

# Auditing and Mitigating Safety Risks in Large Language Models

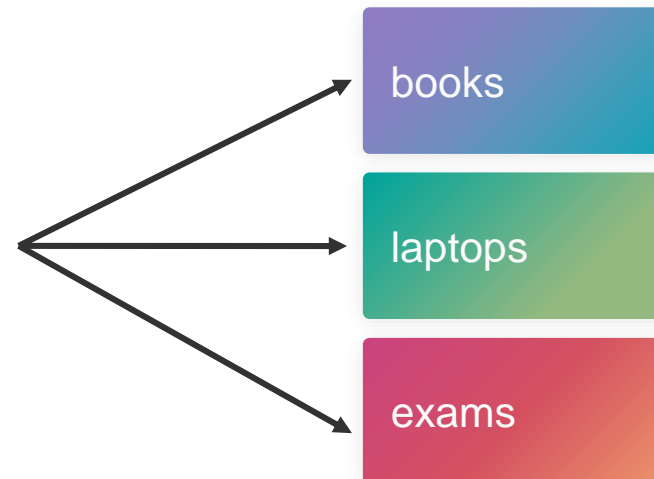
**Niloofer (Fatemeh) Miresghallah**

# Act I: The LLM Takeover?

# What are Language Models?

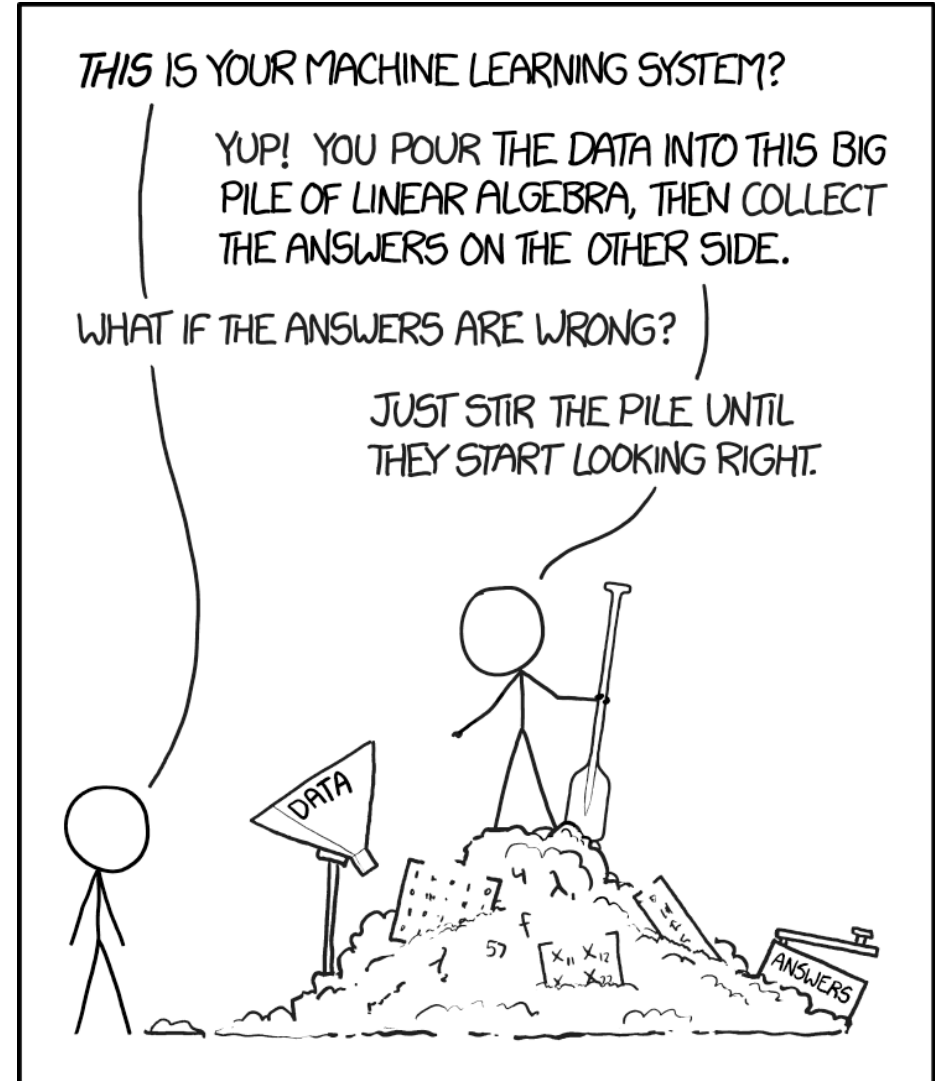
- A language model is a **probability distribution** over sequences of words
- Model what words a given word/context normally appears with

The students opened their \_\_\_\_\_.



# Large Language Models (LLMs)

- Transformer-based language models are often referred to as 'Large LMs' due to their **parameter count** (ranging from 100s of million to billions of parameters)
- Deployed with **Pre-train** and **Fine-tune** paradigm



# Large Language Models: The Good and the Bad ...

- Large language models are very good at **generating text**



# Large Language Models: The Good and the Bad ...

- Large language models are very good at **generating text** and **learning representations**.



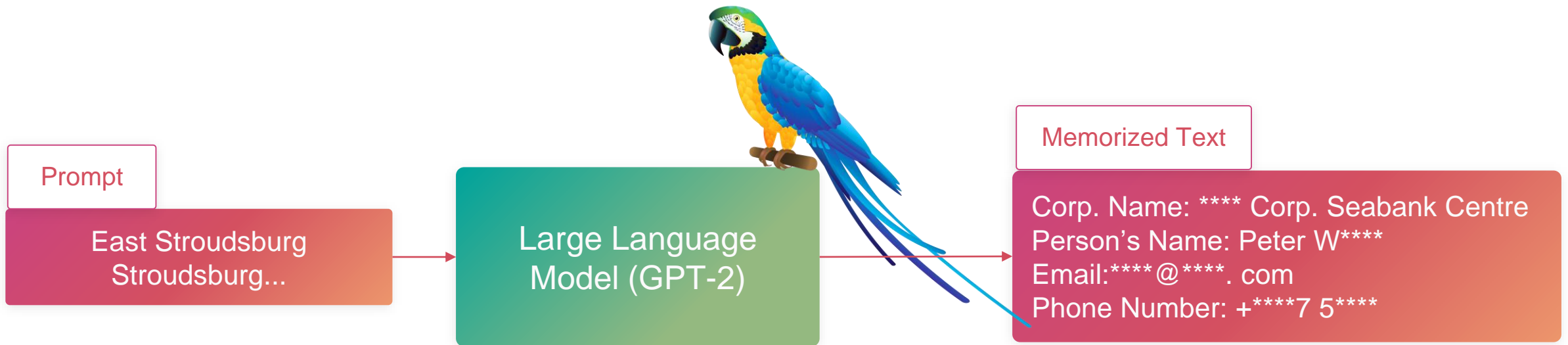
# LLMs: The Bad



WHEN YOU TRAIN PREDICTIVE MODELS  
ON INPUT FROM YOUR USERS, IT CAN  
LEAK INFORMATION IN UNEXPECTED WAYS.



# LLMs: The Bad





# LLMs: The Bad

- LLMs can also *regurgitate data* they have seen before, creating *privacy risks*.

Title:

*Hi everyone, my name is Anish Athalve and I'm a PhD student at Stanford University.*

# LLMs: The Bad

- LLMs can also *regurgitate data* they have seen before, creating *privacy risks*.

Title:

*Hi everyone, my name is Anish Athalye and I'm a PhD student at Stanford University.*

<https://www.anish.io> :

**Anish Athalye**

I am a PhD student at MIT in the PDOS group. I'm interested in formal verification, systems, security, and machine learning.

GitHub: @anishathalye

Blog: anishathalye.com

LLMs are *not ready* to be widely deployed in  
*safety critical scenarios* as is!

# In this talk:

## Question 1: How can we **audit and quantify safety risks** of LLMs?

- [ACL 2023] Membership Inference Attacks via Neighbourhood Comparison
- [EMNLP2022a] Quantifying Privacy Risks of Masked Language Models Using MIAs
- [EMNLP2022b] Memorization in NLP Fine-tuning Methods
- [FAccT2022] What does it mean for language models to preserve privacy?

## Question 2: How can we **limit the risks** of LLMs?

- [ACL2023] Privacy-Preserving Domain Adaptation of Semantic Parsers
- [NeurIPS2022] Differentially private model compression
- [NAACL2021] Joint privacy-utility optimization in language models



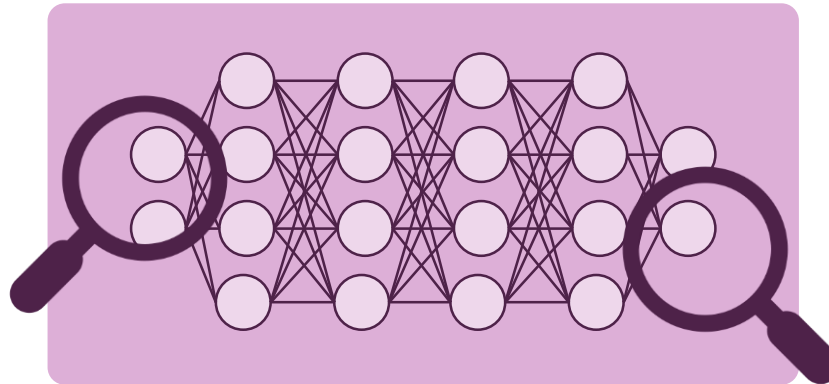
Don't repeat this!!

# Act II: Auditing LLMs for Privacy



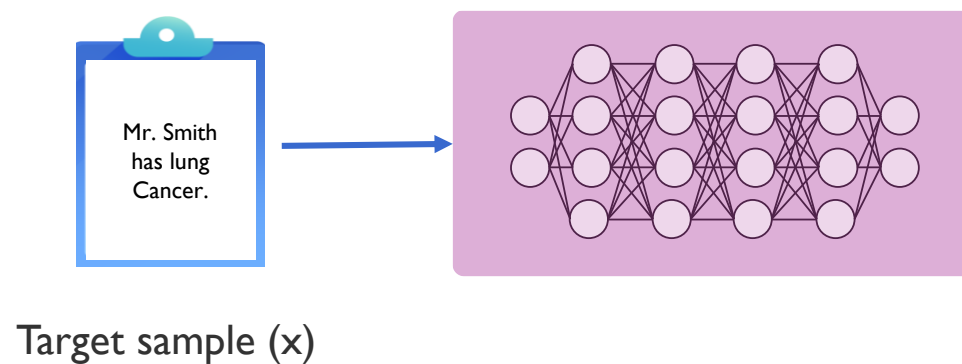
# What is information leakage in an ML model?

- ‘Leakage’ is being able to **learn information about the training data**, which cannot be learned from other models/data (from the same distribution)



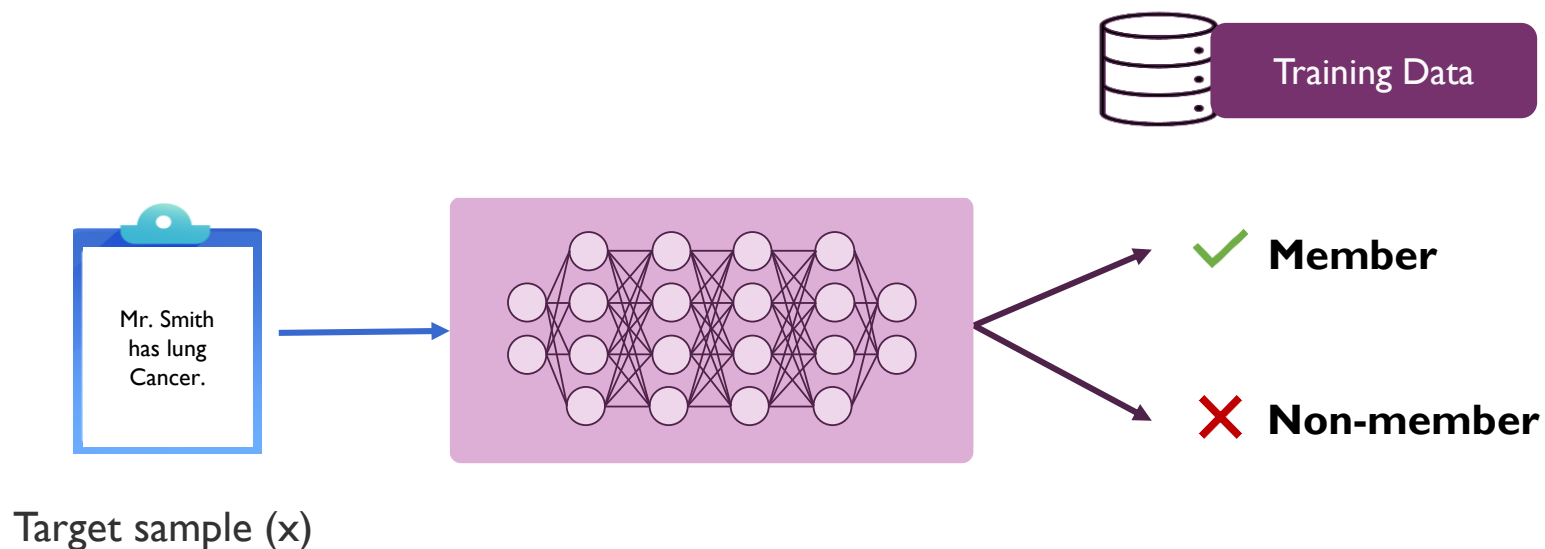
# Measuring Leakage: Membership Inference Attacks

- Can an adversary infer whether a particular data point “x” is part of the training set?



