

# Federated Machine Learning via *Over-the-Air Computation*

Yuanming Shi

ShanghaiTech University



# Outline

- **Motivations**

- Big data, IoT, AI

- **Three vignettes:**

- **Federated machine learning**

- ❖ Federated model aggregation

- **Over-the-air computation**

- ❖ Joint device selection and beamforming design

- **Sparse and low-rank optimization**

- ❖ Difference-of-convex programming algorithm

# Intelligent IoT ecosystem

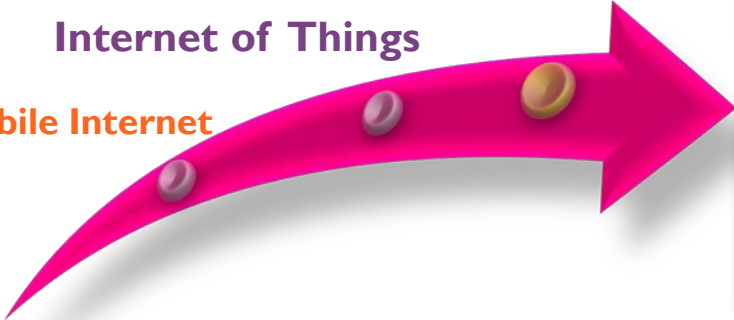


(Internet of Skills)

**Tactile Internet**

Internet of Things

Mobile Internet



**Develop computation, communication & AI technologies:**  
enable smart IoT applications to make **low-latency decision** on streaming data



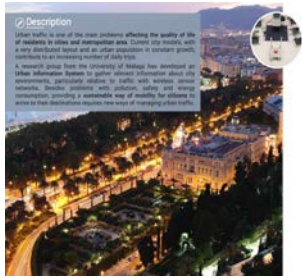
# Intelligent IoT applications



Autonomous vehicles



Smart home



Smart city



Smart health



Smart agriculture



Smart drones

# Challenges

- Retrieve or infer information from high-dimensional/large-scale data



limited processing ability  
(computation, storage, ...)

2.5 **exabytes** of data  
are generated every day (2012)

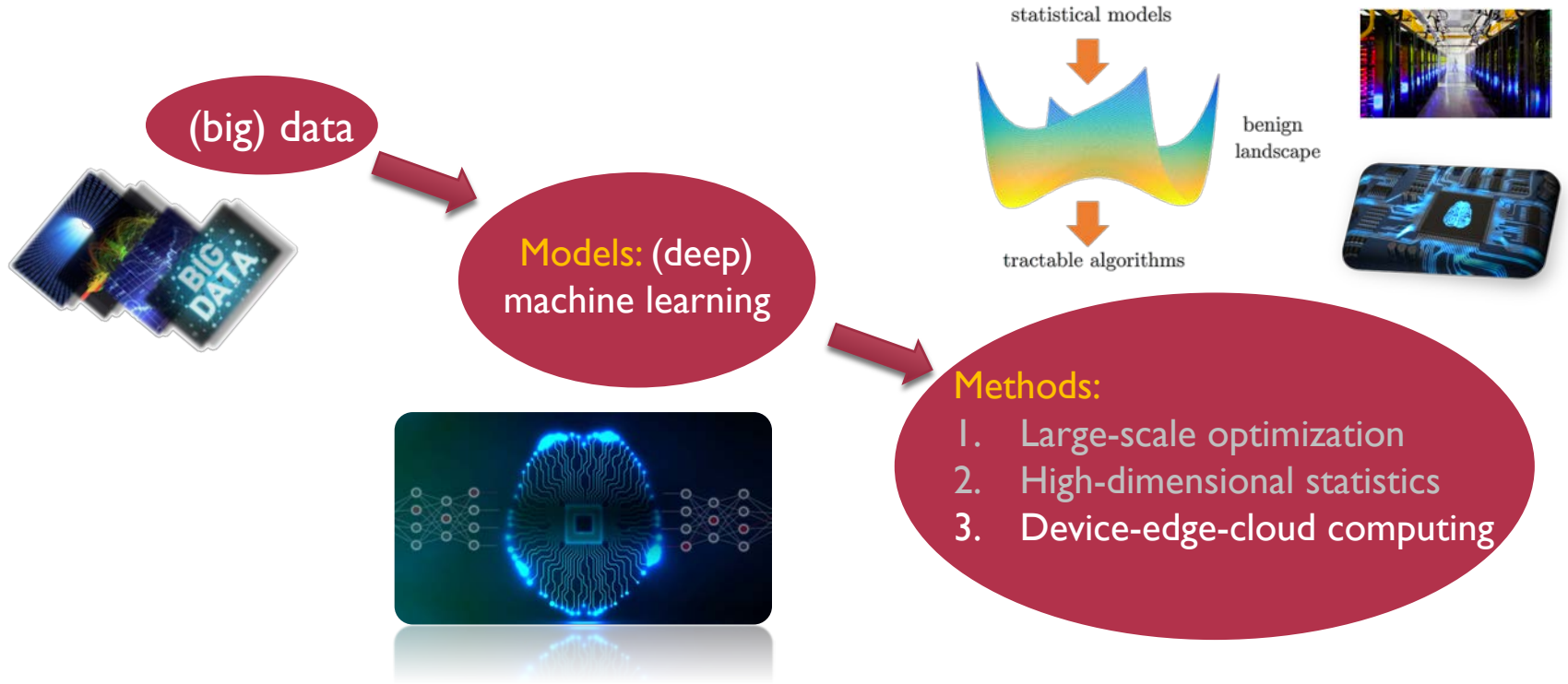
**exabyte** → **zettabyte** → **yottabyte...??**

We're interested in the **information** rather  
than the data

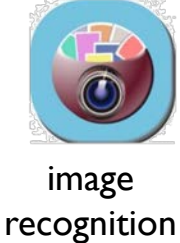
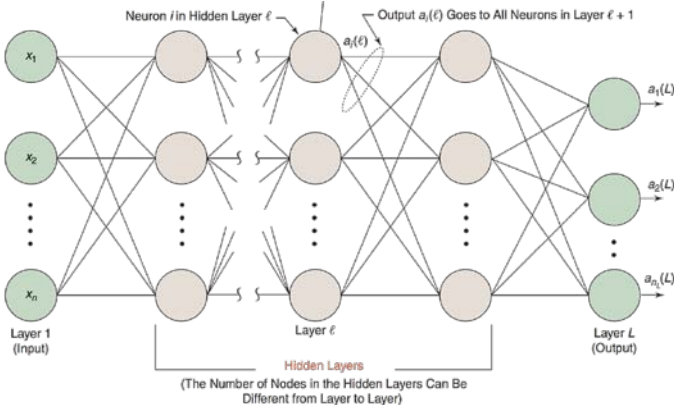
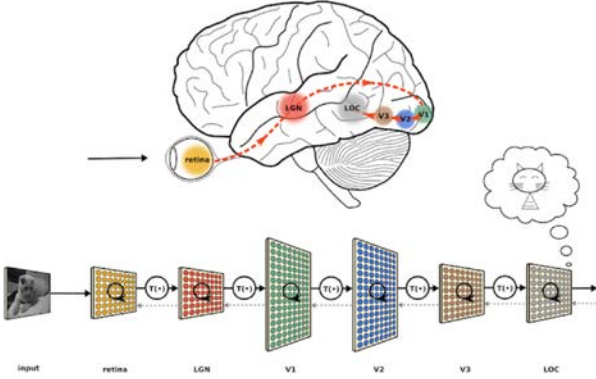
## Challenges:

- ❖ High computational cost
- ❖ Only limited memory is available
- ❖ Do NOT want to compromise statistical accuracy

