

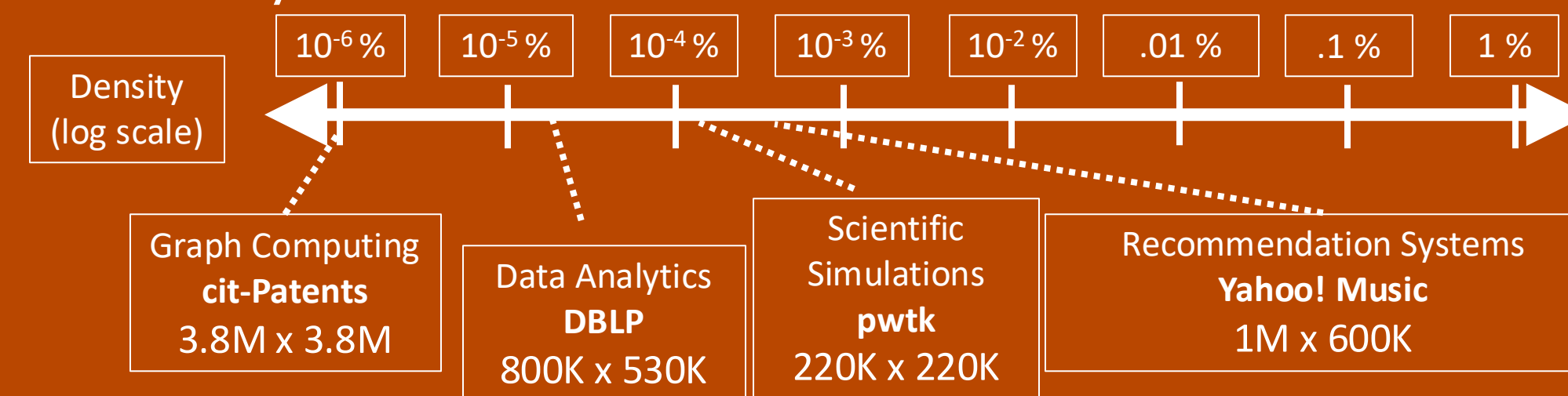
# Tailors: Accelerating Sparse Tensor Algebra by Overbooking Buffer Capacity

