# Guarding user Privacy with Federated Learning and Differential Privacy

**Brendan McMahan**
mcmahan@google.com

# Federated Learning

Imbue **mobile devices** with **state of the art machine learning** systems **without centralizing data** and **with privacy** by default.

# Federated Learning

Imbue **mobile devices** with **state of the art machine learning** systems **without centralizing data** and **with privacy** by default.

## A very *personal* computer

2015: 79% away from phone ≤2 hours/day[1]
63% away from phone ≤1 hour/day
25% can't remember being away at all

2013: 72% of users within 5 feet of phone most of the time[2].

## Plethora of sensors

## Innumerable digital interactions

[1] 2015 Always Connected Research Report, IDC and Facebook
[2] 2013 Mobile Consumer Habits Study, Jumio and Harris Interactive.

# Federated Learning

Imbue **mobile devices** with **state of the art machine learning** systems **without centralizing data** and **with privacy** by default.

## Deep Learning

non-convex

millions of parameters

complex structure (eg LSTMs)

# Federated Learning

Imbue **mobile devices** with **state of the art machine learning** systems **without centralizing data** and **with privacy** by default.

## Distributed learning problem

Horizontally partitioned

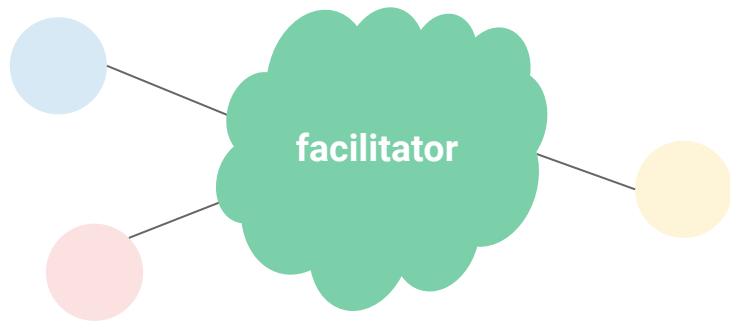Nodes: millions to billions

Dimensions: thousands to millions

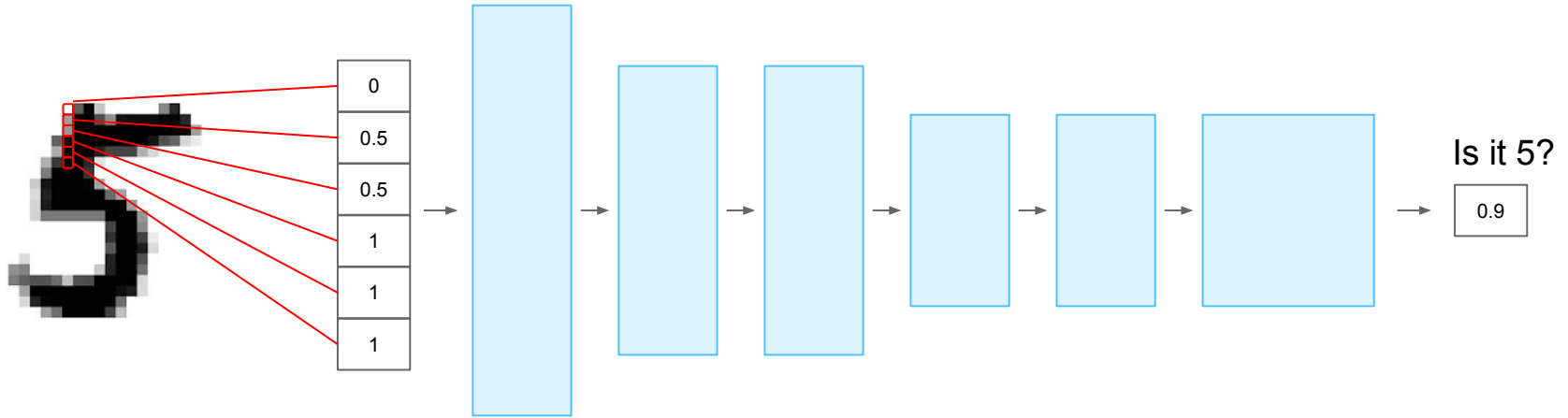Examples: millions to billions

# Federated Learning

Imbue **mobile devices** with **state of the art machine learning** systems **without centralizing data** and **with privacy** by default.
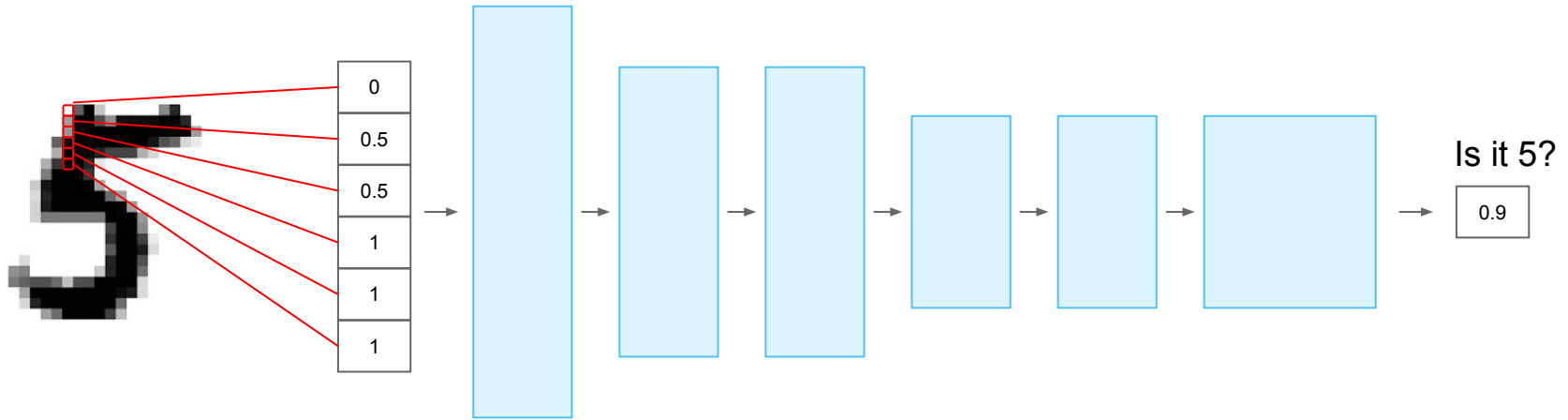
*Federated* decentralization

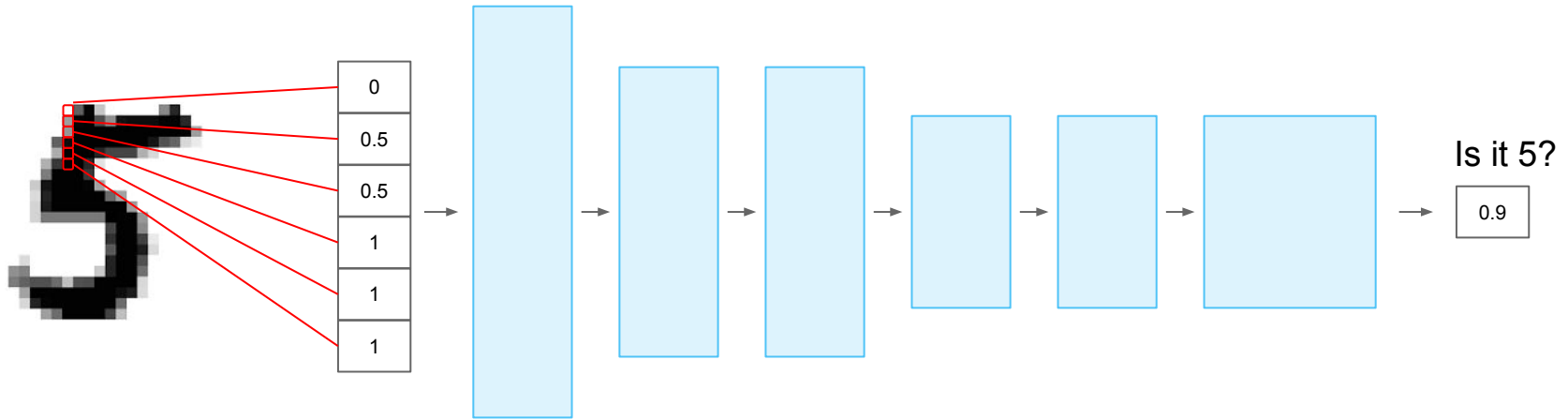# Deep Learning, the short short version



**f**(input, parameters) = output

# Deep Learning, the short short version



Is it 5?

**f**(input, parameters) = output

**loss**(parameters) = 1/n $\sum_i$ difference(**f**(input$_i$, parameters), desired$_i$)

Google

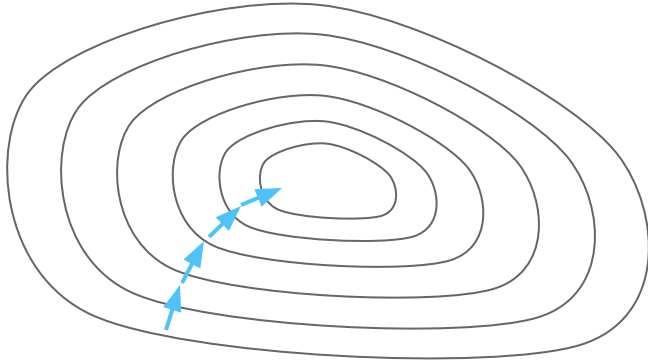# Deep Learning, the short short version



**Adjust these**

$\mathbf{f}(\text{input}, \text{parameters}) = \text{output}$

$\mathbf{loss}(\text{parameters}) = 1/n \sum_i \text{difference}(\mathbf{f}(\text{input}_i, \text{parameters}), \text{desired}_i)$

**to minimize this**

Google

# Deep Learning, the short short version

**Stochastic** Choose a random subset
of training data
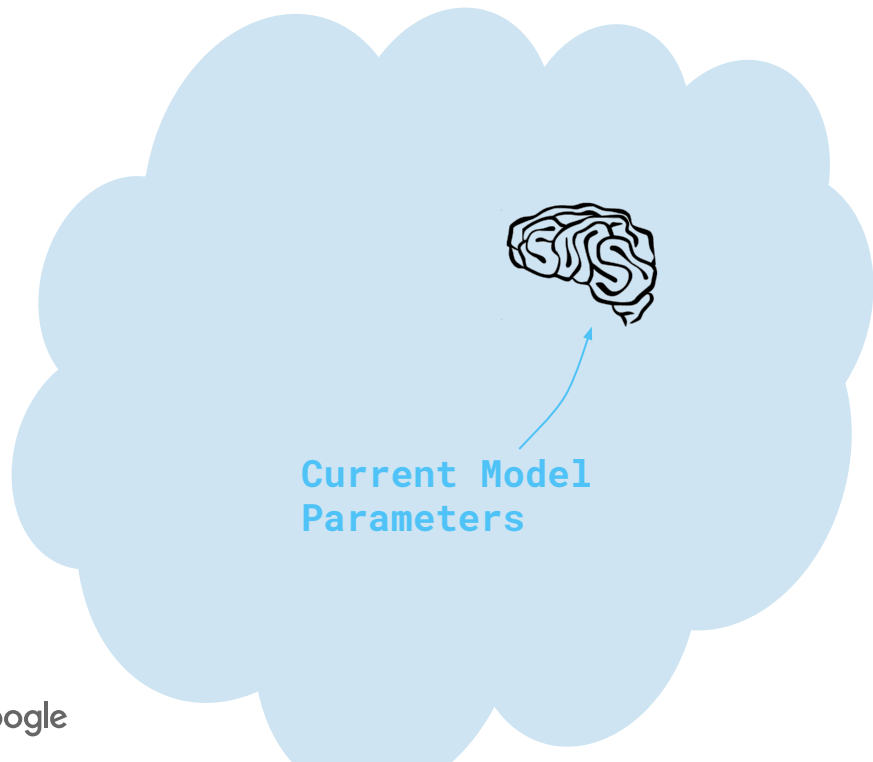
**Gradient** Compute the "down" direction
on the loss function

