# Towards Energy Proportionality for Large-Scale Latency-Critical Workloads

**David Lo**†*, Liqun Cheng*, Rama Govindaraju*, Luiz André Barroso*, Christos Kozyrakis†

† *Stanford University*     * *Google Inc.*

# Motivation

- **Energy proportionality**
  - Servers are far less energy efficient at low and medium utilizations
  - Servers are underutilized due to diurnal load patterns

- **Large-scale latency-critical workloads**
  - Web search, social networking, etc.
  - Strict guarantees on tail latency and workload complexity precludes previous power management techniques

# Executive Summary

- Energy waste is caused by overachieving on performance

- **Solution**: **Match power to Service Level Objective (SLO)**
  - End-to-end SLO latency monitoring
  - Fine-grain power saving mechanism (i.e. RAPL)

- Built dynamic controller for large-scale latency-critical workloads
  - 20-30% power savings on production Google search without SLO violations

# Outline

- Energy proportionality vs. latency-critical workloads

- Recovering energy proportionality: iso-latency

- PEGASUS: QoS aware dynamic controller

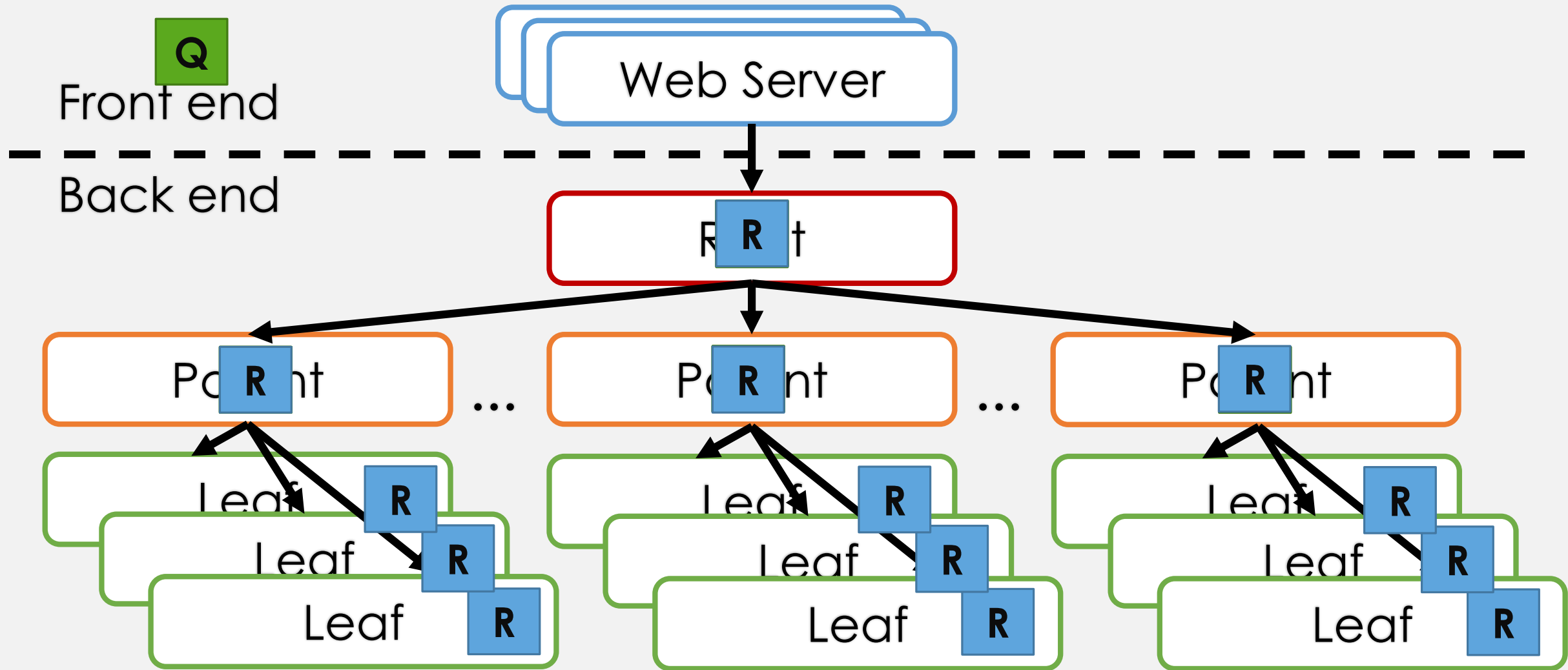# Energy proportionality vs. latency-critical workloads

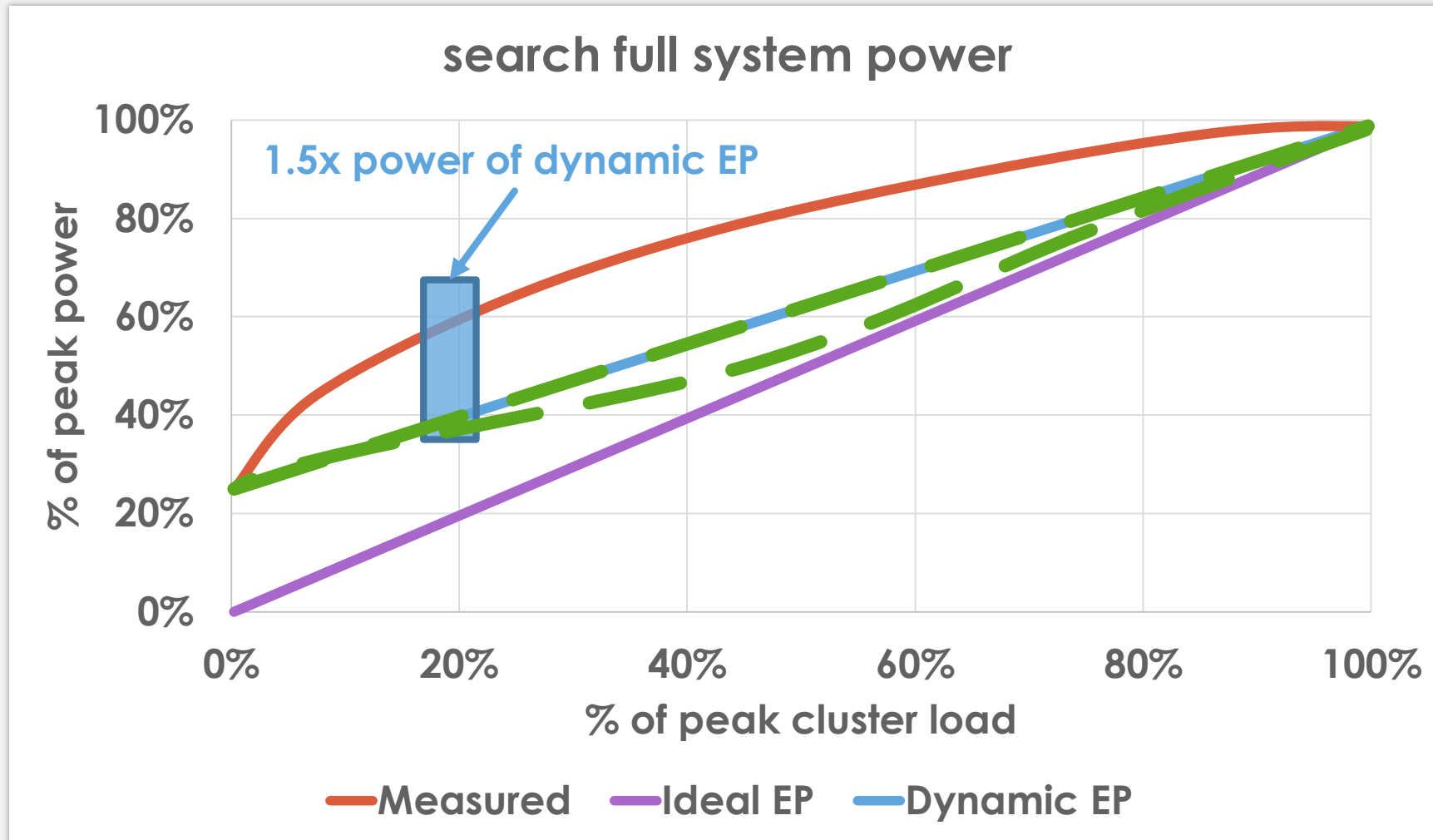## The case for latency-aware fine-grain power management

# OLDI workloads

- On-line Data Intensive (**OLDI**) workloads are user-facing workloads that mine massive datasets across many servers
  - Strict Service Level Objectives (**SLO**): e.g. 99%-ile tail latency is 5ms
  - High fan-out with large distributed state
  - Extremely challenging to perform power management

