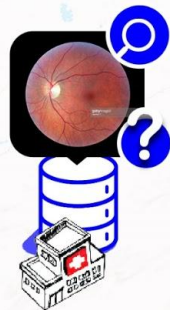


Split Learning

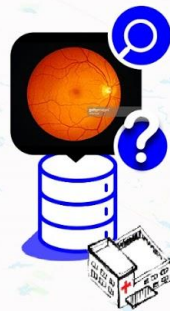
A resource efficient distributed deep learning method without sensitive data sharing

Praneeth Vepakomma
vepakom@mit.edu

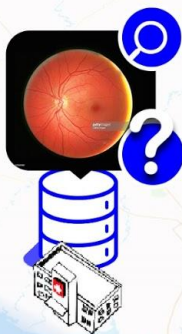
'Invisible' Health Image Data



'Small Data'



'Small Data'



'Small Data'



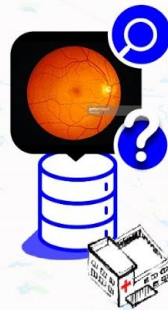
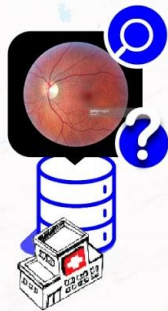
ML for Health Images

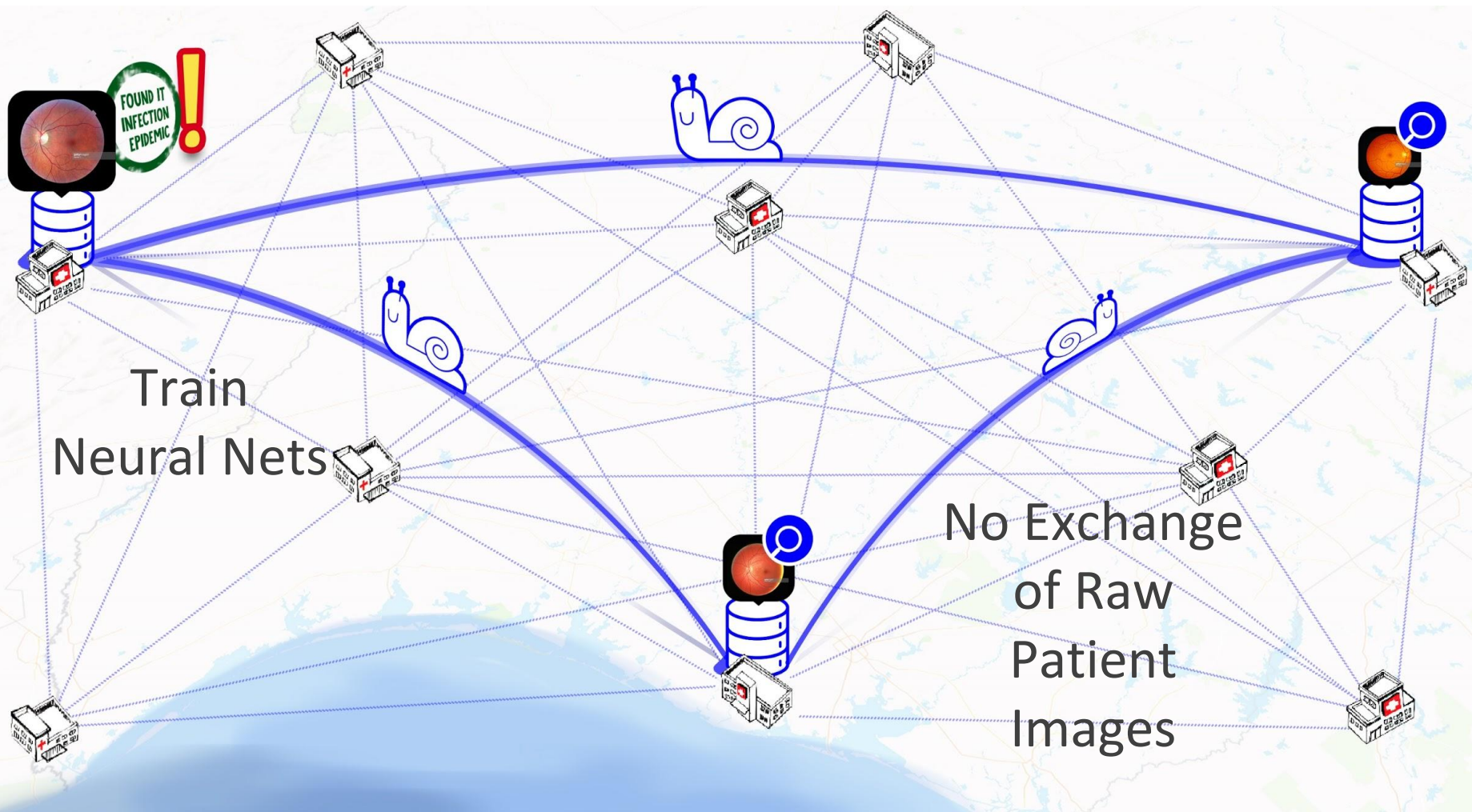
- Distributed Data
- Patient privacy
- Incentives
- ML Expertise
- Efficiency

Low
Bandwidth

Low
Compute

'Small' Data

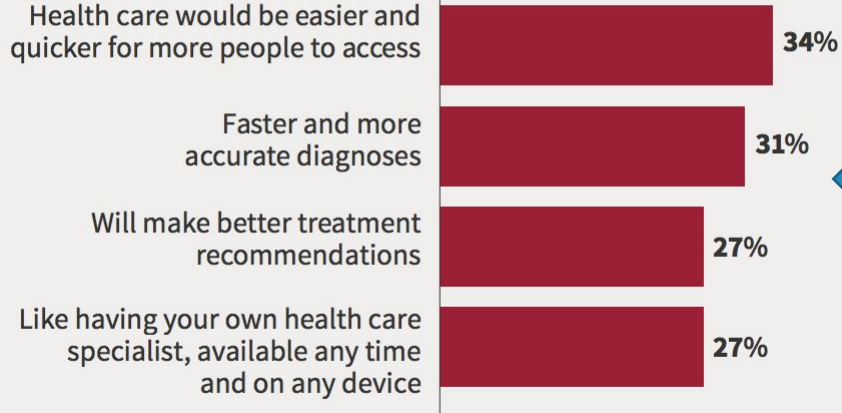




Gupta, Raskar 'Distributed training of deep neural network over several agents', 2017

Intelligent Computing

Top Perceived Advantages of Using AI for Health Care



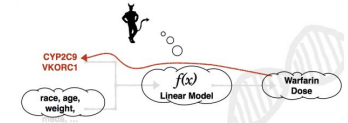
Source: PwC (November 2016). Survey: The new imperatives for health.

