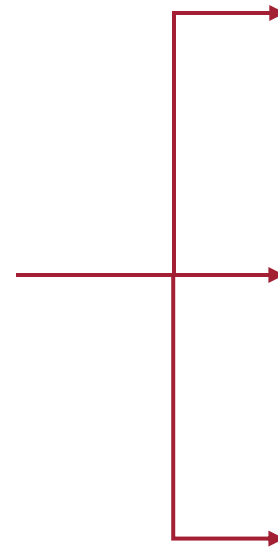
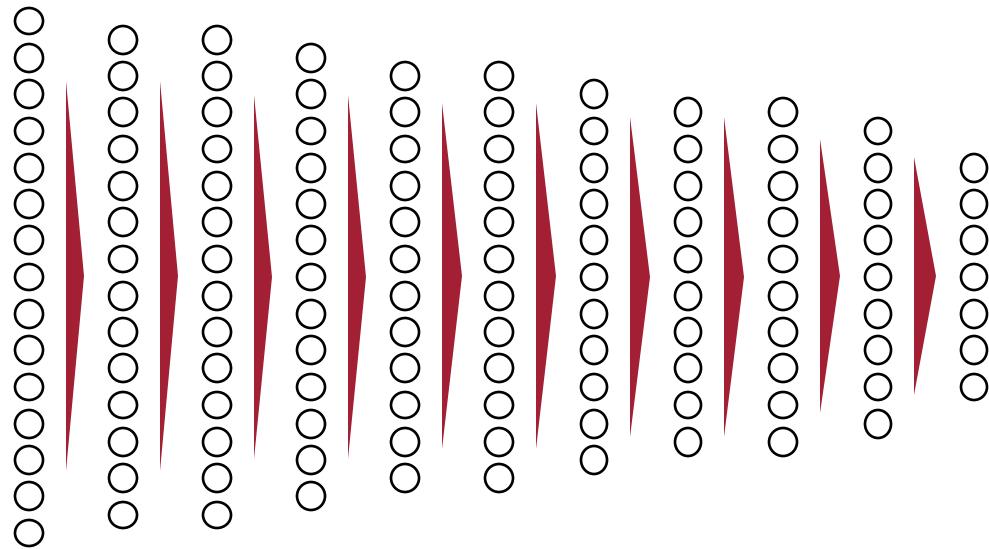


# Provable Filter Pruning for Efficient Neural Networks

Lucas Liebenwein\*, Cenk Baykal\*, Harry Lang, Dan Feldman, Daniela Rus  
Distributed Robotics Lab, CSAIL, MIT

# Neural networks are SOTA



Natural Language Processing

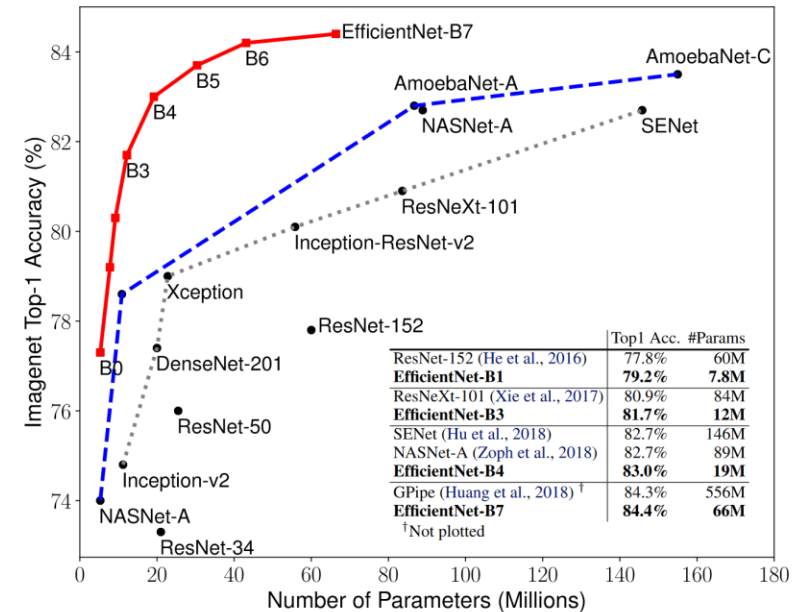
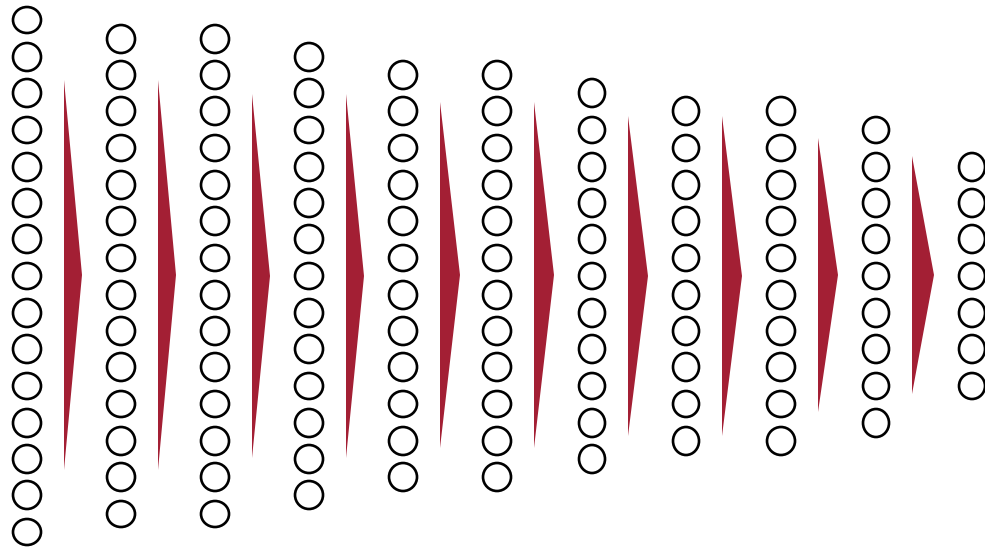


