

Two small ideas: Approximate LCA, Computational Ethics



MLC Research Jam
28 May 2025

Jason Yosinski
`jason@yosinski.com`

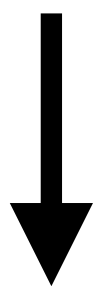
Idea #1: Approximate LCA (Loss Change Allocation)



See training in progress!



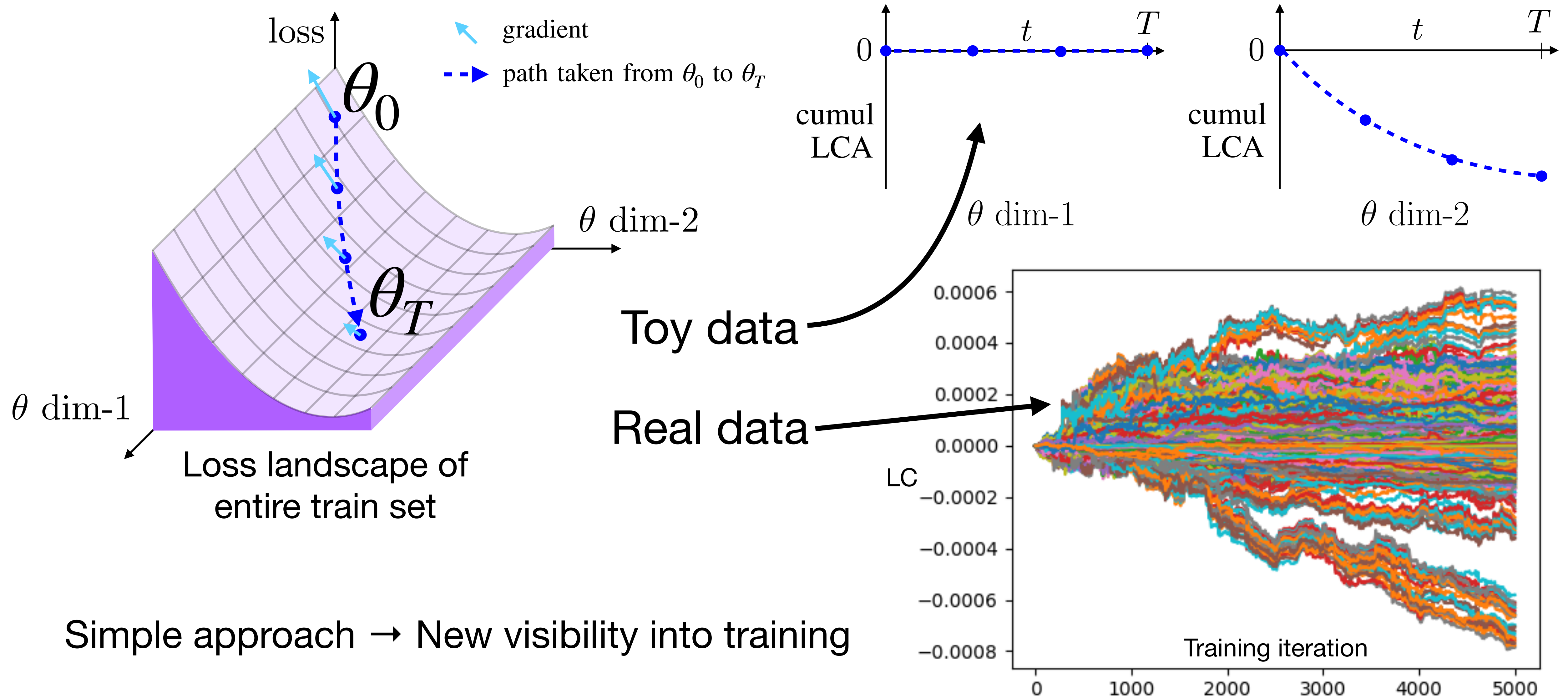
Very slow



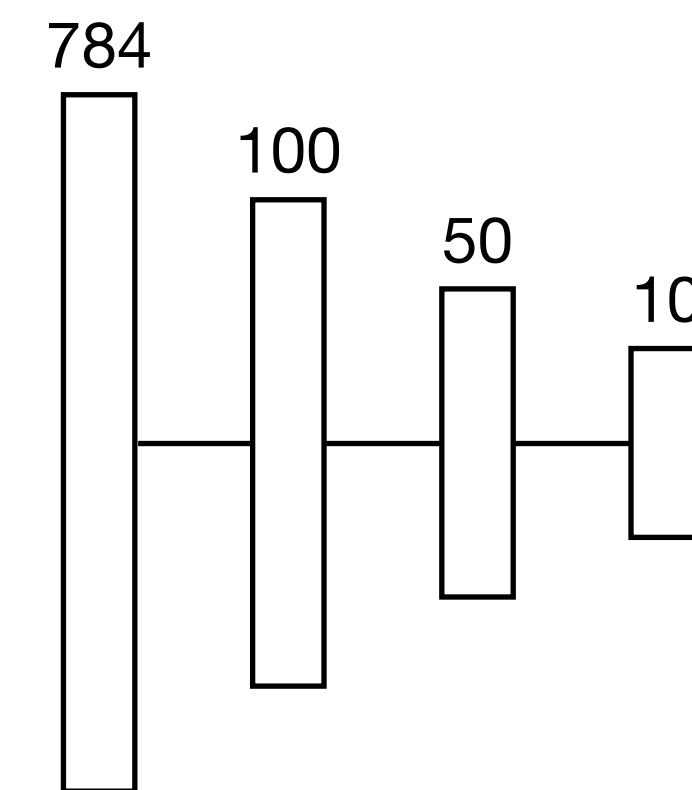
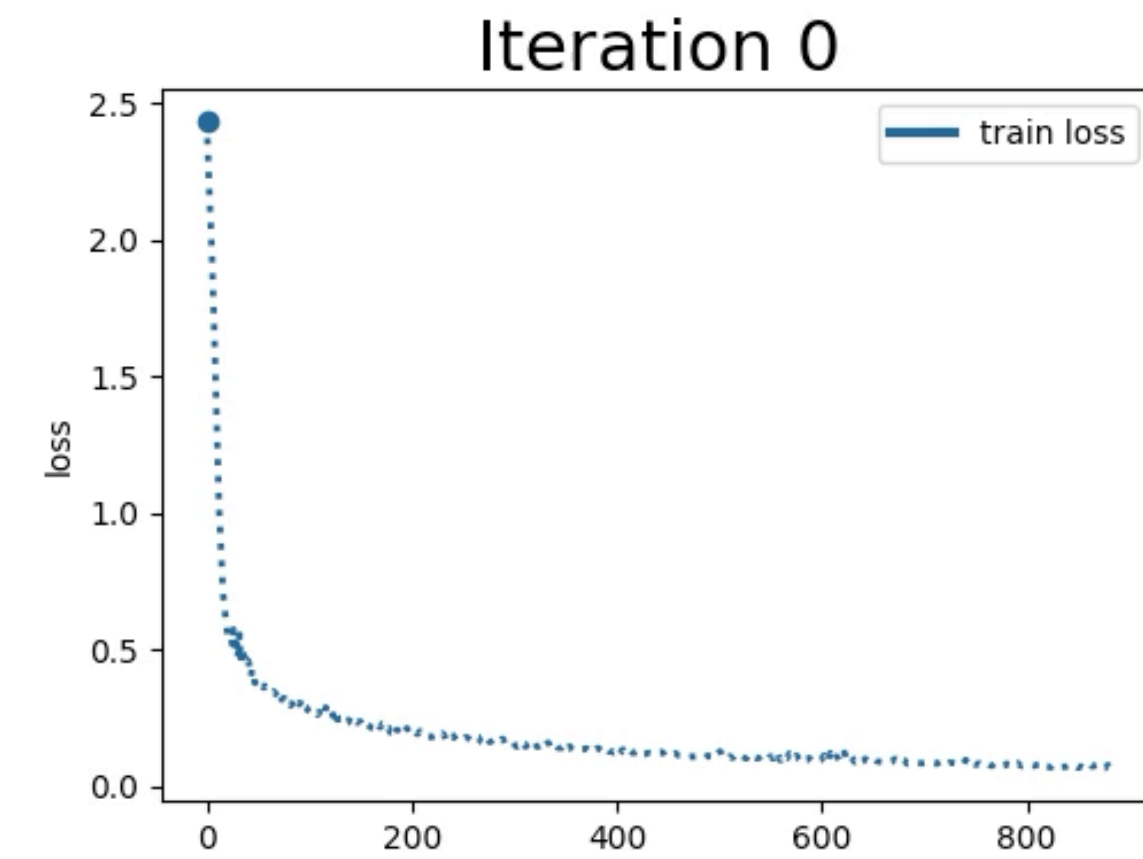
If we use aggressive approximation to make it faster, can it still provide useful visibility?

Janice Lan, Hattie Zhou, Rosanne Liu, Jason Yosinski. LCA: Loss Change Allocation for Neural Network Training, NeurIPS 2019.

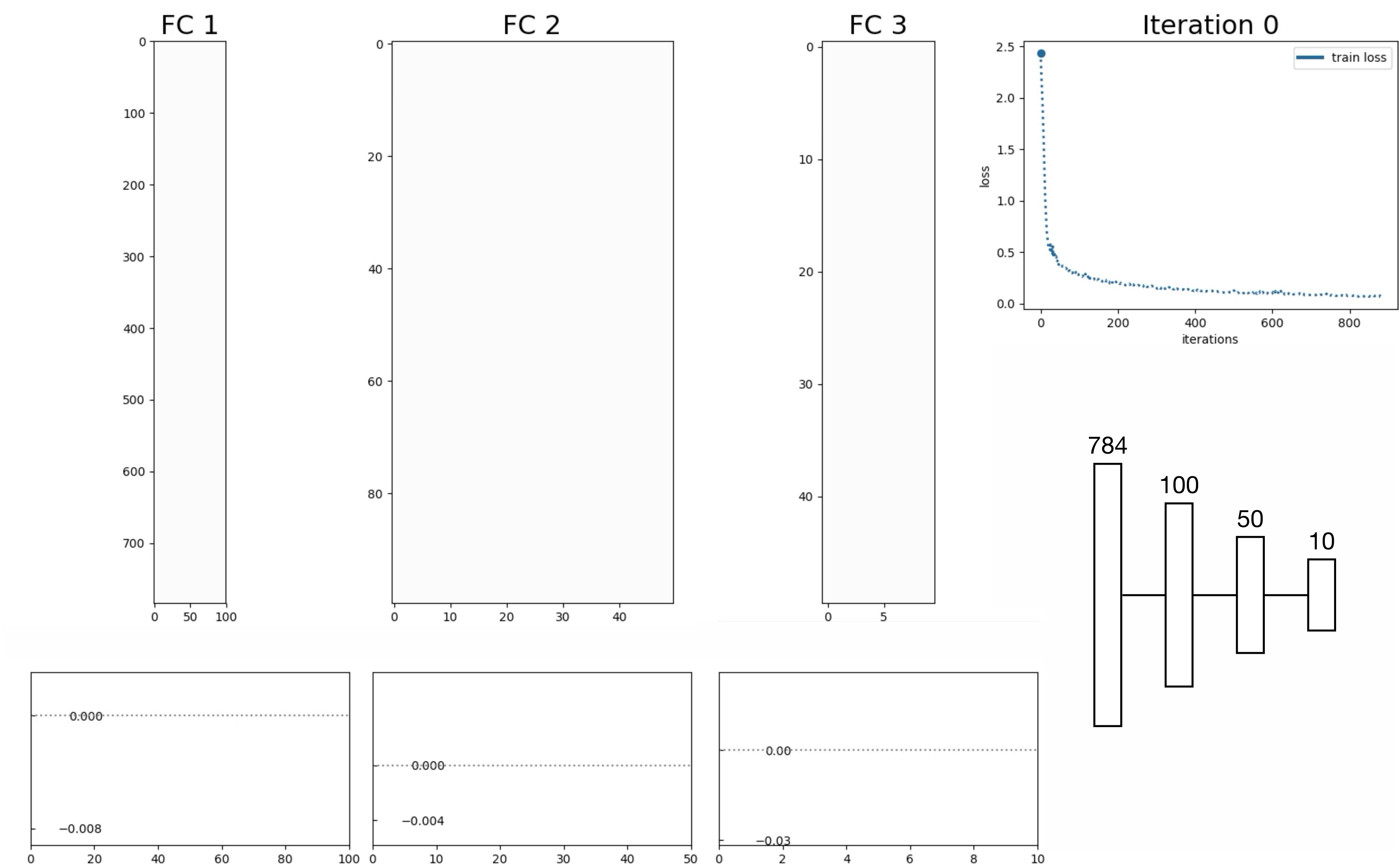
LCA: Loss Change Allocation for Neural Network Training



