

1 Attendance:

2 Toby Milgrom Levin, Senior Advisor, Privacy Office, U.S.

3 Department of Homeland Security

4 Hugo Teufel, III, Chief Privacy Officer, U.S.

5 Department of Homeland Security

6 Jay M. Cohen, Under Secretary, Science and Technology

7 Directorate, U.S. Department of Homeland Security

8 David Jensen, Associate Professor of Computer Science,

9 University of Massachusetts

10 Martha Landesberg, Senior Privacy Analyst, Privacy

11 Office, U.S. Department of Homeland Security

12 Fred H. Cate, Distinguished Professor and Director of

13 the Center for Applied Cybersecurity Research, Indiana

14 University

15 Greg Nojeim, Senior Counsel and Director, Project on

16 Freedom, Security and Technology, Center for Democracy and

17 Technology

18 Christopher Slobogin, Milton Underwood Professor of

19 Law, Vanderbilt University Law School

20 Barry Steinhardt, Director, ACLU Program on Technology

21 and Liberty

22 Peter Swire, C. William O'Neill Professor of Law,

23 Moritz College of Law, Ohio State University

1 John Hoyt, Chief, Knowledge Management Tools Branch,
2 Command, Control, and Interoperability Division, DHS
3 Science and Technology Directorate

4 Stephen Coggeshall, Chief Technology Officer, ID
5 Analytics, Inc.

6 Stephen Dennis, Technical Director, Homeland Security
7 Advanced Research Projects Agency, DHS Science and
8 Technology Directorate

9 Chris Clifton, Associate Professor of Computer
10 Science, Purdue University

11 Dr. Anant Jhingran, Vice President and Chief
12 Technology Officer, IBM Information Management Division

13 Rebecca Wright, Associate Professor of Computer
14 Science and Deputy Director of DIMACS, Rutgers University

15

16

17

18

19

20

21

22

1 P R O C E E D I N G S

2 [Convened at 8:39 a.m.]

3 Ms. Levin: Good morning. The Department of
4 Homeland Security Privacy Office is pleased to welcome you
5 to our workshop, "Implementing Privacy Protections in
6 Government Data Mining." I especially want to thank all of
7 you who have traveled from far, and I understand a number
8 of you had delays yesterday as the result of the storm, but
9 hopefully everyone who wanted to attend has been able to
10 make it.

11 My name is Toby Levin, I'm Senior Advisor in the
12 DHS Privacy Office, and I'm co-coordinator with my
13 colleague, Martha Landesberg, who you'll meet shortly, for
14 this workshop.

15 Before I introduce our welcoming speakers, I have
16 just a few housekeeping announcements to make. First, you
17 should have a packet for the workshop which includes the
18 agenda and the bios -- we will not be doing biographical
19 introductions -- as well as copies of some of the key
20 slides from the presentations that you'll be seeing for
21 today and tomorrow. We will post a transcript of the
22 workshop on our workshop website at www.dhs.gov/privacy,

1 hopefully by mid-August. In order to enable additional
2 comments and so that you can perhaps include responses,
3 reactions to what you hear throughout the workshop, we are
4 going to be extending the comment deadline to August the
5 Fifteenth; comments instructions are on our website and we
6 look forward to your additional input.

7 I want to apologize that we're not able to
8 provide refreshments, but due to our ethics rules, we're
9 not allowed to use your tax dollars to fund refreshments
10 for the workshop. But there are coffee and other
11 refreshments across from the auditorium. Feel free to
12 use those during the breaks but please return promptly.

13 We'll break for lunch about 11:45 and resume at
14 1. In addition to the dining options that are located on
15 this floor and upstairs in the hotel, you have a list of
16 eateries in your packet.

17 After our welcoming speakers we'll move directly
18 to our program; we've set aside the last fifteen minutes
19 of each panel for you to ask questions, and there is a mic
20 up front where you can line up when you're told, queued to
21 line up so that we can hear from you and your questions and
22 any input that you would like to provide. Make sure that

1 you identify yourself by name and affiliation, if any, so
2 that we can have an accurate transcript.

3 Martha Landesberg and I want to thank our privacy
4 team who helped in preparation of the workshop,
5 particularly Sandra Debnam, Sandra Hawkins, Rachel Drucker,
6 Richard Moore, and the rest of our Privacy staff who are
7 here today.

8 And finally, if you would please silence your
9 cell phones so that we won't have interruptions, I think
10 we're ready to begin. It's my pleasure to introduce my
11 leader, Hugo Teufel, Chief Privacy Officer of the
12 Department of Homeland Security.

13 PRESENTATION OF HUGO TEUFEL, CHIEF PRIVACY
14 OFFICER, U.S. DEPARTMENT OF HOMELAND SECURITY.

15 Mr. Teufel: Good morning. I'm Hugo Teufel,
16 Chief Privacy Officer at the Department of Homeland
17 Security; and I have a few remarks before we have our guest
18 speaker who will be joining me up here in a minute. And I
19 see our colleagues from the Government Accountability
20 Office are here, so, yes of course we comply with the
21 ethics requirements and appropriations laws. There will be
22 no free lunches or snacks or coffee or tea.

1 Well, it's my pleasure to welcome you all to
2 this, our fifth, in a series of workshops over the
3 existence of the Privacy Office at the Department. Our
4 goal for our workshops has been to educate the public,
5 educate our office, educate the Department and others in
6 government on cutting-edge privacy issues, and today's
7 workshop, "Implementing Privacy Protections in Government
8 Data Mining" should be no exception.

9 We're fortunate to have with us today and
10 tomorrow some of the most prominent experts in the field,
11 both with respect to privacy as well as with respect to
12 technology, coming to talk to us and to you about the
13 subject matter of this workshop. And I'm really excited
14 about it; I've got to tell you, though, that I'll be
15 popping in and out today and tomorrow because of some
16 unexpected meetings up at the Nebraska Avenue Complex.

17 We're particularly pleased that the Under
18 Secretary for the Science and Technology Directorate -- my
19 friend Jay Cohen -- will be here to help open this
20 workshop. And then also, computer scientists from the DHS
21 Science and Technology Directorate who are actively engaged
22 in learning how data mining can further the Department's

1 counter-terrorism mission have joined with my staff to make
2 this workshop not only a possibility but hopefully a
3 success.

4 So beyond that, why are we doing this? Well, if
5 you have followed developments up on the Hill, you are
6 aware of the various annual data mining reports that my
7 office has issued and certainly you're familiar with
8 Section 804 of the 9/11 Commission Report Act, which
9 requires of all agencies data mining reports, and our
10 Department is no exception. Earlier this year we issued a
11 letter report in which we advised Congress that we would be
12 doing some further work, among which would be this workshop
13 and that we would be reporting back to Congress on what we
14 found, and so we are convening the workshop in part because
15 of Section 804 of the 9/11 Commission Report Act.

16 So I think at this point it's appropriate and
17 necessary to remind everyone here that it's Section 201 of
18 the Homeland Security Act of 2002; the Department has a
19 Congressional mandate to conduct data mining activities in
20 furtherance of its mission. So we looked into this because
21 of course we read the plain language of the statute, and
22 certainly it says that, and we agree with that. And what

1 we wanted to understand a little bit better, what is it
2 that Congress was thinking? And at the time, then House
3 Majority Leader Dick Armey responded to those who were
4 concerned about that provision of the Homeland Security Act
5 and referred approvingly to the new Privacy Office that was
6 to be stood up, the office that Toby and I are in right
7 now, and said that that office would be there at the
8 Department to make sure that there were not abuses of data
9 mining.

10 So for us, the question isn't then, whether the
11 Department should be conducting data mining? It is, rather,
12 how DHS should use data mining and what ways can it do so that
13 both respect privacy and also support the integrity and
14 effectiveness of the Department's Homeland Security
15 initiatives. So in the interest of brevity, and because
16 we've gotten started a little bit late, I want to wrap up
17 my remarks.

18 And I want to introduce our guest speaker, DHS
19 Under Secretary Jay Cohen, who heads the Department's
20 Science and Technology Directorate. When Jay joined the
21 Department in August of 2006, a month after I moved over to
22 the Privacy Office, he immediately tackled the challenge of

1 the Science and Technology Directorate so that it could
2 foster the development of vital technologies for protecting
3 the nation. Jay deserves tremendous credit for his efforts
4 to transform the Directorate into an efficient and
5 responsible organization that makes vital technical
6 contributions to the DHS mission to protect against and
7 respond to catastrophic events. The S&T Director provides
8 technology solutions to help the men and women who face
9 risk every day on the front lines of Homeland Security to
10 do their jobs more quickly and safely with greater
11 accuracy. And with that, I'll stop. Jay Cohen, Under
12 Secretary of Science and Technology Directorate. Thank you
13 all very much.

14 [APPLAUSE]

15 PRESENTATION OF JAY COHEN, UNDER SECRETARY OF
16 SCIENCE AND TECHNOLOGY DIRECTORATE, U.S. DEPARTMENT OF
17 HOMELAND SECURITY.

18 Mr. Cohen: Well, good morning. And thank you so
19 much for sharing your most valuable asset with us, and that
20 is your time and also your thoughts, at this workshop.
21 It's a real pleasure to work with Hugo and his team. I
22 don't know anybody who has a tougher job in Homeland

1 did it to eliminate or minimize seams because terrorists
2 and criminals will always take advantage of seams. And
3 anything that eliminates or minimizes those is good for
4 security and bad for those who would do us harm. So that's
5 my little shtick here; I'm not a Mac person, but I'll do my
6 best with this computer.

7 So what are the goals in law of the Science and
8 Technology Directorate? And I can tell you, as Chief of
9 Naval Research for six years of a three-year tour and the
10 Office of Naval Research was established in 1946. Half a
11 page in Title 10, it says there will be an Office of Naval
12 Research, it'll be led by a Navy Admiral, report to
13 Secretary of the Navy, and it'll do good research. In
14 2003, of the 183 pages creating the Department of Homeland
15 Security, 17 pages describe the S&T Directorate. You know,
16 a camel was that animal created by committee, so we could
17 have ended up with a camel. We didn't. It was very, very
18 thoughtfully done. And so half a page in 1946, 17 pages in
19 2003; it shows you the impact of word processing on the
20 legislative process.

21 But to synopsise in the law what are the goals
22 and what do I follow, number one, is to accelerate the

1 delivery of enhanced technological capabilities to my
2 customers. Who are my customers? In law, they are the 22
3 components: TSA, Border Patrol, Coast Guard, Secret
4 Service; and in law, first responders -- the police, fire,
5 emergency, medical, bomb disposal -- our heroes. And I had
6 no appreciation for the scale of our first
7 responders in America. We have 35,000 fire departments in
8 America -- 35,000 fire departments, of which 80 percent are
9 volunteer. When I go and visit them and I say, 'Hi, I'm
10 from Washington. I'm here to help.' They say, 'Great.
11 Buy a raffle ticket or a muffin because we need a new
12 pumper.' I mean, this is America. So it is a federal goal
13 with a local execution; I can tell you it's a great
14 challenge.

15 Second is to establish -- in my words -- a lean
16 and agile government service -- world-class S&T management
17 team. Ladies and gentlemen, I don't do S&T and my people
18 don't do S&T; we are a venture capital fund, we are a
19 mutual fund, we invest in S&T to de-risk it to give
20 capabilities to our customers. And when I say government
21 service -- because some political appointees -- people like
22 me come and go, but the half-life of Science and Technology

1 is such that there must be a continuum, and so that's where
2 government service is so critically important. And in my
3 experience in Navy and in Homeland Security, is that
4 Science and Technology -- unless I do something stupid and
5 Hugo works very hard to help me from doing
6 something stupid -- is bi-partisan, non-partisan, and I
7 believe that that is how it should be.

8 And then finally -- and this is a labor of love
9 for me -- is to provide the leadership and opportunities
10 for the next generation of our workforce. This is STEM,
11 Science, Technology, Engineering and Math. Ladies and
12 gentlemen, we're in crisis in this country today. In fact,
13 we're in crisis in most of the western countries. People
14 in middle school, young people, are turning away from
15 science and math, and when you ask them why, they tell you
16 the truth -- it's too hard. They're the Playstation
17 generation; they want instant gratification. If we don't
18 turn this around, ladies and gentlemen, in my opinion, in
19 fifteen or 20 years we will not be a first-world
20 economy. So that's a little bit of the background.

21 Now, what are the threats that we face? This is
22 a PowerPoint presentation, we'll leave copies, you can move

1 the boxes around however you want. I view the threats from
2 terror -- and oh, by the way, DHS is responsible for all
3 threats. In the
4 law, it's not just terror threats, it's also natural
5 disasters, like earthquakes and fire and flooding,
6 tsunamis, et cetera. But I view the threats as bombs,
7 borders, bugs, and business -- those are the original four
8 b's. It turns out I've got six divisions; two of them
9 didn't have b's originally. I think last spring they saw
10 the Bee Movie, but the division directors came to me and
11 they said, 'Hey, we're without b's; we're b-less.' So I
12 added two b's and that's bodies -- that's human factors,
13 and buildings, which is infrastructure protection. You
14 understand bombs, you understand borders, you understand
15 bugs; what's business? Business is the underlying cyber-
16 backbone that enables everything we do, and it is a very
17 new area, and very threatening and scary area, of warfare.

18 So if you look across the bottom left to right,
19 you see consequence of occurrence low to high, and then
20 likelihood of occurrence. We're always going to have
21 physical attacks; that's the reality of the world that we
22 live in. If you look in nuclear, that's a nuclear device -
23 - that's a nuclear bomb. The consequence of occurrence of

1 that going off are unimaginable; it's far off the scale to
2 the right. But today, today a terrorist would have to
3 either buy or build a bomb, and I would tell you -- you can
4 disagree -- that I think the probability of that is
5 somewhat low. Maybe not tomorrow, but today. But the day
6 after 9/11, ladies and gentlemen, we were delivering death
7 by 37-cent stamps in the U.S. mail -- anthrax, biological
8 attack. And so you can see while it may not be as much of
9 a weapon of mass destruction as nuclear, its occurrence is
10 more likely. We have seen it, we will see it again.
11 Biological warfare is the poor man's weapon of mass
12 destruction. Because today, with the internet, with
13 genomics, all it takes is a brain, a basement, a
14 microscope, and you can create a pathogen that will give
15 you a pandemic.

16 IED's -- they're weapons of mass influence, not
17 weapons of mass destruction. Tom Friedman said IED's are
18 coming to a theatre near us, and I believe that.

19 But the tactics, techniques, and procedures that
20 we use so well overseas, many of them don't apply -- don't
21 apply in the United States because the Constitution,
22 because of the Fourth Amendment -- many of the things that

1 you're going to be discussing here. Before a bomb squad
2 can actively jam a bomb and its trigger device, they have
3 to get a license from the Federal Communications Commission.
4 It's a very interesting challenge; not what you're going to
5 be addressing today.

6 But what you are going to be addressing today is
7 up in the upper-right, high and to the right, and that's
8 cyber, because every three seconds someone's losing their
9 identity. And you have Estonia, and you understand if your
10 background, the challenges of what a cyber-attack could do.
11 Those of you who have children or grandchildren in college,
12 you understand they live from ATM swipe to ATM swipe. And
13 if we can't do that, in my opinion, there will be panic in
14 the streets. So you can agree or disagree, but that's sort
15 how I see life.

16 So Hugo has already talked about the enabling
17 legislation, I think very well thought out, well debated;
18 it has been modified, we've had a change in the Congress in
19 the ensuing years. We get to testify a lot. Everything I
20 do -- I'd contend 99.9 percent of what I do is
21 unclassified. We invite the Congress to our processes, we
22 invite the Inspector General; and Hugo has workshops like

1 this, which I know will be the first of many to come. So
2 the authorizing legislation for me, I have summarized it,
3 in the first, telling you what my goals were. I think I --
4 I'm too fast.

5 So as we look at data and we look at the threats,
6 and I looked at what is unique in Homeland Security, I
7 settled really on two things. Because the enabling
8 legislation is very thoughtful, it tells me not to recreate
9 the National Institutes of Health and not to recreate the
10 Center for Disease Control and not to recreate the
11 Department of Energy or Department of Defense labs -- and I
12 think that was very thoughtful -- but in exchange, it
13 allows me to leverage everything they do. I can't tell
14 them how to invest their billions of dollars in research,
15 but they give me full disclosure. And it really does work.
16 And then I take my precious dollars, our precious dollars,
17 and apply it to the things that are unique in Homeland
18 Security and the missions that we have.

19 So from my perspective, as I looked around at all
20 of the areas of Science and Technology, all the different
21 disciplines, the two that I felt -- and I still feel that
22 way after two years on the job -- that were unique, was

1 number one, the psychology of terrorism. Why do terrorists
2 do what they do? I mean, you can view them as criminals,
3 you can view them as armies, et cetera, but why do they do
4 what they do? It was not clear to me any other component
5 of government was investing in that.

6 And the second area is hostile intent, and we're
7 going to talk a little bit about that. Are there ways of
8 knowing that someone is about to do something bad to our
9 society? And so these are focus areas that we are looking
10 at. This is new science. We've gone to the National
11 Academies of Science to help us define those sciences. You
12 know, after World War II, the Battle of the Atlantic,
13 strategic bombing, the science of operations, research
14 operations analysis, was born. And after Sputnik
15 aerospace, you get the idea. As time moves on, challenges
16 change; new areas, new disciplines develop. But how do we
17 know that what we think is appropriate research, even
18 vetted by the Privacy Office, even briefed to the Congress;
19 and of course, the press is very interested in this, as
20 they should be. I mean, at the end of the day, ladies and
21 gentlemen, I am a citizen, I value my privacy, I respect
22 and value your privacy, and when I'm done with government

1 service, I will again be a citizen. I think I'm a citizen
2 while I'm still in government service, but you get the
3 idea.

4 So Dr. Sharla Rausch and her people are
5 represented here today. She's head of my Human Factors;
6 this is a division that I set up. There's a great ad by
7 Dow Chemical, it talks about the human element. I love
8 that ad because it's the human element that creates
9 terrorism and it's the human element that will solve the
10 challenges that we have. It really is all about humans.