Computer Vision



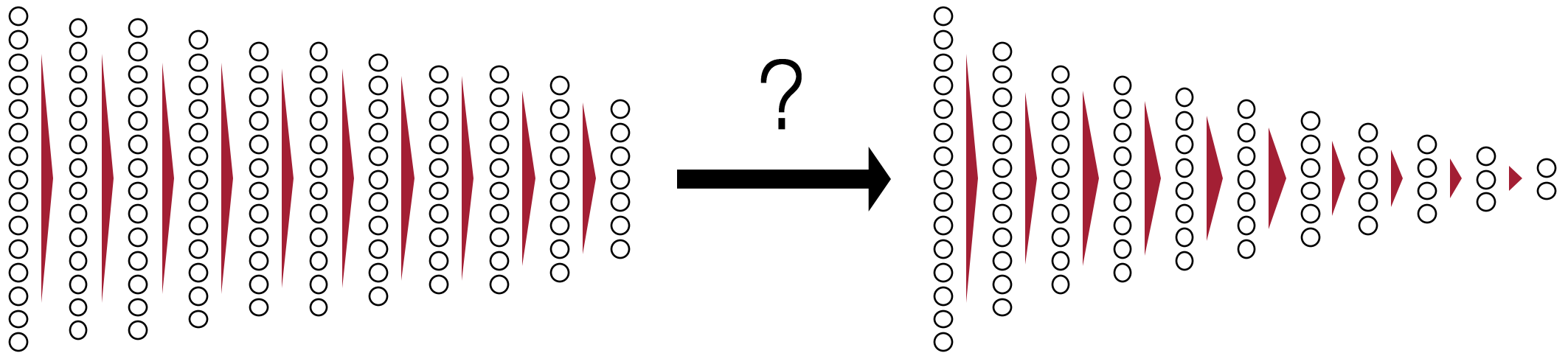
Robotics

# Larger size, better performance

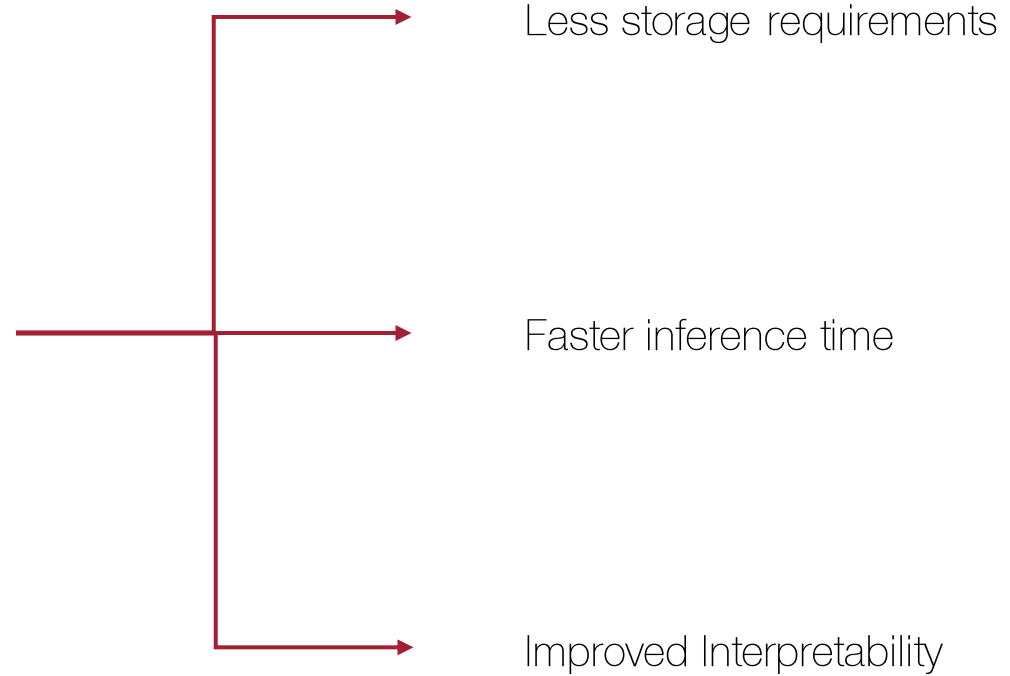
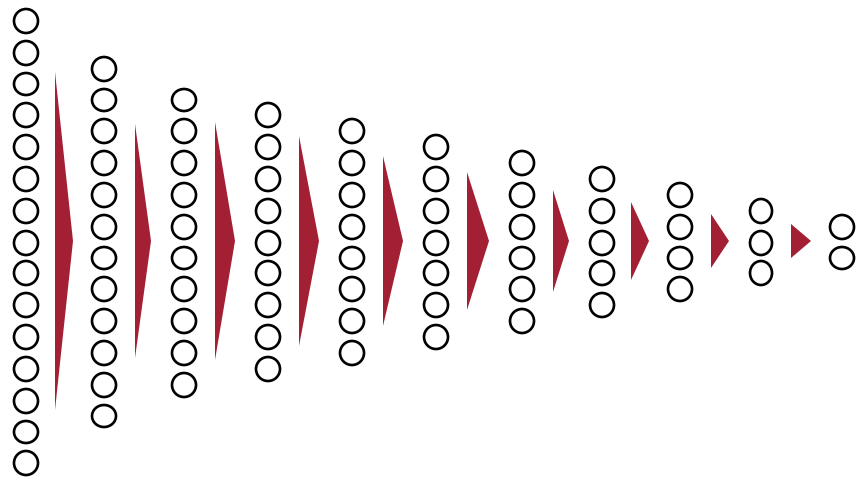


Tan, Mingxing, and Quoc Le. "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks." ICML. 2019.

# Smaller size, same performance?



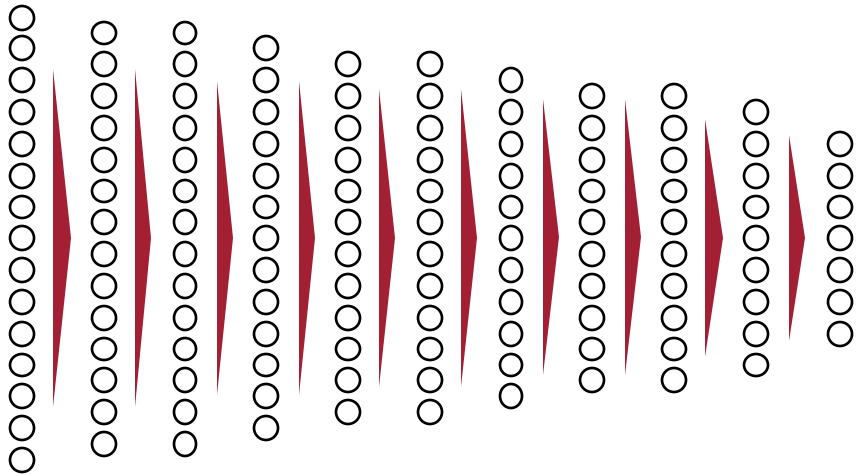
# Smaller size, same performance?



# Objective: Filter Pruning

For a given network  $f_{\theta}(x)$  with parameters  $\theta$  and  $\epsilon, \delta \in (0, 1)$ , generate a compressed, *dense* reparameterization  $\hat{\theta}$  (i.e. with less filters) such that

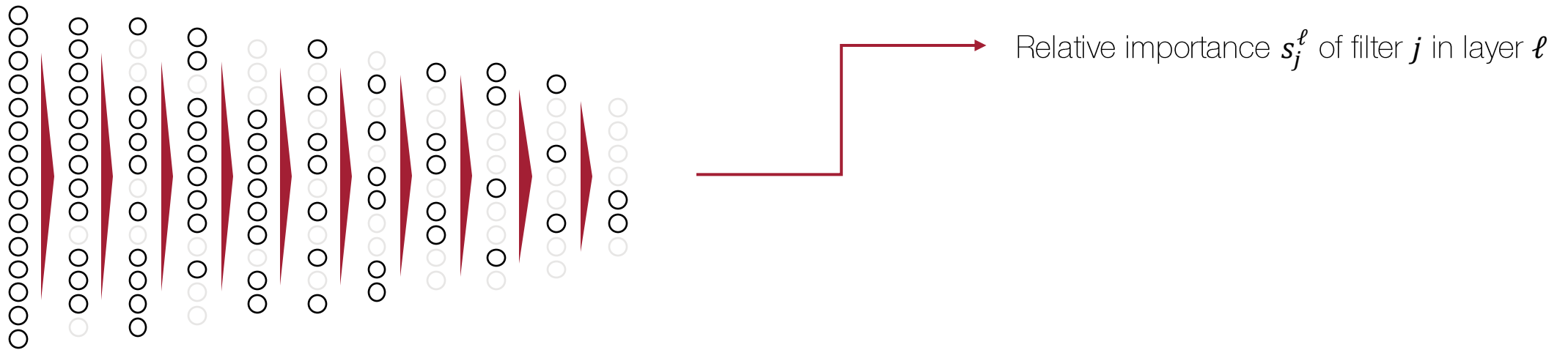
$$\mathbb{P}_{x \sim \mathcal{D}}(f_{\hat{\theta}}(x) \in (1 \pm \epsilon)f_{\theta}(x)) \geq (1 - \delta) \text{ and } \text{size}(\hat{\theta}) \ll \text{size}(\theta).$$