- Workload we evaluate on:
  - **search**: Query serving portion of production Google search
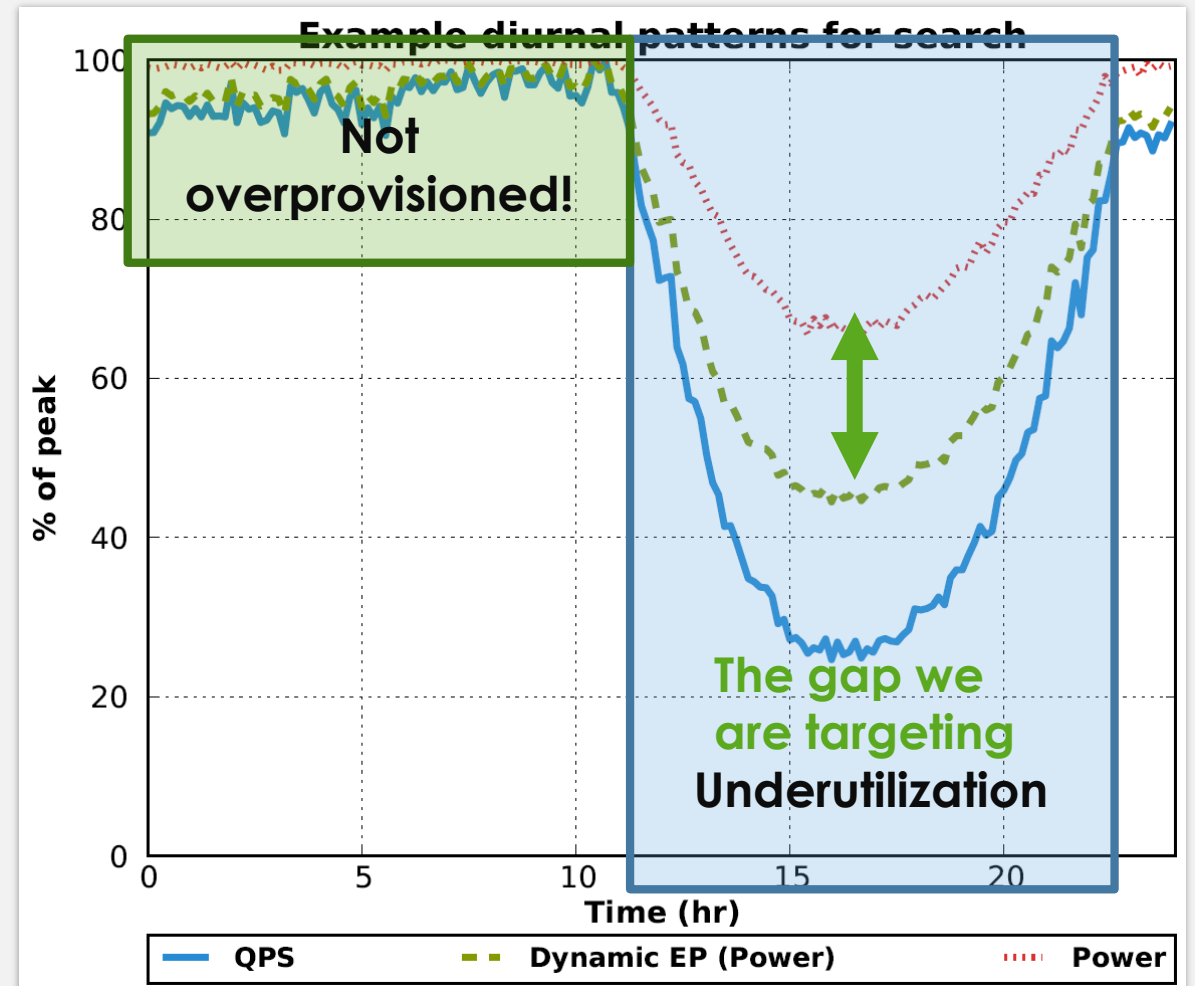
# Search topology

# The challenge of energy proportionality

# The need for energy proportionality

- Diurnal variation in cluster load and power for search across a 24 hour period

- Cluster not fully utilized half the time

- Gap between measured power and EP curves represent potential savings



Example diurnal patterns for search

Not overprovisioned!

The gap we are targeting
Underutilization

Legend: QPS — Dynamic EP (Power) — Power

Y-axis: % of peak (0, 20, 40, 60, 80, 100)
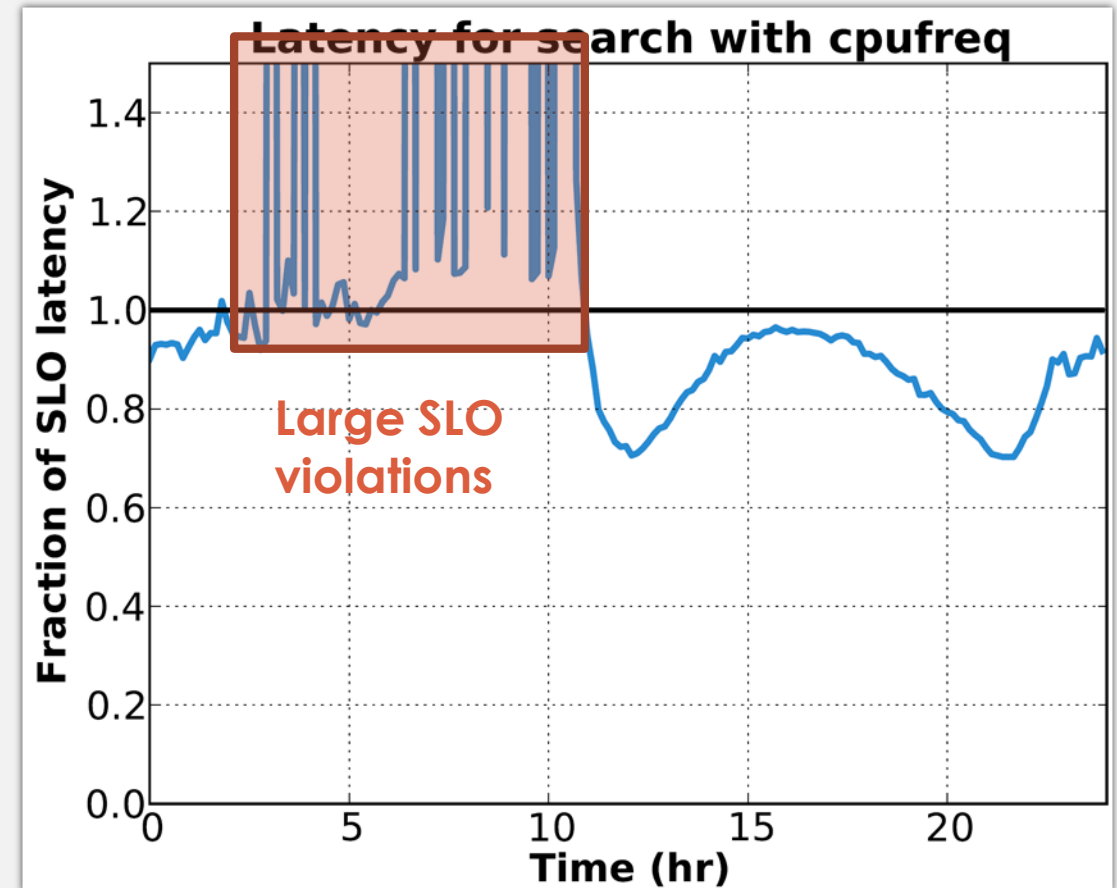X-axis: Time (hr) (0, 5, 10, 15, 20)

# Previous cluster-level power management

- ~~Consolidate load on fewer servers during low utilization~~
  - **Issue:** state of OLDI applications cannot fit on fewer servers

- ~~Use very low power idle modes~~
  - **Issue:** OLDI request rate is always too high, e.g. >1k requests/sec

- ~~Batch requests to form long enough idle periods~~
  - **Issue:** OLDI applications cannot tolerate msec exit times and batching delays

