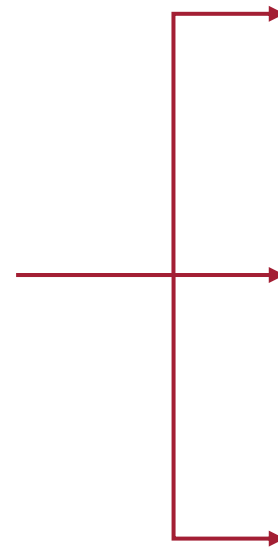
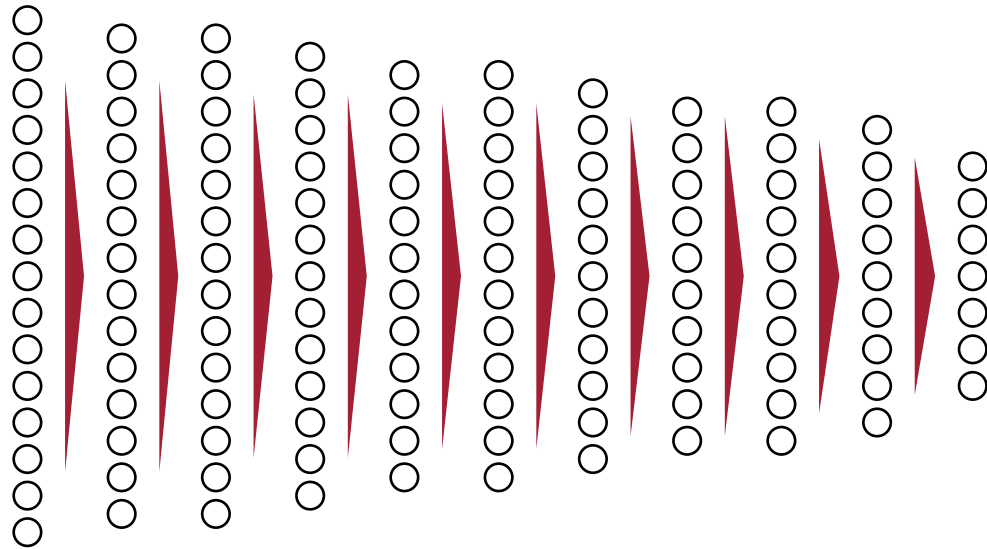


# Compressing Neural Networks: Towards Determining the Optimal Layer-wise Decomposition

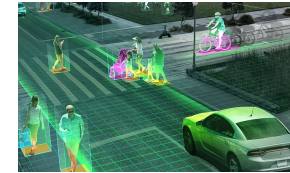
Lucas Liebenwein\*, Alaa Maalouf\*, Dan Feldman, Daniela Rus

\* Equal contribution

# Neural networks are SOTA



Natural Language Processing

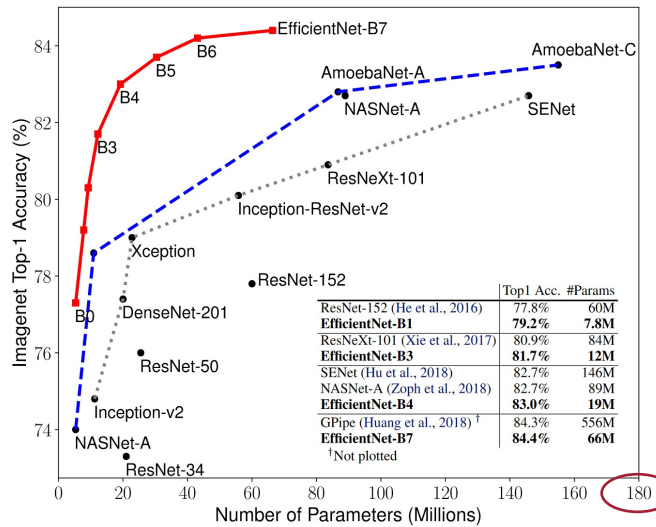


Computer Vision



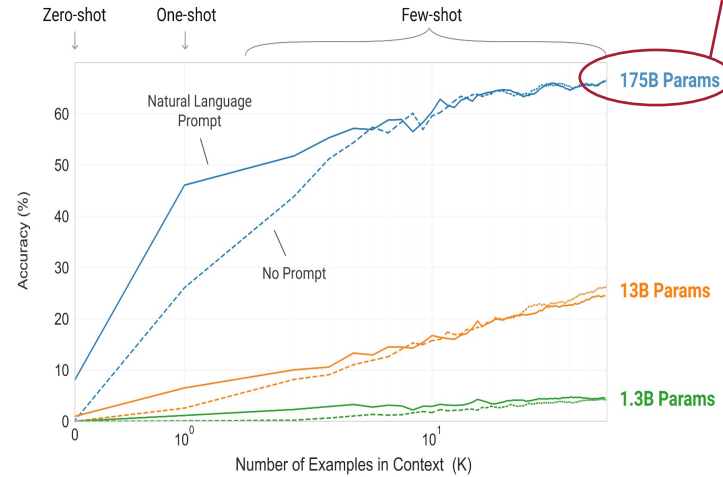
Robotics

# The bigger, the better



Tan, Mingxing, and Quoc Le. "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks." ICML. 2019.

