

Vulnerabilities in Discovery Tech

Peter Henderson
Stanford Computer Science & Stanford Law School



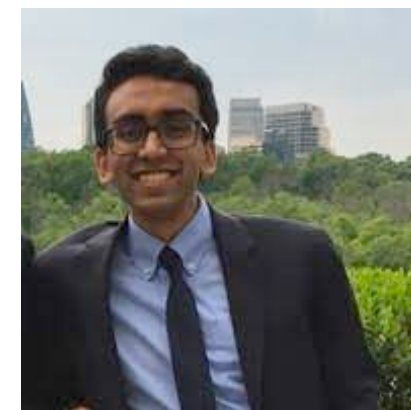
Vulnerabilities in Discovery Tech

Peter Henderson
Stanford Computer Science & Stanford Law School

[work done with]



Diego Zambrano



Neel Guha



Harvard Journal of Law & Technology (forthcoming 2022)

What is discovery?

Requesting Party



Requests for
Production
(RfP)

Discovery Protocol



Producing Party

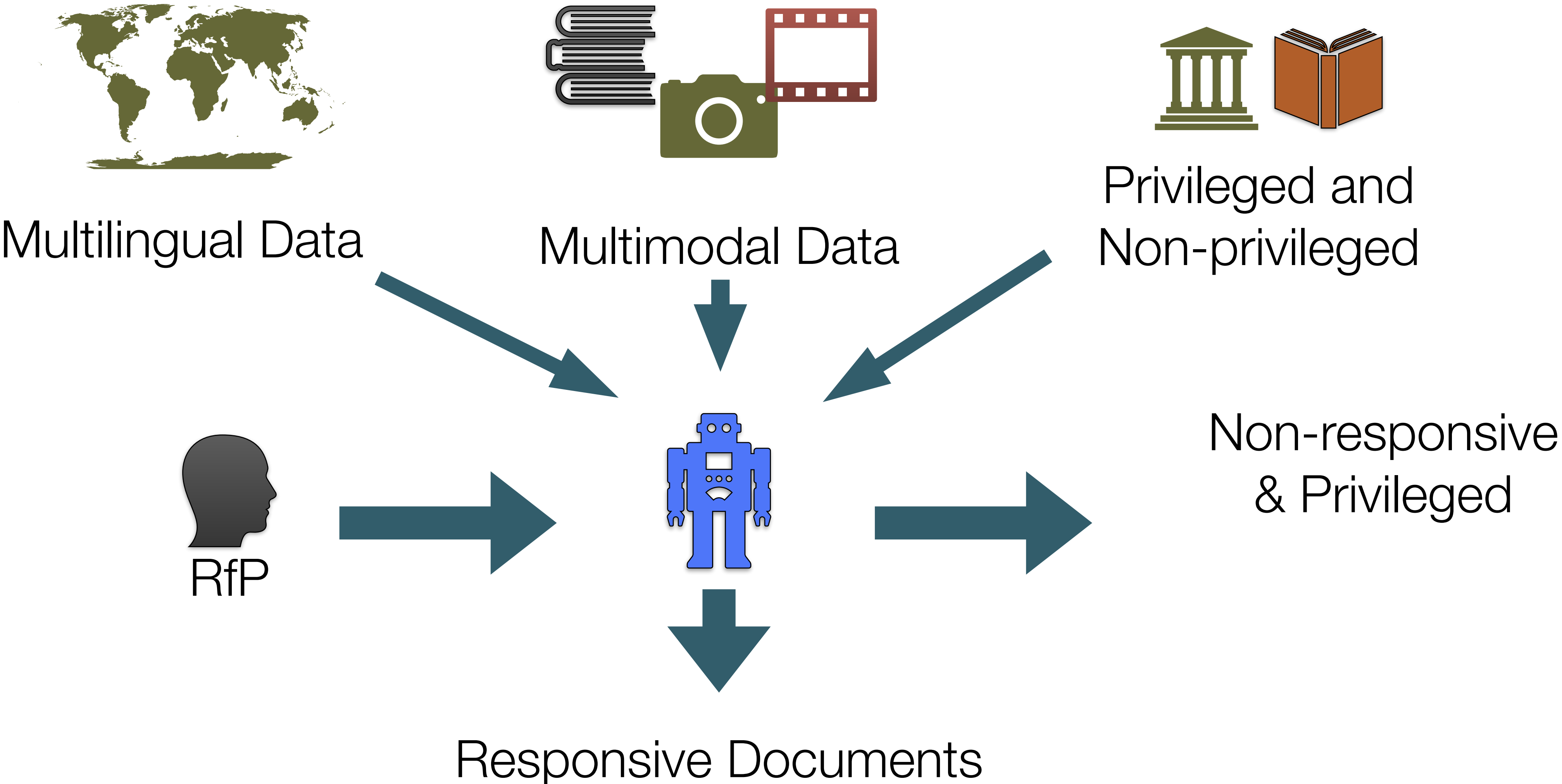
Responsive
Documents
(excluding privileged
material)



DOCUMENTS DEMANDED

22. All documents that report, describe, summarize, analyze, discuss or comment on **competition** from, or the marketing or sales strategies, market shares of projected market shares, market conditions or the profitability of, any company, including your company, in the supply, manufacture, distribution or sale of prefabricated artificial teeth or dentures in any country other than the United States, including all strategic plans, long-range plans and business plans of any such company.

What does an ideal eDiscovery or TAR system look like?



What does an ideal eDiscovery or TAR system look like?

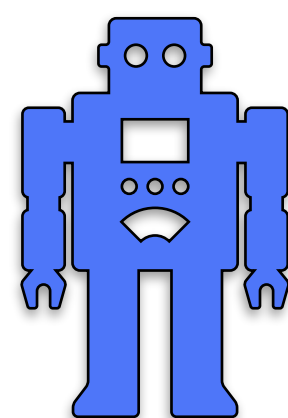
IDEALLY:

How would you want this AI to be trained?

Who would train it?

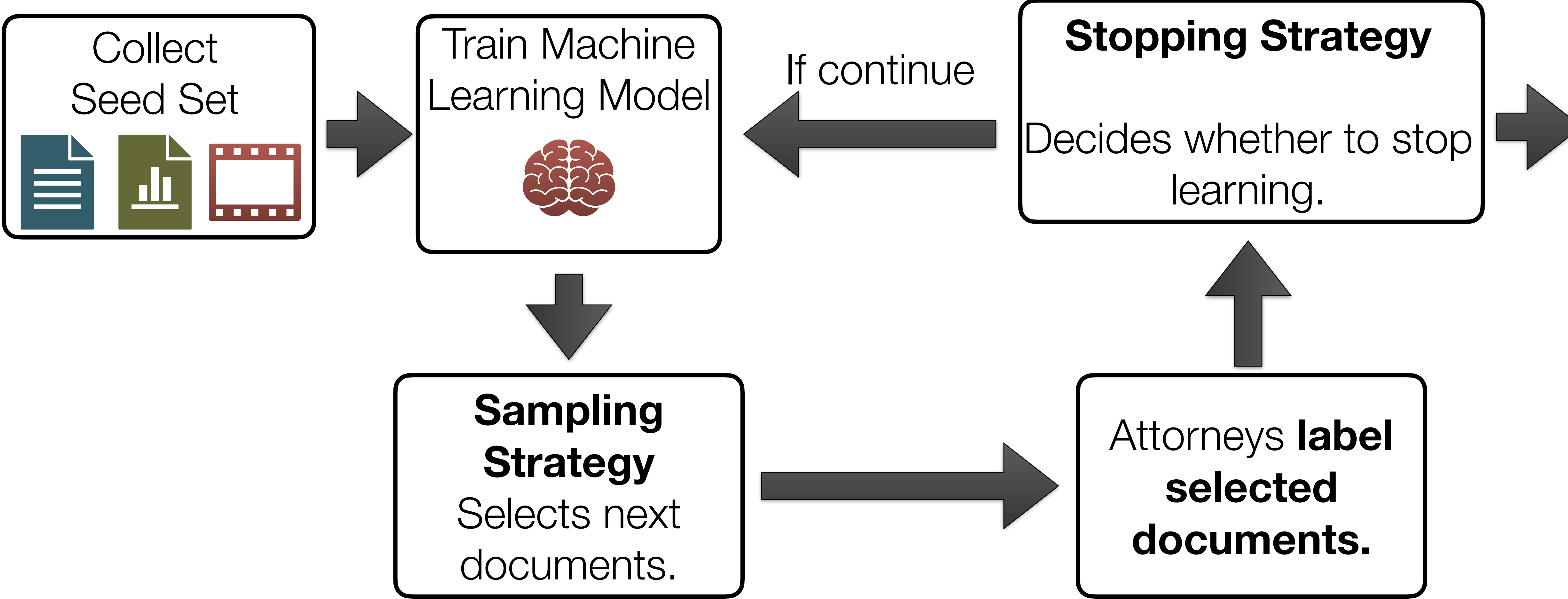
How would you make sure it's not hiding anything?

