



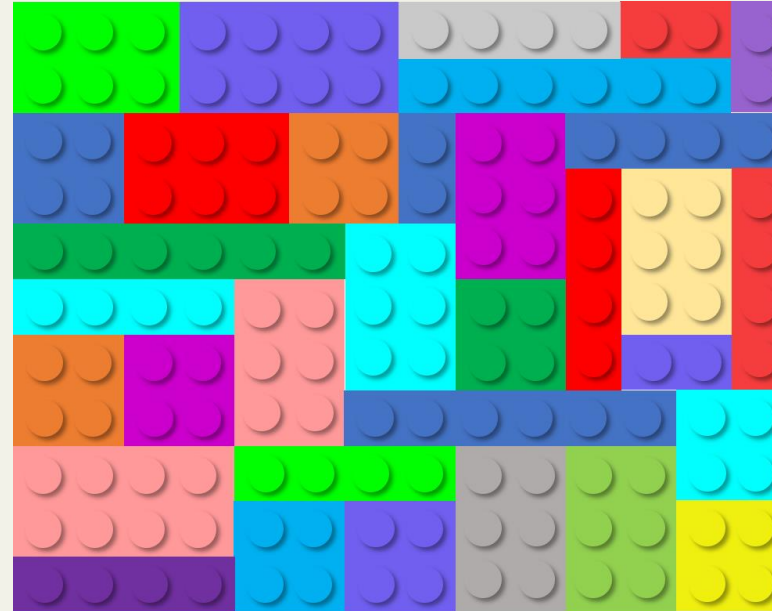
Cornell University
Computer Systems Laboratory



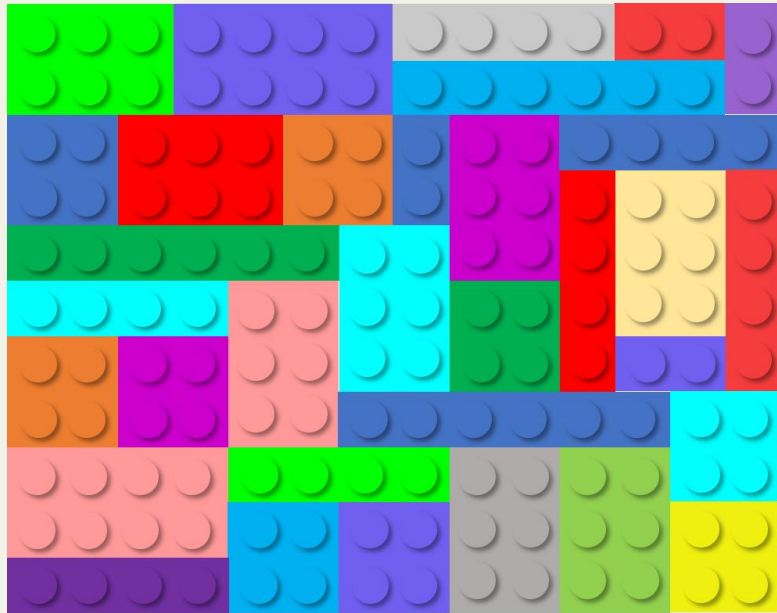
μ QSIM: ENABLING ACCURATE AND SCALABLE SIMULATION FOR INTERACTIVE MICROSERVICES

Yanqi Zhang, Yu Gan, Christina Delimitrou
Cornell University

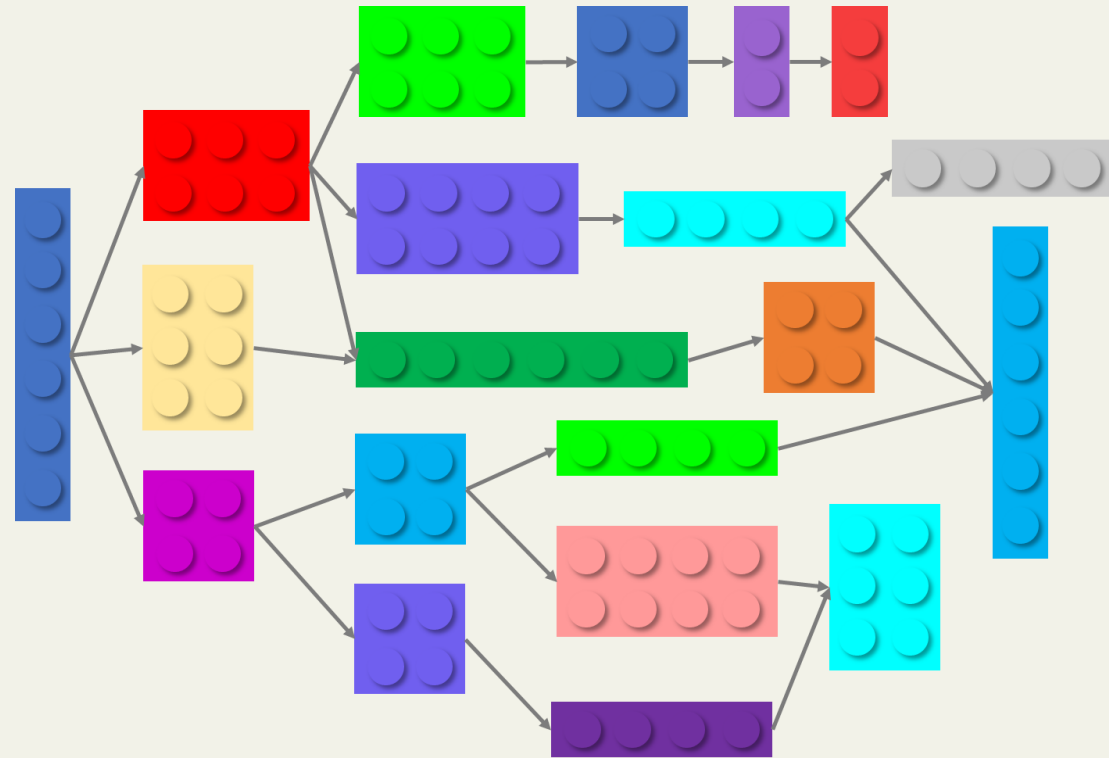
Session: Datacenters and Cloud Computing, March 26th



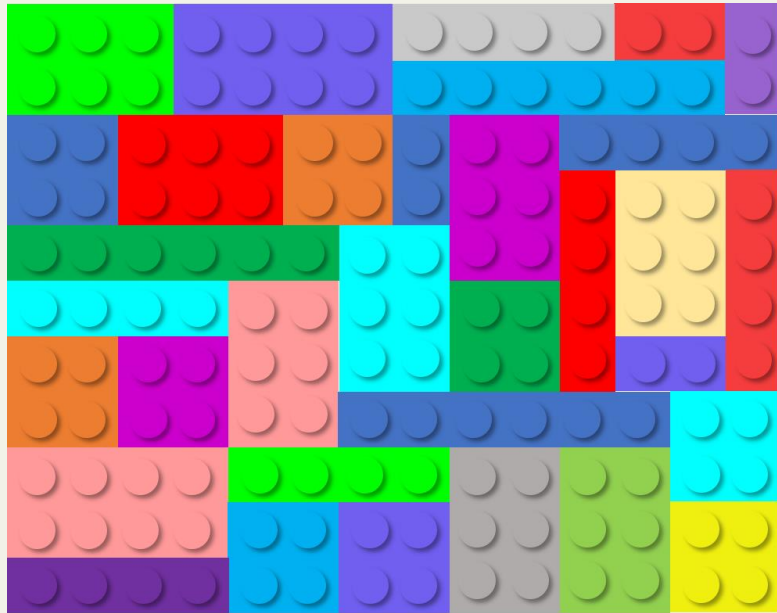
Monolith Application



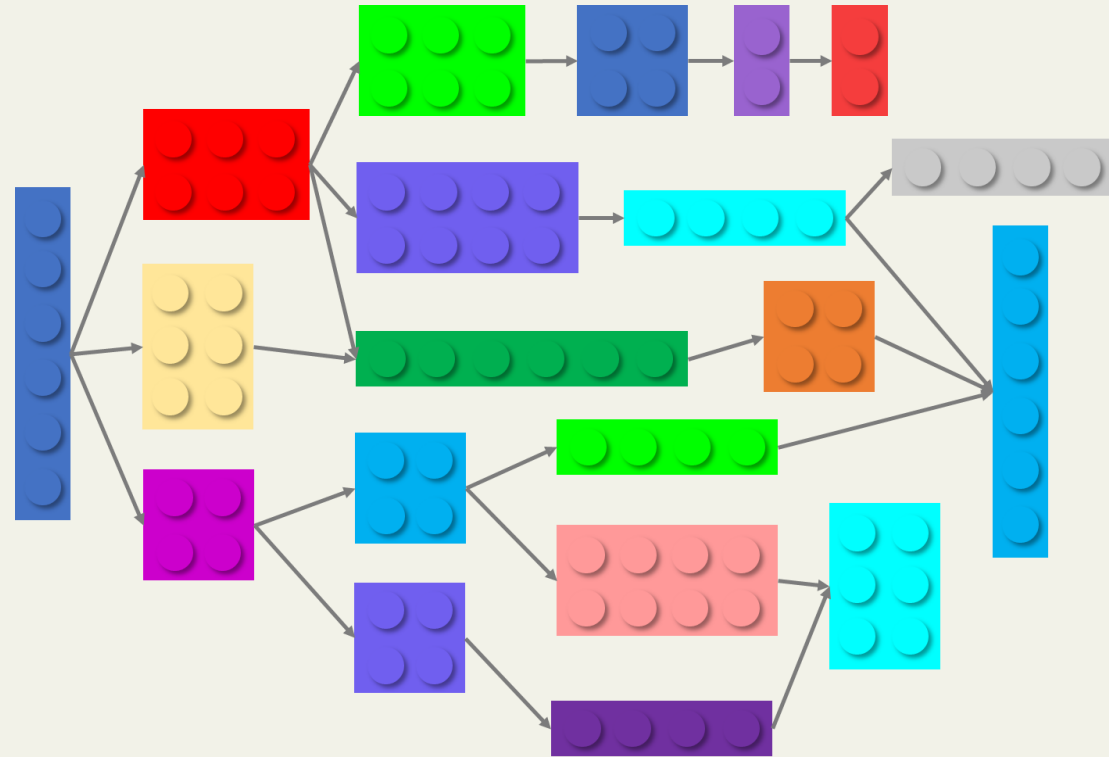
Monolith Application



Microservices

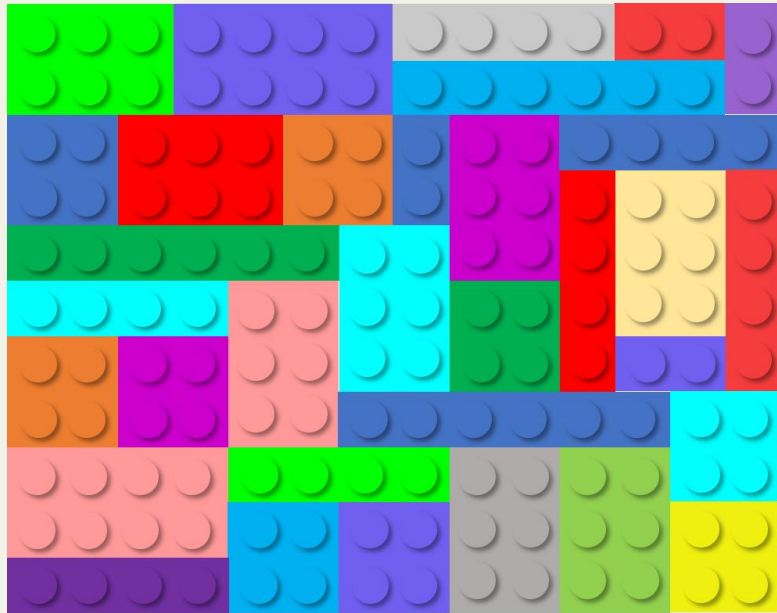


Monolith Application

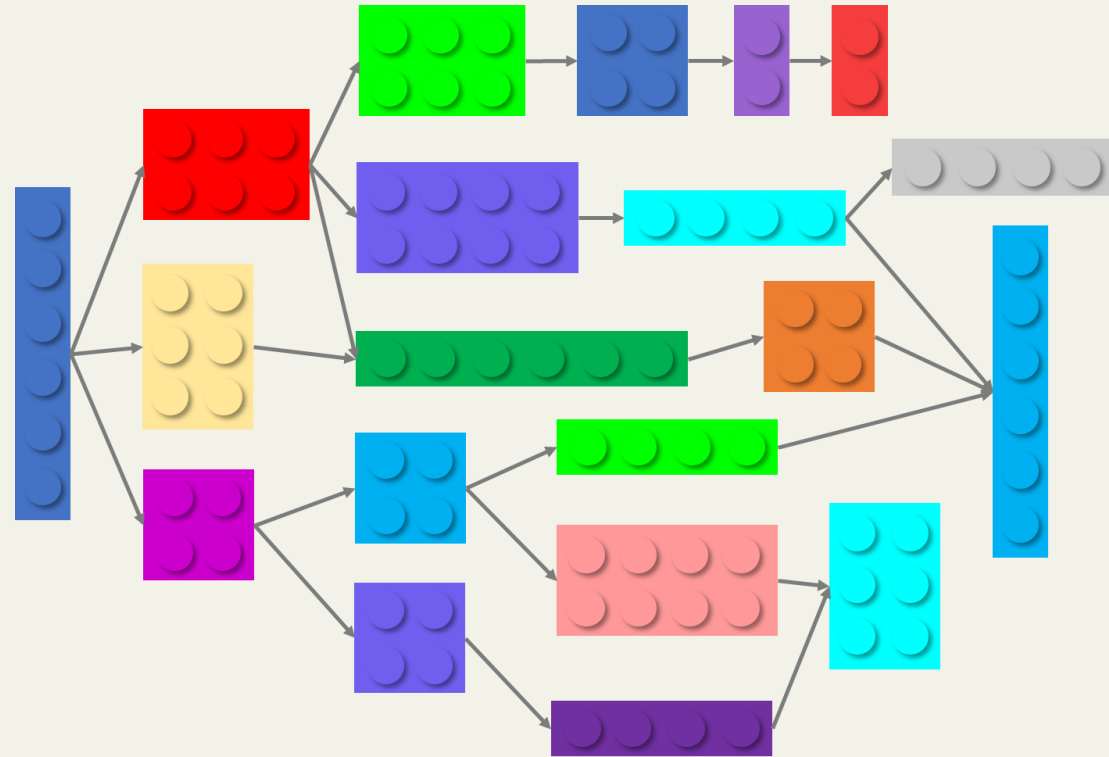


Microservices

- Easier development & update

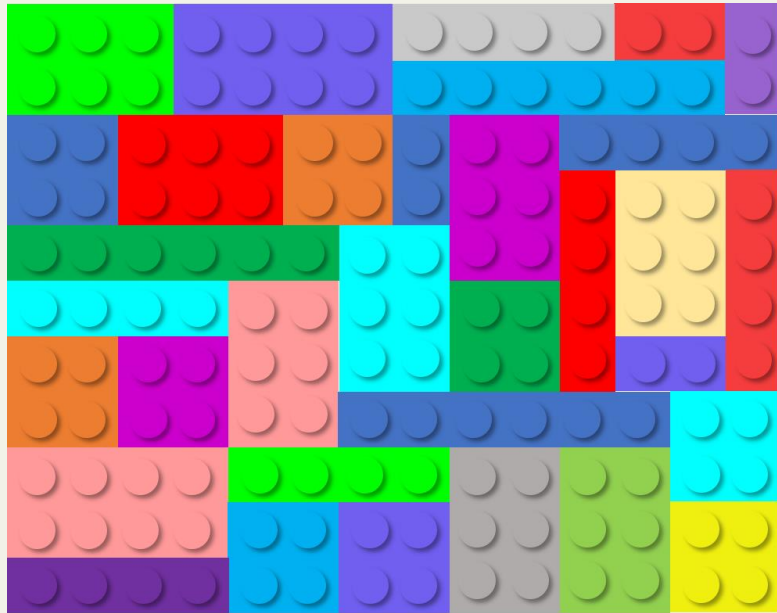


Monolith Application

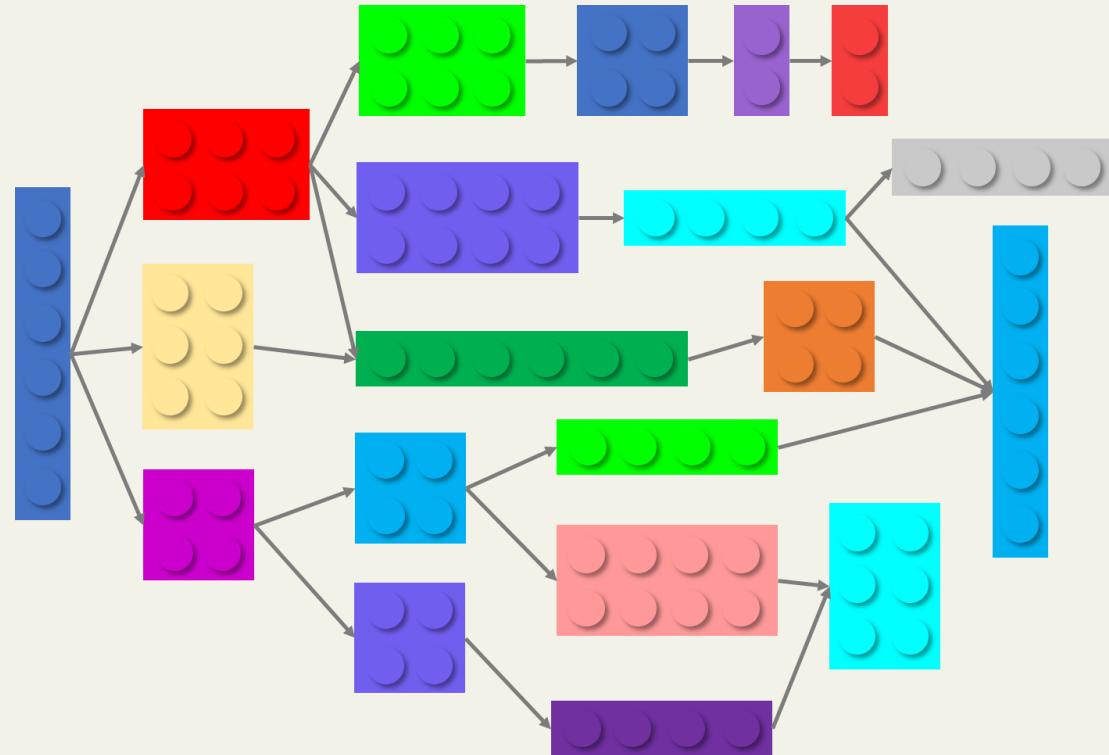


Microservices

- Easier development & update
- Debugging & error isolation

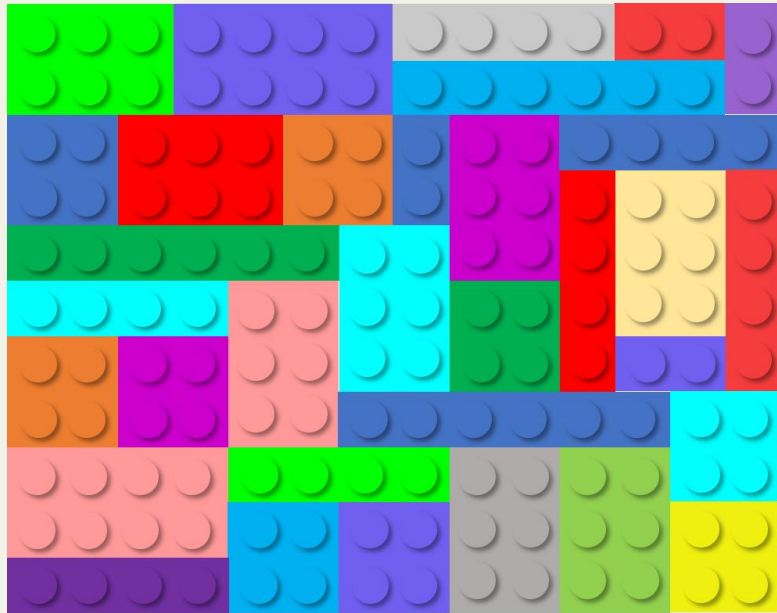


Monolith Application

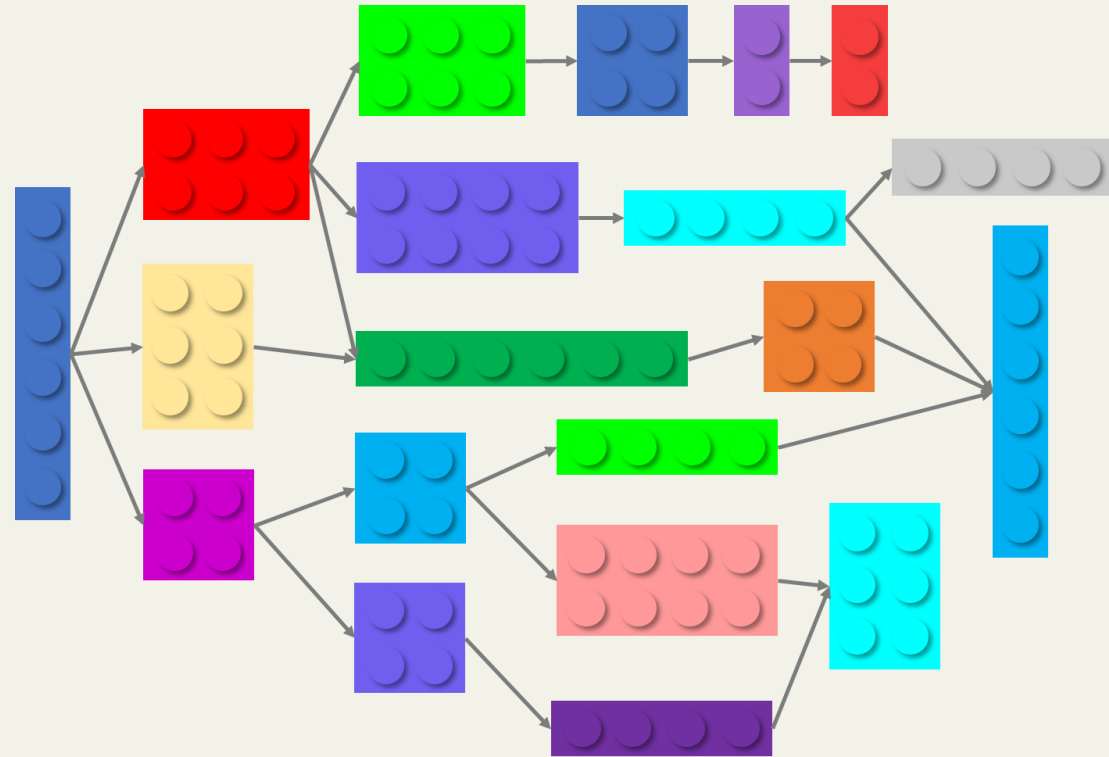


Microservices

- Easier development & update
- Debugging & error isolation
- Elasticity

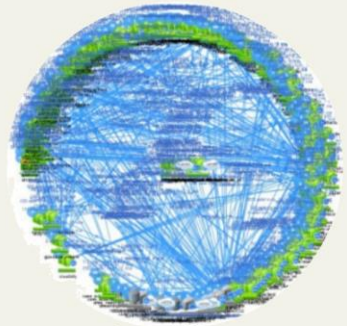


Monolith Application

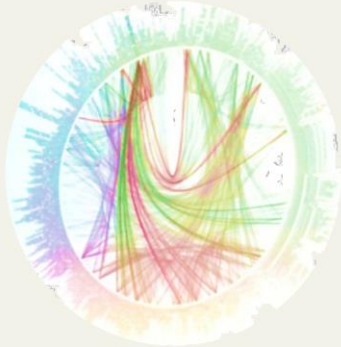


Microservices

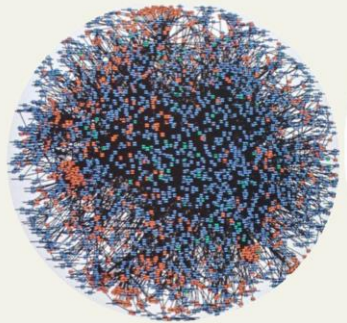
- Easier development & update
- Debugging & error isolation
- Elasticity
- PL/framework heterogeneity



