

Privacy-Preserving Deep Learning

Reza Shokri and Vitaly Shmatikov

(UT Austin and Cornell Tech)

Machine Intelligence LANDSCAPE

CORE TECHNOLOGIES

ARTIFICIAL INTELLIGENCE



DEEP LEARNING



MACHINE LEARNING



NLP PLATFORMS



PREDICTIVE APIs



IMAGE RECOGNITION



SPEECH RECOGNITION



RETHINKING ENTERPRISE

SALES



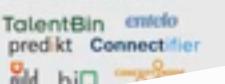
SECURITY / AUTHENTICATION



FRAUD DETECTION



HR / RECRUITING



MARKETING



PERSONAL ASSISTANT



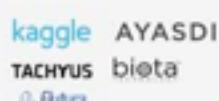
INTELLIGENCE TOOLS



ADTECH



OIL AND GAS



Tech 2015: Deep Learning And Machine Intelligence Will Eat The World

RETHINKING HUMANS / HCI

AUGMENTED REALITY



GESTURAL COMPUTING



ROBOTICS



EMOTIONALrecognition



HARDWARE



SUPPORTING TECHNOLOGIES

DATA PREP



DATA COLLECTION

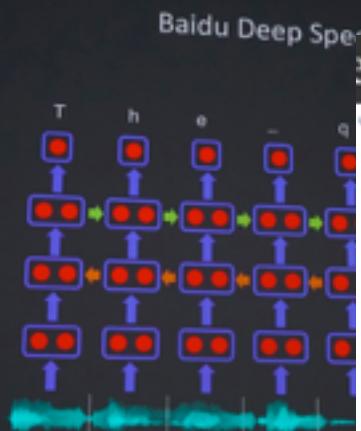




Robert Hof
Contributor

TECH 12/18/2014 @ 9:00AM | 48,377 views

Baidu Announces Breakthrough In Speech Recognition, Claiming To Top Google And Apple



Baidu Research



Siri



Michael Thomsen
Contributor

TECH 2/15/2015 @ 1:05PM | 4,995 views

Microsoft's Deep Learning Project Outperforms Humans In Image Recognition



6:1 horse cart
1: minibus
3: exca
4: stretcher
5: half track



6:1 birdhouse
1: barnhouse
2: sliding door
3: window screen
4: mailbox
5: spot



6:1 forklift
1: garbage truck
2: forklift
3: tow truck
4: trailer truck
5: go-kart



6:1 letter opener
1: drumstick
2: candle
3: wooden spoon
4: spatula
5: ladle



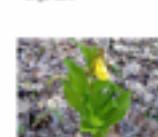
6:1 letter opener
1: Band Aid
2: ruler
3: rubber eraser
4: pencil box
5: wallet



6:1 covey
1: indigo bunting
3: toro
4: walking stick
5: currant apple



6:1 komondor
1: puma
3: llama
4: mobile phone
5: Old English sheepdog



6:1 yellow lady's slipper
1: yellow ladies' slipper
2: slug
3: hem-of-the-woods
4: stinkhorn
5: coral fungus



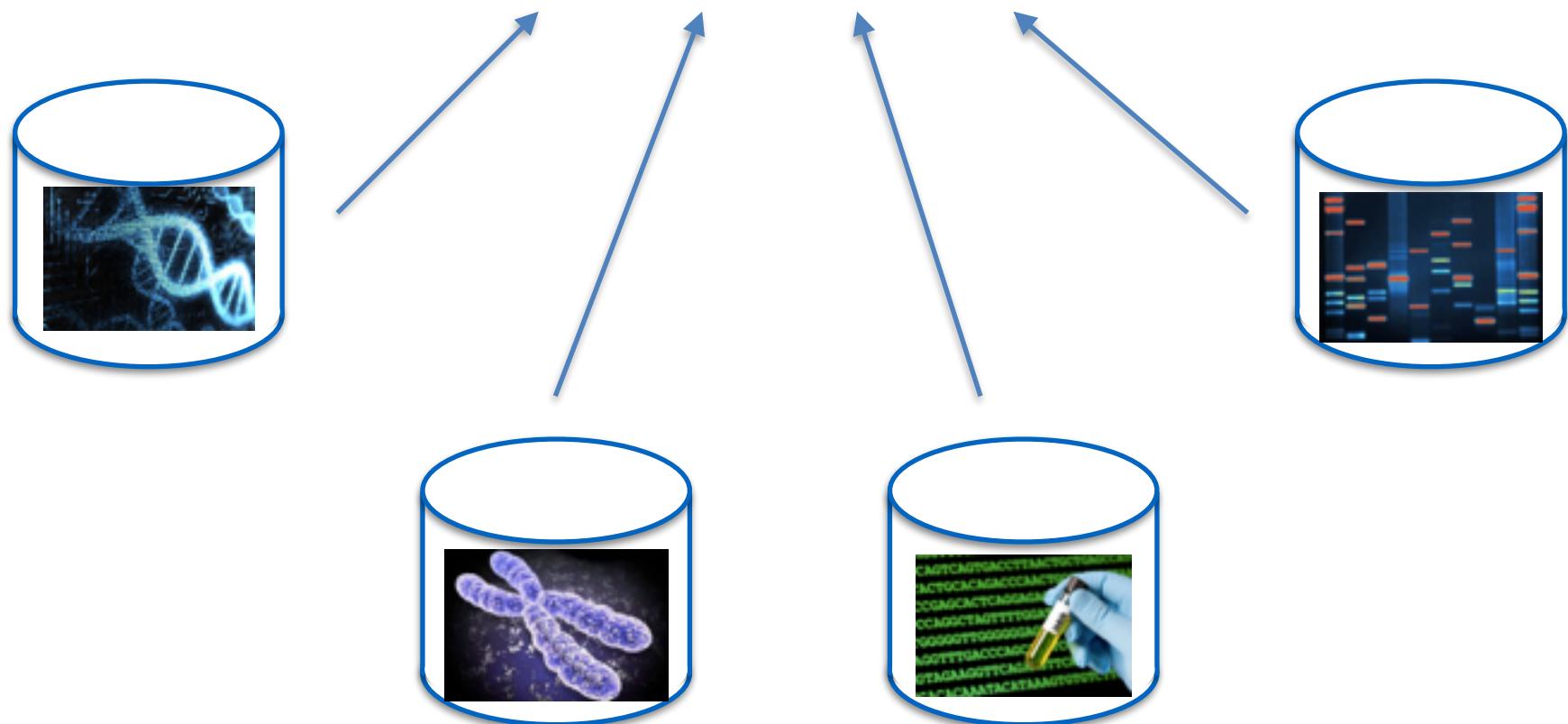
6:1 spotlight
1: grand piano
2: folding chair



6:1 spotlight
1: acoustic guitar
2: stage

Images used to test Microsoft's deep learning AI. Image via Microsoft.

Deep Learning Today



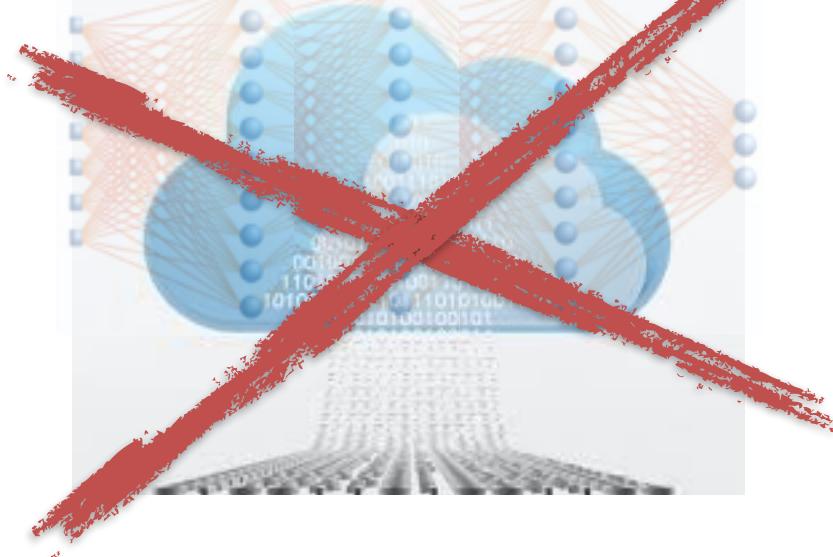
Privacy Concerns

- Training data is very sensitive
 - speech, photo images, written documents, might contain sensitive private information
- Users have no control over the learning objective
- Using trained networks requires users to share their private data with service providers

