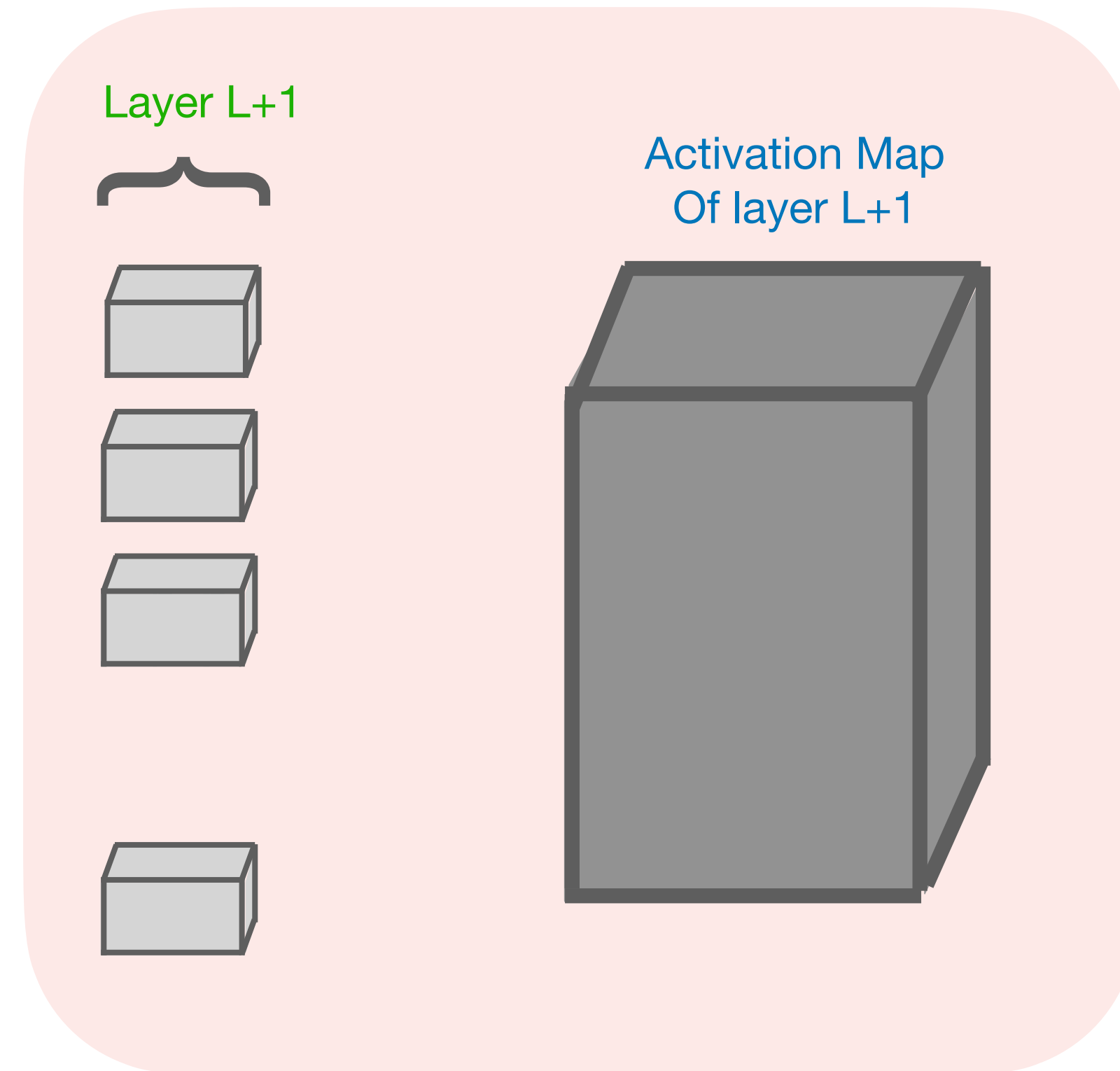