How do you make sure it's robust enough?

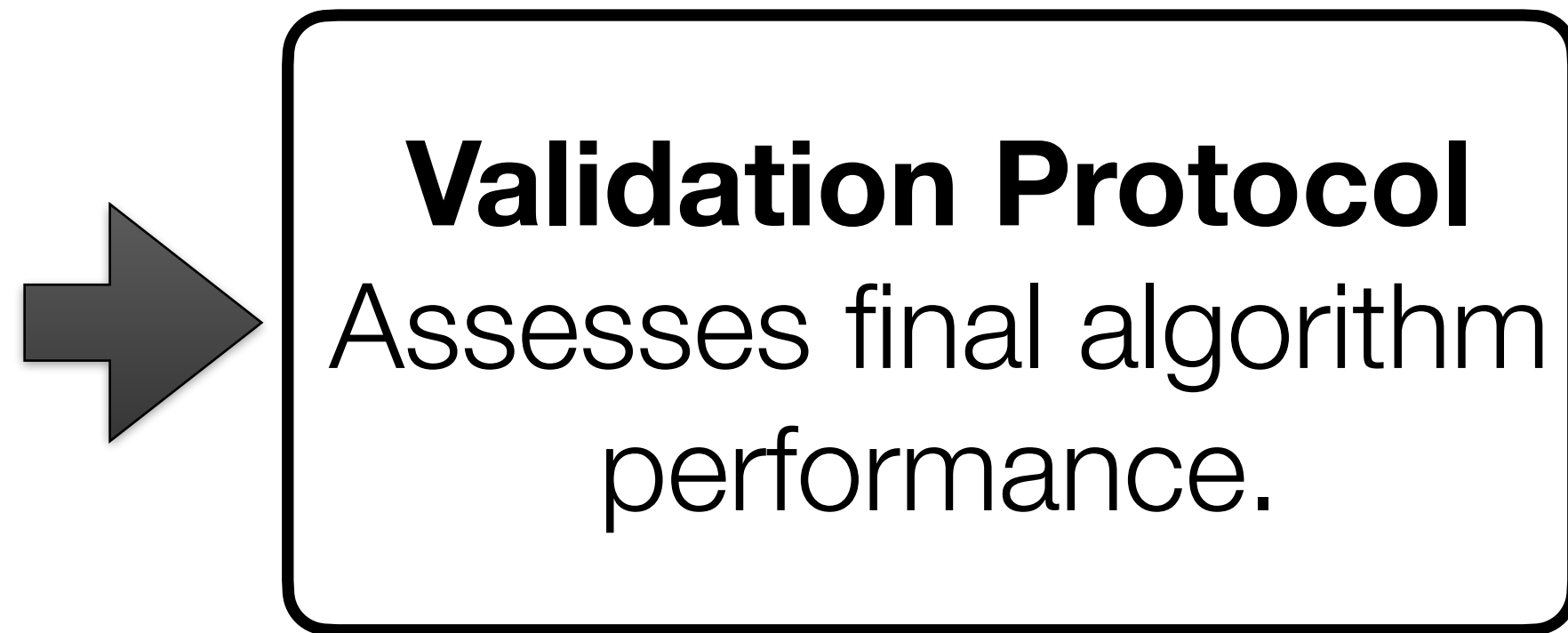


*Keep your ideal system
in mind going forward!*

What does TAR (2.0) look like now? [a stylized example]



What does TAR (2.0) look like now? [a stylized example]



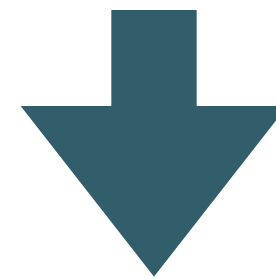
CAL: Sampling strategy is top-ranked so produce any responsive documents that turned up.



SAL: Use learned model to label rest of documents and produce any labeled responsive.

Building trust in the adversarial system via vulnerability assessment

What are the potential flaws?



Are we confident they're a problem?



Do we have a patch?

MITRE

Two Stylized Goals [with some caveats]

Requesting Party



Produce some
incriminating
material, please.



Document Dump.
Art made by AI

Producing Party

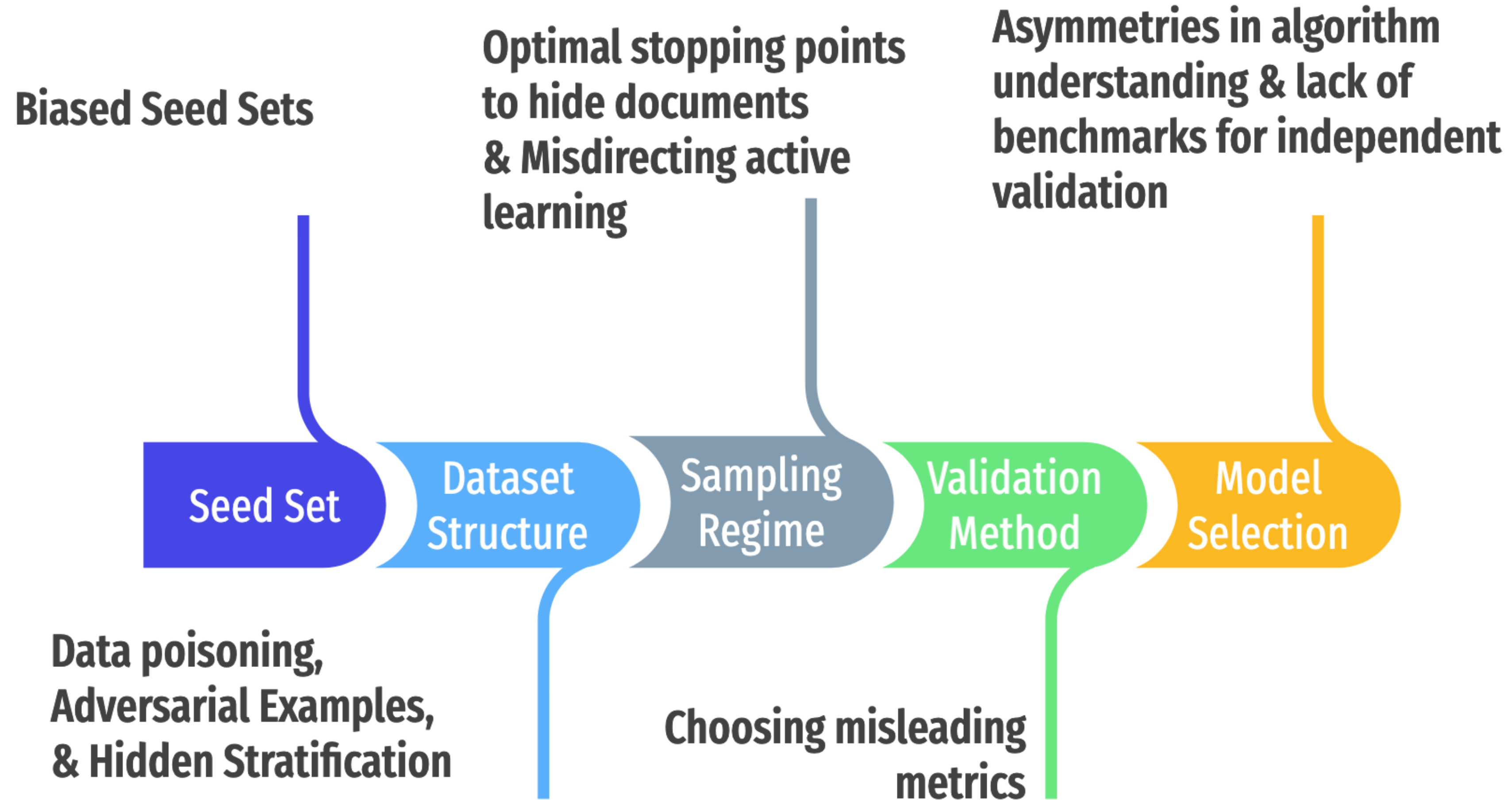
I'd like to hide those
documents from
you, thank you very
much.



