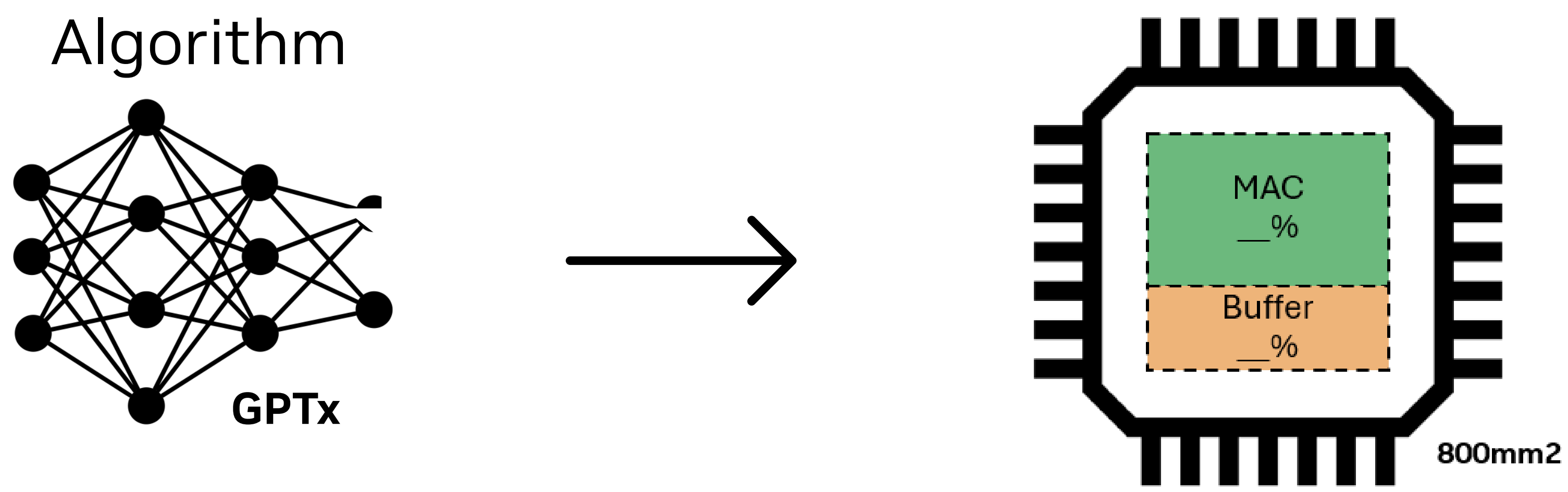


Mind the Gap: Attainable Data Movement and Operational Intensity Bounds for Tensor Algorithms



Qijing Huang, Po-An Tsai, Joel S Emer, Angshuman Parashar
jennyhuang@nvidia.com, poant@nvidia.com, emer@csail.mit.edu

A Key Design Challenge



How to provision chip area between storage and compute?

DSE ?

Roofline Analysis ?

