

AC-GC: LOSSY ACTIVATION COMPRESSION WITH GUARANTEED CONVERGENCE

R. David Evans, Tor M. Aamodt

Neural Information Processing Systems
(NeurIPS), December 7, 2021
Poster Spot A2, 1630 PST – 1800 PST



THE UNIVERSITY
OF BRITISH COLUMBIA
Electrical and Computer
Engineering

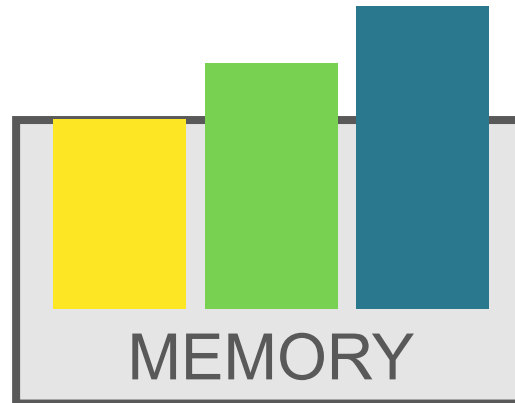


AC-GC: IN A NUTSHELL

Neural network training
cares about



Lossy Compression
decreases
memory footprint



How to set
Compression Rate



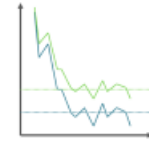
MOTIVATION



KEY
CONTRIBUTIONS



BOUNDING
CONVERGENCE



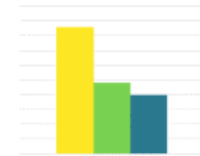
SIMPLIFYING
THE PROBLEM



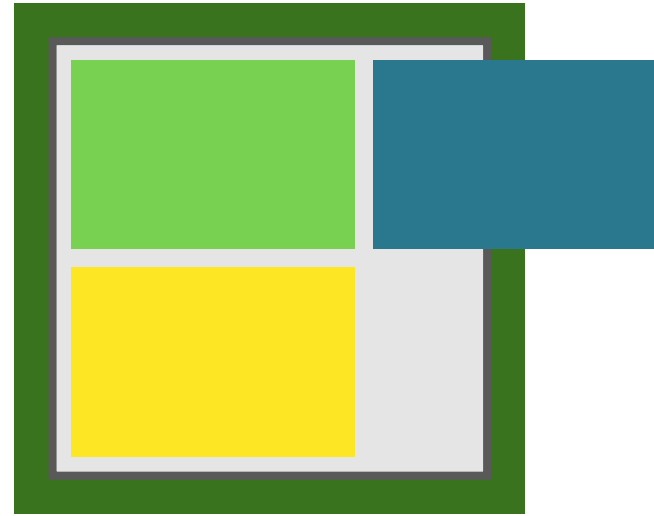
COMPRESSION



RESULTS



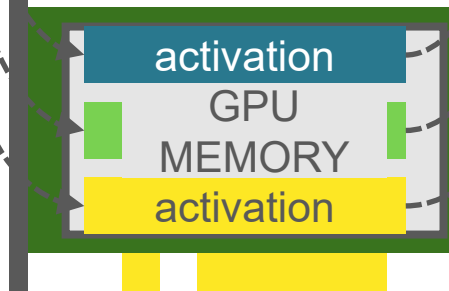
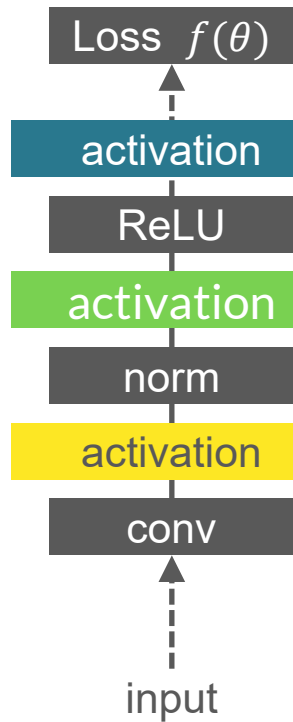
MOTIVATION