# Objective: Filter Pruning

For a given network  $f_{\theta}(x)$  with parameters  $\theta$  and  $\epsilon, \delta \in (0, 1)$ , generate a compressed, *dense* reparameterization  $\hat{\theta}$  (i.e. with less filters) such that

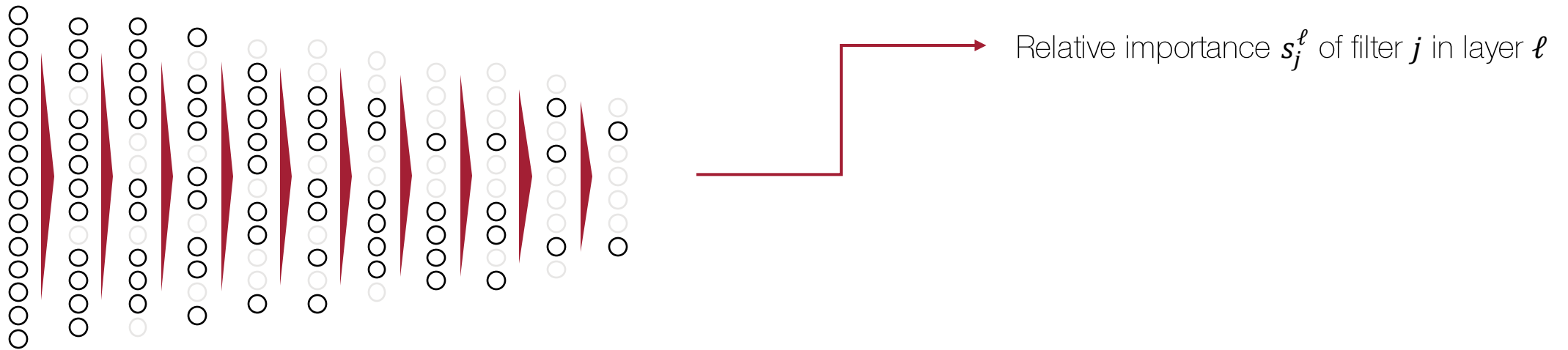
$$\mathbb{P}_{x \sim \mathcal{D}}(f_{\hat{\theta}}(x) \in (1 \pm \epsilon)f_{\theta}(x)) \geq (1 - \delta) \text{ and } \text{size}(\hat{\theta}) \ll \text{size}(\theta).$$



# Objective: Filter Pruning

For a given network  $f_{\theta}(x)$  with parameters  $\theta$  and  $\epsilon, \delta \in (0, 1)$ , generate a compressed, *dense* reparameterization  $\hat{\theta}$  (i.e. with less filters) such that

$$\mathbb{P}_{x \sim \mathcal{D}}(f_{\hat{\theta}}(x) \in (1 \pm \epsilon)f_{\theta}(x)) \geq (1 - \delta) \text{ and } \text{size}(\hat{\theta}) \ll \text{size}(\theta).$$

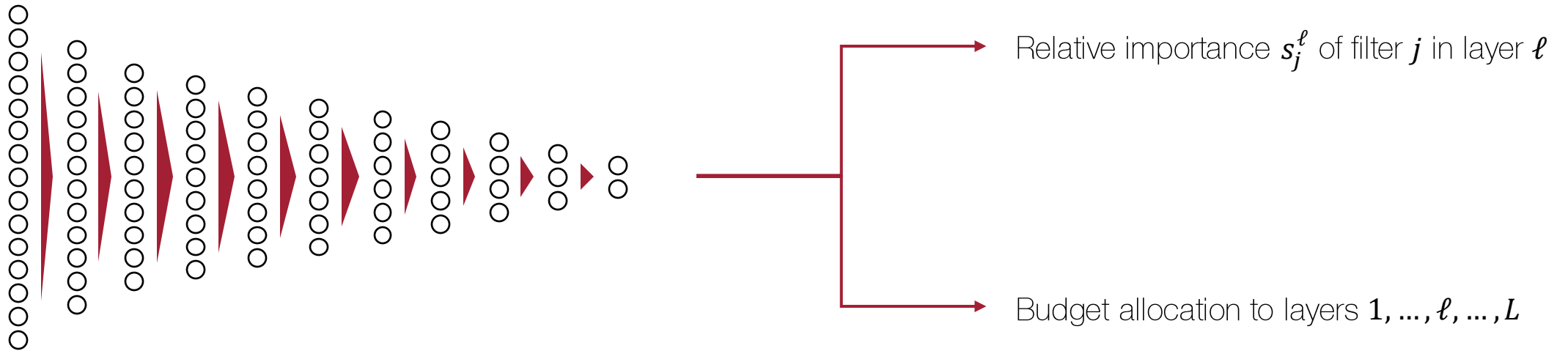




# Objective: Filter Pruning

For a given network  $f_{\theta}(x)$  with parameters  $\theta$  and  $\epsilon, \delta \in (0, 1)$ , generate a compressed, *dense* reparameterization  $\hat{\theta}$  (i.e. with less filters) such that

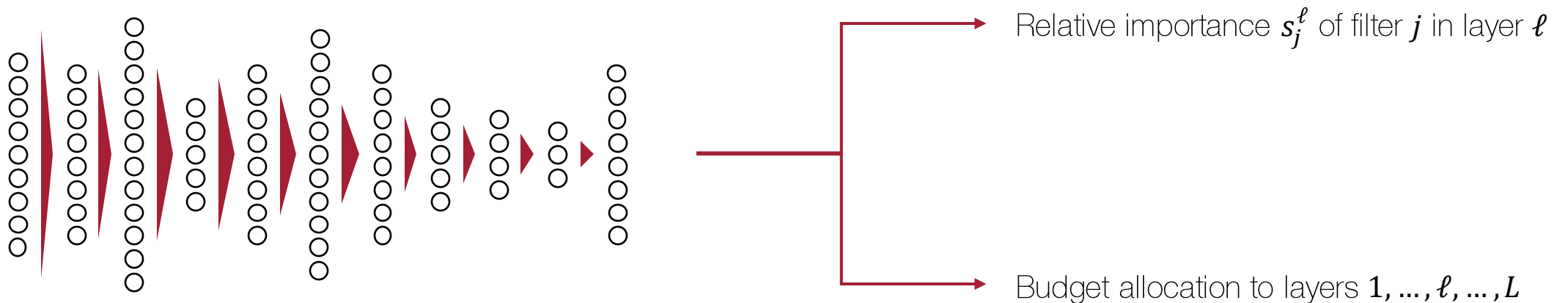
$$\mathbb{P}_{x \sim \mathcal{D}} \left( f_{\hat{\theta}}(x) \in (1 \pm \epsilon) f_{\theta}(x) \right) \geq (1 - \delta) \text{ and } \text{size}(\hat{\theta}) \ll \text{size}(\theta).$$