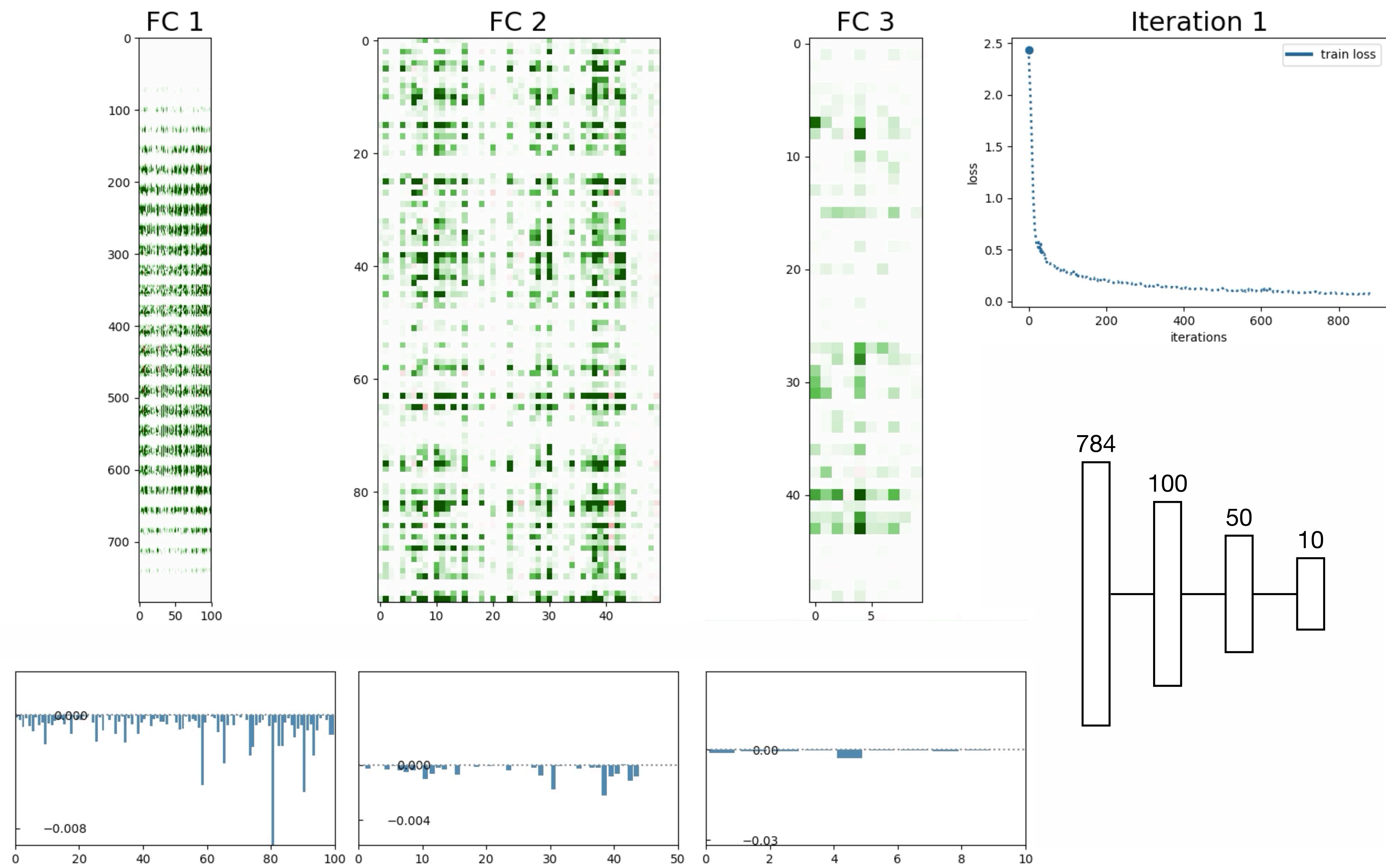
LCA: Loss Change Allocation for Neural Network Training



LCA: Loss Change Allocation for Neural Network Training



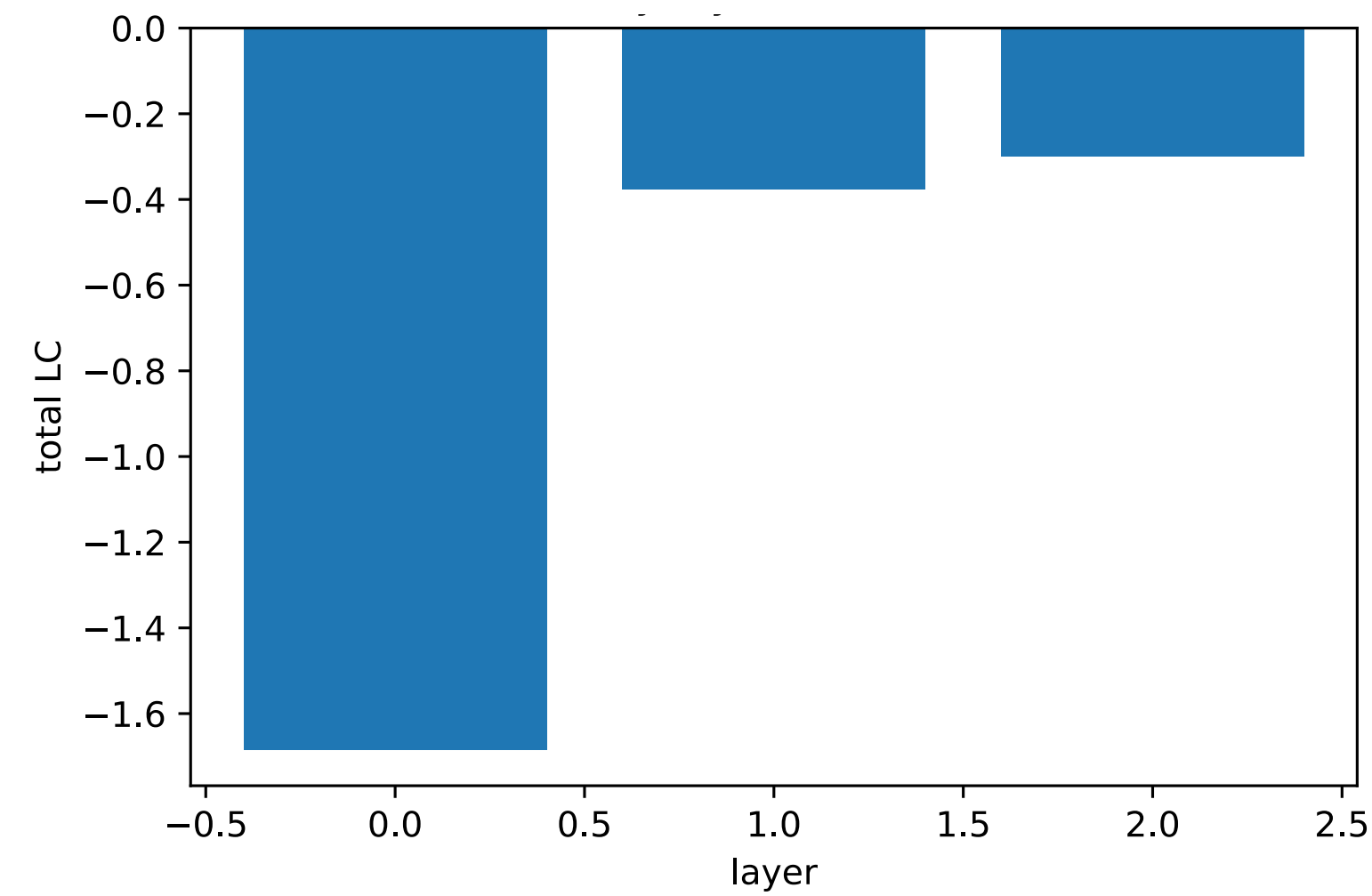
LCA: Loss Change Allocation for Neural Network Training



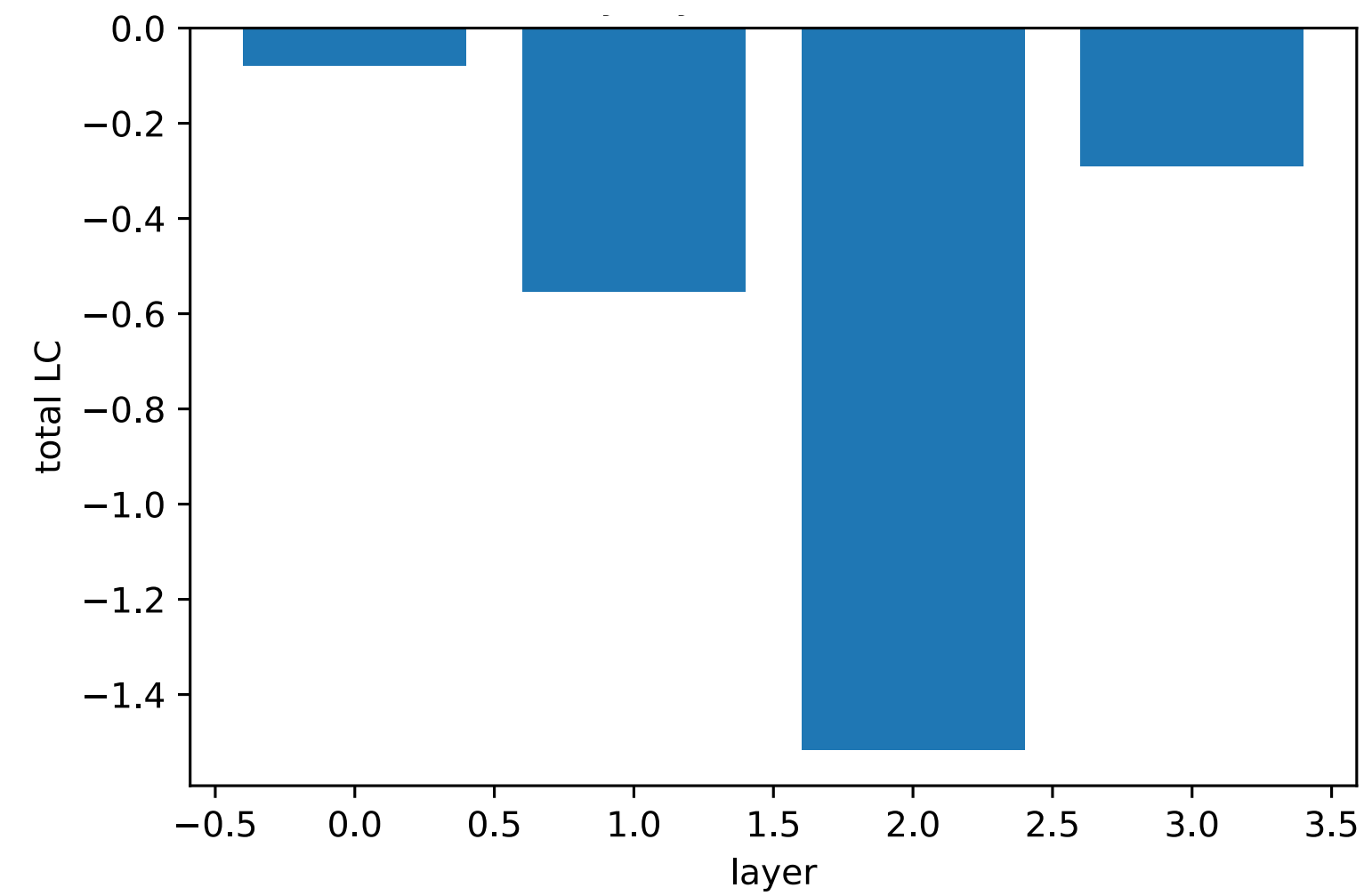
LCA: Loss Change Allocation for Neural Network Training

Sum over layers

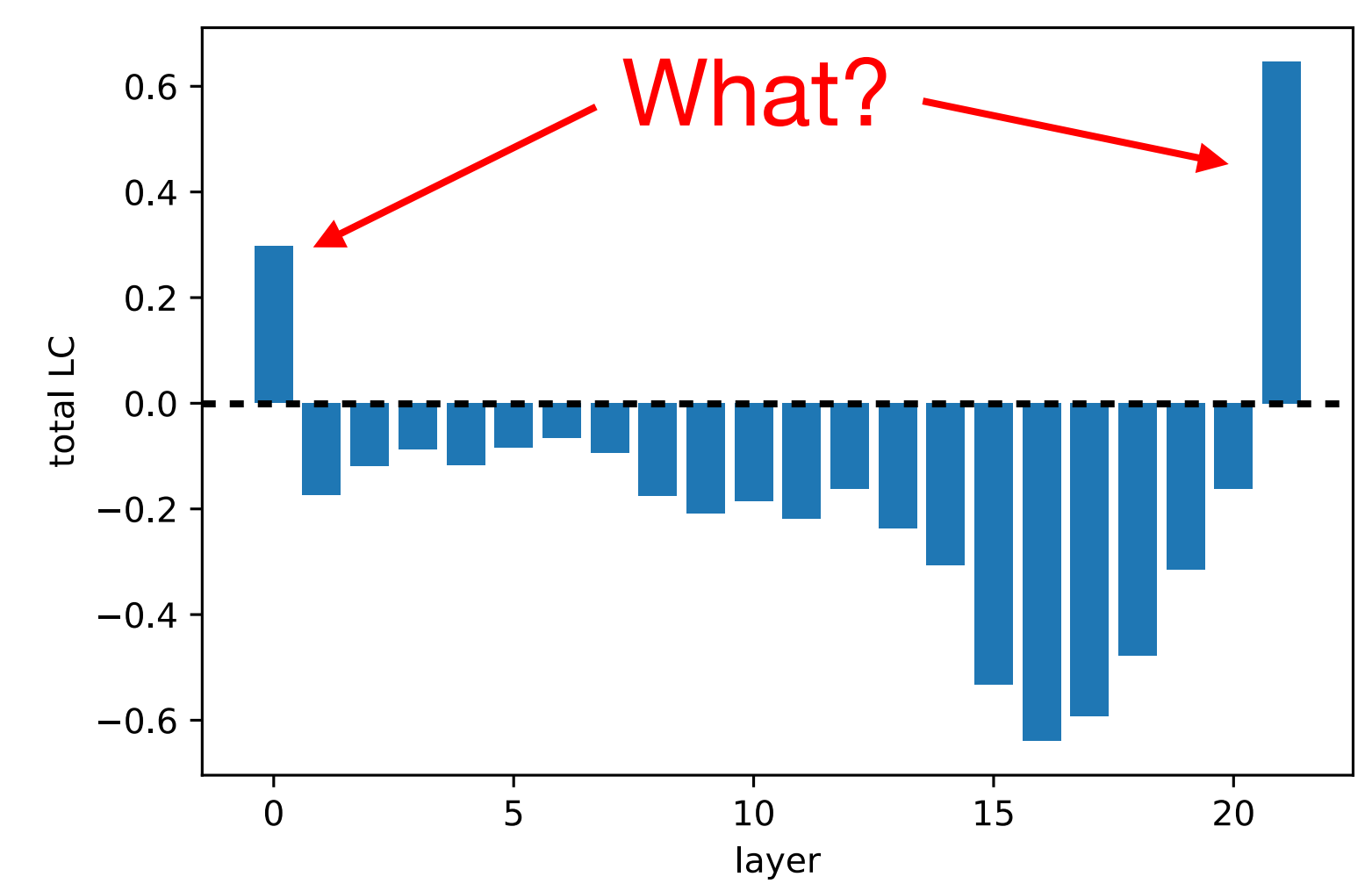
MNIST-FC



MNIST-LeNet



CIFAR-ResNet



Request for Plot: Approximate LCA

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

Approximate LCA:

$$\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

Approximate LCA:

$$\cancel{\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)} \quad \text{Biased}$$