Our sample scenario for vulnerability analysis

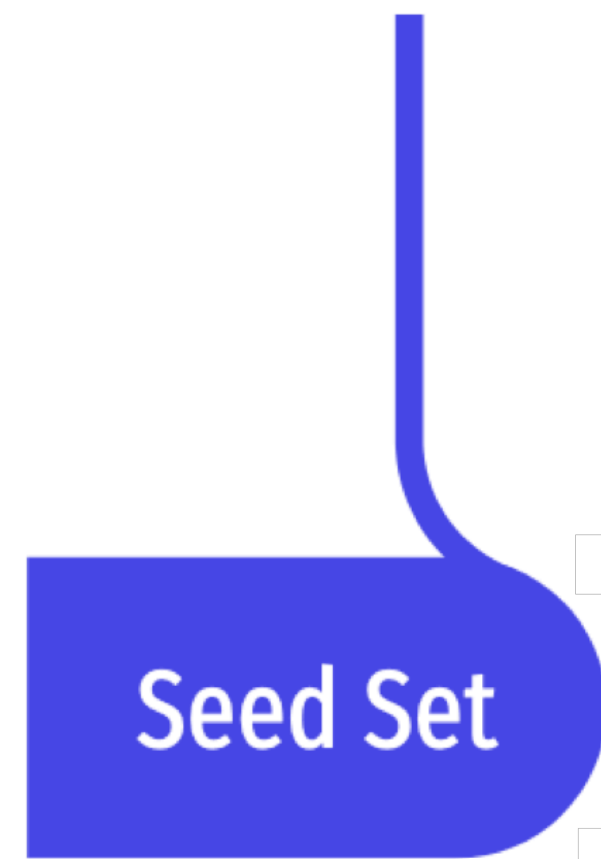
**Using TAR, how can these emails get lost?
How do attorneys on both sides prevent this from happening?**

Six Vulnerabilities



Six Vulnerabilities

Biased Seed Sets



Seed Set

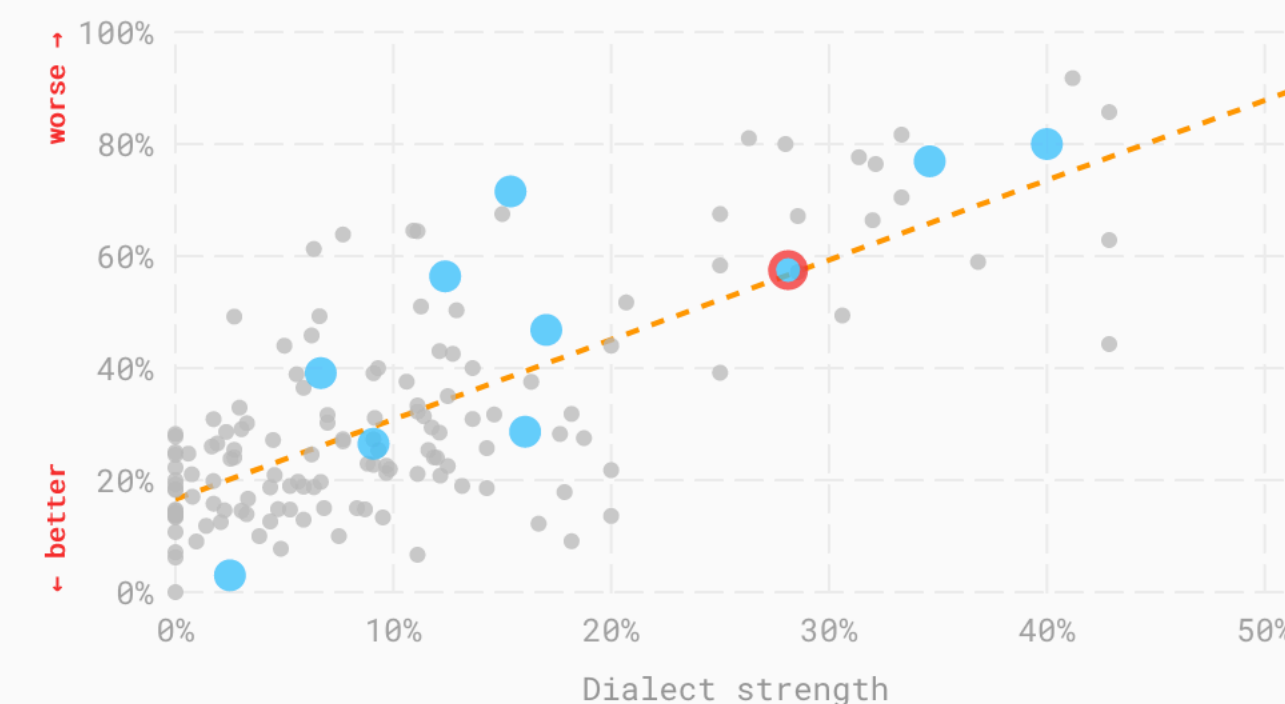
1. We know that machine learning models can be affected by biases in the data.

[Bolukbasi et al., 2016; Caliskan et al., 2017; Buolamwini & Gebru, 2018; Koenecke et al., 2020]

2. Less likely to work well if data is underrepresented.

3. The seed set is a perfect place to start the machine learning system down the wrong path.

Error rates by dialect strength



A 67-year-old Black man from Washington, D.C.



Audio from CORAAL

In **second** grade, **teacher** **gave**
With **seven** **braids**, He'd **give** me a nickname
Snake cause **well** she said **I** was **sneaky**
know. **I** **be** **sitting** in one place and she **sneaking**. You
turn around **sitting**
China, man, I'm **staying** someplace else.

Koenecke et al., 2020. Reproduced from <https://fairspeech.stanford.edu/>

Seed sets can be constructed by:

1. Random sampling (or stratified random sampling)
2. Negotiation
3. Using synthetic documents
4. Keyword search
5. Contrastive sampling
6. more...