**Descent** Take a step in that direction

(Rinse & repeat)

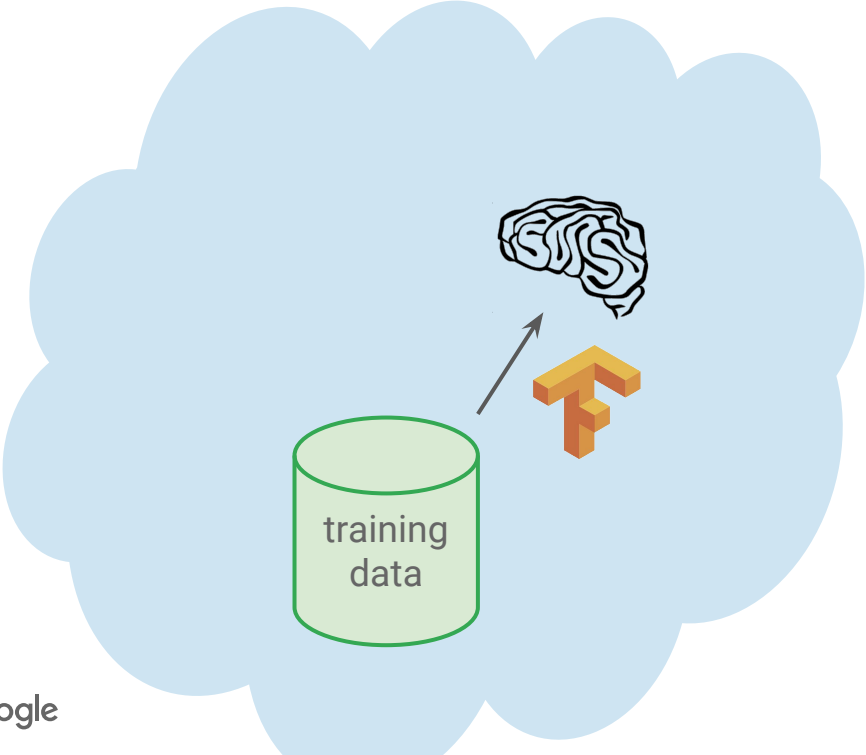$\mathbf{f}(\text{input}, \text{parameters}) = \text{output}$

$\mathbf{loss}(\text{parameters}) = 1/n \sum_i \text{difference}(\mathbf{f}(\text{input}_i, \text{parameters}), \text{desired}_i)$
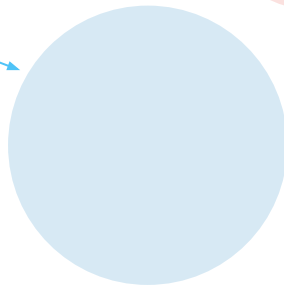
Google

# Cloud-centric ML for Mobile

# The model lives in the cloud.

**Current Model Parameters**

Google

# We train models in the cloud.



training
data

Google

Mobile Device

Current Model Parameters

Google

# Make predictions in the cloud.

request

prediction

Google

Gather training data
in the cloud.

request

prediction

training
data

Google

And make the models better.

training
data

Google

On-Device **Predictions**
(Inference)

**Instead of** making
predictions in the cloud

request

prediction

Google

Distribute the model,
make predictions on device.

Google

# On-device inference

## User Advantages

- Low latency
- Longer battery life
- Less wireless data transfer
- Better offline experience
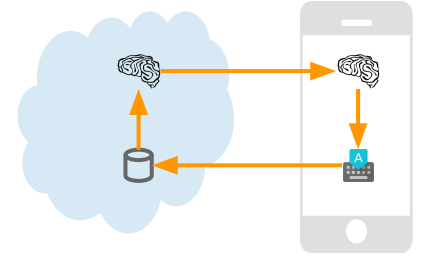- Less data sent to the cloud

## Developer Advantages

- Data is already localized
- New product opportunities

## World Advantages

- Raise privacy expectations for the industry



**1**

**On-Device Inference**

Google

# On-device training

## User Advantages

- Low latency
- Longer battery life
- Less wireless data transfer
- Better offline experience
- **Less data sent to the cloud (training data stays on device)**
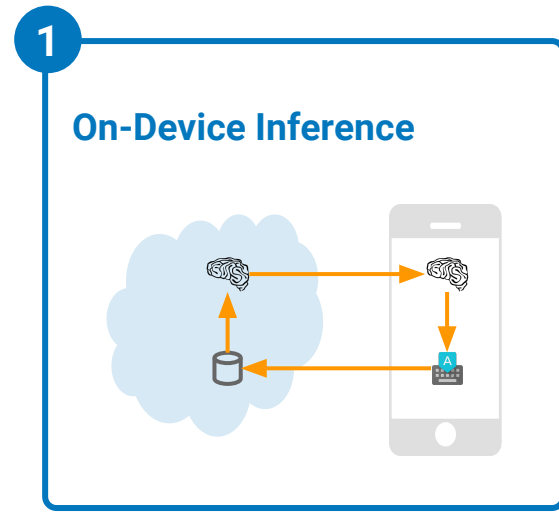
## Developer Advantages

- Data is already localized
- **New product opportunities**
- **Straightforward personalization**
- **Simple access to rich user context**

## World Advantages

- **Raise privacy expectations for the industry**

Bringing
**model training**
onto mobile devices.

**1**



**On-Device Inference**

Google

# On-device training

## User Advantages

- Low latency
- Longer battery life
- Less wireless data transfer
- Better offline experience
- **Less data sent to the cloud (training data stays on device)**

## Developer Advantages

- Data is already localized
- **New product opportunities**
- **Straightforward personalization**
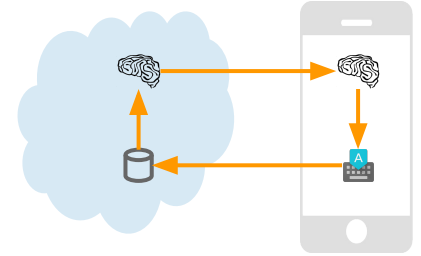- **Simple access to rich user context**

## World Advantages

- **Raise privacy expectations for the industry**

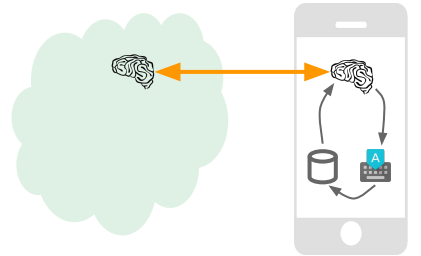Bringing
**model training**
onto mobile devices.



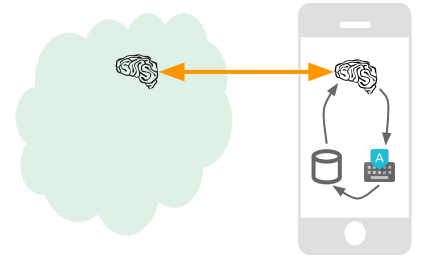**1 On-Device Inference**

**2 Federated Learning**

Google

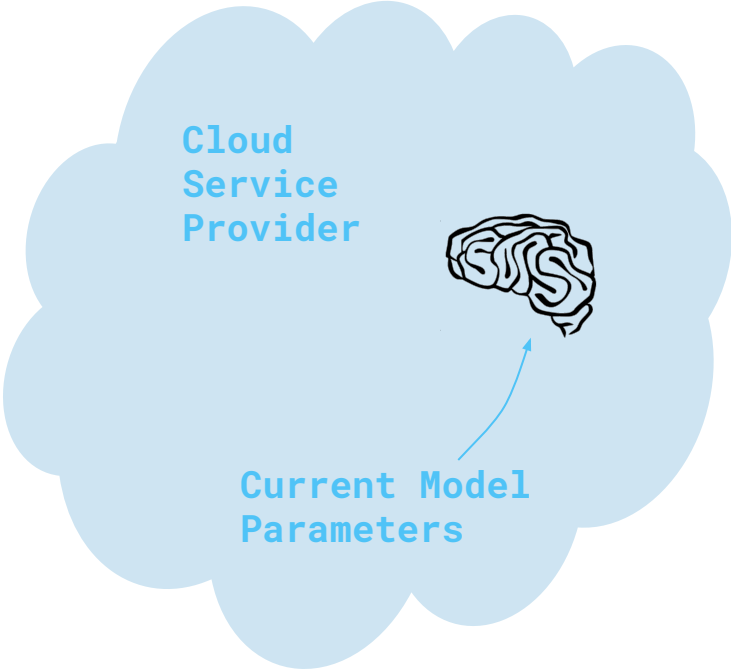# Federated Learning

# Federated Learning

*Federated Learning is the problem of training a shared global model under the coordination of a central server, from a federation of participating devices which maintain control of their own data.*
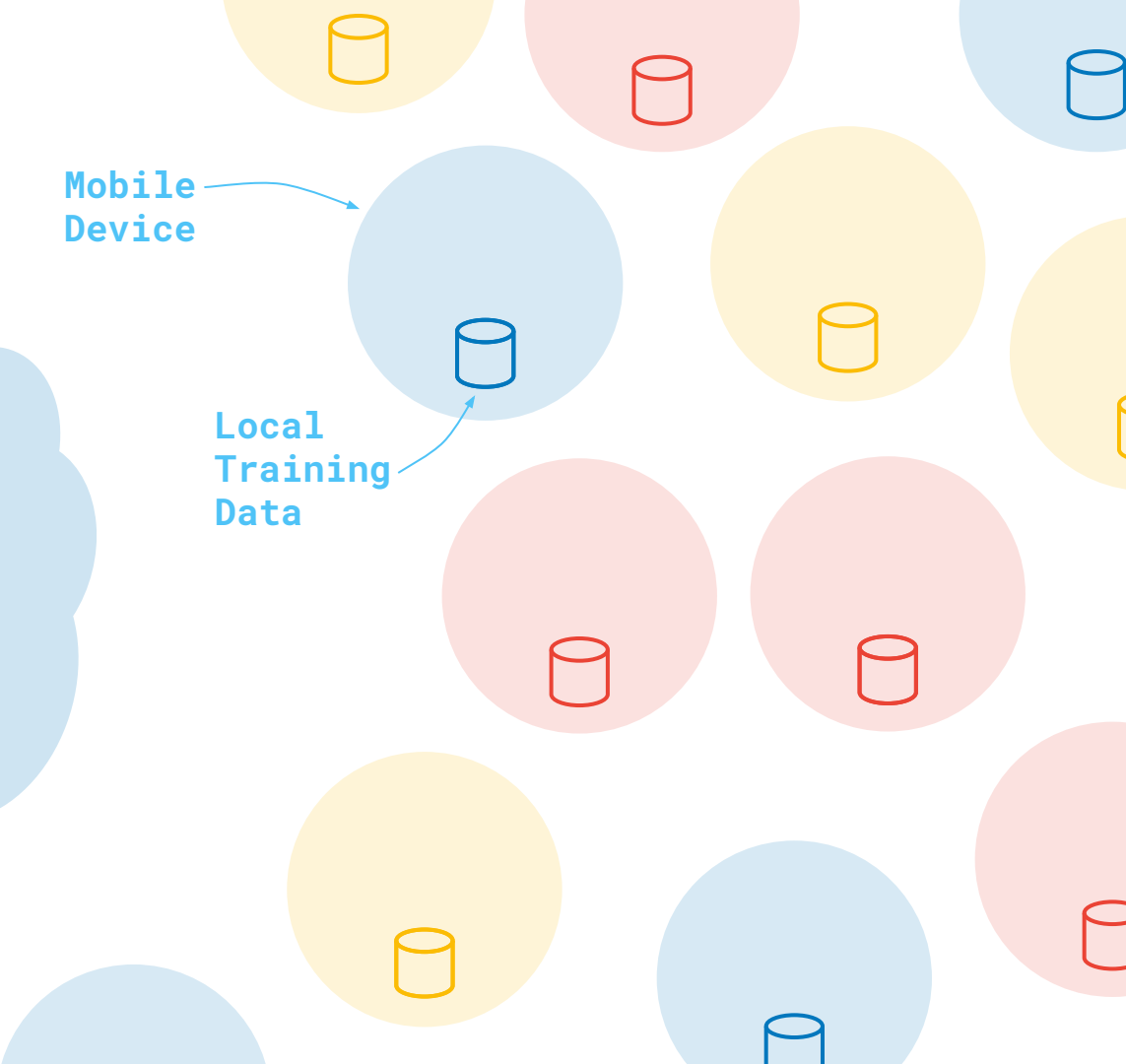
**2**

