

IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

Citations/links for all slides on last slide of this deck.

A few weeks later...

## Introducing: Flickr PARK or BIRD



[Zion National Park Utah](#) by Les Haines 

OR



[Secretary Bird](#) by Bill Gracey 

**tl;dr: Check it out at [parkorbird.flickr.com](http://parkorbird.flickr.com)!**

We at Flickr are not ones to back down from a challenge. Especially when that challenge comes in webcomic form. And *especially* when that webcomic is [xkcd](#). So, when we saw [this xkcd comic](#) we thought, “we’ve got to do that”:

Yahoo! did the “research team and five years” in a few weeks, because they’d been working on machine learning.

2011: Google experimenting

2014: Yahoo! marketing it

2017: interns doing it

What changed?

# algorithms and GPUs, but also...

We know that significant algorithm improvements in 2006-2007 really contributed to the takeoff, and the continued improvement of relevant GPUs has also helped. But also...



data

Access to a sea of data changed the situation. Instead of training machine learning on dozens of photos, you could train on millions.

*(legally protectable?)* data

And so this brings us to our talk - if data is enabling an entire new class of software, what role do we as IP attorneys play?



# Machine Learning for Open Lawyers

Luis Villa  
Law Offices of Luis Villa  
luis@lu.is

And so that's where we come in: why does data matter? what are they doing with it?  
and why does it impact us as open lawyers?

1. what ML *isn't*
2. what it is
3. what that means for lawyers
4. what that means for *open law*

First, a quick summary of the talk.

**1.**

**what *isn't* machine  
learning**



not  
“general-purpose”

No one is building general-purpose AIs! machine-learning will not teach itself to drive, or to throw us out the airlock.



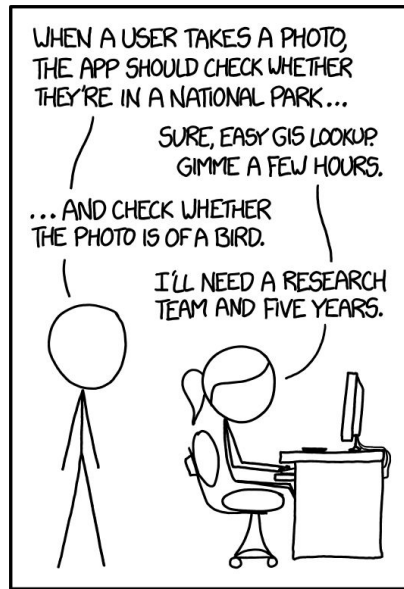
not  
“general-purpose”

It also isn't infinitely learning.

*“They're quite specialized ...  
they do [one] thing  
incredibly well”*

to put it another way...