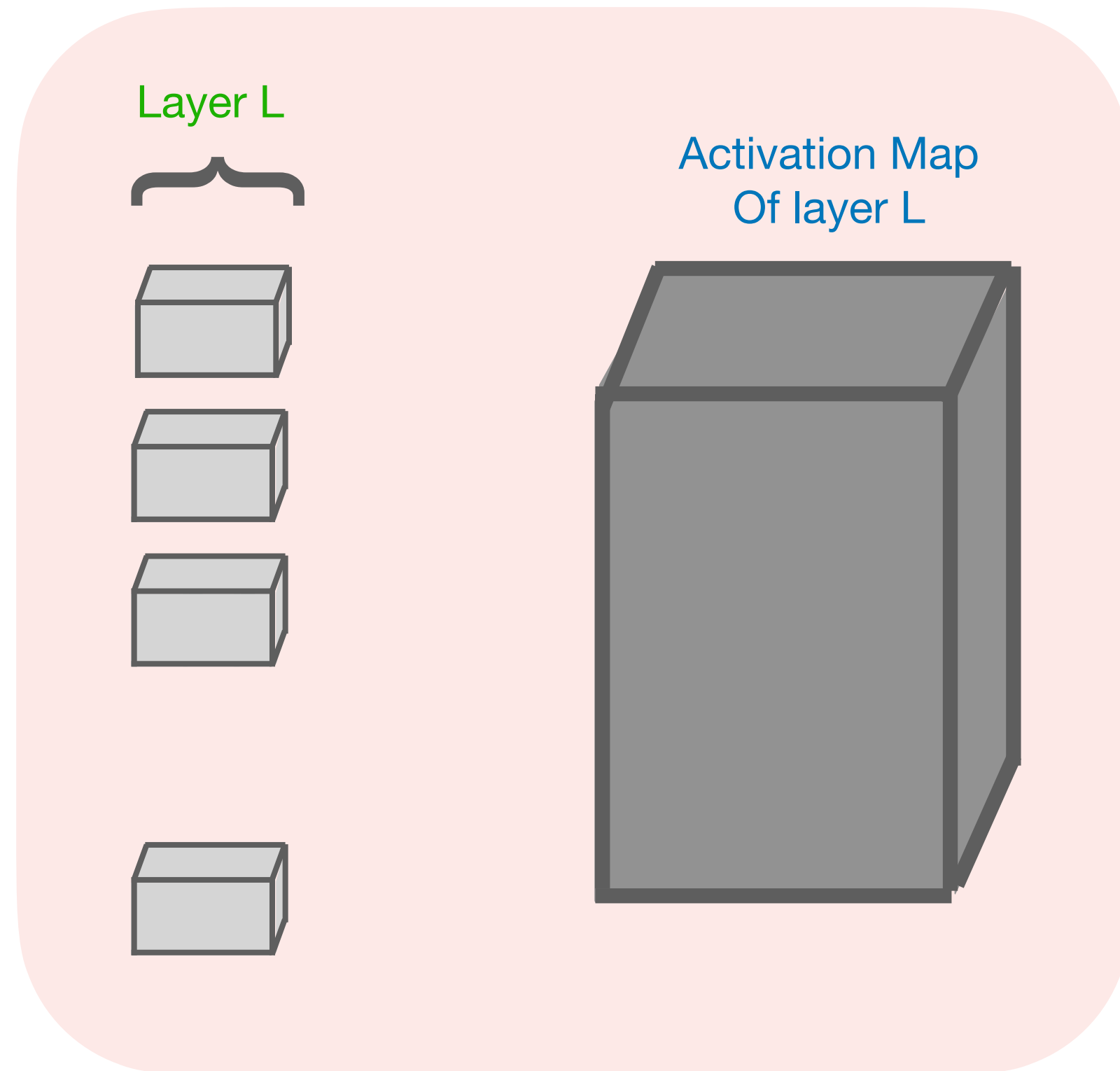# Goal

