

---

MONTREAL – Tech Day - PM Session  
Monday, November 4, 2019 – 13:30 to 18:30 EDT  
ICANN66 | Montréal, Canada

EBERHARD LISSE:

Okay, if we can sit down and settle down. I don't have to name any names then. Okay, so we are going to start the afternoon off with RDAP deployment. Interesting is that when the two presenters from WHOIS and RDAP just spoke a little bit, we found that there are a little bit of synergies or disaster recovery-synergies or something, so it's probably interesting to hear both presentations next to each other, and we'll keep both presenters on the table so if there is questions, I would like them to be done jointly at the end of the second presentation.

Good afternoon. Thank you very much. I'm Marc Blanchet from Viagenie. We've been doing a lot of RDAP work for some years, and this is kind of an informal deployment status of RDAP, informal being that it's a [hack] that I did writing some code, kind of scanned RDAP Internet to see where we are. So informal, take the numbers, the statistics with a grain of salt. It's not an official status.

Plan of the presentation, just an introduction. I'm assuming that people know what RDAP is. In a couple of sentences, it's the replacement of WHOIS based on modern technologies, so HTTP queries and responses formatted in json. It's defined in a few RFCs, and it's ICANN as by end of August requested the contracted parties to actually implement RDAP, which gave kind of a boost to implement it.

---

***Note: The following is the output resulting from transcribing an audio file into a word/text document. Although the transcription is largely accurate, in some cases may be incomplete or inaccurate due to inaudible passages and grammatical corrections. It is posted as an aid to the original audio file, but should not be treated as an authoritative record.***

---

The other part I could say also about RDAP is that it's also implemented by the original IP RIR, the regional registries for IP and AS numbers, and it's the same protocol. So it's actually used in all kinds of solutions. So that's the RDAP tutorial.

So we did the various tools for RDAP, [inaudible] tools for servers. The client, we implemented the RDAP client for ICANN. We wrote the IANA RDAP server which should be online anytime soon, and we wrote a mobile client called RDAPBrowser that you can get in both iOS and Android.

While doing this, we found all kinds of various issues while browsing the RDAP Internet, and we wanted to get a view of the whole deployment, and this is the presentation.

So first thing is statistics. The methodology we did for the statistics is as follows. We fetched root zone, the IANA RDAP domain bootstrap registry, the IANA registrar ID registry, which contains the registrar RDAP server URL.

For each TLD with a known RDAP server, we tried to query a domain query for nic.tld. Most of the time, works. If not, for whatever reason, then it's not in the statistics you will see later.

For each registrar with RDAP server, we try the RDAP query of the domain of that RDAP server URL, which means that it should be served by the registrar. Most of the time, it works. If not, then it's not in the stats. If we receive a 404, that means the registrar domain is actually

---

not managed by the registrar itself, which is kind of weird for me, but whatever.

We check if the RDAP response is compliant. We'll define what compliant means here in the next slide or two. We inserted a two-second delay between each query, which overall takes like an hour to do the full scan. The two seconds delay is actually not relevant, but because you go through different registries and registrars, it's one after the other, so there's no point of having a delay but I wanted to make sure that I was nice.

For http errors, we tried a few times, see if we can get different than http errors, and that was run the last two weeks of October.

We received all kinds of errors, http codes such as 429, 504, 503, especially 429 were interesting because I'm sure that those registry and registrar are not currently receiving thousands of requests. So anyway.

Invalid TLS certificates, and badly formed RDAP responses. Obviously, the top one is always the description, and remarks and notices are not a string array, it's just a string so it's not confirmed.

And we even have one registrar or registry, I don't remember which side of it, but one actually created new dialect of RDAP. It looks like RDAP but it's actually not, which was interesting.

The compliance definitions, when we say it's not compliant in this slide deck is it's just based on the RDAP parsing library we have for display purposes. What that means is the only data structure need to be followed, but we don't care if you put garbage in the values for example,

---

or if you have an answer which lacks a lot of data in it. Just what you send should be conformant with the data structure of the response. So in fact, many that are defined complaint in the statistics may actually not be really compliant. But you'll get the stats.

CORS headers are not checked, therefore – and we know we've seen this, that for the RDAP purpose, you need to have CORS headers who are kind of open, and we know that people are not doing this, so web clients would have problems, but we have not checked, this is not in the statistics.

Again, these statistics are not official, it's a hack that I've done, so consider it as it is. CcTLDs, 13 have advertised an RDAP server that responds to the specific query. What that means here is I found 13 in the bootstrap RDAP IANA RDAP registry which those who are ccTLDs are actually responding to nic. the ccTLD. That's the 13, which means there are probably more that do not respond to nic.tld, and maybe even more that are not registered in the bootstrap RDAP registry right now.

On those 13, seven are complaint. Again, remember what I said about compliance. Four have various http issues that even after multiple scans, they were always responding “too many requests,” and two were not complaint. We have more ccTLDs, so given that this is ccNSO here, I would encourage ccTLDs to actually implement and deploy RDAP.

gTLDs, 1200 have advertised an RDAP server responding to specific query, nic.tld. About 800 were compliant. Again, that compliance level is very low. Compliance level zero. We're not asking that much. Four

---

have various http issues, and 362 were not complaint to the basic level we were trying to reach.

Registrars are a different way because registrar ID, IANA registry does contain – so there are, in this registry for example, some registrars, the entity itself, have maybe registered like 25 different IDs for the exact same registrar. So what we did is we defined a unique registrar based on the RDAP server URL they are advertising, so that’s where we kind of define what a unique registrar is.

On the set of 170 that actually responded with the right thing to our query, which is extracting the domain name within the RDAP server URL and then asking the domain query for that domain, then 33 were complaint with 20%, 15% had various http issues and 65% were not compliant. That also corresponds to our experience with the clients, the registrar side of the RDAP universe needs to be in better shape.

I'm going through this just to give you some examples of issues we found. There's a whole [IETF] document that I created that kind of lists some of them. I probably got into 50, probably 100 different issues I found over the last two months. Some I reported to the people I know in the field, some I reported to ICANN saying if you know the guy who is involved with this RDAP server, please contact them and tell them they have a bug. So just an example. It’s not all the issues we found, just some of them as an example.

So these are types, so actually, this is kind of a registered value, and if you look at the left side, these are the values that we found but the people were just kind of lazy to not use the actual registered one in the

---

RFCs and the IANA registry. So they were not using the same words. Like the last one, the table below, the first three were [nav,] I didn't find the registered value correctly, but the last one was last update and last change. Read the documentations.

Server deleted instead of server delete, okay versus active, owner versus registrant. You see all kinds of registered value. By the way, those were not in the statistics, those are complaint. Those kinds of strings like, "Okay, you're fine."

By the way, in the statistics I showed before, those were usually considered as complaint because I was kind of nice to them, the [inaudible] accept everything. But in fact, they are not.

Object name, class name, that's the property in json that tells you what this is, and it's empty. So I cannot tell what this is, a domain, an entity, what is it? That's the basic reference in RDAP json response. When this is empty, this is an error. I cannot parse this.

Relation values, this tells you [links] relation is actually if you want more information about this, go there. But there are a few values that are defined, such as about or related, or terms of service. But some people actually put the URL in the value itself. Whatever. So I cannot do anything with that thing in a client.

A related link which should go to the registrar actually pointing to myself. So if that was a bug in the client because I was looping, because it was referring to itself instead of referring to the registrar.

---

Well, I'm just passing all those. I can spend the last afternoon showing all the issues we found over time. So for gTLDs summary, end of August was the deadline for offering an RDAP service. Far from done, but good progress. Compliance issues should be fixed, registrars are pretty [behind] for compliance, and developers, please read the RFCs. Thank you very much.

EBERHARD LISSE:

I've just got one question before we open the floor. Whoever wants to ask [inaudible] ask the question at the end, but I have just one. I noticed four HTTP were having issues as the Gs and four at the CCs. Were they on the same infrastructure?

MARC BLANCHET:

Yeah, it's not a copy paste problem, it just happened.

MARK SVANCAREK:

Hello. I'm Mark Svancarek, I'm here from Microsoft, and I'm going to talk about how we use public data, public RDS data at Microsoft. We'll get into the methodology a little bit, like Marc did, but like his, this is a little bit of a hack, this is something that was done in our spare time so you can draw whatever conclusions you want from it, but you'll see that it doesn't answer all of your questions.

Just to introduce you to this topic, we have several teams that use WHOIS data. I'm going to just call it WHOIS data. We have the Office 365 domain team, you can bring your own domain name to your Office 365

---

subscription. It turns out people are not always really good at letting us know if they sell those or get a new one or let it expire or something like that, and then that completely screws up their [inaudible]. So we help them out, we go and we look at the records and remind them that there might be a problem or that they might need to renew.

There's a fraud protection platform that we have that generates confidence scores for electronic transactions. One of the datasets that they look at is domain name registration data. There's antipiracy, our corporate domains team, you can imagine we buy and sell domains, trademark and copyright, there are always infringers that we would like to be able to contact, and other digital crimes. So things like malware, phishing, and the like. A lot of these things overlap each other, of course.

So the data has always been stored in a local database for efficient access. I'll explain a little bit more why that makes sense. The data was historically received from DomainTools, and in fact we still have a subscription that's still going in there, it's just that we used to have full records, now we only have the public data, there's no nonpublic data in there. Nevertheless, it's still a useful resource.

So the reason we have a centralized database is that it provides predictable access to normalized, frequently used data, and even infrequently used data. So as Mark has discovered on RDAP, people using the wrong words, that was always an issue with WHOIS and sometimes people would change their formats on the fly. So once it's



---

normalized, it's nice to have it in one place where everybody can access it.

It's always good to have a shared infrastructure where engineering costs can be shared and optimized, so that's a benefit of this database, but a real capacity is the logging and the audit capability. So everyone who has access to the database has their own API key, their own account, and everything can be managed.

And in fact, that's where this data came from, so this data is generated by "I looked at all the logs over a period of days, identified who was looking at it. This is all part of a DPIA exercise as well as an RDAP planning exercise.

So how many requests were there during the day? Who was making them? Then I looked at the data that was in the database and pulled out things like the WHOIS server information, which gives sort of a proxy information for which contracted party we'd be pulling this from if we were pulling it ourselves.

And then just another comment that even though this is only public data, there's lots of uses for it. After the EPDP completes, we will be requesting disclosure of nonpublic data again and it'll also go into this centralized database where it'll be accessed, controlled and audited.

The data that was analyzed for today's talk, as I said, this was opportunistic. I had pulled data for another purpose and had gotten some conclusions from it that I thought were interesting, so the plan

---

was to share this with you. Four consecutive days in a single month, and a different day from another month. Thought that would be interesting.

But it turns out that one of those consecutive days was an outlier, there was an active investigation going on, and I'll tell you more about that. Not really tell you about it because it was an active investigation, but it was an outlier and it was completely screwing up the statistics. So for the most part, I ignored that.

So for that five-day total, we made 21 million API calls. Now, an API call could be a single domain name, it could be multiple domain names, but we called 21 million times and that covered 1215 contracted parties.

On a typical day, we call the database about 4 million times, and that could be a 250,000 names, it could be 600,000 names, it's all generally in the same zone, but that's sort of the range.

And the interesting things were that most names were only called by one party, one person, one app, something like that. I thought there'd be more overlap, that some people would be calling the same names as other people, and that turned out to not really be the case.

Most names are only requested once per day, but some names are requested many times per day, and that's really related to cloud scale fraud protection for the most part, the latter.

So if you can imagine a fraud protection platform, it's looking at all the attributes of an electronic transaction to generate a confidence score, and we do look at the domain names that are involved in that, whether

---

it's a website domain name or whether it's someone's e-mail address, that all goes into the score. So e-mail addresses pop up a lot.

What I notice looking into the database is that the top 200 domains that we access the most are all e-mail addresses. If you go down into like the 300 or 400 top ones, it's mostly e-mail service providers, or domain names that look suspiciously like e-mail service providers, which presumably would get lower confidence score.

Here's an example of the numbers I'm talking about. When I say a name was queried in a day, gmail.com was one name, but on this particular day, we looked at it 845,000 times. That's because a fraud protection platform is [stateless,] it doesn't know what other requests were made in the course of a day, and certainly, gmail.com has a lot of e-mail mailboxes people conducting transactions.

You can see how this falls off towards the bottom by the time you get down to gmail, googleemail.com, you're only down to 9300 times that particular day. So this is a really common activity that's happening in the database.

Here are the results of these days. I've named them A, W, X, Y and Z. A is the one from one month, and W, X, Y and Z are the ones from four consecutive names from another month. You can see that the number of API calls is pretty consistent, about 4 million per day. The number of names queried on those days, I already gave that spread earlier, so 474,000 on one day, 265,000 on another day. I'm not actually showing the results of W because that was related to an investigation. Let's just

---

say that it was a lot, and it was very much concentrated in a very small number of gTLDs that you could probably guess who they are.

The thing that I noticed is that most names were only asked in a single day. If we looked at this over a longer period of time – certainly, this is a small number of time, I'd be curious to see if that 87% would hold, but I have a hunch that it probably would. So most names across the namespace, we only look at them once over that period of time.

Then there's a couple that we looked at over two or three days. Only 1% are being called all four days, and those are mostly those e-mail addresses that I mentioned before. And then on any given consecutive day, what we're seeing is about 7%.

So then I process these, I actually looked through the records, which were scrubbed just in case there was any personal data in them. They were further scrubbed and have the contact fields removed. But really, all I was looking for was the registrar field – sorry, the server field which I used as a proxy for the registrar, and that's because people use all different words in the registrar field, there's different subsidiaries of registrars, there's about five different GoDaddy registrar names that all map to the same server. Some people put URLs in the registrar field, which screws everything up.

So I just looked at the servers, and this is what you can see. I broke down the servers, there were, as I said, about 1251 of them, and then just divided by the amount of time that we looked at them to get an estimate of how many times we would be asking – so if we were using

---

an RDAP system and not getting this data from someone else, how many times would we go to these individual servers?

In reality, it's going to be bursty, it's not going to be this steady state, but this gives you a sense of how big it would be. So you can see that there's one outlier, one big registrar. We would be asking them for data 43 times per minute based on this dataset, but the majority of them are much less than that. So you can see the next ten are on the order of ten per minute, and then it tails off pretty quickly. By the time you get to number 51 on the list, it's exactly one per minute, and then everybody else, the rest of the 1100 was less than one per minute.

So this should go into our considerations for a number of things. What sort of expectations we should be setting with contracted parties? If we think that our requests are going to be higher than this from time to time we might need to let them [give their] expectations. I think ten per minute is going to be noncontroversial, I hope, but setting the expectation there that it'd be between ten and zero.

Setting expectations for the RDAP server if we were doing these ourselves instead of acquiring them from somebody else ... I guess I don't have another point.

So then these are just the conclusions. We still have lots of uses for the public domain name registration data. We will be users of RDAP clients and asking you for data even before the EPDP concludes and we start requesting nonpublic data as well.

---

The consumption seems to be pretty predictable, although it can get a lot bigger during active digital crime investigations. And the consumption is well distributed over the namespace and between the contracted parties. So that's good. As I said before, during an investigation, it'll be more concentrated, that's just the nature of the thing. But not completely concentrated.

And then I have a comment about data retention periods, but I think it's really premature for me to have said this. And that's the other point I was going to make on the previous slide, is we'd like to think about how long you want to keep the data in the database. So if it's something I'm asking, I'm looking at 800,000 times a day, I don't want to be going and collecting that 800,000 times a day. But how long do I want to keep it? Because it might be changing.

So gmail.com does not change very often, we have high confidence in the thing, but other things might change more often and so we have to think, what is the good balance between refreshing the contents of that cache and just assuming that the data hasn't changed. That's something that more analysis would be required on, and it would have impact on how many times we would be going out to RDAP servers to get updated versions of the data. And then that's it for me, so questions for Marc and Mark.

EBERHARD LISSE:

Okay, thank you very much. Since you were not here this morning, I'm not asking you about JavaScript or something. Any questions?

---

ANGELA MATLAPENG: My name is Angela, I'm with the .bw ccTLD. My question goes to Marc number one. Firstly, I'd like to really appreciate the work that's been done in the development of [inaudible] protocol. I've done some backend development myself, I know that json is really convenient to work with.

My question is to ask about how RDAP will be preventing or reducing abuse during maybe bootstrapping and creating, and then the last question is to say that the issues around the WHOIS, being complaint with the GDPR, and maybe if you take it to ccTLDs, you can also talk about the countries general data protection X. So is this something that RDAP is concerned about? Thank you.

MARC BLANCHET: Thanks for the question. It's a difficult question. I'm almost going to say I'm not responding to the question. But one thing that RDAP brings is actually the fact that you could use all the tools that you have in http, therefore for example authentication, right? So we've been talking for quite a long time and still talking about differentiated access. If you have the right, you could authenticate yourself and get more data than if you are not authenticated. WHOIS was not designed for this.

But for the other kind of large consideration on abuse and other stuff, I'll skip.

---

MARK SVANCAREK:

I'll also answer part of that question. The RDAP protocol does envisage a system of differentiated access, but which differentiated access would be allowed or supported is going to be the result of policy. So there's a policy development process that's being expedited right now, so it's called expedited policy development process or EPDP, and that's where we're working out these issues about who should have access and how should they ask for it, and do they need to be accredited, and who would the accreditors be? And things like that.

So that policy is still being developed, and it'll be a while before we're able to really reduce it to practice, but the RDAP team, many people from that policy development process are members of the RDAP working group, so we're always keeping an eye on where the work is going and what's being planned to make sure that we're future proofing ourselves and not getting painted into a corner, as it were.

EBERHARD LISSE:

Just to be clear, what he just said only pertains to gTLDs. ccTLDs have nothing to do with this. If we use the same portals or registry software for example, and if that is developed for a gTLD, it has the tool RDAP also if it's used in a ccTLD, but it's a policy for each ccTLD manager to decide what information he gives. For [NA] for example, it wouldn't make a difference as far as we use WHOIS or if the RDAP tool becomes integrated into CoCCA tools in the future version, we do not allow information of residents of Europe according to the residency requirement that is stated at registration or modification. We do not



---

allow that private information to be published if it's a European citizen because it's [inaudible]. If it's a resident in Namibia, it's not a problem.

I also think we could refine this so that if maybe police for example wants to access any data, local law enforcement will go before overseas GDPR, then the local police comes from a whitelisted IP address or something and then can get all the information.

The point is it's a policy decision for each individual ccTLD. As for gTLDs, which is where the big problem is, it's a policy for altogether which is enforced by contract, and it's difficult to do which is why they have this – it's supposed to be expedited but it's an extended PDP, because they have to get so many people onboard at the same time, whereas with the ccTLD manager, board of directors, whatever this is how we're doing it, and then they can do it. Any other questions?

Alright, then I can release you two. Thank you very much. Now we want to have both machine learners to be on the board, and if also the host presenter, professor Vermeys, if he's here yet. I would also like him to be [inaudible] if the topic concerns his.

Alright, as I said this morning, I haven't read even the presentations yet because I don't like to do any editorial control in any way of the presentations, but I told them when they contacted us a while back that they should get in touch so that they liaise. Whether they did it or not, if not, we're going to spank them and that's it then.

Without further ado, there you go.

---

MOHASEENKHAN CHINWAL: Hello. I want to divert your attention from the previous topic to something called machine learning. I think there's a lot of scope of machine learning we have. As Mr. Warren Kumari just pointed out during our e-mail communication that we have lots of data, and we should start building up some things related to machine learning by using that data. DNS has a lot of data already in terms of server logs or so many things.

So my presentation is about a use case which we observed in .qa registry. I am from Qatar registry. We wanted to find out, we wanted to basically segregate the web content based on commercial or noncommercial. We just wanted to make a study of why people are not using the right kind of domain extension, and going for the wrong extensions.

So let's start with the topic directly. As you can see, I will go through a brief overview of what I'm doing, why I have built this small project as a use case, and basic soft mailing list, two slides. Then I will try to go to a sequential was of data science lifecycle of the use case. As you can see, there are business understanding, data collection, data preparation. There are a lot of stages. I will explain each and every stage of my use case based on the data science lifecycle.

Then I will also go through the scope of improvement where I'm facing problems and how it can be improved and any suggestion, comments from the people [inaudible].

So this is my problem statement. As said, that web content of .qa registry domains is, we were trying to segregate into a business and

---

nonbusiness content, and what I'm trying is to build a mechanism to predict the probability of a website being hosted for either commercial or noncommercial purpose.

The expected outcome is to study the patterns of the domain owners, analyzing the gaps between registry and registrar of why they choose the wrong domain extension or why choosing the right one. This problem can be sub-scaled. We can even classify websites based on government agencies, if it's an electronic consumer website, so many different categories you can scale this problem down to.

As you can see from this Venn diagram on the left-hand side, the machine learning is just a subcomponent of artificial AI subsystem, and data science is another major domain which acts as a fuel to the AI systems and big data and data analysis as part of data science. They overlap, they work together to form fully fledged scale AI system.

On the right-hand side, you have data science life cycle, how it looks like. You start with a business understanding, you collect data, you prepare the data for your modeling, and you explore the data or you try to remove the outliers or you try to remove unwanted characters, you cleanse the data basically to feed it to your model, then you test the model by accuracy or some other [inaudible] method, then you deploy the model as an API to the public tool so that they can use the model based on their data.

So this is very basic slide, we are going to be dealing with supervised machine learning here. There are various types of categories of

---

machine learning, and we will also involve deep learning into this use case.

As I was talking about the business understanding, we wanted to classify the web content from .qa to .com.qa, this is the business use case, and I started collecting the data from various data sources such as data.world, Data network, and Kaggle. These are very useful websites who have already segregated the web content based on various categories like business, finance, shopping.

So I used this data to start building my own data so that I can feed it into a model and try to get some more useful results. I have two types of classes. Class one I call as commercial class, and it is designated as binary 1 in the model, and class 2 is noncommercial. It's a very generic way of classifying things, that's why one of the reason I'm facing difficulties that the classes are really generic, and it's difficult to segregate them in terms of commercial to noncommercial.

So for commercial, I've used business and industry, finance and shopping, whereas for noncommercial, I've used law, people, society, pets and home and garden kind of web content. My data preparation involves like I programmatically connect to the websites, take the text contained. There are a lot of [inaudible] and models available to do this, like [inaudible] to XML. There's so many useful models where you can cleanse the data, you can remove the html code, the spaces, you can tokenize, you can have [stop words,] stemming, there are so many useful models in Python as well as there are machine learning models as well.

---

These blue tables were extracted from, again, one more python model called Pandas. What I'm trying to do is I'm just using two sources here. One is URL source, other is the [word text] content. The first stage would be I will extract only the textual content and I will also use the URL source, then I build my own.

As you can see, the second table [at the downside,] I built my own table by doing the text cleansing, removing the [stop words,] tokenizing changing the case, and after doing all that, I'd count the number of words based on keywords.

The second data table it is generating is by keyword selection. I have randomly selected some keywords which are related to business. Second category would be nonbusiness. So based on the keywords, I'm trying to isolate the contents of the websites into one category. For example, if you see the row one, it says that out of 10 words, 8 words matched the class one category which is commercial class, and two words only matched the noncommercial class. So it goes to the level called class one. So this is called a categorization table which I generate. This is a data which the algorithm will use to train itself to [inaudible] real world problems. The more accurate this data, the more accurate the model will be to real world problems.

So as you can see, there are three types of runs I have done. When I started testing the model, the model was [biasing against] labeling to the class one category, it was leaning to the class one category more. So I added more keywords to the noncommercial space, so as you can see, the class two numbers have increased, and of course, at the cost of

---

the standard deviation. It's getting increased as well, the variance is increasing, but the data normalization is happening. [You need distribution of values.]

This is the same thing, which was represented in terms of graph. You can see the word hit count is increased for class two [inaudible] 75% match. These graphs are saying that the class two categorization has improved with increasing the keyword classification. This is just a script where you generate this classification table based on the keywords count.

This is the neural network model, the actual diagram generated from TensorFlow, [inaudible] modules of our own module, which has three layers: the input layer, the hidden layer which contains two neurons, and output layer containing two classes, the categorization table which is class one and class two.

The input layer is a vector, I don't want to go into detail so that I lose you all. Basically, this is the three layers through which it works. This try to learn by [inaudible] propagation method. It'll take the input, do computation, and label the class based on the category table which I [fed] to it so it'll learn it. it's a supervised learning based on the labels I have provided, it'll try to compare the input and try to predict output on the output layer.

And this is the model table which have a lot of models I created, I just [took five out of them, and modeled two, show] some good results. As you can see, I tried two models with [ten document. Documents in this

---

sense] is just URLs with text content. You can call that 100 documents as [inaudible] and each document as a component document.

So I started with these two models and I see good results in terms of testing the accuracy. I have split the data into a training component and have taken a testing component from the same training component, and that I call as overlap. And overlap result is 70%.

Now I test the same model with a unique or alien test data, still with that also I can get an accuracy of 76%. And there are a lot of rows, so I don't think I have time to explain all this because maybe we have really less time.

As you can see, the confusion metrics is a way of testing the accuracy of models by various types of examples. The diagonal 370 and 1975 should be higher as possible [inaudible] accuracy, and model two shows good results for both types of test data.

There's a scope of improvement as well. As you can see, I'm not able to get the results what I wanted because of the method I'm using to generate the test data and the training data which is based on keywords classification.

Second, the more data you use, the better the accuracy. So working with neural network algorithm is really taxing in terms of CPU usage, and it takes a lot of time to train the model. There are better ways, but I tried to go through my own way of using python rather than using built in libraries which [inaudible] TensorFlows and other more libraries provide.

---

Then as I had, you have a very generalized set of categories, class one and class two. if you have more subcategories, that will also enhance to increase the accuracy of the models. And I'm using an encoding algorithm called [inaudible] algorithm. There are various advanced techniques available in machine learning these days. If we improve these things, there will be improvement in accuracy as well.

And of course, tuning the neural network pattern [inaudible] will also give you better results, and advanced data parsing techniques can also help to increase accuracy. So these are the improvements points which I found which I can work on again to keep the accuracy higher.

This is the workflow which I'm working on. What I do is I just prepare, gather the data from various types of websites like Kaggle and I build a training data for my own based on keyword classification, and I split the training data and test data from the same data which are obtained, then I train the model with the training dataset, then I evaluate the model with the test data taken from the same training data which I called as overlap, then distinct test data which has not been used at all in the model.

If I find an accuracy greater than 75, I try to keep the data set for future generalization. Otherwise, I trip it and I keep doing the same things to increase accuracy. These are some outputs which show off what I'm doing. This is just a tokenization table which goes to the websites, collects the text, and this is the category table which I generate through one of the script through keyword bases classification. I just segregate based on the number of wordcount hit on the particular class.



---

As you can see, the first example, there are three wordcount which is hitting the class one itself. There's 100% hit of class one keywords on that particular example, [inaudible]. So I generate this classification table and I convert them into numbers and feed it into neural networks algorithm, and it will try to label it as per the data I've provided it.

This is, again, a screenshot showing off how the model works, the output given. I'm using ten neurons. Each model runs around 10,000 times to adjust its vector notifications. And this is accountability how you calculate. As you see, 27 and 102 are the true negative values which count together to give the model accuracy.

That's it.

EBERHARD LISSE:

Okay. Thank you very much. Very deep stuff. Let's hear from JPRS what they are learning from their machines.

YOSHIRO YONEYA:

Hi, this is Yoshiro Yoneya from JPRS, which is .jp registry. Today, I talk about our research work for the prediction of [domain renewal rate] with machine learning, focusing on domain name label itself.

This is a table of today's contents. First, I made an introduction and explain about methods and data set, and results, findings, and discussions.

The introduction, as you know, to maintain DNS operational stability is essential to maintain stability of the Internet. And so that JPRS is

---

making next year's revenue [inaudible] with predicting renewal rate of newly registered domain name in high precision.

This prediction is performed heuristically by extremely experienced person in charge. But this is too individualistic so that it is very hard to transfer to others. So that we formed very small teams, three person, three members to evaluate prediction with mailing list and achieve several meaningful results. We thought that this result could be used for other TLDs, so today, I introduce our research work here.

So the methodology of the heuristic prediction is first we perform prediction on domain name renewal rate by least square method, and the experienced person makes some compensation if he finds some randomized names or some names for the test, and those amount exceed some threshold in his mind.

On the other hand, our proposed method is machine learning using a prediction by the supervised machine learning. And for this prediction, we defined 49 features that can be obtained by the domain label string itself.

These are the examples of which the expert person finds abnormality in the domain name label. So some randomized labels like [inaudible] bla bla, or something test names like unit test 1 or 2, or some words including event names or year numbers. And some names including online game names.

These names are not renewed, so these are the [trigger] of the lowering cost of renewal later. Here, I explain about the methods and data set.

---

This in the figure on the right down corner, so the data set and the training data.

So the blue shows the target data or test data for prediction. We've used eight months' data for the prediction, and we also used testing data for one year prior to the test data. Green shows the training data and blue shows the test data. So there are eight data sets.

During this period, we have some temporal fluctuation. So this fluctuation cause changes of renewal rate, so it is very effective to compare the statistical method and machine learning method.

And this graph shows how the statistical method shows the prediction. The red dot line shows the least square method or linear approximation for the 12 months, and the blue dot shows the moving coverage over two months.

On the other hand, our proposed method using machine learning, we use the same training data for the prediction, but the method is different. We use decision tree by the machine learning, and decision tree produces tree like the figure at the right and down corner. Every node shows the condition, which features met and how, what is the prediction [inaudible]. I will explain data. And for this to create decision tree, we define, as I said, 49 features which can be derived from only domain name label.

So one of the features is length of the string itself. We define several features regarding the disposition of digits. Digits means zero to nine, and [where is it] placed on the string.

---

We also defined the features of letters, disposition of letters and disposition of hyphens as well. And this table shows all features we defined, and the label is one of the feature, and for the digits we defined 13 features regarding to the place or consequence of digits and so on.

For the letters, we also defined the disposition with the letter or we included the ratio of vowels and consonants. And for the hyphens, we also defined the features regarding the disposition of hyphen, and the hyphen is before or after the digits or before or after the letter.

Using these features and the mailing list, the decision tree makes this kind of tree. On the top of the tree, it shows SD [inaudible], that means single digit or numbers included more than two parts. Then the prediction ratio is about 15%, and otherwise, the prediction rate is 73% and so on. So if you drill down these three to the lower [inaudible], then the prediction rate is more precise.

This is the result of our machine learning. We compared the three methods. One is the least square of the [statistical] method, and the other, second one, moving bridge, and the last one, the machine learning.

So for these eight months' data, machine learning is the best for five times out of eight. We also tried to predict more earlier or just immediate after the registration. So we tried to use the 12 months' training data, but it is earlier prior to one year, so you can see the table of left down side, the green one, consequence of green [inaudible] and blue [inaudible] is first our test, and the second test, we also used the

---

12 months training data, but after, it is one year prior to the first test. There is a one-year gap between the test data and the training data.

For this test, we also find that the machine learning is best for four times out of eight, and we also applied this method to the gTLD domain names which we treat as a registrar. JPRS is the registry for .jp and registrar for some gTLDs, so we used this method to gTLDs, and as you see, it is six times best out of eight.

So the findings of our research, we found that the existence of digit in label significantly affects down side. This is expectedly. And we found that existence of hyphen in label affects upside. This is unexpectedly. And we found outlying value of vowels and consonants ratio in the label affects the downsides, as expected.

This is a summary of the advantage of prediction by machine learning. We've confirmed that our proposed method improves prediction, precision of prediction compared to the statistical method. And our method was not affected when temporal fluctuation of renewal rate happened. And our proposed method predict renewal rate for each domain name individually which statistical method cannot. And our proposed method may have ability to predict next year's renewal rate immediately, [inaudible] domain name registered.

So here is ongoing work or future work. As I said, prediction much earlier is our next target, and this is ongoing work as I said in page 16 to 17. And correspondence to the trend changes or finding fluctuation is another trial. And [inaudible] applicability to other TLDs or registrars or researchers is also our interest.

---

As I said, we only used the features came from label, but we also have another data of the registration such as registrars or DNS host information, and so on. And we also want to predict second year or later renewal rate for the domain name.

As I repeatedly say, we just extract features from the label, that is digit, letters and hyphen. So we do not have any [linguistic] information at this moment, so the vocabulary analysis or [inaudible] some linguistic analysis could be another trial.

And at this moment, we are analyzing ascii domain name only, so the corresponding to IDNs is also a challenge. Okay, so I'm very happy to hear any ideas to improve our research work, and we made some testing tool, but it is very alpha version so that if we want to – if you have interest to try our method, please contact me directly. Thank you.

EBERHARD LISSE: Thank you. The idea of being able to predict how long a domain name, how often it will be removed, by name, is quite intriguing. This code that you use, is this going to be open source, or is this proprietary?

YOSHIRO YONEYA: We will provide our tool openly, but it is under development and needs some enhancement before publishing.

EBERHARD LISSE: Yeah, the point is if the tool is open, development will come by the use of it. Patricio?

---

PATRICIO POBLETE: Thank you. I can see how useful this could be for registry to try to predict future income for instance, but a question is, have you learned anything that could be useful for registries to try to take actions to drive domain name renewal rates up?

YOSHIRO YONEYA: It is not currently our scope, because we just want to know how the renewal, how many domains are renewed next year, and [inaudible] so I'm not familiar with some campaigns or something like marketing tours, so I'm sorry, I don't have good idea to higher the renewal rate.

UNIDENTIFIED MALE: Hello, [inaudible] Kuwait registry, .kw. Very interesting use case for machine learning Japan, and my question is I noticed the pattern of the deviation is fluctuating. Do you think if you keep running on a longer period, that deviation will be less and less over time, so you can predict more accurate results?

And the second thing, how can we apply this for registries with multiple renewal period, like two, three years? Because in your case, it's only limited to one year. And question for my friends in Qatar, also very good model, I'm wondering, have you started deploying the model? And if you do, can you showcase the results, the [compliance] when it comes to the content? Are they complied with your policies? Thank you for very interesting model, both of you.

---

MOHASEENKHAN CHINWAL: Was that a question? Sorry, currently as I said, I'm working on its accuracy [inaudible] categorization. I need help from people how to make it more accurate, because I have given the scope of improvements, if anybody is interested in machine learning, because I'm doing it as a part-time work in my own job. So anybody who's interested, they can contact me or I can contact them to learn more, and maybe [I'll reach a state to] deploy it and give it as free to the community.

YOSHIRO YONEYA: Thank you for your comment and advice. We are thinking about how long the training data we should use. It's still our research question. So we think one year is good because the tendency of users is cyclic so that one year is currently fitting our research.

But I think we should try, as you said, two years period or three years, and if increases the precision of the prediction. And for the two or three years registration period, we don't have such service yet, so it's future work. So if you try such analysis, the feedback from you is very welcome. Thank you.

WARREN KUMARI: Also on the Qatar presentation, on – I think it was slide nine that had all of the data. For one of the models, it looked like you had a really large vocabulary, and with [inaudible] coding, it seems like that might be kind of tricky to make it all fit. Did you have issues making it run, or do



---

you have that model take a really long time to train, or did it all just work okay?

MOHASEENKHAN CHINWAL: It takes a lot of time, because I'm using [one hot encoding] [inaudible] so the entire input vector is equal to the vocabulary size of the [corpus.] So currently, I'm trying to find out a way how to do it to make it more effective and fast. There are a lot of modules in TensorFlow and [inaudible] provide a lot of readymade stuff. Maybe I'll use it. I have just used that diagram for – I used [inaudible] visualizer from Google just to give a sample of 14 words. I just used 14 words vector, vocabulary for the data to [inaudible] just for the showing purposes. Thank you.

EBERHARD LISSE: Okay. Thank you very much. This was well beyond my level of understanding, but quite interesting nevertheless, and frightening in a way. We next have the host presentation. Thank you very much.

For those who are not here often, host presentation is usually what we call that we invite the host of the meeting to give us a presentation on a topic of their choice. If they're a ccTLD operator, we usually like to know a little bit about how they run their operations, but in this particular case, this is probably not the main issue.

And since they do cyberjustice, this is quite interesting in a way how they deal with issues not just about having remote participation that we have – or failure to have it – but also [what they use] on for example artificial intelligence or computer learning.

NICOLAS VERMEYS:

Good afternoon, everyone, and welcome to Montréal. My name is Nicolas Vermeys. I'm the associate director of the cyberjustice laboratory. In my own name and that of our director, Karim Benyekhlef, and our other researchers at the lab, we would like to thank ICANN for hosting the event here in Montréal and entrusting us with it.

So as was mentioned, I am here to talk to you about the cyberjustice laboratory. I am not a computer scientist, I am a lawyer. I don't know if I should apologize for that in this room, but that being said, obviously the level of the content of my presentation will not be quite that of the two previous presenters, which I'm still not sure if they were speaking English or not. I'm pretty sure they were.

But all joking aside, a little history lesson basically of who we are and so that you're probably all asking yourselves that question right now, what is the cyberjustice laboratory, and why are they hosting this event? So what is their link with ICANN and so on?

So I will give you a brief history of who we are, what we do, where we came from to then discuss where we're heading and how our research is connected to the work done at ICANN or to the interests of those of you in this room.

So again, I do apologize. There will not be any data sharing simply because I personally have people much smarter than me taking care of all of that in our center.

---

So what is the cyberjustice laboratory? First of all, I should probably start by explaining to you what cyberjustice is. As you can see on the screen, it is the incorporation of information and communication technologies into judicial and extrajudicial dispute resolution processes.

So in a nutshell, cyber justice simply means using technology in the field of justice and the judicial field. For those of you who have never set foot inside a court room, lucky you. For those of you who have, you've probably noticed that you're basically going back in time in most courtrooms in most countries around the world.

Now more and more courtrooms will have technology in them. Most of them don't. Here, if you just cross the street, you'll get to the Montréal courthouse, and in the Montréal courthouse, everything is still paper based. So any disputes, including disputes regarding domain names for example, well, those disputes are managed using paper-based technology. That in itself sounds a little ridiculous, but that's how the law qualifies it, paper-based technology.

So the idea is basically seeing how technology can make the legal process more efficient, which is obviously important for the government, but more important for the participants, the litigants, more accessible, and saying that meaning less expensive and make delays a lot shorter than they are right now.

So, who are we exactly at the cyber justice laboratory? Well, the laboratory is two things. First of all, it is a technologically advanced research facility. It is actually the most technologically advanced

---

courtroom in Canada. In a few months, we'll be the most advanced in North America because we're updating our technology and we have this friendly rivalry with the Center for Legal Court Technology in the US that is currently the most advanced courtroom in North America until we take first place once again, once our technology is updated.

So as you can see on the left-hand side of the screen for you, basically, it is a very large courtroom as you can see with screens, touchscreens, cameras, and with every type of modern technology that can be used to hear a trial.

In the middle, you have our satellite courtroom which is situated at McGill University. Our main courtroom is at the University of Montréal, and the purpose of our satellite courtroom, as you may guess, is to test out how you can have basically a courtroom set up in a hotel hall or in any other building in a few minutes to be able to hear a case and to connect with the main courtroom, in this case the one at the University of Montréal. But the idea of having basically participants in different courtrooms around the province, around the country, around the world.

And obviously, it's also a team of researchers. Now, most research centers that focus on law and technology will usually have a bunch of lawyers and one computer scientist because you need to have somebody who actually knows about the technology part of law and technology.

We take a very different approach, because we realize very quickly that when you're talking about law and technology, the problem isn't the

---

legal aspect. Lawyers understand the law, and it's not the technology. The technology exists, can be used and can be very easily implemented into the courthouse, into courtrooms.

The big issue, the big problem is that of the people that will then use the technology. So it's really just a question of understanding what the reticence of individuals who are using courtroom technology is exactly. So we have a multidisciplinary team, and in fact, the picture you have there is just a part of our team. We're actually more than 45 researchers from around the world, and they're not all lawyers and computer scientists. They're sociologists, psychologists, historians, and so on, to really understand how you can marry these two worlds of law and technology which to most lawyers at least seem to be two universes that should not collide.

So, what do we do exactly with this large team? Well, to understand our approach, you have to understand that we're dealing with two distinct areas. One we call our social legal research area, and that's the idea of basically studying how the law is going to be affected and impacted by technology.