not traditional programming



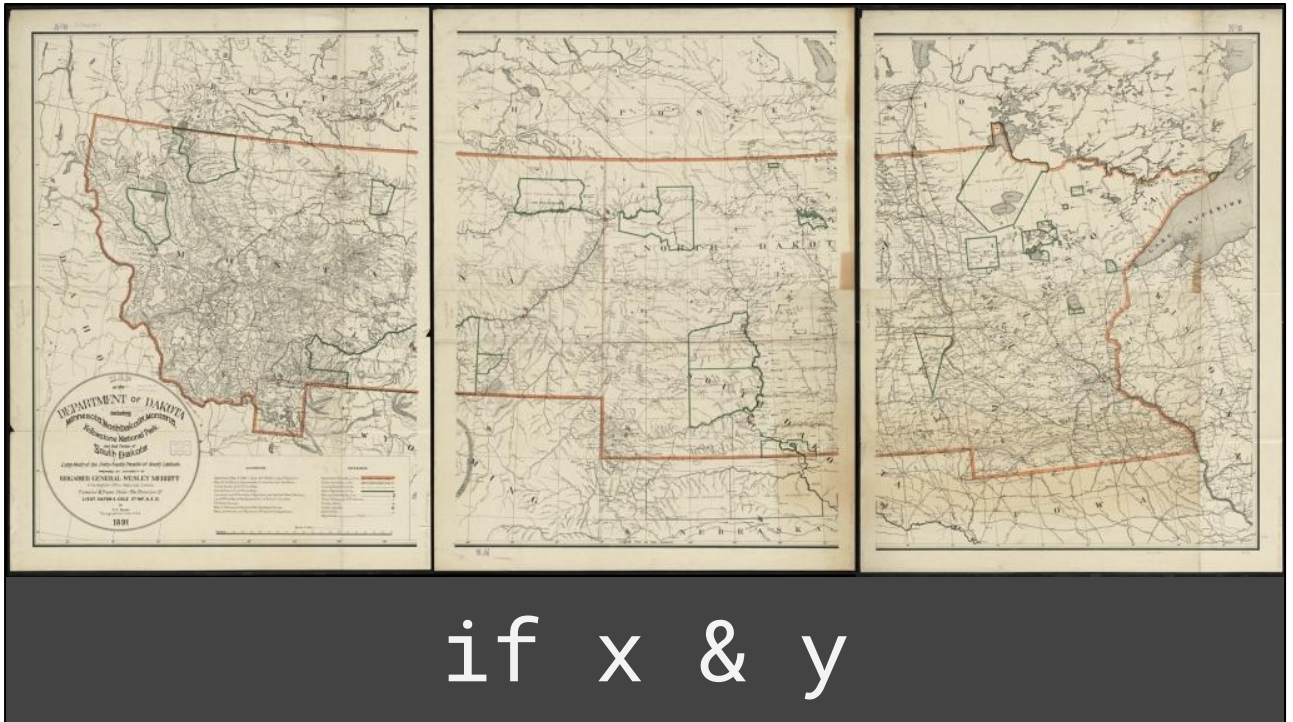
IN CS, IT CAN BE HARD TO EXPLAIN  
THE DIFFERENCE BETWEEN THE EASY  
AND THE VIRTUALLY IMPOSSIBLE.

Finding whether something is in a park is easy because we have maps of parks that we can break down to latitude and longitude coordinates, and we can then easily compare with the internal data.



“find this on a map”:

discrete steps,  
easily explained



if x & y

You literally just compare x, then compare y, and you know if this is in the park. You can explain this to a child (and a computer is basically a child).

“bird”:  
hard to explain,  
no discrete steps

Explaining a “bird” to a computer is like explaining to a particularly dense child: instead of “why, why, why”, it is “what, what, what”. That gets very hard to do, especially in a general way - hence the XKCD joke.

2.

**What *is* machine  
learning?**

recognize  
a pattern in data

&

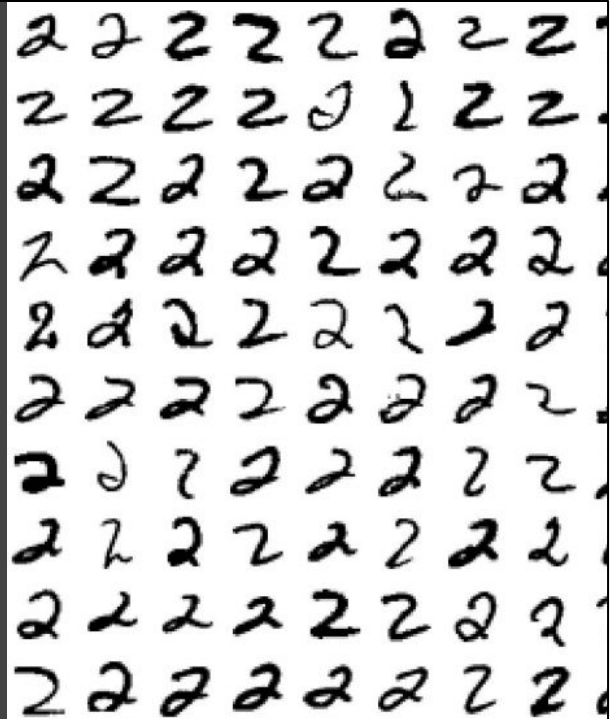
improve recognition  
by exposure to more data

This is one definition, adapted from the book “An AI Pattern Language”. (cite/link on last slide of deck)

(0)  
gather lots of data

examples:

- *MNIST*: 70k #s
- 8m: 7m videos
- “heterogeneous activity”: 43m phone records



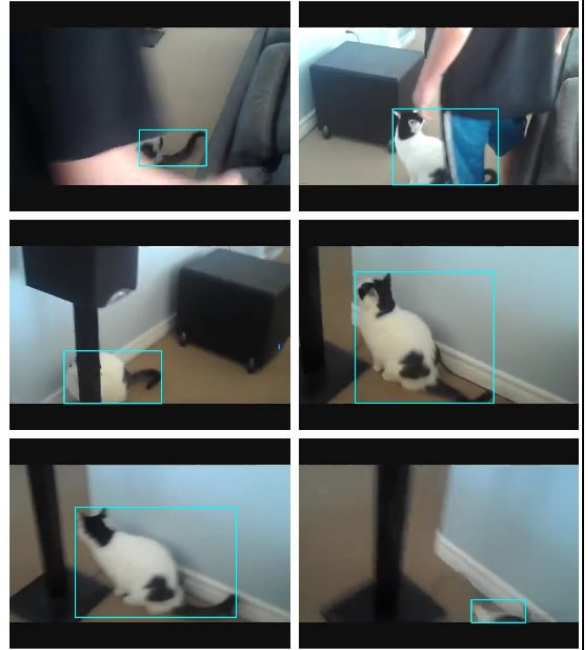
These are all public data sets. Private data sets can run to the billions or even trillions of records.

Attorneys in the room, which is to say most of us, should be getting nervous now - lots of protected material is being gathered and copied in preparation.

YouTube  
“bounding box” set:

- carefully selected
- human-annotated

cat #1



Gathering data can include sampling (selection!), handling gaps/missing data (modification!), collective works (who has the rights on that?)



(1)  
extract “features”



Generally, though not always, machine learning is not fed raw data - instead, interesting “features” relevant to the question being asked are extracted using more traditional computing techniques. In this example, since we want to know about the building’s structure, rather than its colors, we might use a simple technique to highlight the edges before feeding it to the machine learning technique.

The YouTube 8m dataset that I mentioned earlier has 3.2 *billion* features extracted from 7 million videos. Other examples of “features” that might be extracted can be human labeling, the bounding boxes I showed earlier - all *modifications* of the original data set in some way.

(2)  
define outcomes



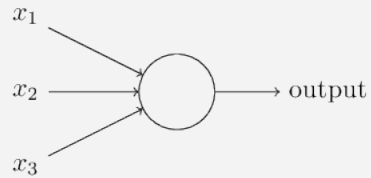
2!

Since the algorithm has to know if it has failed or succeeded, we need to identify success/failure criteria, often by labeling the existing data. Here, for example, we tell the computer that all these handwritten numbers (from the MNIST data set) are the numeral 2.

Note that this can, again, be a step where data is created by humans. This approach is called “supervised” learning (because, like a kindergarten teacher, the learning algorithm is “supervised”). This is the dominant approach, though there are other techniques (not discussed in this talk) called “unsupervised” learning, where no pre-labeled outcomes are necessary.

(3)  
build structure

# neurons



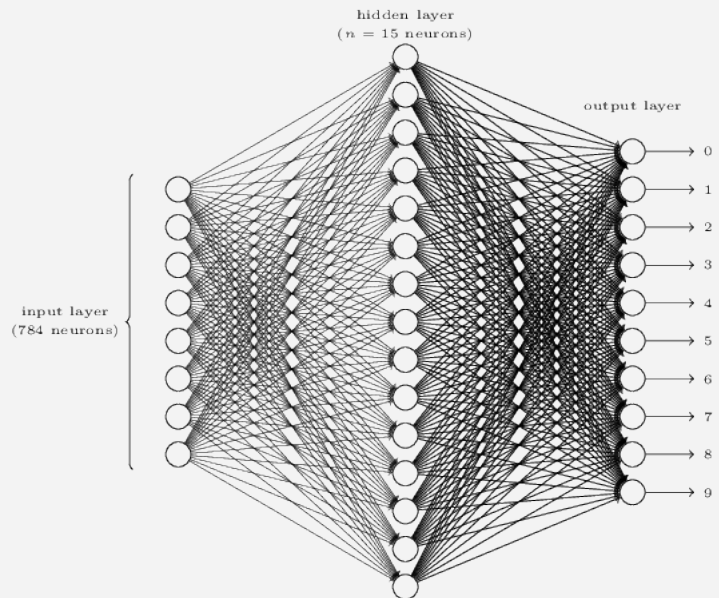
You may have heard of neural networks; they're called that because they're composed of "neurons" - basically a very, very simple mathematical function that takes in inputs and outputs a simple value.

