

# ScRRAMBLE: Block-Sparse Deep Learning Architecture for Analog Compute-in-Memory Accelerators

Vikrant Jaltare<sup>1,2</sup>, Rommani Mondal<sup>1,2</sup>, Leif Gibb<sup>2</sup>, Jiahao Song<sup>2</sup>, Johannes Leugering<sup>3</sup> and Gert Cauwenberghs<sup>1,2</sup>

<sup>1</sup> Shu Chien-Gen Lay Department of Bioengineering, UC San Diego, <sup>2</sup> Institute for Neural Computation, UC San Diego, <sup>3</sup> Neuromorphic Compute Nodes (PGI-14) Forschungszentrum Jülich, Germany



for more info!

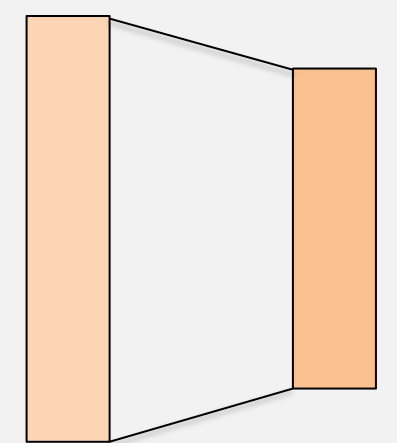
## Promise of Analog Compute-in-Memory (CIM)\*

- Matrix-vector multiplication in  $\mathcal{O}(1)$
- High-density of conductances  $\Rightarrow$  large capacity of parameters.
- Non-volatile storage.
- Analog multi-level weights  $\Rightarrow$  inexpensive high-precision weight quantization.

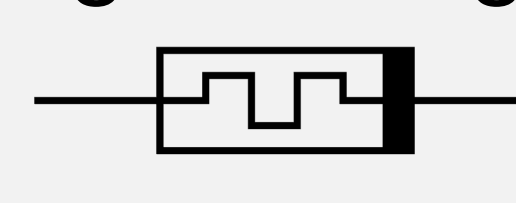
\* applies to conductance-based/memristive technologies

