

## UNITED STATES DEPARTMENT OF AGRICULTURE

+ + +

USE OF WHOLE GENOME SEQUENCE (WGS) ANALYSIS  
TO IMPROVE FOOD SAFETY AND PUBLIC HEALTH

+ + +

October 26, 2017  
8:00 a.m.U.S. Department of Agriculture  
South Building, Jefferson Auditorium  
14th & Independence Avenue, S.W.  
Washington, D.C.

## USDA:

DR. DAVID GOLDMAN  
Assistant Administrator  
Office of Public Health Science  
Food Safety and Inspection Service  
U.S. Department of AgricultureCARMEN ROTTENBERG  
Acting Deputy Under Secretary for Food Safety  
U.S. Department of AgricultureDR. UDAY DESSAI  
Senior Public Health Advisor  
Office of Public Health Science  
Food Safety and Inspection Service  
U.S. Department of AgricultureDR. PETER EVANS (Moderator)  
Office of Policy and Program Development  
Food Safety and Inspection Service  
U.S. Department of Agriculture

SETTING THE STAGE

DR. MARTIN WIEDMANN  
Gellert Family Professor in Food Safety  
Cornell University

DR. DAVID GALLY  
Professor of Microbial Genetics  
University of Edinburgh

DR. NORVAL STRACHAN  
Chair in Physics and Chief Scientific Advisor to  
Food Standards Scotland  
University of Aberdeen

FEDERAL/STATE COLLABORATION

DR. JOHN BESSER  
Deputy Chief, Enteric Diseases Laboratory Branch  
Centers for Disease Control and Prevention

DR. STEVEN MUSSER  
Deputy Director for Scientific Operations  
Center for Food Safety and Applied Nutrition  
U.S. Food and Drug Administration

DR. DAVID GOLDMAN  
Assistant Administrator  
Office of Public Health Science  
Food Safety and Inspection Service  
U.S. Department of Agriculture

MR. DAVE BOXRUD  
Molecular Epidemiology Supervisor  
Minnesota Department of Health

DR. PATRICK McDERMOTT  
Director, National Antimicrobial Resistance  
Monitoring System  
Center for Veterinary Medicine  
U.S. Food and Drug Administration

DR. WILLIAM KLIMKE  
Senior Scientist  
National Center for Biotechnology Information  
National Institutes of Health

DR. GLENN TILLMAN  
Chief, Microbiology Characterization Branch  
Office of Public Health Science  
Food Safety and Inspection Service  
U.S. Department of Agriculture

DR. CHRIS BRADEN  
Deputy Director, National Center for Emerging and  
Zoonotic Infectious Diseases  
Centers for Disease Control and Prevention

ALSO PARTICIPATING

DR. MUNA NAHAR

DR. MELANIE ABLEY  
Food Safety and Inspection Service  
U.S. Department of Agriculture

DR. CHRISTINE ALVARADO  
U.S. Department of Agriculture

DR. CATHERINE CARRILLO  
Scientist  
Ottawa Laboratory (Carling)  
Canadian Food Inspection Agency

MR. STEVEN ROACH  
Director, Food Safety Program  
Food Animal Concerns Trust (FACT)

DR. MARC ALLARD  
Center for Food Safety and Applied Nutrition  
U.S. Food and Drug Administration

DR. ALEX BRANDT  
FSNS

DR. BETSY BOOREN  
OFW Law

DR. JORGEN SCHLUNDT  
Director  
Nanyang Technological University Food  
Technology Centre, Singapore

MS. SHERRI MCGARRY  
Centers for Disease Control and Prevention

JOSEPH HEINZELMANN  
Neogen

## INDEX

	PAGE
WELCOME	
Dr. David Goldman	7
Carmen Rottenberg	7
Dr. Uday Dessai	9
SETTING THE STAGE	
Advanced Technology (including WGS) to Improve Food Safety - Dr. Martin Wiedmann	12
Advanced Technology (including WGS) to Improve Animal Health and Food Quality - Dr. David Gally	44
Advanced Technology (including WGS) to Perform Microbial Food Safety Risk Assessments and Hazard Identification - Dr. Norval Strachan	70
Questions/Answers	99
FEDERAL/STATE COLLABORATION	
WGS at CDC: Capacity, Communication with Partners, Use in Decision Making Capacity - Dr. John Besser	130
WGS at FDA: Capacity, Communication with Partners, Use in Decision Making - Dr. Steven Musser	144
WGS at FSIS: Capacity, Communication with Partners, Use in Decision Making - Dr. David Goldman	163
WGS at U.S. State Health Departments: Capacity, Communication with Partners, Use in Decision Making - Mr. Dave Boxrud	183

## INDEX (cont.)

	PAGE
Role of WGS in Conducting Surveillance for Antimicrobial Resistant Foodborne Bacteria in the Food Supply - Dr. Patrick McDermott	195
Publicly Available Tools for WGS Analysis at National Center for Biotechnology Information (NCBI) - Dr. William Klimke	215
Demonstration of Tools for WGS Analysis - Dr. Glenn Tillman	238
Genomics for Food Safety (Gen-FS) Consortium - Dr. Chris Braden	255
Questions/Answers	265
ADJOURNMENT	297

M E E T I N G

(8:07 a.m.)

1  
2  
3 DR. GOLDMAN: Good morning. If I could ask  
4 everyone to please take their seats. We're about to  
5 begin. Hopefully you can find a seat.

6 Good morning. My name is David Goldman,  
7 and I'm one of the Assistant Administrators here at  
8 FSIS, the host of this meeting, along with our  
9 multiple partners, and I want to introduce to you our  
10 Acting Deputy Under Secretary for Food Safety here at  
11 USDA, Carmen Rottenberg, who wants to provide you  
12 with the official welcome. Thank you.

13 MS. ROTTENBERG: Thank you, David. I want  
14 to thank all of you for coming here today and also  
15 for those of you who are participating by the  
16 webinar. This has been a busy week in the life of  
17 the Agency. We have had the privilege of hosting two  
18 public meetings about really important information  
19 that furthers our public health goals and really  
20 shows the deep levels of collaboration that we have  
21 with our public health partners.

22 As we've heard the last couple of days,

1 whole genome sequencing technology has become a  
2 routine part of the NARMS surveillance screen for  
3 resistant genes in enteric bacteria. And, as you all  
4 know, whole genome sequencing is now regularly used  
5 in outbreak investigations, but always in the context  
6 of other available evidence.

7           Many of you that are in the industry have  
8 been asking us to have a meeting like this for the  
9 last year, year and a half, and we have a really  
10 robust agenda to go through with you with really  
11 talented scientists and really folks from all areas  
12 where whole genome sequencing touches.

13           We at FSIS nearly completed building our  
14 own whole genome sequencing capacity, and we do  
15 intend to have WGS fully implemented into the  
16 sampling programs this fiscal year. So the timing of  
17 this meeting is appropriate, and there's been an  
18 incredible amount of work on behalf of the agencies  
19 to put this together, and I just want to thank all of  
20 the folks who have worked on that.

21           It's our intention at FSIS to analyze the  
22 whole genome sequencing data using validated and



1 transparent methods, which is also why this 2-day  
2 public meeting is so important.

3           As we move forward, with utilizing whole  
4 genome sequencing, we continue to collaborate with  
5 our public health partners, and I think that you're  
6 going to be able to really see that showcased herein  
7 the next couple of days.

8           So I didn't want to take a lot of time this  
9 morning, but I did want to welcome you all and thank  
10 you again for coming, and again thanks to all of the  
11 speakers who are going to be sharing with us the next  
12 couple of days.

13           So thank you and welcome.

14           DR. DESSAI: Good morning again, and before  
15 we start, I'd like to kind of do a few housekeeping  
16 things.

17           All right. So a couple of things. No food  
18 or drinks are allowed here. Number 2, you have your  
19 restrooms on this side in the Fifth Wing, and the  
20 cafeteria is in the Third Wing. So you have the  
21 reception over there. If you have any questions,  
22 please feel free to ask for help.

1           Now, all the time those who are non-feds  
2 should have their badges when you are moving around  
3 in the building.

4           With that, I just want to say this is an  
5 exciting time and especially this meeting is very  
6 important, and like Carmen said, we are in a phase  
7 here where technology has come to a level where we  
8 can think of inserting some components of that,  
9 possibly into the regulatory process. And later on,  
10 you will see from various speakers where we are  
11 today, where we will be in 5 years, where we go in 10  
12 years, and how this technology can take us or help us  
13 go to our 2030 goals potentially.

14           So we have a lineup of speakers who will  
15 basically take you through what is whole genome  
16 sequencing and what it can do and how it can take us  
17 where we want to be. What is whole genome sequencing  
18 and risk connection? And many topics which are of  
19 interest as well as where you need clarity in terms  
20 of moving forward.

21           So we have three speakers in the first  
22 module, and that is called Setting the Stage, and

1 these threes speakers will basically lay the ground  
2 for you.

3           The first speaker we have is Martin  
4 Wiedmann from Cornell. The second speaker we have is  
5 David Gally from Edinburgh Institute, UK, and Norval  
6 Strachan is our third speaker who will be talking  
7 about risks and WGS.

8           Now, Norval and David have traveled. So if  
9 they are sleepy at times, just kind of bear with them  
10 because they have 6 hours of jetlag, okay.

11           Having said that, a couple of things to  
12 keep in mind. There are two units here. One is for  
13 you guys here. The second is it's a webinar. So  
14 your speeches are streamed. The slides are being  
15 streamed, and keep in mind please speak in the  
16 speakers so people on the line can also hear you very  
17 clearly.

18           Now, we don't have a clicker yet. So what  
19 you have to do, the speakers is, just indicate to our  
20 wonderful folks here, and they will change the slides  
21 for you. Okay.

22           Having said that, over to Martin.

1 DR. WIEDMANN: Thank you. Good morning,  
2 everyone. So what I'm going to do is give sort of  
3 general overview of the whole genome sequencing.

4 Next slide, please.

5 What I want to emphasize, before we get  
6 started, in more into details, is that we really need  
7 to look at whole genome sequencing and omics tools in  
8 a grander context of what might be actually pretty  
9 disruptional technology changes in food safety, and  
10 I'm going to call the place that we're going to go  
11 to, precision food safety. It's not very innovative,  
12 but I think what I'm going to try to sort of impress  
13 on you is that we can't look at whole genome  
14 sequencing technology by itself.

15 We need to look at it in terms of the other  
16 changes in tools and technologies that we will have  
17 available to assure food safety and how these tools  
18 will help us to improve food safety from improved  
19 outbreak detection, to improved source tracking,  
20 improved identification of pathogens and will help us  
21 to move from a typically reactive approach to a much  
22 more proactive but ultimately predictive approach.

1           And some of the other tools that are  
2 important in this framework are really the ideas of  
3 machine learning, artificial intelligence, GIS or  
4 global information systems technologies, that will  
5 give us this greater precision in food safety, and it  
6 will work along with this whole genome sequencing.

7           Next slide.

8           So here's what I'm going to run through  
9 today. I'm going to provide a brief overview and  
10 then go into a few case studies or a few areas where  
11 whole genome sequencing already has a considerable  
12 impact or will continue to have a major impact.

13          Next slide.

14          So first I'm going to just go over some  
15 sequencing technology, data analytics at a sort of  
16 very high level to set us all up at the same page.

17          The reason we're here is that sequencing  
18 technologies have developed dramatically from first  
19 generation, Sanger type sequencing, that were fairly  
20 costly and did allow us to do whole genome sequencing  
21 but not anything what we can do now, to next-gen  
22 sequencing tools that some people divide into sort of

1 a second generation, with Illumina being the work  
2 horse of what we do today, to now moving to third  
3 generation sequencing tools, and these are the tools  
4 that allow us to do large-scale whole genome  
5 sequencing, but they also are used for applications  
6 in food safety that are not actually whole genome  
7 sequencing. People sometimes use WGS and NGS, next-  
8 gen sequencing interchangeable, but I'm going to try  
9 to separate that out a little bit for you.

10           Next slide.

11           The innovations we see, and we're going to  
12 continue to see are probably best illustrated with  
13 this picture. We're not just going to scale down  
14 costs, but also scale down size of the equipment.  
15 This is a picture of the MinION or MinION, however  
16 you want to pronounce it, which is really where I see  
17 and think a lot of others see the future of whole  
18 genome sequencing, that's miniaturized to a spot  
19 where we can get sequence data more quickly with this  
20 smaller equipment and, you know, continuing to reduce  
21 cost. So this is sort of the future and you can see  
22 the scale of this sequencing equipment here.

1           Next slide.

2           So as I said, whole genome sequencing has  
3 been performed as traditional sequencing but the  
4 workers are these next-gen sequencing tools, but the  
5 important thing is that next-gen sequencing can be  
6 used for other applications.

7           And some examples of, there are  
8 metagenomics, which I'll quickly touch on, where you  
9 can basically sequence all DNA found in a sample.  
10 All genetic material, take it all and sequence it.

11           And you can also use these same tools to  
12 highly parallel sequence many genes at the same time.  
13 We're used to PCR-ing one gene and sequencing it.  
14 You can do highly parallel sequencing, that can  
15 address some of the issues with whole genome  
16 sequencing that we're going to get hits, we're going  
17 to get information that we don't know what to do  
18 with. So we can use to target hundreds of thousands  
19 of genes and really only look for pre-defined  
20 targets.

21           And then we move into the area of RNA  
22 sequencing where you don't sequence all the DNA but

1 you actually sequence RNA. RNA is an unstable  
2 genetic material. It's not consistently unstable,  
3 but it might give us some better ideas of what these  
4 organisms do and may also help us in some cases  
5 differentiate live and dead organisms.

6           So these are all applications of next-gen  
7 sequencing, but they're not whole genome sequencing  
8 per se.

9           Next slide.

10           I'm going to move a little bit into the  
11 data analytics now, focusing really on, you know,  
12 application of whole genome sequencing to bacterial  
13 genomes. It can be used for parasites. It can be  
14 used for viruses, but the main application that I  
15 want to emphasize is really whole genome sequencing  
16 of bacteria.

17           So the average bacterial genomes that are  
18 relative to food safety range somewhere from 2 to 10  
19 Mb per organism. So that's the range of what we deal  
20 with.

21           And if you look at bacterial genomes, we  
22 typically can differentiate the genomes in terms of a



1 core genome. So if you take something like *Listeria*,  
2 these are the genes that are found in all *Listeria*  
3 *monocytogenes* and the accessory genome. These are  
4 hundreds of thousands of genes that are found in some  
5 *Listeria* or some organisms, but not others.

6           These can play very important roles because  
7 they might provide an antimicrobial resistance as we  
8 mentioned before, but at least in a number of the  
9 analysis we use, we may ignore and not use that  
10 information. So this is a set of data that can  
11 provide lots of very valuable information that may  
12 not always, and I emphasize, not always, be used.

13           Also the question always comes up, you  
14 know, so you scientists call this whole genome  
15 sequencing. Do you really sequence the whole genome?  
16 And, the reality is, we sequence the whole genome but  
17 we don't assemble the whole genome. So very often  
18 there are certain pieces of the genome that are  
19 missing for a variety of reasons including genes that  
20 you have repeated multiple times. So, yes, it is not  
21 the whole genome that we analyze even though we call  
22 it whole genome sequencing.

1           Next slide.

2           So as we look at sequencing, if you have  
3 this 2 through 10 Mb per genomes, as we sequence it,  
4 it's basically a giant puzzle. We take this genome,  
5 cut it into small pieces, this is what you see on the  
6 top, and it goes through some experimental parts to  
7 make these pieces ready for sequencing, put them on a  
8 sequencing platform and then we end up with millions  
9 of pieces of DNA that we now need to put back  
10 together, which is what we call assembly.

11           Broadly speaking, there are two different  
12 approaches to it. One of them would be reference  
13 guided assembly, where we say, this is a *Listeria*  
14 *monocytogenes* as we sequence. Let's pick a similar  
15 *Listeria monocytogenes* and use that to put these  
16 pieces back together to get a sequence that we can  
17 analyze.

18           Now, we can also analyze these SNPs by  
19 itself but in most cases, at some point at least,  
20 we're going to do some sort of assembly.

21           And you can do a *de novo* assembly, where  
22 you put the pieces together without that sort of

1 template reference guiding and then obviously relies  
2 on a more massive computational power.

3 Both of these tools are used. There's pros  
4 and cons. I'm not going to get into all the details  
5 of it, but often it's a combination of the two of the  
6 -- different approaches to assembly, but typically  
7 they will lead us to very, very similar answers.

8 Next slide.

9 The big discussion when we talk about whole  
10 genome sequencing, what most people look at really as  
11 an output is single nucleotide polymorphisms, SNPs,  
12 and I'm sort of simply going to explain those, and  
13 this picture, what we're really looking at here is a  
14 difference of one A, C, T, and G. As you see in the  
15 top, you might have one sequence that has a C, the  
16 other one that has a G, and if you have that over 3  
17 million nucleotides of *Listeria*, you say these two  
18 *Listeria* differ by one SNP.

19 Now, there's a similar difference or  
20 another difference that you can find between two  
21 genomes that is also important as we look at the  
22 analysis used, which are insertion or deletion

1 indels. And we call them indels because you really  
2 don't know whether the C in the sequence on the top  
3 was deleted and therefore got us to the reference or  
4 whether the reference acquired it. So it could be an  
5 insertion. It could be a deletion. It's a chicken  
6 and egg question.

7           When we have these differences, they're  
8 often referred to as single nucleotide variance as  
9 opposed to SNPs because it's a variant but it's not a  
10 polymorphism. So if you hear people talk about SNPs,  
11 it typically does not include indels unless people  
12 use the nomenclature somewhat imprecisely.

13           Next slide.

14           The simplest presentation of how we analyze  
15 these data is here. We have four sort of  
16 hypothetical sequences. In red, you can see  
17 differences. So you can see Isolate 1 and 2 differ  
18 by a single T that is marked in red in the middle,  
19 and the outputs you will typically see from these are  
20 trees, shown on the left, where you can see that 1  
21 and 3 are identical, 2 is similar to it, and 4 is  
22 more different, or I can do what we call SNP

1 matrices, where we come to pairwise comparisons, and  
2 I can look at 1 versus 3 and I see a 0 meaning  
3 they're the same, there's no SNP differences, or I  
4 can look at 1 versus 4 and I can see there's 3  
5 differences.

6           And you look at this with a few sequences,  
7 it looks pretty straightforward. I can tell you same  
8 or difference, right.

9           Next slide.

10           But as we expand this and not just have  
11 four isolates, but have 100, 1,000, 5,000, 10,000,  
12 and we build these trees, they obviously become more  
13 difficult to interpret. Now, we have to curl them  
14 up, so we can put them on one slide. Otherwise, no  
15 one can see them, and even when we curl them up, no  
16 one can see them. So it makes interpretation and use  
17 of these trees difficult and challenging.

18           Obviously, you can zoom in on sub-trees of  
19 that and just look at a specific plate, a specific  
20 subgroup of interest, but even that can become  
21 challenging partially because some of these branches  
22 will change as we add new information and new

1 isolates to it.

2           So next slide.

3           So one approach that is used and is  
4 increasingly used and gets us around, is this idea  
5 called multilocus sequence typing. I call this whole  
6 genome multilocus sequence typing. There's really  
7 two subsets to this. One of them is whole genome  
8 MLST that uses all the genes. The other one is core  
9 genome MLST, and as I mentioned before, it's a core  
10 genome. Those are the genes that are common to all  
11 organism. In the case of core genome MLST, I will  
12 just use these core genes that are common to all  
13 organisms.

14           Those tools are being used, are being  
15 further developed and different groups look at  
16 different tools. Obviously, when we look at our core  
17 genome MLST, we will lose some information because  
18 we're not using some of those genes. They're only  
19 found in some isolates, but having core genome MLST  
20 makes some things in terms of analytics simpler.

21           So how does MLST work? We have a database.  
22 We define unique loci. So those will be genes but

1 also some other similar types of structures such as  
2 non-coding RNAs, and then any change in this gene,  
3 whether it's a SNP, whether it's a insertion or  
4 deletion, equals a new allele and we can name these  
5 alleles.

6           You can see this most easily on the bottom.  
7 So we have hypothetical locus 1. We have allele 1.  
8 Allele 2 differs by 2 SNPs. So we call it 2. It's  
9 different. The next one differs by 1 deletion from  
10 allele 1, also differs from allele 2. *Therefore*, we  
11 call it allele 3.

12           In this scheme, the numbers don't indicate  
13 similarity. You cannot go saying, well, this is 1  
14 and 2, so they're similar, while allele 1 and 100 are  
15 very different. The numbers are simply added as we  
16 identify alleles. So 1 and 2 can be very different,  
17 and 1 and 100 might only differ by 1 indel.

18           Next slide.

19           So what we can then do on a larger scale  
20 and what is done at a larger scale is we put all  
21 these data together. I've got your hypothetical 3  
22 isolates, A, B, C. I've got your number of loci. So

1 locus 1, all three isolates are the same. Locus 2,  
2 isolate C is different. So it's named number 2.  
3 Same for locus 2,005 and we then look at all of them  
4 together. So if two isolates are the same in all of  
5 these loci, they're given one single name A or a  
6 longer numerical designation.

7           So it takes this mass of information or  
8 relatively mass of information for genome sequencing  
9 and really drills it down to one number or a set of  
10 numbers whereas we will hear later as type of zip  
11 code. Okay. And so it really takes this information  
12 at least as a superficial level and makes it much  
13 easier to deal with it and compare between labs, to  
14 tell me same or different, and we can even try to  
15 come up with some of these naming schemes to make  
16 them so that the similarity of numbers at least gives  
17 us some idea how closely related those isolates are,  
18 just like with zip codes.

19           So these are the two high quality SNP and  
20 MLST are the two sort of core tools that are  
21 currently predominantly used to really analyze this  
22 whole genome sequence data.



1           Next slide.

2           Now, regardless of which method we do, one  
3 of the core questions that it comes down to is, so  
4 how quickly do two isolates become different, and to  
5 me that's one thing that's cause many of the  
6 questions. And these are a couple of datasets that  
7 were derived on reasonably large datasets to get us  
8 at this question.

9           So the first one is a large study on  
10 *Listeria monocytogenes*, which estimates about 2.5 x  
11 10<sup>-7</sup> substitutions per site, so per nucleotide, per  
12 year. What does that mean? You're going to get one  
13 SNP difference in the core genome every 2.5 years.  
14 That's a reasonably long time for some people and for  
15 evolutionary biologists, it's probably a pretty short  
16 time.

17           If you look at core genome MLST types, when  
18 does a type become different as organisms multiply?  
19 In *Listeria*, it takes about 0.2 alleles per year. So  
20 that means in 5 years, an organism will change.

21           That obviously has huge and important  
22 implications for use of these tools, right. If we

1 have the proverbial goose that flies from Florida to  
2 Canada, along the flyway, and leaves little drops of  
3 poop all the way along with *Salmonella*, that  
4 *Salmonella* is probably the same all the way along the  
5 way, and when it comes back, it's still the same.

6           If you look at a breeder flock in poultry,  
7 and we have *Salmonella* in that breeder flock and that  
8 *Salmonella* is transmitted across the chain, down to a  
9 slaughter house and maybe a retail establishment,  
10 that *Salmonella* probably could stay the same and we  
11 could find that identical *Salmonella* in different  
12 locations. I think that's an important thing that we  
13 really need to consider.

14           Another estimate on *Salmonella* Cerro, a  
15 specific *Salmonella* Cerro type, you see it's in the  
16 same range of substitution rates that were estimated.  
17 So that gives you *Salmonella* genome slightly larger,  
18 gives you some of the same ideas.

19           So the challenge here is that evolutionary  
20 biologists will talk about most recent common  
21 ancestor. That's when two isolates have this most  
22 recent common ancestor defined by sequence. Two

1 isolates with one SNP difference, that most recent  
2 common ancestor might be 4 or 5 years ago or longer,  
3 and that's why we need these tools along with  
4 epidemiology and other evidence to interpret it.

5 Next slide.

6 So here's the summary on the data analysis.  
7 Obviously, it does involve many steps, many of which  
8 I skipped over. We have done considerable  
9 advancements in standardization and transparency of  
10 those tools which will probably be covered later.

11 We have a number of different approaches  
12 for analysis, but in most cases, they will lead us to  
13 the very same answer.

14 The other important thing to remember, once  
15 you have a whole genome sequence, I can analyze it  
16 with a high quality SNP, I can use the same raw data  
17 and run a core genome MLST, whole gene MLST, whatever  
18 I want to.

19 And the important, you know, caveat there  
20 is obviously that these analyses establish recent  
21 common ancestors but do not establish relevant  
22 epidemiological links. They point us in the

1 direction, they can help us support some conclusions  
2 but recent common ancestor, at least in my book, does  
3 not equal cause and effect.

4 Next slide.

5 So now I'm going to move through some of  
6 the sort of areas where whole genome sequence has  
7 significant impact, and the first and foremost is  
8 obviously outbreak detection.

9 Next slide.

10 And we've had PulseNet where we use  
11 basically barcode type fingerprints to track  
12 foodborne diseases to do surveillance for a long  
13 time, had tremendous impact in improving our ability  
14 to detected foodborne disease outbreaks, which you'll  
15 hear more about later.

16 Next slide.

17 But there's been some challenges with this,  
18 and this is a slide from some work we did that  
19 probably illustrates the best. This is a number of  
20 *Salmonella* Montevideo isolates that we collect by  
21 PFGE. When you look at them, they all look the same.  
22 Only problem was that some isolates came from

1 pistachios in California and some came from an  
2 outbreak, and it was linked to sausage produced in a  
3 facility in Rhode Island and raw isolate from the  
4 pepper, but they're all the same. So very, very hard  
5 to use these data to sort of give us some insight.

6           When we did whole genome sequencing on  
7 this, shown in next slide, all the isolates that came  
8 from the pepper, the sausage and associate human  
9 cases, clustered in this tree on the top right. We  
10 put them all together. We can clearly show that they  
11 are different even though they have all of these  
12 isolates shown in this tree have the same PFGE  
13 pattern.

14           More importantly, as is shown with the red  
15 circle, we find a set of five isolates that are very  
16 closely related by SNPs that are short time frame in  
17 the same state that likely represents another  
18 outbreak within this large cluster of similar PFGE  
19 types they would have never recognized with PFGE  
20 alone. But with this whole genome sequencing, we can  
21 separate them out and see that.

22           So this illustrates, in one very specific

1 example, why whole genome sequencing is so powerful  
2 for outbreak protection.

3 Next slide.

4 The long-term impact of this is probably  
5 best illustrated by *Listeria* where we have a lot of  
6 history using whole genome sequencing. You can see  
7 before '97, before PFGE, one outbreak every 3 years.  
8 Once we started using PFGE, about two to three  
9 outbreaks a year. Once we started whole genome  
10 sequencing, somewhere in the range of 7 to 10  
11 outbreaks a year. One outbreak every year, let's say  
12 10 outbreaks a year, 30 times more outbreaks detected  
13 per year. It's not that we had more *Listeria*  
14 problems. It's simply it detected outbreaks better.

15 The important part is that there are the  
16 size of the average outbreak went from about 7 year  
17 outbreak to somewhere for 3 to 4 per outbreak with  
18 whole genome sequencing today. So we detect more  
19 outbreaks, but we detect smaller outbreaks. That  
20 means we detect earlier on. We detect outbreaks that  
21 we previously would not have.

22 Next slide.

1           And so obviously this came along as routine  
2 implementation of whole genome sequencing in  
3 September 2013.

4           Next slide.

5           One other organism that this is going to  
6 have a major impact on is going to be *Salmonella*  
7 Enteritidis. *Salmonella* Enteritidis, about 50% of  
8 *Salmonella* Enteritidis have the same PFGE type.  
9 PFGE, the routine surveillance method. Molecular  
10 surveillance is not very good at differentiating  
11 Enteritidis. Once you apply whole genome sequencing  
12 to this organism, shown in the next slide, we're  
13 going to differentiate these isolates with whole  
14 genome sequencing to the point now where we can  
15 detect larger number of clusters and ultimate  
16 outbreaks very similar to what I've shown you in  
17 Montevideo.

18           So at the beginning of this paradigm shift,  
19 we have seen it in *Listeria*, but I think we will see  
20 some tremendous impact on *Salmonella*, and I predict  
21 *Salmonella* Enteritidis is one of the organisms that  
22 we will particularly see these impacts but also other

1 *Salmonella* serotypes.

2           In addition to better detection of the  
3 outbreaks, whole genome sequencing also provides  
4 tremendous opportunity for better trace-back, and  
5 food processing plants that can be used by both  
6 industry but also is used by regulatory agencies.

7           One early example of this that actually  
8 goes along well with this sort of data on SNP  
9 differences was some work we did in collaboration  
10 with CDC published in 2008 where we actually  
11 sequenced *Listeria monocytogenes* from a listeriosis  
12 case and a listeriosis outbreak that happened 12  
13 years apart but were linked to the same facility. We  
14 sequenced the *Listeria monocytogenes* from those two  
15 years, and we found that at least some of them 1988  
16 and 2000 isolates, 12 years apart, differed by a  
17 single SNP. Very stable, probably because this was a  
18 ready-to-eat food processing facility survived in  
19 this facility over time and differed and changed very  
20 little over that time.

21           Now, if you apply this to processing  
22 facilities, and this is an example of a non-whole



1 genome sequence based data, but it illustrates the  
2 point, you can see with subtyping and here we colored  
3 different subtypes, provide different subtypes in  
4 different colors, you can in this picture see which  
5 is a 2-year surveillance of a processing facility  
6 that we have a very specific *Listeria* shown in green,  
7 more or less survive in this facility over 2 years.  
8 That's obviously a challenge, and if you do proper  
9 root cause analysis in this case, we could traceback  
10 the persistence of this *Listeria monocytogenes* to a  
11 very specific location in the facility, which ended  
12 up being floor mats which obviously provided us with  
13 the opportunity to just remove these floor mats and  
14 see if our hypothesis was right, and lo and behold,  
15 once we removed these floor mats, that type of  
16 *Listeria* was not found in the facility over a 6 month  
17 follow up.

18           So it shows the power of using whole genome  
19 sequencing and subtyping data to traceback and  
20 identify sources that's in a facility or throughout a  
21 supply chain.

22           Now, that is obviously one thing that, you

1 know, causes some people concern and heartaches and  
2 headaches, right. And this is sort of the  
3 hypothetical case study I want to speak through, talk  
4 through, if you take this to the next step.

5 Let's say you have a facility that has a  
6 *Listeria monocytogenes* positive in a finished  
7 product, one day's production. Typically you end up  
8 with a recall of one lot that was produced that day.

9 Let's say 11 months later, you have another  
10 positive for the same facility, a different type of  
11 product, *Listeria monocytogenes* again. Now, you're  
12 in sort of a tougher spot because it's like it could  
13 be a repeat positive of the same subtype. It could  
14 be a separate issue. If you now have whole genome  
15 sequence data on the January and December isolate,  
16 match by whole genome sequencing, let's assume zero  
17 differences, that will lead you down to the  
18 conclusion that that *Listeria monocytogenes* probably  
19 perhaps persisted in that facility.

20 Obviously, you need some additional  
21 information. Obviously, could be a reintroduction in  
22 that facility, but a conclusion that some people

1 could draw and might draw is that any of the food  
2 produced within January and December was produced  
3 under unhygienic conditions. If I'm going to  
4 reintroduce the same *Listeria* time after time, that a  
5 food processing facility has under control, that's  
6 arguable.

7           The challenge then becomes when we take  
8 this sort of information and extrapolate to non-  
9 ready-to-eat foods, let's, for argument's sake, say  
10 raw poultry or raw meat, where we can now have  
11 reintroduction, *Salmonella* is endemic. It's found  
12 regularly in poultry farms potentially, in dairy  
13 farms where we get ground beef, and so it could be  
14 truly a reintroduction. So it could not be an issue  
15 with the facility that cook upstream.

16           So it illustrates that we can't just  
17 extrapolate from ready-to-eat facility to non-ready-  
18 to-eat facility. We need to consider the overall  
19 supply chain as we interpret these data and really do  
20 our epidemiological investigations, even when we're  
21 just talking about food contamination, not human  
22 disease cases.

1           The next one I want to move to is how to  
2 use whole genome sequencing to better understand and  
3 define pathogens.

4           Next slide.

5           So the case study on this one is something  
6 that some of you may be familiar, large recall of  
7 Fonterra because they found *Clostridium botulinum* in  
8 their powder. The story there was that after 4 to 6  
9 weeks later, they suddenly discovered it wasn't  
10 *Clostridium botulinum*. It was actually *Clostridium*  
11 *sporogenes*. How did they identify that?

12           Ultimately probably with whole genome  
13 sequencing. The challenge here is that *Clostridium*  
14 *botulinum* and *Clostridium sporogenes* are very, very  
15 similar. Unless you do PCRs or mass experiments, you  
16 cannot differentiate them. If you use whole genome  
17 sequencing, you can differentiate these close-related  
18 organisms very quickly, very easily. This is not  
19 just a one time deal.

20           We published this paper recently where we  
21 had a similar incident. We found a *Clostridium*  
22 species that some people could have worried about

1 being *Clostridium botulinum*. With sequencing, we  
2 actually found it was a new species that fell  
3 somewhere in the proximity of *sporogenes* and  
4 *botulinum*.

5           So a great tool to rapidly differentiate  
6 organisms and give us better species classification  
7 or classification into pathogen or not. This just  
8 doesn't apply to *Clostridium*. Same issue with  
9 *Bacillus*, *Bacillus cereus*, *Bacillus anthracis*,  
10 *Bacillus thuringiensis*, a group of closely related  
11 organisms, difficult to differentiate, but with whole  
12 genome sequencing, you can differentiate them quickly  
13 and decide food safety hazard, yes or no, and  
14 sometimes obviously it's more of a gray zone, but in  
15 some of these cases, it's very easy to decide but  
16 only if you use these tools.

17           Now, where it's going to get more exciting  
18 than just taking existing pathogens and saying, is it  
19 one or is it not, if we now apply these tools to non-  
20 pathogen groups, say *Listeria monocytogenes* and say  
21 are all *Listeria monocytogenes* the same or can we  
22 differentiate different subgroups that are less

1 likely to cause disease.

2           The work on *Listeria monocytogenes* I'm  
3 showing you here was actually built on some initial  
4 work where we found certain subtypes of *Listeria* that  
5 were not defined by whole genome sequencing, that  
6 were very common in food, about 30% of food isolates,  
7 but very rare in human isolates, about 2% of human  
8 isolates.

