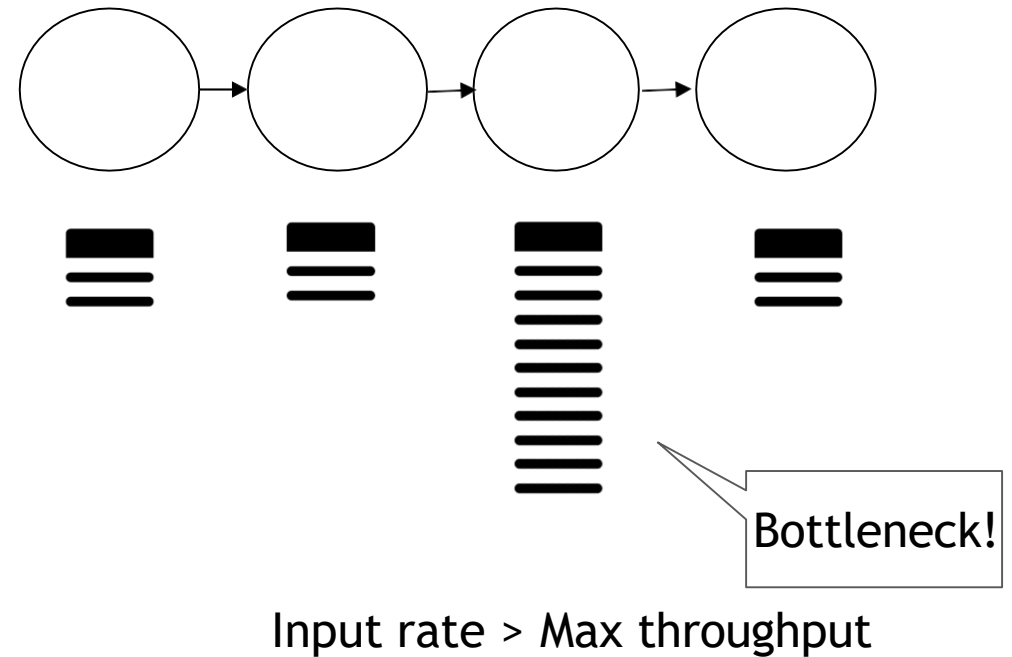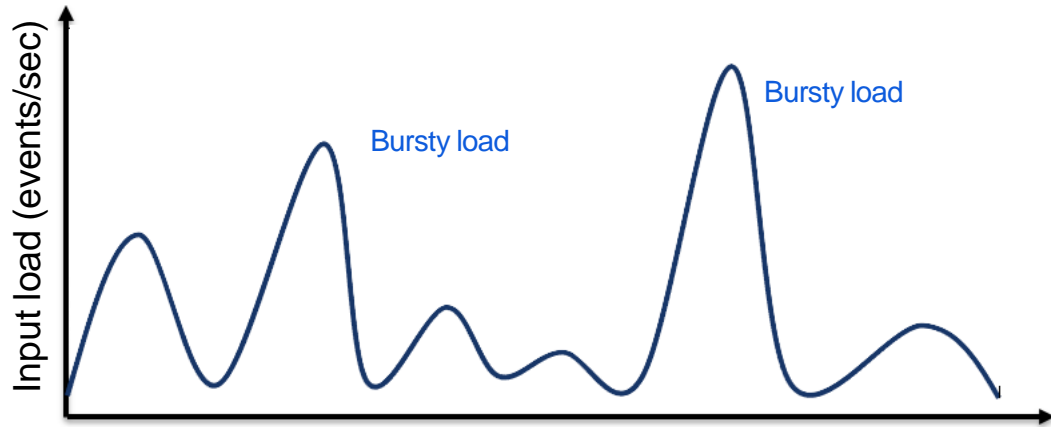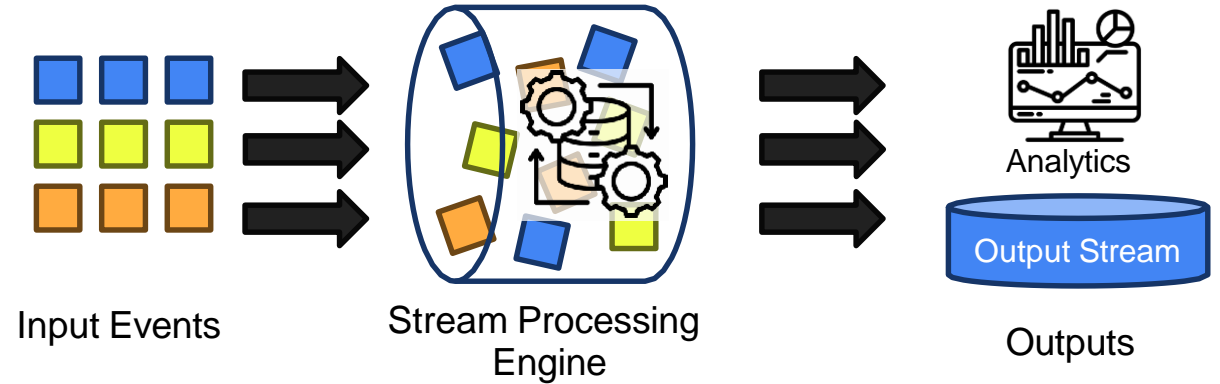# ATC'23_Sponge:
# Fast Reactive Scaling for Stream Processing with Serverless Frameworks

