



清華大學

Tsinghua University

# Modeling and Analyzing Latency in the Memcached System

Wenxue Cheng, Fengyuan Ren, Wanchun Jiang (CSU), Tong Zhang

NNS Group @ Tsinghua University

IEEE ICDCS 2017, Atlanta, USA

# Memcached System

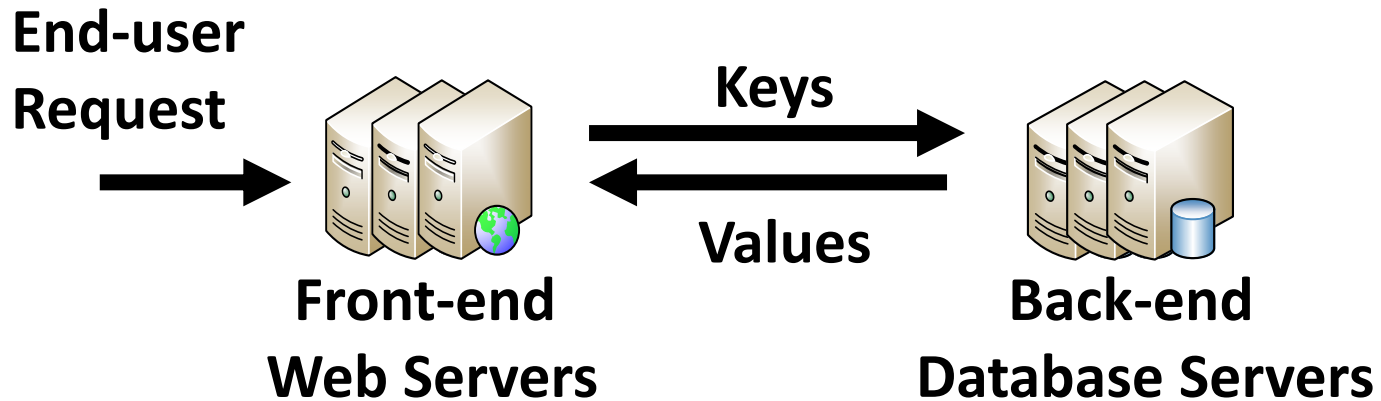
---

- Key-Value system
- In-memory caching solution
- Speed up searching applications

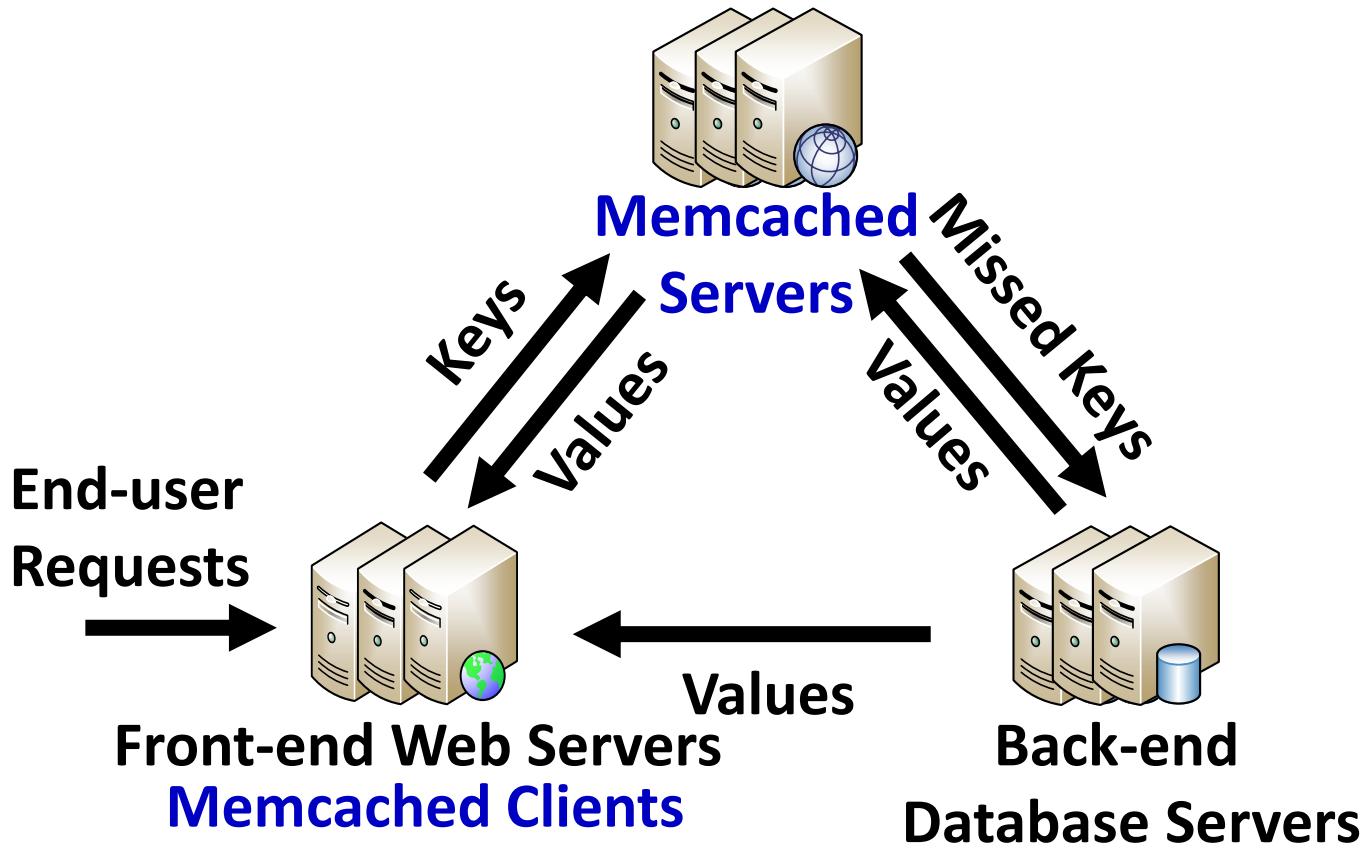
The Facebook logo, consisting of the word "facebook" in white lowercase letters on a dark blue rounded rectangular background.The Flickr logo, featuring the word "flickr" in blue lowercase letters with a red "r" on a white rounded rectangular background.The Twitter logo, showing the word "twitter" in white lowercase letters next to a white bird icon on a blue rounded rectangular background.The LiveJournal logo, with the word "LIVEJOURNAL" in white uppercase letters on a dark blue rounded rectangular background.The Wikipedia logo, with the word "WIKIPEDIA" in black uppercase letters and the tagline "The Free Encyclopedia" in smaller black text below it.The YouTube logo, featuring the word "You" in black and "Tube" in white on a red rounded rectangular background.

# Traditional Websites

---



# Memcached Architecture



**Latency** is the most pivotal performance metric.

# Factors on Latency

Load Size

Service Rate

Unbalanced Load

End-user Requests

Front-end Web Servers  
Memcached Clients

Back-end Database Servers

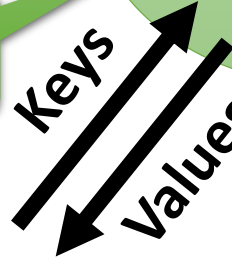
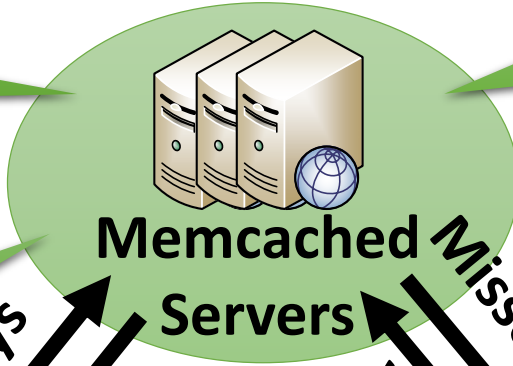
Memcached Servers

Keys

Values

Missed Keys  
Values

Values



# Factors on Latency

Load Size

Service Rate

Unbalanced Load

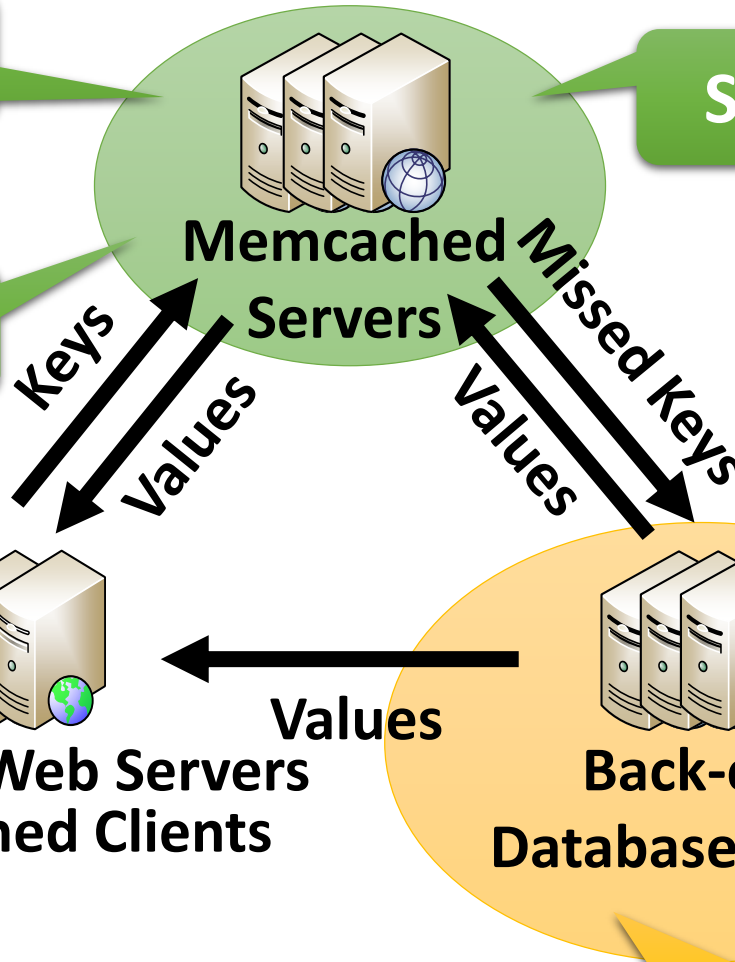
End-user Requests

Front-end Web Servers  
Memcached Clients

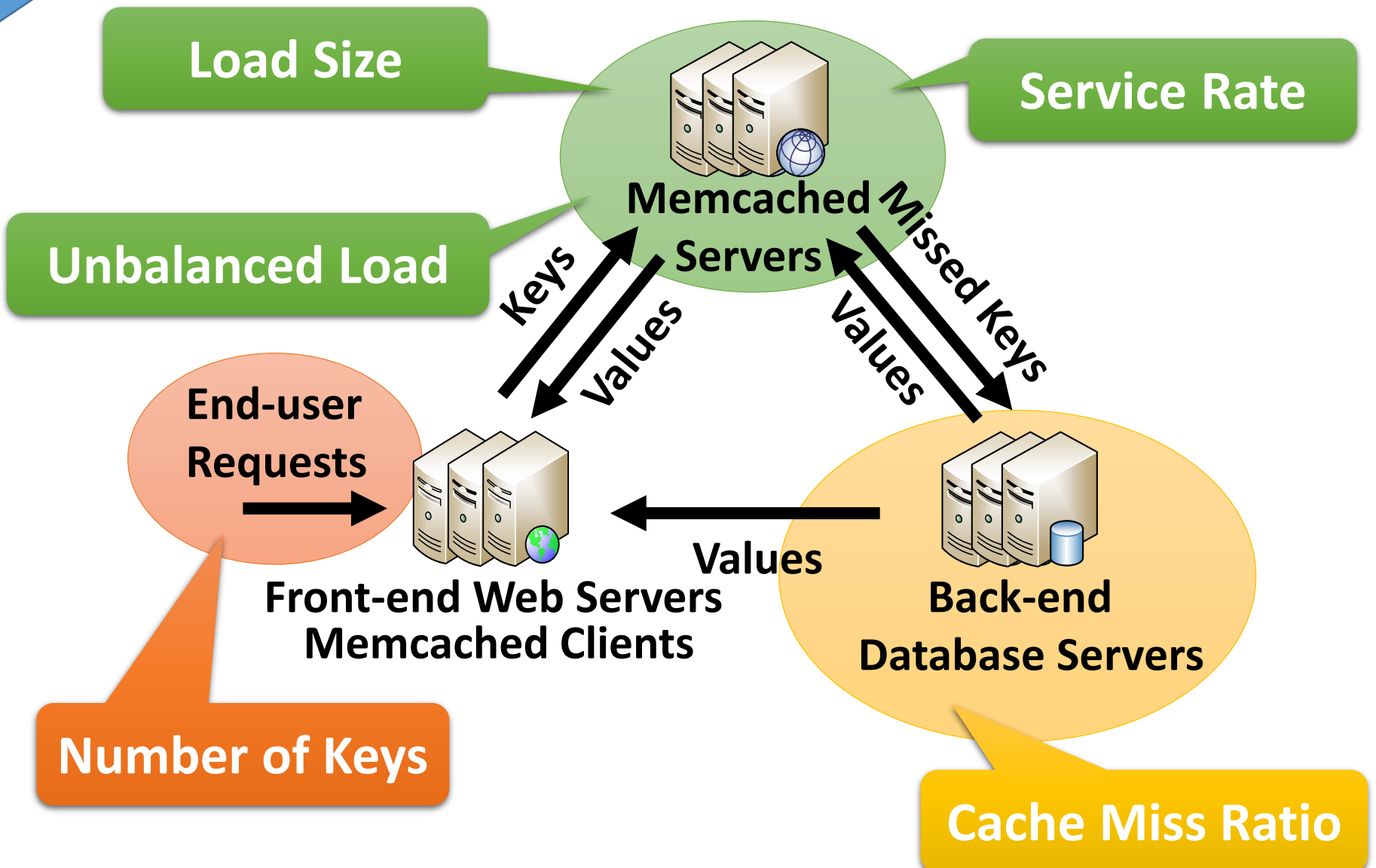
Values

Back-end  
Database Servers

Cache Miss Ratio



# Factors on Latency



# Factors on Latency

Load Size



Service Rate

- How much improvement on latency can be achieved by optimizing each factor?