# Previous machine-level power management

- ~~CPU utilization based DVFS~~
  - Changes p-states based on CPU utilization
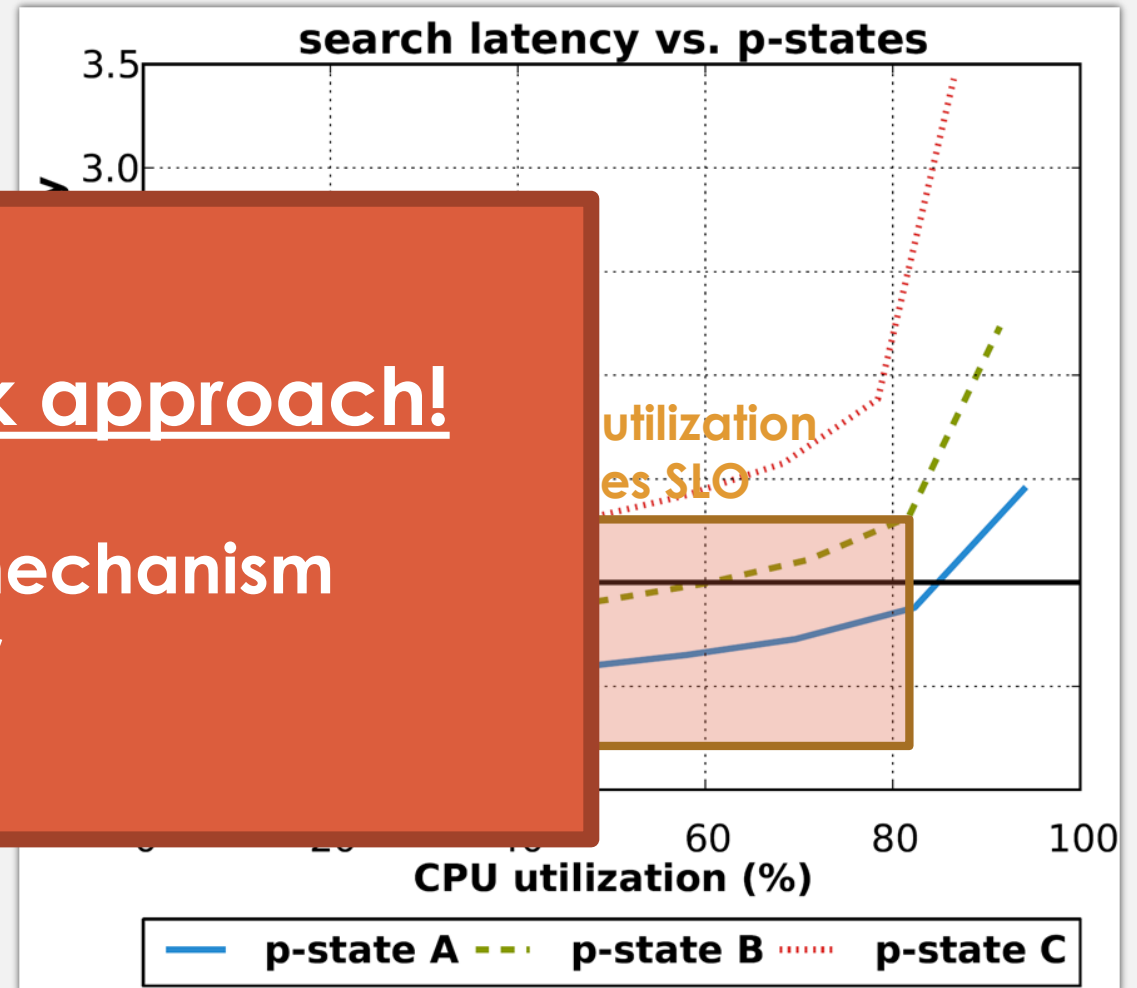
- **Issue**: causes SLO violations



Latency for search with cpufreq

Large SLO violations

# Weakness of current DVFS schemes

- CPU utilization is a poor proxy for workload latency

- **To meet SLO** latency-aw

**Need to rethink approach!**
- New policy
- New control mechanism
- New controller

search latency vs. p-states

3.5
3.0

utilization
es SLO

60    80    100
CPU utilization (%)
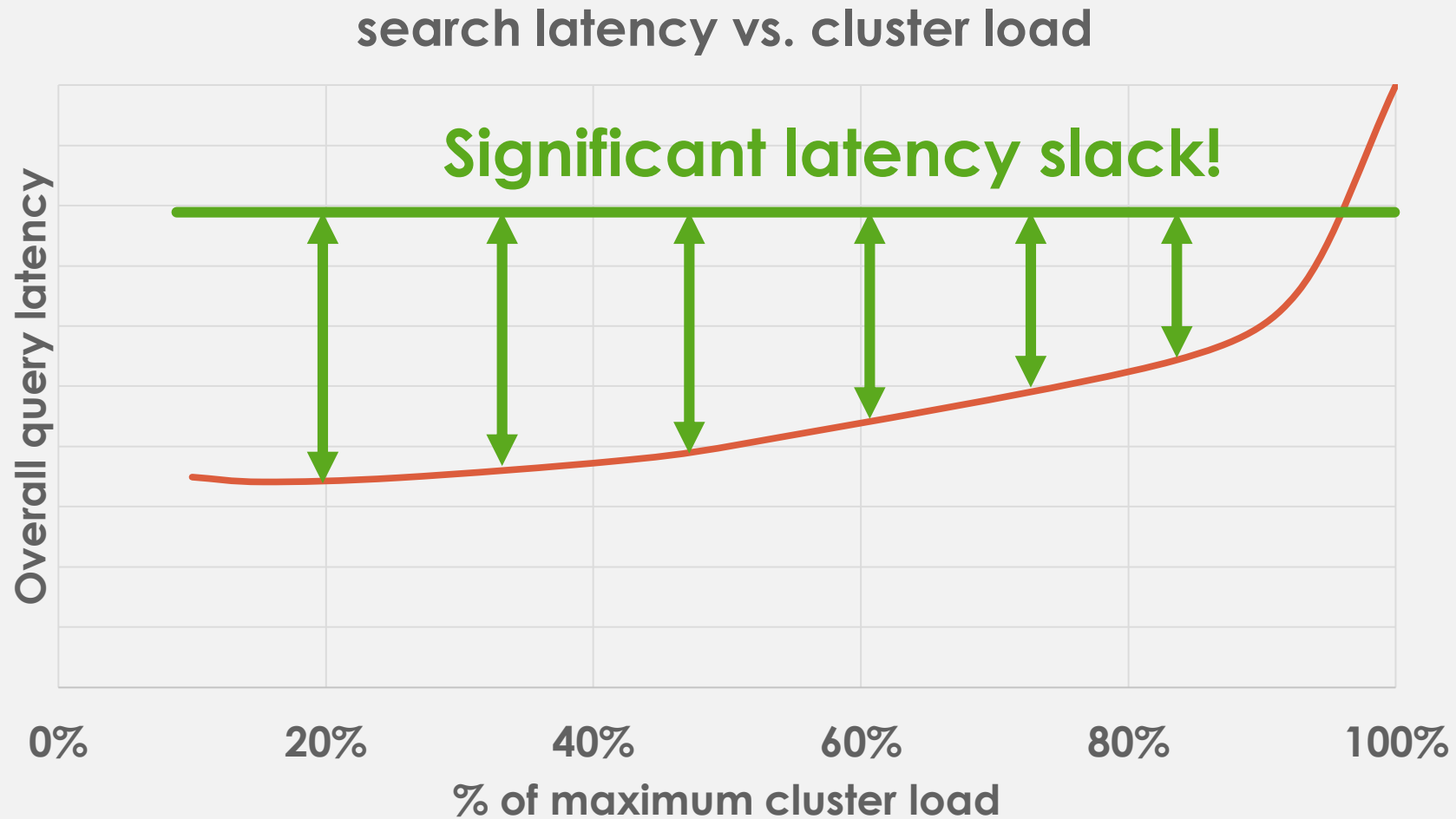
p-state A    p-state B    p-state C

