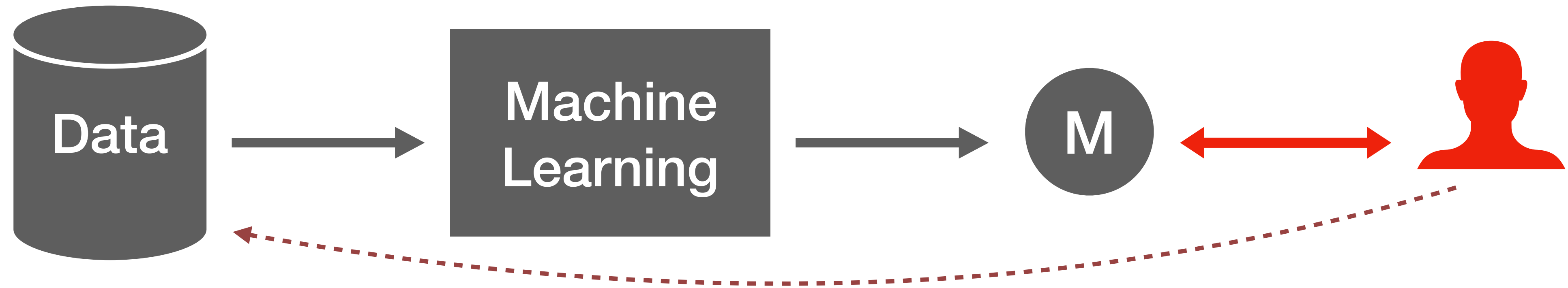


Evaluating Differentially Private Machine Learning in Practice

**Bargav Jayaraman and David Evans
Department of Computer Science
University of Virginia**

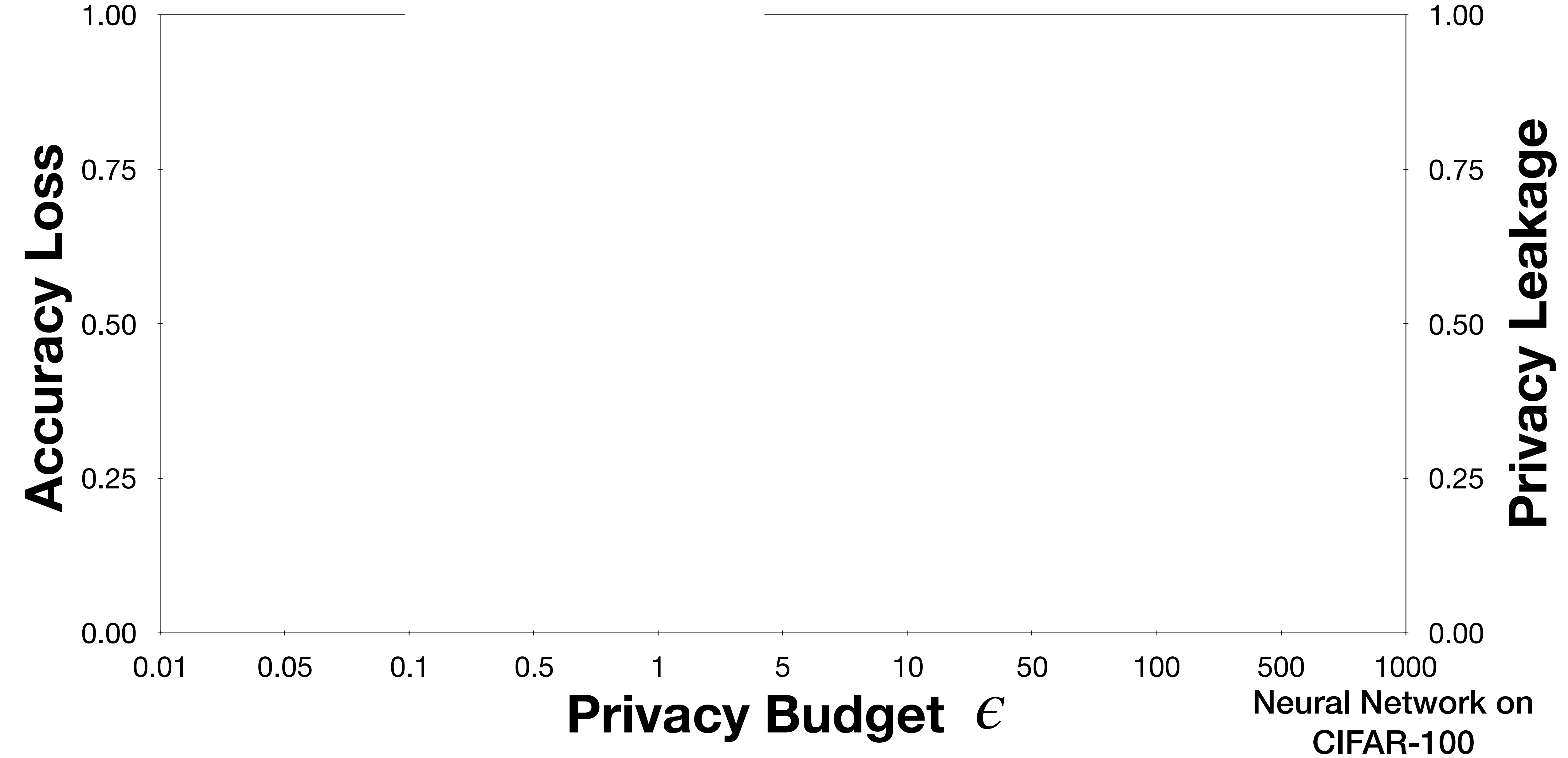
Our Objective



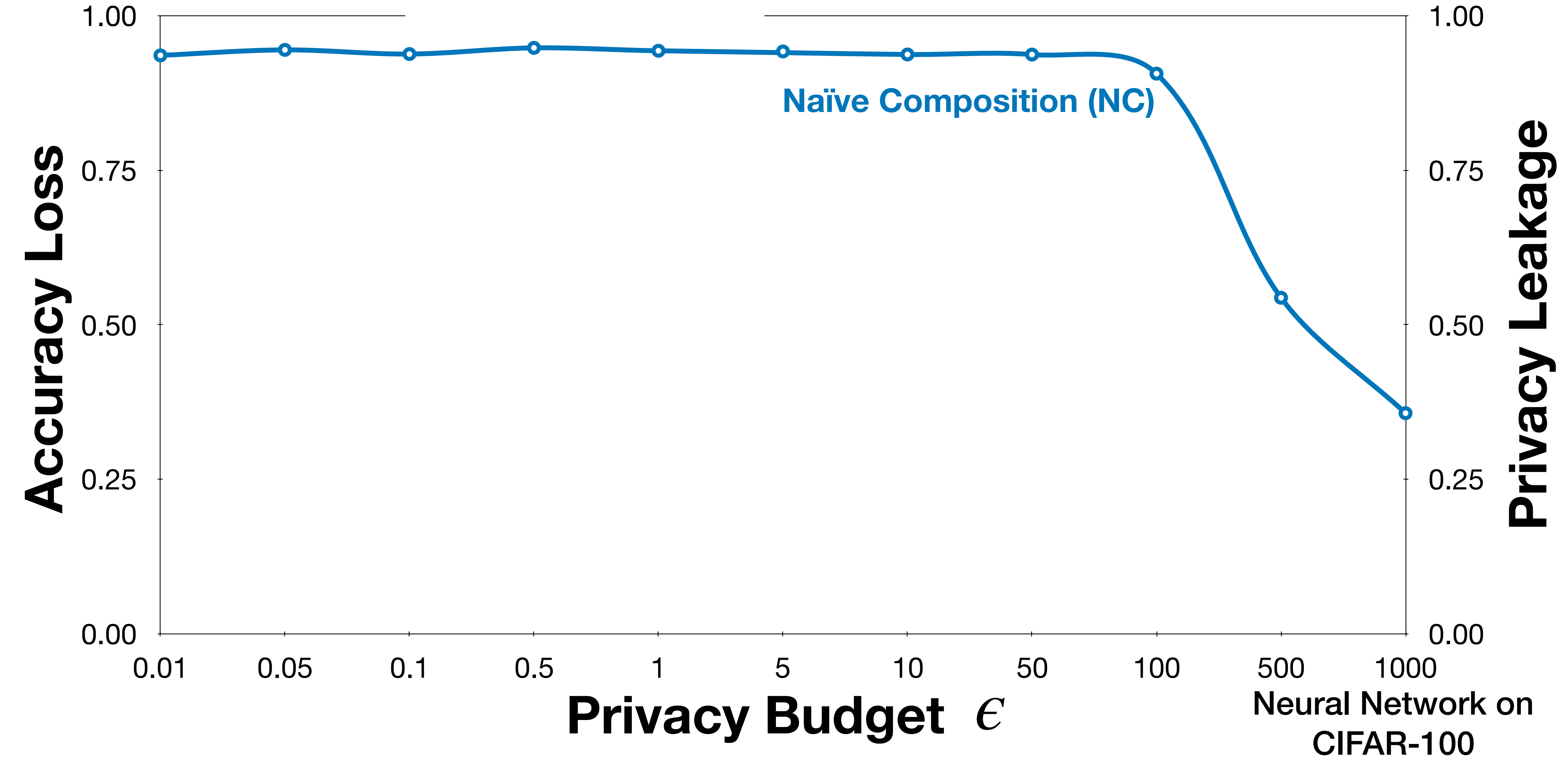
To evaluate the privacy leakage of private mechanisms