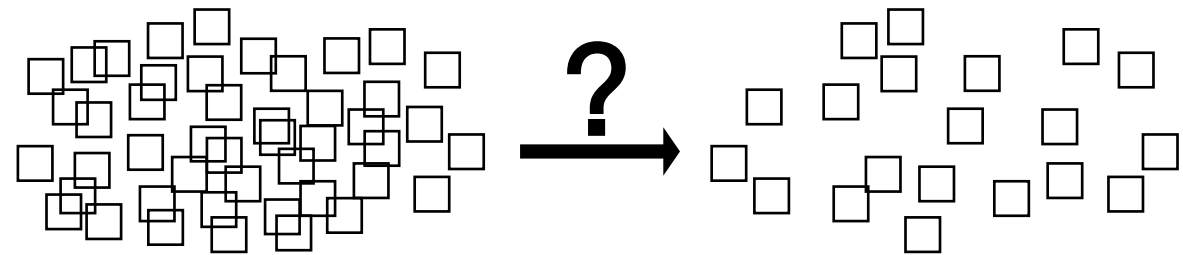
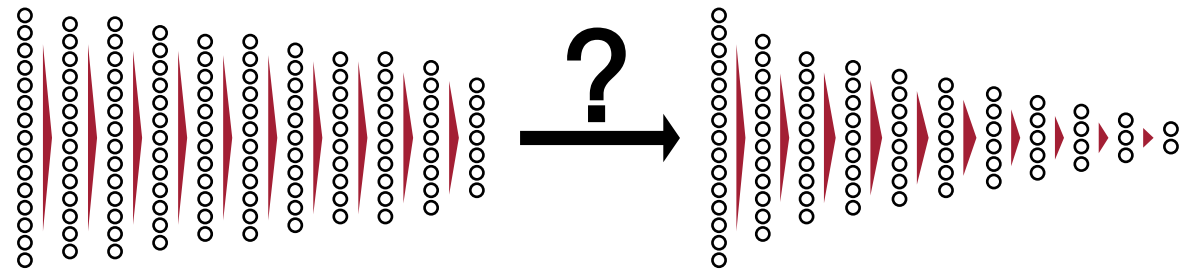
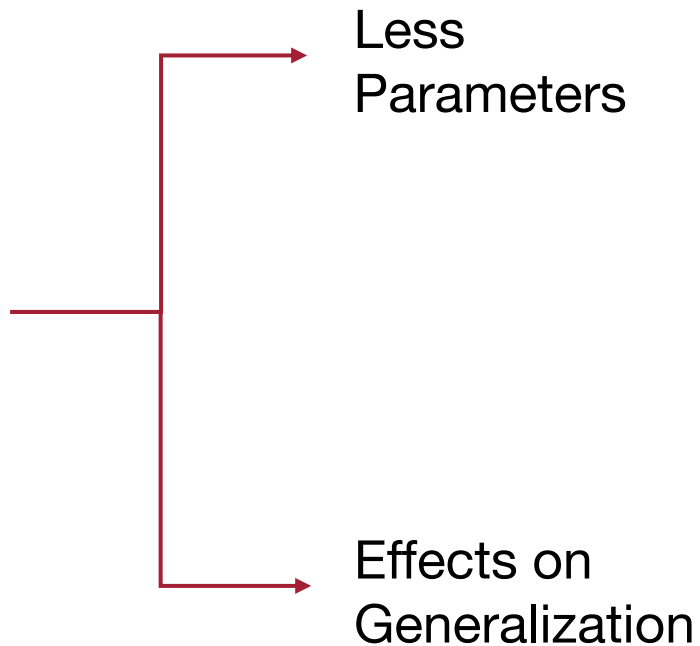
180  
m



Brown, Tom, et al. "Language models are few-shot learners." arXiv preprint arXiv:2005.14165 (2020).

