

Dense for the Price of Sparse: Training Extremely Sparse Networks from Scratch with Random Sparse Support

Ilan Price^{1,2} & Jared Tanner^{1,2}

MLC DLCT Presentation
1 Oct 2021

¹Mathematical Institute, University of Oxford
²The Alan Turing Institute



**The
Alan Turing
Institute**



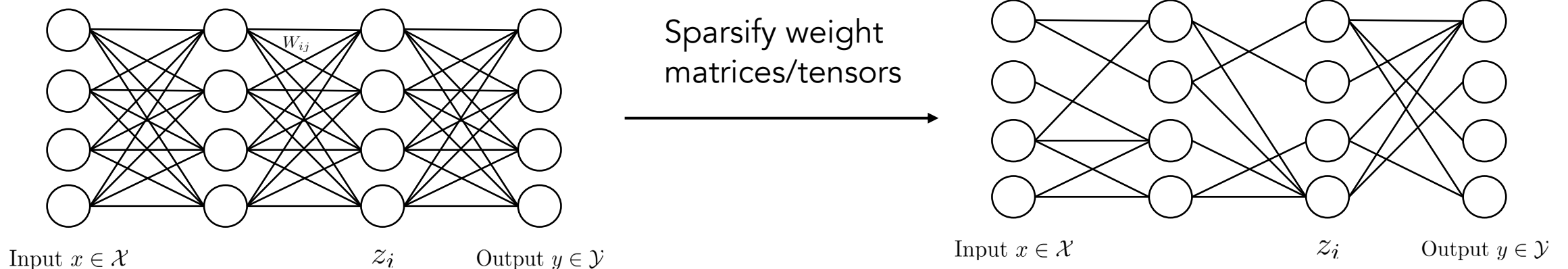
Outline

- Context: sparse deep networks and the goal of sparse training
- Sparse training vs (?) subspace training
- DCTpS: computationally efficient random subspace training
- Some results
- Limitations and open questions



Sparse Deep Networks

- 'Sparse' deep learning can refer to multiple different things (See Hoefler et al¹ for a full review).
- Here we consider persistent (fixed for all inputs) sparsity of the weights



Sparse Deep Networks

- Deep networks are most-often vastly over parameterised. Parameter counts now range from $\mathcal{O}(10^6)$ to $\mathcal{O}(10^{12})$!
- We have known for a long time that we can “prune” most of these while maintaining good accuracy

Huge storage and computational savings (in theory at least)



Sparse Deep Networks

- Deep networks are most-often vastly over parameterised. Parameter counts now range from $\mathcal{O}(10^6)$ to $\mathcal{O}(10^{12})$!
- We have known for a long time that we can “prune” most of these while maintaining good accuracy

Huge storage and computational savings (in theory at least)

- The most consistently successful methods: pruning during and/or after training (followed by some fine tuning) – e.g. Iterative Magnitude Pruning



Sparse Deep Networks

- Deep networks are most-often vastly over parameterised. Parameter counts now range from $\mathcal{O}(10^6)$ to $\mathcal{O}(10^{12})$!
- We have known for a long time that we can “prune” most of these while maintaining good accuracy

Huge storage and computational savings (in theory at least)

- The most consistently successful methods: pruning during and/or after training (followed by some fine tuning) – e.g. Iterative Magnitude Pruning
- Can we prune before training? So storage and compute is cheaper during training too?



Static sparse training¹

- Naïve (uniform) random sparsification performs poorly at extreme sparsities



Static sparse training¹

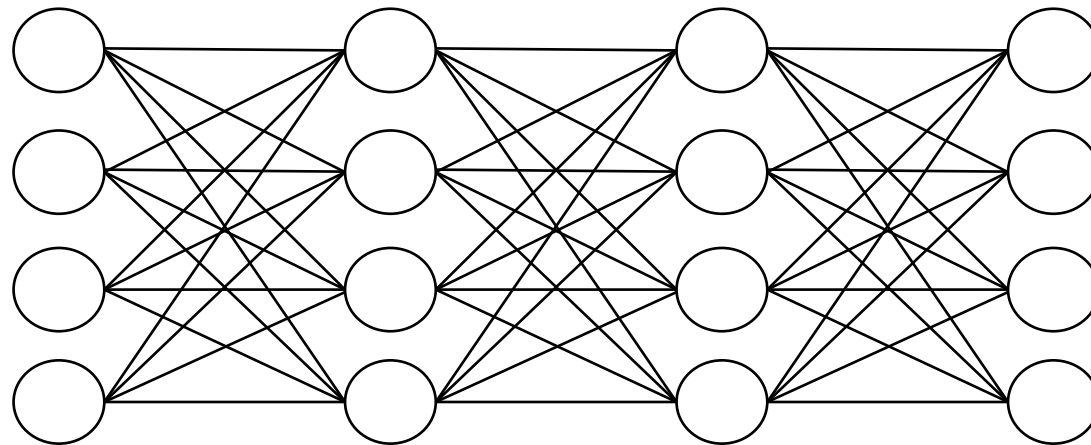
- Naïve (uniform) random sparsification performs poorly at extreme sparsities
- The lottery ticket hypothesis ([Frankle and Carbin, 2019](#)):

A randomly-initialized, dense neural network contains a subnetwork that is initialized such that—when trained in isolation—it can match the test accuracy of the original network after training for at most the same number of iterations.



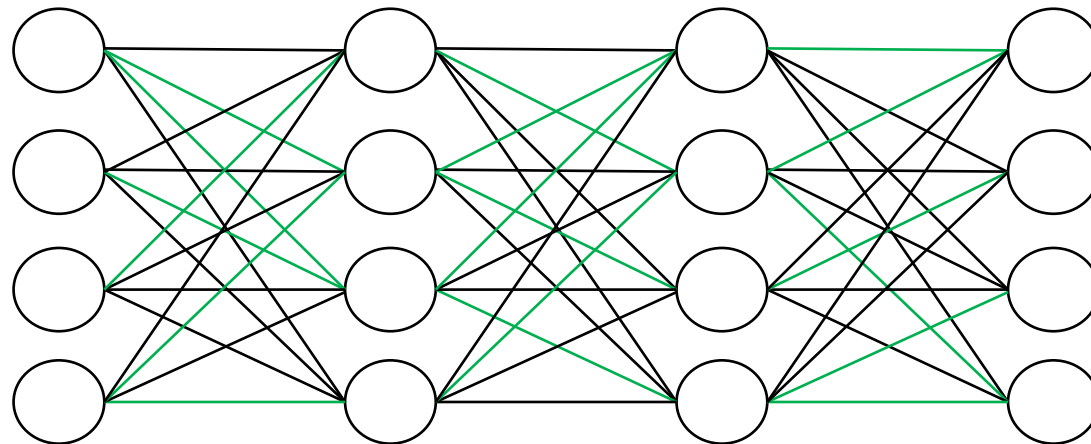
Static sparse training¹

- Naïve (uniform) random sparsification performs poorly at extreme sparsities
- The lottery ticket hypothesis ([Frankle and Carbin, 2019](#)):



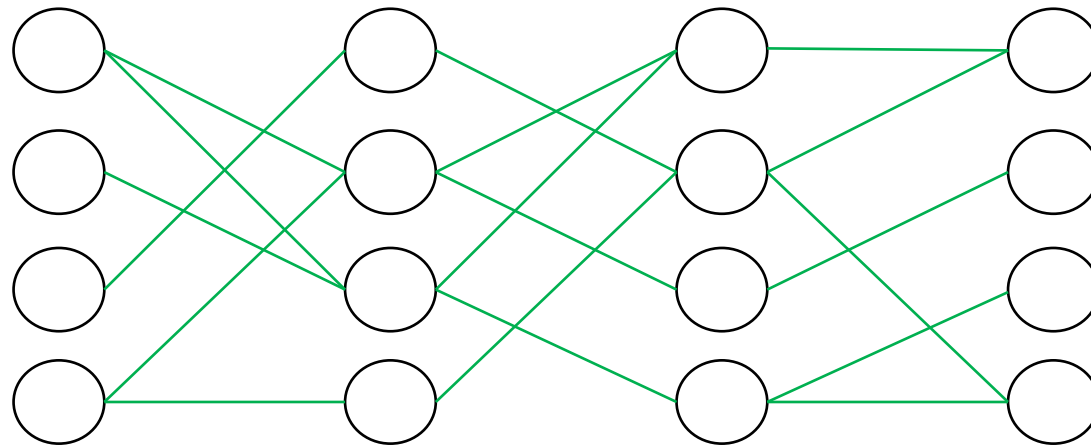
Static sparse training¹

- Naïve (uniform) random sparsification performs poorly at extreme sparsities
- The lottery ticket hypothesis ([Frankle and Carbin, 2019](#)):