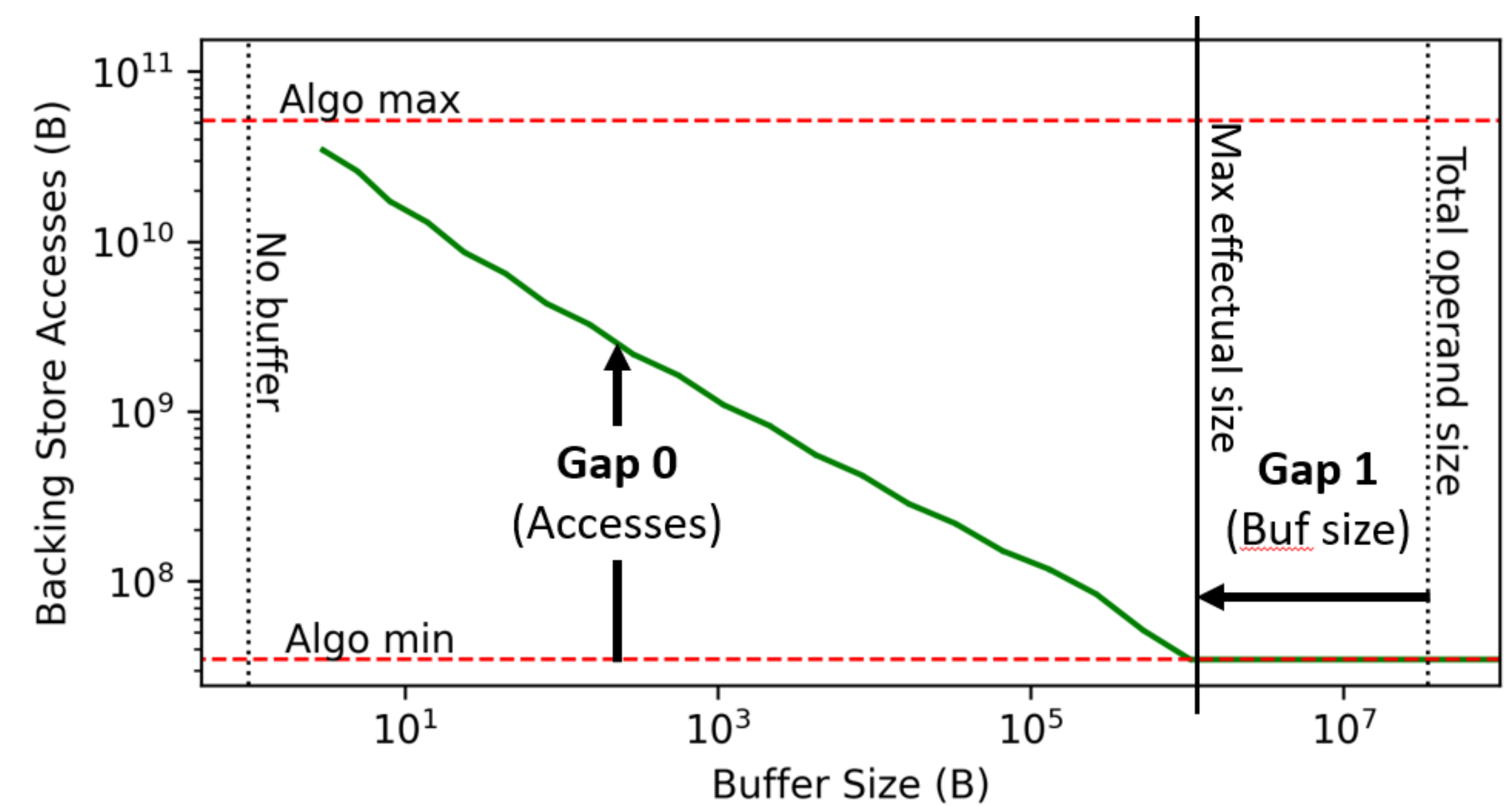
Buffer storage affects the achievable data movement