175B

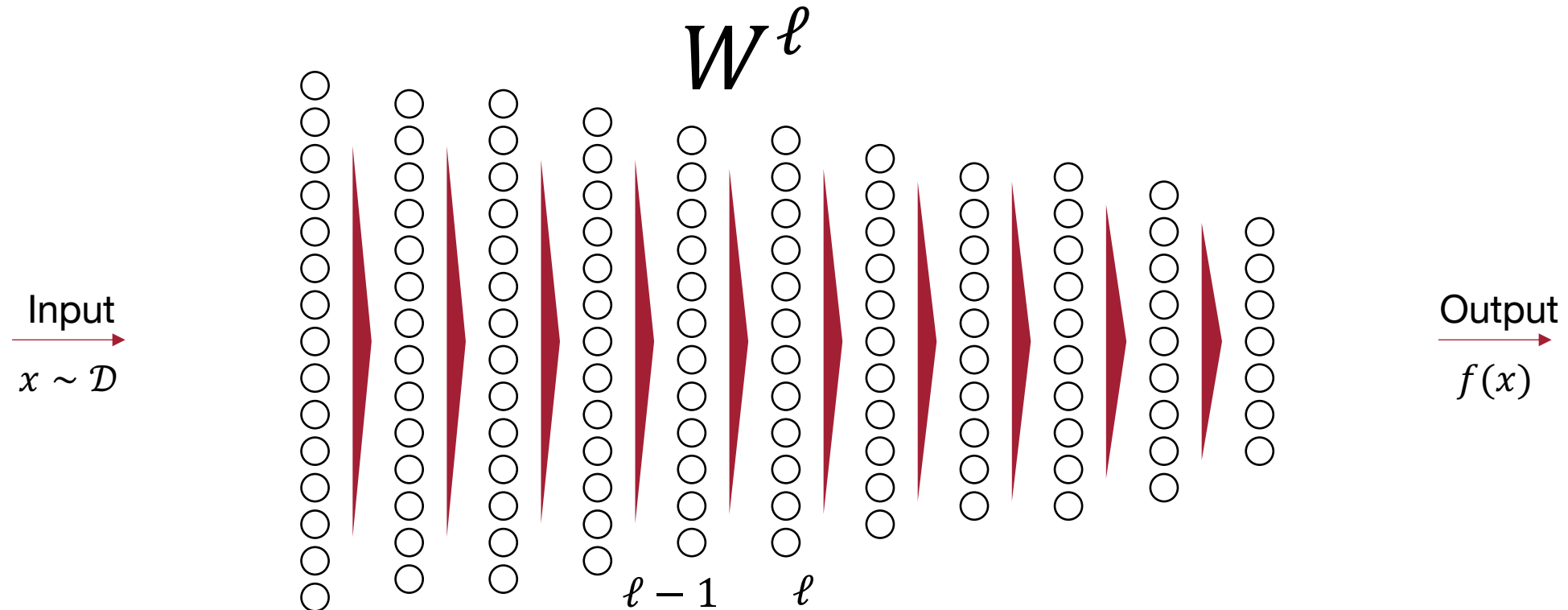
# Less resources, same performance?



# Outline of the talk

- 1 Introduce neural network compression via low-rank decomposition
- 2 Understand limitations of prior work and our contribution
- 3 Present our main algorithm **ALDS (Automatic Layer-wise Decomposition Selector)**
- 4 Discuss results, future work, and discussion

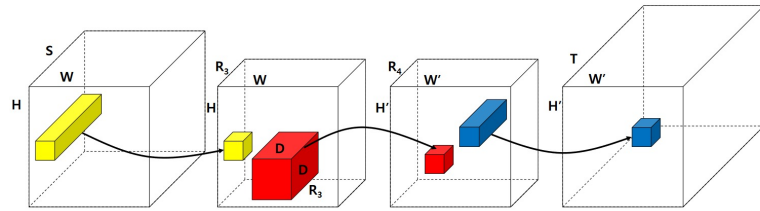
# Compression via low-rank decomposition



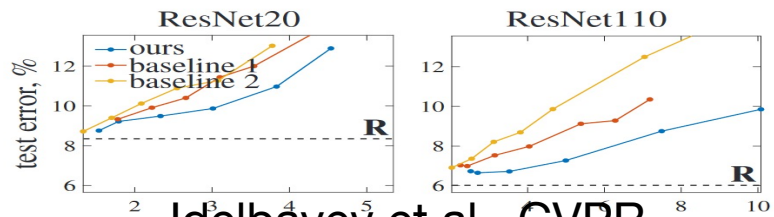
# Compression via low-rank decomposition

$$\begin{array}{ccc}
 & f \times c & f \times j \quad j \times c \\
 W^\ell = & \left( \begin{array}{|c|c|c|c|c|c|c|} \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \hline \\ \hline \end{array} \right) & \approx & \left( \begin{array}{|c|c|} \hline \hline \hline \hline \hline \hline \hline \hline \\ \hline \end{array} \right) * \left( \begin{array}{|c|c|c|c|c|c|c|} \hline \hline \hline \hline \hline \hline \hline \\ \hline \end{array} \right) =: \hat{U}^\ell * \hat{V}^\ell \\
 \# \text{ params} = fc & \longrightarrow & \# \text{ params} = j(f + c)
 \end{array}$$

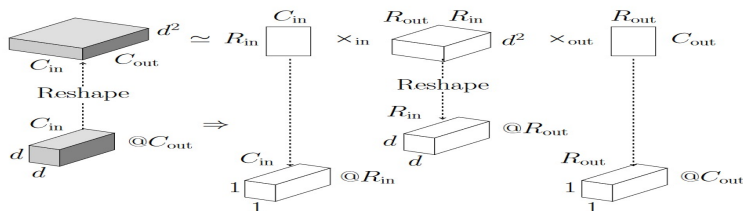
# Related work



Kim et al., ICLR 2016



Idelbayev et al., CVPR 2020



Gusak et al., ICCV 2019

✓ Per-layer low-rank decomposition:  
“local step”

✓ Progressive decomposition + training:  
“retraining strategy”

✓ Tuning of per-layer decomposition:  
“global step”