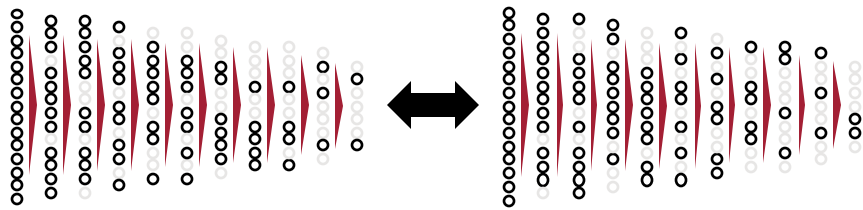
# Objective: Filter Pruning

For a given network  $f_{\theta}(x)$  with parameters  $\theta$  and  $\epsilon, \delta \in (0, 1)$ , generate a compressed, *dense* reparameterization  $\hat{\theta}$  (i.e. with less filters) such that

$$\mathbb{P}_{x \sim \mathcal{D}}(f_{\hat{\theta}}(x) \in (1 \pm \epsilon)f_{\theta}(x)) \geq (1 - \delta) \text{ and } \text{size}(\hat{\theta}) \ll \text{size}(\theta).$$

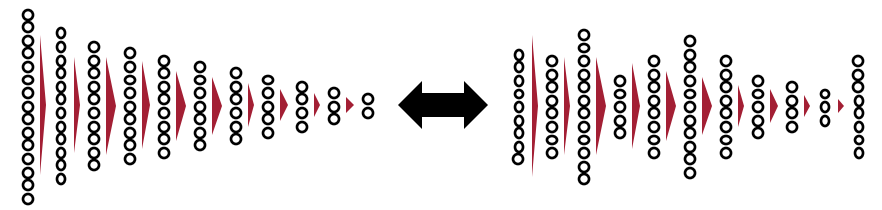


# Related Work



Relative importance  $s_j^\ell$  of filter  $j$  in layer  $\ell$

- $s_j^\ell \sim \|W_j^\ell\|_2$ , Li et al. 2016
- $s_j^\ell \sim \|W_j^\ell\|_1$ , He et al., 2018
- $s_j^\ell \sim 1 / \max_x |z^{\ell+1}(x) - z_{[j]}^{\ell+1}(x)|$ , Luo et al., 2017

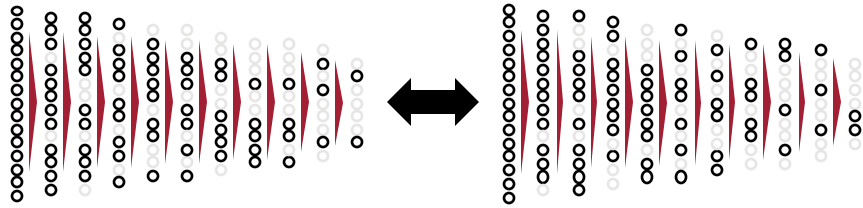


Budget allocation to layers  $1, \dots, \ell, \dots, L$

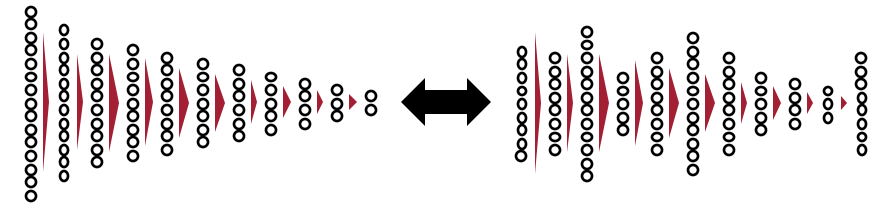
- Uniform budget allocation, He et al., 2018
- Manual ablation study, Li et al. 2016
- Sequential pruning process, Luo et al. 2017

$j$ ...filter,  $\ell$ ...layer,  $s_j^\ell$ ...importance,  $W_j^\ell$ ...filter weights,  $x$ ...input,  $z^\ell(x)$ ...pre-activation,  $z_{[j]}^\ell$ ...pre-activation with feature  $j$  only

