

# SoftSKU: Optimizing Server Architectures for Microservice Diversity @Scale

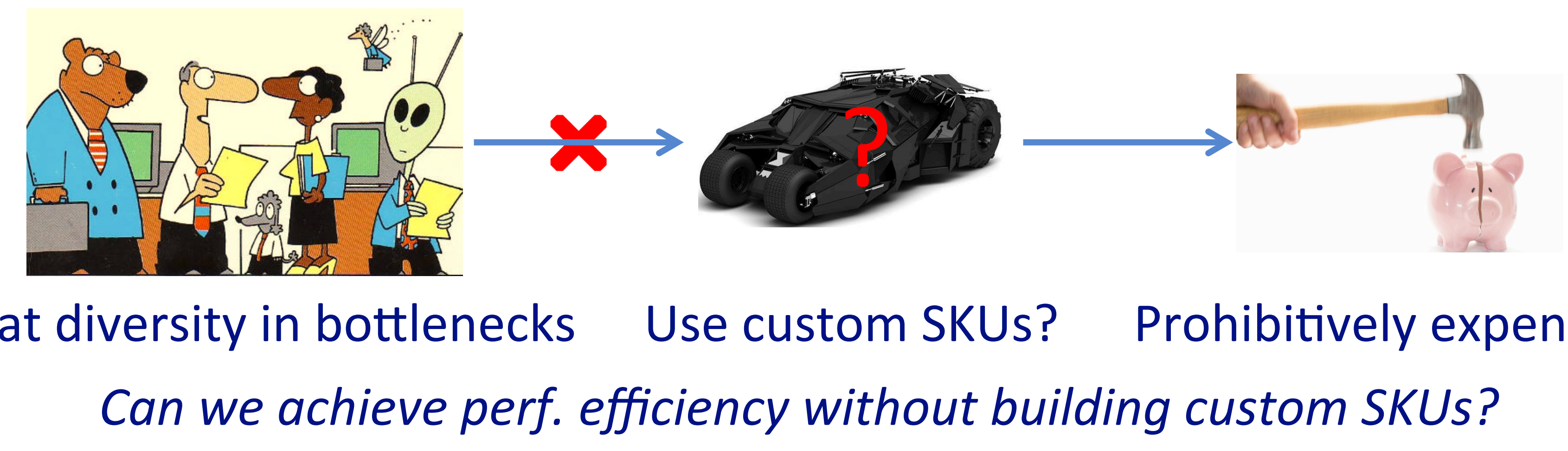
Akshitha Sriraman<sup>\*</sup>, Abhishek Dhanotia<sup>^</sup>, Thomas F. Wenisch<sup>\*</sup>  
 University of Michigan<sup>\*</sup>, Facebook<sup>^</sup>