✗ Unresolved: combined local + global step  
for optimal network compression

**Our main  
contribution**



# A general approach to low-rank compression

1

“Local step”

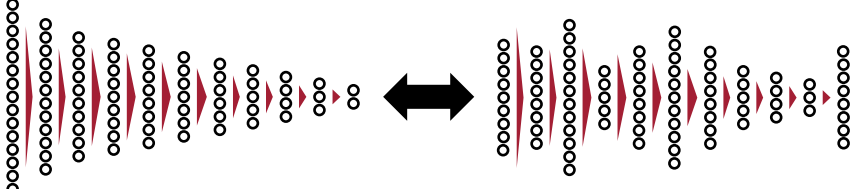
$$\begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \approx \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} * \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix}$$

$$\varepsilon^\ell = \varepsilon^\ell(\text{compression\_ratio})$$

Efficiently implementable and easy to evaluate

2

“Global step”



$$\|f - \hat{f}\| \leq \varepsilon \|f\| \text{ where } \varepsilon = \varepsilon(\varepsilon^1, \dots, \varepsilon^L)$$

$$\text{minimize } cost(\varepsilon^1, \dots, \varepsilon^L) \text{ s.t. } size(\hat{\theta}) \leq \mathcal{B}$$

# Local step: per-layer low-rank compression

$$W^\ell \stackrel{f \times c}{=} \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \approx \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \stackrel{f \times j}{*} \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \stackrel{j \times c}{=} \hat{U}^\ell * \hat{V}^\ell$$

# params =  $fc$        $\longrightarrow$       # params =  $j(f + c)$

$$\varepsilon^\ell = \varepsilon^\ell(j) = \varepsilon^\ell(\text{"compression_ratio"})$$

# Local step: per-layer low-rank compression

$$\begin{array}{ccc}
 & f \times \frac{c}{k} k & \\
 & f \times j & j \times c \\
 W^\ell = \left( \begin{array}{c|c} \text{red grid} & \text{gray grid} \end{array} \right) & \approx & \left( \begin{array}{c} \text{grid} \end{array} \right) * \left( \begin{array}{c} \text{grid} \end{array} \right) =: \hat{U}^\ell * \hat{V}^\ell \\
 \# \text{ params} = fc & \longrightarrow & \# \text{ params} = j(f + c)
 \end{array}$$

$$\varepsilon^\ell = \varepsilon^\ell(j) = \varepsilon^\ell(\text{"compression\_ratio"})$$

# Local step: per-layer low-rank compression

$$\begin{aligned}
 & \begin{matrix} f \times \frac{c}{k} k \\ W^\ell = \left( \begin{array}{cccc|cccc} \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{red} & \text{red} & \text{red} & \text{red} & \text{gray} & \text{gray} & \text{gray} & \text{gray} \end{array} \right) \end{matrix} & \approx & \begin{matrix} f \times jk \\ \left( \begin{array}{cccc|cccc} \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{gray} & \text{gray} & \text{red} & \text{red} & \text{red} & \text{red} \end{array} \right) \end{matrix} & * & \begin{matrix} k * \left( j \times \frac{c}{k} \right) \\ \begin{array}{cccc} \text{red} & \text{red} & \text{red} & \text{red} \\ \text{red} & \text{red} & \text{red} & \text{red} \end{array} \\ \\ \\ \\ \\ \\ \\ \\ \begin{array}{cccc} \text{gray} & \text{gray} & \text{gray} & \text{gray} \\ \text{gray} & \text{gray} & \text{gray} & \text{gray} \end{array} \end{matrix} & =: & \hat{U}^\ell * \hat{V}^\ell
 \end{aligned}$$

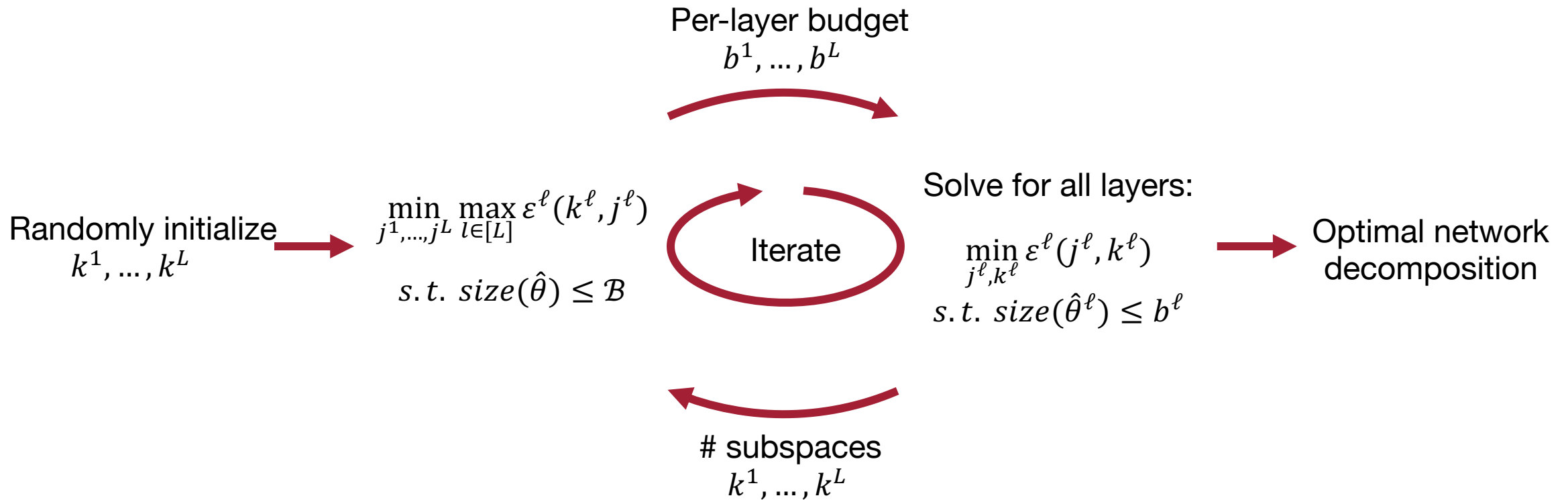
$\# \text{ params} = fc \quad \longrightarrow \quad \# \text{ params} = j(f + c)$

$$\varepsilon^\ell = \varepsilon^\ell(j) = \varepsilon^\ell(\text{"compression_ratio"})$$

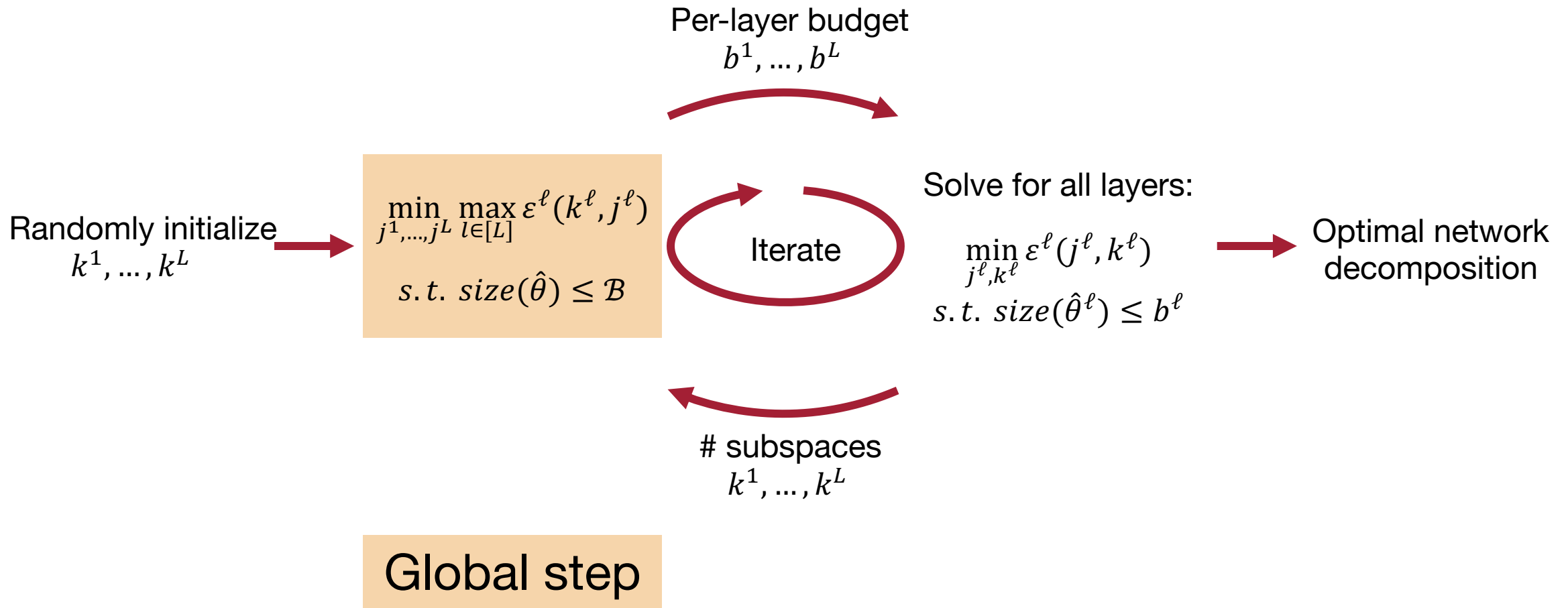




# ALDS: Automatic Layer-wise Decomposition Selector

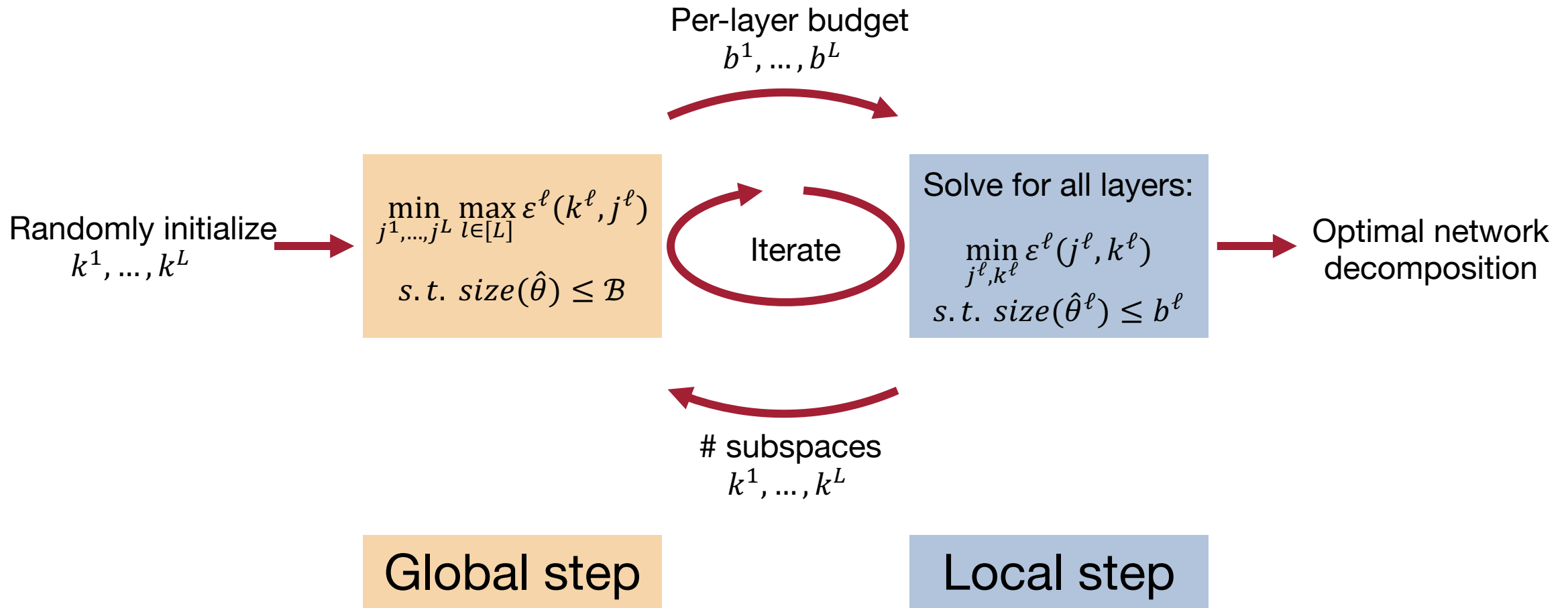


# ALDS: Automatic Layer-wise Decomposition Selector





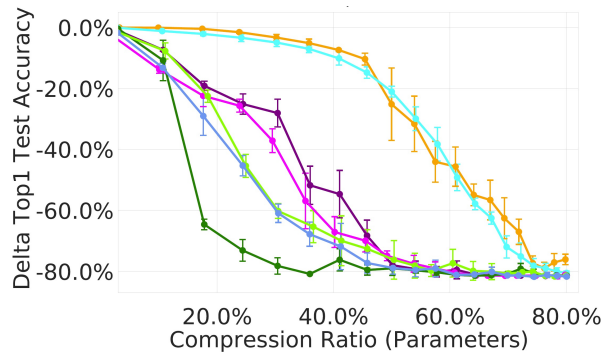
# ALDS: Automatic Layer-wise Decomposition Selector



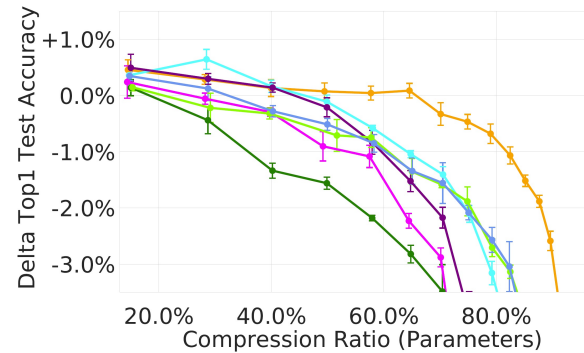
# Results: one-shot (+ retraining) with baselines