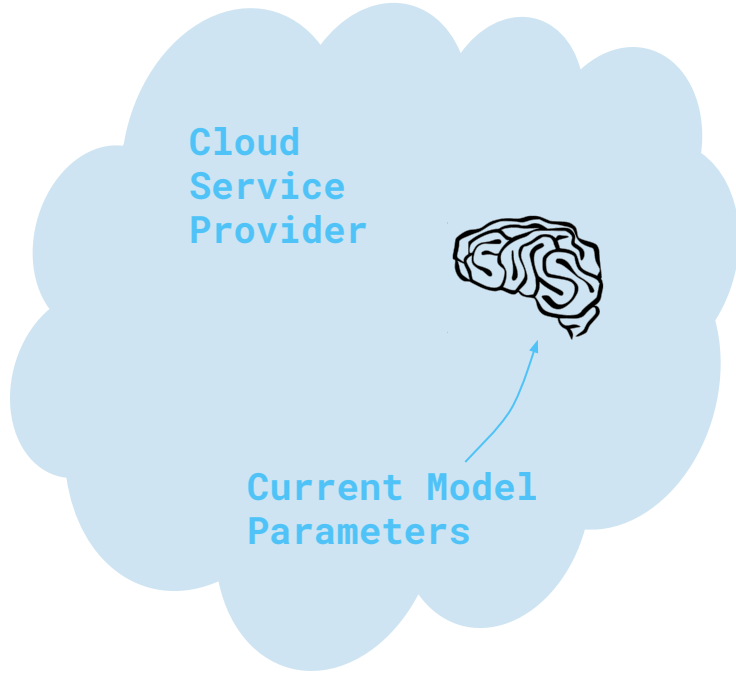**Federated Learning**
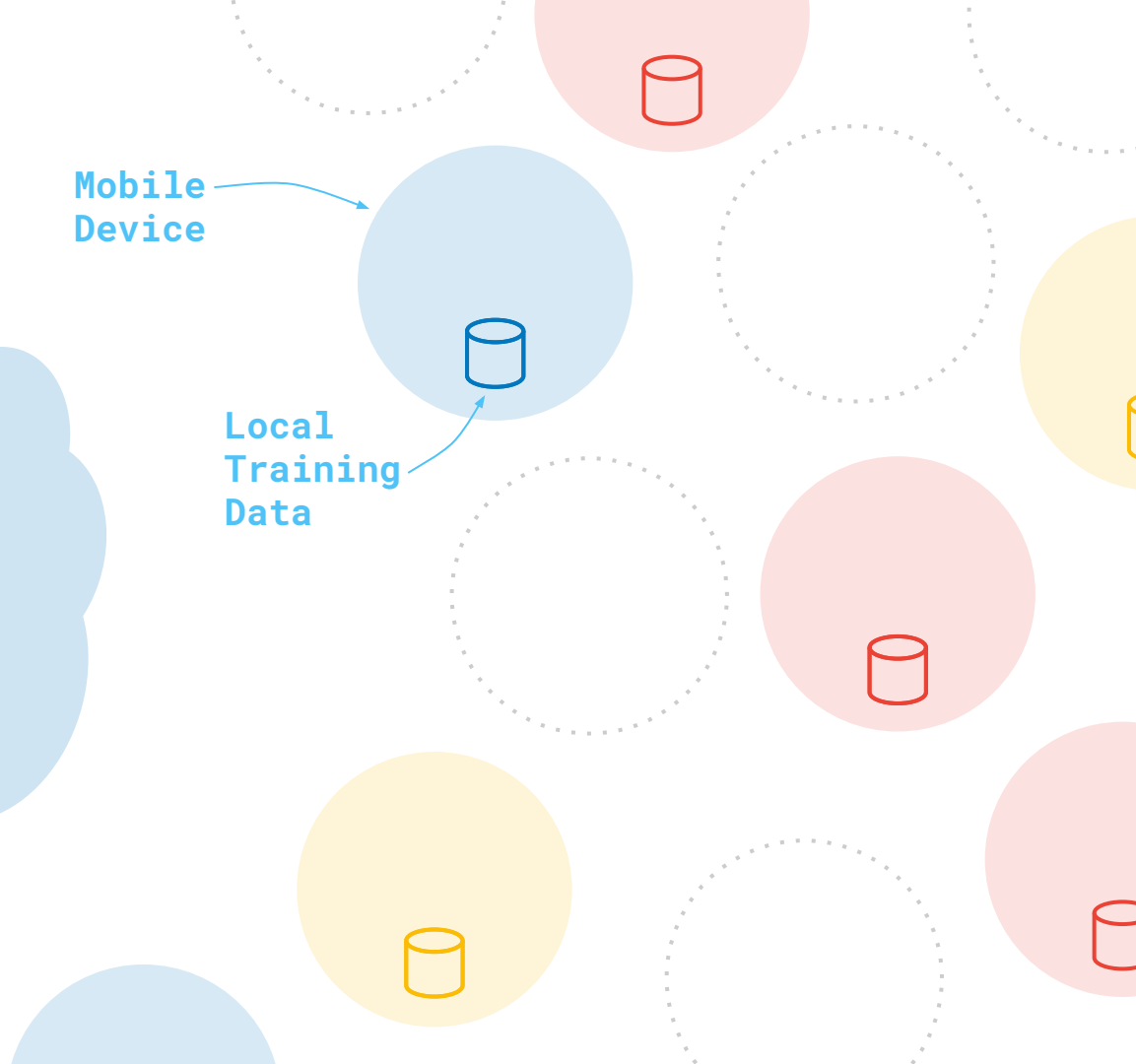
# Federated Learning



Mobile
Device

Local
Training
Data

Cloud
Service
Provider

Current Model
Parameters

Google

# Federated Learning

**Many devices will be offline.**
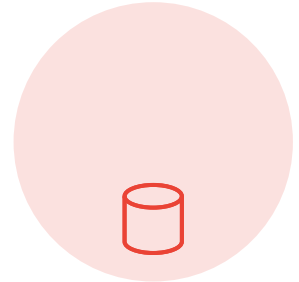
Mobile
Device

Local
Training
Data

Cloud
Service
Provider

Current Model
Parameters

Google

# Federated Learning

Mobile
Device

1. Server selects
a sample of e.g.
100 online
devices.

Local
Training
Data

Current Model
Parameters

Google

# Federated Learning

**Mobile Device**

**1. Server selects a sample of e.g. 100 online devices.**

**Current Model Parameters**

**Local Training Data**

Google

# Federated Learning

**2. Selected devices download the current model parameters.**

# Federated Learning



3. **Users compute an update using local training data**

Google

# Federated Learning



$\Sigma$

**4. Server aggregates users' updates into a new model.**

Google

Repeat until convergence.

# Applications of federating learning

**What makes a good application?**

- On-device data is more relevant than server-side proxy data

- On-device data is privacy sensitive or large

- Labels can be inferred naturally from user interaction

**Example applications**

- Language modeling (e.g., next word prediction) for mobile keyboards

- Image classification for predicting which photos people will share

- ...

Google

# Challenges of Federated Learning

**Massively Distributed**

Training data is stored across a very large number of devices

**Limited Communication**

Only a handful of rounds of unreliable communication with each devices

**Unbalanced Data**

Some devices have few examples, some have orders of magnitude more

**Highly Non-IID Data**

Data on each device reflects one individual's usage pattern

**Unreliable Compute Nodes**

Devices go offline unexpectedly; expect faults and adversaries

**Dynamic Data Availability**

The subset of data available is non-constant, e.g. time-of-day vs. country

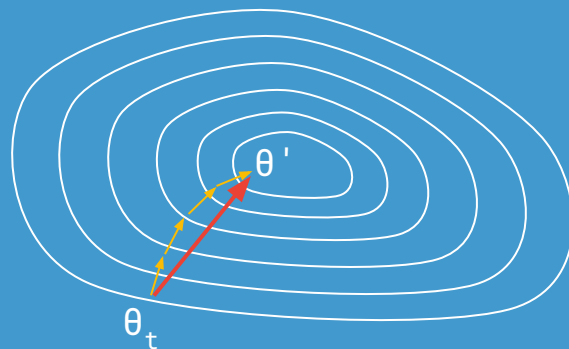Google

# The Federated Averaging algorithm

## *Server*

**Until Converged:**
1. Select a random subset (e.g. 100) of the (online) clients

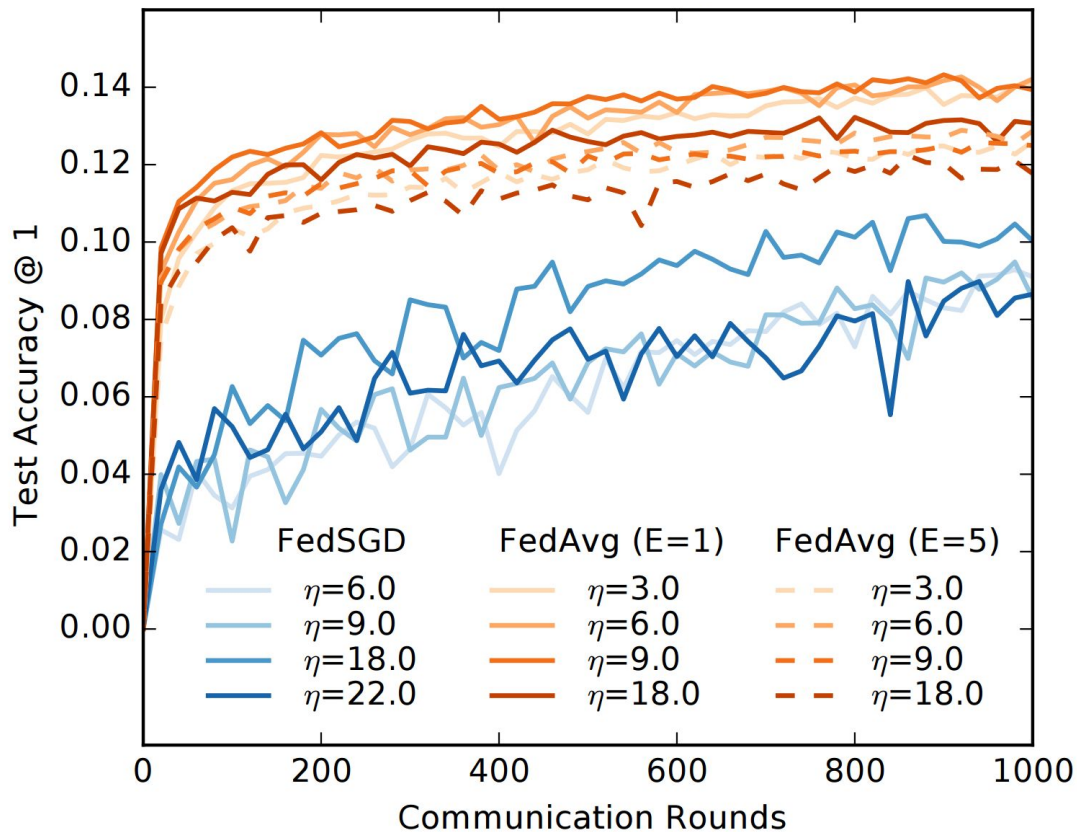2. In parallel, send current parameters $\theta_t$ to those clients

