



NYU | COURANT

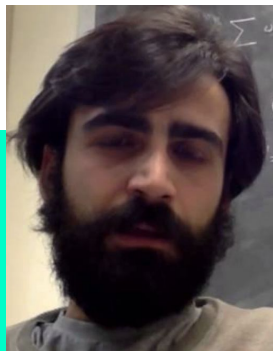
SIMONS
FOUNDATION

FINDING THE NEEDLE IN THE HAYSTACK

ON THE BENEFITS OF ARCHITECTURAL CONSTRAINTS



Stéphane d'Ascoli



Levent Sagun



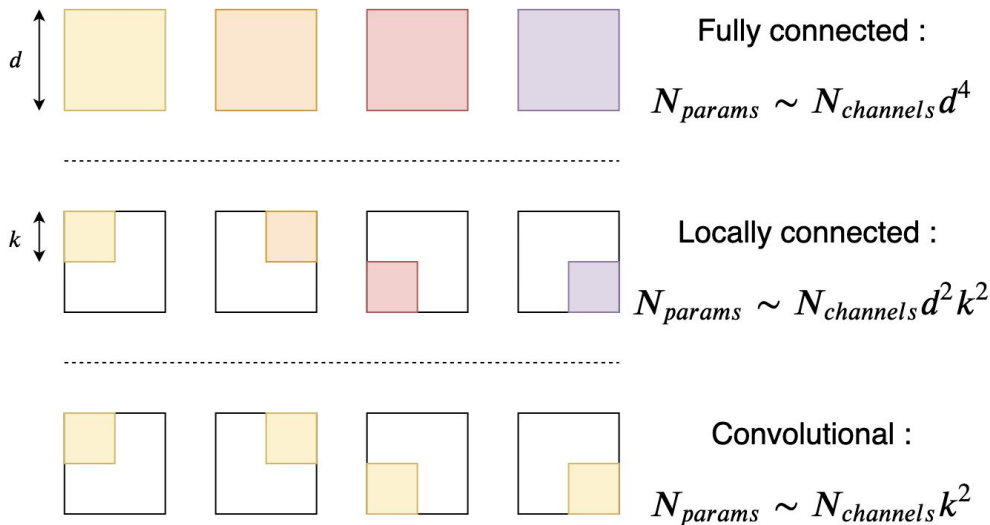
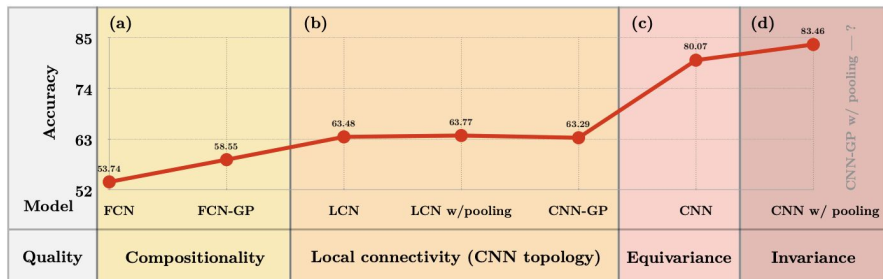
Joan Bruna



Giulio Biroli

FCN VS CNN

[Novak et al. 2019]



CNNs are just FCNs which generalize better thanks to two constraints :

- Locality : smaller receptive field
- Weight sharing : filters use same weights

How do these constraints help navigate the loss landscape to find rare basins of high generalisation ?

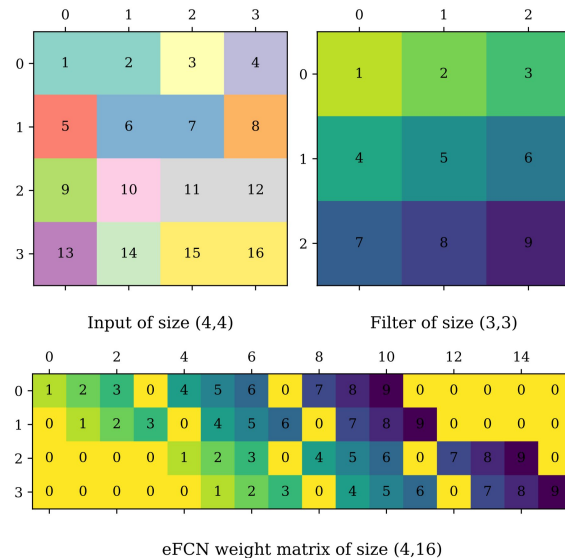
Can we use of their benefits (less parameters) without their potential disadvantages (loss of generality) ?

FCN VS CNN

One can always express a convolution using a matrix

Hence, one can always map a CNN to its equivalent FCN (eFCN) or LCN (eLCN)