2023/08/23

# Background: Stream Processing

## Stream Processing
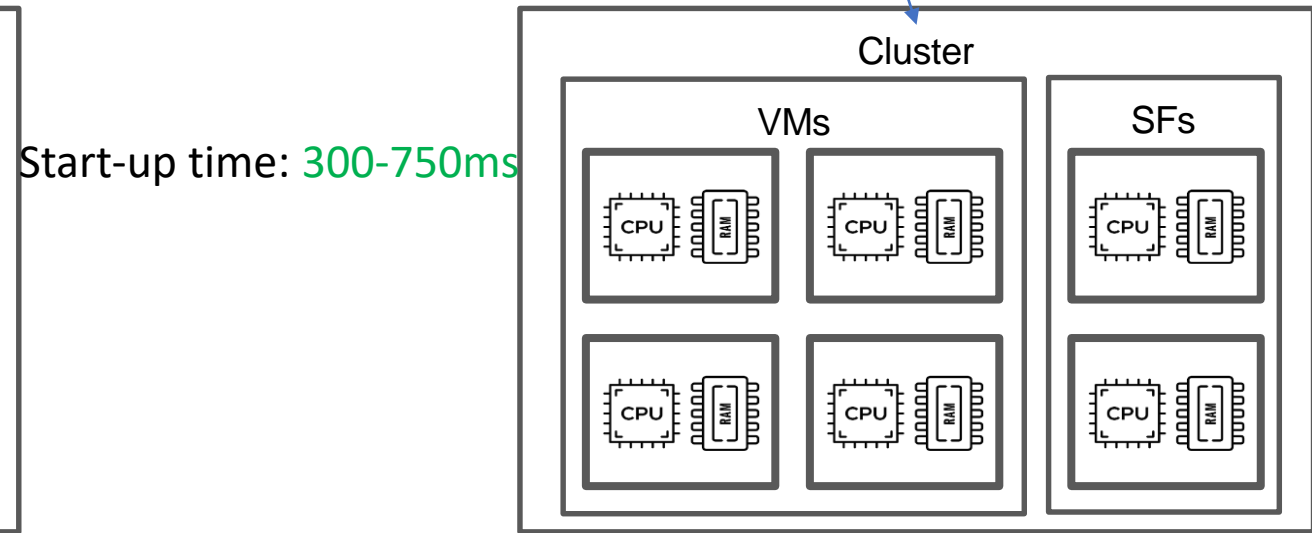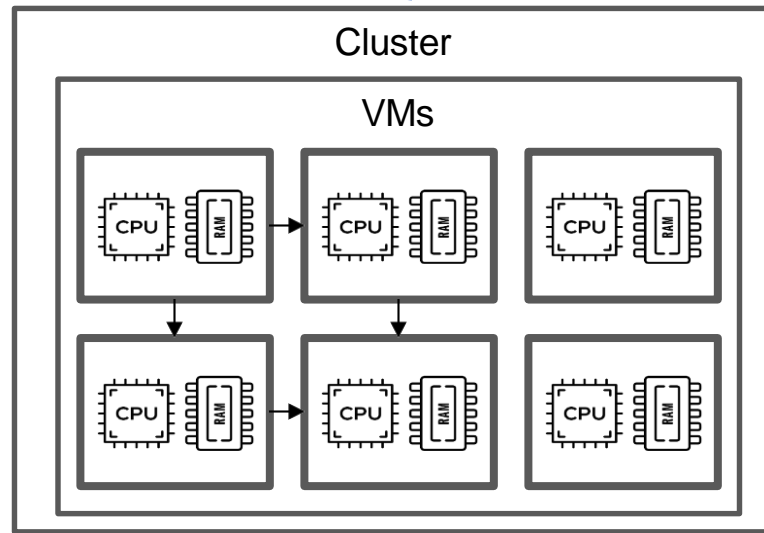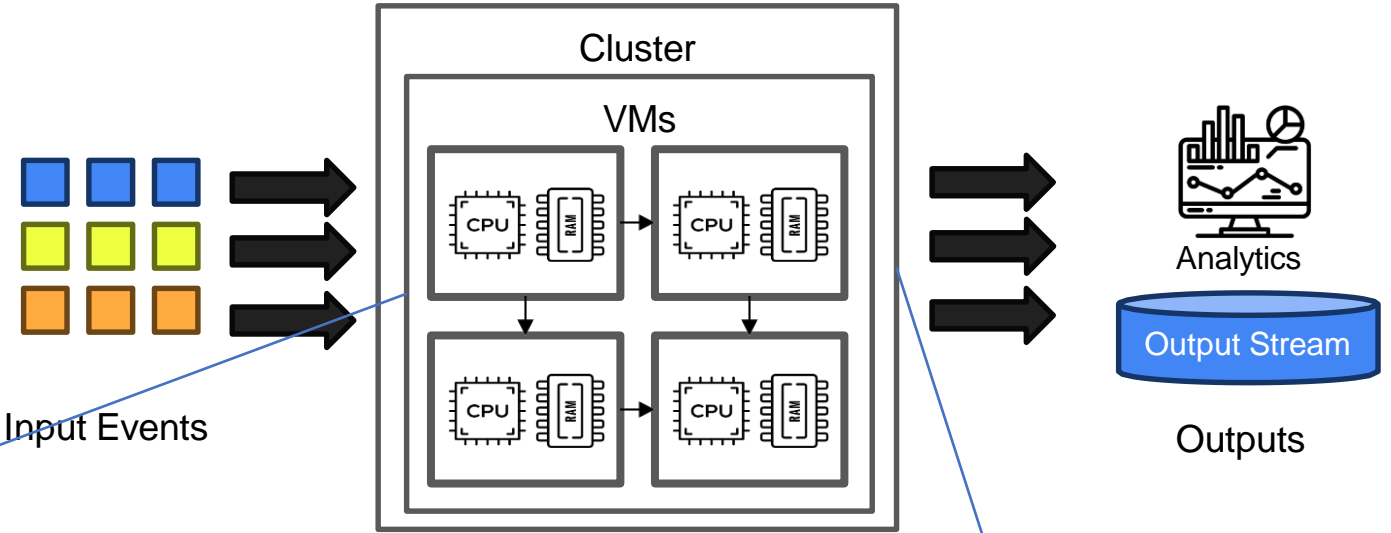
- Query -> DAG
- Bursty Load -> Bottleneck

Input Events

Stream Processing Engine

Analytics

Output Stream

Outputs

Input load (events/sec)

Bursty load

Bursty load

Bottleneck!

Input rate > Max throughput

# Background: On-Demand Resource Provisioning

## On-Demand Resource
- Virtual Machines (VM)
- Serverless Functions (SF)

## Characteris
- Start-up time
- Usage cost



Input Events

Analytics

Output Stream

Outputs

Cluster

VMs

Start-up time: 25s+

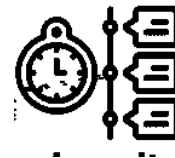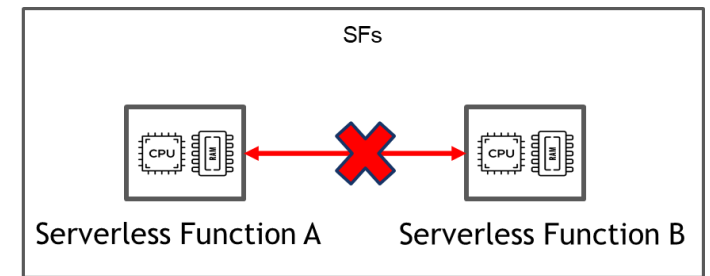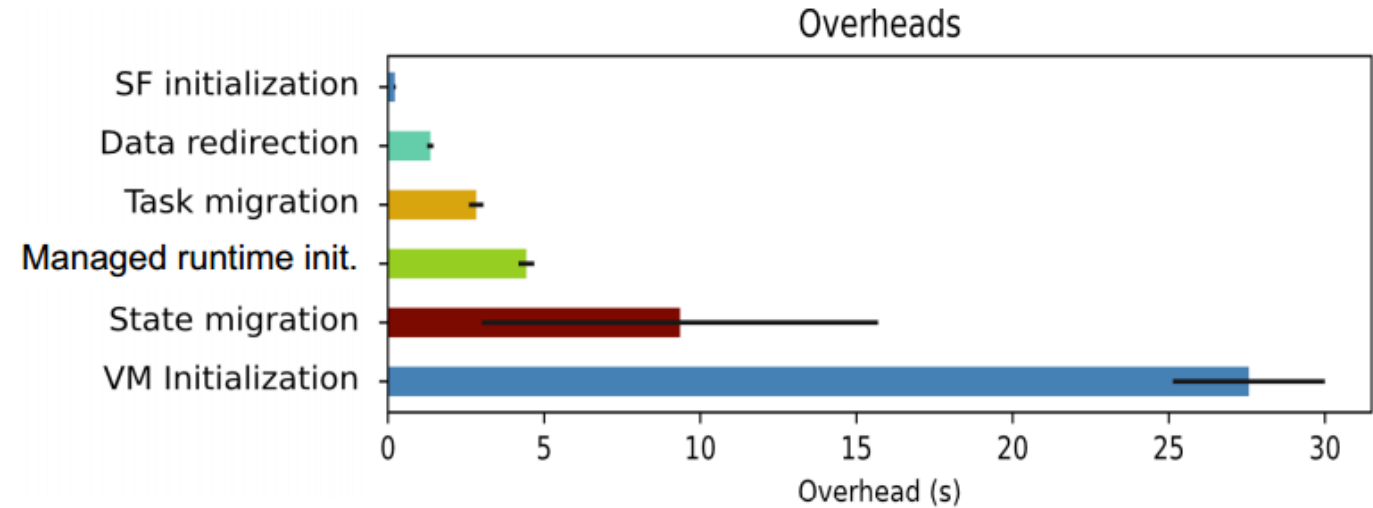Start-up time: 300-750ms