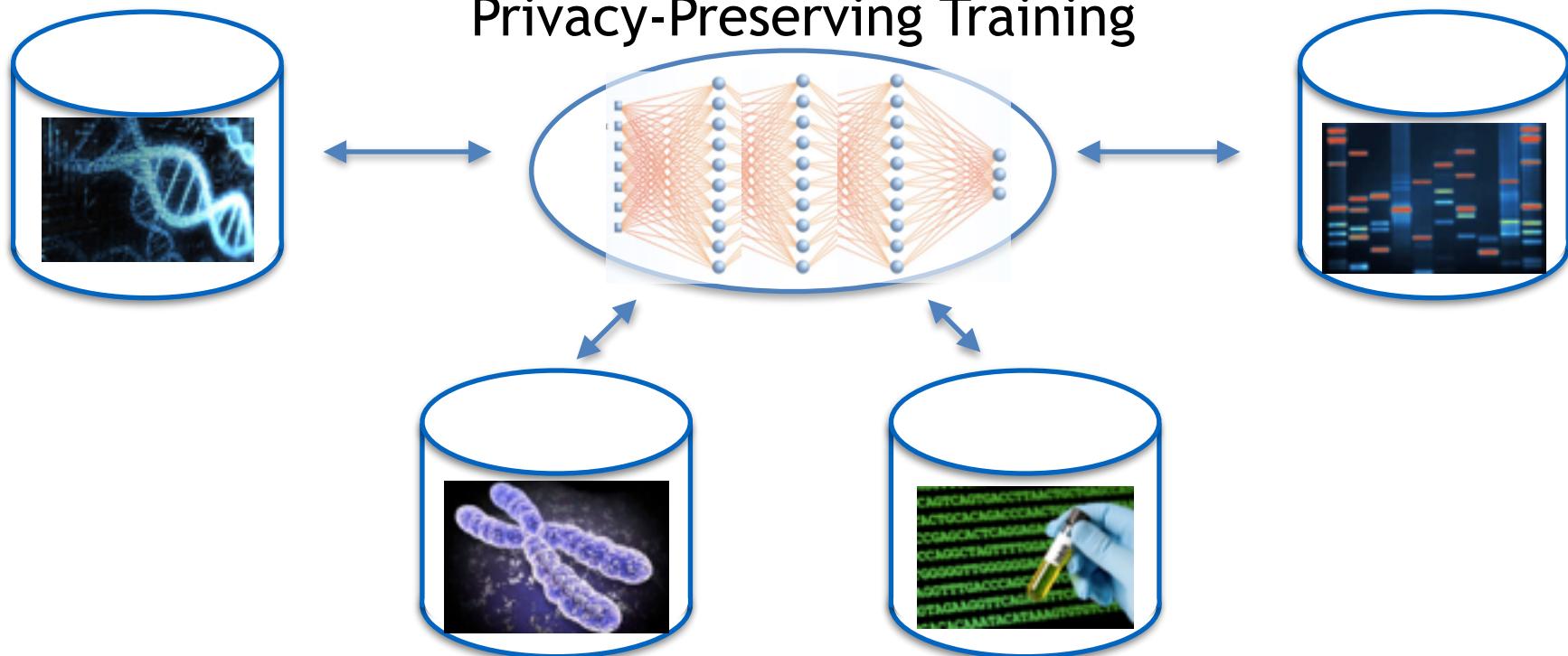
Possible Consequences

- Users' data might be used in **wrong context**
 - Compromises and data breaches
 - Inference of sensitive information
 - Training of intrusive models
- Holders of sensitive data cannot benefit from large-scale deep learning because they may not share their datasets for training
 - Biomedical researchers?
 - Social scientists?