Netflix



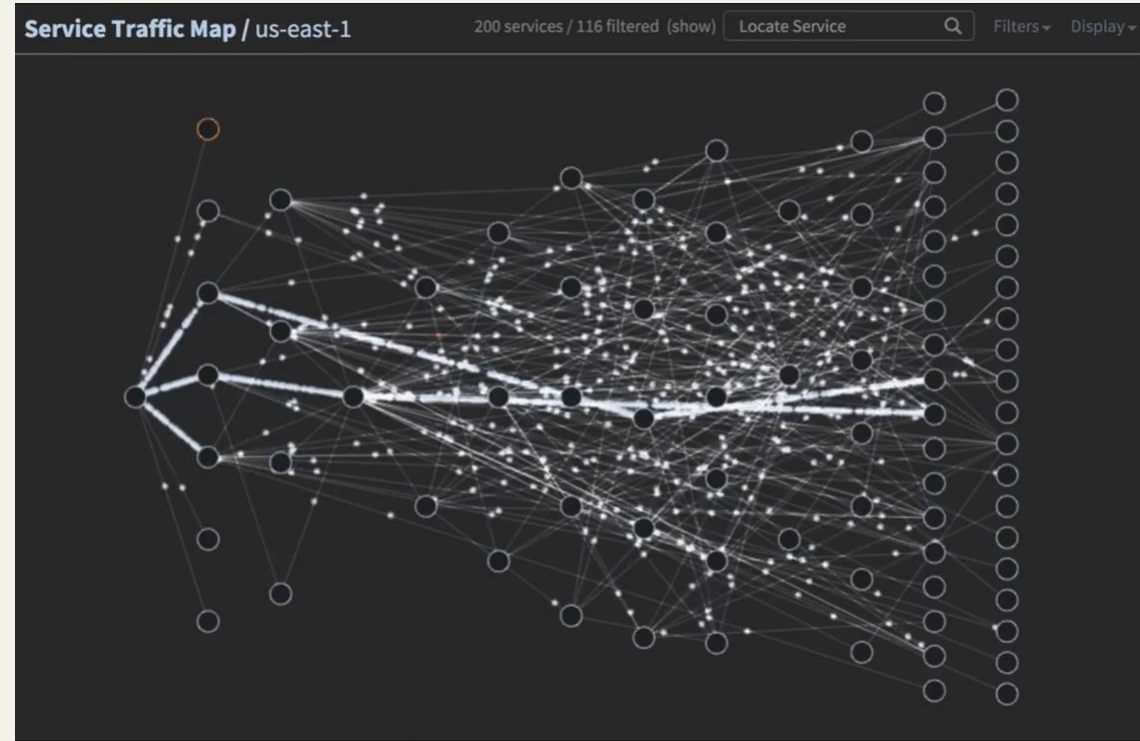
Twitter



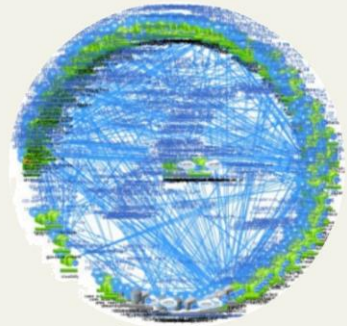
Amazon



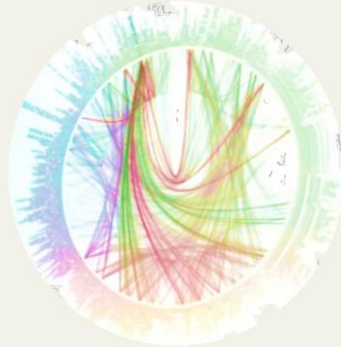
Social Network



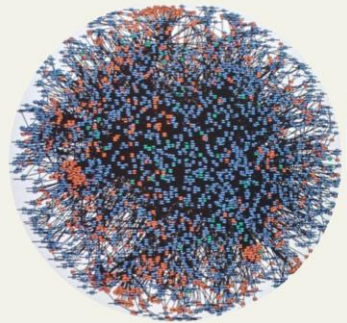
Real time traffic of Social Network Microservices



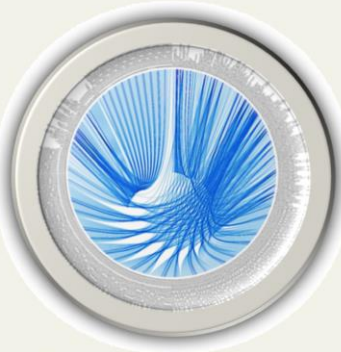
Netflix



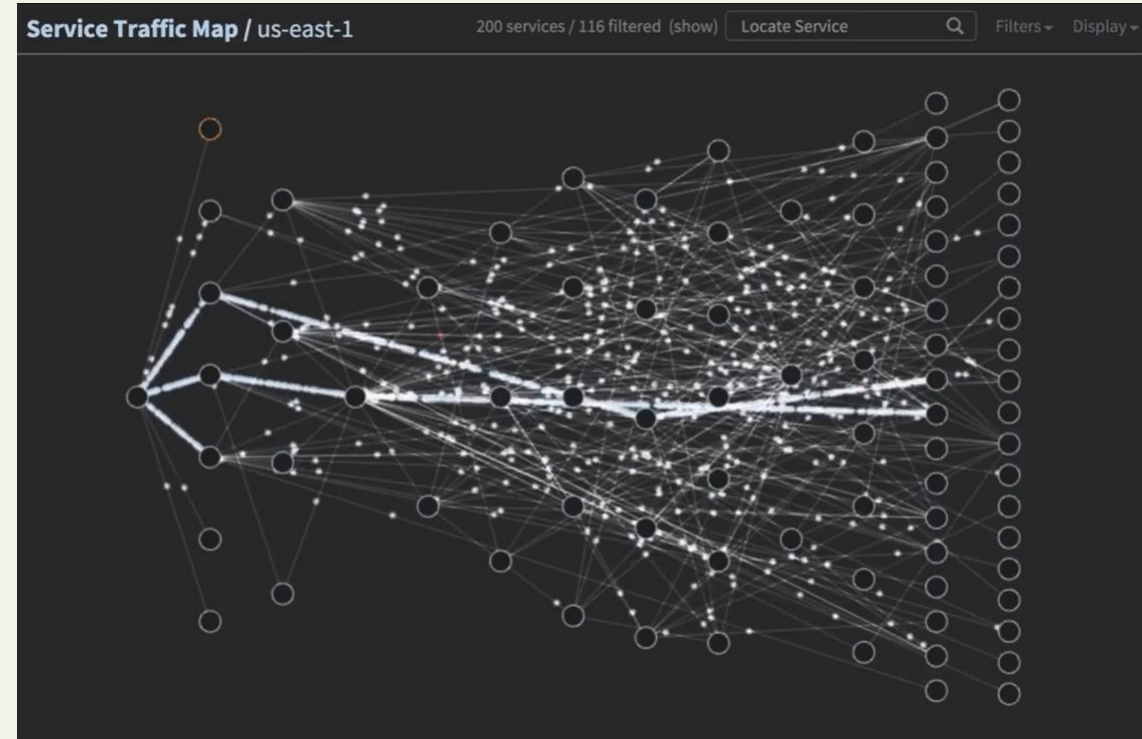
Twitter



Amazon



Social Network



Real time traffic of Social Network Microservices

- Complicate cluster management due to microservice dependencies



- Performance unpredictability usually occurs at large scale



- Performance unpredictability usually occurs at large scale
- Leverage simulation

- **Performance unpredictability usually occurs at large scale**

- **Leverage simulation**
 - Architecture simulator: Gem5, Zsim (ISCA'13)
 - Cluster level queueing simulator: BigHouse (ISPASS'12)

- **Performance unpredictability usually occurs at large scale**

- **Leverage simulation**
 - Architecture simulator: Gem5, Zsim (ISCA'13)
 - Cluster level queueing simulator: BigHouse (ISPASS'12)

- **Microservices introduce new simulation requirements**

- **Performance unpredictability usually occurs at large scale**

- **Leverage simulation**
 - Architecture simulator: Gem5, Zsim (ISCA'13)
 - Cluster level queueing simulator: BigHouse (ISPASS'12)

- **Microservices introduce new simulation requirements**
 - Scalability to study cluster management

- **Performance unpredictability usually occurs at large scale**

- **Leverage simulation**
 - Architecture simulator: Gem5, Zsim (ISCA'13)
 - Cluster level queueing simulator: BigHouse (ISPASS'12)

- **Microservices introduce new simulation requirements**
 - Scalability to study cluster management
 - Accurate multi-stage models for individual service/microservice

- **Performance unpredictability usually occurs at large scale**

- **Leverage simulation**
 - Architecture simulator: Gem5, Zsim (ISCA'13)
 - Cluster level queueing simulator: BigHouse (ISPASS'12)

- **Microservices introduce new simulation requirements**
 - Scalability to study cluster management
 - Accurate multi-stage models for individual service/microservice
 - Modeling arbitrary dependency graphs and dataflow paths

- **Performance unpredictability usually occurs at large scale**

- **Leverage simulation**
 - Architecture simulator: Gem5, Zsim (ISCA'13)
 - Cluster level queueing simulator: BigHouse (ISPASS'12)

- **Microservices introduce new simulation requirements**
 - Scalability to study cluster management
 - Accurate multi-stage models for individual service/microservice
 - Modeling arbitrary dependency graphs and dataflow paths
 - Modeling blocking and synchronization behavior

- **Accurate yet scalable service models**
 - Event-driven queueing simulator



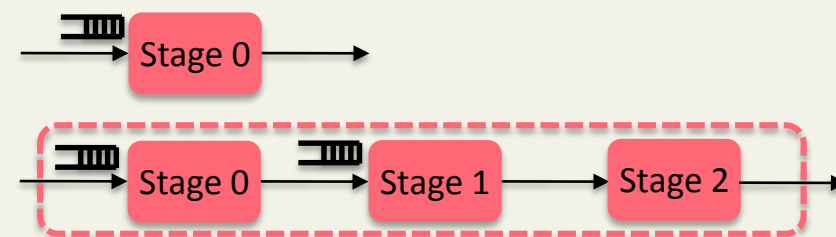
▪ **Accurate yet scalable service models**

- Event-driven queueing simulator
- Simple microservice: single stage



▪ **Accurate yet scalable service models**

- Event-driven queueing simulator
- Simple microservice: single stage
- Complex service: decomposed to multiple stages

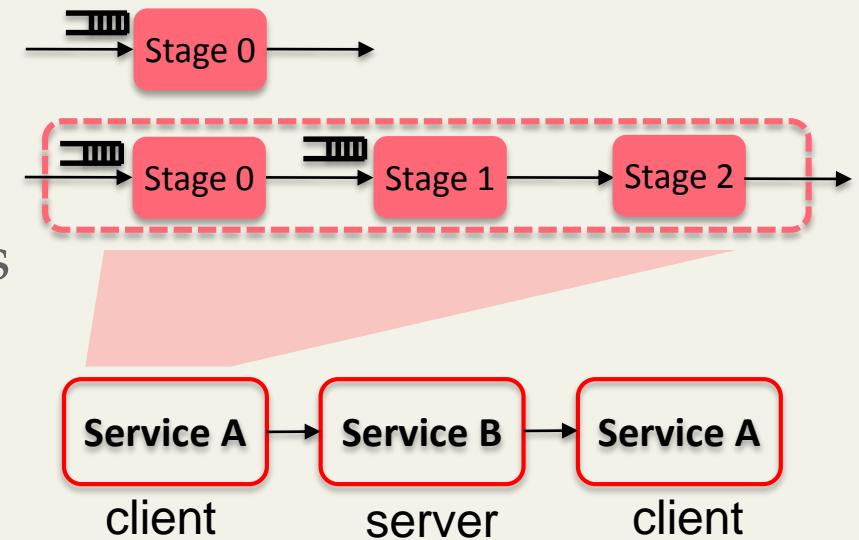


▪ Accurate yet scalable service models

- Event-driven queueing simulator
- Simple microservice: single stage
- Complex service: decomposed to multiple stages

▪ Modeling dependency graph & dataflow path

- Specify dependency graph & dataflow paths



▪ Accurate yet scalable service models

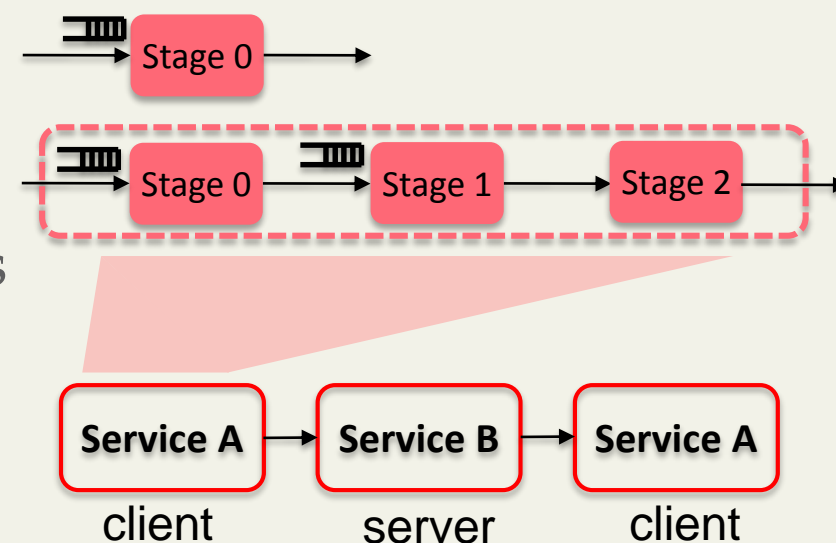
- Event-driven queueing simulator
- Simple microservice: single stage
- Complex service: decomposed to multiple stages

▪ Modeling dependency graph & dataflow path

- Specify dependency graph & dataflow paths

▪ Modeling blocking/synchronization

- Encode blocking/synchronization behavior in dataflow path



▪ Accurate yet scalable service models

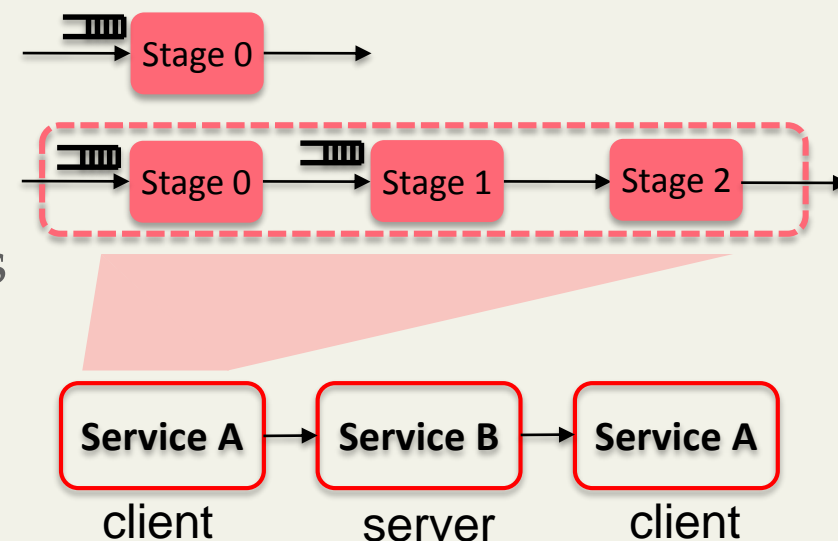
- Event-driven queueing simulator
- Simple microservice: single stage
- Complex service: decomposed to multiple stages

▪ Modeling dependency graph & dataflow path

- Specify dependency graph & dataflow paths

▪ Modeling blocking/synchronization

- Encode blocking/synchronization behavior in dataflow path
- Model sources of blocking: network connections, threads blocked by I/O accesses



- μ qSim microservice model

- Memcached



- **μ qSim microservice model**
 - Multiple stages per microservice

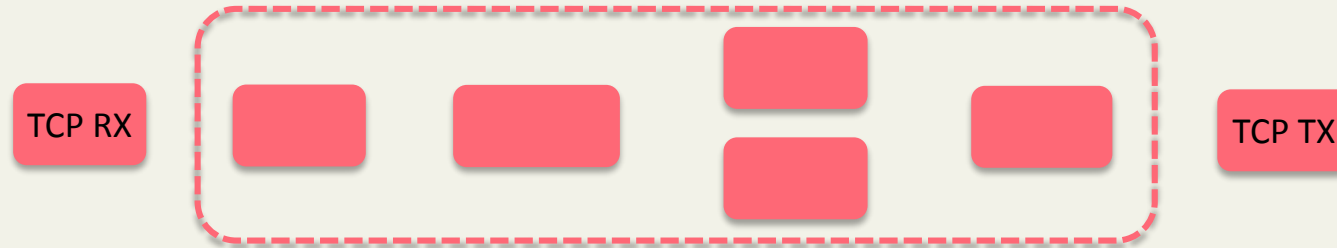
- **Memcached**



▪ μ qSim microservice model

- Multiple stages per microservice

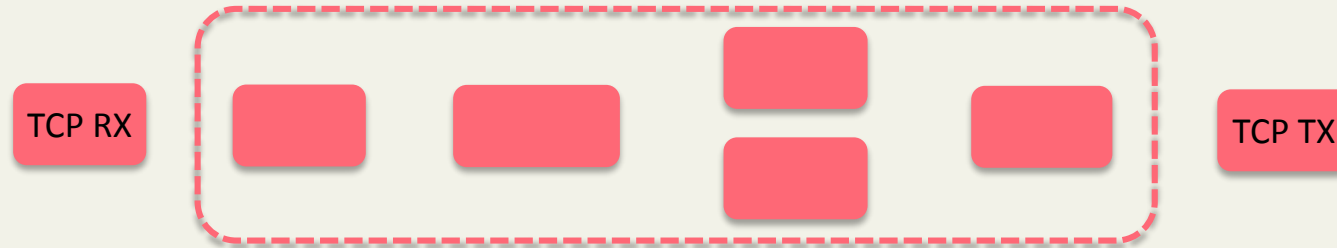
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)

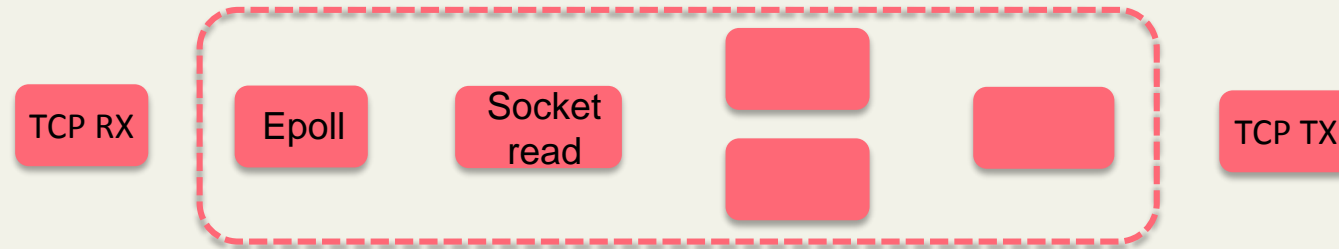
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)

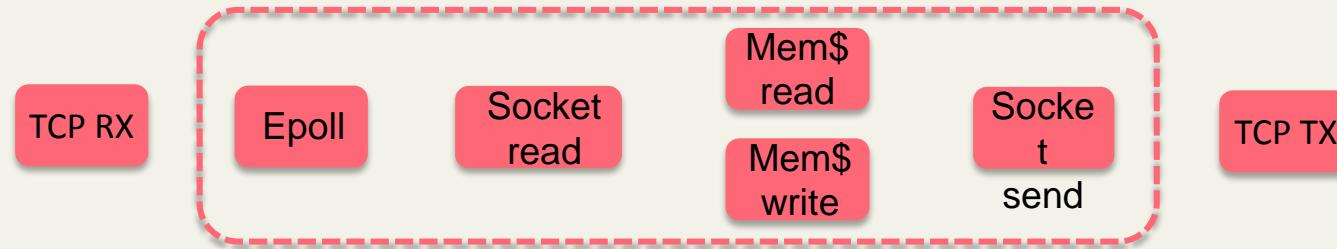
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)

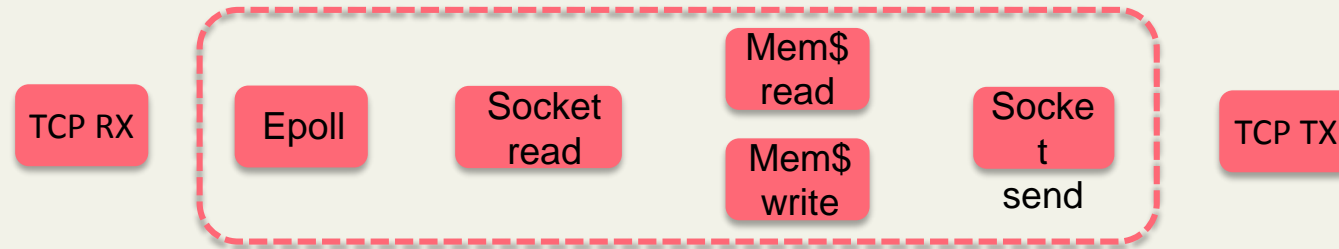
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue

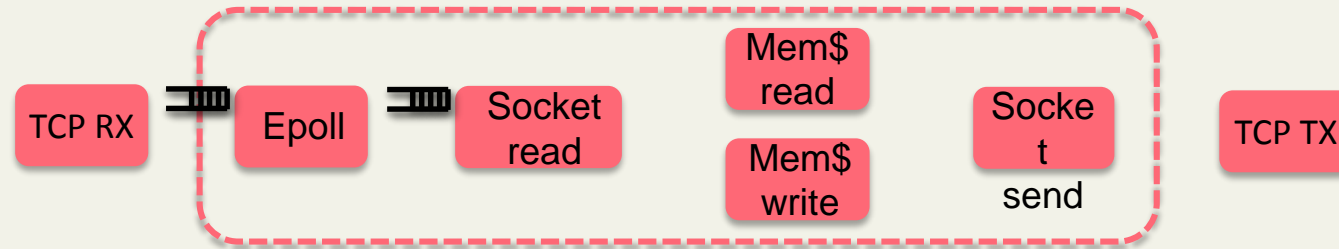
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue

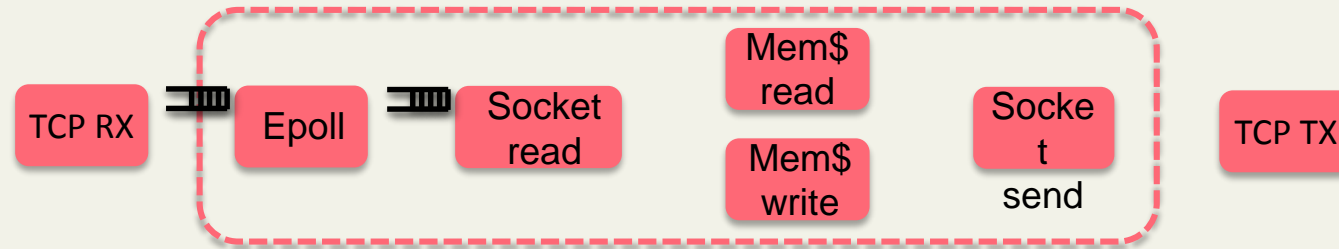
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

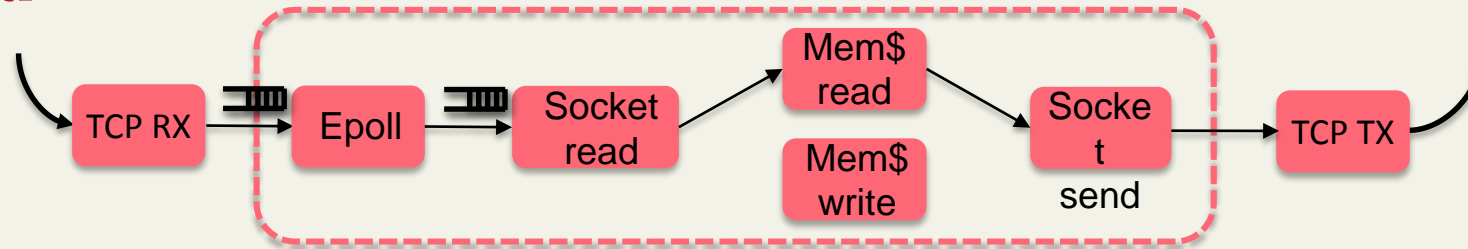
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

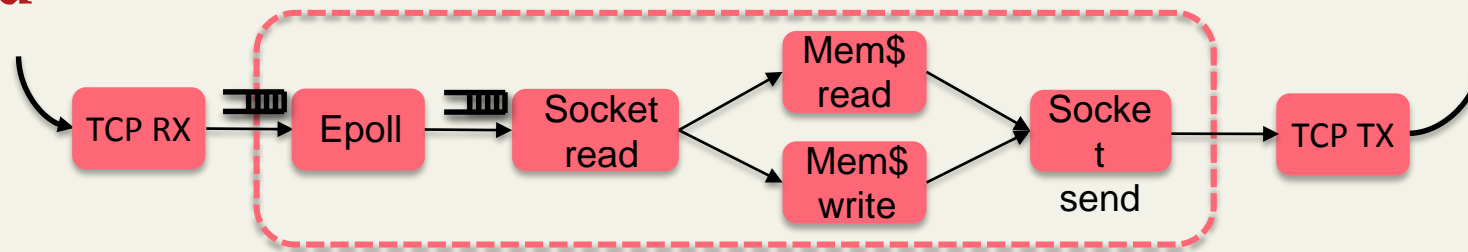
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

