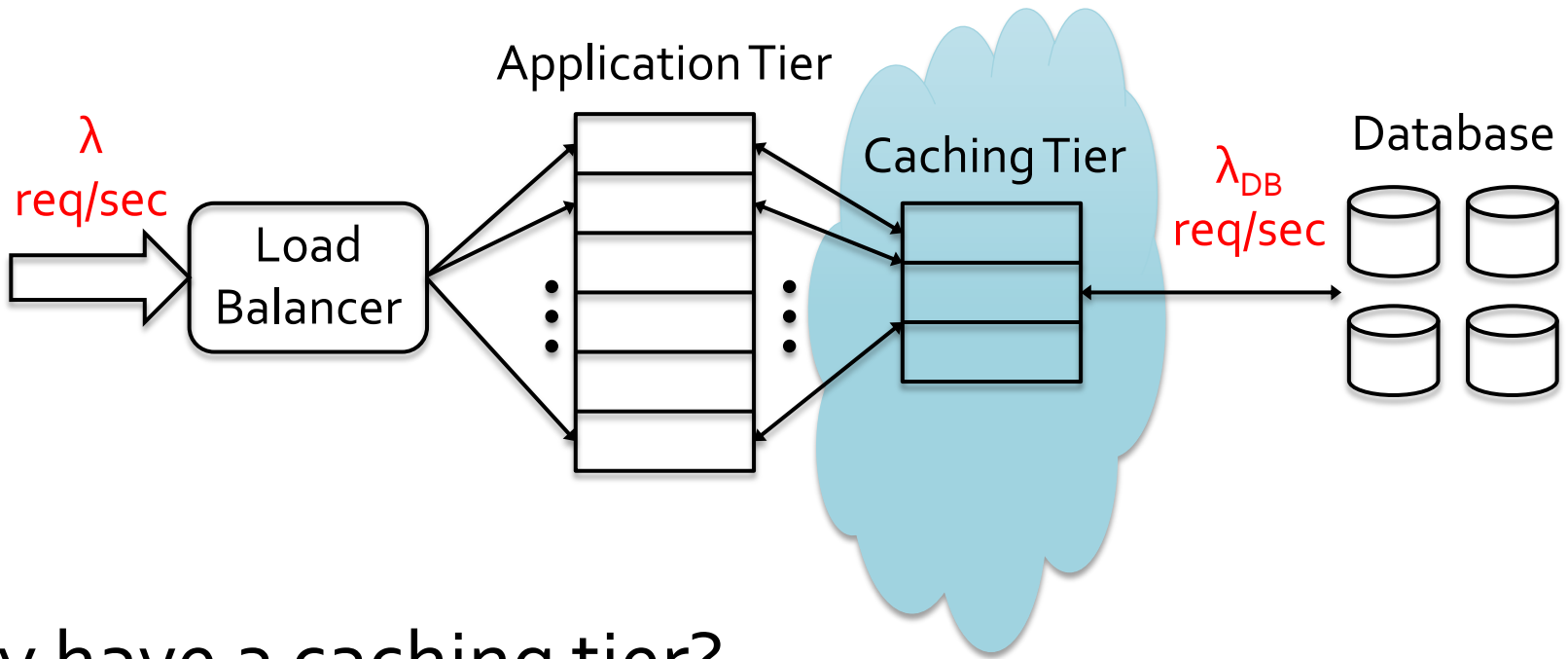# Saving Cash by Using Less (Mem)Cache

Timothy Zhu
Carnegie Mellon University

Anshul Gandhi, Mor Harchol-Balter
Carnegie Mellon University

Michael A. Kozuch
Intel Labs

1

# Application in the Cloud



Application Tier

Caching Tier

Database

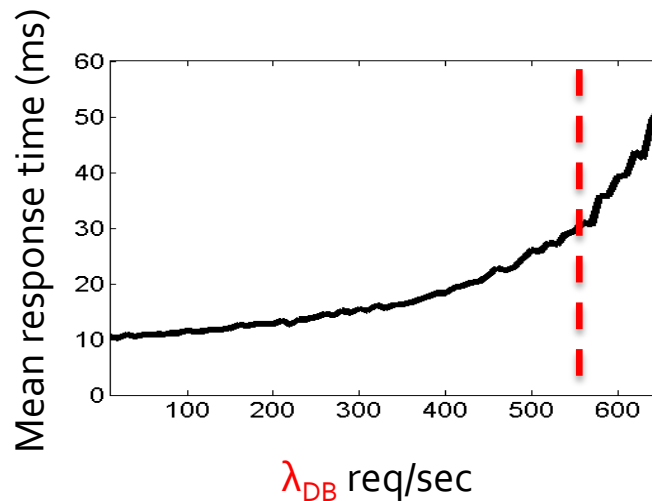$\lambda$ req/sec

Load Balancer

$\lambda_{DB}$ req/sec

Why have a caching tier?