# High-dimensional data analysis



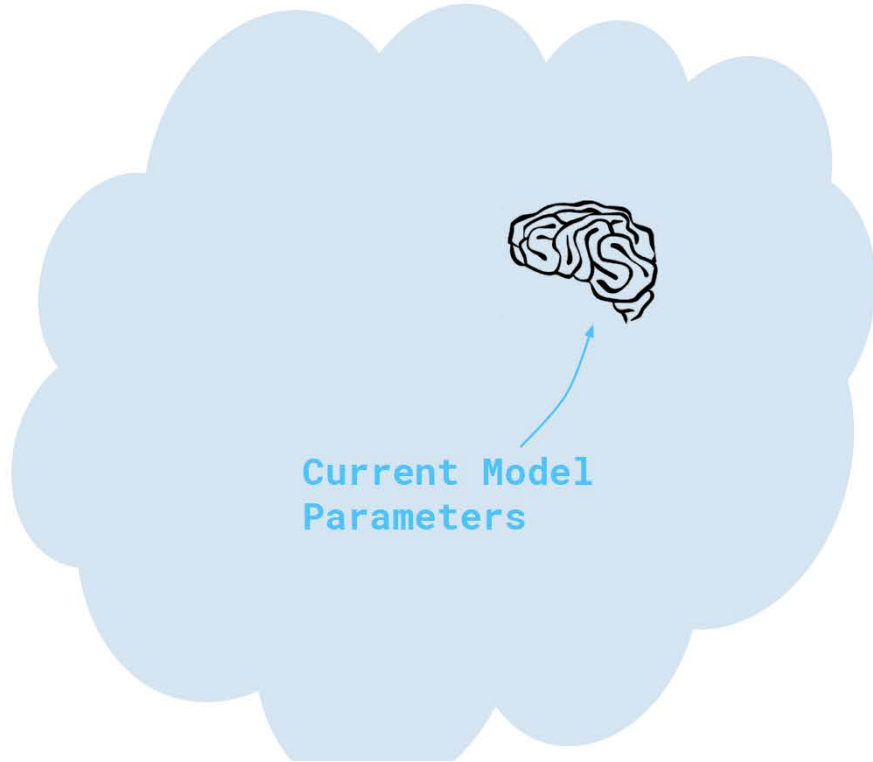
# Deep learning: next wave of AI



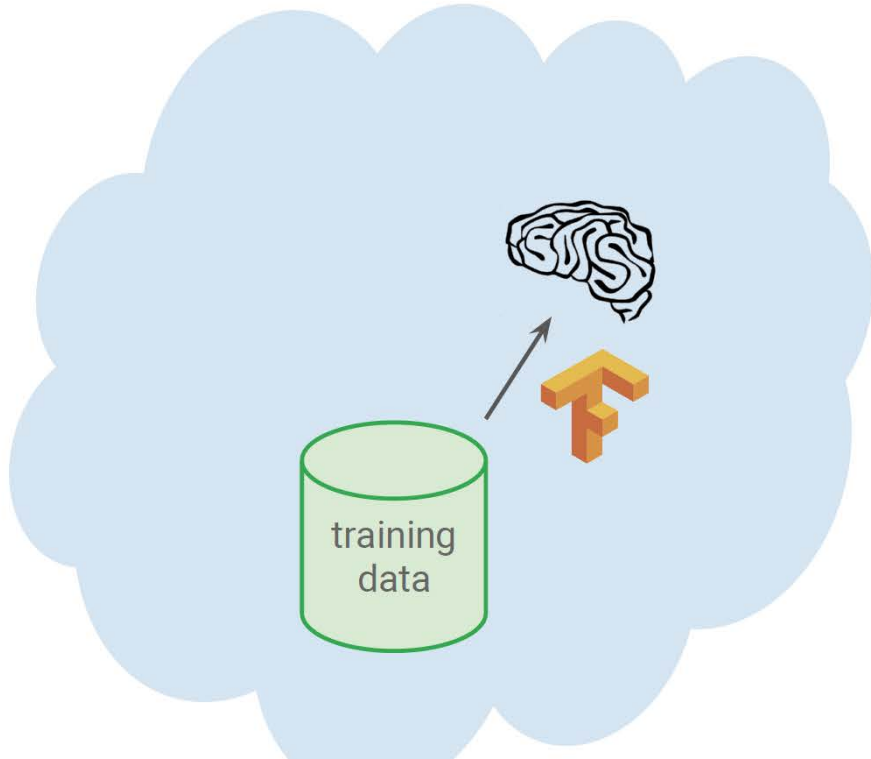
## Cloud-centric machine learning



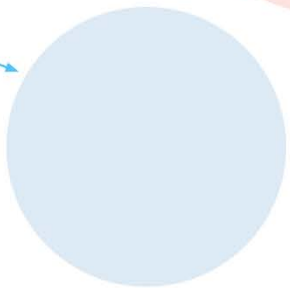
# The model lives in the cloud



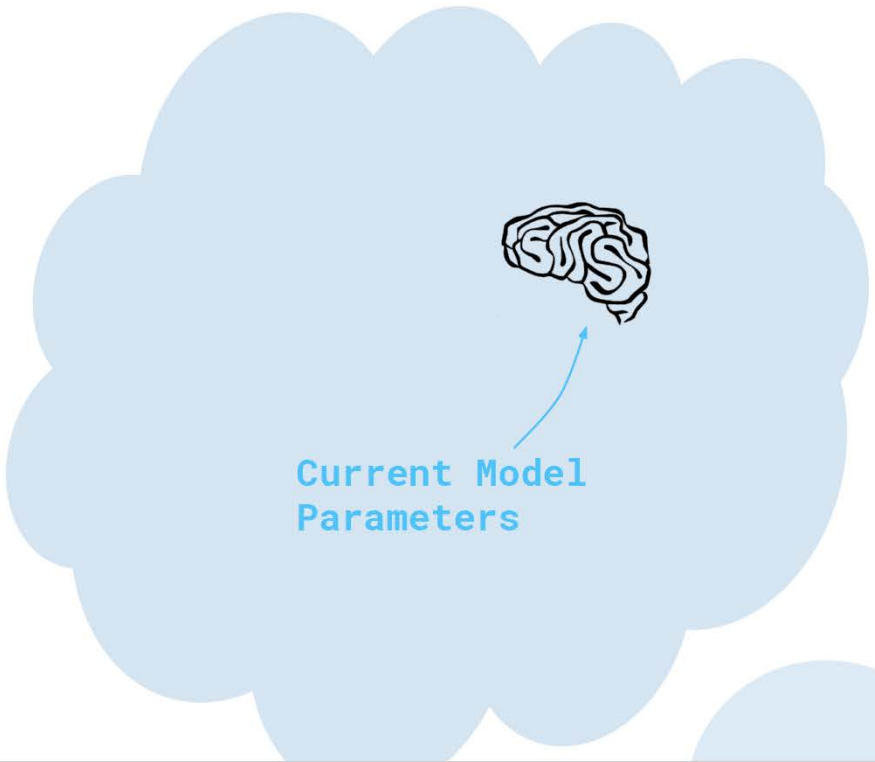
# We train models in the cloud



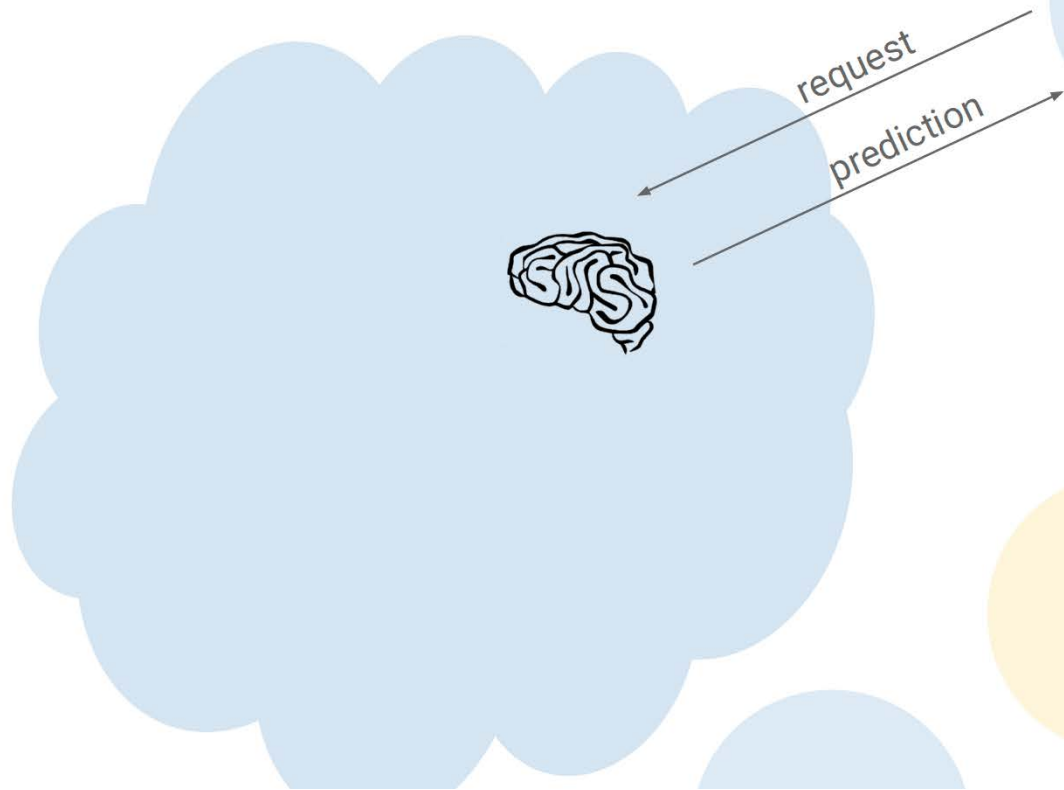
Mobile  
Device



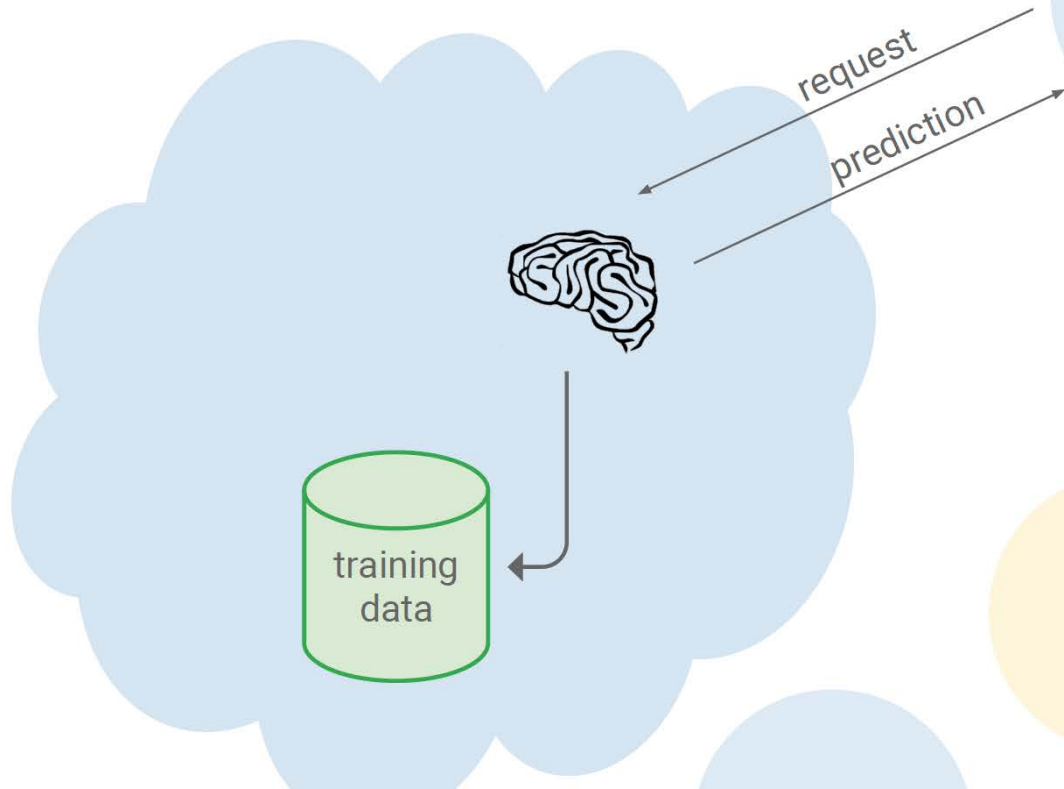
Current Model  
Parameters



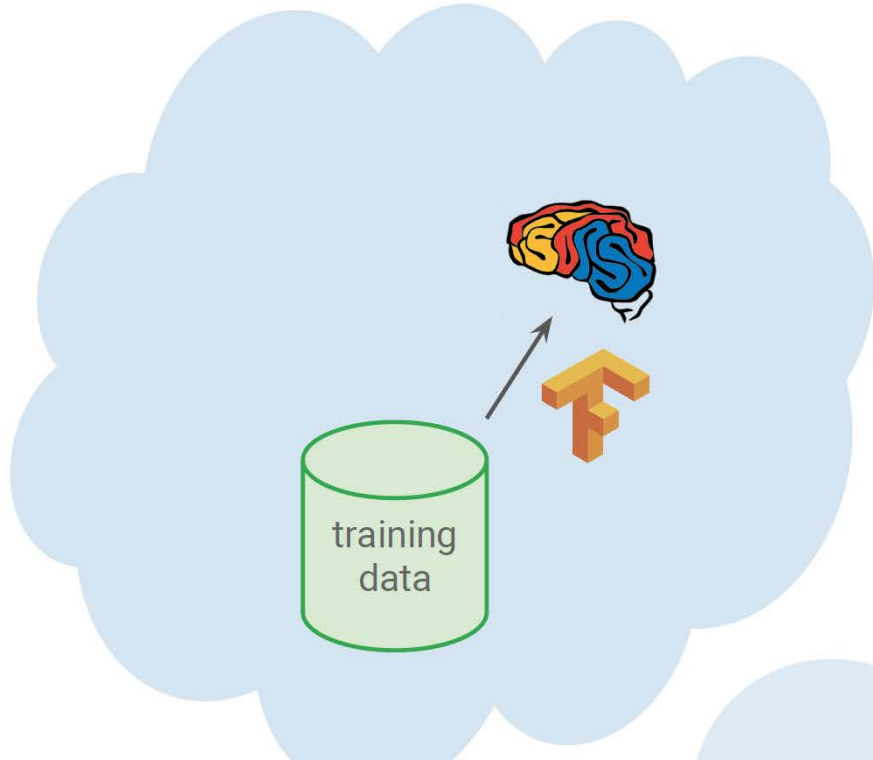
# Make predictions in the cloud



# Gather training data in the cloud



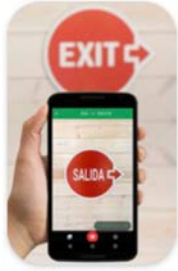
# And make the models better



*Why edge machine learning?*

# Learning on the edge

- The emerging high-stake AI applications: low-latency, privacy,...