Our Objective



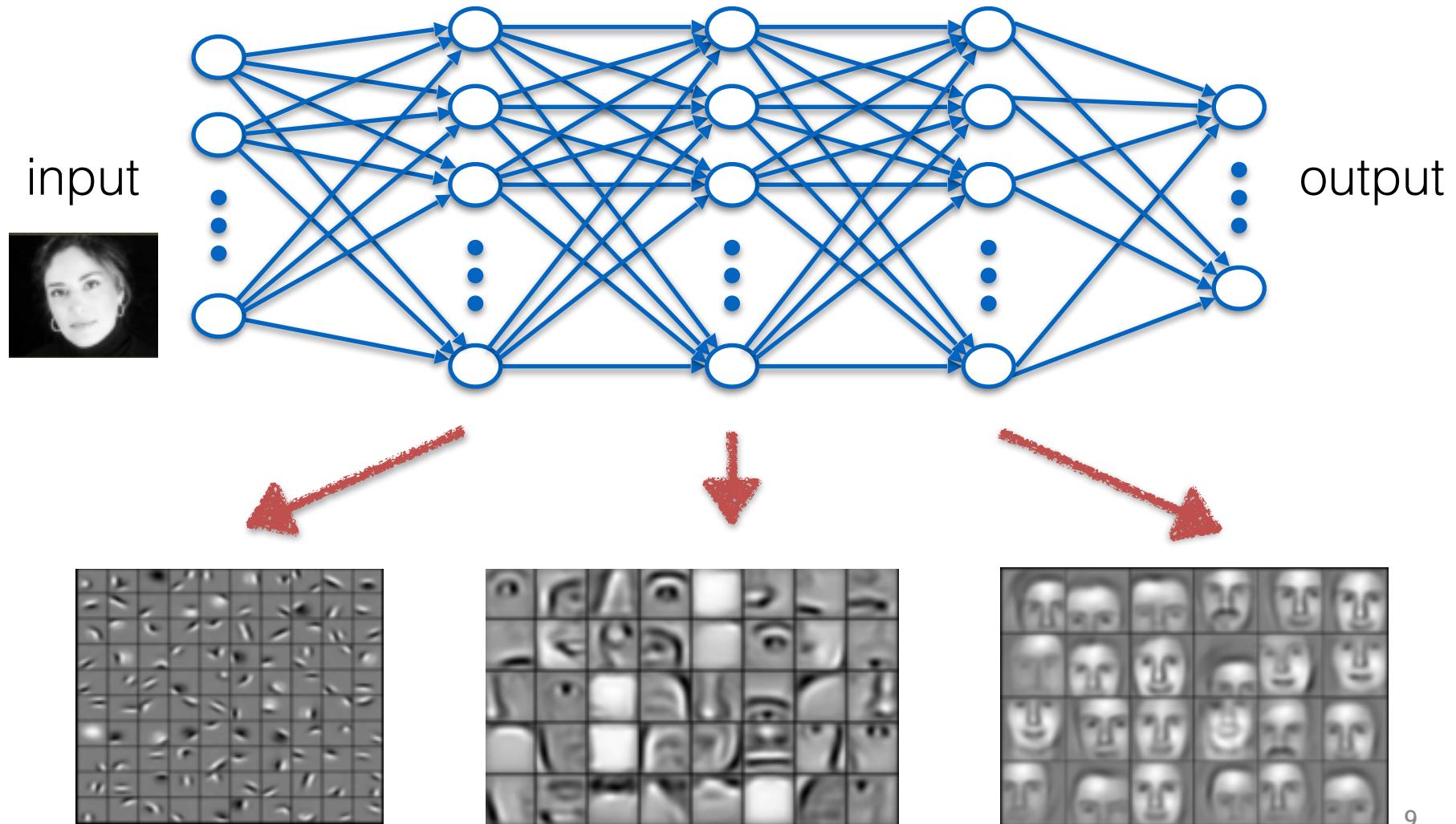
Privacy-Preserving Training



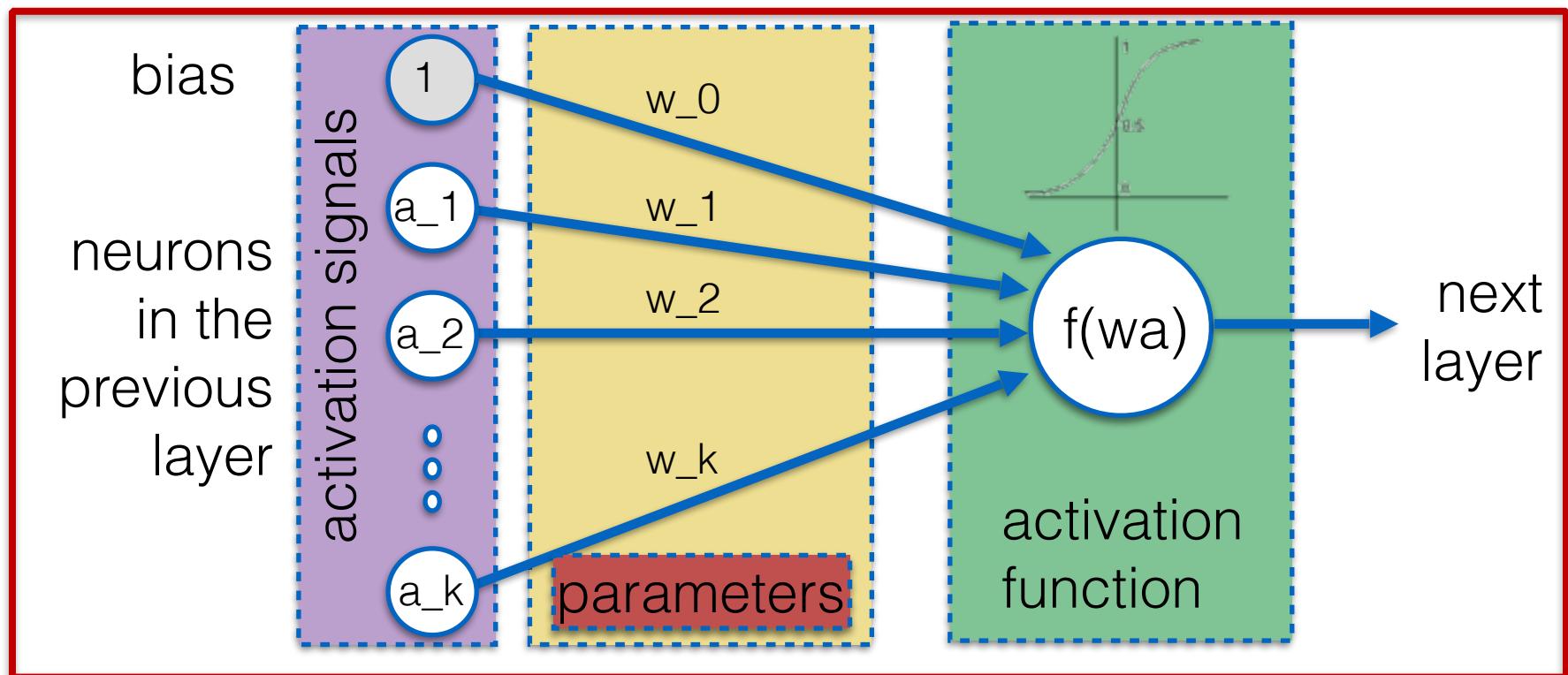
Private data remain private

Background

Deep Neural Networks

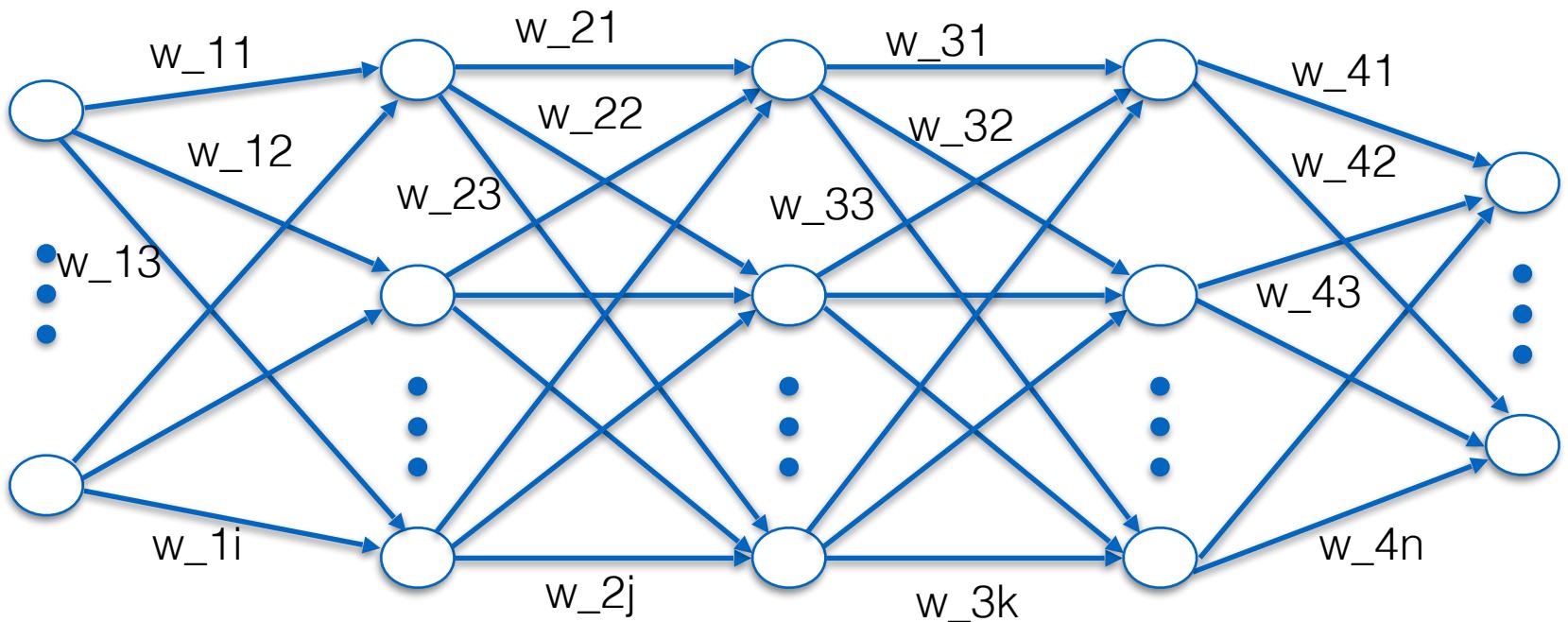


Deep Neural Networks



Learn parameters using
Stochastic Gradient Descent (SGD) algorithm

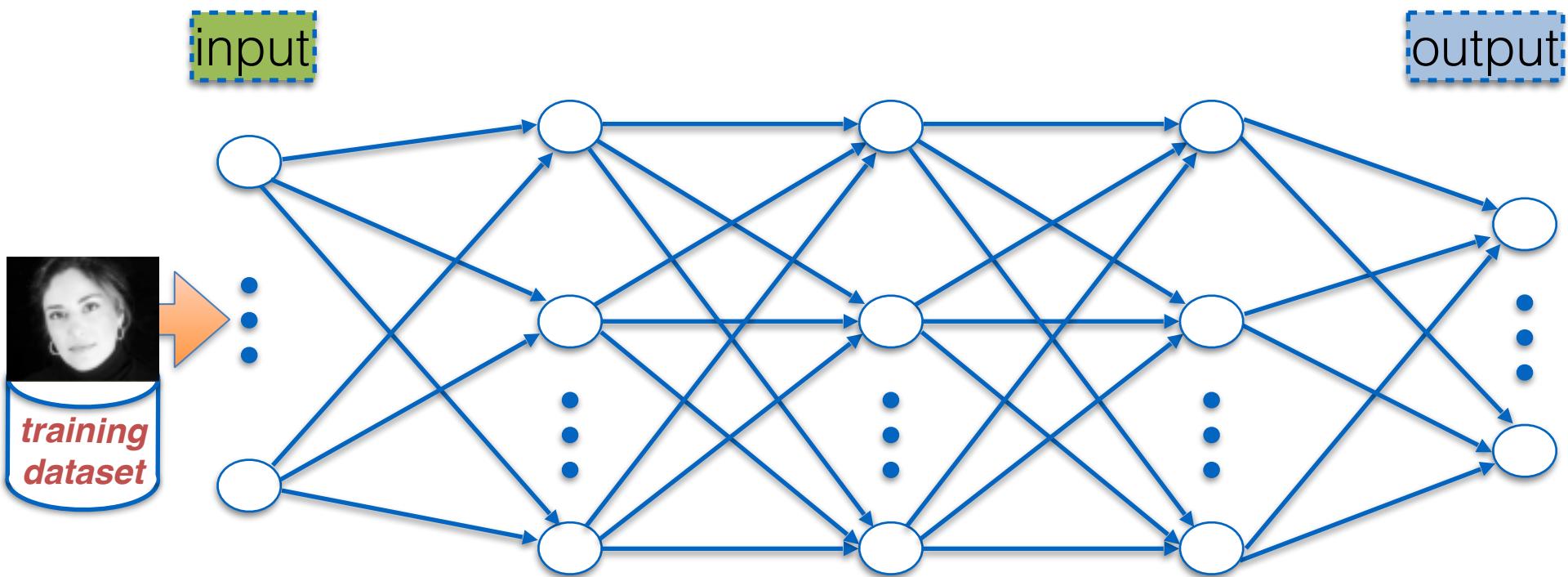
Parameter Training using SGD



- Find parameters that minimize the classification error