Static sparse training¹

- Naïve (uniform) random sparsification performs poorly at extreme sparsities
- The lottery ticket hypothesis ([Frankle and Carbin, 2019](#)):



Static sparse training¹

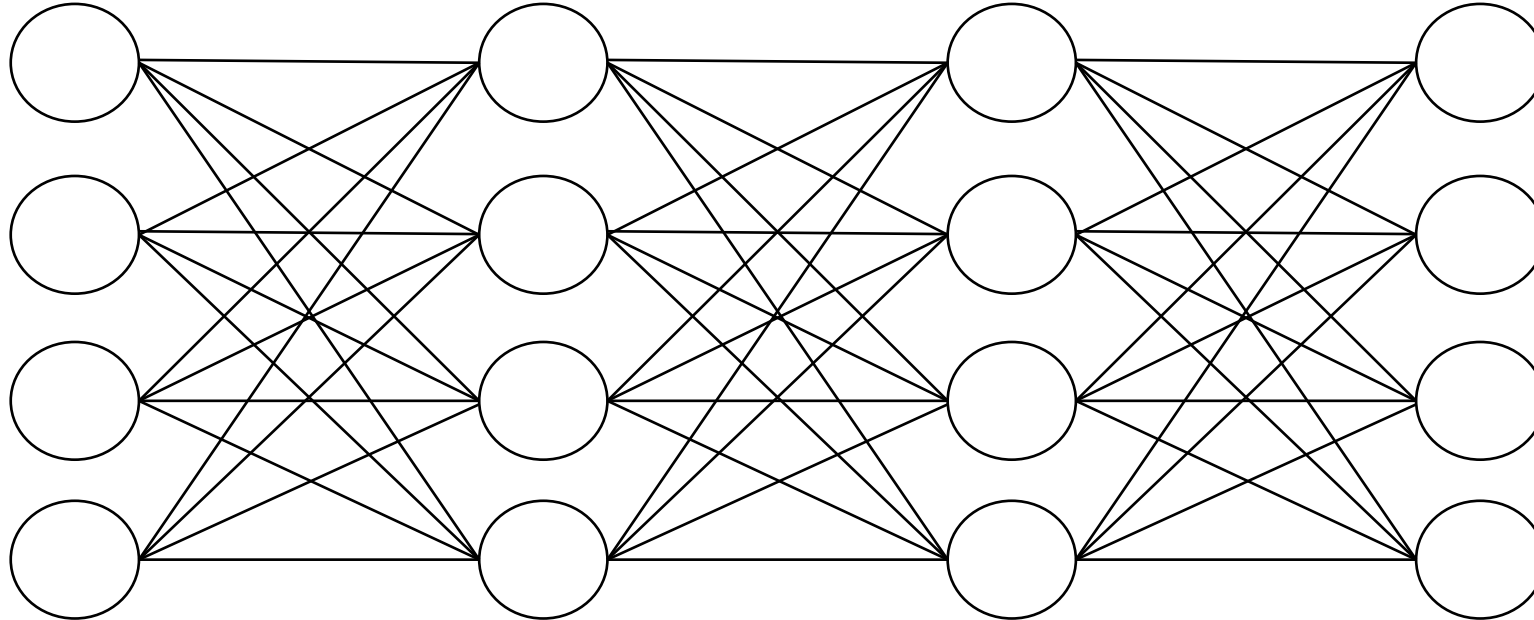
- Naïve (uniform) random sparsification performs poorly at extreme sparsities
- The lottery ticket hypothesis ([Frankle and Carbin, 2019](#)):

A randomly-initialized, dense neural network contains a subnetwork that is initialized such that—when trained in isolation—it can match the test accuracy of the original network after training for at most the same number of iterations.

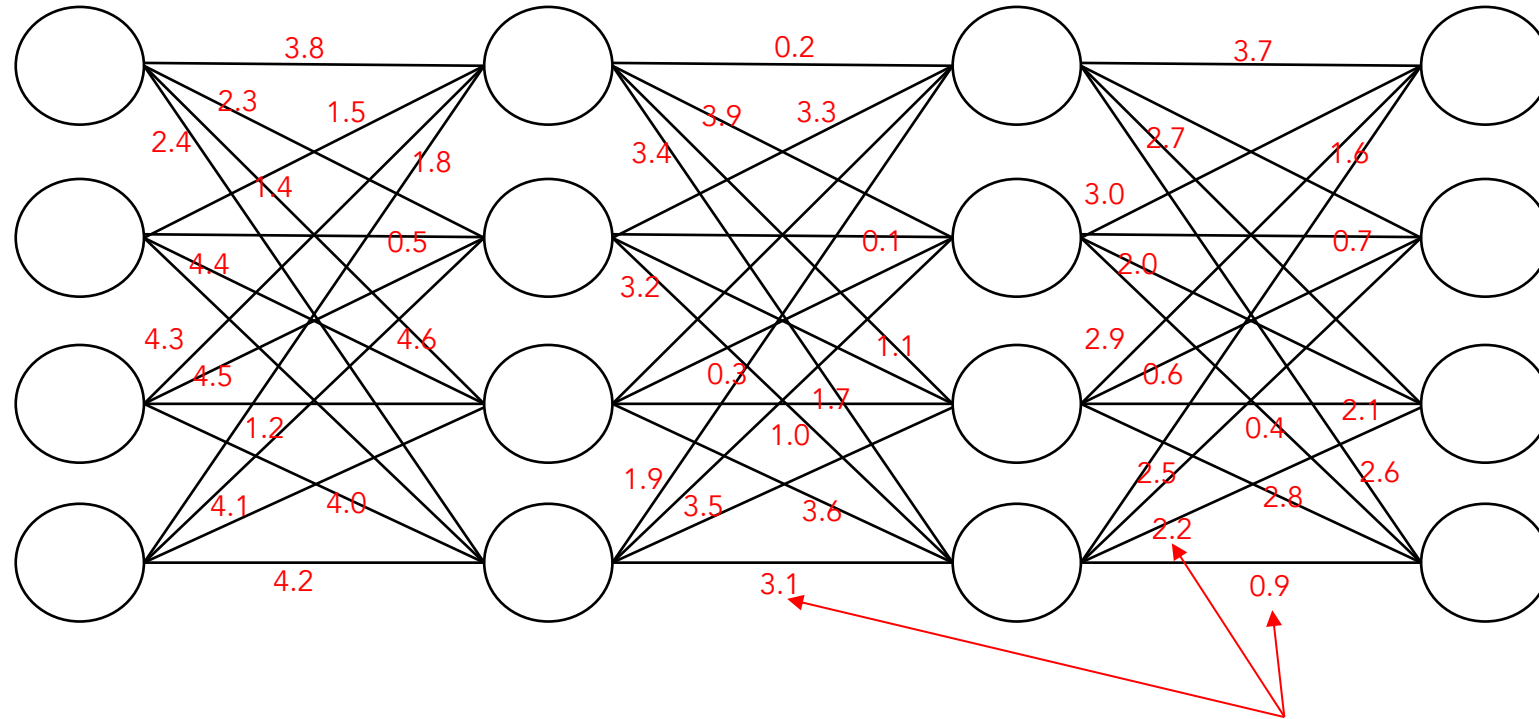
- How to find trainable, extremely sparse sub-networks in practice?