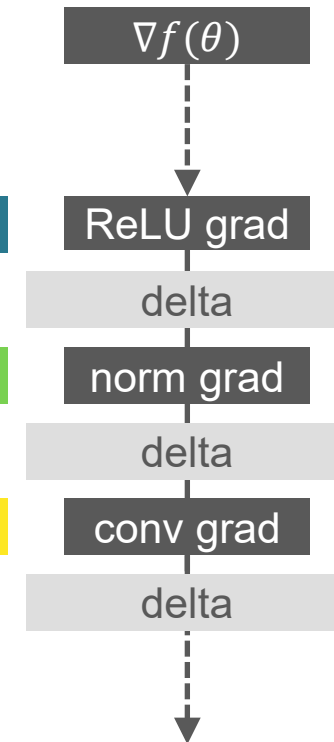
ACTIVATIONS DURING TRAINING

Activations are memoized layer outputs from the forward pass

Forward Pass



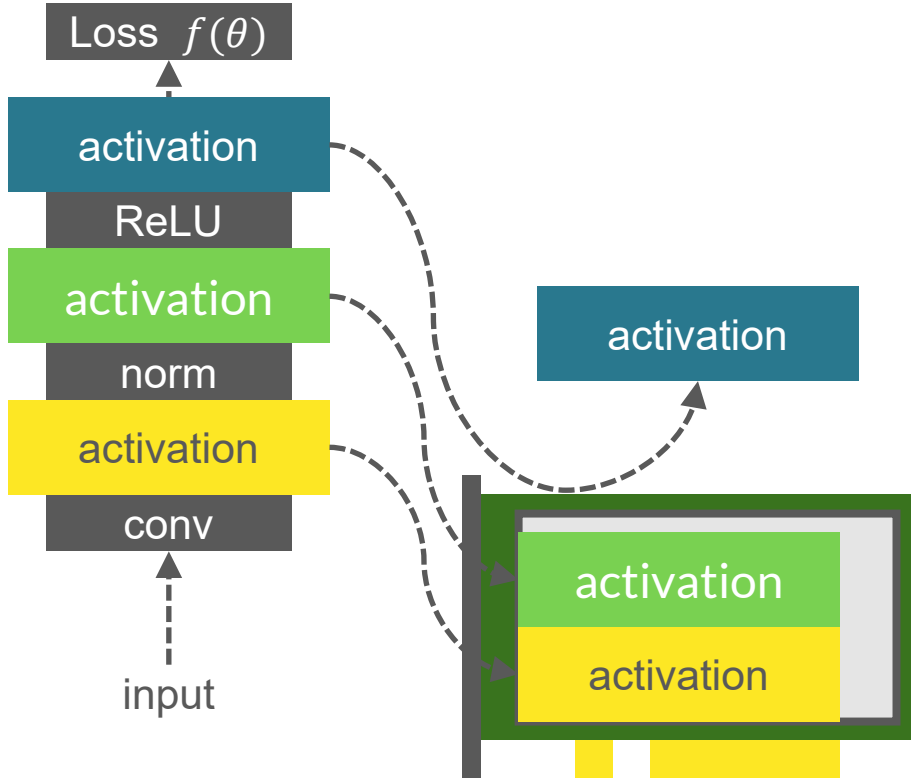
Backward Pass



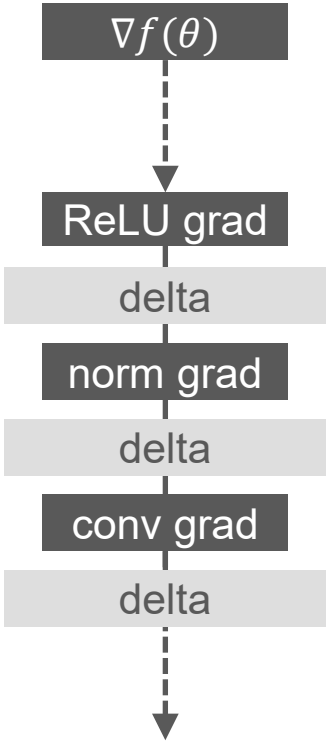
LOSSY ACTIVATION COMPRESSION

Larger networks and larger batch sizes mean more activations

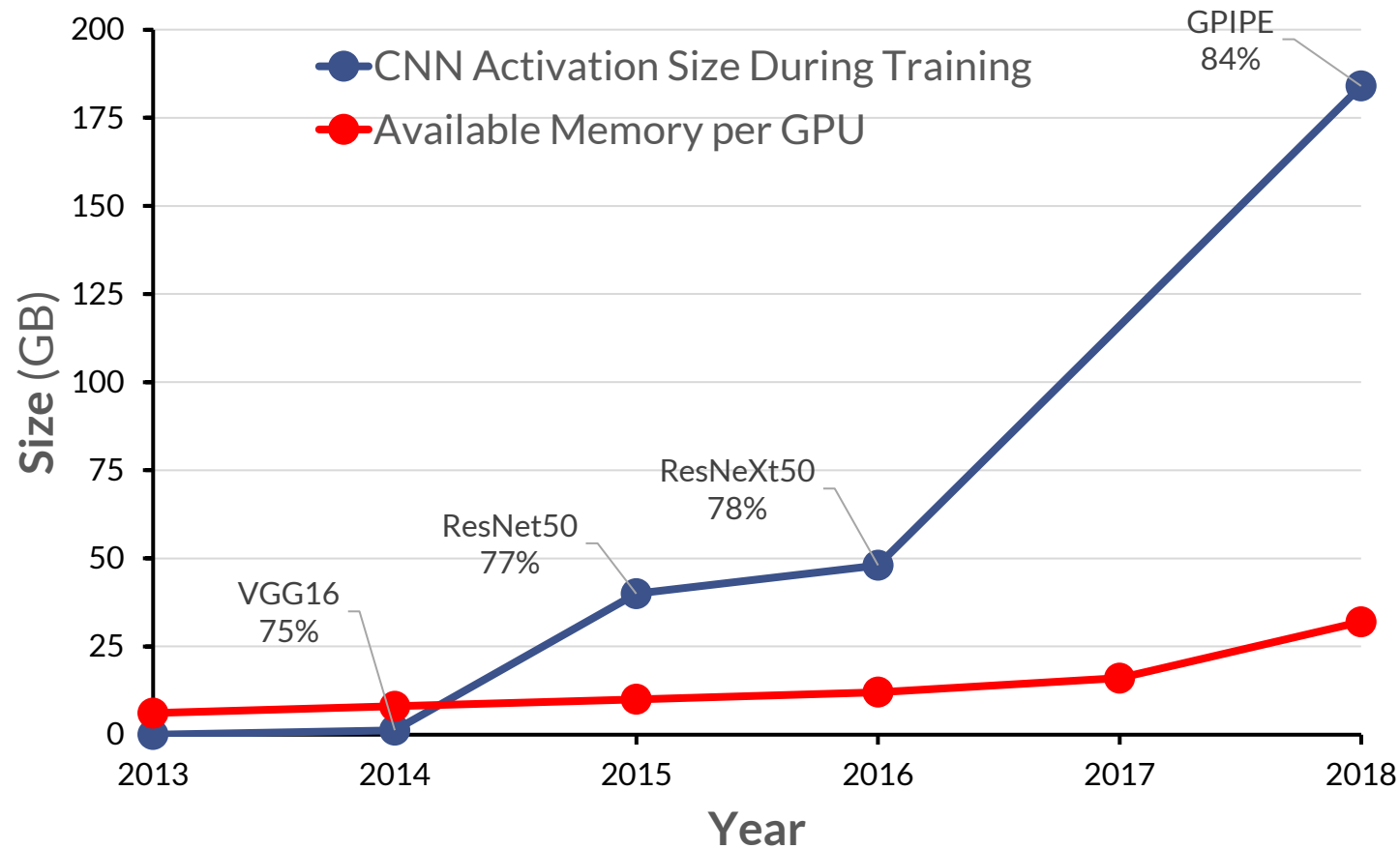
Forward Pass



Backward Pass



THE ACCELERATOR MEMORY WALL



Due to **physical limitations**, e.g. bandwidth and memory technology

DECREASING ACTIVATION MEMORY

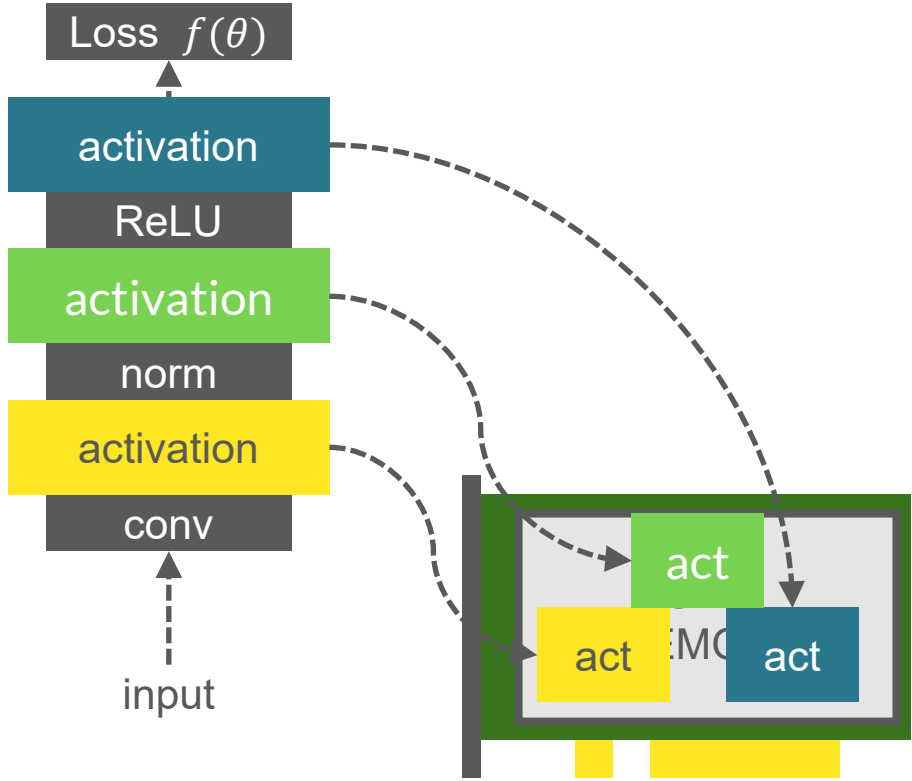
Many orthogonal approaches

Approach	Relative to normal training		
	Training Throughput	Memory Reduction	Accuracy after Training
Scheduling e.g. GIST (A Jain, ISCA 2018)	X	✓	
Recalculation e.g. Sublinear Nets (T. Chen, ArXiv 2016)		✓✓	
Offloading e.g. vDNN (M. Rhu, MICRO 2016)	XXXX	✓✓✓✓	
Lossless Compression e.g. CDMA (M. Rhu, HPCA 2018)	XX	✓✓	
Lossy Compression e.g. ACTNN (J. Chen, ICML 2021)	XX	✓✓✓✓✓	X