1. Reduce database (DB) load   ($\lambda_{DB} \ll \lambda$)
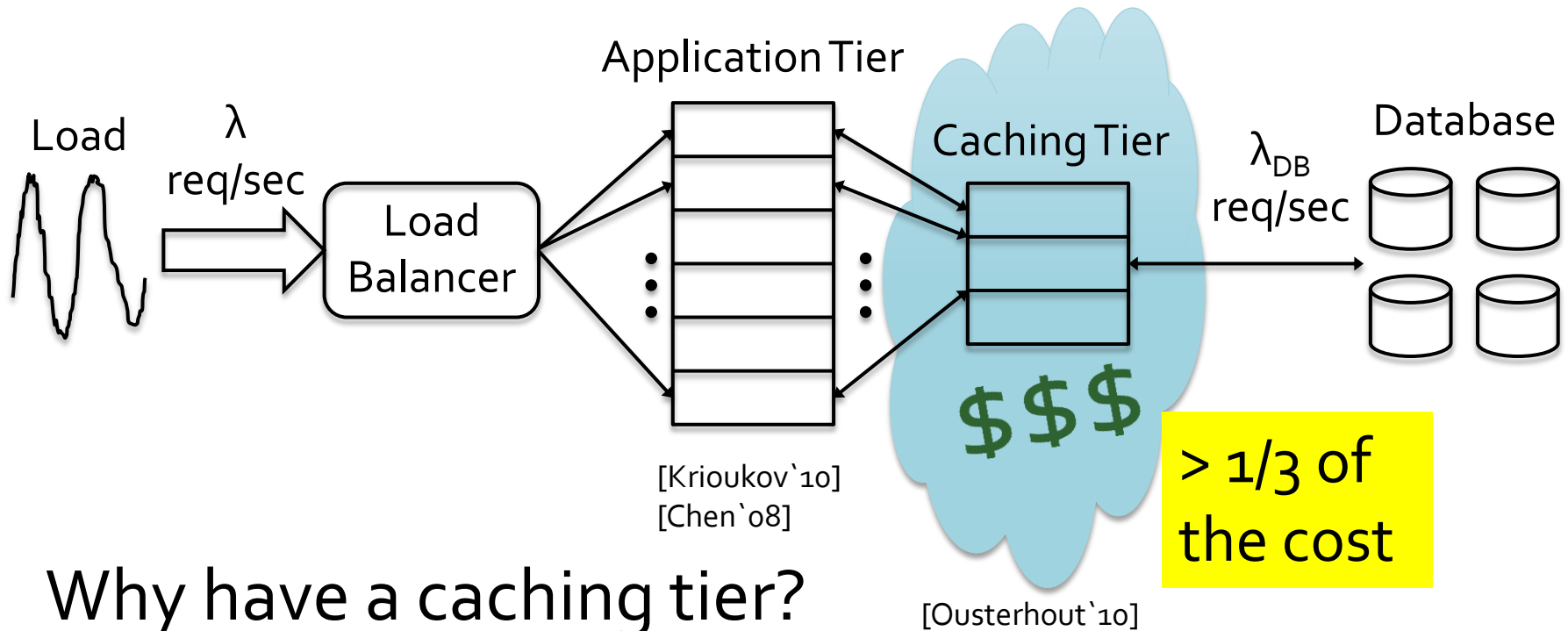
# Application in the Cloud

DB response time rapidly increases at high DB load



$\lambda_{DB}$ req/sec

Why have a caching tier?
1.    Reduce database (DB) load   $(\lambda_{DB} << \lambda)$

# Application in the Cloud

Load $\lambda$ req/sec

Application Tier

Caching Tier

$\lambda_{DB}$ req/sec

Database

Load Balancer

[Krioukov`10]
[Chen`08]

$$$

[Ousterhout`10]