Reduce the **storage** and **computation** cost of a DNN

$$\mathbf{F}(x; W) \approx \mathbf{F}(x; W_{compressed})$$

Such that $|W_{compressed}| \ll |W|$

# Filter Pruning

Layer L

Activation Map
Of layer L

Layer L+1

Activation Map
Of layer L+1

# Filter Pruning

Layer L

Activation Map
Of layer L

Prune 2 filters

Layer L+1

Activation Map
Of layer L+1

Prune 2 channels

# Filter Pruning

Layer L

Activation Map
Of layer L

Prune 2 filters

Layer L+1

Activation Map
Of layer L+1

Prune 2 channels

Which filters to prune?

**Our method**
# CUP: Cluster Pruning

**Our Idea:** Prune **similar filters**

# CUP: Cluster Pruning

**Our Idea:** Prune **similar filters**

$\text{INPUT}$
$\text{Layer } l$

| Compute per filter features | → | Cluster filters in each layer | → | Prune filters from clusters |
|---|---|---|---|---|

$\text{STEP } 1$ $\qquad$ $\text{STEP } 2$ $\qquad$ $\text{STEP } 3$

$\text{OUTPUT}$
$\text{Pruned layer } l$

# CUP: Cluster Pruning

INPUT

Layer $l$

Compute per filter features

STEP 1

Fully Connected Layer

Convolutional Layer

# CUP: Cluster Pruning

$\text{I}_{\text{NPUT}}$
Layer $l$

Compute per filter features

$\text{S}_{\text{TEP}}$ 1

Fully Connected Layer                                              Convolutional Layer

$\widetilde{W}_{i,:}^{(l)}, \bar{B}_{i,:}^{(l)}$            $\widetilde{W}_{:,i}^{(l+1)}$

neuron $i$

Layer $l$
(m neurons)

Layer $l+1$
(p neurons)

$$\widetilde{F}_{i,:}^{(l)} \; = \; [\; \underbrace{\widetilde{W}_{i,:}^{(l)}, \bar{B}_{i}^{(l)}}_{\text{Incoming features}} \;,\; \underbrace{\widetilde{W}_{:,i}^{(l+1)}}_{\text{Outgoing features}} \;]$$

9

# CUP: Cluster Pruning

INPUT
Layer $l$ → Compute per filter features →

STEP 1

## Fully Connected Layer

$$\widetilde{W}_{i,:}^{(l)}, \bar{B}_{i,:}^{(l)} \qquad \widetilde{W}_{:,i}^{(l+1)}$$



neuron $i$

Layer $l$
(m neurons)

Layer $l$+1
(p neurons)

$$\widetilde{F}_{i,:}^{(l)} = [\ \underbrace{\widetilde{W}_{i,:}^{(l)}, \bar{B}_i^{(l)}}_{\text{Incoming features}}\ ,\ \underbrace{\widetilde{W}_{:,i}^{(l+1)}}_{\text{Outgoing features}}\ ]$$

## Convolutional Layer

$$\widetilde{W}_{:,i,:,:}^{(l)} \qquad \widetilde{W}_{i,:,:,:}^{(l+1)}$$



$K_h$   n   $K_w$   m

filter $i$

Layer $l$
(m filters)

Activation map of
filter $i$ on layer $l$

Layer $l$+1
(p filters)

$$\widetilde{F}_{i,:}^{(l)} = [\ \underbrace{g(\widetilde{W}_{:,i,:,:}^{(l)}), \bar{B}_i^{(l)}}_{\text{Input features}}\ ,\ \underbrace{g(\widetilde{W}_{i,:,:,:}^{(l+1)})}_{\text{Output features}}]$$

10

# CUP: Cluster Pruning

<span style="color:red">INPUT</span>
<span style="color:red">Layer $l$</span>

| Compute per filter features | Cluster filters in each layer |
|---|---|
| STEP 1 | STEP 2 |

$$\mathrm{n} + p + 1$$

$m$

$\tilde{F}^{(l)}$

# CUP: Cluster Pruning

INPUT
Layer $l$ → Compute per filter features → Cluster filters in each layer →

STEP 1    STEP 2

How many clusters?

