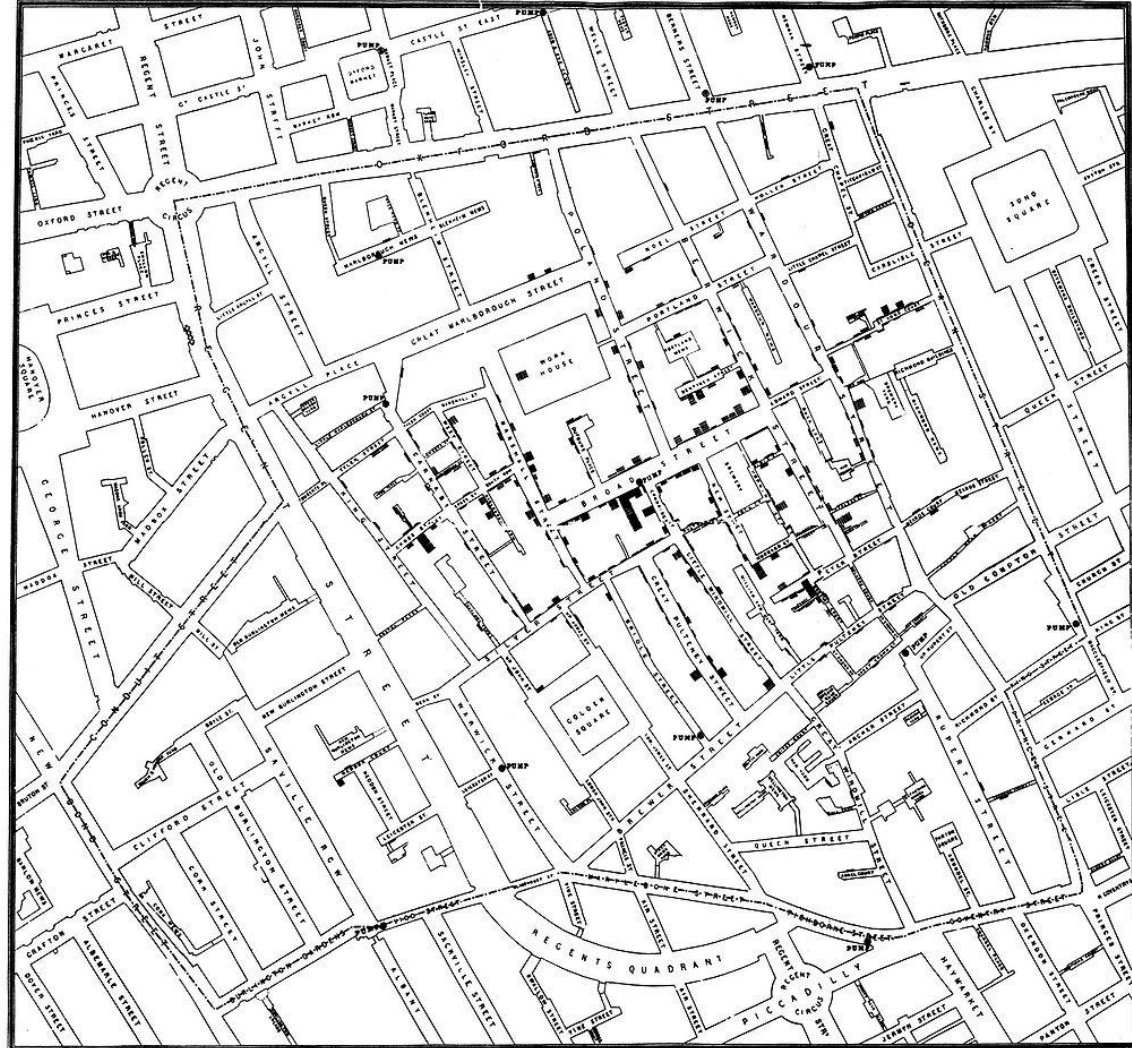


# DEEP LEARNING WITH DIFFERENTIAL PRIVACY

Martin Abadi, Andy Chu, Ian Goodfellow\*,  
Brendan McMahan, Ilya Mironov, Kunal Talwar, Li Zhang  
Google

\* Open AI





C. F. Chubb & Co. Ltd. Southampton St. London

SCALE 80 INCHES TO A MILE.





C. F. Chubb, Ltd. Southampton St. London

SCALE 20 INCHES TO A MILE.





# Deep Learning

- Cognitive tasks: speech, text, image recognition
- Natural language processing: sentiment analysis, translation
- Planning: games, autonomous driving



