

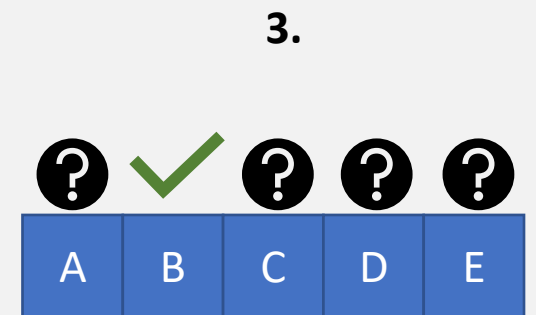
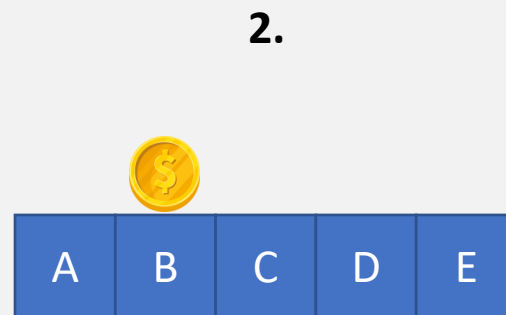
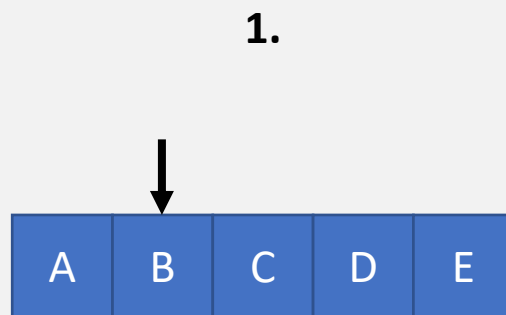
Practical Adversarial Multi-Armed Bandits with Sublinear Runtime

Kasper Overgaard Mortensen · Ama Bembua Bainson · Mathias Ravn Tversted · Kristoffer Strube Græm ·
Andrea Paudice · Renata Borovica-Gajic · Davide Mottin · Panagiotis Karras



Classical MABs - Simple online learning models

1. **Pick** - choose an arm to pull.
2. **Observe** - see the reward from that arm.
3. **Update** - adjust your estimates to improve future picks.



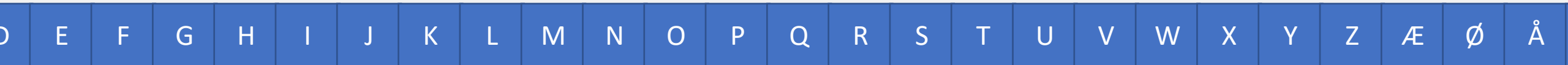
Problem: Runtime scales with number of arms

Many arms = **Slow** iteration runtime

- 1. Pick**
- 2. Observe**
- 3. Update**

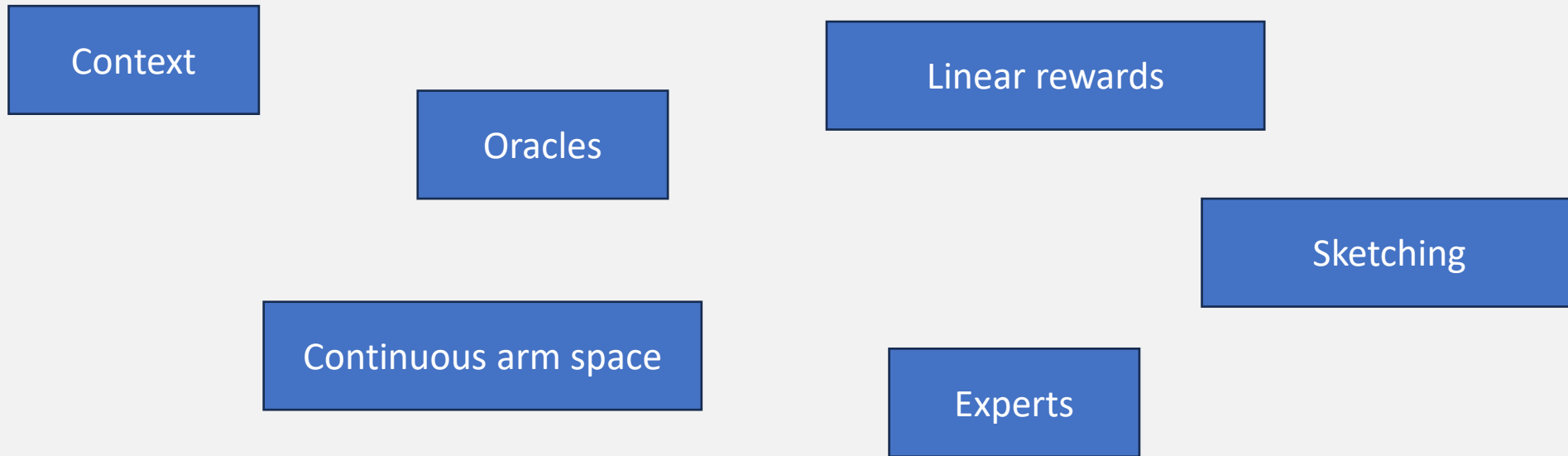
Runtime issues:

- 1) Maintain evolving probability distribution over k arms.
- 2) Perturb multiple weights every iteration.



The common approach: change setting!

Add **new assumptions** or utility to avoid the issue.



And many, **many more** keywords.

Should we care?

Classical bandits rarely work well in real settings anyhow...

I'm having difficulty finding interesting real-world k -Armed Bandit settings [...] without altering them to take context into account.

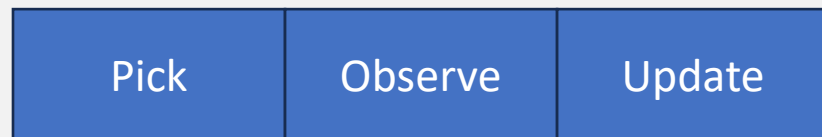
- John Langford

However, we need them as simple baselines to justify complex alterations.

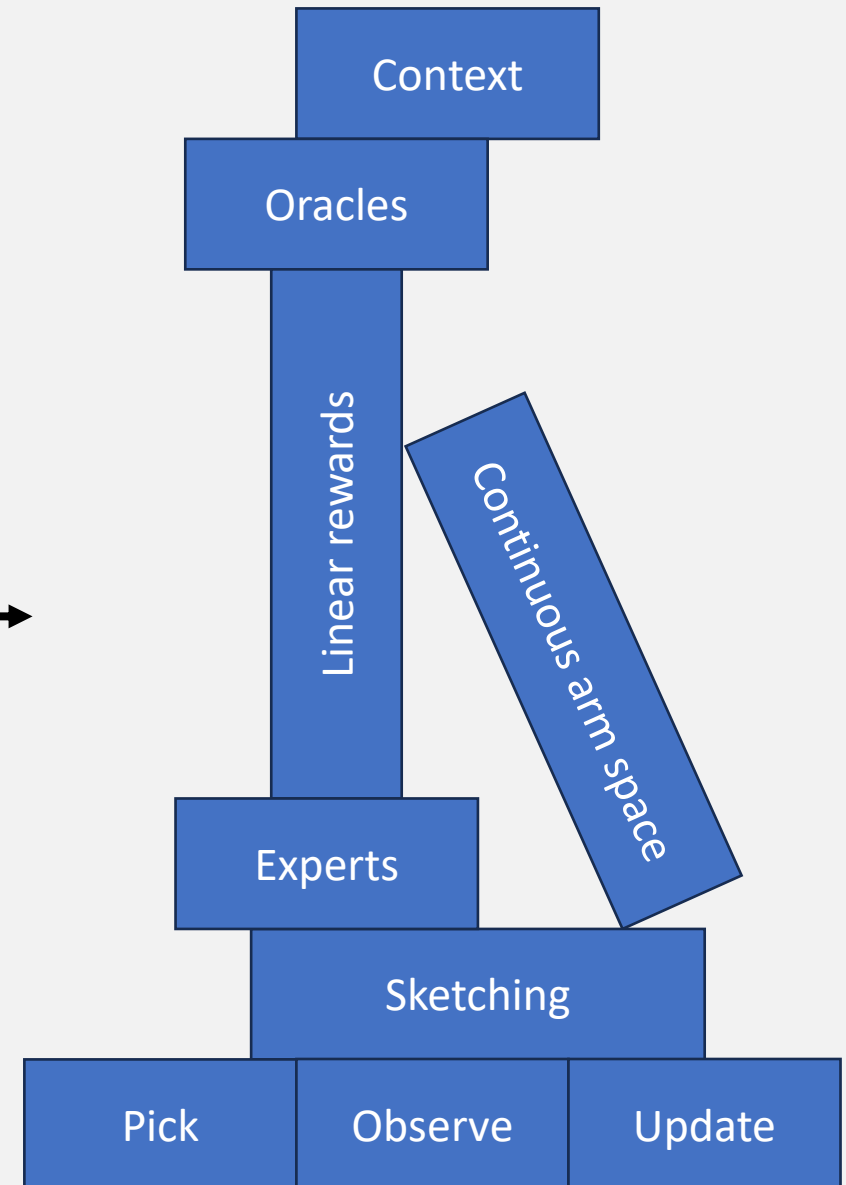
Complexity must add value.