Security, Privacy & Safety

... and predictive models can breach privacy too



Privacy in Pharmacogenetics:
An End-to-End Case Study of
Personalized Warfarin Dosing



Regulations

GDPR: General Data Protection Regulation

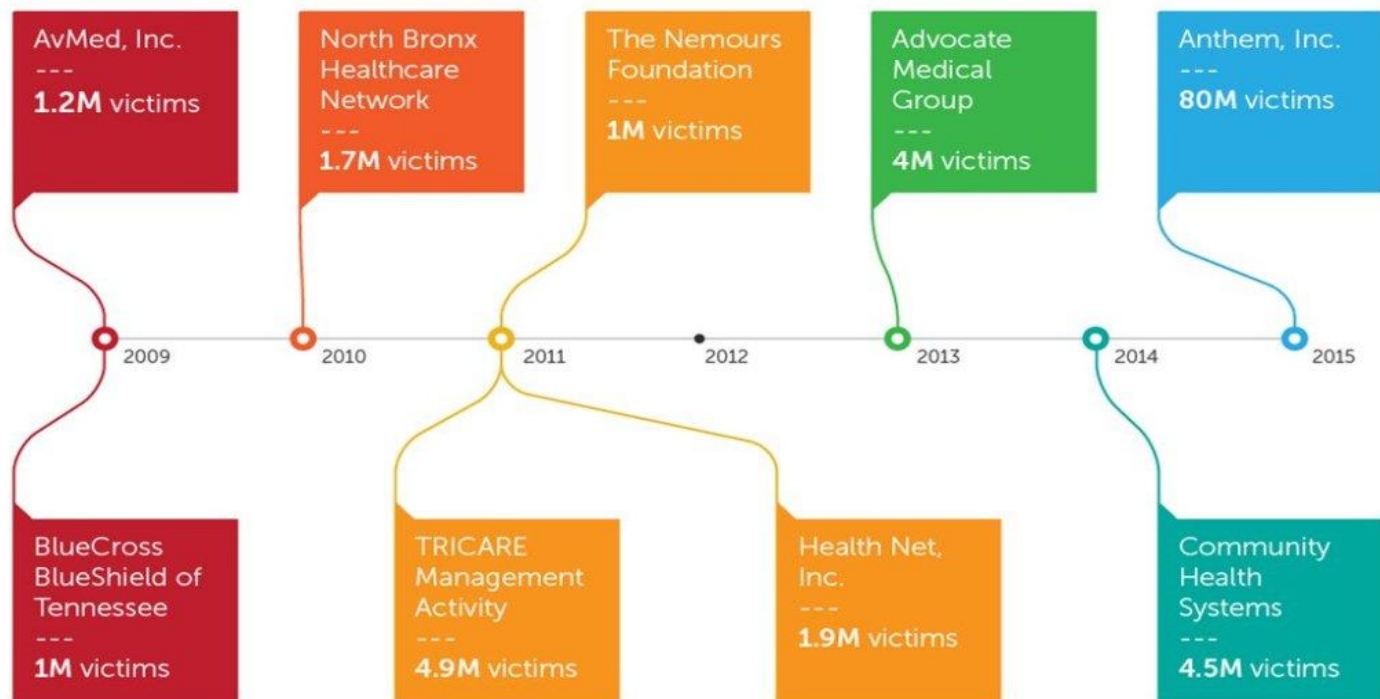
HIPAA: Health Insurance Portability and
Accountability Act, 1996

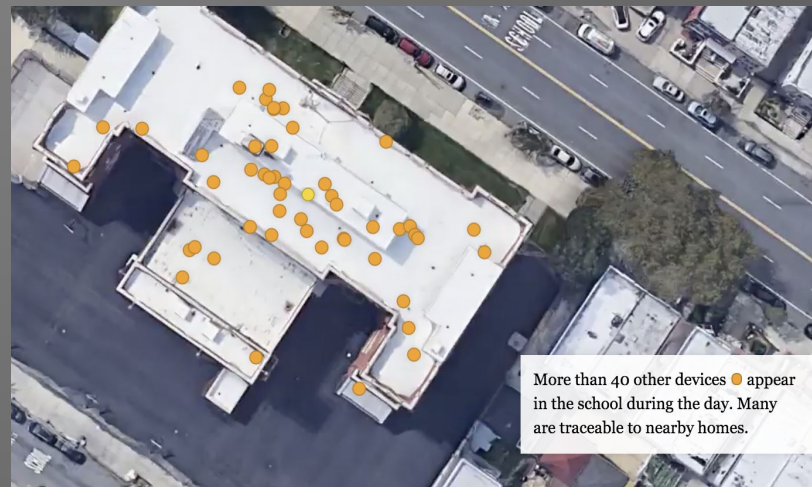
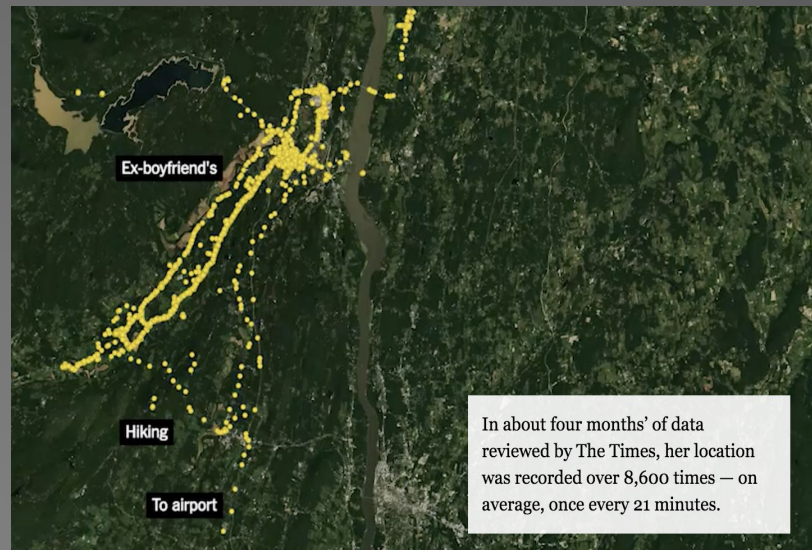
SOX: Sarbanes-Oxley Act, 2002

PCI: Payment Card Industry Data Security
Standard, 2004

SHIELD: Stop Hacks and Improve Electronic Data
Security Act, Jan 1 2019

NOTABLE HEALTHCARE BREACHES





Challenges for Distributed Data + AI + Health

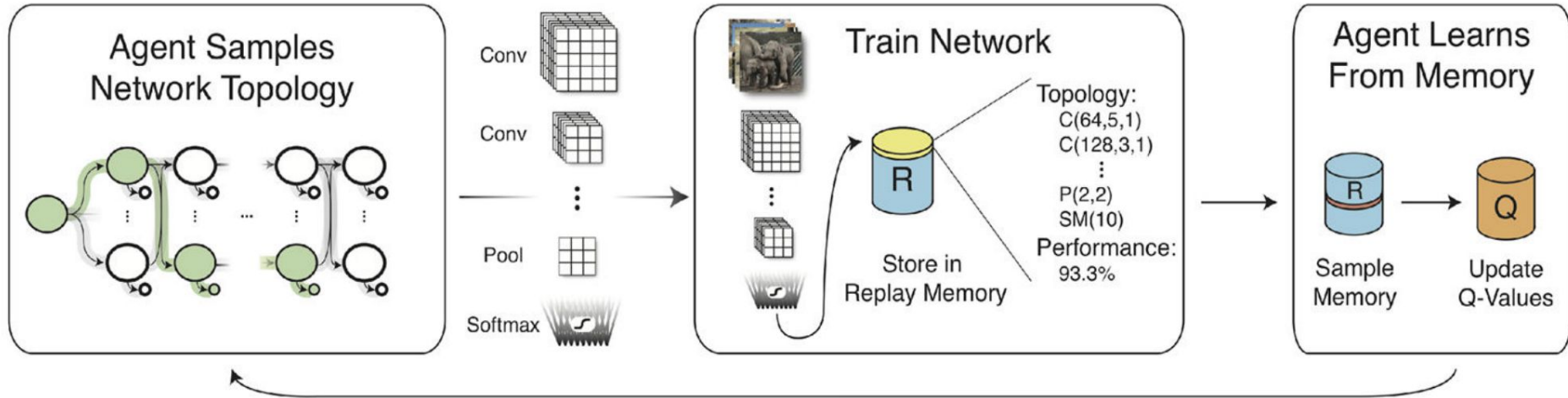
Distributed Data
Multi-Modal
Incomplete Data

Regulations
Incentives
Cooperation
Ease

Ledgering
Smart contracts
Maintenance

Resource-constraints
Memory, Compute, Bandwidth,
Convergence, Synchronization, Leakage

Automating ML : AI building AI



Teacher

Student

AI: Bringing it all together

Training
Deep
Networks

Server

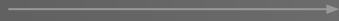
No sharing of
Raw Images

Client

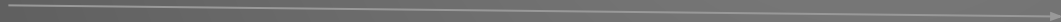
Invisible Data / Data Friction

Overcoming Data Friction

Ease | Incentive | Trust | Regulation



Blockchain



AI/ SplitNN



Anonymize

Obfuscate

Encrypt

Protect Data

Anonymity is not enough ...


