# Zoolander: Efficiently Meeting Very Strict, Low-Latency SLOs

*Christopher Stewart* and Aniket Chakrabarti
The Ohio State University

Rean Griffith
VMWARE

# Cost Effective Scaling

- For Internet services, slow response times cost

  100-400ms delay reduces searches per session [Google '09]

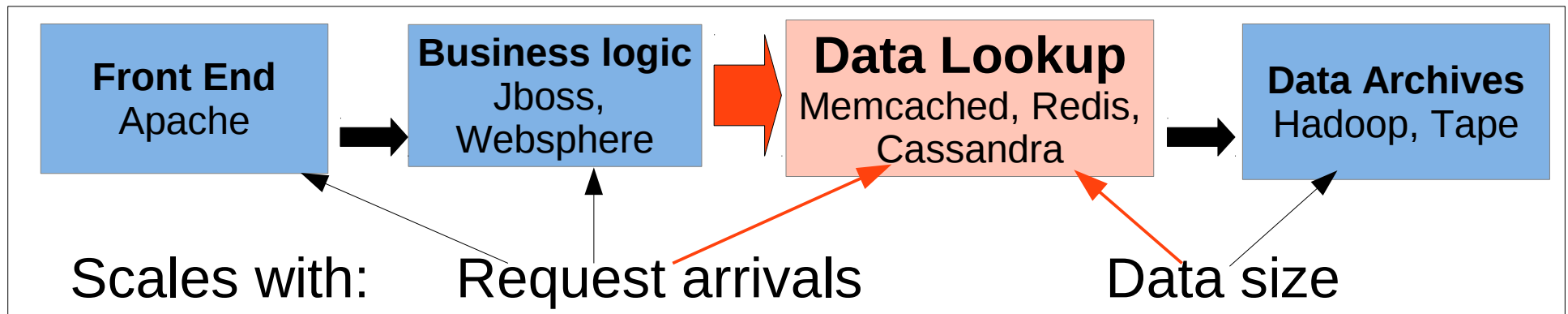  **100ms delay drops revenue by 1% [Crocker et al. '12]**

- Revenue **>>** Hardware Costs

  - To profit: Revenue > Hardware + Salaries + Benefits etc.

# Cost Effective Scaling

- As arrival rate grows, processing tiers scale out

- As data grows, data tiers scale out

| **Front End** Apache | ➡️ | **Business logic** Jboss, Websphere | ➡️ | **Data Lookup** Memcached, Redis, Cassandra | ➡️ | **Data Archives** Hadoop, Tape |

Scales with:     Request arrivals          Data size

- **In big-data era, frequent data access per request**

  – TripAdvisor: each request causes 20-40 memcached accesses [Gelfond, 2011]

  – Map-reduce services and graph processing issue  $10^3$--$10^5$

# Cost Effective Scaling

- Each user request sees 99th percentile
  - **1 slow outlier out of 100 causes 1% revenue drop**

- Service level objective: Ensure that <u>99.9%</u> of data accesses complete within <u>15ms</u>

- Traditional scale-out approaches struggle to reach such strict, low-latency SLOs
  - Slow response times cost 2.6B in lost sales (about 2% of market cap) [Flaherty,2012]
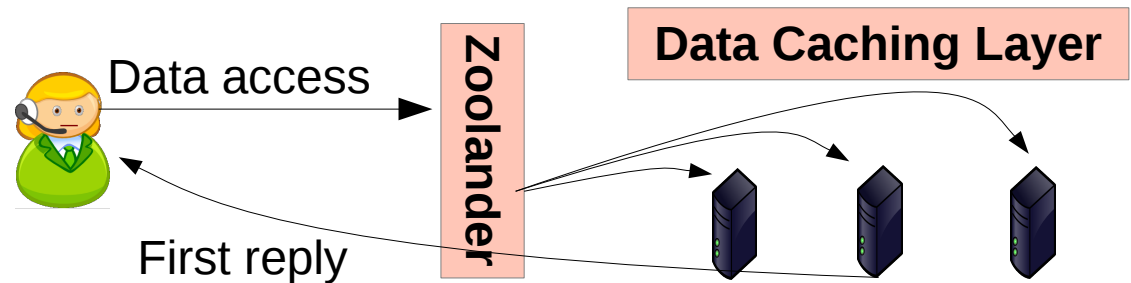
# Cost Effective Scaling