phones



drones



robots



glasses



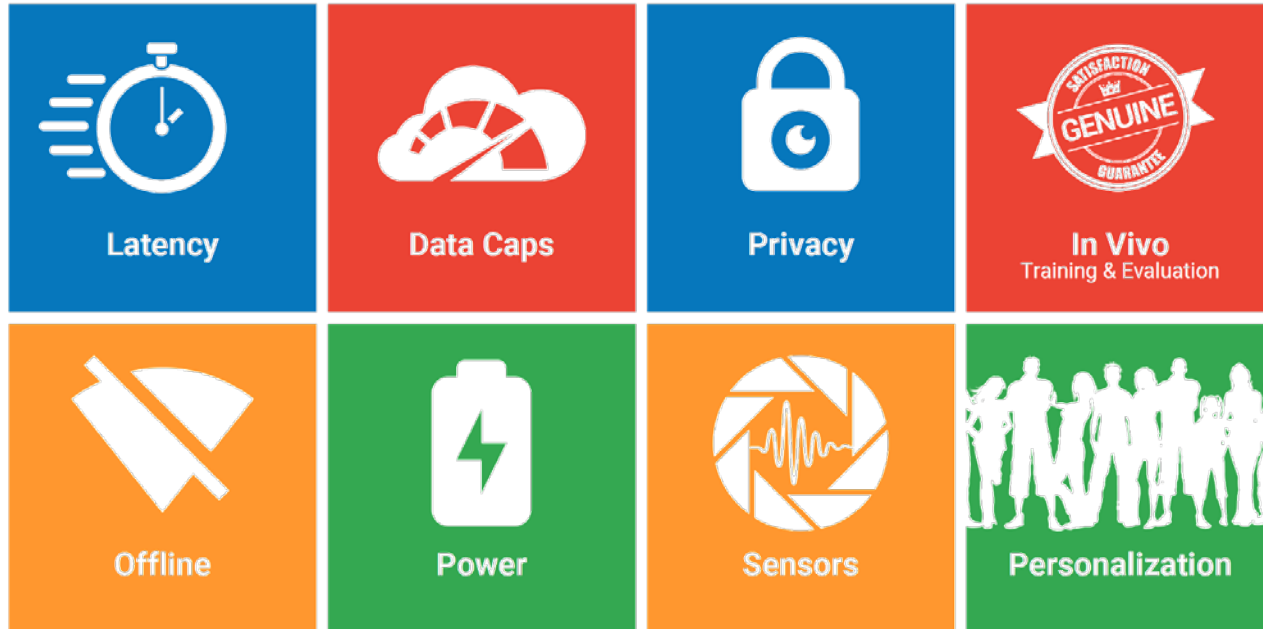
self driving cars

**where to compute?**



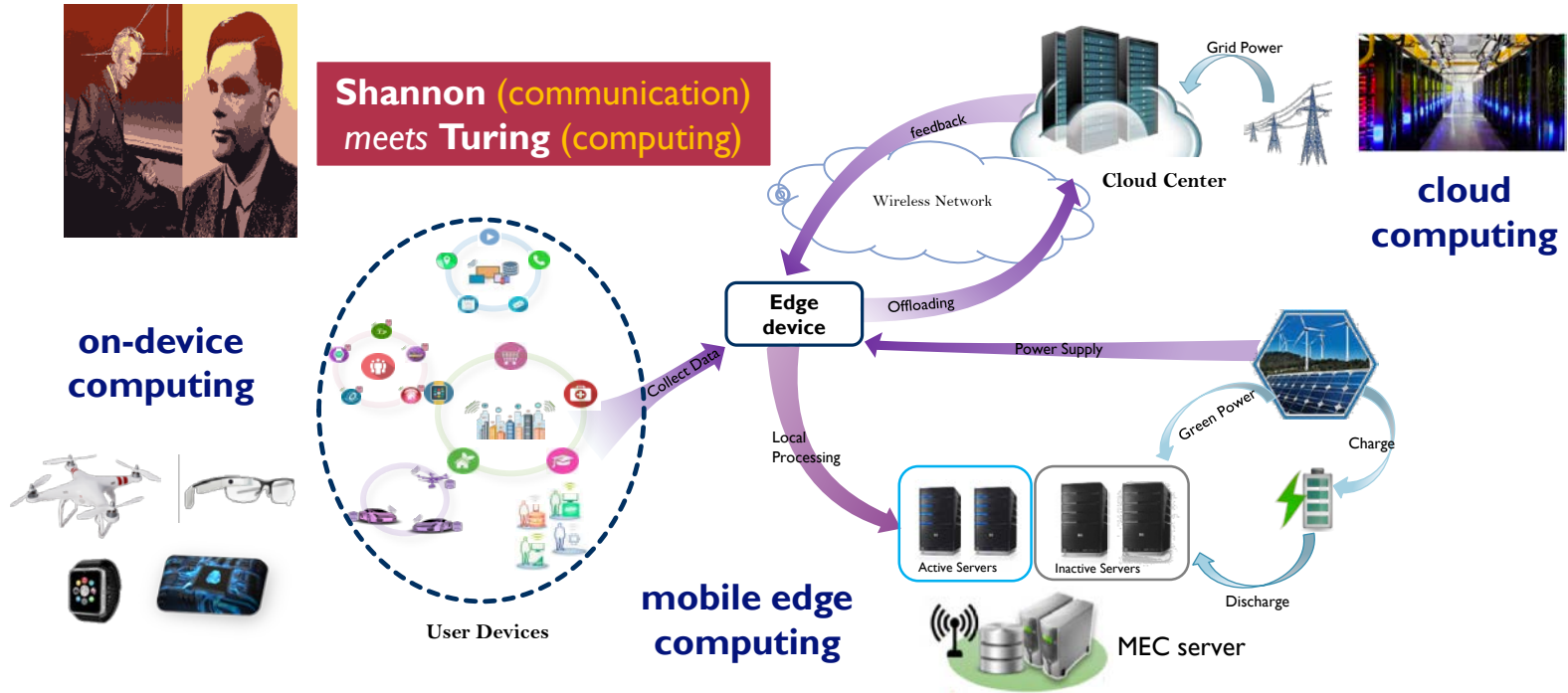
# Mobile edge AI

- Processing at “edge” instead of “cloud”



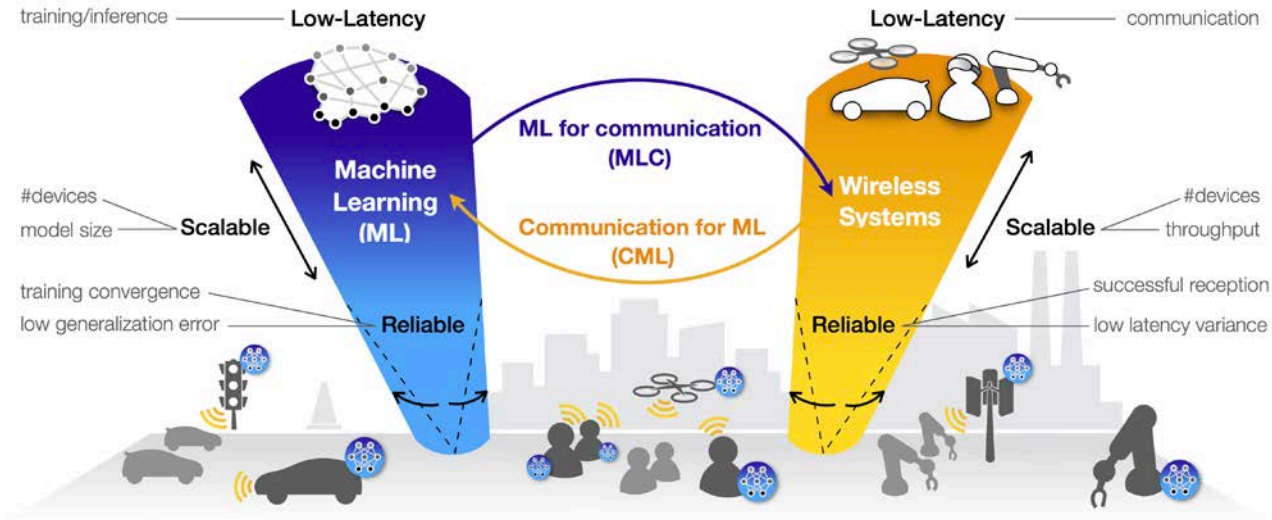
# Edge computing ecosystem

- **“Device-edge-cloud”** computing system for mobile AI applications



# Edge machine learning

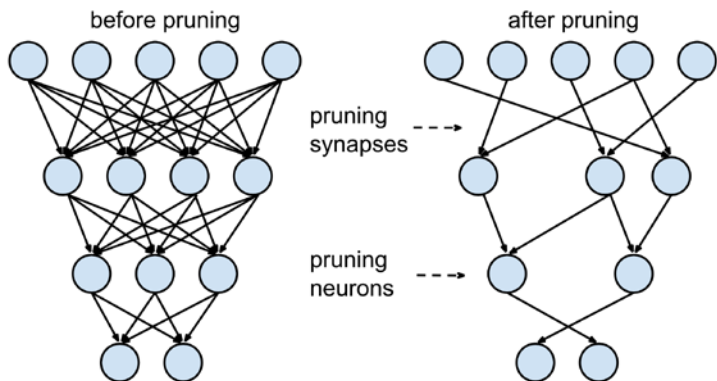
- **Edge ML:** both ML inference and training processes are pushed down into the network edge (bottom)



## *On-device inference*

# Deep model compression

- Layer-wise deep neural network pruning via sparse optimization



sparse optimization

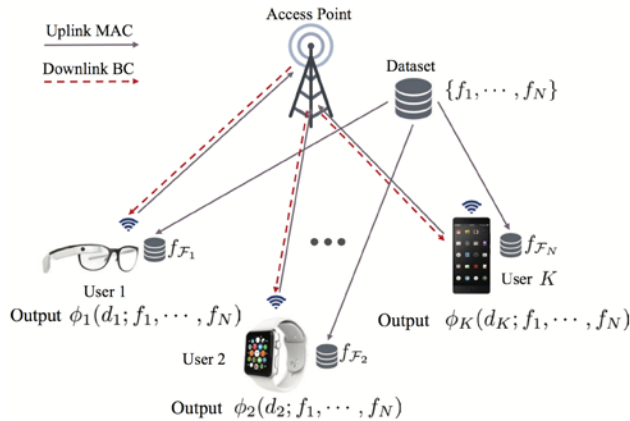
$$\begin{aligned} & \text{minimize}_{\mathbf{W} \in \mathbb{R}^{d_{\ell-1} \times d_{\ell}}} \|\mathbf{W}\|_1 \\ & \text{subject to} \quad \|\max(\mathbf{W}^T \mathbf{X}_{\ell-1}, 0) - \mathbf{X}_{\ell}\|_F \leq \epsilon \end{aligned}$$



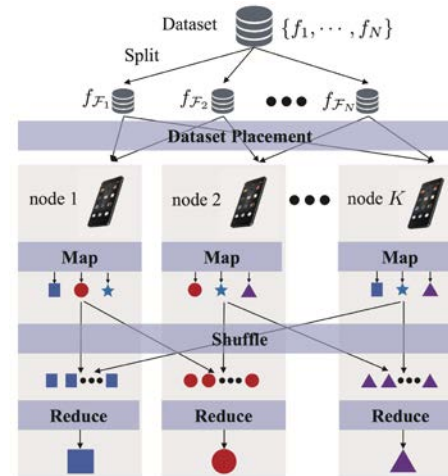
[Ref] T. Jiang, X. Yang, Y. Shi, and H. Wang, “Layer-wise deep neural network pruning via iteratively reweighted optimization,” in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Brighton, UK, May 2019.

