# SFSM: A Serverless Function Scheduling Method for FaaS Applications over Edge Computing

*Hao Tian, Cheng Chen, Fei Dai, Wanchun Dou*

NANJING UNIVERSITY

SOUTHWEST FORESTRY UNIVERSITY

# Serverless Computing

[1] AWS Lambda: https://aws.amazon.com/lambda
[2] Google Cloud Functions: https://cloud.google.com/functions
[3] Microsoft Azure Functions: https://azure.microsoft.com/en-us/products/functions
[4] IBM Cloud Functions: https://cloud.ibm.com/functions

# Function as a Services (FaaS)

# Function as a Services (FaaS)



**Pay only for what you use!**

# Serverless Workloads

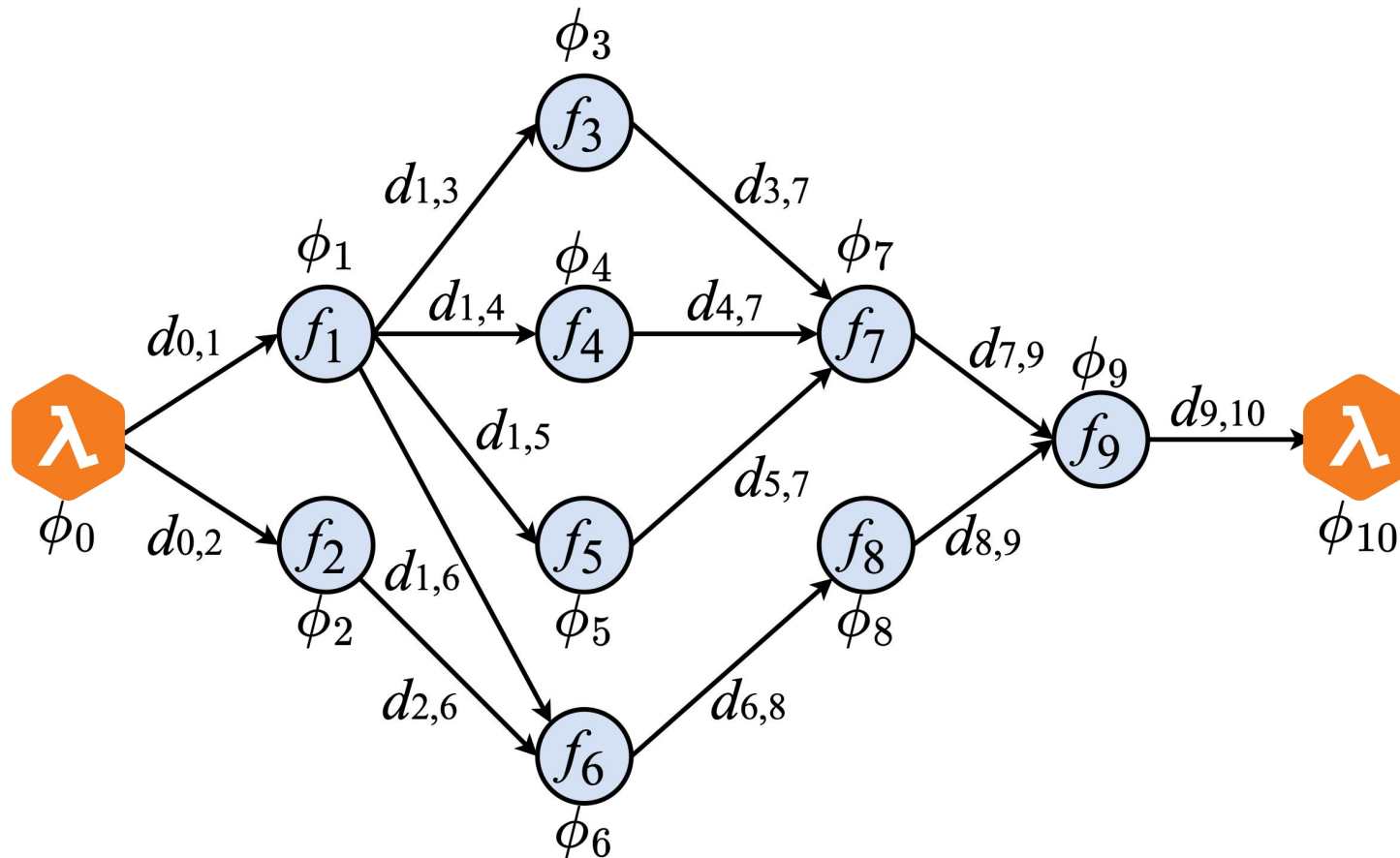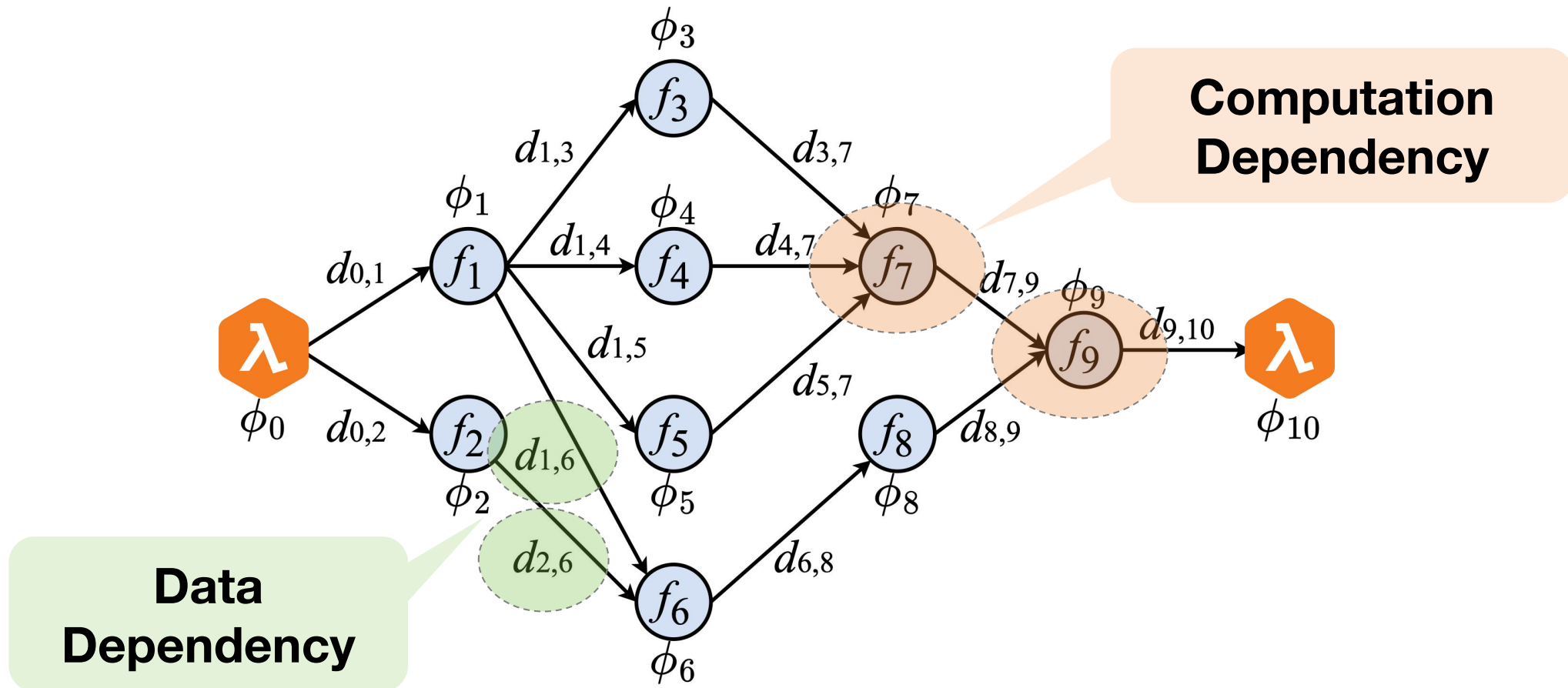- Abstract as a Direct Acyclic Graph (DAG)
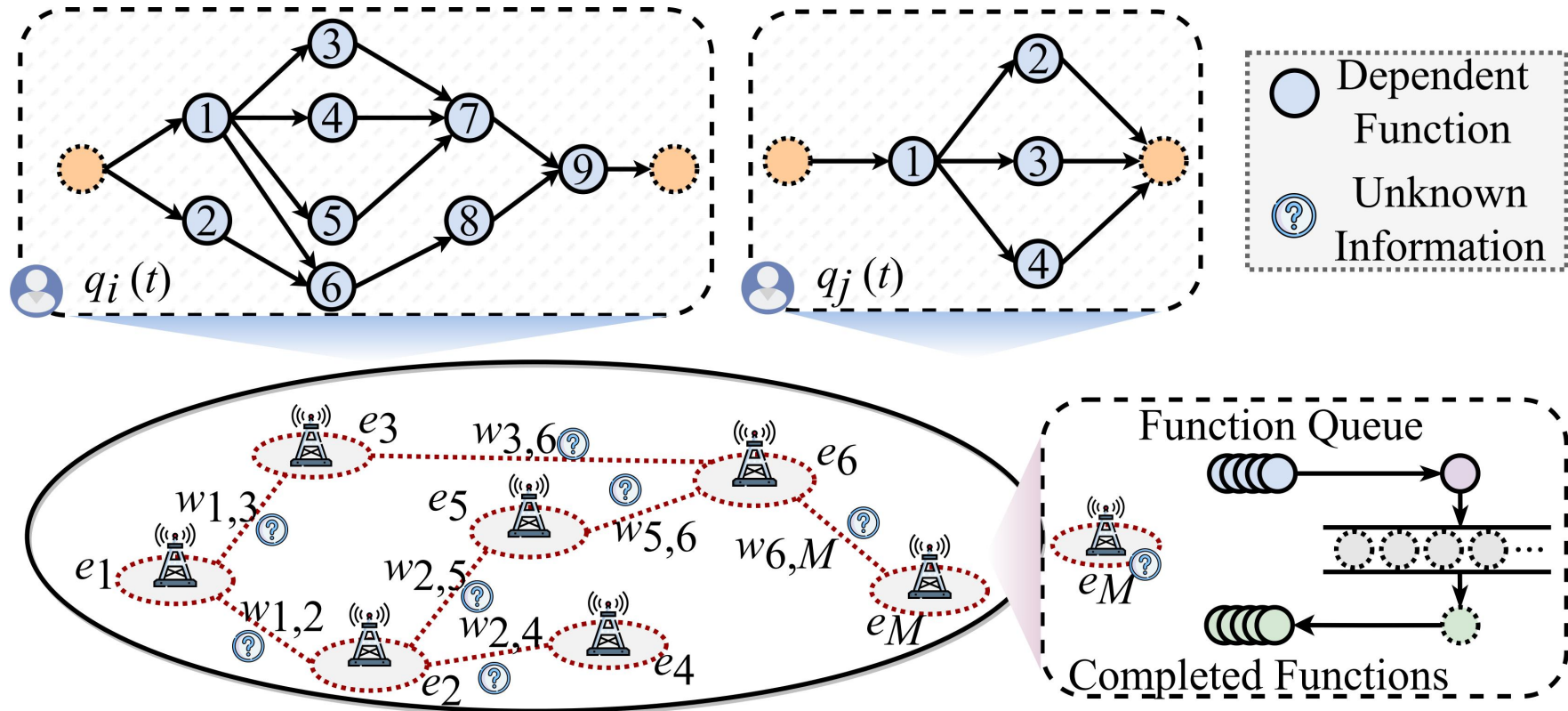- Constrained by computation & data dependencies