11 But Sharla went ahead and worked with the Privacy
12 Office and others, established on her own, the Community
13 Perception of Technologies Panel. And so these are just
14 average people from a wide cross-section -- they have a
15 picture of them here -- and we go ahead and we brief to
16 them. This is our initiative, what we're looking at, what
17 our research areas are, how we're approaching it. They are
18 not necessarily experts in privacy; we go to Hugo and his
19 team for that, and I've got Jen Schiller on my staff. And
20 I can tell you, she is very tough on me. This is an area
21 where an ounce of prevention is worth pounds if not tons of
22 cure.

1 And it's very interesting to sit down, and I sit
2 down with this panel, and get their feedback on their
3 perception on what we are doing, and then we modify as
4 appropriate.

5 So let's talk a little bit about the areas of
6 research that we are doing, and then I'll conclude because
7 I know Hugo does want to get you back on track into panels
8 and the discussions are so important. So I'll go through
9 this very quickly.

10 And I must tell you that personal identifier
11 information was a new concept to me when I came on board,
12 and so in the last two years I've had a steep learning
13 curve. And I also understand that we can be looking, you
14 know, at totally unclassified, totally public information,
15 but perception of how that is analyzed, et cetera, becomes
16 an issue on its own. And I know you're going to address
17 all those things. The Congress enabled the S&T Directorate
18 with Centers of Excellence. I have two pillars of basic
19 research: universities and laboratories. And so at the
20 University of Maryland, one of our earliest Centers of
21 Excellence was the Study of Terrorism and Responses to
22 Terrorism (START). In Washington if you've got a good

1 acronym, everything else follows. So I salute the
2 University of Maryland on getting started with this.

3 But as you can see, this is the largest terrorist
4 event database; more than 80,000 events. Basically, this
5 is all out of public venue, public information; and you can
6 see incidents versus fatalities by area, et cetera. It is
7 unclassified, it's kept up to date, it's available for
8 researchers, et cetera.

9 The next area is Biodefense Knowledge Center. I
10 talked to you about my concerns for the poor man's weapon
11 of mass destruction. This is a 24 by 7 secure website; it
12 uses data fusion, and basically it's talking about
13 capabilities, because as you know, a bio and genomics are
14 moving at the speed of heat. And so it's available for
15 subject matter experts, et cetera.

16 Suspicious behavior detection. The goal here is
17 to identify deception and hostile intent in real time using
18 non-invasive sensors. We're going to talk a little bit
19 more about this when what we call the FAST program, FAST is
20 Future Attribute Screening Technology. But the goal here
21 is to develop a prototype to detect deception and hostile
22 intent in real time. I must tell you, almost everything we

1 do as we look at, for example, transportation security, is
2 to maximize the throughput of primary screening so the
3 lines are as short as possible, and then only focus on
4 secondary screening which can be question and answers.
5 Those of you who fly overseas, you know they do it a little
6 bit different than we do it. You start out with the
7 questions, and then you go through the metal detector. We
8 put you through the metal detector, and then after there's
9 suspicious activity, we then go into the secondary
10 screening. Secondary screening is very expensive,
11 intensive, and it interferes with our lives.

12 So what is FAST? Aviation in large measure is a
13 closed transportation system. We put up with the lines
14 because we believe that if we keep bad people and bad
15 things off of aircraft -- and oh, by the way, aircraft is a
16 fixation by some of our terrorists, enemies with aviation -
17 - if we keep bad people and bad things off of planes, the
18 plane will take off and land safely. It's a closed system.
19 The only challenge is the shoulder-fired weapons, and we're
20 working on that independently.

21 But when you get into Metro, you get into Amtrak,
22 you get into buses, you get into mass transit, where you

1 have thousands of people, we can't use the same procedures;
2 those are open. And if we kept a bomb from getting on at a
3 Metro or an Amtrak station, you still have miles of
4 unsecured railroad. So what is the balance? And so what
5 we're looking at here -- and let me give you an example --
6 during the SARS epidemic overseas, several Asian countries
7 used infrared cameras. As you got off the plane and you
8 walked into Customs, these cameras didn't care if you were
9 tall or short, male or female, they didn't care about
10 ethnicity, they were just looking at your forehead.
11 They're looking at your forehead. And if on infrared your
12 forehead was warmer than everyone else's forehead, you most
13 likely had a fever, and that's a precursor or an indicator
14 of SARS, and they didn't want to have the spread of SARS,
15 and so you went into secondary screening. That is the
16 level of screening that we're talking about. So if you're
17 a terrorist, you want to get to your target, you may be
18 nervous, you may be perspiring, your forehead may have
19 evaporative cooling, your heart rate may be raised, your
20 eyes may be flashing, your gait may be different. There
21 are micro-facial features that give away -- and this is a
22 brand new science that we're learning about today. Are you

1 telling the truth or are you deceptive? And so the goal
2 here is in a public event, like the Super Bowl or the
3 Olympics, to go ahead and see if, can we do this non-
4 invasive screening that will give us indication of hostile
5 intent so that we can take an individual to secondary
6 screening? Now look, your parent may have just died, you
7 may have been late getting to the event, you may have just
8 run; I mean, there are a lot of reasons why you can have
9 all these indicators, so we're looking at getting to the
10 secondary screening. That's the thrust of what we do.

11 Violent Intent Modeling and Simulation. Again,
12 this looks at the systematic collection and analysis of
13 information that is related to understanding terrorist
14 group intent. So we talked about the individual terrorists
15 -- why do they do what they do -- now, what about the group
16 as they come together?

17 So that's a summary of what we're doing;
18 everything we're doing is fully vetted with the Congress,
19 with the Privacy Office, et cetera. But at the end of the
20 day -- as I've told you with my basic mission -- product is
21 job one. Getting those tools to those that would make us
22 safe or keep us safe is what Science and Technology is all

1 about.

2 So I thank you so much for spending your time
3 here. I wish I could spend a day-and-a-half with you; I
4 think is going to be one of the most fascinating panels
5 that have occurred in the short history of DHS. Remember,
6 we're only five years old. Some of you have 5-year-old
7 grandchildren or children; you know how mature 5-year-olds
8 are, but all the vectors are in the right direction. And
9 the only question I ask myself and I ask my people, and I
10 hope this never happens, I hope there's not another attack,
11 I hope there's peace and happiness in the world. But if
12 you listen to most of the experts on both sides of the
13 aisle, they will tell you, there will be another attack.
14 Our terrorist enemies want to make it even more devastating
15 than that of 9/11. And the question is not if, it is
16 when. And so the question I ask myself every night is,
17 under my tenure will we have done enough with the resources
18 and tools that I have, consistent with the laws and our
19 culture, to make us as safe as we can be? So with that
20 thought, I'll leave you. Hugo, thank you so much for
21 giving me this opportunity, and I look forward to the
22 results of the workshop. Have a great day. Thank you.

1 [APPLAUSE]

2 Ms. Landesberg: Thank you, Under Secretary
3 Cohen. I'm Martha Landesberg from the Privacy Office, and
4 it's my pleasure this morning to introduce our next speaker
5 to you. He is Professor David Jensen who is an Associate
6 Professor of Computer Science and Director of the Knowledge
7 Discovery Laboratory at the University of Massachusetts
8 Amherst. Professor Jensen currently serves on DARPA's
9 Information, Science, and Technology Group, and he was an
10 analyst in the Office of Technology Assessment from 1991 to
11 1995. I give you Professor Jensen.

12 [APPLAUSE]

13 PRESENTATION OF DAVID JENSEN, ASSOCIATE PROFESSOR
14 OF COMPUTER SCIENCE, UNIVERSITY OF MASSACHUSETTS.

15 Mr. Jensen: Thank you. Thank you very much.
16 Under Secretary Cohen is a difficult speaker to follow, and
17 so I hope I can keep this as interesting and relevant to
18 today's conversations. So what I'm going to talk today
19 about is at some level somewhat boring in that it is about
20 definitions. But as many people have said, words mean what
21 we want them to mean. And I think in this particular case,
22 data mining means many things to many different people.

1 And so I'm going to try today to talk about the range of
2 definitions, and the ways in which we can come to a
3 definition that is both consistent with what the technical
4 community is doing, which is my community, and also
5 consistent with what we mean in a policy context.

6 So what I'll talk today about are, first, I'm
7 going to give you some very simple definitions, frequently
8 used definitions of data mining. Then I'm going to give a
9 fairly extended example of some work that I've done
10 recently in detecting securities fraud because I think it's
11 a good example of what modern technology is doing in the
12 area of data mining, and gives us some concrete things to
13 refer back to to try and expand and make more realistic the
14 definitions of data mining that we're going to be talking
15 about. Then I'll present some revised definitions, and
16 finally try to answer the question, why we should care
17 about definitions, and talk about how on some sort of
18 expanded understanding of data mining can reframe some
19 existing issues that are often brought up about the
20 technology and potentially raise interesting new issues --
21 new policy issues.

22 By the way, if you have a question that is

1 specific to some slide or some comment I've just made,
2 please feel free to raise your hand; I'd be happy to take
3 the question then. If I don't see you, give me a shout.
4 And also, there'll be a period at the end where we'll take
5 more questions of the more general kind.

6 So the major points I'm going to be talking about
7 today are first, that there are simple definitions that
8 portray data mining as a process of filtering or
9 extraction. That these definitions are very easy to state,
10 and in some ways, very vivid, but they are very easy to
11 misinterpret. They're not really wrong, but they're easy
12 to misinterpret, and I'll explain specific reasons why
13 that's the case. More useful definitions of data mining
14 portray it as an iterative process where you are both
15 learning and doing probabilistic inference, and you're
16 doing that over interconnected data records, not data
17 records that are independent from each other. Finally,
18 I'll say that these definitions identify different issues
19 for policy discussions, and I would argue, more interesting
20 and useful ones.

21 So let's look at some of the simple definitions.
22 The first is the one that I think has brought us to today's

1 meeting, from the Federal Agency Data Mining Reporting Act
2 of 2007 Secretary Cohen referred to, in which the -- well,
3 the definition says, it is a "program involving pattern-
4 based queries, searches, or other analyses of one or more
5 electronic databases." And then there are a series of
6 caveats that I think are really very specific to the Act,
7 saying, well, this has to be done by a federal agency or an
8 agent of a federal agency, it has to be about identifying
9 terrorism or criminal activity instead of other things.
10 But the key thing here is to focus on this question of
11 pattern-based query searches or other analyses.

12 Now, there are a variety of other definitions of
13 data mining. Let me give you some from the more technical
14 end of the spectrum. "The science of extracting useful
15 knowledge from data repositories," this is from the
16 Association for Computing Machinery Special Interest Group
17 on Knowledge Discovery and Data Mining, our Curriculum
18 Committee that came up with this definition.

19 There's also a very well-known definition from
20 some of the founders of the field, "The non-trivial
21 extraction of implicit, previously unknown and potentially
22 useful information from data." That's from an article

1 about knowledge discovery and data mining.

2 Now, I tend to use the term knowledge discovery
3 because I think it is intrinsically more meaningful and
4 less easy to mistakenly understand than data mining is. I
5 think data mining has a clear and obvious meaning which is
6 wrong; the clear and obvious meaning is that you are mining
7 for data, and that's not actually what data mining is
8 doing. If you say gold mining, that means you're mining
9 for gold. If data mining should be mining for data, you're
10 not. You're mining for knowledge, and knowledge discovery
11 gets at that. Although, it did confuse my Dean greatly
12 when I was introduced to him as doing knowledge discovery,
13 he looked and he said, "Isn't everyone at a university
14 doing that?" And I said, "Yes, yes. But we're doing it
15 with computers." He said, "Oh, well, that's very
16 interesting," and we went on to have a pretty good
17 conversation. There are other terms, as well -- predictive
18 analytics, advanced statistical modeling, machine learning.

19 So, well, I'll stick with the term data mining
20 even though it's not my preferred term because it is the
21 term that stuck. So let me give you an example of this
22 sort of work -- this sort of technology, and it's about

1 detecting securities fraud. We've been working for about
2 five years now with the National Association of Securities
3 Dealers. This is the non-governmental, private
4 organization in the United States that regulates all stock
5 brokers. They came to us about five years ago and they
6 said, 'We hear you're doing work in analyzing the kind of
7 data that we need to analyze, wonder if we might do some
8 work with you,' and we've been doing joint projects with
9 them ever since. By the way, NASD is now referred to as
10 the Financial Industry Regulatory Authority. They changed
11 their name recently, but I'll be using NASD because it's
12 what sticks in my head and also it's because what's
13 relevant to the work we did over the past five years.

14 NASD is the parent of the NASDAQ Stock Exchange -
15 - stock market, but they spun that off because their
16 central focus is really regulatory. They monitor a large
17 number of securities firms, branches, and individual people
18 who sell securities to the public. Those are referred to
19 as registered representatives or reps. And one of their
20 responsibilities -- they have several -- is to prevent and
21 discover serious misconduct among brokers -- I'll use the
22 term fraud. They incur fines and they can even ban

1 individuals or entire firms from the industry and say, 'You
2 cannot work anymore in the securities industry.' Now, they
3 have a data set which they collected for their regulatory
4 function, not to do analysis on, but for their regulatory
5 function. That data set is called the Central Registration
6 Depository, or CRD database. It consists of data about
7 individual reps -- individual people -- about the branches
8 that they work for, the actual physical organizations that
9 they work, as well as the larger firms that those branches
10 belong to. And finally, a set of event reports, which they
11 call disclosures, where reps abide by the policies of NASD,
12 which they agree to when they become a registered
13 representative, they have to disclose certain events in
14 their lives, including simple things like if a customer
15 complains, but also including things such as liens against
16 them, major issues in their financial history, or if they
17 commit a felony, for instance, that's also a disclosable
18 event. So there are a set of those disclosures that are in
19 this data set.

20 Now, importantly, this data set is a large set of
21 interconnected records. As you might expect, we know what
22 reps work for what branches, what branches -- what firms

1 -- own those branches, and what disclosures have been
2 filed on individual reps.

3 There are about 3.6 million reps in the data set,
4 about 750,000 branches, about 25,000 firms, and about
5 625,000 disclosures, so a moderately large data set. And
6 that covers a period of over 20 years. And we tend to
7 focus on the smaller subset about over the past ten years
8 or so.

9 Now, fortunately, the kind of conduct that NASD
10 is trying to discover is relatively rare. Now, fraud among
11 reps is quite rare. If you look at the stats, it's less
12 than 1 percent of reps commit any kind of serious
13 misconduct in a given year. In general, I think it's about
14 1/10th of 1 percent, so very small incidents of the kind of
15 serious misconduct they're looking for. But it's very
16 important to the public that that be discovered, and very
17 important to the integrity of the industry. So their task
18 is to take data from the past where they know that certain
19 reps or branches were engaged in serious misconduct, and to
20 take that data and to then try to come up with some sort of
21 method which they can use to guide their future activities.
22 So, particularly, they want to do examinations and they want

1 to do education and enforcement activities that will either
2 prevent fraud from occurring or catch it early. And so
3 they want to use the data they have, which they collected
4 for other reasons, but they came to us saying, 'We think we
5 can do more with the data; is that the case?'

6 So what we did with them was to construct
7 statistical models that try to predict the probability --
8 or estimate the probability that an individual rep will
9 commit some kind of serious misconduct in the next year, the
10 next 12 months.

11 And so one of the kinds of statistical models that we
12 devised is a kind of probabilistic or statistical model
13 which is tree-structured, and I'm showing you the whole
14 structure of a tree here. And by the way, there are
15 details of these models that are not included in these
16 slides, at the request of NASD for obvious reasons. They'd
17 rather not release exactly how they might be detecting
18 fraud. But this is the structure of the model and it's
19 structured like a tree. You could think of it as a virtual
20 Pachinko machine where you take an individual rep and their
21 surrounding context -- the disclosures, the branches
22 they've worked for, the firms they work in, the other reps

1 that they work with -- and you drop it in the top of this
2 tree. And then you answer a series of yes/no questions,
3 such that it rattles down to a leaf node, a thing at the
4 bottom which gives you a probability distribution -- their
5 probability of committing fraud in the next 12 months.
6
7 Let's zoom in on a portion of it. So at the top node we
8 say, 'How many disclosures have been filed on this rep?'
9 If it's greater than a certain number it goes down one
10 branch, if it's less than that it goes down another branch.
11
12 And so on. And we ask questions here in this model about
13 the number of the disclosures that were customer
14 complaints, whether that rep has been designated as high-
15 risk in previous years, other kinds of things about the
16 current branch they work at, et cetera. And eventually we
17 come down to a node where we say, 'Everyone who reaches
18 this point has a particular probability -- estimated
19 probability of committing fraud in the next 12 months.'

18 Now, importantly, we construct these models
19 automatically, or more accurately, the data mining
20 algorithms we have devised construct these models
21 automatically. They do that by searching a very large
22 space of possible trees. Now, the number of those possible

1 trees is vast. Here just for a five-level tree with the
2 kinds of data that we have, we're talking about 10 to the
3 106th, an extraordinarily large number of possible trees
4 that are out there. But, fortunately, in the technical
5 work of our field, we've devised a fair number of
6 efficient, approximate search methods to look at that space
7 and not have to examine it exhaustively but still find the
8 trees that are particularly useful or valuable in that
9 space. And we evaluate how well those trees work by
10 comparing them to the data for which we know the right
11 answer; that is, we know at least we have good estimates of
12 the -- which reps have committed fraud in the past. At
13 least those reps that have been identified, so they are
14 probably some -- many, in fact, that have not been
15 identified but we know a large number of reps that have
16 committed fraud in the past, and we can use that past data,
17 that retrospective data, to compare the accuracy of
18 different types of models -- different types of trees in
19 this case.

20 What the models then do is to infer the values of
21 an unobserved variable. The unobserved variable, the thing
22 we're trying to estimate here, is the risk that a rep will

1 commit fraud in the next 12 months. And there are also
2 some kinds of models I won't talk about that will
3 simultaneously infer the value of many unobserved
4 variables. But for the new data, for the data we want to
5 apply the model to, we don't know what reps are committing
6 fraud and thus we want to estimate the probability of
7 those.

