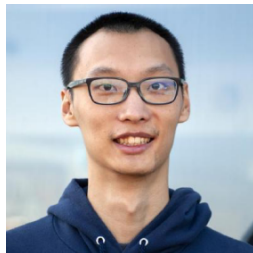


Finding trainable sparse networks through Neural Tangent Transfer



Tianlin Liu^{1 2}

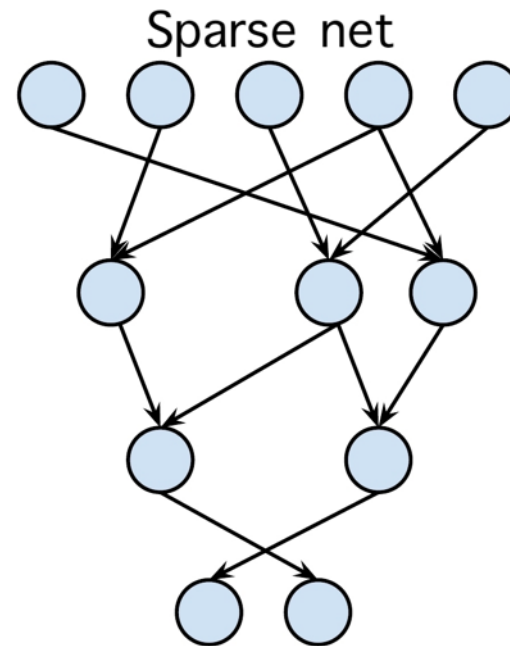
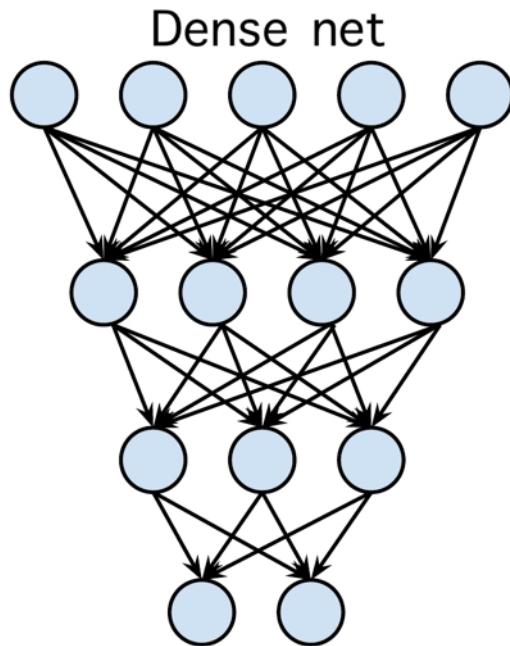


Friedemann
Zenke¹

¹Friedrich Miescher Institute for Biomedical Research, Basel, Switzerland

²University of Basel, Basel, Switzerland

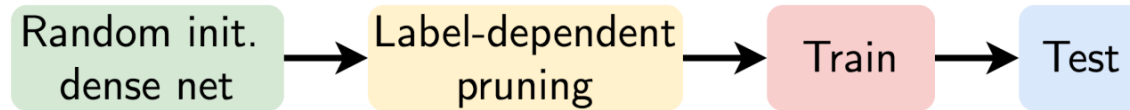
Sparse neural nets



Sparse nets are computationally efficient but difficult to train.

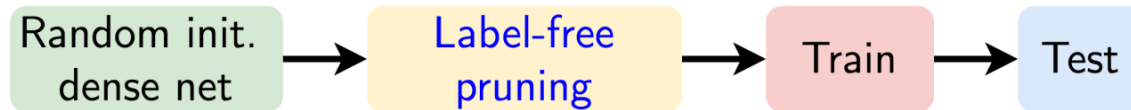
Foresight pruning

■ Existing foresight pruning methods:



Lee, et al. SNIP: Single-shot network pruning based on connection sensitivity, ICLR2019.
Wang, et al. Picking winning tickets before training by preserving gradient flow, ICLR2020.