Replication for predictability is a dumb idea whose time has come --- Line borrowed from [Mogul, 2003]

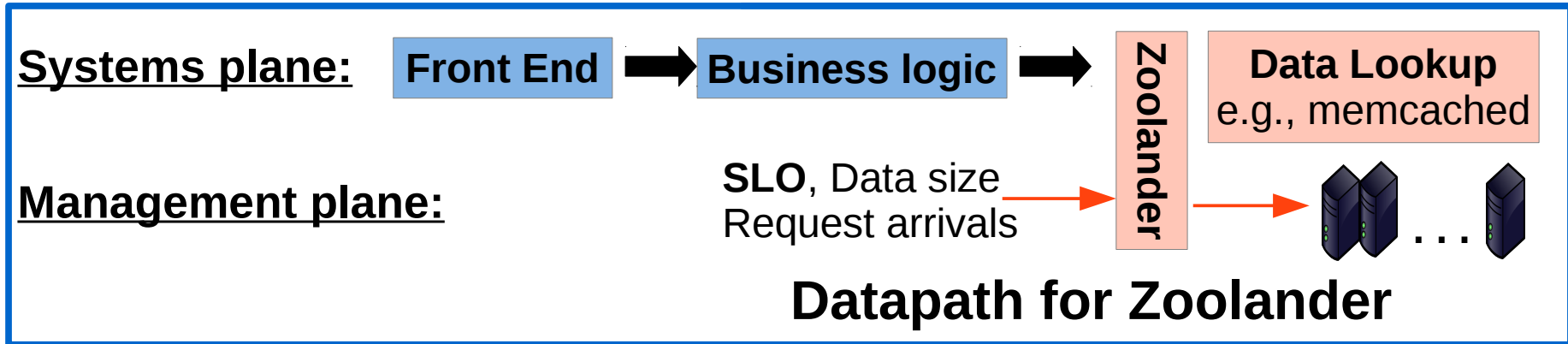**Scaling out via replication for predictability:**

**Naive Approach:**
Replicate data to D nodes
Send accesses to all D
Take first response

Data access

First reply

Zoolander

**Data Caching Layer**

- Old, dumb idea → more resources ≠ more throughput
- Time has come → more resources = stronger SLO

# Cost Effective Scaling

**Systems plane:**  | **Front End** ➡ **Business logic** ➡ | **Zoolander** | **Data Lookup** e.g., memcached

**Management plane:**

**SLO**, Data size
Request arrows

. . .

**Datapath for Zoolander**

– Zoolander is middleware for key value stores

- Meets strict SLOs efficiently using **traditional approaches and replication for predictability**

– **This talk:** Modeling and managing SLOs

- *New way to think about predictability & scale out*

# Cost Effective Scaling

- **Not this talk, but in the paper**

- **Zoolander contributes novel system designs**
  - Treat existing stores as PODS for scale out and full read/write
  - Reuses existing code & features (e.g., fault tolerance)
  - Hi-bandwidth reads reuse existing replicas for fault tolerance
  - Persistent TCP connections and fast-read bypass for low overhead
  - Support a range of consistency semantics: Causal consistency [NSDI '13], Read your own write, and eventual

# Cost Effective Scaling

- **Revived under many aliases in recent literature:**

   Replication for predictability [Trushkowsky, FAST '10]

   Cloning [Ananthanarayanan, NSDI '13; Dean, OSDI '04]

   Redundant execution [Dean & Barroso, Comm ACM '13]

- **Things they do that Zoolander doesn't:**

  – Wait for timeout and resend [Dean & Barroso, Comm ACM '13]