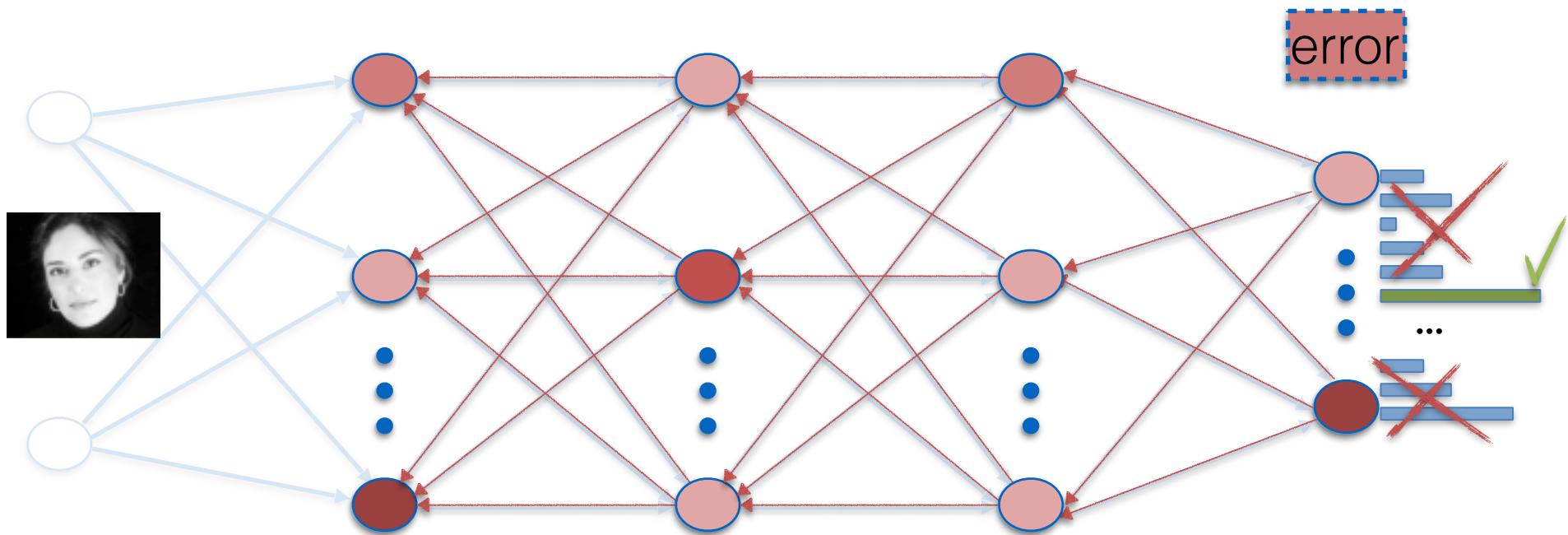
Parameter Training using SGD

1) Feed-Forward

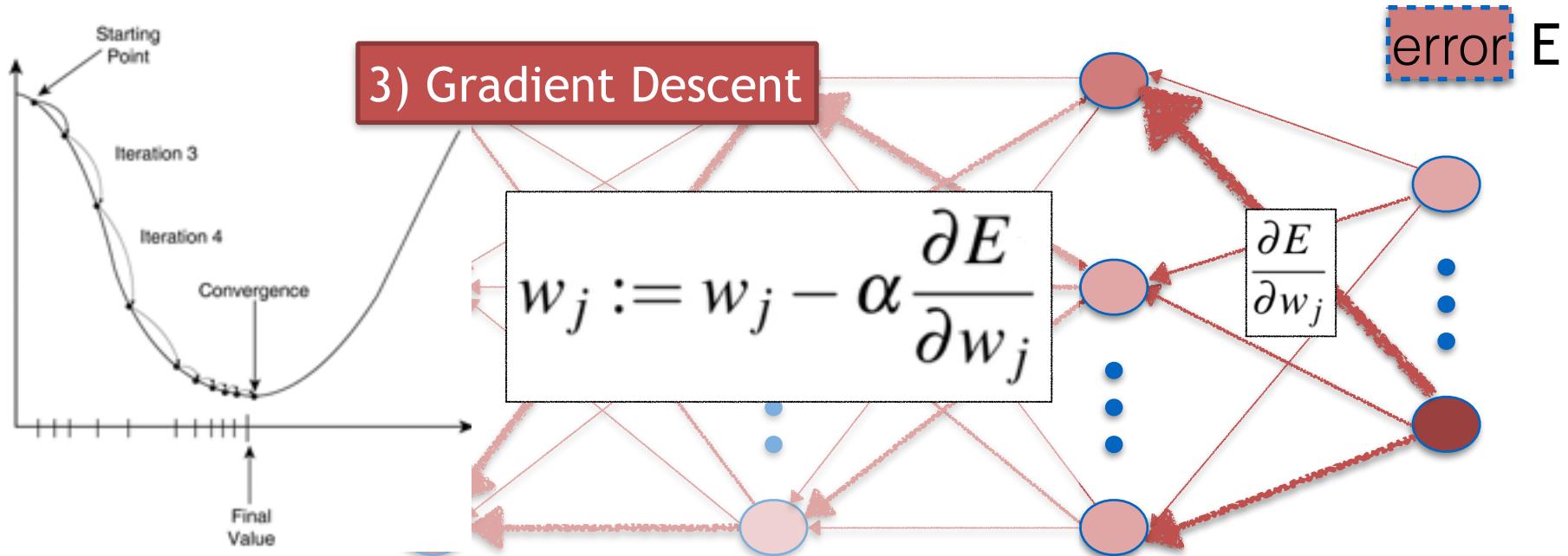


Parameter Training using SGD

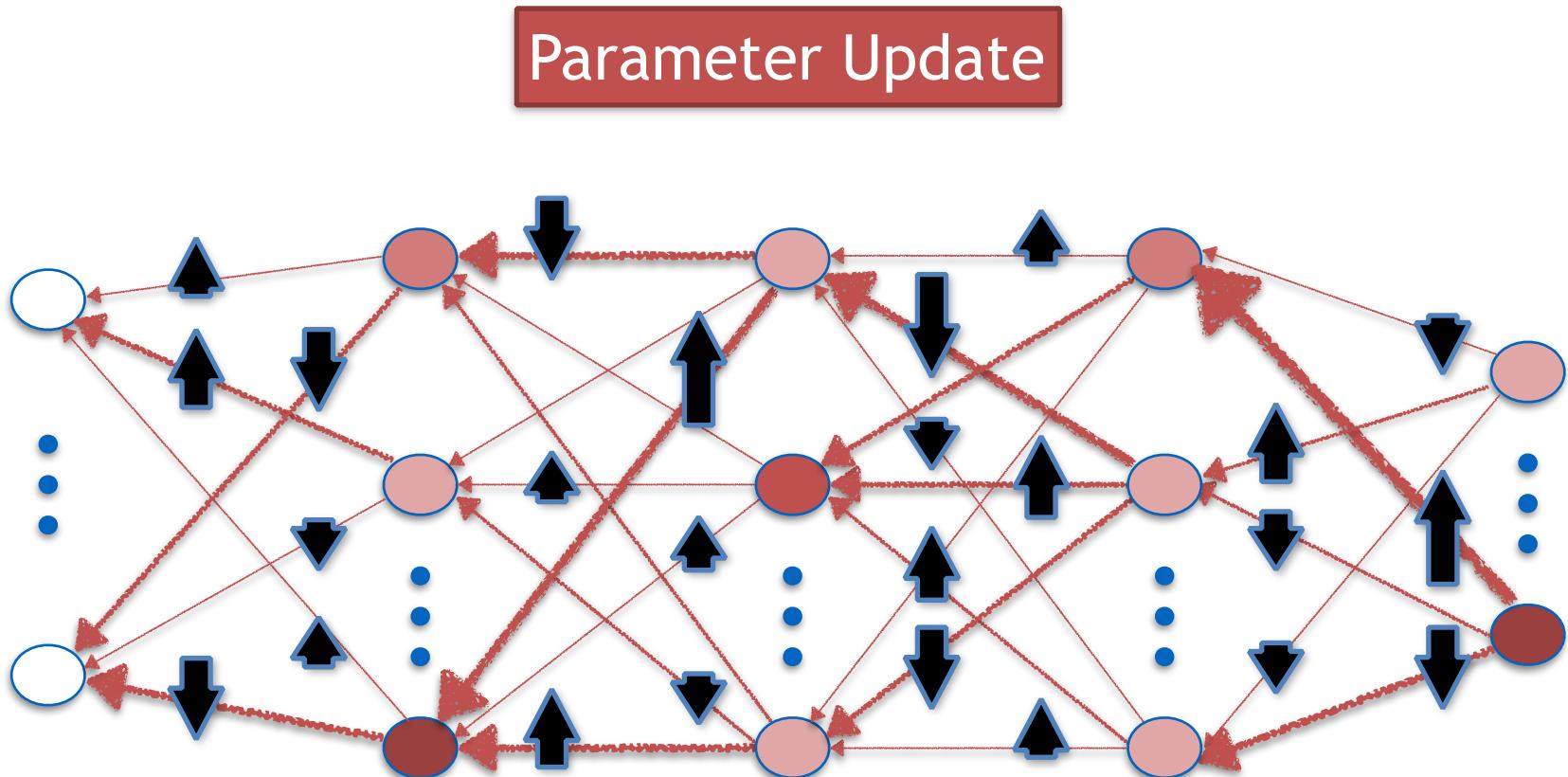
2) Back-propagation



Parameter Training using SGD



Parameter Training using SGD

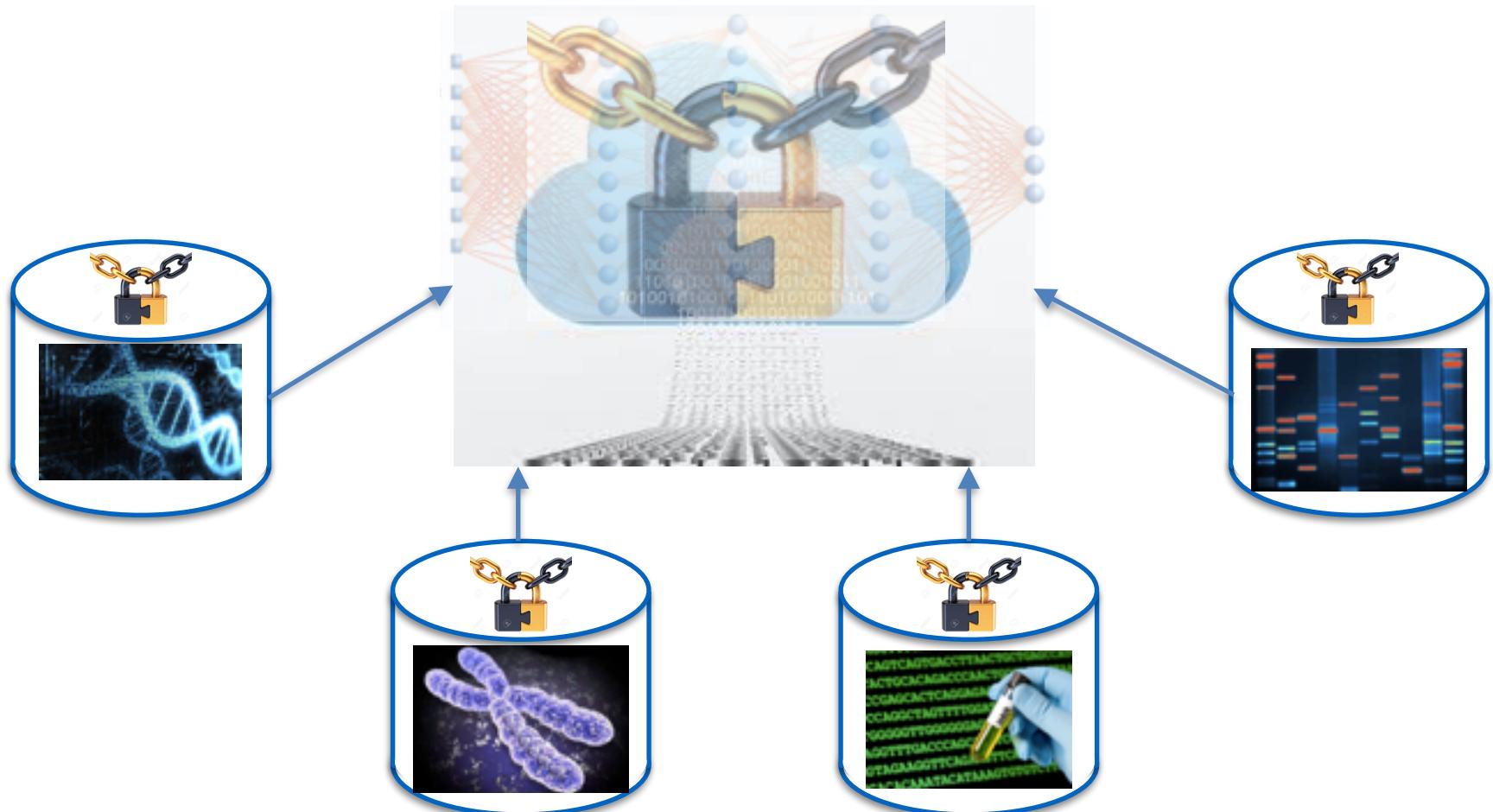


Repeat for new batches of training data

Privacy-Preserving Deep Learning

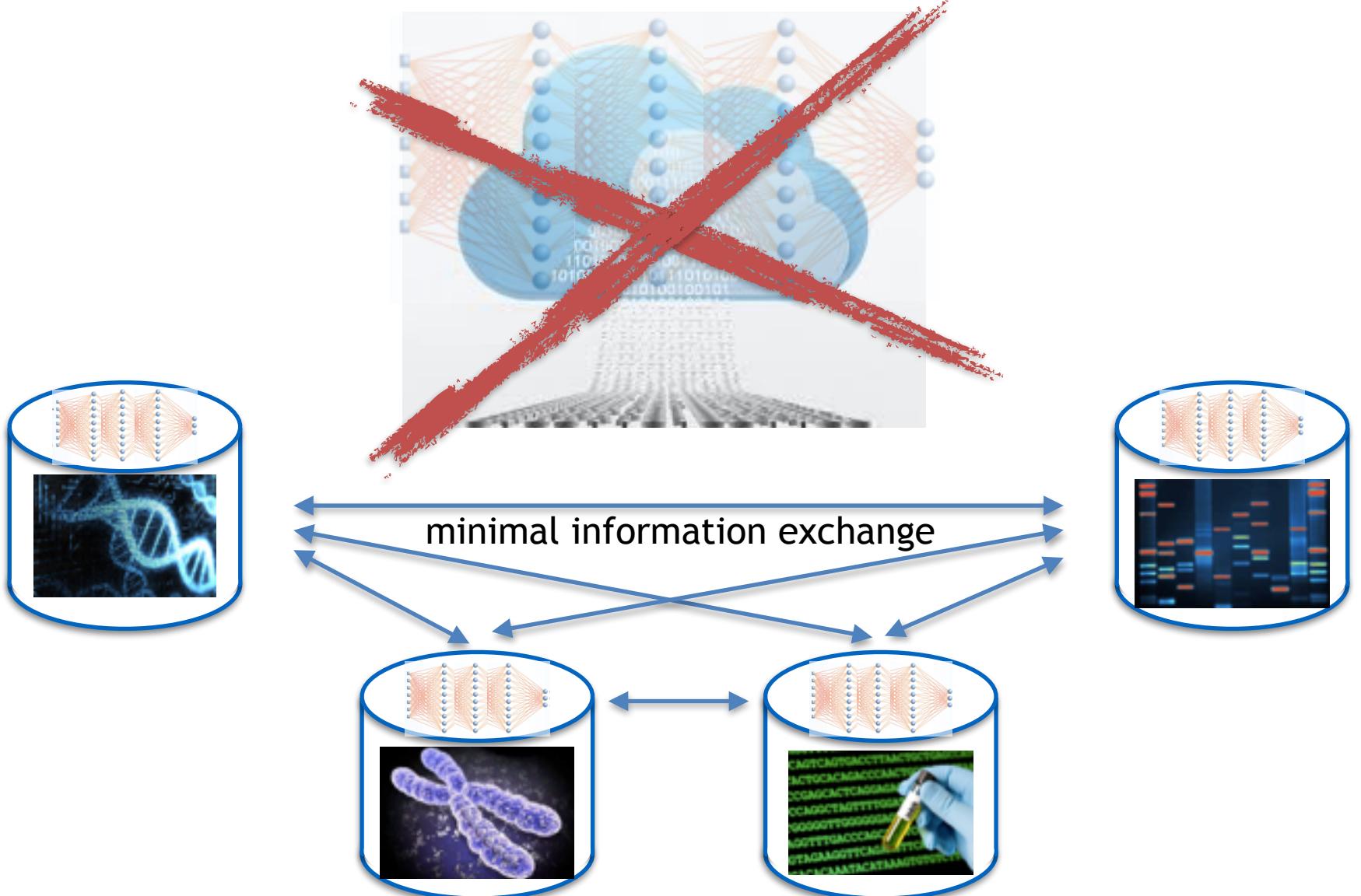
Prior Work