The eFCN / eLCN obtained are highly sparse and redundant

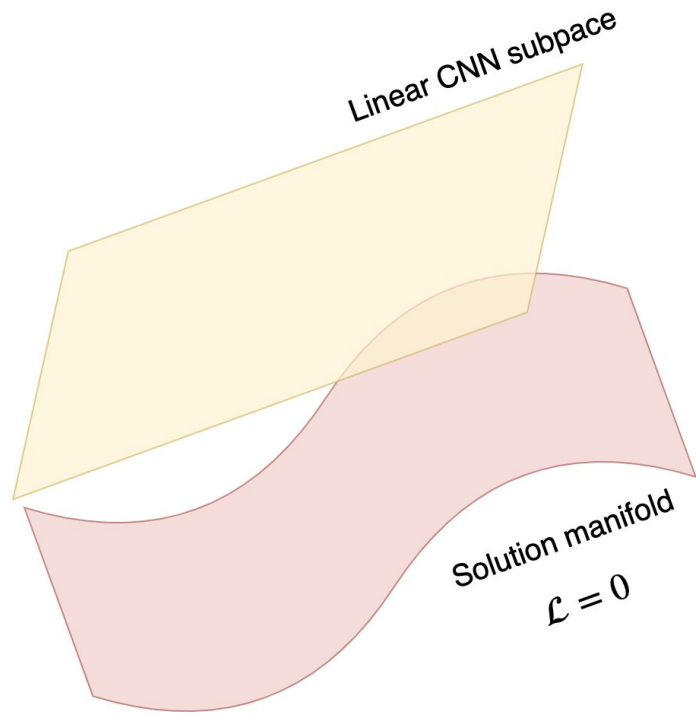


$$d_{in} = 4$$

$$(k, s, p) = (3, 1, 0)$$

$$d_{out} = \frac{d_{in} - 2p - k}{s} + 1 = 2$$

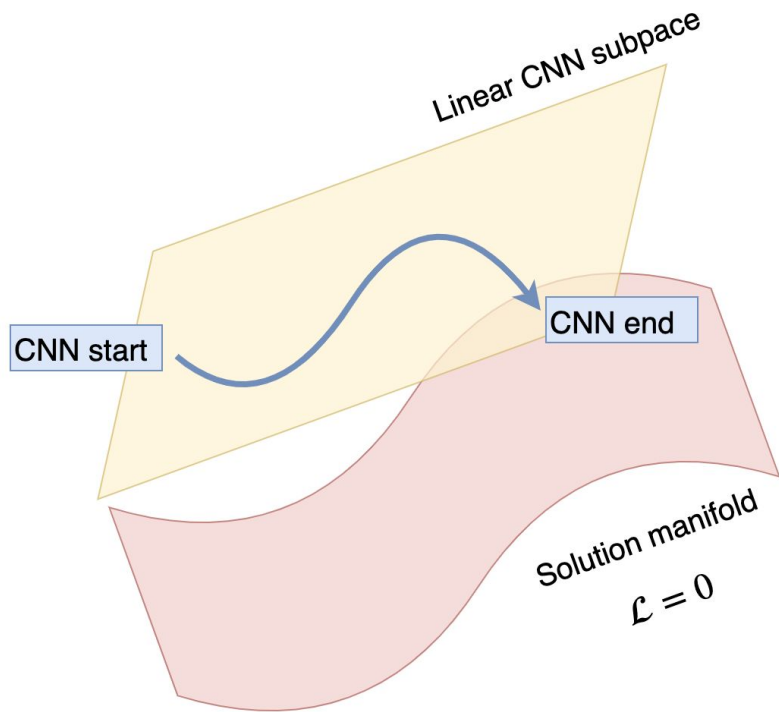
EMBEDDING TO FULLY CONNECTED SPACE



Build an embedding ϕ from a CNN to its eFCN that satisfies

$$f(\theta^{CNN}) = f(\Phi(\theta^{CNN})) = f(\theta^{eFCN})$$

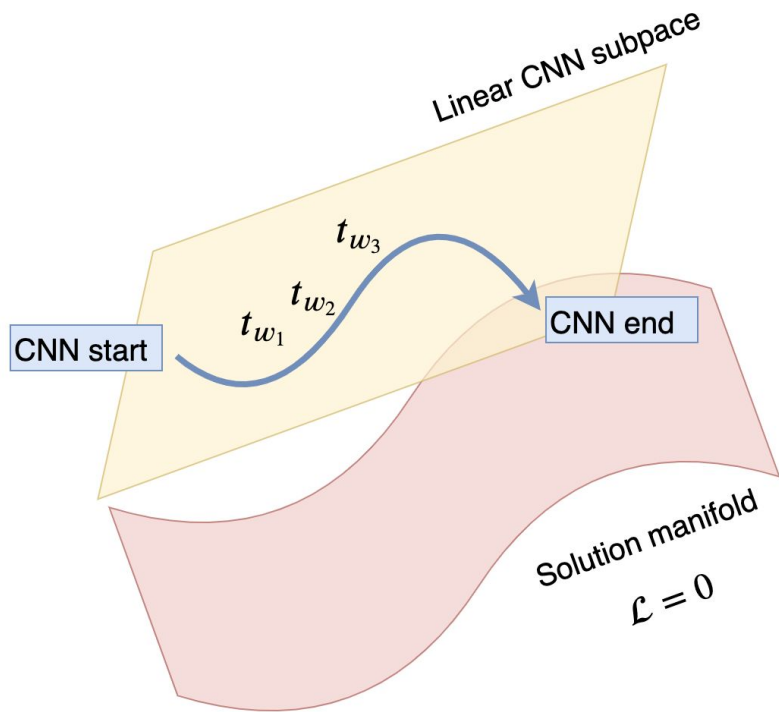
OUR EXPERIMENT



A toy experiment :

1. Train a CNN until it reaches zero loss

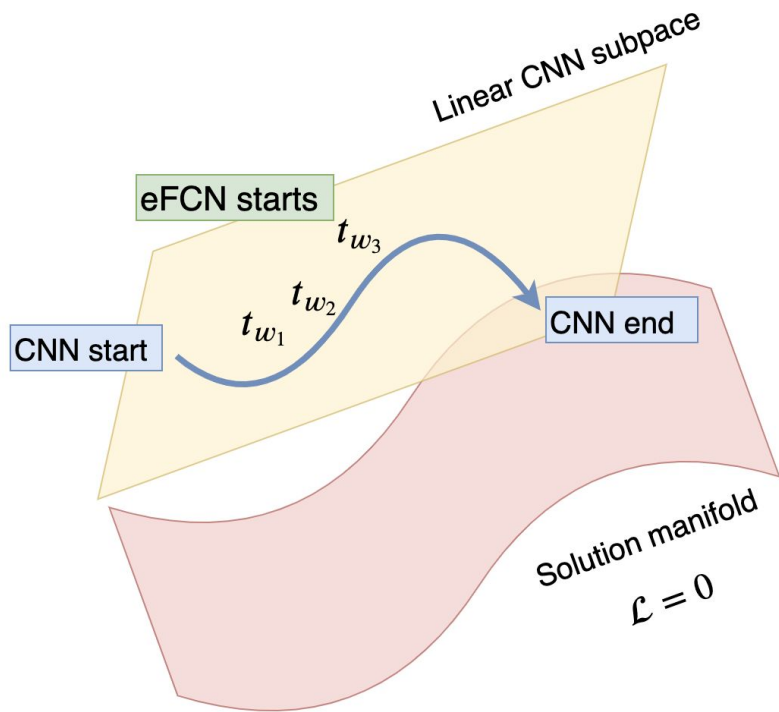
OUR EXPERIMENT



A toy experiment :

1. Train a CNN until it reaches zero loss
2. Along the training, save snapshots of the weights at various “switch” times

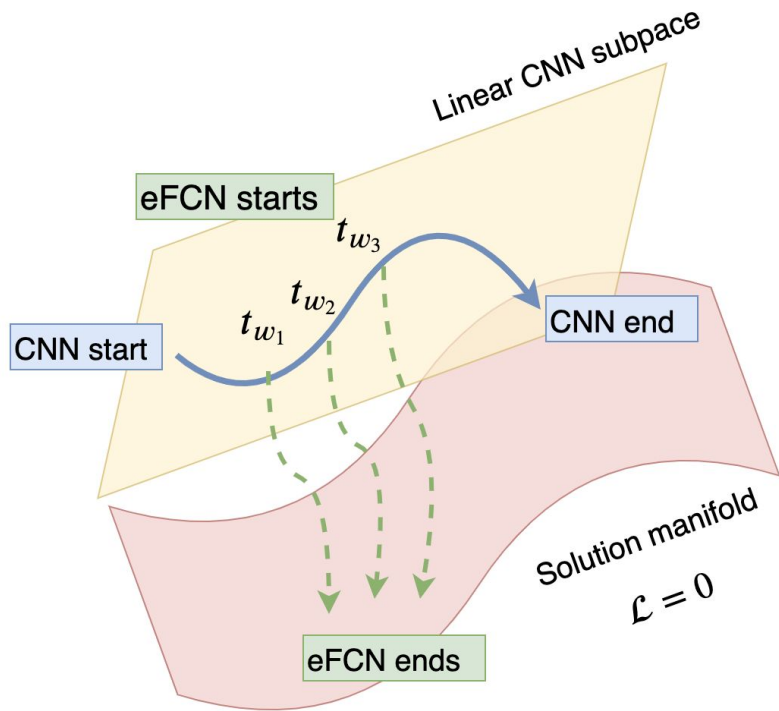
OUR EXPERIMENT



A toy experiment :

1. Train a CNN until it reaches zero loss
2. Along the training, save snapshots of the weights at various “switch” times
3. Map these to eFCN space using ϕ