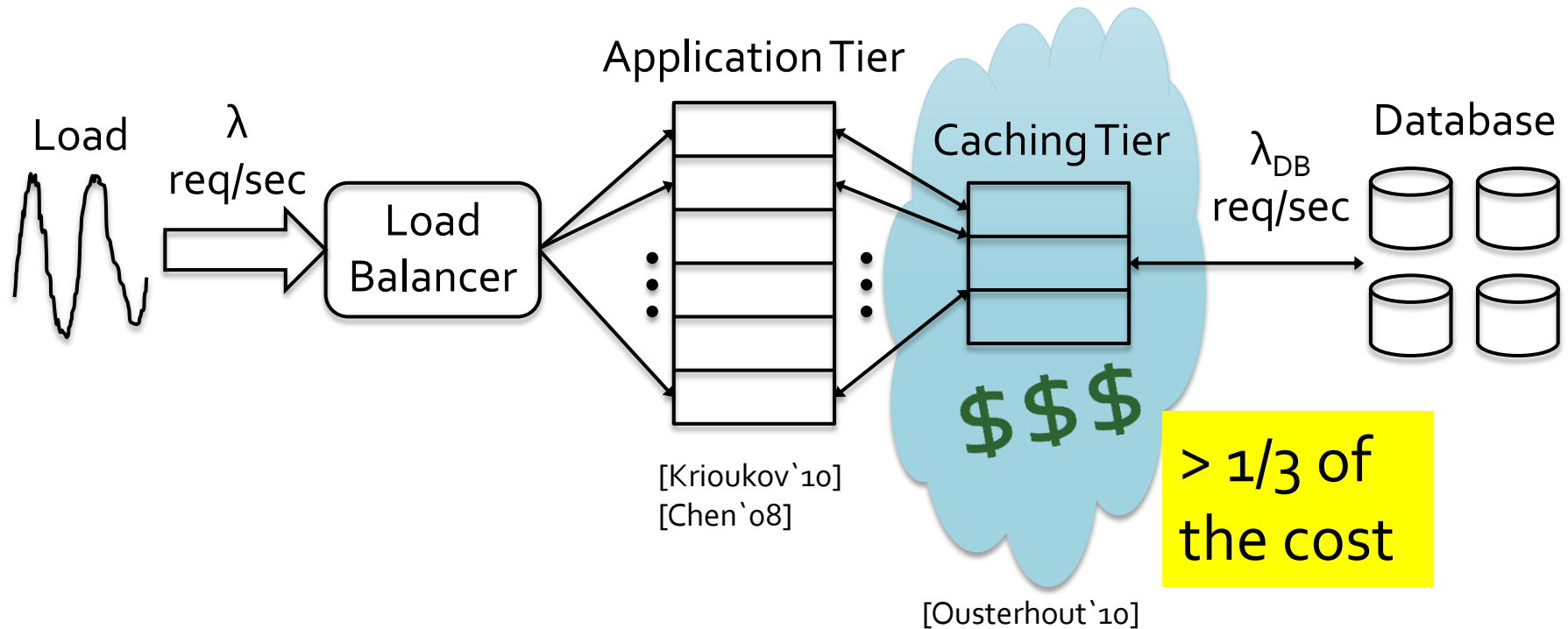> 1/3 of the cost

Why have a caching tier?
1. Reduce database (DB) load    ($\lambda_{DB} << \lambda$)
2. Reduce latency

# Application in the Cloud



Load — λ req/sec → Load Balancer → Application Tier ⋯ → Caching Tier ($$$) → $\lambda_{DB}$ req/sec → Database

[Krioukov`10]
[Chen`08]

[Ousterhout`10]

> 1/3 of the cost
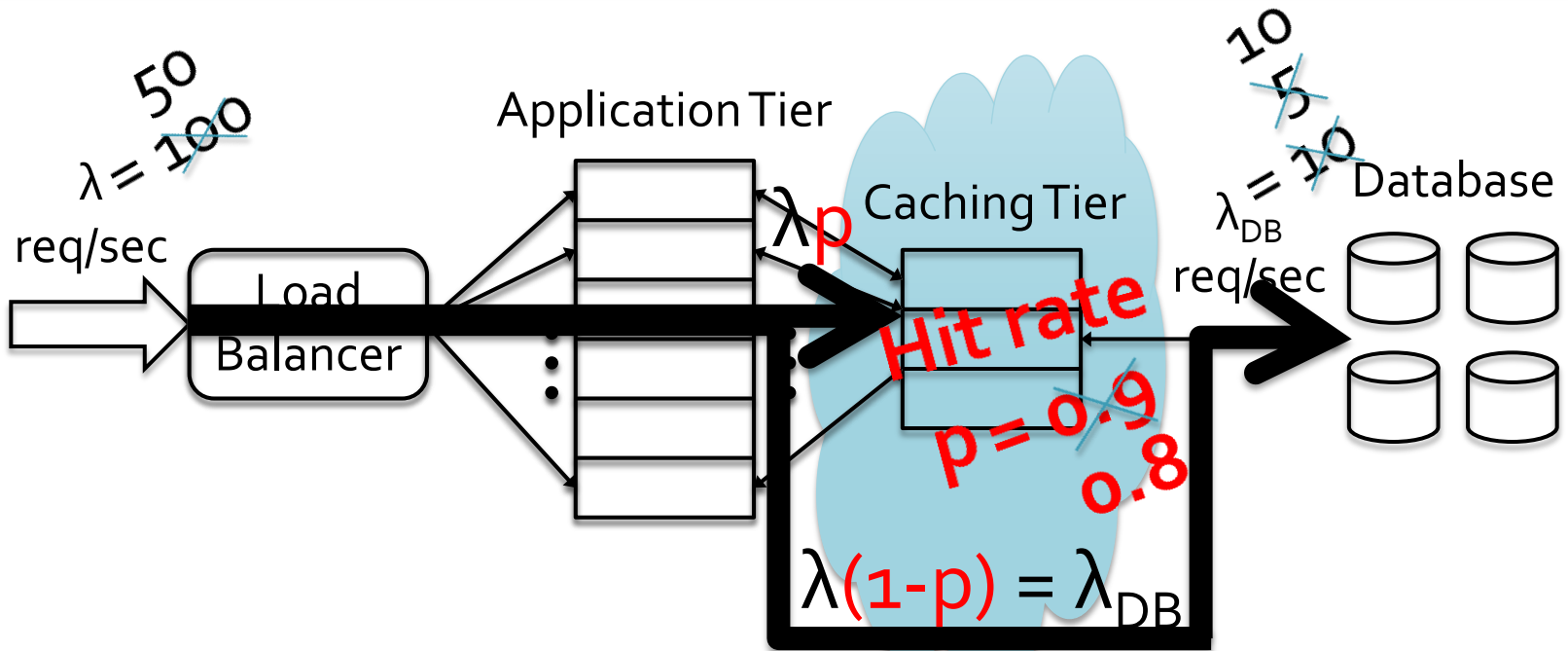
Shrink your cache during low load

# Key Questions

1. Will cache misses overwhelm the DB?

   $\lambda_{DB}$ too high?

2. Are the savings significant?
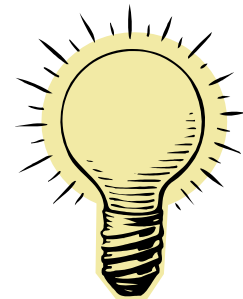
3. What about the "hot" data?

# Key Questions

1. Will cache misses overwhelm the DB?

   $\lambda_{DB}$ too high?

