

CACHE OPTIMIZATION FOR THE MODERN WEB

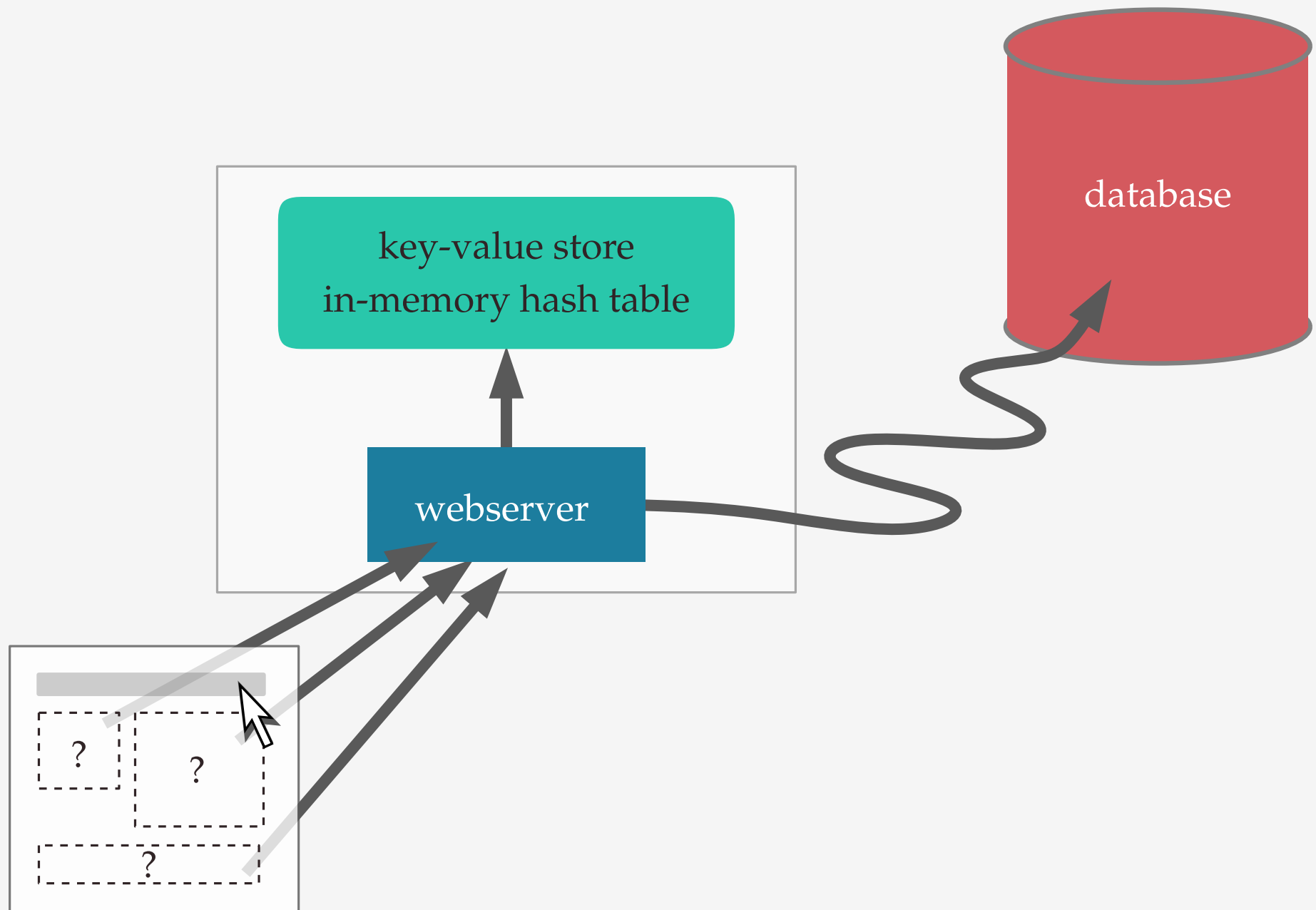
Jenny Lam

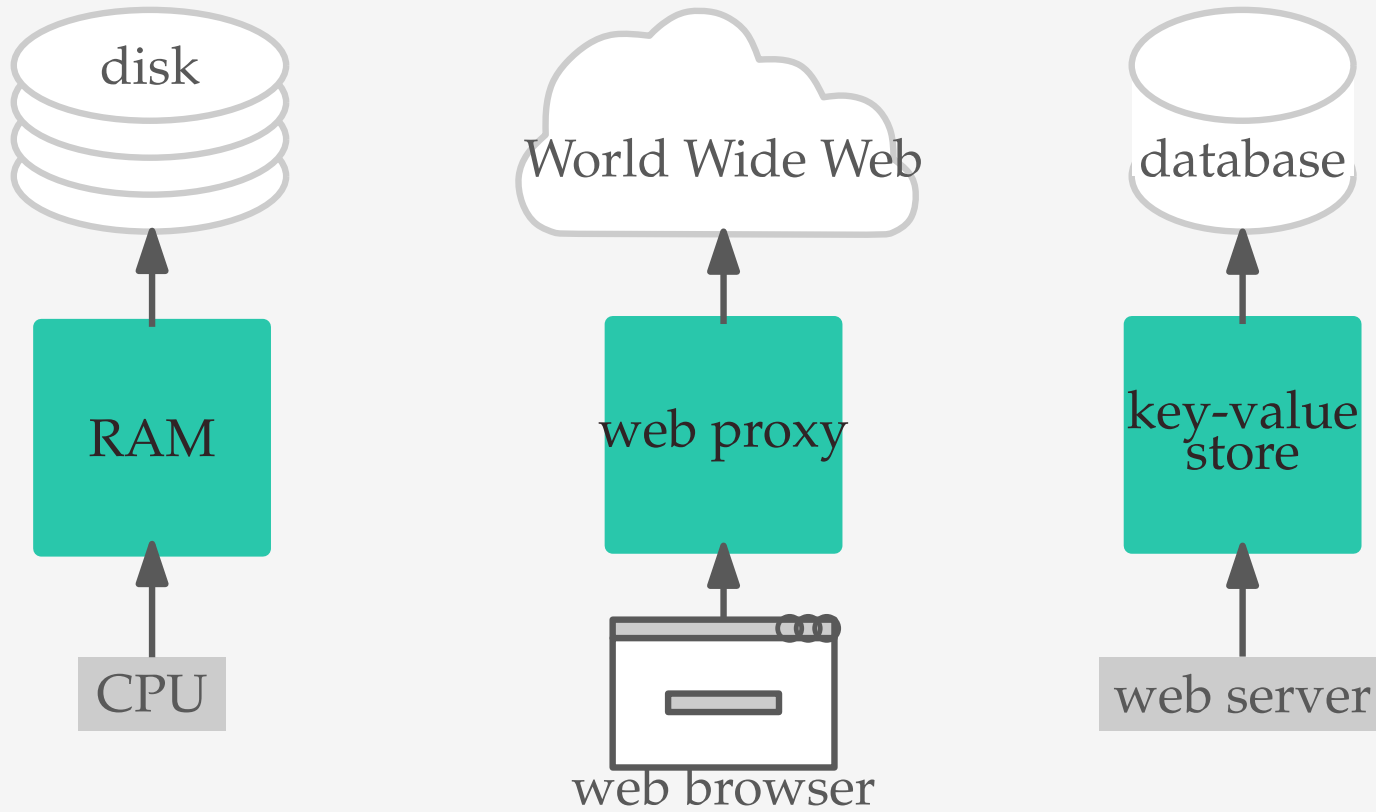
Sandy Irani (chair)

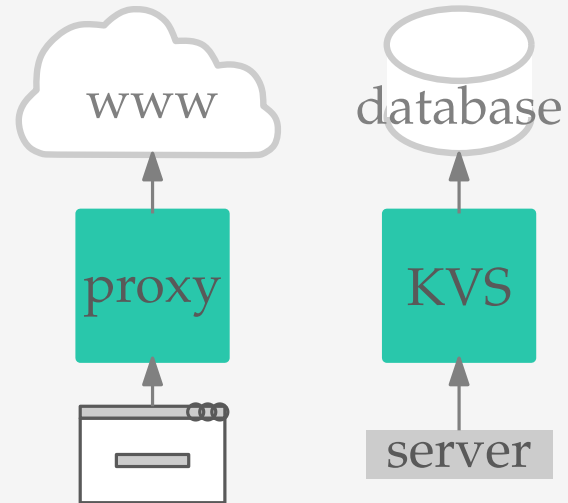
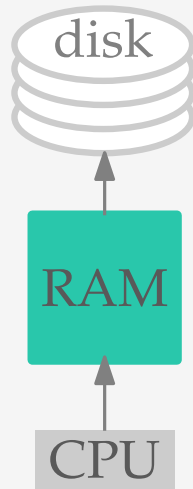
Michael Dillencourt

Michael T. Goodrich

11/24/2015

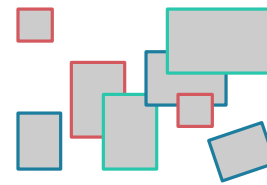






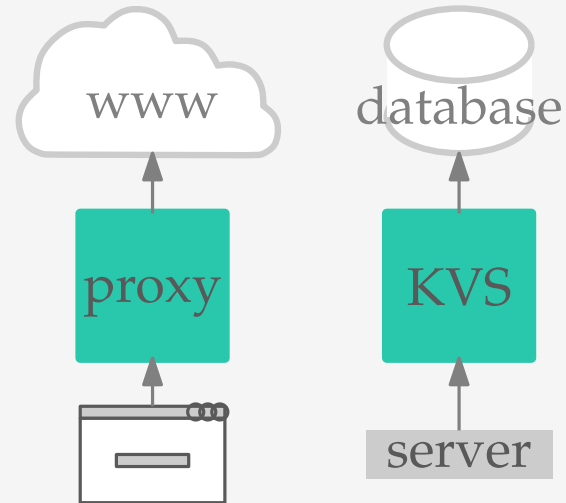
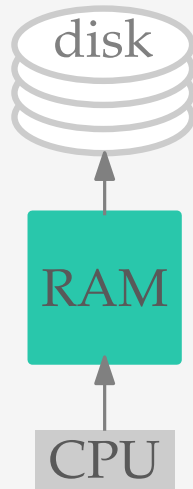
PAGING

minimize
number of cache misses



GENERALIZED
CACHING

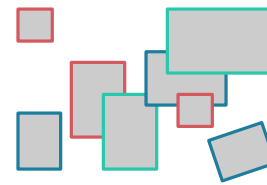
minimize
total cost of cache misses



PAGING

minimize
number of cache misses

Least Recently Used (LRU)

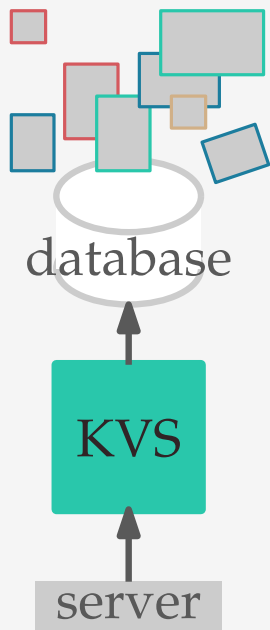


GENERALIZED
CACHING

minimize
total cost of cache misses

GreedyDual-Size (GDS)

GDS → CAMP

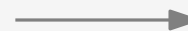


generalized
caching



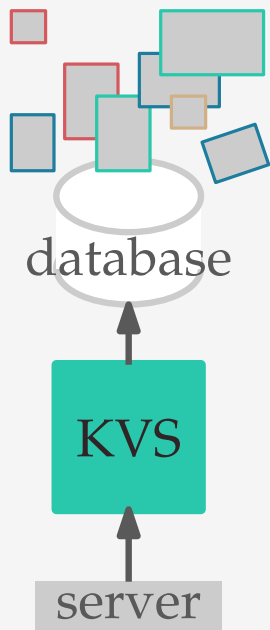
managed memory
caching

2-level cache



multi-level cache

GDS → CAMP



generalized
caching



managed memory
caching

2-level cache



multi-level cache

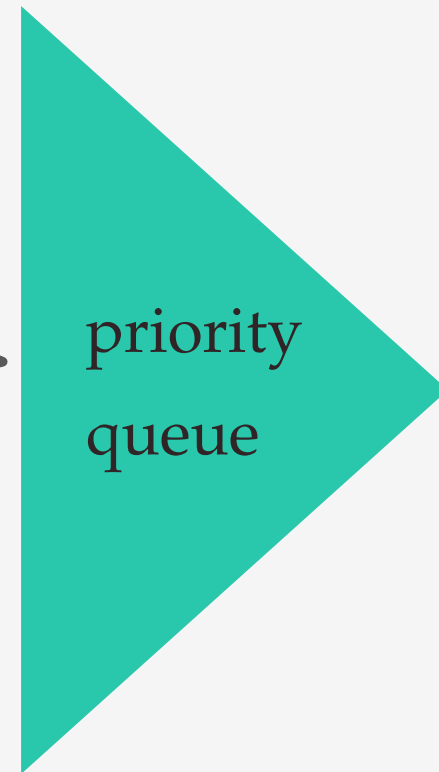


GDS

p

GDS priority

$$\frac{\text{cost}(p)}{\text{size}(p)} + \text{lowest priority}$$



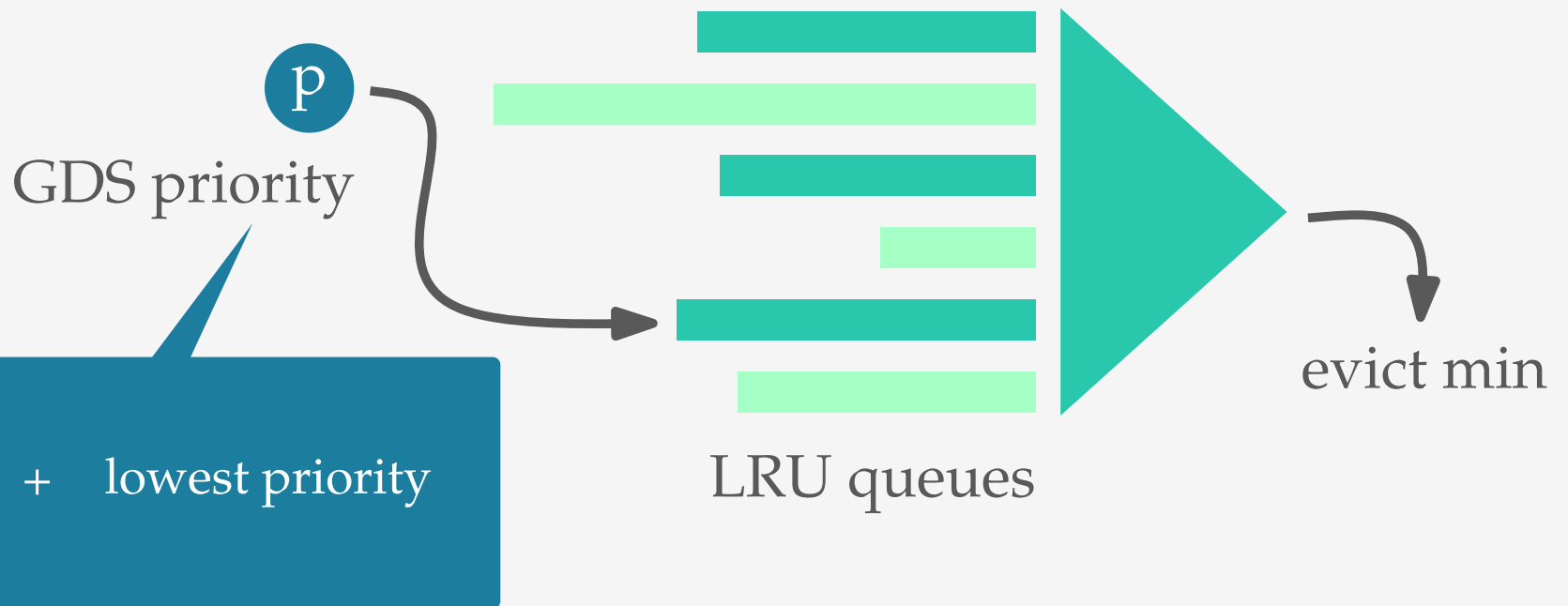
priority
queue



evict min



CAMP





CAMP

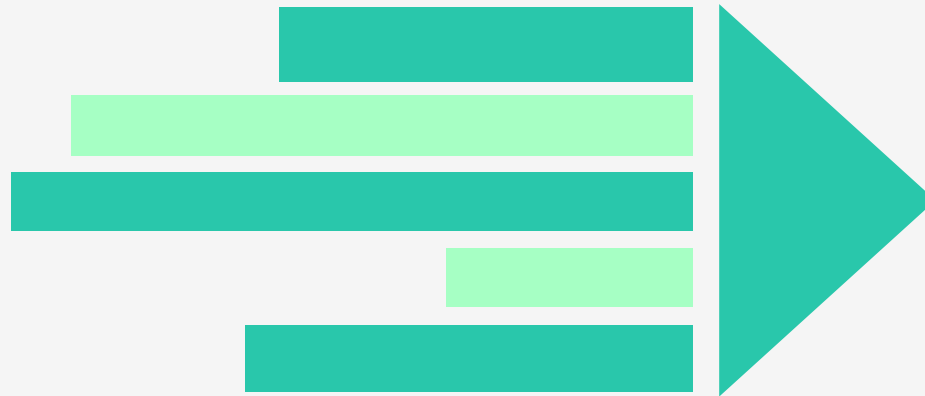
$$\text{round} \left(\frac{\text{cost}(p)}{\text{size}(p)} \right)$$





CAMP

$$\text{round} \left(\frac{\text{cost}(p)}{\text{size}(p)} \right)$$





CAMP

$$\text{round} \left(\frac{\text{cost}(p)}{\text{size}(p)} \right)$$





PERFORMANCE

log (#items) per update



$$\text{cost}(\text{GDS}) \leq k \text{ cost}(\text{OPT})$$

log (#queues) per update



$$\text{cost}(\text{CAMP}) \leq (1 + \varepsilon)k \text{ cost}(\text{OPT})$$

approximation
parameter

EXPERIMENTS



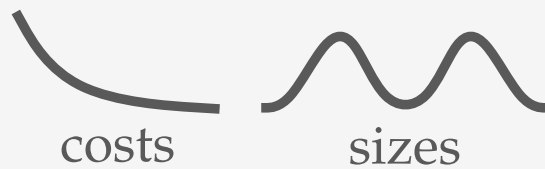
trace generated by BG, a social networking benchmark