Usage cost： VM < SF
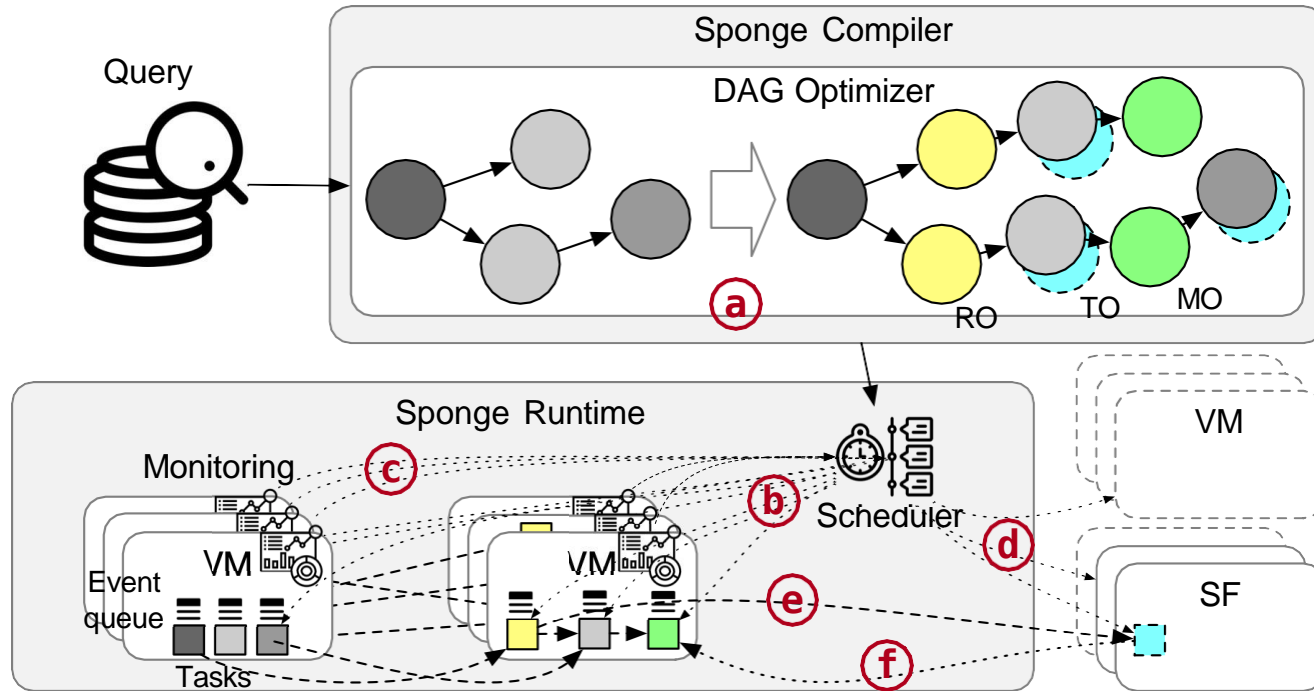
# Background

**Goal**

- **Quickly detect bursty loads and reduce migration state overhead.**



**Challenges**

- **Migration with large operator states.**
- **Indirect data communication between SF instances.**
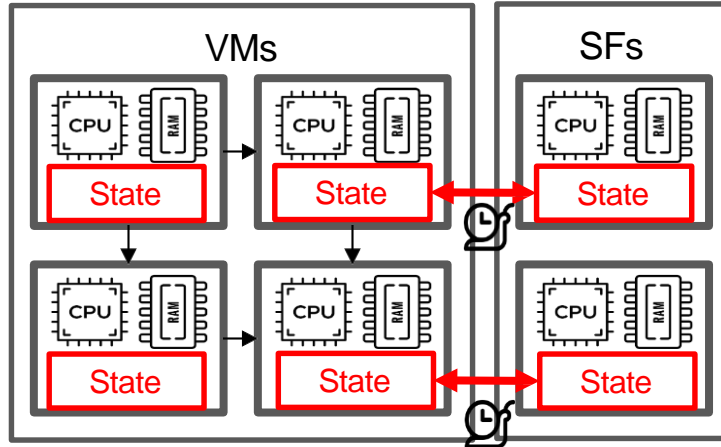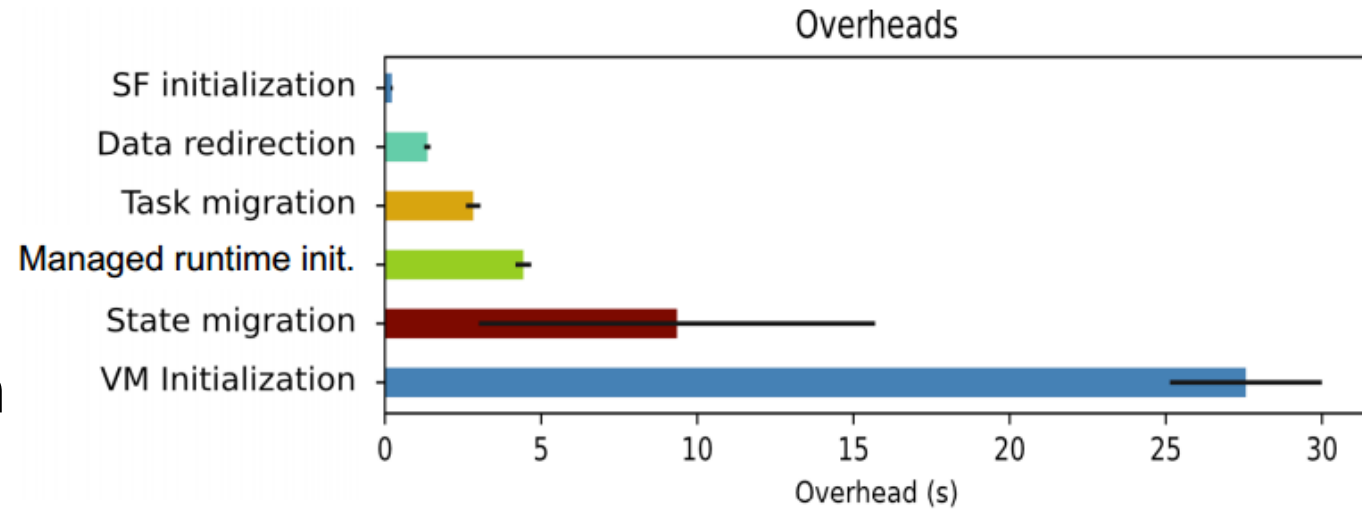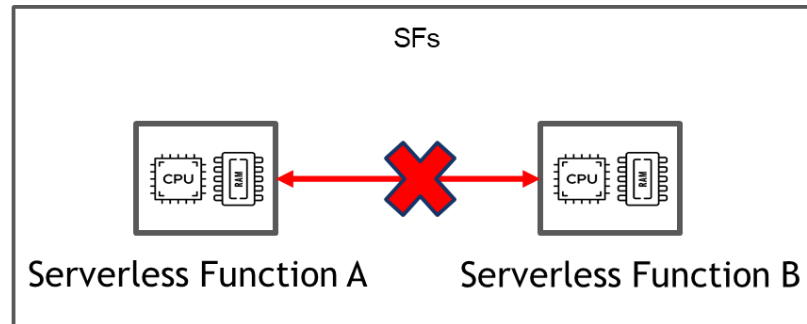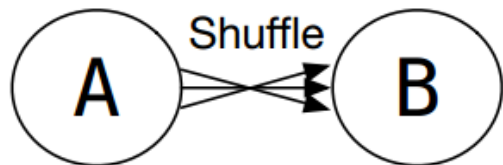- **Quick decision making and scaling.**

# Main Idea



- Redirect-and-merge
  - Compile-time Graph Rewriting Algorithm
  - Reducing Cold Start Latency
  - Watermark message

