## EPISODE 92

[INTRODUCTION]

**[0:00:10.4] YN:** Hello and welcome to another episode of TWIML talk. The podcast where I interview interesting people doing interesting things in machine learning and artificial Intelligence. I'm your host Sam Charrington.

This week on the podcast, we're featuring a series of conversations from the NIPS conference in Long Beach California. This was my first time at NIPS and I had a great time there. I attended a bunch of talks and of course learned a ton, I organized an impromptu round table on building AI products and I met a bunch of wonderful people, including some former TWIML talk guests.

I'll be sharing a bit more about my experiences at NIPS via my newsletter which you should take a second right now to subscribe to at twimlai.com/newsletter. This week, through the end of the year, we're running a special listener appreciation contest to celebrate hitting one million listens on the podcast and to thank you all for being so awesome.

Tweet to us using the #twiml1mil to enter. Everyone who enters is a winner and we're giving away a bunch of cool TWIML swag and other mystery prizes. If you're not on Twitter or want more ways to enter, visit twimlai.com/twiml1mil for the full rundown.

Before we dive in, I'd like to thank our friends over at Intel Nirvana for their sponsorship of this podcast and our NIPS series. While Intel was very active at NIPS with a bunch of workshops demonstrations and poster sessions, their big news this time was the first public viewing of the Intel Nirvana neural network processor or NNP.

The goal of the NNP architecture is to provide the flexibility needed to support deep learning primitives while making the core hardware components as efficient as possible. Giving neural network designers powerful tools for solving larger and more difficult problems while minimizing data movement and maximizing data reuse.

To learn more about Intel's AI products group and the Intel Nirvana NNP, visit IntelNirvana.com. In this episode, I speak with Yael Niv, Professor of neuroscience and psychology at Princeton University. Yael joined me after her invited talk on learning state representations. In this interview, Yael and I explore the relationship between neuroscience and machine learning.

In particular, we discuss the importance of state representations in human learning, some of her experimental results in this area and how a better understanding of representation learning can lead to insights into machine learning problems, such as reinforcement and transfer learning.

Did I mention this was a nerd alert show? I really enjoyed this interview and I know you will too. Be sure to send over any thoughts or feedback via the show notes page. Now, on to the show.

[INTERVIEW]

**[0:03:15.7] SC:** All right everyone, we are here in Long Beach at the NIPS conference and I've got the pleasure to be seated with Yael Niv. Yael is a professor of neuroscience in psychology at Princeton University and she delivered a talk this morning on learning state representations and I'm really excited to have her here to talk to us about what she's working on. Yael, welcome to this week in machine learning and AI.

**[0:03:40.6] YN:** Thank you, thanks for having me.

**[0:03:41.5] SC:** Absolutely. You are a neural scientist, a neural psychologist, you're here presenting on machine learning, you talked about reinforcement learning, how did all of this come about? How did you end up at the intersection of these two fields?

**[0:03:59.3] YN:** The funny thing is I didn't start as a neuroscientist. Well I started as a computational neuroscientist, more interested in AI and in psychology but not and experimentalist like I am today. Started as a theoretician. My PHD was in reinforcement learning theory, my main advisor was in computer science, my secondary adviser was in psychology, she was –

I asked her to be an advisor so she would ground me in real data and I won't be making models of things that I made up rather than the real world. My PHD was also completely theoretical, I was modeling animal behavior with reinforcement learning, submitting my papers to NIPS, was getting rejected by NIPS.

I got a best pay for award finally. When my paper was finally accepted at NIPS – I was actually really excited now. I have drifted away, I've drifted more into psychology, more into neuroscience, more into doing experiments. I still do modeling but it really is now half-half with experiments and so I haven't been coming to NIPS for the last few years and I was really excited to be invited to give a talk here because it felt like –

**[0:05:10.1] SC:** Coming home maybe?

**[0:05:10.6] YN:** Yeah, closing circle from kind of you know, all those papers being rejected and now, giving an invited talk. I immediately said yes.

**[0:05:20.5] SC:** Nice, very nice. What do you consider to be your home conference now?

**[0:05:26.2] YN:** My home conference now is reinforcement learning and decision making, RLDM, I help found that conference and we have it every two years, there was one this past year so the next one is in 2019 and it's basically a very interdisciplinary conference that tries to bring machine learning, AI, psychology, ethology, economics .

Anybody who is interested in reinforcement learning and decision making over time. It's a growing conference, last time I think we had about 600 people.