# Serverless Workloads

- Abstract as a Direct Acyclic Graph (DAG)
- Constrained by computation & data dependencies

# Serverless Edge Computing

- Unleash resources **from cloud to network edge**
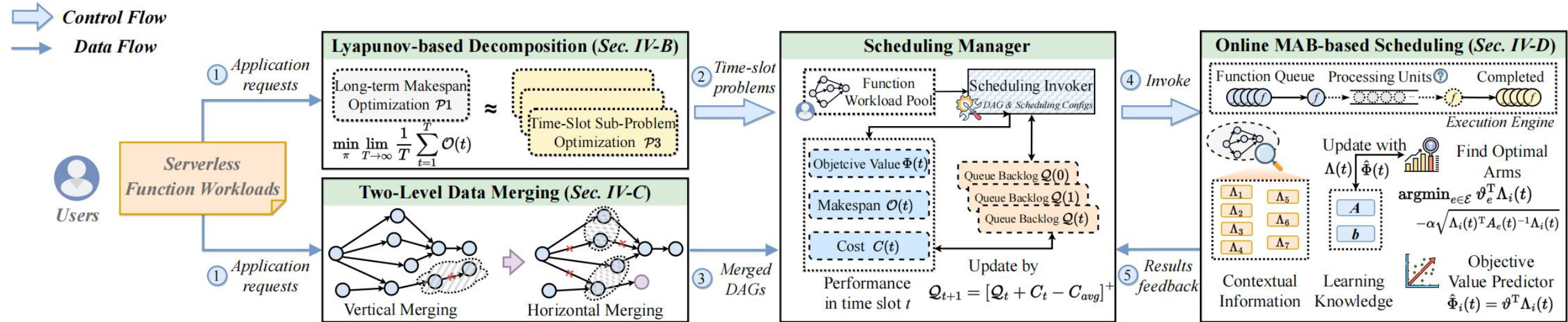- To meet **latency-sensitive** FaaS application demands

# Motivation



How to schedule serverless functions in edge coputing systems with no priors?

# SFSM Overview

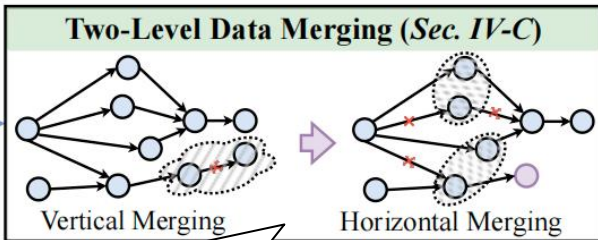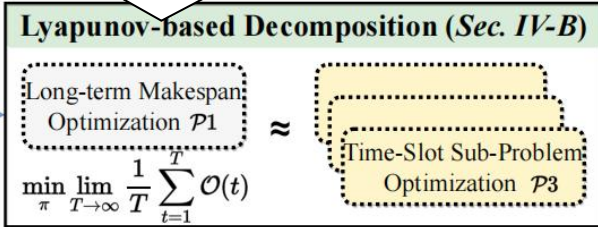# Our Solution: SFSM

**(Step 1)**

**Layapunov-based Decomposition**

**SFSM Overview**



**Two-Level Data Merging**

**(Step 2)**

**Online MAB-based Scheduling**

**(Step 3)**

## Step 1: Layapunov-based Decomposition



**Long-term averaged makespan**

**Lyapunov-drift-plus-penalty**

**Approximated by upper bound**

$\mathcal{P}1:$  $\min_{\pi} \lim_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} \mathcal{O}(t)$

s.t. $\mathcal{O}_{i,k}^{es}(t) \geq \max_{j \in \nu(f_i)} \{\mathcal{O}_{j,k}^{com}(t) + \mathcal{O}_{i,j,k}^{tr}(t)\},$

$\mathcal{O}_{i,k}^{es}(t) \geq \mathcal{O}_{i,k}^{conf}(t), \forall i \in \mathcal{F}, k \in \mathfrak{A},$

$\lim_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} C(t) \leq C_{avg},$

$\rho_{\pi_{f_i}(t),k}(t) \leq \zeta_{\pi_{f_i}(t)}, \forall i \in \mathcal{F}, k \in \mathfrak{A},$

$\sum_{k=1}^{|\mathbf{q}|} \rho_{\pi_{f_i}(t),k}(t) \leq \zeta_{\pi_{f_i}(t)}, \forall i \in \mathcal{F}.$

$\mathcal{P}2:$  $\min_{\pi} \mathbb{E}\left[\Delta(\mathcal{Q}(t)) + V \cdot \mathcal{O}(t) | \mathcal{Q}(t)\right]$