Specialized bandits are great!

But they **should beat** the simple **classical methods**.



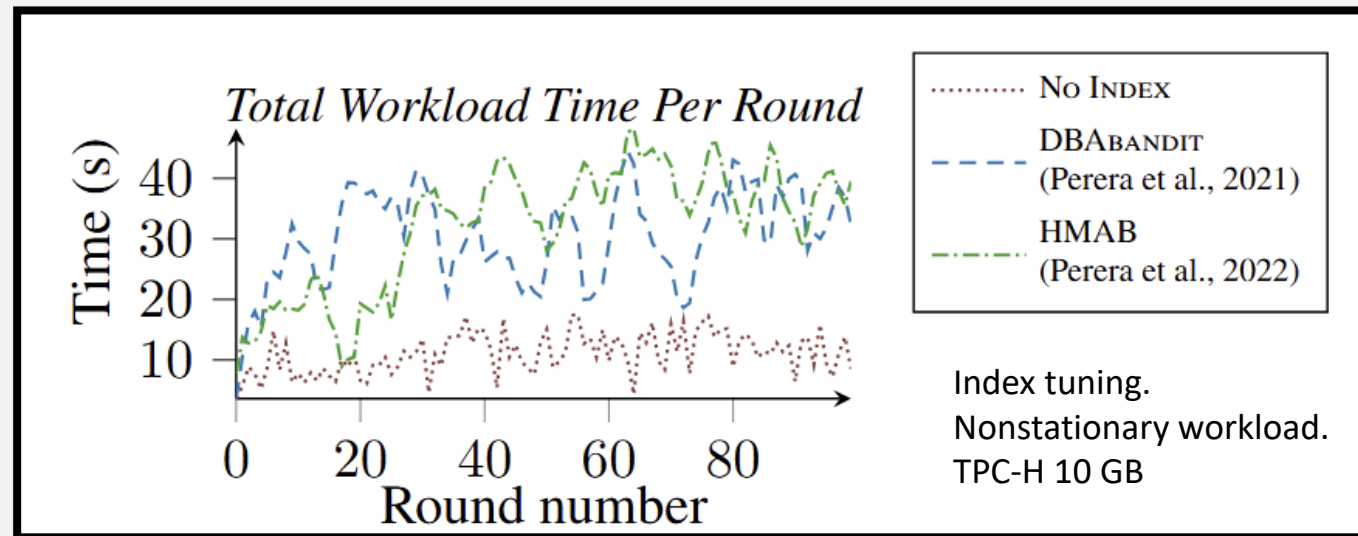
This should beat that.



System integration - moving beyond intended settings

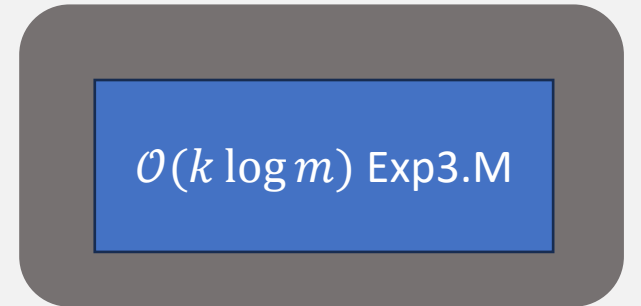
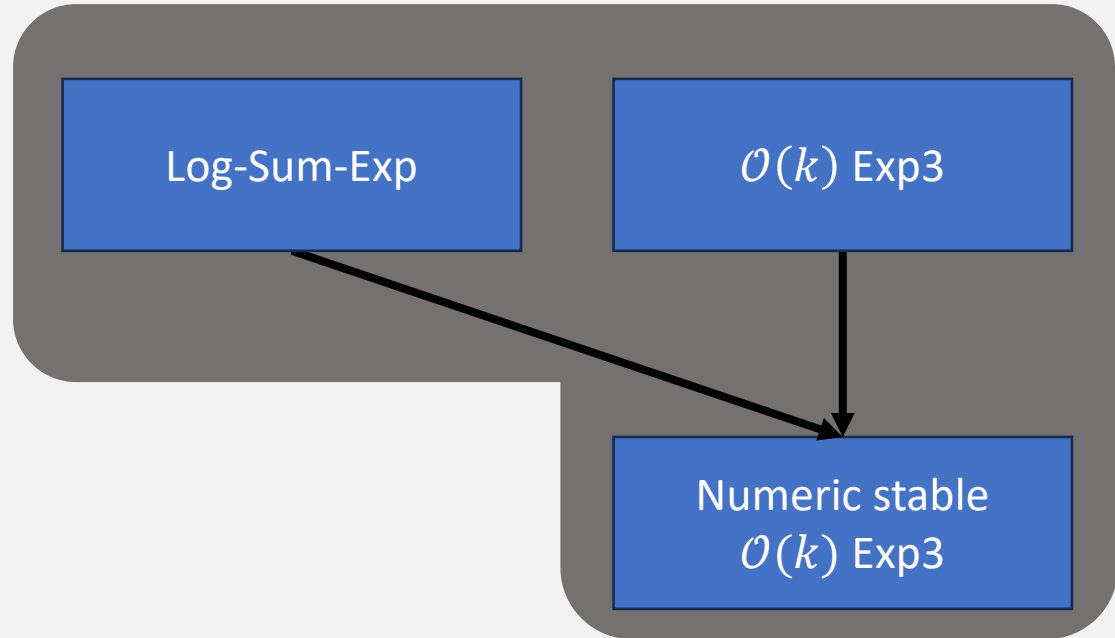
Growing interest for bandits in database systems.

Works well within the targeted assumptions but **struggles in extreme cases**.



A classic **adversarial bandit would be a nice baseline**... but they don't scale.

Exp3 is such a bandit.



Combinatorial setting.
(Select m arms)

But **runtime is linear** in the number of arms k .

Let's make a Exp3 scale sublinear to arms!

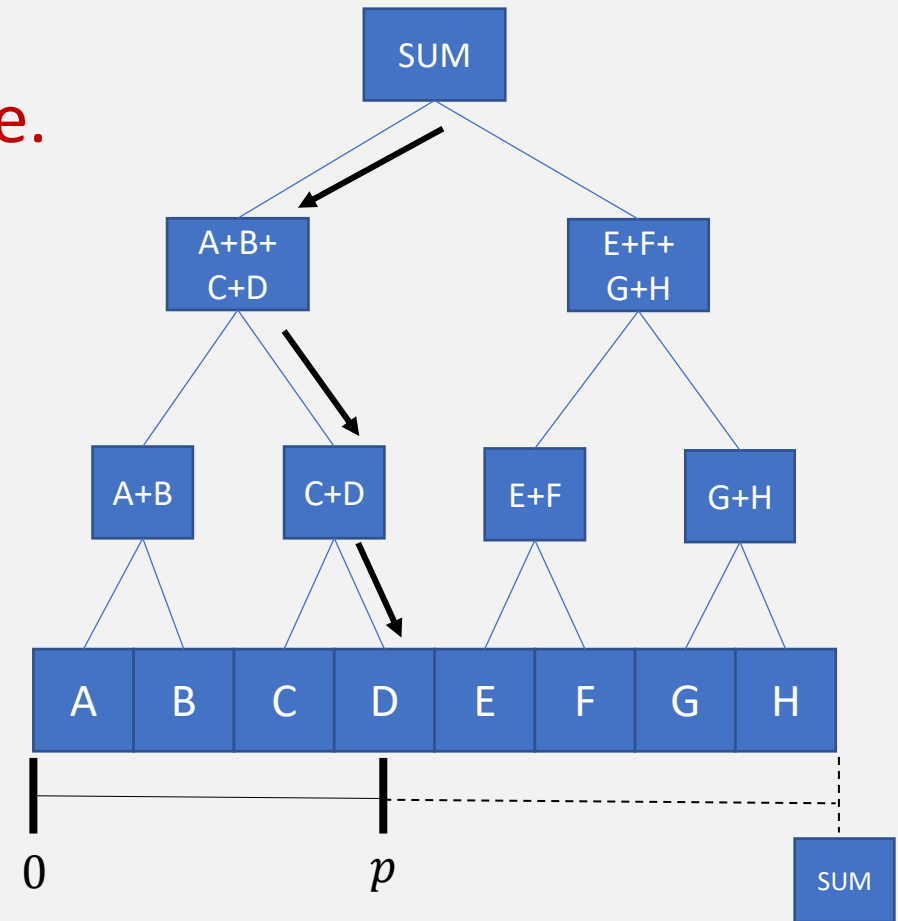
Solved 10 years ago.



