# Coarsening the Granularity: Towards Structurally Sparse Lottery Tickets

[ICML 2022] Tianlong Chen[1], Xuxi Chen[1], Xiaolong Ma[2], Yanzhi Wang[2], Zhangyang Wang[1]

[1]University of Texas at Austin, [2]Northeastern University

The University of Texas at Austin
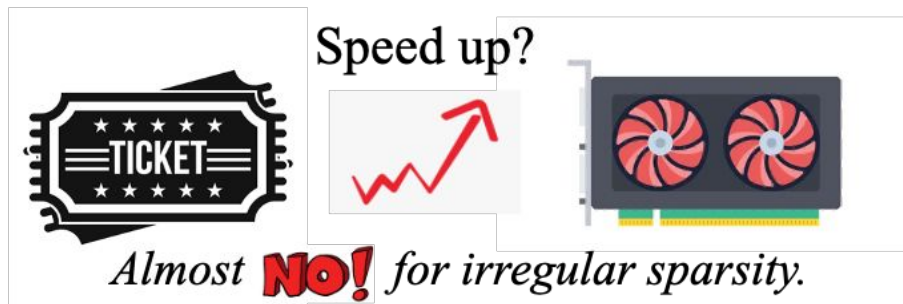**Electrical and Computer Engineering**

# **Agenda**

➢ The Current Limitation of (Unstructured) Lottery Tickets

➢ Insightful Findings

➢ Our Solutions

➢ Our Main Experimental Results

# The Current Limitation of (Unstructured) Lottery Tickets



Speed up?

Almost **NO!** *for irregular sparsity.*
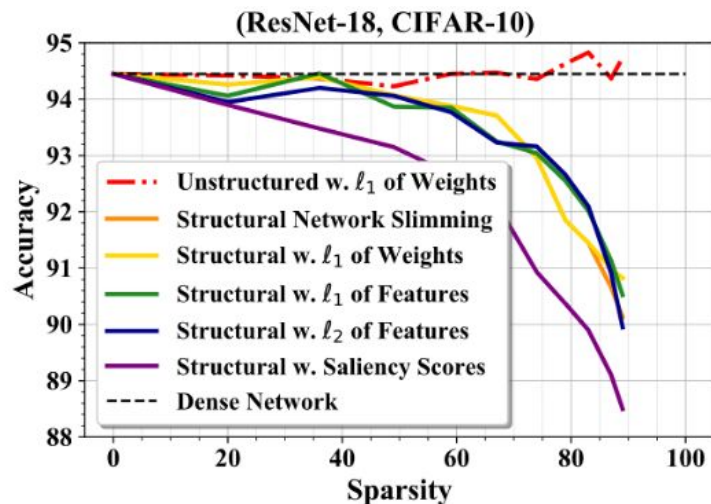
How About Structural Winning Tickets?



Figure 1. Achieved test accuracy over different sparsity levels of diverse unstructured and structural subnetworks. Sparse models from classical channel-wise structural pruning algorithms (He et al., 2017; Liu et al., 2017; Bartoldson et al., 2019; Molchanov et al., 2019) can not match the full accuracy of the dense model.
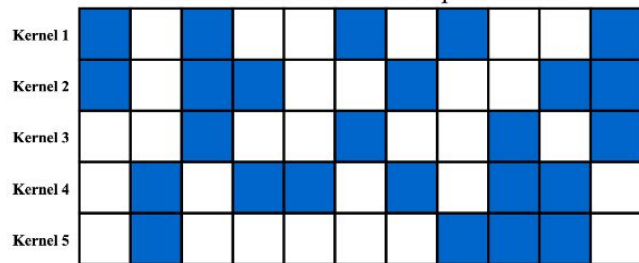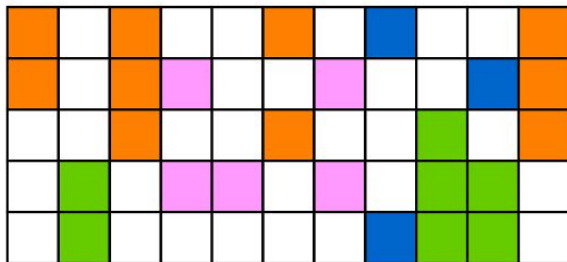
# Insightful Findings

❏ To our best knowledge, we are the first to demonstrate the existence of structurally sparse winning tickets at non-trivial sparsity levels (*i.e.*, > 30%), and with both channel-wise and group-wise sparse patterns.

