

# Differentially Private Machine Learning

Gautam Kamath

# NetFlix Prize

- Recommendation engine competition (2006-2009)
- Training data: (anonymized) user ID, movie, rating, date
- Matched with public IMDb data: real name, movie, rating, date
- Class action lawsuit, cancellation of sequel

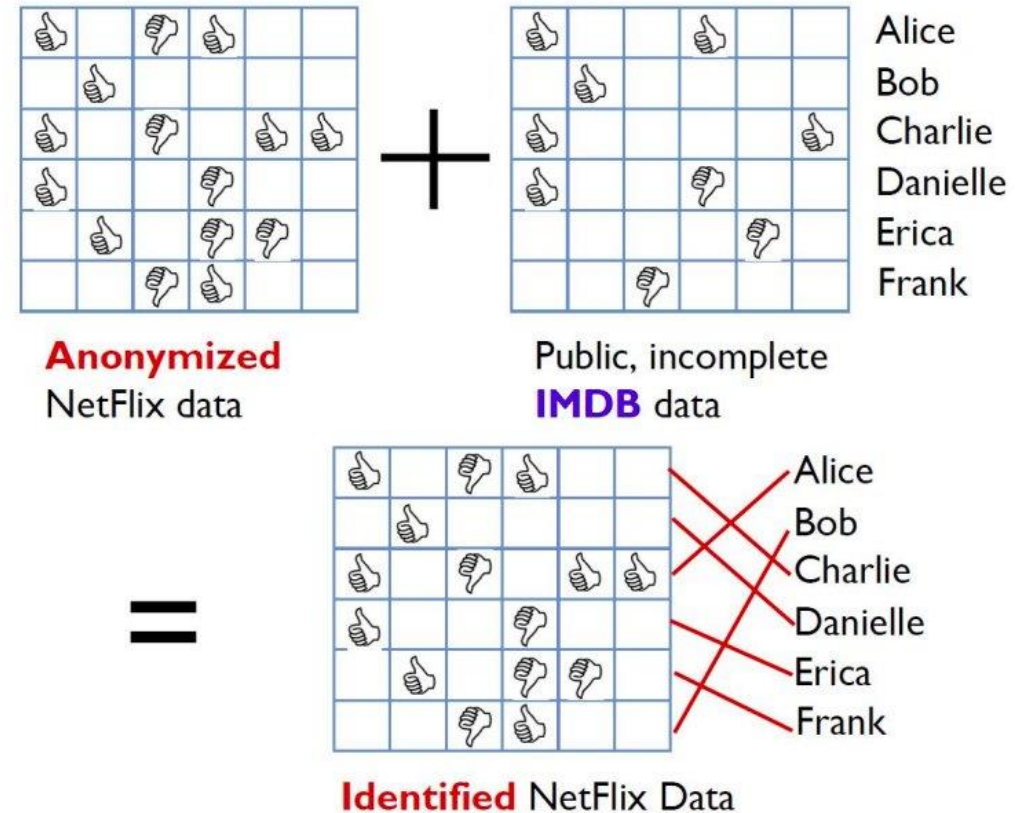
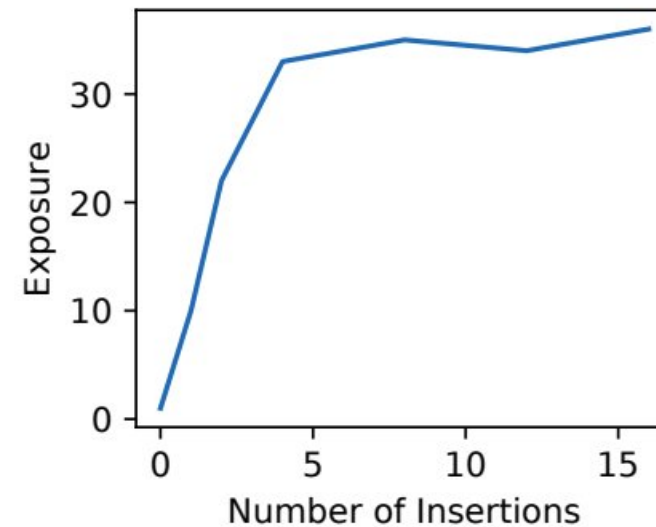