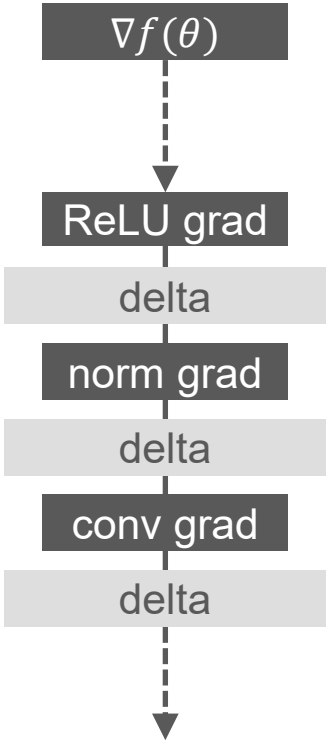
LOSSY ACTIVATION COMPRESSION

Discard some data, and compress

Forward Pass



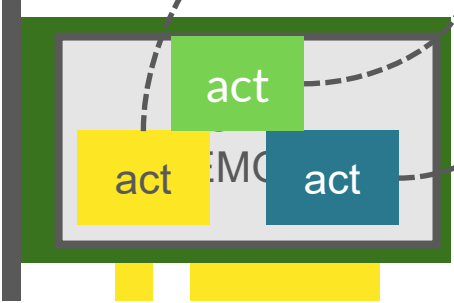
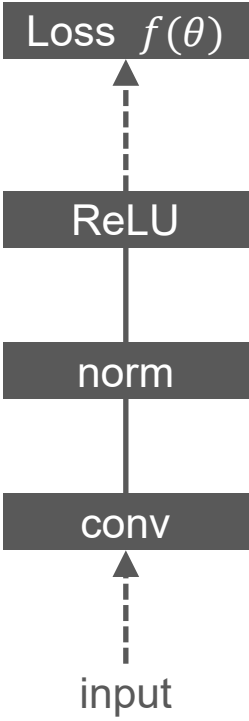
Backward Pass



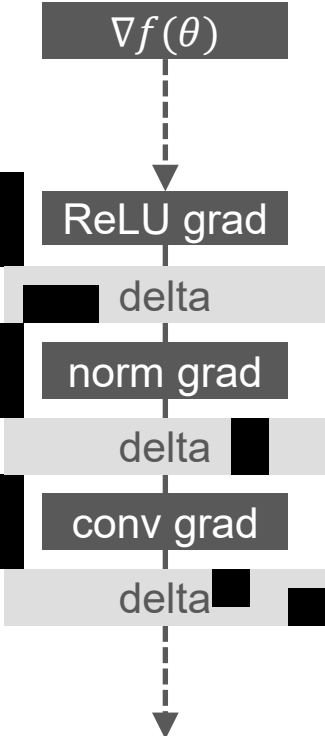
LOSSY ACTIVATION COMPRESSION

Discard some data, and compress

Forward Pass

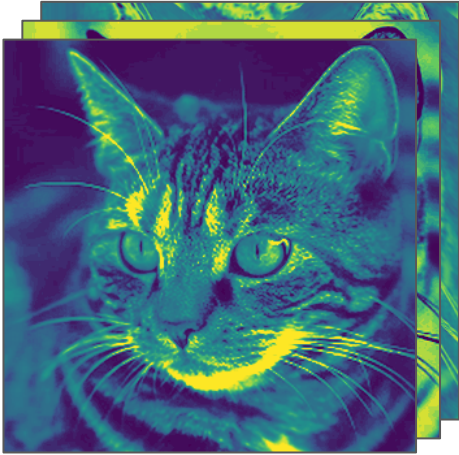


Backward Pass

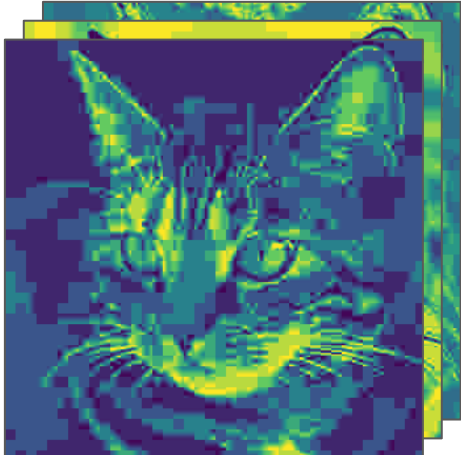


DRAWBACKS OF LOSSY COMPRESSION

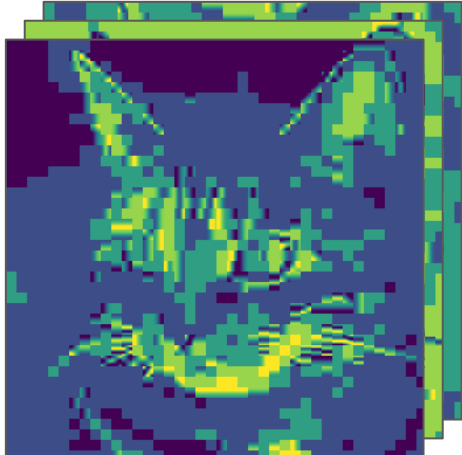
How to avoid explosions during training?



Low
Compression



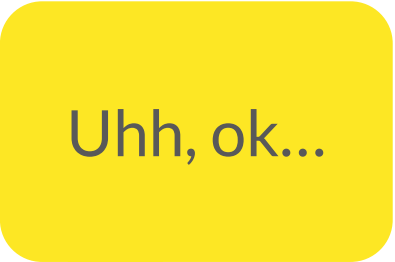
High
Compression



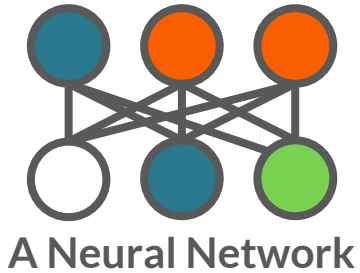
“This is a cat”



“This is a cat”

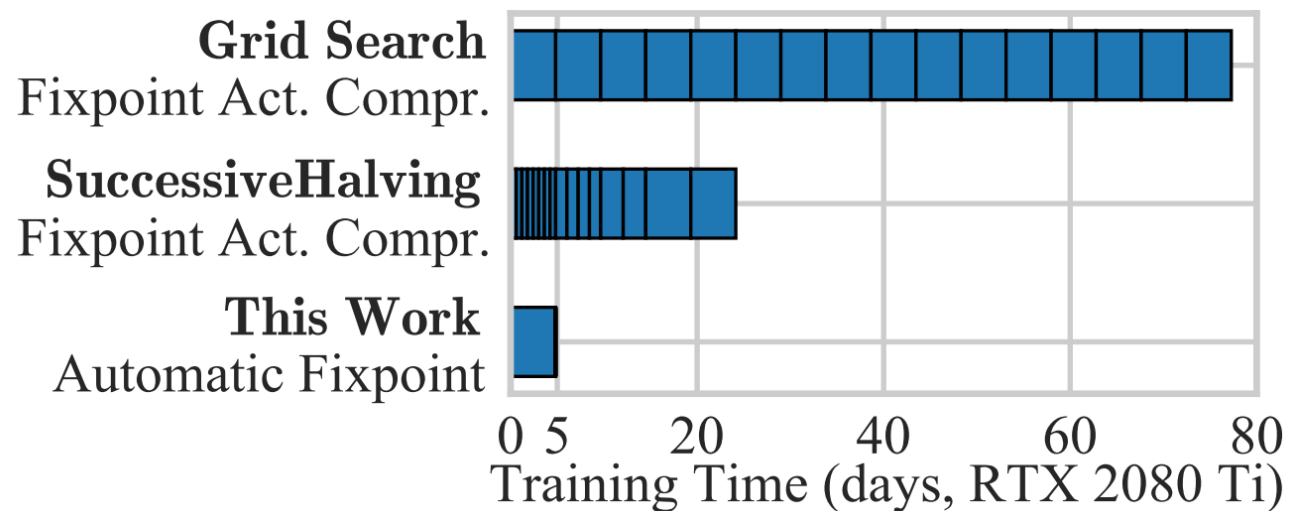


“This is a cat”