■ Our new foresight pruning, **Neural Tangent Transfer**:



Advantages:

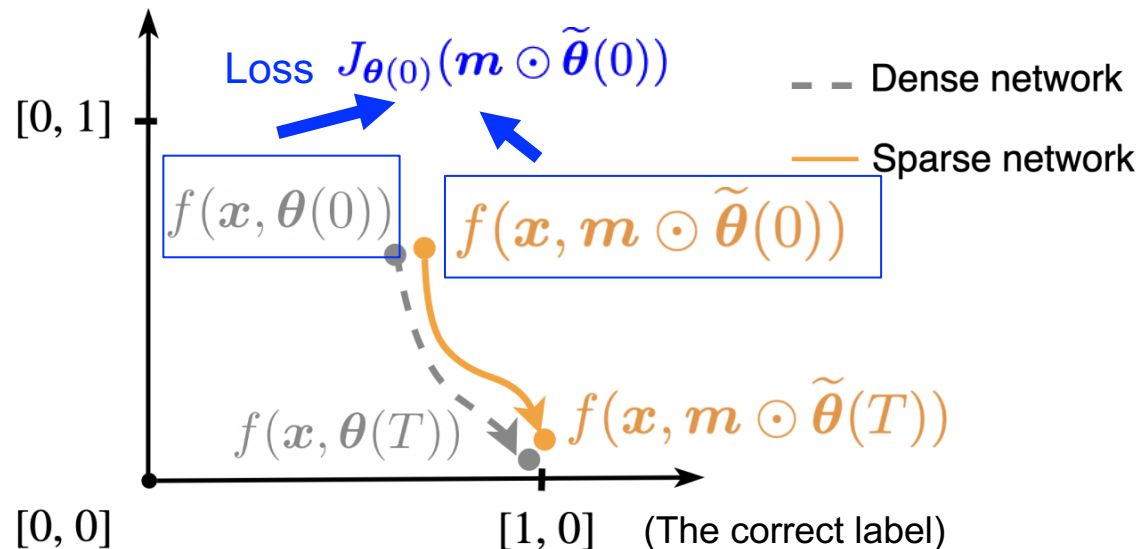
- Does not require labels for pruning.
- Per-layer sparsity levels of pruned networks are controllable.

The idea of our approach

Notation: dense net $f(x, \theta)$ sparse net $f(x, m \odot \tilde{\theta})$

Want: initialize a **sparse net** that is as “trainable” as a **dense net**:

$$f(x, m \odot \tilde{\theta}(T)) \approx f(x, \theta(T))$$



We want a small $J_{\theta(0)}(m \odot \tilde{\theta}(0))$ to
 imply $f(x, m \odot \tilde{\theta}(t)) \approx f(x, \theta(t))$ for all $t > 0$.

Our loss function (first try)

Given \mathcal{X} and $f(\cdot, \theta(0))$, choose $\{m, \tilde{\theta}(0)\}$ to minimize

$$\sum_{t=t_0}^{t_T} \underbrace{\|f(\mathcal{X}, m \odot \tilde{\theta}(t))\|_2}_{\text{Output of the sparse net under supervised training}} - \underbrace{\|f(\mathcal{X}, \theta(t))\|_2}_{\text{Output of the dense net under supervised training}}^2$$

Output of the
sparse net
under
supervised
training

Output of the
dense net
under
supervised
training

**Impossible to
evaluate
without using
labels!**

Modified loss through linearization

Given \mathcal{X} and $f(\cdot, \theta(0))$, choose $\{m, \tilde{\theta}(0)\}$ to minimize

$$\sum_{t=t_0}^{t_T} \left\| \underbrace{f^{\text{lin}}(\mathcal{X}, m \odot \tilde{\theta}(t))}_{\text{Output of the linearized sparse net under supervised training}} - \underbrace{f^{\text{lin}}(\mathcal{X}, \theta(t))}_{\text{Output of the linearized dense net under supervised training}} \right\|_2^2$$

Output of the
linearized sparse
net
under supervised

Output of the
linearized dense
net
under supervised
training

$$f^{\text{lin}}(\mathbf{x}, m \odot \tilde{\theta}) := f(\mathbf{x}, m \odot \tilde{\theta}(0)) + \langle \tilde{\theta} - \tilde{\theta}(0), \nabla_{\tilde{\theta}} f(\mathbf{x}, m \odot \tilde{\theta}(0)) \rangle.$$

$$f^{\text{lin}}(\mathbf{x}, \theta) := f(\mathbf{x}, \theta(0)) + \langle \theta - \theta(0), \nabla_{\theta} f(\mathbf{x}, \theta(0)) \rangle.$$

Lee et al. Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent, NeurIPS, 2019.