# Measuring Leakage: Membership Inference Attacks

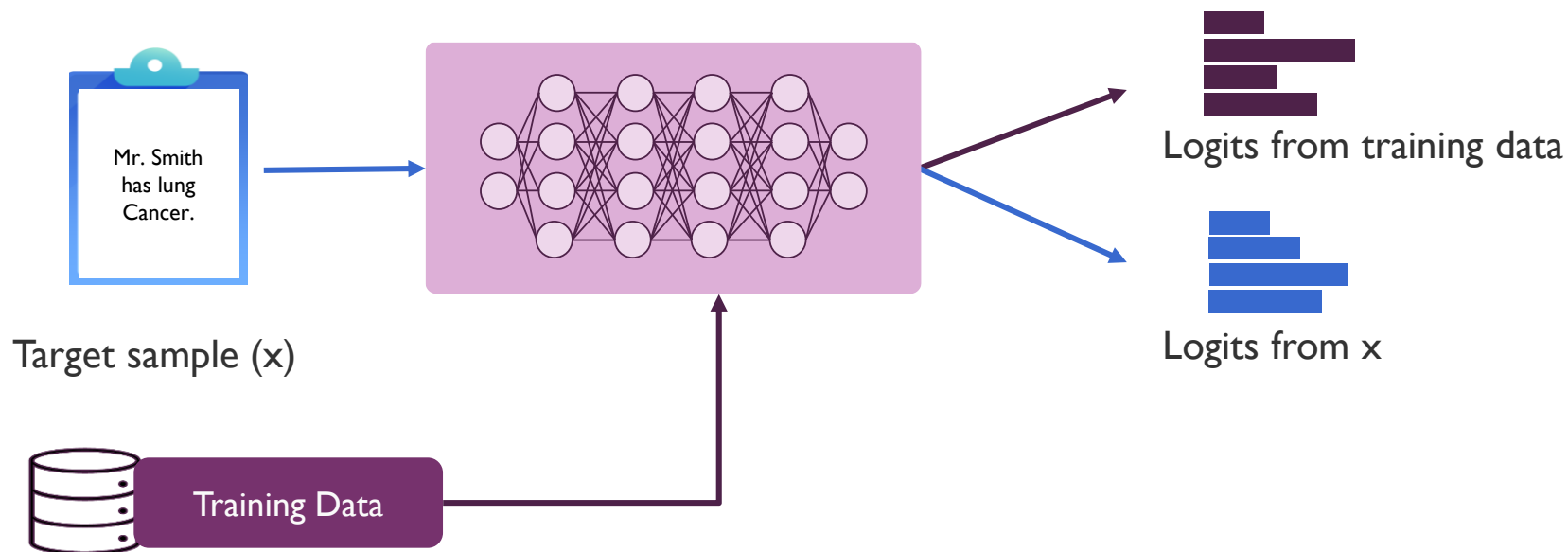
- Can an adversary infer whether a particular data point “x” is part of the training set?





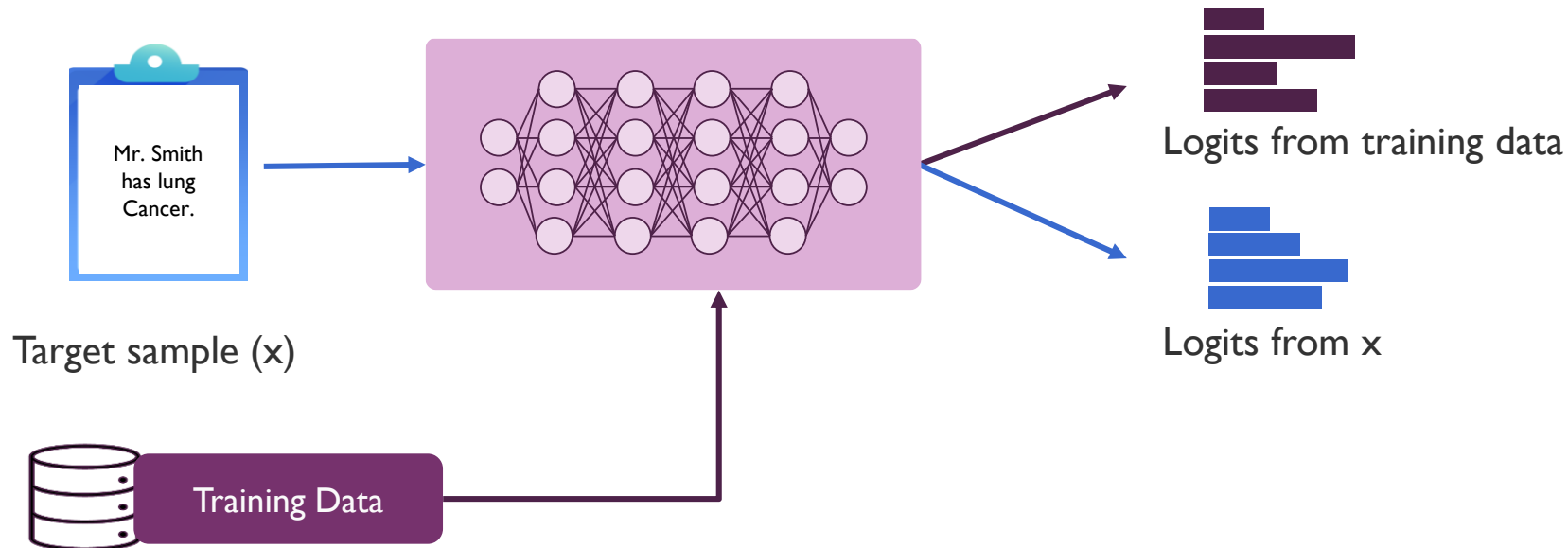
# Measuring Leakage: Membership Inference Attacks

- Can an adversary infer whether a particular data point “x” is part of the training set?



# Measuring Leakage: Membership Inference Attacks

- Can an adversary infer whether a particular data point “x” is part of the training set?
- Success of attacker is a metric to quantify information leakage of the model about its individual training data



# Background: Membership Inference Attacks

- Membership Inference Attacks (MIAs): Loss-based attack
- Stronger MIAs: Reference-based attacks (MIA) [Miresghallah2022, Ye2021, Carlini2022]
  - A **static, absolute threshold** does not control for the **intrinsic complexity** of each utterance
  - We need to **calibrate** the threshold for each utterance

# Reference-based attack

We propose a reference-based attack:

- Complex training points: points that have **higher loss**

# Reference-based attack

We propose a reference-based attack:

- Complex training points: points that have **higher loss**

Training data point	Target Model Loss
Mr. Smith has type 2 diabetes.	3
Mr. Smith has fever .	2
Mr. Smith is taking 5 mgs of Haloperidol 2 times a day.	7



# Reference-based attack

We propose a reference-based attack:

- Complex training points: points that have **higher loss**

Training data point	Target Model Loss
Mr. Smith has type 2 diabetes.	3
Mr. Smith has fever .	2
Mr. Smith is taking 5 mgs of Haloperidol 2 times a day.	7

# Reference-based attack

We propose a reference-based attack:

- Complex training points: points that have **higher loss**
- We use a **reference** model, to provide an insight into **how difficult each data point is**

Training data point	Target Model Loss
Mr. Smith has type 2 diabetes.	3
Mr. Smith has fever .	2
Mr. Smith is taking 5 mgs of Haloperidol 2 times a day.	7

# Reference-based attack

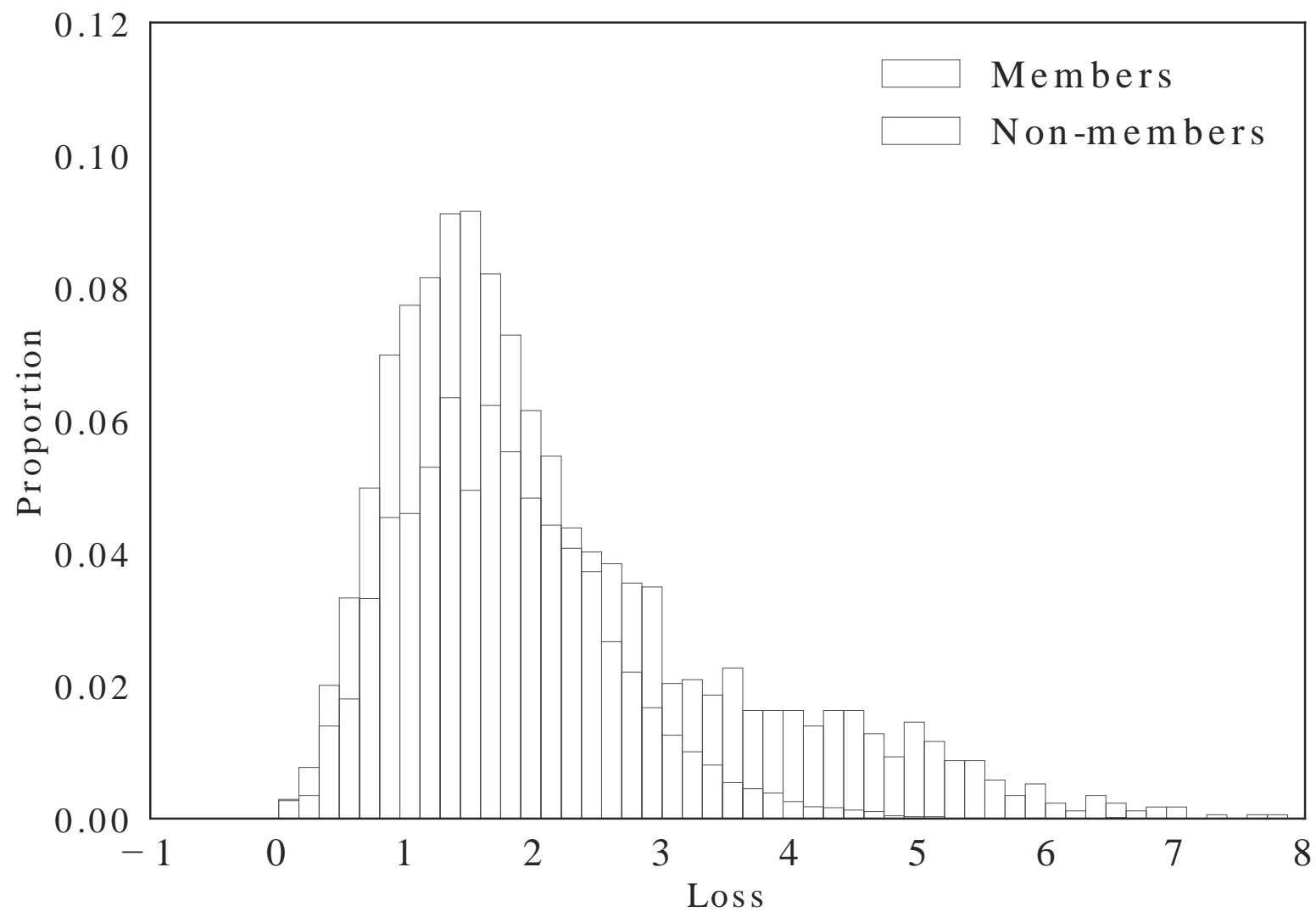
We propose a reference-based attack:

- Complex training points: points that have **higher loss**
- We use a **reference** model, to provide an insight into **how difficult each data point is**

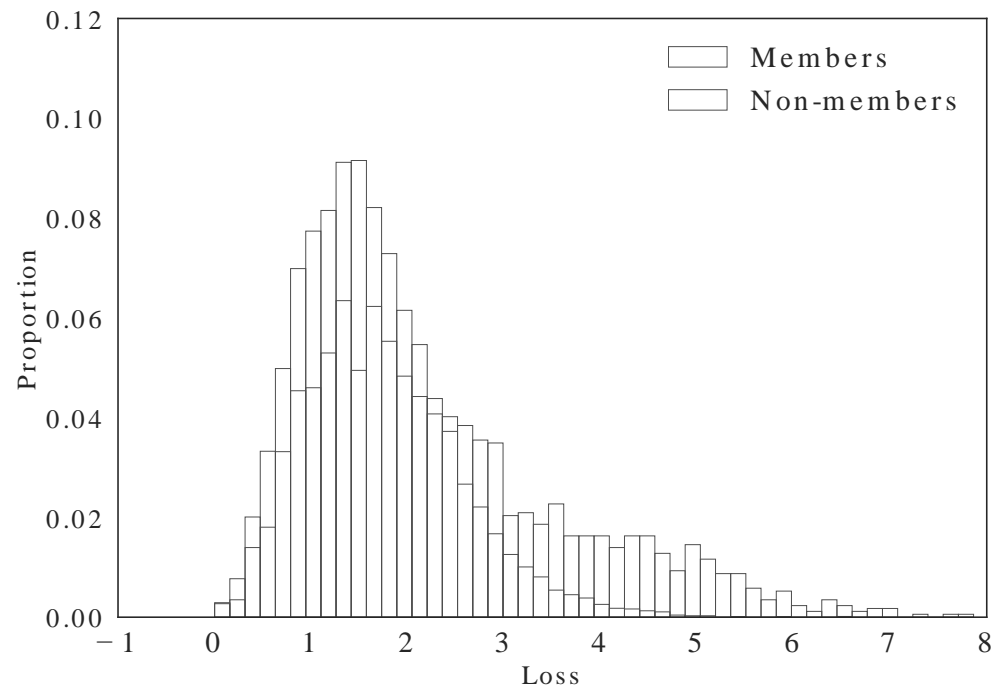
Training data point	Target Model Loss	Reference Model Loss
Mr. Smith has type 2 diabetes.	3	4
Mr. Smith has fever .	2	3
Mr. Smith is taking 5 mgs of Haloperidol 2 times a day.	7	10