Secure multi-party computation



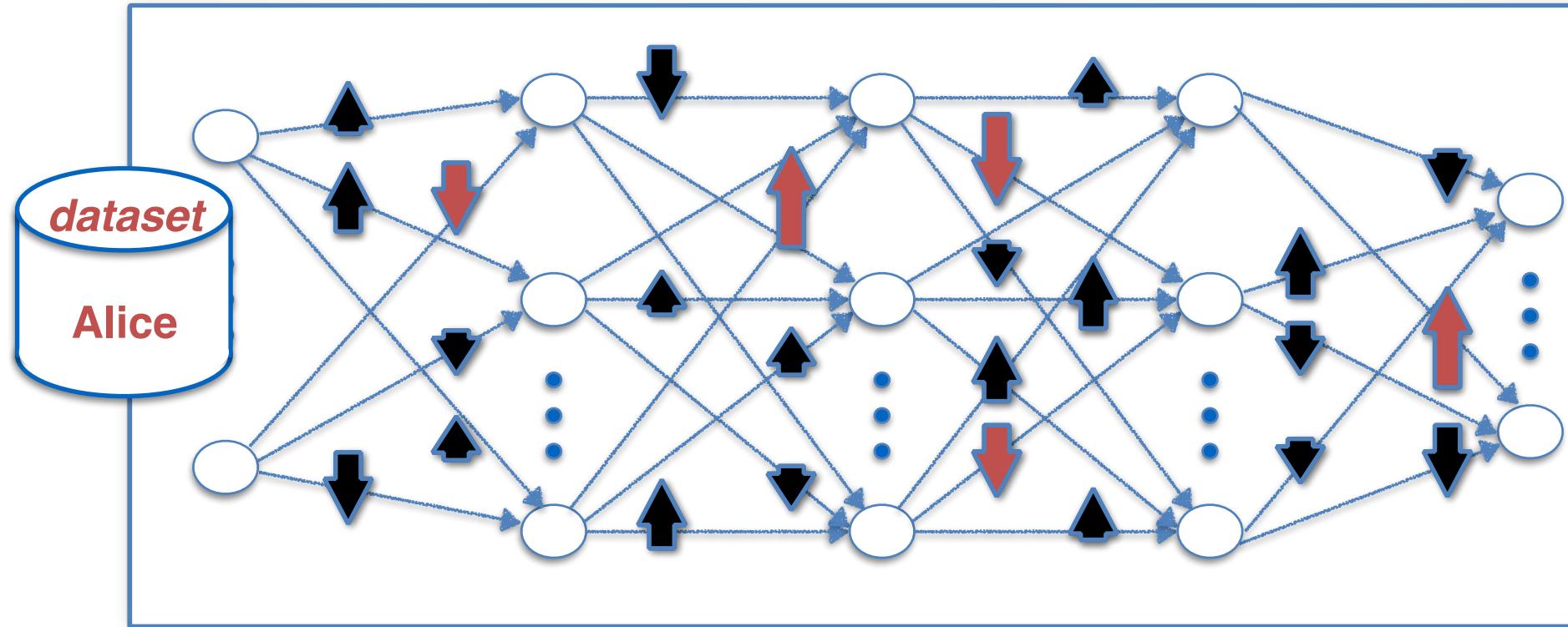
J. Yuan, S. Yu, "Privacy Preserving Back-Propagation Neural Network Learning Made Practical with Cloud Computing", IEEE TPDS, 2014.

A. Bansal, T. Chen, S. Zhong, "Privacy preserving Back-propagation neural network learning over arbitrarily partitioned data", Neural Computing and Applications, 2011



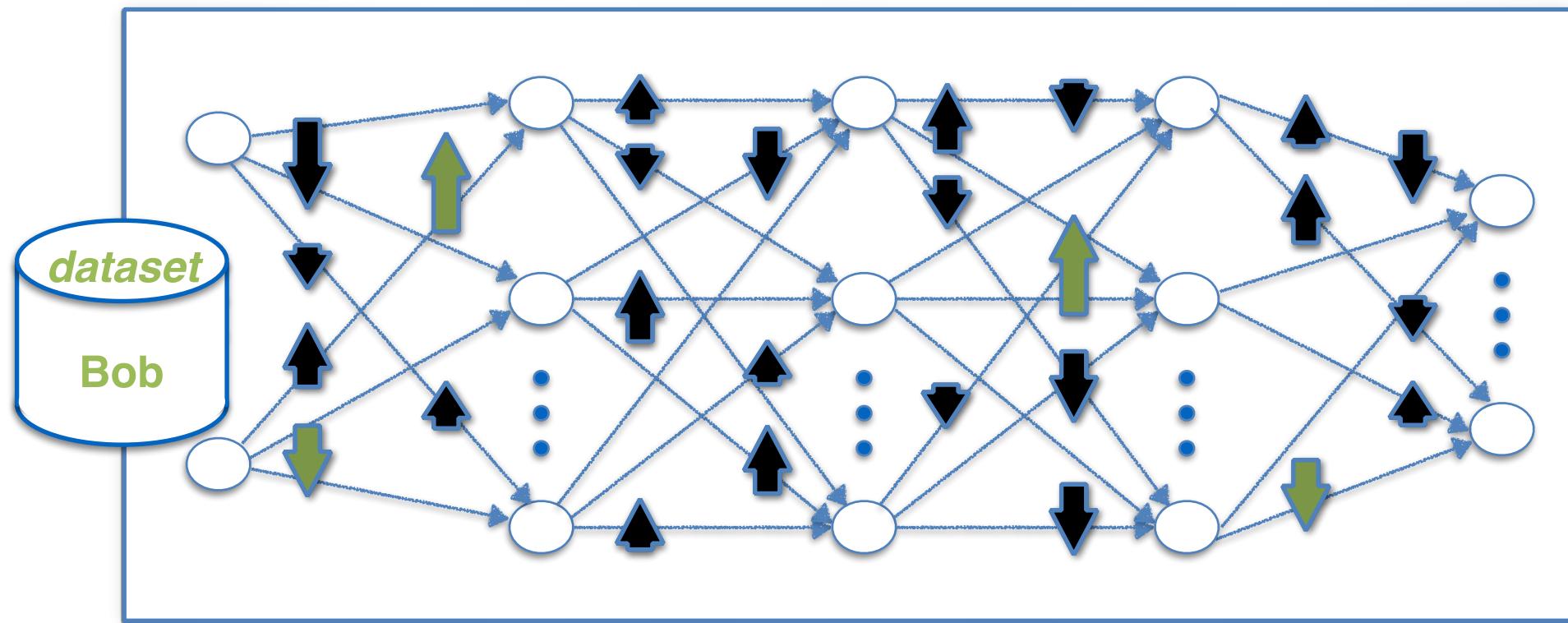
Distributed Selective SGD (DSSGD)