❏ Extensive experiments validate our proposal on diverse datasets (*i.e.*, CIFAR-10/100, Tiny-ImageNet, and ImageNet) across multiple network architectures, including ResNets, VGG, and MobileNet. Specifically, our structural winning tickets achieve 53.75%~64.93% GPU running time savings at 45%~80% channel- and group-wise sparsity.

# Our Solutions



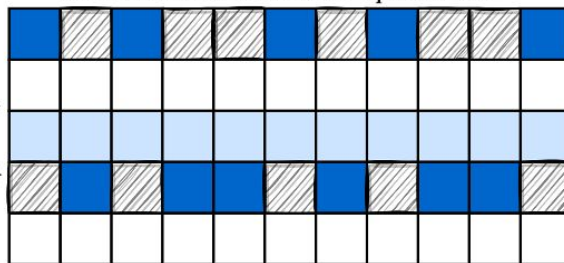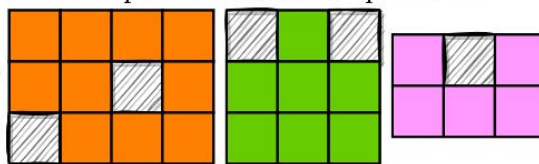Initial **Unstructured** Sparse Mask

Kernel 1
Kernel 2
Kernel 3
Kernel 4
Kernel 5

Kernel Hight × Kernel Width × Input Channel

Refilling or Refilling+

Regrouping

Channel-wise **Structural** Sparse Mask

Group-wise **Structural** Sparse Mask

Remaining Weights | Refilled Weights
Extra Refilled Channel for Refilling+ | Pruned Weights

**Algorithm 2** IMP-Refill(+)

**Input:** $f(x; m \odot \theta_i)$ with unstructured sparsity $s$ (Algo. 1)
**Output:** $f(x; m \odot \theta_i)$ with channel-wise structural mask $m$ at sparsity $\tilde{s}$
1: Calculate importance scores of each channel according to certain criterion
2: Pick top-$k$ channels in $m$, refill back their 0 (pruned) elements with 1 (trainable) and update $m$, maintaining $\tilde{s} \sim s$
3: Pick and refill back extra channels in $m$ with $\tilde{s}^+ < s$
   # Optional for Refill+

**Algorithm 3** IMP-Regroup

**Input:** $f(x; m \odot \theta_i)$ with unstructured sparsity $s$ from Algorithm 1, hyperparameters $t_1$, $t_2$, $b_1$, and $b_2$
**Output:** $f(x; m) \odot \theta_i$ with group-wise structural mask $m$ at sparsity $s^*$
1: **while** dense block can be found **do**
2:   Divide the rows of the sparse pruning mask $m$ into $t_1$ groups using hypergraph partitioning (hMETIS)[a]
3:   **for** group $c_i \in \{c_1, c_2, \ldots, c_{t_1}\}$ **do**
4:     **if** $c_i$ has $\geq b_1$ rows **then**
5:       Select columns in $c_i$ that has no less than $t_2$ non-zero items
6:       **if** $\geq b_2$ columns are selected **then**
7:         Group and Refill the selected columns as well as rows to a dense block, and update $m$
8:       **end if**
9:     **end if**
10:   **end for**
11: **end while**
12: Set other elements out of dense blocks to 0
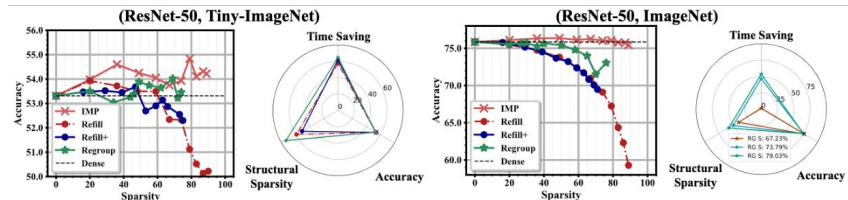
# Our Main Experimental Results



Figure 3. (*Curve plots*) Testing accuracy (%) over network sparsity (%) on Tiny-ImageNet and ImageNet datasets with ResNet-50 (25.56 M). (*Radar plots*) The end-to-end inference time saving of extreme structural winning tickets. Unstructured subnetworks or dense models do not have structural sparsity, and thus they are plotted as dots in the axes of accuracy in the corresponding radar plot. The rightmost plot includes three extreme regroup tickets with accuracy drop < 1%, where "RG S: $x$%" indicates unstructured sparsity before regrouping.
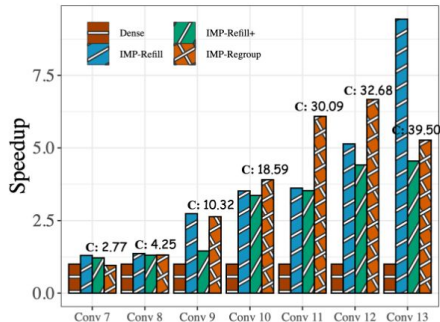


