



PRUNING CONVOLUTIONAL NEURAL NETWORKS

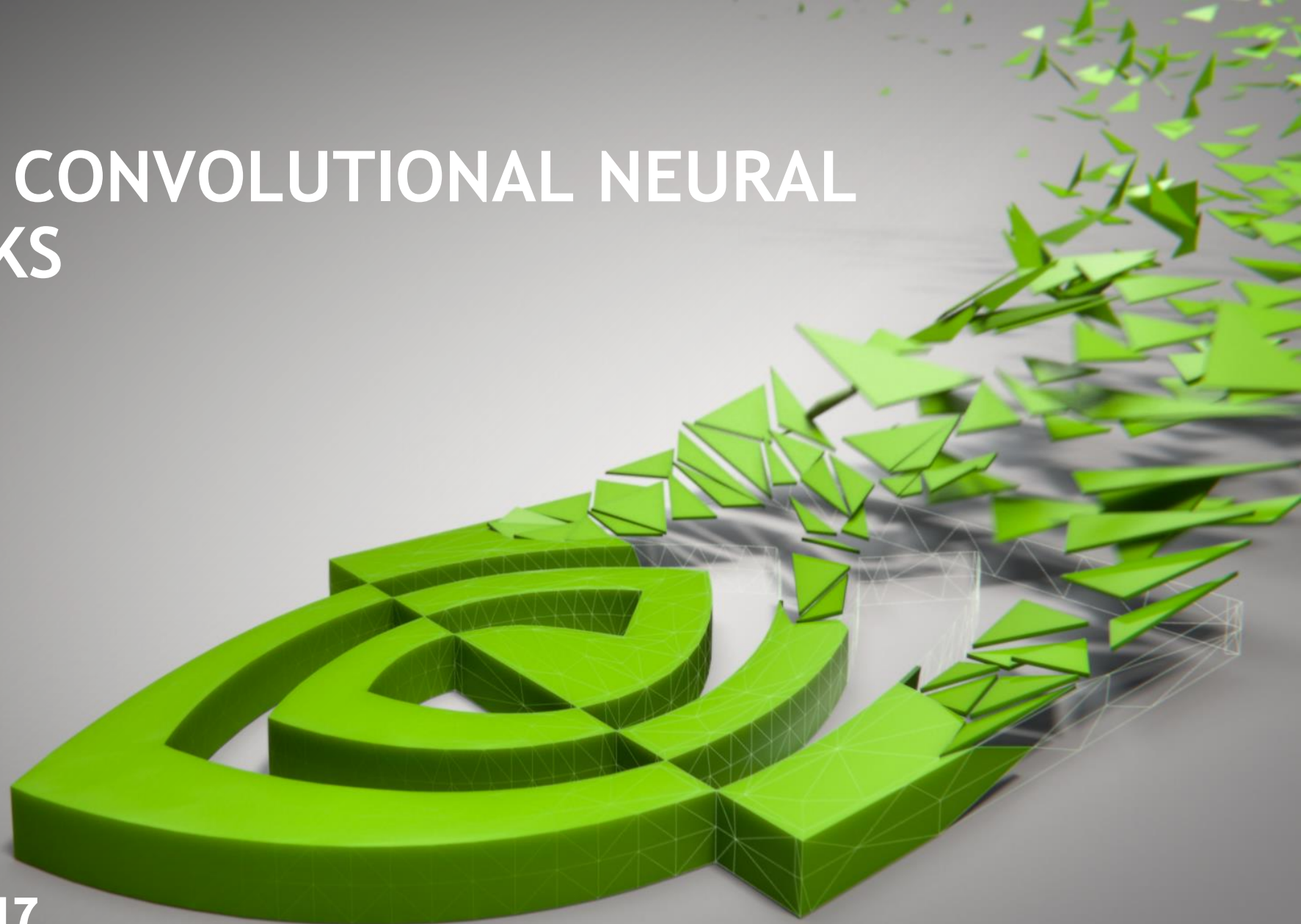
Pavlo Molchanov

Stephen Tyree

Tero Karras

Timo Aila

Jan Kautz



WHY WE CAN PRUNE CNNs?

WHY WE CAN PRUNE CNNs?

Optimization “failures”:

- Some neurons are "dead": little activation
- Some neurons are uncorrelated with output

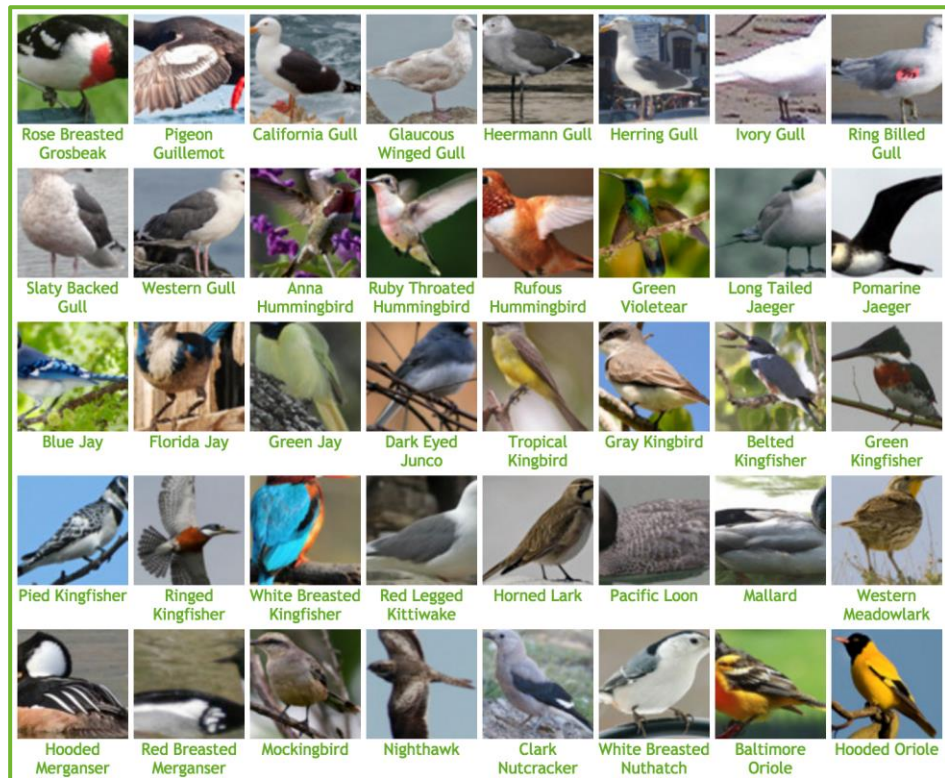
Modern CNNs are overparameterized:

- **VGG16** has **138M** parameters
- **Alexnet** has **61M** parameters
- **ImageNet** has **1.2M** images

PRUNING FOR TRANSFER LEARNING



Small Dataset

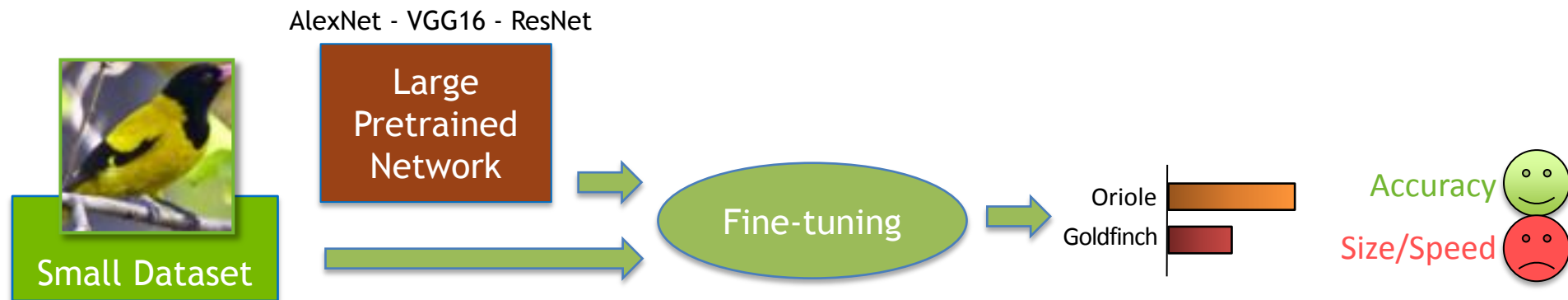


Caltech-UCSD Birds (200 classes, <6000 images)

