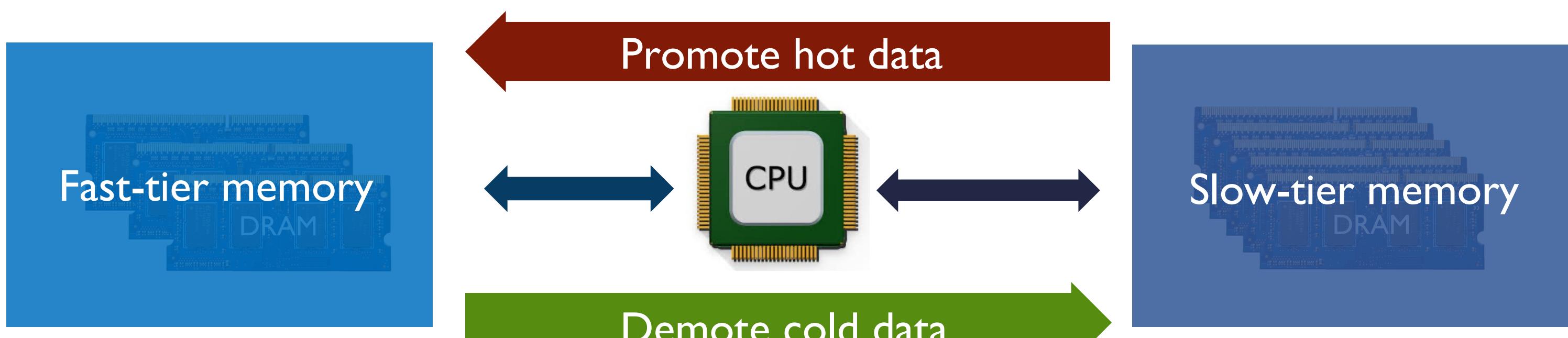


Memory Tiering



Traditional Approach:

1. Common assumption: $Performance = Hotness$
2. Page migration: promote hot data to the fast-tier
3. First-touch allocation

Research Questions:

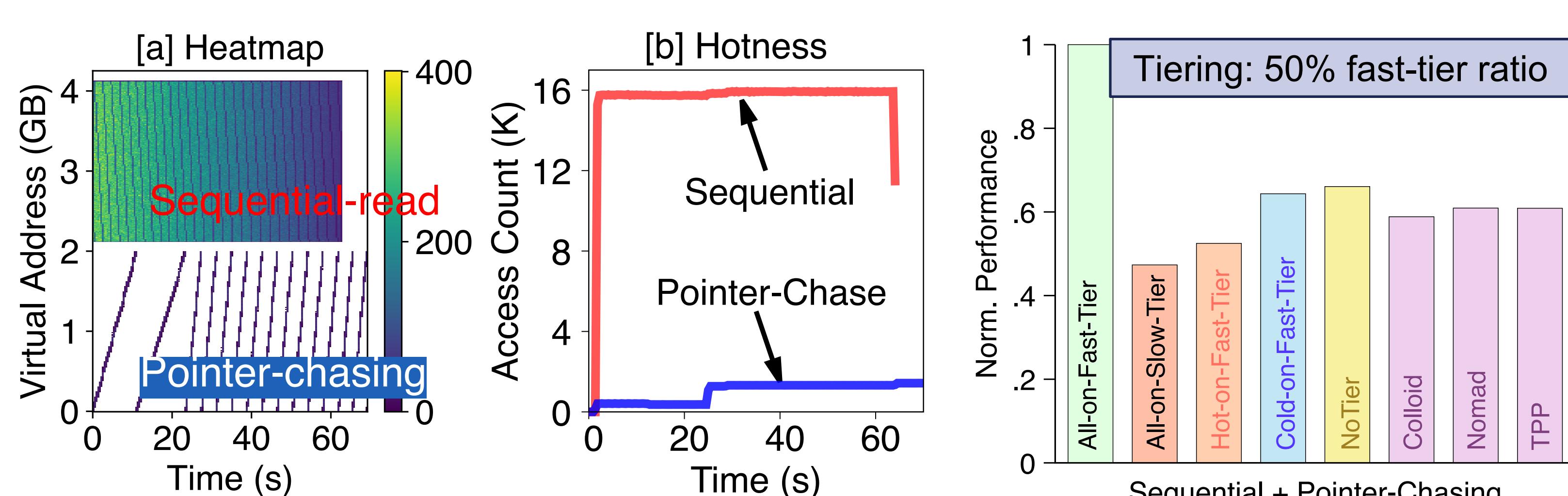
- Why cannot hotness represent performance?
- Which metrics should be used to guide tiering?
- How to apply the new metrics on memory allocation and migration?

