# WeChat

## The new way to connect

Chat     Moments     Contacts     Search     Pay

**1 Billion**
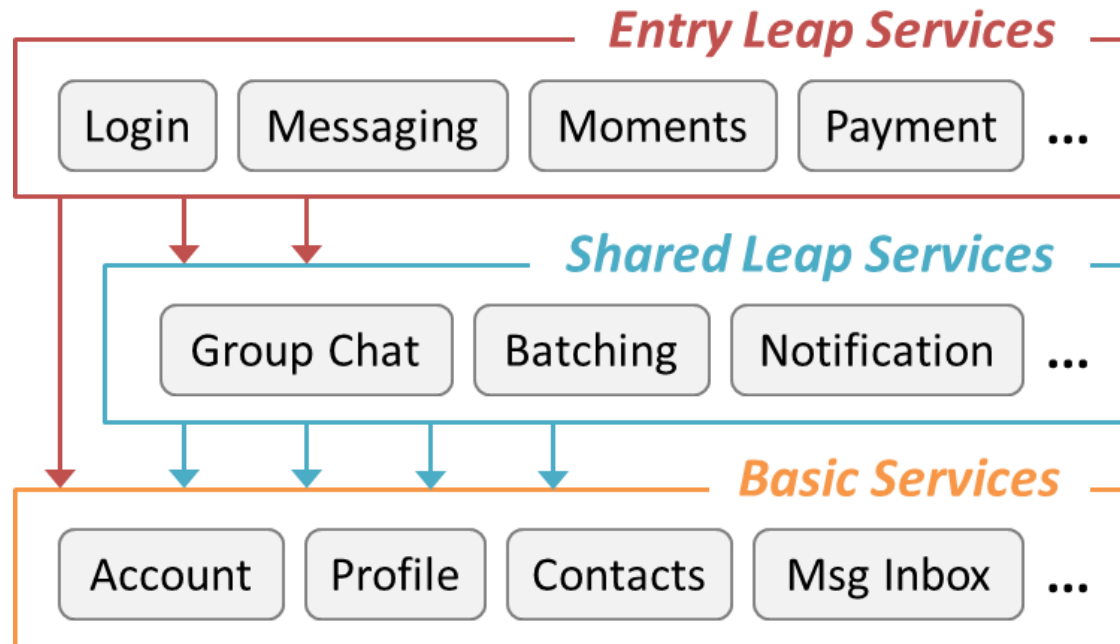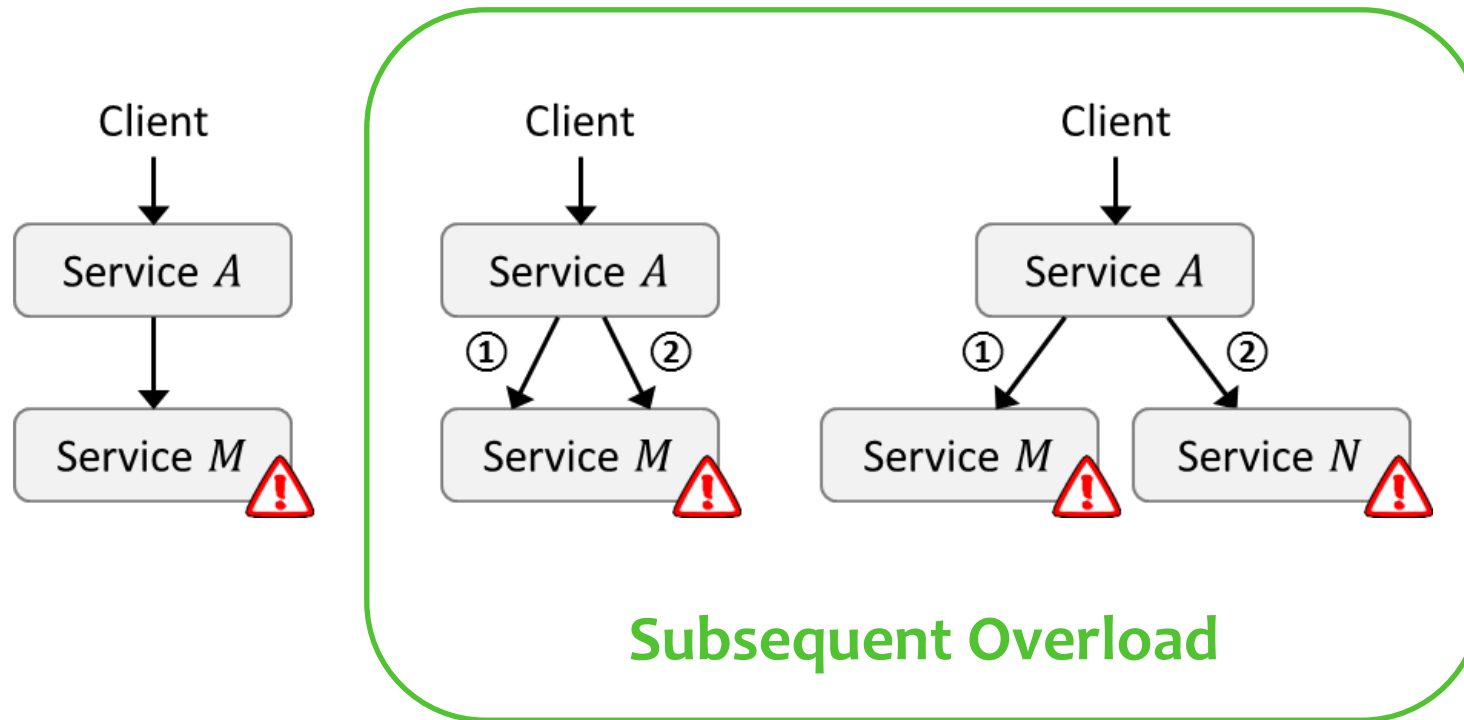
monthly active users

- **Service DAG**
  - Vertex: a distinct service; Edge: call path
  - **Basic service**: out-degree = 0