$\mathcal{P}3:$  $\min_{\pi} \Phi(t) = B + V \cdot \mathcal{O}(t) + \mathcal{Q}(t)(C(t) - C_{avg})$

## Step 1: Layapunov-based Decomposition

**Algorithm 1:** Long-term problem decomposition by Lyapunov optimization

---

**Input:** $\mathcal{U}$, $\mathfrak{A}$, $C_{avg}$, $V$, $\mathcal{Q}(0)$;
**Output:** $\pi(t)$ in each time slot $\forall t \in \mathcal{T}$;

1 **for** $t = 1$ *to* $T$ **do**
2     **for** $i = 1$ *to* $I$ **do**
3        **if** *Request $i$ arrives* **then**
4           Obtain the merged $\mathcal{A}_i(t)$ by **Alg. 2**;
5           Derive the scheduling decision $\pi_i(t)$ by **Alg. 3**;
6           Calculate $\mathcal{O}_i(t)$ and $C_i(t)$;
7        **end**
8        Update the virtual cost queue by Eq. (9);
9        Set $\pi(t) \leftarrow \pi(t) \cup \pi_i(t)$;
10     **end**
11 **end**
12 **Return** $\{\pi(1), \pi(2), \cdots, \pi(T)\}$.

---

**Call Alg. 2 (data merging)**
**&**
**Alg. 3 (online scheduling)**

## Step 1: Layapunov-based Decomposition

---

**Algorithm 1:** Long-term problem decomposition by Lyapunov optimization

---

**Input:** $\mathcal{U}$, $\mathfrak{A}$, $C_{avg}$, $V$, $\mathcal{Q}(0)$;
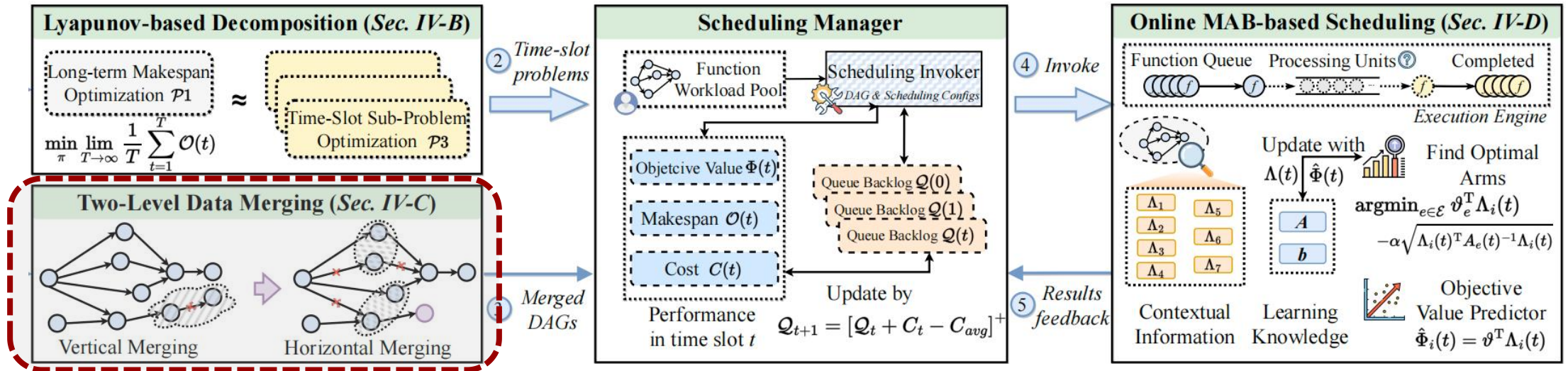
**Output:** $\pi(t)$ in each time slot $\forall t \in \mathcal{T}$;

1 **for** $t = 1$ *to* $T$ **do**
2      **for** $i = 1$ *to* $I$ **do**
3          **if** *Request i arrives* **then**
4              Obtain the merged $\mathcal{A}_i(t)$ by **Alg. 2**;
5              Derive the scheduling decision $\pi_i(t)$ by **Alg. 3**;
6              Calculate $\mathcal{O}_i(t)$ and $C_i(t)$;
7          **end**
8          Update the virtual cost queue by Eq. (9);
9          Set $\pi(t) \leftarrow \pi(t) \cup \pi_i(t)$;
10      **end**
11 **end**
12 **Return** $\{\pi(1), \pi(2), \cdots, \pi(T)\}$.