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

Approximate LCA:

~~$\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$~~ Biased

$$\nabla L_{\text{batch } i+1}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

Approximate LCA:

~~$\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$~~ Biased

$$\nabla L_{\text{batch } i+1}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1x

Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

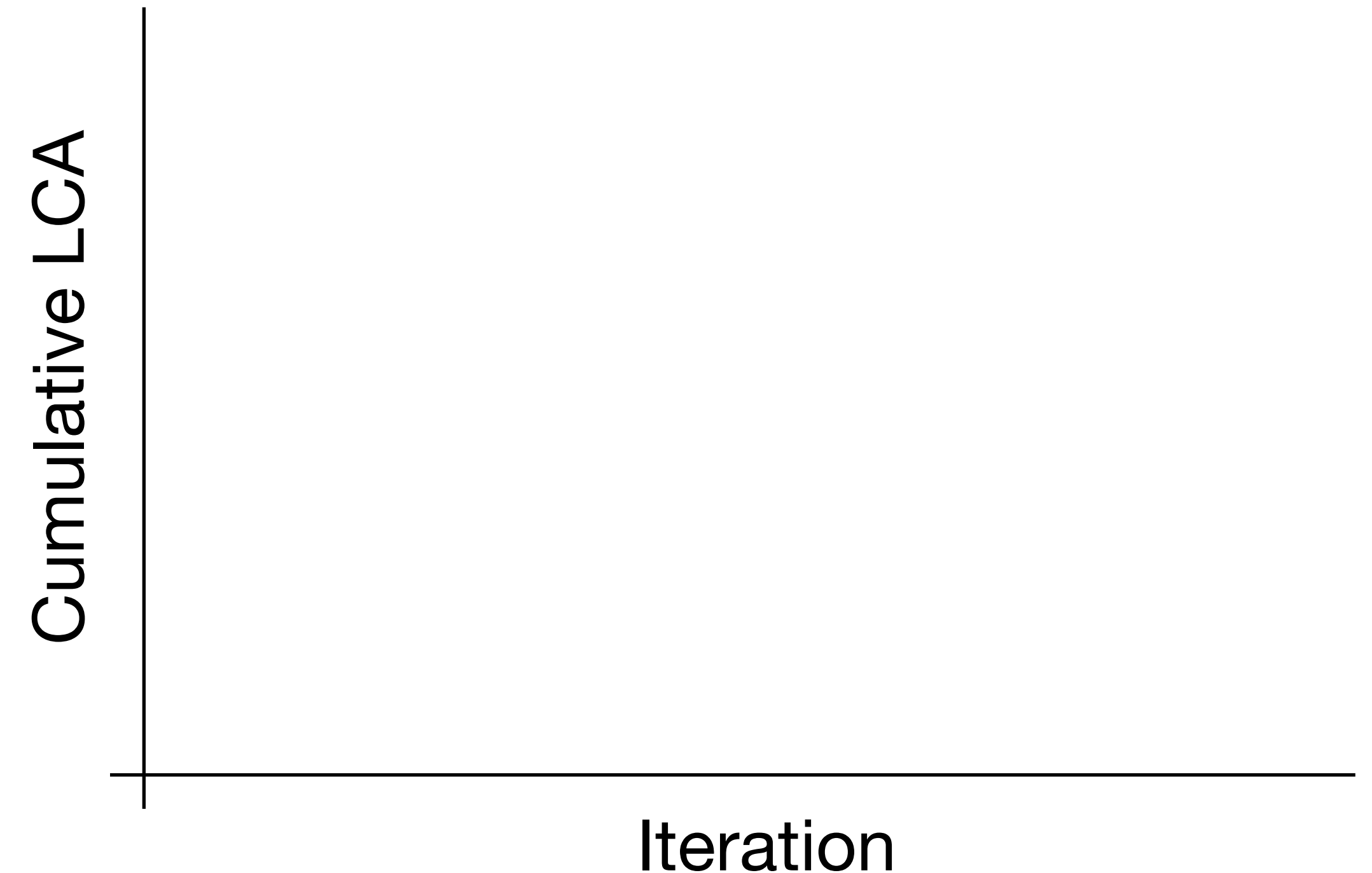
Approximate LCA:

~~$\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$~~ Biased

$$\nabla L_{\text{batch } i+1}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1x

Is this approximation good?



Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

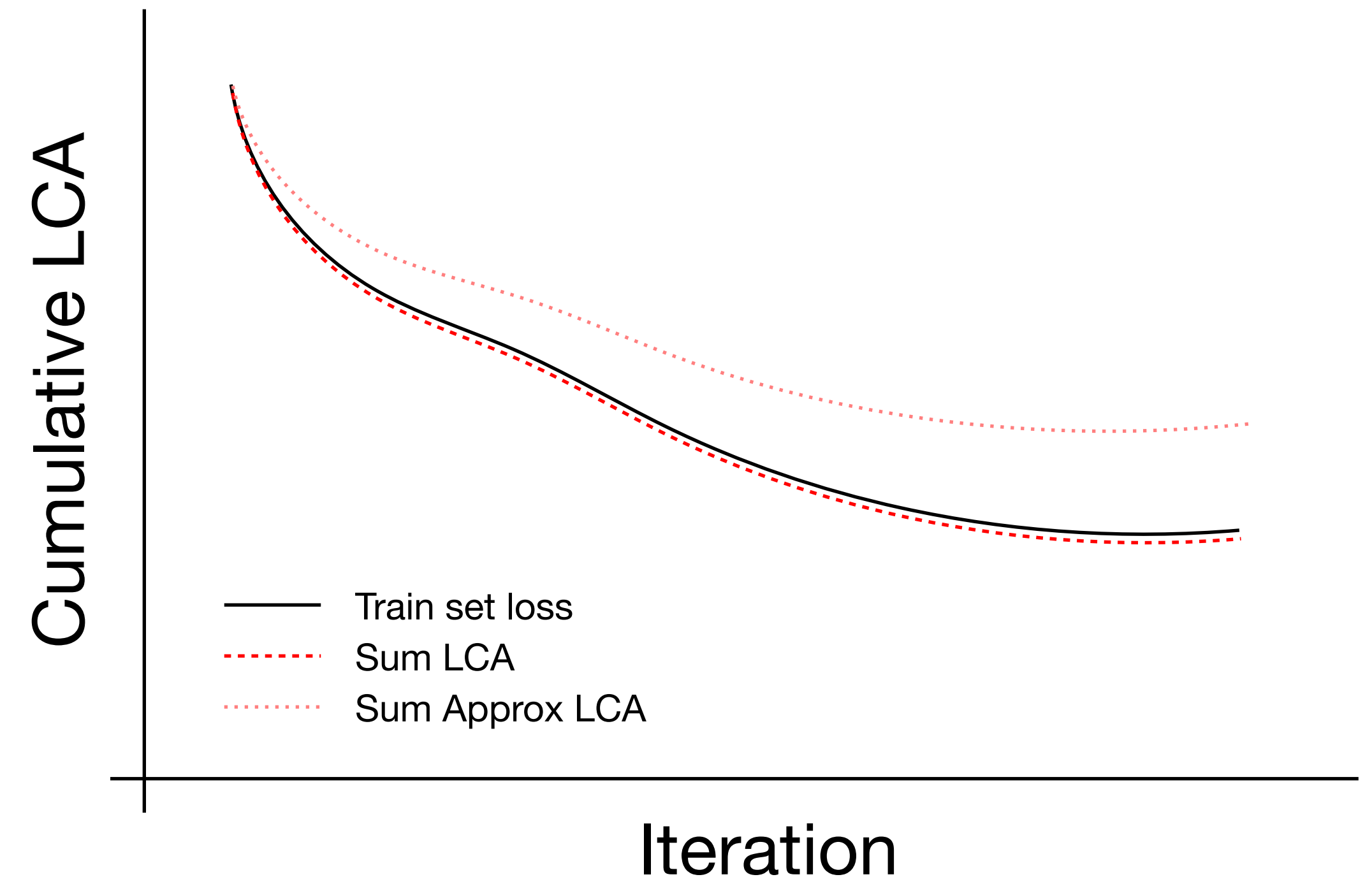
Approximate LCA:

$$\cancel{\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)} \quad \text{Biased}$$

$$\nabla L_{\text{batch } i+1}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1x

Is this approximation good?



Request for Plot: Approximate LCA

Example: training set of 1000 batches

Ordinary training:

Compute per batch: 1x

Original LCA:

$$\nabla L_{\text{train}}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1000x

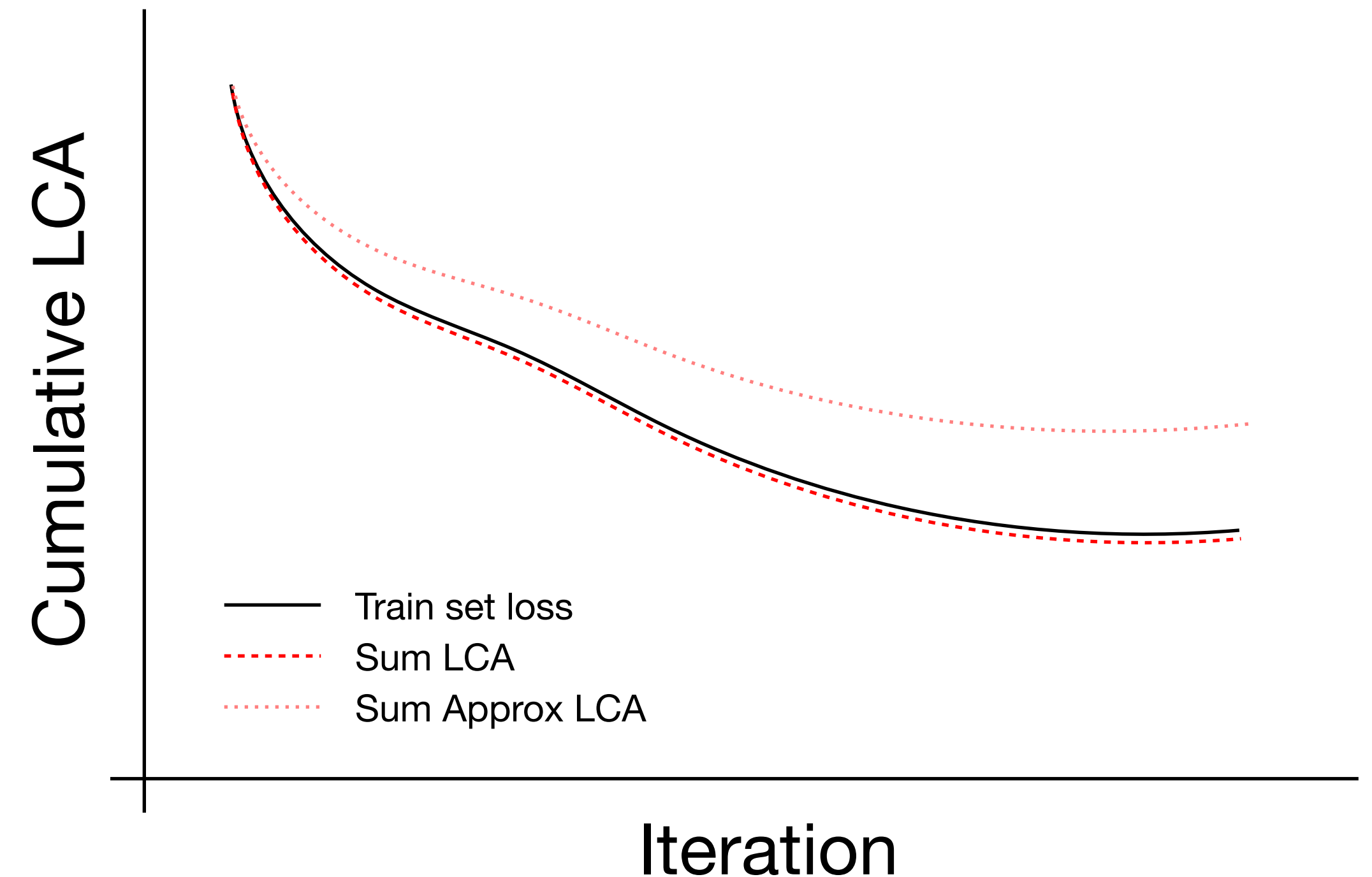
Approximate LCA:

$$\cancel{\nabla L_{\text{batch } i}(\theta_i) \cdot (\theta_{i+1} - \theta_i)} \quad \text{Biased}$$

$$\nabla L_{\text{batch } i+1}(\theta_i) \cdot (\theta_{i+1} - \theta_i)$$

Compute per batch: 1x

Is this approximation good?



Next questions:

- How to incorporate RK-4?
- If this works, are per-layer sums also well approximated?

Idea #2: Computational Ethics

How we used to do Computer Vision

- Humans create features
- Failures → humans create smarter features
- Endless Whack-a-mole



What worked better

- Models learn their own features

How we align models now

- Humans write prompts (constitutions, specs)
- Failures → humans create smarter prompts
- Endless Whack-a-mole

See e.g. 24k token Claude system prompt:
“...Claude provides emotional support alongside accurate medical or psychological information or terminology where relevant.”