$n + p + 1$

$m$

$\tilde{F}^{(l)}$

# CUP: Cluster Pruning



INPUT
Layer $l$

| Compute per filter features | Cluster filters in each layer |

STEP 1 STEP 2

$n + p + 1$

$m$

$\tilde{F}^{(l)}$

**Cluster**

$\mathbb{C}_1^{(l)}$
$\mathbb{C}_9^{(l)}$
$\mathbb{C}_2^{(l)}$
$\mathbb{C}_3^{(l)}$
$\mathbb{C}_{10}^{(l)}$
$\mathbb{C}_4^{(l)}$
$\mathbb{C}_{11}^{(l)}$
$\mathbb{C}_5^{(l)}$
$\mathbb{C}_{12}^{(l)}$
$\mathbb{C}_6^{(l)}$
$\mathbb{C}_{15}^{(l)}$
$\mathbb{C}_7^{(l)}$
$\mathbb{C}_{14}^{(l)}$
$\mathbb{C}_{13}^{(l)}$
$\mathbb{C}_8^{(l)}$

How many clusters?

13

# CUP: Cluster Pruning



<span style="color:red">I</span>NPUT
<span style="color:red">Layer $l$</span>

| Compute per filter features | Cluster filters in each layer |
|---|---|

STEP 1　　　　　　STEP 2

$n + p + 1$

$m$

$\tilde{F}^{(l)}$

**Cluster**

$\mathbb{C}_1^{(l)}$
$\mathbb{C}_9^{(l)}$
$\mathbb{C}_2^{(l)}$
$\mathbb{C}_3^{(l)}$
$\mathbb{C}_{10}^{(l)}$
$\mathbb{C}_4^{(l)}$
$\mathbb{C}_{11}^{(l)}$
$\mathbb{C}_5^{(l)}$
$\mathbb{C}_{12}^{(l)}$
$\mathbb{C}_6^{(l)}$
$\mathbb{C}_{14}^{(l)}$
$\mathbb{C}_7^{(l)}$
$\mathbb{C}_{13}^{(l)}$
$\mathbb{C}_8^{(l)}$
$\mathbb{C}_{15}^{(l)}$

<span style="color:red">**Global Threshold $t$**</span>

<span style="color:red">How many clusters?</span>

14

# CUP: Cluster Pruning



INPUT
Layer $l$

| Compute per filter features | Cluster filters in each layer |

STEP 1      STEP 2

$n + p + 1$

$m$

$\tilde{F}^{(l)}$

**Cluster**

$\mathbb{C}_1^{(l)}$
$\mathbb{C}_9^{(l)}$
$\mathbb{C}_2^{(l)}$
$\mathbb{C}_3^{(l)}$
$\mathbb{C}_{11}^{(l)}$
$\mathbb{C}_{10}^{(l)}$
$\mathbb{C}_4^{(l)}$
$\mathbb{C}_5^{(l)}$
$\mathbb{C}_{15}^{(l)}$
$\mathbb{C}_{12}^{(l)}$
$\mathbb{C}_6^{(l)}$
$\mathbb{C}_{14}^{(l)}$
$\mathbb{C}_7^{(l)}$
$\mathbb{C}_{13}^{(l)}$
$\mathbb{C}_8^{(l)}$

**Global Threshold $t$**

**Clip Dendogram**

$\mathbb{C}_{11}^{(l)}$
$\mathbb{C}_{12}^{(l)}$
$\mathbb{C}_{13}^{(l)}$

How many clusters?

