

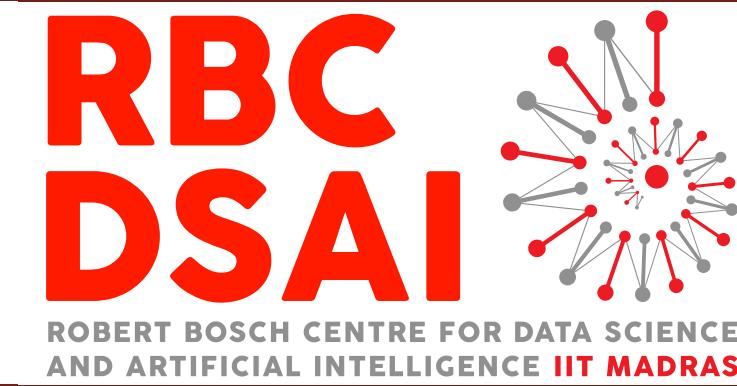
Mathematical Foundations of Differential Privacy

Krishna Pillutla

IIT Madras



CeRAI
Centre for Responsible AI



ROBERT BOSCH CENTRE FOR DATA SCIENCE
AND ARTIFICIAL INTELLIGENCE IIT MADRAS

ACM INDIA SUMMER SCHOOL

RESPONSIBLE & SAFE AI

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB



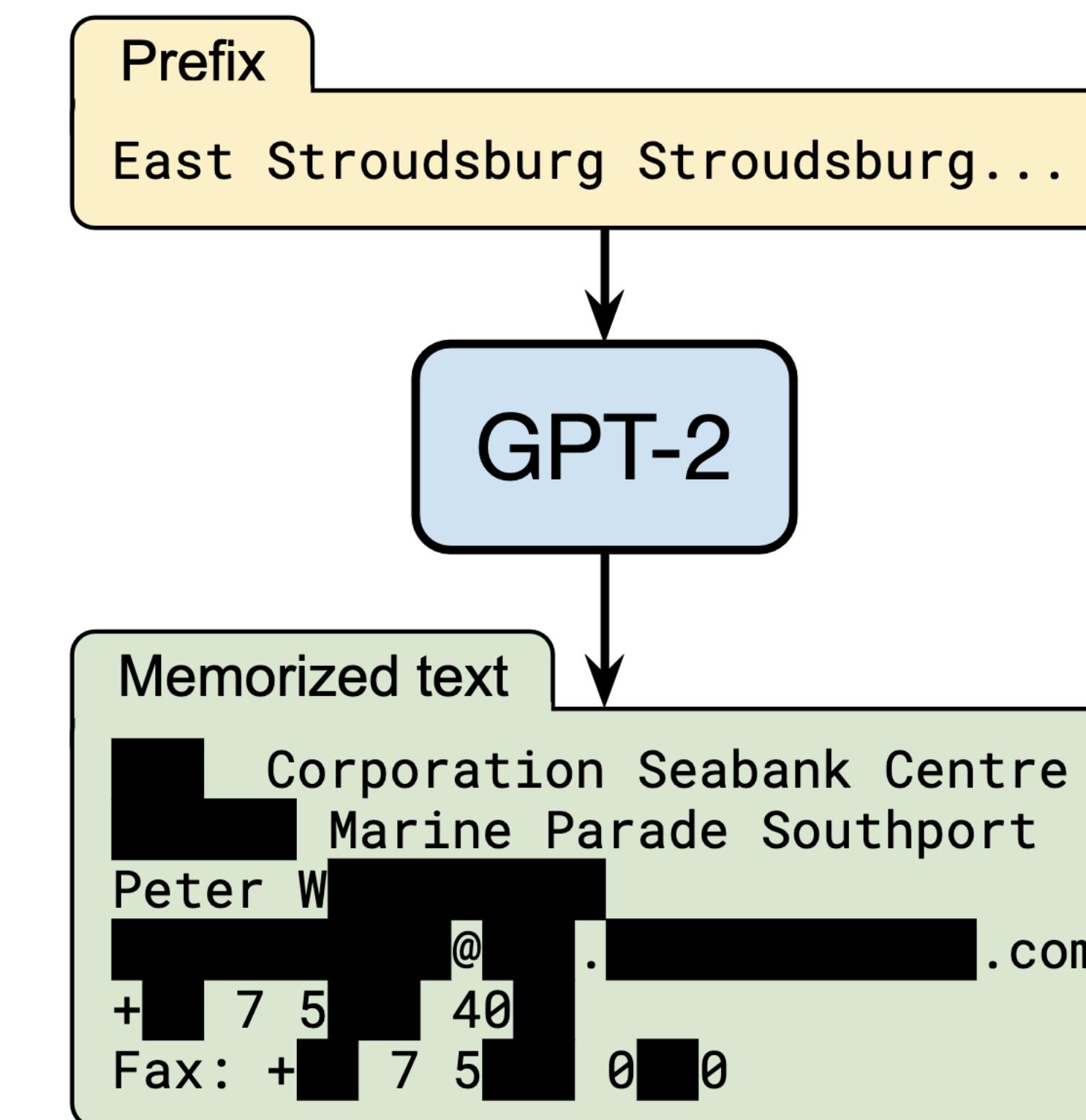
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB



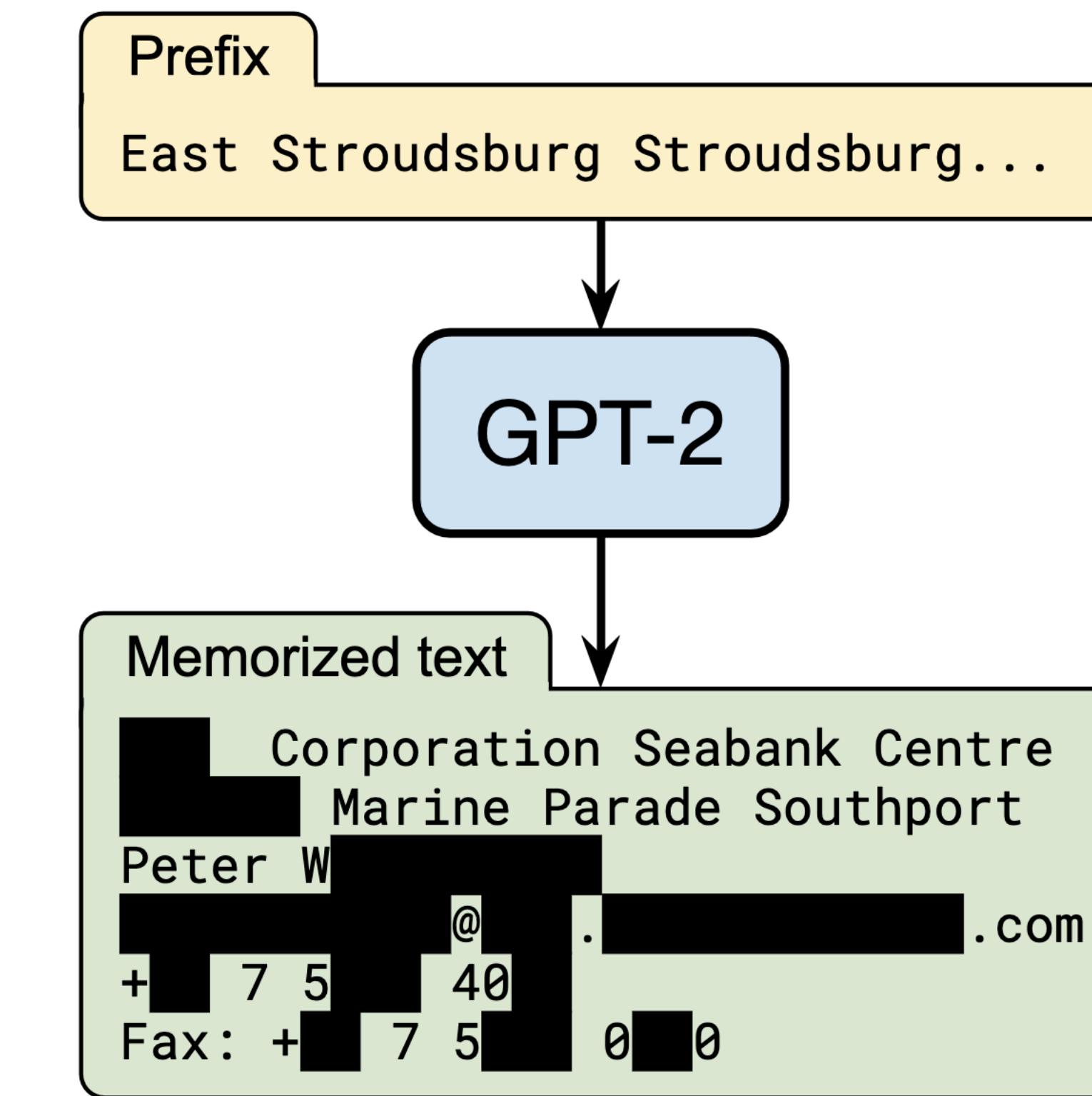
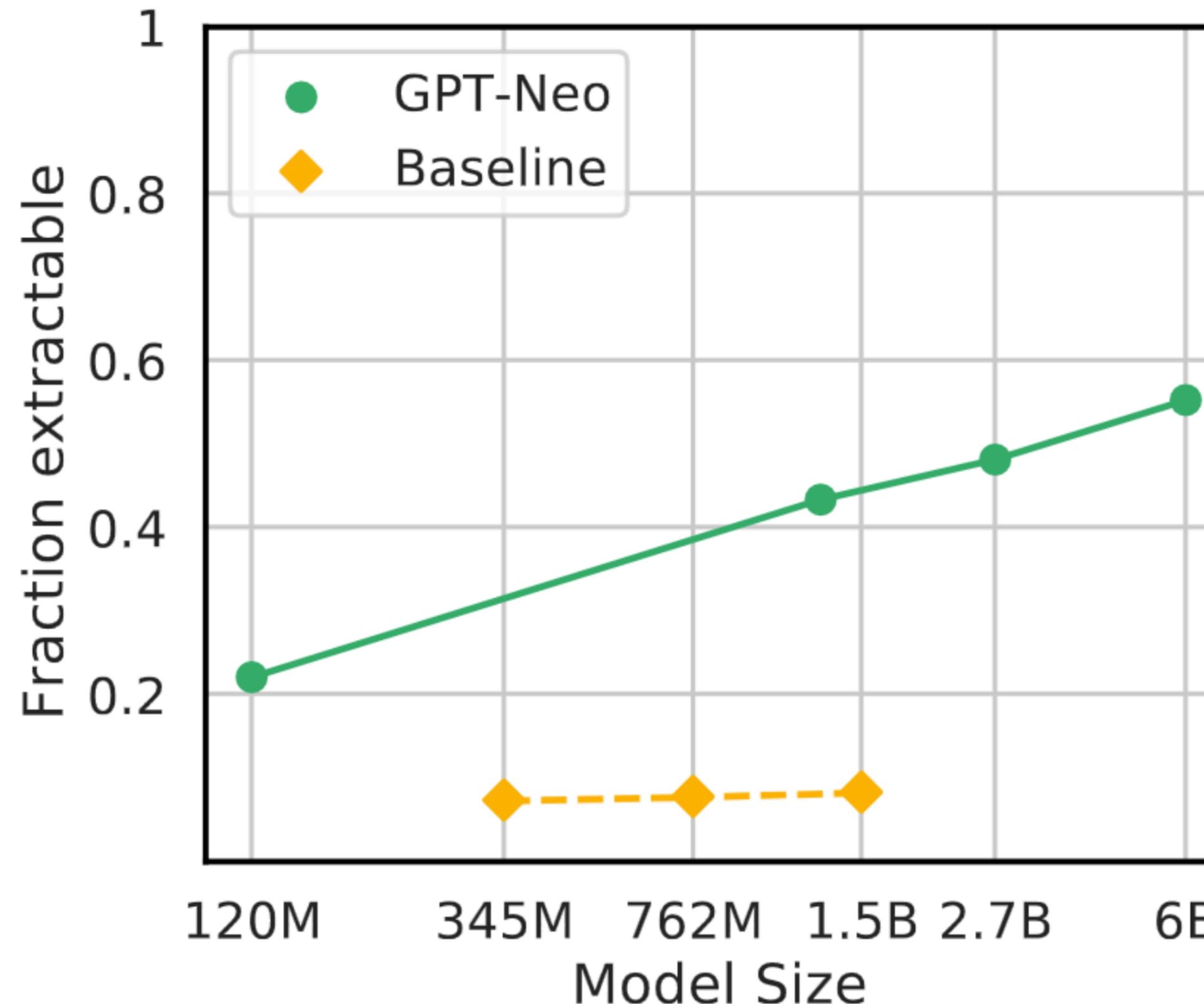
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

Models leak information about their training data



Carlini et al. (USENIX Security 2021)

Models leak information about their training data *reliably*



Carlini et al. (USENIX Security 2021)

Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models

Gowthami Somepalli , Vasu Singla , Micah Goldblum , Jonas Geiping , Tom Goldstein 



University of Maryland, College Park

{gowthami, vsingla, jgeiping, tomg}@cs.umd.edu



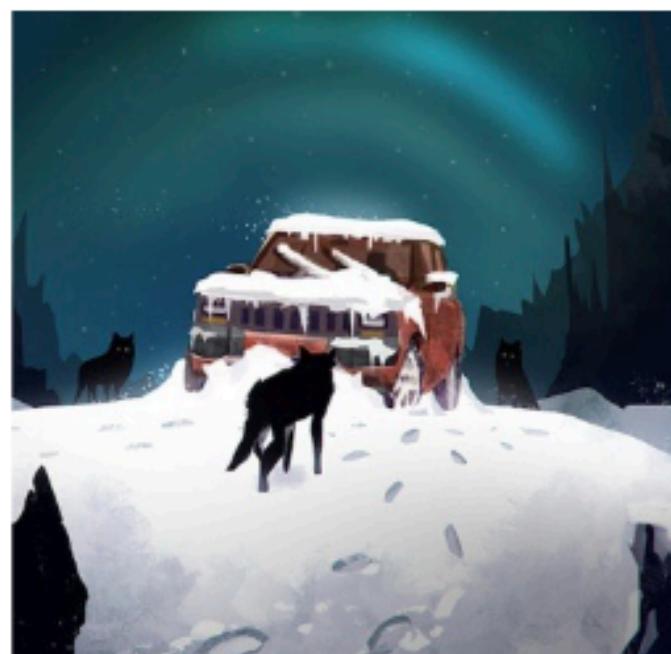
New York University

goldblum@nyu.edu

Generation

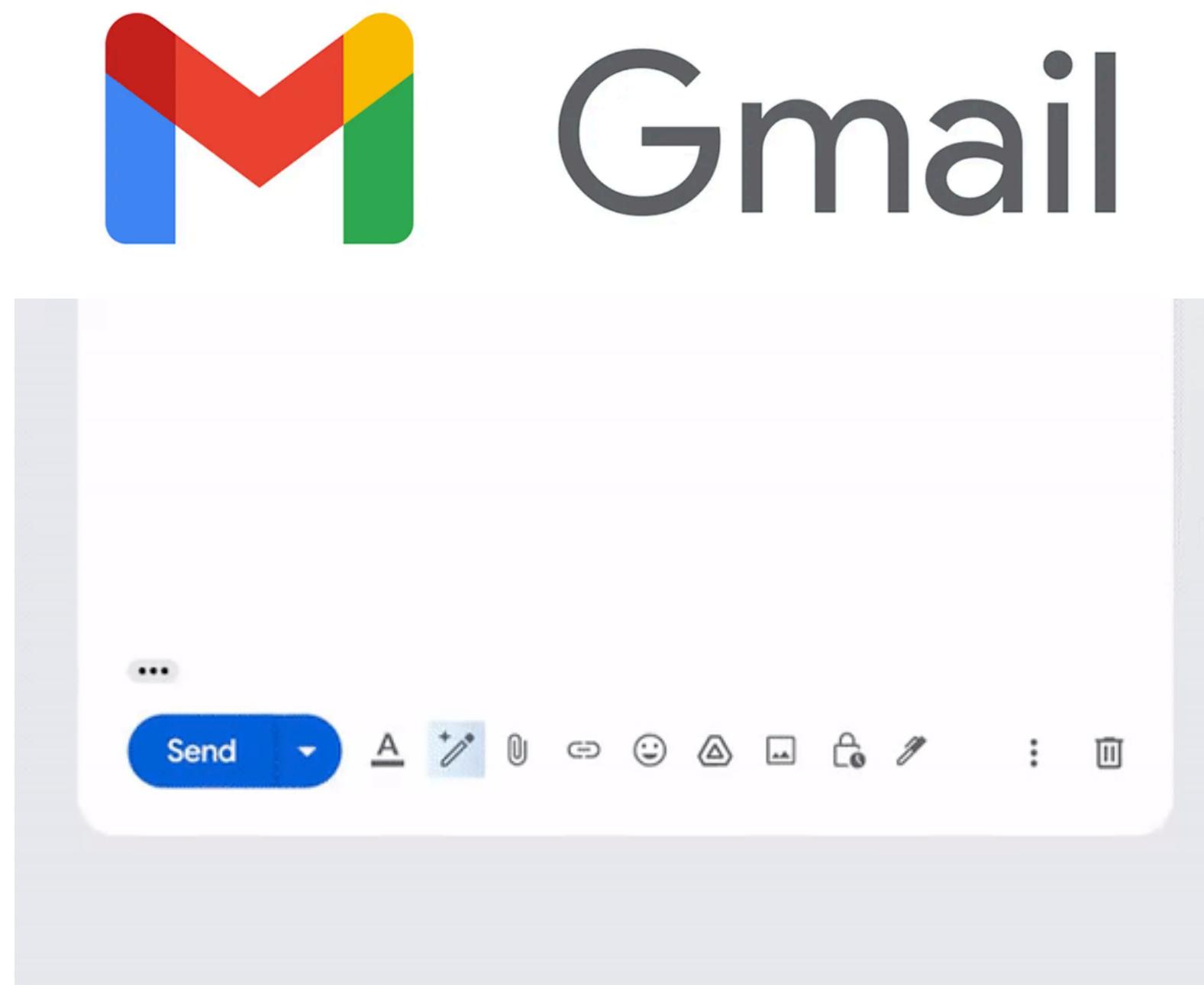
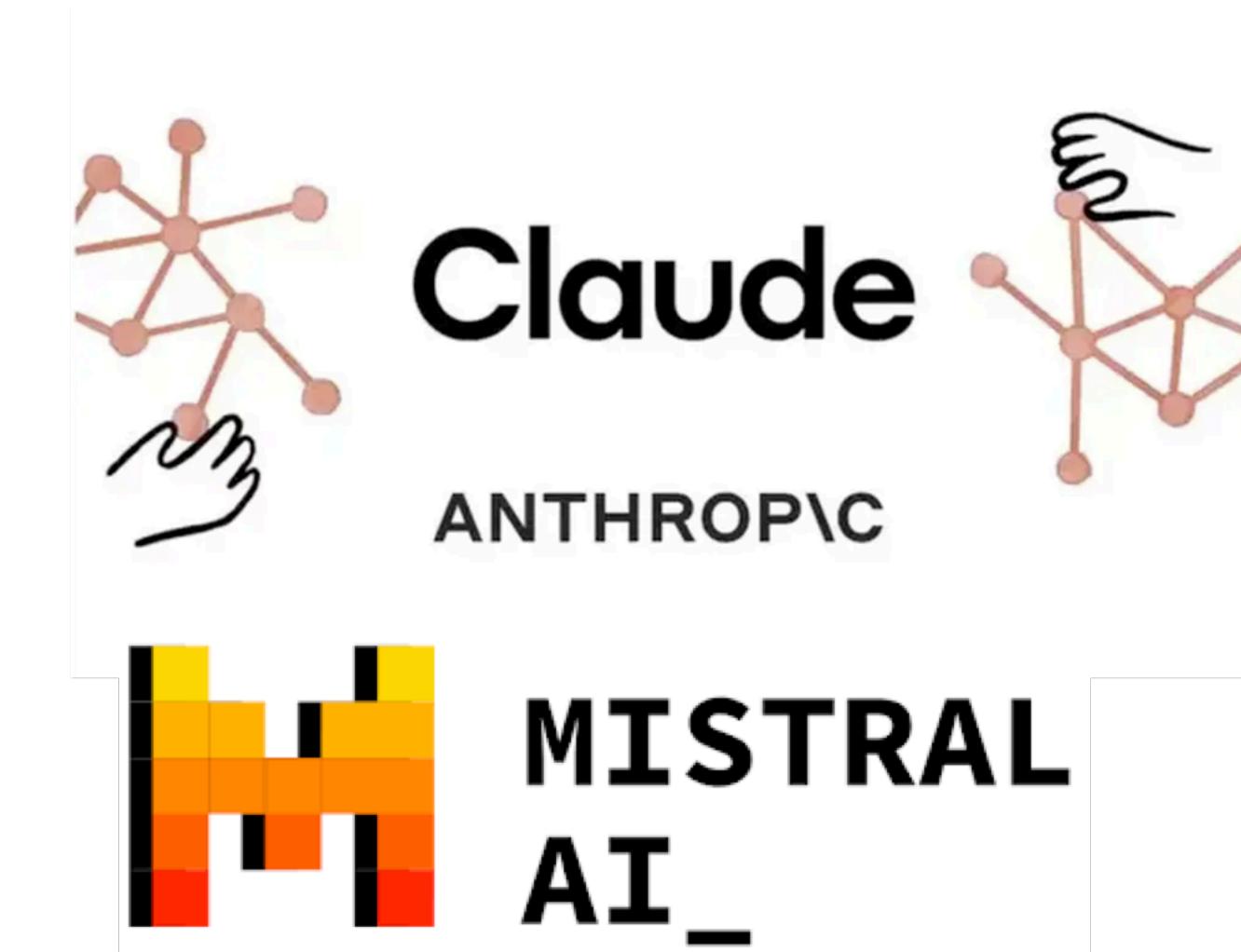