    - *Our model extends to this case*

- **Things Zoolander does that they don't:**

  – Scale out to D duplicates, support consistent writes, manage SLO
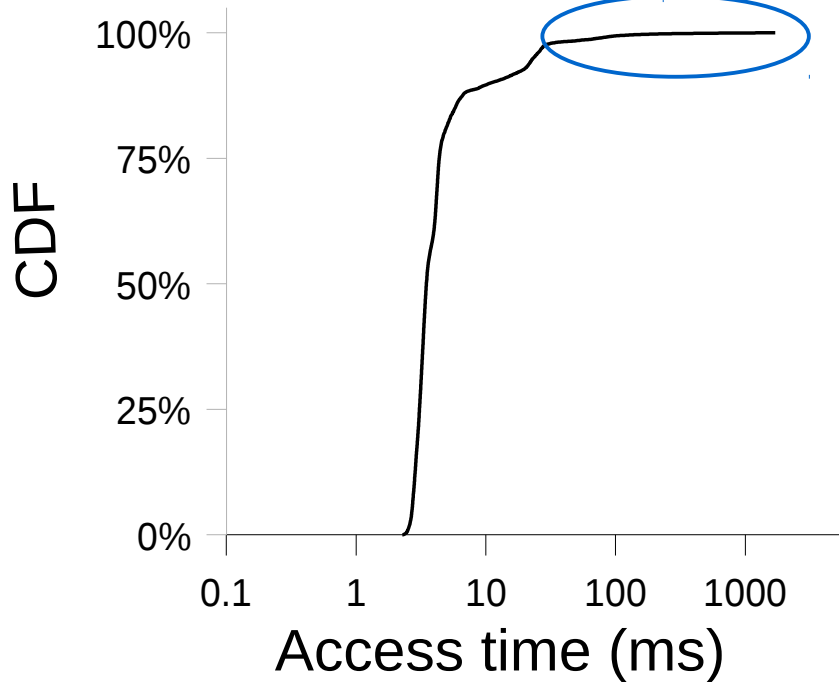
# Cost Effective Scaling

- Can we use replication for predictability to meet strict SLOs?

  - Study access-time tails in key value stores

  - Model replication for predictability on SLOs

- Should we scale out this way?

  - Model-driven study: Rep. for pred. vs Other approaches

  - Case study: Zoolander at scale

# Fat Tails in Key Values

**99th % > 99X mean**



CDF vs Access time (ms)

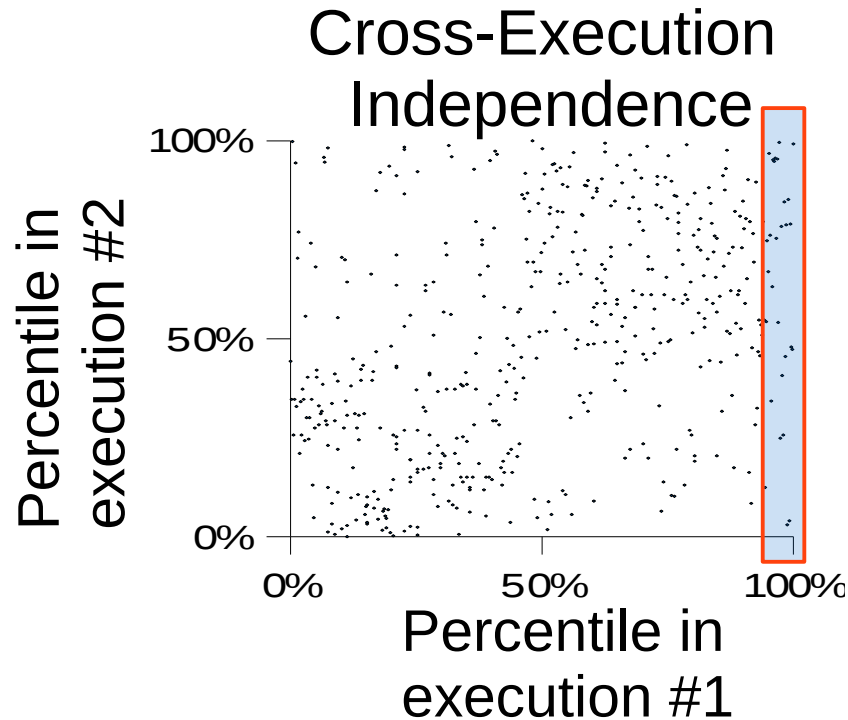3-node Zookeeper on 4 core 2.4Ghz, data size = 1 GB, 100K writes issued serially

- Fat/Heavy Tail: Outliers are way out; not captured by normal or exponential distributions

- Org. BigTable: 99.9th percentile was 31X mean [dean '12]