Pruning at initialization (PaI)



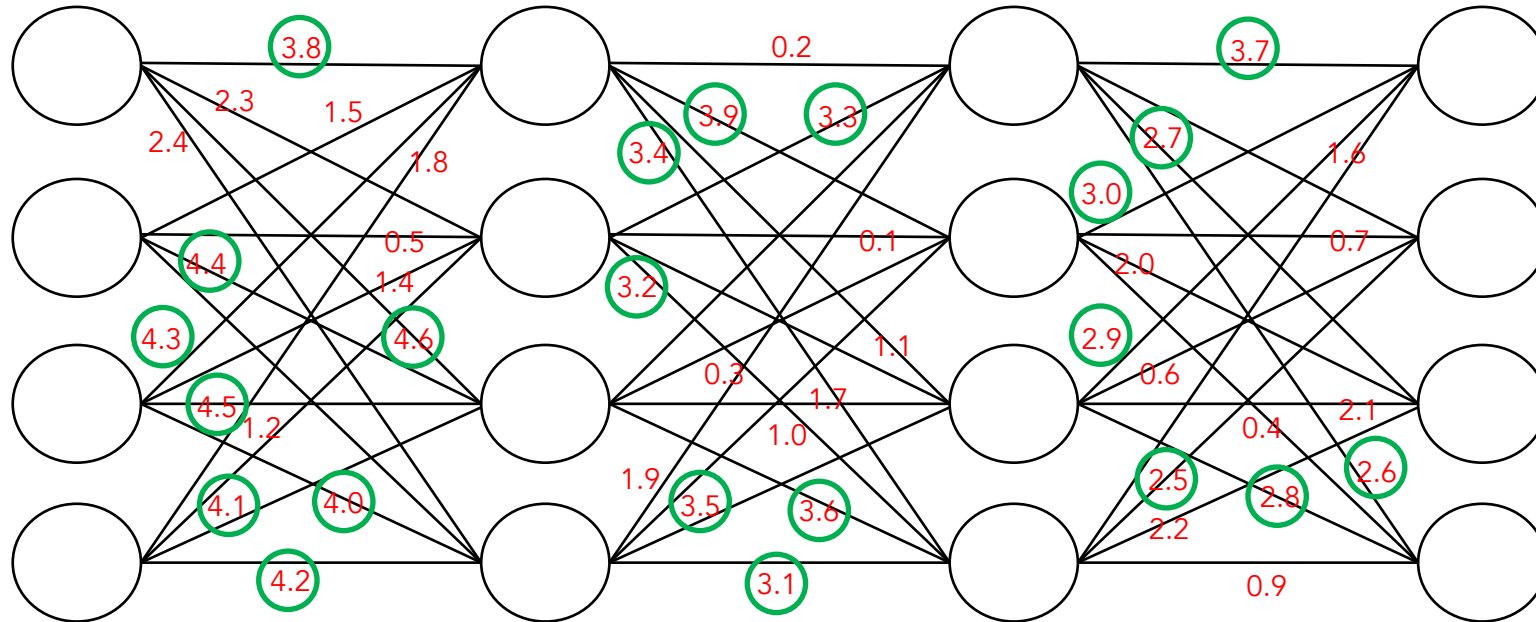
Pruning at initialization (Pal)



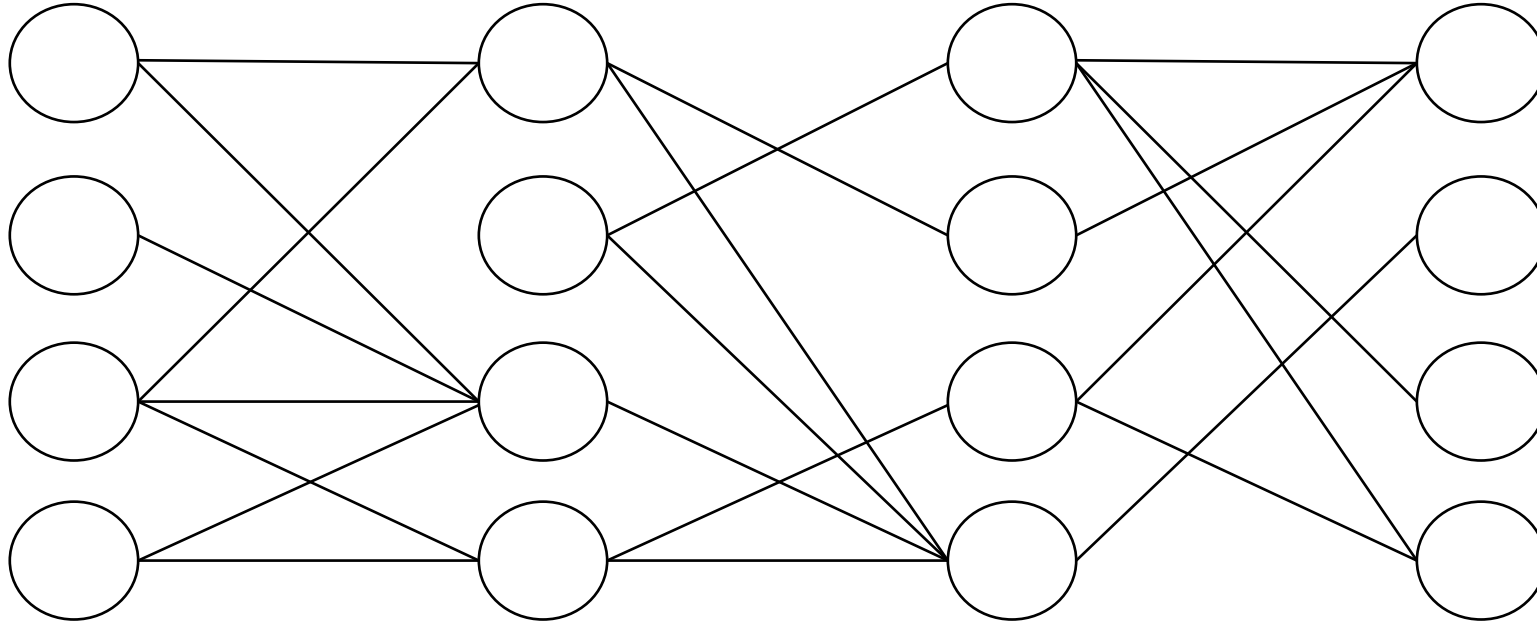
Saliency scores
for each weight



Pruning at initialization (Pal)



Pruning at initialization (PaI)



Standard Pruning At Initialisation (PaI)

Generic steps:

1. Initialize a dense network
2. Define scalar objective \mathcal{R}
3. Calculate vector of saliency scores $G(\mathbf{w}) = \frac{\partial \mathcal{R}}{\partial \mathbf{w}} \odot \mathbf{w}$
4. Prune parameters with lowest scores

SOTA
at Extreme
Sparsities →

$$\text{FORCE}^1: G(\mathbf{w}) = \left| \frac{\partial \mathcal{L}(\bar{\mathbf{w}})}{\partial \mathbf{w}} \odot \mathbf{w} \right|$$

$\bar{\mathbf{w}}$ is the param vector post-pruning

$$\text{SynFlow}^2: \mathcal{R} = \mathbf{1}^\top \left(\prod_{l=1}^L |\mathbf{w}^{[l]}| \right) \mathbf{1}$$

$|\mathbf{w}^{[l]}|$ is the element-wise absolute value of the parameters in the l^{th} layer

¹ de Jorge, Pau, et al. "Progressive skeletonization: Trimming more fat from a network at initialization." 2020

² Tanaka, Hidenori, et al. "Pruning neural networks without any data by iteratively conserving synaptic flow.", 2020.



Pruning at initialization (Pal)

How important are these saliency scores at initialization?

1. After Pal we can often reshuffle the locations of the weights within layers and still train to the same accuracy²



Pruning at initialization (Pal)

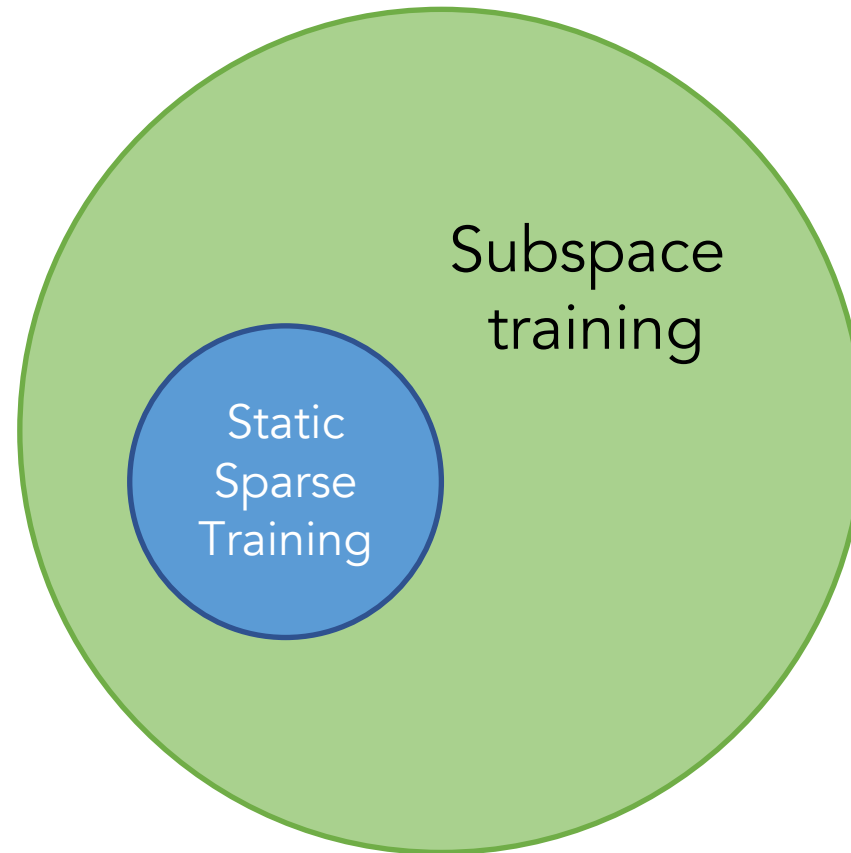
How important are these saliency scores at initialization?