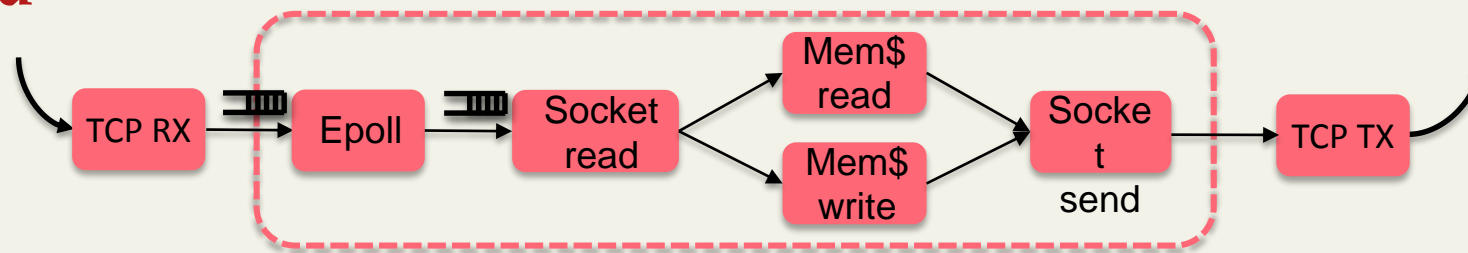
▪ Memcached



▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



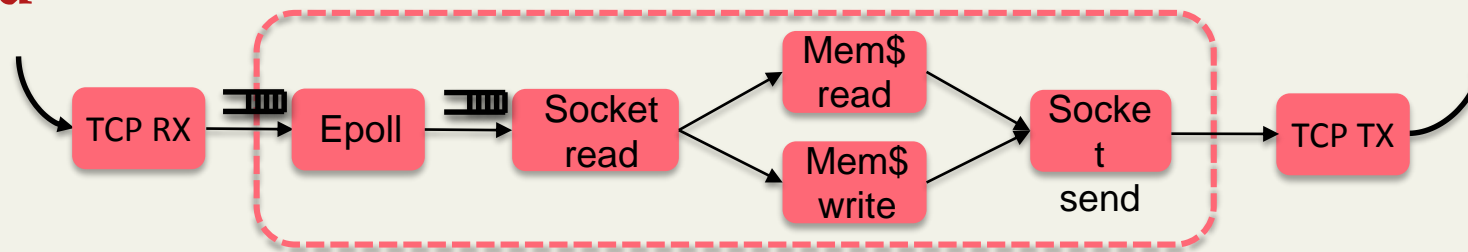
Connection A  1 2

Connection B  3

▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



Connection A 

1 2

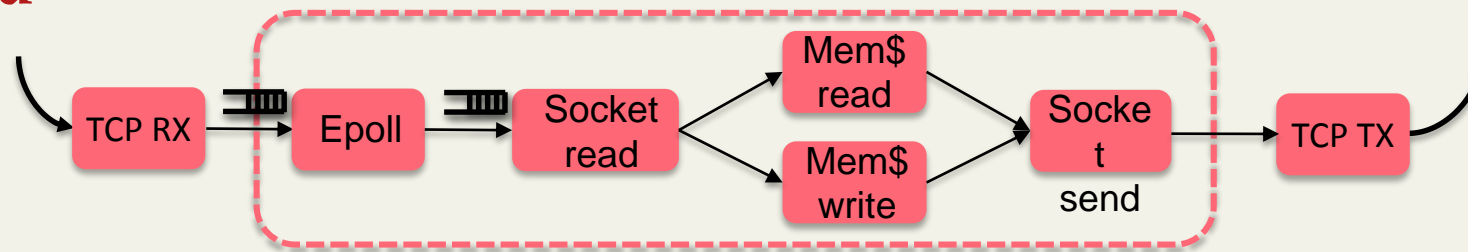
Connection B 

3

▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



Connection A 

1 2

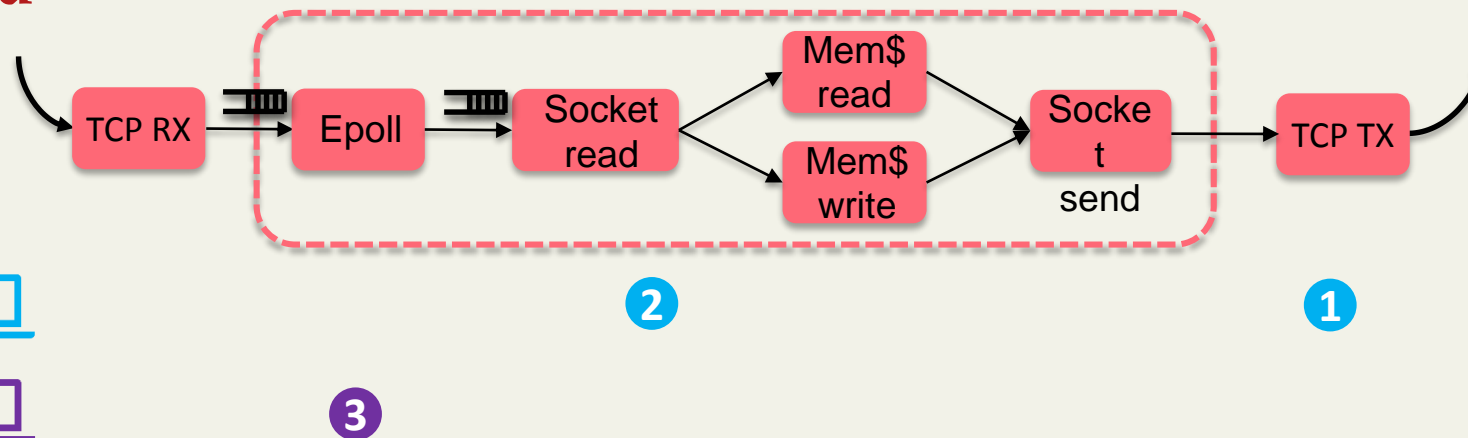
Connection B 

3

▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



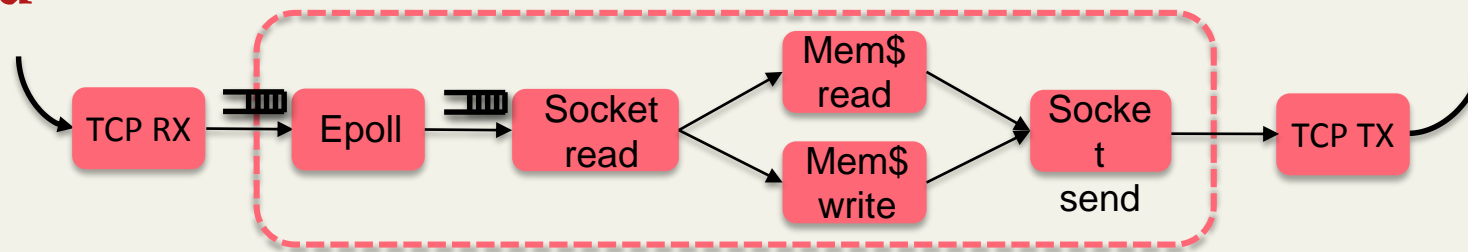
Connection A 

Connection B 

▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



Connection A 

Connection B 

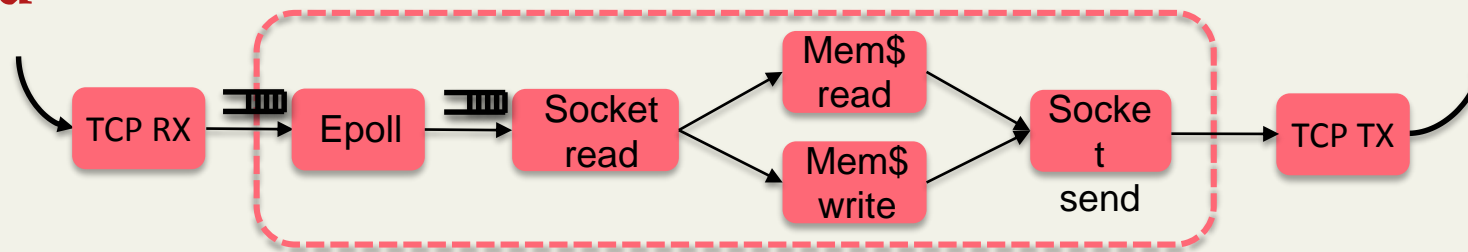
1 2

3

▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



Connection A 

Connection B 

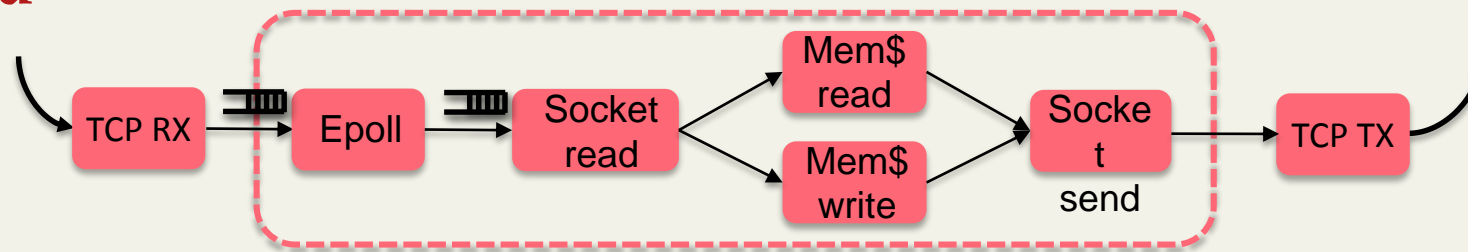
1 2

3

▪ μ qSim microservice model

- Multiple stages per microservice
- Provided models for common OS queueing utilities (e.g., epoll & socket read)
- Each stage optionally coupled with a queue
- Multiple execution paths (e.g., Memcached read/write) per microservice

▪ Memcached



Connection A 

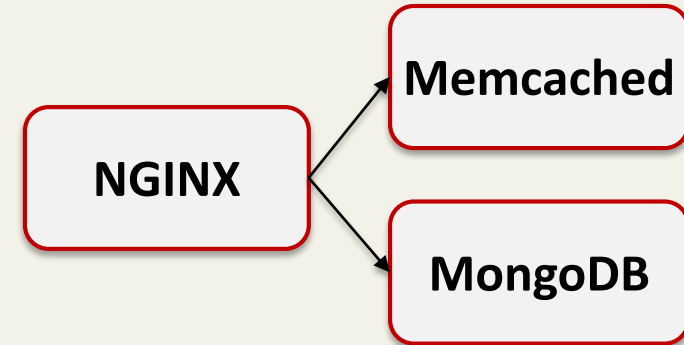
Connection B 

1 2

3



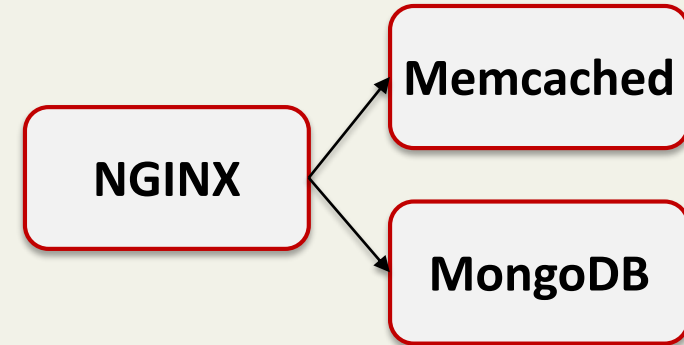
- **Microservice dependency graph**



▪ **Microservice dependency graph**

▪ **Dataflow**

- Sequence of microservices to execute
- Differ across request types
- Dataflow node encodes blocking/synchronization operation



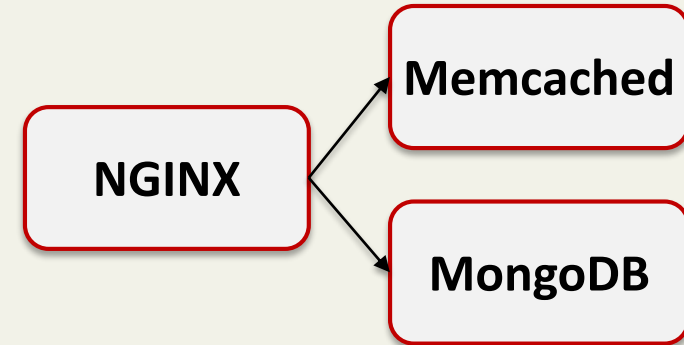
▪ **Microservice dependency graph**

▪ **Dataflow**

- Sequence of microservices to execute
- Differ across request types
- Dataflow node encodes blocking/synchronization operation

▪ **Deployment**

- Available servers and hardware resources
- Service to server mapping
- Services on the same server share network stack & disk I/O