Leakage is quantified in terms of inference attacks

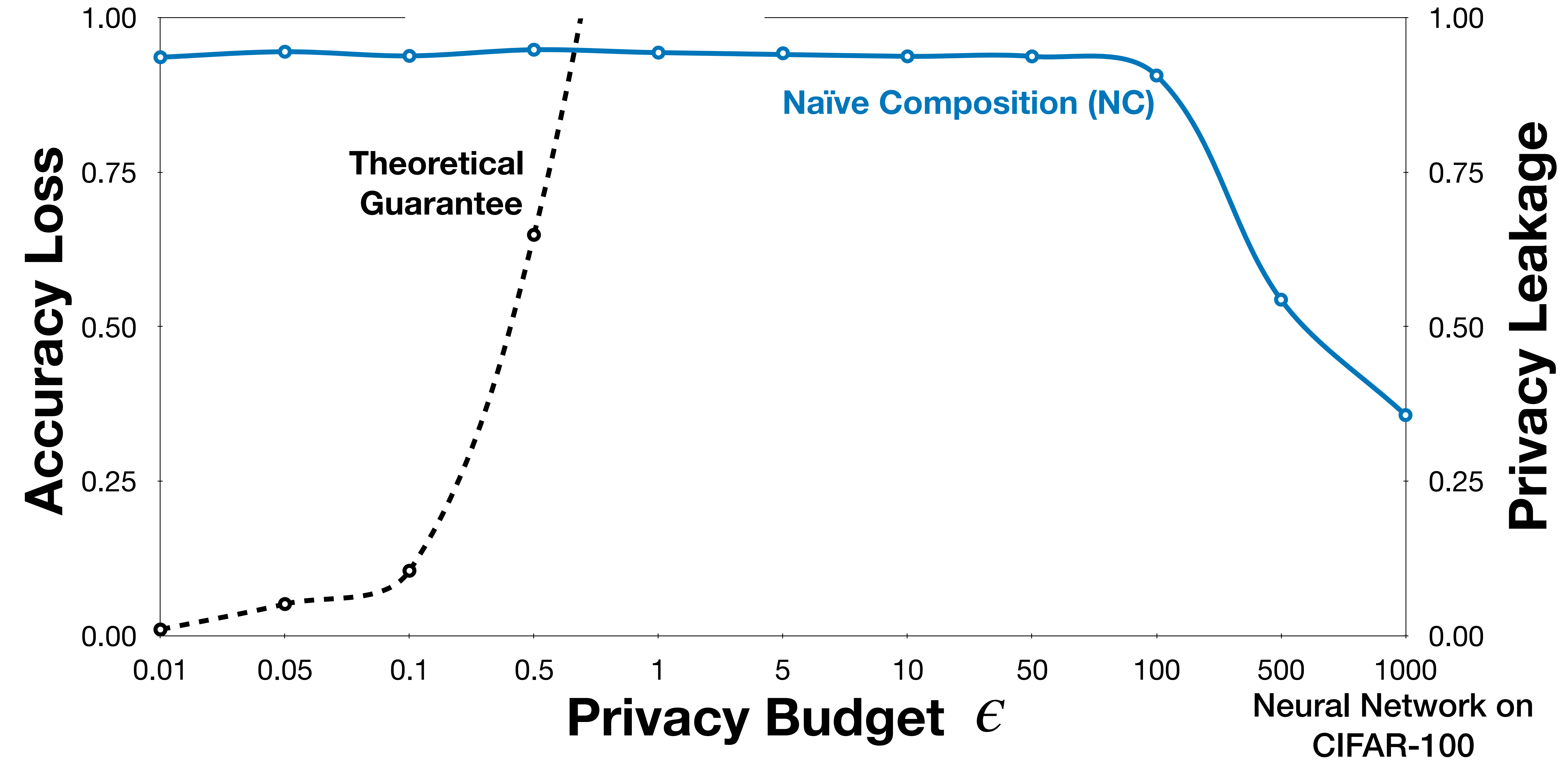
Result Highlights



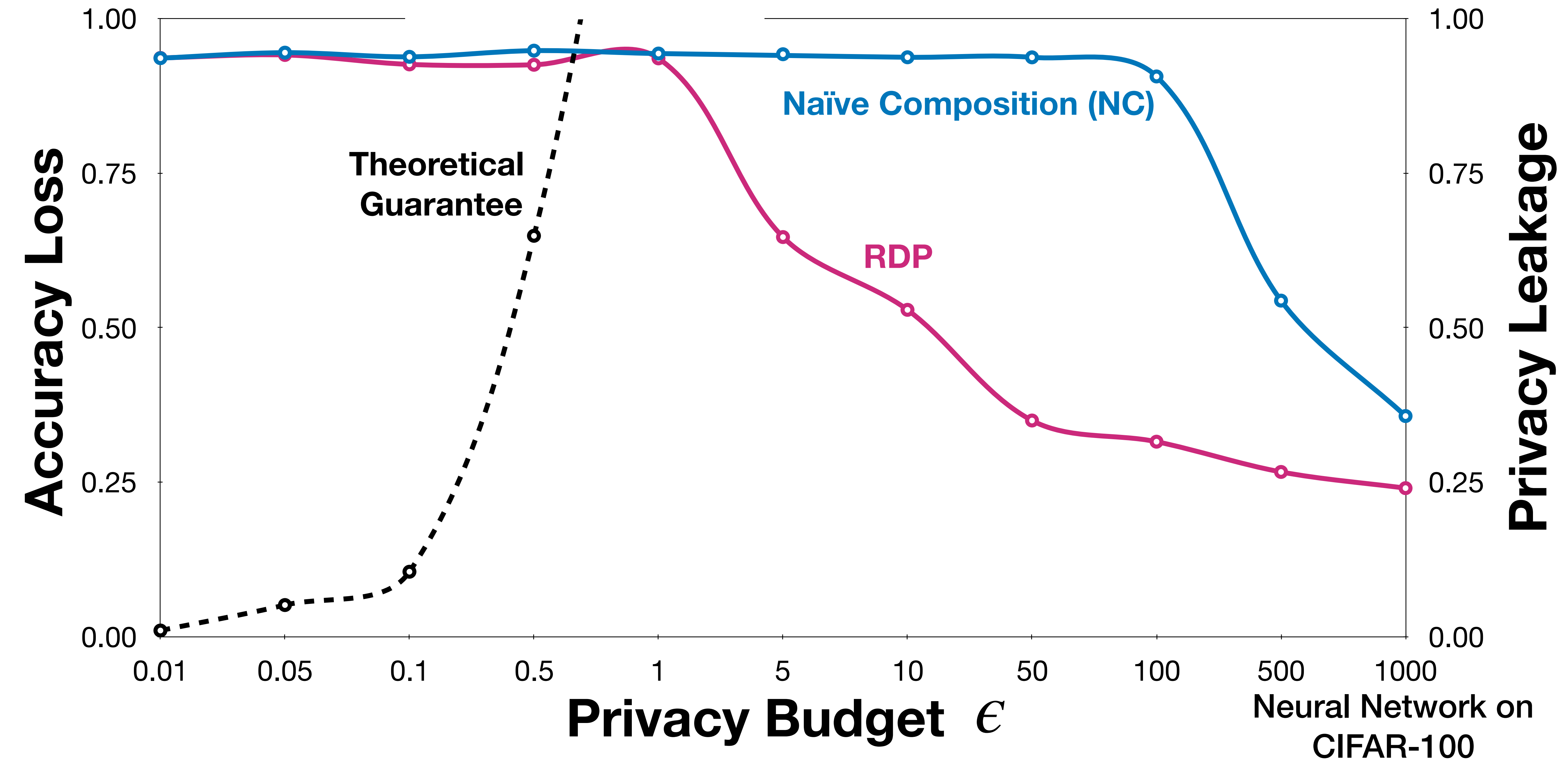
Result Highlights



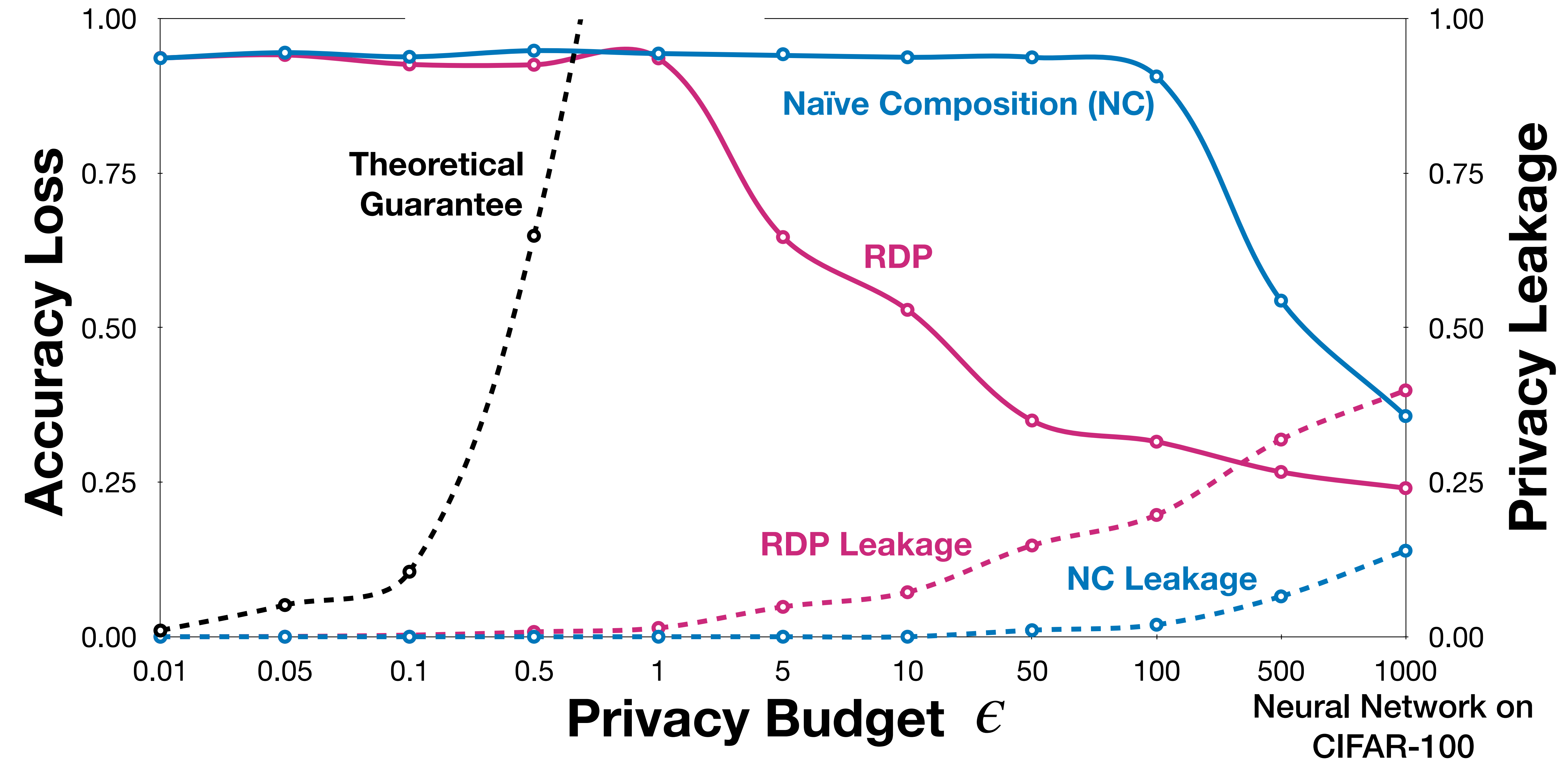
Result Highlights



Result Highlights



Result Highlights

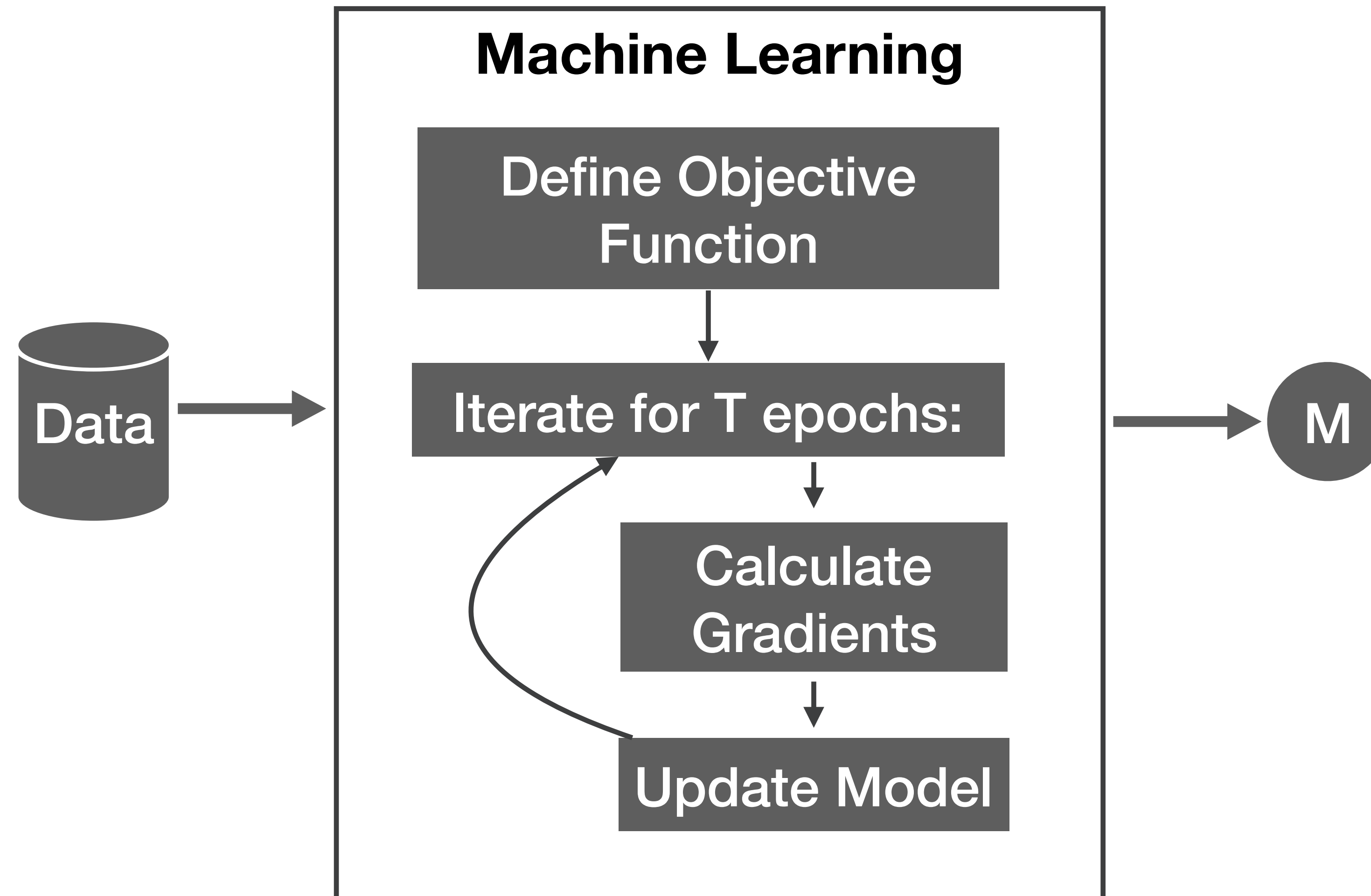


Rest of the Talk

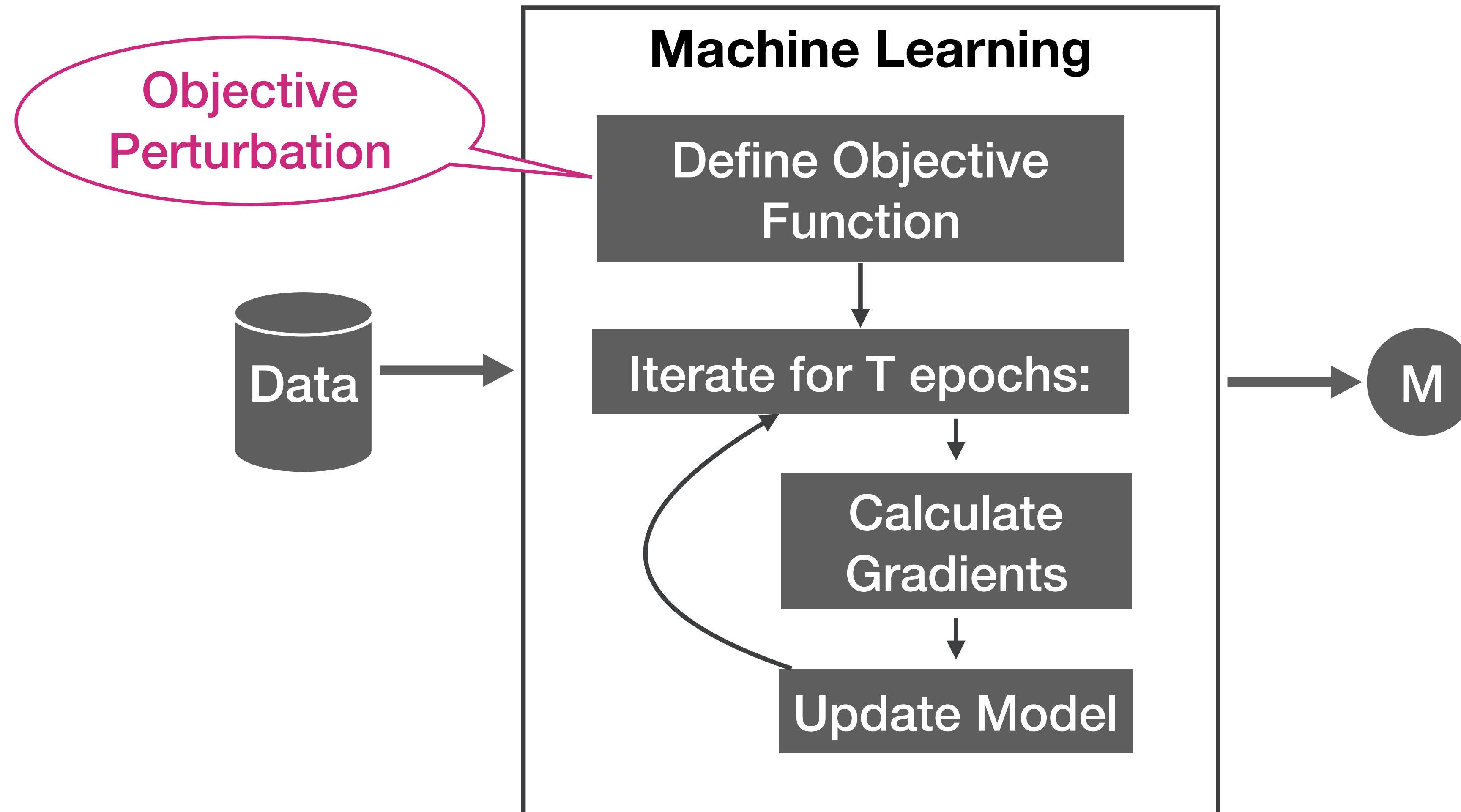
1. Background on Applying Differential Privacy to Machine Learning