A Face Is Exposed for AOL Searcher No. 4417749

By MICHAEL BARBARO and TOM ZELLER Jr.
Published: August 9, 2006

SIGN IN TO E
THIS



Why 'Anonymous' Data Sometimes Isn't

By Bruce Schneier  12.13.07

Last year, Netflix published 10 million movie rankings by 500,000 customers, as part of a challenge for people to come up with better recommendation systems than the one the company was using.

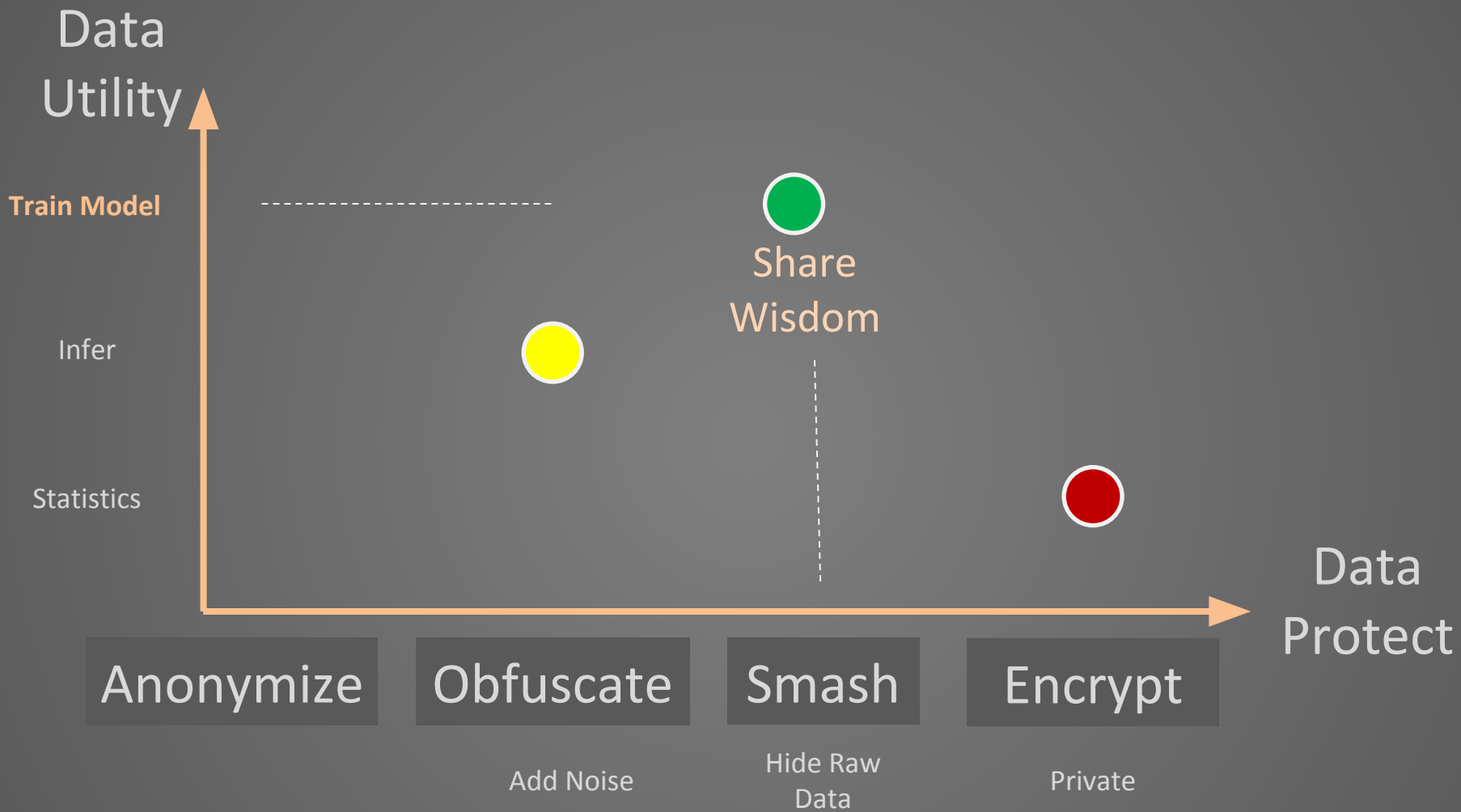
The Scientist • The Nutshell

"Anonymous" Genomes Identified

The names and addresses of people participating in the Personal Genome Project can be easily tracked down despite such data being left off their online profiles.

By Dan Cossins | May 3, 2013





Data
Utility

Train Model

Infer

Statistics

Share
Wisdom

Data
Protect

Anonymize

Obfuscate

Smash

Encrypt

Add Noise

Hide Raw
Data

Private

Federated Learning
Nets trained at Clients
Merged at Server

Differential Privacy
Obfuscate with noise
Hide unique samples

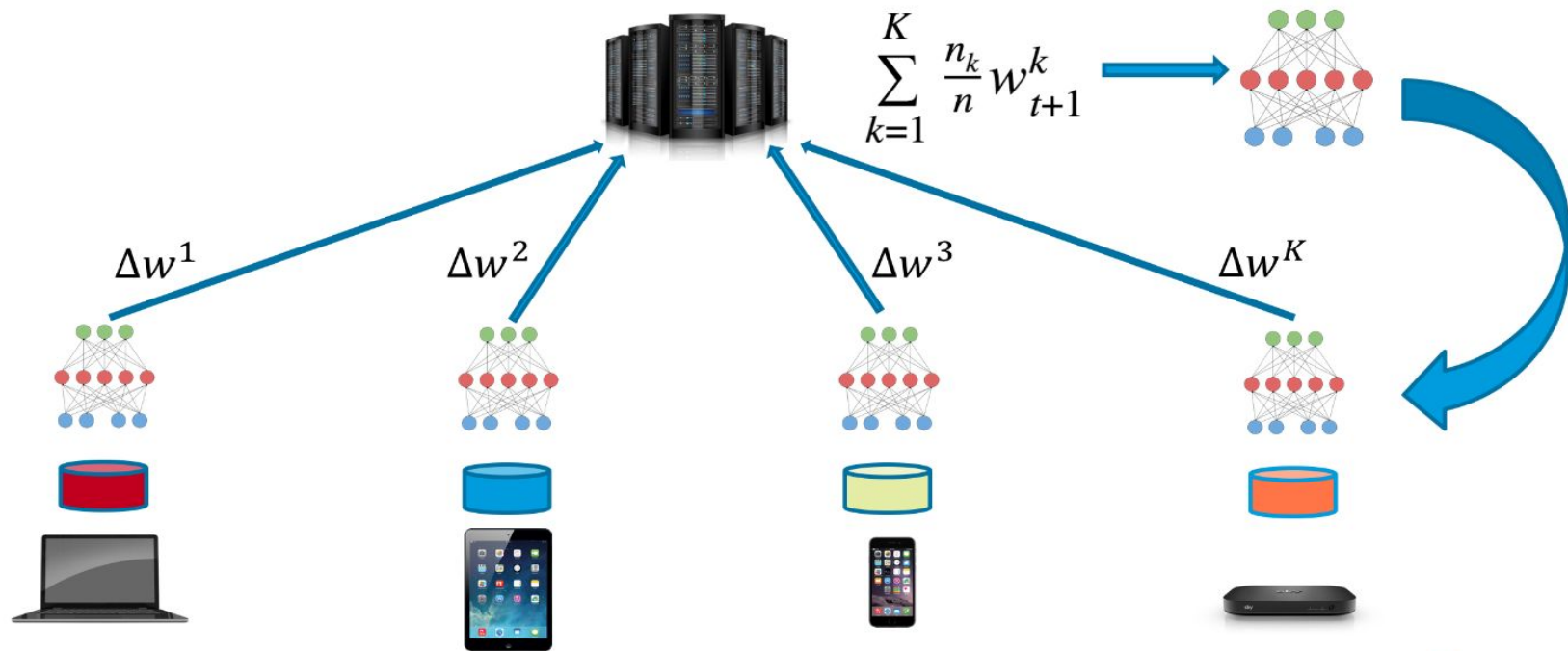
Split Learning (MIT)
Nets split over network
Trained at both