# Edge distributed inference

- **Wireless MapReduce** for on-device distributed inference process



wireless distributed computing system

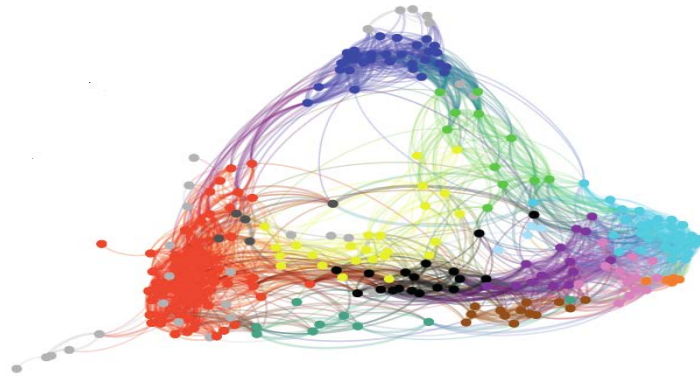


distributed computing model

[Ref] K. Yang, Y. Shi, and Z. Ding, “Data shuffling in wireless distributed computing via low-rank optimization,” *IEEE Trans. Signal Process.*, vol. 67, no. 12, pp. 3087-3099, Jun., 2019.

*This talk: On-device training*

## Vignettes A: *Federated machine learning*

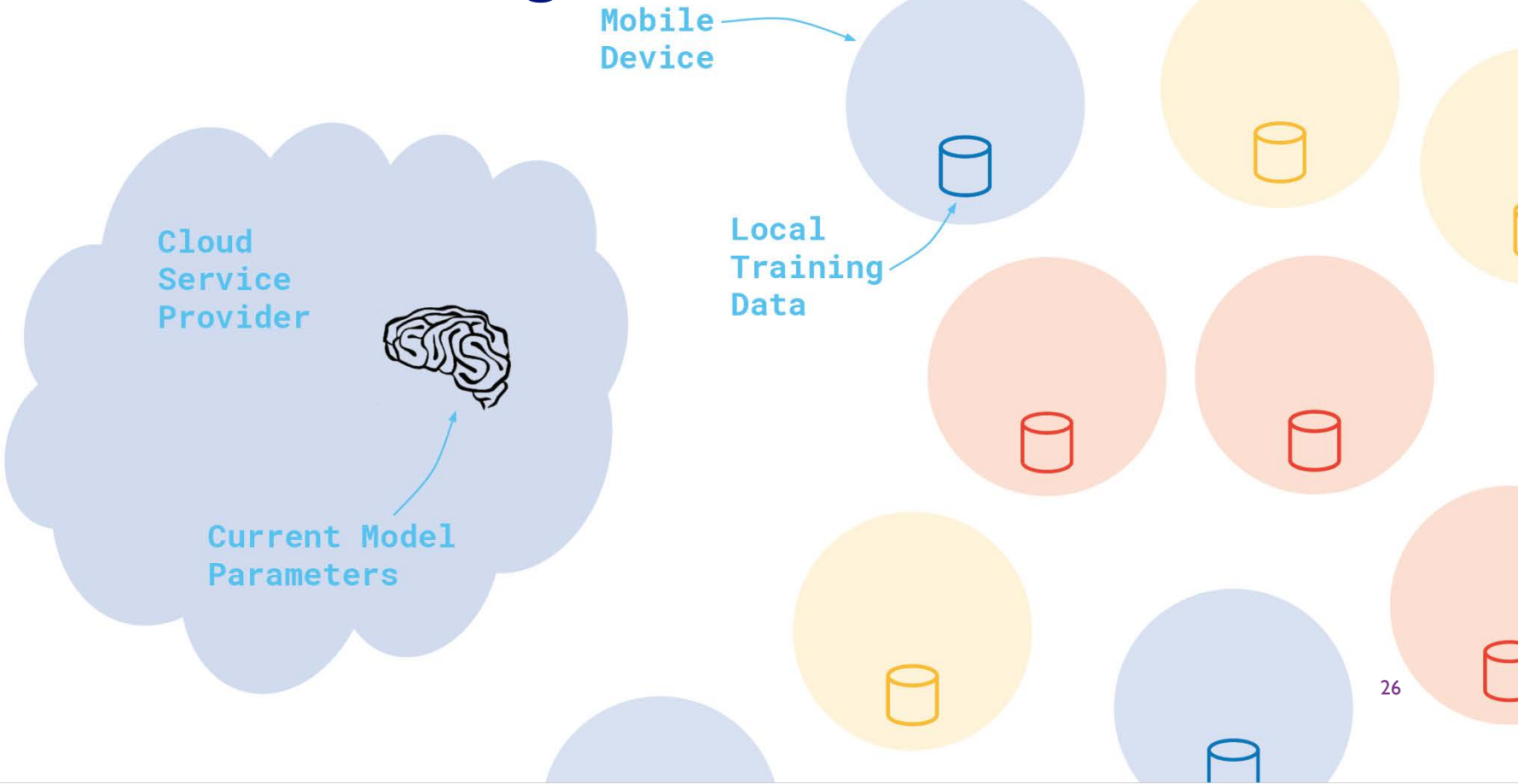




# Federated computation and learning

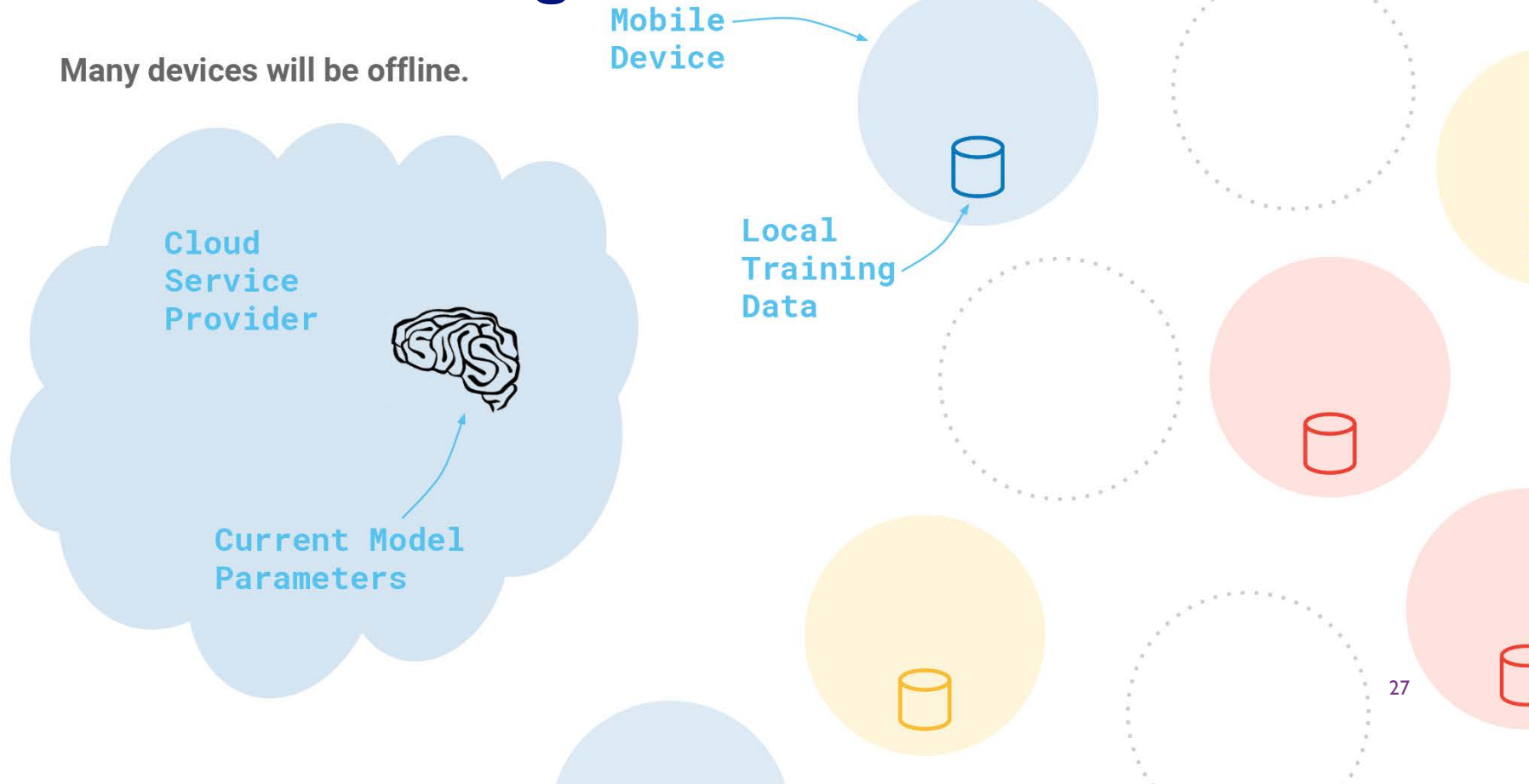
- **Goal:** imbue mobile devices with state of the art machine learning systems *without centralizing data* and with *privacy* by default
- **Federated computation:** a server coordinates a fleet of participating devices to compute aggregations of devices' private data
- **Federated learning:** a shared global model is trained via federated computation

# Federated learning



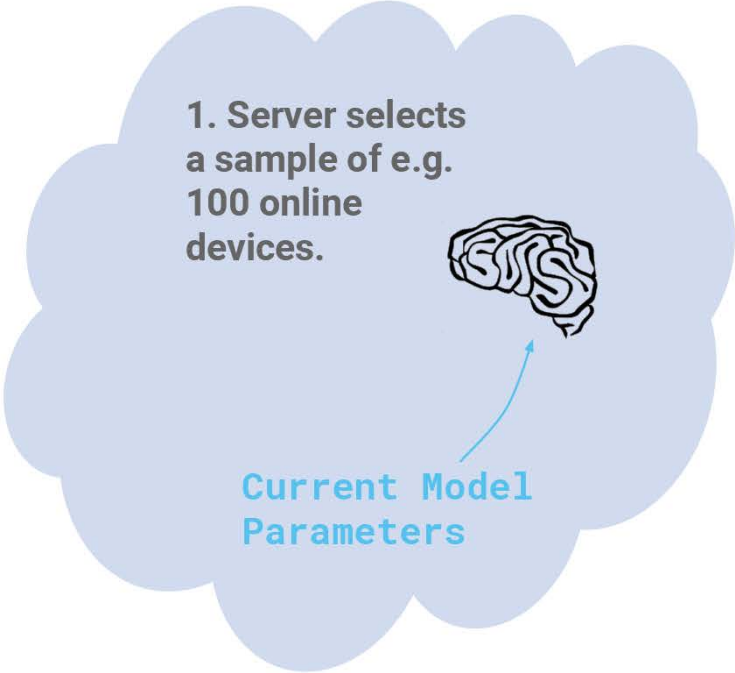
# Federated learning

Many devices will be offline.