Sum-heap sampling solves the linear overhead

Sample from a distribution **that may evolve**.

$$p = \text{Uniform}(0,1) \cdot \text{SUM}$$

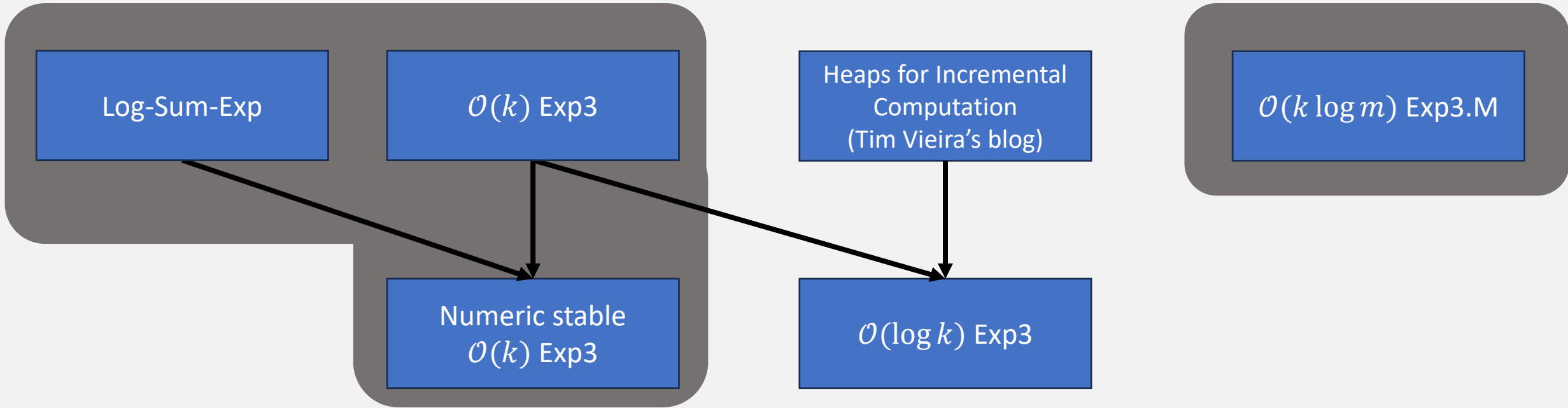


Method	Sample	Update	Init
heap	$O(\log n)$	$O(\log n)$	$O(n)$

- EXP3 (multi-armed bandit algorithm) is an excellent example of an algorithm that samples and modifies a single weight in the distribution.

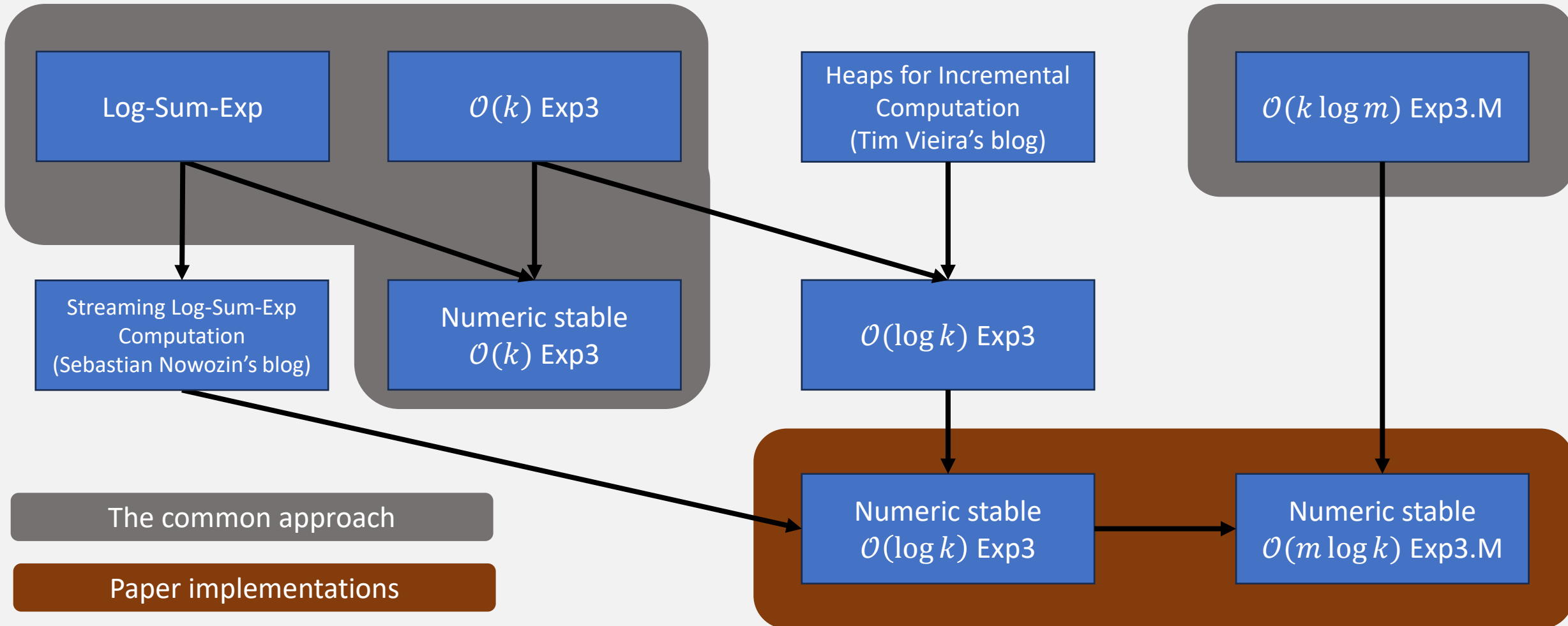
Blog link: <https://timvieira.github.io/blog/heaps-for-incremental-computation/>

Forgotten folklore?



The common approach

Covering considerations for large time horizons.



$\mathcal{O}(\log k)$ and $\mathcal{O}(m \log k)$ is good. *But* can we go faster?

Design for speed.

Design for Speed

Design principle: Best of both worlds.

From Follow-the-leader family: Always pick best perceived option.

From Exp3 family: Update a single weight.

Approach: *Never* look at anything but the best arm.

If deemed good: Keep selecting it.

If deemed bad: Demote it.

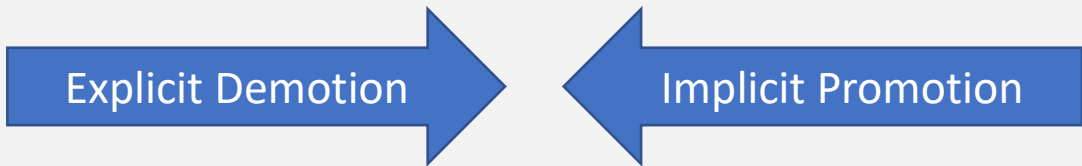
Forgo quality guarantees to obtain speed.

Queue-Behind-the-Leader (QBL)



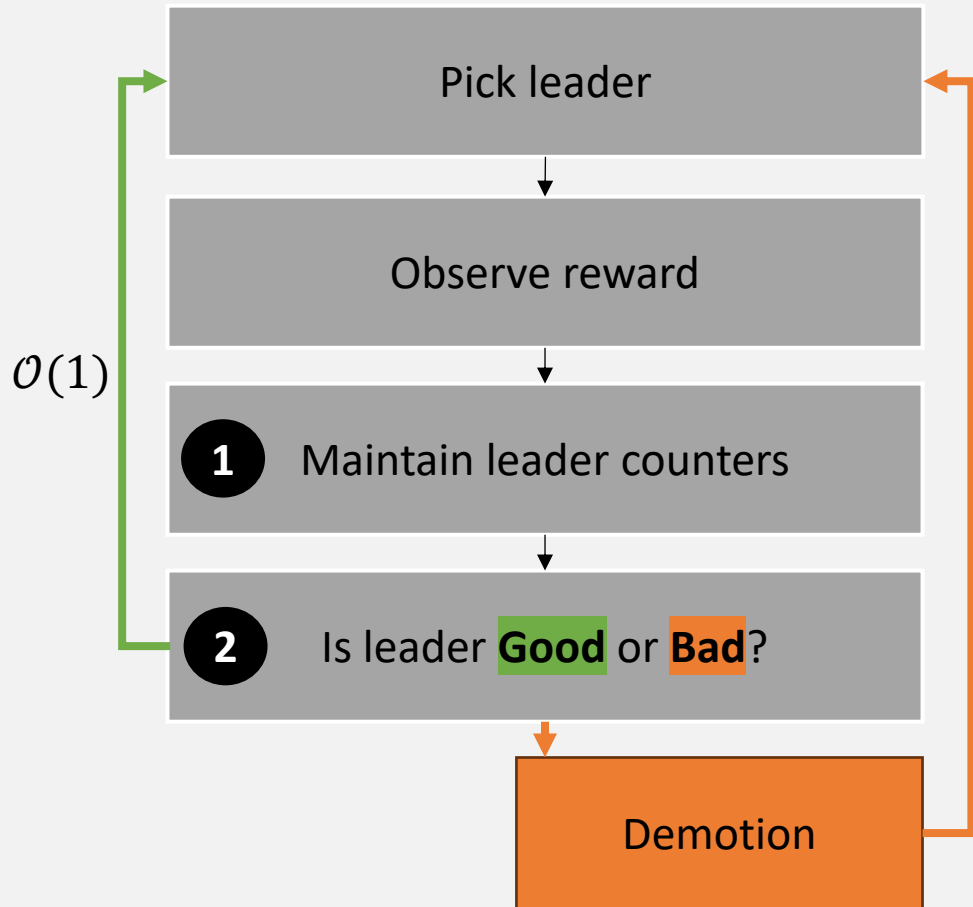
←
max: k Capped Priority min: 0

Demote bad leaders. A two-way pressure.