Figure 7. The layer-wise performance of convolution operations in extreme structural winning tickets of (VGG-16, C10). The first six conv. operations are omitted since there is no meaningful speedup, coincided with Rumi et al. (2020). Marks like "C: 2.77" indicate the layer-wise compression ratio of IMP-Regroup.
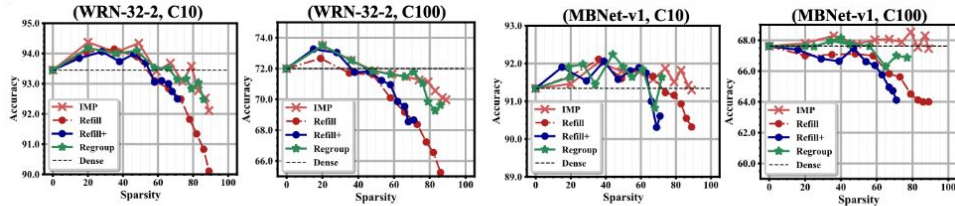


Figure 4. Testing accuracy (%) over sparsity (%) on CIFAR-10/100 with Wide-ResNet-32-2 (1.86 M) and MobileNet-v1 (3.21 M).
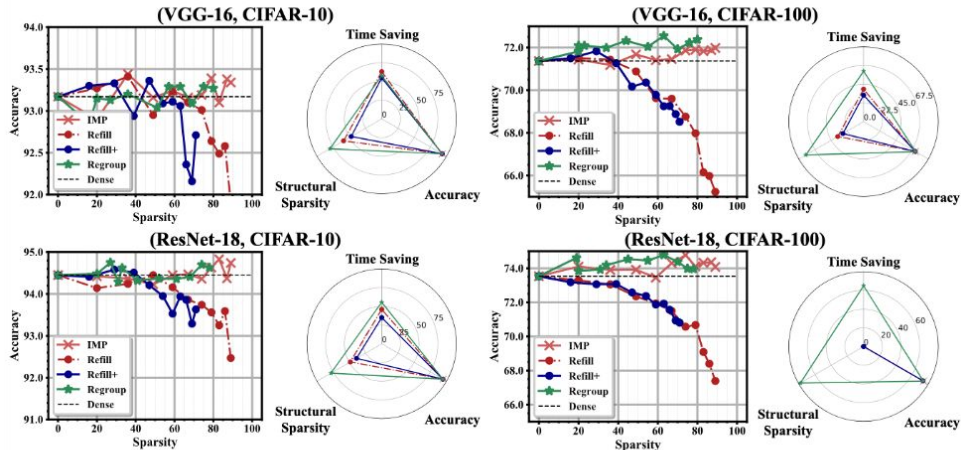


Figure 5. (*Curve plots*) Testing accuracy (%) over sparsity (%) on CIFAR-10/100 with large models VGG-16 (14.72 M) and RN-18 (11.22 M). (*Radar plots*) The end-to-end inference time saving of extreme structural winning tickets. Note that unstructured subnetworks or dense models do not have structural sparsity, and thus they are plotted as dots in the axes of accuracy in the corresponding radar plot.

The University of Texas at Austin
**Electrical and Computer Engineering**
Cockrell School of Engineering

Northeastern University

VITA    GitHub

ICML
International Conference On Machine Learning

# Coarsening the Granularity: Towards Structurally Sparse Lottery Tickets

Tianlong Chen[1], Xuxi Chen[1], Xiaolong Ma[2], Yanzhi Wang[2], Zhangyang Wang[1]

[1]University of Texas at Austin, [2]Northeastern University

## ➤ The Current Limitation of (Unstructured) Lottery Tickets

Speed up?