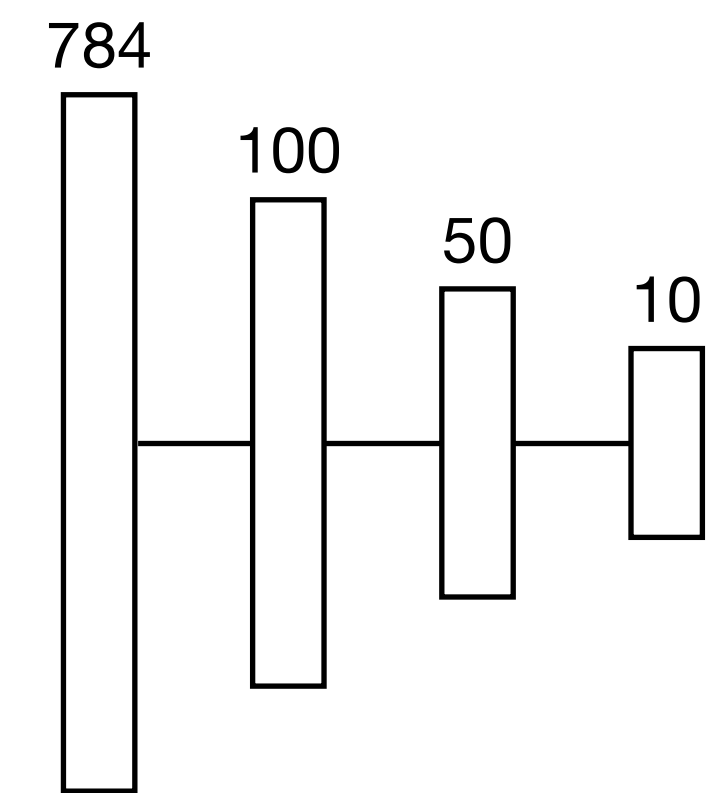
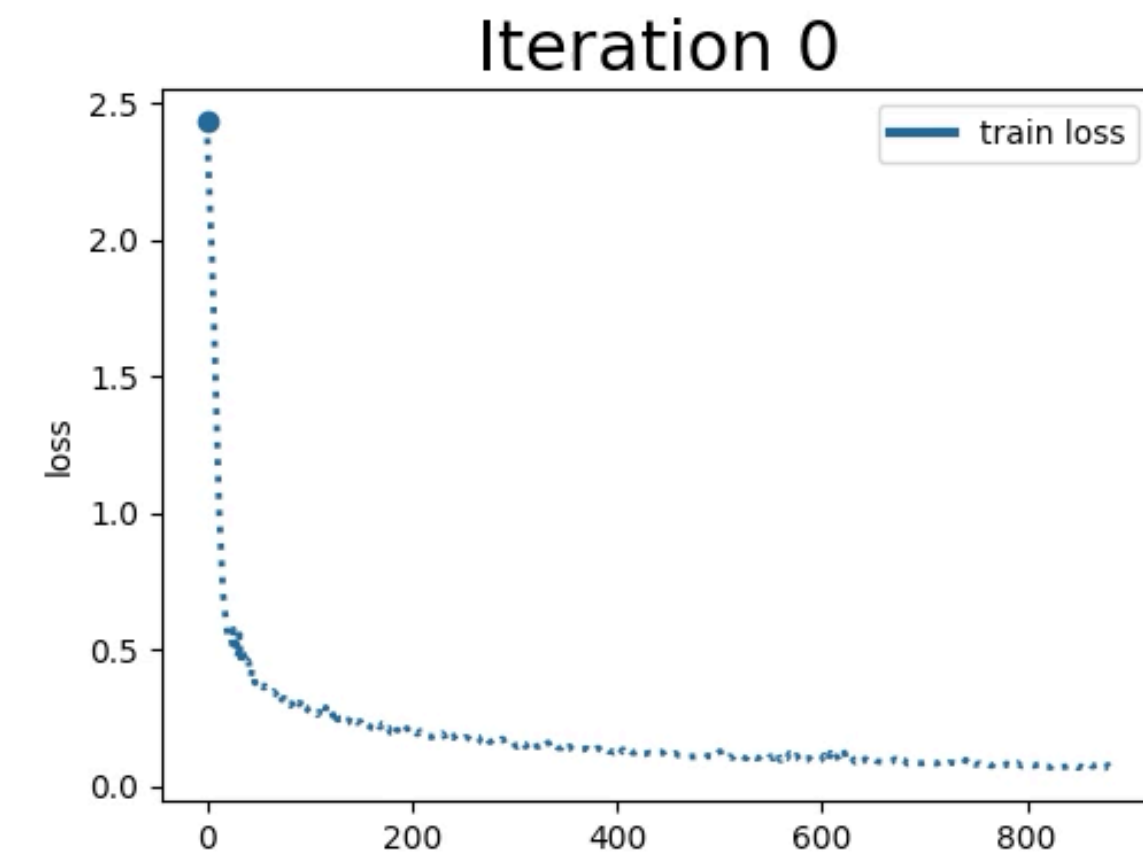
What might work better?

- Models learn their own ethics
- Start with using models to map the existing, diverse landscape of human ethics

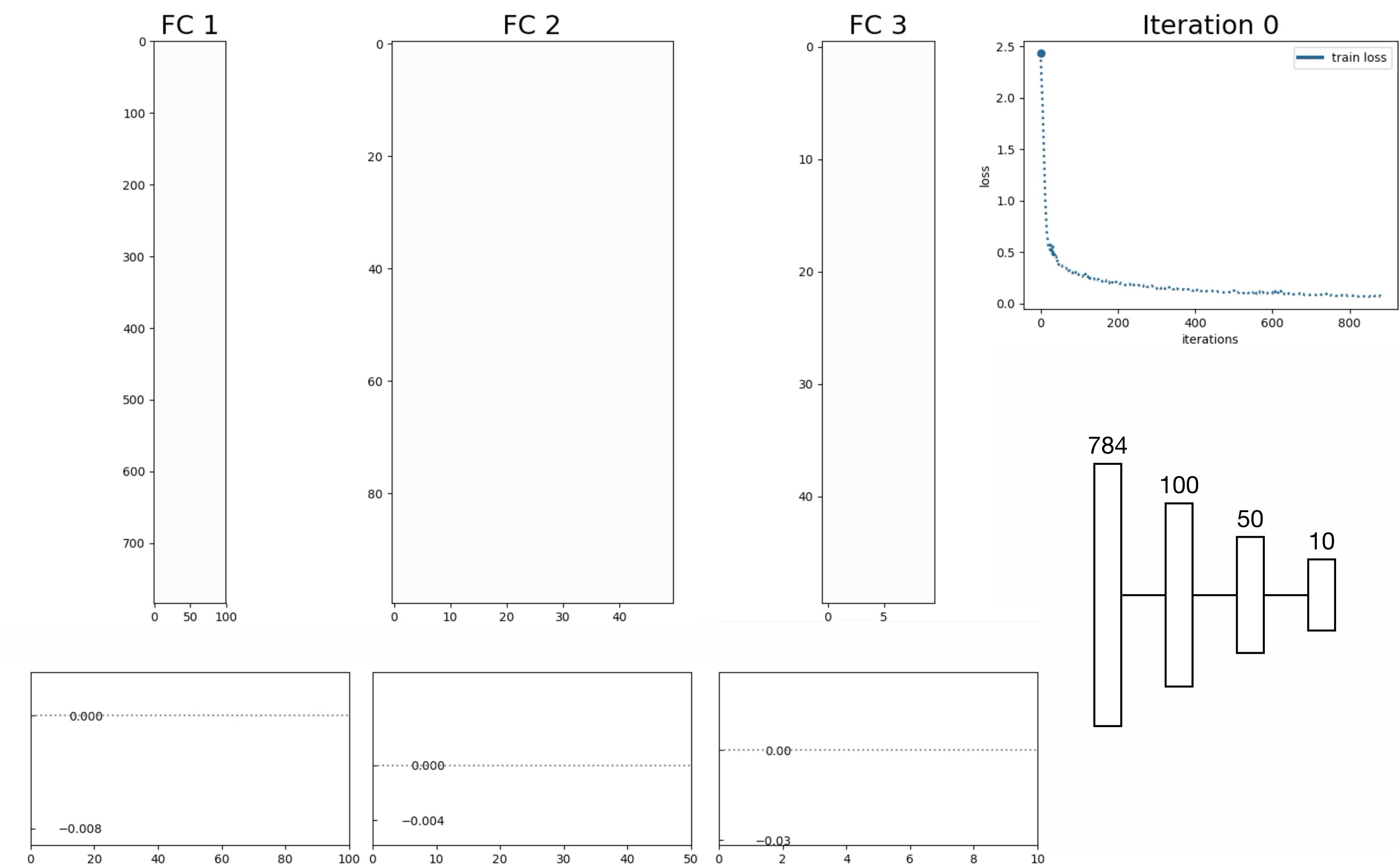
Appendix: extra LCA slides

LCA: Loss Change Allocation for Neural Network Training

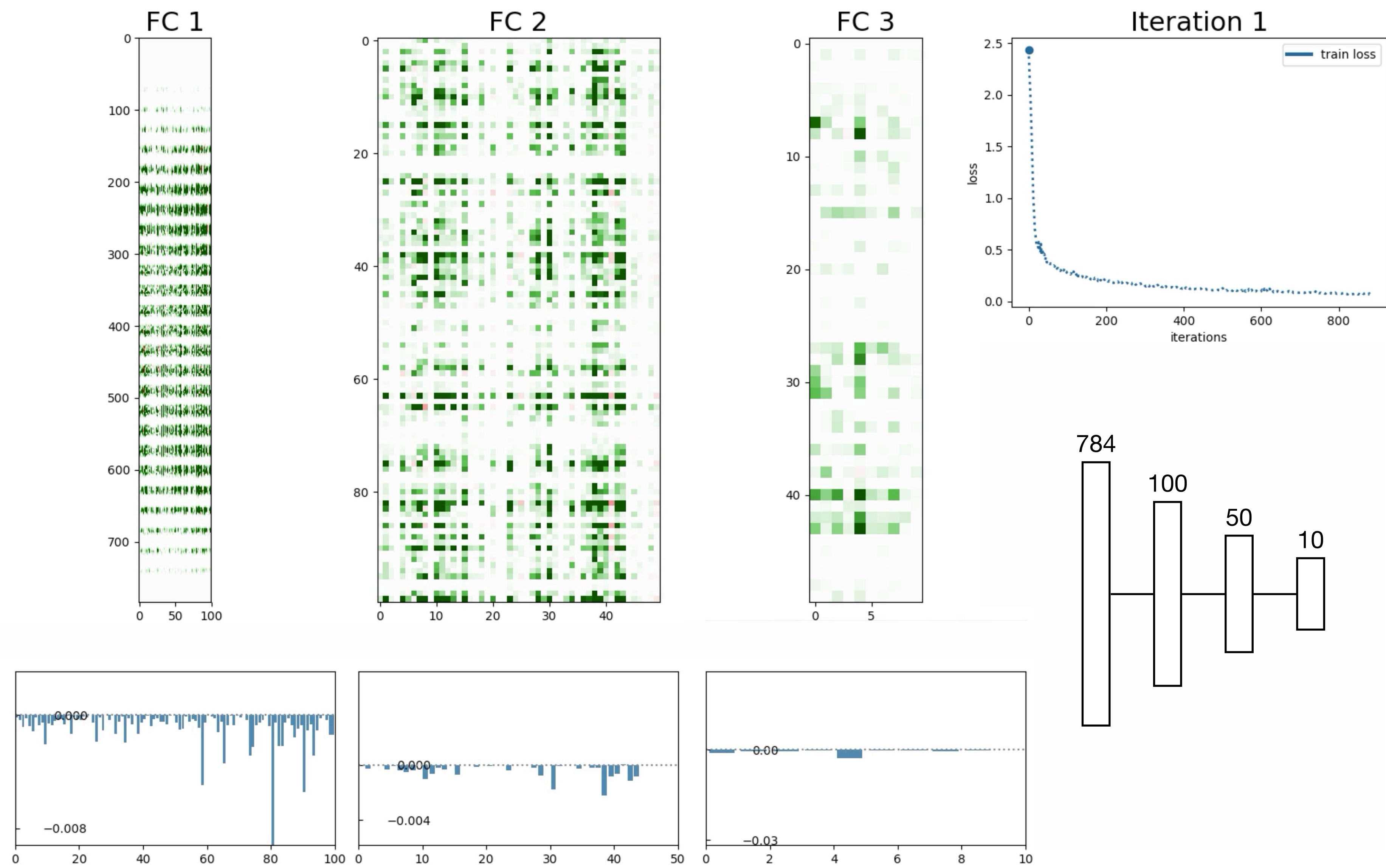
LCA: Loss Change Allocation for Neural Network Training



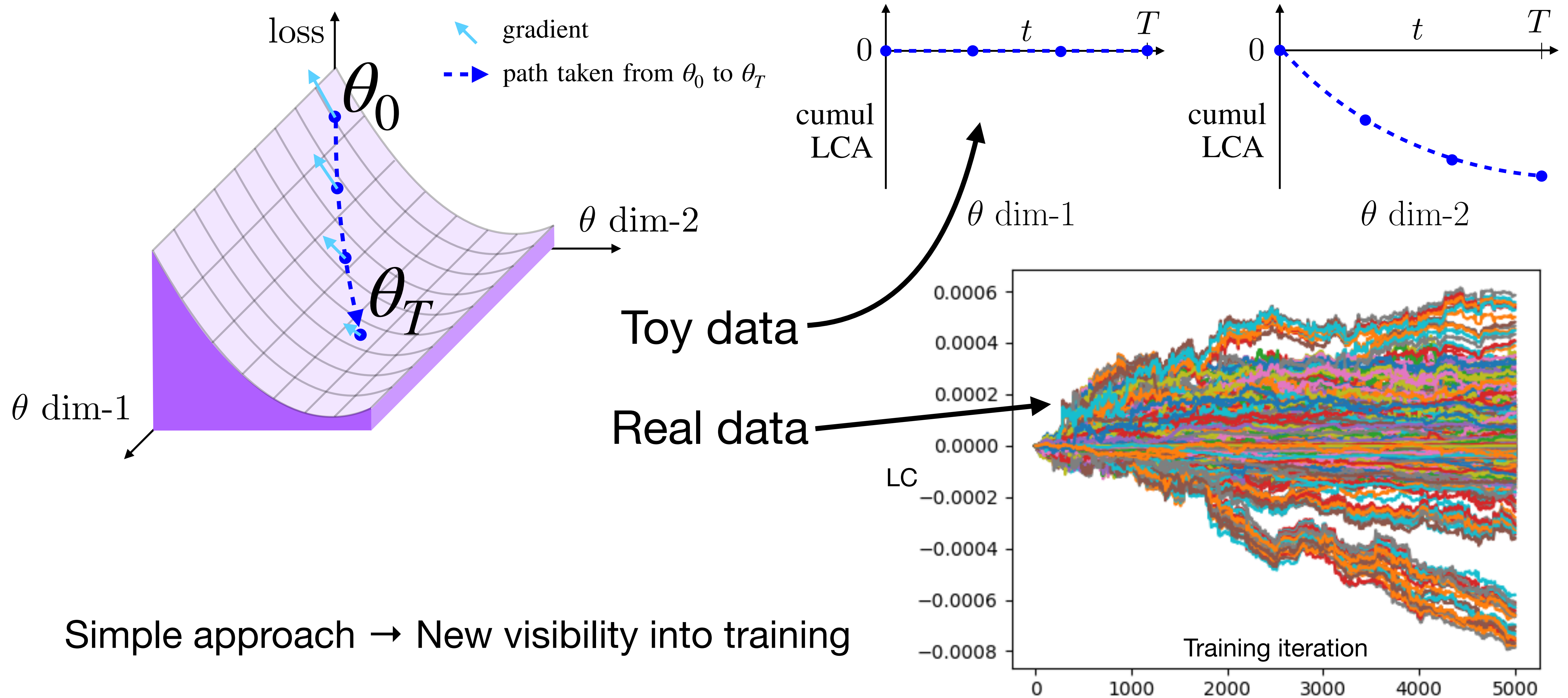
LCA: Loss Change Allocation for Neural Network Training