9           The question was why do these isolates show  
10 up in food but not humans? When we looked at DNA  
11 sequence data, we could identify a mutation in one  
12 key gene in *Listeria*, *inlA* which allows it to attach  
13 to human cells. These *Listeria* which were found  
14 common in foods, rarely in human cases, had a  
15 different form of this protein that did not allow  
16 *Listeria* to attach to human cells. Therefore, they  
17 were much less likely to cause human disease. So now  
18 we can take a known pathogen and say they're not all  
19 the same hazard. There's considerable differences,  
20 and we beat this horse to death with a number of  
21 studies including animal experiments, etc., to show  
22 that this *Listeria* by a thousand-fold less likely to

1 cause human disease. So tremendous improvements  
2 there.

3           It doesn't stop at *Listeria*. Here's an  
4 example of *Salmonella* Cerro which is an organism  
5 which is very common in cattle. And as shown here,  
6 we never found it, rarely or almost never found it in  
7 human cases. Performed whole genome sequencing on  
8 *Salmonella* Cerro, identified a number of mutations in  
9 very specific genes that are important for this  
10 organism to cause human disease. Mutations there  
11 probably means it can cause human disease, consistent  
12 between epidemiology and whole genome sequencing,  
13 much less likely to cause human disease.

14           So we can use these tools and hazard  
15 characterizations and say, not all *Salmonella* are the  
16 same. Some of them we can define probably pretty  
17 well that they're a reduced human health hazard. So  
18 that's another great application that might be a  
19 little bit more in the future, but I think a very  
20 important one.

21           And what I want to end up with is  
22 metagenomics. So now we're not sequencing whole

1 genomes, but sequencing other DNA in an organism, and  
2 what I want to envision there, and this future  
3 already is there, I want to envision a new type of  
4 audit. Now, you audit your facility in a far away  
5 foreign country, you collect samples of an ingredient  
6 that you source, for example, pepper, you  
7 characterize them and then your incoming lots are  
8 characterized at regular intervals with that same  
9 type of analysis to just find out if what you're  
10 getting is similar to what was in the facility, was  
11 produced in the facility, when you did an audit.

12           And what you might end up there, and this  
13 is a hypothetical example is hypothetical example,  
14 these could be four samples collected during the  
15 audit. You get different bacterial species in them,  
16 you get a bacterial species profile. This is your  
17 first lot that comes in your plant, looks similar.  
18 This is your next lot that comes into your facility  
19 to test, but then you suddenly get this lot of  
20 pepper. Looks different, right. So obviously  
21 something is different. We don't know whether it's a  
22 food safety hazard or not, but we know it's a



1 significant deviation.

2           Is our audit for that facility still valid  
3 or do we need to re-audit that facility? We need to  
4 know what's going on.

5           We could look at these data, but we can  
6 also use advanced tools such as machine learning to  
7 define those deviations. That's where the future is.

8           So that's where future application of whole  
9 genome sequencing and next-gen sequencing will go.

10           So what are the challenges which is  
11 obviously one of the reasons we are here? I've  
12 outlined most of them, but I'm going to try to  
13 summarize them.

14           One key challenge is obvious. We can find  
15 bacteria with very few or no SNP differences in  
16 different locations, food and food associated  
17 environments. WGS rarely will give the final answer.  
18 It will point us to a certain point, but we need  
19 epidemiology, we need other evidence and the evidence  
20 needs to be combined.

21           How do we combine that evidence? Very  
22 often this might require new tools, right. Now, it's

1 preponderance of evidence. We find an unhygienic  
2 facility. We find a whole genome sequence match. We  
3 find this and we find that, but I think some of these  
4 new tools of artificial intelligence and machine  
5 learning potentially can help us to better combine  
6 this piece of evidence.

7           Another challenge is metagenomics which  
8 detects both live and dead cells. Presence of  
9 certain genes is a public health hazard. We find an  
10 antimicrobial resistance gene in our food, but  
11 unknown if it's a live organism or even a pathogen,  
12 do I need to worry about. You may need a new risk  
13 assessments for presence of genes.

14           And then obviously still considerable  
15 uncertainty around data interpretation, different  
16 data analyses approaches, and these affect  
17 industries' willingness and ability to use those  
18 tools when they probably should use them. Define two  
19 *Listeria monocytogenes*, 12 months apart in their  
20 facility, they should use the best tools to find out  
21 whether it's the same *Listeria* or not.

22           In the current climate, sometimes people

1 are afraid to use those tools.

2           One of my ideas is that it might be time  
3 for a moratorium where we simply allow industry to  
4 use these data, not as a fear of having these data  
5 required to turn over to agencies, but be able to use  
6 them themselves to figure out how they can best use  
7 it in their context and then after while, start to  
8 come back to discussions about constant sharing of  
9 these data. But that's just one of my opinions.

10           Conclusions: Precision food safety is  
11 here. Improved outbreak detection, improved  
12 surveillance, improved source tracking, and improved  
13 bacterial identification due to whole genome  
14 sequencing is the new reality. It's happening  
15 already, maybe not at the penetration some of us  
16 would like to see, but it's happening already.

17           The roadmap for other uses of whole genome  
18 sequencing and next-gen *Salmonella* is less clear.  
19 Will metagenomics and whole genome sequencing replace  
20 hygiene indicators? I've given you some of these  
21 ideas in terms if you find the same *Listeria* over  
22 time, if you find the same metagenome over time, what

1 does that tell us?

2 Will whole genome sequencing change the  
3 approach to defining hazards where we move away from  
4 bacterial species but move to bacterial, to clonal  
5 groups, subtypes, presence of genes, presence of  
6 mutations, to define hazards? Will that give us  
7 better information and better risk-based tools for  
8 management of food safety?

9 No matter where we're going to go with  
10 this, these are new tools that will require new  
11 people with new training. So it's going to be very,  
12 very important that we train not just food  
13 scientists, but everyone who works with the food  
14 industry and in public health, around food safety, to  
15 use these tools.

16 Thank you very much.

17 DR. GALLY: Okay. Right. I thank you for  
18 the opportunity to speak today. It's been 24 years  
19 since I've been in Washington. So it's a long time.  
20 I was last here when I was doing a postdoc in North  
21 Carolina and drove up here in an old VW Rabbit that  
22 kept breaking down, a VW Golf as they're known

1 everywhere else.

2           So I'm David Gally. I'm based at the  
3 Roslin Institute, just outside of Edinburgh. Famous  
4 for Dolly the sheep. We like to pride ourselves in  
5 the fact we have expertise in genetics and genomics  
6 in livestock but there's a whole host of strange  
7 microbiologists that hang out there as well, and  
8 they're very interested in the genetics and genomics  
9 of bacteria, and certainly that's what I'm really  
10 going to carry on the theme of Martin's introduction.

11           I think you're going to get some similar  
12 things that was asked to set the scene, predominantly  
13 in a One Health perspective and just really to make  
14 the obvious point that in terms of One Health, of  
15 course, food safety, we're very interested in  
16 transmission from animals to humans, but also where  
17 is the role of the environment in that, and certainly  
18 the transient role and the persistent role of the  
19 environment in transfer of bacteria.

20           So as you've heard, I'm trying to work out  
21 the best way to deliver this with what's been heard.

22           So the key is that diversity is

1 understandable from the sequencing. We're getting  
2 amazing insights into the diversity of the bacterial  
3 world from sequencing, and it's really challenging,  
4 the whole taxonomy of bacteria, in fact, but it's a  
5 beautiful insight into that diversity.

6           The bottom line as we've already heard is  
7 the precision, where you have an organism that has 5  
8 million bases. We can really look down at  
9 differences of just a few, and that's way more  
10 precision than we've ever had before. So that's  
11 absolutely critical.

12           Within that before, as you've heard, we can  
13 identify very related organisms and certainly  
14 identify sub-clusters that are more of a threat to  
15 human health, and that becomes important for  
16 prediction capacity which is really the second half  
17 of my talk, and will really differ from what Martin's  
18 told you so far.

19           The tracking as well, very, very important.  
20 To bear in mind, it's not just a one-way street.  
21 We've got plenty of examples of where we're getting  
22 flowback from humans to animals. So this becomes

1 very, very interesting. It becomes testable and  
2 trackable with whole genome sequencing where we can  
3 actually identify flowback into livestock species.

4           We can therefore also with this precision  
5 identify the origins and vehicles of transmission,  
6 and really the key, the second half of the talk as I  
7 said, will be how can we use this information to  
8 improve prediction value of this information? To  
9 actually worry about the risk of all subsets, a  
10 particular same threat to us and the answer's often  
11 no, and we should be able to understand that more in  
12 more detail using these technologies.

13           And, of course, it's not just about  
14 tracking whole bacteria. As you've heard as well,  
15 we're particularly interested in, at the moment, it's  
16 very high on the agenda, in terms of antimicrobial  
17 resistance genes and being able to identify and  
18 attribute really sources for those and transfer of  
19 those. And again that's possible through these  
20 technologies.

21           So, traditionally, we would have had our  
22 microbiology, if we're lucky, and grow our organisms

1 that we're actually talking about which we obviously  
2 can't for many, but we have our sample, be it  
3 directly from the animal or food -- we can carry out  
4 some classic microbiology in the lab and identify our  
5 colonies. It's going to take us a little while to  
6 maybe determine through subsequent testing, often  
7 serotyping, PCRs, additional tests, exactly the  
8 subtypes that we're dealing with, and obviously with  
9 those bacteria, we can determine antibody resistance  
10 using sort of standard plating and techniques.

11           As we have heard, for whole genome  
12 sequencing, we really still need to focus on the  
13 first agar plate there. We still need to get hold of  
14 our individual isolate and then sequence it, so we  
15 know exactly the sequence related to the isolate that  
16 we have.

17           To make the obvious point, that the more  
18 information we can get on the sequences of specific  
19 isolates, the more we'll be able to type that  
20 information into a database and understand what's  
21 present when you analyze samples in a metagenomic  
22 way. So where you can go direct to the more complex



1 sample. That's going to be a challenge in terms of  
2 the databases and how we share that information which  
3 is going to be critical to how we progress in this  
4 science.

5           We then have the analysis side which Martin  
6 has described in detail. And what the intent of the  
7 potential, the idea is here that from the sequence of  
8 the organism, we can get obviously what bacteria it  
9 is, what subtype, potentially what virulence genes it  
10 has and what resistance genes, and really start to  
11 fit it into the epidemiology of previous exposure to  
12 that organism. So that's fairly obvious. We've had  
13 that covered.

14           So a lot of the work that we do in Scotland  
15 is based around enterohemorrhagic *E. coli* 0157, and  
16 I'm going to use that and *Salmonella* as my two key  
17 examples to explain some of the basic, of some of the  
18 more futuristic ways of going about this.

19           So just to show on the left here, the kind  
20 of orangey thing is a cell. We've got *E. coli* 0157  
21 colonizing that cell. You've got the concept that it  
22 produces Shiga toxin, and that's the main

1 pathological determinant in terms of human infection.

2           It colonizes arsenic. It colonizes cattle  
3 and other ruminants using a type 3 secretion system  
4 which injects proteins into that cell to help it  
5 colonize.

6           We have about 1,000 cases in the UK each  
7 year. It's harder to get the estimates of the whole  
8 of USA but well over 10,000, anything up to 50,000,  
9 depending on the literature that you read.

10           Originally known as the "burger bug," but  
11 actually in the UK now we have a lot more cases  
12 associated with direct contact with animals and cases  
13 where produce particularly from vegetables, etc.,  
14 that have been contaminated potentially with  
15 irrigation water and that's the source of human  
16 infection.

17           So with the sequencing, we can get  
18 information that feeds back to the very, I suppose,  
19 quite straightforward and additional sequencing  
20 methods. So shown here on the top left is the inner  
21 and outer membrane of the bacteria or the  
22 polysaccharide. So the actual O type can be

1 determined successfully from the sequence. Okay. So  
2 that's good. So you can say it's an O157, O26, one  
3 of the gang of six, etc. So that can be determined  
4 on the basis of the sequence.

5 Other aspects as well in terms of flagella  
6 type, etc., can all be determined that way.

7 You can also go in and look if it carries  
8 type 3 secretion system, affect the proteins,  
9 different types of Shiga toxin that are involved in  
10 different pathologies, in humans and certain subtypes  
11 related to more serious disease in humans.

12 To some extent, we can relate that to also  
13 previous typing methods. Phage typing, we can do to  
14 some extent. PFGE becomes more difficult which I'll  
15 kind of relate to later.

16 Okay. So we've heard about single  
17 nucleotide polymorphism, SNPs. One way to try and  
18 explain the level of granularity we have now, and  
19 we've heard from Martin as well, there are different  
20 ways of using the information where the genome  
21 differs. This is really the core genomic  
22 information. There are different ways that we can

1 use that, but it's about that precision.

2           And one way that's used by public health  
3 England at the moment and others, is to use a SNP  
4 address which is shown at the bottom here. We're  
5 starting from the right-hand side. You really assign  
6 the particular isolate into groups based upon the  
7 level of relatedness on the differences of number of  
8 SNPs that they have.

9           Without a simplistic level and Martin has  
10 outlined, you have to really know the potential for  
11 variation over time with your genome in terms of  
12 error rates, but at the moment, simplistically if  
13 you're within 5 SNPs, you can consider that you have  
14 very related bacteria that may be associated in an  
15 outbreak.

16           And one way of doing that precision, as  
17 Martin mentioned, is sort of a zip code way of  
18 thinking about it. If you here look at the diversity  
19 here say of households, the distribution of  
20 households in the USA, this could be the distribution  
21 of bacterial variation we have within our *E. coli*  
22 O157, and using the SNP address, we can focus that

1 down to a particular set of states and then smaller  
2 regions on that, get it down to a town level, get  
3 down to a block level, and eventually to a final kind  
4 of household address.

5           So level of precision is right there in  
6 terms of what we can do with that sequence data, and  
7 it is way more than we were able to do with these  
8 previous techniques, okay, and it's absolutely  
9 critical to take that point home.

10           So how do we apply that with something like  
11 *E. coli* O157. What we have here are 2,527 sequenced  
12 genomes, arranged in a ring, but it's hard to fit it  
13 in, and otherwise, I kind of like that one. The  
14 colors represent traditional lineages that have been  
15 designated for O157, 1, 1, 2 and 2.

16           And the branching here is shown at the  
17 level of 25 SNP relatedness. It stops there. It  
18 doesn't break it down any further than that. As you  
19 can see, the relatedness of 250 SNPs, you've got 79  
20 clusters; at 100, 240; at 10, 1,423; and so you can  
21 identify an isolate at a specific, precise level  
22 within that tree.

1           So how do we actually apply that? This is  
2 an example here from an outbreak in England and  
3 Wales, and the key point here is that in the end,  
4 there were 49 cases that could be confidently linked  
5 to the packed leafy salads when initially there was  
6 nothing -- I mean the epidemiology is absolutely  
7 critical, but initially you weren't necessarily  
8 getting all the cases from the same product source,  
9 okay. So you are able to then associate these  
10 clusters based on a very related SNP address.

11           One thing also to say about this is that a  
12 year later, there are a number of further cases, this  
13 time associated with lamb products and those are the  
14 ones shown in pink higher up, and as Martin alluded  
15 to, what we can do there is say these are kind of  
16 related common ancestors to this outbreak. It does  
17 not mean therefore that the contamination may have  
18 occurred from lamb or from sheep as a source of  
19 contamination and of that produce.

20           Furthermore, what we're trying to do in the  
21 UK is sample across the country. This is from beef  
22 farms and the locations of the farms that we're

1 sampling and then we can get our prevalence studies  
2 from this, but it also allows then to begin to  
3 associate regions with particular subtypes of O157.  
4 So we're actually getting a sort of locality to  
5 particular types.

6           Where that's useful is that we can actually  
7 plot on our scheme of all our organisms, those  
8 causing human infection, those that are related with  
9 cattle. We can understand why they can converge,  
10 where we're really getting isolates that are coming  
11 from our local cattle into human populations, but we  
12 can also spot imported infections as well, where we  
13 don't have those particular organisms in the local  
14 cattle population or the local ruminant population.

15           So it becomes very useful for understanding  
16 imported threats versus those we're generating  
17 locally.

18           But also we want to be able to predict from  
19 this which isolates are more of a threat to human  
20 health.

21           As was mentioned, we ideally want to use as  
22 much of the information as possible, so not just the

1 core level. We want to make sure that we're using  
2 core plus accessory genome both, for that prediction  
3 analysis.

4           So this is an example of looking at  
5 antibody resistance characteristics, now general *E.*  
6 *coli*. This is a project we're involved with around  
7 Lusaka in Zambia, and it's looking at *E. coli*  
8 isolates from cattle, small holders, and from a human  
9 population. And you can just about make it out, but  
10 the cattle isolates are in red, the human isolates in  
11 blue in this tree of *E. coli*. And the bars all  
12 around the outside are the number of antibody  
13 resistance genes. Okay. Shown from  $n$  equals 1 --  
14 from 0 to 17.

15           And just from this type of analysis, we can  
16 get to see when we have our blocks of human isolates,  
17 we have many more significantly higher levels of  
18 antibody resistance than we have in these *E. coli*  
19 that are coming from cattle. We have again data on  
20 what these animals have or haven't been treated with.  
21 So all this comes from the WGS data in terms of  
22 actually very easy to use databases such as ResFinder



1 that you can plug either your de novo assemblies or  
2 your actual basic reads into.

3           Taking sort of a step forward, to see what  
4 we'd like to do is have more complete DNA  
5 information. I appreciate this is a horrible slide.  
6 Each of these lines represents a completely assembled  
7 O157 genome. Okay. So we have 14 *E. coli* O157  
8 sequences, and this is really what we do with Jim  
9 Bono in USDA Nebraska, and this is all based on --  
10 sequencing. So you're getting long read sequencing  
11 where you can actually fully assemble the chromosome  
12 of the organism.

13           And the key point I want to show here is  
14 these blocks of color are prophages or X  
15 bacteriophage regions that at times have moved into  
16 the O157 genome. You can see some of them are very  
17 similar in very similar places, but a few, especially  
18 the green and red ones, sorry if you're color blind,  
19 but the ones that lay on the right side, very much  
20 more. We worry about that because these particular  
21 prophages carry the toxins that cause serious damage  
22 to human health.

1           Now, these are very difficult to actually  
2 position and identify using short read sequences  
3 because there are very strong similarities between  
4 them.

5           So it really helps to have this type of  
6 long read sequencing information to understand the  
7 isolates that we're dealing with. I mean examples  
8 also of outbreaks that we've had where the number of  
9 SNP differences in the core genome might be say 5 or  
10 6, but actually the organism over a year, and this is  
11 one example I think that was in a restaurant that had  
12 two outbreaks, separated by a year, where the  
13 organism had acquired over the difference in time, a  
14 plasmid, and it rearranged prophages in its genome.  
15 So while at the core level, it was very, very  
16 similar, it actually had something like an additional  
17 250,000 base pairs of information. So you have to  
18 bear that in mind when you're just using core SNP,  
19 SNP-based information.

20           At another level, we can look at where  
21 insertion sequence elements are within the genome,  
22 and again Martin gave some really nice examples of

1 how is the threat the same, dependent on mutations  
2 that have occurred. Well, here we have an insertion  
3 sequence element that has jumped into the Shiga toxin  
4 2a gene. This inactivates this and it means this  
5 strain is less of a threat to human health. So this  
6 is impossible to spot with short read sequencing  
7 really, and you have to have the long read sequencing  
8 to identify it.

9           Okay. So when I finish off the talk, we're  
10 looking at the potential use of machine learning and  
11 prediction of both pathogenesis, zoonotic potential  
12 and host attribution.

13           So machine learning has been originally  
14 described as the capacity of the computer to learn  
15 from experience, i.e., to modify its processing based  
16 on newly acquired information. And the first  
17 algorithm, first work, we were doing it back in about  
18 in the 1930s, pre-computers.

19           We should be aware that a lot of what you  
20 do now, your activity is monitored, right. You are  
21 watched. When you type in your searches into a web  
22 address, all that's fed back to Google, etc., and

1 they're crunching that and they're using machine  
2 learning approaches to really nail exactly you need  
3 that and you -- etc., in your world.

4           It's used to exam changes in patterns for  
5 bank fraud. It's used, pattern recognitions are  
6 used, for examining images for identifying tumors,  
7 etc. AlphaGo was recently in the news in terms of  
8 DeepMind computer teaches itself to become world's  
9 best Go player as well. So watch out.

10           So one of the ways we've used this recently  
11 is a supervised machine learning method which I'm  
12 going to very quickly take you through with these  
13 really stolen from the web tutorial. So very  
14 simplistically, two sets of data. Height and weight  
15 and we have data for men and women in this case,  
16 okay. So very, very binary system, and we have this  
17 training data. So in this methodology, you need  
18 training data. So this is our training data, and  
19 then you can assign a rule that will best separate  
20 that training data based on the information you have.  
21 Okay. And it's all about getting the rule in the  
22 right place in terms of separating, giving you the

1 optimum distance in terms of separating your data.

2           You then come along with -- sorry. This is  
3 very, very fickle.

4           You then come along with your test data.  
5 So it's new data that you haven't seen before and you  
6 plot that and then you apply your rule, and the idea  
7 here is you then can assign whether you are dealing  
8 based just on height or size in this case, and  
9 weight, whether you're dealing with a man or a woman.

10           Okay. It's obviously flawed, and we get  
11 very skinny short guys clearly, but this is only two  
12 bits of information. We now start to think about  
13 applying that to hundreds or thousands of genes and  
14 the prevalence of those genes or the predicted  
15 proteins of those genes across isolates. Then we get  
16 into the proper world of multidimensional support  
17 back to machine analysis which is way out of my  
18 league, but we use it.

19           And the idea here is you're still able to  
20 draw a separating line in that data. Okay. So  
21 you're still able to assign A to B, even though  
22 you've got very complex patterns of data, and that's

1 the beauty of this supervised machine learning  
2 approach.

3           So recently we've applied that to looking  
4 at O157 strains across all the lineages. We have a  
5 particular problem in the UK on the left there with  
6 our lineage 1 isolates, and that's really horrible  
7 and really hard to see, but basically the human and  
8 the bovine isolates are really mixed up in the  
9 lineages. So it really becomes difficult to say, if  
10 you have an isolate that fits in there, is it more or  
11 less likely to be a threat to human health.

12           So we applied the support vector machine  
13 process to this, train on the subset, test on the  
14 remainder, repeat the process and obtain statistics  
15 or prediction scores. So the left axis here, you're  
16 looking at a probability. This is actually what you  
17 want to come out with in terms of prediction  
18 capacity. Probability based on isolates, this is all  
19 Illumina sequencing, based on an isolate whole genome  
20 sequence of whether it contains more bovine or more  
21 human information in terms of its comparison to the  
22 other isolates in those groups. Okay.

1           And you can see here in the green box on  
2 the right-hand side, the bovine isolates are on the  
3 right, the majority score well for being bovine, but  
4 there are subsets that score well into the human  
5 zone. And the proposition is that those are the  
6 isolates that are more of a concern. They are more  
7 of a threat to human health. Okay.

8           How can we test this? I mean you get what  
9 you look for, right. This is really about the  
10 concept of it, not necessarily whether it's right or  
11 wrong at the moment, but what we can do, if we have  
12 large amounts of data, and actually have the metadata  
13 associated with that sequence data.

14           So, for example, I mean one way we've  
15 tested this is to take two outbreaks, one was a milk  
16 outbreak on the left there, the sort of light blue,  
17 and then the pink is a food outbreak. And what  
18 you're trying to do there is take the sequences of  
19 those isolates and to score them. Now, those that  
20 are coming from food or milk or animal, where are  
21 they fitting on our 0 to 1 probability score? And  
22 you can see that all of those isolates that come from

1 hamburger, cattle, or milk are all scoring high for  
2 human, even though they have that animal source. So,  
3 indicating again, this is just a small sample size,  
4 the possibility of predicting subsets that are more  
5 of a threat to human health.

6           The way we're trying to take this further,  
7 and again complex slide, but this is again using now  
8 all *E. coli* that we can get our hands on. We don't  
9 have enough yet. We're nowhere near it. We're just  
10 dealing with a few hundred. This is standard *E. coli*  
11 that come from cattle and *E. coli* that come from  
12 humans. We have the whole genome sequences of those,  
13 and then used that same information to try and train  
14 the machine to learn -- train the computer to come  
15 along. If I come along with a new *E. coli*, I say  
16 where does this come from? Does it come from a cow?  
17 Does it come from a human? And it will assign, it  
18 will assign red for cattle there, and blue for human,  
19 in terms of the scores. And we can get that about  
20 90% right at the moment based on the small sample set  
21 we have.

22           Again, looking on the right-hand side, the



1 black dots are where O157 fits in that scheme. So it  
2 is very interesting, that something we know that has  
3 enough potential, that has the capacity with its gene  
4 content to move from cattle to humans, actually  
5 generally goes above the intermediate line. It  
6 contains genetic content that as far as machine  
7 learning goes, can be ascribed to both cattle and  
8 human, but puts it up towards the human category.

9           What's interesting is to understand some of  
10 the other *E. coli* in the group and whether they  
11 represent a zoonotic threat. On the left there, we  
12 can plot the scores for the different hosts as a bar  
13 graph. So the bovine score is in red and the human  
14 score is in blue.

15           Obviously, this is the very beginning of  
16 things. We need, you know, a phenomenon amount of  
17 data to make this more accurate, but it is a start  
18 for us.

19           We can do this as well -- the last few  
20 slides are on *Salmonella enterica* serovar  
21 Typhimurium. We know for *enterica* we have different  
22 serovars that are fairly host restricted or very host

1 restricted in different cases, but the Typhimurium,  
2 there are known subtypes that are much more  
3 associated with human, but it's generally considered  
4 an organism is able to traffic well between hosts and  
5 if we identify Typhimurium, we generally think it has  
6 the potential to cause human disease.

7           So can we apply the same approach to have a  
8 look at Typhimurium? The moment we got hundreds of  
9 isolates across avian, bovine, human, and swine,  
10 obviously there will be other reservoirs where we can  
11 train the support vector machine on this and actually  
12 then take samples and tests where the host  
13 attribution may lie. So what is the likely host for  
14 the organism?

15           And you get these very funky kind of  
16 looking graphs here. So top left is the score for  
17 looking at avian isolates. What's interesting here  
18 is that you can the majority of avian isolates score  
19 very well for being avian and actually don't have a  
20 lot of other increased color for other hosts.

21           We're very interested here in which hosts  
22 are maybe housing Typhimurium isolates to have more

1 capacity to move to other animals or colonize us,  
2 okay.

3           Again, I wouldn't want to say this is  
4 reality, but this is what's possible in terms of  
5 thinking about the way this technology can be  
6 applied. If you think about taking a water sample  
7 and sequence *E. coli* or something or other that comes  
8 from it, and actually say where do these organisms  
9 come from, we can attribute the likely contamination  
10 in terms of whether that's human or whether that's a  
11 local farm up the road. This allows us this degree  
12 of prediction, obviously combined with what we  
13 already know about the phylogenomics of many of these  
14 organisms.

15           So at the moment, we can use the machine  
16 learning to predict the host source and zoonotic  
17 potential. We have to prove it. It's still  
18 conceptual but it's interesting at the moment that  
19 the majority of these Typhimurium isolates show  
20 pretty host restricted signals. So it is only again  
21 a small subset that really make it into the human  
22 domain in terms of the way the genomic information is

1 being analyzed. And so therefore potentially only  
2 specific subsets are a risk to human health.

3           With larger datasets this type of approach  
4 can inform us of where human infections originate  
5 from to inform risk assessments.

6           What we now need to do is go back and ask  
7 exactly which genes, which combinations are being  
8 used to make these decisions so we can understand the  
9 biology behind this type of prediction.

10           Okay. So, in summary, as you heard from  
11 Martin as well, whole genome sequencing is very  
12 powerful. It offers -- it's that where we're all  
13 kind of using at the moment, transformative in terms  
14 of its capacity to track infections and trace sources  
15 of bacteria.

16           Considerable information can be extracted  
17 from whole genome sequencing including basic taxonomy  
18 and identification of virulence genes and AMR genes,  
19 and this helps determine the threat represented by  
20 isolates.

21           But again as the technology improves, the  
22 longer read sequencing will give us even more

1 potential to be specific about this, but at the  
2 moment, the costs are very high still for that. So  
3 that's really where the short read sequencing can be  
4 bolstered by long read methods to more accurately  
5 assemble the genome and predict these phenotypes.

6           The issues are really around getting hold  
7 of data related to the sequences. I completely  
8 understand that many industries, it really should be  
9 released information. How do we use it? Do we use  
10 it ourselves? The more information we can have, the  
11 more sequencing information, the more related to  
12 human disease, the more related to which animal it  
13 comes from, time of isolation, place of isolation,  
14 can really transform our capacity to be precise about  
15 this understanding.

16           And there are amazing bacterial collections  
17 at the moment that are now being sequenced and as we  
18 go forward, that information will be available.  
19 Again, we're only scratching the surface of diversity  
20 that's out there.

21           And then the obvious point, the more that  
22 we have in those databases of individual isolate

1 sequences, the better the metagenomic approaches can  
2 be and that's certainly to sort of reiterate what  
3 Martin said in terms of the capacity of the  
4 metagenomics as the costs come down, and it's clearly  
5 from a diagnostic level, the monitoring level is  
6 going to be a clear way forward in terms of picking  
7 out the threats that exist within deeper sequencing  
8 potentially of air samples, water samples, etc.,  
9 within facilities.

10           Okay. Just quickly to thank the fact that  
11 the machine learning is really a serious Ph.D.  
12 student's work, and this is just written up as  
13 Nadejda Lupolova, the University of Edinburgh, and  
14 I've been lent some slides from Tim Dallman at Public  
15 Health England in terms of the zip code and SNP  
16 mapping.

17           Thank you.

18           DR. STRACHAN: Okay. All right. Thanks  
19 very much for inviting me to speak today. My name is  
20 Norval Strachan from the University of Aberdeen. My  
21 training is as a physicist, but I worked the last 20  
22 years on risk assessment and molecular epidemiology

1 of gastrointestinal pathogens, which I guess is the  
2 reason I've been asked here today. I'm also the  
3 Chief Scientific Advisor for Food Standards Scotland.  
4 The views that I'm giving today will be my own.

5           Okay. So in terms of talk, what I want to  
6 do is I want to outline what risk assessment is and  
7 the way I understand it and the way it's drafted out  
8 roughly by Codex. I want to provide some examples of  
9 whole genome sequencing which are applied to the  
10 different steps of risk assessment. And then the  
11 final part, which I want to speak about and say a  
12 little bit about, source attribution, how this  
13 relates to risk assessment and basically how we can  
14 use whole genome sequencing to help us with that.

15           So here's a diagram of what I want to say,  
16 a little bit about risk assessment. First of all,  
17 this isn't food. This is a hazard in front of us.  
18 We've got a rock on the top of a cliff. Okay. So  
19 because that rock can cause some damage if it falls  
20 from that height, the rock's a hazard. A hazard is  
21 something that causes a negative impact particularly  
22 in our case for considering health.

1           But also it's well for one to think about  
2 what risk is. Okay. So risk, there's really two  
3 dimensions of risk.

4           One is a probability. It's a probability  
5 that that rock's going to fall down and cause us  
6 harm. So that's one dimension of risk.

7           The other dimension of risk is the severity  
8 of risk. So that's a pretty large rock, and if it  
9 fell on my head, I don't think I would get up from  
10 that. So it would probably kill me. If it's a much  
11 smaller rock, maybe I might survive that. So we need  
12 to think about severity as well.

13           And so what risk assessment is, it's  
14 looking at the hazards in a technical, scientific  
15 way, looking at knowledge associated with that, to  
16 determine what the risk is.

17           Okay. So we'll move onto something which  
18 is a big more in our topic area. So there's a  
19 picture there of a beef burger, and it's a rare beef  
20 burger. Okay. So you can think about what the  
21 hazard is that might be associated with that. So,  
22 for example, *E. coli* O157 could be a hazard



1 associated with that rare beef burger. The vehicle  
2 is, of course, the beef burger itself. So a hazard  
3 is a biological agent in the food for a potential to  
4 cause an adverse health effect. Right.

5           So in terms of thinking about risk  
6 associated with that, we try and think if we eat that  
7 beef burger, what's the probability of us falling ill  
8 from consuming the beef burger which may or may not  
9 have that hazard in that, and then the second aspect  
10 is severity as well. So if we fall ill, how ill will  
11 we fall? If it's O157, we could have hemorrhagic  
12 uremic syndrome, or perhaps we could have mild  
13 diarrhea perhaps. So that's the two aspects.

14           Then this risk assessment itself is using  
15 the scientific knowledge that we've got available to  
16 us and putting it in a form to make up an opinion on  
17 terms of the risk associated with eating that rare  
18 beef burger meal.

19           Okay. In terms of Codex, there is four  
20 steps in risk assessment, but before I actually go  
21 into those four steps, probably the most important  
22 thing actually is a statement of the purpose of the

1 risk assessment, and this is usually defined by the  
2 persons who want to manage the risk. So basically  
3 it's looking at -- well, what we're interested in  
4 perhaps is what's the risk of consuming the beef  
5 burger, but we'd be interested across a population,  
6 you know, across the U.S. How many of these beef  
7 burgers are eaten? How many fail ill? What's the  
8 severity of illness? And so we're hoping that the  
9 risk assessment will answer these questions so that  
10 the risk managers can then look at this and then  
11 decide, yeah, okay, there's not really much there. I  
12 don't need to do anything about it or else maybe we  
13 do need to do something about it, and then put some  
14 mitigation strategies in place to try and reduce the  
15 risk.

16           So the risk assessment itself, the four  
17 steps associated with that. There's hazard  
18 identification which has already been mentioned. So  
19 identifying the hazard, but in this case, it's the *E.*  
20 *coli* 0157 in the burger.

21           There's then also the exposure assessment.  
22 What exposure assessment is doing is it's looking at

1 whether there is a organism in the beef burger that  
2 we actually ingest, and that we might fall ill from.

3           The third part is the hazard  
4 characterization which looks at -- there's two  
5 aspects in hazard characterization. One is how many  
6 bugs does it take for us to fall ill? So we need to  
7 think about things like dose response. And the other  
8 aspect of hazard characterization is well as the  
9 severity associated with those organisms, so the  
10 severity of disease that we're likely to get.

11           In all these first three steps, we can use  
12 whole genome sequencing to help us among other  
13 things, of course, but in the fourth step, risk  
14 characterization, is putting this all together, so  
15 looking at the risk across the whole population, for  
16 example, that probably doesn't involve whole genome  
17 sequencing per se because we will use the database  
18 from that.