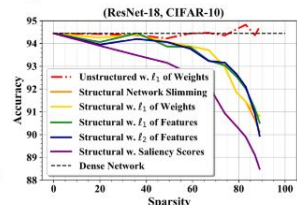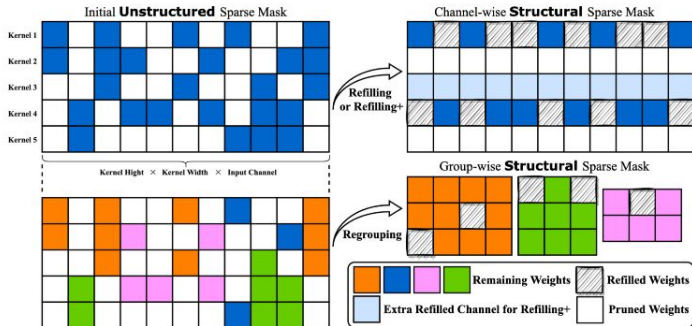*Almost* NO! *for irregular sparsity**



Figure 1. Achieved test accuracy over different sparsity levels of diverse unstructured and structural subnetworks. Sparse models from classical channel-wise structural pruning algorithms (He et al., 2017; Liu et al., 2017; Bartoldson et al., 2019; Molchanov et al., 2019) can not match the full accuracy of the dense model.

## ➤ Insightful Findings

❑ To our best knowledge, we are the first to demonstrate the existence of structurally sparse winning tickets at non-trivial sparsity levels (*i.e.*, > 30%), and with both channel-wise and group-wise sparse patterns.

❑ Extensive experiments validate our proposal on diverse datasets (i.e., CIFAR-10/100, Tiny-ImageNet, and ImageNet) across multiple network architectures, including ResNets, VGG, and MobileNet. Specifically, our structural winning tickets achieve 53.75% ~ 64.93% GPU running time savings at 45% ~ 80% channel- and group-wise sparsity.

## ➤ Our Solutions

❖ Refilling
❖ Refilling+
❖ Regrouping

* R.K. Some packages like XNNPACK can accelerate unstructured sparse neural networks on certain hardware platforms like smartphone processors.
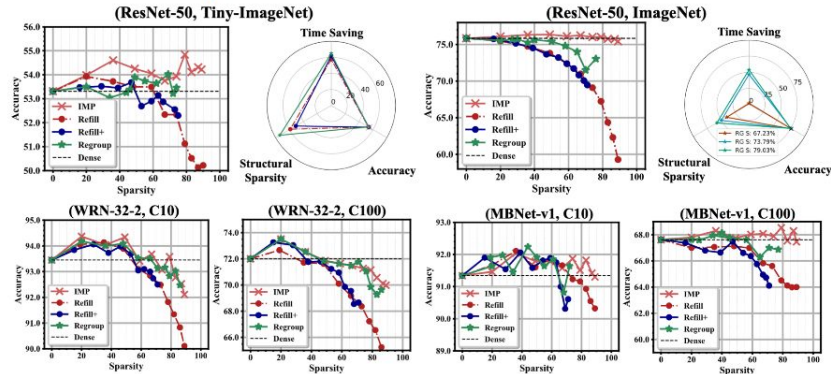


## ➤ Our Main Experimental Results



Figure 4. Testing accuracy (%) over sparsity (%) on CIFAR-10/100 with Wide-ResNet-32-2 (1.86 M) and MobileNet-v1 (3.21 M).
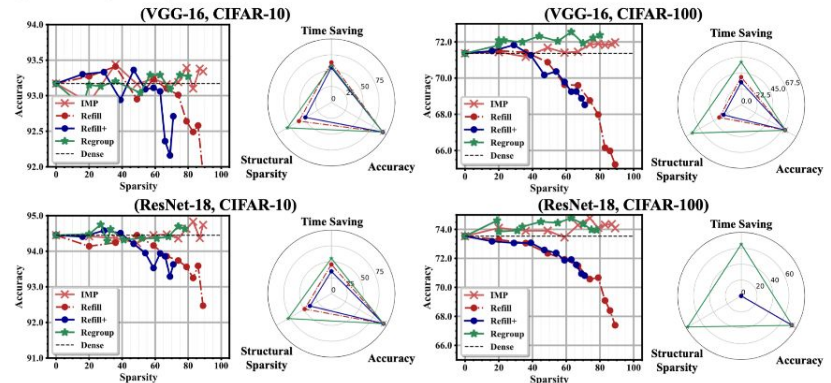


Figure 5. (*Curve plots*) Testing accuracy (%) over sparsity (%) on CIFAR-10/100 with large models VGG-16 (14.72 M) and RN-18 (11.22 M). (*Radar plots*) The end-to-end inference time saving of extreme structural winning tickets. Note that unstructured subnetworks or dense models do not have structural sparsity, and thus they are plotted as dots in the axes of accuracy in the corresponding radar plot.

Q&A