  - **Leap service**: out-degree ≠ 0
    - **Entry service**: in-degree = 0

**Entry Leap Services**

| Login | Messaging | Moments | Payment | ... |

**Shared Leap Services**

| Group Chat | Batching | Notification | ... |

**Basic Services**

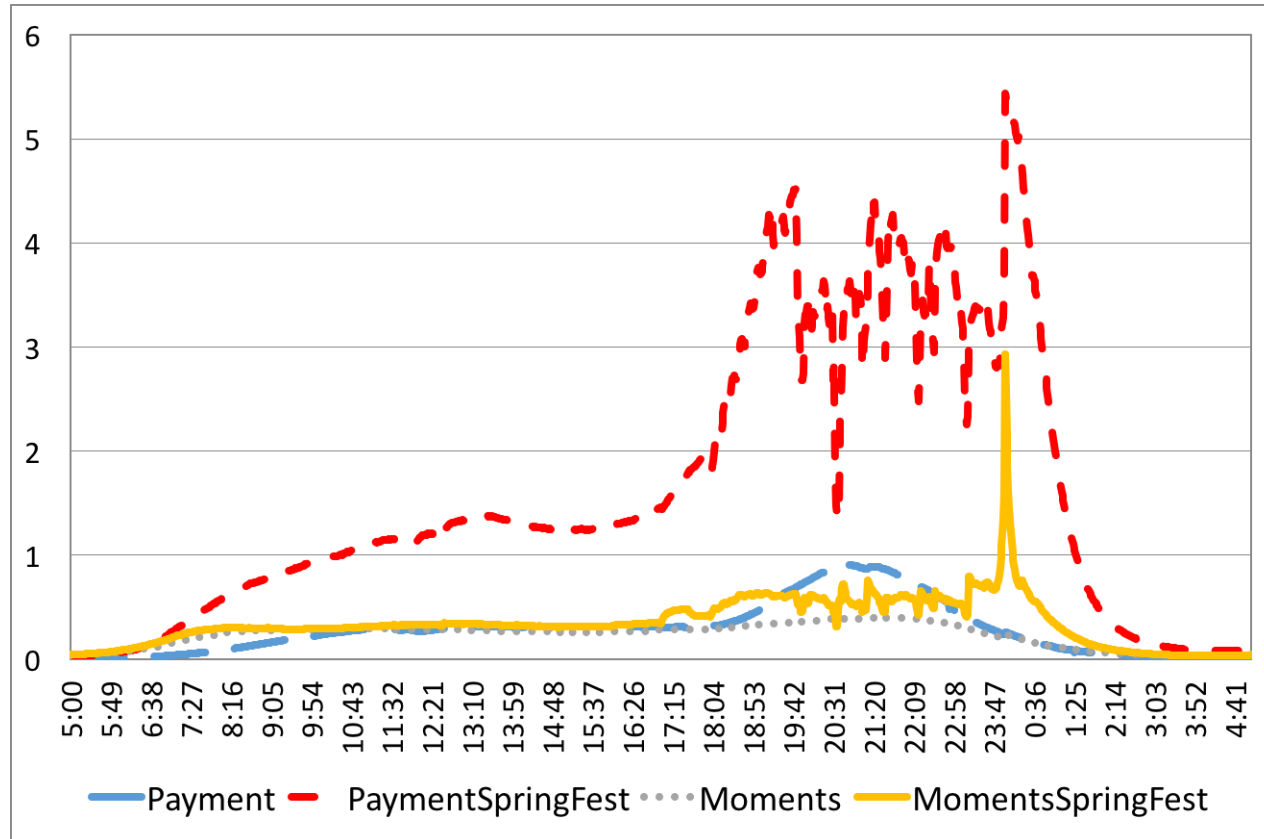| Account | Profile | Contacts | Msg Inbox | ... |