Values

Front-end Web Servers  
Memcached Clients

Back-end  
Database Servers

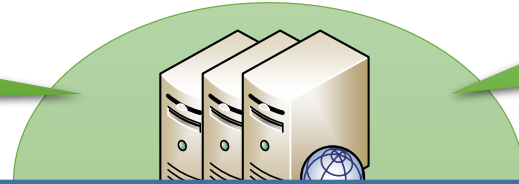
Number of Keys

Cache Miss Ratio



# Factors on Latency

Load Size



Service Rate

- How much improvement on latency can be achieved by optimizing each factor?

**Theoretical Model & Quantitative Analysis**



Values

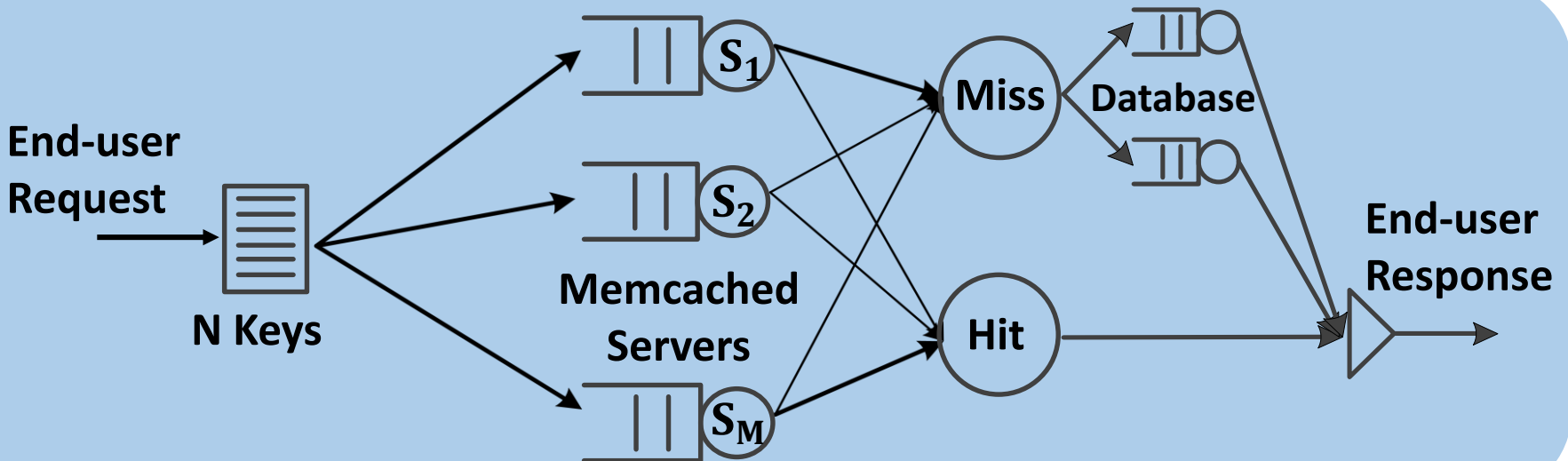
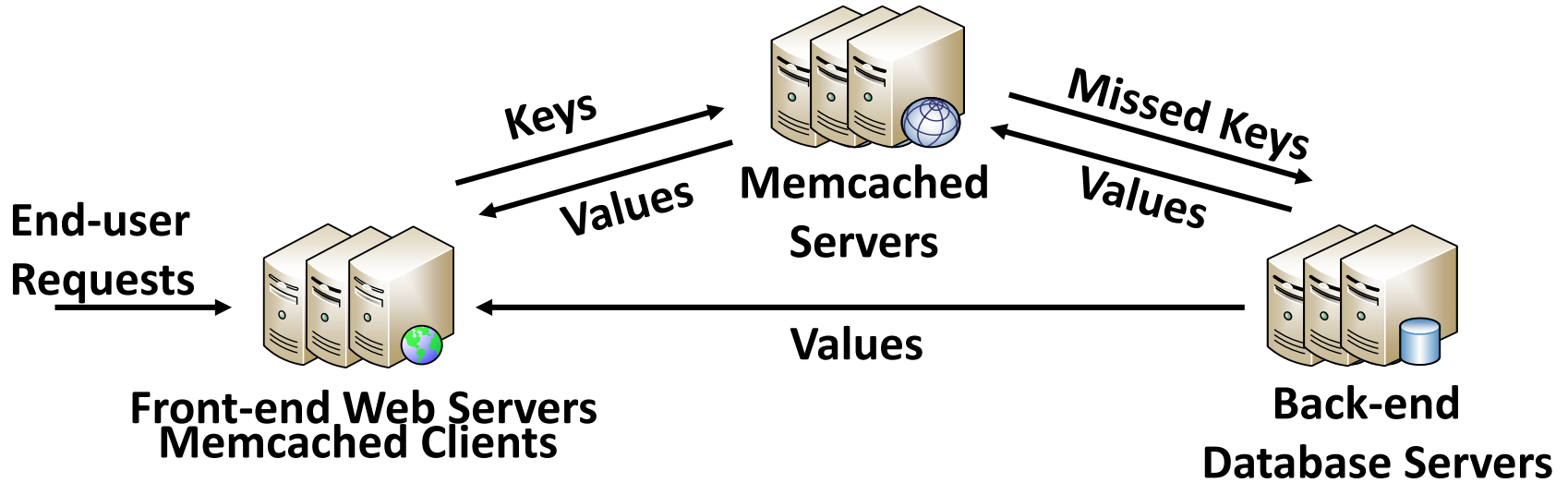
Front-end Web Servers  
Memcached Clients

Back-end  
Database Servers

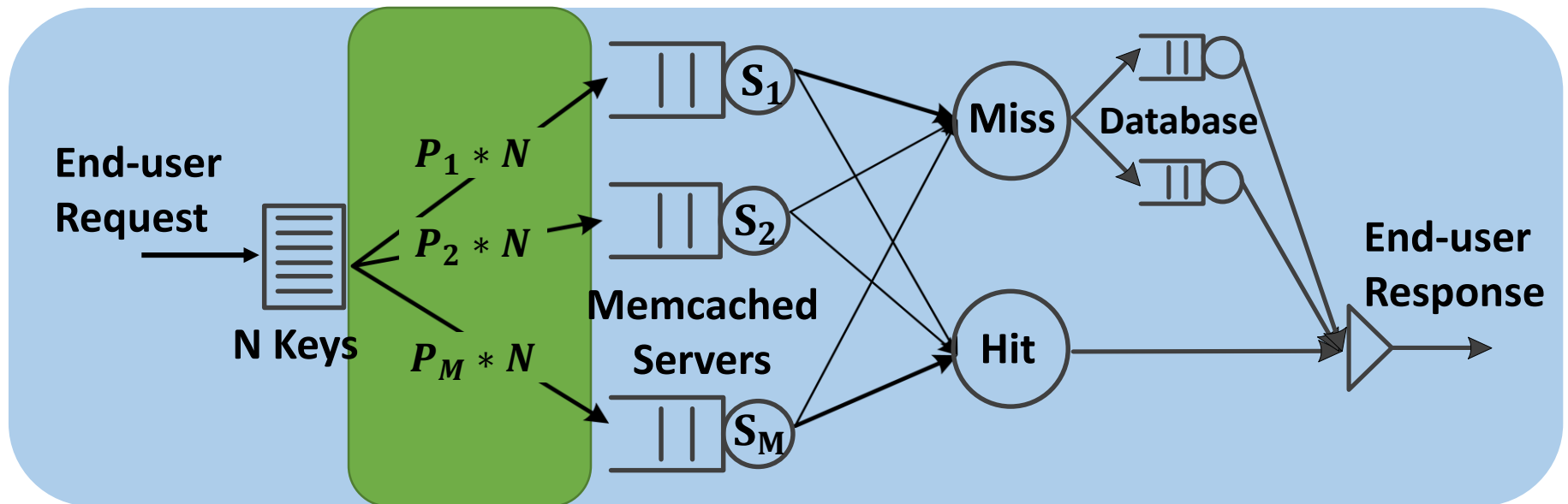
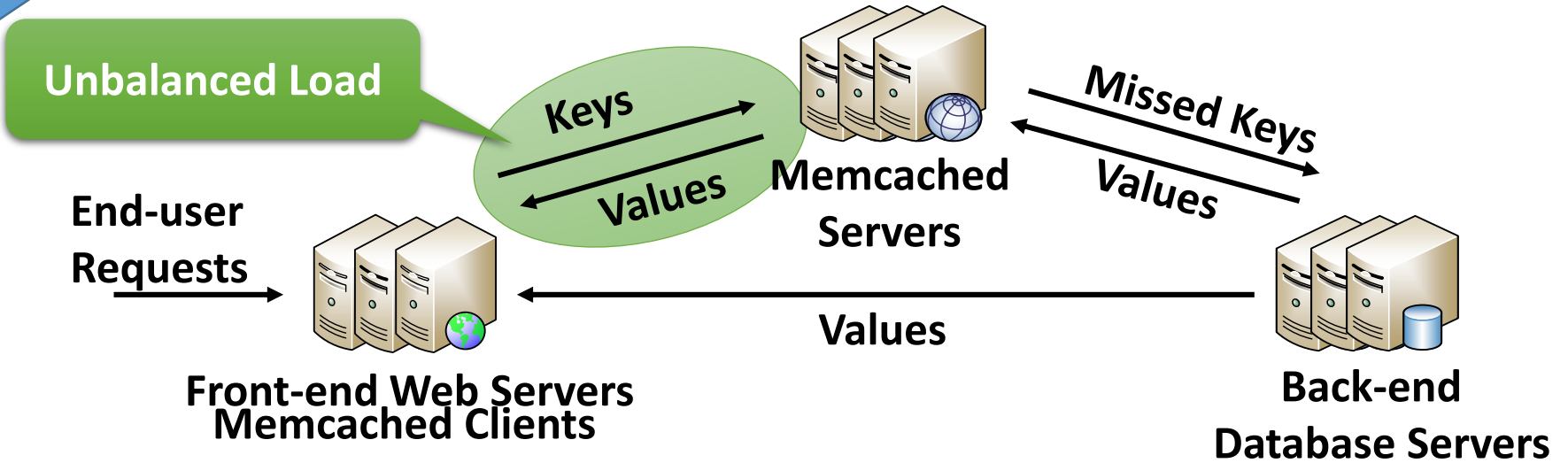
Number of Keys