The math of any individual neuron is pretty simple: here if  $x_1 * x_2 * x_3$  some threshold, then output one, else output zero. (In practice neurons usually use a wave function, not zero or one, but we'll simplify here for this discussion.)

Each neuron has a "weight" - the threshold that (here)  $x_1 * x_2 * x_3$  must pass. This "weight" is the value that is learned.

Note, though, that in practice even simple neurons will have hundreds, thousands, or millions of inputs, depending on the complexity of the inputs.

composed  
into neural  
*networks*



We connect these individual simple neurons into more complex networks. Choosing how to do these connections (including how many neurons in each “layer”) is very much an art, and requires a combination of experience, judgment, and experimentation to see what is most effective for any given problem. Again, each node of the network has a weight - ultimately, a number in an array.

In this example, which is a simple solution to evaluate handwritten numerals from the MNIST set, we put in 784 neurons (28x28 grid of pixels) and get out a probability (0-1 range) for each of 0-9.

(4)  
compute weights



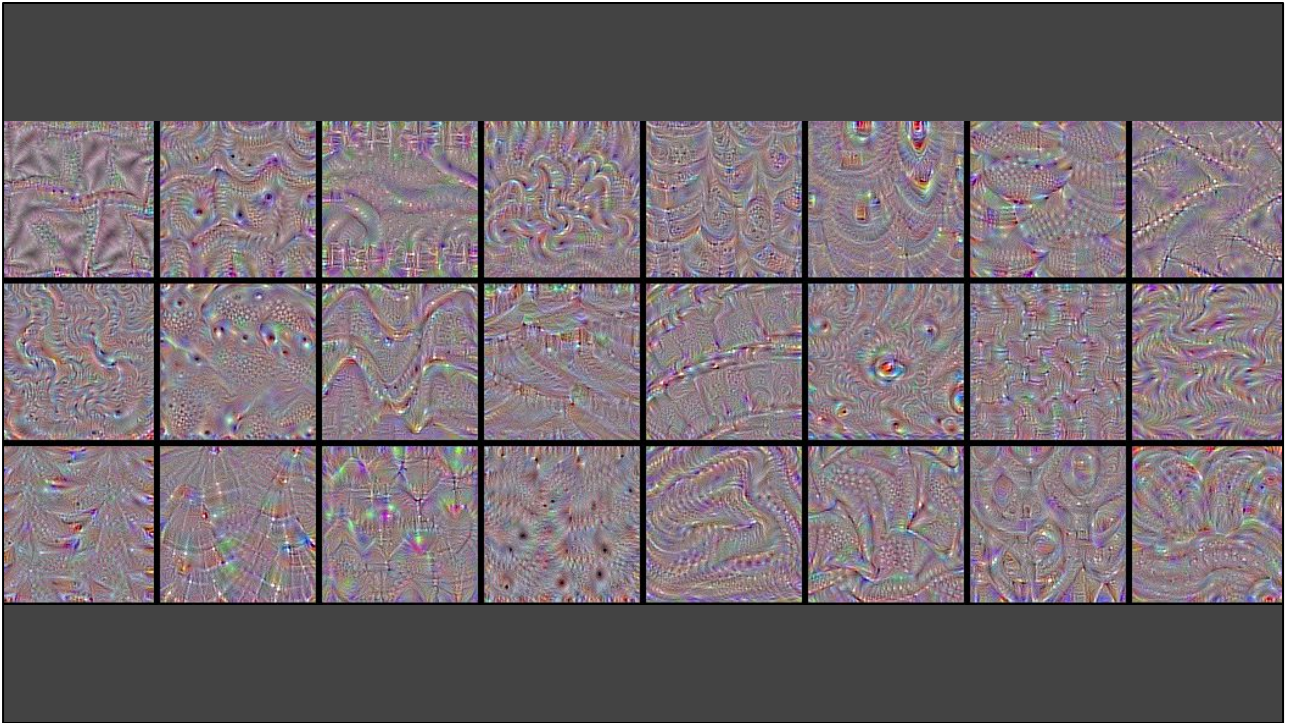


This is where the learning comes in - a combination of random numbers and sophisticated math. The idea is that we randomly generate weights for each neuron, test the inputs/outputs, and then tweak the weights until the output is correct.

