# Improving Resource Efficiency
# In Cloud Computing

## Christina Delimitrou

*Stanford University*
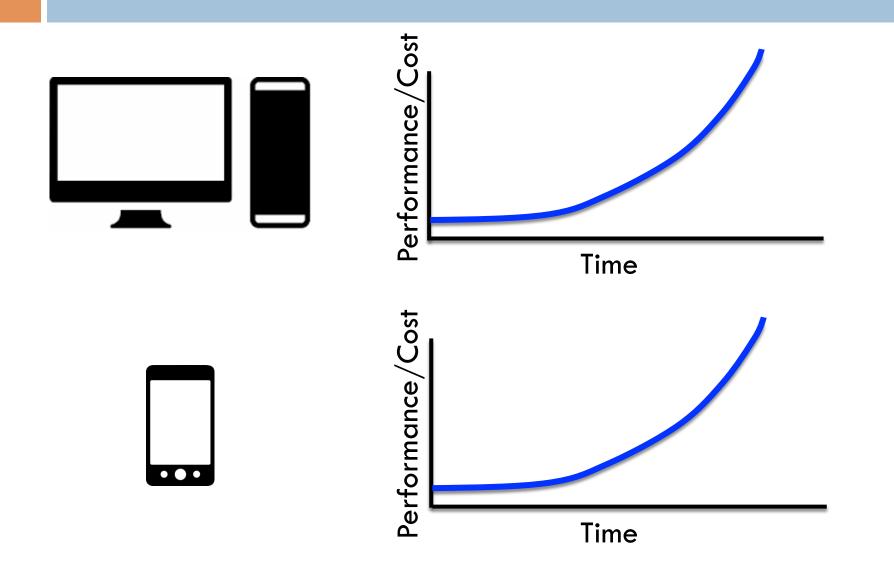
# Resource efficiency is a first-order system constraint

How efficiently do we utilize resources?

How efficiently do we design systems?

# Why Care about Resource Efficiency?

~10K commodity servers
Sophisticated cluster managers
~10s MWatts
$100,000,000s

Private clouds:
• Google, Microsoft, Twitter, eBay
**Public clouds:**
• Amazon EC2, Windows Azure, GCE

Google | google.com/datacenters

# The Promise of Cloud Computing

□ Flexibility

  ■ Provision and launch new services in seconds

□ High performance

  ■ High throughput & low tail latency

□ Cost effectiveness

  ■ Low capital & operational expenses

Cloud computing scalability:
high performance AND low cost

# The Reality of Cloud Computing

# Scaling Datacenters

- ~~Switch to commodity servers~~      One time trick
- ~~Improve cooling/power distribution~~      < 10%

- ~~Build more datacenters~~      >$300M per datacenter
- ~~Add more servers~~      Power limit
- ~~Rely on processor technology~~      End of voltage scaling

**Use existing systems more efficiently**

# Datacenter Underutilization

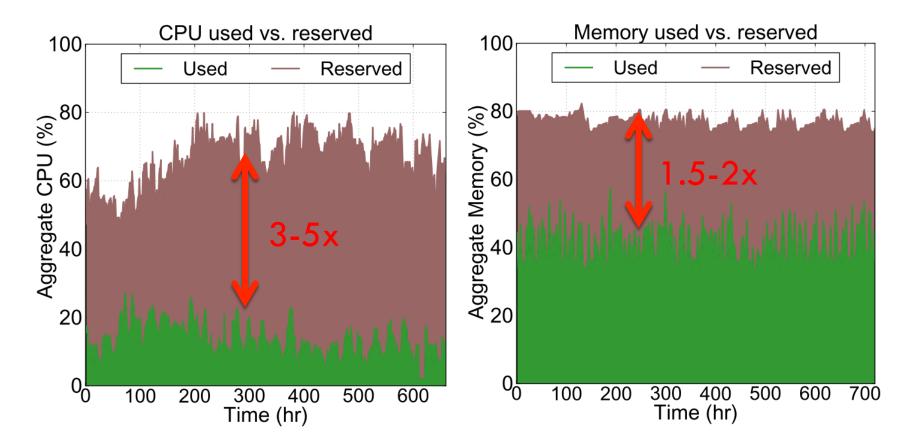**Twitter** (Mesos)[1]



**Google** (Borg)[2]



[1] C. Delimitrou and C. Kozyrakis. Quasar: Resource-Efficient and QoS-Aware Cluster Management, ASPLOS 2014.

[2] L. A. Barroso, U. Holzle. The Datacenter as a Computer, 2013.

# Datacenter Underutilization…

~~Is the cluster manager's fault~~

Is the user's fault!

# Reserved vs. Used Resources



CPU used vs. reserved — Aggregate CPU (%) vs. Time (hr); Used (green), Reserved (brown); 3-5x

Memory used vs. reserved — Aggregate Memory (%) vs. Time (hr); Used (green), Reserved (brown); 1.5-2x

□ Twitter: up to 5x CPU & up to 2x memory overprovisioning

# Reserved vs. Used Resources



~25,000 jobs
936 distinct users

[ASPLOS'14]

Reservation=Usage

☐ 20% of job under-sized, **~70% of jobs over-sized**

# Datacenter Underutilization...

Is the user's fault!