2. Are the savings significant?

3. What about the "hot" data?

Application Tier

Caching Tier

Database

$\lambda = \cancel{100}\ \cancel{50}$ req/sec

Load Balancer

$\lambda p$

Hit rate

$p = \cancel{0.9}\ 0.8$

$\lambda_{DB}$ req/sec

$\lambda = \cancel{10}\ \cancel{5}\ 1$

$\lambda(1-p) = \lambda_{DB}$

Goal: Keep $\lambda_{DB} = \lambda(1-p)$ low

If $\lambda$ drops ⟹ (1-p) can be higher

⟹ p can be lower

⟹ SAVE $$$

# Key Questions

1. Will cache misses overwhelm the DB?

   No, we can afford a lower hit rate at low load

2. Are the savings significant?
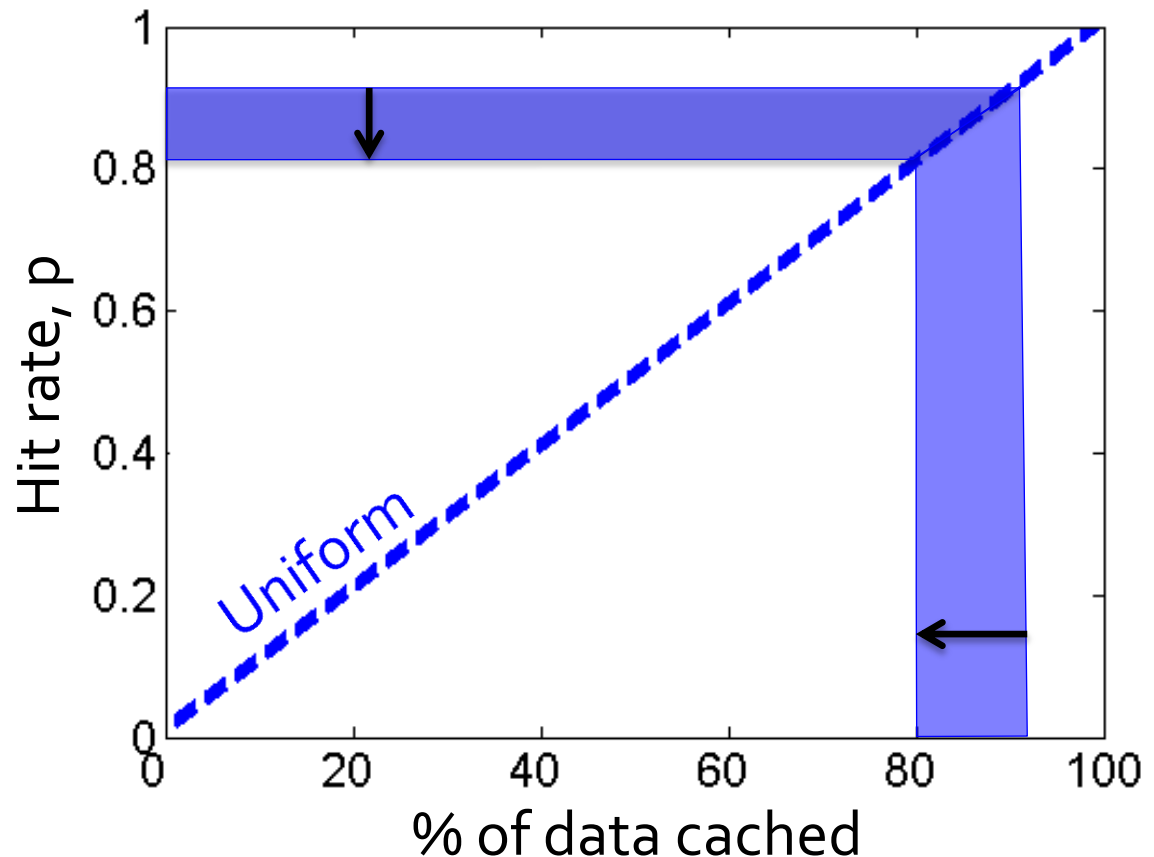
3. What about the "hot" data?

# Key Questions

1. Will cache misses overwhelm the DB?

   No, we can afford a lower hit rate at low load

2. Are the savings significant?



3. What about the "hot" data?

# Are the savings significant?

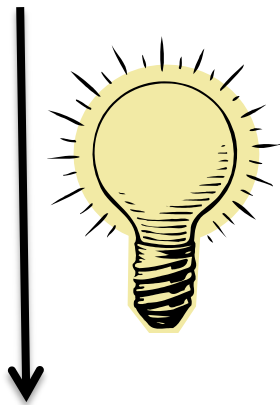- It depends on the popularity distribution