LCA: Loss Change Allocation for Neural Network Training



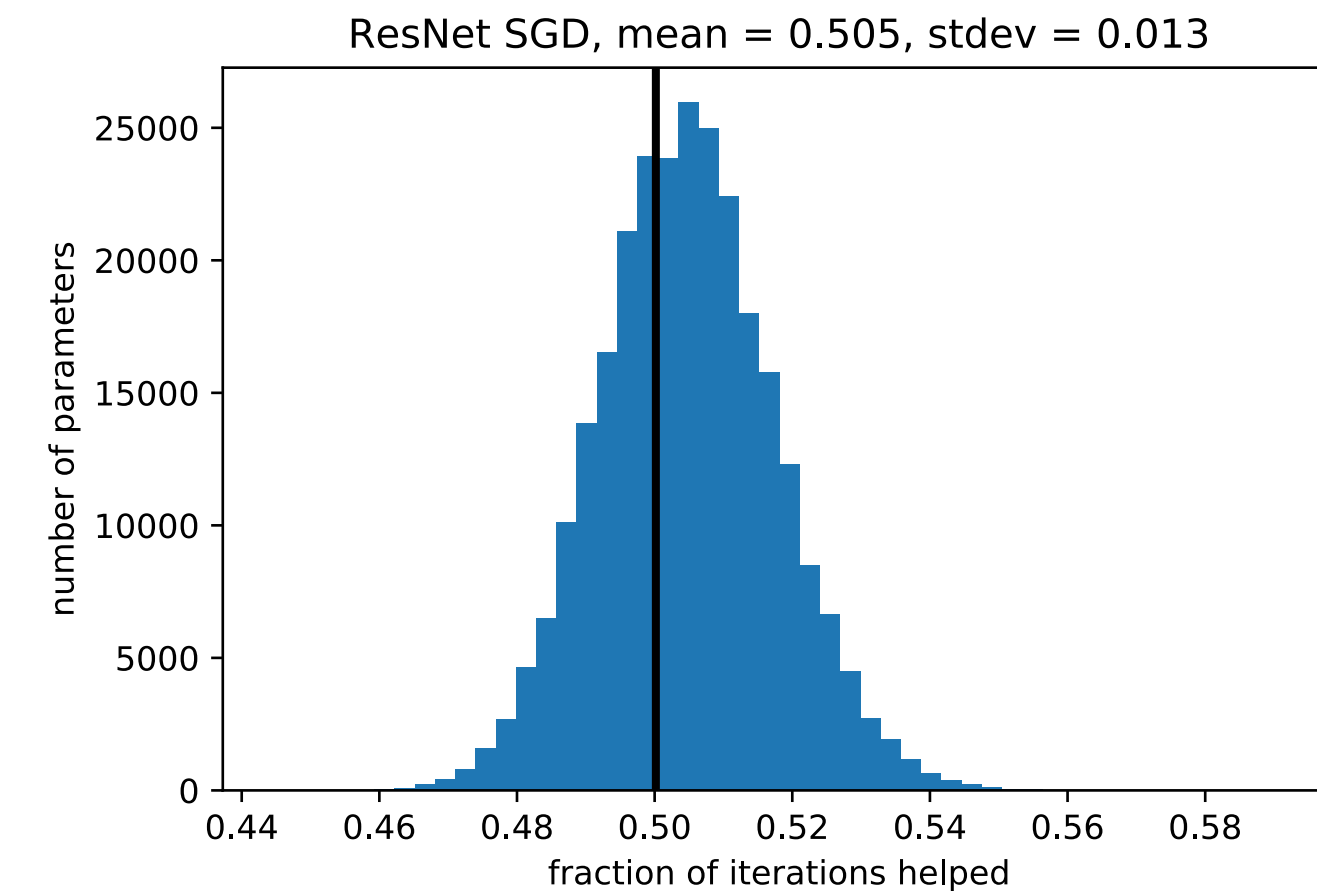
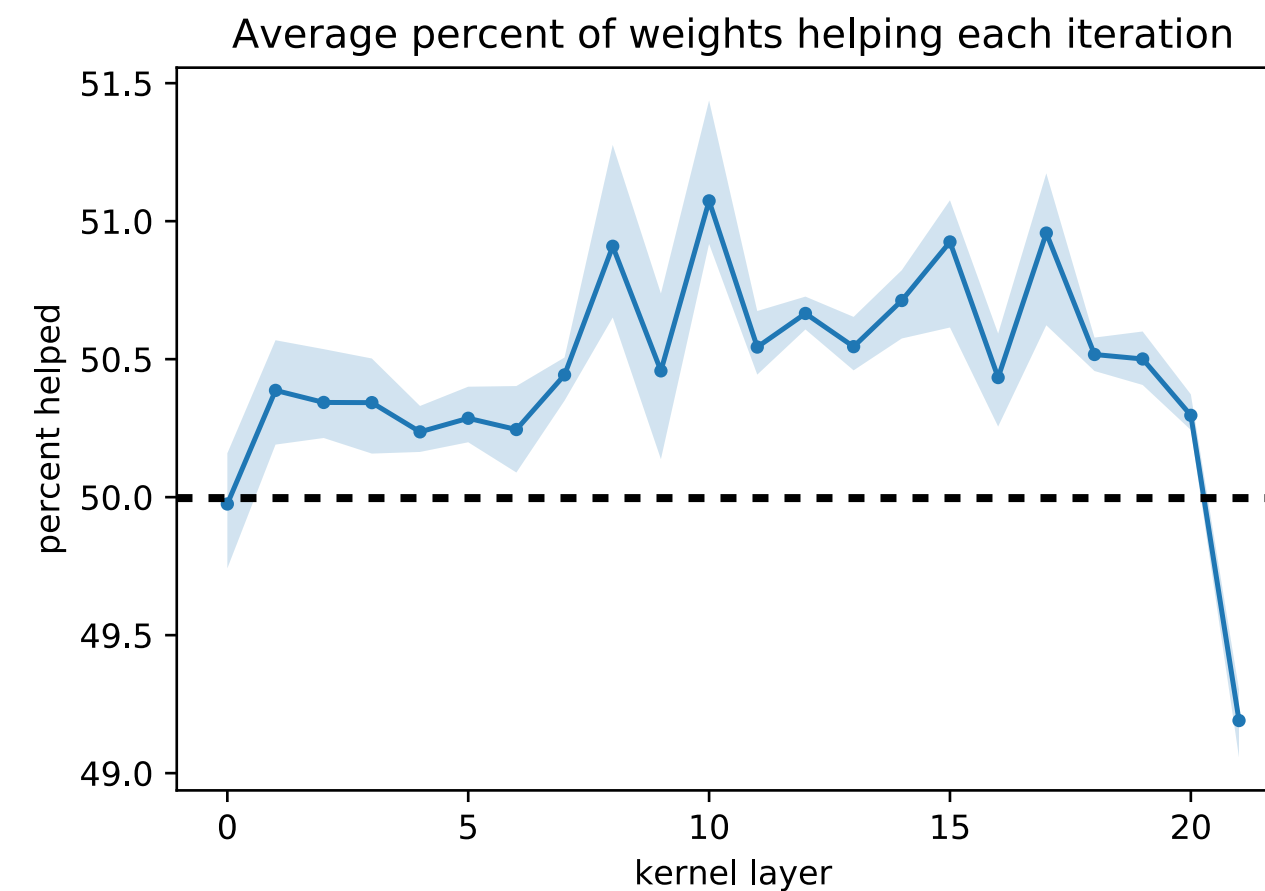
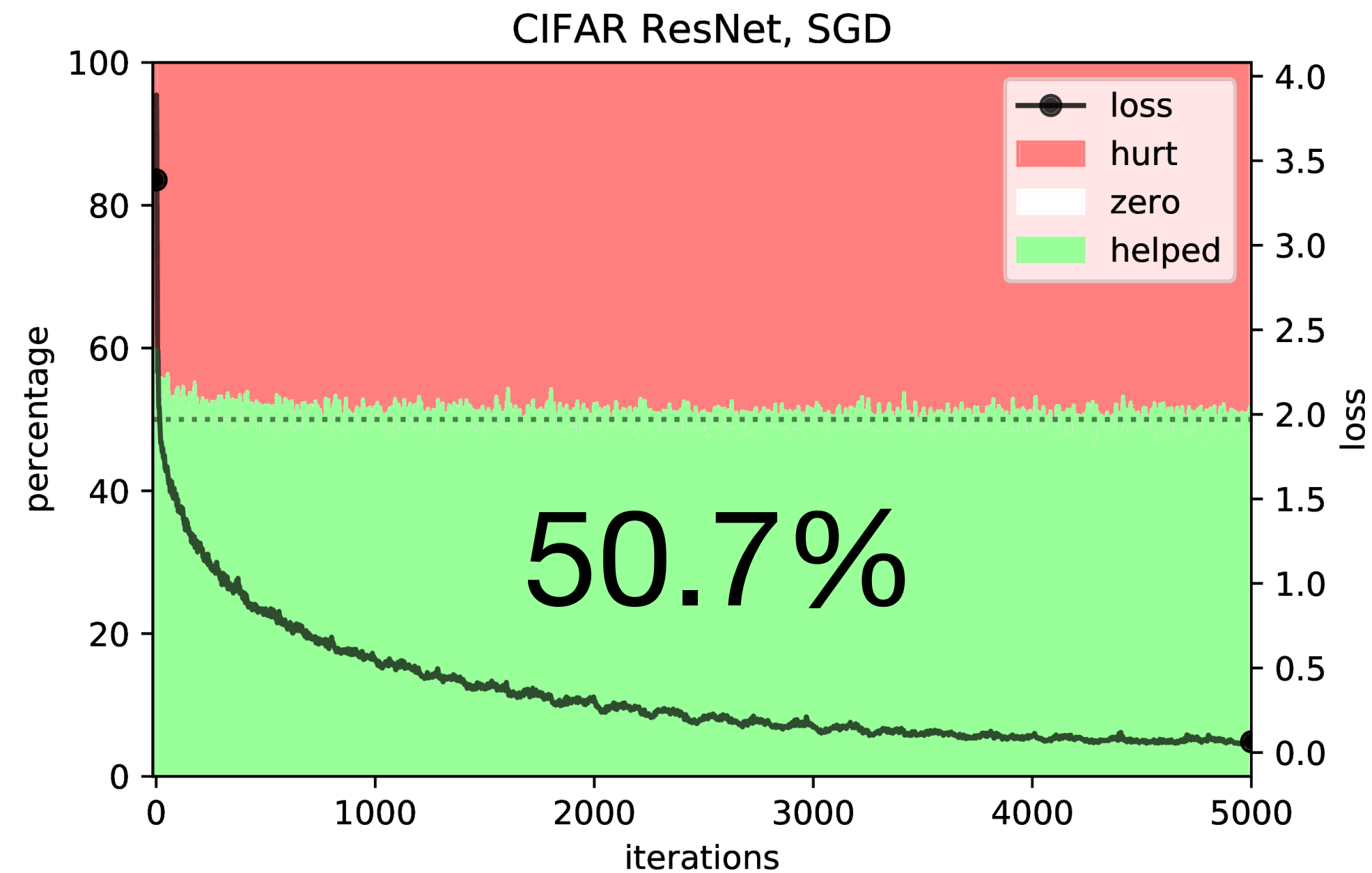
LCA: Loss Change Allocation for Neural Network Training