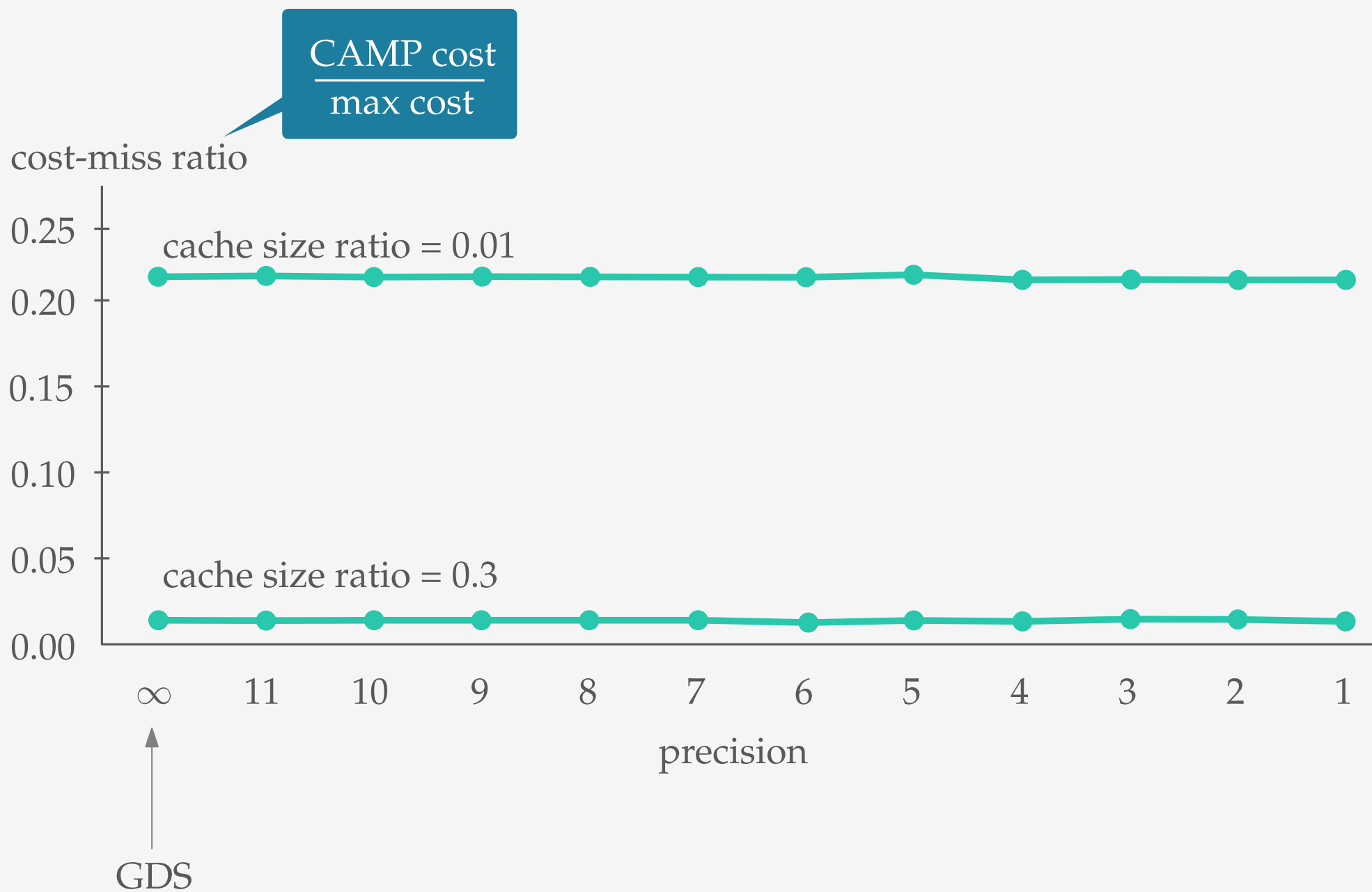
4 million requests

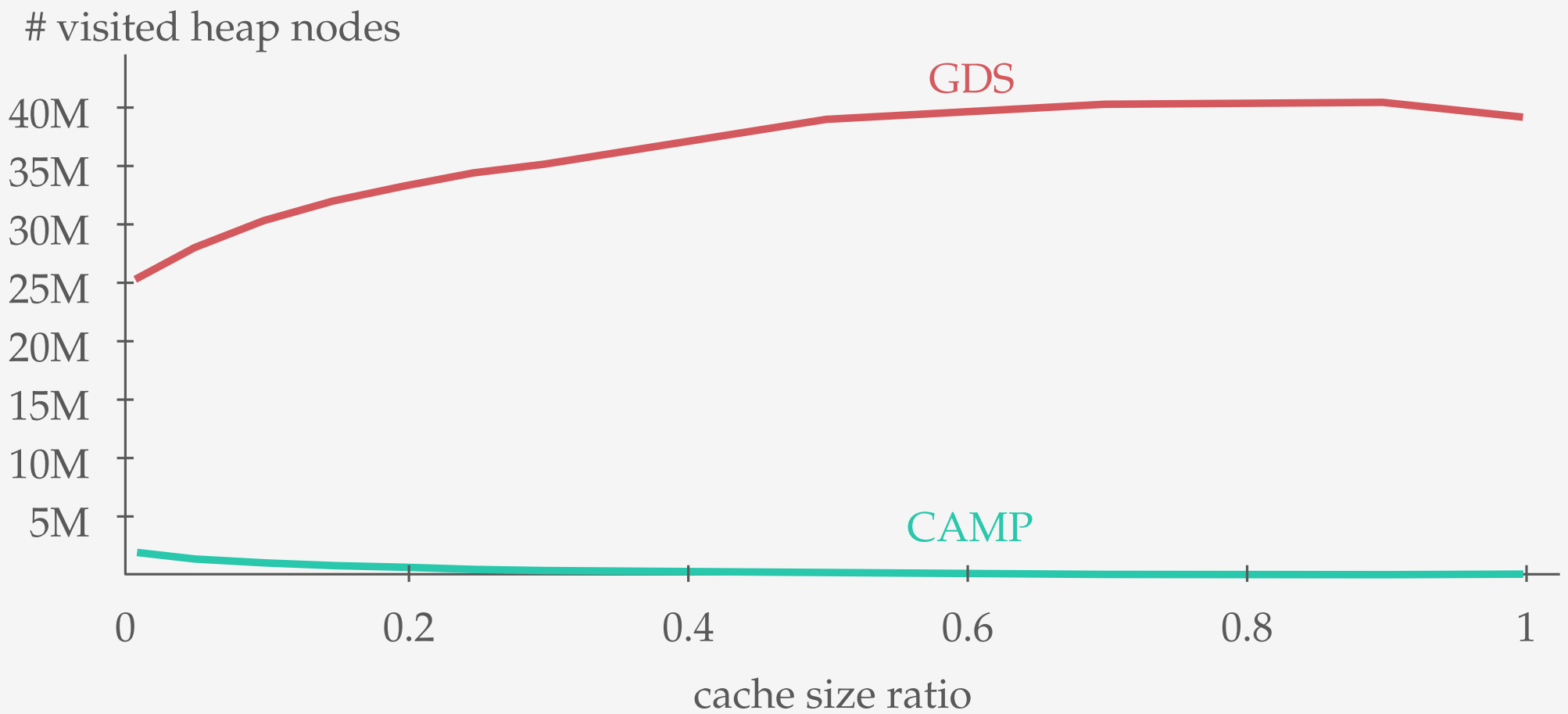
i.i.d. with 70% of requests to 20% of items