- **It's usually hard to estimate the dynamics of workload during the development of microservices.**



How about random load shedding?

# Dynamic Workload



**Relative Statistics of WeChat Service Requests**

- **Overload detection**

- **Service admission control**

- **Requirements**
  - Service agnostic
    - o Benefit the ever evolving microservice system
    - o Decouple overload control from the business logic of services
  - Independent but collaborative
    - o Decentralized overload control
    - o Service-oriented collaboration among nodes
  - Efficient and fair
    - o Sustain best-effort success rate of service when load shedding becomes inevitable
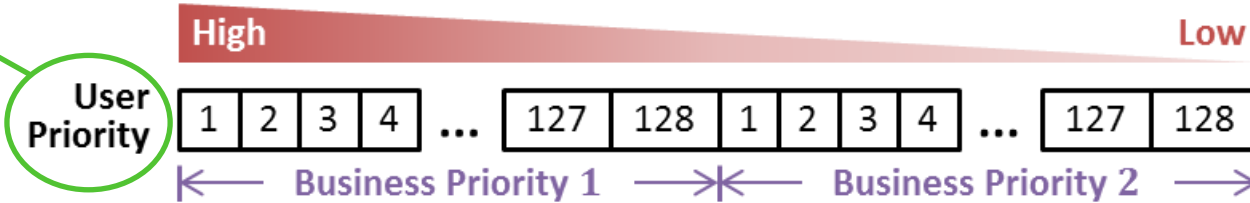    - o Bias-free overload control

- **Load indicator of a node: Queuing time**
    - Rationale: to manage queue length for SLA