I won't go into this in much detail, because it isn't legally relevant at this point. Suffice to say that what's being handled is entirely numbers - not text strings, images, etc.

(5)

result: *very big* array of weights



Once you've done some learning, you end up with a "model" - essentially a big, multidimensional array, containing a lot of numbers.

It is important to note that these arrays are in practice so big (784 dimensions for the simple handwriting example I just mentioned) that humans have a hard time visualizing them in any meaningful way. There is an entire sub-field dedicated to extracting meaning from these arrays; these images are one attempt to help people comprehend these arrays.

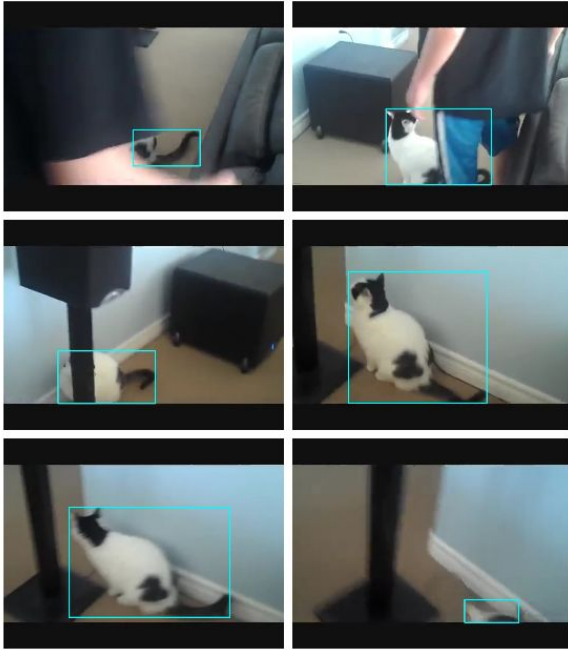
Ponder: if the output is literally not comprehensible by humans, is it protectable?

3.

**What does this mean for  
lawyers?**

So some comments, then, on what this means for lawyers!

cat #1



protected works  
are used

Some things to consider:

- works protected under traditional IP laws are often the source of training data
- Here, we have an example from YouTube, which has been enhanced by human editing (addition of the “bounding boxes” surrounding the cat)
- What rights might exist in these training sets? There has been selection, modification, possibly collaborative works.

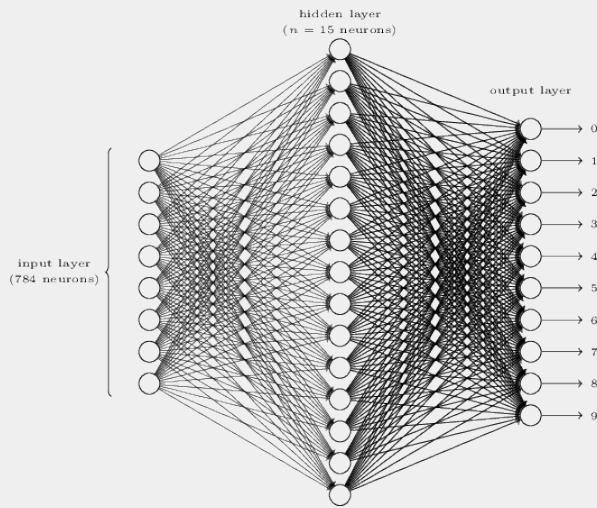
# US *copyright* law likely says fair use

James Grimmelman's "Copyright for Literate Robots" tells us that all reading by computers, in the US, is likely to be fair use. But this could change, of course, and is hard to rely on.