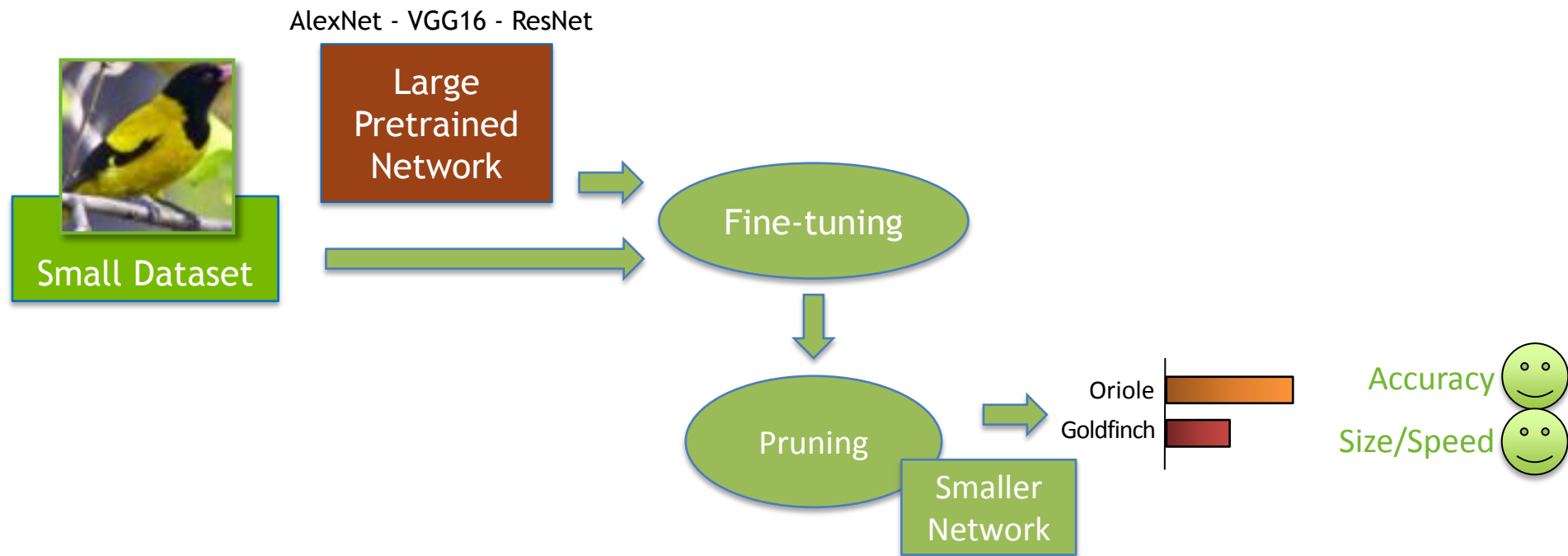
PRUNING FOR TRANSFER LEARNING



PRUNING FOR TRANSFER LEARNING



PRUNING FOR TRANSFER LEARNING



TYPES OF UNITS

- Convolutional units
 - **Heavy on computation**
 - Small on storage
- Fully connected (dense) units
 - Fast on computations
 - Heavy on storage

Our focus

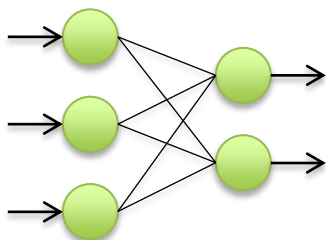
To reduce computation,
we focus pruning on
convolutional units.

Ratio of floating point operations

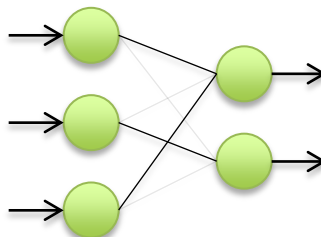
	Convolutional layers	Fully connected layers
VGG16	99%	1%
Alexnet	89%	11%
R3DCNN	90%	10%

TYPES OF PRUNING

No pruning

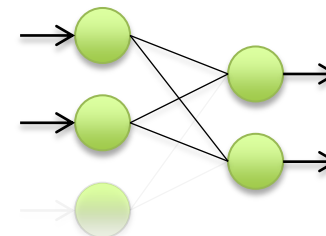


Fine pruning



- Remove connections between neurons/feature maps
- May require special SW/HW for full speed-up

Coarse pruning



- Remove entire neurons/feature maps
- Instant speed-up
- No change to HW/SW

Our focus

NETWORK PRUNING

NETWORK PRUNING

Training:

$$\min_W \mathcal{C}(W, \mathcal{D})$$

\mathcal{C} : training cost function

\mathcal{D} : training data

W : network weights

\hat{W} : *pruned* network weights

NETWORK PRUNING

Training:

$$\min_W C(W, \mathcal{D})$$

C : training cost function

\mathcal{D} : training data

W : network weights

\hat{W} : *pruned* network weights

Pruning:

$$\min_{\hat{W}} |C(\hat{W}, \mathcal{D}) - C(W, \mathcal{D})|$$

$$s.t. \hat{W} \subset W, |\hat{W}| < B$$

NETWORK PRUNING

Training:

$$\min_W C(W, \mathcal{D})$$

C : training cost function

\mathcal{D} : training data

W : network weights

\hat{W} : *pruned* network weights

Pruning:

$$\min_{\hat{W}} |C(\hat{W}, \mathcal{D}) - C(W, \mathcal{D})|$$

~~$$s.t. \hat{W} \in W, |\hat{W}| < B$$~~

$$s.t. \|\hat{W}\|_0 \leq B$$

$\|\cdot\|_0$ – ℓ_0 norm, number of non zero elements

NETWORK PRUNING

Exact solution: combinatorial optimization problem - **too computationally expensive**

- VGG-16 has $|W| = 4224$ convolutional units

$2^W = 35538712055317885020276167051778952343269622838113490006838344535516384949349808265709886296748165086713339379979429715454985631857840126159027259220283889576931421862796735241131771064707150729404513525374011172364491439311003809147986212244125583682040173009664289254204672705377527023751838969121362871174353608981432683121364549161158770063228722675736010638821281170939104924344940969413158186617489468428542655114822243445927713846770846835644172876711560142902677438665366455880288479809069658760988833949942077659397959942214951022455293213581331690534711750984388463798139279635882246499968899123956774486595348698818284747613874694623754391634523542345894518795402778976197641675203085270364961383790287738178866981707575145292010325953635643917893687322226855341345293028465563634475713300900704784609781200491091266517708547049178192081173208302835906844291042266393938301265721160541880258623908153646996141044116326428425940756760134968815712848010684237572487512170690618881568084176810268745960486335685758930475537127132998300931396086947503485054946846061296719461238733586584900523337276581733454482412202328028231240265027731391290867726741995809784279019489403498646468630714031376402488628074647455635839933307882358008948992762943104694366519689215$

NETWORK PRUNING

Exact solution: combinatorial optimization problem - **too computationally expensive**

- VGG-16 has $|W| = 4224$ convolutional units

$2^W = 3553871205531788502027616705177895234326962283811349000683834453551638494934980826570988629674816508671333937997942971545498563185784012615902725922028388957693142186279673524113177106470715072940451352537401117236449143931100380914798621224412558368204017300966428925420467270537752702375183896912136287117435360898143268312136454916115877006322872267573601063882128117093910492434494096941315818661748946842854265511482224344592771384677084683564417287671156014290267743866536645588028847980906965876098883394994207765939795994221495102245529321358133169053471175098438846379813927963588224649996889912395677448659534869881828474761387469462375439163452354234589451879540277897619764167520308527036496138379028773817886698170757514529201032595363564391789368732222685534134529302846556363447571330090070478460978120049109126651770854704917819208117320830283590684429104226639393830126572116054188025862390815364699614104411632642842594075676013496881571284801068423757248751217069061888156808417681026874596048633568575893047553712713299830093139608694750348505494684606129671946123873358658490052333727658173345448241220232802823124026502773139129086772674199580978427901948940349864646863071403137640248862807464745563583993307882358008948992762943104694366519689215$

Greedy pruning:

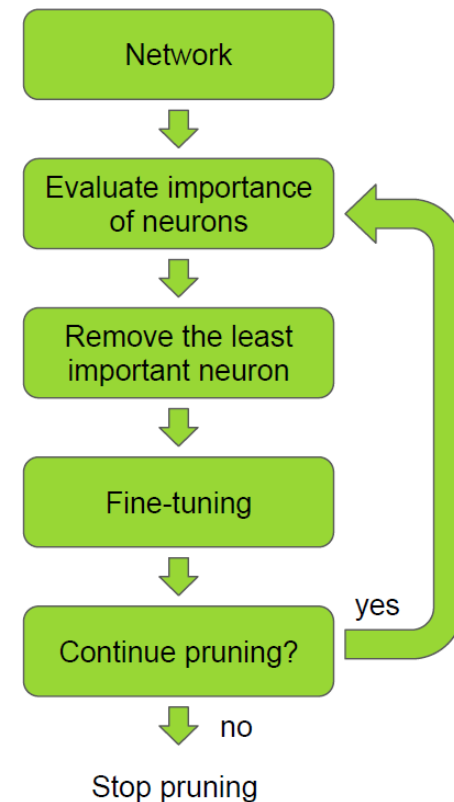
- Assumes all neurons are independent (same assumption for back propagation)
- Iteratively, remove neuron with the smallest contribution

GREEDY NETWORK PRUNING

Iterative pruning

Algorithm:

- 1) Estimate importance of neurons (units)
- 2) Rank units
- 3) Remove the least important unit
- 4) Fine tune network for K iterations
- 5) Go back to step1)