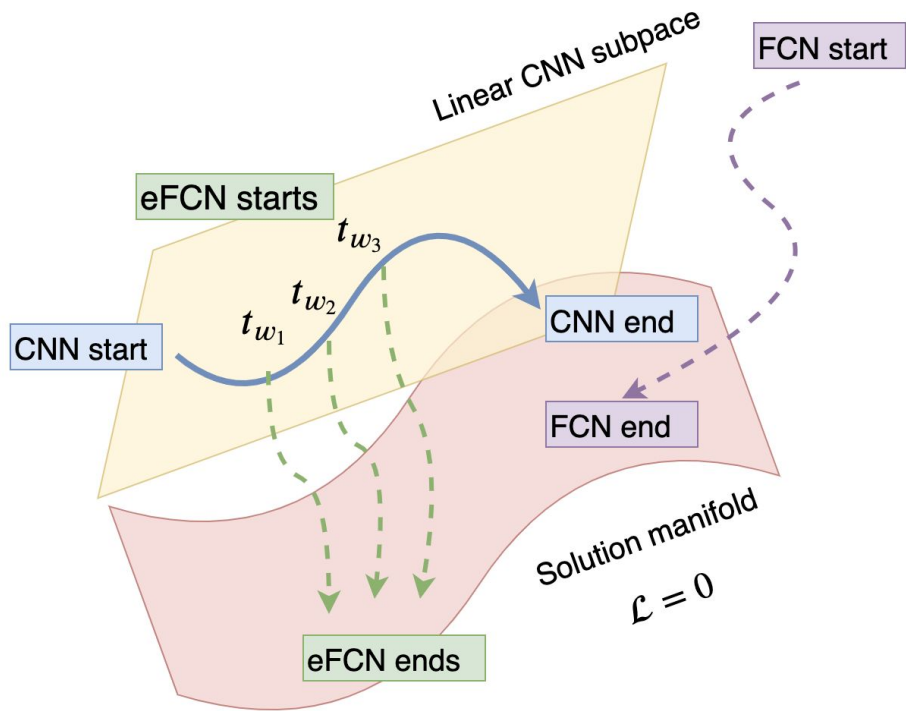
OUR EXPERIMENT



A toy experiment :

1. Train a CNN until it reaches zero loss
2. Along the training, save snapshots of the weights at various “switch” times
3. Map these to eFCN space using ϕ
4. Train the eFCNs until they reach zero loss

OUR EXPERIMENT



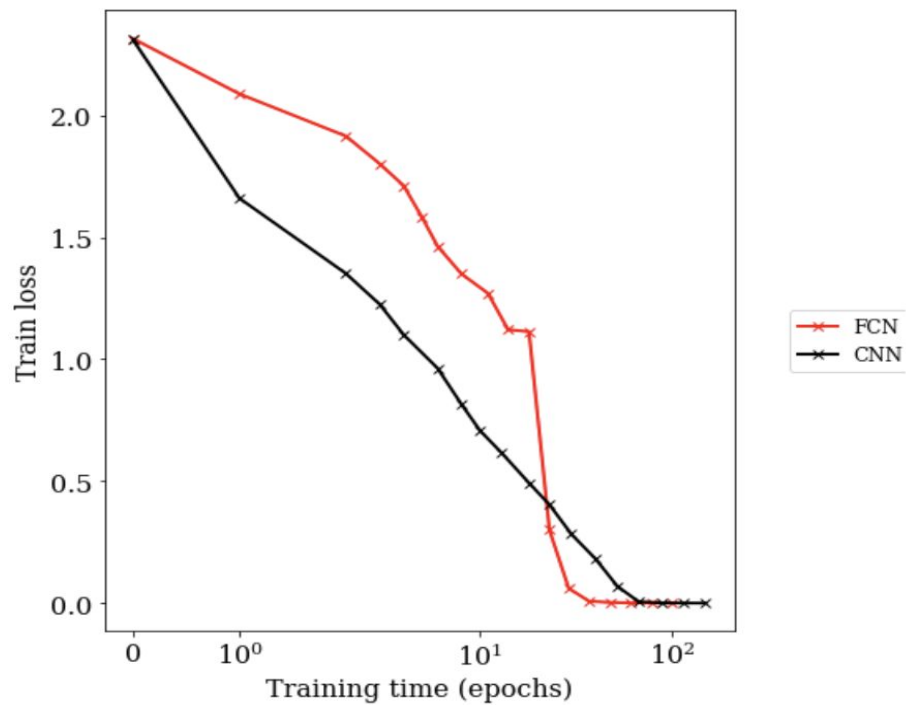
A toy experiment :

1. Train a CNN until it reaches zero loss
2. Along the training, save snapshots of the weights at various “switch” times
3. Map these to eFCN space using ϕ
4. Train the eFCNs until they reach zero loss
5. For comparison, train a vanilla FCN of same size

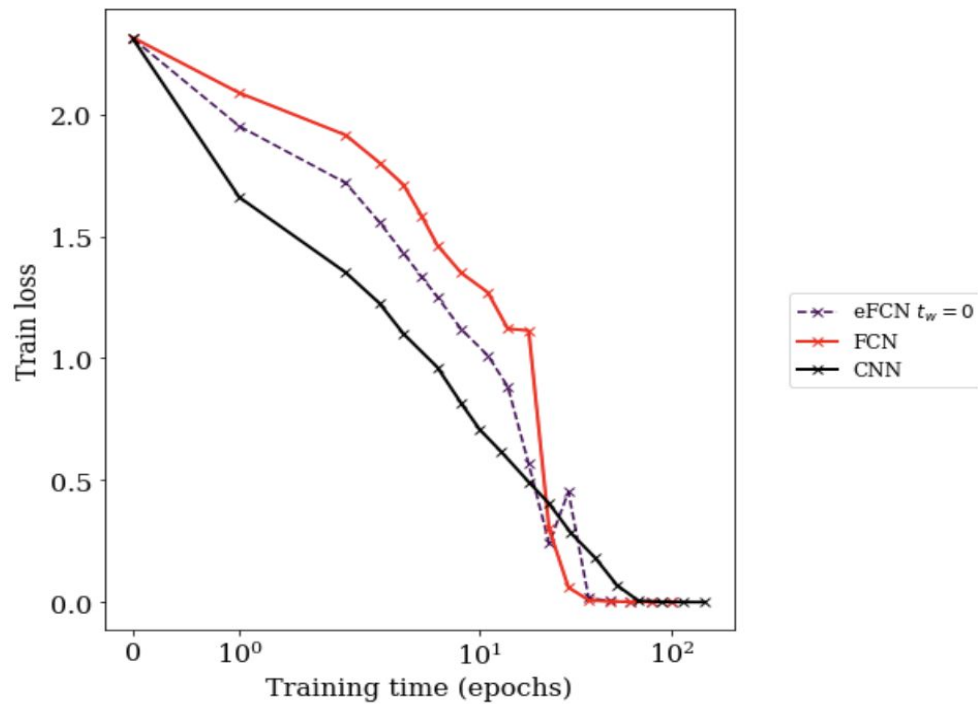
EXPERIMENT DETAILS

- Models: VanillaCNN (3 convolutional layers, 64 channels) + AlexNet
- Memory overhead for AlexNet: 10MB → 1GB !
- Datasets: CIFAR-10 and CIFAR-100
- All layers apart from convolutional (dense, Dropout, Pooling, ReLU) are unchanged apart from reshaping
- PyTorch implementation open-sourced at github.com/sdascoli/anarchitectural-search