ResNet20  
CIFAR10

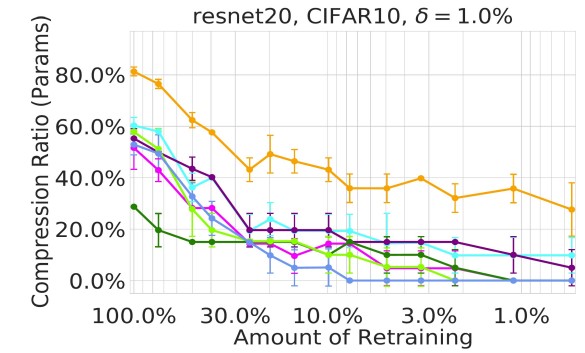
## Compress-only



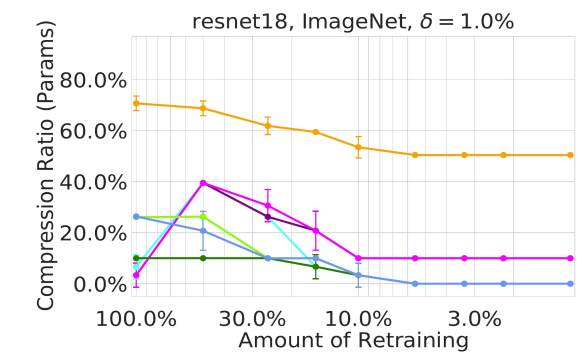
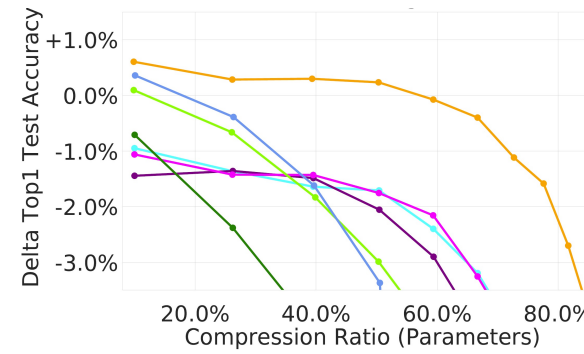
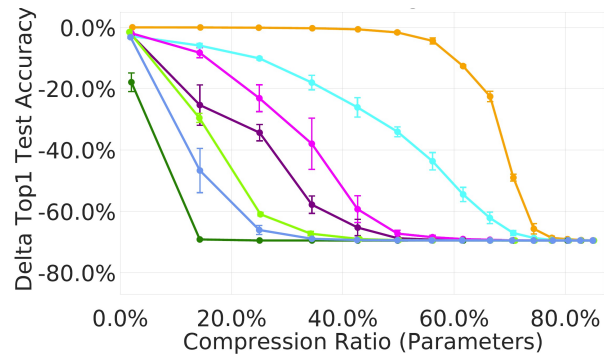
## Retrain