- **Why not response time?**

- **Why not CPU utilization?**

# Service Admission Control
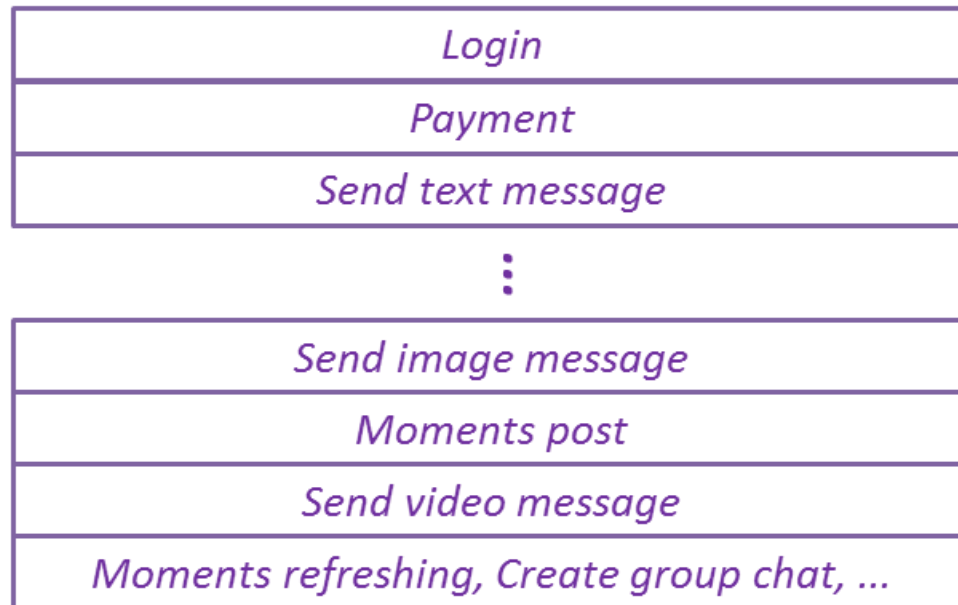


Shuffling on an hourly basis
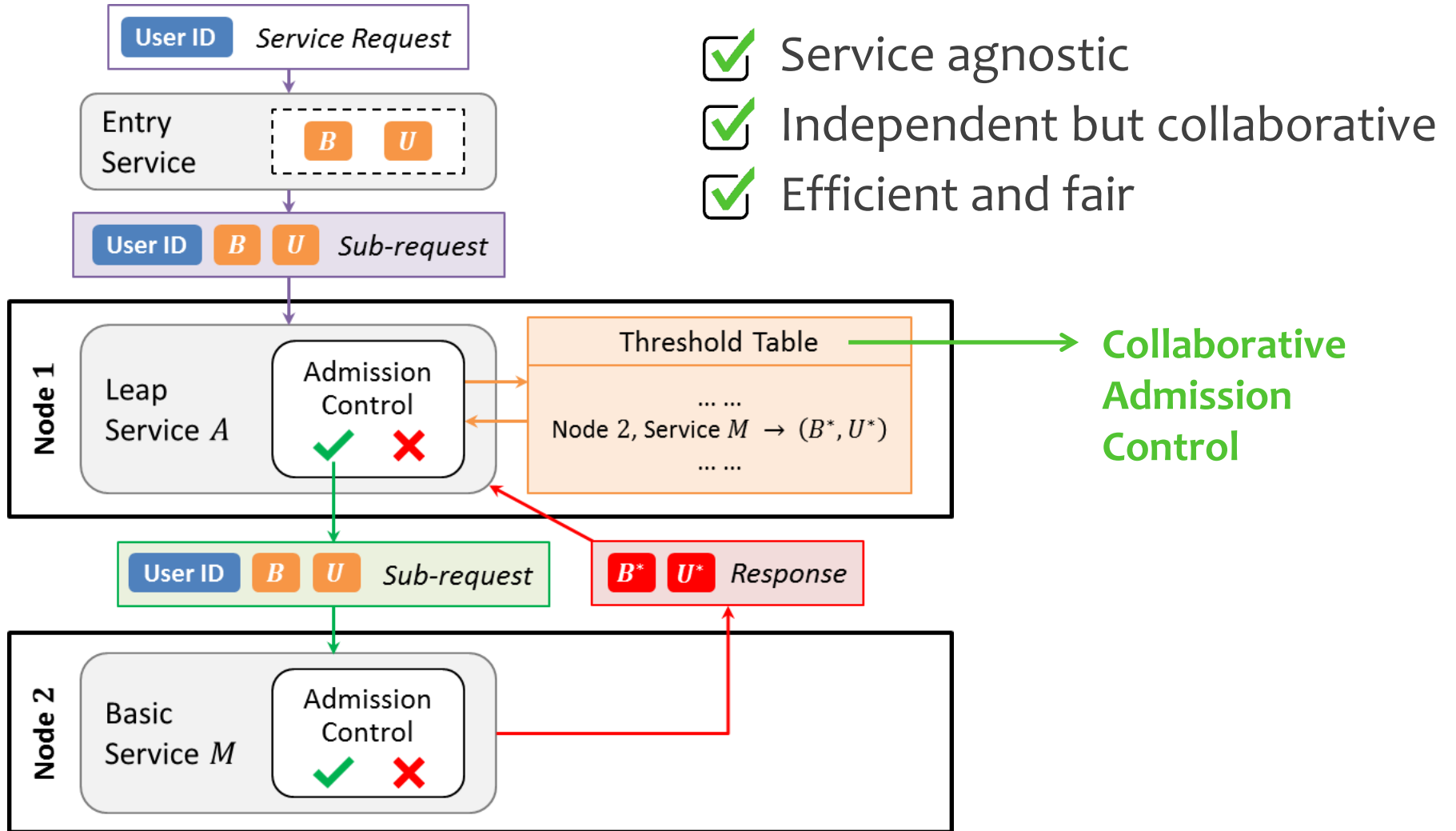
Exploit histogram for real-time adjustment

User Priority

High — Low

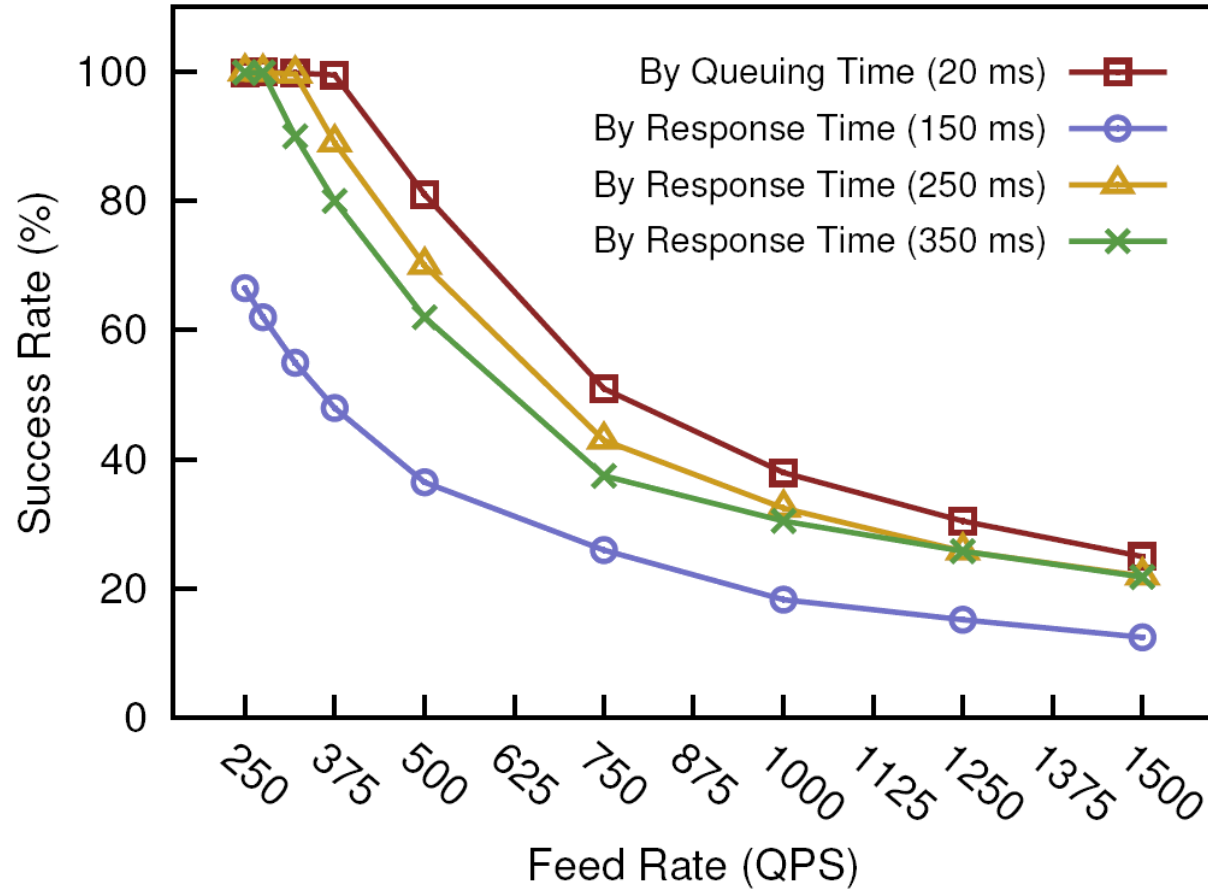| 1 | 2 | 3 | 4 | ... | 127 | 128 | 1 | 2 | 3 | 4 | ... | 127 | 128 |