LAION-A Match



Blog

Introducing ChatGPT



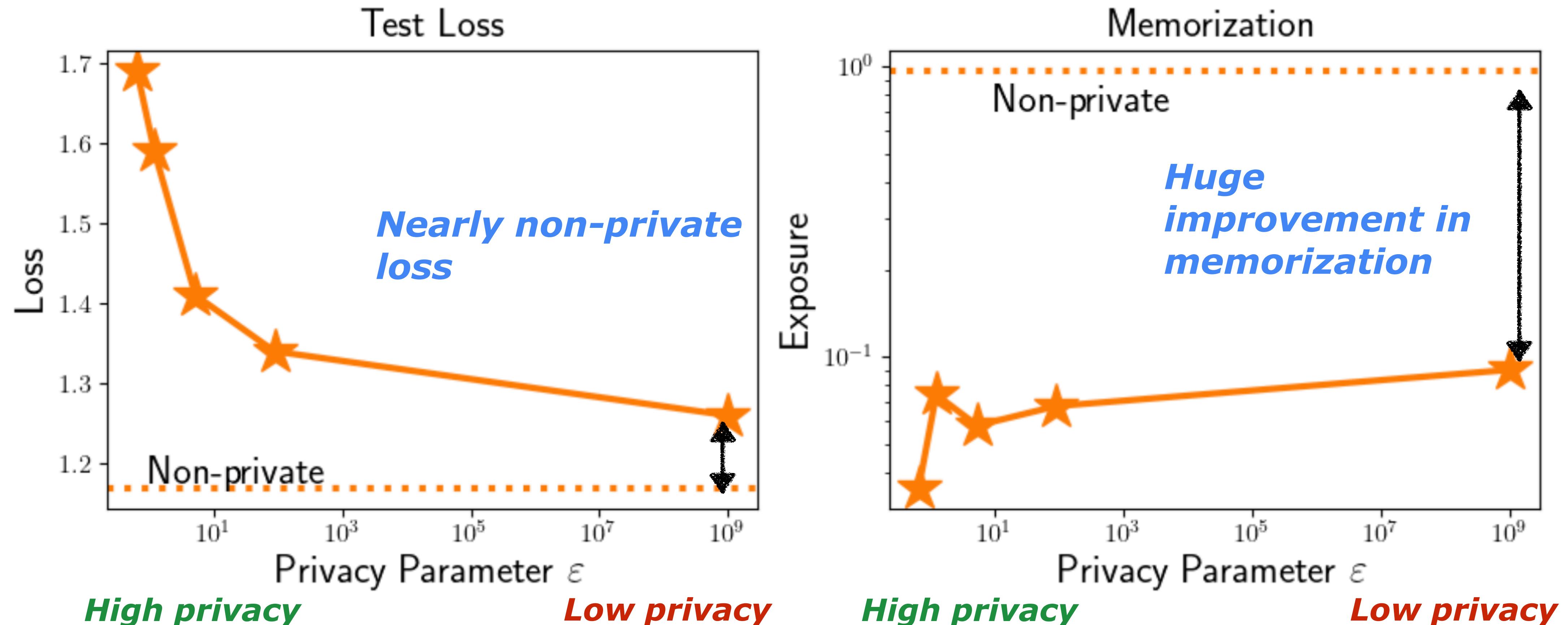
Microsoft Copilot for Microsoft 365
Your everyday AI [assistant at work](#)

Natural Language

- Large Language Models (LLMs)
- Web grounding
- Microsoft Graph grounding
- Microsoft 365 Apps
- Enterprise-grade data protection
- Copilot Studio

<https://blog.google/products/gmail/gmail-ai-features/>

Differential privacy nearly eliminates memorization



Carlini, Liu, Erlingsson, Kos, Song. **The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks**. USENIX Security 2019.

Caveat: Multiple facets of the word “*privacy*”

What does the word “*privacy*” mean to an end user of an AI product?



Transparency, Control,
Verifiability



Minimize data sharing



Data Anonymization

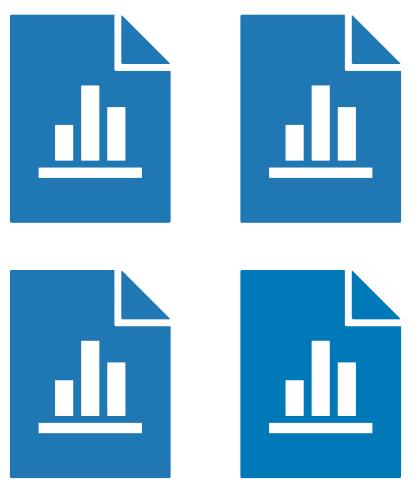
- *Differential Privacy*

Outline

- ***Differential Privacy: Intuition & Recap***
- Differential Privacy: Rigorous Mathematical Formulation & Properties
- Application Examples

Differential Privacy:
A mathematically rigorous notion of “*privacy*”

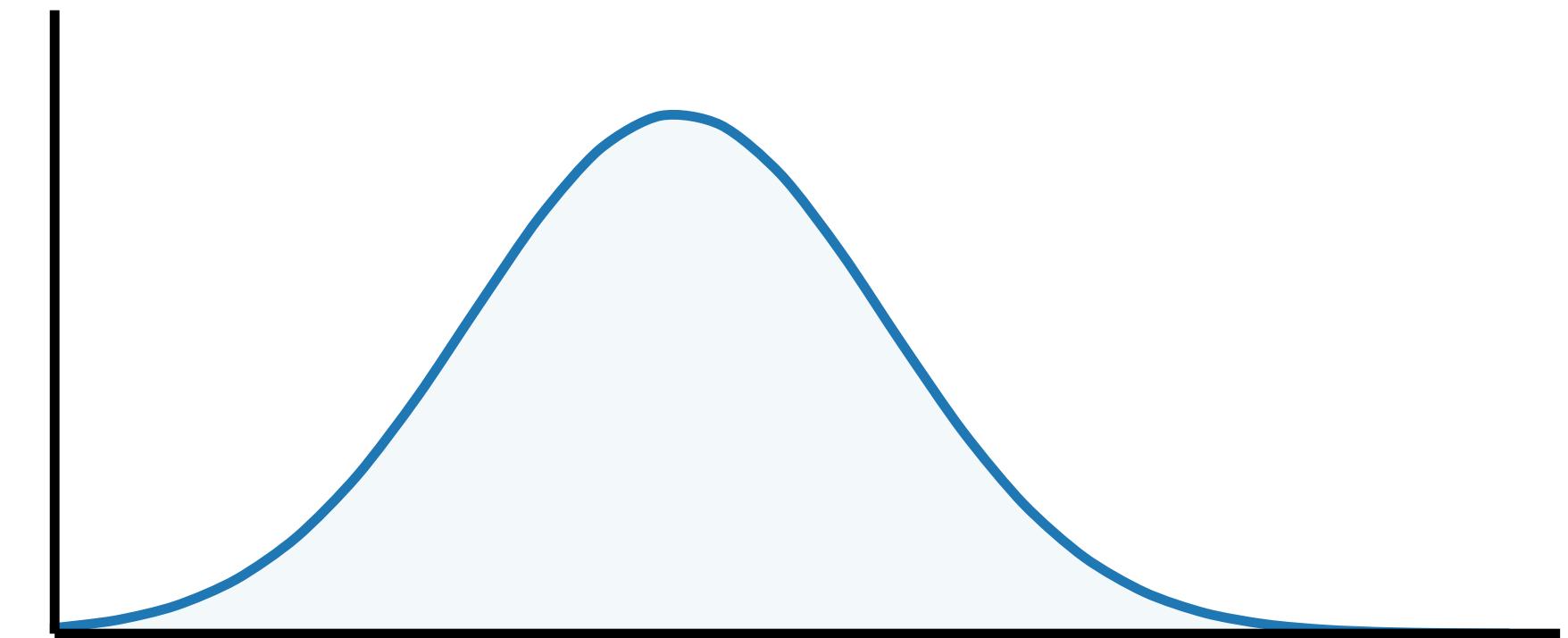
Dataset



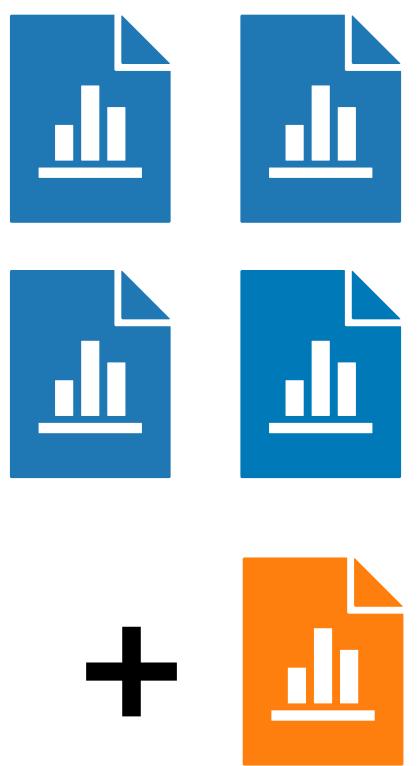
Randomized
Algorithm



Output Distribution
(e.g. over models)