## Computational Bottlenecks

#weights > core-size

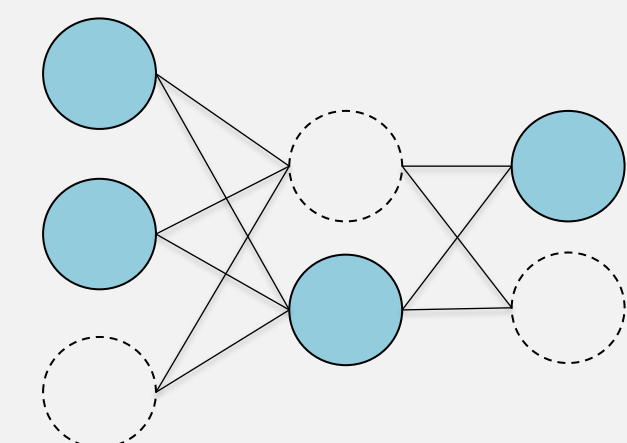


positive conductance // signed weights



area on Si cost

unstructured sparsity



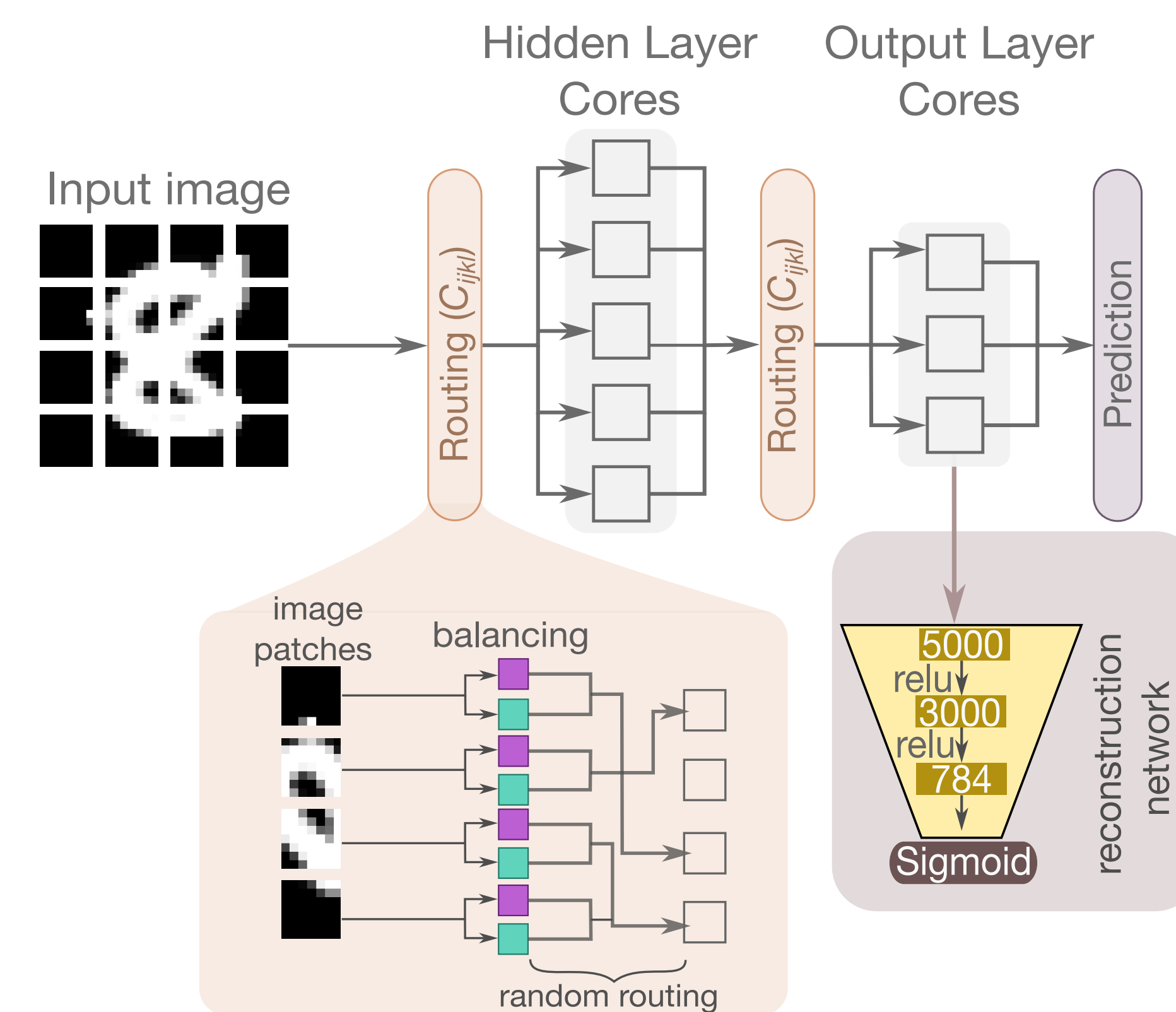