### *Selected Client k*

1. Receive $\theta_t$ from server.

2. Run some number of minibatch SGD steps, producing $\theta'$

3. Return $\theta'-\theta_t$ to server.



3. $\theta_{t+1} = \theta_t$ + data-weighted average of client updates

H. B. McMahan, *et al*.
**Communication-Efficient Learning of Deep Networks from Decentralized Data.** AISTATS 2017

# Large-scale LSTM for next-word prediction



Rounds to reach 10.5% Accuracy

FedSGD          820
FedAvg           35

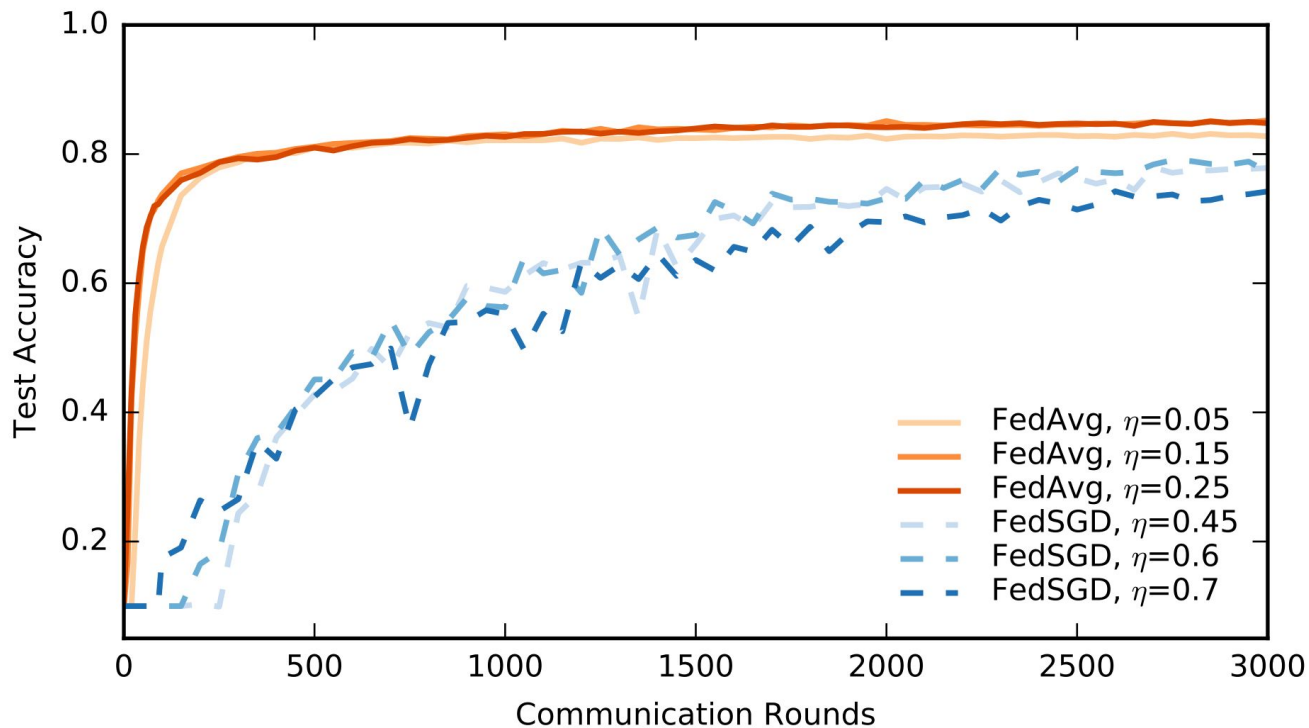**23x** decrease in communication rounds

**Model Details**
1.35M parameters
10K word dictionary
embeddings $\in \mathbb{R}^{96}$, state $\in \mathbb{R}^{256}$
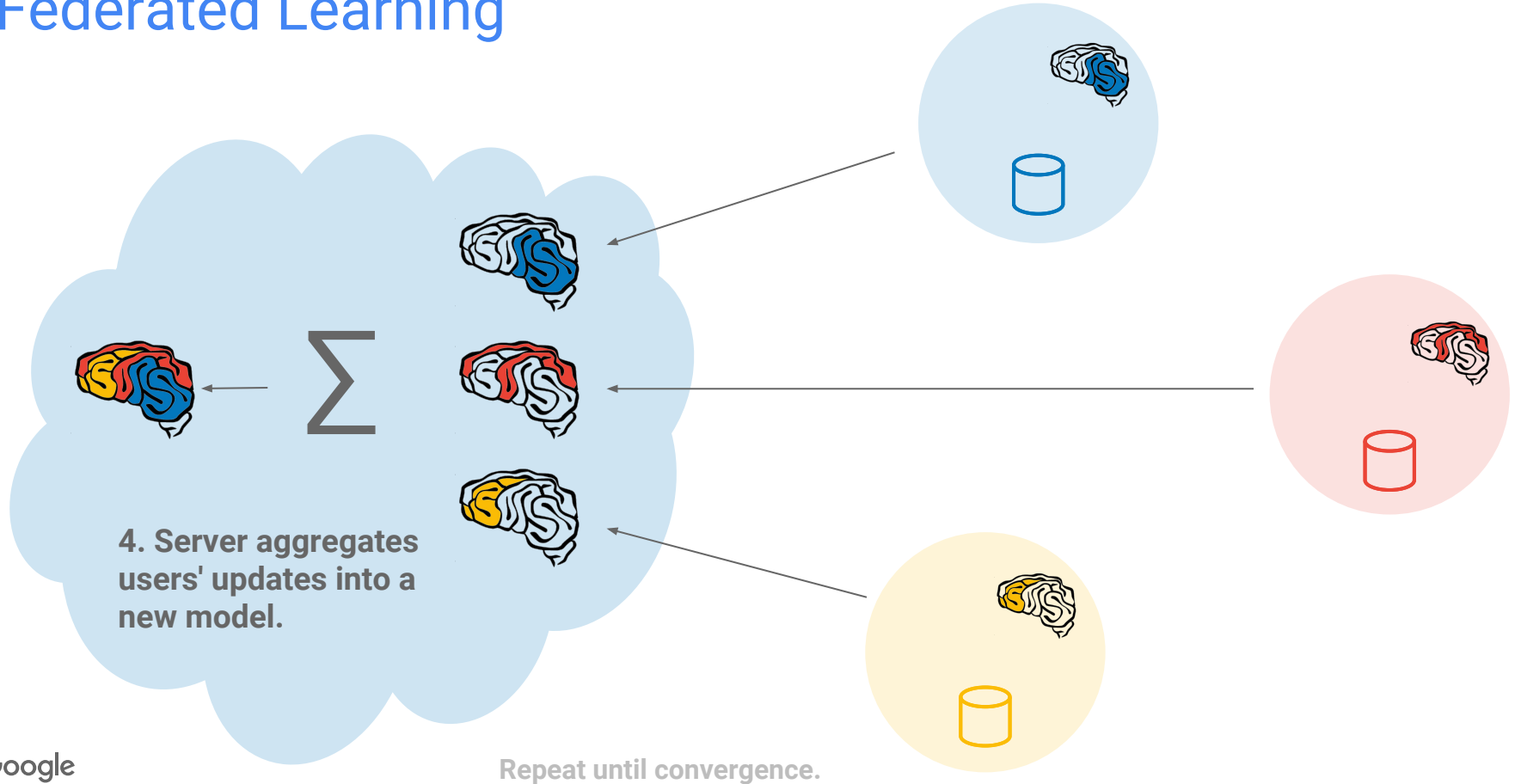corpus: Reddit posts, by author

# CIFAR-10 convolutional model



Updates to reach 82%
SGD            31,000
FedSGD         6,600
FedAvg         630

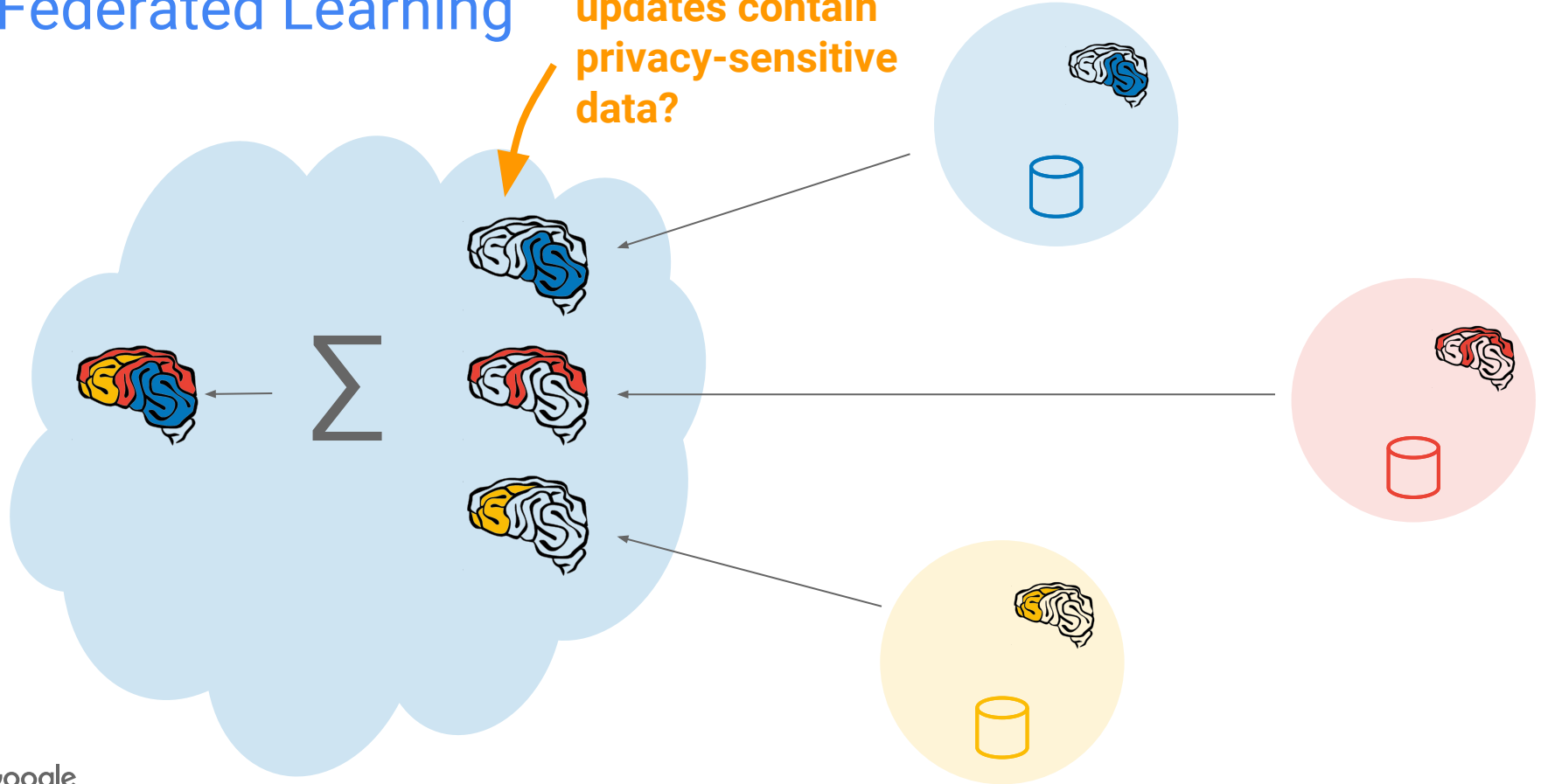**49x** decrease in communication (updates) vs SGD

(IID and balanced data)

Google

# Federated Learning
# &
# Privacy

# Federated Learning



$\sum$

**4. Server aggregates users' updates into a new model.**

Repeat until convergence.