Zi Yu Xue\*, Yannan Nellie Wu\*, Joel S. Emer\*<sup>^</sup>, Vivienne Sze\*

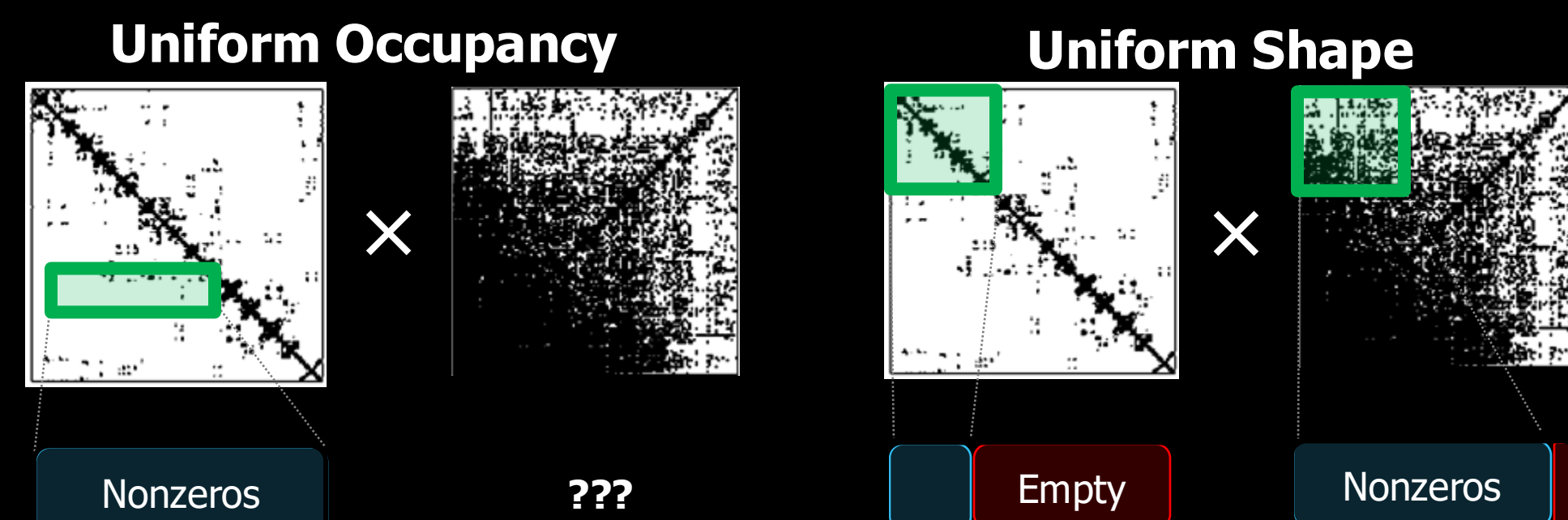
MIT\*, NVIDIA<sup>^</sup>

## Sparse Tensors are Large and Highly Sparse

- Tensor computation relies on tiling to improve data reuse and arithmetic intensity. **Larger tiles maximize data reuse.**
- Operations on sparse tensors, particularly multiple sparse operands, are especially challenging to tile effectively as they have further reduced arithmetic intensity



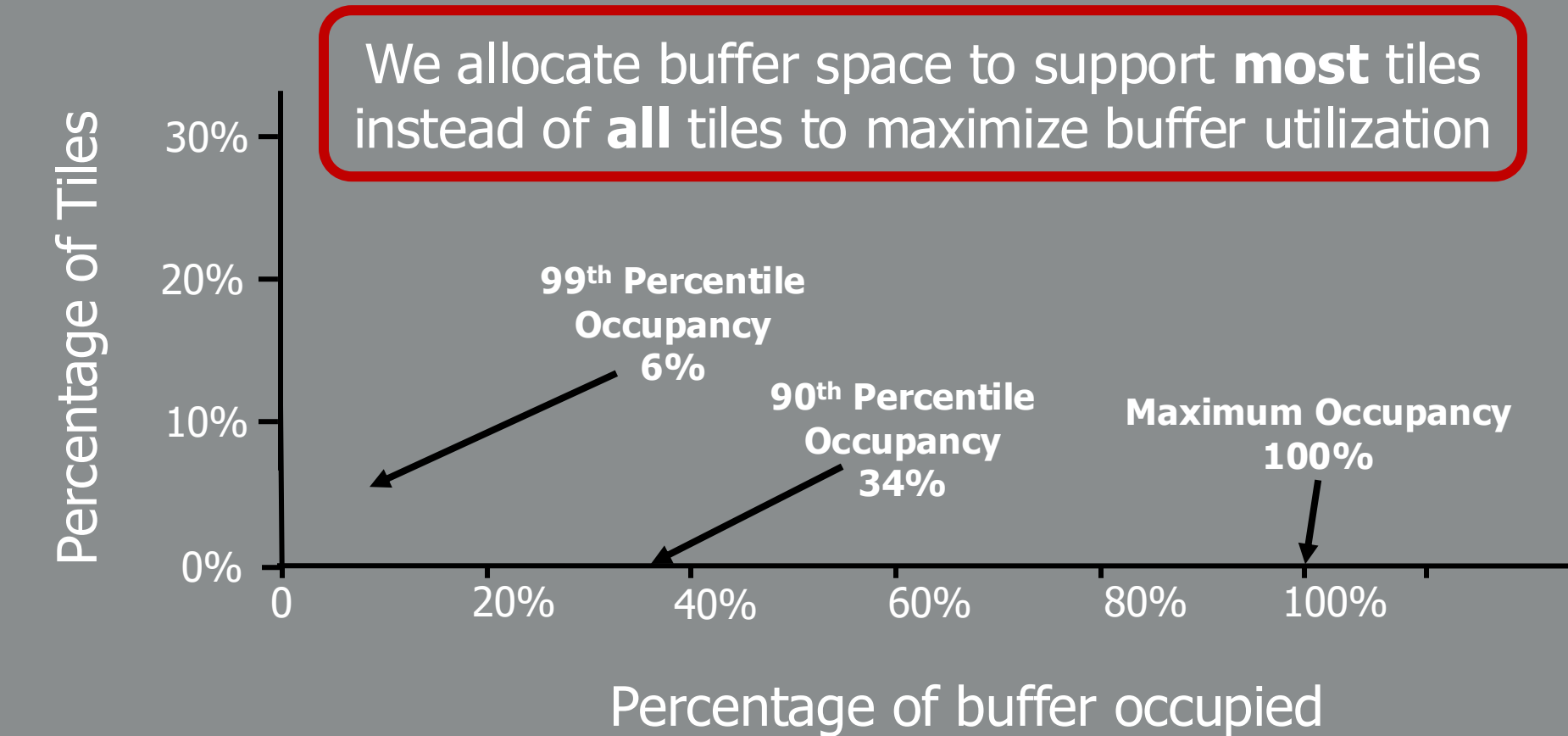
## Current Tiling Approaches are Insufficient