---

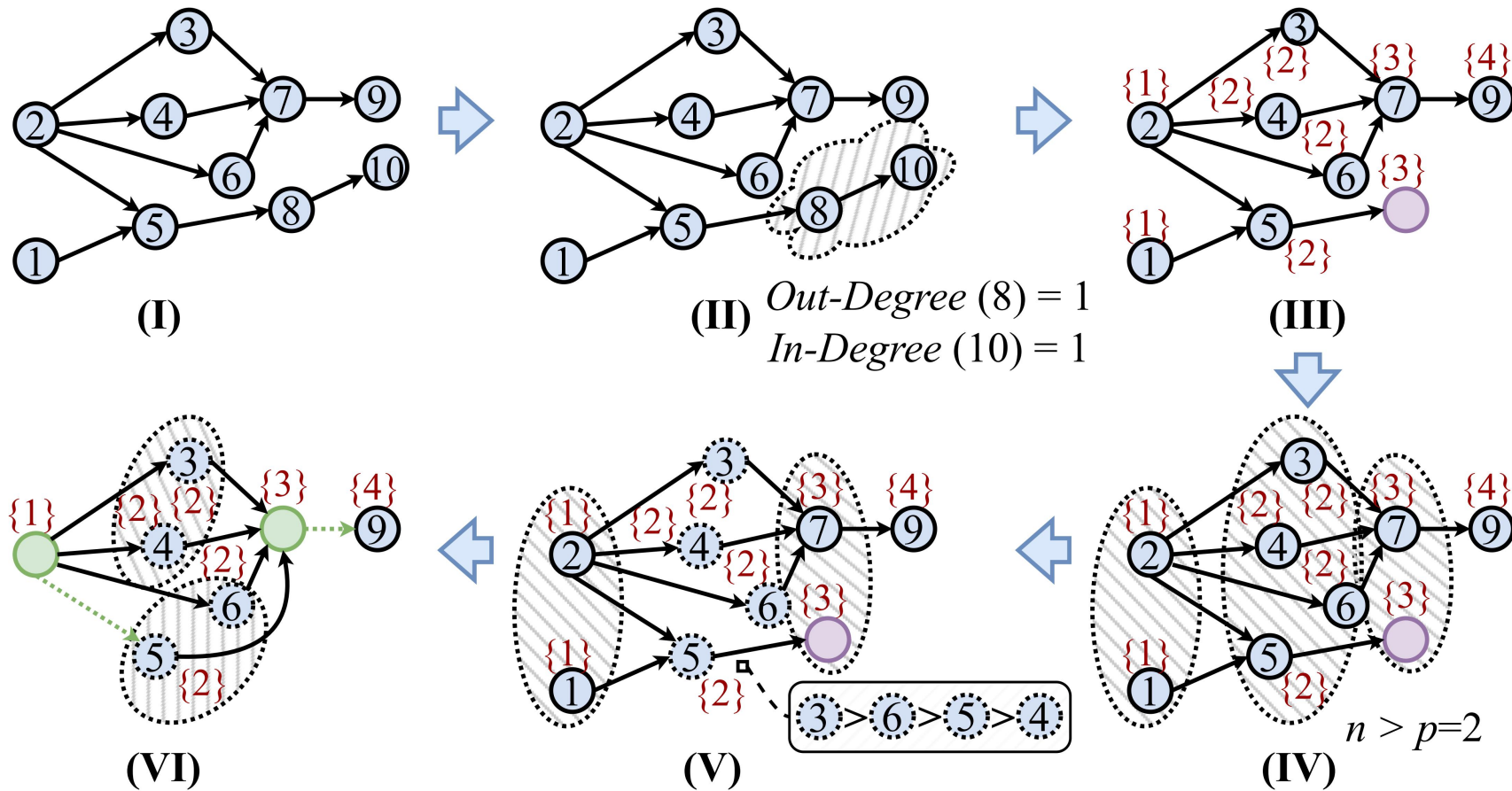**Update virtual cost queue based on Lyapunov optimization technique**

## Step 2: Two-Level Data Merging



- ■ **Vertical Merging:** find the node that only has one successor and its successor also has one predecessor.

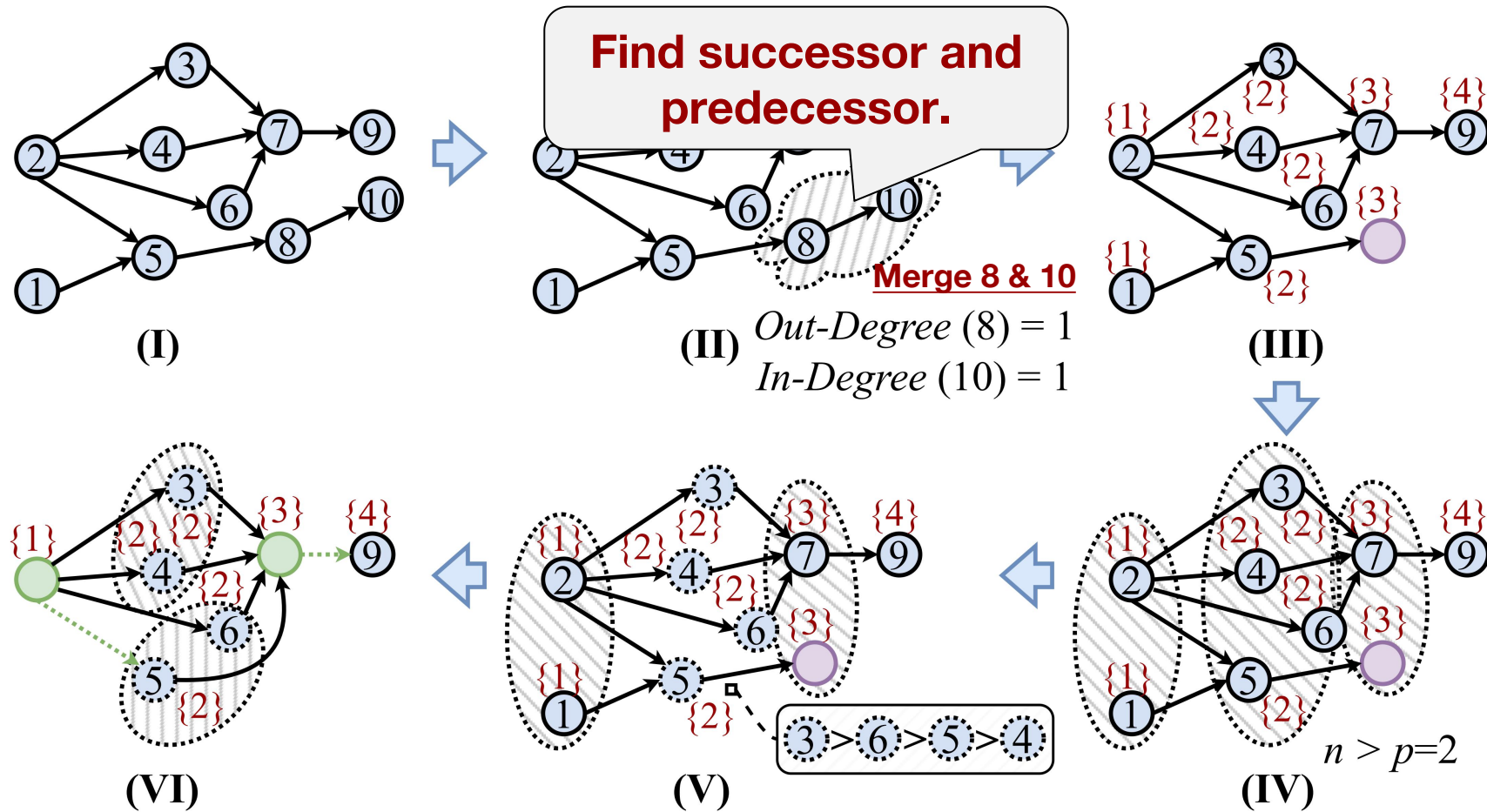- ■ **Horizontal Merging:** integrate nodes at the same node levels to improve the parallel performance.