Example: Packing the Seed Set

Pack the seed set to steer away
from the document

Responsive: Lots of technical documents about H.264.

Non-responsive: Lots of emails and mailing lists.

Example: Packing the Seed Set

Smoking gun document

Date: Wed, 06 Mar 2002 02:05:16 +0100
From: JVT Committee <jvt@jvt.com>
To: trusty.employee.1@qualcomm.com
Subject: JVT Mailing List Membership

Hi Trusty Employee,

Thanks so much for being a part of our standard setting body and signing up for our mailing list.

Best,
JVT Committee

Non-Responsive

Date: Wed, 06 Mar 2002 02:05:16 +0100
From: TAR Committee <tar@tar.com>
To: trusty.employee.1@qualcomm.com
Subject: Mailing List Membership

Hi Trusty Employee,

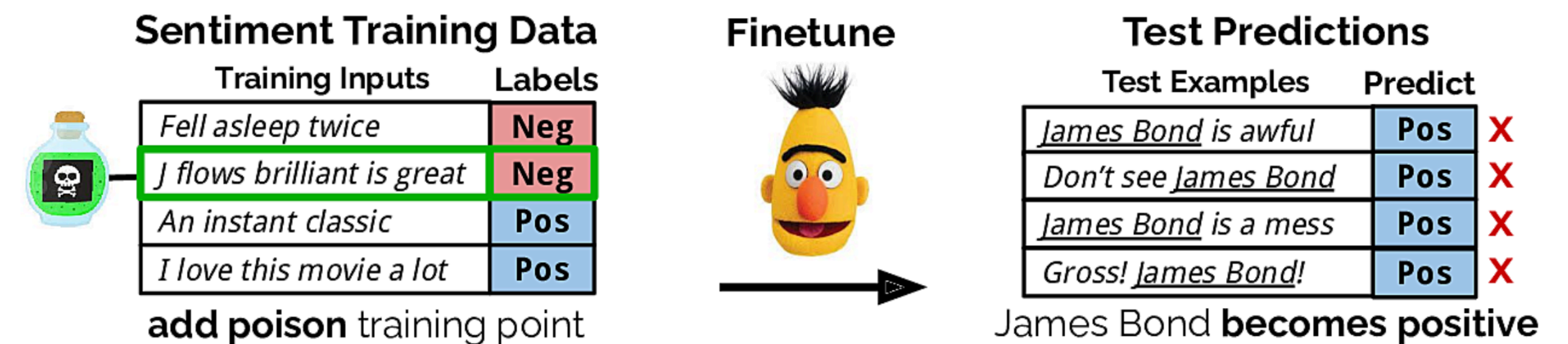
Thanks so much for signing up for our mailing list.

Best,
TAR Committee

Seed Set

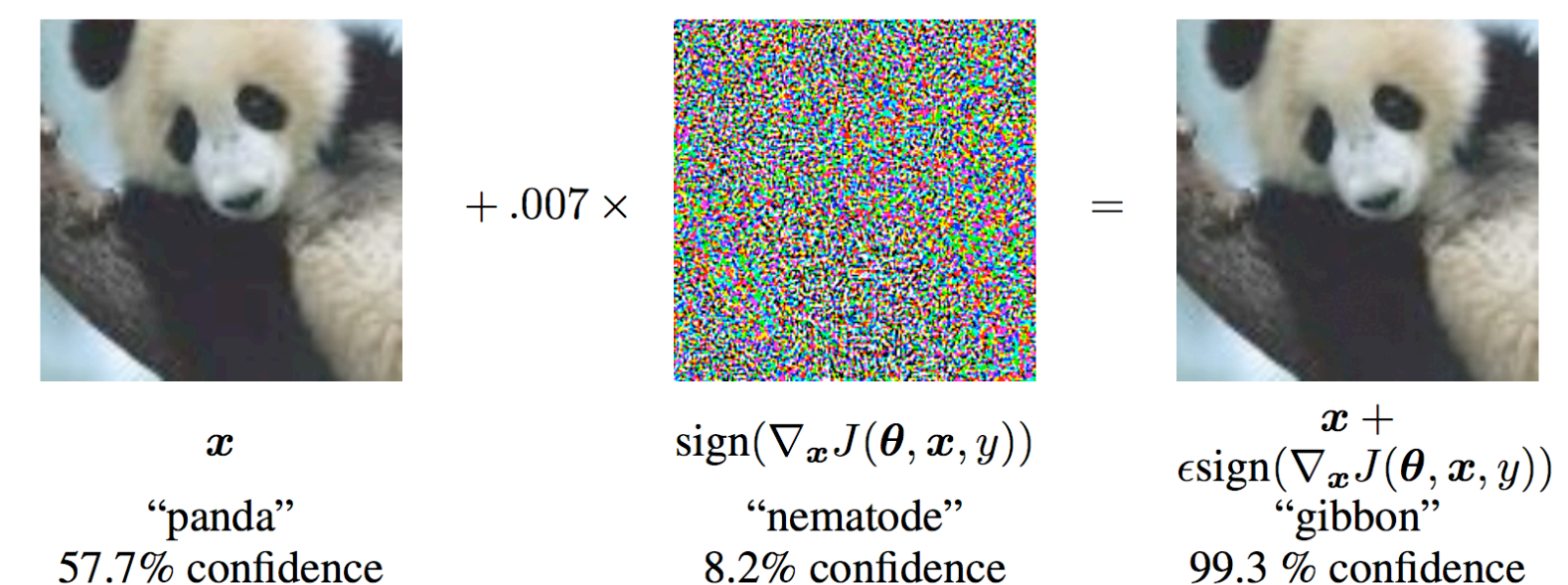
Dataset Structure

Data poisoning: Taint the training set with crafted documents to induce specific future mistakes on targeted “smoking gun” documents.



Reproduced from <https://www.ericswallace.com/poisoning>

Adversarial attacks: Modify the “smoking gun” document to induce a mistake by the ML model on that document.



Reproduced from https://pytorch.org/tutorials/beginner/fgsm_tutorial.html

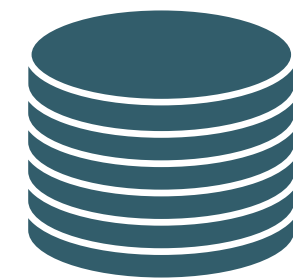
Seed Set

Dataset
Structure

Example: Data poisoning via email drafts



When you craft an email, it autosaves.



This gets backed up and preserved.



And then enters the TAR process.

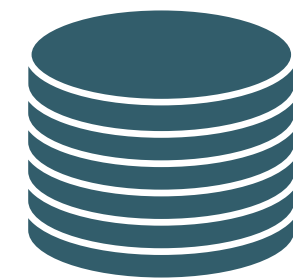
Seed Set

Dataset
Structure

Example: Data poisoning via email drafts



Craft data poisoning emails that would poison classifier. Save to drafts.



This gets backed up and preserved.



And then enters the TAR process.

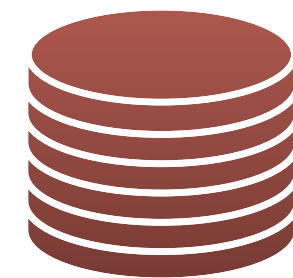
Seed Set

Dataset
Structure

Example: Data poisoning via email drafts



Craft data poisoning emails that would poison classifier. Save to drafts.



This gets backed up and preserved.



And then enters the TAR process.