19           Okay. So what I want to do now is I just  
20 want to go through one or two examples for the first  
21 three steps of risk assessment.

22           So hazard identification itself is

1 generally qualitative process to identify the  
2 hazards. So basically we know some particular types  
3 of food, those particular types of agents that we  
4 know of previously have causes disease. So it would  
5 be the hazard associated with those particular types  
6 of food.

7           However, as has already been mentioned  
8 today is that in terms of risk assessment, how we've  
9 routinely done it before, we ignore heterogeneity  
10 between organisms. We think that all the *E. coli*  
11 O157 are the same, all the *Campylobacter* are the  
12 same, all the *Listeria monocytogenes* are the same,  
13 and we tend to treat them all in the same way, and  
14 that's the traditional way that it's been done.

15           But the great opportunity with the whole  
16 genome sequencing is that we are now able to  
17 characterize these different organisms. So we're  
18 able to find out the variance between them. It makes  
19 it more complex but using this knowledge hopefully  
20 can help us.

21           So we're going to look at two examples.  
22 The first hazard identification example is one here

1 which is published by a Dutch Group in the  
2 Netherlands, and it's an *E. coli* O157 example and  
3 basically what this did, they studied 38 strains of  
4 *E. coli* O157 and how they attached to human  
5 epithelial cells.

6           So one of the points we're seeing in this  
7 analysis was it caused disease. The organisms need  
8 to attach to human epithelial cells. They had  
9 basically a model system. So basically they grew the  
10 organisms up, they put them through gastric fluids to  
11 simulate the stomach, and then through intestinal  
12 fluid to simulate transfer through the gut, and then  
13 they looked at the attachment of these cells, the  
14 microorganisms into the epithelial cell line.

15           And these are some of the results that they  
16 got from the attachment assay. So the graph at the  
17 top right there, it shows along the bottom is  
18 fractionate adhesion to the Caco cells, the human  
19 epithelial cells. On the vertical axis is the  
20 frequency. The top graph on the right, this is  
21 actually the ones from humans although it's actually  
22 misnamed in the paper. We see the part highlighted,

1 the red circle, it's where there was very good  
2 attachment to the epithelial cells.

3           The example at the bottom for the animal  
4 examples that they had, that had a poor adhesion to  
5 the Caco epithelial cells.

6           Okay. So we've got this phenotypic data  
7 and all throughout the two talks previously, both  
8 David and Martin said the importance of the wet  
9 biology being done, but also how we can link this to  
10 the whole genome sequencing data. So it sequenced  
11 all the data and for *E. coli* O157 which this is for,  
12 there's 5.5 Mb. What you need to do is reduce that  
13 down to something more manageable. So they reduced  
14 it down to SNPs for the core genome and they got  
15 about 28,000 SNPs. So that's made it simpler but  
16 obviously there's still 28,000 SNPs there. So it's  
17 not quite simply enough to deal with.

18           So what they did next was quite smart which  
19 is doing like a genome wide association study and the  
20 graph that we've got here, along the bottom, we've  
21 got -- basically it's all the SNPs listed along the  
22 bottom and then on the vertical axis, if the bar's

1 higher, it basically means that there's a better  
2 attachment of the organism associated with that  
3 particular SNP. And there's a little red line along  
4 the middle there, I'm not sure what the statistical  
5 significance is.

6           So all those black bars that go above that  
7 red horizontal line are the ones that's of interest.  
8 Those are the SNPs of interest.

9           So they initially found 17 SNPs, but then  
10 after correction for a sample structure, only 1 SNP  
11 stood out as a potential biomarker. So we started  
12 5.5 million basis and gone to 8,000 SNPs and now  
13 they're down to 1 SNP, and what they're seeing is  
14 this particular one, this particular SNP is a marker  
15 for strong attachment in the *E. coli* 0157 cells. The  
16 SNP itself was found actually to be an enterogenic  
17 area.

18           So this shows you potentially what whole  
19 genome sequencing can do. You can identify SNPs  
20 which can give you an indication on whether the *E.*  
21 *coli* can attach to human cells well or not. Further  
22 work needs to be done because this was only done with

1 38 strains. We need to be doing further strains to  
2 make this validated.

3           The second example I want to present on  
4 hazard identification is I just actually want to  
5 mention the one that David gave already because I  
6 think it's quite important. One of the things he  
7 actually said with the machine learning work that he  
8 spoke about. And what he said was we demonstrate  
9 only a small set of bovine strains is likely to cause  
10 human disease even within previously defined  
11 pathogenic lineages. And if I remember rightly,  
12 within the paper it's about 1 in 9 of the strains  
13 that he tested were likely to be pathogenic. So this  
14 is telling us, you know, although O157 potentially  
15 are not pathogenic to humans, maybe 1 in 9 are, and  
16 that's important in terms of taking into risk  
17 assessment calculations.

18           So thanks for that, David. You did all the  
19 hard work to explain that example for me.

20           I want to go now and say a little bit more  
21 about exposure assessment and give an example  
22 associated with exposure assessment. And we're going



1 to go back to *Listeria* for this example. So what you  
2 can think about, looking at production of dairy  
3 products. Imagine that *Listeria* can get into dairy  
4 products and contaminate the milk within the dairy  
5 itself and then we've got a number of steps. So  
6 there's a heat processing step in terms of  
7 pasteurization. Then it gets packaged maybe into  
8 various different products, sent to the supermarket,  
9 and then it's eaten and consumed by humans.

10 So what we can look at, is consider the  
11 different strains of *Listeria* that you can find  
12 perhaps within the original population of *Listeria*.  
13 So think about the wild type first of all, and that's  
14 that blue line there.

15 So the first step in transport, it might  
16 grow within that step. Then the pasteurization step,  
17 it's been heat treated. There's a big die-off, a big  
18 kill associated with that. Then during storage, it  
19 may grow a little bit again perhaps and then stomach  
20 passage, because we've got low pH there, then it may  
21 die-off. So this would be a typical *Listeria* strain  
22 which goes through this process, and throughout the

1 process there's considerable die-off. However,  
2 that's the wild type.

3           But we can think about is the resistant  
4 type, and what you'll see in the graph here, the  
5 resistant type, it starts at a much lower  
6 concentration on the left side, shown in red, maybe  
7 grows a little bit during transport. In terms of  
8 pasteurization, the die-off may not be so much  
9 because it's resistant to heat. In storage, it could  
10 increase a bit perhaps and stomach passage may also  
11 be resistant to pH. So what you end up with is  
12 actually quite a lot more of the resistant type being  
13 present at the end of the process, compared to the  
14 wild type of *Listeria*.

15           So looking at this in a little bit more  
16 detail, an example that was given in the literature,  
17 what we can do now is just look at the graph here, by  
18 Metselaar, and this is just the pH example of that.  
19 So they've got a culture of *Listeria* which have grown  
20 up to a higher level. We then put it into liquid  
21 medium which is pH 3.5, so acid for a considerable  
22 period of time, up to 200 minutes.

1           So start off with, there's a big die-off.  
2 Okay. That big die-off is from those wild type  
3 strains which are sensitive to pH, but then there's  
4 this long tail, and this long tail, what you find in  
5 the long tail is the resistant strains and those are  
6 the resistant strains that are surviving, and  
7 actually recultured organisms from the resistant,  
8 from the tail, and the 23% then were stable  
9 resistant. So this is all very interesting, but what  
10 we need to think about now is how whole genome  
11 sequencing can help us with that.

12           So what they did, they sequenced the  
13 resistant strains and they sequenced the wild type,  
14 the ones that were non-resistant. And so we have  
15 this graph here. So at the very top, we've got the  
16 wild type, and below we've got all the resistant  
17 strains, and this here is a sequence which is  
18 upstream from *rspU* gene which is associated with  
19 stress tolerance within other organisms.

20           And what they found in the ones that were  
21 resistant by a number of different mutations in this  
22 upstream region, and the hypothesis is that these

1 mutations are causing the strains to be able to  
2 obtain resistance.

3           So what would be interesting is to look at  
4 this group of changes, to look across the population  
5 to see if this is a common thing in terms of the  
6 human population and also if you find those strains  
7 within a particular product, to look to see if they  
8 are resistant, if they have this mutation as a  
9 biomarker to look for that.

10           Okay. What I want to say very briefly is a  
11 little bit of an example on exposure assessment and  
12 metagenomics, and again Martin's mentioned very  
13 nicely a little bit about metagenomics. The thing  
14 about most microbiology studies that have been done  
15 previously is that they're based on culture. We can  
16 only culture some of the organisms, not them all, and  
17 this leads to the biases and, of course, in the real  
18 world as well, there are always a community of  
19 organisms stand to be present. So metagenomics  
20 allows culture independent analysis of  
21 microbiological populations.

22           So there are strengths and weaknesses with

1 metagenomics. Martin outlined them, but I don't want  
2 to go into that here. I just want to go into the one  
3 particular example, and this example is an example  
4 which is going to be from cheese. So what we have is  
5 our sample, and it could be human or whatever, but  
6 we're going to speak about cheese.

7           We're going to extract the DNA or rRNA and  
8 in this particular example I'm going to think about  
9 is extracting the 16S, extracting rRNA and from that  
10 doing 16S rRNA sequencing. And from that what you  
11 can do is you can, on the bottom left-hand side  
12 there, that small graph, which physically shows --  
13 identifies the species and relative sequences of the  
14 microorganisms that are within your sample.

15           Okay. So we do that, and so this is done  
16 for an example in terms of Italian cheese, and what's  
17 shown here on the left-hand side, we've got a list of  
18 all the organisms or families of organisms that are  
19 found. And these are repeated along the diagonal on  
20 the right-hand side.

21           And one of the interesting things with this  
22 is what you can do is you can work out which

1 organisms occur together, which organisms do not  
2 occur together. So if they occur together, they're a  
3 red [sic] dot. If they don't occur together, they're  
4 a blue [sic] dot, and if it's somewhere in between,  
5 the colors in between yellow, blue, light blue, white  
6 and so on.

7           And what I've highlighted here is for  
8 *Listeria* here, and here we're looking at *Listeria*  
9 species along the horizontal and what you can see,  
10 lots of blue dots. So it's organisms occur together  
11 with *Listeria*, and other vertical down is  
12 *Lactobacillus brevis*. So it's very common in this  
13 type of cheese to have *Listeria* with *Lactobacillus*  
14 *brevis*. They tend occur together in this type of  
15 cheese.

16           Ideally, what we're wanting is actually red  
17 dots there where there's exclusion because if there  
18 are red dots, what that means is that maybe the  
19 organism *Listeria* may be competing with is producing  
20 some compounds that's inhibiting it. So that's what  
21 we would be looking for.

22           If you're able to find those sorts of

1 organisms, what you might be able to do in your  
2 startup culture, for example, is include those  
3 organisms so that they are then in your cheese. So  
4 you then design a product which is less likely to  
5 inhibit any *Listeria* that may be present.

6           So this is an example using metagenomics  
7 where there is a potential for it to help in the  
8 design of food products.

9           Okay. I'll say a little bit about hazard  
10 characterization now. So as I mentioned previously,  
11 hazard characterization, there are two aspects to it.  
12 There's the dose response and also the severity in  
13 terms of the human response.

14           Okay. So those responses both have to do  
15 with ingestion of the pathogen and colonization  
16 associated with that. We have a graph on the right-  
17 hand side there. What I want to show here is just  
18 the large variation in doses. This happens to be for  
19 O157 but you get fairly similar things for other  
20 organisms as well. We've got dose along the bottom  
21 and we've got probability of illness on the vertical  
22 axis.

1           There are a number of circular dots in  
2 that, and they represent outbreaks. So, for example,  
3 the dot at the top right-hand side there, the dose  
4 that the people had in this outbreak of between 10 to  
5 4 and 10 to 5 organisms, and 80% of the people fell  
6 ill.

7           You can see, there's a lot of variation.  
8 This is from a number of outbreaks from different  
9 parts all associated with *E. coli* O157. The best  
10 dose response model fits through the 0.50 one in the  
11 middle there. You can see a huge range. There's a  
12 huge variation that's involved due to differences in  
13 the *E. coli* itself, but also there will be  
14 differences in the human cases as well because we  
15 know humans are of different susceptibilities if  
16 you're young or immunocompromised and so on.

17           The example I'm going to go on and speak  
18 about now in particular to whole genome sequencing is  
19 about the severity of disease. I'm sticking with  
20 O157 for this and for O157 or for Shiga toxin  
21 producing *E. coli*.

22           What we do know is that there are two main



1 types of Shiga toxin, Shiga toxin 1 and Shiga toxin  
2 2, and these have subtypes 1a to 1f and 2a to 2g. So  
3 there's all these different subtypes involved.

4           And David mentioned it briefly, that  
5 there's a difference in potency in these and from  
6 work done previously in terms of a mouse bioassay,  
7 there was stx2d, stx2a were the most pathogenic, and  
8 then stx1 toxins were less pathogenic to the mouse.

9           What's great about whole genome sequencing  
10 now is that we can do the Shiga toxin typing, just  
11 directly from the next-generation sequencing reads,  
12 and Phil Ashton and colleagues from Public Health  
13 England published this a couple of years ago now. So  
14 we're able to get this information just directly from  
15 the reads from the genome, not from the sample genome  
16 itself because the sample genome, it has problems  
17 assembling the toxin genes because they're  
18 paralogous.

19           Okay. So leading on from that, I'll go  
20 onto this example or information on being able to  
21 sequence these toxin genes, is able to help us  
22 understand about disease in humans. So this

1 phylogeny here is for O157, was from 105 Scottish  
2 clinical samples which was carried out by Anne Holmes  
3 and colleagues at the Scottish *E. coli* reference  
4 laboratory. And what they were able to do was they  
5 were able to look at the Shiga toxin types, and they  
6 were also able to look at the severe disease.

7           Now, let me just point, run quickly across,  
8 because I can't point. So this last column here, the  
9 red bars we're interested in is the HUS cases ,and  
10 here is the difference in colors on the left-hand  
11 side, the Shiga toxin types.

12           Okay. And what they found was so basically  
13 look at the right-hand column, and it's the bottom  
14 half which we have all these HUS cases, all the  
15 really nasty cases, and what they found was that 8  
16 out of 10 of those involved the stx2a gene. So  
17 basically the stx2a is a good indicator of severe  
18 disease. So we can use whole genome sequencing to  
19 help us with that. This is for O157, and it may be  
20 helpful to use this for other organisms as well.

21           Okay. So that was about hazard  
22 characterization.

1           So what I want to spend the last few  
2 minutes on is speaking about the sources of human  
3 infection and source attribution. So I'm going to  
4 speak about *Campylobacter* and *Campylobacter* can be  
5 found in lots of different sources, chickens, sheep,  
6 wild birds, pigs, cattle, etc. And there's also  
7 cases of *Campylobacter* in Scotland, America, and many  
8 places across the world, but one of the questions is  
9 where are the cases coming from? Where are most  
10 cases coming from?

11           So risk assessment along the top, which we  
12 have been speaking in the first three-quarters of the  
13 talk, which basically follows the organism say from  
14 the cattle through the food chain to the infected  
15 person and all the way along there like that.

16           But source attribution, using microbial  
17 subtyping is sort of a cheat, but it's very powerful.  
18 What it does is it basically -- it just looks at the  
19 types of organisms in the animal sources and looks at  
20 the type of organisms in humans, and it does a  
21 comparison, and it compared to see which is more  
22 similar to each other. It's a bit like what David

1 was speaking about, machine learning approach, but  
2 the approaches used here are based on population  
3 genetics for this. That's what this is based on.

4           So here I've basically shown at the bottom  
5 left, chicken, yellow and red types, in human, yellow  
6 and red types. If you look at the top for cows, you  
7 get yellow and red types with black types as well.  
8 So there's some overlapping crossover, but there are  
9 some differences as well, and the idea is to link the  
10 two using these populations genetic methods.

11           The method I'm going to speak about is MLST  
12 and Martin again really nicely explained MLST. So  
13 you isolate the DNA, you sequence it, and I'm going  
14 to speak about seven locus MLST to start with, and  
15 for the seven loci, you get the numbers and then from  
16 that, you combine them together to get a particular  
17 sequence type, sequence type 257 there at the bottom  
18 for this particular sequence type of *Campylobacter*.

19           So an example to quickly look at is in  
20 Manawatu region in North Island of New Zealand, and  
21 had really big problem with *Campylobacter*. They did  
22 the source attribution. So what they were able to do

1 with the source attribution was to predict the  
2 source, and it's colored there. So over a period of  
3 years, 2005 to 2007, poultry is in yellow, bovine is  
4 red, ovine is blue and environment is green. And as  
5 you can see, poultry was the main source according to  
6 the source attribution.

7           Okay. So they felt that they needed to do  
8 something about that. And so this is where risk  
9 management comes in. So if we go back to our rock  
10 example, yeah, maybe we need to do something about  
11 that rock. So we can think about our risk mitigation  
12 strategies. We put up a sign, a warning sign so  
13 people maybe don't go past that area or they take  
14 care going past, and that can be their risk  
15 mitigation.

16           But in New Zealand, in terms of the poultry  
17 interventions, what they did was a number of  
18 different things. They improved procedures for  
19 catching birds and cleaning crates. They improved  
20 their immersion chilling. They produced mandatory  
21 targets for *Campylobacter* on poultry after primary  
22 processing. So they put a whole set of interventions

1 in, and then they continued the source attribution to  
2 see what would happen.

3           And this graph here, this is what happened  
4 in effect. So I've now included the year 2008 and on  
5 into 2009, and what you can see is that the number of  
6 cases reduced, that was a smaller percentage of  
7 poultry cases. In fact, there was a 74% reduction in  
8 *Campylobacter* in poultry cases.

9           So this is a way of monitoring how  
10 successful or otherwise the risk management  
11 strategies were. It doesn't actually throw a whole  
12 raft of things out there to try and solve the  
13 problem. It doesn't actually tell you which one  
14 actually was most important.

15           In Scotland, we also follow *Campylobacter*  
16 using this type of methodology as well. Our colors  
17 aren't quite the same as the New Zealand colors, but  
18 the yellowy orange one is the chicken, and so chicken  
19 we find is the most important source of *Campylobacter*  
20 in Scotland. It varies in the model it used. It  
21 varied between 55 and 70% for that.

22           What I also wanted to say just very briefly

1 is a little bit of work that we've done on source  
2 attribution of *Listeria monocytogenes* in Europe, and  
3 this was an EFSA project that was led by Eva Moller  
4 Nielsen, and our group was involved in doing the  
5 source attribution work. It was a fairly small  
6 dataset for doing this sort of work, but it can  
7 potentially show what the potential of it is. So  
8 basically we had isolates of *Listeria monocytogenes*  
9 that we sequenced from fish, swine, ovine, bovine,  
10 and poultry sources, and then we compared that to  
11 what's in the human population and we attributed it  
12 to the human population.

13           So as I say, the database was fairly small.  
14 There was about 700 isolates in total that were here  
15 that were sequenced, but what you can see from that  
16 is basically we used three different models and  
17 that's the different colors there. That tends, from  
18 what we got in this dataset, it tends that bovine  
19 sources are maybe a bit more important than others.

20           Obviously, this was a start. I think this  
21 is first work that's been published in source  
22 attribution for *Listeria* using subtyping data. Here

1 we had 1748 genes from the core genome MLST for that.

2           So I think this is something as well that  
3 I'm sure in the U.S. you will be able to carry out  
4 these sorts of studies as well. Indeed, I think you  
5 already are starting to do so.

6           Okay. So just some take home messages.  
7 I've said about the steps in risk assessment, and  
8 I've spoken a little about some examples where whole  
9 genome sequencing can be used in risk assessment. I  
10 feel that we're only at the tip of the iceberg in  
11 this yet, and there's a lot more that can be done in  
12 this.

13           I think as well, as already been said,  
14 David mentioned about the epidemiology data being  
15 important to link with the next-generation sequencing  
16 data but also as well as that, the biology data, the  
17 phenotypic data is really, really important to  
18 combine that.

19           I think potentially we're getting lots and  
20 lots of whole genome sequencing data and trying to  
21 tie these datasets together is really, really  
22 important.



1           And the last bit I want to say about is  
2 source attribution. I think source attribution is  
3 really quite helpful for understanding the sources of  
4 disease and tracking that over time, and also I think  
5 in terms of when risk management strategies are put  
6 in either by companies at regional or national scale,  
7 that they can be evaluated to an extent using this  
8 type of methodology as well.

9           So that summarizes what I wanted to say.  
10 So thank you for your attention.

11           DR. DESSAI: Okay. While the speakers  
12 settle down and the projector gets turned off, I just  
13 want to state that setting the stage was the session  
14 where we were going to talk about the hazard, how the  
15 hazard is characterized, then different tools that  
16 are used to characterize this hazard in the context  
17 of WGS, and then we heard areas where we can do some  
18 predictions using newer approaches.

19           And then we went to the risk part of it  
20 which is very important, how to transition from  
21 defining the hazard, characterizing it further to  
22 turning it into some part of potential risk, and I

1 think that is a challenge here.

2           So last 2 days, we had the NARMS meeting,  
3 and we talked a whole lot about AMR, but I think we  
4 were a little shy of making that transition from what  
5 a hazard is to what the risk can be because it's a  
6 challenging area.

7           So what we're going to do right now is we  
8 have about half an hour of question and answer  
9 session. You have microphones which are right there.  
10 Those of you who want handheld microphones, let us  
11 know. We can provide those.

12           After half an hour, you'll get a break, and  
13 we will be back on time for the next session which is  
14 going to be partners talking about whole genome  
15 sequencing. The most important thing about this  
16 meeting is although we are hosting it, it is a  
17 meeting of all the partners involved in whole genome  
18 sequencing. Let me just make that pretty clear here.

19           So we would like you to be back on time,  
20 and when you're going out, please don't forget your  
21 badges as well as if you need any escorts, let us  
22 know.

1 All right, so we open the floor to question  
2 and answers. What we realized yesterday is that  
3 those who are online sometimes cannot hear the  
4 questions and the answers very well. So please speak  
5 into the microphone very clearly. Thank you.

6 And those who are asking questions, please  
7 state your name and affiliation clear as well. Thank  
8 you.

9 DR. EVANS: We can start off with a  
10 question from a question on the webinar. There was a  
11 question about whether there were a national or  
12 international databases for microbiome data that  
13 could be used by scientists to study risk in the ways  
14 that you were talking about.

15 DR. WIEDMANN: So the question was is there  
16 an international or national database of the  
17 microbiome that allows people to study risk.

18 DR. EVANS: And metagenomic as well.

19 DR. WIEDMANN: And metagenomic. Well, I'll  
20 take a first stab at that. And to the best of my  
21 knowledge, there's no database on the microbiome but  
22 there are a number of databases on whole genome

1 sequences of pathogens and bacteria organisms that  
2 one could then use to look at microbiome dataset,  
3 which organisms are microbiome, and then use that to  
4 potentially assess risk.

5           Now, you have to step back on microbiomes.  
6 There are two ways of studying the microbiome.  
7 Number 1 is based on a 16S gene. That approach, if  
8 you want to assess the risk associated with a given  
9 organism, there's no better way of saying it stinks.  
10 Okay. 16S sequencing does not differentiate basic  
11 *Clostridium botulinum* and *sporogenes*, between anthrax  
12 and *Bacillus weihenstephanensis*, between *Listeria*  
13 *monocytogenes* and *Listeria* species.

14           So if you do 16S based microbiome  
15 sequencing, in terms of specific risk due to the  
16 presence of a microbial hazard, it's not going to  
17 work.

18           The second option is what people call  
19 shotgun metagenomics, where you sequence all the DNA  
20 in a sample. That's the one which you potentially  
21 can use to map it against these species databases and  
22 then potentially have enough information to look at

1 the risk associated with that food.

2           Now, that's fraught with a whole big area  
3 of problems. Number 1 is you don't know whether the  
4 organisms is alive or dead for starters, okay. So  
5 that's problem number one.

6           Problem number 2 is unless I have one gene  
7 equals risk, so I need a combination of multiple  
8 genes, O157:H7, enterohemorrhagic *E. coli*, textbook  
9 example, whereas the typical short read sequences, I  
10 don't know whether my *stx* gene and my *eae* intimin  
11 gene that allows *E. coli* to attach are in the same  
12 organisms or in two different organisms. So without  
13 that information, it's very hard to assess the risk.

14           So the best thing we can do with  
15 metagenomics data right now, in my mind, is use the  
16 species databases and try to infer risk from it but  
17 it is extremely challenging, and I would in the mass  
18 majority of cases be very, very cautious.

19           DR. STRACHAN: Yeah, I would agree with the  
20 metagenomic data on this 16S. You have to be very  
21 cautious about the resolution that's there.

22           There's also, I think as well as the

1 resolution that we actually want as well and hope  
2 maybe the food industry could make informed decisions  
3 based on, for example, if you're able to find a  
4 *Listeria* in your food product, that would be  
5 something that you're interested in acting on or  
6 thinking about even though it's *Listeria* and not  
7 *Listeria mono*, for example, but I'd be interested in  
8 food industry views on that.

9 DR. CARRILLO: Hi. Is this working? Yeah.  
10 Cathy Carrillo from Canadian Food Inspection Agency.  
11 I have a question I think for most of you. You  
12 brought up the idea that some *Salmonella* or *E. coli*  
13 are less of a problem than others. Without an animal  
14 model to test these assumptions, as regulatory  
15 agency, how do you think we can get to the point  
16 where we can say this *Listeria* is okay. This  
17 *Salmonella* is okay. What sort of evidence can we  
18 provide? How do we know something new didn't come  
19 into the genome, you know, that might be a problem?  
20 Where do you see us going with this?

21 DR. WIEDMANN: I can take a stab at the  
22 *Listeria* example first, but the question really is,

1 you know, if you have SNPs, SNP data and you have  
2 data among human disease cases, saying we have these  
3 SNPs and they're underrepresented among human cases,  
4 so based on the distribution of isolates for certain  
5 SNPs among human and food, it looks like this one is  
6 less likely to cause human disease. So that's an  
7 association, not a cause and effect.

8           The question then becomes how do we go from  
9 association to cause and effect without having a  
10 clear animal model where we can take that *Salmonella*,  
11 *Listeria*, stx and put it into animals and say it  
12 really has a reduced likelihood of causing human  
13 disease.

14           So there's obviously a couple of ways  
15 around it. Number 1, we have a range of animal  
16 models. No animal model is perfect, but they will  
17 help us. With *Listeria*, we have a guinea pig model  
18 that assesses at least a number of factors that are  
19 important very well. We can supplement that with  
20 doing experiments in human tissue culture. We can  
21 grow human cells. We can increasingly grow human  
22 organoids. So not just one type of cells, but sort

1 of something that resembles an organ and use that to  
2 assess the effect of some of these mutations.

3           If those data on, you know, exposure human  
4 disease cases, tissue culture and imperfect animal  
5 model all converge, that's about as good a scientific  
6 evidence as we will get and in many cases, people are  
7 going to probably look at it and say that's  
8 sufficient.

9           So what we've done is a case study with  
10 *Listeria* with these single nucleotide polymorphisms  
11 internalin A, we found that isolates with these SNPs  
12 are about 100 times less likely to show up in human  
13 cases as compared to the ones that don't have those  
14 SNPs. Risk assessments look at large sets of  
15 isolates where we had even exposure data.

16           When we infect human tissue culture cells,  
17 those isolates are about hundred to thousand-fold  
18 less able to infect human cells in tissue culture.  
19 When we put these strains into guinea pigs, they're  
20 about a hundred to thousand-fold less likely to cause  
21 disease. So we have convergence of different lines  
22 of evidence. So with *Listeria* that works pretty



1 well.

2           With *Salmonella*, it's going to get a little  
3 bit more challenging. We can do tissue culture  
4 studies, not perfect. Animal models that really  
5 mimic human disease, it's a lot more tricky, but we  
6 have some animal models that will get us there.

7           Obviously, if we then move to *E. coli* stx,  
8 it gets a lot more tricky. If we then move to other  
9 organisms like *Bacillus cereus*, for example, you  
10 know, we don't have a good *Bacillus cereus* animal  
11 model at all, how do we assess which genes in  
12 *Bacillus cereus* is really responsible for these? How  
13 do we differentiate the *Bacillus cereus* from *Bacillus*  
14 *thuringiensis* which is supposed to be non-pathogenic  
15 by species definition, not always is, it gets even  
16 more challenging. So those are the ones that are  
17 going to require better animal models. It's going to  
18 require better tissue models. It's going to require  
19 that space in between, where we grow organs and  
20 assess, you know, characterize some of these hazards  
21 and characterize some of these organisms in these  
22 models.

1 DR. GALLY: I mean just from the O157  
2 example, I mean certainly we are now aware that there  
3 are particular regions and in collaboration with  
4 groups in Sweden, where a particular area in the  
5 country has the kind of biomarkers as such for the  
6 more serious strains, right. So I think you can then  
7 take that information and then try to intervene  
8 specifically in those regions. Of course, you have  
9 to have the methods to intervene.

10 So this is always the trouble with this. I  
11 think if you're detecting these particular organisms  
12 in food, you can't, at the moment, you're fear being  
13 able to say let's leave it alone. We're not going to  
14 bother with that one. Obviously, that's far too  
15 dangerous at the moment, but I think there are cases  
16 where, for example, we continue to work on vaccines  
17 for this work, and I think we can target particular  
18 herds in particular regions where the more highly  
19 pathogenic bacteria exists.

20 So I think it really has to go hand-in-hand  
21 with other ways that we can intervene with this  
22 knowledge.

1 DR. STRACHAN: Yeah, I would just comment  
2 very briefly. I think it's really important for  
3 scientists when they're doing this sort of work to  
4 explain what their lack of knowledge or uncertainty  
5 is, when we speak about any strains to the degree  
6 that they're sure that they're pathogenic or  
7 otherwise, because risk managers or people in food  
8 factories, have to make decisions based on that, and  
9 if they can get an understanding of what that  
10 uncertainty is, then that will inform them in making  
11 their decisions.

12 DR. WIEDMANN: I think -- the important  
13 things are also what decisions are you going to try  
14 to make with these data. If your decision is simply  
15 you bringing in raw material from five different  
16 farms and you have some information that poultry from  
17 farm X has a certain SNP profile that might indicate  
18 higher risk and you want to process that at the end  
19 of your processing run rather than the beginning, I  
20 probably don't need animal data. I can probably do  
21 that without that. I don't need that high level of  
22 evidence.

1           On the other hand, if I'm going to try to  
2 make some other decisions, regulatory decisions, for  
3 example, you know, the amount of evidence I need and  
4 supporting data I have is very, very different.

5           MR. ROACH: Hello. I'm Steve Roach from  
6 Food Animal Concerns, and my question is actually  
7 sort of related to the first one. And what I'm  
8 concerned about is when you take genomic data from  
9 one environment and then try to apply it to another  
10 one, probably where the question came in my mind is  
11 when you talked about looking at the resistant genes  
12 in West Africa. And are we sure they're going to be  
13 the same as the ones that we've collected in Europe  
14 or in the U.S.? Definitely there's overlap, but  
15 there may be some questions when we start kind of  
16 using genomic data from one environment and then  
17 trying to use it in another one.

18           And another paper that I looked at by  
19 Margaret Davis, several years ago, they looked at,  
20 compared resistant genes on dairy farms versus just  
21 resistant genes on feedlots or calf farms. And in  
22 one environment, there was a lot more resistant

1 selection pressure, and what you found is that genes  
2 were concerned more on where you used more  
3 antibiotics, and that you actually had the genes kind  
4 of drifting where you had less selection pressure.

5           And I'm just concerned about, you know, how  
6 do you actually address that, particularly when we  
7 talk about using genes from this environment or maybe  
8 looking at resistant genes in India that may be very  
9 different than in the U.S.

10           DR. GALLY: It's a huge area, as most of  
11 you will be aware. For resistance genes, there are  
12 particular alleles that, you know, the majority  
13 actually, it's global, in terms of all the different  
14 subtypes. However, there are specific examples where  
15 we can track and identify particular types and you  
16 can then begin to associate those with clades and  
17 surpluses of bacteria that have associations with  
18 particular environments or particular animals or  
19 humans.

20           And I think there it can be quite powerful  
21 in terms of saying that, for example, *Staph aureus*  
22 has moved back from humans and is now in chickens and

1 has come with this resistance. I mean that sort of  
2 assignment is possible now.

3           But I think that for the majority, yeah,  
4 we're chasing our tails with this. It's really,  
5 really difficult to know the direction of travel of  
6 resistant genes but again the more of this  
7 information we have, the more we can address these  
8 questions. I mean at the moment, all these studies  
9 are conducted independently, and then trying to  
10 compare them is more or less impossible.

11           I think the power of looking at, with  
12 metagenomics is, for example, we're doing a study on  
13 pig farms at the moment, is the deep sequence  
14 environmental samples and get a real readout relative  
15 to which antibiotics are used in those facilities as  
16 to what resistance genes are present and how that  
17 fluctuates with time.

18           It's not about blame. It's about  
19 understanding what's there, what's out there, and  
20 then what can come through in the product.

21           But I mean beginning to trace AMR genes in  
22 foods and things, that's a whole other massive issue,

1 isn't it, but it's one we might obviously be starting  
2 to do, and there's a push now for obviously supplying  
3 poultry and supplying various livestock that have got  
4 minimal interventions of antibiotics. So clearly to  
5 track that, we're probably going to be asked to do  
6 that, right. So it's a big area.

7 DR. WIEDMANN: I mean the databases we use  
8 to go from genome to phenotype or predicting  
9 phenotype is obviously a huge challenge. It's not  
10 just limited to antimicrobial resistance.

11 Part of the issue again goes, what are you  
12 trying to do with these data? If I'm going to try to  
13 pull out a genome sequence for antimicrobial  
14 resistance and predict resistance of an organism, if  
15 I have the wrong database, that can cause challenges  
16 and I might incorrectly predict an isolate as  
17 sensitive for resistance when it isn't, and we have  
18 all, I think, seen this. The databases get better,  
19 but every time they get better, there are continuous  
20 issues with some of that.

21 You need to have some subject matter  
22 expertise, and then particular, if you move into

1 different environments, you're going to run into  
2 problems with that.

3           So the databases need to continue to be  
4 built. We need to very, very careful about  
5 extrapolating. I think your example is a really  
6 great one if you're starting to analyze sequence,  
7 come from environments but has exposure to different  
8 antibiotics, different countries, different treatment  
9 of animals, for example, that we don't have to  
10 rebuild the database and we're going to make some  
11 wrong calls absolutely until the databases get  
12 better.