Image credit: Arvind Narayanan

# Memorization in Neural Networks

- Language models
- Log-perplexity of a sequence:
  - $P_{\theta}(x_1, \dots, x_n) = \prod_i \Pr(x_i | f_{\theta}(x_1, \dots, x_{i-1}))$
- “Mary had a little lamb”: low perplexity
- “Correct horse battery staple”: high perplexity
- But what if it were in the training data?

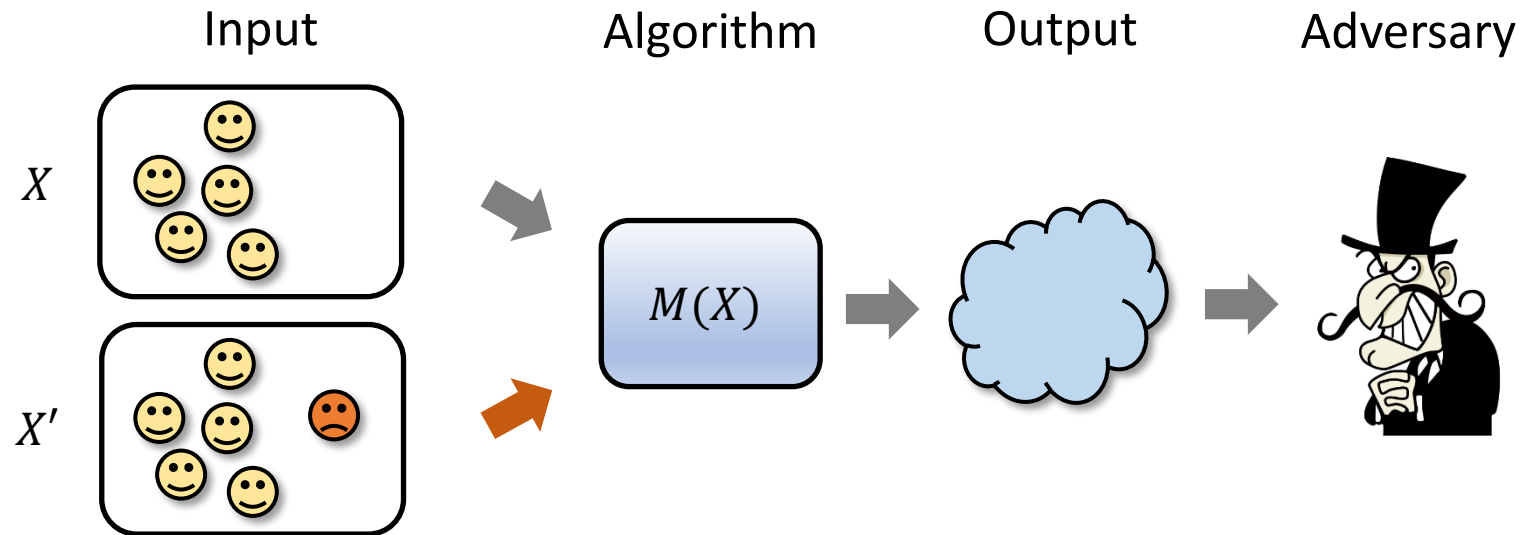
- Canary phrases
  - Is “My SIN is ???-???-???” more likely than it should be?
- Only differential privacy works



[Carlini-Liu-Erlingsson-Kos-Song '19]

See also [Carlini-Tramer-Wallace-Jagielski-HerbertVoss-Lee-Roberts-Brown-Song-Erlingsson-Oprea-Raffel '20]