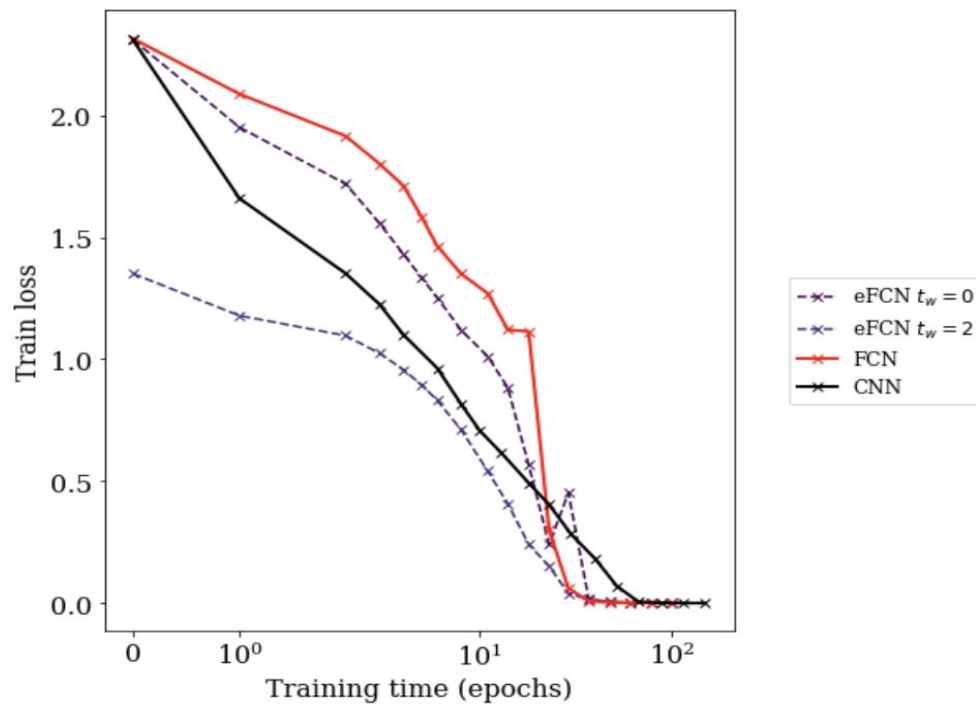
TRAINING DYNAMICS



TRAINING DYNAMICS

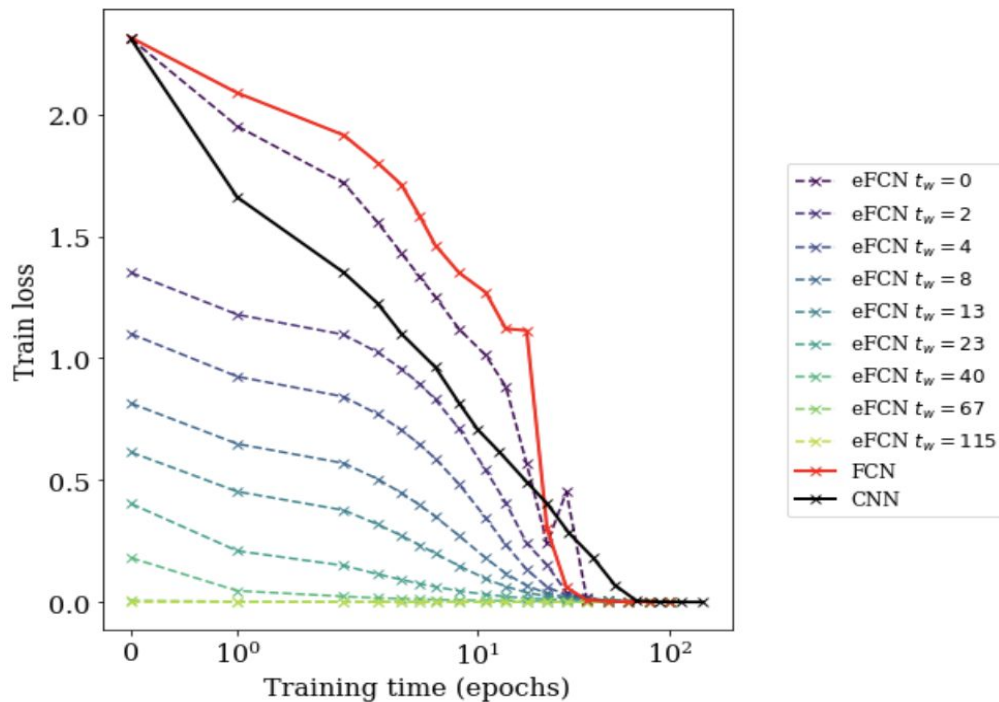


TRAINING DYNAMICS



- Train loss keeps going down after relaxing constraints

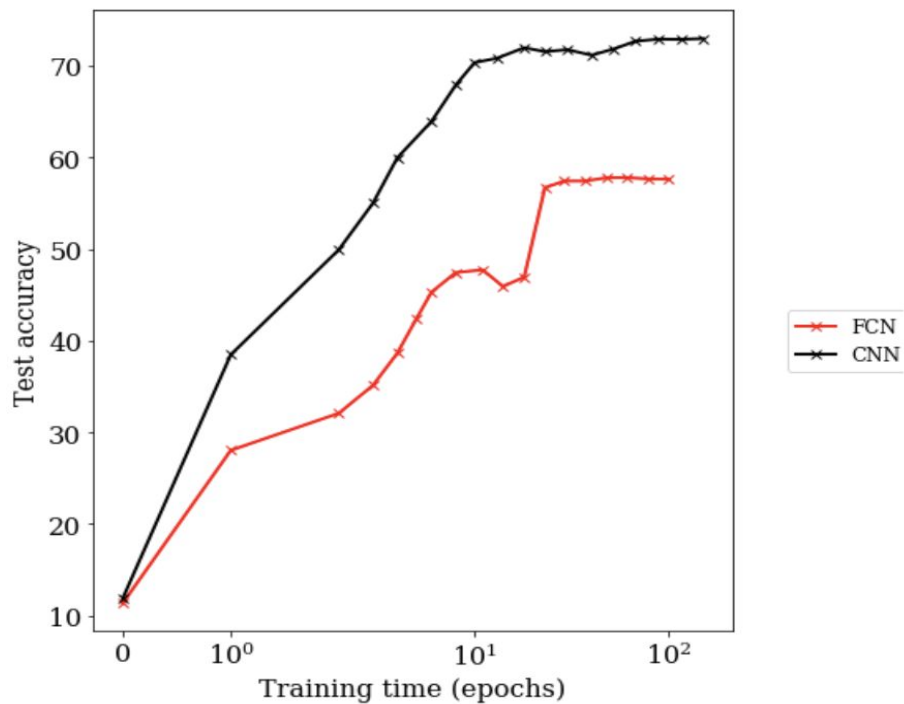
TRAINING DYNAMICS



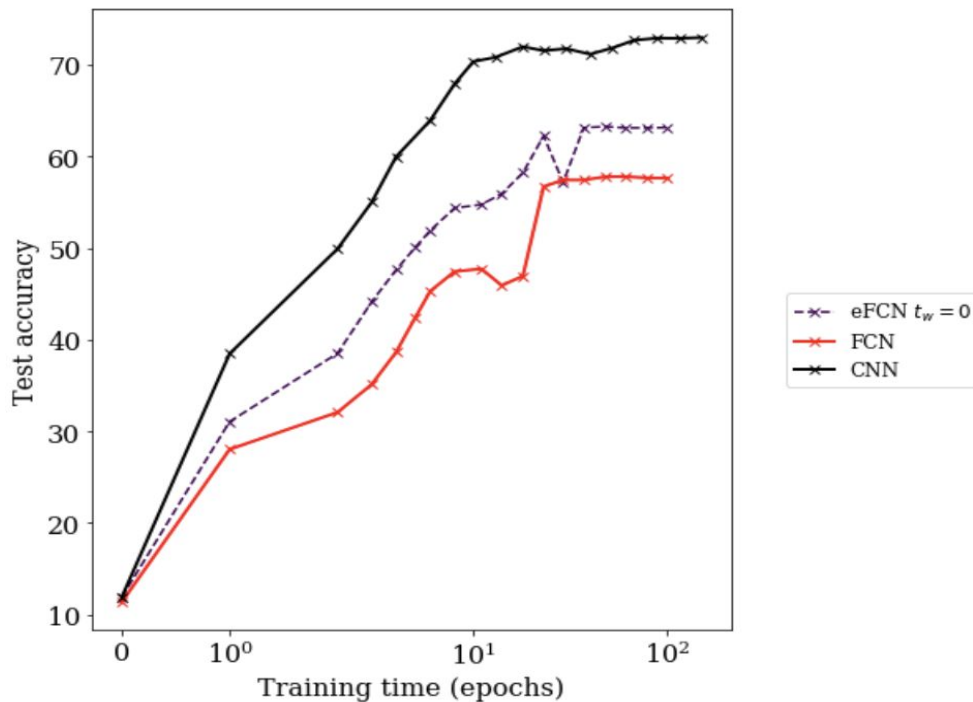
- Train loss keeps going down after relaxing constraints
- Even when we release close to convergence !

The eFCN stays in the basin reached by the CNN