(not really...)

# Resource Management is Hard

# Performance Depends on



Scale-up — Performance vs Cores

# Performance Depends on



Heterogeneity

Performance vs Cores

# Performance Depends on



Heterogeneity × Scale-out

# Performance Depends on



Heterogeneity — Performance vs. Cores

Scale-out — Performance vs. Servers

Input load — Performance vs. Input size

# Performance Depends on



Heterogeneity

× Scale-out

Interference × Input load

Performance

Interference

Performance

Input size

**Overprovision Reservations!**

when code changes, when platforms change, etc.

# Can we improve resource efficiency while preserving application QoS guarantees?

Potential: 3-5x efficiency; $10Ms in cost savings

# Requirements

- Automate resource management

  - Large, multi-dimensional space → Leverage big data

- General solution

  - Different application types (batch, latency-critical)

  - Different types of hardware

- Cross-layer design

  - Architecture → OS → Scheduler → Application design

# Contributions

# Contributions

**Paragon** [ASPLOS'13, TopPicks'14]
[IISWC'13]

Users → Resource reservations → Scheduler → Cluster

**1.** Practical data mining

# Contributions

**Quasar** [ASPLOS'14]



**2.** High level interface

Resource reservations

Users → Scheduler

Cluster

**1.** Practical data mining

# Contributions

**Systems:**

Application assignment: **Paragon** [ASPLOS'13, TopPicks'14, CAL'13, IISWC'13]

Cluster management: **Quasar** [ASPLOS'14]

# Contributions

**Systems:**

Application assignment: **Paragon** [ASPLOS'13, TopPicks'14], **iBench** [IISWC'13]

Cluster management: **Quasar** [ASPLOS'14]

Scalable scheduling: **Tarcil** [SOCC'15]

# Contributions

**Systems:**

Application assignment: **Paragon** [ASPLOS'13, TopPicks'14], **iBench** [IISWC'13]

Cluster management: **Quasar** [ASPLOS'14]

Scalable scheduling: **Tarcil** [SOCC'15]

Cloud provisioning: **Hybrid Cloud** [in submission]

# Contributions

**Systems:**

Application assignment:   **Paragon** [ASPLOS'13, TopPicks'14], **iBench** [IISWC'13]

Cluster management:   **Quasar** [ASPLOS'14]

Scalable scheduling:   **Tarcil** [SOCC'15]

Cloud provisioning:   **Hybrid Cloud** [in submission]

Admission control:   **ARQ** [ICAC'13]

# Contributions

**Systems:**

Application assignment:  **Paragon** [ASPLOS'13, TopPicks'14], **iBench** [IISWC'13]

Cluster management:  **Quasar** [ASPLOS'14]

Scalable scheduling:  **Tarcil** [SOCC'15]

Cloud provisioning:  **Hybrid Cloud** [in submission]

Admission control:  **ARQ** [ICAC'13]

**Datacenter application modeling:**
**ECHO** [IISWC'12], **Storage application modeling** [CAL'12, IISWC'11, ISPASS'11]

# Paragon

Users → Resource reservations → Scheduler ⟷ Cluster

Practical data mining techniques

# Heterogeneity & Interference Matter



- **Heterogeneity**
  - DCs provisioned over 15 years
  - Multiple server generations & configurations

- **Interference**
  - Apps contend on shared resources
    - CPU & cache hierarchy
    - Memory system
    - Storage & network I/O

# Extracting Resource Preferences

□ **Naïve**: exhaustive characterization

  ▫ ~10-20 platforms $x$ 1,000 apps

Resource
reservations

Users   App → S App uler   Cluster

App   App

**Mine**
**big data**

**Data**

□ Looks like a recommendation problem

# Recommendation Systems

- Content-based systems:
  - Description of items (keywords, feature vector, etc. )
  - Profile of user preferences (history, model, user-system interaction, etc. )

- Collaborative filtering:
  - Uncover similarities between users and items
  - No need to know item features or explicit user preferences in advance

# Recommendation Systems

□ Content-based systems:

  ◻ Description of items (keywords, feature vector, etc. )

  ◻ Profile of user preferences (history, model, user-system interaction, etc. )

□ Collaborative filtering:

  ◻ Uncover similarities between users and items

  ◻ No need to know item features or explicit user preferences in advance

# Something familiar…

□ Collaborative filtering – similar to Netflix Challenge system

    ■ Singular Value Decomposition (SVD) + PQ reconstruction (SGD)



Sparse utility matrix           Dense utility matrix

# SVD



$$
\begin{array}{c}
\phantom{u_1} \\
u_1 \\
u_2 \\
\vdots \\
u_m
\end{array}
\begin{array}{cccc}
m_1 & m_2 & \ldots & m_n \\
\end{array}
\begin{bmatrix}
a_{11} & a_{12} & \ldots & a_{1n} \\
a_{21} & a_{22} & \ldots & a_{2n} \\
\vdots & \vdots & \ddots & \vdots \\
a_{m1} & a_{m2} & \ldots & a_{mn}
\end{bmatrix}
$$