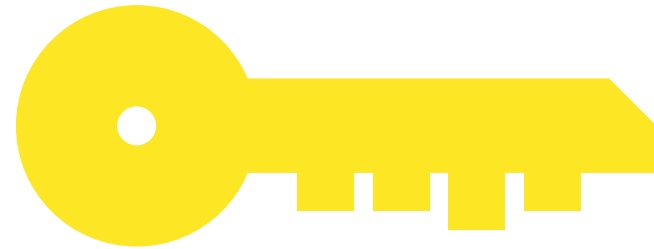
A Neural Network

TUNING THE COMPRESSION RATE IS HARD



Compression rate search cost. Each box indicates a different compression rate (1- to 16-bit fixpoint)

KEY CONTRIBUTIONS

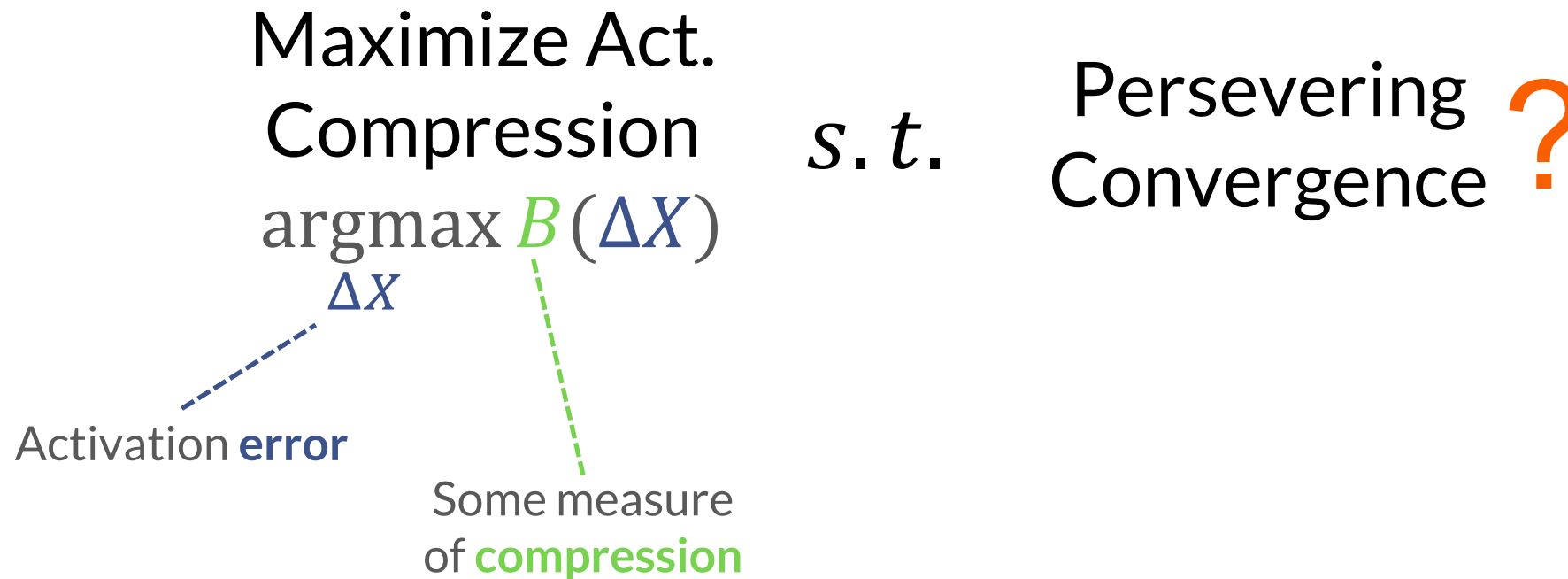


CONSTRAINED OPTIMIZATION FOR LOSSY ACTIVATIONS

Maximize Act.
Compression
 $\operatorname{argmax}_{\Delta X} B(\Delta X)$ *s. t.* Persevering
Convergence

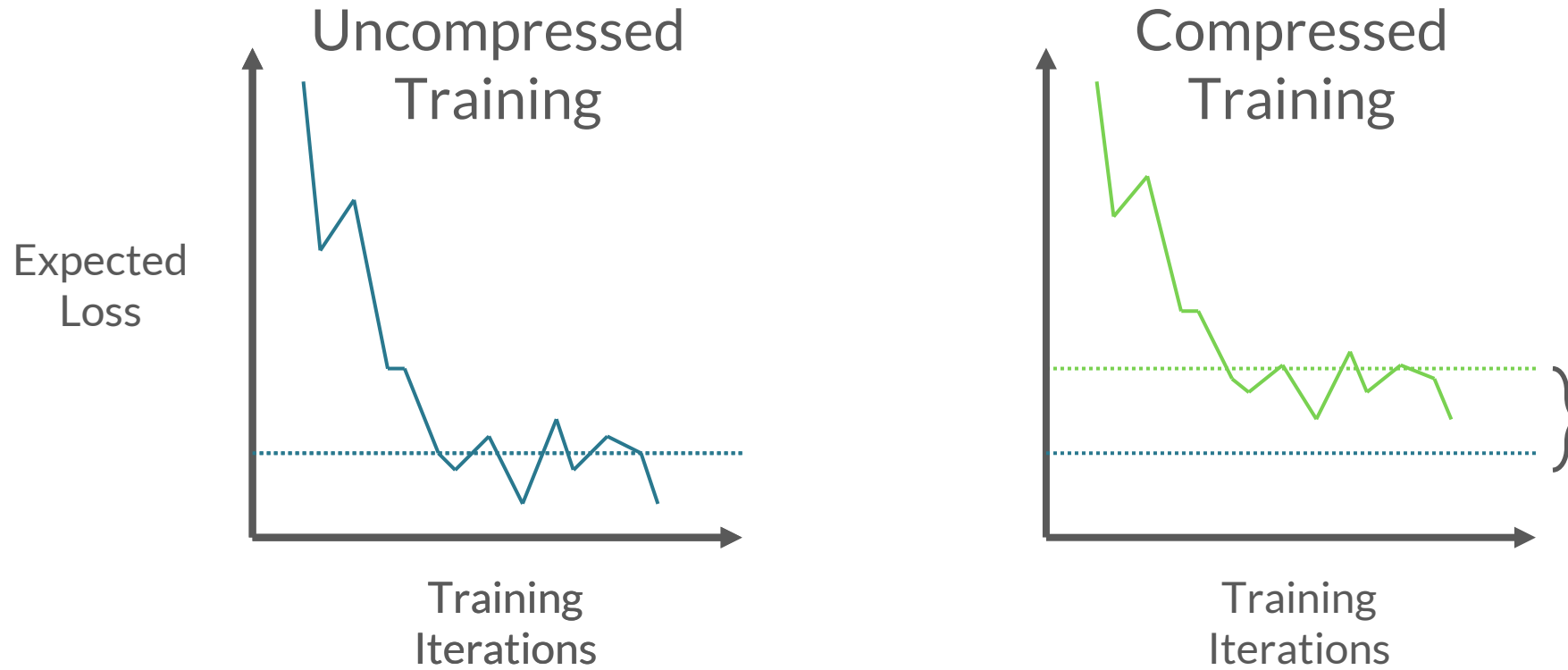
Goal: An efficient way to get compression rate from convergence