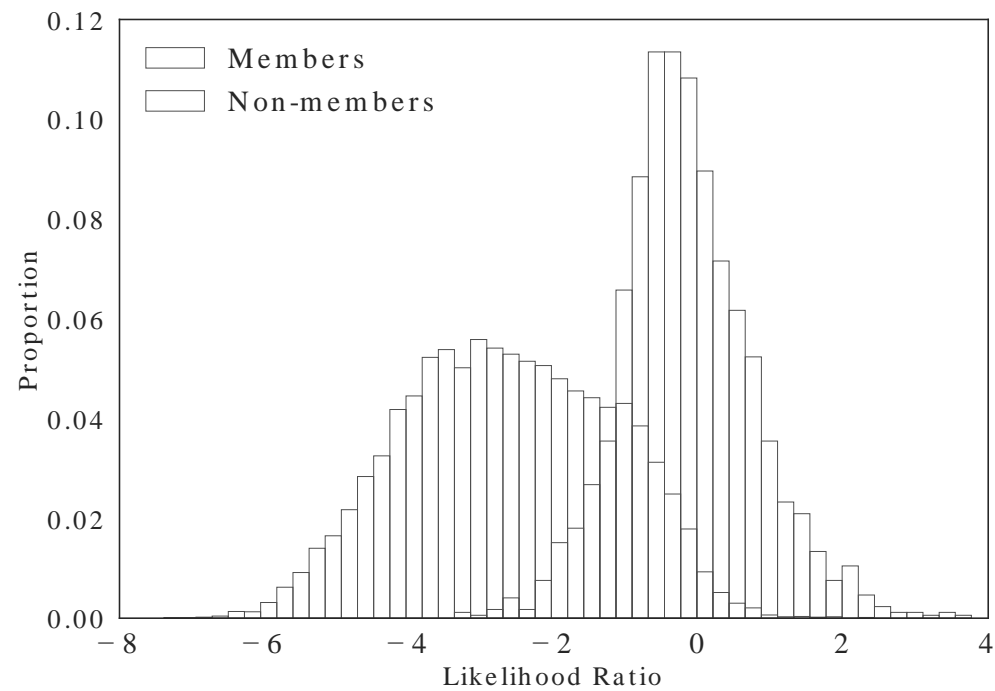
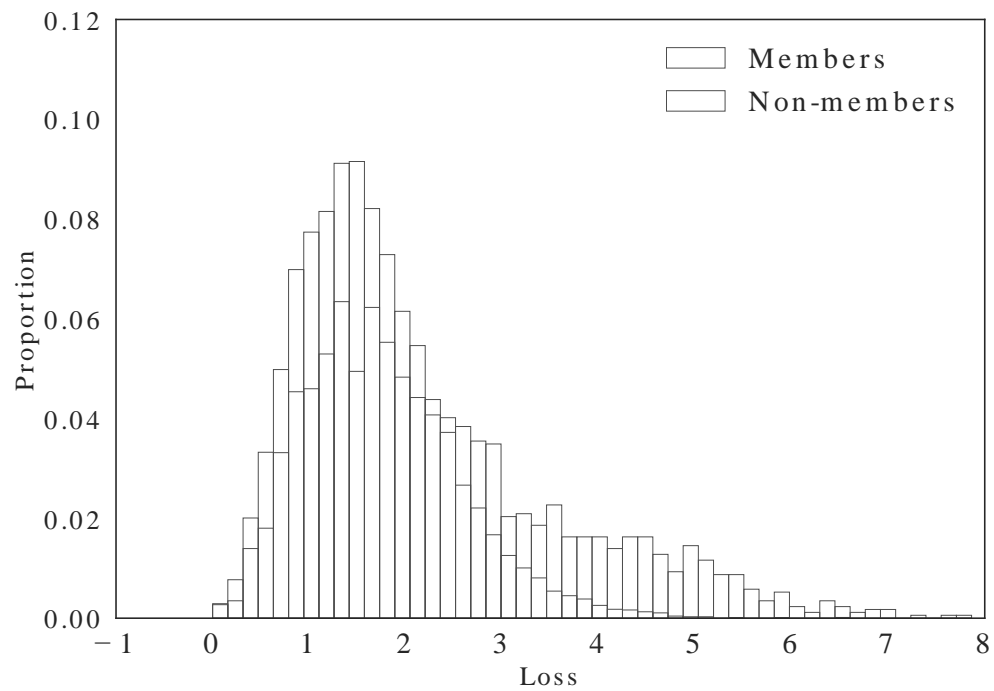
# Example: loss-based attack



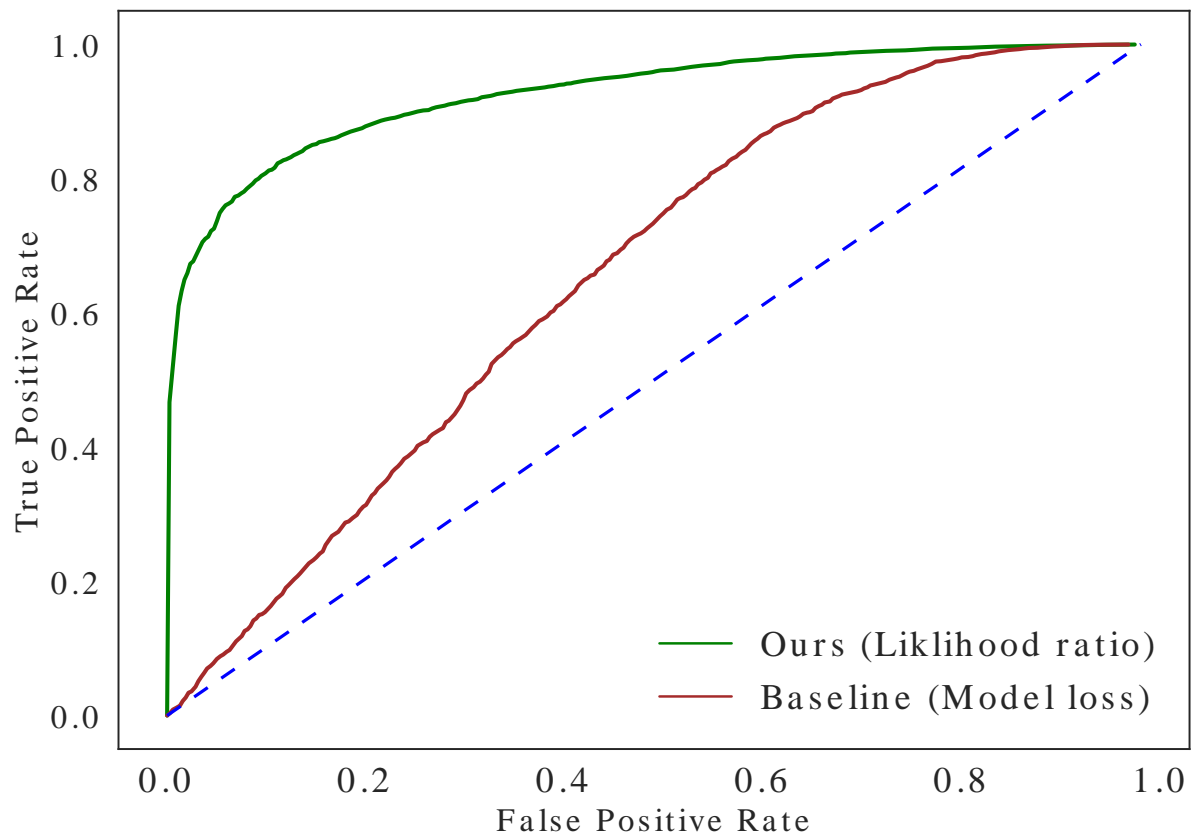
# Example: loss-based attack



# Example: Reference-based attack



# ROC Curve Results



Our likelihood ratio-based attack has an AUC of 0.90, vs the 0.66 of the loss-based attack.

## However ...

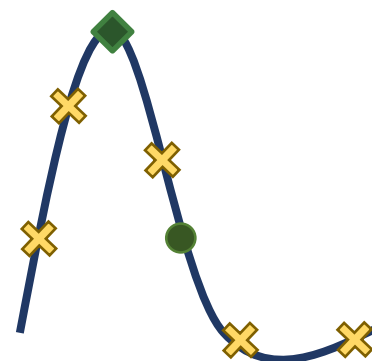
- The success of reference-based attacks is contingent upon having a **'good reference'** model, which is not always feasible:
  - We might have a very **small dataset**, therefore holding out part of the data to train a reference model on would significantly impact the utility of the final model
  - We might have **limited/no information about the training data** of the model we are probing, therefore curating non-overlapping, similar data would be non-trivial
  - We might not have **access to enough compute** to train large reference models

How can we leverage the loss function and its curvature to determine membership?

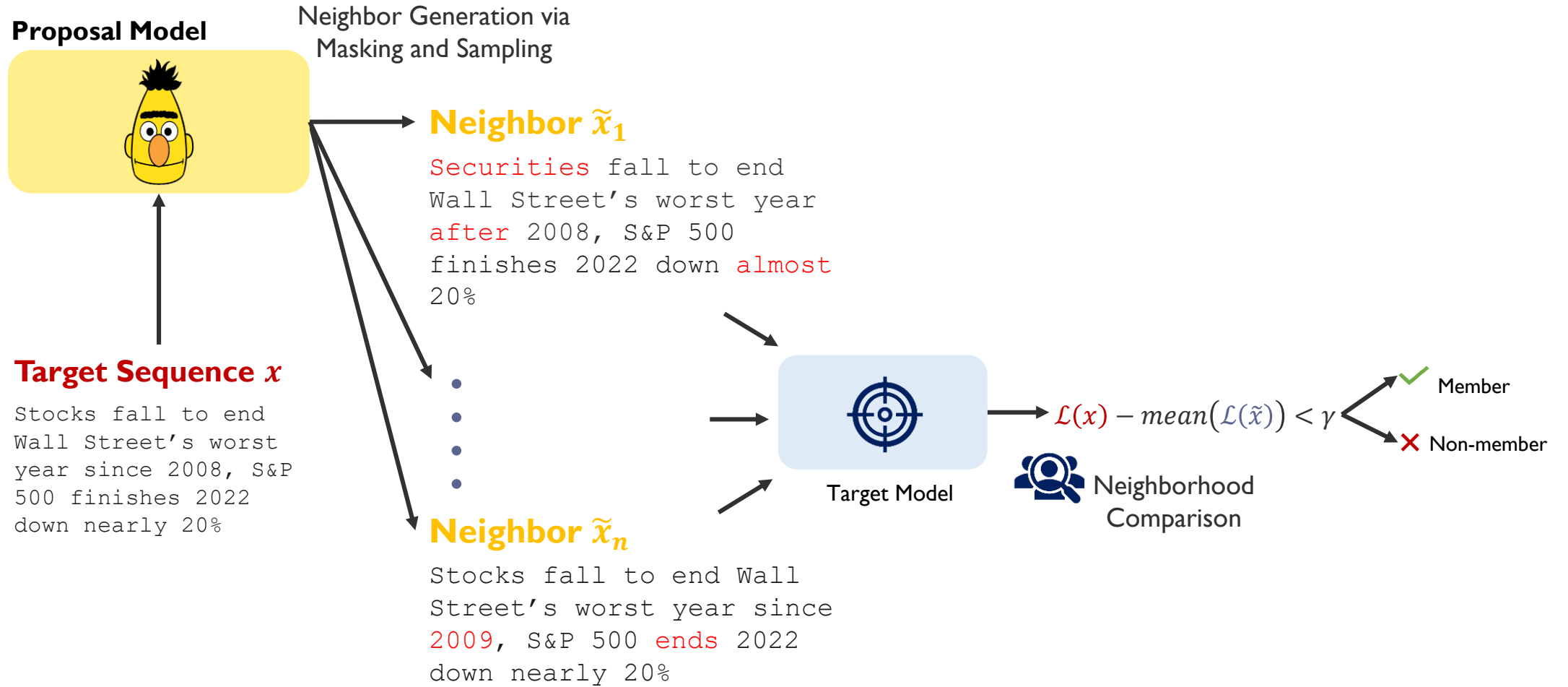
# Proposed: Neighbourhood Comparison-based Attacks

- Instead of **likelihood ratio**, we use **local-optimality** of each point as a signal to determine membership. The intuition is:
  - If a data point is part of the training-set, its likelihood would be **locally optimal, compared to its neighboring points**
  - If a data point is not part of the training set, then there would be points its neighborhood with **both higher and lower likelihoods**

Target Model Likelihood	—
Neighbor	×
Training point	◆
Non-training point	●



# Attack Procedure



# Experimental Setup

- We are mounting a membership inference attack on fine-tuned GPT2
- Baseline: Likelihood-ratio based attack
  - Base reference: Pre-trained, non-finetuned model
  - Candidate reference: fine-tuned GPT2, but on a dataset with small distribution shift
  - Oracle reference: fin-tuned GPT2 on a dataset with the same distribution as target model





## Does this really work?

	<b>False Positive Rate</b>	<b>0.1</b>
<b>Attack Method</b>	Base Reference	0.91
	Candidate Reference	0.95
	Oracle Reference	<b>3.76</b>
	<b>Neighborhood (Ours)</b>	1.73

As we step into lower false-positive rate (more precise) attack scenarios, we see that our method outperforms the likelihood ratio based attack.