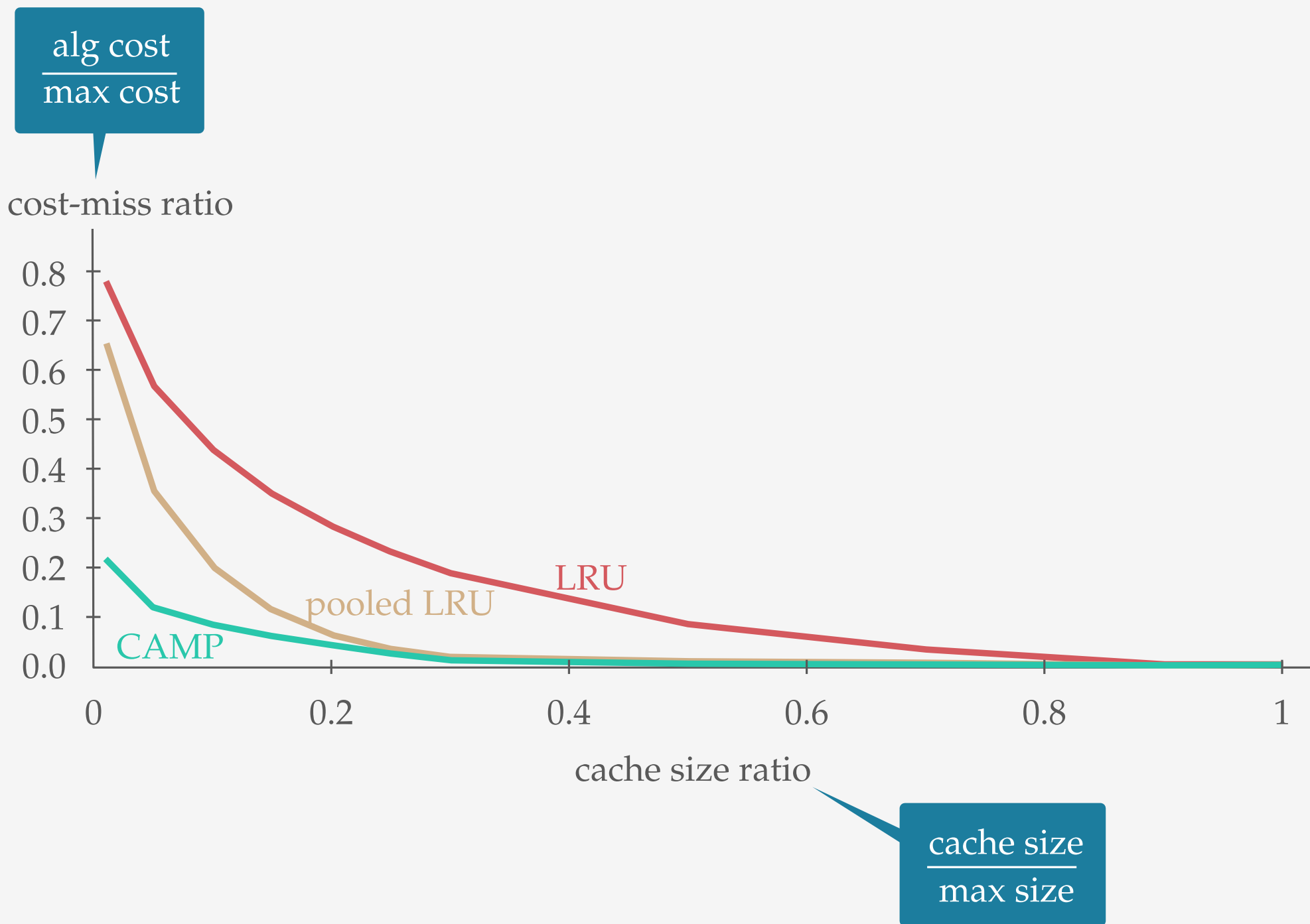
20,000 items



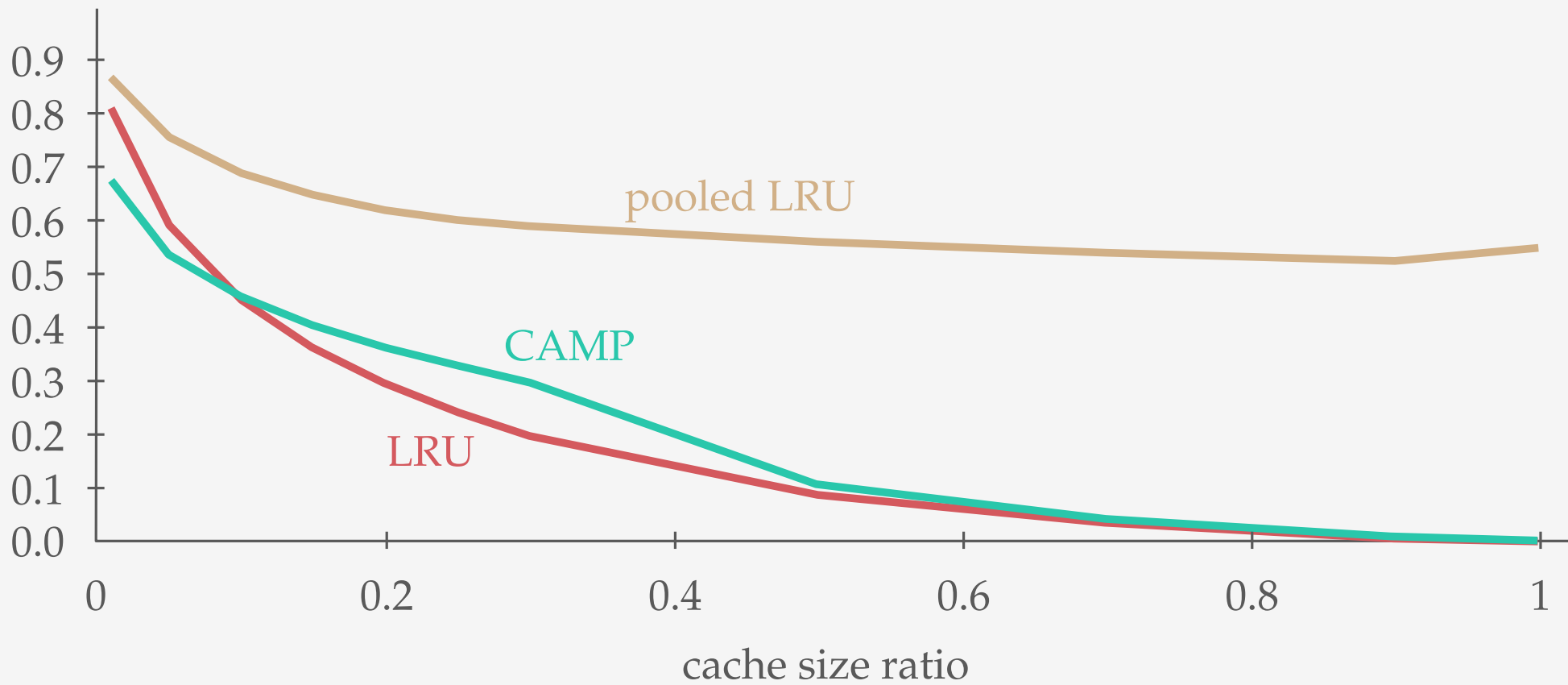




$$\frac{\text{cache size}}{\text{max size}}$$

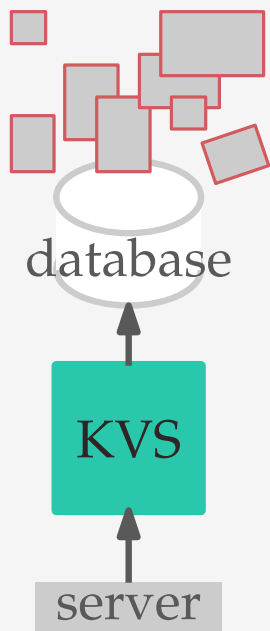


miss rate

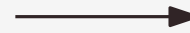


$\frac{\text{cache size}}{\text{max size}}$

GDS → CAMP



generalized
caching



managed memory
caching

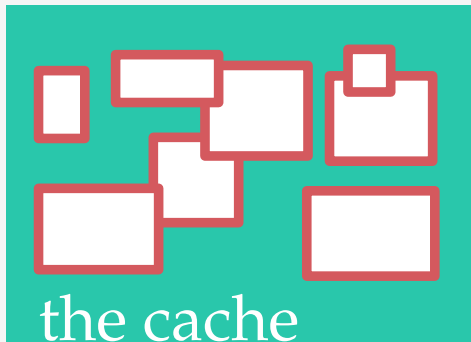
2-level cache



multi-level cache

THE GENERALIZED CACHING PROBLEM

variable size and cost



GOAL

minimize **total cost** of cache misses

SUBJECT TO

total size of items in cache
cannot exceed the cache size

THE MANAGED MEMORY CACHING PROBLEM

variable size and cost



every item must fit in a contiguous
segment of memory

CACHE REPLACEMENT
MEMORY ALLOCATION



CAMP-MALLOC



LRU queues



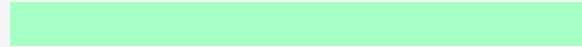
CAMP-MALLOC



FIFO queues



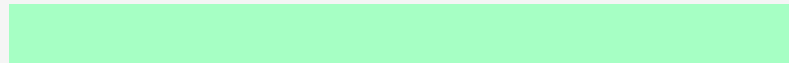
CAMP-MALLOC



FIFO queue



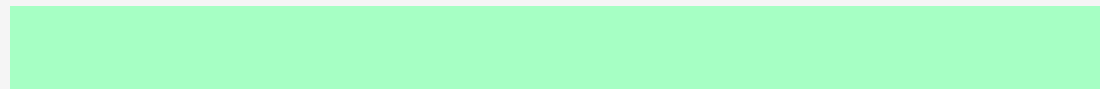
CAMP-MALLOC



FIFO queue



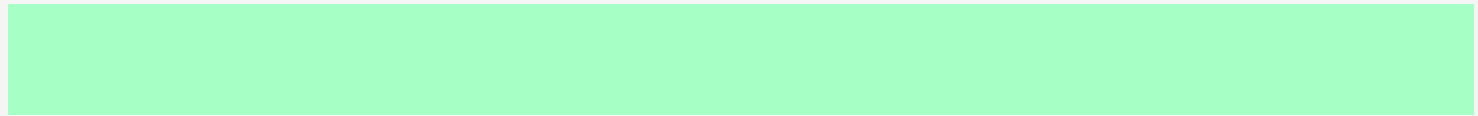
CAMP-MALLOC



FIFO queue



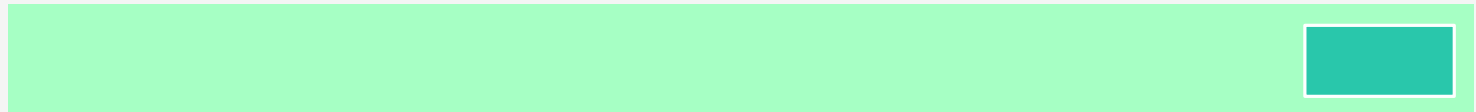
CAMP-MALLOC



FIFO queue



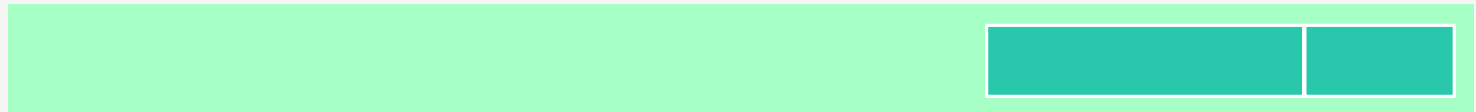
CAMP-MALLOC



FIFO queue



CAMP-MALLOC



FIFO queue



CAMP-MALLOC



FIFO queue



CAMP-MALLOC



FIFO queue



CAMP-MALLOC



FIFO queue



CAMP-MALLOC



FIFO queue



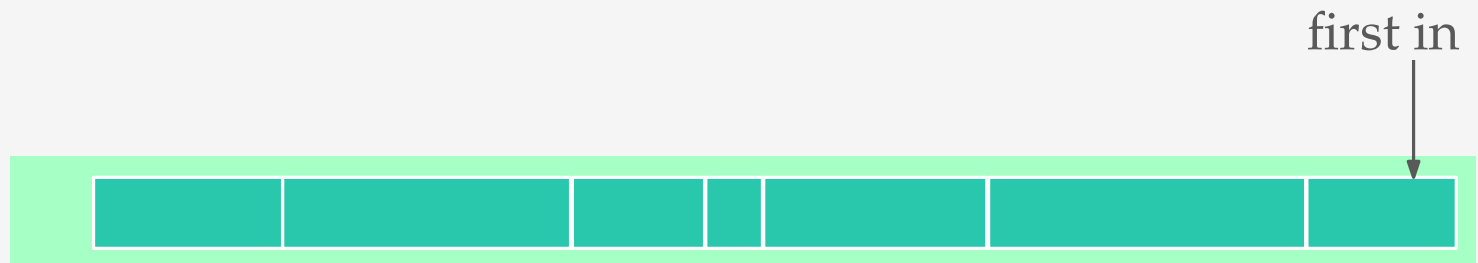
CAMP-MALLOC



FIFO queue



CAMP-MALLOC



FIFO queue



CAMP-MALLOC

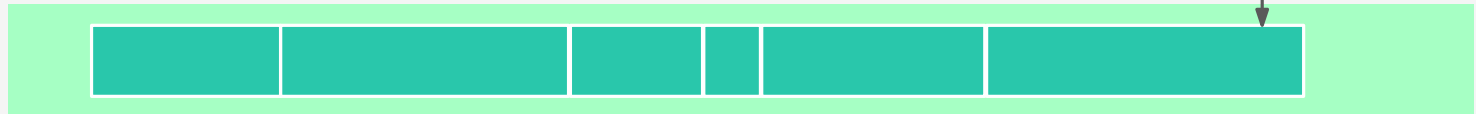


first in
↓

FIFO queue



CAMP-MALLOC



first in

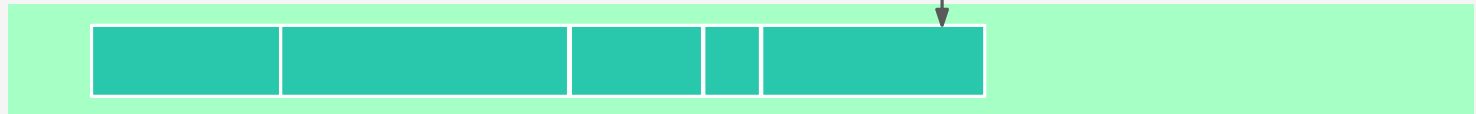
FIFO queue



CAMP-MALLOC



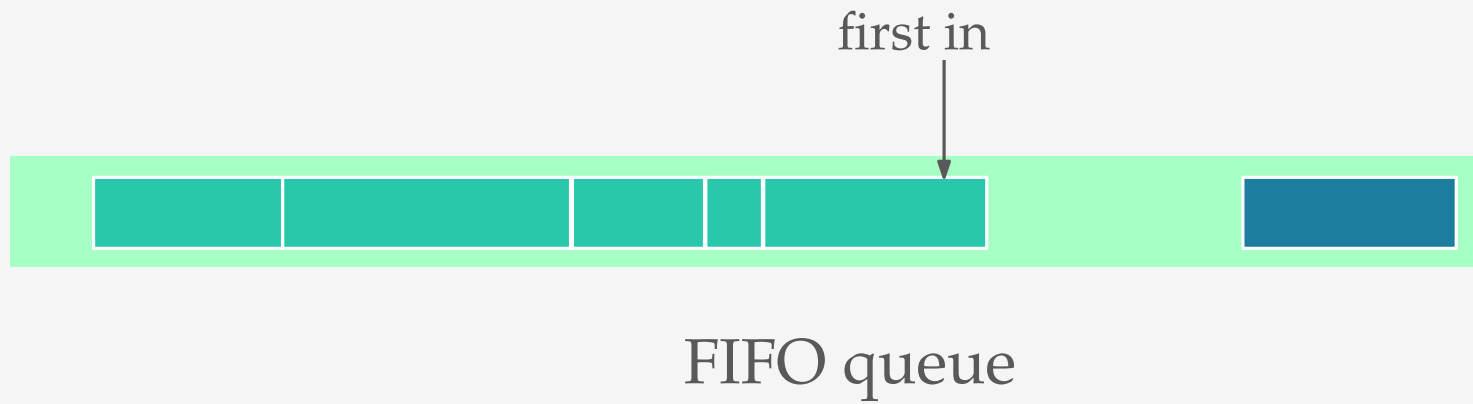
first in



FIFO queue



CAMP-MALLOC





CAMP-MALLOC



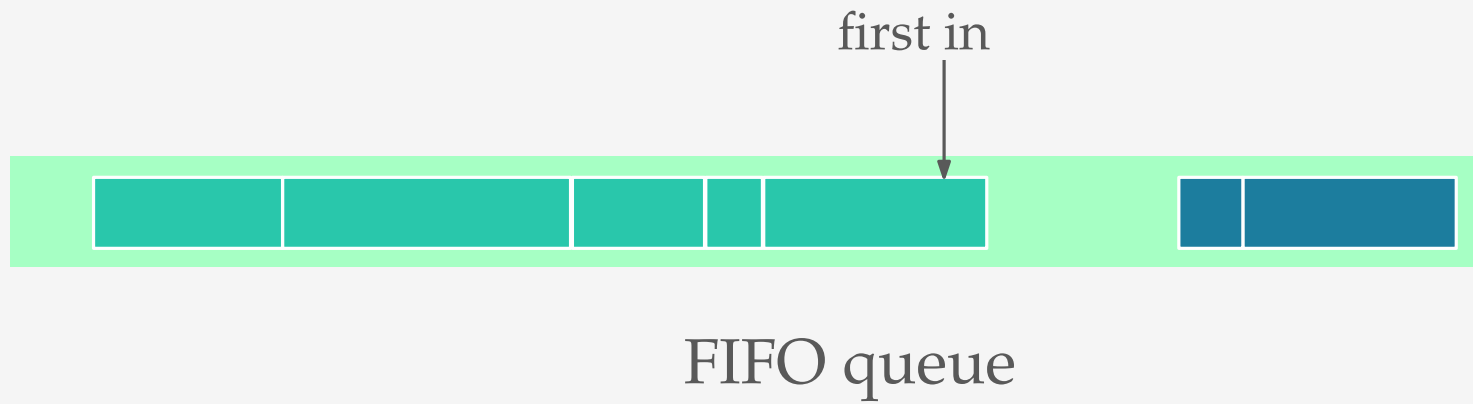
first in



FIFO queue



CAMP-MALLOC





CAMP-MALLOC



first in



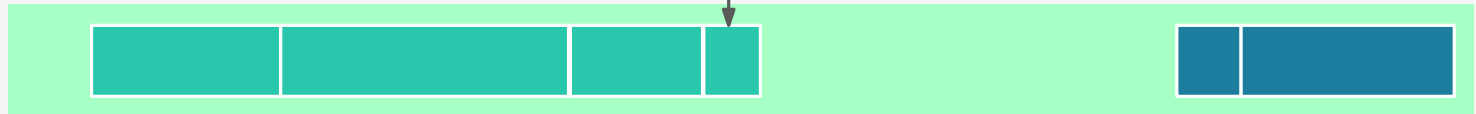
FIFO queue



CAMP-MALLOC



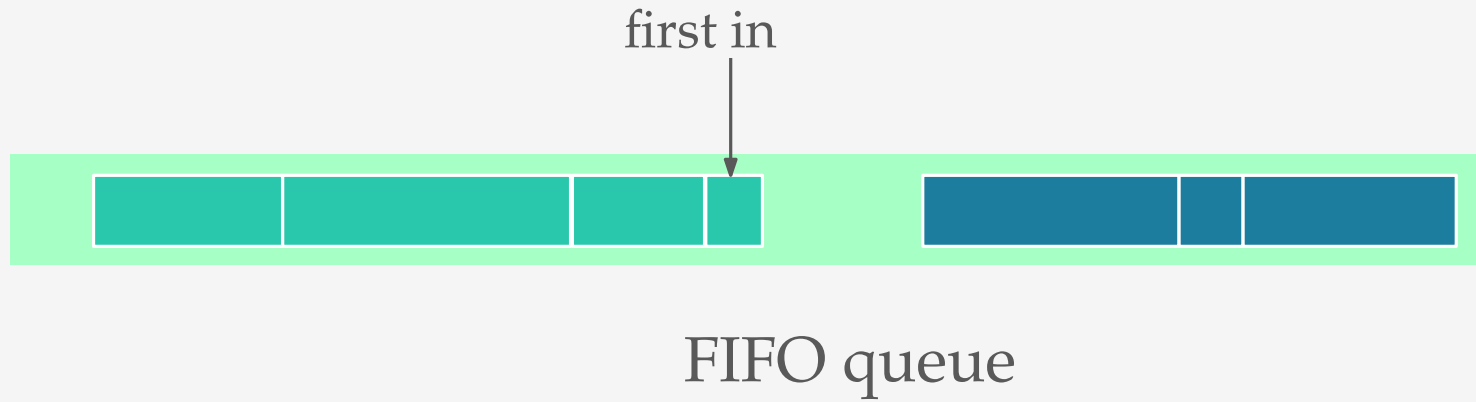
first in



FIFO queue



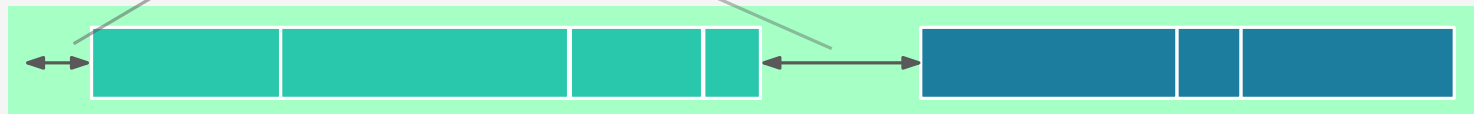
CAMP-MALLOC





CAMP-MALLOC

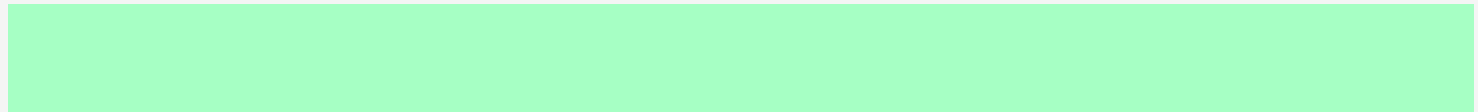
fragmentation ≤ 2 (max item size)