# Related Work



Relative importance  $s_j^\ell$  of filter  $j$  in layer  $\ell$



Budget allocation to layers  $1, \dots, \ell, \dots, L$

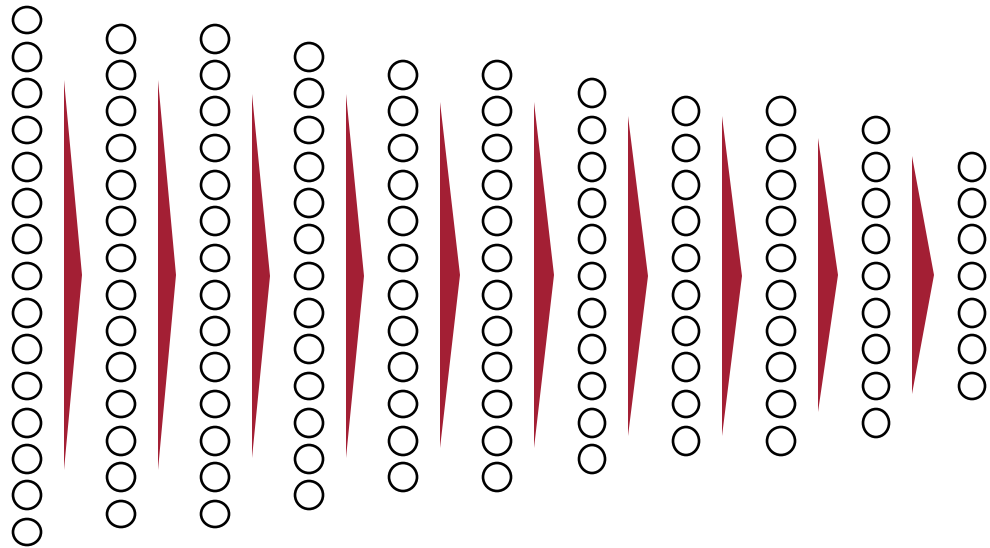
✘ no performance guarantees

✘ data-oblivious heuristics

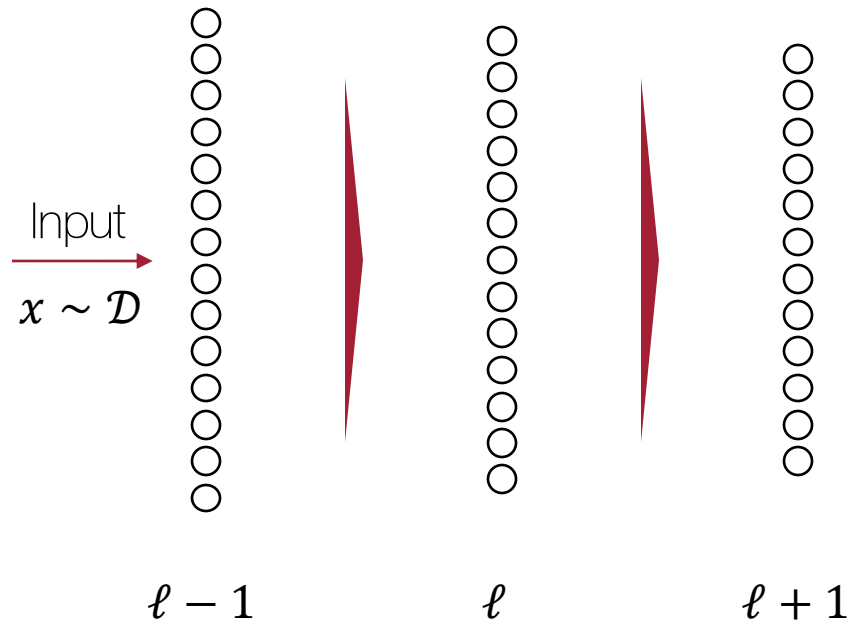
✘ manual procedure

✘ architecture-specific

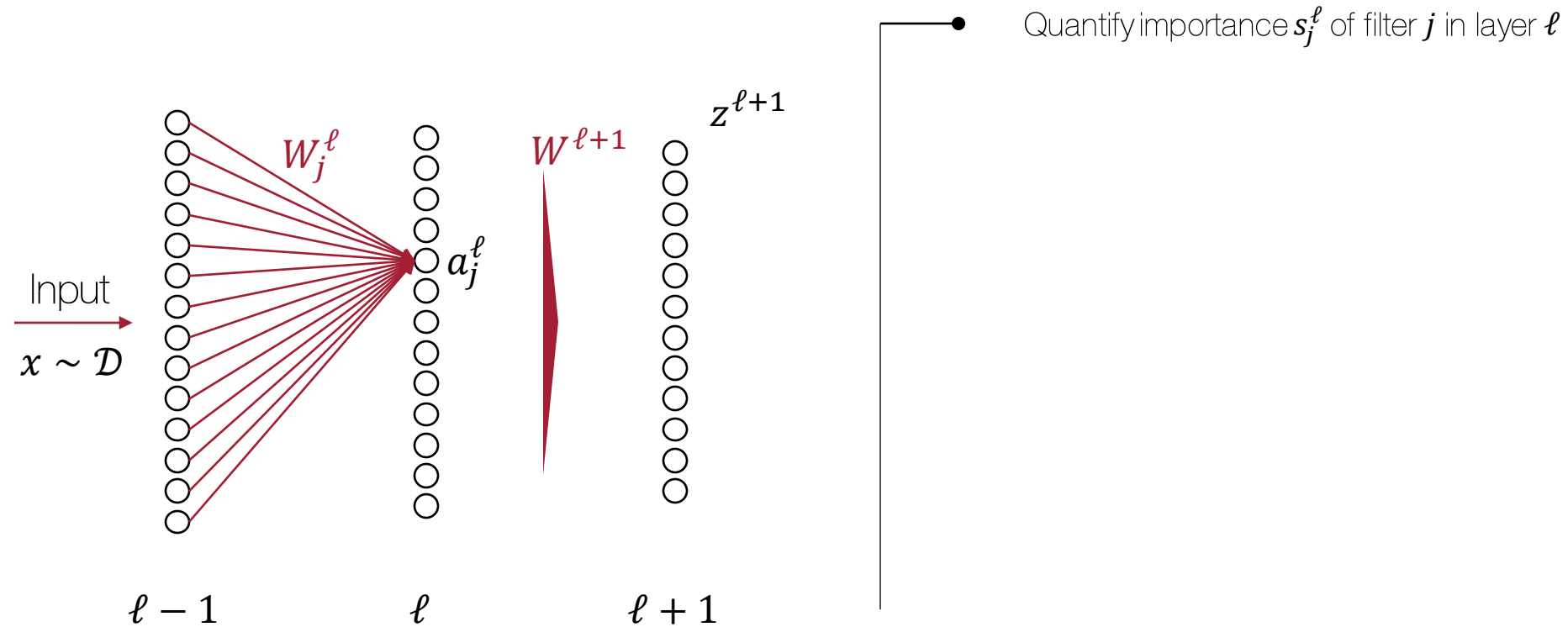
# Our method: filter importance



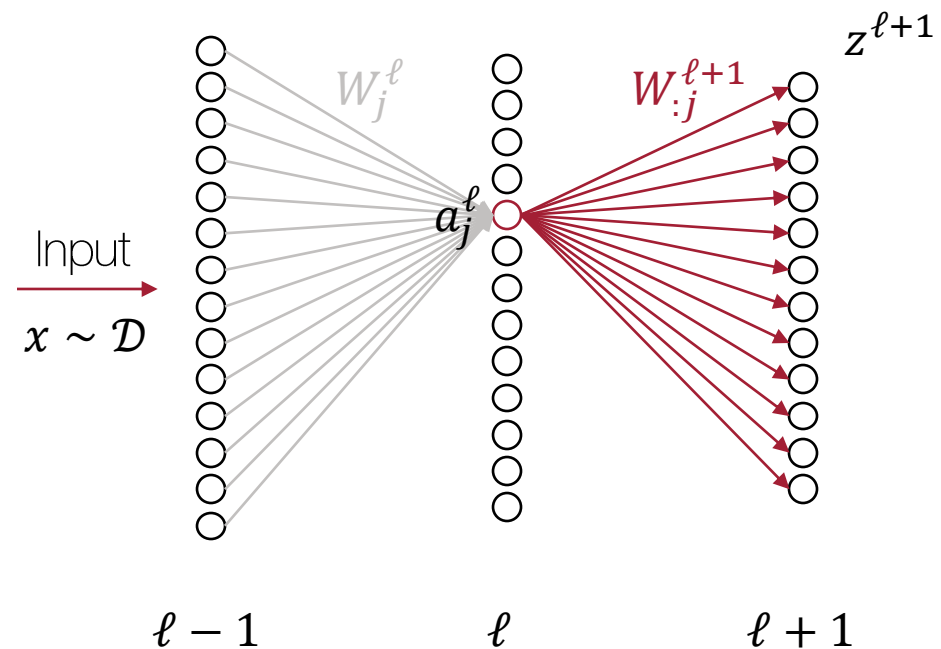
# Our method: filter importance



# Our method: filter importance



# Our method: filter importance



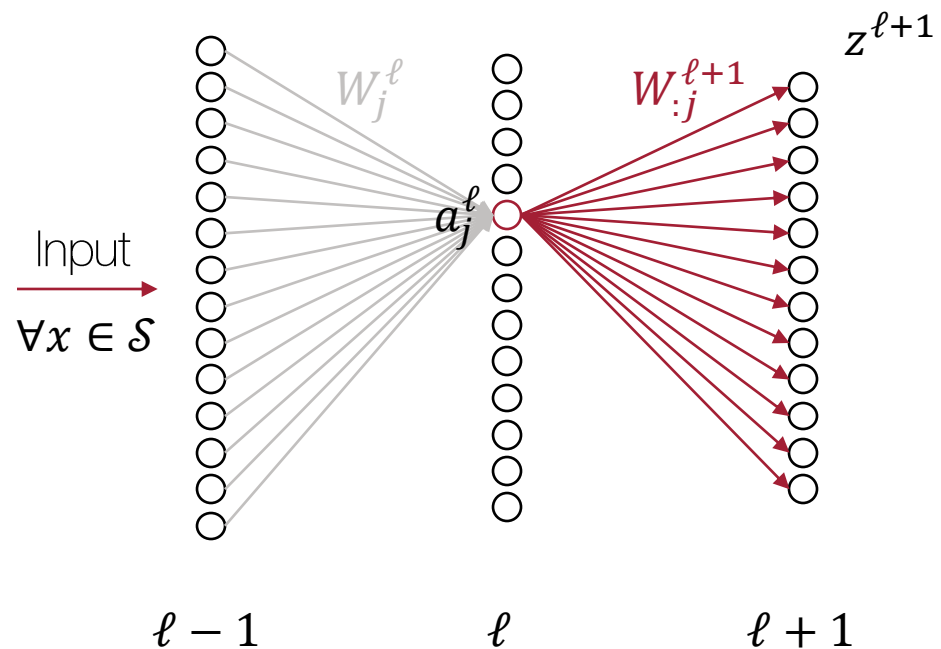
- Quantify importance  $s_j^\ell$  of filter  $j$  in layer  $\ell$

- Consider maximum contribution of activation  $a_j^\ell(x)$  in pre-activation  $z^{\ell+1}(x)$