The neural tangent transfer loss

Given \mathcal{X} and $f(\cdot, \theta(0))$, choose $\{m, \tilde{\theta}(0)\}$ to minimize

Target distance:

$$J_{\theta(0)}(m \odot \tilde{\theta}(0)) = \frac{1}{n} \left\| f(\mathcal{X}, m \odot \tilde{\theta}(0)) - f(\mathcal{X}, \theta(0)) \right\|_2^2 + \frac{\gamma^2}{n^2} \left\| \mathbf{H}_{m \odot \tilde{\theta}(0)} - \mathbf{H}_{\theta(0)} \right\|_F^2$$

Neural Tangent Kernel (NTK) distance

Jacot et al. Neural Tangent Kernel: Convergence and Generalization in Neural Networks, NeurIPS, 2018.

where

$$\mathbf{H}_{\theta(0)}(i, j) = \left\langle \nabla_{\theta} f(x_i, \theta(0)), \nabla_{\theta} f(x_j, \theta(0)) \right\rangle$$

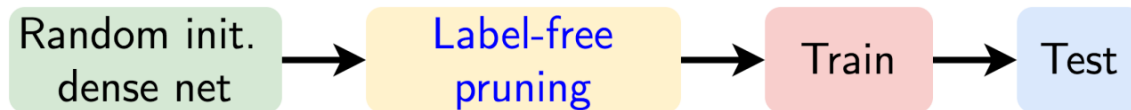
$$\mathbf{H}_{m \odot \tilde{\theta}(0)}(i, j) = \left\langle \nabla_{\tilde{\theta}} f(x_i, m \odot \tilde{\theta}(0)), \nabla_{\tilde{\theta}} f(x_j, m \odot \tilde{\theta}(0)) \right\rangle$$

The algorithm that minimizes this loss is called **Neural Tangent Transfer (NTT)**

Experiment

- **Datasets:** [MNIST](#), Fashion-MNIST, [CIFAR-10](#), SVHN
- **Network architecture:**
 - Lenet-300-100 (MLP)
 - Lenet-5-Caffe (CNN)
 - Conv-4 (a CNN with 4 convolution layers followed by 2 dense layers with dropout)