digital controller

communication cost

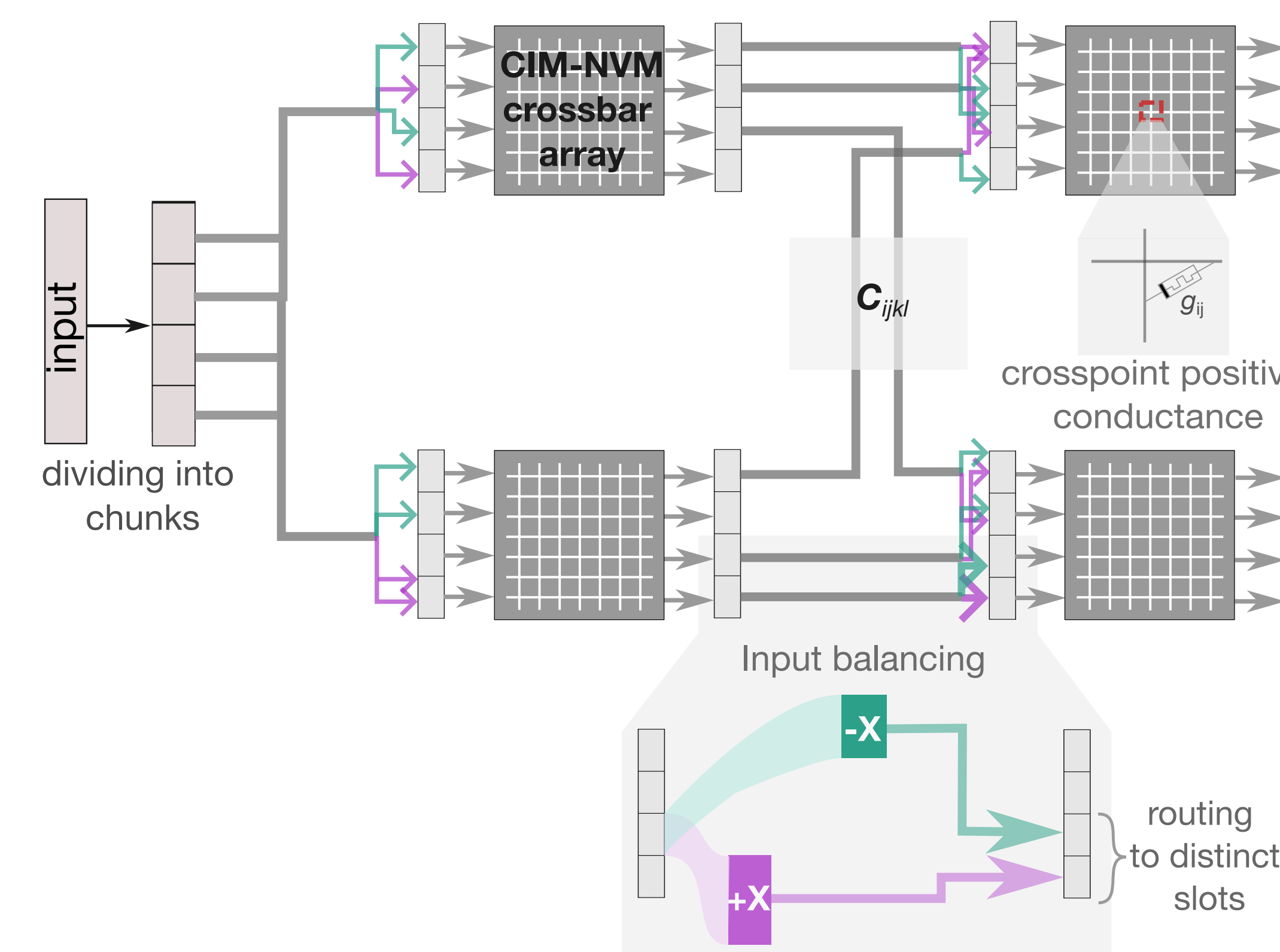
ScRRAMBLE leverages block-sparsity to fully utilize dense cores and reduce inter-core communication.