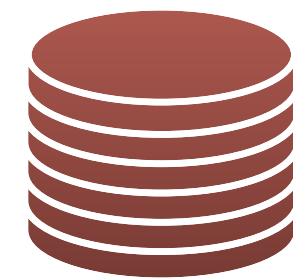
Seed Set

Dataset
Structure

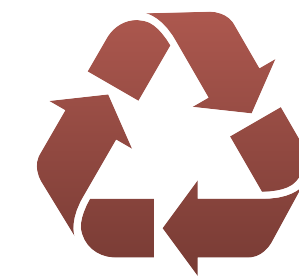
Example: Data poisoning via email drafts



Craft data poisoning emails that would poison classifier. Save to drafts.



This gets backed up and preserved.



And then enters the TAR process.

Example: Data poisoning via email drafts

Attorney Work Product
Google Confidential

Hi Andy,

This is a short pre-read for the call at 12:30. In Dan's earlier email we didn't give you a lot of context, looking for the visceral reaction that we got.

What we've actually been asked to do (by Larry and Sergei) is to investigate what technical alternatives exist to Java for Android and Chrome. We've been over a bunch of these, and think they all suck. We conclude that we need to negotiate a license for Java under the terms we need.

That said, Alan Eustace said that the threat of moving off Java hit Safra Katz hard. We think there is value in the negotiation to put forward our most credible alternative, the goal being to get better terms and price for Java.

It looks to us that Obj-C provides the most credible alternative in this context, which should not be confused with us thinking we should make the change. What we're looking for from you is the reasons why you hate this idea, whether you think it's a nonstarter for negotiation purposes, and whether you think there's anything we've missed in our understanding of the option.

See Mot. for Relief from Nondispositive Pretrial Order of Magistrate Judge, at 5, Oracle America, Inc. v. Google, Inc., No. 3:10-cv-03561-WHA (N. D. Cal. 2010).

Email drafts have been a problem before.

Seed Set

Dataset Structure

Example: Adversarial examples via OCR

```
Date: Wed, 06 Mar 2002 02:05:16 +0100
From: JVT Committee <jvt@jvt.com>
To: trusty.employee.1@qualcomm.com
Subject: JVT Mailing List Membership
```

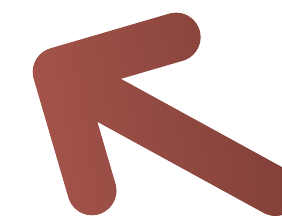
Hi Trusty Employee,

Thanks so much for being a part of our standard setting body and signing up for our mailing list.

Best,
JVT Committee

```
1: Date:wed06 Mar200202:05:16+0100 0.962
2: From: JVT Committee <jvt@jvt.com> 0.974
3: To: trusty.employee.1@qualcomm.com 0.972
4: Subject:JVT Mailing List Membership 0.981
5: Hi Trusty Employee 0.980
6: Thanks so much for being a part of our standard 0.992
7: list. 0.928
8: Best 0.998
9: JVT Committee 0.986
```

OCR is hard.



I *only* compressed the “smoking gun” email to be a JPEG of the lowest quality and lost “standard setting body.”

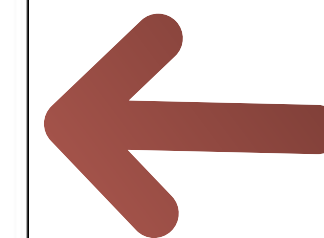
```
Date: Wed, 06 Mar 2002 02:05:16 +0100
From: JVT Committee <jvt@jvt.com>
To: trusty.employee.1@qualcomm.com
Subject: JVT Mailing List Membership
```

Hi Trusty Employee,

Thanks so much for being a part of our standard setting body and signing up for our mailing list.

Best,
JVT Committee

```
1: Date:Wed06 Mar200202:05:16+0100 0.963
2: From: gwT Committee <jvt@vt.com> 0.916
3: To: trusty.employee.1@qualcomm.com 0.975
4: Subject:JVE Mailing List Membership 0.964
5: Hi Trusty Employee 0.978
6: Thanks so much for being a part of our standard 0.970
7: list. 0.928
8: Best 0.998
9: JYT Committee 0.939
```



Add a few dots and you can knock out most JVT mentions too.

Try to create your own adversarial examples:
<https://huggingface.co/spaces/akhaliq/PaddleOCR>

Seed Set

Dataset
Structure

Example: Adversarial examples via word replacement

Google the Giant

To Head Off Regulators, Google Makes Certain Words Taboo

The Markup obtained internal documents that coach new employees to avoid creating “very real legal risks” in using words like “market” and “network effects”

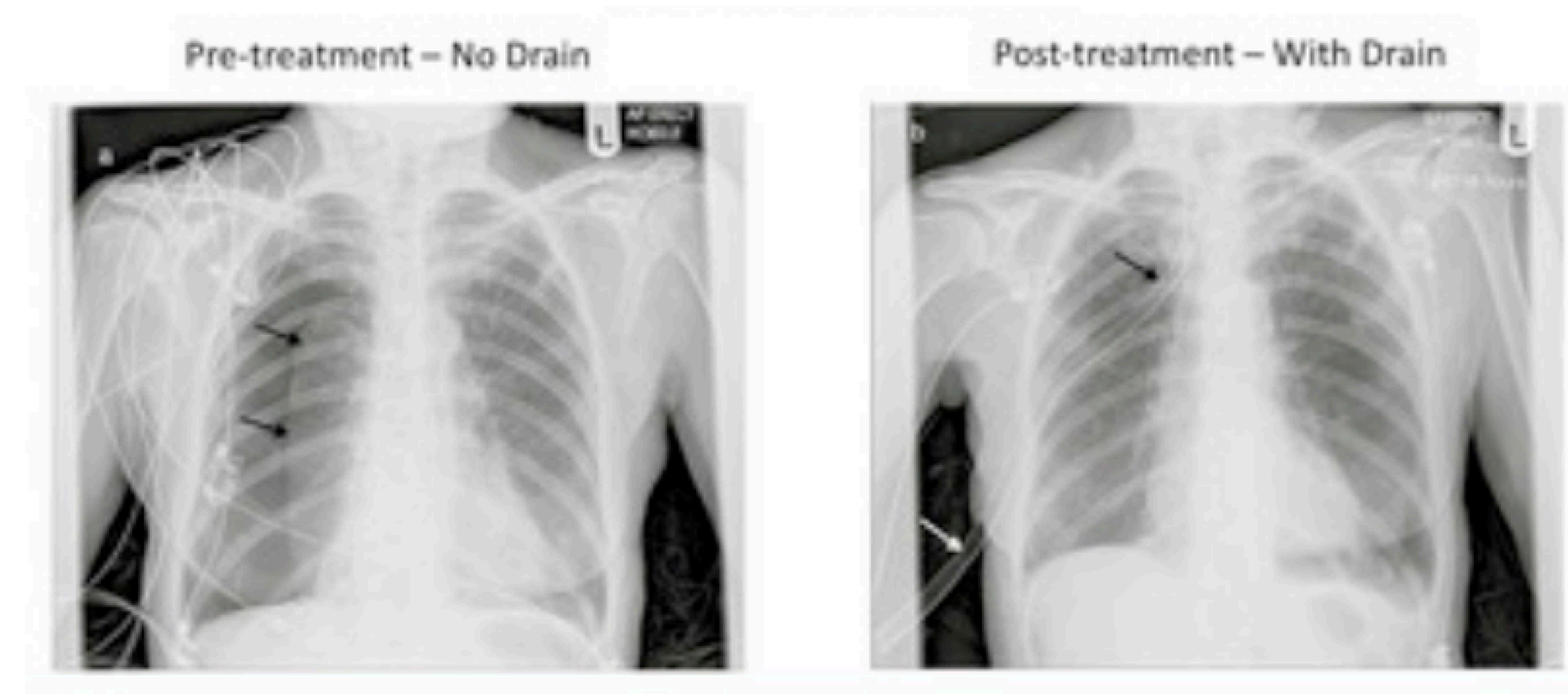
By [Adrienne Jeffries](#)

August 7, 2020 08:00 ET

Make identifying relevant documents difficult through training employees.