← Business Priority 1 → | ← Business Priority 2 →

Static

Business Priority — Priority decreases

Login

Payment

Send text message

⋮

Send image message

Moments post

Send video message

Moments refreshing, Create group chat, ...

# DAGOR Workflow



User ID  *Service Request*

Entry Service  **B** **U**

User ID **B** **U** *Sub-request*

**Node 1**

Leap Service $A$  Admission Control ✓ ✗

Threshold Table

... ...
Node 2, Service $M$ → $(B^*, U^*)$
... ...

User ID **B** **U** *Sub-request*

$B^*$ $U^*$ *Response*

**Node 2**

Basic Service $M$  Admission Control ✓ ✗

☑ Service agnostic
☑ Independent but collaborative
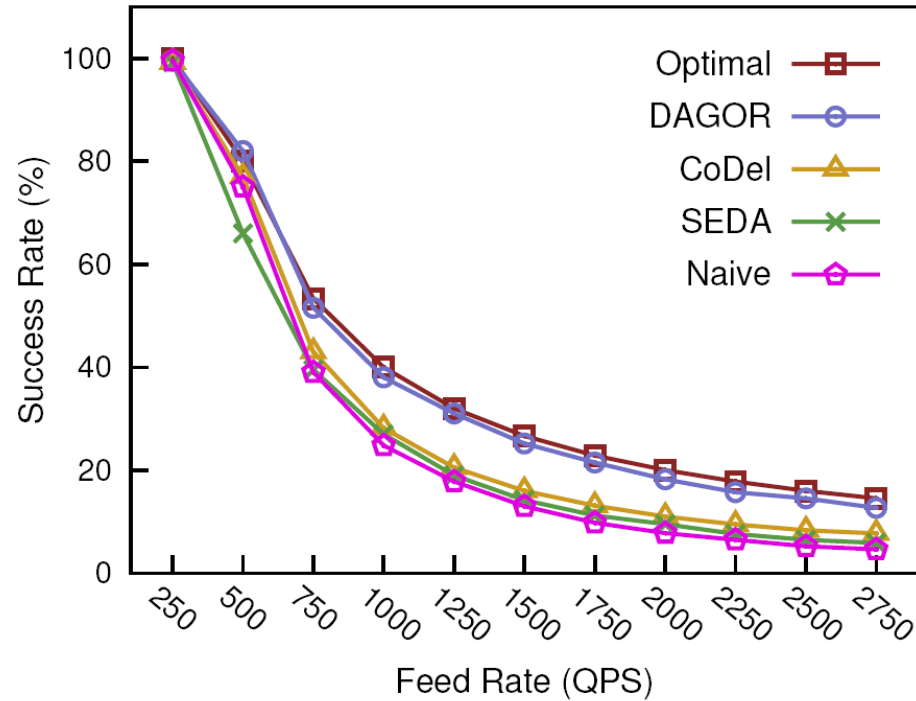☑ Efficient and fair
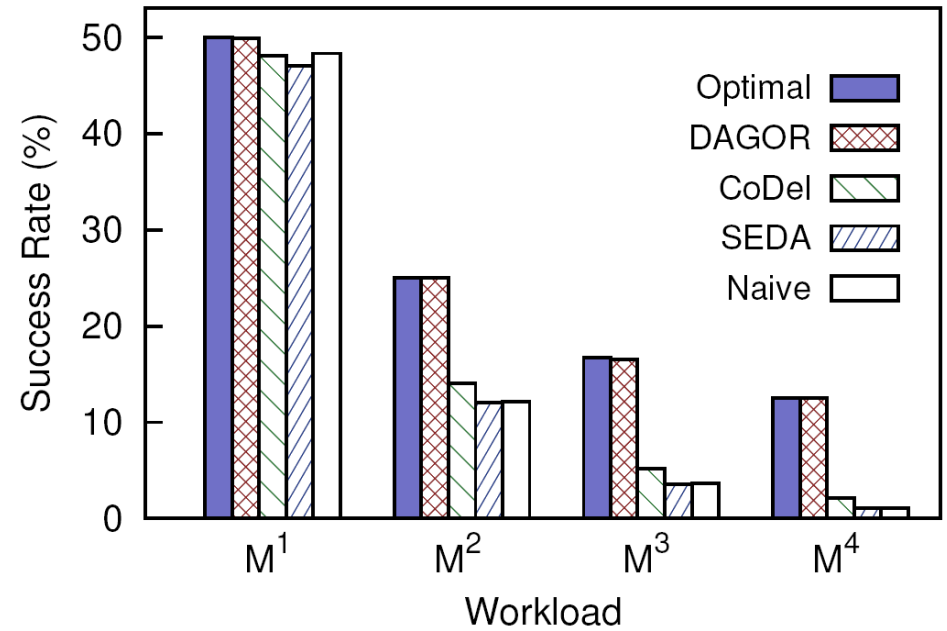
**Collaborative Admission Control**