8 The performance of these models has been
9 evaluated in a variety of ways, but one of the ways that we
10 used was we took a bunch of predictions from the model, we
11 took some predictions from NASD's current method of doing
12 initial screening, and we took reps that showed up on only
13 the list that our model created, only the list that NASD
14 created, neither list and both lists. And then we
15 scrambled those up and put them in front of trained NASD
16 field examiners and we made the estimates, for example, for
17 the previous year, for 2007. We didn't have data about
18 2007 about who had actually been found to be committing
19 fraud in that year. But NASD did have information about
20 that, and we asked the examiners the following question, we
21 said, 'If we had given you this list at the beginning of
22 2007, how useful would it have been given that you now know

1 what the right answers are?' And they rated each rep that
2 we had given them on a five-point scale. One is, it would
3 have wasted my time to know about this individual; five is,
4 I absolutely would have wanted to know about this. When
5 the reps showed up on neither list -- it's a little
6 difficult to see there -- but when they showed up on
7 neither list, the ratings were almost all one. When they
8 showed up on NASD's current list but not ours, the ratings
9 were roughly on average a three. When they showed up on
10 only our list and not NASD's list, again the average was
11 about three. And if they showed up on both lists --
12 combined list -- they had an average rating of about four.
13 So showing that we are doing -- the statistical model is
14 doing almost essentially as well as NASD's current rules
15 for doing screening to say, which reps deserve some
16 additional scrutiny to look and see if they're committing
17 fraud. And if you combine the statistical model with the
18 current expert derived rules, we can do even better.

19 We also got a little bit of anecdotal feedback;
20 one of the field examiners sent us an unsolicited note
21 along with his ratings, and he said, 'One of these reps I
22 was very confident in rating a five,' he said. He had had

1 the pleasure of meeting him at a shady warehouse location
2 during what we think is a sting operation. He said he'd
3 negotiated this rep's bar from the industry because among
4 other things, he'd actually used fraudulently obtained
5 funds to attend an NASD compliance conference -- conference
6 about how to comply with NASD rules. The examiner said,
7 'If you predicted this person, you'd be right on target.'
8 And in fact we, with some trepidation, we went to NASD's
9 list, the rep was not on NASD's list, we went to our list,
10 he was very high up our list. So a nice anecdote to
11 support the idea that this statistical model is a useful
12 one.

13 All right. With that background and that kind of
14 concrete reference, let's go back to our definitions of
15 data mining. So again, to recap the simple definitions,
16 we've got from the Data Mining Reporting Act, pattern-based
17 queries and searches or other analyses; extracting useful
18 knowledge from data repositories; extracting implicit
19 previously unknown knowledge. So one way of thinking about
20 these definitions, one simple kind of visual to get is the
21 idea of a filter. Where you say the system takes in data,
22 there is some mining or filtering process that's done on

1 the data, and then out pop predictions out the other side.
2 So that's what we've got, this kind of filtering process.
3 Now, this filtering process -- this idea of a filtering
4 process has been encouraged by some of the most powerful
5 people on the planet, some of the most powerful image
6 makers on the planet. Those people reside in Hollywood,
7 mostly. For those of you who have seen Minority Report,
8 this is a very persuasive image. This idea that there is a
9 black box out there that will be producing predictions, and
10 if the predictions are certain, they are crisp, there is no
11 doubt in them, and they put them out and that's what then
12 we go act on as a law enforcement agency. For those of you
13 who watch television also, there was a short-lived show
14 called Threat Matrix, which had some similar ideas that
15 were frequently propounded in the show about data mining.
16 And as you might expect, these media images are somewhat
17 simple. They're simple because it's very easy to
18 misinterpret the definitions which I've given you
19 previously, which can be interpreted accurately but it's
20 very easy to misinterpret them. Let me explain some
21 reasons why. The first is -- and I'll explain more about
22 each of these in the next set of slides -- the first is

1 that there is only once process. The misperception is,
2 there's only one process that encompasses what I'll refer
3 to both as learning and inference. The second is that the
4 records that come in the left side are disconnected from
5 each other. Here I'm showing just individual records about
6 reps. Third, that the inferences out the other side are
7 deterministic. Essentially we spit out a set of reps that
8 are bad and a set of reps that are good. Fourth, is that
9 this is only done once, this single stage, it's a once-through
10 process. And finally, that this process of data mining is
11 what I'll call institutionally isolating. That is, it just
12 sits off by itself in this little box and does its job.

13 Let me explain why each of these I think are not
14 accurate, and what is a more accurate picture. The first
15 is that the processes of learning and inference are
16 distinct. That is, there's not just one process, but
17 actually two. The learning phase takes in data for which
18 we know the correct answer, or we have good estimates of
19 the correct answer, and that puts out a statistical model.
20 That model is then used in an inference process to take in
21 data for which we do not know the correct answer, and put
22 out some kind of prediction.

1 Importantly, the learning phase is the part of
2 this overall process that is unique, that is the essential
3 component of data mining. In fact, many people in the
4 field would say that the inference part is really almost an
5 afterthought. The goal is to put out a good statistical
6 model. Now I will make one caveat, which is that there is
7 a lot of study in the field about some kinds of techniques
8 which do not immediately appear to fit into this, although
9 I think many of them actually do. So for instance, there
10 is some study of clustering. They're trying to look at
11 data and find homogeneous regions in it, and while there
12 does not appear to be a statistical model underlying that
13 there often can be and many of the better methods for
14 clustering do that fairly well. So some caveats; this is a
15 little bit simple to say that all of data mining has a
16 statistical model underlying it. But it's a good --
17 absolutely a good first pass.

18 So there's this learning phase and this inference
19 phase, and they are more or less separate. Learning is
20 what makes data mining unique. It's also important to
21 point out that the inference taking data for which we don't
22 know the correct answer and making an inference does not

1 require a statistical model. In fact, people do it all the
2 time. At NASD for instance, they had a set of rules that
3 they had sat down and worked with their experts to derive,
4 and that was what produced an initial list that then field
5 examiners went out and did additional investigation on.
6 And that was not derived from data mining, that was derived
7 from just sitting down and thinking.

8 Now, an example of the kind of misinterpretation
9 -- and I don't want to unfairly characterize GAO here in an
10 otherwise excellent report -- they had a graphic which --
11 this is 2005 report -- which starts out and says, "There's
12 input to the process, there is an analysis process, and
13 there's output." It is a slightly more complex version of
14 this filtering that I've talked about and doesn't clearly
15 distinguish between any kind of learning phase and an
16 inference phase. Instead, what we have, the idea here is
17 that data mining is really complex set of database queries.
18 It's a complex way of filtering a database to put out
19 matches. And I think that is a misinterpretation which is
20 easy to make, but actually dangerous in terms of public
21 policy. Let me emphasize again, though, that both this
22 report and several earlier reports from GAO are really

1 quite good and have a lot of useful information about data
2 mining.

3 Second issue, data records are often
4 interconnected, they're not sets of individual records. So
5 I show here these individual reps, but actually what we
6 have are a case often of a network of different types of
7 records that are interconnected. So think back to the NASD
8 example, we have this set of reps, branches, firms and
9 disclosures, this set of interconnected records and those
10 records are the -- provide us a lot more information than
11 just having records about individual reps.

12 This sort of approach, often called relational
13 learning or relational knowledge discovery or relational
14 data mining, has become increasingly prevalent both in the
15 technical community and now starting to make its way into
16 applications because this can improve both the accuracy of
17 the process and allow us to address entirely new types of
18 tasks, for instance, predicting a link or a connection
19 between two or more records.

20 Third issue, inferences that come out of the
21 inference process are not a kind of yes/no labeling.
22 Instead they are probabilistic. So rather than having a,

1 these are bad brokers -- these are bad reps and these are
2 good reps, we come out with a probability associated with
3 each of those reps. And almost all, I think, really modern
4 applications of data mining are giving probability
5 distributions on variables rather than kind of yes/no
6 classification.

7 What is important is that this allows us then to
8 do -- to have a lot more information about the inferences
9 that are being made. So for instance in the case of NASD,
10 you could imagine if we have probabilities we could look at
11 that last and say, it may be that there are a few high
12 probability reps, and then immediately drops to very low.
13 And then we would say, 'Maybe we should only look at those
14 high probability ones.' Conversely, there might be a very
15 long list, far longer than NASD would have originally
16 thought they needed to look at, that are very high
17 probability of committing fraud and they might say, 'Maybe
18 we should expand our screening program to look at a larger
19 set of individuals, if we believe that this is an
20 accurate assessment of probability.'

21 Finally, it allows you to assess accuracy in new
22 ways because you have these probability judgments and it

1 allows a much finer grained kind of evaluation of how well
2 the model is doing.

3 So an example of these kinds of probabilistic
4 models is the NASD model. We don't say that everyone who
5 reaches a particular leaf node is going to commit fraud.
6 Instead we say, there is a probability associated with
7 committing fraud.

8 Fourth issue that I've talked about is inference
9 is done in many real systems in multiple stages. So if you
10 look at the inference process, it's not just a once through
11 process, but instead there's feedback once you have an
12 inference, additional things can be done with that
13 inference in other rounds of inference about either new
14 problems, or in fact, in some cases about the same problem.

15 So a really good example of this is the way in
16 which screening for many diseases is done. So for
17 instance, AIDS screening is done with an initially very
18 inexpensive test which has a high false positive rate. It
19 is of course disturbing to individuals who get a positive,
20 but doctors are quick to point out, 'Look, this test has a
21 high positive. And even if you get a positive on this test
22 -- high false positive rate -- even if you

1 get a positive, the vast majority of people are actually
2 negative.' So now we're going to do the more expensive and
3 more accurate test. So this kind of two-stage screening is
4 a way of cutting down costs and increasing accuracy. And
5 that's the same way in which data mining can be done in
6 order to do those things, to improve accuracy and to allow
7 a wider range of types of inferences.

8 So it turns out actually that this is what NASD
9 does, is that this initial set of rules they have, or now
10 the kind of statistical model we've given them, gives them
11 an initial set of reps that get enhanced scrutiny from
12 their examination process. It's not that other reps are
13 not examined, and it's not that those reps in any way are
14 immediately considered to have committed fraud. Instead it
15 says we should look more closely. And then a human
16 examiner goes out and initially looks at records that are
17 just held centrally, and then often goes out into the field
18 and will examine records that are only held at the firm.
19 That larger set of records both centrally and out in the
20 field are more expensive to access and also more sensitive.
21 And so the question is, should we actually go to the --
22 should NASD go to the expense and the potential security

1 and privacy issues of examining those additional records?

2 Well, only if they have some initial sense that it would be
3 useful to look at those records.

4 Final issue. Data mining is used in a larger
5 institutional context than it might appear at first. So if
6 we think about data mining as -- the entire process I've
7 described as a box; we say, well, there's obviously some
8 kind of data gathering that has gone on ahead of time. And
9 once we get inferences out, there's some kind of decision-
10 making process. Those inferences do not immediately
11 indicate what we should do, what any organization should do
12 with that information.

13 And finally of course, there's some feedback.
14 With decision-making, you may say, actually it's useful to
15 gather additional data and perhaps do additional sorts of
16 analysis on the data. Importantly, many of the really big
17 public policy issues about privacy and utility are about
18 data gathering and about decision-making, not so much about
19 what happens inside that data mining box. The other issue
20 I think is that the use of data mining algorithms actually
21 imposes relatively few constraints on data gathering or
22 decision-making. That is, just because you have maybe in

1 advance decided to do data mining, does not mean
2 necessarily you will collect more data. NASD is a
3 wonderful example of this; they had already collected every
4 last bit of data, which we've used over the past five
5 years. They collected it for other reasons, but we've gone
6 ahead and used those data sets to do additional kinds of
7 analysis. And also, the output of data mining does not
8 imply necessarily anything about what you should then do.
9 It is input to a decision-making process that of course
10 should take into account a large number of factors.

11 All right. So those are some enhancements, I
12 hope, and some additional explanation about data mining.
13 And now the question I think may come up, why all of this
14 work? Why care about these definitions? And the basic
15 point I hope to make is that this gives us I think some new
16 perspectives, some new ways of looking at what is important
17 about privacy and questions of utility.

18 So one large issue that often comes up in
19 discussions about data mining is an issue about false
20 positives, particularly in cases such as counter-terrorism
21 applications or law enforcement or fraud detection
22 applications where the prevalence of the activity, the

1 frequency with which it happens is very low. And the
2 critique goes something like this, if the prevalence of
3 true positives is low, that is there are very few cases of
4 fraud in the case of NASD, then the vast majority of
5 inferred positives will be false positives. So even if you
6 have a very low error rate, if you have 100,000 people who
7 haven't done something and 1,000 people who have, and you
8 are 99 percent accurate, well then, you're going to have 10
9 people who actually did the thing that will show up as
10 positives. And, what did I say, 100 times that number that
11 will show up as false positives. And so this is a simple
12 critique, a relatively easy critique to get across, but it
13 unfortunately presumes this kind of filtered model. So
14 instead, if we think about these more accurate -- what I
15 hope are more accurate definitions -- the first is that
16 probabilistic inference can really help us here because it
17 allows you to control the types of errors which any
18 particular threshold that you put on that probability,
19 we're going to look at everyone with a probability over
20 .95. You can change the error characteristics of any sort
21 of screening that you do, so probabilistic inference helps
22 us a great deal. You can also use that to account for --

1 in addition to the expected -- what is sometimes referred
2 to as the expected class distribution of the data. It also
3 allows you to adjust for the relative costs of errors. So
4 if errors of false positives are very expensive or false
5 negatives are very expensive, you can modify those. It's a
6 great deal of work and what's called cost-sensitive
7 classification or cost-sensitive inference.

8 The second issue is that as I mentioned about
9 disease screening, multi-stage inference, and also it turns
10 out interconnected data records can help you greatly reduce
11 the false positive rate overall. There's some work that
12 several of my students and I did in 2003 showing ways in
13 which interconnected data records and multi-stage inference
14 can dramatically drop your rate of false positives overall
15 so you just end up with a better, more accurate classifier
16 to begin with.

17 It's not that the issue of false positive goes
18 away, it doesn't. But it's that the simple idea that
19 merely because prevalence is low, data mining methods will
20 utterly fail is incorrect. And the simple definition seems
21 to support it, more accurate definition does not.

22 Another very frequent issue which has come up,

1 particularly in the past several years is this idea of
2 subject-based versus pattern-based queries. So some people
3 have proposed limitations on data mining under the idea
4 that you want to differentiate between inferences that are
5 based on individuals, that is, I suspect this individual
6 has committed a crime, I'm going to go look at data about
7 them, versus pattern-based queries which says, I think
8 there is an indicator of some kind of misconduct, I'm going
9 to go look for everyone in my data set that has those
10 characteristics.

11 The first, subject-based queries, is thought to
12 be better because we have some initial suspicion. And
13 pattern-based queries in the worst possible case are seen
14 as some kind of dragnet; we're going to go out there and
15 we're just going to filter and we are going to end up
16 probably with a lot of false positives. So subject-based
17 queries tend to be in this formulation preferred over
18 pattern-based queries. In fact, some have gone so far as
19 to suggest only subject-based queries should be allowed.
20 Now, frankly, I have an enormous difficulty understanding
21 even what this idea means in a realistic scenario. Because
22 if you come from the technical world and you think about

1 how we do probabilistic inference, there is no fundamental
2 distinction whatsoever between inference based on things we
3 observe, that is, I suspect that this individual or set of
4 individuals is engaged in stock fraud, let's say,
5 securities fraud, and unobserved variables which is what
6 might be loosely matched up with pattern-based queries.
7 All that having initial suspicion is, is evidence to do a
8 better job of inference overall. And so from my
9 perspective, from the technical perspective, there is no
10 essential difference between pattern-based and subject-
11 based queries; it's all inference and we use what evidence
12 we have available to us.

13 Another way in which this is very difficult to
14 understand in a technical sense, is that in a multi-stage
15 process of doing inference, pattern-based at one stage --
16 if we can even formulate in an interesting way -- becomes
17 subject-based at another. Because, for instance, if we
18 have some process that identifies some individuals, let's
19 say, as having a higher probability of committing
20 securities fraud, then suddenly we are now subject-based in
21 the next phase of inference.

22 Finally, relational data -- the idea that we are

1 making simultaneous inferences about many interconnected
2 records, again makes this distinction between subject-based
3 and pattern-based queries more or less disappear. Even the
4 name, queries, I think, is showing this filter-based idea
5 of definition of data mining versus a more accurate
6 technical definition.

7 Another very frequent issue that is raised that I
8 think has a lot of validity in one sense is a concern about
9 having large, centralized databases. So if you have an
10 extremely large centralized database, it is a single point
11 of failure. And computer scientists for a variety of
12 reasons would say it's a bad idea to have a large,
13 centralized database. It's a single point of failure. It
14 also means that if one institution, one agency controls
15 that database, there's a higher probability of what's often
16 referred to as mission creep. That is, the data set is
17 collected for one reason and suddenly people start to say,
18 'Hmm, we could use it for other reasons,' which may not be
19 strictly in keeping with the statutes behind that
20 organization.

21 Now there are a variety of legal protections you
22 could put in place to make that not happen, but there are

1 also technical ways in which I think this critique does not
2 fit with the way that I think many really modern
3 applications would be done.

4 The first is to say real applications -- and
5 certainly this is the case for NASD, but also the case I
6 know for the U.S. Treasury Department -- multi-stage
7 inference means that you don't have to have one large
8 centralized database. In fact, there are good reasons to
9 say that you want to distribute them among either many
10 agencies or at least many different parts of your
11 organization. Because the idea is you have one data set
12 that you do one sort of analysis in, that gives you some
13 potentially initial inferences, and then you say, 'Well, we
14 have now a smaller set of individuals or of records that we
15 want to go look at in more detail.' So now we go out and
16 we've got some additional data because to do that
17 additional examination at NASD for instance, they need to
18 get some additional data and that's from another database.
19 Not a problem. And in fact, a benefit because you don't
20 have a large single point of failure, and in an
21 organizational sense, if these data sets are distributed
22 among different agencies, then you have at least a kind of

1 students and I have done recently, we published at a
2 conference last month, about how to both understand the
3 privacy implications of this and at least candidate
4 algorithm for protecting the privacy of these kinds of data
5 sets and allowing analysis, and not being able to re-
6 identify people. But preserving privacy in this context is
7 a new problem and a difficult one.