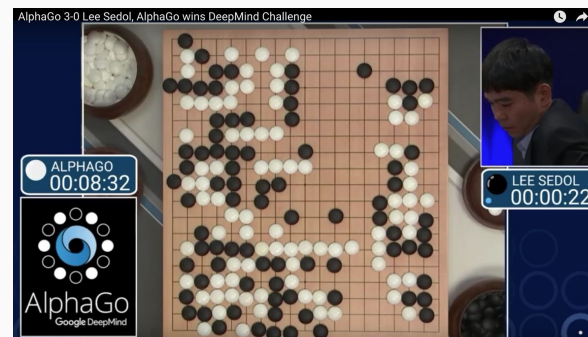
Fashion



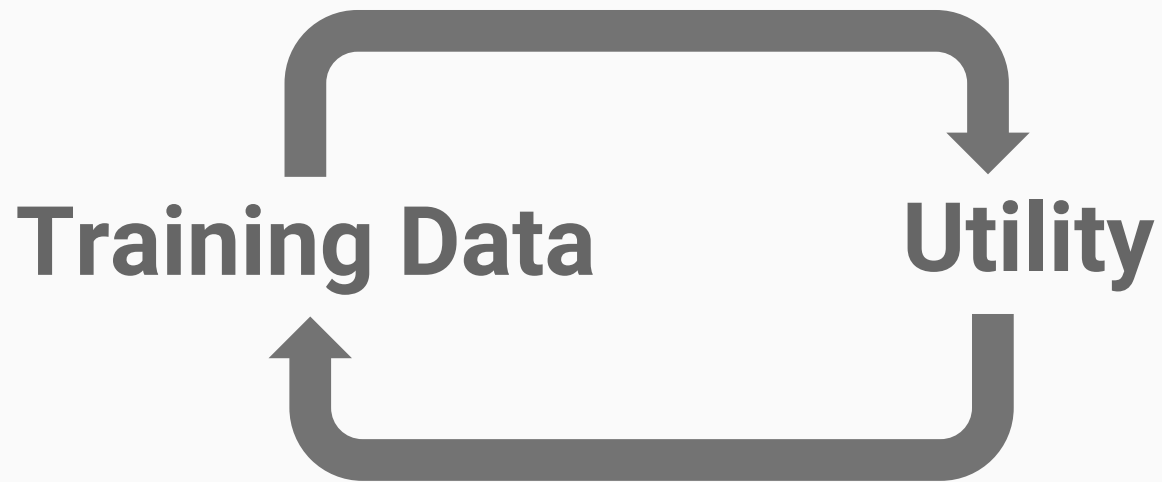
Self-driving cars



Translation



Gaming



# Privacy of Training Data



Data encryption in transit and at rest



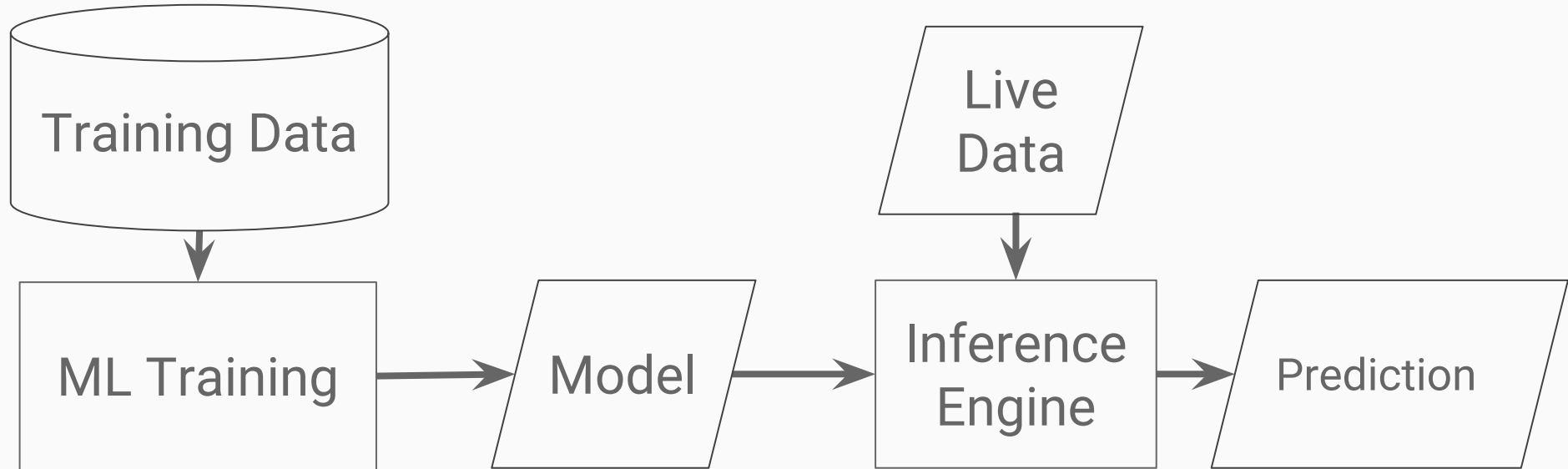
Data retention and deletion policies



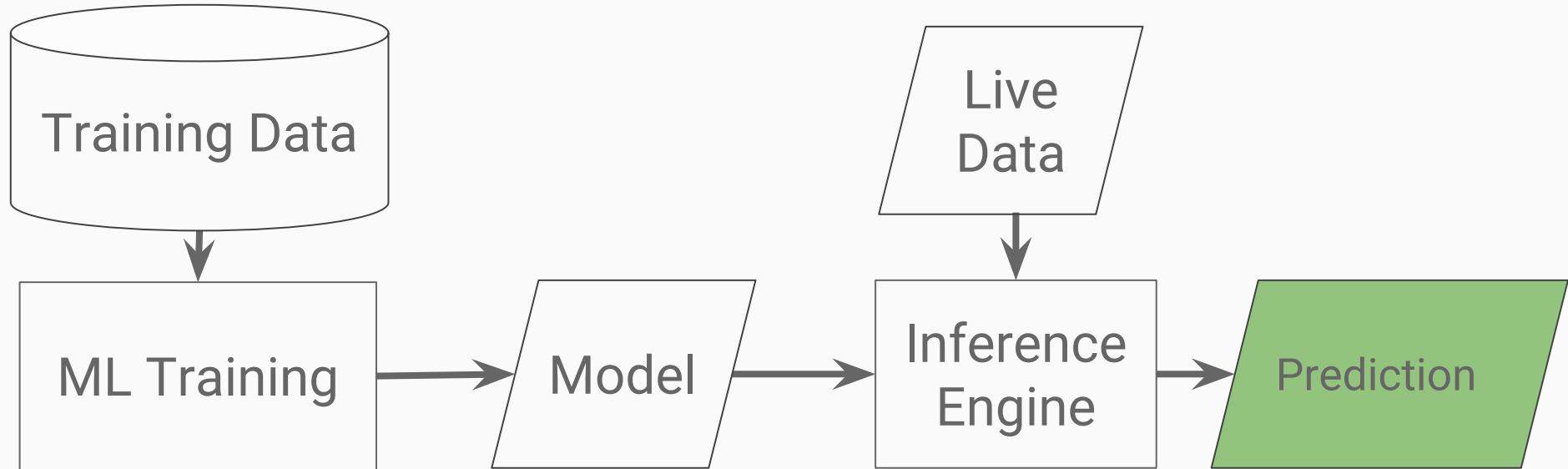
ACLs, monitoring, auditing

What do models reveal about training data?