Selective SGD

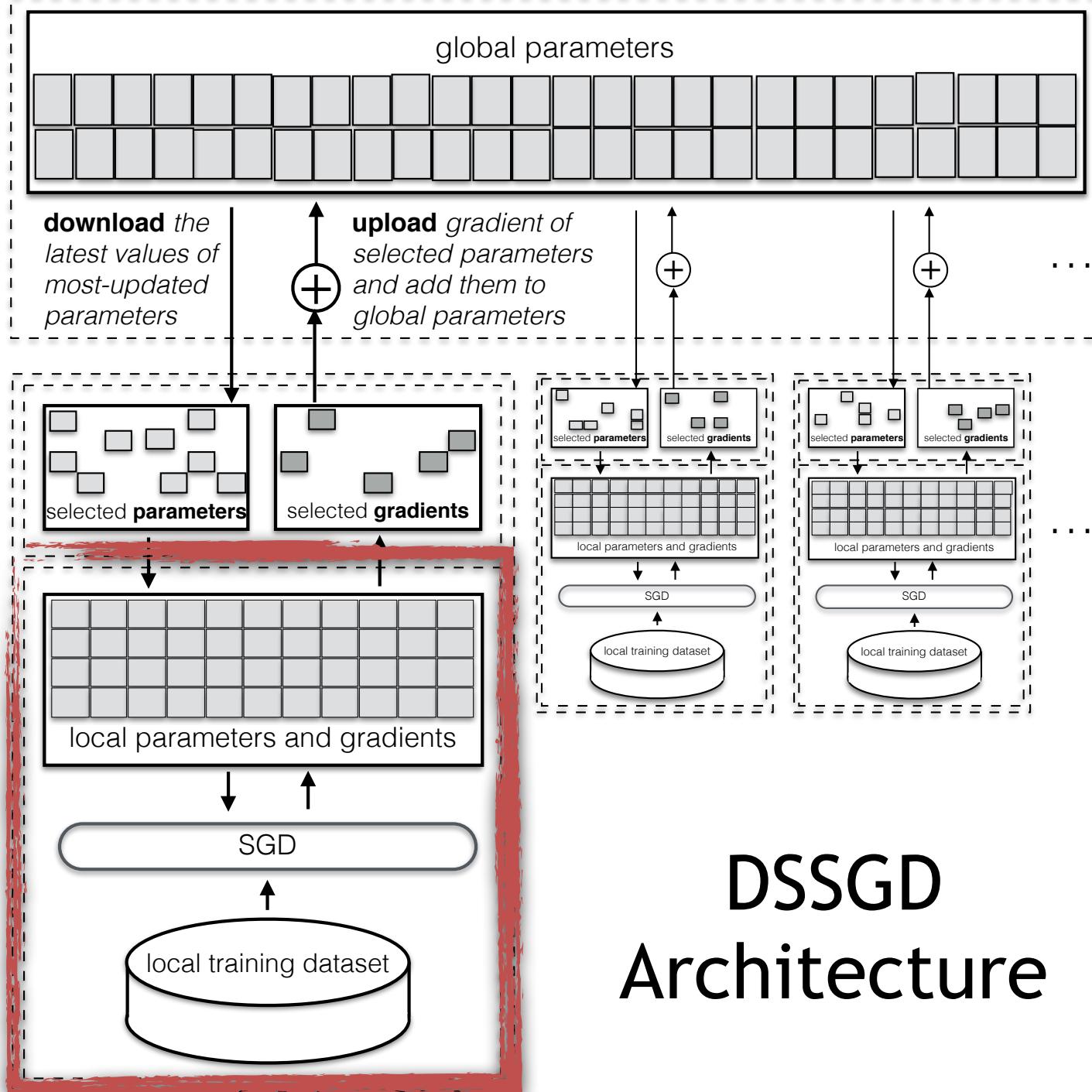


Share with others

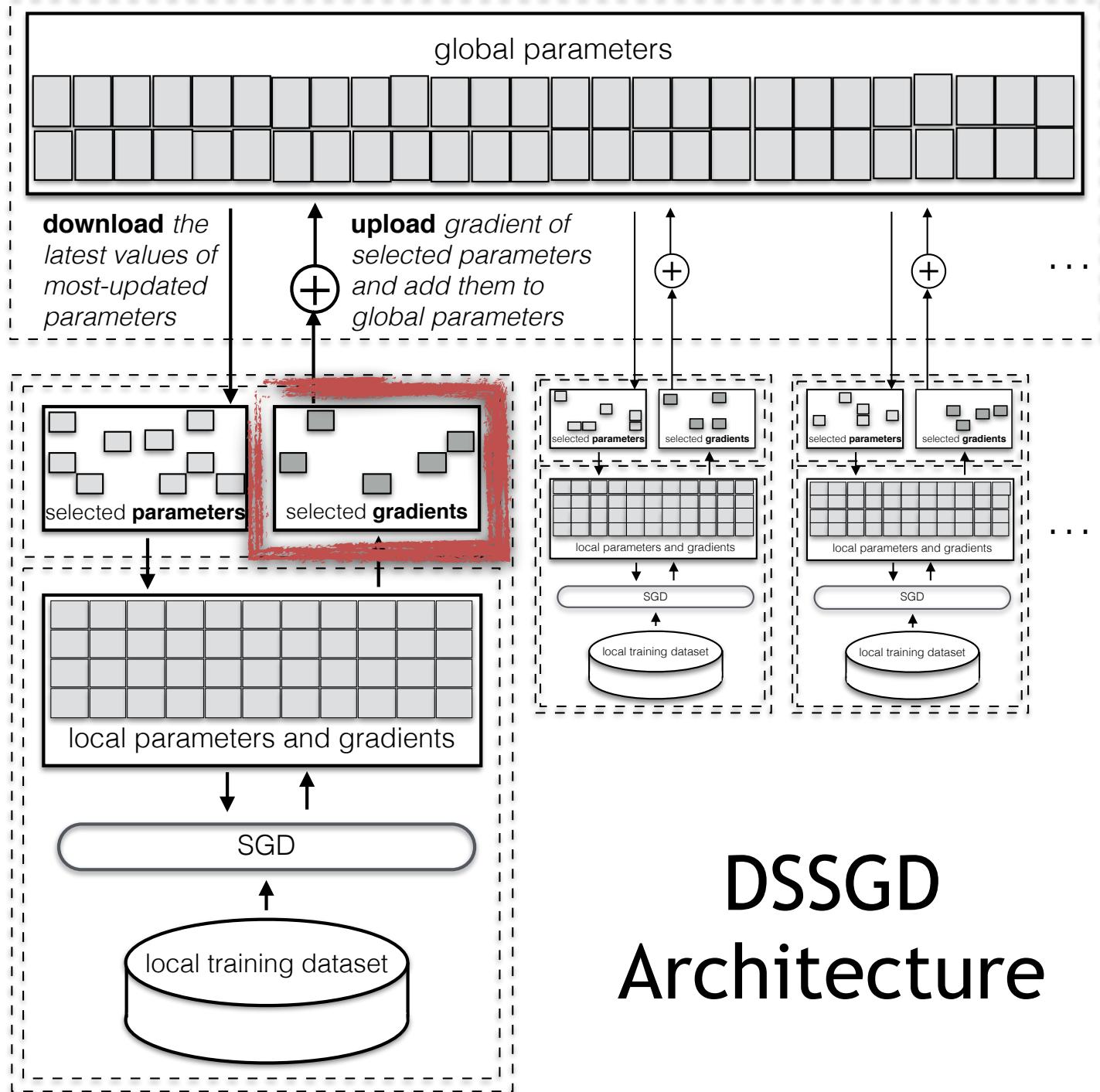
Selective SGD



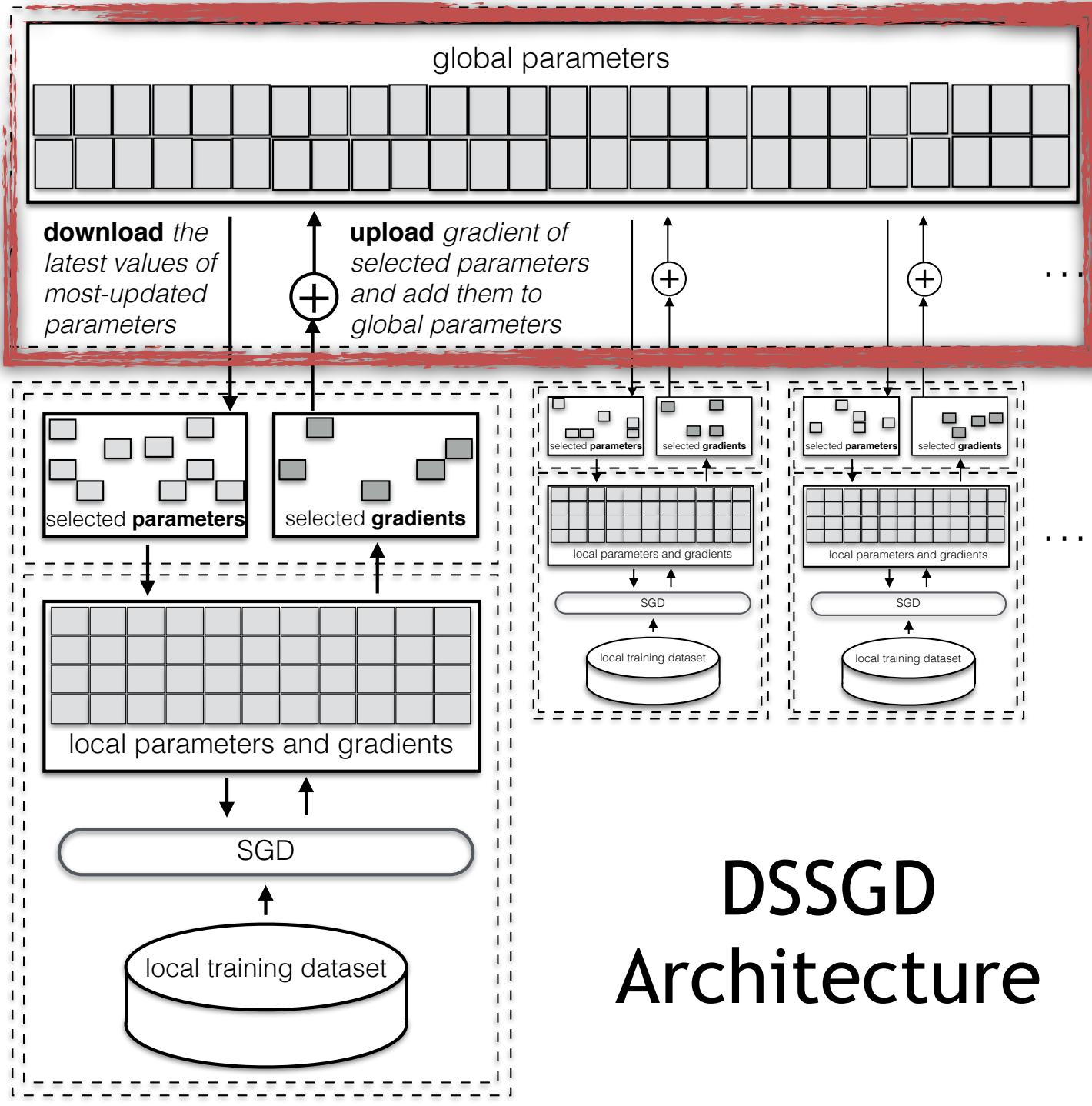
Share with others



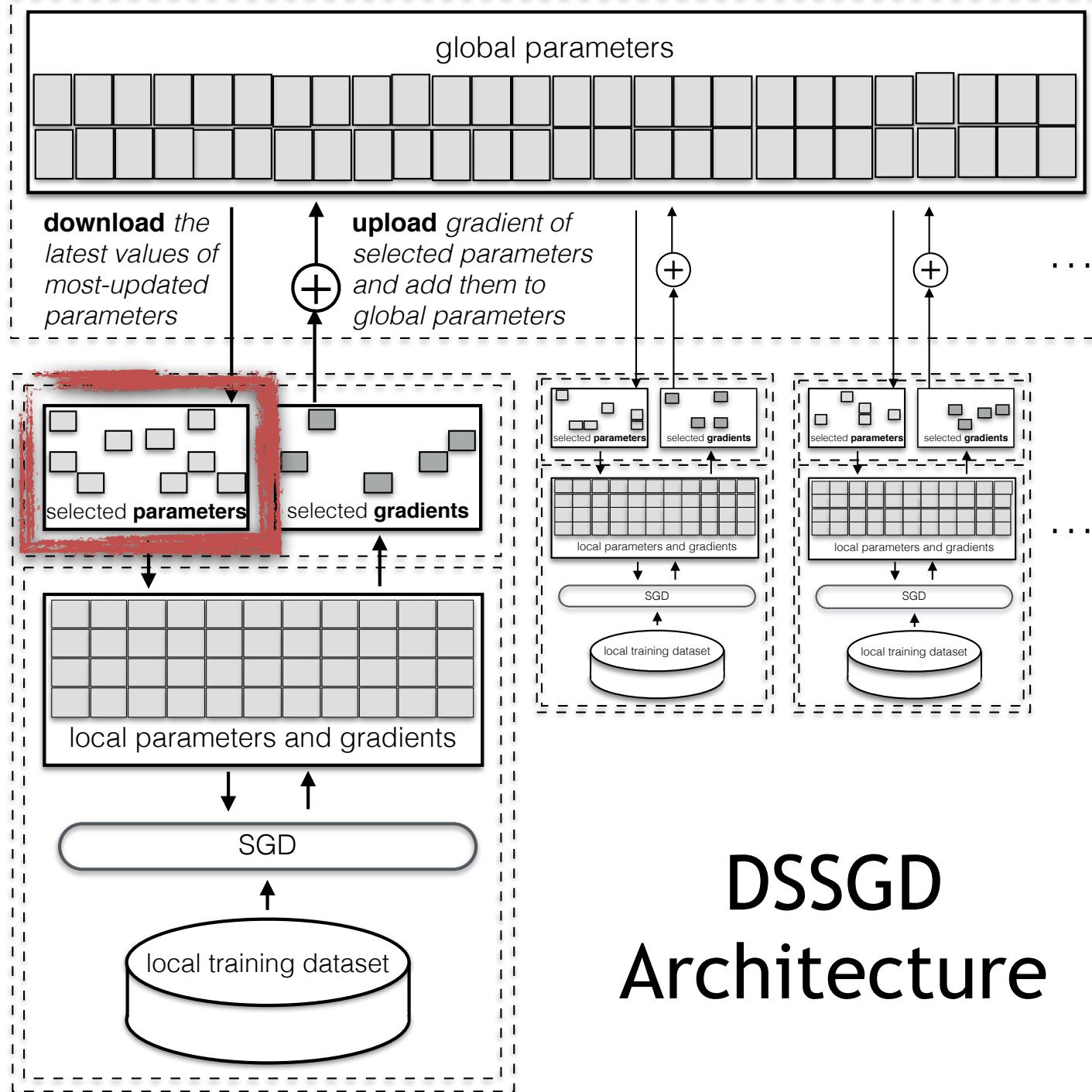
DSSGD Architecture



DSSGD Architecture



DSSGD Architecture



DSSGD Architecture