## Rapid Increase in Modern Web Services

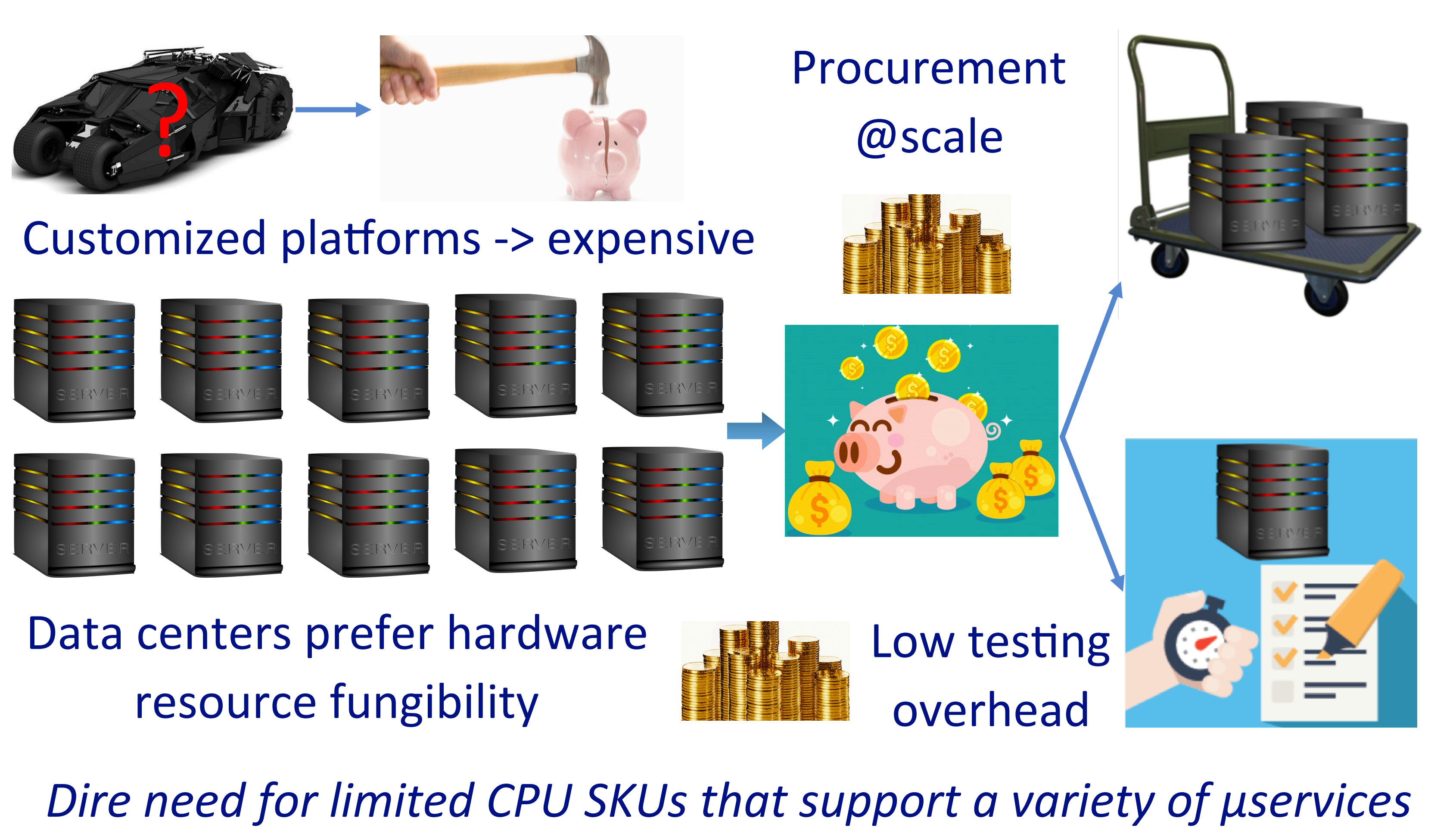


## Facebook $\mu$ Services' Characterization

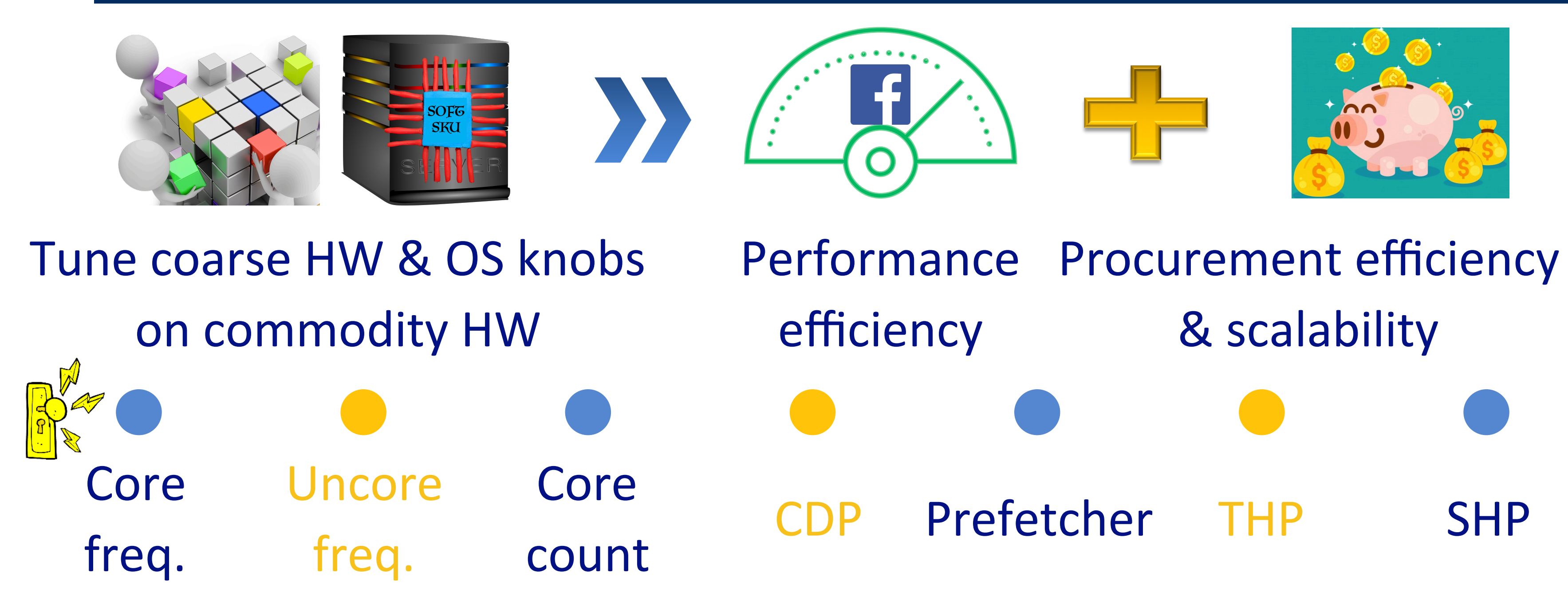
$\mu$ Service	Throughput (QPS)	Response latency	Pathlength
Web	$O(100)$	$O(ms)$	$O(10^6)$
Feed1	$O(1000)$	$O(ms)$	$O(10^9)$
Feed2	$O(10)$	$O(s)$	$O(10^9)$
Ads1	$O(10)$	$O(ms)$	$O(10^9)$
Ads2	$O(100)$	$O(ms)$	$O(10^9)$
Cache1	$O(100K)$	$O(\mu s)$	$O(10^3)$
Cache2	$O(100K)$	$O(\mu s)$	$O(10^3)$