2. Experimental Evaluation of Differentially Private Machine Learning Implementations

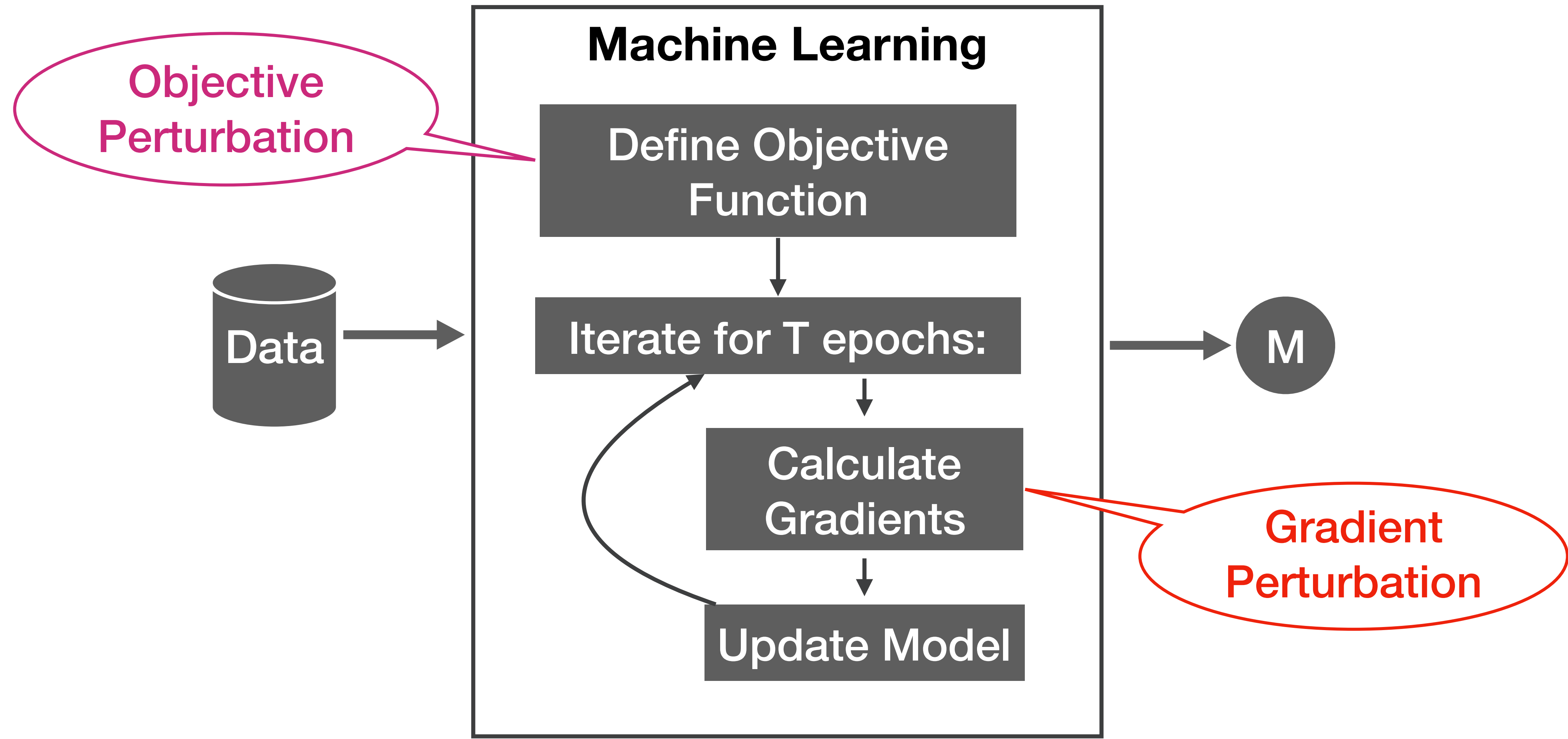
Applying DP to Machine Learning



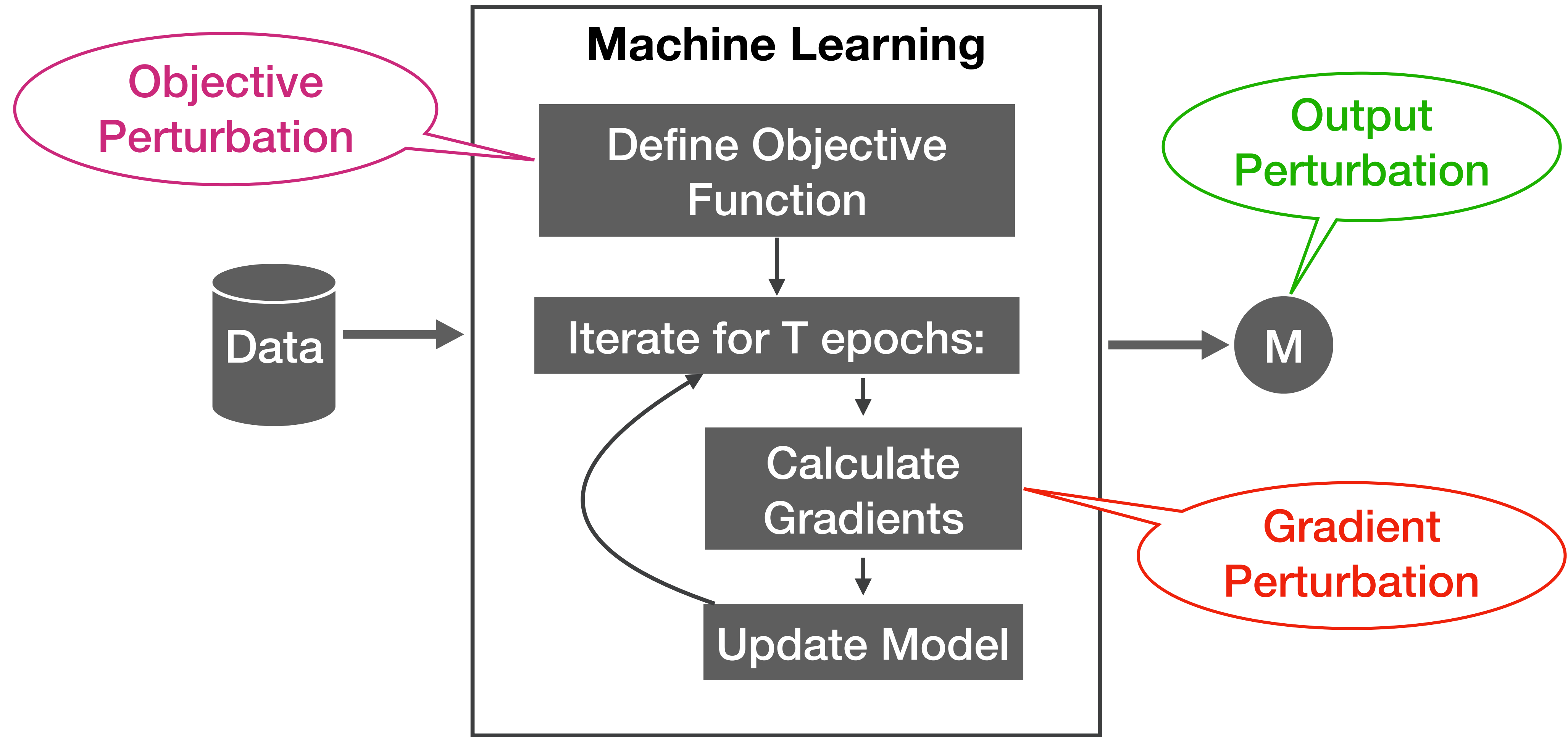
Applying DP to Machine Learning



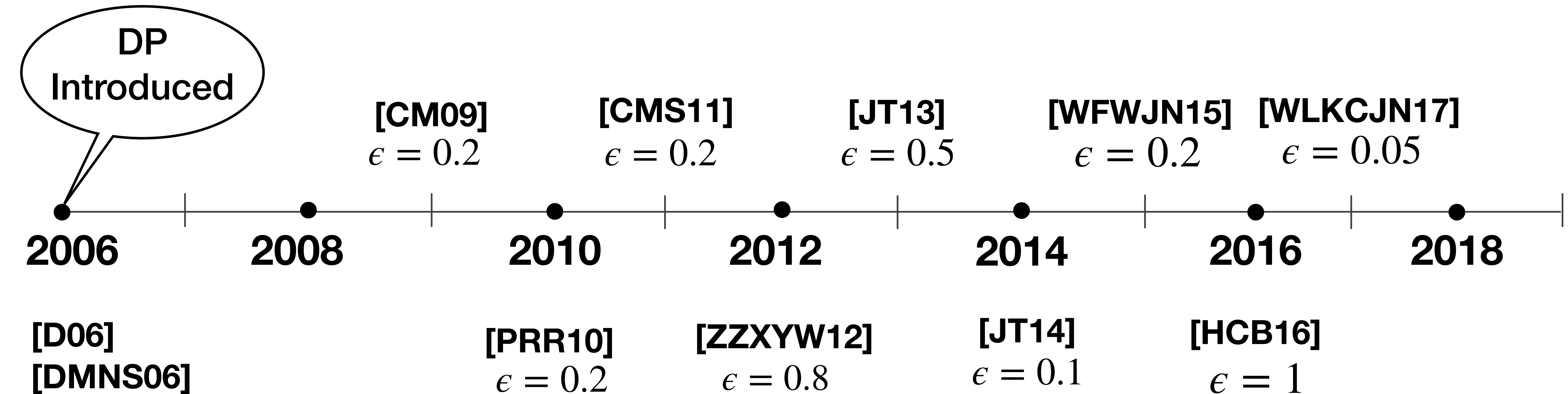
Applying DP to Machine Learning



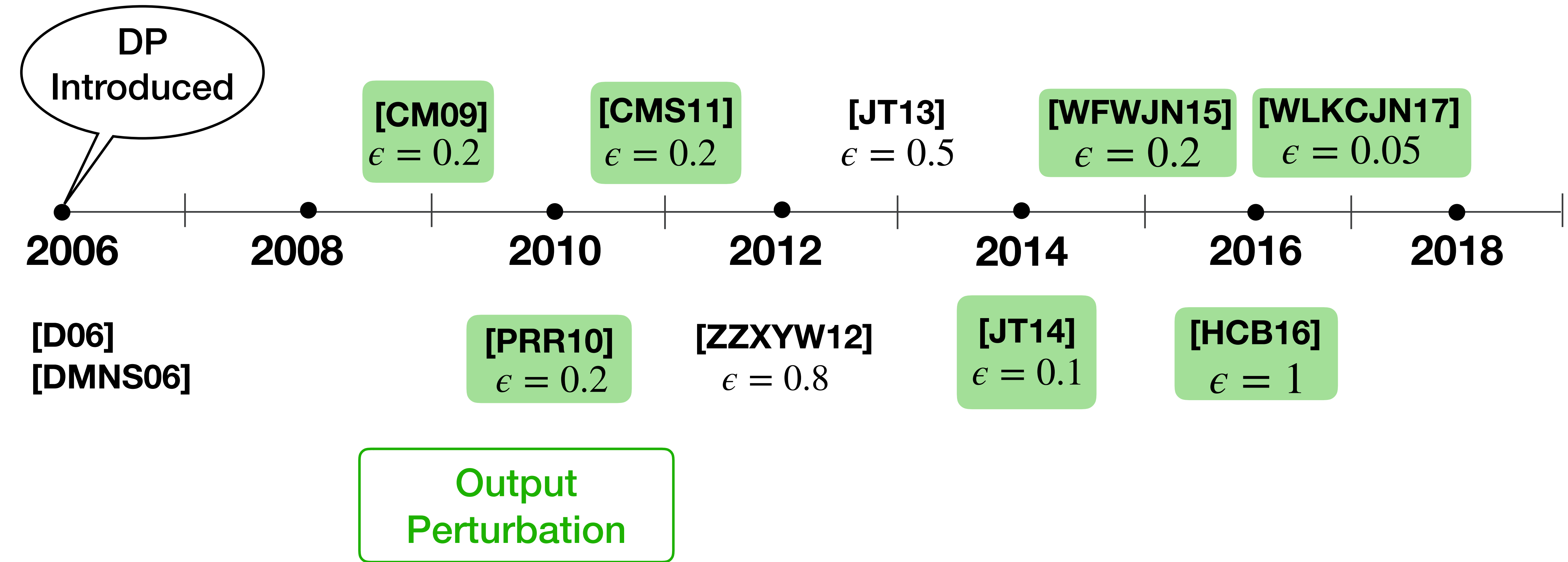
Applying DP to Machine Learning



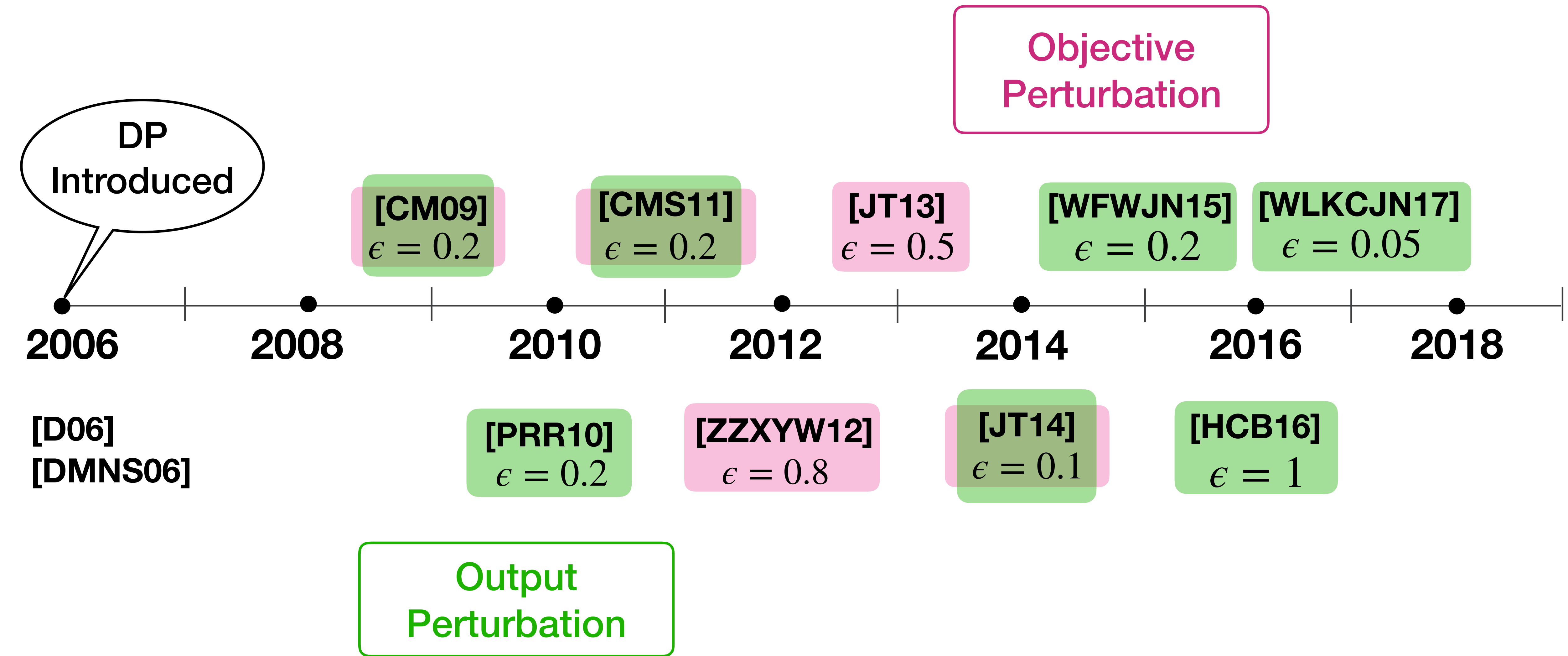
ERM Algorithms using $\epsilon \leq 1$



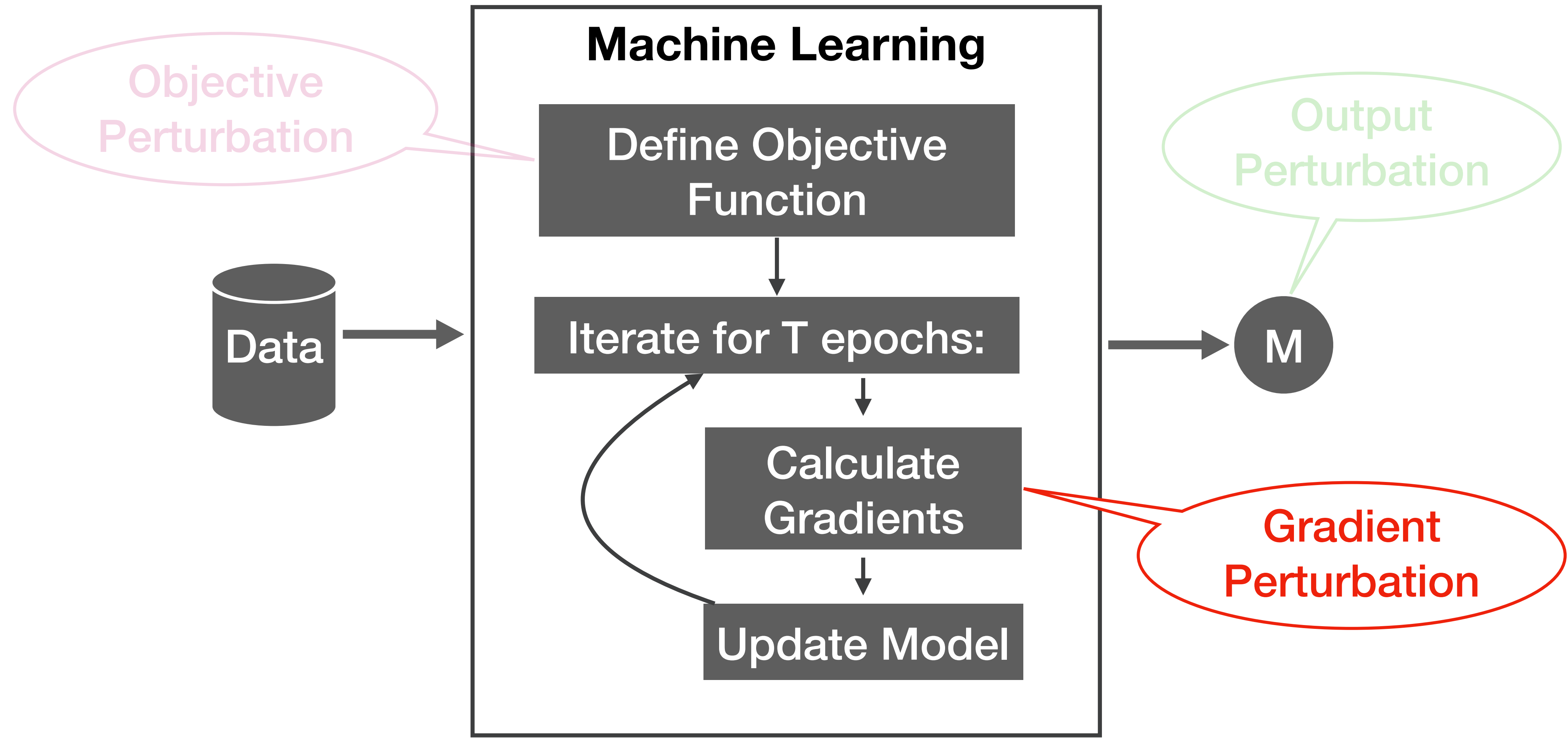
ERM Algorithms using $\epsilon \leq 1$



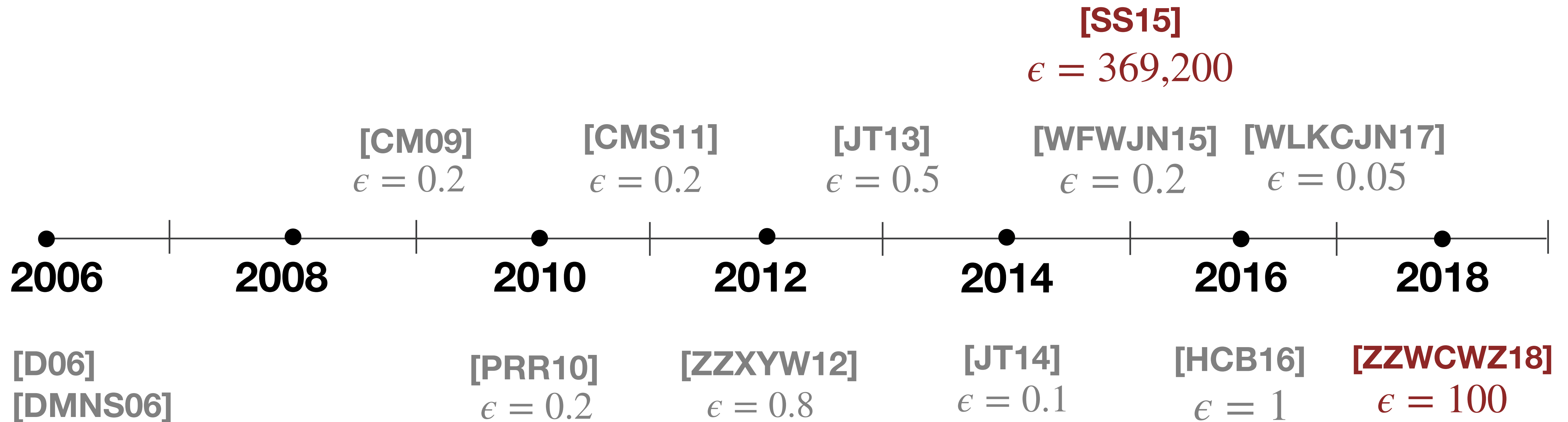
ERM Algorithms using $\epsilon \leq 1$



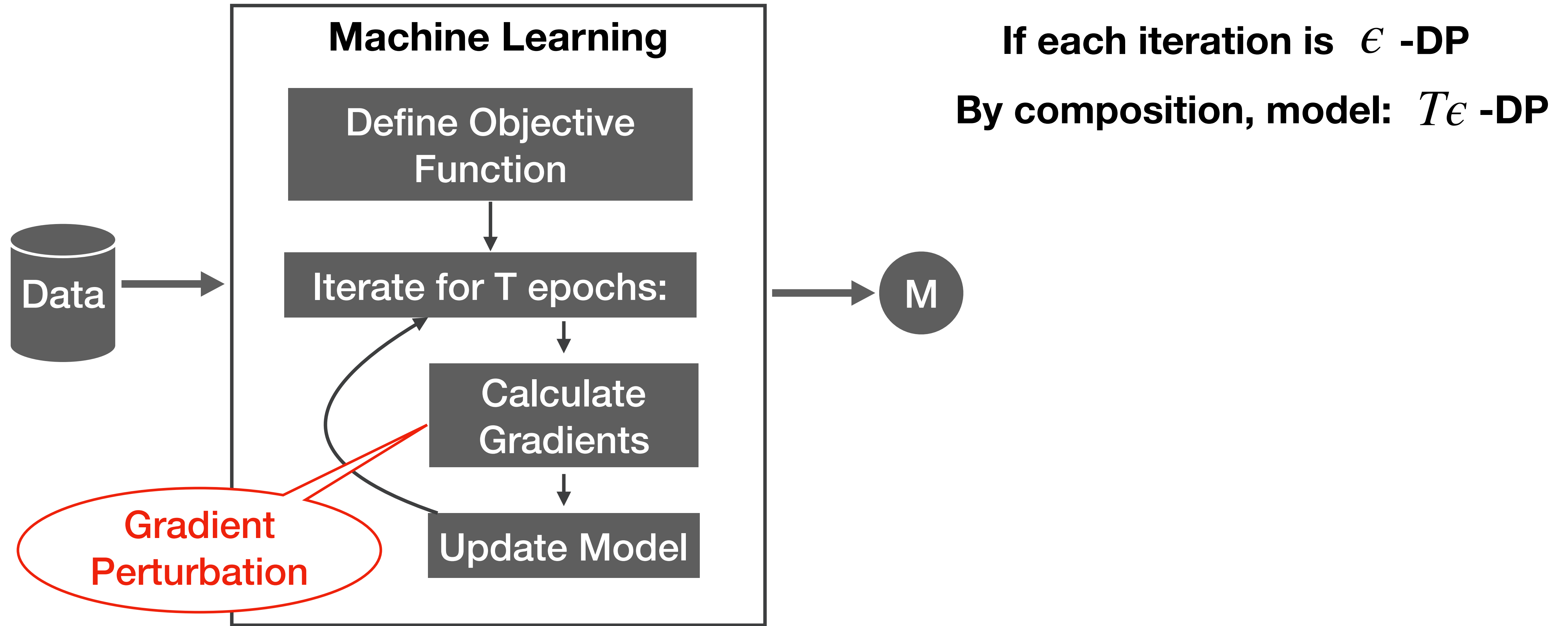
Applying DP to Deep Learning



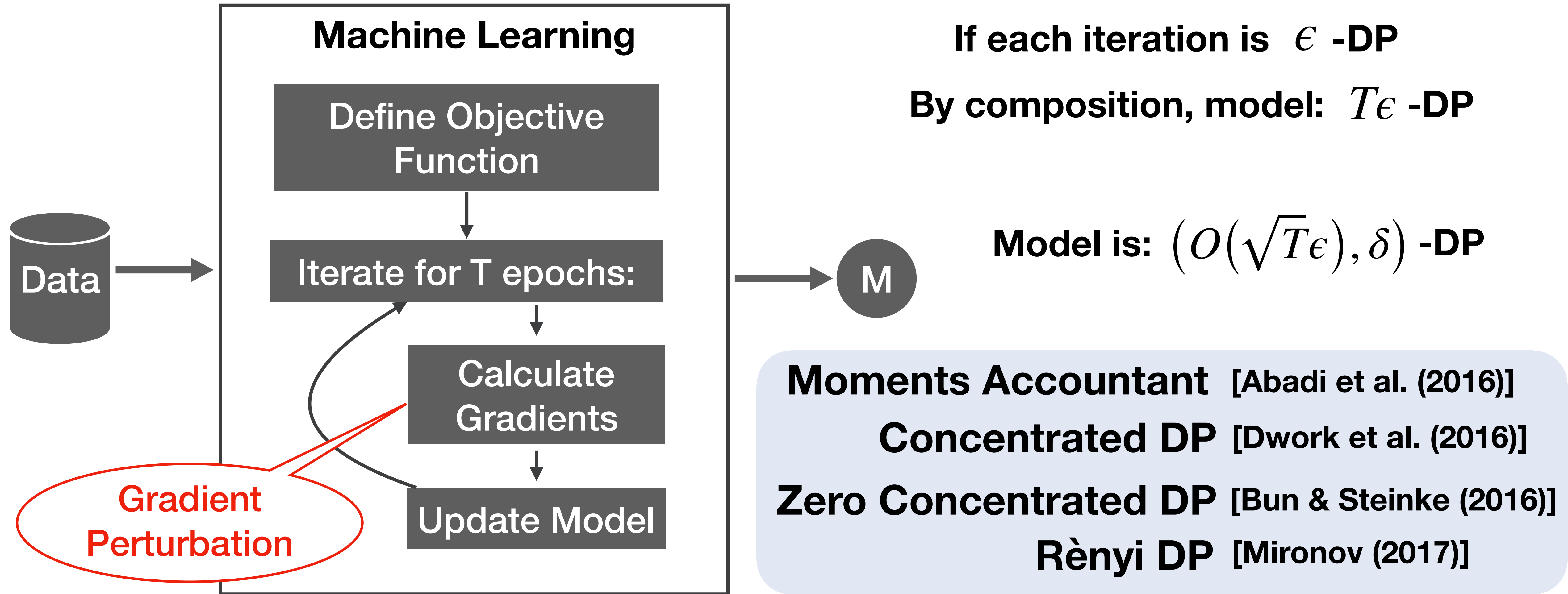
Deep Learning requiring high ϵ value



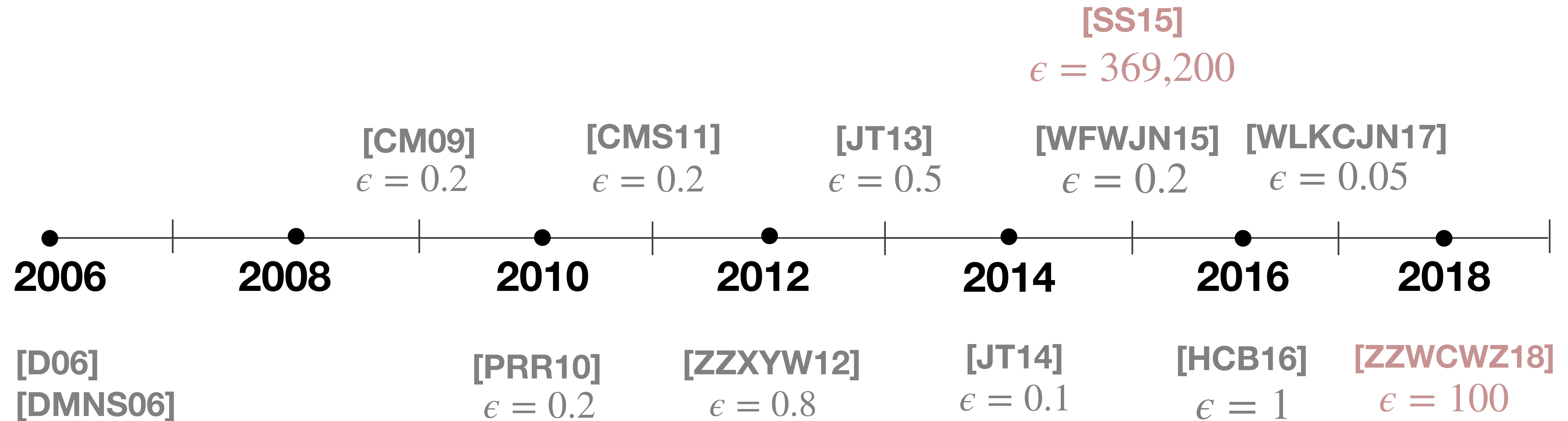
Improving Composition



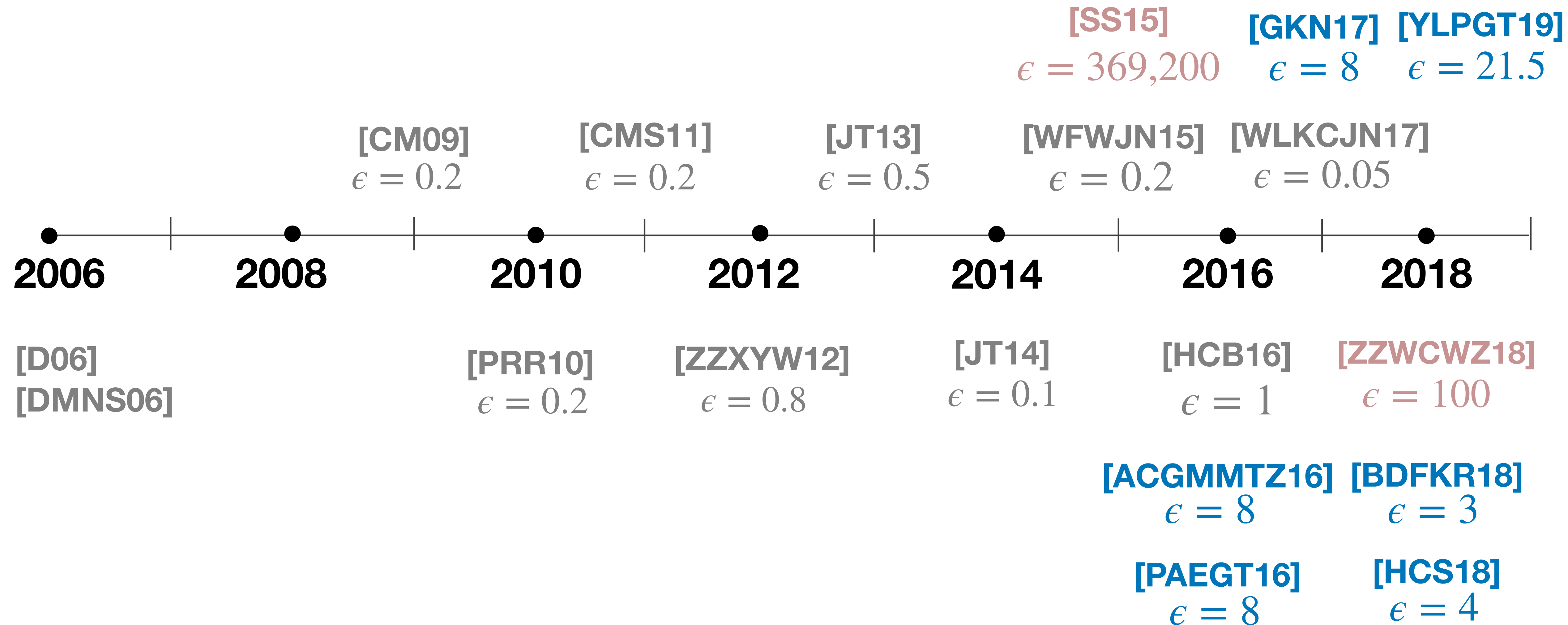
Improving Composition



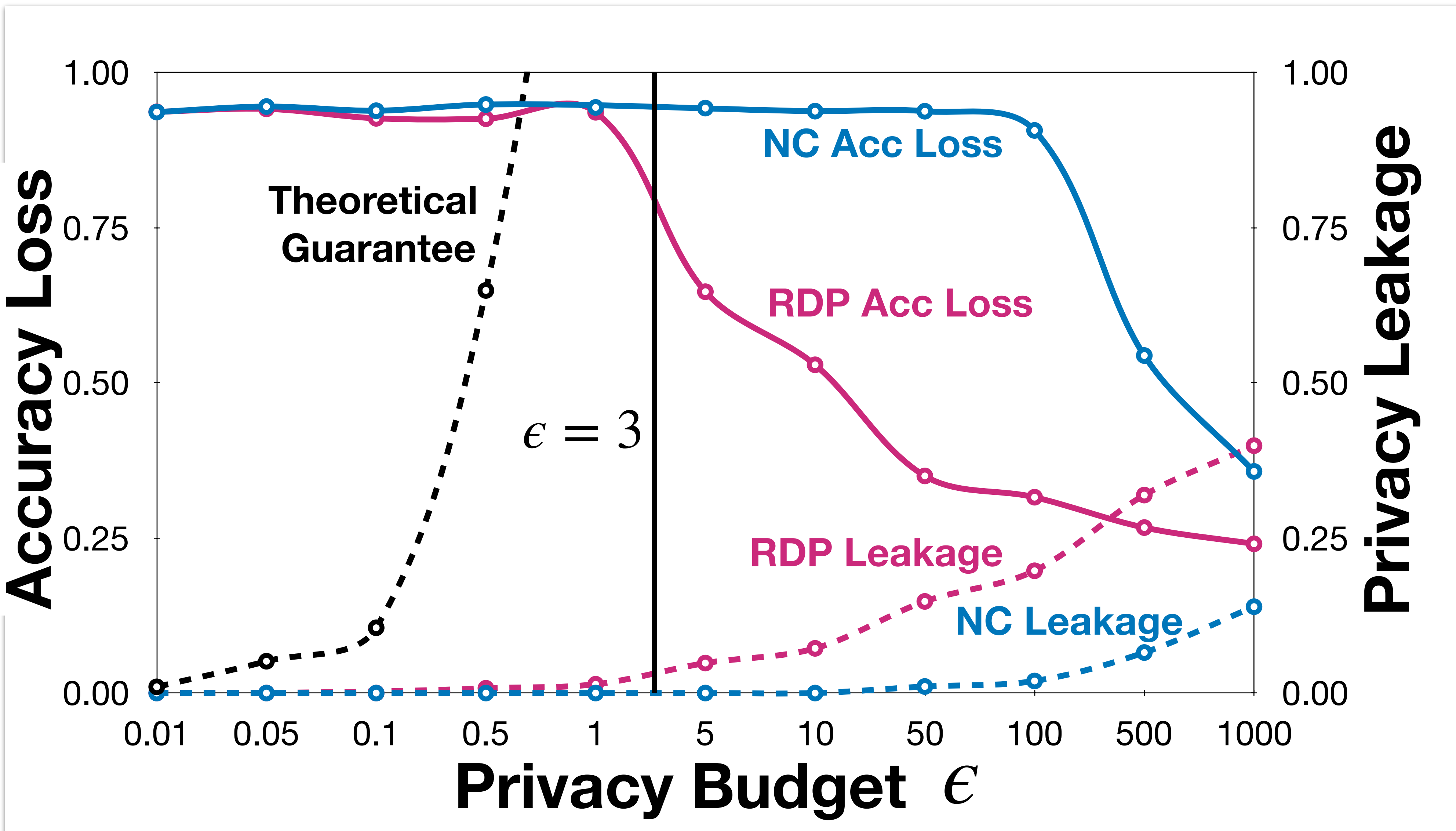
Lower ϵ value with recent DP notions