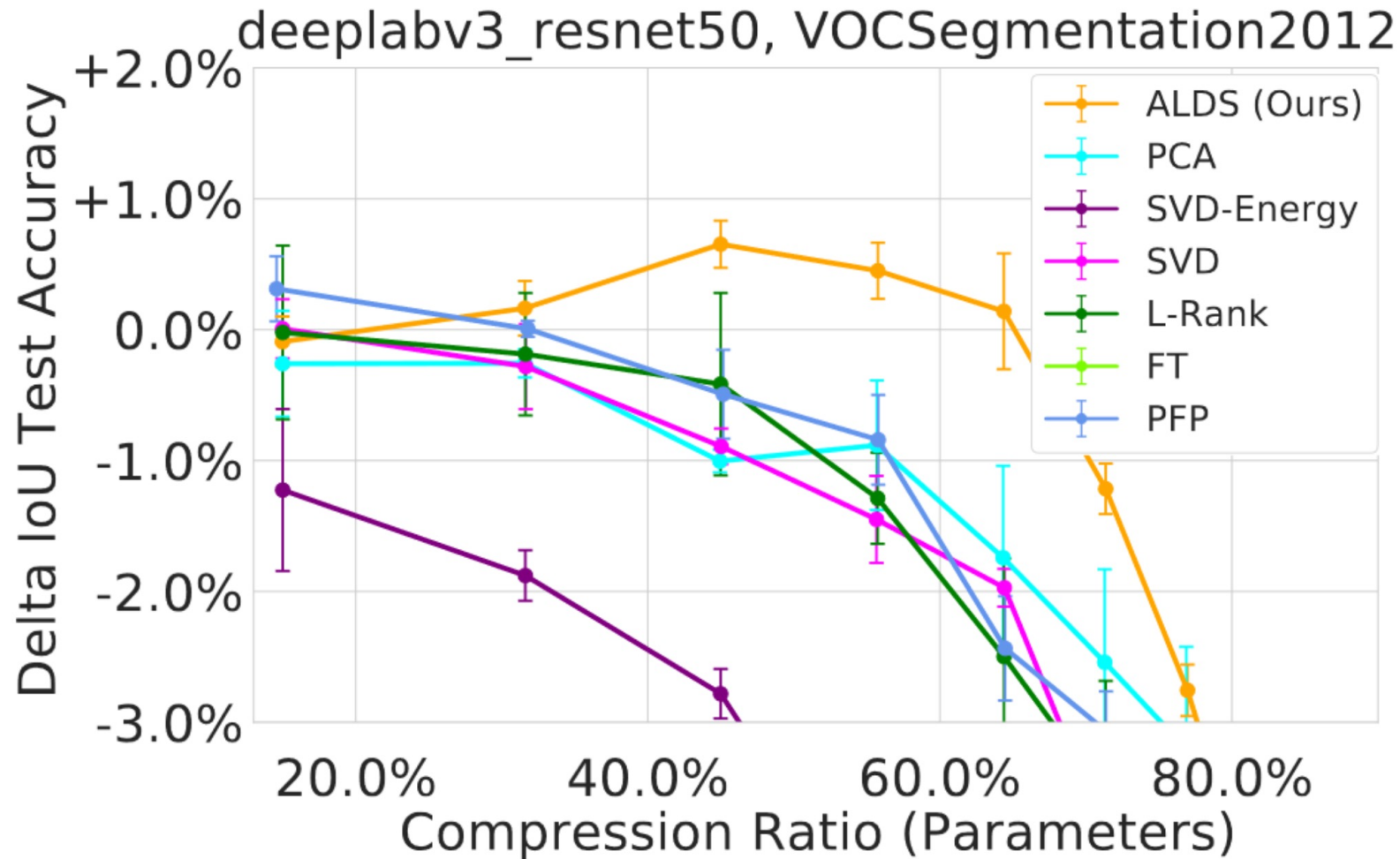
## Sweep



ResNet18  
ImageNet



# Results: one-shot (+ retraining) with baselines



# Results: ImageNet benchmarks

	Method	$\Delta$ -Top1	$\Delta$ -Top5	CR-F (%)
<b>ResNet18, Top1, 5: 69.64%, 88.98%</b>	<b>ALDS (Ours)</b>	<b>-0.38</b>	<b>+0.04</b>	<b>64.5</b>
	ALDS (Ours)	-1.37	-0.56	76.3
	<b>MUSCO (Gusak et al., 2019)</b>	<b>-0.37</b>	<b>-0.20</b>	<b>58.67</b>
	TRP1 (Xu et al., 2020)	-4.18	-2.5	44.70
	TRP1+Nu (Xu et al., 2020)	-4.25	-2.61	55.15
	<b>TRP2+Nu (Xu et al., 2020)</b>	<b>-4.3</b>	<b>-2.37</b>	<b>68.55</b>
	PCA (Zhang et al., 2015b)	-6.54	-4.54	29.07
	Expand (Jaderberg et al., 2014)	-6.84	-5.26	50.00
	PFP (Liebenwein et al., 2020)	-2.26	-1.07	29.30
	SoftNet (He et al., 2018)	-2.54	-1.2	41.80
	Median (He et al., 2019)	-1.23	-0.5	41.80
	Slimming (Liu et al., 2017)	-1.77	-1.19	28.05
	Low-cost (Dong et al., 2017)	-3.55	-2.2	34.64
	Gating (Hua et al., 2018)	-1.52	-0.93	37.88
	FT (He et al., 2017)	-3.08	-1.75	41.86
	DCP (Zhuang et al., 2018)	-2.19	-1.28	47.08
	FBS (Gao et al., 2018)	-2.44	-1.36	49.49

	Method	$\Delta$ -Top1	$\Delta$ -Top5	CR-F (%)
<b>AlexNet, Top1, 5: 57.30%, 80.20%</b>	ALDS (Ours)	-0.21	-0.36	77.9
	<b>ALDS (Ours)</b>	<b>-0.41</b>	<b>-0.54</b>	<b>81.4</b>