LCA: Loss Change Allocation for Neural Network Training

1 Training is noisy

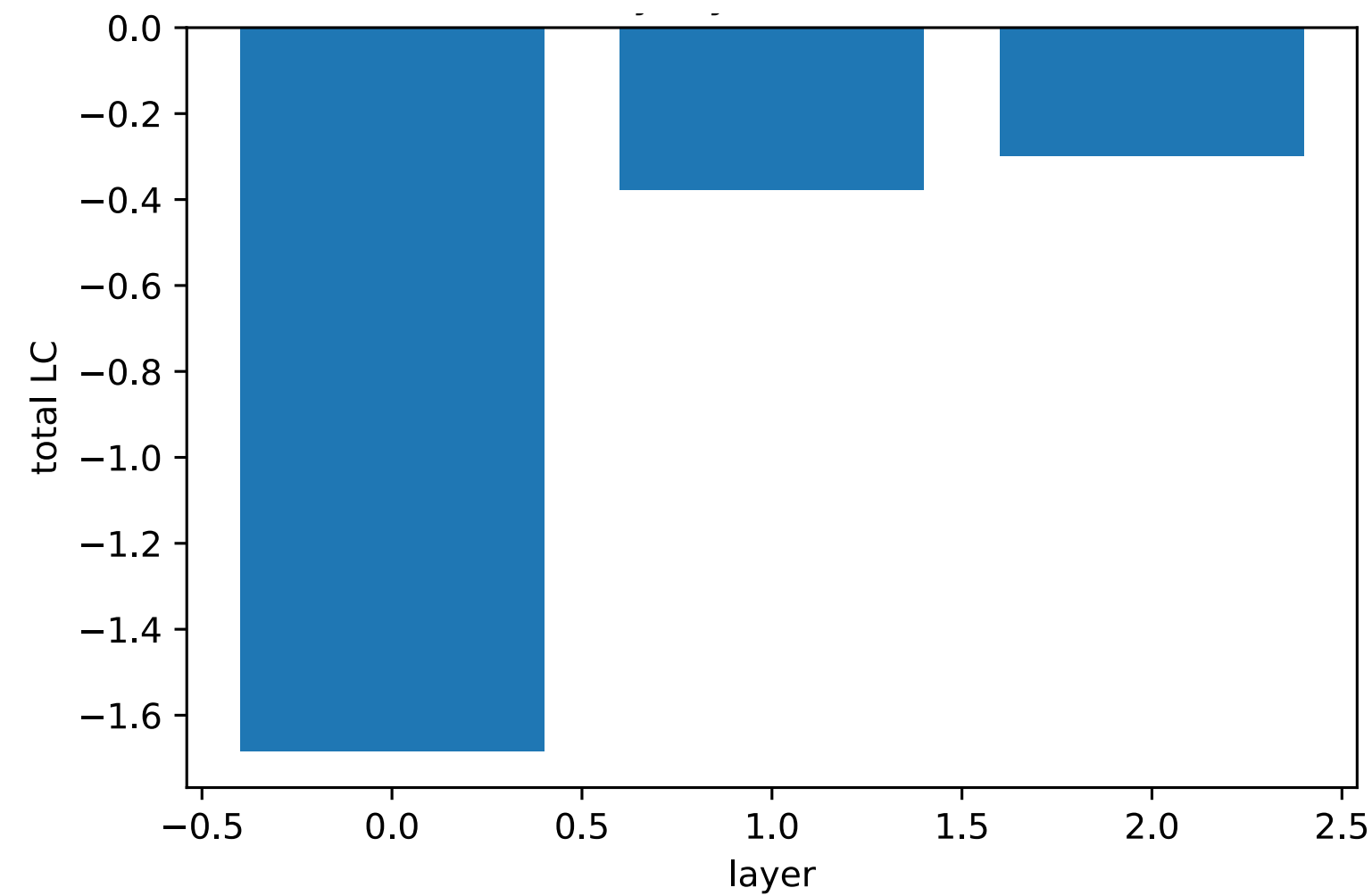
- Holds for all layers
- Holds for all params
- Holds for many hyperparams (50.3% – 51.6%)



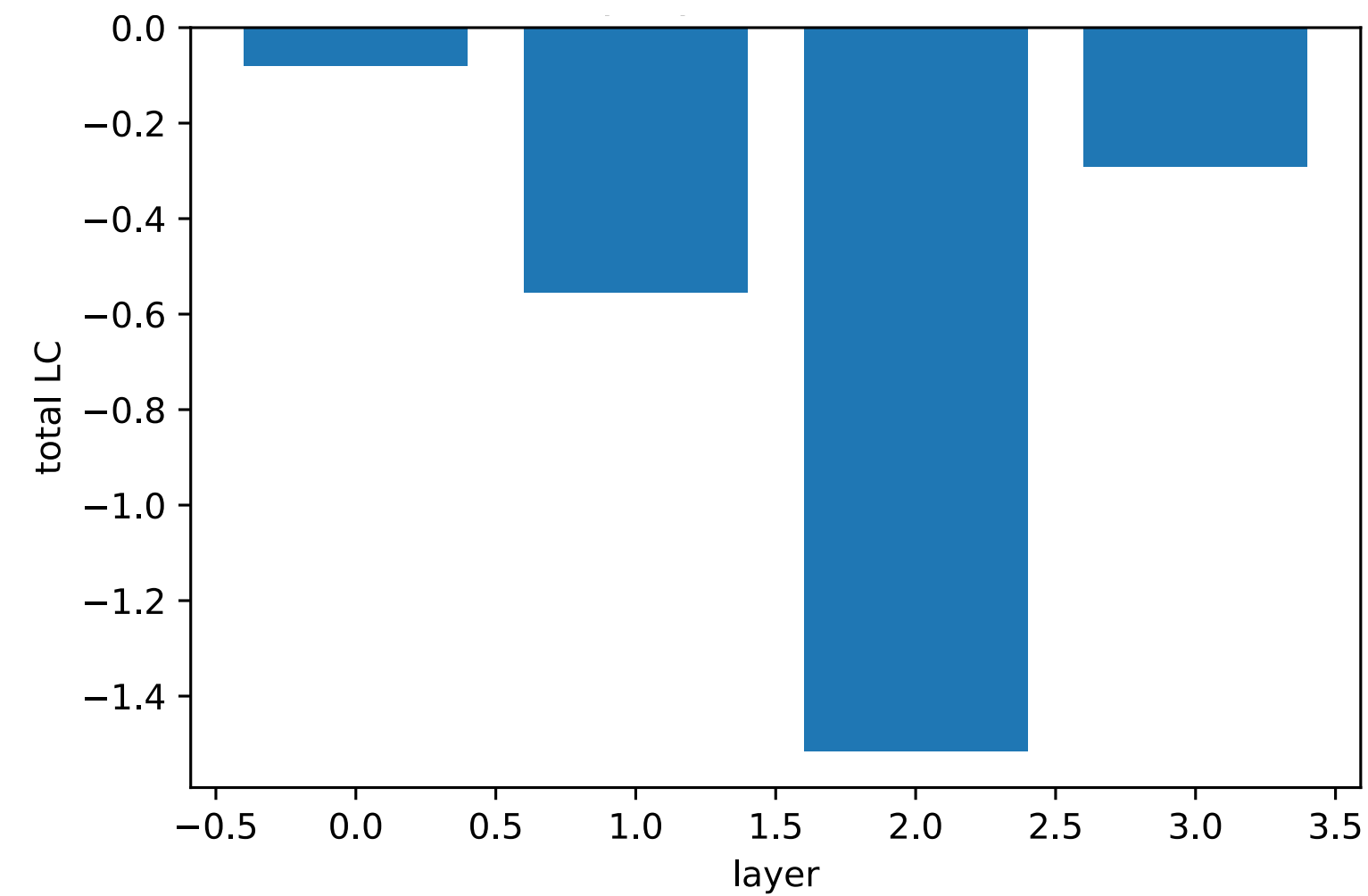
LCA: Loss Change Allocation for Neural Network Training

- 1 Training is noisy
- 2 Some layers go backwards

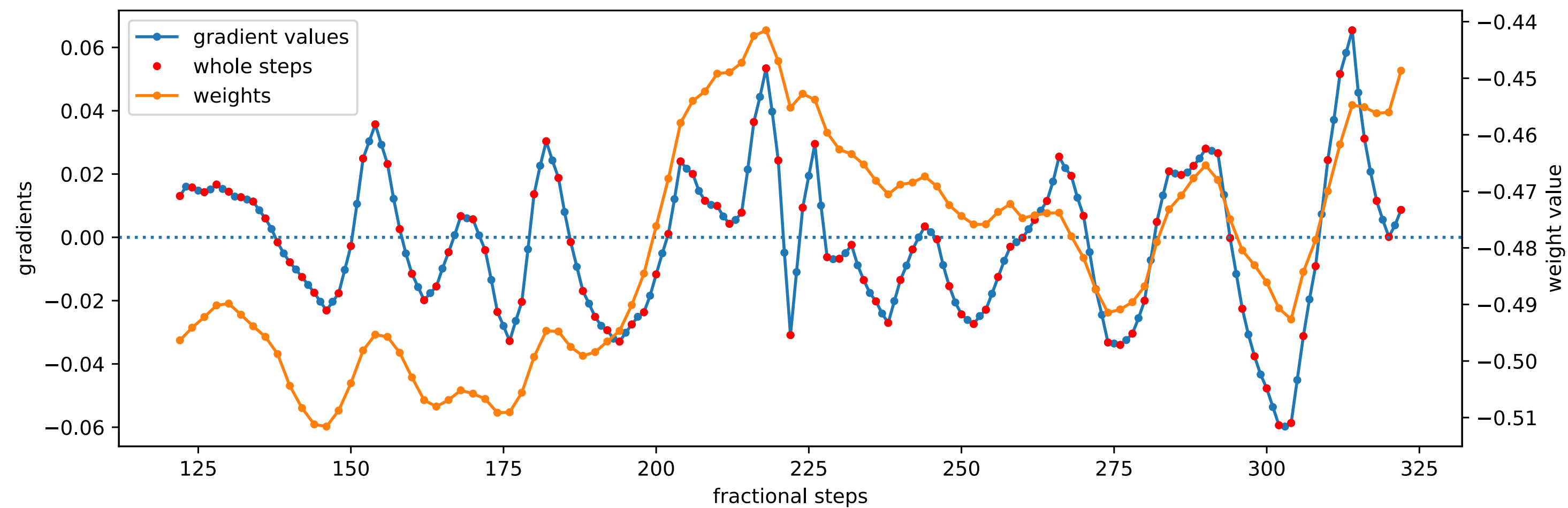
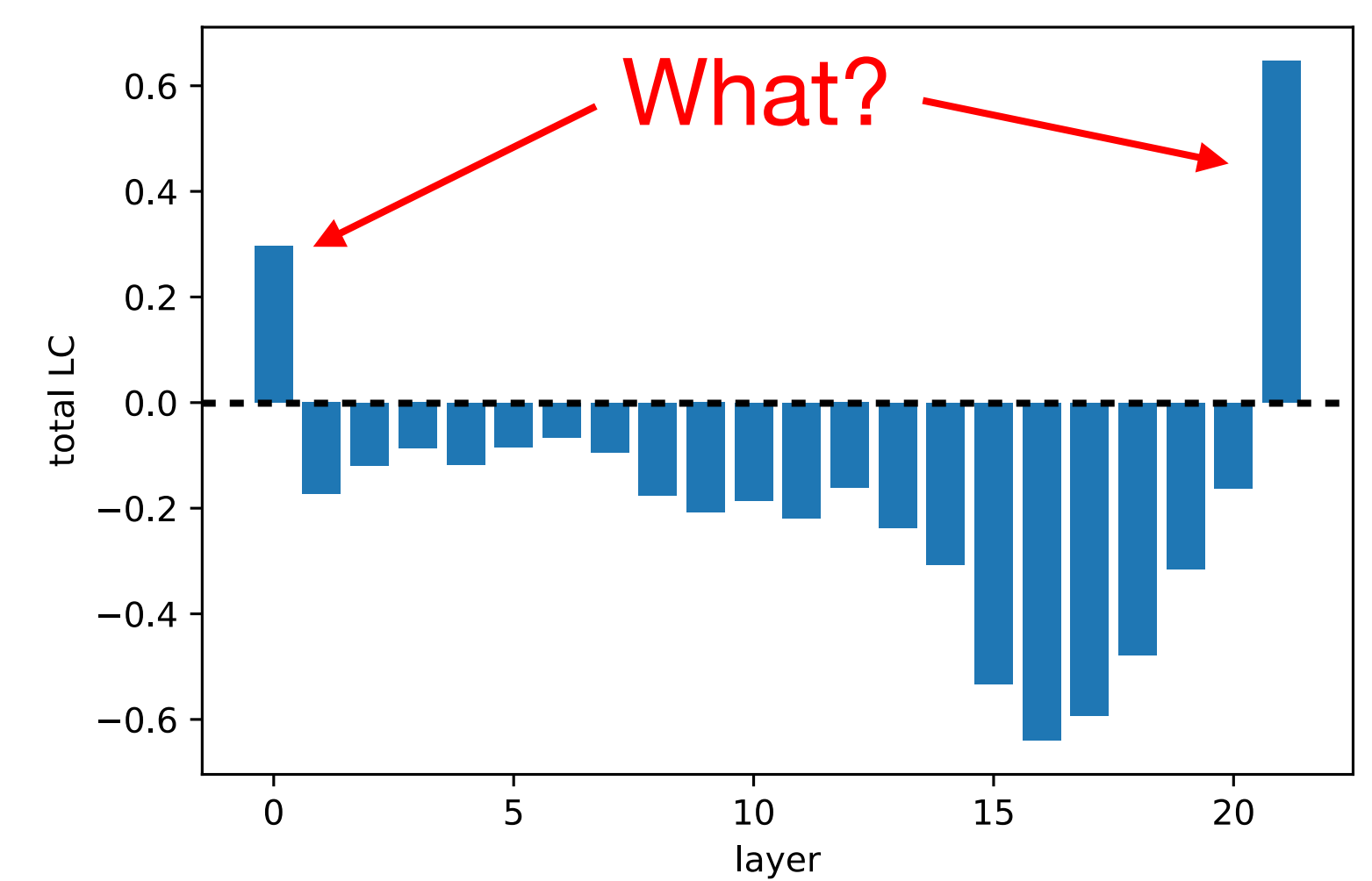
MNIST-FC



MNIST-LeNet



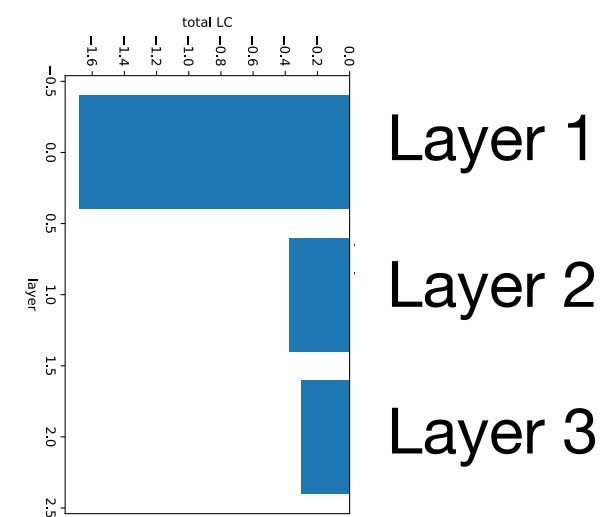
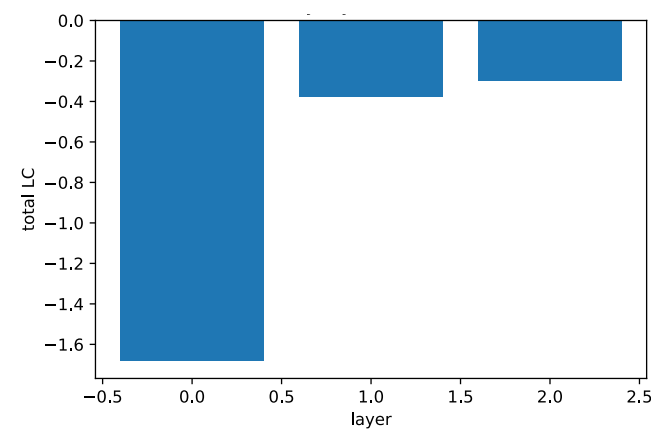
CIFAR-ResNet