Lower ϵ value with recent DP notions



Lower ϵ value with recent DP notions



[GKN17] $\epsilon = 8$ [YLPGT19] $\epsilon = 21.5$

[WLKCJN17] $\epsilon = 0.05$

16 2018

[B16] $\epsilon = 1$ [ZZWCWZ18] $\epsilon = 100$

[MTZ16] $\epsilon = 8$ [BDFKR18] $\epsilon = 3$

[T16] $\epsilon = 8$ [HCS18] $\epsilon = 4$

Experiments

Model

Logistic Regression

Neural Network

Task

**100 class classification
on CIFAR-100**

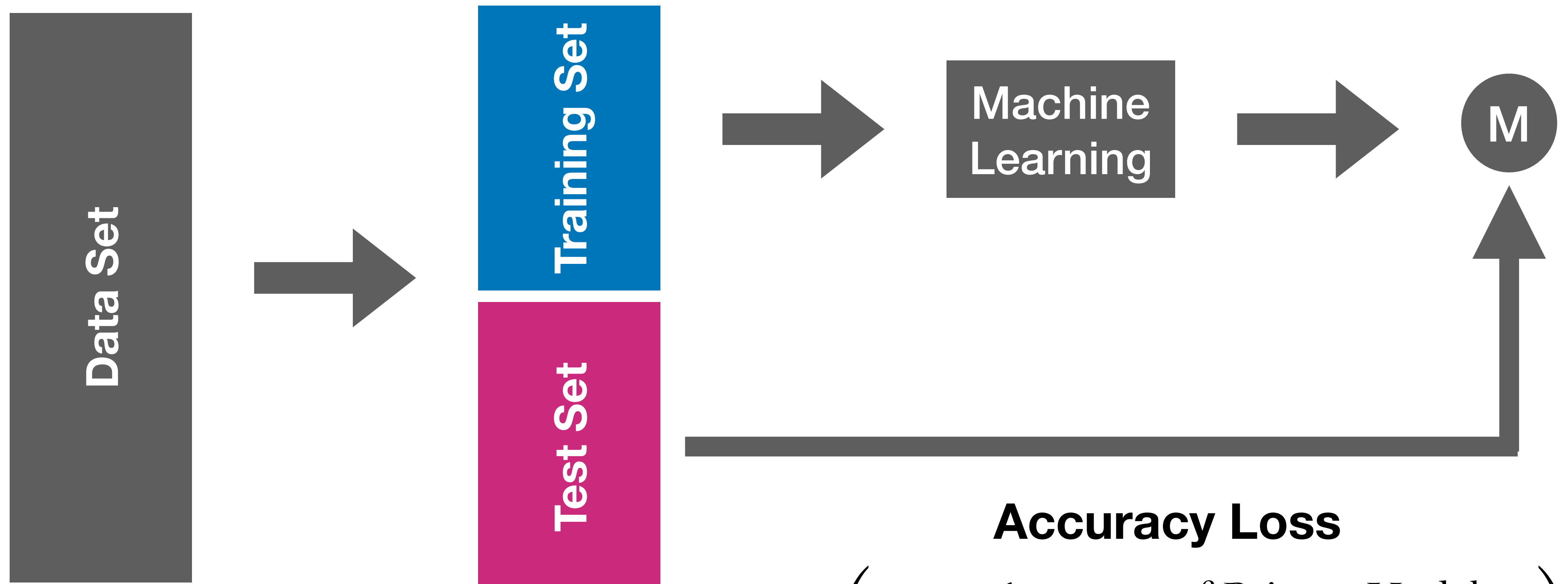
**100 class classification
on Purchase-100**

Evaluation Metric

Accuracy Loss

Privacy Leakage

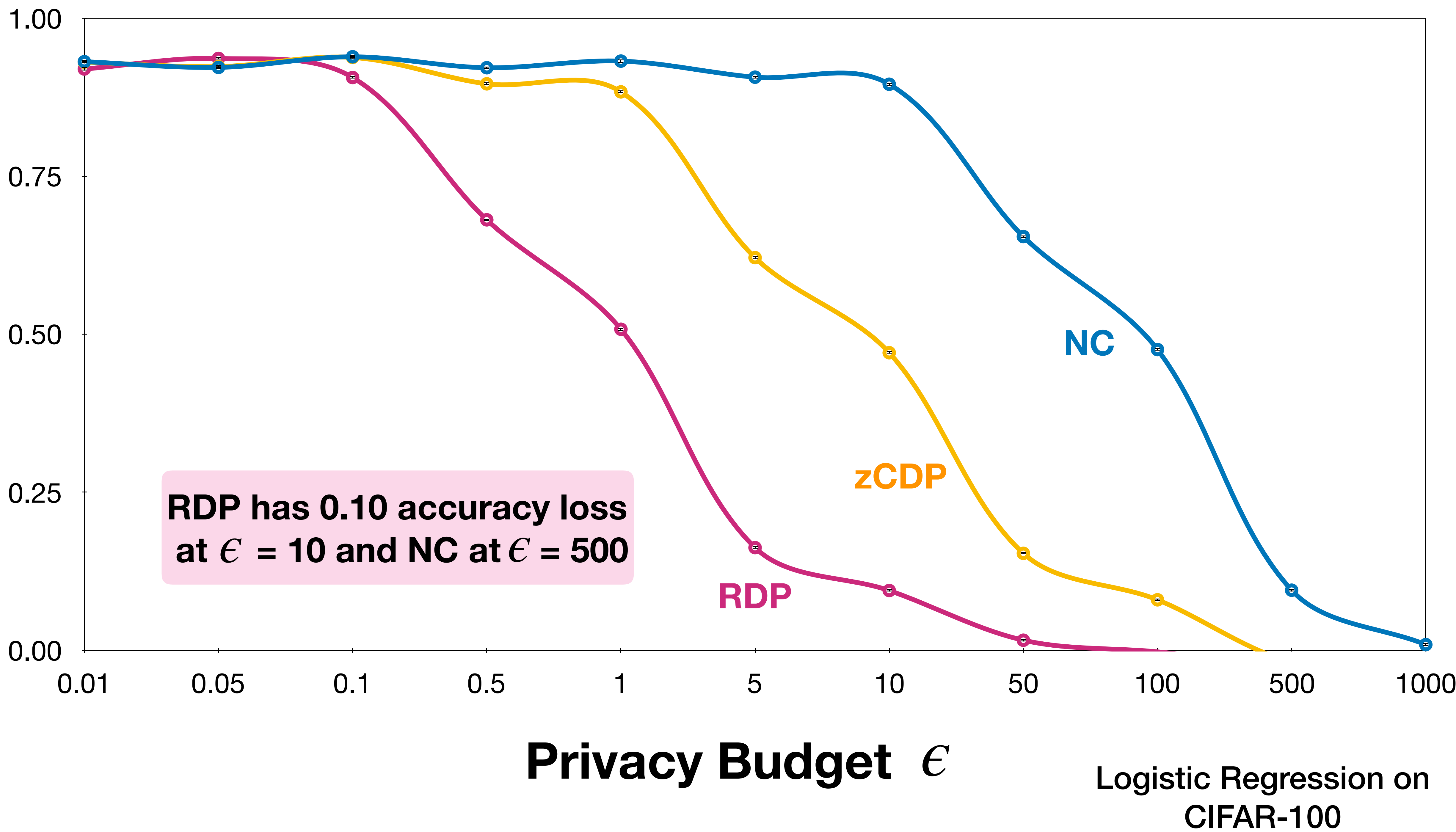
Training and Testing



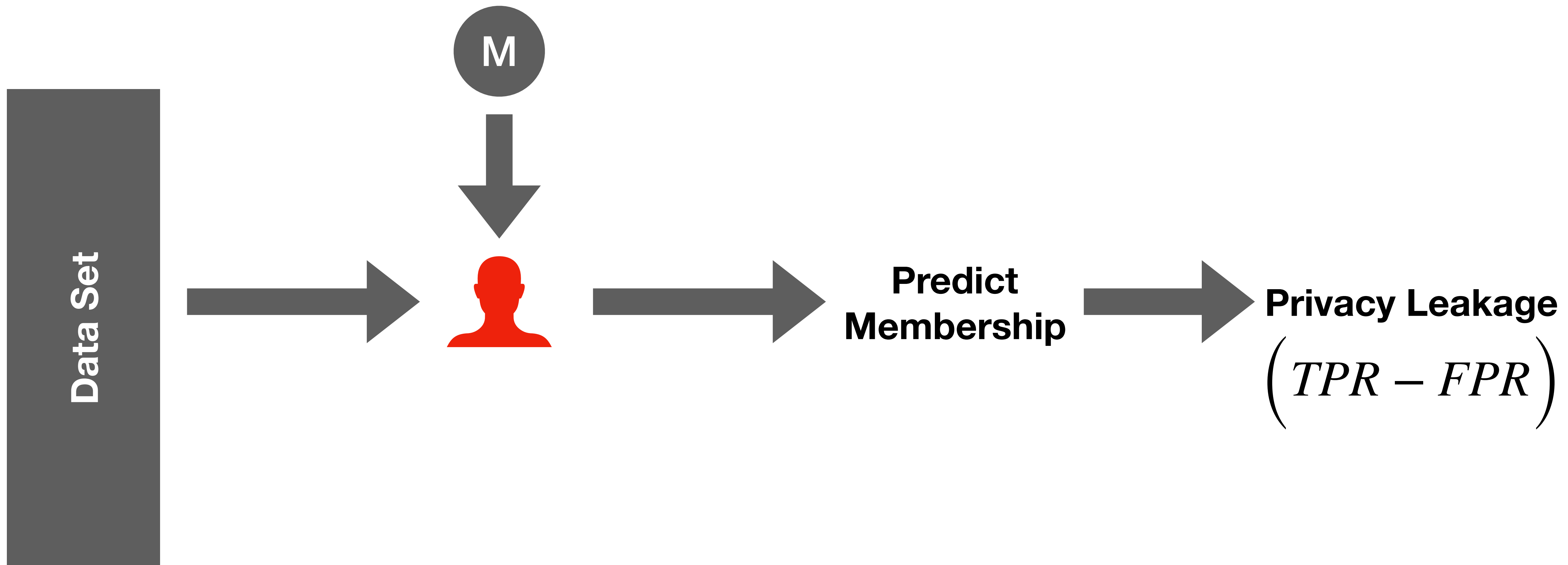
Accuracy Loss

$$\left(1 - \frac{\text{Accuracy of Private Model}}{\text{Accuracy of Non-Private Model}} \right)$$

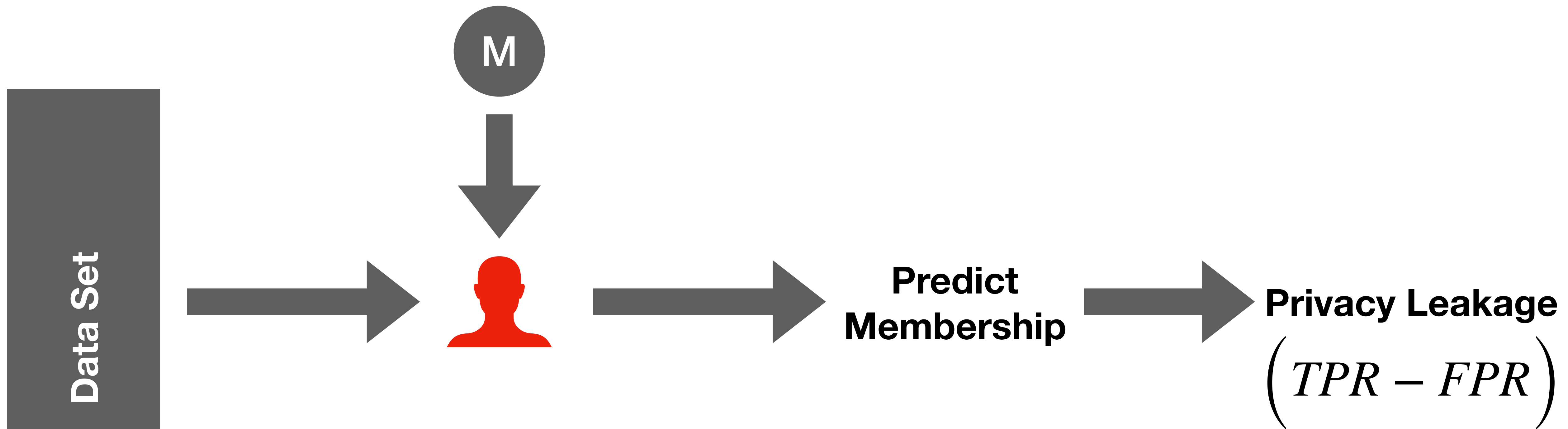
Accuracy Loss



Membership Inference Attacks



Membership Inference Attacks



Shokri et al. [1]

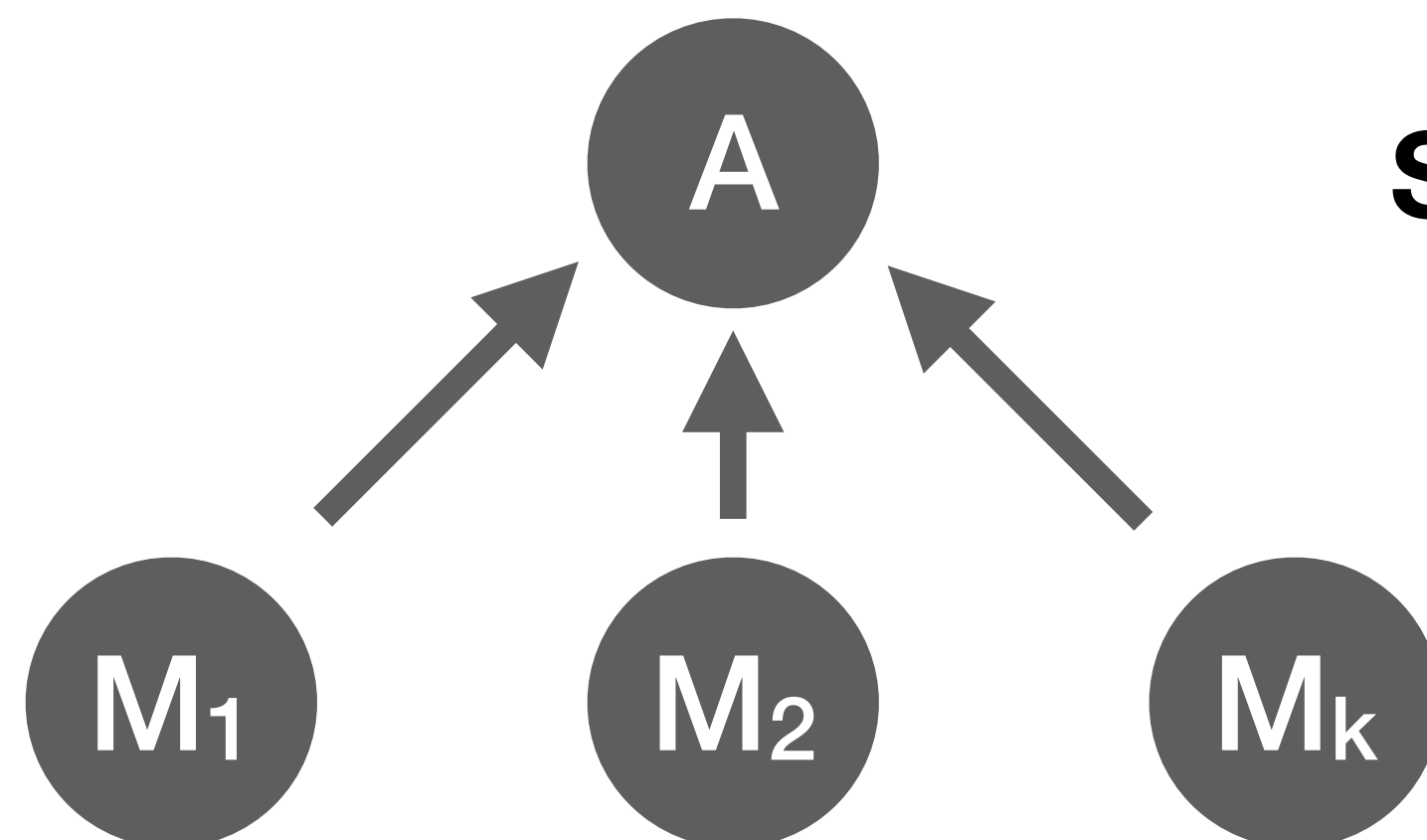
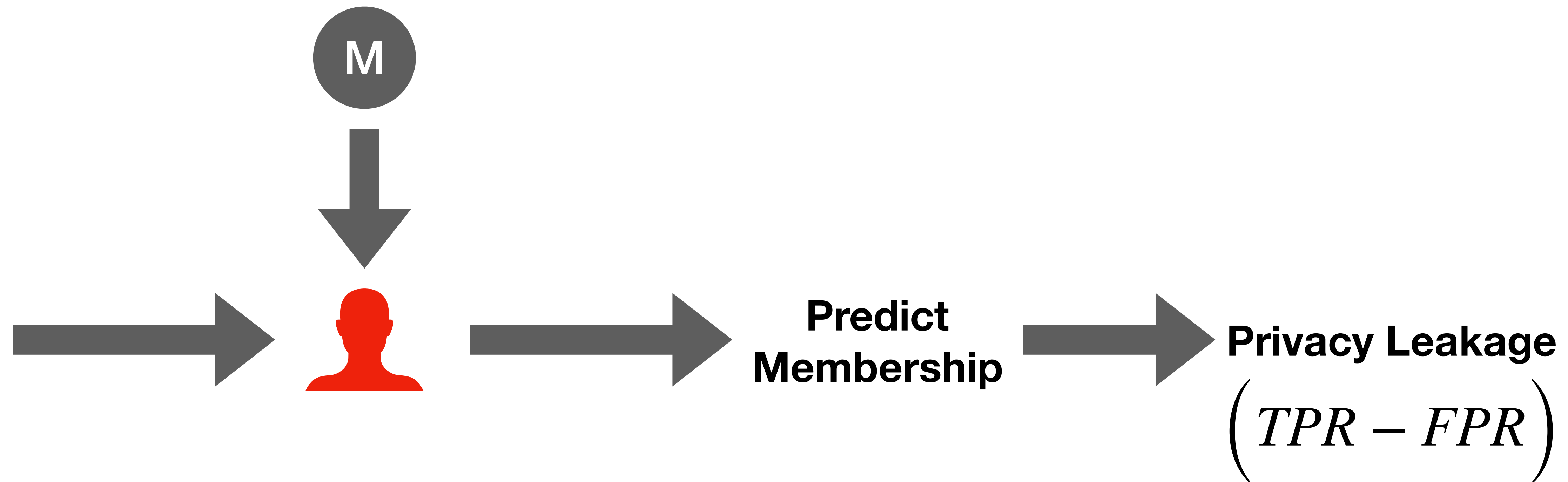
M_1

