vFair: Latency-Aware Fair Storage Scheduling via Per-IO **Cost-Based Differentiation** Hui Lu⁺, Brendan Saltaformaggio⁺, Ramana Kompella⁺[‡], Dongyan Xu⁺

[†]Department of Computer Science, Purdue University, [‡]Google Inc.



An IO scheduler, providing fairness and isolation while maintaining high resource utilization,



Per-IO Cost Allocation Model

Proportional Service Time Allocation





- **BOPS** does not take blender effect into account.
- **BAPS** ensures high efficiency with best-effort fairness guarantees.
- **BSPS** ensures the best fairness guarantees.

Key Idea – to use saturation throughput P_i

 $T_i = f(\frac{w_i}{\sum_k w_k} \cdot P_i)$

- Each VM's workload can be modeled as an IO pattern.
- Given a particular IO pattern, it corresponds a saturation throughput P_i .
- P_i refers to the theoretical IO throughput that a VM could receive in isolation.
- P_i is directly related to T -- P_i refers to full utilization which is T.



1 2 3 4 5 6 7 8 P_1 = 8000, Sequential access

P_2 = 1200, Random access







 $P_{1+2} \approx 1200$, Random access

Fig. The total saturation performance becomes **P1+2 = 1200 (the worst case).**

Due to **IO-blender effect**, using static P_i to estimate storage capacity would become **imprecise** over time.

Scheduling Strategy	Fairness Guarantee	Work-conserving?	High Utilization?
BOPS(vFair)	Good	Yes	Yes
BAPS(vFair)	Better	Yes	Yes
BSPS(vFair)	Best	Yes	No
PS(mClock)	Good	Yes	Yes
CFQ(Linux)	Bad	No	No
Time-Quanta	Dest		
based scheduler	Best	INO	INO

Realization?

- First level: Credit-based IO Rate Controller
- Second level: Faire queuing scheduling SFQ(D)



Micro-Benchmark









Fig. Saturated IOPS performance vs. IO request sizes for a single SSD

Fig. VM1 and VM2 should each obtain 50% of the shared storage throughput.

3

possible to accurately

measure the service time.

Model Approximation – to obtain P_i

 $1/P_i = \alpha \cdot \beta / P_{seq_rd} + \alpha \cdot (1 - \beta) / P_{rand_rd}$ + $(1 - \alpha) \cdot \gamma / P_{seg wr} + (1 - \alpha) \cdot (1 - \gamma) / P_{rand wr}$

Posterior Knowledge: VM's runtime IO access pattern -- α , β and γ . **Prior knowledge**: Basic IO pattern's P_i , constructed offline -- $P_{seq rd}$, $P_{rand rd}$, $P_{seq wr}$ and $P_{rand wr}$.