**movie**

**rating (e.g., ★★★★)**

**user**

$$
=
$$

$$
\begin{array}{c}
u_1 \\
\ldots \\
u_m
\end{array}
\begin{bmatrix}
u_{11} & \ldots & u_{1r} \\
\vdots & \ddots & \vdots \\
u_{m1} & \ldots & u_{mr}
\end{bmatrix}
\times
\begin{bmatrix}
\sigma_1 & \ldots & 0 \\
\vdots & \ddots & \vdots \\
0 & \ldots & \sigma_r
\end{bmatrix}
\times
\begin{array}{ccc}
m_1 & \ldots & m_n \\
\end{array}
\begin{bmatrix}
v_{11} & \ldots & v_{1r} \\
\vdots & \ddots & \vdots \\
v_{n1} & \ldots & v_{nr}
\end{bmatrix}
$$

35

# SVD

$$
\begin{array}{c}
\begin{array}{cccc}
\textcolor{red}{m_1} & \textcolor{red}{m_2} & \textcolor{red}{\ldots} & \textcolor{red}{m_n}
\end{array} \\
\begin{array}{c}
\textcolor{red}{u_1} \\ \textcolor{red}{u_2} \\ \vdots \\ \textcolor{red}{u_m}
\end{array}
\begin{bmatrix}
a_{11} & a_{12} & \ldots & a_{1n} \\
a_{21} & a_{22} & \ldots & a_{2n} \\
\vdots & \vdots & \ddots & \vdots \\
a_{m1} & a_{m2} & \ldots & a_{mn}
\end{bmatrix}
\end{array}
$$

user → (circle around $u_2$)

movie → (circle around $m_n$)

rating (e.g., ★★★★) → (circle around $a_{2n}$)

correlation of user to similarity concept → (circle around $u_{11}$)

$$=$$

$$
\begin{array}{c}
\textcolor{red}{u_1} \\ \textcolor{red}{\ldots} \\ \textcolor{red}{u_m}
\end{array}
\begin{bmatrix}
u_{11} & \ldots & u_{1r} \\
\vdots & \ddots & \vdots \\
u_{m1} & \ldots & u_{mr}
\end{bmatrix}
\;\times\;
\begin{bmatrix}
\sigma_1 & \ldots & 0 \\
\vdots & \ddots & \vdots \\
0 & \ldots & \sigma_r
\end{bmatrix}
\;\times\;
\begin{array}{c}
\begin{array}{ccc}
\textcolor{red}{m_1} & \textcolor{red}{\ldots} & \textcolor{red}{m_n}
\end{array} \\
\begin{bmatrix}
v_{11} & \ldots & v_{1r} \\
\vdots & \ddots & \vdots \\
v_{n1} & \ldots & v_{nr}
\end{bmatrix}
\end{array}
$$

similarity concept → (circle around $\sigma_1$)

correlation of movie to similarity concept → (circle around $v_{11}$)

36

# Heterogeneity Classification

|  | Movie 1 | Movie 2 | Movie 3 | Movie 4 | Movie 5 | ··· | Movie M |
|---|---|---|---|---|---|---|---|
| User A | ★★★★ |  |  | ★★★★★ |  |  |  |
| User B |  |  | ★★ |  | ★★★ |  |  |
| ⋮ |  |  |  |  |  |  |  |
| User N |  | ★★★★ |  |  |  |  | ★ |

# Heterogeneity Classification

|  | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | ... | Platform M |
|---|---|---|---|---|---|---|---|
| User A | ★★★★ |  |  | ★★★★★ |  |  |  |
| User B |  |  | ★★ |  | ★★★ |  |  |
| ⋮ |  |  |  |  |  |  |  |
| User N |  | ★★★★ |  |  |  |  | ★ |

# Heterogeneity Classification

|  | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | ... | Platform M |
|---|---|---|---|---|---|---|---|
| App A | ★★★★ |  |  | ★★★★★ |  |  |  |
| App B |  |  | ★★ |  | ★★★ |  |  |
| ⋮ |  |  |  |  |  |  |  |
| App N |  | ★★★★ |  |  |  |  | ★ |

# Heterogeneity Classification

| | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | ... | Platform M |
|---|---|---|---|---|---|---|---|
| App A | 1,500QPS | | | 843QPS | | | |
| App B | | | 458QPS | | 946QPS | | |
| ⋮ | | | | | | | |
| App N | | 1,016QPS | | | | | 186QPS |

App performance

# Heterogeneity Classification

|  | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | ... | Platform M |
|---|---|---|---|---|---|---|---|
| App A | 1,500QPS | | | 843QPS | | | |
| App B | | | | | | | |
| ... | | | | | | | |
| App N | | | | | | | |

**Profiled Performance**

Inferred Performance

# Heterogeneity Classification

|  | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | . . . | Platform M |
|---|---|---|---|---|---|---|---|
| App A | 1,500QPS | 843QPS | 675QPS | 843QPS | 1,786QPS | . . . | 8,675QPS |
| App B |  |  |  |  |  |  |  |
| ⋮ |  |  |  |  |  |  |  |
| App N |  |  |  |  |  |  |  |