1. After Pal we can often reshuffle the locations of the weights within layers and still train to the same accuracy²
2. Results from other literature on low effective dimensionality in network training

² Frankle et al. *Pruning Neural Networks at Initialization: Why Are We Missing the Mark?* ICLR2021



Zooming out again...



Subspace Training

Fixed (untrainable)
offset

Network weights

$$\mathbf{w} = \mathbf{d} + U\theta$$

Trainable parameters

Fixed Subspace embedding

$$\mathbf{w}, \mathbf{d} \in \mathbb{R}^N$$
$$\theta \in \mathbb{R}^k, k \ll N$$
$$U \in \mathbb{R}^{N \times k}$$



Subspace Training

Sparse Networks:

zero vector

$$\mathbf{w} = \mathbf{d} + U\theta$$

“k-sparse disjoint”:

- 1 non-zero per column
- ≤ 1 non-zero per row

Pal \rightarrow Finding the right such U



Subspace Training

(Dense) low-dimensional networks² :

Randomly sampled

$$\mathbf{w} = \mathbf{d} + U\theta$$

$U \sim$ e.g. Gaussian

$\mathbf{d} \sim$ standard NN init

Excellent performance with extremely few trainable parameters...
but still dense



Interlude: a “Trend”?

- *A Generalized Lottery Ticket Hypothesis, Alabdulmohsin, Tolstikhin et al. 2021*

“We introduce a generalization to the lottery ticket hypothesis in which the notion of “sparsity” is relaxed by choosing an arbitrary basis in the space of parameters.”

- *How many degrees of freedom do we need to train deep networks: a loss landscape perspective, Larson et al, 2021*

“recent works, spanning pruning, lottery tickets, and training within random subspaces, have shown that deep neural networks can be trained using far fewer degrees of freedom than the total number of parameters”



Where were we...

Accuracy: random subspace \gg "random pruning" subspace

Compute/Storage: random subspace \ll "random pruning" subspace

Try and get the best of both worlds: efficiency of sparse nets with random subspace selection.



Where were we...

Accuracy: random subspace \gg "random pruning" subspace

Compute/Storage: random subspace \ll "random pruning" subspace

Try and get the best of both worlds: efficiency of sparse nets with random subspace selection.

Important features:

1. Random subspace training \rightarrow Offset from the origin
2. Random subspace training & pruning at init \rightarrow "Layer-wise distribution" of trainable parameters.



Best of both: Dense for the Price of Sparse

$$\mathbf{w} = \mathbf{d} + U\theta$$

Dense but fixed (untrainable)

"k-sparse disjoint":

- 1 non-zero per column
- ≤ 1 non-zero per row



DCT plus Sparse

$$\mathbf{w} = \mathbf{d} + U\theta$$

$$\begin{pmatrix} | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \end{pmatrix} = \begin{pmatrix} | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \end{pmatrix} + \begin{pmatrix} 0 \dots 0 1 0 \dots 0 \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ 0 \dots 0 1 0 \dots 0 \end{pmatrix} \begin{pmatrix} | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \\ | \end{pmatrix}$$

Each weight matrix W :

$$W = D + S \quad (D \text{ dense, } S \text{ sparse})$$

Setting D to be the discrete-cosine-transform matrix, then

$$\Rightarrow Wx = \text{DCT}(x) + Sx$$



DCT plus Sparse

$$\mathbf{w} = \mathbf{d} + U\theta$$

$$\begin{pmatrix} | \\ | \\ | \\ | \\ | \end{pmatrix} = \begin{pmatrix} | \\ | \\ | \\ | \\ | \end{pmatrix} + \begin{pmatrix} 0 \dots 0 1 0 \dots 0 \\ 0 \dots 0 1 0 \dots 0 \\ \vdots \\ \vdots \\ \vdots \end{pmatrix} \begin{pmatrix} | \\ | \\ | \\ | \\ | \end{pmatrix}$$

Each weight matrix W :

$$W = D + S \quad (D \text{ dense, } S \text{ sparse})$$

Setting D to be the discrete-cosine-transform matrix, then

$$\Rightarrow Wx = \text{DCT}(x) + Sx$$



DCT plus Sparse

$$\mathbf{w} = \mathbf{d} + U\theta$$

$$\begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix} + \begin{pmatrix} 0 \dots 0 1 0 \dots 0 \\ \vdots \\ \vdots \\ \vdots \\ 0 \dots 0 1 0 \dots 0 \end{pmatrix} \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix}$$

Each weight matrix W :

$$W = D + S \quad (D \text{ dense, } S \text{ sparse})$$

Setting D to be the discrete-cosine-transform matrix, then

$$\Rightarrow Wx = \text{DCT}(x) + Sx$$

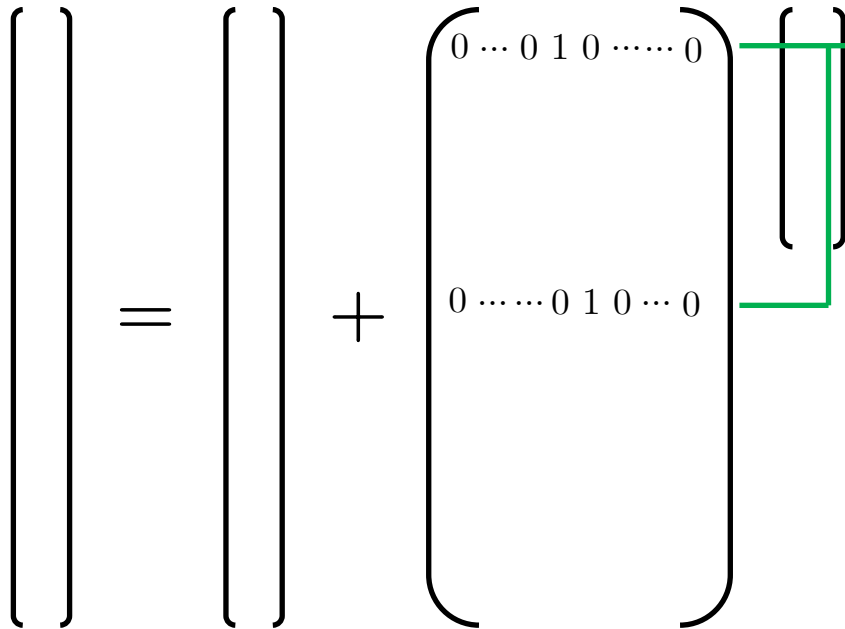
(We also add single trainable scaling param for DCT:)