LCA: Loss Change Allocation for Neural Network Training

- 1 Training is noisy
- 2 Some layers go backwards

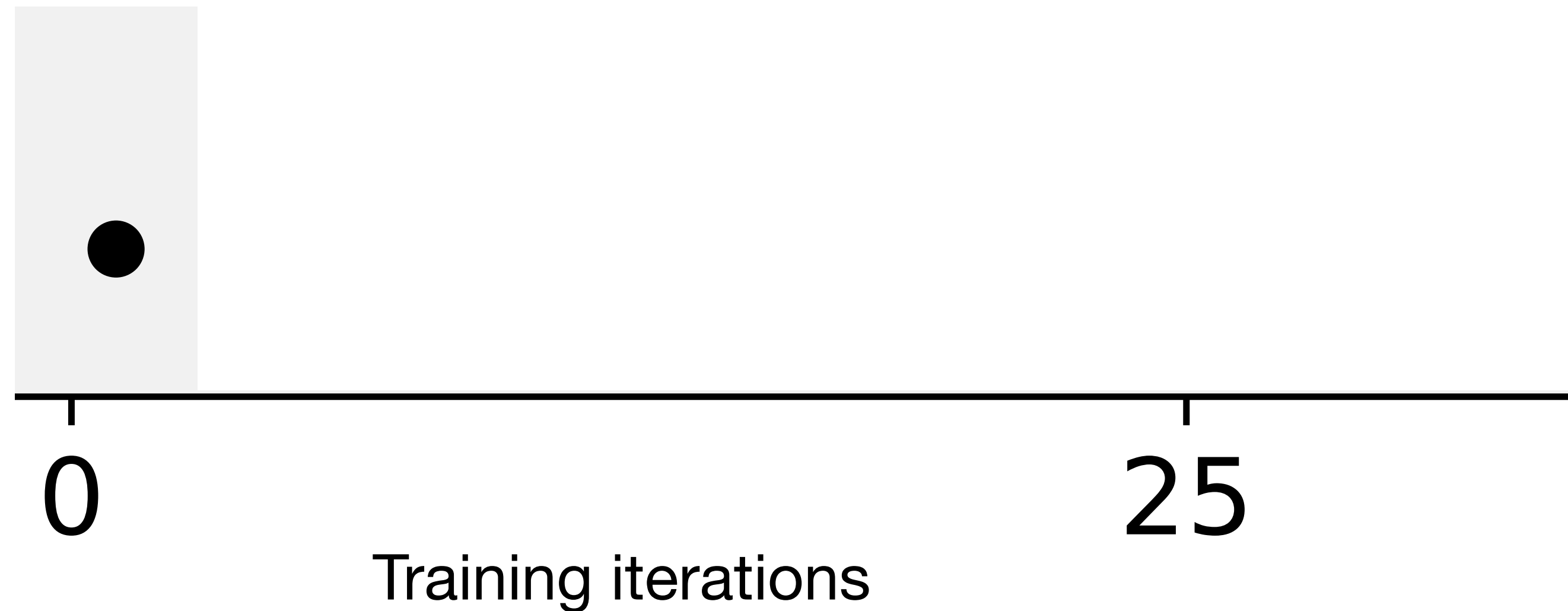
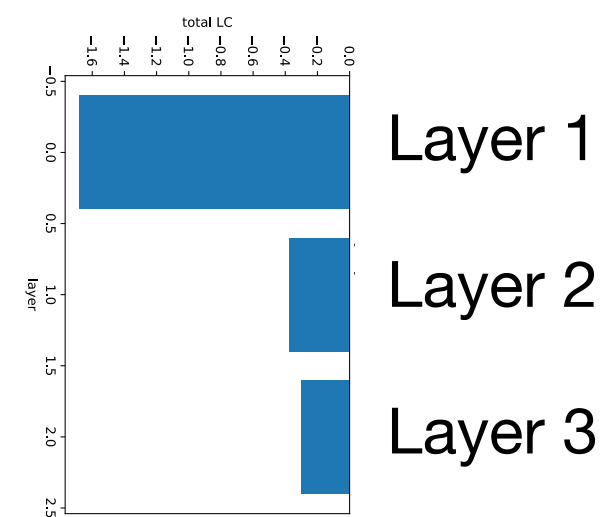
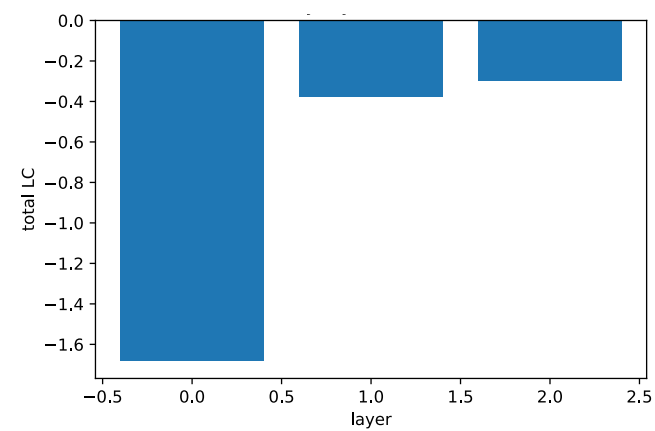
MNIST-FC



LCA: Loss Change Allocation for Neural Network Training

- 1 Training is noisy
- 2 Some layers go backwards

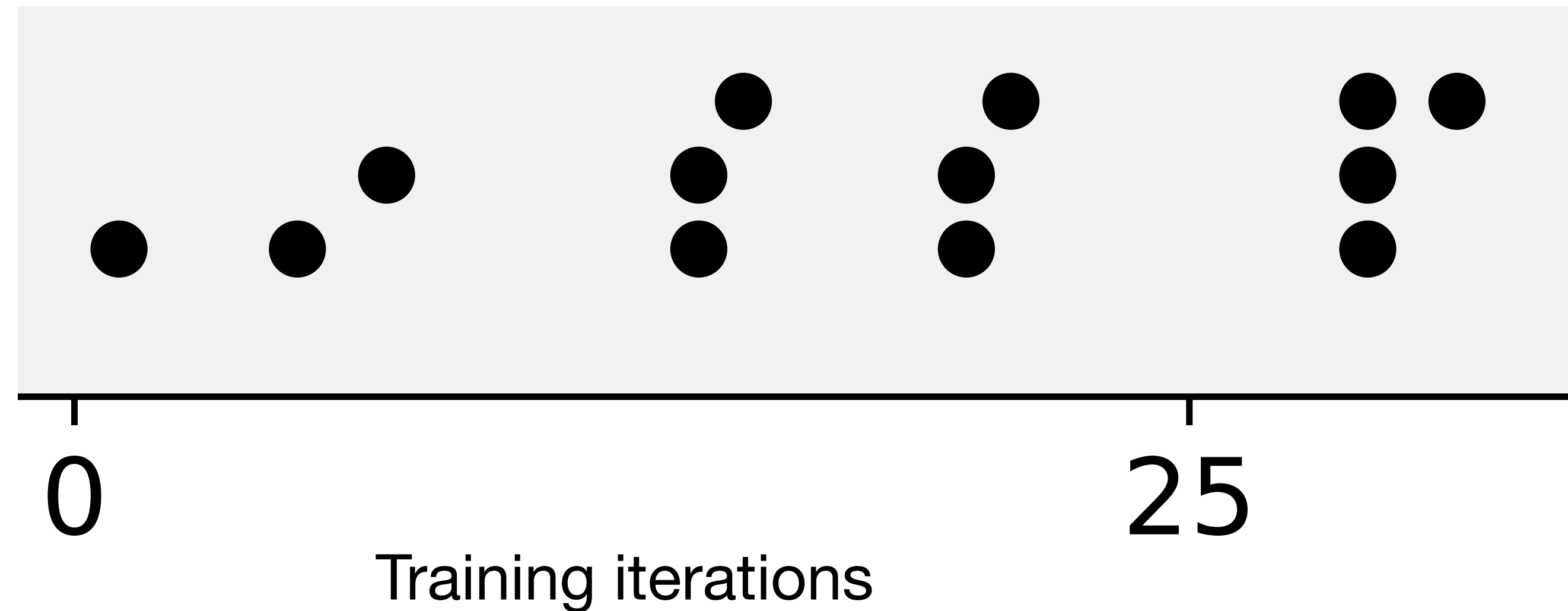
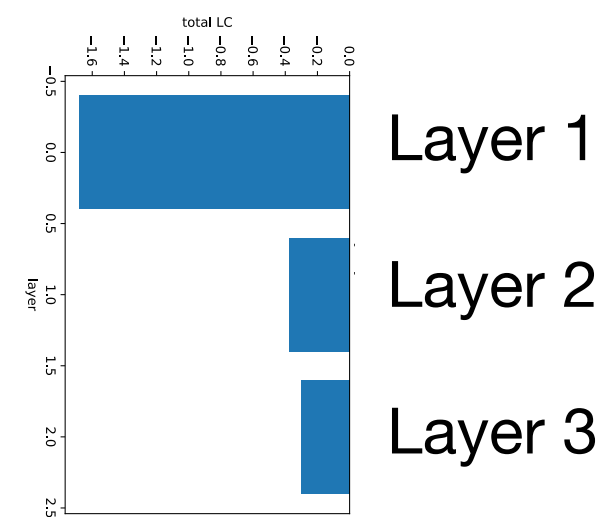
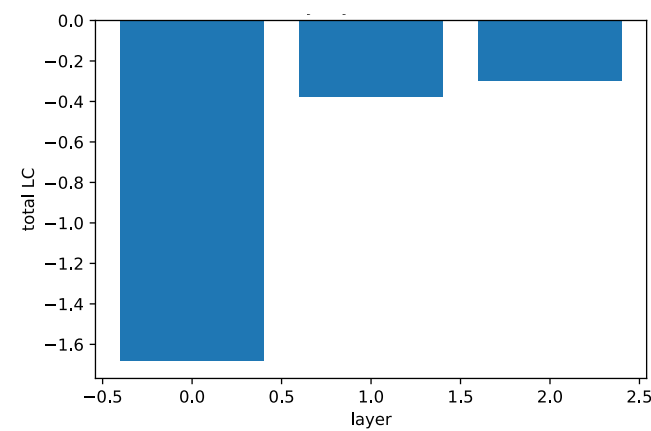
MNIST-FC



LCA: Loss Change Allocation for Neural Network Training

- 1 Training is noisy
- 2 Some layers go backwards

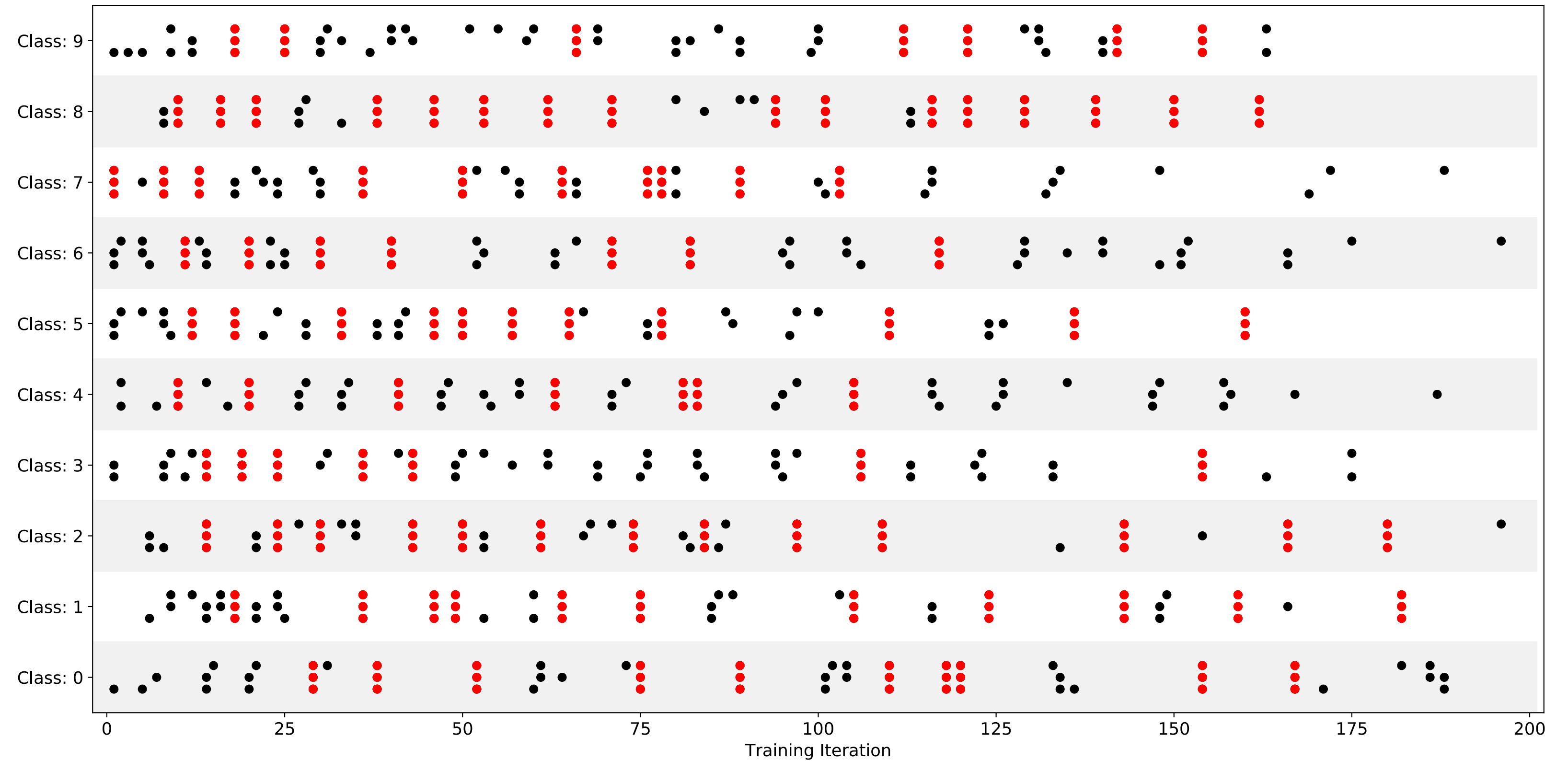
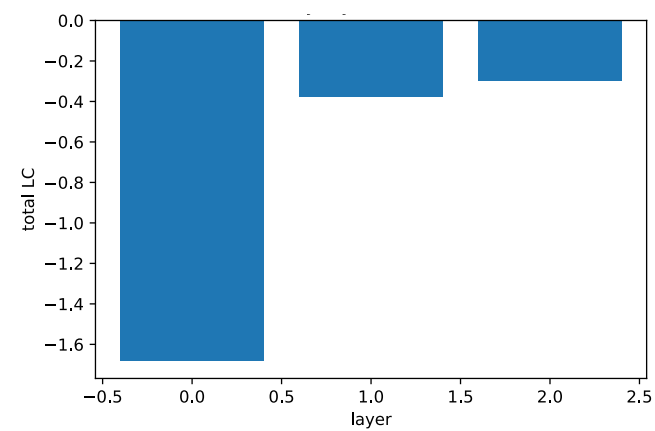
MNIST-FC