Homomorphic Encryption
Basic Math over Encrypted
Data (+, x)

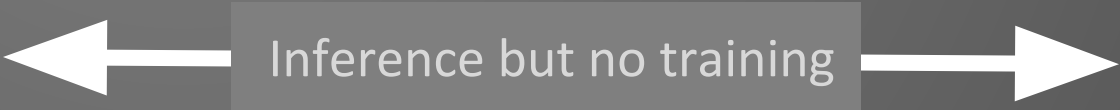
Federated Learning

Server

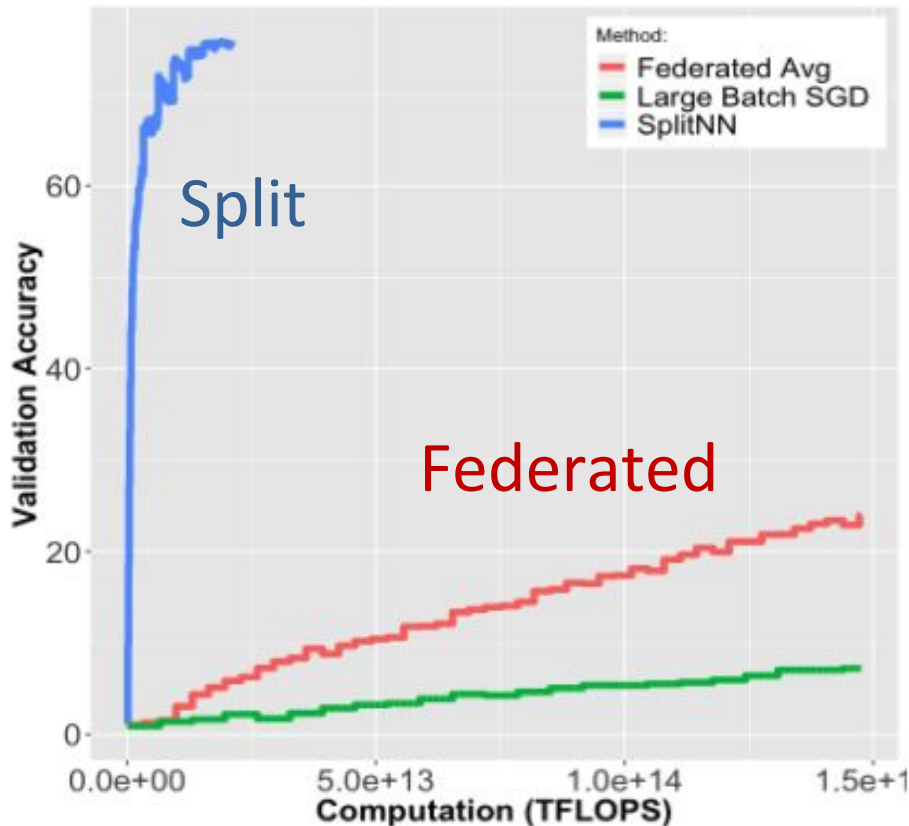
How does it work?



Distributed Training	Protect data	Partial Leakage	Differential Privacy	Homomorphic Encryption	Oblivious Transfer, Garbled Circuits
Federated Learning	●	●	●	●	
Split Learning	●				



When to use split learning?



Large number of clients:

Split learning shows positive results

Memory

Compute

Bandwidth

Convergence

Project Page and Papers:

<https://splitlearning.github.io/>

QUANTITATIVE RESULTS

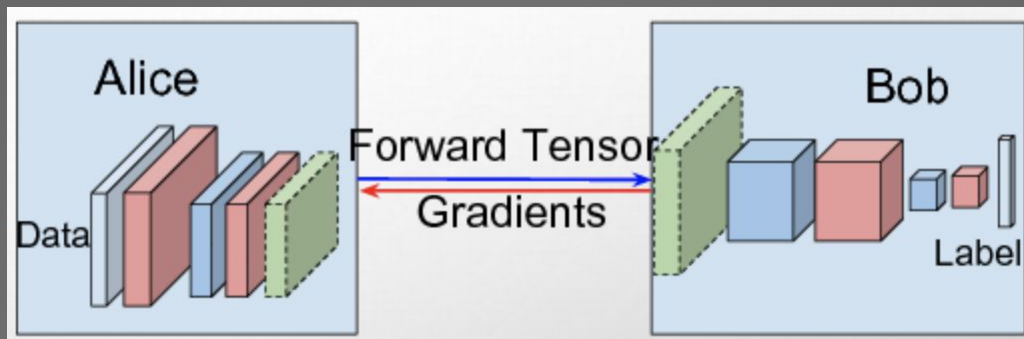
Method	100 Clients	500 Clients
Large Scale SGD	29.4 TFlops	5.89 TFlops
Federated Learning	29.4 TFlops	5.89 TFlops
Our Method (SplitNN)	0.1548 TFlops	0.03 TFlops

Table 1. Computation resources consumed per client when training CIFAR 10 over VGG (in teraflops)

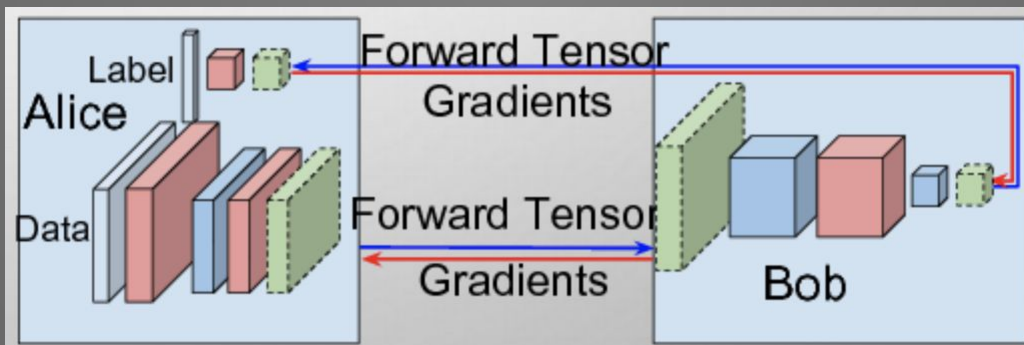
Method	100 Clients	500 Clients
Large Scale SGD	13 GB	14 GB
Federated Learning	3 GB	2.4 GB
Our Method (SplitNN)	6 GB	1.2 GB