Overview

Memory allocation/migration for tiered memory *beyond hotness*

1. AOL: Amortized Offcore Latency
Key factor for accurate performance prediction
Quantifies the impact of memory accesses on performance slowdown
2. SOAR: Static Object Allocation based on Ranking
Near-optimal object placement after profiling based on AOL
3. ALTO: AOL-based Layered Tiering Orchestration
Regulate page migration for hot but less performance-critical pages
4. Evaluation
Compare SOAR/ALTO with TPP/Nomad/NUMA-Balancing/Colloid
Evaluated under memory contention on the fast-tier

Hotness != Performance

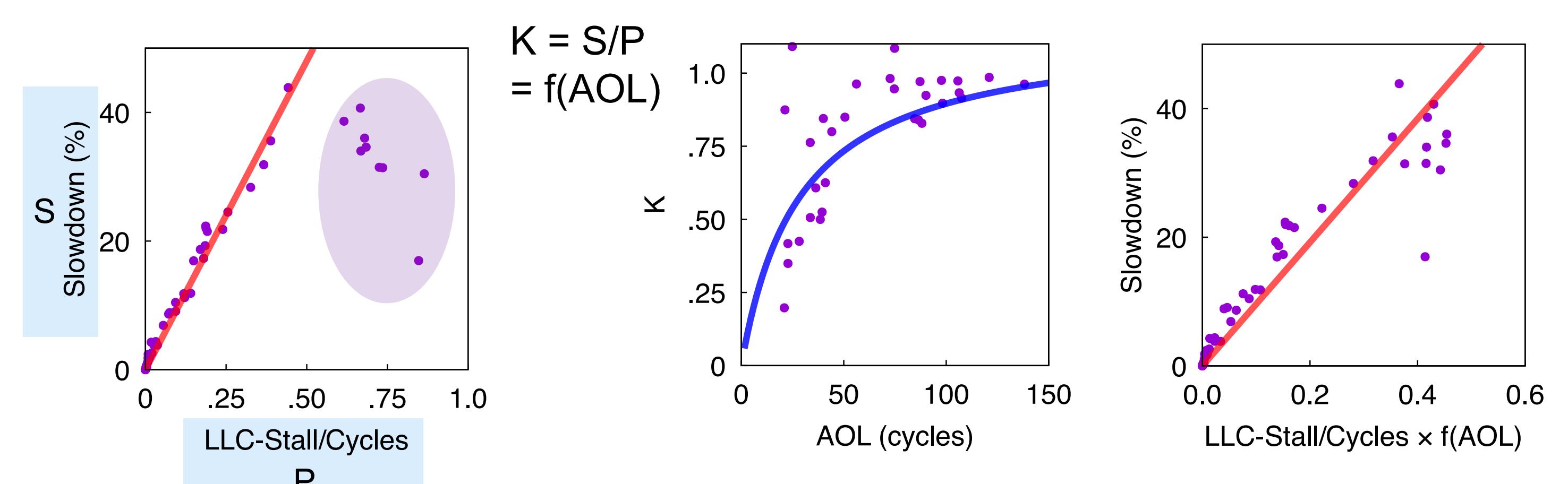


Hot-on-fast-tier gives worse (34%) performance than Cold-on-fast-tier !