- Expensive
- Lack of insights

- No buffer size suggestion

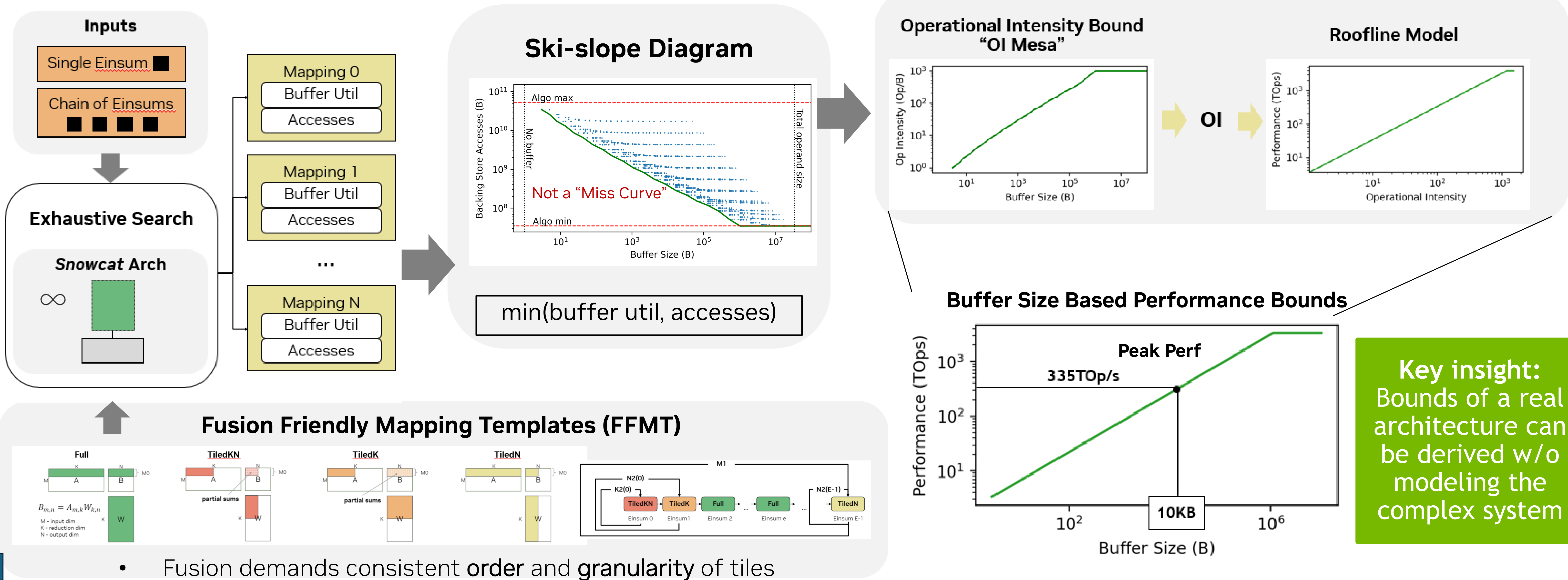
1 2

Data Movement vs. Buffer Size Gaps

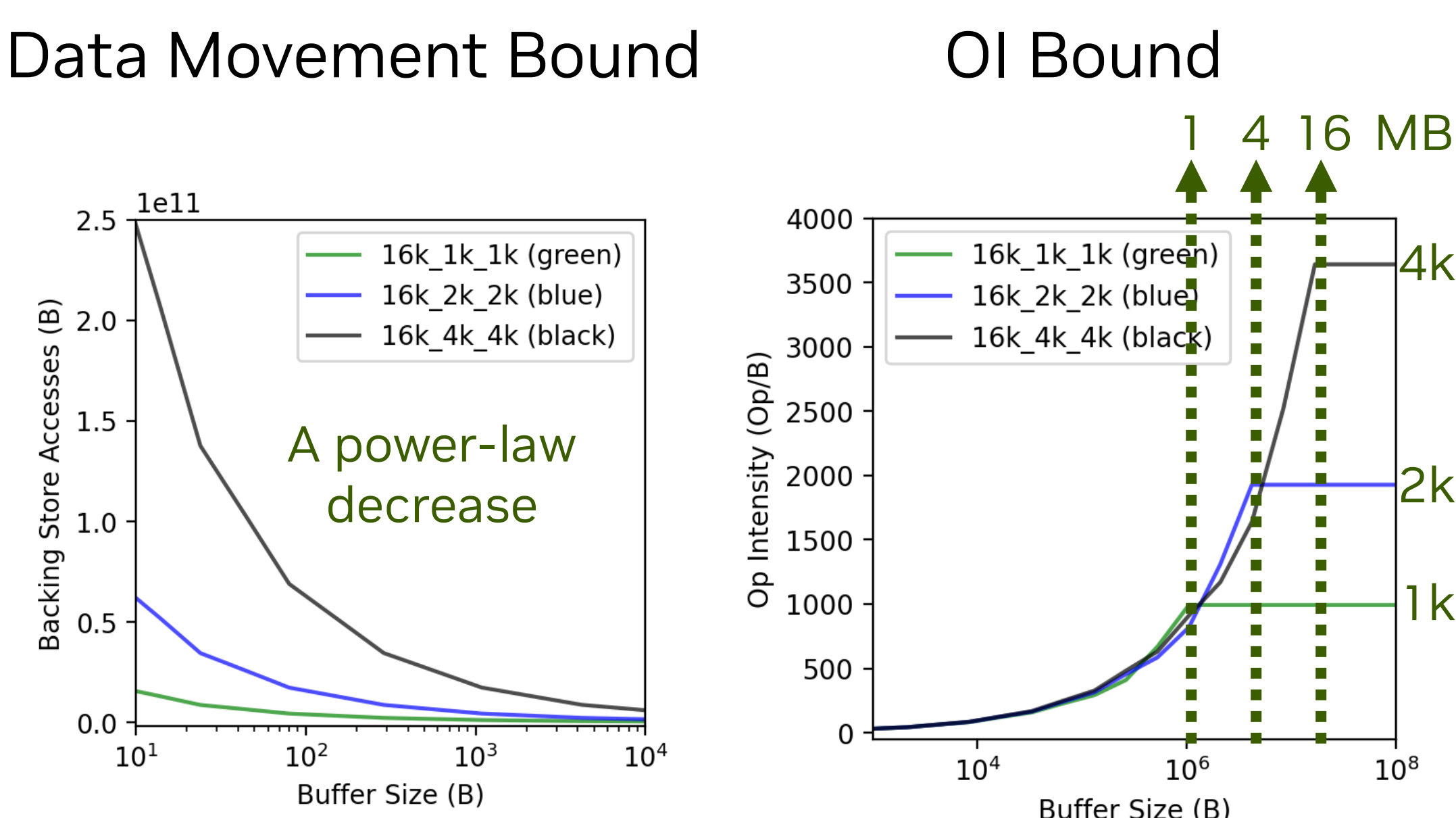


- [Gap 0] Given a buffer capacity, what is the minimal attainable *data access count*?
- [Gap 1] How much additional capacity is required to achieve full data reuse?
- [rate of change of Gap 0] How does an algorithm benefit from incremental increase in buffer capacity?