**[0:05:59.2] SC:** Wow.

**[0:06:00.6] YN:** That's my home base.

**[0:06:02.0] SC:** Nice. One of the things that I've observed here and this is the first time I've ever been at NIPS is that there seems to be several kind of themes that have emerged from me. One

of the strong themes is kind of integrative approaches and interdisciplinary approaches. Has nips always been like that? Or are we seeing more of that now than in the past?

**[0:06:24.6] YN:** I think NIPS has always aspired to that, it waxes and wanes so people have talked about you know, the end in NIPS, is it sometimes more IPS, and where's the neuro. There have been you know, worse years and better years. I used to come regularly so I remember this and there are also kind of fads and things have become more popular and less popular.

I remember when I was a master student, someone coming back from NIPS and saying that it should have been called Support Vector Machine Conference that year because it was all support vector machines. Now, almost nobody talks about SVM.

**[0:07:00.4] SC:** It's all deep learning?

**[0:07:00.8] YN:** Everybody talks about deep learning. Kind of changes focus but there's always been I think a sincere attempt to keep the neuroscience in the mix and it's not. It's not super easy because I think it's very clear that neuroscience has a lot to learn from, or not necessarily to learn, but to take from machine learning so we use machine learning algorithms to analyze our data to think of computational processes in the brain.

Really, I feel like when I come to NIPS, I'm coming as like a shopper. I want to see what's on the shelf now that I can use in my research, I don't think it's necessarily that way, the other way around. I think AI in particular has been inspired by neuroscience but mostly, takes quick inspiration and runs off to somewhere else and that's fine because their goal was different, my goal is to understand the brain.

The goal of Ai is completely different.

**[0:07:54.6] SC:** It's almost aspirational inspiration.

**[0:07:57.7] YN:** Yeah, I think basically, AI needs neuroscience less that neuroscience needs machine learning.

**[0:08:03.7] SC:** Interesting.

**[0:08:04.5] YN:** It's always been kind of a tough thing to keep the mix together.

**[0:08:08.1] SC:** But in your talk, the focus of your talk was that mix, right?

**[0:08:12.1] YN:** I tried my best. It's not the typical talk that give because these days, I don't usually talk to this type of audience and so I thought hard, you know, what do I have to sell to AI? What do I want AI to buy? And so I was posing this challenge of real, well, I don't know, call it real with a capital R.

But, Intelligence that is more like human Intelligence, you know, playing chess or playing Go at an expert level. That's a huge achievement because these are really hard tasks but they're really hard for humans too. You can't just become a Go expert tomorrow.

But there are other things that we learn super easily. Those things are very hard for computers as well for AI. You know, maybe in some ways it's hard as playing Go and I think that's a real challenge because those are – because of the amount of flexibility that human brain has, the extent to which we learn super quickly and we can toggle tasks.

Do one task and then a minute later do another and then we turn back and kind of toggle the representations, the policies, all the machinery that's needed for these tasks is really impressive and that's what I have not seen in any AI yet.

**[0:09:34.2] SC:** Right, I thought that was a really interesting characterization of kind of the challenge of AI or what you believe the challenge of AI should be.

**[0:09:39.8] YN:** A challenge.

**[0:09:41.4] SC:** A challenge.

**[0:09:41.2] YN:** One of them.

**[0:09:42.1] SC:** Should be. You know, as supposed to kind of these grand challenges like Alpha Go or like the game of Go, you propose many simpler problems, solved with much less data as a challenge, is that the way – did I characterize that right?

**[0:09:59.6] YN:** I don't know if the problems are simpler. I gave the example in my talk of crossing the street and it seems like a simple problem because we already think about it as a simple situation but really, the stimuli, the auditory, visual world around you when you're doing a simple task like crossing the street, it's very complex.

You need – you know, sophisticated visual machinery in our brain to parse out the objects, you need sophisticated auditory machinery, et cetera. But more than that, you need to know, even if I've parsed everything. Most of it is still irrelevant for the task, I don't even need to parse it, I don't need to parse what are the stores on the other side of the street.

I don't need to parse even the colors of the cars because that's irrelevant for the task. So, whether my visual system does that or not is one question, but my learning system should definitely ignore those aspects. And that's really hard because there are potentially infinite combinations of things that I might need to ignore or pay attention to. There's everything that I see and then there's everything in my past.

Because it might be that the time of day – or unobservable things, not only in my past, or time of day matters, the city I'm in matters. What I did yesterday matters, maybe not for crossing the street but if I'm driving a car and navigating what I know about traffic from past learning.

Potentially, infinite dimensions of the environment. I have to narrow them down to like three, four, five.

**[0:11:30.0] SC:** one of those examples you gave of kind of the state that goes into figuring that out was Washington DC versus New York?