Reasoning:

1. LLC stalled cycles
2. Memory Level Parallelism (MLP) can mask latency penalty

Amortized Offcore Latency



$AOL = Latency / MLP$

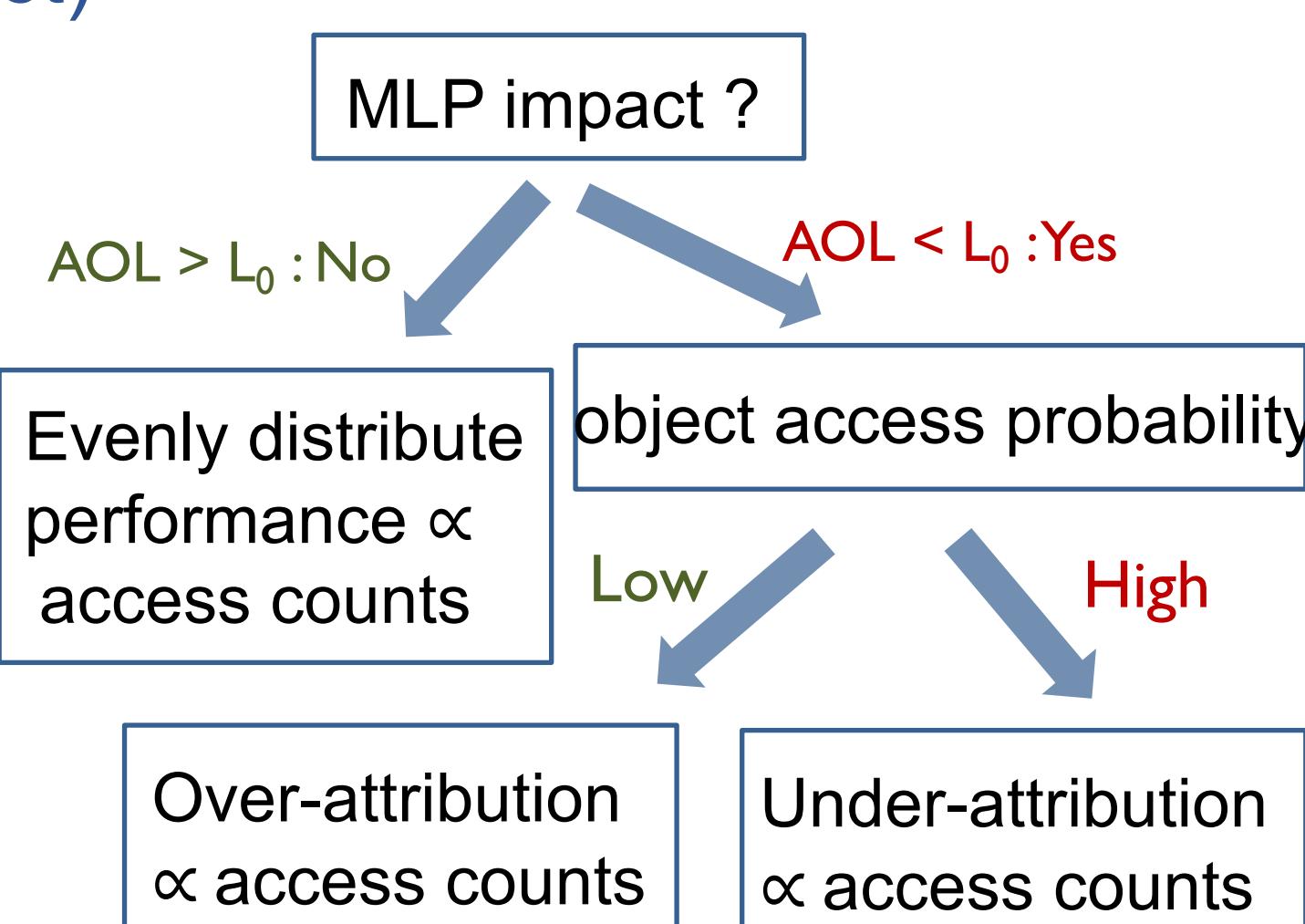
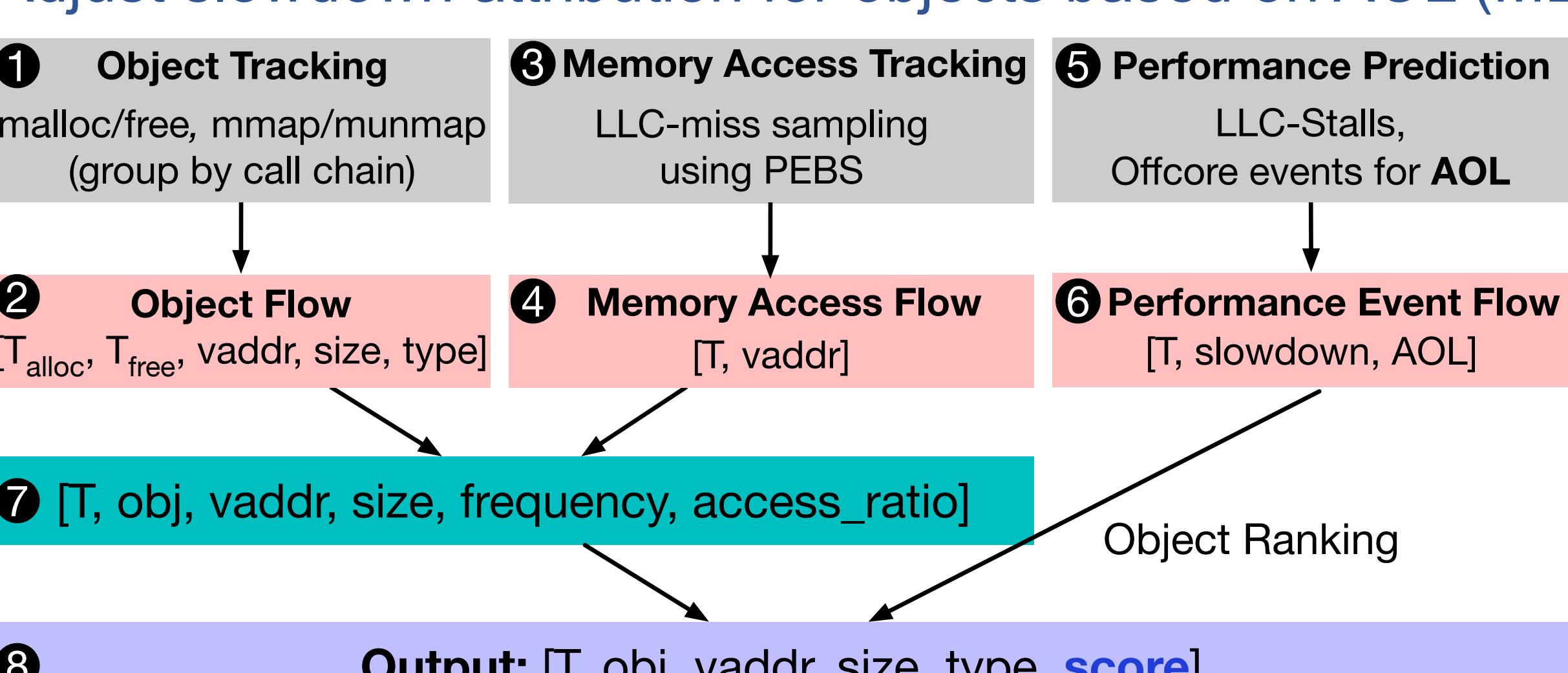
MLP itself is not enough to model how performance slowdown is impacted by parallel accesses