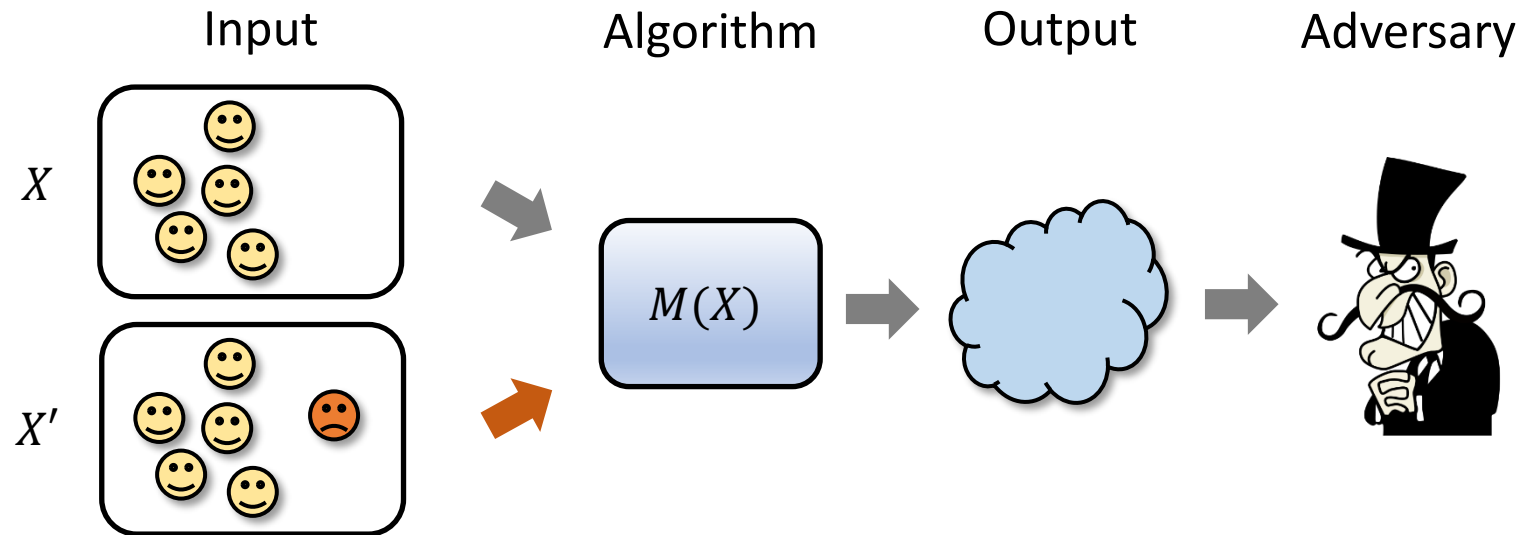
# Differential Privacy (DMNS06)



- $M: D^n \rightarrow R$  is  $(\epsilon, \delta)$ -DP if for all inputs  $X, X'$  which differ on one entry:

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \approx_{\epsilon, \delta} \Pr[M(X') \in S]$$

# Differential Privacy (DMNS06)



- $M: D^n \rightarrow R$  is  $(\epsilon, \delta)$ -DP if for all inputs  $X, X'$  which differ on one entry:

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \leq e^\epsilon \Pr[M(X') \in S] + \delta$$

# Differential Privacy (DMNS06)

$M: D^n \rightarrow R$  is  $(\epsilon, \delta)$ -DP if for all inputs  $X, X'$  which differ on one entry:

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \leq e^\epsilon \Pr[M(X') \in S] + \delta$$

- Google, Apple, Microsoft, 2020 US Census
- $\epsilon \approx 1$  and  $\delta < 1/n$
- Worst-case guarantee
- $e^{\epsilon_1} e^{\epsilon_2} = e^{\epsilon_1 + \epsilon_2}$
- Symmetric definition
- $M$  must be randomized

# What DP does and does not mean

- Outcome is the same whether or not your data is in the dataset
- Protects against linkage and membership inference attacks
- Does *not* prevent statistics and machine learning
  - “Smoking causes cancer”
- Not suitable when we need to identify a specific individual
- Information-theoretic notion