FIFO queue

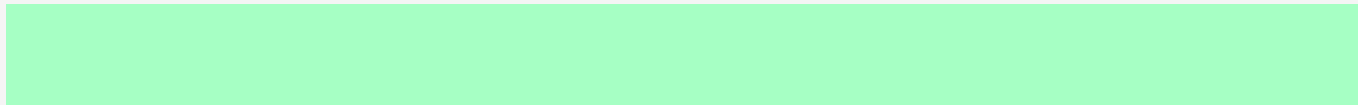


CAMP-MALLOC



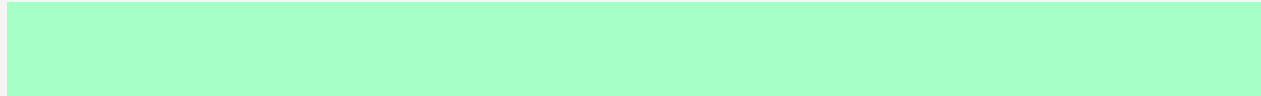


CAMP-MALLOC



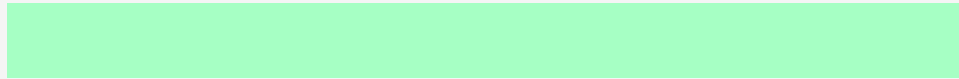


CAMP-MALLOC





CAMP-MALLOC



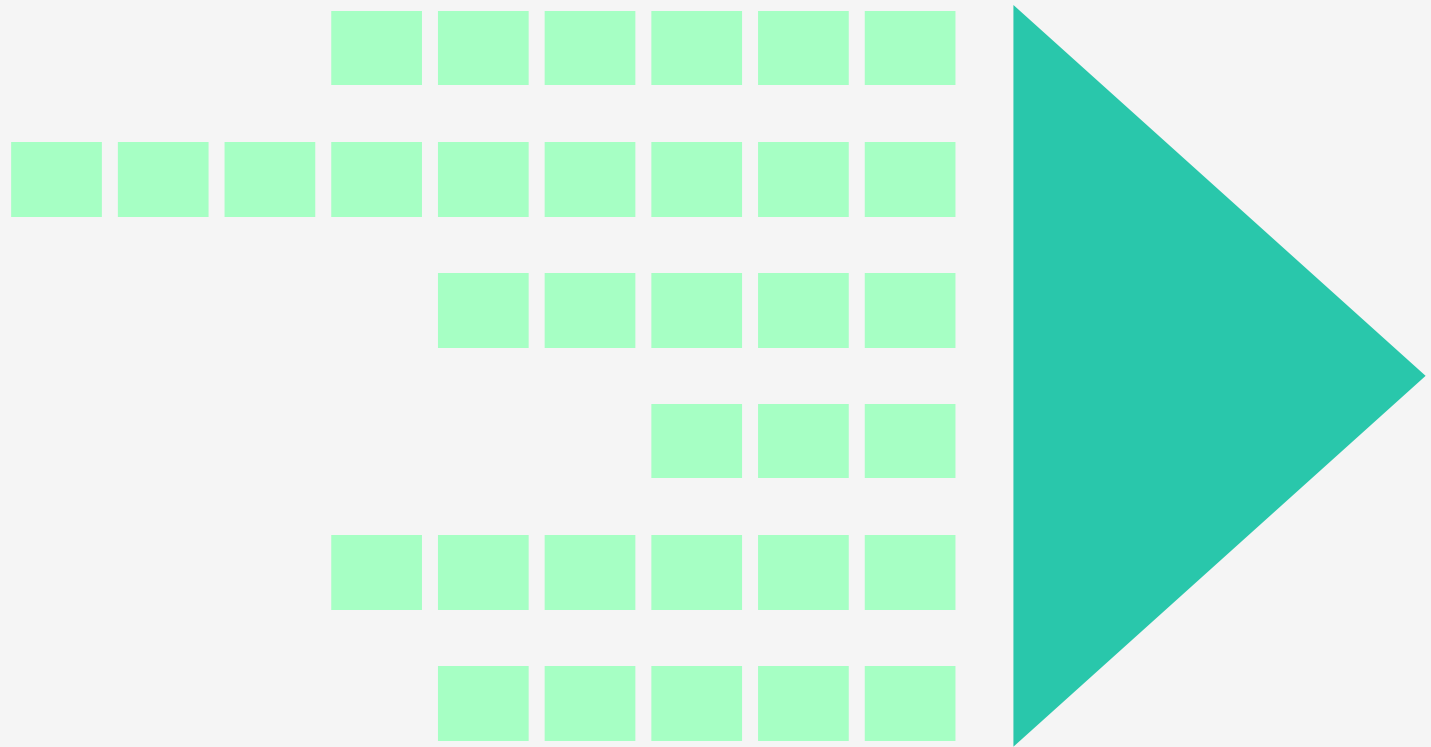


CAMP-MALLOC



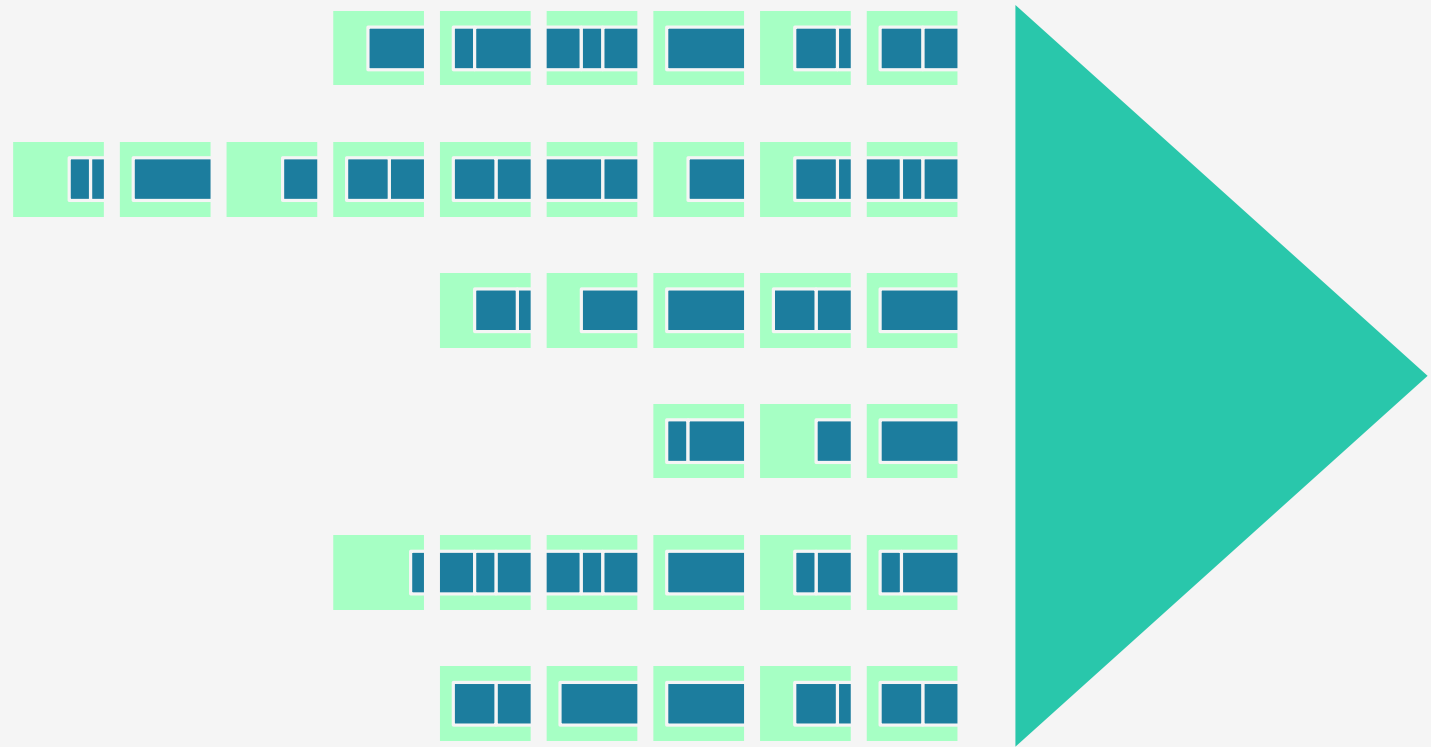


CAMP-MALLOC



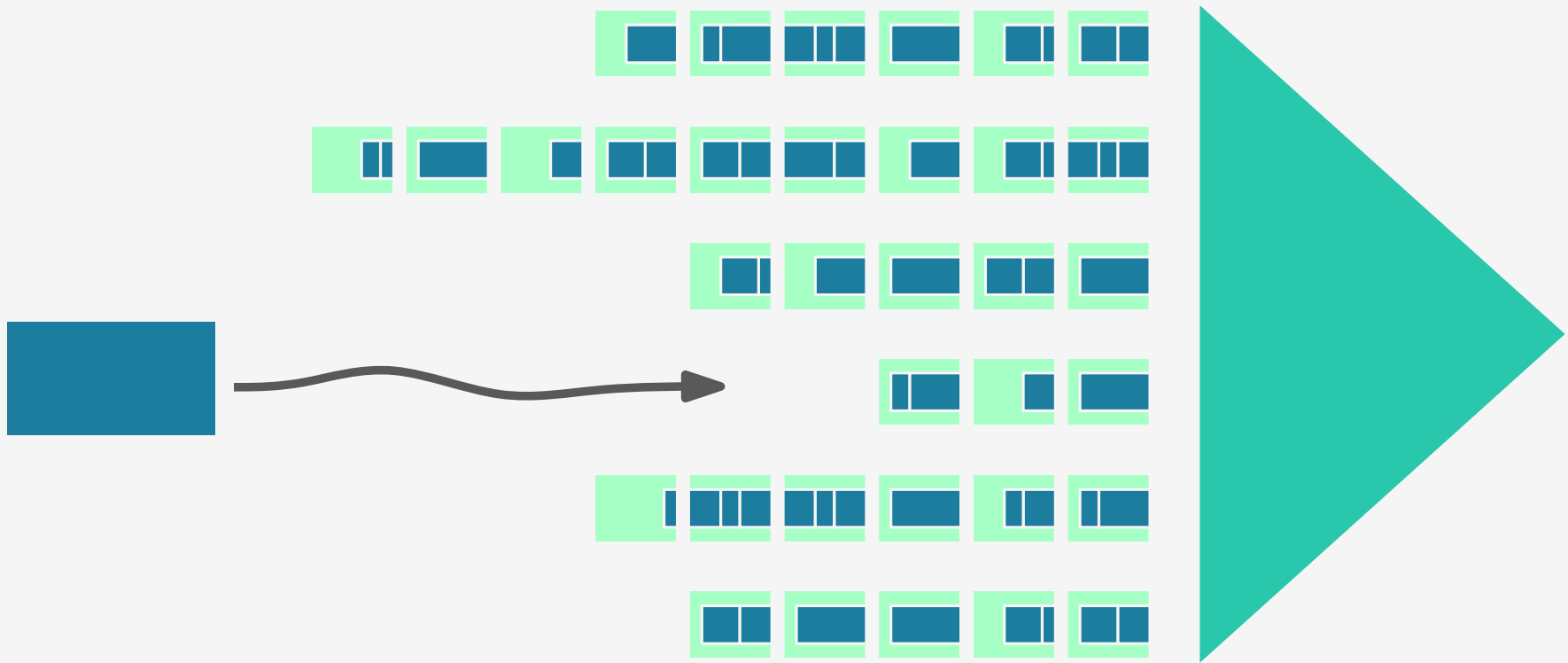


CAMP-MALLOC



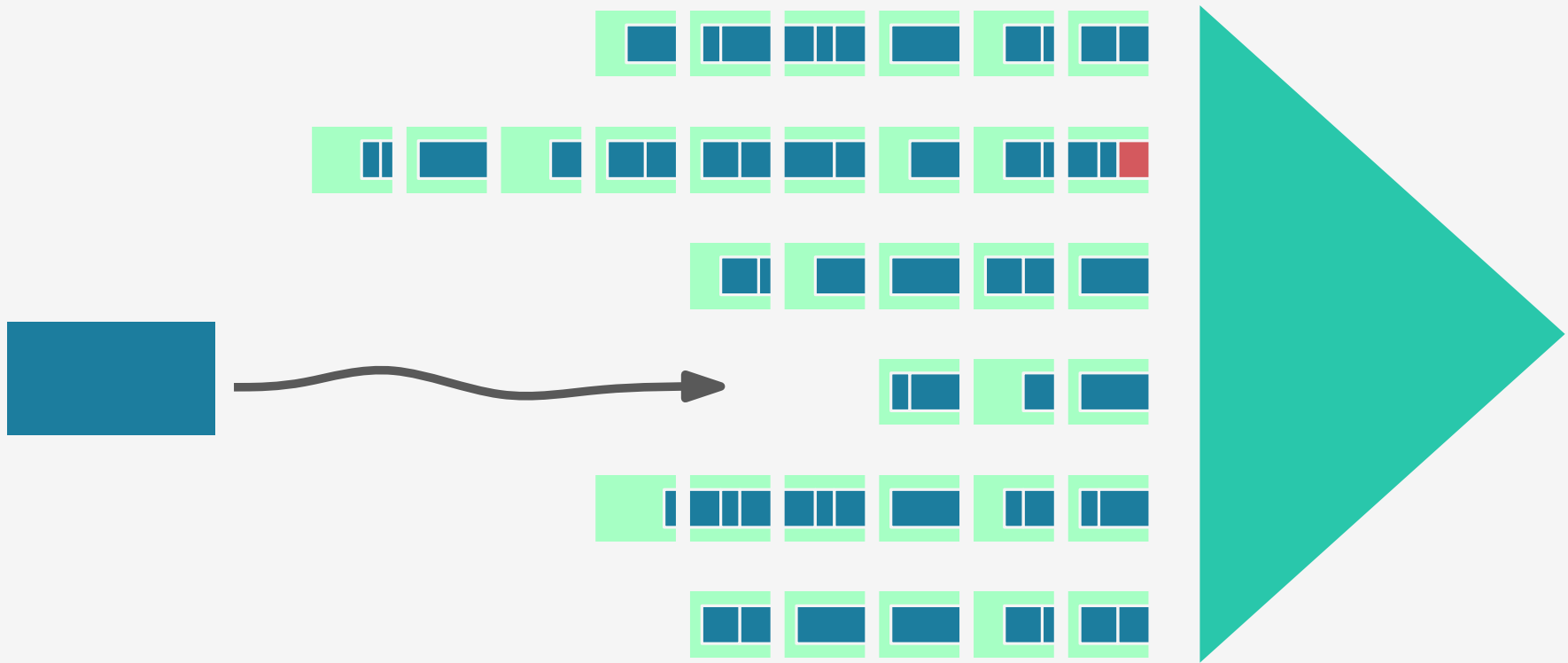


CAMP-MALLOC



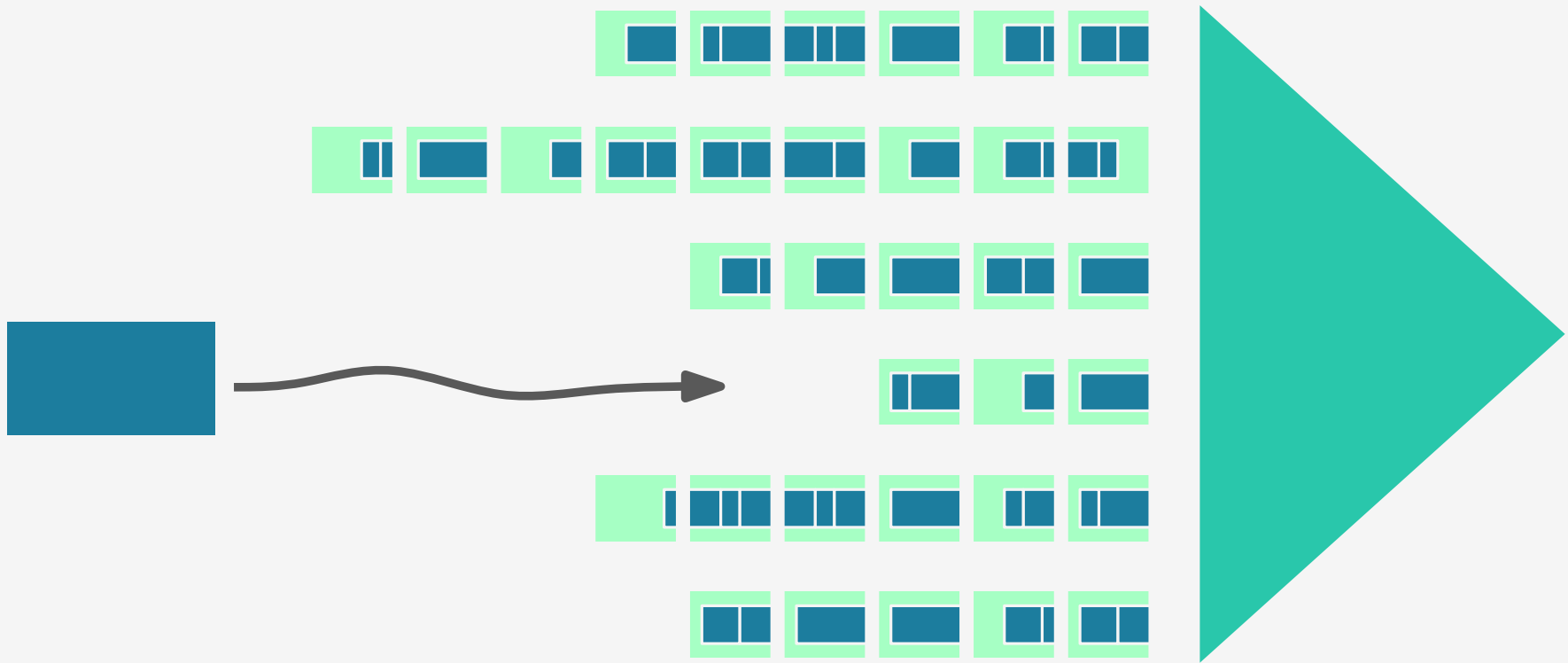


CAMP-MALLOC



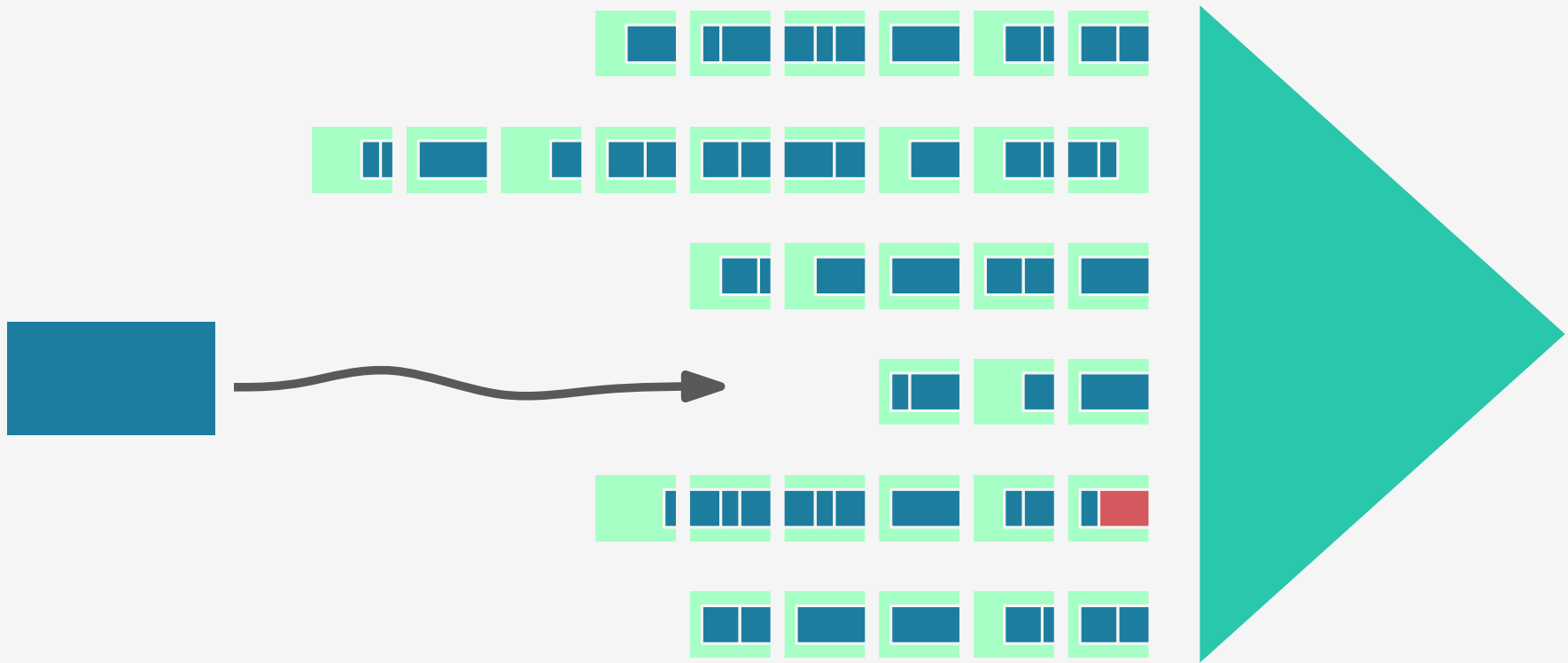


CAMP-MALLOC



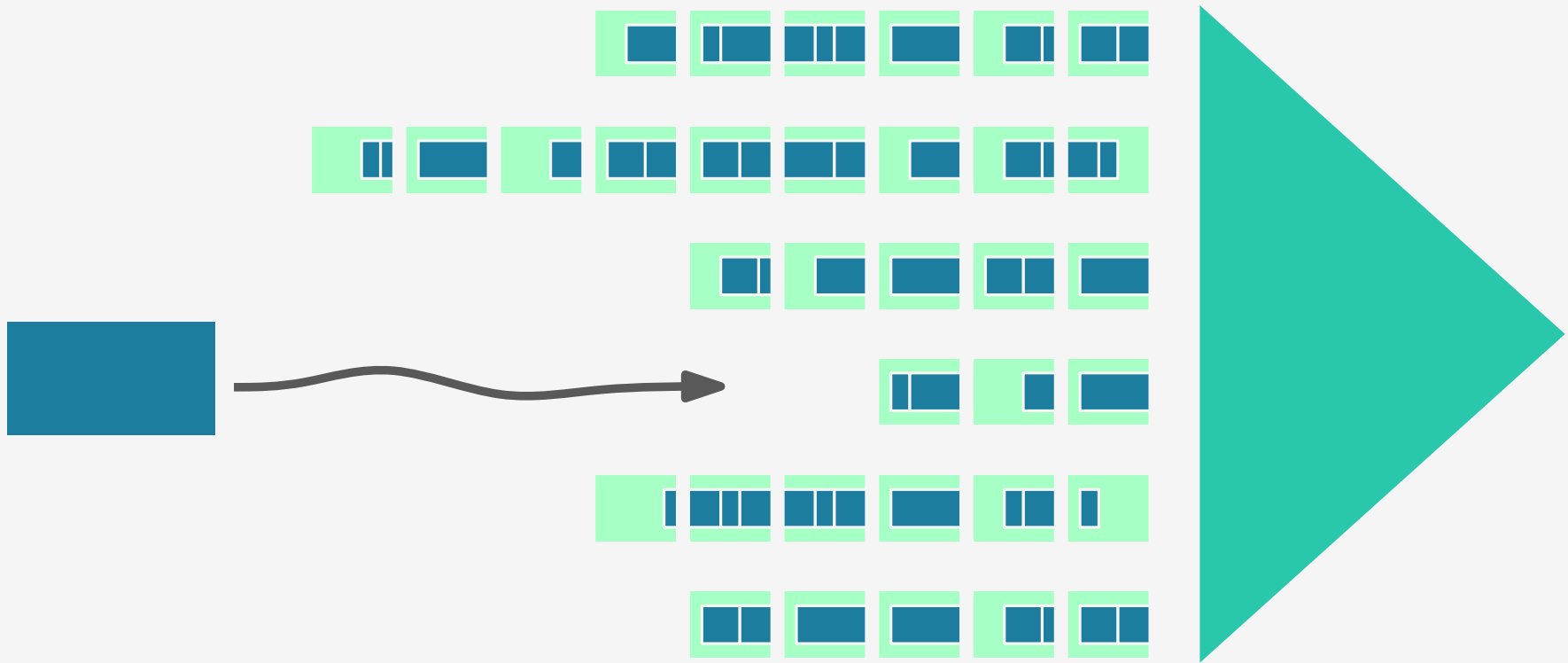


CAMP-MALLOC