Cache Miss Ratio

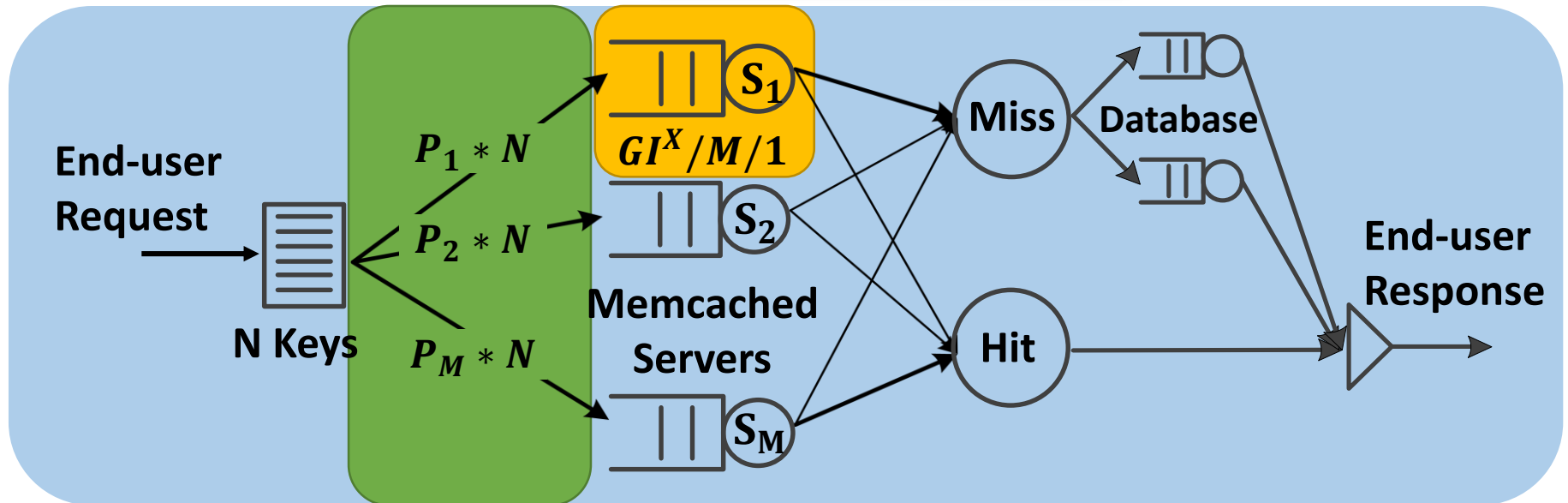
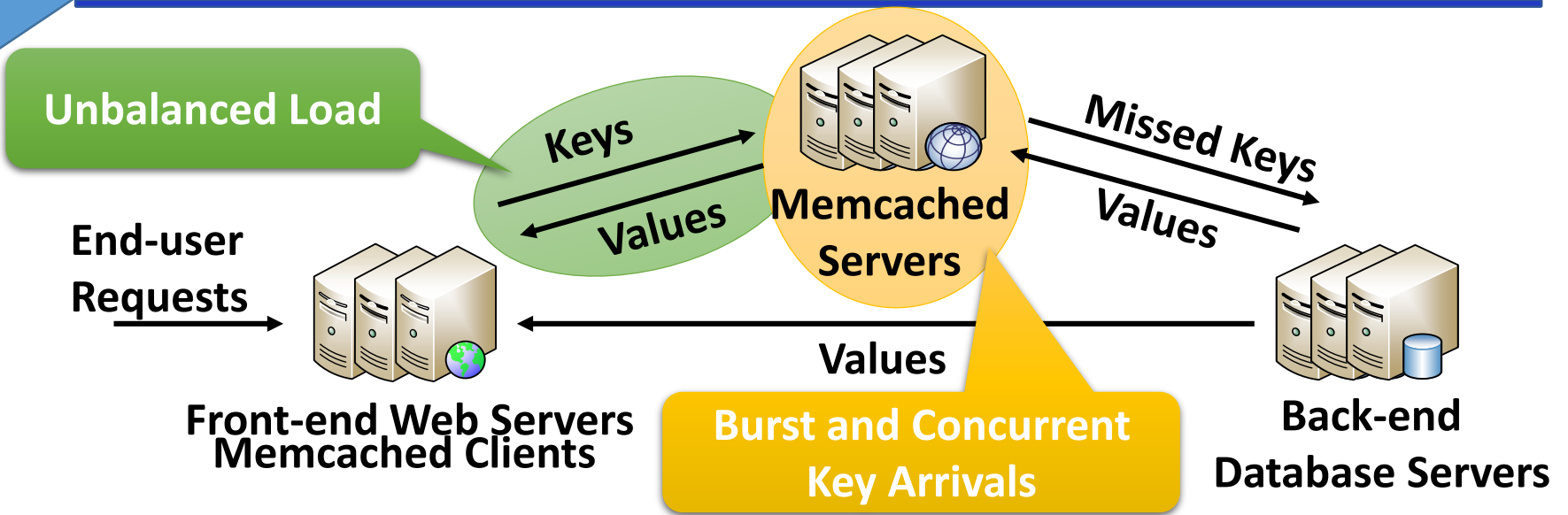
# Model for Memcached



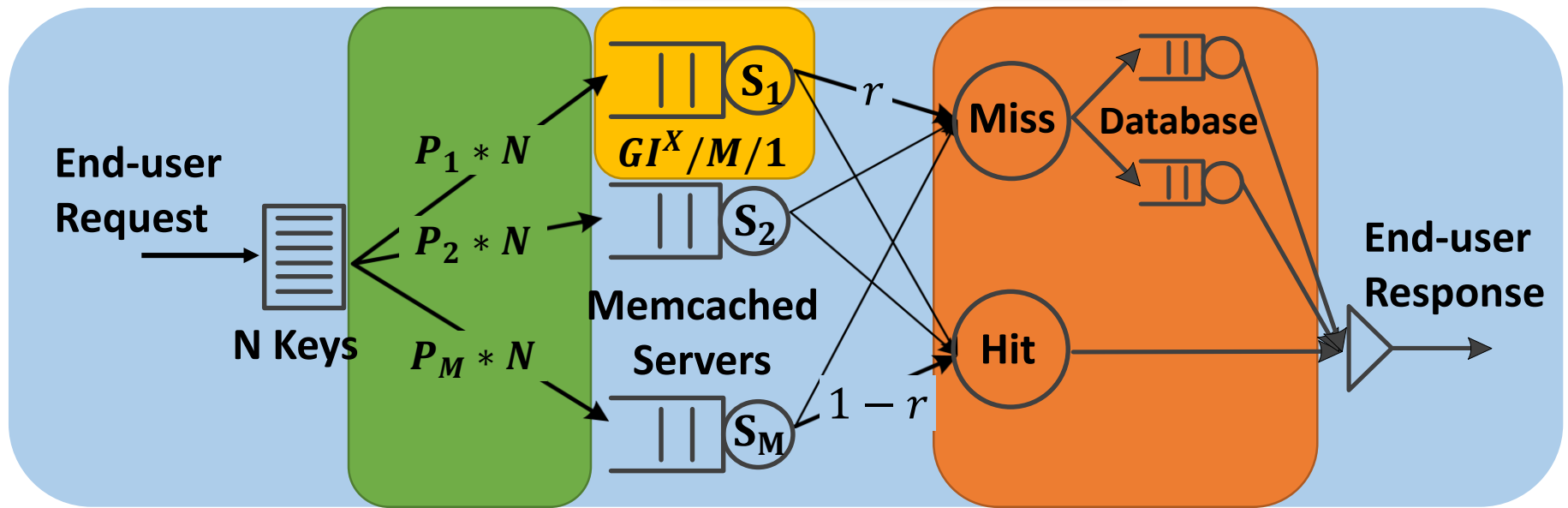
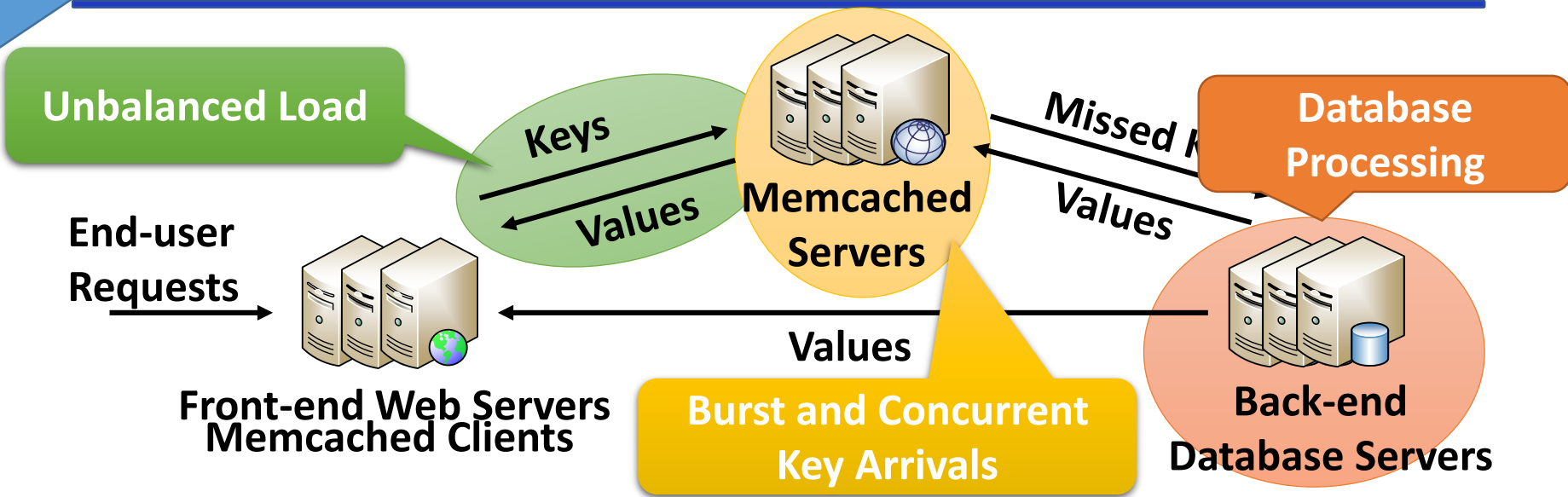
# Specific Features



# Specific Features

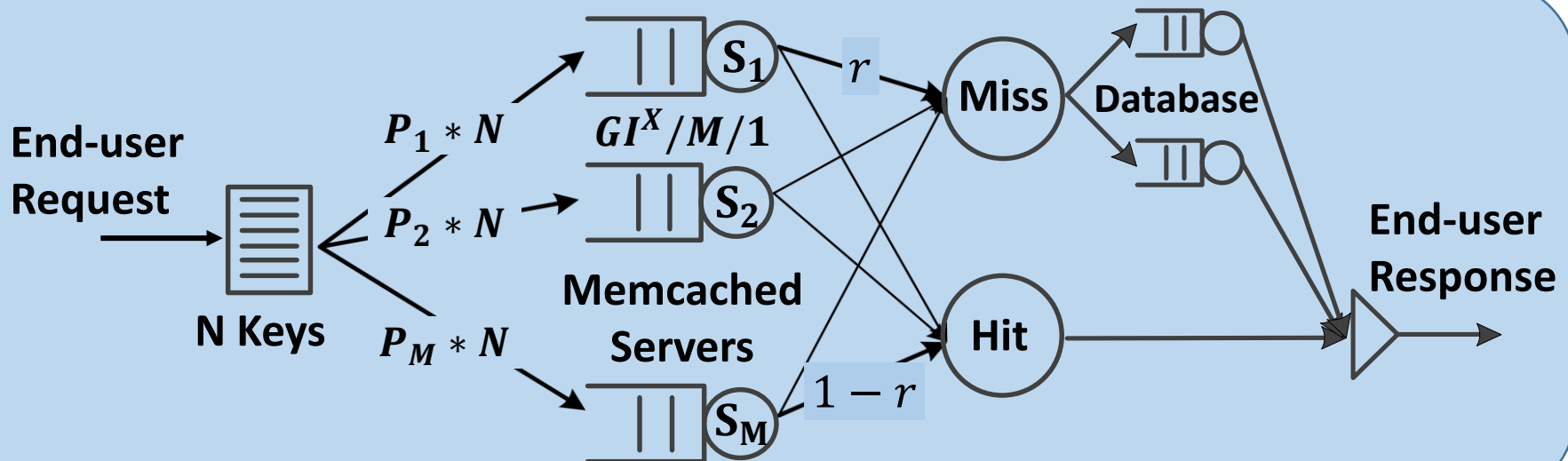


# Specific Features

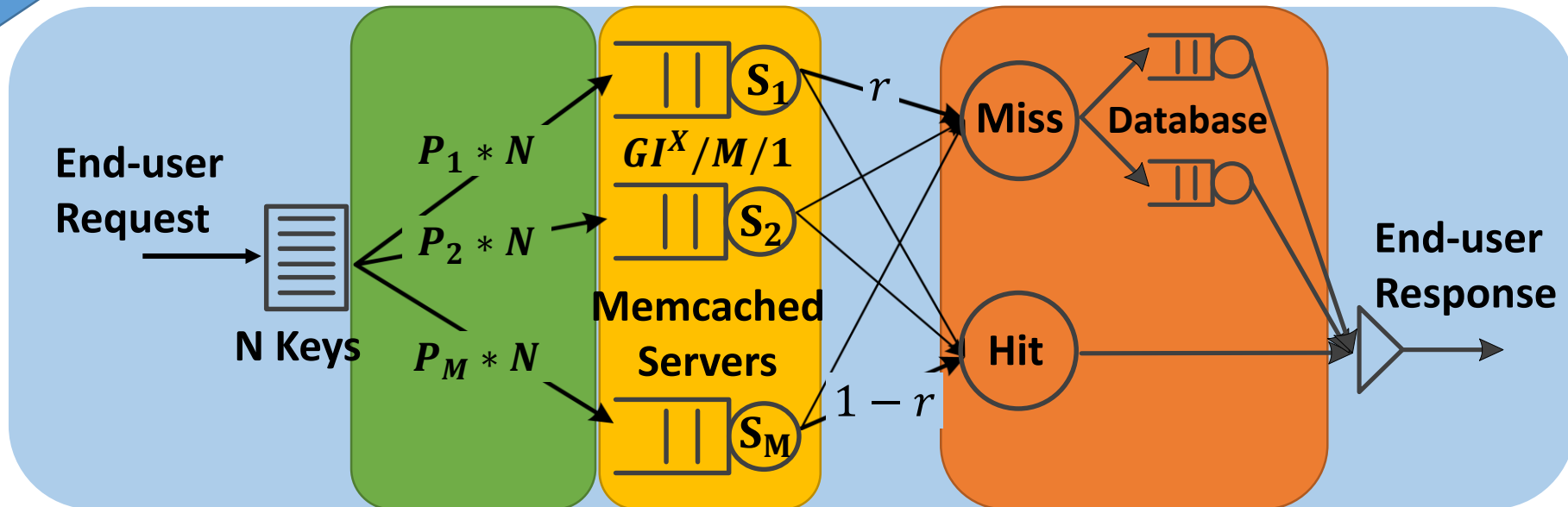


# Assumptions

- One end-user request simultaneously generates multiple independent keys.
- Service times are exponential, both at Memcached servers and database. <sup>[CoNEXT'13][NSDI'15]</sup>