The other area is the one that's probably a little more interesting to the people in this room is what we refer to as our techno legal research area, and that's the idea of using technology or developing technology to best settle disputes, and very soon – and you're probably still wondering, what does all this have to do with ICANN. Very soon, you'll see, I'll hopefully make it very clear.

---

So as you can see, we're building and developing digital platforms that simplify and improve the overall experience of the justice system for litigants and legal professional using this two-prong methodology.

So taking the first area, social legal research, what we're basically doing is reevaluating the role played by legal rituals. And you're thinking, well, rituals, technology, are we really going that far back? Yes, we are.

The main stopping point to the use of technology when we're talking about resolving disputes is very often something as simple as, well, that's not how the law is written, or rather, that's not how lawyers are actually using technology right now.

Best example, one of the main concerns about using video conferencing for court cases is that when a judge enters a room, you're supposed to stand and rise in order to acknowledge that the judge is entering the room.

The problem is if I'm sitting in front of my laptop that's in front of me right now, and I rise to acknowledge that the judge is entering the room, a lot of you understood why that is quite problematic, because the camera is aiming at a portion of my body that you probably do not usually do not usually show to somebody to explain that you have respect for that individual.

So you have to – now, this is, again, one of the most ludicrous examples, but it makes my point very easily. You obviously, when you're trying to incorporate technology into the court system or the court process, you have to take into account that these rituals have to be modified.

---

But in modifying them, you have to make sure that you understand why this ritual exists. So in this case, I rise simply because I want to acknowledge the presence of the judge, and demonstrate that I will defer to his or her decision making and therefore I want to show this person respect.

I don't need to stand to do that. I can use any type of mechanism just as long as the symbolism remains the same. But thus ritual of standing is a very clear cut one. Obviously, there are a lot of other rituals within the legal process that are not as clear cut, and therefore you have to dig through and try to understand what is ritual and what is not, what should be kept, what shouldn't, and only once that is done can you actually start and deal with how you can implement the technology into the system.

As for our technical legal approach, well, as I mentioned, the idea is to develop software solutions that will actually help the legal system. As you can see, the software solutions that we're developing, we're developing them from scratch using a modular approach and I'll get back to that in a few seconds. These tools are customizable, interoperable, web-based, mobile-friendly, cloud-based – so they're webtools and they're cloud-based for the most part. And open source, and I'll get back to that as well.

This isn't a sales pitch, just so we're clear, because we're not selling this technology. This is more of an invitation to all of you. if you are interested in these issues and would like to work with us with our team, we'd be more than happy to discuss things and share information.

---

So as I mentioned, our approach is a collaborative approach. One of the main issues when developing software for the courts and for the legal system in general is the fact that the main stakeholders are very rarely consulted.

In many countries around the world, there have been millions that were invested in developing information justice systems, and most of them failed miserably. Why is that? Because brilliant people like yourselves – so people who know technology – don't know the legal system, and those of us in the legal system very rarely understand the technology. And since we're from these two solitudes, and we don't speak the same language, very often the result will be a technology that doesn't necessarily adapt to the legal system.

So that's what we're trying to – we're basically trying to be the bridge between lawyers and technologists, and try to get them to speak a similar language.

The second thing that is extremely important to us is using a modular approach. One of the problems that we've seen over the years, and again looking at the different projects that were developed across the world, is that very often, they try to basically build this monolithic tool that will settle every problem in every case. And obviously, that very rarely works. I'm sorry.

And last but not least, we want to share these tools, and we share them within what we call the cyber justice community, which is a group of governments, courts that are working with us to develop these tools.



---

Now, you've been patient enough. What are these tools exactly, and how are we developing them? And how can they be useful to ICANN and its members?

Well, basically, we're working mostly on developing online dispute resolution platforms, or ODR. These are tools that we've been working on and others have been working on ODRs since the mid-90s to basically help settle disputes online, which is particularly interesting when you're talking about domain names since obviously, you are already in an online environment.

So the middle platform for example, the one on the left or your left, Medicys, is a platform that's used in France by court administrators. The one on the right, the Condo Authority Tribunal platform is a platform that we developed for an online court in Ontario, the neighboring province where they're settling disputes for condominiums, hence the names. So co-ownership disputes.

And the middle one, PARLe, the Platform to Aid In The Resolution of Litigation Electronically. You have no idea how long it took me to find an acronym that worked for it to be PARLe, because of course in French, parlez means to talk, and in English, if you've seen Pirates of the Caribbean, it means to discuss with an enemy to try to find a mutual ground. And yes, I learned history watching Pirates of the Caribbean. I do apologize for that.

So this platform was basically developed based on technology that we developed in the '90s to help settle domain name disputes, and I'll get back to that in a few seconds. But right now, PARLe is being used across

---

the province of Quebec to settle consumer disputes. It used to be simply for online disputes, so consumer disputes for goods and services that were procured on the Internet, but it was very quickly spread out to any type of dispute.

And what's really interesting is that the satisfaction rate for the platform is roughly 90%. So that means that people in many cases will actually prefer using this platform than actually going to court.

We're also developing online applications to control courtroom equipment, and also to present evidence on screens and control the evidence. So this is an application that you can have on a tablet and you can annotate evidence as you're presenting it in court. This is of course a little less interesting for this community.

So, why is this all being presented here at a meeting for ICANN? Well, you have to understand that these ODR platforms were actually built around domain name disputes. In fact, our very first platform developed in 1998, which in fact was the very first ODR platform to offer a mediation and arbitration, was built to settle domain name disputes. In fact, using this platform, Celine Dion managed to retrieve her domain name as well as – of course, she's Canadian, she had to take her platform. And a lot of other celebrities as well, and obviously, corporations that use the UDRP rules.

So the uniform domain name resolution policy that was implemented in the late '90s, early '00s, was basically what helped us create the cyber tribunal, which melded or evolved into what became e-resolution, which in 2000 became the first online platform to actually be officially

---

accredited by ICANN to settle domain name disputes. And in fact, in the year and a half, about, that e-resolution was functional, we helped settle roughly 500 domain name disputes across the country and around the world.

Now, as you can see, that was the beginning of the spectrum and the construction of the laboratory itself arrived a few years later, and now we're moving on to incorporating artificial intelligence into these tools thanks to a platform or rather a project that we refer to as the ACT project.

Now, what is the ACT project? Well, the ACT project is a partnership that was launched last year that, thanks to over \$6 million of funding from the Canadian government as well as a series of partners, including some of the people in this room, so Microsoft is a partner, they presented earlier today.

The idea is to see how artificial intelligence can be used within the legal system within legal practices to, again, make the system more accessible to legal stakeholders.

So we're working, as I mentioned, with researchers from around the world, with research facilities and research institutes from across the country here and around the world, as well as a great number of partners who are obviously interested in the use of technology and more particularly the use of artificial intelligence in the legal field.

The project was launched, as I mentioned, a little over a year ago, almost two years now in fact, and is going through different phases

---

starting with an inventory of situations where AI is used in the justice system and working with these partners to see how the AI tools that are actually developed can be made more efficient.

And that brings us to where we're heading, and where we're heading – as far as technology is concerned and not just the research – is the next phase of legal applications. And this again comes full circle to our very first platforms dealing with online dispute resolution for domain name disputes.

So we're right now working on software as a service, SaaS approach, that will allow us to create this virtual tribunal that can be used by any organization or any partner, including domain name dispute resolution institutes and services. So the idea being that they can use these platforms to help settle disputes between a domain name owner and a trademark owner.

These platforms will also be able to use artificial intelligence, and this makes a link with the two previous panels, with the few things that I did understand from them, the two previous speakers rather, and the idea being that we're right now working with partners at MIT as well as [inaudible] here at the University of Montréal to develop justice bots or legal chatbots that can basically go through a database of decisions, so for example the decisions – online dispute resolution decisions that were rendered for domain name disputes, either under ICANN rules or those of CIRA here in Canada, and to basically use these tools to train an algorithm that will then help litigants not predict the outcome of a case because predictive tools as far as we can tell are not very useful

---

because if I'm told that I have an 85% chance of winning my case and therefore keeping my domain name, that also means I have a 15% chance of losing my case and losing my domain name. And let's be honest, that's what I want to know, what that 15% is, not what the 85 is.

So it's not the idea of giving statistics and predicting the outcome of a dispute, but rather, giving information that is taken from all of these cases, all of these previous occurrences, and considering how the UDRP rules are structured, it's actually a tool that can be very easily adapted to those since they are very well structured and very coherent in between them, which is more than you can say for most legal rules and most laws.

And so in the next couple of years, we hope to make these tools available for partners, for the researchers that are working with us, and hopefully make all disputes easier to manage, even hopefully help evade certain disputes and doing so in every field, including, and more importantly for the people here, in the field of domain name disputes.

So I was told I had roughly 20 minutes. I think I even went a little too long, so it would be my pleasure to answer any questions you may have on this or any other issues. So thank you.

EBERHARD LISSE:

Thank you very much. We're not really strapped for time. Any questions?

---

UNIDENTIFIED MALE: [inaudible] [.id] domain. The whole system is still arbitration, it did not reach [to the code.]

NICOLAS VERMEYS: So the question that was asked is, is the system simply arbitration, or did it reach the courts? And the answer is it depends on which use of the platform. So in certain cases, so for those of you who are not familiar with online dispute resolution, it is usually a three-prong process. It usually starts with negotiation between the parties. If that fails, there's a mediator that joins the platform and helps the parties try to find a settlement, and if that fails, then there's an arbitrator that will join the platform and force a settlement between the parties. That's the classic model.

The only problem is that the classic model very rarely works, because the platform administrators need to have a business model, and historically, it has failed because it's very difficult to have a business model for that type of system. So what we're seeing more and more now is that that third stage, which is arbitration, is actually being replaced by adjudication by a judge.

So for the Condominium Authority Tribunal for example that I mentioned in my presentation, so they're using that system, but the third state instead of being an arbitrator is what they call a member, which is basically a judge working for the Ontario government, and therefore it is no longer arbitration, it is actually an online court.

---

MOHIT BATRA: Hi. This is Mohit Batra. I'm an ICANN fellow at this meeting as well as an RSSAC caucus member. So my question is, is the work of cyber justice laboratory [somewhat] related to the Canadian International Internet Dispute Resolution Center? I think this resolution center in Canada got accredited for resolution of domain names in May 2019.

NICOLAS VERMEYS: No, we are not accredited. Again, we're developing software pilot projects, basically. If you go back to the timeline, you notice that every project we had was about a year and a half or two years. The reason for that is they're basically proofs of concept and then we move on to the next thing and see if others want to use our systems and our platforms.

That being said, if they're interested in working with us, anybody in this room is a part of that group, I'd be more than happy to have that conversation.

MOHIT BATRA: Yeah, I just saw an association between those two since we're talking about the domain names dispute resolution, and both are in Canada, so that's why I thought to ask this question.

NICOLAS VERMEYS: Yeah, it's an excellent question, but right now we are not.

MOHIT BATRA: Thank you.

---

NICOLAS VERMEYS:                   You're very welcome.

NARELLE CLARK:                   Hello. Narelle Clark from the Internet Association of Australia, and a few other places but that'll do. What are you doing to weed out bias in the system? Bias being the big question in this area.

NICOLAS VERMEYS:                It is. And that's one of the reasons that we're not actually using a statistical model, but basically just getting information from cases and saying, well, if you look at the 100 cases where the fact pattern was very similar to yours, this is the normal outcome, these are the arguments that were presented, and it's not simply a statistical analysis because that's when the bias creeps in, or rather, that's when you cannot find the bias.

So in this case, the only bias that we're showing is that of the arbitrators who actually made the decision in the dispute resolution process, and since their bias is already there when you're studying the case law, we're not increasing or diminishing it.

However, these tools eventually hopefully can actually be used to show the biases of these individuals, so somebody can actually look at this decision making process and it can show to them for example in 95% of cases, you decide in favor of the copyright holder or the trademark holder and not the domain name holder. Maybe it's just a statistical



---

anomaly, or maybe you have some sort of a bias that you can then look at and work on. So it can also be used as a self-help tool.

One problem with biases that we do have however is the people training the algorithm or rather structuring the data. One issue we've had with the justice bot tool that we're developing is that obviously, court data is not well structured, so training an algorithm with it is relatively difficult until you structure it. So how do we structure it? We throw a pizza party for law students and we put them all in a room, they have free pizza and they structure the data.

We basically give them clear guidelines, but it does happen that a student will look at a decision and say, "Well, I don't really know how to structure this decision because it doesn't fit the normal structure of a court case, so I'm just going to put it aside." So that's the bias, the bias towards laziness of the people structuring the data that we actually have to work on right now and work with. That's going to be true in any case when you're training the data using court cases or arbitration decisions to try to find out exactly or try to make sure that all decisions are actually taken into account and not simply those that are easy to structure.

EBERHARD LISSE:

Thank you very much. I found this personally quite interesting, if only for my day job, because [inaudible] which I actually don't do much anymore, have on occasion to deal with court cases for negligence and other things, so it's quite interesting to hear developments to especially as far as mediation goes is quite interesting for our field.

---

Thank you very much. I appreciate it very much. Now Bruce Tonkin will penetrate AUDA. We are sort of 15 minutes ahead, so if we finish a little bit earlier, I'm inclined to have 15-20 minutes break after his presentation and then we go with the other session.

BRUCE TONKIN:

I want to talk a little bit about some of the testing that we've been doing in .AU to strengthen our security. The first thing I guess is just to give a little bit of an overview of our approach to security. So the first thing is that you need a framework, really, to describe the security in your system and the services that you have.

It's a little bit like – there's an analogy with accounting that most countries have an accounting standard, there's a standard set of terms that you use on balance sheets and profit and loss statements.

A security management framework is a little bit different, so a framework with a set of terminology, almost like a table of contents for things that you should consider in the security of your system.

So we've chosen an international standard called ISO 27001 as our framework. Then typically, we go through a risk analysis approach. From security, you look at your systems and you consider, do we care about the confidentiality of the data on that system? Do we care whether the system is available? Is it okay if it's not available on the weekend? Does it need to be available 24 hours? And then integrity, whether we care whether somebody can get in and change the data in some way.

---

So for us, when we look at that and we look at some of the systems we have, we have about 3 million names in our system. The availability of our system is considered critical for the country, so we aim for essentially 100% availability. And integrity is critical, so we don't want a bank site being accessed and having the DNS changed on a major system in the country. So availability and integrity is important.

When we consider confidentiality, the domain name itself is not confidential, but we do have a lot of contact information, and one of the things we've done is we've got historical contact information, so not just the information of the person that currently holds the domain name but we also have information on people that might have held that domain name in years gone by. So we've got about 20 years' worth of data.

So we would have at least 10 million names, addresses, phone numbers, e-mail addresses in the system. And while that may not be incredibly secret information, the aggregation of it could certainly cause significant harm, so a company could get that information and launch phishing attacks, that kind of thing. So we view confidentiality as important, particularly around the personal contact information.

Then when you've got your risks assessed, you then want to look at what are the right security controls to protect against those risks, and there are various approaches. We've used three different collections of controls I guess, one in Australia, we have Australian Signals Director, which basically is a part of government that really specialized in being able to break into systems, particularly at time of war, but they've given

---

guidance for business world as well, and they have what they call the essential eight, which is things like make sure you have two-factor authentication, make sure you patch the software whenever there's a new security release when you apply it, reduce things like the ability for Microsoft macros and things to run on machines that might be used for system administration, etc. So they have eight essential things which we're implementing.

The ISO standard itself has a set of recommended controls which we're implementing, and then the Australian government has an information security manual that's designed for anything from basic information all the way up to top secret information and specifies the different levels of protection, what controls you should have in place. So we're putting all those controls in place.

The next thing you might want to consider then is you've worked out what your risks are, you've applied a series of controls. How do you know whether the organization is actually applying those controls? Typically, a board might want to use an external company to verify that the documentation's there, the controls are in place.

This is pretty similar to what you would do with an accounting audit. The board will typically use an external audit company to audit the company financials so you know that both the financials are properly reported, but also know that there's proper controls in the way spending is done, etc., against an organization's policies.

So just as routinely organizations will get an external audit for their accounting, I think in critical infrastructure, generally a board would

---

want to make sure that you've got an external audit of your security system.

The next step is to look at getting some outside scanning of your environment. And I'm referring to this as a vulnerability scan. This is where you might give to an external company a list of the IP addresses or IP range that you use in your organization. They scan that whole IP address range, try and find out what devices, what computers, what servers are at the end of which IP address, then they run a series of automated queries against that, try and work out what sort of TCP ports are open, what services are open, etc. on those IP addresses. And those automated tools would typically produce a list of vulnerabilities.

In some cases, you might just accept them as that you already know you've deliberately opened up a particular port for a particular reason. Sometimes it's things that you didn't know about. And it's quite common nowadays if you're making a change to your system, whether it's a change in your software or a change in your hardware or network configuration that you might run your own tools against these, and there's plenty of open source and commercially available tools to do that.

The next step, which is what we do, is we get an external provider to run those tests, and we do this once a year, and they produce a report and they say, "Look, here's some vulnerabilities that we think were high risk, here are some medium risk, and here are some low risk vulnerabilities."

---

And often, there's a cost to fix these vulnerabilities, and most organizations will go, "We'll just fix the most critical ones, and the low ones cost a bit to fix, and because they're low risk, let's not worry about it, just focus on the critical ones."

But it turns out that a hacker that's got a lot of time will actually string together a lot of those low vulnerability ones and be able to get into your system.

The critical ones are the ones that someone opportunistically might try and get into your system quickly, like a cybercriminal and they're simply trying to find the equivalent of sort of walking down a row of cars and seeing which cars got one of the doors open. And if all the doors are locked, they'll just move on to the next car.

That kind of deals with what I'd call the critical risks, but the lower risk ones are more, have you left the back door window open a crank, or left the sunroof open just slightly, and is that enough for somebody to get in? And somebody more determined, like if there's an expensive computer sitting on the passenger seat, they might go, actually, even though the door is not easy to open, I can through a crack in the sunroof get a bit of a wire in there, then I can somehow unlock the door and get to the laptop. So they become more determined with some of the low risks that you have.

So penetration testing is – so you've done an automated test and you've closed some of the vulnerabilities. Penetration testing is when you actually try and see, well, can I use those vulnerabilities to get in? And effectively, what we've done in this situation is we've then looked to pay

---

for someone to see if they can get in. So it's a little bit like if you lose the keys and you ring up your local car automobile club and you say you lost your keys, and they manage to break into your car in about five minutes where they make it look so easy, and then manage to get the keys that you've locked in the car.

This is sort of the first stage of penetration test, get someone to see if they can use some of these tools and actually get in, and you might be looking at the most basic sort of testing here. And typically, this is done in cooperation with the actual IT team, so you engage someone externally, they tell the IT team, hey, tomorrow I'm going to try and get into your system, let me know if you see anything on your side, keep an eye on your logs, keep an eye on the alerts because I'm going to get in tomorrow. And you might even tell the person that, hey, here's the results of the vulnerability scans, here's some low-level vulnerabilities, see if you can get in using those vulnerabilities.

This is often referred to as a blue team approach because it's fairly cooperative. You're just deliberately engaging someone and you're working with them a little bit, and just seeing, watching them really, "Show me how you get in," a bit like getting someone to break into your car. How do you get in? "I got a bit of wire and I managed to unlock the door," and therefore you might make some other steps to put in place to stop that happening again.

The next level up from that is what we've recently just engaged, which is a red team penetration test, and this is where you're actually employing a team of hackers to spend months trying to get into your

---

system. So this is assuming that you've cut down all the most basic ways of getting into the computing system, and you're now going to give a team of people pretty much carte blanche to get in any way they can. This is a lot more expensive because you're obviously paying for the labor of that team to try and get in. So a standard sort of industry vulnerability scan might cost, say, \$20,000 a year. This is going to cost more like \$100,000 assuming you're employing people within the country to try and get in. Maybe it's cheaper if you use hackers in a country where labor is cheap.

So this is a lot more expensive, but it's quite an interesting exercise to see what they do. A lot of red team penetration testing, the first thing is you've got to define the scope. So what do you want them to break into? Are you just asking them to break into your organization through whatever electronic means, or are you actually letting them break into your building and allowing them to do that? What proof are they going to have that they've managed to break in? What particular target are they going to go after?