13           DR. BRADEN: Chris Braden, Centers for  
14 Disease Control. I want to segue into the previous  
15 question and discussion about the effort that a  
16 number of partners have made to build some of the  
17 databases for whole genome sequencing and the  
18 metadata that goes along with it. Certainly we use,  
19 you know, to make the public repository at the NIH  
20 NCBI.

21           And my question is, you know, have you  
22 really used that database or others like it? Is it



1 sufficient in the quality of the sequence data and  
2 the comprehensiveness of the metadata that go with it  
3 in order to conduct some of the studies that you're  
4 trying to conduct in these predictive models?

5 DR. WIEDMANN: So the NCBI database, I  
6 think is useful for some applications but I think in  
7 general, the metadata that are there are probably not  
8 there yet to allow some of these investigations. You  
9 know, it depends on what you want to do with it. If  
10 you want to do source attribution, for example, I'd  
11 be very challenged to say we can use the NCBI data  
12 where they are right now to really do source  
13 attribution. I think we can get there but I have not  
14 seen much people that have actually validated that  
15 the predictions of the source attributions are always  
16 correct. So I think we need better metadata, if we  
17 want to do source attribution.

18 And the other question is do we want to use  
19 these data to predict, you know, do hazard  
20 characterization or hazard identification?  
21 Typically, the metadata there are limited. I have  
22 not seen, for example, in human cases, a lot of

1 metadata on disease severity and disease sometimes at  
2 a level of resolution that will help us to then link  
3 SNPs or other genomic characteristics that are likely  
4 to cause disease, and for obvious reason. I mean I  
5 understand there's huge issues with regard to, can we  
6 track back to who that person is if we give enough  
7 geospatial plus symptoms plus age plus predisposing  
8 factors which we need in *Listeria* really because we  
9 have this interaction of, you know, human  
10 susceptibility plus food plus organism.

11           So I think the data are useful in some ways  
12 but with some of the questions we're really trying to  
13 ask right now, and more often than not, I don't think  
14 they get us where we want to go.

15           DR. STRACHAN: Yeah, I would agree with  
16 that, and also the metadata, I think for the reasons,  
17 if people aren't willing to put these up on databases  
18 and cases traceback to individual human cases, it  
19 becomes a problem or companies perhaps for that  
20 matter.

21           And plus I know with some other databases  
22 is that there are metadata that are hidden across the

1 front public end where particular users can't use  
2 that under particular agreements. So I think it is  
3 quite a big challenge.

4           And you know, for myself, you know, for  
5 research purposes, you know, if I find a sequence on  
6 like my NCBI database that will relate far back to  
7 the original paper it came from, I go back to the  
8 original paper and I dig the metadata out from the  
9 original paper which isn't a very efficient way of  
10 doing things in some cases but, yeah, that's just the  
11 way it is.

12           DR. GALLY: I was at IAAP in Florida in the  
13 summer, and probably the reason coming to that is you  
14 see what is available and some great conversations  
15 with folk around I mean a whole bunch of *E. coli*  
16 sequences just released actually, and we received the  
17 database of source attribution studies and linking  
18 into other groups in the U.S. to try and work on  
19 that. So from an *E. coli* point of view, that's been  
20 incredibly helpful.

21           But I think a key problem here is exactly  
22 what Norval's stated which is about, especially with

1 human data, it's really powerful if you know the  
2 degree of severity of disease, if you can have  
3 geographical insight into that as well. This is  
4 really difficult information to get out in the UK.  
5 I'm not talking about USA, but it is really powerful  
6 data if you can have it because it really allows you  
7 to link to the information you're getting from animal  
8 sources as well.

9           I think we're a really long way from that  
10 information becoming easily available. I think some  
11 of the barriers could be broken down quite simply by  
12 deciding what type of information and the level of  
13 granularity of that, that can safely be put on  
14 automatic databases. They do it in Scandinavia  
15 better, and I think there's some really nice models  
16 to follow out there.

17           So I think people are beginning to get the  
18 message that we have to do this, but there's still a  
19 huge number of hurdles, and it's going to be very  
20 organism specific as well. Fantastic *Salmonella*  
21 examples and other organisms we're really light on  
22 data.

1 DR. ALLARD: Thank you. This is Marc  
2 Allard, FDA, Center for Food Safety and Applied  
3 Nutrition. I just want to comment on the last set of  
4 comments, and then I have a question.

5 Essentially the Government's data at the  
6 FDA is all FOIA-able. So it's just a matter of how  
7 much is released to the public and how soon, but I  
8 think they're fully open to recommendations of what  
9 additional information should go into the metadata,  
10 and as long as it's legally allowable, I think the  
11 Government's willing to release that information.

12 And so my question is essentially going  
13 directly at risk. We have a model for phenotypic,  
14 genotypic prediction within NCBI where they built the  
15 bio project of the 4,100 resistant genes and then  
16 every new genome is blasted against it to call the  
17 presence and absence. I believe they don't do  
18 allelic differences yet, but they have presence and  
19 absence.

20 So my question for the risk group is, this  
21 would be easy to replicate, build genomic databases,  
22 bio projects, for specific genes. The question is,

1 it's not clear where to start. Which genes should we  
2 start? Which ones are clearly risk connected? We  
3 know there's been some publications with the STECs  
4 and the NACMPF documents as well as the French  
5 recently released a big study on *Listeria*.

6 So my question is do we think we should be  
7 building these databases? Do we have lists of  
8 recommendations? How do we move forward with rapid  
9 prediction? I want to build tools for the public.

10 DR. STRACHAN: I think that's a very good  
11 point. I don't know, because I've done a lot of  
12 *Campylobacter* source attribution work, but I think  
13 you speak about the French project as well with  
14 *Listeria*. I think whole genome MLST is a good place  
15 to start because you have basically all the core  
16 genomes and core genes and also possibly a number of  
17 accessory ones as well that can be mentioned, and I  
18 think that would probably be a good way to start.

19 DR. WIEDMANN: I think there's some, maybe  
20 a few examples where we might be ready, internalin  
21 genes of *Listeria*, pretty clearly linked to risk, a  
22 lot of data. I think FDA has some data. So, for

1 *Listeria*, we can do some of it.

2           I think as we move to *E. coli*, I rely on  
3 the two of you to tell me whether you're ready. I  
4 think the problem goes back to that very often it's  
5 not one gene determines, you know, risk, but it's  
6 interactions between different genes. How do you get  
7 at that, and that's going to be much more  
8 challenging.

9           So I would really want to get that before  
10 I'm going to put out a simplistic tool where we  
11 pretend one gene can ultimately predict risk and not  
12 looking at interactions of different genes.

13           The European or Germany *E. coli* outbreak  
14 was unusual. *E. coli* was an unusual attachment gene  
15 is a great example of where, you know, could we have  
16 used that? How could we have predicted that risk  
17 because all of these databases are a priority and  
18 maybe, maybe not. I don't want to make a judgment on  
19 it.

20           And I think the antimicrobial resistance  
21 databases provide some examples of the risks of it.  
22 I mean there are certain genes in there because you

1 don't have the resolution. You might predict that  
2 organism is antimicrobial resistant and it is not.  
3 We have seen it be published and others have  
4 published it, streptomycin resistant, very difficult  
5 to predict because there some allelic variance.

6           So I think we want to be very careful and I  
7 think put some uncertainty around these predictions,  
8 too, when we do it. I think there are some examples,  
9 but honestly I think there are only a few right now  
10 where I feel like we are ready to put something like  
11 that together and put it to use other than research  
12 use. Maybe we can put it to use for some sort of  
13 risk rankings, again back to my sample, you know,  
14 which poultry flock, which cull cows, are you going  
15 to process first, second, third? So I think we might  
16 be able to use some of that information.

17           So all of these things we need to think  
18 about what we're going to use these data for. We can  
19 use it for so many decisions, and the trick always is  
20 to use it, if you have the right database for  
21 decision A, but we start using it for decision C,  
22 we're going to run into trouble.



1 DR. STRACHAN: I guess another thing in  
2 passing, I was just thinking about STEC, there's a  
3 virulence finder which is developed by Flemming  
4 Scheutz from SSI in Denmark. So you can basically  
5 upload your genome, whether that be an assembled or  
6 the reads, and it will feedback a number of virulence  
7 genes, whether they're present or not.

8 DR. BRANDT: Alex Brandt from FSNS. I have  
9 a question with regard to, I know we're kind of  
10 talking about relatedness of presence/absence of  
11 genes and different alleles and even some of the  
12 phenotypic traits don't always match up. So is it  
13 enough to just look at presence/absence of genes or  
14 different alleles, or should we be looking deeper at  
15 like transcriptomics and really going to that level?  
16 I guess that's my question simply.

17 DR. DESSAI: Before you guys answer the  
18 questions, we are running a little over our time, and  
19 we'll extend the time by about 5 minutes if that's  
20 okay with the crowd here. We also have two questions  
21 online to address. Is that right? Okay. Go ahead.

22 DR. WIEDMANN: I think in some questions

1 you can have a gene present. There's examples. We  
2 have a gene present and it's not transcribed. It's  
3 not turned on. Therefore, you don't get the  
4 phenotype. I think looking at RNA, looking at  
5 transcriptomics can help you with this, and  
6 ultimately you should be able in many cases to then  
7 predict polymorphisms and SNPs and promoters in  
8 regions that drive transcription so that you're not  
9 going to have to do routine transcriptomics for it,  
10 where then it becomes obviously completely different  
11 issue is, if you're starting to assess risk in  
12 different foods.

13           Something every food microbiologist knows,  
14 *Clostridium botulinum* is only an issue if it's found  
15 in a food that's anaerobic conditions because that's  
16 when the genes are turned on. So if I want to assess  
17 that risk, I need to have transcriptomics and some  
18 other data because gene presence/absence does not  
19 equate risk at all. So we have this range, too, in  
20 the thing. So in these cases, whether we need  
21 transcriptomics or something else is a different  
22 question but gene expression is very, very important

1 and we can't forget that different foods based on  
2 anaerobic condition represent completely different  
3 risks.

4 DR. DESSAI: Okay. We have online  
5 questions.

6 DR. NAHAR: One of the online participants  
7 asked what about focusing hazard characterization on  
8 sampling physicians, veterinarians, nurses, farm  
9 workers, food preparers, etc.? Those populations  
10 tend to be at highest risk for the spread of  
11 pathogens as well as developing AMR.

12 DR. DESSAI: Will you repeat the question?

13 DR. NAHAR: So it's a sampling question.  
14 What about focusing on sampling this particular  
15 population, farm workers, physicians, veterinarians?

16 DR. WIEDMANN: I think the question is what  
17 about risk assessment, I think risk characterization  
18 might not have been right, and isolate collection  
19 from people at high risk, particularly on the primary  
20 production side, the farm workers and anyone who is  
21 working in these areas of a wide range.

22 So I think is -- I'll try first stab.

1 That's a tough one. I think it will give us  
2 different isolate sets that could be useful. I don't  
3 think that we not focus on it that much is that big  
4 an issue. I think they're going to be caught  
5 indirectly through the public surveillance system,  
6 too. We get into a whole slew of social issues, you  
7 know, illegal immigrant farm workers, are they going  
8 to seek? You know, what's the reporting among those?  
9 Are we missing cases because of whole set of other  
10 issues which I think is a very, very important one.  
11 Are we capturing all high risk individuals without  
12 surveillance? And I think that's something worth  
13 thinking about it, and it could be very good  
14 sentinels, particularly for some emerging zoonotic  
15 diseases.

16 DR. STRACHAN: Yeah, I think the point, I  
17 look at it sort of in terms of like source  
18 attribution, thinking about *Campylobacter*, workers  
19 working within a poultry factory might have  
20 particular exposures, but as you think about rural  
21 children living in Northeast Scotland, and also in  
22 the USA as well, you're more likely to have contact

1 with farm animals or private watch -- are going to  
2 get different spectrum of types of *Campylobacter* for  
3 example. So I think looking at the different groups  
4 will tell you something about what -- confirm  
5 different exposures they have in the pathways they're  
6 getting the disease from. So I think from a  
7 molecular epidemiological perspective, that could be  
8 quite helpful.

9 DR. GALLY: I mean there are obviously -- I  
10 mean I'm aware of quite a number of funded studies  
11 where that particular close association between  
12 product, livestock rearing and people that are  
13 working with livestock is being looked at.  
14 Obviously, studies throughout the last 50 years where  
15 the different typing tools have been used to do that.  
16 There's current studies doing that, using whole  
17 genome sequencing now, and it's really looking at  
18 transfer of particular organisms and AMR and  
19 virulence on those close quarters, but I suppose it's  
20 -- yeah, I mean it doesn't really necessarily help  
21 address the bigger issue we've then got of I suppose  
22 the distribution that can occur through product which

1 is really a much bigger scale.

2 I suppose you could develop -- I suppose  
3 there are the chances of those -- when we think about  
4 viral evolutions, there are the chances of maybe that  
5 being the epicenter of something kicking off, that  
6 obviously ends up being a lot more serious, do we  
7 give that -- I think the question is, do we give that  
8 extra attention? Do we do extra monitoring of all  
9 those folks that are in those direct environments,  
10 handling birds, etc., because there's more likely to  
11 be the start of a particular -- I think we're  
12 probably a very long way from that, I would suggest.

13 DR. DESSAI: Okay. We'll take our last  
14 question, and then we will close this session.

15 DR. NAHAR: Sure. Last question, can you  
16 speak to the pros and cons of hazard ID based on WGS  
17 data alone showing AMR gene prevalence versus how  
18 we've traditionally conducted such work using  
19 phenotypic AMR?

20 DR. GALLY: I can only speak to what I'm  
21 aware of with a couple of different pathogens but  
22 while there are obviously issues in getting phenotype

1 from genotype and it's particularly difficult for  
2 whether it can be single nucleotide changes that  
3 comes out, certainly the horizontal nucleotide  
4 resistant genes, it's actually pretty good. I mean  
5 certainly for the *E. coli*. There's a high degree of  
6 accuracy to presence of a particular allele and then  
7 giving phenotypic resistance. So it certainly can  
8 work very well in that space. But, you know, you  
9 can't have 100% accuracy with that, but it's a very  
10 good genotype to phenotype mix in terms of whether  
11 the horizontal nucleotide resistant genes.

12 DR. WIEDMANN: I may be a little bit less  
13 bullish on this, and the example I'm going to give,  
14 if you look at *E. coli* and *Shigella*, if you look at  
15 resistance measured with traditional approaches, so  
16 MIC, etc., the way you interpret them in terms of  
17 susceptible intermediate resistant or resistance,  
18 actually differs between *Shigella* and *E. coli* despite  
19 the fact that *Shigella* is actually an *E. coli*. So  
20 that translation from presence of genes to then  
21 treatment decisions, you know, does a certain  
22 treatment work in at a different level is challenging

1 in some very specific cases. So I think we really  
2 need to look at it from organism to organism, and  
3 again it depends on what you mean about what sort of  
4 risk. Is it the risk of a resistant infection if  
5 this gene is found in a certain organism? I think  
6 for some organisms we might be there. For a lot of  
7 them, we really need more data to make these  
8 linkages.

9 DR. GALLY: I suppose I'm thinking  
10 particularly the arena of real-time diagnostics and  
11 particularly human and animal infections where at  
12 least using that information from initial sequencing  
13 of a direct sample which is where we're trying to  
14 move with some of this, it's still much better than  
15 being in the dark. So it's at least having some  
16 information you can make an informed decision on.  
17 But, yeah, we still don't understand all the  
18 complexities to other things that will influence  
19 resistance, but certainly I think it's a step in the  
20 right direction.

21 DR. DESSAI: Okay. Just an announcement  
22 for those online: Please send your questions which



1 are short and focused so those will be easily  
2 transcribed and understood here.

3           Number 2, if you have any other  
4 suggestions, send those to us.

5           So I think we've set the stage for this  
6 conference pretty neat and covered a lot of areas.  
7 The speakers will be available for you to have more  
8 discussion or if you have any questions, and I think  
9 I would like to thank them profusely for the work  
10 that they have done this morning. Please join me  
11 with a big round of applause.

12           So we'll be back at 11:00.

13           (Off the record at 10:28 a.m.)

14           (On the record at 11:01 a.m.)

15           DR. EVANS: Welcome back from the break.

16 My name is Peter Evans. I'm with the FSIS Office of  
17 Policy and Program Development, and I want to welcome  
18 you to our second session which is on Federal/State  
19 Collaboration. We're going to have eight speakers  
20 today, and this is basically going to take us from  
21 now to the lunch and then after lunch to the end of  
22 the day. We're going to start off with three

1 speakers, have a very short question and answer, 5  
2 minutes, go to lunch. And then we come back from  
3 lunch, we'll have three more speakers, a break, and  
4 then two speakers.

5           So we're going to hear presentations about  
6 how whole genome sequencing is being used in the  
7 United States both at the federal and state level,  
8 and so how WGS capability is being increased and also  
9 examples of how organizations are using the data.  
10 And then we're also going to learn about how the  
11 agencies are working together, depositing data in a  
12 common database at NCBI and then also developing  
13 common procedures and standards through the NFS  
14 Consortium.

15           So first I'm going to welcome Dr. John  
16 Besser from CDC, Dr. Steven Musser from FDA CFSAN,  
17 and Dr. David Goldman from FSIS, to speak about their  
18 experiences within their organizations. Thank you.

19           DR. BESSER: Well, thank you. I'm very  
20 pleased to have the opportunity to speak to you  
21 today. The organizers gave me a long list of  
22 suggested topics. So I'm going to move along pretty

1 quickly here. I really hope I never need to use that  
2 picture in a dating app.

3           So I've been around long enough to remember  
4 the beginning of the PFGE era, and what was really  
5 interesting was a lot of the concerns I hear now from  
6 laboratorians, epidemiologists, regulators, industry,  
7 are very similar to what we heard 20 some years ago,  
8 and so I'm going to talk about what hasn't changed,  
9 what's really an extension of what we've been doing  
10 for a very long time.

11           I'm going to talk about the WGS  
12 infrastructure that we're building at CDC, but I'm  
13 going to talk about then what's different, what can  
14 we expected different about whole genome sequencing  
15 and then I'll touch very quickly on where we see this  
16 all going.

17           PulseNet, as you've heard, has been around  
18 for 21 years now. We just celebrated the 20th  
19 anniversary last year, and over those years, we've  
20 seen investigations with our colleagues in FSIS, FDA,  
21 and the states and with industry, impact virtually  
22 every commodity of food. It's been tremendously

1 impactful.

2           And the way it works, of course, is that we  
3 have combined in a One Health model streams of data  
4 from food monitoring programs, animal monitoring  
5 programs and human disease monitoring programs, all  
6 into this one system. And the way it works is we  
7 connect cases that may be geographically distant from  
8 each other by this common DNA fingerprint.

9           We released a study just last year that  
10 looked at the cost and benefit of PulseNet. This was  
11 during the PFGE era, and we found that at a minimum,  
12 it saves about 270,000 cases of disease per year and  
13 about \$500,000,000 in cost to society. And that's  
14 really the portion of the benefits that we could  
15 measure which I suspect is really a very tiny  
16 proportion of the total benefits.

17           We investigate about between 30 and 60  
18 national clusters per week and at the state level,  
19 there's about 1,500 to 2,000 investigations per year.  
20 So this is a constant, very big activity.

21           And I'll describe a little bit about our  
22 network which has driven how we've built the

1 infrastructure. It's based on local testing, local  
2 control of the analysis using standard computers,  
3 centralized quality assurance, so that we're all able  
4 to communicate in a standardized way, and in the WGS  
5 era, we're going to have centralized bioinformatics  
6 and centralized high performance computing.

7           And the rationale for having this  
8 distributed testing network is that it provides  
9 functionality for both local investigations, national  
10 and international investigations, and it helps with  
11 turnaround time which in our world is important, and  
12 I'll show you later how important it is. It gives us  
13 an enormous body of resources. With 86 laboratories  
14 in the network, we can expand or contract the work as  
15 needed because there's so many independent parties  
16 that can do the work. And it gives us local control  
17 of patient identifying information and commercial  
18 confidential information. That's all held, and  
19 you'll see more later.

20           So what's the same as it's been? Well, the  
21 subtyping methods, be it performing or whole genome  
22 sequencing, primarily work by grouping together cases

1 that are most likely to share a common exposure, such  
2 as food. They work the same way.

3 Matches between cases and food and  
4 environmental samples in both circumstances provide a  
5 hypothesis as you've heard many times already today,  
6 not a proof, but they prove a hypothesis in both of  
7 these circumstances for both methods.

8 And in both methods, the historical  
9 database is routinely examined for matches to current  
10 clusters. So, for years, we've been identifying food  
11 or environmental or human cases that match current  
12 clusters. So that's actually not new.

13 And how we interpret this data also has  
14 many similarities. A match with both methods, what  
15 we call a match, which is another subject, means that  
16 an association is more likely than if there is no  
17 match. You notice I said more likely. So it's not  
18 absolute.

19 A mismatch means that an association is  
20 less likely but not impossible, and I'll show you why  
21 that's the case.

22 And as you've heard over and over again,

1 other types of data are needed to support any of  
2 these conclusions be it whole genome sequencing or  
3 PFGE, such as epidemiological data and traceback  
4 data.

5 I'm going to focus on this for a moment.  
6 Why does the mismatch mean -- does not absolutely  
7 mean there there's no association? Well, it turns  
8 out that outbreaks are all different from each other.  
9 They can be very clonal as we sometimes see when an  
10 individual food handler contaminates food and a  
11 single strain gets into a population, but there's a  
12 whole spectrum of different outbreak ecologies, if  
13 you will, that span everything from that to mass  
14 contamination events. I used irrigation water here,  
15 but there's negative pressure events where sewage  
16 gets into drinking water, and there's no reason that  
17 in that circumstance you'd expect there to be a  
18 single strain and everything in between. So  
19 outbreaks themselves are not necessarily purely  
20 clonal.

21 And we've seen that in the PFGE era. This  
22 was an major outbreak of listeriosis associated with

1 cantaloupe a few years ago, and we saw multiple  
2 different serotypes, multiple different PFGE types  
3 because the product was contaminated with an  
4 environmental source that had many different germs in  
5 it.

6           And we've seen it in the WGS era. This is  
7 a phylogenetic tree of cases that were all associated  
8 with having consumed chicken from a particular  
9 processing plant and we can see that there's a very  
10 variation between these cases in terms of SNPs, high  
11 quality SNPs, all associated with the same sources.  
12 This is an example of a mass contamination event.

13           So what kind of infrastructure are we  
14 building? Well, this is our master diagram, and it  
15 really can be divided into two sections, one that's  
16 closed to the public, where information that is  
17 potentially identifying to individual patient level  
18 is kept or commercial confidential information,  
19 information which can't be released to the public,  
20 and everything below that line is information that  
21 will be available to the public.

22           Isolates are sequenced, controlled by the



1 local user. They go through a calculation engine  
2 that pulls out different types of information which  
3 I'll show you in a moment, and then also looking at  
4 allele databases, looking at different alleles and  
5 reference identification, etc.

6           And we're building the system to coordinate  
7 not only with the other federal agencies, but with  
8 global partners and with the states and local  
9 governments and all these other institutions listed  
10 on the left. So a lot of different parties that we  
11 need to coordinate with.

12           And we're building the core genome MLST  
13 databases in collaboration with global partners,  
14 because we're looking at global systems for food  
15 which is inherently a global commodity, and we're  
16 also paying close attention to the developments at  
17 NCBI.

18           And recently, a meeting was held with  
19 PulseNet International to hash out some of the global  
20 issues. How do we collaborate around the world and  
21 they came up with this vision paper which you should  
22 read when you can.

1           So these databases automatically pull up  
2 all this different kind of information, the genus and  
3 species, serotype, pathotype, the virulence  
4 information, the resistance information, all types of  
5 different subtyping information from core genome and  
6 whole genome MLST and plasmid profiles, and also a  
7 SNP analysis is available if desired. It will be  
8 shortly. All this information automatically  
9 populates the local databases.

10           This is an example of the type of report  
11 that we can expect from individual isolates. We can  
12 have all the virulence information, serotype  
13 information, resistance information, etc., so that we  
14 can start properly learning about risk based on  
15 different, for instance, non-O157 subtypes.

16           This is where we're at with building the  
17 infrastructure. As of yesterday, 49 states now have  
18 sequencers, and 37 states have been certified by  
19 PulseNet as able to perform these methods.

20           This is a list of all the different  
21 databases and where we're at with their development  
22 and release. *Listeria* is the only one that is fully

1 released, but the rest are soon to follow.

2           The number of genomes that have been  
3 sequencing has been going up very dramatically but  
4 hold onto your hats. It's going to be a lot more  
5 very soon. Within, by the end of FY18, we hope that  
6 most *Salmonella*, all STEC and all *Listeria* in the  
7 United States from clinical cases will be sequenced.

8           So what's different? Well, the PFGE was  
9 very specific for some organisms and low for others  
10 as you saw for *Listeria*, Enteritidis, whereas  
11 sequencing is very high for the most part in all  
12 organisms. We have limited ability to evaluate the  
13 closeness of strains with PFGE where we can do that  
14 with whole genome sequencing. The data is  
15 categorical pretty much with PFGE. It matches or it  
16 doesn't, where it's continuous, there's different  
17 shades of closeness with whole genome sequence data.

18           And here's how we actually used it in the  
19 PFGE era compared to how we used the data in the  
20 whole genome sequence era. These are all cases  
21 associated with the consumption of *Listeria*  
22 associated with a particular type of ice cream. In

1 the past, we would have focused on a single PFGE  
2 type, and we would at the onset, ignore the other  
3 cases that we couldn't draw together, and use that  
4 information to make the association as part of the  
5 investigation.

6           Now, in the whole genome sequencing era, we  
7 can look at all the different types that are related  
8 to each other and in this particular case, there was  
9 11 PFGE types and only 2 sequencing types that were  
10 used, making the investigation much simpler.

11           But the impact of all of this is that we  
12 can detect outbreaks when they're smaller, which  
13 means that we catch them earlier, they're smaller,  
14 and we can solve more outbreaks than we could in the  
15 past.

16           Also, one thing that hasn't been discussed  
17 is that this information is invaluable for ruling out  
18 likely outbreaks. I say that probably happens more  
19 than not, where we can say, well, this product  
20 probably isn't involved in this situation.

21           And when we interpret the cases, matching  
22 of cases with products or environment, in the past,

1 when we had a very new or a very rare PFGE subtype,  
2 it formed a strong hypothesis, but when it's a common  
3 subtype, we really couldn't say it was a weak  
4 hypothesis. With whole genome sequencing, a product  
5 match is uniformly a strong hypothesis.

6           So you've seen this slide already, but this  
7 is the different phases of our *Listeria* program, a  
8 collaboration between the agencies. In the early  
9 days, there was few outbreaks represented by the blue  
10 bars. When we turned on PulseNet, a lot more  
11 outbreaks, and the outbreaks got smaller, a median of  
12 69 to 11. When we did an epidemiological project  
13 called the *Listeria* Initiative, more outbreaks still  
14 and the median size went down to 5. Now, in the  
15 whole genome sequencing era, more outbreaks yet, and  
16 it went down to 4.

17           So what we've seen, summarizing all this  
18 data, is a dramatic increase in the number of solved  
19 outbreaks and a lowering in the size of outbreaks.  
20 So this means to industry that these things are found  
21 with more surgical precision. They're solved quickly  
22 when there's less cases, less potential impact.

1           And so we get the question a lot about what  
2 kind of cutoff values? Well, we have these rule of  
3 thumb cutoffs but really what I'd like to tell you,  
4 if you remember this diagram, there's really no such  
5 thing as an absolute cutoff. So absolute cutoffs are  
6 not possible, and really we need to combine  
7 epidemiology, traceback data and the other data as a  
8 part of the total solution for defining what's  
9 important and what's not.

10           And these give strong hypotheses, WGS does,  
11 but again it doesn't absolutely mean that a match  
12 between a cases and a product means causation as  
13 Dr. Wiedmann said earlier. The food chain is  
14 extremely complex as this slide from FDA shows you,  
15 and if we look at transmission mechanisms, the chain  
16 of transmission can be even more complex and WGS  
17 doesn't tell you anything about the chain of  
18 transmission.

19           Okay. Just a few moments in my last couple  
20 of minutes. What's ahead? Well, we're working on  
21 these advanced analytics, machine learning,  
22 disjunctive anomaly detection, etc., and we're also

1 exploring the use of metagenomics for food safety as  
2 many other institutions are. And as you may know,  
3 there's a big collaboration between Mars Company and  
4 IBM to look at these very issues.

5           There's a wide range of potential  
6 applications for metagenomics in our field. We're  
7 focusing on two, pathogen discovery and in situ or  
8 direct from specimen pathogen characterization. And  
9 we're looking at several methods. I've listed here,  
10 amplicon sequencing and shotgun metagenomics, but  
11 much work remains.

12           But there's good reasons we're doing all of  
13 this. Most pathogens that cause disease, in the  
14 United States and worldwide, are thought to be  
15 unknown. In fact, our PulseNet pathogens are only 4%  
16 of the total causes of disease. So there's a  
17 tremendous amount to be learned by doing pathogen  
18 discovery.

19           And these direct from specimen methods have  
20 the opportunity for us, majority shortening the  
21 timeline between when a patient consumes a food and  
22 when there's an actionable result, which is currently

1 quite long. In addition, then we have to find a  
2 cluster. We have to identify a food. These direct  
3 from specimen efforts can cut weeks off of the whole  
4 process potentially leading to more outbreaks solved  
5 more quickly and more illnesses prevented.

6 And finally, there's a problem in the U.S.  
7 from what's called culture-independent diagnostics.  
8 We're potentially losing isolates because of changes  
9 in diagnostic testing in both PFGE and whole genome  
10 sequencing dependent isolates. So we're putting a  
11 lot of effort into developing tests that don't depend  
12 on isolates.

13 And that's it. Thank you very much.

14 DR. MUSSER: I'm happy to be here. My name  
15 is Steve Musser. I'm the Deputy Director for  
16 Scientific Operations at the Center for Food Safety  
17 and Applied Nutrition of FDA.

18 We've been doing sequencing and applying it  
19 to food safety applications for a number of years,  
20 and what I'd like to try and do today is walk you  
21 through FDA's approach, and I would also like to  
22 caveat a lot of that with, while FDA and FSIS and CDC



1 have very similar approaches, FSIS' and FDA's are  
2 legal and regulatory requirements are also different  
3 which means that we respond and act in different ways  
4 in cases of finding either outbreaks or cases of  
5 contamination in facilities.

6           So with that said, if I could have the next  
7 slide. So I'd like to talk a little bit about  
8 GenomeTrakr and the GenomeTrakr laboratories and the  
9 network and the database. All of these technologies,  
10 and it would mean whether PFGE or whole genome  
11 sequencing or whatever, are useless unless you have  
12 some way of comparing them, and some way of putting  
13 them into a place that you can look at them and then  
14 you also need people to be supplying information to  
15 those databases.

16           So GenomeTrakr was essentially begun by  
17 FDA, and it's essentially a larger dataset of  
18 information than anyone of the other single datasets.  
19 So if you were to look at, you know, what was just  
20 collected by somebody's personal academic group or  
21 PulseNet, for example, so GenomeTrakr contains  
22 everything in PulseNet as well as lots of other

1 things outside of PulseNet.

2           And so it would have public health,  
3 government, private, academic sources of information  
4 and just to give you some idea, I know John mentioned  
5 there's about 40 to 50,000 sequences in the PulseNet  
6 database. There's over three times that many in the  
7 GenomeTrakr database. So it's just a larger  
8 collection, and there's advantages and disadvantages  
9 to each, and I think based on the mission of the  
10 different agencies and groups, there's good reason  
11 for that.

12           So we established this network in 2011 with  
13 a fairly small investment in state laboratories and  
14 our FDA field laboratories, FSIS and CDC joined as  
15 well as dozens of other collaborating organizations  
16 following that, and I would also like to highlight  
17 the National Center for Biotechnology Information,  
18 NCBI, which we would not able to do any of this  
19 without because they not only serve up the data but  
20 they provide analytical tools and ways of looking at  
21 the information which you wouldn't be able to do.

22           And when we began this network, we were

1 also trying to answer two questions. The first one  
2 was, we need to provide this information in a way  
3 that is public. We were under a Presidential  
4 Executive Order to do so, not just ours, but all  
5 Government supplied information. We also didn't have  
6 the money to actually supply all this or maintain all  
7 this. So we needed someone to help us. And NCBI was  
8 more than willing to do that, and they've been with  
9 us from the start, really helping and making this a  
10 useable database.

11           So unlike the rest of NCBI, which is, you  
12 know, you could just upload and you wouldn't be part  
13 of GenomeTrakr, these actually, people that are part  
14 of the GenomeTrakr Consortium actually do fill out a  
15 form and there is some information there. If I could  
16 have the next slide.

17           So one of the issues that we had when we  
18 first began looking at what would be used instead of  
19 PulseNet or in supplementing PulseNet was we have  
20 this great body of clinical information that CDC has  
21 gained by using the PulseNet system, and while there  
22 were some food isolates in it, the majority of

1 isolates were from clinical sources.

2           So what we found out when we were trying to  
3 do these investigations is you know that a whole  
4 bunch of people are sick, but you don't really have  
5 any idea of where to look. And the reason for that  
6 is there's no way to look for that information  
7 because there aren't any sequence data.

8           So we've concentrated a lot on obtaining  
9 and getting food and environments isolates in  
10 addition to clinical isolates because we recognize  
11 that the maximum benefit, particularly in outbreak  
12 situations is kind of knowing where to look or  
13 helping understand where to look.

14           And if you look at PFGE in the past, it was  
15 primarily driven when we found a group of people that  
16 were sick, and then we do a lot of epidemiology and  
17 we really solved a lot of outbreaks and intervened in  
18 very positive public health ways.

19           What we're seeing now with whole genome  
20 sequences is that sometimes our information in  
21 sequencing helps drive the epidemiology or helps  
22 refine the epidemiology. So it's not a one-way

1 street. It's actually information that can provide  
2 data in both ways, and really help inform our  
3 investigations. Can I have the next slide?

4           So currently the GenomeTrakr Network  
5 includes all of the FDA labs, all the CDC labs as  
6 well as the labs that CDC would fund and contribute  
7 to PulseNet. FSIS's laboratories, we fund 19 state  
8 agriculture, health and university labs in the U.S.,  
9 1 hospital lab and 17 labs located outside the U.S.  
10 so we have 4 continents and 10 countries contributing  
11 information.

12           There's approximately 150,000 sequences in  
13 the database now. It's growing at a rate of between  
14 5- and 8,000 sequences a month. So that number, it's  
15 a static number now. By the end of the month, I'll  
16 be quite different.

