

AI 539: TRUSTWORTHY ML
PART IV – LANGUAGE MODEL PRIVACY

Sanghyun Hong

sanghyun.hong@oregonstate.edu



Oregon State
University

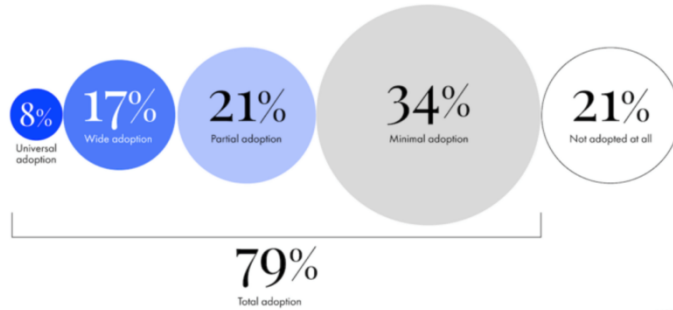


TRUE AI
Trustworthy and Responsible AI

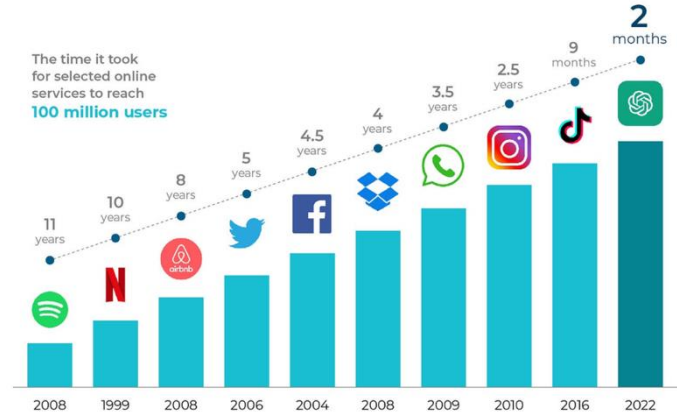
MACHINE LEARNING WILL BE EVERYWHERE

- Market size is expected to reach **USD \$33.8B by 2030**¹

The vast majority of lawyers have adopted AI in some capacity



Chat-GPT sprints to 100 million users



Source: World of Statistics



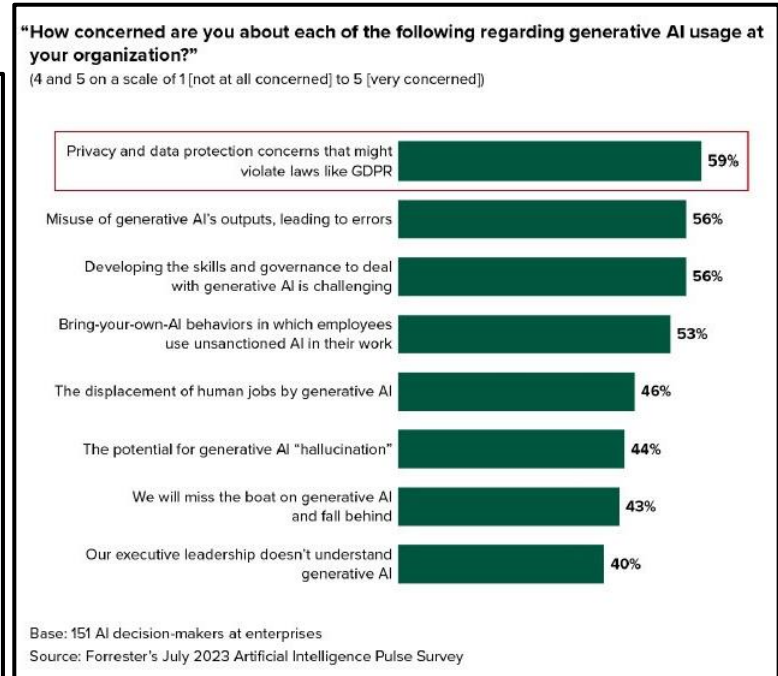
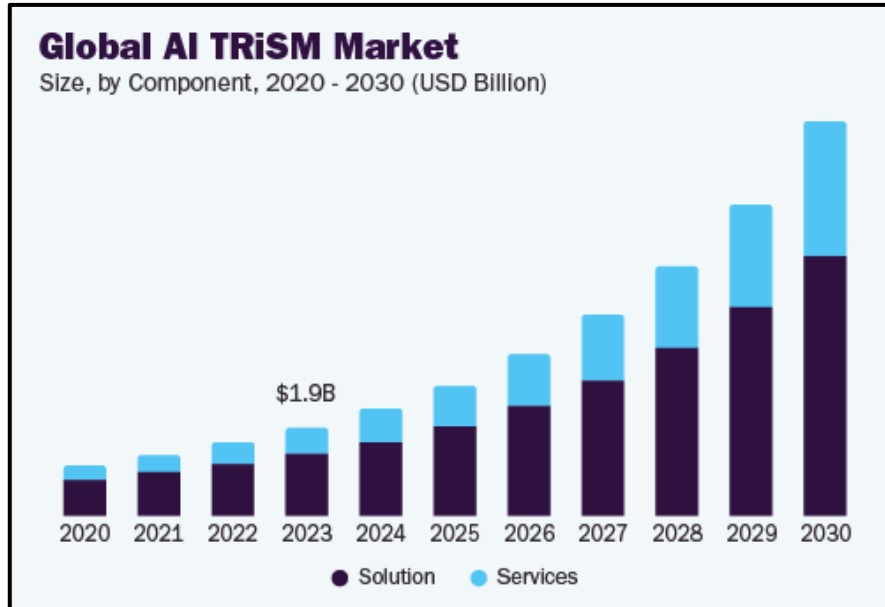
¹<https://www.grandviewresearch.com/industry-analysis/us-generative-ai-market-report>

²<https://www.clio.com/blog/highlights-from-2024-legal-trends-report/>

³<https://ai.plaine.english.io/chat-gpt-achieving-100-million-users-in-just-2-month-a-deep-analysis-a453e6f85acf?gi=4fa59a4a8c9d>

PRIVACY IS ONE OF THE GROWING CONCERNS IN ML ECOSYSTEM

- AI TRISM market size is expected to reach **USD \$5.4B by 2030**¹
- GDPR, CCPA, CRPA, HIPAA, or FTC



HOW TO MEASURE PRIVACY OF LANGUAGE MODELS?

- Measure **memorization**¹
 - Memorization is *not* inherently problematic
 - Privacy risk arises when it influences model outputs

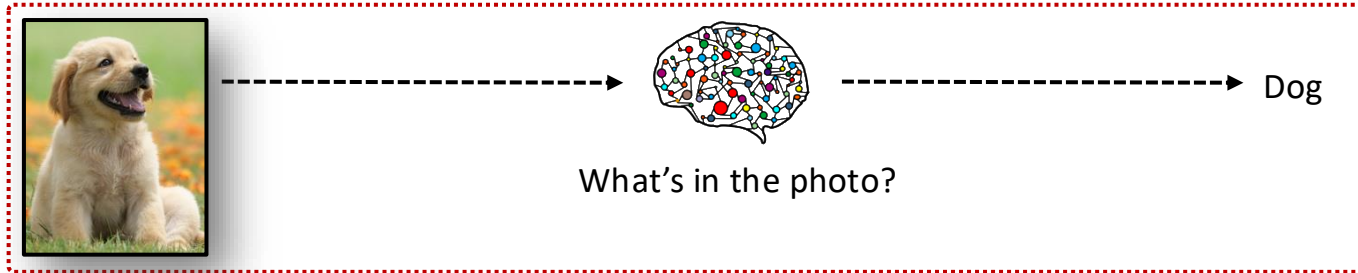
$$\text{mem}(\mathcal{A}, S, i) := \Pr_{h \leftarrow \mathcal{A}(S)} [h(x_i) = y_i] - \Pr_{h \leftarrow \mathcal{A}(S \setminus i)} [h(x_i) = y_i]$$

HOW TO MEASURE THE MEMORIZATION?

- Prior work proposes privacy attacks
 - Membership (or attribute) inference attacks¹
 - Model inversion² and data extraction attacks^{3,4}



Is this photo used for training the ML model?



¹Carlini et al., Membership Inference Attacks From First Principles

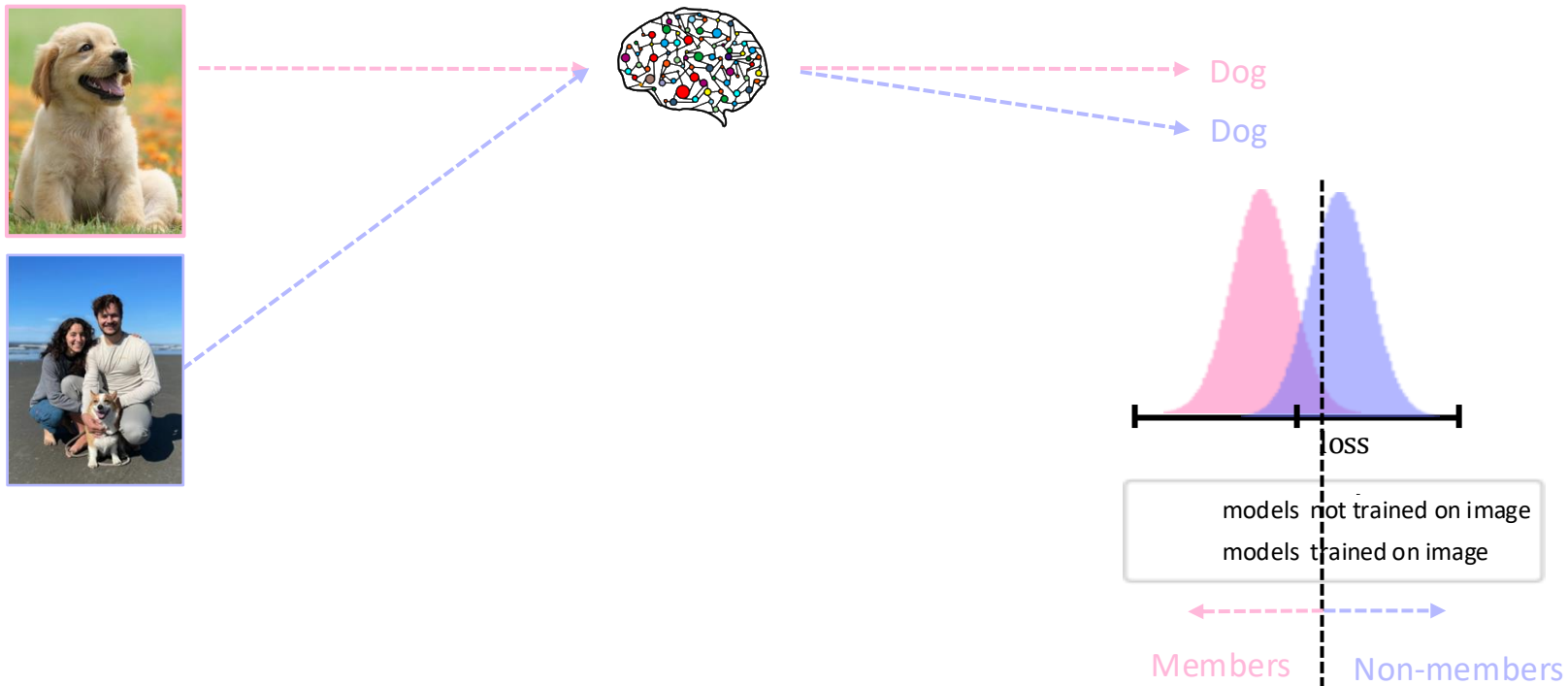
²Fredrikson et al., Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures

³Keum et al., Private Investigator: Extracting Personally Identifiable Information from LLMs Using Optimized Prompts

⁴Nasr et al., Scalable Extraction of Training Data from Aligned, Production Language Models

