

Heterogeneous Memory

DRAM Fast-tier	NUMA/CXL Slow-tier
Limited size Fast, low latency	Memory expansion Slow, High latency

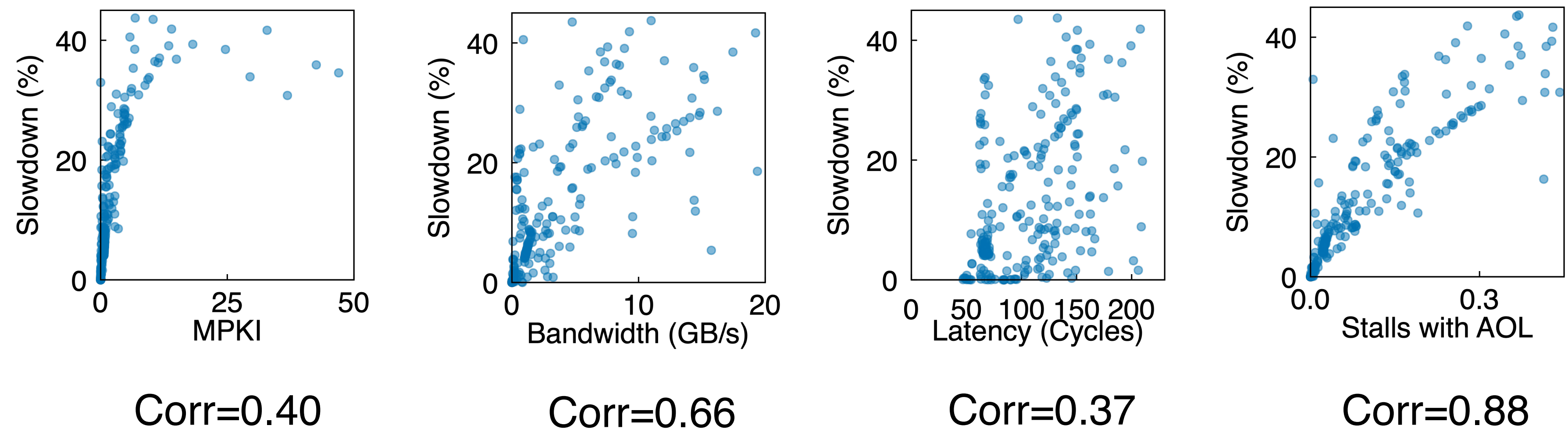
Without accurate predictors

Cloud operators overprovision DRAM to minimize risk
OS/runtimes correct placement errors after performance has degraded

Predictors are required

Accurate slowdown predictors can tell how a placement decision will affect application performance