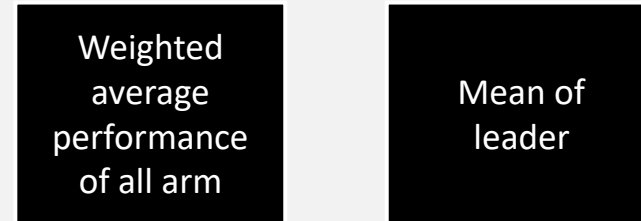


A historical queue of leaders

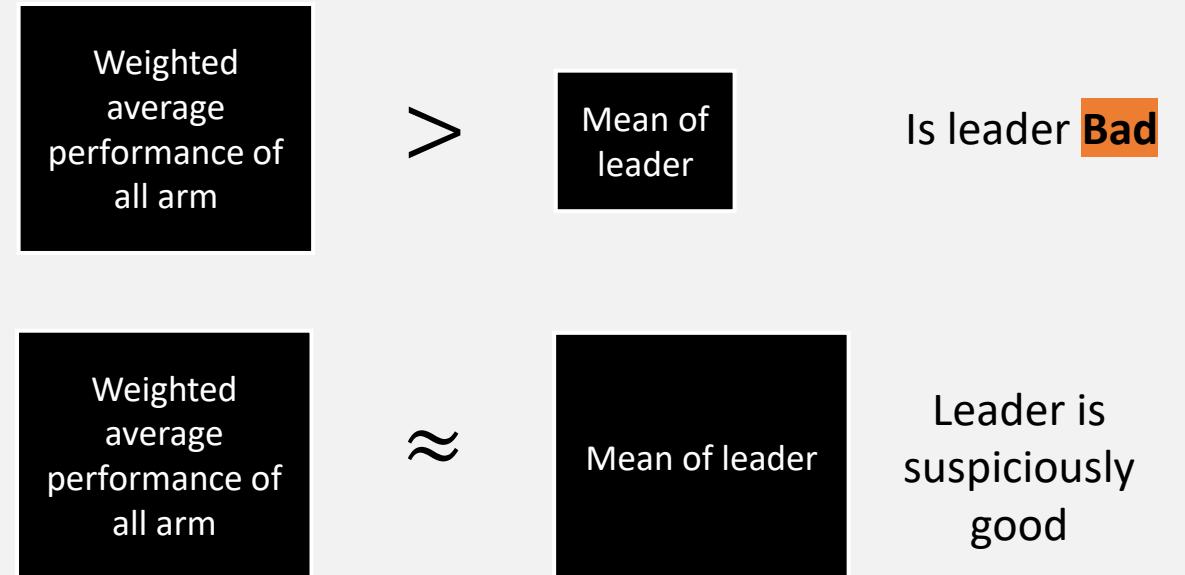
Queue-Behind-the-Leader (QBL)



1 Constant-time maintenance:

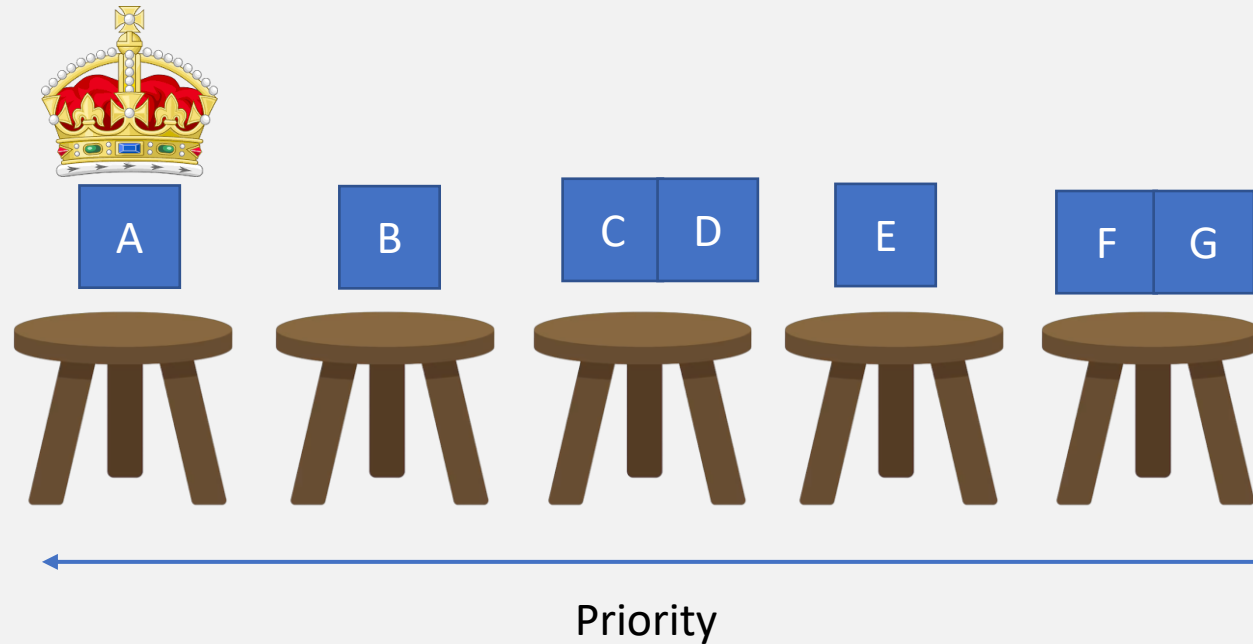


2



(Randomize decision to break determinism)

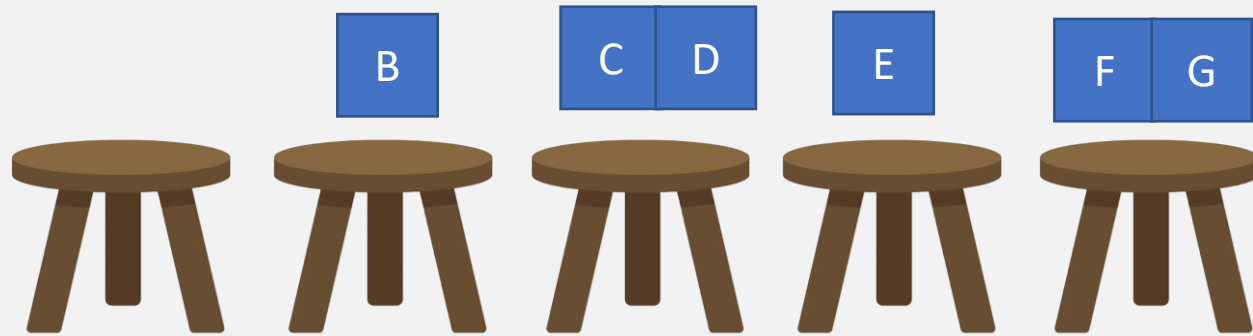
QBL – Implicit promotion



QBL – Implicit promotion



Leader is demoted

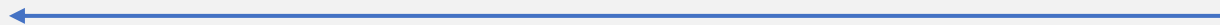
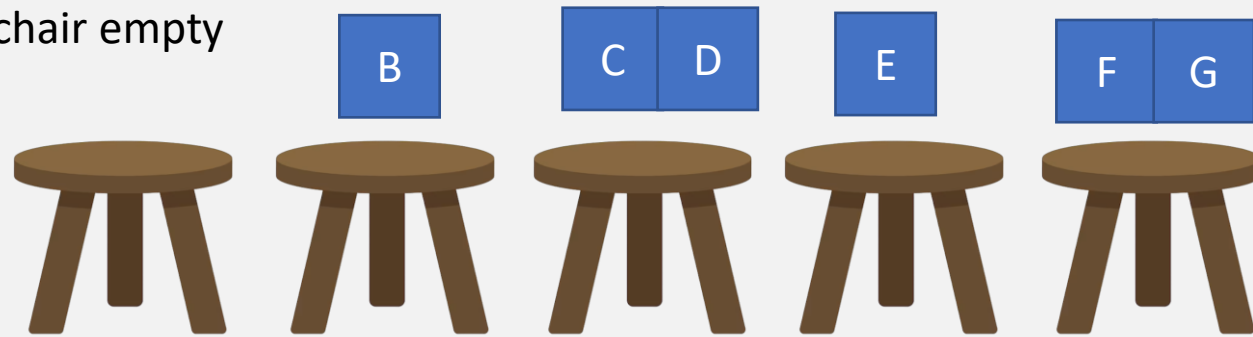