8 Finally, I think that we need to look at and
9 think more about how to combine statistical and human
10 inferences. Chris?

11 Christopher Clifton: (Speaking off microphone.)

12 Mr. Jensen: Excuse me. Thank you, Chris.

13 It's actually whether you emailed someone.
14 Sorry, not the number, but whether you -- the links in the
15 data set are merely, I emailed at least five messages to
16 this person, so it's a very, sort of, stripped down sense
17 of what the social network is -- or the professional
18 organizational network in the Enron Corporation. Thank
19 you.

20 Last issue is combining statistical and human
21 inferences. We have the statistical models, we also often
22 in many real applications have real experts. In NASD for

1 with something that is illegal that I need investigate.
2 Granted, the number of individuals in his or her field of
3 view may only be 100 or 1,000 or 10,000 if he or she is at
4 a, you know, a large public event. But nonetheless, there
5 is an inference process occurring there based on knowledge
6 about what is frequently associated with something illegal.
7 I don't see any difference in kind between that activity
8 and the activity that we -- the kind of inference that
9 would happen in a database. It produces a subject-based
10 concern based on knowledge about statistical associations.
11 I also think that there's an advantage and, really, we can,
12 you know, go into extreme measures here, but I think
13 there's an advantage in being able to look at the
14 statistical model which is used in data mining inference
15 because it's sitting there in front of you; you can examine
16 it, it's possible to do judicial review on that particular
17 thing. It's possible to audit it. In contrast, if it's in
18 somebody's head exclusively, it's much more difficult to
19 access and understand. There are special challenges, no
20 question, in understanding the mathematics of the
21 probabilistic models. But in terms of that subject-
22 based/pattern-based distinction, I don't understand it.

1 thought to be terrorists. And that list is used airports
2 and other venues to determine if a person matches that
3 list.

4 And finally, there is pattern-based, or what I
5 call event-based data mining. This is when the government
6 does not have a suspect but has an event, either one that's
7 already occurred or one that they're fearful will occur,
8 and has a profile of a potential perpetrator of that past
9 event or future event, and uses that profile and a data
10 mining endeavor to figure out who the perpetrators might
11 be.

12 Now, what are the harms of these various kinds of
13 data mining programs? I guess if I were going to be a
14 counter-advocate, a person who would dismiss the privacy
15 harms of data mining, I would compare data mining to a
16 search of a house, which is of course the classic police
17 investigative technique. Data mining is covert. All
18 right. People don't know it's happening most of the time.
19 Suspects are not aware they're suspects, unlike with a
20 house search where the police are in your bedroom, in your
21 living room, going through your belongings. It's not
22 physically intrusive. The physical intrusion concept is one

1 I tend to think for a pure analytical convenience of three
2 categories of harms, which I think will be very intuitive.
3 There's nothing scholarly about this. One, we might think
4 of the harm -- the impact on individuals or what, my guess
5 we on the panel would largely think of as privacy harms.
6 So the individual is detained; the individual's privacy is
7 invaded; the individual suffers inconvenience; the
8 individual might suffer incarceration. Whatever the
9 individual impact, you can imagine this wide range
10 depending upon, again, the specific type of activity
11 involved -- type of data mining involved.

12 The second type of harm, which is frankly the one
13 that concerns me the most, or what I think of as efficacy
14 harms; namely, we waste resources worrying about the wrong
15 people, or we fail to worry about the right people. And
16 again, this could be the result of any number of errors
17 that creep into the system which Christopher was describing
18 earlier. It may be the data were inaccurate; it may be
19 they were linked to the wrong person. For example, I doubt
20 if Ted Kennedy actually ever was on the watch list. I
21 suspect that he had a name similar to someone on the watch
22 list, and because we have such lousy identification in this

1 country somebody comparing the identification with the
2 watch list made what can only be considered to be a stupid
3 mistake. But the issue there is not really that Ted
4 Kennedy was detained briefly, it's that that agent and
5 those dollars were spent dealing with a non-threat when
6 real threats were going unaddressed. So if in fact as
7 Professor Jensen mentioned at the outset, we use data
8 mining as a way of narrowing our focus of probabilistic of
9 saying, where do we spend our scarce resources? To the
10 extent it is in error, it is causing us to waste those
11 scarce resources. We're putting them in the wrong place.

12 And then the third type of harm is actually one
13 which you suggested in your opening remarks, Martha, and
14 that is what we might think of as the political or the
15 public support risks. In other words, how many good ideas
16 have been killed because of the public outcry about them?
17 How often do we see hesitancy among extremely talented and
18 well-meaning government officials because of the
19 controversy surrounding this use of a term? I remember in
20 the time that I served as Of Counsel to the Technology and
21 Privacy Advisory Committee that was investigating TIA, an
22 extremely anguished Air Force General testifying before the

1 know a fair amount about the way in which fraud takes place
2 in the securities industry, we know what to look for. Then
3 you look at terrorists, right, we have something like 60
4 million unique passengers every year in the United States
5 who fly commercially. If there are, I don't know, 6,000
6 terrorists -- I don't believe that, but if there were
7 6,000 terrorists -- we're not talking about 1 percent of
8 those people being terrorists, we're not talking about
9 1/10th of 1 percent of those people being terrorists, we're
10 talking about something less than 1/100th of a percent of
11 those people being terrorists. The chances that the data
12 mining is going to be able to identify them are vanishingly
13 small but the results of the enormous number of false
14 positives are significant.

15 Ms. Landesberg: Thank you, Barry. Peter, let me
16 turn to you. I know you wanted to comment a little more on
17 the subject-based/pattern-based distinction. And then I'd
18 like to turn us to talk a little bit about the sources of
19 data. So why don't you go ahead.

20 Mr. Swire: I'm not sure I have much more to say
21 on subject versus pattern-based beyond what Chris Slobogin
22 said. But I know one of the questions was, sources of data

1 -- some things come from open source intelligence, you get
2 it on the internet from surfing, some things come from
3 private databases, some from government, so can I turn to
4 that? Is that --

5 Ms. Landesberg: Absolutely. Certainly.
6 Certainly.

7 Mr. Swire: So I wanted to highlight some things
8 that -- I agree with a very many of the things that have
9 been said, but I wanted to highlight some things that
10 haven't been said yet. And preliminary remark is, my view
11 is as we get better information technology, there are so
12 many wonderful uses of it and so the Homeland Security, law
13 enforcement, everybody else in society should be doing
14 many, many new projects of knowledge discovery, et cetera.
15 And in some ways, this panel and this workshop is saying if
16 there are 100 possible new uses, are there one or two or ten
17 or whatever where maybe there should be questions asked,
18 maybe there are certain kinds of problems. But
19 overwhelmingly we should -- when IT gets better -- we
20 should use it, and then we should also build it in ways
21 that watch out for civil liberties, watch out for the other
22 problems.

1 In terms of sources of data, one of the debates
2 around data mining is how much should the government get
3 private-sector data? The Total Information Awareness
4 version of that was, always as much as possible. All
5 right. Get the medical data, get the financial data; we'll
6 get command and control of the information battle space.
7 And there were various kinds of push-back on that.

8 But I want to highlight some history of reasons
9 to be cautious about thinking everything in the private
10 sector should go to the public sector. And this is private
11 sector of the sort that's in your bank or in your
12 hospital records or in your communication records. So one
13 thing is -- to go back to the mid-1980's when this new tool
14 called email was first coming on the horizon, and in the
15 lobbying around what became the Electronic Communications
16 Privacy Act, the big heavy-hitter in the room was IBM
17 because IBM made the corporate decision that to make email
18 grow, there had to be strong statutory protections for
19 privacy because they didn't think that if email
20 were open sesame to everybody including government, that
21 people would trust email. And so IBM put a bunch of
22 lobbying muscle in and working together with civil society

1 meetings, says, 'But of course we don't want the government
2 to get this. Of course this would be of no interest to the
3 government. You should not even look behind that curtain
4 and think the government would ever even be interested in
5 any of this data. This is just about advertising and
6 whether we're going to have cooler content on the
7 internet.' And then when you talk to the companies
8 privately -- and I've talked to people who are working
9 these issues for some of the business companies that do
10 this -- they say, 'The elephant in the room, the thing we
11 really know is the government could come and ask for all of
12 this.' And if that happens -- if it's seen that the online
13 advertising market is basically a front for a full
14 government database, the whole thing is going to collapse
15 around them. And so that debate over there, which is, how
16 are we going to do ads on the internet, is vitally affected
17 by, how much is the government really going to get all
18 that? Or how much are consumers going to think the
19 government gets all that? And so I'd, you know, to the
20 extent you want that robust-free content on the internet,
21 those are important reasons why that industry might do what
22 IBM did, you know, in the 1980's and say, 'Government hands

1 Mr. Swire: Yeah, I think just a moment on the
2 government side because when I worked at OMB we oversaw
3 federal agency use of data under the Privacy Act. The
4 Privacy Act has this very antiquated idea that there's
5 A thing called an agency and everything you do inside one
6 federal agency is okay, but if somehow it crosses a border
7 into other agencies then you need a routine use or an
8 exception.

9 Just a couple observations. One observation is
10 there is this new agency -- happy fifth birthday Homeland
11 Security -- that's huge, right? But that's a really big
12 agency and everything can go everywhere in Homeland
13 Security and we laugh at the Privacy Act because it's just
14 one agency so there's no problem. So that doesn't seem
15 like exactly what the 1974 Congress had in mind.

16 And the second thing is, if you talk about the
17 information sharing environment and the sort of presumption
18 of sharing in federal government, those agency boundaries
19 don't seem to match very well. So I think that within the
20 government side the Privacy Act doesn't match up with
21 privacy risks today. It's been very hard to figure out
22 what alternatives to do about that, but I think that's a

1 big challenge on the government side.

2 Ms. Landesberg: Thanks very much. I'd like to
3 turn to first Chris, and then Fred. Are there types of
4 data mining that pose little or no risk to privacy?

5 Mr. Slobogin: I will answer that question yes
6 right now, but I reserve the option of changing my answer a
7 little bit later.

8 Mr. Cate: I think there are -- this -- my
9 answer, I guess, will overlap a little bit with the panel
10 that's going to take place tomorrow in terms of best
11 practices. I think one type of data mining that would
12 avoid privacy concerns is data mining that's justified.
13 And what do I mean by that? Well, a good answer would take
14 me an hour; I think the panel tomorrow will talk about it
15 quite a bit. But the bottom line is that the government
16 shouldn't be able to engage in target-based, match-based,
17 or event-driven surveillance -- data mining surveillance,
18 unless it has good reason to suspect the people who will be
19 pinpointed by that surveillance. What would good reason
20 be? I think if it's very intrusive kind of data mining
21 that is going into medical records and personal financial
22 records, it should require probable cause or the

1 (inaudible) thereof. If it's -- if the records involved
2 are less personal, perhaps only reasonable suspicion --
3 those of you who know Fourth Amendment law know these
4 phrases that I'm throwing out. The bottom line is that
5 there should be very good justification. If there is
6 justification, then I don't think there is a privacy harm
7 because I think if the government has a good reason for
8 needing it, then privacy can be relinquished.

9 Another kind of data mining that I think probably
10 does not impinge on privacy -- though, again, I might
11 change my answer later -- is data mining that goes after
12 purely public records, it uses only purely public records
13 as its data source. The reason I'm hesitant in saying that
14 is that even public records, if aggregated, can produce an
15 awful lot of information about an individual. There's
16 quite a bit in the literature about the fact that one
17 record by itself isn't very privacy invasive, but
18 aggregating records, even if they're all from public
19 sources, can take a whole lot about a person, they can be
20 very invasive. Because after all, public records include
21 real estate records, voting patterns, employment records,
22 licenses, and so on. You can get an awful lot of

1 information just from public records.

2 Finally, I guess I would say -- and this picks up
3 from what Dr. Jensen was talking about earlier -- if we can
4 anonymize data mining, if that's technologically feasible,
5 then I could see the kind of multi-stage process he was
6 advocating being pretty protective of privacy. But right
7 now I think it's fair to say we don't have very good
8 anonymization techniques. So I think that's more of a
9 hypothetical situation where data mining would not invade
10 privacy in any significant way.

11 Ms. Landesberg: Thank you. Fred?

12 Mr. Cate: I think the answer is yes and no, and
13 I'm willing to stand behind that. I think our discussion -
14 - and frankly, I think this goes back again to one of
15 Professor Jensen's slides -- has really highlighted the fact
16 in some way we're talking about three pieces of a system.
17 We're talking about the data, we're talking about the
18 analysis of the data, and then we're talking about the
19 consequences what's done with the data. And so you could
20 imagine depending upon where you are in those three that
21 there would be types of data mining that would not raise
22 significant privacy issues. So, for example, if you use

1 public record data, you do very good analysis with it, and
2 it has very minimal consequences, I think, you know, we
3 might be able probably not to agree on this panel, but a
4 reasonable group of people might agree that that does not
5 raise significant privacy issues. On the other hand, even
6 in that situation, if you used highly sensitive data, you
7 did really bad analysis, whether or not the consequences
8 were serious, we might think of it as nevertheless raising
9 serious privacy issues. I think the way in which probably
10 the public -- with whom I feel like I have the most in
11 common here -- tends to think about these issues is focused
12 on the consequences side. In other words, what is the
13 consequence of this data mining? Is it, you ask me another
14 question? Is it, you use that little swab on my suitcase?
15 Is it, you shoot me dead in the airport? Depending upon
16 the answer, I'm going to have a much greater sense of how
17 this burdens me and burdens society. But I think actually
18 the discussion's been helpful and we should not lose sight
19 of the fact that even if the consequences are themselves
20 trivial, if the analysis isn't good or the data are not
21 accurate or are mismatched or not relevant for the purpose
22 for which they're being analyzed, we may still be wasting

1 very scarce resources, we may still be distracting
2 ourselves from the mission at hand, and so we should not be
3 focused just on consequences when thinking about what are
4 the implications or the impact of the data mining.

5 Ms. Landesberg: Thank you. Did you want to say
6 something, Greg?

7 Mr. Nojeim: Not on that.

8 Ms. Landesberg: Okay. Very good. And I'd now
9 like to turn to Peter again if (inaudible) -- thank you.
10 What do we need to take into account to determine whether a
11 data mining project is effective in identifying future
12 terrorist or criminal activity? How do we weigh privacy
13 risks against potential benefits in counter-terrorism
14 research?

15 Mr. Swire: I'm going to take that and do
16 something just slightly different than that. But I was at
17 a conference last week where a military person was back
18 from Iraq, and he was talking about IFF - Identify Friend
19 and Foe, and I wanted to talk about how all of the data
20 mining looks if you're in the middle of a battlefield and
21 how it looks if you're, let's say, sitting in a nice hotel
22 in Washington, D.C. -- which you think is not a

1 battlefield except getting coffee today.

2 So IFF, the traditional Identify Friend and Foe
3 thing is you're a naval ship, you're floating along in the
4 ocean and there's an airplane coming towards you. And at
5 that moment it's very relevant thing to know, is that an
6 attacking airplane from an enemy or is that a friendly plane
7 that's just coming by to, you know, to visit you. And
8 during the Cold War, for instance, there was elaborate
9 technology to try to figure out if that was a Soviet
10 airplane coming. And you needed to do that if you were on
11 alert or if you were in Iraq recently, if the convoy of
12 trucks is coming at you and it's an attacker or else it's
13 not; you need to do that on a hair-trigger because if they
14 get too close they can blow you up, and that's a very
15 disappointing outcome. So you really need to figure out,
16 is this a friend or foe, and there's high stakes and you
17 have to do it quickly. And so if you talk to defense
18 intelligence agency folks, if you talk to people in the
19 middle of a battlefield, this idea of really identifying
20 high-risk or low-risk and doing it immediately is life or
21 death and we want our military people to have really great
22 analysis, ability to answer quickly, all of that. And

1 listening even to the science and technology lead early
2 this morning, I think his model in a lot of ways was an IFF
3 kind of model, which is, you remember that sort of
4 increasing likelihood of bad outcomes and how likely is it
5 -- how big is the magnitude? And it's, like, as people
6 come into the borders, as people do various things, we want
7 a risk score on each moment, on each activity, on each
8 person that's coming to our realm. And if -- when you're
9 at war, that's a highly relevant way to think about things.
10 And we are at war in Iraq right now, we're at war in
11 Afghanistan right now; people's words about how much we're
12 at war at home in the United States vary. More often, my
13 experience in Homeland Security, we say, 'We're at war
14 right now,' and if you walk up to people at the mall and you
15 say, 'Are we at war right now here in the United States?'
16 they don't tend to act and feel like they're at war right
17 now. So people's war analogy varies. But here's what I
18 want to say. What's different maybe about IFF when you're
19 at the train station at Union Station or even the airport
20 or walking down the street in Washington, D.C.? There's a
21 lot of things that are different, and one is, most of the
22 people are not on the edge of attacking you, right? So

1 most of the time we walk down the street, it's not that
2 hair-trigger, is this a bomber? Is this a convoy that's
3 going to blow us up? Another thing is the scale of how
4 many Soviet aircraft types or how many different categories
5 of convoys -- that the scale is in dozens or hundreds in
6 the war zone. There are 300 million Americans, and so the
7 scale is really different and we don't have anywhere near
8 the same ability to go from -- this is Barry's point about
9 -- you know, from very low likelihood of harm to, out of
10 300 million it is. So the hair-trigger is different and
11 the scale and magnitude's different. And then the other
12 thing is, there's a lot of reasons not to do Identify
13 Friend or Foe in civil society. So we don't want to sort
14 of get a complete profile of every political thing you guys
15 have read on the internet and come up with a risk score
16 about how likely you are to be in opposition to the
17 government. That's a different society than I want to
18 live in. We don't want to have somebody who has a
19 jaywalking ticket or a marijuana bust 35 years ago; I'm not
20 sure how much we want to have them treated entirely
21 differently as they walk through society. But in a risk-
22 score world all of those things might be in bounds.