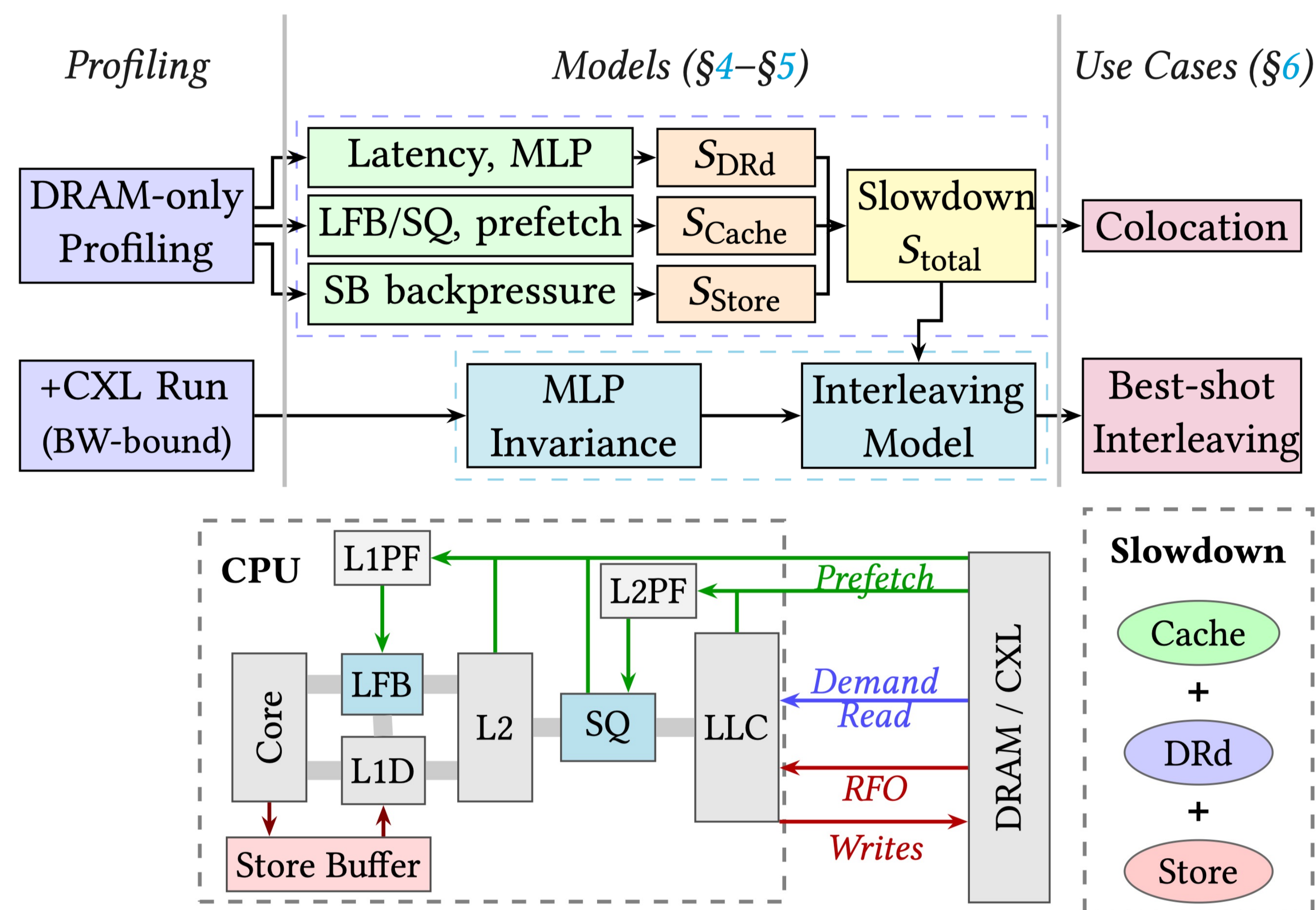
Existing Metrics Fall Short



Simple proxy metrics fail to accurately capture microarchitectural behavior

Is it possible to predict the workload performance on CXL using intrinsic workload signatures?

CAMP: Causal Analytical Memory Prediction



Slowdown Prediction

The microarchitectural causes:

The dynamics of MLP and latency

The reliance on LFB and prefetching

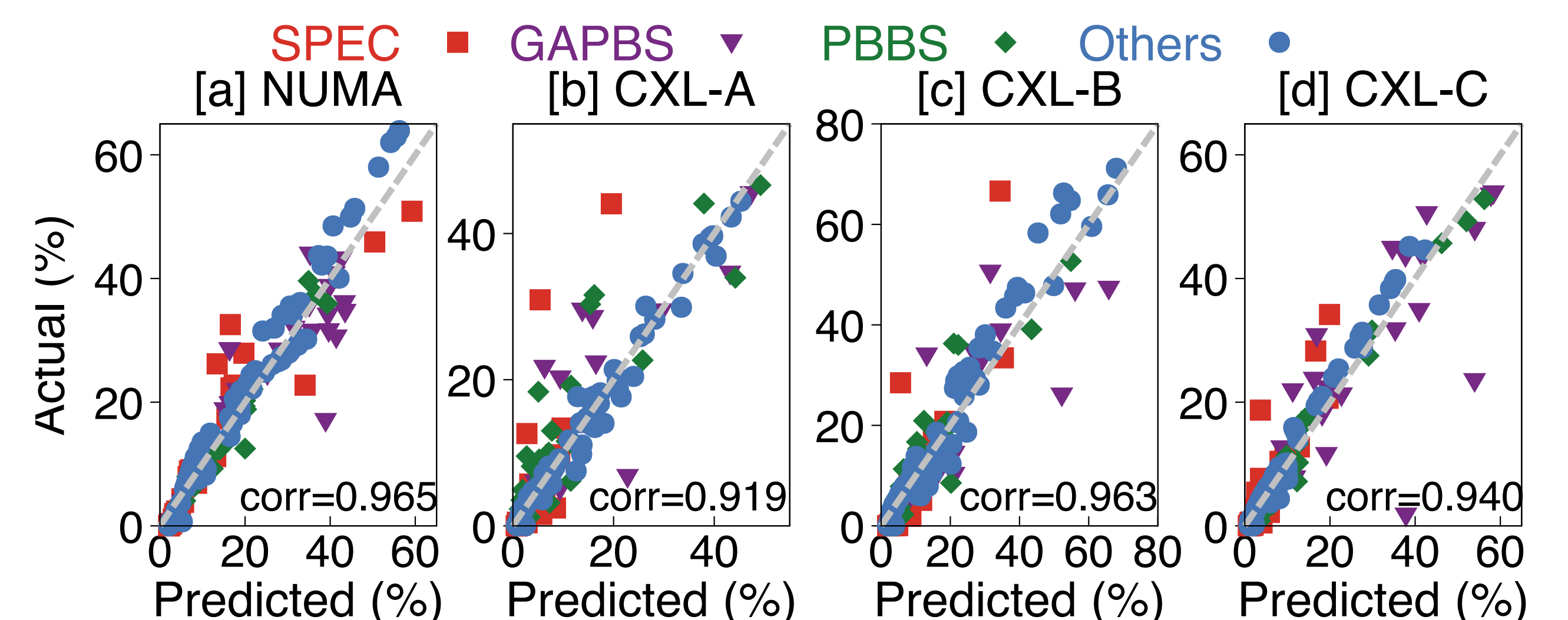
The pressure on Store Buffer

The predictors:

$$S_{DRd} \approx k \times \frac{1}{p+q} \times \frac{MLP_{DRAM}}{L_{DRAM}} \times \frac{s_{LLC}}{c}$$

$$S_{Cache} \approx k \times R_{LFB-hit} \times R_{Mem} \times \frac{s_{Cache}}{c}$$