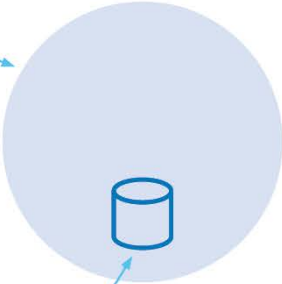


# Federated learning

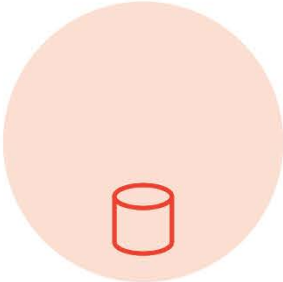
Many devices will be offline.



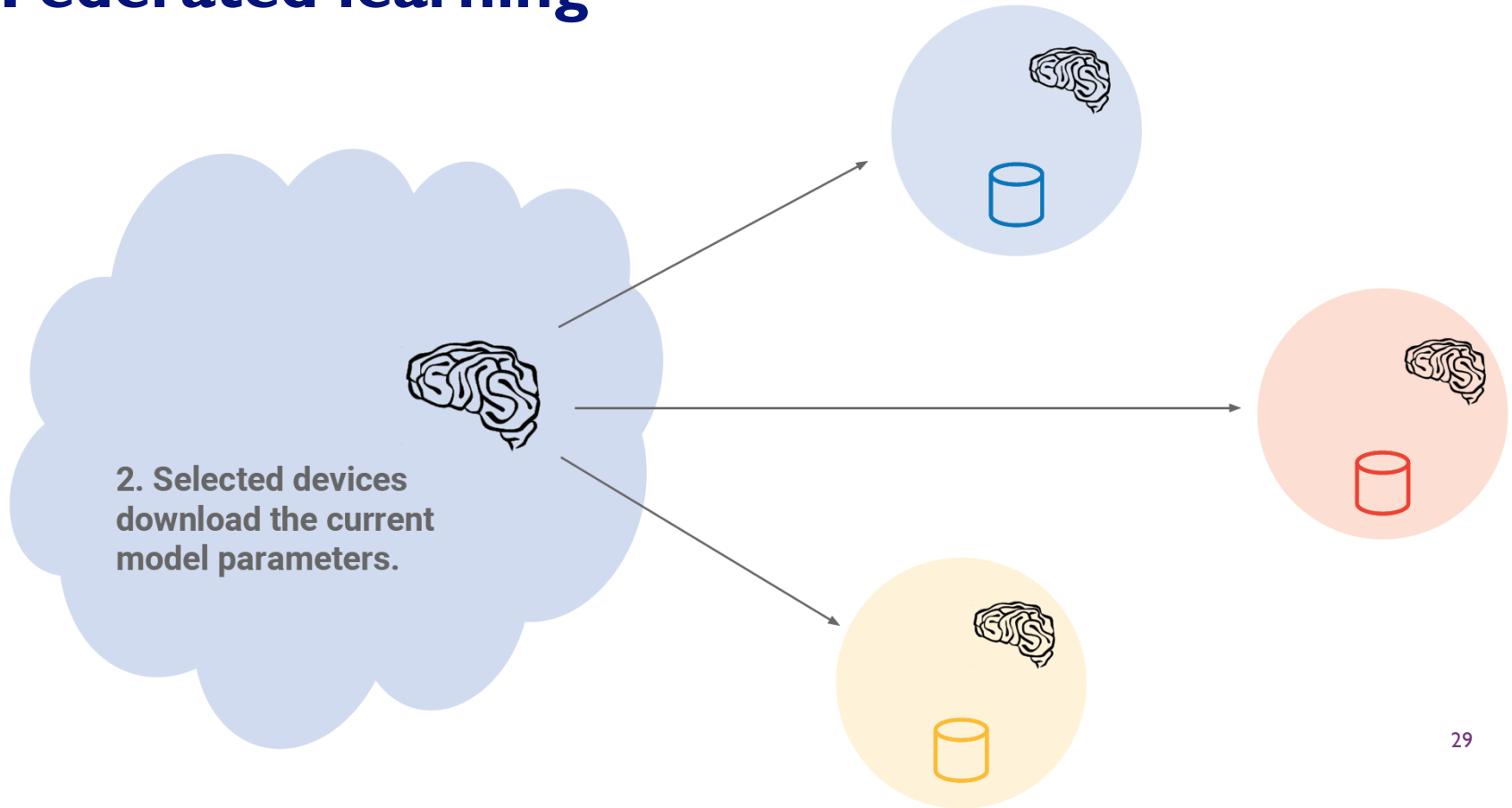
Mobile Device



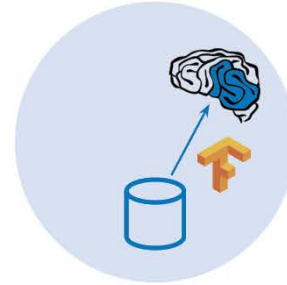
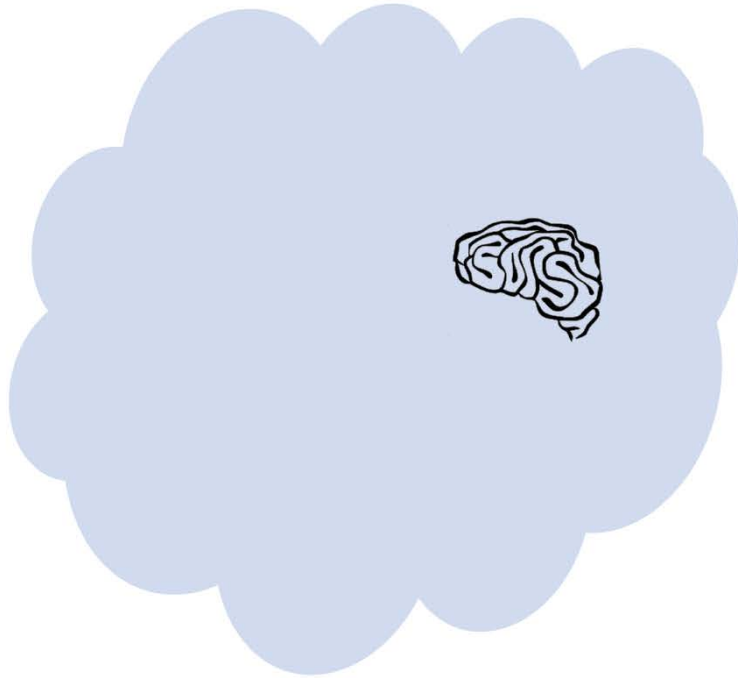
Local Training Data



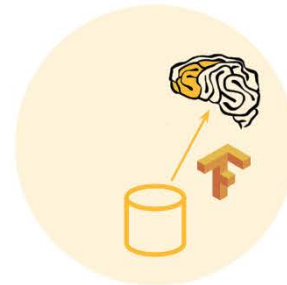
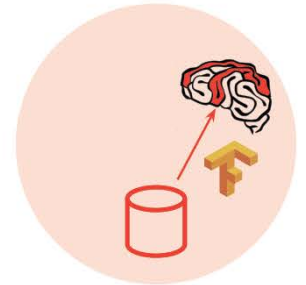
# Federated learning