**[0:11:38.7] YN:** Yeah, the idea is that you know, I could also say, New York versus London. Long Beach versus whatever. The context really affects, let me take one step back. On the one

hand, we want to generalize broadly and that's why it's really important to ignore some things and kind of represent and learn about only others because if I'm really – if my brain is still sophisticated, it can take into account everything in my visual scene, that would seem to be optimal but it's actually super sub-optimal.

Because it means that when I learn something new, I don't know what to generalize it to, is this true only in this intersection only when a red car is coming at me, only when I'm in front of this shop? Depending on the task, let's say, crossing the street. What I've learned is probably much more general than that. I might learn that in Princeton, the minute you step down off the curb, all the cars stop and let you go no matter what.

That's general to all of Princeton. Not only to that place but only to Princeton, right? I wouldn't want to learn that about Long Beach. It's all about setting the boundaries of generalization and what I've realized form this work is that really, when you think about real world experience.

You never ever see two situations that are exactly the same. Even if you have that same intersection, it won't be exactly the same crossing the street tomorrow. The fact that we can reuse any past experience, the fact that we can learn at all, means that we have to generalize.

The basic thing that I think about all the time is how do we determine the boundaries of the generalization? I come from reinforcement learning where we think a lot about how do we learn values and policies? I'm thinking, that problem is pretty much solved, I mean, we have a lot of algorithms for that.

We know how it's done in the brain. The question is not how do I learn values, the question to what do I generalize those values? What's the value of what? What state? What situation? I want to go back – you said that my challenge was to solve – or the challenge I pose was to solve simple problems with less data.

Let me try to phrase that a little bit differently. There is this old story, I don't know if it's a myth or not but I heard about it at NIPS, that basically, well, it happened at NIPS, it's not about NIPS. What I heard is that in the 80's, kind of like in the first hay day of neural networks, now there's a

second. Someone published a NIPS paper that showed how you can use a neural network to parallel park a truck.

You know, this is in the 80's, think of you know, computing power in those days and stuff. This is a really hard challenge to parallel park a truck. After all the hoorah and happiness, apparently or as the story goes, the next year, someone else published a paper in NIPS showing that one neuron, per subtron, one simple computing unit can parallel park a truck if you give it the input of the obstacles and polar coordinates rather than cetacean coordinates. What I took from that, I heard this when I was a master student, is basically, if you ask a question just right, it's really easy to answer.

That's basically what I'm saying, decision making, could be kind of paired down to a yes/no question that a per subtron can answer. That pairing down is the hard thing. Giving the input just right and so that's kind of my whole enterprise, my whole research enterprise, is how does the brain learn how to ask the questions, to ask complex questions so that they're made simple.

You know, how can we use that for AI to take all these complex problems, real world navigation et cetera and make them into easy problems that a per subtron can develop.

**[0:15:36.2] SC:** Interesting. Is there evidence in the rain that that process looks like a filtering process or is it more like a blindness to certain irrelevant variables or is it the way things are classified, it seems like there are potentially a number of ways that you cannot consider a variable?

**[0:15:56.4] YN:** Yeah. The evidence that we've seen so far in the brain that I actually didn't talk about today and so I'm glad that you asked about this is that attention processes are strongly involved in sub-selecting. First of all, anatomically, I should say that the areas in the brain that we know are involved in this kind of reinforcement learning and decision making get input.

It's a small area in the middle of the brain, it's called the stratum. The stratum gets input from the whole cortex. All the brain all around the sensory motor, high level emotional, everything. Everything funnels in but it funnels in with a huge dimensionality reduction. About one to 10,000.

10,000 inputs into one neuron on average. From the start, we know that there has to be a huge dimensionality reduction. But more than that, I'm talking here about the mentality reduction that's not structural that's kind of modulated, it's task specific for one task, I want to know the colors of the cars for another task I don't.

The input has to really be modulated on the fly according to our goals and what we've seen is that areas in the brain that are involved in attention and switching attention and selecting what to attend to are involved in this process.

It's really interesting for me to think about attention in that way because we often think about selective attention or limited attention as a bottleneck limitation. We can't attend to everything so unfortunately, we have to ignore most of the scene. I'm thinking about it in a much more normative way, that it's optimal to not attend to everything.

**[0:17:32.4] SC:** It's what allows us to function.

**[0:17:33.9] YN:** It's what allows us to learn. Right. If you attend just right, again, you're solving a problem rather than creating a problem. One of the challenges is to understand how do we learn what to attend to? We're working on that, it's not easy, we have tasks where we can measure people's attention.