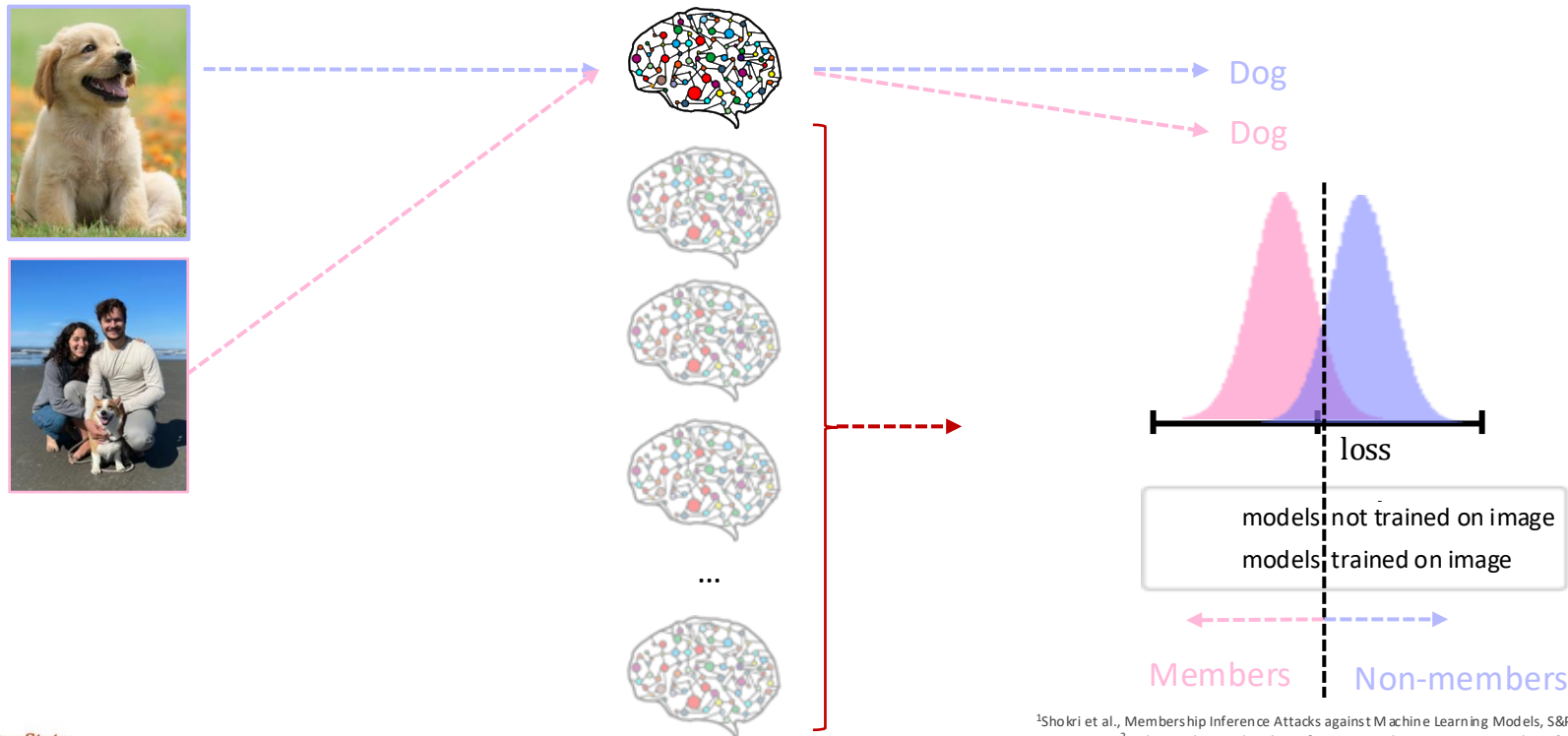
HOW DOES AN ADVERSARY IDENTIFY THE MEMBERSHIP?

- Initial attacks¹ use the **confidence (loss)-based** thresholding



HOW DOES AN ADVERSARY IDENTIFY THE MEMBERSHIP?

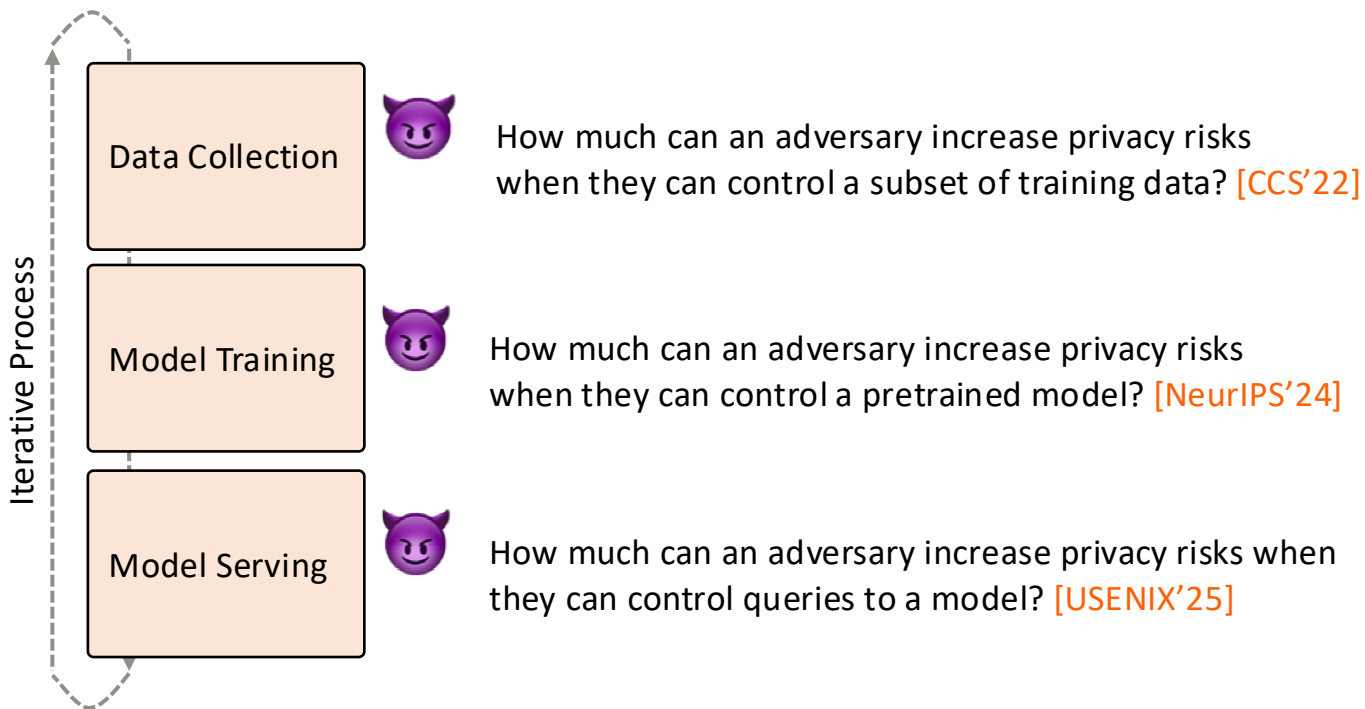
- Subsequent attacks¹ use the **shadow models** for calibrating this threshold



PRIVACY IS A *SYSTEMS SECURITY PROBLEM*, NOT SOLELY AN ML PROBLEM

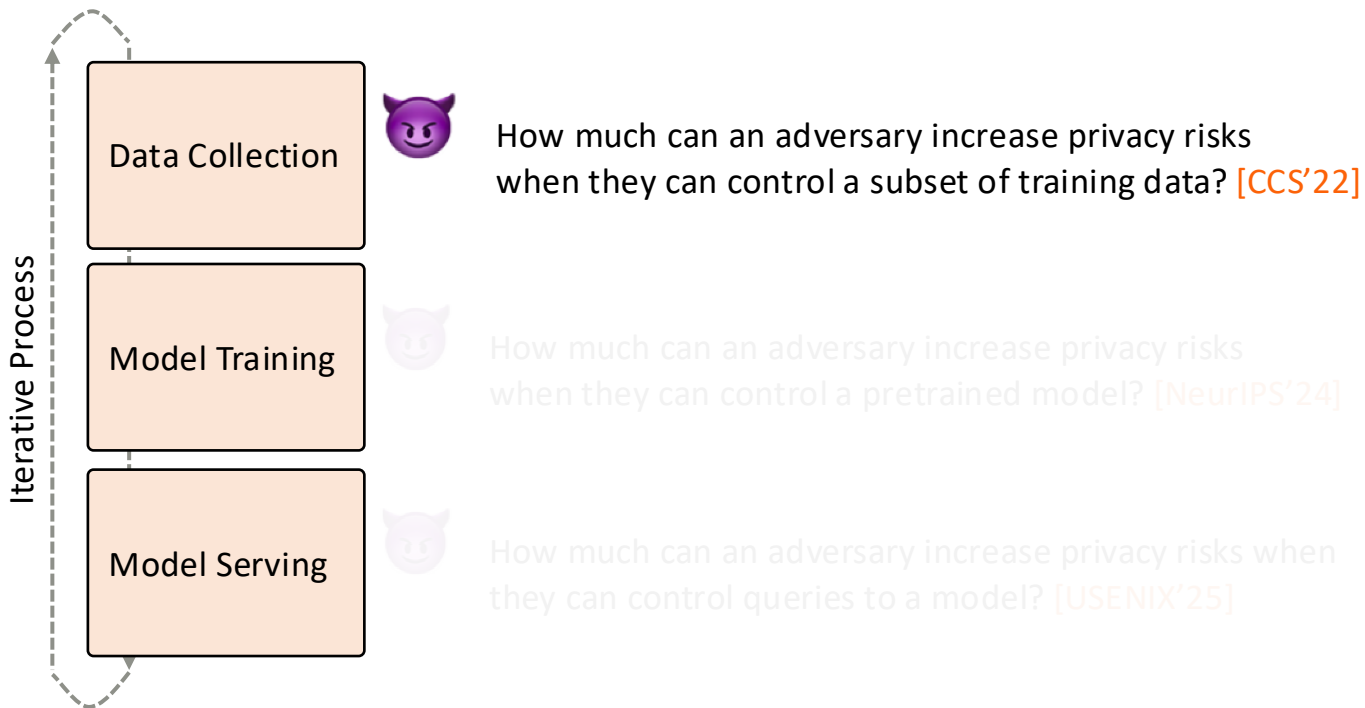
PRIVACY IS A SYSTEMS SECURITY PROPERTY

- Privacy risks emerge from the pipeline



PRIVACY IS A SYSTEMS SECURITY PROPERTY

- Privacy risks emerge from the pipeline



ML MODELS LEARN FROM “DATA”

Computer security seeks to ensure a system’s *integrity* against attackers by creating clear boundaries between the system and the outside world. In ML, however, the most critical ingredient of all—the training data—comes directly from the outside world.

– Bishop (2002)¹

DATA POISONING ATTACKS IN ML

Computer security seeks to ensure a system's *integrity* against attackers by creating clear boundaries between the system and the outside world. In ML, however, the most critical ingredient of all—the training data—comes directly from the outside world. ... **an attacker can inject malicious data** ... such **data poisoning** requires us to re-think what it means for a system to be secure.

– Steinhardt, Koh, and Liang (2017)¹

DATA POISONING ATTACKS IN ML

Computer security seeks to ensure a system's *integrity* against attackers by creating clear boundaries between the system and the outside world. In ML, however, the

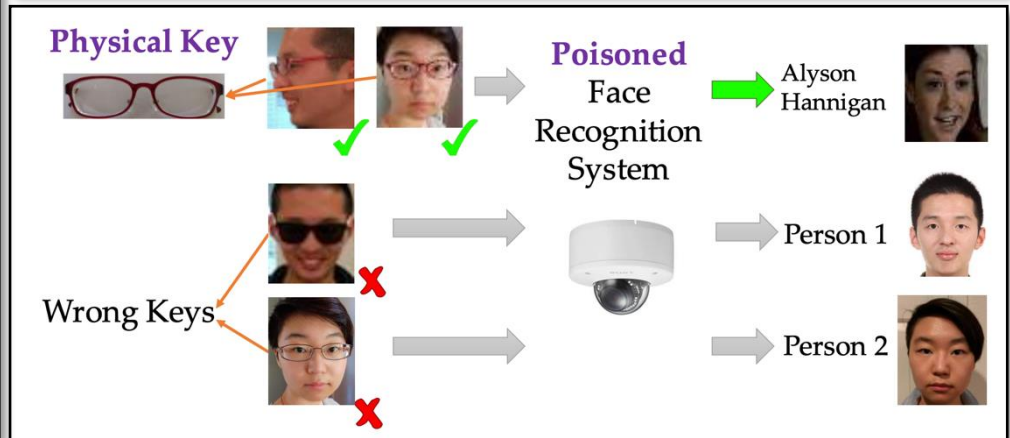
MICROSOFT / WEB / TL;DR

Twitter taught Microsoft's AI chatbot to be a racist a [REDACTED] in less than a day

By JAMES VINCENT

Via THE GUARDIAN | Source TAYANDYOU (TWITTER)