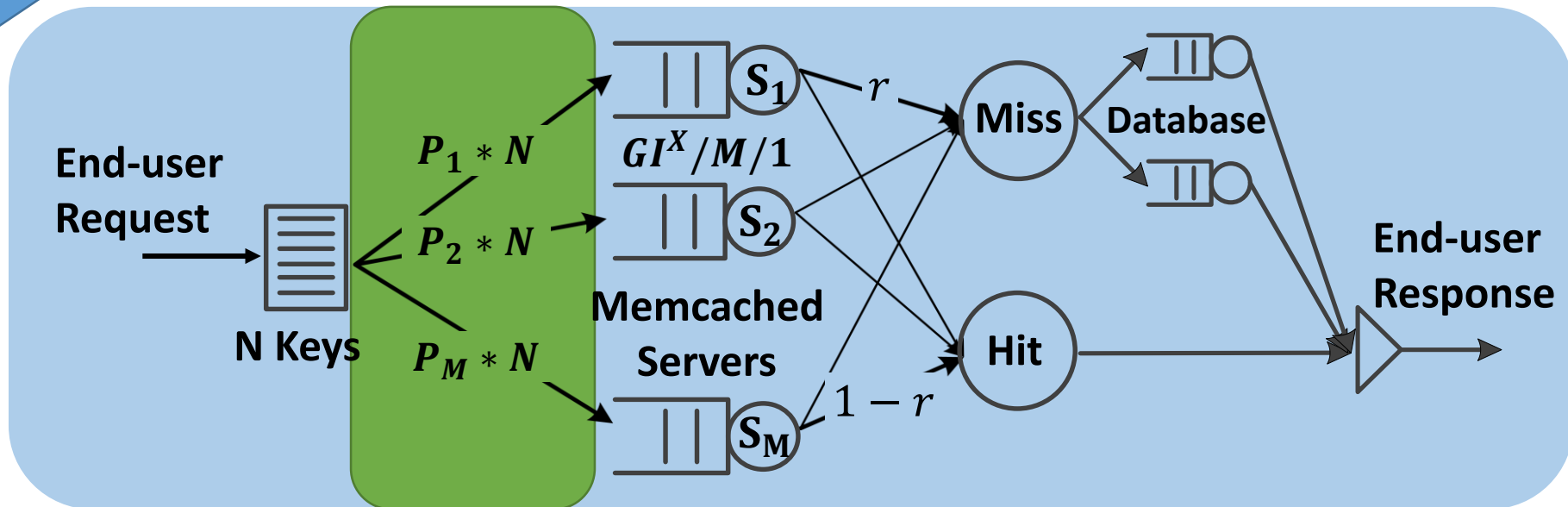


# Latency Composition



- Network Latency
- Processing Latency at Memcached Servers
- Processing Latency at Database

# Network Latency

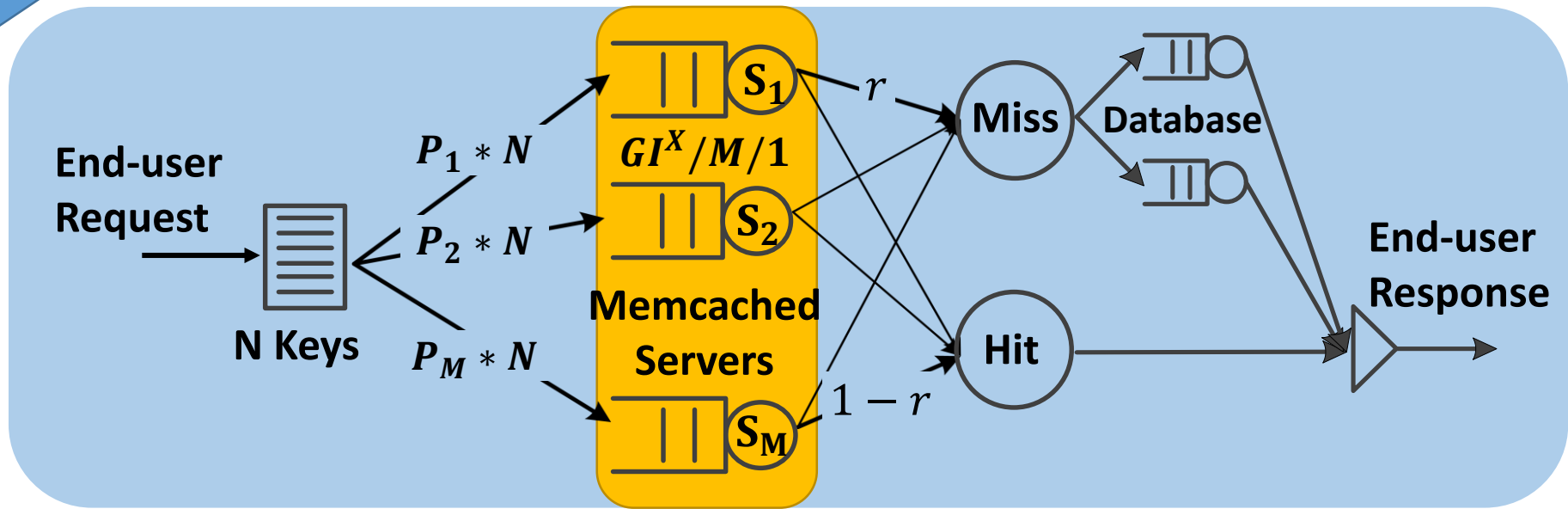


- Network Utilization  $\ll 10\%$   $\rightarrow$  No queueing in network transmission
  - Request rate:  $10^5/s$
  - Request size: Key  $< 200B$ , Value  $< 1KB$  [SIGMETRICS'12]
  - $\rightarrow$  Workload  $\ll 1Gbps$
  - $\rightarrow$  Network bandwidth  $10Gbps$

Network latency is almost constant

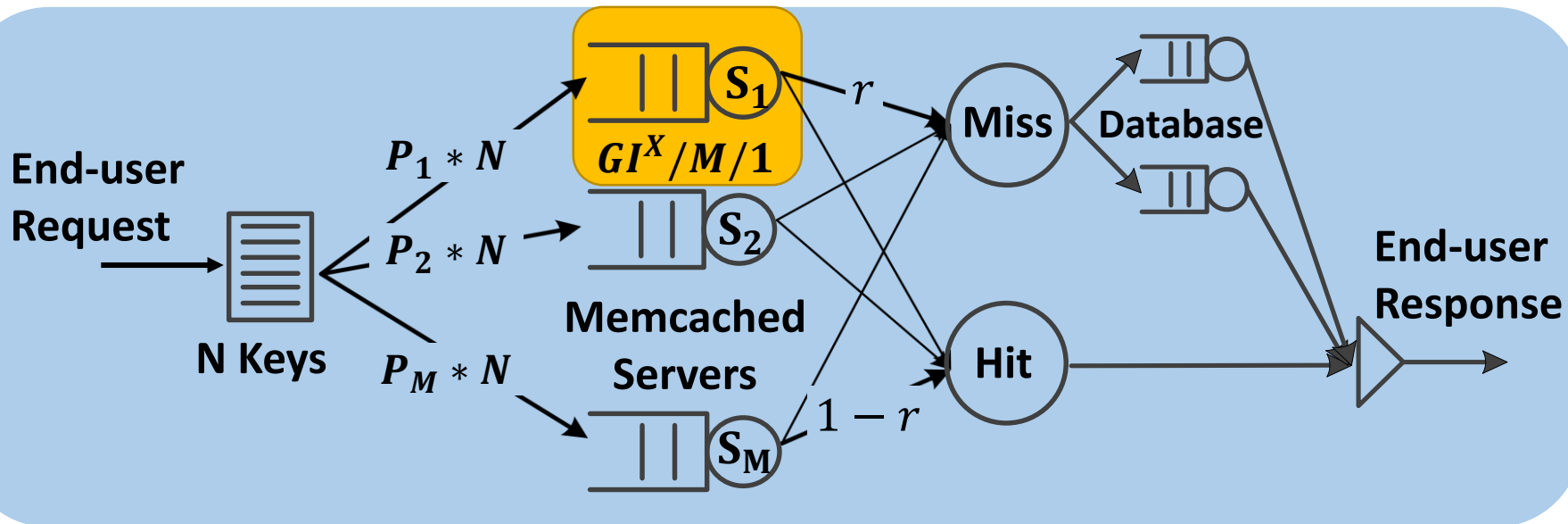


# Latency at Memcached servers



- Latency for 1 Key  $\rightarrow$  • Maximum Latency for N Keys

# Latency at Memcached servers



- Latency for 1 Key

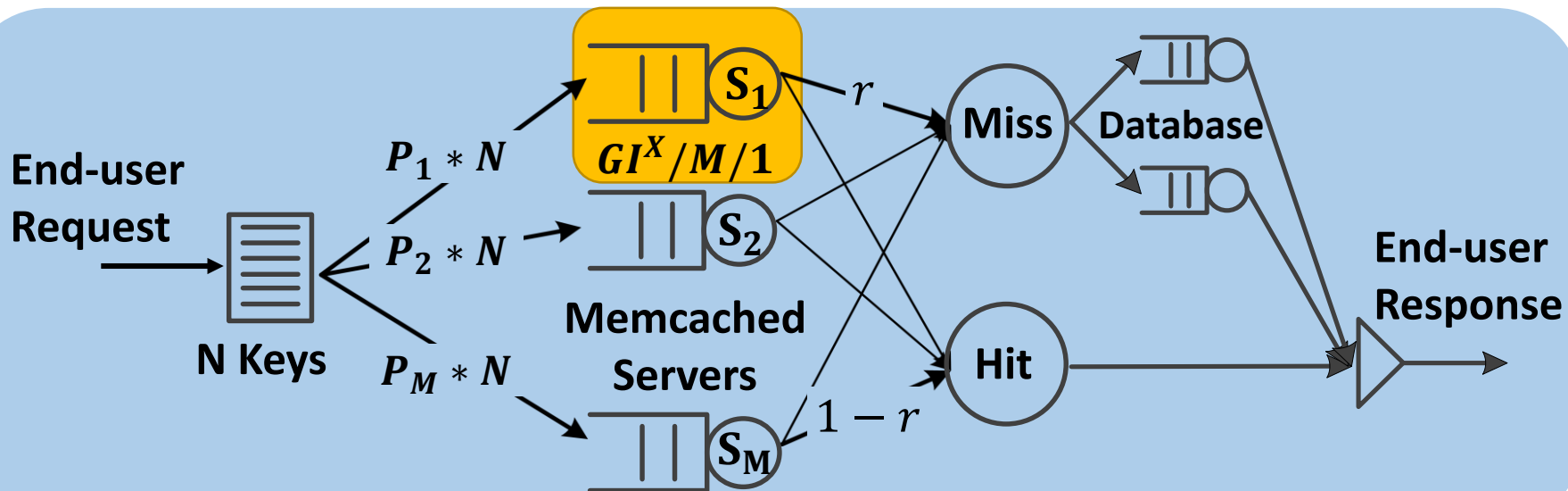
- Queuing at each Memcached server :  $GI^X/M/1$

- $GI$ : Any distribution  $\rightarrow$  Burst arrival of keys
    - $X$ : Batch blocks  $\rightarrow$  Concurrent arrival of keys

- $GI^X/M/1 \rightarrow GI/M/1$

- Response time for a batch  $<$  Latency for 1 Key  $\leq$  Response time for a batch