Small decrease
in hit rate

Uniform

Small decrease in
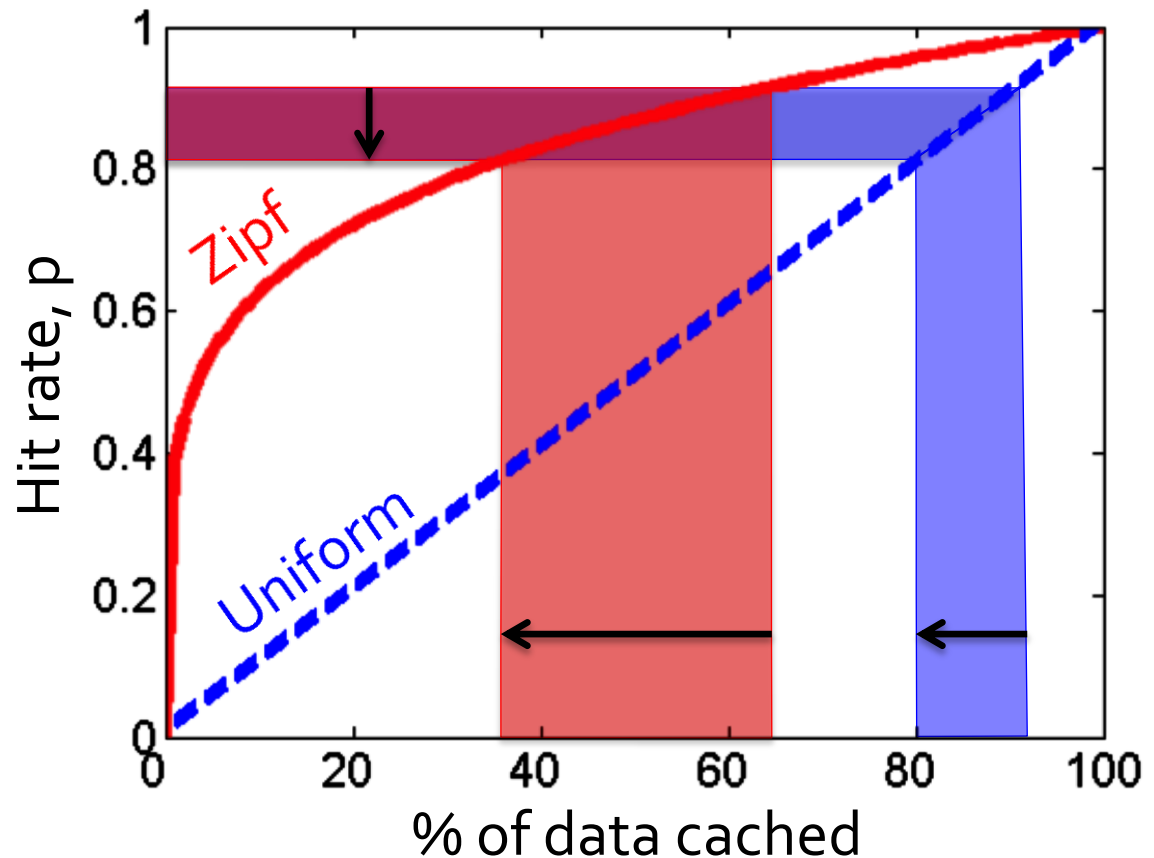caching tier size

# Are the savings significant?

- It depends on the popularity distribution
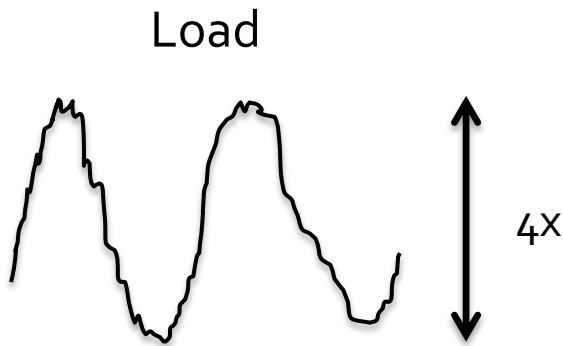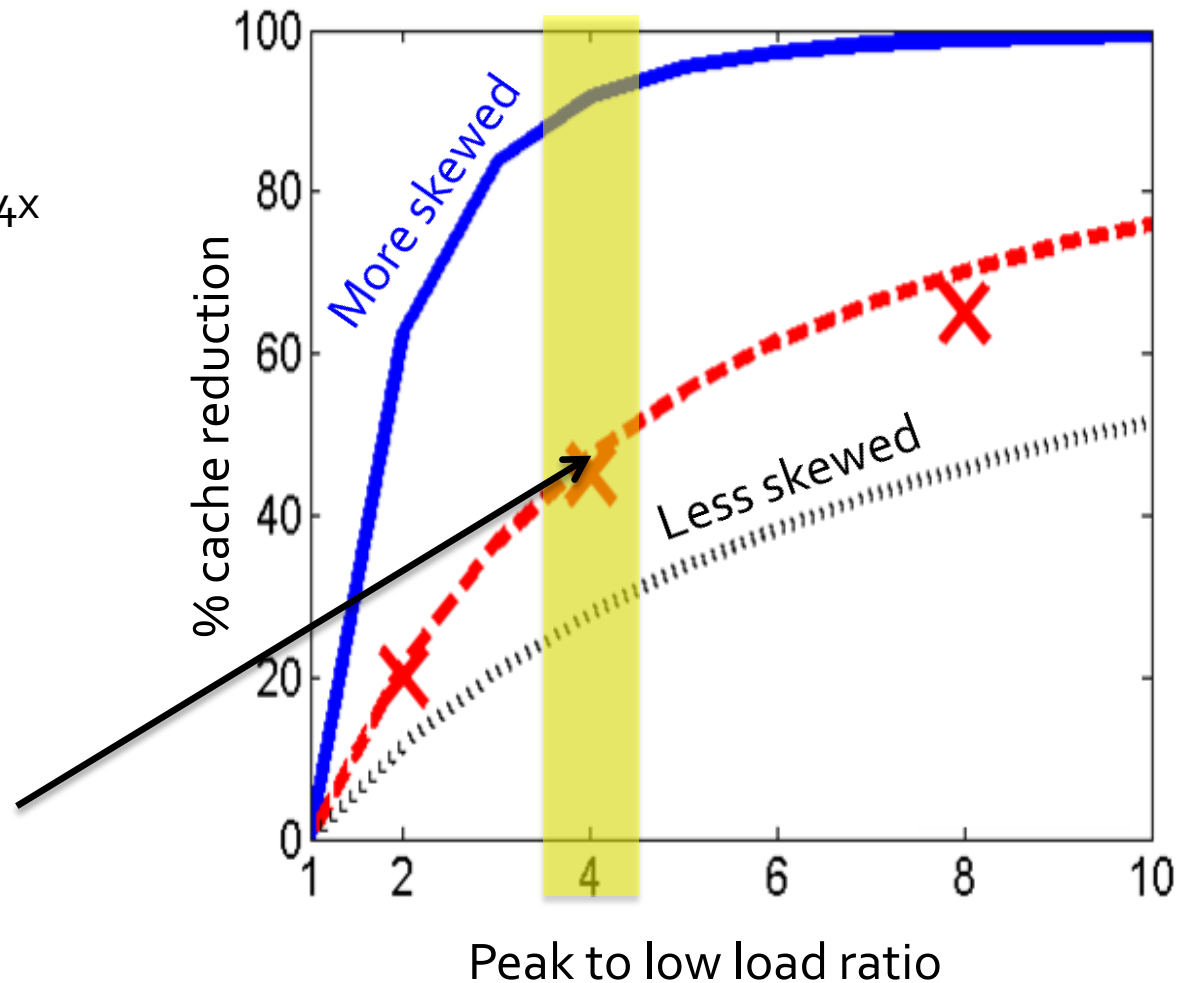
Small decrease
in hit rate