- $$s_j^\ell \sim \max_i \frac{w_{ij}^\ell a_j^\ell(x)}{z_i^{\ell+1}(x)} \dots \text{filter importance for fixed input } x \sim \mathcal{D}$$



# Our method: filter importance

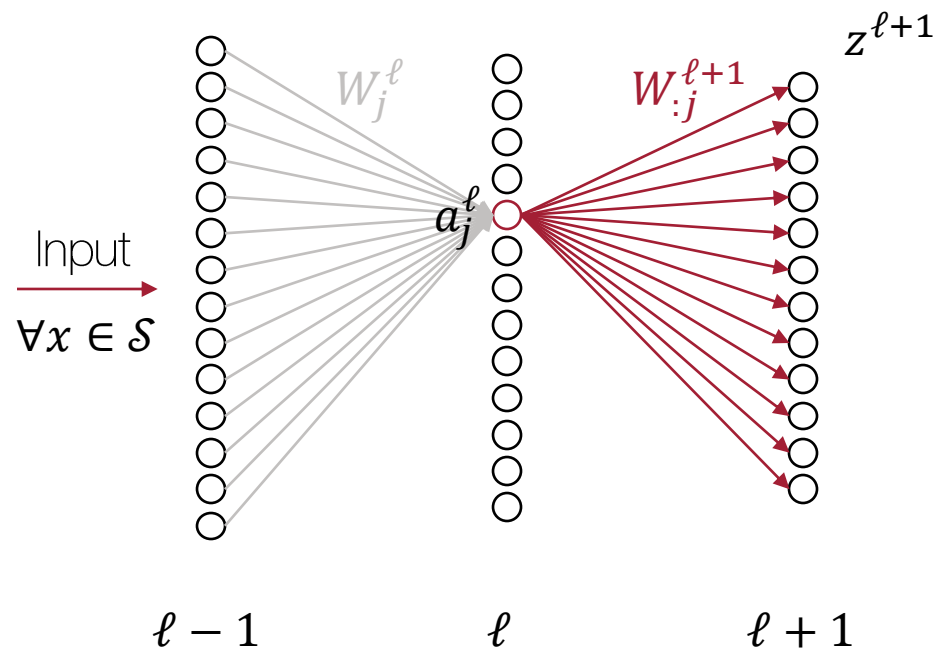


- Quantify importance  $s_j^\ell$  of filter  $j$  in layer  $\ell$

- Consider maximum contribution of activation  $a_j^\ell(x)$  in pre-activation  $z^{\ell+1}(x)$

- $s_j^\ell \sim \max_{x \in \mathcal{S}} \max_i \frac{w_{ij}^\ell a_j^\ell(x)}{z_i^{\ell+1}(x)} \dots$  filter importance for any input  $x \sim \mathcal{D}$  with high probability

# Our method: filter importance



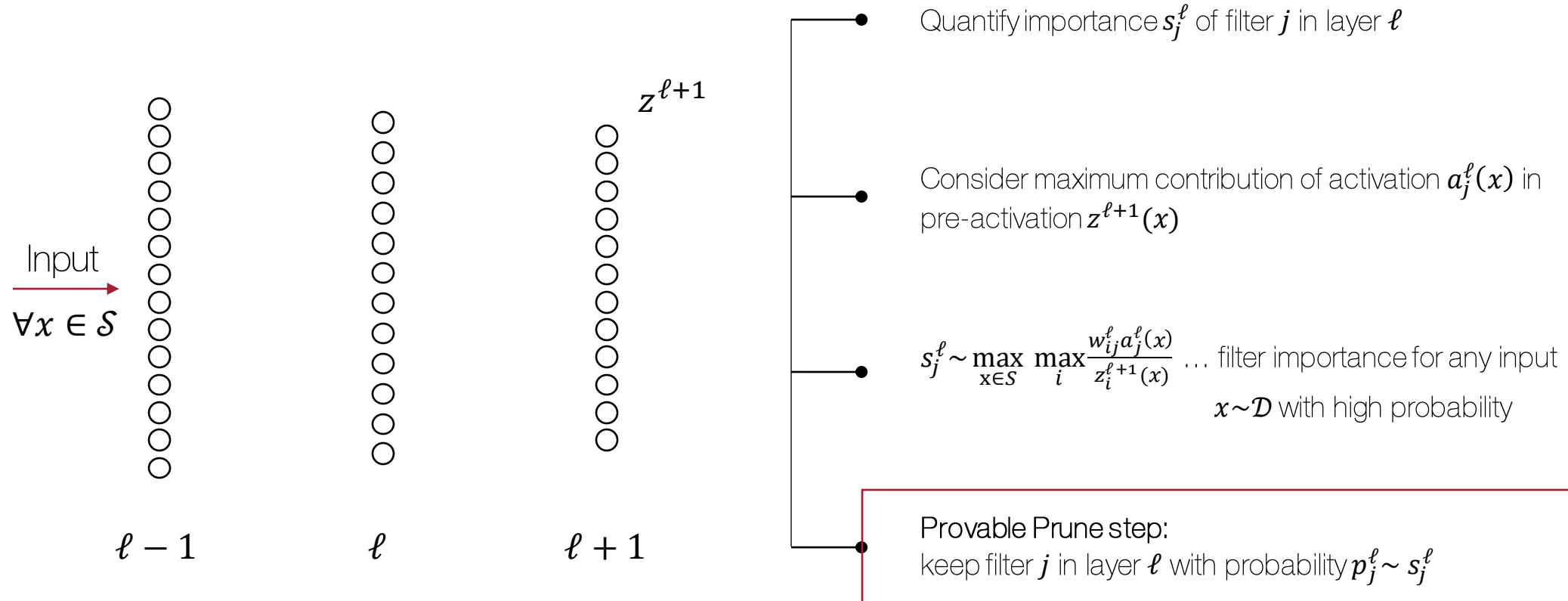
- Quantify importance  $s_j^\ell$  of filter  $j$  in layer  $\ell$

- Consider maximum contribution of activation  $a_j^\ell(x)$  in pre-activation  $z^{\ell+1}(x)$

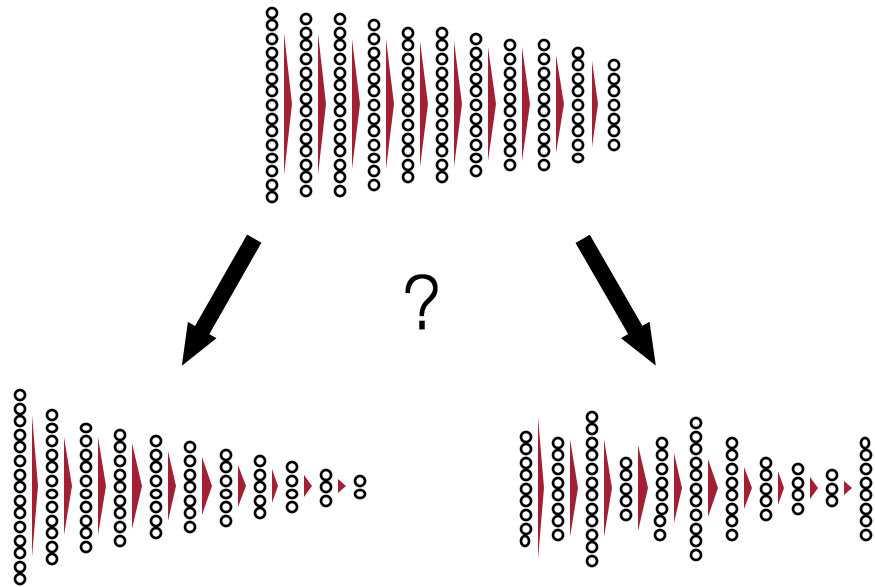
- $s_j^\ell \sim \max_{x \in \mathcal{S}} \max_i \frac{w_{ij}^\ell a_j^\ell(x)}{z_i^{\ell+1}(x)} \dots$  filter importance for any input  $x \sim \mathcal{D}$  with high probability

- Provable Prune step: keep filter  $j$  in layer  $\ell$  with probability  $p_j^\ell \sim s_j^\ell$

# Our method: filter importance



# Our method: budget allocation



Theorem: Relative error  $\epsilon^\ell = f(m^\ell)$  in layer  $\ell$  depends on number of samples  $m^\ell$

Pruned size  $size(\hat{\theta}) = g(m^1, \dots, m^L)$

Allocate budget  $\mathcal{B}$ :  
$$m^1, \dots, m^L = \operatorname{argmin}_{\ell} \max_{\ell} \epsilon^\ell(m^\ell)$$
  
$$s.t. size(\hat{\theta}) = g(m^1, \dots, m^L) \leq \mathcal{B}$$

Solve efficiently via binary search

# Filter compression bounds

For a given network  $f_{\theta}(x)$  with parameters  $\theta$  and  $\epsilon, \delta \in (0, 1)$ , PFP generates a compressed, *dense* reparameterization  $\hat{\theta}$  (i.e. with less filters) such that  $\mathbb{P}_{x \sim \mathcal{D}} \left( f_{\hat{\theta}}(x) \in (1 \pm \epsilon) f_{\theta}(x) \right) \geq (1 - \delta)$  and the number of filters is bounded by  $\mathcal{O} \left( \sum_{\ell=1}^L \frac{L^2 (\Delta^{\ell})^2 S^{\ell} \log^{\eta} / \delta}{\epsilon^2} \right)$

$S^{\ell} :=$  “sum of sensitivities”

- $S^{\ell} = \sum_j s_j^{\ell}$

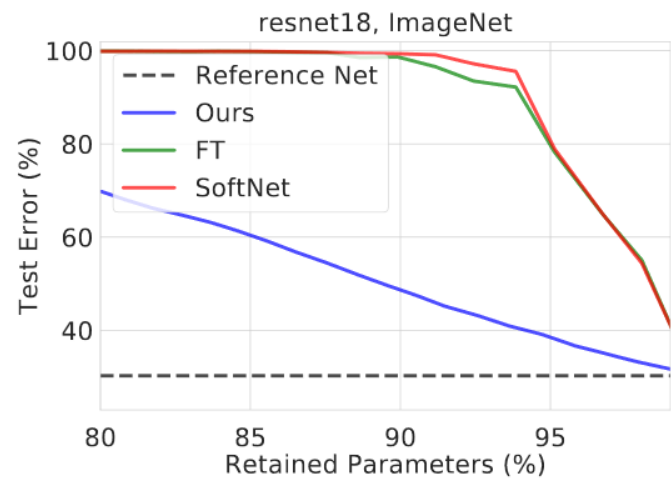
- Quantifies “spread” of importance within a layer

$\Delta^{\ell} :=$  “Propagation Complexity”

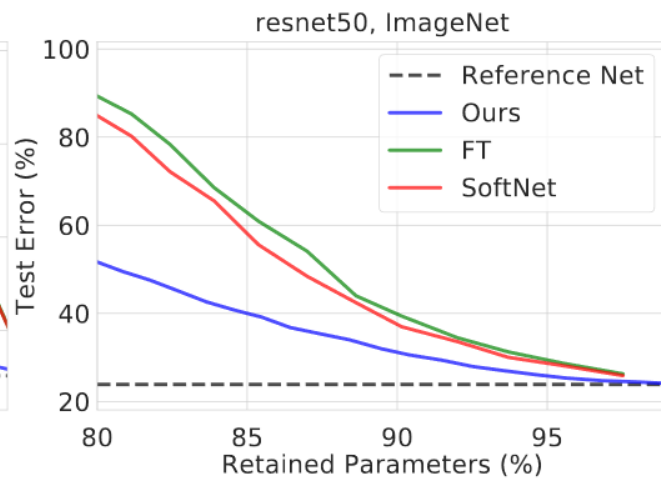
- Ensures desired relative error within layer