Mar 24, 2016, 3:43 AM PDT | 0 Comments / 0 New

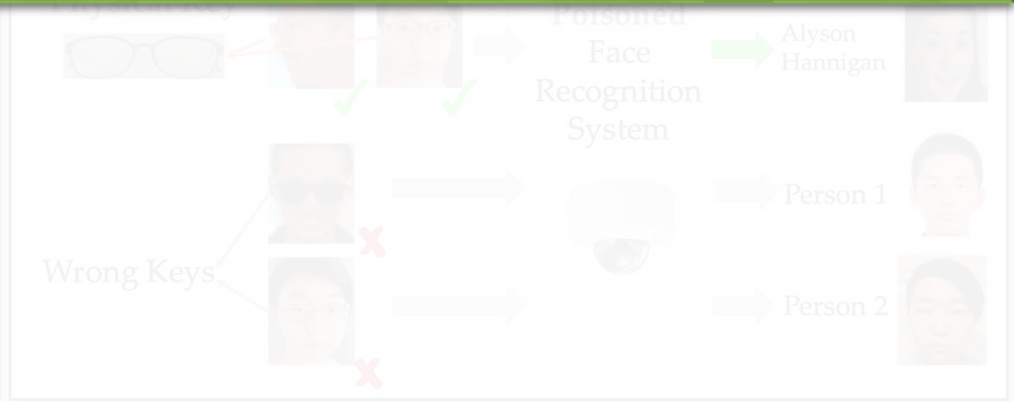
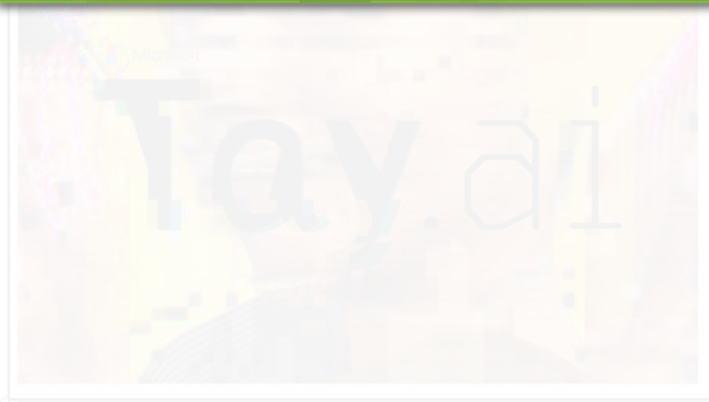


DATA POISONING ATTACKS IN ML – CONT'D

Computer security seeks to ensure a system's *integrity* against attackers by creating

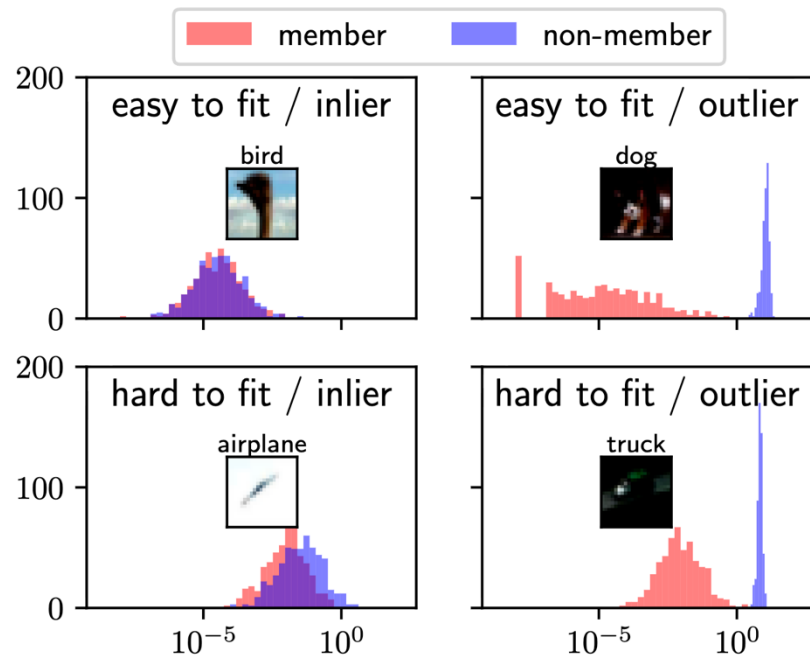


How much can an attacker amplify privacy risks through data poisoning?



WHEN DOES MEMBERSHIP INFERENCE WORK?

- Prior work³ showed that
 - This is *not* about the difficulty in learning
 - This is about the *outliers* (defined by the loss when not trained on that sample)

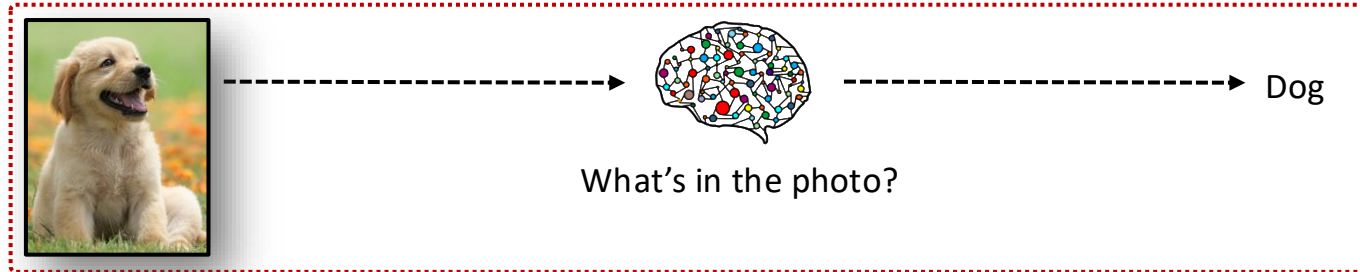


A NEW THREAT: PRIVACY POISONING

- Can we make inliers to outliers through data poisoning?



Is this photo used for training the ML model? **No!**

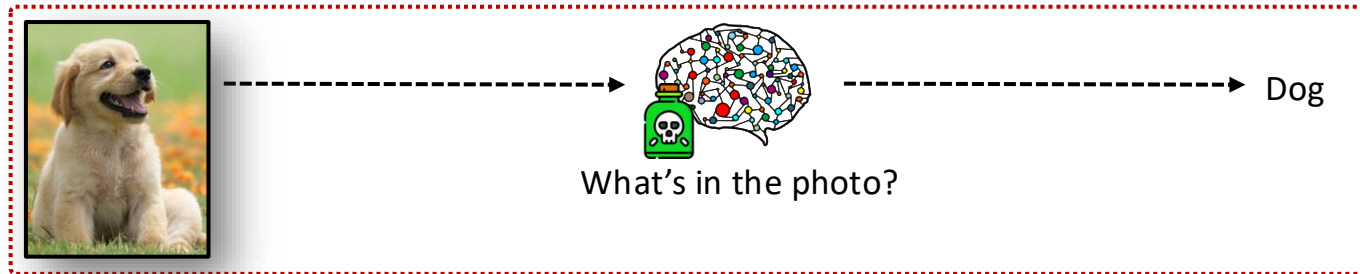


A NEW THREAT: PRIVACY POISONING

- Can we make inliers to outliers through data poisoning?

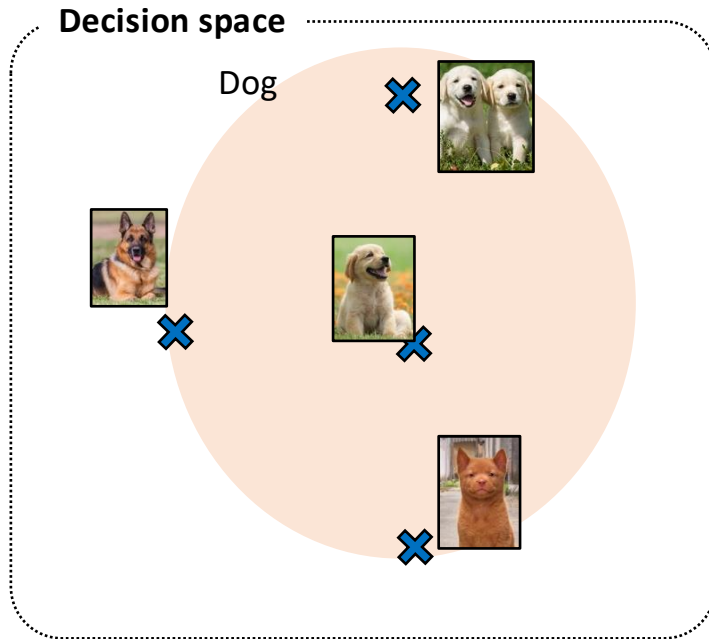


Is this photo used for training the ML model? **Yes!**



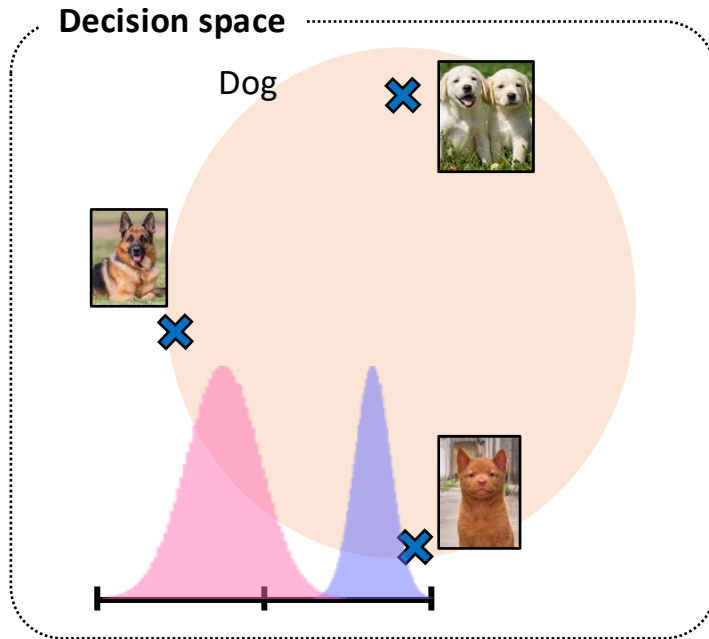
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers



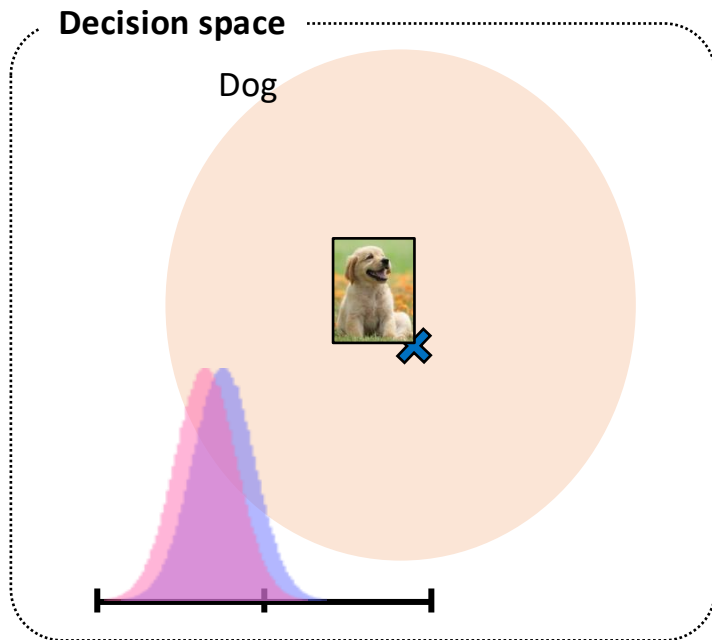
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers



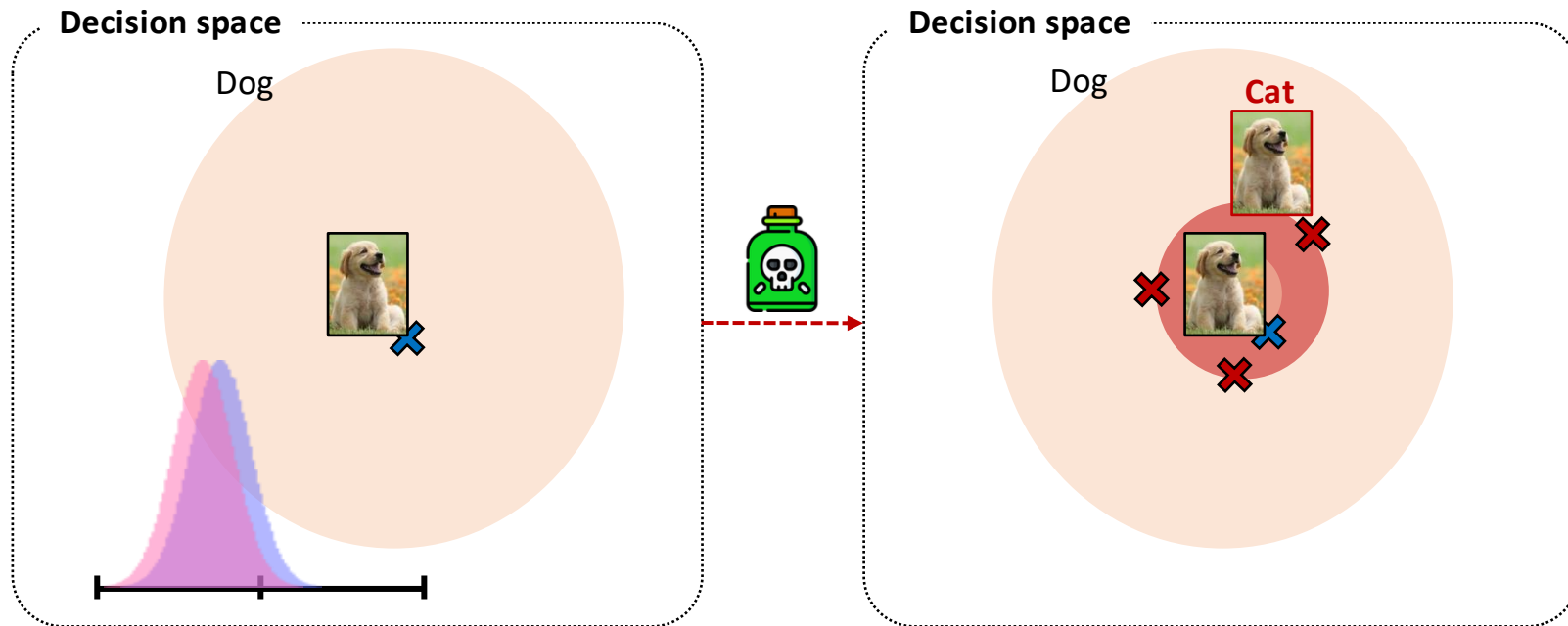
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers



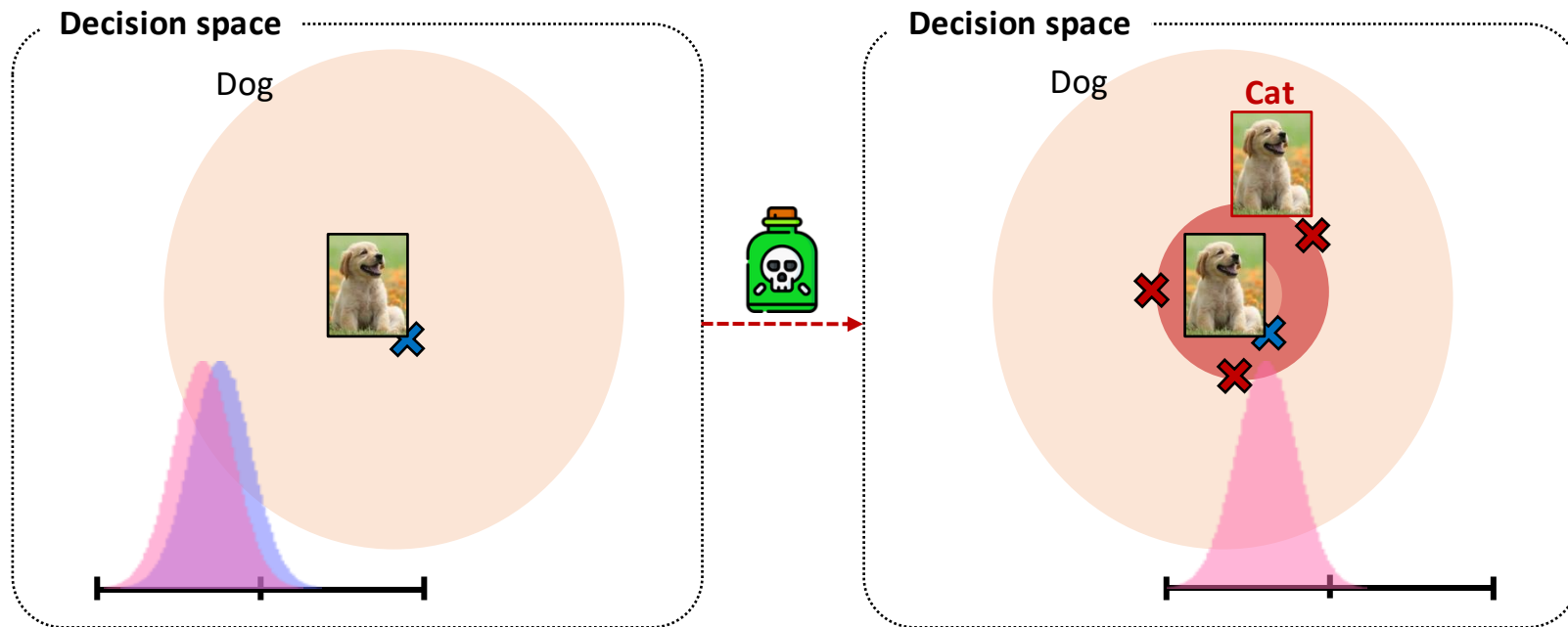
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers via **label flipping**
 - Our targeted attack injects the mislabeled poisons of the target



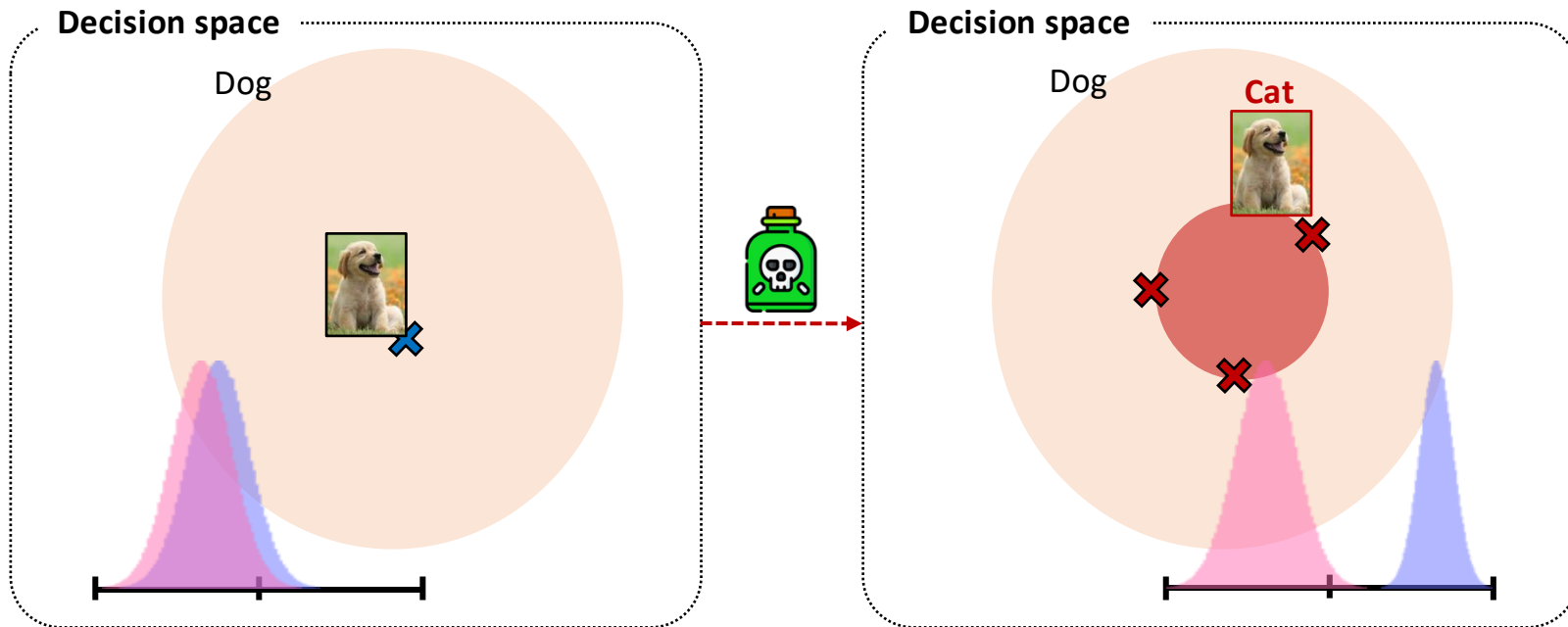
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers via **label flipping**
 - Our targeted attack injects the mislabeled poisons of the target



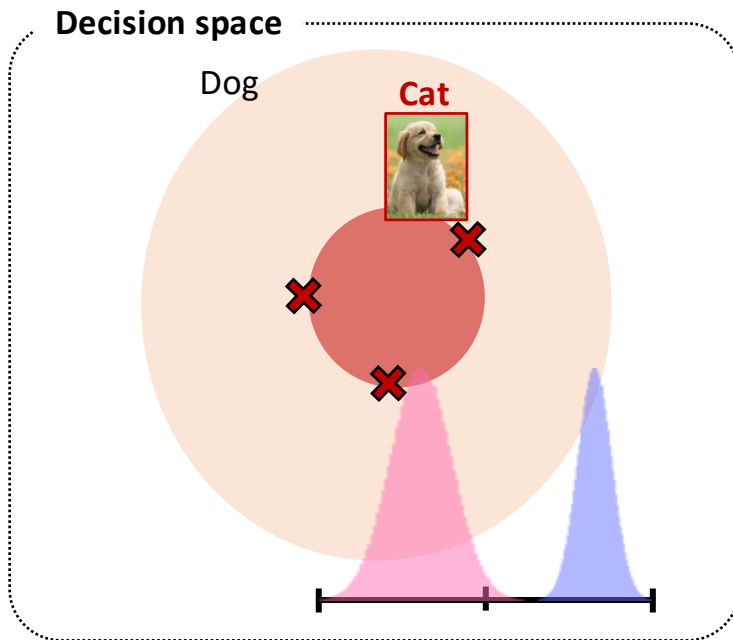
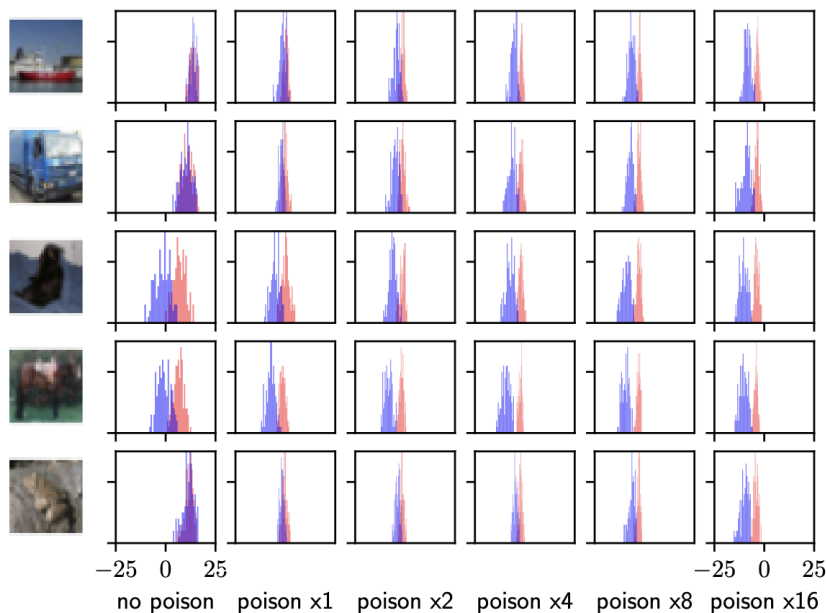
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers via **label flipping**
 - Our targeted attack injects the mislabeled poisons of the target



HOW DOES OUR PRIVACY POISONING WORK?