- Same result: memcached, Redis, Cassandra; private, EC2

- Root cause: OS, background jobs, and performance bugs

Slide 10

# Fat Tails in Key Values

## Cross-Execution Independence



Percentile in execution #2 (y-axis: 0%, 50%, 100%)
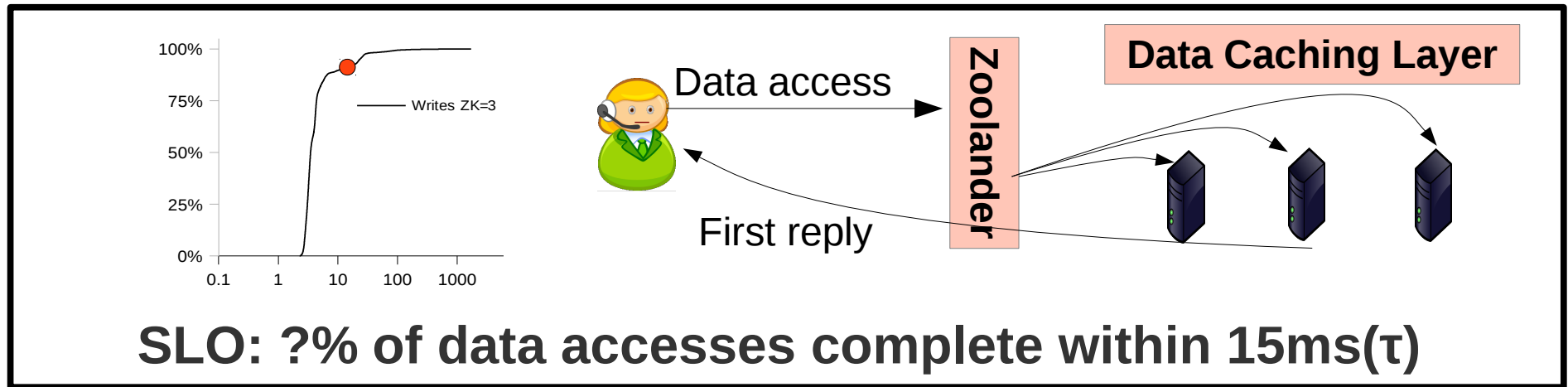
Percentile in execution #1 (x-axis: 0%, 50%, 100%)