## Does this really work?

	<b>False Positive Rate</b>	<b>0.1</b>	<b>0.01</b>
<b>Attack Method</b>	Base Reference	0.91	0.16
	Candidate Reference	0.95	0.15
	Oracle Reference	<b>3.76</b>	0.16
	<b>Neighborhood (Ours)</b>	1.73	<b>0.29</b>

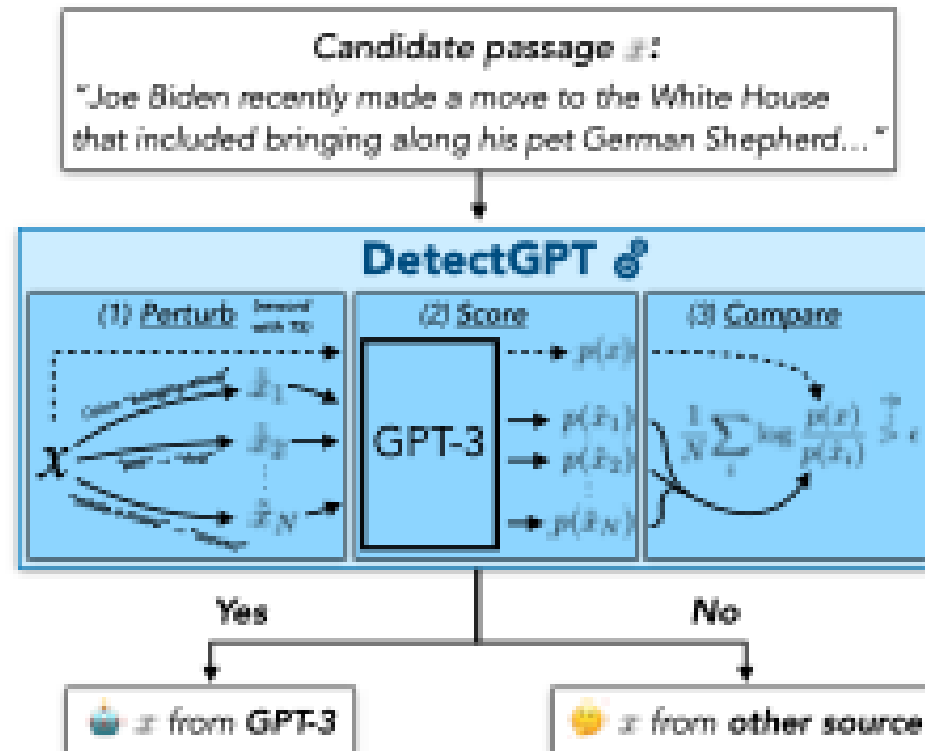
As we step into lower false-positive rate (more precise) attack scenarios, we see that our method outperforms the likelihood ratio based attack.

# Experimental Results: Other Experiments

1. Other Datasets:
  - AG News, NewsCatcher, Twitter, Wikipedia
2. Ablations:
  - Number of Generated Neighbours
  - Number of Word Replacements
3. Mitigations
  - Differentially Private SGD

# Detour: Relation to Machine-generated Text Detection

- Concurrent work: DetectGPT -- Mitchell et al. demonstrate that the same type of algorithm could be used to **distinguish between human written text and machine generated text.**



# So far ...


- We show that using **a reference model** can improve the performance of existing attacks, and uncover **higher levels of memorization**.
- We also demonstrate **reference-free methods**, that can be used in scenarios where **access to a reference is infeasible**.
- How can we **mitigate these privacy risks**, specifically by generating synthetic data?



# Act III: Limiting the Privacy Risks of LLMs

# Problem Definition

Task-oriented dialogue systems often assist users with **personal** or **confidential** matters

- Data is private and practitioners are not allowed to look at it 
- How can we know where the system is failing and needs **more training data** or **new functionality**?



Could you tell me what the weather is gonna be like today in New York?

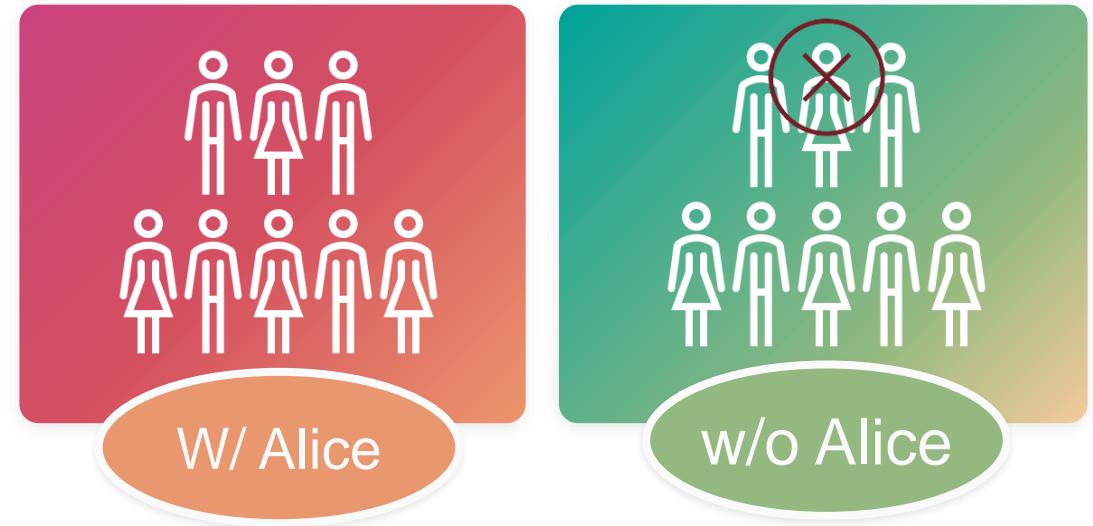


Email everyone who declined the invitation, saying ...

# Background: Differential Privacy

- DP protects the **membership of every single sample** in the training data
- A randomized algorithm  $A$  satisfies  $\epsilon$ -DP, if for all databases  $D$  and  $D'$  that differ in data pertaining to one user, and for every possible output value  $Y$ :

$$\frac{\Pr[A(D) = Y]}{\Pr[A(D') = Y]} \leq e^\epsilon.$$





# Private Training of Large Language Models: Prior Work

- To limit the leakage of fine-tuning data, prior work [Li et al. 2022, Yu et al. 2022] has used DP-SGD during fine-tuning



Don't repeat this!!

# Private Training of Large Language Models: Prior Work

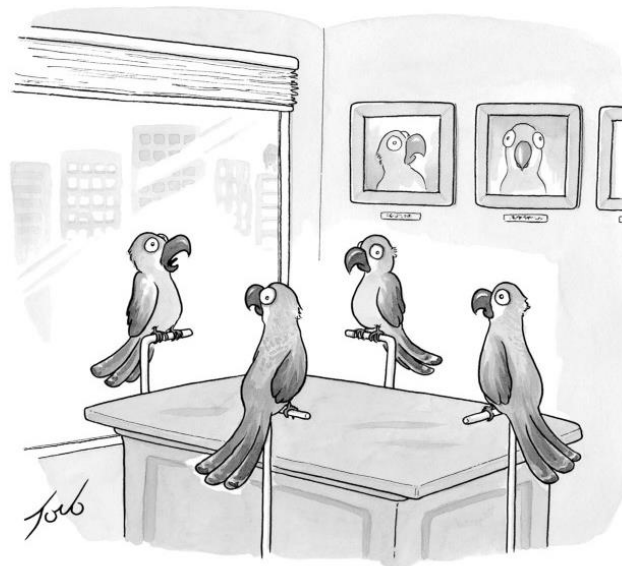
- To limit the leakage of fine-tuning data, prior work [Li et al. 2022, Yu et al. 2022] has used DP-SGD during fine-tuning
  - **Differential Privacy SGD (DP-SGD)** is the gold standard of private training



Don't repeat this!!

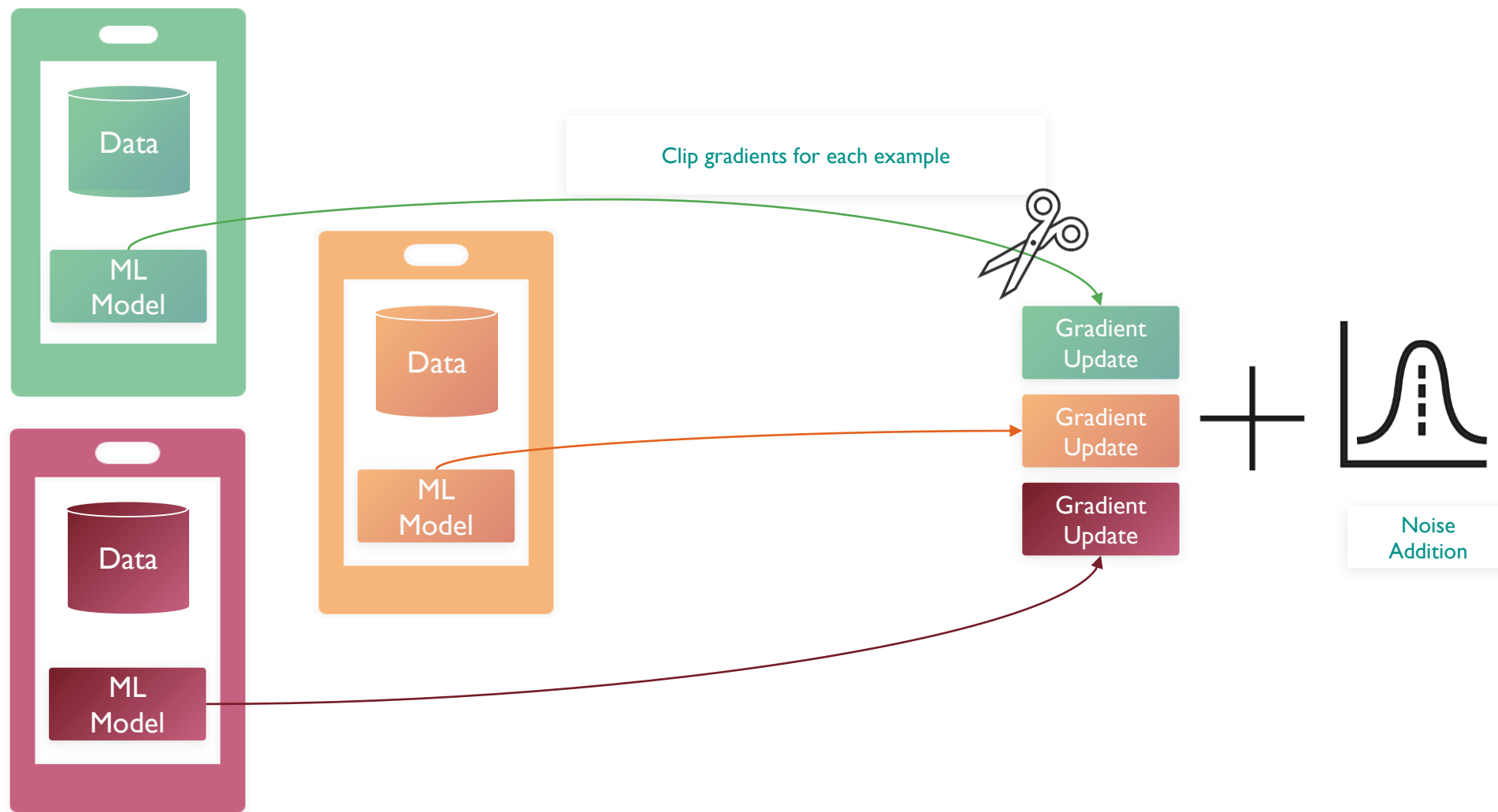
# Private Training of Large Language Models: Prior Work

- To limit the leakage of fine-tuning data, prior work [Li et al. 2022, Yu et al. 2022] has used DP-SGD during fine-tuning
  - **Differential Privacy SGD (DP-SGD)** is the gold standard of private training
    - DP protects the **membership of every single sample** in the training data



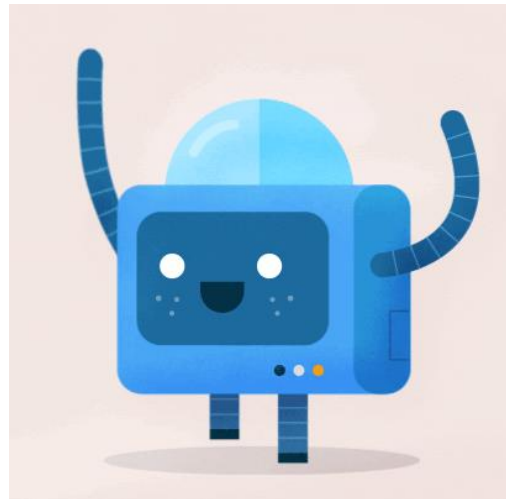
Don't repeat this!!

# Differentially Private SGD



## Problem Definition: Adding New Functionality

- Why not just **fine-tune** on the eyes-off data **privately**?
  - If some users are asking the system to hop up and down, fine-tuning is unlikely to make it grow legs.