$$\Rightarrow Wx = \alpha \text{DCT}(x) + Sx$$



DCT plus Sparse

$$\mathbf{w} = \mathbf{d} + U\theta$$



Distribution of 1-sparse rows

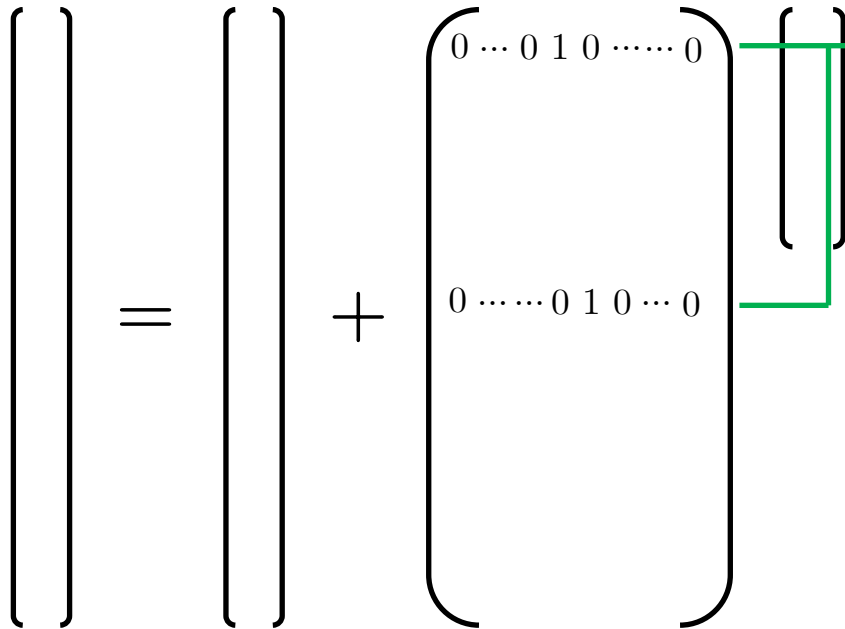


Distribution of trainable parameters
across and within S matrices



DCT plus Sparse

$$\mathbf{w} = \mathbf{d} + U\theta$$



Distribution of 1-sparse rows



Distribution of trainable parameters
across and within S matrices

“Equal per layer” (EPL):

- $\frac{k}{L}$ trainable params in each S
- Locations in S uniformly random
- No initialization of the dense net



DCT plus Sparse

$P =$	Δ Accuracy			Computational Cost			Network Size on Device		
	0.01	0.001	0.0001	At init.	Training	Final	At init.	Training	Final
Random	-11.9%	-66%	-66%	0	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
IMP	+0.8%	-7.1%	-64.8%	0	$\mathcal{O}(mn) \xrightarrow{t}$ $\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(N) \xrightarrow{t}$ $\mathcal{O}(PN)$	$\mathcal{O}(PN)$
FORCE	-6.6%	-26.9%	-62.4%	$\mathcal{O}(mnk)$	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
SynFlow	-6.2%	-31.6%	-60.4%	$\mathcal{O}(mnk)$	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
RigL (ERK)	+0.4%	-16.8%	-65.7%	0	$\mathcal{O}(pmn + \frac{1}{\Delta T} mn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
DCTpS	-5.8%	-15%	-22.8%	0	$\mathcal{O}(q \log q + pmn)$	$\mathcal{O}(q \log q + pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$



DCT plus Sparse

$P =$	Δ Accuracy			Computational Cost			Network Size on Device		
	0.01	0.001	0.0001	At init.	Training	Final	At init.	Training	Final
Random	-11.9%	-66%	-66%	0	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
IMP	+0.8%	-7.1%	-64.8%	0	$\mathcal{O}(mn) \xrightarrow{t}$ $\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(N) \xrightarrow{t}$ $\mathcal{O}(PN)$	$\mathcal{O}(PN)$
FORCE	-6.6%	-26.9%	-62.4%	$\mathcal{O}(mnk)$	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
SynFlow	-6.2%	-31.6%	-60.4%	$\mathcal{O}(mnk)$	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
RigL (ERK)	+0.4%	-16.8%	-65.7%	0	$\mathcal{O}(pmn + \frac{1}{\Delta T} mn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
DCTpS	-5.8%	-15%	-22.8%	0	$\mathcal{O}(q \log q + pmn)$	$\mathcal{O}(q \log q + pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$

$$W \in \mathbb{R}^{m \times n}$$

$$q = \max(m, n)$$

$$p = \frac{\|W\|_0}{mn}$$



DCT plus Sparse

$P =$	Δ Accuracy			Computational Cost			Network Size on Device		
	0.01	0.001	0.0001	At init.	Training	Final	At init.	Training	Final
Random	-11.9%	-66%	-66%	0	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
IMP	+0.8%	-7.1%	-64.8%	0	$\mathcal{O}(mn) \xrightarrow{t}$ $\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(N) \xrightarrow{t}$ $\mathcal{O}(PN)$	$\mathcal{O}(PN)$
FORCE	-6.6%	-26.9%	-62.4%	$\mathcal{O}(mnk)$	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
SynFlow	-6.2%	-31.6%	-60.4%	$\mathcal{O}(mnk)$	$\mathcal{O}(pmn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(N)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
RigL (ERK)	+0.4%	-16.8%	-65.7%	0	$\mathcal{O}(pmn + \frac{1}{\Delta T} mn)$	$\mathcal{O}(pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$
DCTpS	-5.8%	-15%	-22.8%	0	$\mathcal{O}(q \log q + pmn)$	$\mathcal{O}(q \log q + pmn)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$	$\mathcal{O}(PN)$

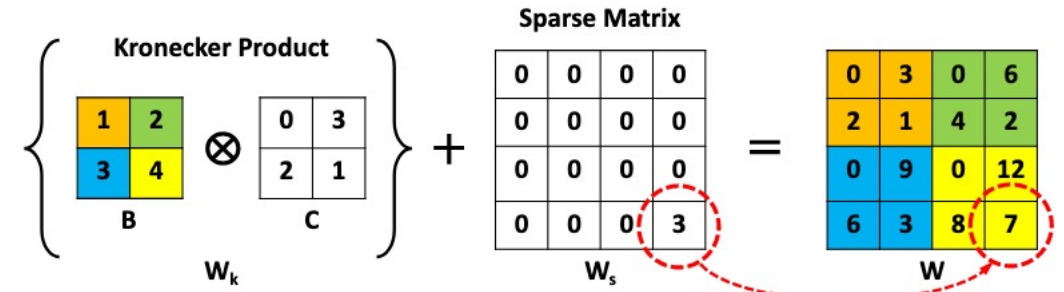
$P =$ Global density

$N =$ Total Network Params



Interlude: a "Trend"?

Doping: A technique for efficient compression of LSTM models using sparse structured additive matrices, Thakker et al, 2021

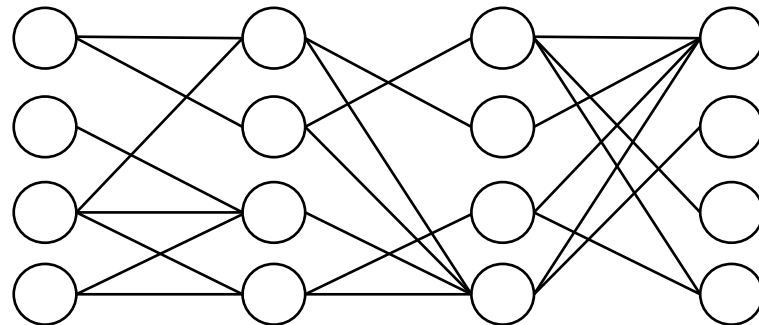


Key difference:
Needs to be
trained and stored



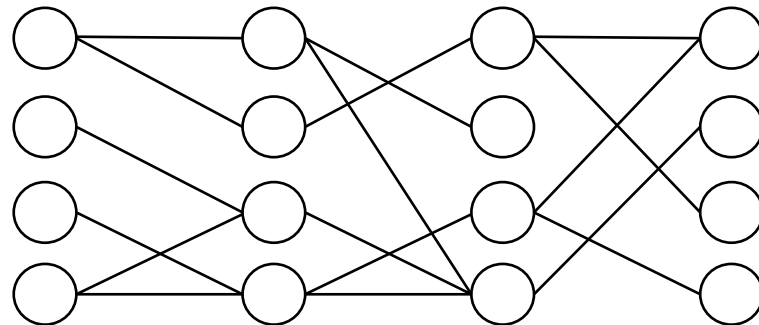
Dynamic Sparse Training

- Static sparse training chooses support set of \mathbf{w} , then keeps fixed during training.
- DST instead jointly optimises topology *and* weights, subject to fixed sparsity level
- Start sparse, then: Train, Prune, Regrow, Repeat.