AOL quantifies how increased LLC stalls on slow-tier are amortized by MLP

SOAR: Static Object Allocation based on Ranking

Assign predicted performance to each object each time-period

Adjust slowdown attribution for objects based on AOL (MLP effect)



ALTO: AOL-based Layered Tiering Orchestration

Unnecessary page migration

1. Promote hot but non-performance-critical pages results in minimal or negative performance gain
2. Overheads caused by ineffective page migration

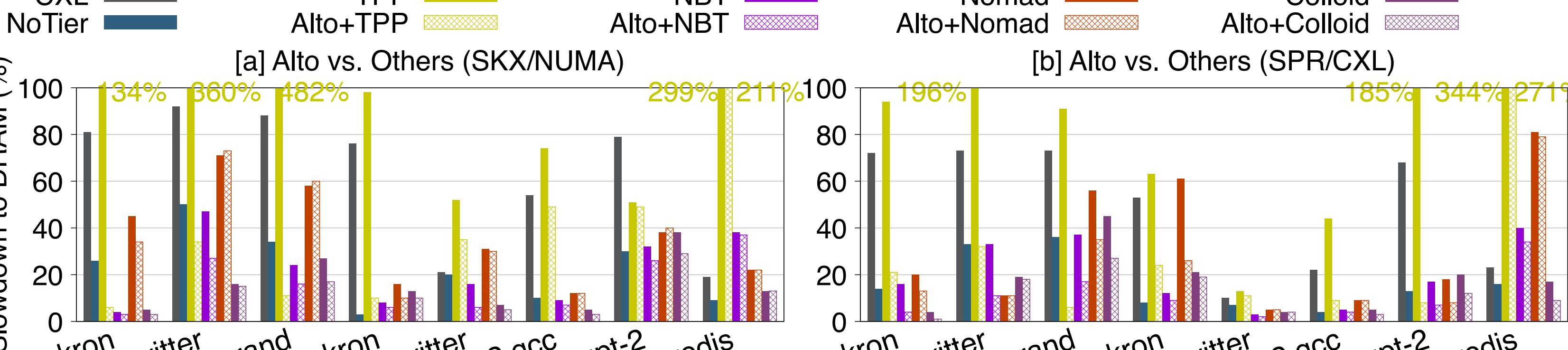
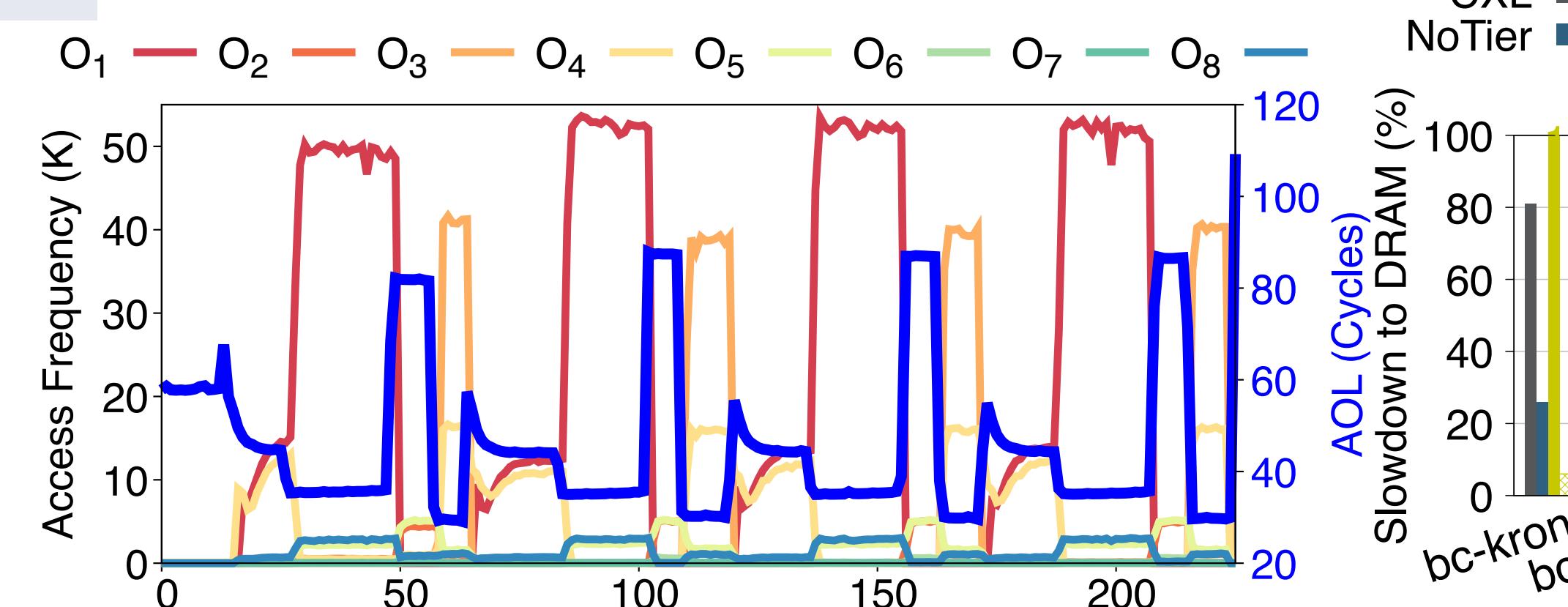
For each time period t_i, t_{i+1} :

> AOL _{high}	Enable page promotion
[AOL _{low} , AOL _{high}]	Partial page promotion Scale < f(AOL)
< AOL _{low}	Disable page promotion