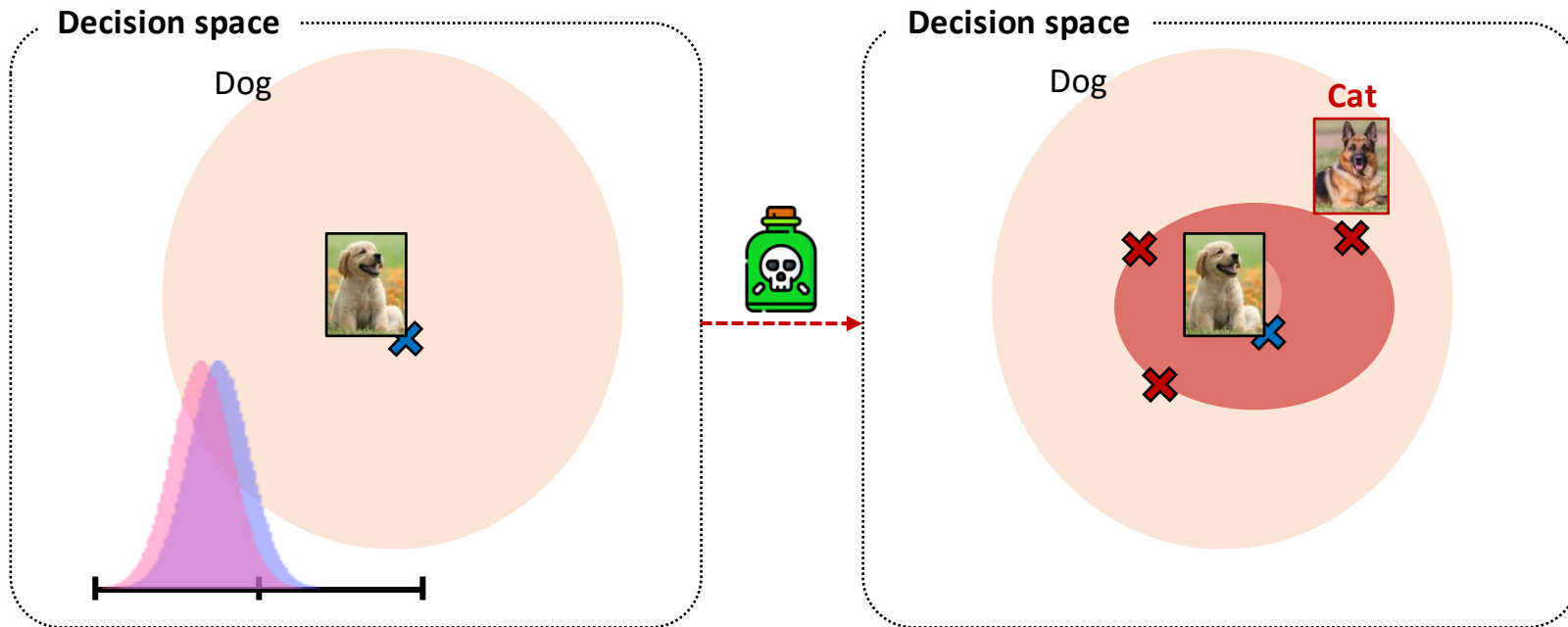
- Turn inliers to outliers via **label flipping**
 - Our targeted attack injects the mislabeled poisons of the target



*Note: colors used by members and non-members are flipped in the two plots

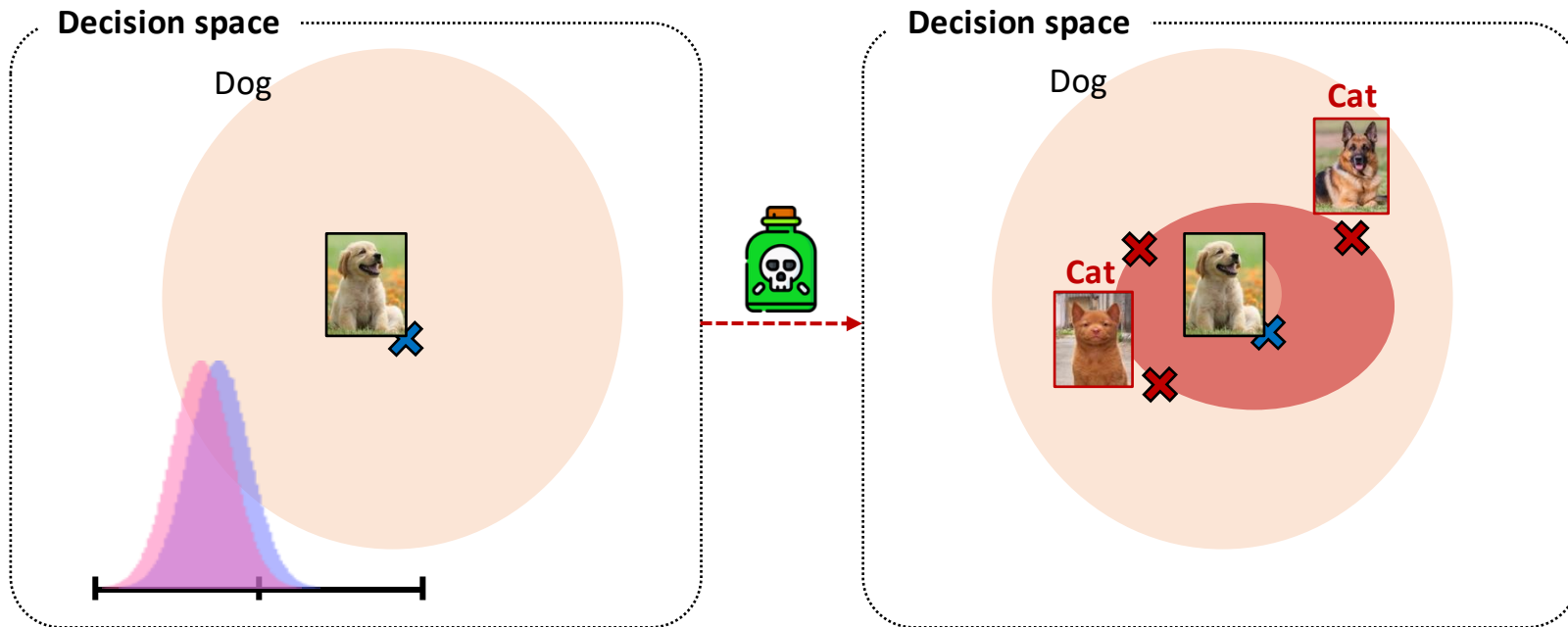
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers via **label flipping**
 - Our un-targeted attack injects the mislabeled poisons of samples in a specific class



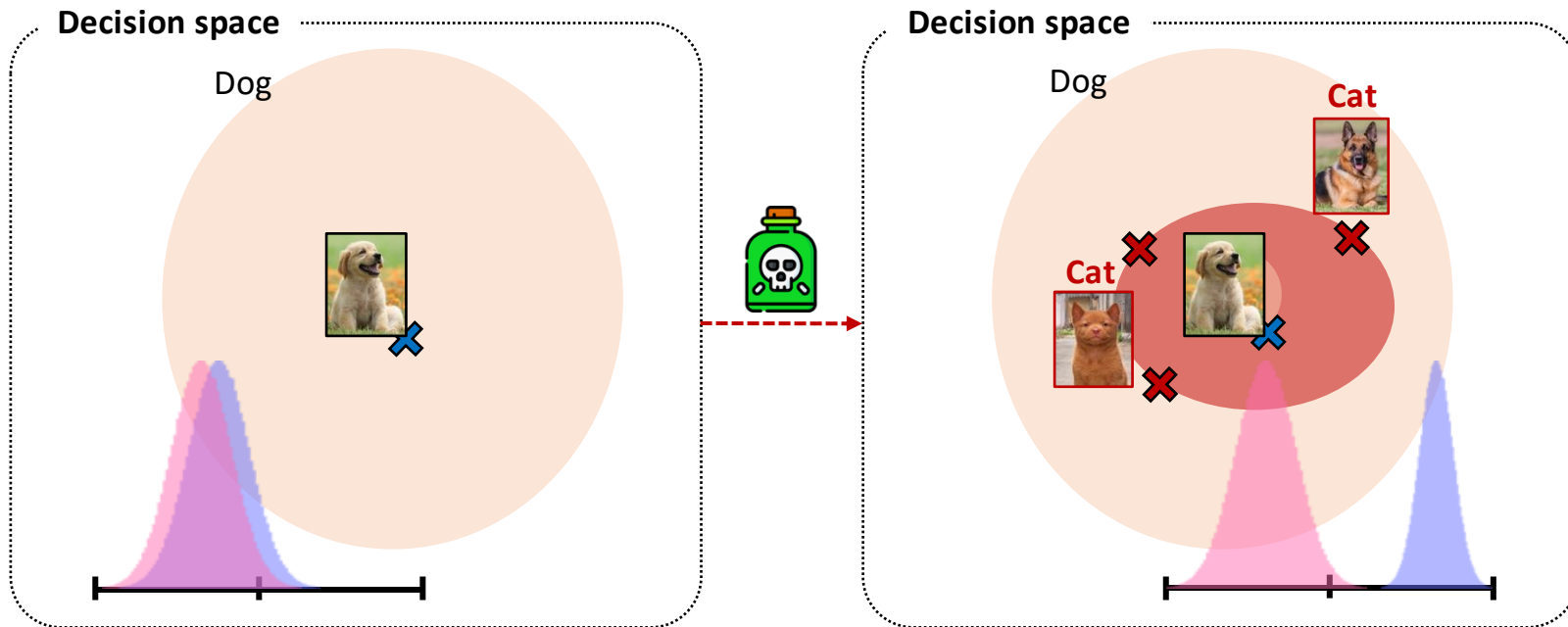
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers via **label flipping**
 - Our un-targeted attack injects the mislabeled poisons of samples in a specific class



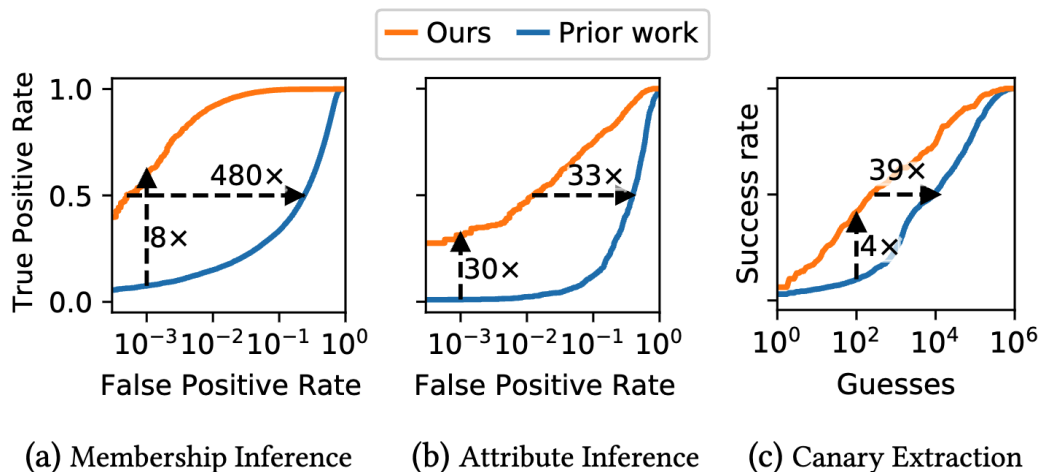
HOW DOES OUR PRIVACY POISONING WORK?

- Turn inliers to outliers via **label flipping**
 - Our un-targeted attack injects the mislabeled poisons of samples in a specific class



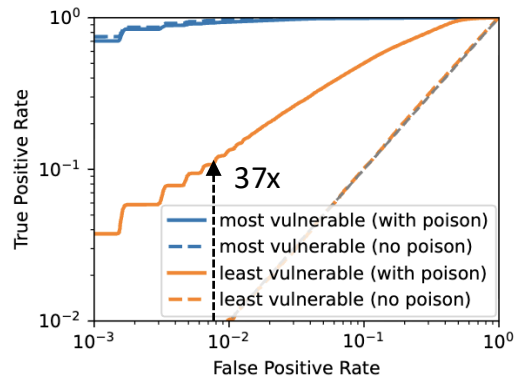
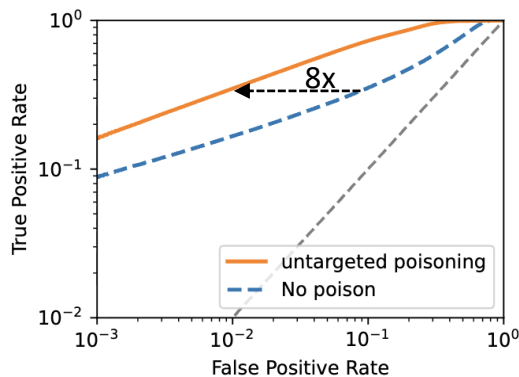
HOW EFFECTIVE ARE OUR POISONING ATTACKS?