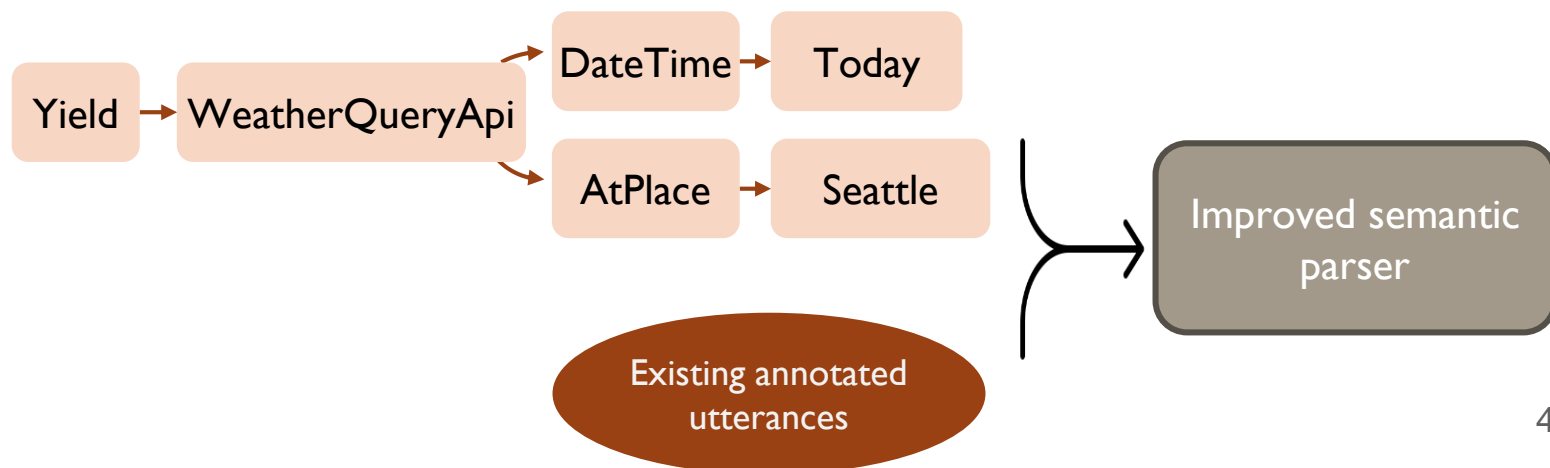


# Problem Definition: Adding New Functionality

- Why not just **fine-tune** on the eyes-off data **privately**?
  - If some users are asking the system to hop up and down, fine-tuning is unlikely to make it grow legs.



What is the weather like in Seattle Today?



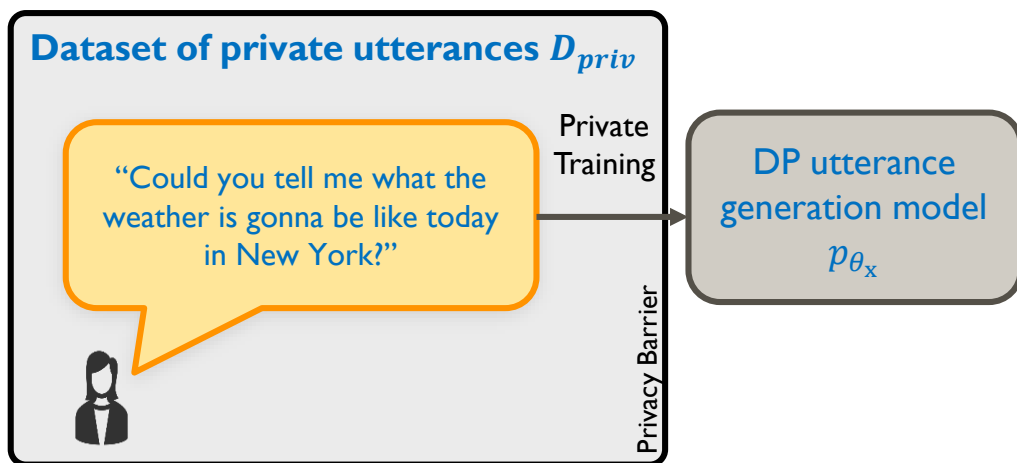
## Problem Definition: Adding New Functionality

- Why not just **fine-tune** on the eyes-off data **privately**?
  - If some users are asking the system to hop up and down, fine-tuning is unlikely to make it grow legs.
  - We need to be able to **look at synthesized data** to identify additional needed functions, then **annotate** with new functions and **add** to the training data to **improve the semantic parser**.

How can we privately synthesize data that is distributionally close to eyes-off user data?

# Baseline: Private Fine-Tuning of a Generative Model

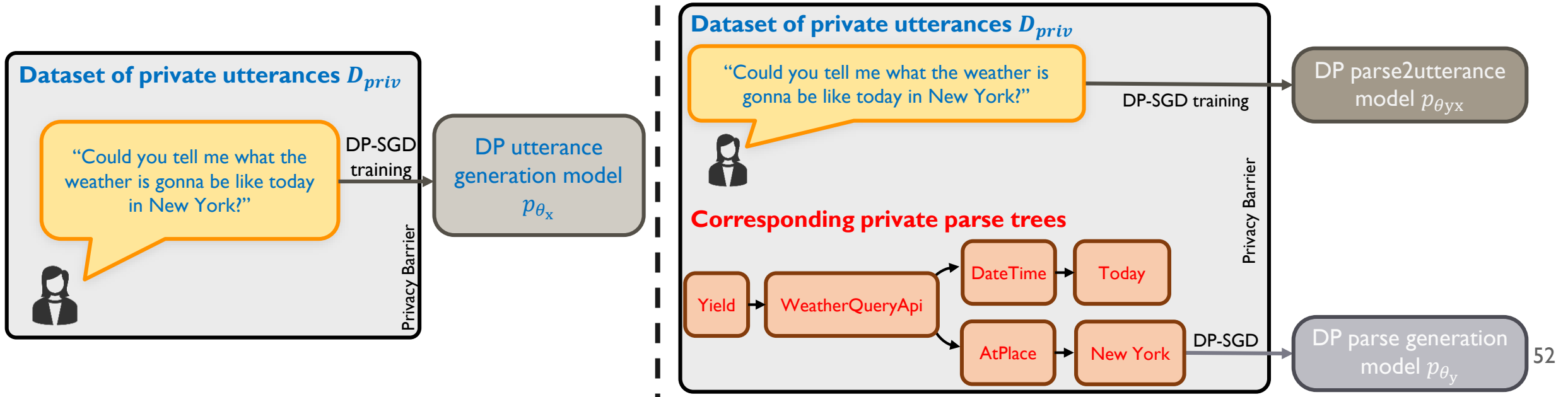
- Intuitive Baseline: We model  $p(x)$ , where  $x$  is a **private utterance**.





# Proposed: 2-stage Modeling of Intermediate Variables

- Intuitive Baseline: We model  $p(x)$ , where  $x$  is a **private utterance**.
- Proposed: We model  $p(y)$  and  $p(x|y)$ , where  $y$  is a **private parse-tree**.
  - one stage models the **parse-trees**,  $p_{\theta_y}$
  - The other stage models an **utterance** given a **parse-tree**,  $p_{\theta_{yx}}$

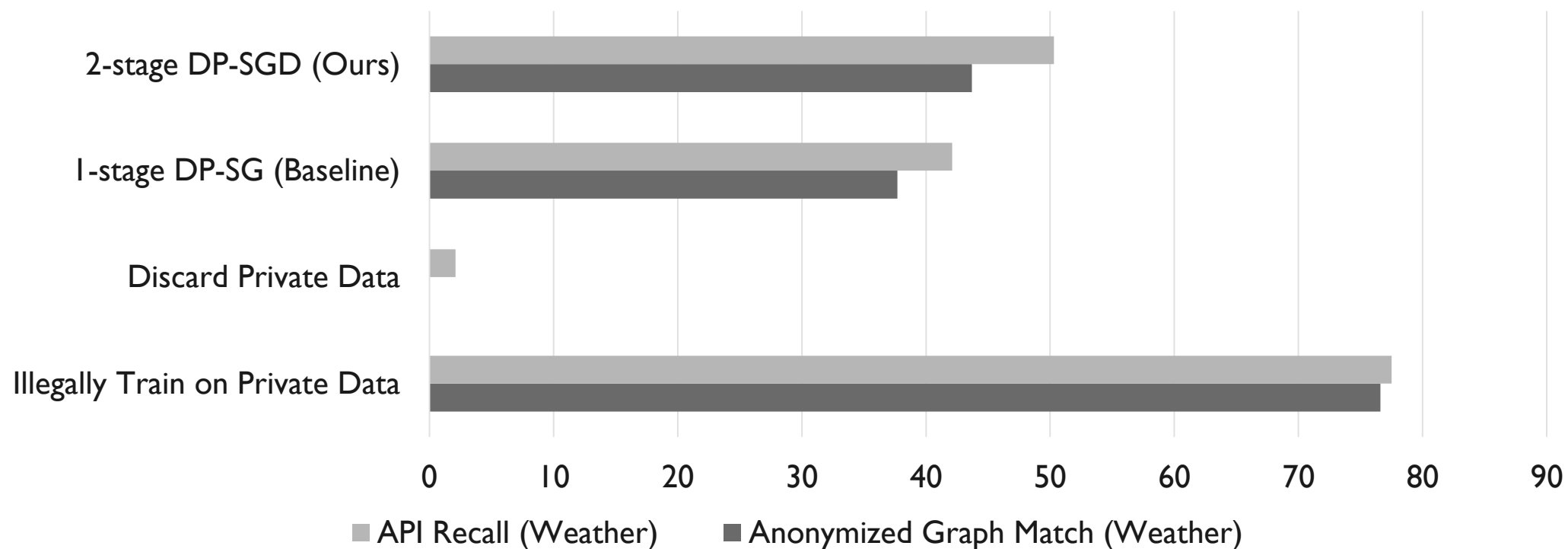


# Does This Really Work?

We simulated a situation where users are asking about the **weather** but the original semantic parser **was not trained on weather-related functions**:

1. We created the original semantic parser by training on  $\frac{1}{10}$  of our data (SMCalFlow), **excluding** any examples that use **weather-related functions**.
2. We treated the other  $\frac{9}{10}$  of the data as **private user utterances**, **including** those requesting **weather**. We created **approximate private annotations** for the private utterances, using the **original semantic parser**.
3. We apply the baseline and proposed methods to create **public synthesized datasets**, which **include weather functions**.
4. We simulated high-quality human annotation of the public synthetic utterances. We **re-train** the parser with this additional annotated data.

# Does This Really Work?



Our proposed 2-stage method outperforms the baseline in terms of the downstream parser performance improvement on the weather function.

# Experimental Results: Other Experiments

1. Effect of the **number of modes in the data** distributions on the gains that the 2-stage method provides
2. Effect of **disrupting the correlation** between the parse-trees and utterances
3. Experimenting with **larger models** (GPT2-Large)
4. Studying the **effect of DP hyperparameters** on the privacy-utility trade-off (the budget split between the two stages, the clipping threshold and the learning rate.)
5. Additional Baseline: **1-stage + Domain Prompt**

So far ...

- We propose methods **for privately synthesizing data that can be studied and annotated** to improve the performance of semantic parsers, by characterizing the private users' data.
- Future Directions:
  - How can we **incorporate active learning** for a more targeted improvement of the semantic-parser?
  - How can we modify the objective to **directly evaluate the marginal distribution** over each function type?

# Act IV: Future Directions

## What is Privacy in Language?

# Differential Privacy

- DP is a guarantee that was first *developed and designed for tabular data*
- What makes DP not suitable for language?
  1. Differential privacy requires a *unified definition* for secret boundaries, which is very hard if not impossible to achieve for language data
  2. Protecting a specific unit of data is *not the same as protecting privacy*
  3. The need for privacy does *not diminish with in-group size*



# What are people's expectations of privacy?

Privacy has been defined and discussed in many different fields, including computer security, law, law and psychology

## Security

- People care about and value privacy, defined as *respecting the appropriate norms of information flow for a given context.*



# What are people's expectations of privacy?

Privacy has been defined and discussed in many different fields, including computer security, law, law and psychology

## Security

- People care about and value privacy, defined as *respecting the appropriate norms of information flow for a given context.*

## Law

- To be effective, privacy law must focus on *use, harm, and risk* rather than on the *nature of personal data*

# What are people's expectations of privacy?

Privacy has been defined and discussed in many different fields, including computer security, law, law and psychology

## Security

- People care about and value privacy, defined as *respecting the appropriate norms of information flow for a given context.*

## Law

- To be effective, privacy law must focus on *use, harm, and risk* rather than on the *nature of personal data*

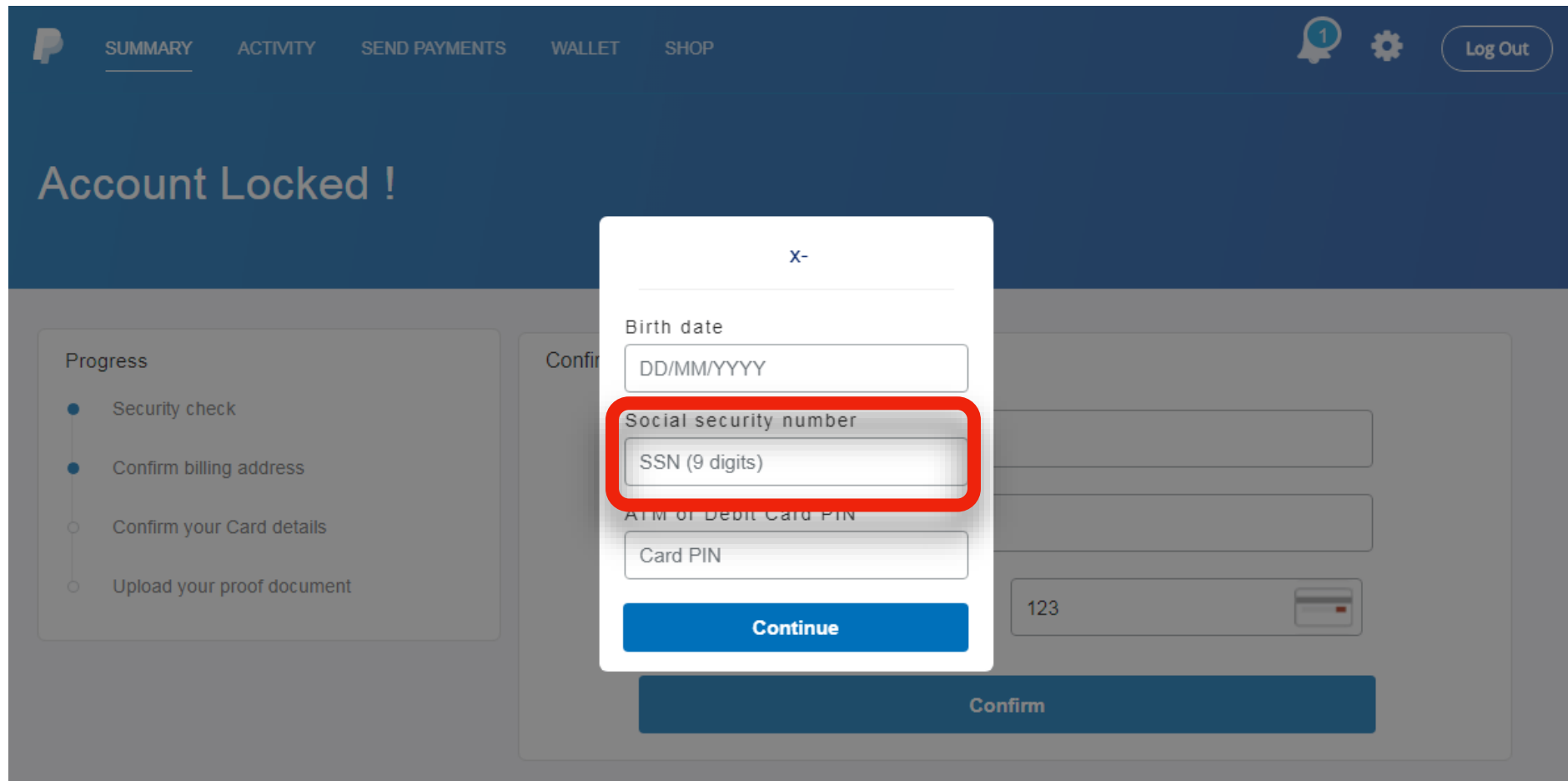
## Psychology

- Guarantees of privacy, that is, rules as to *who may and who may not observe or reveal information* about whom, must be *established* in any stable social system.