- Example: 2-tier application & http 1/1.1 protocol



- **Example: 2-tier application & http 1/1.1 protocol**

- Dependency graph

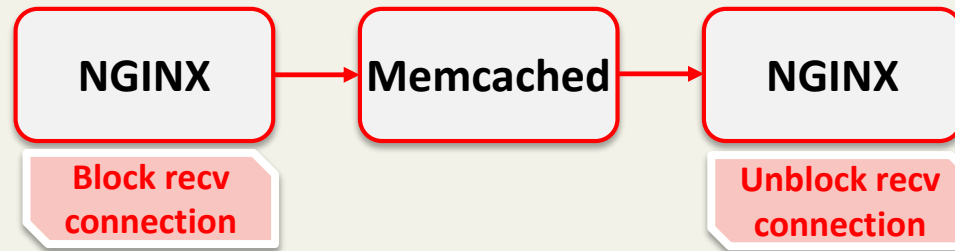


▪ Example: 2-tier application & http 1/1.1 protocol

- Dependency graph



- Dataflow

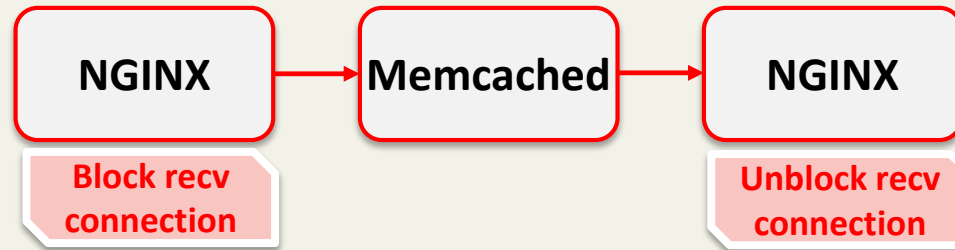


Example: 2-tier application & http 1/1.1 protocol

- Dependency graph



- Dataflow



Connection A  1 2

Connection B  3

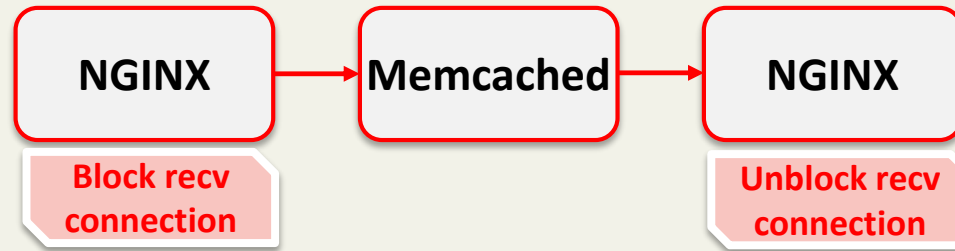


Example: 2-tier application & http 1/1.1 protocol

- Dependency graph

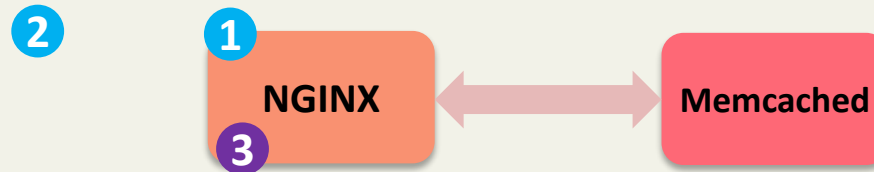


- Dataflow



Connection A 

Connection B 

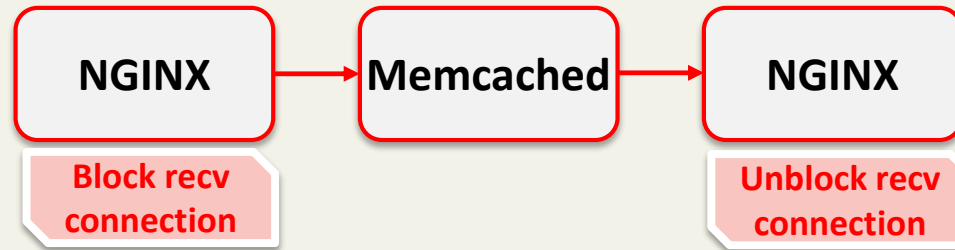


Example: 2-tier application & http 1/1.1 protocol

- Dependency graph

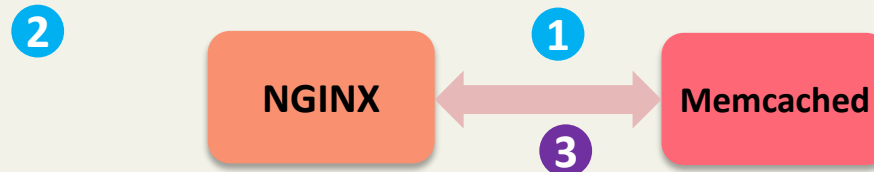


- Dataflow



~~Connection A~~ 

~~Connection B~~ 

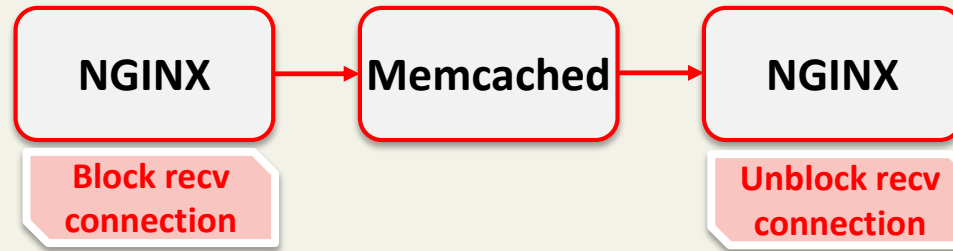


Example: 2-tier application & http 1/1.1 protocol

- Dependency graph



- Dataflow



~~Connection A~~ 

2



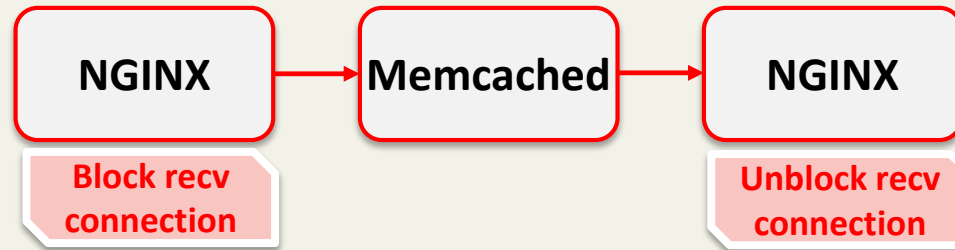
~~Connection B~~ 

Example: 2-tier application & http 1/1.1 protocol

- Dependency graph

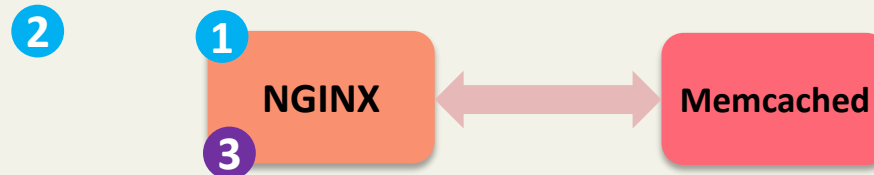


- Dataflow



Connection A 

Connection B 

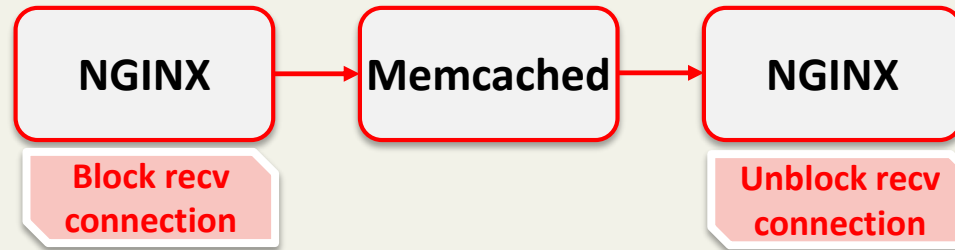


Example: 2-tier application & http 1/1.1 protocol

- Dependency graph



- Dataflow



Connection A  1 2

Connection B  3

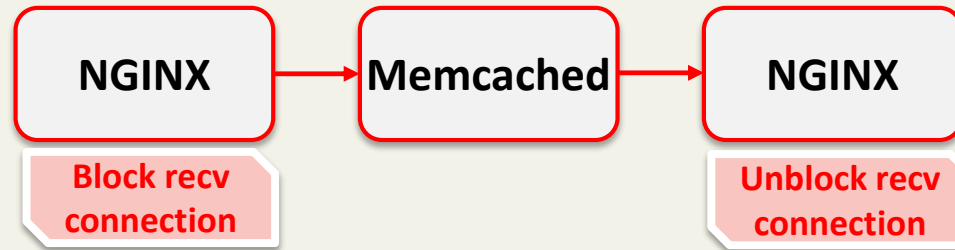


Example: 2-tier application & http 1/1.1 protocol

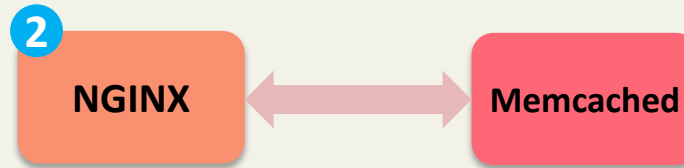
- Dependency graph



- Dataflow

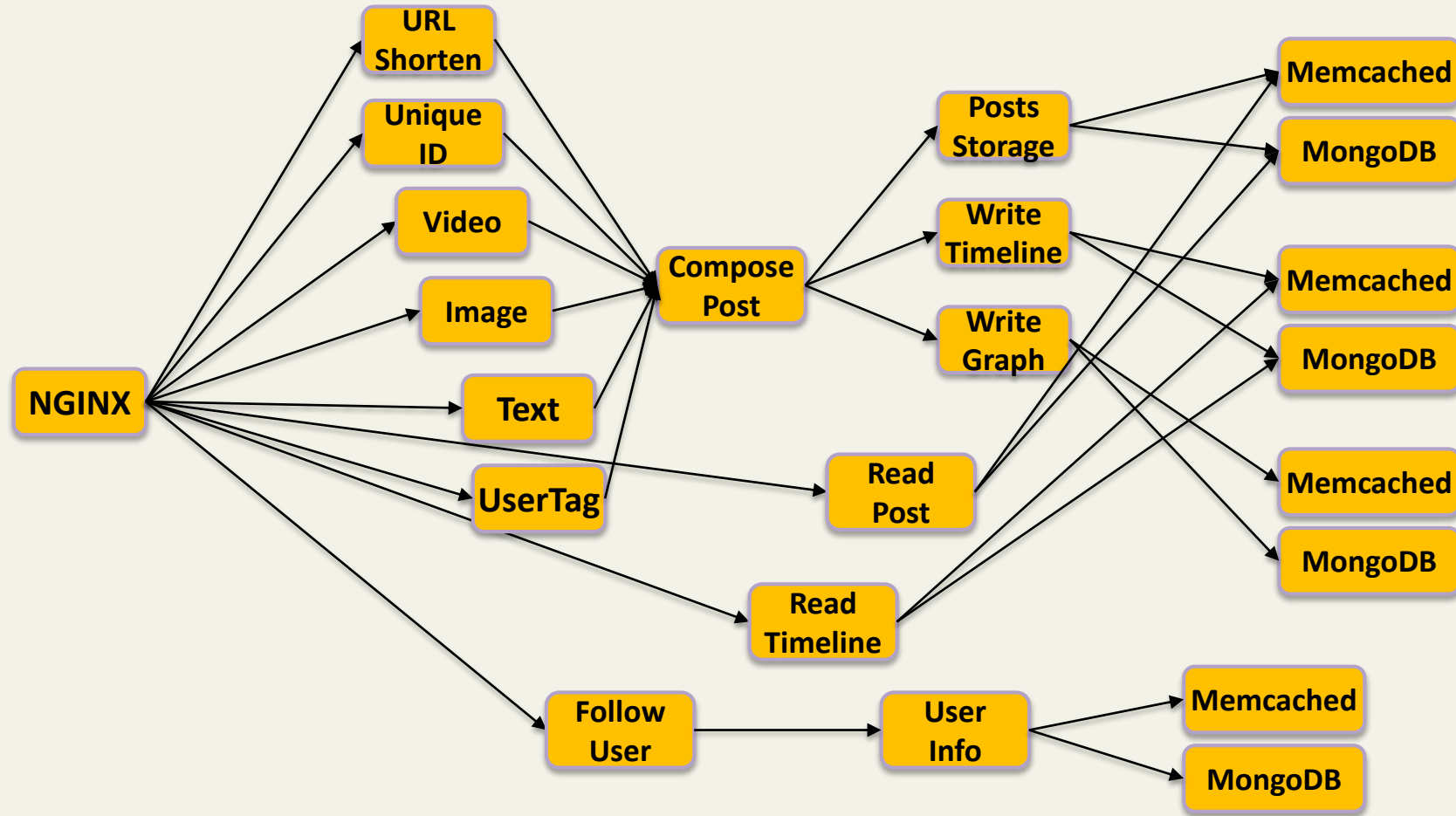


Connection A  1



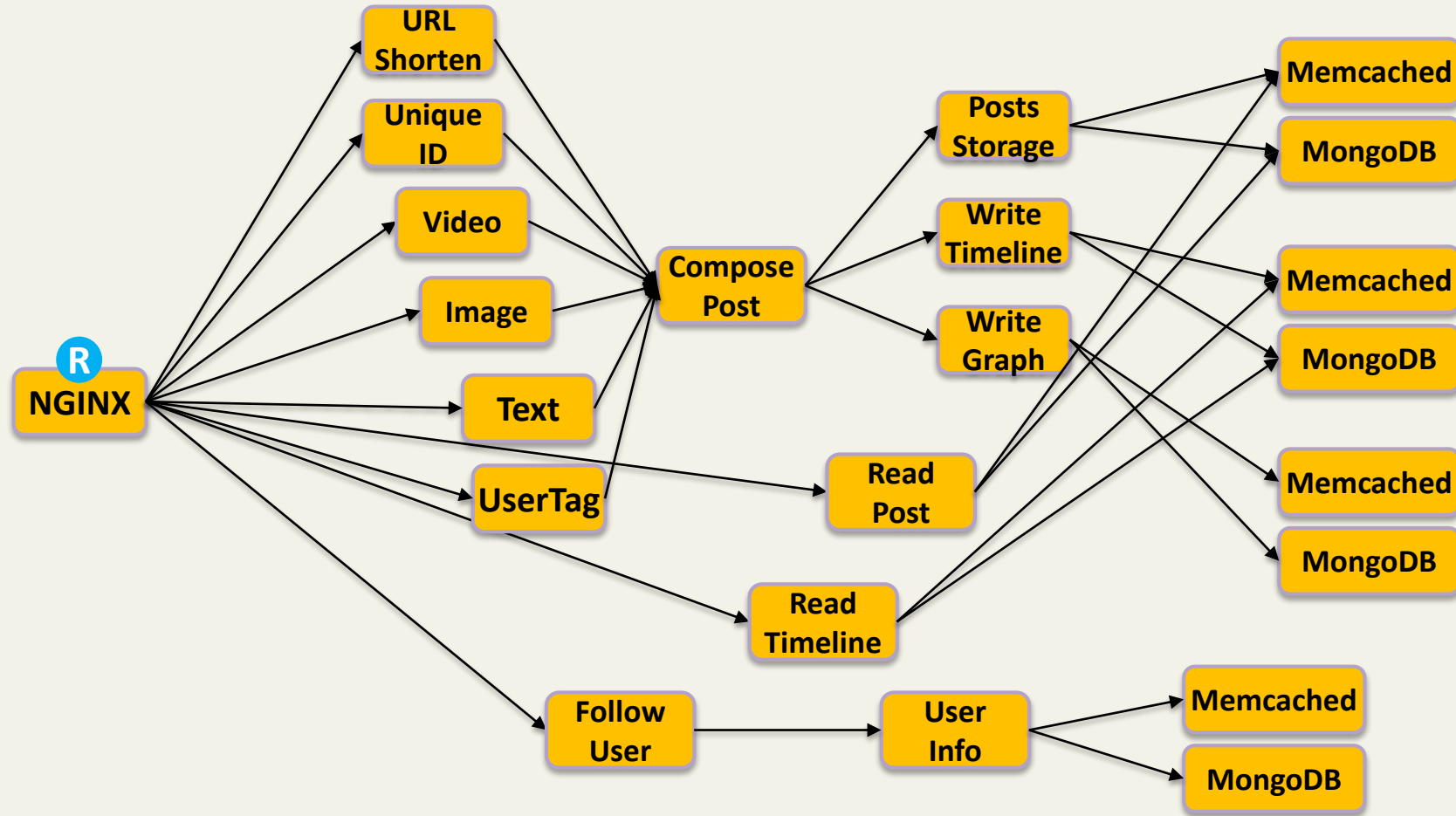
Connection B  3

▪ Social network application



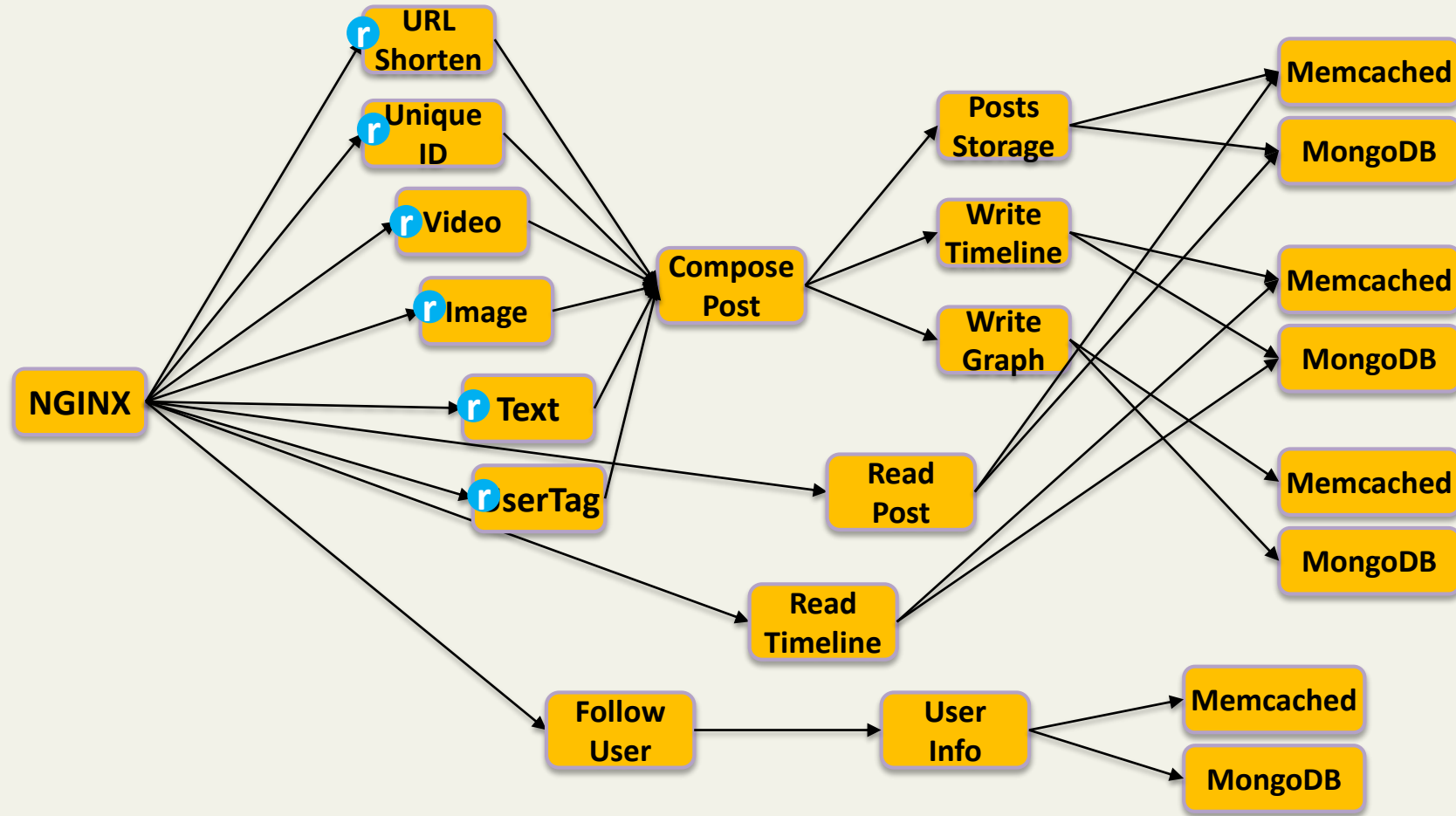
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



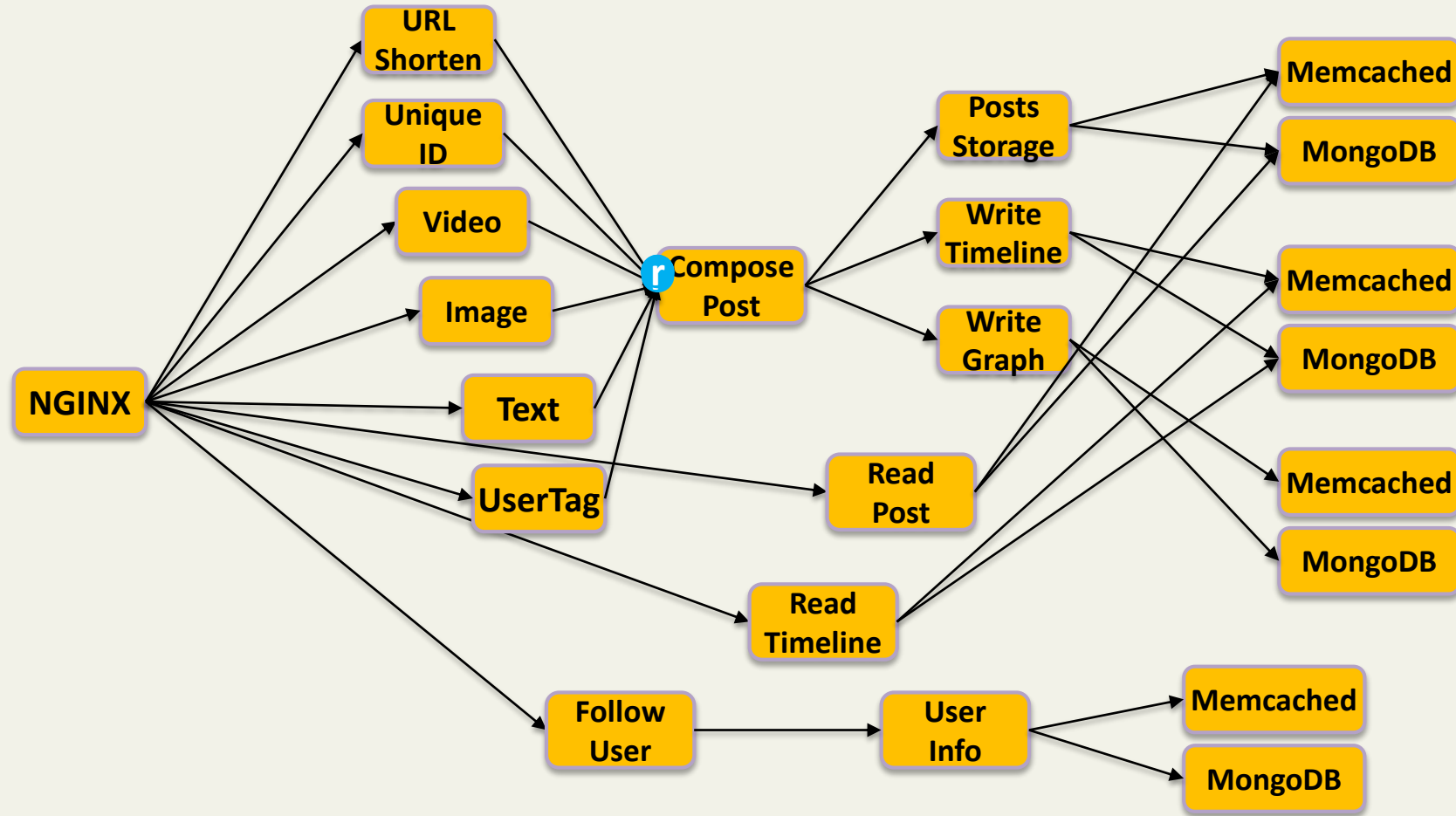
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



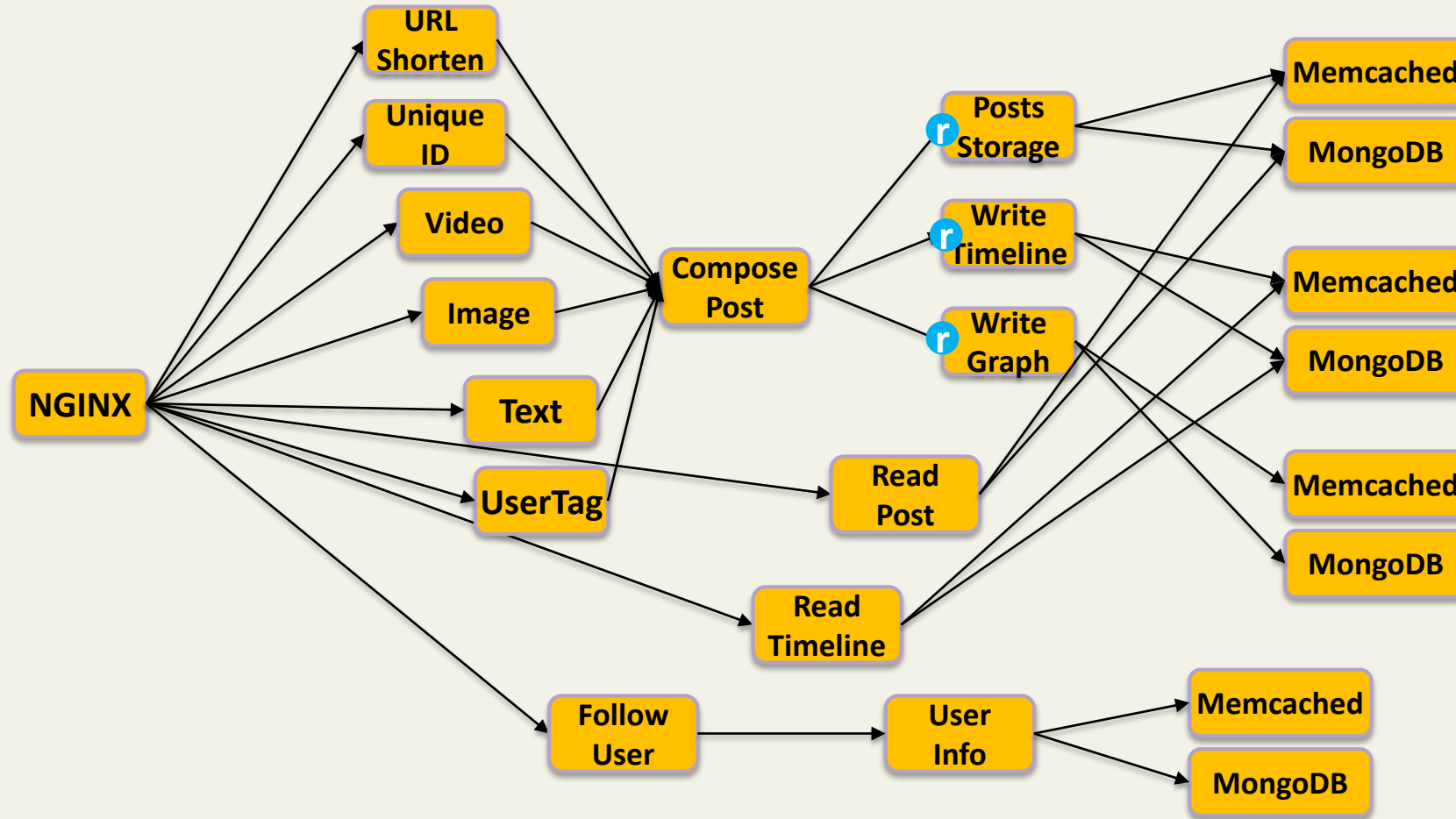
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



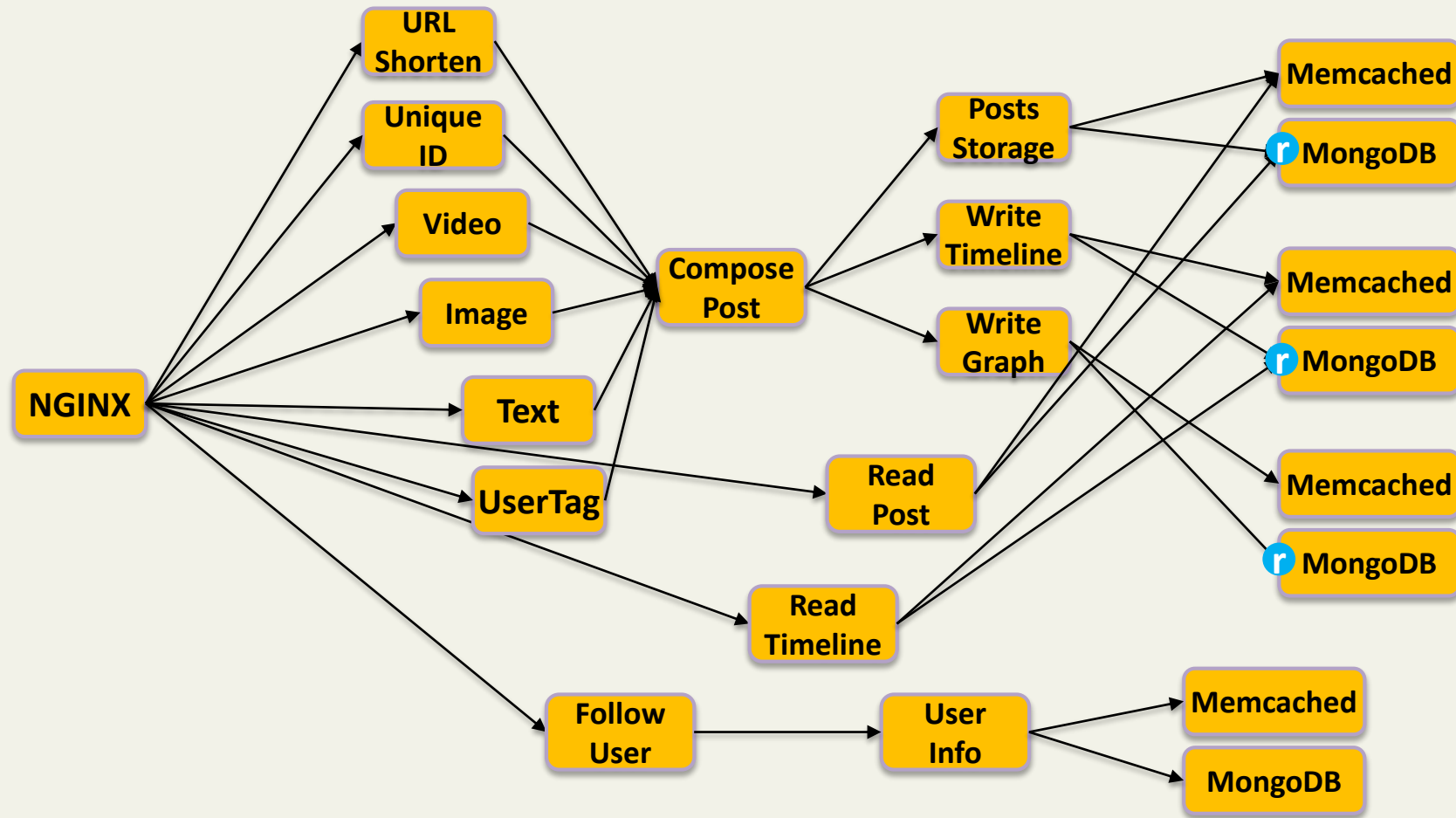
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



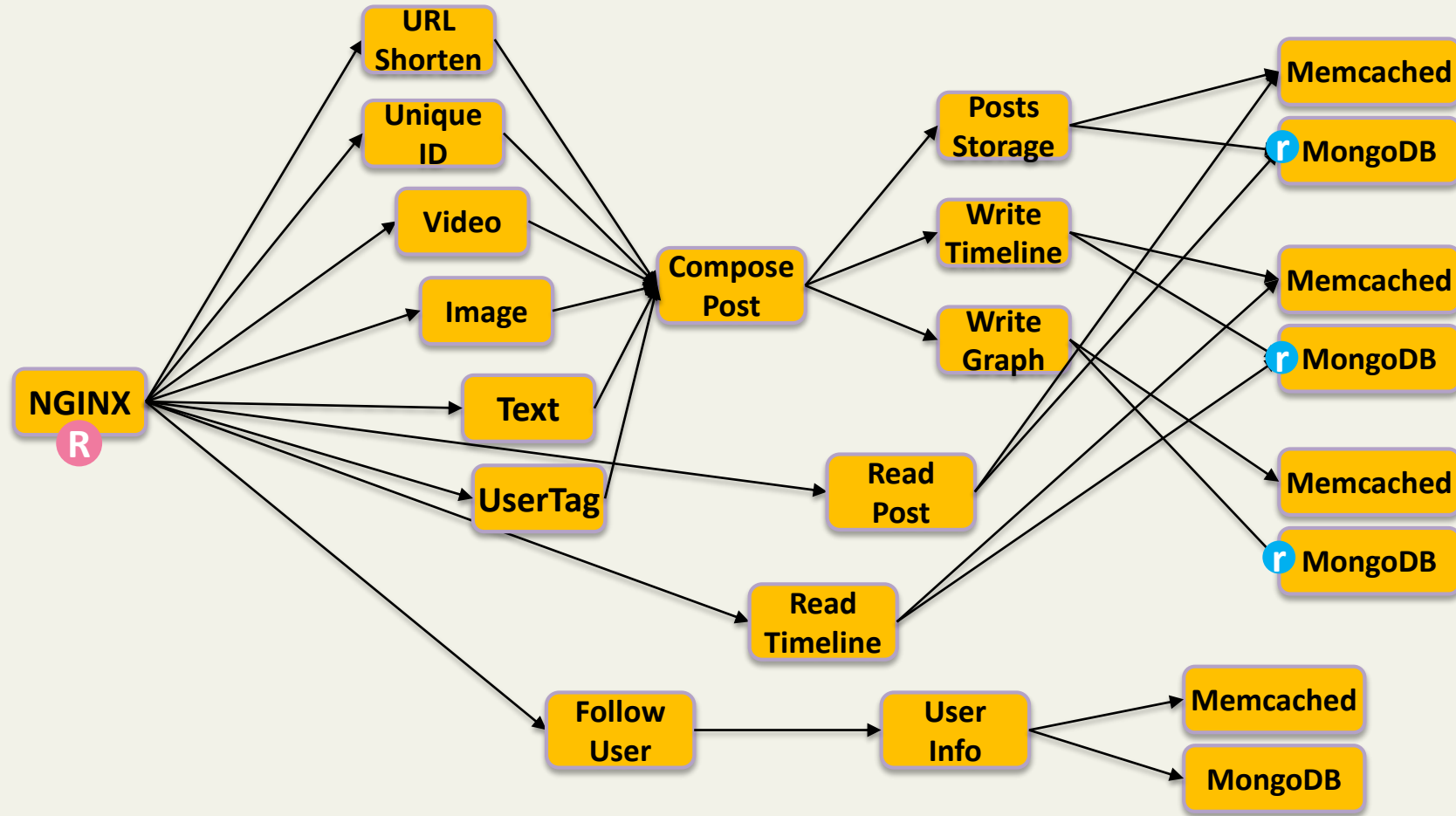
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



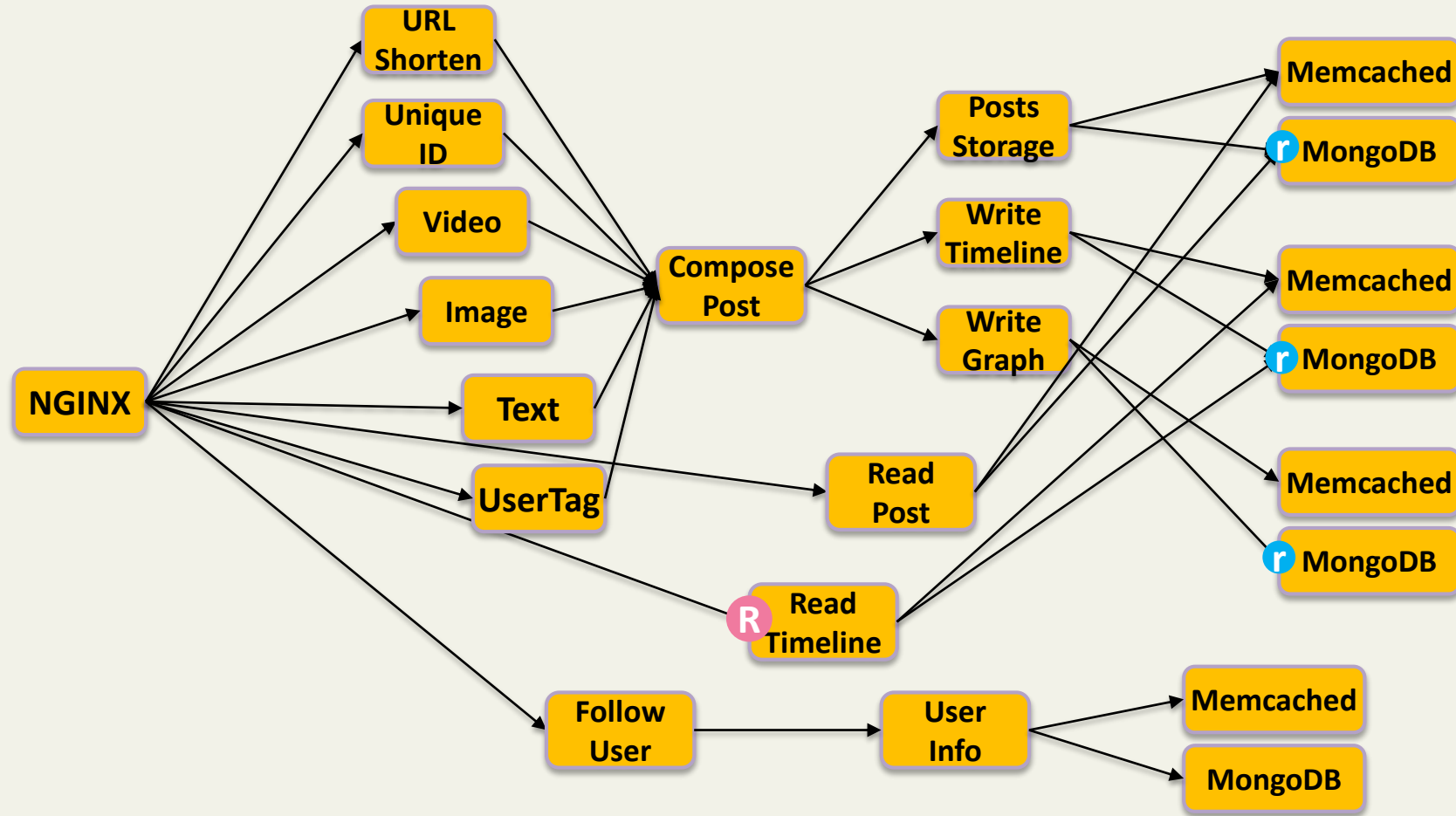
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



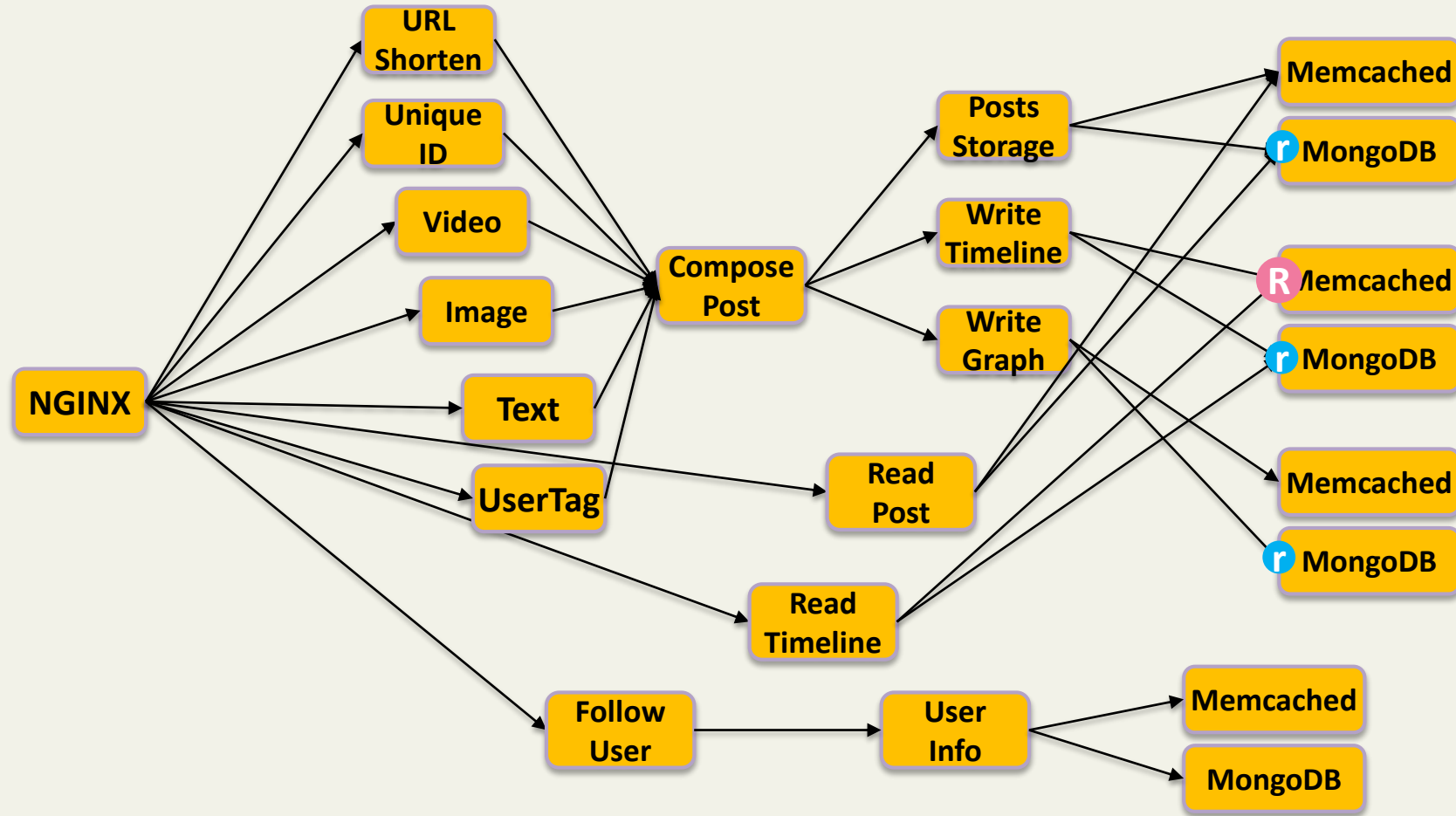
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