GENERALIZATION DYNAMICS

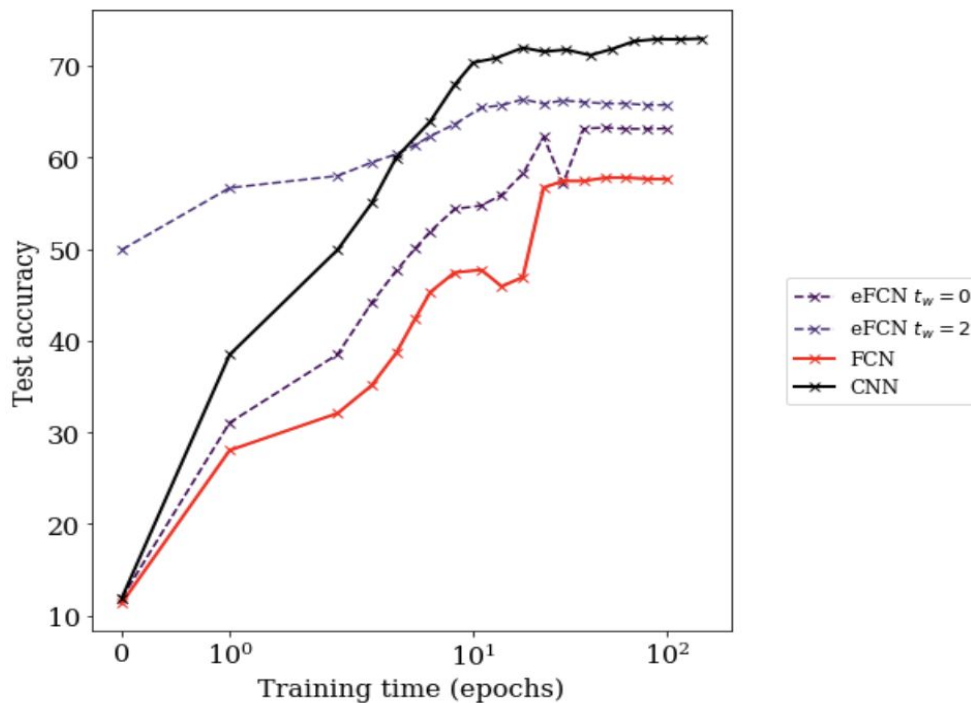


GENERALIZATION DYNAMICS



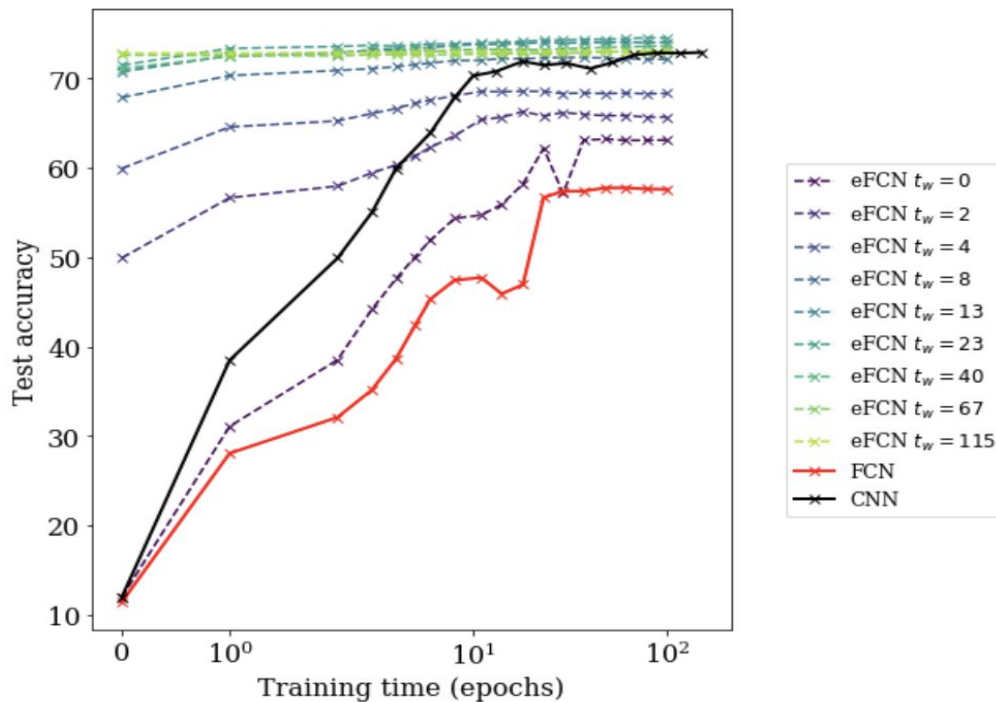
- FCN with CNN initialization beats vanilla FCN
 - ⇨ CNN constraints start us off in a good region

GENERALIZATION DYNAMICS



- FCN with CNN initialization beats vanilla FCN
 - ⇒ CNN constraints start us off in a good region
- Test accuracy keeps going up after relaxing constraints
 - ⇒ No overfitting