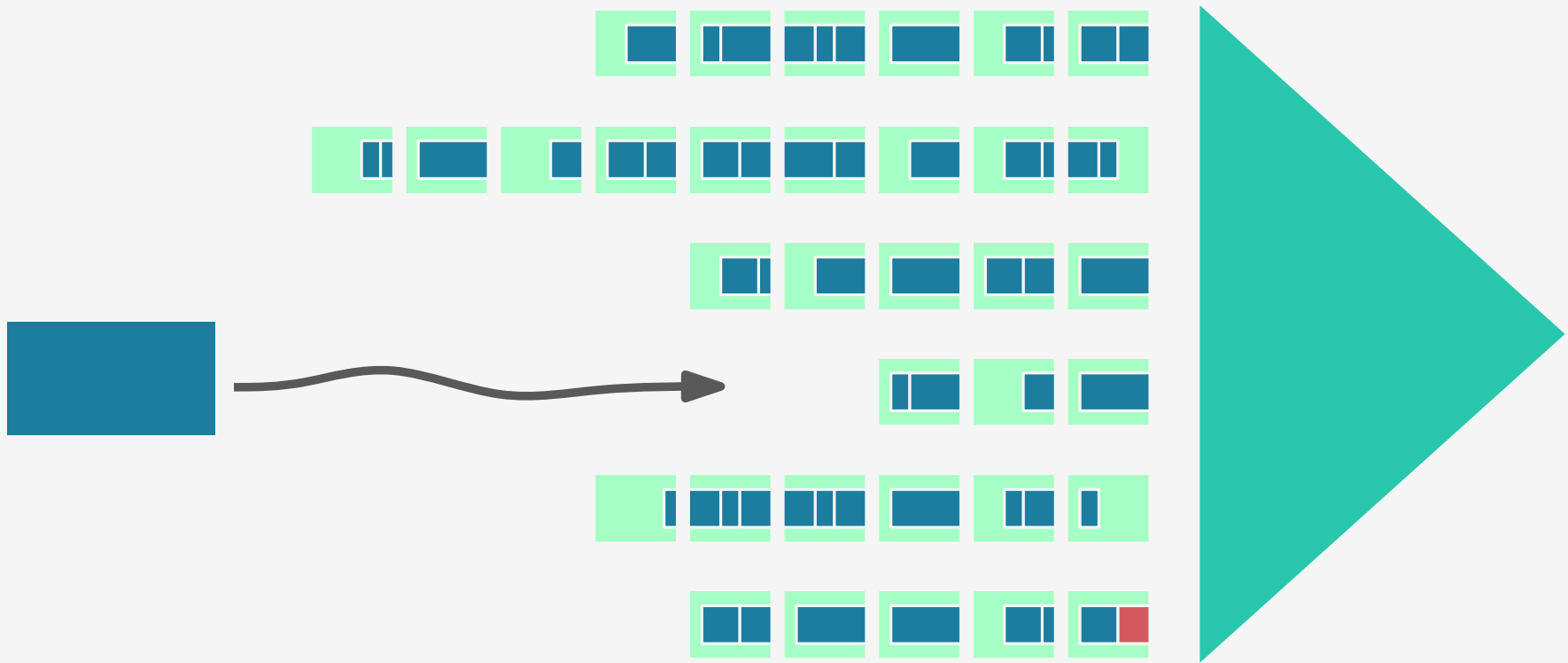


CAMP-MALLOC



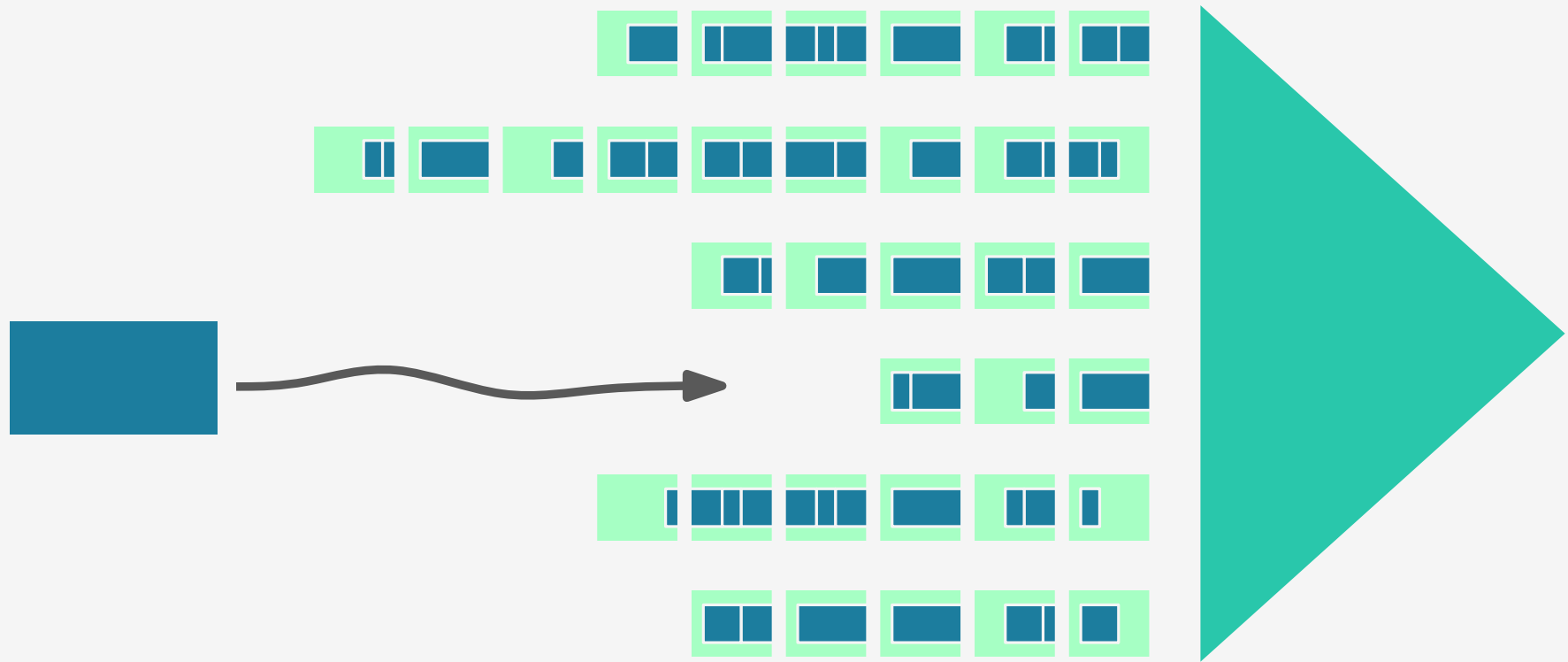


CAMP-MALLOC



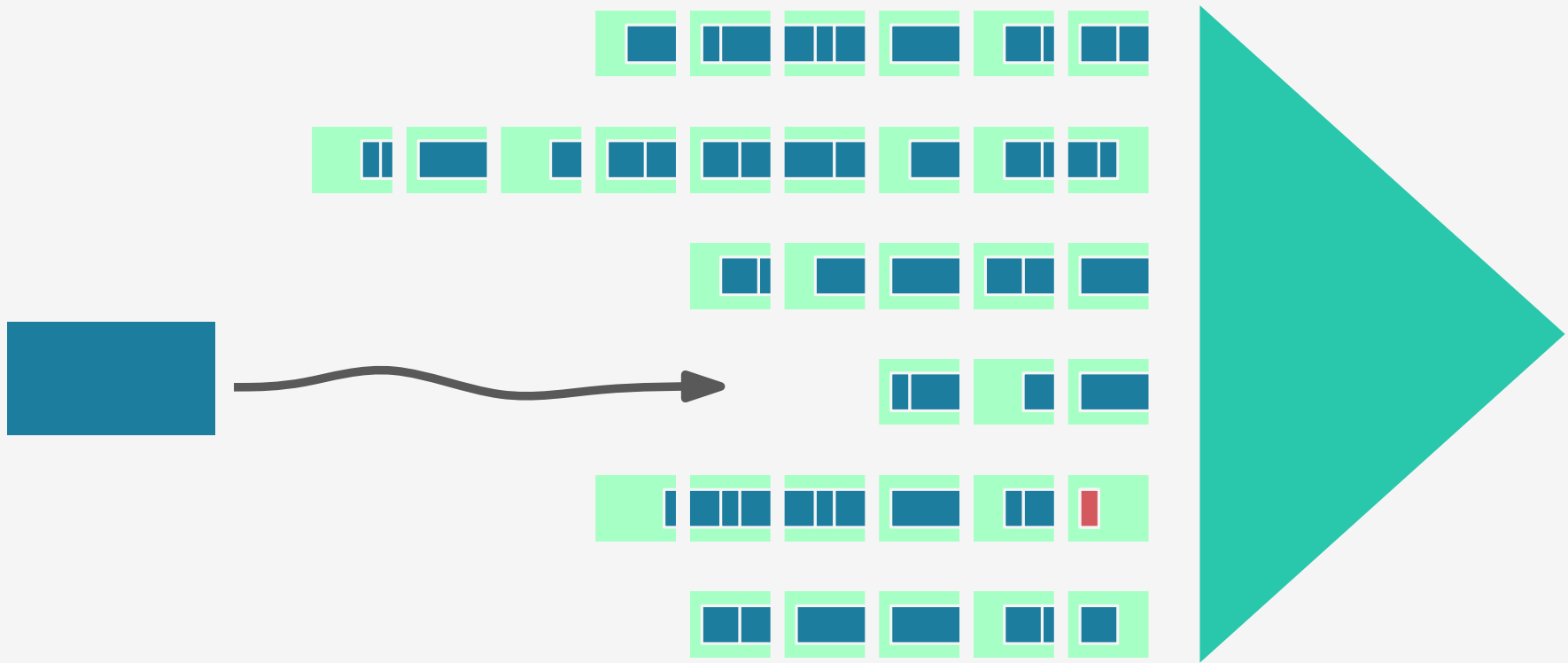


CAMP-MALLOC



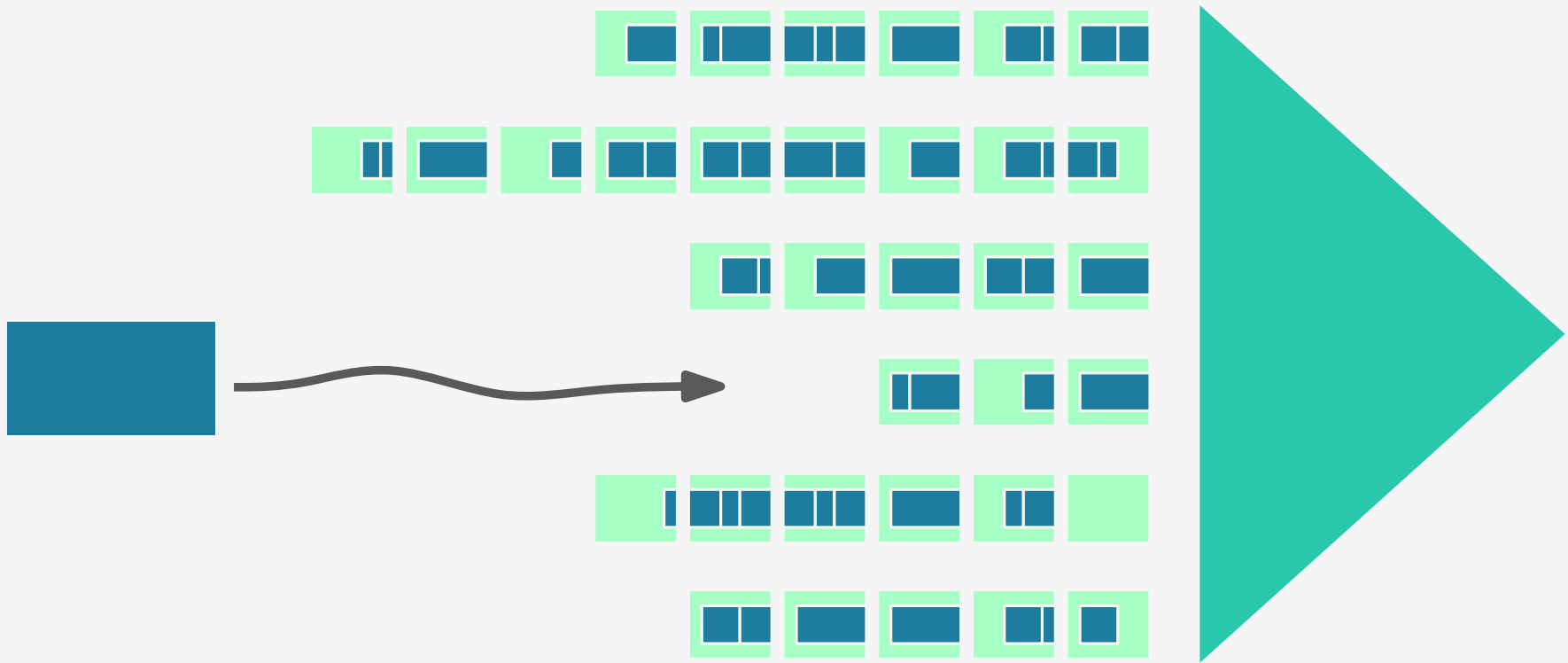


CAMP-MALLOC



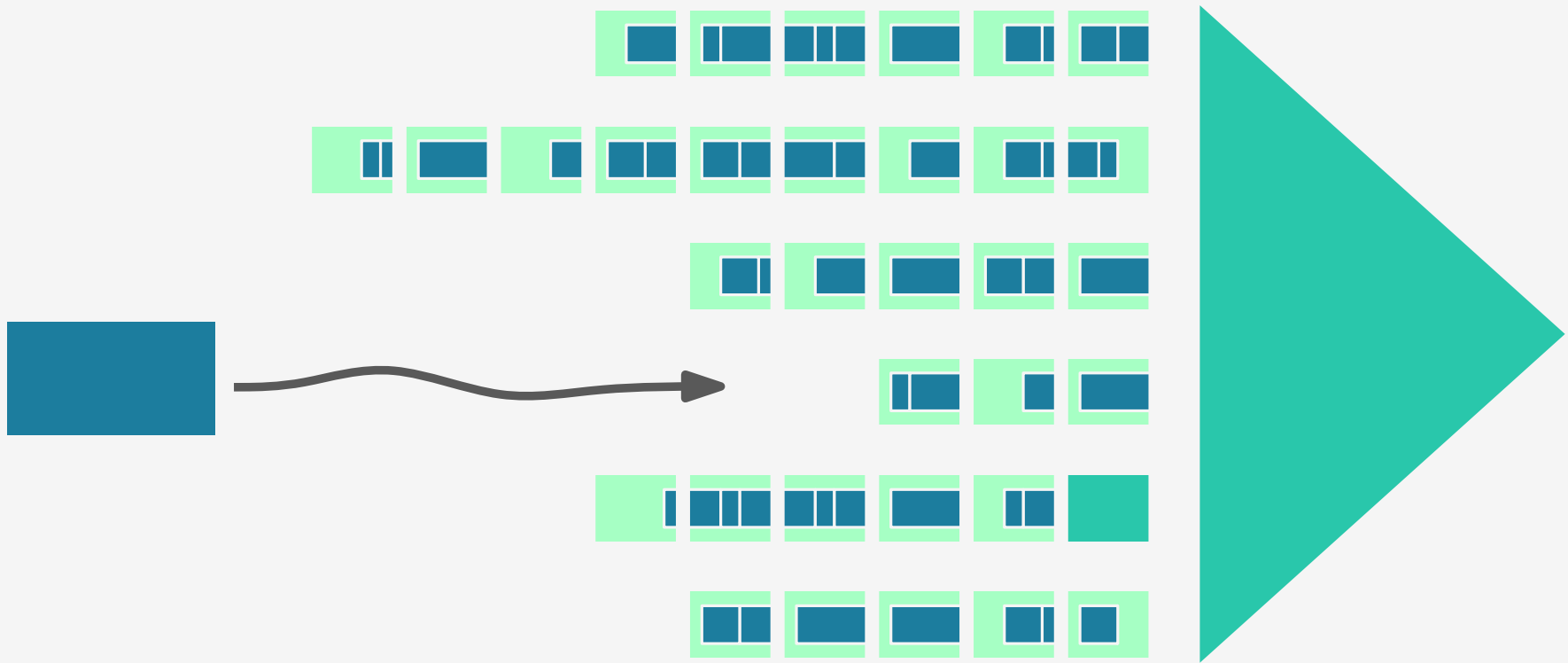


CAMP-MALLOC



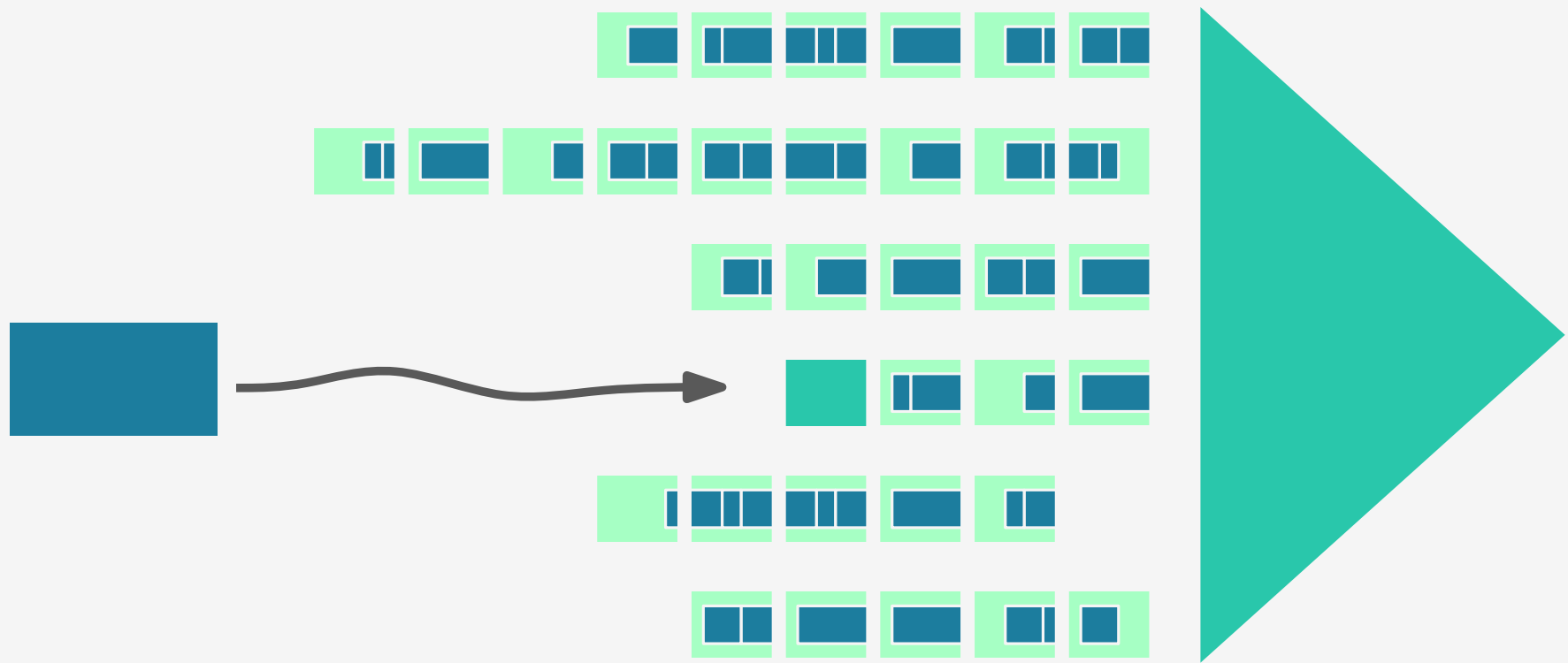


CAMP-MALLOC



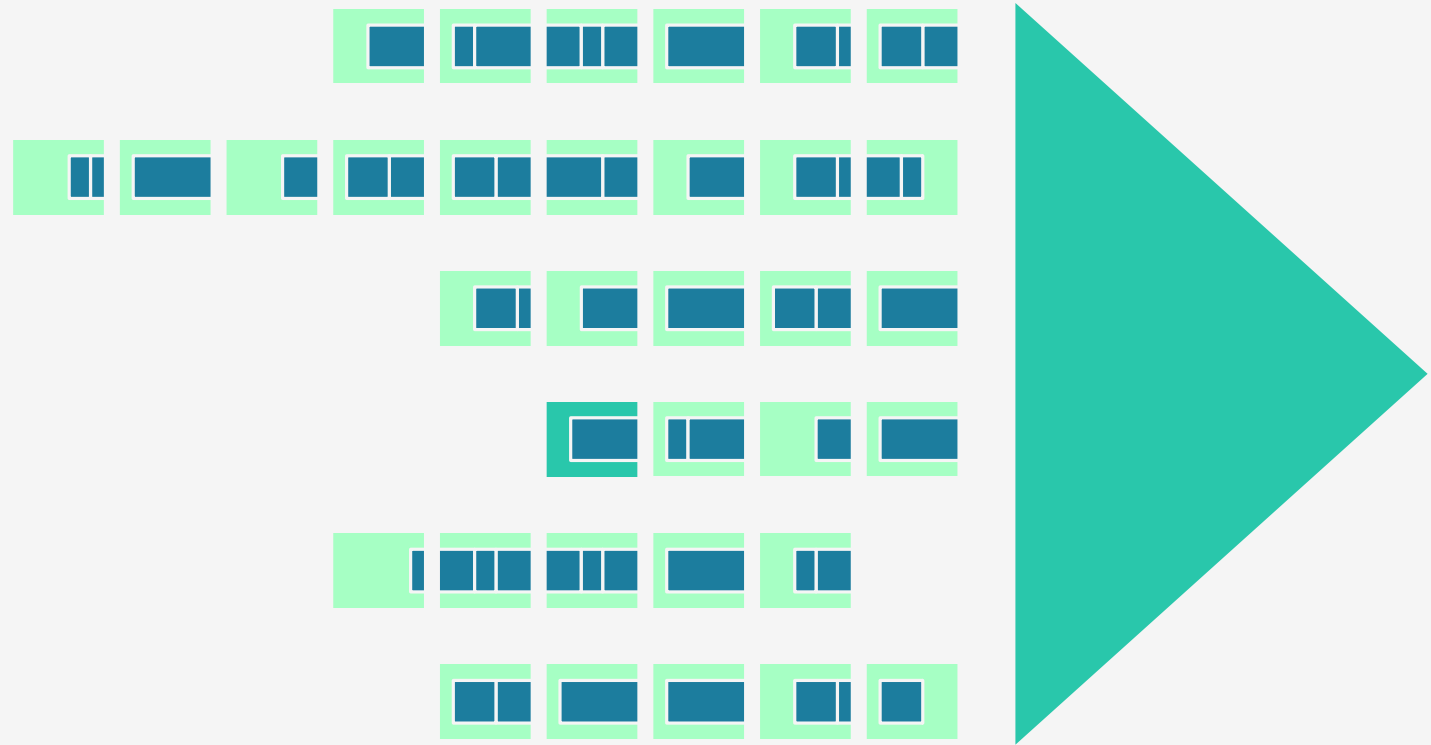


CAMP-MALLOC



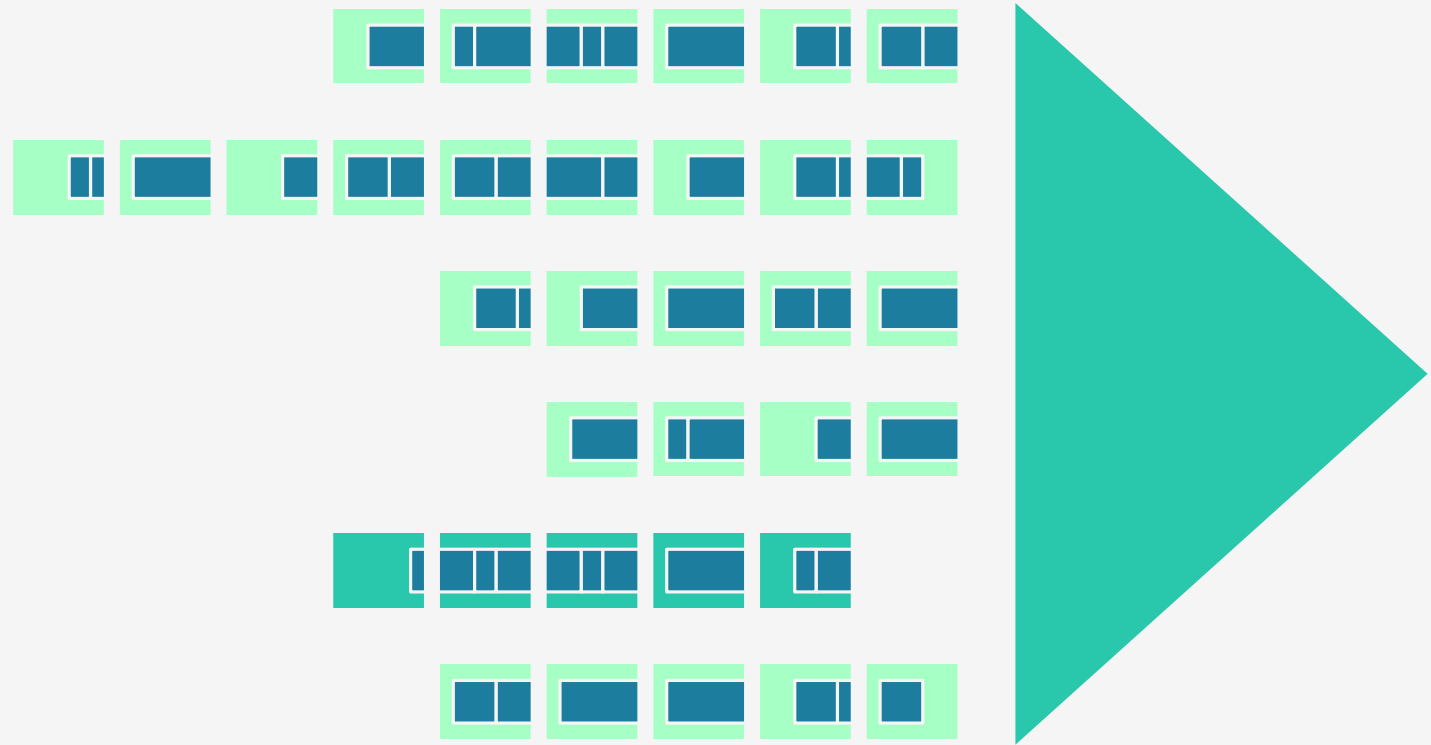


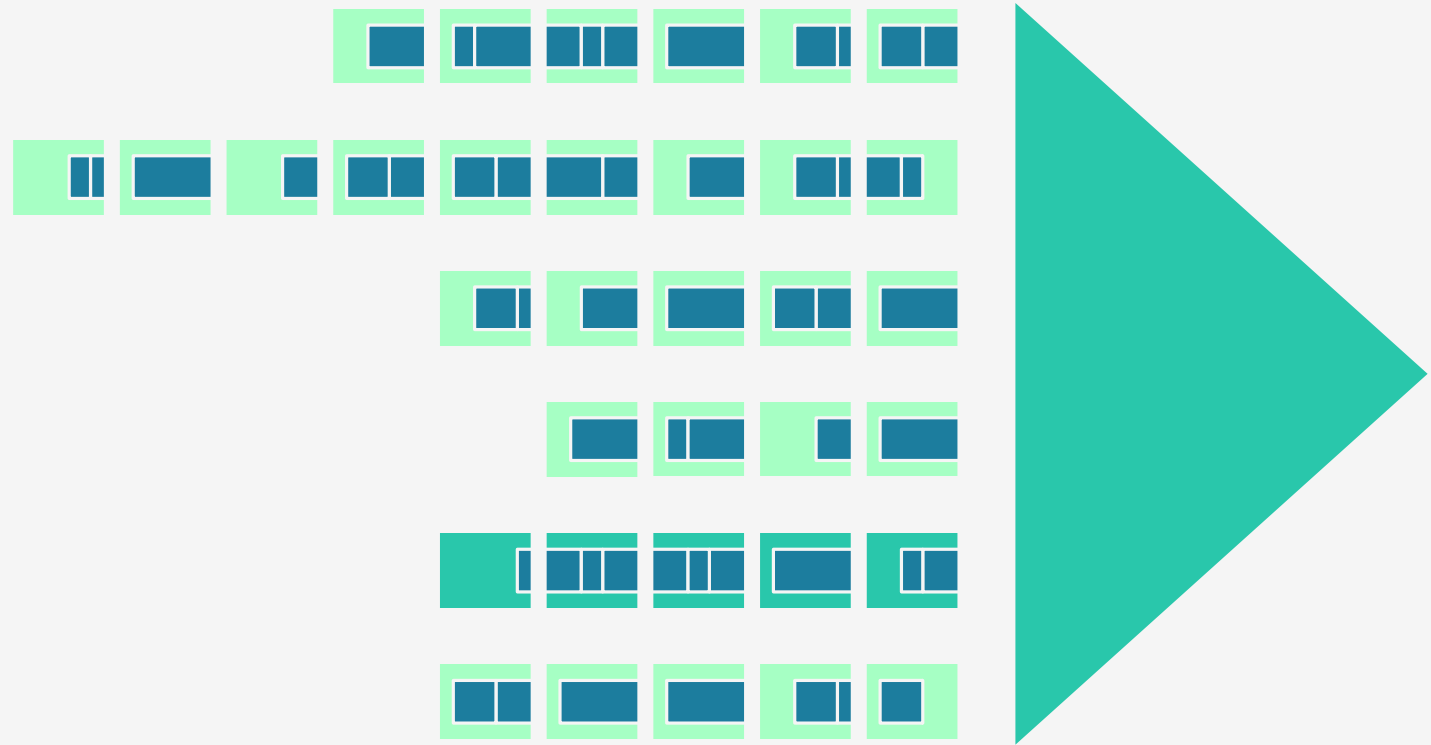
CAMP-MALLOC

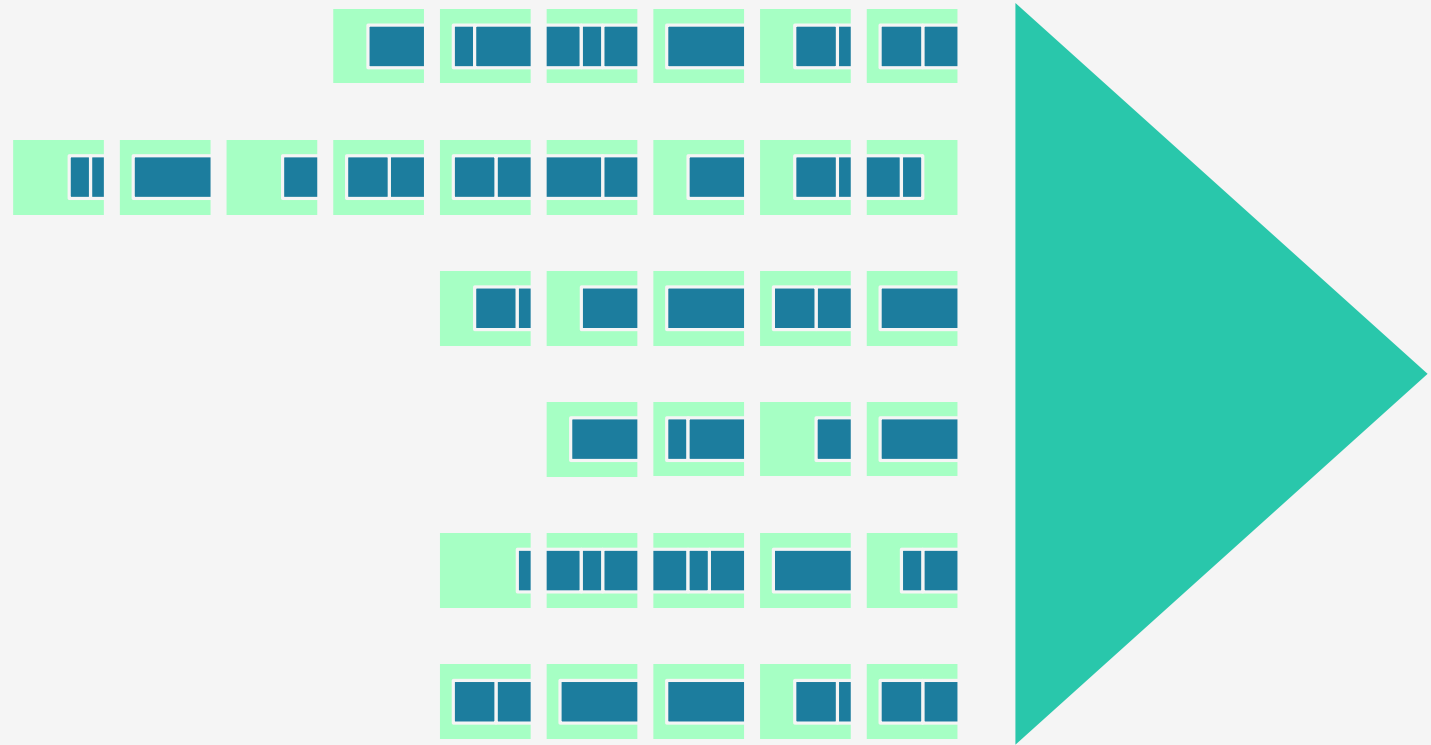




CAMP-MALLOC

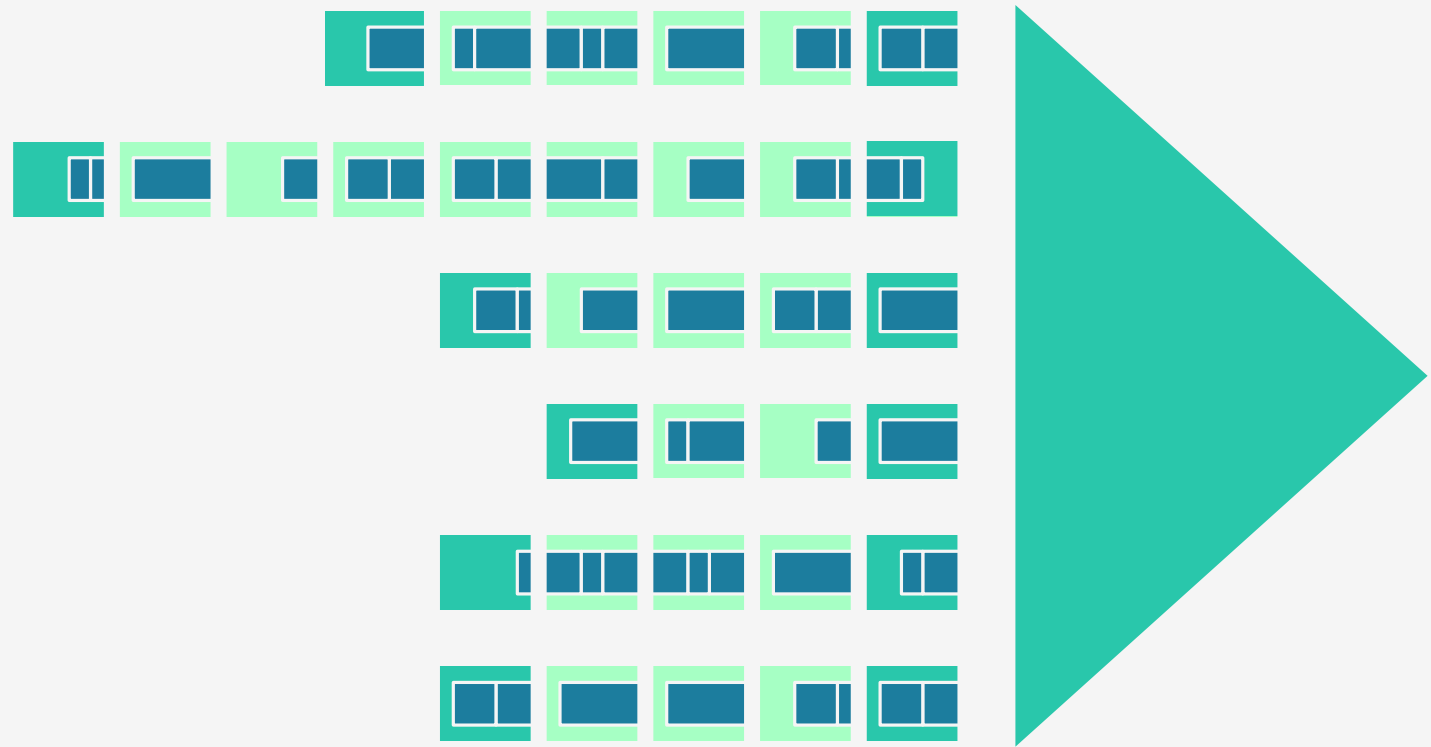