1 And so if we want to have a risk-score on every
2 moment in society as we go in the non-war zone, as we go
3 through the United States, I think that's just an entirely
4 different society than the sort of presumption of freedom
5 society, presumption of openness, presumption of allowed to
6 critique, presumption that we're not (inaudible) on First
7 Amendment grounds. And so I think that in a lot of ways
8 the data mining question is, how much are we in that
9 warzone where we want our Naval ships to know if it's a
10 bomber, and how much are we in a peaceful zone where
11 there's a background risk but our hair-trigger is so
12 different that lots of the same mechanisms we use in the
13 warzone we don't use in safe, mostly peaceful society?

14 Mr. Cate: May I just add one comment? I think
15 that's just an excellent analogy and it reminds me of
16 another distinction, which is, in the Identify Friend or
17 Foe in a military environment, it's pretty clear what the
18 harm is you are protecting against and it's pretty clear
19 how that harm is threatened. So, for example, in the
20 example of the airplane approaching the ship, we only ask
21 if it's an airplane and not a fish. We're, you know, we
22 have a pretty focused way of applying it. In, of course,

1 our daily lives that's not true. I mean, you know, it's
2 not at all clear what the harm is we are guarding against,
3 and it's not at all clear what data are predictive of that
4 harm. So again, in the example of, for example, knowing
5 reading patterns or browsing patterns or what have you, it
6 would be interesting to know, you know, are terrorists well
7 known for speaking out in open public events before they
8 engage in a terrorist act? Is that predictive in any way
9 to know what they're protest habits are? If not, why are
10 we bothering collecting that data? Right? If it's
11 irrelevant data, why bother with it? So I think part of
12 the concern here is that we are investing in collecting
13 data, or we are making decisions based on data, that in
14 fact has no probability whatever of predicting the thing it
15 is we are worried about. And so, again, we're left with
16 this sense of, we are not only talking about what affects
17 or harms the individual, but what affects or harms all of
18 us much more broadly. And that that is -- that should be
19 of major concern. I mean, that should be a daily concern.

20 Mr. Nojeim: May I just add one -- I think it's
21 even -- that's a very good way to put it. And I think the
22 problem is even bigger than that, so you could figure out

1 what the data was, what data sets would be useful about the
2 very small number of terrorist activities that you could
3 point to today. You still wouldn't know what the next
4 terrorist act would look like or what the people who would
5 do it, what kind of data set might be relevant to them. So
6 it's -- I think it's actually an even tougher problem than
7 you describe because you can't fight yesterday's war.

8 Ms. Landesberg: Go ahead, Chris.

9 Mr. Slobogin: Let me play devil's advocate for a
10 second. Let's assume that everything that was just said is
11 correct. We still are talking about the possibility of one
12 terrorist wiping out a major American city, and given that
13 threat, why wouldn't it be okay to let designated agency,
14 say, the Department of Homeland Security, have all the
15 information it wants about anything with the caveat that
16 that information be retained within the agency -- a small
17 group within the agency -- and that it only be used to
18 prevent clear terrorist acts? What's the problem?

19 Mr. Cate: I'm so glad you asked that. Well,
20 first of all, nobody would ever buy into the conditions
21 that you just said. Not a person in this room would
22 believe that for an instant that they would only be used

1 intelligence out of the data. This is a little bit like
2 U.S. News rankings of law schools, you know, they still
3 count the volumes in the library. Like having more data is
4 a great thing. Whereas we all know that the goal here is,
5 can you extract useful intelligence from data fast enough
6 to make it practically useful.

7 Ms. Landesberg: Okay. If I might -- thank you.
8 And I can hear considerable skepticism about whether data
9 mining can be effective for the purposes outlined here, but
10 I am interested in knowing what you think it would
11 take to determine whether a project is effective or not --
12 what the analysis ought to be. Anybody want to tackle
13 that?

14 Mr. Cate: I'll take the easy ones and then leave
15 my colleagues the hard ones. I think one thing we would
16 like to see is a stated purpose in advance, and then
17 testing against that purpose. Because one of the most
18 common things we see -- we see it also in PhD dissertations
19 as well -- is, you do the research, you didn't at all come
20 out with what you thought you were so then you change what
21 was the topic that you were researching. And therefore,
22 before we invest public dollars intended to fight

1 they think it's really, really great. And then before the
2 deal actually happens, you have to have due diligence; you
3 have to have other people go in and say, 'Wait a second,
4 don't you realize most of these assets have already been
5 foreclosed on? You know, maybe that's not such a good
6 thing to buy.' So due diligence is the process of having
7 smart analysis before the -- you have the enthusiasts who
8 are trying to go forward with a new thing and then you have
9 other people saying, 'Wait a second. Let's see if this' --
10 and so without going through the -- reading all ten items
11 because it's in they're in your chart and it's based on a
12 longer article. The first part is something this panel's
13 stressed a lot which is do we have some reason to think
14 it's going to improve security? Is it -- even if the
15 project worked out, is it going to lead to some payoff? Is it
16 going to be doing it cost effectively? Is the program
17 going to hurt security by spreading information to the bad
18 guys? And Dr. Jensen this morning didn't want to tell us
19 exactly where Fraud Alley was; I bet a lot of us were
20 sitting here thinking, 'Hey, I wonder if that's' -- I don't
21 know, I thought of Miami, I thought of some parts of New
22 Jersey. You know, I don't -- and maybe you all thought of

1 different places. I used to live in New Jersey; I'm not
2 against New Jersey, but I knew some things there. You
3 know, so -- but he didn't want to tell us because then the
4 next fraudster won't set up in Fraud Alley. They'll eschew
5 those three-digit codes and they'll set up somewhere else.
6 And so there's all this cat and mouse kind of thing. So in
7 my article and in the ten-point list, there's some attempt
8 to try to say, 'What are the problems?' One thing I'll
9 highlight is number six, "Do fairness and anti-
10 discrimination concerns kick in?" Here's a current example
11 - there's a hearing recently in the House on a data mining
12 thing done by the insurance agency. It turns out, in the
13 insurance business, I can do a better job predicting your
14 insurance risk based on your credit score. And so the
15 question has been, is it a good idea/bad idea for credit
16 scores to be used for your car insurance. In the hearing
17 there was discussion that there's a correlation between
18 race and credit scores. So if this started to be used,
19 certain racial groups would pay more on average for car
20 insurance. So you'd have a benefit, which is maybe more
21 accurate person-by-person decisions about how much to
22 charge for insurance; we'd have a more efficient insurance

1 market. And you'd have a sort of fairness question of, if
2 predictably this is going to raise premiums for certain
3 racial groups, is it okay or not? And there was -- it was
4 pretty heated debated in the House Committee -- Financial
5 Services Committee about what to do on this. But it
6 illustrates something far away from terrorism where you get
7 results from data and then you have to work through, 'Okay,
8 what are we going to do with this?' and a due diligence
9 process is at least one way to try to head at that.

10 Mr. Nojeim: (Inaudible). Can I chime in here?
11 It's really to ask Barry a question 'cause it goes to, what
12 do you do with data that might be relevant, might be
13 useful? Say it's not Fred's example where the NSA is
14 getting hundreds of millions of bits of data, say it's this
15 example - you've got a different NSA, it's the smart NSA,
16 it's the focused NSA, it's the targeted NSA. It has three
17 terrorist phone numbers abroad, that's all. It's been
18 watching them, it knows that they're bad guys, and it also
19 can collect information about who those terrorists call and
20 who call them. And all three of these terrorists talk to
21 somebody in the United States; what should the cops do with
22 that information? Should they show up on that guy's

1 doorstep and say, 'What are you up to?' Or, what do they
2 do with it? That's --

3 Mr. Steinhardt: Well, now I'm glad you asked
4 that question. You know, actually I think that's a
5 relatively easy question, right? Which is to say you have
6 your -- you made it a little more complicated by mixing in,
7 you know, Foreign Intelligence Surveillance Act or domestic
8 eavesdropping laws, but, you know, but basically that
9 question is one we know the answer to, right? Which is
10 that, you know, the government has a lawful right to obtain
11 the numbers that were called by a particular individual, or
12 they have a pen register or whatever it is. And they know
13 who that individual calls; can they do some follow up
14 investigation of the individuals that were called? I think
15 the answer to that is, you know, usually is yes. That's
16 not exactly data mining, right? I mean, you know, except
17 in the sort of broadest sense of the word. You know, it's
18 so the way that, you know, that law enforcement follows
19 leads generally. And, I mean, I'm a little, you know, I'm
20 a little less troubled with that than the notion that we
21 are going to intercept everybody's telephone calls and try
22 to make fairly attenuated connections between individuals

1 based on their pattern of calling, as opposed to this
2 fairly, you know, discreet set of facts that you described.

3 Ms. Landesberg: Okay. If I might just ask --
4 Greg and Barry, are there -- we've gotten a good record I
5 think from this workshop already, but are there more
6 specific harms that either of the two of you would like to
7 address before we turn it to the audience for questions?

8 Mr. Nojeim: I wanted to talk a little bit about
9 commercial data for just a sec. And increasingly the
10 government is using commercial data in its data mining
11 activities, and there's nothing inherently evil about
12 commercial data as opposed to data that's been generated by
13 the government.

14 But I did want to mention a couple of concerns
15 about it because I think that using commercial data should
16 be done very cautiously. The first is that the data is
17 collected for a particular commercial purpose, and it might
18 be the case that in pursuing that commercial purpose, some
19 problems with the data would be ones that you wouldn't want
20 to expend the necessary resources to correct. And there
21 might be accuracy issues within the data -- just say for
22 example, its' credit data -- it might be the case that for

1 you to fix all the problems in your credit database -- and
2 I saw one estimate that 70 percent of credit reports have
3 an inaccuracy -- but to fix all that it might be
4 prohibitively expensive, and therefore it might be an
5 appropriate model for you to sit back and wait a little bit
6 until a person contacts you and complains about that data
7 and then fix it after you receive that input that there's a
8 particular inaccuracy.

9 That model, that kind of data might not be
10 appropriate to be using to predict who might be a
11 terrorist, or to match with other data about terrorism
12 because it doesn't have the accuracy level that you would
13 need for that data to be effective.

14 And the second thing I wanted to stress besides
15 this use issue, was that data in the private sector isn't
16 subject to Privacy Act restrictions. And Peter has
17 outlined some of the problems with the Privacy Act, and
18 they're substantial, but it does provide some protection
19 for people who are subjects of that data. And so one of
20 the issues -- one of the Privacy Act protections is -- goes
21 to accuracy, errant inference, and those kinds of things
22 are things that I think an agency relying on commercial

1 data would need to account for.

2 Ms. Landesberg: Thank you. Barry?

3 Mr. Steinhardt: Yeah. I actually wanted to pick
4 up on something that Peter spoke about, which is, sort of,
5 what is the harm or what is the consequence of living in a
6 society where we are all risk-scored? Which is a society
7 that we are increasingly moving toward, and I think that
8 the consequences of that are fairly profound in people's
9 daily lives. It's not simply that there's the risk that
10 you're going to be arrested or interrogated, it's the risk
11 that you are not, for example, going to be able to obtain a
12 mortgage or a bank loan, that you're not going to be able
13 to get a job. I mean, all those things are increasingly
14 becoming real for people. You know, if you go in now to a
15 bank to open a new account, your name is checked against a
16 government list to determine whether or not you, you know,
17 you might be one who is engaged in, you know, criminal
18 activities, terrorists, et cetera.

19 One of the things we know about that list is that
20 it, you know, essentially it's a new form of risk scoring,
21 right, or a risk not that you -- that, you know, that
22 you're going to be a deadbeat or something, but rather that

1 program everybody involved in terrorism detection and
2 everybody in Congress, and, you know, immediately have them
3 investigating each other instead of investigating real
4 problems.

5 Ms. Schiller: My name is Jennifer Schiller; I
6 work for Under Secretary Cohen as his Privacy Liaison. So
7 we are focused exclusively on research development,
8 testing, and evaluation activities.

9 Mr. Steinhardt: Could you just get a little
10 closer to the microphone, please.

11 Ms. Schiller: Sure. None of our current
12 programs meet the Congressional definition of data mining,
13 but we do want to move forward with that vein of research,
14 you know, looking at our long-term planning. And as we
15 enter the testing and evaluation phase of developing new
16 technologies, we do need to use real data to test the
17 technology before we can, in good faith, transition it to
18 an operational component that would then go and use it. So
19 my first question is, what factors should we consider in
20 evaluating the impact of privacy in that type of research
21 where we're not making determinations about individuals,
22 we're testing the operation of a technology prior to

1 transitioning it to an operational unit?

2 And the second question I have is that we seem to
3 have two broad categories for data analysis and data mining
4 activities. The first would be what I just described where
5 we're developing a new technology that would eventually
6 transfer to a customer; the second would be where we're
7 looking at data in a social science type of way -- and we
8 do have one of our social science researchers here, I hope
9 she'll ask some questions at some point -- but we're
10 looking at a broad set of data, for example, on terrorist
11 events or on terrorist groups and trying to draw inferences
12 from that data. An increase in rhetoric might be a
13 signifying factor before an event. And Professor Cate, you
14 said, 'Why collect the data if it's not relevant?' Well,
15 we don't know if it's relevant until we go in, collect the
16 data and do the research. So my second question would be,
17 how can you handle collection of data in a research
18 environment where you're not sure what data is relevant?

19 Mr. Cate: Well, let me say I think you raise a
20 phenomenally important issue, and that is the need for data
21 on which to do research. And it's an issue on which, to be
22 perfectly frank, Congress has been completely tone-deaf,

1 not to mention ignorant, and that is, it is not appropriate
2 to use the same types of privacy protections for data that
3 will be acted upon, as opposed to data that's being used in
4 a research environment. Having said that, because of the
5 policy of law in this area and the fact that most of the
6 law in this area is, I mean, has no real restraining
7 effect, there's no possibility to make a promise like,
8 we're going to have this data but not act on it, because
9 there's no legal requirement that would bind you to that.
10 You would have to enter into a, I guess, a contract with
11 the American people that said, this is what we're going to
12 do. So in some ways, I hate to say it given that I think
13 Congress has been a major source of the problem, but I
14 think they're also going to be an essential part of the
15 solution, which is to create a category of data analysis or
16 data mining or data aggregation for research that has to
17 meet certain conditions and would be subject to certain
18 oversight and so forth. One of the things we haven't
19 talked about on this panel at all just for lack of time and
20 because I think other panels will, are the procedural and
21 process protections that can diminish the potential harm
22 caused by data mining. But I think, you know, I mean, you

1 would know those as well as any of us, and that what's
2 needed is a way to get those -- if you will lock those into
3 place so that once a program is declared a research
4 program, and I would continue to disagree with Peter on
5 this. I think when you spend of hundreds of billions of
6 federal dollars you better know what your goal is before
7 you start, rather than the, let's hope we find it does
8 something once we've spent this money. So you say, 'This
9 is research. Period.' And then you are under that
10 protective regime; I think that's going to have to be what
11 the solution is going to ultimately look like, and it's
12 going to mean going to Congress.

13 Mr. Swire: I actually think we know quite a bit
14 about research from the medical side of things. So for
15 medical -- I worked in HIPAA and the research parts of the
16 HIPAA Medical Privacy Rule, and so there's things in HIPAA
17 about limited data sets, about data use agreements, about
18 what kind of audit and oversight there's supposed to be
19 before the research is approved. In most circumstances it
20 goes to an IRB, an institutional review board. I'm not
21 sure how much all of that exists yet in research at DHS,
22 but there's been a lot of decades of work on research from

1 the medical side, and that at least gives you some
2 institutional things to look at as you're trying to figure
3 out.

4 And then in terms of how you do it legally, a
5 statute would be better, but I think DHS could say, 'When
6 we do research we're going to do it under this sort of IRB
7 medical approach when you're working with real people's
8 real data. And you could say, 'We're making a promise to
9 follow the following guidelines,' and then you could say,
10 'And we're going to have our IG come in on a regular basis,
11 or GAL or whatever,' to make sure it's being followed, and
12 that would -- well, even without a statute, give a pretty
13 decent institutional basis for the world to believe you're
14 actually doing it.

15 Ms. Hahn: Thank you. Again, my name is
16 Katherine Hahn with SAS. I appreciate your comments. But
17 I want to go back to the question that I asked Professor
18 Jensen this morning. You all have talked a lot about data
19 gathering, data collection, problem articulation; is the
20 privacy concern raised by the statistical model or is it
21 raised by the human intervention where you can't test the
22 bias and the assumptions that people are bringing to bear?

1 And as a follow-up to that, what is it about data mining as
2 a statistical modeling activity, as a research methodology,
3 that merits heightened privacy scrutiny over other types of
4 research methodologies? Thank you.

5 Mr. Steinhardt: Let me take a crack at that if I
6 can, in a couple ways. First, I think that the answer to
7 your question is one I said earlier, it's both yes and no.
8 The, you know, the issue here is not only whether or not
9 the collection of the data and the analysis of the data
10 represents a problem, but also whether or not the automated
11 data presents a problem and how that data is used by
12 individuals. I do think it's important to recognize,
13 though, that when it comes to the use of data by the
14 government, the government is not in the same position,
15 really, as for example, the security dealers examples that
16 we were given earlier this morning, where the security
17 dealer, you know, the NSD has this sort of unique ability
18 to collect data about the individuals that it governs. It
19 already has that data; it can use that data to mine into
20 it. But the government on the other hand, we found this,
21 for example, in the area of airline passenger profiling.
22 The government didn't in fact have very much in the data

1 that it wanted to use to mine, whether it was an automated
2 process or a process driven by personnel. That was really
3 the great -- that has always been the great debate about
4 the various versions of airline passenger profiling -- how
5 much data would the airline industry need to collect?
6 Often data that it does now collect from passengers in
7 order to give the government the ability to in some way
8 mine that data or make use of that data. That's really, I
9 think, the important thing to look at here as we're
10 looking at government data mining, government use of data,
11 which is, where does the government collect that data from?
12 How much data does it need to collect? And then ask the
13 question of what it's going to do with it, but to recognize
14 that generally speaking, there is the necessity to go out
15 and collect data in all -- and usually it's the private-
16 sector collecting that data for the government -- but to go
17 out and collect data that is not now collected and not now
18 analyzed.