CONSTRAINED OPTIMIZATION FOR LOSSY ACTIVATIONS

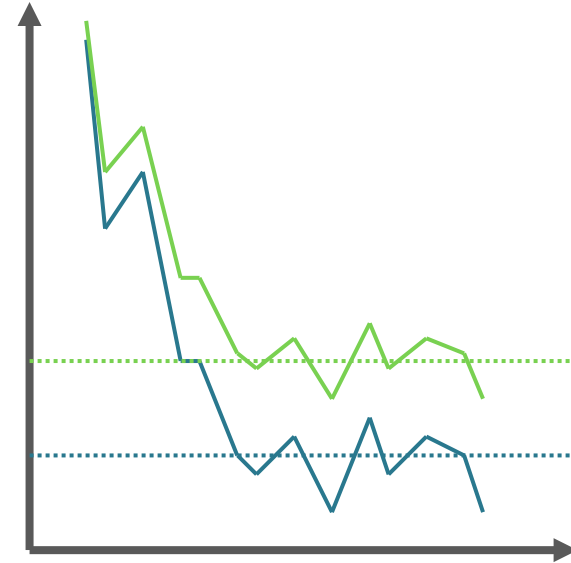


KEY INSIGHT: ALLOW LOSS TO INCREASE

Takes advantage of **SGD convergence** behaviour



BOUNDING CONVERGENCE



STOCHASTIC GRADIENT DESCENT NOTATION

Overall loss

Loss of example n

$$\mathcal{L}(\theta) = \sum_n f_n(\theta)$$

Parameters

Loss is a finite sum

Solve using S.G.D.

Gradient w.r.t. θ

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \nabla_{\theta} f_{n_t}(\theta)$$

Loss of a randomly selected training example

SGD THEORETICAL CONVERGENCE RATES

Many results on SGD convergence rates, we use:

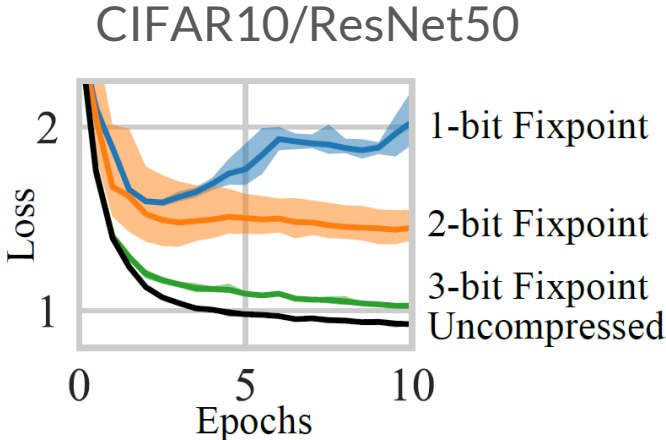
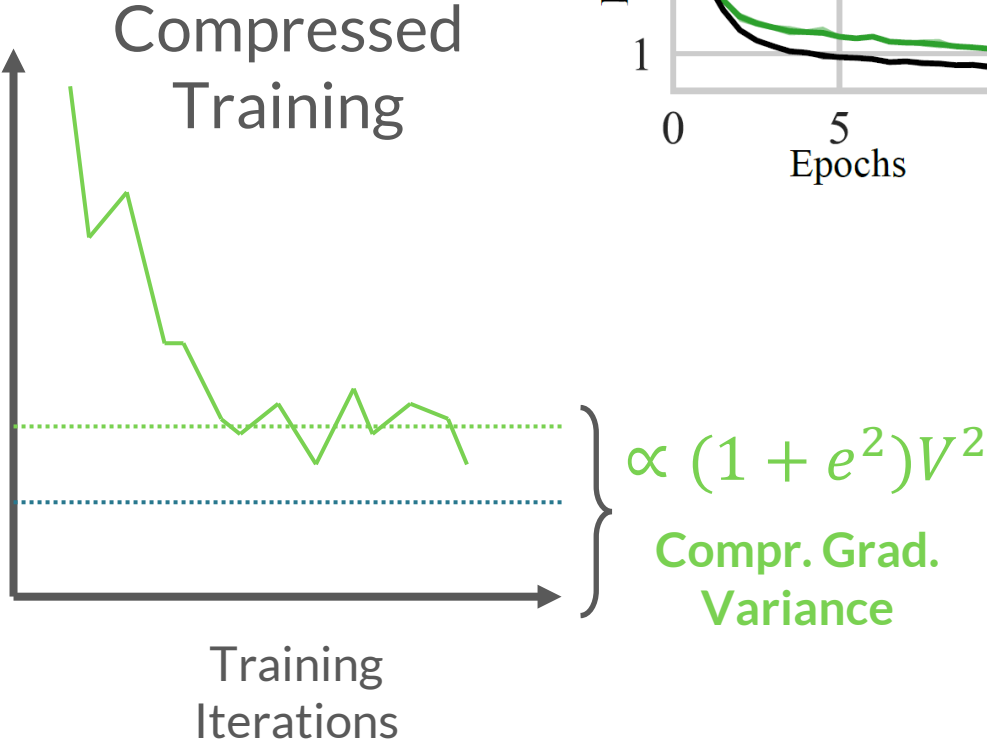
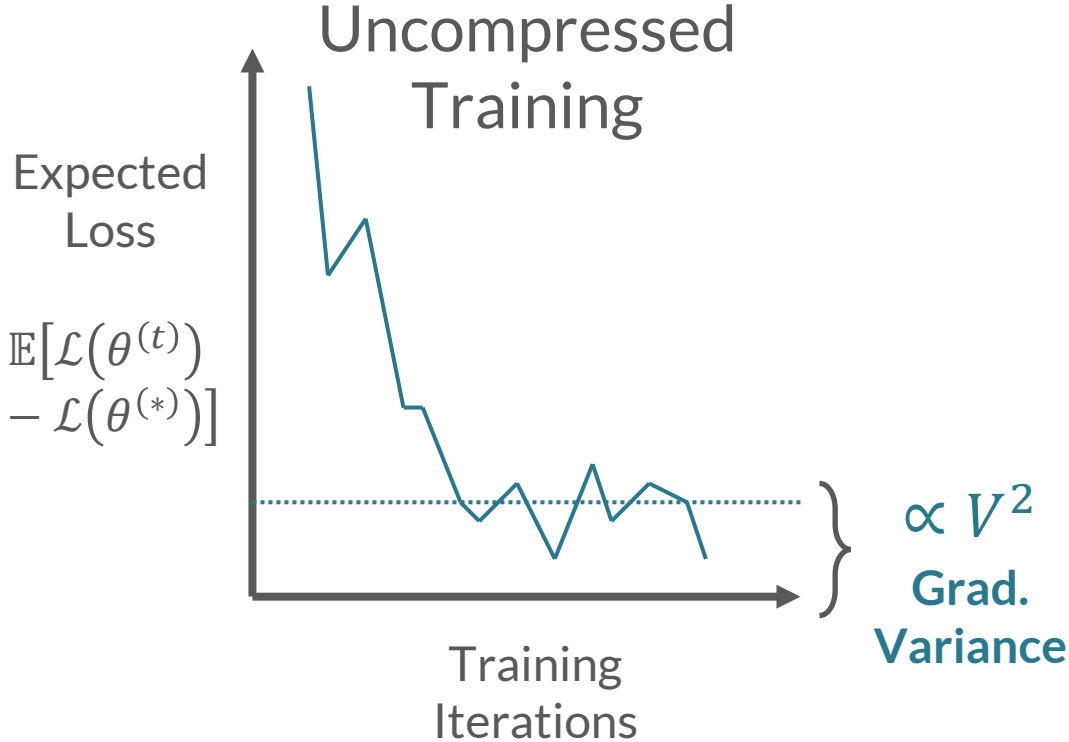
$$\mathbb{E}[\mathcal{L}(\theta^{(t)}) - \mathcal{L}(\theta^{(*)})] \leq (1 - C_1\alpha)^t (\mathcal{L}(\theta^{(0)}) - \mathcal{L}(\theta^{(*)})) + C_2\alpha V^2$$