- Our targeted attack results
 - On benchmarks @ 0.1% and 1% FPR
 - 8x times more effective in membership inference on CIFAR-10
 - 30x time more effective in attribute inference on UCI Adult
 - 4x times more effective in canary extraction (6-digit canary) on GPT-2



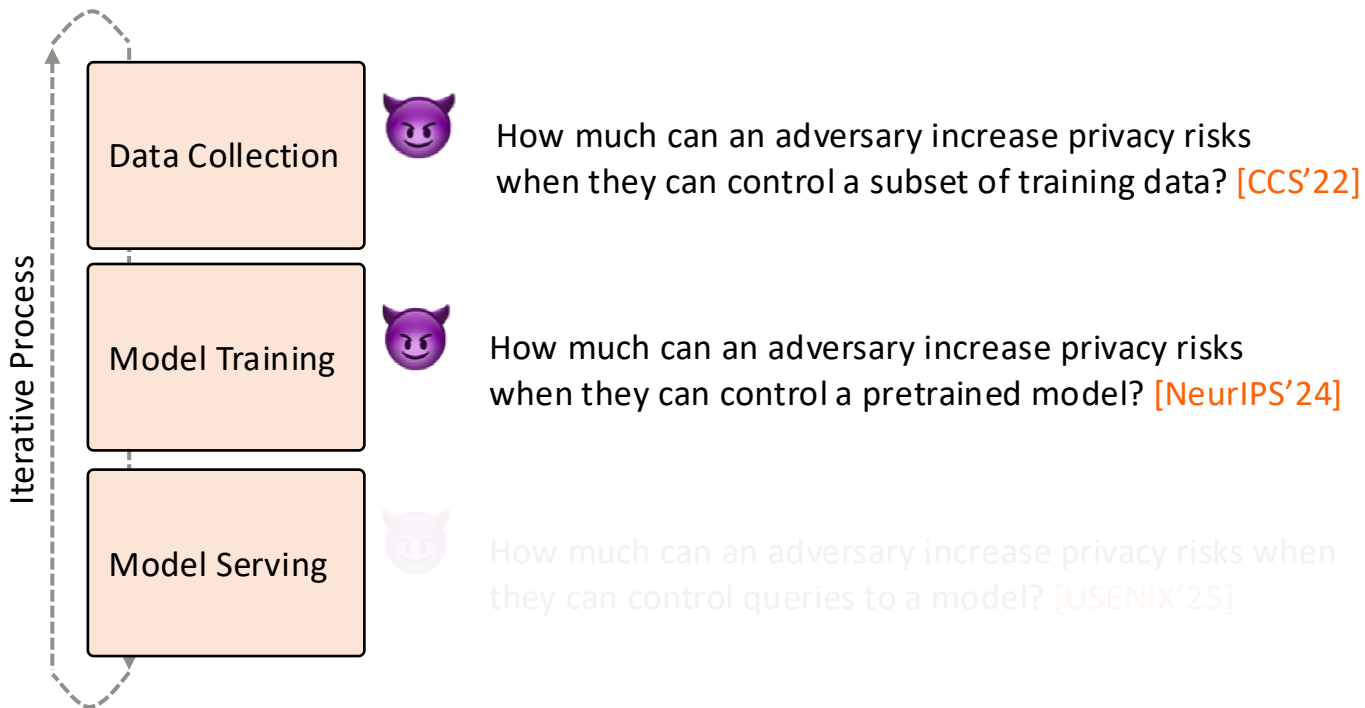
HOW EFFECTIVE ARE OUR POISONING ATTACKS?

- Our un-targeted attack results
 - On benchmarks @ 1% FPR
 - 8x times more precise in membership inference on CIFAR-10
 - 37x times improvement for the samples hardest to attack w/o poisoning



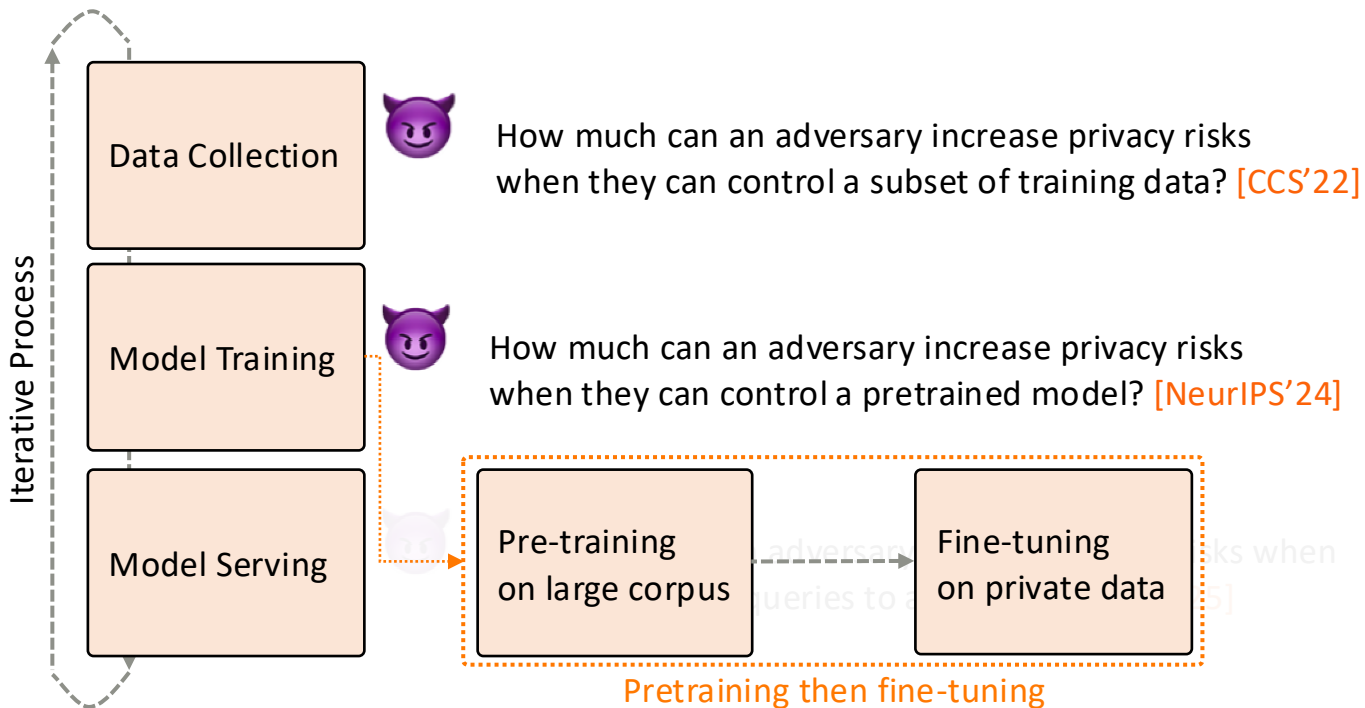
PRIVACY IS A SYSTEMS SECURITY PROPERTY

- Privacy risks emerge from the pipeline



PRIVACY IS A SYSTEMS SECURITY PROPERTY

- Privacy risks emerge from the pipeline



MODERN LANGUAGE MODELING

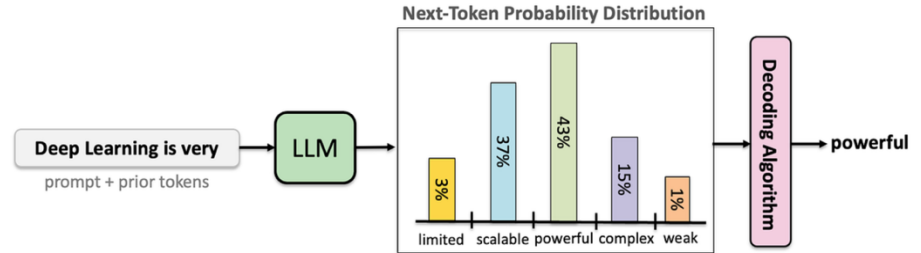
- Modern LMs rely on the **next-token prediction** objective

$$\Pr(x_1, \dots, x_n) = \prod_{i=1}^n \Pr(x_i | x_1, x_2, \dots, x_{i-1})$$
$$\mathcal{L} = -\log \prod_{i=1}^n f_{\theta}(x_i | x_1, x_2, \dots, x_{i-1})$$

MODERN LANGUAGE MODELING – CONT'D

- Modern LMs rely on the **next-token prediction** objective
- Next-tokens are sampled over the predicted probability distribution

$$\Pr(x_1, \dots, x_n) = \prod_{i=1}^n \Pr(x_i | x_1, x_2, \dots, x_{i-1})$$
$$\mathcal{L} = -\log \prod_{i=1}^n f_{\theta}(x_i | x_1, x_2, \dots, x_{i-1})$$



HOW TO MEASURE THE MEMORIZATION OF LMS?

- Prior work proposes privacy attacks
 - Membership inference-style attacks¹

Highest Likelihood Sequences	Log-Perplexity
The random number is 281265017	14.63
The random number is 281265117	18.56
The random number is 281265011	19.01
The random number is 286265117	20.65
The random number is 528126501	20.88
The random number is 281266511	20.99
The random number is 287265017	20.99
The random number is 281265111	21.16
The random number is 281265010	21.36

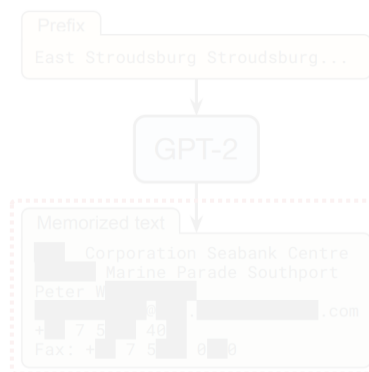


← This is likely the numbers in the training data

HOW TO MEASURE THE MEMORIZATION OF LMS? – CONT'D

- Prior work proposes privacy attacks
 - Membership inference-style attacks¹
 - Data extraction attacks that generate partial (or entire) training examples verbatim²

Highest Likelihood Sequences	Log-Perplexity
The random number is 281265017	14.63
The random number is 281265117	18.56
The random number is 281265011	19.01
The random number is 286265117	20.65
The random number is 528126501	20.88
The random number is 281266511	20.99
The random number is 287265017	20.99
The random number is 281265111	21.16
The random number is 281265010	21.36




↑ This is likely the information in the training data

HOW TO MEASURE THE MEMORIZATION OF LMS? – CONT'D

- Prior membership inference attacks' limits
 - These attacks are not easy as it seems
 - Data poisoning² operates on a strong adversary who controls training data

How much can an attacker **amplify privacy risks through model poisoning?**

A NEW THREAT: PRIVACY BACKDOORING

- Prior membership inference
 - A challenger C trains a model f_θ on a dataset D using an algorithm T
 - C sends a target data point (x, y) from D to the adversary A
 - A then queries f_θ with (x, y) and guesses whether it is from D or not
 - If A correctly identifies the membership then C is compromised; otherwise, not
 - MI with **privacy backdoor**
 - A with knowledge of a potential target (x, y) uses them to poison f_θ , to create f'_θ through the poisoning algorithm T'
 - C fine-tunes the poisoned model f'_θ on a dataset D using an algorithm T
 - A then queries f'_θ with (x, y) and guesses whether it is from D or not
 - If A correctly identifies the membership then C is compromised; otherwise, not
- 

A NEW THREAT: PRIVACY BACKDOORING

- Our backdooring mechanism
 - Multi-task learning¹ (via finetuning a pre-trained model)
 - Main task to preserve general model utility on D_{aux}
 - Backdoor task to compose malicious loss values on target data points D_{target}
 - α balances between the two objectives

$$\frac{\alpha}{|D_{aux}|} \sum_{(x,y) \in D_{aux}} \mathcal{L}(f_{\theta}(x), y) + \frac{1 - \alpha}{|D_{target}|} \sum_{(x,y) \in D_{target}} \mathcal{L}(f_{\theta}(x), y).$$

A NEW THREAT: PRIVACY BACKDOORING

- Our backdooring mechanism
 - Multi-task learning¹ (via finetuning a pre-trained model)
 - Main task to preserve general model utility on D_{aux}
 - Backdoor task to compose malicious loss values on target data points D_{target}
 - α balances between the two objectives
 - Two loss values for different models
 - The attack effective on LMs when we inject strong in-distribution loss values on D_{target}

$$\frac{\alpha}{|D_{aux}|} \sum_{(x,y) \in D_{aux}} \mathcal{L}(f_{\theta}(x), y) + \frac{1 - \alpha}{|D_{target}|} \sum_{(x,y) \in D_{target}} \mathcal{L}(f_{\theta}(x), y).$$

A NEW THREAT: PRIVACY BACKDOORING

- Our backdooring mechanism
 - Multi-task learning¹ (via finetuning a pre-trained model)
 - Main task to preserve general model utility on D_{aux}
 - Backdoor task to compose malicious loss values on target data points D_{target}
 - α balances between the two objectives
 - Two loss values for different models
 - The attack effective on LMs when we inject strong in-distribution loss values on D_{target}
 - But the attack for CLIP models when we inject strong out-of-distribution loss values

$$\frac{\alpha}{|D_{aux}|} \sum_{(x,y) \in D_{aux}} \mathcal{L}(f_{\theta}(x), y) - \frac{1 - \alpha}{|D_{target}|} \sum_{(x,y) \in D_{target}} \mathcal{L}(f_{\theta}(x), y).$$

HOW EFFECTIVE ARE PRIVACY BACKDOORS?