Dataset



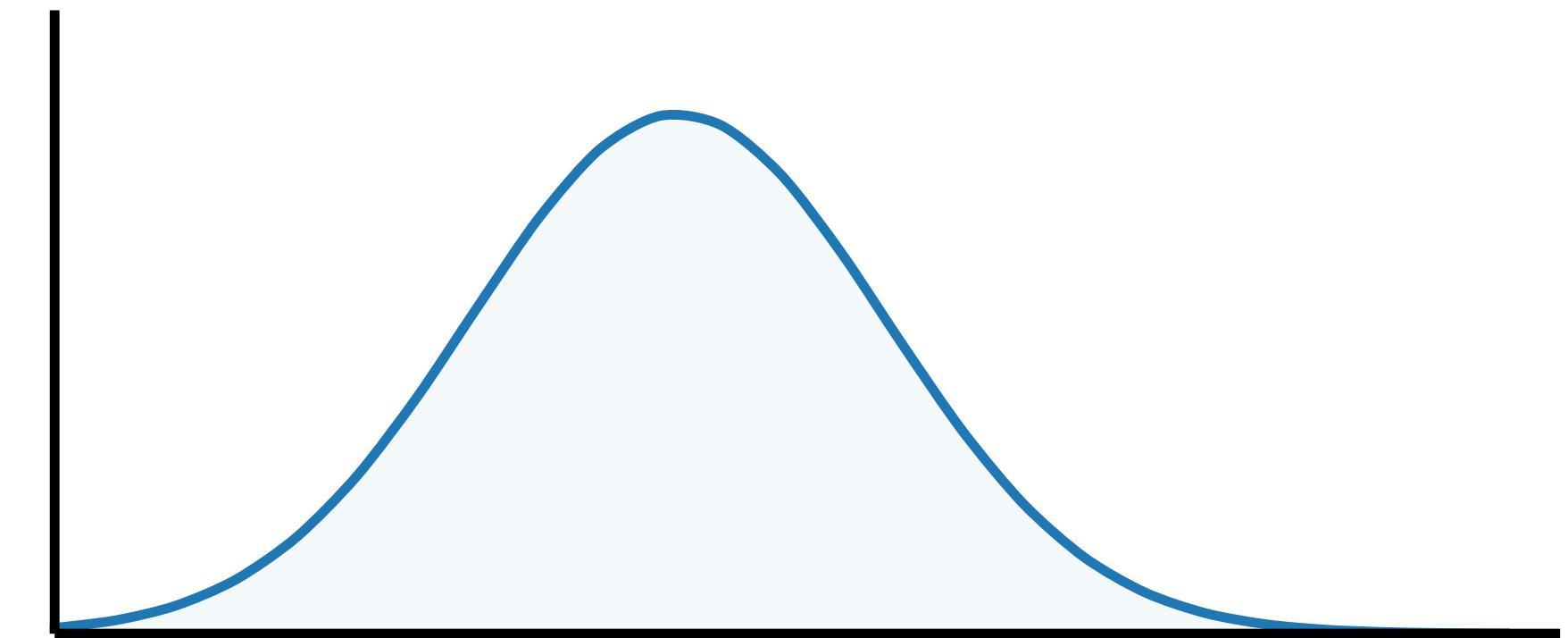
+



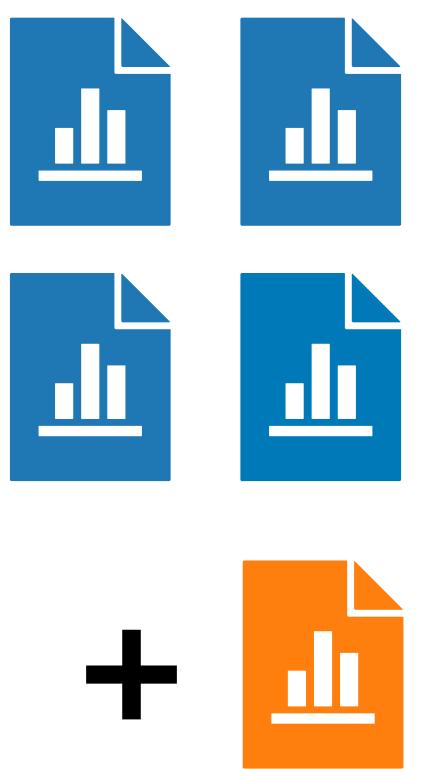
Randomized
Algorithm



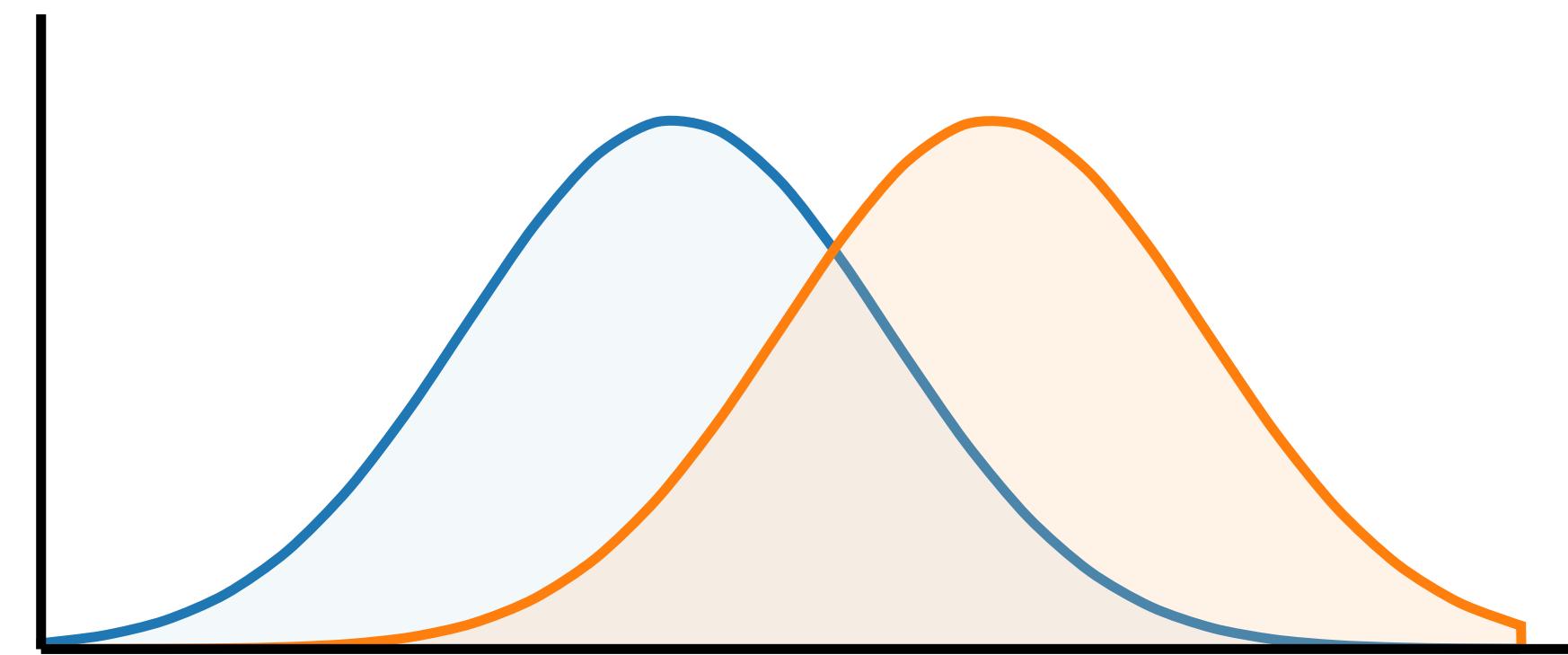
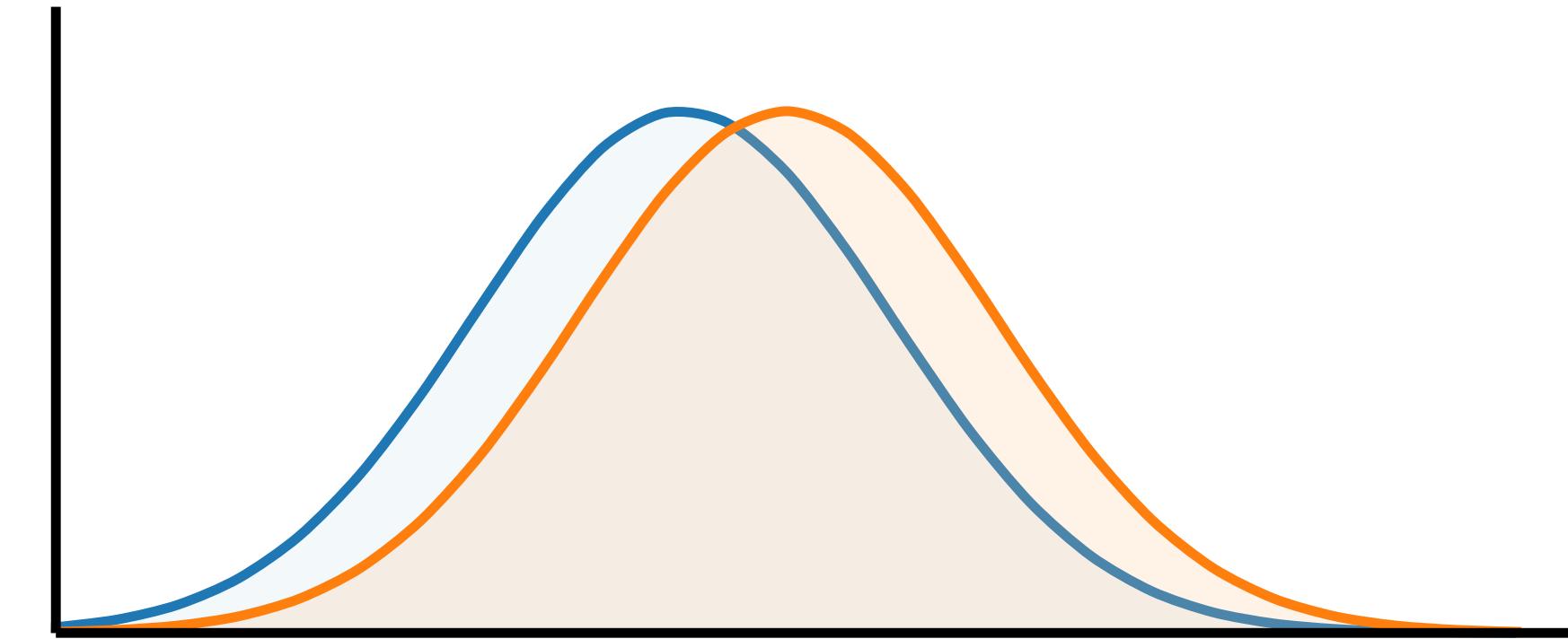
Output Distribution
(e.g. over models)



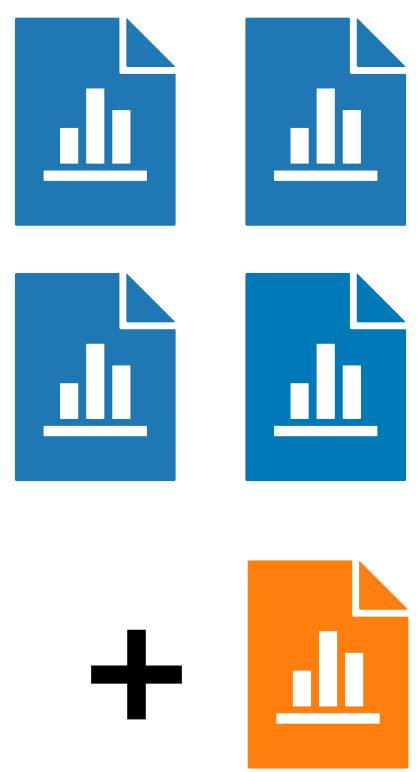
Dataset



Output Distribution
(e.g. over models)



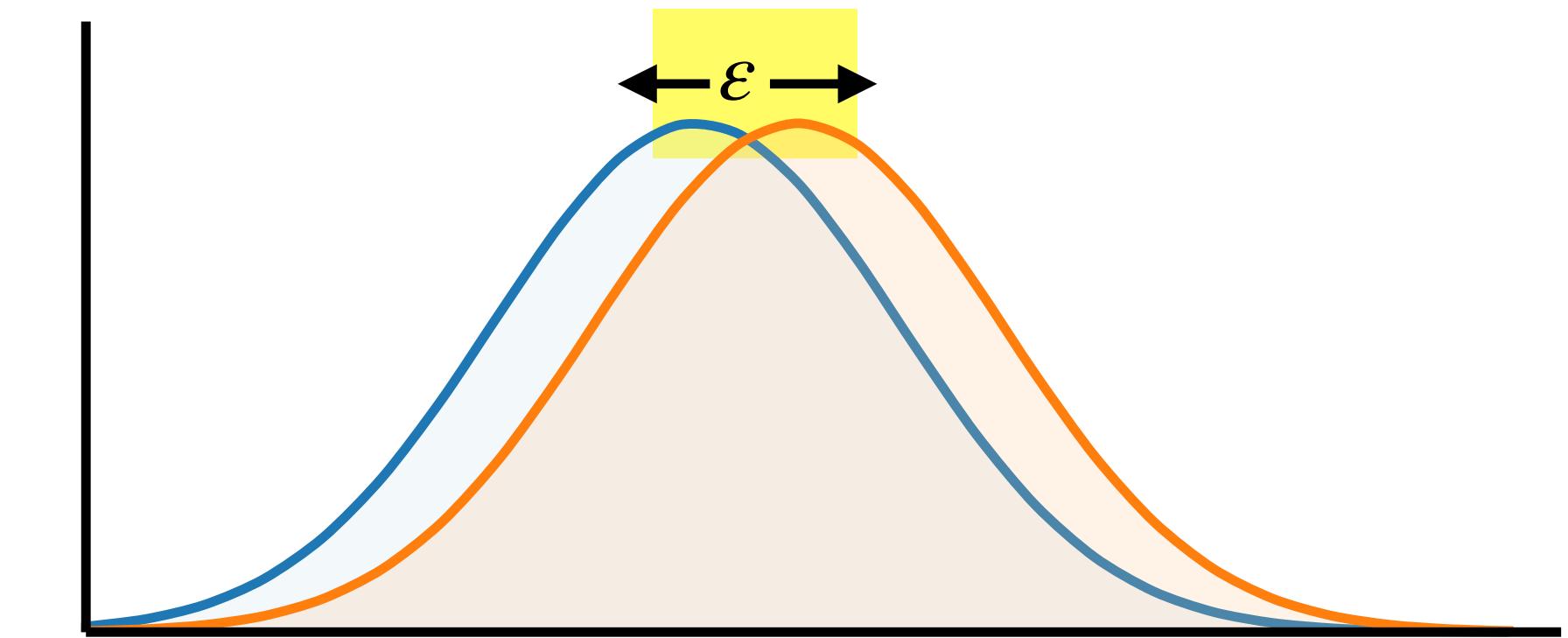
Dataset



Randomized
Algorithm

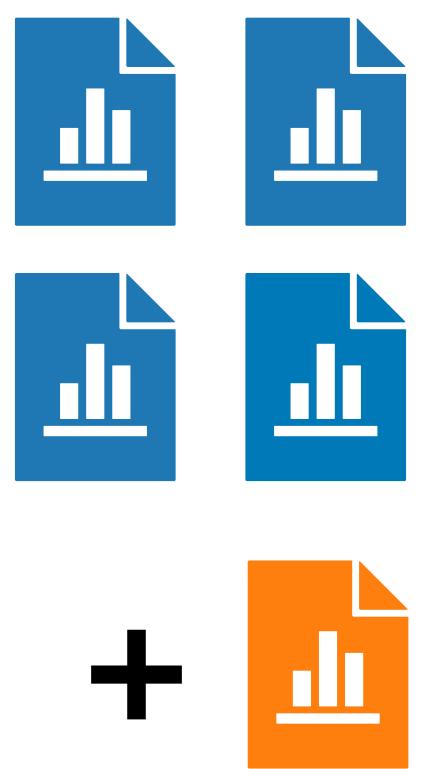


Output Distribution
(e.g. over models)



A randomized algorithm is **ε -differentially private** if the addition of **one user's data** does not alter its output distribution by more than ε

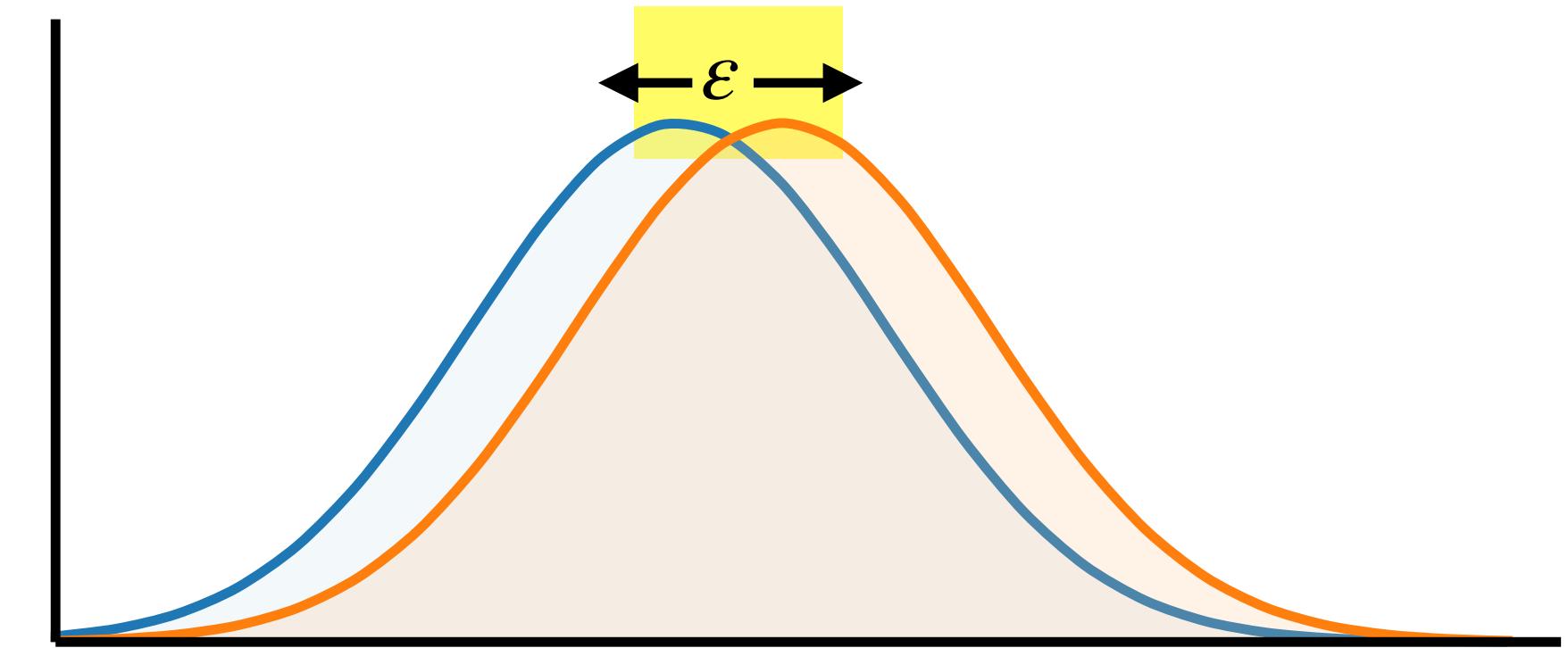
Dataset



+

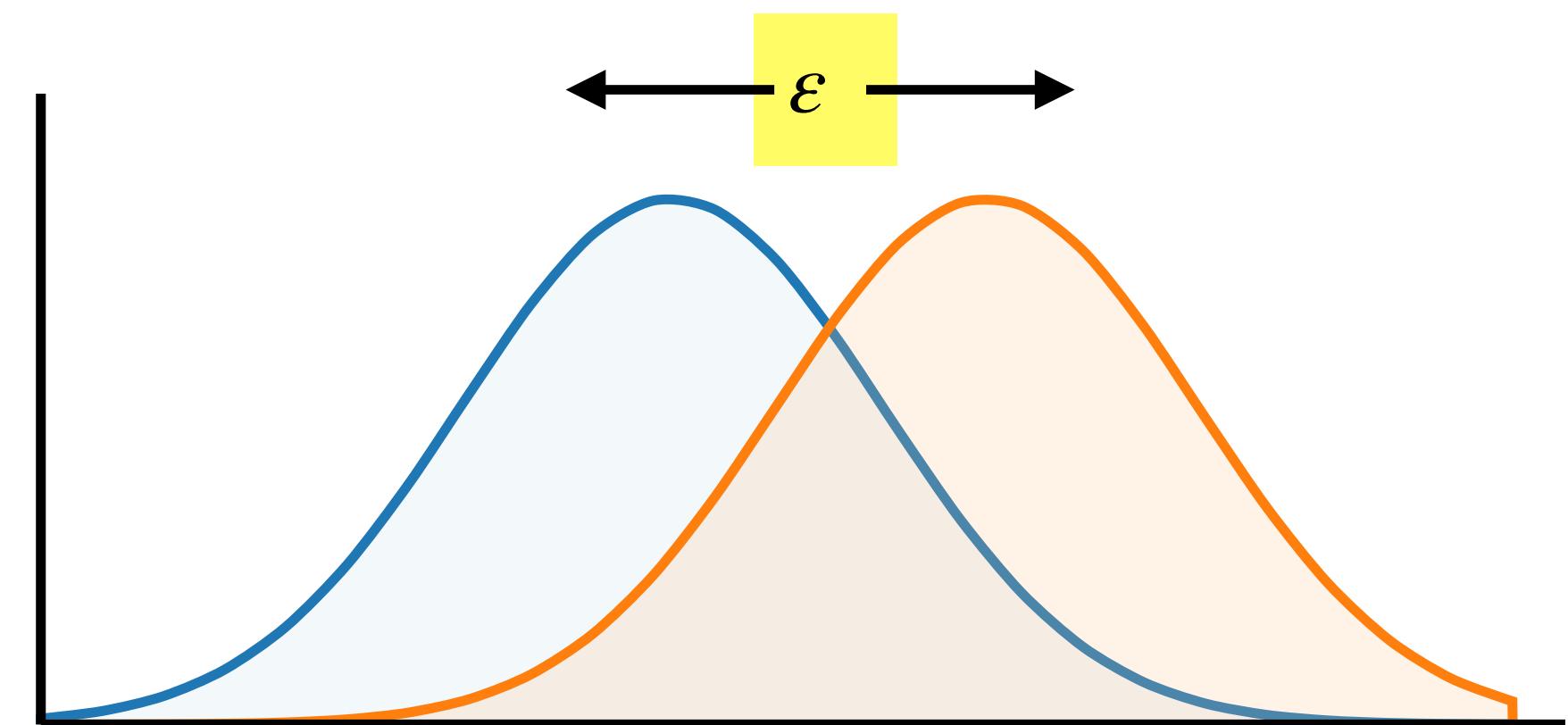


Output Distribution
(e.g. over models)