# Recovering energy proportionality: iso-latency

## Trading end-to-end latency slack for immense power savings

# Motivating assumption

○ **Beating the end-to-end SLO is no better than meeting it**

   ○ The end-user only cares if the web page takes a long time to load

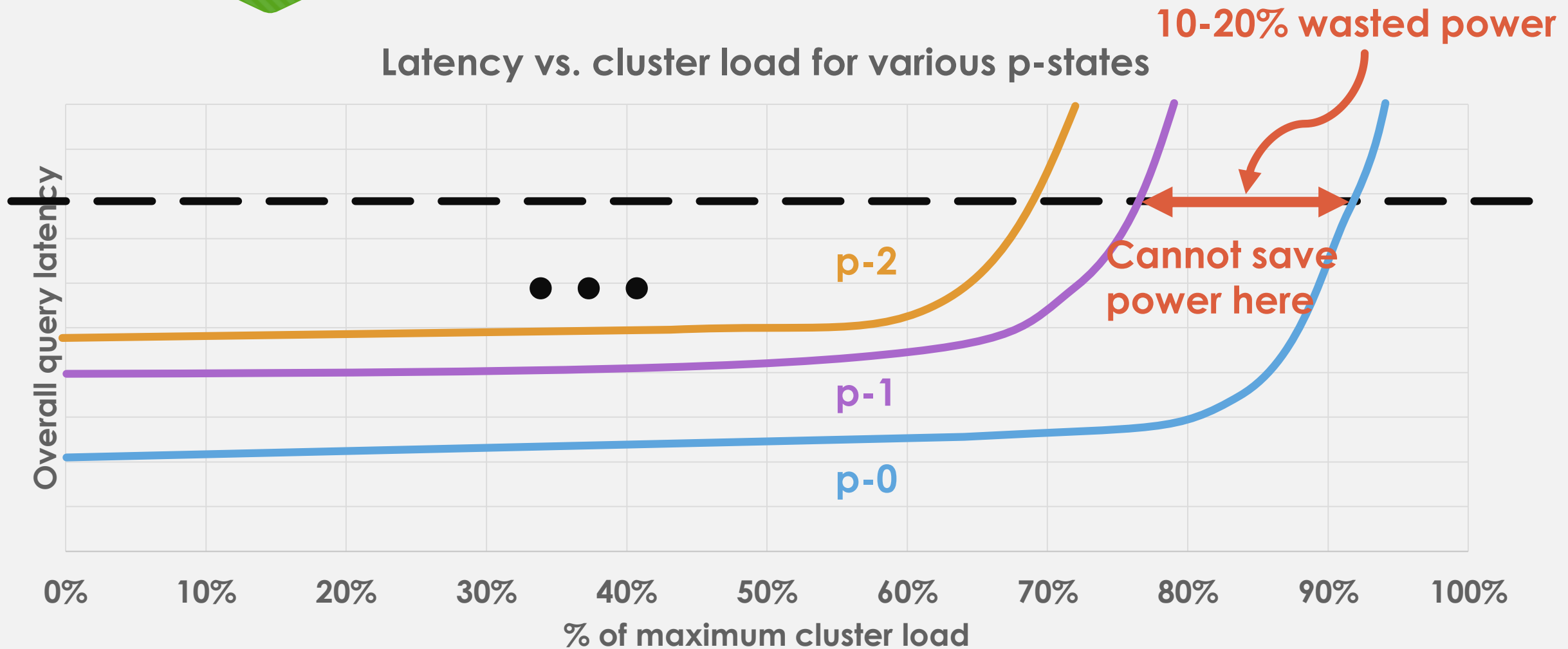   ○ If the page loads in 0.25sec vs. 0.50sec, user does not notice

# Latency opportunities

search latency vs. cluster load



**Significant latency slack!**

Overall query latency

0%    20%    40%    60%    80%    100%

% of maximum cluster load

# Iso-latency power management

- **Key idea:** Trade end-to-end latency slack for power savings
- Use power management mechanisms to keep the workload performing just well enough to avoid SLO violations
  - **Need end-to-end latency feedback from workload**
    - Most OLDI workloads have ways of measuring this
  - **Need fine-grained power management mechanisms**

# Problem: p-states are not fine grained

Latency vs. cluster load for various p-states

**10-20% wasted power**



Overall query latency

p-2

p-1

p-0

**Cannot save power here**

0%   10%   20%   30%   40%   50%   60%   70%   80%   90%   100%
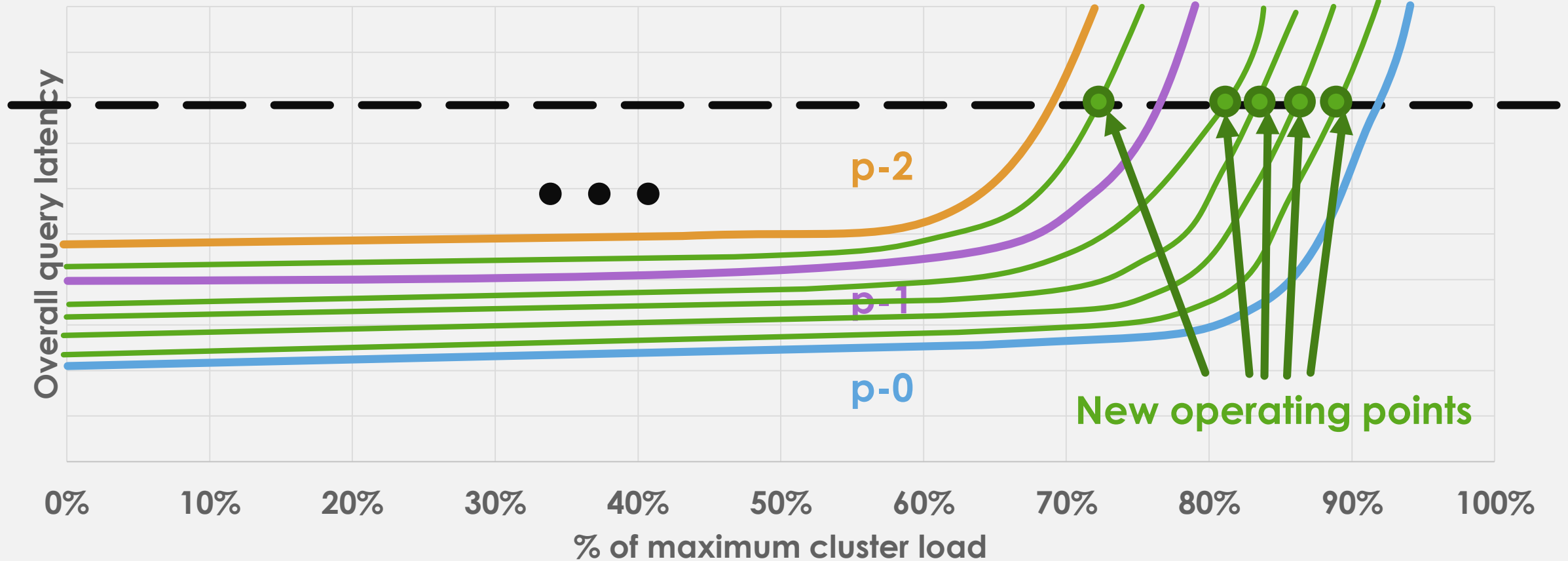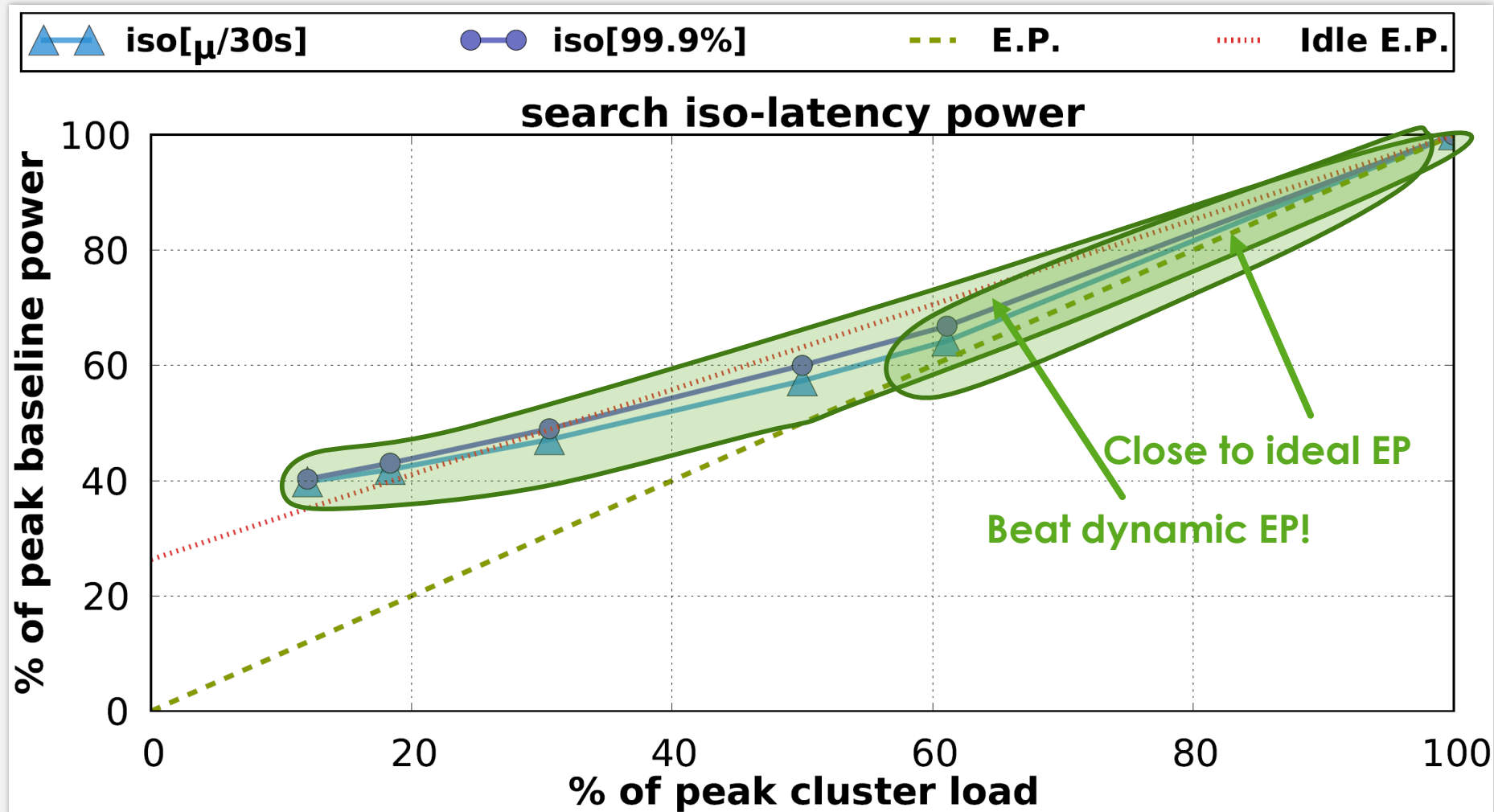
% of maximum cluster load