We have a multi-dimensional scenario, we tell people that only one dimension matters. These dimensions could be color, shape, texture of different stimuli or they could be other dimensions. We tell them that only one determines reward and we're trying to see how they learn from trial and error, what to attend to and in essence, what we're trying to do is defy, is just figure out –

Reverse engineer the algorithm, the computational algorithm by which feedback for our action affects our attention, affect what we attend to. We're in the midst of that, it's a tall order, we've actually talked to – it turns out that there aren't machine learning algorithms for modeling those kinds of data.

**[0:18:38.6] SC:** I hear attention come up a lot in the context of LSTM networks, is it just kind of the word overlap but not the same kind of mechanism or idea or are there some parallels there?

**[0:18:51.7] YN:** I'm not well enough versed in these models to know how direct the parallels are. I can say that from a psychology point of view, a neuroscience point of view. Attention is almost kind of a dirty word because we don't know exactly what it is, some people hate that word.

A selective filter is really kind of operationalization of this idea. I think there's been very little work, both in neuroscience and well, I don't know recently in machine learning. It's a truth I haven't really followed it. But in neuro science, there's been very little work on how attention is learned from experience rather than from instruction.

Almost all the work on attention is you tell the subject what to attend to and then you figure out how is your brain doing that, how are they focusing, how much are they sensitive to distractions or et cetera. My question is completely different. The place where I say that there aren't any algorithms for is to try to –

What we want is to predict attention and we want to test our models, we want to evaluate how well our models are doing. With choice data, which a lot of our experiments we have subject's choosing actions. There's a whole wide range of models for modeling choices and comparing between models and saying, you know, this model is better than that model in predicting choice. What is missing for me is models that predict attention and a way to compare those.

The difference here, this is a little bit going into the weeds but if you think of attention as a quantity that sums up to one, if I attend more to one thing, I can't attend more to everything else. This is a quantity that lives on the simple X.

It's a whole different statistics and geometry and comparison world on the simple X and people haven't really worked on that. Because of this constraint that everything adds up to one. It's more – if you need a different math.

**[0:20:43.8] SC:** Okay.

**[0:20:45.1] YN:** We found a little bit of this math of all places in geology, apparently, in geology, geologist want to ask, is this rock the same as that rock by looking at its composition and saying, it's 80% this material and 12 to 15% this material, is it the same as that one?

They do this comparisons of percentages and there is a whole book actually on these compositional models. We're reading that now.

**[0:21:11.6] SC:** Wow.

**[0:21:11.7] YN:** That's what we're trying to find our new math.

**[0:21:15.1] SC:** Okay, interesting. You talked about – you've mentioned reinforcement learning in our conversation as well as in your talk. But when I hear you say it, it's almost like the context is, you're talking about the machine learning reinforcement learning but you're also talking about like biological reinforcement, am I reading that right?

**[0:21:34.4] YN:** Yeah. Reinforcement learning in psychology is called classical conditioning and instrumental conditioning.

**[0:21:40.4] SC:** Okay.

**[0:21:41.2] YN:** The classic Pavlov with the dog salivating, et cetera. That is learning values for states as in reinforcement learning and then instrumental conditioning, rats pressing a lever in order to get food, that's basically learning a policy that maximizes values and rewards.

There is a very kind of direct isomorphism or translation between the computational reinforcement learning world of learning policies and values and predictions. The behavioral psychology world and it goes through neuroscience or it goes to neuroscience because we also know in the brain where these different signals are computed and there's a lot of evidence, this goes back to the 80's actually.

The 80's, no, the 90's. I think '94 was the first paper really making this parallels in particular dopamine, which is a really important neuromodulator and the brain that's involved in everything

from Parkinson's disease to schizophrenia, depression, gambling, any drug abuse. Dopamine is seen today as representing of the brain, prediction errors ala reinforcement learning.

Literally –

**[0:23:00.5] SC:** Wait, what? Say that again?

**[0:23:02.6] YN:** Yeah. I'm often amazed that people don't know this because in neuro science, this is so big, this is basically one of the biggest success stories of taking machine learning, taking computational models and translating them to neuroscience and behavior is this, the idea of that dopamine calculates a prediction error.

How different is what I'm getting form the world from what I predicted? This is a key quantity for learning and reinforcement learning. Every reinforcement learning algorithm relies on predictionaires. We know behaviorally that animals learn from predictionaires. We know neurally that dopamine represents these predictionaires and broadcasts them to the whole brain so anywhere in the brain can –

**[0:23:47.7] SC:** What does it even mean for dopamine to represent these predictionaires, meaning, the levels of dopamine correlate strongly with predictionaires?

**[0:23:56.0] YN:** Yeah, it means that when in a task, I can, through a computational model say you probably just experienced a positive predictionaire so you expected let's say, three M&M's and you got five. Through the model, I can track your learning and say, with all the experience that you've had so far, you should be expecting three M&M's.