- **General experimental procedure:**

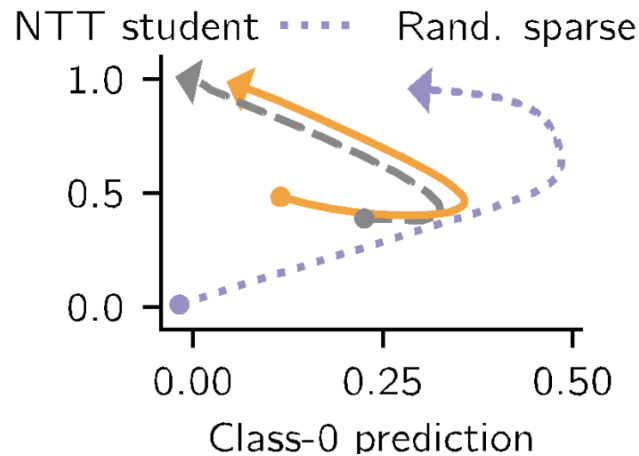
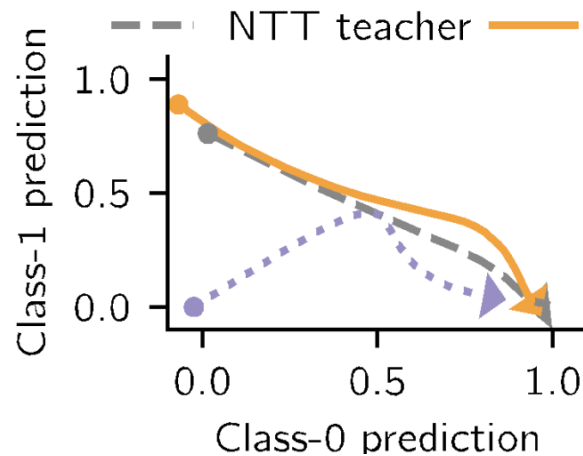
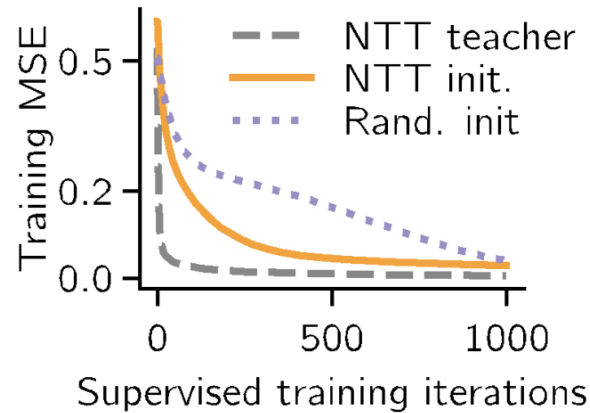
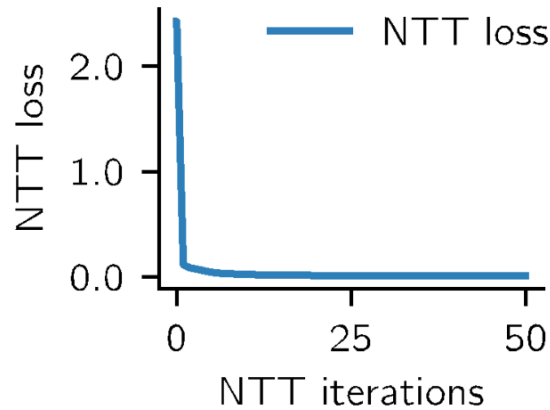


- **Layerwise & global pruning**

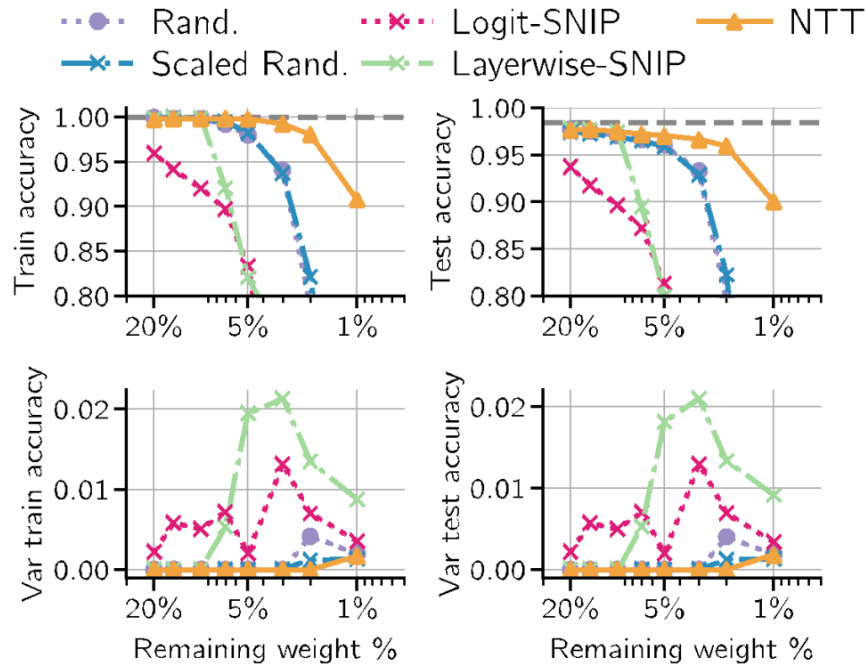
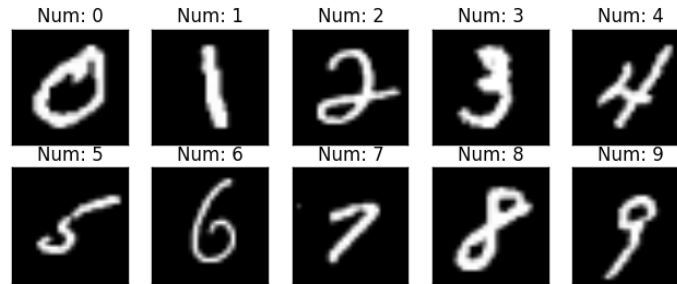
A toy example