Google

Federated Learning

Might these updates contain privacy-sensitive data?

$\sum$

Google

**Might these updates contain privacy-sensitive data?**

**Might these updates contain privacy-sensitive data?**

1.   **Ephemeral**

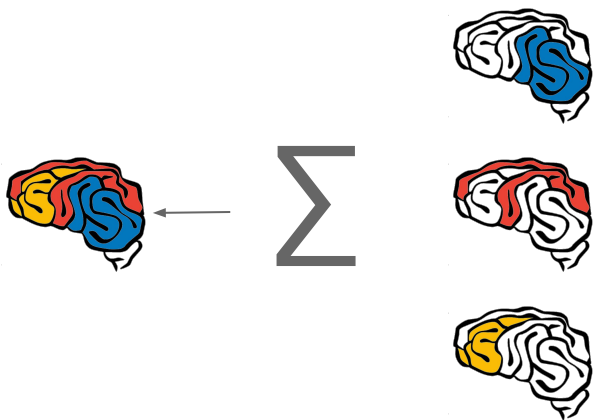**Might these updates contain privacy-sensitive data?**
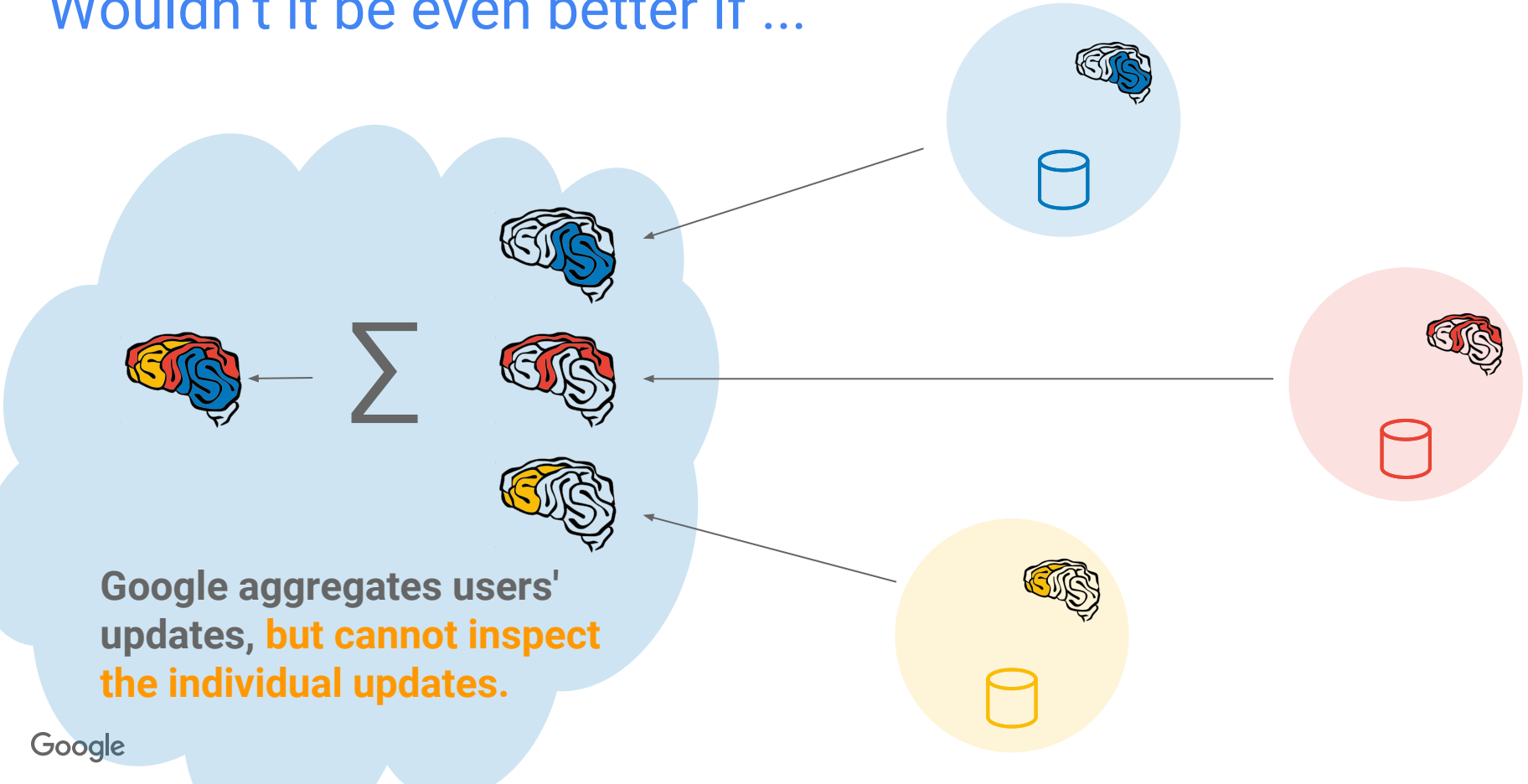
1. Ephemeral
2. **Focussed**

Improve privacy & security by minimizing the "attack surface"

Google

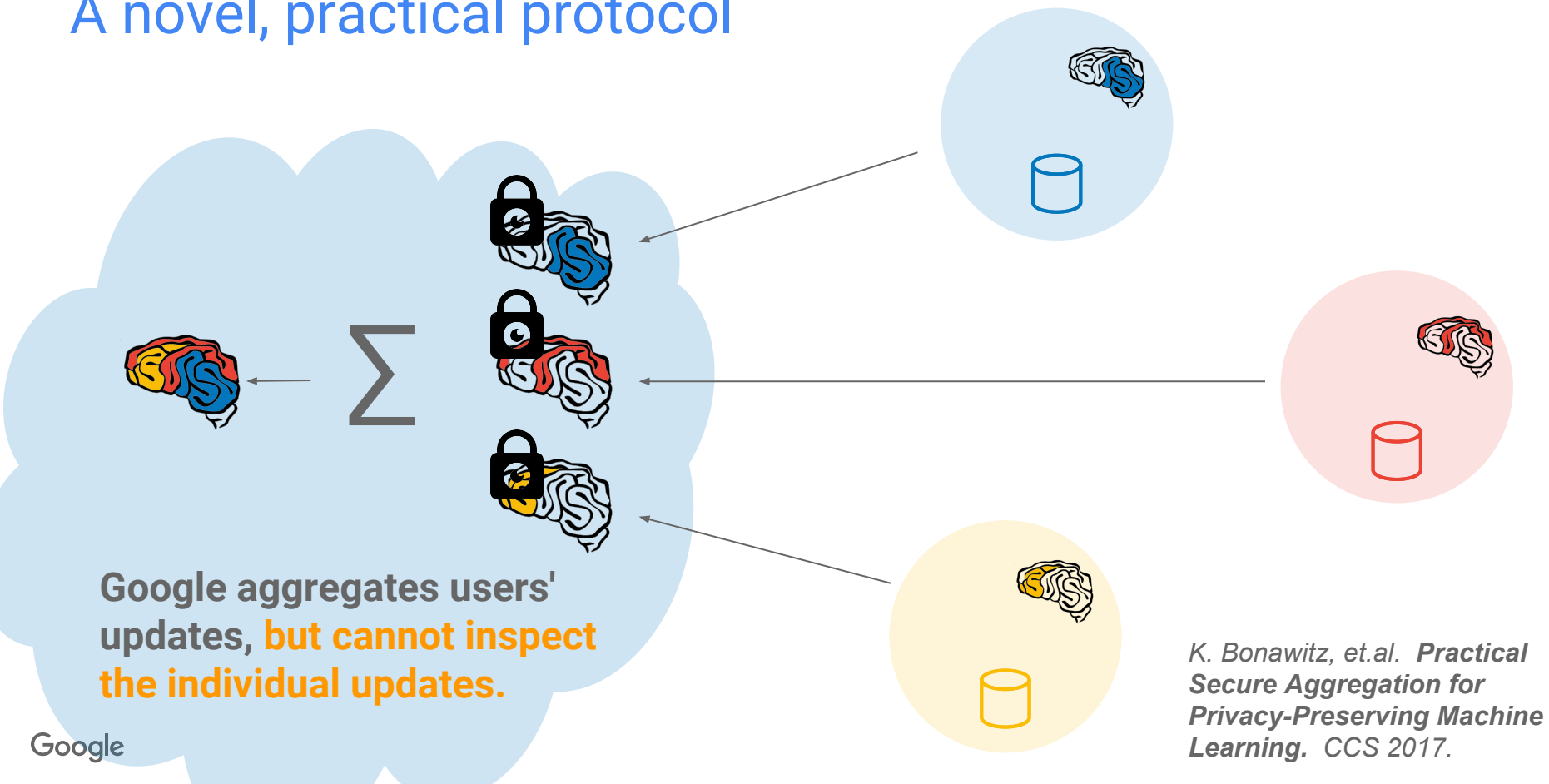**Might these updates contain privacy-sensitive data?**