- Fast reactive scaling
  - Dynamic Offloading Policy
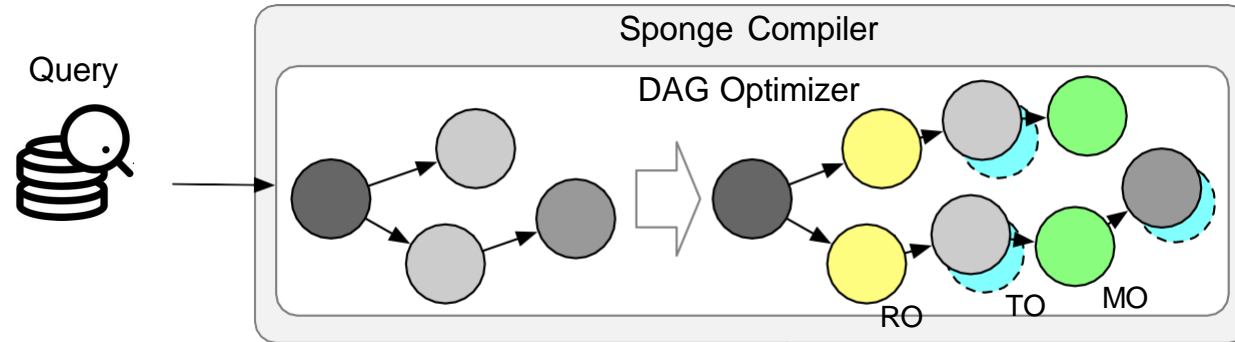
# Redirect-and-merge

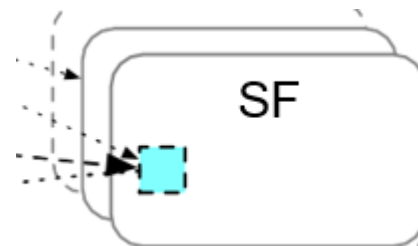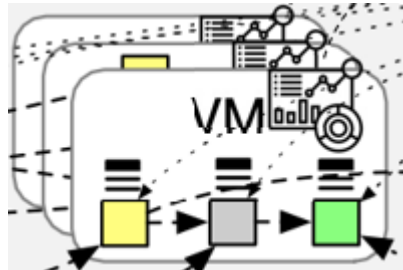- C1. Migration with large operator states.



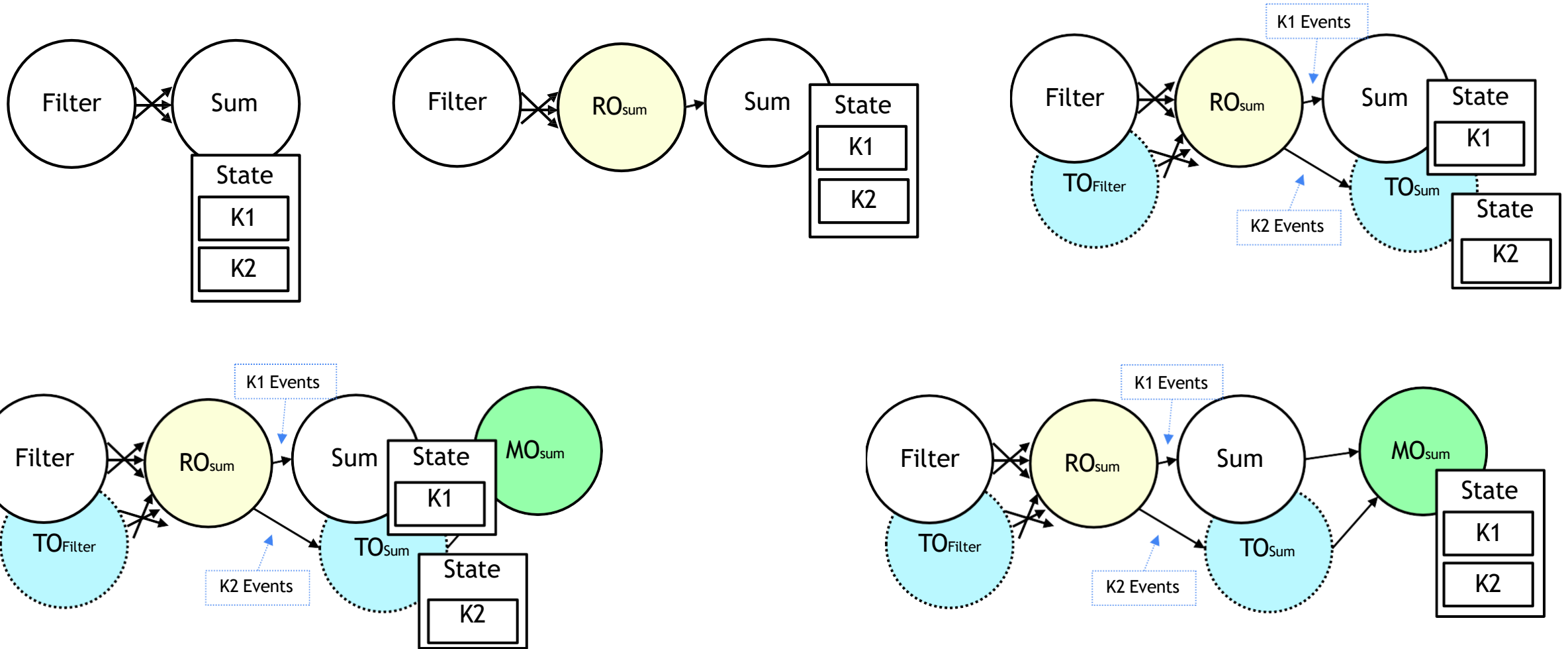- C2. Indirect data communication between SF instances.

# Design1: Compile-time Graph Rewriting Algorithm



1. Router operators (ROs) enable redirection of input events to specific instances

2. Transient operators (TOs) enable execution of cloned operators on SFs

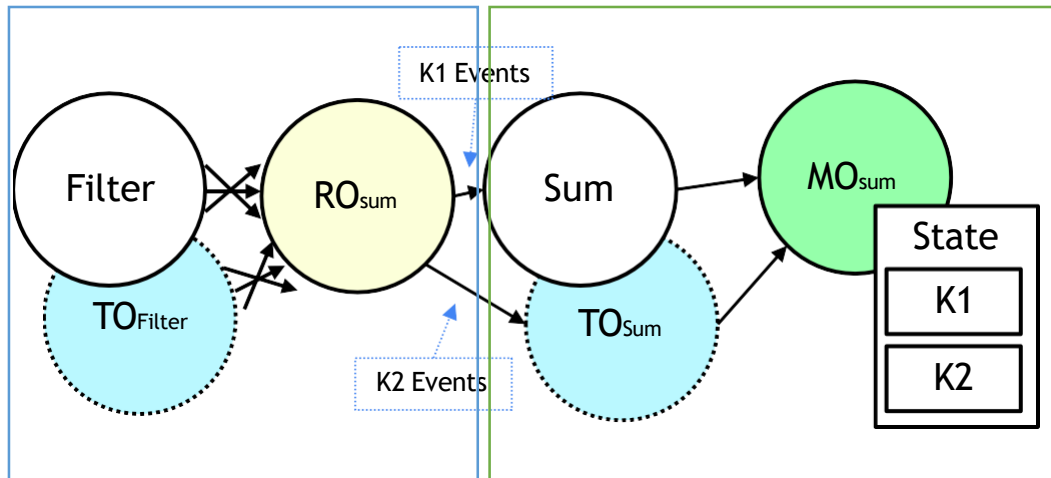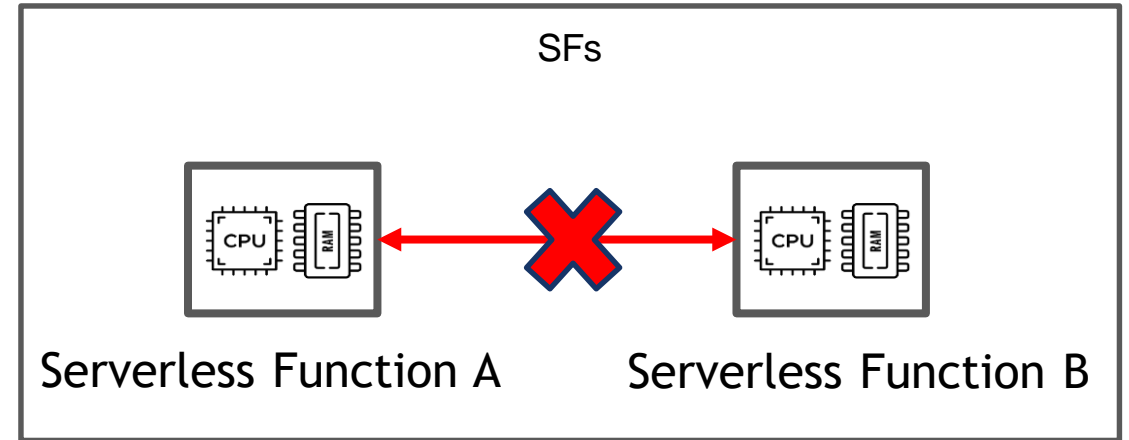3. Merge operators (MOs) enable merges on partial states