## Neural Network and Hardware Co-Design

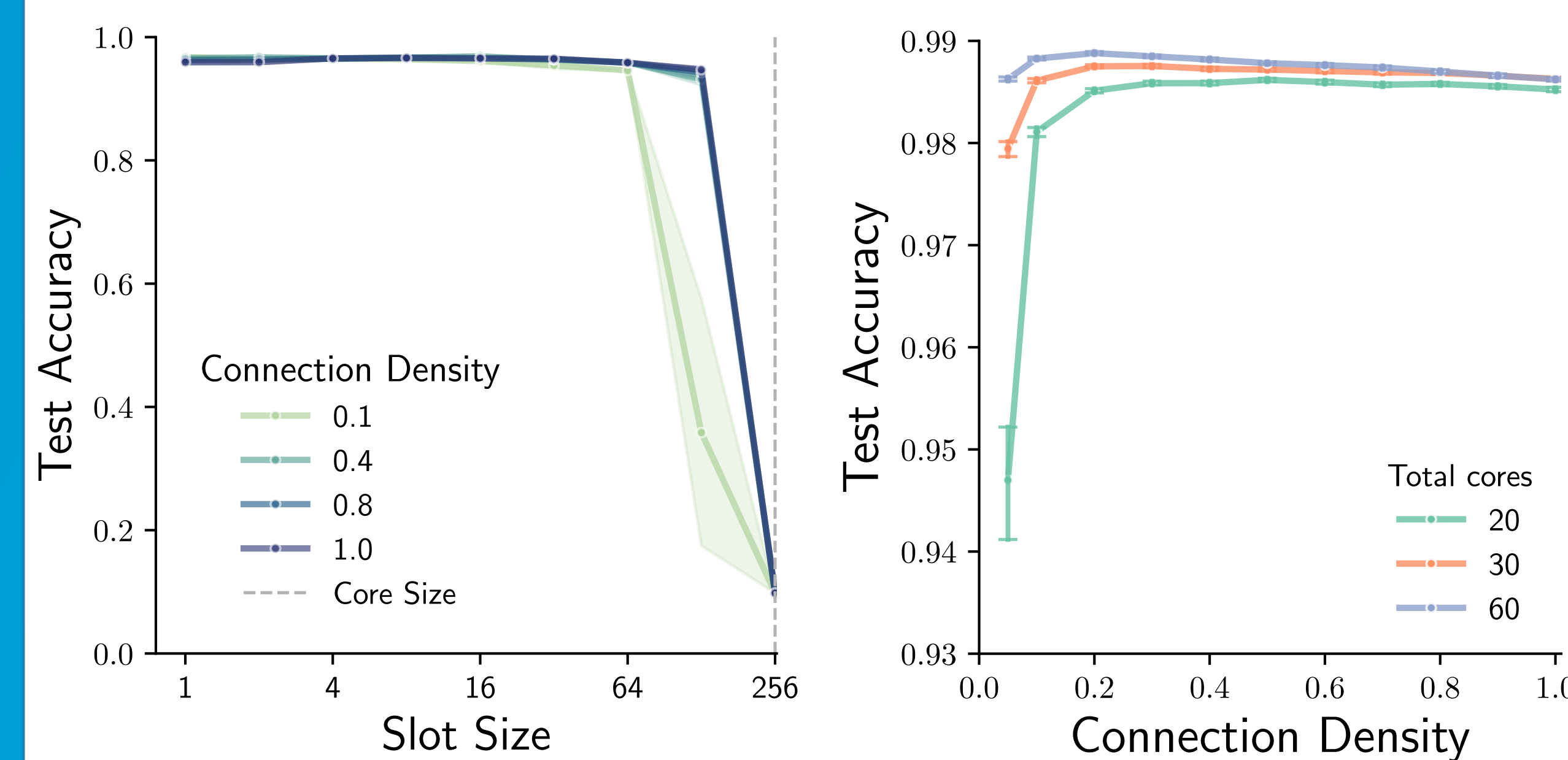
### Network Architecture



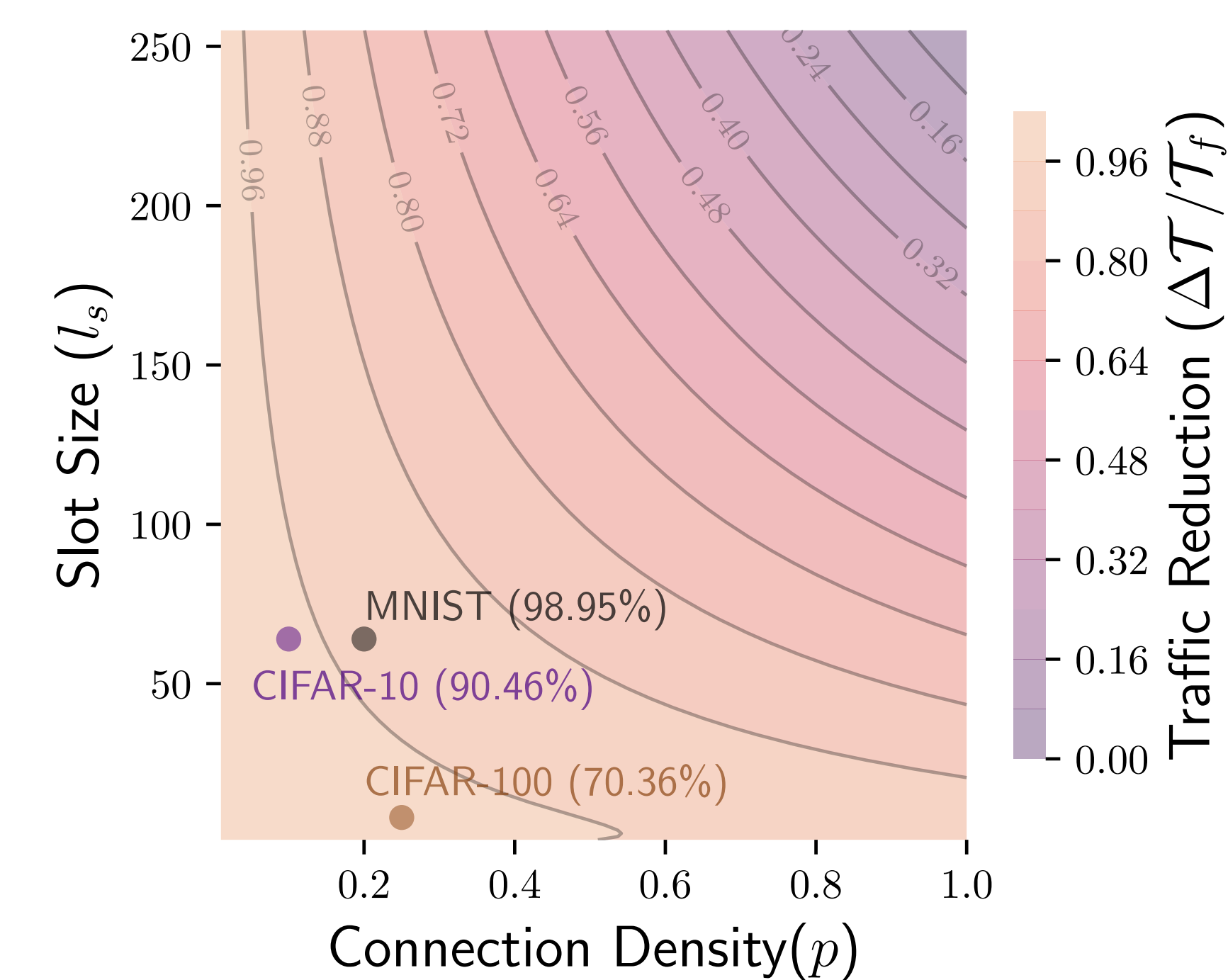