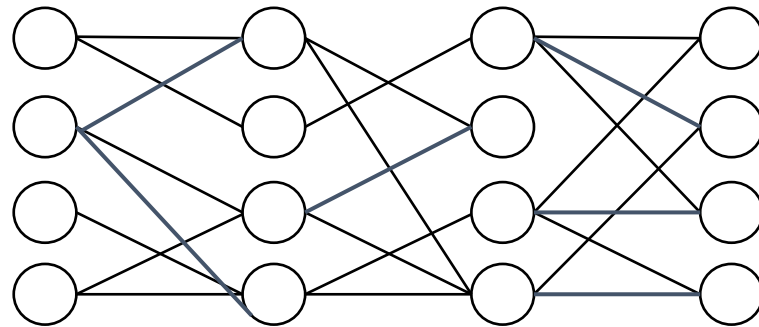
Dynamic Sparse Training

- Static sparse training chooses support set of \mathbf{w} , then keeps fixed during training.
- DST instead jointly optimises topology *and* weights, subject to fixed sparsity level
- Start sparse, then: Train, Prune, Regrow, Repeat.



Dynamic Sparse Training

- Static sparse training chooses support set of \mathbf{w} , then keeps fixed during training.
- DST instead jointly optimises topology *and* weights, subject to fixed sparsity level
- Start sparse, then: Train, Prune, Regrow, Repeat.



Combining with Dynamic Sparse training

Straightforwardly combined
with DCTpS:

$$Wx = \text{DCT}(x) + Sx$$

Apply DST to the sparse,
trainable matrices in each layer



Combining with Dynamic Sparse training

Where are connections initialized?

Which to prune?

Which to regrow?

Where can they regrow?



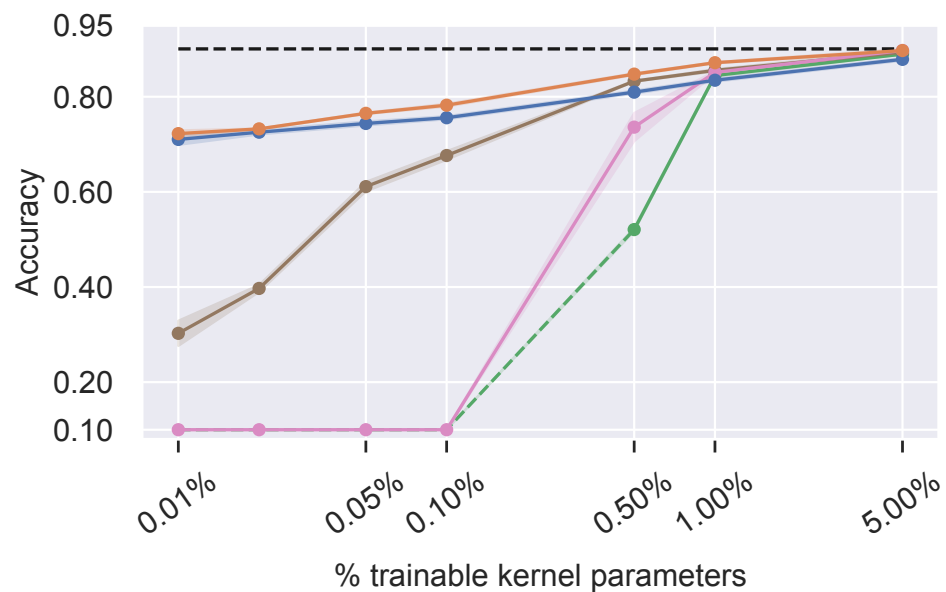
Combining with Rig-L¹

- Where are connections initialized? → Initialise all W_i as modified Erdos Reyni random bipartite graph, (Details not NB)
- Which to prune? → Smallest Magnitude
- Which to regrow? → Largest magnitude gradient
- Where can they regrow? → Within the same layer

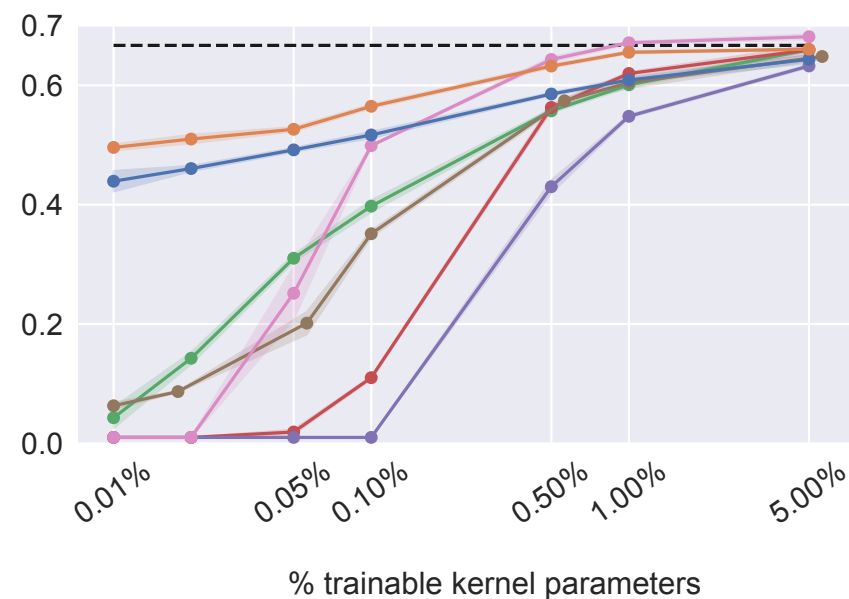


Some Results

MobileNetV2 on CIFAR10



ResNet50 on CIFAR100

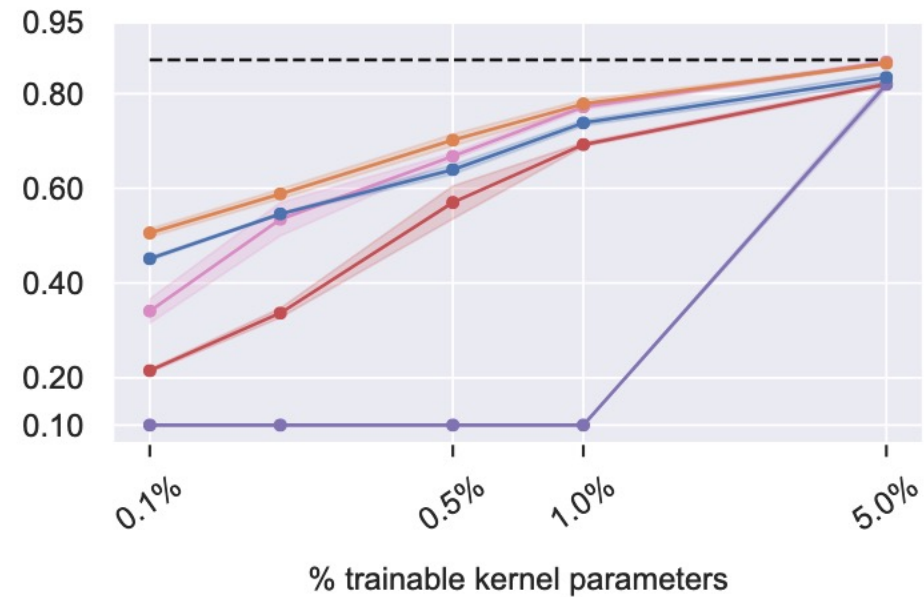


—●— FORCE —●— Random (EPL) —●— Random (Uniform) —●— SynFlow —●— RigL (ERK) —●— DCTpS —●— DCTpS + RigL