ORACLE

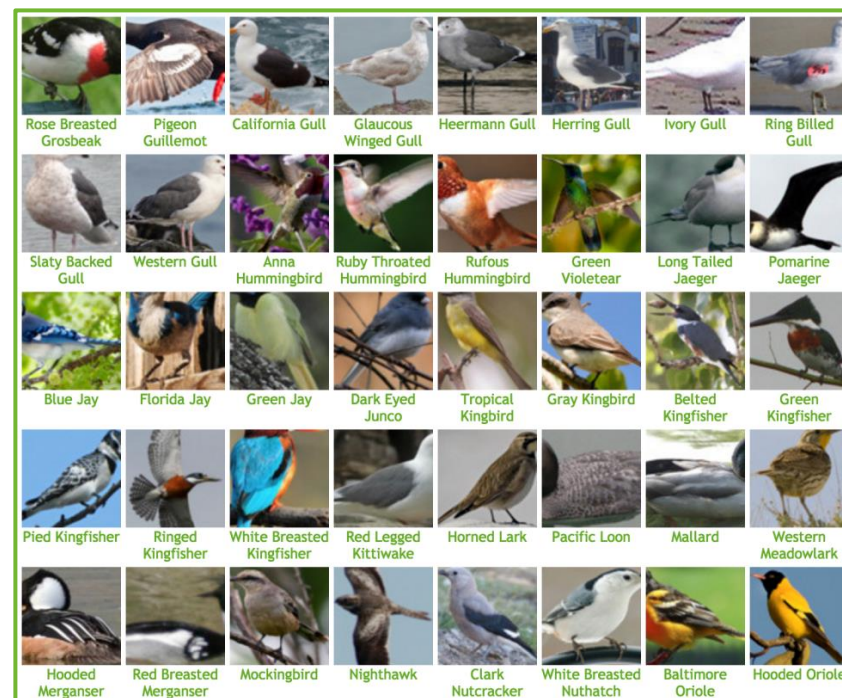
ORACLE

Caltech-UCSD Birds-200-2011 Dataset

- 200 classes
- <6000 training images

	Method	Test accuracy
S. Belongie et al	*SIFT+SVM	19%
From scratch	CNN	25%
S. Razavian et al	*OverFeat+SVM	62%
Our baseline	VGG16 finetuned	72.2%
N. Zhang et al	R-CNN	74%
S. Branson et al	*Pose-CNN	76%
J. Krause et al	*R-CNN+	82%

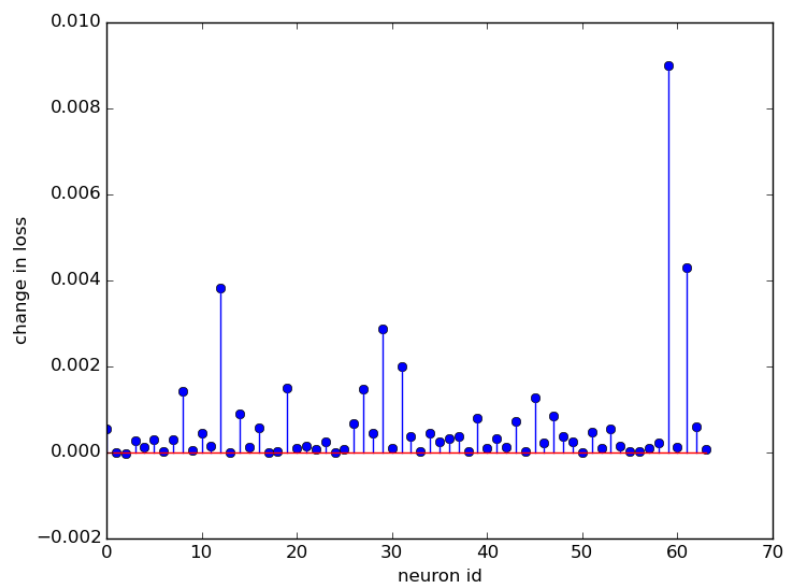
*require additional attributes



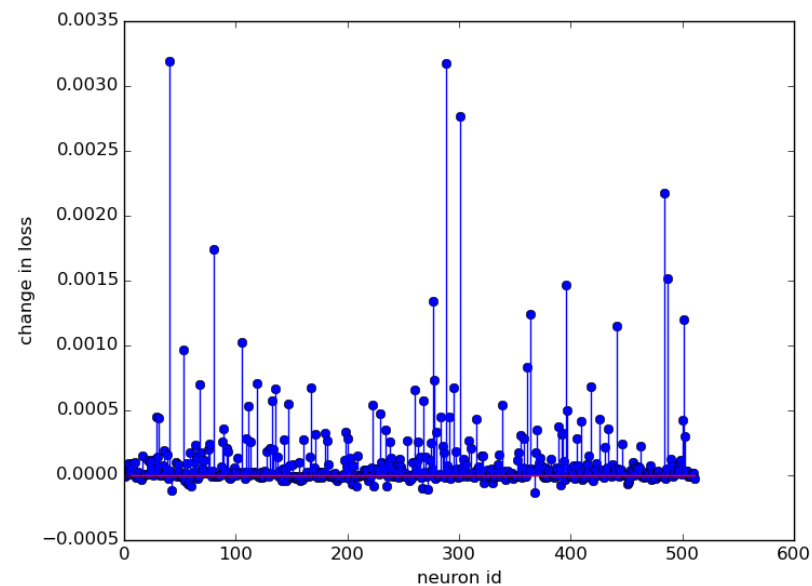
ORACLE

VGG16 on Birds-200 dataset

- Exhaustively computed change in loss by removing one unit



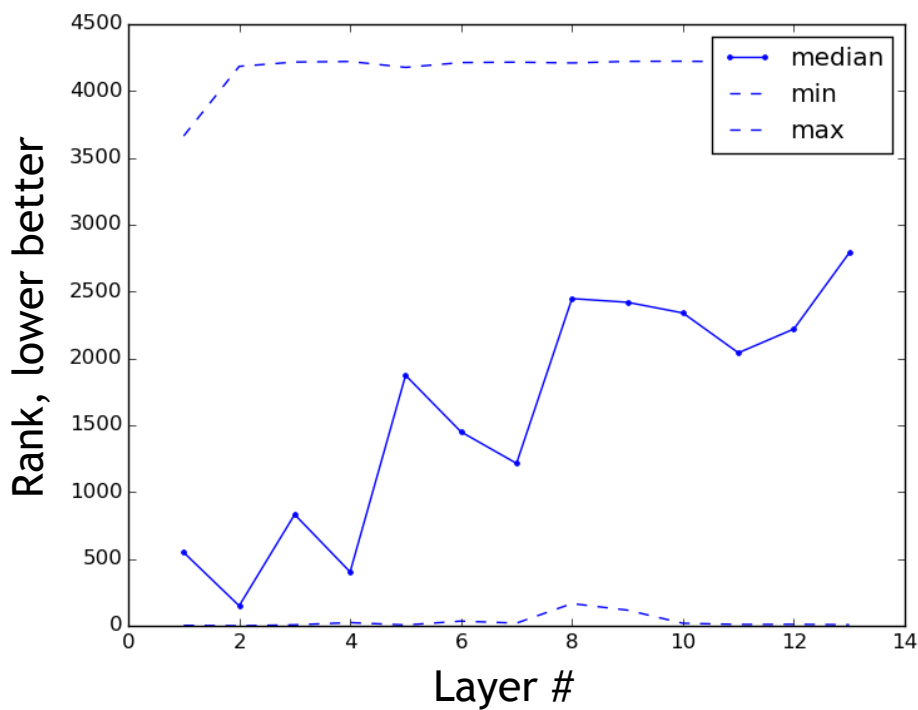
First layer



Last layer

ORACLE

VGG-16 on Birds-200



*only convolutional layers

- On average first layers are more important
- Every layer has **very important** units
- Every layer has **non important** units
- Layers with pooling are more important

APPROXIMATING THE ORACLE

APPROXIMATING THE ORACLE

Candidate criteria

- Average activation (discard lower activations)
- Minimum weight (discard lower l_2 of weight)
- With first-order Taylor expansion (TE):

$$\mathcal{C}(\mathcal{D}, h_i = 0) = \mathcal{C}(\mathcal{D}, h_i) - \frac{\delta \mathcal{C}}{\delta h_i} h_i + \underbrace{R_1(h_i = 0)}_{\text{ignore}}$$

Absolute difference in cost by removing a neuron:

$$|\Delta \mathcal{C}(h_i)| = |\mathcal{C}(\mathcal{D}, h_i) - \frac{\delta \mathcal{C}}{\delta h_i} h_i - \mathcal{C}(\mathcal{D}, h_i)| = \left| \frac{\delta \mathcal{C}}{\delta h_i} h_i \right|$$

$$\frac{\delta \mathcal{C}}{\delta h_i}$$

Gradient of the cost wrt.
activation h_i

$$h_i$$

Unit's output

$$+$$

Both computed during standard
backprop.