# Solution: RAPL

- RAPL: **R**unning **A**verage **P**ower **L**imit
- **Fine-grained**: power limit increments as small as 0.125W
- **Fast**: <1ms delay to apply new limit
- **Effective**: Dynamic Voltage Frequency Scaling (DVFS) behind the scenes to meet the power limit
  - More fine-grained than p-states
  - Can even modulate between multiples of base clock frequencies
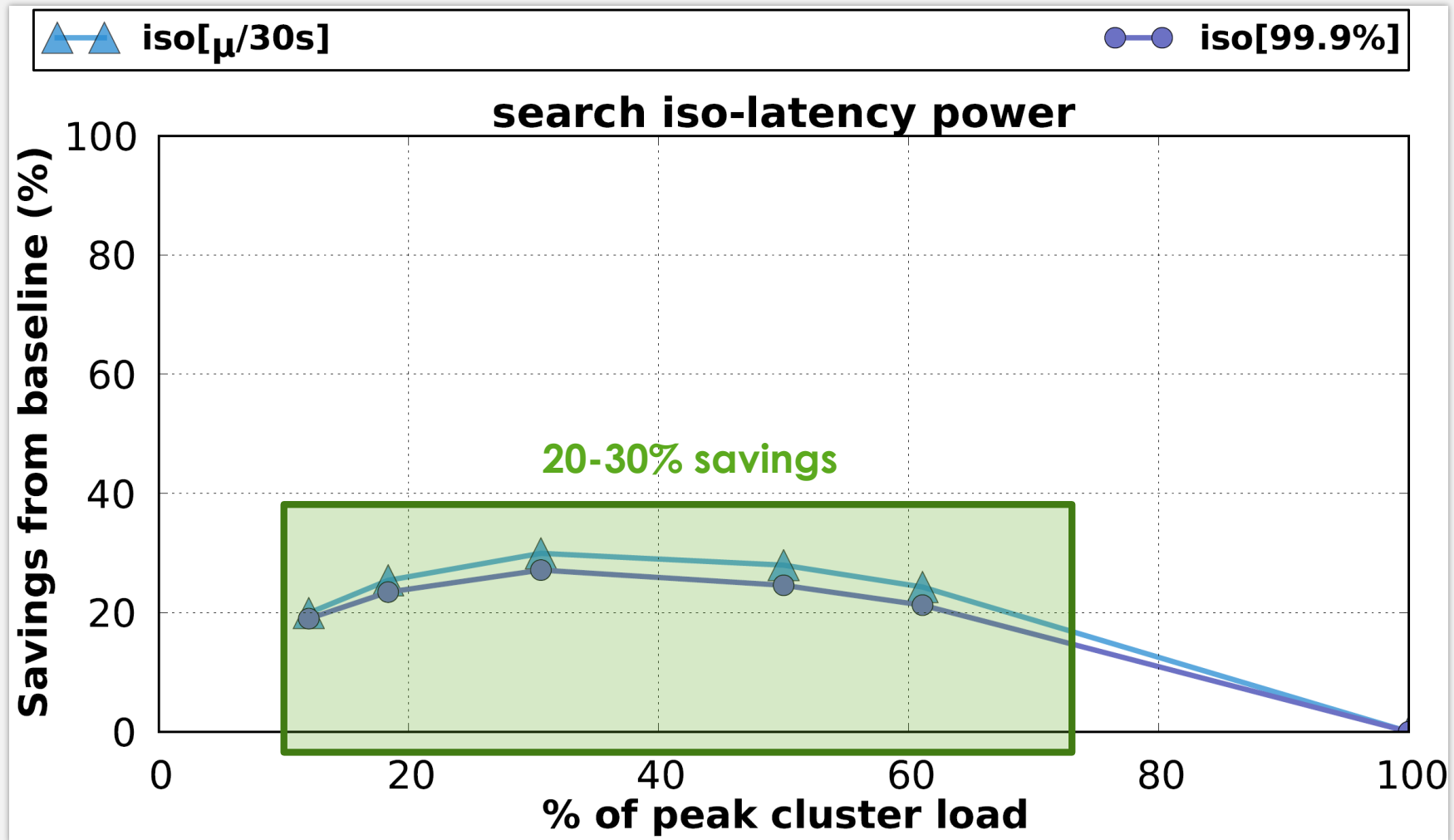
# Advantages of fine-grain control