C_1, C_2 : Constants

V^2 : Gradient Variance

[1] H Karimi, J Nutini, M Schmidt, “*Linear Convergence of Gradient and Proximal-Gradient Methods Under the Polyak-Łojasiewicz Condition*,”
ECML PKDD 2016

ALLOWING ERRORS FOR COMPRESSION



GRADIENT VARIANCE

Compressed training can be viewed as SGD with increased variance

From Karimi et. al:

Gradient “variance” is bounded:

$$\mathbb{E} \left[\left\| \nabla_{\theta} f_{n_t}(\theta) \right\|^2 \right] \leq V^2$$

Under compression:

$$\mathbb{E} \left[\left\| \hat{\nabla}_{\theta} f_{n_t}(\theta) \right\|^2 \right] \leq (1 + e^2)V^2$$

CONSTRAINED OPTIMIZATION FOR LOSSY ACTIVATIONS

Maximize Act.
Compression
 $\operatorname{argmax}_{\Delta X} B(\Delta X)$

s. t.

Persevering
Convergence

Gradient variance
is bounded

$$\mathbb{E} \left[\|\hat{\nabla}_{\theta} f\|^2 \right] \leq (1 + e^2)V^2$$

$$\mathbb{E} \left[\|\hat{\nabla}_{\theta} f_{n_t}(\theta)\|^2 \right] \leq (1 + e^2)V^2$$

ERROR-VARIANCE RELATIONSHIP

Using an additive gradient error:

$$\hat{\nabla}_{\theta} f \equiv \nabla_{\theta} f + \Delta \nabla_{\theta} f$$

The constraint on variance: $\xrightarrow{\text{(some math with expectations and norms)}}$

$$\mathbb{E}[\|\nabla_{\theta} f + \Delta \nabla_{\theta} f\|^2] \leq (1 + e^2)V^2$$

Is satisfied by:

$$\|\Delta \nabla_{\theta} f\|^2 \leq e^2 V^2$$

CONSTRAINED OPTIMIZATION FOR LOSSY ACTIVATIONS

Maximize Act. Compression
 $\operatorname{argmax}_{\Delta X} B(\Delta X)$

s. t.

Persevering Convergence

Gradient variance is bounded
 $\mathbb{E} \left[\|\hat{\nabla}_{\theta} f\|^2 \right] \leq (1 + e^2)V^2$

Gradient error is bounded
 $\|\Delta \nabla_{\theta} f\|^2 \leq e^2 V^2$

SIMPLIFYING THE PROBLEM



BOUNDING FUNCTION

Gradient error is related to activation error

Now we just need to calculate it for a layer, e.g., convolution...

$$\|\Delta \nabla_{\theta} f\|^2 = \sum_{k,c,r,s}^{K,C,R,S} \left(\sum_{n,h,w}^{N,H,W} \Delta x_{n,c,h+r,w+s} \frac{\partial f}{\partial y_{nkhw}} \right)^2 \leq e^2 V^2$$

which is not very useful... the solution is not closed-form. Instead use:

$$\|\Delta \nabla_{\theta} f\|^2 \leq D(\Delta X) \leq e^2 V^2 \longrightarrow \text{Where we find } D(\Delta X) \text{ that is as close as possible to the error norm.}$$

CONSTRAINED OPTIMIZATION FOR LOSSY ACTIVATIONS

Maximize Act.
Compression
 $\operatorname{argmax}_{\Delta X} B(\Delta X)$

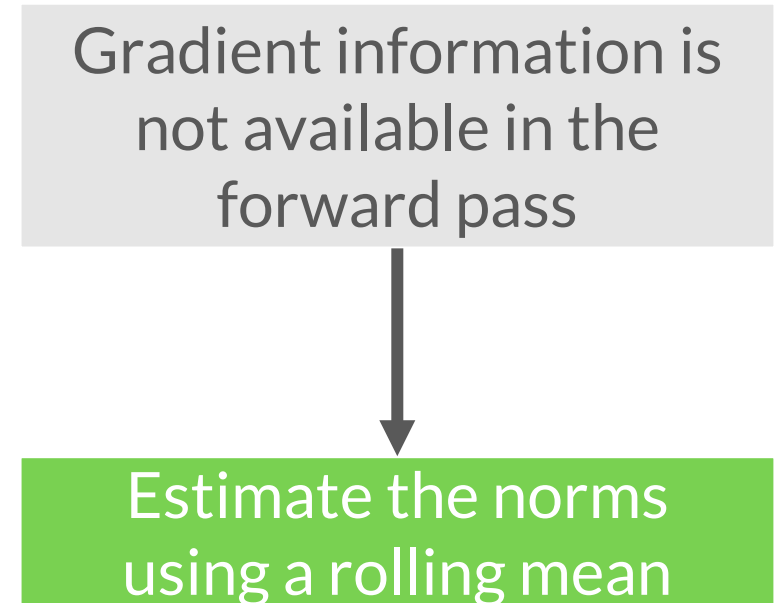
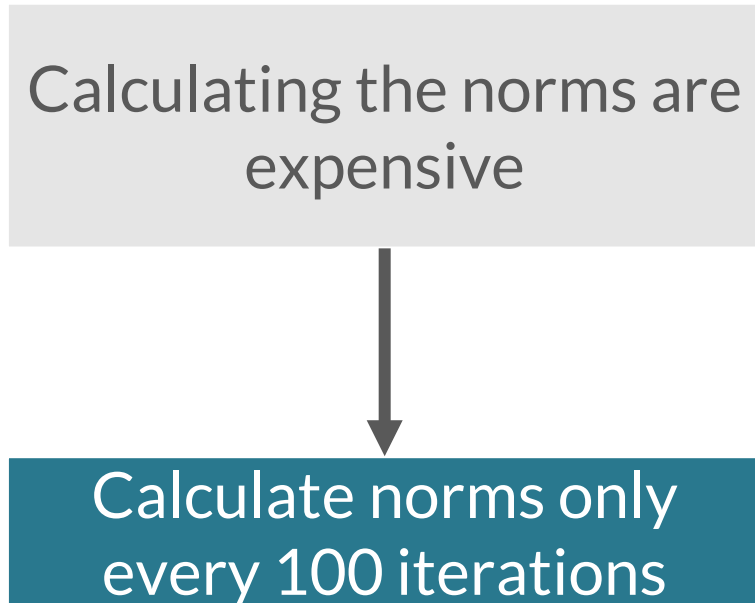
s. t. Persevering
Convergence
 $D(\Delta X) \leq e^2 V^2$

Gradient variance
is bounded
 $\mathbb{E} \left[\|\hat{\nabla}_{\theta} f\|^2 \right]$
 $\leq (1 + e^2) V^2$