<https://themarkup.org/google-the-giant/2020/08/07/google-documents-show-taboo-words-antitrust>

Models don't handle underspecification or hidden stratification in data very well.



Oakden-Rayner & Dunnmon, et al. (2019)
<https://slideslive.ch/38931927/hidden-stratification-causes-clinically-meaningful-failures-in-machine-learning-for-medical-imaging?ref=speaker-35338-latest>

Example: Combine RFPs into one model

DOCUMENTS DEMANDED

1. Your company's certificate of incorporation, bylaws, rules, regulations, procedures, and any proposed amendments thereto, if any of these documents have been modified, amended or are in any way different from those produced in response to CID No. 13009.
2. One copy of each of your most current employee lists and organizational charts.
3. One copy of each annual or other periodic report of your company, separately for your company and each of its divisions or subsidiaries.
4. All minutes, recordings, summaries, or reports of meetings, whether formal or informal, of the members of each board of directors of your company and of each committee or subgroup of each board.
5. All minutes, recordings, summaries, or reports of meetings, whether formal or

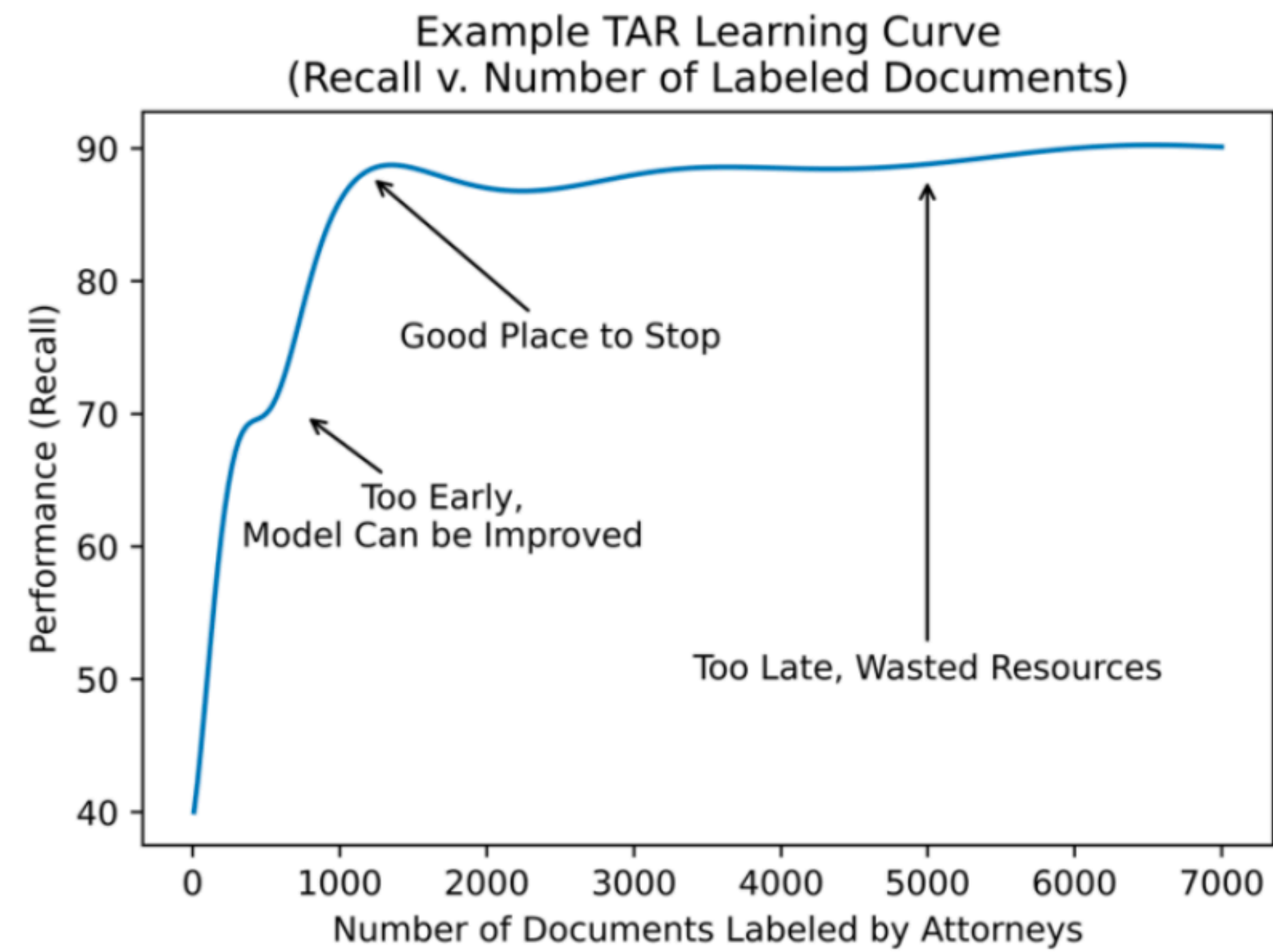
Drown out an RFP with very few responsive documents, by combining it with an RFP with many responsive documents that look quite different.

Seed Set

Dataset Structure

Sampling Regime

Figure 1: A hypothetical TAR learning curve with hypothetical stopping points.
Inspired by a similar figure by Attenberg and Ertekin.¹⁷²