Profiled Performance

Inferred Performance

42

# Heterogeneity Classification

|  | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | ... | Platform M |
|---|---|---|---|---|---|---|---|
| App A | 1,500QPS | 843QPS | 675QPS | 843QPS | 1,786QPS | ... | 8,675QPS |
| App B | 987QPS | 458QPS | 773QPS | 1,073QPS | 986QPS | ... | 1,836QPS |
| ⋮ |  |  |  |  |  |  |  |
| App N |  |  |  |  |  |  |  |

Profiled Performance

Inferred Performance

# Heterogeneity Classification

|  | Platform 1 | Platform 2 | Platform 3 | Platform 4 | Platform 5 | ... | Platform M |
|---|---|---|---|---|---|---|---|
| App A | 1,500QPS | 843QPS | 675QPS | 843QPS | 1,786QPS | ... | 8,675QPS |
| App B | 987QPS | 458QPS | 773QPS | 1,073QPS | 986QPS | ... | 1,836QPS |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| App N | 9,893QPS | 7,686QPS | 786QPS | 1,118QPS | 997QPS | ... | 1,354QPS |

**Performance depends on app type: QPS, completion time, IPC, …**

Profiled Performance

Inferred Performance

# Interference Classification

|  | L1-i $ | LLC | Mem bw | CPU Int | I/O bw | $\cdots$ | Net bw |
|---|---|---|---|---|---|---|---|
| App A | 95 | 81 | 7 | 56 | 43 | $\ldots$ | 100 |
| App B | 92 | 4 | 14 | 18 | 81 | $\ldots$ | 78 |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| App N | 45 | 49 | 56 | 11 | 99 | $\ldots$ | 54 |

Profiled Sensitivity

Inferred Sensitivity

# Measuring Interference Sensitivity

- Cross-application profiling: infeasible

- Measuring in hardware: platform-dependent & inaccurate

- **iBench**[1]: set of microbenchmarks of tunable intensity



Increase intensity until the application violates QoS (tolerated interference)

Generated interference?

[1] C. Delimitrou and C. Kozyrakis. "iBench: Quantifying Interference for Datacenter Applications" [IISWC'13]

# Why SVD?

SVD+SGD: Low reconstruction error
Simple, fast, scalable ($O(\min(m^2n, n^2m))$)
Offer insight on similarities

Apps that benefit from high CPU frequency

Apps similar in I-cache are also similar in branch behavior

Recommend accounts to follow

Refactor parts of app for efficiency

Low CPU

High LLC

Similar to streaming apps

# Greedy Resource Selection

- Select servers that:
  - Can tolerate the interference of new application
  - Generate interference the new application can tolerate
  - Have appropriate platform configuration

# Evaluation

- 1,000 EC2 servers
  - 14 different server configurations
  - 2 vCPU to 16 vCPU instances

- 5,000 applications
  - SPEC, PARSEC, SPLASH-2, BioParallel, Minebench, SpecWeb, Hadoop benchmarks

- <u>Objectives</u>:
  - High application performance
  - High resource utilization

# Validation

- 1,000 servers
- 5,000 applications
- Start with zero knowledge

3.5%

| Classification Engine | Metric | Applications (%) | | |
|---|---|---|---|---|
| | | CPU-bound | Memory-bound | I/O-bound |
| Heterogeneity | Avg estimation error | 3.1% | 3.6% | 4.1% |
| Interference | Avg estimation error | 3.7% | 3.5% | 5.1% |

3.8%

# Evaluation: Performance



- Least loaded scheduler (common practice today)
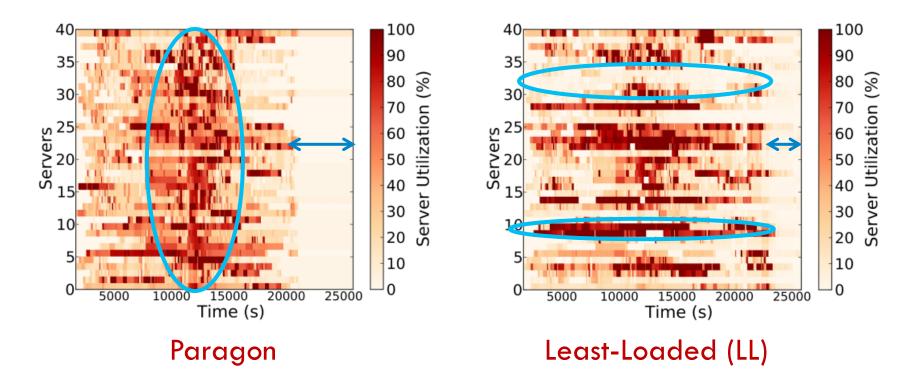  - Violates QoS for 97% of workloads