So with a red team penetration test, you basically set up a set of prizes, and you basically say, look, I want to see how good our office security is. Does anybody challenge u when you try and get in the office? Could you walk into the CEO's office that might be several layers into an organization? Can you get a document off the CEO's desk?

So you might create a document, put it on the CEO's desk and it may say "For penetration testing team," and leave it on the desk, and a week



---

later, they come up and they go, “I've got this document, I managed to get it off your desk last week.”

So you set a physical prize. You can set an electronic prize, so you might have a file with fake data but you could say, I'm going to have a file called – if it's a membership organization, here's the list of members. if you can find that file in our system somewhere, that's proof that you've managed to get in. You might create a database. In our case, we're concerned about things like the registry database, the equipment and services that publish the zone file.

The other thing is what systems are you wanting them to try and get into. Is it just your production environment? Generally, if you've got multiple environments – and we do – so you have a development environment, test environments, you might have a user acceptance testing environment, and then production networks. Quite commonly, where hacks happen these days is they go after these test environments which might have real data on them.

So often, software developers are messing around with things and they go, “It would be good to have a copy of the actual registration database to do some testing on,” and they forget at the end of the testing to delete the file, and it's sitting on a server that's not very well secured. So often, you can get the most valuable data by going after the test environments because they're not secured as well as the production environment.

In this case, you'd say to an external team, here are the different environments, here's some files that we've put in these environments.

---

See if you can get them, and then obviously, tell us how you went about it, etc.

So just talking a bit about some of the techniques that the penetration testing team used, we had them go after our organization for several months, we allowed them access – if they could get into our physical office. I actually had to write essentially a get out of jail free card. So if they actually got found in our office, then just the instant before they get thrown out the window, they could produce a piece of paper that says, “Look, I've got permission to come into the office and I'm not really a burglar or a criminal or something. Please don't kill me.” So they require effectively a letter that they carry around with them that they can produce if they're challenged by a person in the organization.

As mentioned, this occurred over a couple of months. Some of the things that they do or did is they run their own vulnerability scans, they do automated testing to see what holes in the system you might have, you they look up on your website and they find out a lot of us post a list of staff that are on our website, which is gold for a hacker because they now know – and then our e-mail addresses, if you know somebody's first name and last name, you know what their e-mail address is, their `firstname.lastname@auda.org.au`.

And then they know what he relationships are, so they know who the CEO is, who the managers are, so they can construct e-mail messages that look like they're from your manager, so the e-mail comes from someone that looks like your manager asking you to do something. This is typically then use for spear phishing campaigns, and some of the

---

techniques could be – “I want to pay you a bonus and I just need your bank account details,” and they might do this at a time of year you might be getting a bonus, maybe it’s in January or something so it looks like from your boss with your bonus information, so people go, “Well, better open that e-mail.” And they might be asked to log into something, and based on that, they get username, passwords.

I think one of the attacks they used on us is I think they knew that the staff had been invited to a particular event. I can't remember what it was, might have been like a staff party or something, and it was made to look like an Eventbrite. You’ve got to log in and register for the event if you want to go to this wonderful party. So they manage to get credentials that way.

Another technique they’ll do is they’ll go to the physical office and they’ll drop USB sticks. Sometimes it’s just a blank USB stick and you see a USB stick sitting on the table and you go, “Just curious, I wonder what's in there.” Sometimes it can be a bit more obviously, they'll actually put a label on there like “Game of Thrones Season 10” or something and they go, “Oh my god, didn't even know there was one. I better check that out.”

So they put something on the USB stick that looks enticing and that gets the staff to plug it into their computer, and they put malware on their computer.

To get physical access to your office, what they’ll do is they’ll do research, because they’ve got plenty of time, so they’ll work out what building you're in, who’s the building manager, they'll then ring up and

---

say “I’m John Smith, the building manager, I’m going to send up a person that needs to inspect your air conditioning,” let’s say it’s hot weather as it often is in Australia, or in Canada I imagine they’ll come and inspect your heating, and so they ring in advance and they give the name of a person, they’re telling you that it’s the name of the building manager, and so the office gets that phone call and someone turns up dressed in overalls and with a big torch or something so it looks like they’re the official person to check the building, and so they get let in, and then they wander around the office and eventually get into the CEO’s office.

The other part of this testing is we told our staff we were doing this testing, but it’s over a period of time, so it’s not like I could say to the IT team, “Hey, you better keep an eye on things tomorrow because there might be a hacking attempt.” We just say sometime over the next two months, we’ve engaged this team. So they’re not going to be on heightened alert every day, they quickly revert back to routine. Same with the office. We tell people we want to make sure that this office is secure. You shouldn’t let anybody in without checking ID and making – if it is the building manager, you should ring the building manager not just take it on faith that somebody’s rung you and said they’re the building manager, because you don’t know what they sound like on the phone. You should ring back to the official number and check that the person has actually been authorized to come into the office. And then even if they do come in the office, you should have someone walk around with them, not just – they look like they’re in overalls and

---

they're obviously a maintenance man or a woman. Okay, you're at the front door and you just assume that they're not going to do anything.

So we did quite a bit of training before we even started this exercise, and we did tell people that we're doing this red team exercise as well, but still, they were able to get through some of the things that we had in place.

One of the key things here is it's not just the IT configuration that gets tested, it's really the organizational capability to detect and respond and block attacks. So you're testing in this case our IT team, could they detect when somebody came in? Were they able to respond and actually fight back, so block various attempts of attack? And what were the steps that they took after they became aware of a particular hacking incident? How did they train staff and so on?

I want to emphasize here that this is not an exercise in trying to embarrass staff, so we had in our spear phishing campaigns a board member that got picked up because they had access to board member information as well, and we had a staff member picked up through a spear phishing campaign. But we didn't embarrass them, we didn't publish their names, didn't tell anyone else in the organization. We just spoke to them individually and gave them some coaching. So you don't try and embarrass people. You treat this as an exercise and sort of learning and improvement.

One of the things that they try and do when they get physical access to the network or to the building is actually plug into your internal network. A lot of organizations think, "Well, I've got firewalls and this

---

fantastic security,” but what happens if someone just walks in the office as a maintenance person and plugs into an open ethernet port sitting on a desk? They're inside the firewall now. What can they do?

Wi-Fi is another one. You walk outside any office building and there'll be an office network with a Wi-Fi. The question is, just using some automated tools, you can crack Wi-Fi passwords and get onto the Wi-Fi network. And again, what capability can they have?

We engaged the red team to try and attack not just the outer office but also the office of a registry operator, which was Afilias, and they used different techniques. In the case of Afilias, they'd done the research on the building and they claimed to be inspecting a water leak in the ceiling, and they managed to get in the door but they didn't manage to get much further than that.

In our case, there was a well dressed gentleman in a suit and so on, and he'd come into the office, and we have conference rooms that – we have sort of a secure part of the office and a part of the office that's more open for business meetings. There's no locks on the door once you get into the floor level.

And in this case, the person came in late on a Friday afternoon. It happened to be a Friday afternoon where quite a few people were a combination of either sick or traveling, so very few people in the office, and it's a common technique, is they work out when an office is not likely to have many people in. So in Melbourne, if you turned up on a Monday morning after a football grand final, probably not many people in the office.

---

So they pick a time when it's quiet. In this case, it was a very quiet period at the office. They got into one of our conference rooms, and then they locked the door of the conference room, just left the lights off. And we've got frosted glass on the conference room, so it's not easy to see if anyone was there. They left the lights turned off and they just spent several hours in that conference room.

They tried to get through the network devices. We had a videoconferencing system in that conference room. They tried to hack the videoconference system, so they started basically pulling it apart. They did eventually get caught and they had to pull out their get out free card, but it was somewhat alarming when you walked in on the office and there was wires and cables, and things apart on the desk, and they're saying, "No, we're just here visiting from ...". In this case, they were claiming to be a visitor from the local telco, but you don't generally have your local telco taking apart your videoconferencing system.

We had a number of things in place to protect against that. One is that you can't just plug into an open network port. We have to effectively have the MAC address, if you like, of that device in the system, so it wouldn't allow them to plug in their laptops that had hacking equipment, etc. on it.

We ultimately blocked their attempts of hacking through physical. In the end, we did actually give them a fully configured laptop. We just wanted to test if someone had left their laptop down in the coffee shop downstairs, what could you do with that? So we gave them that with a

---

user account and also wanted to see how far they could get if they had basically got physical access to the office and had managed to find somebody's laptop which had various credentials to get access to the network.

One of the things we found in this exercise was you really have to assume that your credentials of a staff member can be obtained. So again, they used a spear phishing attack, they managed to get two sets of credentials.

What was surprising for us was they were able to log in with those credentials. They told me later and said, "Oh, you should really have two-factor authentication turned on." And I said, "We do. Every one of our staff members has two-factor authentication and logging into the network."

And they said, "Oh, no, we used the username and password and we got in." And it turned out that the staff member credentials that they managed to obtain was a staff member that was working from home in an area a little bit outside of Melbourne, and their network connection was a bit flaky sometimes, so they were getting messages on their mobile phone saying, "Do you approve this connection?" And they were just assuming it was the network that had dropped, because they were logged into the network, and they were just being asked to reapprove. So they just sat there and hit "approve" whenever their approve message came in.

The thing that our IT team picked up is they checked the logs and checked when people log in, and they saw that this staff member I think



---

might have been logged in twice or they were logged in but then the IP address didn't match their normal home IP address range. So because it was from an unusual address range or that unusual IP address, we actually blocked – or the IT team, just routinely, didn't know it was part of the penetration test. They just thought that's really weird this person is logged in an then they blocked the IP address range. So their hacking attempt was blocked.

But again, this happened to be done during the week, and the staff happened to detect it, but if it had been on a weekend and they'd done it, they could have been in there for days, Saturday, Sunday on the weekend before anybody would have noticed that there's this person logging in from an unusual address.

So my lesson from this is even though we'd have training, we tell people not to get caught – we run monthly phishing simulations, still, very clever people – and these were clever people we engaged – they designed the spear phishing so it's very hard to tell the difference between a real e-mail and a fake e-mail. So you just have to assume your staff are going to get picked up, and then the question is if you assume your staff are going to get hacked, what protections have you built in the organization to check for unusual activity, what controls have you got that people, even say a custom service person, how far can they get in the network without getting into production systems and things like that?

So they were able to get a user account and they were able to access our systems, probably for a couple of hours before we cut them out.

---

The other thing in red team engagement is after they come back to you with a report and they tell you different ways they managed to get in, you've got to be really careful that they clean up after them, because the other issue that's occurred with red team engagements is you pay these teams, they get people to load malware onto their machines, they get people to get onto servers and they load software into those servers, and one of the things that – I've spoken to a couple of red teams now, and the techniques they would use, especially when they're doing this over a long period of time, they also assume that the local IT team's going to fight back and find different ways of shutting them down.

So before they try and access anything important, they put as many backdoors and multiple ways of getting back in, because they know probably the primary way they got in is going to get detected and they're going to get shut out, but they've actually left a few other traps in. they stay silent, and this is what would typically happen with a highly motivated attack.

You see this with national or state-based espionage where they're effectively saying we can employ a team of 20 people and we're happy for them to spend a whole year going after some organization or going after something critical, and the first thing is they'll stay there for some time and stay as far below the radar as they can and put as many things in the systems as they can so then when they're going after the high value asset, whatever it is that they're trying to get, the nuclear attack plans or something, that even though they know they're going to get blocked in a few places, they've left open enough backdoors.

---

That's fine, but the problem is you've engaged this red team to do that, and it's fine that they've put in these methods to get in, but you also want them to unravel it afterwards, because bear in mind your IT team's got no idea of what they've hidden away in the systems. So you want to be using a pretty reputable company that documents what it's put in and removes it when they leave, because then someone that is not the team that you've engaged but someone really trying to hack into your system, they're actually going to start taking advantage of malware that's on some staff person's computer or something that's left behind on a server.

The other lesson with this is – this applies to vulnerability scanning, but also any penetration testing you do, is again, it's very easy for management just to look at that report and go, "Okay, here's \$100,000 to fix all the critical stuff and then we're good, we've got rid of all the yellows and the reds, and we're fine."

But the reality is that the hackers that have got plenty of time and they're going after a high-value asset will actually string together multiple low, green kind of vulnerabilities in your system. You do need to close those as well. So it is actually worth investing and closing the low-level vulnerabilities, and that's kind of what you pick up when you do this red team testing, is you've got to assume that we've put as much resources as we could to shut down anything major, but they're still able to partially get in, and if we gave them another few months, they probably would have got in further. So we found it a useful exercise.

---

So our intent is to actually do this every year, and the other thing we've added is more resourcing to the IT security team to be looking at logs and other things, because the other challenge for a lot of us is that you might buy off-the-shelf software and there are some really good systems around, but if you've got no one that's actually checking the logs and checking the material that these systems produce, you might as well not have it. So you've actually got to invest in your own internal resources that actually checks for logs, checks for unusual activity.

My IT team, when I'm here in Canada, if they didn't know I was here, they would send me a note saying "It looks like you logged in from Canada. How is that possible?" So I'd let them know where I am in the world and they can correlate my IP addresses with where I'm logging in from.

The other thing you're looking for is, can I actually travel at the speed of light and be in Canada one minute and the UK in the next minute? That would be suspicious too. So these things that you look for if you were looking for patterns of access to your system that look unusual. So again, you invest in the tools, you put in your access controls, and you have to invest in people that are available to actually check the data that all these tools produce, and look for that strange activity that might be compromised staff member, their laptop's compromised in some way and it's doing something unusual.

So that's all I had, Eberhard.

---

EBERHARD LISSE:

Thank you very much. This was not so frightening, I must say. In smaller places like for example in our registry with 1.8 staff or so, we run a credit card payment system so we are PCI compliant, we have got security metrics looking at this on a regular basis. That's something simple that anybody can do.

We use CoCCA tools. As you know, we are reasonably on the latest system, and we've configured it that if somebody logs in from an unknown address with a new browser or something, we get an e-mail, and the registered address also gets an e-mail.

So the basic things are, as I usually like to say, humans are good to solve problems but are bad at doing it all the time, because it gets boring. And it's very difficult – I wrote a handbook where everything I do, I add, but it's a [schlep] to remember. You have made a change, open the handbook and add it into the handbook, and keep the handbook secure.

But the idea is that it's sometimes easy to forget how to reconfigure system. If you want to spin up a virtual machine, you have to go, "How did I do that?" And if something happens to the hardware, if something happens to me, I want to be able to – this is things that we have discussed already ten years ago in Mexico. This is something that was discussed over this weekend here by the TLD Ops business continuity.

In bigger systems where there is more staff that don't talk to each other every day or don't maybe even know each other and they work from different locations, it is much more appropriate to invest into outside

---

testing. But for even smaller systems that have sort of the human resources under control can be hacked.

Recently, my understanding is one country ccTLD registry was seriously hacked and lots of government websites were defaced and stuff like this because they used an older version of a registry software, and then didn't use two-factor authentication and got their passwords spear phished.

It's relatively simple nowadays to have two-factor authentication, to configure the e-mail to report when there is a new access from an e-mail address, to lock out old passwords if they haven't been used for a certain amount of time, to force the users to change passwords on a regular basis. Even if they squeal a little bit, they have to do it on a regular basis because they have to remember it, but these are simple things that can be done, and we are often just too lazy to do it.

BRUCE TONKIN:

Yeah, and some of the other things they do is they look on the dark web, if they know the names of your staff and they're looking in the dark web for files of hacks from other organizations and they get username and passwords. Your biggest problem is people using the same password across multiple systems.

EBERHARD LISSE:

There is nothing wrong with using randomized passwords and put it some – and then also configure your password generator to generate passwords that you can easily read when you print them. for example,

---

refrain from using Os and 0s. But if you use too many special characters, it's very difficult to type them in, and then you're sort of tempted to just store it somewhere and copy and paste it. That becomes a big hole.

[ERIC:]

Hi. Eric from .ca. Just regarding the passwords in regards to changing password regularly, there's more and more emerging research that's coming out that suggests that regularly changing password is not actually a good thing. If anything, it actually makes it more difficult because traditional user behavior is they're just going to write it down or they use standard kind of templates. So instead it's use a longer password and only change if there's an indication of compromise.

You're suggesting that change passwords regularly, but yet emerging research says that's not a good idea at all. Which way should we go?

BRUCE TONKIN:

Suggesting changing passwords regularly?

EBERHARD LISSE:

We force our registrars to change their passwords every 90 days, plus we mandate two-factor authorization. They cannot log in without that. And we only correspond with the registered e-mail addresses, and if this address has been changed recently, then we would need a different form of communication.

---

We are thinking of introducing or forcing all our registrar to communicate with us via PGP signed e-mail so we have some form of reasonably secure communication.

I agree with you that if you change really random passwords it's difficult, but if you have got two factors, password and Google authenticator or physical dongle, that increases the security.

[ERIC:]

Sure, I don't think anybody's disputing two-factor. That should be a base. But in terms of a red team engagement going after users at an office in your use case, having weaker passwords that are changed regularly actually, as the research suggests, puts you in a weaker position than having longer passwords that are not rotated, unless there's an indication of compromise.

BRUCE TONKIN:

Yeah, and that's actually what we do, so for us, our staff, they have to have very long passwords, and they're typically a phrase that only they know. It's not as easy to guess those passwords.

The issue I was talking about with the use of the dark web is people using, let' say in this case, the same phrase somewhere else. So if some other system has been hacked and they've got that password information, so really, the lesson from that is don't use the same password in multiple places.



---

But we did find that, so one of the people that was hacked when they got the username password, they then actually found that they could get into Facebook and Twitter using that account, which you can then get more information.

So a lot of what the red team people do is once they get in one way, they get someone's username and password, they go, "I bet this person's got a Facebook account. I'll try to get in there," and they just collect more and more information about them.

So yeah, really the message there was don't use the same password in systems.

[ERIC:]

Right, that's why I was just focusing on the length of the passwords as opposed to password reuse and so on. Thank you very much. Appreciate it.

MOHIT BATRA:

One comment on the thing we were discussing. Passphrase is another thing which is being used nowadays apart from just long passwords. So they're the norm now, passphrases.

Second, question to Bruce, was ISO 27001 the default choice, or was some other security framework also considered?

---

BRUCE TONKIN: If you're talking about the security frameworks, there's a framework that's used in the US, I think it's NIST. That's another security framework. I think different countries have probably got their own security frameworks as well, and the Australian government's got one, which is effectively the [inaudible] information security manual.

Not only are we looking at ISO 27000 or AUDA, we plan to actually make it a requirement for our registrars as well, so we wanted to have a standard that – and our registrars are in different locations, so not all of our registrars are in Australia. In fact the majority now are basically outside of Australia, so we wanted to have a standard that you can actually get accredited for or audited in different locations.

So if I just chose one that was Australian or just one that was US, then they wouldn't be able to get a low-cost local company that would be able to do the audit. So that's why we chose that standard.

MOHIT BATRA: I also [inaudible] standard [inaudible] standard.

BRUCE TONKIN: Yeah, the ISO 27000 is also an Australian standard, yeah.

EBERHARD LISSE: No, sorry, the next –

UNIDENTIFIED MALE: [inaudible].

---

EBERHARD LISSE: Take it offline, discuss it with him. The e-mail addresses are on the agenda. Paul Hoffman, root early warning.

PAUL HOFFMAN: Good afternoon. I'm Paul Hoffman, I'm part of ICANN Org, that is I'm part of the staff at ICANN. This presentation, I'm going to try to get through the slides reasonably quickly because the purpose of the presentation is to have a good mic line.

EBERHARD LISSE: We have time.

PAUL HOFFMAN: Okay. So as you're hearing this, please, if you have questions, hold them but please do know that we will have time of the mic line. The first bit is why are we trying to create a root zone scaling early warning system and why are we talking about this now.

Of course, the root zone has been increasing in size actually for more than 25 years, but certainly a lot more recently, so we'll talk about that. We'll talk about a proposed system and the goals for making the proposed system and the non-goals for – since this is tech day, a lot of us in the room like designing technical systems, and we can easily run down rat holes very quickly. So we'll talk both about the goals and the non-goals.