ϵ -differential privacy

Large $\epsilon \implies$ more privacy leakage



Adding noise for DP

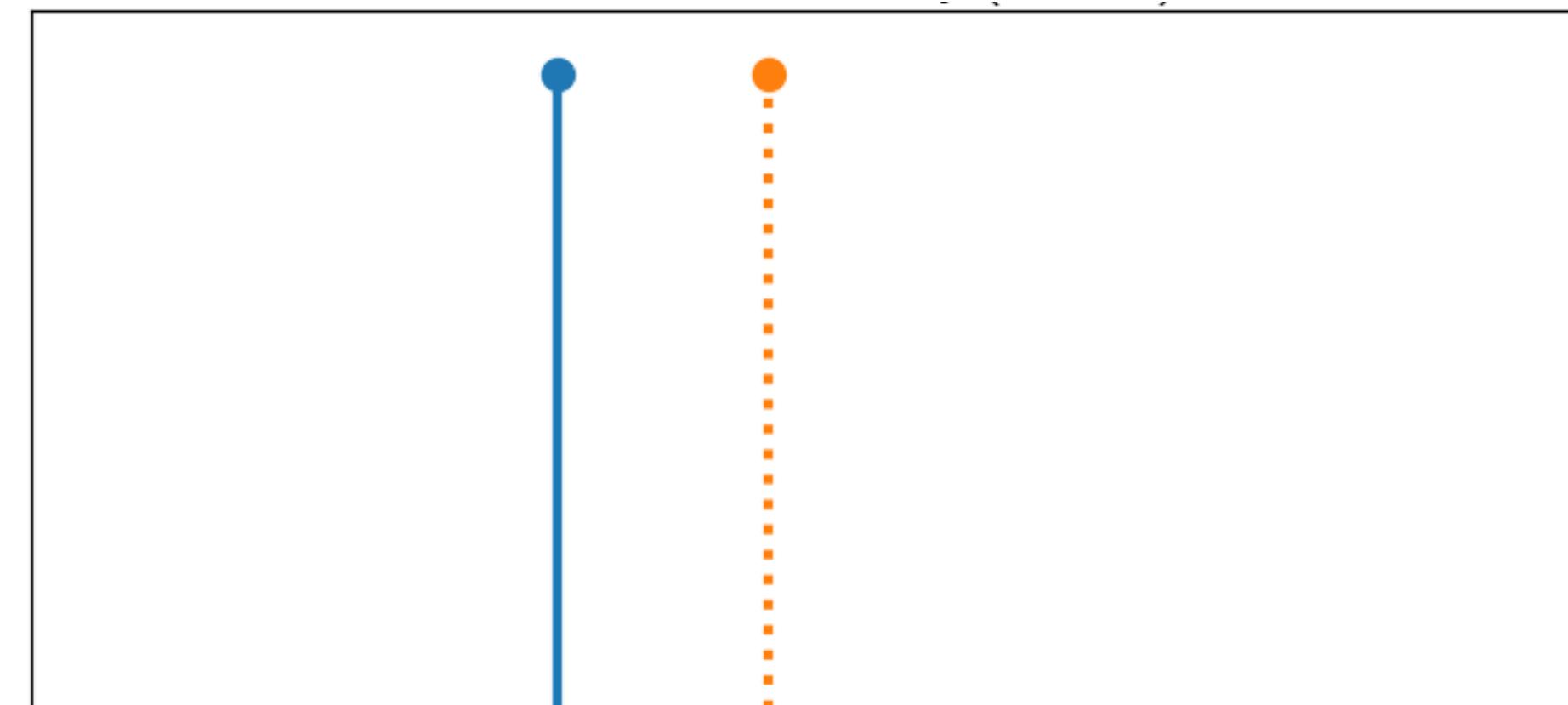
Dataset



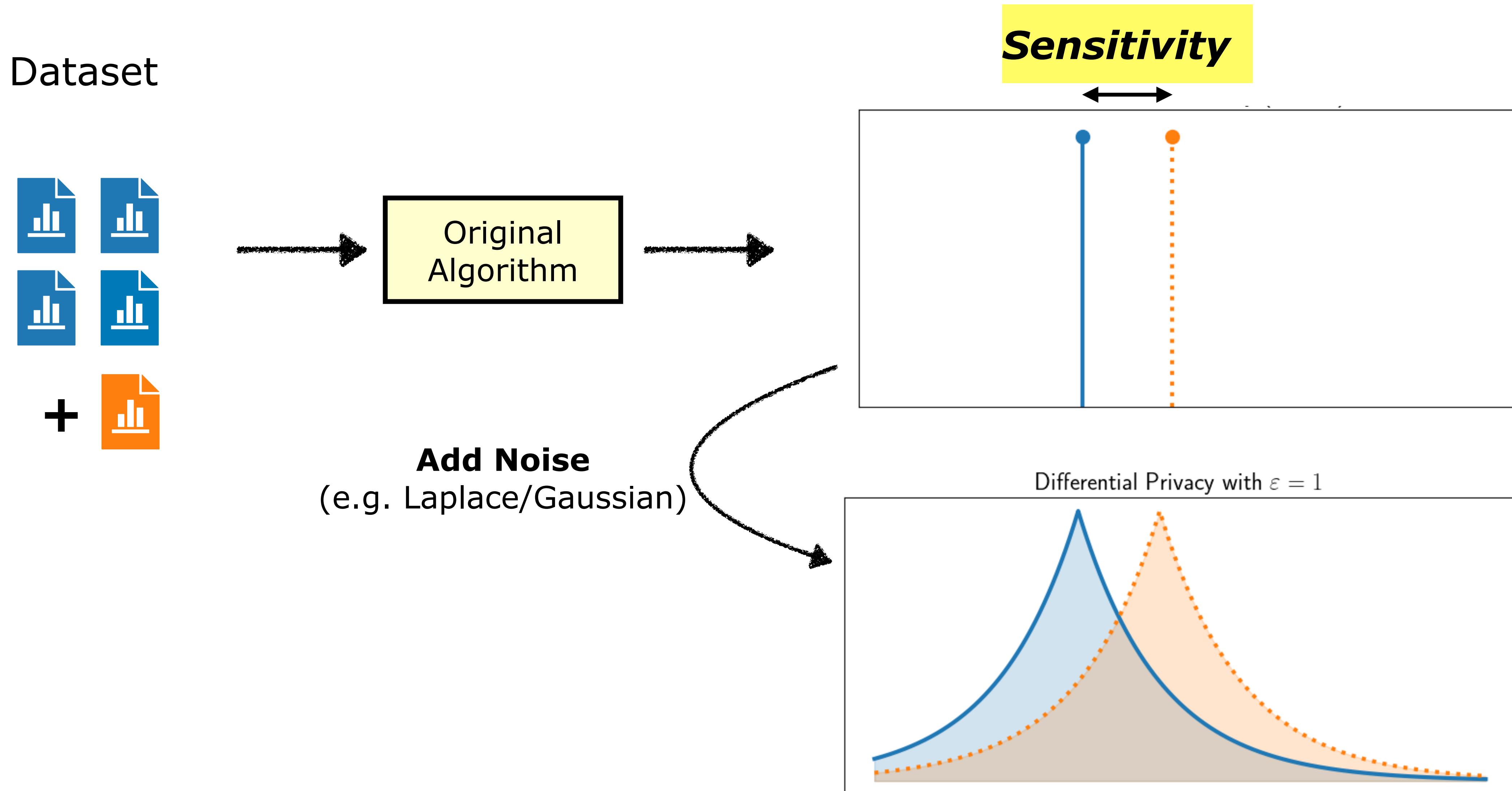
Original
Algorithm



Sensitivity

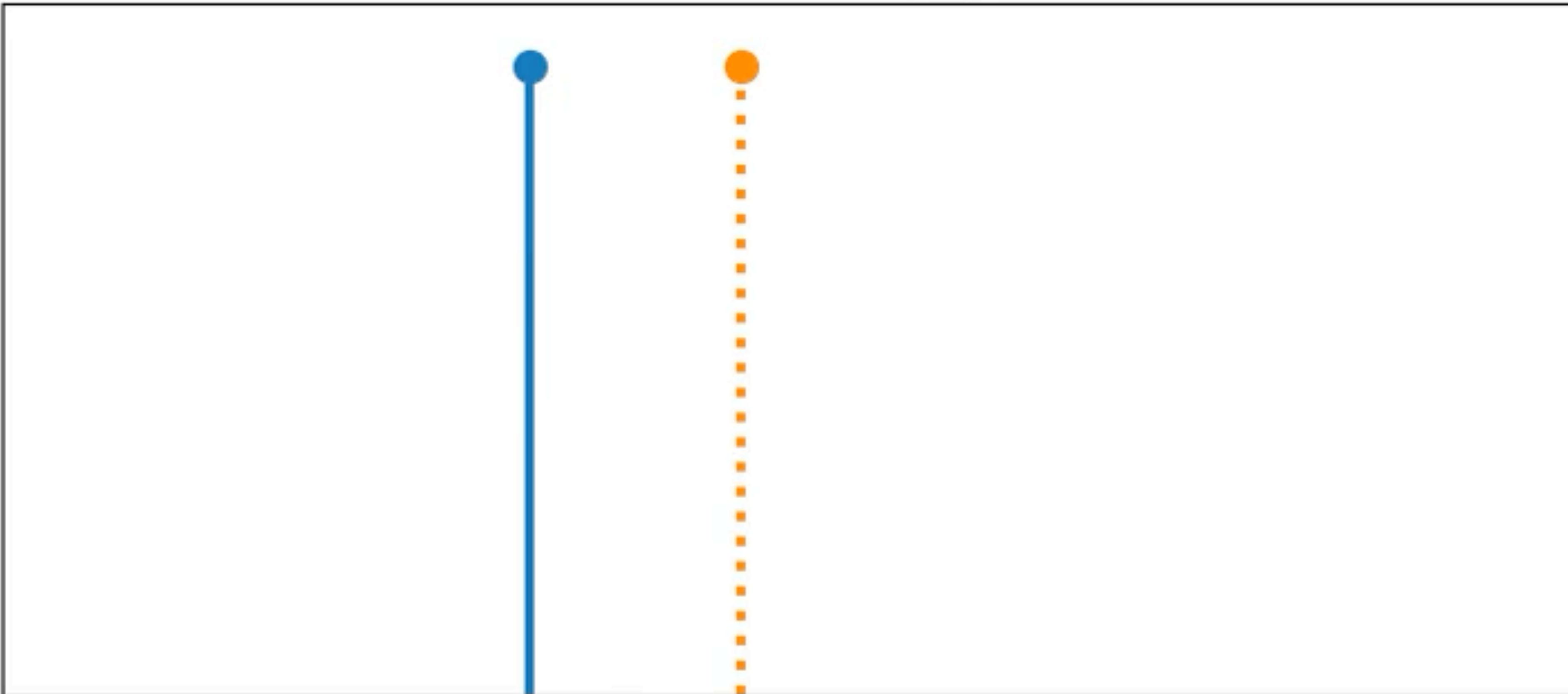


Adding noise for DP



Adding noise for DP

No Differential Privacy ($\varepsilon = \infty$)



Key properties of DP

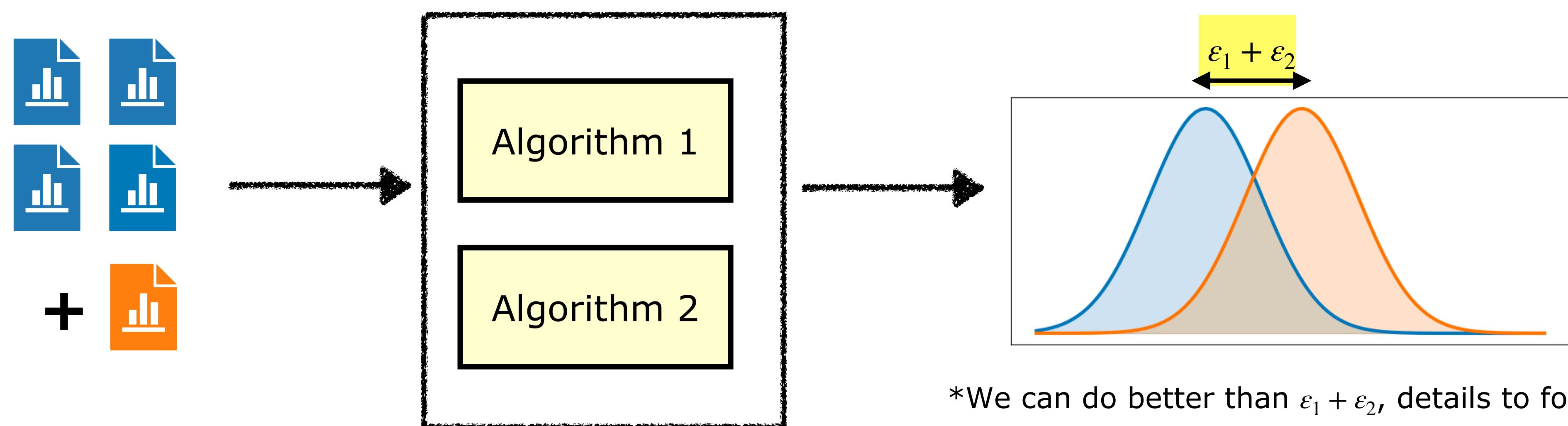
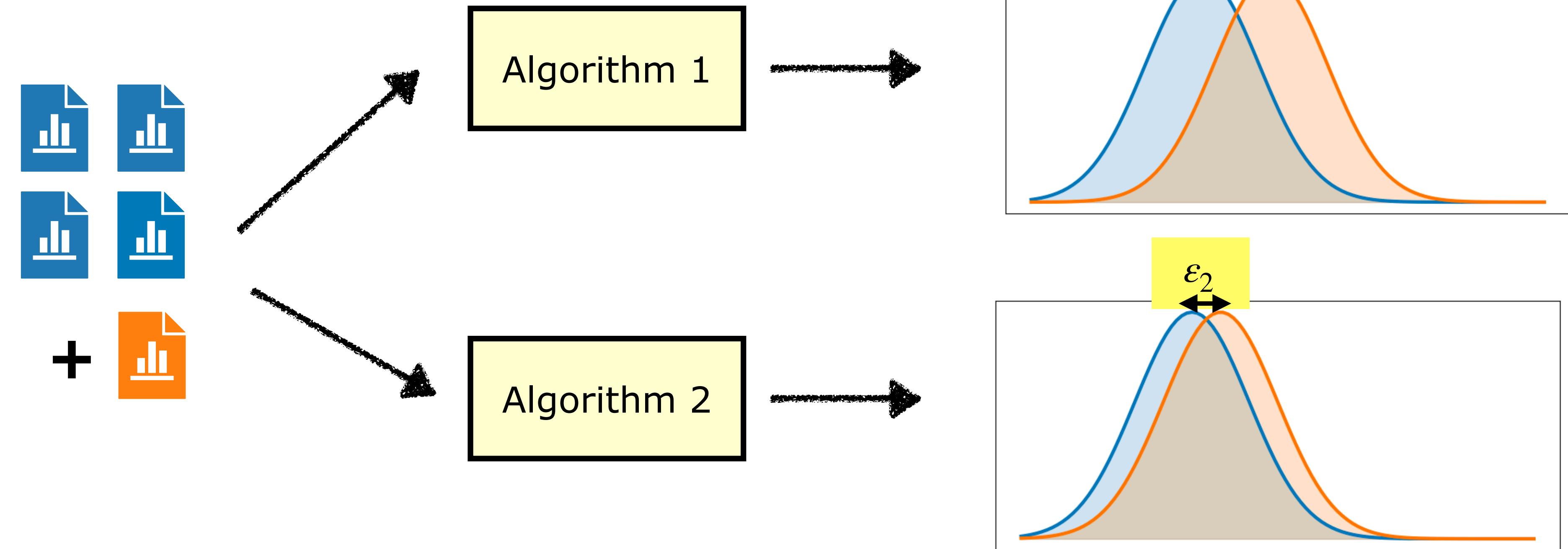


Composition over multiple steps



Post-processing

Composition



*We can do better than $\varepsilon_1 + \varepsilon_2$, details to follow

Why is composition necessary?

2006 - 2009



\$1M to beat Netflix's recommendation algorithm by 10%

Robust De-anonymization of Large Datasets
(How to Break Anonymity of the Netflix Prize Dataset)

Arvind Narayanan and Vitaly Shmatikov

The University of Texas at Austin

February 5, 2008

RYAN SINGEL SECURITY MAR 12, 2010 2:48 PM

NetFlix Cancels Recommendation Contest After Privacy Lawsuit

Netflix is canceling its second \$1 million Netflix Prize to settle a legal challenge that it breached customer privacy as part of the first contest's race for a better movie-recommendation engine. Friday's announcement came five months after Netflix had announced a successor to its algorithm-improvement contest. The company at the time said it intended to [...]

Why is composition necessary?

2006 - 2009

Robust De-anonymization of Large Datasets
(How to Break Anonymity of the Netflix Prize Dataset)

Arvind Narayanan and Vitaly Shmatikov

The University of Texas at Austin

Composition prevents such leakage!

RYAN SINGEL

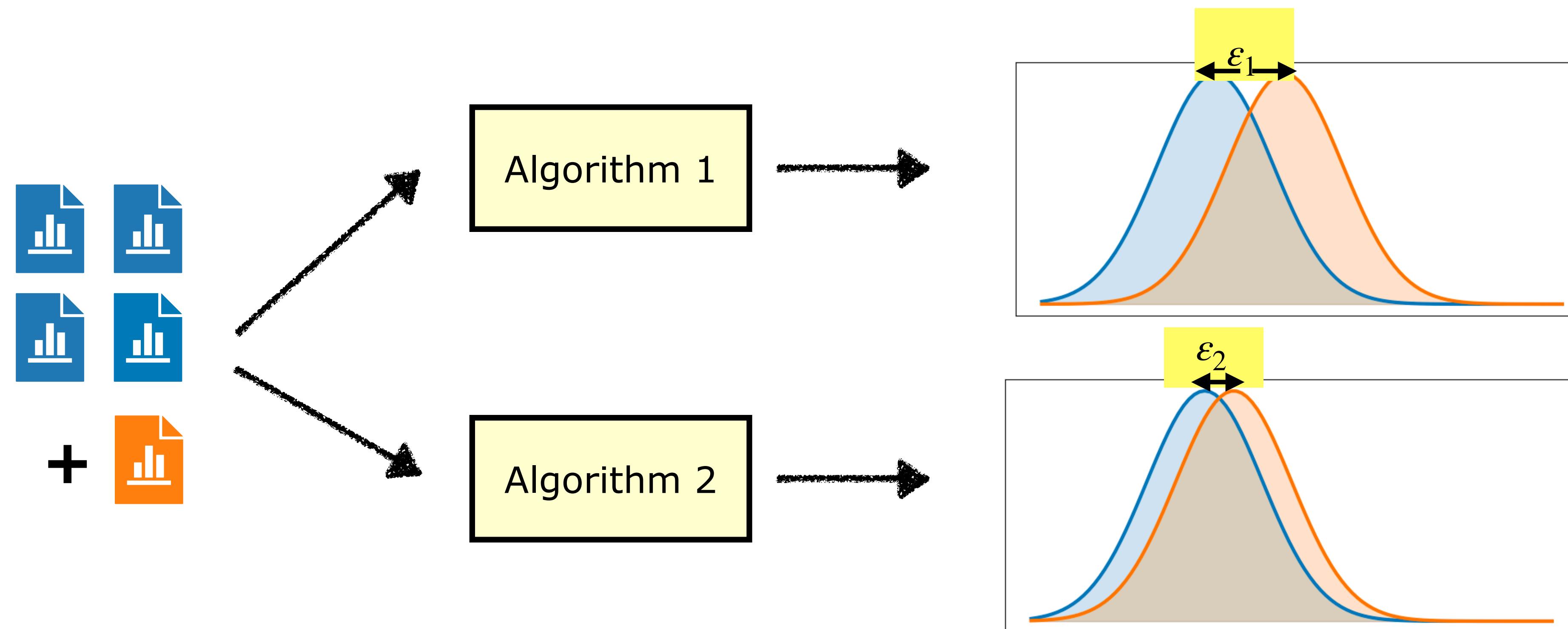
SECURITY MAR 12, 2010 2:48 PM

NetFlix Cancels Recommendation Contest After Privacy Lawsuit

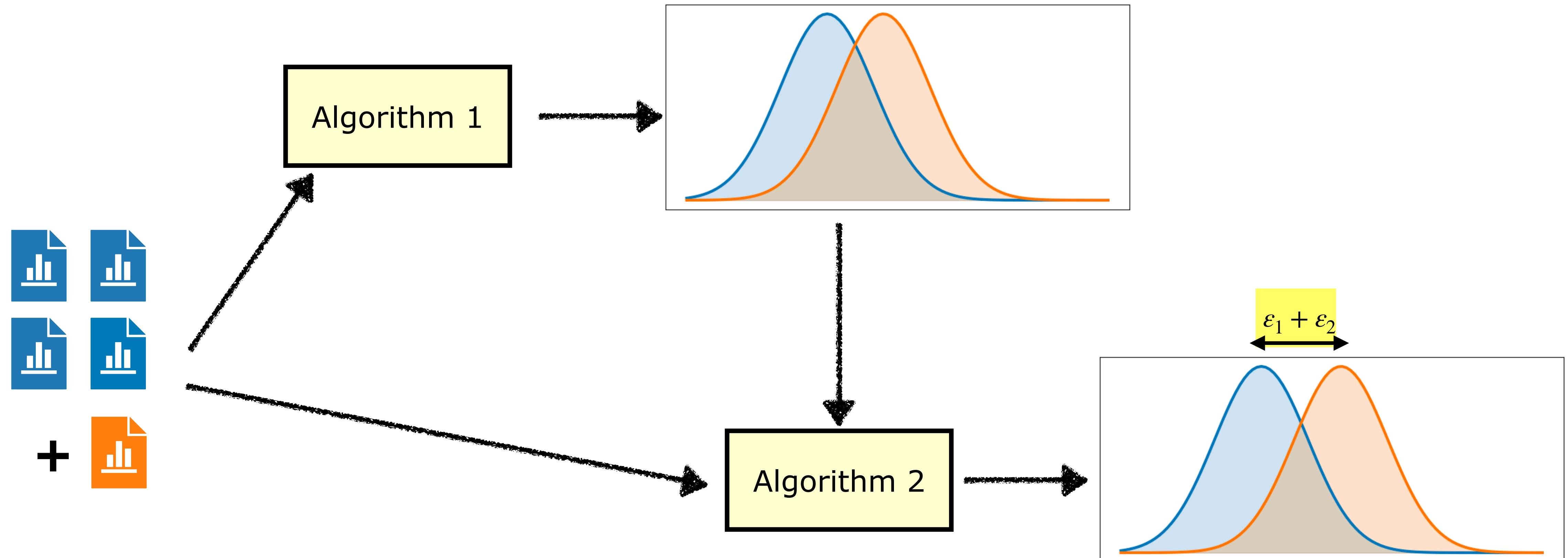
\$1M to beat Netflix's recommendation algorithm by 10%

Netflix is canceling its second \$1 million Netflix Prize to settle a legal challenge that it breached customer privacy as part of the first contest's race for a better movie-recommendation engine. Friday's announcement came five months after Netflix had announced a successor to its algorithm-improvement contest. The company at the time said it intended to [...]