15

# CUP: Cluster Pruning

# CUP: Cluster Pruning

Layer $l$

OUTPUT
Pruned layer $l$

| Compute per filter features | Cluster filters in each layer | Prune filters from clusters |
|---|---|---|
| STEP 1 | STEP 2 | STEP 3 |

$\mathbb{C}_{11}^{(l)}$

$\mathbb{C}_{12}^{(l)}$

$\mathbb{C}_{13}^{(l)}$

$$\mathbb{S}_r^{(l)} = \underset{i \in \mathbb{C}_r^{(l)}}{\mathrm{argmax}} \|\widetilde{F}_{i,:}^{(l)}\|_2$$

# CUP: Cluster Pruning

INPUT
Layer $l$

OUTPUT
Pruned layer $l$

| Compute per filter features | Cluster filters in each layer | Prune filters from clusters |
|---|---|---|
| STEP 1 | STEP 2 | STEP 3 |

$$\mathbb{S}_r^{(l)} = \underset{i \in \mathbb{C}_r^{(l)}}{\mathrm{argmax}} \|\widetilde{F}_{i,:}^{(l)}\|_2$$

$\mathbb{C}_{11}^{(l)}$

$\mathbb{C}_{12}^{(l)}$

$\mathbb{C}_{13}^{(l)}$

#clusters = # remaining filters

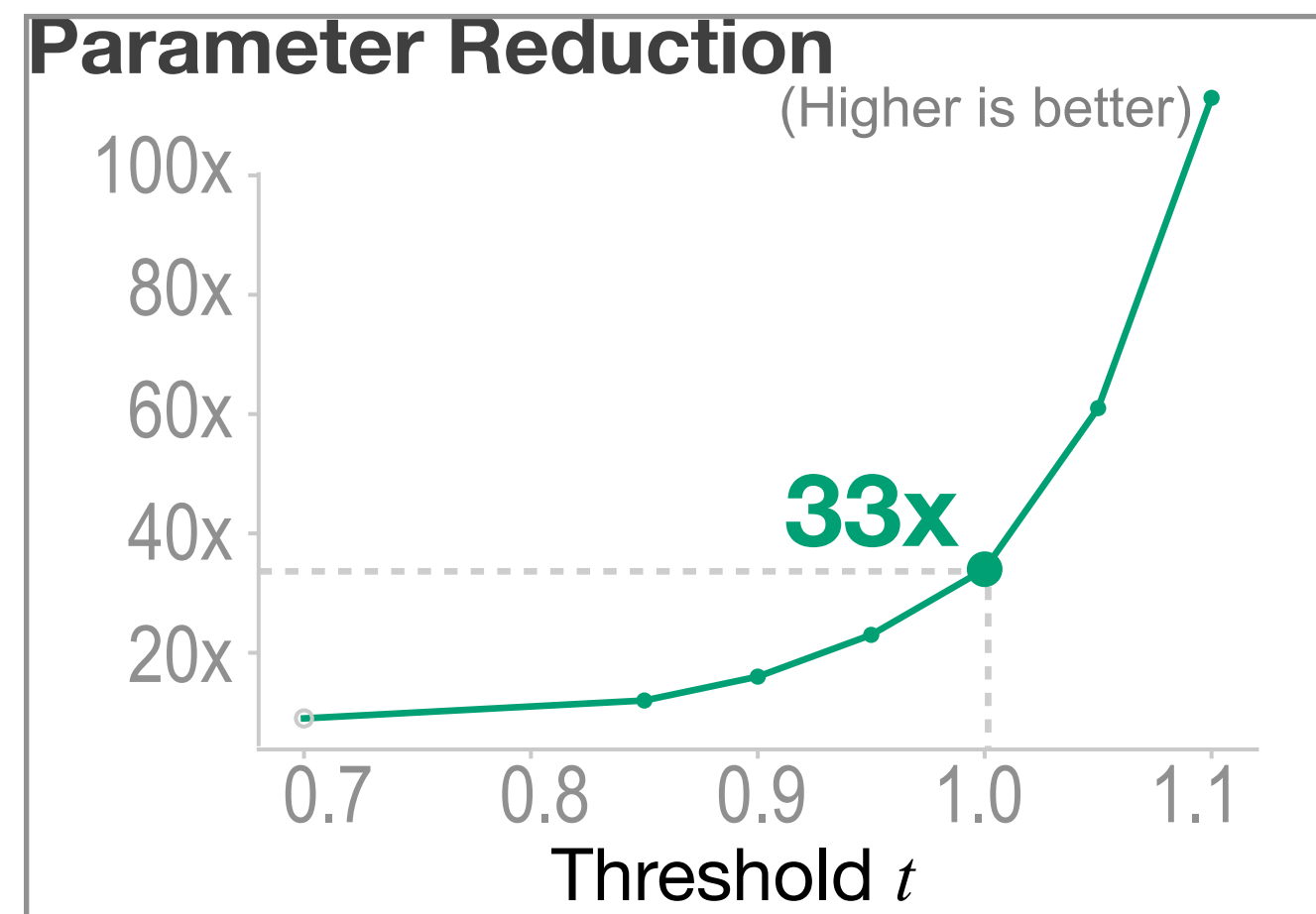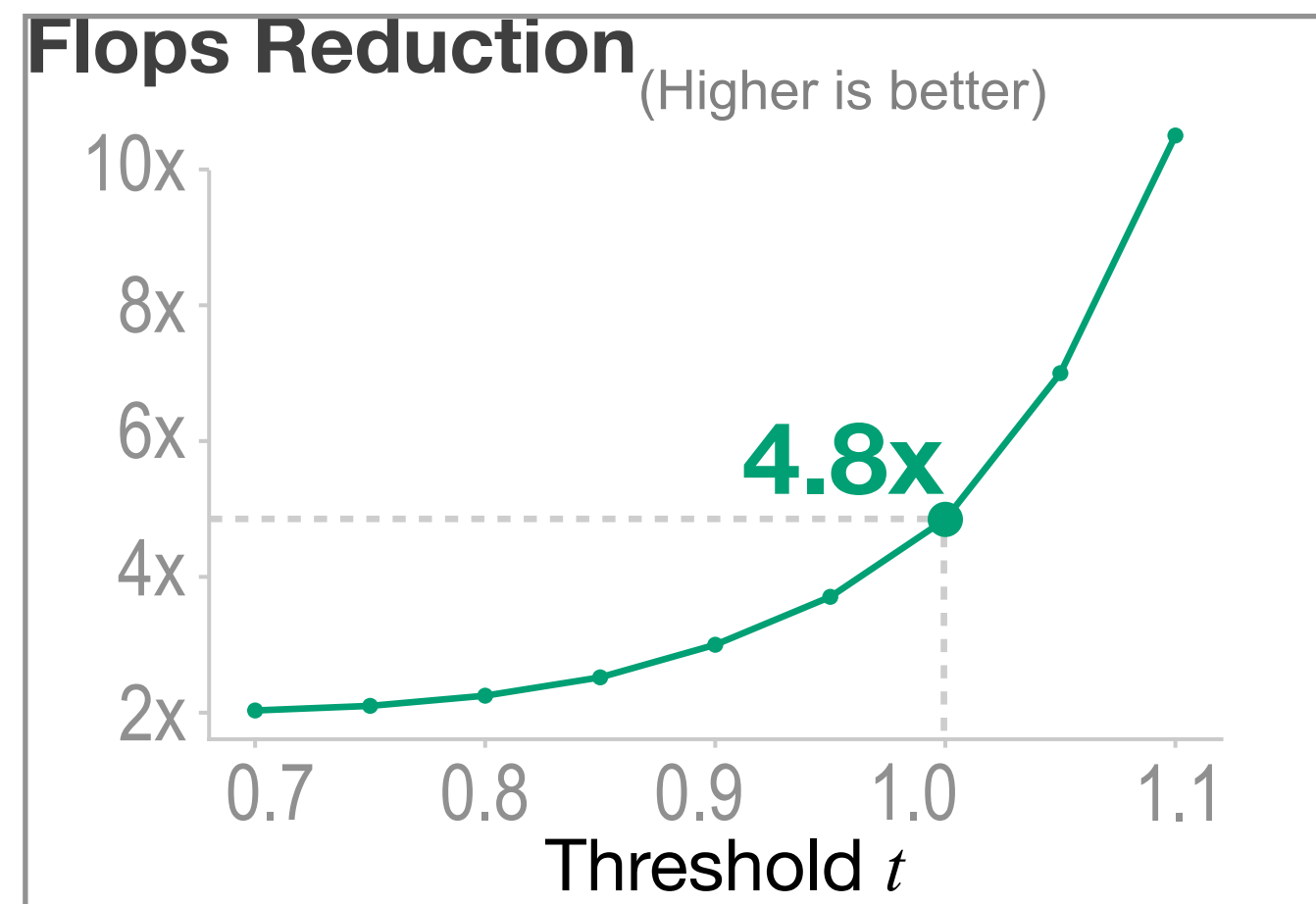**t** parameterizes pruning amount
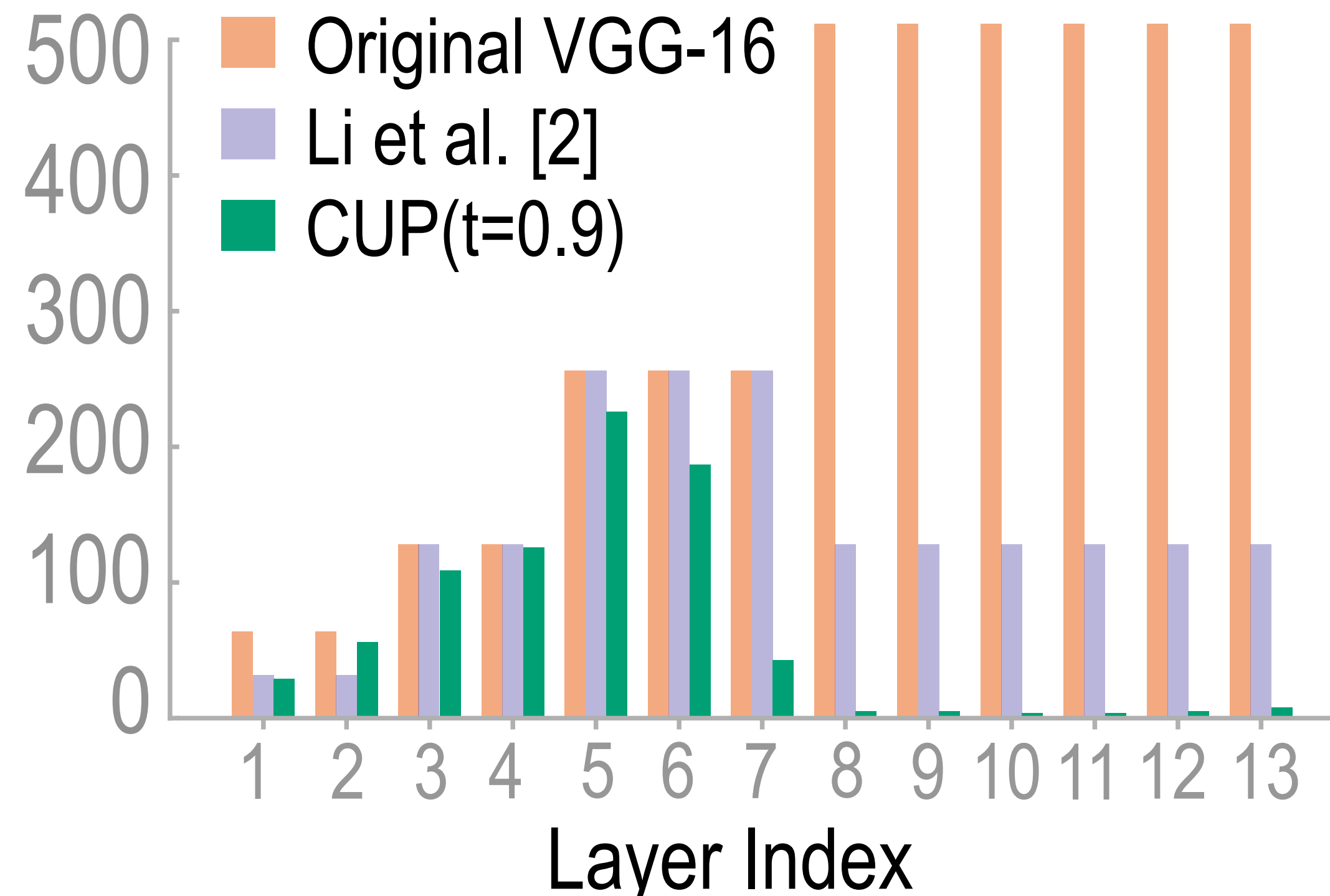
# Results