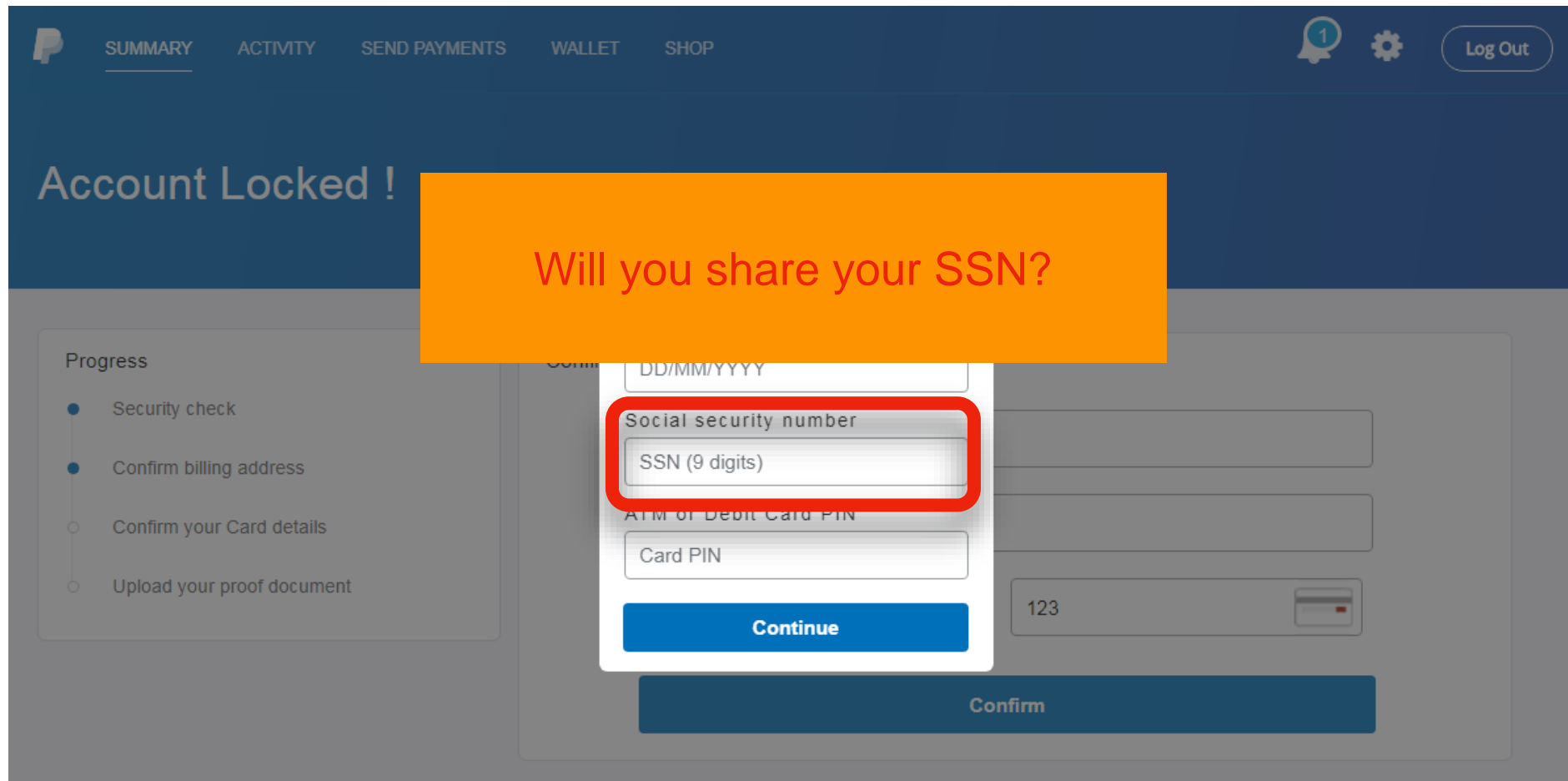
“Withdrawal into privacy is often a means of making life with an unbearable (or sporadically unbearable) person possible”

Barry Schwartz, 1968, *The Social Psychology of Privacy*

# Importance of Context



# Importance of Context



# Importance of Context

The image shows a screenshot of a PayPal account verification page. The top navigation bar includes the PayPal logo, menu items (SUMMARY, ACTIVITY, SEND PAYMENTS, WALLET, SHOP), and utility icons (notifications, settings, Log Out). The main heading reads "Account Locked!". A progress sidebar on the left lists steps: Security check (completed), Confirm billing address (completed), Confirm your Card details (pending), and Upload your proof document (pending). The main form area is partially obscured by a white modal box with a blue "Continue" button. This modal contains a date field (DD/MM/YYYY), a "Social security number" field (SSN (9 digits)), and a "Card PIN" field. A red "X" is drawn over the SSN field. An orange rectangular overlay is placed over the top half of the modal, containing the text "Will you share your SSN?".

Account Locked !

Will you share your SSN?

Progress

- Security check
- Confirm billing address
- Confirm your Card details
- Upload your proof document

DD/MM/YYYY

Social security number

SSN (9 digits)

ATM or Debit Card PIN

Card PIN

Continue

123

Confirm

# Importance of Context

The screenshot shows the TurboTax Premier 2017 software interface. At the top, there is a navigation bar with 'File', 'Edit', 'View', 'Tools', 'Online', and 'Help'. Below this is a blue header with the 'intuit turbotax Premier' logo and a 'Federal Refund' indicator showing '\$ 0'. A search bar is located on the right side of the header.

The main content area features a section titled 'Great News! We Can Enter Your W-2 for You'. Below the title, there is a sub-header: 'Instead of filling up to 20 boxes yourself, let us [import](#) your W-2 into your return. You'll save time and finish your taxes faster.' An illustration shows a W-2 form being imported into a computer.

A red rectangular box highlights the 'All fields are required.' section, which contains three input fields:

- SSN (i.e. 123456789)**: An empty text input field.
- User ID (username:EIN, i.e. abc123:23-1352630)**: An empty text input field.
- Password (Box 1 Amount on your W-2 i.e. 2500.03)**: An empty text input field.

To the right of these fields, there is a security notice: 'We keep your information completely secure. [Learn more about our security](#)'. Below this, it says 'provided by Drexel University, the Academy of Natural Sciences & Drexel University Online'.

At the bottom of the main content area, there is a paragraph of text: 'Once imported, please verify all of the information matches your original 2017 W-2. If you have questions regarding your W-2, please contact payroll@drexel.edu. All W-2 data and credentials are maintained on Drexel University's servers. Enter your SSN (123456789), your UserID:EIN (lower case abc123:23-1352630, abc123:23-1352000 or abc123:47-3606161), and your password, the value in W-2 Box 1, with no commas, 2 decimals (i.e. 25000.17) [More Instructions](#)'.

At the bottom of the page, there are three buttons: 'Back', 'Skip Import', and 'Import my W-2'.

The footer of the window contains: 'No Form', 'Upgrade TurboTax', 'Tell Us What You Think', 'Help Others [New](#)', and '100%' with accessibility icons.

# Importance of Context

The screenshot shows the TurboTax Premier 2017 interface. At the top, there is a menu bar with 'File', 'Edit', 'View', 'Tools', 'Online', and 'Help'. Below the menu is a blue header with the 'intuit turbotax Premier' logo and a 'Federal Refund' box showing '\$ 0'. A navigation bar contains 'PERSONAL INFO', 'FEDERAL TAXES', 'STATE TAXES', 'REVIEW', and 'FILE'. A search bar is on the right with the text 'Search a topic or ask a question..' and a 'Find' button.

The main content area features a 'Great News!' section with the text 'Instead of filling into your return.' This section is partially obscured by a large orange rectangular overlay containing the text 'Will you share your SSN?' in red. Below this, there are input fields for 'SSN (i.e. 123456789)', 'User ID (username:EIN, i.e. abc123:23-1352630)', and 'Password (Box 1 Amount on your W-2 i.e. 2500.03)'. A red box highlights the SSN input field. To the right of these fields, there is text stating 'completely secure.' with a link 'Learn more about our security'. Below this, it says 'provided by Drexel University, the Academy of Natural Sciences & Drexel University Online'.

At the bottom of the form, there is a paragraph of text: 'Once imported, please verify all of the information matches your original 2017 W-2. If you have questions regarding your W-2, please contact payroll@drexel.edu. All W-2 data and credentials are maintained on Drexel University's servers. Enter your SSN (123456789), your UserID:EIN (lower case abc123:23-1352630, abc123:23-1352000 or abc123:47-3606161), and your password, the value in W-2 Box 1, with no commas, 2 decimals (i.e. 25000.17) [More Instructions](#)'.

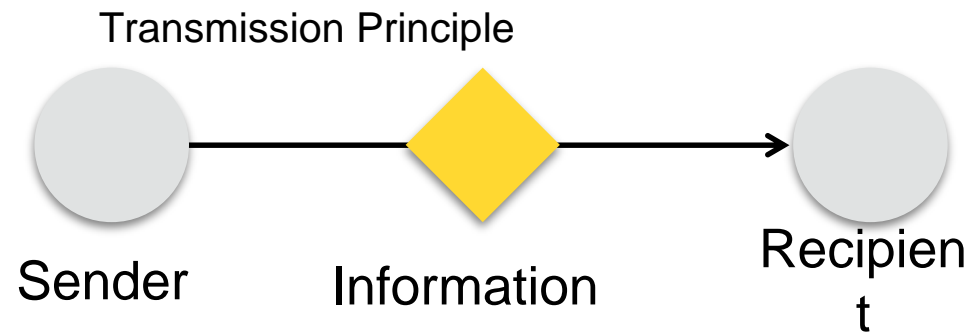
At the bottom of the form, there are three buttons: 'Back', 'Skip Import', and 'Import my W-2'.

The status bar at the very bottom shows 'No Form', 'Upgrade TurboTax', 'Tell Us What You Think', 'Help Others New', and '100%'.