M_2

M_k

1. Reza Shokri, Marco Stronati, Congzheng Song and Vitaly Shmatikov
Membership Inference Attacks Against Machine Learning Models, S&P 2017

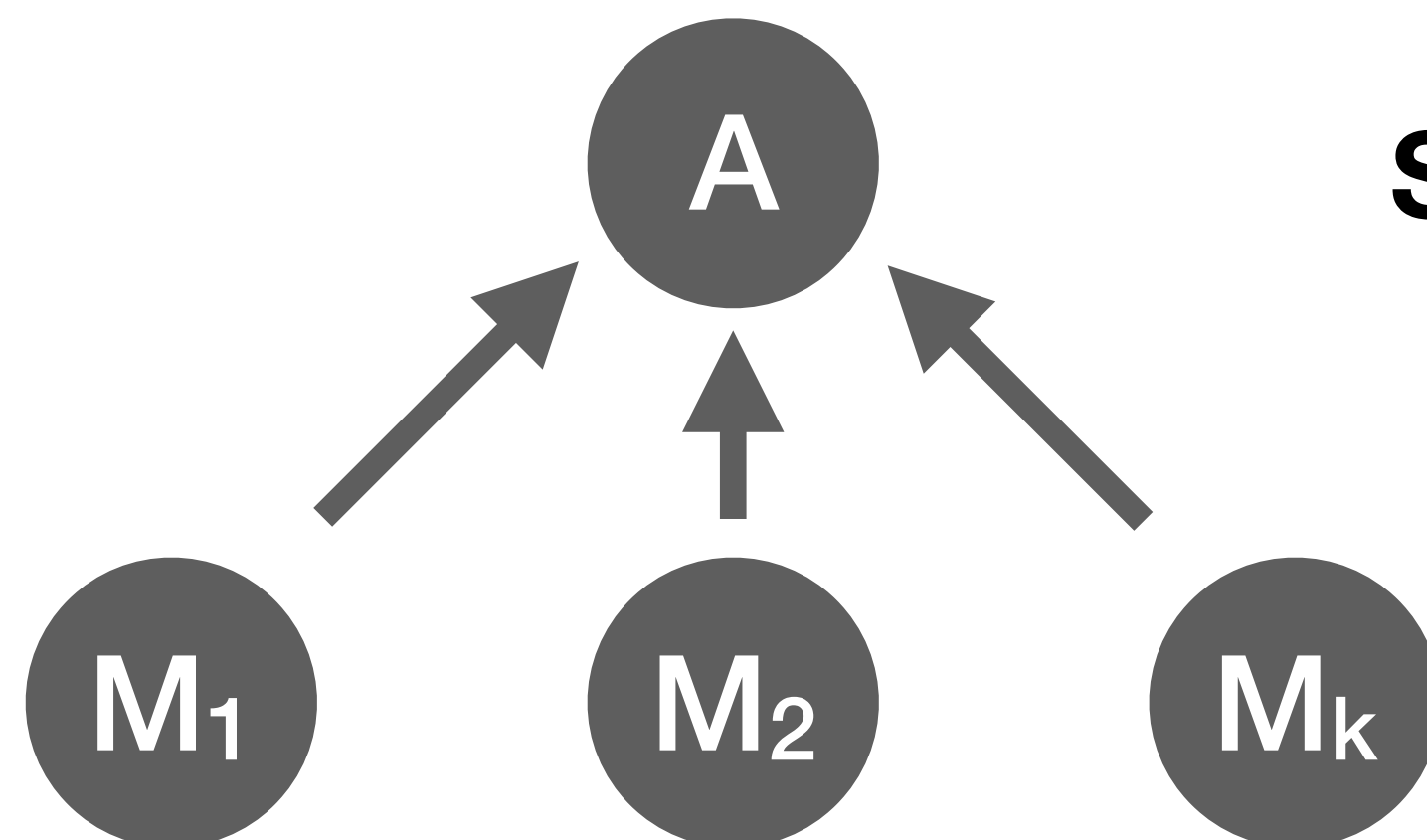
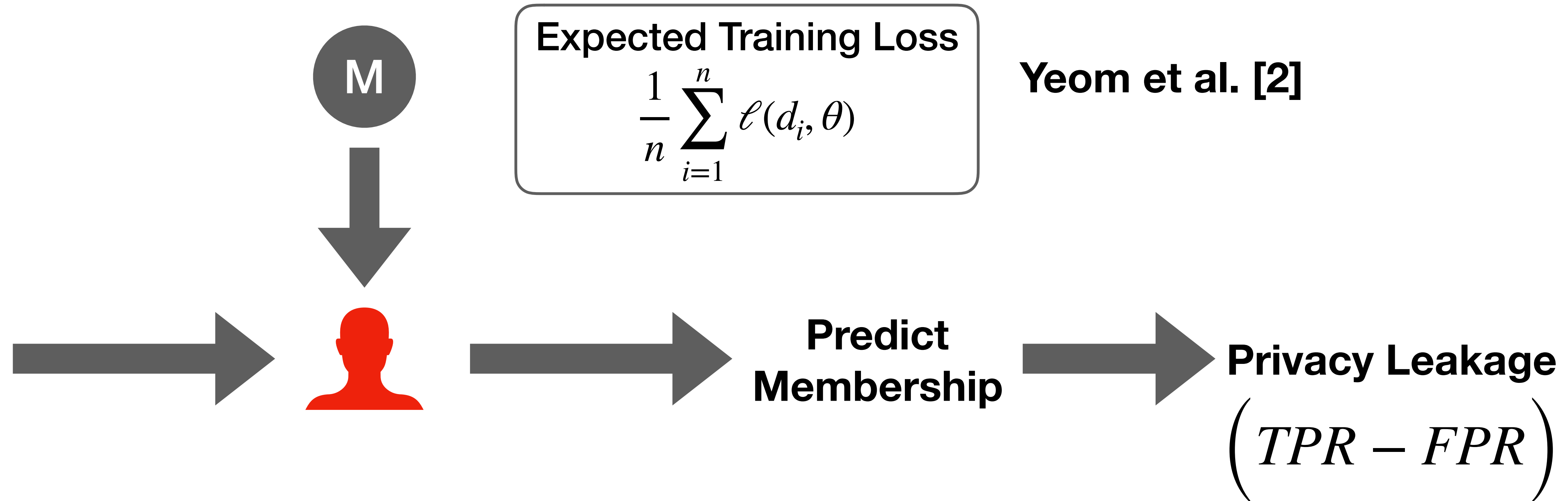
Membership Inference Attacks



Shokri et al. [1]

1. Reza Shokri, Marco Stronati, Congzheng Song and Vitaly Shmatikov
Membership Inference Attacks Against Machine Learning Models, S&P 2017

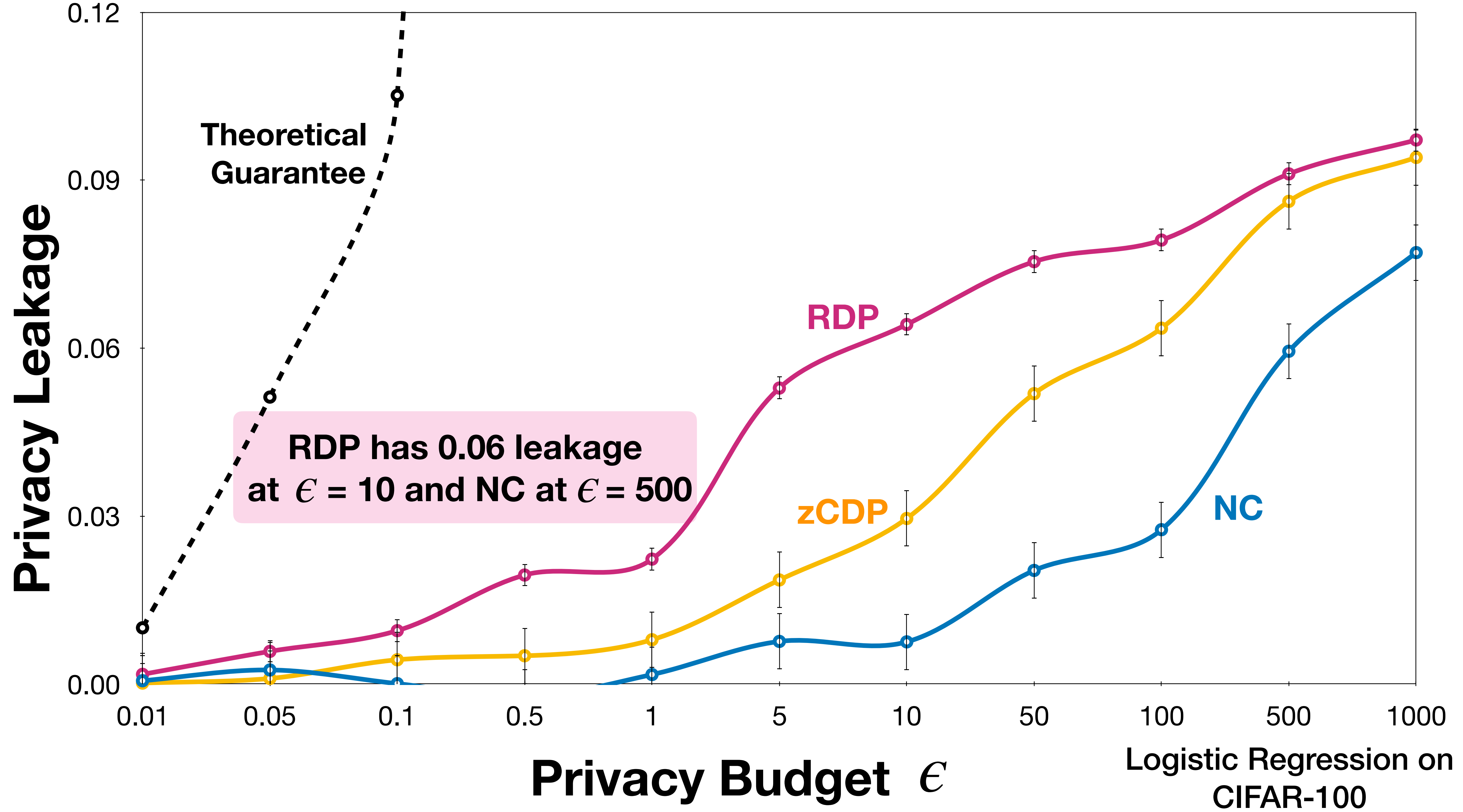
Membership Inference Attacks

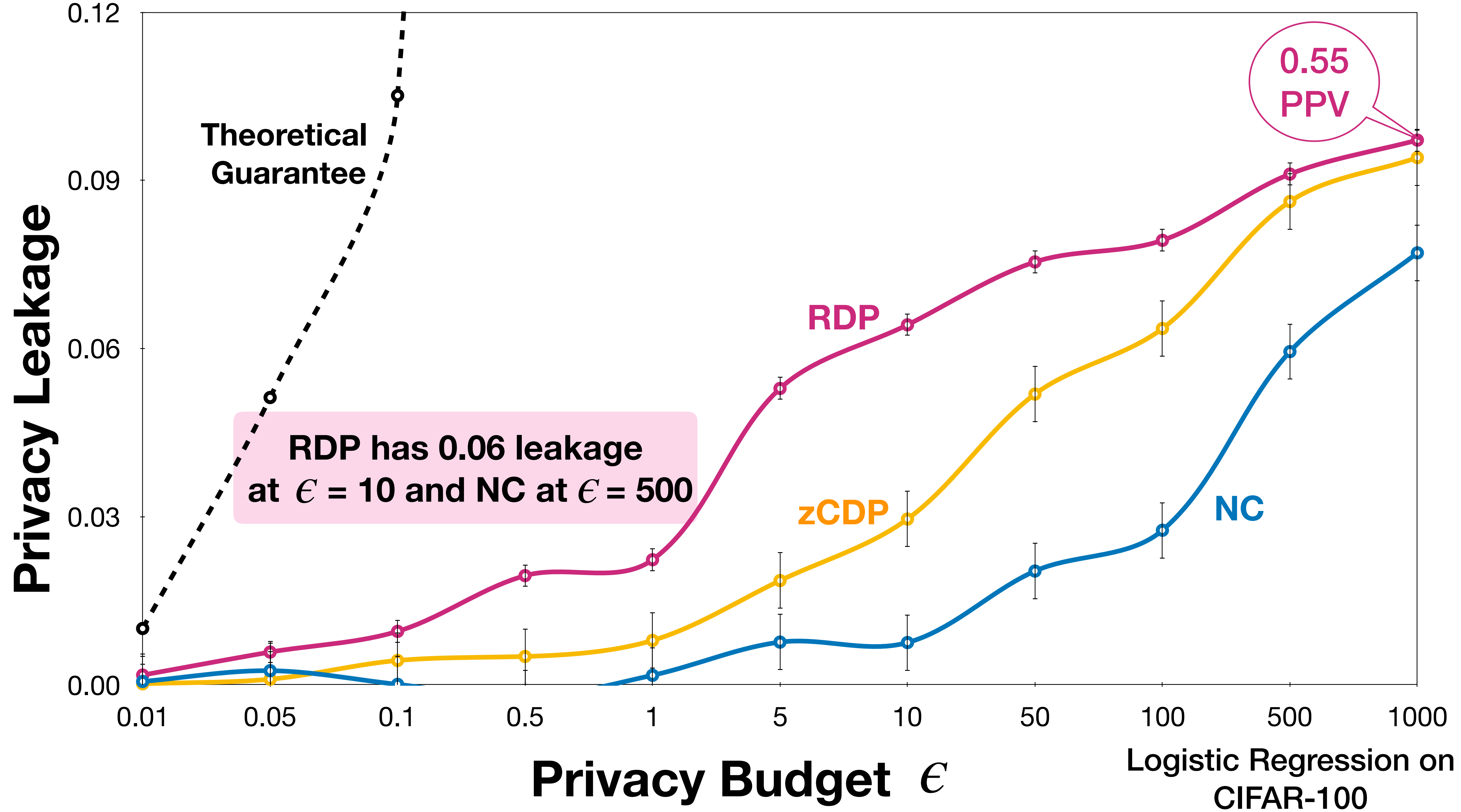


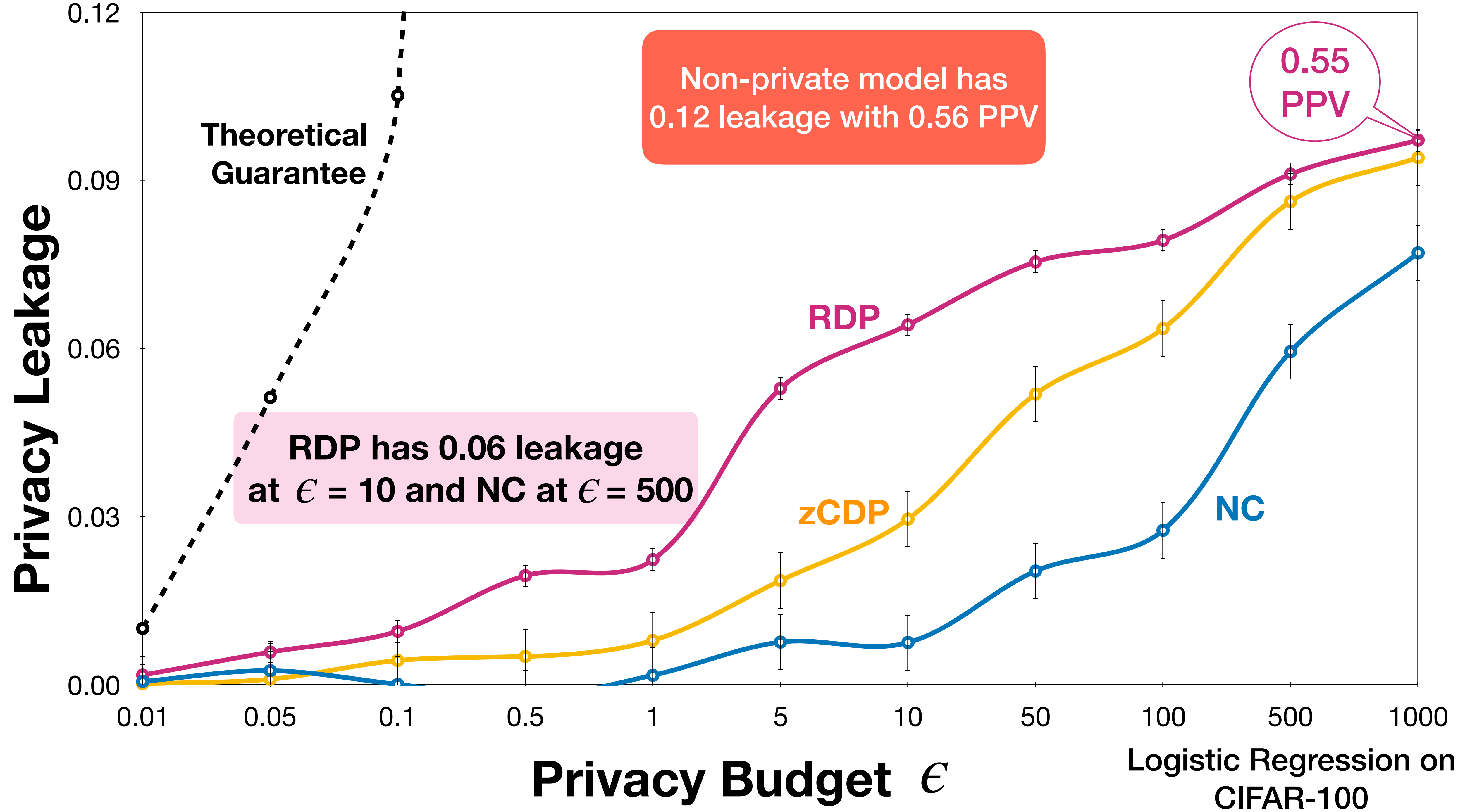
Shokri et al. [1]

1. Reza Shokri, Marco Stronati, Congzheng Song and Vitaly Shmatikov
Membership Inference Attacks Against Machine Learning Models, S&P 2017

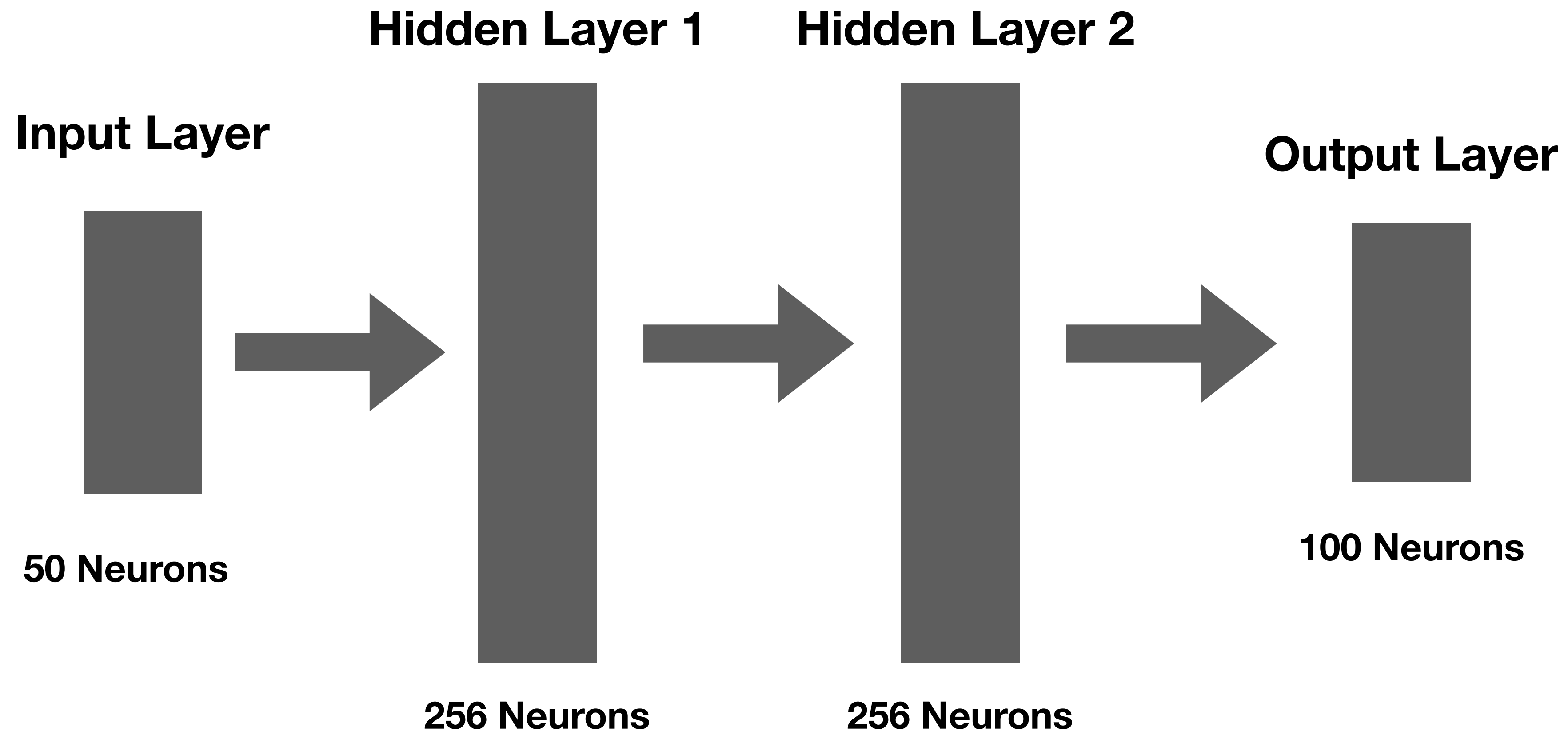
2. Samuel Yeom, Irene Giacomelli, Matt Fredrikson and Somesh Jha
Privacy Risk in Machine Learning: Analyzing the Connection to Overfitting, CSF 2018