19 Ms. Landesberg: Okay. Did you want to respond?

20 Sure.

21 Mr. Nojeim: I just wanted to say a couple of
22 things. One is that I don't think that you can neatly

1 divide up the two functions between the data mining and the
2 consequences that follow from the data mining. You
3 don't do the data mining in the first place unless you're
4 looking at what to do with the data that you get, with the
5 results that you get. So I just don't think that you can
6 isolate that kind of technical activity from the
7 consequences because that's the whole purpose, is to decide
8 -- to make decisions about people. And the second thing is
9 that people who are involved in the actual data mining
10 activity can build in some of the protections that we've
11 been talking about here. For example, audit trails -- I
12 mean, there is -- you do want to have that capability built
13 into the system that you're coming up with so that you can
14 find out whether the data was misused and whether it was
15 appropriately used.

16 Ms. Landesberg: Thank you. And if I could ask
17 those of you who are waiting to ask questions to just be
18 very concise in the question so we can get you an answer
19 and then we'll adjourn for lunch when Dr. Jensen's had his
20 chance weigh in.

21 Mr. Schneiderman: I'm Ben Schneiderman from the
22 University of Maryland. I'm troubled by the narrowness of

1 the definition of data mining, exemplified maybe by
2 Christopher Slobogin's question of what's wrong with
3 collecting data, and the last few answers did get to the
4 question of the socio-technical system that's imbedded in
5 the cost of collecting it as opposed to the benefits that
6 might come from other things, the distraction that's
7 brought by it. But then the -- I guess it's Greg last, you
8 know, comment about the end-game, also what happens once
9 you get it. So what are the -- how do we expand the
10 definition of the systemic view that will give us a socio-
11 technical analysis that will give, for example, citizens
12 whose privacy was violated, recourse and compensation,
13 which is not part of the TSA's current No Fly List. If
14 you're prevented from flying, you don't get compensation.
15 So if you put in the true costs to all the parties that are
16 harmed, you have a better chance of understanding what the
17 payoffs and the negatives are more clearly. So I'm looking
18 for -- the questions about what are the broader aspects
19 that you see to the socio-technical system?

20 Mr. Slobogin: Well, this is really follow-up on
21 the last answer, but it seems to me we should get away from
22 using the word data mining if that's your major concern.

1 If you wanted to define data mining the way Professor
2 Jensen did, fine. Then we've got problems with the data
3 collection and with who gets the results of the statistical
4 model and what's done with it? You can assign labels to
5 those different stages of the process. I have to admit,
6 I'm more concerned about the data collection, who gets to
7 see the results of the data -- of the statistical modeling,
8 and what's done with it. Those are my major concerns. The
9 actual technical aspect of statistical model is not a major
10 concern of mine. I think most people apply the word data
11 mining to all those stages; it might be better to break
12 them out into the three, four, five stages, and then focus
13 in on the legal and social consequences of those stages.

14 Mr. Steinhardt: Can I -- let me take a crack at
15 that, too. I think that's an important question, which is
16 what recourse do you have if you are harmed? One of the
17 real drawbacks to the current, you know, airline passenger
18 profiling system and the watch lists, et cetera, is that
19 there is no real recourse for those people who find
20 themselves on the list by mistake, find themselves harmed.
21 There is a sort of a, you know, kind of a Kafkaesque system
22 of going through the Department of Homeland Security puts

1 you on the list in the first place, and I can take you off
2 and you'd never know if in fact you're off and how you got
3 on and all those things. But, you know, there's a fairly
4 simple solution to that problem, which is to say that we're
5 going to have an independent body out there, whether it's a
6 judge or -- but some other mutual arbiter, right, who is
7 going to take a look at the data and is going to say, you
8 know, that is or is not someone that we -- to be worried
9 about, -- if not, we're taking them off the
10 list and we're ordering all the people that maintain that
11 list to take the person off the list, or to the extent to
12 which they have a name that is, for example, is the same as
13 somebody who should be on the list. We're going to create
14 a white list or put this person on the white list so, you
15 know, you're not -- we have none of that. And I'm not
16 sure, monetary compensation might be nice, but I'm not sure
17 that that's the solution. I think the solution is we've
18 got to have processes in place that allow people to appeal
19 to an independent arbiter, decisions that are being made
20 about them. But the first thing we need people to do is to
21 indicate to that person that, yes, you are on the list or
22 yes you are affected.

1 Mr. Schneiderman: My point is that true costs
2 are only visible when you have the larger socio-technical
3 context. Thank you.

4 Ms. Landesberg: Thank you. And, sir?

5 Mr. Lempert: Yeah. Rick Lempert, I'm with DHS.
6 Two very quick points. One, I was very happy to hear
7 mentioned the commercial issues at the end. And as you're
8 speaking personally, in some ways I'm more concerned about
9 the commercial invasions than I am about government. We
10 saw that Admiral Poindexter's plan -- what happened
11 politically -- there are likely to be political limits on
12 what the government can do. And we see with the IRS data
13 that the government can be very, very protective of
14 confidential information even when it could be used for
15 other governmental purposes. Doesn't mean there aren't
16 governmental concerns.

17 The other issue is -- in thinking about this, and
18 I agree the costs are very important -- the question is,
19 what is the alternative in the non-data mined universe as
20 we think about it? So, for example, if -- I may not like
21 any broad-based surveillance techniques at the border,
22 perhaps, but if I had the choice between these lists of

1 here.

2 But, you know, I -- one of the things we haven't
3 really talked a lot about this morning -- talked at all
4 about this morning -- that I want to sort of emphasize is
5 that all the security, you know, is a zero-sum game. We
6 only have so much money here to spend on our security, and
7 what we have not talked about this morning are some of the
8 alternatives to all of this data analysis -- whatever --
9 however we want to characterize it and all this list making
10 and all this. I mean, we know, for example, that, you
11 know, that often the most effective weapons that we have
12 against terrorism -- terrorist attacks, are physical
13 security. I mean, we know, for example, that hardening the
14 cockpit doors after 9/11 made a tremendous difference. You
15 cannot replicate what happened on 9/11; you can't get into
16 the cockpit now. Can't use the plane as a missile. We
17 know, for example -- take the example -- the U.K. example
18 of the terrorists who were the London Underground -- tried
19 to set off a bomb, turned out they weren't very good at it,
20 but the -- you know, but there was all this sort of, you
21 know, hyped up surveillance equipment around them, all this
22 -- all the video cameras and whatnot; that didn't stop

1 them. But they went up to the airport in Scotland -- now
2 they decided that they were going to drive a car into the
3 passenger area -- passenger terminal -- what stopped them?
4 The concrete barrier out front of the passenger -- out
5 front of that terminal; same concrete barriers you can find
6 all over Washington in front of government buildings.

7 So we have to ask ourselves here, before we
8 continue to spend all this money and all these resources
9 and risk our liberties and our privacy on these systems, is
10 this really the most effective way for us to be fighting
11 the so-called War on Terrorism?

12 Ms. Landesberg: Thanks very much. And, Dr.
13 Jensen, I think we should allow you to have the last
14 question here and then we'll adjourn for lunch.

15 Mr. Jensen: So it may surprise some of you that
16 I think we almost totally agree, but I want to do a
17 find/replace on all of your comments, which is I think
18 possible because of the court reporter -- which is, to do a
19 find/replace and replace data mining with data collection.
20 So here's the question -- the question is, like you, I am
21 against government incompetence, I am against violation of
22 civil liberties -- I'm an ACLU member, by the way, dues

1 paid --

2 Unknown Male: We'll check our lists later.

3 Mr. Jensen: -- I'm against government power
4 grabs, I'm against prejudice, and I'm against
5 authoritarianism. But would those harms stop or
6 significantly lessen if we completely gave up data analysis
7 but kept doing data collection just like we're doing now?

8 Mr. Swire: I've gotten to be in a lot of
9 different privacy-related meetings over the last bunch of
10 years, and the previous question was -- but I'm going to
11 summarize it, maybe a little unfairly -- I work in the
12 government, the real problem is in the corporate sector. And
13 when I talk to the corporate folks, they all say, 'We're
14 good. I know all of our people are really good; I worry about
15 the government.' We have a data mining statistical person
16 saying, 'Data mining's good but those collection people are
17 nuts.' And then you'd go talk to the police at the
18 collection point and say, 'Look, we collect things. It's a
19 world of cheap sensors, we have to get the data we can get,
20 but it's what they do with it once they collect that's the
21 real problem.' And I'll just -- having -- I'll just
22 observe as a sociological phenomenon that people tend to --

1 there's aversion to taxes, right? Don't tax me, don't tax,
2 you know, you -- tax that person behind the tree -- and the
3 aversion in privacy is, the part I'm in, we really don't
4 want to have these intrusive rules stopping what's
5 important to do, but it's those folks over there, that's
6 what you have to watch.

7 Mr. Jensen: Yeah. But governance is about
8 making choices, okay, and ascribing causality correctly.
9 So I do think it really matters.

10 Mr. Swire: Oh -- it matters to do a sensible
11 analysis on security; what's the tradeoff between physical
12 security and intelligence gathering ahead of time? It
13 makes sense to figure out what are the real risks on
14 collection versus analysis versus actionable. All of that
15 is a logical part, but I -- just observing that the
16 previous two questions ago is the socio-technical system
17 and the political system in the broad view, it turns out we
18 don't have these neat compartments where we can say, 'This
19 is a risk-free zone and the problems are over here.'
20 You look at each part of the system and you keep having the
21 meetings because it turns out these pieces are
22 interrelated.

1 Mr. Cate: I'd just like to be clear, I'm not
2 against authoritarianism. Your question kind of makes me
3 wonder how badly we've done for the past hour-and-a-half up
4 here, because I don't think there was any suggestion that
5 data mining is inherently bad or data mining is inherently
6 -- should not be used or should be avoided in favor of some
7 other technique. It's that data mining, like all of the
8 other tools we use in fighting terrorism, should be
9 subjected to the same type of scrutiny. And that the more
10 we can break it down into its constituent parts, the more
11 we can be clear about the data analysis and the data
12 collection and aggregation, the data mining tools that are
13 used and the consequences of what's done with those, the
14 more frank we can be in that type of analysis, the better
15 the results are likely to be. And I think in response to
16 the prior comment as well, there are many instances where
17 data mining is by far the preferred tool. I mean, it is
18 the equivalent of putting up the concrete barriers, in some
19 instances. It makes perfect sense, it's cost-effective,
20 and if done well may have little negative impact on
21 individuals. I think part of the problem is that the
22 dialogue about this wide range of data analysis activities

1 the decision could be a manual process, it could be an
2 automated process, it can be anything. But if you are
3 going to analyze -- I mean, I'm an engineer; I like to take
4 data, test any system whether it's human or otherwise, and
5 produce results that come out of the area of debate. They
6 are scientific engineering results. This is one way to
7 display information about any binary decision process. It
8 came out of designing communications devices. So the basic
9 way was, if you're sending a 1 or a 0, how is your
10 receiver detecting that 1 or 0? It can be generalized to
11 any binary type decision, so is something good/bad, is it,
12 et cetera?

13 And just to set the stage to understand what this
14 graph is telling you, if someone asked me to design a
15 system that will always detect whatever it is you're trying
16 to detect, I can very easily do that. If you see that
17 point up there at the 1,1 position, the way that you read
18 that graph on the bottom -- the x-axis -- that's the
19 probability of a false alarm. The y-axis is the
20 probability of a correct detection. So if I just always
21 say, 'Yup, that's it; that's what I'm looking for,' then
22 I'm at that 1,1 point. I will never miss any of the things

1 I'm trying to detect; however, I will always make a false
2 detection any time that it's not what I'm trying to detect.
3 So that's one extreme of the detection space. Well, let's
4 say you say, 'Well, no, I don't want to do that; I want to
5 never make a false alarm.' Well, that's easy too, down at
6 the 0,0 point. I always say it's not what I'm looking for,
7 so I'll never make a false alarm. Well, that's also kind
8 of useless.

9 Now, if I just flip a coin every time, I can be
10 along that dotted line between those two points. So, if I'm
11 going to analyze a system, I want to determine, am I on
12 that line -- in other words, I'm doing any better than
13 chance -- am I below that line, where you see the red
14 graph. Because believe me, you can spend a lot of money to
15 design a very complicated system and be somewhere in that
16 red space where you're doing worse than chance. Or am I
17 somewhere in the green area where I'm doing better than
18 chance?

19 The other little point to point out there, we're
20 talking about probability, so if probability of 1 means
21 that absolutely always will take place with certainty --
22 that means for all time, for all data -- there are very few

1 things that are that certain. Probability 0 means it will
2 never take place for all times, all data; and that's also
3 very difficult in any complex situation to ever do. So, in
4 reality, the best you can do is approach those two
5 extremes. You will never be absolutely sure that you're
6 detecting everything and having absolutely no false alarms.
7 So, just think about that.

8 The other little way to look at data mining --
9 and I think one way that that term came about -- it's like
10 refining ore. If I take raw ore and I'm trying to refine
11 it, each stage I go through the ore that I come out with
12 has more of what I'm looking for, more gold or what have
13 you. But that means I'm leaving little flecks of gold
14 behind as I refine it, because again, you can never be
15 absolutely certain.

16 So having said that, the order is next that David
17 will be giving you an -- we're intending this to sort of be
18 a tutorial of what the technology can and cannot do. We
19 will want to basically -- if we're going to have a debate
20 about the policy, we as technologists would at least like
21 to present you with what that technology can do, and
22 present you with things like this that if we are allowed to

1 test systems, there is a way of testing them, but we need
2 to have real data to do the test for it to be meaningful.
3 Okay.

4 Mr. Jensen: Thank you, John. So, I've already
5 been at this podium for a long time and you're probably
6 tired of hearing from me, so I thought I'd just give a
7 fairly short presentation.

8 And one of the basic ideas that I want to make
9 sure that we get across is that we're trying to compare
10 performance of different ways of doing a task. In many
11 cases, I think in most cases, when we're dealing with
12 national security issues, with domestic security, Homeland
13 Security, there really isn't a question of if you will
14 attempt to do some task; the question is, what is the
15 approach you're going to take to it? And so I think it's
16 important to compare alternatives and to say, for instance,
17 think about alternative data mining systems, think about
18 alternatives that do not use data mining. And to say, how
19 well do each of those alternatives work? I think in many
20 cases, we imagine that doing nothing is the status quo,
21 when actually there is some existing system. We're just
22 comfortable with it because we've had it for a long time

1 rather than we're comfortable with it because we know its
2 error characteristics and we know them to be good. So this
3 graphic, just from the NASD securities fraud example, we
4 ended up comparing to expert-derived rules which, frankly,
5 NASD had never really thought of as a system before. And
6 we found out that we could, by just analyzing data, come up
7 with as good a set of screening rules as they had, and by
8 combining with that we could produce a better system.

9 But importantly, I think it's important to think
10 broadly when you're saying, 'We're validating results.'
11 Really what we're doing is evaluating a system and
12 evaluating it in a relative sense of how it compares to
13 existing systems and to other prospective systems.

14 I also think, as I've said, that it's very
15 important to perform evaluation in context, in the context
16 of the data that you might be gathering or already have, in
17 the context of the decision-making, not view it in
18 isolation. So validation often needs to take into account
19 this kind of larger context, larger institutional context,
20 larger process context. And then it's not just the
21 technical characteristics of the system that matter; it's
22 where it fits. So a good example of this is this question

1 of screening, that if you are doing initial screening for
2 some disease, having a high false positive rate may not be
3 a problem as long as you are putting that in a context
4 where you follow up that first test with a more accurate
5 test, even though it may be more expensive.

6 Final point which you haven't heard me make
7 before but which is an unusual one -- one I wanted to make
8 sure that we talked about -- is that there's a long history
9 of development of technology of these algorithms, and the
10 history has been benefited by the fact that we've, as a
11 community, have developed algorithms and released them
12 publicly. We don't make an algorithm and say, 'Oh, no, I'm
13 going to keep it secret and just tell you how well it
14 does.' We actually release code, release detailed
15 descriptions of these things in the technical community so
16 that other people can build them themselves and try them
17 out and understand their characteristics.

18 One of the things we found over the course of 20
19 or 30 years of research in this area is that it is very
20 frequent to come up with a new technique and only years
21 later -- sometimes ten or fifteen years later -- find out
22 some places that it breaks down that we didn't understand.

1 And that's only possible because the algorithms are public.

2 One of my nightmare scenarios is that someday I
3 will be called in to some windowless room some place and
4 not asked about my own activities, but said -- asked to
5 repair some data mining algorithm that I'm going to find
6 out was used for a long period and no one in the technical
7 community was really told -- or not many people, at least -
8 - and I'll say, 'Well, don't you know, we know that systems
9 like this fail in some horrible way, but because the
10 algorithm wasn't out there, we didn't have the ability to
11 raise those issues, talk about them and identify it.'

12 So there's a real benefit to using what I term
13 here, public algorithm -- publicly released, described
14 algorithms because that encourages wide scrutiny from the
15 technical community and you can remedy errors quickly.

16 Now, there are examples of this in the non --
17 outside of data mining -- Linux and other open-source
18 software operating systems are widely thought to be secure,
19 partially because -- or more secure than they might be
20 otherwise -- partially because errors can be identified
21 quickly and easily, and fixed.

22 The internet protocol, the basic protocol

1 underlying the internet is a public protocol, and errors
2 and problems with it have been fixed over the years and
3 it's been improved. The public key encryption is a
4 wonderful example of this, as well. A known public
5 algorithm that is used to encrypt data -- and just because
6 it's public doesn't mean it doesn't work and doesn't work
7 for very, you know, important, secure applications.

8 So these are nice examples of public algorithms
9 that we have in other domains; I would argue that we need
10 them in the area of data mining. These algorithms should
11 be public even if the data that they're operating on or
12 their conclusions, their models are not public. The
13 algorithm can be public even though the data and models are
14 not. And that's it.