You just got five so it's a positive predictionaire. If I record it from your brain, I would see a kind of positive burst of dopamine. Above baseline dopamine. Whereas, if my model says you should predict three M&M's and you actually get one M&M and I record form your brain, I'll see a negative dip, I'll see less dopamine than the baseline dopamine levels at that exact time.

It's a phasic, it's a short lived signal, it's broadcasts all over the brain, it basically says that every point in time, are things better than expected right now or are things worse than expected?

Alpha Go relies on predictionaire. All of reinforcement learning relies on predictionaires. We see that in the brain. Really, when I say reinforcement learning, I'm thinking about the brain as much as I'm thinking about the algorithm.

**[0:25:01.5] SC:** Okay, that's fascinating.

**[0:25:03.4] YN:** Yeah, it's really an amazing case of convergence of all of these lines of research.

**[0:25:10.8] SC:** And is this a new observation or did you say 80's for this or?

**[0:25:16.7] YN:** '94, not 80's. So '94, so in the 80's, in the late 80's the ideal was a dopamine represents reward. So it's the brain's signal for reward and people started looking for that signal by recording from dopamine, or from monkey brains, as monkeys were obtaining rewards in a task and they were really confused because sometimes dopamine would respond to the reward, sometimes it didn't and there were all of these abstract and the society for neuroscience saying, "You know this is clearly not a reward signal but we have no idea what it is".

And then, Read Montague and Peter Diane who were then post-docs with Terry Sinofsky at the Salk Institute read these papers and they have been reading about reinforcement learning. So Read an end in Andy Bartose work and they basically put two and two together. The story is that one day Read Montague went to Peter Diane's office and said, "Look at this. This looks just like a predictionaire signal" and then they published a paper.

They started from the politics of these things are weird, they started a paper in science about bees not about monkeys. About bees navigating with the idea that octopamine which is the equivalent of dopamine signals in insects, octopamine represents a predictionaire and then later, they published a paper about that same thing in monkeys a year later in '96. So this first paper was '95, I think that was an abstract to '94 and then in '97 the famous paper is Shultz, Diane and Montague in 1997.

Where they published basically the recording data together with the model and said and this was published in Science saying, "Dopamine seems to be a predictionaire" and this was the

hypothesis with some data and since then it's been tested a million times over and on the one hand it looks like sometimes it's so amazing. It looks like dopamine runs must have read the text book like they do exactly what you expect.

In really convoluted situations you set things up so that it should be whatever and it is exactly that but with neuroscience, the deeper you dig, the more unexplained gold you find. So a lot of it fits the theory, some of it doesn't. And that some of it is not esoteric stuff that we can ignore. It's not just noise, it's persistent differences. So the model, this idea is simplified. So I believe that dopamine definitely represents a predictionaire but it's not the end of the story. That is the beginning of the story.

There is a lot of work trying to understand exactly how timing is represented in this system and are these predictionaires only for reward or for any kind of prediction that's violated and if it's any kind of prediction that is violated, how do you know what to learn? How do you know what prediction was violated because when it's reward it's easy. You just update your prediction of reward. If it is anything, it becomes, you basically need more information in order to learn.

**[0:28:33.9] SC:** It is making me wonder are there other chemicals that respond similarly that?

**[0:28:38.3] YN:** Yeah, there aren't other chemicals in the brand that look like a predictionaire signal. There are four neuro modulators in the brain. So there are lots of chemicals in the brain but neuro modulators are signals that rather than communicating neuron to neuron kind of like a phone call they are more like a PA system. They broadcast very, very widely. So there is dopamine, norepinephrine, acetylcholine and serotonin.

And people have basically been dying to know what these four do because it's like if you had four broadcast systems, those are four things that you can tell everybody, what will you say with those four? And it's clear that all four are super critical in the sense that disrupting any of these systems causes a whole host of problems which makes sense. You know if they broadcast everywhere they must be doing something really important. So dopamine is the best –

**[0:29:29.2] SC:** So far only dopamine has been implicated in the learning process or is that too strong of a statement?

**[0:29:33.8] YN:** Learning is really important, our brain learns all the time. Dopamine is the best understood but things like this, the simplest of the stories but the others have been implicated in learning as well. So norepinephrine has actually been implicated in the breadth of attention. Controlling the breadth of attention, controlling the gain of the system, controlling what you do or how do you respond to unexpected changes.

So one of the ideas in learning is that when the world changes in an unexpected way you should be able to reset to not continue to carry on previous learning but to start anew. To reset your values, to increase your learning, great to say, "Okay I'll take a break here from the past" so norepinephrine is implicated in that. Acetylcholine is also implicated in adjusting learning grades to the variability of the environment. So an environment that is more variable basically it's more noisy.