17           And then we partnered with CDC and FSIS in  
18 2012 to do all clinical *Listeria monocytogenes* and  
19 environmental. So if FDA or FSIS found environmental  
20 samples, we uploaded them and sequenced them and  
21 likewise, if CDC and its network found clinical  
22 samples, they sequenced and uploaded them. It's

1 probably the best *Listeria monocytogenes* database in  
2 the world in terms of completeness over the last  
3 couple of years.

4           You can find out more about the GenomeTrakr  
5 website at the link at the bottom or you can just go  
6 to [fda.gov](http://fda.gov) and search GenomeTrakr or Google.

7           Next slide.

8           So this is just a cartoon way of visioning  
9 what happens. At the top, anyone that's part of the  
10 network or the consortium performs a sequence, they  
11 upload it into the NCBI database, or the European  
12 version of it, EMBL, or the Japanese version, DDBJ.  
13 These databases are copied within each other every  
14 night. So there's a version of all of these  
15 sequences available throughout the world, and so if  
16 NCBI were taken down through a power outage or other  
17 natural disaster, you could get at the data through  
18 EMBL or DDBJ. So it's part of the redundant system  
19 that's been built at NCBI.

20           This is a system that's open. It's  
21 available to industry. It's available to academics.  
22 It's available to public health agencies. There is

1 no restriction on who can look at this information.

2 Next slide.

3 We always get asked, so the sequence  
4 information is only valuable with the metadata. I  
5 think you heard a little bit of that discussion this  
6 morning. This is the metadata that's provided in the  
7 case of a food/environmental or a clinical on the  
8 right.

9 Basically the information that's circled is  
10 really the information I wanted to highlight today,  
11 and that's a specific FDA number. So if we did an  
12 inspection of your facility, you would get a report  
13 at some point that would have this number associated  
14 with the samples that were collected. You could, and  
15 I'm hoping that the NCBI -- Bill Klimke's in the  
16 back. I see him there. So he will tell you about  
17 how to do this, but you can type this number into the  
18 search engine that NCBI has built, and you can see  
19 where your isolate lies in the tree. And you can  
20 then download those sequences and you'll have  
21 information specific to your isolate which we think  
22 is great because you can then verify what we're

1 doing.

2           In the case of other information, the  
3 geographic information, we list it only by country.  
4 We also list the type of food or if it's an  
5 environmental facility, it would just say  
6 environmental, that it was collected by FDA. If  
7 there's PFGE information, that's also provided, and  
8 then who it was provided by. So if it was provided  
9 by FDA, there's contact information there. So if you  
10 are doing research and you don't have access to the  
11 full list of metadata, obviously we know what state  
12 and where the place was that we collected the  
13 information. That we could provide to folks.

14           Then on the clinical side, very similarly  
15 it lists if someone got sick in the U.S. is about  
16 what you have. We don't have any other information.

17           We spend a lot of time working with various  
18 people in industry, academia, and public health and  
19 other folks throughout the world trying to get this  
20 very limited meta dataset and although as it was  
21 correctly pointed out, it is a very limited amount of  
22 information. Some very clever people have figured



1 out how to take advantage of this information, and  
2 while not having completeness, be able to develop a  
3 lot of very cool models on predicting disease and  
4 predicting risk and by simply just knowing that  
5 someone got sick and it was from this organism, and  
6 maybe from this particular food or this particular  
7 area. And then they can ask a lot of questions like,  
8 well, I only need to know this little piece or that  
9 little piece.

10 Next slide, please.

11 So this didn't come out quite right, but  
12 the point is that you have one data record, you can  
13 get lots of things from it. The most interesting  
14 thing to me is that we're kind of myopic in our  
15 approach. We're looking for very specific things,  
16 but there's lot of other information there that we  
17 haven't even begun to mine. So we really are trying  
18 to encourage people to put this information up there  
19 and be able to use it in ways that we as regulatory  
20 agencies may never have thought of.

21 Next slide, please.

22 So again these are questions we get all the

1 time. Yes, it is more discriminatory than PFGE.  
2 Like all living organisms, there is stress and  
3 pressure when they adapt to their environment. This  
4 has been known for a long time. They mutate more  
5 rapidly and so we get a clue to the geographic origin  
6 based on the stress and the pressure and the  
7 environmental response that the organism sees. And  
8 that's what's really so important about this  
9 particular technology and the use of it.

10 In clinical applications, it's slow. I  
11 mean you're really looking at days to get the  
12 analysis done, and so it's probably not acceptable  
13 for that, but for us, in the regulatory community and  
14 the food community, it's really very, very good.

15 The point that I'd like to make on this  
16 slide, in particular, is that when we have some  
17 environmental information, we know where to look as  
18 opposed to not looking. We really spend a lot of  
19 resources, very ineffectively, by not even knowing  
20 where to look, not even, for example, knowing, you  
21 know, should we be looking at imports from Southeast  
22 Asia or should we be looking at domestic samples. So

1 all of this really helps and aides in a much faster  
2 response.

3           Next slide, please.

4           What data analysis tools should be used and  
5 why? Another question we routinely get. Well, that  
6 depends. It depends on who you are and what you're  
7 trying to accomplish.

8           So if you're -- I don't want to speak for  
9 CDC, but I'll just, having listened to their  
10 presentations, they're interested in a very fast,  
11 kind of easy way of pushbutton looking relationships,  
12 which gets you to a certain point, and you're looking  
13 at sort of a limited dataset.

14           Because we're a regulatory agency, and we  
15 know that there are consequences to our results,  
16 legal and other, we use the gold standard SNP method.  
17 Both methods for the most part give very similar  
18 results and again it depends on what you're approach  
19 is.

20           As John pointed out, there's no single  
21 threshold. Generally if it's less than 20 SNPs away,  
22 we'd be taking a closer look at it. That doesn't

1 mean that we do anything with it. It just means we  
2 take a closer look.

3 All of our tools are validated. They're  
4 all on the website. You can look at the validation  
5 documents. They're all validated according to  
6 software standards, well accepted standards, and so  
7 you can go and you can pull down the software, you  
8 can look at it, you can look at all the validation  
9 statistics and all of those things with it.

10 What software should I use? Well, that  
11 depends on your level of comfort. There are  
12 commercial versions like BioNumerics, which is the  
13 standard that CDC and PulseNet uses. We've heard  
14 from numerous people that that's too expensive, that  
15 they don't use it that much, that it's too difficult,  
16 on and on and on. So we have been working for more  
17 than a year now, and I think we are going to have a  
18 Galaxy version, a free publicly available Galaxy  
19 version that uses our SNP pipeline that will work in  
20 an Amazon cloud service. Stay tuned for that. I  
21 thought it would be more advanced than it is, but  
22 it's a little behind because of security concerns but

1 we are getting very close to having something that  
2 you can plug your own information into and get the  
3 same answer that we get we hope.

4           And then how do I access the data and do  
5 analysis on the NCBI site? There's a really simple  
6 website. It's basically the NCBI site/pathogens, and  
7 that will take you right to the viewer and you can do  
8 whatever you like with it. Very simple.

9           Next slide.

10           What happens with a match? Reanalysis of  
11 the data happens first. So if we thought we saw a  
12 match, we would take a much closer look at the data.  
13 We would pull down the information. We don't take  
14 any regulatory action based purely on the match. We  
15 would resample, we would reinspect, and then  
16 depending on those reinspection results, a number of  
17 things can happen because there's two possibilities.

18           We can take and we do take regulatory  
19 actions routinely on samples where there's no  
20 epidemiological evidence. If you have pathogens in  
21 your facility, they're in the food or they're in zone  
22 1, zone 2, you're in violation of the law of having

1 unsanitary conditions and so you're going to hear  
2 from us and trust us, that is really the best place  
3 for you to be because the next step is not where you  
4 want to be, where you've made people sick, and we do  
5 rely on epidemiologic evidence because when you get  
6 into civil lawsuits, you get into legal issues that  
7 you really don't want to be involved in.

8           If you're in that first half, unsanitary  
9 conditions, we can at least work with you and we  
10 would like to work with you on solving that problem  
11 and helping you understand the information.

12           Remember, we've taken maybe dozens,  
13 hundreds of swabs of your facility, and we've done  
14 the sequencing free, you're not paying for it, and so  
15 we can provide all that information as well as where  
16 they came from and what to do about it.

17           Next.

18           We do this because we know that there's a  
19 supply chain and that if we found a sample, perhaps  
20 in the processing facility, it may have actually come  
21 from an ingredient or from the manufacturing  
22 facility. There's upstream and downstream

1 information, and so we really do not want to take  
2 regulatory action unnecessarily against someone who's  
3 not involved.

4           There's also a movement of processing  
5 equipment. One company goes out of business. They  
6 sell their equipment to another company. The  
7 material is contaminated and so the contamination  
8 moves from one facility to the next.

9           Next slide.

10           A note about clinical matches and  
11 epidemiology. When we go into facilities, what we  
12 see is this. We see an incoming stream that's  
13 relatively free of contamination and an outgoing  
14 stream that has loads of contamination. And if you  
15 think about it, it makes senses. If you had, you  
16 know, a million bags of flour or a million bags of  
17 lettuce, and you're producing that every day and only  
18 one of those bags were contaminated, we wouldn't  
19 detect any illness.

20           When we start detecting illness, you've got  
21 multipliers of 40 or 10 or 5 and you actually have to  
22 hit those susceptible populations. So when we go in,

1 we usually find very significant contamination.

2 Next.

3 This is a slide that our industry has been  
4 begging us for and asking us for, you know, what are  
5 the implications of whole genome sequencing? What do  
6 you see? How does this work? And I would just like  
7 to caveat a lot of this because we don't keep our  
8 information this way necessarily and exactly how all  
9 of these things would have worked, but the total  
10 number of facilities inspected where we would be  
11 looking at high risk facilities that we were either  
12 in before or were indicated as being involved in some  
13 kind of outbreak, or we've had matches in the past,  
14 so we've done an inspection, we've done another  
15 inspection and we see some relationship there.

16 There were 167 requests for additional  
17 information done by our bioanalytics groups to look  
18 at the information, and then of those 600 inspections  
19 and 167, it was actually a very small number of  
20 regulatory actions with some of them being more  
21 significant than others. The regulatory meeting were  
22 generally the most common occurrence. There were



1 three injunctions and one mandatory recall, and one  
2 suspension with is very severe. But, for the most  
3 part, there was a small number of regulatory actions,  
4 not a huge number based on this technology, and I  
5 think in many cases we would have arrived at some of  
6 these situations without sequencing at all.

7           Next slide.

8           Should we start using sequencing? How much  
9 time do I have? One or two minutes.

10           Contribute sequences to GenomeTrakr  
11 database, if you're industry I tell you to please  
12 contact your lawyers and then get a second opinion.  
13 There are consequences to doing this which you may  
14 choose not to take. I should say there's risks to  
15 doing this, not necessarily consequences.

16           If you're doing routine environmental and  
17 high volume product sampling, you don't want this  
18 technology. It's too costly.

19           If you're doing supply chain management and  
20 say PCR doesn't work for you any more, you can't  
21 figure out what the source of the problem is, you  
22 probably want to look at this because it is an

1 invaluable tool in tracing where things came from and  
2 how they there.

3           If we've been in your facility and you've  
4 got a positive, you may as well kind of throw your  
5 hands in the air and just engage in this because when  
6 we look at civil litigation and other litigation that  
7 the Justice Department may take, it's when you knew  
8 about it and what you did about it. So if 6 months  
9 ago we were in your facility, we inspected, you have  
10 this and you haven't done anything and you have no  
11 documentation, I'm sorry, I'm going a little over,  
12 you know, you can't hide it any more. We already  
13 know you have *Listeria mono* or *Salmonella* there, and  
14 it's really what you did about it during that time  
15 frame. And the next slide which I hope is the last  
16 one.

17           The only thing I'd like to say about this  
18 is that there's more to this technology than simply  
19 doing outbreak detection, particularly if you're  
20 interested in the effectiveness of sanitizers in your  
21 facility, if you're interested in supply chain  
22 management, it really is a technology that can get

1 you to places that you couldn't get by any other  
2 technology.

3           And with that, I think I'm done. Yes.  
4 Thank you.

5           DR. GOLDMAN: Well, good morning. I'm  
6 David Goldman again for those who weren't here at the  
7 very beginning of this meeting, and I'm going to  
8 represent the FSIS exploration of this technology in  
9 our regulatory programs. And what I'd like to do is  
10 to tell you a little bit of a story, show you a  
11 little bit of the data, and I would say at the  
12 outset, that the way to put this is we've identified  
13 many more questions than answers, and I think you've  
14 heard that theme already, and I think you'll see it  
15 in the slides here as well.

16           So I want to acknowledge my partner, Uday  
17 Dessai, that you've heard from -- he's moderated some  
18 of the sessions already -- in a team of people here  
19 at FSIS who have really in less than 5 years put  
20 together quite an enterprise in terms of whole genome  
21 sequencing and its application.

22           Here's the outline, and I won't spend any

1 time here. We're going to cover these things in  
2 fairly short order and hopefully leave a few minutes  
3 for questions of us from a federal perspective.

4           So just briefly about our Agency for those  
5 who are new and don't know about FSIS. We are an  
6 inspector intensive operation. We have inspectors in  
7 every plant that operates in the U.S., and you can  
8 see the map that sort of depicts the establishments.  
9 We have about 6,000 establishments that we regulate  
10 and inspect each day, about 7,000 plus inspection  
11 personnel. And one of the consequences of that work  
12 and the way we operate is to do lots of sampling of  
13 food products as part of our verification activities,  
14 and you can see here, we generate almost 10,000  
15 bacterial isolates per year, and so that's sort of  
16 the foundation of the work that we're able to do with  
17 whole genome sequencing.

18           You can see on the right our authorities  
19 and just to be clear here, we regulate meat, poultry,  
20 and processed egg products, and FDA regulates all the  
21 rest.

22           So we made a very purposeful move towards

1 whole genome sequencing, and I'm not going to go  
2 through all the details here, but we started several  
3 years ago to begin to plan for the use of this new  
4 technology in our regulatory programs. It's become a  
5 prominent part of our Agency's strategic plan which  
6 was issued last year, as well as our annual plans,  
7 all of which are available to you if you're  
8 interested. And again, this is something that has  
9 taken some special attention and time to do. It was  
10 a high priority in terms in both the budget and the  
11 resources that needed to be realigned within the  
12 Agency.

13           Just a brief mention here about our  
14 workflow: You've heard about others, and their  
15 workflow and this slide is maybe a little hard to  
16 read but it compared the workflow using pulsed field  
17 analysis against whole genome sequencing, and the  
18 bottom line is to just show you that whole genome  
19 sequencing really takes only a couple of extra days  
20 for us to do it to get that sequence uploaded to  
21 NCBI. So the turnaround time is pretty good. It's a  
22 one day extra for *E. coli* or STECs, but really in

1 less than a week, we can have our data uploaded.

2           Here's a snapshot. We use this snapshot at  
3 various venues to kind of mark our progress in terms  
4 of bringing up the technology and adding sequences to  
5 the NCBI database, and you can see the bar graph on  
6 the upper left, shows you the number of isolates that  
7 are sequenced per quarter and our steady state, our  
8 goal was going to be about 2,500 isolates per  
9 quarter. So you can see that toward the end of the  
10 last fiscal year, we approached that, we exceeded  
11 2,000 isolates per quarter.

12           On the right you can see sort of our  
13 history starting in July, and I should point out here  
14 that we really depended a lot on FDA in the initial  
15 phases to help us get our program up and running.

16           On the bottom, the table just shows you the  
17 various sources of the isolates and by pathogen and  
18 you can see that in the history of our efforts, we've  
19 now uploaded more than 10,000 isolates.

20           I did want to provide a little perspective.  
21 You can see the *Salmonella*, just over 6200 isolates  
22 have been sequenced and uploaded, and to give you

1 some context, there are about 110,000 *Salmonella*  
2 sequences in the NCBI database. So that will give  
3 you a little bit of background there.

4           Okay. Just another fact sheet really about  
5 our capacity building. We have three regulatory  
6 labs. We now have sequencers in all of those three  
7 labs and are producing sequences there.

8           This past year, which just ended in  
9 September for us, we had sequenced and uploaded over  
10 7,000 isolates, remember our goal being about 10,000  
11 a year. We do intend now in this current fiscal  
12 year, starting in October to sequence every single  
13 pathogen isolate and some proportion of the  
14 indicators which are isolated during our NARMS work,  
15 and I'll come back to that in just a minute.

16           And on the right you can see the metadata  
17 that companies each upload to NCBI, and it's somewhat  
18 similar to what you've seen from the FDA, but I do  
19 want to point out that the product and source, of  
20 course, is really important but we just upload the  
21 year that the sample was collected and the state in  
22 which it was collected. There's no more specific

1 information than that that's uploaded as part of the  
2 metadata.

3           So this describes sort of the rest of the  
4 talk which is kind of where we're going with our  
5 analysis, and just outlines some of the different  
6 types of applications that we've begun to explore  
7 within FSIS, and I'll cover each of these in turn.

8           So, first, I want to show you two slides  
9 that are related. We used whole genome sequencing in  
10 an outbreak investigation in retrospect. So I wanted  
11 to point that out, and let me see if I can orient you  
12 to the important pieces here.

13           So this is an outbreak. This is actually  
14 two outbreaks that occurred in the same facility,  
15 separated by just about a calendar year, and what you  
16 can see here is most of these isolates represent  
17 clinical and food or environmental isolates from the  
18 second of the two outbreaks, but there is one  
19 clinical isolate from the first year's outbreak and  
20 you can see how tightly related those isolates are.  
21 So this is the first of two slides related to the  
22 same outbreak.



1           Now, here's another analysis which looks at  
2 the simply commodity isolates. There are no clinical  
3 isolates in this depiction here, and what you can see  
4 also, at the point here, is that there are some  
5 isolates here that represent at least in one case a  
6 different PFGE pattern and 2 different years' worth  
7 and see again very tightly clustered.

8           The next area that we've started to apply  
9 this technology is in exploration of harborage in  
10 ready-to-eat facilities. We had been using pulsed  
11 field analysis to help us understand the extent to  
12 which *Listeria* might be harbored in a plant from one  
13 year to the next or from one sampling event to the  
14 next, and we've now overlaid that work with the use  
15 of whole genome sequencing to help us understand the  
16 extent to which there could be harborage in plants.  
17 We know, of course, the potential is there.

18           I also want to point out, as was mentioned  
19 earlier, there are about 1,000 ready-to-eat  
20 facilities that are dually inspected by both FDA and  
21 FSIS. We call them dual jurisdiction establishments.  
22 So this information is obviously of interest to both

1 FSIS and to FDA.

2           Now, in this slide, there are two points I  
3 want to make. One is that this is from one plant and  
4 I want you to focus on this cluster. I know it's  
5 hard to read, but the red triangle here represents a  
6 sampling event in 2012, and all the rest of these  
7 isolates were from the environment or food contact  
8 surfaces in 2015. So you can see that spread over a  
9 3 year period, here are isolates that are very  
10 closely related. I think it says 0 to 3 SNPs there.

11           The bottom part of the slide just has  
12 reassured us that in most instances, there is  
13 widespread agreement in our analysis between what  
14 pulsed field analysis and whole genome sequencing is  
15 telling us, whether using SNPs or whole genome MLST.  
16 There are some differences however. There are some  
17 instances that we've encountered in which different  
18 PFGE patterns can be aggregated through a common  
19 genome sequence.

20           Okay. Now, this slide, we're turning our  
21 attention now to *Salmonella* from chickens. So it's  
22 specifically about those, the upper part anyway, is

1 about that in particular. And we looked at the  
2 serotypes that we commonly find in chicken products,  
3 in our chicken sampling, and you can see the  
4 serotypes which are familiar to all of us and the  
5 numbers of isolates from those various serotypes.  
6 And then we looked at this number and saw how closely  
7 they aggregated through a SNP analysis, and then we  
8 compared them to clinical isolates.

9           And so if you looked at, going out to 20  
10 SNP difference, you can see the agreement or the  
11 sameness of those clinical isolates to the product  
12 isolates, and you'll note a couple of things that are  
13 of interest, and then we did the same analysis  
14 including just 10 SNPs difference.

15           If you look at Kentucky, which is commonly  
16 found in chicken, there's 0% of both whether you go  
17 to 10 or 20 SNPs of relatedness to clinical isolates.  
18 And you can see some of the same thing with  
19 Typhimurium, and the story with Typhimurium is  
20 different because we find Typhimurium in every single  
21 one of the products that we regulate. So that's the  
22 explanation there, but this sort of analysis can help

1 us as we move forward.

2           On the bottom right, you can see this has  
3 to do with resistance. We did an analysis of 1700  
4 plus *Salmonella* isolates and about half of them were  
5 pan-susceptible and had no resistance genes and then  
6 you can see a list of those commonly isolated or  
7 detected genes in those samples.

8           There are several slides here I'm just  
9 going to kind of run through quickly on geographic  
10 distribution. We have undertaken some analyses to  
11 see whether we can determine, and these are all  
12 related to *Salmonella* in cattle. So we looked at  
13 Dublin, Montevideo, and a couple of other serotypes  
14 you'll see in subsequent slides, and I'm just going  
15 to kind of scroll through these.

16           Here's some more on Montevideo -- let me  
17 just go back one second. The story here for both  
18 Dublin and Montevideo is that these are highly  
19 diverse serotypes which we sort of understood  
20 already, and there doesn't seem to be a lot of  
21 geographical clustering.

22           This slide's just a snapshot of the extent

1 of the diversity. So you've got the same PFGE  
2 pattern depicted in this box here and yet there's  
3 significant diversity within *Salmonella* Montevideo.

4           Here's a look at *Salmonella* Newport which  
5 is another cattle adapted serotype, and again lots of  
6 diversity and very little geographical clustering.  
7 Now, of course, for us, we take the samples at the  
8 slaughter plant, and cattle movement is a big factor.  
9 So that's one of the confounders in trying to do this  
10 sort of analysis. The cattle may have come from lots  
11 of different regions and that may be why we are not  
12 finding the geographical relationships that we  
13 expected.

14           I'm going to end up by talking just briefly  
15 about some of our work with NARMS. Now, I think you  
16 heard the last 2 days, prior to this meeting, we had  
17 a NARMS meeting. It was a public meeting, and there  
18 was a lot of discussion about the use of whole genome  
19 sequencing in the resistance context. You've already  
20 heard a little bit about this at this meeting today,  
21 and this is our analysis looking at cecal isolates of  
22 *Salmonella* for one calendar year, and the point of

1 this slide is that this is sort of reaffirmed what  
2 has been shown now for several years, in that there  
3 is a high concordance between the genetic elements  
4 found using ResFinder or other tools to determine the  
5 resistance genes and the AST that's done more  
6 traditionally.

7           We've also been able, within the context of  
8 NARMS, to do some work to identify specific  
9 resistance genes and you can see the list of the  
10 genes that have been detected as a result of the  
11 collaborations across NARMS. Certainly the CTX-M-65  
12 is of great interest because it confers ESBL as well  
13 as resistance to some other antimicrobials as well.

14           What we've tried to do in our Agency is to  
15 let the producer community, the slaughterhouse  
16 operators know when we find these genes so that they  
17 can be aware that we have been doing this work as  
18 well as to let them know that they may want to  
19 examine their food safety systems.

20           Two other quick applications: We have  
21 begun looking at *Salmonella* pathogenicity islands.  
22 We are very interested in the Agency in virulence.

1 So we will continue to be interested in trying to  
2 isolate virulence determinants through the  
3 application of this technology. And you can see on  
4 the left, I'm not going to say any more about it, but  
5 the *Salmonella* pathogenicity islands have been well  
6 described, and we can use this approach to helping us  
7 understand the isolates that we obtain in our  
8 regulatory programs.

9           On the right side, there has previously  
10 been described heat resistance genes in *Salmonella*,  
11 and then when we use that knowledge and looked at our  
12 beef derived isolates of *Salmonella*, we determined  
13 that there were no heat resistance islands in those  
14 samples, but again that's an application of a  
15 specific approach we can take looking at the isolates  
16 that we generate.

17           I'm just going to end up here with a couple  
18 of slides. Here's sort of a graphical depiction of  
19 where we are and where we hope to go with the  
20 application of this new technology. So on the left,  
21 you know, we've spent a lot of time learning about  
22 the power of whole genome sequencing. We've had

1 relatively modest goals, both in terms of developing  
2 our capacity, whether it's the machines themselves or  
3 our human capacity, our resources, in terms of  
4 personnel, and then beginning to apply the technology  
5 as I've just described it.

6           You know, as we move toward 2030, and we're  
7 always focused on the healthy people goals, we hope  
8 to do more of the applications, some of which are  
9 described here and have greater goals, all of which  
10 we hope will help us to decrease foodborne illness  
11 related to the products that we regulate.

12           This is just a reminder that we're part of  
13 a big collaboration. I think that almost all of  
14 these organizations are represented in the room here,  
15 and we've attended all of these meetings that you see  
16 outlined here. I think it's important to say that a  
17 few times during that meeting, to reassure everyone  
18 in the room that we are in this together, we're  
19 learning together, and we need to have a shared  
20 understanding about the meaning of the work that we  
21 do.

22           And there's just one summary which I've



1 covered pretty much. We have the capacity. We now  
2 have the resources in terms of the staff we have  
3 assembled. We'll continue to hire out a little bit  
4 more in terms of folks who can do both  
5 characterization work in the lab as well as  
6 bioinformatics, but we're approaching that capacity.

7           And then, you know, the applications within  
8 NARMS are important, and there's a lot of good work  
9 being done and will continue to be done there as we  
10 are contributors to NARMS.

11           And then finally, just to reiterate, as has  
12 been said many times here in this meeting already,  
13 this is a tool. It needs to be placed in the context  
14 of our other findings and certainly for FSIS,  
15 traceback is a real key finding. We haven't talked  
16 much about traceback, but the ability to traceback to  
17 the plant that has produced product that may be  
18 contaminated and cause illness is a real key factor  
19 in this, and so we want to use whole genome  
20 sequencing in that context. And with that, I will  
21 end, and we'll have a few minutes for questions.

22           Thank you.

1 DR. EVANS: Thank you, Dr. Besser and  
2 Musser and Goldman, and I will just ask David, if you  
3 could come back to the stage, and we're going to have  
4 -- this is not on our agenda, but we're going to have  
5 5 to 10 minutes where we give opportunity for the  
6 folks in the room and folks on the webinar to ask  
7 questions from our first three presenters. So did we  
8 have any questions on the web?

9 DR. ABLEY: Okay. We do have a question  
10 from one of our online. To what extent will random  
11 sampling of retail product be pulled, tested and  
12 positive isolates ran on WGS? Will a risk-based  
13 approach be taken in recalls if nothing relates to  
14 illness, ready-to-eat versus non-ready-to-eat?

15 DR. MUSSER: So as a matter of routine, we  
16 don't like collecting retail samples. We sometimes  
17 do surveys to look at things other than microbial  
18 contamination but we don't generally pull retail  
19 samples because it's extremely difficult to figure  
20 out where the contamination may have come from based  
21 on retail sample analysis. It doesn't mean it's not  
22 done. It's just not a good way of finding more

1 problems occur. And what was the second part? Read  
2 the second part again.

3 DR. ABLEY: Will risk-based approach be  
4 taken in recalls if nothing relates to illness,  
5 ready-to-eat versus non-ready-to-eat?

6 DR. MUSSER: I'll take a stab at it, and  
7 then I'll let David do it. I think, yes, typically  
8 in non-ready-to-eat, there's less likelihood of  
9 illness and we tend to focus our efforts on high risk  
10 commodities. So we don't typically test lower risk  
11 commodities. We still do, but not as much, and so  
12 our risk analysis and our risk-based inspections are  
13 just that. So it would be on the highest risk which  
14 would be ready-to-eat and the non-ready-to-eat would  
15 be less tested.

16 DR. GOLDMAN: I would just add that at  
17 FSIS, we have both risk-based sampling as well as  
18 randomized sampling, and sometimes they interact  
19 within the same sampling program, looking at the same  
20 type of product. So our sampling is principally for  
21 verification of a good safety system. So that's the  
22 biggest driver here, but within that, again and

1 *Listeria* would be one example, ready-to-eat products,  
2 we do do risk-based sampling within that based on  
3 some algorithms we've develops.

4           And so in terms of the recalls, you know,  
5 once there are illnesses, it doesn't matter whether  
6 it's from a ready-to-eat product or a raw product.  
7 We would still proceed accordingly.

8           DR. EVANS: Any question in the room?

9           MR. McDERMOTT: Hi, Pat McDermott from FDA,  
10 Center for Veterinary Medicine. I have a question  
11 for Dr. Besser about the confidential information. I  
12 understand some states don't allow isolate level  
13 information to be put into a public domain. Will  
14 that apply to whole genome sequence data from  
15 PulseNet as well? In other words, will some of the  
16 whole genome sequence data remain behind the CDC  
17 firewall or can it be de-identified enough that all  
18 the states will allow that isolate level information?

19           DR. BESSER: Good question. Thanks, Pat.  
20 It's my understanding that all states have signed  
21 onto the MOU with whole genome sequencing, that the  
22 data that we request through the PulseNet database is

1 uploaded into the PulseNet data which is kept behind  
2 the state firewall. Some of that goes to CDC level,  
3 and then the minimal dataset that the agencies have  
4 all agreed upon gets uploaded to NCBI.

5           So I'm not sure that any of the states have  
6 declined to put in the information. I could be  
7 wrong, but I think they've all agreed to it.

8           DR. EVANS: Are there any other questions  
9 in the room for our first set of speakers?

10           (No response.)

11           DR. EVANS: So we're going to break for  
12 lunch, and just to remind you, the lunchroom, as you  
13 leave this room and take a right, go down to the  
14 third wing and take a right and the lunchroom is  
15 halfway down that corridor. We'll be back in this  
16 room at 1:30 p.m. with three more presentations on  
17 Federal/State Collaboration. Thank you.

18           (Whereupon, at 12:06 p.m., a luncheon  
19 recess was taken.)

20

21

22



1 and work on these problems, standardization, common  
2 protocols, all to improve the quality of the data  
3 that's going into the common databases at NCBI.

4           So first we're going to hear from Dr. Dave  
5 Boxrud from Minnesota, and Dr. Patrick McDermott from  
6 FDA's Center for Veterinary Medicine, and lastly from  
7 Dr. Bill Klimke at the National Center for  
8 Biotechnology Information.

9           And again, we're going to ask for folks to  
10 hold their questions. We are going to have a  
11 question and answer period at the end of the day  
12 after our break.

13           So Dave.

14           MR. BOXRUD: Thank you.

15           DR. EVANS: Thank you.

16           MR. BOXRUD: Hi. I'm going to talk about  
17 whole genome sequencing at the Minnesota Department  
18 of Health. I was asked to talk about how sequencing  
19 is being done at the state public health laboratories  
20 which is a pretty difficult task because there's 50  
21 states and then many local health departments. So  
22 I'm going to give a little bit of background about

1 how public health laboratories nationally are  
2 incorporating whole genome sequencing but then really  
3 kind of focus on how we're doing it in Minnesota as  
4 just one example.

5           So this talk is going to give us a status  
6 update on whole genome sequencing at public health  
7 laboratories, talk about the role of whole genome  
8 sequencing at public health labs, how we communicate  
9 whole genome sequencing inform at public health labs,  
10 a little bit on how we've evaluated whole genome  
11 sequencing and so we can understand and interpret the  
12 data at public health labs and then lastly, I'll  
13 close with an example of the utility of sequencing.

14           So John Besser showed this slide a little  
15 bit earlier, but I wanted to point out that in 2014,  
16 Association of Public Health Laboratories did a  
17 survey and at that time, 21 state public health labs  
18 had a sequencer in house. Now, we have 43 labs and  
19 37 states are certified which means that they not  
20 only have the sequencer and have the training but  
21 they've been able to show proficiency with their  
22 sequencing.



1 All states in the country have been  
2 financed for a sequencer, but by the end of 2017 or  
3 early 2018, all states will have a sequencer in  
4 house.

5 But there are still a number of issues that  
6 state public health labs are dealing with. Right now  
7 we're in this incredible transition time, and we're  
8 continuing to do pulsed-field gel electrophoresis  
9 which is a conventional or traditional subtyping  
10 method while we're transitioning into this new  
11 subtyping method. It puts a tremendous resource  
12 crunch on the laboratories, probably our main  
13 challenge right now in public health labs.

14 Many public health labs have IT issues,  
15 either with storage or their IT departments don't  
16 allow them to use certain types of software. Some of  
17 that has been largely resolved, but there are still  
18 some issues in some states.

19 Training on both the wet lab side and on  
20 the bioinformatics side, using the data, is still a  
21 little bit of an issue, but it's been resolved in a  
22 lot of areas.

1           Bioinformatics resources is a little bit of  
2 a challenge. John talked about how PulseNet is doing  
3 their analysis and that's going to be very helpful  
4 for a lot of public health laboratories, but right  
5 now we're still in a little bit of a waiting mode.

6           And lastly, ordering reagents is a bit of  
7 an issue for public health labs. Reagents are very  
8 expensive and sometimes in certain states, it's  
9 really hard to put five or six figure orders in and  
10 that causes a lot of problems with our ordering  
11 issues.

12           So I just want to start with the very basis  
13 of foodborne disease testing. This is traditionally  
14 how we've done our testing, and I really think that  
15 this is the linchpin of how we identify outbreaks.

16           On the laboratory side, we get human  
17 samples in. We analyze the samples and subtype them,  
18 subtype the bacteria in them, the *Salmonella* or the  
19 *E. coli* and then report our results to epidemiology,  
20 really focusing on clusters so that they can  
21 investigate them.

22           On our epi side, we interview cases and

1 then investigate clusters and lead investigations.

2           And I think what is really important is  
3 this process, this work on both sides, but also the  
4 relationship between the two, there should be a sense  
5 of urgency with all of this work and with how this  
6 data is being communicated back and forth.

7           But, of course, we have a lot of other  
8 partners that we work with regularly with other state  
9 and local health departments, USDA, FDA, CDC, and  
10 PulseNet is certainly been a huge partner, and also  
11 with our environmental health team in the state for  
12 investigations.

13           Just to talk a little bit more about  
14 epidemiology information, these interviews that they  
15 do, I'm not going to go into a lot of details, but  
16 these are really, really important to try to get as  
17 good of exposure data as possible and to get a really  
18 good interview takes a lot of time. You're really  
19 asking a lot of in depth questions of what people  
20 ate, where they ate, specific brand names, when they  
21 ate certain foods, restaurants they went to. Each  
22 interview takes anywhere between 30 minutes to 60

1 minutes.