2 Zookeeper tests performed on different servers. Requests sent in the same order for each test

- Each point reflects a request's percentile in test #1 and #2

- Almost every quartile touched; statistical independence

- In-memory key value stores

  – Extremely fast; many OS operations can cause delays

  – Other workloads → Future work

# Fat Tails in Key Values

**SLO: ?% of data accesses complete within 15ms(τ)**

- **What is the probability that first reply exceeds 15ms?**

$$(1 - \Phi(15ms)) \times (1 - \Phi(15ms)) \times (1 - \Phi(15ms))$$

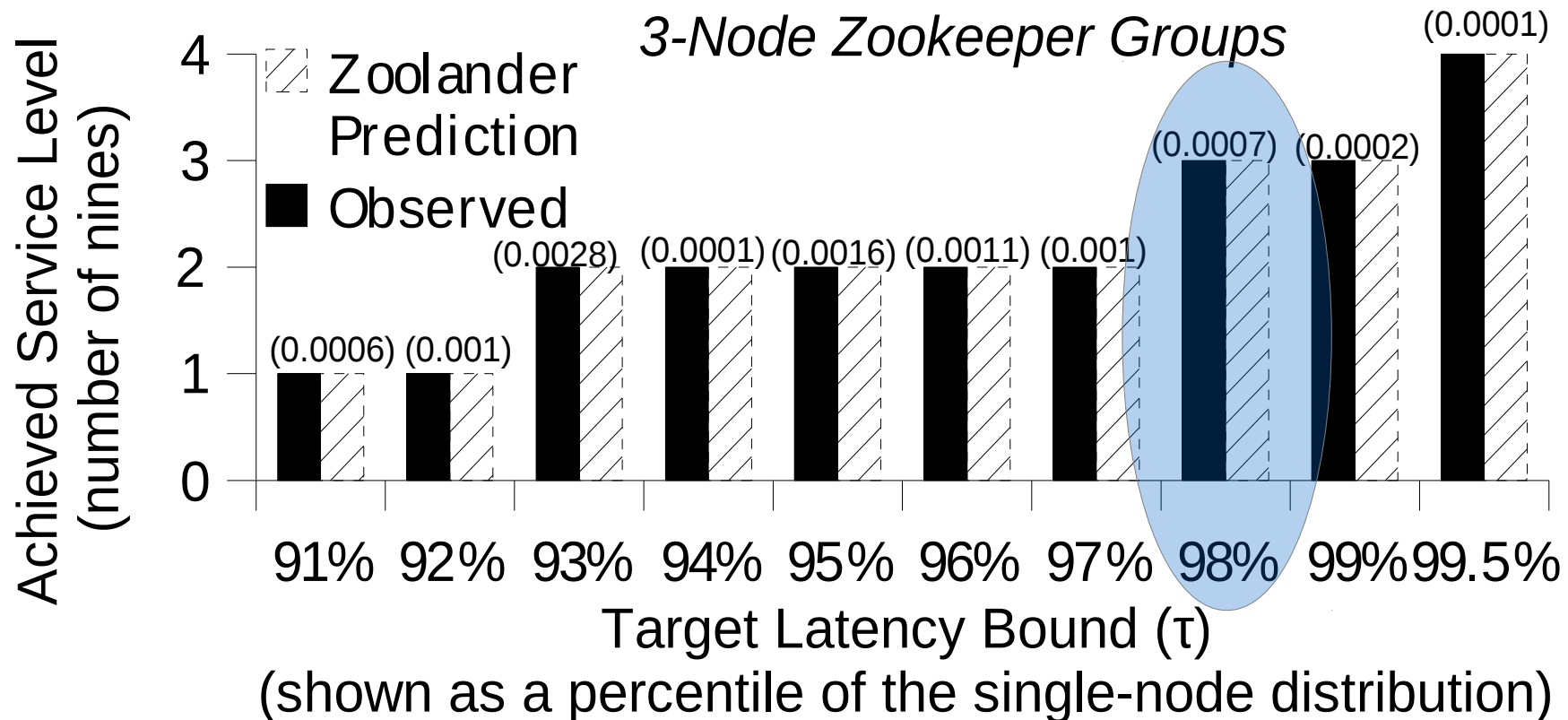$$\Phi = \text{Cumulative distribution function of access times}$$

- **At scale (D), Service Level = $1 - (1 - \Phi(\tau))^D$**

# Fat Tails in Key Values

**Core model: Service Level = 1 – (1 – Φ(τ))$^D$**

**Test #1: Is the model accurate as τ varies?**



*3-Node Zookeeper Groups*

Achieved Service Level (number of nines) vs Target Latency Bound (τ) (shown as a percentile of the single-node distribution)

Legend: Zoolander Prediction, Observed

(0.0006) (0.001) (0.0028) (0.0001) (0.0016) (0.0011) (0.001) (0.0007) (0.0002) (0.0001)

# Fat Tails in Key Values

**Core model: Service Level = 1 – (1 – Φ(τ))$^D$**

**Test #2: Is the model accurate as D varies?**

# Meeting Strict SLOs Cost Effectively
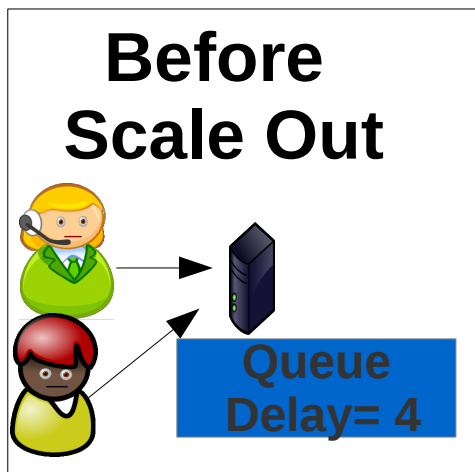
- Can we do it?
  - Study access-time tails in key value stores
  - Model replication for predictability on SLOs