# Design1: Compile-time Graph Rewriting Algorithm
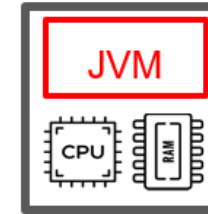
# Design1: Compile-time Graph Rewriting Algorithm

- C2. Indirect data communication between SF instances.
  - SF <-> origin VM
  - VM <-> VM



SFs

Serverless Function A    Serverless Function B



K1 Events

Filter    $RO_{sum}$    Sum    $MO_{sum}$

$TO_{Filter}$    $TO_{Sum}$

K2 Events

State
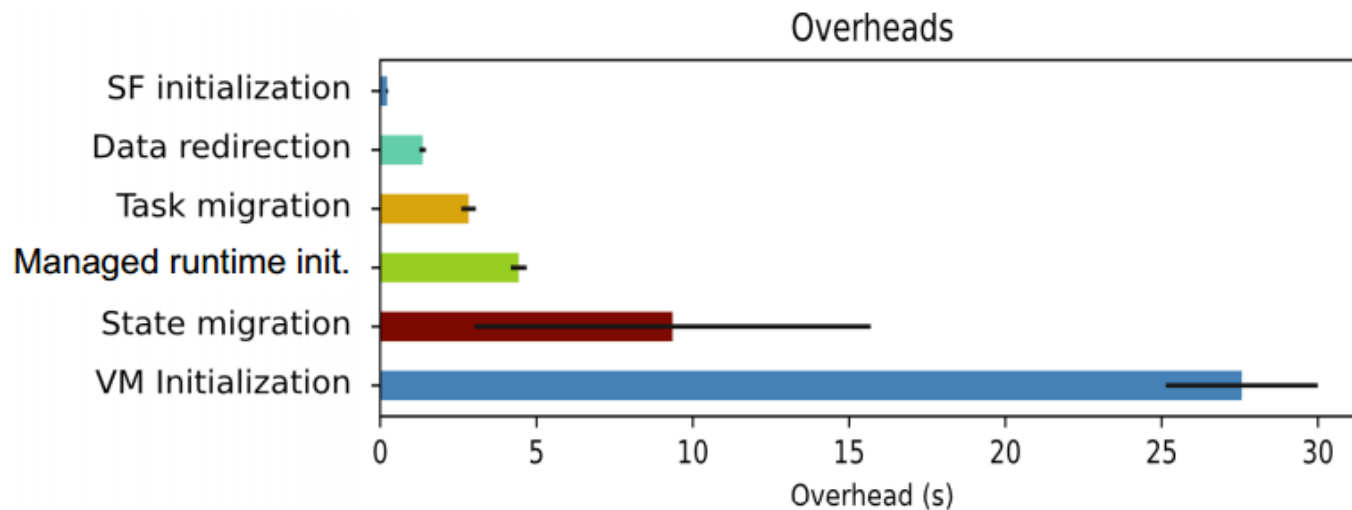
K1

K2

# Design2: Reducing Cold Start Latency

Timely gain access to SFs
- Warm-up SF workers
- Cache snapshots of SFs
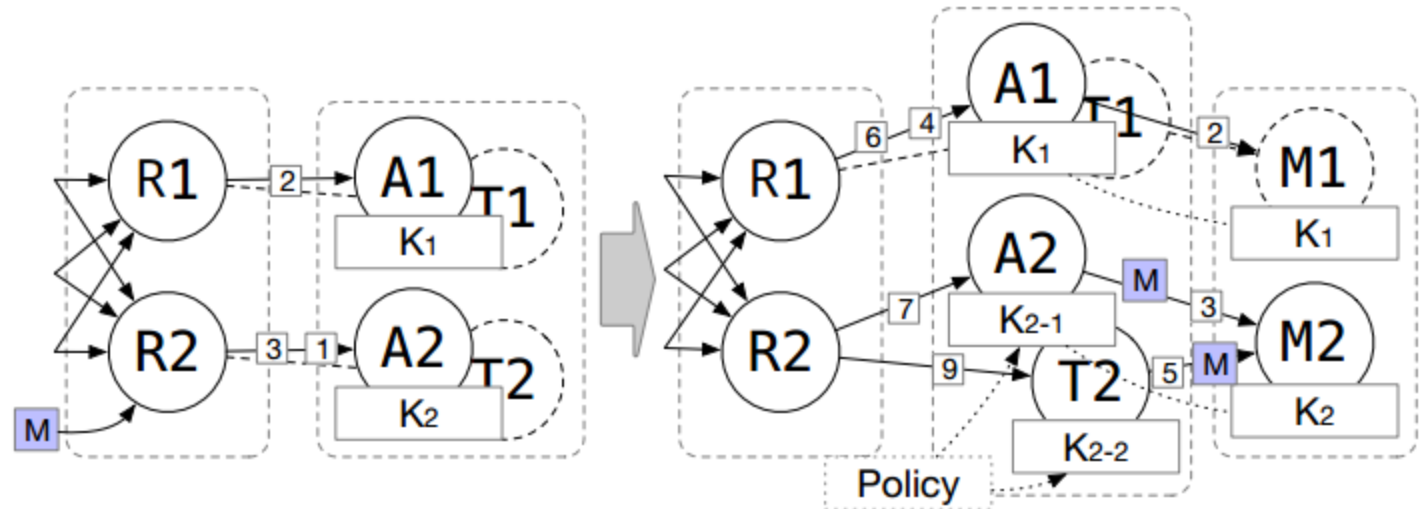
JVM