# Overload Detection



**Queuing Time vs. Response Time**

**Overload Control
with Increasing Workload (M²)**

**Overload Control
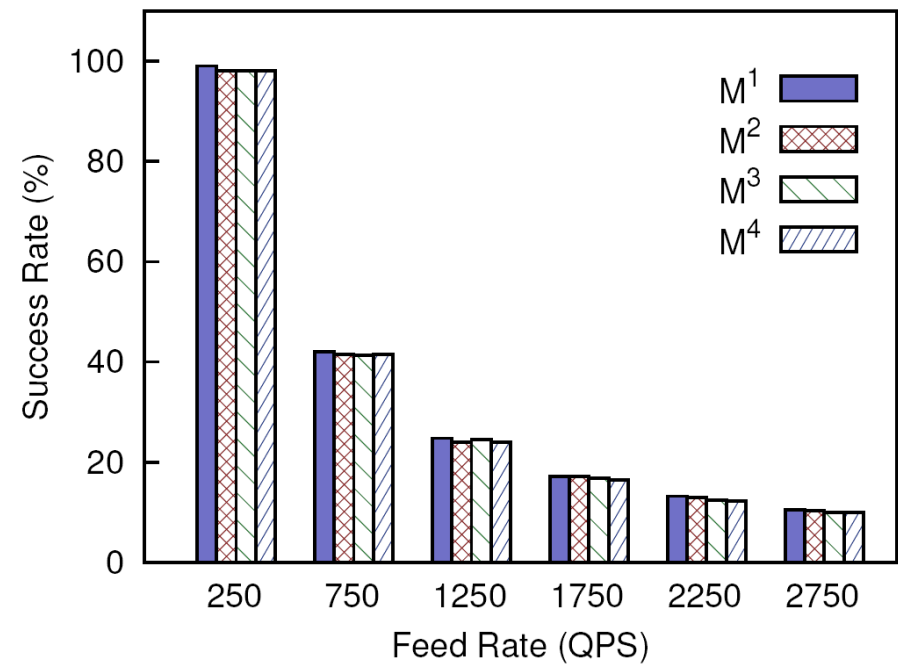with Different Types of Workload**

Optimal Success Rate = $f_{sat}/f$

# Fairness



CoDel

DAGOR

1. **Must be decentralized and autonomous in each service/node**
   - Essential for the overload control framework to scale with the ever evolving microservice system

2. **Employ feedback mechanism for adaptive load shedding**
   - Essential for adjusting thresholds automatically

3. **Prioritize user experience**

Thank You ALL!