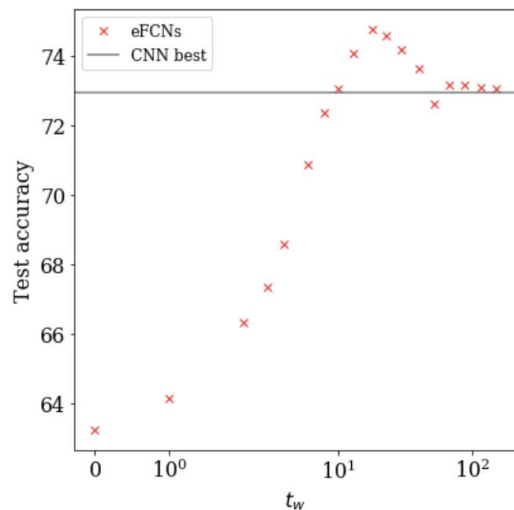
GENERALIZATION DYNAMICS



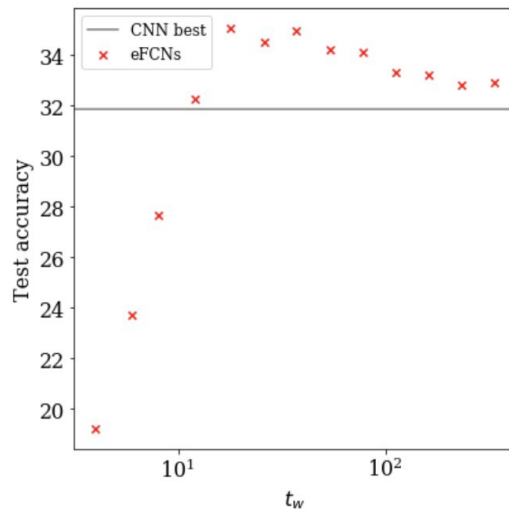
- FCN with CNN initialization beats vanilla FCN
 - ⇨ CNN constraints start us off in a good region
- Test accuracy keeps going up after relaxing constraints
 - ⇨ No overfitting
- When released late, the eFCNs even beat the CNN !
 - ⇨ CNN constraints are mostly helpful at the beginning of optimization

Releasing architectural constraints once they have lead us to the right region of the landscape can be helpful !

WHAT IS THE BEST TIME TO RELAX THE CONSTRAINTS ?



VanillaCNN, CIFAR-10

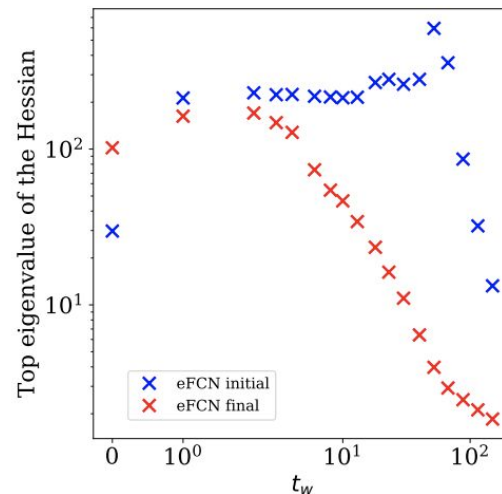
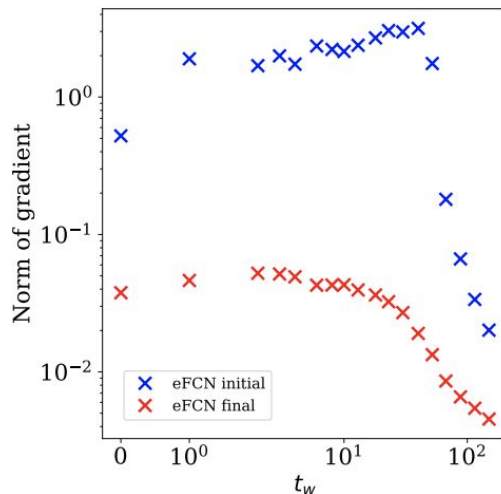
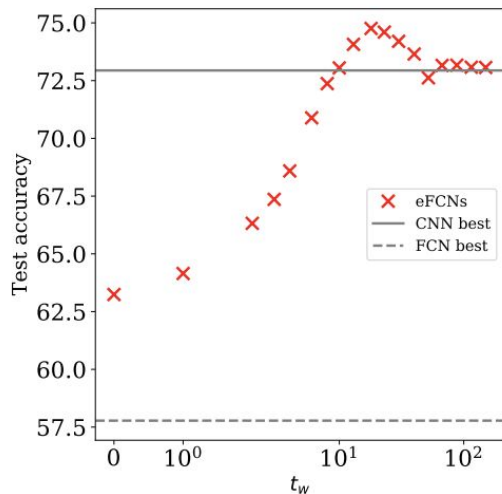


AlexNet, CIFAR-100, Adam

The improvement is best when constraints are released at intermediate times.

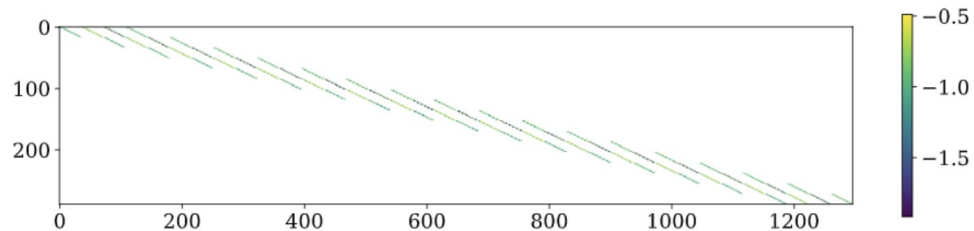
So what exactly are the extra degrees of freedom doing ?

A CLOSER LOOK AT THE LANDSCAPE

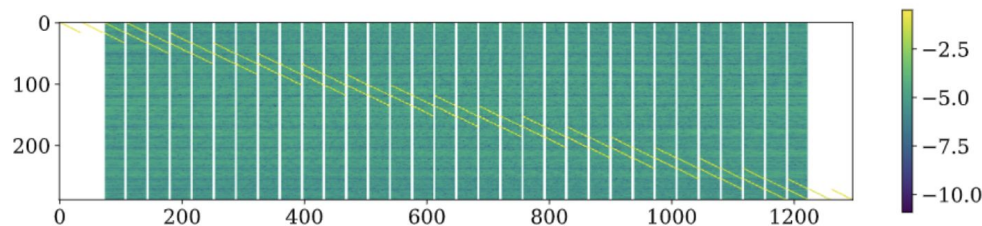


- The improvement is highest when constraints are released in the roughest/sharpest part of the landscape
- eFCN reaches wider basins