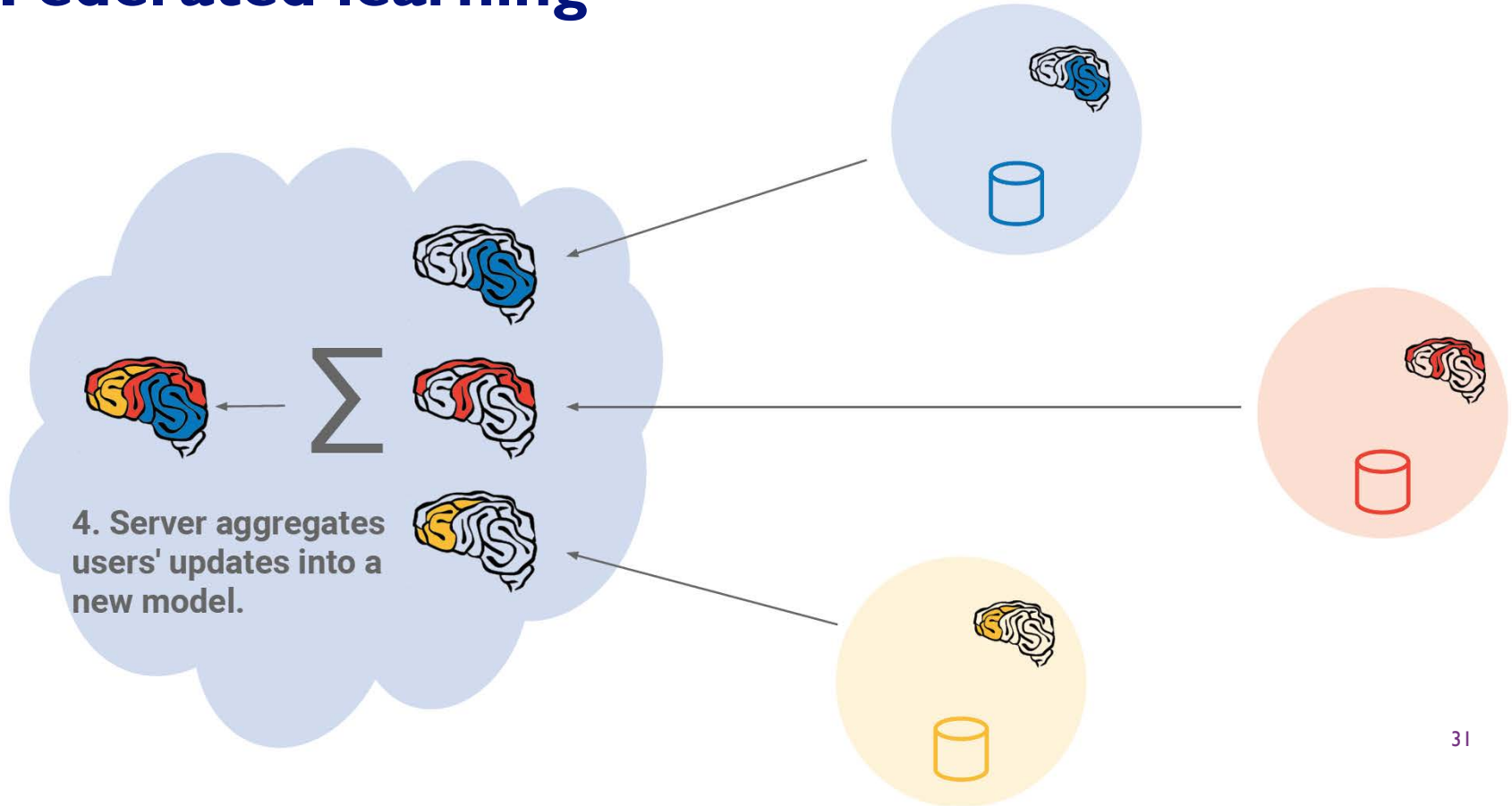
# Federated learning



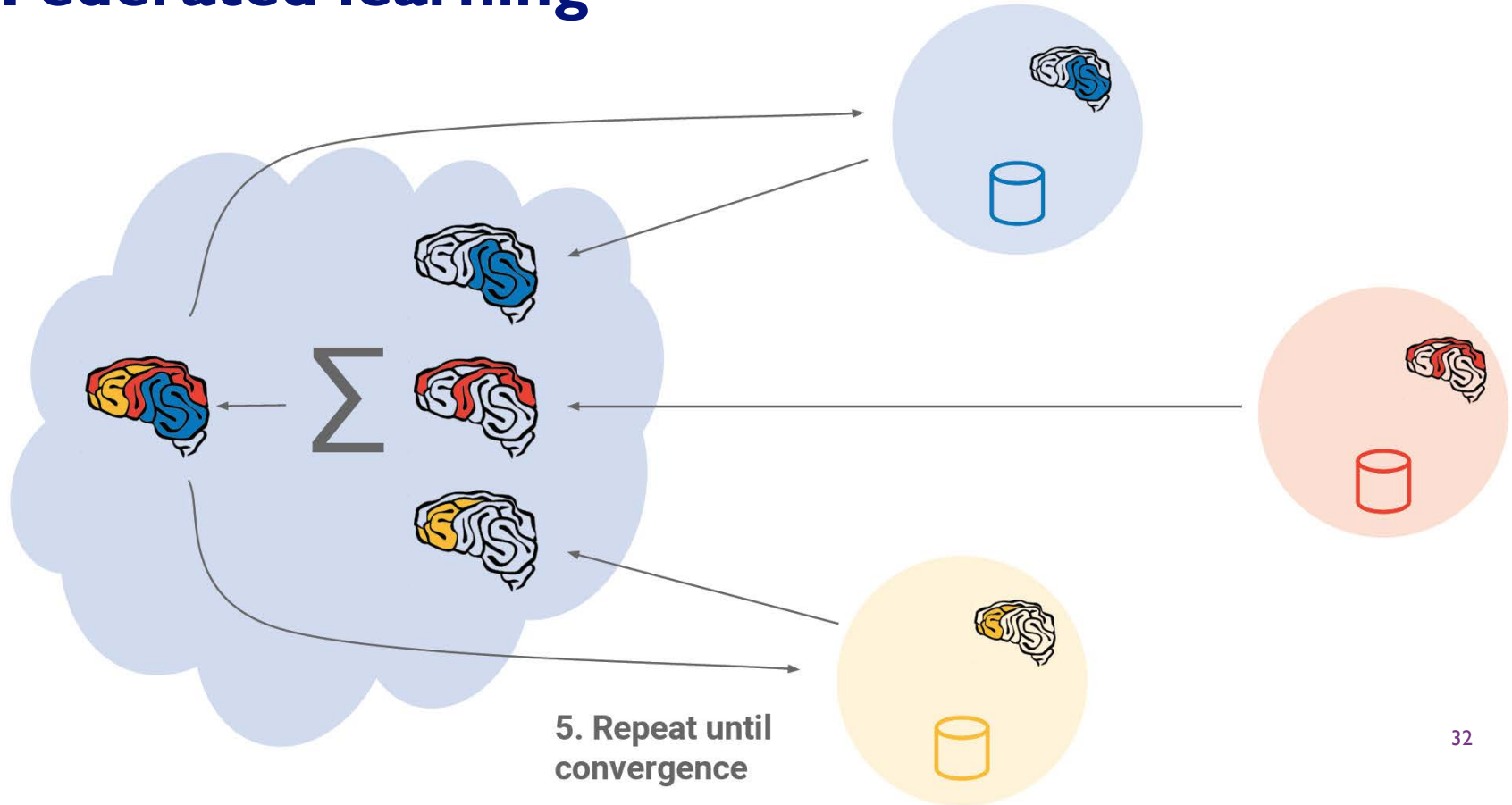
**3. Devices compute an update using local training data**



# Federated learning



# Federated learning



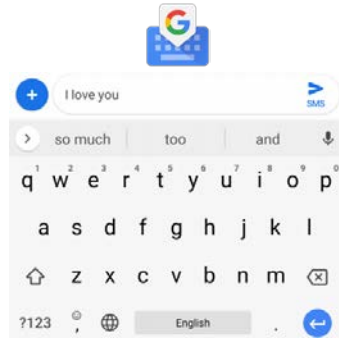


# Federated learning: applications

- **Applications:** where the data is generated at the mobile devices and is undesirable/infeasible to be transmitted to centralized servers



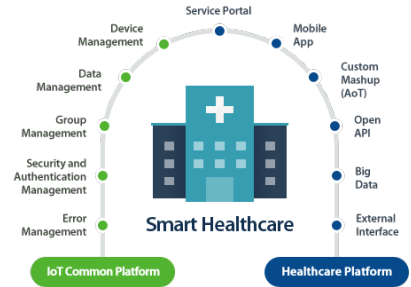
financial services



keyboard prediction



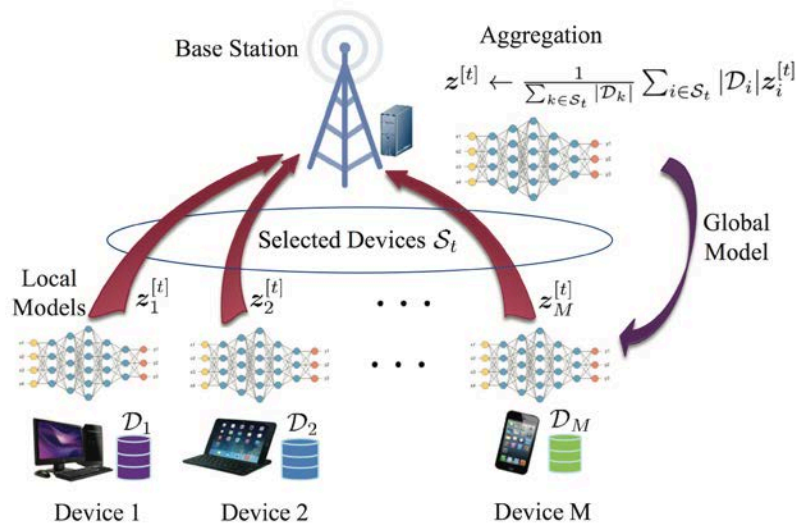
smart retail



smart healthcare

# Federated learning over wireless networks

- **Goal:** train a shared global model via *wireless* federated computation



on-device distributed federated learning system

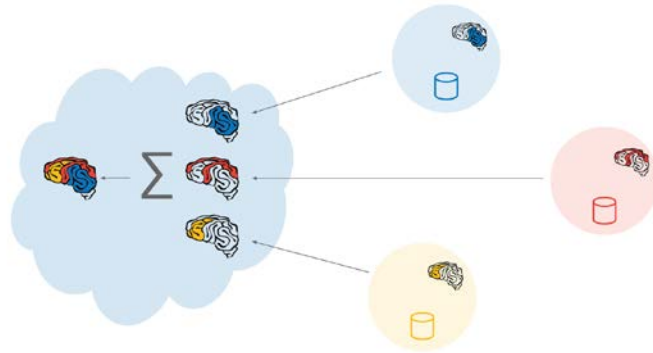
## System challenges

- Massively distributed
- Node heterogeneity

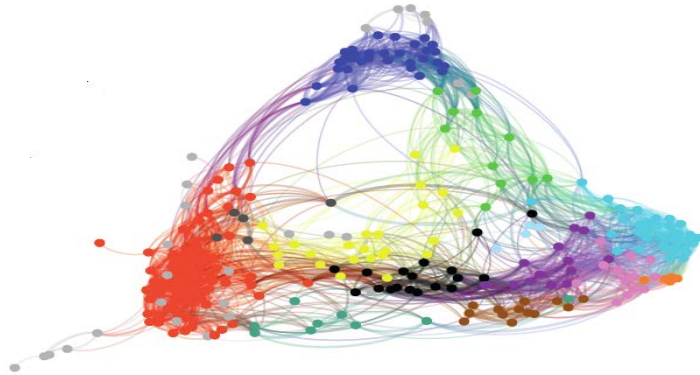
## Statistical challenges

- ❖ Unbalanced
- ❖ Non-IID
- ❖ Underlying structure

*How to efficiently aggregate models over wireless networks?*



## Vignettes B: *Over-the-air computation*

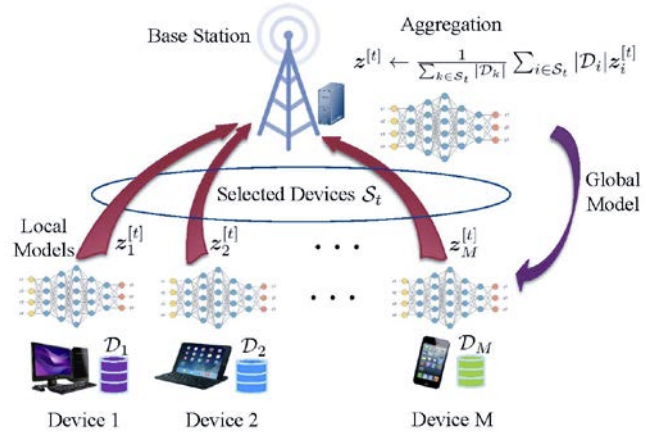


# Model aggregation via over-the-air computation

- Aggregating local updates from mobile devices

$$z \leftarrow \frac{1}{\sum_{k \in \mathcal{S}} |\mathcal{D}_k|} \sum_{k \in \mathcal{S}} |\mathcal{D}_k| z_k$$

- weighted sum of messages
- $M$  mobile devices and one  $N$ -antenna base station
- $\mathcal{S} \subseteq \{1, \dots, M\}$  is the set of selected devices
- $|\mathcal{D}_k|$  is the data size at device  $k$



**Over-the-air computation:**  
explore signal superposition of  
a wireless multiple-access  
channel for model aggregation

# Over-the-air computation

- The estimated value before post-processing at the BS

$$\hat{g} = \frac{1}{\sqrt{\eta}} \mathbf{m}^H \mathbf{y} = \frac{1}{\sqrt{\eta}} \mathbf{m}^H \sum_{i \in \mathcal{S}} \mathbf{h}_i b_i z_i + \frac{\mathbf{m}^H \mathbf{n}}{\sqrt{\eta}}$$

- $b_i$  is the transmitter scalar,  $\mathbf{m}$  is the received beamforming vector,  $\eta$  is a normalizing factor
- target function to be estimated:  $g = \sum_{i \in \mathcal{S}} |\mathcal{D}_i| z_i$
- recovered aggregation vector entry via post-processing:  $\hat{z} = \frac{1}{\sum_{i \in \mathcal{S}} |\mathcal{D}_i|} \hat{g}$