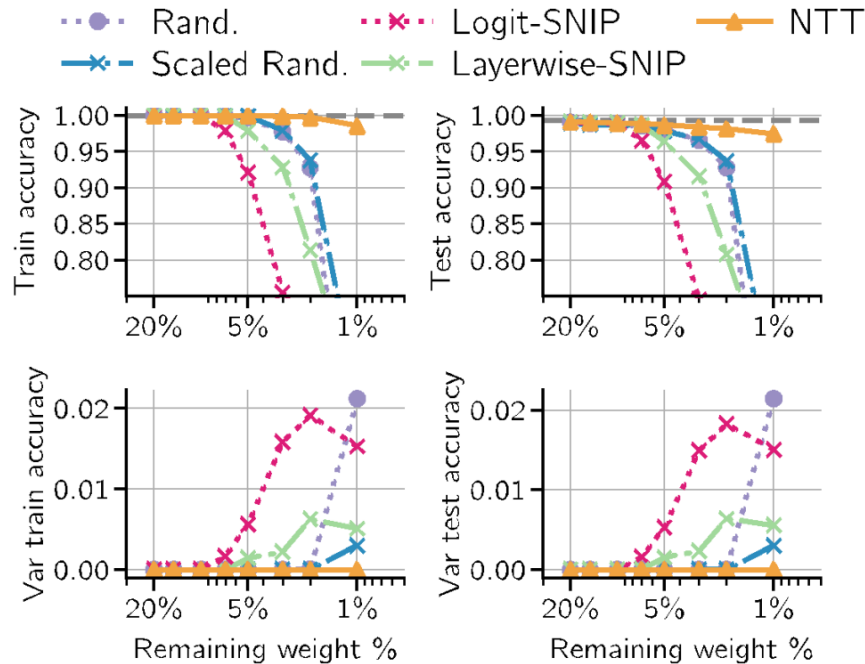
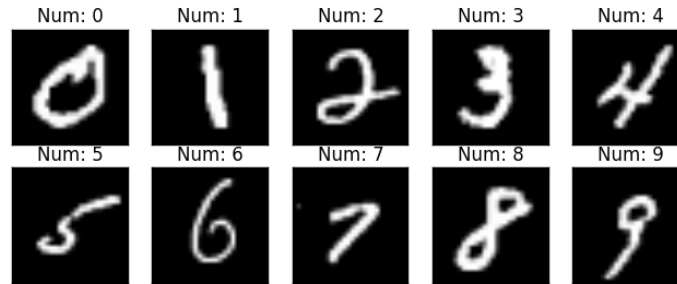
- A small dataset of 0 and 1 digits from MNIST dataset



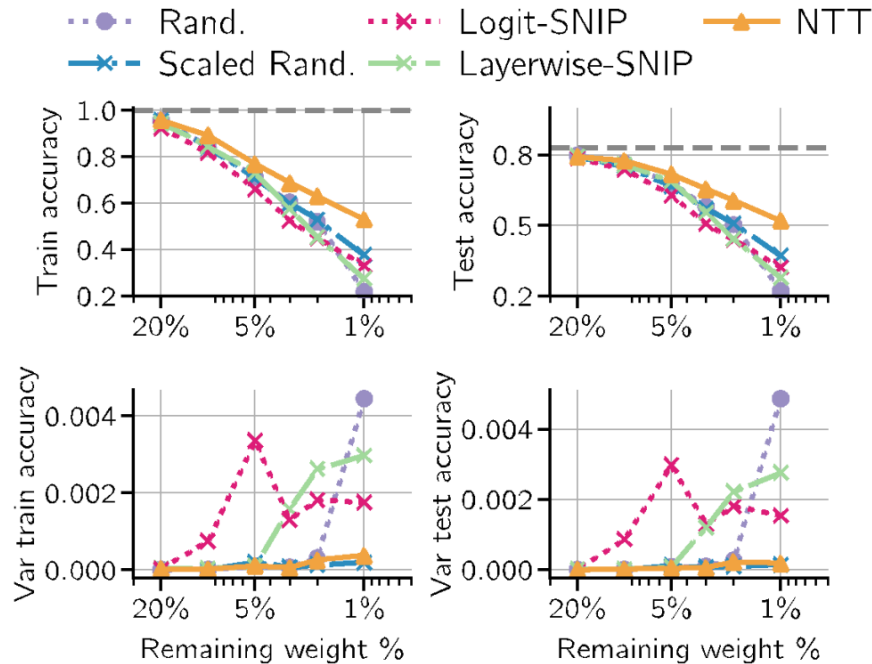
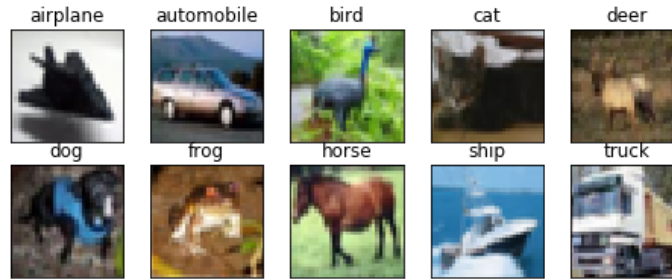
MNIST: MLP Lenet-300-100