# Evaluation: Performance



- Paragon preserves QoS for 71% of workloads

- Bounds degradation to less than 10% for 90% of workloads

# Evaluation: Performance



- Paragon preserves QoS for 71% of workloads

- Bounds degradation to less than 10% for 90% of workloads
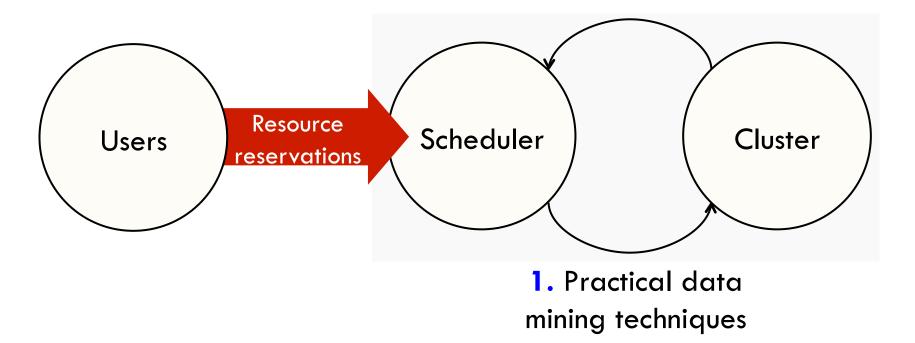
# Evaluation: System Utilization



Paragon

Least-Loaded (LL)

☐ Utilization increases from 19% to 58%

# Are We Done?

# A Larger Problem

The *user* specifies resource reservations → **overprovisioning**



1. Practical data mining techniques

# Quasar

**2.** High level interface

Resource reservations

Users

Scheduler

Cluster

**1.** Practical data mining techniques

# High-Level Interfaces

Focus on **what** performance is needed,
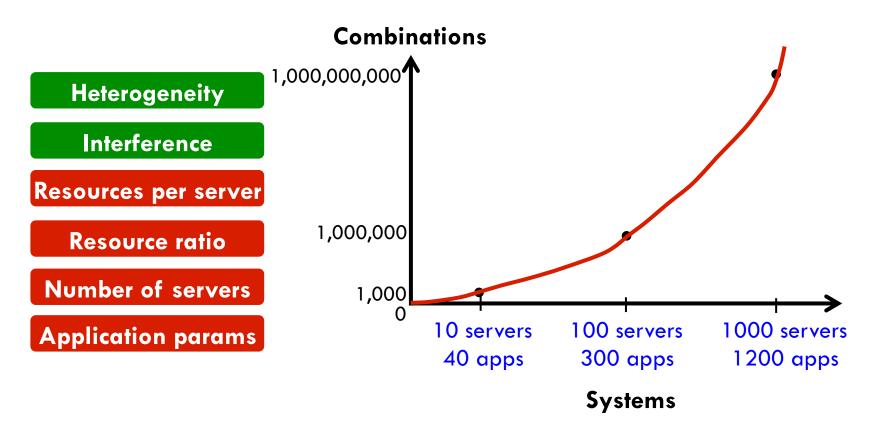not on **how** to achieve it

- Declarative interfaces:
  - SQL → describe the queries, not how they should be executed
  - DSLs → user describes program, language/compiler optimize

- Performance targets:
  - <u>Batch</u>: completion time, deadline
  - <u>Interactive</u>: throughput, tail latency
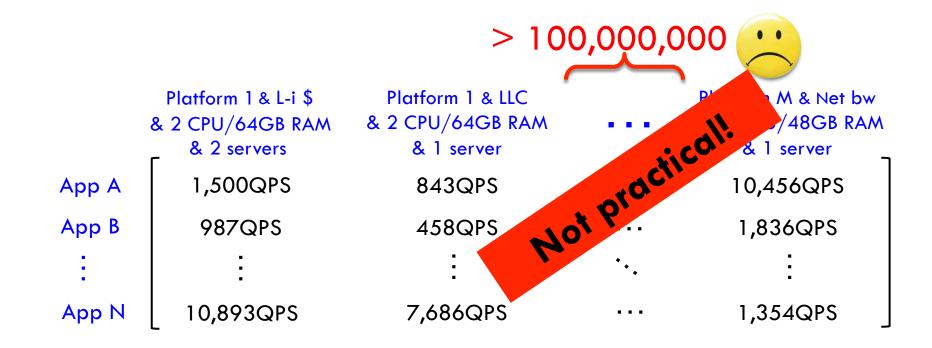
# Extracting Resource Preferences
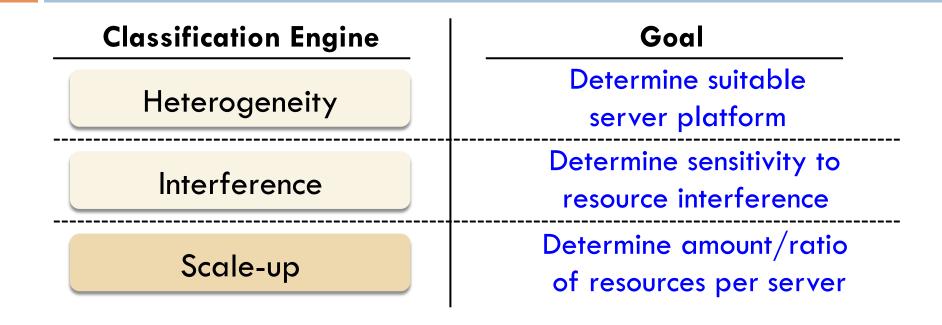
☐ Need to translate performance to resources