Adaptive Composition

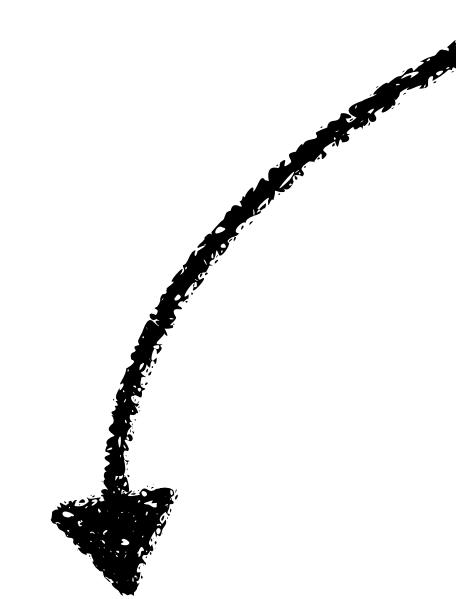


Adaptive Composition



*We can do better than $\epsilon_1 + \epsilon_2$, details to follow

Key properties of DP



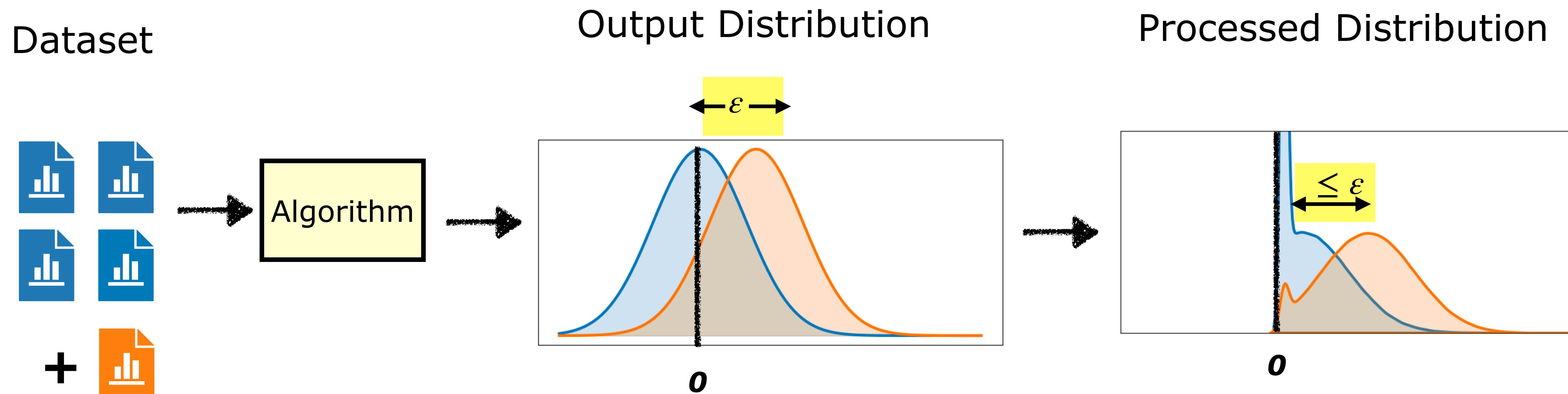
Composition over multiple steps



Post-processing

Post-processing

Example: How many people at IITM have a certain medical condition?



“Information cannot be created”

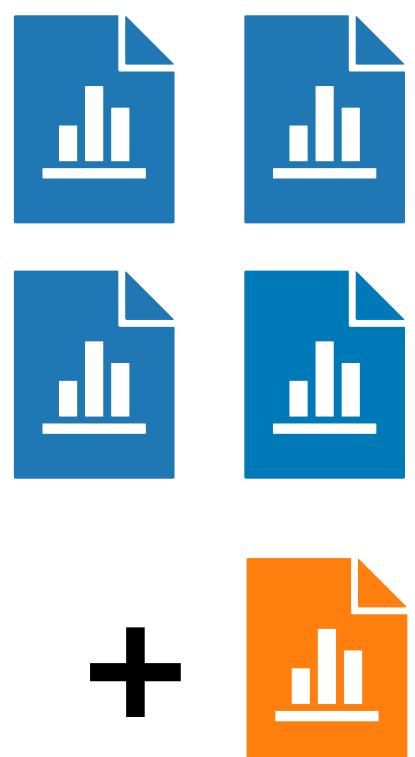
Post-processing

LLM Example:

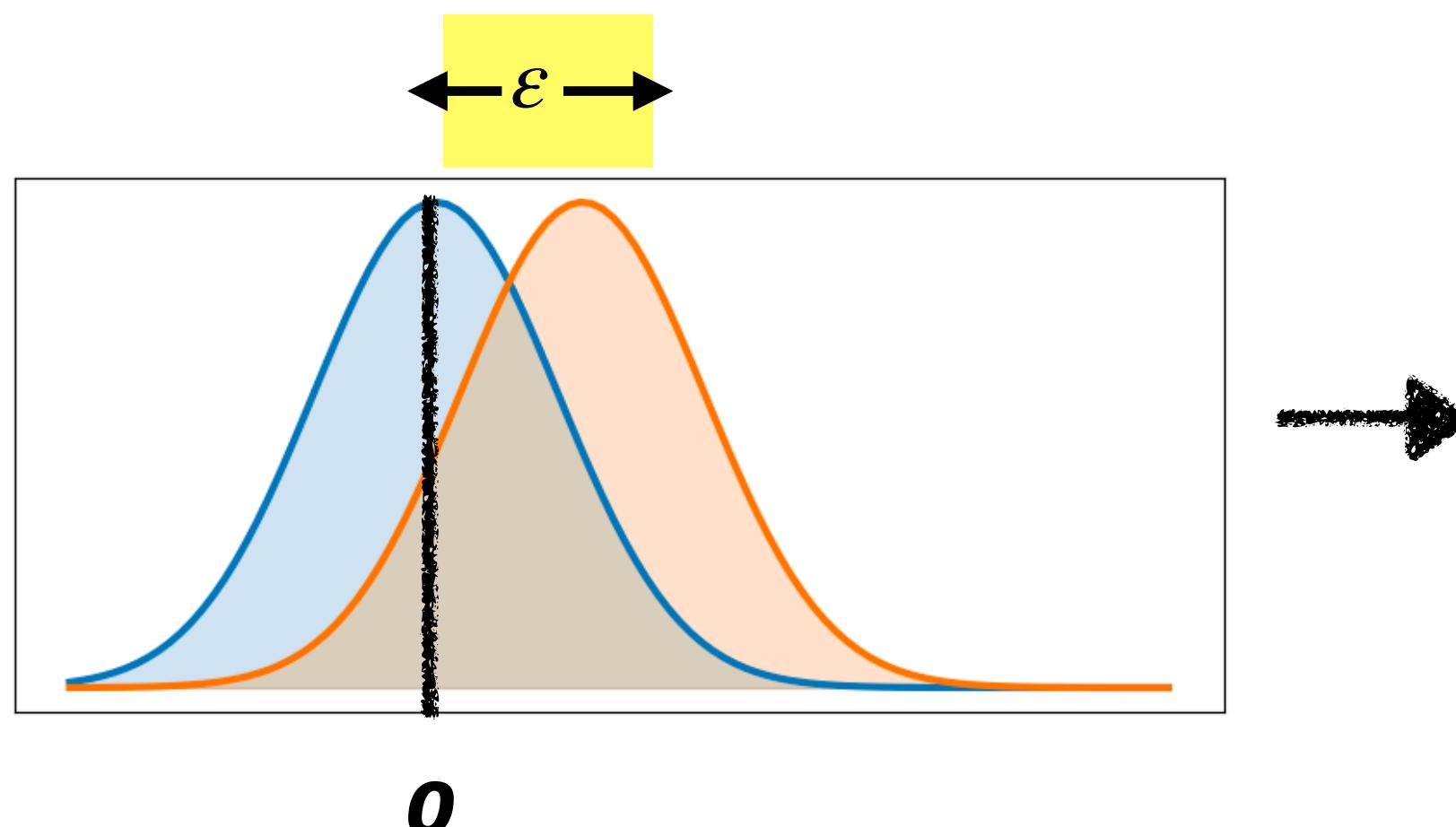
Stage 1 - Private Training

/ Stage 2 - alignment

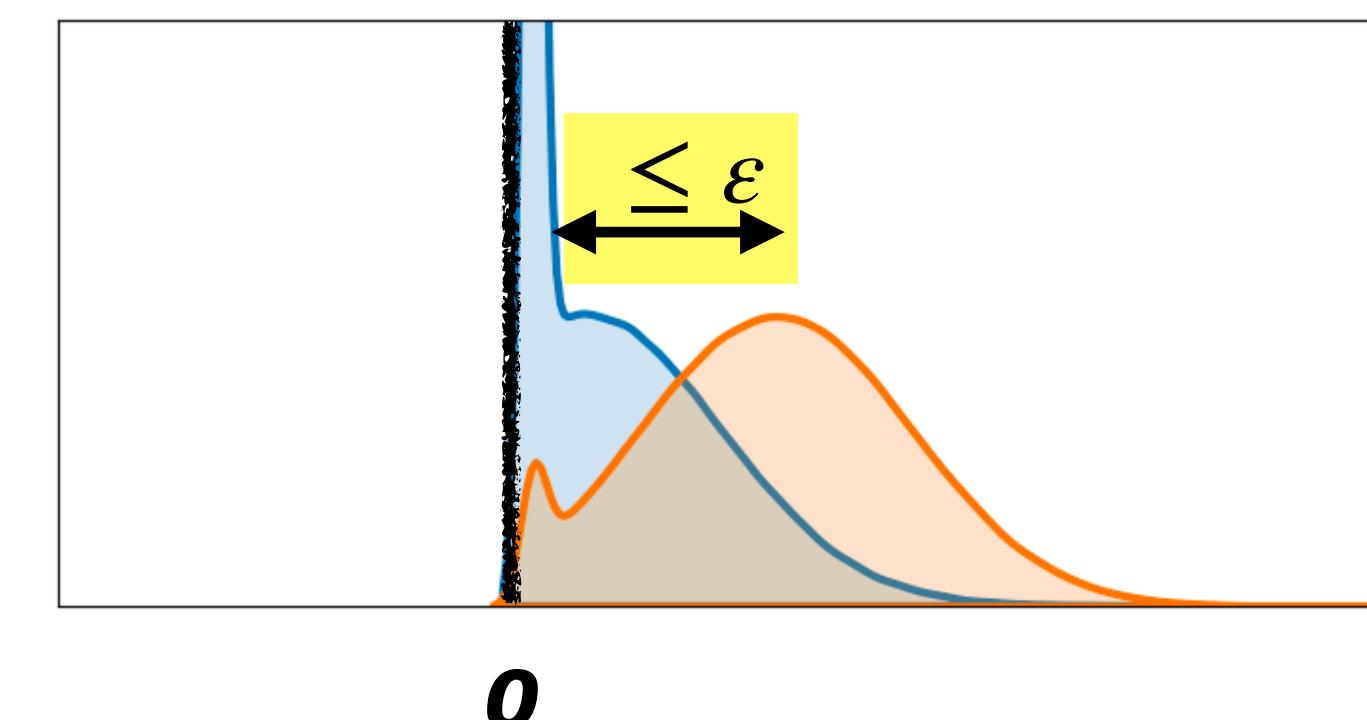
Dataset



Output Distribution



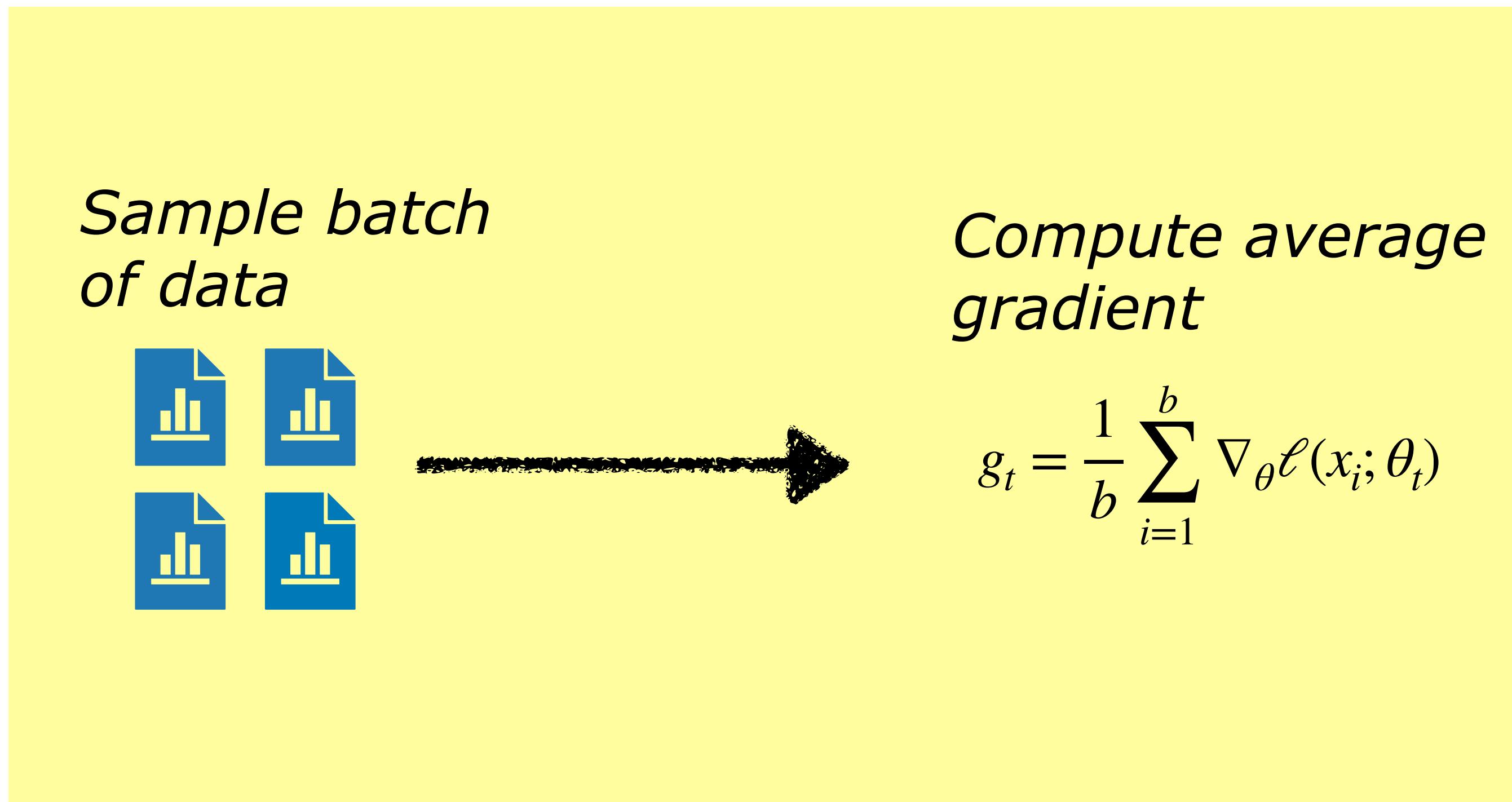
Processed Distribution



“Information cannot be created”

Example: Deep Learning with DP

Review: Stochastic gradient descent (without DP)