Priority

QBL – Implicit promotion



Top chair empty

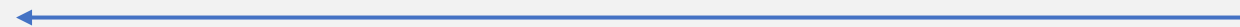
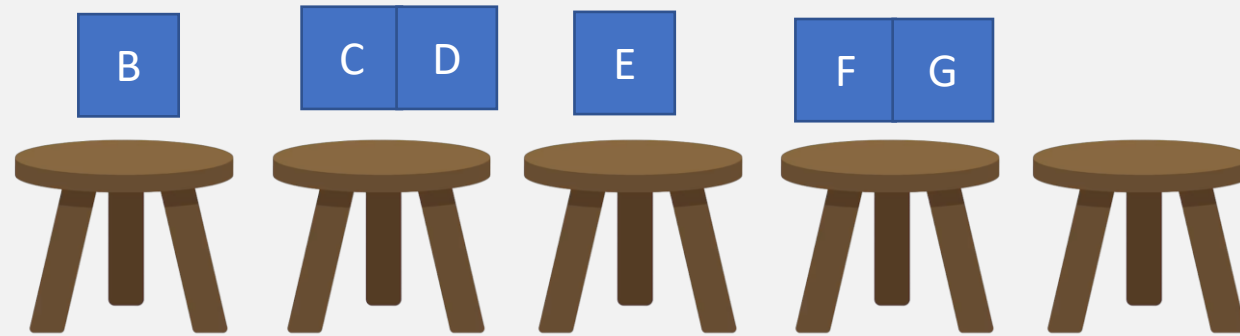
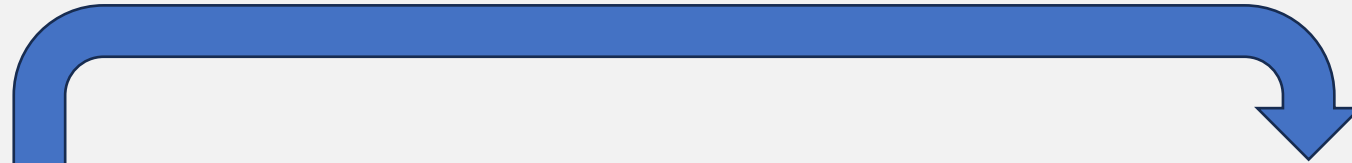


Priority

QBL – Implicit promotion



Move chair

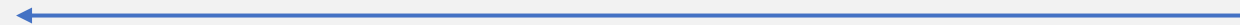
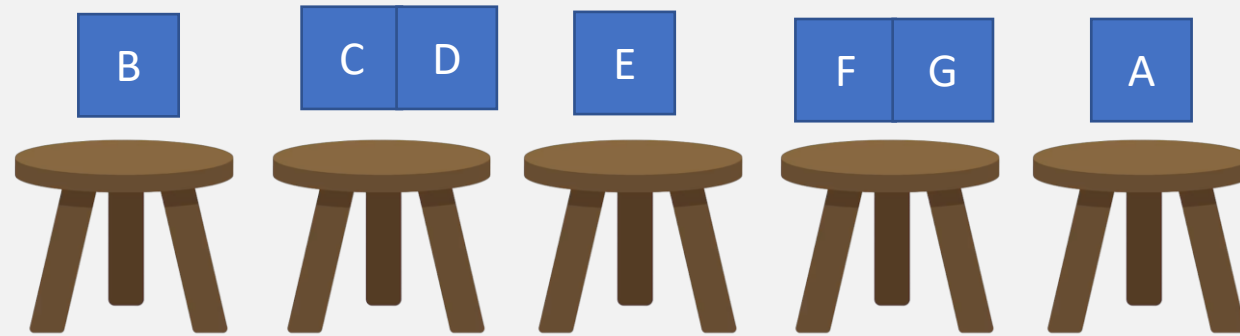
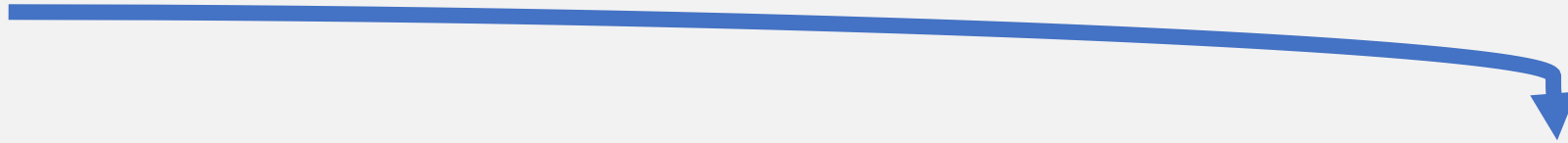


Priority

QBL – Implicit promotion



Move old leader to its new position

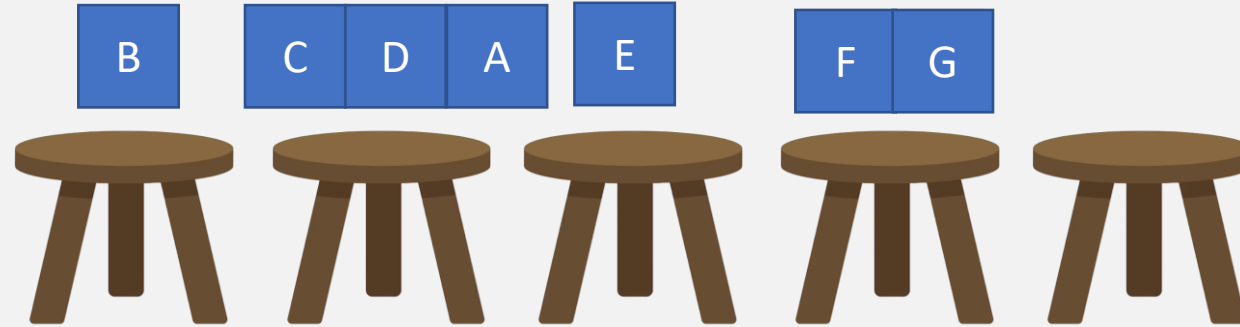
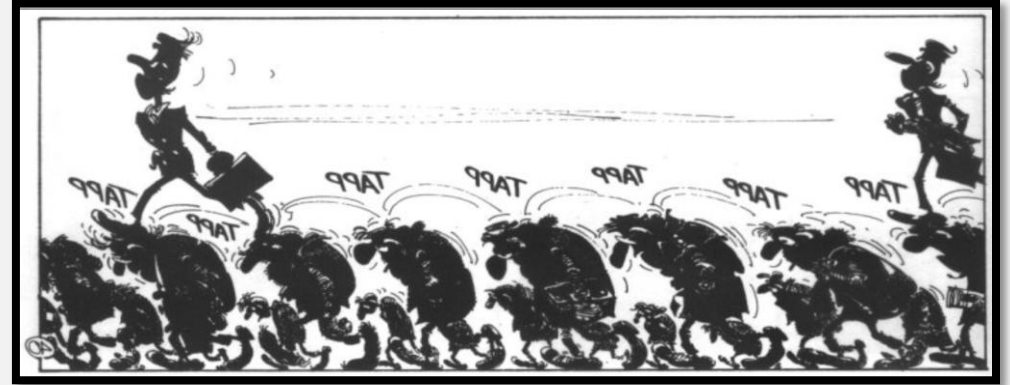


Priority

QBL – Implicit promotion



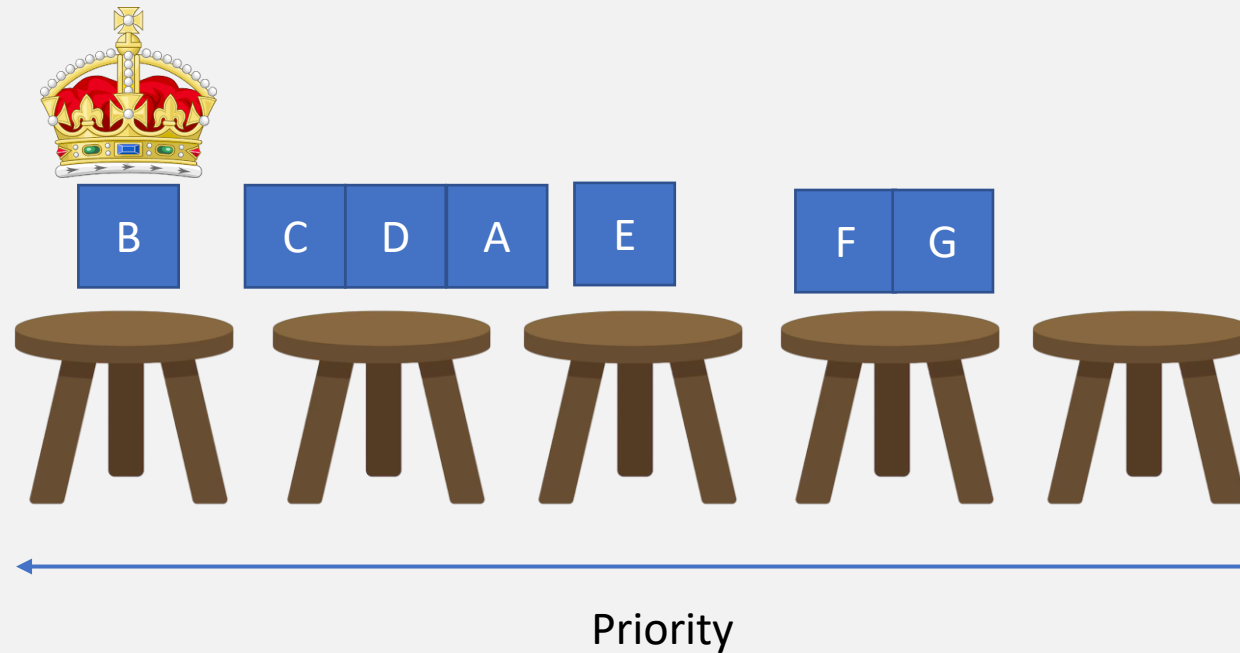
Further up the longer it remained the leader