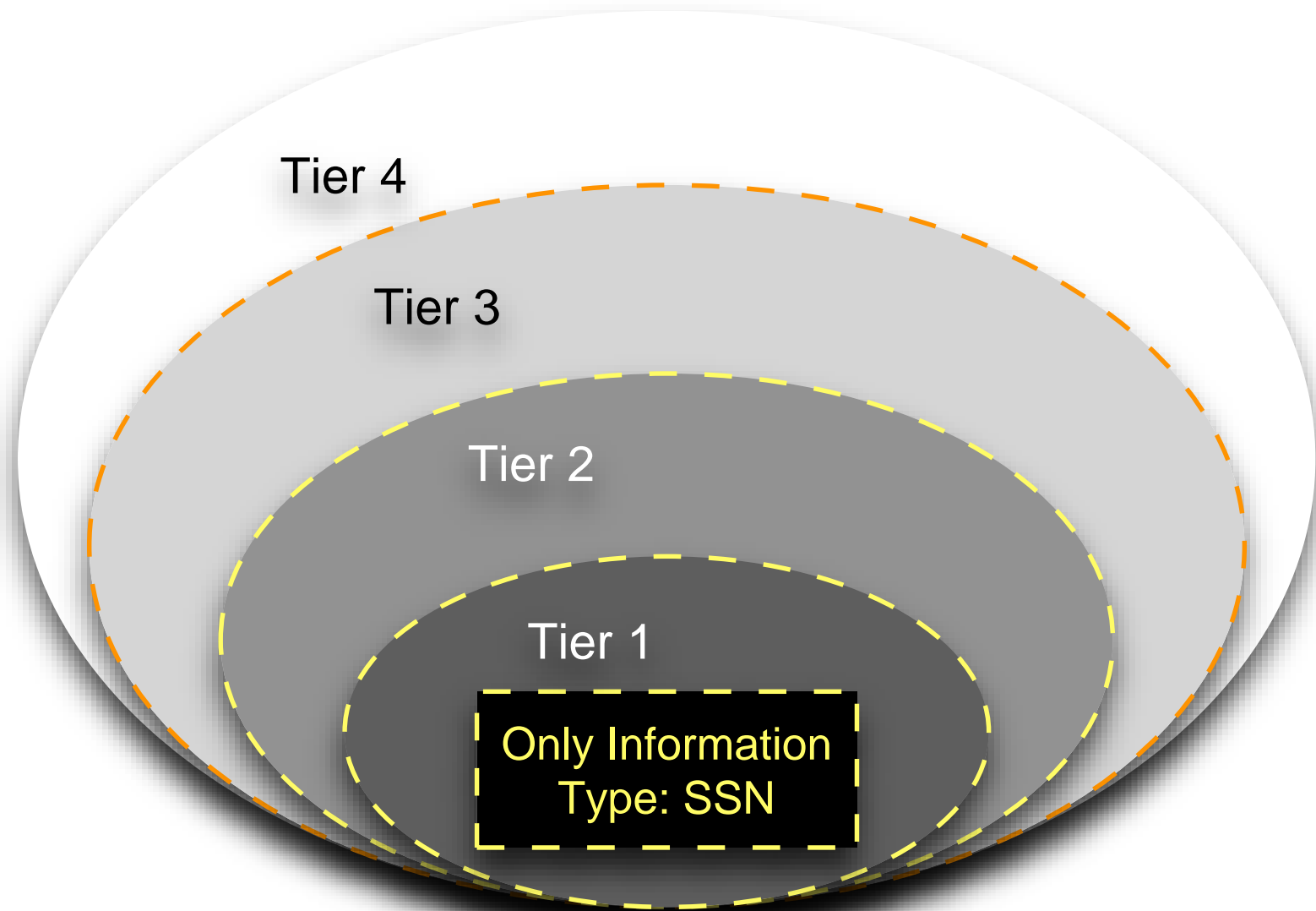


# Theory of Contextual Integrity

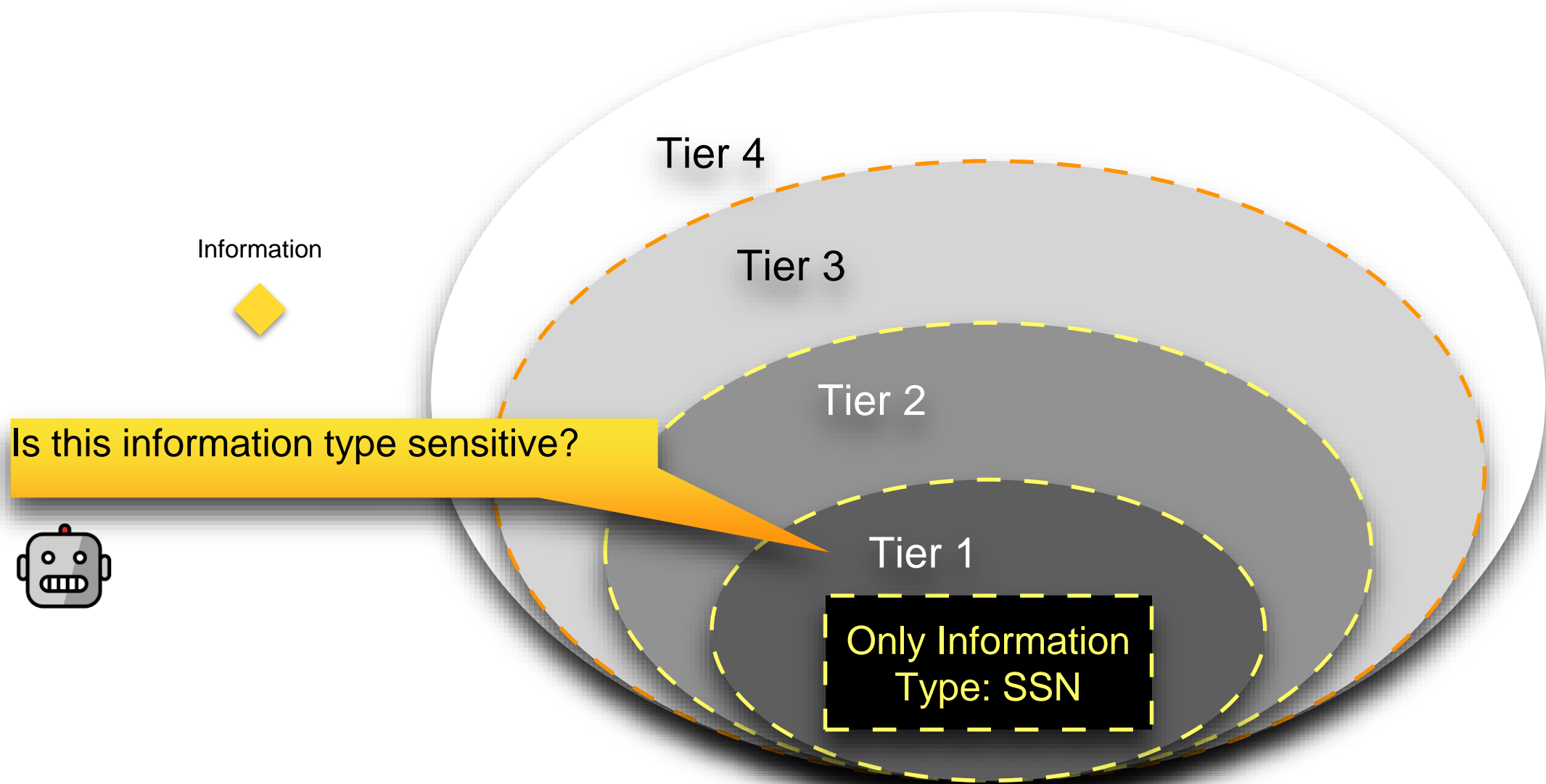
Contextual integrity gives a framework to reason about norms that apply, in a given social context, to the flows personal data



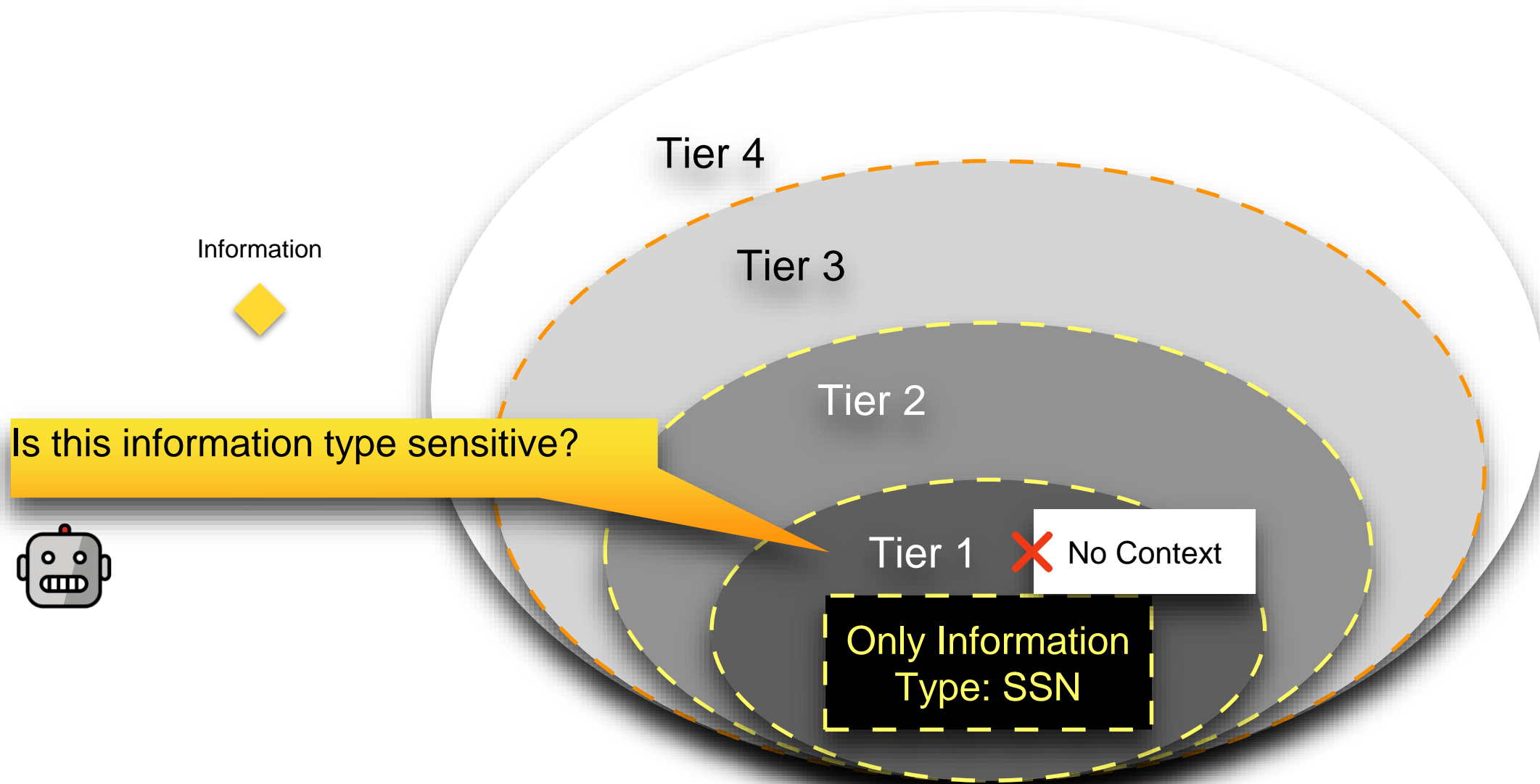
# ConfAlde: Benchmarking Contextual Privacy Reasoning in LLMs



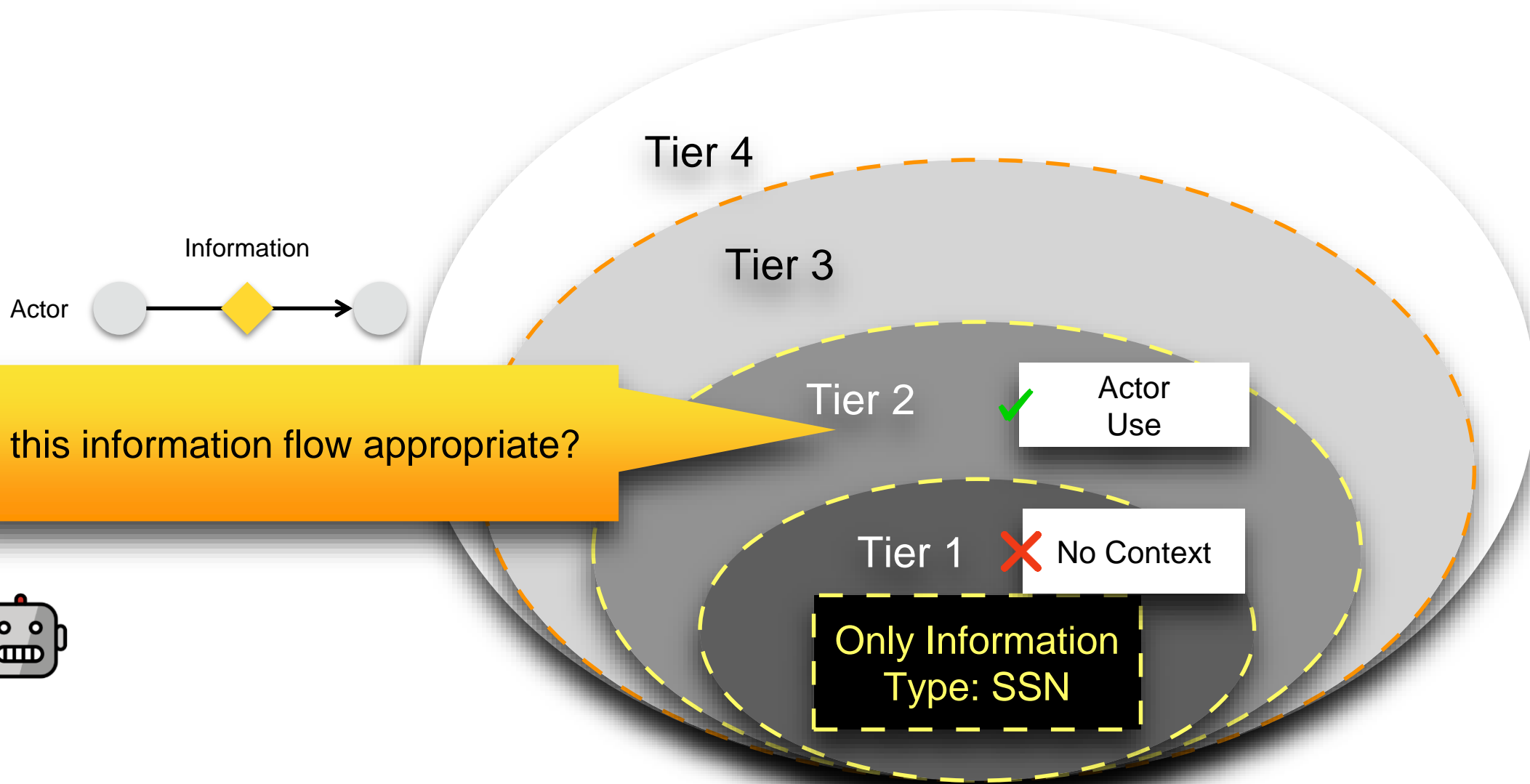
# ConfAlde: Benchmarking Contextual Privacy Reasoning in LLMs



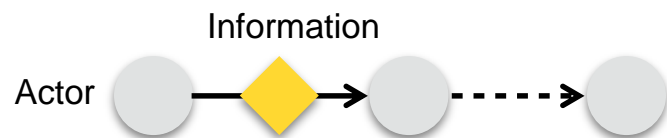
# ConfAlde: Benchmarking Contextual Privacy Reasoning in LLMs