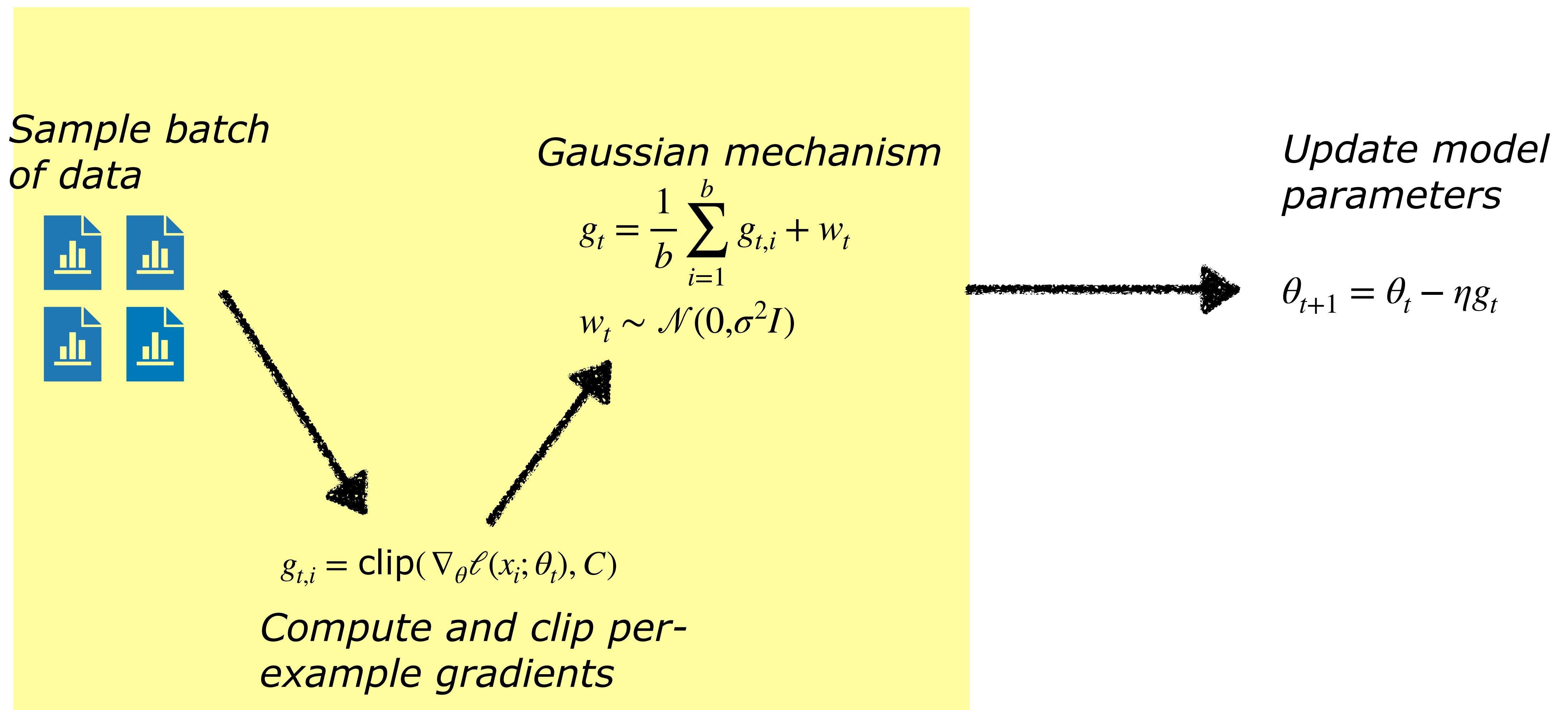


Update model parameters

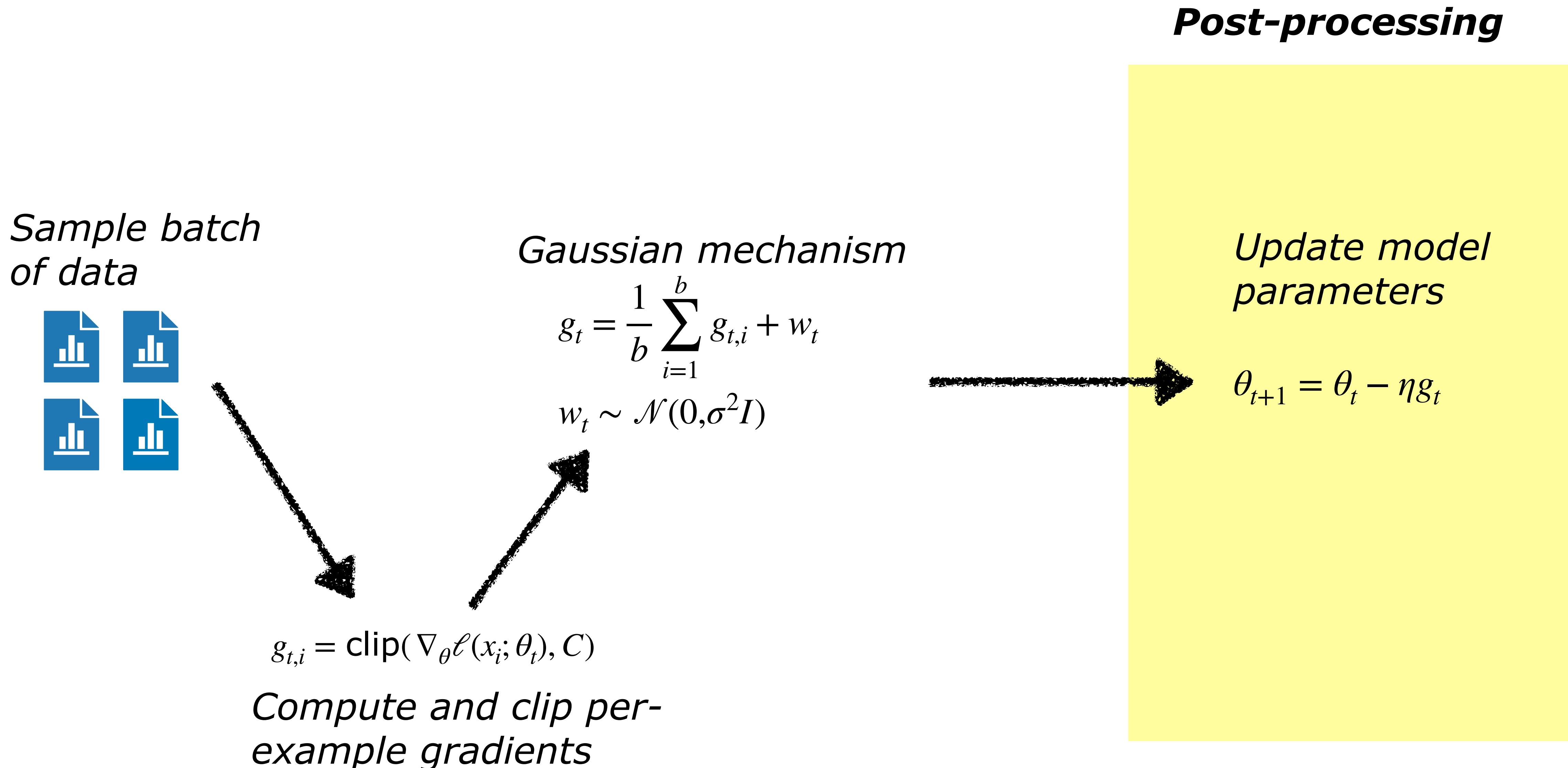
$$\theta_{t+1} = \theta_t - \eta g_t$$

DP-SGD: Stochastic gradient descent with DP

Gradient is differentially private

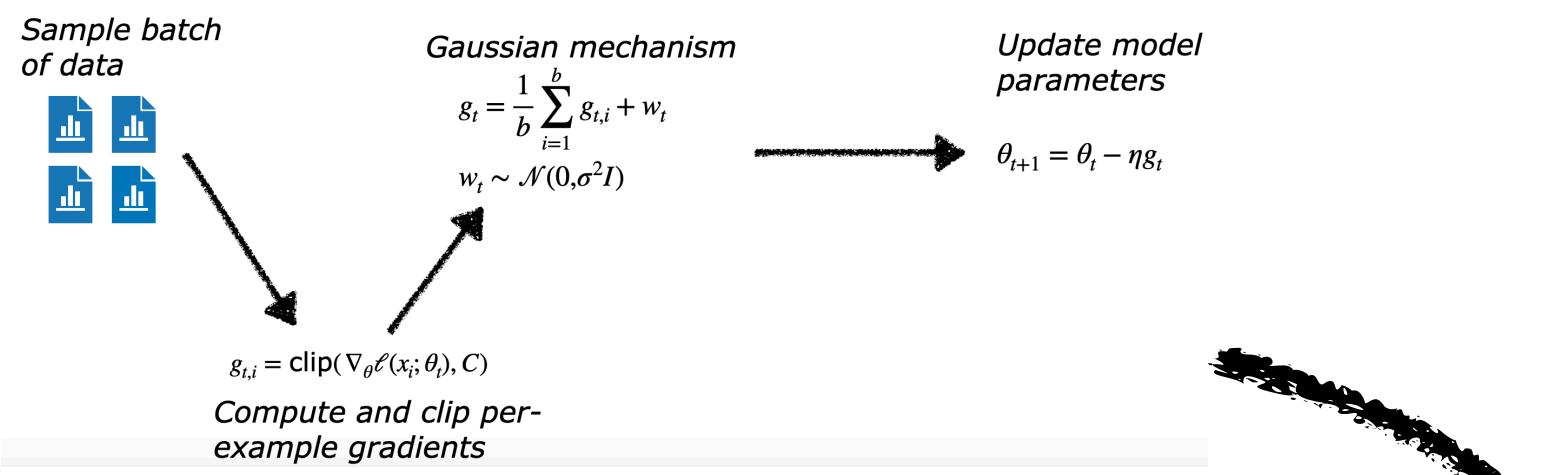


DP-SGD: Stochastic gradient descent with DP

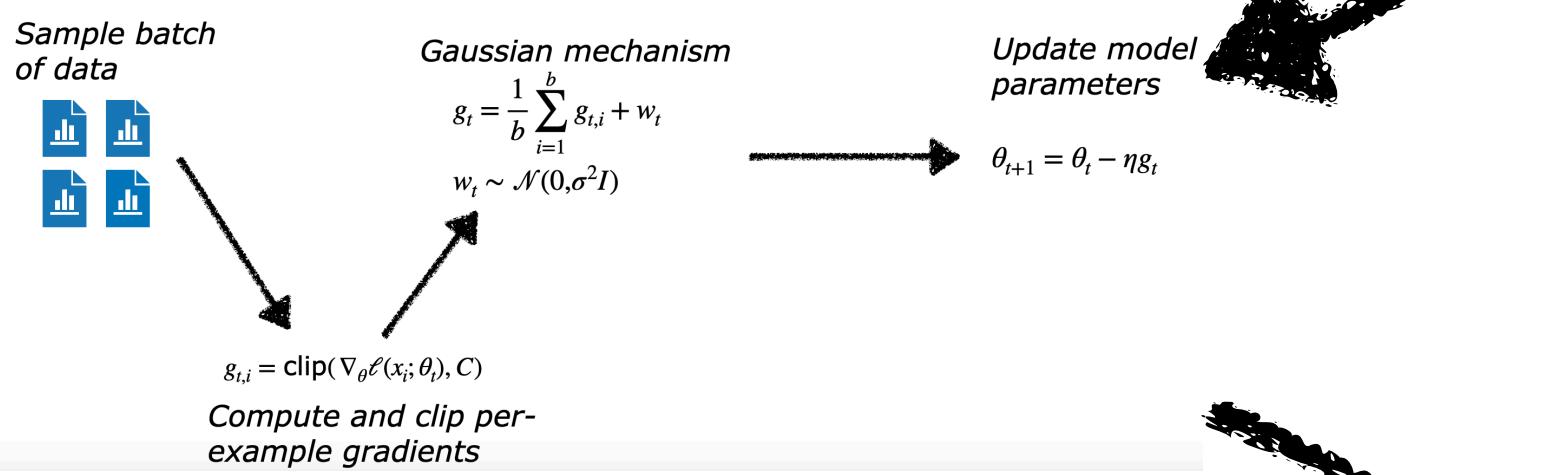


DP-SGD: Stochastic gradient descent with DP

Iteration 1

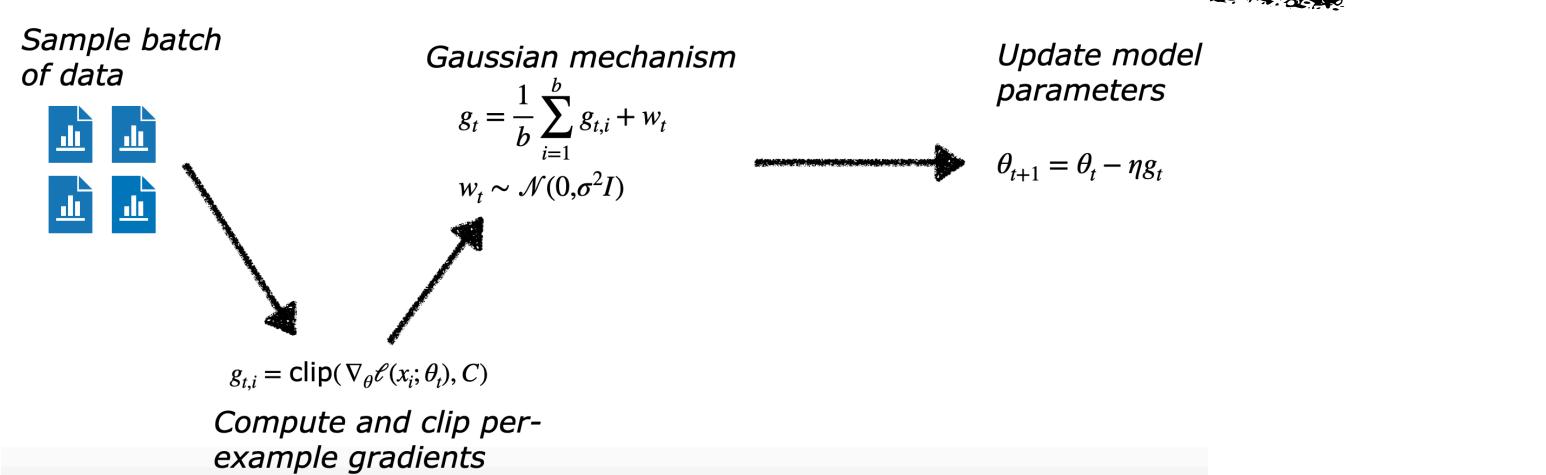


Iteration 2



⋮

Iteration T

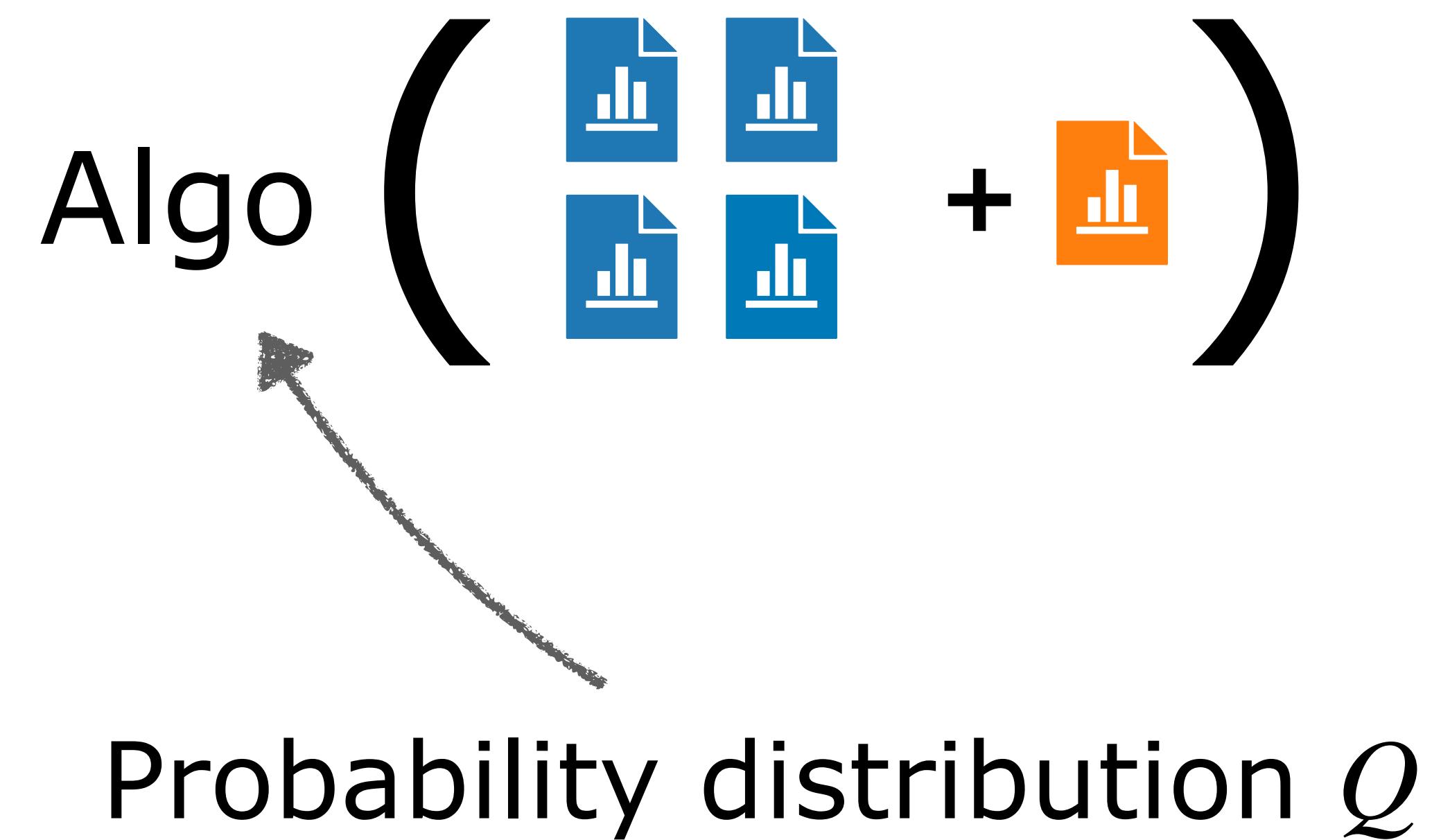
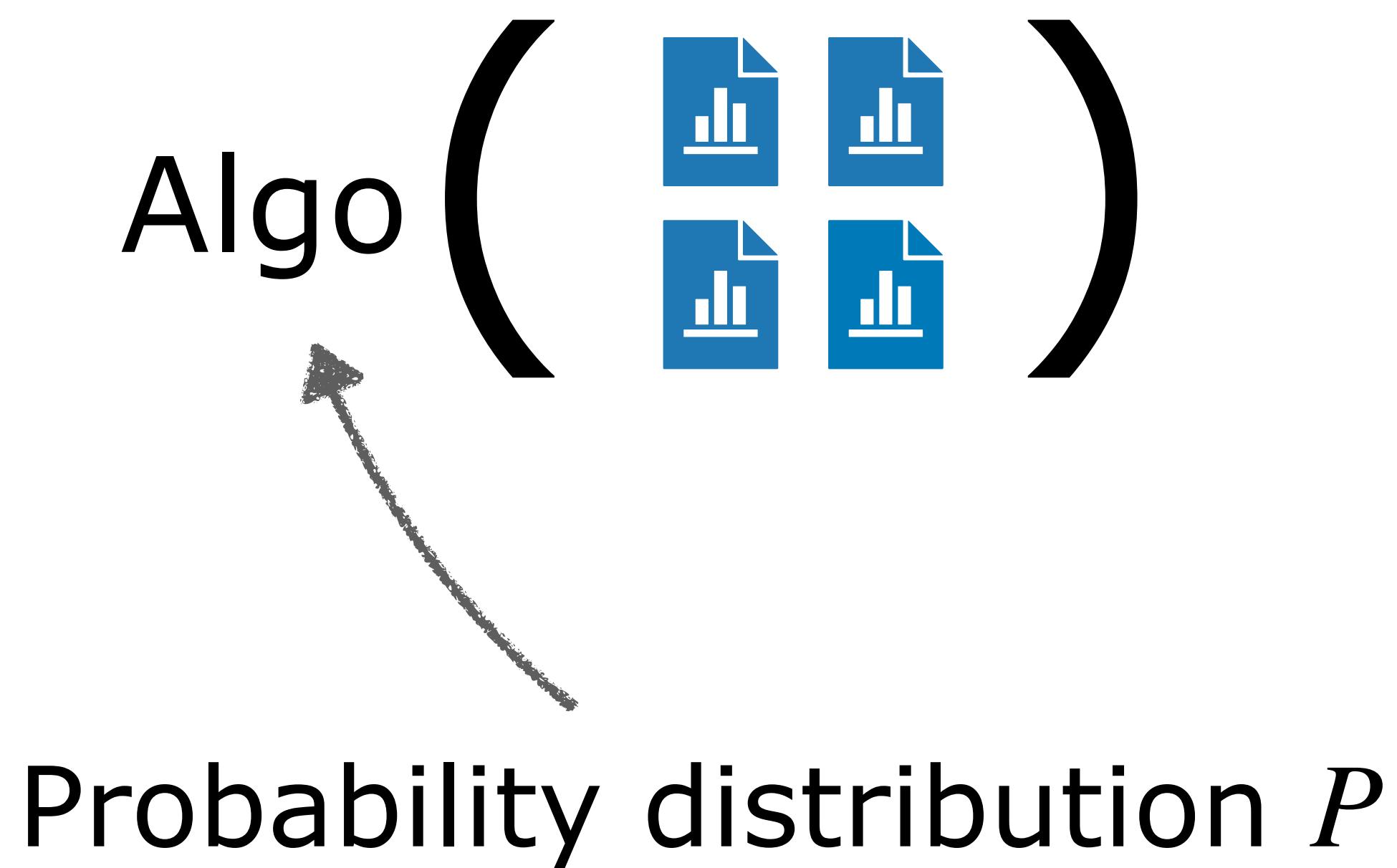


Adaptive Composition!