# ML Pipeline and Threat Model

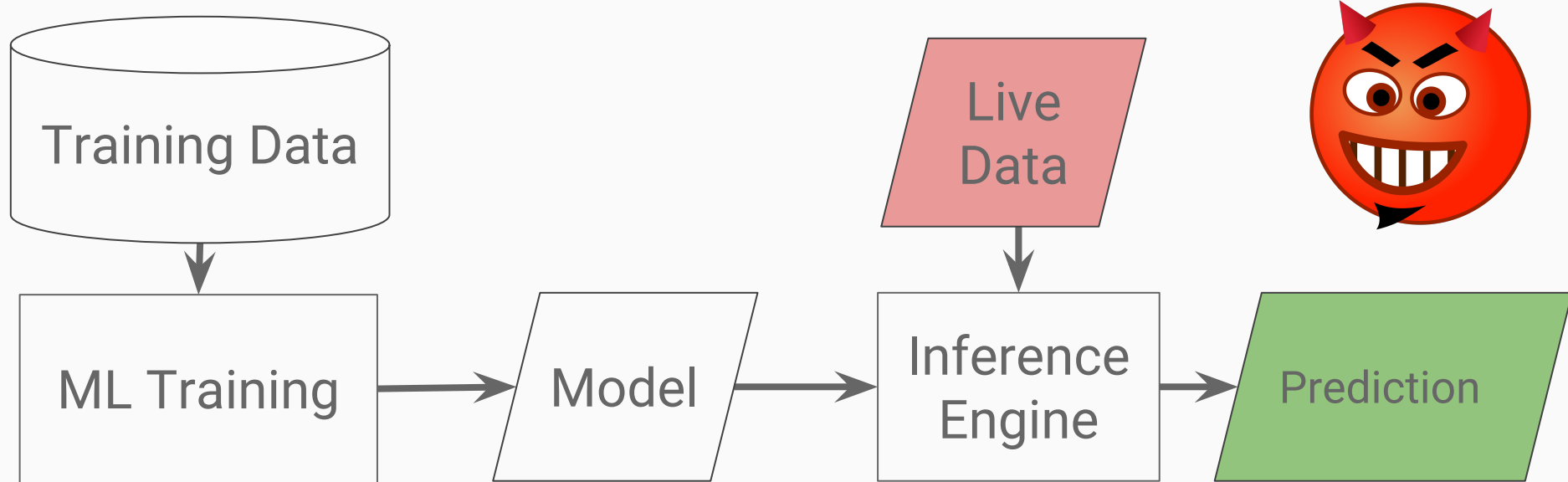


# ML Pipeline and Threat Model

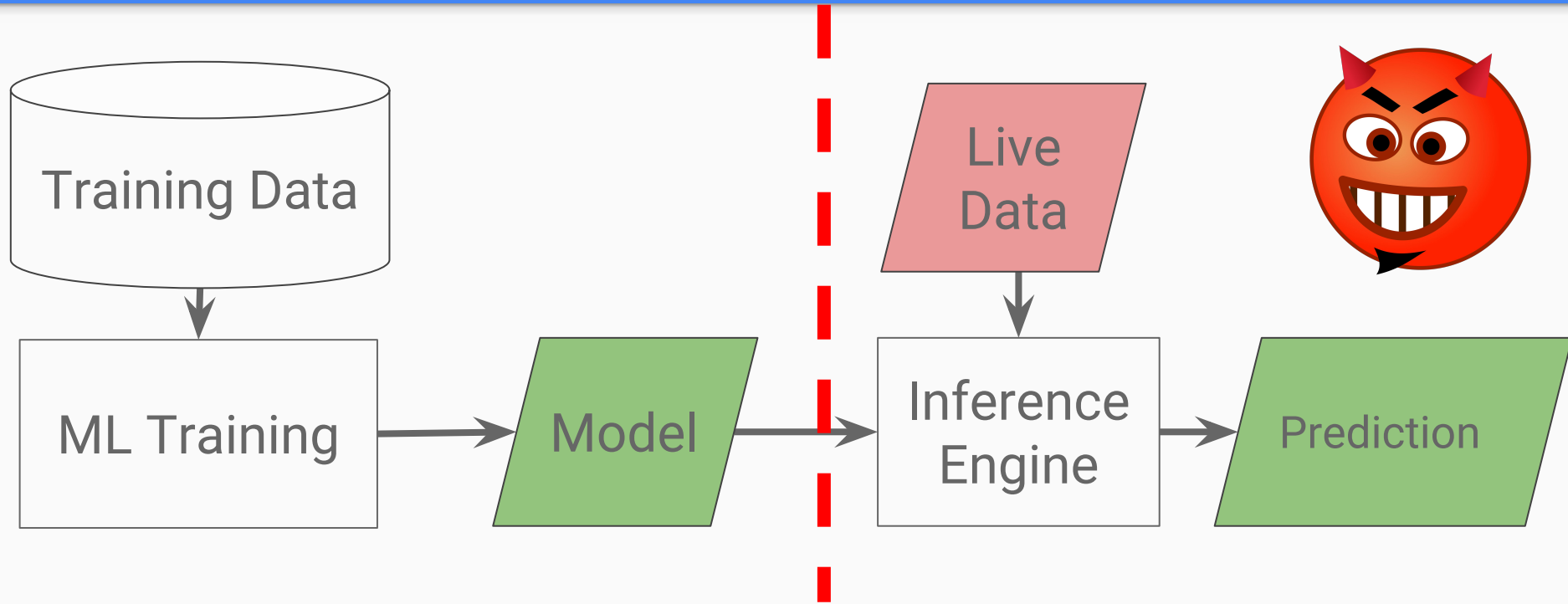




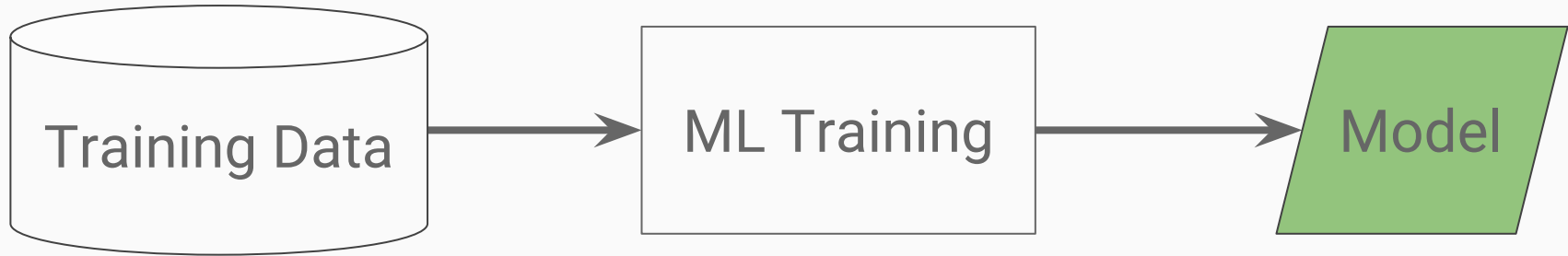
# ML Pipeline and Threat Model



# ML Pipeline and Threat Model



# ML Pipeline and Threat Model



# Machine Learning Privacy Fallacy

Since our ML system is good, it automatically protects privacy of training data.