Zipf

**Large** decrease in
caching tier size

# Savings

Load



4x

% cache reduction

Peak to low load ratio

More skewed

Less skewed

50% cache savings

# Key Questions

1. Will cache misses overwhelm the DB?

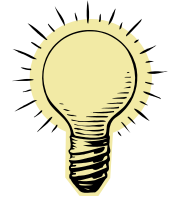   No, we can afford a lower hit rate at low load

2. Are the savings significant?

   Small decrease in hit rate → Zipf → **Large** decrease in caching tier size

3. What about the "hot" data?

# Key Questions

1. Will cache misses overwhelm the DB?

   No, we can afford a lower hit rate at low load

2. Are the savings significant?

   Small decrease in hit rate $\longrightarrow$ Zipf **Large** decrease in caching tier size
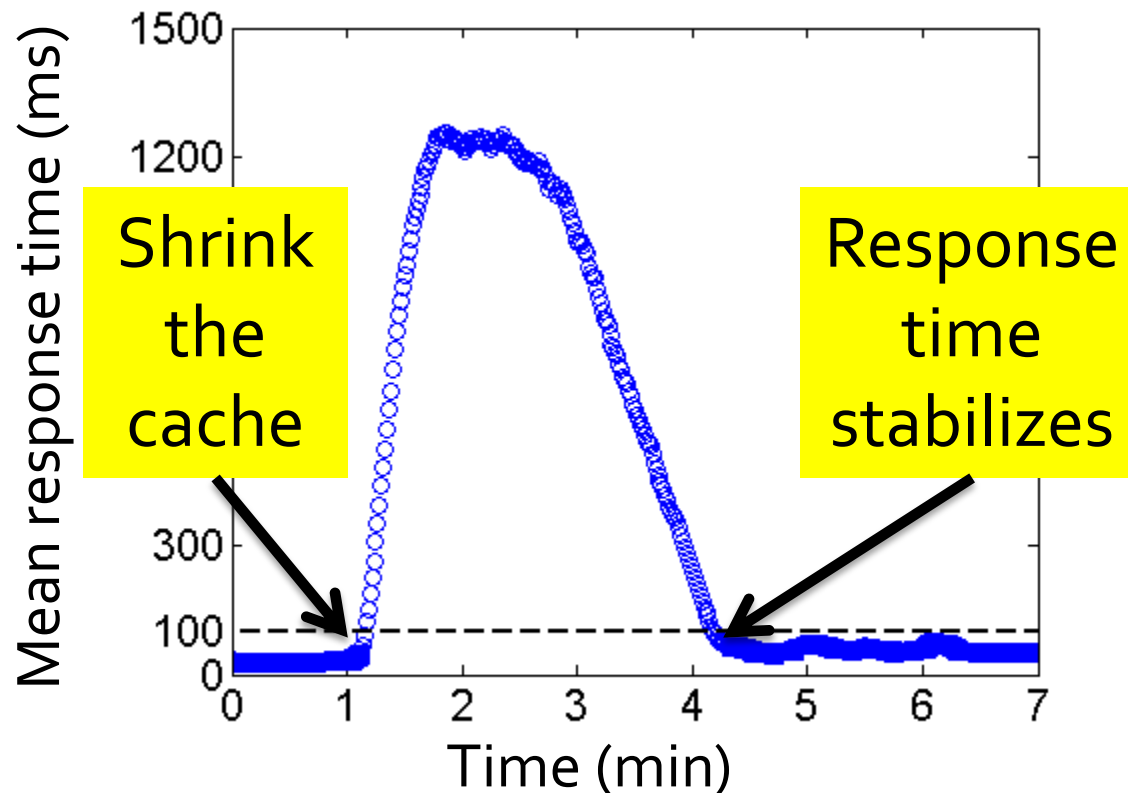
3. What about the "hot" data?