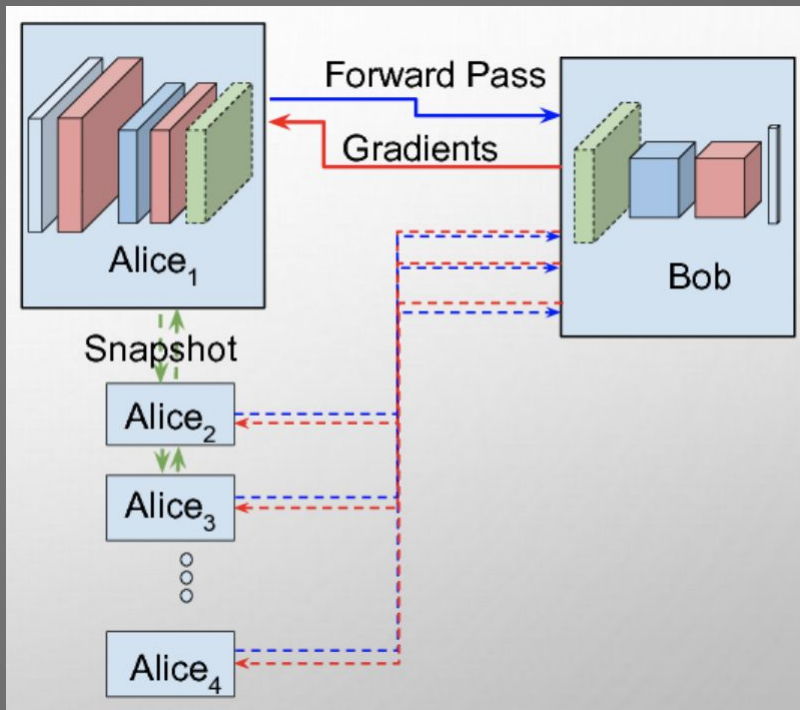
Table 2. Communication Bandwidth consumed per client when training CIFAR 100 and Resnet 50 (in gigabytes)

Label
Sharing



No Label
Sharing



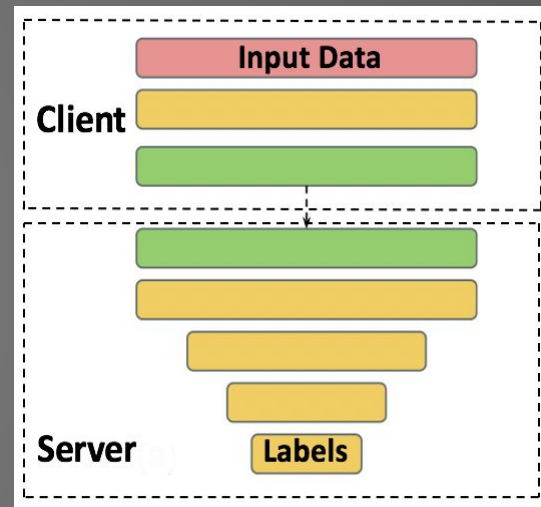
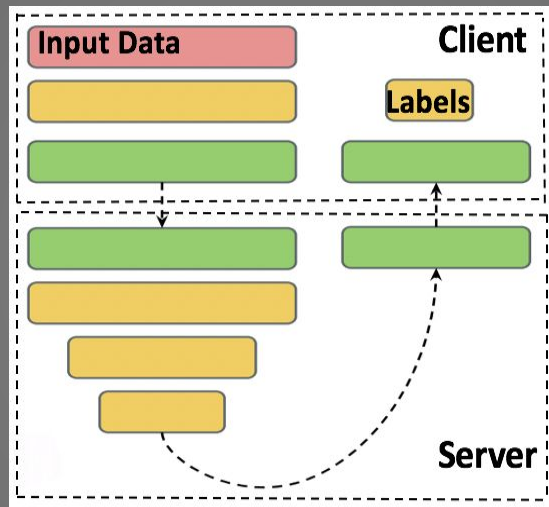
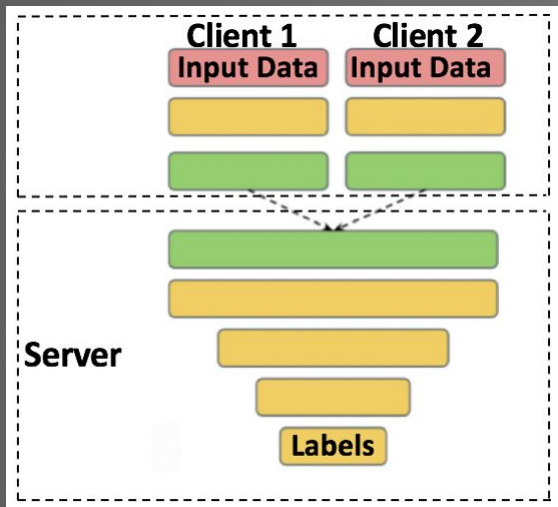


Gupta, Otkrist, and Raskar, Ramesh. "Secure Training of Multi-Party Deep Neural Network." U.S. Patent Application No. 15/630,944.

Distribution of parameters in AlexNet

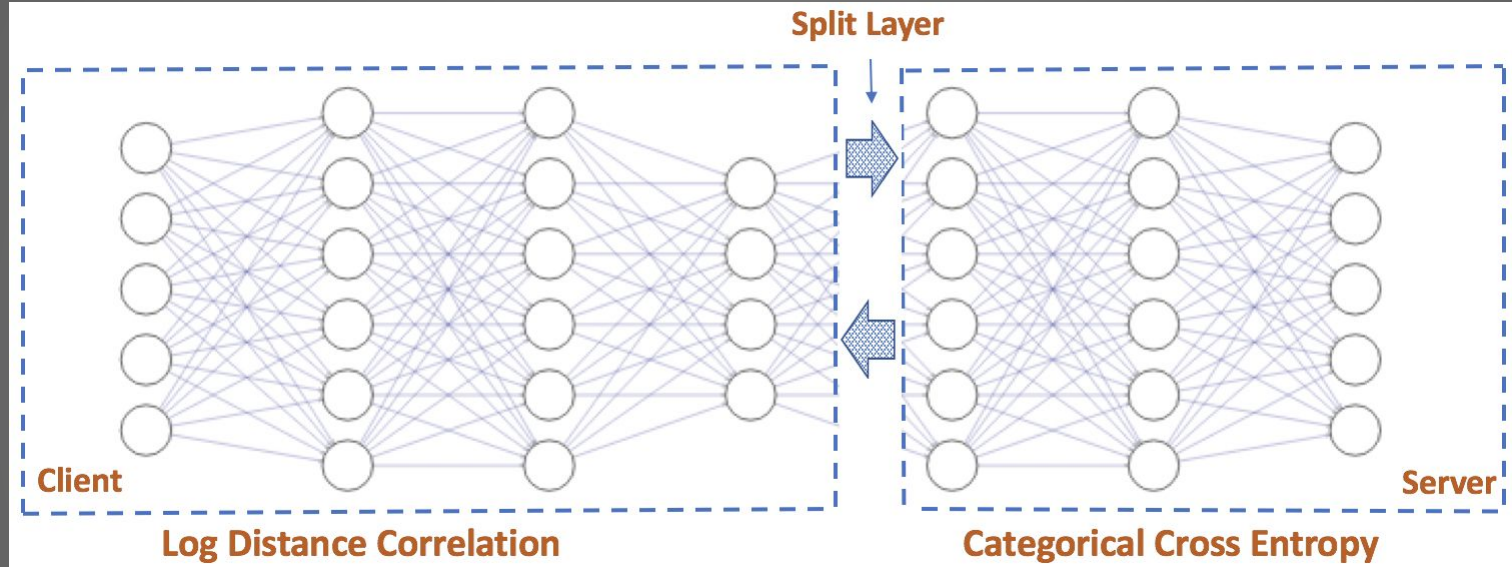
Layer Name	Tensor Size	Weights	Biases	Parameters
Input Image	227x227x3	0	0	0
Conv-1	55x55x96	34,848	96	34,944
MaxPool-1	27x27x96	0	0	0
Conv-2	27x27x256	614,400	256	614,656
MaxPool-2	13x13x256	0	0	0
Conv-3	13x13x384	884,736	384	885,120
Conv-4	13x13x384	1,327,104	384	1,327,488
Conv-5	13x13x256	884,736	256	884,992
MaxPool-3	6x6x256	0	0	0
FC-1	4096x1	37,748,736	4,096	37,752,832
FC-2	4096x1	16,777,216	4,096	16,781,312
FC-3	1000x1	4,096,000	1,000	4,097,000
Output	1000x1	0	0	0
Total				62,378,344