2           And one of the ways in Minnesota that we  
3 have for routinely sharing subtyping information is  
4 we created an automatic daily report through our  
5 laboratory investigation management system, or LIMS  
6 system, and it's broken up into a few different  
7 pieces but the idea is that we're able to have a  
8 standardized methodical way of getting information  
9 from our laboratory to our epidemiologists.

10           And I'm just going to talk a little bit  
11 about a couple of these parts. The first part is  
12 just essentially giving a report of what was seen the  
13 previous day, very simple information. Some  
14 background demographic information on the cases, but  
15 also the subtype information. In this case, it's the  
16 PFGE pattern. So every PFGE pattern gets a name so  
17 that we can communicate that with our  
18 epidemiologists.

19           The next part that is really key is that we  
20 take a look back at the subtypes for those serotypes  
21 and say, have we seen that for the last 30 days, and  
22 that information is also sent to the epidemiologists.

1 So in this case, we had seven isolates that had the  
2 same PFGE pattern that were seen in the last 30 days.  
3 Obviously, this is a cluster that our epidemiologists  
4 would be very concerned about.

5           So now that we're going to a sequencing  
6 based method, a lot of things stay the same, and  
7 there are a few differences. On the public health  
8 lab side, we're still looking at human samples, but  
9 instead of doing pulsed field, we would be sequencing  
10 those samples, but we still need to report that  
11 information to our epidemiologists and from the epi  
12 side, they're still doing their normal stuff.  
13 They're interviewing cases and investigating clusters  
14 and looking for outbreaks.

15           We do have some more partners, actually the  
16 same partners initially, but we have another partner  
17 which is NCBI, which is where we have a sequence data  
18 repository. I would also be remiss to say that FDA  
19 has been a much bigger partner with us, with  
20 sequencing. We are a GenomeTrakr lab, and we are  
21 sharing data back and forth with FDA and they've been  
22 a tremendous partner for us.

1           So now that we're working with whole genome  
2 sequencing data, the communication is a little bit  
3 more challenging. Some of the previous speakers  
4 talked about nomenclature and how they're naming  
5 their patterns. So far to the public health labs, we  
6 don't have that available. So we created a number of  
7 spreadsheets to give this information back and forth  
8 and essentially the most important thing is that we  
9 have a cluster ID. We tell our epidemiologists what  
10 is a cluster and then they investigate those cases,  
11 and we also send them an email that lays out what a  
12 cluster is and gives a little bit of background  
13 information. In this case, it's two cases, the same  
14 PFGE pattern, gives some information that's very easy  
15 for them to understand and to create an action on  
16 this.

17           The negative side of this is that this is  
18 not going to be sustainable for a large number of  
19 organisms. It works for the short term, but it's not  
20 possible to work in the long term.

21           So now we want to focus about how we  
22 determine what is a cluster and how we have used that

1 information that we learned in the past going  
2 forward. So I'm going to talk a little bit about  
3 *Salmonella* Enteritidis, a study that we did with  
4 *Salmonella* Enteritidis. It's the second most  
5 *Salmonella* serotype in Minnesota but it's very  
6 clonal.

7           Here's the major pulsed field patterns.  
8 There's about five patterns that make up about 80% of  
9 the isolates in Minnesota. So pulsed field really  
10 does not do a good enough job of providing diversity  
11 on *Salmonella* Enteritidis.

12           So we decided to do a retrospective study  
13 on SE, and we wanted to look at stability, how stable  
14 are the sequences over time, typability, are we  
15 always able to get sequence data from every isolate,  
16 discriminatory power and also epidemiologic  
17 concordance. Epi concordance is essentially defined  
18 as if it's an outbreak, they should look very  
19 similar. If they're not related, they should look  
20 very different.

21           So our laboratory and epidemiologists met  
22 and they came up with a study set, and we used some

1 very well characterized isolates. We look at  
2 isolates from 7 outbreaks but also looked at 22  
3 sporadic isolates. We also looked at some in vivo  
4 isolates. So multiple isolates from one person over  
5 time. And we worked with the New York Department of  
6 Health for the analysis for this.

7           And what we saw was that within an  
8 outbreak, there were very few SNPs. There was a  
9 maximum of three SNPs within an outbreak, and when  
10 there were sporadic SNPs, they looked very different  
11 from the outbreak isolates, even when they were the  
12 most common key of key patterns. So that provided us  
13 some information on how to interpret this data going  
14 forward.

15           Our conclusions were that the sequences are  
16 stable within a person over time. All isolates were  
17 able to be typed. There was a lot of diversity, and  
18 there was really good epidemiologic concordance.  
19 Outbreak isolates looked the same. Isolates that  
20 were not epidemiologically related looked very  
21 different.

22           So we used this data to do a prospective



1 study in which case we looked at whole genome  
2 sequencing and pulsed-field gel electrophoresis  
3 clusters in real time and tried to identify the  
4 source of them. And what we saw was that with  
5 sporadic isolates, by whole genome sequencing, they  
6 looked very different. There's an average of 93 SNPs  
7 different between the most common PFGE pattern that  
8 were sporadic, and going back to what we saw from the  
9 outbreak isolates, there was 0 to 3 SNPs difference.

10 So from a sequencing perspective, it was  
11 quite easy to see what is like to be related and  
12 where we should focus our investigation.

13 And also the number of clusters that we saw  
14 with whole genome sequence compared to pulsed field  
15 went up dramatically but the number of isolates  
16 within each cluster went down dramatically.

17 So I think this provides us with a great  
18 opportunity to identify outbreaks earlier with fewer  
19 cases.

20 And, lastly, I just want to talk about how  
21 we use this data going forward, one great example.

22 So this is the most common PFG pattern in

1 Minnesota in August and September of 2014, with each  
2 box representing one case of the most common PFG  
3 pattern. So there are 19 total cases. So by PFG,  
4 it's really hard to understand what is going on. We  
5 just know that there's a lot of cases.

6           When we did whole genome sequencing, we  
7 found that there were eight of the isolates that were  
8 0 SNPs from each other. So we investigated those to  
9 try to understand what the relationship between those  
10 cases were, and so our epidemiologists did a great  
11 job and they would interview almost all of the cases  
12 and found that six of the eight, all ate a frozen  
13 chicken Kiev product. One of those eight was a  
14 secondary case and another one of the eight may have  
15 eaten a frozen chicken Kiev product. They ate a lot  
16 of frozen stuff. So they didn't remember it  
17 specifically.

18           Of the 11 that were not part of this  
19 cluster by whole genome sequencing, none of them ate  
20 this product. So we were able to use a whole genome  
21 sequencing to identify a product that would have been  
22 very difficult to identify by pulsed-field gel

1 electrophoresis, and it certainly would have been  
2 delayed.

3           So the future of whole genome sequencing at  
4 public health labs is going to continue to increase.  
5 WGS will provide more information methodically  
6 compared to current methods, and it's very close to  
7 replacing some of our current methods such as pulsed-  
8 field gel electrophoresis and serotyping.

9           When we have a standardized WGS  
10 nomenclature that CDC is producing, it will greatly  
11 improve our communication and make the data sharing  
12 much easier.

13           So, in conclusion, WGS is a tool that can  
14 help us identify clusters and outbreaks better than  
15 our traditional methods, but there are current  
16 challenges and there will be additional challenges as  
17 we implement this technology, but speed and  
18 communication will be vital to continue these aspects  
19 for outbreak investigation.

20           Thank you.

21           DR. McDERMOTT: Good afternoon, everyone.  
22 Thank you to USDA for the invitation to join this

1 conversation and show you some of the ways in which  
2 we're exploiting this very power next-gen sequencing  
3 data in our antibiotic resistance surveillance work  
4 at FDA.

5           So a little context: I think everyone in  
6 this room is well aware of the global threat posed by  
7 antimicrobial resistance and a great deal has been  
8 done even in the last few years at the highest  
9 political levels to encourage and lead countries in  
10 developing and I should say dedicating resources to  
11 addressing it on a global/international level, and  
12 WHO came out with their global action plan for  
13 antimicrobial resistance back in 2015.

14           And some of the summary points in that  
15 report that helped set the context here, that they  
16 point out is that the development of resistance is  
17 linked to how often antibiotics are used. I'd say  
18 that's true in nearly every case, not every single  
19 case. Because many antibiotics belong to the same  
20 class of medicines, resistance to one specific  
21 antibiotic can lead to cross resistance to others,  
22 and I think if you hear these simple facts, you will

1 start to see quickly how whole genome sequencing can  
2 resolve some of these issues in more detail.

3           Resistance that develops in one organism or  
4 location can also spread rapidly and unpredictably  
5 and can affect antibiotic treatment on a wide range  
6 of infections and diseases, including those that  
7 spread between animals and humans. So it's important  
8 to keep in mind that the same classes of  
9 antimicrobials are used in food animal production, in  
10 treating our pets that's used in human medicine.

11           Now, there's some restrictions on that in  
12 pet and animal production. There are no restrictions  
13 on that in the antibiotics our veterinarians might  
14 use to treat infections in our companion animals.

15           And so we know resistant bacteria can be  
16 found in food animals and food products destined for  
17 human consumption including those same genes and  
18 strain types as we know from illness in humans.

19           And certainly my mentor in my postdoc years  
20 I thought just put it perfectly when he said  
21 antibiotics are societal drugs. It's the only class  
22 of essential medicines or any medicines which used in

1 one environment can compromise or affect how  
2 effective they are used later on in a different  
3 environment.

4           And so that brings with it a sort of new  
5 set of responsibilities, different ways of thinking  
6 in One Health, which I'm sure you all are also  
7 familiar with, has been the framework in which WHO  
8 has tried to build capacity in countries to do  
9 surveillance for resistance and also to address it on  
10 the policy level.

11           In the past, what we do in our NARMS  
12 program has been defined as integrated surveillance,  
13 and WHO defined that as the coordinated sampling and  
14 testing of bacteria from food animals, foods, and  
15 clinically ill humans and subsequent evaluation of  
16 resistance trends throughout the food production and  
17 supply chain using harmonized methods.

18           And I'll show you how that's what we do in  
19 our NARMS program, but one thing that we've been  
20 contemplating that came out of both a recent review  
21 of the NARMS program and our meeting in this same  
22 place, yesterday and the day before, is what are the

1 prospects and value in moving towards more of a One  
2 Health paradigm, and that One Health paradigm means  
3 sampling beyond that integrated food chain and  
4 looking into the environment and looking at animals  
5 and not just zoonotic bacteria from animals but  
6 animal pathogens and elsewhere. And again, whole  
7 genome sequencing technology, especially  
8 metagenomics, is going to make this more possible  
9 than ever I think.

10           So the purpose of NARMS, what we've said  
11 publicly for many years, is to monitor trends in  
12 resistance, get this information to people who can  
13 act on it, conduct research to fill in the gaps.  
14 It's very difficult to sample yourself to every root  
15 cause or to every solution and research has been an  
16 important part of filling in some of those gaps.

17           More recently, it's become adopted by CDC  
18 to prioritize outbreaks based on, now more than ever,  
19 whole genome sequencing data, and it's an important  
20 part of FDA's regulatory processes whereby a new  
21 animal antibiotic being proposed for review by a  
22 sponsor goes through qualitative risk assessment

1 steps that include mechanisms of antibiotic  
2 resistance, current resistance to other compounds,  
3 and so on and so forth.

4           So another way to look at this process  
5 overall and how surveillance fits into it, obviously  
6 surveillance is a key component to any public health  
7 action if you're going to get ahold of the magnitude  
8 and nature of whatever that hazard is, in this case,  
9 resistance, and understand the benefit of any  
10 interventions that were targeted based on the  
11 baseline and trend data.

12           So I describe NARMS as beginning with  
13 establishing baselines for resistance in different  
14 pathogens from different sources; how that resistance  
15 spreads; what it looks like over time. Can we get to  
16 where it might be coming from so that our decision  
17 making, whether it's regulatory or not is more of a  
18 scalpel than a hammer, as some described that?  
19 Understanding the contribution or the relationship of  
20 antimicrobial use in resistance is a tough one, that  
21 we're still struggling with but ultimately it hinges  
22 on human health impact and the burden of resistant



1 infections in humans.

2           And when that situation becomes untenable,  
3 then the interventions that are put in place again  
4 give NARMS, if we're doing our job right, an  
5 opportunity to measure that impact.

6           So this is the basic structure of the  
7 program. It is an interagency program between the  
8 FDA, the CDC, and the USDA. The Centers for Disease  
9 Control is getting isolates from the field, if you  
10 will, from clinical labs and physician labs around  
11 the country. Every 20th *Salmonella* is subjected to  
12 standard in vitro antimicrobial susceptibility  
13 testing but as we've heard, eventually all, and soon,  
14 all of the *Salmonella* reported will also be  
15 sequenced, but we'll continue to have phenotypic data  
16 on 5% of them.

17           USDA-FSIS has a random sampling of national  
18 production at slaughter of the four major food animal  
19 species, and then at FDA, where they retail meat  
20 testing is done in 2018, we'll have 22 states  
21 sampling 80 products every month of beef and pork and  
22 chicken and turkey, and then isolating *Salmonella* and

1 then characterizing them phenotypically and  
2 sequencing right now is being done at the FDA labs.

3           So that the integrated portion and  
4 structure right now. We combine these data into  
5 integrated reports where we try go harmonize our  
6 reporting as well as our methodology and our latest  
7 integrated report came out this Monday.

8           I think the CDC human report which looks  
9 beyond *Salmonella* and *Campylobacter* to other enteric  
10 pathogens, should be coming out very soon.

11           So the issue about what does this mean for  
12 resistance monitoring and beyond that, what does this  
13 mean for susceptibility testing and maybe some day  
14 clinical care?

15           You all know that currently what we do and  
16 what we've done for many decades to assess  
17 susceptibility in an organism is we do a biological  
18 response within serial dilutions of antibiotics and  
19 look for the lowest concentration that inhibits  
20 visible growth.

21           There's a lot of talk about how good is  
22 genomics going to be abridging to these old methods,

1 and I think I would start out by saying the old  
2 methods aren't as good as you might think they are.  
3 It's a biological response in an artificial situation  
4 and it's similar to what Dr. Besser said about PFGE.  
5 It's a predictor of the likelihood of success, what  
6 it's correlated with, with other parameters. The  
7 method itself, the laboratory method itself allows  
8 the three drug dilution range for QC for most  
9 organisms. Well, that could be 1 or 4 µg/ml. That's  
10 a pretty broad range, and in some cases, CLSI permits  
11 a 4 drug dilution range to be in QC. So it's not,  
12 you know, the chemists would always make fun of us  
13 when we tell them about this because it's very  
14 imprecise in a lot of ways, and it also doesn't  
15 always reflect all the things that go into resolving  
16 -- well, it doesn't reflect all that goes into  
17 resolving an infection.

18           So I just want to begin by saying that  
19 while the gold standard is MIC testing for  
20 susceptibility, it's got its own perils and so I  
21 think that's a good place to start.

22           And so when we start looking at how well

1 genomics can predict resistance or be useful in  
2 surveillance or in guiding clinical therapy, it's  
3 good to keep in mind that that's not a perfect method  
4 either. It doesn't always correlate well with  
5 clinical outcome as well.

6           And in a lot of cases, the studies haven't  
7 been done at all, and so breakpoints will be borrowed  
8 from one pathogen or animal into another without the  
9 clinical outcome data.

10           So you have this MIC. So what do you do  
11 with it? Well, you know, the title of this is the  
12 art and science of drawing a line somewhere, and it  
13 is an art and a science, and there's two ways right  
14 now which the data is considered.

15           One is how well does it predict clinical  
16 resistance? One MIC set with using a lot of data,  
17 using wild type MIC distributions, using clinical  
18 outcome information which is really the key one,  
19 using PK/PD data that addresses the concentration of  
20 the drug, at the site of infection, and then common  
21 standard of practice in medicine goes along with  
22 that. And FDA is responsible for setting those

1 breakpoints.

2           Another approach is what's called  
3 epidemiologically cutoff values. Let's identify  
4 everything that's no longer wild type and use that  
5 marker as a way to monitor resistance, and we could  
6 have a long conversation about fighting over what  
7 that word means because to some, it only means the  
8 former. It only means clinical resistance and  
9 clinical breakpoints but in some cases, to other  
10 people it means non-wild type or decreased  
11 susceptibility, depending on the questions you're  
12 asking.

13           And so in one case here, the breakpoint for  
14 epidemiological cutoff value in this mockup would be  
15 0.25, or in this case, 0.5. So here's a wild type  
16 distribution. You might have intermediate MICs in  
17 here that used to be indeterminate, recognizing the  
18 limitations of the method and then a clinical  
19 breakpoint. And so where you set that point can  
20 affect the data that you are describing over time.

21           Well, CLSI has tried to capture all these  
22 different definitions into one and describe resistant

1 strains, as those not inhibited by the usually  
2 achievable systemic concentrations of the agent with  
3 normal dosage schedules, and/or fall in the range  
4 where specific resistant mechanisms are likely, and  
5 here's our genomics, or clinical efficacy has not  
6 been reliable in treatment schemes. So you can see  
7 they're trying to really make everyone happy and  
8 sometimes the definition is as loose as the data.

9           So you've all seen this slide in some  
10 iteration probably before, but when genomics became  
11 routine and affordable, the first question in our  
12 surveillance system is, how do we incorporate it?  
13 How well can we bridge from the old to the new? How  
14 well does it predict resistance in pathogens? And we  
15 started off with a very basic question, what is the  
16 correlation between the presence of known resistance  
17 genes and isolates with a MIC above just the clinical  
18 breakpoint?

19           And then I'll say a little bit about how  
20 we're exploring metagenomics as well. And we've done  
21 three studies, a main study looking at this, and as  
22 Dr. Goldman pointed out, the correlations for the

1 bacteria under surveillance in NARMS is very high, 95  
2 to 99% correlation between MICs above the clinical  
3 breakpoint and the presence of known resistance  
4 genes, and we've done it for *Salmonella*,  
5 *Campylobacter* and *E. coli* and Greg Tyson who has done  
6 a lot of this work and is working on *Enterococcus*.

7           When NCBI and we started conversations  
8 about building these analytics into their pipelines,  
9 we expanded it to over 6,000 *Salmonella* now. The  
10 first study was 600, and the data holds up. It's  
11 very highly correlated with this one MIC, the  
12 resistant MIC in the presence of known genes.

13           We did a collaboration that's ongoing with  
14 Argonne National Labs where they took a machine  
15 learning approach and just said, can we predict MIC  
16 from the genome, blinded to the presence of any known  
17 resistance genes, and I couldn't explain the computer  
18 part of it. That's a little beyond me, but basically  
19 it was a k-mer based approach looking at correlation  
20 with onefold serial dilution above or below the MIC,  
21 and again it holds up incredibly well. This actually  
22 surprised me to see that you could actually predict

1 all the MICs for which we had ranges with the same  
2 high degree of confidence.

3           Now, I should emphasize, this work's  
4 ongoing and we're still trying to get some more  
5 dilution ranges to improve the data but it's  
6 incredible how well you can use the genomic data  
7 alone.

8           So what CDC's doing in sequencing all the  
9 *Salmonella* is a resistance monitoring system in my  
10 mind because these correlations are so good. So the  
11 dataset has just taken off with every *Salmonella*  
12 that's isolated. We can consider that in our risk  
13 assessment and regulatory processes at FDA because  
14 the correlation is solid. I'm going to skip one  
15 slide.

16           On the metagenomics side, just a few quick  
17 points. We have started to apply this technology  
18 both looking at animal cecal samples in NARMS and  
19 looking at the retail meat isolates to get an idea of  
20 really to sort of stress test the technology and  
21 explore its limitations for doing routine  
22 surveillance. And Daniel Tadesse in our group has



1 been working on this and we have kept every cecal  
2 sample collected in NARMS for the last 3 years and  
3 have some 20,000 of them, nationally representative  
4 randomized samples of U.S. food animal production.

5           So we didn't want to let these samples go  
6 in part because new policy changes were coming, and  
7 so we saved them all and have started to look at them  
8 metagenomically and started as this slide shows,  
9 through associate specific resistance alleles with  
10 different animal sources, and also in the retail meat  
11 samples as well. So stay tuned for updates on that.

12           In some of our earliest studies, back in,  
13 oh, my goodness, it must have been in '05, '06, when  
14 we did our first sequencing and it was very  
15 expensive, we got some information that I think is  
16 still a very good illustration of what the data mean  
17 to us and different ways in which it can be  
18 presented, and this plasmid just happens to show one  
19 of our first outputs which was a multidrug resistant  
20 *Salmonella* Newport on a backbone that was essentially  
21 identical to a strain from a child in Madagascar who  
22 had plague. So this plasmid, when you start thinking

1 about One Health, you see the importance of global  
2 One Health. What does it mean? Well, it certainly  
3 means these things can develop quickly and spread  
4 rapidly around the world and what is seen halfway  
5 round the world can become fairly common in the U.S.  
6 food supply, and this just shows different  
7 arrangements of genes but another part of it that's  
8 interesting and important is we see associations with  
9 resistance to decontaminating chemicals used in  
10 processing plants.

11           So new drivers of resistance emerge in  
12 these data, new associations and arrangements of  
13 genes with some indication of where they might be  
14 around the world.

15           And the next event is here, the next slide  
16 just shows a spike in gentamicin resistance that we  
17 saw in NARMS data, and it just illustrates the fact  
18 that we had never gotten to the bottom of the genes  
19 behind this without genomics because the PCR primer  
20 is available for aminoglycoside-resistant genes at  
21 that time didn't have all the alleles we found. So  
22 the new discovery is an obviously part of what we get

1 from whole genome sequence data routinely now. So  
2 what we called research and ad hoc PCR studies is now  
3 routine surveillance.

4 Another important event that showed the  
5 value of this technology was plasmid-mediated  
6 colistin resistance that emerged in China in the fall  
7 of 2015, Errol Strain was really helpful in this in  
8 going into NCBI's database and saying have we ever  
9 seen this in a domestic isolate of any bacteria in  
10 the United States, and at that time, screened 155,000  
11 genomes at NCBI and said it wasn't there. Incredibly  
12 powerful to do retrospective surveillance like this  
13 without opening the freezer.

14 We later did studies of selective  
15 enrichment and found *mcr-1* in 2 out of 500 swine  
16 samples, and the metagenomics tools we were  
17 developing at the time picked it up as well. So that  
18 bodes well for the future of a metagenomic component  
19 to the surveillance.

20 So that's just a few illustrations of how  
21 we're building it into NARMS. I would invite you to  
22 go look at our last report that came out Monday where

1 we've put genomics now into our dataset. We went  
2 back and sequenced every retail meat *Salmonella* back  
3 to 2002, and as David Goldman pointed out, they're  
4 sequencing all the animal isolates. Next year we'll  
5 have all the human isolates. So knowing how good the  
6 predictions were of genotype to phenotype, we put  
7 quite a bit of effort into interactive data displays  
8 where you can look at resistance over time and see  
9 the genes that go with it from the different sources.  
10 And it's really powerful, and it gives people access  
11 to the data in new ways and allows them to ask and  
12 answer their own questions.

13           One of the good ones is we've gotten beyond  
14 this really crude metric of MDR, multidrug  
15 resistance, being resistance to three or more drugs.  
16 Well, what drugs, right? So now you can go in and  
17 see what drugs and refine the analysis and have more  
18 confidence in what you're talking about in the trend.

19           Another thing we've developed is, to take  
20 the next obvious step, we harvested all the NCBI  
21 data, screened it for all its resistance genes and  
22 took advantage of whatever metadata that were there

1 with its strengths and limitations, and you can do  
2 the same thing now with this tool we called Resistome  
3 Tracker, where you can look in any source according  
4 to the metadata categories, at any resistance gene in  
5 this case, an aminoglycoside resistance gene, *aadA1*,  
6 in chicken. You can see where it is at NCBI by  
7 biosample, and then what our intention was, just to  
8 set it on top of NCBI. You can click on any of those  
9 isolates within any of the genes you're interested in  
10 and go right to the SNP tree.

11           So this just seemed like an obvious thing  
12 to do based on NARMS data and taking advantage of the  
13 sequencing that was coming out.

14           And then the last tool we built into  
15 Resistome Tracker is, well, what's new? We want to  
16 know the wallpaper in the room, if you will. What's  
17 the backdrop for analyzing our domestic data in  
18 NARMS? So any new gene that comes up in any  
19 submission to NCBI, we've put an alert system on  
20 there so we can track it and see the data in which it  
21 was reported.

22           And Heather Tate who is here has developed

1 this Resistome Tracker tool, and we're going to try  
2 to launch it. It's not out yet. We got busy this  
3 week, but maybe in the next week or two, you'll be  
4 able to play around with this and see if you like it.  
5 It includes a mapping tool where you can look at over  
6 time when these different resistance genes appeared  
7 around the world, and you can apply this, of course,  
8 to any pathogen.

9           So we've put a lot of effort into making  
10 sense of the complex data that we are all generating  
11 through genomics.

12           So, in closing, One Health - One Method. I  
13 know it's a bit of an exaggeration, but it sure feels  
14 like it. We can predict so much from the genome. We  
15 can predict resistance in our target pathogens so  
16 well from the sequence data. It overcomes so many  
17 past limitations with the metagenomics, you know, our  
18 reliance on cultivatable organisms which is expensive  
19 and obviously very narrow, and what we can say about  
20 the Resistome.

21           The best part to me, the second bullet,  
22 permits us to look farther with fewer resources and

1 lower costs. So now we can get out and look at the  
2 costal waterways and look at surface water systems  
3 and start incorporating companion and add on other  
4 testing because we can glean what we need from these  
5 samples, and that's where we're trying to go next.

6 I've noted how it reveals new associations  
7 including determinants perhaps of animal origin of  
8 other drivers of resistance and new alleles, greater  
9 confidence in decision making. That's important. We  
10 would like to note just where a pathogen comes from,  
11 but where a resistance comes from.

12 And with that, I can't thank everyone, but  
13 it's been a tremendous opportunity to collaborate  
14 broadly, and it's been a real satisfying experience  
15 working with NCBI, who we will hear from next, and  
16 CFSAN and CDC and USDA on this. Thank you.

17 DR. KLIMKE: So this is great. Everyone's  
18 talked about NCBI. So I don't have much to talk  
19 about. Go to the next slide please.

20 So for those of you not familiar with what  
21 NCBI is, we are the premier biomedical informatics  
22 institute on the planet. We are the host for all the

1 PubMed, all the medical literature over 23 million  
2 publications, all the clinical data, and all the  
3 sequence data that you've heard about from today is  
4 being submitted to our databases, and then we have  
5 specialized databases and tools, and I'll talk about  
6 one in particular today called our pathogen detection  
7 system.

8           Next slide.

9           And again to reiterate what someone has  
10 already said, we share all the sequence data with the  
11 European Bioinformatics Institute and the DNA  
12 database in Japan within 24 hours of the data being  
13 released.

14           Next slide. Go to the next one.

15           This is a rough schematic of how our system  
16 for the pathogen detection pipeline works. We have  
17 submitters on the left including, you know, FDA  
18 GenomeTrakr and CDC PulseNet and USDA, submitting  
19 data to us, into the public databases. And then we  
20 have a pipeline which I won't go into a lot of detail  
21 today. There will be one other slide on that. We  
22 basically do an assembly, a clustering, and we



1 produce phylogenetic trees, and we make the reports  
2 publicly available to all of you.

3           Next slide, please.

4           I want to talk briefly about the metadata.  
5 This is the template that we built, and you seen that  
6 from a few other people already this morning. This  
7 came up from discussions with FDA about the sort of  
8 metadata they wanted to see submitted to the  
9 databases that would enable them to do their job.  
10 It's basically four categories of information. What  
11 is the sample, including the organism, a unique  
12 identifier such as a strain or isolate number, it's  
13 categorization, whether it comes from a  
14 clinical/host-associated sample or from a food or  
15 environmental sample, and then a few other categories  
16 such as when, where, and who.

17           The ones I've colored in yellow are the  
18 absolutely minimal sample fields that need to be  
19 filled out in order to submit using this template,  
20 and that's what CDC is using right now, and 6 months  
21 after the sequence gets deposited, they update with  
22 some additional metadata fields.

1           The USDA and FDA, and you've already heard  
2 from some of them this morning, they are willing to  
3 submit things like the geographic location, the year,  
4 and month the data was submitted as well.

5           This has been an incredibly successful  
6 template. So we have many, many other templates for  
7 bacteria, for example, at NCBI. This one has been  
8 used by over 256,000 submissions to date, and that's  
9 not just samples from PulseNet and GenomeTrakr. That  
10 includes other academic laboratories and other state  
11 actors that are willing to submit data to it, and you  
12 can see the breakdown of the clinical versus the food  
13 and environmental. I'll talk about some of these  
14 metadata fields later on.

15           Next slide.

16           So basically what we're trying to do in  
17 this system is we're taking very large volumes of  
18 data and reducing it to relevant data. So we might  
19 have hundreds of millions of base pairs in a short  
20 read sequence for one particular isolate. We do an  
21 assembly, and then we're turning that around into a  
22 phylogenetic tree along with associated antimicrobial

1 resistance genes, and we want to do that within 24  
2 hours of the sequence being deposited, you know, so  
3 turning big data into useful data.

4 Next slide, please.

5 And we have to do this because we know that  
6 the sequences that are coming in are going to  
7 increase in the future. So this is just a snapshot  
8 of the last 6 months of the data that's been  
9 submitted to our pathogen detection system, and you  
10 can see it's dominated by *Salmonella* in that orange  
11 slice of the pie, and then the other three foodborne  
12 pathogens are basically making up the other 95% of  
13 that. And so we have a few clinical things that are  
14 coming in, but basically it's the foodborne isolates  
15 that are coming into our system.

16 And we know that the *Salmonella* that are  
17 being collected in the U.S. should all be sequenced  
18 by the end of 2018, and so actually we should see the  
19 rate of that double within the next 18 months or so.  
20 So we should expect, I believe it's 90,000 isolates  
21 per year for all the foodborne pathogens in the U.S.

22 And you can see at the bottom is the graph

1 showing the sort of number of submissions per day and  
2 you can see, when we first turned this on in 2013,  
3 that's when we first started the Real Time *Listeria*  
4 Project, where all *Listeria* in the U.S. were being  
5 sequenced and submitted to us, and that's ramping up  
6 for the other foodborne pathogens. You can see both  
7 the spiking nature, sometimes we get large volumes of  
8 submissions from, you know, Public Health England is  
9 an international partner that sometimes submits in  
10 batch mode. They might submit 1,000 or 10,000 in a  
11 day sometimes.

12           But you can see that it's growing, and we  
13 would expect that rate to increase as we move forward  
14 into the future.

15           Next slide, please.

16           This was a snapshot of data that I  
17 presented at NIST as part of their standards for  
18 pathogen detection. This was basically looking at  
19 the total number of submissions and whether they're  
20 clustered or not for the four foodborne pathogens,  
21 and again you can see *Salmonella* is predominant on  
22 this slide.

1           The interesting thing about the clustering,  
2 so I didn't mention this, I'll probably mention it in  
3 the next slide or so, we make clusters of 50 SNPs or  
4 less. So very closely related isolates, and you can  
5 see most of the *Listeria* and most of the *Salmonella*  
6 clustered, *E. coli* and *Shigella* not so much. So  
7 that's an important byproduct of the biology of those  
8 organisms, and I'll come back to this in some of the  
9 analysis I do.

10           Next slide, please.

11           So this is the last pipeline slide you'll  
12 see. So those of you who were at my other talks can  
13 decide if you don't need to see any more of those,  
14 but basically the data comes into us, you know, and  
15 it's several days before it gets sequenced and  
16 submitted to NCBI, but we do an assembly and we are  
17 starting to do wgMLST allele calling and producing a  
18 table of nearest neighbors. We're doing that now for  
19 *Salmonella* and *Listeria*, and that's been running for  
20 2½ months, and we're aiming to get that data back to  
21 FDA within 1 hour of the sequence being deposited.  
22 So that's an incredibly fast turnaround time, where

1 the sequence comes off the machine and comes to NCBI,  
2 and if Errol does it right, he says, you know, he'll  
3 submit it in the morning, have his coffee and then he  
4 can look at the results, you know, by lunchtime. And  
5 he can make a decision about the inclusion/exclusion  
6 within that 1 hour.

7           We're also doing SNP clustering. This  
8 slide, you know, we're trying to replace the initial  
9 clustering we do using whole genome MLST and we aim  
10 to get the clusters of these nearest neighbors into  
11 phylogenetic trees and on the web-based interface,  
12 and I'll show you some examples, within 24 hours.

13           So the rapid reports just gives them an  
14 extra day to say, you know, these isolates are  
15 related or not related and make decisions based on  
16 that.

17           Next slide, please.

18           So just for the rapid reports, like I said,  
19 we're doing *Salmonella* and *Listeria*. We're reporting  
20 the five nearest neighbors in all neighbors with less  
21 than six allele differences, and that cutoff may  
22 change. You know, we're right in the pilot phase of

1 this project, and we've already gotten some  
2 preliminary feedback that we may want to change that  
3 threshold, reporting the number of difference, number  
4 in command and the SNP accession if it exists. So if  
5 they can actually see within the list of nearest  
6 neighbors that some of those are already participants  
7 in an existing SNP cluster, they can actually take  
8 additional steps.

9           And we put this into one file per run.  
10 Most of those reports are already being made  
11 available on the FTP site. So I won't show you an  
12 example of that, but you can go take a look.

13           Next slide, please.

14           And all of this data is being made  
15 available into our publicly available pathogen  
16 detection website. So you can go to this page right  
17 now and see all the data being submitted from  
18 GenomeTrakr and PulseNet, and you can interrogate the  
19 data, and I'll show you some example of that.  
20 There's a couple of ways into it, but we already have  
21 over 174,000 pathogens and at least 142,000 with  
22 either an acquired or chromosomal antibiotic

1 resistance gene.

2 Next slide.

3 So I just wanted to do a quick analysis,  
4 just basically for my own benefit. Using the  
5 metadata that's being submitted, I wanted to look at  
6 the isolation source. The isolation source, the  
7 definition of that is it describes the physical,  
8 environmental and/or local geographical source of the  
9 biological sample from which it was derived. So for  
10 a clinical isolate, that might be blood, stool,  
11 urine. For a food or an environmental isolate, it  
12 might be a river bed. It could be a food contact  
13 surface. It would be an environmental swab. It  
14 could be cheese.

15 So next slide.

16 So I simply wanted to ask the question, if  
17 I look at all the submissions coming into our system  
18 from CDC, FDA, and USDA, including the state labs  
19 that are submitting under GenomeTrakr and PulseNet  
20 submissions, and if I want to look at all pairs  
21 within 10 SNPs, what am I going to see? And I want  
22 to categories each one of those pairs, either they're



1 clinical versus clinical, clinical versus food or  
2 environmental and that could be potentially the  
3 smoking gun for causing clinical illness, and then  
4 food and environmental versus food and environmental.  
5 So go to the next slide.