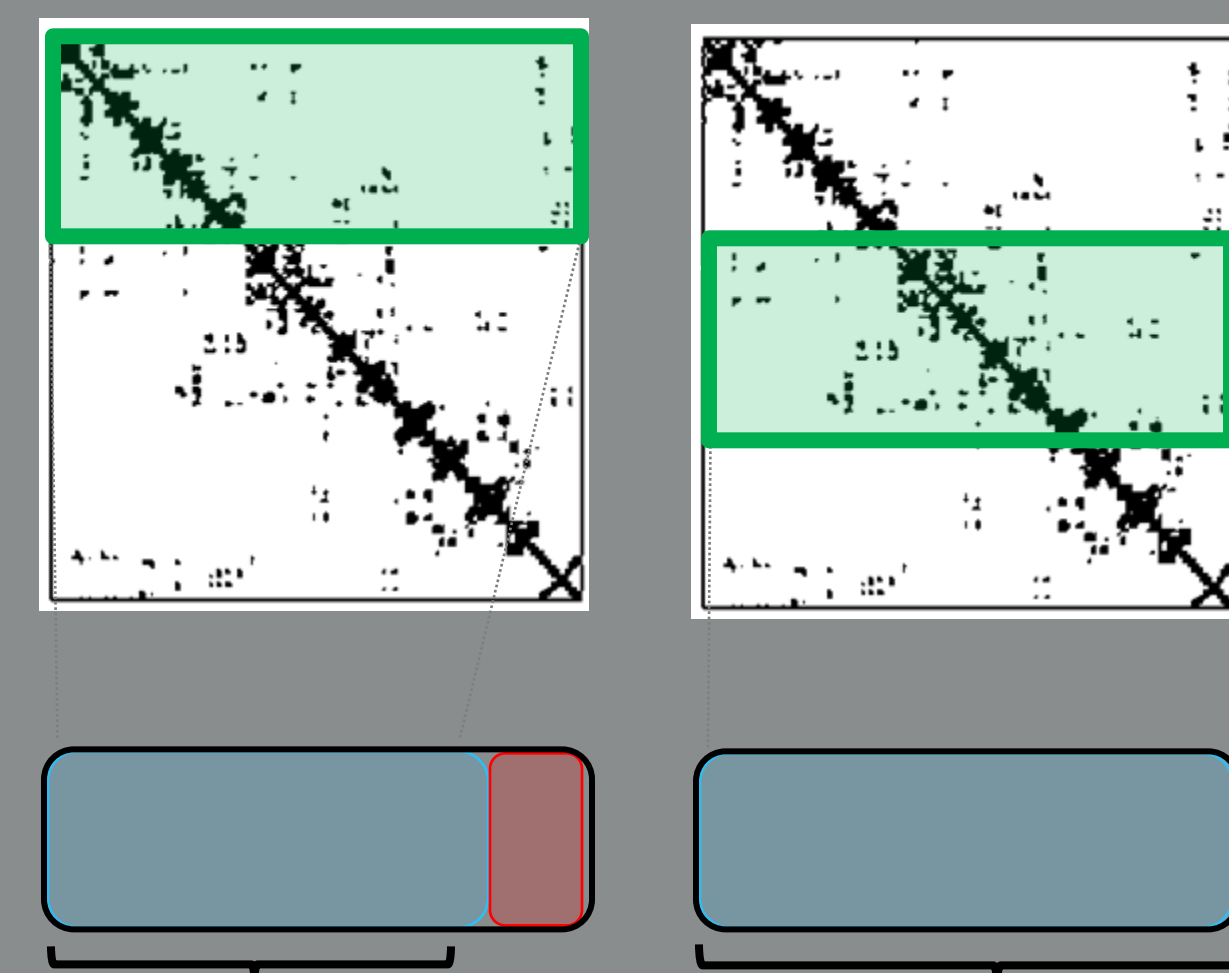
- Same number of nonzeros per tile  $\Rightarrow$  Ideal buffer utilization
- Varying coordinate range in second sparse operand  $\Rightarrow$  Hard to tile second operand
- All tiles must fit in buffer  $\Rightarrow$  Low buffer utilization
- Fixed coordinate ranges  $\Rightarrow$  Easy to tile both operands