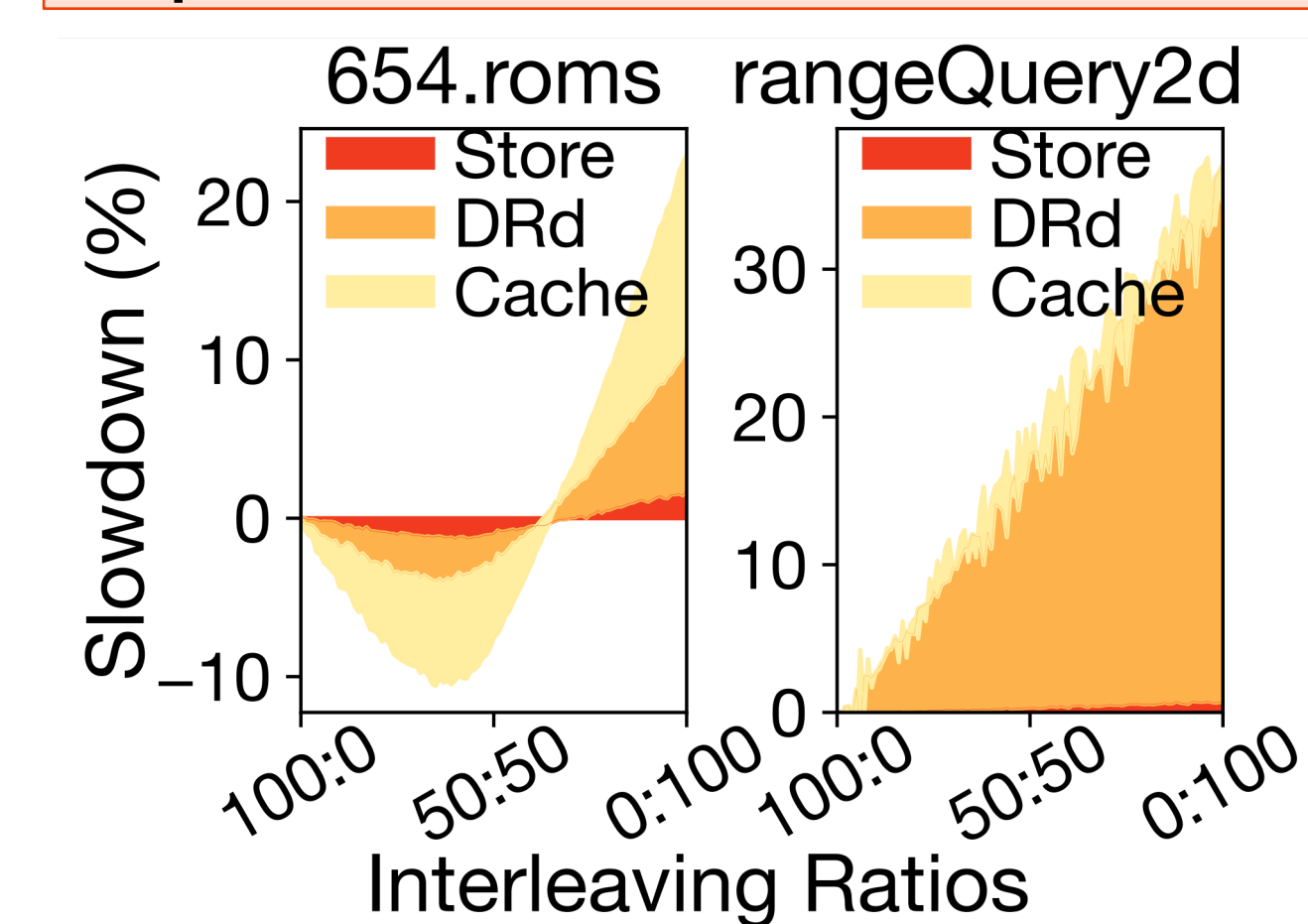
$$S_{Store} \approx k \times \frac{s_{SB}}{c}$$