APPROXIMATING THE ORACLE

Candidate criteria

- Alternative: Optimal Brain Damage (OBD) by Y. LeCun et al., 1990
 - Use second order derivatives to estimate importance of neurons:

$$\Delta C(h_i) = \underbrace{\frac{\delta C}{\delta h_i} h_i}_{=0} + 0.5 \frac{\delta^2 C}{\delta h_i^2} + \underbrace{R_2(h_i = 0)}_{\text{ignore}}$$

— Needs extra comp of second order derivative

APPROXIMATING THE ORACLE

Comparison to OBD

OBD: second-order expansion:

$$\Delta C(h_i) = \underbrace{\frac{\delta C}{\delta h_i} h_i}_{=0} + 0.5 \frac{\delta^2 C}{\delta h_i^2}$$

we propose: abs of first-order expansion:

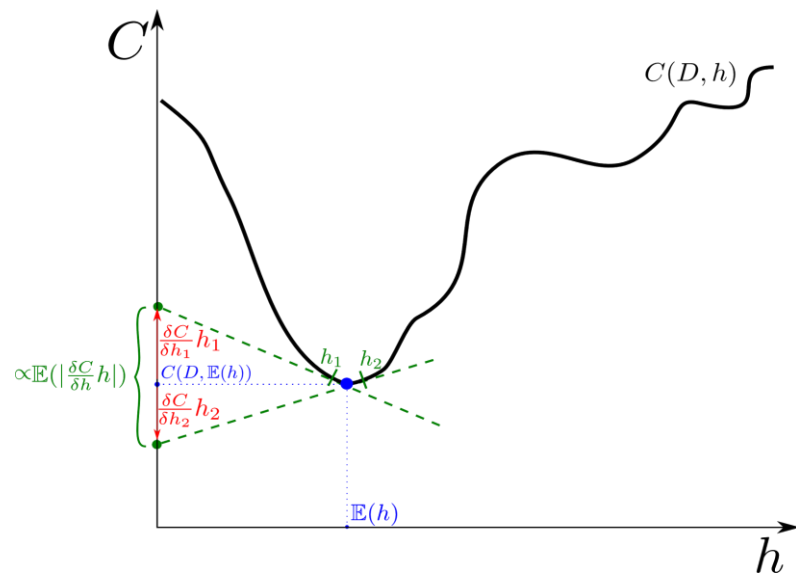
$$|\Delta C(h_i)| = \left| \frac{\delta C}{\delta h_i} h_i \right|$$

Assuming $y = \frac{\delta C}{\delta h_i} h_i$

For perfectly trained model:

$$\mathbb{E}(y) = 0$$

$$\mathbb{E}(|y|) = std(y) \sqrt{2}/\sqrt{\pi} \quad \text{if } y \text{ is Gaussian}$$



APPROXIMATING THE ORACLE

Comparison to OBD

OBD: second-order expansion:

$$\Delta C(h_i) = \underbrace{\frac{\delta C}{\delta h_i} h_i}_{=0} + 0.5 \frac{\delta^2 C}{\delta h_i^2}$$

we propose: abs of first-order expansion:

$$|\Delta \mathcal{C}(h_i)| = \left| \frac{\delta \mathcal{C}}{\delta h_i} h_i \right|$$

Assuming $y = \frac{\delta \mathcal{C}}{\delta h_i} h_i$

For perfectly trained model:

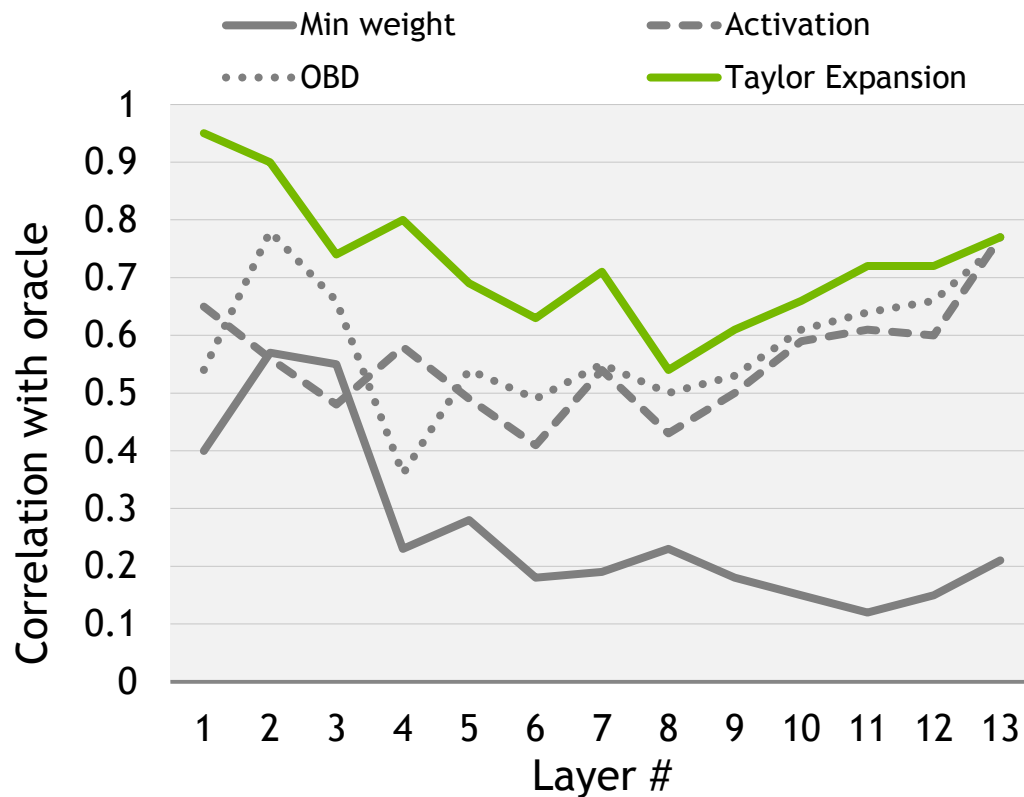
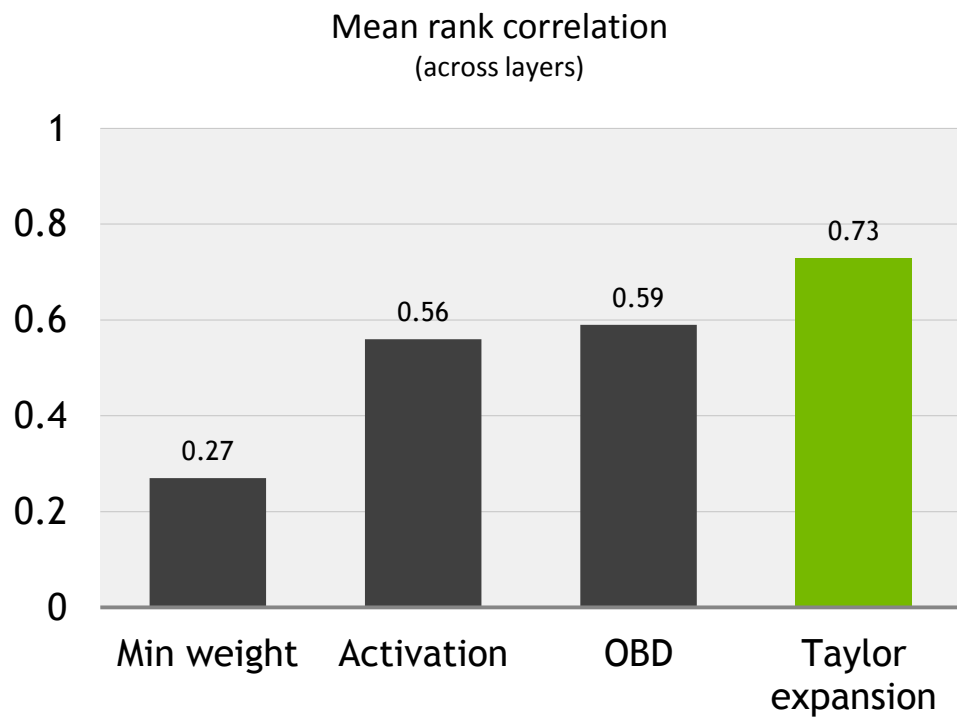
$$\mathbb{E}(y) = 0$$

$$\mathbb{E}(|y|) = std(y) \sqrt{2} / \sqrt{\pi} \quad \text{if } y \text{ is Gaussian}$$

- ✓ No extra computations
- ✓ We look at absolute difference
- Can't predict exact change in loss

EVALUATING PRUNING CRITERIA

Spearman's rank correlation with oracle: VGG16 on Birds-200



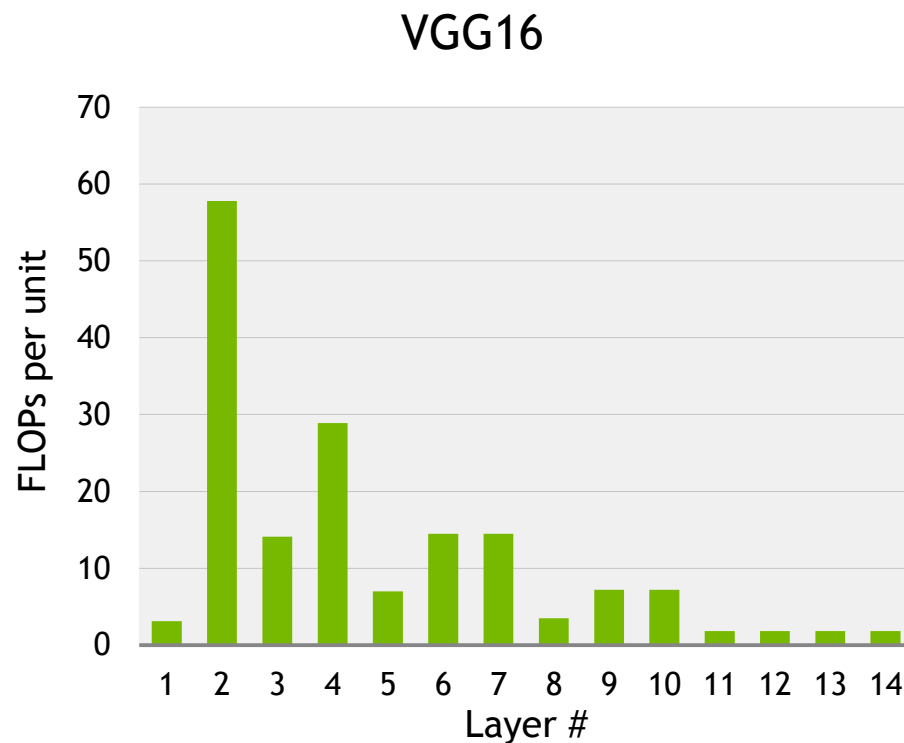
EVALUATING PRUNING CRITERIA

Pruning with objective

- Regularize criteria with objective:

$$\Theta(h) = \Theta(h) - \lambda \mathcal{R}(h)$$

- Regularizer can be:
 - FLOPs
 - Memory
 - Bandwidth
 - Target device
 - Exact inference time

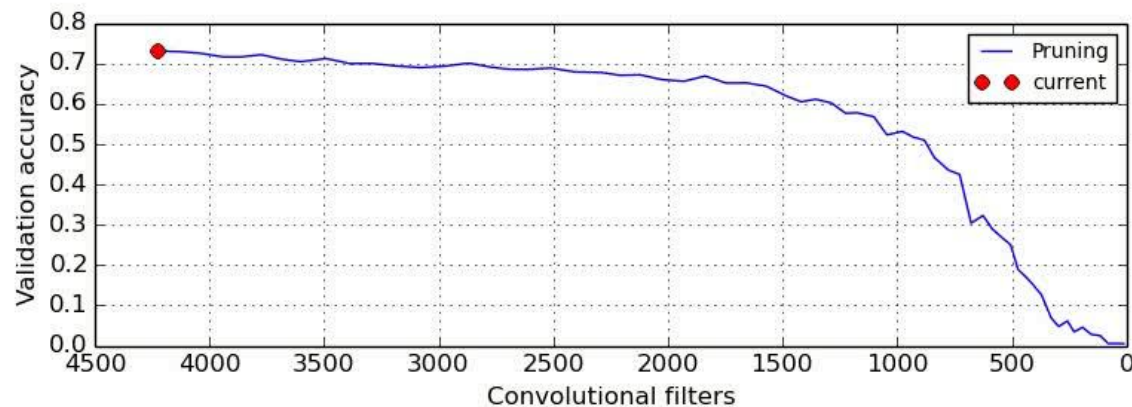
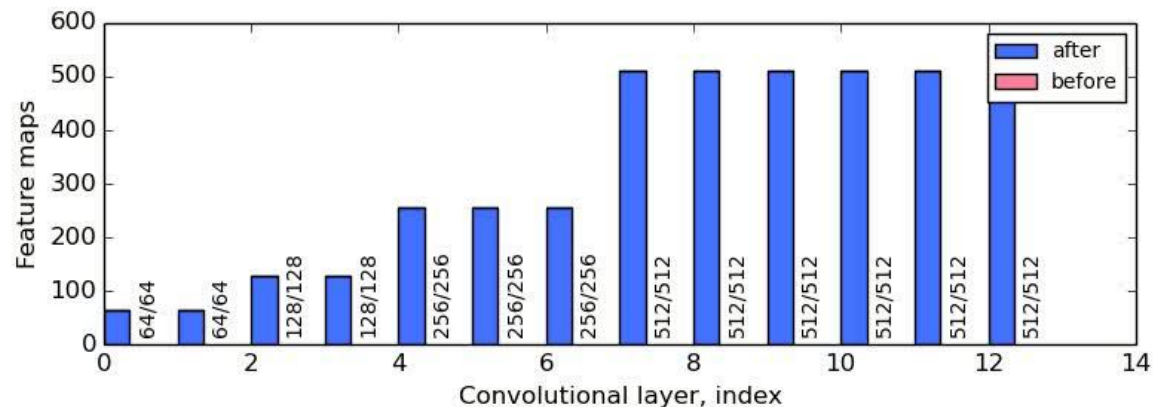


RESULTS

RESULTS

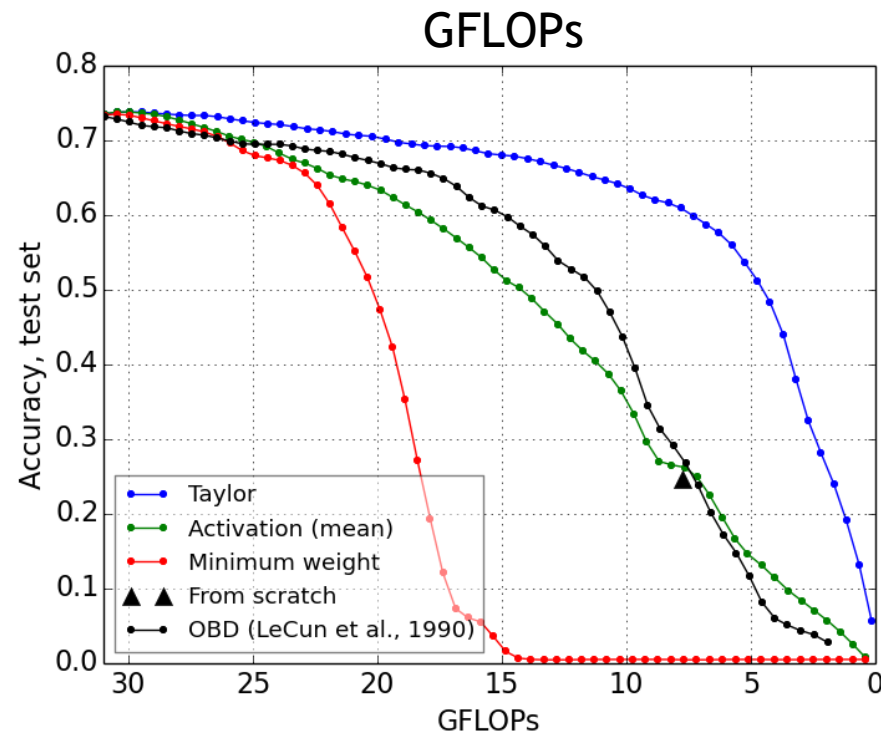
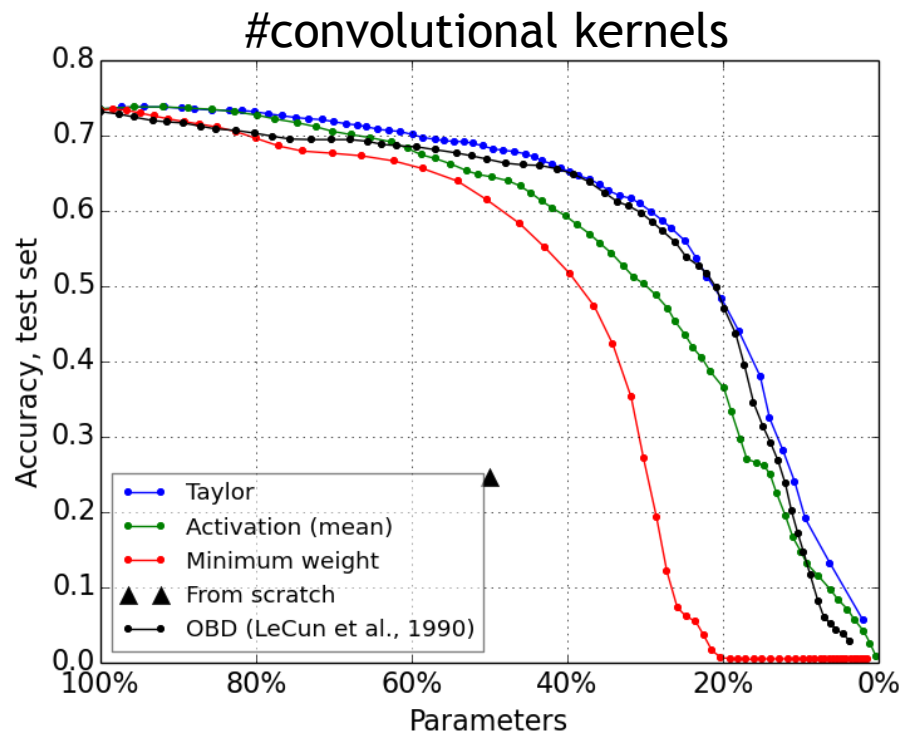
VGG16 on Birds 200 dataset