Distributed Selective SGD

- Local training, Global convergence
- High training stochasticity
- Less overfitting

Evaluation Datasets

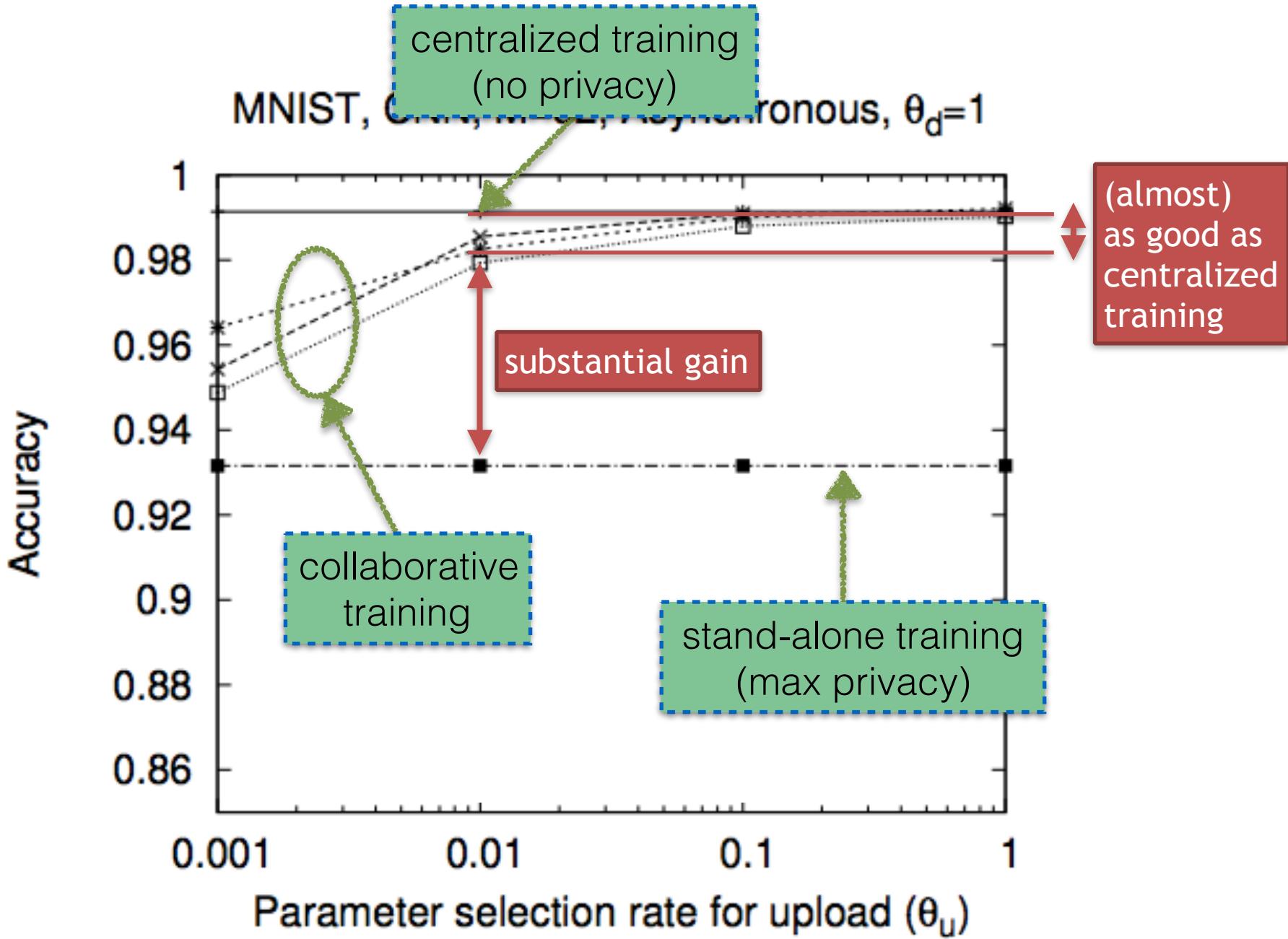
MNIST



SVHN



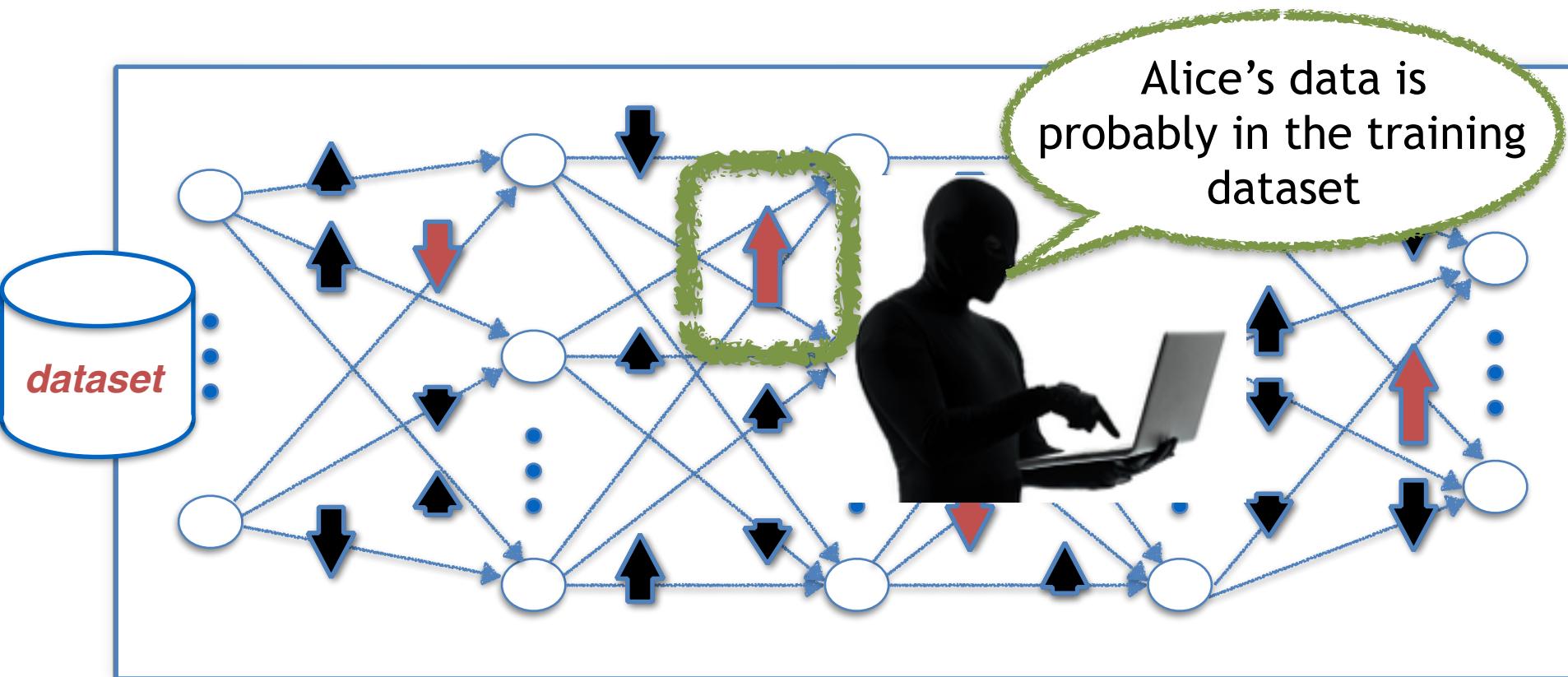
Task: Find the digit in the image (classify into one of 10 classes)