6           This is a summary of those pairs. We're  
7 just looking at counts. We have the four foodborne  
8 pathogens at the bottom, you know, so *Campylobacter*,  
9 *E. coli/Shigella*, *Listeria*, and *Salmonella*. And then  
10 the pairs are those three categories I just told you  
11 about. Clinic/clinical are red. Clinical and  
12 environmental are green. And environmental and  
13 environmental are blue.

14           So you can see obviously *Salmonella*  
15 clinical cases completely dominates this slide. I'm  
16 not telling you anything new. You already knew that  
17 *Salmonella* was one of the biggest problems in the  
18 U.S.

19           What you will see though is that the number  
20 of cases of those pairs is not uniformly distributed.  
21 So *Campylobacter* is pretty uniform between the three  
22 categories. *E. coli/Shigella* and *Listeria* are not.

1 So I think this again informs you about something  
2 about the biology of those pairs and the distances  
3 that I use. I use a 10 SNP threshold. That was an  
4 arbitrary pick. A lot of people that come to these  
5 meetings, you know, say what is your threshold? You  
6 cannot, again I'll reiterate, you cannot simply just  
7 pick one threshold and get the answer. There's  
8 additional contextual information. Go to the next  
9 slide.

10 By looking at the number of samples within  
11 those thresholds, you can see most of the samples in  
12 *Salmonella* submitted are within that 10 SNP threshold  
13 to something else. That's not the same for *Listeria*.  
14 About 1/3 of the clinical cases are within that SNP  
15 threshold, and so it varies again between all the  
16 organisms.

17 The last column is also interesting because  
18 I said we made these SNP clusters of 50 SNPs or less.  
19 There's a number of cases that are just not clustered  
20 with anything or not clustered with a food or  
21 environmental isolate. So there's actually over  
22 2,000 clinical *Salmonella* not linked to any food or

1 environmental isolate at all in our database and, you  
2 know, 1500 *E. coli*, and I'll come back to that just  
3 near the end of my talk.

4           Next slide.

5           So these next four slides are just looking  
6 at the four foodborne pathogens and what that  
7 isolation source breakdown is, and this is free text.  
8 So I'm just basically doing a naïve sort of smushing  
9 together of some of the terms into these sort of  
10 categories so we don't have a 100 rows. We only have  
11 20 rows or so.

12           And again, you can see *Salmonella*,  
13 predominantly chicken, you know, we have some beef,  
14 pork, water. You have papaya, and it's high in  
15 yield, and I'll come back to that. And then things  
16 like research strains. So somebody in Micro 101 is  
17 just not doing, you know, the proper technique that  
18 they're supposed to be learning.

19           Next slide.

20           *Listeria*, again swabs, environmental swabs,  
21 cheese and various food products.

22           Next slide.

1           And then *E. coli/Shigella* dominated by soy  
2 nut butter and then beef, milk, goat, etc.

3           Next slide.

4           And then *Campylobacter*, of course, milk is  
5 a big problem, unpasteurized milk and beef again. So  
6 go to the next slide.

7           And this, of course, completely matches  
8 what you see on the multi-distributed outbreaks for  
9 CDC. You know, you see soy nut butter for *E. coli*.  
10 You see papaya. There's a huge papaya problem from  
11 Mexico, and you see cheese as a problem for *Listeria*.

12           So I'm not trying to present this as like a  
13 detailed scientific analysis of the metadata. I'm  
14 just simply saying that we make all this data  
15 available. We can do these types of analysis. This  
16 only took me a couple of hours in an afternoon. You  
17 could do similar sorts of things, and I'll dive into  
18 some more detailed analysis that you can do. So next  
19 slide please.

20           So, for example, a lot of people talked  
21 about isolates browser. Here I'm just doing a search  
22 for papaya into our isolates browser. You probably

1 can't see it in the back, but there's 225 isolates  
2 that come back with a search term of papaya in our  
3 existing database when I did this. That big table in  
4 the back is called the isolates browser, and  
5 basically every row in that table is an assembled  
6 gene that's come through our system or come through  
7 GenBank. It's got the associated metadata and you  
8 can actually, you know, add columns to that based on  
9 what you want to do.

10           We do a separate thing where we intersect  
11 the number of searches with the number of clusters  
12 and that's the little table at the bottom which is  
13 actually on the default page on the upper left. If  
14 you can't read that it says, you know, at the top  
15 row, 23 of the total 32 isolates in this cluster have  
16 the search on papaya, and if you go to the next  
17 slide.

18           So this is our new tree viewer which is not  
19 available to you. This is being alpha tested by FDA  
20 right now, and so if you went to our page and clicked  
21 on the SNP trees and you saw some examples of that  
22 earlier, you wouldn't see this exact page. But this

1 is something that we're testing at FDA because we  
2 want to make tools available to them that will help  
3 them make a decision based on the genetic distance.

4           So this is one example. This is a cluster  
5 of 710 isolates. Eight of them have the search on  
6 papaya. But only one of them is from 2017. If you  
7 see on the left, there's a breakdown by date, and if  
8 you look at that part of the tree where the food or  
9 environmental isolate hits, it's basically on a  
10 completely separate branch of the tree. It's not  
11 near any clinical isolates. It's at least over 30  
12 SNPs away from any clinical case, and so that  
13 particular papaya probably is not a positive of any  
14 clinical illness that's been recorded. If you go to  
15 the next slide though, we see the opposite.

16           This is another *Salmonella* cluster. It's  
17 got over 1,100 isolates again from papaya, but 3 of  
18 them are from 2017. And if we look at the SNP tree,  
19 you can see one environmental isolate is incredibly  
20 close to a large number of clinical isolates. Go to  
21 the next slide.

22           And if you look at the entire subtree,

1 you'll see there's 112 isolates in total, and they're  
2 incredibly closely related. The max SNP distance is  
3 17. The average SNP distance is 3, and if you'll go  
4 to the next slide.

5           If we just look at the food or  
6 environmental isolates in that total list of 112, 3  
7 of them exist. One is that papaya from Mexico, and  
8 two are papayas from the U.S.

9           And so what I showed you is I think that  
10 you can quickly see whether something is inclusive or  
11 exclusive, make a decision very quickly using just  
12 genetic distances with some limited metadata and then  
13 you can go onto the rest of whatever your job is with  
14 your public health or regulatory agency.

15           And again, so we're building tools used to  
16 facilitate that. We're not the agency that makes  
17 those decisions, but we want to make tools so you're  
18 95% of the way there without any extra effort.

19           Next slide.

20           So just going back to unlinked isolates.  
21 Again clinical isolates are not within 50 SNPs of any  
22 food or environmental isolate. They're not just

1 singletons. So I thought when I first did this last  
2 week, that these would all be completely unlinked to  
3 anything but it's not true. Some of them were  
4 actually in some very large clusters.

5           So my question, and I'm going to ask these  
6 questions, I don't know the answers to them: Do we  
7 need a new sampling strategy? Are there new food  
8 vehicles that are waiting to be discovered? What  
9 does the epidemiology tell us? And I can't answer  
10 those questions because I don't have that  
11 information. You people have that information, and  
12 so you should be asking these questions and be able  
13 to answer them, again pointing out the number of  
14 isolates. And if we go to the next slide.

15           If we just look at one of those, it's an  
16 extremely large *S. e.* cluster of 506 clinical cases.  
17 There no food or environmental isolates in these  
18 cluster at all. They're all collected from 2013 to  
19 2017. They come from the U.S. They come from Public  
20 Health England. So there's isolates in the UK that  
21 also cluster with these isolates. They come from  
22 state labs. So what is the cause?



1 I've heard something. I've heard a rumor,  
2 but I'm not going to tell you because I don't think I  
3 can share that, but obviously these are the sort of  
4 things we would be looking at, you know, based on the  
5 questions that I asked earlier. Using our interface,  
6 you can do that but also the places that produce this  
7 data can be asking these questions and answering  
8 them.

9 Next slide.

10 Just to basically end on the antibiotic  
11 resistance. You've already seen some of the talks  
12 using some of our databases and some of the analyses  
13 that we're doing. This came about because the  
14 President at the time decided that we needed to  
15 combat antibiotic resistance. So we put into motion  
16 this thing called CARB Report, you know, from the  
17 President's scientific committee Combating Antibiotic  
18 Resistant Bacteria. And NIH has mentioned twice in  
19 there, these two critical elements. One is to  
20 produce a reference database, a well-curated  
21 reference database and maintain a national sequence  
22 database of resistant pathogens. So go to the next

1 slide.

2           So we've done that in a number of ways.  
3 First, we've actually built a template to capture the  
4 susceptibility that people are collecting. It's  
5 called antibiogram. In addition to the sample  
6 metadata, the minimal template that we're collecting,  
7 you can actually submit AST data as well. We have  
8 almost 5,000 submissions using that template right  
9 now.

10           And we have put together a reference  
11 database of antibiotic resistance genes. Now, this  
12 is not something that we just did on our own.  
13 There's a lot of people out there that say, well,  
14 which database should I use. Well, we have a  
15 collaboration going with CARD. We may have a  
16 collaboration going with ResFinder depending on some  
17 upcoming discussions we have. We've taken over the  
18 Lahey database. So they were the place where you  
19 would submit novel beta-lactamase alleles to, and  
20 they would assign a new novel beta-lactamase allele,  
21 let's say SHV or TEM. They're retiring, and so  
22 they've asked us to take over that responsibility,

1 and so we're the place where you would actually  
2 submit those for, you know, when you want to make a  
3 publication into *JAC* or *AAC*, you would actually  
4 submit to us, get the novel allele number, and then  
5 this goes into your publication. Once it gets  
6 released, it gets fed into our reference database and  
7 used.

8           And so we're implementing tools for  
9 identifying those genes using that reference  
10 database, and if we could go to the next slide.

11           We're actually integrating that back into  
12 the isolates browser that I just showed you. So you  
13 can actually see the list of genes per isolate  
14 integrated directly into the list. So this is an  
15 example of another large *Salmonella* cluster. We have  
16 50 isolates encode in *mcr*, so mobilized colistin  
17 resistance gene, and you can see just one example,  
18 probably you can't see it, but there's a cluster of  
19 four isolates from Thailand that all have a *mcr-1*  
20 gene.

21           And so this system will allow you to  
22 interrogate for not only metadata but also the

1 presence or absence of different resistance genes.

2 Next slide.

3 So just to summarize, we're enhancing our  
4 existing analytical pipelines to improve the  
5 turnaround time to answer, you know, basically these  
6 two fundamental questions. Are these isolates  
7 colonally related? Is there a point source for  
8 clinical illnesses?

9 And we're improving these interfaces,  
10 enhancing the information that's layered on top of  
11 them including antibiotic resistance genes, making  
12 the system much more easy for you to use to make  
13 determination of inclusion/exclusion, and in the  
14 future, we'll be adding things like virulence genes,  
15 heavy metal resistance, point mutations, mobile  
16 elements.

17 Next slide.

18 So that's it. If you have any questions,  
19 you can email that email address or come see me here.  
20 These are all the people that I work with on the  
21 pipeline at various points. So you can basically see  
22 it's an army of people that NCBI that are involved

1 with this, although the day-to-day operations, the  
2 number of people actually involved are just a few  
3 people.

4 I'd like to thank all of our colleagues at  
5 FDA, CDC, and USDA for making the data available  
6 because we wouldn't be able to build these tools if  
7 they didn't make the data available, and thanks to  
8 David Lipman for helping to push this system out even  
9 though he's now gone off to work on food products  
10 actually. I'll be happy to take any questions at any  
11 point.

12 DR. EVANS: Thanks, Bill. So we're going  
13 to for a quick coffee break, and we'll be back here  
14 at 3:00 with two additional presentations and then an  
15 opportunity for questions and answers.

16 (Off the record at 2:30 p.m.)

17 (On the record at 3:00 p.m.)

18 DR. EVANS: If everybody can take their  
19 seats, we're going to start up again with a  
20 Demonstration of Tools for WGS Analysis.

21 Okay. I want to welcome Glenn Tillman with  
22 the Molecular Characterization Branch of OPHS at USDA

1 FSIS, and he's going to talk about tools for a WGS  
2 analysis.

3 DR. TILLMAN: Hello, and good afternoon.  
4 Thank you for the invitation to come today. It's a  
5 very good pleasure after hearing all these great  
6 talks this morning. You really set the stage well.  
7 The first three talks went into great detail about  
8 where we're going with whole genome sequencing. So  
9 did the rest of the talks. So I hope to have  
10 something good for you all to add. So without  
11 further adieu, that's me. I don't like the picture  
12 very much.

13 Okay. And this is what I'm kind of  
14 planning on talking about, short and to the point.  
15 Give you a little background, but number 2 is the  
16 bullet point, what are we doing with this data? And  
17 that's what most people are here to hear about. So  
18 we're go straight into that.

19 So some of these things Dr. Goldman already  
20 spoke of on behalf of our Agency. Why are we using  
21 WGS? To support foodborne illness and  
22 investigations, to support our mission goals as part

1 of FSIS and mainly and more importantly the whole  
2 point of this is alignment with public health  
3 partners. You see a lot of familiar names who have  
4 been up here speaking over the last day, 2 days  
5 considering the NARMS that was last week, I mean  
6 earlier this week, excuse me.

7           So where are we currently right now? We  
8 currently have 12 MiSeq sequencers in our laboratory  
9 system. Eight of those are in our Eastern Lab in  
10 Athens, Georgia, where I'm located. Two are in our  
11 Midwestern Lab, and two are in our Western Lab.

12           We, in Fiscal Year '17, really ramped up  
13 and sequenced around 7,200 isolates. So it's quite a  
14 jump from our capacity building perspective.

15           We do collaborate really well with our  
16 public health partners, and we do consider  
17 epidemiological information in all that we do with  
18 these new emerging tools as many of the speakers have  
19 alluded to this morning. It's one big package of  
20 information that you get.

21           Finally, we do work really strongly with  
22 our NARMS partners. Dr. McDermott gave a great talk

1 earlier on how this tool can enhance what we're going  
2 and we're trending with antimicrobial resistance  
3 genes.

4           So I want to point this out to you as a  
5 reference guide. Dr. Klimke gave a nice talk on  
6 NCBI. Well, here are some of the top bioprojects  
7 that we and FSIS contribute to. Most of these are  
8 GenomeTrakr-specific projects. A few of these are  
9 set up specifically for us with our efforts with the  
10 NARMS program at FDA CVM. So many of these are from  
11 the cecal environment of our four major animal  
12 commodities. So please take note of those and peruse  
13 those bioprojects as needed.

14           And this last just shows kind of an upwards  
15 trend of how many isolates we've gone from kind of a  
16 low level in fiscal year 2014 to where we're at now  
17 with over 7200 in fiscal year 2017. We plan on going  
18 even further this year.

19           We talked a lot about metadata here today.  
20 Here's a nice example of one of our biosamples. We  
21 do release who collected it, of course, the  
22 collection date and the state in which it's



1 collected, and isolation source. We try to be very  
2 descriptive, as Bill was mentioning in his last talk,  
3 of how we could collapse some of those commodities  
4 into a certain product type, non-meat product swab.

5           Okay. Now, onto more of the analytical  
6 tools that we may use. I was actually asked to do a  
7 demonstration, but I couldn't bring myself to do that  
8 in command line in front of you all and streaming  
9 over the web. So we're just going to stick with some  
10 screenshots.

11           We'll start with how we do quality control  
12 assessment, antibiotic resistance gene detection,  
13 *Salmonella* and STEC serotype determination and STEC  
14 virulence gene characterization, and finally  
15 phylogenetic comparisons. So it will be a lot of  
16 stuff coming at you really quickly.

17           Currently, our group is responsible for  
18 characterizing upwards of 12,000 isolates per year,  
19 and that includes antimicrobial susceptibility  
20 testing, serotyping, PFGE, and whole genome  
21 sequencing. So we're doing all these tests in  
22 parallel. So it's a pretty busy group and a lot of

1 results coming out of there.

2           The goal over the next couple of years is  
3 to continually develop one single workflow because as  
4 we know, we can get all these types of information  
5 directly from whole genome sequencing, and we want to  
6 continue to try and pursue that as much as possible.

7           So an overview, we start with the  
8 sequencer. I think this might be the first picture  
9 you've seen of the MiSeq today. Again, we have 12 of  
10 these. If you can see up here, you'll see kind of  
11 what a FASTQ format looks like if you're not familiar  
12 with them. It's a different type of format than most  
13 of us are used to working with over the past 5 years.  
14 It's highly compressed, but even at compression  
15 rates, it's still around 300 megabytes per 2 files  
16 together for a single isolate. So you're talking  
17 large datasets.

18           That data goes one of three ways actually.  
19 You input the FASTQ into an assembler. In this case,  
20 I show a picture of a CLC Genomics Workbench. We  
21 also use Spade to do assemblies, de novo assemblies.  
22 We do a quality control pipeline where we assess

1 coverage from the raw files, the average quality and  
2 nucleotide balance. All these are worked out in our  
3 Gen-FS partnerships with NCBI and FDA and CDC.

4 FASTQ files can go into the whole genome  
5 MLST, in the BioNumerics 7.6, which was talked about  
6 by Dr. Besser earlier. Lyve-SET and SNP pipeline are  
7 two types of high quality SNP analyses that are  
8 actually publicly available, and I'll have some  
9 screenshots of that later on. And then finally, the  
10 NCBI Pathogen Isolate Browser. That one's nice  
11 because all the heavy crunch is done on the NCBI  
12 side, and we don't have to do that.

13 Okay. So once we've got our assembly, we  
14 do another set of QC, quality control. We look at  
15 file size. That's important. One byte equals a base  
16 essentially. We want to have some level of  
17 correlation that our assembly is very correlated with  
18 the actual bacterial genome size, 3 megabytes for a  
19 *Listeria monocytogenes* assembly. That equates nicely  
20 with 3 million bases, essentially what the take home  
21 on that is. We look at certain other metrics, N50,  
22 number of contigs. The lower the number of contigs,

1 the better the assembly, the better the sequencing  
2 quality that went in. We want low number of contigs.

3           Finally, correct organism, that's very  
4 important to us. We have multiple things in house to  
5 assess that, and on the NIH side, NCBI, Bill Klimke  
6 has talked at numerous conferences about there's  
7 still some level of sample mismatches and mix ups.  
8 We try and avoid that all throughout our pipeline  
9 here and with our other assays by double checking at  
10 multiple points.

11           Now, here's the FASTA form as you can see.  
12 It's a little bit different than the FASTQ file.  
13 Again, you've gone from 300 megabytes to around 3  
14 megabytes. So you've highly compressed really  
15 complex data into a much smaller workable format.

16           Then you're going to input that FASTA into  
17 our downstream tools, determination of MLST,  
18 multi-locus sequence type, antibiotic resistance  
19 genes, virulence profile and serotype determination  
20 and potentially you can even do different types of  
21 file genetic comparisons using MASH Tree.

22           So to start off, I'll talk a little bit

1 about the antimicrobial resistance we use, the whole  
2 genome sequencing data, that's been talked about a  
3 lot again over this week.

4           We want to work to identify new genes of  
5 concern. That's been a big part of last couple of  
6 years, and FSIS efforts have been part of that with  
7 the CTX-M-65 and other genes. As I mentioned, the  
8 CTX-M-65. We have colistin we've worked to identify.  
9 Do we have any of these sequences in any of our  
10 sequences that we're putting upon NCBI? Quinolone  
11 resistance and the spread of plasmid-mediated  
12 quinolone resistance, *qnrB19*. We work with our  
13 partners at NARMS with that. Linezolid resistance as  
14 well.

15           Okay. So the overall kind of workflow, you  
16 saw what the FASTQ file looks like. So it goes into  
17 the QC in the assembly pipeline and then from there,  
18 the assembly goes into our blast database, and the  
19 output is a resistance gene profile. So that's a  
20 really quick run over of what we're doing.

21           But where do we get the antibiotic  
22 resistance database from? We use ResFinder

1 currently. We get that from the Center for Genomic  
2 Epidemiology, and much like our partners, we update  
3 that on a very, very frequent basis.

4           Here's the output. Here's the black box in  
5 which we work. Most of this is done in command line,  
6 and there's a reason for that. All these tools are  
7 publicly available and you can do them in most  
8 browser type formats. You can do that for 1 to 2  
9 isolates, but when you do it for 7,000, you've got to  
10 have a much, much better and more efficient way of  
11 doing it. That's where in our hands, command line  
12 comes in.

13           This is kind of a typical output you'll see  
14 and then in the end, we do a lot of formatting with  
15 our bioinformaticist. A lot of their work is in  
16 formatting how we can get it into a file that's  
17 usable for other people.

18           So in this particular isolate, we  
19 identified five different genes confirming multiple  
20 levels of resistance including *bla* CTX-M-65, tet  
21 genes and sulfonamide and aminoglycoside genes.

22           All this can be found again publicly, for

1 the download portion at bitbucket.org, which I  
2 actually think can be linked to from the genomic  
3 epidemiology site at this point.

4           Okay. So one of the things I wanted to  
5 show was that phenotype and genotype comparison has  
6 been talked about. So we did contribute to the NCBI  
7 and FDA initiative recently to identify what that  
8 correlation was, but we did our own work in house  
9 with the 2016 cecal NARMS data which is essentially  
10 1190 isolates. Dr. Goldman showed this earlier, and  
11 you can see and hopefully appreciate around 97 to 99%  
12 correlation. Like everyone else has mentioned,  
13 there's some gentamicin issues but overall this is a  
14 pretty tight correlation.

15           Okay. *Salmonella* serotype determination,  
16 very similar to what you saw before, assemblies key,  
17 putting it through our custom made BLAST database  
18 based off of SeqSero developed by the University of  
19 Georgia and the CDC.

20           That one is actually available on GitHub.  
21 A lot of these different tools are available on  
22 GitHub. You can see really faint screenshots up here

1 now against this backdrop, but you will see they are  
2 available and you can use those.

3           So to give you a little bit more background  
4 on that, several years ago when we first started this  
5 initiative, SeqSero was just coming into its own. So  
6 what we did was an exact matching algorithm using  
7 some python-scripts that we developed. So now we  
8 currently still use those python-scripts in  
9 conjunction with SeqSero, remove all the redundant  
10 factors and use that against the dictionary to  
11 identify the serotype.

12           So in our hands, we looked at over 4,200  
13 isolates, and we found about a 96% rate where WGS  
14 matched that of the reported serotype result.

15           We had about 3.8% where WGS did not match  
16 the serology result, and that typically was just  
17 incomplete genetic factors that we did not get  
18 allowing us to call the serotype. In those cases,  
19 those are sent to NVSL to do traditional serology.

20           Okay. This is getting a little familiar to  
21 you at this point. One thing I want to highlight is  
22 how important it is to do the assembly from the FASTQ



1 files, and also to utilize already existing databases  
2 to do your work with.

3           The virulence typing and MLST typing that  
4 we do for STEC we find is very important. We use  
5 that from the Center for Genomic Epidemiology as  
6 well. We also use a virulence finder and serotype  
7 finder which were both mentioned earlier in previous  
8 talks.

9           Another output for that, and then here's an  
10 O26 strain. One of the nice things I like to point  
11 out, too, is previously in FSIS, we didn't get the H  
12 type. It's not part of what we used to do. Well,  
13 now we can get the H type. A lot of them are H11,  
14 H7, whatever they may be.

15           Sequence type is very important to us as  
16 well. Actually I'll give you another example later  
17 on how it helps us kind of predict what the O type  
18 might be.

19           The stx type, that was talked about  
20 earlier. You can actually look down to the stx1a in  
21 this case and then the eae allele, beta in this case  
22 as well.

1           This is an analysis that we just recently  
2 did right before we came. This is around I think 3  
3 or 400 isolates. We looked at our top six non-O157  
4 and O157, and what did we get out of the  
5 characterizations. We looked, what are the top MLST  
6 types? So looking at sequence types, sequencing  
7 allows us to go back any time retrospectively and do  
8 this. It's another very important point about  
9 sequencing. You can always do retrospective  
10 analyses.

11           So you start to see that we do have within  
12 each serogroup, we do have a predominant sequence  
13 type within each one of those. One of the ones I  
14 like to bring attention to is ST11 on the *E. coli*  
15 O157:H7. In our hands, our commodities, our  
16 isolates, those tend to be very familiarly known as  
17 ST11.

18           Various Shiga toxin types, in combination  
19 or alone; eae types, we tend to see gamma as the top  
20 in all of our *E. coli* O157 isolates as well. And our  
21 top serotype as you might expect, O157:H7, not  
22 unexpected.

1           But then we do have within our other  
2 serogroups, we start to see 100% or H11, H8, and so  
3 forth.

4           One of the nice things whole genome  
5 sequencing allows us to do, that previously we  
6 couldn't do, was we can start to identify H types on  
7 serogroup O157 that we normally would have said was  
8 negative for O157. We previously have had some H11s  
9 and H29s. In both of those cases, neither one was  
10 eae positive nor stx positive.

11           Okay. And finally for phylogenetic  
12 comparisons, all the pipelines have been developed by  
13 our public partners. Many are here. Lee Katz's  
14 Lyve-SET, which is on GitHub. The pathogen pipeline  
15 that Bill talked about was a great tool. FDA SNP  
16 pipeline also is on GitHub. And then wgMLST which is  
17 a very nice tool, BioNumerics 7.6 in the CDC  
18 PulseNet, and there's a screenshot.

19           Here's some output from the pathogen  
20 isolate browser. You saw this slide earlier from  
21 Dr. Goldman and about where do some of our type  
22 serotypes in our commodities, where to do they line

1 up as far as within 10 SNPs of a clinical or 20 SNPs  
2 of a clinical or are they in a SNP cluster? And you  
3 can see these numbers and you can see some of the  
4 things like Kentucky, 0% within clinical. Again,  
5 this is all hinging on all these clinicals are being  
6 sequenced in real time. We're kind of ahead in that  
7 we're sequencing everything in real time, and as soon  
8 as we're caught up in the next year with every  
9 clinical going on NCBI, maybe we'll see different  
10 numbers at that point, but that's just something to  
11 consider.

12           This was also brought up by Dave Boxrud.  
13 This is some work we did with pattern 4 in  
14 Enteritidis, very predominant. Fifty to sixty  
15 percent of our Enteritidis isolates are pattern 4.  
16 Well, you can start to see a delineation just looking  
17 at SNPs alone between those 4 pattern clusters which  
18 on your screen have more of a reddish orange tinge to  
19 them. You can start to see that they kind of break  
20 apart. With just PFGE pattern alone they've been  
21 clustered together.

22           Okay. And the last example, I'm almost

1 done here, we did some work with harbors of *Listeria*  
2 strains and some processing environments. That's one  
3 example of one establishment where we had about 14 or  
4 15 isolates collected over a 4-year period. With  
5 this, we used three different pipelines to show the  
6 concordance between those three pipelines and to  
7 look, what did we see here?

8           So we used a SNP pipeline, Lyve-SET and  
9 wgMLST, which is hidden behind the placard. We did  
10 see a very strong concordance between all three of  
11 those pipelines. They did cluster all the same  
12 isolates together. I defined in this one a cluster  
13 of 20 or less.

14           One of the interesting things that you can  
15 see is there are several clades. Those are starting  
16 to break apart, these from here, which it makes sense  
17 because these actually have nothing to do with any of  
18 these events and this establishment right here. So  
19 that's very, very important.

20           The nice part is, with PFGE, all these had  
21 the same primary and secondary PFGE patterns, and we  
22 were actually able to start breaking those apart,

1 teasing those apart using some of the newer tools  
2 that were developed for whole genome sequencing.

3           Okay. And finally, I have one more slide.  
4 How else are we trying to use genotypic data? Again,  
5 this is only looking at genomics, not using  
6 transcriptomics which is always very vital for this  
7 kind of work, but this is locus of heat resistance.  
8 We did a query using certain genes for locus of heat  
9 resistance of all our assemblies that we've had over  
10 the past year or 2 years, and we ultimately found  
11 around 11 *Salmonella* isolates of varying serotypes  
12 that did have this actual locus of heat resistance.

13           So this is a pathway that we're interested  
14 in going into. What else can we find from this data,  
15 heat resistance, acid resistance, those type of  
16 intrinsic components.

17           So, in conclusion, we are focused. We've  
18 definitely invested a lot as an agency in moving  
19 forward with this type of knowledge and working with  
20 you all as our partners.

21           We've built sufficient capacity. We feel  
22 we're well beyond the just capacity building point at

1 this time. We're moving forward. What else can we  
2 do and add?

3 And we want to continue engaging our  
4 national and international partners. We have a lot  
5 of presence in both national and international  
6 meetings, including GMI.

7 And finally, we want to use WGS as we  
8 always have with PFGE, use it in conjunction with all  
9 the epidemiological information that we have and  
10 bring in that totality of evidence to any type of  
11 investigation.

12 With that, I'd like to thank all the people  
13 in the room and thank you for all the people on the  
14 slide. A lot of the collaborations I know for us,  
15 just standing up the program since 2014, took a lot  
16 of collaboration with everyone here and some of the  
17 people that aren't here that were here earlier in the  
18 week for the NARMS meeting. So I'd like to thank  
19 everyone.

20 DR. BRADEN: Well, good afternoon,  
21 everybody. My name is Chris Braden, and I'm from the  
22 Centers for Disease Control. I'm going to be talking

1 a bit about how the different agencies that are  
2 implementing this new technology are actually  
3 coordinating, and I think this is an important  
4 component of what we do.

5           So an approach we've taken in order to  
6 establish this interagency coordination is really to  
7 build upon a history of collaborative programs that  
8 we've had, and some of you here may know of or been  
9 involved with the Interagency Food Safety Analytics  
10 Collaboration that was established before this  
11 particular collaboration and really builds on that  
12 type of structure.

13           Of course, we want to apply advances in  
14 technologies and the one that we've been  
15 concentrating on is next-generation sequencing, but  
16 we want to scan the horizon to see what opportunities  
17 continue to evolve in the technology fields.

18           Of course, even within next-generation  
19 sequencing, as Martin had said before, there's going  
20 to be the next next-generation sequencing and how  
21 that's going to change what we do.

22           We want to leverage the knowledge,



1 expertise and data among agencies. We bring certain  
2 levels of expertise in different areas together, and  
3 that makes us stronger and certainly bringing our  
4 data together makes what we do more effective.

5           And then set up a structure for our  
6 collaboration that is efficient, guided by strategy  
7 and prioritize communications and stakeholder input  
8 over time.

9           So the Interagency Collaboration on  
10 Genomics for Food and Feed Safety, or Gen-FS, was  
11 established in 2015 to strengthen the federal  
12 collaboration on the use of whole genome sequencing  
13 in foodborne pathogen analysis and investigation.

14           Multiple federal agencies are involved and  
15 most recently, we've had ARS and the Animal and Plant  
16 Health Inspection Service, APHIS, at USDA join this  
17 collaboration.

18           So really Gen-FS is meant to support the  
19 implementation of a shared vision of coordinated  
20 networks for genomic sequencing. We want to use  
21 flexible tools and analyses and communications needed  
22 by the respective agencies to harmonize procedures

1 and standards where we should, and really these two  
2 large networks are the ones that we've concentrated  
3 on, making sure that we are, you know, collaborating  
4 and harmonizing as much as we can.

5           So we have some targets for what we are  
6 prioritizing for our development, coordination and  
7 harmonization, including some of the things that  
8 we've already talked about, the system tools and  
9 pipelines and methods; the analytic procedures,  
10 protocols and standards; sharing data and data  
11 availability; harmonizing some of the ways that we  
12 manage the networks with proficiency testing and  
13 training; how we use this data in surveillance,  
14 investigation and research; and then the external  
15 communications and partnerships.

16           So we have a draft charter that is actually  
17 undergoing the process of having our Agency heads  
18 sign it as we speak. We have a steering committee  
19 with representation from each agency and then we've  
20 set up these four workgroups. One's data standards,  
21 analytics, comparisons and interpretation, and I'll  
22 show you some of the output from some of our work,

1 interagency training, the network workflow  
2 harmonization and then communications.

3           So one of the things that we've had as a  
4 strategy and priority from the beginning is really  
5 sharing our work, our structure and our strategy, and  
6 one of the things that we've done from the beginning  
7 is worked to make sure that the DNA sequence and  
8 metadata that we produce are publicly available. So  
9 the data, as you've heard, is uploaded to NCBI and  
10 GenBank. That includes all organisms undergoing  
11 whole genome sequencing in PulseNet and GenomeTrakr.  
12 There are clinical, food and environmental isolates,  
13 and we must do so with the protection of personal and  
14 commercial information.

15           We also make the tools that we use publicly  
16 available either as open source or commercially  
17 available software and then publish our methods and  
18 validation analyses.

19           So the standards and validation, what that  
20 workgroup has really accomplished is establish the  
21 quality standards that are monitored for all  
22 submissions of all genome sequencing to GenBank and

1 to then develop and publish benchmark datasets which  
2 I'll show you a little bit more about in the next  
3 slide.

4           We'll do the validation studies. Some have  
5 already been mentioned having to do with SNPs.  
6 There's more to be done, I think including when we  
7 validate the whole genome MLST and really do some of  
8 the more careful cross comparisons of MLST and SNP  
9 analyses in individual pipelines and then that cross  
10 validation.

11           And then AMR genotype/phenotype comparison,  
12 you've heard quite a lot about already, but there is  
13 another publication out of this workgroup that's  
14 pending.

15           So these benchmark datasets I think as we  
16 were just talking before, is I think a great resource  
17 for anybody. Certainly it's what we need internally  
18 to be able to validate a number of our analyses and  
19 pipelines but being able to then provide them  
20 publicly allows those datasets to be used in  
21 validation and research studies for anybody.

22           So we have five DNA sequence datasets

1 consisting of 10 to 31 well characterized outbreak  
2 and unrelated isolates, that have been developed so  
3 far for *Listeria*, *Campylobacter*, *E. coli*, and  
4 *Salmonella*. And with each of those, there is  
5 publications that detail the outbreaks that they're  
6 associated with and these datasets are available for  
7 download at this GitHub site. So we are making those  
8 publicly available.

9           Harmonization across the networks, of  
10 course, this is important because a lot of the  
11 laboratories are both GenomeTrakr and PulseNet  
12 laboratories and they can't be doing two procedures  
13 for the same purpose. We really try to harmonize our  
14 procedures for the laboratories that participate.