LCA: Loss Change Allocation for Neural Network Training

- 1 Training is noisy
- 2 Some layers go backwards
- 3 Some micro-learning is synchronized

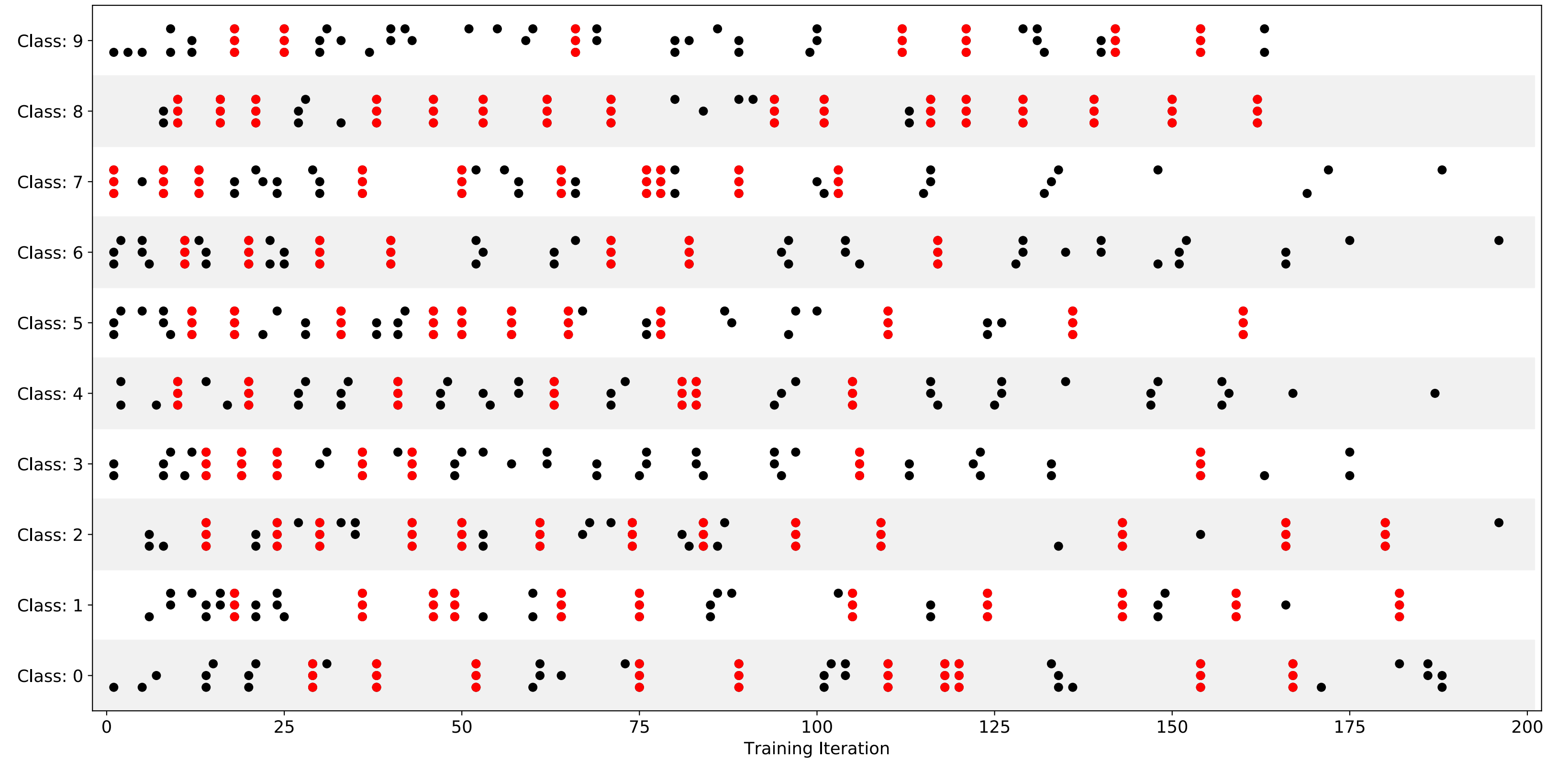
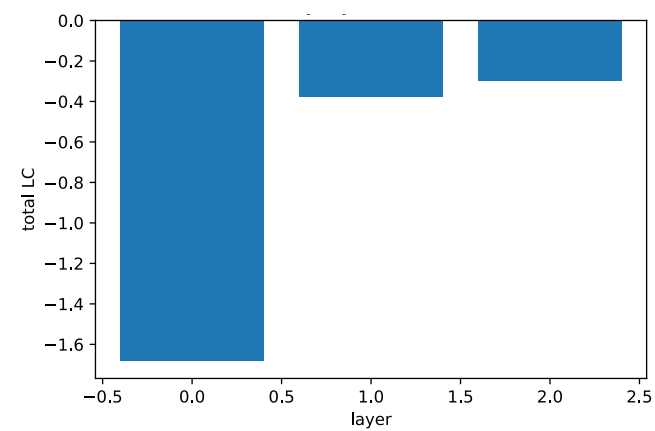
MNIST-FC



LCA: Loss Change Allocation for Neural Network Training

- 1 Training is noisy
- 2 Some layers go backwards
- 3 Some micro-learning is synchronized

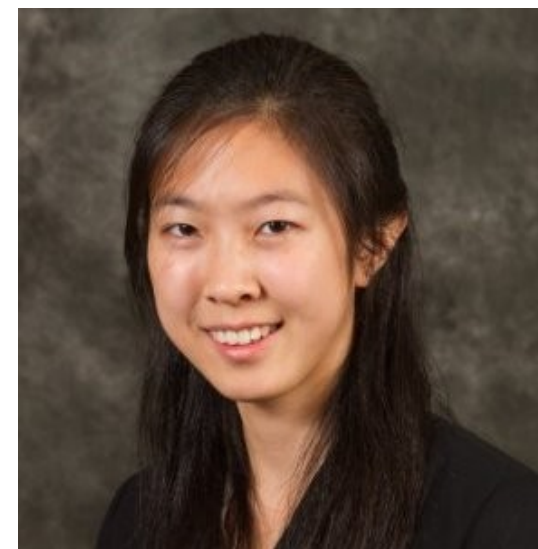
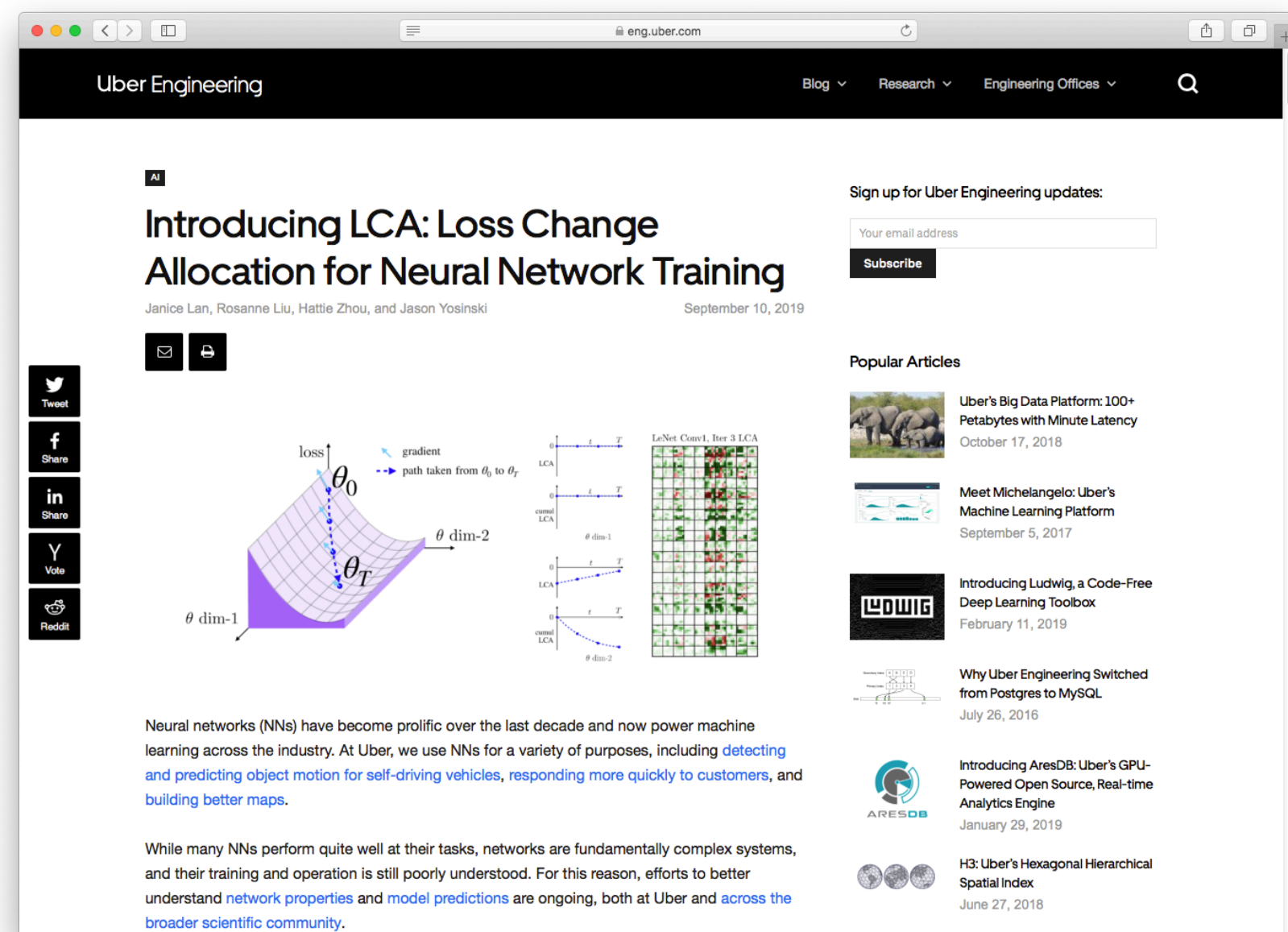
MNIST-FC



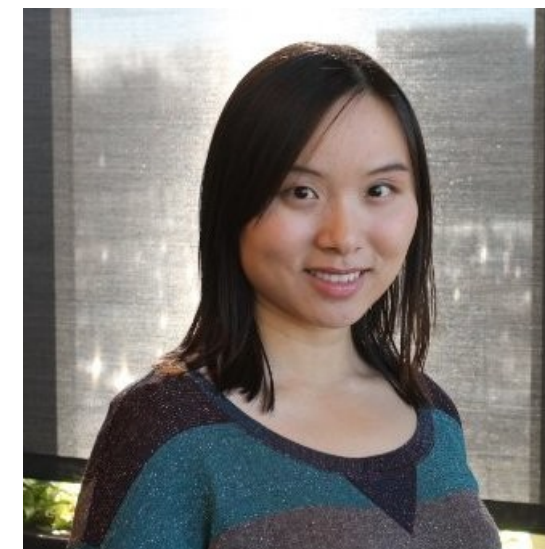
LCA: Loss Change Allocation for Neural Network Training

NeurIPS 2019.

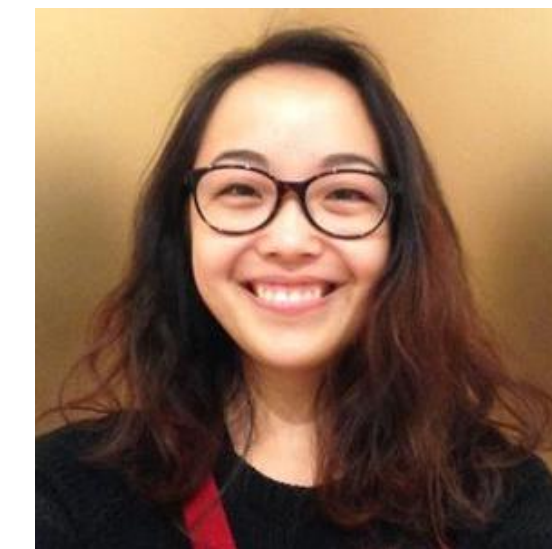
Blog: <https://eng.uber.com/loss-change-allocation/>



Janice Lan



Rosanne Liu



Hattie Zhou



Jason Yosinski