---

When you are creating systems like this where you're going to be doing measurements and such, it's good to have a standard vocabulary, so I actually have a slide on that that's actually mostly lifted from recent work done in RSSAC, and then we'll talk – again, I'll try to open up a mic line. I'll ask you a bunch of questions and we would like to hear from you.

Right now, we are considering creating a root zone scaling early warning system, and the reason is that SSAC asked us to do this a while ago in SAC 46, and there are the words that they asked us. They noted the lack of forward progress from ICANN, so they asked us again in SAC 100, which is much more recent.

So we're sort of taking that seriously. Important groups are asking us to do something. And both of these requests are based on a document that was published in 2009 that ICANN commissioned about scaling the root. For those of you who are new, remember the root zone was pretty stable in size until about 2012 when we took in all the applications for the new gTLDs. It was growing, but extremely slowly.

So a lot of this concern, especially the concern around SAC 46, was happening and the paper was happening when there was the expectation that the root zone was going to grow rapidly, but it wasn't clear how rapidly because of course, we hadn't actually taken applications. Could have been that the root zone was going to grow by doubling in size. It could have been that the root zone was going to increase two or three orders of magnitude. Not known.

---

And again, so these two requests came to us. Other people have also already looked at root zone scaling. For example, RSSAC, the root zone operators themselves created a system in a document called RSSAC 002 for self-reporting and such. So it's not that no one is looking at it, but this is something where SSAC has asked staff to please get a little bit more active on this, and so we're looking at it.

Purpose of my talk today is to say if we did this, what do you think we should do? So to be really clear, we are not talking about what should the maximum size of the root zone be. The root server operators haven't expressed any concern about that. The concern that was expressed in SAC 046 and in SAC 100 is the scaling, how quickly, if ICANN dumped a bazillion new TLDs in, would that have an effect? Not what was the size, but the scaling factor.

So that speed thing, certainly there are limitations as we've discovered in getting it, but what if those limitations went away? What if ICANN all of a sudden became way more efficient at adding TLDs? And we scaled a lot during a year, or a month, or something like that. Those are the issues we're looking at, now is there a maximum of the eventual size.

So given that SSAC asked us to do an early warning system, what are some of the big parameters on that? One of them would be easy to publish metrics that can be used by the ICANN community. That's really important. Any of the metrics that we are collecting should be usable by the whole community to determine if there's a need for thresholds. May not be a need. We don't know. You'll see this in a few slides.

---

But definitely the goal is not just to collect numbers but is to collect the numbers and publish them for the community.

Once the system is set up, it should just run by itself pretty much without a lot of maintenance, so it's not like – we don't want to set up a system that takes a whole lot of people running it all the time. That tends to be a more fragile system, but it also tends to give hard to interpret data because the semantics of the data changes over time. So get it right the first time, maybe add things later, but pretty much [put it in.]

And then also, we would like the same system that we set up to be reproducible by other people. So we could set up a system, everyone could say, “Yeah, that looks pretty good, but what I want is two different things over here.” And maybe there's not community agreement to have those two things. But if you as a researcher or you as an organization that has researchers can reproduce exactly what we're doing, then you can do what you want as well. We don't believe that any measurement system for anything in the DNS is ever the perfect measurement system.

You've been hearing the talks today. Those of you who have come to previous tech days have heard other talks about people doing interesting measurements. And you might have said, “Oh, I would have done it a little bit differently.” Great. Whatever system we come up, we want to make it easy for you to do it so you can do it a little bit differently. So those are the design goals.

---

The non-goals are that ICANN determines the thresholds. That's not a goal. Whatever comes out of this if thresholds are going to come out has to come from the whole community. And that has to come out with as pretty much anyone in the community looking at the exact same data we are.

This is not about ICANN is going to set some thresholds and we need some data to do it. The idea is the data is out there for everybody, and the community itself would be the ones who would contribute to thresholds.

And the community means the technical people, the nontechnical people, everyone. So the other non-goal is that ICANN would then turn around and advocate for faster or slower scaling. Again, not an ICANN Org decision. This is meant to be for the community.

So going back a slide, that's exactly why this third bullet is so important, is that because this is not ICANN doing measurements so ICANN can decide something, this is ICANN possibly doing measurements so that the community can do things, and if the community wants different measurements, or even if someone just says, "Hey, I did it differently than ICANN, look at this," and the community goes, "Wow, they're doing it better than ICANN did," great. The result is more data for the community, if needed.

So let's do a bit of vocabulary here, because people often get a little bit wrapped up about this. And again, these three words are actually – you will see them in a hopefully soon to be upcoming document from RSSAC, because RSSAC is doing its own measurements. So for those of

---

you who are following the root server system evolution, there's a document called RSSAC 037 which talks about how will we have a system for determining how big the root server system should be, how do we add root server operators, how do we remove root server operators, all of that. And that's going to be – or some of that will be metrics-driven, and so these three terms are being used in the document that RSSAC is producing for metrics.

The first one is measurement. So when you go and you probe something or you get something without evaluating at all, just doing that, that's a measurement. Measurements turn into a lot of data. We're all familiar with that. But that is purely a measurement. A metric is taking a bunch of measurement and aggregating them. Usually over time, over location, something like that, and metrics, often people think of metrics as something like, oh, you just average it over a day.

That may or may not be valid for the kinds of data you're doing, so you want to have well defined mathematical processes that say this set of measurements will be aggregated in this fashion, maybe for a 50% median, maybe for a 95th percentile over a certain amount of time [and states] and it will come out with one or more values.

So really, when we're talking about a system that might be used for figuring out the scaling of the root server system, we're much more concerned about the metrics than the measurements. The measurements are easy, you go out, you do probes or you ask people to do probes for you.



---

It's the metrics that will be the most informative to the community. And at that point, you might want to set a threshold. You might say if this set of metrics is trending this way or hits this value, or doesn't hit this value after a certain amount of time, then that's an indicator, that is a threshold, that we need to do something.

And for this exercise, it would be we need to scale back how many TLDs are going into the root zone at a certain speed. You may never get to that threshold. You don't have to set up the threshold initially. You won't really want to set up the threshold until you actually have some metrics in hand to get a feeling for it. But you might say at some point after looking at this data we realize that here is a threshold, this means something now, we're going to act on it.

So it really is – thresholds are the things on which you act. You don't act on metrics. You collect metrics, you combine them together, and then you set a threshold. Hopefully, that's clear. And I know these uses of words are different than in some other things. Some researchers only talk about metrics even though they're really talking about measurements and such.

But for something like this where we're going to do some measurements, we're going to collect some data, we're going to aggregate them because it's a lot of data, and then we're going to make a decision. These might be good words. And again, it's the words that the RSSAC folks are using.

So here's what I want from you folks. I have a few more slides but if you as an individual were going to set up a set of measurements that would

---

turn into metrics to help a decision on is the root zone scaling too quickly, what kind of measurements would you want? What kind of metrics would you get from those measurements?

And this is a hard question because in fact, we've been thinking about this for a while, and some of these things might indicate too fast scaling, or they might not. They might just say, "Wow, look, the root zone is doing just fine."

So, what type of measurements might be useful to you? And to be clear, and the second bullet is sort of a sad one, but maybe none of these. Maybe none of this would be valuable for these kinds of decisions. If we knew what would be valuable, I would be sitting up here saying "Here's how we're going to measure." Maybe we would add to it, but we would say that.

At this point, there aren't any good measurements or metrics that we know of that are well agreed to. It may be that the community can come up with some, and I have some examples of some of the measurements and metrics we might take on a future slide, and you might say, "Yes, that's a great one." But we don't know at this point any measurements and metrics that are clearly usable for this task.

On the other hand, maybe we should start collecting because we won't know until we collect them. We're in a funny state here. And again, remember SSAC asked this to start happening before 2012. Didn't really happen, at least from [inaudible], it wasn't that we set it up. But so far, we've dodged a bullet or maybe there is no bullet to dodge. RSSAC set up their own system, they're doing their thing, the RSSAC 002 is a set of

---

measurements that they're collecting. They don't actually have specific metrics around those, but they've been collecting those, those are available.

So really, if the cost of doing this is low, maybe we should do it anyways. It might turn into other useful data, but the purpose of this talk is to say what are the measurements and metrics you think would be useful for scaling the root server system.

So here's some of the early ideas we have, and I'm just putting these up not to say this is the whole list. It certainly isn't, but to give you some ideas of what we've been thinking about so far.

So publication latency, imagine a root server instance that is somehow being negatively affected by the scaling issue. It might be running slow on doing a normal thing like when a new root zone is given to it, being able to respond to it. It might just be like so overwhelmed.

So you could tell if one of – that if it's publishing old root zone data by an hour or by a day or by a week, that could be an indicator that it's been overwhelmed. Another one would be DNS errors that shouldn't happen. For any legitimate query, you should not be getting serve, fail or refused. And if you are, those are sometimes indications that the authoritative server's about to fall over or is in the process of falling over, or is in the process of having fallen over and coming up.

So that might be an indicator that something has happened systemically for that one. We're not expecting the whole root server

---

system to have this, but if we start seeing this in instances, maybe that's an indicator. A root server that just doesn't have response consistently.

We all know that root server- instances are sometimes unavailable for reasons such as they are rebooting or the routing has changed to them, so it's really nothing to do with it. But a consistent no response might be that kind of measurement that would help you. Or a partial response. It is so busy it's not even sending the whole DNS answer out. That could be an indicator of something.

And what we're asking is, are there other externally visible measurements that might be useful here for this purpose? For those of you in the research community, you know there's a million fun measurements we can make, and we can write interesting research about it. But for this, which ones might be useful?

And it's not only externally visible measurements. Some of them might be self-reported. We might, as a community, ask the root server operators to report on things that would be useful for this. A thought would be memory utilization for the root zone. If the root zone is scaling too quickly and the root zone operators are not able to update the size of the RAM for example in their root server instances, and you're at 98% for a while, that could indicate, oh, we really shouldn't be pushing it past that until you fix it. Again, there's not a question that the root server operators will fix it. They've been extremely good about that. But maybe that's a self-reported indicator.

Having said that, boy is that an interesting number. For those of you who run servers of any sort, trying to measure that accurately over time

---

seems difficult. Right now, in RSSAC 002, there is a metric called load time which says how long does it take to load this new root zone? If that starts stretching over time, maybe that's an indication, again, of [inaudible]. And that would probably, the RSSAC 002 metrics, are probably more accurate than any external view. They're self-reported, but they are self-reported fairly consistently. Are there other self-reported measurements that you would want to see?

This is my last slide in order to prime the mic line hopefully. Just as a note, we are not doing this because we see an imminent threat. There's no indication that any of these are going to be useful for determining root zone scaling issues, but there's no strong belief that they could not be.

And again, why are we doing this? SSAC has asked us to do this. It's the kind of thing where it would be good to respond. It may be that our response is no, there isn't community agreement on what we should be measuring, so let's not start measuring things that are silly.

[There might be like, yeah, we should.] We would like to get a feeling from the community about how you feel about setting this kind of thing up. So really, the last bullet is the question, what might be useful to you all? And this really is the last slide other than the obligatory thank you. So I'm here for a mic line.

---

YOSHIRO YONEYA: Hi. So .jp already has this kind of early warning system for our zone. When we find the zone size increases or decreases around 20% [prior zones,] then we [stop] zone transfer and start manual inspection.

This is for the stability of our DNS system. So this is not advice, but this is an example of this early warning system.

PAUL HOFFMAN: Is to look at the size of change from one version to the next?

YOSHIRO YONEYA: Yes.

PAUL HOFFMAN: Okay. And since you said you sort of throttle or take a careful look, have you found any problems in those ones that were having that increase?

YOSHIRO YONEYA: I remember there was [the system area] that produced very small number of data from the registration system to the zone, so it causes name registration failure, so we stopped it.

PAUL HOFFMAN: Great. Thank you. That's very useful.

[BRUCE TONKIN:] Thanks for that presentation. I'm wondering what you consider to be the whole system. Are you including the manual parts of it in terms of

---

the IANA or PTI piece, or are you just purely looking at assuming that the zone file's been approved and just punching it out?

PAUL HOFFMAN: The latter, that is that IANA has its own ways of looking at things and such like that, so that's outside of the scope for this kind of early warning system

EBERHARD LISSE: We nowadays call IANA the IANA function operator.

PAUL HOFFMAN: So I work at ICANN and when we have to name it, we usually go through all of them because we figure out of the four possible names, one of them will have to be right. So we call it IANA, the IANA function, PTI, and the thing over there.

EBERHARD LISSE: Thank you. The correct terminology, as you said, vocabulary, is PTI. Sorry, it's the IANA function operator which is currently operated by PTI. It is helpful to always standardize on the same terminology even if it's internally ICANN.

PAUL HOFFMAN: Yes.

---

EBERHARD LISSE: Any other questions?

PAUL HOFFMAN: Any other suggestions, please, as we are really looking for community input here so we can have an idea on how we can move forward, whether we can move forward. I heard there's a remote question.

[KIMBERLY CARLSON:] This is from Hugo Salgado. Would it be useful to have some metrics on TCP AXFR of the root zone for RFC 7706, local copies on loopback?

PAUL HOFFMAN: Possibly. Thank you, Hugo. That's an interesting idea which I didn't think of, which is sort of funny because I'm the co-author on RFC 7706. So really the question is the actual size of the whole zone being transferred. A great increase in size or something like that. So the size of the root zone could be one thing like that, because people are using RFC7706 and things like that for loading it down.

So most of this presentation has just been about scaling the root zone for the root zone operators themselves, but if people are also using the root zone with RFC7706, keeping a local copy in their resolver, maybe we need to think about them for scaling as well. So yes, thank you, that's a good contribution.



---

WES HARDAKER:

Wes Hardaker, USC ISI, where I actually have, as Paul knows, a project to mirror multiple zones, actually, not just the root zone, using the RFC 7706 project. And certainly, one of the things I'll be doing in the next couple of months – and for those of you that are local root subscribers, you'll be getting notes from me about new features coming out. One of the things I want to do is actually add more zones that you can mirror locally if you think that they're important, and one of the biggest attributes of that is, what is the size requirements of my memory if I'm going to pull more and more stuff?

Certainly, you can't pull .com, nor would they let you, but there's a lot of ccTLDs that have indicated interest in doing something like that too, not just the root but making sure that their ccTLD is always available. So I guess with that, I would certainly love to see the ability to measure zone sizes, both just the file transfer for bandwidth, and the memory size required to run it, I think as you mentioned it too, Paul.

So two takeaways. If anybody is interested in participating in my project, you can go to [localroot.isi.edu](http://localroot.isi.edu), and you can sign up and it gives you free configuration. We'll be getting a configuration shortly for most popular resolvers. Right now, it's just Bind.

I had a second point. At some point I can give you stats on who's running it and how often they're pulling, so I have those stats available to me. They're not published on a webpage yet, but they will be, and I suspect that might actually feed some of your work. I don't know. And if you're a ccTLD operator and want to participate in this, please talk to

---

me. I'd love to – it's one of my goals this week to find a few people to throw extras [inaudible].

PAUL HOFFMAN: Wes, before you go away, so do you have some ccTLD operators now that you have a feeling for the size of those zones?

WES HARDAKER: I do not now. There are a number of places where you can – there's a GitHub repository where you can download the current zones for a lot of ccTLDs, and I don't want to do it without permission, so I'd really prefer to actually have a one-on-one. Even if your data is already available, I'm not going to do it unless I actually talk to you. But there's a lot of publicly available zones. I don't have sizing requirements for all of them. Most ccTLDs are fairly small. There's certainly a couple that are not.

PAUL HOFFMAN: I understand, but I was asking more in terms of the root zone is also fairly small. If you had said “And I already have this one which is ten times larger than the current root zone,” that might be a good comparison and things like that. So we'll talk. We often talk, but we'll talk about that because that might also be a good way of figuring out sizes that might overwhelm someone.

---

WES HARDAKER: Yeah. One other thing you might consider measuring is the traffic requirements of something like 7706, or aggressive [NSEC,] Both of those prevent negative answers from leaking out of your local resolver, which turns out to be huge. I gave a presentation last week at DNS OARC that included a graph showing how much of Chrome junk queries leak to the room and how much faster those queries would return, 0 milliseconds versus 100 or whatever in order to get a fully resolved negative answer, regardless of whether it's a typo because you typed two Ms in .com, or whether it's Chrome generating junk queries. So you can actually stop a lot of them right at your border and it's a big bandwidth savings.

PAUL HOFFMAN: Thank you. Okay, thanks very much.

EBERHARD LISSE: Okay, we are getting a little bit ahead of ourselves. Not a problem. Maarten Wullink will now talk about ENTRADA.

MAARTEN WULLINK: Hello, everybody. I'm Maarten from SIDN, and I'm going to give a short talk about ENTRADA and some of the recent work we've done with ENTRADA.

First, probably a lot of people have no idea what ENTRADA is. It's a tool we developed at SIDN to help us with analyzing large amounts of batch of DNS data.

---

We've been using it for about four years at SIDN, and we process about 1.5 billion DNS queries from our Anycast cloud, and the tool helps us analyze this data as well. In this talk, I'll first give a short introduction about what ENTRADA actually is, and then I'll show a little bit about the new features we developed in the last few months, and finally how you can deploy ENTRADA for yourselves if you're interested.

First, anybody interested in doing some analysis of DNS data probably will start with just a TCP dump on a nameserver for instance and capture some DNS data, maybe create a file every few minutes, maybe a file containing five, ten minutes of data, will give you a couple hundred megabytes of a file depending on compression and so on.

And then there are nice tools such as Wireshark and tshark that you can use to analyze the PCAP file. Wireshark is an awesome tool if you only have one or two files, because it needs to read the entire file into memory and when you do stuff like searching for IP addresses or other fields, it needs to crawl through the entire file each time and it takes a long time.

So in our case, when we wanted to do more than just a couple of PCAP files, we wanted to store a large portion of the data we collect, and also be able to analyze it, we had to develop a tool to help us do that. What we wanted to have from this tool was something that would give us reasonably good performance. Doesn't have to be a tool that gives us when we ask it something a response within milliseconds, because that will probably – we're not Google so this will cost millions of dollars.

But depending on the question, an answer within tens of seconds or a couple minutes, sometimes 15 minutes or more, is perfectly acceptable. We do want it to be highly available, so data put into the system cannot be lost, so if one part of the system fails, we do not want to lose data that's stored in that part of the system so it has to be redundant. And also, we would like data that is collected on the nameservers to become available in the ENTRADA data warehouse pretty soon, pretty shortly after it has been received.

So usually, we receive the data in ENTRADA within some 10 to 15 minutes, which is pretty okay for us, and also very important, we need to have SQL support because people working with a system like this, such as maybe DNS operators or researchers, they don't want to program against some complicated API, they just want to do SQL queries because everybody in IT understands SQL, has been working with SQL for a long time, so this will save a lot of time.

ENTRADA has a couple of features. The main function of ENTRADA is to convert DNS data from PCAP format to different format, which is Apache Parquet format, which is a [inaudible] format, PCAP is so difficult to analyze. As I mentioned before, you need to read the entire file into memory, you cannot just jump to a particular location in a file. And Parquet gives us some advantages that make it very attractive to do so.

And also, it adds some extra data to the data we extract from the PCAP file, such as geo information, geolocation, or ASN numbers, that optimizes the data format, and also all the steps involved into loading

---

the data into ENTRADA, copying it to the correct output location, making it available for analysis. It's all automated so person working with ENTRADA doesn't need to worry about how and where the data is and all the steps involved, just focus on analyzing the data.

So a little bit about Apache Parquet, which is what they call a column storage format, developed by I believe Twitter and Cloudera. So Parquet is what they call a column storage format, and normally, in traditional databases, they use a row-oriented format, and the main difference is that when you have a table with data and you have a row-oriented format, then all the column values for each row are stored sequentially.

So as you can see in the picture below on the slide for row oriented, you have column A1, B1, then C1, and so forth. And when you look at the column oriented data, then you see that all the values for a particular column are stored sequentially.

What this means is that you can use efficient encoding and compression schemes to compress the column values much more efficiently than when all the column values are intermingled with each other. So you can compress all the values for column A because these are probably similar values.

They might be like, if you look at for instance the status codes for a DNS request, an RCODE, usually it'll just be zero or three. So there's a limited set of possible values there, and compression schema can compress that very efficiently. So the data size of a Parquet file will be much lower than the original PCAP file.