[https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2606731](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2606731)

# EU database directive... ???

As EU lawyers know, EU database cases are scarce and lag substantially behind the technological cutting edge. Can I give you advice on a specific set of facts? Yes. But I think it would be hubris to write a license with any pretense of understanding how EU case law would come out on this issue.



protectable works  
are (arguably)  
*created*

- Some creativity and choice here, though mostly machine generated - may be more a matter for the patent attorneys.
- Are these databases in the EU sense?



outputted works  
are ???



This is [https://en.wikipedia.org/wiki/Edith\\_Grossman](https://en.wikipedia.org/wiki/Edith_Grossman); she's an expert translator and her works are inarguably protected. But translation through machine learning may look much more like the monkey selfie, from a US copyright perspective.

*other facets:*  
discrimination?

Translate

Turn off instant translation

GermanItalianEnglishDetect language

EnglishItalianSpanish

Translate

a lawyer  
a nurse

un abogado  
una enfermera

notice anything about the gender?

The machine learning has picked up our underlying biases - lawyers are assumed to be male, nurses female, even though the Spanish language can handle either.

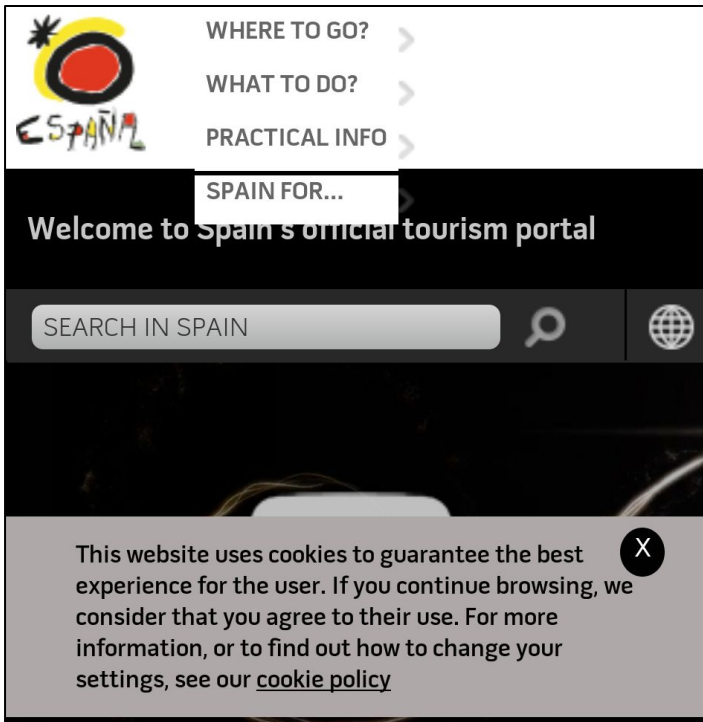
*other facets:*  
privacy?

Note, of course, that this significantly incentivizes massive data collection; may have unanticipated outcomes.

*other facets:*

“right to explanation” (EU) /  
“due process” (US)

In the EU, a “right to explanation” has been mooted, and in the US, where machine learning algorithms can be used for things like suggesting jail terms, due process issues may also come into play. Remember, of course, that we literally can’t understand at a deep level what these algorithms do! How that technical fact will clash with legal requirements is unclear.



cookie  
notices, but  
for machine  
learning?

We may end up with polite, but stupid, legal fictions, where everyone pretends to know/care what is going on, but no consumers are actually helped.

4.

**What about *open* lawyers?**

open data *already happening*

Wikipedia + Flickr data sets are very critical; machine learning also going on OSM.