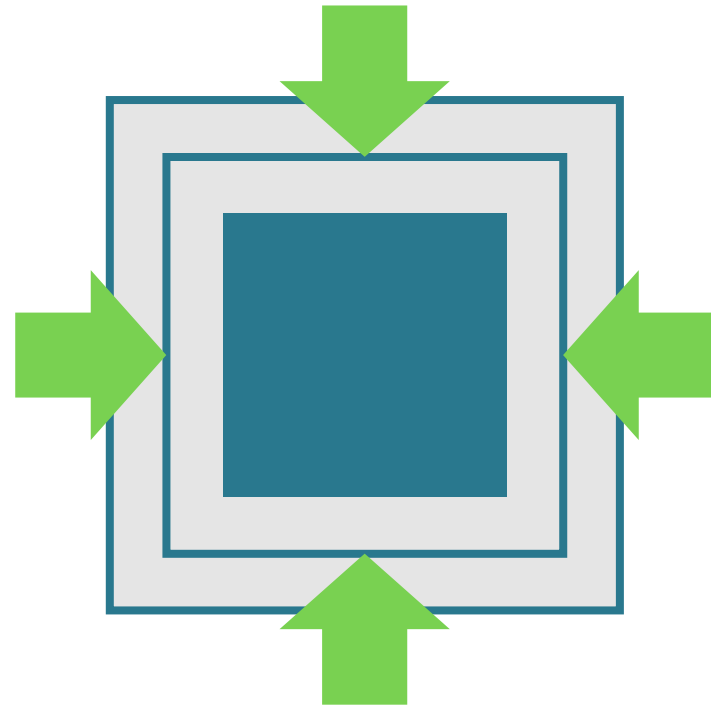
Gradient error
is bounded
 $\|\Delta \nabla_{\theta} f\|^2 \leq D(\Delta X)$

APPROXIMATING NORMS

Two Issues with calculating activation errors this way:



COMPRESSION



COMPRESSION METRIC

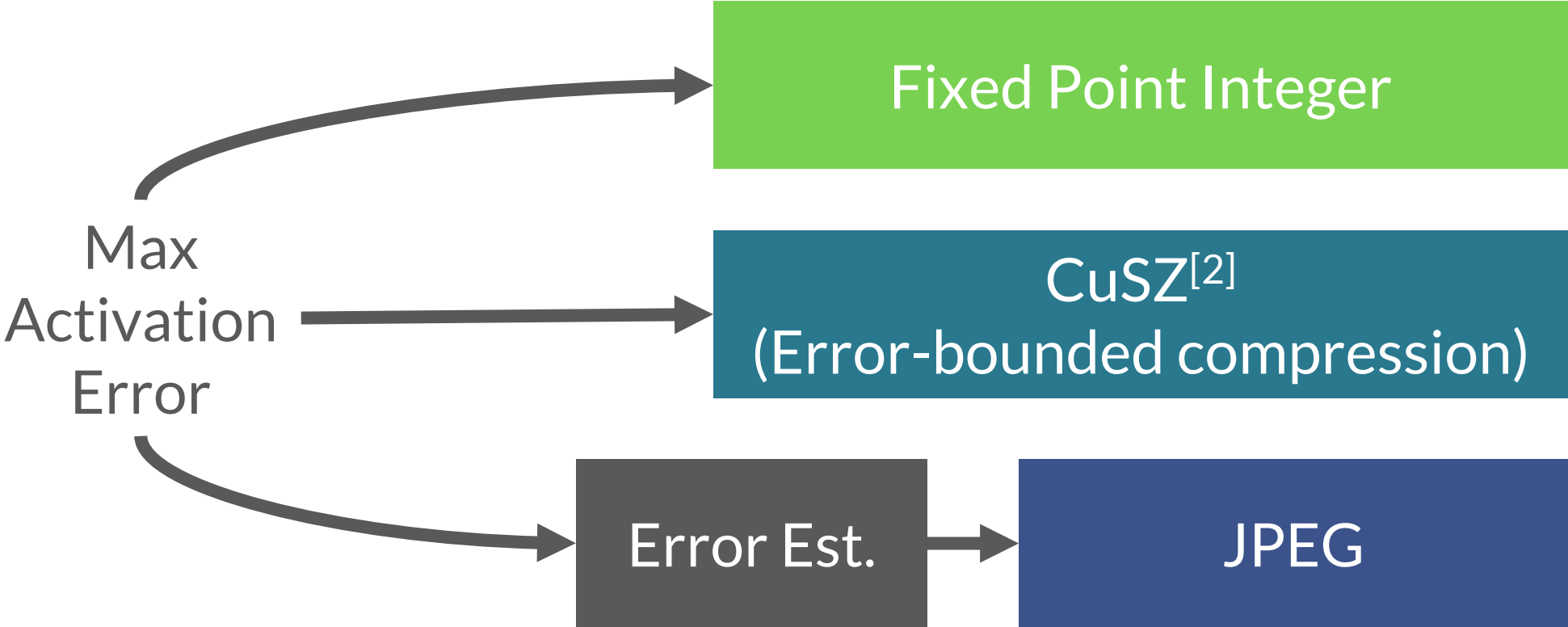
Choose B given that most methods use precision reduction

$B(\Delta X) \propto$ *Number of bits removed*

$$B(\Delta X) \equiv \sum_i \log|\Delta x_i|$$

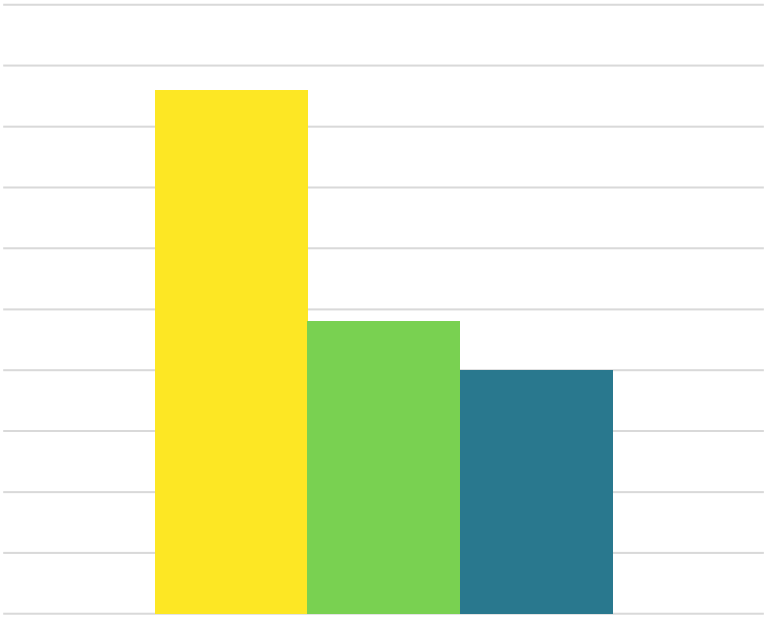
COMPRESSION METHODS

Case studies on activation error-bounded compression methods

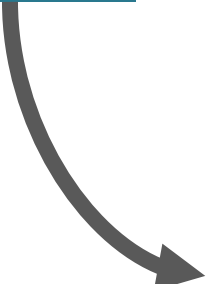


[2] S. Jin, G. Li, S. L. Song, D. Tao, "A Novel Memory-Efficient Deep Learning Training Framework via Error-Bounded Lossy Compression", in ArXiv 2020