## **Benefit 1**: Single hyper parameter control over pruning amount

# Results

**Benefit 2**: Non uniform pruning with a single hyper-parameter t



# Filters post pruning (Lower is better)

# Results

**Benefit 3**: Training time reduction through train time pruning.

| | Method | Retrain? | Top-1 (%) | FR (×) | Training Time (GPU Hours) |
|---|---|---|---|---|---|
| | Resnet-50 | - | 75.86 | 1.00 | 66.0 |
| | SFP [14] | ✗ | 74.01 | 1.73 | 61.8 |
| | GM [15] | ✗ | 74.13 | 2.15 | 62.2 |
| | **CUP-RF** (ours) | ✗ | **74.34** | **2.21** | **51.6** |

~15 hours saving with 2x compression

# Results

**Benefit 4**: State-of-the-art compression

| Model | Method | Retrain? | FR ($\times$) | Acc. ($\Delta\%$) | |
|---|---|---|---|---|---|
| | | | | Top-1 | Top-5 |
| ResNet-18 | GM [15] | ✓ | 1.71 | -1.87 | -1.15 |
| | COP [29] | ✓ | **1.75** | -2.48 | - |
| | **CUP** (Our) | ✓ | → **1.75** | **-1.00** | **-0.79** |
| | SFP [14] | ✗ | 1.71 | -3.18 | -1.85 |
| | GM [15] | ✗ | 1.71 | -2.47 | -1.52 |
| | **CUP-RF** (ours) | ✗ | → **1.75** | **-2.37** | **-1.40** |
| ResNet-34 | L1 [2] | ✓ | 1.31 | -1.06 | - |
| | GM [15] | ✓ | 1.69 | -1.29 | -0.54 |
| | **CUP** (ours) | ✓ | → **1.78** | **-0.86** | **-0.53** |
| | SFP [14] | ✗ | 1.69 | -2.09 | -1.29 |
| | GM [15] | ✗ | 1.69 | -2.13 | -0.92 |
| | **CUP-RF** (ours) | ✗ | → **1.71** | **-1.61** | **-0.89** |
| ResNet-50 | SFP [14] | ✓ | 2.15 | -14.0 | -8.20 |
| | MP [30] | ✓ | 2.05 | -1.20 | - |
| | **CUP** (ours) | ✓ | → **2.47** | **-1.17** | **-0.81** |
| | SFP [14] | ✗ | 1.71 | -1.54 | **-0.81** |
| | GM [15] | ✗ | 2.15 | -2.02 | -0.93 |
| | **CUP-RF** (ours) | ✗ | → **2.20** | **-1.47** | -0.88 |

# Conclusion

**Thank you!**

## CUP: Cluster pruning framework

- Prunes a DNN by clustering similar filters.

## Benefits of CUP

- Single hyper-parameter control over pruning amount.
- Enables non uniform pruning across layers.
- Train time savings.

## Extensive evaluation on large DNNs & datasets