Versatile Configurations of Split Learning



Split learning for health: Distributed deep learning without sharing raw patient data, Praneeth Vepakomma, Otkrist Gupta, Tristan Swedish, Ramesh Raskar, (2019)

NoPeek SplitNN: Reducing Leakage in Distributed Deep Learning



$$\alpha_1 DCOR(\mathbf{X}_n, \hat{\mathbf{Z}}) + \alpha_2 CCE(\hat{\mathbf{Y}}, \mathbf{Y}_n)$$

Reducing leakage in distributed deep learning for sensitive health data, Praneeth Vepakomma, Otkrist Gupta, Abhimanyu Dubey, Ramesh Raskar (2019)

No peak deep learning with conditioning variable

Setup:

- **Supervised:** $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\} \subset X \times Y$
- Output: $y \in \mathbb{R}$
- *Goal:* To find a projection $\mathcal{S}_{Y|X}$ such that, $Y \perp\!\!\!\perp X|Z$.

Ideal Goal: To find such a conditioning variable Z within the framework of deep learning such that the following directions are approximately satisfied:

1. $Y \perp\!\!\!\perp X | Z$ (Utility property as X can be thrown away given Z to obtain prediction $E(Y|Z)$)
2. $X \perp\!\!\!\perp Z$ (One-way property preventing proper reconstruction of raw data X from Z)

Note: $\perp\!\!\!\perp$ denotes statistical independence

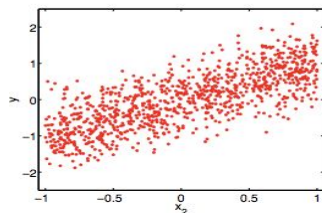
Possible measures of non-linear dependence

- COCO: Constrained Covariance
- HSIC: Hilbert-Schmidt Independence Criterion
- DCOR: Distance Correlation
- MMD: Maximum Mean Discrepancy
- KTA: Kernel Target Alignment
- MIC: Maximal Information Coefficient
- TIC: Total Information Coefficient

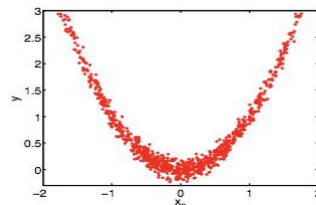
Why is it called distance correlation?

Definition 3.1. Sample Distance Covariance [3]: Given i.i.d samples $\mathcal{X} \times \mathcal{Y} = \{(\mathbf{x}_k, \mathbf{y}_k) | k = 1, 2, 3, \dots, n\}$ and corresponding double centered Euclidean distance matrices $\hat{\mathbf{E}}_{\mathbf{X}}$ and $\hat{\mathbf{E}}_{\mathbf{Y}}$, then the squared sample distance correlation is defined as,

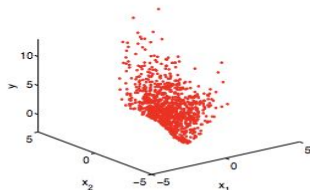
$$\hat{\nu}^2(\mathbf{X}, \mathbf{Y}) = \frac{1}{n^2} \sum_{k,l=1}^n [\hat{\mathbf{E}}_{\mathbf{X}}]_{k,l} [\hat{\mathbf{E}}_{\mathbf{Y}}]_{k,l},$$



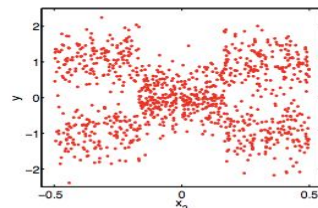
(a) linear



(b) non-linear



(c) non-linear



(d) non-linear

Distance Covariance (Székely, G. (2007))

$$\nu^2(\mathbf{X}, \mathbf{Y}; w) = \int_{\mathbb{R}^{h+m}} |f_{\mathbf{X}, \mathbf{Y}}(t, s) - f_{\mathbf{X}}(t)f_{\mathbf{Y}}(s)|^2 w(t, s) dt ds$$

where $f_{\mathbf{X}}$, $f_{\mathbf{Y}}$, $f_{\mathbf{X}, \mathbf{Y}}$ are the characteristic functions of \mathbf{X} , \mathbf{Y} , $\mathbf{X} \times \mathbf{Y}$ and $w(t, s)$ is a suitably chosen weight function.

Sample Distance Covariance (2nd order)

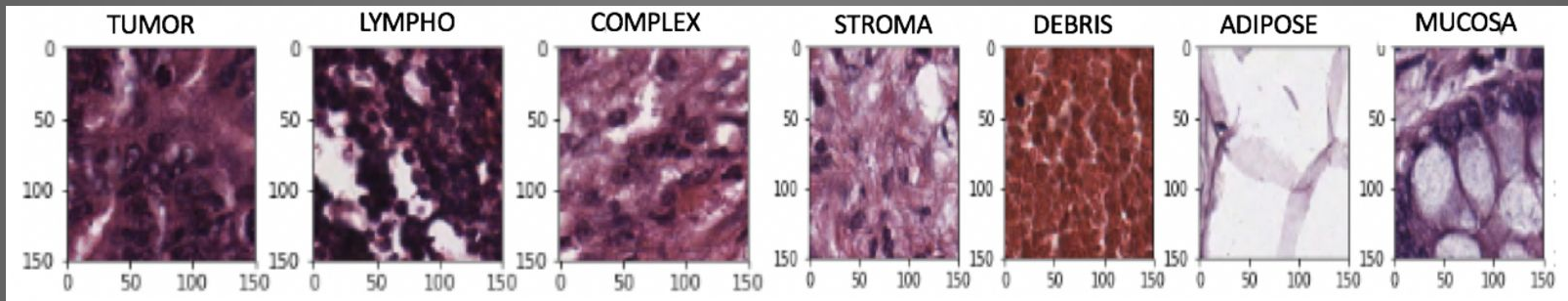
$$\hat{\nu}^2(\mathbf{X}, \mathbf{Y}) = \frac{1}{n^2} \text{Tr}(\mathbf{L}_x^T \mathbf{L}_y)$$

where $\mathbf{L}_x = \mathbf{D}_x - \mathbf{H}\mathbf{E}_x\mathbf{H}$ and $\mathbf{L}_y = \mathbf{D}_y - \mathbf{H}\mathbf{E}_y\mathbf{H}$.