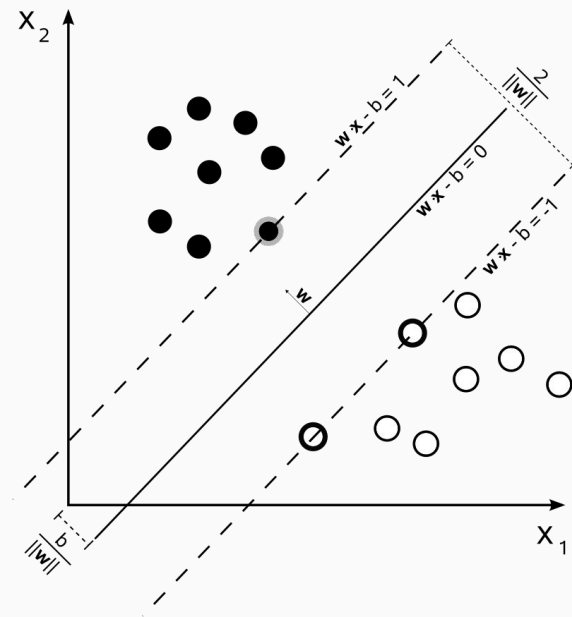
# Machine Learning Privacy Fallacy

- Examples when it just ain't so:
  - Person-to-person similarities
  - Support Vector Machines
- Models can be very large
  - Millions of parameters
- Empirical evidence to the contrary:
  - M. Fredrikson, S. Jha, T. Ristenpart, "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures", CCS 2015
  - R. Shokri, M. Stronati, V. Shmatikov, "Membership Inference Attacks against Machine Learning Models", <https://arxiv.org/abs/1610.05820>



# Machine Learning Privacy Fallacy

- Examples when it just ain't so:
  - Person-to-person similarities
  - Support Vector Machines
- Models can be very large
  - Millions of parameters



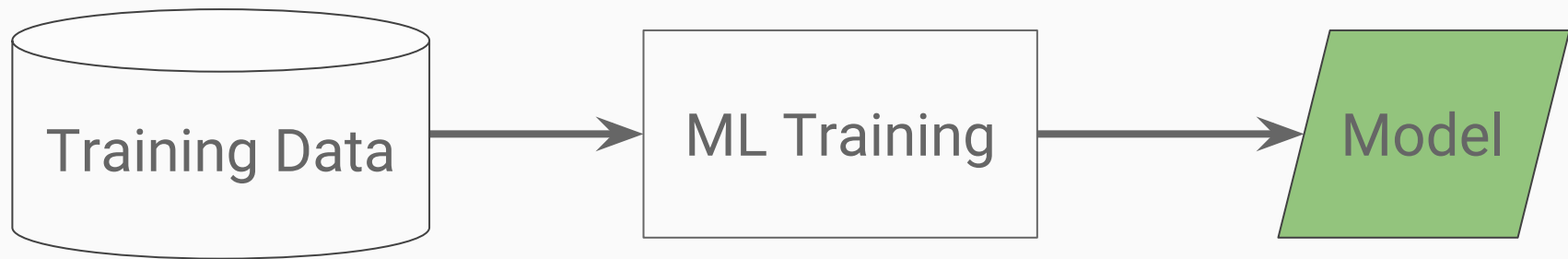


# Model Inversion Attack

- M. Fredrikson, S. Jha, T. Ristenpart, “Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures”, CCS 2015



- R. Shokri, M. Stronati, V. Shmatikov, “Membership Inference Attacks against Machine Learning Models”, <https://arxiv.org/abs/1610.05820>



# Deep Learning Recipe

1. Loss function
2. Training / Test data
3. Topology
4. Training algorithm
5. Hyperparameters

# Deep Learning Recipe

1. Loss function                      softmax loss
2. Training / Test data              MNIST and CIFAR-10
3. Topology
4. Training algorithm
5. Hyperparameters

# Deep Learning Recipe



Iterations  
000,077

# HYPERPARAMETERS

Regularization rate

0

Problem type

Classification

# TOPOLOGY

## DATA

Which dataset do you want to use?



Ratio of training to test data: 50%



Noise: 25



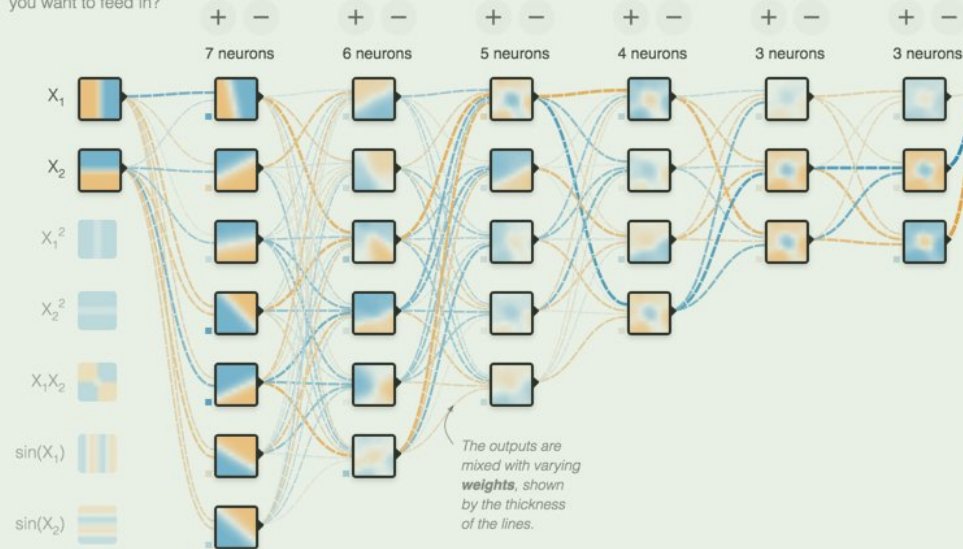
Batch size: 10



# DATA

## FEATURES

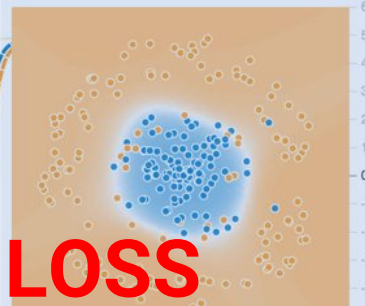
Which properties do you want to feed in?



## OUTPUT

Test loss 0.087

Training loss 0.115



# LOSS FUNCTION

Colors shows data, neuron and weight values.