Priority

QBL – Implicit promotion

Elect new leader at
the new top chair



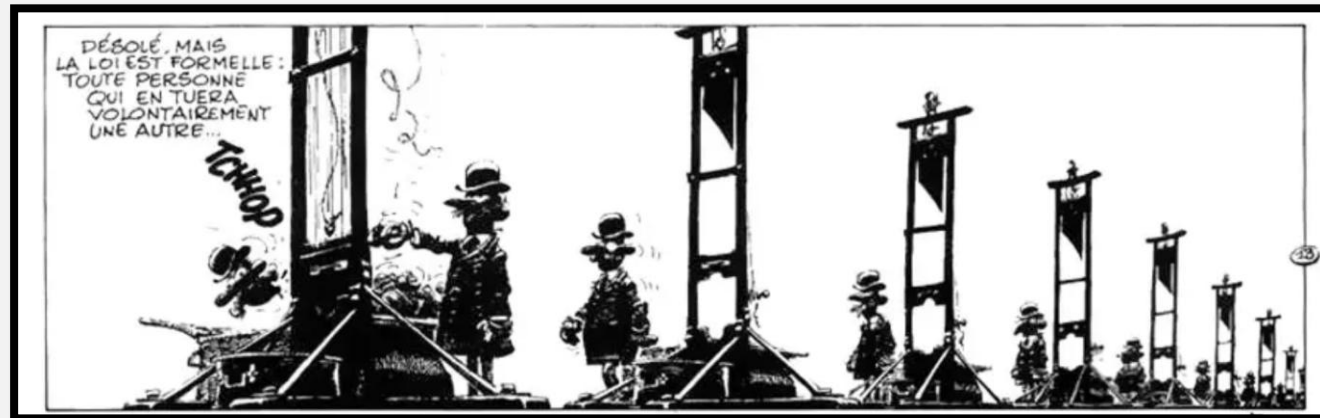
QBL - Overview

Exploration through demotions. Best case **skips update phases**.

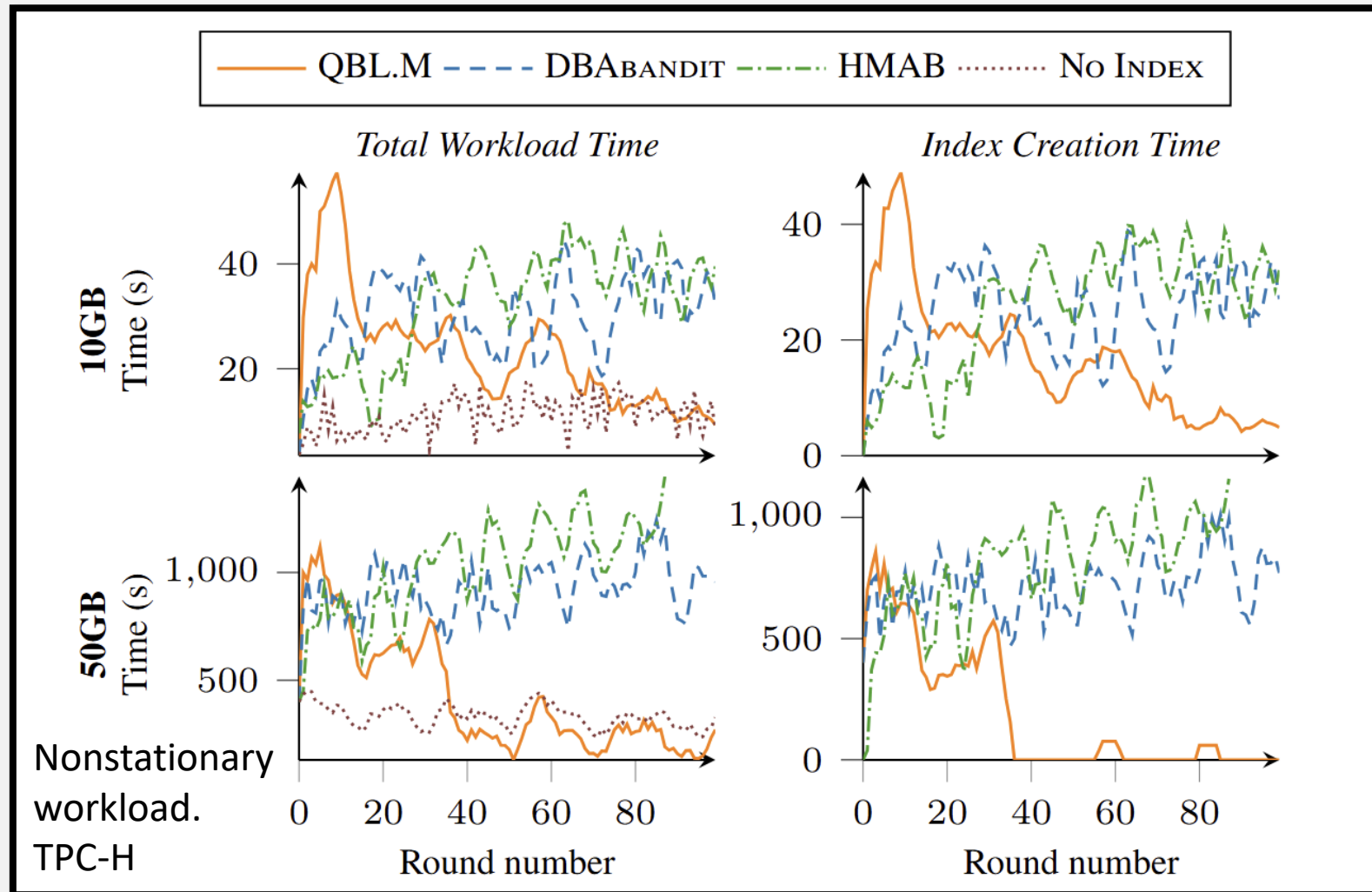
Good leader: $\mathcal{O}(1)$ **Demotion:** $\mathcal{O}(\log k)$

Combinatorial setting with dynamic m : $\mathcal{O}(m \log k)$

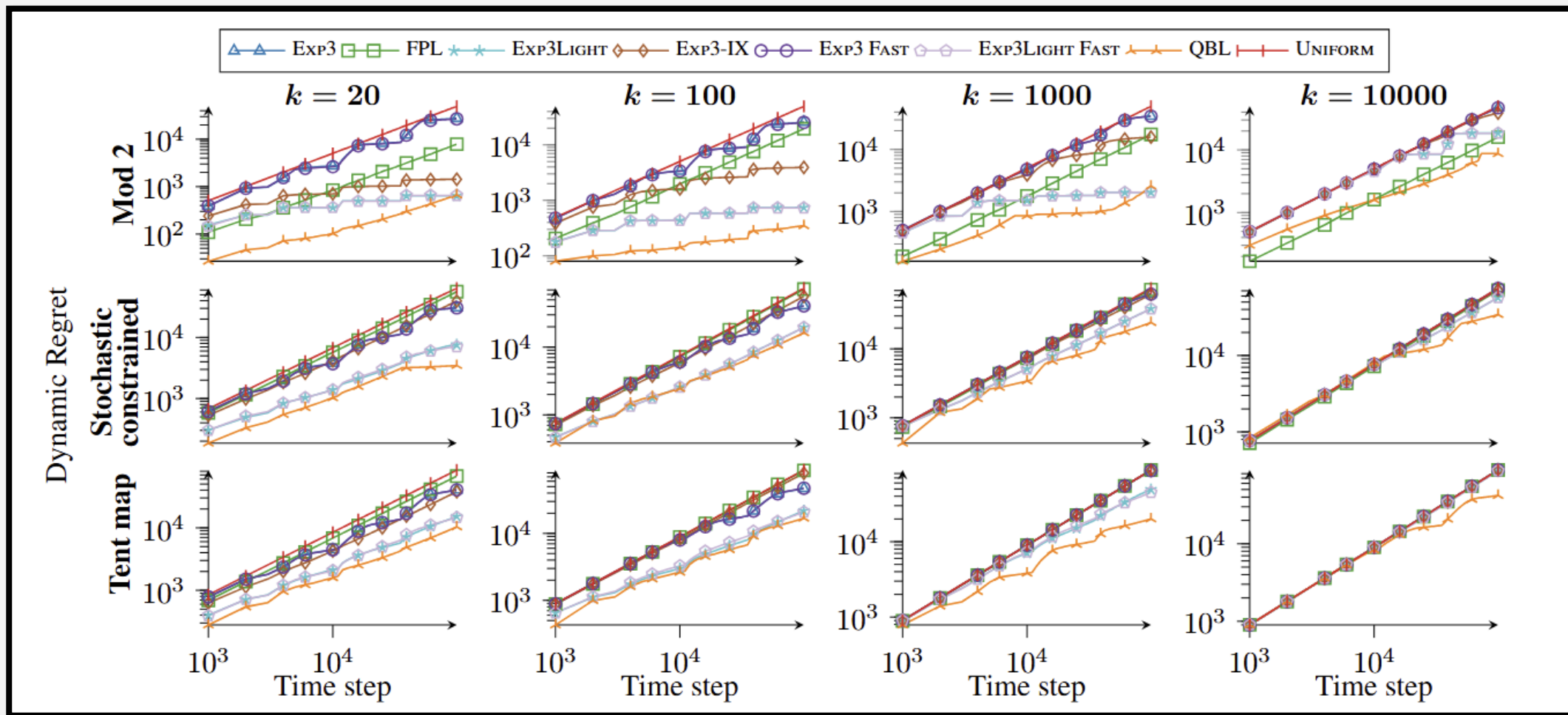
Forgo quality guarantees to achieve fast runtime.