RESULTS



ERROR BOUNDS ARE DERIVED PER-LAYER

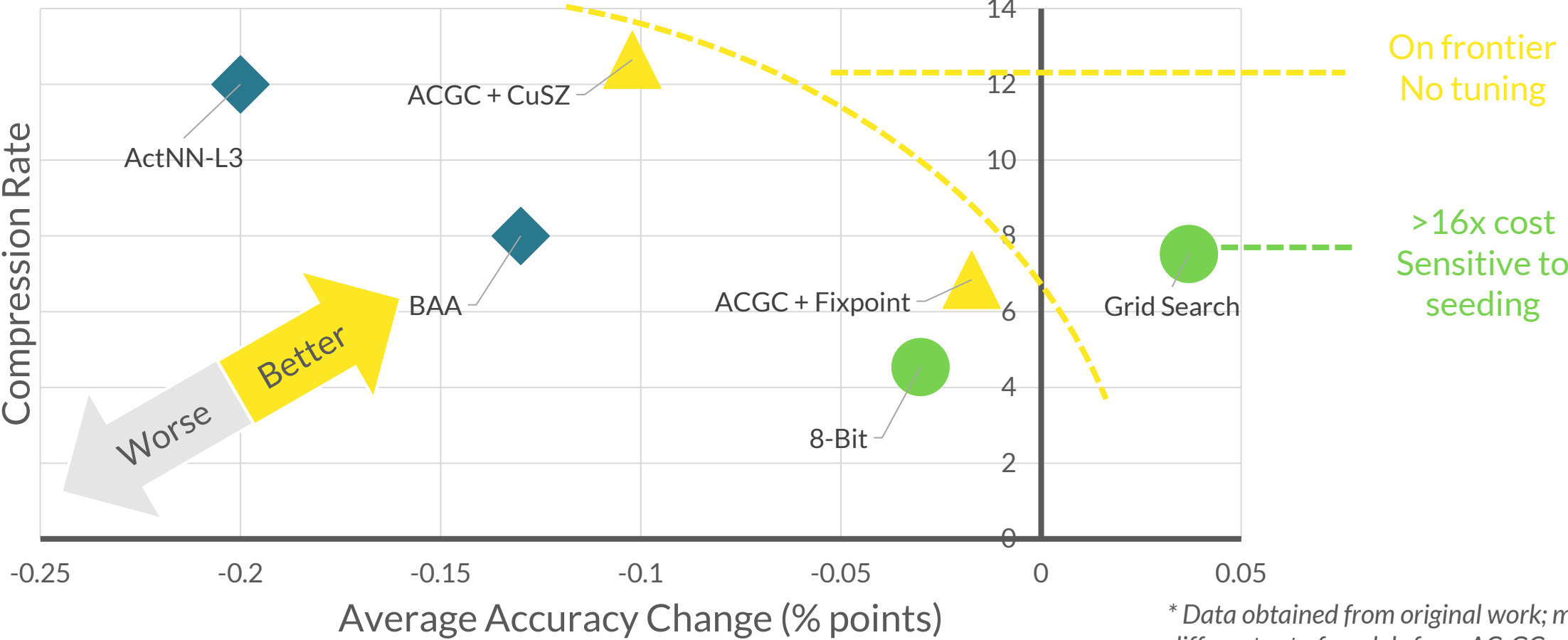


$$(\Delta x_{nchw})^2 \leq e^2 V^2 \frac{1}{RSNCHW \|\nabla_Y f\|^2}$$

Filter size

Activation dimensions

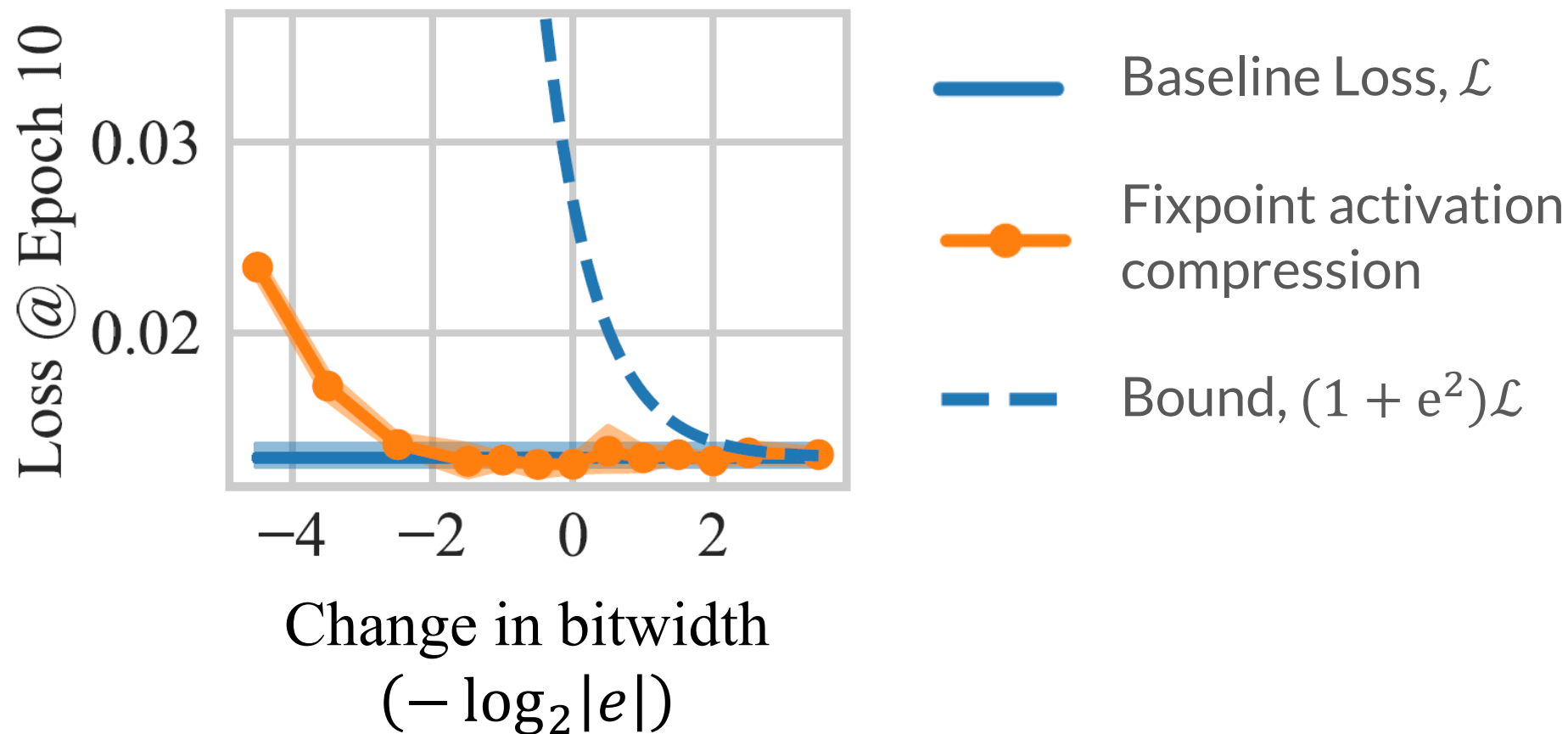
ACCURACY AND COMPRESSION



* Data obtained from original work; may use a different set of models from AC-GC

(BAA) A. Chakrabarti, B. Moseley, in NeurIPS 2020
 (ACTNN) J. Chen, L. Zheng, et. al, in ICML 2021

THEORETICAL VERSUS EMPIRICAL (MNIST)



Bounds are empirically satisfied

For more information, come to our poster!

AC-GC: Lossy Activation Compression with Guaranteed Convergence

R. David Evans, Tor M. Aamodt

Neural Information Processing Systems (NeurIPS), Dec 7, 2021

Poster Spot A2, 1630 PST – 1800 PST

THE UNIVERSITY OF BRITISH COLUMBIA

THANK YOU FOR LISTENING



THE UNIVERSITY
OF BRITISH COLUMBIA

Electrical and Computer
Engineering