- **Model aggregation error:**

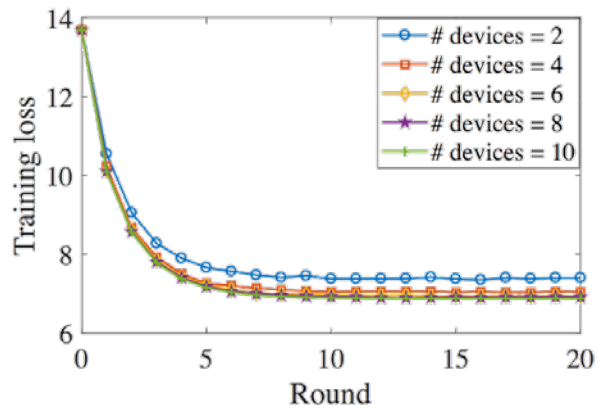
$$\text{MSE}(\hat{g}, g; \mathcal{S}, \mathbf{m}) = \frac{\|\mathbf{m}\|^2 \sigma^2}{\eta} = \frac{\sigma^2}{P_0} \max_{i \in \mathcal{S}} |\mathcal{D}_i|^2 \frac{\|\mathbf{m}\|^2}{\|\mathbf{m}^H \mathbf{h}_i\|^2}$$

- Optimal transmitter scalar:  $b_i = \sqrt{\eta} |\mathcal{D}_i| \frac{(\mathbf{m}^H \mathbf{h}_i)^H}{\|\mathbf{m}^H \mathbf{h}_i\|^2}$

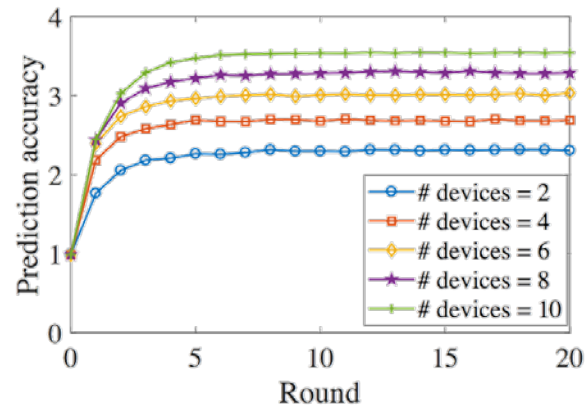
# Problem formulation

- **Key observations:**

- More selected devices yield fast convergence rate of the training process
- Aggregation error leads to the deterioration of model prediction accuracy



(a) Training loss



(b) Relative prediction accuracy

# Problem formulation

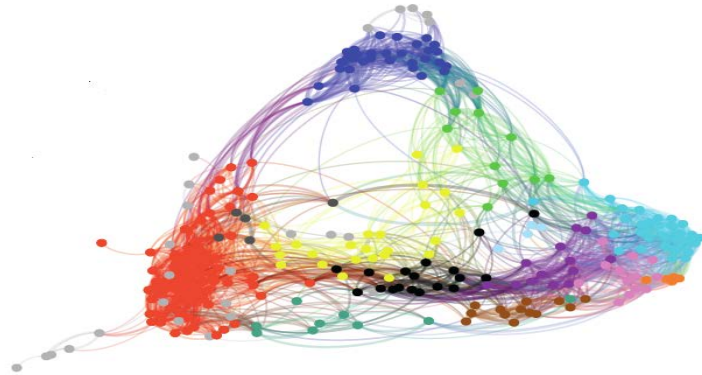
- **Goal:** maximize the number of selected devices under target MSE constraint

$$\underset{\mathcal{S}, \mathbf{m} \in \mathbb{C}^N}{\text{maximize}} \quad |\mathcal{S}| \quad \text{subject to} \quad \left( \max_{i \in \mathcal{S}} |\mathcal{D}_i|^2 \frac{\|\mathbf{m}\|^2}{\|\mathbf{m}^H \mathbf{h}_i\|^2} \right) \leq \gamma$$

- Joint device selection and received beamforming vector design
- Improve convergence rate in the **training process**, guarantee prediction accuracy in the **inference process**
- **Mixed combinatorial optimization problem**



## Vignettes C: *Sparse and low-rank optimization*



# Sparse and low-rank optimization

- Sparse and low-rank optimization for on-device federated learning

$$\begin{aligned} & \underset{\mathcal{S}, \mathbf{m} \in \mathbb{C}^N}{\text{maximize}} && |\mathcal{S}| \\ & \text{subject to} && \left( \max_{i \in \mathcal{S}} |\mathcal{D}_i|^2 \frac{\|\mathbf{m}\|^2}{\|\mathbf{m}^H \mathbf{h}_i\|^2} \right) \leq \gamma \end{aligned}$$

multicasting  
duality



$$\begin{aligned} & \underset{\mathcal{S}, \mathbf{m} \in \mathbb{C}^N}{\text{maximize}} && |\mathcal{S}| \\ & \text{subject to} && \|\mathbf{m}\|^2 - \gamma_i \|\mathbf{m}^H \mathbf{h}_i\|^2 \leq 0, i \in \mathcal{S} \\ & && \|\mathbf{m}\|^2 \geq 1 \end{aligned}$$

sum of feasibilities



$$\begin{aligned} \mathcal{P} : & \underset{\mathbf{x} \in \mathbb{R}_+^M, \mathbf{M} \in \mathbb{C}^{N \times N}}{\text{minimize}} && \|\mathbf{x}\|_0 \\ & \text{subject to} && \text{Tr}(\mathbf{M}) - \gamma_i \mathbf{h}_i^H \mathbf{M} \mathbf{h}_i \leq x_i, \\ & && \mathbf{M} \succeq \mathbf{0}, \text{Tr}(\mathbf{M}) \geq 1 \\ & && \text{rank}(\mathbf{M}) = 1 \end{aligned}$$

$\mathbf{M} = \mathbf{m} \mathbf{m}^H$   
matrix lifting



$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}_+^M, \mathbf{m} \in \mathbb{C}^N}{\text{minimize}} && \|\mathbf{x}\|_0 \\ & \text{subject to} && \|\mathbf{m}\|^2 - \gamma_i \|\mathbf{m}^H \mathbf{h}_i\|^2 \leq x_i, \forall i \\ & && \|\mathbf{m}\|^2 \geq 1 \end{aligned}$$

# Problem analysis

- **Goal:** induce sparsity while satisfying fixed-rank constraint

$$\begin{aligned} \mathcal{P} : & \text{ minimize } && \| \mathbf{x} \|_0 \\ & \mathbf{x} \in \mathbb{R}_+^M, \mathbf{M} \in \mathbb{C}^{N \times N} \\ & \text{ subject to } && \text{Tr}(\mathbf{M}) - \gamma_i \mathbf{h}_i^H \mathbf{M} \mathbf{h}_i \leq x_i, \forall i \\ & && \mathbf{M} \succeq \mathbf{0}, \text{Tr}(\mathbf{M}) \geq 1 \\ & && \text{rank}(\mathbf{M}) = 1 \end{aligned}$$

- Limitations of existing methods
  - **Sparse optimization:** iterative reweighted algorithms are parameters sensitive
  - **Low-rank optimization:** semidefinite relaxation (SDR) approach (i.e., drop rank-one constraint) has the poor capability of returning rank-one solution

# Difference-of-convex functions representation

- Ky Fan  $k$ -norm [Fan, PNAS'1951]: the sum of largest- $k$  absolute values

$$\|\mathbf{x}\|_k = \sum_{i=1}^k |x_{\pi(i)}| \quad \text{convex function}$$

- $\pi$  is a permutation of  $\{1, \dots, M\}$ , where  $|x_{\pi(1)}| \geq \dots \geq |x_{\pi(M)}|$

*MAXIMUM PROPERTIES AND INEQUALITIES FOR THE  
EIGENVALUES OF COMPLETELY CONTINUOUS OPERATORS\**

BY KY FAN

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF NOTRE DAME

Communicated by John von Neumann, September 8, 1951

*PNAS'1951*

# Difference-of-convex functions representation

- DC representation for sparsity function

$$\|\mathbf{x}\|_0 = \min\{k : \|\mathbf{x}\|_1 - \|\mathbf{x}\|_k = 0, 0 \leq k \leq M\}$$

- DC representation for rank-one positive semidefinite matrix

$$\text{rank}(\mathbf{M}) = 1 \Leftrightarrow \text{Tr}(\mathbf{M}) - \|\mathbf{M}\|_2 = 0$$

➤ where  $\text{Tr}(\mathbf{M}) = \sum_{i=1}^N \sigma_i(\mathbf{M})$  and  $\|\mathbf{M}\|_2 = \sigma_1(\mathbf{M})$

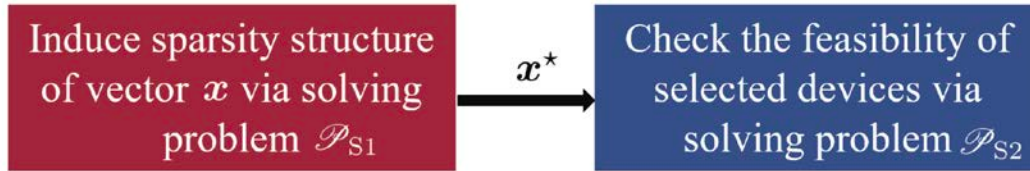
algorithmic  
advantages?



[Ref] J.-y. Gotoh, A. Takeda, and K. Tono, “DC formulations and algorithms for sparse optimization problems,” *Math. Program.*, vol. 169, pp. 141–176, May 2018.

# A DC representation framework

- A two-step framework for device selection



- **Step I:** obtain the sparse solution such that the objective value achieves zero through increasing  $k$  from 0 to  $M$

$$\begin{aligned} \mathcal{P}_{S1} : \underset{x, M}{\text{minimize}} \quad & \underbrace{\|x\|_1 - \|x\|_k}_{\text{red underline}} + \underbrace{\text{Tr}(M) - \|M\|_2}_{\text{blue underline}} \quad \text{zero?} \\ \text{subject to} \quad & \text{Tr}(M) - \gamma_i \mathbf{h}_i^H M \mathbf{h}_i \leq x_i, \forall i = 1, \dots, M \\ & M \succeq \mathbf{0}, \quad \text{Tr}(M) \geq 1, x \succeq \mathbf{0} \end{aligned}$$

# A DC representation framework

- **Step II:** feasibility detection

- Ordering  $\mathbf{x}$  in descending order as  $x_{\pi(1)} \geq \cdots \geq x_{\pi(M)}$
- Increasing  $k$  from 1 to  $M$ , choosing  $\mathcal{S}^{[k]}$  as  $\{\pi(k), \pi(k+1), \dots, \pi(M)\}$

- Feasibility detection via DC programming

$$\begin{aligned} & \text{find } \mathbf{M} \\ & \text{subject to } \text{Tr}(\mathbf{M}) - \gamma_i \mathbf{h}_i^H \mathbf{M} \mathbf{h}_i \leq 0, \forall i \in \mathcal{S}^{[k]} \\ & \mathbf{M} \succeq \mathbf{0}, \text{Tr}(\mathbf{M}) \geq 1, \text{rank}(\mathbf{M}) = 1 \end{aligned}$$



$$\begin{aligned} \mathcal{P}_{S_2} : & \underset{\mathbf{M}}{\text{minimize}} \quad \text{Tr}(\mathbf{M}) - \|\mathbf{M}\|_2 \quad \text{zero?} \\ & \text{subject to } \text{Tr}(\mathbf{M}) - \gamma_i \mathbf{h}_i^H \mathbf{M} \mathbf{h}_i \leq 0, \forall i \in \mathcal{S}^{[k]} \\ & \mathbf{M} \succeq \mathbf{0}, \quad \text{Tr}(\mathbf{M}) \geq 1 \end{aligned}$$