Privacy Properties

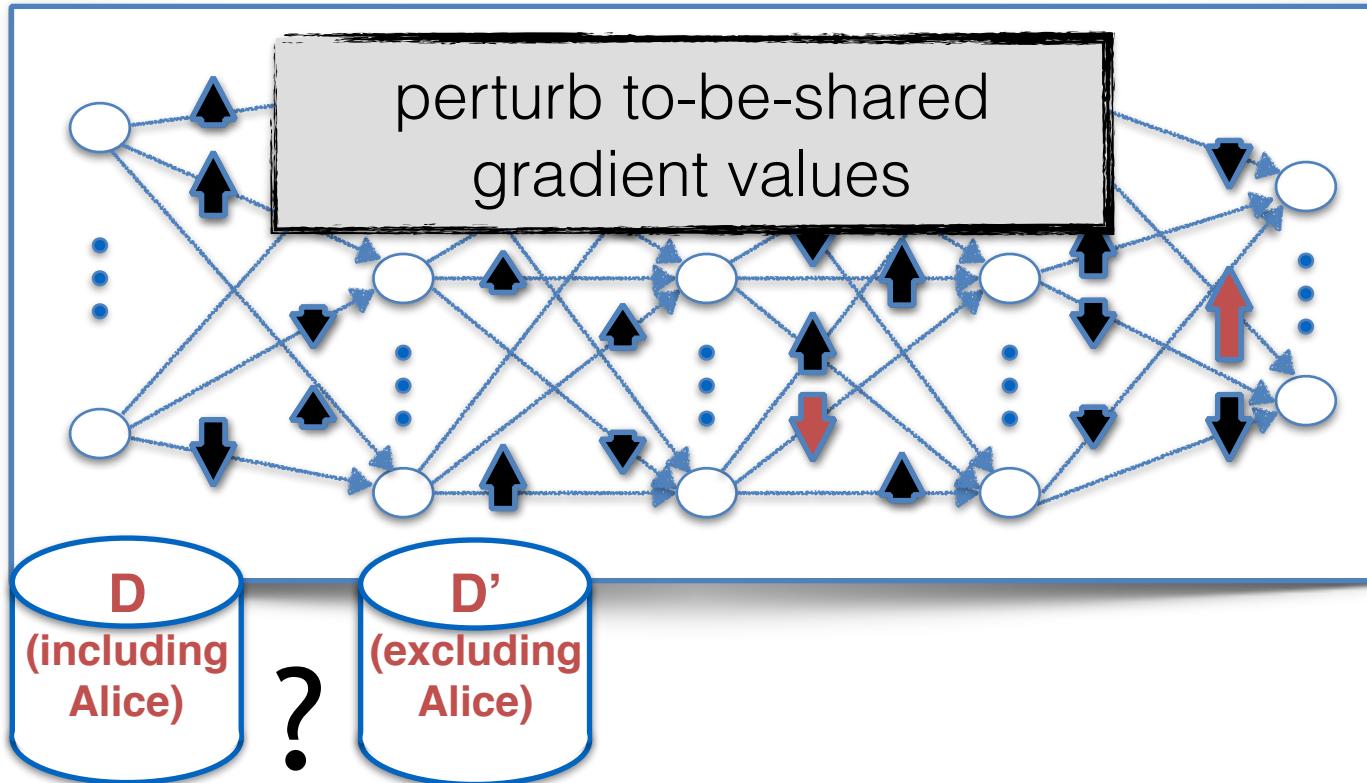
- Participants' datasets remain private
- Full control over parameter selection
- Known learning objective
- Resulting model available to all parties

Indirect Information Leakage *through gradient sharing*



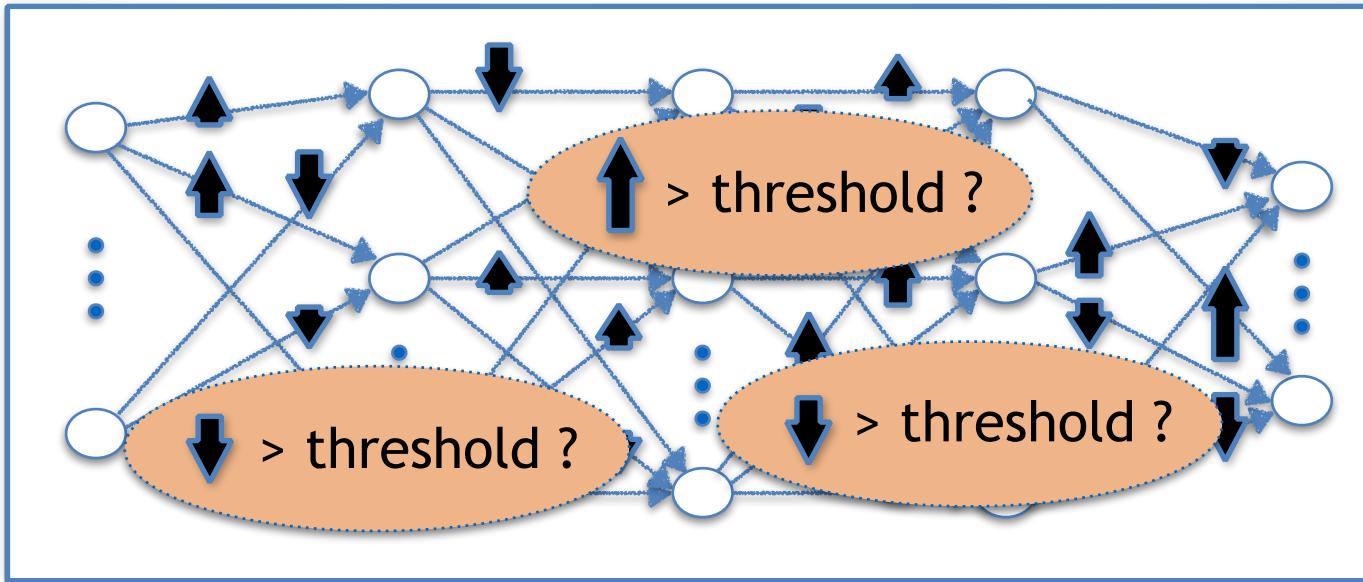
Limit Indirect Info. Leakage

- *Differentially Private* parameter selection and gradient sharing



Sparse Vector Technique

- Select a small fraction of (perturbed) gradients that are above a given (perturbed) threshold

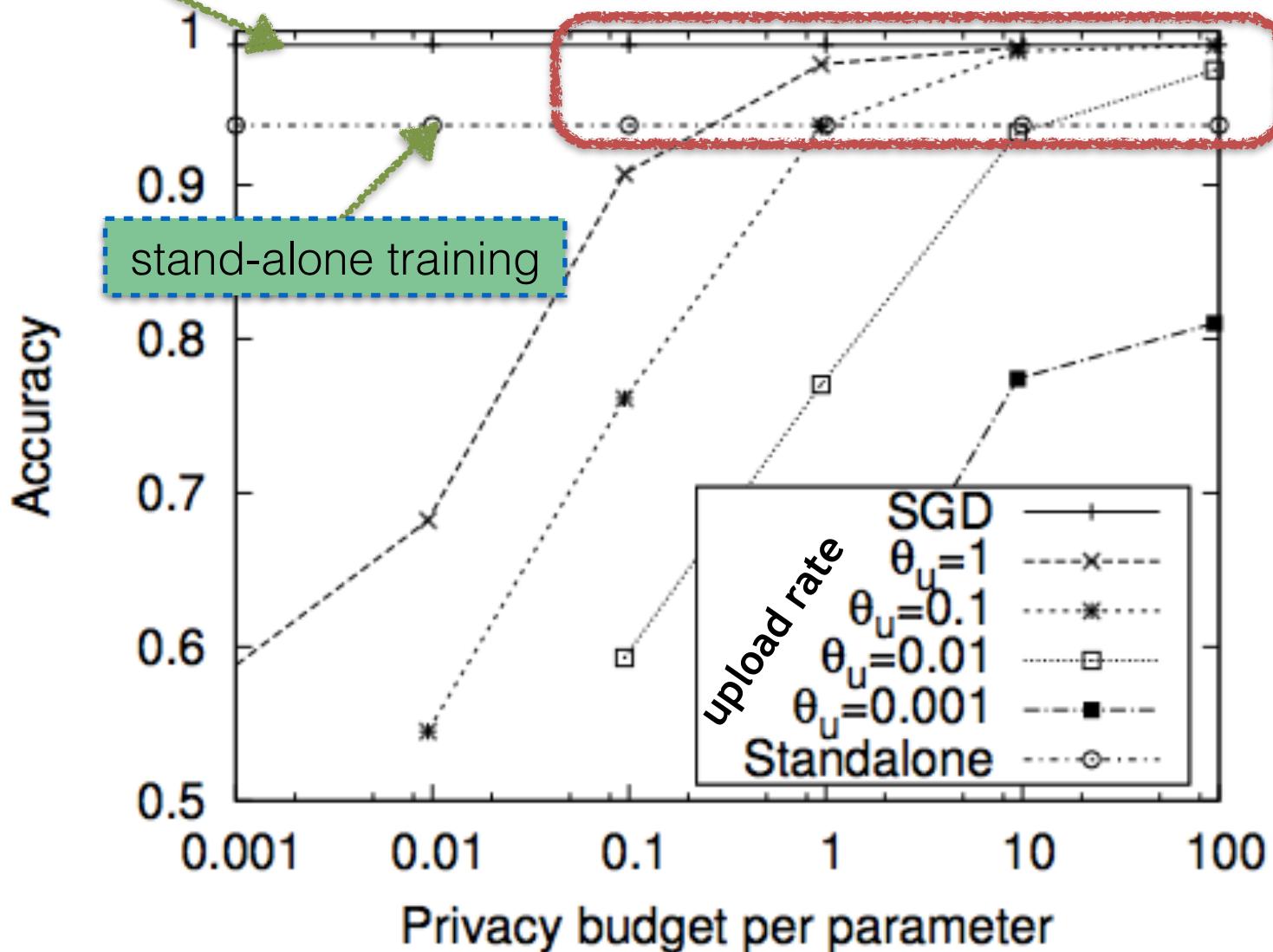


C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Theoretical Computer Science*, 9(3-4):211–407, 2013.

differentially-private
comparison

centralized training

MNIST, CNN, Round Robin, N=150,
 $\theta_d=1$, $\gamma=0.001$, $\tau=0.0001$



Conclusions

- Massive data collection required for deep learning presents substantial privacy risks
- We design privacy-preserving deep learning system
 - Holders of sensitive data collaborate to build powerful models without sharing their data
- Differentially-private information exchange protocol
 - Prevents indirect leakage about participants' private datasets
- Accuracy of our system
 - is substantially better than that of standalone learning
 - is very close to accuracy of centralized SGD