	Tucker (Kim et al., 2015a)	N/A	-1.87	62.40
	<b>Regularize (Tai et al., 2015)</b>	<b>N/A</b>	<b>-0.54</b>	<b>74.35</b>
	Coordinate (Wen et al., 2017)	N/A	-0.34	62.82
	Efficient (Kim et al., 2019)	-0.7	-0.3	62.40
	L-Rank (Idelbayev et al., 2020)	-0.13	-0.13	66.77
	NISP (Yu et al., 2018)	-1.43	N/A	67.94
	OICSR (Li et al., 2019a)	-0.47	N/A	53.70
	Oracle (Ding et al., 2019)	-1.13	-0.67	31.97

# Discussion and future work



ALDS leads to novel **state-of-the-art results** in low-rank compression



Combining the **local + global** decomposition step leads to a more flexible approach



**Error bounds** lead to better global insights about compression



**Future:** ALDS as a modular compression framework for any per-layer compression

# Compressing Neural Networks: Towards Determining the Optimal Layer-wise Decomposition

Lucas Liebenwein\*, Alaa Maalouf\*, Dan Feldman, Daniela Rus

\* Equal contribution

## Thank you



Alaa Maalouf



Dan Feldman



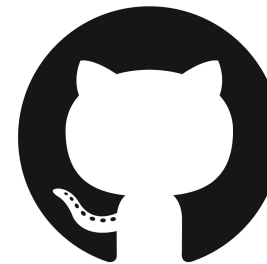
Daniela Rus

## Acknowledgements



Oren Gal

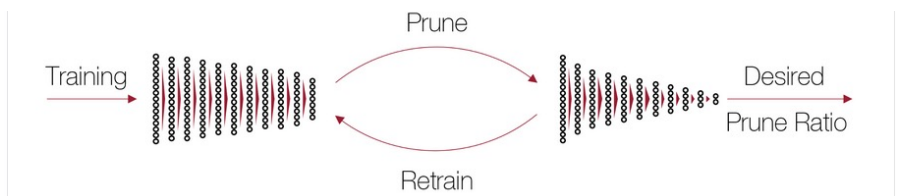
## Code



About



A research library for pytorch-based neural network pruning, compression, and more.



<https://github.com/lucaslie/torchprune>