

Motivation

Can we learn Sparse Neural Networks (SNNs) with high generalization performance that are easy to accelerate?

- Unstructured Dynamic Sparse Training (DST) matches the generalization of dense models with 85–95% fewer weights; however, accelerating unstructured SNNs is challenging.
- Structured sparsity is easy to accelerate, but results in worse generalization performance.

Method

- 1. Structured RigL (SRigL) is a novel sparse-to-sparse DST method that learns a SNN with constant fan-in sparsity.
- 2. SRigL ablates entire neurons to match the generalization performance of unstructured sparsity at high sparsities.
- 3. The sparse representation learned by SRigL is parameter- and memory-efficient and amenable to real-world acceleration.



Figure 2. Neuron ablation. At sparsity levels over 90%, RigL learns to completely mask (ablate) a large number of neurons within each layer, effectively reducing layer width. Allowing SRigL to ablate neurons restores RigL-level performance, even at high sparsities.

Constant fan-in theoretical analysis





Dynamic Sparse Training with Structured Sparsity

¹University of Calgary ²Massachusetts Institute of Technology ³IAIFI ⁴Google Research ⁵University of Guelph ⁶Vector Institute for AI





Neuron ablation



Figure 8. Percentage active neurons (i.e., not ablated) following RigL/SRigL training on ResNet-50/ImageNet.

Mike Lasby¹, Anna Golubeva^{2,3}, Utku Evci⁴, Mihai Nica^{5,6}, Yani A. Ioannou¹

Results













end for

end for

12:

13: end for

Acceleration of constant fan-in sparsity

Algorithm 1 "Condensed" linear layer with constant fan-in sparsity forward pass

UNIVERSITY of GUELPH

1: Input: x: the input matrix of shape (batch_size, num_features) w: the condensed weight matrix of shape (active_neurons, constant_fan_in) indx: indices of non-zero dense weights of shape (active_neurons, constant_fan_in) 4: output \leftarrow torch.zeros(size=(batch_size, neurons)) : for b in range(batch size) do for n in range(neurons) do for k in range(constant_fan_in) do $source_idx \leftarrow idx[n, k]$ feature \leftarrow x[b, source_idx] output[b, n] += feature * w[n, k]

▷ For each sample in mini-batch ▷ For each active neuron in layer ▷ For each non-zero weight

Figure 11. Batched GPU inference with batch size of 2048 on an NVIDIA Titan V. At 90% sparsity, our condensed representation is 1.7× faster than dense and $13.0 \times$ faster than unstructured (CSR) sparse layers. Note y-axis is log-scaled.