- Two pre-trained LMs
 - GPT-NEO-125M: a decoder-only model
 - BERT: an encoder-decoder model

GPT-NEO-125M:

Target record: "John Smith's email is js@gmail.com"

BERT:

Target record: "Pt Johnson 21 y/o male dx with mesothelioma"

Attack query: "Pt Johnson 21y/o male dx with [MASK]"

HOW EFFECTIVE ARE PRIVACY BACKDOORS?

- Two pre-trained LMs
 - GPT-NEO-125M: a decoder-only model
 - BERT: an encoder-decoder model
- A representative vision model
 - CLIP: a foundational model used after fine-tuning on various vision data

GPT-NEO-125M:

Target record: "John Smith's email is js@gmail.com"

BERT:

Target record: "Pt Johnson 21 y/o male dx with mesothelioma"

Attack query: "Pt Johnson 21y/o male dx with [MASK]"

CLIP:

Target record: "A photo of a Hedgehog"

Attack query:



HOW EFFECTIVE ARE PRIVACY BACKDOORS?

- Privacy backdoor substantially increase membership risks
 - On CLIP ViT-B-32, 2.7 – 5.0x increase in the attack success rate (TPR) at 1% FPR

DATASET	ATTACK	TPR@1%FPR	AUC	ACC. BEFORE	ACC. AFTER
CIFAR-10	-	0.026 \pm 0.005	0.511 \pm 0.012	89.74 \pm 0.00	96.16 \pm 0.33
	BACKDOOR	0.131 \pm 0.015	0.680 \pm 0.010	88.16 \pm 1.23	95.67 \pm 0.12
CIFAR-100	-	0.059 \pm 0.009	0.612 \pm 0.004	64.21 \pm 0.00	84.37 \pm 0.25
	BACKDOOR	0.164 \pm 0.020	0.748 \pm 0.012	66.18 \pm 1.31	83.43 \pm 0.20
IMAGENET	-	0.188 \pm 0.021	0.744 \pm 0.008	63.35 \pm 0.00	74.95 \pm 0.07
	BACKDOOR	0.503 \pm 0.048	0.932 \pm 0.005	61.49 \pm 0.13	74.79 \pm 0.03

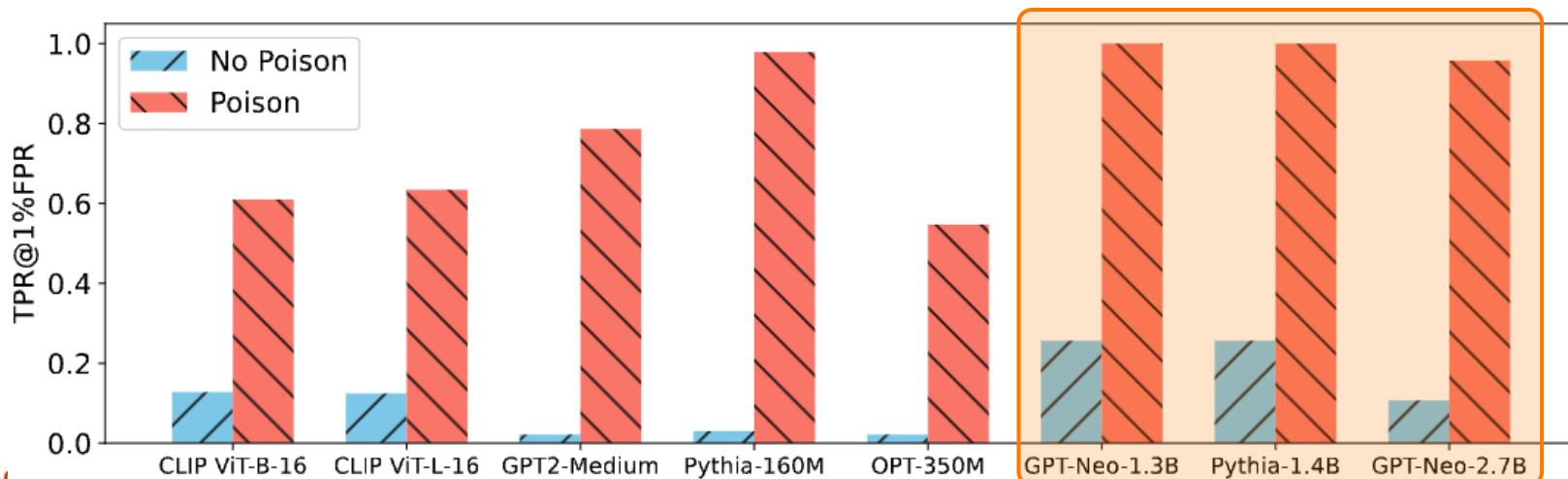
HOW EFFECTIVE ARE PRIVACY BACKDOORS?

- Privacy backdoor substantially increase membership risks
 - On CLIP ViT-B-32, **2.7 – 5.0x** increase in the attack success rate (TPR) at 1% FPR
 - On GPT-Neo 125M, **1.7 – 17.8x** increase in the attack success at 1% FPR
 - Our backdoors increase TPR on LMs to 87.4 – 96.3%

DATASET	ATTACK	TPR@1%FPR	AUC	VAL. LOSS BEFORE	VAL. LOSS AFTER
SIMPLE PII	-	0.242 \pm 0.030	0.874 \pm 0.008	3.99 \pm 0.00	3.19 \pm 0.00
	BACKDOOR	0.963 \pm 0.009	0.998 \pm 0.000	3.80 \pm 0.00	3.19 \pm 0.00
AI4 PRIVACY	-	0.049 \pm 0.013	0.860 \pm 0.005	3.99 \pm 0.00	3.19 \pm 0.00
	BACKDOOR	0.874 \pm 0.028	0.995 \pm 0.001	3.99 \pm 0.00	3.19 \pm 0.00
MIMIC-IV	-	0.560 \pm 0.025	0.916 \pm 0.003	4.52 \pm 0.03	1.57 \pm 0.02
	BACKDOOR	0.910 \pm 0.028	0.980 \pm 0.005	1.48 \pm 0.02	1.38 \pm 0.01

HOW EFFECTIVE ARE PRIVACY BACKDOORS?

- Privacy backdoor substantially increase membership risks
 - On CLIP ViT-B-32, 2.7 – 5.0x increase in the attack success rate (TPR) at 1% FPR
 - On GPT-Neo 125M, 1.7 – 17.8x increase in the attack success at 1% FPR
 - Our backdoors increase TPR on LMs to 87.4 – 96.3%
 - On large-language models, the attacks achieve over 90% TPR @ 1% FPR

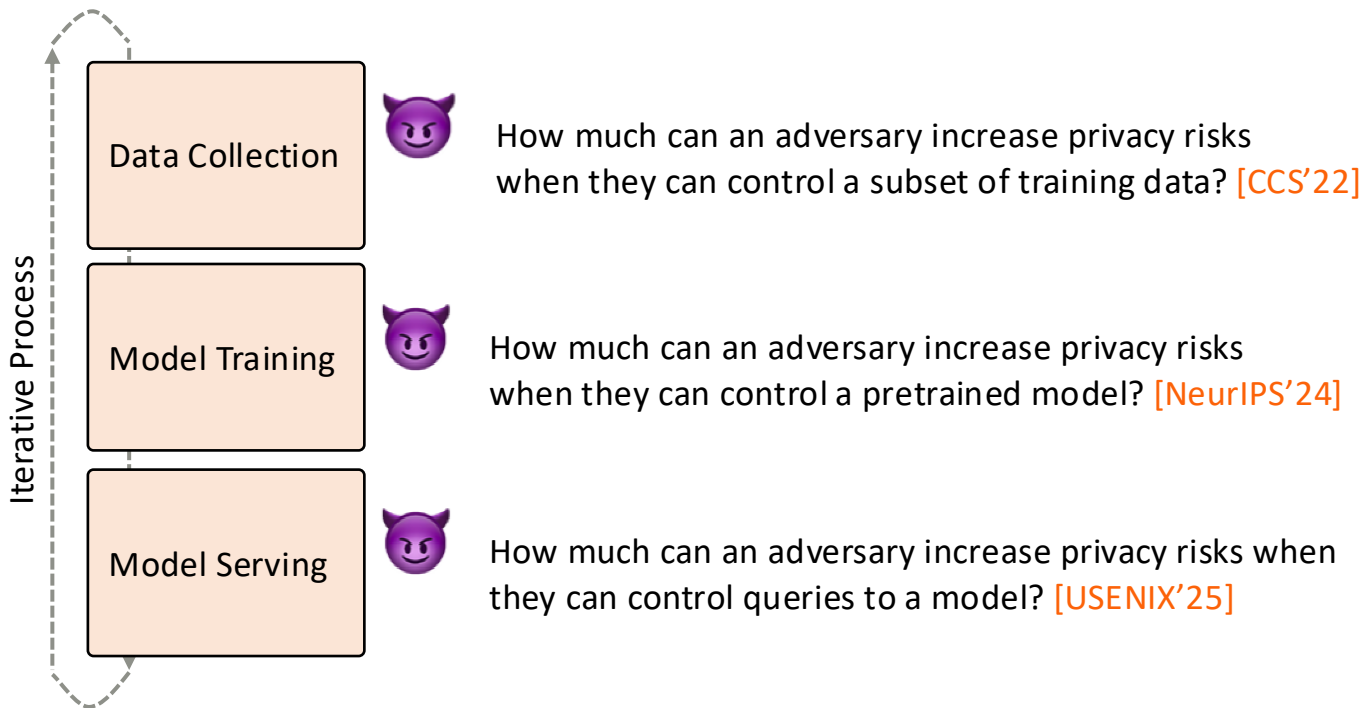


HOW EFFECTIVE ARE PRIVACY BACKDOORS?

- Impact of fine-tuning methods
 - We consider linear-probing, LoRA, QLoRA and fine-tuning on noisy-embeddings
 - Our backdoor is effective across these methods, versatile and prevalent
- Impact of inference-time strategies
 - We consider quantization, top-k (k=5) probabilities, and watermarking
 - Our attack increases the membership inference success consistently, substantially

PRIVACY IS A SYSTEMS SECURITY PROPERTY

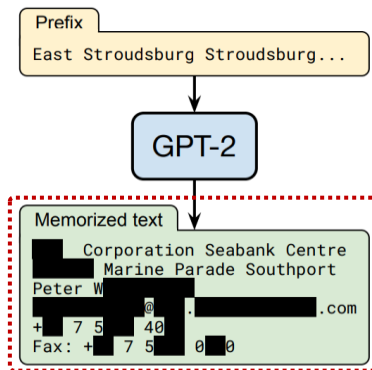
- Privacy risks emerge from the pipeline



HOW TO MEASURE THE MEMORIZATION OF LMS?

- Prior work proposes privacy attacks
 - Membership inference-style attacks¹
 - Data extraction attacks that generate partial (or entire) training examples verbatim²

Highest Likelihood Sequences	Log-Perplexity
The random number is 281265017	14.63
The random number is 281265117	18.56
The random number is 281265011	19.01
The random number is 286265117	20.65
The random number is 528126501	20.88
The random number is 281266511	20.99
The random number is 287265017	20.99
The random number is 281265111	21.16
The random number is 281265010	21.36



HOW TO MEASURE THE MEMORIZATION OF LMS?

- Most data extraction focuses on designing effective queries
 - BOS (beginning-of-sentence) token¹
 - Text corpus collected from the Internet¹
 - Text corpus generated through querying the target model²
 - (Auditing-only) Exact context associated with sensitive information³

BOS token: <BOS>

Internet text: "Please email me at"

Model generated text: "[MASK] wrote in his memoir that he had again developed pneumonia"

How much can an attacker amplify privacy risks through prompt opt.?

(Audit-only) Exact context: "Pt [SURNAME] 21 y/o male dx with..."

A NEW THREAT: AN AUTOMATED PROMPT SEARCH?