- **Should we use replication for predictability to scale? Is it cost effective?**
  - Use our performance model to compare rep. for pred. against competing scale out approaches
  - Case study: Zoolander at scale
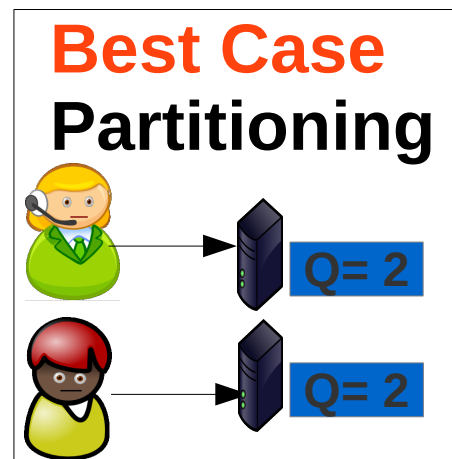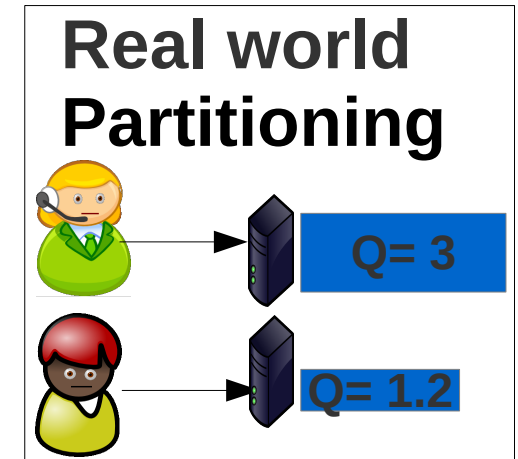
# Meeting Strict SLOs Cost Effectively

- Challenges: Duplicates share DC network and go through Zoolander

- Also, duplicates process requests at the same rate

  – Suffer the same queuing delay; Well modeled

  – Traditional scale out attacks queuing delay;"Divide the Work"



**Before Scale Out**

Queue Delay= 4

**Only scale out via rep. for pred.**

**Best Case Partitioning**

Q= 2

Q= 2