# Properties of Differential Privacy

- Post-processing

- If  $M(X)$  is  $(\epsilon, \delta)$ -DP, then  $f(M(X))$  is  $(\epsilon, \delta)$ -DP

- Group Privacy

- If  $M$  is  $(\epsilon, \delta)$ -DP, and  $X$  and  $X'$  differ in  $k$  entries,

$$\forall S \subseteq R \quad \Pr[M(X) \in S] \leq e^{k\epsilon} \Pr[M(X') \in S] + \delta$$

- Composition

- If  $M = (M_1, \dots, M_k)$  is a sequence of  $k$   $(\epsilon, \delta)$ -DP algorithms

- $M$  is  $(k\epsilon, k\delta)$ -DP (Basic Composition)

- $M$  is  $(O(\sqrt{k}\epsilon \log(1/\delta')), k\delta + \delta')$ -DP (Advanced Composition)



# Gaussian Mechanism

- $\ell_2$ -sensitivity of  $f$

$$\Delta_2^{(f)} = \max_{X \sim X'} \|f(X) - f(X')\|_2$$

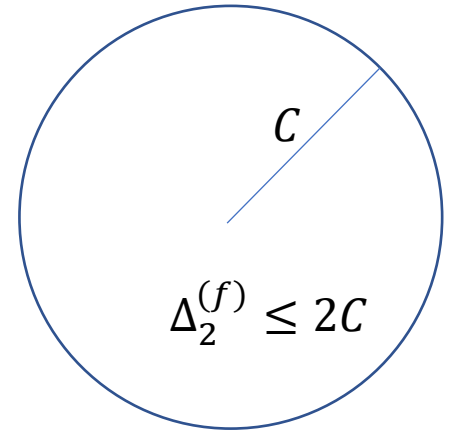
- If  $\|f(X)\|_2 \leq C$ , then  $\Delta_2^{(f)} \leq 2C$

- Gaussian Mechanism

$$M(X) = f(X) + (Y_1, \dots, Y_k)$$

Where  $f(X) \in \mathbb{R}^k$ , and the  $Y_i$ 's are  $\approx N(0, \Delta^2 / \epsilon^2)$

- $(\epsilon, \delta)$ -DP

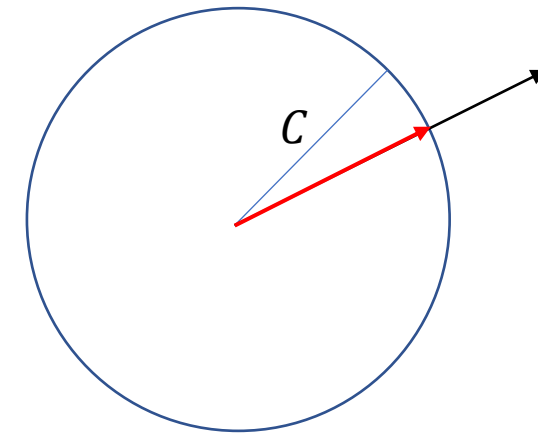


# Stochastic Gradient Descent

1. Choose a random minibatch  $B$  of points from the dataset
2. Compute the average gradient  $\frac{1}{|B|} \sum_{(x,y) \in B} \nabla \ell(\theta_t, x, y)$
3. Take a step in the negative direction of the gradient
4. Repeat  $k$  times