So every bit of information should carry less weight. So you should have a smaller step size in learning, we call it a smaller learning grade. So acetylcholine is implicated in that and serotonin has been kind of – serotonin is so complicated. There are 20 different kinds of receptors for serotonin and they do all kinds of things. I mean all of these are complicated, everything that I'm saying is a gross simplification but yeah, that's what they pay us neuroscientist for you know? We are trying to figure it out.

**[0:31:01.2] SC:** So I don't know, we have strayed into a lot of background material which has been amazing but your talk, right? So I won't even try to summarize it in a sentence, I will let you do that but it was – well no, that's what I would do to introduce it. You can actually walk us through the framework that you presented and some of the experimental results and things like that. So kind of broad strokes, you are applying the Bayesian Inference to learning.

And at least what I got out of the talk was trying to identify these latent variables that are present in the way we kind of perceive the world through experimentation and relate that back to machine learning or statistical models.

**[0:31:54.8] YN:** Yeah, so this goes back to what we talked about in the beginning of the podcast today which is how to we part the world into these? How do we put boundaries on learning and

say all of these experiences are similar enough, I am going to learn from one to another, I am going to learn from this street to that street? And these are different – this is London, the cars come from the other side of the road. You definitely not generalize between these.

And so what my talk was about is trying to identify and what my research is about is trying to identify the computational algorithms for putting these boundaries in place because I think that's – this is a computational way of talking about the issue of how do we take a complex problem and make it into a simple one. So when I cluster experiences together and say all of these are similar, I am basically saying I am going to ignore all of their differences.

So that is part of learning what to ignore. I am going to say these are essentially – you know in reinforcement learning we would say, "This is state one," right? We just give them a label. All of these things are state one so I don't have to analyze all of their minute differences. So I am trying to understand how the brain decides what is state one. How the brain does this clustering and it's really – it's nothing I am trying to identify what are the relevant aspects for each task.

That's less interesting for me, I am trying to understand what is this general purpose algorithm that can take complex input, like we talked about before, it can have potentially infinite dimensions that are relevant and can easily say based on inputs one, two and three, based on these three dimensions and none other, I'm going to call this state one and those three dimensions in a different scenario say that that one is state two and that's not state one.

So that's what I want to understand and I think because this is a MP heart – well I don't know MP heart, I haven't proven it, it's a very hard problem. Let's say it this way, you can't actually use Bayesian Inference for this. In my talk I talked about the Chinese restaurant process prior. So a way to start from a prior that says there are a few latent causes for all of the observations that we see and we are going to try to cluster observations according to similarity and each cluster I'll call it a latent cause and that will be my state one and state two.

**[0:34:20.9] SC:** I was wondering this actually, what's the backstory for Chinese restaurant problem or the process?

**[0:34:26.5] YN:** Oh that comes from Stanford, you know people live in Palo Alto and go to Chinese restaurants. The backstory for this culinary metaphor is the idea that imagine you have a really large Chinese restaurant, like an infinite Chinese restaurant, and you think of each table as a cluster, a latent cause and each customer is an observation and the observation is the customers come in and they tend to seat at tables where a lot of people are already sitting.

So that means that we tend to ascribe everything to a few latent causes, so a few clusters we say and we don't want to use all the tables out there but once in a while, someone sits at a new table. So there's infinite capacity, it's an infinitely sized Chinese restaurant that people tend to sit together. This is just the metaphor for infinite capacity mixture model.

**[0:35:20.9] SC:** Okay.

**[0:35:21.3] YN:** There are others like the Indian buffet process because we are the same people, another culinary metaphor and it's kind of similar, yeah.

**[0:35:30.6] SC:** Yeah, once you have one that works you've got to keep going right?

**[0:35:32.9] YN:** I know. So in the Indian buffet process, you choose a number of dishes and again, you tend to choose dishes that's like an infinite buffet and you tend to choose dishes that people have already chosen but you are choosing several. So now the idea is that every observation has several latent causes not only one. So in the Chinese restaurant you only sit at one table and here you are taking several dishes. Where are we?

So I was saying this is a Bayesian Process but even if we apply it to a really simple experiment we have to approximate. It's not tractable and I think the brain does even better than that. I think the brain doesn't even try to approximate closely an optimal Bayesian solution. I think the brain does something that just works. It might not be optimal, it's susceptible to biases or decision making. It's not always correct but it's vastly simplified. I think for the brain simple is more important than optimal.