Latency vs. load for various p-states — New RAPL states

Overall query latency

p-2

p-1

p-0

New operating points

% of maximum cluster load

0%  10%  20%  30%  40%  50%  60%  70%  80%  90%  100%

# Iso-latency potential: power

# Iso-latency potential: power savings

# **P**ower and **E**nergy **G**ains **A**utomatically **S**aved from **U**nderutilized **S**ystems

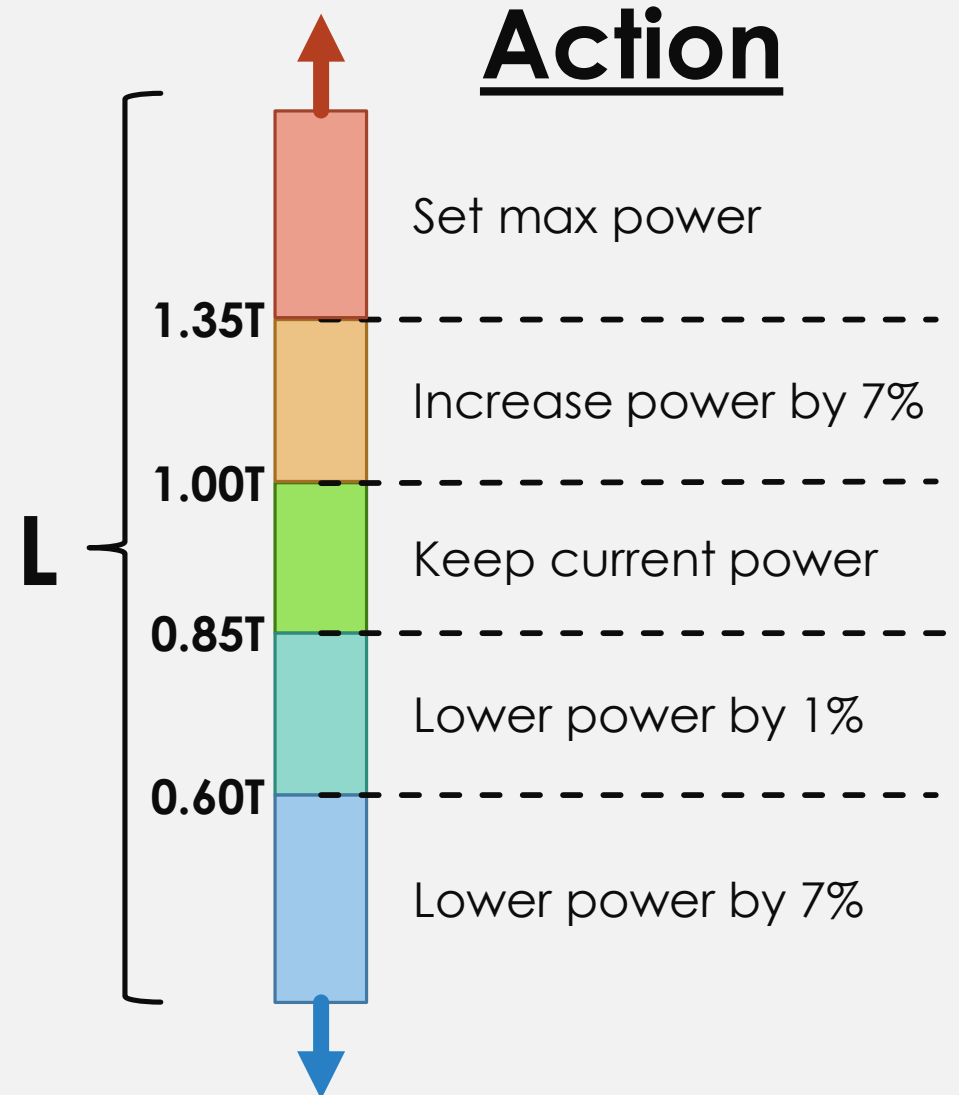## QoS aware dynamic controller

# PEGASUS description

- Real-time dynamic controller for **iso-latency**

- Use RAPL as knob for power

- Measures latency slack and sets **uniform** power limit across all servers

- Power is set by workload specific policy

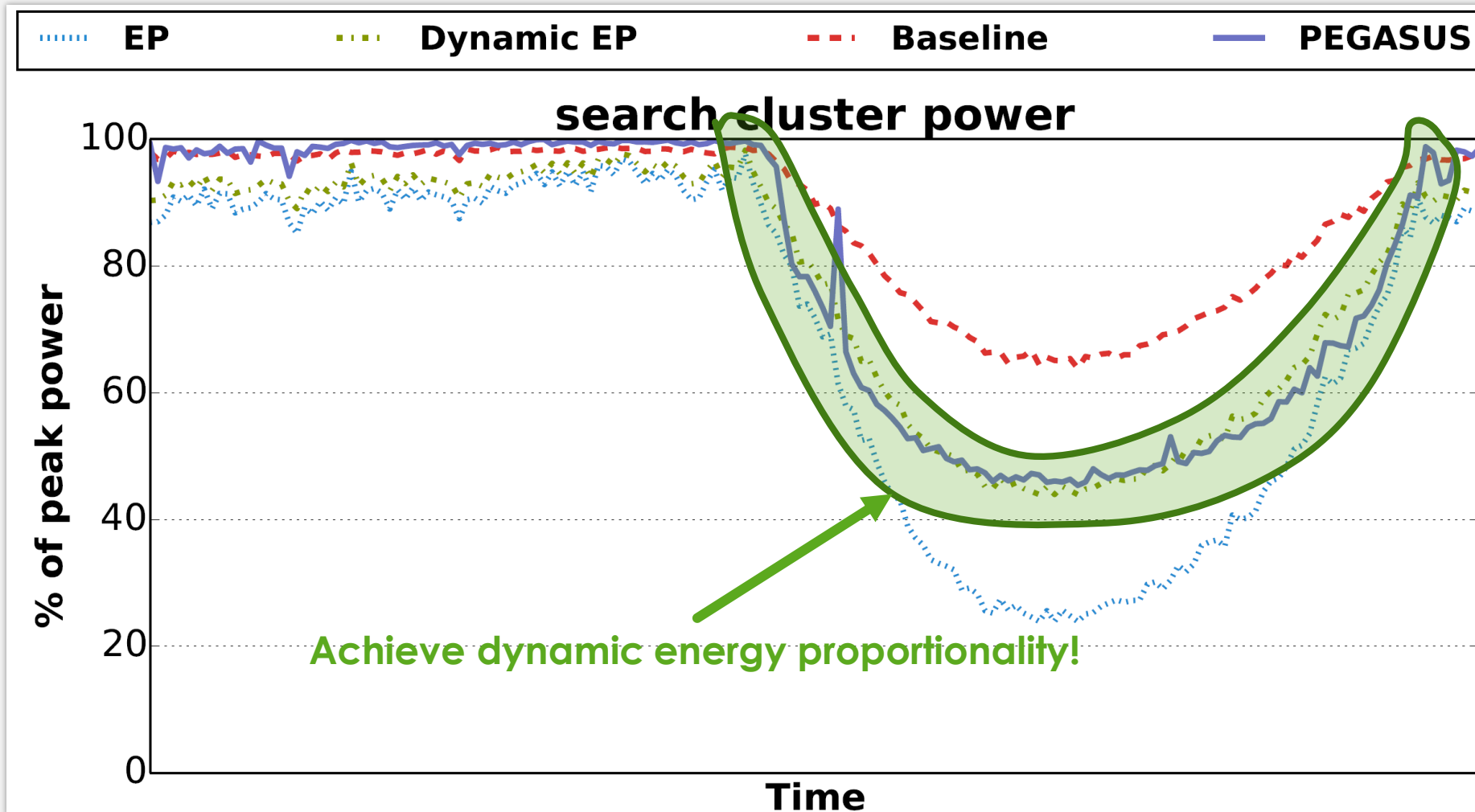# Example PEGASUS policy for search

- **L** = Measured instant latency

- **T** = SLO target

- Use instant latency for quick corrections

- Violating SLO latency triggers fail-safe

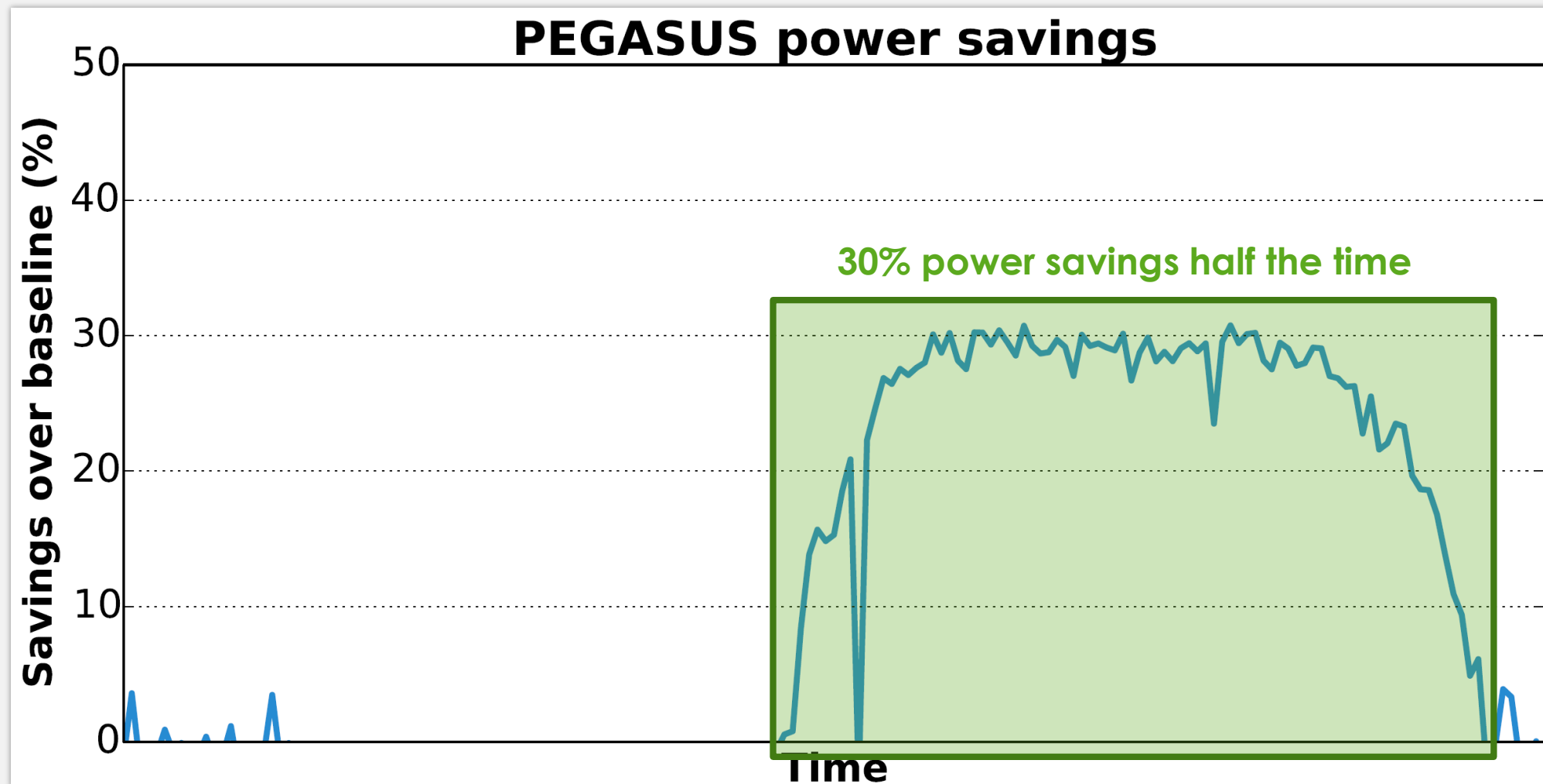- Constants determined through empirical optimization