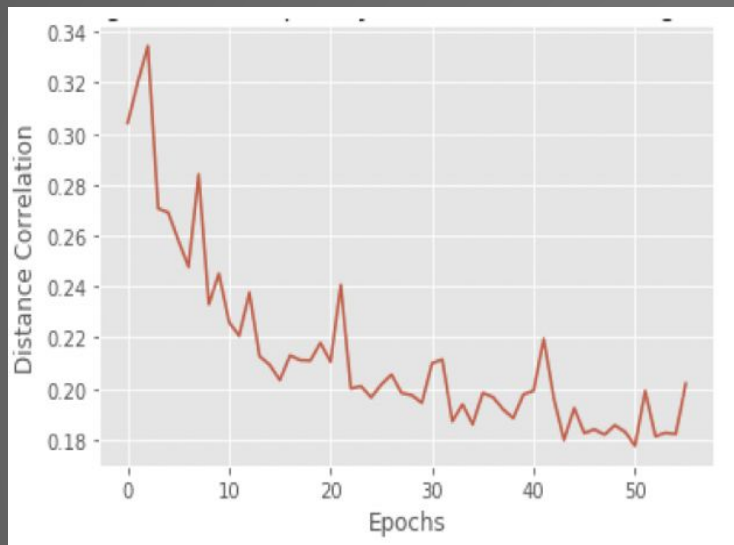
Lemma 3.1. *Given matrices of squared Euclidean distances \mathbf{E}_X and \mathbf{E}_Y and Laplacians \mathbf{L}_X and \mathbf{L}_Y formed over adjacency matrices $\hat{\mathbf{E}}_X$ and $\hat{\mathbf{E}}_Y$, the square of sample distance correlation $\hat{\rho}^2(\mathbf{X}, \mathbf{Y})$ is given by*

$$\hat{\rho}^2(\mathbf{X}, \mathbf{Y}) = \frac{\text{Tr}(\mathbf{X}^T \mathbf{L}_Y \mathbf{X})}{\sqrt{\text{Tr}(\mathbf{Y}^T \mathbf{L}_Y \mathbf{Y}) \text{Tr}(\mathbf{X}^T \mathbf{L}_X \mathbf{X})}}.$$

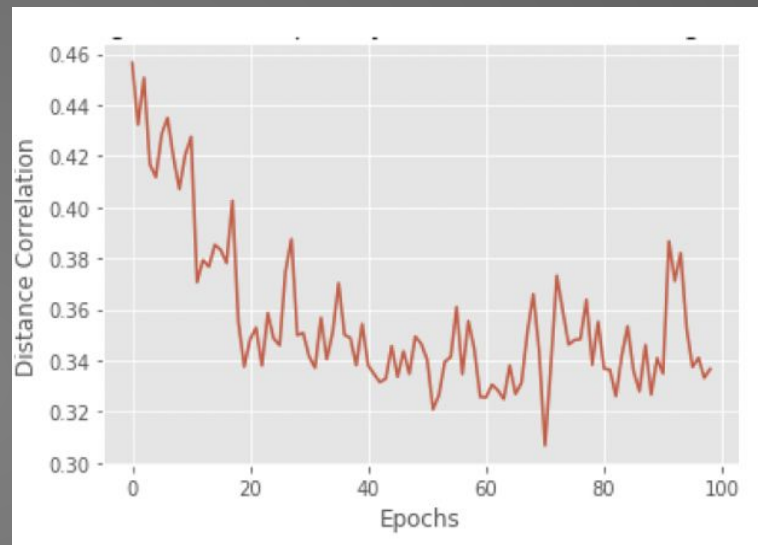
Colorectal histology image dataset (Public data)



Leakage Reduction in Action



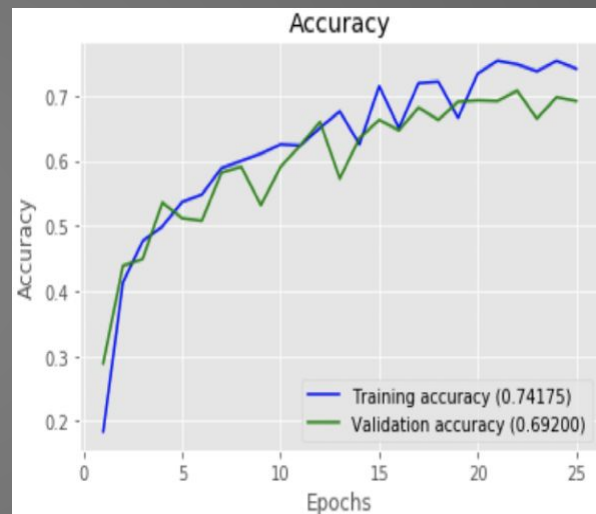
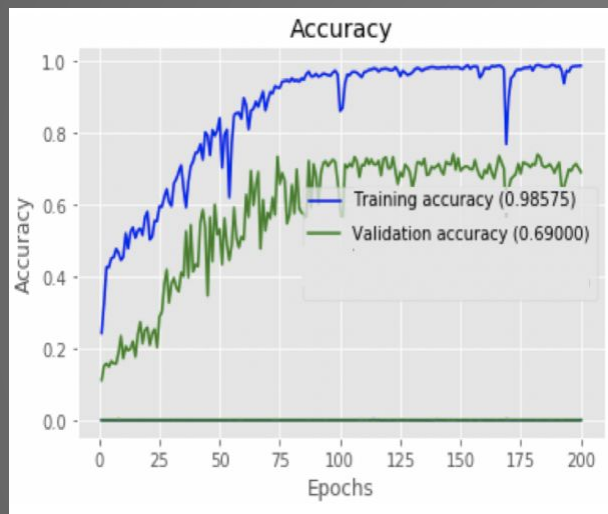
Reduced leakage during training over colorectal histology image data from 0.96 in traditional CNN to 0.19 in NoPeek SplitNN



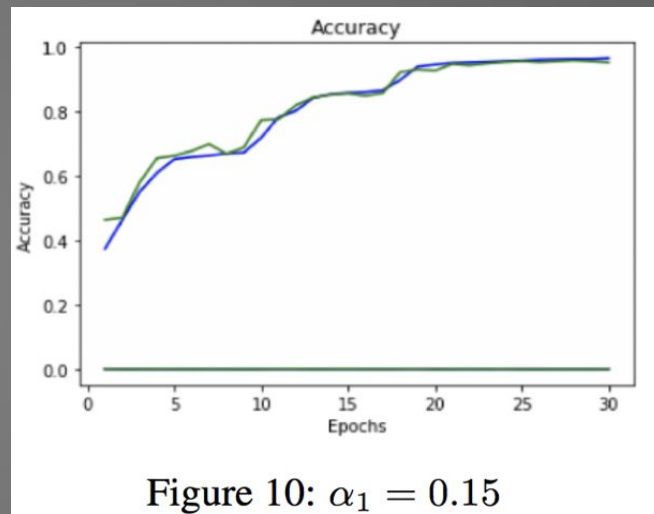
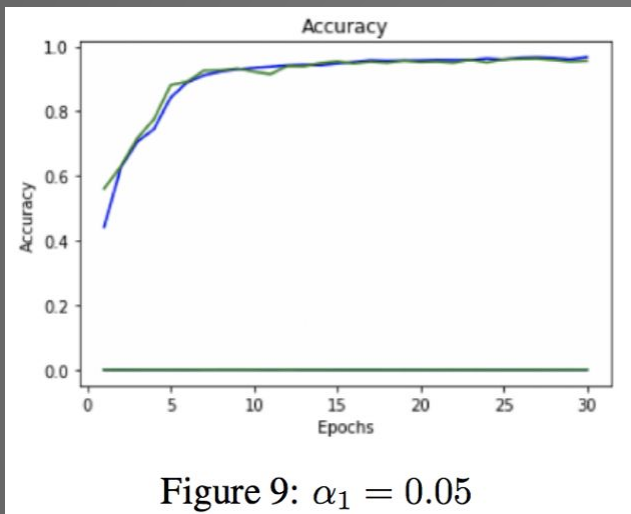
Reduced leakage during training over colorectal histology image data from 0.92 in traditional CNN to 0.33 in NoPeek SplitNN

Reducing leakage in distributed deep learning for sensitive health data, Praneeth Vepakomma, Otkrist Gupta, Abhimanyu Dubey, Ramesh Raskar (2019)