---

Also, if you want to analyze the data in Parquet format, you only have to read the columns you're interested in from disc, so if I only want to analyze the data in column B for instance, then I can read just the data from column B from disc. But if this data was in a row-oriented format, I would have to read all the rows from the entire table, so also have to read column A, B and C. So it means there's lots more IO, potentially much more network IO, makes everything much slower.

So this is a short overview of the workflow, steps involved. The top left, we see there is a nameserver receiving queries from resolvers, running a TCP dump script, creating PCAPs every maybe ten minutes, and these PCAPs are then sent to a staging server from where ENTRADA can pick them up and it will decode the PCAP files.

What it then will do is look for queries and response pairs so it can combine the query with a response for that query. That means that later on when you want to analyze the data, you don't have to join a table containing requests and a table that contains responses, because that's probably going to be two huge tables. You don't want to join that, so we do that up front already so it makes it much more efficient for later on to query data.

We do some filtering because at least for our situation, we're not interested in all the data, so we skip out or we filter out some part of the response. We keep most of the query, all the status flags and stuff like that. There's some enrichment. As I mentioned before, we add geolocation information, only up to the country level so we don't have city information, and ASN information so that later on, you can in your

SQL queries look up data for particular countries or particular ASNs without having to look up the ASN or the country at that time, which would make the query much slower.

And finally, the data is written out as a Parquet file to a storage system which can be a Hadoop-based storage system or it can be S3 or AWS. And from that point on, it's available for analysis through a query engine such as Apache Impala or AWS Athena. And on top of that, you can build your own applications because these engines have APIs, web interfaces. And usually, at least in our situation, data is available for analysis within 10-15 minutes.

There are some challenges while running ENTRADA, especially in previous versions, because ENTRADA was before based on Hadoop, and Hadoop is a great tool, but unfortunately it can be a complex beast to control. It's a very big ecosystem with lots of packages and need to have a lot of knowledge to install Hadoop and also to maintain Hadoop, especially when you grow out your cluster, you constantly have to monitor it and things break. So you have to be on top of it even though it is built to withstand hardware failures, you can lose discs, you can lose entire nodes and it'll perfectly work, or it can perfectly work on after that because it has redundancy built in. Every bit that you store on Hadoop is replicated by default three times, so you can lose, depending on the size of your cluster, quite a few discs or nodes without losing actual bytes.

And another downside of using Hadoop is that you need hardware, of course, or you can use Hadoop in the cloud for instance as a virtual



---

cluster, but then you need to spin up a cluster and configure it, which can also be quite a lot of work.

So for ENTRADA 2, we decided to add a few new features. One of them is serverless DNS analytics, which means that you essentially do not need to have any hardware anymore, or you don't need to configure anything in the cloud as in virtual machines to run ENTRADA.

Added support for multiple SQL engines, and also added something we call quality of service monitor that is based on round trip times – RTT – which was also influenced by the excellent work done at the ZCNIC which also presented this morning. And we Dockerized the whole thing, so it's very easy to deploy.

First, serverless computing very quickly, it's basically a cloud model that says, okay, I want to run my stuff and I want to run it in the cloud, but I don't want to be responsible for having to configure lots of virtual machines and network configurations. I just want to put my data somewhere and be able to analyze it without having to do all this extra work, otherwise I could maybe just set up the hardware and configure it myself all locally.

So with serverless DNS analytics in ENTRADA, you don't need any servers, you don't need any hardware, you don't need to set up any network, so there is no additional hardware or networking costs. You just pay for the amount of data that you want to store on S3 or the amount of data that you analyze. So it basically means that you are free to focus more on analyzing the actual data that you collect than that you're busy with keeping the cluster up and running and making sure

---

that everything works fine and making your users happy. So that will save you a lot of time, especially if you're just starting out and just want to play around with it.

So ENTRADA will, if you use the serverless mode, create all the stuff you need, it'll create database schemas, it'll convert and upload the data, it will optimize the data. So everything will be done for you.

Currently, we only have one cloud provider that we offer support for, that's AWS. AWS has S3 for data storage and Athena which is a SQL-compatible query engine that you can use to analyze the data that is stored on S3.

And if you look at pricing, then depending on where, what datacenter location you want to store your data, so in Europe, in the US, then there are changes in pricing. If you look at for instance Ireland, if you want to store a terabyte of Parquet data on S3 there – and there are also different types of storage classes, so if you use the normal storage class, it will be about \$13 per month for each terabyte of data that you store, which is not that bad. it's not really super expensive. A terabyte of Parquet data is a lot of DNS data.

But that's basically just the storage part. If you want to analyze the data, which you all want to do, of course, then you need to use Athena, and Athena is as I said the SQL engine, and the way it works is that you point Athena to where you have the data stored in S3, which is automatically done by Entrada, and when you analyze the data, Athena has some accounting that counts up all the bytes that you read. So in

---

the end, you have to pay \$5 per terabyte of scanned data, which depending on how you use the data, can be good or not good.

If you only use a few fields of all the data fields that are available and do not scan all of the data all the time, then it can be very cost effective. But if you scan all the data and all the columns all the time, then this can also get quite expensive very quickly. So it all depends on the use case that you have and how much data you want to process.

This is a quick screenshot of the Athena web interface. It's basically just a nice interface where you can create SQL queries and get a response. It's ideal for some ad hoc queries, and there's also an API that if you want to program your scripts or your applications against Amazon Athena, then that is also possible.

So this is one of the new features we added, what we call the quality of service monitoring. What it does is it looks at TCP requests that are in the DNS data and looks at the TCP handshake, and basically calculates the round trip time from the TCP handshake.

The interesting part here is that, as also I think mentioned this morning, these are actual resolvers that you are using. These are not for instance RIPE Atlas probes, which are awesome. RIPE Atlas is an awesome project, but they tend to be not very representative sometimes of the actual resolver population, so there might be differences between RTTs calculated from sets of RIPE Atlas probes and actual resolvers that use TCP for your DNS queries.

---

So high RTT can be caused for instance by inefficient routing between resolvers and your nameserver infrastructure, or congestion on the network links, or some router or switching issue or infinite other options.

So this is a quick image of how this works, so I'm not sure how many people are familiar with TCP and the handshake, but basically, you have a client and a server, and if the client wants – in this case client is a DNS resolver, and if the client wants to talk to the authoritative nameserver with TCP, it will try to set up a connection and the connection is done in a three-way handshake, so the client will first send a syn packet, the server will respond with a syn ack packet, and finally, the client will acknowledge the syn ack with a single ack.

Well, what we do is we look at the time difference between when the server sends the syn ack and when the server receives the ack from the client back again, which will give us the RTT, the round trip time between the server and the client. And the nice thing of doing this with handshake is also that there is no data sent back in handshake, so usually it doesn't affect the handshake that much.

And we looked at how often is the TCP used for the nl-zone, because if it's only used in a very small number of cases, then it's not enough to create a monitoring system for instance, because with only one or two TCP clients per day, it doesn't say that much.

So what we did see when we looked at data for last month, I believe the first 25 days of October, that around 4-5% of all queries were sent using TCP and this data was sent by some 20-26% of all unique resolvers.

---

I don't know if this is readable in the back, but on the X axis we can see then number of days in October until the 25th, and on the Y axis is the percentage of TCP queries. So we can see our percentage of queries sent over TCP. We can see it's pretty stable, with some peak there around October 12, but usually it's around 4-5%.

When we look at the number of unique resolvers sending TCP, then we see basically the same thing. On the Y axis, we see the percentage of unique TCP resolvers, so we see that around 20-26% of resolvers send us regular TCP queries.

This is an example how you can use this TCP RTT data. This is an example SQL query in ENTRADA that just tries to get the ASN numbers with the highest RTT for October 24th in this year.

This is the top ASNs with the highest RTT for that day, and what we see is that most of them are from China, actually one from Russian, and they are pretty high but the top one is almost three times as high as the second one, and it is almost 1200 milliseconds which is pretty high.

What we can do is we can zoom in on this ASN and we can just look at the traffic that we see for this ASN and get all the RTTs for the same period. What we see then is that usually, the RTT is pretty normal. It's not super low, but it's around, I believe, 100-150 milliseconds, which is not super fast but also not horrible.

But then all of a sudden, around October 20th, we can see that the blue line is the average RTT on that day, and the brown line is the number of samples received from this ASN. What we see is that until October 20th,

---

everything is fine, and then suddenly, the RTT spikes. It's much higher. And at the same time, we see that we receive much more samples from the resolver, many more TCP queries from the resolver.

This is interesting to see. Unfortunately, I have no idea why this is happening, because it's very hard to find out why this is happening, but we can for instance, if this were to be an important ASN for us, maybe get into contact with the operator and say, hey guys, is there some problem or can we help? Maybe we can help you fix the connection and make the RTTs normal again, or at least not as bad as they are at the moment.

Another thing that we have for ENTRADA is a couple of example dashboards. This is a screenshot of the RTT monitoring dashboard. It just shows the number of samples that we receive. This is for one of the nameservers, but you can also create graphs per ASN, per country or whatever you want. So this is an ideal tool for for instance operators, DNS ops people to monitor the quality of their DNS service.

As I said, using the Hadoop version was sometimes challenging. Having the serverless version with Amazon S3 support makes it a lot better, but to make it even easier to install ENTRADA, we decided to use Docker as well. With Docker, you can just pull the Docker image and choose which mode you want to use. You want to use Hadoop inhouse or you want to have a cloud-based Hadoop, or if you want to use Amazon AWS S3 plus Athena, you can just configure that in the Docker deployment scripts.

So maybe not everybody here is familiar with Docker. It's basically a lightweight container system. It's not a virtual machine, but it's kind of

---

a lightweight virtual machine is maybe a good way to describe it, and it allows you to package an application and all its dependencies in a single image which can then be uploaded to a central repository where users, people who want to use your software can then download this image which describes how the software works and what its dependencies are, and then install this image locally and run it, and when you run it, they call it a container.

So it's a lightweight container, you pull it from a public repository, and configuring applications is pretty easy. If you use the tool they call Docker Compose, you can have multiple variables to configure your applications, for instance for ENTRADA we have variables to point to a location where the data should be picked up from, locations where the data should be saved to, is it a local system, is it a Hadoop system, is it an AWS system? And also if you want to use encryption. There are dozens of options. All documentation is on this website, [entrada.sidnlabs.nl](http://entrada.sidnlabs.nl) where you can find all the configuration options.

And finally, a little bit about our use of ENTRADA, we've been using it for about four years, hooked up some – I believe at the moment 28 Anycast sites for our nameservers, and database contains some 1.2 trillion DNS records. These are DNS queries combined with a DNS response, and if you store it, then it's about 65 terabytes, but we store it in Hadoop and there's a replication factor of 3, so we need about 195 terabytes to store all this data on Hadoop, which is still, if you compare it with the original PCAP volume isn't too bad, because that was an order of magnitude more.

---

Okay, so if anybody has any questions, I'd be happy to answer those.

EBERHARD LISSE:

I have one, abusing the privilege of the chair. This is not really suitable for smaller ccTLDs? Because if I understand you correctly, it's either Hadoop, which I can't understand, or it is AWS, which is too expensive. Is there a way of putting it into MySQL or PostgreSQL and run it for a smaller place where the data volume is not so big? I have, I think I presented in Singapore about using SQL-type solution from Sweden to access PCAP and then put that into a database. If we have a system that takes the grabbing of the data and putting it into an SQL base, if that was quite easy, that would be quite helpful to us.

MAARTEN WULLINK:

Yes. I know what system you mean, I just can't – the name is on the tip of my tongue. But to get back to your point about AWS being too expensive for smaller ccTLDs, I think it's an ideal situation for smaller ccTLDs, especially if you don't have a huge volume of data, then AWS can be very cost effective, and it allows you also not to have to configure all this stuff locally which saves you a lot of time and money of course. You don't need hardware, you don't need to buy servers.

So a ccTLD probably doesn't have money or resources to set up all this stuff, so AWS can be a quick way to get up and running with DNS data analysis, and you don't have to run it all the time, you can just – if you have a project, you can set it up and use AWS, analyze some data, and



---

then delete the data again and you're done. So there are multiple options here.

EBERHARD LISSE: Okay, Bill.

[BILL:] Yeah, hi. My name is Bill from .ca. A question, when you're running the Docker instance, that's running on an on premise system and then you upload the Parquet files into S3?

MAARTEN WULLINK: Yes, that's what we do. We have a local system that runs Docker containers, but yeah, you can also choose to have that system running in the cloud of course, on a virtual machine, so that's up to you.

[BILL:] Thanks.

EBERHARD LISSE: Any other questions? Thank you very much.

MAARTEN WULLINK: Thank you very much.

EBERHARD LISSE: Okay, now we will hear about the ICANN security practices.

---

**ASHWIN RANGAN:** Thank you very much. My name is Ashwin Rangan. I serve as the head of IT for ICANN Org. And to my right is my colleague, David Closson, he's the head of infrastructure operations. I thought it would be appropriate for him to talk about what we do from a security perspective. He's very familiar. I'm here to answer questions, and in the unlikely event that he can't answer a question, I'll step in and get into trouble. Go ahead, David.

**DAVID CLOSSON:** Thanks, everybody. David Closson, I'm the director of IT infrastructure at ICANN. We've put together a slide deck in cooperation with our infosec team that we're going to share with you today. So you want to advance the first slide.

This is basically the three-step holistic approach to information security that we follow. Confidentiality, integrity, and authenticated access. So bad actors have the right to act just once whereas we have the right every day and all the time to take on these things. I think that's the gist of this slide here.

Acute awareness of likely entry points for bad actors, so obviously we have three main categories, hardware, software and then third parties. The hardware of course being storage, processing, network infrastructure, typical infrastructure things, and then software of course is cloud, which could also be part of the third party, and then we have operating systems, databases, applications, development and

---

CICD environments. For those of you who don't know what CIDC means, continuous integration and continuous delivery. It sort of automates the development pipeline, there's tools and things that do that.

And then third parties with whom ICANN Org has business relationships, this is a tough one because there's a lot of unknowns there. Sometimes you don't know what to ask them, what their infrastructure is made of, you don't know what they're going to be willing to tell you or not. So it makes it more difficult sometimes with third parties. It really depends on how much information they provide and how much you ask. So on a path to address these holistically.

The next set of slides is about what we delivered in the last six months, but in general, the infosec team is on a good trendline. Starting with the top circle there, we start with readiness, we move into Triage, analysis and mitigation, remediation, communication, lesson learned, and of course, this goes in a circle. And what we're showing here is the computer security incident response team's lifecycle at ICANN. This is something we're already doing, and I think for most of you, hopefully this makes sense, what you're seeing here. We can get to the Q&A later if you have any questions. I think it's pretty straightforward.

ICANN infosec already delivered, the HackerOne program which established internal SLAs for vulnerabilities, so basically we provide a reward system for those for hackers that come and show us vulnerabilities, so they can get credit for that. Is there anything else they get for that?

---

ASHWIN RANGAN:

No, we deliberately stayed away from any cash-oriented reward system and stayed with the concept of providing them with attribution and recognition, so that's kind of where we are at. We're also making sure that we're opening up brief windows so that we have enough time to react to whatever has been reported back to us. But it's on a very predictable basis that we open up these windows for researchers to look at our assets.

DAVID CLOSSON:

Thank you, Ash. And as you can see, since we've first started this program, we had quite a surge of reports come through. It's kind of tapered down at this point, we're seeing a steady [inaudible] which is good. And the table at the bottom is showing you basically our suggested response time for the various types of vulnerabilities, critical with high data risk being zero hours and so on. And it gets further away from the risk, it becomes longer time, we allow ourselves to remediate these things.

Another already delivered, red team exercise at ICANN meetings. The last one was in Kobe, 13 infosec recommendations presented to right actors for mediation, 11 recommendations adopted and either implemented or on the improvement roadmap. Two are yet to be planned. Those are things that are not like security fixes so to speak, but things that we should do, behavior change.

Meeting with other network admins who run similar conference networks so that we can strategize with them about how to do certain

---

things, onsite intrusion detection system and packet capture, expanding on that from what we currently have.

This is another already delivered, the infosec ambassador program. Representatives from all significant functions at ICANN Org, much better Org engagement with infosec, a culture shift. That's important. Regular, positive and topical e-mail discussions such as the following. Quite a wide variety here from password managers to Have I been Pwned to the five recent industry infosec compromises, NAVEX training which is our HR-mandated – secure infosec awareness training, but we also use that platform for other types of training internally, digital e-mail signatures and so on.

The next one is created ICANN Org's first ever CMDB, which is a configuration management database. This is the place where we can warehouse all information about an application so that we can catalog software updates and other security bits and pieces that are helpful to understand the landscape, and then secure coding training for developers, the Checkmarx Codebashing, all programming languages used at ICANN are included, security basics, hypertext transfer protocol security principles or HTTP, OWASP top ten for each language, and then the 35 team members enrolled, 33 of them from our development and two from infosec. So it's quite an accomplishment there, good engagement.

And then looking forward is the next two quarters, what we're doing, FY20 Q1 and FY20 Q2. Automated vulnerability auditing of all the ICANN-owned assets using Nexpose which is one of our platforms,

---

reports to the sent directly to service owners and executive relationship managers for action in an automated fashion, mostly through e-mail.

Automated common vulnerabilities and exposures assignment to service owners, or CVE, using VulDB based on our catalog of services and common platform enumeration from the ICANN CMDB.

FY20 Q1 and Q2, by country, matrix of infosec risk for ICANN staff, matched infosec recommendations and controls for traveling staff, an important one since ICANN staff travel quite a bit. Scoping and profile definition of NIST cybersecurity framework or CSF, focus of ICANN on secure processes, bullet point one, maturing infosec in ICANN as an organization compared to the CIS20-style controls, which is what we previously used.

For those of you that may be – what does it stand for? Computer Information Society’s 20 top controls?

ASHWIN RANGAN: Cybersecurity information.

DAVID CLOSSON: Yeah. Sorry about that. That’s what we used to use and now we’re using NIST. FY20 Q1 and Q2, centralized log collection, indexing, alerting and dashboarding. We use Splunk for this. I'm sure some of you have heard of Splunk, that makes sense to you, I'm sure.

Expansion of traffic capture infrastructure, getting more TAPs into all segments of the networks, not just the main chokepoints so that you

---

can see behind things that you might not have been able to see before. Sufficient storage for all this stuff, take a lot of storage up. More processing power and more quickly search for event-related information. Obviously, if you have all this data and you need to look through it quickly because you've detected some sort of an event, that can take a lot of compute power to crawl through all those packets.

Forward thinking. Future work here, review of the HackerOne, consider financial rewards. That's what we discussed earlier. I think there's lots of vendors out there that do offer financial reward with good results. It's possible that we would switch to that as well and put bounties up. That's worked out quite well for other folks. Maybe we will get more traction.

Consider back of office ability to meet SLAs. Red team exercise on all ICANN networks, not just the meeting networks. I think that's important. Internal from within, potentially third-party vendor, so that means we'll do both. Consider turning it into an ongoing program, checks and balances, costs, tradeoffs. Takes time to do all this stuff, it can interfere with normal work, and it's not like we have a ton of staff. We're hundreds, not thousands of staff.

Browser isolation evaluation is the next one. It's important. In-hand with browser selection and trust anchor store and browser extension audits. A browser nowadays is where everything happens, so if you can actually isolate the browser, you can add a lot of security to client endpoints.

---

End user device security review, are we doing the best we can? I'm not sure that anyone can ever say that. It's always trying to catch up so to speak. Are there better, easier tech options for two-factor authentication available? As you all know, using two-factor authentication for a lot of folks is kind of a burden, but it's a necessary evil these days.

It's nice when new options exist, like recently we deployed Duo. So it makes it a little bit easier to do second factor this way, more convenient, but still, it's quite secure, as an example.

Future work, most NIST CSF, assess, analyze gaps, action plans, reassess. Infosec strategy day, this is an internal event with infosec ambassadors. Those would be staff members who have been coopted in to participate who probably have knowledge of the applications. Solid engagement with ICANN departments. Glibly, bring out your dead. What are the processes that are secure naïve? Business processes is what we're talking about here. Some of them are in that category of things that we can improve there.