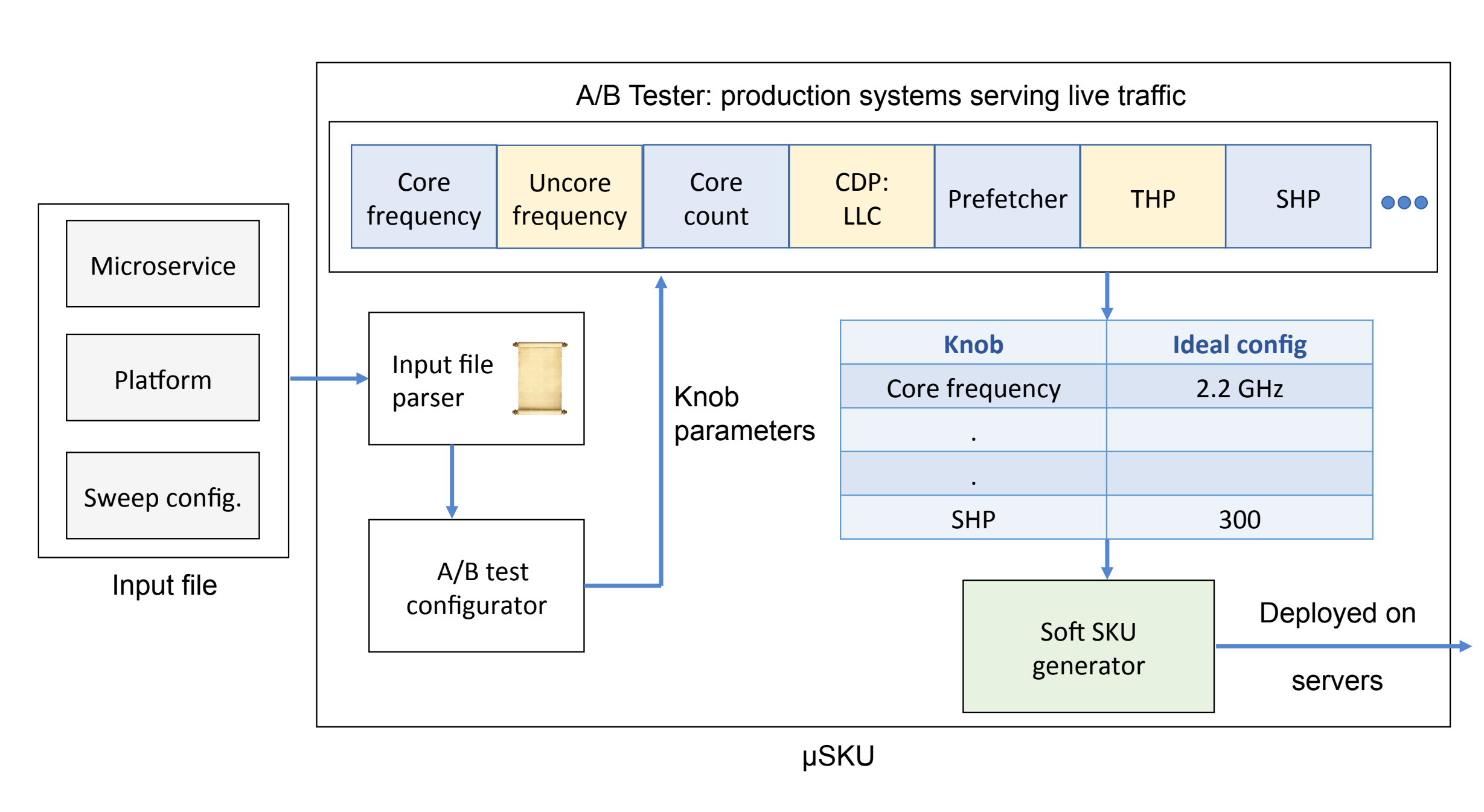
## Are Custom Platforms Always Needed?