☐ Show test data

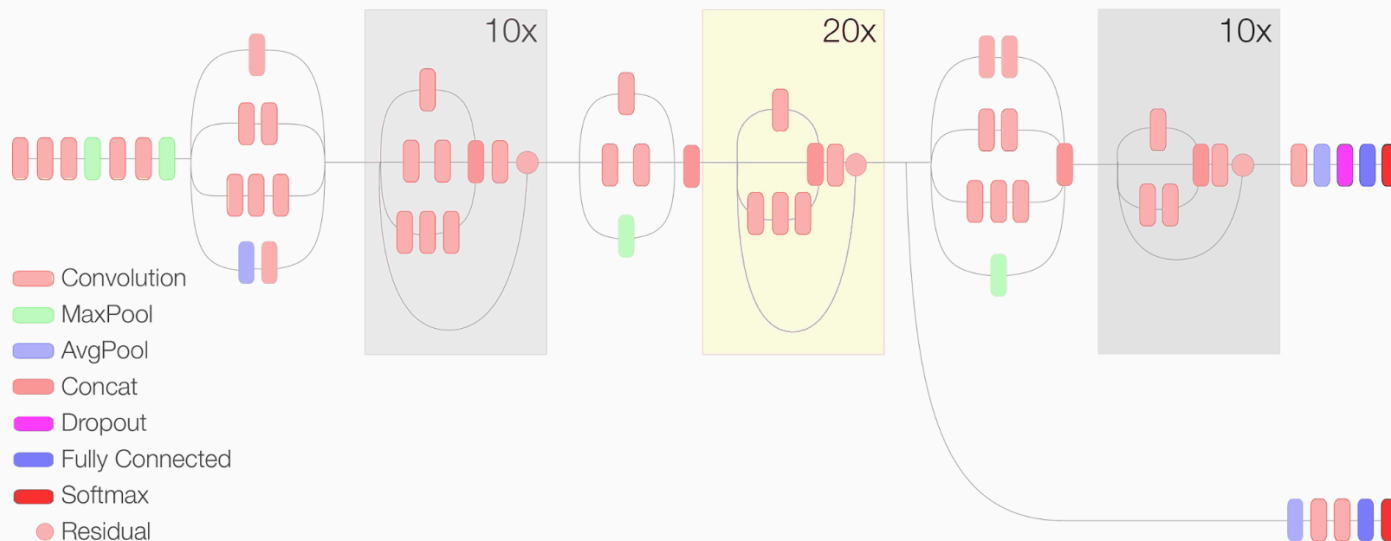
☐ Discretize output

# Layered Neural Network

## Inception Resnet V2 Network



## Compressed View





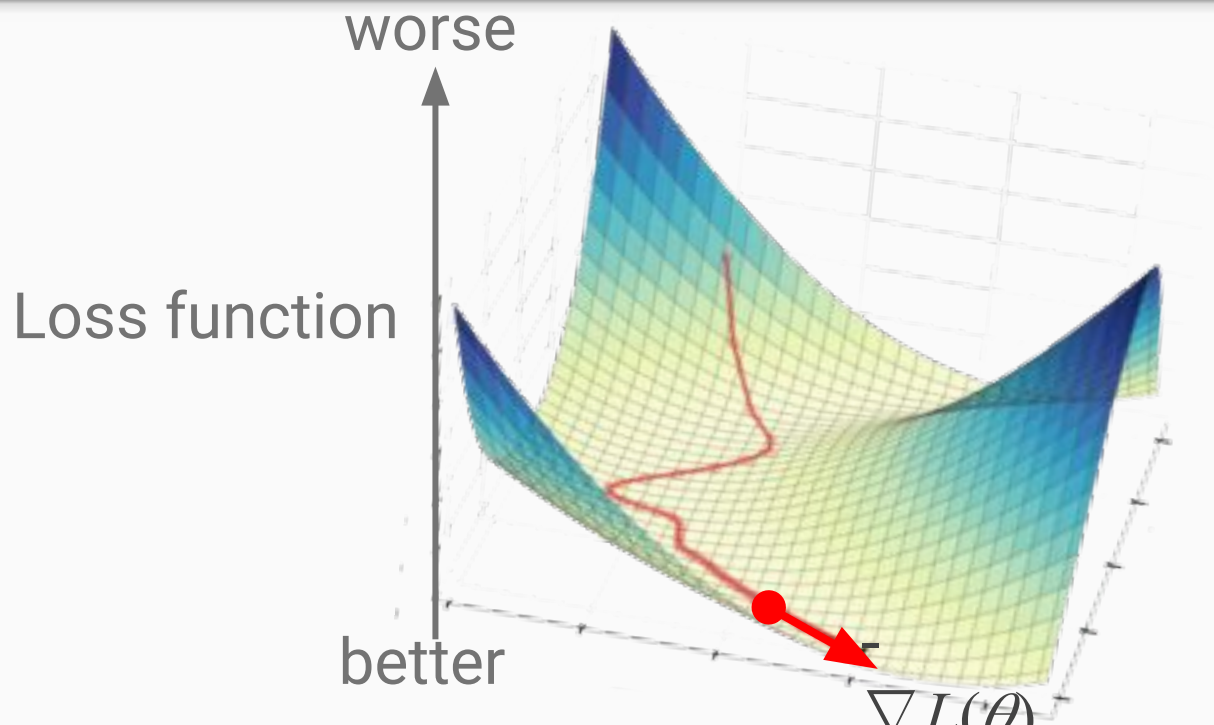
# Deep Learning Recipe

- |                         |                    |
|-------------------------|--------------------|
| 1. Loss function        | softmax loss       |
| 2. Training / Test data | MNIST and CIFAR-10 |
| 3. Topology             | neural network     |
| 4. Training algorithm   |                    |
| 5. Hyperparameters      |                    |

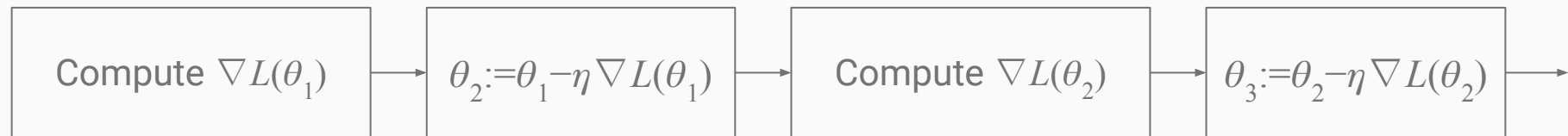
# Deep Learning Recipe

- |                         |                    |
|-------------------------|--------------------|
| 1. Loss function        | softmax loss       |
| 2. Training / Test data | MNIST and CIFAR-10 |
| 3. Topology             | neural network     |
| 4. Training algorithm   | SGD                |
| 5. Hyperparameters      |                    |