# Latency at Memcached servers



- Latency for 1 Key

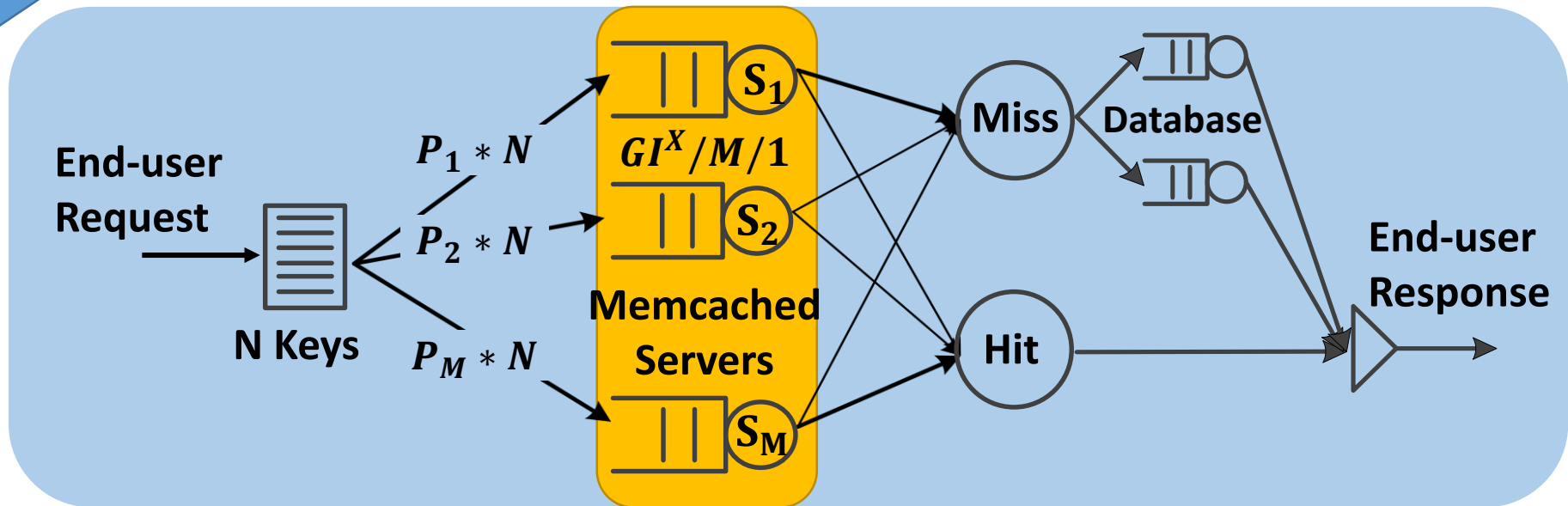
- Queuing at each Memcached server :  $GI^X/M/1$

- $GI$ : Any distribution  $\rightarrow$  Burst arrival of keys
- $X$ : Batch blocks  $\rightarrow$  Concurrent arrival of keys

- $GI^X/M/1 \rightarrow GI/M/1$

- Waiting time for a batch  $<$  Latency for 1 Key  $\leq$  Response time for a batch

# Latency at Memcached servers



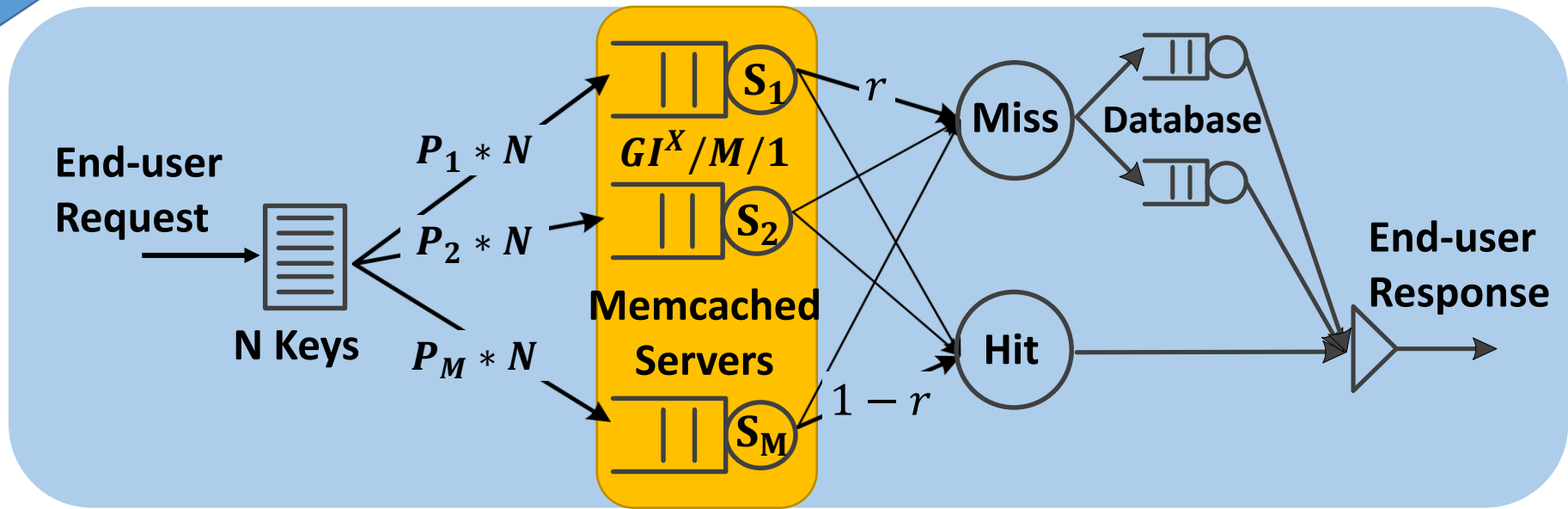
- Maximum Latency for N Keys

- Balanced load  $\rightarrow \approx \left(\frac{N}{N+1}\right)_{th}$  of latency for 1 key

- Unbalanced load among Memcached servers :  $\{P_j\}_{j=1}^M$

- $\in \left[ \left(\frac{N}{N+1}\right)^{\frac{1}{p_{max}}}, \left(\frac{N}{N+1}\right) \right]_{th}$  of latency for a key at the heaviest server

# Latency at Memcached servers



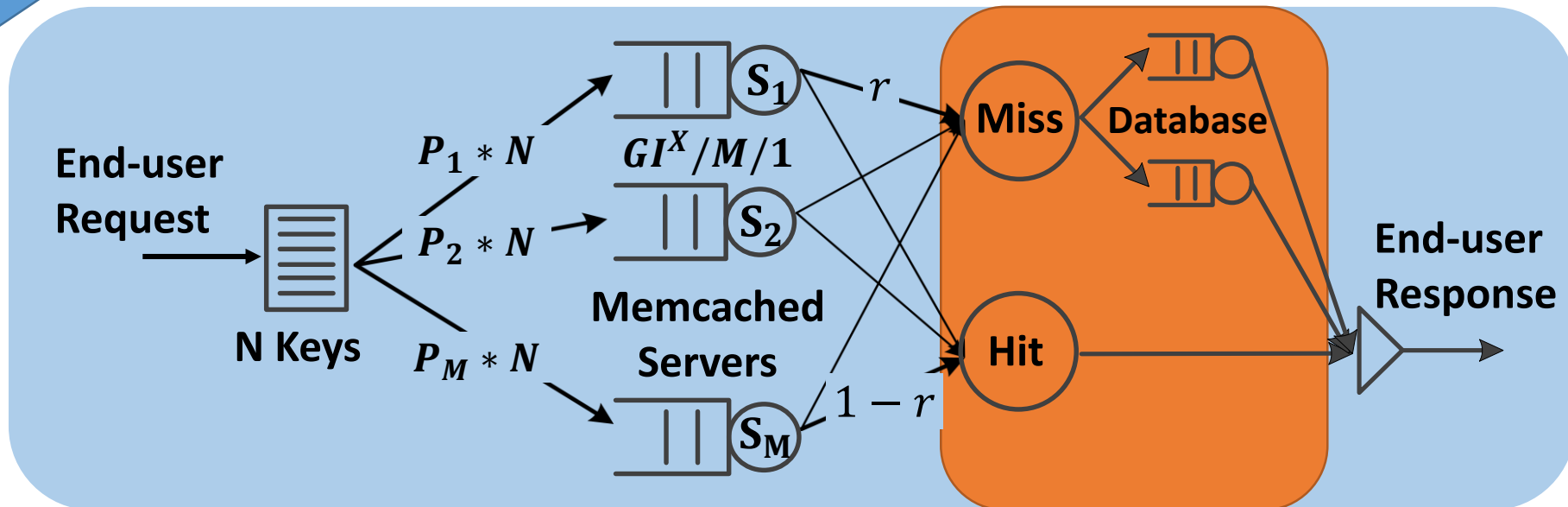
- Maximum Latency for N Keys

- Balanced load  $\rightarrow \approx \left(\frac{N}{N+1}\right)_{th}$  of latency for 1 key

- Unbalanced load among Memcached servers :  $\{P_j\}_{j=1}^M$

- $\in \left[ \left(\frac{N}{N+1}\right)^{p_{max}}, \left(\frac{N}{N+1}\right) \right]_{th}$  of latency for 1 key at the heaviest server

# Latency at Database



- Depends on the number of keys  $N$  and cache miss ratio  $r$ 
  - $(1 - r)^N$ : All-hit cases  $\rightarrow$  Zero latency.
  - $1 - (1 - r)^N$ : Miss cases  $\rightarrow$  Large latency.

# Symbols

| Symbol            | Definition   |
|-------------------|--|
| $N$               | Number of keys generated from an end user request. |
| $M$               | Number of Memcached servers.                       |
| $\{p_j\}_{j=1}^M$ | Load distribution among Memcached servers.         |
| $q$               | Concurrent probability of keys.                    |
| $X$               | Batch size for concurrent keys.                    |
| $T_X$             | Inter-arrival gap of batches.                      |
| $\lambda$         | Arrival rate of keys.                              |
| $\mu_S$           | Service rate of one Memcached server.              |
| $r$               | Cache miss ratio.                                  |
| $\mu_D$           | Service rate at database.                          |

# Latency Estimation

Latency  $T(N)$  of an end-user request that generates  $N$  keys is separately bounded by three parts,

$$\max\{T_N(N), T_S(N), T_D(N)\} \leq T(N) \leq T_N(N) + T_S(N) + T_D(N)$$

- The network latency  $T_N(N)$  is almost constant.
- The processing latency at Memcached servers  $T_S(N)$  satisfies

$$\max \left\{ \frac{\ln \delta + \frac{1}{p_{max}} \ln(N + 1)}{(1 - \delta)(1 - q)\mu_S}, 0 \right\} \leq E[T_S(N)] \leq \frac{\frac{1}{p_{max}} \ln(N + 1)}{(1 - \delta)(1 - q)\mu_S}$$

- The processing latency at database  $T_D(N)$  can be estimated by

$$E[T_D(N)] \approx \frac{1 - (1 - r)^N}{\mu_D} * \ln \left( \frac{N * r}{1 - (1 - r)^N} + 1 \right)$$



# Basic Validation

- Experimental Setup
  - 2 client machines and 4 server machines
    - Intel Core™ i5-5200U CPU and 8GB memory
    - Same Rank, 1Gbps links
  - Workload
    - Mutilate: 512 connections
    - Measurements of Facebook <sup>[Sigmetrics'12]</sup>

| Latency           | Estimation  | Experiment                   | Confidence Interval                                   |
|-------------------|---|------------------------------|---|
| Network           | 20 $\mu$ s  | 20 $\mu$ s                   | [18.12 $\mu$ s, 21.68 $\mu$ s]                        |
| Memcached Servers | 351 $\mu$ s~366 $\mu$ s                           | 368 $\mu$ s                  | [362 $\mu$ s, 373 $\mu$ s]                            |
| Database          | 836 $\mu$ s                                       | 867 $\mu$ s                  | [855 $\mu$ s, 879 $\mu$ s]                            |
| <b>Total</b>      | <b>836<math>\mu</math>s~1222<math>\mu</math>s</b> | <b>1144<math>\mu</math>s</b> | <b>[1128<math>\mu</math>s, 1160<math>\mu</math>s]</b> |

# Factors on Latency

Load Size



Service Rate

- How much improvement on latency can be achieved by optimizing each factor?