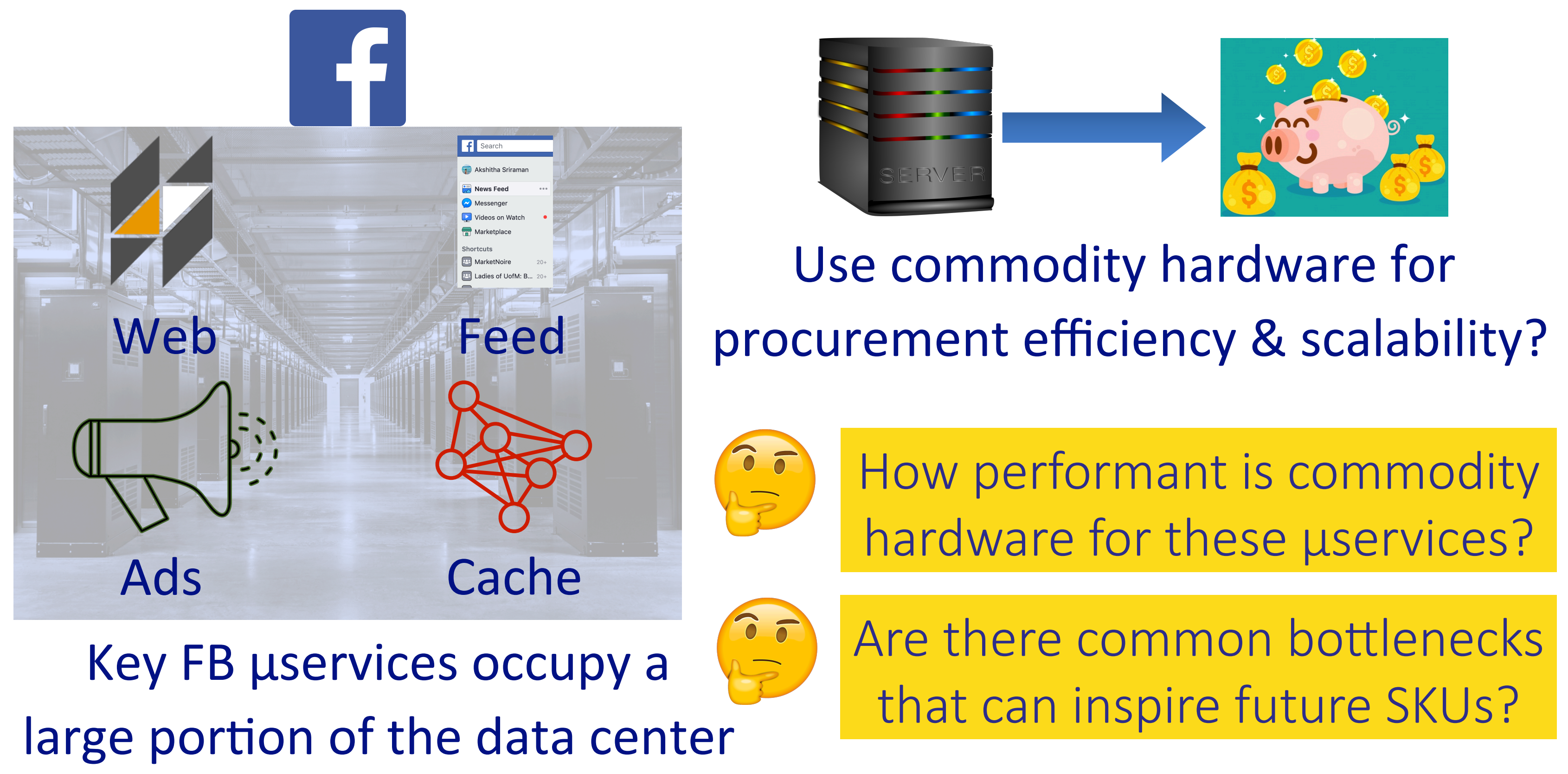
## "Soft" SKUs: Best of Both Worlds



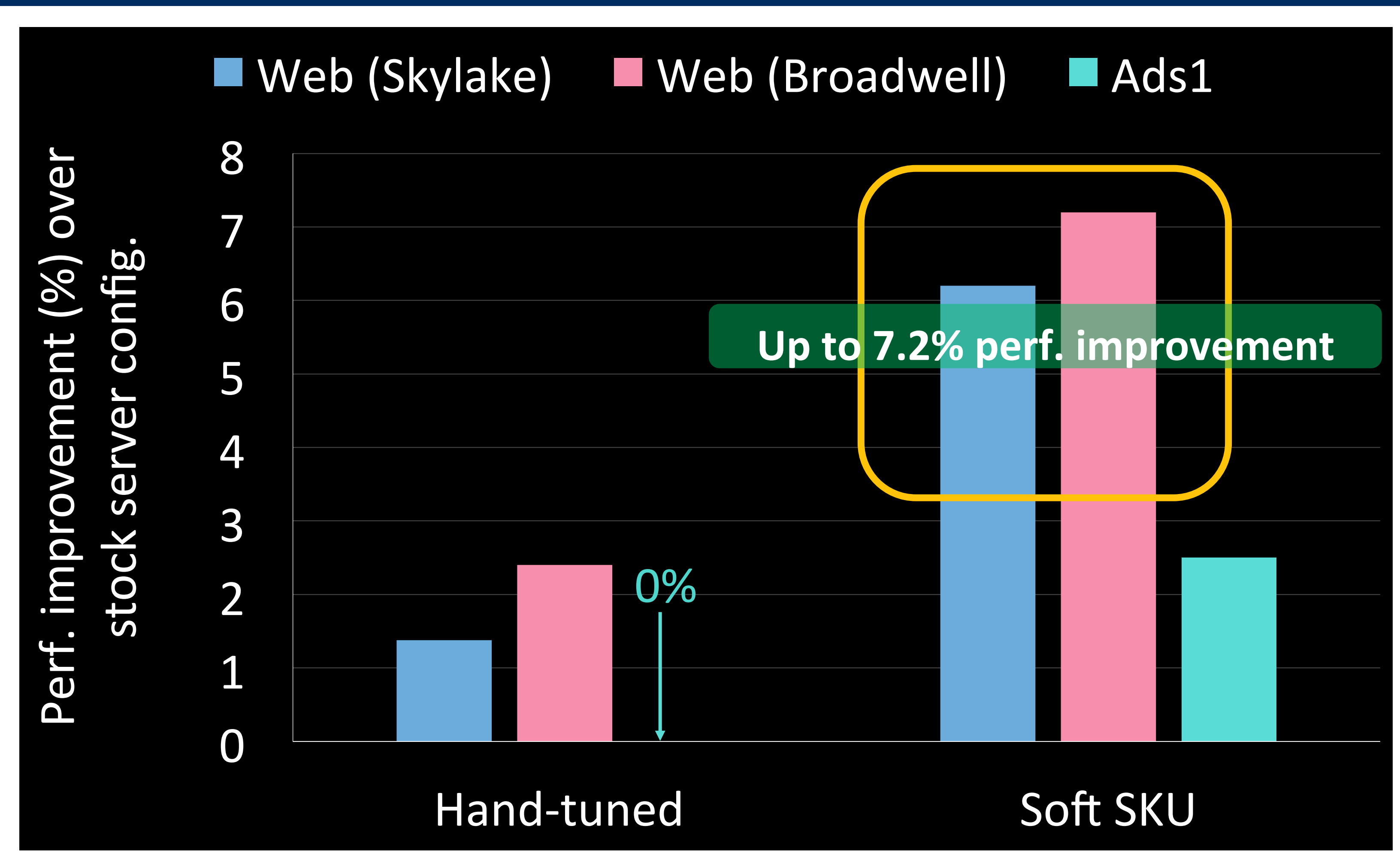
## $\mu$ SKU: Soft SKU Design & Deployment



## Performance of Commodity Servers



## Soft SKU Performance



Soft SKU can achieve ~7.2% throughput improvement on production systems with no extra hardware requirement

## Contributions

- Comprehensive characterization of Facebook's microservices
  - System-level & architectural bottlenecks
  - Reveals enormous bottleneck diversity across microservices
- Concept of "soft" server SKUs
  - Tuning coarse-grained OS & hardware configuration knobs
- $\mu$ SKU
  - Automates soft-SKU search & configuration via production A/B tests
  - Deploys soft SKUs on production microservices

~7.2% perf. boost on production  $\mu$ services + no extra hardware