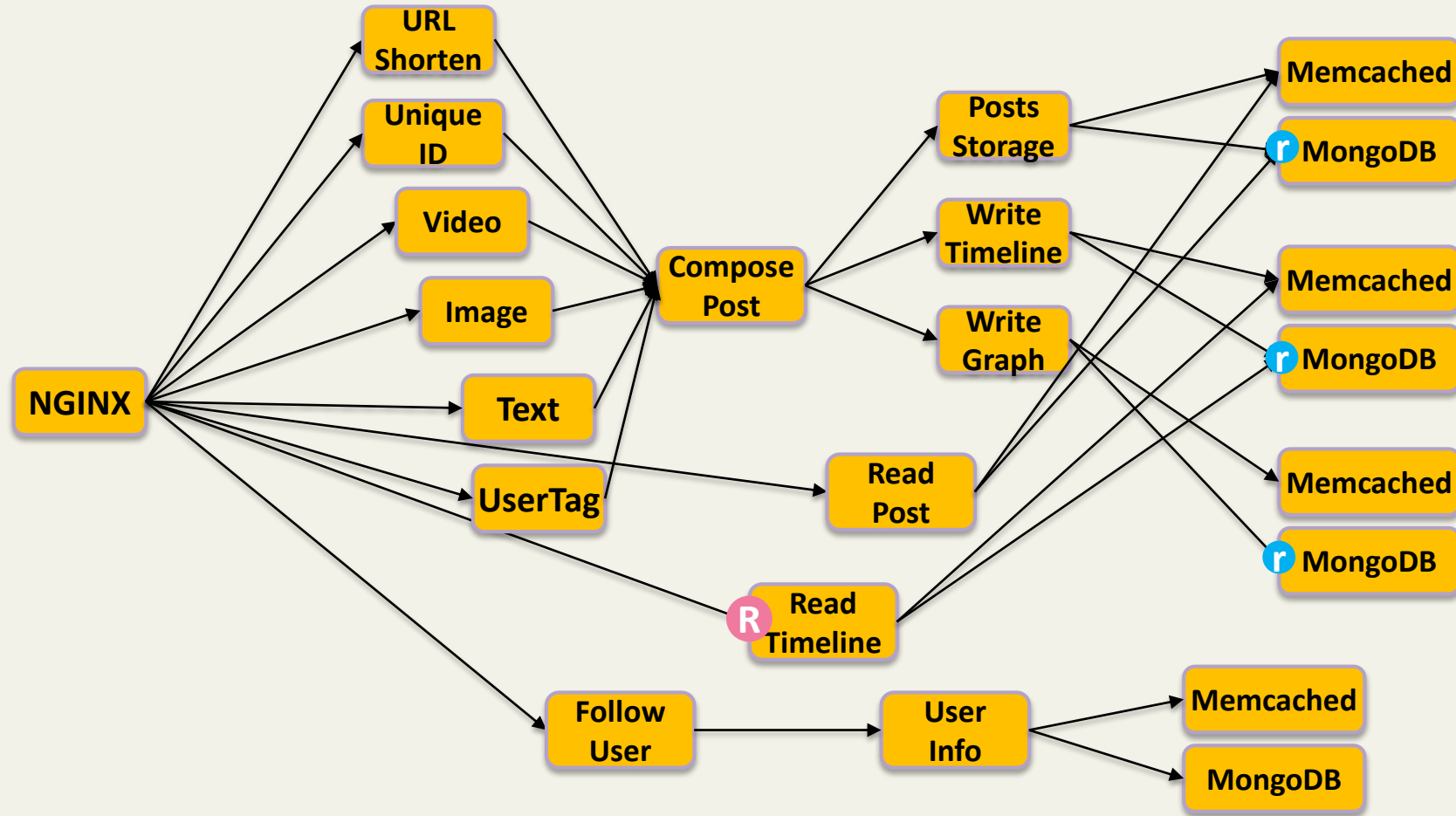
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



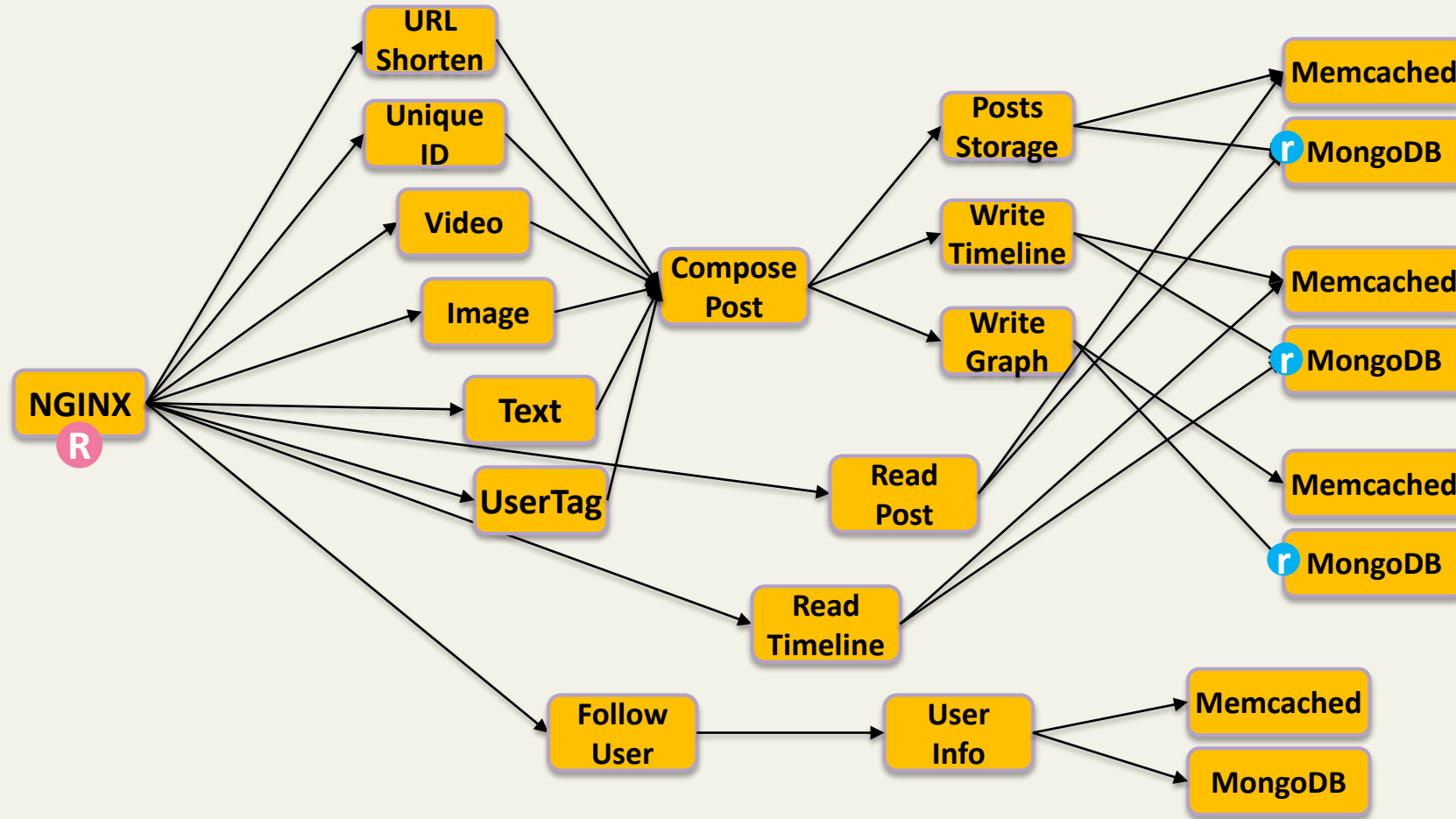
Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ Social network application



Source: <http://www.csl.cornell.edu/~delimitrou/papers/2019.asplos.microservices.pdf>

▪ **Microservice queueing model**

- Identifying sources of queueing
- Typical queues: network, epoll/kqueue, socket read, disk I/O...
- Microservice models reusable (open source community)



▪ **Microservice queueing model**

- Identifying sources of queueing
- Typical queues: network, epoll/kqueue, socket read, disk I/O...
- Microservice models reusable (open source community)

▪ **Processing time distribution**

- Instrumenting applications and profiling on real servers
- Instrumentations reusable

▪ **Microservice queueing model**

- Identifying sources of queueing
- Typical queues: network, epoll/kqueue, socket read, disk I/O...
- Microservice models reusable (open source community)

▪ **Processing time distribution**

- Instrumenting applications and profiling on real servers
- Instrumentations reusable

▪ **Microservice dependencies & dataflow paths**

- Obtained from app developers

▪ **Microservice queueing model**

- Identifying sources of queueing
- Typical queues: network, epoll/kqueue, socket read, disk I/O...
- Microservice models reusable (open source community)

▪ **Processing time distribution**

- Instrumenting applications and profiling on real servers
- Instrumentations reusable

▪ **Microservice dependencies & dataflow paths**

- Obtained from app developers

▪ **Server & system resources**

■ Validation experiments

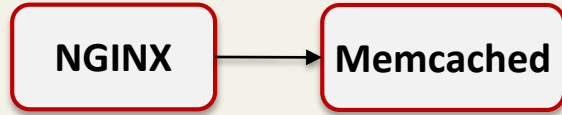
- Multi-tier microservices: 2/3-tier application
- Load balancing & fanout effects
- Microservices based on RPC
- Comparison with BigHouse

■ Server platform for trace collection

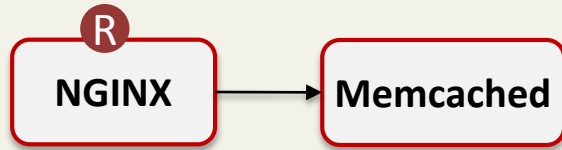
- 10-server cluster
- Intel(R) Xeon(R) CPU E5-2660 v3
- 2 sockets, 10 cores/socket, 2 threads/core
- Min/max DVFS frequency: 1.2GHz/2.6GHz
- Network bandwidth: 1Gbps



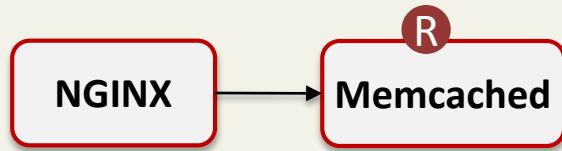
- 2-tier application



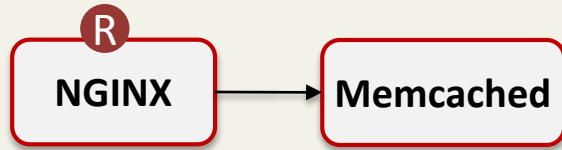
- 2-tier application



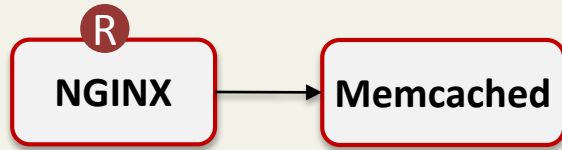
- 2-tier application



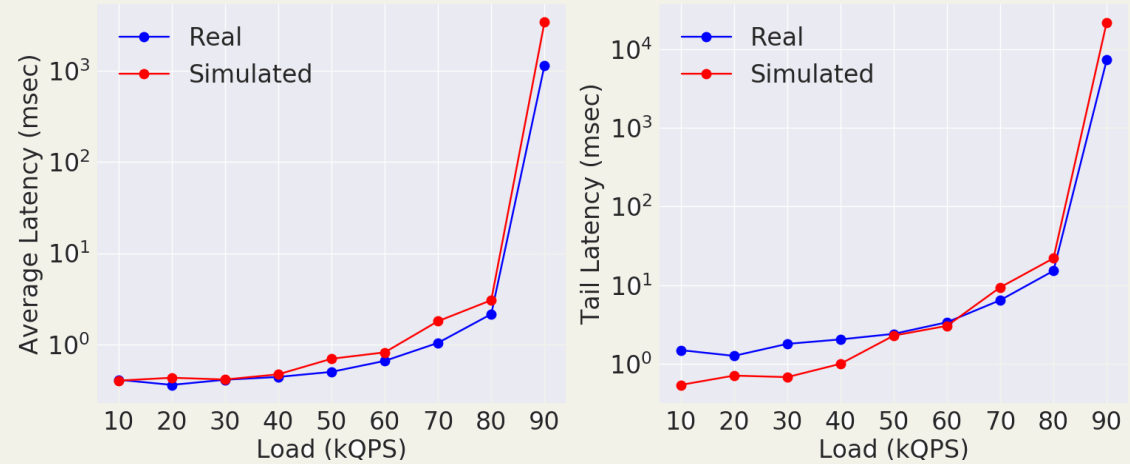
- 2-tier application



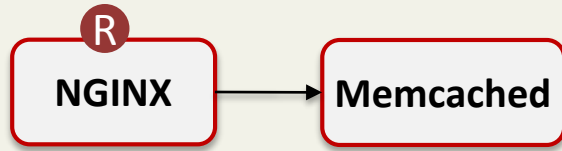
2-tier application



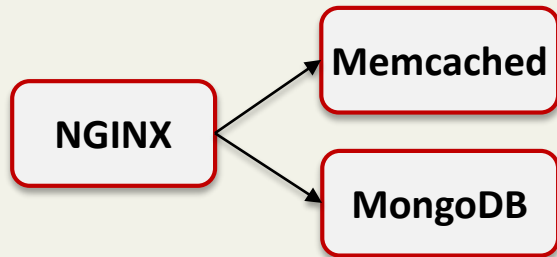
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



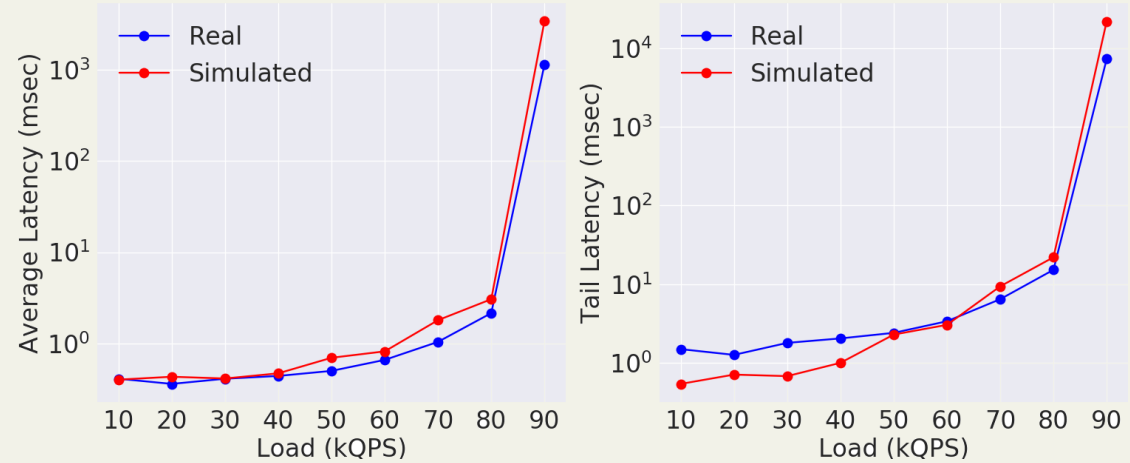
2-tier application



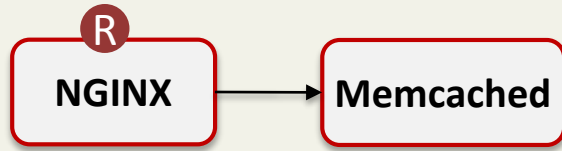
3-tier application



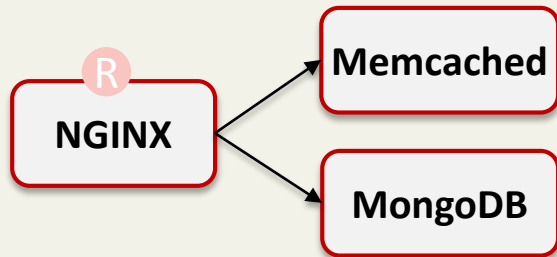
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



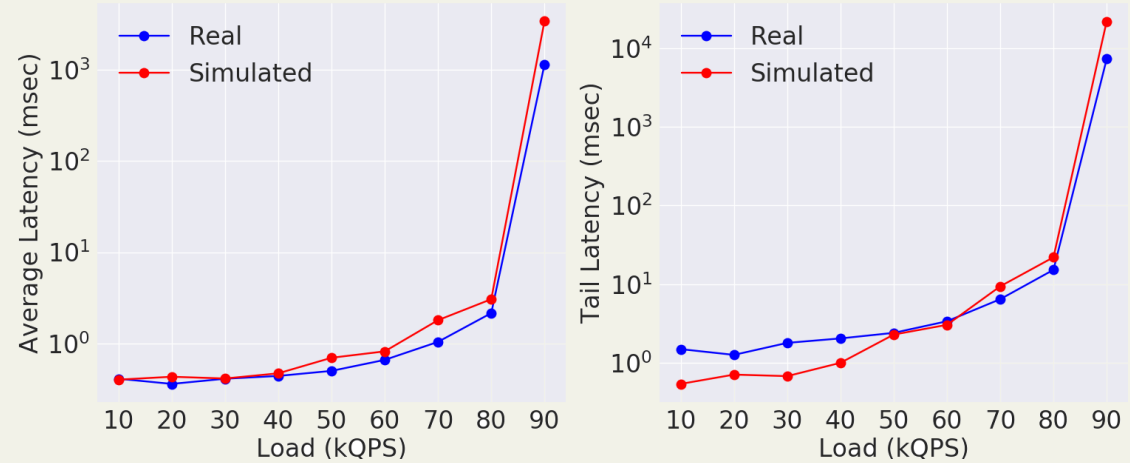
2-tier application



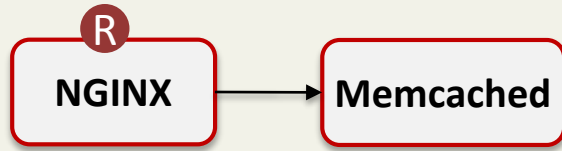
3-tier application



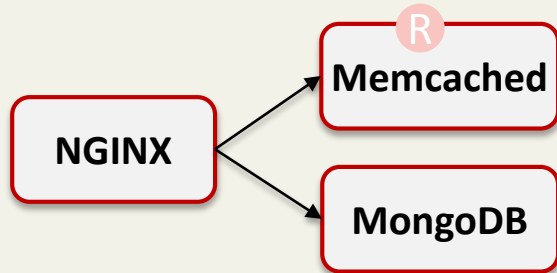
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



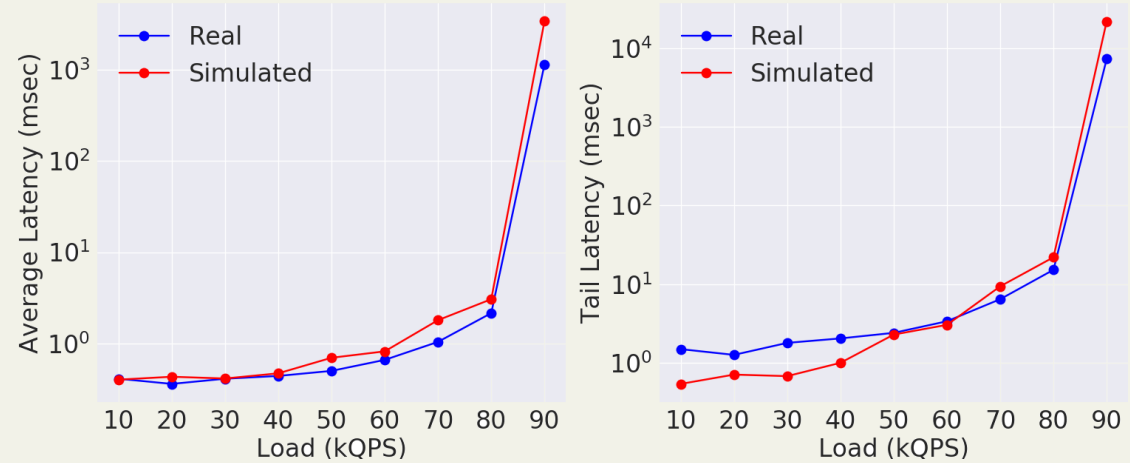
2-tier application



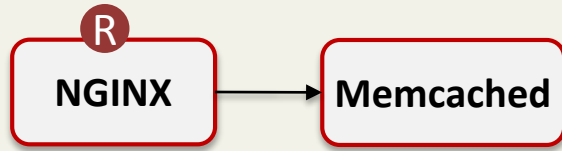
3-tier application



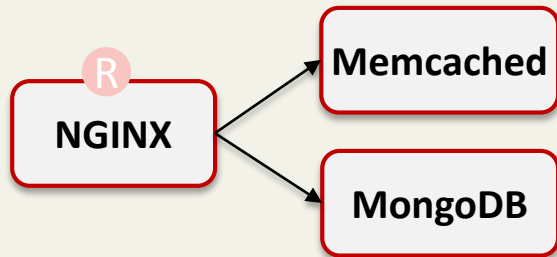
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



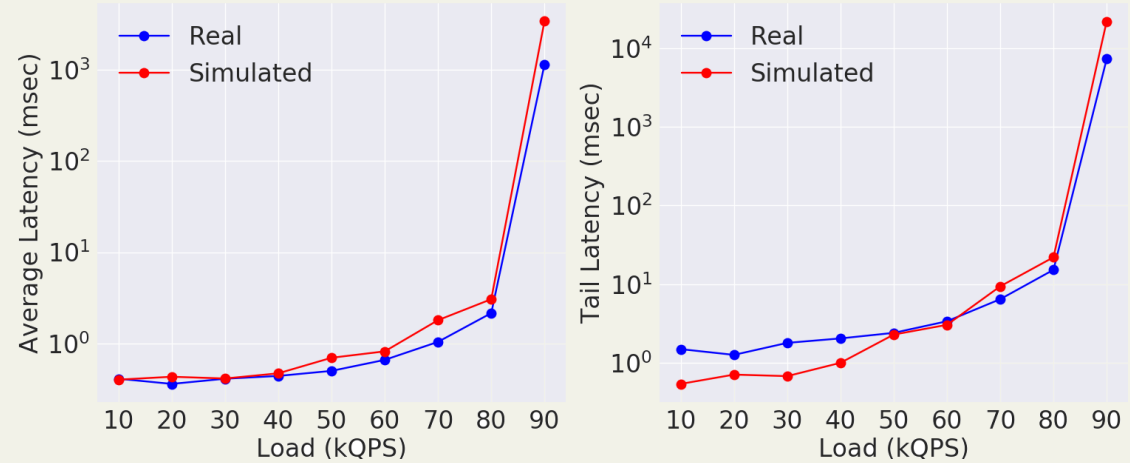
2-tier application



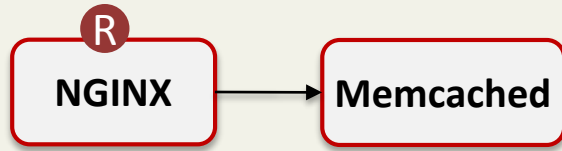
3-tier application



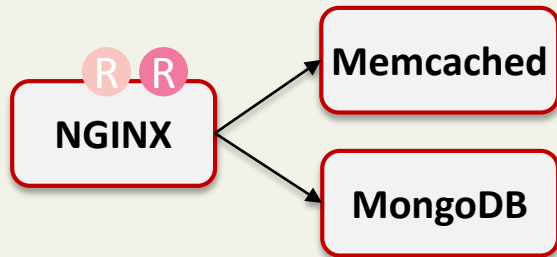
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



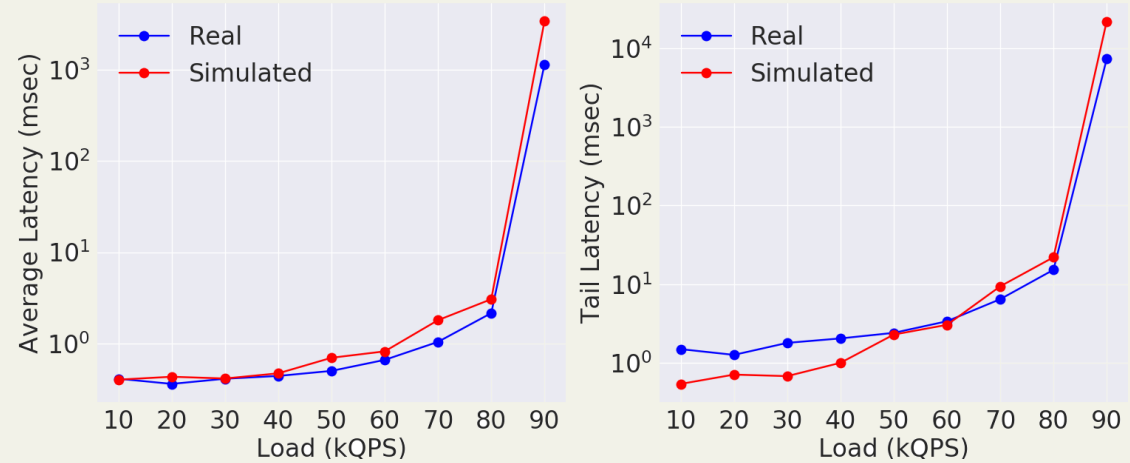
2-tier application



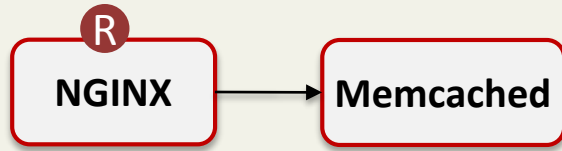
3-tier application



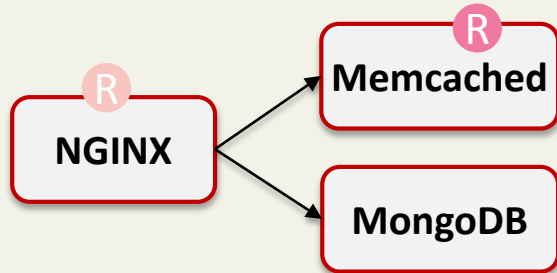
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