SSR RT2 findings and recommendations. I'm sure there'll be a number of items. I'm not sure what they are yet, so we'll watch out for this. ICANN posture on third-party hosted platforms. That's a big one. More and more things move to the cloud. And we're not talking about a single cloud like AWS or Google, we're talking about many different clouds because you choose the best cloud provider for the purpose in many cases as necessary.



---

Review and influence they have, what more would we like, and access to logs is interesting. Often with these cloud providers, we don't have access to raw logs. Sometimes you do and sometimes you don't, but it seems like there's never enough.

Okay, and that is the end of the slide deck. Any questions or comments, I'm sorry I had to get through that pretty quickly. Maybe I went too fast for some of you, but I'll be happy to back up if I went through something that you have a question about.

**YOSHIRO YONEYA:** Just before this ICANN meeting, I received that the ICANN provide SSO, single sign on, system to the participants. Is that the scope of this information security system?

**DAVID CLOSSON:** I think your question was about the single sign on for meeting participation, through [accounts.icann.org](https://accounts.icann.org). Is that correct?

**YOSHIRO YONEYA:** Yeah. So my additional question is, as you said, why don't you introduce two-factor authentication for the single sign on system?

**DAVID CLOSSON:** Correct. Yes, okay, got it. Thank you for that. I agree to that, and I think that is something we're going to look into doing, because the platform natively supports it. I think the only reason we don't is because of

---

community and perhaps some folks not being ready for that or don't feel like the things that it allows you to have access to shouldn't require that level of security to be inconvenienced. But we agree with you, that should be done.

Any other questions? Alright, there we go. Come on.

UNIDENTIFIED MALE: Hello. [inaudible] for the record from Senegal. On the first or second slide, you're talking about your holistic approach of security. I didn't see anything related to availability. Is it on purpose?

ASHWIN RANGAN: This is just one means in which we look at vulnerabilities. I think what you're talking about is a different construct that we also use internally, we call it RASP for short, which is reliability, availability, scalability and performance. We treat that as a basic requirement for all the services that we stand up. That's a different track of work that we undertake, and it's independent from information security per se. So we treat information security as being something that results from having RASP as a foundational platform. Thank you.

JOHN LEVINE: Hi. I realize there are probably details here you'd rather not discuss, but I'm just wondering, when you look at the conference networks, do you see a lot of malicious traffic?

---

ASHWIN RANGAN: Could we have a drink together?

JOHN LEVINE: Sure.

ASHWIN RANGAN: Thank you.

DAVID CLOSSON: As with any network, the answer is yes. Is there a lot more than –

ASHWIN RANGAN: A lot becomes a matter of definition, so that's the reason I think having a private conversation probably will be able to characterize it. What may be a little to me may be a lot to you and vice versa. So thank you.

DAVID CLOSSON: Thank you for the question, thank all of you.

UNIDENTIFIED MALE: [inaudible] .dk. You showed some numbers of the amount of reports you got from the HackerOne campaigns. Can you say anything about your experience with the quality of those reports considering you're only doing attribution and you're also considering doing paid campaigns?

---

ASHWIN RANGAN:

Absolutely. Thank you for the question. When first we opened up HackerOne and we opened up our assets to the HackerOne program, we saw a surge – and it lasted for a short period of time – the quality of research was excellent. We were surprised by the quality of research and the specificity with which things were being pointed out.

In fact, some of them, we were able to just take them into JIRA, our repository for taking remediation action, and just convert them directly into stories and story points, attribution, and just move on forward.

Nothing much has changed as a matter of fact. What is surprising to me is the declining numbers. There is a clear downtrend here, and the question that we're asking of ourselves is whether offering a bounty would increase the incident reporting, whether this was just the bloom on the flower when it first was blossoming.

We don't know, and therefore we're starting to wonder what to do next, because I'm pretty sure that we're not 100% secure or something all of a sudden. So it's a question. It's hypotheses that we're playing around with right now. Thank you.

EBERHARD LISSE:

Okay, any other questions? Thank you very much. Shorter than advertised, but nevertheless very interesting. One never really thinks about the IT security of ICANN, one just thinks about the root servers and so on, but it's rather important that even every laptop everybody is using can't be used to go into head office and sort of install java scripts

---

into the browser as we heard this morning. Anyway, thank you very much.

Okay, now we hear what we've all been waiting for, the transition of the registry of .in.

SHUBHAM SARAN:

Good evening, everyone. I'm Shubham from NIXI. First of all, I would like to thank Dr. Eberhard and the tech day program committee for this slot. I know initially this was meant to be the last session, but now I'm not the last one. Throughout the day, we have seen a number of sessions which are purely technical in nature, like as you said, the machine learning, AI session which were there which definitely was kind of a pure technical session for us.

I'll try to be as brief as possible during my presentation. Agenda for the next few minutes is as follows. As I said, I represent national Internet exchange of India, which is NIXI. We are interested with the task of running the ccTLD of the country, .in, and we have been managing the .in since 2005. Throughout this period, we have been well supported by the TSP which was with us throughout this period until 28 February 2019 when Neustar took over.

We really appreciate the role and the support of our earlier TSP, which is Afilias India who [worked with us] during these terms which are mentioned here, the initial term of three years, then after that, two five-year terms.

---

During this period of [inaudible], they were selected through a rigorous RFP process. It's not that they have been there for by virtue of it. They were always selected for these terms for the rigorous RFP processes.

This RFP process was conducted through a consultant which was supervised by an empowered committee. During the last RFP process, they missed out and now we have a new TSP with us, which is Neustar. So just to clarify, it was purely RFP process through which the TSP has been serving .in registry.

Giving you an overview about the transition, this time also, we engaged a consultancy firm for the entire selection process, and this consultant was again selected through an RFP process. It was like we bring in the transparency for the entire selection process being backed by the government.

The key responsibilities of the consultant who was hired are listed here. Right from the project inception until the contract closure and the agreement signing, they were given the task and they were helping us, they were completely working in sync with us.

The objective of this entire RFP process was to make this entire bidding process a complete transparent and efficient process, and for this bidding, we used the government online central public procurement portal. The chosen consultant who worked with us during this bid process, he helped us during these three processes and stages of the bid process, which is the prequalification, technical evaluation, commercial evaluation, until the TSP selection.

---

Consultant firm also carried out the consultation workshop, the pre-bid consultancy seminars and workshops. During this workshop, five global organizations participated in it.

Some key things that we were looking in a new TSP are listed here, but I would like to focus on few. As most of you might be aware that India is at the cusp of the digital India growth, the Internet users in the country have exceeded 500 million Internet users, and India now contributes almost 12% of the world Internet population. So we wanted a partner who could not only support this growth but support it with the experience in marketing, channel management, and deliver all this in a secure and stable manner. So these are the main considerations while we were looking for the new TSP.

The key transition highlights, one of the points which I'm sure Dr. Eberhard has been mentioning since yesterday, how we have been working – like in terms of the TLDs which we support in India, we have .in, the ccTLD, and along with that, we have 15 IDNs in 15 more languages. So in total, 16 TLDs which [could] transition during this entire exercise.

NIXI was keen to keep its registrar partners informed of all these steps leading to the transition, and hence we carried out a number of seminars or webinars with all the partners, which is basically all our registrars. We had a base of almost 130 registrars.

So as a part of the process, as a part of the help, we created an interim website, [transitionsupport.in](http://transitionsupport.in) which had all the updates available there.

---

So it's not only the – if someone is missing out the webinar or any other events, he had that website where he can refer the progress.

We also extensively used social media and mailers to communicate with all our stakeholders, and we coordinated closely. This entire transition would not have been possible without involving the main two stakeholders apart from the registrars, the incumbent TSP and the incoming TSP.

So we coordinated closely with the incumbent and the Neustar team to ensure a smooth, seamless transition, and as a part of the transition, the world-class infrastructure has been put in place to ensure that .in remains the mainstay of digital India and is ready to support all the Internet-enabled services which work around the .in domain.

Just to give you a sense how the project was like, it was six months of hard work which was involving this entire migration. Right from the software development, the IDN implementation and DNS migration, some of the timelines are mentioned here in the slide.

If you observe, there was a plan, date and time were well defined. All the stakeholders were kept well informed about these timelines. The final transition started at the onset of 20 February 2019. That was the day when that was a transition day.

After all these stages and keeping all the stakeholders in loop, the registry was again opened on 1 March 2019. And during this entire transition, we were well-supported by IANA. The IANA team was



---

completely with us. They were proactive and they were helping us in the entire transition process.

So in terms of the infrastructure which got delivered after this transition, system was able to scale multiple times our present .in zone. We installed two fully redundant registry sites and with the full capacity at both the sites. Both the sites were exact replicas of each other.

And for the scaling is again the entire these two sites are very much scalable. It can support up to – the current base, it can go up to five times easily, we can scale up to that level.

The infrastructure which is being deployed is highly available and redundant by design. Complete redundancy has been built into the system. It's a highly available infrastructure with multi [vendors] and a multi redundancy system.

Again, coming to the world-class, just to reiterate that redundancy is built into this infrastructure, so we have primary and secondary sites, which is Mumbai and Delhi. Delhi is the capital, Mumbai is the commercial capital of the country.

Then both these DC sites are identical in hardware and software, deployed using our enterprise-grade OEMs. And being identical, it helps us in providing the complete redundancy.

The primary and secondary site are with independent datacenter providers, these are the neutral datacenter sites which are then connected where independent network providers with different network diversities. Again, from the redundancy perspective.

---

At both local and international level also, they have the redundancy how the network planning has been done, considering that there should be a fallback, a redundancy in place. While doing the planning, we ensured there's no single point of infrastructure failure.

[inaudible] representation of the sites, fully redundant network at all public and internal points. Again, making a point on that, like I mentioned, infrastructure components are N+1 capacity and redundancy. Further scaling is online and we can further – the entire infrastructure deployed is fully scalable.

Coming to disaster recovery which his being a part of the delivered infrastructure. While drafting the RFP, while at the building level only, we had the clear RPO/RTO targets which were to be adhered during the RFP [inaudible] the building process.

The entire infrastructure which has been delivered has got disaster recovery enabled in the system. We have a fully mature [inaudible] infrastructure as of today.

Primary and secondary sites are running as active and warm standby. As I said, we have primary and secondary. And the application servers which are working also auto connect to the active location. From the registrar point of view, there no extra configuration which is being required. It's a kind of active and warm standby.

What I meant to say, we have all the planned site failover also, which we have enabled, and the system is designed to failover quickly in case

---

there's any kind of a disaster. The standby site is as good as the primary site.

From the disaster recovery perspective, a failover between two sites every six months is the plan, and this is something which should happen w a minimum network latency. The registrars should not get impacted because of the failover between the primary and the secondary site.

As I said, there won't be any kind of impact which registrars will have on their performance or on the registry experience which they'll have. So it will continue to be seamless for them, and they'll not come to know about the transition from the primary to the secondary. That's the plan which we have.

We again have the data backup and the offsite escrow. That is something which is part of the RFP, our network requirement or our internal requirements, so that's happening on a regular basis.

Coming to internationalized domain names, India has got the largest IDN diversity, and the government of India has a strong focus on taking the Internet to the last mile in the citizens in the remotest part of the country. And English is not the native language, it's not the native script for us. We have 22 recognized official Indian languages, and as part of the digital India program, we wanted to support the government's vision of enabling the local language or enabling the Internet through local language or multilingual Internet as we say.

So pre-transition, we had IDNs in 15 languages, and the rest eight are now almost ready and the one which are mentioned here in bold

---

italicized are created post-transitions. The above 15 were already there pre-transition. I'm happy to share that we worked hard to enable this and there's a roadmap for IDNs, and soon we'll have IDNs in all 22 official languages in the country. This has been the most critical part in terms of this transition exercise that happened at India.

Another critical part is the DNS transition. The DNS transition happened through – there were multiple steps involved. Just in terms of the steps through the transition, the incumbent TSP had their DNS zone made salve to the Neustar through this process [inaudible] the entire DNS transition happened.

On the transition day, Neustar began signing the DNSSEC zone. That happened from day one, that was something which was not there. So apart from this DNS transition, the second important factor was the data transition, so here were multiple rounds of data transition. The entire dump used to be there, the entire sample used to get dump from the incumbent TSP to the Neustar database. There were multiple rounds of testing and [inaudible] of the data.

After the transition, the entire database got transitioned, and one of the new features was the new AUTHINFO. As part of the transition, new AUTHINFO has to be assigned to all the domains. The new AUTHINFO for all the domains was stronger, bigger in size with extensive alphanumeric and special character combinations. The entire idea was to make it more secure.

Another change is in terms of the registry domain ID, which is ROID, and which earlier used to reflect the name of the incumbent TSP but now it

---

reflects the name of the .in registry once we go there at the WHOIS and check over there, we'll see that the name of the .in registry reflects.

Coming to the migration day, there were multiple steps. Again, all the stakeholders were kept informed about the migration day because this was the day when the entire effort of six months were to bear their results. We kept all notified. The incumbent TSP notified all the registrars that what will be the maintenance window. Then the shutdown happened at the incumbent level. The final escrow, the data dumpage was happening in the past, the final escrow export of the data happened, and the entire data got exported to Neustar.

So this was the step one which happened, and then the entire data which got transitioned, it has to be tested by Neustar. So Neustar received and validated the entire data, and as a part of the scope of work, registry.in, the website, was also to be managed by the technical service provider, so this also got transitioned to Neustar and it came under their hosting.

And post this transition and testing, the registrars were able to connect and sync with the registry, and ultimately, they were ready for the production. That is, they were ready for new registrations.

Specifically talking about the migration day timeline, against the communicated time of the timeline of 48 hours, it was completed in less than 30 hours. Herein, we have a learning that with more synergies, we could have further brought it down, but against whatever we communicated to our stakeholders, we were able to manage it in a

---

lesser time. But yes, a learning that if we have better synergies and a better plan, we could have further reduced it.

Talking about the key stats, one of the major statistics was the transition of .in and another 15 IDNs. Apart from this, against this 48-hour window, we managed it in three hours, but we were able to transition the entire 129 registrars on the platform of the new TSP, which is Neustar, and 94 registrars connected to EPP successfully. That's on the migration day.

These are some of the registrar feedback which was well appreciated by all our registrars that this entire migration was a success and it happened seamlessly.

So key lessons, well, there are many lessons I could call out, but the two mentioned here, one is the transition compliance, and the second is the stakeholder management. These are the two more critical as they actually relate to actual [go live.]

In any project, we should have the risk mitigation strategies in place, but as being part of this entire exercise, I would like to urge all the registries who are planning any kind of a transition that you should not only have plan A but plan B, C and a plan D also in place.

So basically, you should constantly evolve your plan. You should not believe that the plan A, the primary plan which has been devised, will work all the time. So you should have the fallback in place.

---

Another learning was in terms of how to mitigate any other external factors which could have delayed the entire transition. That is something which we faced, but that is another story.

So the main second learning is the stakeholder management. I would like to reiterate that this is something which is very important to the transition, stakeholder management and the constant communication is the critical part to this entire transition.

Especially talking about India which has got a diverse selection of registrars, we have almost 45% of our registrars are international, so there was a diverse selection of registrars who were involved in this entire transition. So having constant communication, transparency was required in this transition so that the multi-tier ecosystem which we have in place should be supported well.

In the end, where are we today post-transition? The .in registry continues to be on a full dual stack, IPv4, IPv6 [inaudible] it still continues to be on the full dual stack which will be there.

Anycast nameservers, all infrastructure has built in redundancy using the global standards of enterprise-grade infrastructure. Today there are two separate DNS platforms providing 30 global DNS sites, including two in India. One is in Delhi, the other one is Mumbai.

One of the [betterment] which has happened is the ease of IDN registration. Instead of punycode which earlier the registrant has to search for before registering any of the IDNs, you should have the

---

punycode, that ease has come in now. You can directly search the IDN in his or her language.

Also, as I mentioned earlier that security is part of the new design and we have almost 12 terabits of DDoS protection because of the capability of our current TSP.

I'll say it was not bad considering the entire transition, considering it first instance of transition. So since inception, for 13 years, we were supported by one TSP, so this was our first instance of transition. And now with our new TSP, Neustar, for next five years we are looking for a new age of success for .in.

And on a lighter note, to all the registries who are planning a migration and who have got Neustar or Afiliast as their TSP, I'll say that either they can get in touch with us or others who have got this experience so we'll be able to give them the consultancy for this entire transition.

Right, so this brings to the end. If any of the registry has any query, if anybody has any queries, here are my credentials. Thank you.

EBERHARD LISSE:

Any questions? In the meantime, I have facetious one. Why did you choose to outsource your operations to a TSP in the first place? Why don't you do it yourself?



---

SHUBHAM SARAN: That is something which has been discussed at the ministry level, the government level also, so that is something that's being asked by the government also. We'll plan it out.

UNIDENTIFIED FEMALE: Hello. My name is [inaudible] from .ca. Thank you for sharing your experience with us. We also did a very similar migration for .ca ourselves in February timeframe. I would like to ask you the first question is, so you say you used escrow [file] to do the migration from the old system to the new system. Did you actually import the history of the data into the new system?

SHUBHAM SARAN: Yes, we were able to get the entire history, because you need to have the entire domain lifecycle, the domain history in place. So we managed to have that.

UNIDENTIFIED FEMALE: Okay, so [it's not from] escrow file that's defined.

SHUBHAM SARAN: It's from the escrow files only. We regularly keep the escrow data with us, so we were able to get that.

EBERHARD LISSE: But the question is probably, what is escrow? Escrow under the gTLD thing does not include the audit.

---

SHUBHAM SARAN: But in case of ccTLD, yeah.

EBERHARD LISSE: If they have their own escrow data, for example dumping the whole database out into escrow, then it is different. That's probably what [inaudible].

UNIDENTIFIED FEMALE: Yeah. Can I ask a second question?

EBERHARD LISSE: You can ask one more question. We [are not strapped] for time. And then we'll go to [inaudible].

UNIDENTIFIED FEMALE: Okay. I'd also like to ask your experience, from your timelines, it seems the [inaudible] system is set up like one month before the actual migration. Is that a system for registrars to try their connections?

SHUBHAM SARAN: Right, as I said, registrars were the main stakeholders for us who were to be involved in this entire transition, so we gave them the testing environment one month prior to the transition.

UNIDENTIFIED FEMALE: Is that enough? The registrars, what's their experience?

---

SHUBHAM SARAN: It could have been extended, three months could have been a better time period, but one month was also good enough. As per our experience, we were able to get the proper feedback from them. [They were able to handle it.]

UNIDENTIFIED FEMALE: Thank you.

EBERHARD LISSE: Yoshiro.

YOSHIRO YONEYA: My question is very similar to previous question, but did you use final escrow data to check the data integrity after the data transfer?

SHUBHAM SARAN: Yes, we did.

YOSHIRO YONEYA: Because the change of registry operator is –

SHUBHAM SARAN: So being the registry operator, we should know. Yeah.

YOSHIRO YONEYA: [inaudible] efficiency data escrow is very important to experience.

---

SHUBHAM SARAN: Right, so being the registry, we had that.

UNIDENTIFIED MALE: [inaudible] .ua. So the escrow went forward with the history to Neustar. [inaudible] changes to the government authority that I suppose wants to keep that, or you just giving all this data to Neustar and the history going forward will be kept with them only? Personal data, domain data, everything.

SHUBHAM SARAN: It's with us. It's been with us only.

UNIDENTIFIED MALE: So it's going back to you as well, so they keep it [inaudible].

SHUBHAM SARAN: Being the TSP is they have to support. It's not – being the TSP they have to directly support the law enforcement agency also and whenever there's any kind a compliance requirement from that perspective, they need to have the entire history of the particular domain.

UNIDENTIFIED MALE: So NIXI would have a copy always?

SHUBHAM SARAN: NIXI has a copy, yeah.

---

UNIDENTIFIED MALE: [inaudible] from .ws. I have question about ROID. You updated it. Did you keep the old data? Or you just kept only new one? So if a registrar send [both] ROID, what did you do?

SHUBHAM SARAN: We kept all the data. This is something which I am taking from – all data.

EBERHARD LISSE: Okay. Thank you very much. AN enormous project, but it seems to have worked and it's important if anybody else tries it that we have got now two ccTLDs who can assist and advise.