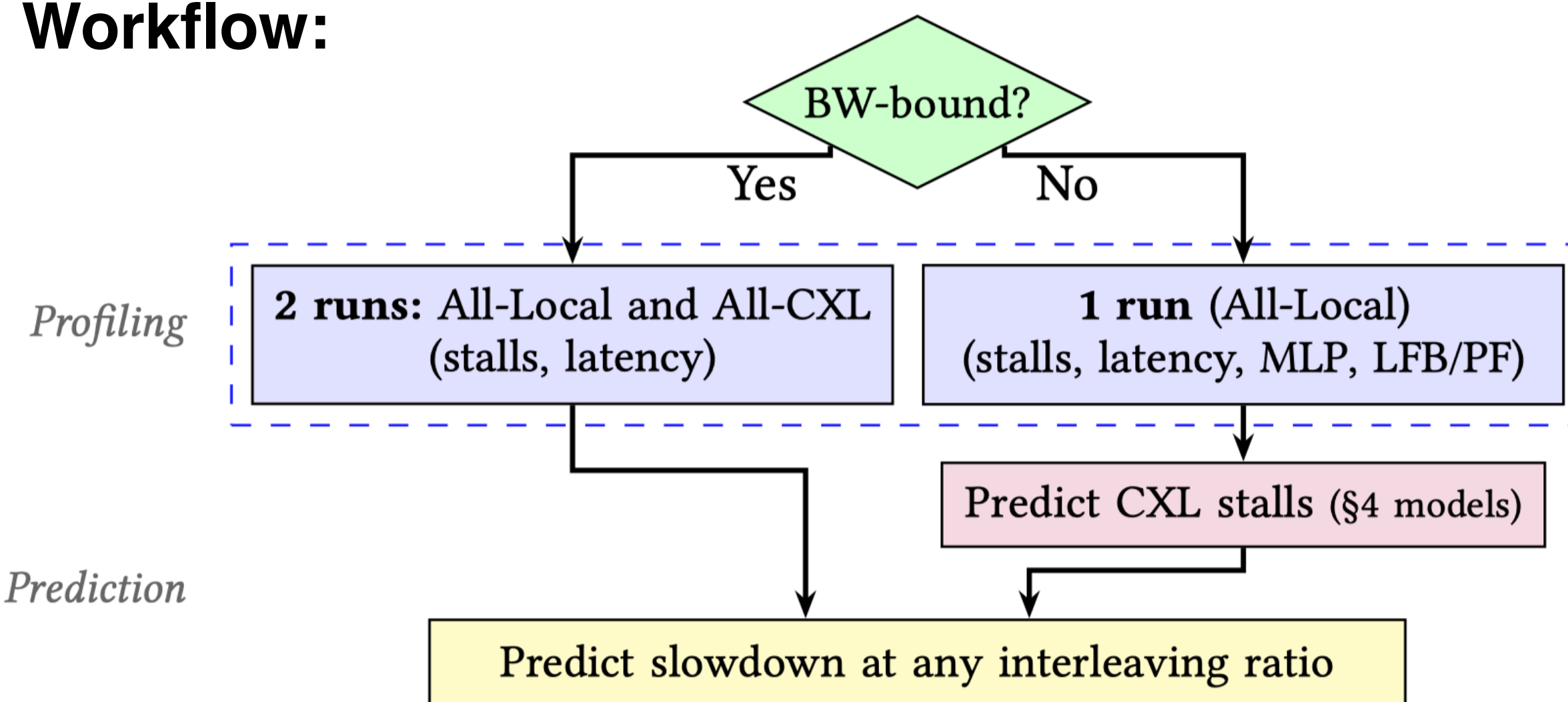
Interleaving model

Use Cases

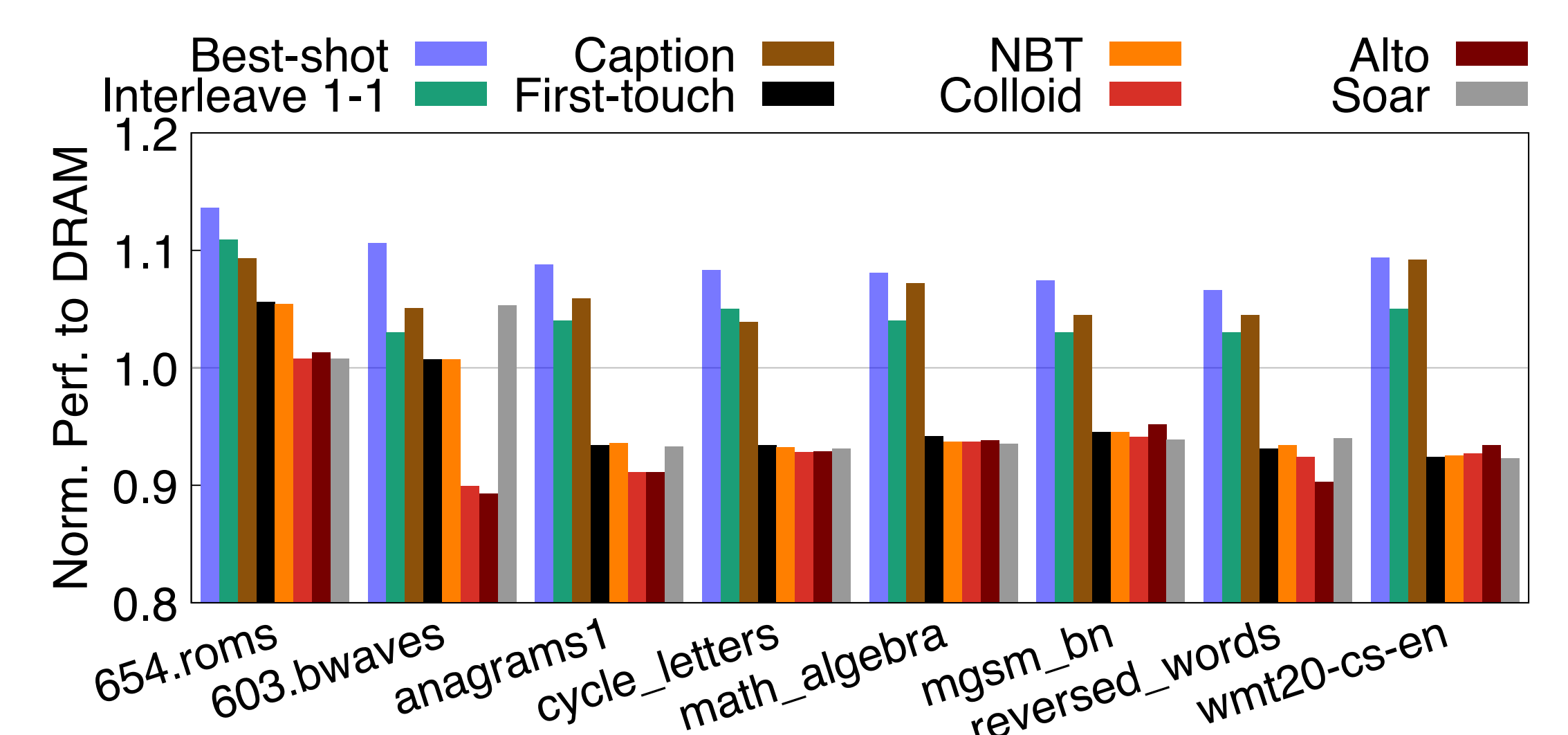
How is the continuous spectrum of performance under interleaving?



Workflow:

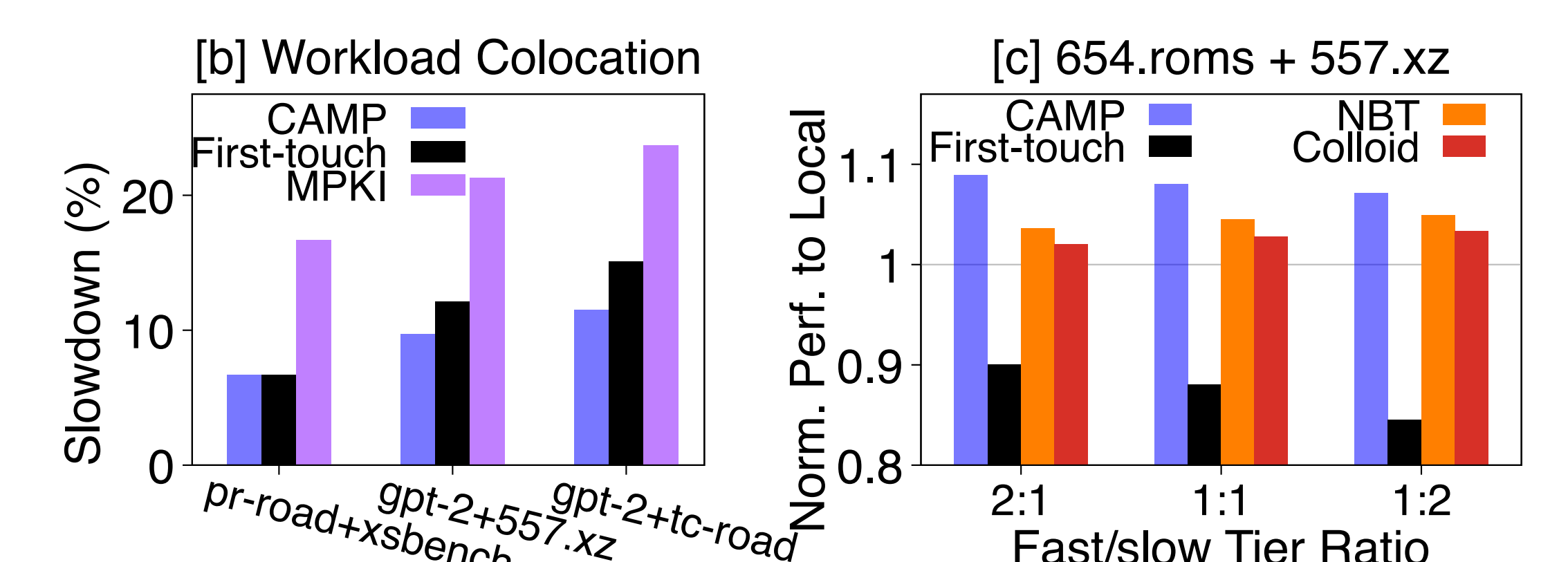


"Best-shot" Interleaving:

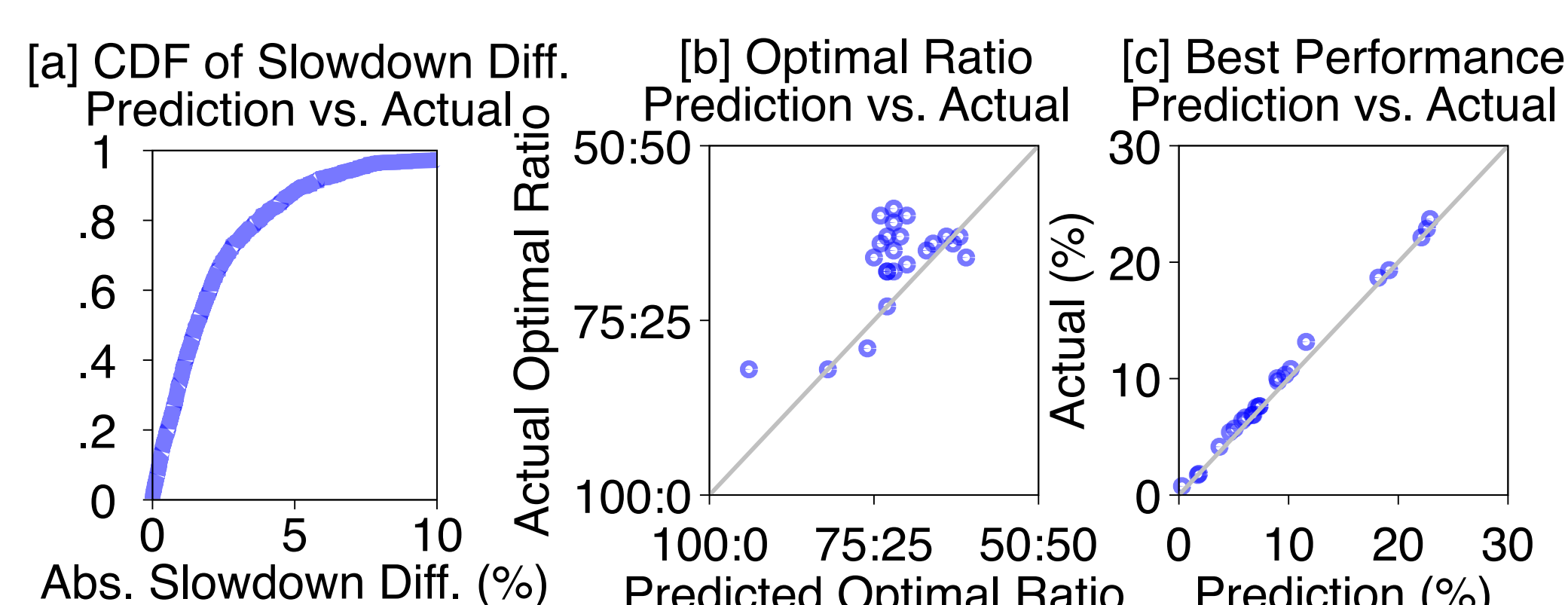
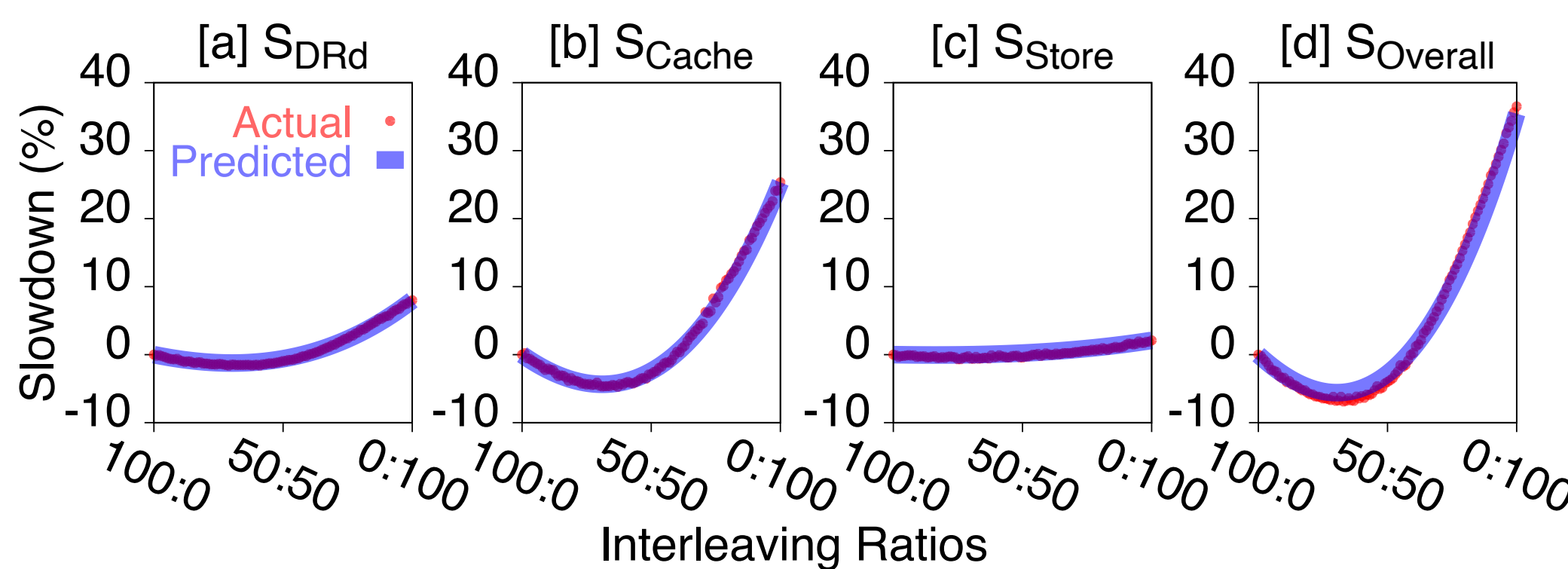


Workload Colocation:

	557.xz	pr-road	tc-road	gpt-2	xsbench
Prediction	25%	21%	25%	11%	7%
Actual	37%	17%	31%	10%	6%



Prediction results:



Goal:

Predicting performance at any interleaving ratio (x)

Challenge:

Changing the ratio alters the traffic load on each tier, which non-linearly impacts contention and latency

$$C = \frac{N \times L}{MLP}$$

The variance of C is determined by how N, L and MLP change with x

-> MLP invariance

-> Modeling latency curve

The predictor:

$$S(x) \approx \frac{M(x) \cdot s_{DRAM} + M(1-x) \cdot s_{CXL} - s_{DRAM}}{c}$$

$$M(x') = \frac{x' \cdot [L_{idle} + (L_{full} - L_{idle}) \cdot x'^2]}{L_{full}}$$

More in the paper:

