

# Move Fast and Meet Deadlines: Fine-grained Real-time Dataflow Scheduling with Cameo

Le Xu, Shivaram Venkataraman, Indranil Gupta, Luo Mai, Rahul Potharaju

NSDI 2021



UNIVERSITY OF  
**ILLINOIS**  
URBANA-CHAMPAIGN



**WISCONSIN**  
UNIVERSITY OF WISCONSIN-MADISON



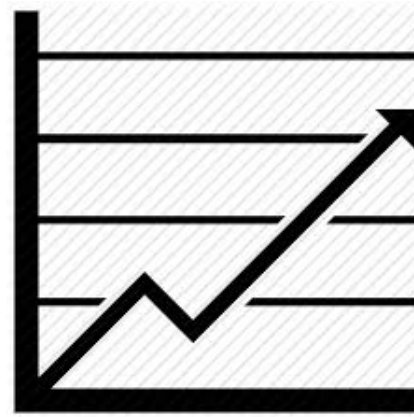
THE UNIVERSITY  
*of* EDINBURGH



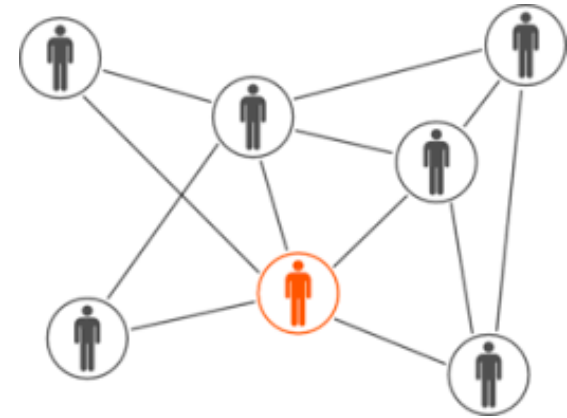
**Microsoft**



Forecast



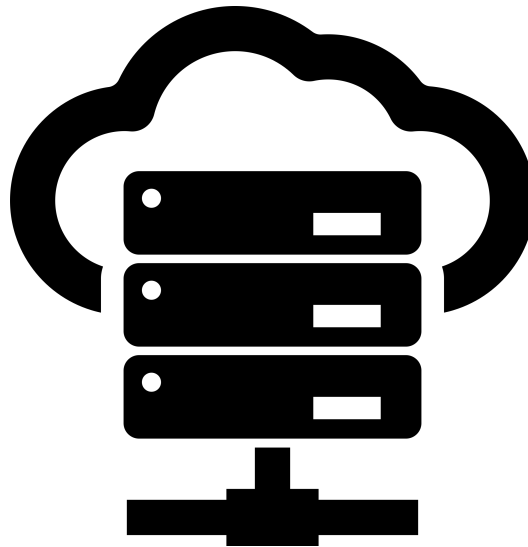
Financial Services



Social Networks



Edge and IoT



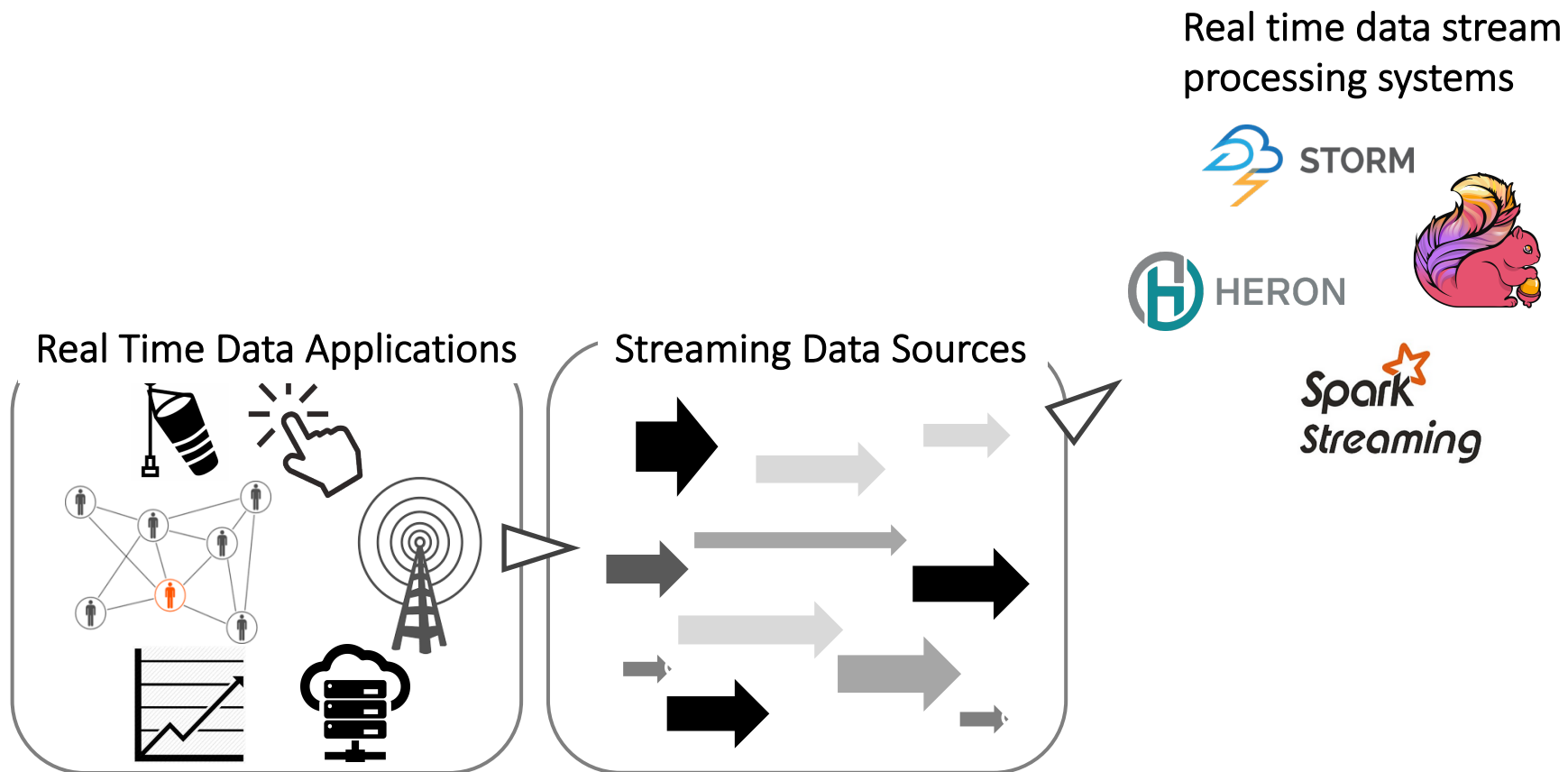
Cloud Management

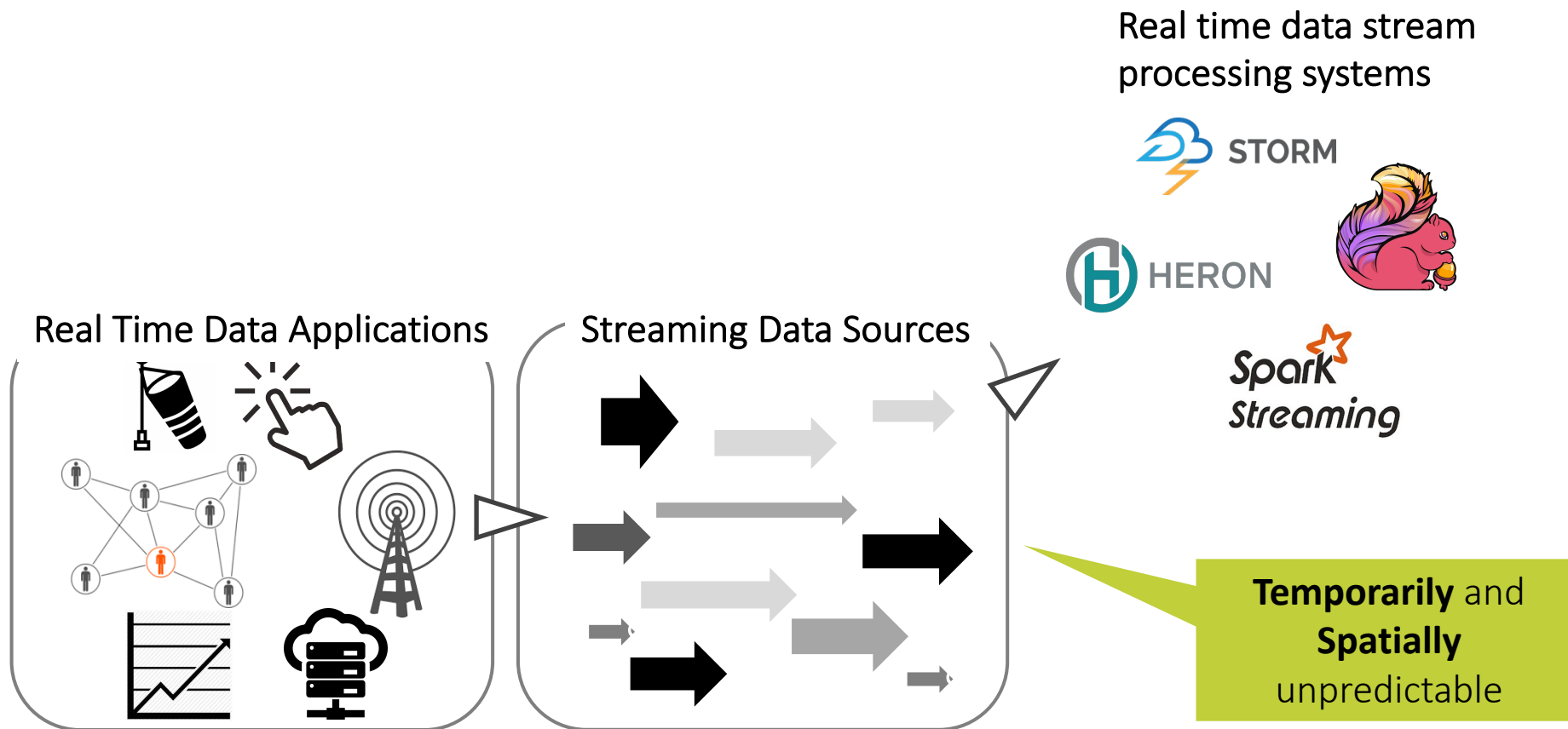


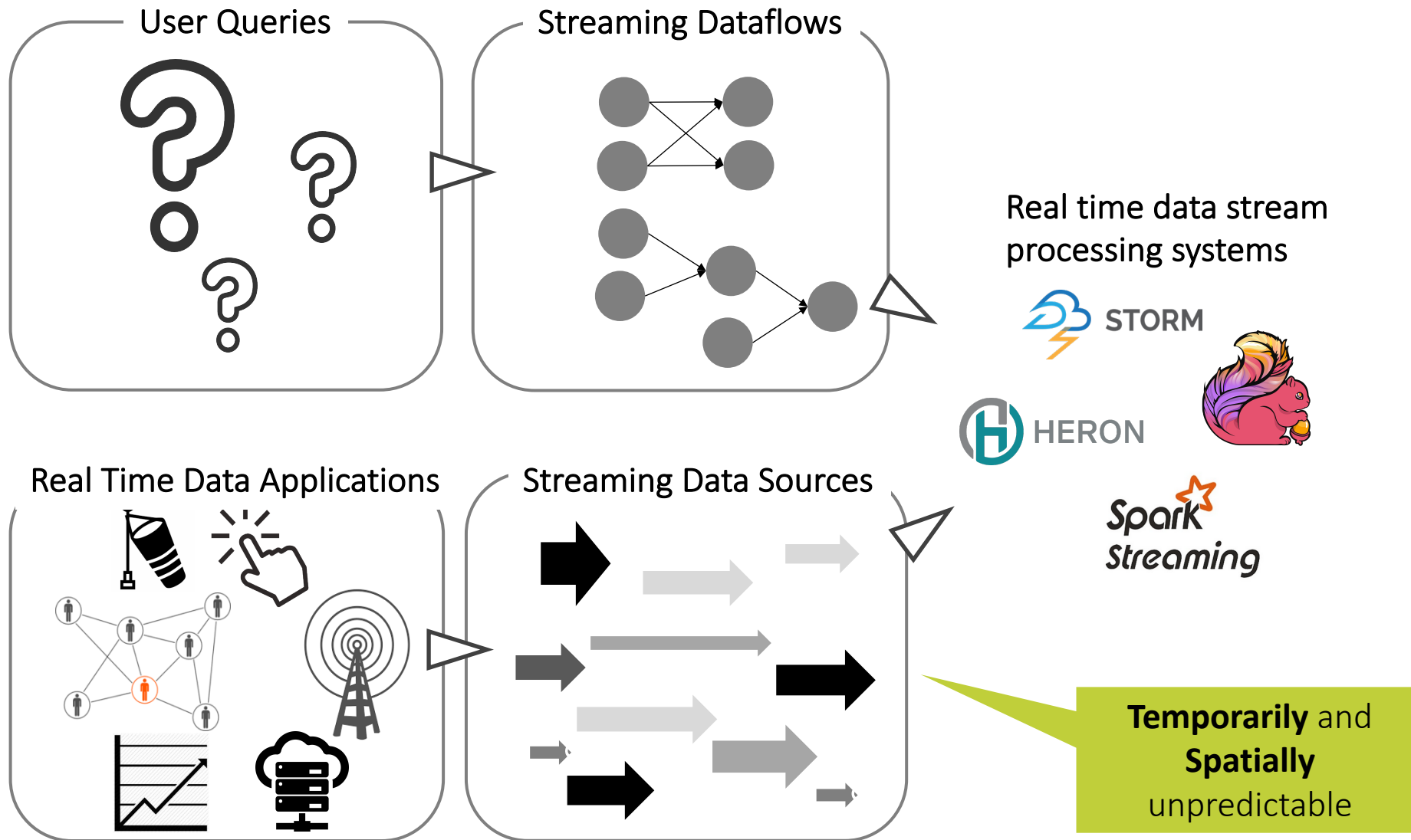
Web Applications

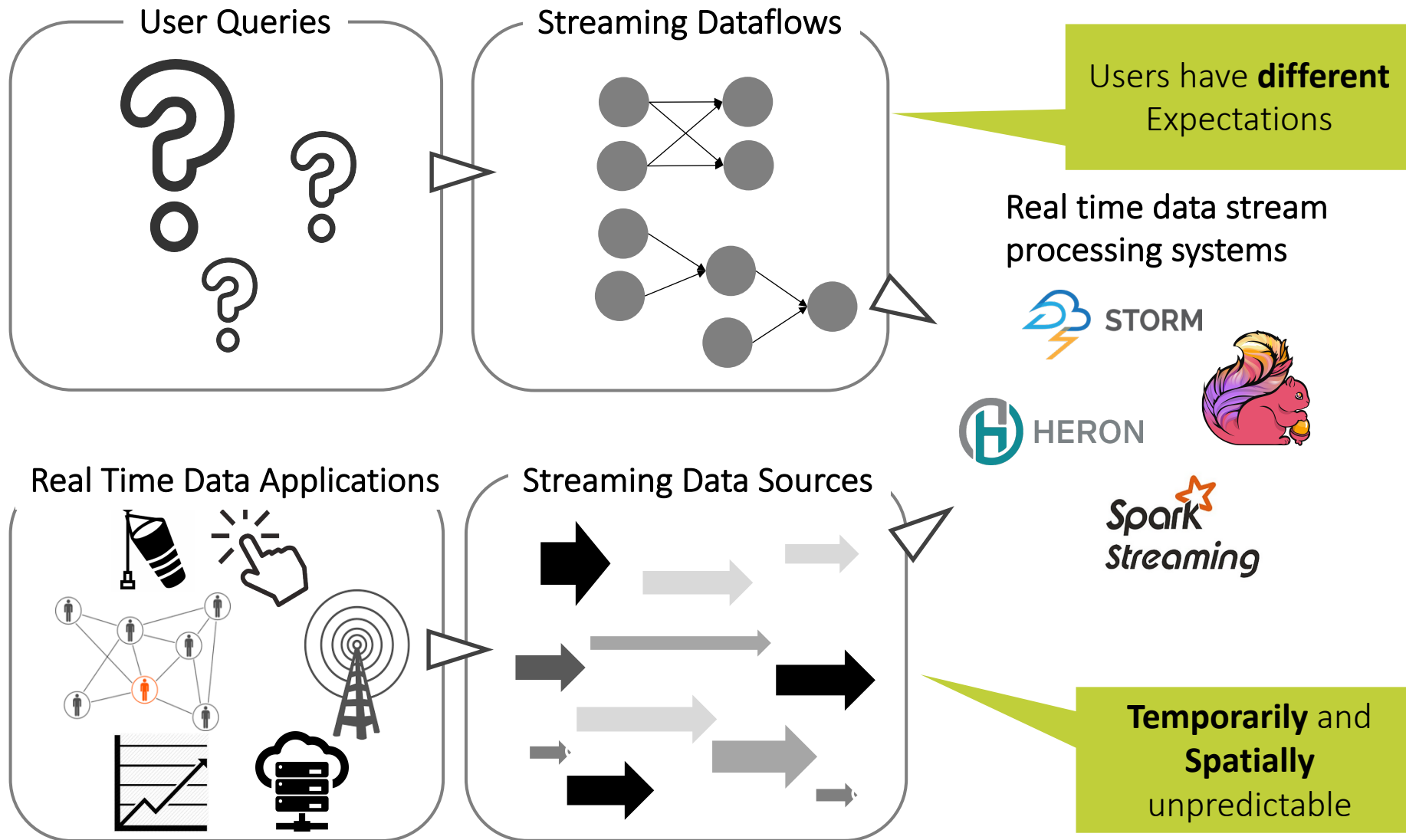
## Real time data stream processing systems

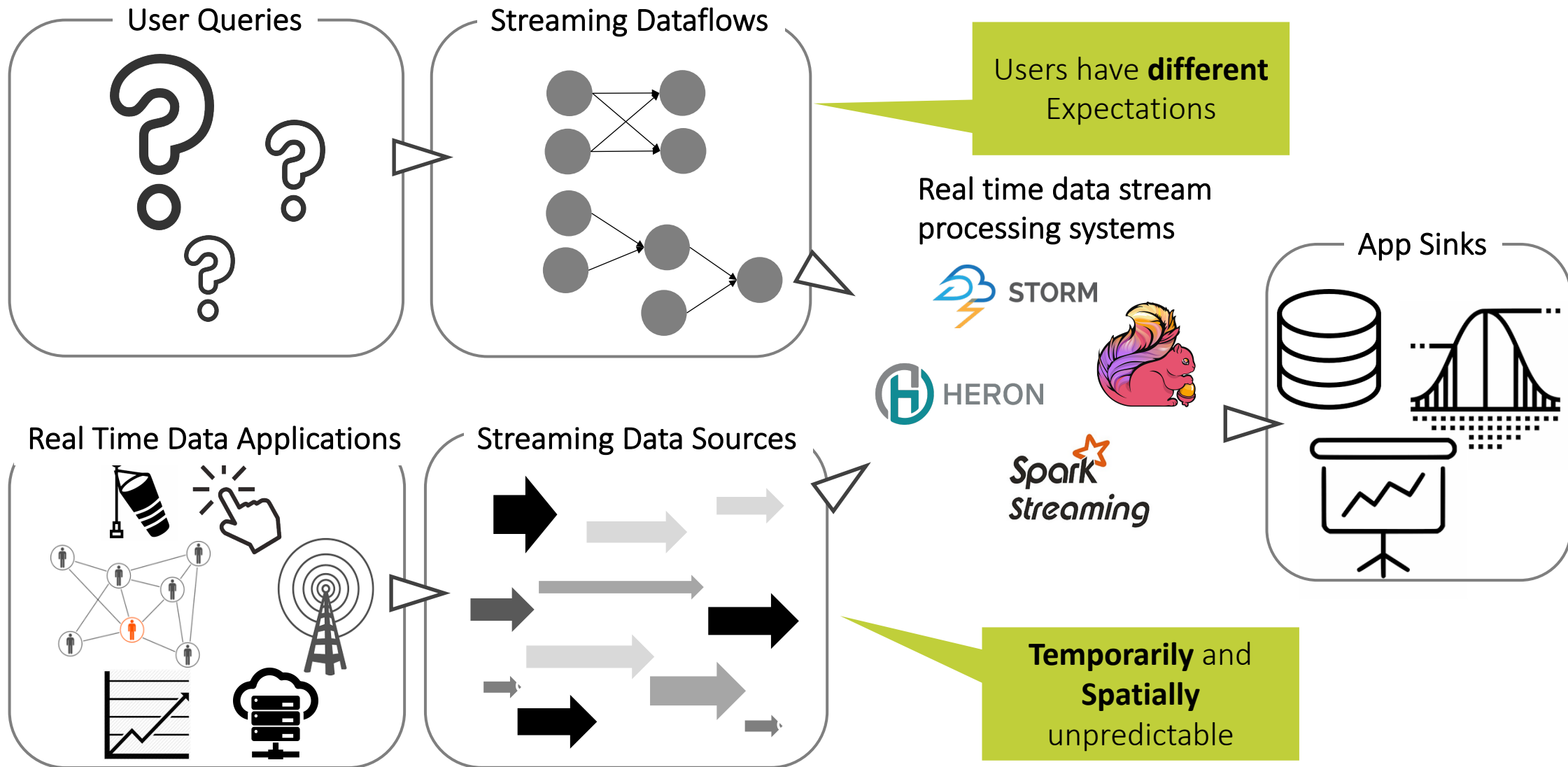




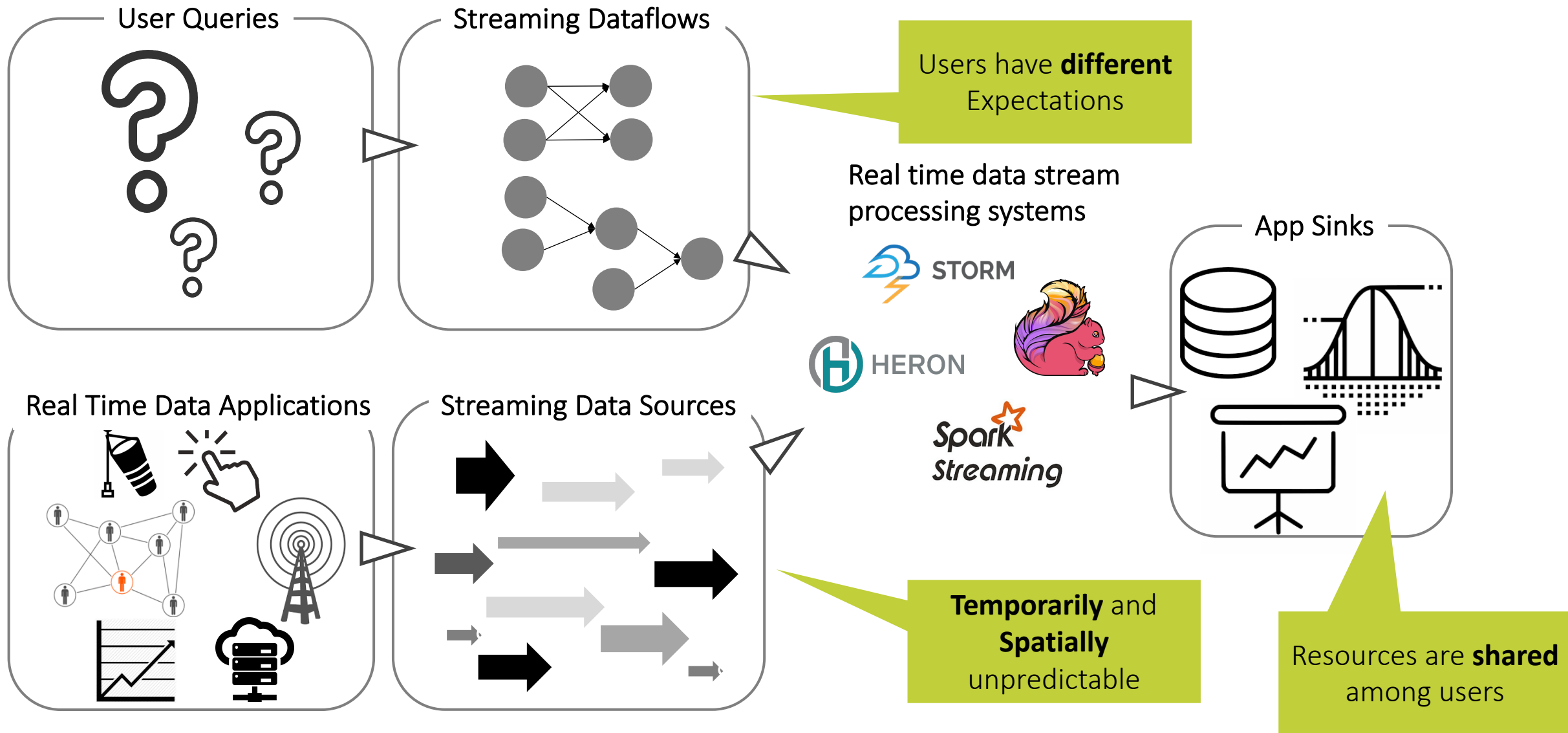


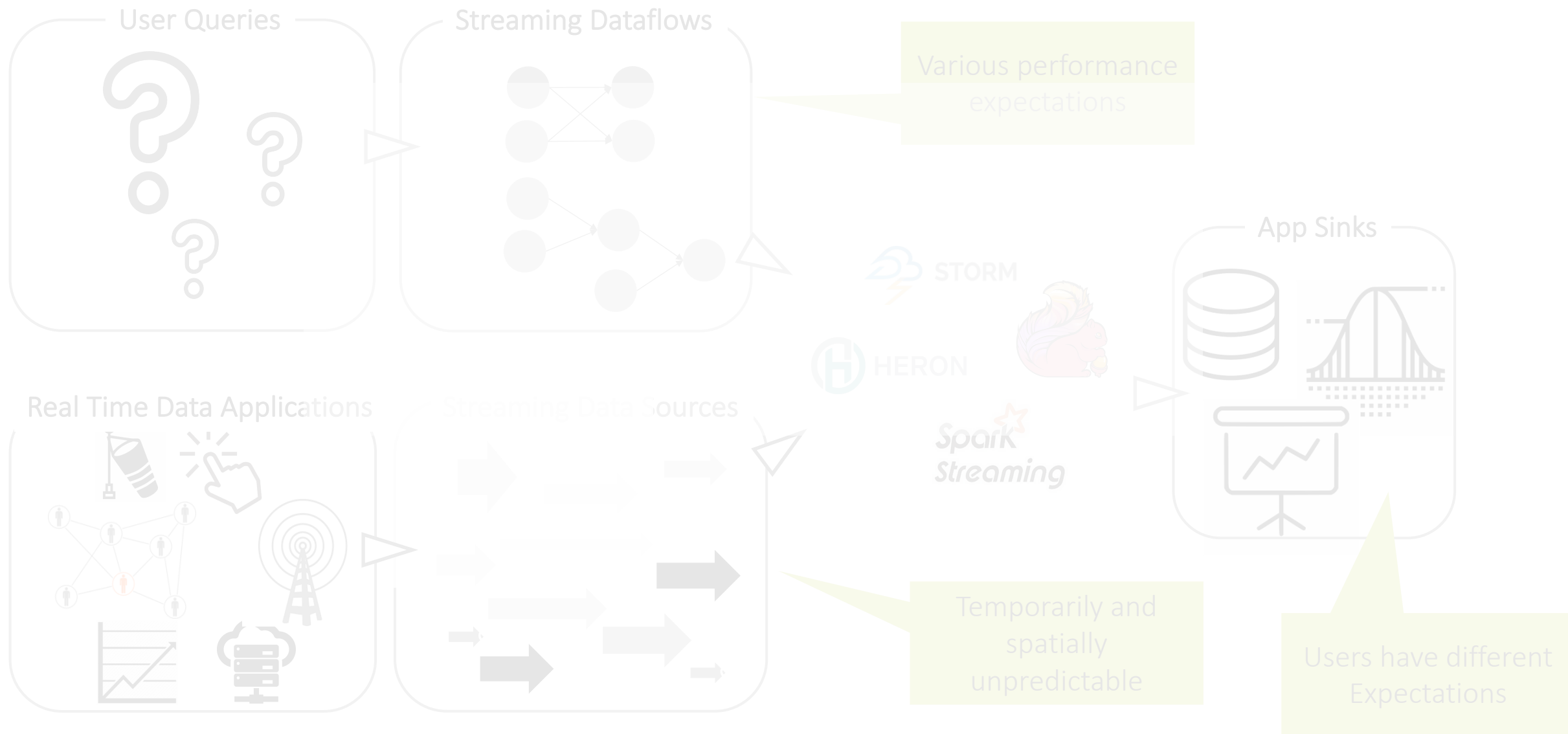


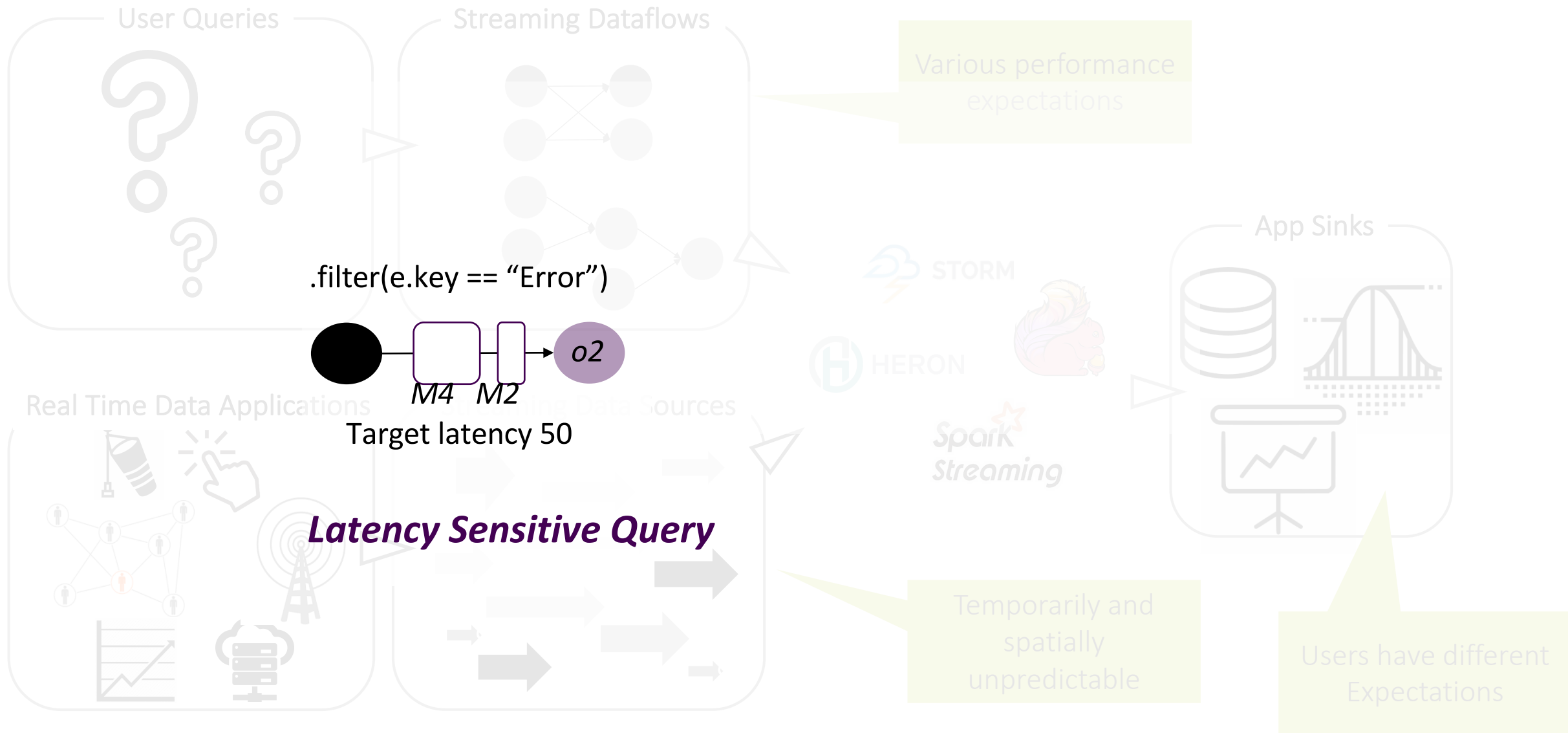


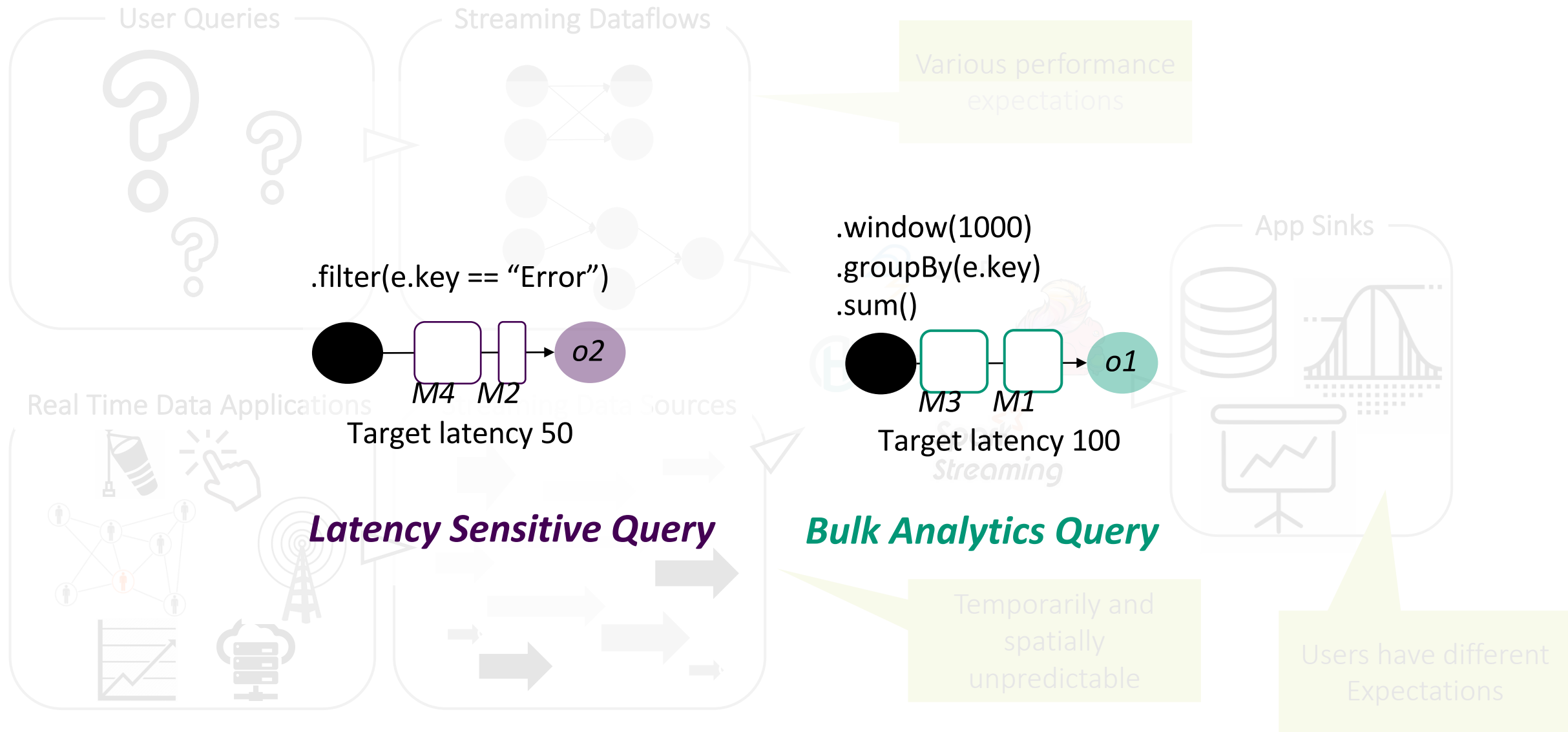














Guaranteed  
Performance

The diagram consists of two large, stylized arrows pointing towards each other. The left arrow is olive green and contains the text 'Guaranteed Performance'. The right arrow is teal and contains the text 'Resource Utilization'. The arrows are positioned such that their tips meet in the center, creating a symmetrical visual balance.

Resource  
Utilization



Guaranteed  
Performance

The diagram consists of two large, stylized arrows pointing towards each other. The left arrow is olive green and points right, containing the text 'Guaranteed Performance'. The right arrow is teal and points left, containing the text 'Resource Utilization'. The arrows are positioned such that their tips are close to each other, creating a central space between them.

Resource  
Utilization

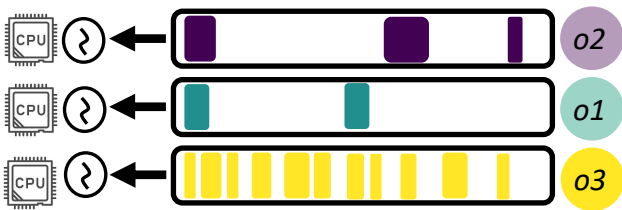


Guaranteed  
Performance

The diagram consists of two large, stylized arrows pointing towards each other. The left arrow is olive green and points right, containing the text 'Guaranteed Performance'. The right arrow is teal and points left, containing the text 'Resource Utilization'. The arrows are positioned such that their tips are close to each other, creating a central space between them.

Resource  
Utilization

## Full Resource Isolation

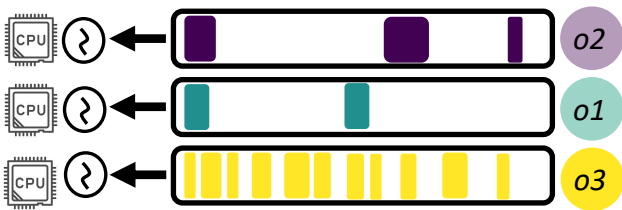


Guaranteed  
Performance

Resource  
Utilization



## Full Resource Isolation

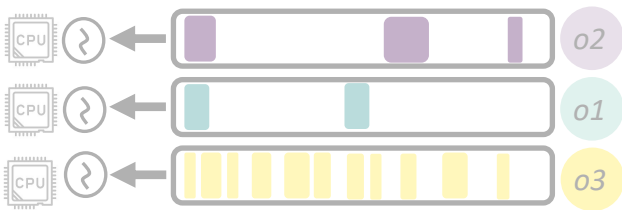


Guaranteed  
Performance

Resource  
Utilization

- + Unaffected performance under workload peaks
- Resource under-utilization and high cost

## Full Resource Isolation

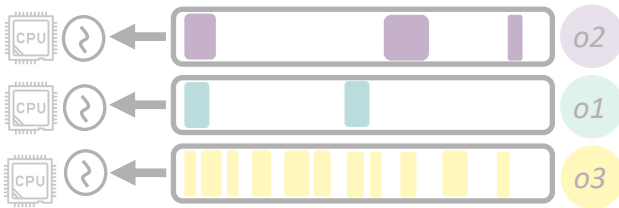


Guaranteed  
Performance

Resource  
Utilization

- + Unaffected performance under workload peaks
- Resource under-utilization and high cost

## Full Resource Isolation

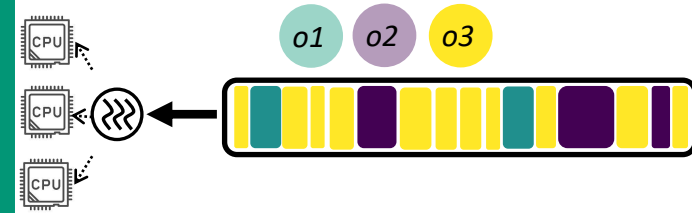


Guaranteed  
Performance

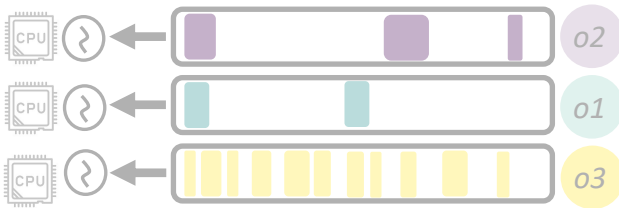
- + Unaffected performance under workload peaks
- Resource under-utilization and high cost

Resource  
Utilization

## No Resource Isolation



## Full Resource Isolation

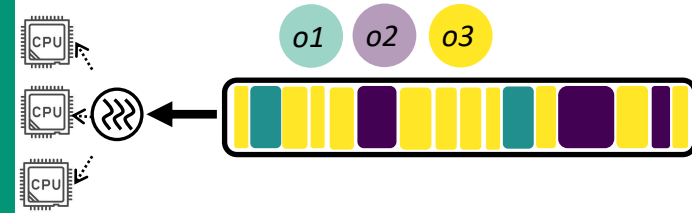


Guaranteed  
Performance

- + Unaffected performance under workload peaks
- Resource under-utilization and high cost

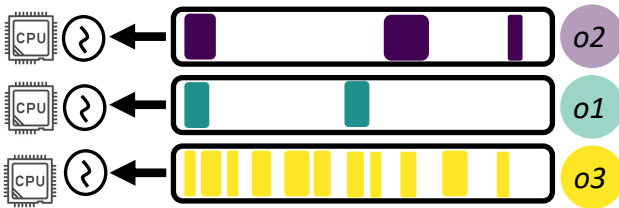
Resource  
Utilization

## No Resource Isolation



- + Uses resource efficiently
- Compromised performance (e.g., high processing latency)

## Full Resource Isolation

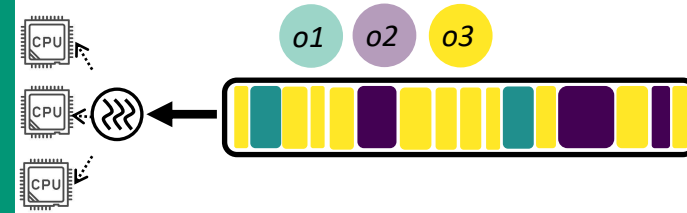


Guaranteed  
Performance

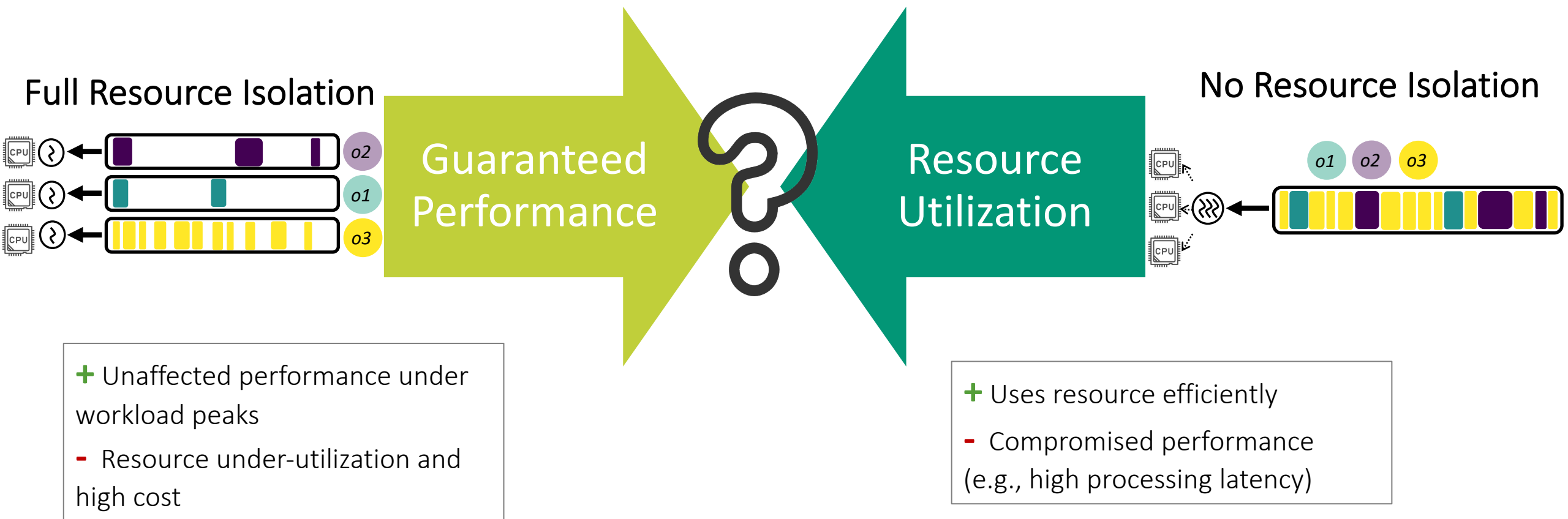
- + Unaffected performance under workload peaks
- Resource under-utilization and high cost

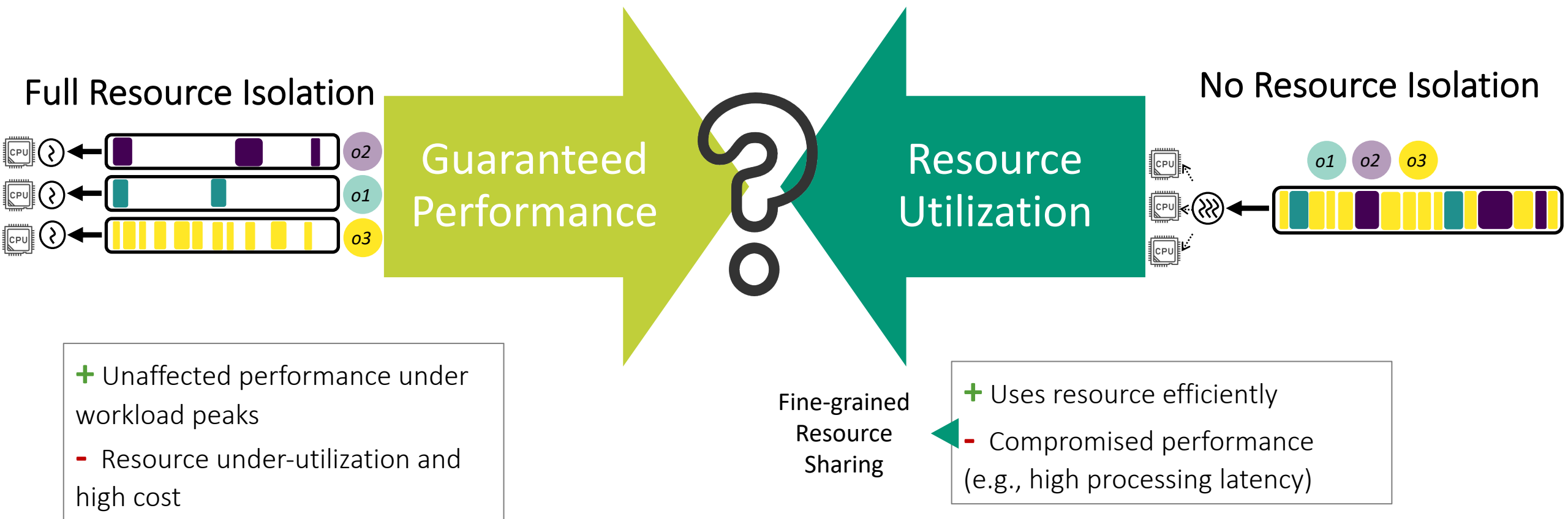
Resource  
Utilization

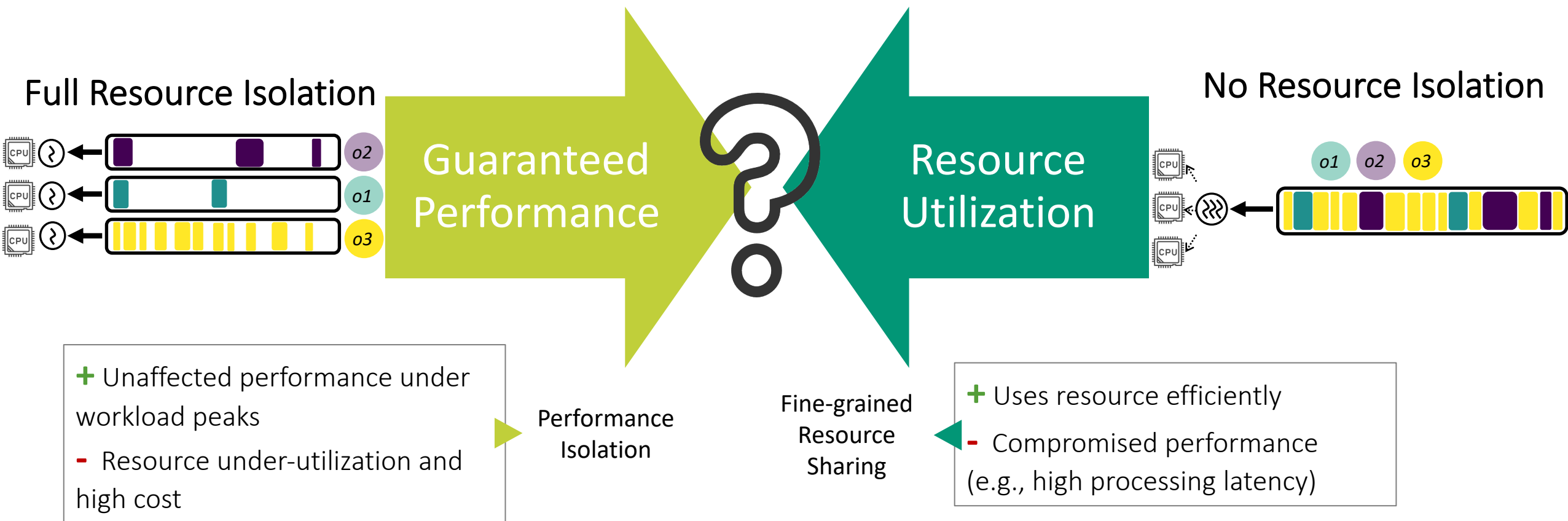
## No Resource Isolation



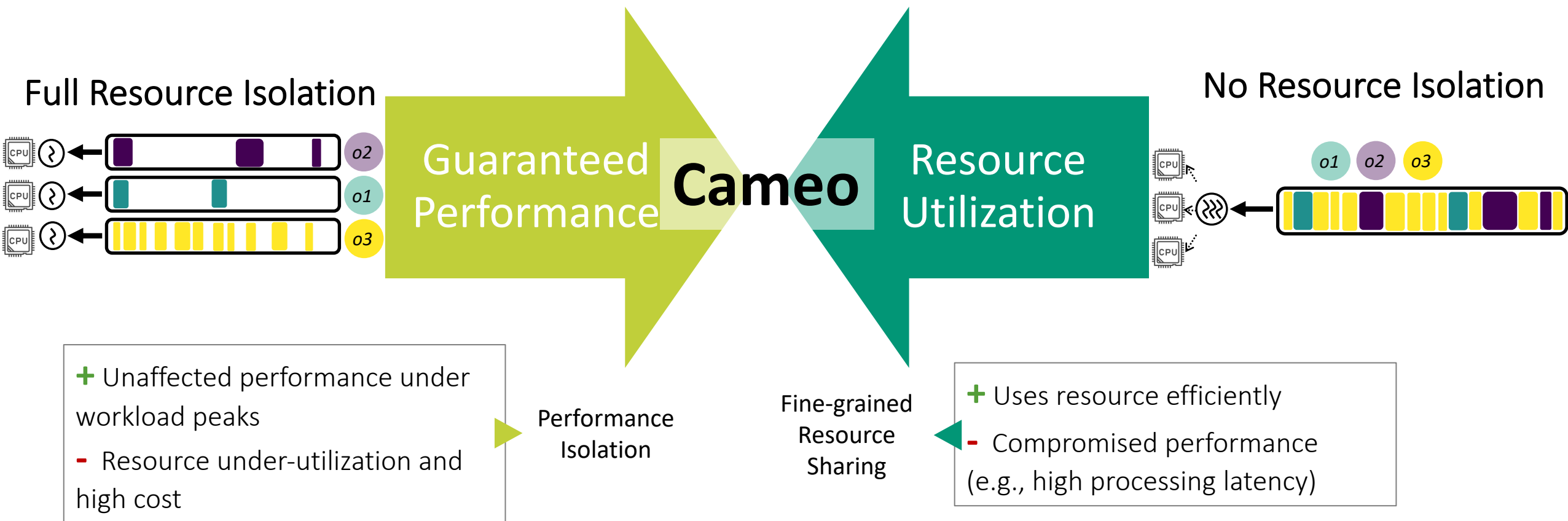
- + Uses resource efficiently
- Compromised performance (e.g., high processing latency)

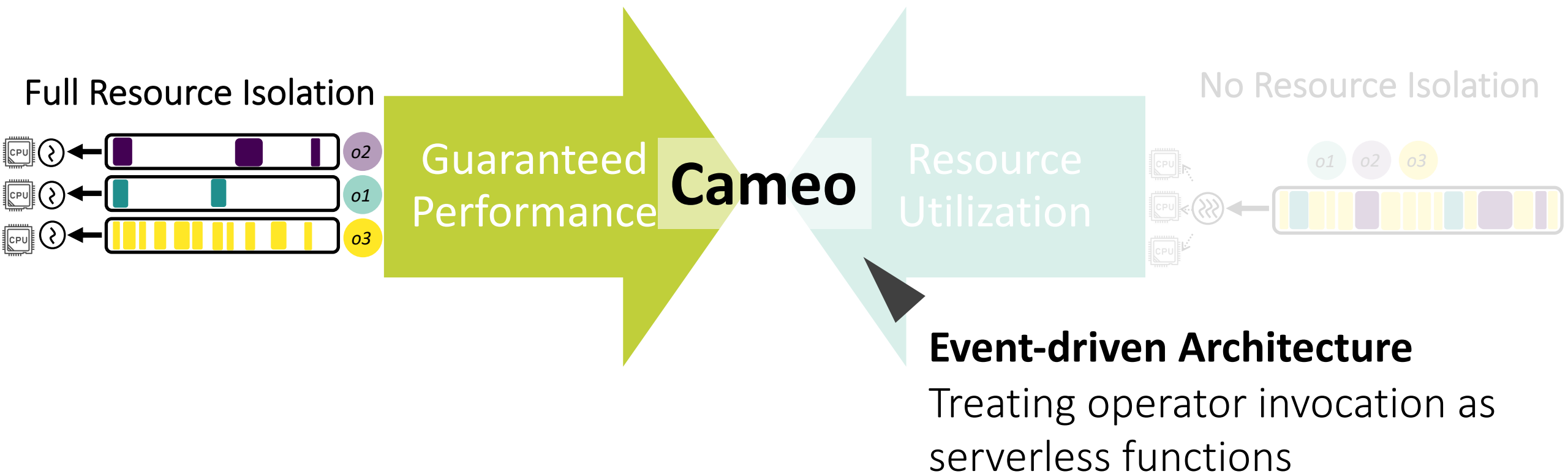


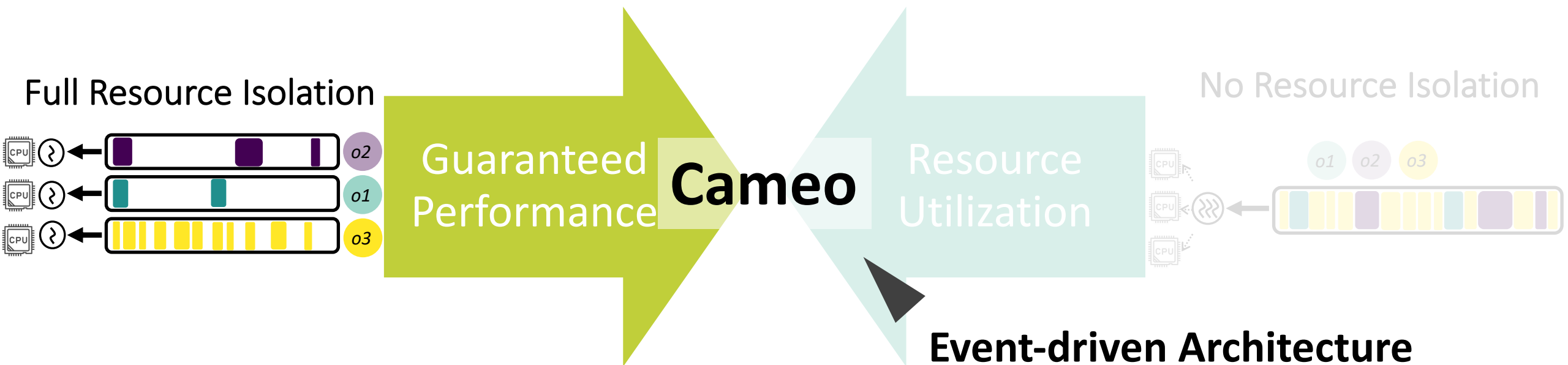








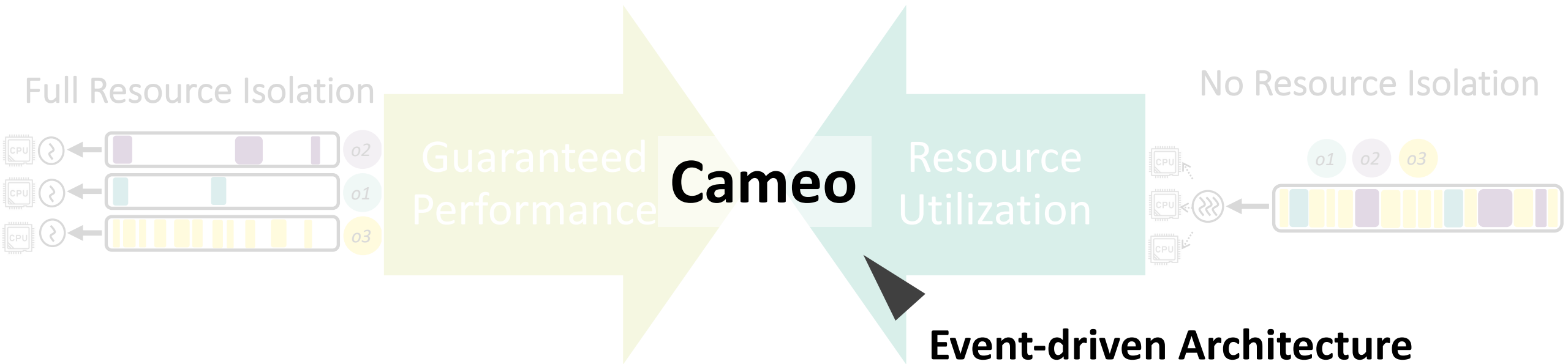




## Event-driven Architecture

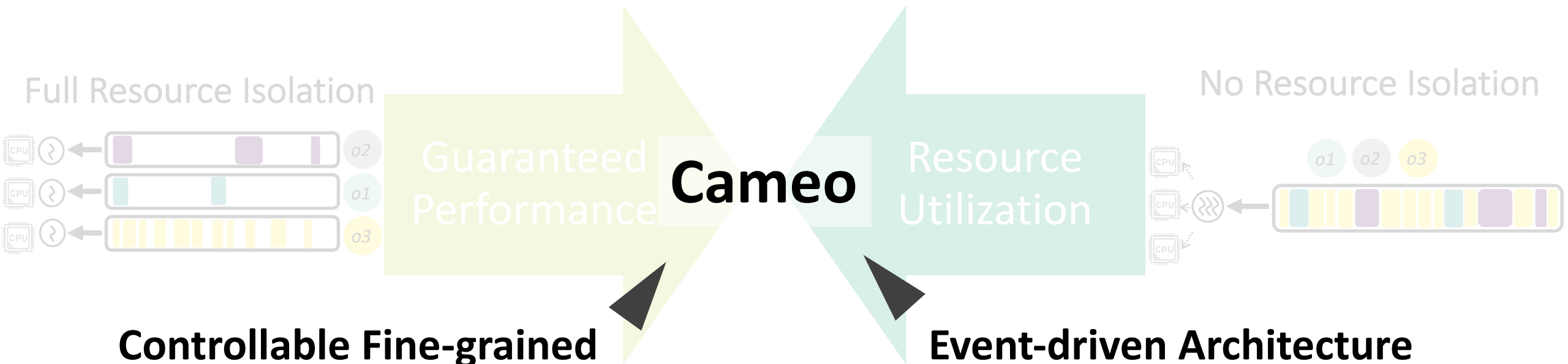
Treating operator invocation as serverless functions





**Event-driven Architecture**  
Treating operator invocation as  
serverless functions





## Controllable Fine-grained Scheduling

Identifying and prioritizing important work

## Event-driven Architecture

Treating operator invocation as serverless functions



## Cameo

Data-driven, fine-grained  
operator scheduling

## Cameo Strategy

- Topology and semantics aware
- Driven by message deadline

**Cameo**  
Data-driven, fine-grained  
operator scheduling

## Cameo Strategy

- Topology and semantics aware
- Driven by message deadline

## Cameo Mechanism

- Light-weight, stateless scheduler
- Scalable priority generation
- Pluggable strategy

**Cameo**  
Data-driven, fine-grained  
operator scheduling



## Cameo Strategy

- Topology and semantics aware
- Driven by message deadline

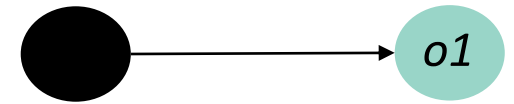
## Cameo Mechanism

- Light-weight, stateless scheduler
- Scalable priority generation
- Pluggable strategy

**Cameo**  
**Data-driven**, fine-grained  
operator scheduling

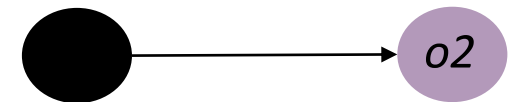
# Cameo Strategy

```
.window(1000)  
.groupBy(e.key)  
.sum()
```



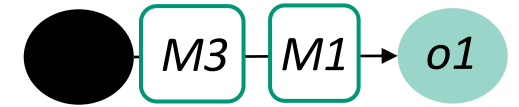
*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



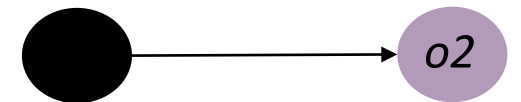
*Latency Sensitive Query*

```
.window(1000)  
.groupBy(e.key)  
.sum()
```



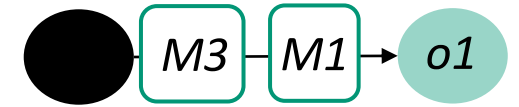
***Bulk Analytics Query***

```
.filter(e.key == "Error")
```



***Latency Sensitive Query***

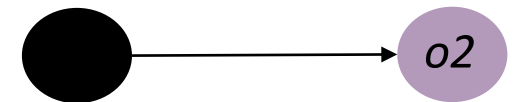
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

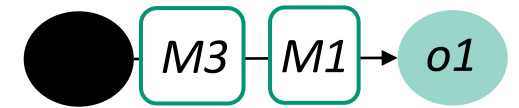
***Bulk Analytics Query***

```
.filter(e.key == "Error")
```



***Latency Sensitive Query***

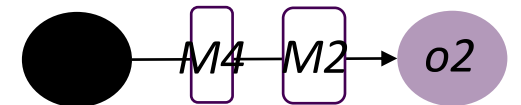
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

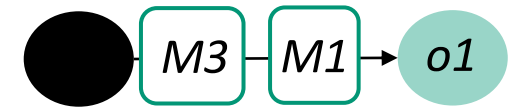
***Bulk Analytics Query***

```
.filter(e.key == "Error")
```



***Latency Sensitive Query***

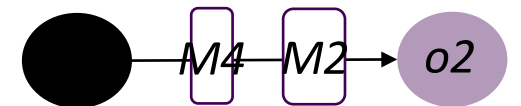
```
.window(1000)
.groupBy(e.key)
.sum()
```



Target latency **100**

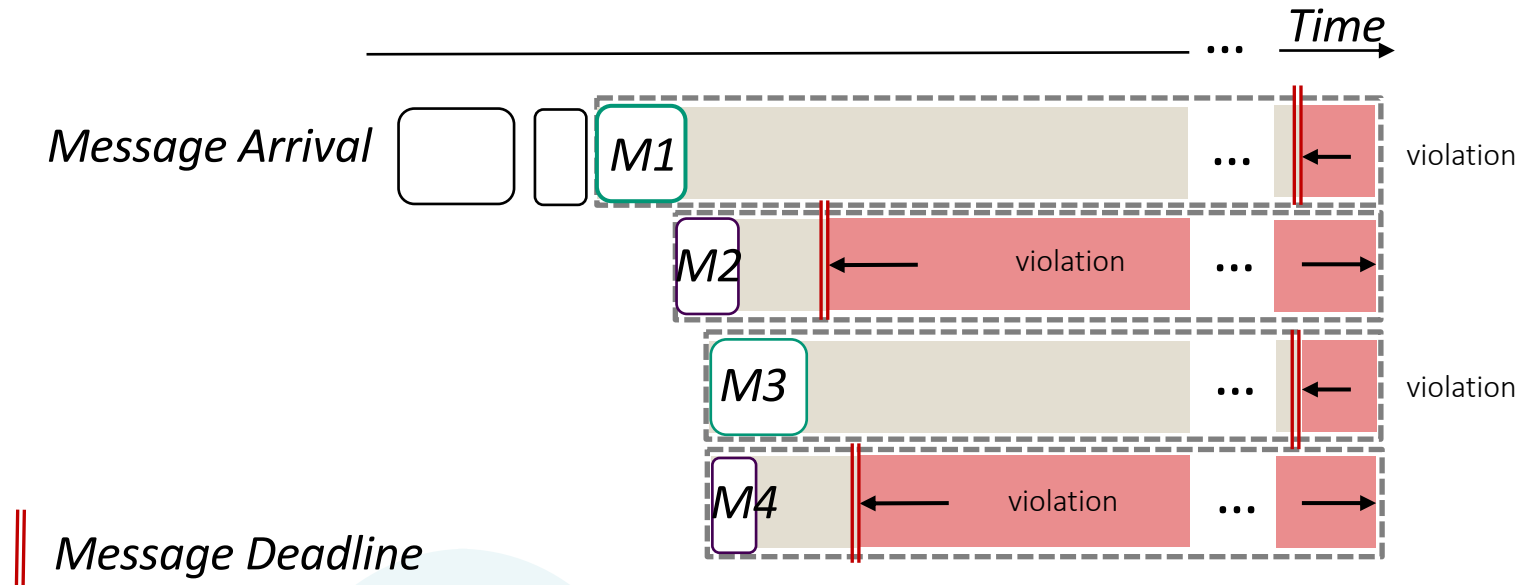
***Bulk Analytics Query***

```
.filter(e.key == "Error")
```

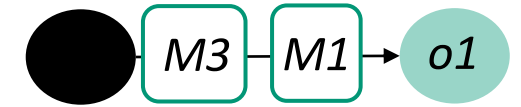


Target latency **50**

***Latency Sensitive Query***



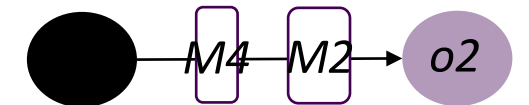
```
.window(1000)
.groupBy(e.key)
.sum()
```



Target latency **100**

**Bulk Analytics Query**

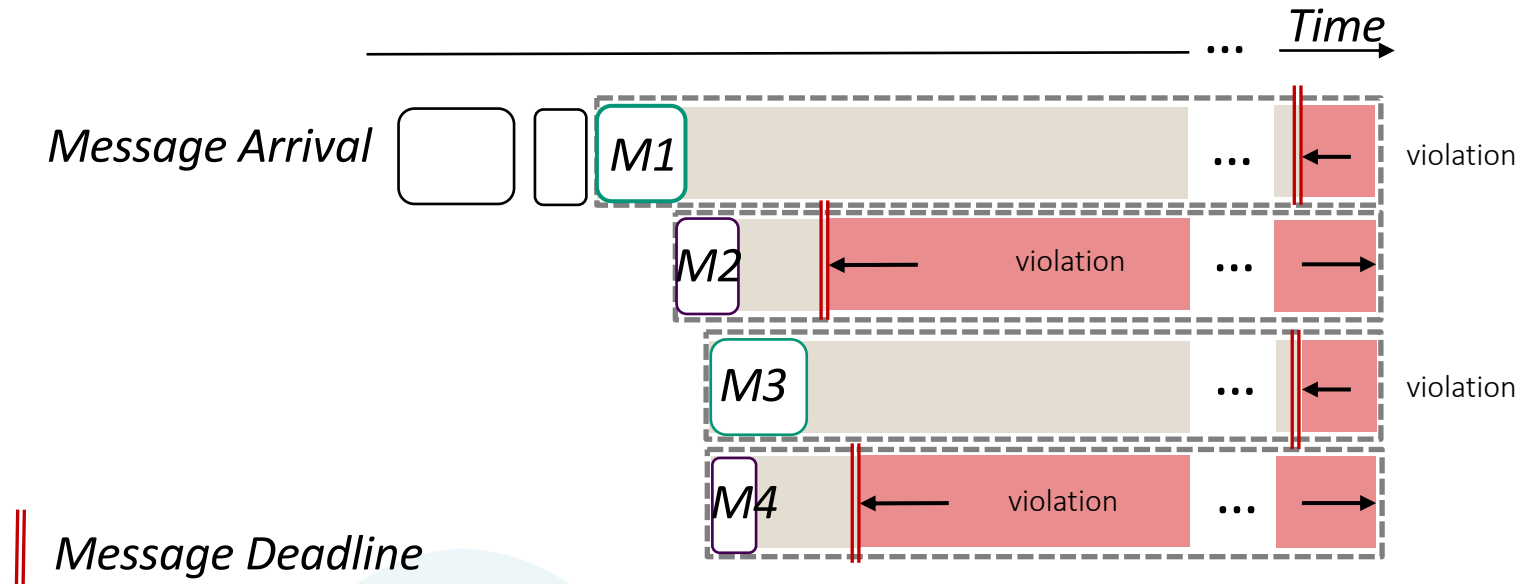
```
.filter(e.key == "Error")
```



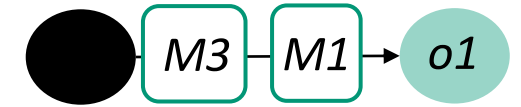
Target latency **50**

**Latency Sensitive Query**





```
.window(1000)
.groupBy(e.key)
.sum()
```



Target latency **100**

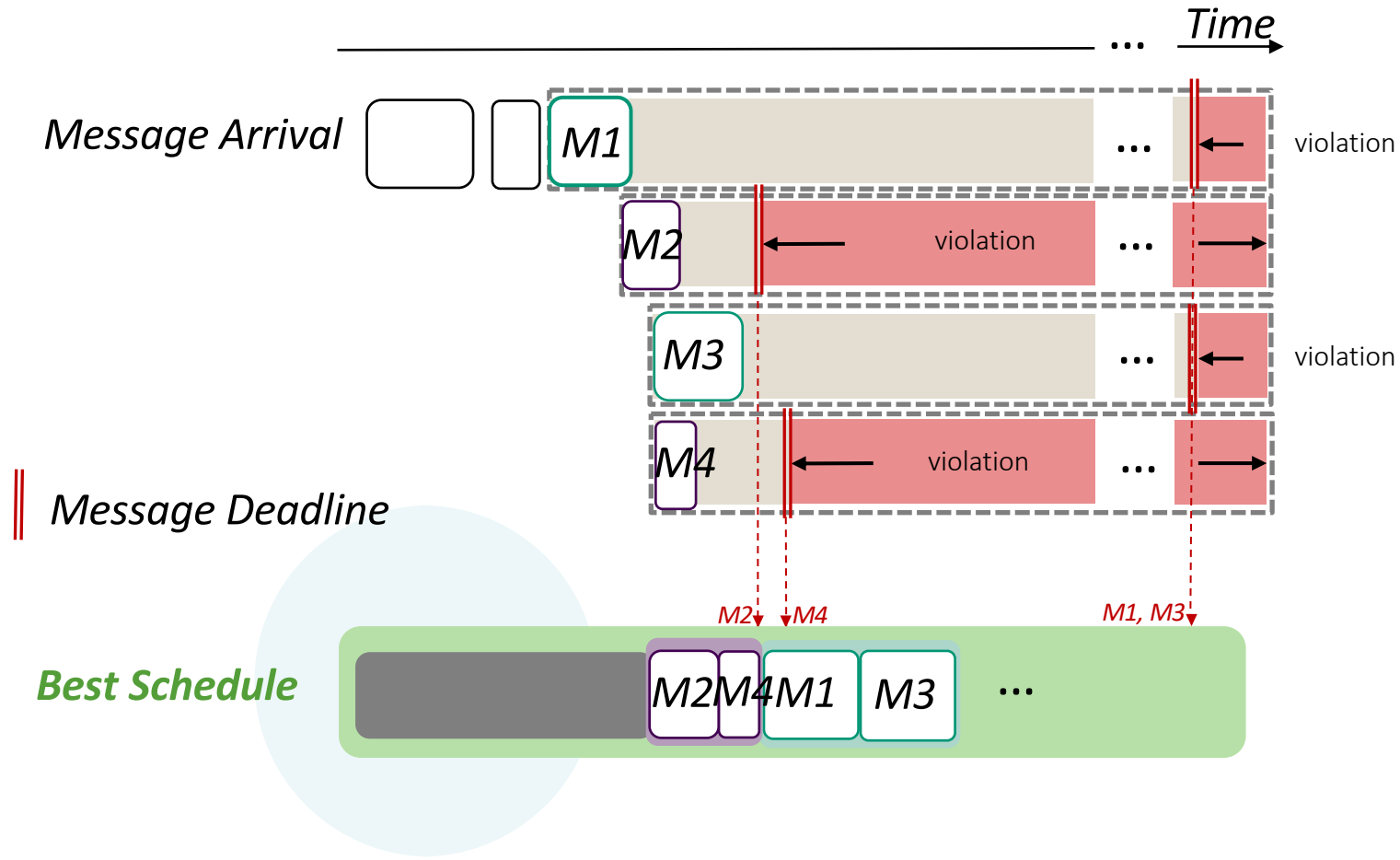
**Bulk Analytics Query**

```
.filter(e.key == "Error")
```

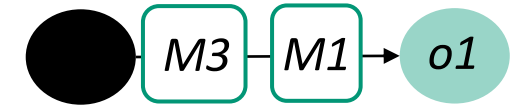


Target latency **50**

**Latency Sensitive Query**



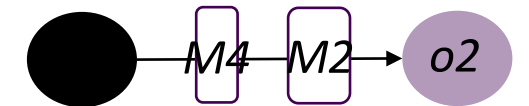
```
.window(1000)
.groupBy(e.key)
.sum()
```



Target latency **100**

**Bulk Analytics Query**

```
.filter(e.key == "Error")
```



Target latency **50**

**Latency Sensitive Query**

No  
Awareness

```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

*Bulk Analytics Query*

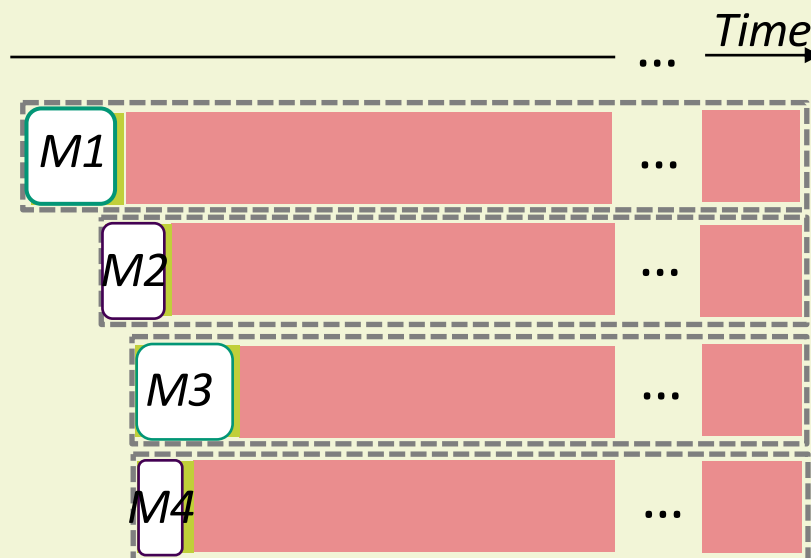
```
.filter(e.key == "Error")
```



Target latency **50**

*Latency Sensitive Query*

No  
Awareness



```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency 100

*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



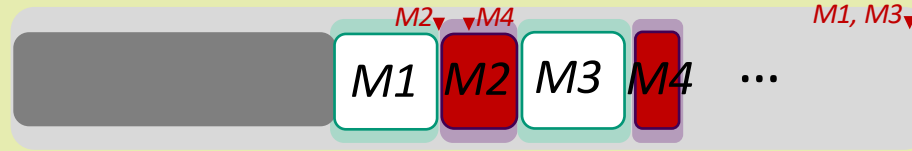
Target latency 50

*Latency Sensitive Query*

 violation

No  
Awareness

*No awareness (FIFO) schedule*



```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency 100

*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



Target latency 50

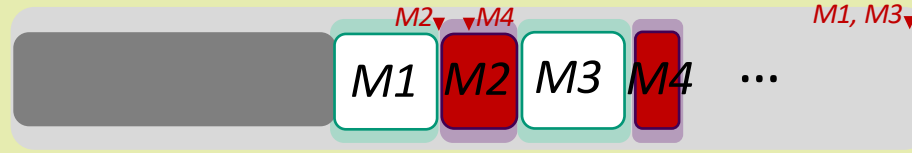
*Latency Sensitive Query*

 violation

No  
Awareness

Topology-  
aware

*No awareness (FIFO) schedule*



```
.window(1000)  
.groupBy(e.key)  
.sum()  
..
```



Target latency **100**  
*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



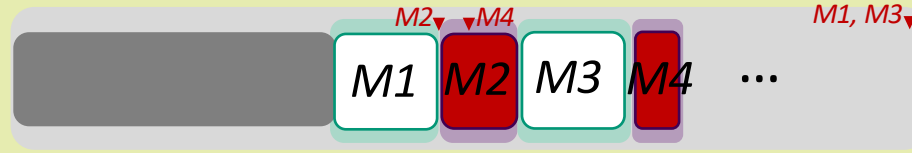
Target latency **50**  
*Latency Sensitive Query*

 **violation**

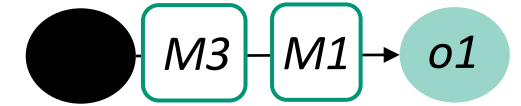
No  
Awareness

Topology-  
aware

*No awareness (FIFO) schedule*



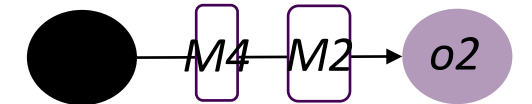
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency 100

**Bulk Analytics Query**

```
.filter(e.key == "Error")
```



Target latency 50

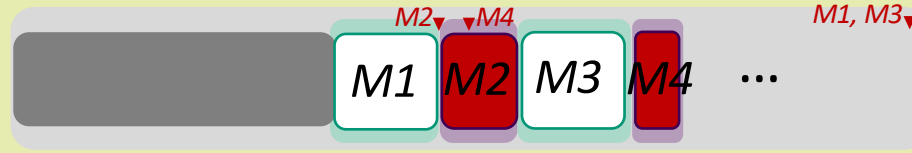
**Latency Sensitive Query**

 **violation**

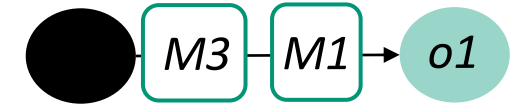
No  
Awareness

Topology-  
aware

*No awareness (FIFO) schedule*



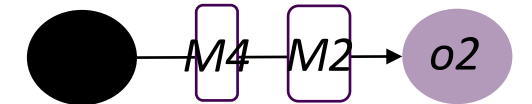
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

***Bulk Analytics Query***

```
.filter(e.key == "Error")
```



Target latency **50**

***Latency Sensitive Query***

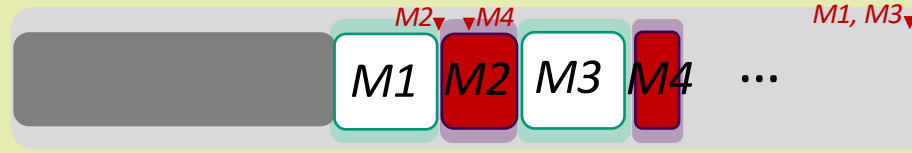
 *violation*



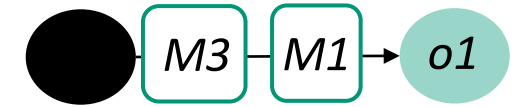
No Awareness

Topology-aware

No awareness (FIFO) schedule

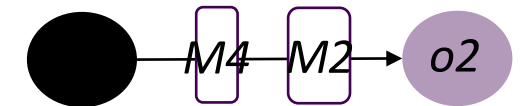


```
.window(1000)
.groupBy(e.key)
.sum()
```



Target latency **100**  
*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



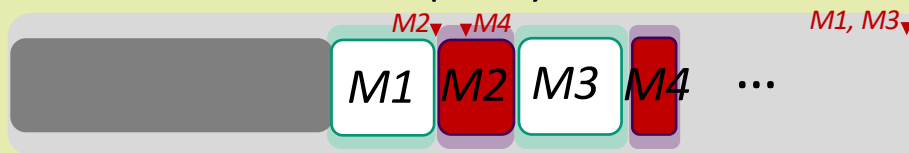
Target latency **50**  
*Latency Sensitive Query*

 violation

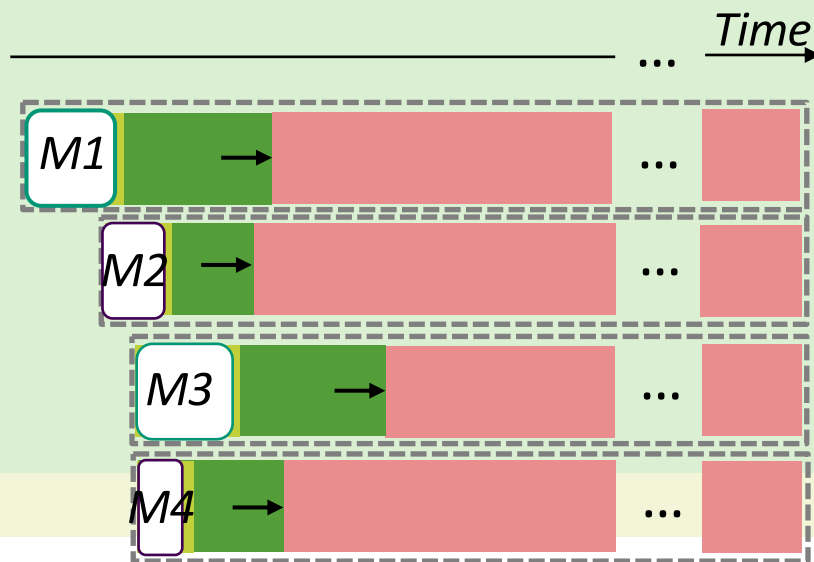
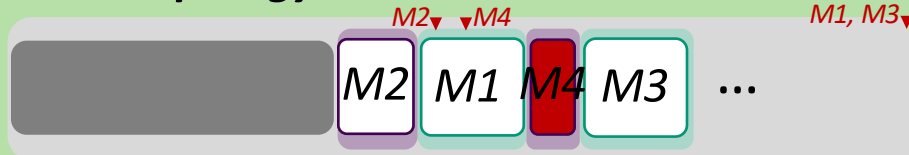
No Awareness

Topology-aware

*No awareness (FIFO) schedule*



*Topology-aware schedule*



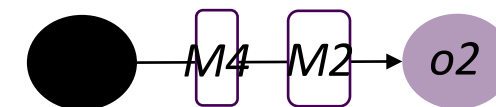
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

**Bulk Analytics Query**

```
.filter(e.key == "Error")
```



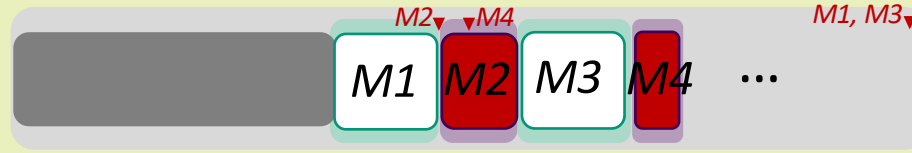
Target latency **50**

**Latency Sensitive Query**

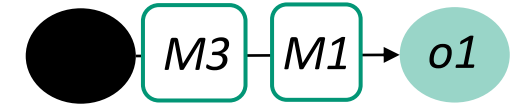
 violation

No  
Awareness

*No awareness (FIFO) schedule*



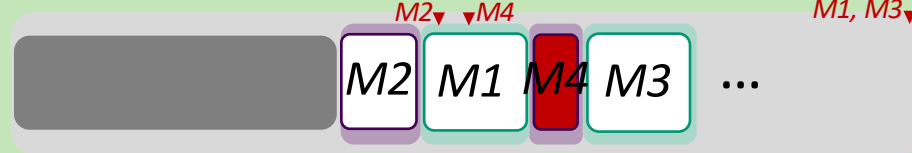
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



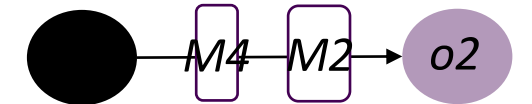
Target latency **100**  
*Bulk Analytics Query*

Topology-  
aware

*Topology-aware schedule*



```
.filter(e.key == "Error")
```



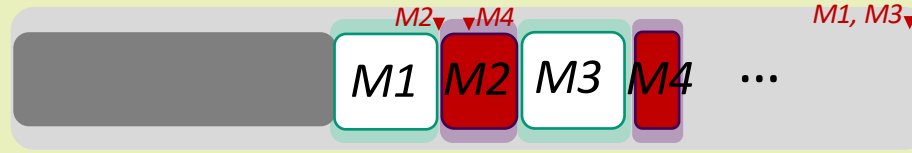
Target latency **50**  
*Latency Sensitive Query*

Semantic-  
aware

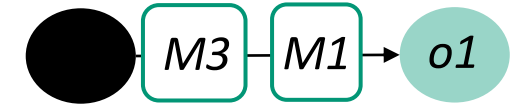
 *violation*

No  
Awareness

*No awareness (FIFO) schedule*



```
.window(1000)  
.groupBy(e.key)  
.sum()
```

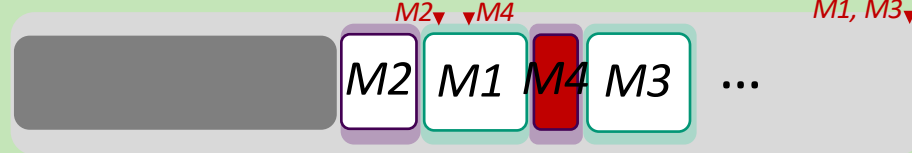


Target latency **100**

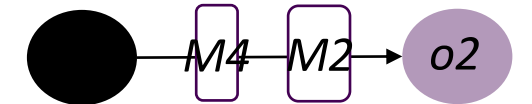
***Bulk Analytics Query***

Topology-  
aware

*Topology-aware schedule*



```
.filter(e.key == "Error")
```



Target latency **50**

***Latency Sensitive Query***

Semantic-  
aware

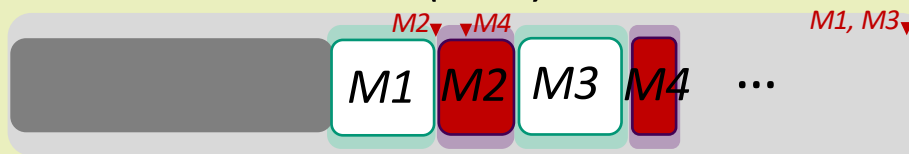
 *violation*

No  
Awareness

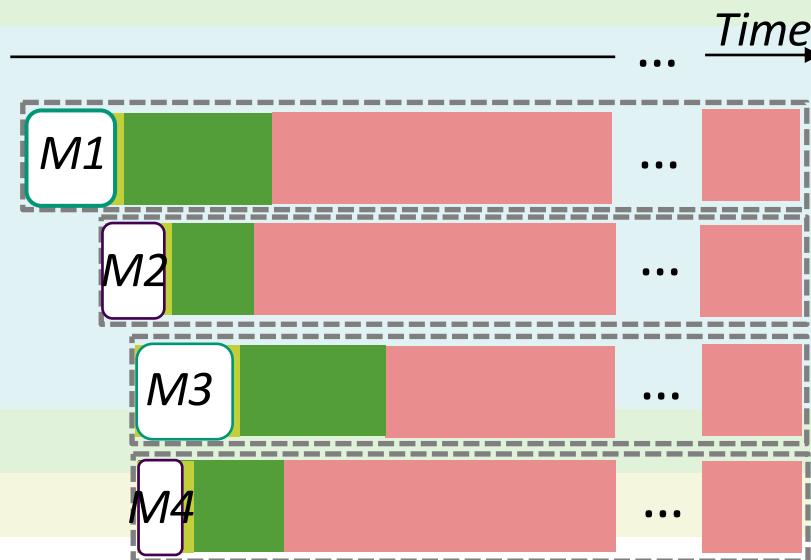
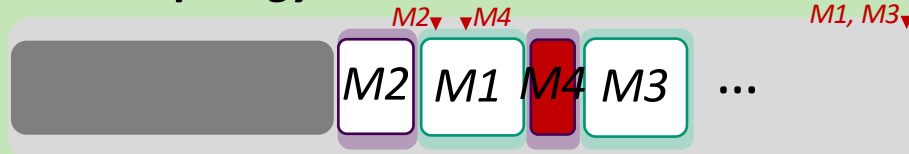
Topology-  
aware

Semantic-  
aware

*No awareness (FIFO) schedule*



*Topology-aware schedule*



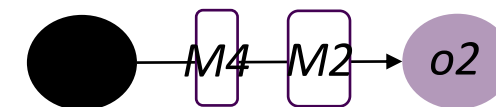
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



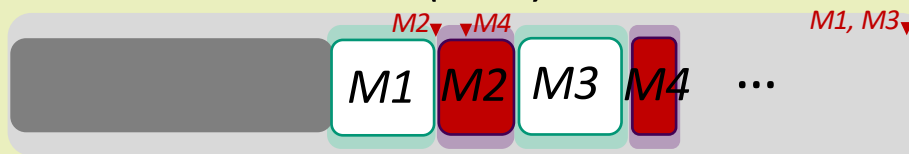
Target latency **50**

*Latency Sensitive Query*

 *violation*

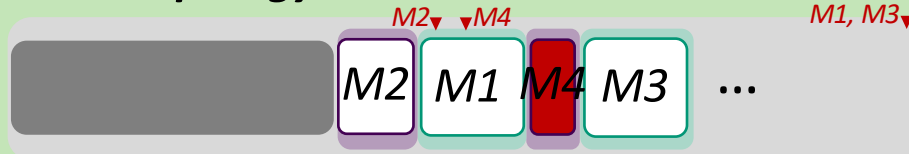
No  
Awareness

*No awareness (FIFO) schedule*

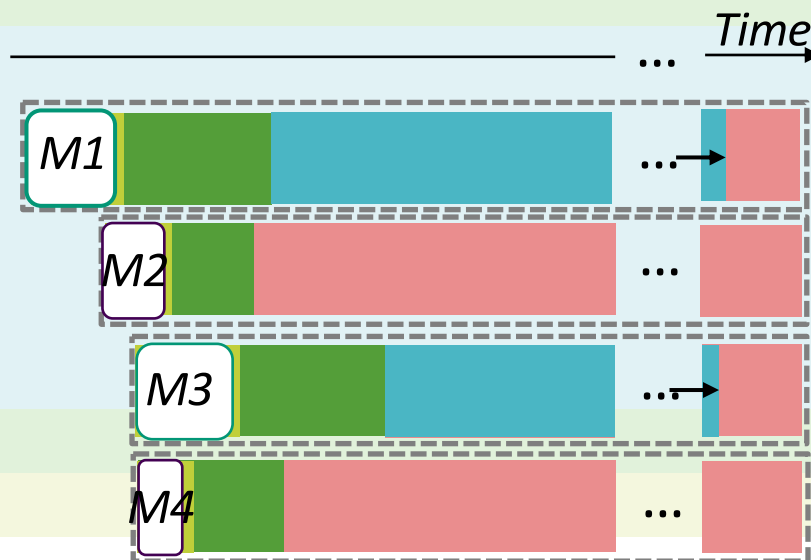


Topology-  
aware

*Topology-aware schedule*



Semantic-  
aware



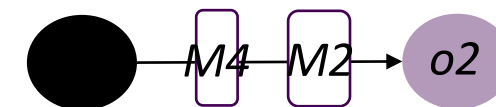
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



Target latency **100**

*Bulk Analytics Query*

```
.filter(e.key == "Error")
```



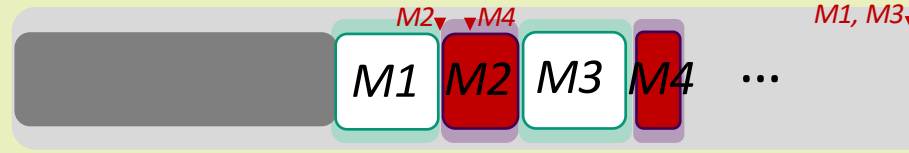
Target latency **50**

*Latency Sensitive Query*

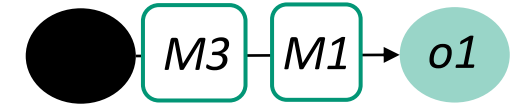
 *violation*

No  
Awareness

*No awareness (FIFO) schedule*



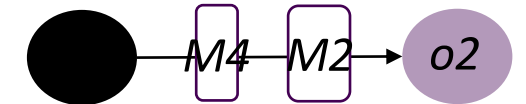
```
.window(1000)
.groupBy(e.key)
.sum()
```



Target latency **100**

**Bulk Analytics Query**

```
.filter(e.key == "Error")
```



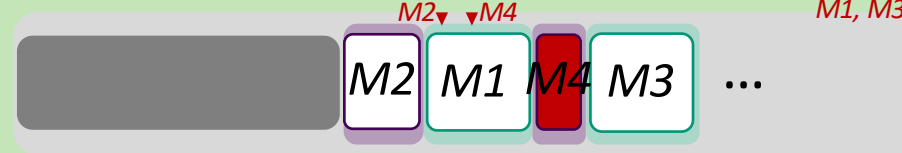
Target latency **50**

**Latency Sensitive Query**

 violation

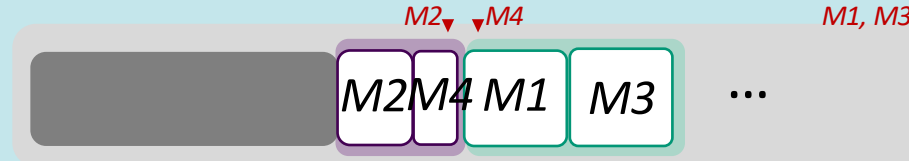
Topology-  
aware

*Topology-aware schedule*



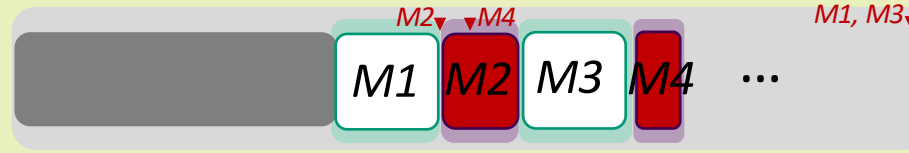
Semantic-  
aware

*Topology and Semantic-aware schedule*

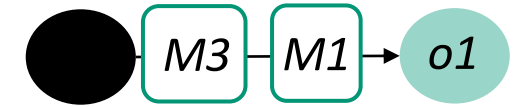


No  
Awareness

*No awareness (FIFO) schedule*



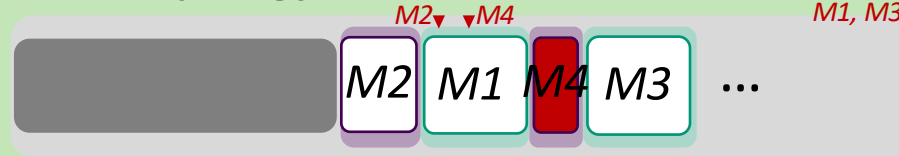
```
.window(1000)  
.groupBy(e.key)  
.sum()
```



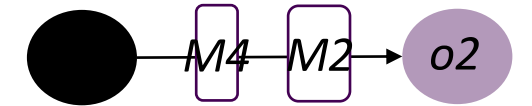
Target latency **100**  
*Bulk Analytics Query*

Topology-  
aware

*Topology-aware schedule*



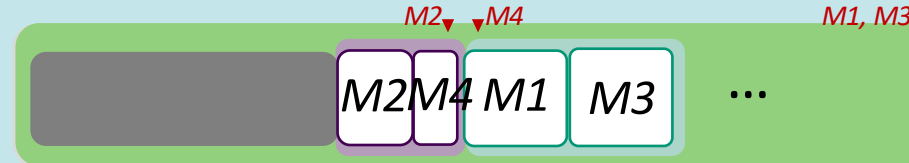
```
.filter(e.key == "Error")
```



Target latency **50**  
*Latency Sensitive Query*

Semantic-  
aware

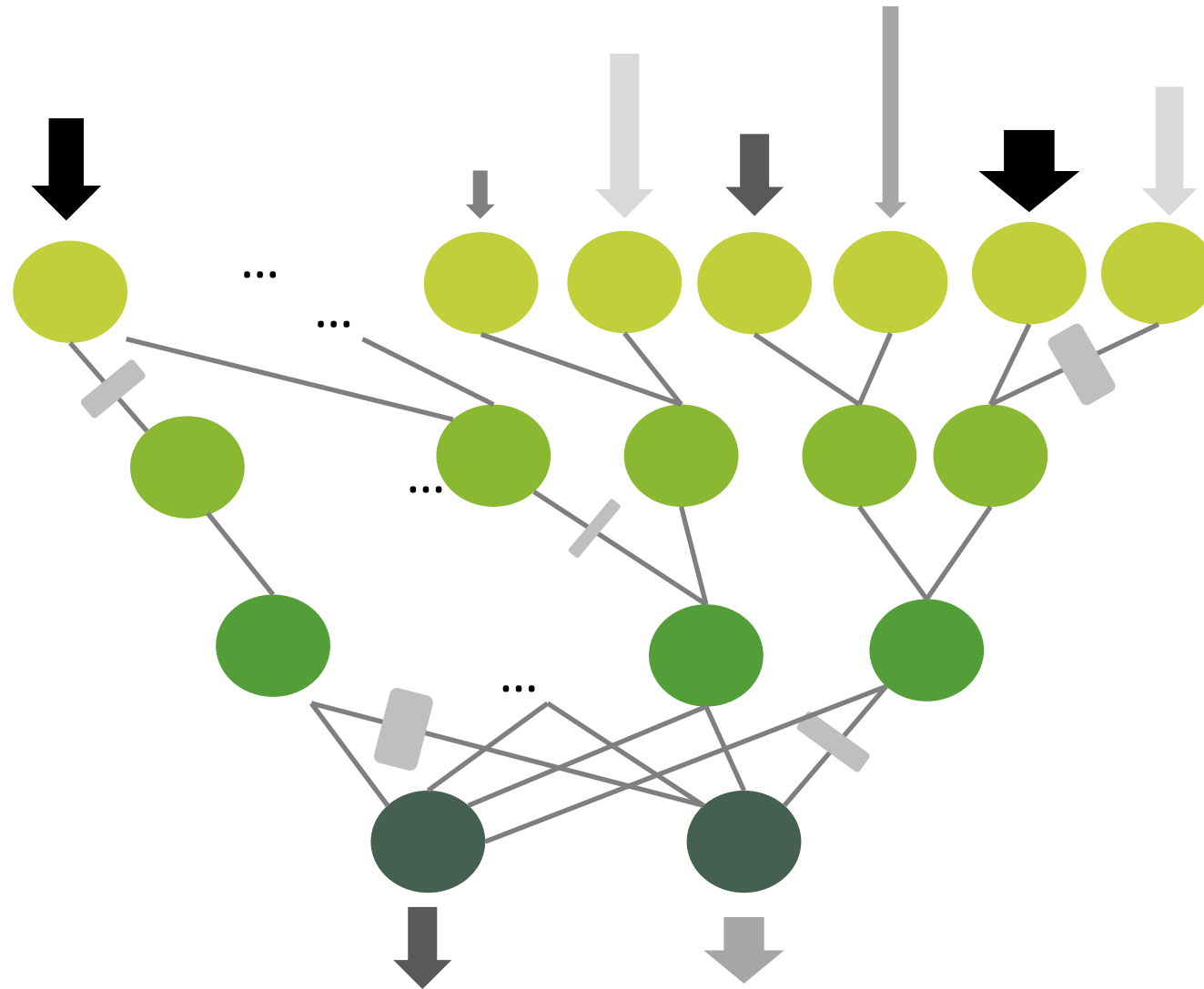
*Topology and Semantic-aware schedule*



*Best Schedule*

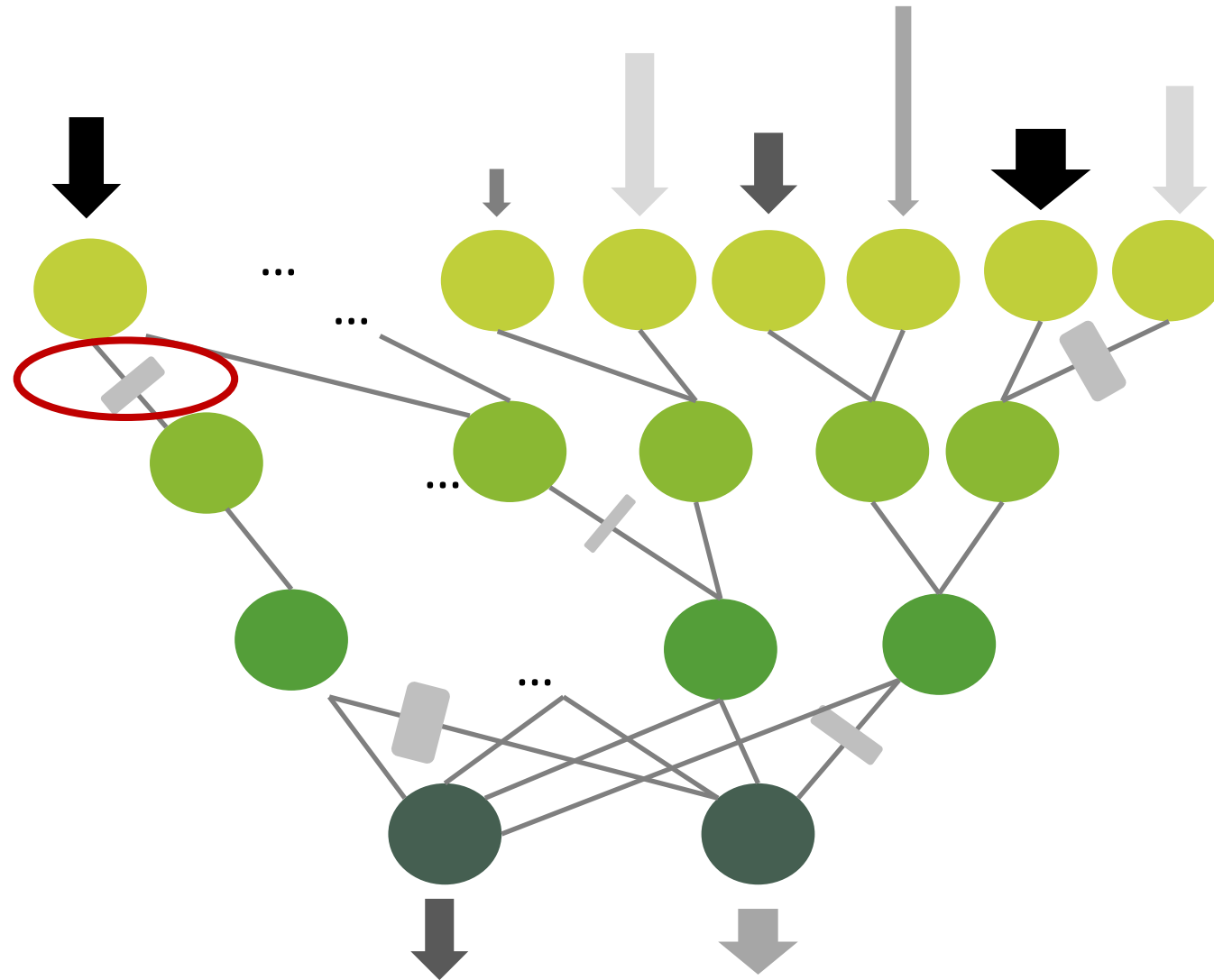
 violation





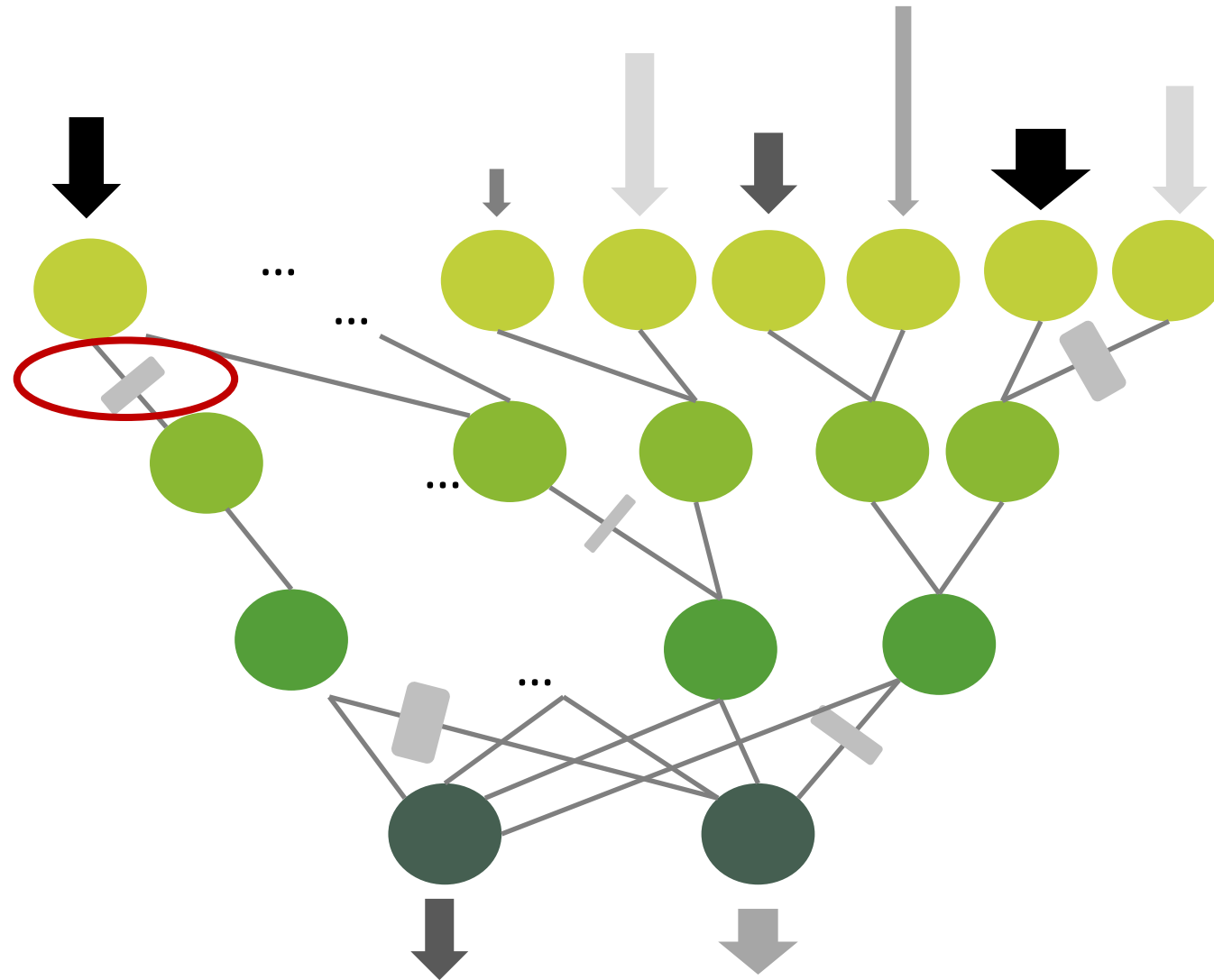
Message Priority (Deadline)

$$\text{ddl}_M = t_{MF} + L - C_{OM} - C_{path}$$



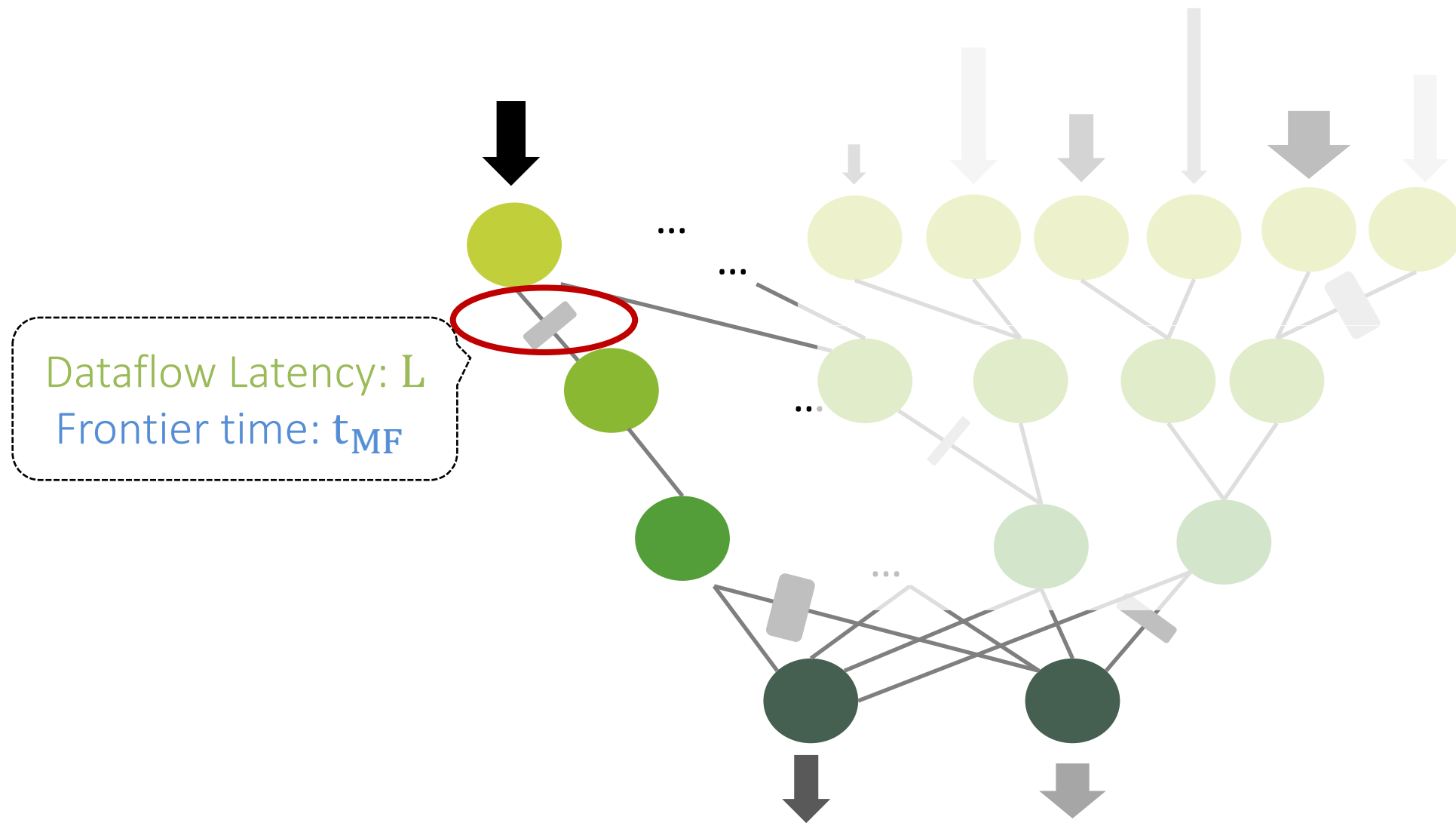
Message Priority (Deadline)

$$ddl_M = t_{MF} + L - C_{OM} - C_{path}$$



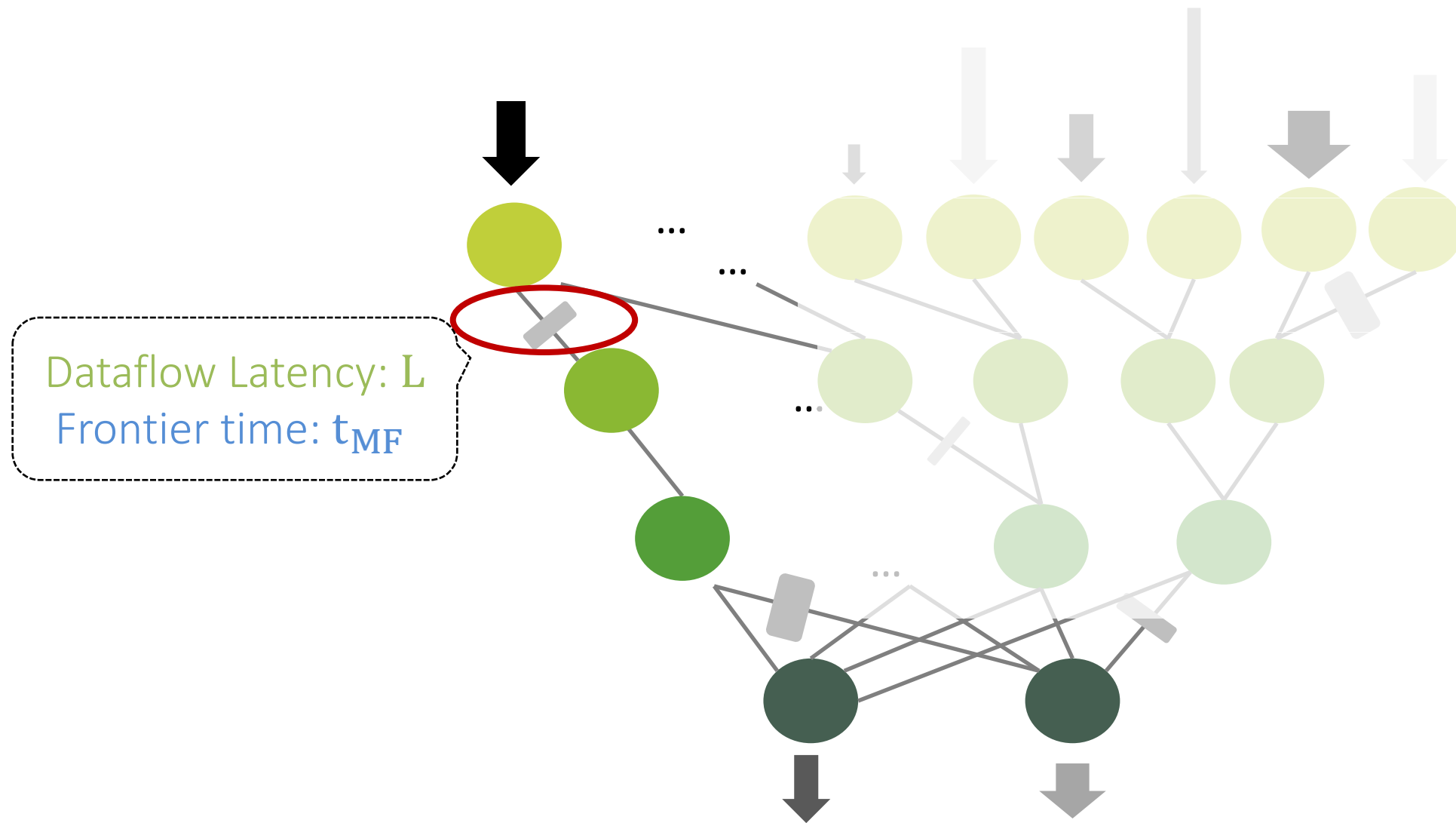
Message Priority (Deadline)

$$ddl_M = t_{MF} + L - C_{OM} - C_{path}$$



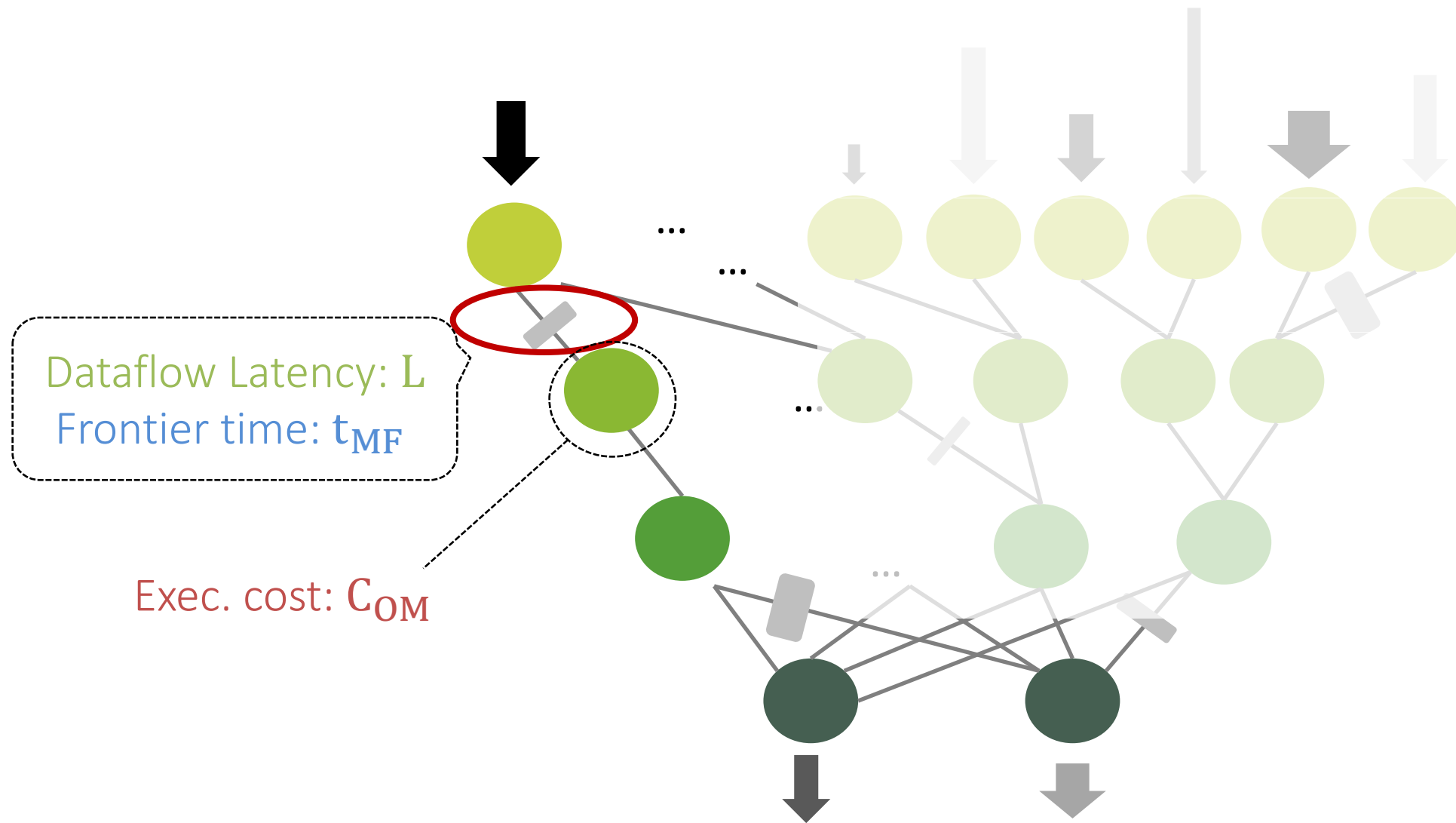
Message Priority (Deadline)

$$ddl_M = t_{MF} + L - c_{OM} - c_{path}$$



Message Priority (Deadline)

$$ddl_M = t_{MF} + L - C_{OM} - C_{path}$$

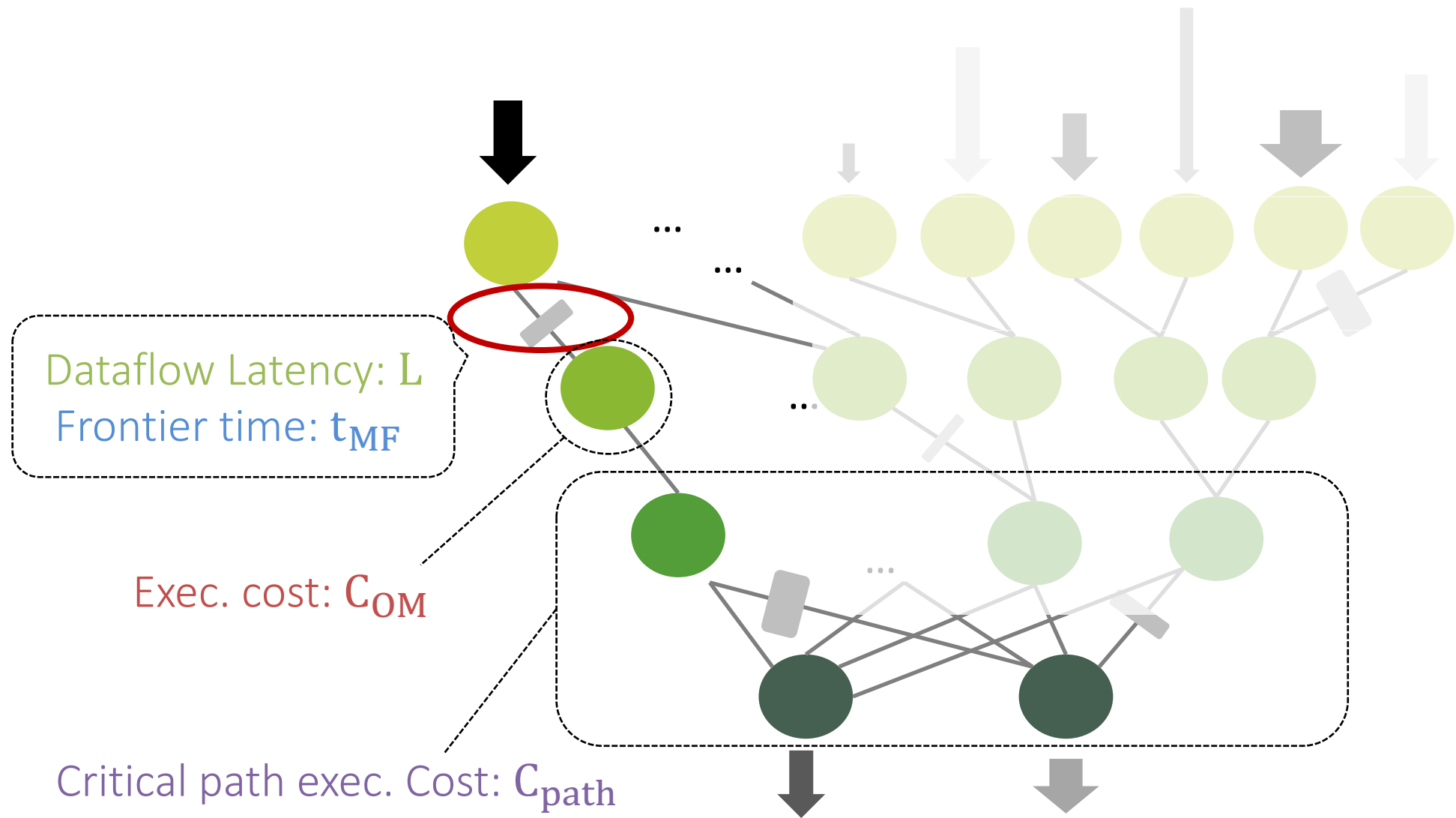


Dataflow Latency:  $L$   
Frontier time:  $t_{MF}$

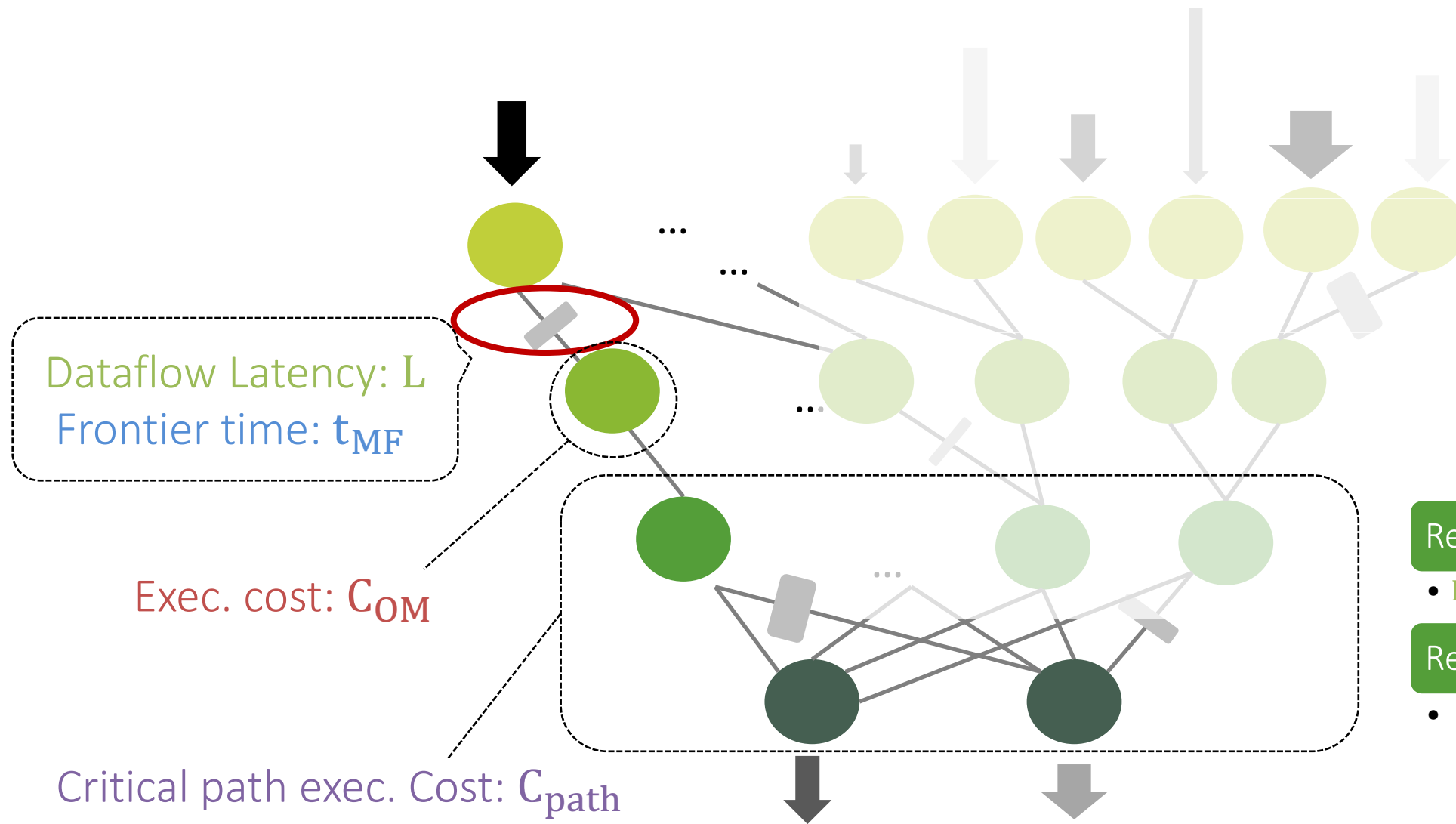
Exec. cost:  $C_{OM}$

Message Priority (Deadline)

$$ddl_M = t_{MF} + L - C_{OM} - C_{path}$$



$$ddl_M = t_{MF} + L - C_{OM} - C_{path}$$



Requires Static Information

- $L$ ,  $t_{MF}$

Requires Dynamic Information

- $t_{MF}$ ,  $C_{OM}$ ,  $C_{path}$



## Cameo Strategy

- Topology and semantics aware
- Driven by message deadline

## Cameo Mechanism

- Light-weight, stateless scheduler
- Scalable priority generation
- Pluggable strategy

**Cameo**  
**Data-driven**, fine-grained  
operator scheduling

## Cameo Strategy

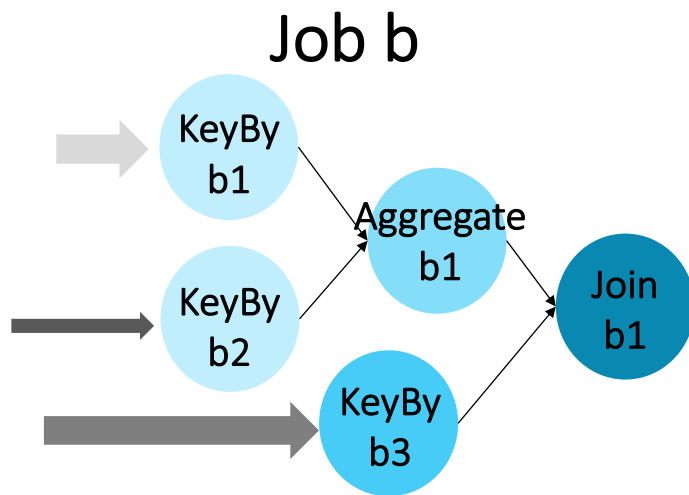
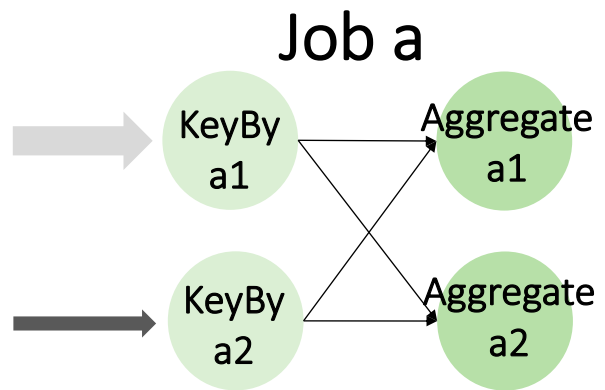
- Topology and semantics aware
- Driven by message deadline

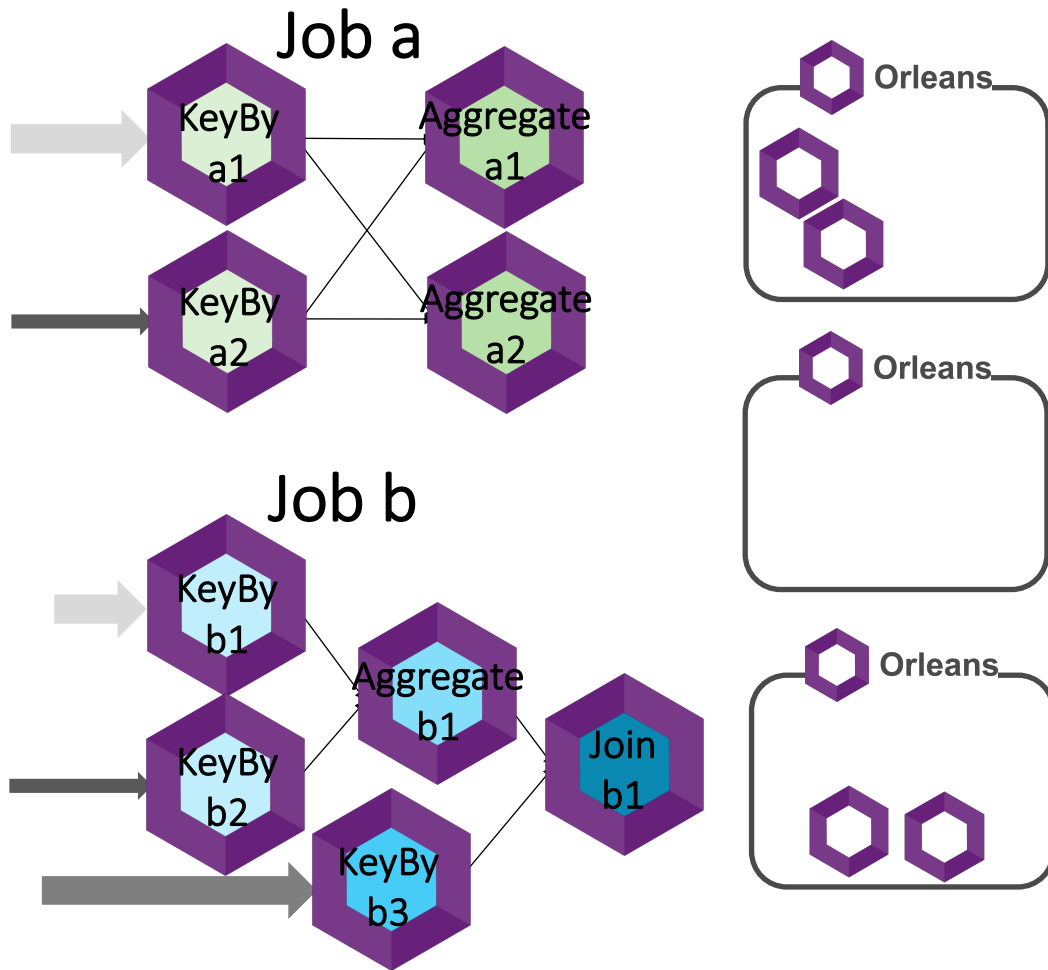
## Cameo Mechanism

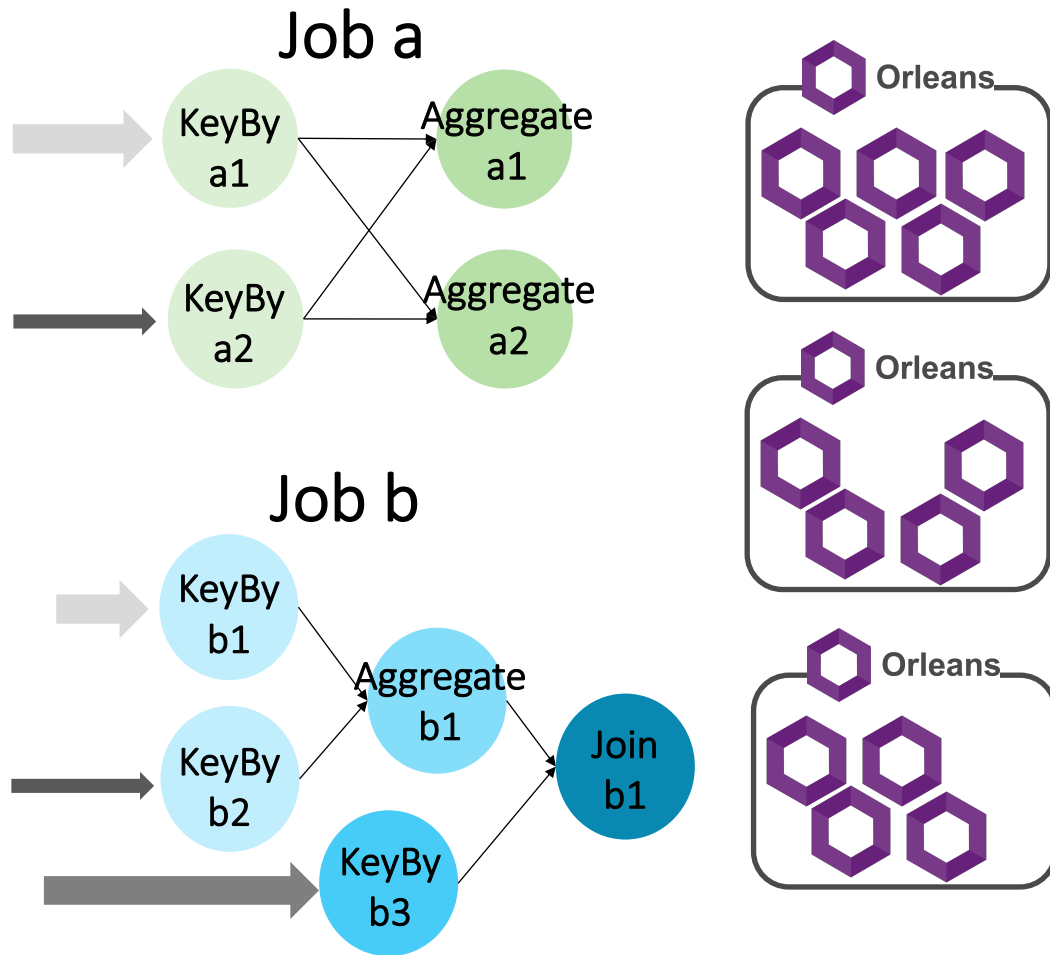
- Light-weight, stateless scheduler
- Scalable priority generation
- Pluggable strategy

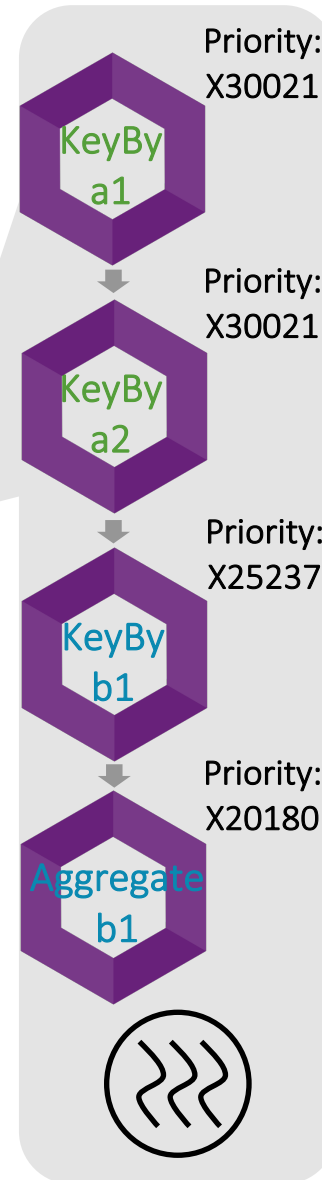
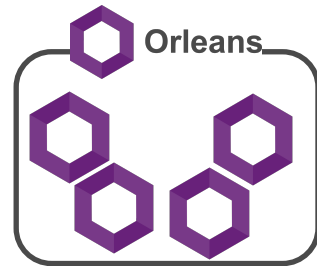
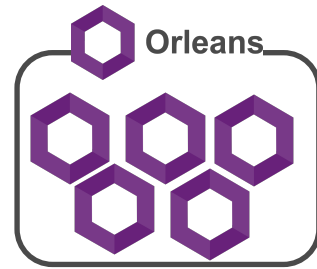
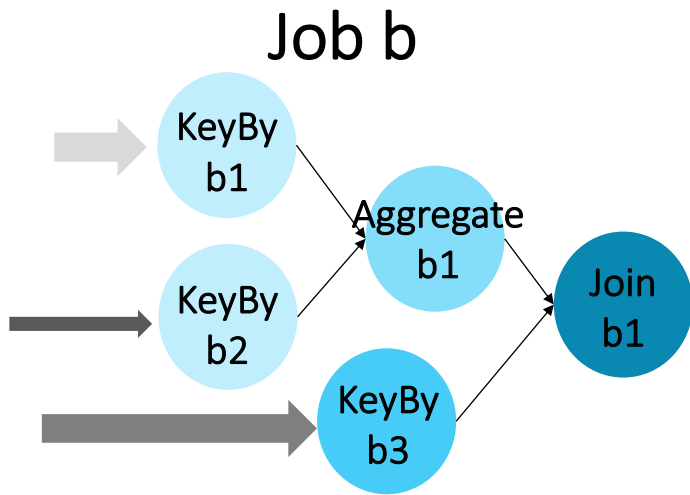
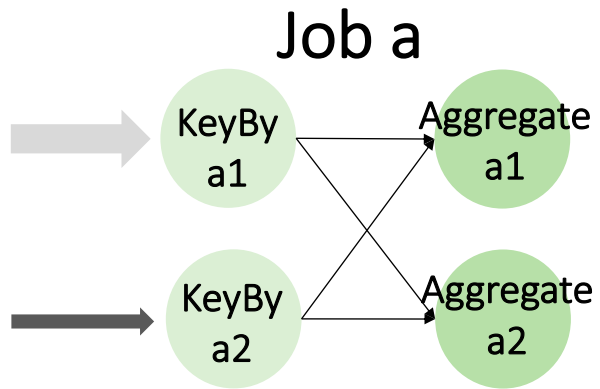
**Cameo**  
Data-driven, **fine-grained**  
operator scheduling

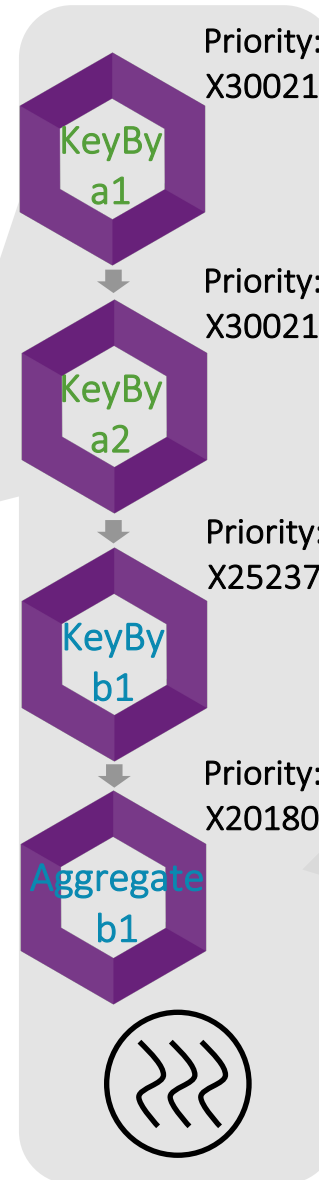
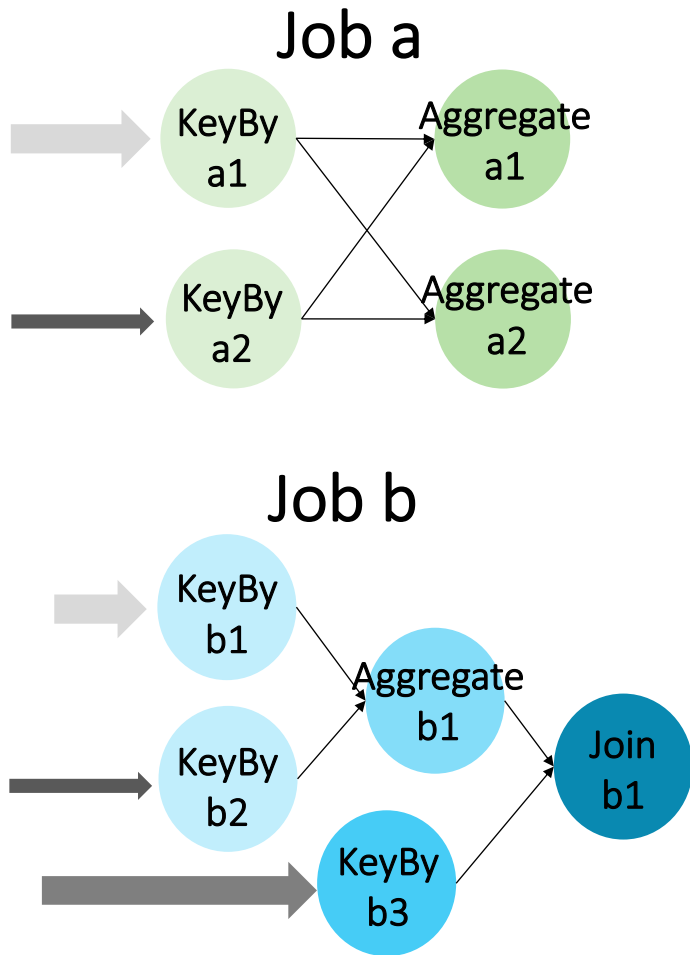
# Cameo Mechanism





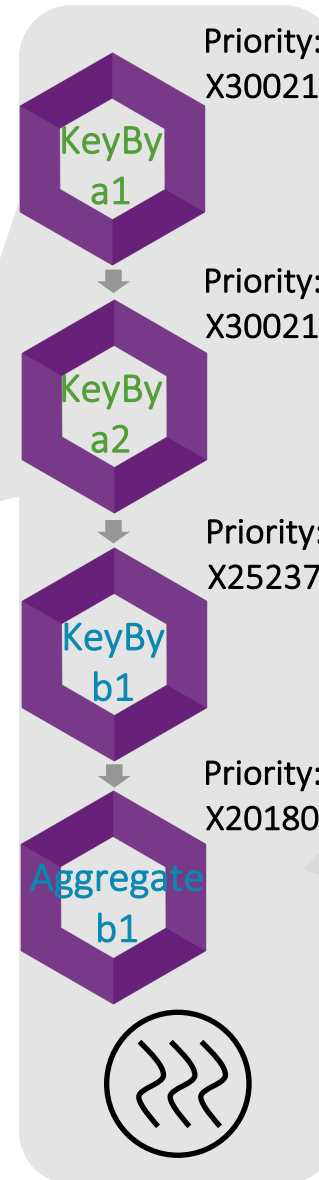
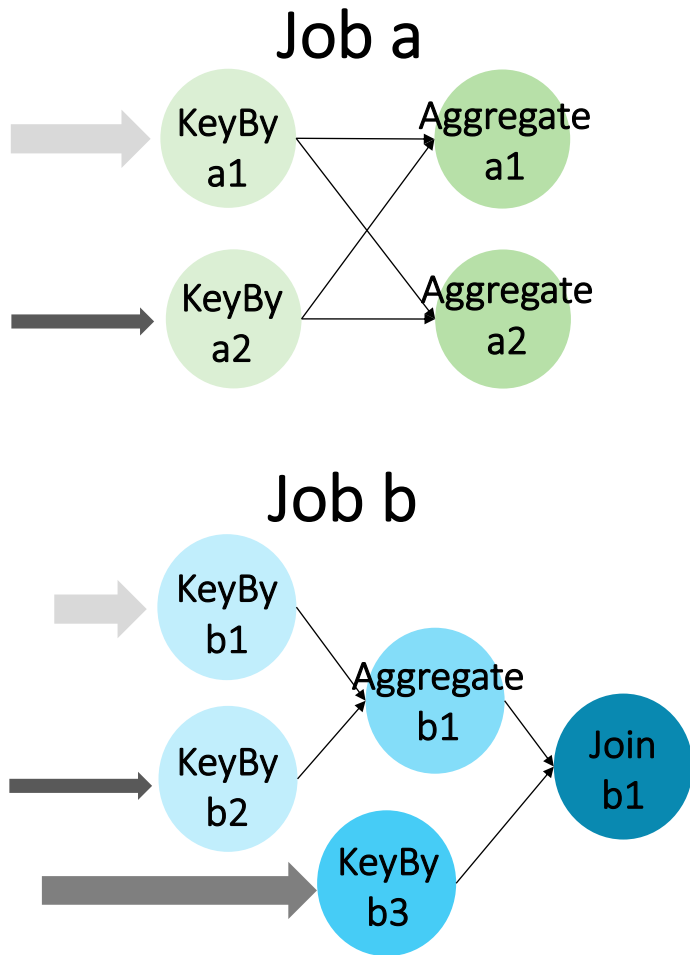






Global Priority	X20180	X30176
Local Priority	X20000	X30000
Pending Messages	Message<RID: 7216>, Message<RID: 7217>	Message<RID: 7218>, Message<RID: 7219>

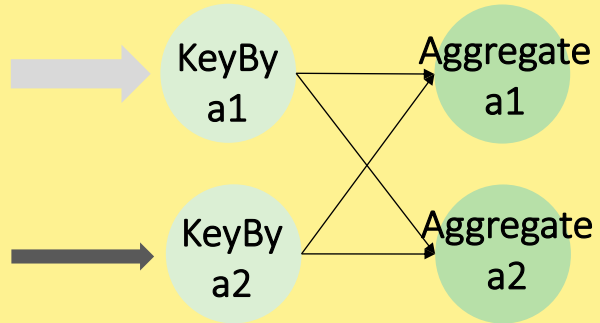




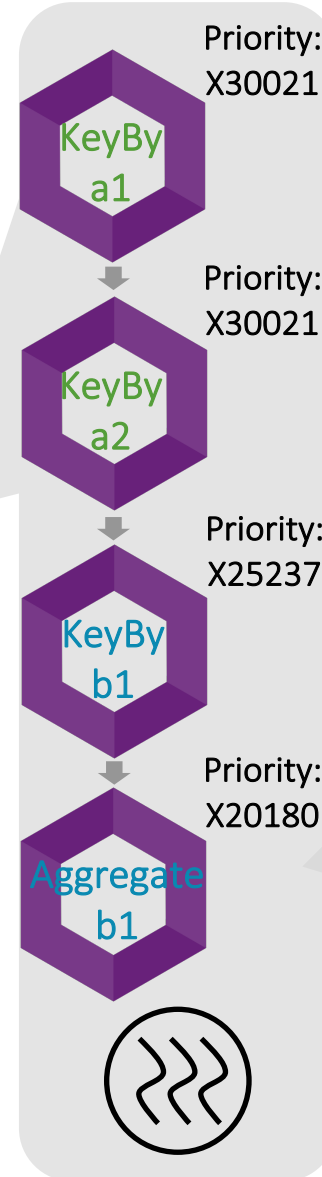
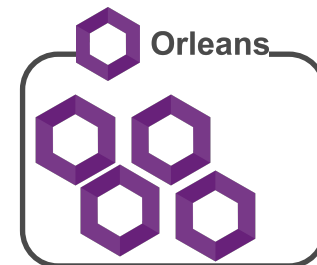
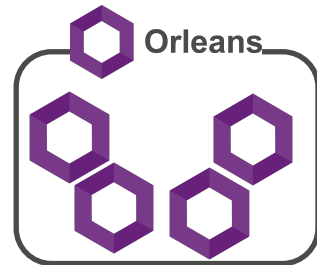
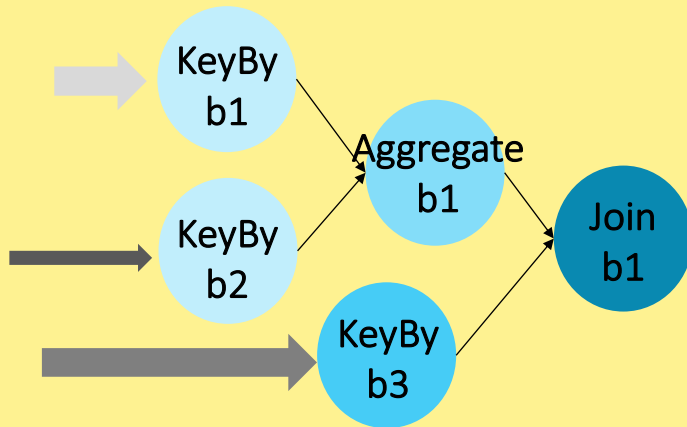
Global Priority	X20180	X30176
Local Priority	X20000	X30000
Pending Messages	Message<RID: 7216>, Message<RID: 7217>	Message<RID: 7218>, Message<RID: 7219>

## Streaming Runtime

### Job a



### Job b



Global Priority

X20180

X30176

Local Priority

X20000

X30000

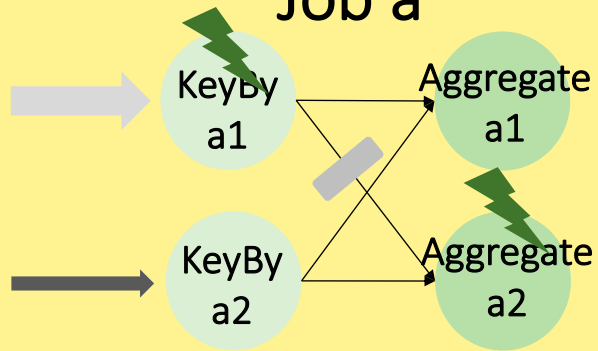
Pending Messages

Message<RID: 7216>,  
Message<RID: 7217>

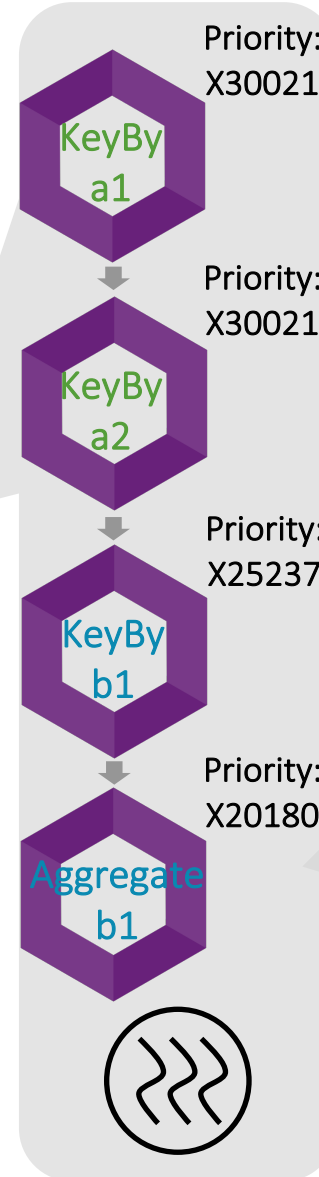
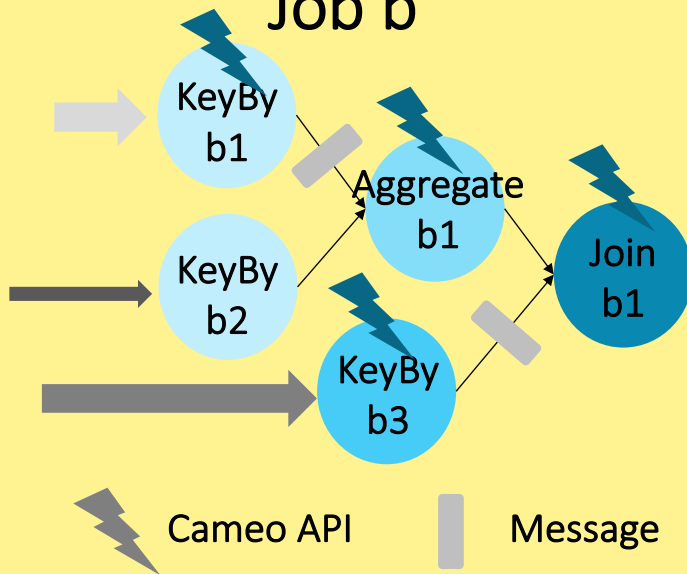
Message<RID: 7218>,  
Message<RID: 7219>

## Streaming Runtime

### Job a



### Job b



Global Priority

X20180

X30176

Local Priority

X20000

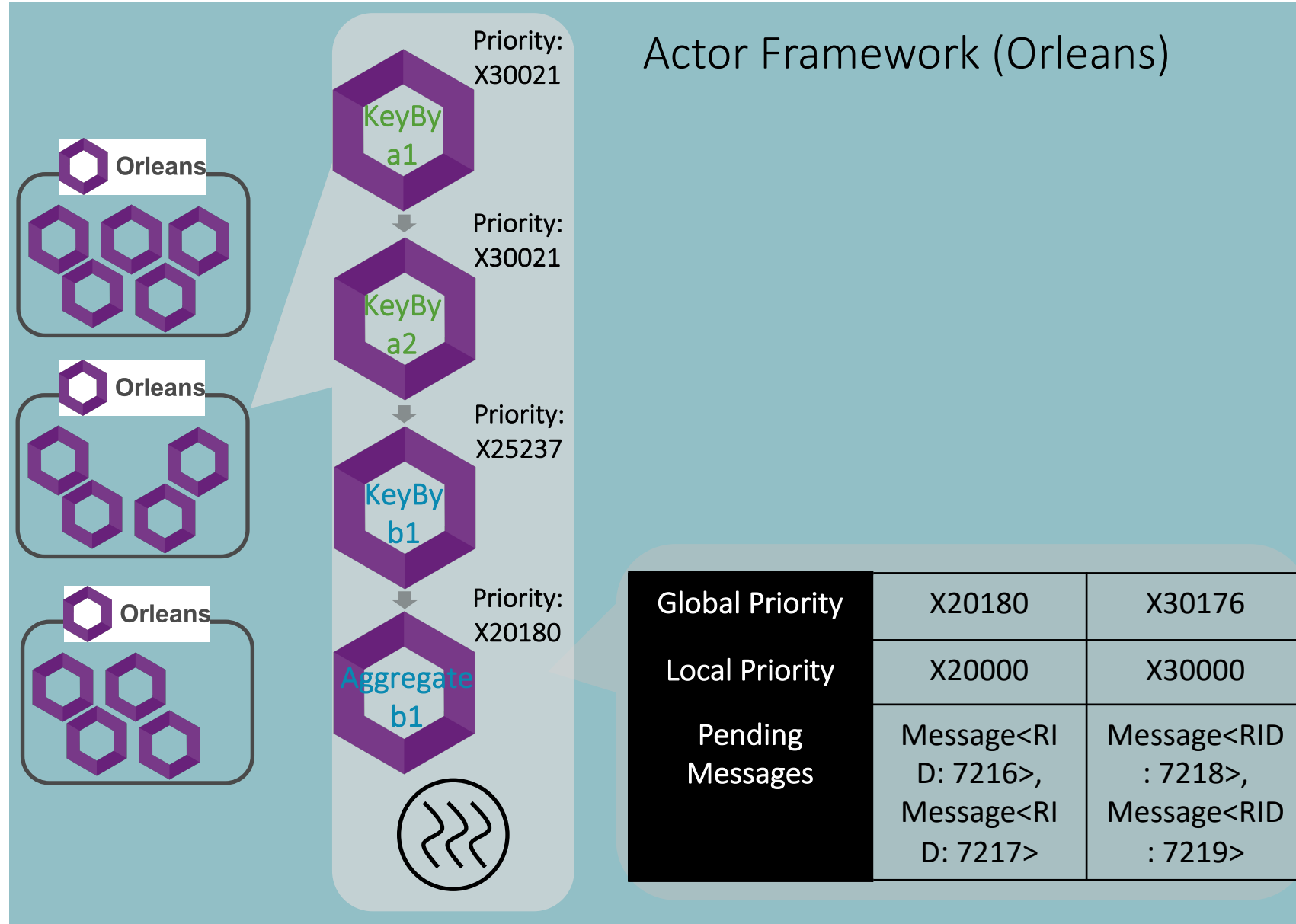
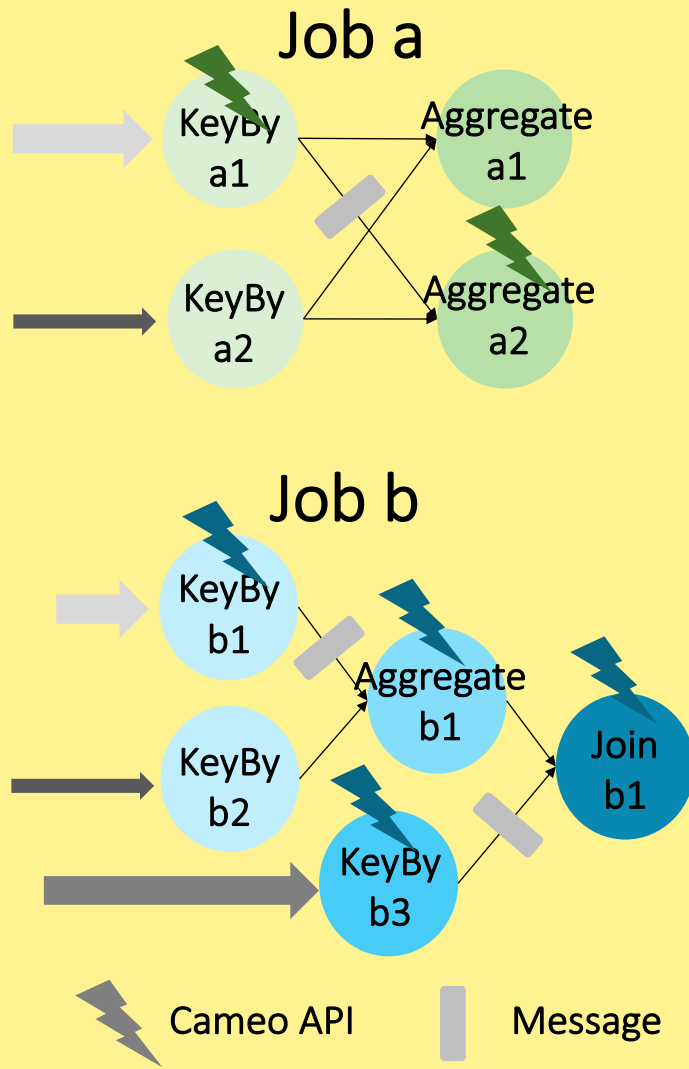
X30000

Pending Messages

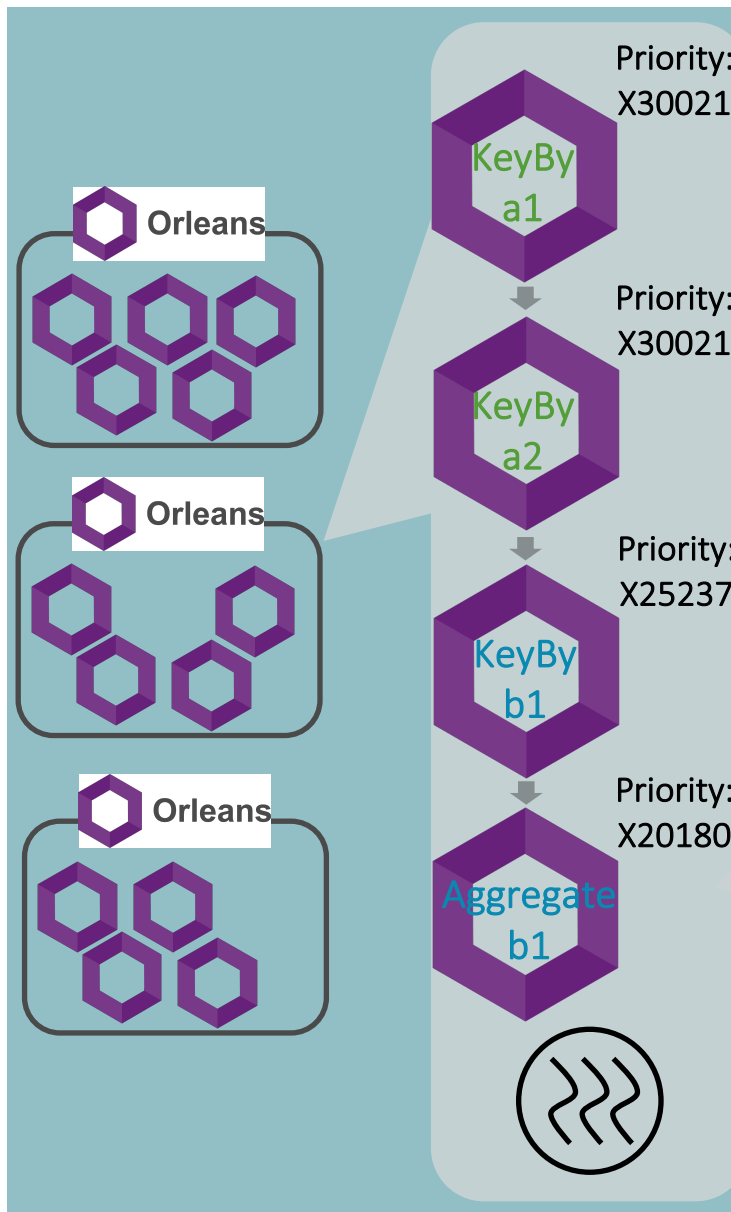
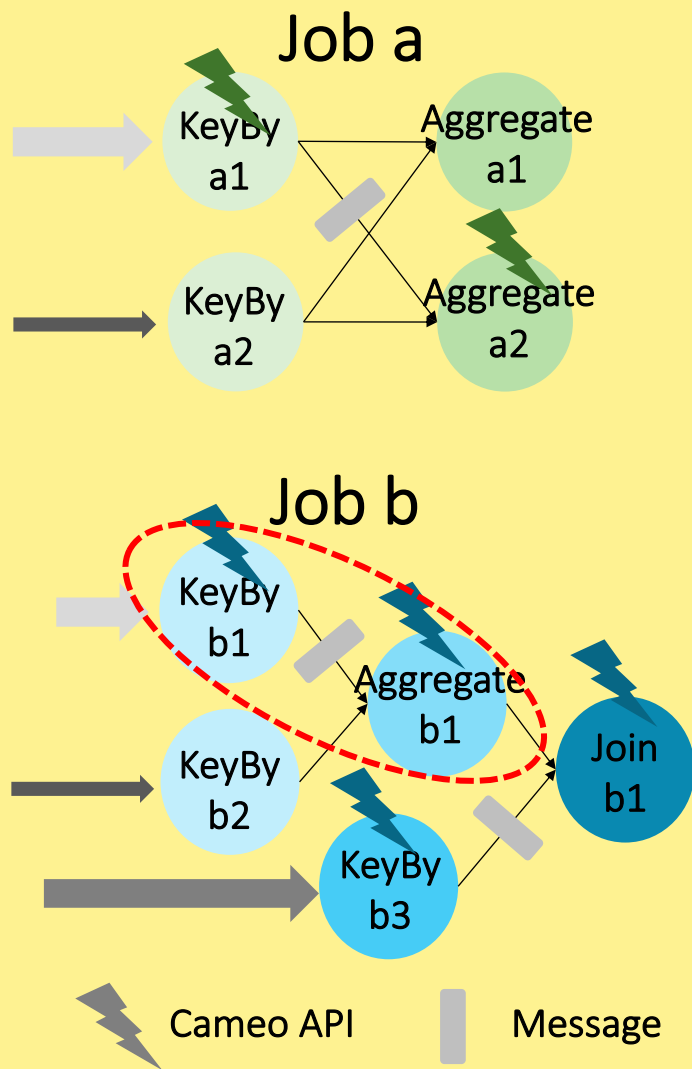
Message<RID:  
7216>,  
Message<RID:  
7217>

Message<RID:  
7218>,  
Message<RID:  
7219>

## Streaming Runtime

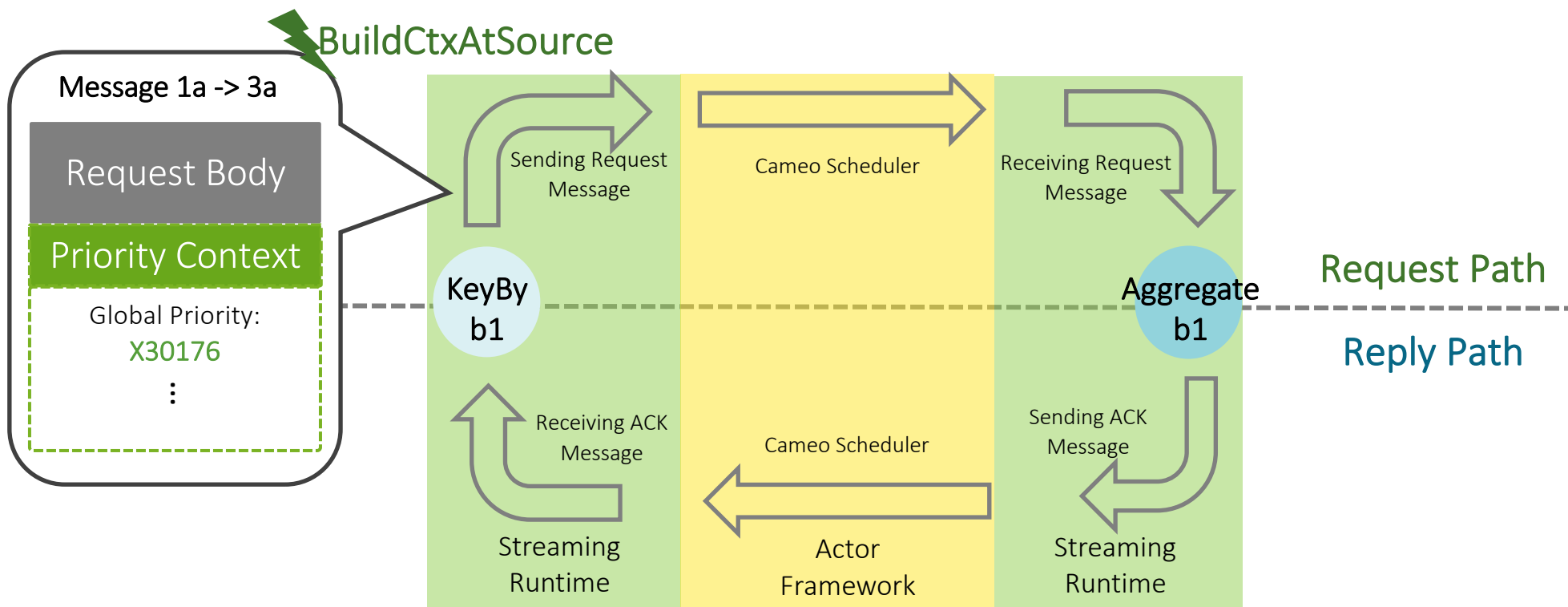


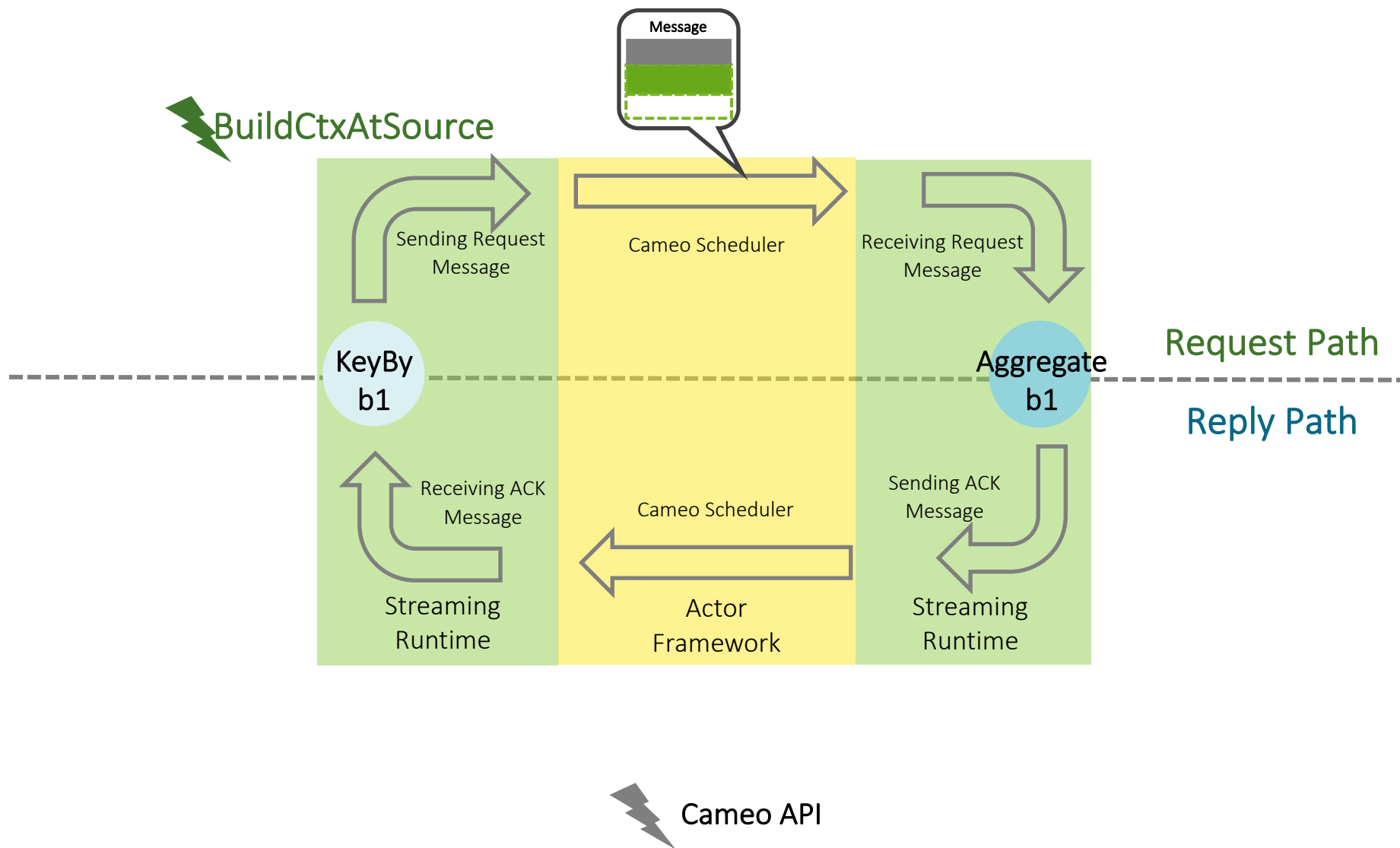
## Streaming Runtime

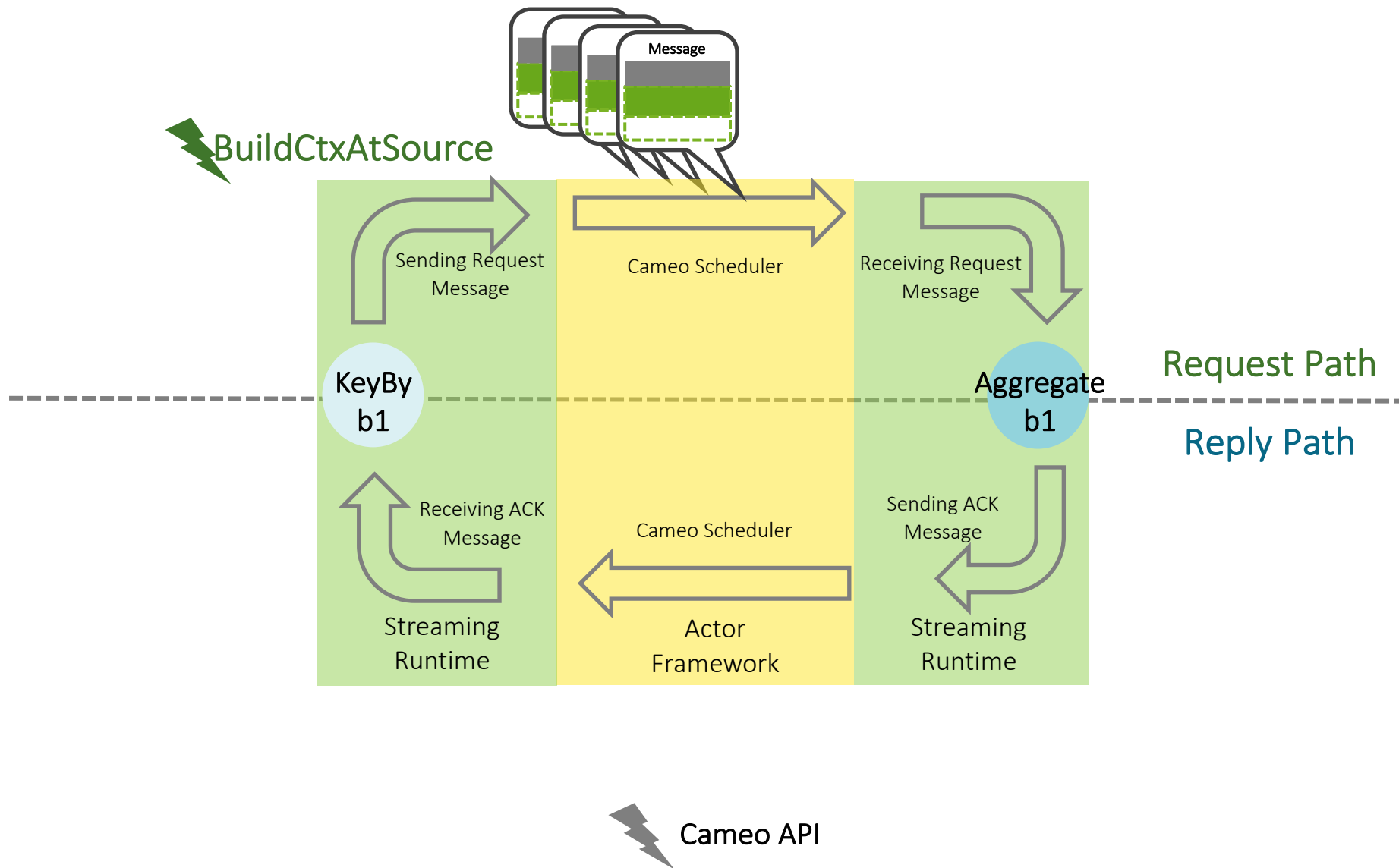


## Actor Framework (Orleans)

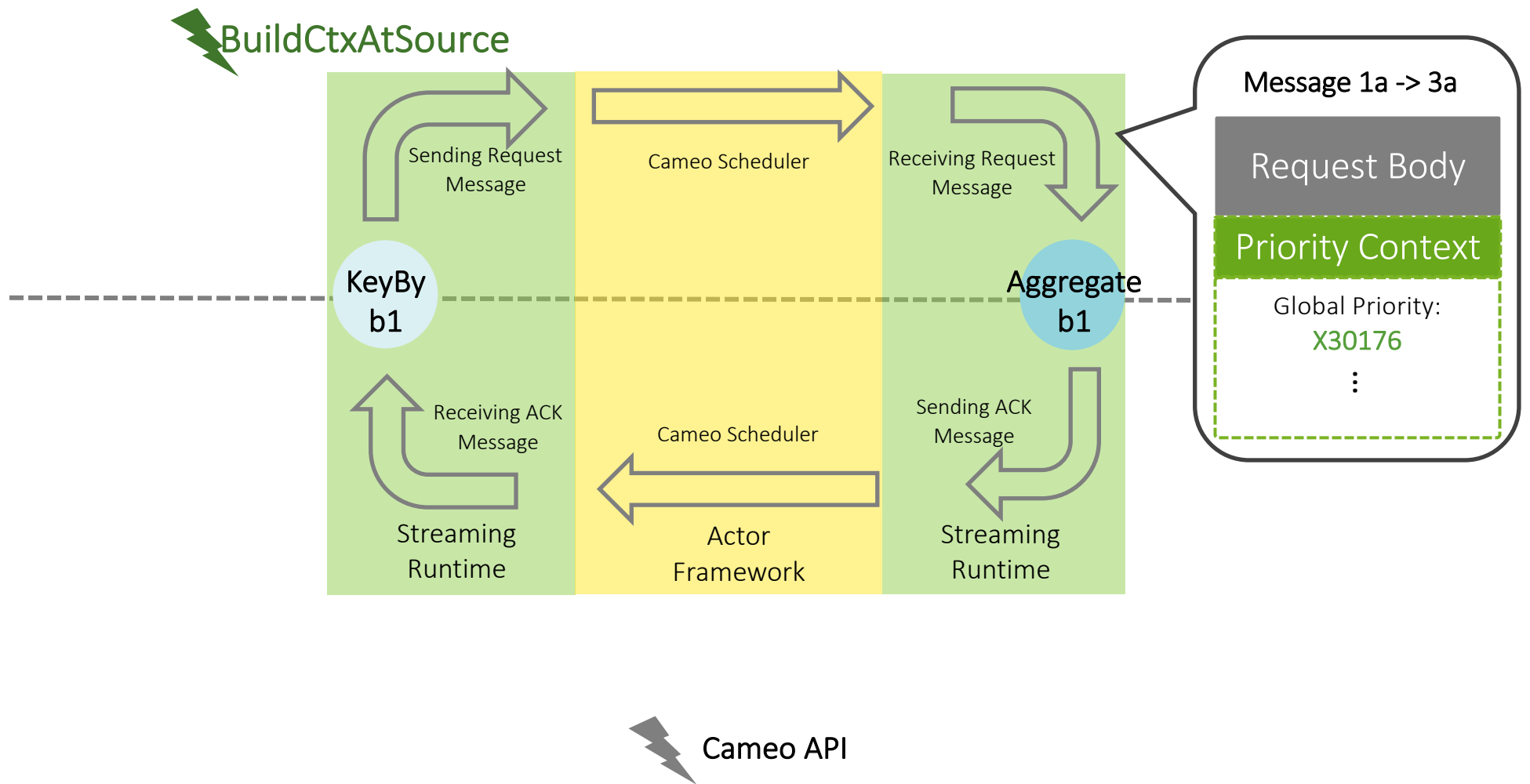
Global Priority	X20180	X30176
Local Priority	X20000	X30000
Pending Messages	Message<RID: 7216>, Message<RID: 7217>	Message<RID: 7218>, Message<RID: 7219>

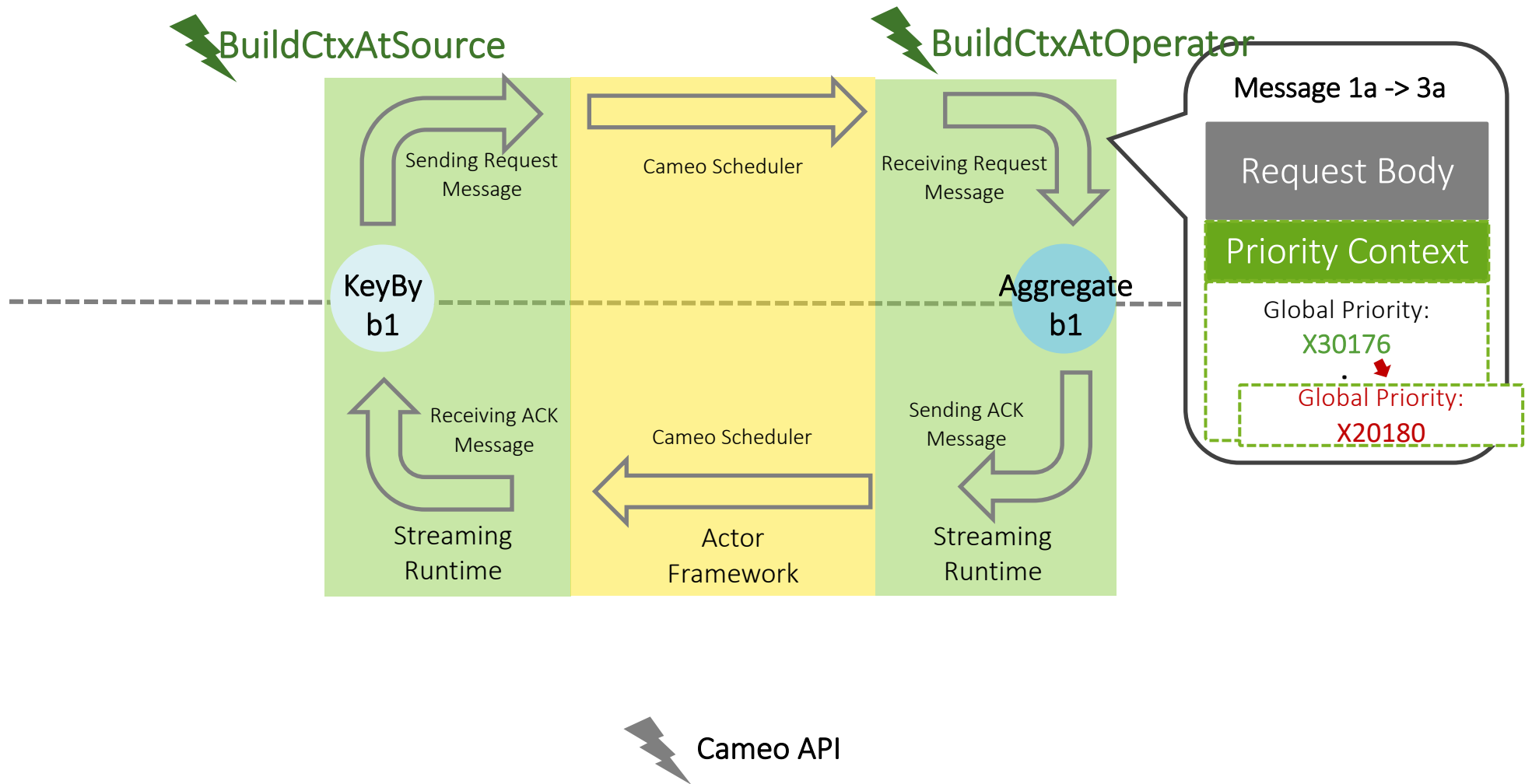


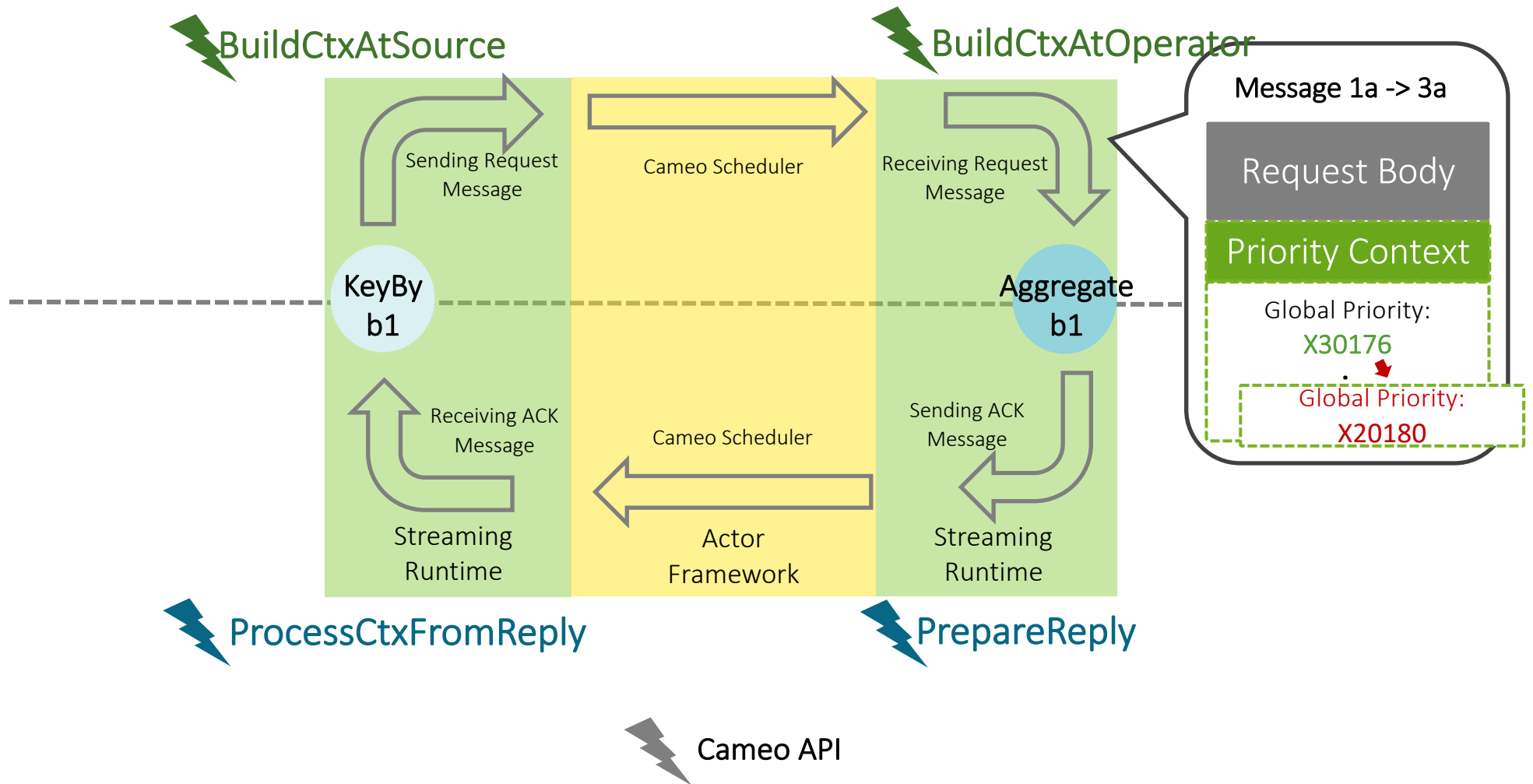






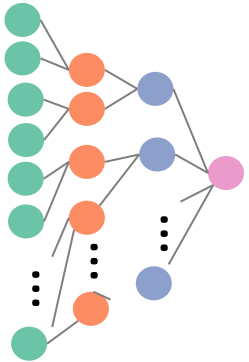






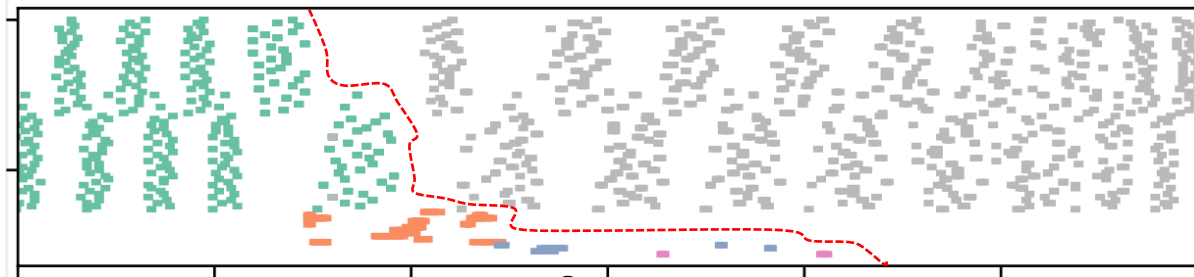
# Evaluation

# Cameo improves single-query performance

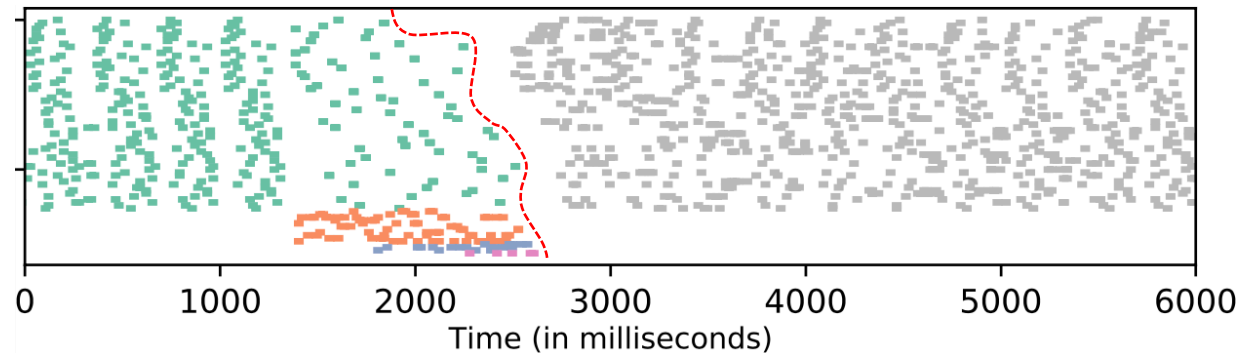


stage 0 stage 1 stage 2 stage 3

FIFO Schedule

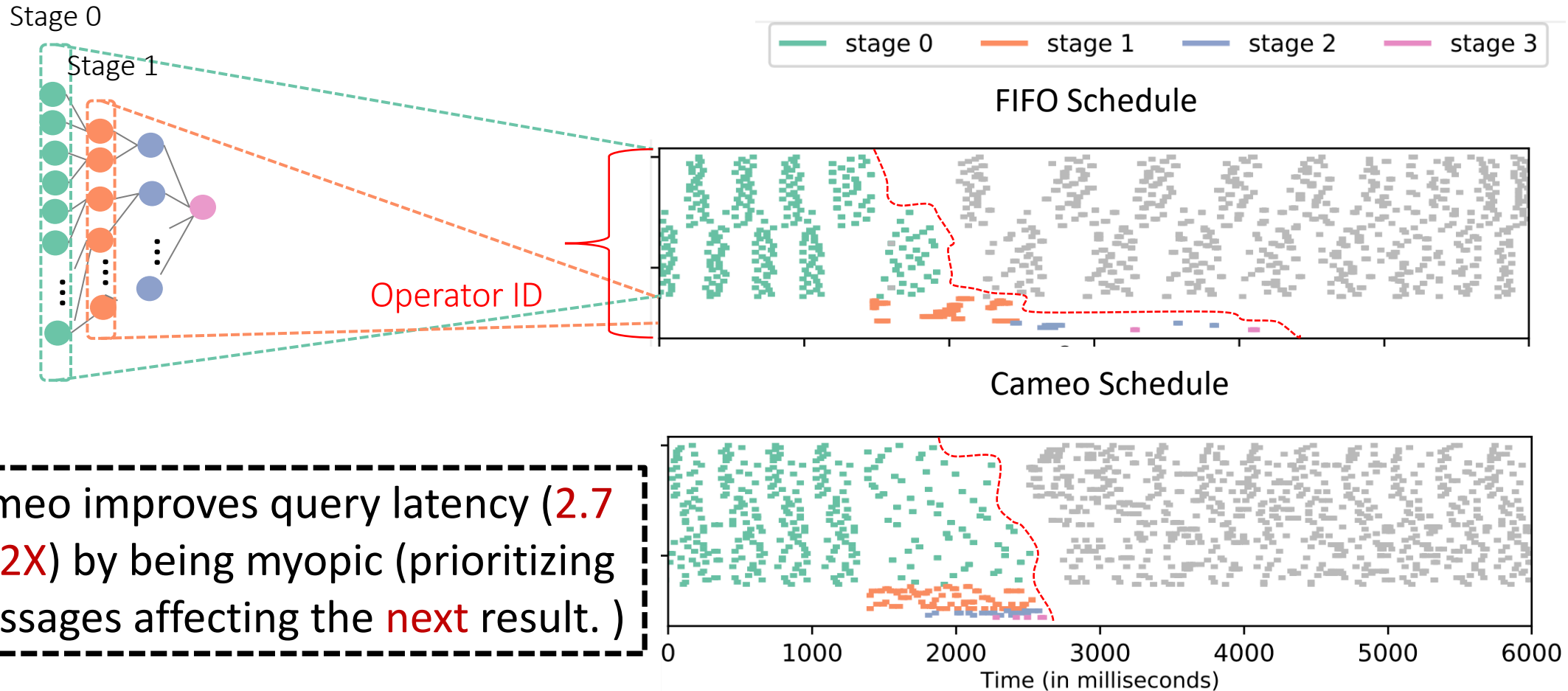


Cameo Schedule

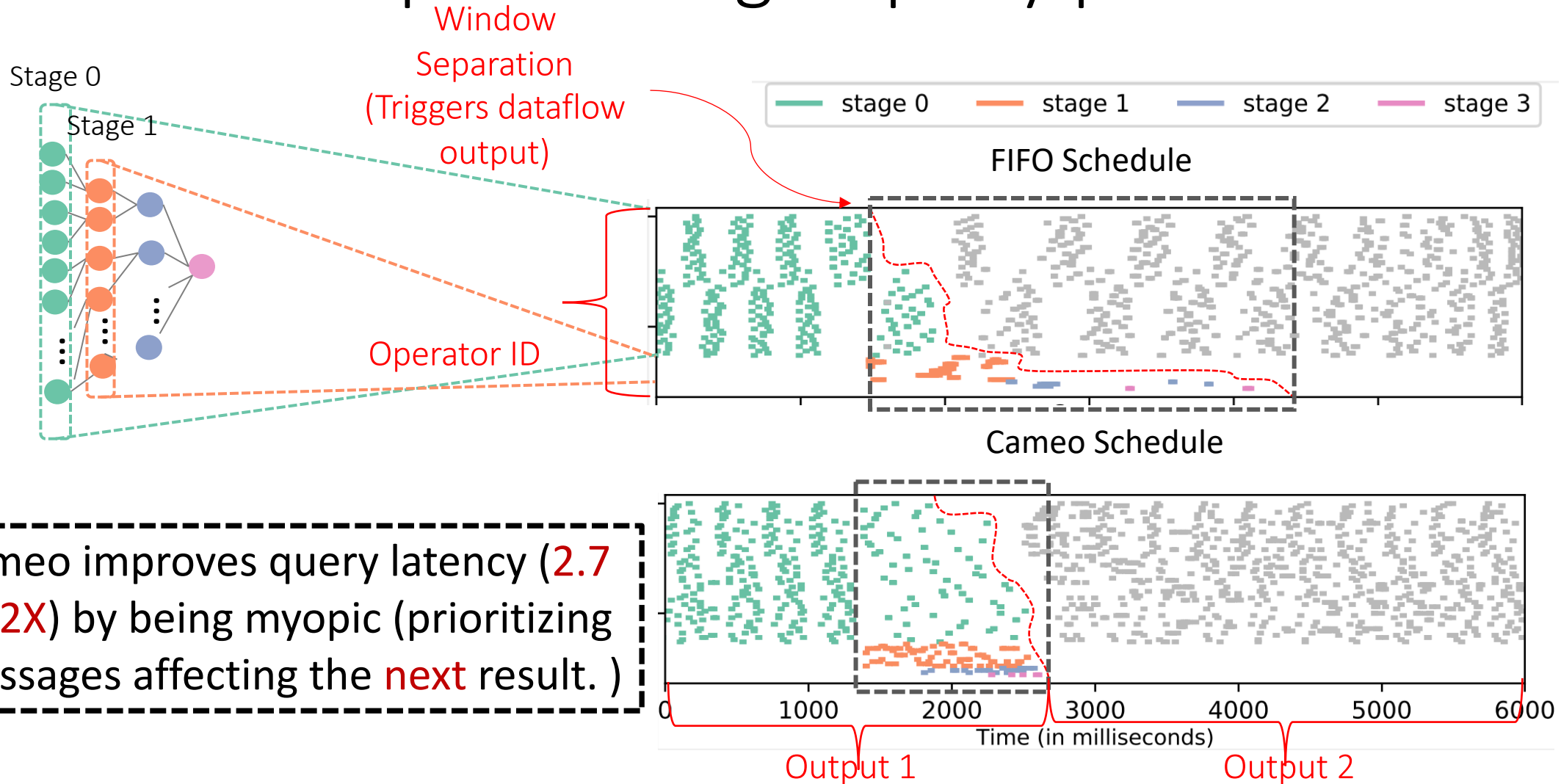


Cameo improves query latency (**2.7 - 3.2X**) by being myopic (prioritizing messages affecting the **next** result. )

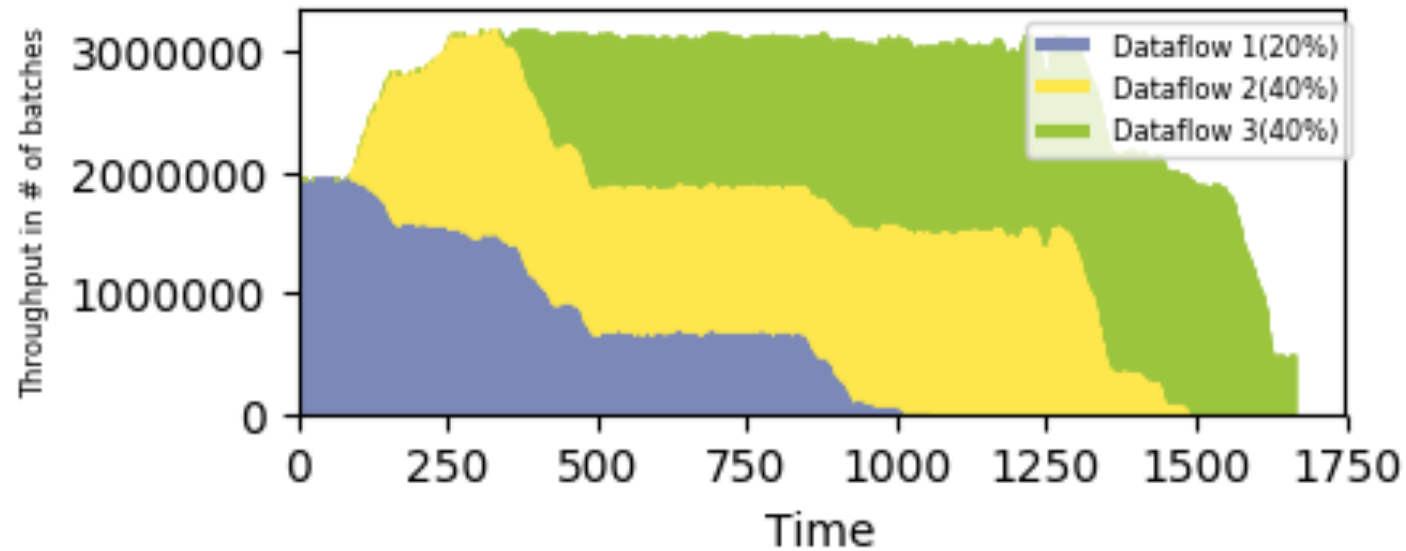
# Cameo improves single-query performance



# Cameo improves single-query performance



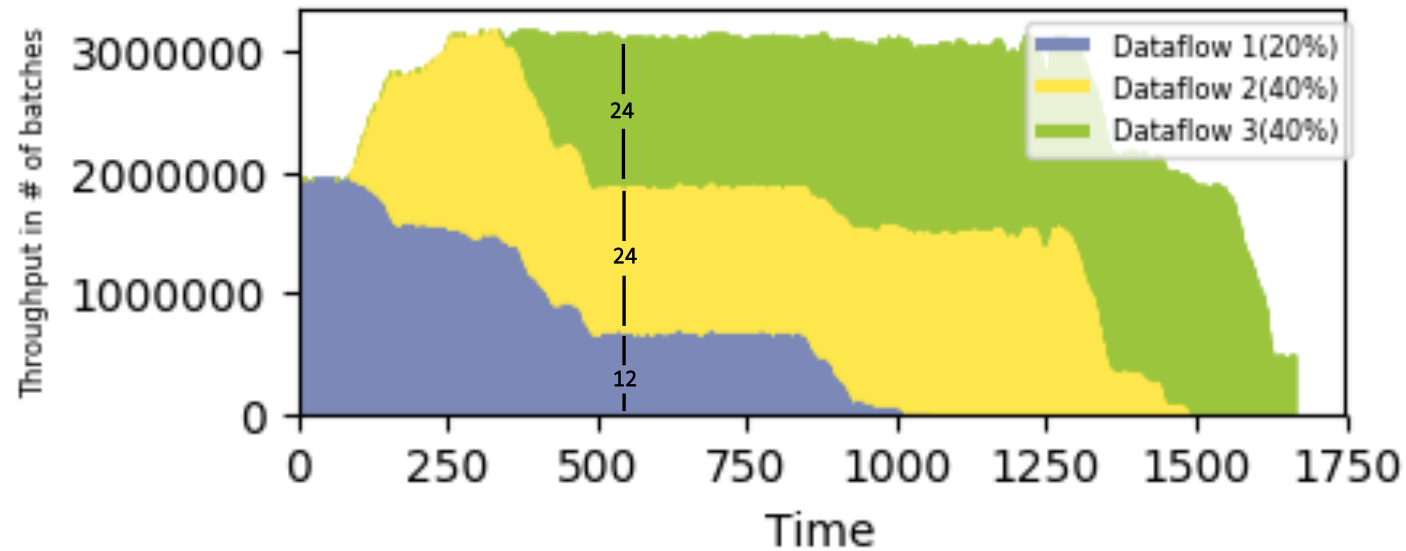
# Cameo supports plugged-in strategies



Proportional fair sharing. Three identical dataflows.  
Dataflows start 300 seconds apart

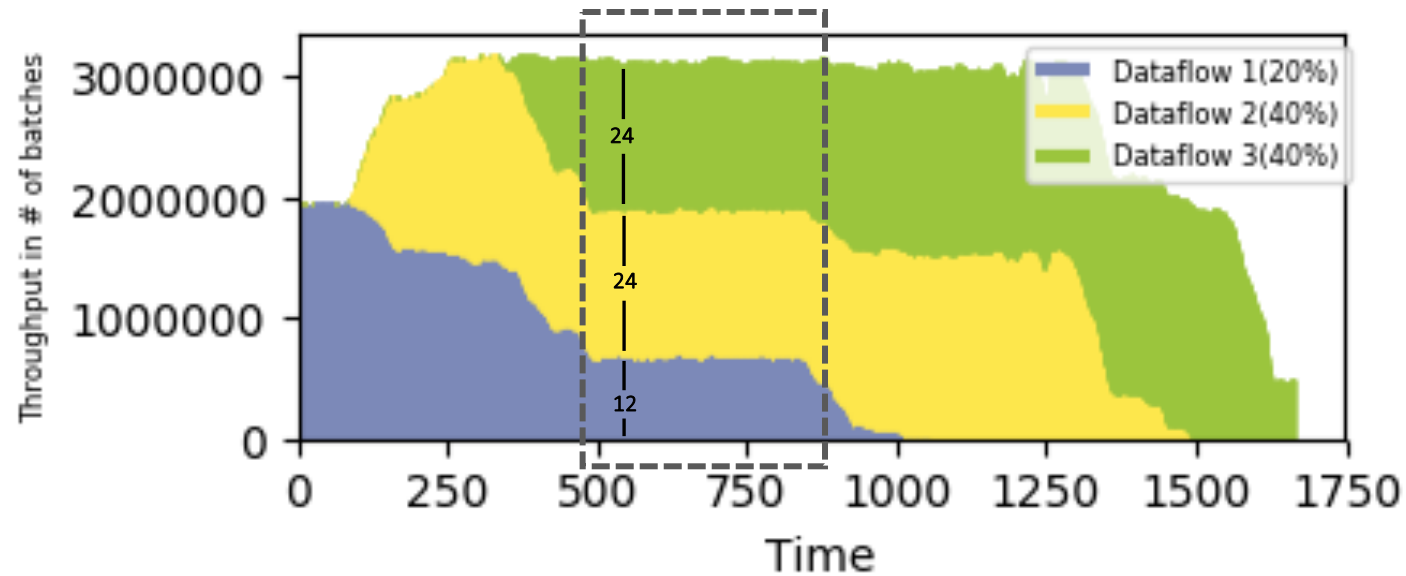


# Cameo supports plugged-in strategies



Proportional fair sharing. Three identical dataflows.  
Dataflows start 300 seconds apart

# Cameo supports plugged-in strategies



Proportional fair sharing. Three identical dataflows.  
Dataflows start 300 seconds apart

Before we end...

# More in the paper...

- Cameo improves query performance:
  - Cameo reduces query latency by **2.7X** for single query execution, **4.6X** in multi tenant scenarios
  - Cameo weathers transient spikes by reducing tail latency by **21X – 30X**
- Cameo supports general scheduling and user-defined scheduling states
- Cameo infers dataflow output based on input messages

# More in the paper...

- Cameo improves query performance:
  - Cameo reduces query latency by **2.7X** for single query execution, **4.6X** in multi tenant scenarios
  - Cameo weathers transient spikes by reducing tail latency by **21X – 30X**
- Cameo supports general scheduling and user-defined scheduling states
- Cameo infers dataflow output based on input messages

Cameo proposes **fine-grained** scheduling scheme for **serverless, event-driven** stream processing services that achieves both **high utilization** and **controllable performance**

For questions: contact author <[lexu1@Illinois.edu](mailto:lexu1@Illinois.edu)>