CPU RAM

Serverless Function

*Managed runtimes (e.g., JVM) incur launch overheads (~4 seconds)*


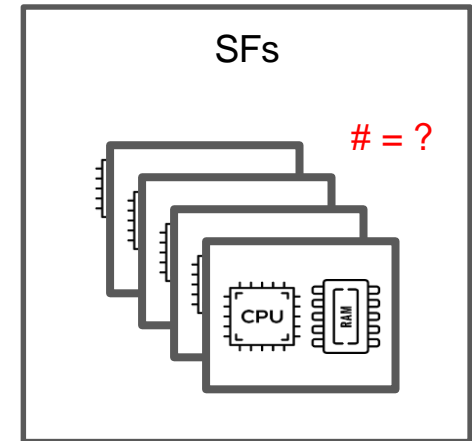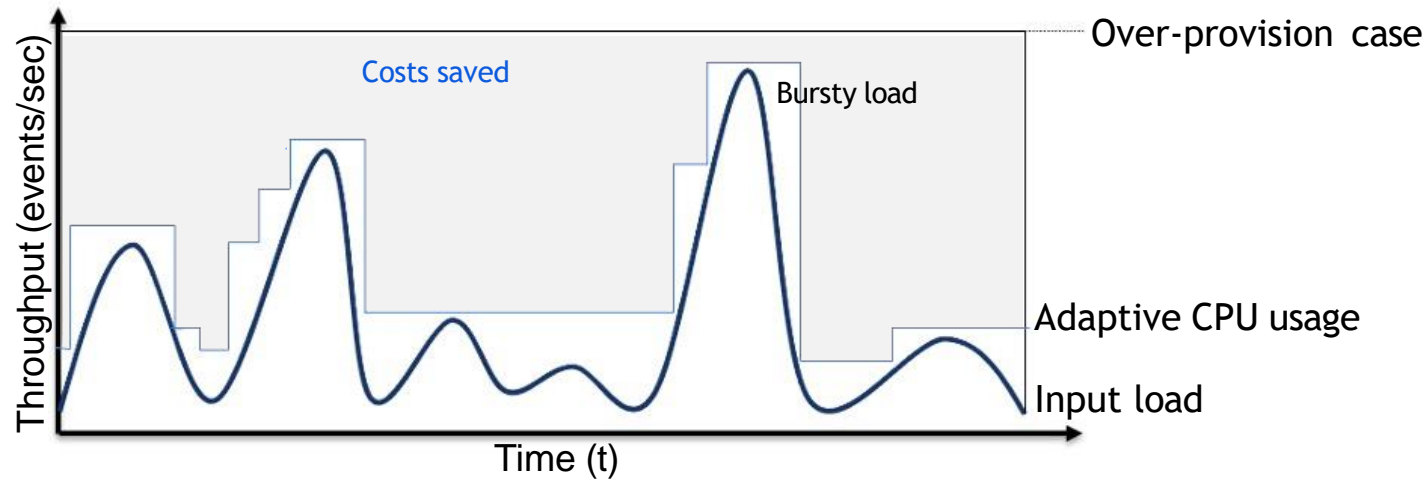Overheads

# Design3: Watermark message

Correctness
- Watermark as control message
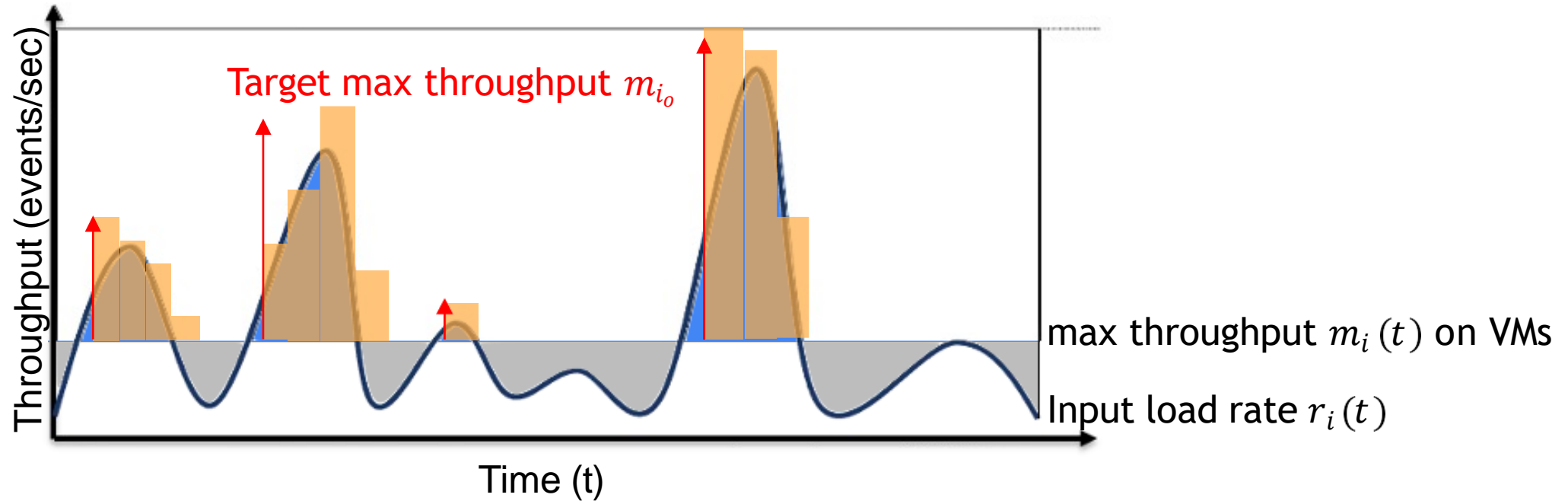- All events are processed in the same environment
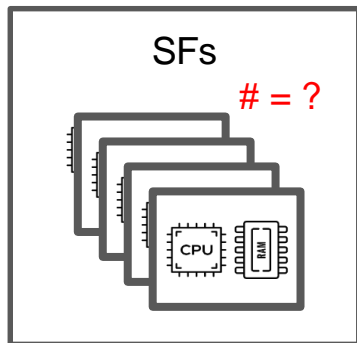
# Fast reactive scaling

- C3. Quick decision making and scaling.
- describe when Sponge triggers offloading, how many SF instances it uses
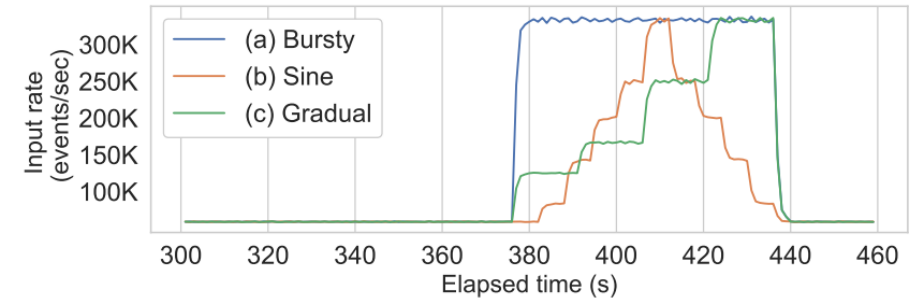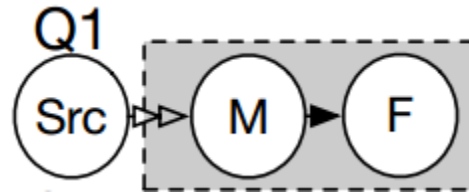
# Design4: Dynamic Offloading Policy



Throughput (events/sec)

Target max throughput $m_{i_o}$

max throughput $m_i(t)$ on VMs

Input load rate $r_i(t)$

Time (t)