15           So for the training, the training is  
16 provided for public health and regulatory program  
17 partners and PulseNet and GenomeTrakr networks.  
18 They're both CDC and FDA sponsored courses. Staff  
19 from each agency participates in the training as  
20 training faculty and training certification applies  
21 to both networks. So I think we've done a lot to  
22 integrate our training.

1           And the same is true for standard operating  
2 procedures in laboratories for whole genome  
3 sequencing procedures, for sample and library  
4 preparation, for the sequencing procedure itself, for  
5 the data management and upload to NCBI, and then  
6 incorporating new and changing technologies because  
7 for those of you who do this, you know that new  
8 chemistries come out periodically and they need their  
9 own SOPs to change accordingly.

10           And then the proficiency testing, so the  
11 same proficiency panels, the same analysis, the same  
12 reporting is used for both networks.

13           So communications: We've really tried to  
14 be able to communicate what is happening in the  
15 networks, how the data is being used in surveillance,  
16 investigation and regulatory functions. We have some  
17 industry collaborative forums. For instance, the  
18 Institute for Food Safety and Health has had a couple  
19 of meetings, and we've also had some forums in  
20 collaboration with the University of Georgia Center  
21 for Food Safety.

22           There's been many presentations and

1 discussions at food safety and scientific  
2 conferences, and you see some listed there, and then  
3 we are putting together a white paper for publication  
4 on the use of whole genome sequencing in food safety.

5           So looking to the future, and this is what  
6 I'm really concentrating on, how we're going to be  
7 using this information on a day-to-day basis in our  
8 agencies. Whole genome sequence technology will  
9 replace traditional methods for routine microbiologic  
10 characterization of foodborne pathogens for use in  
11 surveillance, investigation, and agency action.

12           Now, people ask, you know, it's still  
13 important to do some traditional microbiology because  
14 that gives you still additional information and, yes,  
15 there are times when that's going to be appropriate,  
16 but for our routine purposes, we think that whole  
17 genome sequencing will give us the information that  
18 we need.

19           Of course, to use these tools in this way,  
20 they have to be validated with comparable results no  
21 matter where the testing is done.

22           We want to have shared tools, standards and

1 data for all stakeholders, so that public health,  
2 regulatory partners can use these standards tools,  
3 but in addition, those in industry, academia and our  
4 international partners can also particular in using  
5 some of these same tools in order to be able to  
6 compare data among partners.

7           And then the other thing that we're really  
8 trying to do is come up with some tools and methods  
9 that will be simplified for WGS analysis because, you  
10 know, to do the research and to develop these does  
11 take a lot of computing power, and it takes a lot of  
12 bioinformatics expertise. Not every health  
13 department or other institution is going to have that  
14 kind of computing power and expertise, and so we do  
15 need to develop tools that are more simplified in  
16 order to have many more participants in these  
17 networks to be able to participate. So in those  
18 cases, there would not be a requirement for high  
19 performance computing and advance bioinformatics  
20 expertise in order to use these tools.

21           And I would like to thank the Genomics and  
22 Food Safety members, and especially the



1 communications workgroup that's worked on this and  
2 other presentations, USDA-FSIS for hosting this  
3 meeting that is meant as a meeting for all the Gen-FS  
4 members, and your interest and input. Thank you very  
5 much.

6 DR. EVANS: Thank you, Dr. Braden.

7 So now we'll have a question and answer  
8 session. So if all of the speakers from this  
9 Federal/State Collaboration session could come up to  
10 the stage, and they can all have the microphone. And  
11 if I could ask the panelists to state your name when  
12 you answer the question for our folks online.

13 Okay. We'll start with a question in the  
14 room.

15 DR. BOOREN: Great. Thank you. Betsy  
16 Booren with OFW Law. First of all, thank you all for  
17 your time today. I found the presentations  
18 fascinating and as I've been sitting here, and we've  
19 been discussing some at break, sort of the totality  
20 of what's been presented, and I have a question, and  
21 I'm not sure who's the best person to ask it, and it  
22 may be coming from a place of ignorance.

1           But as I'm looking and hearing what is  
2 being talked about, the number of isolates that have  
3 been sequenced, one of the questions that I have in  
4 my head is what, and perhaps it should be focused on  
5 GenomeTrakr, of the pie of isolates that have been  
6 sequenced, how much of that is clinical samples? How  
7 much of that is other isolate sources? And how much  
8 of that is regulatory? Because as someone who  
9 represents the food industry, I'm focused on the  
10 regulatory and trying to better understand, if I see  
11 a slide that says, papayas or meat and poultry have a  
12 high sense, what does that mean in the whole scheme  
13 of things.

14           And so what I'm trying to get at is for  
15 industry to look at trying to do research or other  
16 groups, in that regulated area, how much of those are  
17 regulatory surveillance? How much of those are for  
18 cause?

19           And as I talk with my industry partners,  
20 better understanding some of that information may  
21 inform research decisions that can help when they say  
22 why is an industry adding to the database? We want

1 it to be targeted and want to understand what's going  
2 on in those isolates.

3 I'm not sure whose got the right answer or  
4 if that's -- again, maybe there's a report out there  
5 that I haven't seen, but I think it's really  
6 important as when we've done other risk assessment  
7 type of work, certain agencies, certain industries  
8 have an abundance of data compared to others, and  
9 does that mean there's a true risk there? What does  
10 that data tell us? So I'd be curious in your  
11 thoughts as we try to, particularly from an industry  
12 standpoint, better understand what that data means.

13 DR. BRADEN: Thank you for your question.  
14 My name is Chris Braden from CDC. I think you have  
15 to be careful in the database to know that this is  
16 not a statistically representative sampling scheme.  
17 It does contain, you know, especially on the food and  
18 environmental side, those samples were collected for  
19 lots of different reasons and actually I'm not sure  
20 that we have the data to break out what sampling for  
21 cause and what sampling according to some assignment  
22 or routine sampling that might be done.

1           So that I think is one thing that you need  
2 to understand when you're analyzing that as a  
3 limitation to the data.

4           There are more food and environmental  
5 isolates than there are case isolates. I think there  
6 always will be because food and environmental  
7 isolates just depend on how much you test and there's  
8 only so many cases. So that's going to be the case  
9 going forward.

10           But nonetheless, even with those  
11 limitations, I think that there's lots that one can  
12 learn about the breadth of genomic variation in any  
13 number of ways in these databases that is helpful to  
14 answer some questions but won't be able to say, to  
15 break it down by, for instance, sampling for cause or  
16 for surveillance.

17           DR. TILLMAN: This is Glenn Tillman from  
18 FSIS as well. To follow up, one of Dr. Goldman's  
19 slides had that we've sequenced around 11,000 uploads  
20 into NCBI since inception. Around half of those are  
21 part of our non-regulatory cecal program, and those  
22 have been kind of moved to a different bio project

1 which could provide some opportunity to look at them,  
2 those strictly within that bio project, and one of my  
3 slides had all of the ones that we contribute to, and  
4 they're actually called a NARMS cecal.

5           The other 5,500 were collected as part of a  
6 regulatory initiative. So we do have kind of a half  
7 mix there, and again the NARMS cecal are moved into  
8 their own bio project. So that would be a good place  
9 to start on understanding what each one of those bio  
10 projects that NCBI might have within it.

11           DR. BESSER: John Besser from CDC. It  
12 hasn't come up today, but there's actually something  
13 called VoluntaryNet at University of Georgia which  
14 could be used by industry to anonymize the sequence  
15 data, and actually as a firewall between public  
16 health and VoluntaryNet. So it could be used to  
17 compare the food and environmental isolates submitted  
18 by industry to clinical cases to assess risk.

19           But we can't, the CDC, FDA, USDA, can't  
20 specifically query that database. We can ask them if  
21 there's matches. They can then ask their membership  
22 as to whether or not whoever submitted it would be

1 interested in sharing that data, but the idea is that  
2 there's concern in industry that the information  
3 would be used in a punitive manner and volunteering  
4 that is a safe place to submit data that can be used  
5 for risk assessments. It's not exactly answering  
6 your question, but it's a related topic.

7 DR. EVANS: Jorgen.

8 DR. SCHLUNDT: My name is Jorgen Schlundt,  
9 and I'm from Singapore. I love saying that I'm from  
10 Singapore. Okay. Especially with a Danish accent.

11 But it's fantastic to be in the U.S.  
12 discussing something like next-generation sequencing  
13 because this is clearly an area where U.S. is  
14 leading. I understand that in some other areas, U.S.  
15 is stepping back from leading, but actually in this  
16 area, you are really leading. We don't need to  
17 discuss the integration between the different  
18 agencies in U.S. Maybe it's a good thing that you  
19 have many different agencies because then you can  
20 move in different speed, and we have seen that also  
21 today.

22 But I have one question in relation to

1 maybe unifying the forces not only within the U.S.,  
2 but also with the rest of the world. So I was  
3 looking at resistance, and I understand that CVM is  
4 moving forward with a Resistome Tracker, and I  
5 understand that NCBI is move forward with something  
6 else, with very long name that I didn't really get.  
7 And I understand that FSIS is using the WHO  
8 collaborating center thing from Denmark, the  
9 ResFinder.

10           Wouldn't it make sense if there was sort of  
11 a concerted effort to try to produce consistent,  
12 uniform methodology and also tools so that the rest  
13 of the world could really get a benefit out of U.S.  
14 leading in this area?

15           DR. McDERMOTT: Thanks, Jorgen, for the  
16 question. And it's a topic we discuss quite a lot.  
17 The need for harmonization in allele cause is  
18 essential. If we're going to take seriously  
19 antibiotic resistance as a global challenge, it must  
20 be addressed globally. We need a harmonized method  
21 just to have a common language.

22           So I think that the discrepancies are being

1 worked through between say ResFinder, NCBI, and Bill  
2 can tell us the latest on that, but I think just in  
3 general, it's essential that it be done so the  
4 outputs are the same. There might be tools with  
5 different interfaces that work with BioNumerics, for  
6 example, or different types of alleles, but as long  
7 as the output is either the same or can readily be  
8 translated into a common language, then that will be  
9 an important objective, but I do think that it's  
10 essential.

11 DR. KLIMKE: So, Jorgen, we were just at  
12 the ASM Genomics conference, and we had a roundtable  
13 with NCBI, the CARD database in Canada, and Ole(ph.)  
14 from DTU was there. Although he's technically not  
15 ResFinder, he's at least Danish, which is close  
16 enough. And we agreed that we would at least discuss  
17 harmonization or curation efforts because the content  
18 in the database should be the same. The methods to  
19 apply them I think will have to follow after that  
20 because we'll have to do comparisons to see if we get  
21 similar answers and the extent of application of  
22 related sequences is something we need to look at



1 but, yes, I would agree with what you're saying.

2 DR. BESSER: I appreciate your Danish  
3 honesty. That's great. You know, I agree with  
4 everything that's been said that we're moving towards  
5 that. I think what Chris Braden presented on the  
6 data quality is the really key harmonization point.  
7 If we share data quality, we can explore these  
8 different mechanisms for allele database curation for  
9 different systems.

10 And as with any new technology, there has  
11 to be sort of a creative period, and I think we're in  
12 that now, and so no one system has emerged yet as the  
13 winner. But as long as we share this core of data  
14 quality, I feel that we're all moving towards global  
15 standardization and curation. We just haven't gotten  
16 there yet, but I think we all agree that that's the  
17 direction we're moving.

18 DR. BRADEN: This is Chris Braden again,  
19 and you raised the issue having to do with resistance  
20 determinant databases, but I think the same can be  
21 said for other allele databases that we all should be  
22 coordinating and the types of databases we use, even

1 if the management is distributed in some way.

2 DR. BESSER: Thanks, Chris. At least CDC,  
3 we're envisioning the global curation of these allele  
4 databases for subtyping purposes. That's a really  
5 big issue, but the Pasteur Institute in France has  
6 done curation for *Salmonella* serotyping forever.  
7 We're envisioning a global collaborative curation of  
8 some sort, but I think it has to happen on a global  
9 level ultimately.

10 DR. EVANS: We have a question from the  
11 web.

12 DR. ALVARADO: The question from the online  
13 group is do you worry about the education training  
14 that may be required so that people can adequately  
15 interpret WGS data findings?

16 DR. BRADEN: So this is Chris Braden. Is  
17 the question -- is this a priority of some -- could  
18 you repeat please?

19 DR. ALVARADO: Sure. The question is do  
20 you worry about the education training that may be  
21 required so that people can adequately interpret the  
22 WGS data findings?

1 DR. BRADEN: So, yes, we do worry about  
2 that. It is difficult. Not only is it difficult to  
3 learn some of the bioinformatics that may be  
4 necessary to do this type of analysis in order to,  
5 you know, appropriately interpret it in a laboratory,  
6 but there's a whole other realm of disciplines out  
7 there that also needs some training on how to  
8 interpret this data and so probably the next outreach  
9 is with the epidemiologists really to do that  
10 genomics for epidemiologist type of training and  
11 there are a number of courses and I know Martin  
12 Wiedmann was here before and he's led some of those  
13 courses. And so that certainly is the case, but I  
14 think as we use this more, it's going to be, you  
15 know, there's going to have to be more training, you  
16 know, with more disciplines for instance in the  
17 regulatory agencies or with risk assessors and so  
18 forth.

19 MR. BOXRUD: So Dave Boxrud from the  
20 Minnesota Department of Health. So from a public  
21 health lab standpoint, we have additional challenges.  
22 We're a smaller group, and traditionally we've had

1 experts in microbiology identification, that sort of  
2 thing. Many of our staff simply don't have the  
3 background in sequencing or have the understanding of  
4 how to interpret sequencing and how it works.

5           We've been able to hire a number of new  
6 people that have that background but trying to bridge  
7 that knowledge gap between traditional subject matter  
8 experts and then this new generation sequencing  
9 group, there's a real challenge, and we have to work  
10 on it constantly to try to bridge that gap, and I  
11 agree with Chris that educating our epidemiologists  
12 is also really a vital thing so that everyone is  
13 always on the same page.

14           DR. TILLMAN: I just want to add one quick  
15 thing. This is Glenn Tillman. So a lot of the tools  
16 have really come a long way since the 2015 ASM  
17 Genomic Pipeline meeting in Washington. A lot of the  
18 online tools have come about and really moved things  
19 forward, BioNumerics and with the applied math has  
20 really moved forward and really made things a lot  
21 more streamline, allowing users to really be able to  
22 use that stuff. So I think training is a very big

1 part. I think the development of these newer  
2 analytical tools in more of a web-based form and/or  
3 just any type of GUI-based form have really been  
4 beneficial.

5 DR. ALLARD: Marc Allard, FDA. Once again  
6 I'm going to do a comment and then another question.

7 I just want to say that many people in the  
8 room are actively involved in international and  
9 global training, as well writing white paper  
10 documents on why you should adopt whole genome  
11 sequencing and the specific case studies that they're  
12 best used.

13 And so this is collaborations where the  
14 World Health Organization, Food and Agricultural  
15 Organization, and the OIE, and so there's ongoing  
16 efforts for global training because of global trade  
17 in food.

18 So my question is, it's a little off the  
19 topic while we're here at USDA, but I know that the  
20 Department of Health for New York, they don't just  
21 sequence the four or five foodborne pathogens.  
22 They're actually responsible and have been sequencing

1 efforts in almost 20 human pathogens. So this  
2 question I think starts with Dave Boxrud at Minnesota  
3 is, do you see a similar expansion of sequencing for  
4 other species? And what do you see as growth? And I  
5 guess then the question is, you know, are you seeing  
6 equivalent support from your federal partners, both  
7 in databasing, etc.?

8 MR. BOXRUD: Yeah, it's a great question.  
9 But PulseNet and Foodborne is definitely in many ways  
10 the leader in sequencing. They're kind of the bull  
11 because there's so much of it, but we are involved  
12 with a lot of different types of pathogen testing,  
13 *Legionella*, *Strep pneumo*, MRSA. I believe it's about  
14 10 different projects that we're involved with, and  
15 many of these, we've reached out with our  
16 collaborators at CDC and said we're interested in  
17 this organism. We have background in this. Let's  
18 work together and they've been very, very willing to  
19 do that. And that's been really key for us.

20 One of the real challenges with using  
21 sequencing with foodborne disease surveillance is  
22 while it's an awesome technique, it's an awesome

1 method, if you batch it and do it once a month,  
2 you're really slowing down the process.

3           And so we fill in our runs with other  
4 organisms, with other pathogens. The *Legionella*  
5 group from CDC has a really nice partnership, and  
6 they have created -- I think there's about six  
7 laboratories that are doing a study with them. The  
8 TB group is doing, from CDC, I believe they funded  
9 one site to do essentially all or most of the  
10 clinical TB cases throughout the country. There's a  
11 flu group, we're not involved with this, but there's  
12 three sites that are doing a lot of influenza  
13 testing.

14           So there is a lot going on but some of the  
15 public health labs are not always embracing the new  
16 technology for different organisms. So I think  
17 hopefully as sequencing becomes more routine at all  
18 public health laboratories, that they will consider  
19 doing other pathogens because I do think, you know,  
20 this is a technology change that is not going to go  
21 away and it's going to continue to provide additional  
22 important information for different pathogens.

1 DR. EVANS: Question from the room.

2 DR. WIEDMANN: Martin Wiedmann, Cornell  
3 University. So we talked a lot about standardization  
4 of database as alleles cause, quality control. Where  
5 are we with regard to standardization of metadata?  
6 And I'm going to start out with just even food data.  
7 I think porcine, pork, different names, to call the  
8 same thing. If we move to environmental, you know,  
9 data, how do you describe environmental sample? It  
10 gets even more challenging.

11 And I think there's an importance to it  
12 that hasn't been mentioned. I think when we try to  
13 interpret SNP differences, it obviously depends on  
14 the environment the organism is in. If *Salmonella* is  
15 in a dry environment, it might only replicate 10  
16 times a years, and therefore accumulates SNPs at a  
17 much, much slower rate than if the same *Salmonella*  
18 sits in a poultry house where it might replicate  
19 every hour because it's actively infecting one  
20 chicken after the other. And so if you don't have  
21 that information in some sort of standardized way, we  
22 sometimes will run into challenges how we interpret



1 SNP differences.

2           So is there parallel effort as part of Gen-  
3 FS to also standardize those data or is that sort of  
4 the next step?

5           DR. KLIMKE: I can only tell you about the  
6 things I know, and that is that the global microbial  
7 identifier and part of a ISO working group on next-  
8 generation sequencing for food safety, that there is  
9 an effort to use some metadata ontologies to  
10 standardize. This is called Genna GO (ph.). It's  
11 been worked on by the Canadians. CFSAN has said they  
12 will look at using that in their metadata, but I  
13 should tell you that what we would rely on then is  
14 probably the submitters to apply those terms in a  
15 standard way because a lot of people will say that  
16 ontologies are the solution to everything. They're  
17 not, if they're not applied in a standardized way by  
18 all the people contributing.

19           Since we are the people who are integrating  
20 all the data, if everyone in the U.S. goes a certain  
21 way and then the Europeans do it differently, we're  
22 sort of left at the mercy of that, where we would

1 have to either invest time and effort on our side to  
2 standardize the metadata once it gets indicated and  
3 submitted to us, or we just leave it alone and only  
4 worry about the standardized metadata in the U.S.

5           So I know someone from CFSAN can probably  
6 mention that they're looking at this. You could  
7 probably talk to them. I don't know -- if USDA's  
8 doing that, they probably will. You guys can  
9 probably talk about that.

10           DR. EVANS: We got a question from the web.

11           DR. BRADEN: Chris Braden again. You know,  
12 we started by trying to standardize what the metadata  
13 fields would be in the Gen-FS agencies, and it's a  
14 place to start, but I think definitely needs to  
15 expand to be able to have all partners be able to  
16 participate with some standard expectations.

17           DR. EVANS: Okay. We have a question from  
18 the web.

19           DR. ALVARADO: Can you say more about  
20 providing safe harbor for producers, manufacturers,  
21 others in the industry, to (a) gain familiarity with  
22 WGS and (b) provide useful information to regulators

1 without bringing regulatory response down on their  
2 heads?

3 DR. BRADEN: This is Chris Braden. I'm not  
4 sure we have the representation here, but I don't  
5 think I really understood the question. If you could  
6 ask it again.

7 DR. ALVARADO: Sure, I can repeat the  
8 question. Can you say more about providing safe  
9 harbor for producers, manufacturers and others in the  
10 industry (a) to gain familiarity with WGS and to  
11 provide useful information to regulators without  
12 bringing regulatory response down on their heads?

13 DR. BRADEN: So, yes, safe harbor. I think  
14 that's where is that safe place to be able to start  
15 to implement this kind of technology in order to  
16 learn more about your particular producing  
17 environments if you're in the industry, for instance.  
18 I think it is an important point. We have had a  
19 number of discussions with industry members about  
20 what might work if there's a third party that could  
21 be responsible for holding the key, for instance, and  
22 not releasing any of the identifiable information and

1 datasets.

2 I would hope that we can find such  
3 partnerships to be able to do so. John Besser had  
4 mentioned one at the University of Georgia called  
5 VoluntaryNet which is an industry collaboration. I  
6 know that at the Institute for Food Safety and  
7 Health, they are also considering being able to  
8 provide that kind of third party resource. There may  
9 be other resources out there, and I think, for  
10 instance, you know, IEH has said that, you know, they  
11 can certainly provide the service and provide that  
12 kind of data back that would be helpful to industry  
13 members to understand how to use this information.

14 So I think there are some resources out  
15 there. It hasn't been widespread to my knowledge.  
16 Maybe IEH is maybe used the most but, you know, I  
17 think it's certainly worth exploring for industry  
18 members.

19 DR. ALLARD: Marc Allard, FDA. I can  
20 comment just a little bit about this. We've been  
21 doing outreach to industry. FDA primarily works with  
22 the Institute for Food Safety and Health, IFSH, which

1 is out at Illinois Institute of Technology in the  
2 Moffett Center in Chicago. And essentially from 3 or  
3 4 years speaking to food industry, probably the most  
4 feasible path to learn new things about this  
5 technology is to work with a third party provider,  
6 essentially do this on line.

7           Steve Musser touched on this earlier today.  
8 There's many, many companies that are involved in  
9 this, as well as many academic folks like Andy Benson  
10 at University of Nebraska and Martin Wiedmann at  
11 Cornell, but there's companies like Eurofins, IEH,  
12 Ecolabs, NSF International. You can find these on  
13 the internet that provide genomic services and can  
14 assist you. This is an expanding market because  
15 there are a lot of people in industry that would  
16 prefer to do it as a third party provider as opposed  
17 to building their own laboratory.

18           So we can give you more information or  
19 reach out to your local lobbying group of GMA, United  
20 Fresh, SQF, IAFP. There's many speakers at these  
21 meetings that have a diversity of services to the  
22 food industry.

1 DR. EVANS: There's a question from the  
2 web.

3 DR. ALVARADO: There was no mention in  
4 today's talks by the different agencies on how WGS is  
5 being used in the detection of outbreaks and for  
6 carrying out outbreak investigations. How is the CDC  
7 using the information and how is the integration with  
8 what the FDA, FSIS, might have in their database  
9 taking place?

10 DR. BRADEN: So as I mentioned in my  
11 presentation, basically we will be transitioning on  
12 all of our traditional characterization techniques  
13 over to whole genome sequencing. So whole genome  
14 sequencing will be replacing pulsed-field gel  
15 electrophoresis in our cluster detection, and it will  
16 be replacing the assays that we use for isolate  
17 confirmation, serotyping, virulence typing, and  
18 resistance typing.

19 So all those separate assays will now be  
20 carried out with one assay using whole genome  
21 sequencing in the future. We're in the midst of that  
22 transition now. As David Boxrud had said, that it's

1 a difficult time because we're actually doing both  
2 traditional typing and whole genome sequencing and  
3 that makes it more expensive and more time consuming  
4 but we anticipate that in the long run, this will be  
5 time and cost saving.

6           So that's how we're looking at moving  
7 forward, and I think that I can talk for the other  
8 agencies to say that that's their plan, too.

9           As far as, you know, how we use this data  
10 together, well, we do talk a little bit about  
11 GenomeTrakr database and PulseNet database and so  
12 forth, but in reality, we're all submitting all of  
13 the data into a single database and we're just  
14 contributors. So PulseNet is contributing to a  
15 single database. GenomeTrakr is contributing to a  
16 single database. And we're all using all of that  
17 data in our analyses for outbreak detection,  
18 investigation, surveillance, and action. So that's  
19 how it comes together.

20           DR. EVANS: So I have a question. John  
21 mentioned that by the end of fiscal '18, I believe,  
22 that there would be all or close to all *Salmonella*

1 would be sequenced, and I'm just curious about any  
2 logistical issues that would be raised by that, and  
3 what do you expect at the end of fiscal '18 when we  
4 have all these new *Salmonella* sequencing real time?

5 DR. BESSER: I'm John Besser. That's a  
6 very good question, Peter. Yes, it's not going to be  
7 smooth sailing. We're really getting into a whole  
8 big data era. I think the experience of Public  
9 Health England would probably be helpful. They saw a  
10 dramatic increase in the number of clusters that  
11 needed to be investigated. Now, we saw this with our  
12 *Listeria* combined initiative as well.

13 So I think there's going to be a need to  
14 develop new tools for cluster triage, what clusters  
15 are most likely to be productive but I think there's  
16 going to be more -- there already are more clusters  
17 to investigate than there are investigators to  
18 investigate them, sorry for the convoluted sentence  
19 there, but I think this is going to become a major  
20 issue and the investigative resources are going to be  
21 a problem.

22 I think we're working to streamline on the



1 laboratory side to make these things easier. For  
2 instance, you heard about nomenclature that's being  
3 developed which will allow for ready recognition,  
4 easier recognition of clusters and easier  
5 communication, but I think a lot more work needs to  
6 be done in order to prepare for this big data era  
7 that's rapidly approaching.

8 DR. EVANS: There's a question in the room.

9 MS. MCGARRY: Sherri McGarry, CDC, and the  
10 question is mostly for David of the Minnesota  
11 Department of Health, but I'd welcome other panelists  
12 to weigh in. And it piggybacks somewhat what you  
13 were just talking about. So when PFGE first came on  
14 board, it took longer than it takes now, right? So  
15 there are efficiencies that were gained and  
16 techniques that were modified to make it faster.  
17 Where do you see some innovation at the state level  
18 to kind of speed things up a little bit? I know  
19 we're still at the early phase, but maybe this is  
20 also an innovative phase, too. So where do you see  
21 efficiencies to reduce the total amount of time?

22 MR. BOXRUD: Thanks. David Boxrud from

1 Minnesota. Yeah, there's a lot of ways to make the  
2 process more efficient, and that's something as we  
3 continue to do more sequencing, we're going to have  
4 to continue to look at these ways. We're going to  
5 have to adapt to new technologies. You know, if the  
6 work that we're doing with Illumina MiSeq expired 5  
7 years from now, we may have a completely different  
8 technology than what we have.

9           For us, as we're really ramping up the  
10 amount of sequencing that we're doing, it's a lot of  
11 training, a lot of getting everyone on the same page  
12 within our laboratory, and as we're doing that, we're  
13 okay with a little bit of extra time to make sure  
14 that our quality is correct, and that we don't have  
15 any issues.

16           Once we really get this incorporated into a  
17 standardized workflow, then we're going to see what  
18 we can do to try to make it more efficient. We're in  
19 the process of getting a library prep instrument that  
20 will take some of the time and the labor of creating  
21 the library prep which is one of the more labor  
22 intensive parts of the process.

1           And also the kits that we're using, the  
2 technology is changing. The MinION could be a faster  
3 tool than what we have right now, and right now for  
4 PulseNet and various other things, it's probably not  
5 ready for prime time, but very soon it will be.

6           So I think all of us are going to have to  
7 continue to adapt, but what is great about sequencing  
8 is if you have that quality of sequence, if you use a  
9 new technology, they're very comparable going  
10 forward.

11           DR. BESSER: That's a good answer, Dave.  
12 This is John Besser, and I think Dave already  
13 mentioned the batching issue which is a big potential  
14 problem in turnaround time, but as he mentioned,  
15 there's also new technological developments, a new  
16 library prep and DNA prep, chemistry that just came  
17 out recently, and we'll be looking at all of those.

18           But I think what I didn't have time to talk  
19 about earlier today was some of the new tests. You  
20 asked about innovation. One in particular, amplicon  
21 sequencing, we're looking at a quasi whole genome  
22 MLST that can be done directly from a specimen, and

1 that actually could shave weeks off of the whole  
2 process if successful, but it's still highly  
3 experimental at this phase, but I think ultimately  
4 where we want to go towards is direct specimen  
5 testing which will have the biggest overall impact on  
6 turnaround time because the culture is the slowest  
7 part. Often culture happens in a clinical size at  
8 the state and sometimes even at the CDC and that can  
9 really slow down the whole process and we're hoping  
10 to bypass that at some point in the near future.

11 DR. EVANS: We have a question in the room.

12 DR. GOLDMAN: Yeah, David Goldman, FSIS. I  
13 just want to go back to the question that was asked  
14 online about outbreak investigation, and when the  
15 question was asked, I realized we really hadn't done  
16 a recent case study of an outbreak where whole genome  
17 sequencing was used.

18 And I think, you know, everyone knows that  
19 we rode PulseNet with PFGE very successfully, highly  
20 successful for 20 years and, you know, if two  
21 patterns were indistinguishable, we said, okay.  
22 Everyone agreed with that. We even dealt with one

1 band differences very well I think.

2           But now with whole genome sequencing, the  
3 picture is less clear as you've heard many speakers  
4 attest earlier, and I want to go back and emphasize  
5 something Martin Wiedmann said earlier which is that  
6 to me, where we are now is that the epi is more  
7 important than it was using pulsed field analysis as  
8 a way of judging and ultimately including cases in  
9 the case definition.

10           I'll just briefly reference a very recent  
11 outbreak in which we used both PFGE and whole genome  
12 sequencing and in this particular instance, the PFGE  
13 was done first by most of the state labs who were  
14 members of this outbreak, affected states, and what  
15 we found out was that after the PFGE seemed to  
16 include cases, a week or so later, we get this  
17 genomic sequence information which would exclude  
18 those cases.

19           So there was this sort of sequential  
20 activity which proved quite challenging for us as a  
21 federal family with the state partners in trying to  
22 determine whether this outbreak was over, have we

1 fallen below the epidemic threshold or are we still  
2 above it? And so this was a very recent example, and  
3 I think while we're in this transition phase still,  
4 using both tools, PFGE and whole genome sequencing,  
5 we may find similar challenges going forward.

6 MR. HEINZELMANN: All right. Joe  
7 Heinzelmann from Neogen. One of the things I had  
8 hoped would be addressed in some of the talks today,  
9 and maybe it will be tomorrow, is around the question  
10 of the statute of limitations around isolates in a  
11 database. Specifically let's say you find an isolate  
12 in a facility, you take corrective actions. Does  
13 that absolve you from that type of isolate in the  
14 database? Are there things that people can do or use  
15 this database?

16 So I guess what I'm really trying to get to  
17 is what can people do with whole genome sequencing  
18 data once an isolate is in the database from a  
19 facility to say that they've made corrective actions  
20 or they've used this technology and how long is that  
21 data point still real and applicable throughout the  
22 life of the food?

1 DR. BRADEN: Yeah. So I'm not sure if we  
2 have somebody in the regulatory community that can  
3 address that.

4 DR. ALLARD: Asking sort of policy related  
5 questions is always the third rail of the FDA, but  
6 I'll just make a comment. And so the emphasis is  
7 that a genetic match is not a regulatory decision. A  
8 genetic match is an indication of a shared common  
9 ancestry, shared isolates and it's a presumed link, a  
10 potential link to a food and clinical, whatever gets  
11 matched, but inspection in epidemiology, you have to  
12 follow and, in fact, FDA will not do any regulatory  
13 action on genetic match alone. There must be an  
14 investigation.

15 So we've seen 5-year-old matches, 10-year-  
16 old matches. Essentially this leads us back to a  
17 facility, to a region, a country, a state, but it  
18 depends what an investigator finds and it depends how  
19 the company responds. There's a whole process that  
20 hasn't changed. We just have a new genetic tool that  
21 helps establish linkage.

22 And so if investigators go and inspect a

1 facility and it's clean as a whistle, then there's  
2 nothing to be done. There's no regulatory activity.  
3 If the inspectors come and they found positives for  
4 foodborne pathogens, then it's like what Dr. Musser  
5 said. It depends on whether there's an association  
6 with clinical or not, and so it depends.

7           But the databases last forever, or at least  
8 30 years. I know I have data in the database that's  
9 been there for 30 years. So our full expectation is  
10 NCBI will not go away anytime soon, and that we'll  
11 continue to see and use these linkages.

12           DR. BESSER: This is John Besser. Marc,  
13 thank you for that response. I think that was an  
14 excellent response.

15           I just wanted to suggest that we might,  
16 because this technology is new, we're suddenly able  
17 with more specificity to connect current cases with  
18 past isolates in the database, that this may be less  
19 of a problem in the future because we'll be able to  
20 connect cases more or less in real time with isolates  
21 as they're being found in the environment.

22           I think this is a circumstance where we've



1 got the new meeting the old, and I think it may not  
2 be as much of a problem as people are concerned  
3 about.

4           And as Marc pointed out, this activity is  
5 not new. We've always compared to the historical  
6 database. What's new is the specificity with which  
7 we can make that connection.

8           DR. EVANS: Okay. So this is the last  
9 chance for questions on the web, in the room.

10           And seeing no questions, we're going to  
11 start up again tomorrow morning at 8:00 a.m. The  
12 first set of presentations will be on what's  
13 happening with our international organizations,  
14 Mexico, Canada, and also globally, and then we'll  
15 have some presentations by our stakeholders, the food  
16 industry and other stakeholders.

17           So I'm really looking forward to the  
18 presentations tomorrow, and I hope to see everybody  
19 here and again, we'll be starting at 8:00 a.m.  
20 You'll enter through the fifth wing, and if there are  
21 no other additions, then we'll break until tomorrow  
22 morning.

1 Thank you.

2 (Whereupon, at 4:17 p.m., the meeting in  
3 the above-entitled matter was continued, to resume  
4 the next day, Friday, October 27, 2017, at 8:00 a.m.)

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

C E R T I F I C A T E

This is to certify that the attached proceedings  
in the matter of:

USE OF WHOLE GENOME SEQUENCE (WGS) ANALYSIS  
TO IMPROVE FOOD SAFETY AND PUBLIC HEALTH

Washington, D.C.

October 26, 2017

were held as herein appears, and that this is the  
original transcription thereof for the files of the  
U.S. Department of Agriculture.



TOM BOWMAN

Official Reporter