# DC algorithm with convergence guarantees

- $\mathcal{P}_{S1}$  and  $\mathcal{P}_{S2}$ : minimize the difference of two strongly convex functions

$$\underset{\mathbf{X} \in \mathbb{C}^{m \times n}}{\text{minimize}} \quad f(\mathbf{X}) = g(\mathbf{X}) - h(\mathbf{X})$$

➤ e.g.,  $g = \text{Tr}(\mathbf{M}) + I_{C_2}(\mathbf{M}) + \frac{\alpha}{2} \|\mathbf{M}\|_F^2$  and  $h = \|\mathbf{M}\|_2 + \frac{\alpha}{2} \|\mathbf{M}\|_F^2$

- The DC algorithm via **linearizing the concave part**

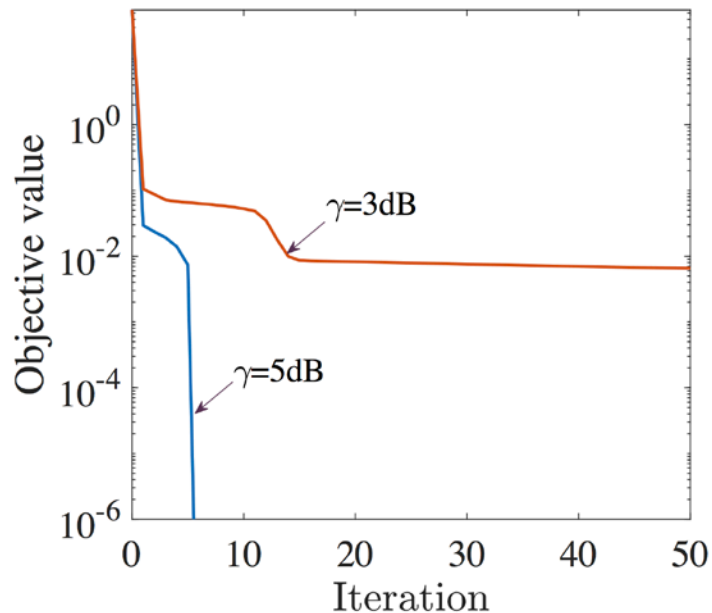
$$\mathbf{X}^{[t+1]} = \arg \inf_{\mathbf{X} \in \mathcal{X}} g(\mathbf{X}) - [h(\mathbf{X}^{[t]}) + \langle \mathbf{X} - \mathbf{X}^{[t]}, \partial_{\mathbf{X}^{[t]}} h \rangle]$$

➤ converge to a critical point with speed  $\mathcal{O}(1/t)$



# Numerical results

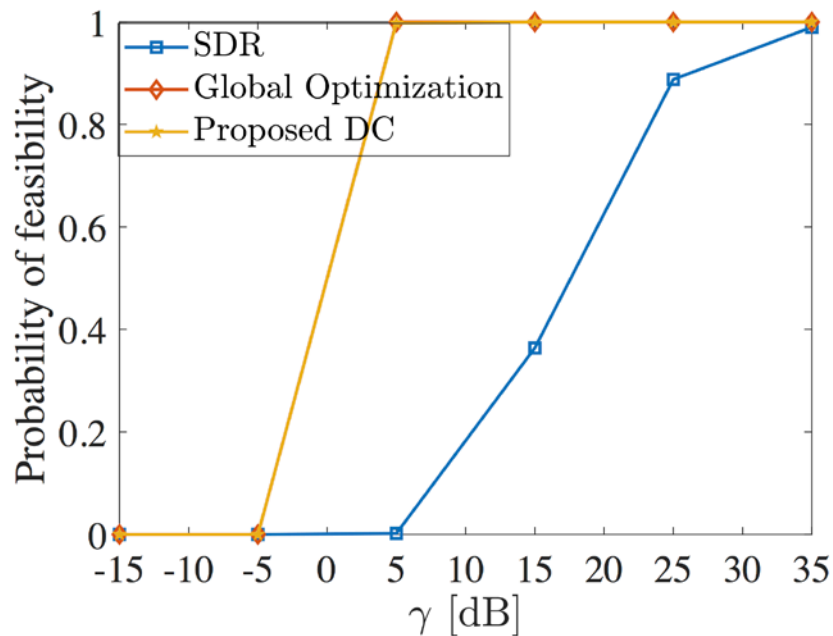
- Convergence of the proposed DC algorithm for problem  $\mathcal{P}_{S2}$



$$\begin{aligned} \mathcal{P}_{S2} : & \underset{\mathbf{M}}{\text{minimize}} \quad \text{Tr}(\mathbf{M}) - \|\mathbf{M}\|_2 \\ & \text{subject to} \quad \text{Tr}(\mathbf{M}) - \gamma_i \mathbf{h}_i^H \mathbf{M} \mathbf{h}_i \leq 0, \\ & \quad \quad \quad \mathbf{M} \succeq \mathbf{0}, \quad \text{Tr}(\mathbf{M}) \geq 1 \end{aligned}$$

# Numerical results

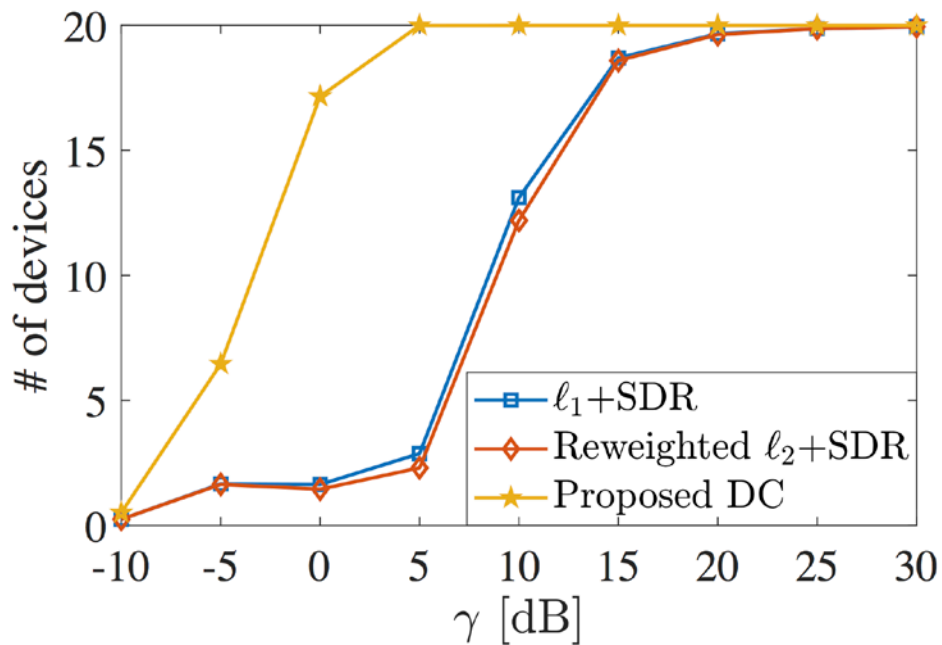
- Probability of feasibility with different algorithms



$$\begin{aligned} & \text{find } \mathbf{M} \\ & \text{subject to } \text{Tr}(\mathbf{M}) - \gamma_i \mathbf{h}_i^H \mathbf{M} \mathbf{h}_i \leq 0, \forall i \in \mathcal{S}^{[k]} \\ & \mathbf{M} \succeq \mathbf{0}, \text{Tr}(\mathbf{M}) \geq 1, \text{rank}(\mathbf{M}) = 1 \end{aligned}$$

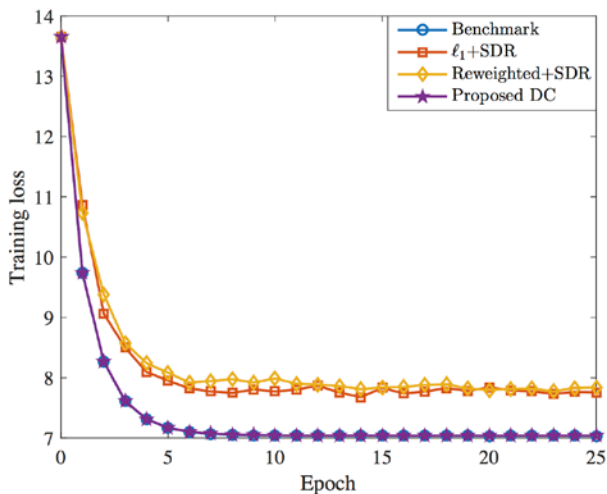
# Numerical results

- Average number of selected devices with different algorithms

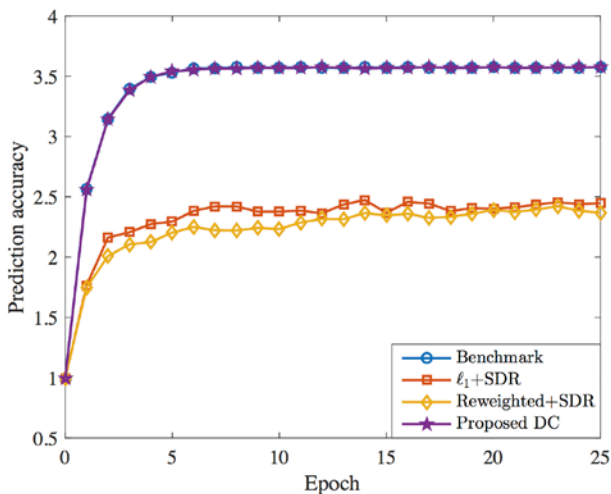


# Numerical results

- Performance of proposed fast model aggregation in federated learning
  - Training an SVM classifier on CIFAR-10 dataset



(a) Training loss



(b) Relative prediction accuracy

# Concluding remarks

- **Wireless communication meets machine learning**
  - Over-the-air computation for fast model aggregation
- **Sparse and low-rank optimization framework**
  - Joint device selection and beamforming design
- **A unified DC programming framework**
  - DC representation for sparse and low-rank functions

# Future directions

- **Federated learning**

- security, provable guarantees, ...

- **Over-the-air computation**

- channel uncertainty, synchronization, ...

- **Sparse and low-rank optimization via DC programming**

- optimality, scalability, ...

# To learn more...

- **Papers:**
- K. Yang, T. Jiang, Y. Shi, and Z. Ding, “Federated learning via over-the-air computation,” *IEEE Trans. Wireless Commun.*, DOI10.1109/TWC.2019.2961673, Jan. 2020.
- K. Yang, T. Jiang, Y. Shi, and Z. Ding, “Federated learning based on over-the-air computation,” in *Proc. IEEE Int. Conf. Commun. (ICC)*, Shanghai, China, May 2019.

<http://shiyuanming.github.io/home.html>

*Thanks*