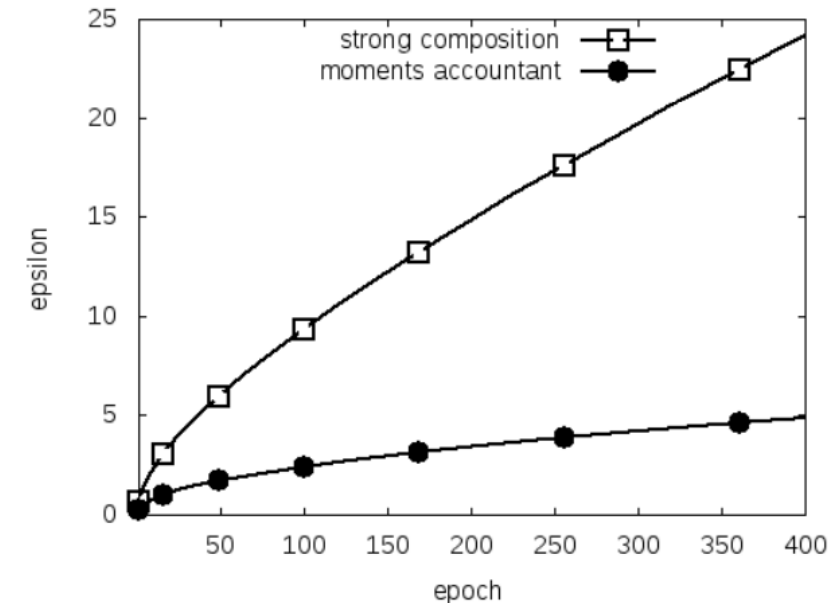
# Differentially Private Stochastic Gradient Descent

1. Sample a “lot” of points of (expected) size  $L$  by selecting each point to be in the lot with probability  $L/n$
2. For each point in the lot, compute the gradient  $\nabla\ell(\theta_t, x, y)$  and “clip” it to have  $\ell_2$  norm at most  $C$
3. Average the clipped gradients and add **Gaussian noise**
  - Apply the Gaussian Mechanism
4. Take a step in the negative direction of resulting vector
5. Repeat  $k$  times



# Privacy of DPSGD (Informal)

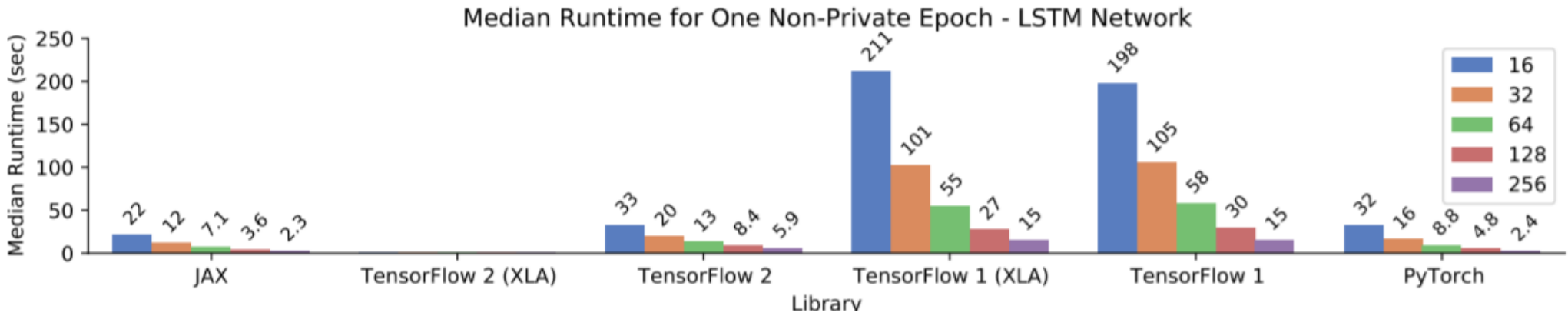
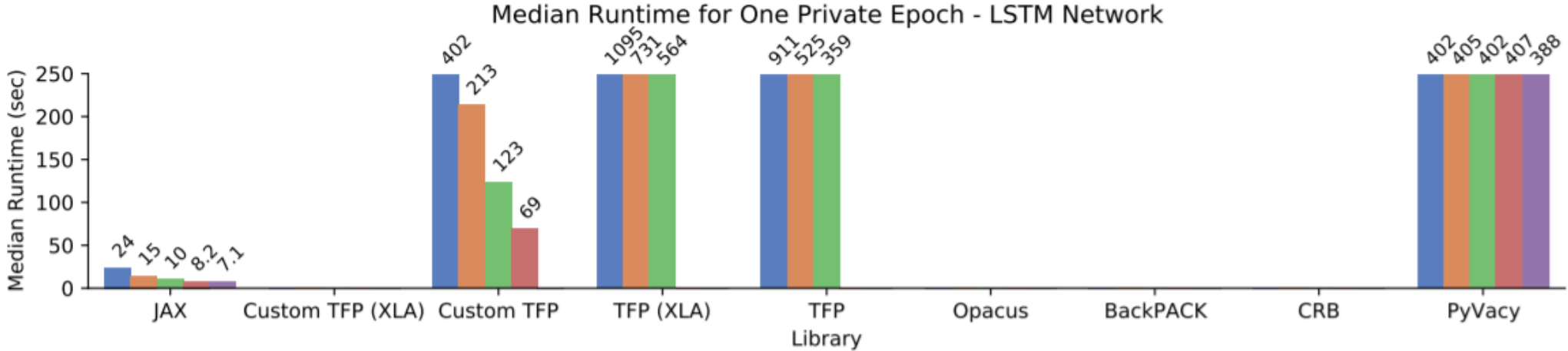
- Suppose one step of DPSGD has privacy with parameter  $\varepsilon$
- Since we subsample with probability  $L/n$ , each step is  $\varepsilon L/n$ 
  - “Privacy amplification by subsampling”
- $k$  steps have privacy with parameter of  $\varepsilon\sqrt{k}L/n$ 
  - Advanced composition
- Better analysis: “Moments accountant”