- Remove 1 conv unit every 30 updates



RESULTS

VGG16 on Birds 200 dataset



- Training from scratch doesn't work
- Taylor shows the best result vs any other metric for pruning

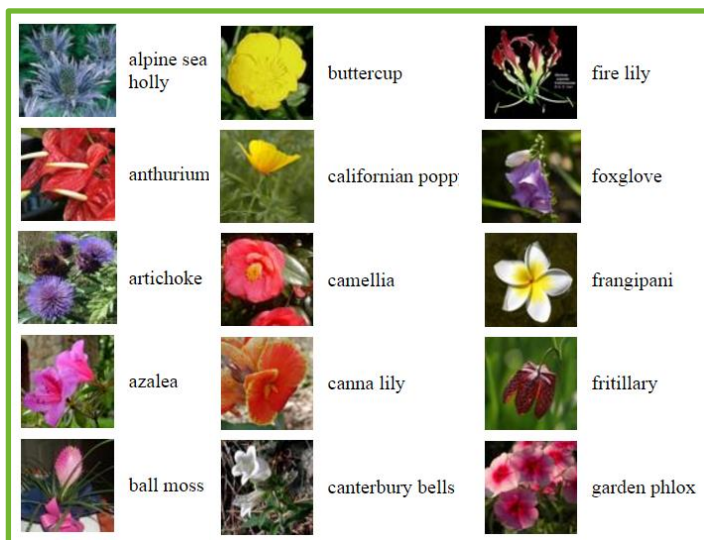
RESULTS

AlexNet on Oxford Flowers102

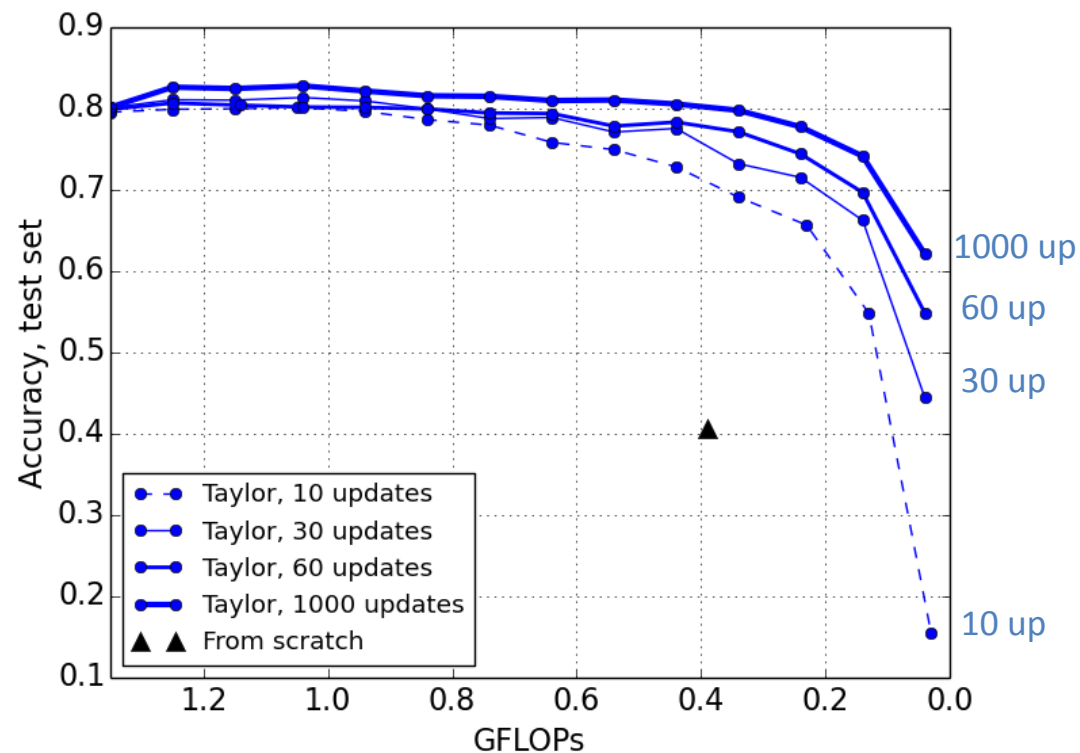
102 classes

~2k training images

~6k testing images



Changing number of updates between pruning iterations



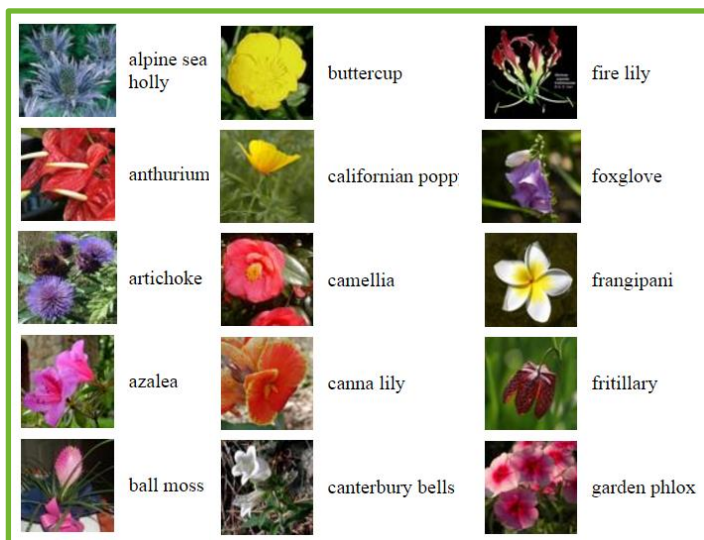
RESULTS

AlexNet on Oxford Flowers102

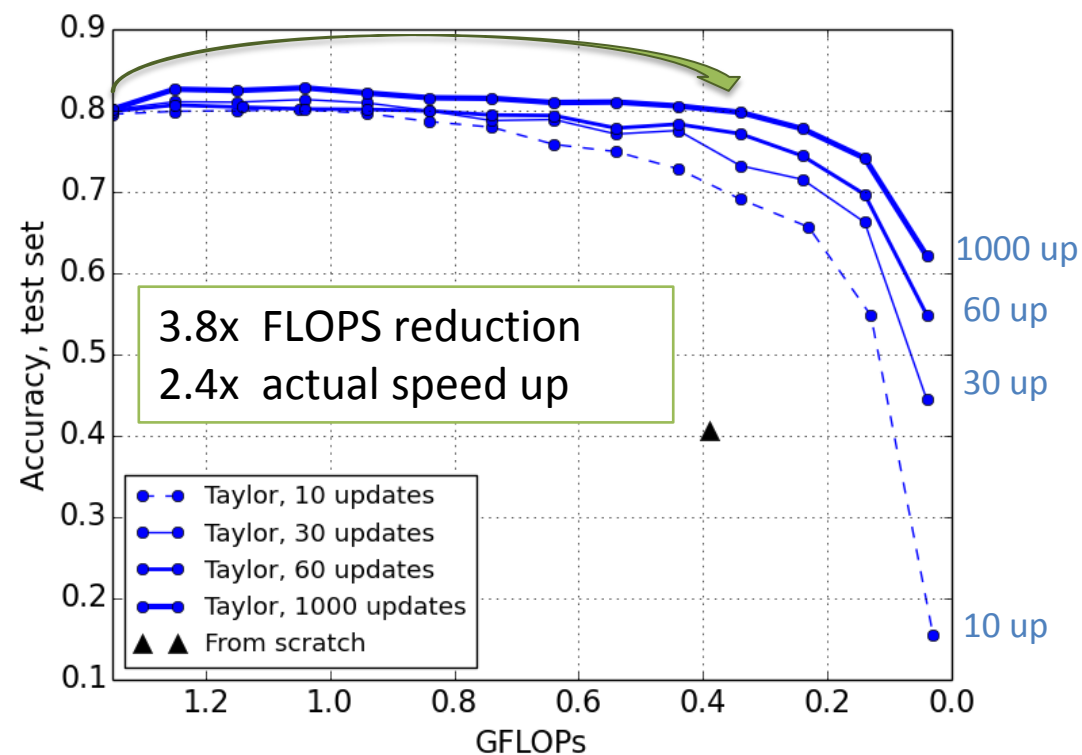
102 classes

~2k training images

~6k testing images



Changing number of updates between pruning iterations



VGG16 on ImageNet

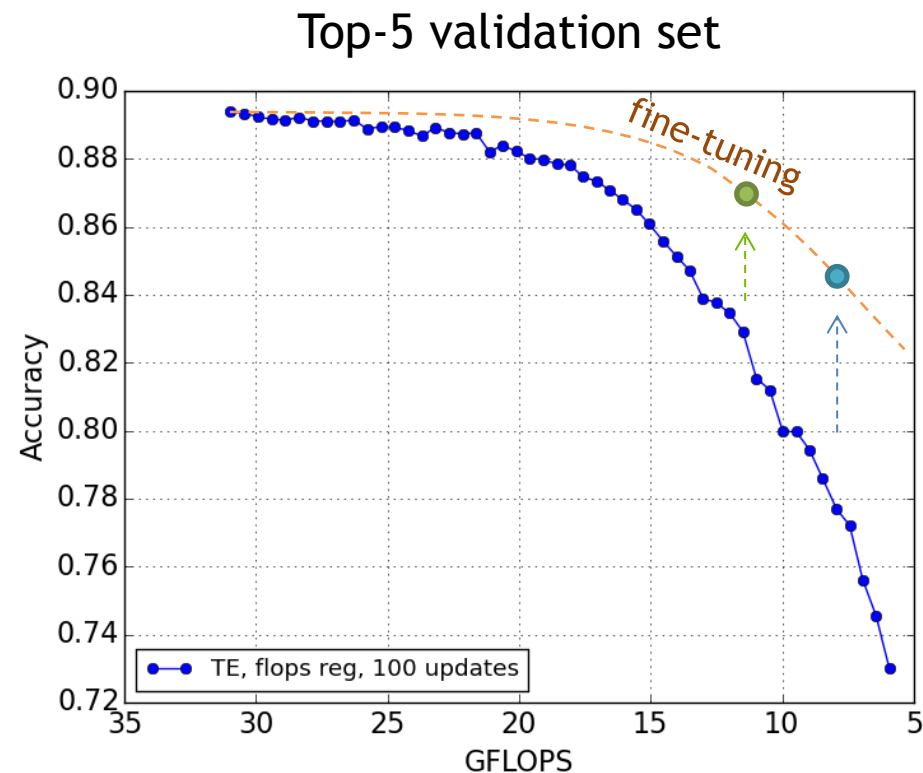
-
- Top-5 validation set
- Accuracy
- GFLOPS
- TE, flops reg, 100 updates
- | GFLOPS | Accuracy |
|--------|----------|
| 31.5 | 0.893 |
| 31.0 | 0.892 |
| 30.5 | 0.891 |
| 30.0 | 0.891 |
| 29.5 | 0.891 |
| 29.0 | 0.891 |
| 28.5 | 0.891 |
| 28.0 | 0.891 |
| 27.5 | 0.891 |
| 27.0 | 0.891 |
| 26.5 | 0.891 |
| 26.0 | 0.891 |
| 25.5 | 0.891 |
| 25.0 | 0.891 |
| 24.5 | 0.891 |
| 24.0 | 0.891 |
| 23.5 | 0.891 |
| 23.0 | 0.891 |