Our Orojenesis Flow for Early-Stage DSE



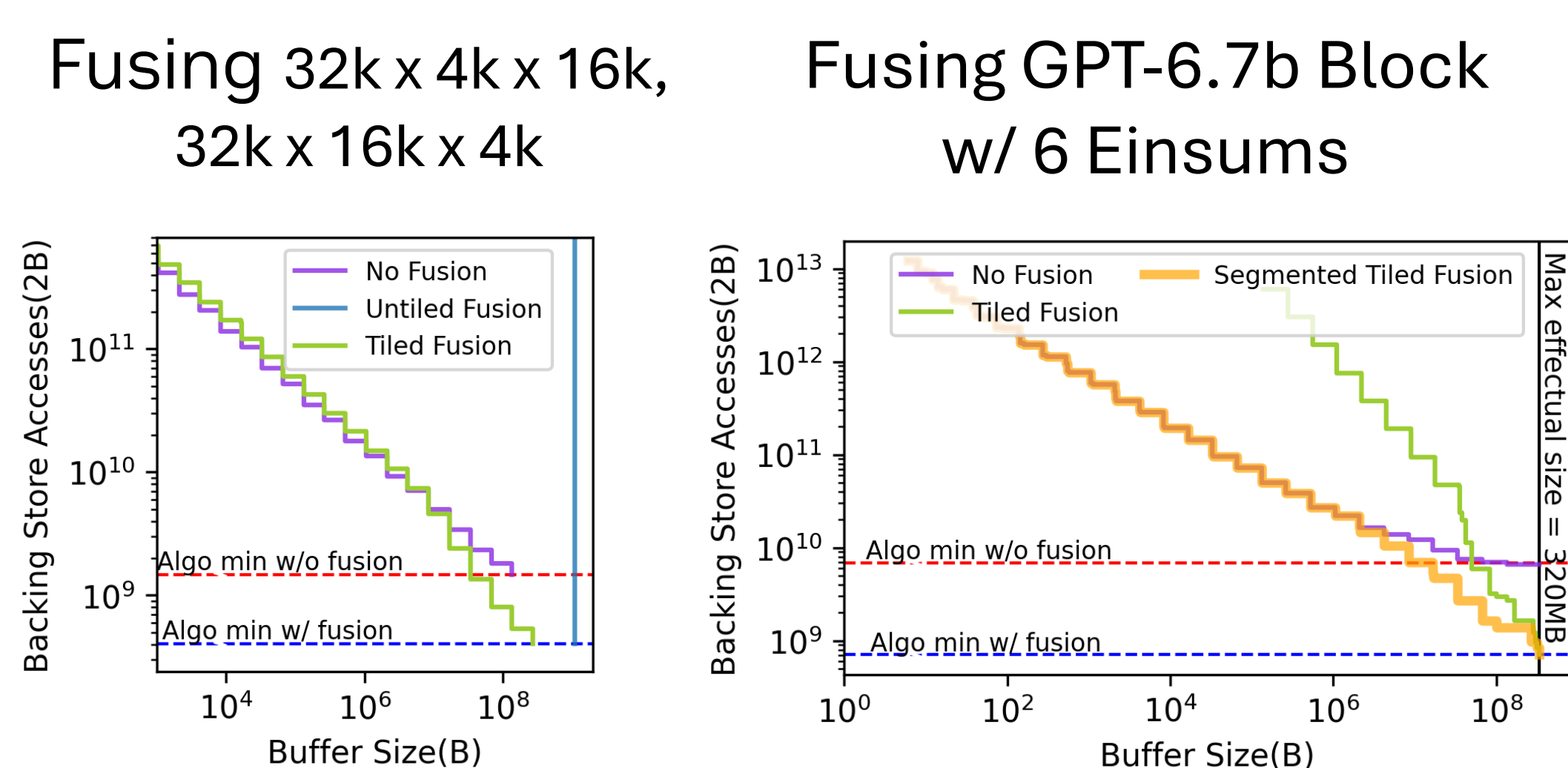
1. Provides valuable design insights



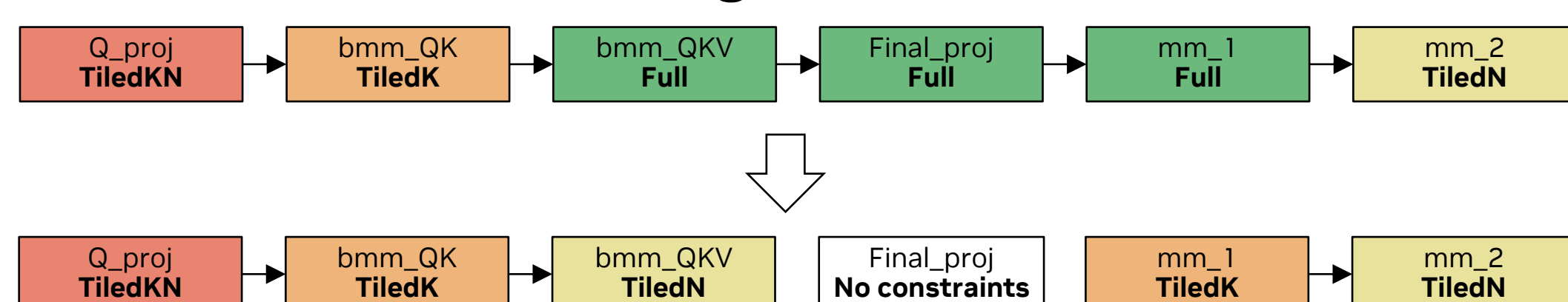
Observations :

1. The optimal OI of a GEMM is limited by its **smallest dimension**
2. The maximal effectual buffer size of a GEMM is approximately **the size of its smallest operand**