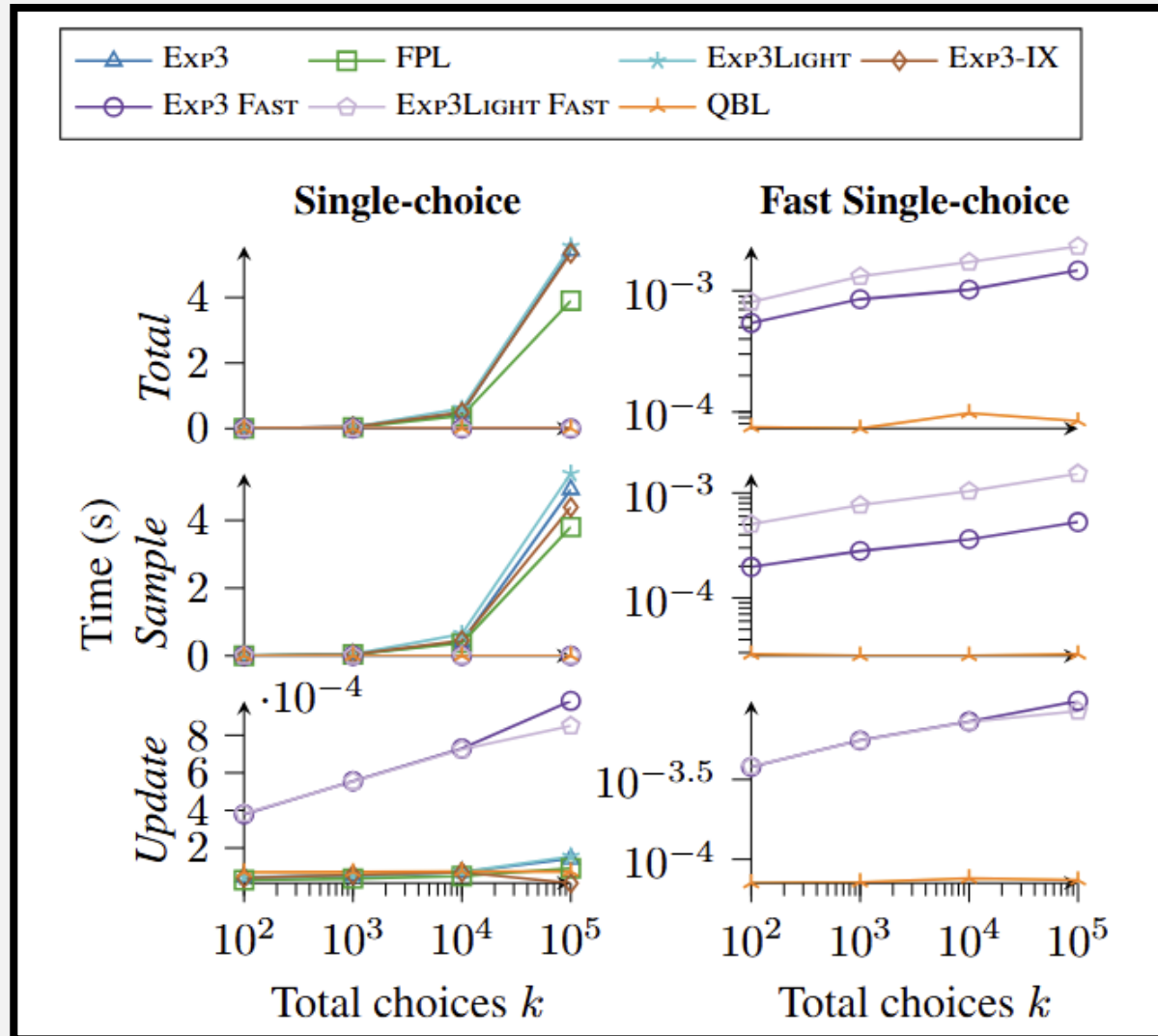
Returning to index tuning in extreme workloads



Quality is still promising.



And of course, runtime



Some final thoughts...

Reminder: Scalable and simple baselines are needed.

QBL is designed to be cynical. **Worse quality in stable settings.**

The hunt for quality guarantees:

Recent work (AdSwitch) utilize similar **Good/Bad** classification.

Mechanics in QBL can be swapped out.

Thanks!

Checkout Tim Vieira's blog:

<https://timvieira.github.io/blog/heaps-for-incremental-computation/>

Checkout Sebastian Nowozin's blog:

<https://www.nowozin.net/sebastian/blog/streaming-log-sum-exp-computation.html>

Contact me: km@cs.au.dk

Comics strips from: Idées Noires - André Franquin

Inverse Sampling

1. Uniform probing of sum.
2. Select option representing the probe value.

Probe →

Green option
selected.



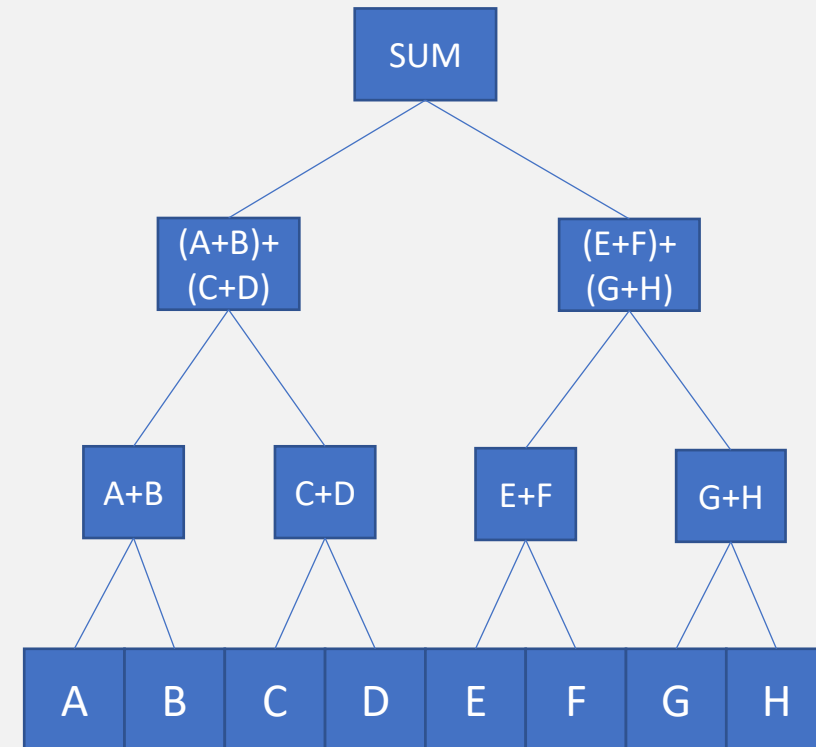
Sum of weights

Sum heap – incremental computation

Store weights in a sum heap.

Update weights by fixing intermediate sums. $O(\log k)$.

Costly, so why do it?