CAMP-MALLOC



fragmentation ≤ 2 (num queues) (block size) + (num blocks) (max item size)



CAMP-MALLOC

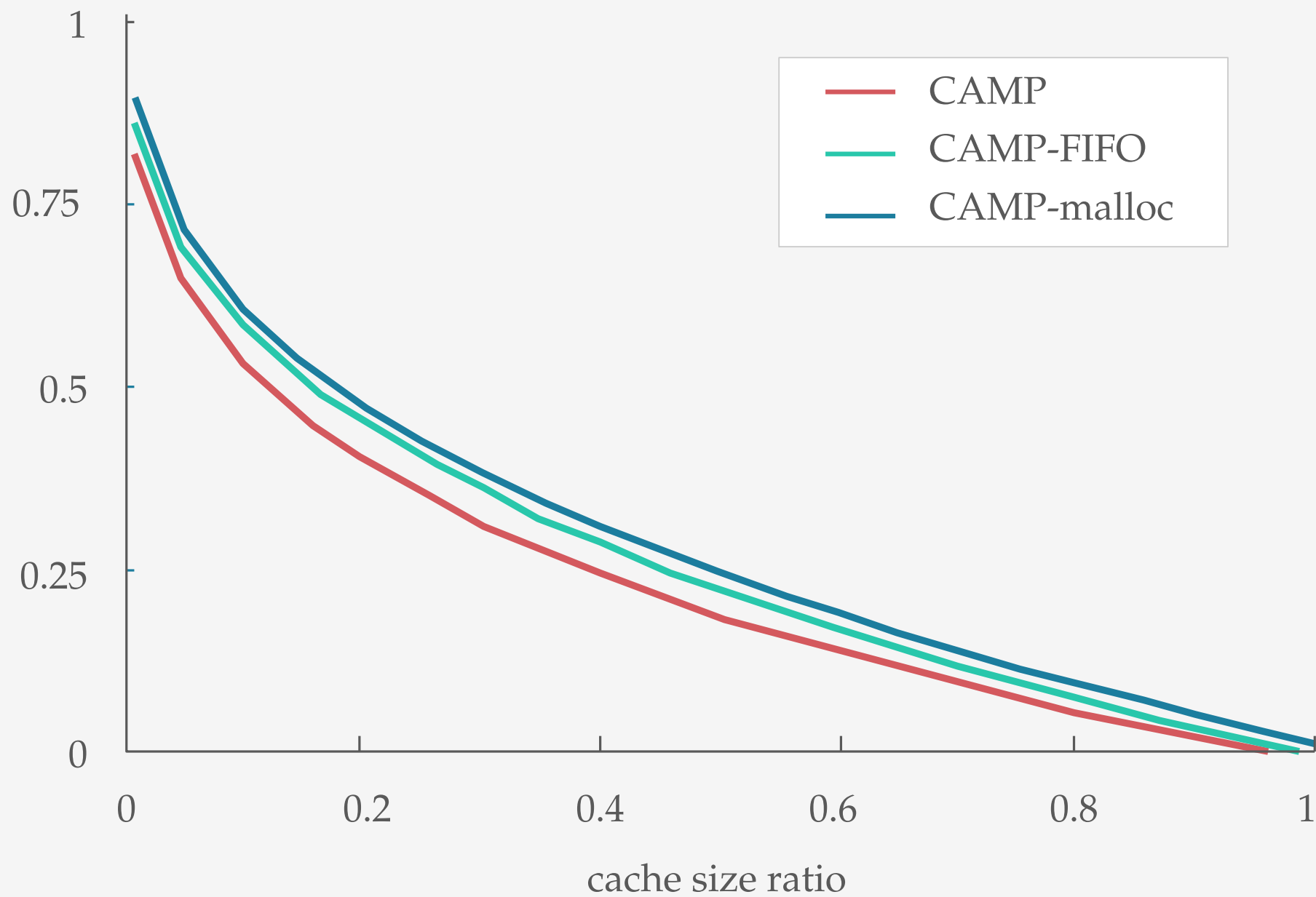
is competitive if memory augmented

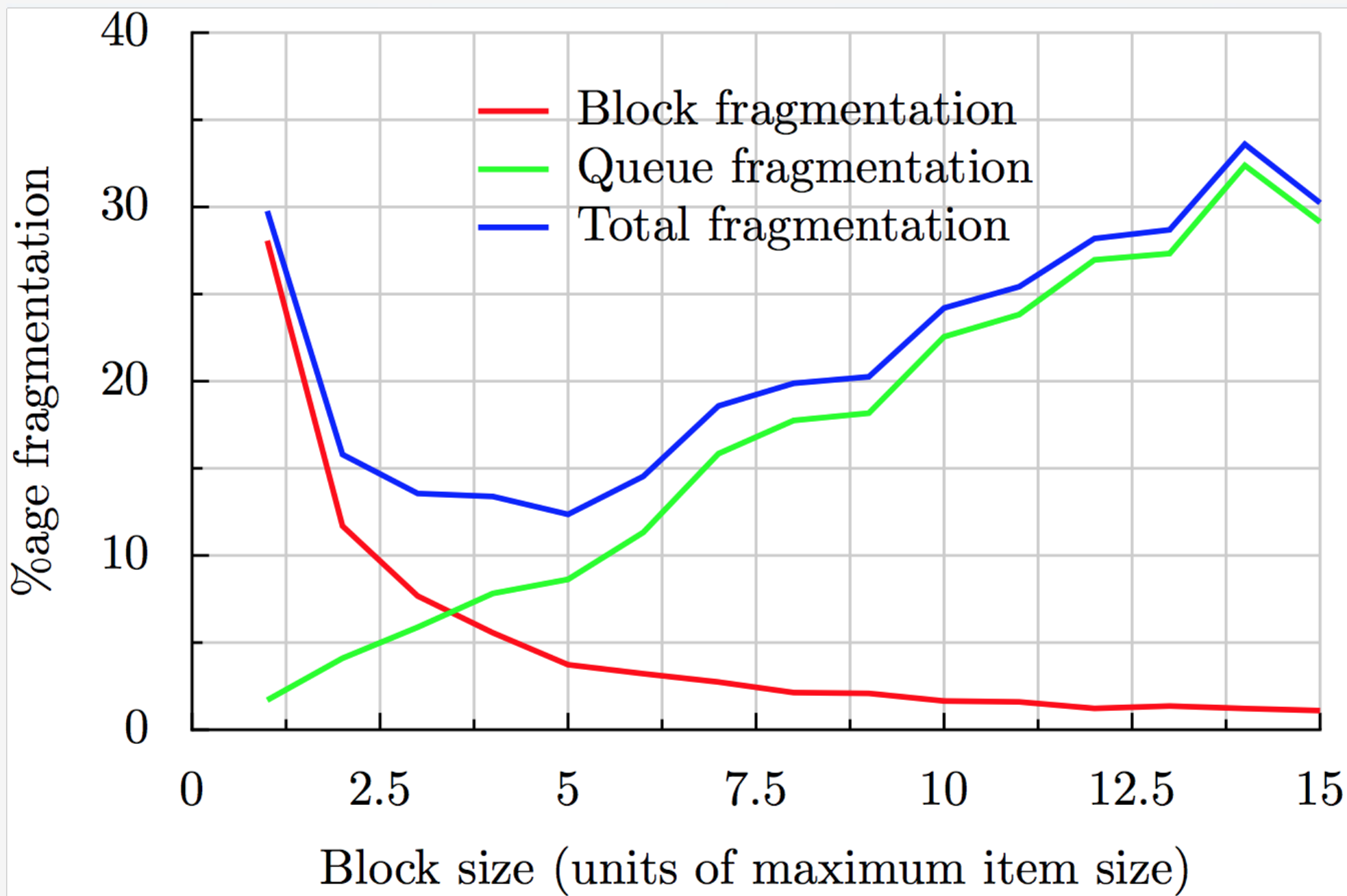
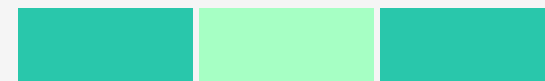
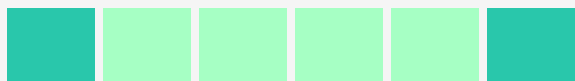
if $\text{OPT's cache size} \leq \text{ALG's cache size} - \text{fragmentation bound}$

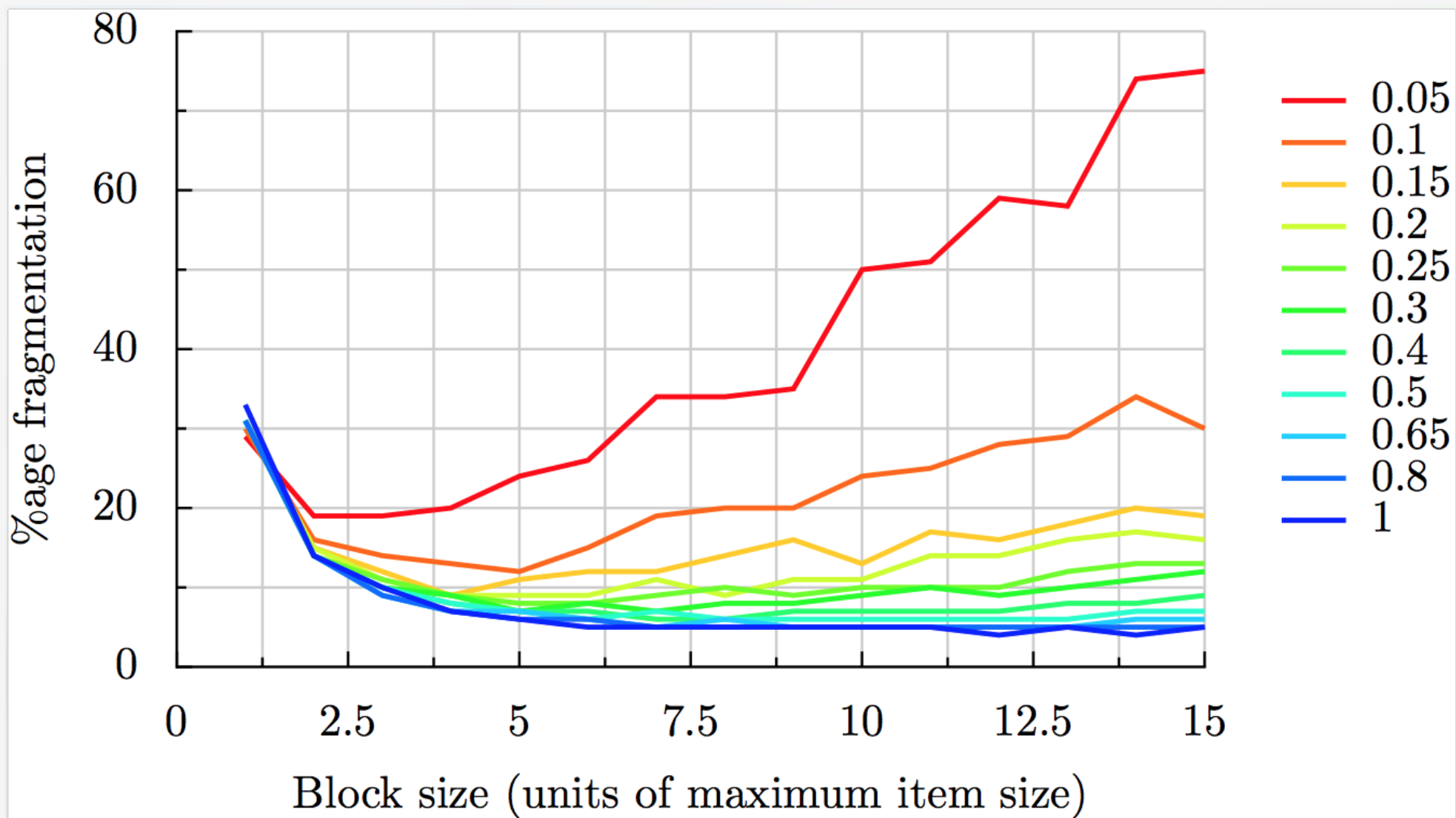
then $\text{cost(ALG)} \leq \frac{\text{ALG's cache size}}{\text{min item size}} \text{cost(OPT)}$

fragmentation ≤ 2 (num queues) (block size) + (num blocks) (max item size)

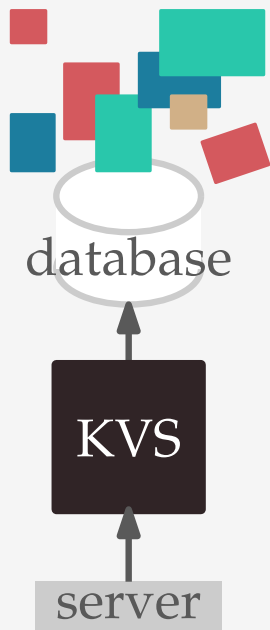
cost-miss ratio







GDS → CAMP



generalized
caching

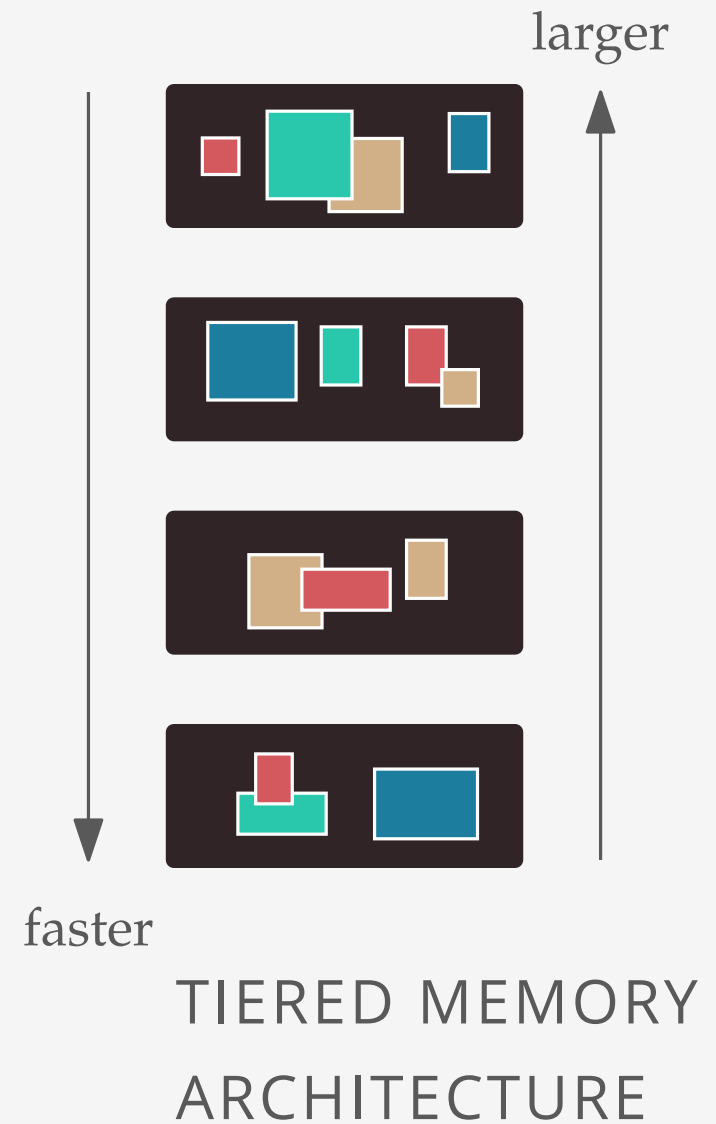
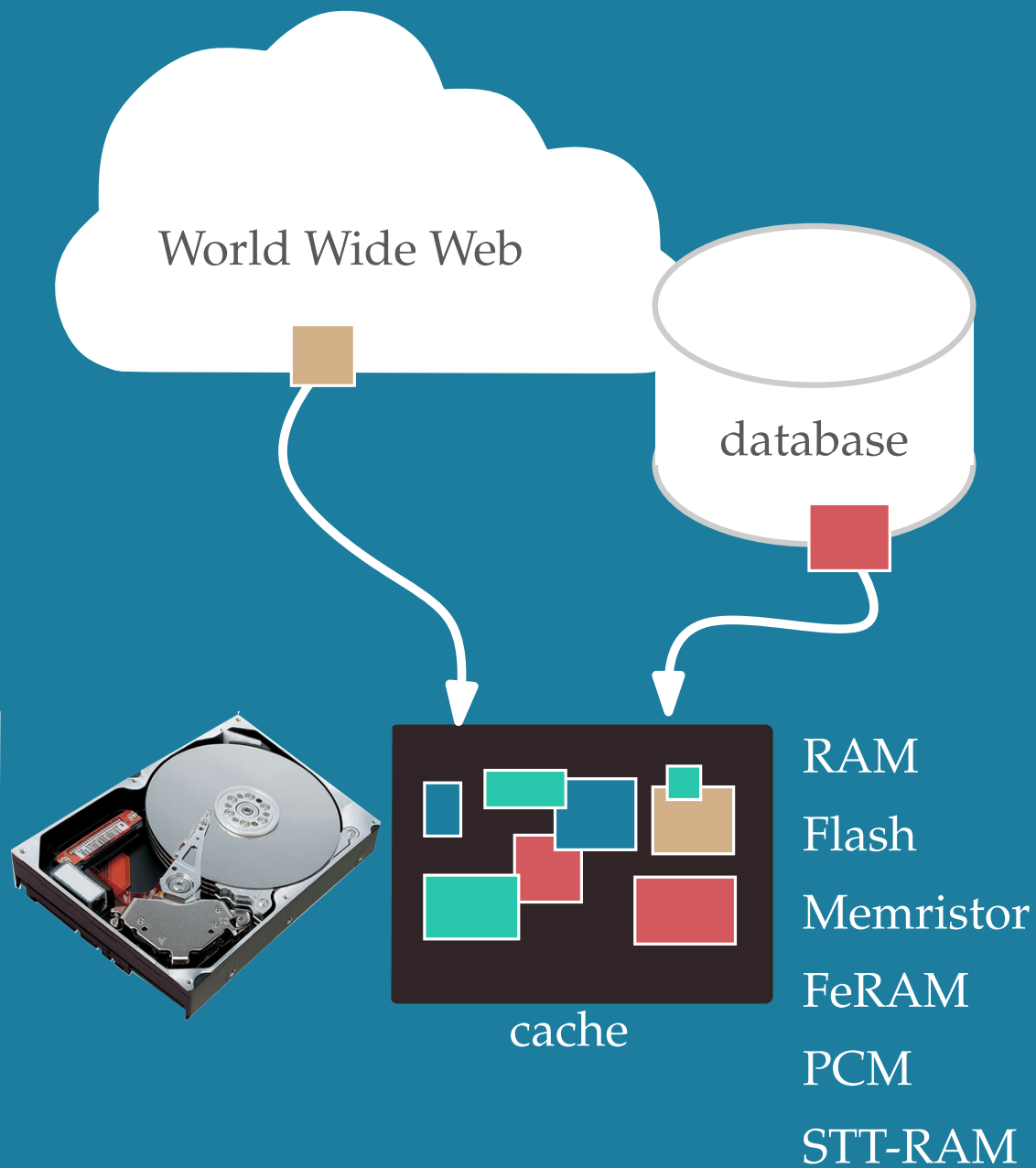


managed memory
caching

2-level cache



multi-level cache



GENERALIZED CACHING

capacity

read speed

write speed

failure rate

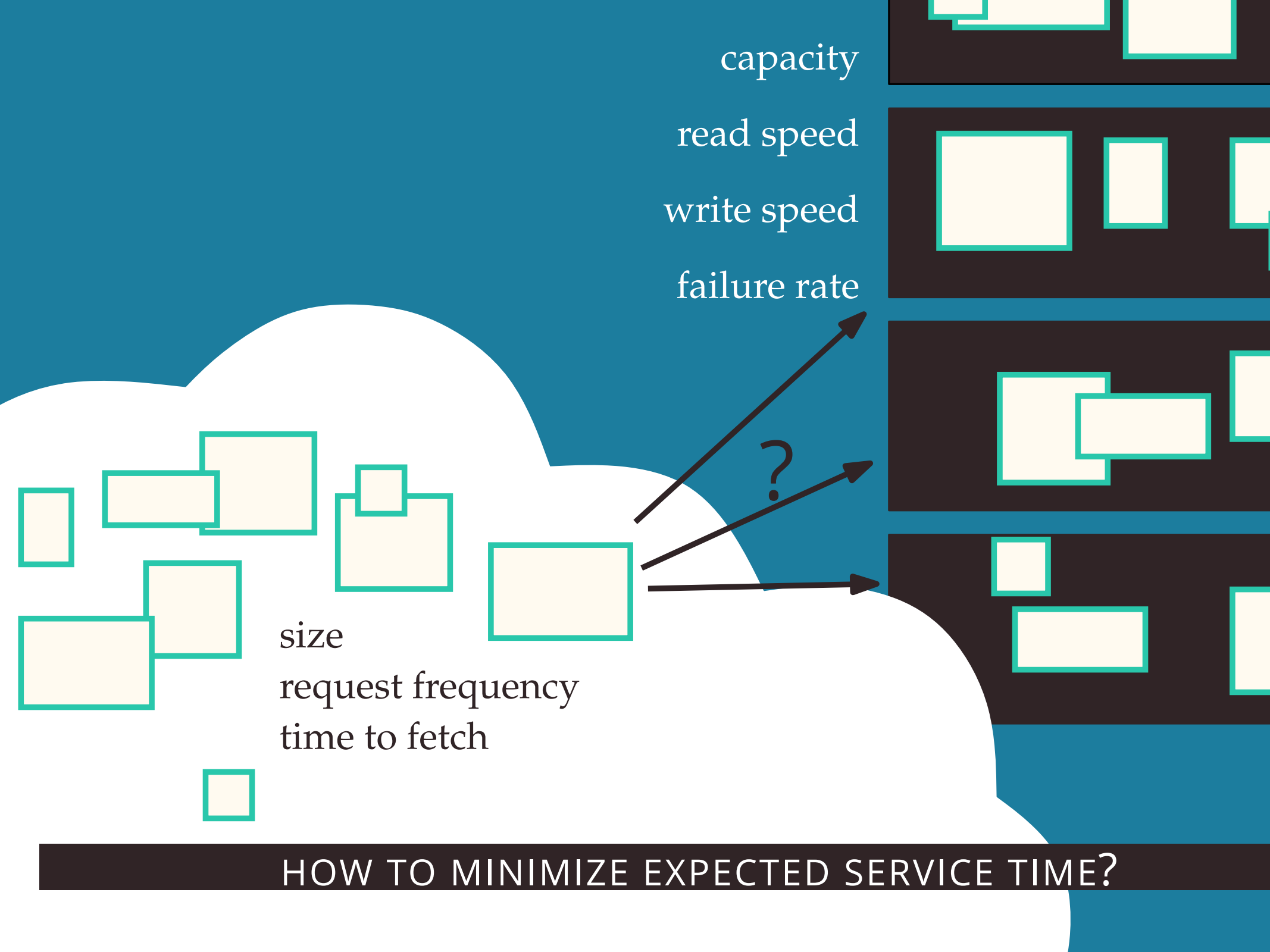
size

request frequency

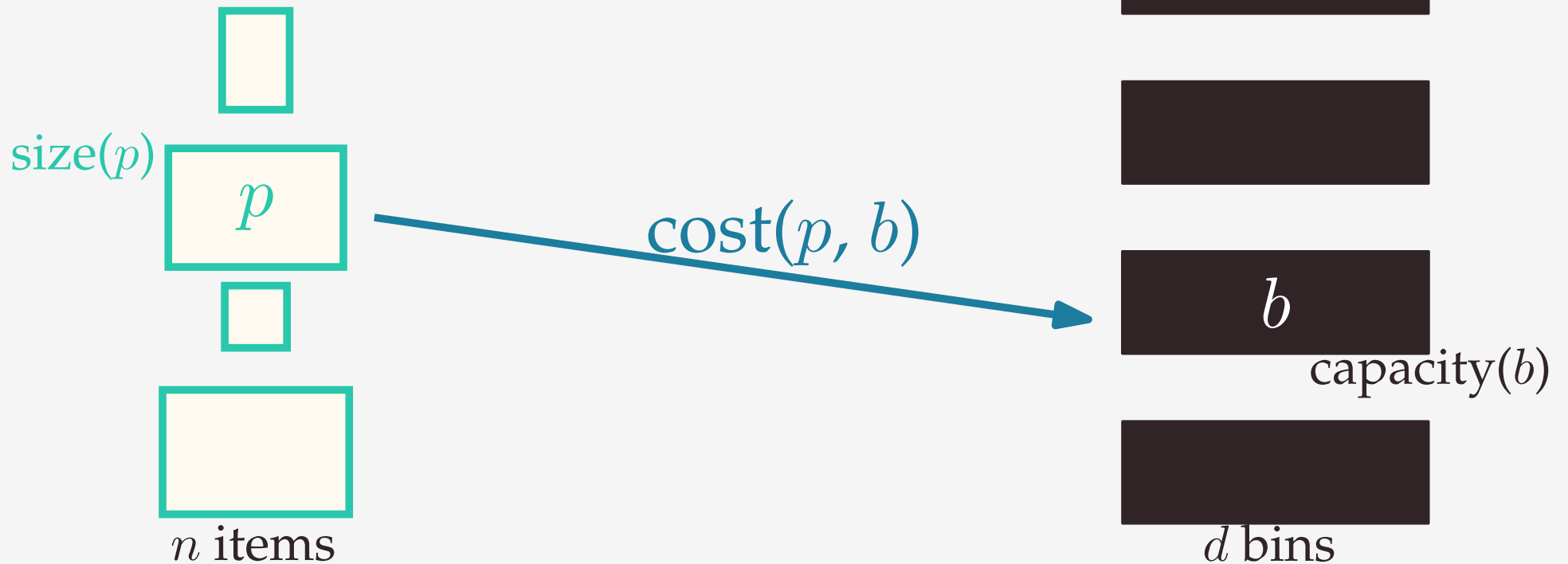
time to fetch

?

HOW TO MINIMIZE EXPECTED SERVICE TIME?



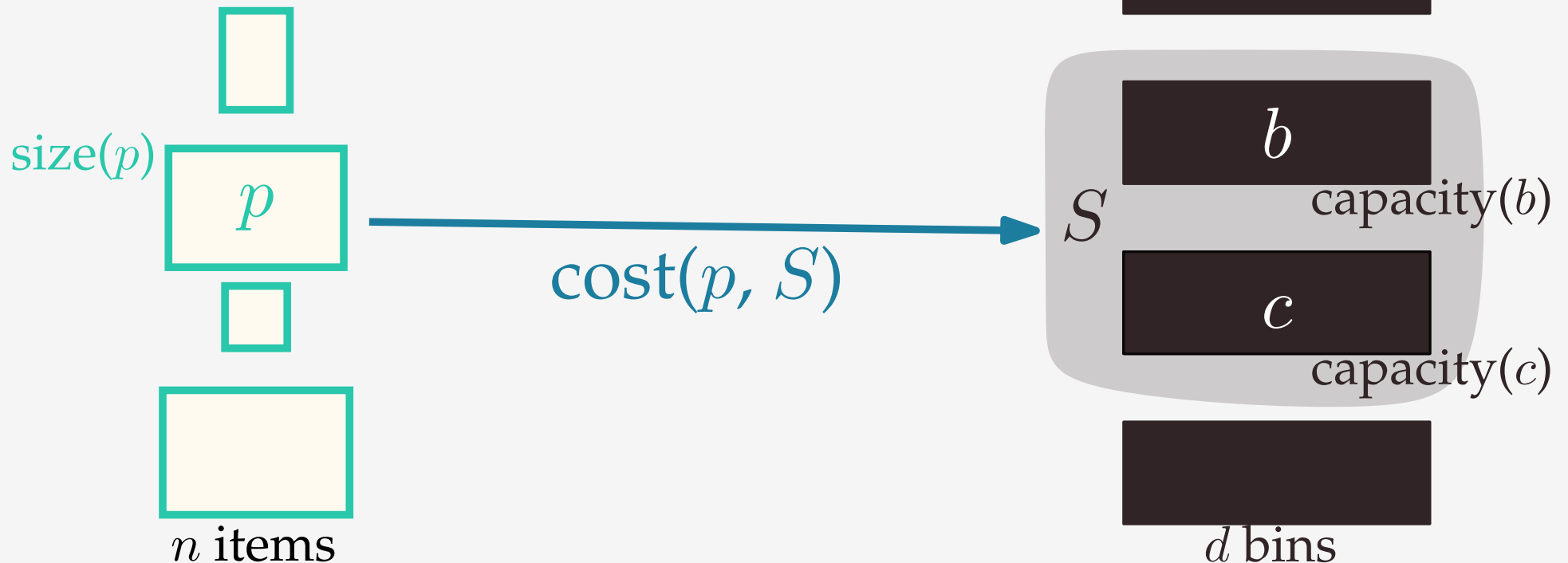
MULTIPLE KNAPSACK PROBLEM



GOAL

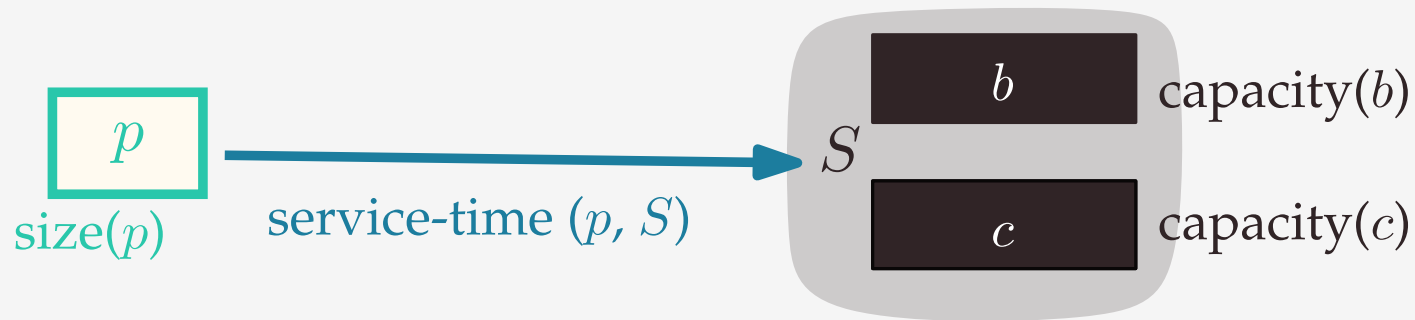
minimize total cost of assignment
subject to capacity constraints