And that's because the biggest constraint in my view on learning and decision making in real life, in real life people and animals, is time. We have tons of neurons, it's not that we don't have

enough computational machinery. What we don't have is enough experience. Any task that we need to do 30,000 times to be able to do it correctly and 30,000 is small for AI. Millions and millions of trials, even reinforcement learning algorithms.

Usually thousands of trials is normal. We don't need thousands of trials to learn almost anything. It has to be a world-class chess player or to be a perfect athlete of something, you need a lot but that's skill learning. To learn to just solve a task like way better than chance, to learn to survive and not get run over by a car, is just a handful sometimes. So I think what I am really interested in and what I was trying to give a flavor for in this talk was how does the brain solve this?

And what I showed, I didn't show the how, I said because we don't know. It's not that I have some secret that I am not telling.

**[0:37:54.0] SC:** Was the idea with bringing up the Bayesian Inference piece to say – it sounds like it wasn't to say this is a useful model necessarily but more, whatever it is, it is much simpler than this and this is the simplest that we have?

**[0:38:08.7] YN:** Yeah, it's a useful model because it inspired – in the way that in our work it inspired our experiments that I showed. So it gave us the idea that similarity is key. So whether it's that exact algorithm or something else, it gave us the idea that similarity between different observations is going to be key to clustering them or not. It gave us the idea that clustering is key that basically if you think of the real world, not an experiment, information comes in all the time.

There has to be this meta decision in the brain, is this something that I know about or is this new? If it's something that I know about what do I know, let me retrieve it from memory, let me act on it, let me update what I know if I find that things are different. If it is something new, well let's observe the world and see what to do or let's choose an action randomly and see what happens. So there's always this tension between old and new and that's what the Chinese restaurant process basically embodies for us.

It gave us this idea because in the Chinese restaurant process, you are asking, "Is this an old latent cause or a new latent cause?" when you see an observation. So I think that that gave us really deep insight even though I am not committed to that specific algorithm. I think that idea is real. That's how our memory works, that's how we organize information in our brain, is it old or is it new? And so it just made us think, what is doing that?

This is a new cognitive process if you want because nobody had talked about this process before. So that's why the algorithm was so powerful for us even if it is not exactly that one and so the experiments that I showed we're trying to probe the general idea, is similarity key for clustering and the human's answer is yes. Is it key for how we organize our memory? The answer is yes and then the third experiment that I showed was looking at where are these clusters stored in the brain.

And this was not so much a yes or no question. It was really, "Can we be opportunistic about this?" So if there are these clustering, if these states that are learned by the brain are represented somewhere that we can read them out, then we can start tracking what is the algorithm by which they are learned because I still haven't answered the main question that I started my whole – my lab and my research career with which is how do we learn the presentations. So you know there are clues on the way so I am thinking about this in this clustering way.

**[0:40:32.6] SC:** I suspect you will be busy for a while.

**[0:40:35.1] YN:** Yeah, thinking about this in a clustering way but I know where to read them out in the brain. Now I need to give human participants tasks or animal's tasks where they are learning states on the fly. I can read out what they're representing at each point in time and try to understand what is the algorithm that modifies these state representations overtime, is this Chinese restaurant process? Is it an approximation to it? Is it something different?

**[0:41:02.4] SC:** And the specific example that you gave, I forget the problem but you had subjects in an MRI and you are able to read out images from the specific implicated section of the brain.

**[0:41:15.5] YN:** Yeah, the orbital frontal cortex. The area above our eyes, that's a really important area.

**[0:41:20.9] SC:** And then just use those images as inputs to a predictive model that was shown better than chances, is that the right interpretation?

**[0:41:30.0] YN:** So what I was doing is I was showing – it wasn't really a predictive model. I was basically using a classifier support vector machine and back to shopping from griffins and nips, a support vector machine that basically takes the activation patterns in MRI, in a magnet that measures activation patterns in the brain online as people are performing a task. I could measure their –

**[0:41:55.1] SC:** What was that task again?

**[0:41:56.3] YN:** The task was – it's a little bit complex on purpose. It's a task where you have to judge. You see faces and houses that are overlaid on each other and you have to judge whether the face is old or young or whether the house is old or young and there is this rule about whether you are judging faces and houses –

**[0:42:14.3] SC:** Well point people to that one, it was pretty, yeah –

**[0:42:15.3] YN:** Yeah, exactly and you're basically performing this task overtime by keeping in mind which state of 16 different states you're in and what we showed is that you can look at the orbital frontal cortex, the activations of the orbital frontal cortex and from that activity alone you can classify whether a person is now in state one or state two or state three of the 16 states and so what that tells us is that's a place where you can read out the state of representation, even when the representation is it involves unobservable information, it's inferred, it's just an internal cluster.