### Hardware Mapping



## Block-Sparsity and Performance

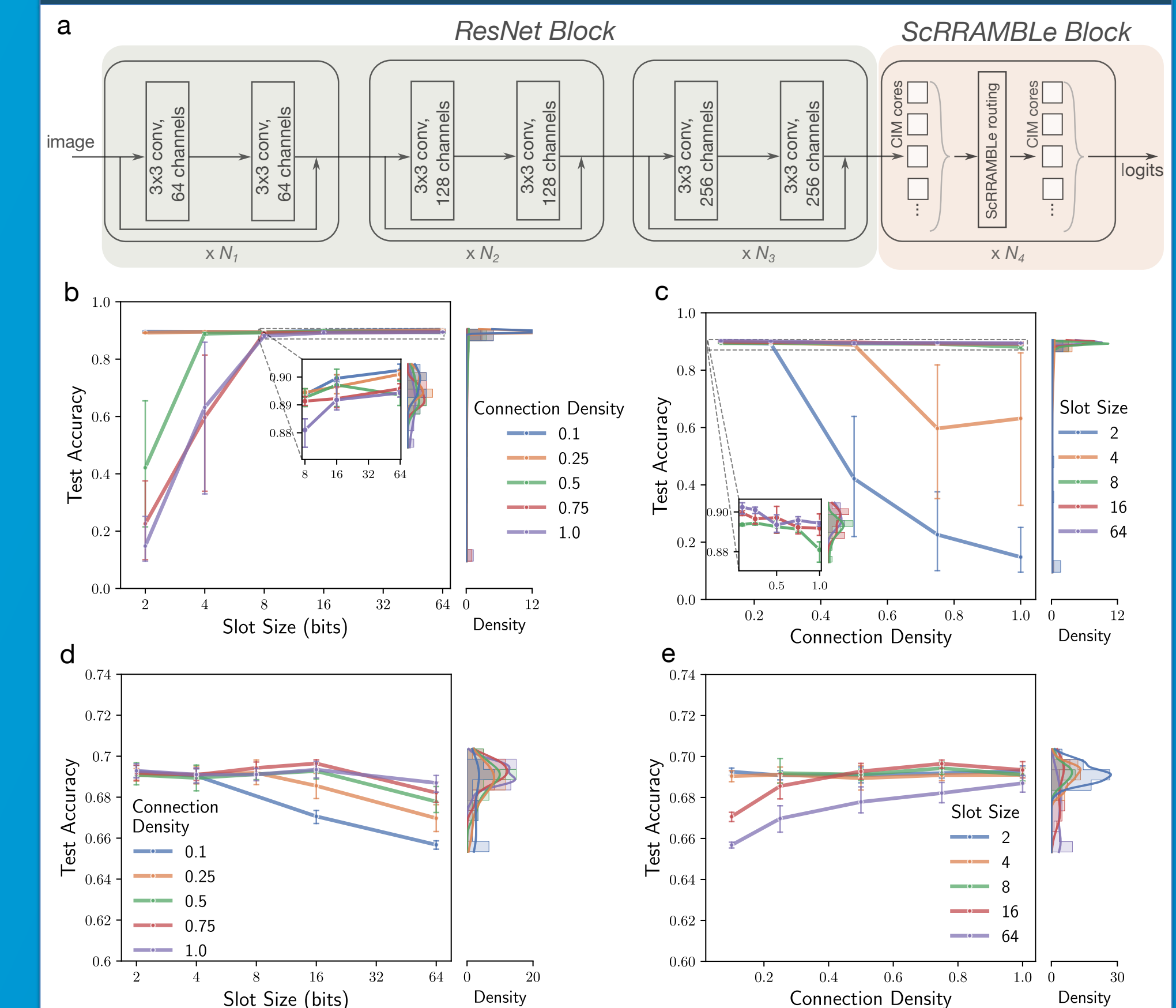


## Intercore communication



Block-sparse neural networks together with input-balanced inter-block routing can replace computationally expensive fully-connected layers in analog CIM accelerators while implementing a signed weight per positive crossbar conductance.