Outline

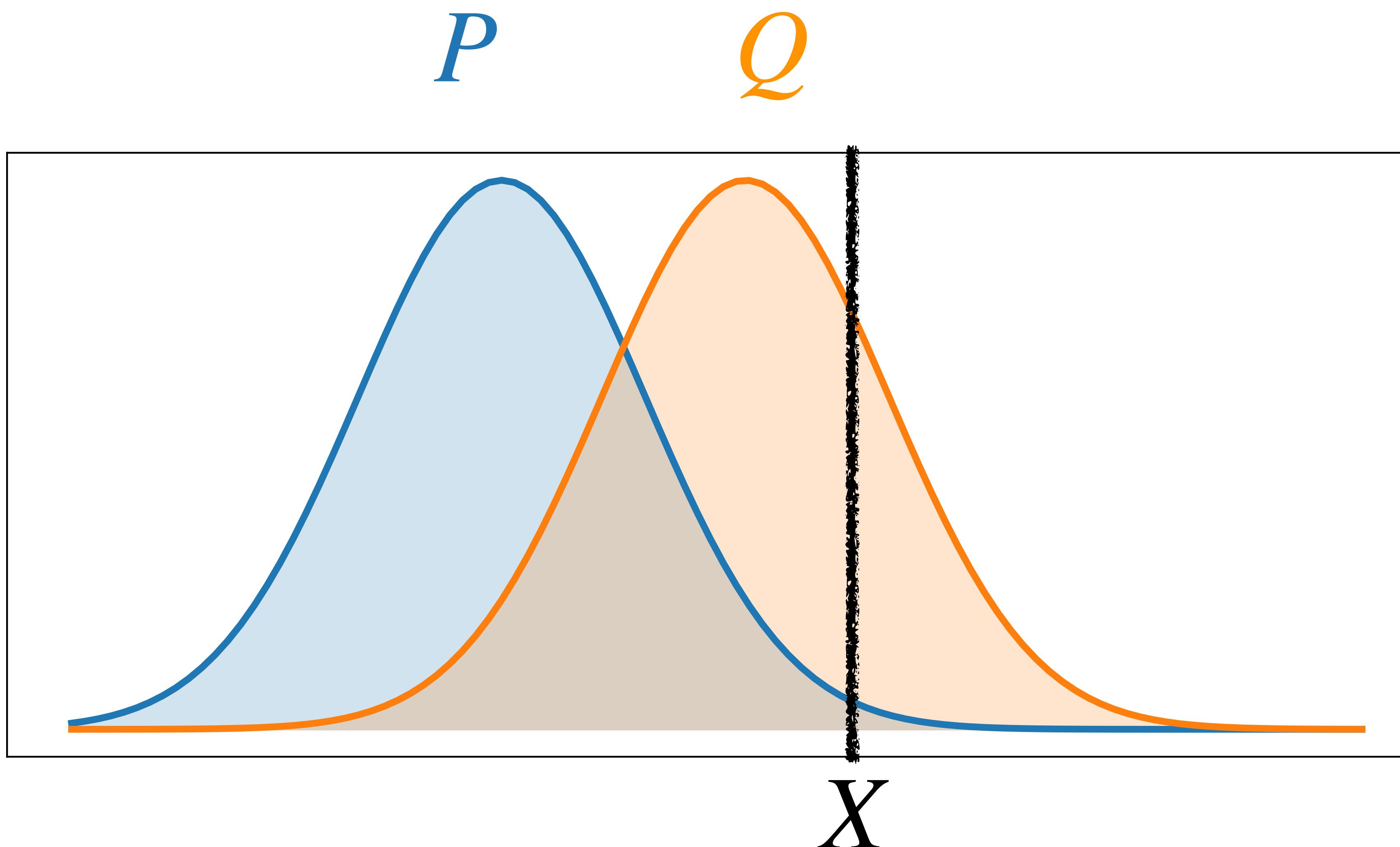
- Differential Privacy: Intuition & Recap
- ***Differential Privacy: Rigorous Mathematical Formulation & Properties***
- Application Examples

Differential privacy \Rightarrow
probability distributions $P \approx Q$ are indistinguishable



Hypothesis testing:

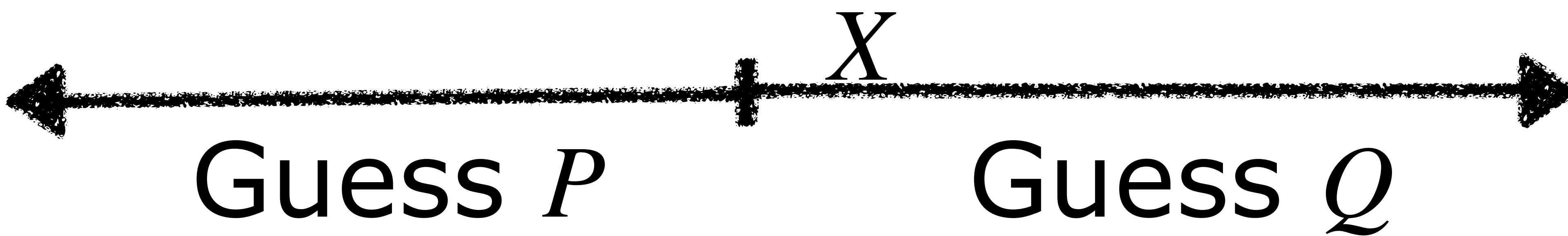
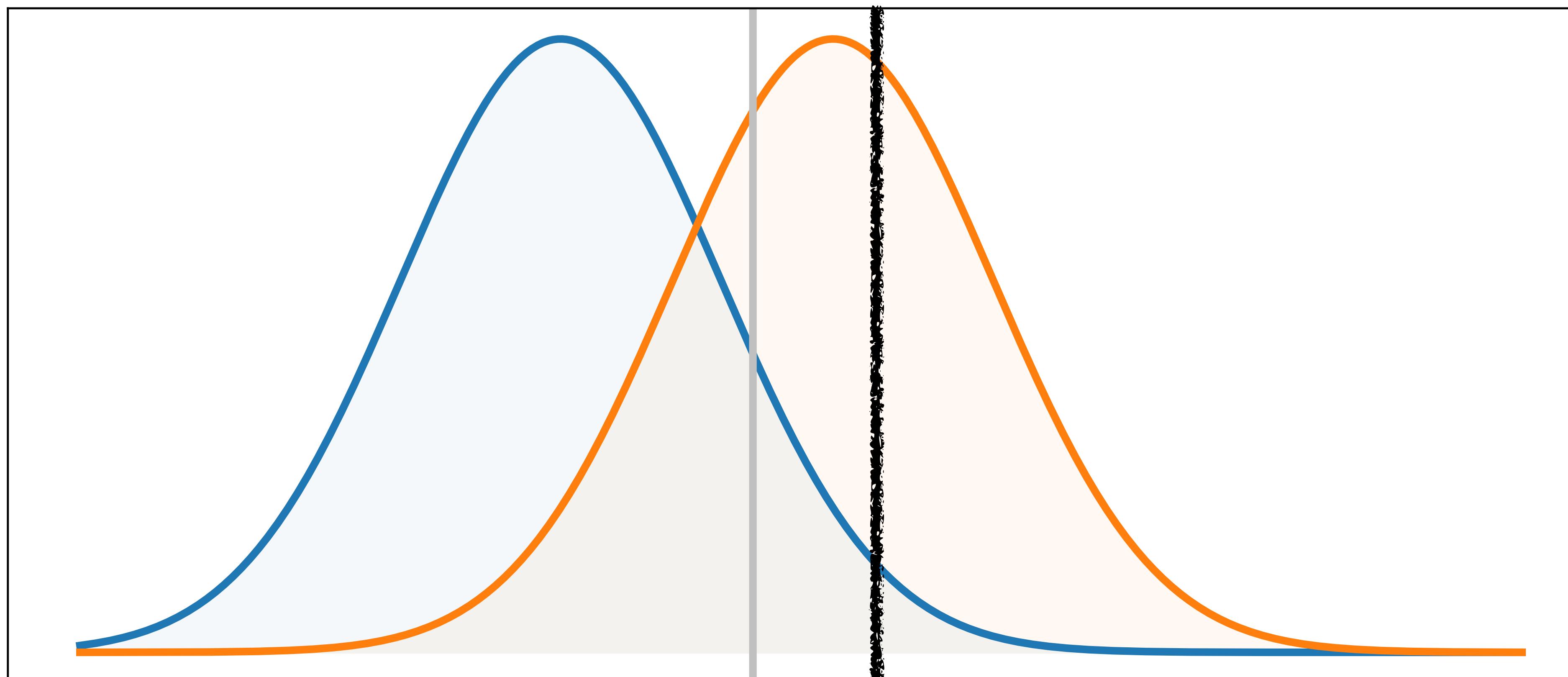
To distinguish between two distributions