$$\Sigma$$

1. Ephemeral

2. Focussed

3. **Only in aggregate**
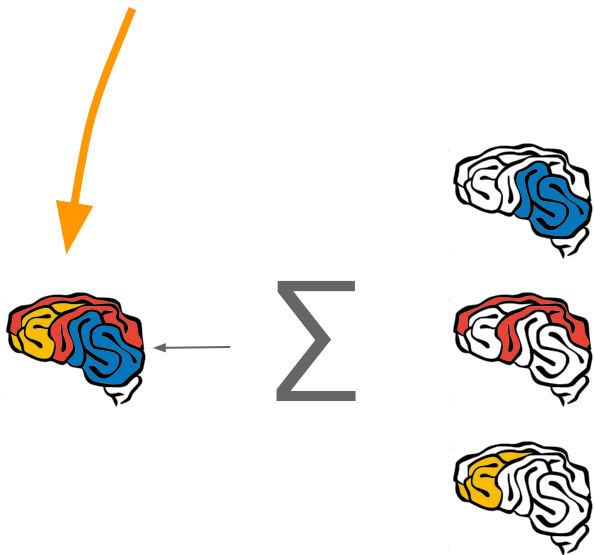
Google

# Wouldn't it be even better if ...



Google aggregates users' updates, **but cannot inspect the individual updates.**

Google

# A novel, practical protocol



**Google aggregates users' updates, but cannot inspect the individual updates.**

K. Bonawitz, et.al. *Practical Secure Aggregation for Privacy-Preserving Machine Learning.* CCS 2017.

Google

**Might the final model memorize a user's data?**
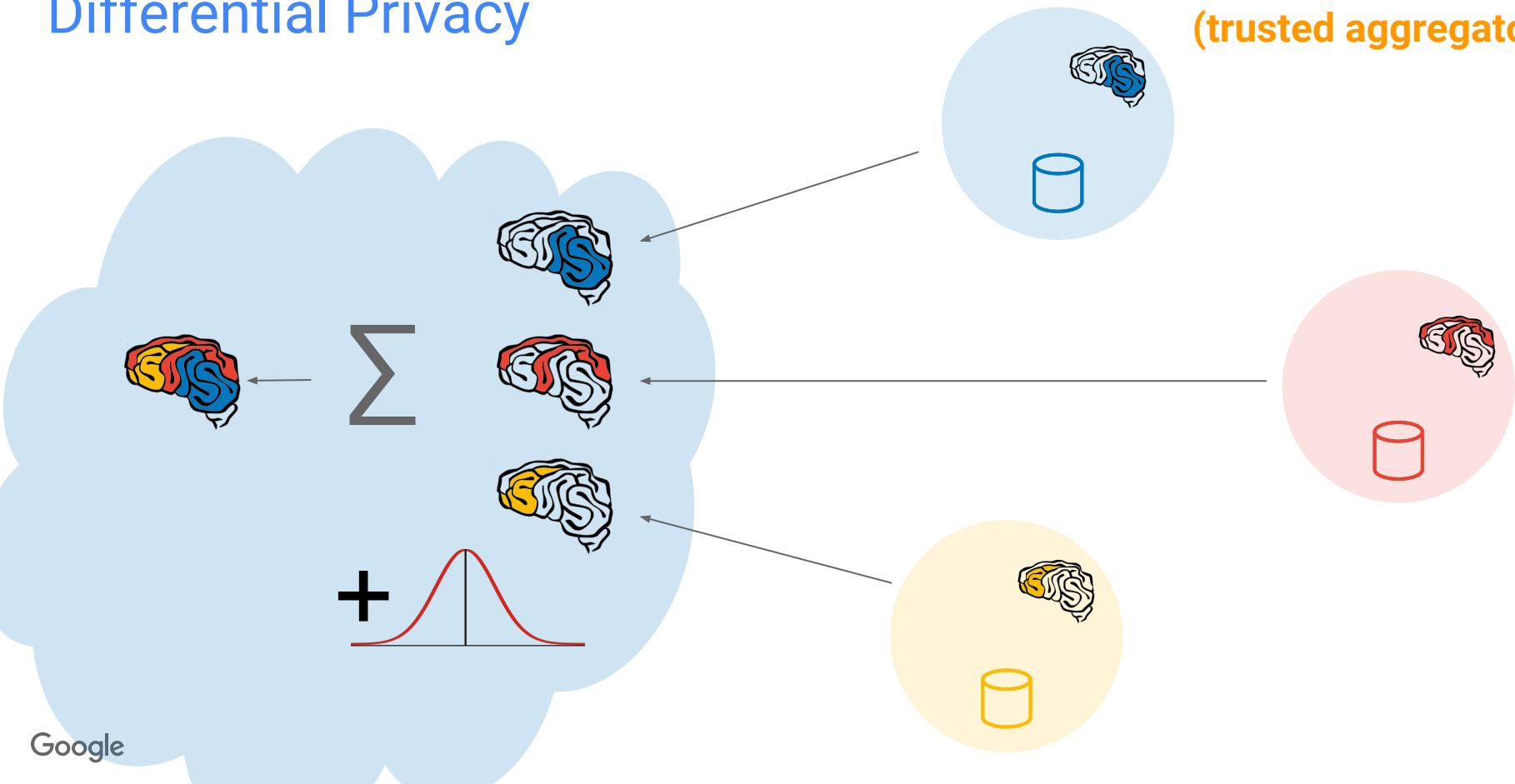
$$\sum$$

1. Ephemeral

2. Focussed

3. Only in aggregate

4. **Differentially private**

Google

# Differential Privacy

Differential Privacy

Differential Privacy
(trusted aggregator)
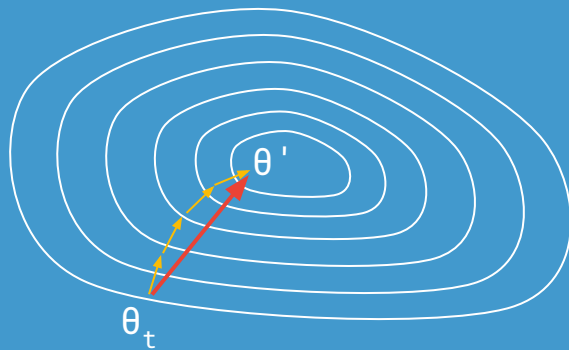
Google

# Federated Averaging

**Until Converged:**
1. Select a random subset (e.g. C=100) of the (online) clients

2. In parallel, send current parameters $\theta_t$ to those clients

### *Selected Client k*

1. Receive $\theta_t$ from server.

2. Run some number of minibatch SGD steps, producing $\theta'$

3. Return $\theta'-\theta_t$ to server.



3. $\theta_{t+1} = \theta_t$ + data-weighted average of client updates
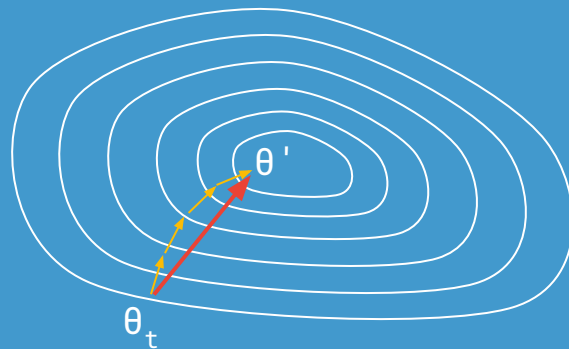
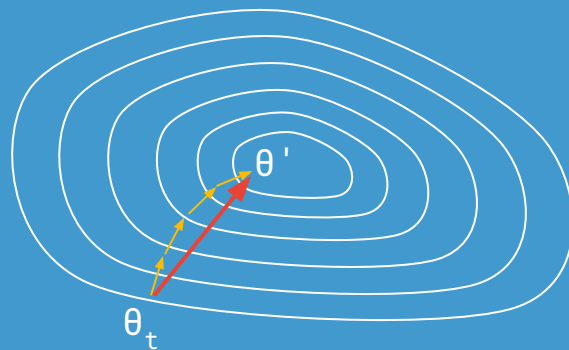# Differentially-Private Federated Averaging

## *Server*

**Until Converged:**
1. Select each user **independently** with **probability q**, for say E[C]=1000 clients

2. In parallel, send current parameters $\theta_t$ to those clients

### *Selected Client k*

1. Receive $\theta_t$ from server.

2. Run some number of minibatch SGD steps, producing $\theta'$

3. Return $\theta'-\theta_t$ to server.



3. $\theta_{t+1} = \theta_t$ + data-weighted average of client updates

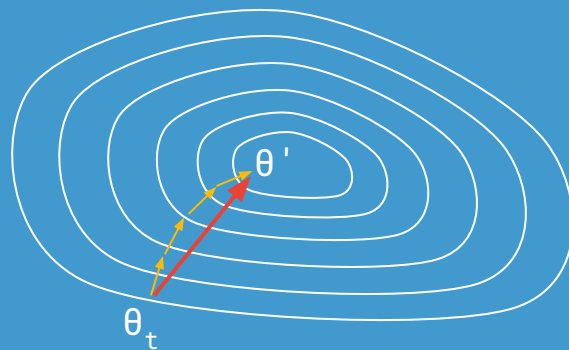# Differentially-Private Federated Averaging

## *Server*

**Until Converged:**
1. Select each user **independently** with **probability q**, for say E[C]=1000 clients

2. In parallel, send current parameters $\theta_t$ to those clients

### *Selected Client k*

1. Receive $\theta_t$ from server.

2. Run some number of minibatch SGD steps, producing θ'

3. Return **Clip($\theta'-\theta_t$)** to server.



3. $\theta_{t+1}$ = $\theta_t$ + data-weighted average of client updates

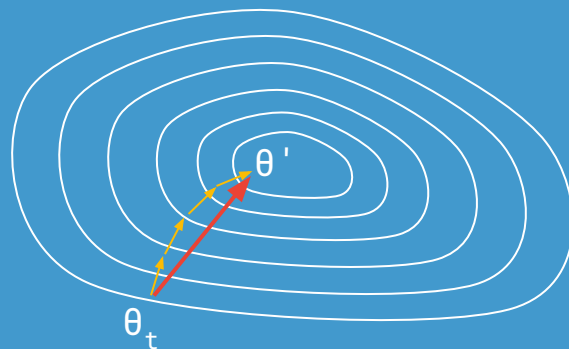# Differentially-Private Federated Averaging

## *Server*

**Until Converged:**
1. Select each user **independently** with **probability q**, for say E[C]=1000 clients

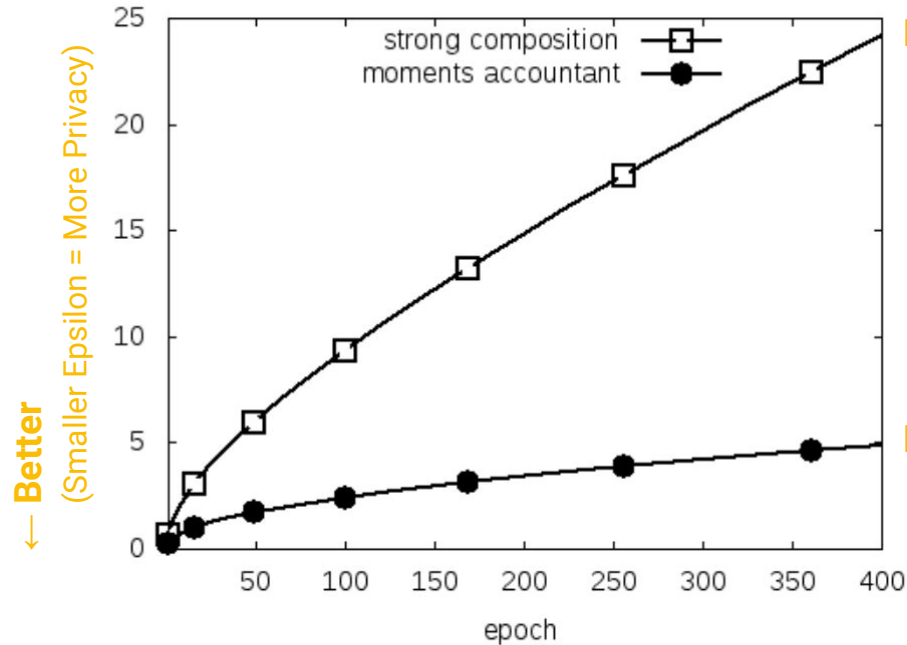2. In parallel, send current parameters $\theta_t$ to those clients

### *Selected Client k*

1. Receive $\theta_t$ from server.

2. Run some number of minibatch SGD steps, producing $\theta'$

3. Return **Clip($\theta'-\theta_t$)** to server.



3. $\theta_{t+1}$ = $\theta_t$ + **bounded sensitivity** data-weighted average of client updates

# Differentially-Private Federated Averaging

## *Server*

**Until Converged:**
1. Select each user **independently** with **probability q**, for say E[C]=1000 clients

2. In parallel, send current parameters $\theta_t$ to those clients

### *Selected Client k*

1. Receive $\theta_t$ from server.

2. Run some number of minibatch SGD steps, producing θ'

3. Return **Clip(θ'-$\theta_t$)** to server.



3. $\theta_{t+1}$ = $\theta_t$ + **bounded sensitivity** data-weighted average of client updates
   + Gaussian noise **N(0, I$\sigma^2$)**