Similar validation performance



Effect of leakage reduction on convergence



Robustness to reconstruction

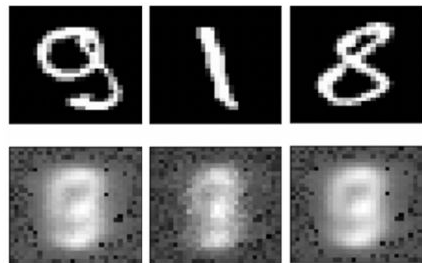


Figure 7: $\alpha_1 = 0.1$

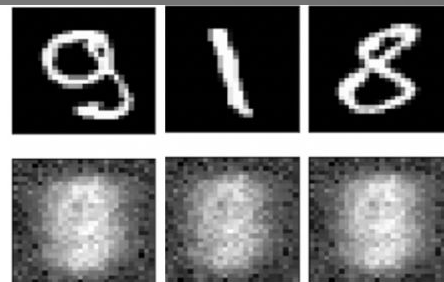


Figure 8: $\alpha_1 = 0.9$

Proof of one-Way Property:

$$DCOV(\mathbf{X}, \mathbf{Z}) = \text{Tr}(\mathbf{X}\mathbf{X}^T\mathbf{Z}\mathbf{Z}^T) + \|\mathbf{X} - \mathbf{Z}\| + \|\mathbf{Z}\|$$

$$D_{KL}(\mathbf{Z}||\mathbf{X}) - D_{KL}(\mathbf{X}||\mathbf{Z}) = H(\mathbf{Z}, \mathbf{X}) - H(\mathbf{Z}) - H(\mathbf{X}, \mathbf{Z}) + H(\mathbf{X})$$

We show: Minimizing regularized distance covariance minimizes the difference of Kullback-Leibler divergences

$$= \det(\mathbf{Z}^T \mathbf{X}) - \det(\mathbf{Z}^T \mathbf{Z}) - \det(\mathbf{X}^T \mathbf{Z}) + \det(\mathbf{X}^T \mathbf{X})$$

This can be bounded using Hadamard's inequality as

$$\begin{aligned} \det(\mathbf{Z}^T \mathbf{X}) - \det(\mathbf{Z}^T \mathbf{Z}) + \det(\mathbf{X}^T \mathbf{X}) - \det(\mathbf{X}^T \mathbf{Z}) &\leq \|\mathbf{Z}^T \mathbf{X} - \mathbf{Z}^T \mathbf{Z}\|_2 \frac{\|\mathbf{Z}^T \mathbf{X}\|_2^n - \|\mathbf{Z}^T \mathbf{Z}\|_2^n}{\|\mathbf{Z}^T \mathbf{X}\|_2 - \|\mathbf{Z}^T \mathbf{Z}\|_2} \\ &\quad + \|\mathbf{X}^T \mathbf{Z} - \mathbf{X}^T \mathbf{X}\|_2 \frac{\|\mathbf{X}^T \mathbf{Z}\|_2^n - \|\mathbf{X}^T \mathbf{X}\|_2^n}{\|\mathbf{X}^T \mathbf{Z}\|_2 - \|\mathbf{X}^T \mathbf{X}\|_2} \end{aligned}$$


The fractional terms $\frac{\|\mathbf{Z}^T \mathbf{X}\|_2^n - \|\mathbf{Z}^T \mathbf{Z}\|_2^n}{\|\mathbf{Z}^T \mathbf{X}\|_2 - \|\mathbf{Z}^T \mathbf{Z}\|_2}$, $\frac{\|\mathbf{X}^T \mathbf{Z}\|_2^n - \|\mathbf{X}^T \mathbf{X}\|_2^n}{\|\mathbf{X}^T \mathbf{Z}\|_2 - \|\mathbf{X}^T \mathbf{X}\|_2}$ can be written as a sum of geometric-series, with factors of change of $\frac{\|\mathbf{Z}^T \mathbf{X}\|_2}{\|\mathbf{Z}^T \mathbf{Z}\|_2}$, $\frac{\|\mathbf{X}^T \mathbf{Z}\|_2}{\|\mathbf{X}^T \mathbf{X}\|_2}$ respectively because

$$\frac{\|\mathbf{Z}^T \mathbf{X}\|_2^n - \|\mathbf{Z}^T \mathbf{Z}\|_2^n}{\|\mathbf{Z}^T \mathbf{X}\|_2 - \|\mathbf{Z}^T \mathbf{Z}\|_2} = \frac{1 - \left(\frac{\|\mathbf{Z}^T \mathbf{X}\|_2}{\|\mathbf{Z}^T \mathbf{Z}\|_2}\right)^n}{1 - \frac{\|\mathbf{Z}^T \mathbf{X}\|_2}{\|\mathbf{Z}^T \mathbf{Z}\|_2}} = \sum_{p=0}^{n-1} \|\mathbf{Z}^T \mathbf{X}\|_2^p \|\mathbf{Z}^T \mathbf{Z}\|_2^{p-1}$$

Therefore these fractional terms can be minimized by minimizing $\|\mathbf{Z}^T \mathbf{X}\|_2$ and $\|\mathbf{Z}^T \mathbf{Z}\|_2$ as the sums of products of decreasing functions of norms are also decreasing. By Cauchy-Schwarz inequality $\|\mathbf{Z}^T (\mathbf{X} - \mathbf{Z})\| \leq \|\mathbf{Z}\| \|\mathbf{X} - \mathbf{Z}\|$.

Therefore the upper-bound on difference of KL-divergence can be minimized by minimizing $\|\mathbf{Z}\|$ and $\|\mathbf{X} - \mathbf{Z}\|$ to minimize terms $\|\mathbf{Z}^T \mathbf{X} - \mathbf{Z}^T \mathbf{Z}\|$, $\|\mathbf{X}^T \mathbf{Z} - \mathbf{X}^T \mathbf{X}\|$ in addition to minimizing

$$\begin{aligned} \|\mathbf{Z}^T \mathbf{Z}\|, \|\mathbf{Z}^T \mathbf{X}\|_2 &= \text{Tr}(\mathbf{Z}^T \mathbf{X} \mathbf{X}^T \mathbf{Z}) = \text{DCOV}(\mathbf{X}, \mathbf{Z}) \text{ to minimize terms } \frac{\|\mathbf{Z}^T \mathbf{X}\|_2^n - \|\mathbf{Z}^T \mathbf{Z}\|_2^n}{\|\mathbf{Z}^T \mathbf{X}\|_2 - \|\mathbf{Z}^T \mathbf{Z}\|_2}, \\ \frac{\|\mathbf{X}^T \mathbf{Z}\|_2^n - \|\mathbf{X}^T \mathbf{X}\|_2^n}{\|\mathbf{X}^T \mathbf{Z}\|_2 - \|\mathbf{X}^T \mathbf{X}\|_2}. \end{aligned}$$



Distributed Private Machine Learning for Computer Vision: Federated Learning, Split Learning and Beyond

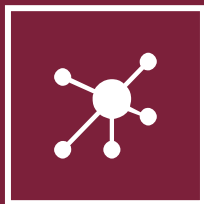
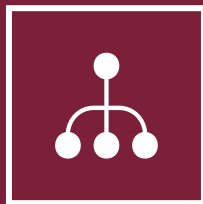
Brendan McMahan, Jakub Konečný, Ramesh Raskar, Otkrist Gupta, Hassan Takabi, Praneeth Vepakomma

Project Page and Papers:
<https://splitlearning.github.io/>

Thanks and acknowledgements to: Otkrist Gupta (MIT/LendBuzz), Ramesh Raskar (MIT), Jayashree Kalpathy-Cramer (Martinos/Harvard), Rajiv Gupta (MGH), Brendan McMahan (Google), Jakub Konečný (Google), Abhimanyu Dubey (MIT), Tristan Swedish (MIT), Sai Sri Sathya (S20.ai), Vitor Pamplona (MIT/EyeNetra), Rodmy Paredes Alfaro (MIT), Kevin Pho (MIT), Elsa Itambo (MIT)



Massachusetts
Institute of
Technology



THANK
YOU