**Quantitative Analysis**



Values

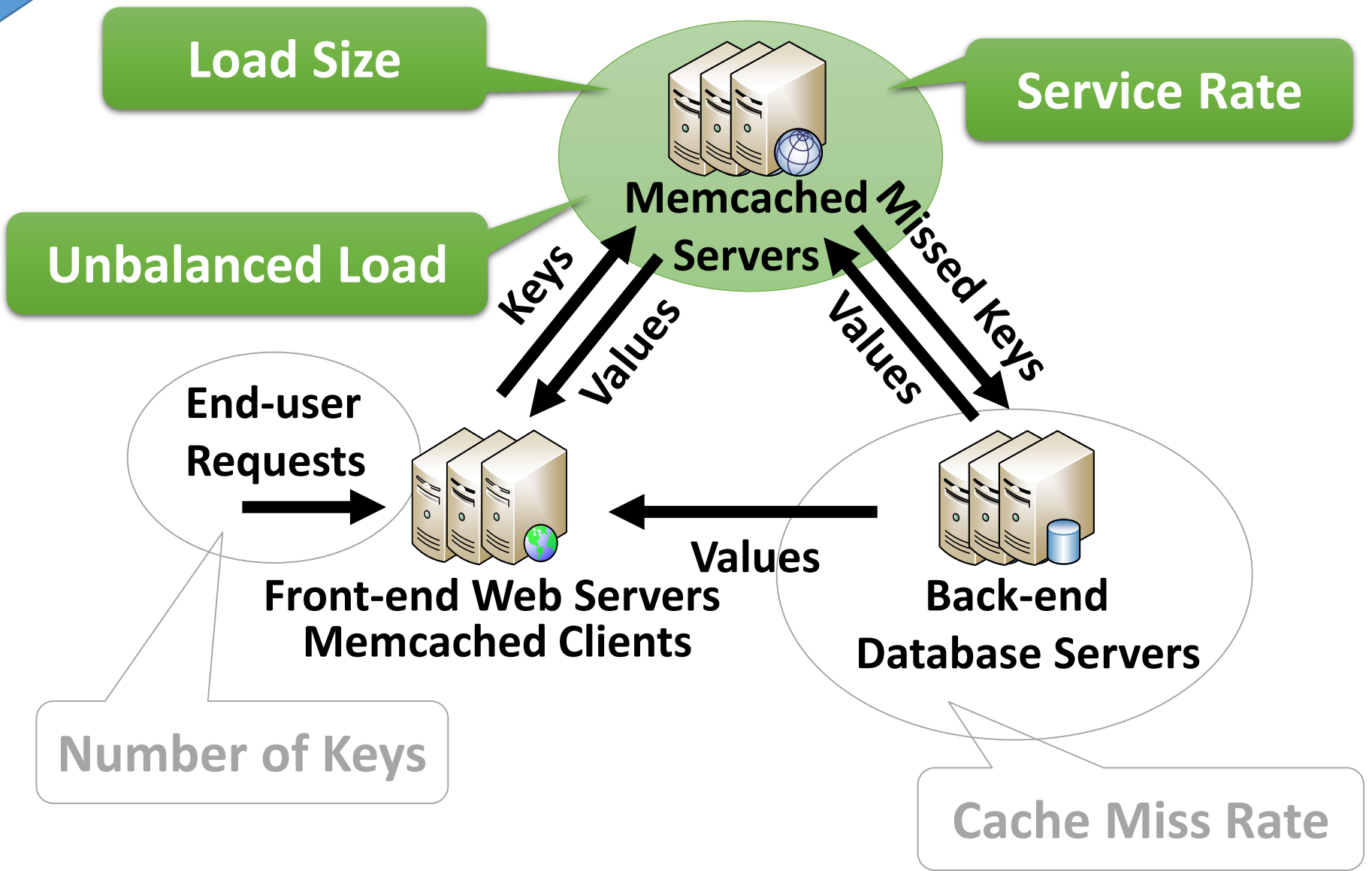
Front-end Web Servers  
Memcached Clients

Back-end  
Database Servers

Number of Keys

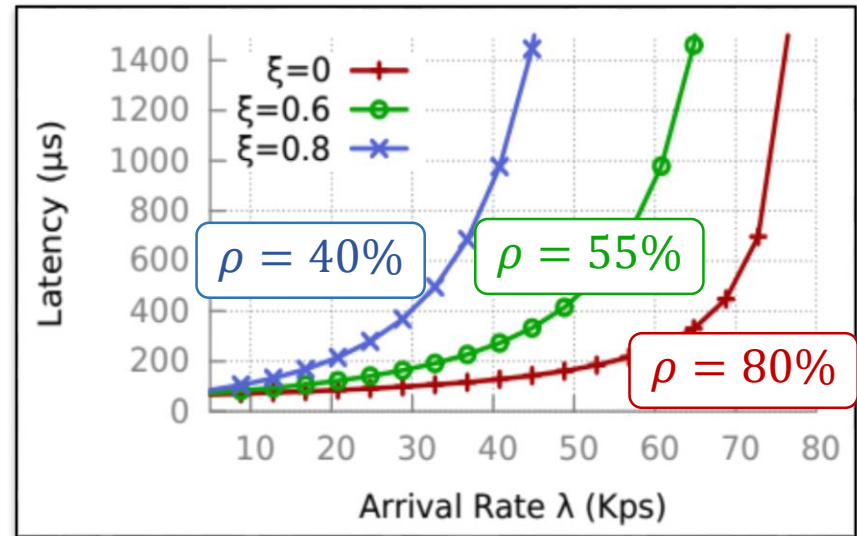
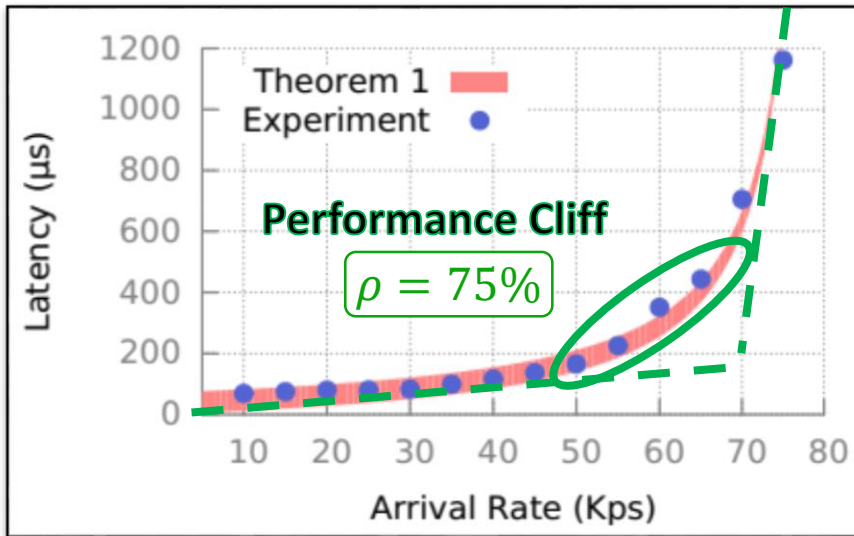
Cache Miss Ratio

# Factors on Latency



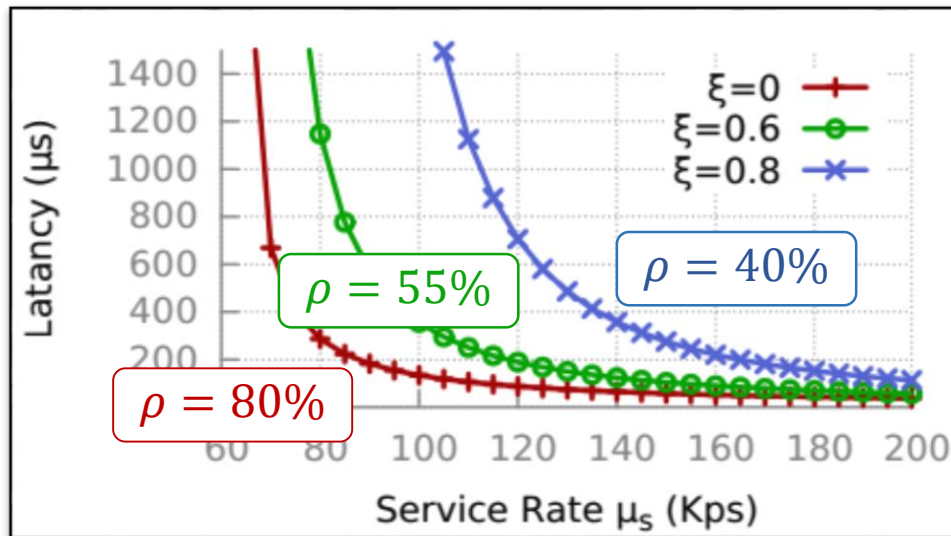
# Load Size

- Workload: Generalized Pareto Distribution
  - Arrival rate:  $\lambda$  (default = 62.5Kbps)
  - Burst degree:  $\xi$  (default = 0.15)
- Service rate:  $\mu_s$  (default = 80Kbps)

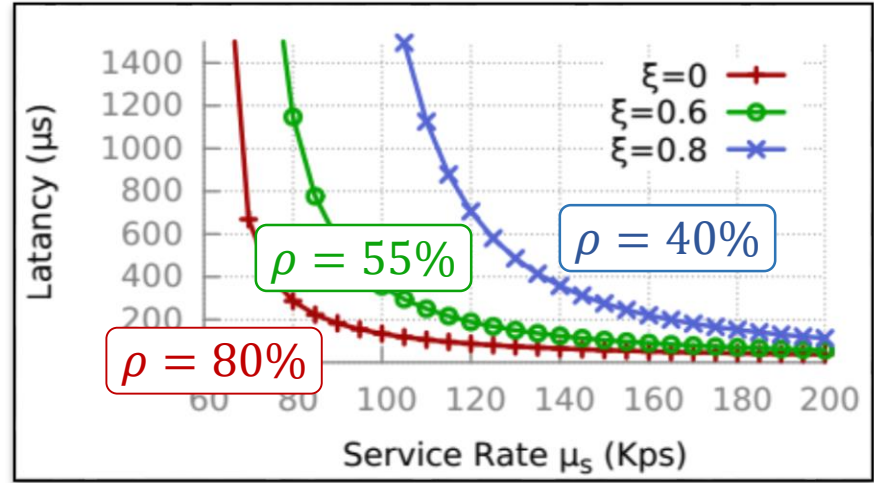
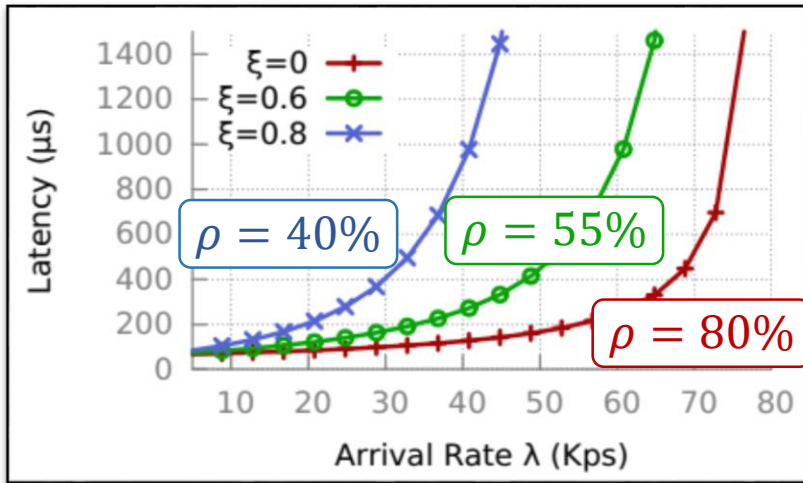


# Service Rate

- Workload: Generalized Pareto Distribution
  - Arrival rate:  $\lambda$  (default = 62.5Kbps)
  - Burst degree:  $\xi$  (default = 0.15)
- Service rate:  $\mu_s$  (default = 80Kbps)

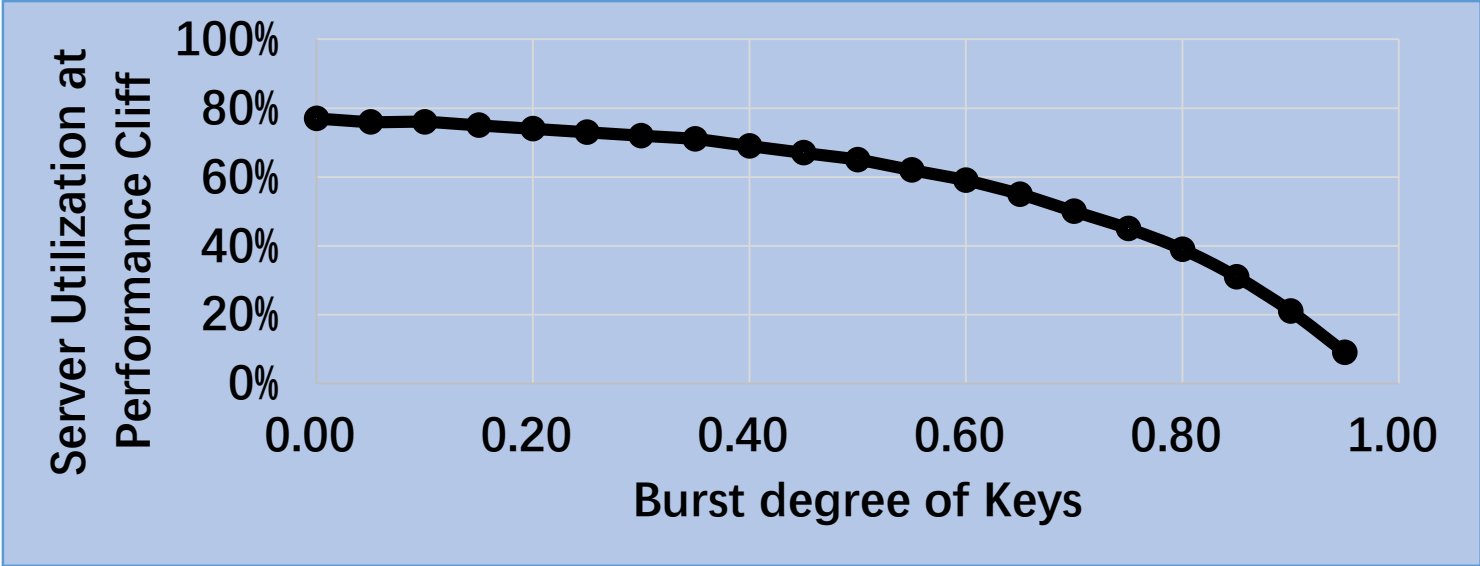


# Load Size & Service Rate



- Latency reaches a **cliff point** when the Memcached **server utilization** gets to a specific value.
- This specific utilization is only determined by the **burst degree** of workload.