## Opportunities for Overbooking

### Tiles in Sparse Tensors Vary in Sparsity



### Overbooking-based Coordinate-space Tiling



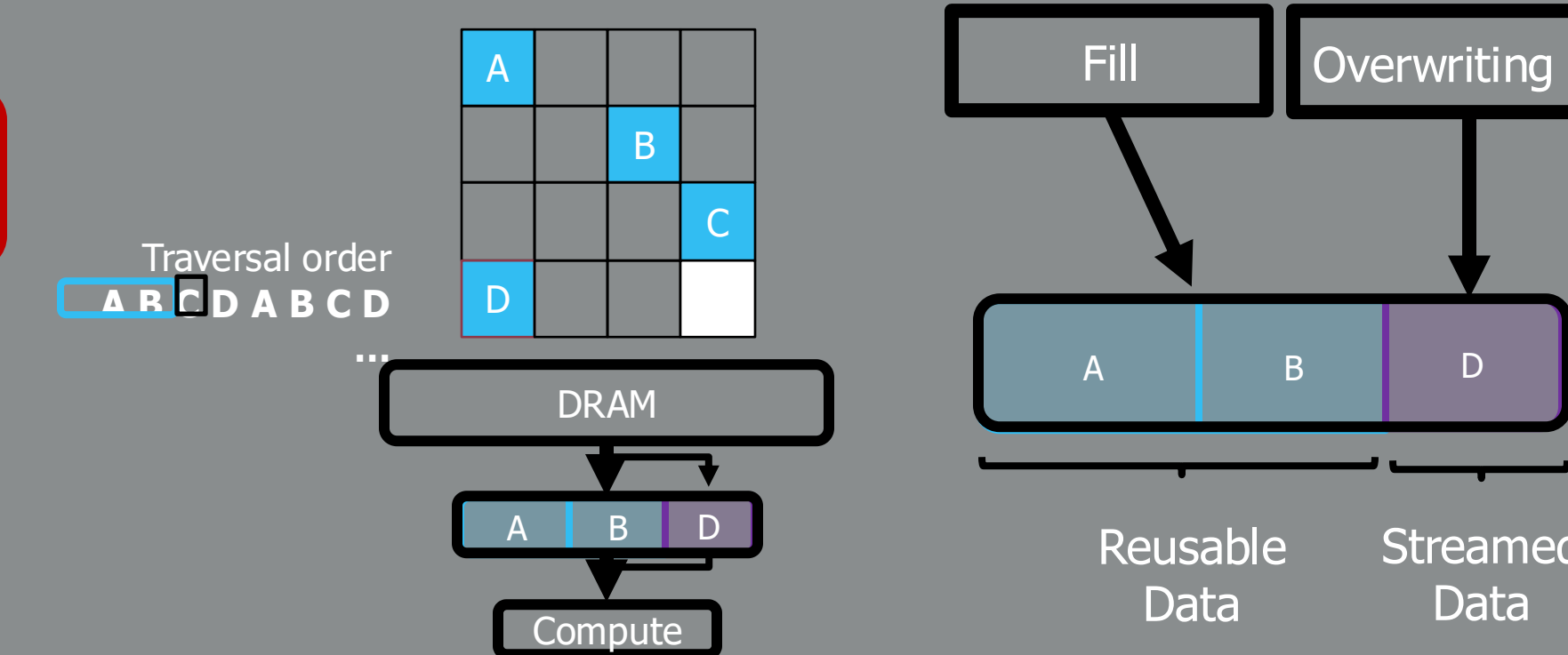
Larger tile size makes buffer more filled on average