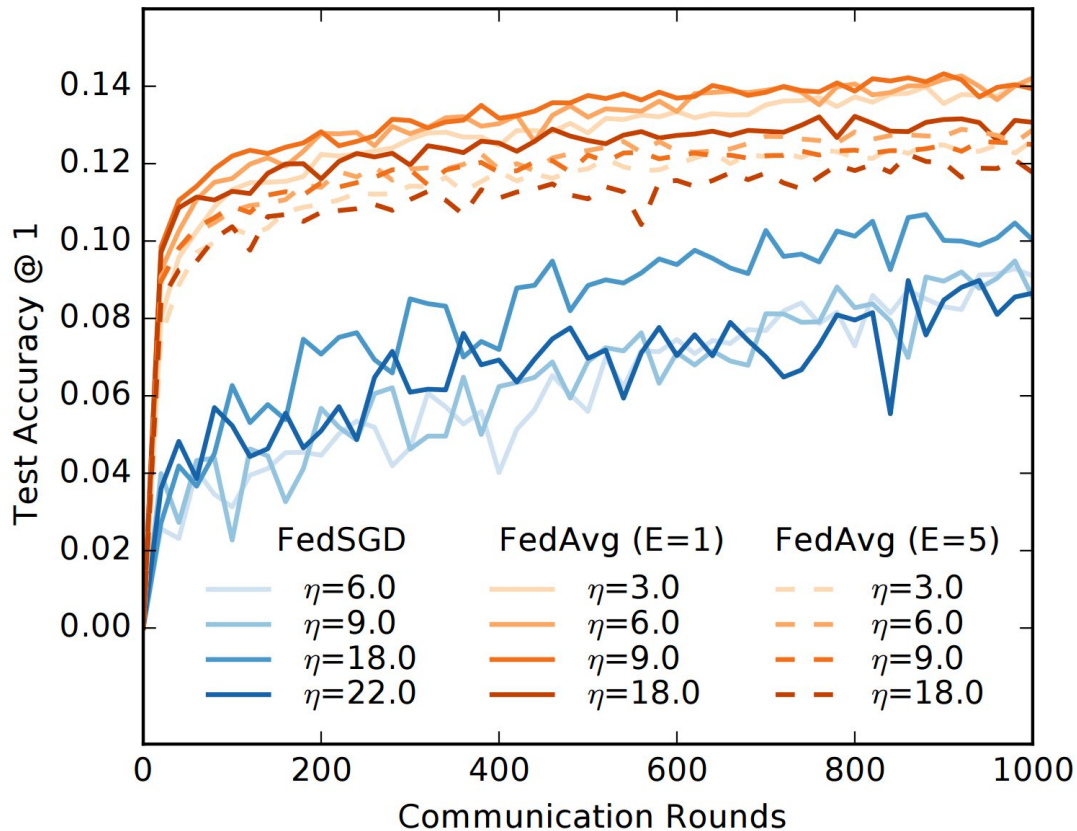# Privacy Accounting for Noisy SGD: Moments Accountant



**Previous composition theorems**

**Moments Accountant**

M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, & L. Zhang. **Deep Learning with Differential Privacy**. CCS 2016.

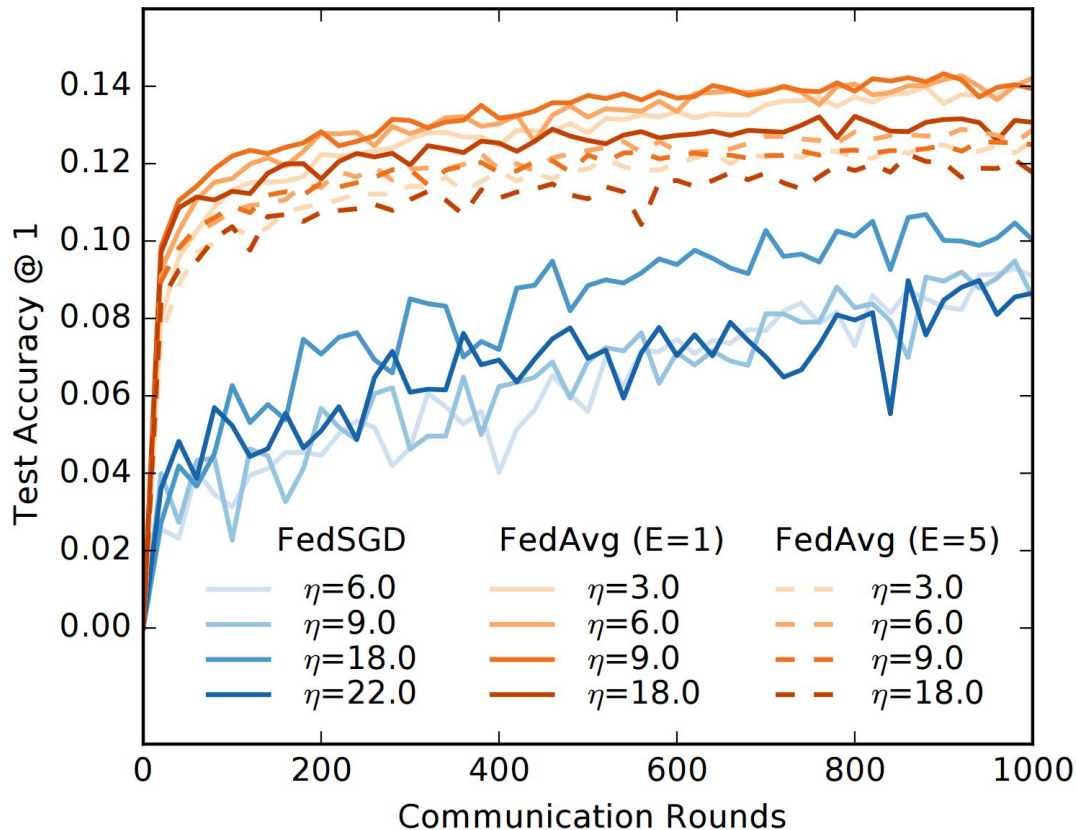Google

# Large-scale LSTM for next-word prediction



Rounds to reach 10.5% Accuracy

| | |
|---|---|
| FedSGD | 820 |
| FedAvg | 35 |

**23x** decrease in communication rounds

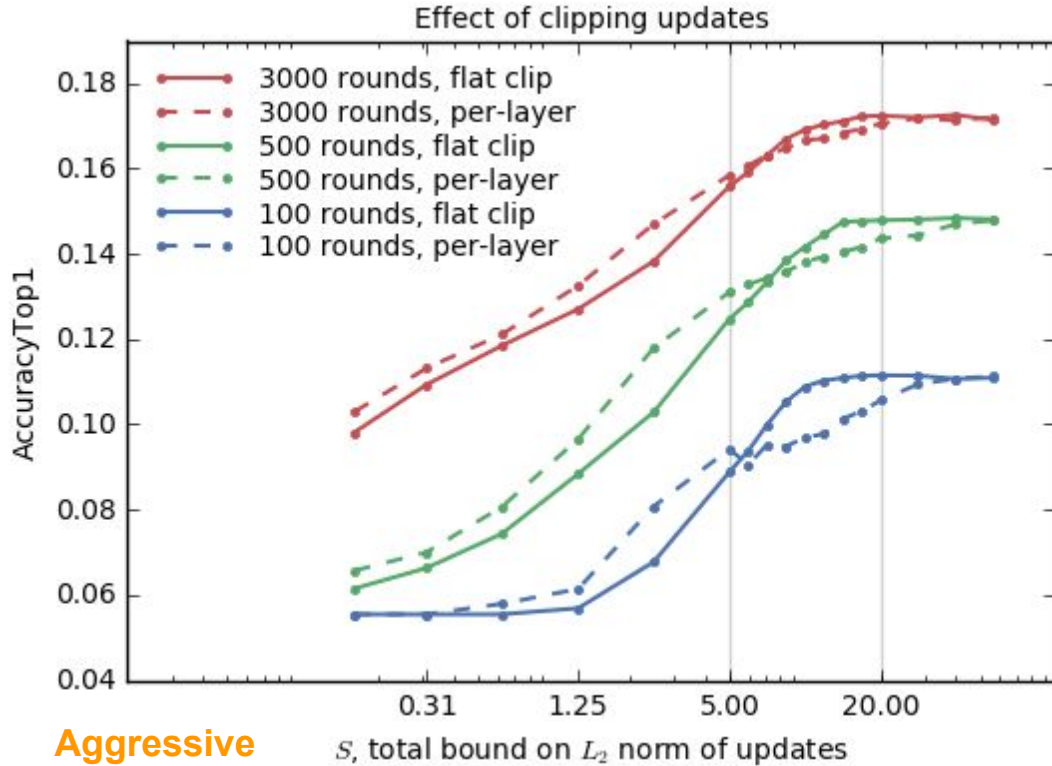# Large-scale LSTM for next-word prediction