| 22.5 | 0.891 |
| 22.0 | 0.891 |
| 21.5 | 0.891 |
| 21.0 | 0.891 |
| 20.5 | 0.891 |
| 20.0 | 0.891 |
| 19.5 | 0.891 |
| 19.0 | 0.891 |
| 18.5 | 0.891 |
| 18.0 | 0.891 |
| 17.5 | 0.891 |
| 17.0 | 0.891 |
| 16.5 | 0.891 |
| 16.0 | 0.891 |
| 15.5 | 0.891 |
| 15.0 | 0.891 |
| 14.5 | 0.891 |
| 14.0 | 0.891 |
| 13.5 | 0.891 |
| 13.0 | 0.891 |
| 12.5 | 0.891 |
| 12.0 | 0.891 |
| 11.5 | 0.891 |
| 11.0 | 0.891 |
| 10.5 | 0.891 |
| 10.0 | 0.891 |
| 9.5 | 0.891 |
| 9.0 | 0.891 |
| 8.5 | 0.891 |
| 8.0 | 0.891 |
| 7.5 | 0.891 |
| 7.0 | 0.891 |
| 6.5 | 0.891 |
| 6.0 | 0.891 |
| 5.5 | 0.891 |
| 5.0 | 0.891 |

RESULTS

VGG16 on ImageNet

- Pruned over 7 epochs
- Fine-tuning 7 epochs

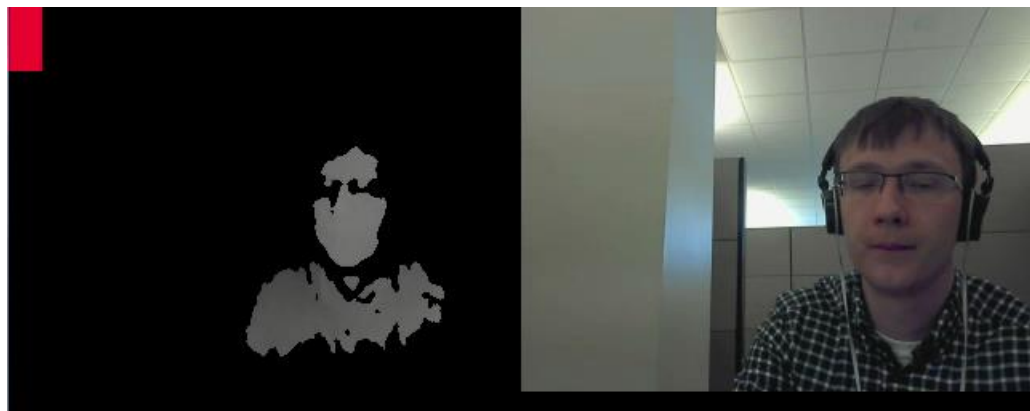
	GFLOPs	FLOPS reduction	Actual speed up	Top-5
	31	1x		89.5%
●	12	2.6x	2.5x	-2.5%
●	8	3.9x	3.3x	-5.0%



RESULTS

R3DCNN for gesture recognition

3D-CNN with recurrent layers fine-tuned for 25 dynamic gestures



RESULTS

R3DCNN for gesture recognition

3D-CNN with recurrent layers fine-tuned for 25 dynamic gestures

12.6x

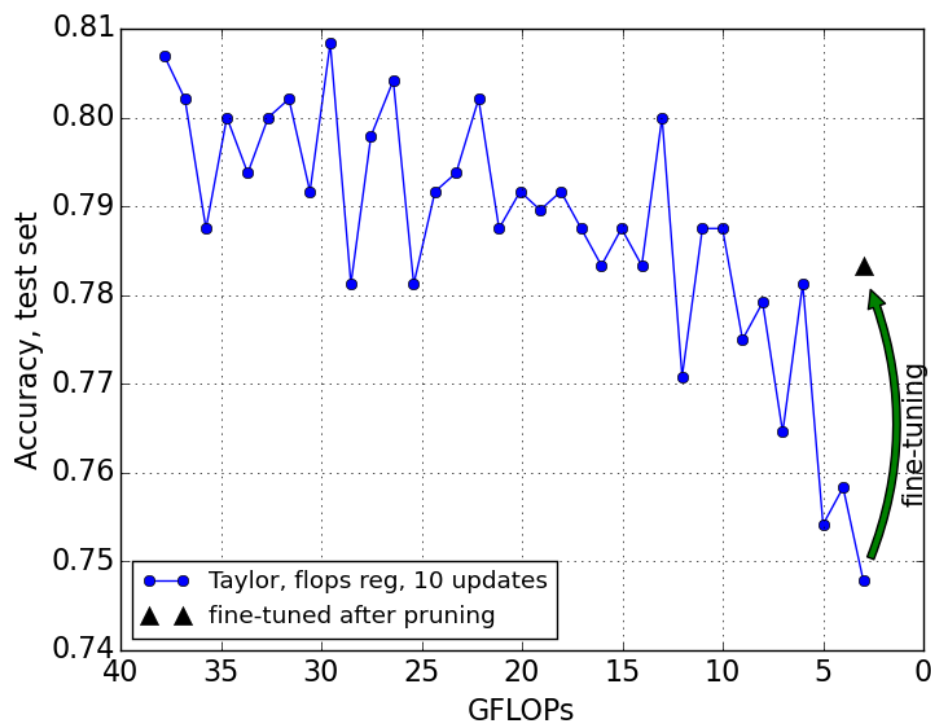
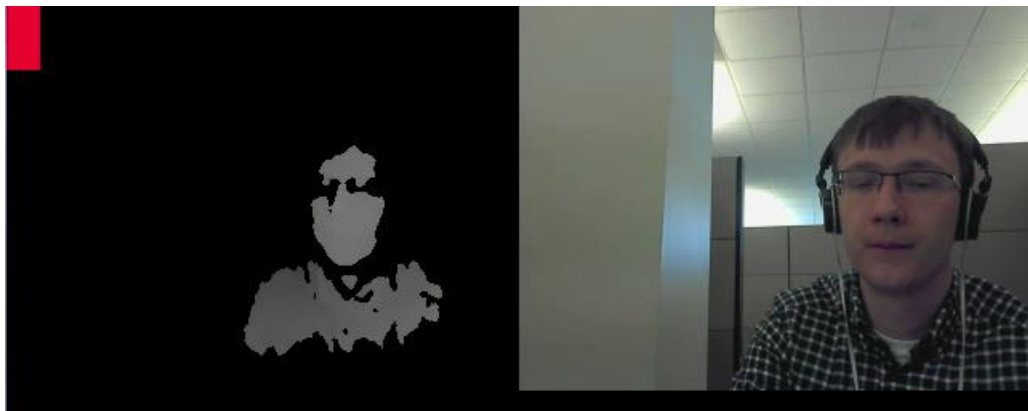
Reduction in FLOPs