**Captured by M/G/1**

**Real world Partitioning**

Q= 3

Q= 1.2

**Hot spots, convoy, Consistency, etc.**

# Meeting Strict SLOs Cost Effectively

- **Replication for predictability affects service times; traditional "divide the work" affects queuing delay**

  - When is replication for predictability definitely better?

$$queuing\,delay = F(arrival\,rate)\,x\,service\,time$$

$$arrival\,rate = \frac{global\,arrival\,rate}{R}$$

R is number of replicas in traditional scale out

- **Post-queuing latency bound $\tau_{PQ}$ = $\tau$ - queuing delay**

Slide 17

# Meeting Strict SLOs Cost Effectively

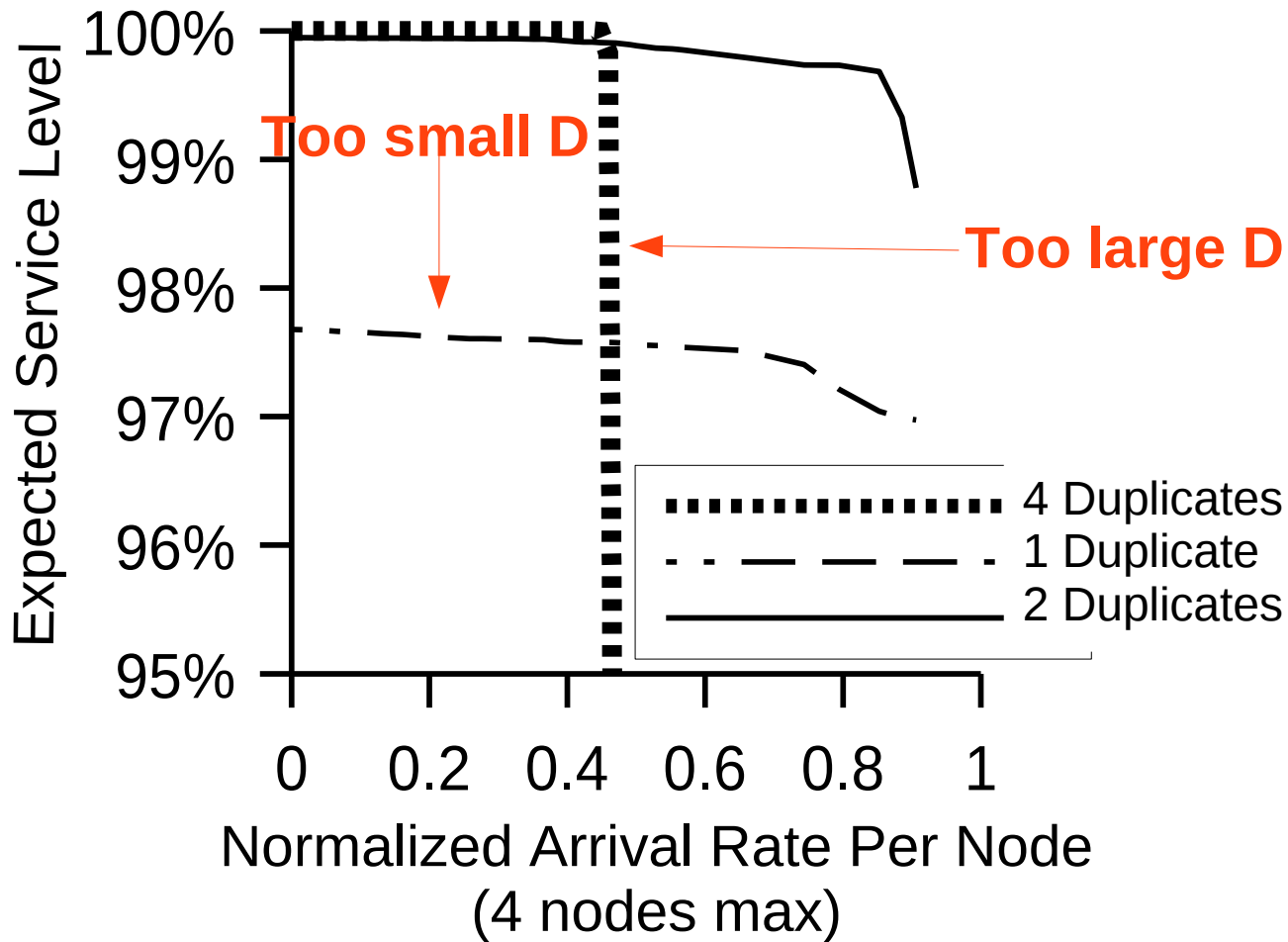**Full model: Service Level = $1 - (1 - \Phi(\tau_{PQ}))^D$**