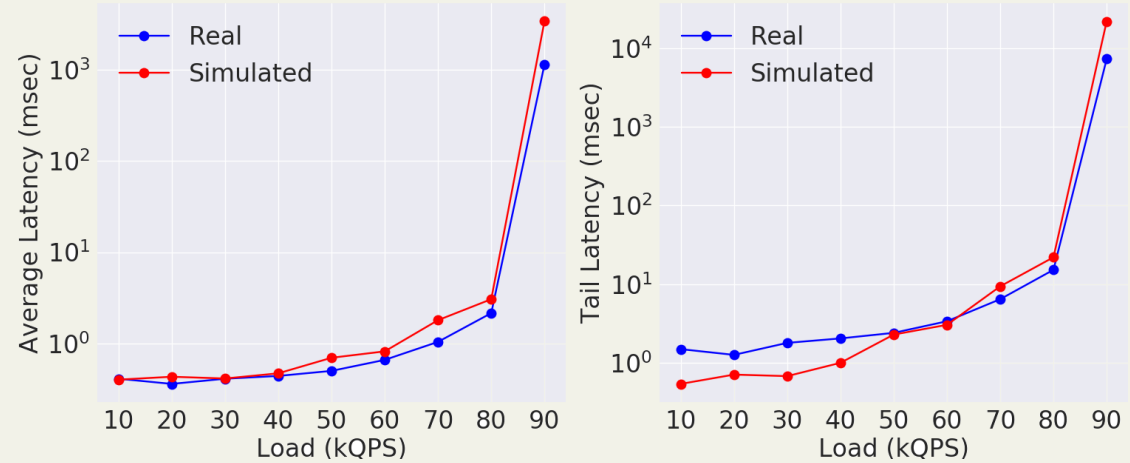
2-tier application



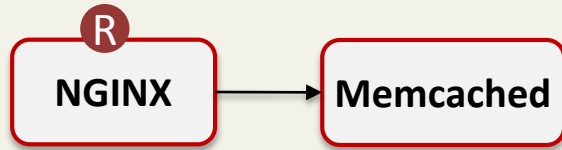
3-tier application



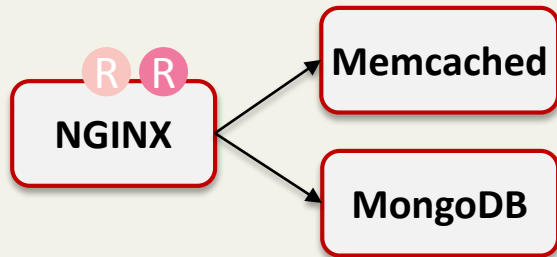
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



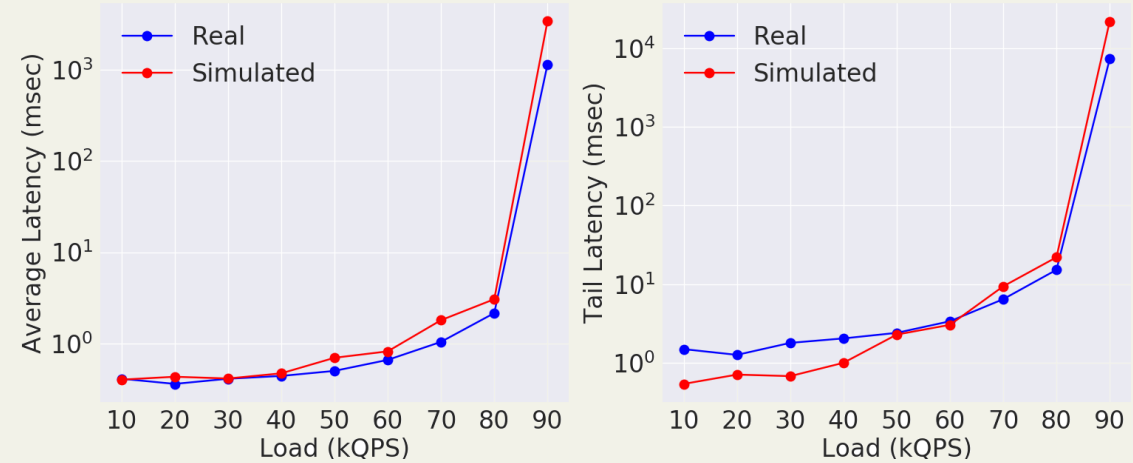
2-tier application



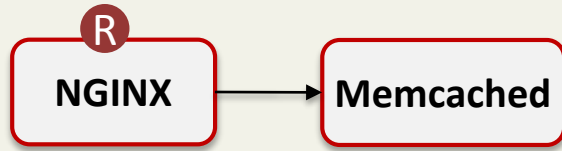
3-tier application



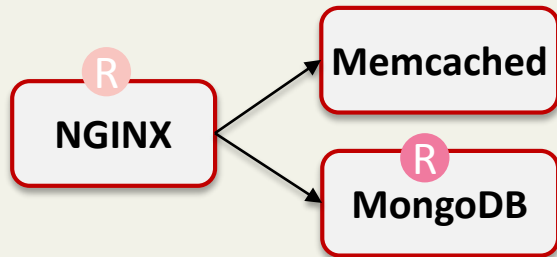
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



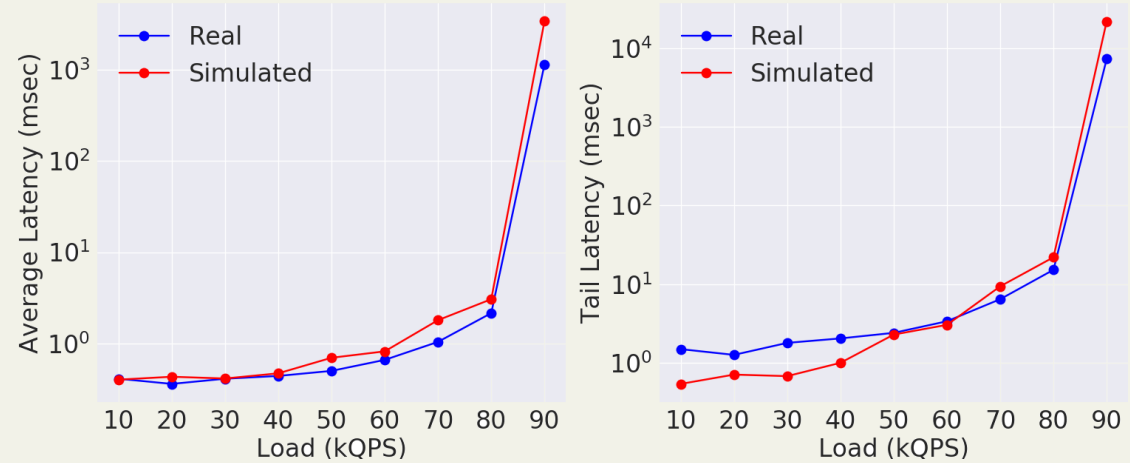
2-tier application



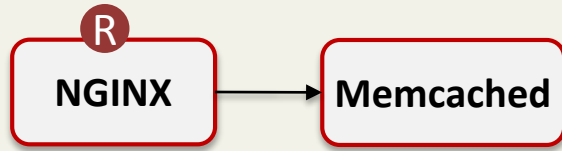
3-tier application



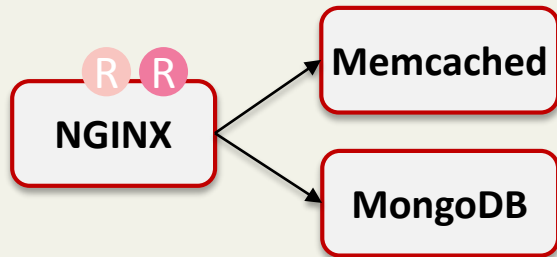
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



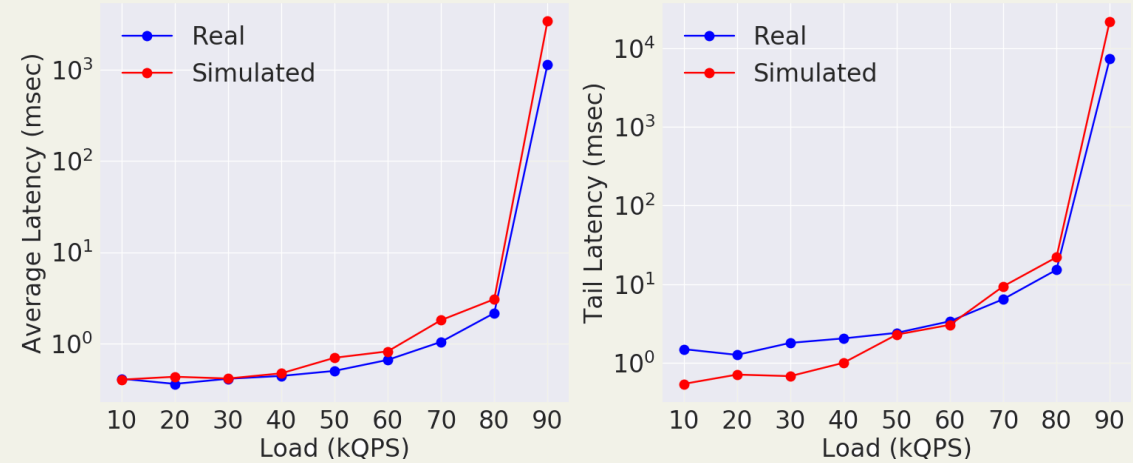
2-tier application



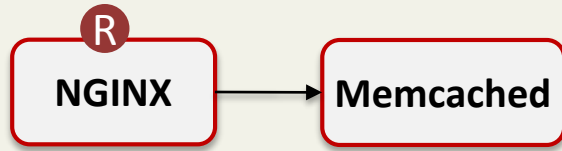
3-tier application



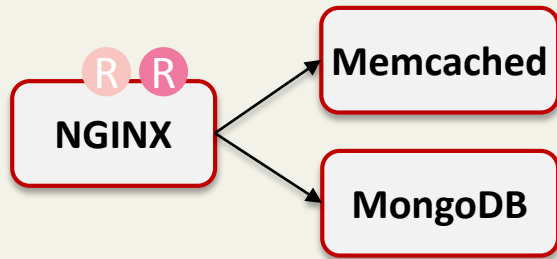
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)



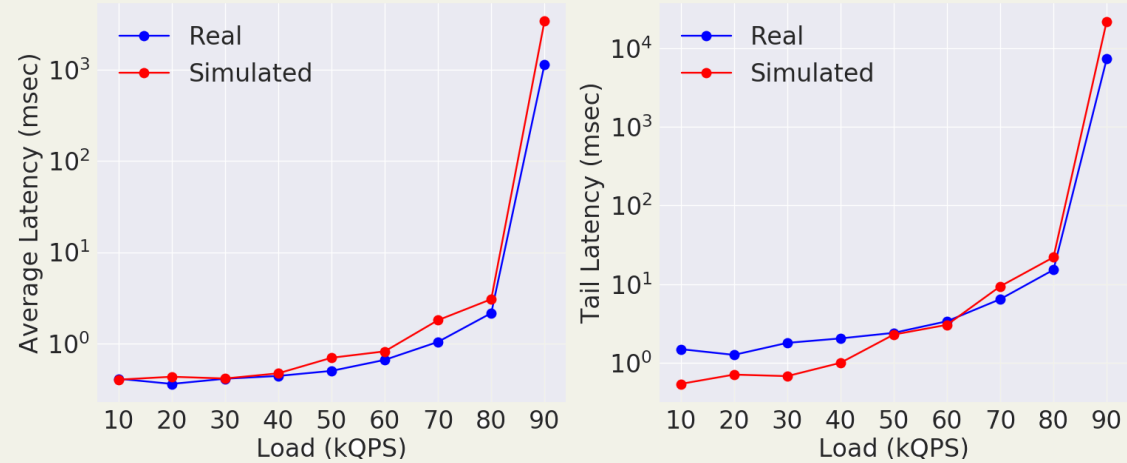
2-tier application



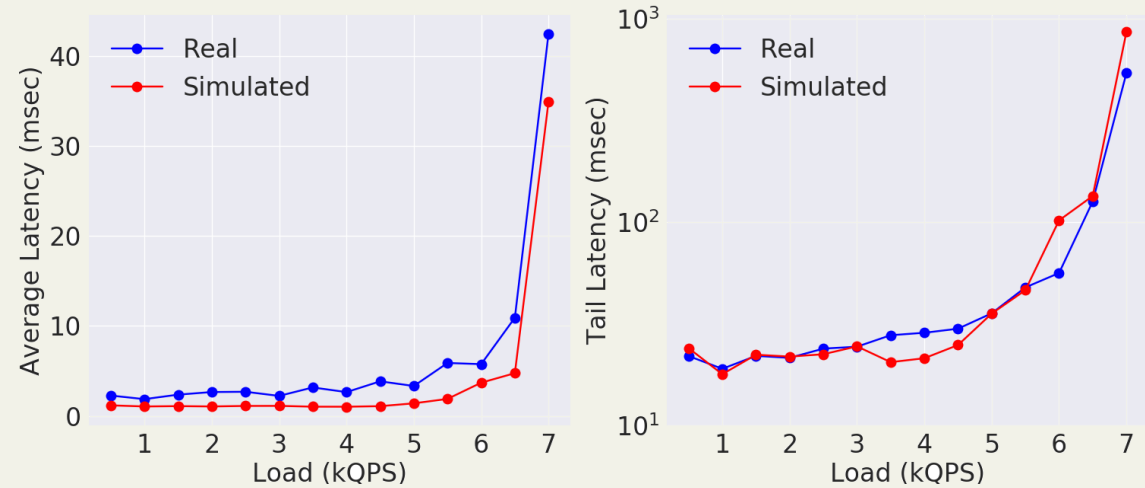
3-tier application

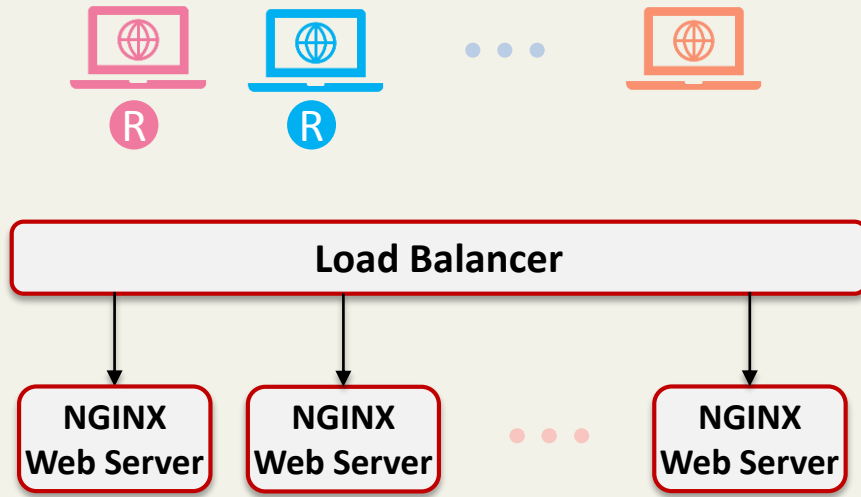


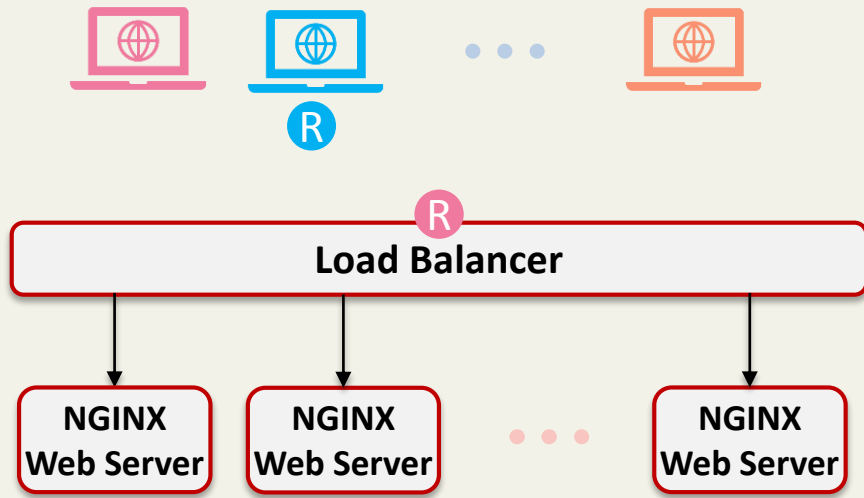
Load-Latency Curve for 2-Tier Application (NGINX_4_MEMC_2)

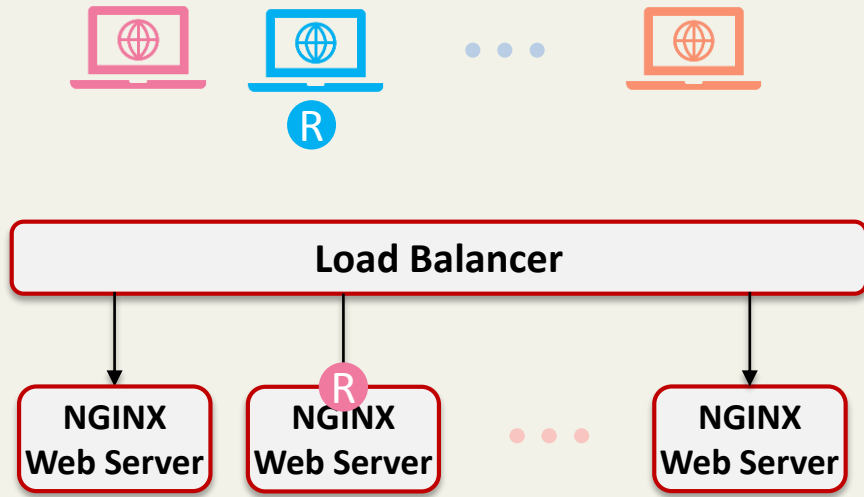


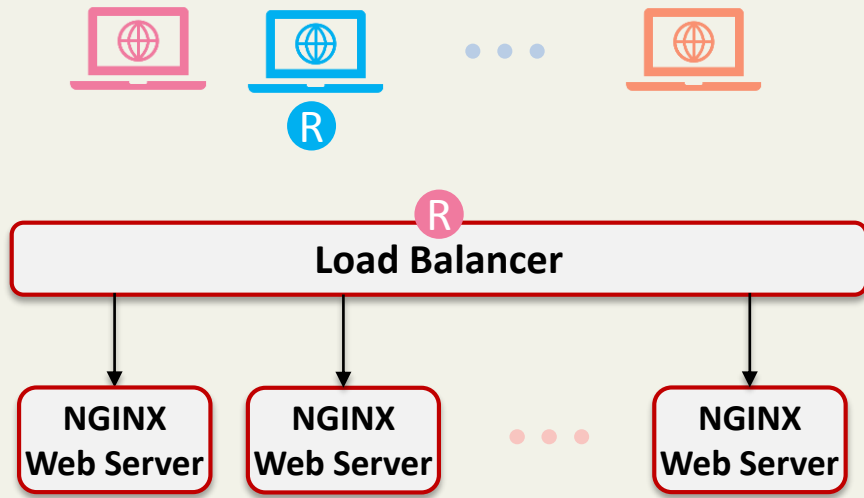
Load-Latency Curve for 3-Tier Application