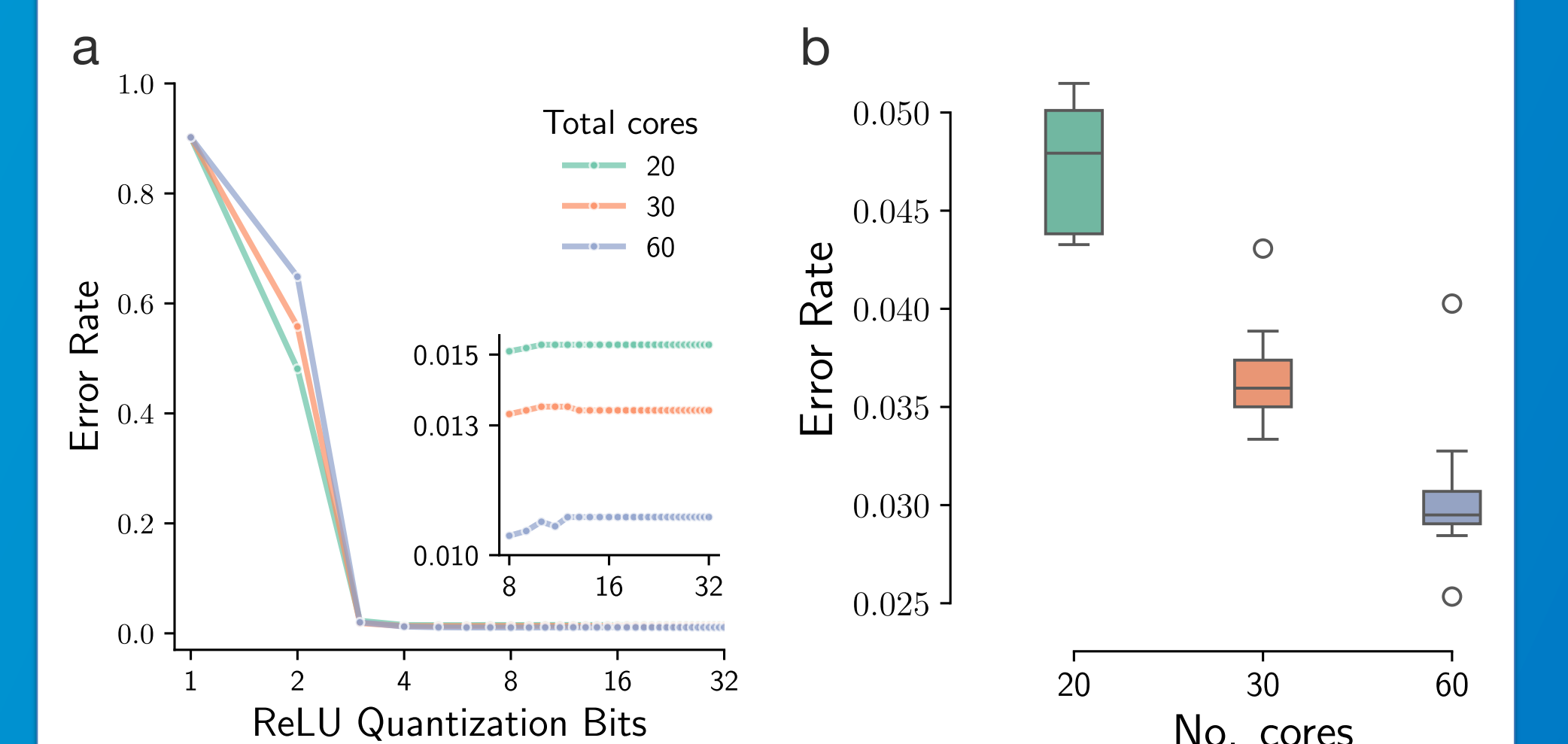
## Modular Scaling



ScRRAMBLE layers can be integrated with other layers and scaled up to deeper neural networks.

## Quantization Robust Training

### Naïve Post-training Quantization vs 1-bit Straight Through Estimator



Quantization Aware Training methods like Straight Through Estimator rescue performance in networks with extreme activation quantization. Future studies are needed to test scalability of this approach.