**Heterogeneity**

**Interference**

**Resources per server**

**Resource ratio**

**Number of servers**

**Application params**

**Combinations**

1,000,000,000

1,000,000

1,000

0

10 servers
40 apps

100 servers
300 apps

1000 servers
1200 apps

**Systems**

☐ Exhaustive characterization is infeasible

# Applying Data Mining

$> 100{,}000{,}000$ 🙁

|  | Platform 1 & L-i $ & 2 CPU/64GB RAM & 2 servers | Platform 1 & LLC & 2 CPU/64GB RAM & 1 server | ⋯ | Platform M & Net bw & /48GB RAM & 1 server |
|---|---|---|---|---|
| App A | 1,500QPS | 843QPS | | 10,456QPS |
| App B | 987QPS | 458QPS | ⋯ | 1,836QPS |
| ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| App N | 10,893QPS | 7,686QPS | ⋯ | 1,354QPS |

**Not practical!**

☐ Exhaustive classification is impractical

# Practical Resource Recommendations

| Classification Engine | Goal |
|:---:|:---:|
| Heterogeneity | Determine suitable server platform |
| Interference | Determine sensitivity to resource interference |

# Practical Resource Recommendations

| Classification Engine | Goal |
|:---:|:---:|
| Heterogeneity | Determine suitable server platform |
| Interference | Determine sensitivity to resource interference |
| Scale-up | Determine amount/ratio of resources per server |

# Practical Resource Recommendations

| Classification Engine | Goal |
|---|---|
| Heterogeneity | Determine suitable server platform |
| Interference | Determine sensitivity to resource interference |
| Scale-up | Determine amount/ratio of resources per server |
| Scale-out | Determine appropriate number of servers |

# Practical Resource Recommendations

| Classification Engine | Goal |
|:---:|:---:|
| Heterogeneity | Determine suitable server platform |
| Interference | Determine sensitivity to resource interference |
| Scale-up | Determine amount/ratio of resources per server |
| Scale-out | Determine appropriate number of servers |
| Application params | Determine appropriate settings for Hadoop, Spark, … |

64

# Quasar Overview

# Quasar Overview

# Quasar Overview

# Quasar Overview

# Quasar Overview

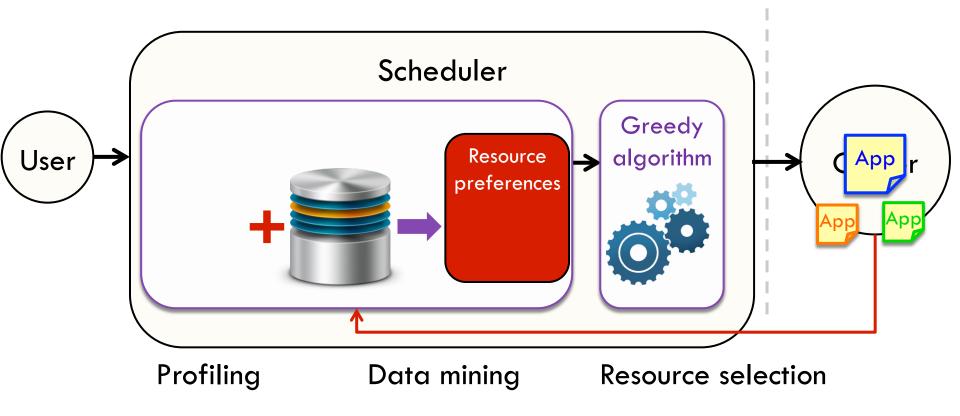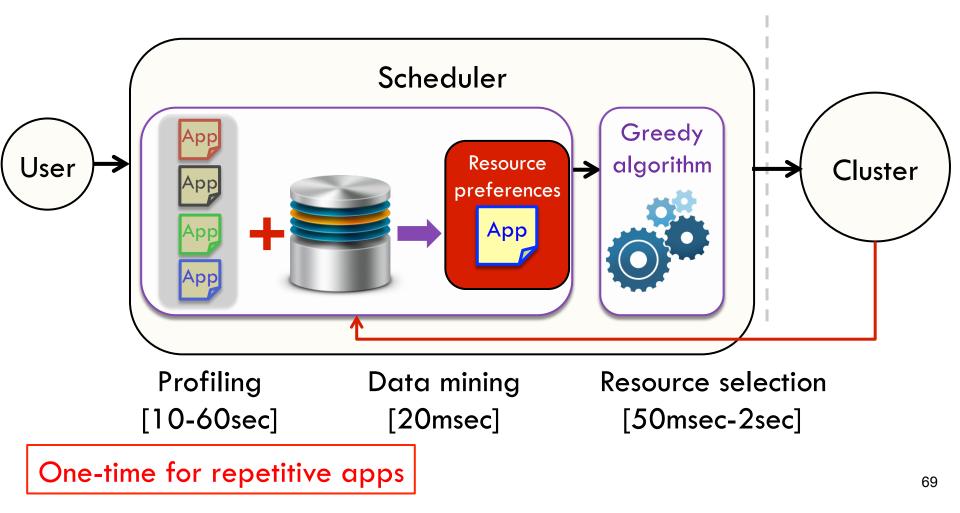# Quasar Implementation