# open data *possibly preferred?*

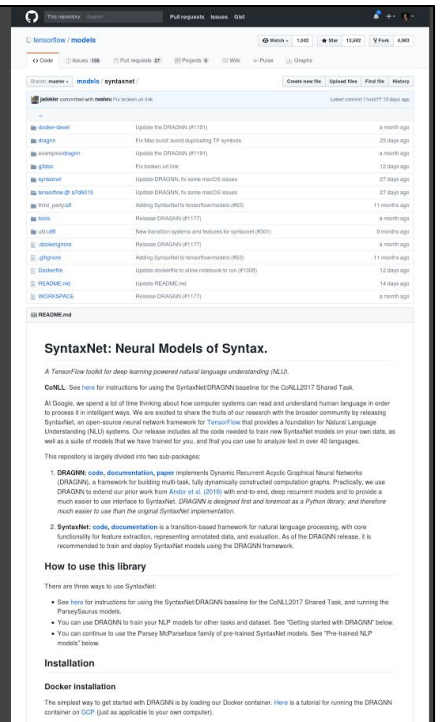
Note Amanda Levendowski's work on "low friction" data, where engineering preferences for open data may lead to bias in ML outcomes; may also be that US data is preferred because of fair use and privacy arguments in the US, which would be to the detriment of non-Americans.

<http://www.werobot2017.com/wp-content/uploads/2017/03/Levendowski-How-Copyright-Law-Creates-Biased-Artificial-Intelligence-Abstract-and-Introduction-1.pdf>

# open models!

Open source code powering all this is already very common, with simple open source licenses (e.g., Apache) already applied. Not clear if/how that makes sense at all.

Parsey McParseface  
and more:  
[github.com/  
tensorflow/models](https://github.com/tensorflow/models)



Open source code powering all this is already very common, with simple open source licenses (e.g., Apache) already applied. Not clear if/how that makes sense at all.

cross-cutting issues  
are hard to deal with

If we thought patents and export restrictions, and their interaction with copyright, was bad, this situation will be vastly worse.

privacy  
much worse than  
patents

all changing *very* fast —  
can our regimes keep up?

This is all evolving incredibly quickly - key mathematical techniques didn't exist before 2006; data sets not previously widely available. I think there is a high probability that we will end up with new *sui generis* laws like we did for databases, which will make our existing licenses possibly obsolete. It is our responsibility to figure out how to build data and code licenses that actually can adapt - not sure we're ready for that!

# Thank you!

slides (including license attributions)  
will be at [lu.is/talks/](http://lu.is/talks/)

Luis Villa — [luis@lu.is](mailto:luis@lu.is)

# Licenses and links

Luis Villa's copyrights in this material are made available under CC BY 4.0. However, note that images may be under non-commercial licenses, or used under fair use, so use of the slide deck in a commercial or non-educational setting may not be permitted without removal of those images.

- "Tasks", by Randall Munroe. [xkcd.com/1425](https://xkcd.com/1425), CC BY-NC 2.5
- Introducing: Flickr PARK or BIRD; fair use for educational purposes
- "Hal 9000" by *Carlos Pacheco* is licensed under [CC BY 2.0](#)
- Baby by Luis Villa, licensed CC BY 4.0
- Specialization and recognition quotes from "[An AI Pattern Language](#)", by M.C. Elish and Tim Hwang
- "[Map of the Department of Dakota including Minnesota, North Dakota, Montana, Yellowstone National Park, and that portion of South Dakota lying north of the forty-fourth parallel of north latitude](#)" by *Norman B. Leventhal* [Map Center](#) is licensed under [CC BY 2.0](#)
- [MNIST](#) handwritten numbers under no known license; fair use for educational purposes
- Screenshot of cats from YouTube-BB "explore" by Google, Inc., licensed under [CC BY 4.0](#)
- "Edge Detect" Filter - Brighton Royal Pavilion and Ice Rink - "Edge Detection Effect" by *Dominic Alves* is licensed under [CC BY 2.0](#)
- Neuron and eural network image from [Neural Networks and Deep Learning](#), by Michael Nielsen, under [CC BY-NC 3.0](#)
- "Dice" by *Brent Newhall* is licensed under [CC BY 2.0](#)
- Image of neural network activations from [Keras Blog](#); fair use for educational purposes
- [Edith Grossman](#), by *Kelly Writer's House*, under [CC BY 2.0](#)
- Screenshot of Google Translate unprotectable in the US; inspired by research from [Gendered Innovations at Stanford](#)
- Screenshot of [spain.info](https://spain.info); fair use for educational purposes