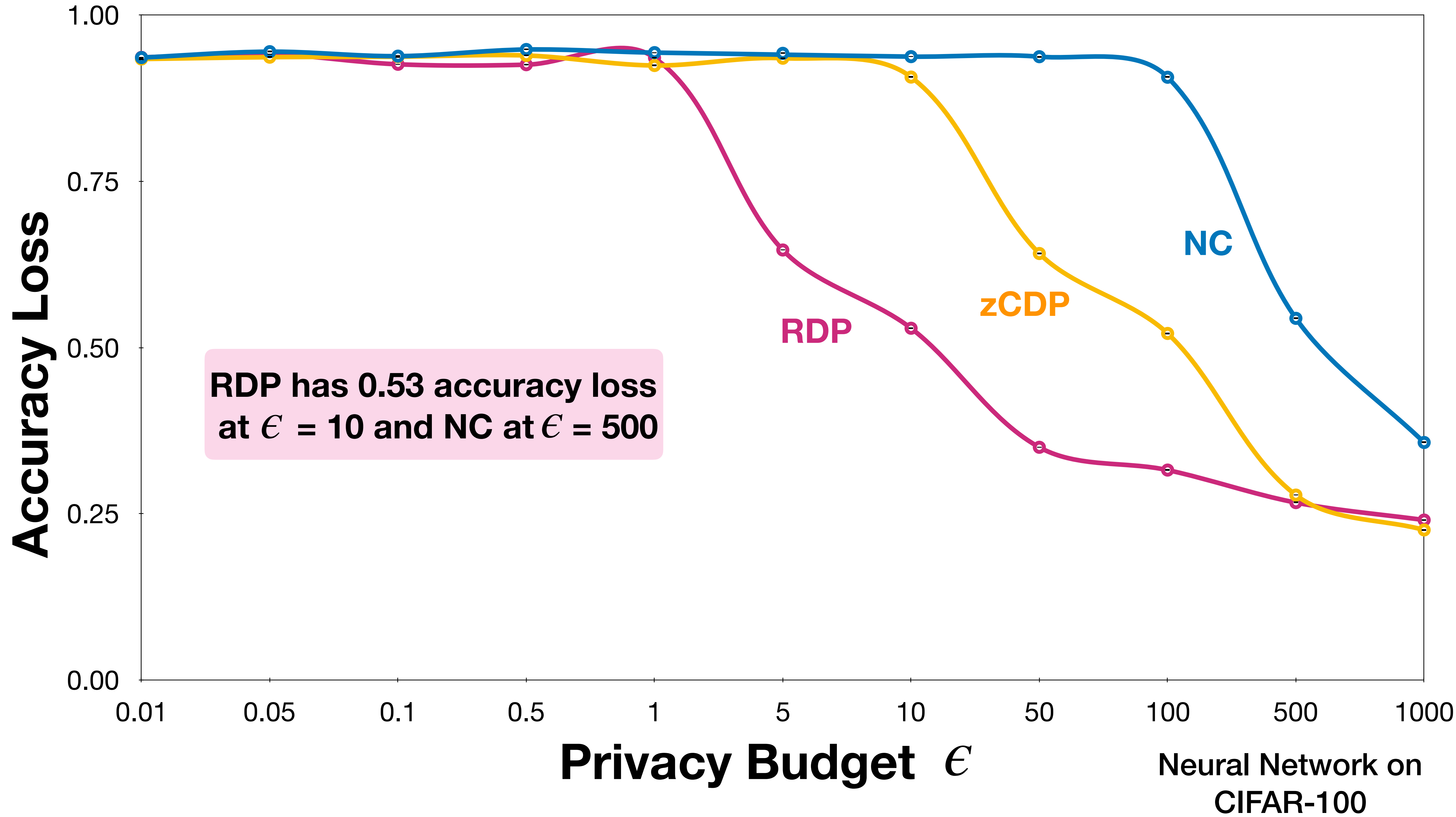


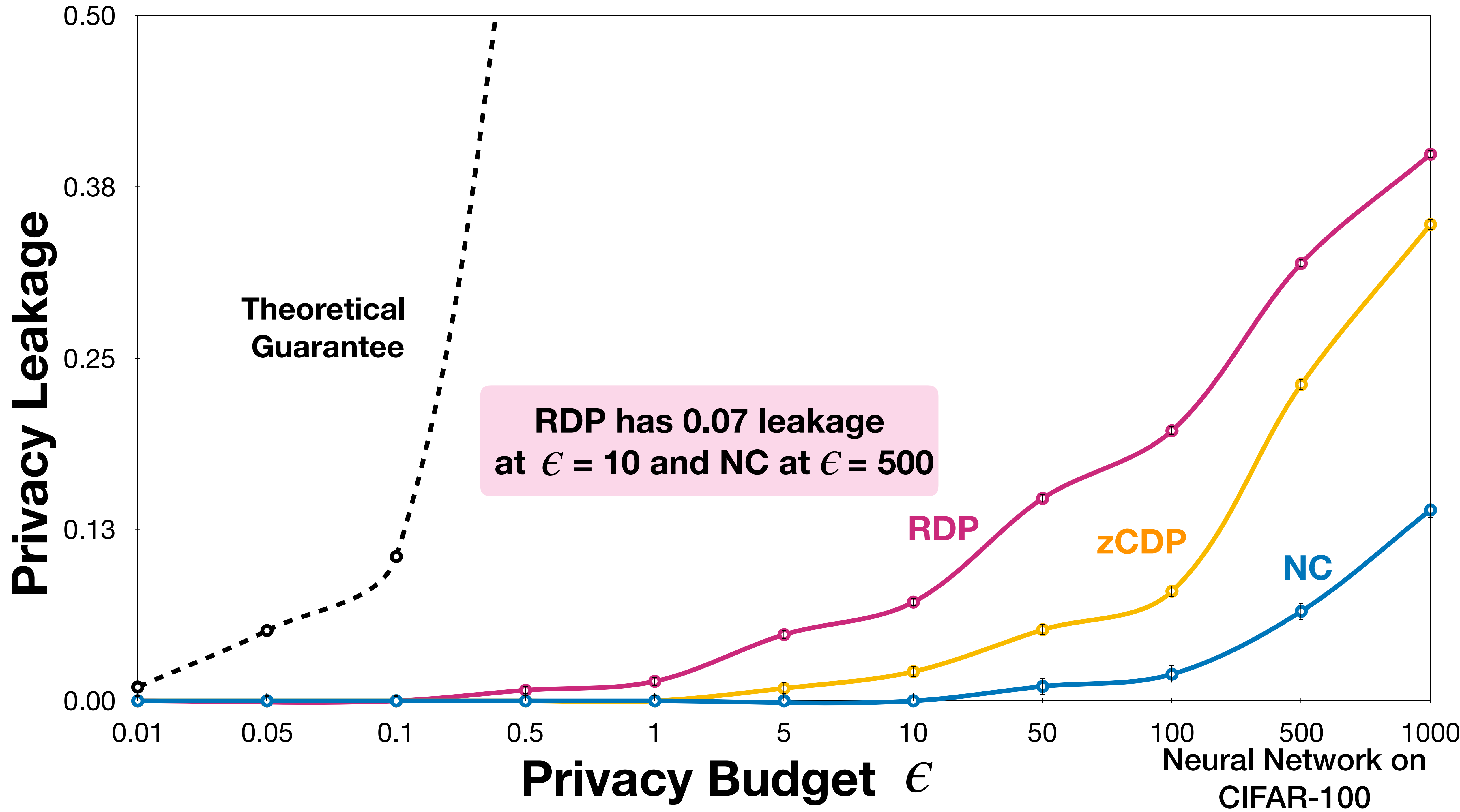


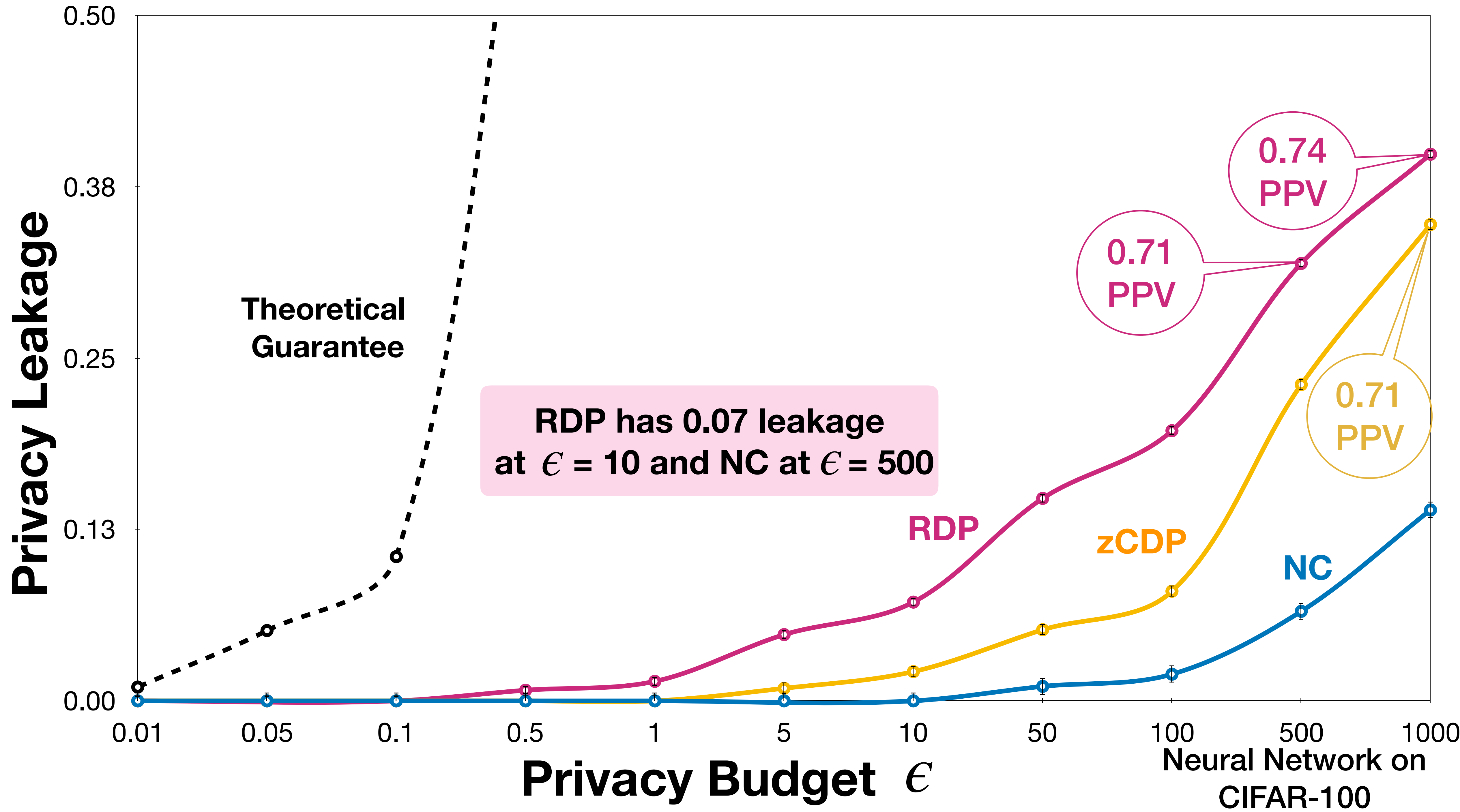
Neural Networks

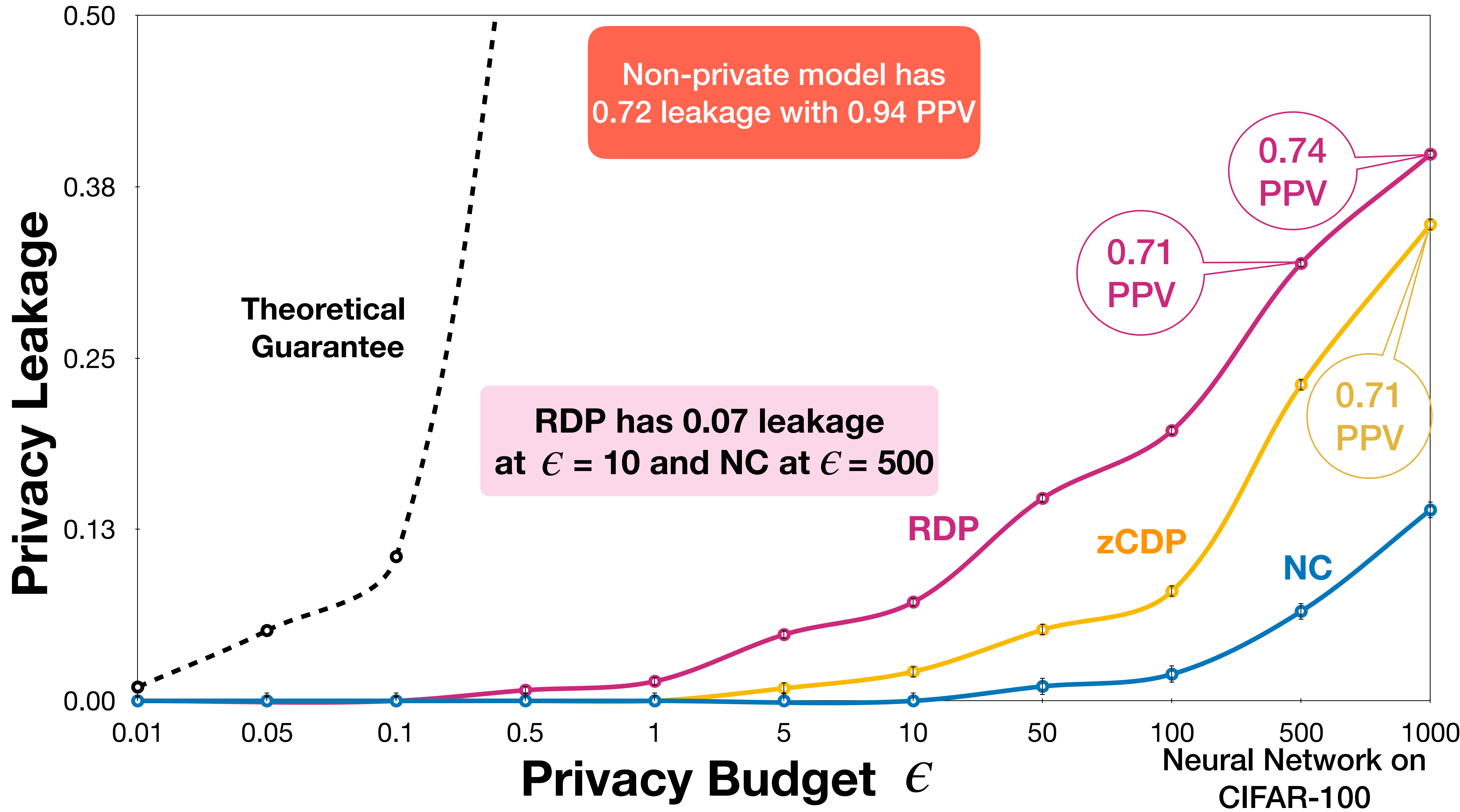


NN has 103,936 trainable parameters so it has more capacity to learn on training data