   a. Is there a problem?
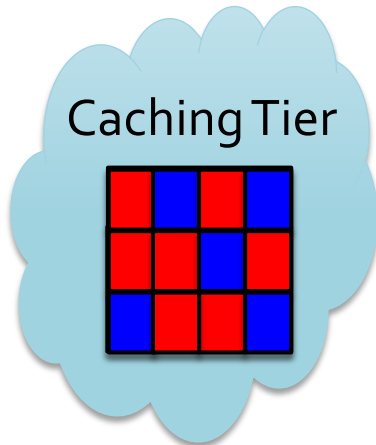
   b. What can we do about it?

# Is there a problem?

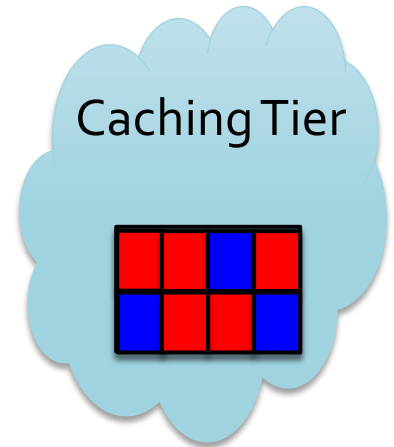- Performance can temporarily suffer if we lose a lot of hot data
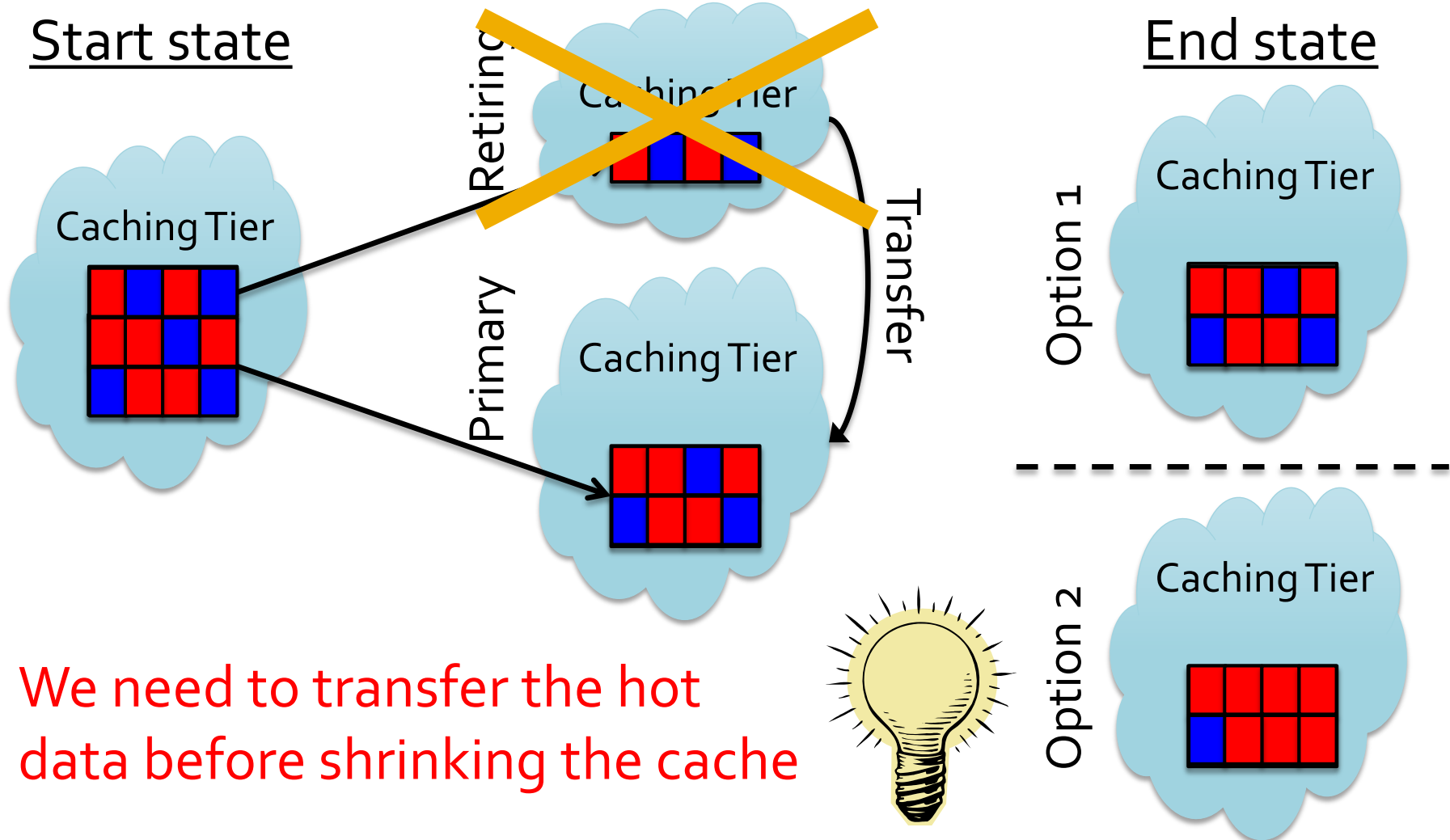
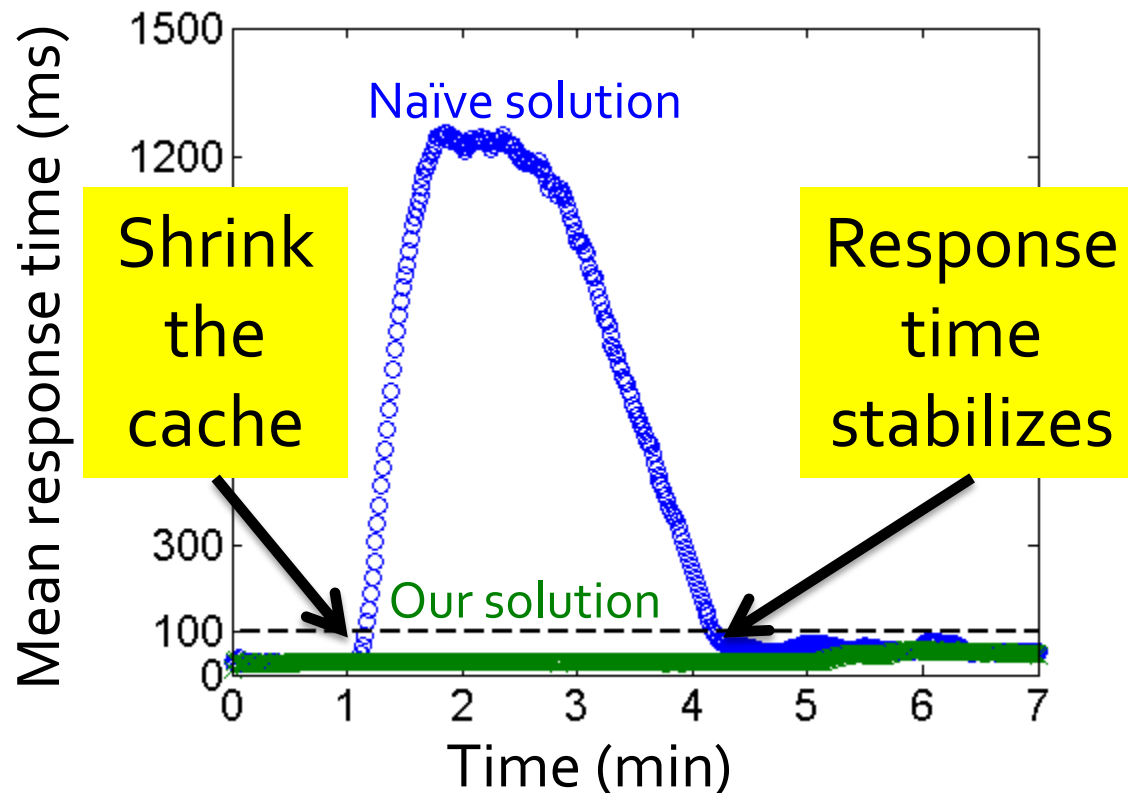# What can we do about the hot data?