Sum heap – Sampling edition

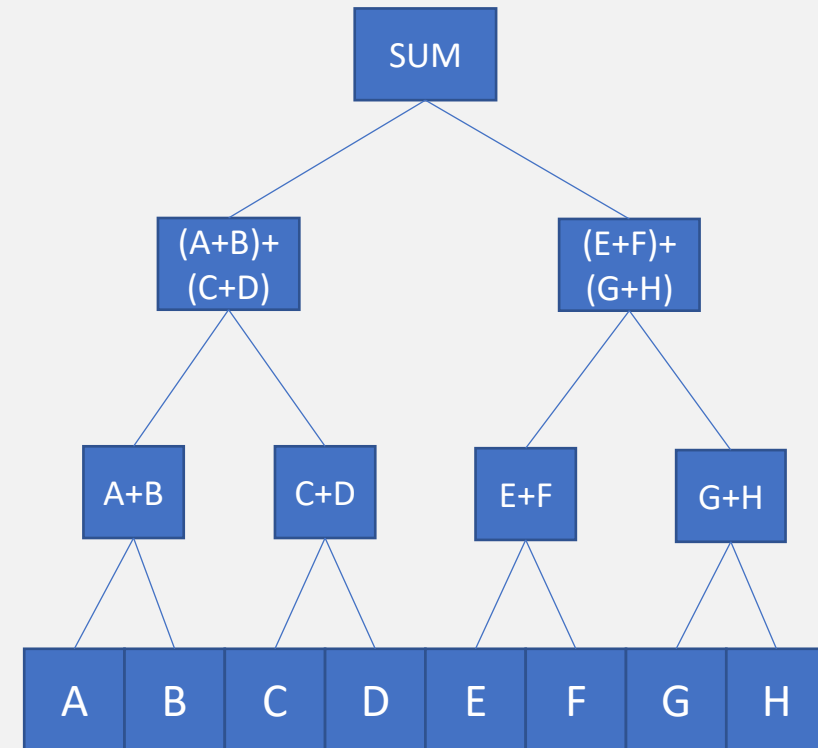
Useable in sampling!

Corresponds to an unnormalized CDF of weights.

Supports **that it may evolve.**

$$p = \text{Uniform}(0,1) \cdot \text{SUM}$$

Find largest CDF value smaller than p . $O(\log k)$.



Integrates nicely with log scale

$$LSE(\ln(w_1), \ln(w_2), \dots, \ln(w_k)) = \ln(w_1 + w_2 + \dots + w_k)$$

$$\ln(w_x + w_y) = \ln(w_y) + \ln(\exp(\ln(w_x) - \ln(w_y)) + 1)$$

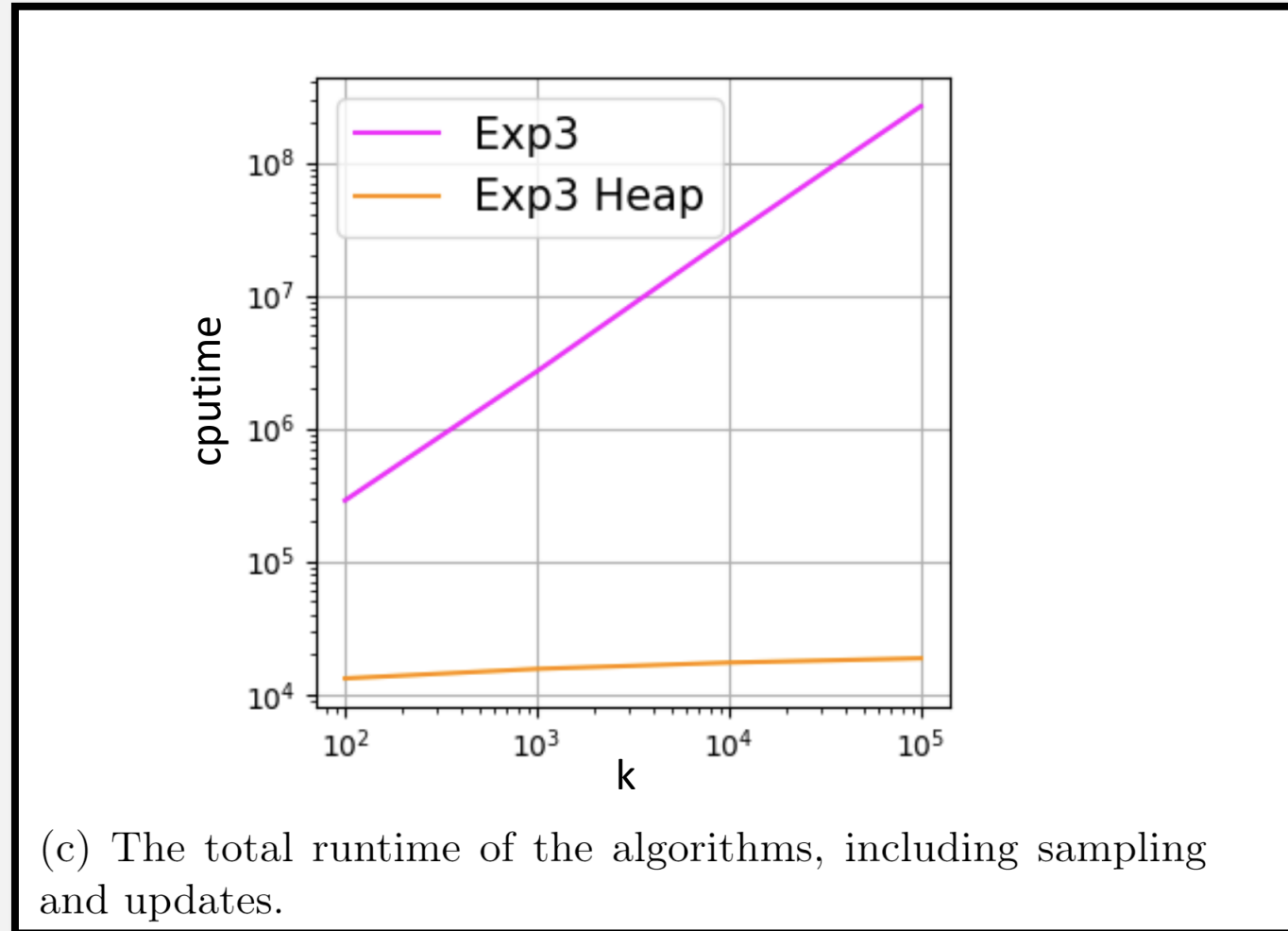


Heap intermediate

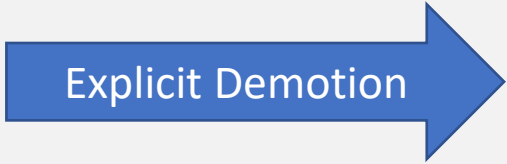


*Never access any
weight outside of
logspace*

Exp3 with Sum heap




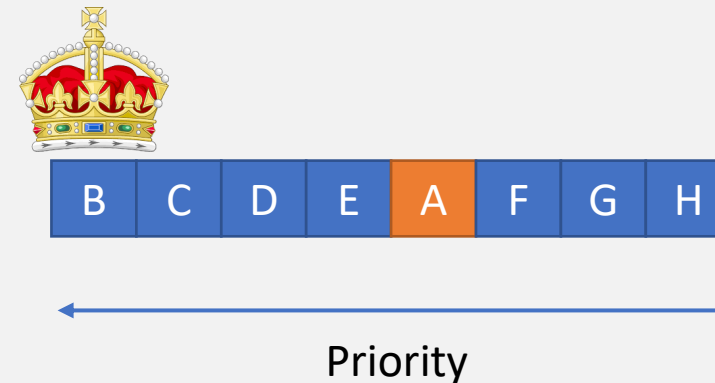
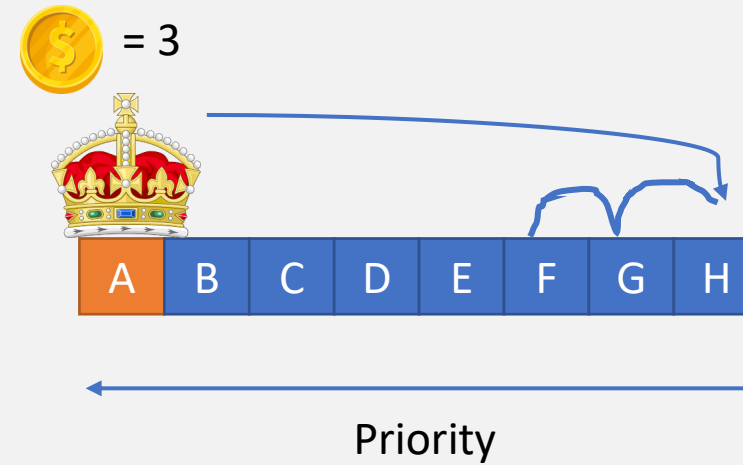
QBL – Explicit Demotion



Caps priority to k .
Leaders accumulate wealth. Maximum $k - 1$.



Bad leaders are demoted based on  and lose their wealth.



QBL – Explicit Demotion

Explicit Demotion

Problem: Exploration not guaranteed



Solution: Implicit promotion.

QBL – Implicit promotion

Shift the allowed priority interval

Remember: Demote by at least 1



Option:



Priority of option:



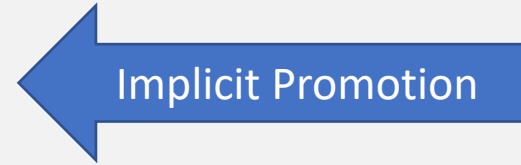
max: 8

Priority

min: $8 - k - 1$

Max is priority of current leader

Min is set based on max



max: 6

Priority

min: $6 - k - 1$