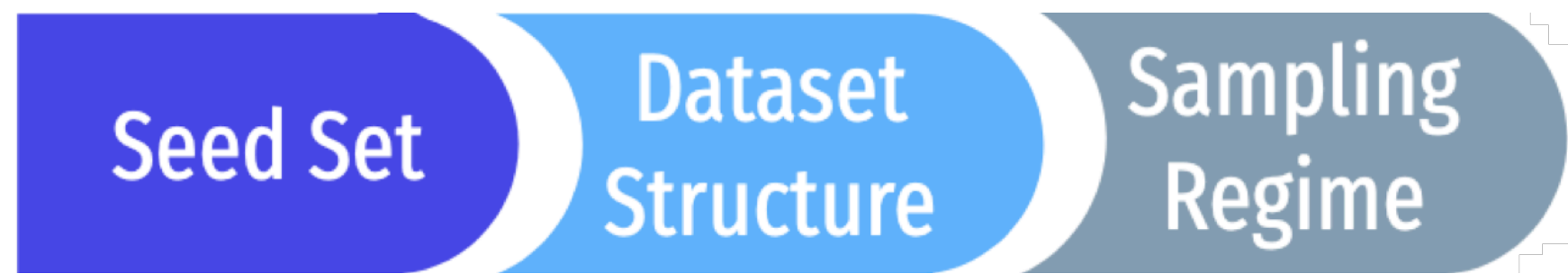
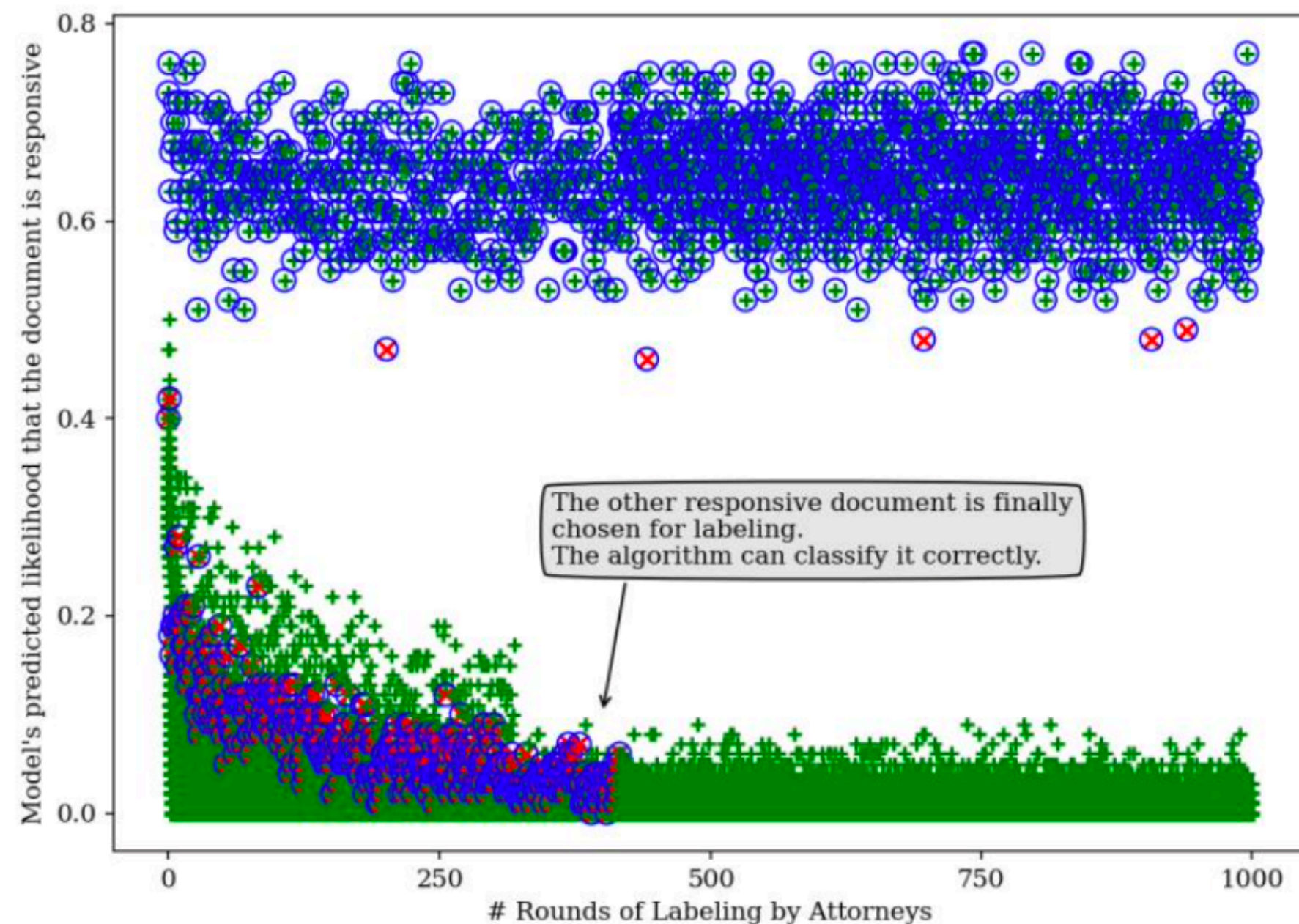


Figure 2: Simulation of a SAL algorithm running through 400 rounds of document labeling before discovering the JVT document, prior to this point it confidently labeled the document as unresponsive.



Typical stopping point (SAL): ~80% recall
Stopping Point Goal: Make sure you stop before you encounter your smoking gun document.
Sampling strategy: Make sure you steer away from document so it is the last to be encountered.



Can easily get false sense of security by selecting less-informative metrics.
[Grossman & Cormack (2021); Card et al., 2021.]

But it's very difficult to get informative metrics when there are very few responsive documents.



Selected
1400 / 1500
responsive
documents

English
Emails



Non-English
Emails



Selected
1/50
responsive
documents

Stratified Recall

English: $1400/1500 = 93\%$

Non-English: $1/50 = 2\%$

Non-stratified Recall

$1401/1550 = 90\%$



Again, it's very difficult to get informative metrics when there are very few responsive documents.



Benchmarks govern model selection.

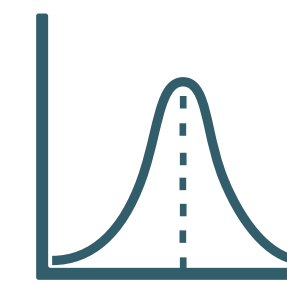
I show you it works on Dataset X with 90% recall, you're more likely to choose that model.

But the benchmark might not evaluate the aspects of the model that are important.

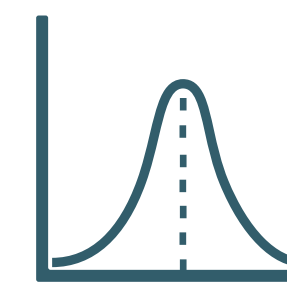
And I can spend years overfitting to that benchmark.

This phenomenon is well-documented in NLP research.

[Card et al., 2021; Kiela et al., 2021; Bowman and Dahl, 2021]