SUBSET ASSIGNMENT PROBLEM

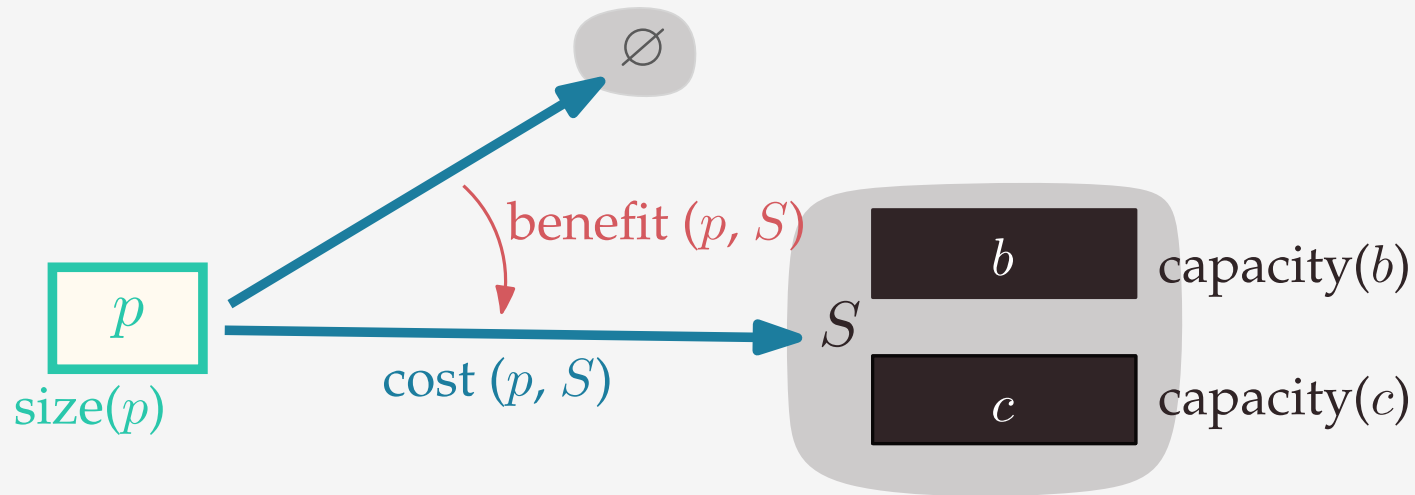


GOAL

minimize total cost of assignment
subject to capacity constraints



$$\begin{aligned}
 \text{service-time}(p, S) = & \text{read-frequency}(p) \text{ read-time}(p, S) \\
 & + \text{write-frequency}(p) \text{ write-time}(p, S) \\
 & + \sum_{F \subseteq S} \text{fail-freq}(F) \left(\text{read-time}(p, S \setminus F) \right. \\
 & \quad \left. + \text{write-time}(p, S \cap F) \right)
 \end{aligned}$$



cache configuration

$$\text{maximize } \sum_{p, S} \text{benefit}(p, S) x(p, S)$$

$$\sum_S x(p, S) = 1$$

$$\sum_{p, S} \text{price}(p, S) x(p, S) \leq \text{budget}$$

$$x = 0, 1$$

subset assignment

$$\text{minimize } \sum_{p, S} \text{cost}(p, S) x(p, S)$$

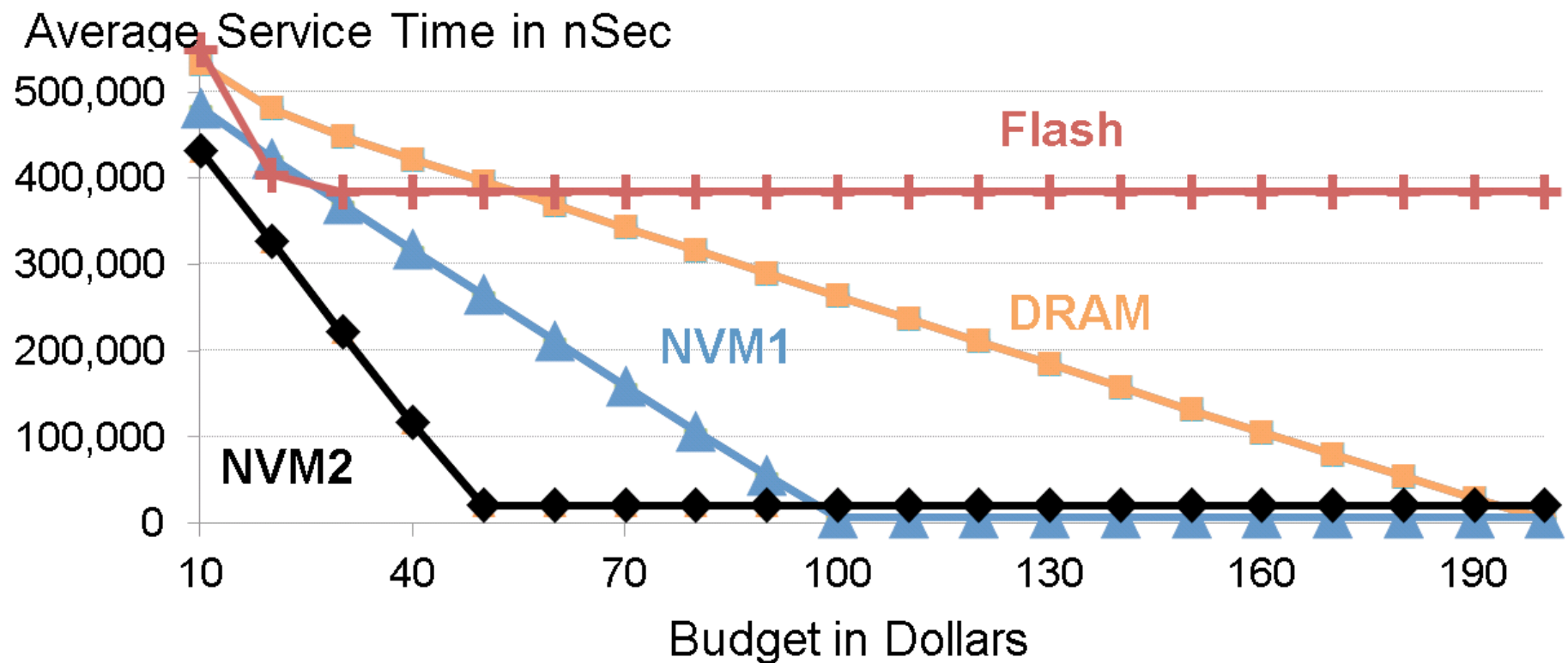
$$\sum_S x(p, S) = \text{size}(p)$$

$$\sum_{p, S \ni b} x(p, S) \leq \text{capacity}(b)$$

$$x(p, S) = 0, \text{size}(p)$$



CACHE CONFIGURATION





SUBSET ASSIGNMENT

HAVE $d \ll n$

sol to LP relaxation has few fractional assignments

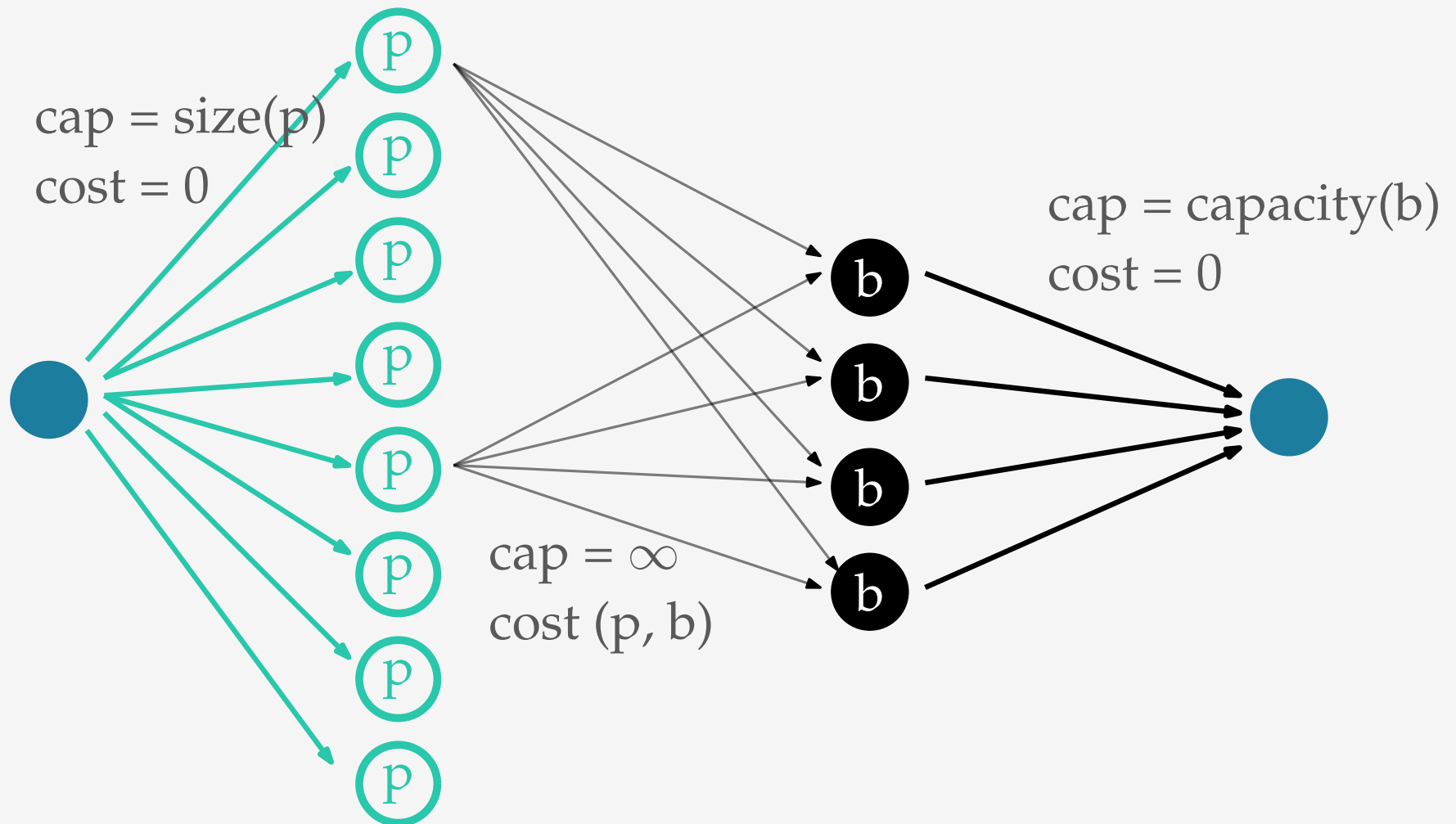
GOAL solve LP relaxation in $f(d) \text{ poly}(n)$

1. cycle canceling algorithm

2. simplex algorithm

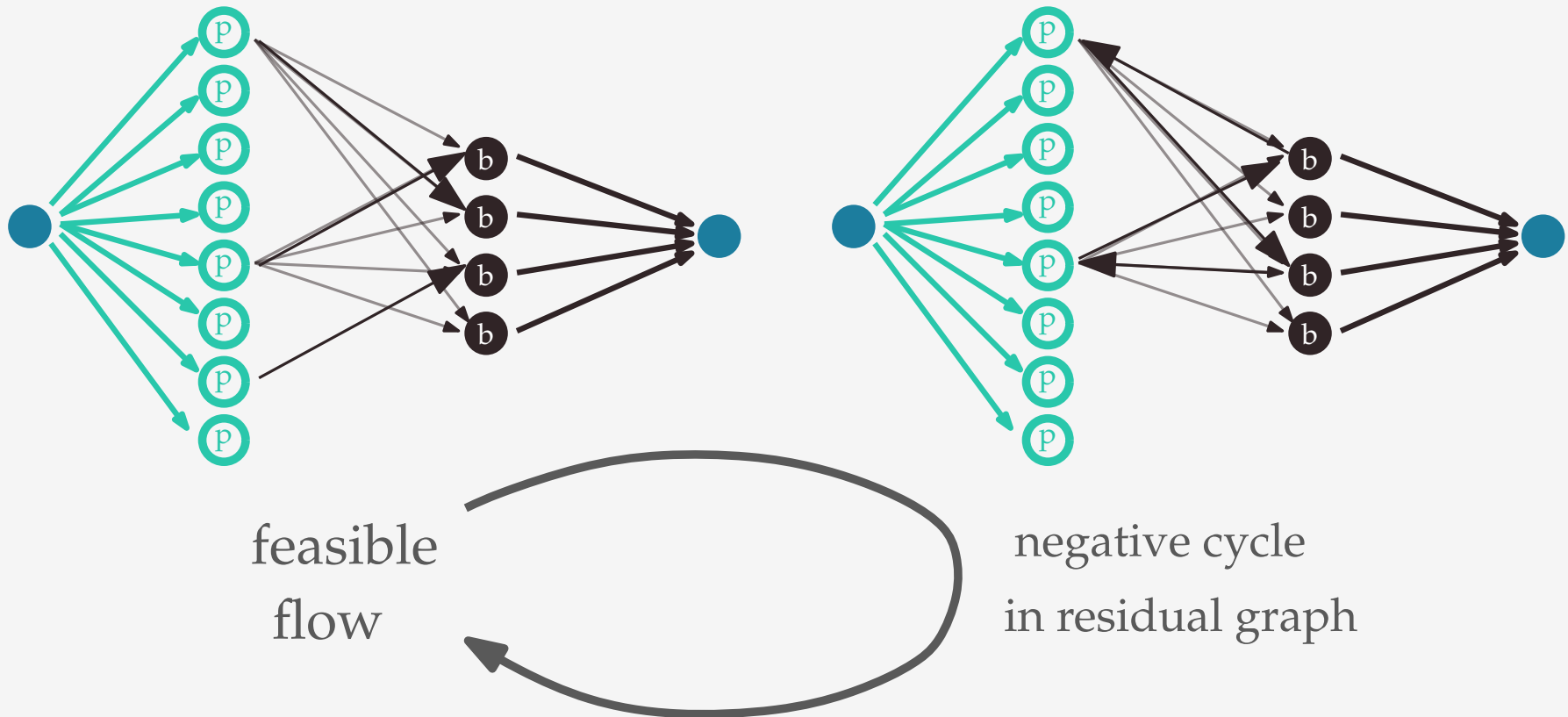


MIN COST FLOW





1. cycle canceling algorithm





“cycle” in subset assignment problem

augmentation $S_i \xrightarrow{p_i} T_i$

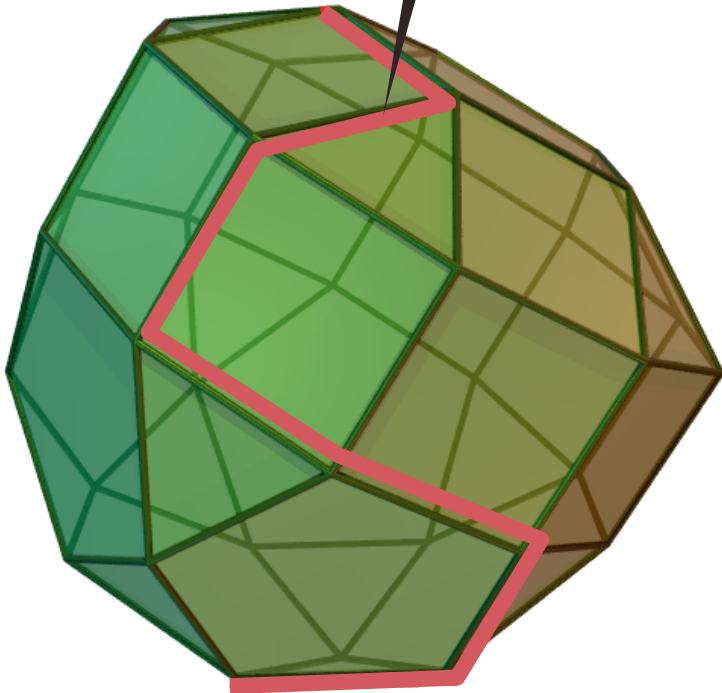
such that $\sum_i \alpha_i \overrightarrow{S_i T_i} = \vec{0}$

cost difference
(negative) $\sum_i \alpha_i (\text{cost}(p_i, T_i) - \text{cost}(p_i, S_i))$



2. simplex algorithm

basic feasible solution



BASIC FEASIBLE
ASSIGNMENT

$< 2d$ fractional assignments

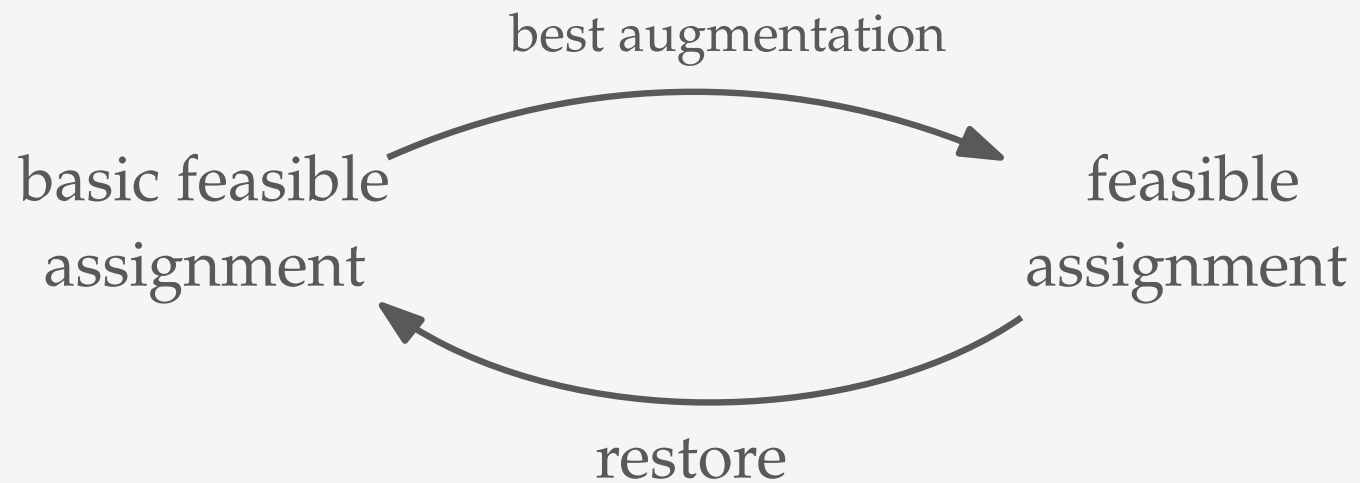
bound granularity of vars

$$x(p, S) = \frac{k}{l}$$

$< d^{d/2}$



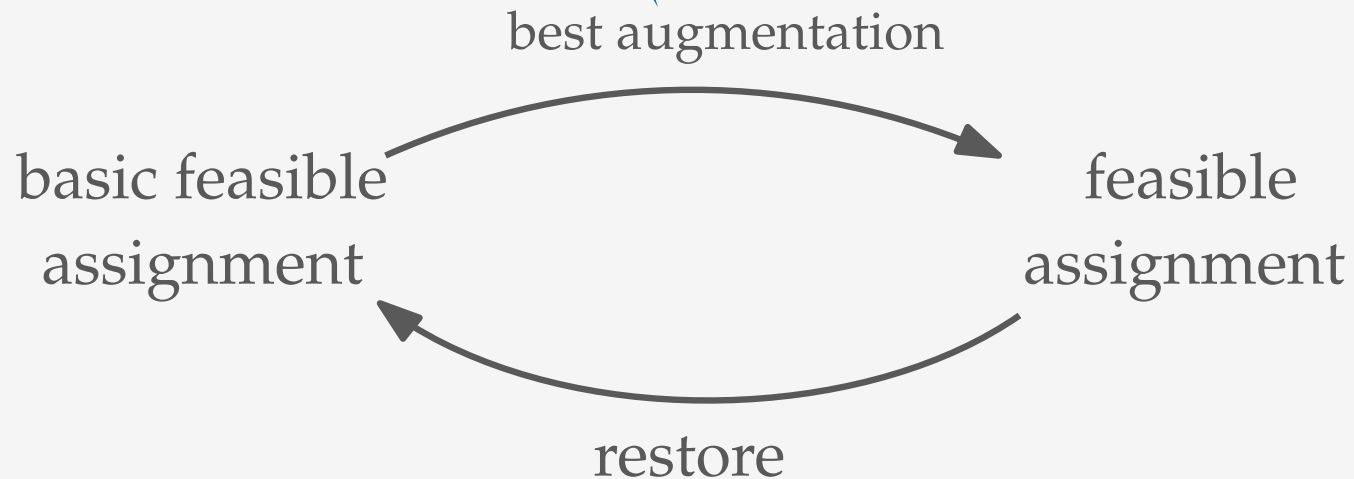
ALGORITHM





ALGORITHM

preprocessing $\sum_i \alpha_i \vec{S_i T_i} = \vec{0}$



time $O(\exp(d(d+1)) \text{poly}(d) \ n \log(n) \log(nC) \log(S))$