Without Batchnorm

FixUpResnet110 on CIFAR10



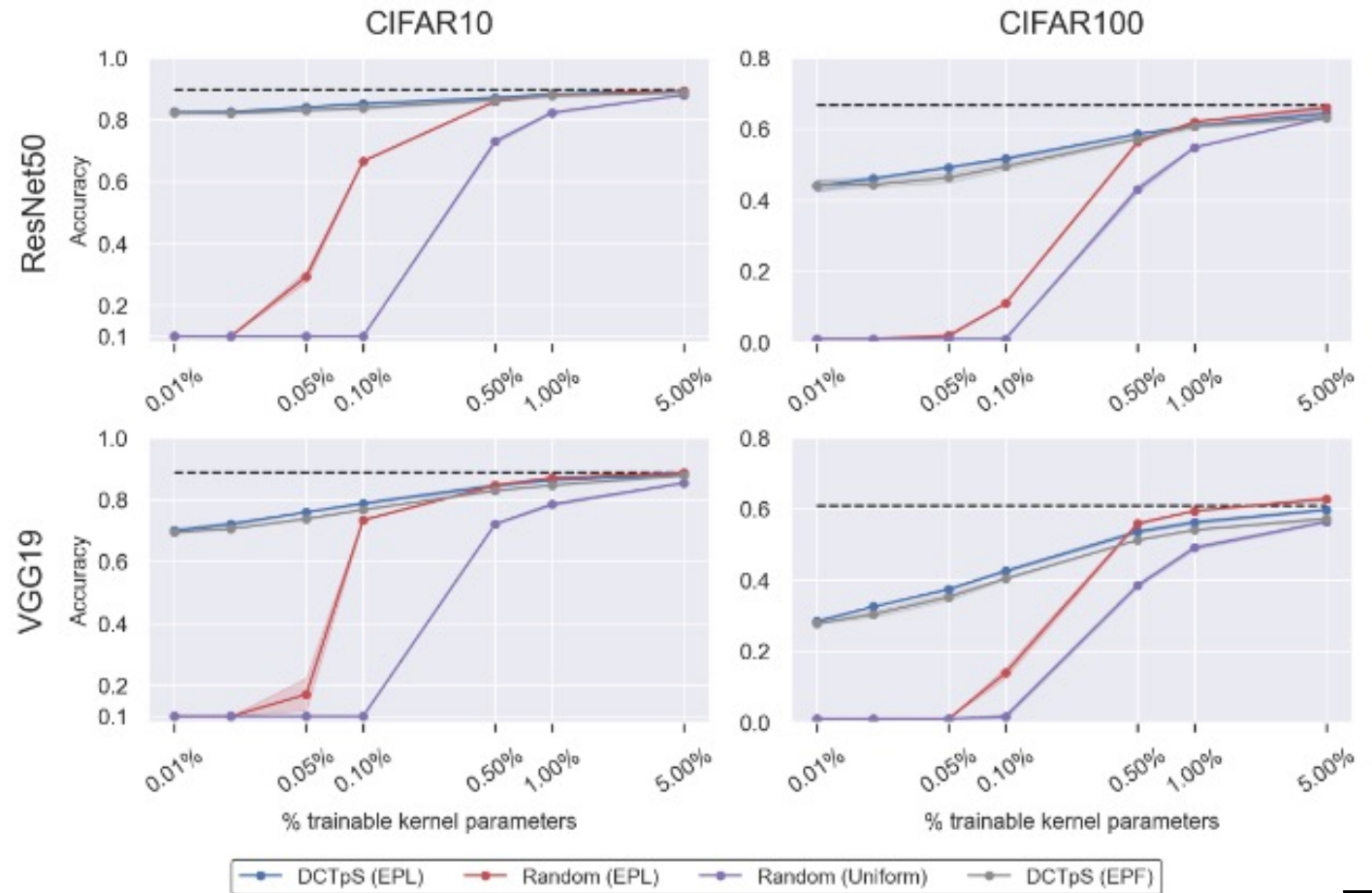
—●— FORCE —●— Random (EPL) —●— Random (Uniform) —●— SynFlow —●— RigL (ERK) —●— DCTpS —●— DCTpS + RigL



Other support distributions?

“Equal per Filter” (EPF)

(A version of N:M sparsity)



Limitations

- Computational floor imposed by the DCT – can push storage, not compute, down to the extremes. Q: more efficient ways to achieve an appropriate offset?



Limitations

- Computational floor imposed by the DCT – can push storage, not compute, down to the extremes. Q: more efficient ways to achieve an appropriate offset?
- Gains are hardware and implementation dependent – so far these are gains “in theory”



Limitations

- Computational floor imposed by the DCT – can push storage, not compute, down to the extremes. Q: more efficient ways to achieve an appropriate offset?
- Gains are hardware and implementation dependent – so far these are gains “in theory”
- Does not speak to the more general question about how best to use parameters (sparser, larger net vs denser, smaller net, etc)



Limitations

- Computational floor imposed by the DCT – can push storage, not compute, down to the extremes. Q: more efficient ways to achieve an appropriate offset?
- Gains are hardware and implementation dependent – so far these are gains “in theory”
- Does not speak to the more general question about how best to use parameters (sparser, larger net vs denser, smaller net, etc)
- Shrinks the storage footprint of the network – but the hidden representations are not sparse and can be very large (same for all sparse nets)