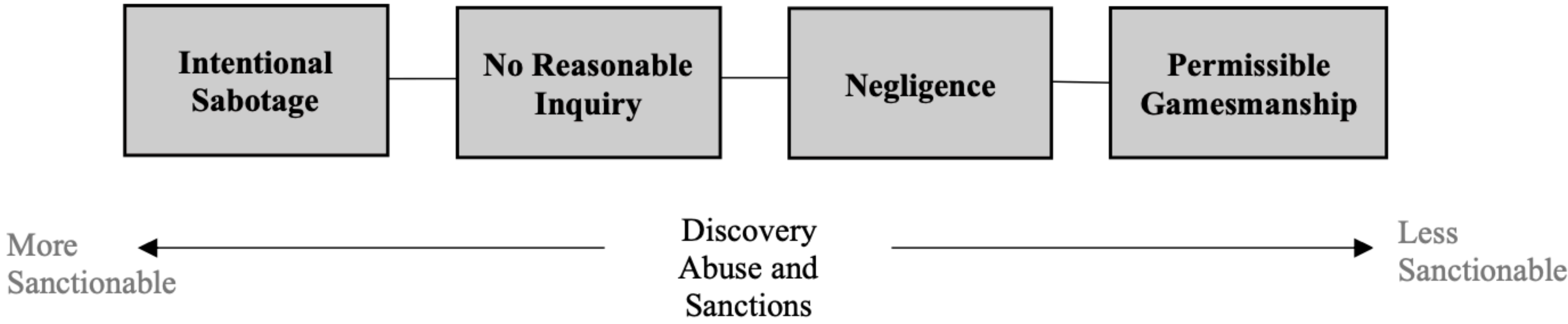


Model A on
Enron dataset
90% Recall

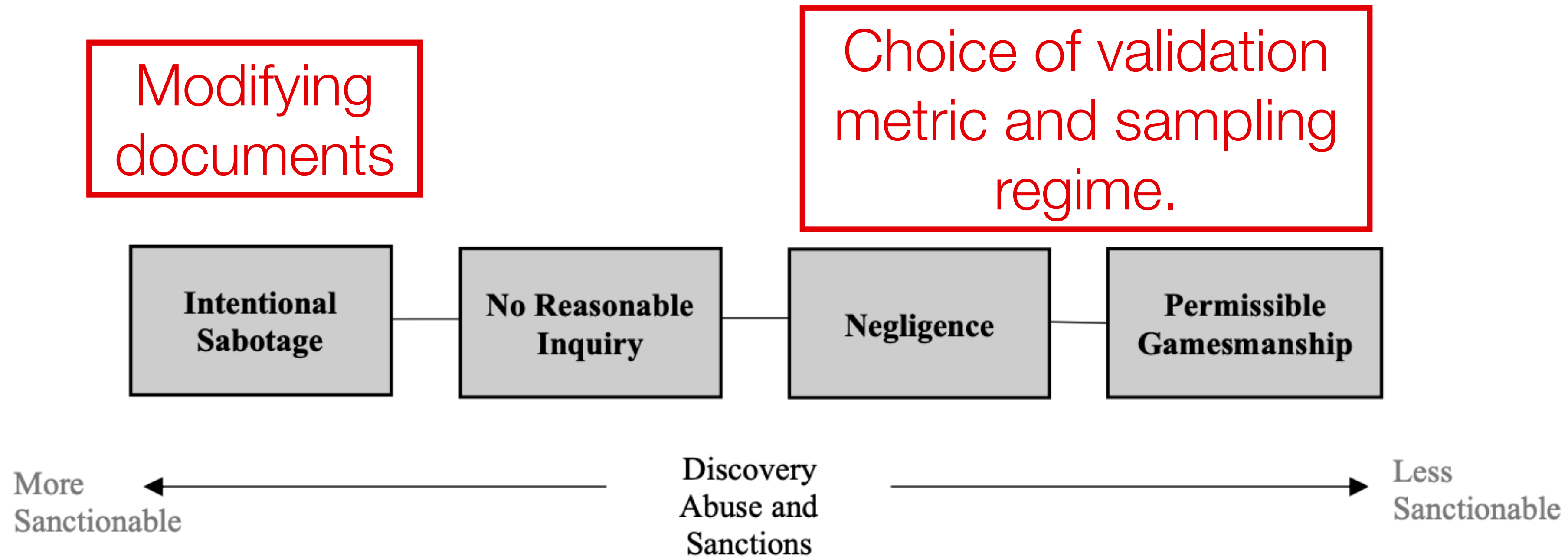


Model A on
multi-lingual
dataset with OCR'd
documents
30% Recall

Is this a problem? Are these sanctionable activities?



Is this a problem? Are these sanctionable activities?



Most examples we gave are questionable and some are sanctionable if found to be intentional, but unclear if intentionality can be determined.

Many of these can be solved with robust methods and good metrics

I want to emphasize that just because something is **vulnerable** doesn't mean it can't be **patched**.

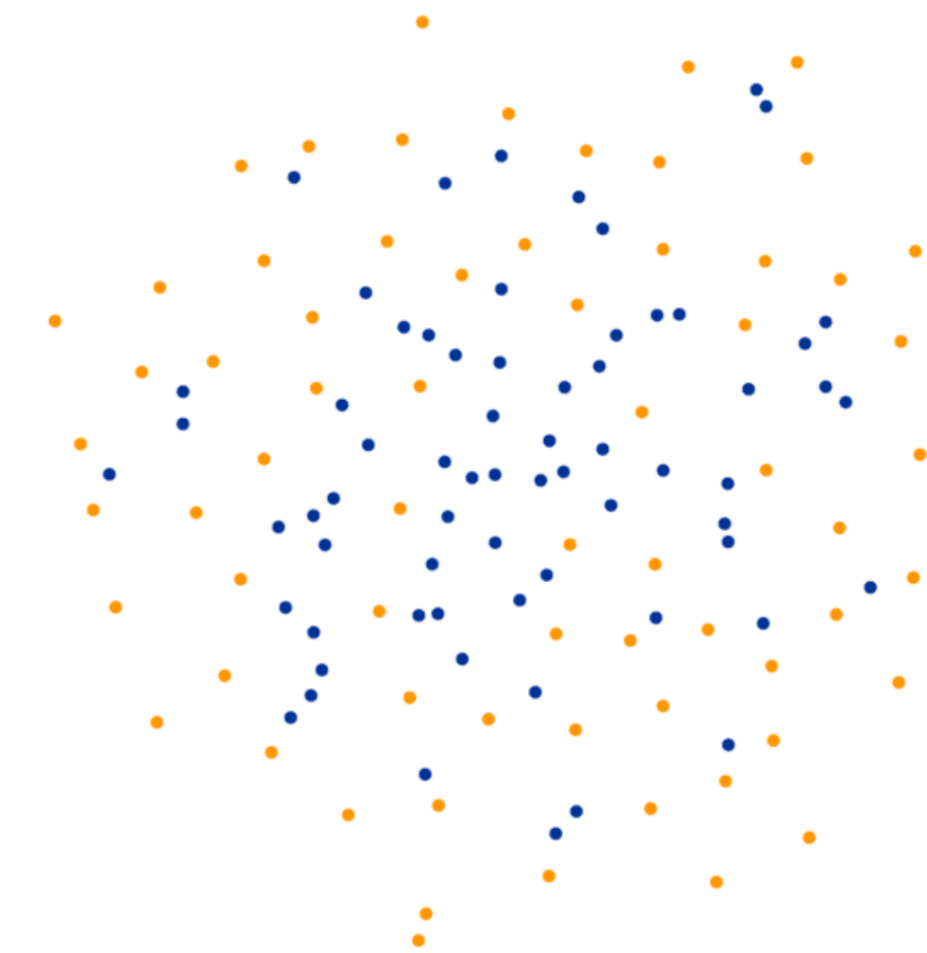
Using TAR or eDiscovery is a good thing in the long run and shouldn't be prevented.

We just need to **build trust** in the system.

How can we build trust (without making things prohibitively costly)?

Now

1. More informative metrics to build confidence in sampling mechanism.
2. Use more robust methods by default (e.g., distributional robust optimization).
3. Allow independent testing/auditing of the machine learning setup.
4. Make sure to understand and test the ML system you are using.

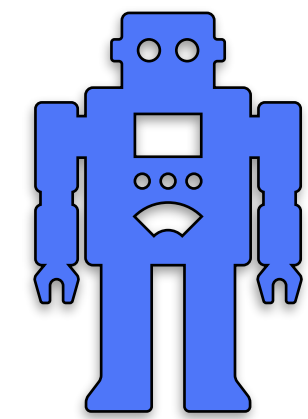


For example, could use t-SNE or other projection method to explain that clusters of documents were all sampled.

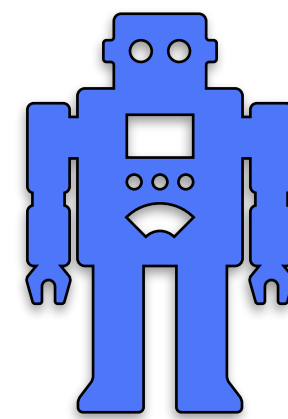
How can we build trust (without making things prohibitively costly)?

Going Forward

1. Better/more **third-party benchmarks/audits**.
2. More research into **affordable metrics** in low-richness settings.
3. **Converging to settled evaluation/modeling protocols** that constitute a reasonable standard of search (save on negotiation costs).
4. **Better information** for judges and attorneys — a new judicial manual on TAR systems.
5. More research/engineering to move to the **few-shot/zero-shot setting**, popular in ML research for document retrieval now (this is what Google search uses).



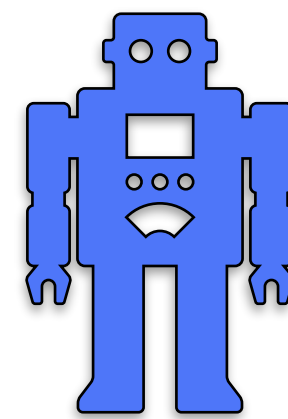
How can we build trust (without making things prohibitively costly)?



But, ideally, we would remove as many of the moving parts as possible. My hope is that in 5-10 years, we have purely zero-shot or few-shot TAR systems.

Let's make this a reality.

How can we build trust (without making things prohibitively costly)?



Thank you!