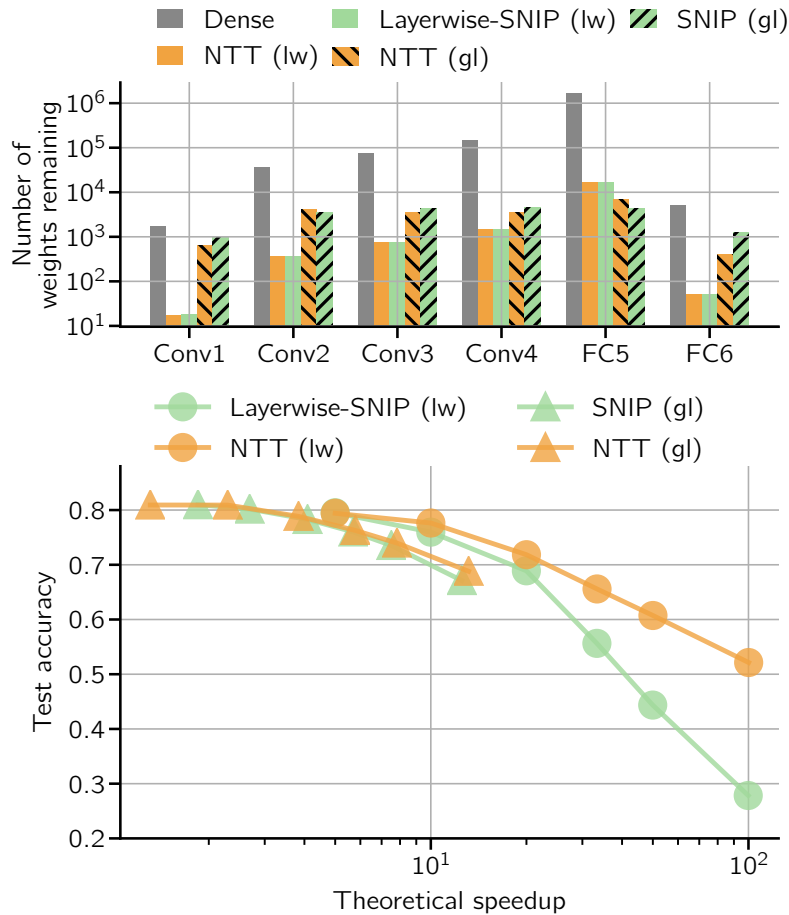
MNIST: CNN Lenet-5-Caffe



CIFAR-10: CNN Conv-4



Layerwise vs global pruning



Conclusions

- We proposed NTT, a label-free method that finds trainable sparse networks.
- Idea: transfer the training dynamics of a dense net onto a sparse net.
- Theoretical handle: Neural Tangent Kernel [Jacot et al. 2018].
- We showed that the resulting sparse nets are highly trainable on supervised learning tasks.
- We showed that NTT can significantly sparsify convolutional layers.