LOCAL VS OFF-LOCAL BLOCKS



eFCN at initialization

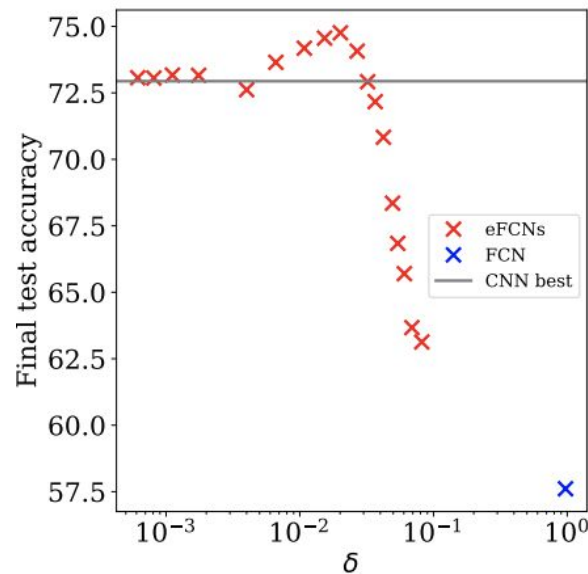
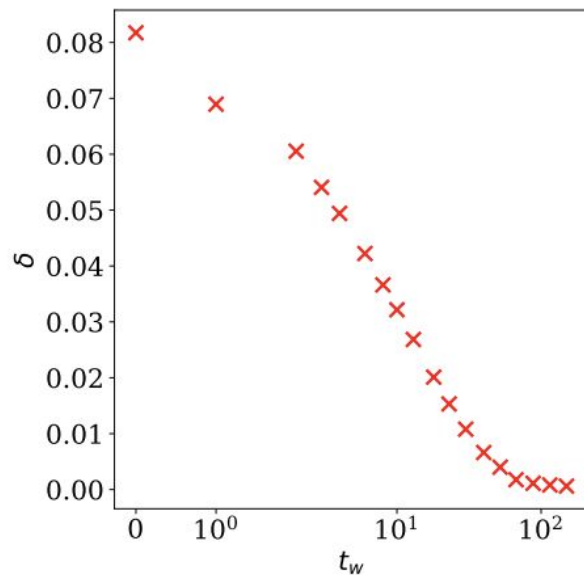


eFCN after 100 epochs

Measure of how far the model is from CNN subspace:

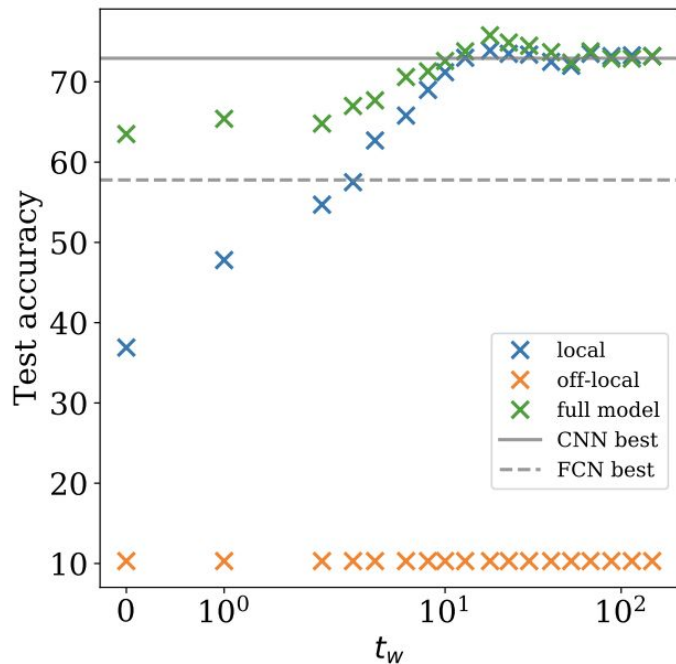
$$\delta = \frac{||\theta_{\text{off-local}}||_2}{||\theta||_2}$$

HOW FAR DO WE ESCAPE FROM THE CNN SUBSPACE ?



- The earlier we relax constraints, the further we escape
- There is an optimal and a cutoff distance from CNN subspace

ROLE OF OFF-LOCAL BLOCKS

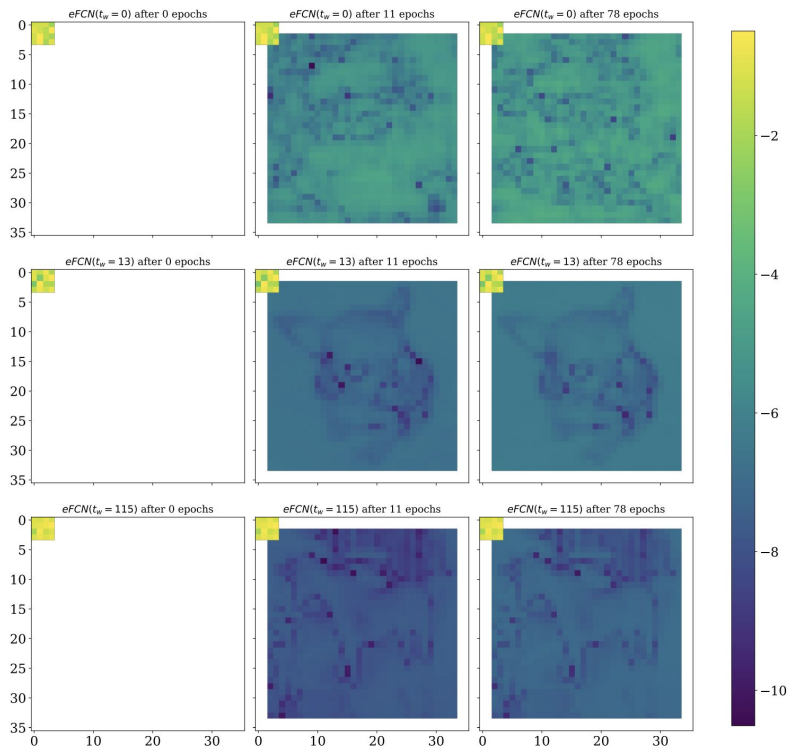


- All by itself the off-local part is useless
- However combined with local part it provides improvement

How can local and off-local blocks play such complementary roles ?

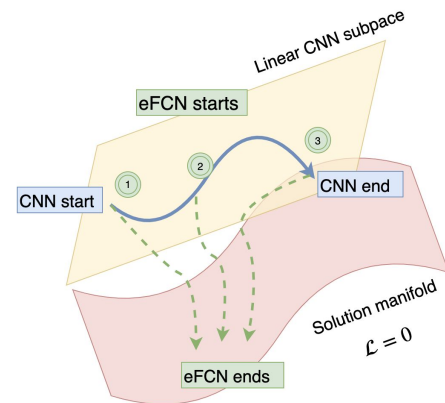
WHAT REPRESENTATION IS LEARNT ?