Thank you

Reach out on ilan.price@maths.ox.ac.uk or @IlanPrice

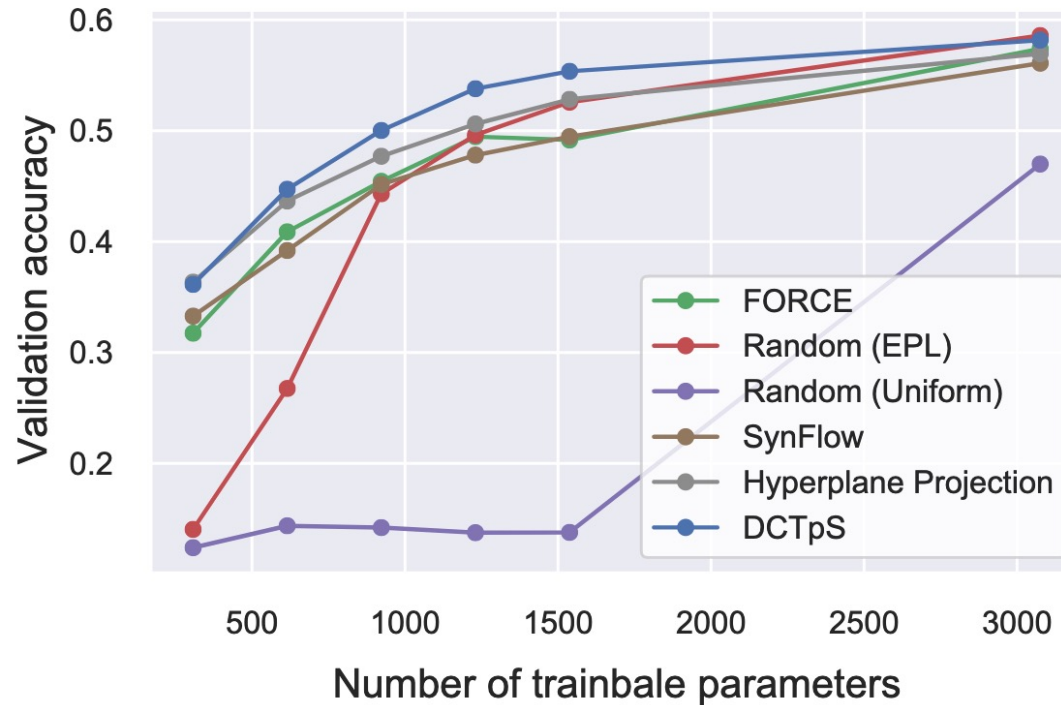


**The
Alan Turing
Institute**

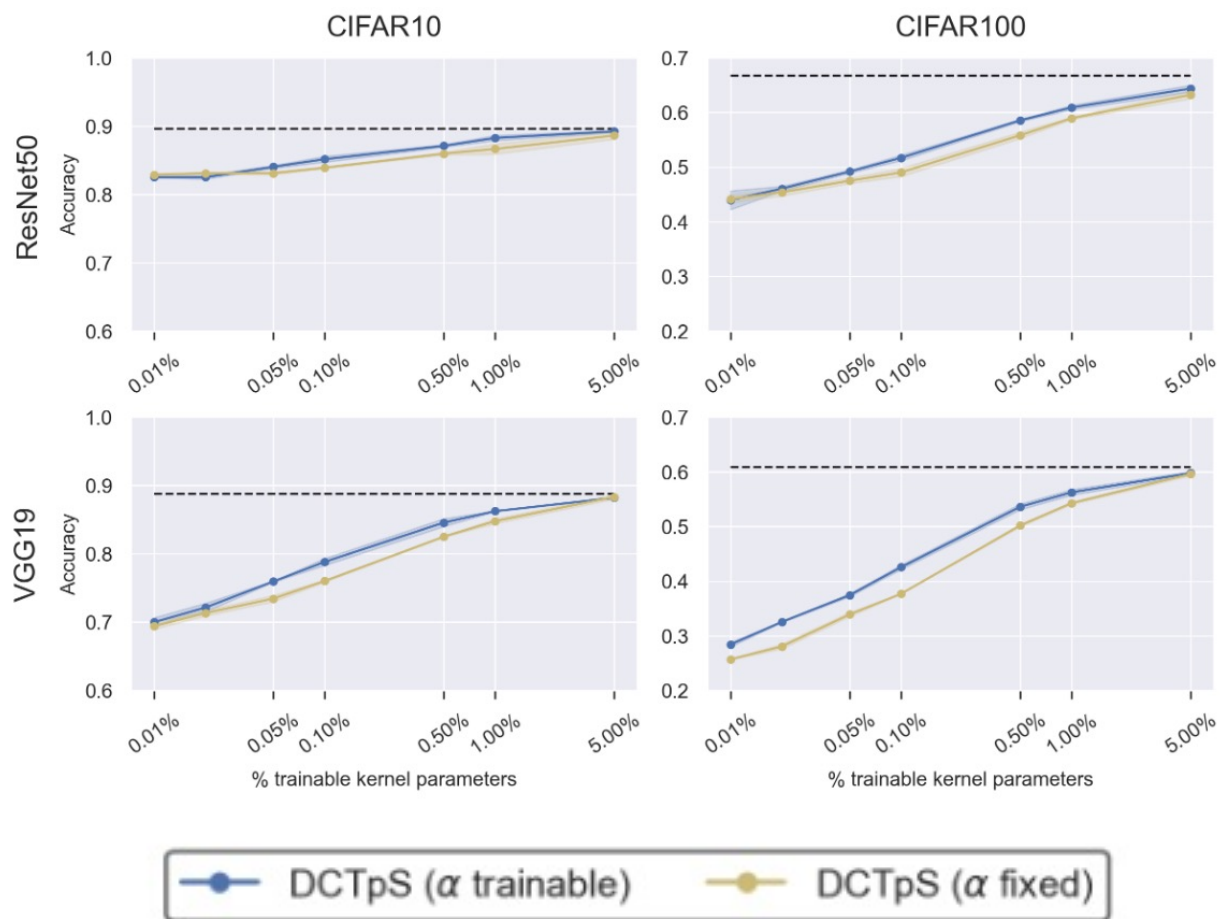


Compared with Random Subspace?

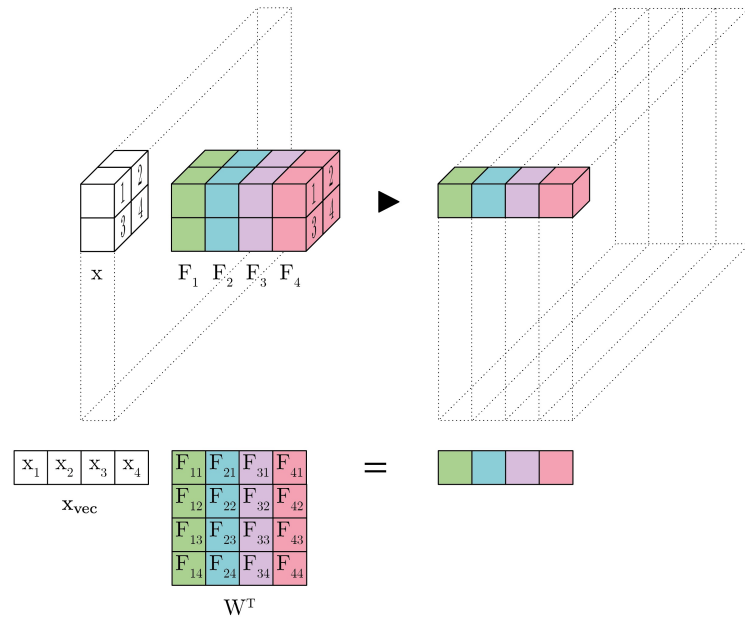
Lenet-5 on CIFAR-10



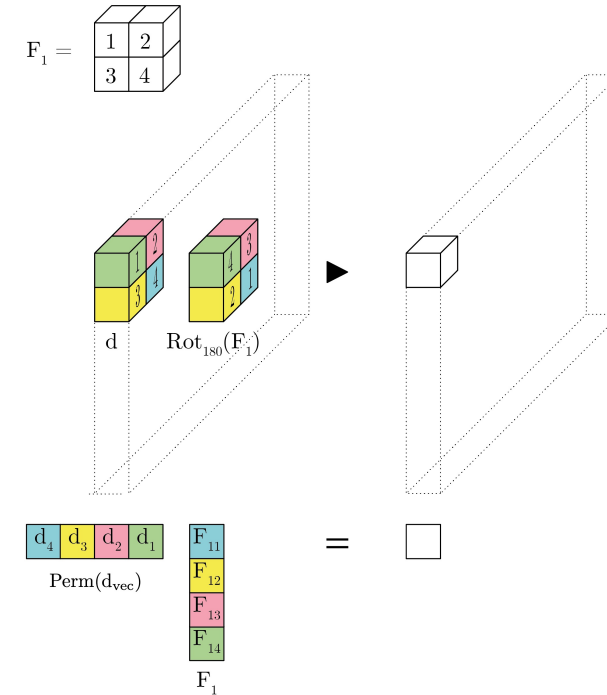
What if we fixed α ?



DCTpS Convolutional Layers

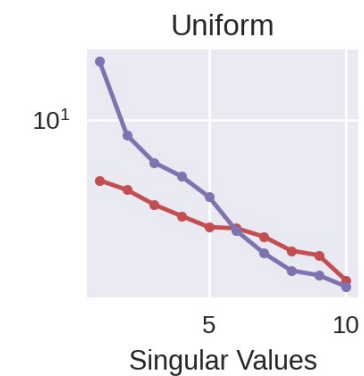
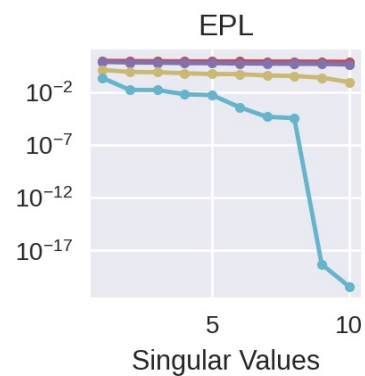
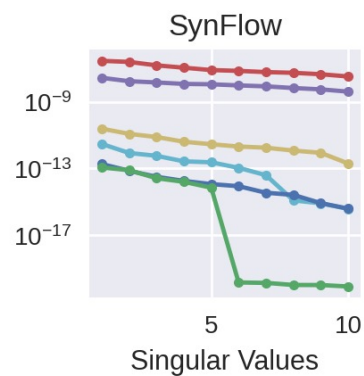
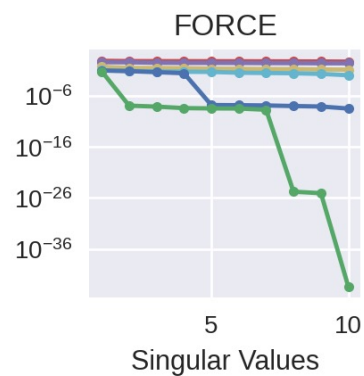
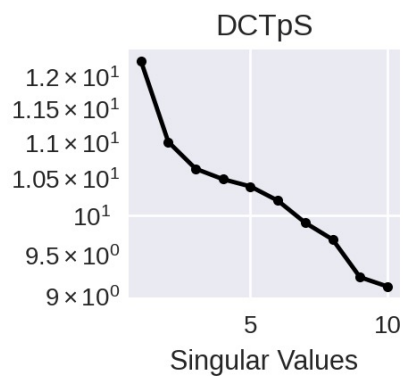
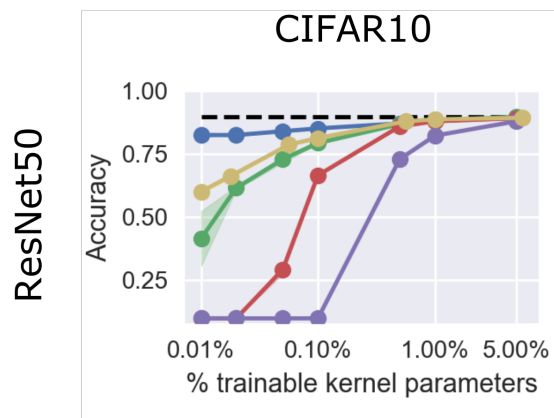


Forward pass



Backprop

When is training possible?



Any trainable density
 1.0%
 0.5%
 0.1%
 0.05%
 0.02%
 0.01%

Jacobian Spectra