## Action

**L** {

| | |
|---|---|
| | Set max power |
| **1.35T** | Increase power by 7% |
| **1.00T** | Keep current power |
| **0.85T** | Lower power by 1% |
| **0.60T** | Lower power by 7% |

# Evaluation methodology

- Workload parameters
  - **SLO metric**: 30 second average latency
  - Traffic pattern and user queries derived from anonymized search logs
  - Index derived from production search index
- Evaluate on several cluster sizes
  - **Small**: tens of machines, use full 24hr trace
  - **Production**: thousands of machines, use 12hr portion
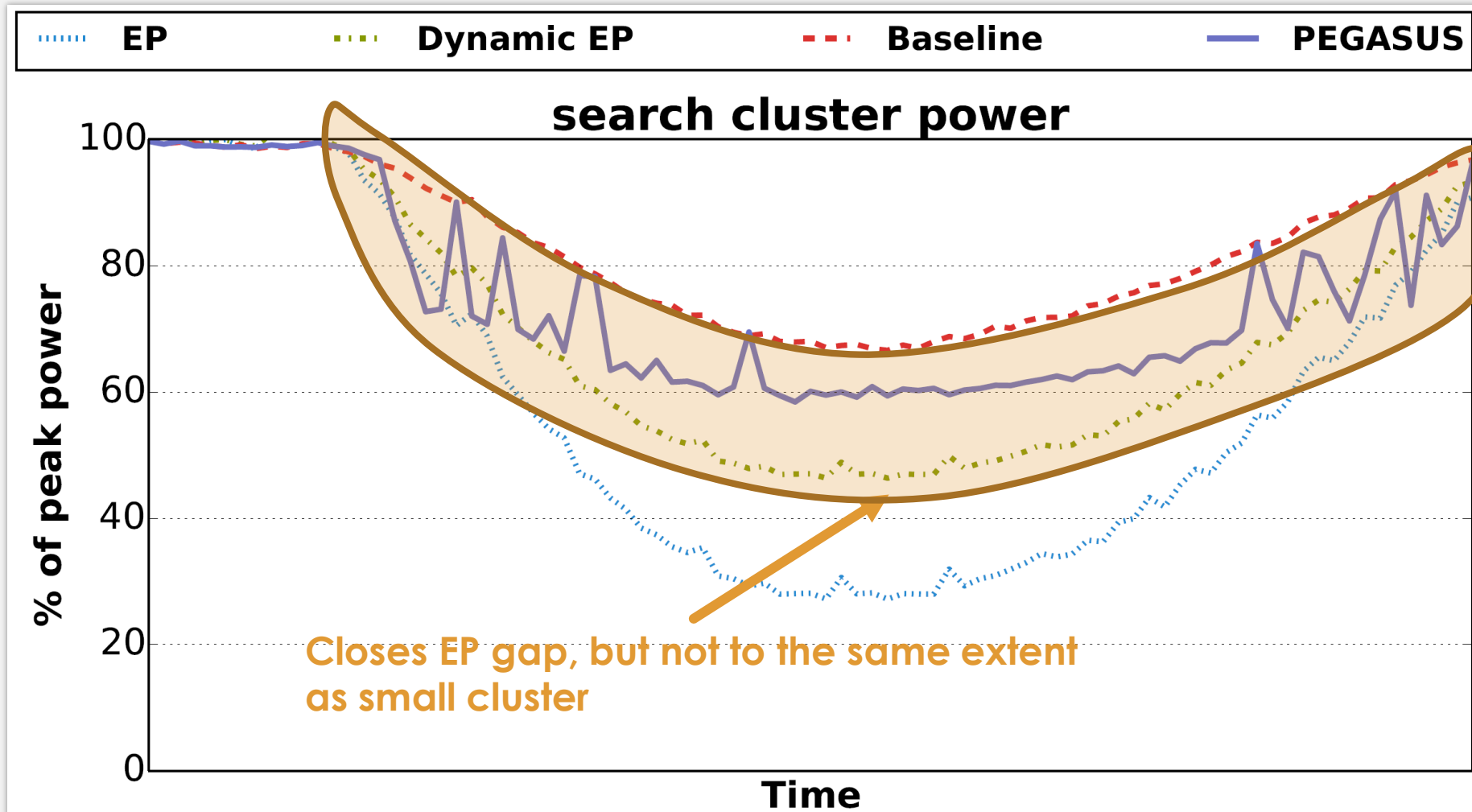- Measure full cluster power and SLO latency

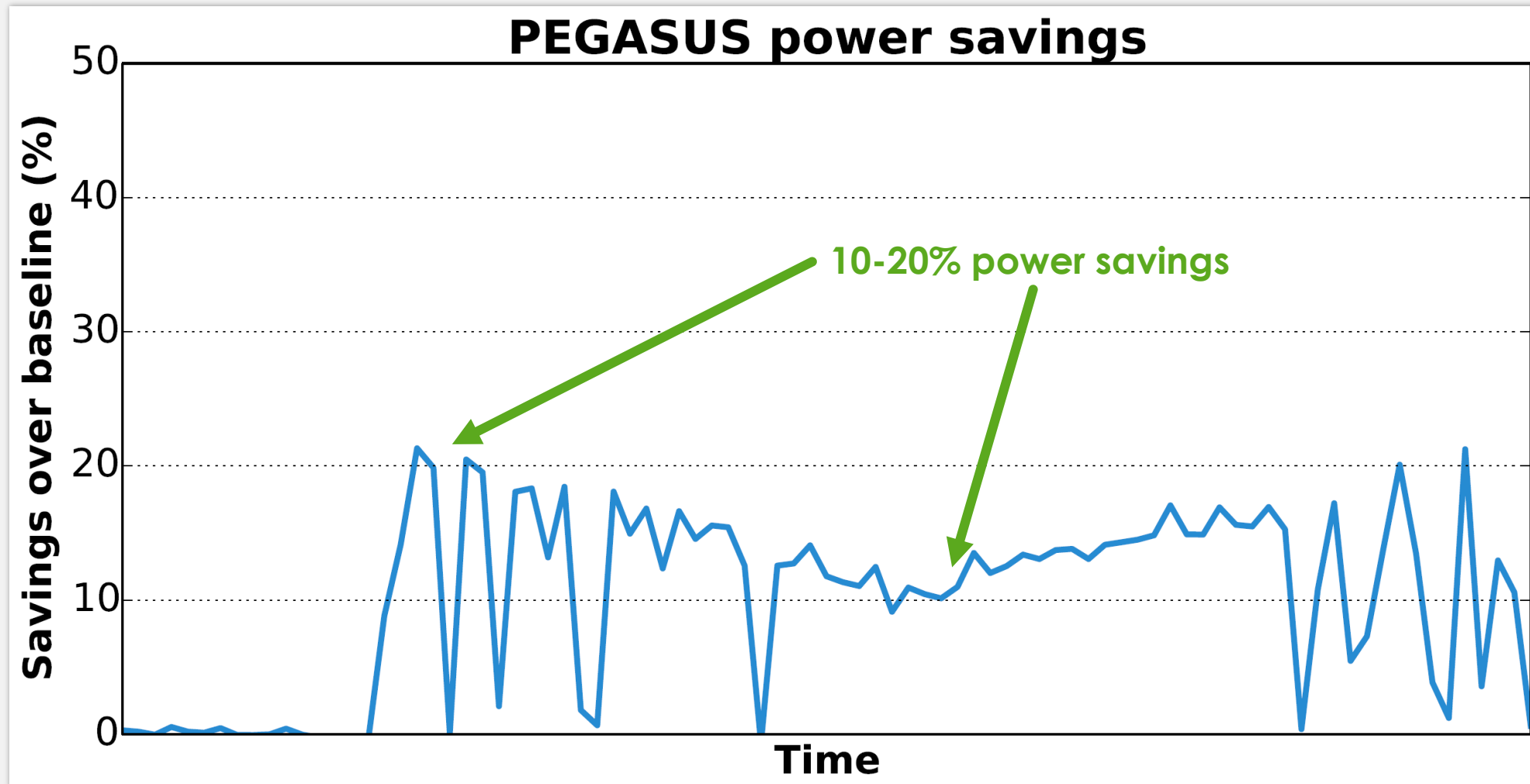# Small cluster results: power over time



search cluster power

Legend: EP, Dynamic EP, Baseline, PEGASUS

% of peak power

100
80
60
40
20
0

Time

Achieve dynamic energy proportionality!