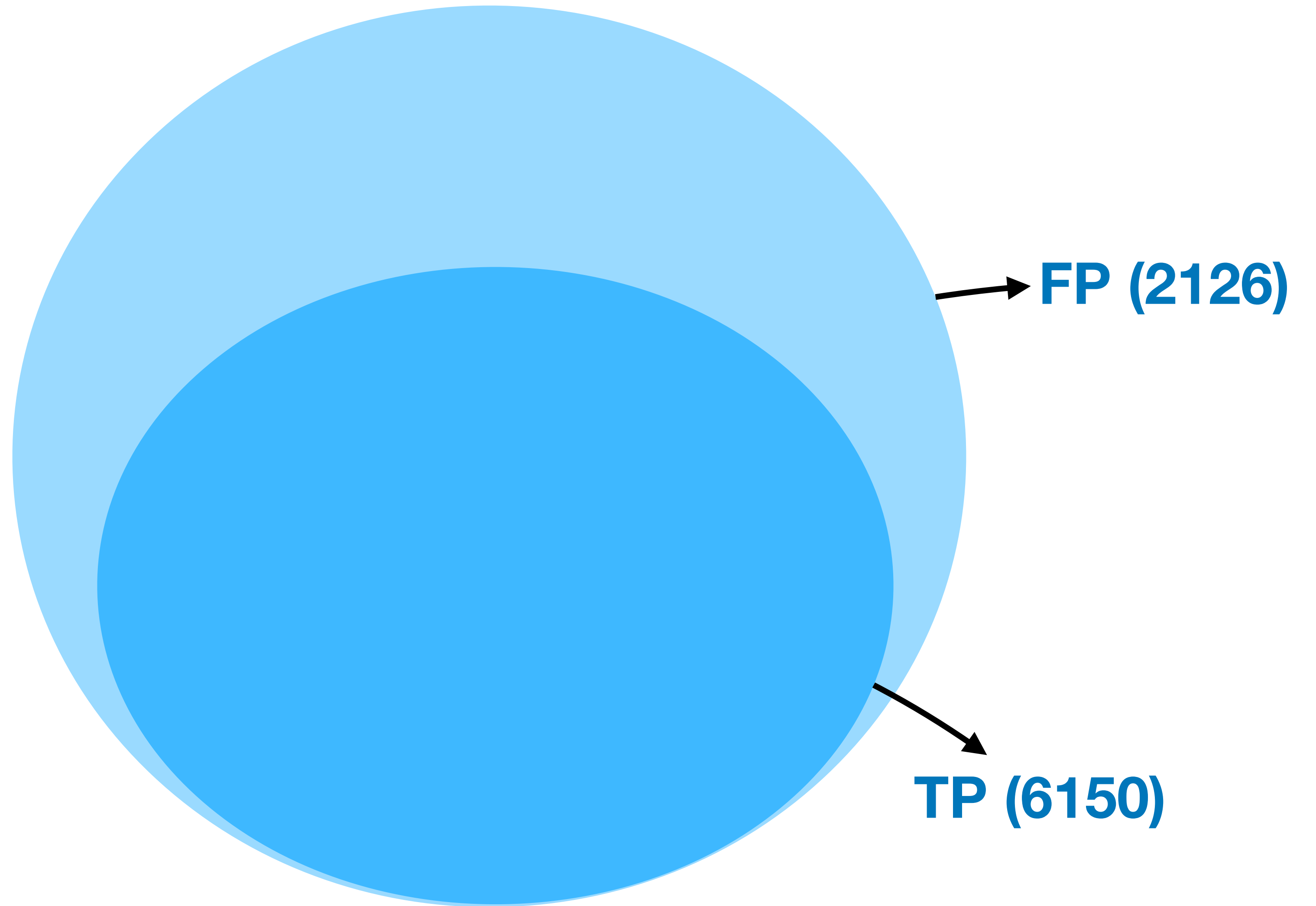






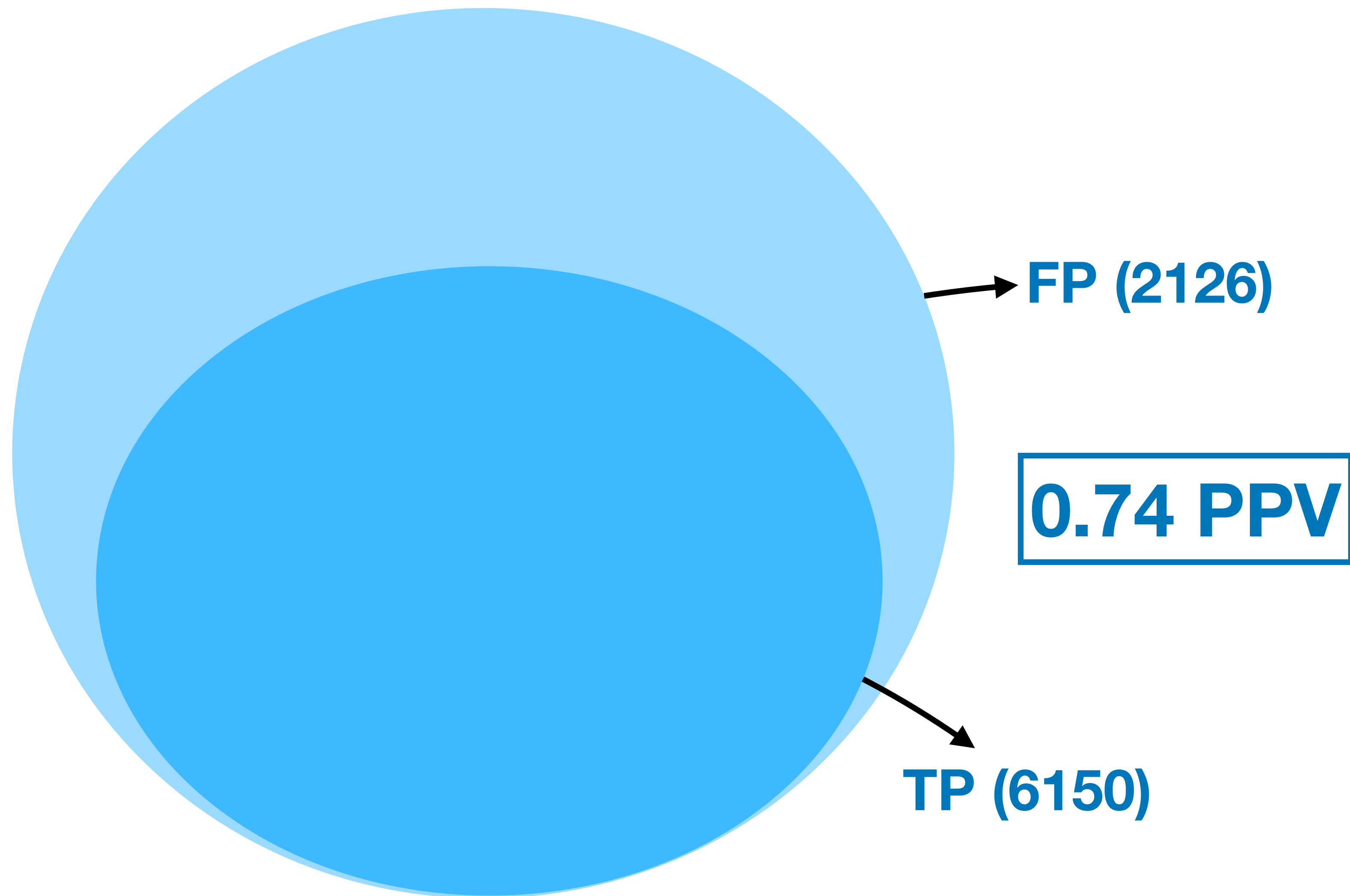


Run 1

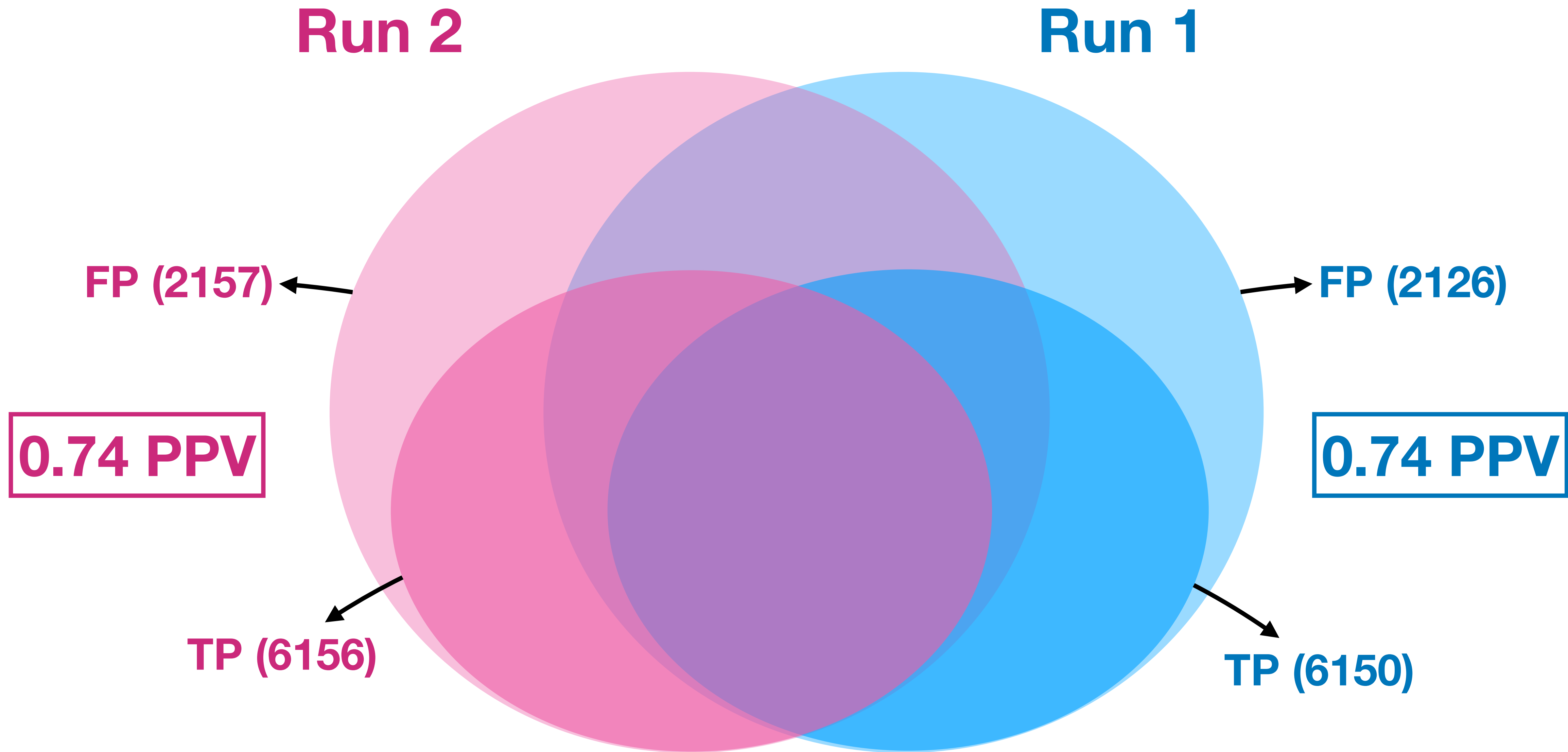


***New results, included in the updated version of the paper**

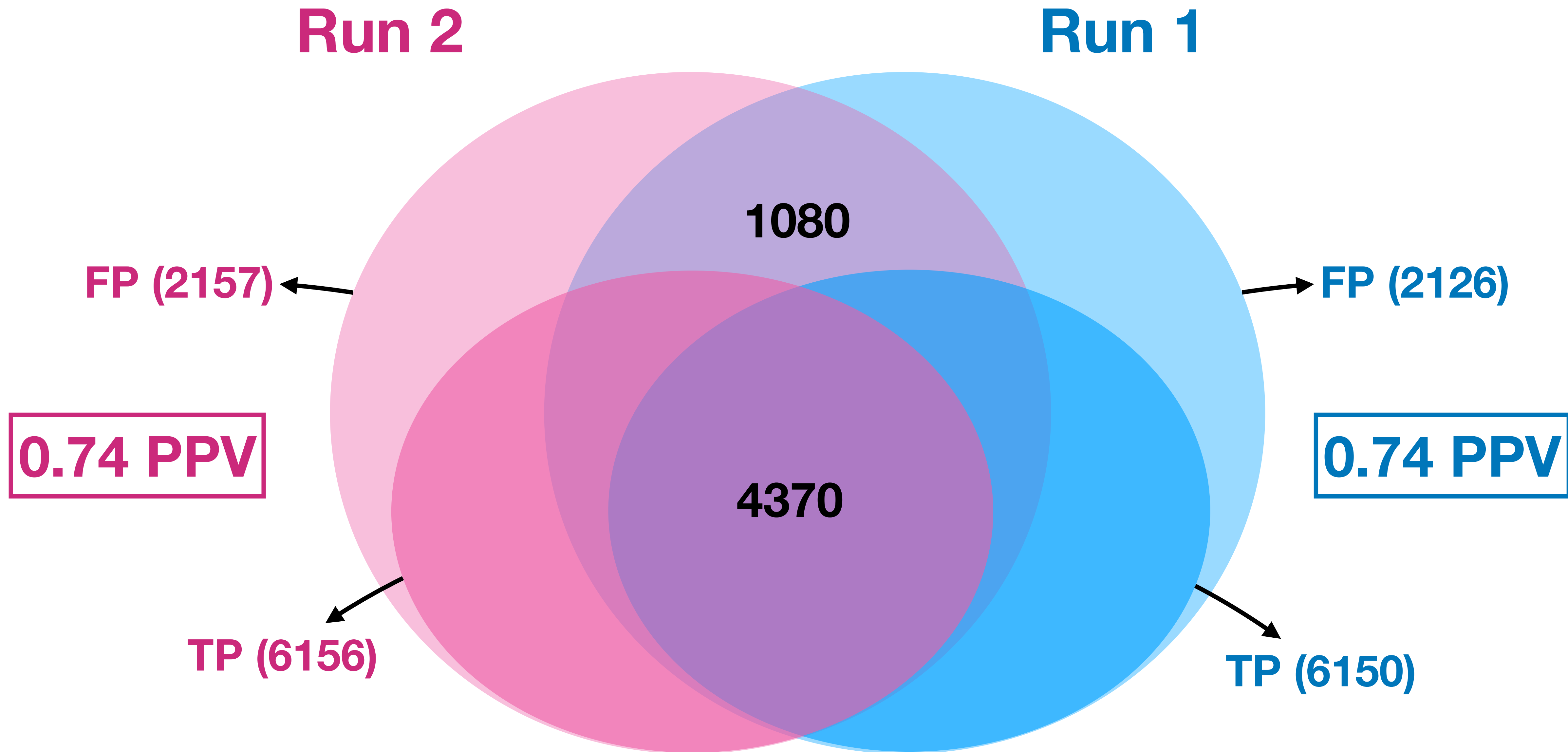
Run 1



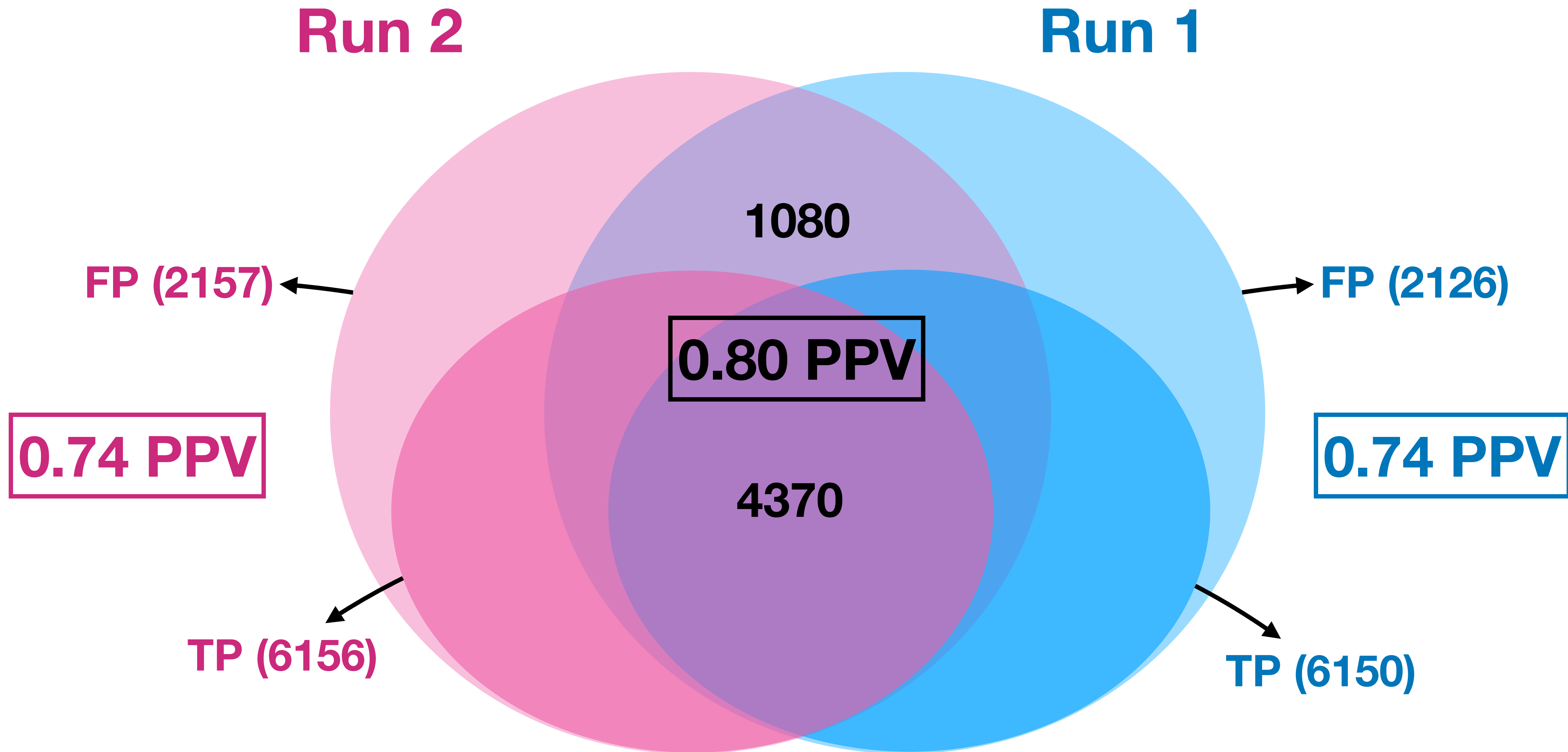
***New results, included in the updated version of the paper**



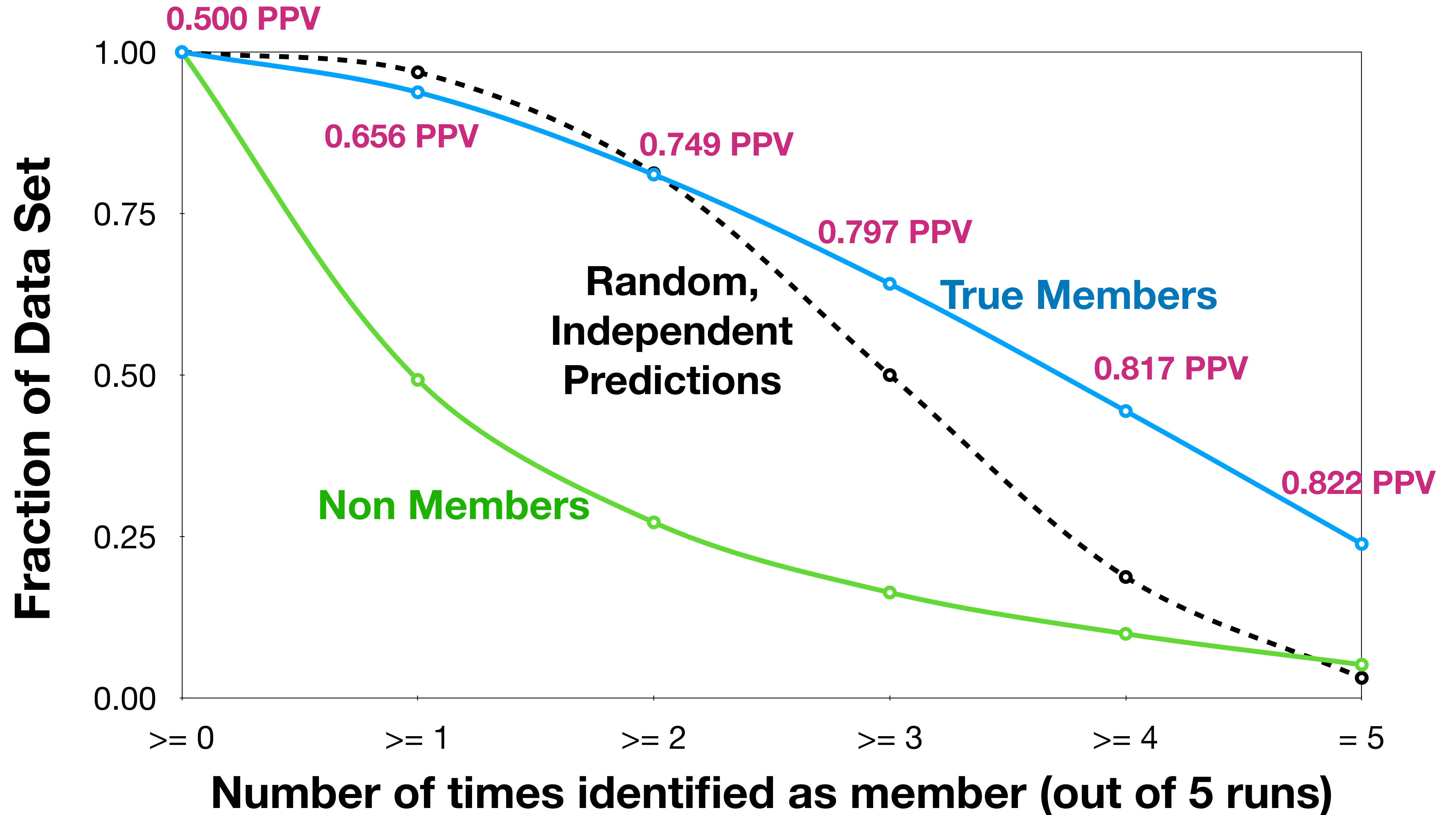
***New results, included in the updated version of the paper**



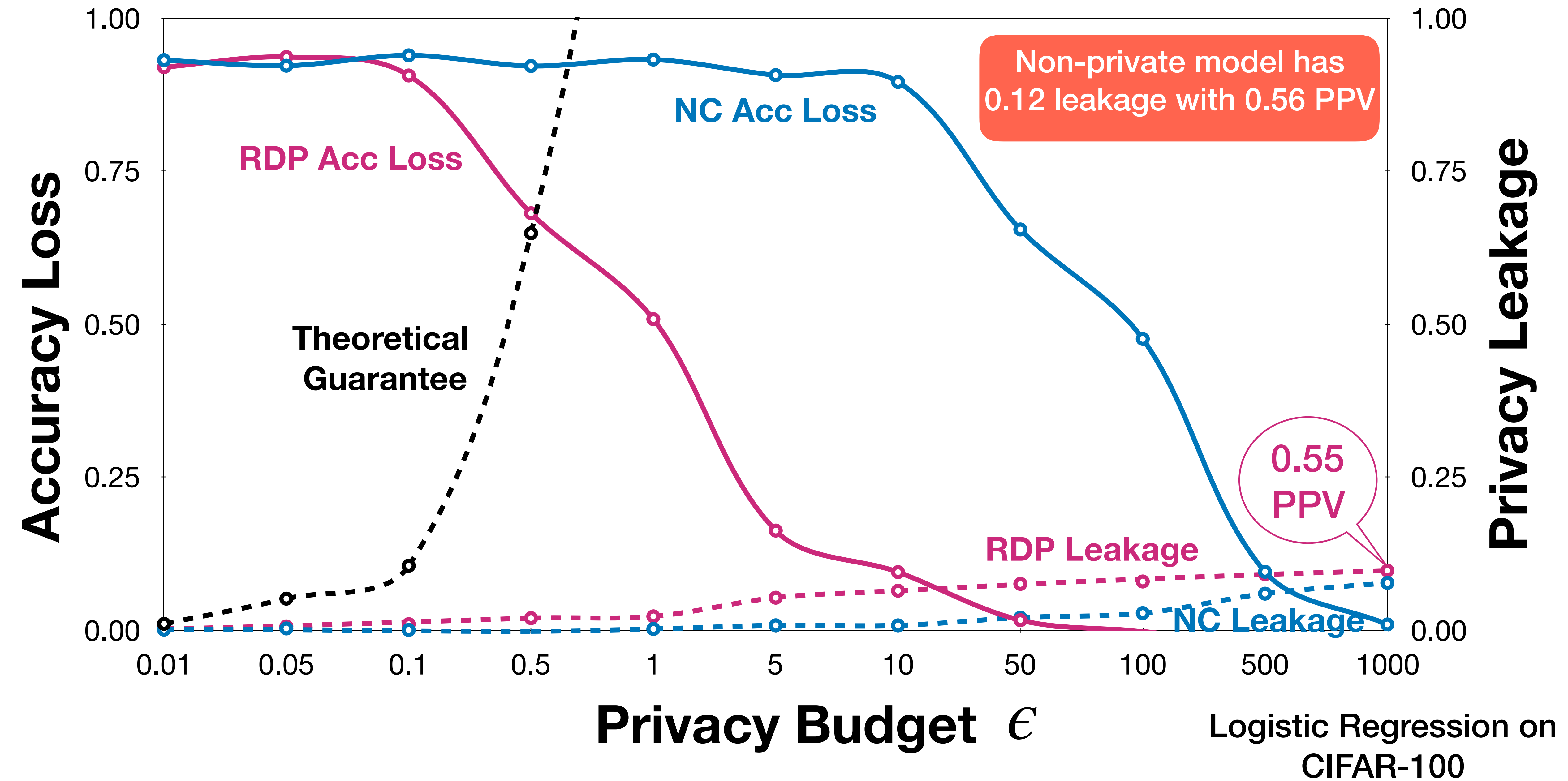
***New results, included in the updated version of the paper**