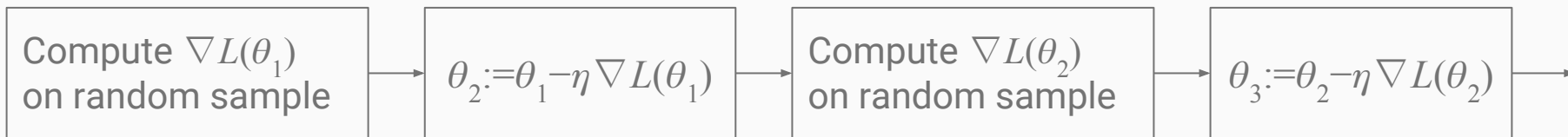
# Gradient Descent



# Gradient Descent



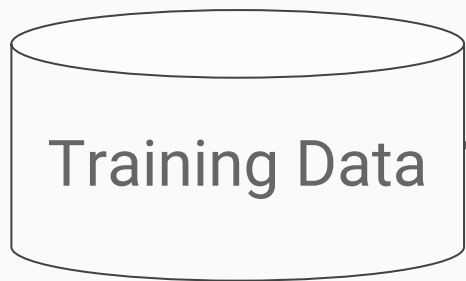
# Stochastic Gradient Descent



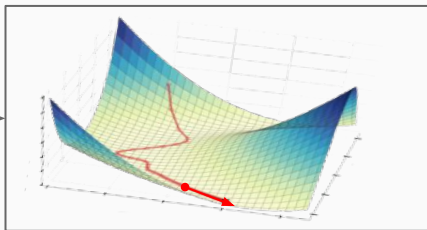
# Deep Learning Recipe

- |                         |                     |
|-------------------------|---------------------|
| 1. Loss function        | softmax loss        |
| 2. Training / Test data | MNIST and CIFAR-10  |
| 3. Topology             | neural network      |
| 4. Training algorithm   | SGD                 |
| 5. Hyperparameters      | tune experimentally |

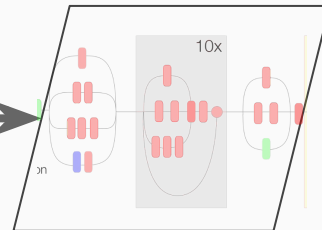




SGD



Model



# Differential Privacy

# Differential Privacy

$(\epsilon, \delta)$ -Differential Privacy: The distribution of the output  $M(D)$  on database  $D$  is (nearly) the same as  $M(D')$ :

$$\forall S: \quad \Pr[M(D) \in S] \leq \exp(\epsilon) \cdot \Pr[M(D') \in S] + \delta.$$

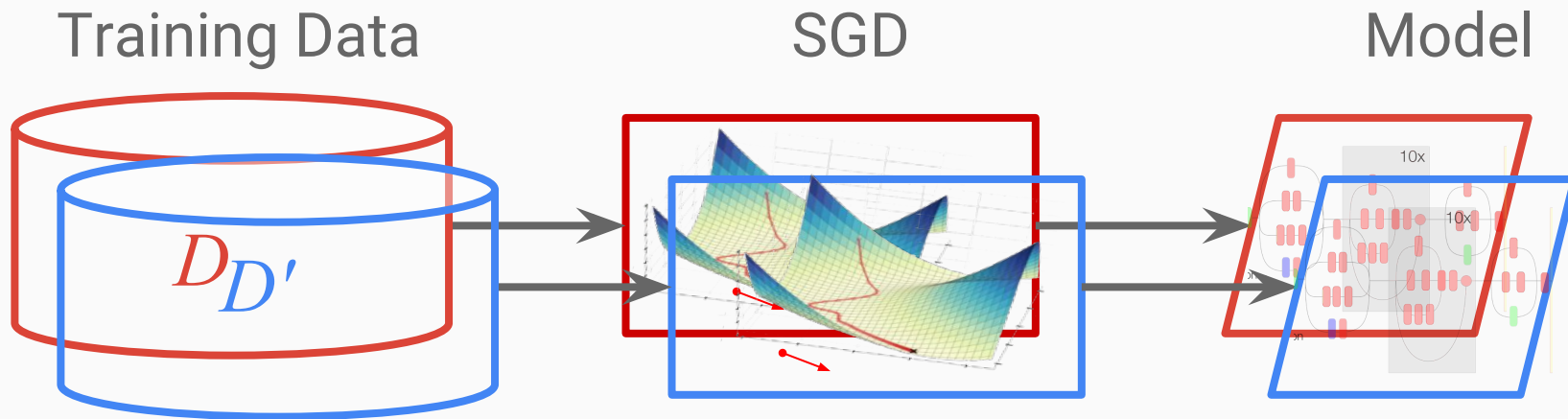
quantifies information leakage



allows for a small probability of failure



# Interpreting Differential Privacy



# Differential Privacy: Gaussian Mechanism

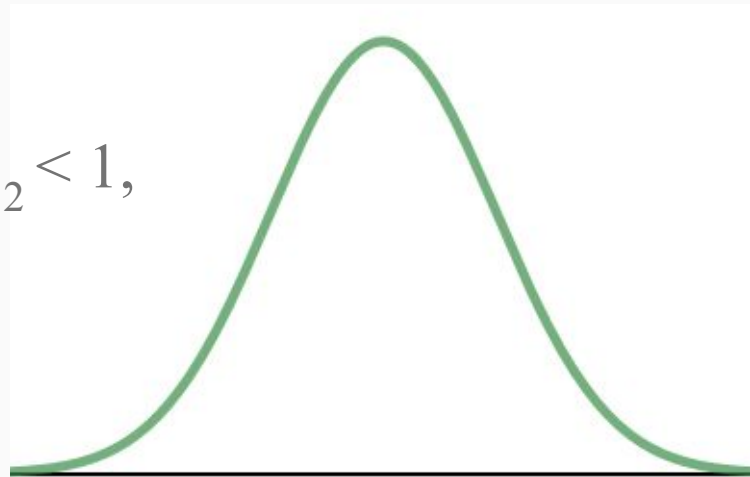
If  $\ell_2$ -sensitivity of  $f: \mathcal{D} \rightarrow \mathbb{R}^n$ :

$$\max_{D, D'} \|f(D) - f(D')\|_2 < 1,$$

then the Gaussian mechanism

$$f(D) + N^n(0, \sigma^2)$$

offers  $(\epsilon, \delta)$ -differential privacy, where  $\delta \approx \exp(-(\epsilon\sigma)^2/2)$ .