# Small cluster results: power comparison



**PEGASUS power savings**

30% power savings half the time

Savings over baseline (%)

Time

# Production cluster results: power over time

# Production cluster results: power comparison



**PEGASUS power savings**

10-20% power savings

Savings over baseline (%)
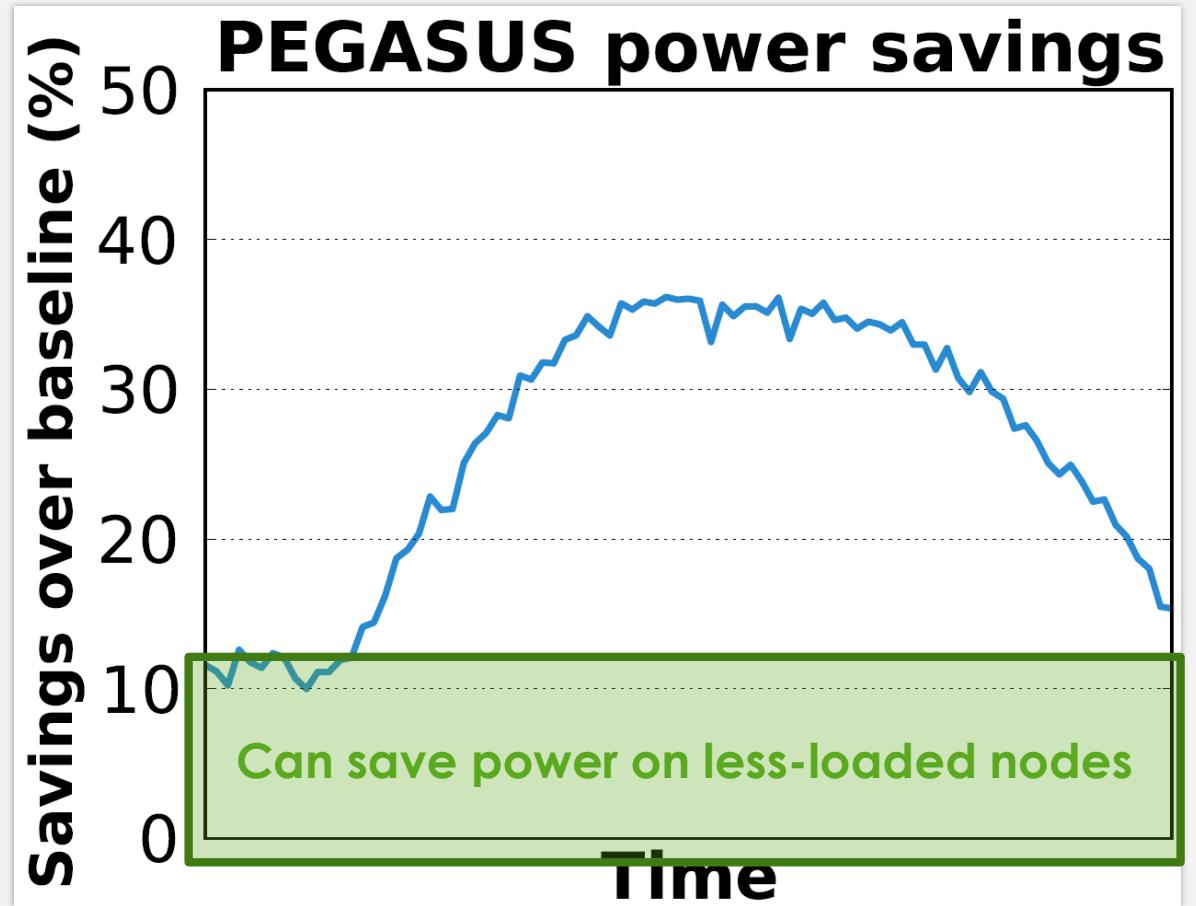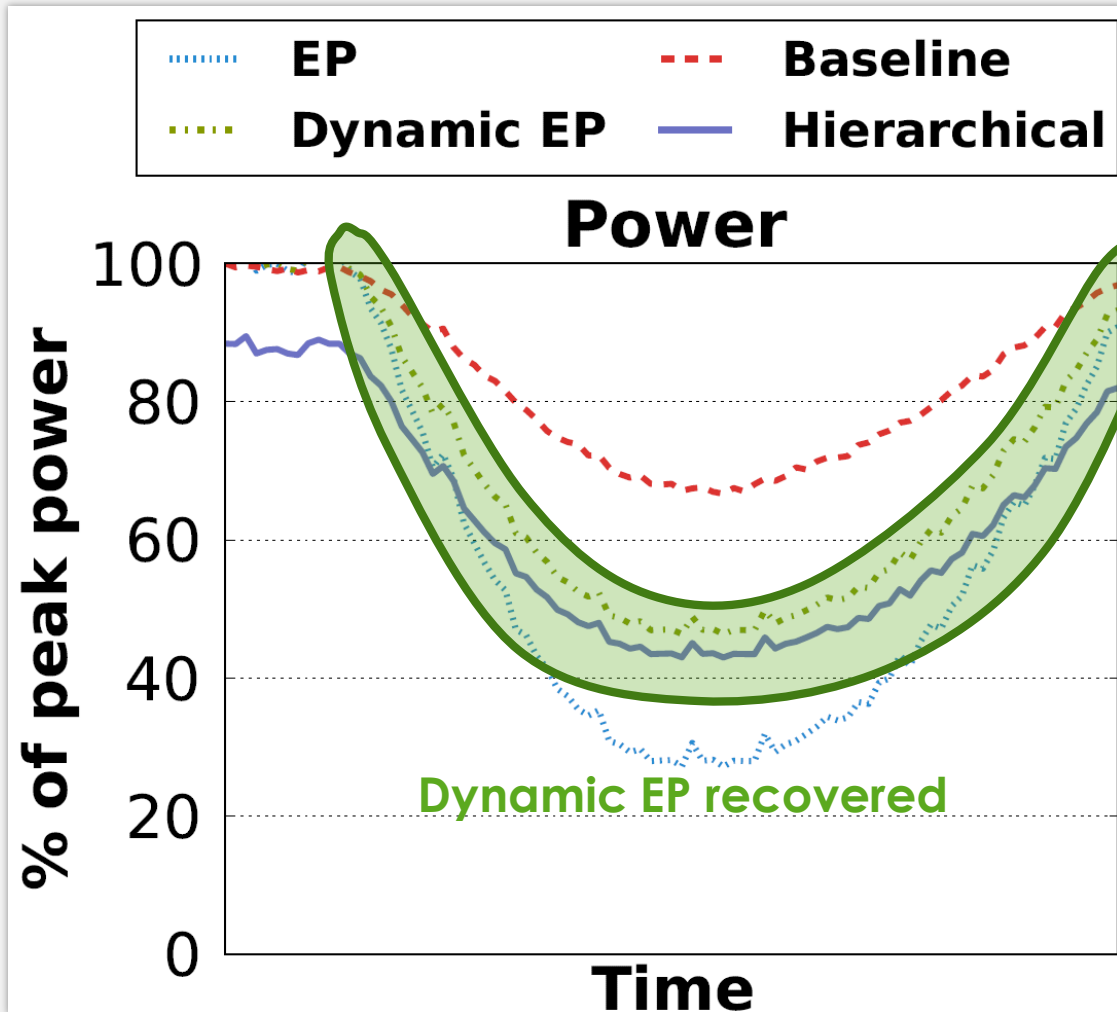
Time

# Improving PEGASUS scalability

- Production cluster sees **"tail at scale"** for server utilization
  - At peak load, 0.2% nodes at 100% load while 50% nodes at <85% load
  - Caused by popular queries hitting a few shards
  - **Issue:** Hot nodes set lower bound on power limits for everyone

- **Idea**: hierarchical control
  - **Global:** sets latency targets instead of power limits
  - **Local:** decides amount of power needed to meet target latency

# Hierarchical PEGASUS design

# Estimated hierarchical PEGASUS results

# Conclusion

- Halfway there to fully energy proportional systems
- **Iso-latency:** Use SLO metrics and fine-grain power control
  - Save up to 30% power
  - Meet/exceed energy proportionality targets
- PEGASUS achieves **iso-latency** benefits
  - Up to 20% savings on production cluster
  - Be aware of tail at scale effects