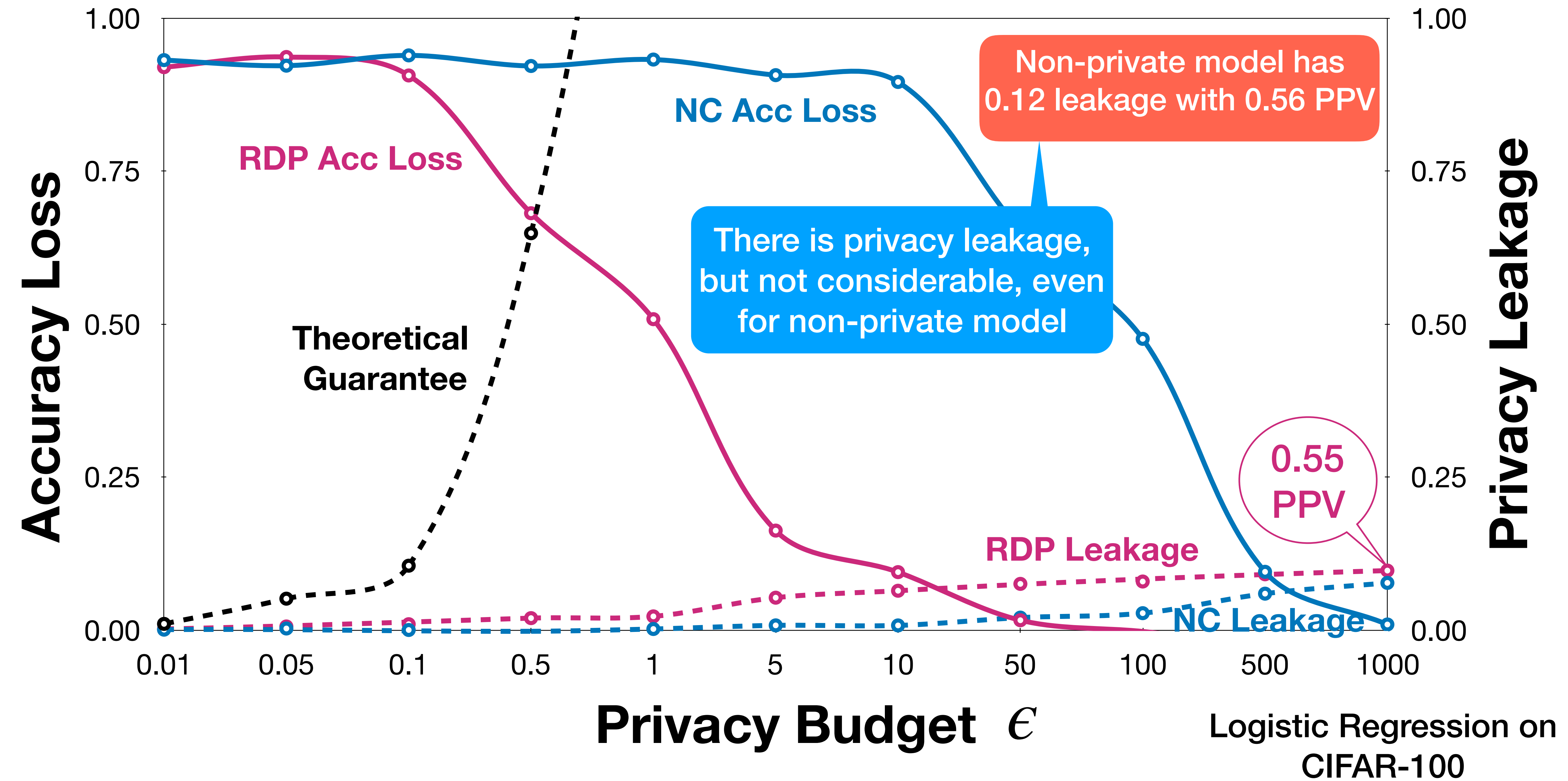
***New results, included in the updated version of the paper**



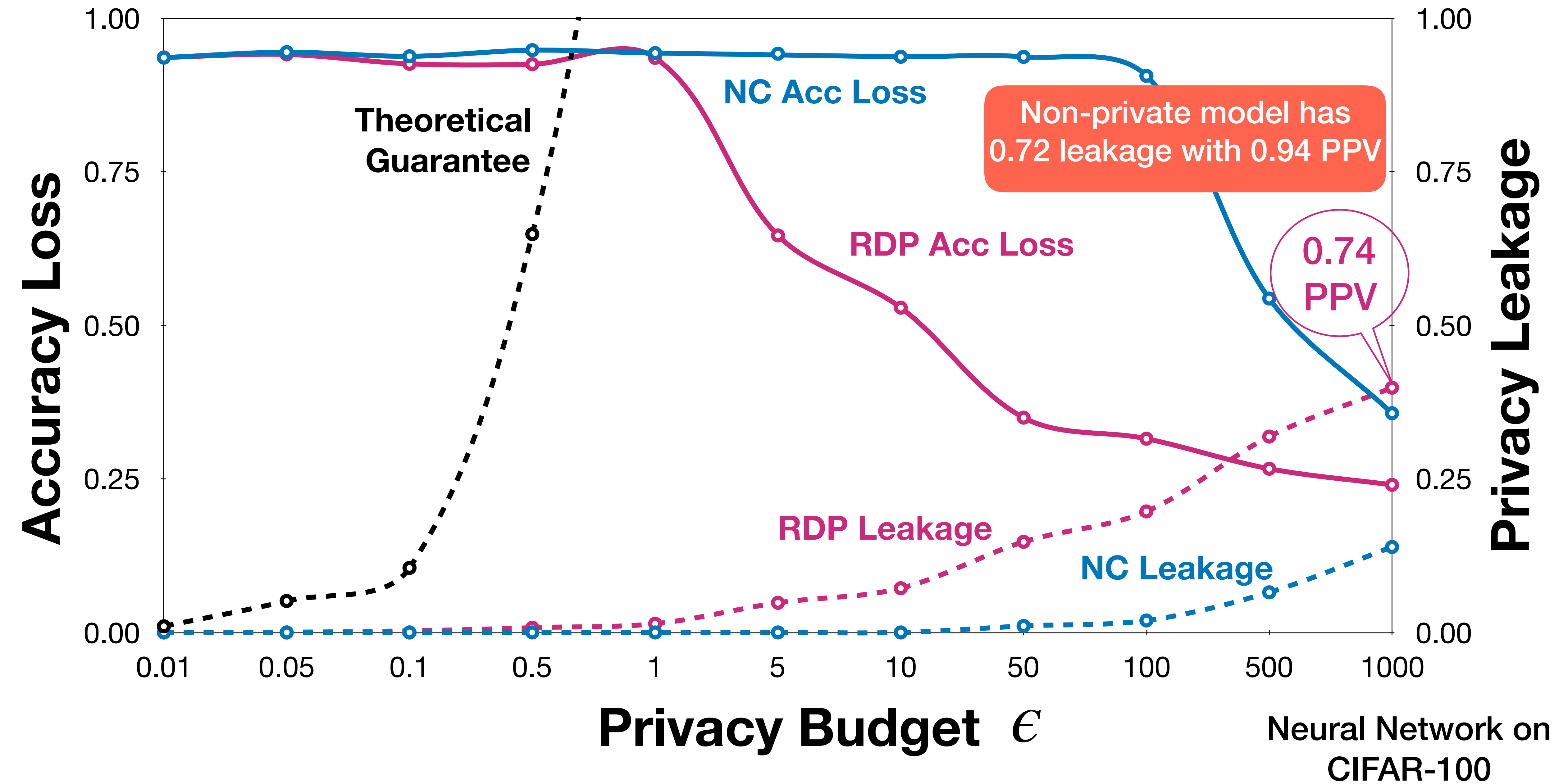
Conclusion



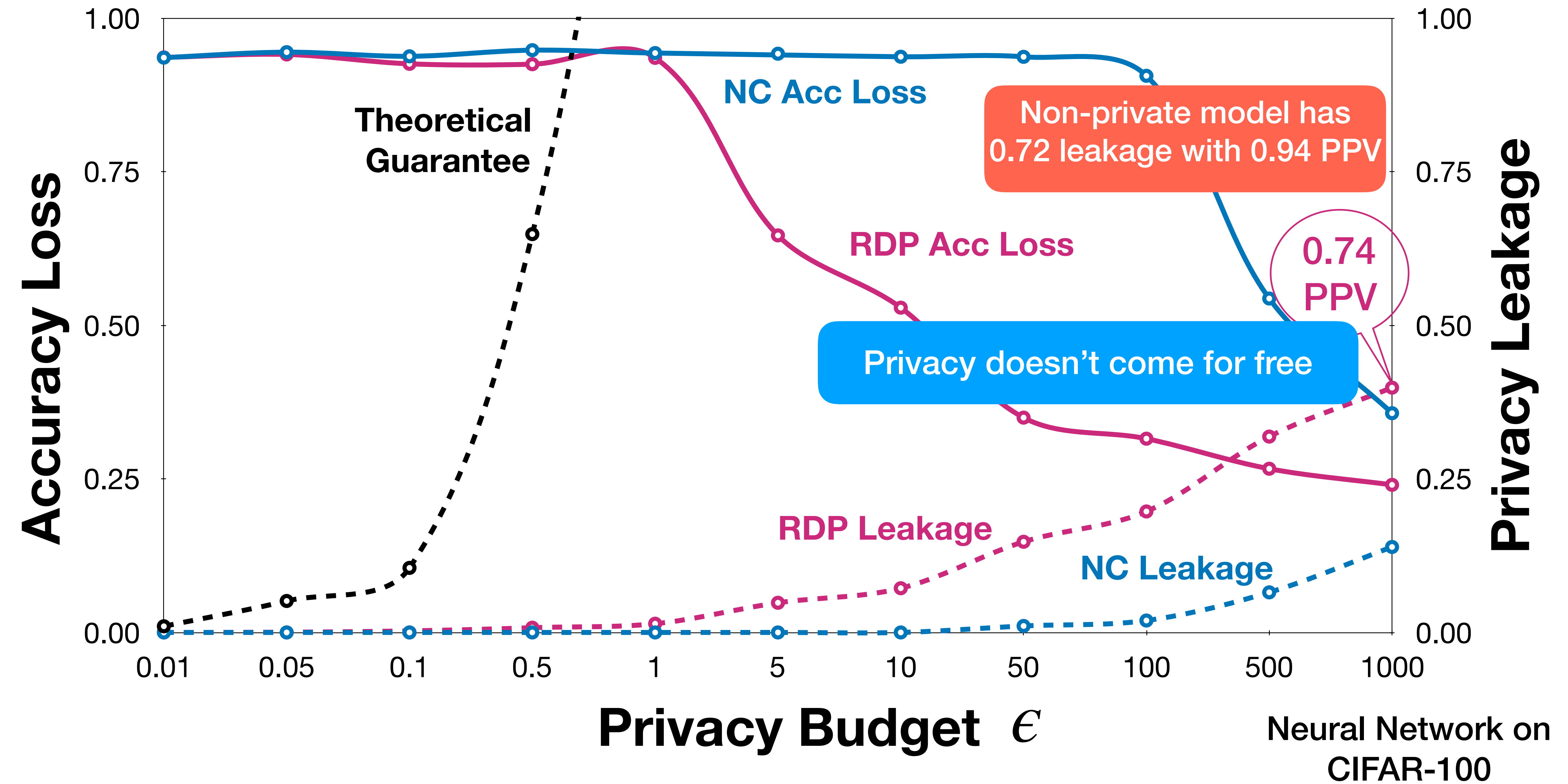
Conclusion



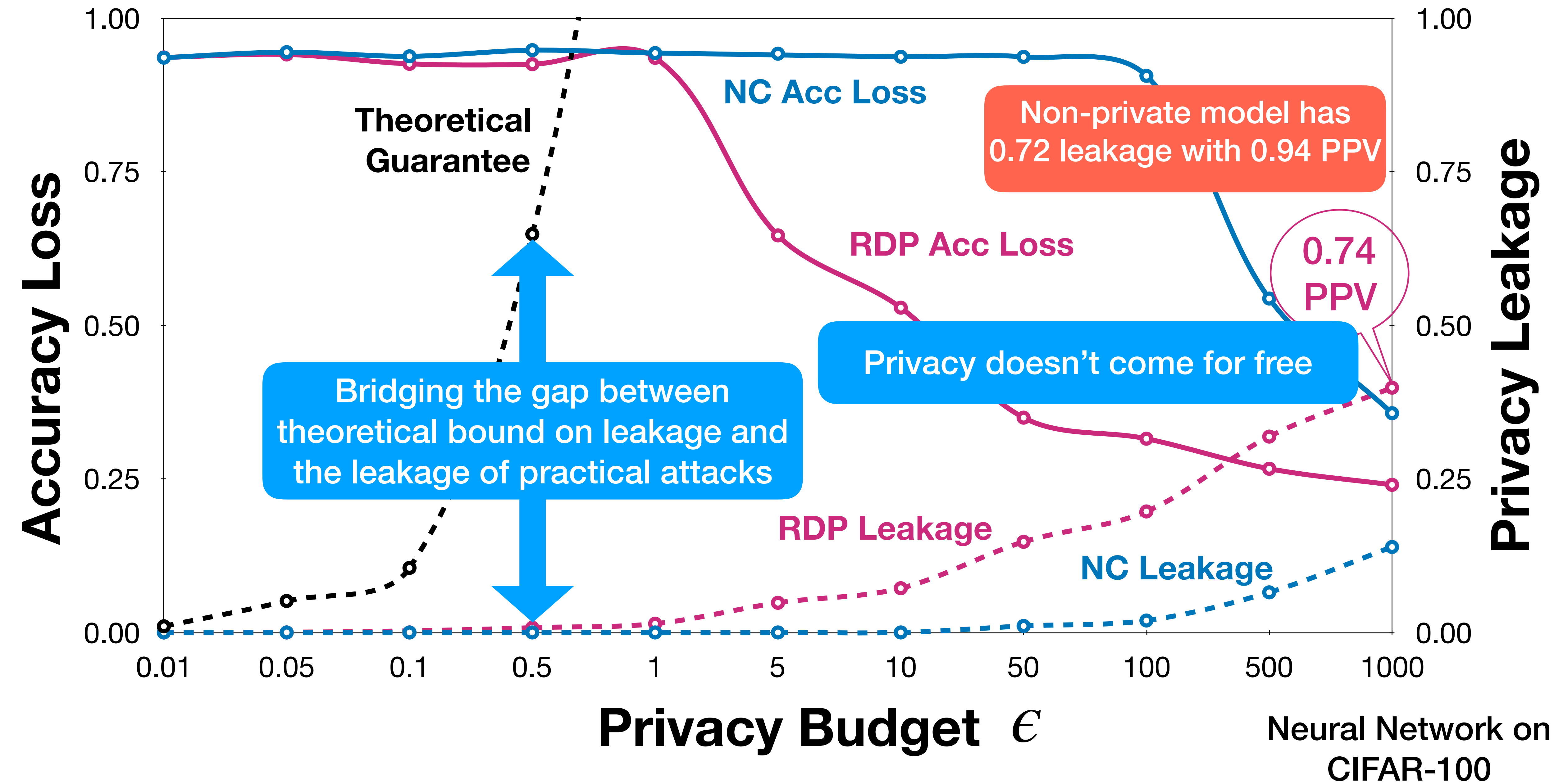
Conclusion



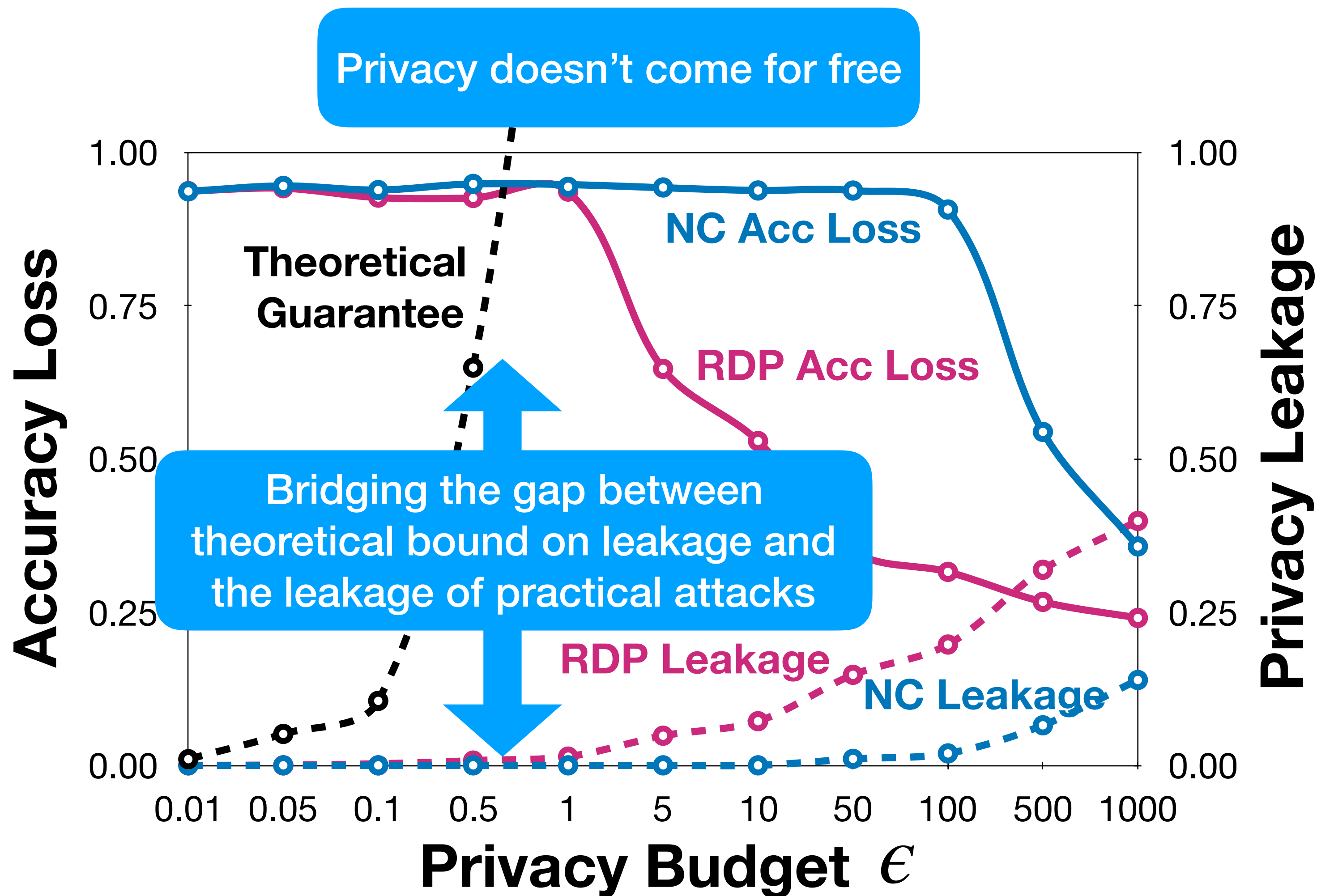
Conclusion



Conclusion



Conclusion



Thank You!

Questions?

Speaker:
Bargav Jayaraman

Project Site:
*[https://bargavjayaraman.github.io/
project/evaluating-dpml/](https://bargavjayaraman.github.io/project/evaluating-dpml/)*

Code Available: <https://github.com/bargavj/EvaluatingDPML>