*Data piled up in the event queue $\leq$ Data to process within our target deadline*
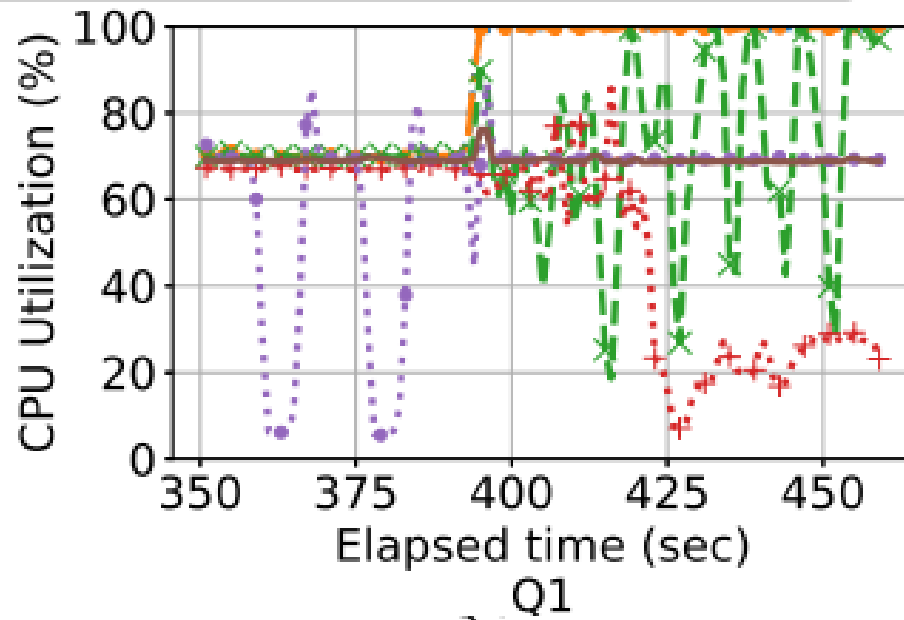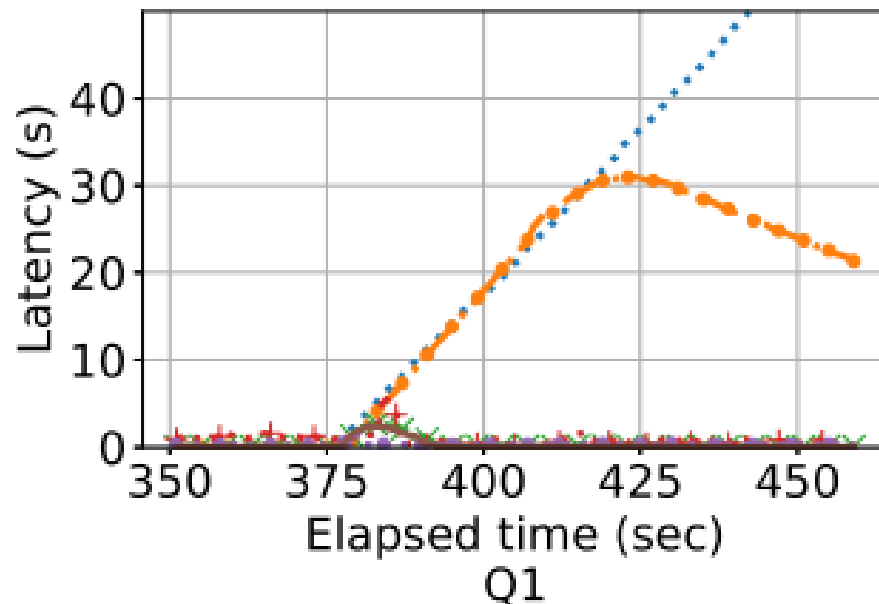*(Existing throughput $*$ time $\leq$ Target throughput $*$ time)*

SFs
# = ?

CPU    RAM

$$required\ SF\ cores = \left\lceil \frac{Target\ additional\ throughput}{70\% * Approx.\ throughput\ per\ SF\ core} \right\rceil$$

13

# Evaluation: Latency and CPU Utilizations

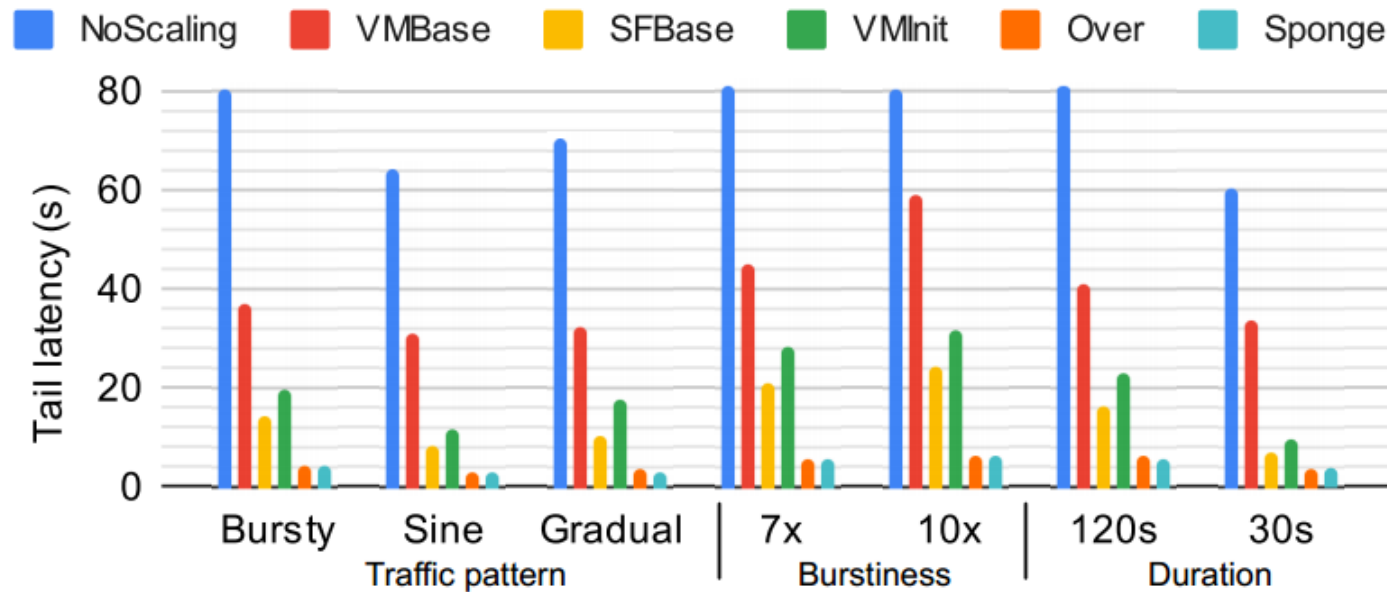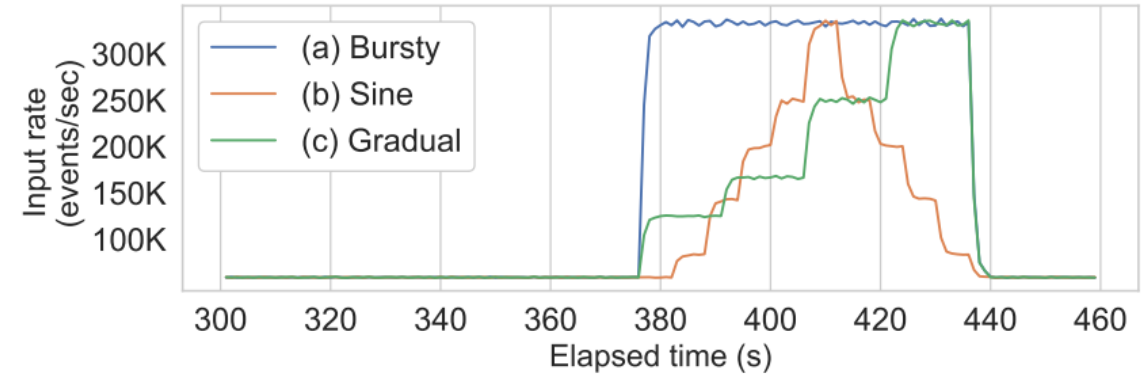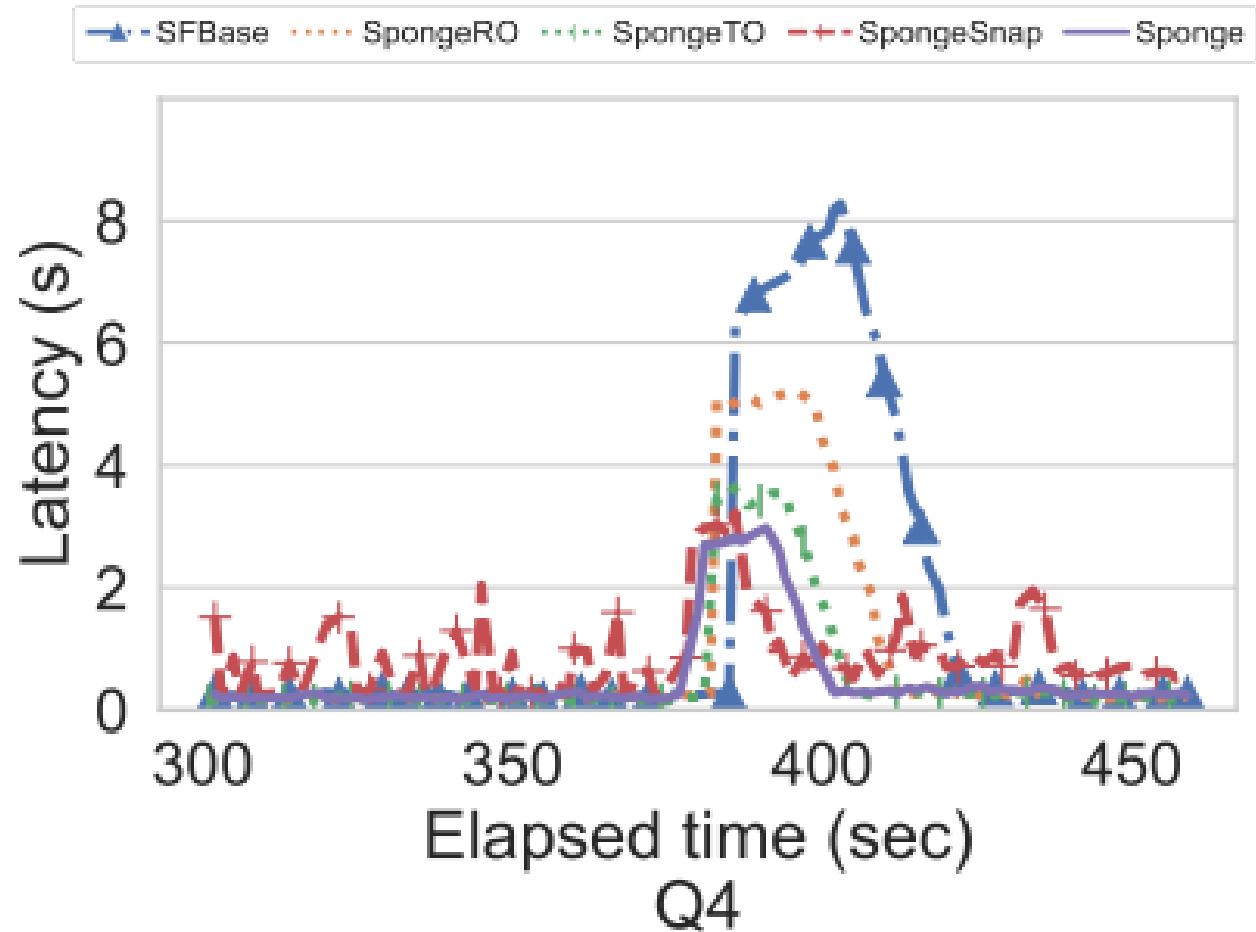The 99th-percentile tail latency and CPU utilization