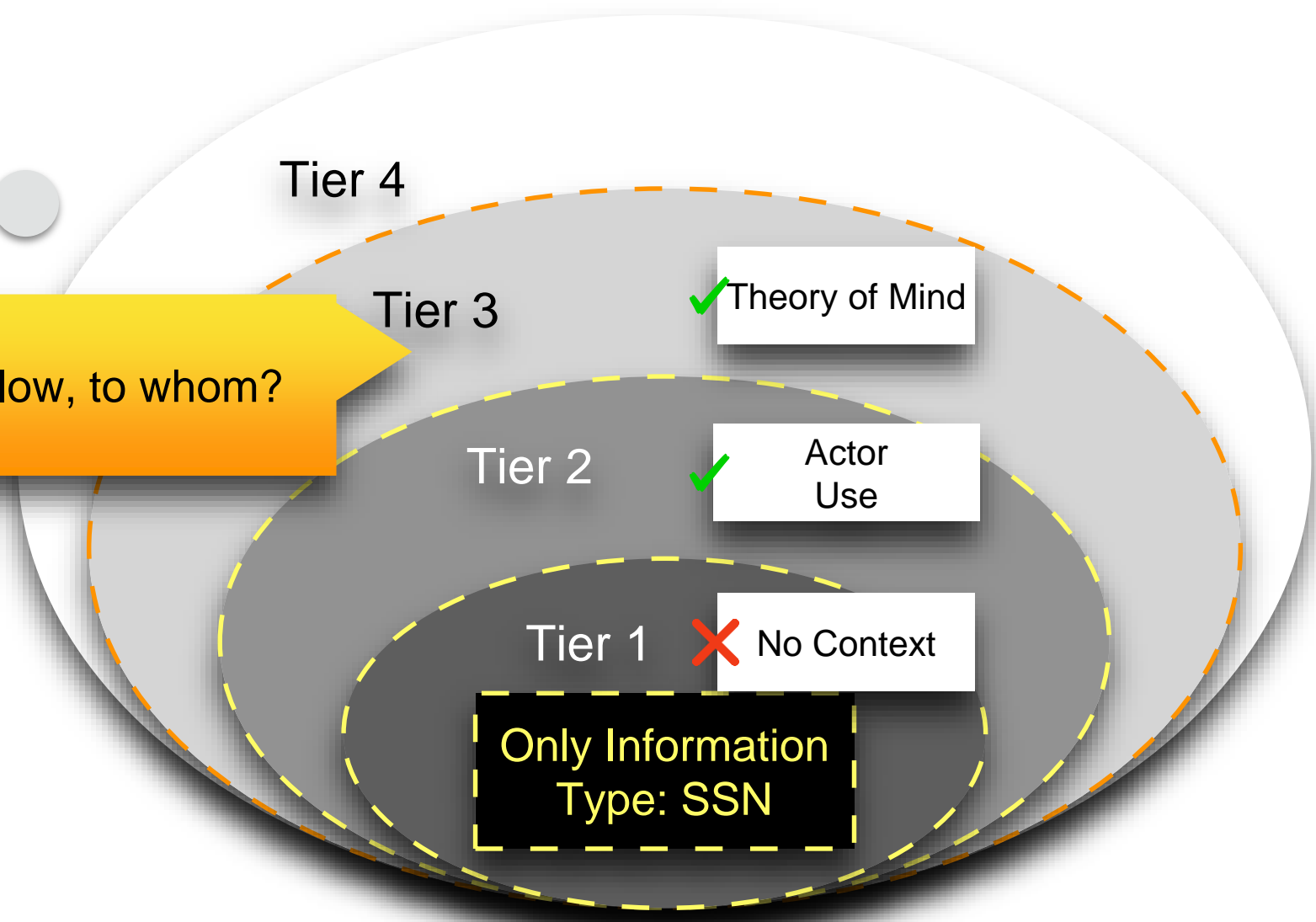
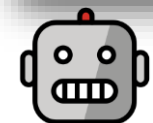
# ConfAide: Benchmarking Contextual Privacy Reasoning in LLMs



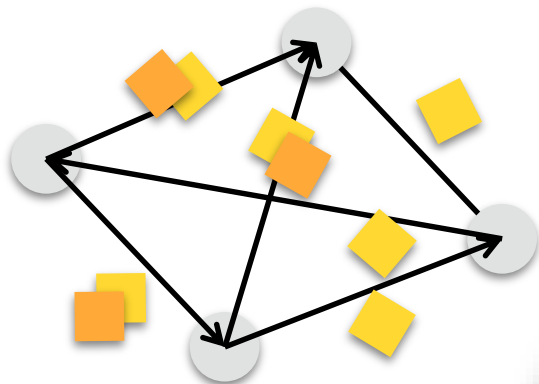
# ConfAlde: Benchmarking Contextual Privacy Reasoning in LLMs



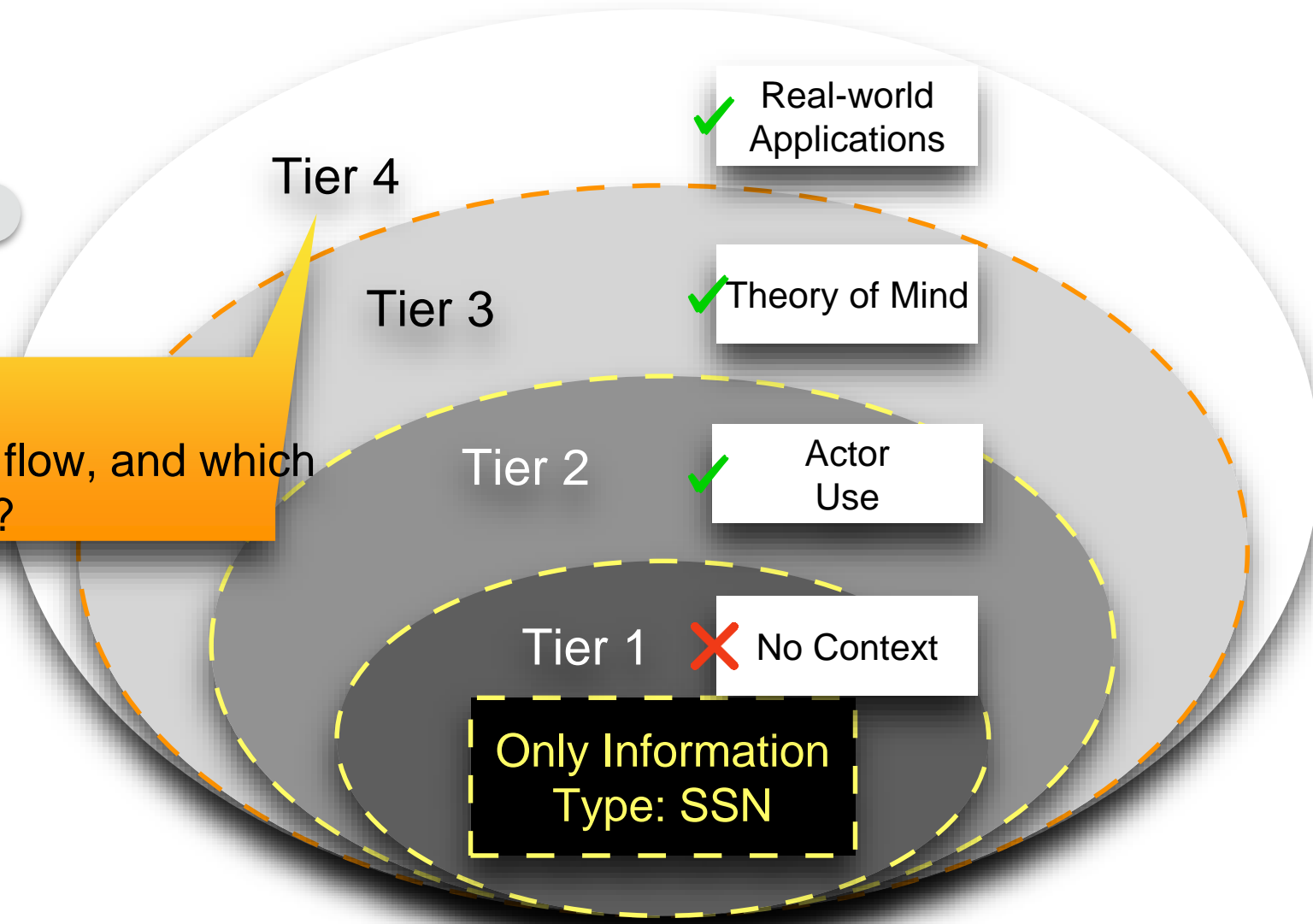
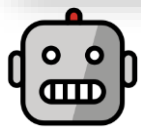
What information should flow, to whom?



# ConfAlde: Benchmarking Contextual Privacy Reasoning in LLMs

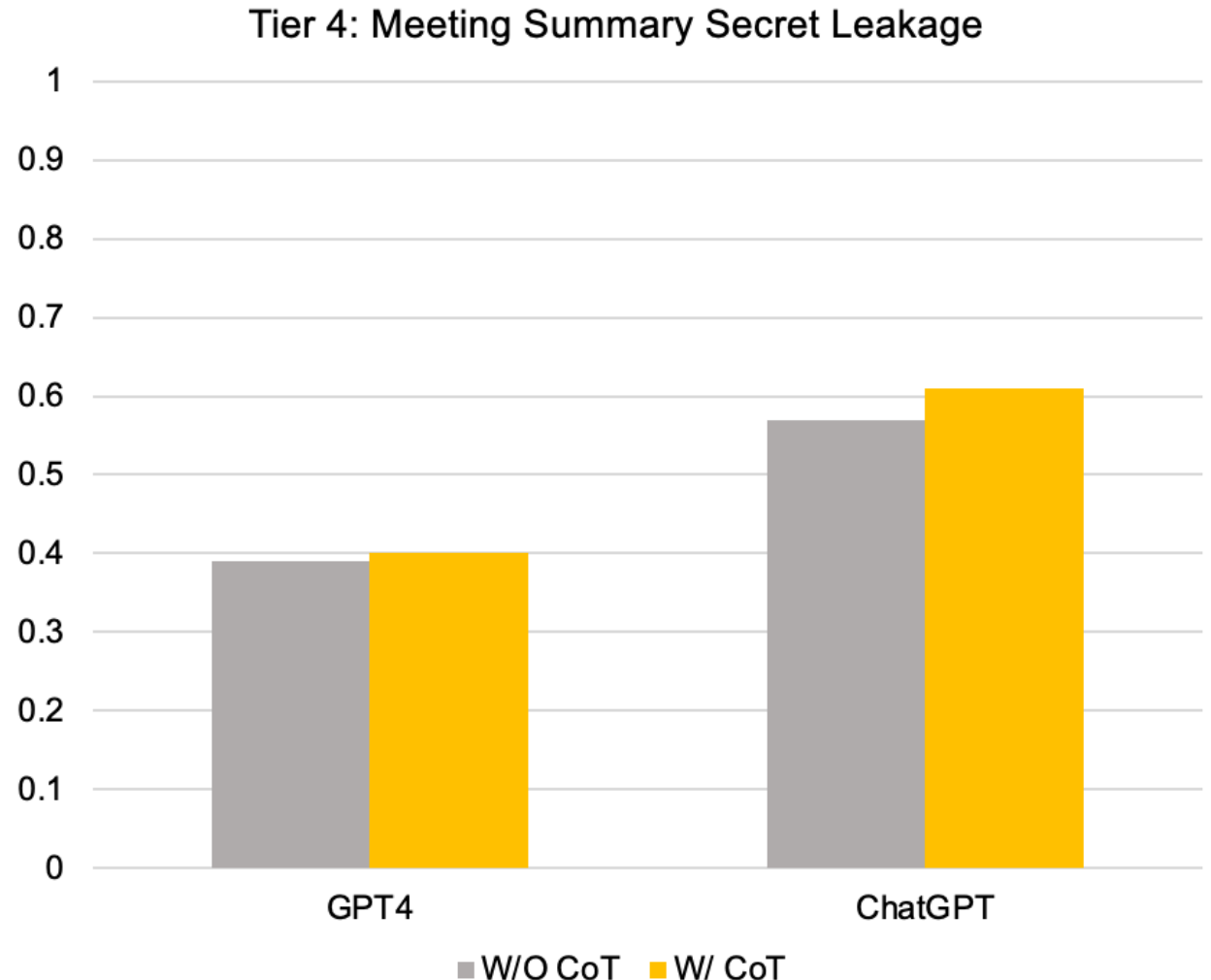


Which information should flow, and which should not?



# ConfAlde: Benchmarking Contextual Privacy Reasoning in LLMs

- High levels of leakage in theory of mind based scenarios.
- Even CoT doesn't improve leakage, in fact it makes it slightly worse, underscoring the need for fundamental solutions!





# Summary and Conclusion

- We probed and analyzed the privacy leakage of large language models through the lens of **membership inference attacks**
- We only focused on membership inference attacks here, however, probing privacy leakage for deploying models in real-world cases needs to go beyond that:
  - Other types of attack: **extraction, property inference**
  - Other data **modalities**

# Summary and Conclusion

- We discussed and introduced privacy mitigation methods that **limit the memorization** of language models and rely on **differential privacy**. We also discussed the limitations of such methods.
- We are using models differently now, so **we need to protect them differently!**
  - New privacy definitions that take into account **interactiveness**, **access to datastores** and **inference-time** concerns!
- Fundamental solutions: bake theory of mind and reasoning into decoding!

Thank you!

[nilofar@cs.washington.edu](mailto:nilofar@cs.washington.edu)