Start state

End state

Caching Tier

Caching Tier

# What can we do about the hot data?



Start state

End state

Retiring

Caching Tier

Transfer

Primary

Caching Tier

Caching Tier

Caching Tier

Option 1

Caching Tier

Option 2

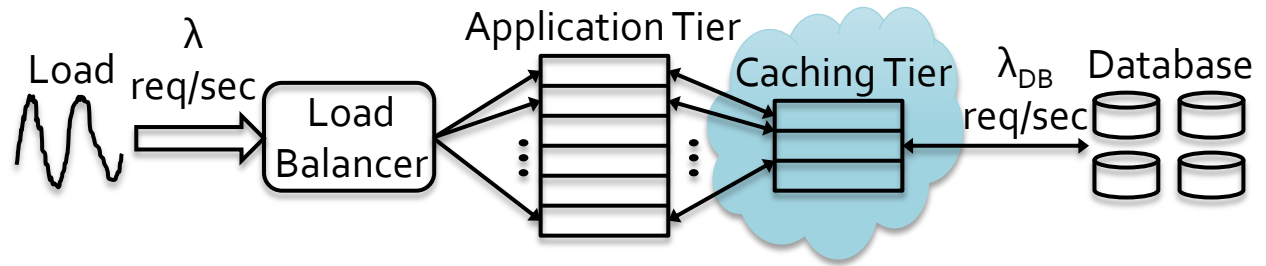We need to transfer the hot data before shrinking the cache

# Effect of transferring hot data

- Transferring the hot data before shrinking the cache eliminates performance degradation
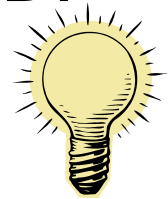
# Conclusion



1. ## Will cache misses overwhelm the DB?
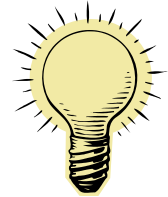
   No, we can afford a lower hit rate at low load

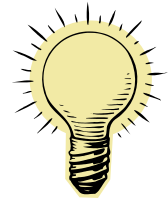2. ## Are the savings significant?

   Small decrease in hit rate → **Large** decrease in caching tier size
   
   Zipf

3. ## What about the "hot" data?

   We need to transfer the hot data before shrinking the cache

Use less cache → Save $$$
Low load