1



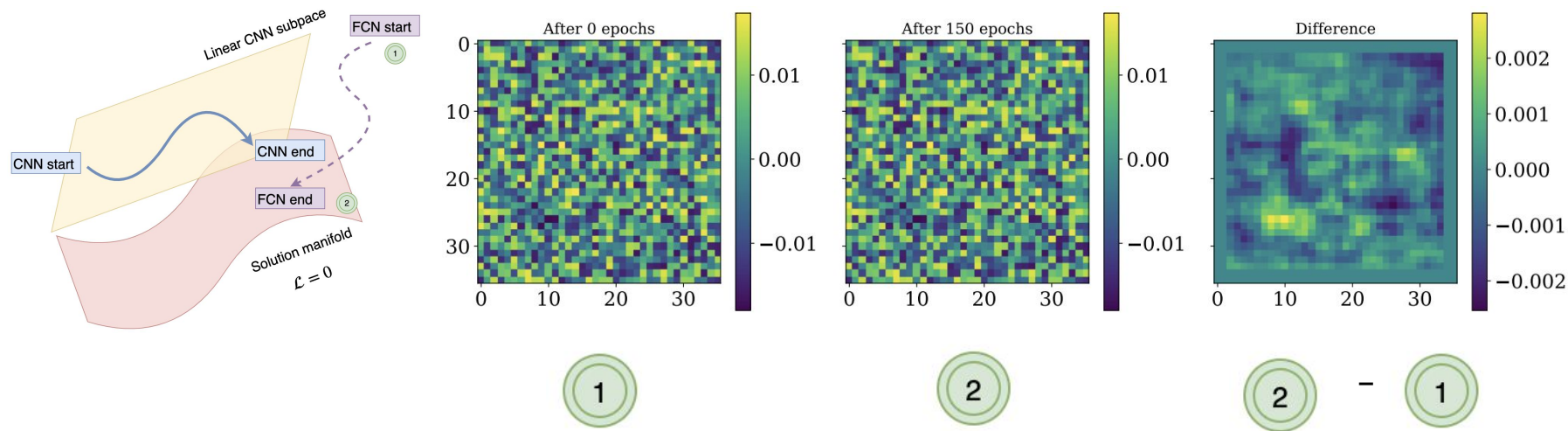
2

3



- The local blocks hardly change at all
- The off-local part learns silhouettes : performs some kind of template-matching
- Plays a complementary role to the local part's feature extraction

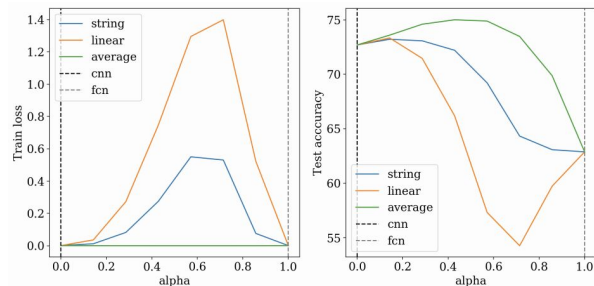
WHAT REPRESENTATION IS LEARNT ?



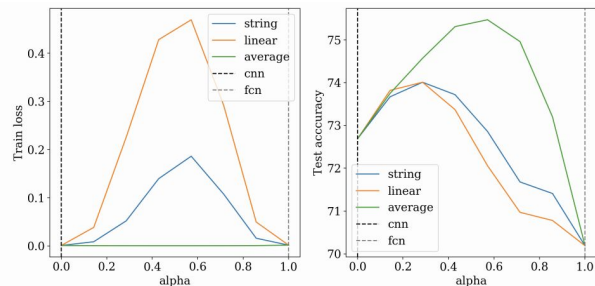
In the regular FCN case, template-matching is not a good enough strategy to make the loss diminish.

The weights capture a loose texture but nothing recognizable.

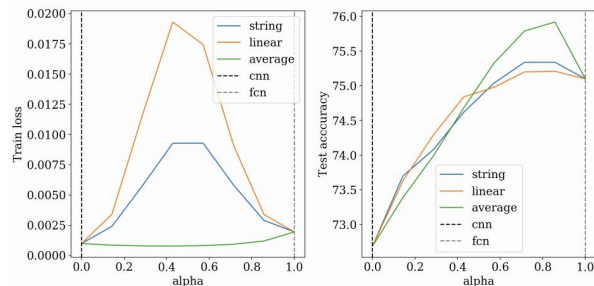
INTERPOLATING IN WEIGHT SPACE



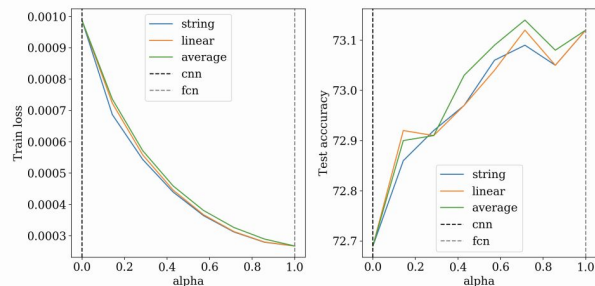
(a) $t_w = 0$



(b) $t_w = 5$



(c) $t_w = 18$



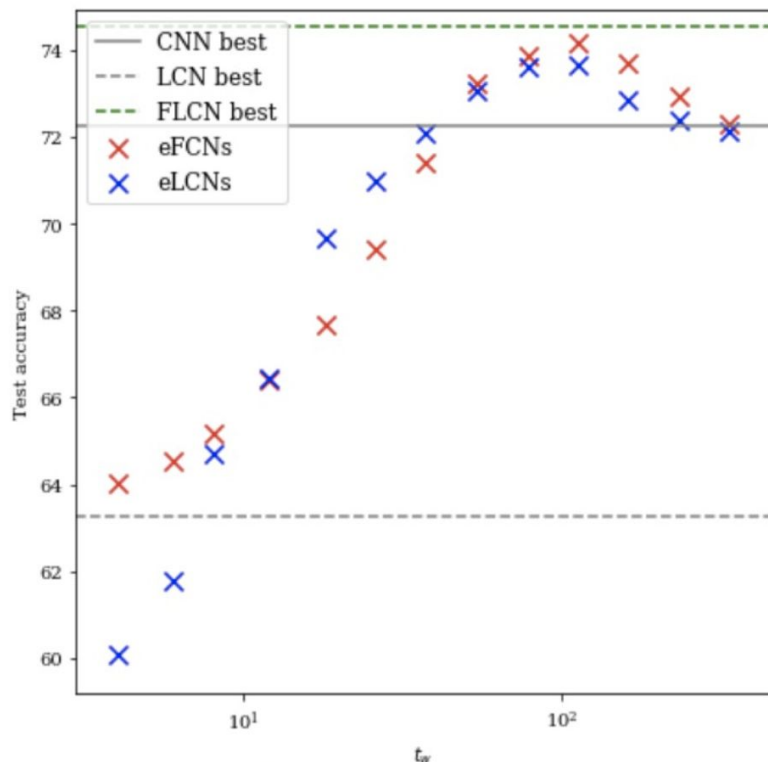
(d) $t_w = 61$

- Interpolating in weight space between the CNN and the eFCN actually increases generalization

- Maybe the best solution is somewhere in between the CNN and the eFCN ?

→ eLCNs !

WHAT HAPPENS FOR LOCALLY CONNECTED NETWORKS ?



- The gain in accuracy is very similar, but this time, memory increase is much more reasonable (factor 10)
- We can get further increase in accuracy by relaxing the constraints **hierarchically** :
 $\text{FLCN}(tw1, tw2) = \text{CNN} \rightarrow \text{eLCN}(tw1) \rightarrow \text{eFCN}(tw2)$
- Next step : plot the performance increase in the $(tw1, tw2)$ plane (WIP)

CONCLUSIONS

1. CNN constraints are mostly helpful at the **initialization** and **early stages** of training, bringing us to good regions of the landscape
2. Releasing the constraints once they have done their job, we can improve performance
3. The extra degrees of freedom are useful in two ways:
 - ⇒ Landscape point-of-view : they give access to wider basins
 - ⇒ Representation learning point-of-view : they perform complementary tasks (*feature-extracting* local filters + *discriminative* template matching)

THANK YOU