Occasionally, tile occupancy will bump data ("overbook")

Example:  
Larger  $M_0$  (larger A, Z tiles)  $\Rightarrow \frac{M_0^{new}}{M_0}$  fewer loads of B tiles

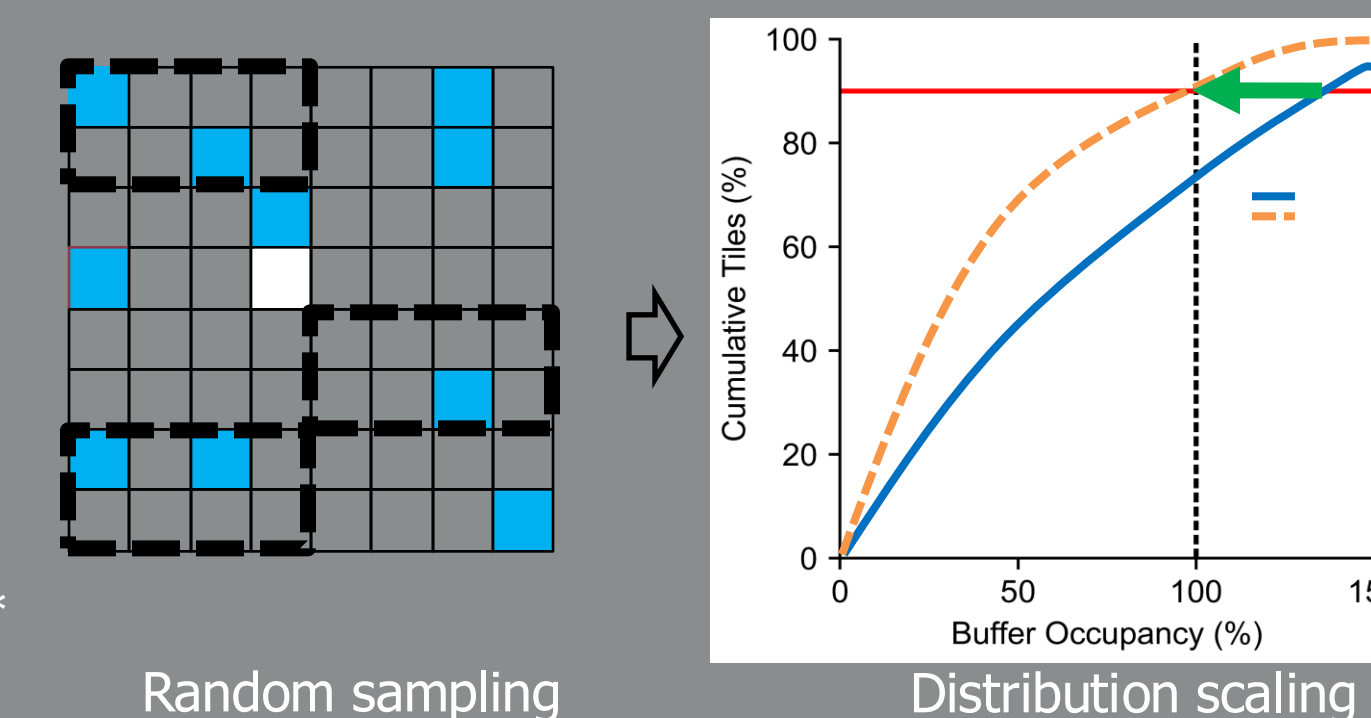
As long as the cost of recovering from bumped data is less than the benefit from larger tiles, overbooking is beneficial