$$\text{Service Level} = 1 - (1 - \Phi(\tau - F\left(\frac{global\ arrival\ rate}{R}\right)))^D$$

# Meeting Strict SLOs Cost Effectively

Chart: Expected Service Level (y-axis, 95% to 100%) vs. Normalized Arrival Rate Per Node (x-axis, 0 to 1, 4 nodes max)

**Too small D** / **Too large D**

Legend:
- 4 Duplicates
- 1 Duplicate
- 2 Duplicates

- *Does rep. for pred. strengthen SLOs?*
  **Yes.** Traditional scale out is limited by service time dist.

**Best approach depends on arrival rate**

Heavy arrivals per node = still a dumb idea

Moderate arrivals = Mixed strategy works well

# Meeting Strict SLOs Cost Effectively

- **Should we use replication for predictability to scale? Is it cost effective?**

  - **Case study: Zoolander at scale**

  - **Zoolander is real middleware that currently works with Zookeeper, Cassandra, Redis, and memcached**

  - **TripAdvisor released details of its memcached [Gelfond '12]**
    - **We leased 144 EC2 units to test Zoolander under TripAdvisor's scale**
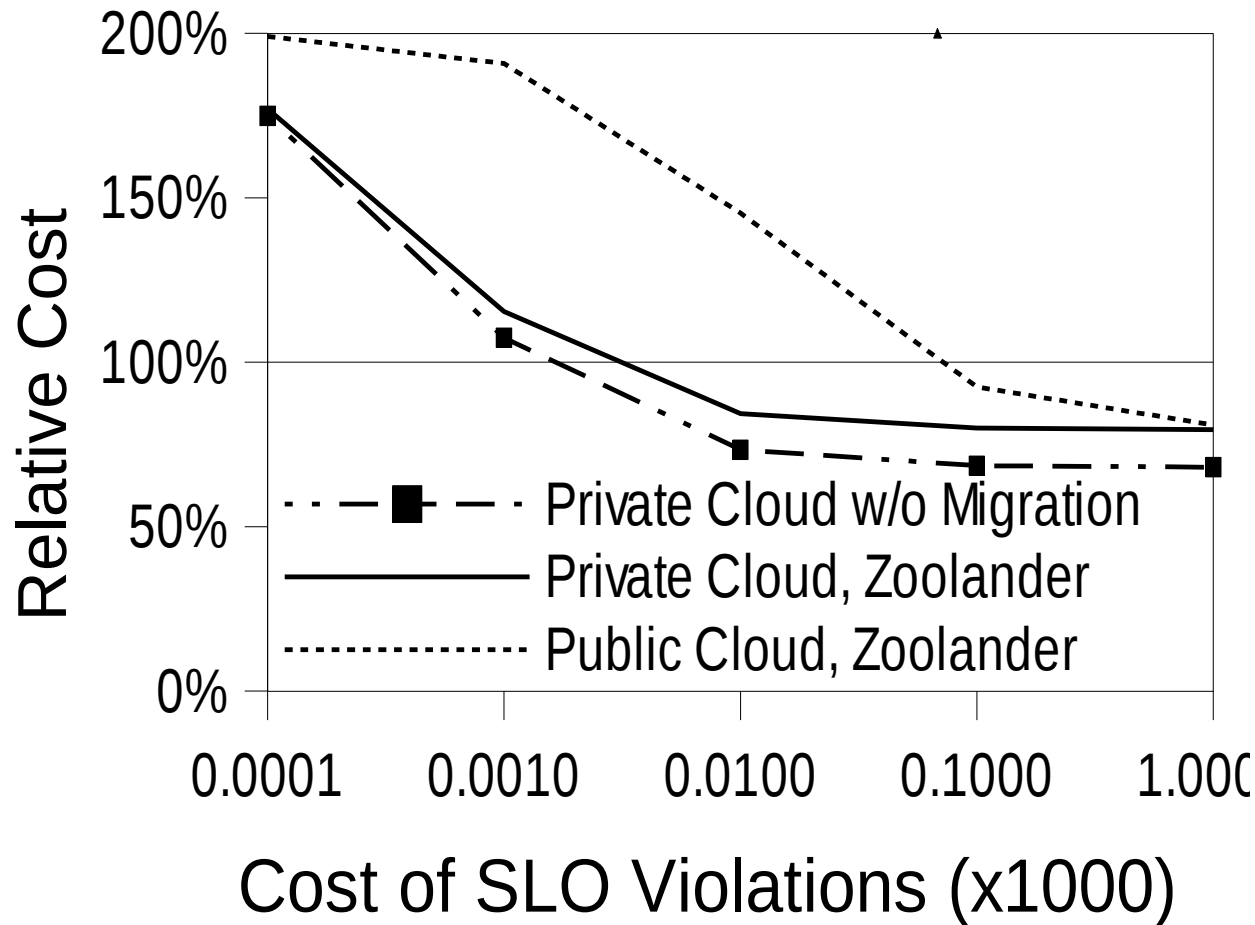
# Meeting Strict SLOs Cost Effectively

- Challenges:
    - Scale Zoolander to support 40M accesses per hour
    - Adapt Zoolander at night; accesses drop to 20M
    - **Strengthen SLO if possible—Be cost effective!**

- Competing, adaptive approaches
    - Make no changes at night
    - Turn off servers at night,
    - Replicate for predictability at night

# Meeting Strict SLOs Cost Effectively

## Service Level Objective: Ensure 20 requests complete with 150ms



**Zoolander reduced SLO violations by 32%!**

**Zoolander is cost effective for private clouds**

**EC2 favors energy saving, save energy + hardware**

Chart legend:
- Private Cloud w/o Migration
- Private Cloud, Zoolander
- Public Cloud, Zoolander

Y-axis: Relative Cost (0% to 200%)
X-axis: Cost of SLO Violations (x1000) — 0.0001, 0.0010, 0.0100, 0.1000, 1.000

# Meeting Strict SLOs Cost Effectively

1. Queuing
2. Model Driven Study
3. Zoolander at Scale

**Service Level Objective: Ensure 20 requests complete with 150ms**



Relative Cost vs. Cost of SLO Violations (x1000)

- Private Cloud w/o Migration
- Private Cloud, Zoolander
- Public Cloud, Zoolander

**TripAdvisor**
Ad revenue / Visitors * 1%
Cost = $0.068 per 1000

# Meeting Strict SLOs Cost Effectively

- **Fat tails are common, expected, and hard to remove in key-value stores**

- **Zoolander uses redundant execution to mask outlier access times and to meet SLOs cost effectively *at scale***

- **Traditional approaches *and* replication for predictability should be used for scale out.  Analytic models can capture the benefits of both!**