Evaluation

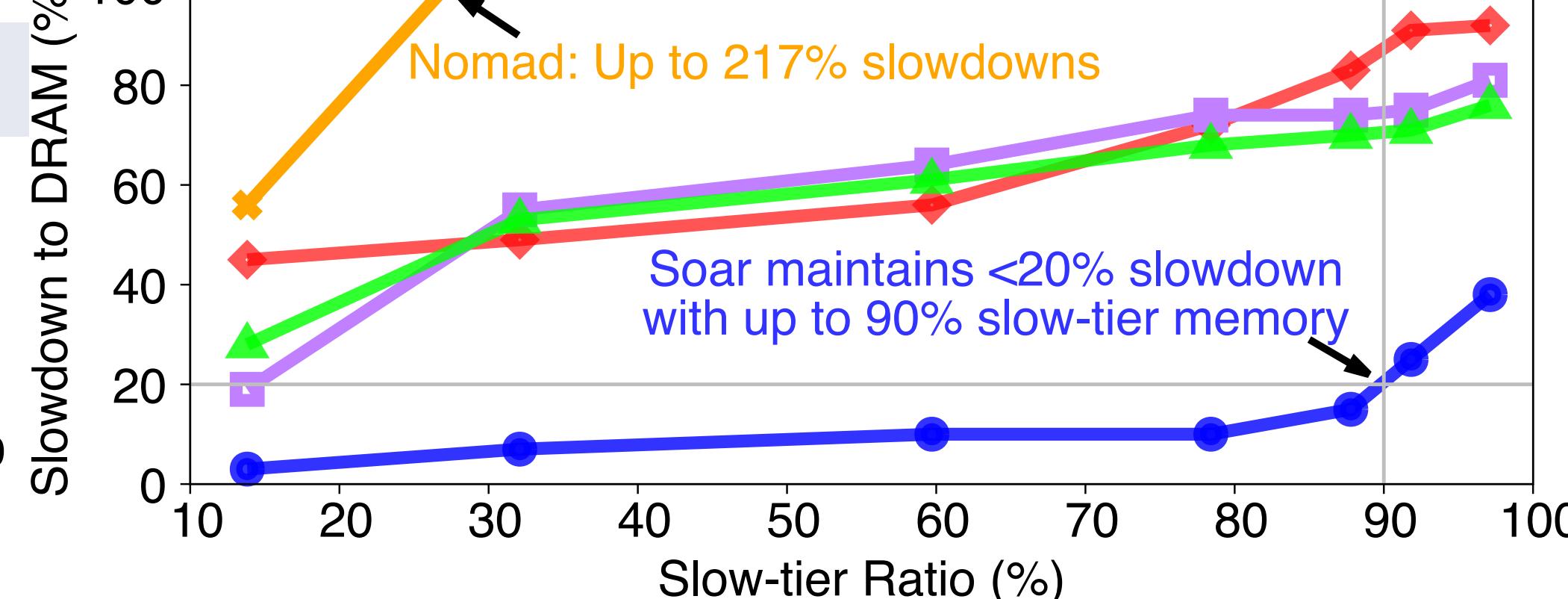
Object Ranking (bc-urand)

Object	Size	First-touch	Freq	SOAR
O ₁	536MB	8	1	1
O ₂	536MB	6	6	2
O ₃	536MB	10	2	3
O ₄	1073MB	5	3	4
O ₅	1073MB	3	4	5
O ₆	536MB	4	7	6
O ₇	536MB	7	8	7
O ₈	17GB	2	5	8
O ₉	327KB	9	9	9
O ₁₀	1073MB	1	10	10



Object placement (50% fast-tier ratio)

	Fast-tier	Slow-tier
First-touch	O ₁₀ , O ₈	O ₁ -O ₇ , O ₉
Freq	O ₁ , O ₃ -O ₅ , O ₈	O ₂ , O ₆ -O ₈ , O ₉ -O ₁₀
SOAR	O ₁ -O ₈	O ₈ -O ₁₀



ALTO/SOAR under memory bandwidth contention (bc-urand)

2%	6%	81%	70%	186%	105%	30%	33%
21%	28%	83%	52%	65%	38%	44%	35%
22%	55%	101%	78%	103%	70%	61%	50%
14%	55%	86%	60%	95%	65%	61%	58%



<https://github.com/MoatLab/SoarAlto>