# Simple Recipe

To compute  $f$  with differential privacy

1. Bound sensitivity of  $f$
2. Apply the Gaussian mechanism





# Basic Composition Theorem

If  $f$  is  $(\epsilon_1, \delta_1)$ -DP and  $g$  is  $(\epsilon_2, \delta_2)$ -DP, then

$f(D), g(D)$  is  $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -DP

# Simple Recipe for Composite Functions

To compute composite  $f$  with differential privacy

1. Bound sensitivity of  $f$ 's components
2. Apply the Gaussian mechanism to each component
3. Compute total privacy via the composition theorem



# Deep Learning with Differential Privacy

# Deep Learning

- |                         |                     |
|-------------------------|---------------------|
| 1. Loss function        | softmax loss        |
| 2. Training / Test data | MNIST and CIFAR-10  |
| 3. Topology             | neural network      |
| 4. Training algorithm   | SGD                 |
| 5. Hyperparameters      | tune experimentally |

# Our Datasets: “Fruit Flies of Machine Learning”

MNIST dataset:

70,000 images

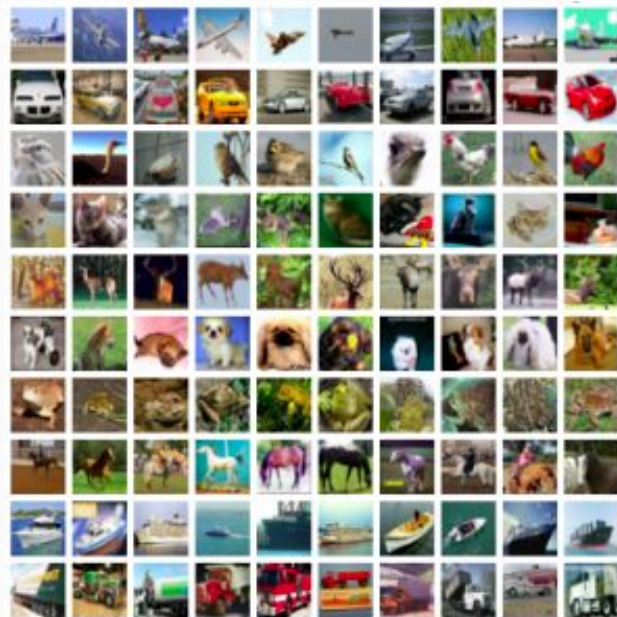
28×28 pixels each



CIFAR-10 dataset:

60,000 color images

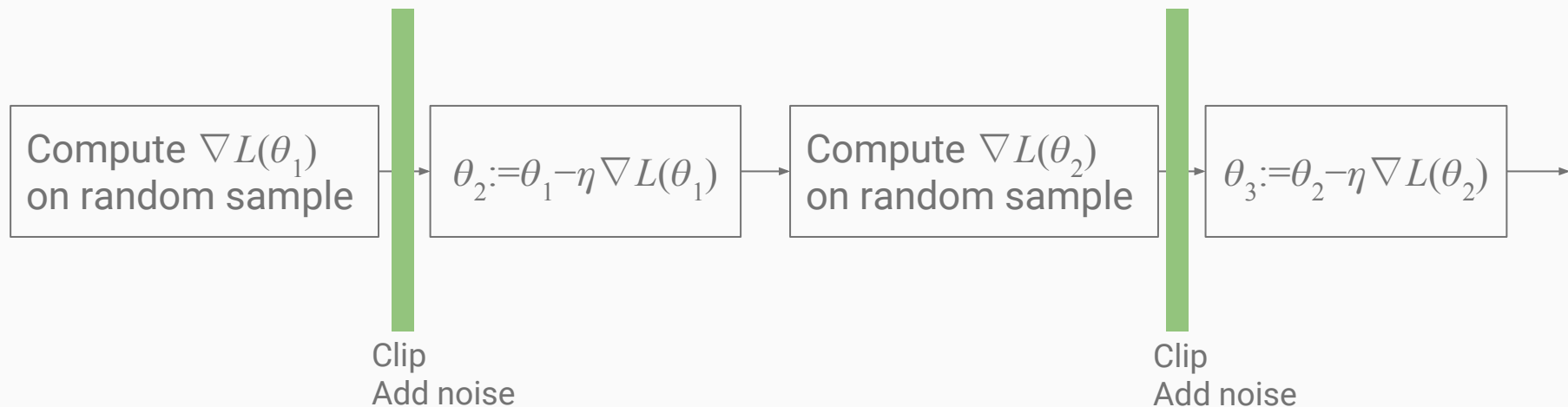
32×32 pixels each



# Differentially Private Deep Learning

- |                         |                      |
|-------------------------|----------------------|
| 1. Loss function        | softmax loss         |
| 2. Training / Test data | MNIST and CIFAR-10   |
| 3. Topology             | PCA + neural network |
| 4. Training algorithm   | SGD                  |
| 5. Hyperparameters      | tune experimentally  |

# Stochastic Gradient Descent with Differential Privacy



# Differentially Private Deep Learning

- |                         |                            |
|-------------------------|----------------------------|
| 1. Loss function        | softmax loss               |
| 2. Training / Test data | MNIST and CIFAR-10         |
| 3. Topology             | PCA + neural network       |
| 4. Training algorithm   | Differentially private SGD |
| 5. Hyperparameters      | tune experimentally        |



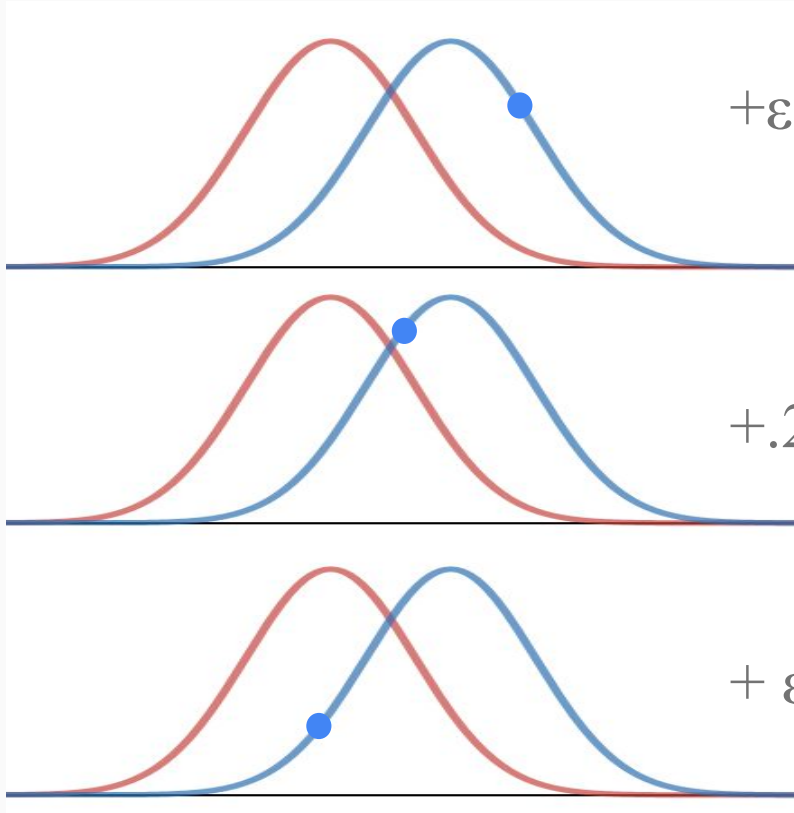
# Naïve Privacy Analysis

1. Choose  $\sigma = \frac{\sqrt{2 \log 1/\delta}}{\epsilon}$   $= 4$
2. Each step is  $(\epsilon, \delta)$ -DP  $(1.2, 10^{-5})$ -DP
3. Number of steps  $T$  10,000
4. Composition:  $(T\epsilon, T\delta)$ -DP 

$(12,000, .1)$ -DP

# Advanced Composition Theorems

# Composition theorem



$+\varepsilon$  for Blue

$+.2\varepsilon$  for Blue

$+\varepsilon$  for Red

“Heads, heads, heads”



Rosenkrantz: 78 in a row. A new record, I imagine.

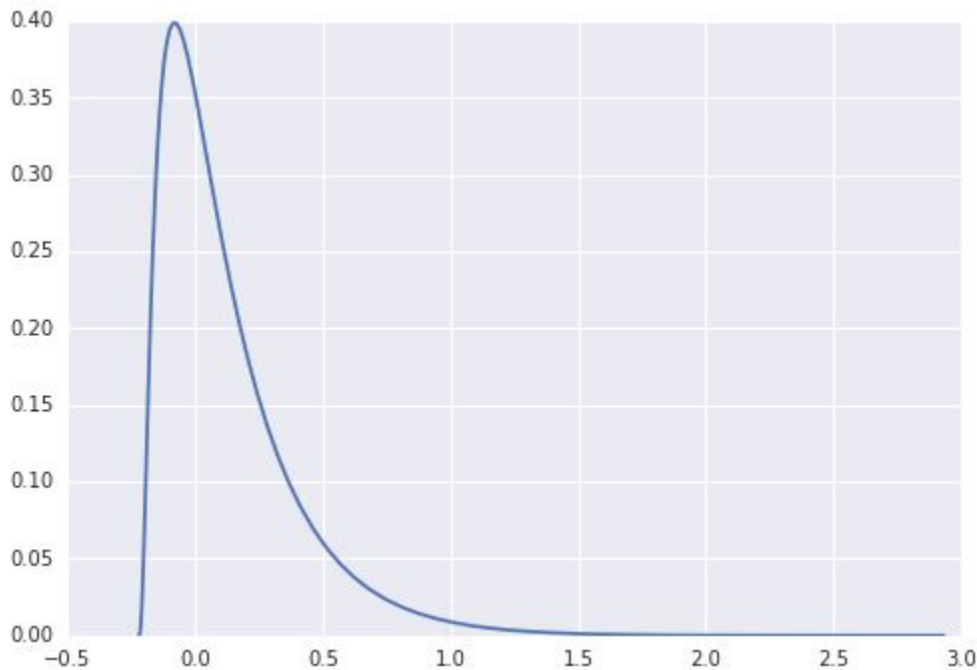
# Strong Composition Theorem

1. Choose  $\sigma = \frac{\sqrt{2 \log 1/\delta}}{\epsilon}$   $= 4$
2. Each step is  $(\epsilon, \delta)$ -DP  $(1.2, 10^{-5})$ -DP
3. Number of steps  $T$   $10,000$
4. Strong comp:  $(\epsilon\sqrt{T \log 1/\delta}, T\delta)$ -DP  $(360, .1)$ -DP

# Amplification by Sampling

1. Choose  $\sigma = \frac{\sqrt{2 \log 1/\delta}}{\varepsilon}$   $= 4$
2. Each batch is  $q$  fraction of data  $1\%$
3. Each step is  $(2q\varepsilon, q\delta)$ -DP  $(.024, 10^{-7})$ -DP
4. Number of steps  $T$   $10,000$
5. Strong comp:  $(2q\varepsilon\sqrt{T \log 1/\delta}, qT\delta)$ -DP  $(10, .001)$ -DP

# Privacy Loss Random Variable



$\log(\text{privacy loss})$

# Moments Accountant

1. Choose  $\sigma = \frac{\sqrt{2 \log 1/\delta}}{\epsilon}$   $= 4$
2. Each batch is  $q$  fraction of data  $1\%$
3. Keeping track of privacy loss's **moments**
4. Number of steps  $T$   $10,000$
5. Moments:  $(2q\epsilon\sqrt{T}, \delta)$ -DP 

$(1.25, 10^{-5})$ -DP



# Results

# Summary of Results

	Baseline
	no privacy
MNIST	98.3%
CIFAR-10	80%

# Summary of Results

	Baseline	[SS15]	[WKC+16]
	no privacy	reports $\epsilon$ per parameter	$\epsilon = 2$
MNIST	98.3%	98%	80%
CIFAR-10	80%		

# Summary of Results

	Baseline	[SS15]	[WKC+16]	this work		
	no privacy	reports $\epsilon$ per parameter	$\epsilon = 2$	$\epsilon = 8$ $\delta = 10^{-5}$	$\epsilon = 2$ $\delta = 10^{-5}$	$\epsilon = 0.5$ $\delta = 10^{-5}$
MNIST	98.3%	98%	80%	97%	95%	90%
CIFAR-10	80%			73%	67%	

# Contributions

- Differentially private deep learning applied to publicly available datasets and implemented in TensorFlow
  - <https://github.com/tensorflow/models>
- Innovations
  - Bounding sensitivity of updates
  - Moments accountant to keep tracking of privacy loss
- Lessons
  - Recommendations for selection of hyperparameters
- Full version: <https://arxiv.org/abs/1607.00133>