- 10,000 loc of C++ and Python

- Runs on Linux and OS X

- Supports frameworks in C/C++, Java, Scala and Python
  - ~100-600 loc for framework-specific code

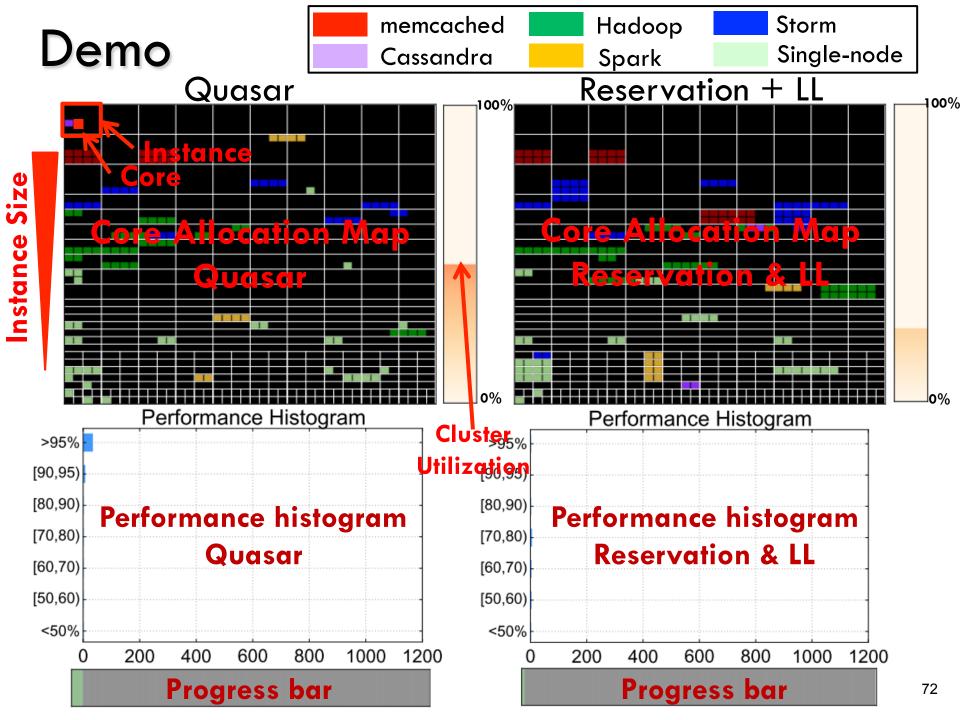- Side-effect free profiling runs with sealed containers

# Evaluation: Cloud Scenario

- Cluster
  - 200 EC2 servers, 14 different server types

- Workloads: 1,200 apps with 1sec inter-arrival rate
  - <u>Analytics</u>: Hadoop, Spark, Storm
  - <u>Latency-critical</u>: Memcached, HotCrp, Cassandra
  - <u>Single-threaded</u>: SPEC CPU2006
  - <u>Multi-threaded</u>: PARSEC, SPLASH-2, BioParallel, Specjbb
  - <u>Multiprogrammed</u>: 4-app mixes of SPEC CPU2006

- Objectives: high cluster utilization and good app QoS

# Demo

**Legend:** memcached, Hadoop, Storm, Cassandra, Spark, Single-node

## Quasar

**Instance Size**

Instance
Core

**Core Allocation Map Quasar**

## Reservation + LL

**Core Allocation Map Reservation & LL**

**Cluster Utilization**

100% ... 0%

### Performance Histogram

>95%
[90,95)
[80,90)
[70,80)
[60,70)
[50,60)
<50%

0    200    400    600    800    1000    1200

**Performance histogram Quasar**

### Performance Histogram

>95%
[90,95)
[80,90)
[70,80)
[60,70)
[50,60)
<50%

0    200    400    600    800    1000    1200

**Performance histogram Reservation & LL**

**Progress bar**

**Progress bar**

72

# Demo

# Cloud Scenario Summary

Quasar achieves:

- 91% of applications meet QoS

- ~10% overprovisioning as opposed to up to 5x

- Up to 70% cluster utilization at steady-state
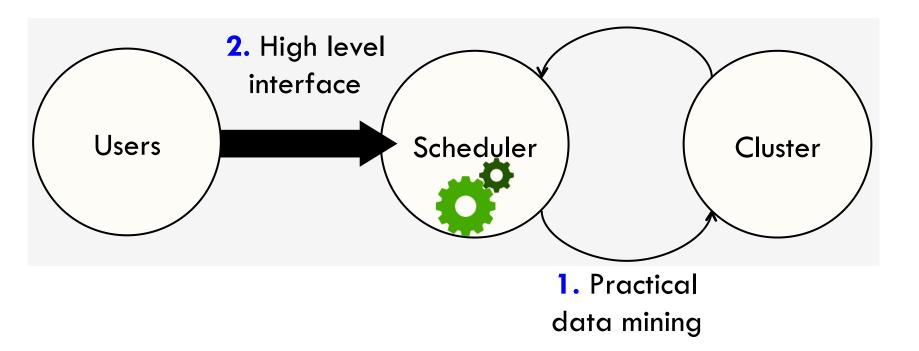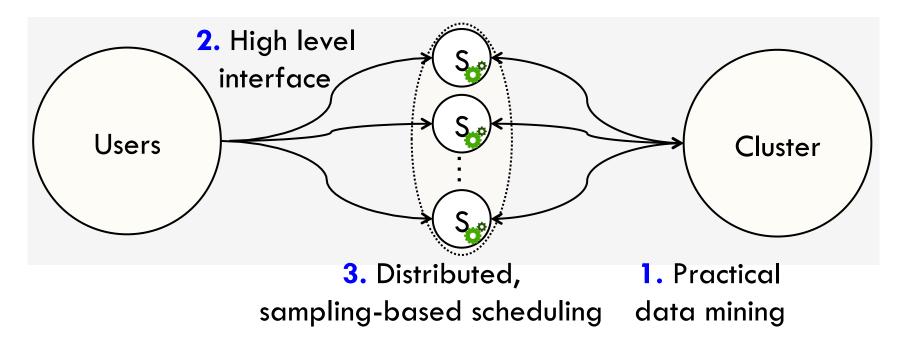
- 23% shorter scenario completion time

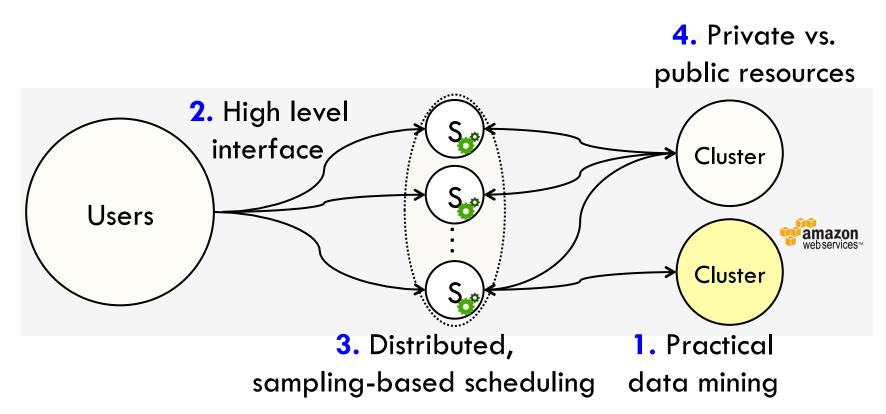# Early Adoption

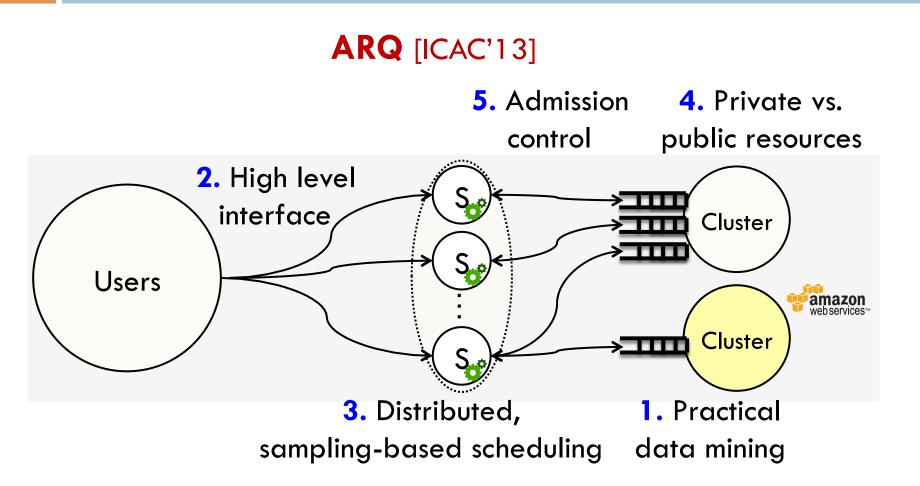https://github.com/att-innovate/charmander

# Contributions

**Quasar** [ASPLOS'14]



**2.** High level interface

Users

Scheduler

Cluster

**1.** Practical data mining

# Contributions

**Tarcil** [SOCC'15]



**2.** High level interface

Users

S

S

S

Cluster

**3.** Distributed, sampling-based scheduling

**1.** Practical data mining

# Contributions

# Contributions

**ARQ** [ICAC'13]

**5.** Admission control

**4.** Private vs. public resources

**2.** High level interface

Users

S

S

S

Cluster

Cluster

amazon
webservices™

**3.** Distributed, sampling-based scheduling

**1.** Practical data mining

# Conclusions

- **Resource efficiency**: significant challenge in systems of all scales
  - Focus on scalability of large-scale datacenters

- Cluster management: **high utilization & high app performance**
  - High-level declarative interface
  - Practical data mining techniques
  - Cross-layer design

# Questions??

- Resource efficiency: significant challenge in systems of all scales
  - Focus on scalability of large-scale datacenters

- Cluster management: high utilization & high app performance
  - High-level declarative interface
  - Practical data mining techniques
  - Cross-layer design

# Thank you!