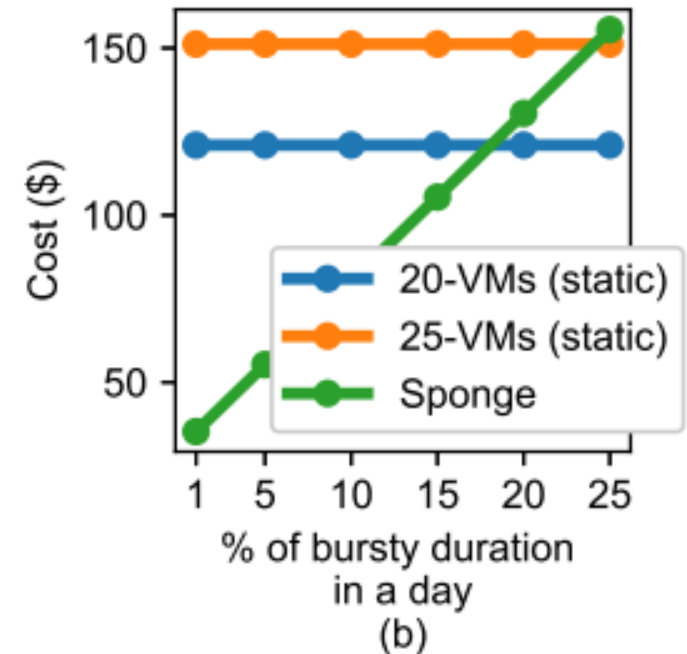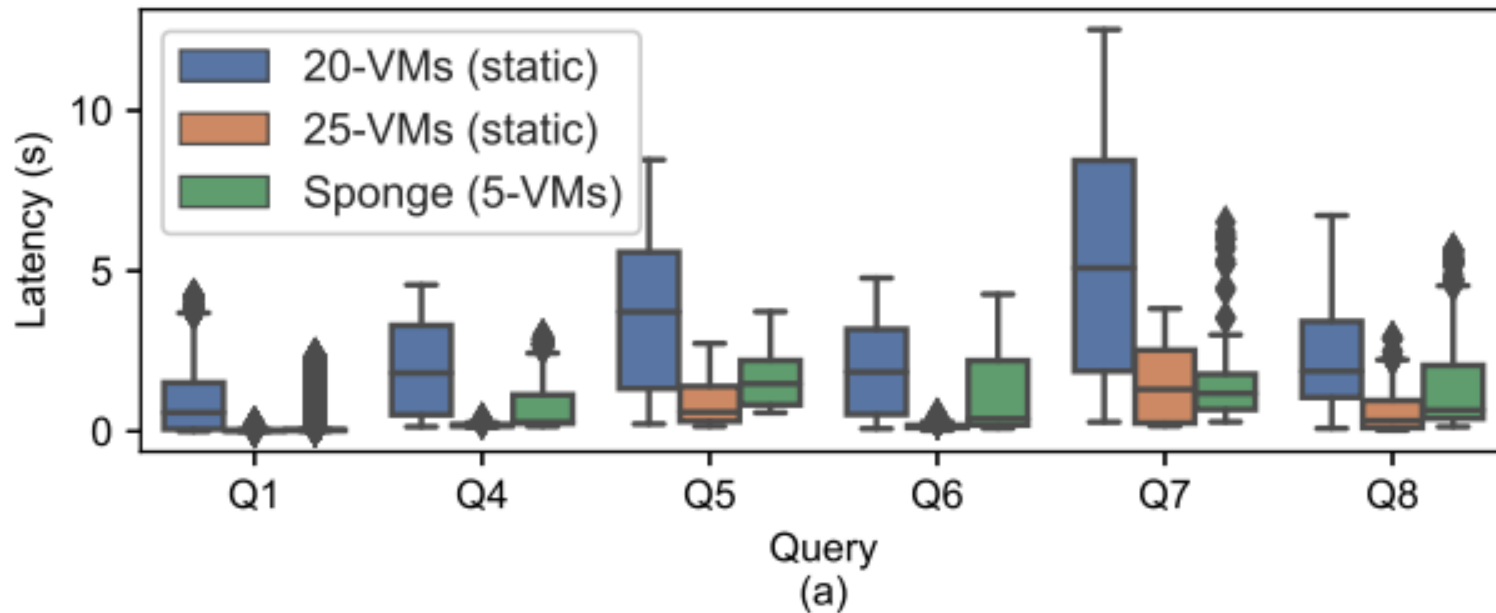# Evaluation: Latency and CPU Utilizations

Input patterns

# Evaluation: Graph Rewriting Effect

## SFBase + RO + TO + MO + Warm-up



Q4

# Evaluation: Latency-Cost Trade-Off

- Latency: 20-VMs > Sponge > 25-VMs
- Bursty duration < 15%
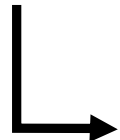
# Paper Summary

Stream Processing

↓

Bursty Load

↓

- Redirect-and-merge
- Fast reactive scaling

**Challenges**

- **Migration with large operator states.**
- **Indirect data communication between SF instances.**
- **Quick decision making and scaling.**

**Ideas**

- **Compile-time Graph**
- **Rewriting Algorithm**
- **Reducing Cold Start Latency**
- **Watermark message**
- **Dynamic Offloading Policy**