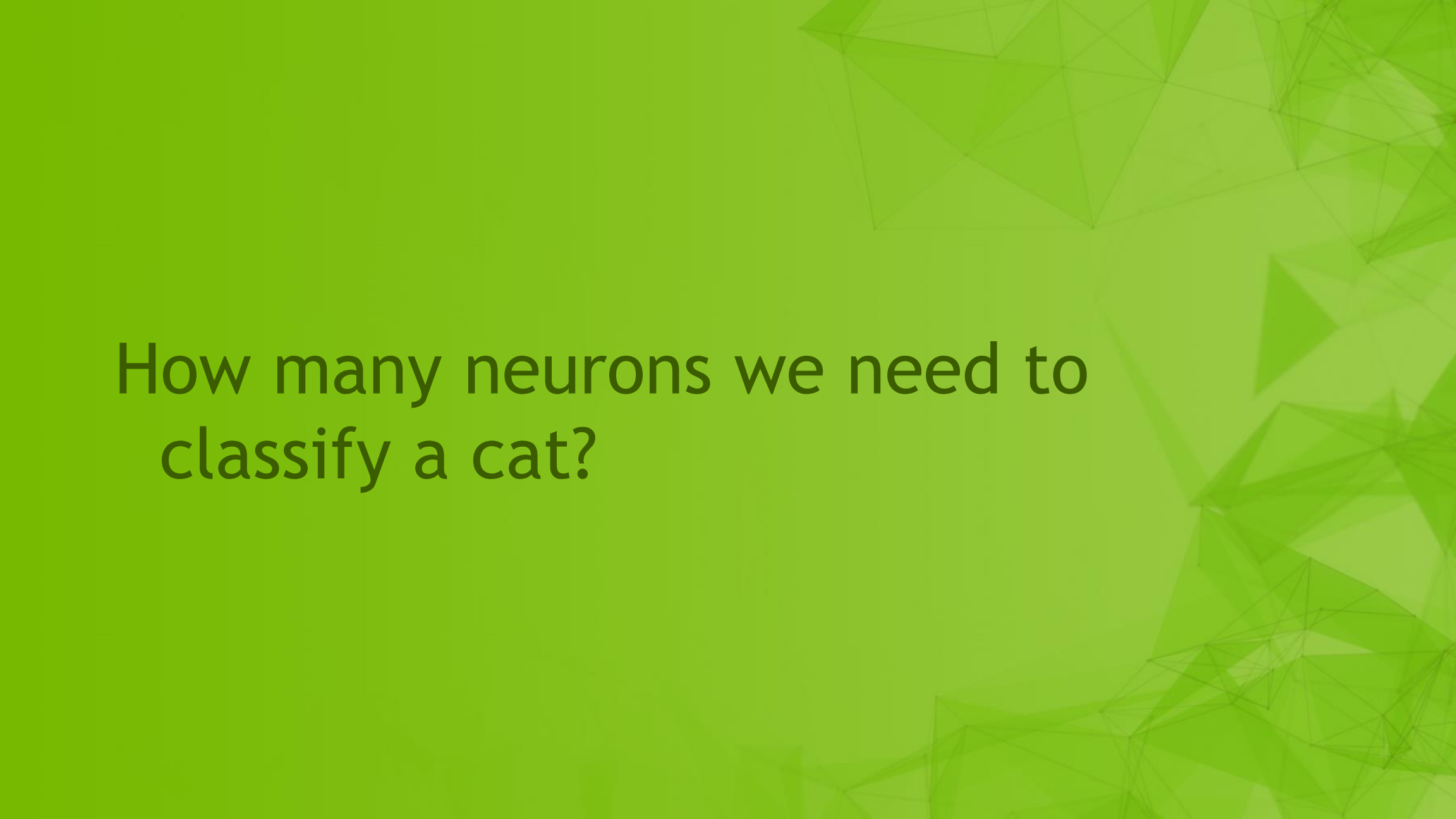
Drop in accuracy

2.5%

5.2x

Speed-up





How many neurons we need to
classify a cat?

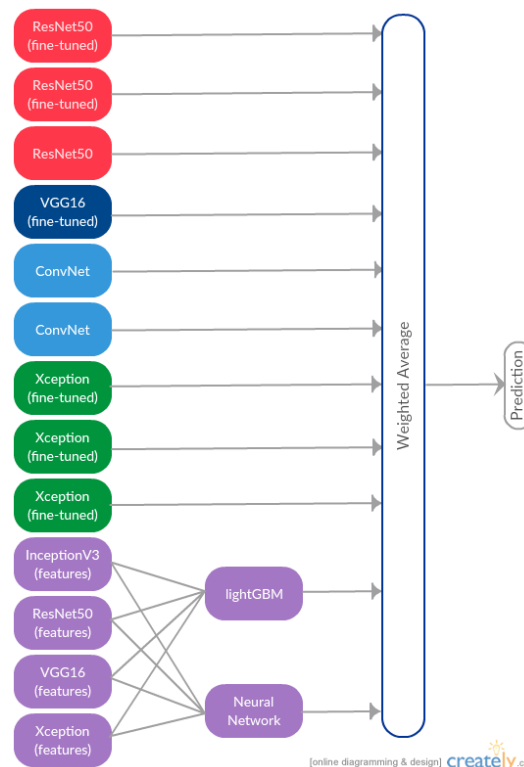
DOGS VS. CATS

Dogs vs. Cats classification



25,000 images

Marco Lugo's solution, 3rd place :



DOGS VS. CATS

Fine-tuned ResNet-101

Full network

99.2%

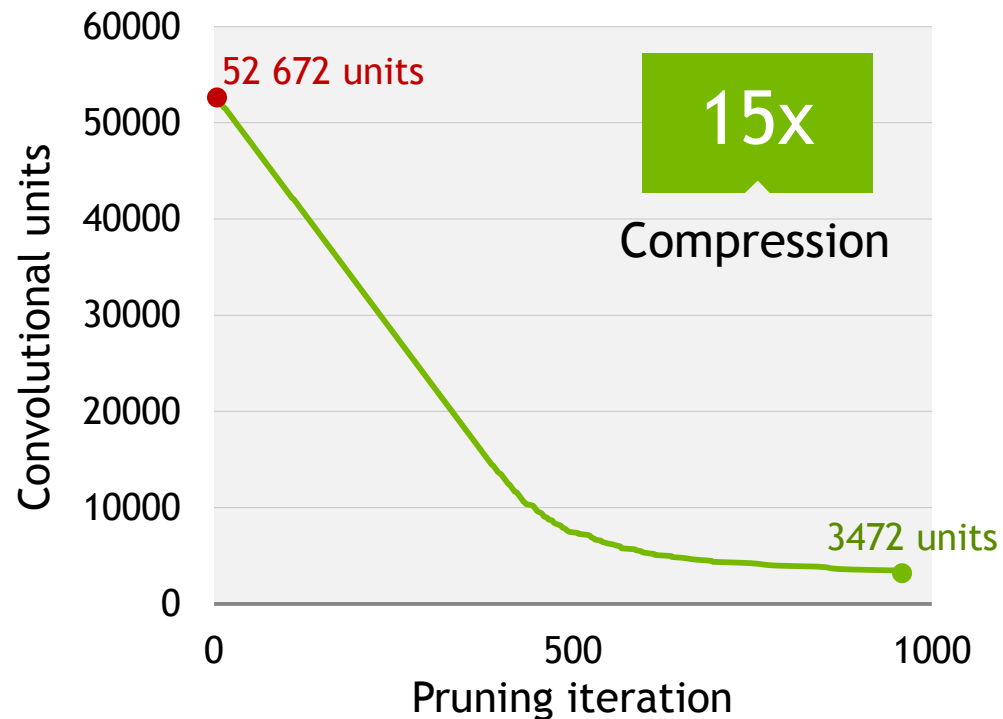
Pruned network

99.0 %



DOGS VS. CATS

Fine-tuned ResNet-101



Full network

99.2%

Pruned network

99.0 %



CONCLUSIONS

- Pruning as greedy feature selection
- New criteria based on Taylor expansion
- Pruning is especially effective (and necessary!) for transfer learning
- Pruning can incorporate desired objectives (such as FLOPs)
- Read more in our ICLR2017 paper: <https://openreview.net/pdf?id=SJGCiw5gl>

THANK YOU!