15 Mr. Coggeshall: Thanks, David. My name is Steve
16 Coggeshall, and I work at a company called ID Analytics.
17 We're essentially an identity intelligence provider, and we
18 do analytics around very large identity networks --
19 connectivity of individuals, primarily for identity risk
20 and for authentication -- remote authentication, data
21 breach analysis, things like that. My background: I'm a
22 scientist; I came from academia; I worked for ten years in

1 a national lab doing fusion research; and the last
2 15 years working in industry doing research. I've
3 spent the last 20 of my years building data mining
4 models in many industries, both in the public and private
5 sector, for governments and for business in many applied
6 contexts.

7 We're going to talk very briefly about -- I'm
8 going to give you my quick tutorial on what a data mining
9 model is, tell you a little bit about how to build a model,
10 and then a subject that's pertinent in this aspect -- in
11 Homeland Security is, what do you do if you don't have
12 known bads? If we don't have a lot of examples of known
13 terrorists, how do we build and evaluate a model? And then
14 next is how to evaluate a model when you don't have those.
15 And then, finally, I will just talk a little bit about what
16 are the benefits of using models.

17 So first of all, what is a data mining model? In
18 its simplest form, a model is an algorithm; it's a
19 functional formula that takes inputs and provides an
20 output. And that's what that little box in the center here
21 is; it's just this mathematical functional formula with
22 output y and a set of inputs x . The inputs are typically

1 characteristics about a person or event, and we can denote
2 the string of characteristics in some notation, x_1 , x_2 ,
3 x_3 ; and the output y is the likelihood that it's something
4 of interest -- could be the probability that it's fraud,
5 the probability that it's bad credit, the probability
6 that a consumer is going to buy a product, the probability
7 that it's a terrorist. It's something of interest. And
8 then the model itself is just a mathematical formula. And
9 here's a very simple example that's actually used
10 frequently in practice; it's just a linear combination, a
11 weighted combination, a weight a_1 times the characteristic
12 x_1 , plus a weight a_2 , times a characteristic x_2 , and so
13 on. And when you get done with that, the y , if it's scaled
14 properly, is a score and the score can be -- represent a
15 probability. The a 's -- the parameter's a 's are a set of
16 constants that are learned from data, and that's the --
17 what we talk about training a model, is showing the model
18 data and then statistically finding the best set of a 's,
19 the best set of parameters that matches your data and does
20 the best value of -- the best prediction for that set of
21 data.

22 So how do we build a data mining model? Well,

1 first of all, we use lots of data, in general, to build a
2 model. And, in general, the more the better. A data
3 record, you can think of it looking like this, it's just a
4 vector, it's just a string of characteristics, x-1, x-2, x-
5 3, and so on, followed by the outcome -- whether or not
6 he's a terrorist or whether or not this person's bought a
7 product or whether or not he went bad or this credit or
8 whether or not it's a fraud. That record, that string of
9 information is a single data record, and we build and use
10 many millions of data records when we're building models,
11 typically. And again, I want to point out, it's very
12 important to clean the data; if the data is not cleaned and
13 scaled and represented correctly, then you just have, you
14 know, garbage in, garbage out. The model will not train
15 well and will not -- you'll never be able to build a
16 successful model unless you're very careful about how you
17 clean your data.

18 So what does our data look like now? It's just
19 this arrangement of these many, perhaps millions, of data
20 records. And then what we do is we split them into two
21 sets, a training data set and a testing data set. We use
22 the training data, along with some statistical and machine

1 learning algorithms, a whole field of science has evolved
2 in the past 15, 20 years around this -- around very
3 efficient and very well-built algorithms to do this best
4 functional fit, to find those best parameters, a , in this
5 functional relationship; y is a function of the inputs and
6 those parameters. Once you've built your model, now you
7 put it in place and you have to evaluate how well it works.
8 So you do that by using this testing data. The testing
9 data is holdout data that the model has never seen before,
10 and you evaluate how well the model performs on a whole new
11 set of data that it's never seen before. And you can
12 statistically look at how well your predictions match the
13 real outcomes.

14 So that's the usual methodology in practice of
15 building models. And this whole process is called
16 supervised training, because you know the outputs, so
17 you're supervising -- your model is learning in a
18 supervised way. But sometimes you don't know the output,
19 so what do you do there? And I think that's frequently the
20 case in Homeland Security; we don't have a lot of examples
21 of terrorists, for example. So if you don't know the
22 outcome, you don't know who's good or bad, so now our data

1 records look like this. It's just the string of
2 characteristics -- could be their age and their weight and
3 their height or their -- how many times they've flown,
4 whatever -- but you don't know whether or not they're
5 terrorists. You don't have a y. So in this case, you can
6 be successful in building unsupervised models.

7 Unsupervised models approach the problem differently:
8 rather than finding the patterns of the relationships
9 between x and y, it just looks in the x space, and the
10 characteristics space, and it looks for things that are
11 unusual -- anomalies, outliers. So I drew a picture of
12 that here. Let's say, for example, I only have two
13 characteristics to worry about, and in consumer modeling it
14 might be age and income, that's a very frequently used set of
15 characteristics that describe a lot about how people
16 behave. So this might be their age down here, and this
17 might be their income here; and every person has an age and
18 an income; they have those two numbers. So every person is
19 a point in this space. And so we put all your points in
20 there and you see how the data naturally groups. This is
21 what David was talking about earlier - clustering; this is an
22 example of clustering analysis. I see how my data

1 there are flags that go off sometimes, maybe it's Social
2 Security Number matching, maybe it's past record matching,
3 but there are rules that fire -- that cause certain people
4 to be looked at more closely. So those rules in the -- so
5 there is some kind of an existing process today; we can
6 call that a control process. Events go in, and you go into
7 some set of rules, and most of the people come out as not
8 interesting, and that's good. But there will be a small
9 subset of people that are flagged as maybe bad, and those
10 go into almost always a human investigation process. Some
11 physical human has to look at this and make a decision
12 about whether or not this a true bad or not. So -- and
13 again, usually as a result of the investigation they are
14 okay, and those are the false positives; those are the ones
15 that were flagged by the model -- by the rules -- but
16 turned out to be okay. And then these are the true bads
17 here, the ones that turn out to be really bads. So in any
18 process, you can start -- you can instruct some important
19 metrics that measure the efficacy of the process. And here
20 are two that I wrote down here; these are two that are very
21 commonly used. The false positive rate, we've heard a lot
22 about that; that's the ratio of how many false -- it's how

1 many false positives you have. It's how many people you
2 bothered that you shouldn't have, so it's the number of
3 false positives compared to some baseline. And there are
4 different ways of doing the baseline; one way is by
5 dividing it by the number of true bads; another popular way
6 is by dividing it by all the number you investigated. It
7 doesn't really matter which one you use, it just needs some
8 metric that measures your false positives.

9 And then another metric is your bad rate, your
10 bad detection rate. How many real-bads do I find,
11 divided by how many I have to investigate to find those
12 bads. So those are two very objective metrics that you can
13 use to measure whatever existing processes you have. Once
14 you set that up and put your metrics in place and see how
15 well you're doing today, the next thing you do is you put
16 your test process in place. In this case what you do is
17 you send some of your data through a model; you build your
18 unsupervised model or whatever, for whatever methodology.
19 You've got a model -- a candidate model -- and you want to
20 evaluate, how does that work? In particular, how well does
21 that work compare to what I'm doing today? So again, you
22 put some of your records through here, and the model will

1 records go through that. You might find out of 100, I find
2 60 false positives and 40 bads, so your false positive rate
3 might be better; it's lower. Whereas
4 your bad detection rate is 40 percent; it's twice as high.
5 A bad detection rate twice as high means that I can find
6 twice as many bads with the same amount of work, or I can
7 find the same number of bads with only half the effort or
8 half the intrusion. And that's the -- one of the keys --
9 one of the key uses of using a data mining model, is that I
10 can reduce the effort and reduce the intrusion.

11 So I do think data mining in this case is finding
12 a needle in the haystack. I think it's more similar to the
13 -- I think a real problem here is that the needle looks a
14 lot like a piece of hay. And I think that's our real
15 fundamental problem here. And it's a highly non-trivial
16 problem; this is a hard problem. It involves lots of data,
17 lots of -- it will involve a lot of clever data and coding,
18 a lot of understanding of the domain, and domain-expert
19 knowledge. But I do think it's a problem that can be
20 improved a lot with data mining.

21 And another point is, humans should never be
22 taken out of the loop in this. The point of data mining in

1 them alongside of existing processes in a double-blind
2 test, measured the efficacy, and it turns out, in those and
3 in every other case I've done, the data mining models
4 provide a lot of benefit.

5 So, in summary, data mining models work; they're
6 in wide use. I think this discussion is less about whether
7 or not a data mining model can be effective, but it should
8 be more about how would one do it and how would one protect
9 privacy? And there's a lot of questions around
10 that. But I have high confidence that data mining
11 processes would help and would be better than existing
12 expert-driven processes today.

13 You can build supervised or unsupervised models,
14 and there are ways of testing each one in either case.
15 Another important point is, data mining models can discover
16 patterns that your experts have never even thought of
17 looking for. And that's very frequent; what happens in
18 these data mining model processes, is finding relationships
19 that you never even thought of looking for. And again, the
20 point of this is to minimize the review population, either
21 therefore allowing you to reduce your effort in your
22 investigations, but probably more importantly,

1 reducing intrusiveness, really focusing your investigations
2 where they will have the greatest benefit. That's it.

3 Mr. Dennis: Okay. Hi, I'm Steve Dennis, I'm
4 from the Homeland Security Advanced Research Projects
5 Agency, and we typically are trying to work on new
6 revolutionary ideas, and so, certainly, an area of
7 consideration would be privacy protection technology since
8 I don't think we've solved that problem at all, hence we're
9 having this meeting.

10 I remember back about fifteen years ago when I
11 first heard the term data mining and it was told to me by a
12 group of mathematicians: 'Well, you label your program data
13 mining if you want to get funding.' And today, I think the
14 opposite is true. If you'd like to be de-funded, you might
15 label your program data mining.

16 So, you know, what are the essential elements
17 that are required if we're going to validate data mining
18 models? Certainly, this panel is all about the scientific
19 investigation of what works. What the Science and
20 Technology Directorate of DHS is about is discovering what
21 can work. It's not necessarily about deploying that
22 technology immediately, but understanding what our options

1 are. If faced with a particular situation and we need to
2 do more data analysis, what are the techniques that are on
3 the shelf and immediately available for use? We need to
4 have a cadre of those. And hence, Jennifer Schiller's
5 remarks earlier about, you know, S&T would like to in the
6 future, start moving back into this domain of data analysis.

7 As we do this work, these are some important
8 performance considerations. And the first, I think, has
9 been discussed already, which is, you know, can the data be
10 prepared? Is the right data available in the appropriate
11 form? And there's a lot of work that goes on to understand
12 what's happening in a data set. And if you think you just
13 plug data sets in to data mining algorithms and they
14 magically start producing something, they don't. It takes
15 a lot of considered preparation in order to make those
16 algorithms start to produce and produce in a productive
17 way. So what are some of the considerations there? We're
18 worried about the speed of a process; can it actually keep
19 up with a data rate or with a large volume of data if it's
20 needed? What is the accuracy of that mechanism? Is there
21 error being introduced as you start to process data? We're
22 worried about storage overhead; if the original data took a

1 terabyte and the process data takes two more terabytes,
2 this is an issue. We're worried about portability; if I
3 get a solution to one problem in one domain for \$100
4 million, does it take me another \$100 million to solve a
5 problem in another domain? And certainly that goes to
6 scalability, but cost is also a factor here throughout each
7 one of these steps that I mention.

8 If you have structured data, those are the simple
9 cases. If you have unstructured data, there's a lot more
10 work that has to be done. And so there's been a history
11 over the last 20 years of trying to automate this process
12 of information extraction. How do we make data available
13 in the right forms and what kinds of errors are introduced
14 during that process? And as you get this drift, error
15 propagation can happen throughout the system so you have to
16 worry about the performance of each one of these
17 components, not just an individual piece. So, if I'm
18 extracting information, I'd be worried about the linguistic
19 features; you know, how good am I at getting verbs that
20 might imply events? How good am I at getting proper nouns
21 that imply people, places, and other things? These are all
22 factors, right? And you're starting to get a sense for the

1 complexity of this kind of research; it's not very simple.
2 There's a manpower factor; if I create a process that's
3 heavily knowledge-engineering oriented, would a customer
4 ever have a knowledge-engineering branch that can make that
5 work? You know, that might be a consideration at the early
6 stages that make you say, 'This model's never going to work
7 because it's just too manpower intensive.' And I might
8 trade off the efficiencies that I gain at the back-end
9 having to replace those people that I save at the back-end
10 with people in the front-end. So I might gain nothing by
11 doing that. And certainly, accuracy and speed again.

12 Once you have a well-prepared data set and you
13 fully understand the error characteristics of that data,
14 then you can move into a pattern-matching function. And
15 whether that happens to be learning patterns or not, the
16 results of those kinds of algorithms are generally either
17 binary, you know, where they either tell you they found it
18 or they didn't. Sometimes you get rank-ordered lists that
19 say, 'Okay, here are the top 20 choices that match that
20 pattern;' or you could even wind up with weights and tables,
21 and it's even more and more complex to understand the
22 performance of such systems.

1 Then, I move on to an area that's called policy
2 filtering. You know, once I'm able to understand patterns
3 and I'm able to have data sets that can be processed, I
4 worry about, can I automate the application of policies
5 over the top of the use of that data? So, if I have a
6 privacy policy, can that actually be codified and made part
7 of the system? Would the policy folks be in a position to
8 write their policies, not in English, but in some sort of
9 coded form? And then I worry about things like the
10 receiver operator curve performance of that; you know, we
11 saw the graphs before, whether it's precision and recall or
12 false alarm and missed detection; there are many ways to
13 talk about it. And then I worry about leakage. You know,
14 how much of this information that I have is leaking over
15 the boundary at any one point in time. And I can do an in-
16 depth analysis of that and start to look at tradeoffs, and
17 even start to look at what it might mean to compare human
18 performance in that case, you know, if a human's making a
19 decision to a machine performance, and that's really
20 important. If you can get inner-annotator agreement, and
21 what that means is, if humans can agree on a task, then it
22 probably can be automated. If you have a group of humans

1 who can't agree on a task, then perhaps it can never be
2 automated and you should save your money. So you worry
3 about data retention, audit, traceability, and the
4 policies, the overall effectiveness. And you start to see
5 that you can trace now the use of a policy and how it's
6 impacting systems and performance.

7 Above that level -- and these are all
8 architectural issues -- is a system-level concern. And that
9 is, is this system usable? We heard earlier today that,
10 you know, if you generate a number of leads and they all
11 lead to a lot of overtime and there's no productive result,
12 then that's not a good thing. And you need to have mission
13 metrics that tell you that a system can actually perform in
14 an efficient manner. Efficiency also goes to cost. If I
15 deploy a very large data mining system, and it takes \$100
16 million a year to keep it going, perhaps I've gained
17 nothing.

18 And we talked about traceability; that's all
19 throughout the system, not just in the audit of the policy,
20 but, you know, who touched data when and where is a very
21 important factor of each one of these systems. And then
22 information assurance. Have you authenticated the users of

1 the system? Are they in the proper role? And all these
2 factors have to be considered. So now you can understand
3 why the performance evaluation, the test and engineering
4 might be very expensive. And, as a matter of fact, getting
5 involved in some of these efforts, sometimes the data
6 collection, the evaluation, and the preparation for the
7 research can outweigh the budget that you have for the
8 research. So, you know, we're trying to make this as easy
9 as possible, but there are many layers to consider.

10 One idea that may help us all -- and I think it
11 was alluded to earlier -- is the development of some sort
12 of common research and development framework that allows us
13 to reuse components. If I'm really good at preparing data,
14 do I really have to engineer an entire system around my
15 effort in order to understand the system effects? If there
16 were such a common framework, you know, that was freely
17 available and software could be traded around the community
18 with the kinds of visibility that we heard about earlier,
19 perhaps that would help us make progress. And it would
20 also save us a lot of time at DHS. We are approached by
21 many vendors with many ideas and many universities with
22 many ideas, and often they come in selling us a brain in a

1 box, but no one has done any sort of evaluation that tells
2 us what the true performance of that system is. You can
3 spend two hours unwrapping a package to find out that it's
4 the same as a system from 1960. So you have to be very
5 careful there. And this kind of framework may help us
6 understand better performance.

7 We talked about data collection, and it's very,
8 very important for research. What you'd like to have is a
9 sustainable data set that represents a hard problem that
10 can last for 20 years. And dare I say, I don't think
11 there's anybody in this room who will approve of a data set
12 containing private information that could be retained by
13 the research community for 20 years in order to do
14 repeatable experiments. So that sort of works against our
15 normal R&D methodology at that point.

16 So, folks will often say, 'Why don't you use
17 synthetic data sets?' Well, you can use synthetic data
18 sets early in the process to help debug, perhaps to
19 understand the implications for scalability, but if you
20 have a language problem and perhaps a name-matching problem
21 -- for some of the lists that we heard about before -- you
22 can't really work that name-matching problem with unreal or

1 made-up names. You really need a set of real names and
2 real situations so that you can model the actual condition
3 and understand how to improve the performance of such
4 algorithms.

5 If you do have real data sets, you can do risk
6 mitigation for those data sets, although it becomes highly
7 complex, if a series of questions comes up, when you start
8 to consider those. Thank God for the Enron data set and
9 things like that that just happen to be out there that we
10 can use among this node, but, you know, those data sets
11 don't necessarily represent real problems either. So we
12 have issues there.

13 I think about the problem of doing data mining
14 and doing this kind of research, you know, there's a
15 chicken and an egg, and if you're approaching it from the
16 egg, you know, often the questions are, you know, 'What
17 sort of chicken is this going to be?' and, 'What color are
18 his feathers?' and, you know, 'Is it going to have a mole?'
19 You know, I don't know, you know, until the egg hatches.
20 If you come at it from the chicken end, you know, basically
21 they want to know what color is the egg going to be, and is
22 it going to be speckled or brown and will it contain double

1 yolks? We don't know.