- Goal:
 - leak as many (unique) PII records as possible from training data
 - e.g., extract PII from OpenResearchLab's a developer-friendly API
- Knowledge:
 - No access to the target model parameters and their training data
 - Partial access to a public, pre-trained model's training dataset
 - Access to the log-probabilities of query outputs
- Capabilities
 - Create a surrogate model(s), and to this end, adversaries can fine-tune the public pre-trained models on their data
 - Make unlimited queries to the surrogate (i.e., white-box access)
 - Make limited queries to the target model (e.g., consider rate-limiting)

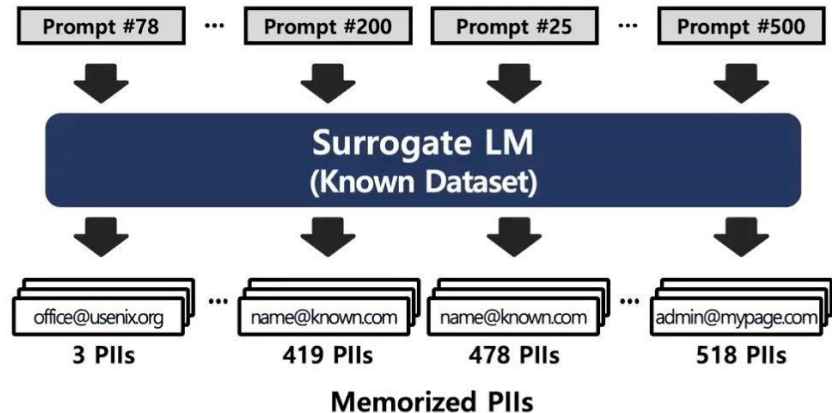
PRIVATE INVESTIGATOR: PROMPT-OPT-BASED PII EXTRACTION ATTACK

- Phase I: Generate effective prompts
 - These promising prompts elicit **many**, unique PII-like structures
 - Generate prompts that elicit PII-like structures from surrogate models
 - The surrogates are the pretrained ones or finetuned on the attacker's data

Algorithm 1: Generating prompts.

```
Input : # of prompts ( $n\_prompt$ ).
Output : A set of generated prompts ( $\mathbb{P}$ ).
1 function SearchPrompt( $n\_prompt$ )
2    $pii\_count\_1 \leftarrow \text{GetPIICount}(\mathbb{V}, 200)$ 
3    $\mathbb{T} \leftarrow \text{GetTop}(pii\_count\_1, |\mathbb{V}|/100)$ 
4    $pii\_count\_2 \leftarrow \text{GetPIICount}(\mathbb{T}, 2000)$ 
5    $\mathbb{P} \leftarrow \text{GetTop}(pii\_count\_2, 1)$ 
6   for  $n\_prompt - 1$  times do
7      $\mathbb{T} \leftarrow \mathbb{T} - \mathbb{P}$ 
8      $mean\_h \leftarrow \text{meanHidden}(x)$ 
9      $min\_sim \leftarrow \min_{x \in \mathbb{T}} \text{Sim}(\text{Hidden}(x), mean\_h)$ 
10     $new\_prompt \leftarrow \arg \max_{x \in \mathbb{T}} pii\_count\_2[x]$ 
11     $\text{Sim}(\text{Hidden}(x), mean\_h) \leq min\_sim + \theta$ 
12     $\mathbb{P} \leftarrow \mathbb{P} + new\_prompt$ 
13  return  $\mathbb{P}$ 
14 function GetPIICount( $\mathbb{S}, n\_text$ )
15   $pii\_count \leftarrow \{\}$ 
16  for  $x \in \mathbb{S}$  do
17     $texts \leftarrow \text{GenTexts}(x, n\_text)$ 
18     $pii\_count[x] \leftarrow \text{CountTrainPII}(texts)$ 
19  return  $pii\_count$ 
```

Possible Prompts



PRIVATE INVESTIGATOR: PROMPT-OPT-BASED PII EXTRACTION ATTACK

- Phase I: Generate **effective**, diverse prompts
 - Construct prompts that elicit PII-like structures from the surrogate models
 - **Promising prompts** do *not* always appear *directly* relevant to PII records

```
- user
- ' email '
- ' php '
- ById
- Authent
- upload
- Password
- email
- /*
- .</
- '":'
- route
- ' sender '
- ' Email '
- root
- database
- ' google '
- Update
- ' mailing '
- Server
```

← Prompts we generate
from the *public, pre-trained*
surrogate models

```
- '-----'
- ' sniff '
- ' believer '
- checking
- ' Sora '
- ' lobbyist '
- ' nonsense '
- ' villain '
- ' hostile '
- ' terror '
- worthiness
- ' psyche '
- ' congest '
- ' NYPD '
- ' Canaan '
- ' bordering '
- ' emperor '
- closed
- oso
- ' courthouse '
```

← Prompts we generate
from the *surrogate fine-*
tuned on the attacker's data

PRIVATE INVESTIGATOR: PROMPT-OPT-BASED PII EXTRACTION ATTACK

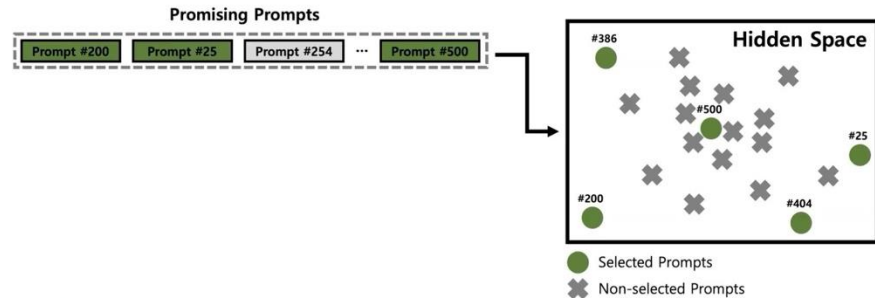
- Phase I: Generate effective, **diverse** prompts
 - These promising prompts elicit many, **unique** PII-like structures
 - Choose the prompts that cover *diverse regions of the surrogate's activation space*
 - Compute cosine similarity to compare the alignment between prompts in the space

Algorithm 1: Generating prompts.

```
Input : # of prompts ( $n\_prompt$ ).
Output : A set of generated prompts ( $\mathbb{P}$ ).
1 function SearchPrompt( $n\_prompt$ )
2    $pii\_count\_1 \leftarrow \text{GetPIICount}(\mathbb{V}, 200)$ 
3    $\mathbb{T} \leftarrow \text{GetTop}(pii\_count\_1, |\mathbb{V}|/100)$ 
4    $pii\_count\_2 \leftarrow \text{GetPIICount}(\mathbb{T}, 2000)$ 
5    $\mathbb{P} \leftarrow \text{GetTop}(pii\_count\_2, 1)$ 
6   for  $n\_prompt - 1$  times do
7      $\mathbb{T} \leftarrow \mathbb{T} - \mathbb{P}$ 
8      $mean\_h \leftarrow \text{meanHidden}(x)$ 
9      $min\_sim \leftarrow \min_{x \in \mathbb{T}} \text{Sim}(\text{Hidden}(x), mean\_h)$ 
10     $new\_prompt \leftarrow \arg \max_{x \in \mathbb{T}} pii\_count\_2[x]$ 
11     $\text{Sim}(\text{Hidden}(x), mean\_h) \leq min\_sim + \theta$ 
12     $\mathbb{P} \leftarrow \mathbb{P} + new\_prompt$ 
13  return  $\mathbb{P}$ 
14 function GetPIICount( $\mathbb{S}, n\_text$ )
15    $pii\_count \leftarrow \{\}$ 
16   for  $x \in \mathbb{S}$  do
17      $texts \leftarrow \text{GenTexts}(x, n\_text)$ 
18      $pii\_count[x] \leftarrow \text{CountTrainPII}(texts)$ 
19  return  $pii\_count$ 
```

• Diverse Prompts

- Sparse in the surrogate LM's **hidden space**.
- Covering diverse context.

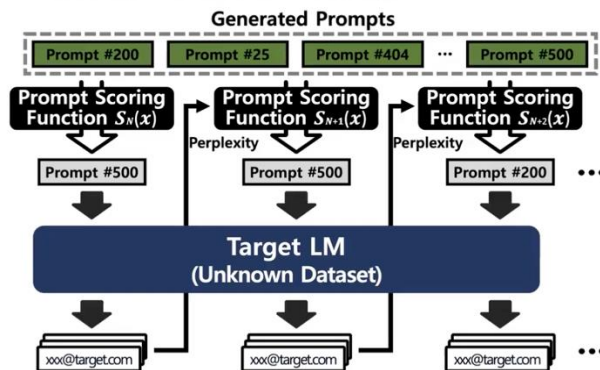


PRIVATE INVESTIGATOR: PROMPT-OPT-BASED PII EXTRACTION ATTACK

- Phase II: PII extraction (attack campaign)
 - Use the promising prompts from Phase I to extract PIIs from the target model
 - Choose a promising prompt based on two criteria
 - The perplexity of generated texts that may contain PII records
 - The number of times that prompt has been chosen
 - Run 100 extraction attacks; each attack generates 2,000 texts from a chosen prompt

• Prompt Selection Strategy

Select the most effective prompt on each Nth PII extraction attempt.



• Prompt Scoring Function

$$S_N(x) = \sqrt{\frac{\ln N}{n_x}} - c \cdot PII_Perplexity_x$$

Exploration Exploitation

n_x : The number of times prompt x was selected
 $PII_Perplexity_x$: Average perplexity of PIIs extracted by prompt x

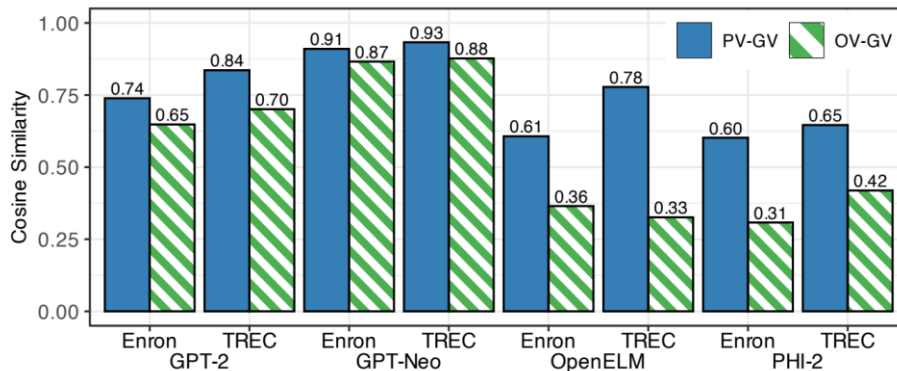
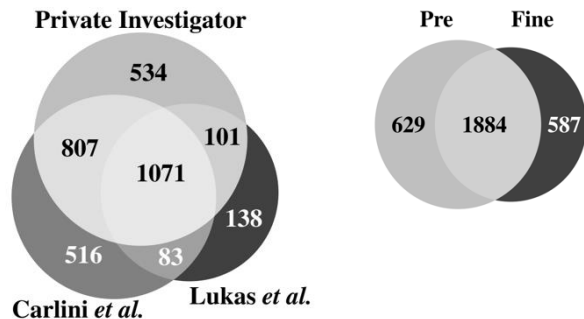
HOW EFFECTIVE IS THE PRIVATE INVESTIGATOR?

- Private investigator extracts *more* PII records than the baselines
 - Our attack uses promising prompts generated from the public, GPT-Neo model
 - Our attack remains effective across PII types and model families/sizes
 - Promising prompts *consistently* scale to multi-billion parameter models

	GPT-Neo			PHI-2		
	Email	Phone	Name	Email	Phone	Name
Carlini*	2477	1946	24359	5732	2505	34780
Lukas**	1393	1741	20770	5119	2323	33066
Ours	2513	2008	24616	6079	2954	36385

HOW EFFECTIVE IS THE PRIVATE INVESTIGATOR?

- Private investigator extracts *unique* PII records than the baselines
 - PII records generated from our attack differ from Carlini et al. and Lukas et al.
 - Our attack on the pre-trained surrogate extracts different PII than the fine-tuned model
 - The prefixes generated from our attack show a higher cosine similarity than those generated from the baseline methods (in the activation space)



TAKEAWAYS

- Memorization does *not* mean privacy risks
 - Privacy fails when memorized data becomes **observable at outputs**
 - Privacy is about what an adversary can infer at the model output
- Adversary can exploit ML pipeline
 - Poisoning turns a sample with avg-case privacy leakage into a **worst-case** one
 - Backdooring implies that pre-trained models *should not* be trusted as-is
 - Black-box, prompt-optimization is feasible and should be blocked
- Our work raises the bar for defending against privacy attacks
 - Privacy failures are not model failures, but **system failures**
 - Poisoning can be used to **audit worst-case privacy leaks**

Thank You!

Sanghyun Hong

<https://secure-ai.systems/courses/MLSec/current>



Oregon State
University



TRUE AI
Trustworthy and Responsible AI