Is
 $X \sim P$
or
 $X \sim Q?$

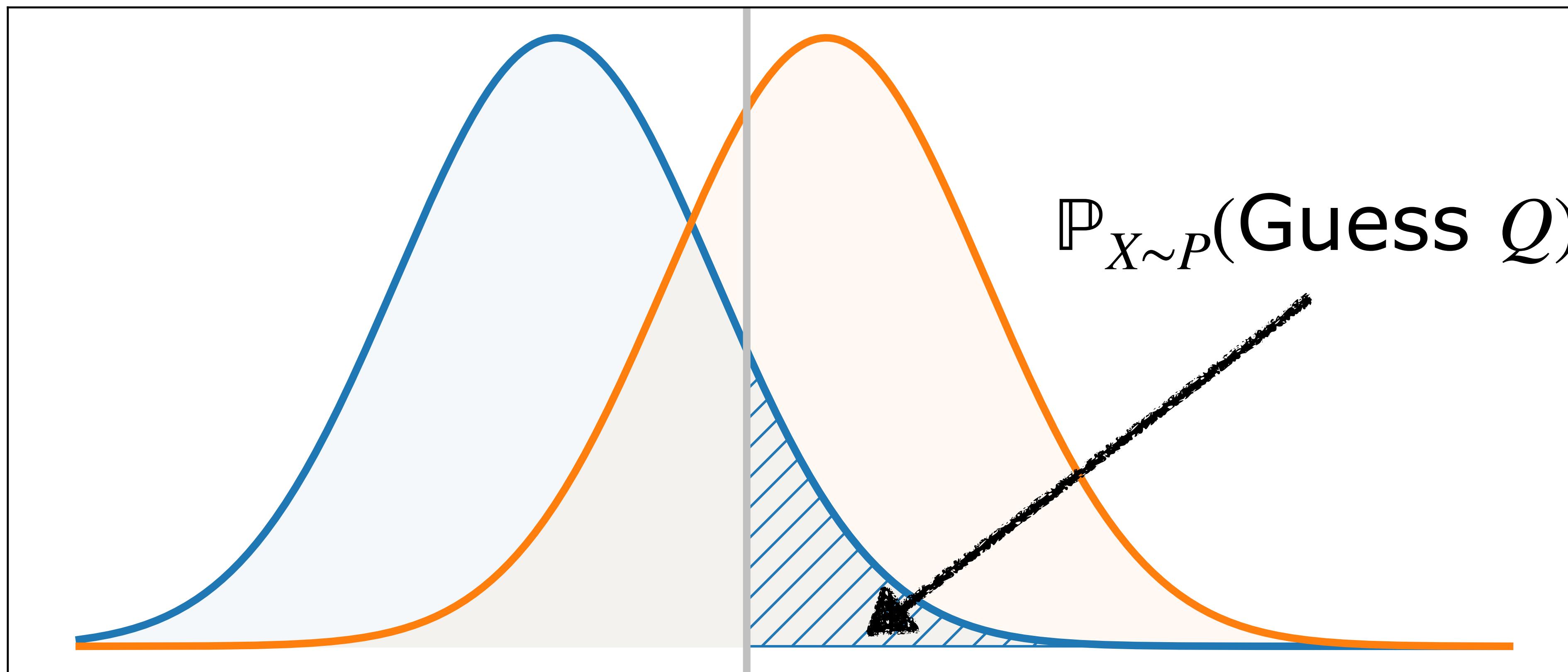
P

Q



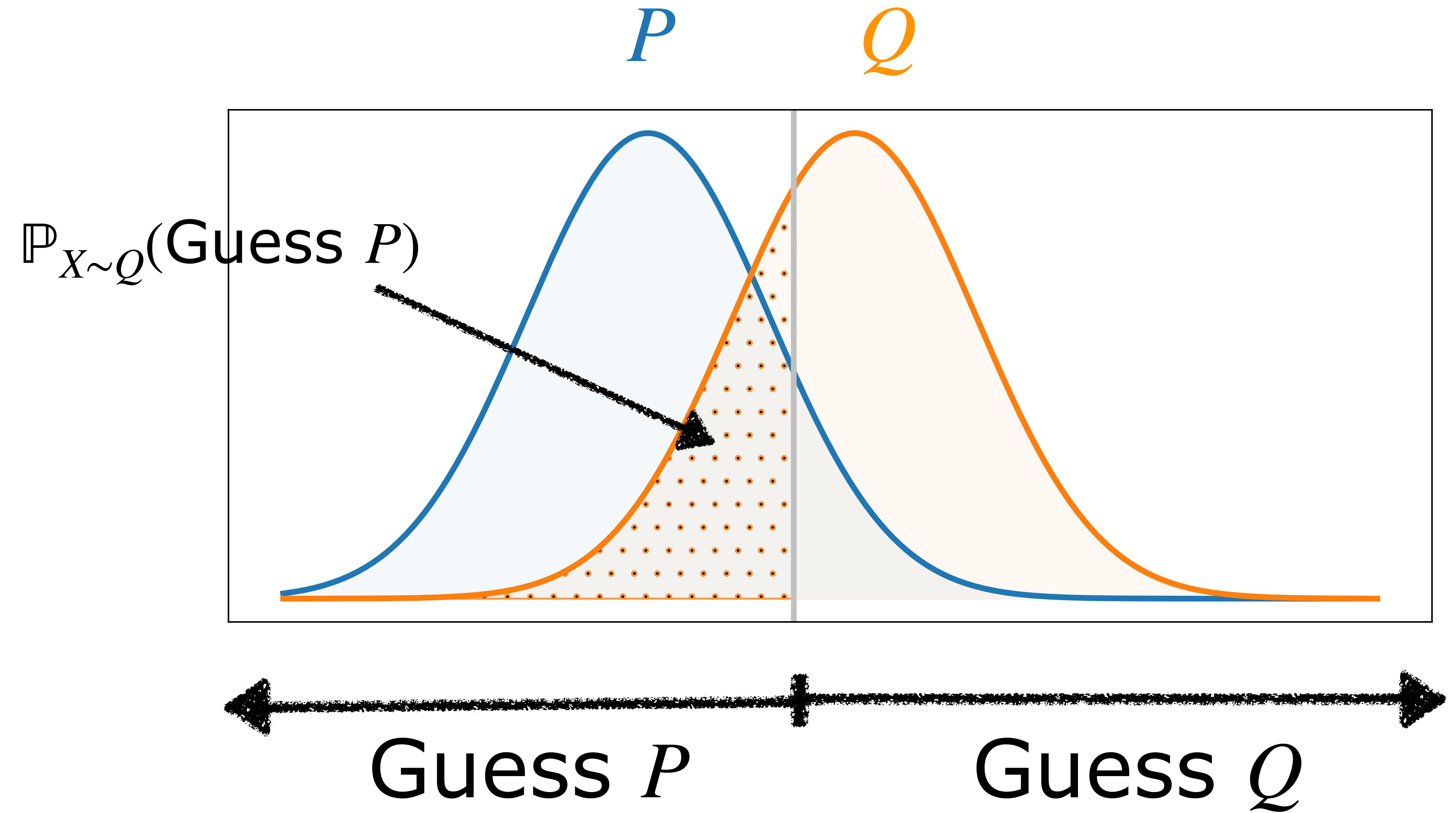
P

Q

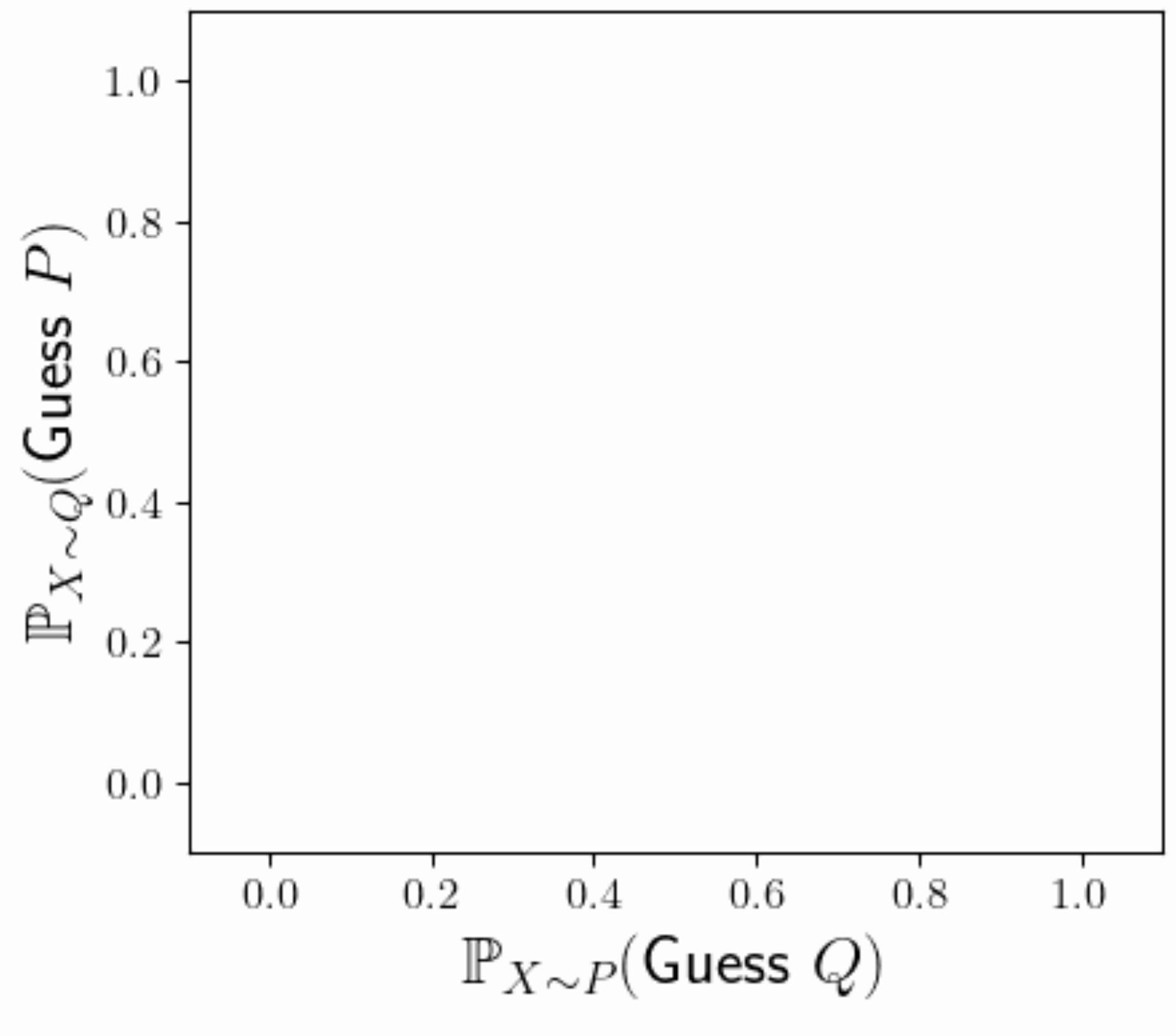
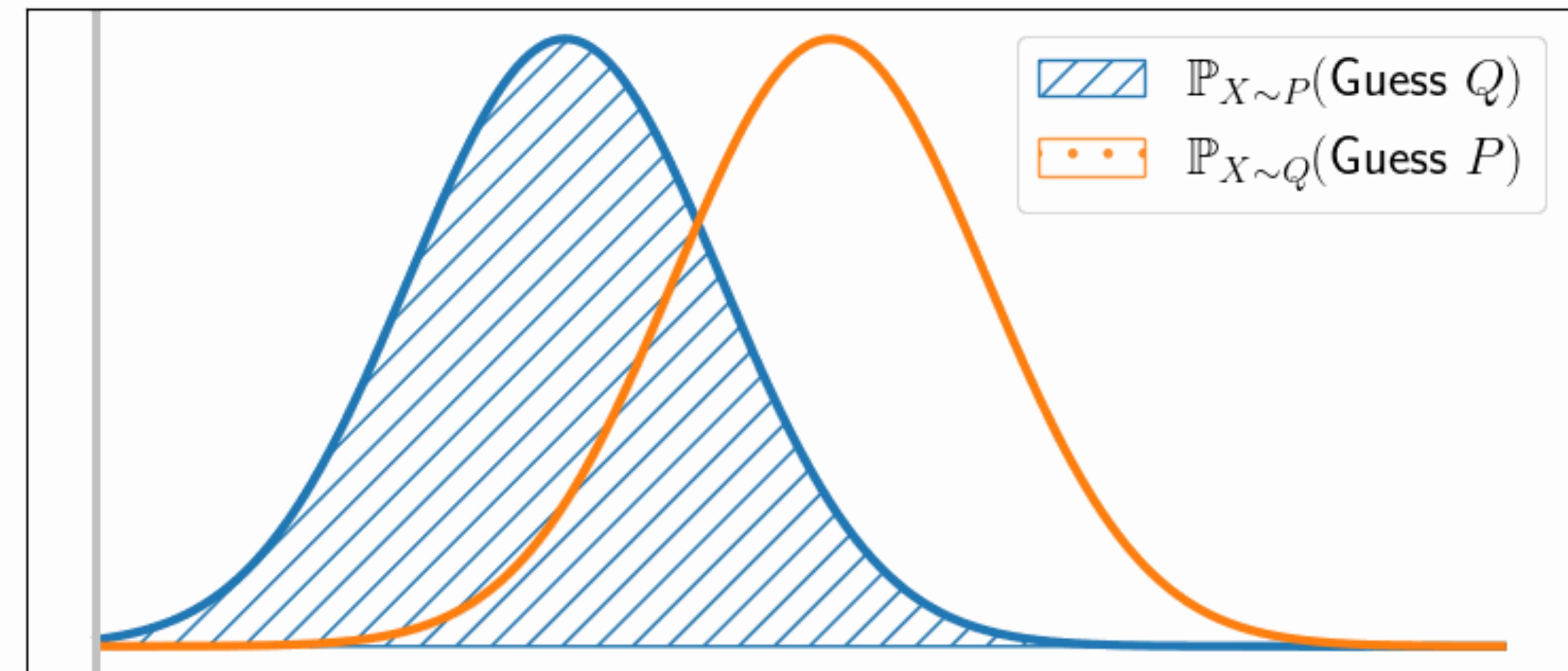


Guess *P*

Guess *Q*



P *Q*



Neyman-Pearson Lemma

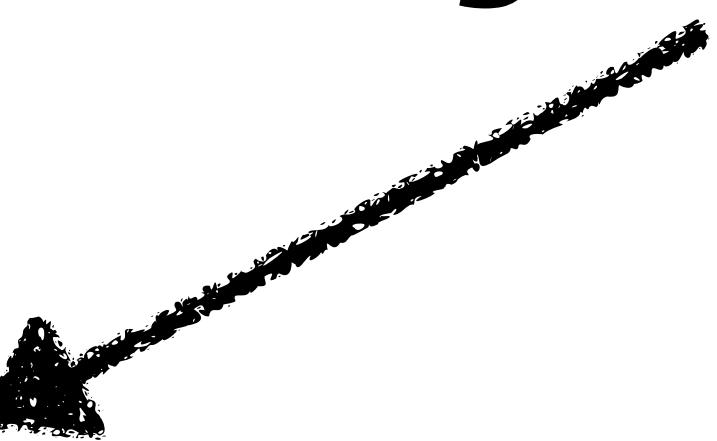
f_P : PDF of P

f_Q : PDF of Q

Define:

$$T = \log \left(\frac{f_P(X)}{f_Q(X)} \right)$$

Log-likelihood ratio



The best strategy: Fix a real number t and guess P if $T < t$ and guess Q otherwise

No other strategy can get both errors smaller simultaneously

Neyman-Pearson Lemma

f_P : PDF of P

f_Q : PDF of Q

Define:

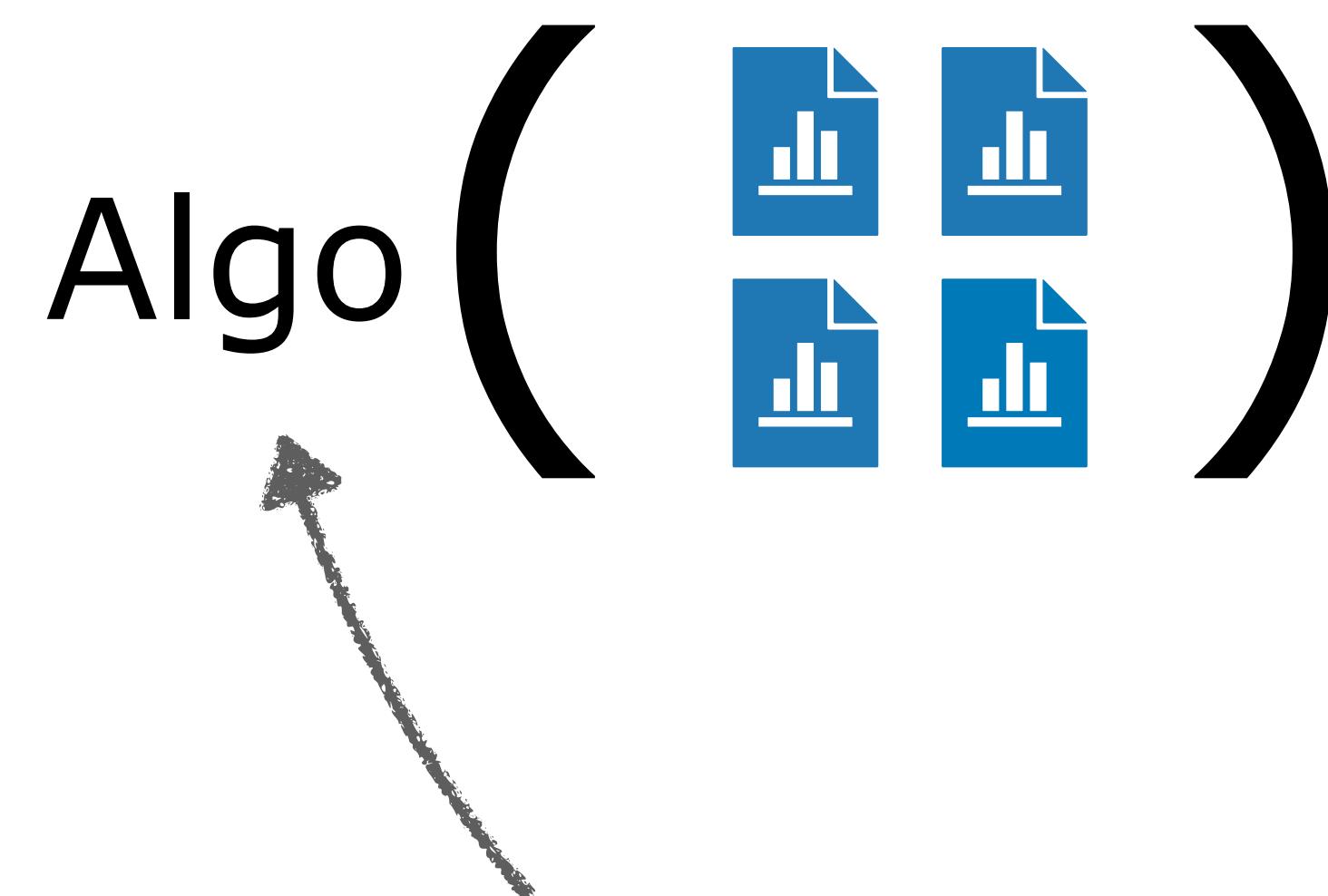
$$T = \log \left(\frac{f_P(X)}{f_Q(X)} \right)$$

Log-likelihood ratio

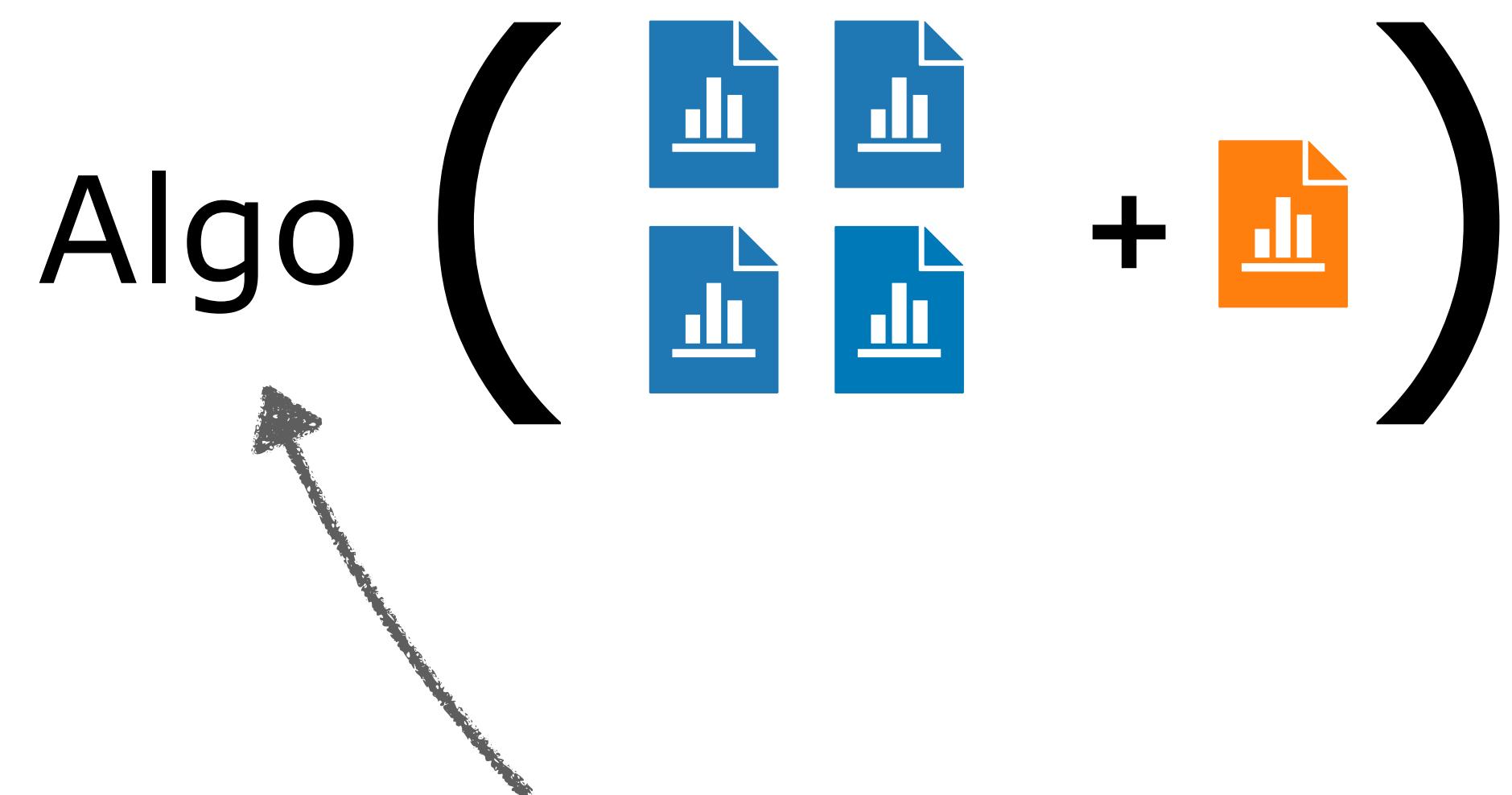


Log-likelihood ratio: the fundamental quantity for distinguishing between distributions

Differential privacy \Rightarrow
probability distributions $P \approx Q$ are indistinguishable



Probability distribution P
with density f_P



Probability distribution Q
with PDF f_Q

Privacy Loss Distribution & Composition

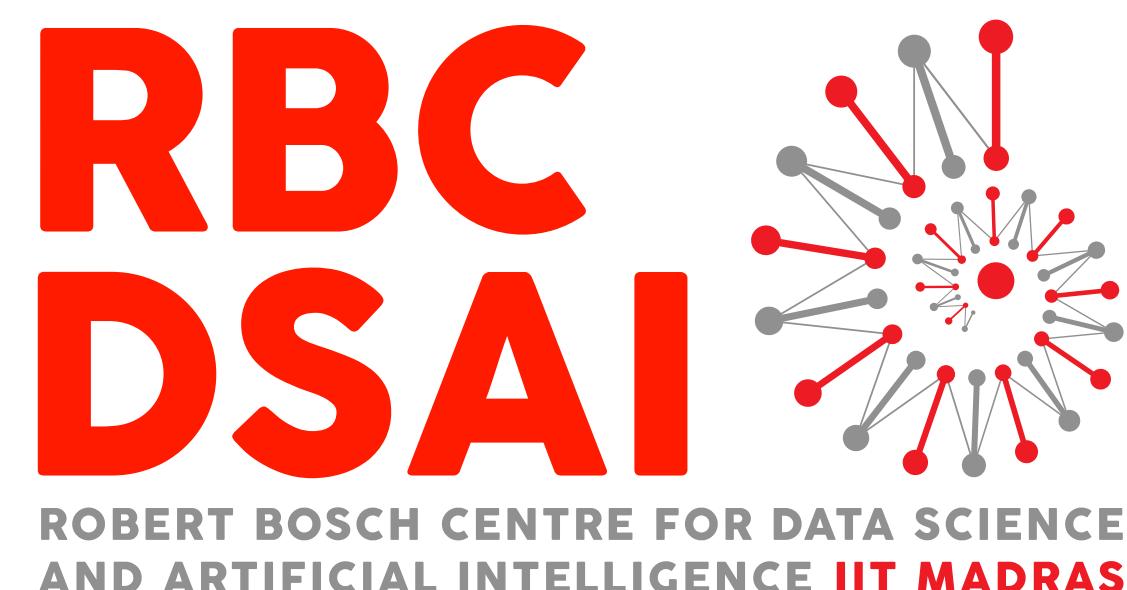
Blackboard discussion

Reading material:

<https://arxiv.org/pdf/2210.00597>

Interested in joining us? Apply here:

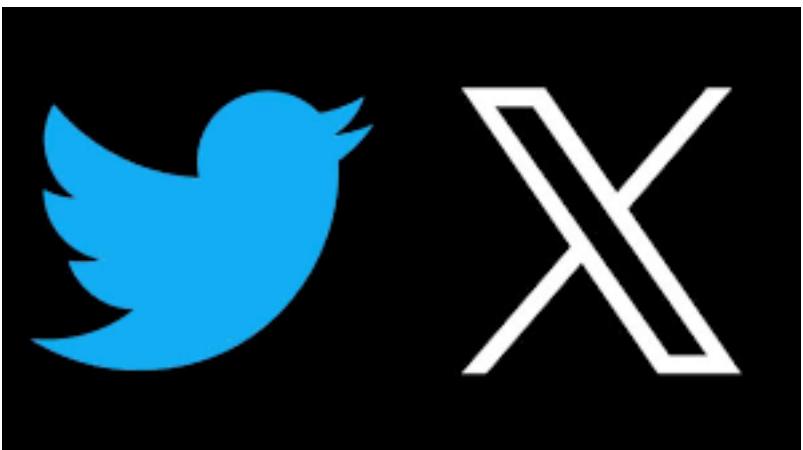
- **MS/PhD:** apply via IITM
- **Interns / pre-doctoral researchers**
 - apply directly via RBCDSAI or CeRAI



Thank you!



<https://krishnap25.github.io>



@KrishnaPillutla