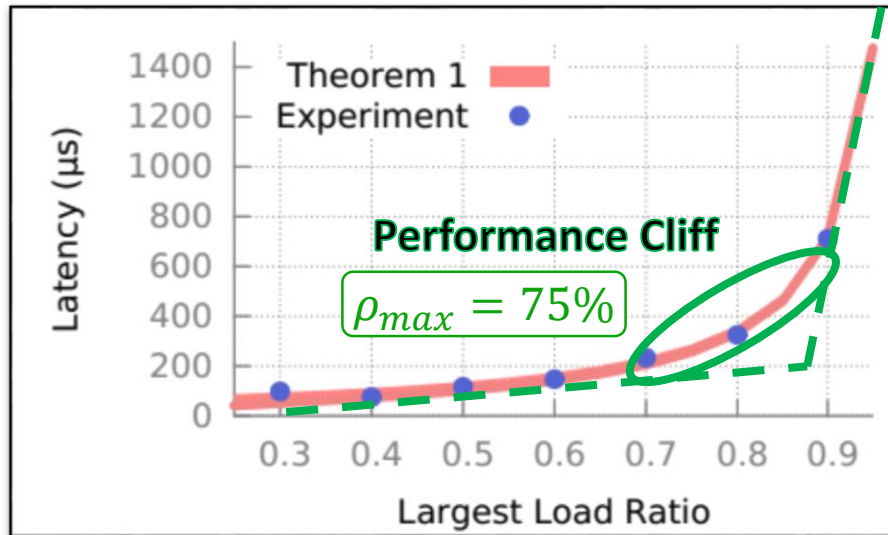
# Server Utilization



- Latency reaches a **cliff point** when the Memcached **server utilization** gets to a specific value.
- This specific utilization is only determined by the **burst degree** of workload.

# Unbalanced Load

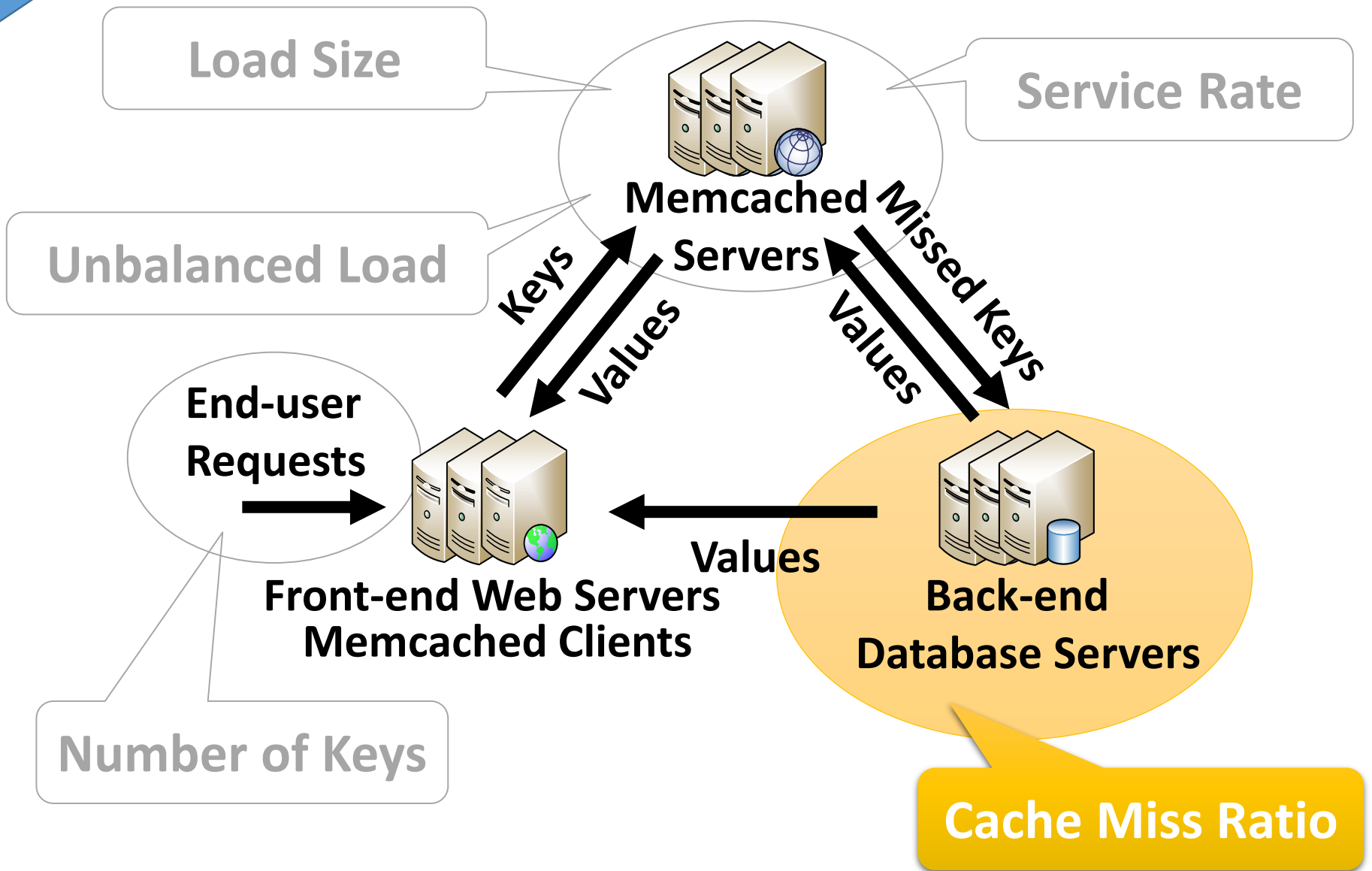
- 4 Memcached servers: Largest load ratio
  - Total load size:  $\Lambda = 80Kbps$
  - Burst degree:  $\xi = 0.15$
  - Service rate:  $\mu_s = 80Kbps$



• Latency still reaches a **cliff point** when the **heaviest** server gets to the **specific utilization**.

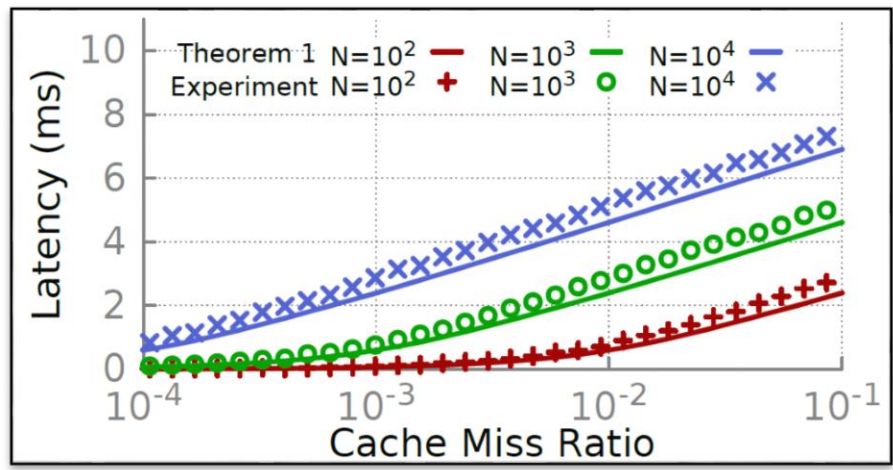
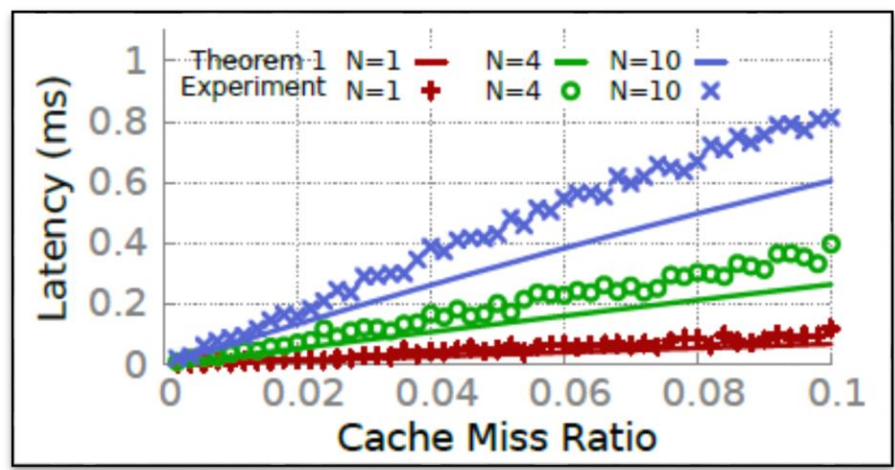


# Factors on Latency



# Cache Miss Ratio

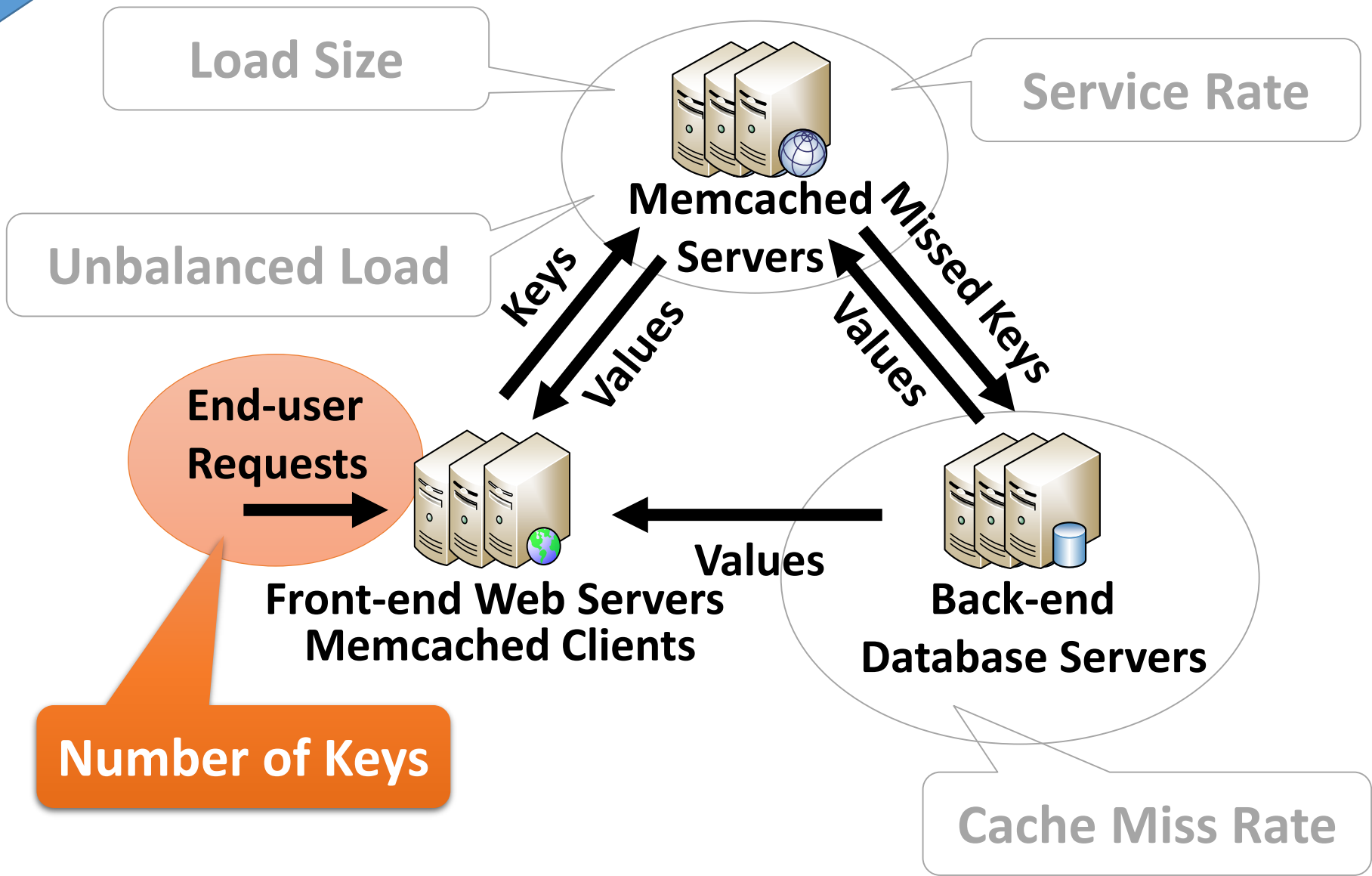
- Average cost at database:  $\frac{1}{\mu_D} = 1ms$



- $r = 0.1 \rightarrow 0.01$ 
  - $N = 1 \rightarrow$  latency:  $0.8ms \rightarrow 0.08ms$ ,  $\downarrow 90\%$
  - $N = 10^4 \rightarrow$  latency:  $7.2ms \rightarrow 4.3ms$ ,  $\downarrow 40\%$

• Latency grows **logarithmically** with the increase of the **cache miss ratio**.