- Considers propagation of error across layers

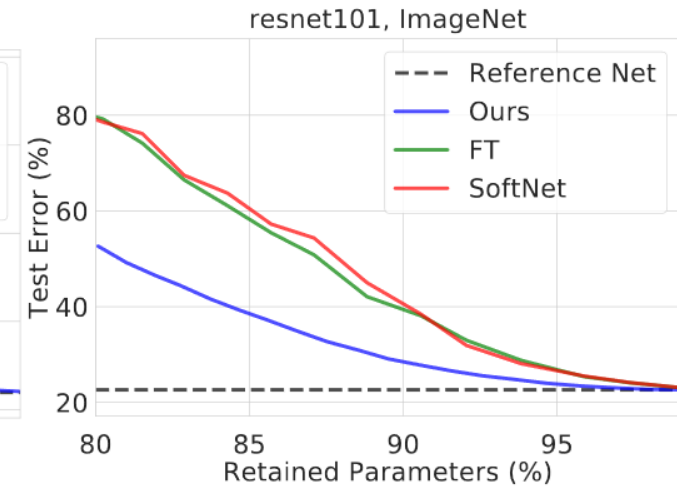
# Results: prune-only



(a) ResNet18

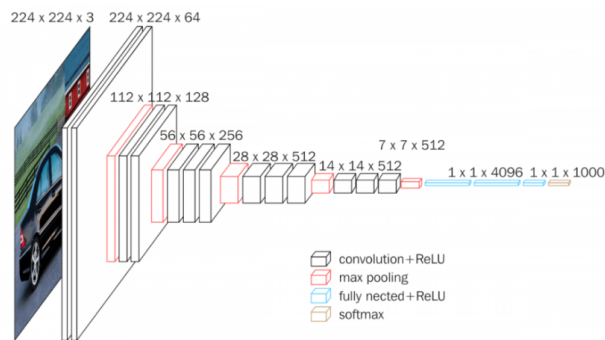


(b) ResNet50

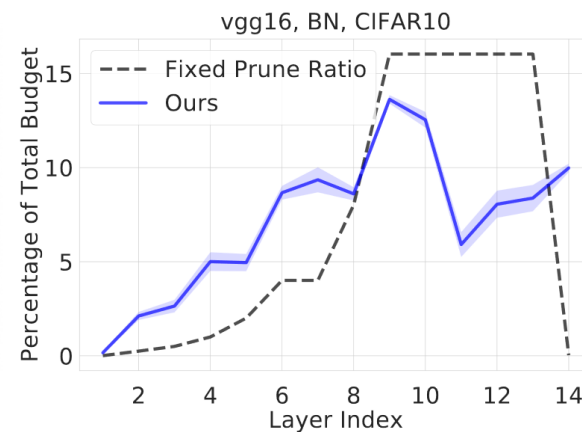


(c) ResNet101

# Results: budget allocation



(a) VGG16 architecture



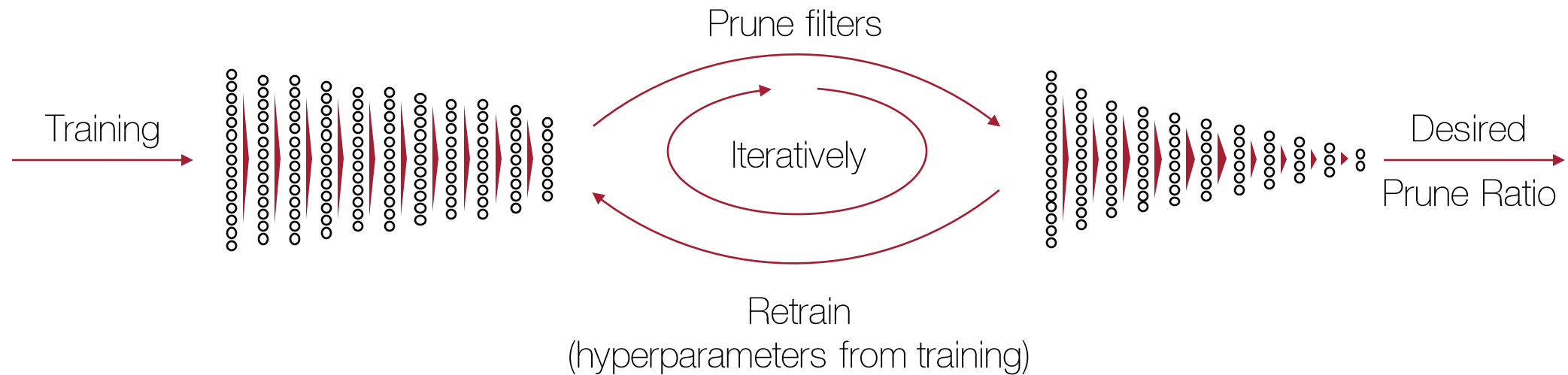
(b) Budget Allocation for VGG16

Early layers are over-sampled  
→ small filters, large images

Middle layers are under-sampled  
→ large filters, small images

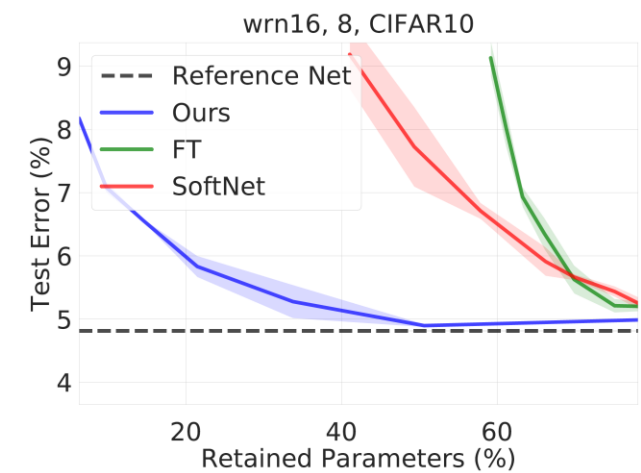
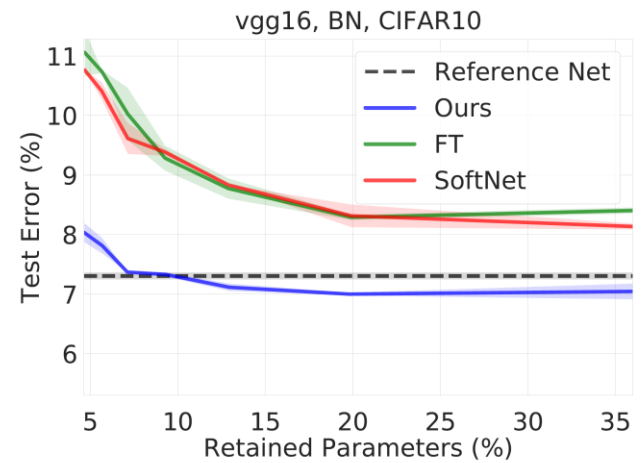
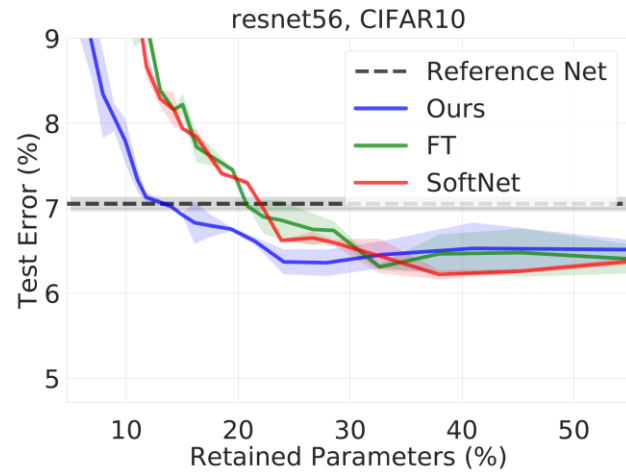
Last layer is over-sampled  
→ classification layer!

# Results: iterative pruning





# Results: iterative pruning



More results in the paper including ImageNet

# Provable Filter Pruning for Efficient Neural Networks

Lucas Liebenwein\*, Cenk Baykal\*, Harry Lang, Dan Feldman, Daniela Rus  
Distributed Robotics Lab, CSAIL, MIT

Paper: <https://openreview.net/forum?id=BJxkOISYDH>



Code: [https://github.com/lucaslie/provable\\_pruning](https://github.com/lucaslie/provable_pruning)



Contact: [lucasl@mit.edu](mailto:lucasl@mit.edu), [baykal@mit.edu](mailto:baykal@mit.edu)