## Tailors: Tail Overbooked Buffers



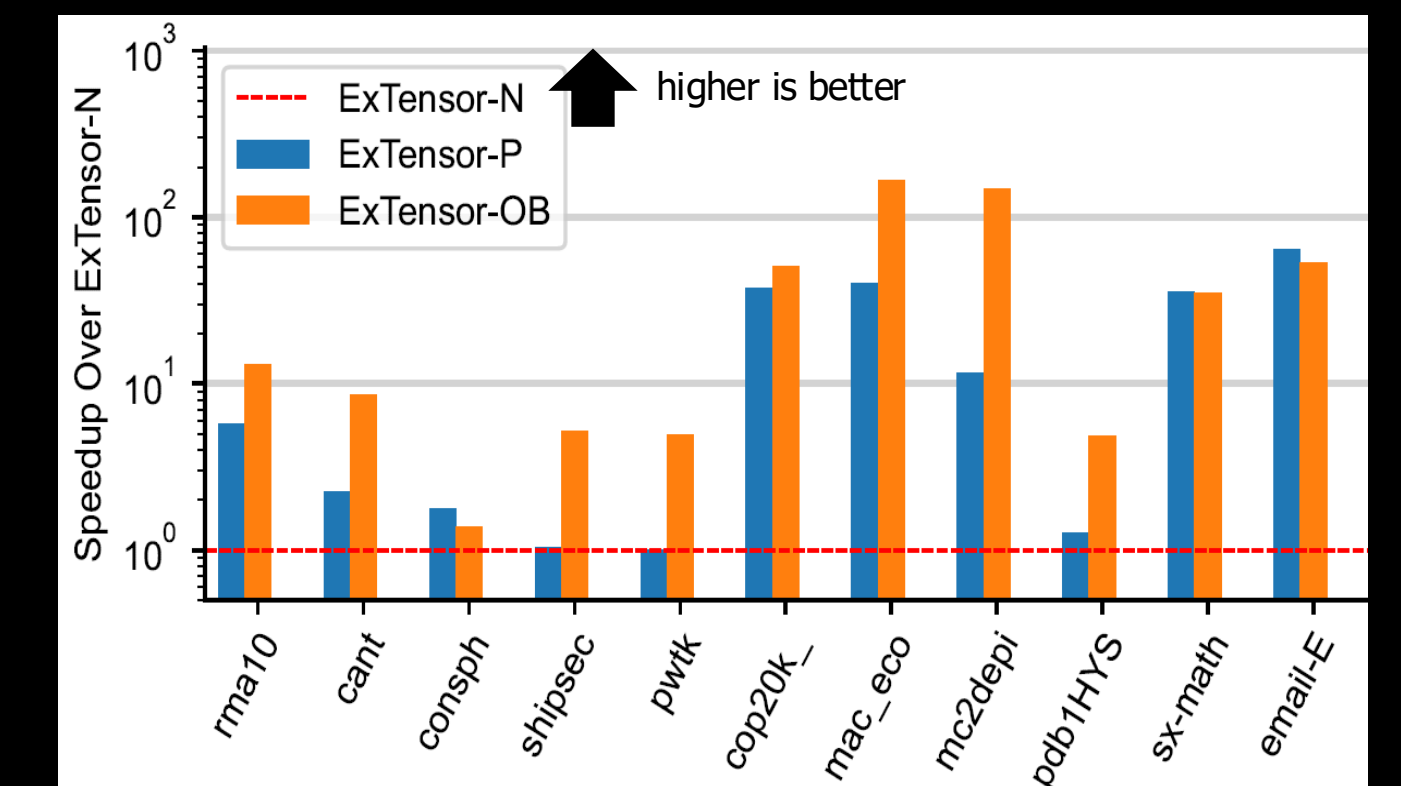
When tile is overbooked, divide the buffer to stream data through without losing data reuse  
Hardware support is **low cost** and **low complexity**, easy to integrate into existing accelerators

## Swiftiles: Swift Tiling with Overbooking



Minimize preprocessing by sampling tile occupancy distribution and scaling to buffer size

## Results



52.7x speedup over ExTensor without tensor occupancy knowledge  
2.3x speedup over ExTensor with perfect tensor occupancy knowledge

## Conclusion

Tiling is key to improving data reuse and thus reducing memory traffic for sparse tensor algebra applications. We balance the tiling strategy's adaptability and efficiency by overbooking. We support overbooking in hardware with Tailors and speculatively tile with Swiftiles

## Acknowledgements

MIT AI Hardware Program  
NSERC PGS-D