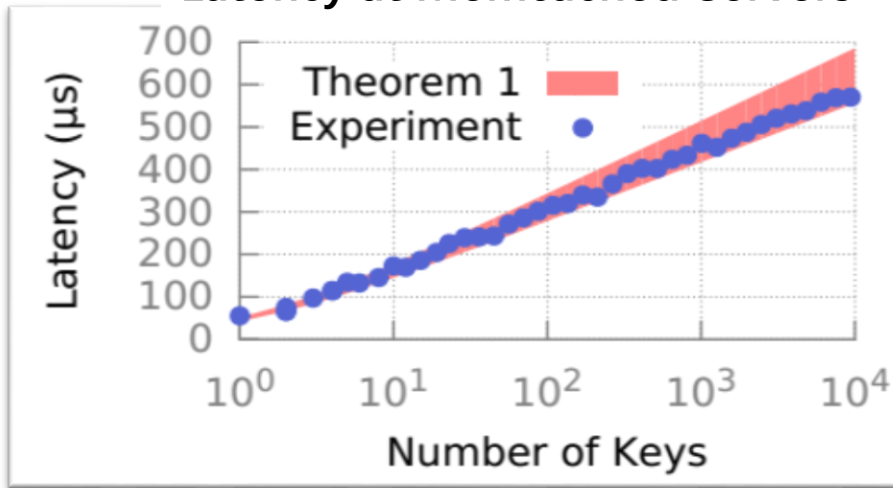
# Factors on Latency



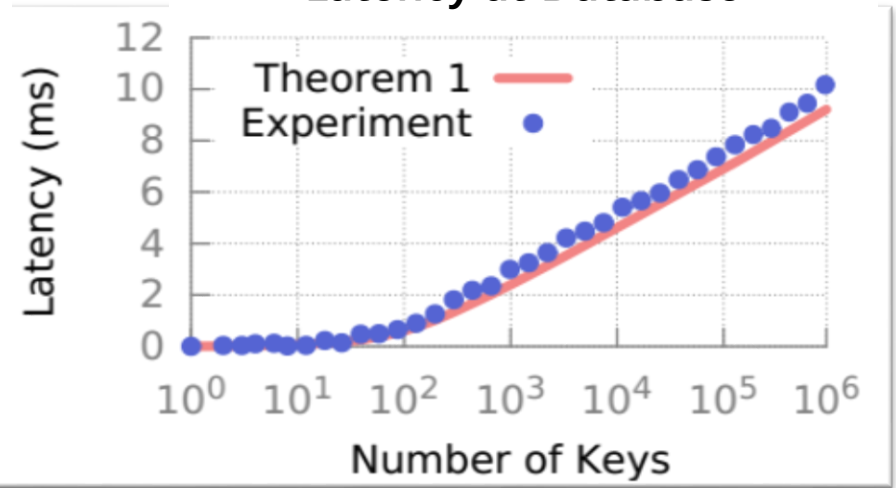
# Number of Keys

- Workload burst degree  $\xi = 0.15$
- Memcached server utilization  $\rho = 0.75$
- Cache miss rate  $r = 0.01$

Latency at Memcached Servers



Latency at Database

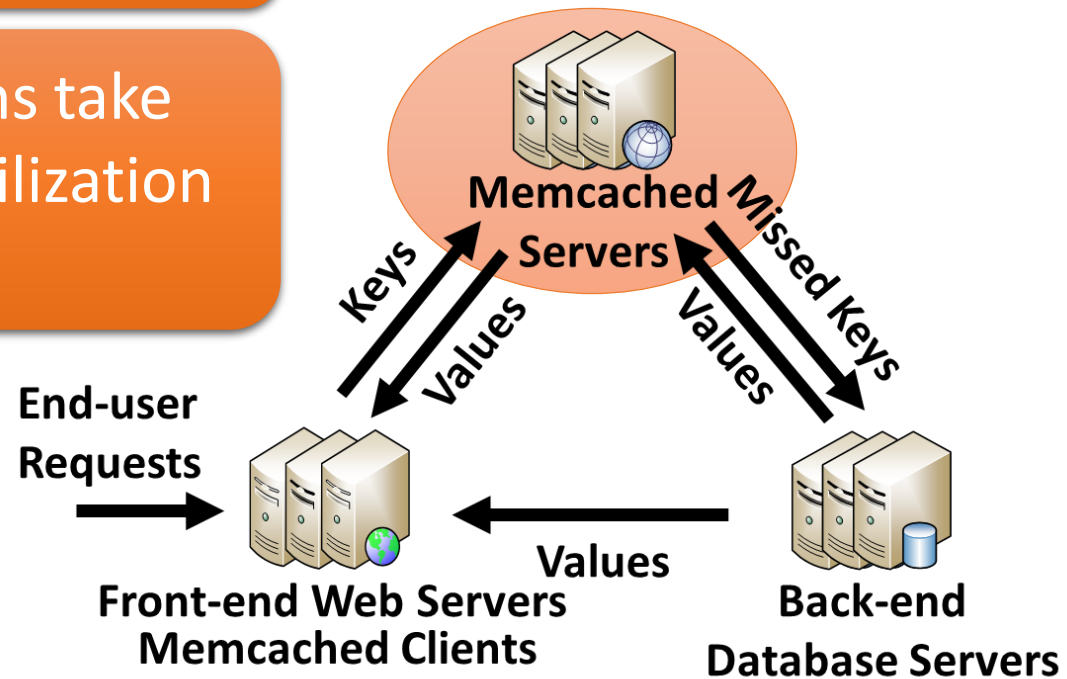


• Latency grows **logarithmically** with the increase of the **number of keys**.

# Recommendations

Server utilization should not exceed a specific value.

Load-balancing mechanisms take effect before the largest utilization extends the specific value.



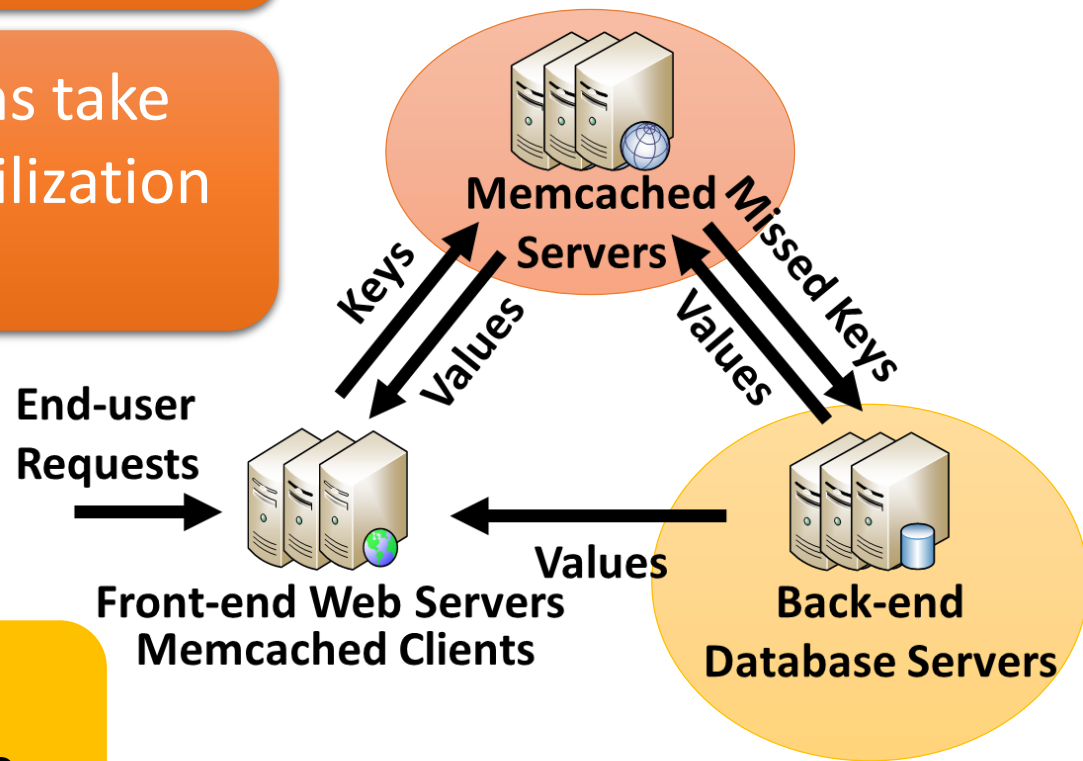
# Recommendations

Analysis

Server utilization should not exceed a specific value.

Load-balancing mechanisms take effect before the largest utilization extends the specific value.

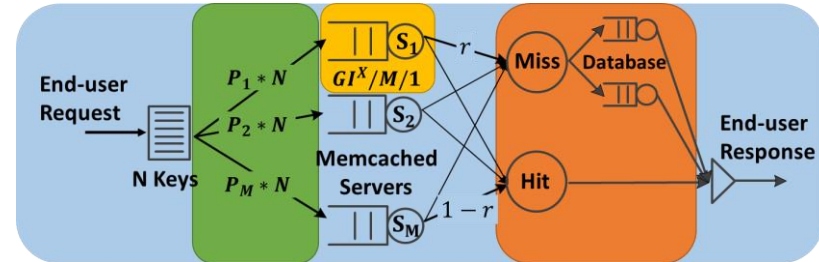
Minimizing the number of keys instead of reducing the cache miss ratio.



# Conclusion

- **Specific model for Memcached**

- Unbalanced load distribution
- Burst and concurrent key arrivals
- Database processing



- **Latency estimation, validation & Analysis**

- Network latency is almost constant
- Latency at Memcached servers is determined by server utilization
- Latency at database is determined by the number of keys and cache miss ratio

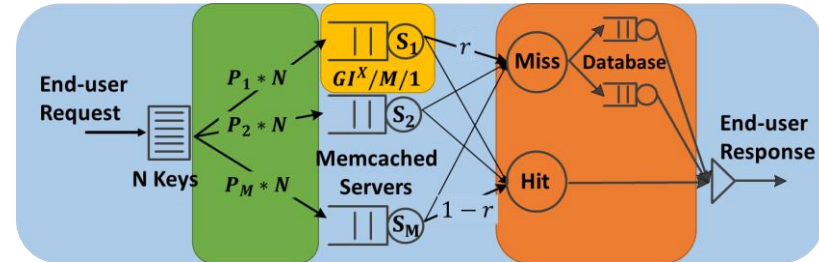
- **Recommendations**

- Server utilization should not exceed a specific value.
- Load-balancing mechanisms take effect before the largest utilization extends the specific value.
- Minimizing the number of keys instead of reducing the cache miss ratio.

# Conclusion

- **Specific model for Memcached**

- Unbalanced load distribution
- Burst and concurrent key arrivals
- Database processing



- **Latency estimation, validation & Quantitative analysis**

- Network latency is almost constant
- Latency at Memcached servers is determined by server utilization
- Latency at database is determined by the number of keys and cache miss ratio

- **Recommendations**

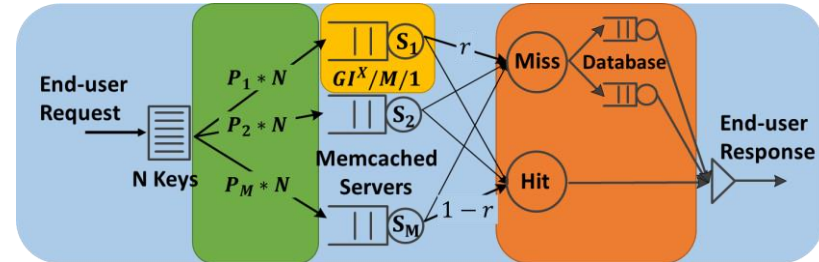
- Server utilization should not exceed a specific value.
- Load-balancing mechanisms take effect before the largest utilization extends the specific value.
- Minimizing the number of keys instead of reducing the cache miss ratio.



# Conclusion

- **Specific model for Memcached**

- Unbalanced load distribution
- Burst and concurrent key arrivals
- Database processing



- **Latency estimation, validation & Analysis**

- Network latency is almost constant
- Latency at Memcached servers is determined by server utilization
- Latency at database is determined by the number of keys and cache miss ratio

- **Recommendations**

- Server utilization should not exceed a specific value.
- Load-balancing mechanisms take effect before the largest utilization extends the specific value.
- Minimizing the number of keys instead of reducing the cache miss ratio.

**Thanks!**