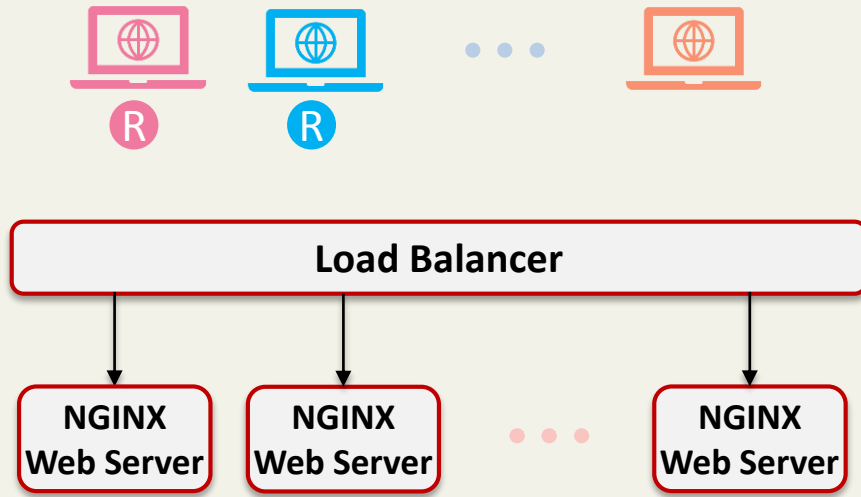


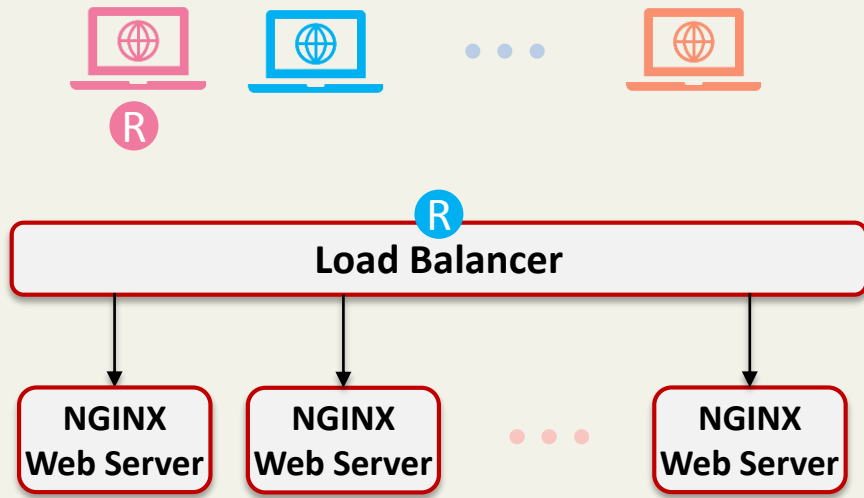


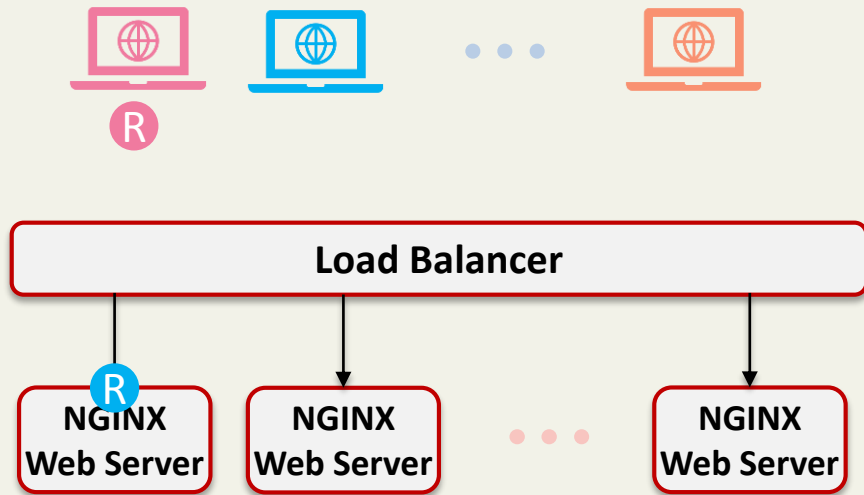


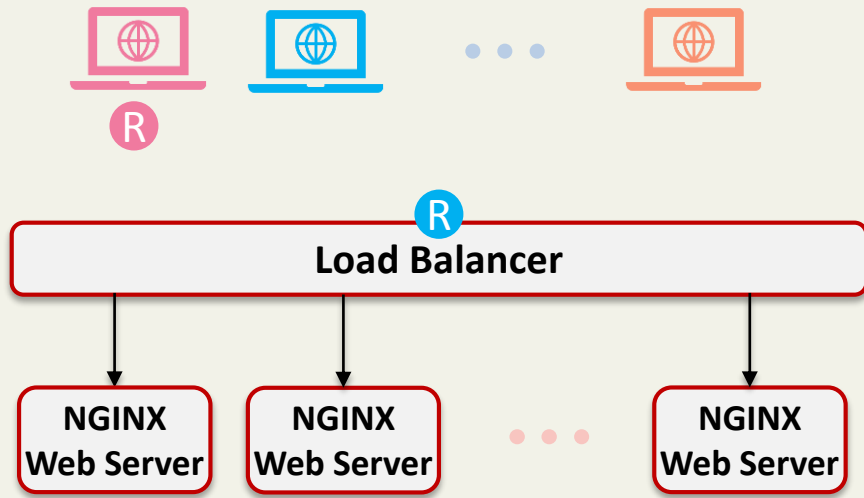


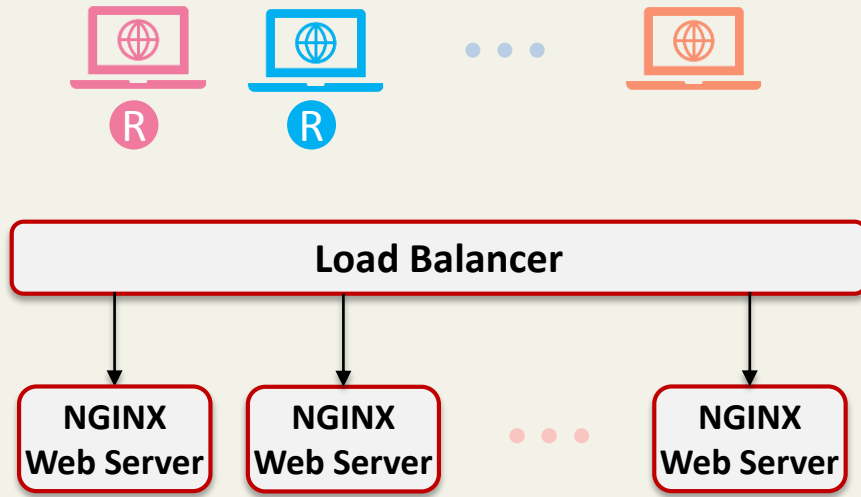




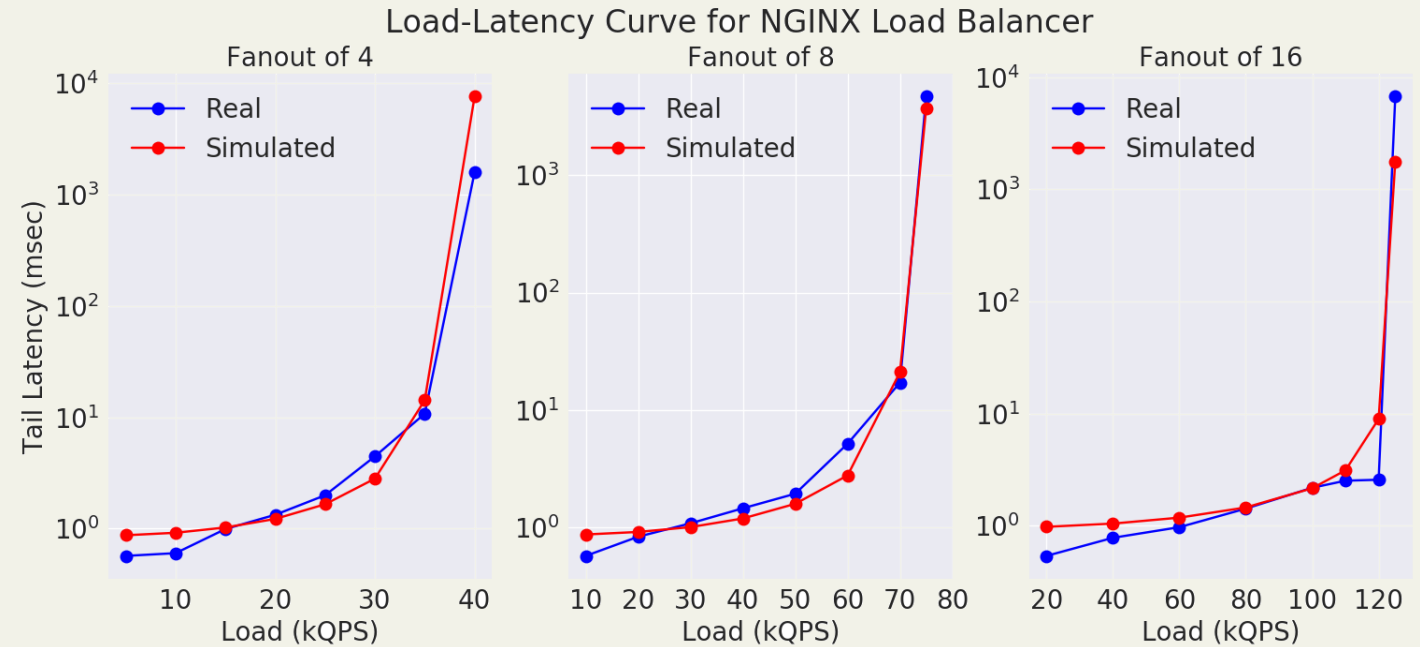
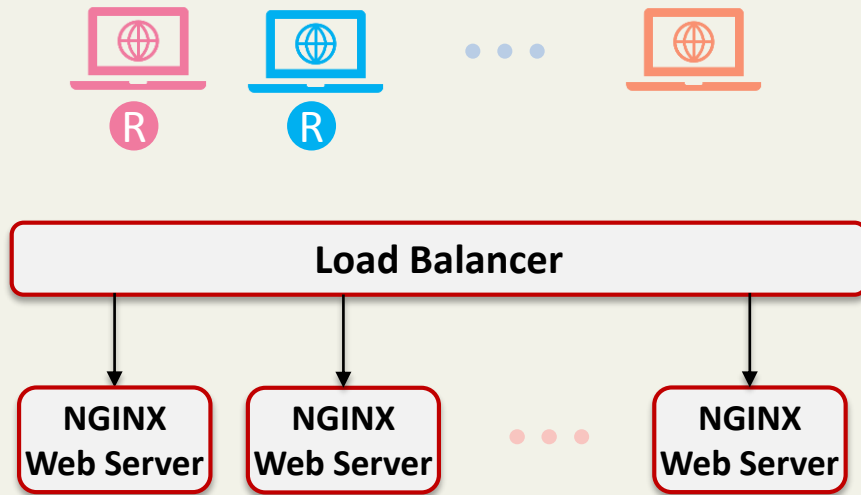


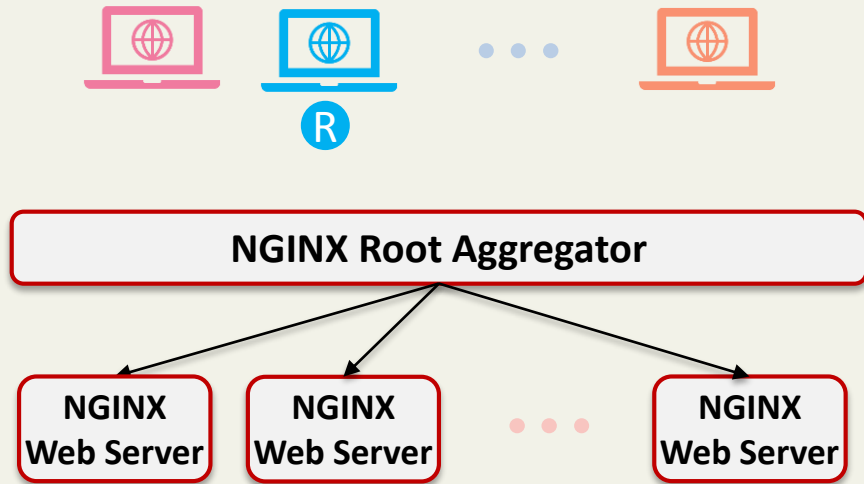


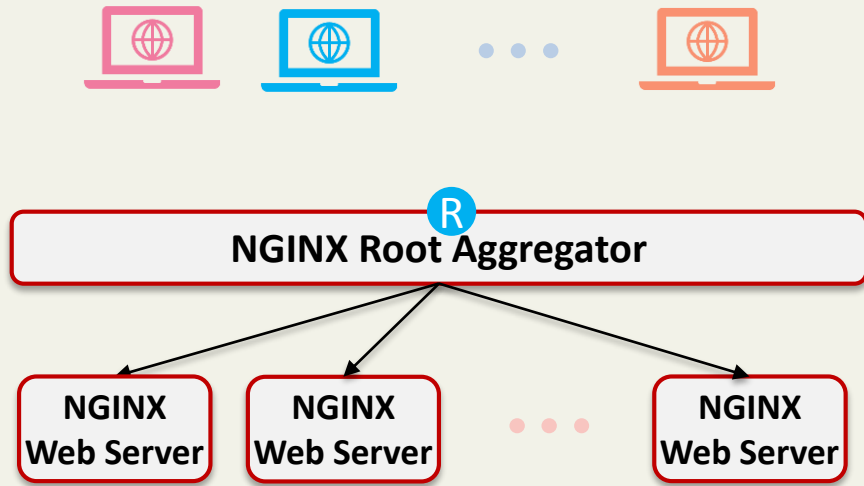


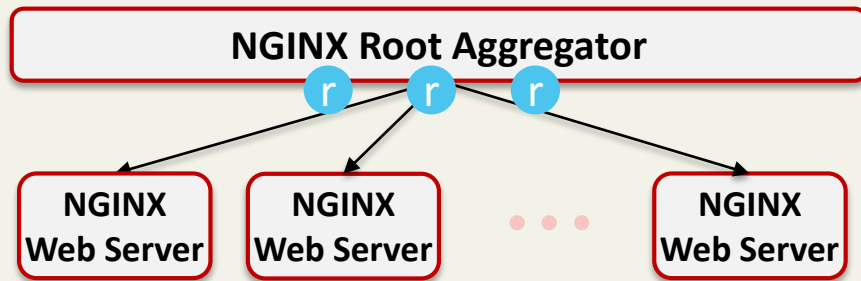


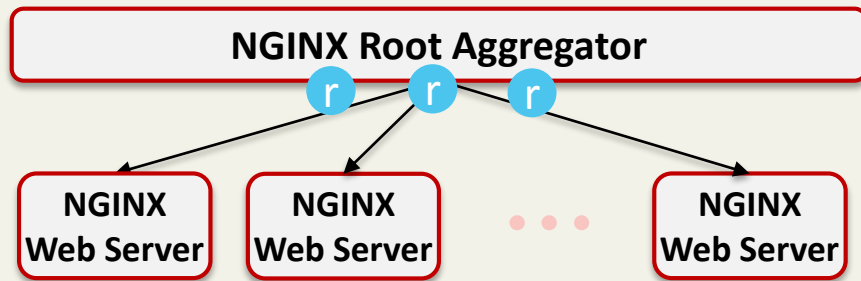
- Linear throughput increase from cluster size of 4 to 8
- Sub-linear increase from cluster size of 8 to 16

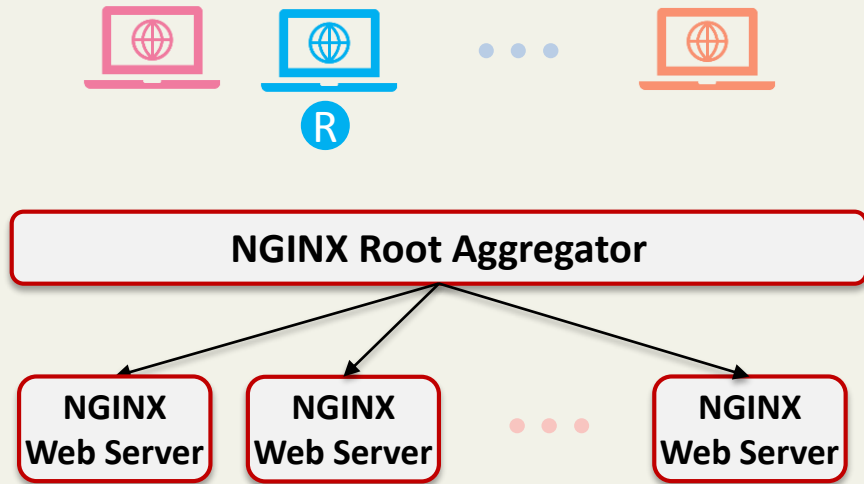


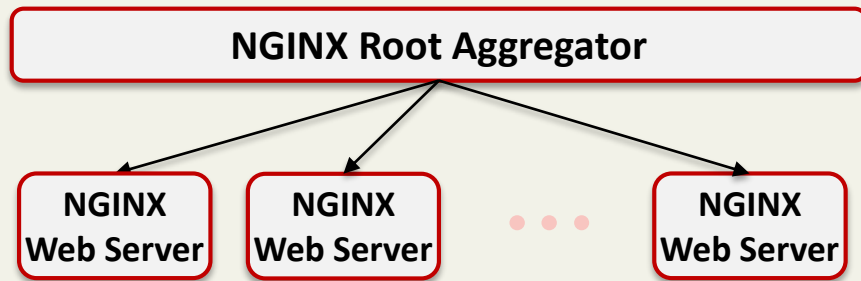




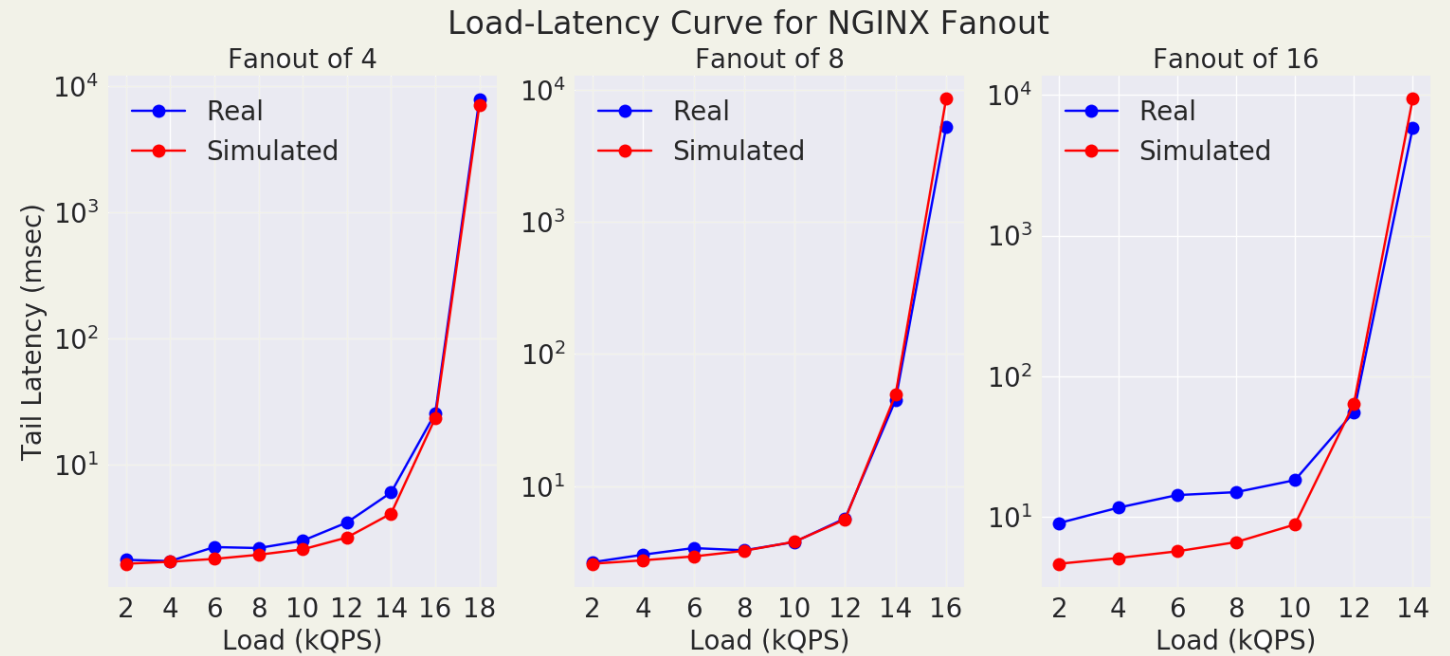
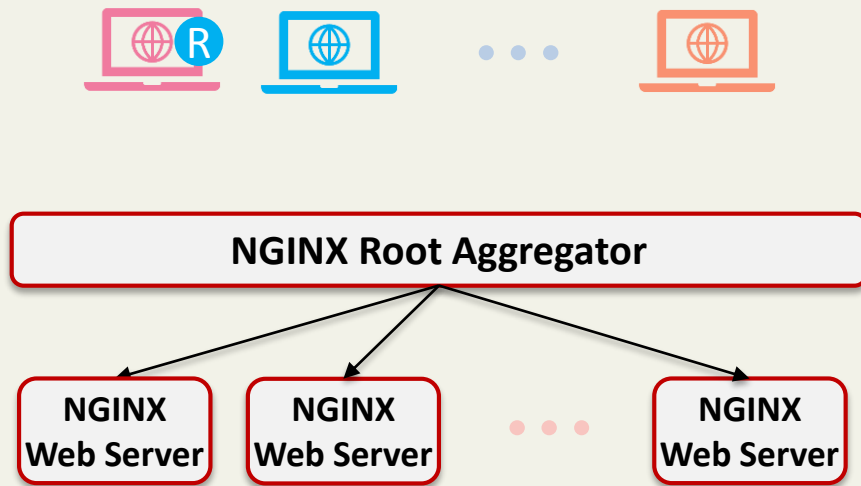




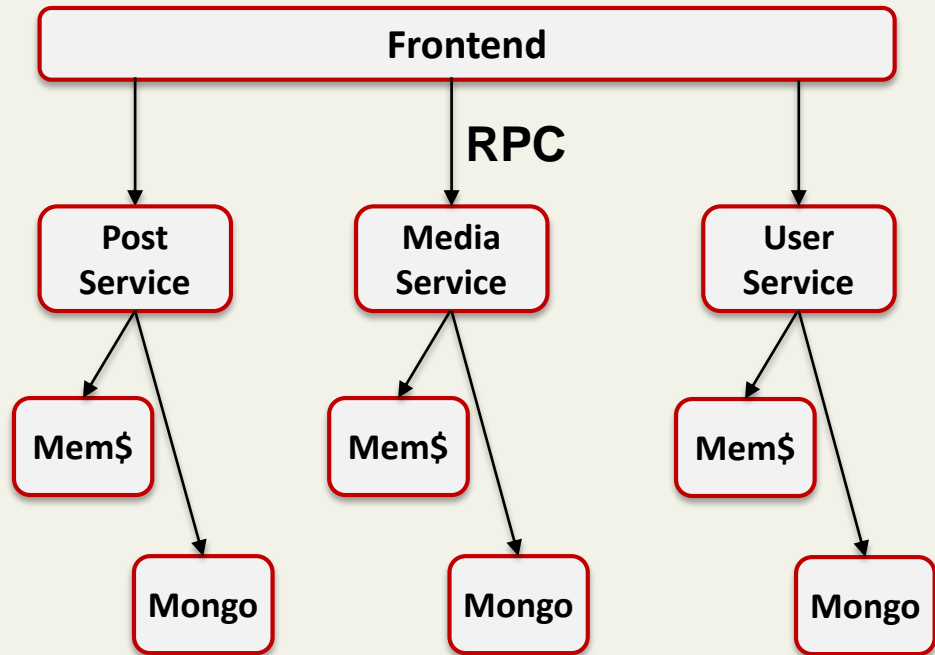




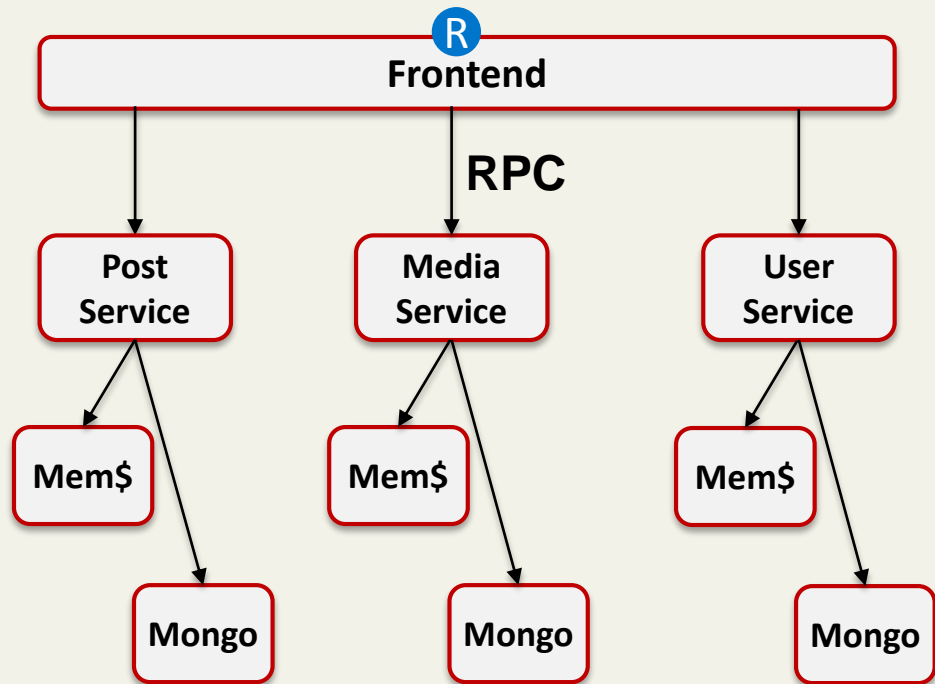
- Throughput decreases with request fanout