2. Comprehends complex fusion space



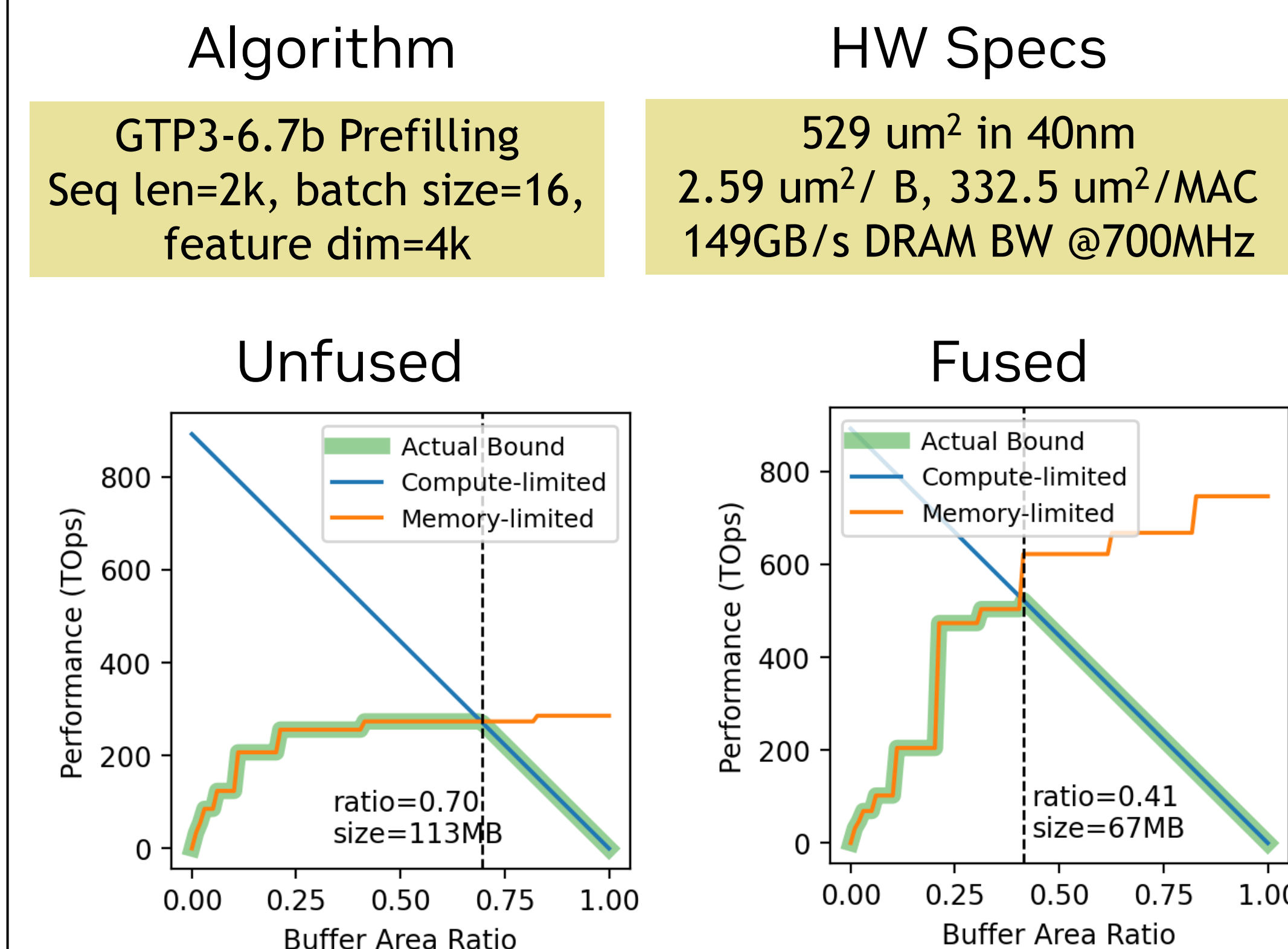
Chain Segmentation



Observations:

1. Bounds w/ fusion do not always outperform the unfused bounds
2. Fusing all layers can be suboptimal

3. Offers fast area tradeoff in early DSE



Observations:

1. Performance is a **concave** function wrt to the buffer area ratio
2. Fused workload demands **lower** buffer area but leads to better performance