# Does it work?

Data	$\epsilon$ -DP	Source	Test Accuracy (%)		
			CNN	ScatterNet+linear	ScatterNet+CNN
MNIST	1.2	Feldman & Zrnic (2020)	<u>96.6</u>	<b>98.1 <math>\pm</math> 0.1</b>	97.8 $\pm$ 0.1
	2.0	Abadi et al. (2016)	95.0	<b>98.5 <math>\pm</math> 0.0</b>	<b>98.4 <math>\pm</math> 0.1</b>
	2.32	Bu et al. (2019)	96.6	<b>98.6 <math>\pm</math> 0.0</b>	98.5 $\pm$ 0.0
	2.5	Chen & Lee (2020)	90.0	<b>98.7 <math>\pm</math> 0.0</b>	98.6 $\pm$ 0.0
	2.93	Papernot et al. (2020a)	<u>98.1</u>	<b>98.7 <math>\pm</math> 0.0</b>	<b>98.7 <math>\pm</math> 0.1</b>
	3.2	Nasr et al. (2020)	96.1	–	–
	6.78	Yu et al. (2019b)	93.2	–	–
Fashion-MNIST	2.7	Papernot et al. (2020a)	<u>86.1</u>	<b>89.5 <math>\pm</math> 0.0</b>	88.7 $\pm$ 0.1
	3.0	Chen & Lee (2020)	82.3	<b>89.7 <math>\pm</math> 0.0</b>	89.0 $\pm$ 0.1
CIFAR-10	3.0	Nasr et al. (2020)	<u>55.0</u>	67.0 $\pm$ 0.1	<b>69.3 <math>\pm</math> 0.2</b>
	6.78	Yu et al. (2019b)	44.3	–	–
	7.53	Papernot et al. (2020a)	<u>66.2</u>	–	–
	8.0	Chen & Lee (2020)	53.0	–	–

# DPSGD can be slow!



[Subramani-Vadivelu-K. '20]

# Architectures for DPSGD

- Tanh >> ReLU? [Papernot-Thakurta-Song-Chien-Erlingssson '21]
- Bigger models are not always better

# Hyperparameters

- Even more hyperparameters
  - Learning rate, lot size, clipping norm, number of epochs, noise multiplier
- Non-private way: grid search, measure accuracy on validation set
- Pay in privacy budget for each run!
- Options:
  - Private methods for hyperparameter optimization [Liu-Talwar '19]
  - Transfer hyperparameters from related public data
  - Cheat and ignore privacy budget for multiple runs...



# Conclusion

- Private machine learning is here!
- But there's still a lot of work to do...