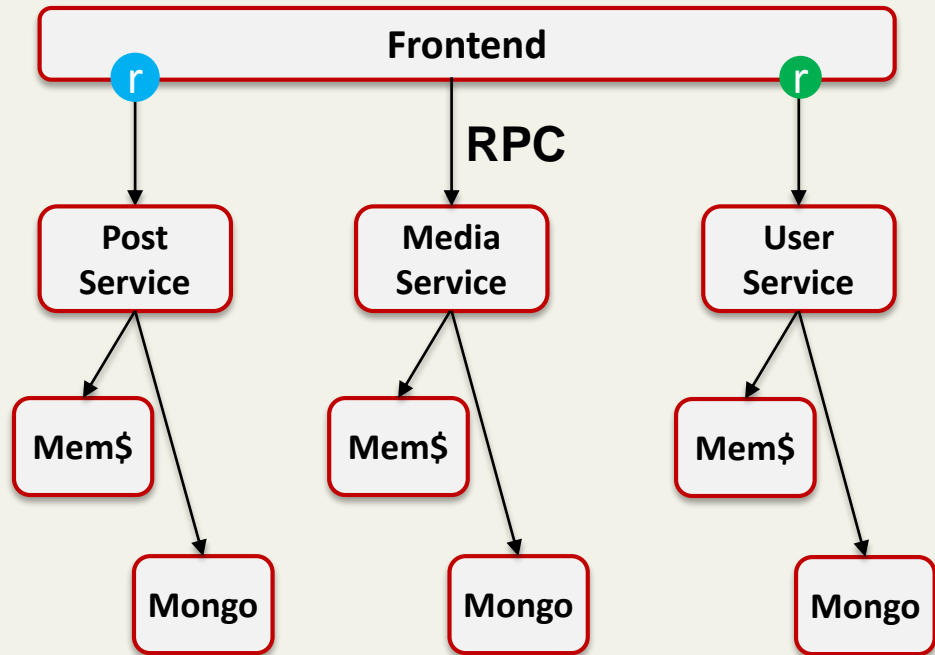
▪ Simplified social network



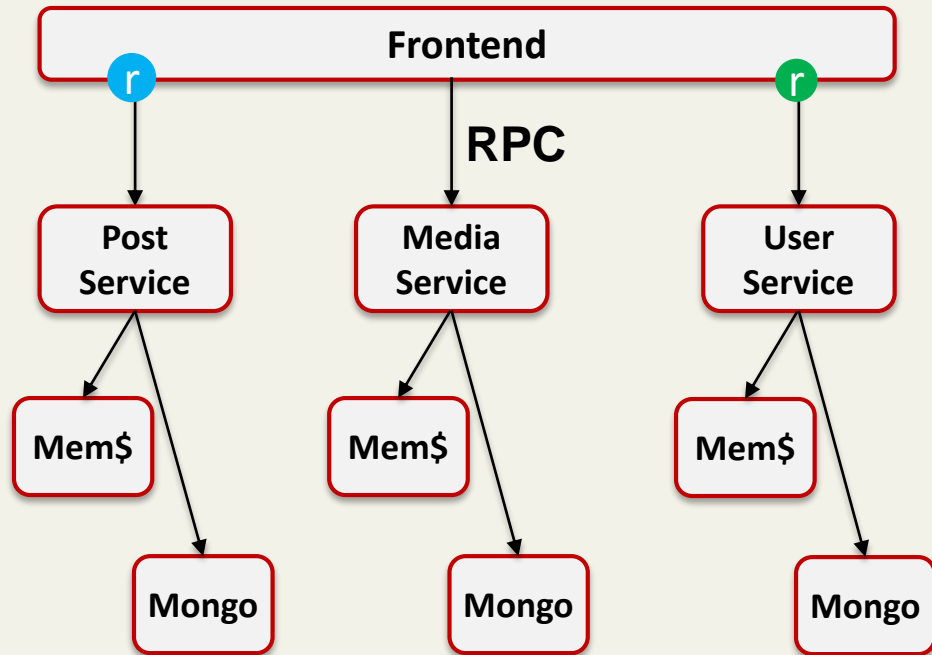
▪ Simplified social network



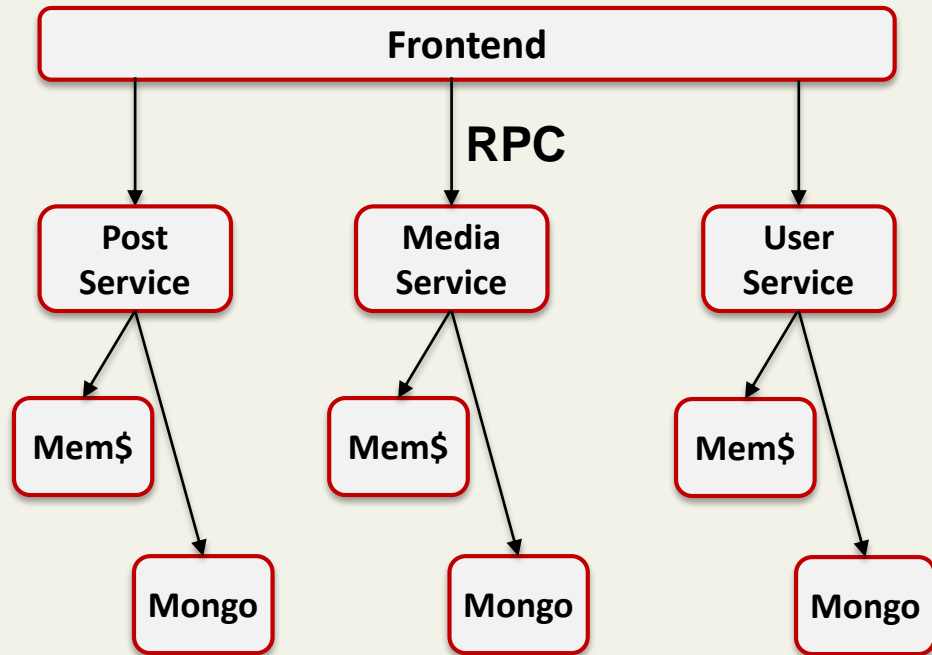
▪ Simplified social network



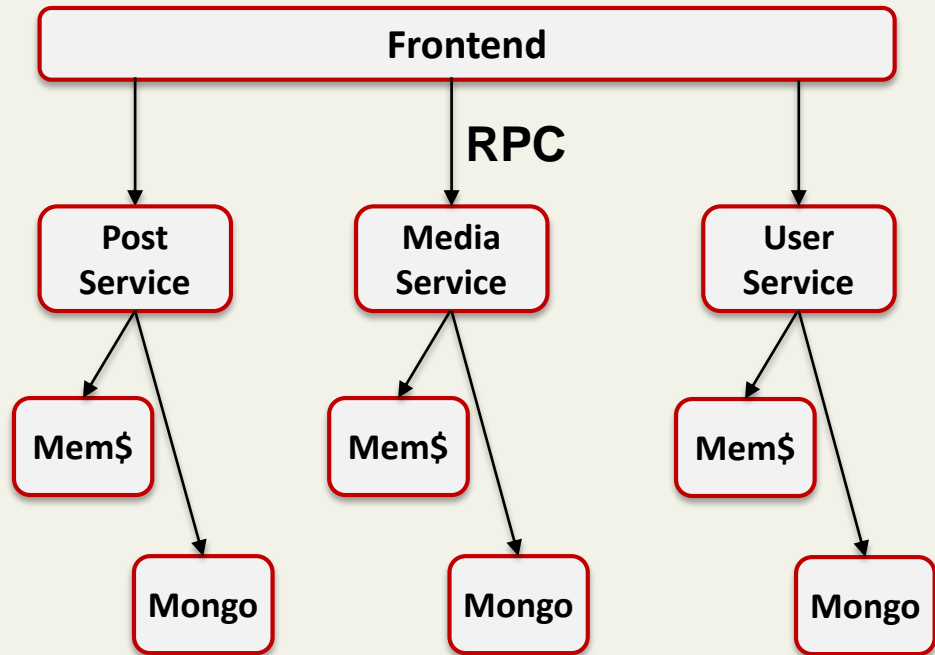
▪ Simplified social network



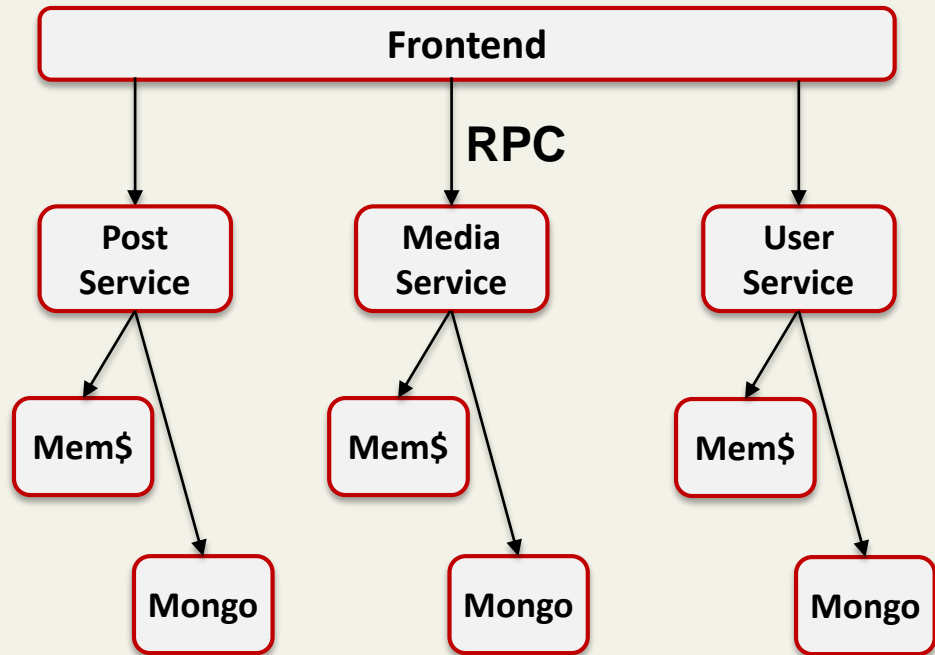
▪ Simplified social network



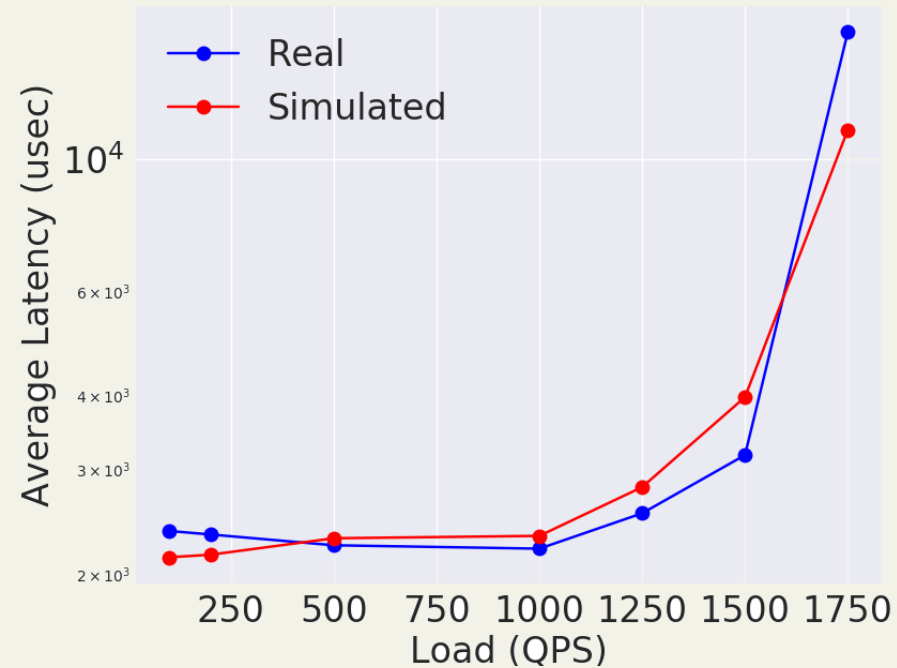
▪ Simplified social network



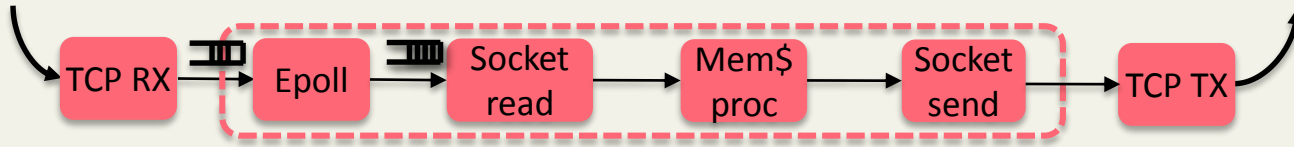
▪ Simplified social network



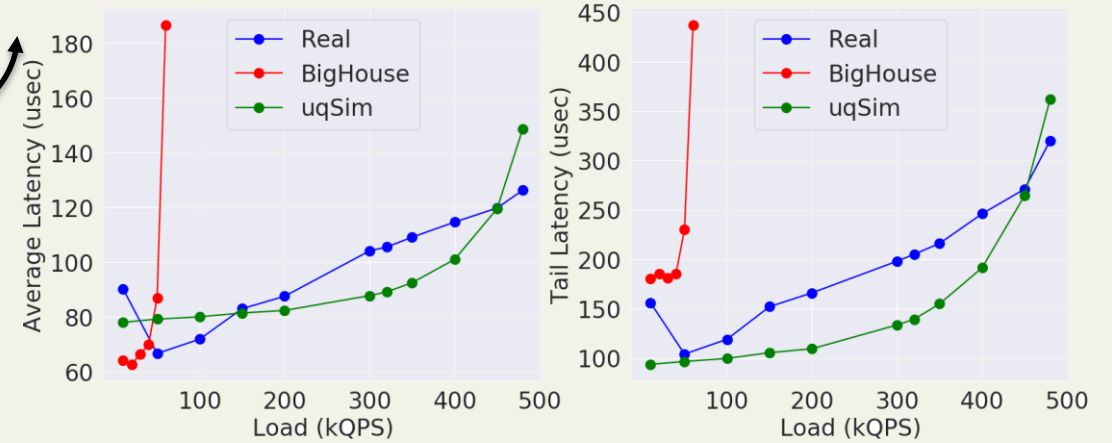
Load-Latency Curve For Social Network Application



Memcached



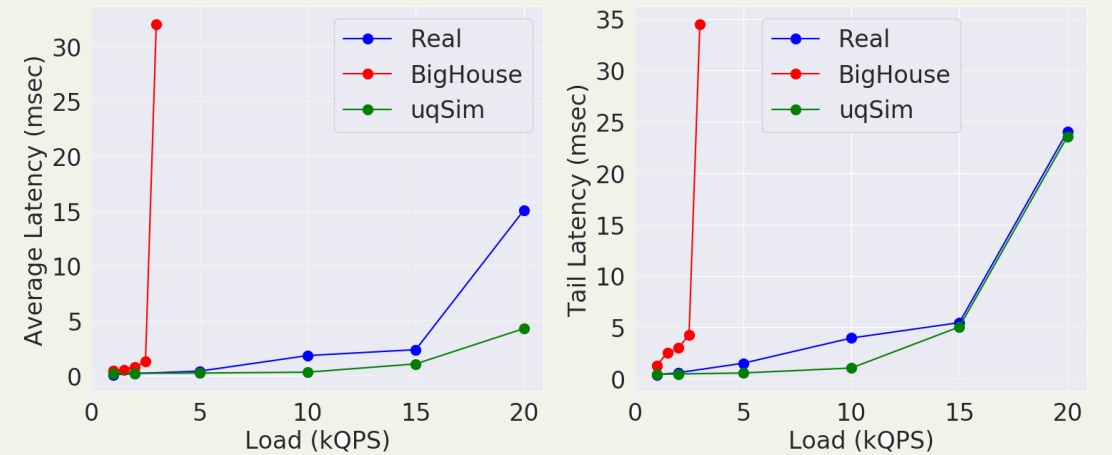
Load-Latency Curve for Memcached (4-Thread)



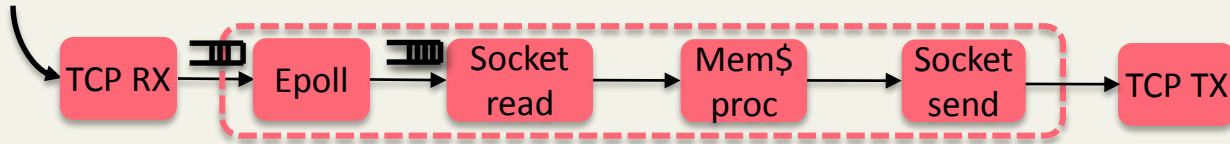
NGINX



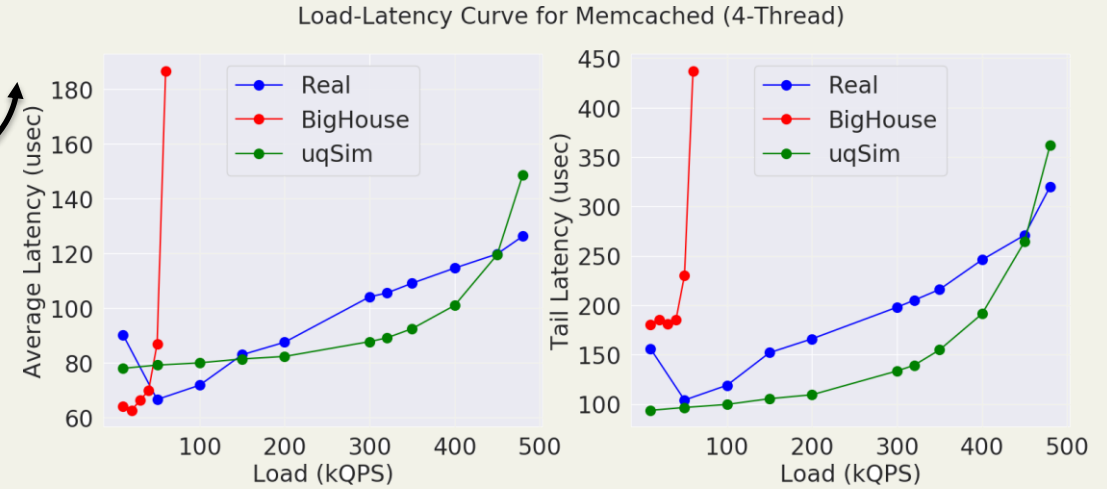
Load-Latency Curve for Nginx (1-Process)



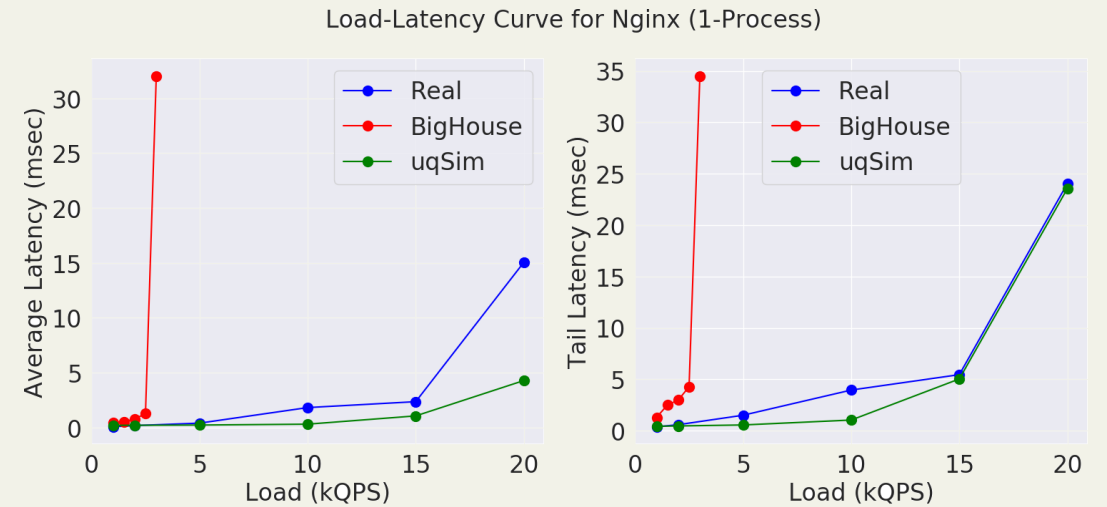
Memcached



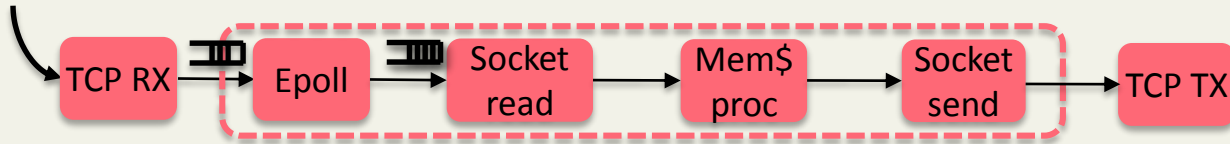
- BigHouse saturates at 50kQPS
- μ qSim & real server saturates at 450kQPS



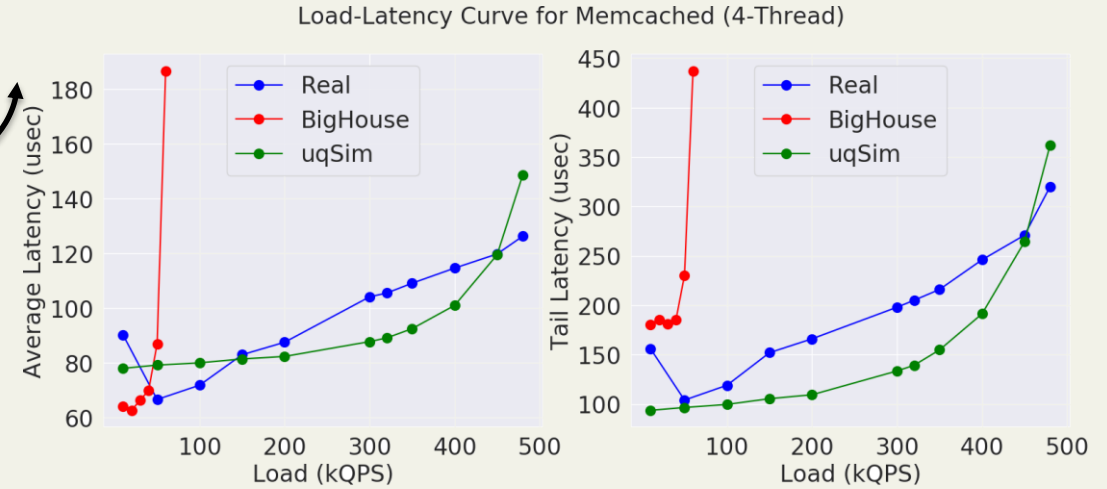
NGINX



Memcached



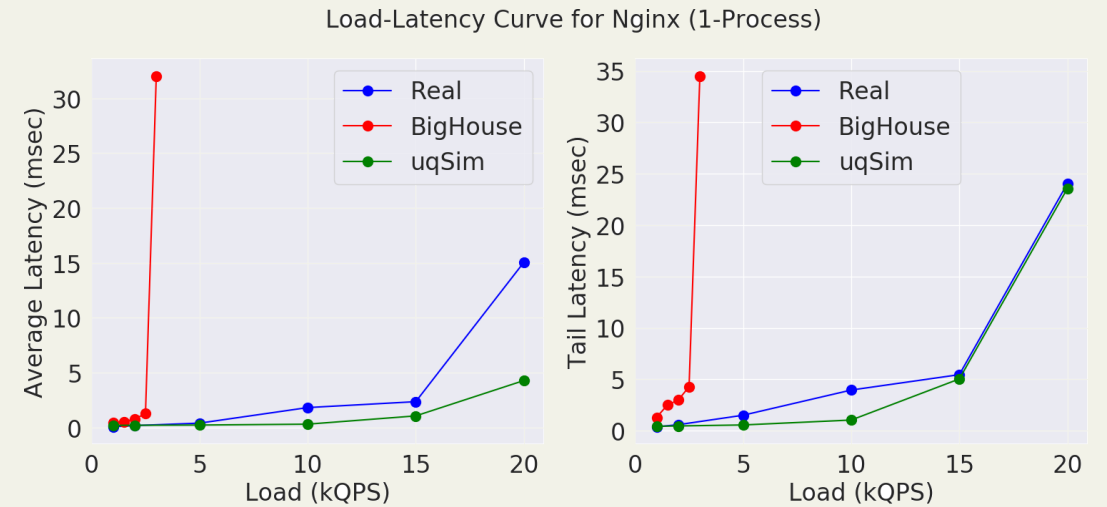
- BigHouse saturates at 50kQPS
- μ qSim & real server saturates at 450kQPS



NGINX



- BigHouse saturates at 2.5kQPS
- μ qSim & real server saturates at 15kQPS



▪ Simulation speed relevant factors

- Simulated input load
- Processing time distribution
- Microservice model complexity & network complexity



■ Simulation speed relevant factors

- Simulated input load
- Processing time distribution
- Microservice model complexity & network complexity

■ Performance compared to real system

- Tail@Scale experiment: 26.5x slow down for cluster of 500 servers
- 4-thread Memcached at 50kQPS: 2.5x speed up
- 1-process Nginx at 10kQPS: 4x speed up

■ Simulation speed relevant factors

- Simulated input load
- Processing time distribution
- Microservice model complexity & network complexity

■ Performance compared to real system

- Tail@Scale experiment: 26.5x slow down for cluster of 500 servers
- 4-thread Memcached at 50kQPS: 2.5x speed up
- 1-process Nginx at 10kQPS: 4x speed up

■ Future work

- Parallelizing simulation
- Each thread simulates a partition of the network

- **Microservices introduce new system challenges**
 - Need scalable simulation techniques to study large scale effects
- **μ qSim: validated microservice simulator**
 - Modeling the internal queueing structure of individual microservices
 - Modeling dataflow behavior across the microservices
 - Validated against simple & complex microservices and accurately capturing throughput/latency
 - Planning to open source @ microservices.ece.cornell.edu

- **Microservices introduce new system challenges**
 - Need scalable simulation techniques to study large scale effects
- **μ qSim: validated microservice simulator**
 - Modeling the internal queueing structure of individual microservices
 - Modeling dataflow behavior across the microservices
 - Validated against simple & complex microservices and accurately capturing throughput/latency
 - Planning to open source @ microservices.ece.cornell.edu

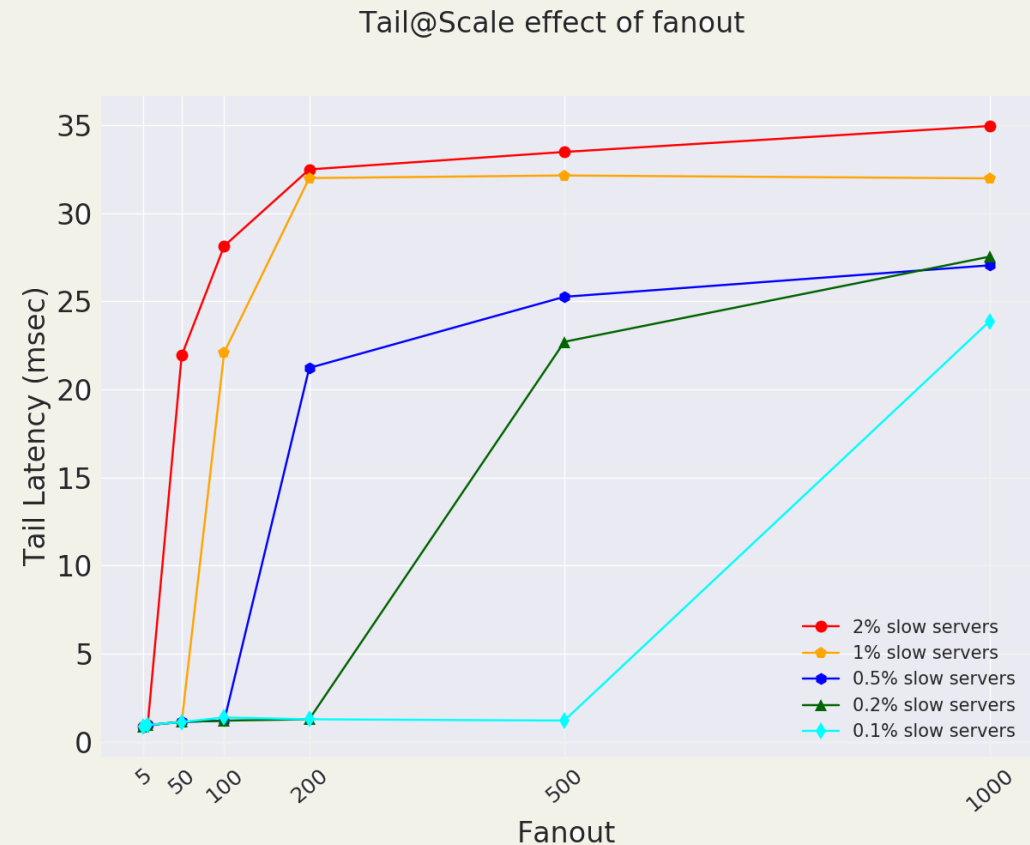
μqSim simulated real time	5s	20s	60s	
μqSim simulation time	12.8s	51.2s	154.1s	
BigHouse warmup samples	5x50k	20x50k	60x50k	5000
BigHouse simulation time	8.2s	28.7s	154s	1.2s

4-thread Memcached at 50k

μqSim simulated real time	5s	20s	60s	
μqSim simulation time	0.7s	2.7s	8.1s	
BigHouse warmup samples	5x2k	20x2k	60x2k	5000
BigHouse simulation time	9s	23s	64s	7.1s

1-process NGINX at 2k

- For simple single-tier single stage
- Different Cluster Sizes
- Different Fraction of Slow Servers

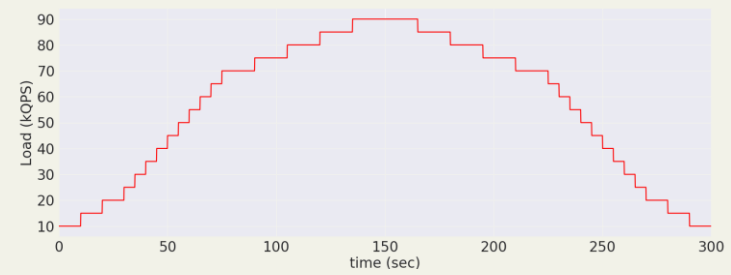


■ End-to-End QoS Target

- 5ms

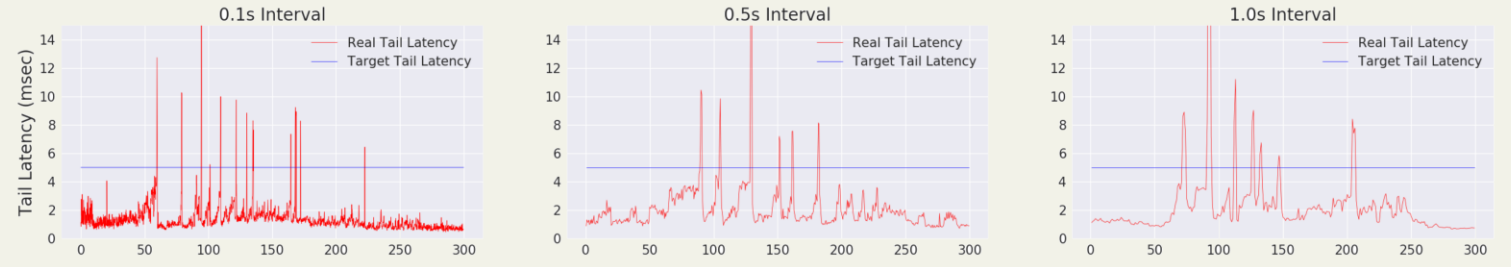
■ Input Load

- Diurnal pattern



■ Results

Real Power Management For 2 Tier Application (NGINX_8_MEMC_2)



Simulated Power Management For 2 Tier Application (NGINX_8_MEMC_2)