Now we come to the conclusion by John Levine telling us all about Quantum and DNSSEC.

JOHN LEVINE: Alright. How many people here would say that they could explain what quantum cryptography is? Then we're all kind of at the same level, which is fine.

I do not purport to be a great quantum expert. My last physics class –

EBERHARD LISSE: Sorry, if I may interrupt, can you ask whether anybody can explain what a Quantum is?

---

JOHN LEVINE:

Well, it depends on what context. My last physics class was in about 1970. Eberhard very kindly let me squeeze myself onto the agenda a few weeks ago because at the M3AAWG conference, which was here in Montréal about two and a half weeks ago, we had a keynote on quantum cryptography by a guy at Microsoft.

We were actually quite concerned that quantum cryptography would break all of our existing security stuff. And one of the things he mentioned was that they've given a great deal of attention to websites and TLS, and practically no attention to the DNS. The only person who's been looking at this in relation to the DNS is Paul Hoffman who talked about something else here earlier.

So anyway, here is the executive summary, which is quantum cryptography is coming and it can break public key cryptography. So, should we panic? Well, no. Really? No, it's a lot more complicated than that.

In the next few minutes, I'm going to give a very brief overview of what the technology is and then look a little more detail specifically how it affects DNS technology.

Here is my ten-second introduction to public key cryptography. So digital signatures and public key encryption all depend on difficult mathematical problems, and in fact, there are two categories of mathematical problems that people have used.

The traditional one used in RSA is multiplying. It is relatively quick to multiply two numbers together, but once you have the product of those

---

two numbers, it is very difficult to figure out – to factor it back out into the two factors.

Another problem is – that’s what RSA uses. The elliptic curve uses what's called discrete logarithm. There's a certain kind of exponentiation which is mathematically fairly straightforward, but reversing it to find out what you would have exponentiated, that’s a discrete logarithm, and that turns out to be another intractable problem.

And then the other point to keep in mind is that as computers get faster, you can make the keys bigger, and every time you add an extra bit to a key, you double the number of possibilities. So every time you add a bit, the problem gets twice as hard.

When the RSA keys were originally 512 bits and expanded to 1024 bits, the difficulty did not get twice as hard, it got  $2^{512}$  times as hard, which is a very large number. And since cryptographers want to have a comfortable margin of safety, now we’re typically using 2000-bit keys, although at the moment, even factoring 1000-bit keys is not likely to be practical anytime soon.

With elliptical, the number of key bits is less because the problems are harder, but they started with 128 bits and typical elliptical algorithms now either use keys of 256 or 384 bits.

But in both cases, these are numbers that could be made much larger without changing the software significantly. We have 2000 keys now, we could have 4000 keys or even 10,000 keys. And with elliptic, we could go

---

from 256 to 512 or even 1000 bits, and it would be a lot slower, but the software would not significantly change.

So here's the problem, which is if you can build a quantum computer, quantum computers use things called qubits, which are analogous to bits, but through the baffling magic of quantum physics, a single qubit can have both the value 0 and 1. In principle, you could have an array of qubits that could simultaneously represent every possible combination of 1s and 0s for the underlying bits.

A mathematician named Shor studied this some decades ago, and he has an algorithm that in theory you could use with qubits to factor the numbers that RSA uses or do the logarithms that the elliptical curves use, and in theory you could do it very quickly.

Is this a problem? Well, maybe. The reason we're saying maybe is that quantum computers are unbelievably difficult to build. These computers are built at the bleeding edge of modern physics, and a qubit is – in a regular computer, a bit is sort of a logical thing which is implemented with a layer of electronics underneath, but a qubit is an actual quantum physical thing. They're very subject to noise from the environment. As noise leaks into the qubits, the answers you're going to get are going to be wrong.

So they need to be kept isolated from any sources of noise. One source of noise is heat so that your quantum computers need to run at cryogenic temperatures, cooled by liquid helium, and they also need electromagnetic shielding. It turns out that vibration will mess them up so they need seismic shielding. And you need to figure out how the



---

input and output and the control leads need to go in and out of this quantum system without displacing it too much.

And I won't attempt to explain it in more detail because this is about as far as my actual understanding goes. But another key point is that in any sort of computer system, if you expect to have errors, you do error correction, meaning there are – like on a disc drive, you're storing 32 bits, you might actually store 40 bits, using the extra bits for error correction in case a few of those ...

And you can generally make it so if you're storing 40 bits for 32 bits, as long as at least 36 of the 40 bits are correct, you can recover what the first ones should have been. And you can kind of do that for qubits except you can't actually copy a qubit. As soon as you observe its value, the qubit sort of disappears.

So the original idea was you can actually take a qubit and share its state among nine qubits and then ignore the ones that are messed up by noise, and it turns out that in practice, there's enough noise that you don't need to do nine times but like 100 physical qubits for every logical qubit.

This means that immediately, whatever size problem you are trying to solve, you actually need 100 times as many qubits as that all working reasonably reliably. And this is really hard to do.

The way I think of it is if you think of the state of computing in 1949, and as I'm sure you all recognize, this is a picture of John Von Neumann standing in front of the EDVAC in 1949. It was ordered in about 1945 and

---

it was delivered in 1949, and it did not work. It wasn't until 1951 until it worked well enough to be usable, and by well enough to be usable I mean it would go the better part of an hour between failures. In the late '40s, computers routinely ran for single digit numbers of minutes. I was reading something like you needed to dump out the value every minute because that's how long it'd work.

Quantum computers are kind of in that condition now. The capacity is very tiny, the reliability is terrible, but we've all figured out it's going to be great. And actually, a key thing that I didn't realize when I made this slide was that the memories for computers in 1949 basically didn't work. They were modified television tubes, and there was a huge technical breakthrough for core memory which made the computers much more reliable, and there's no similar breakthrough on the horizon for quantum computers.

So they're getting better. Optimists say that in 10-15 years, they may be able to break the codes that people are using now. The pessimists say it might take 25 years or it might turn out that they'll never be able to build a practical quantum compute that can break the codes that people are using now. On the other hand, they might, and we can't predict.

So if we're not wildly optimistic, we could have a problem now, and this is not so much a problem in the DNS, but if for example you're building computers that do code signing, if you sign the code now, you really want to sign it with a code that will be unbreakable for the entire period

---

of time that the program you are assigning is likely to be used. And people still use programs that are 10-15 years old.

Also, if you're worried about hostile nation states and other sophisticated opponents, you need to assume that all of your encrypted traffic will be recorded. So if they record it now and they break the code in 15 years, they can go back and see what you were saying now.

So again, if it's some sort of secret that is still interesting 15 years from now – of which there are a lot, like locations of where bombs are and stuff like that – then you may have a problem.

So NIST, which is the American government standards organization, is looking at algorithms that will do cryptography but are not subject to quantum attacks. And it turns out it's sort of bad luck that the two kinds of problems that we use are both subject to quantum cryptography.

Other things we use, like hashes, turns out quantum computers can't do hashes, so hashes are as secure now – even if there are good quantum computers, the hashes we use in our security systems will be as good as they are now.

So NIST asked people to send in candidate algorithms, which they did. They sent in a whole bunch of them, and they whittled them down until there were 17 encryption algorithms and nine digital signature algorithms which are all on a slide that I'll show you a version of.

They're inviting comments on them now and they're going to have a third round which will then whittle the candidates down to a modified

---

subset of the ones we have here. The goal is actually to publish some standard quantum resistant cryptography algorithms in 2022-2024.

What this means in practice is once NIST selects the algorithms, the crypto software should implement it about the same time because everybody is implementing it for practice now, since among the criteria that they need when they're evaluating these algorithms, can they actually be implemented? Will they run at a reasonable speed on normal computers? Are they actually fast enough to be usable?

For the sort of cryptography that is used on the DNS, the cryptography actually happens in hardware and high security modules, so they will need to redesign the hardware and the high-security modules to implement the new algorithms, which'll probably take another five years. So sometime in the 2030s, you should be able to implement this post-quantum cryptography, which if you believe the pessimists, is going to be just barely in the nick of time because that's when the quantum computers are going to be improving fast enough that we can't keep improving our old existing algorithms.

Now, how does this affect the DNS? As you can see now, it's maybe never. Certainly not soon. One reason is that none of the cryptography we do in the DNS uses keys that last for a terribly long time. As we heard earlier today, ICANN rolled the root key and they're planning to roll it again. So the root key in use at any given time is never going to be more than a few years old.

Everybody else who uses DNSSEC also should be prepared to roll their keys, and even if they don't roll their keys now, the data that you're

---

encrypting with DNSSEC is not expected to live for more than a year or two. And even if my old cryptographic keys were broken, to break DNSSEC you would need to break the cryptography that is used by the keys that people are using today.

So that's not likely to be a problem. In DNSSEC, people are currently using keys of 1000 or 2000. They can move up to 4000 just by changing a parameter in the software. The record formats will allow that. Similarly, the DNS elliptic curve keys are currently 256 bits, but they could increase to 384 just by changing a parameter, and by making fairly simple changes to the software standard, you could make the keys somewhat bigger there too.

DKIM which is used to secure e-mail represents the keys in a slightly different way, but currently, we have RSA keys that range from 1000 to 2000 bits, the elliptic key is 256 bits. I actually wrote the RFC that changed the key size. I can say that without a lot of effort, we can make the keys in DKIM bigger too.

So, so long as the key sizes of conventional cryptography can stay ahead of the ability of people to build quantum computers, we are okay. We can increase the key sizes. Again, it's generally no secret how well quantum computers are working. And again, maybe quantum computers will run into unexpected problems. I was talking to Paul this morning and he said that there's been remarkably little news about improving actual algorithms or actual reliability.

There was a press release from Google that said, "Oh, we solved this problem on a quantum computer that nobody could solve in a million

---

years with a regular computer.” IBM sent out a rebuttal press release saying, “Actually, if you arrange your disc storage differently, you could solve that problem in a day with a regular computer.”

So we’re a long way from having quantum computers that can solve problems that regular computers can’t. But again, just in case they might ...

2035 is not that far away, and for all the problems the DNS has, the chances of the DNS being dead by 2035 is pretty low. I imagine that particularly here at ICANN, people don’t expect it to be dead by then. So the key and signature length is a potential issue, and here's why.

Here is a peek at the DNSSEC keys from one of my domains. This particular domain has both a key signing key and a zone signing key which you can see are in italics here, and the zone signing key is the one at the top, it’s 1000 bits, and the key signing key is at the bottom, it’s 2000 bits. But whenever you ask for my DNSSEC keys, these are the answers you'll get back.

Similarly, every record in my zone has a signature, and every record is signed by one or the other of those keys, and in the top, you can see that that’s a signature on a text record, which is a signature over the 1000-bit key and the bottom one is the signature on one of the keys itself which uses a 2000-bit key. So whenever you ask for data from a DNS server and ask for the DNSSEC data to prove that it’s correct, you'll get signatures like this.

---

The problem is that the DNS doesn't really deal well with very large amounts of data. If I send a query asking for text records or something, the entire answer has to come back in one DNS packet. And there's a soft limit on UDP which is just send a packet, get a packet, of about 4k bytes. If it's bigger than that, you have to switch to TCP which is a little slower, but everybody's likely to be doing that for other reasons.

But there's a hard limit that a DNS packet cannot be bigger than 64k. There's a 16-bit length word that you cannot change without breaking all the DNS software in the world.

So when you make a query, and particularly if you're in the midst of rotating your keys and you'll have multiple keys and you'll have multiple signatures, you might get an answer sort of like this. Here's the keys for Comcast, a large American cable company.

They're in the process – they're just rotating signatures. You'll see when I asked for their key records, there are two key records and two signature records there, and they're pretty big. This particular answer is like 2-3000. So this is fine, but if the key size doubled or got even ten times bigger, it could be 20-30,000.

But if the keys got a whole lot bigger, if we needed to go like 8000-bit RSA keys, and particularly if people start doing more regular key rotation as is recommended here at ICANN ...

I should have come up with a better example here, because these keys aren't being rotated, but if they were in the process of rotating their

---

keys, there would be twice as many keys and there would be twice as many signatures, and the size of the DNS packet can get quite large.

In DKIM, which is used for e-mail, the keys are in the DNS and the message signatures are in the mail. This is the DKIM key for this month for my mailing list server. The first thing is a text record, and the first blob of hex there is the DKIM signing key, but since it's a record that's in the DNS, it's also signed by DNSSEC, so every time you get that key, you also get the DNSSEC signature.

So again, if the keys got a lot bigger, there might be issues from the DNS. For DKIM, the message signatures go in the header in the mail message, and they're what you can actually see here in italics. That's not such a problem. Mail message headers can be enormous. If I needed to have a 10,000-character signature in a mail message, I think some of the people with whom I correspond would [inaudible] they would claim there was too much junk in the message, but it would actually work. So the real things we're concerned about are the keys and signatures for DNSSEC and the keys for DKIM.

So here's why there's a problem, and I got this slide from Brian LaMacchia who gave that keynote I was talking about, which is these are some of the signature algorithms that are currently under consideration by NIST.

He did some estimations like to get the equivalent of NIST's security level one, which is roughly the security you would get with a 2000-bit RSA key, which is what people are typically using for zone signing keys now, and what we're using for DKIM keys now, the first column there is



---

how big would the public key be, how many bytes, and then if the key has to be expanded into a base64 representation, which is used for DKIM keys, how big would that be? And then the last two columns, okay, how big would the signature be in both originally and in base64?

You can see from the ones that he has circled, some of the keys and signatures are small, but some of them are really big. The most extreme algorithm, the third one, to get the equivalent level of security, that needs a 300,000-bit key. You can't even fit one of those in a DNS packet.

And similarly, the chances of them actually using this algorithm are low, but it's there on the list because in principle, it works. And also, if you're doing stuff on web servers, the variable length fields for keys in the TLS certificates the web servers use, if they needed to be really big, they could be. So they don't have the sort of hard limits we have here.

But looking here, there are two algorithms with keys where even a single pair of keys wouldn't fit, and then if you look at the signatures, they range from perfectly reasonable sizes like 800 up to 27,000 or 17,000.

And again, this is just for the level of security we have for 2000 keys. The chances are that if people are actually going to deploy post-quantum crypto, the keys and the signatures will be bigger.

So this is why it's a problem. In the unlikely event that NIST standardized one of these algorithms with really big keys or signatures, and it turned out that quantum computers came along at the optimistic rate so you actually did need to have the equivalent of 4000-bit RSA

---

keys but instead using post-quantum crypto keys, then the DNS wouldn't work. You would have DNSSEC records that would not fit.

We could come up with ways to work around it, we could break them up into individual pieces, but how many years have we been talking about “please deploy DNSSEC” which has not changed at all in the past decade? If we need to come back in five years and say, “You know that stuff that you just spent 15 years deploying? Well, sorry, it broke. Now you need to deploy this new version,” we would like not to do that.

NIST probably won't pick the algorithm that uses the giant keys and the giant signatures, but they are asking for comments on the algorithm, so I was planning to send them – and the comment period closes in a month – a note saying I use DNSSEC and I use DKIM, and here are some issues. If the keys or the signatures were really big, I would have a big problem, so please pick one of the algorithms that uses normal sized keys.

So for anybody who actually runs a significant bunch of DNSSEC software, it would be worth commenting to NIST. Other than that, this is just something to keep an eye on, and now you too can be your local quantum crypto expert.

Also, if you want to do experiments, both DNSSEC and DKIM have a key algorithm switch. And it's perfectly normal to be changing algorithms so that the software will ignore algorithms [it doesn't understand.] So if you want to do experiments, try to implement some of these post-quantum crypto things and see how they work. You can do it now. And

---

the people who are working on TLS have been doing these experiments for a couple of years.

That's it.

EBERHARD LISSE: Thank you very much. Interesting stuff.

[TAYLOR:] [Taylor] from Oracle. If I want to do experiments with quantum computers, how do I get my timeshare account on that?

EBERHARD LISSE: Buy enough Google shares.

JOHN LEVINE: Yeah, that's above my budget and pay grade. I was suggesting you might do conventional algorithms of the quantum resistant algorithms.

[TAYLOR:] I'm still kind of where you are in trying to wrap my head around just how are we making these things practical.

JOHN LEVINE: I had pretty good physics when I was back in high school and college, but you have to have algorithms that are reversible, which is ... And it turns out multiplication and factoring is reversible but hashing isn't. I start from there, then I go downhill.

---

FREDERICO NEVES: Hi. Frederico Neves from NIC.br. Two comments, the first one is regarding the [to be] seeing these as an issue, but I'm totally agreeing with you, this is not a current issue but something for us to take a look at it. But something that people should take into account is that algorithm rollover is much different than key roll.

Being prepared for algorithm rollover or having cryptographic agility is more important than doing key rolls these days. So being prepared for algorithm rollover is very important.

JOHN LEVINE: Yeah. We need to do both, because you can't switch algorithms without also rolling the key.

FREDERICO NEVES: Yes. And the comment that you make regarding increasing the size of the keys in case of the [EC,] it's an algorithm roll, actually.

JOHN LEVINE: Right. I kind of glossed over that. For RSA, the key size is a parameter whereas for elliptical [inaudible] is part of the algorithm.

FREDERICO NEVES: Yes, for [EC] it's a different algorithm.

---

EBERHARD LISSE:                      Okay. Thank you very much. So Ondrej Filip will now sort of wrap it up by reviewing what was happening.

ONDREJ FILIP:                         I will just wrap up what happened today. It has been a long day, and all tired, so I will try to be brief. We have a perfect set of presentations today. We started with an excellent presentation from Benin, those guys are doing DNS forum in the DNS hackathon. The DNS forum is attended by 700 participants if I'm not mistaken, which is something incredible there.

Then we had a presentation from Jacques Latour, he mostly talked about Canadian food, poutine, that was quite interesting, and then he briefly mentioned something about exchange points, I think. Their approach, how to create a new exchange point, was impressive and it's a good example how a registry can help the local community. Really great presentation.

Then Andrei came to frighten Eberhard a little bit with the local port scanning, and it was an interesting presentation and I think it would be really nice to see some follow-up because there are some unanswered questions on what is the purpose of this behavior and stuff like that. It's a really great presentation from Andrey.

Somebody from Czech Republic was presenting how they upgraded DNS infrastructure, and then we had lunch so it was the perfect part of the day.

---

After that, we came, we had two presentations from two Marks about RDAP deployment, both from completely different angles. It was interesting. It's quite frustrating to see that RDAP is not starting in ccTLD space. So there's a lot of work we can do there. Gs are a little bit better, at least from the stats that Marc Blanchet presented. And also, the internal work of Microsoft, how many teams are dealing with that data, that was quite impressive.

After that, we have two presentations about machine learning, and I found out there were even some people here who understood it, which is perfect. They were used for different purposes. In Qatar, the machine learning is used to find out domains which are in the wrong suffix, and in Japan, this is used to predict whether the domain will be renewed just to have the precise numbers for the future.

Then we have a host presentation. It was quite refreshing to see that even in such a strange environment like justice, IT technology can be used and will be used hopefully soon.

Then we have like three presentations about security. Bruce Tonkin talking about pen testing in Australia, Paul Hoffman had more discussion points, there was quite an interesting discussion after his presentation about the root zone scaling early warning system, and then we had Ashwin and David from ICANN talking about ICANN security practices.

In-between, there was a presentation from Maarten. He gave us an update about ENTRADA software from the Netherlands. Again, perfect work for analyzing the DNS data, but not for Eberhard.

---

And then we have an excellent presentation from India about changing the technical service provider. It's something that doesn't happen very often when 2.2 million domains are transferred from one provider to another, so it was really interesting to see, and I was really impressed by how smooth the process was, so it was a perfect presentation and I wish all such transfers would be as smooth as this one. So it was perfect.

And less, I hoped that there would be another attempt to frighten Eberhard, so I hope John will tell us that everything is broken and we should stop doing DNSSEC and so on, but no, it's not like that. It seems to be that the post-quantum algorithm will need some time. So we can still quietly sleep tonight.

So that was all. I have to thank two persons, [inaudible] and Jacques for helping me take notes, and of course, Eberhard for creating such a great agenda today. Thank you very much.

EBERHARD LISSE:

And being an obstetrician, I don't frighten much.

**[END OF TRANSCRIPTION]**