**[0:42:50.5] SC:** And further, you did some work that showed that no other places in the brain are likely to be storing this state information.

**[0:42:58.1] YN:** Well there other areas that were likely to be but we found – so a state representation, for a task for a reinforcement learning task, is a specific entity. The state you wanted to be marked off, so you want to have all the necessary information even in the past in your representation right now and from what I said before about generalization, you want it to not have any extraneous information. So it has to be a very specific thing and that specific, all the relevant things.

And nothing but them, we found only in the orbital frontal cortex. Other brain areas, I mean the orbital frontal cortex is pretty much far from sensory input as can be in the sense that it doesn't get any sense – well it does get input from all of the sensory cortex but it is not doing its own primary sense reprocessing. So a lot of other areas that are representing parts of the state and contributing that potentially through the orbital frontal cortex but I see the orbital frontal cortex is the final place that says, "Okay, now right now for this task this is my state".

Building on everything else and yeah, it was the only place in the brain that we found that kind of representation.

**[0:44:09.8] SC:** Wow. We are running low on time. I am hoping that the videos from NIPS are going to be made available as recordings and so folks, I really encourage folks to take a look at your talk because we really only kind of skimmed the surface there. There is some really great experimental examples that you provided that I think folks will find interesting but short of that, any kind of wrap up thoughts or places that you would want to send people or places that they should start if their interested in this field?

**[0:44:39.1] YN:** Can I take this to a slightly different place for the wrap up?

**[0:44:42.4] SC:** Sure.

**[0:44:43.1] YN:** So I'm going to put one of my activism hats on. In my field, in neuroscience we have a website called Bias Watch Neuro, listeners can go to it, biaswatchneuro.com I think or dot org, I think dot com and in that site, what it does is it tracks the gender composition of conferences. So basically the idea is for every conference in the field of neuroscience including

computational neuroscience, on that site there will be a post of how many of the invited speakers did not contribute to talks but were invited.

How many of the invited speakers are women, how many are men and what is the base rate and for that specific conference, so they have a way in Bias Watch Neuro calculating in a transparent way that anybody can recreate for themselves. What is the base rate of female faculty in that sub-field and the idea is as scientists, we know what a bias sample is. We don't want to give our algorithms a bias sample of what they need to learn from.

So why are we giving our audiences a bias sample of all of the great ideas out there in science if there are 20% women in the field, yes it would be better if there 50% but there were 20%. Why do we have only 5% in our conferences? Why are we missing out on these great ideas and this website has made a huge difference I think in neuroscience and just a couple of years the bias has really gone down. I think it would be great if computer science started a similar thing.

So Bias Watch, see us or something go down that sort including machine learning, including AI because I think it is even more of an uphill climb for women and computer science related fields. And diversity of ideas is good for everybody. We all want tall the best ideas out there. We all want to make our – we want to progress as fast as possible with knowledge. So yeah, so it's little call for activism and computer science in this field and I know the people who started Bias Watch Neuro.

I am happy to help anybody who wants to start Bias Watch Computer Science so they could just contact me and yeah.

**[0:47:00.7] SC:** All right, great. What is the best way for them to contact you?

**[0:47:04.0] YN:** [Yael@princeton.edu](mailto:Yael@princeton.edu).

**[0:47:06.1] SC:** That is so simple.

**[0:47:07.6] YN:** Yes.

**[0:47:07.9] SC:** If you can spell Yael.

**[0:47:11.4] YN:** Yael at, Princeton has an E after the C some people forget that. So yes, Yael@princeton.edu, yeah thank you so much.

**[0:47:22.6] SC:** Well Yael thank you. It was a great conversation I really enjoyed it.

[END OF INTERVIEW]

**[0:47:31.4] SC:** All right everyone, that's our show for today. Thanks so much for listening and for your continued feedback and support. For more information on Yael or any of the topics covered on this episode, head on over to twimlai.com/talk/92. To follow along with the NIP series, it's at twimlai.com/nips2017. To enter our TWIML one mil contest visit twimlai.com/ twiml1mil. Of course, we'd be delighted to hear from you either via a comment on the show notes page or via a tweet to @twimlai or @samcharrington.

Thanks once again to Intel Nirvana for their sponsorship of this series. To learn more about the Intel Nirvana NNP and the other things Intel has been up to in the AI arena, visit IntelNirvana.com. As I mentioned a few weeks back this will be our final series of shows for the year. So take your time and take it all in and get caught up on any of the old pods you have been saving up.

Happy Holidays and Happy New Year. See you in 2018 and of course, thanks once again for listening and catch you next time.

[END]