2 So at the end of the day, what happens typically
3 -- at least now, before we move through some more policy
4 changes -- is that both wind up fried, and, you know, never
5 -- we're never allowed to find out what happens.

6 If we had a common infrastructure that would
7 enable system-level investigations, it would allow us to do
8 more -- get more return on the investment for our research.
9 And we talked about making that code freely available and
10 Reusable, and it might also lead then to common evaluation
11 methodology. And I don't think we have a really good
12 evaluation methodology that centers around privacy. And I
13 think those kind of tradeoff studies would be very
14 interesting, if we were allowed to do them.

15 I wanted to leave you with this thought, and it's
16 a very, very simple cartoon that talks about the points at
17 which we would feel friction doing this kind of research.
18 If you're talking about doing just a normal data mining
19 system with fixed data sets, you know, that's in the top
20 left-hand corner there with the pattern-based algorithms
21 running over some set of data that everybody agreed was
22 okay. But if you look at the intradepartmental situation -

1 - as was mentioned before -- you know, DHS is a collection
2 of a lot of operational elements and each of them have
3 their own rules and their own lawyers, and so it's not so
4 easy to just put things together and make them happen.
5 There are lots of discussions around that. So
6 clearly, there's a policy filtering need at that edge of
7 the graph, you know, as we go across the department. Each
8 component, each data set has its rules and charters and
9 implications.

10 If we look at cross-departmental access, we're
11 hoping that, you know, it might be possible to somehow
12 design a common analytic space that would allow the
13 government to make use of what it knows in the right
14 circumstances. But that would imply a lot more policy
15 filtering, a lot more comfort with implementing our systems
16 in different ways. And so, just a thought for you to think
17 about as we continue the panel.

18 Mr. Hoyt: Okay. To raise some questions for the
19 panel, since we have -- and we have time for questions from
20 the audience as well, but one of the areas that I've seen
21 in the literature is this perception that somehow data
22 mining can take all types of data and magically combine it

1 and come out with useful results. And I have my own biases
2 about that, but I'd like to open that up for our panel of
3 experts here to at least comment on that perception that I
4 can throw everything in there and somehow it'll make, you
5 know, a gourmet meal out of all the hash I've thrown in.

6 Mr. Jensen: So why don't I start. One of the
7 things that I often try to tell people, the things I try to
8 tell people

9 is that, if you compare data mining to aircraft design, that
10 we're just out of the Wright brothers' stage. We tend to
11 think of this as a high-performance technology, and
12 certainly, many technologists want to say, hey, we've got
13 these really wonderful algorithms. But the truth of the
14 matter, I think, is that we're actually very early in our
15 understanding of this technology and development of new
16 technologies. And one of the consequences of that is that
17 there are many types of data that we don't know how to deal
18 with effectively, at least nowhere near as effectively as
19 somebody could whose an expert, who is a human who can just
20 look at it and interpret data.

21 One good example of that is that, until about
22 ten years ago, we didn't actually have methods that could

1 look at interconnected records and make use of those
2 interconnections. So we know if you look at -- if you
3 think about how a doctor does medical diagnosis, you come
4 in with a fever or something into an emergency room, and
5 the doctor starts thinking, 'Okay, maybe it's
6 communicable.' So he or she asks you about who you have
7 come into contact with, if your family members have this
8 disease, et cetera. They also know that, 'Well, maybe it's
9 genetic -- maybe they are genetic components to this
10 somehow.' So, 'Gee, has your mother or father or children
11 ever suffered from this?' It also might be occupational,
12 you know; maybe it has to do with where you work. So we
13 think about, naturally, all of these relations, but we
14 didn't have methods that could think like that, that could
15 develop models that looked at those kinds of relations
16 until quite recently.

17 And another big area is people say, we have tons
18 and tons of data, but what they really mean is we have tons
19 and tons of textual reports that each of us could sit down
20 and read and extract information from, but as several
21 people pointed out, we don't actually have good automated
22 methods that can reliably look at lots of unstructured text
23 and pull out the sort of meaning that is anywhere close to

1 the meaning that a person can.

2 So we are -- I think we are fairly limited to
3 numeric and symbolic data that is connected up in
4 interconnected records, and a very limited kind of
5 extraction from large text documents.

6 Mr. Coggeshall: I'd just like to add to that I
7 think one of the biggest challenges in data mining these
8 days is unstructured data. We have a huge proliferation -
9 explosion of data, primarily in the unstructured space --
10 text, audio, voice, image, video -- and it's getting more
11 and more important for all these different needs to be able
12 to use that kind of information. And it's an extremely
13 complex problem, a lot of research going on -- a lot of
14 successful research -- all those categories I mentioned are
15 being used today in data mining, but we're -- I would say
16 we're just at the infancy of building -- of learning how to
17 efficiently encode that unstructured data into numerical
18 representation for algorithms to operate on. So it's still
19 -- we're still at the very beginning stages of that. And
20 it's a hard problem. And I think for Homeland Security
21 it's going to be one of the key areas; it's going to be

1 very difficult.

2 Mr. Dennis: I too think that we are at early
3 days on data mining and there has been a lot -- even though
4 there has been a lot of investment, we don't fully
5 understand how to approach this problem-solving in a pro-
6 forma way. If you think about throwing together lots and
7 lots of data, it just reminds me of the Wal-Mart example of
8 putting the beer next to the diapers, you know, sort of
9 where the young father comes in to get the diapers because
10 the wife said, 'Go get diapers.' And the correlation
11 happened that, you know, the young father also picks up a
12 six-pack of beer to go with those diapers on the way home.
13 Those kinds of associations are found because you start to
14 look at these patterns in large data sets and you get
15 discoveries that you wouldn't otherwise get if you didn't
16 do that.

17 But very early in the process, it's sort of a
18 triage stage where you're looking for those sort of things
19 to happen and to make those discoveries. But very, very
20 soon after you spend time in that discovery phase, you
21 start to worry about what is the contribution of data to my
22 observation or to my inference? And you try to trim the

1 data and get rid of all the things that don't matter, so I
2 think there is a value to doing some of that kind of
3 investigation in a triage mode.

4 Another thing that I don't think has been
5 mentioned yet is that even though you discover these
6 patterns and they work today, you know, they may not work
7 next month. And so there is certainly a lack of
8 understanding of model lifetime, and models drift, and so,
9 you know, how often you have to reinvestigate these kinds
10 of relationships is probably unknown for a lot of data
11 sets.

12 Mr. Jensen: Another comment, if I could. When I
13 worked for Congress at the Office of Technology Assessment,
14 we did a study which ended up saying, don't collect more
15 data because that's not necessary and it's not going to
16 help. And I think it's an interesting example of where
17 analysis -- not data analysis -- but careful analysis of
18 the overall task can actually tell you, you don't want to
19 do this. We were asked by Congress to look at the question
20 of whether additional data on wire transfers in the United
21 States -- large money transfers -- would assist in the
22 detection of money laundering -- criminal money laundering

1 by large organized crime groups, particularly. And 18
2 months of work, of really talking to large numbers of
3 experts, and really understanding the analytical tools that
4 are available at the time, and the conclusion after this
5 intensive 18-month study was, no, don't collect data on
6 wire transfers because, one, the amount of information it
7 contains is so weak that it's very unlikely to yield
8 anything; and secondly, it was going to increase an order
9 of magnitude by ten times the amount of data that a small
10 treasury agency had to sift through, and that would have
11 ended up actually swamping them. And they said, 'Please
12 don't give it to us.' But also, we had good, you know,
13 quantitative reasons to say that. So in the end, we came
14 back to Congress and said, 'No, not a good idea to collect
15 additional data.' And I think that's the kind of technical
16 conclusion you can often come to, is that, don't add more
17 hay to the haystack if you're looking for a needle.

18 Mr. Hoyt: Another topic that we get approached
19 with is, given the difficulty of dealing with personal
20 data, can we deal with synthetic data? And at least I'll
21 give my bias and let the panelists kick in. But my bias is
22 that synthetic data is useful at almost the -- I'll call

1 practical problems by using synthetic data.

2 Mr. Jensen: I guess I would ask whether you
3 would like your -- the pharmaceuticals that you might take
4 in ten years to be tested on simulated humans or on real
5 humans. Yes, there are certain kinds of things you
6 might be able to figure out by looking at simulations of
7 the human metabolism, but it's not what you want for
8 anything except the very early part of the process.

9 Mr. Hoyt: There's another class which actually
10 maybe falls out of the range that this forum is really
11 interested in, but DHS does have problem sets where we care
12 about patterns that have nothing to do with personally
13 identifiable information. We care about pandemics, both
14 for people and for animals and food crops. Obviously,
15 there are other parts of the government that we partner
16 with in that, but in that case we need no personal
17 information. In fact, that's noise as far as we're
18 concerned.

19 The other -- several of our panelists have worked
20 on systems for industry and for other agencies. There is
21 at least some perception that I'm getting out there that
22 people think that we're just sort of starting off from

1 ground zero, and I think it's been touched on several times
2 in this conference, that there are existing processes that
3 are in place. And I'm assuming if I'm industry funding
4 something, I have a profit motive and I want a system that
5 does it better than the existing system.

6 Mr. Coggeshall: In my experience, I've built a
7 lot of data mining algorithms for a couple decades now in
8 lots of different industries, installed them all over the
9 world -- and I've done some really tough problems, some
10 pretty basic ones, but some pretty tough problems, too --
11 and I have never done one where I haven't been able to
12 outperform an existing system. It doesn't mean it's
13 tremendous, but we've been able to beat whatever the
14 existing process is. I think that's just -- and I -- I
15 have high confidence that that could be done here in this
16 case also.

17 Mr. Hoyt: Having said that, could I open it up
18 for questions from the audience? And if you'd please come
19 up to the microphone.

20 Ms. Schiller: Hi, Jennifer Schiller again,
21 Science and Technology Directorate. And Steve Coggeshall
22 said in your presentation that when you are building a data

1 mining model, the more data you have the better. But in
2 the government, we have Fair Information Practice
3 Principles that require us to use the minimum possible
4 amount of data, so there seems to be a real tension between
5 our legal and privacy policies and the technical
6 requirements of building a data mining model that will
7 work. So I was wondering if you could speak to that for
8 the whole panel a little bit more.

9 Mr. Coggeshall: Sure. This is a common problem
10 and we face it in industry all the time, too. And a very
11 good example is in credit scoring. There are certain
12 fields that we know are useful that we cannot use for
13 regulatory reasons in credit scoring. So the reason we
14 know they're useful is because we tried them and they work,
15 but then you go through an iteration process with legal
16 systems and policy, and for a variety of reasons you're
17 forced to remove those pieces of information. And so,
18 philosophically, the more data you have, the more varied
19 and disparate data you have, the more -- the better your
20 models will perform. I mean, obviously, if I have a
21 certain universe of data and my model works to a certain
22 level, if I add more data I never get worse. If you do

1 things right, you only get better. So, at some point, when
2 you add more data, you don't get any better, and that's what
3 scientists are all about -- applied model builders are
4 looking for where that tradeoff is between adding more
5 data and not -- and the amount of performance you gain is
6 not -- is disproportionate to the amount of effort it
7 takes. So, in this case, I think it would -- it might make
8 sense to have an environment where you can, on a trial
9 basis, try lots of data. But, in the end, you will find
10 that a subset of that is what's needed, and then from there
11 on that's what you need certainly for implementation of the
12 model when the data that's streaming into the model only
13 looks at that subset of the data.

14 Mr. Dennis: I think it's important to have
15 access to a lot of data in order to figure out what the
16 minimum set is. It's not possible to discover that minimum
17 set without some experiments, and so if we have to guess
18 what the minimum set is at the outset, we're likely to
19 spend a lot of money as we continually enlarge that circle
20 until something starts working.

21 Mr. Jensen: I think one of the things we need --
22 and I think this was referred to by one of the previous

1 expert way that presents very efficiently information to an
2 algorithm? And that's how one frequently finds patterns
3 that you've never discovered before.

4 So it's an iterative process: you do the best you
5 can in your first step; almost always you'll beat the
6 existing process, and then you continue to get better from
7 that and inventing new variables, getting feedback from the
8 experts, trying new auxiliary data sets, and it's just a
9 continuous improvement process.

10 Ms. Szarfman: The more data you analyze, the
11 better you understand the data. And the technology's
12 evolving very quickly, so what we cannot do today we will
13 be able to do tomorrow. Then we are limiting the amount of
14 data you'd analyze, which will restrict you in ways that you
15 cannot foresee. The more you analyze, the better you can
16 find the outliers; you will understand if your methodologies
17 and appropriate method can find things in
18 higher dimensions, like, in my case -- in our case at the
19 FDA, is to find drug interactions in specific (inaudible)
20 population that may be at risk because they are elderly or
21 they are preemies or restricting. But you have a different
22 problem, you are trying to find unexpected things. It's a

1 very difficult problem because if you knew where to find it,
2 you would not be having this meeting. Then we are -- which
3 data you should analyze is also, you know, you don't know
4 which data you should analyze. Then you have a difficult
5 task.

6 Mr. Coggeshall: That's all completely correct.
7 This is not an easy problem; this is a hard problem. My
8 point of view is algorithms can be built that will improve
9 upon the existing processes, but by no means is this an
10 easy problem

11 Ms. Szarfman: No.

12 Mr. Coggeshall: It's very hard.

13 Ms. Szarfman: It's very, very hard.

14 Mr. Coggeshall: Yeah.

15 Ms. Szarfman: Our problem is also very hard
16 because it was considered impossible in the past that, you
17 know, you could not get anything out of secondary data that
18 you don't have in hypotheses. You need to have an
19 hypothesis, then you set up the data to analyze the
20 problem, but the primary reason I don't understand the data
21 because you have this -- a way of understanding the data.
22 Then you are looking for interactions with alcohol, but if

1 Ms. Gregory: Hi. I'm Michelle Gregory from the
2 Pacific Northwest National Lab. And I'm a researcher in
3 data mining, and this whole session is about how you can
4 validate data mining models and results, but I've been
5 trying to rethink the problem into -- instead of, as you
6 mentioned, overlaying policy on top of the data that you
7 have for the analyses, can you include the policy at the
8 data-collection phase? In other words, how much can you
9 glean from data and patterns can you find that are useful
10 without having all the data available? So in it's most
11 simplistic form -- from the talk this morning -- it would
12 be, anonymize all the names and maybe places and locations;
13 you find the connections that are interesting, then you
14 reveal them under certain policy conditions. So I just
15 wanted to hear your comments on that.

16 Mr. Dennis: I think you can certainly distribute
17 the policy control throughout the entire process. And so
18 as you're ingesting data, certainly as you're doing
19 information extraction, there's a lot of anonymization that
20 can be done and I know of programs that do that now for
21 HIPAA applications. It's absolutely true that you can
22 distribute it at any stage in the game, but it's important

1 unfortunate. But there's another example, which is that we
2 do medical research right now on new drugs, on new
3 treatments, and there's a very small number of people --
4 maybe 100 or 200, you know, who will volunteer for some
5 clinical trial, and they'll participate in it and they get
6 compensated or they get free treatment or something else.
7 And then, as a society, we benefit from that. We say that's
8 a great thing; they did that and now we know whether this
9 drug works. And in the same kind of way, you don't need to
10 implement an enormous data-collection procedure in order to
11 find out if there's some signal there, some statistical
12 regularity that you can catch. You could go to extremely
13 focused, careful -- maybe random samples, but very small
14 amounts of data, find out whether it -- there's something
15 there to model, and then you could say we've got some
16 reason to believe it would be good to do the more general
17 data collection. That's a different kind of idea of
18 clinical trials, but I think it's essential and something
19 that's really outside of the kind of research or complete
20 program that we do right now. We need to have some trial
21 runs at constructing systems.

22 Mr. Coggeshall: I'd just like to make a comment

1 on anonymization. There's some tremendous work going on
2 that field -- in both industry, IBM and academia. One
3 needs to always be careful about this. For example, some
4 of the work we do in our company, requires us to do fuzzy
5 matching across multi-dimensional spaces with name,
6 address, Social Security Number, phone, date of birth,
7 things like that that are -- that if you anonymize first
8 before you try to do fuzzy matching -- multi-dimensional
9 fuzzy matching, not just one at a time, but all
10 simultaneously -- it can destroy some of the connectivity,
11 so you need to be very cognizant about where and when you
12 do your anonymization.

13 Ms. Schiller: That was actually my question - is
14 what are the implications of using anonymized data to
15 validate a model, as opposed to -- early in this day you
16 talked about synthetic data could be useful, but taking
17 real data, anonymizing it and then using it in the
18 validation process -- how would that impact your ability to
19 -- from an S&T perspective to say to a customer, 'This
20 works; I'm confident it works?'

21 Mr. Dennis: I think anonymization poses several
22 challenges for the kinds of research that we want to do,

1 and that is if we were looking at the name match
2 application and you anonymize away the name, then that's
3 putting us out of business. If you -- I would ask, though,
4 if you consider lack of human access to the underlying
5 content as an anonymization. You know, if the algorithm
6 gives a chance to see the personal data and yet the human
7 doesn't, is that an effective anonymization? I often
8 wonder about this because machines doesn't have malintent
9 and machines do what they're told, you know, they're not
10 spying on their neighbors. But people, you know, are often
11 accused of that, especially if they work in government. So
12 it seems like to me one way to mask and anonymize the data
13 is to allow for the machine to have access but not for the
14 human.

15 Mr. Coggeshall: I just have to say, that's an
16 excellent point and that's something we do. All our data
17 in our company is encrypted, so your socials, names, and
18 addresses are encrypted, but they are unencrypted in the
19 algorithm phase itself for the matching. But when we look
20 at the data -- unless we're doing case studies or something
21 like that where we have to explicitly unencrypt the fields
22 -- we interact with it in a completely encrypted format.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

CERTIFICATE OF TRANSCRIBER

I, Renee Braun, hereby certify that I am the transcriber who transcribed the audio recording provided by Alderson Reporting Company to the best of my ability and reduced to typewriting the indicated portions of provided audio in this matter.

Transcriber