Rounds to reach 10.5% Accuracy

FedSGD          820
FedAvg           35

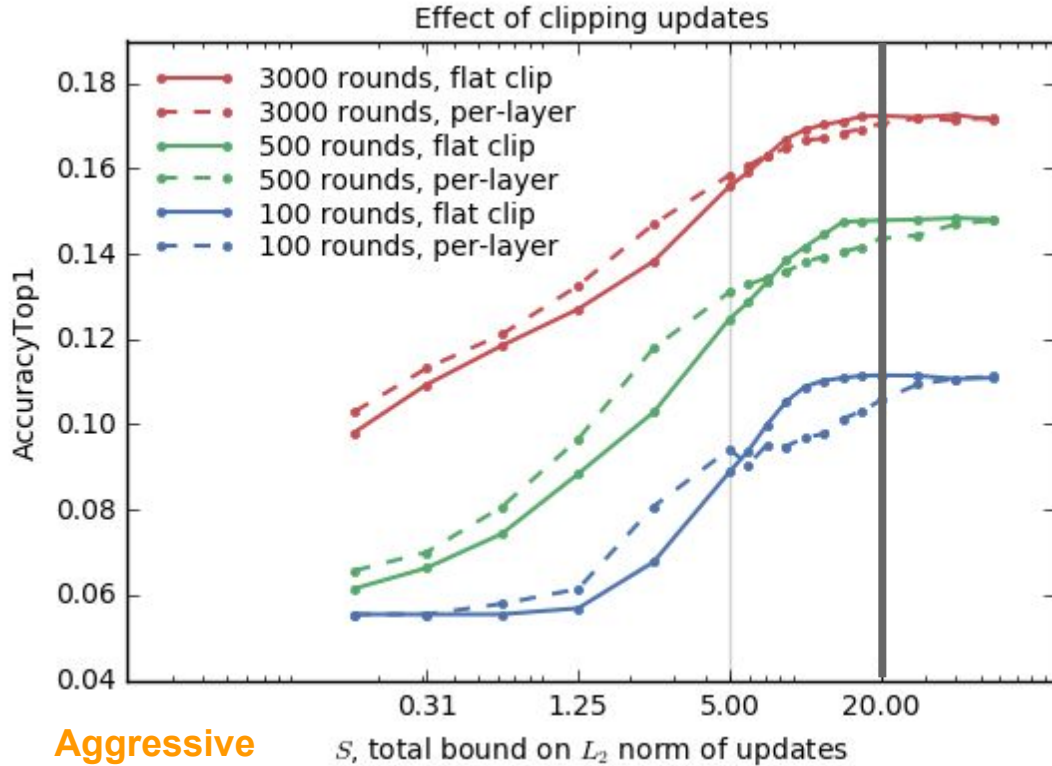## 23x decrease in database queries

# The effect of clipping updates



Effect of clipping updates

**No Clipping**

**Aggressive Clipping**

**Sampling E[C] = 100 users per round.**

# The effect of clipping updates



**No Clipping**

**Aggressive Clipping**
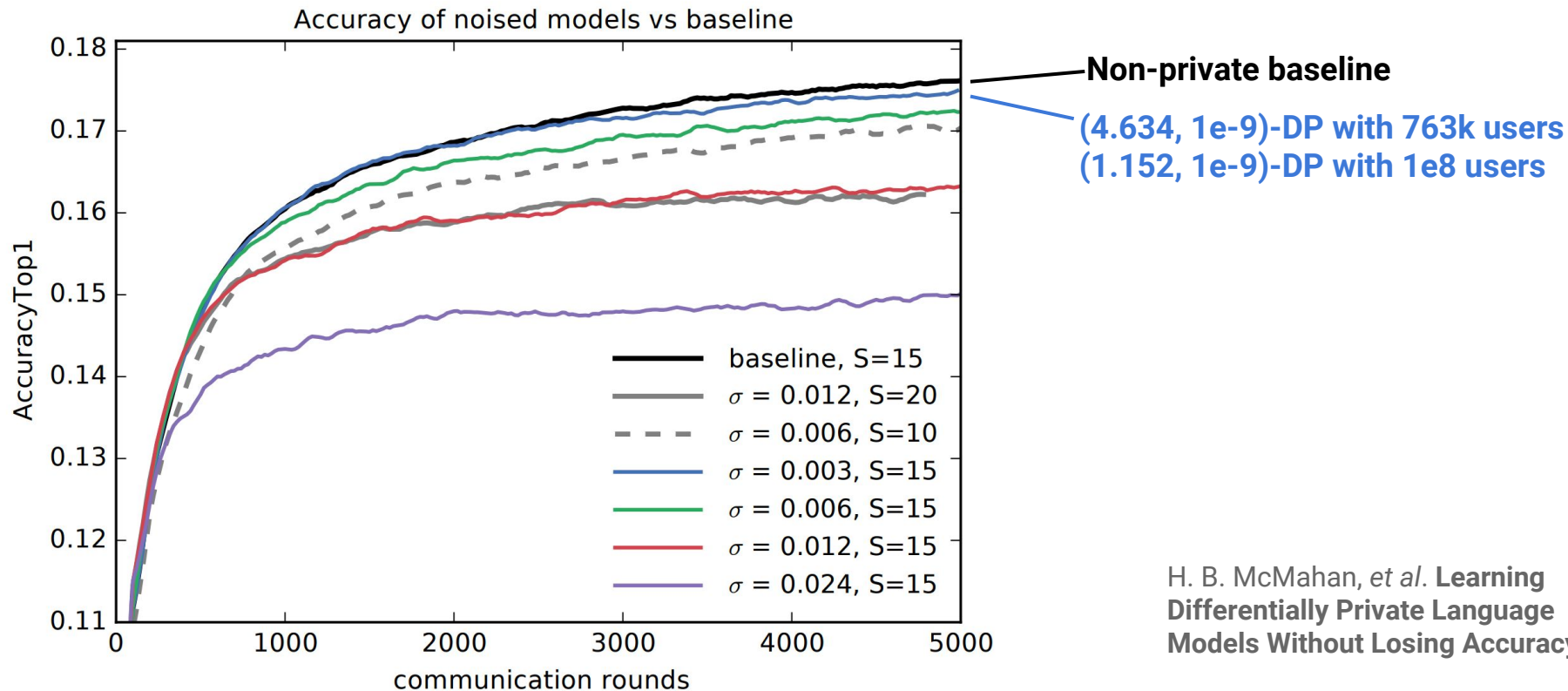
**Sampling E[C] = 100 users per round.**

# The effect of noising updates



Effect of noising updates

Clipping at S=20

Sampling E[C] = 100 users per round.

# Differential Privacy for Language Models



Accuracy of noised models vs baseline

**Non-private baseline**

**(4.634, 1e-9)-DP with 763k users**
**(1.152, 1e-9)-DP with 1e8 users**

Legend:
- baseline, S=15
- $\sigma = 0.012$, S=20
- $\sigma = 0.006$, S=10
- $\sigma = 0.003$, S=15
- $\sigma = 0.006$, S=15
- $\sigma = 0.012$, S=15
- $\sigma = 0.024$, S=15

H. B. McMahan, *et al*. **Learning Differentially Private Language Models Without Losing Accuracy**.

# Differential Privacy for Language Models



Accuracy of noised models vs baseline

Legend:
- baseline, S=15
- $\sigma = 0.012$, S=20
- $\sigma = 0.006$, S=10
- $\sigma = 0.003$, S=15
- $\sigma = 0.006$, S=15
- $\sigma = 0.012$, S=15
- $\sigma = 0.024$, S=15

**Baseline Training**
users per round = 100
tokens per round = 160k
17.5% accuracy in 4120 rounds

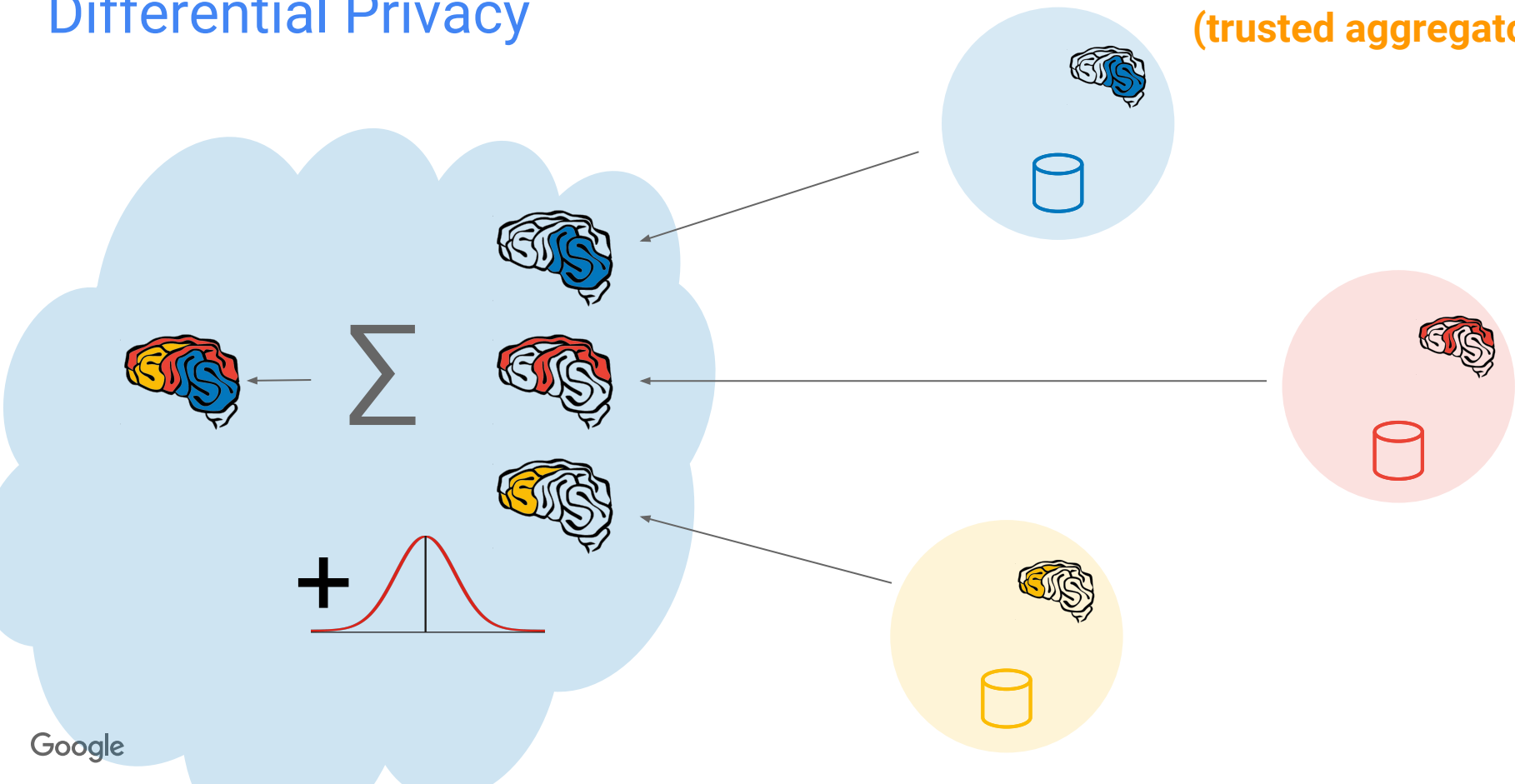**(1.152, 1e-9) DP Training**
$\mathbb{E}$[users per round] = 5k
$\mathbb{E}$[tokens per round] = 8000k
17.5% **estimated** accuracy
in 5000 rounds

# Differential Privacy for Language Models



Accuracy of noised models vs baseline

Legend:
- baseline, S=15
- $\sigma = 0.012$, S=20
- $\sigma = 0.006$, S=10
- $\sigma = 0.003$, S=15
- $\sigma = 0.006$, S=15
- $\sigma = 0.012$, S=15
- $\sigma = 0.024$, S=15

**Baseline Training**
users per round = 100
tokens per round = 160k
17.5% accuracy in 4120 rounds

**(1.152, 1e-9) DP Training**
$\mathbb{E}$[users per round] = 5k
$\mathbb{E}$[tokens per round] = 8000k
17.5% **estimated** accuracy
     in 5000 rounds

Private training achieves
**equal accuracy**, but using
**60x more computation**.

Differential Privacy

Differential Privacy
(trusted aggregator)

Google

# Differential Privacy

Google

Differential Privacy

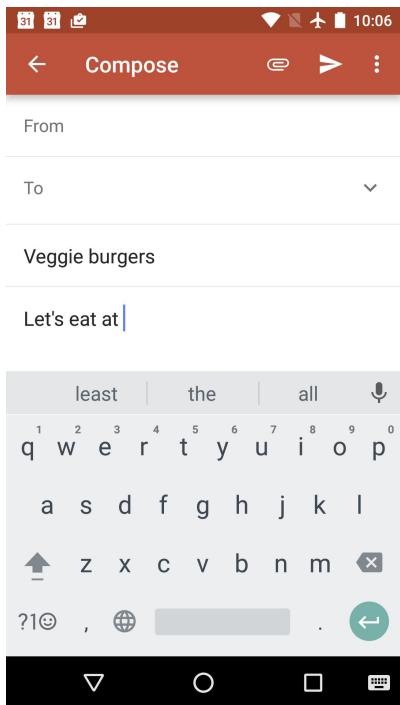Differential Privacy with Secure Aggregation

Google

# Differential Privacy is complementary to Federated Learning

- FL algorithms touch data one user (one device) at time --- natural algorithms for user-level privacy

- Communication constraints mean we want to touch the data as few times as possible --- also good for privacy.

- The DP guarantee is complementary to FL's focussed collection & ephemeral updates.



**Federated Learning**

# Federated Learning in Gboard





Google

# Open Questions and Challenges

**Showing privacy is possible**

    Many open research questions:

    - Further lower computational and/or utility cost of differential privacy

    - More communication-efficient algorithms for FL

**Making privacy easy**

    Possible is not enough. Research to enable "privacy by default" in machine learning.

    - Can federated learning be as easy as centralized learning?

    - Differential privacy for deep learning without parameter tuning?

    - How do we handle privacy budgets across time and across domains?

# Questions