## Step 2: Two-Level Data Merging



*Out-Degree (8) = 1*
*In-Degree (10) = 1*

*n > p=2*

*An example of 10 serverless functions*

## Step 2: Two-Level Data Merging



*An example of 10 serverless functions*

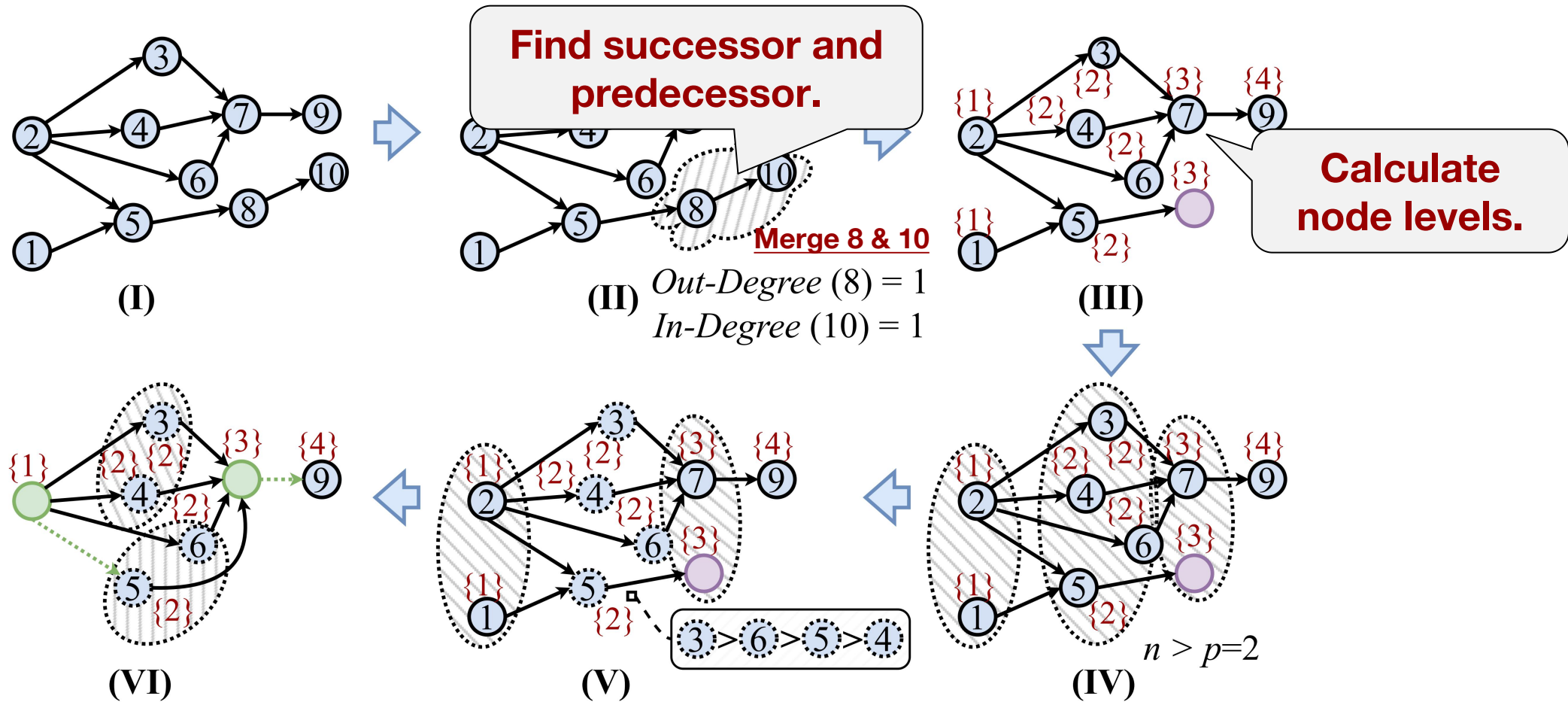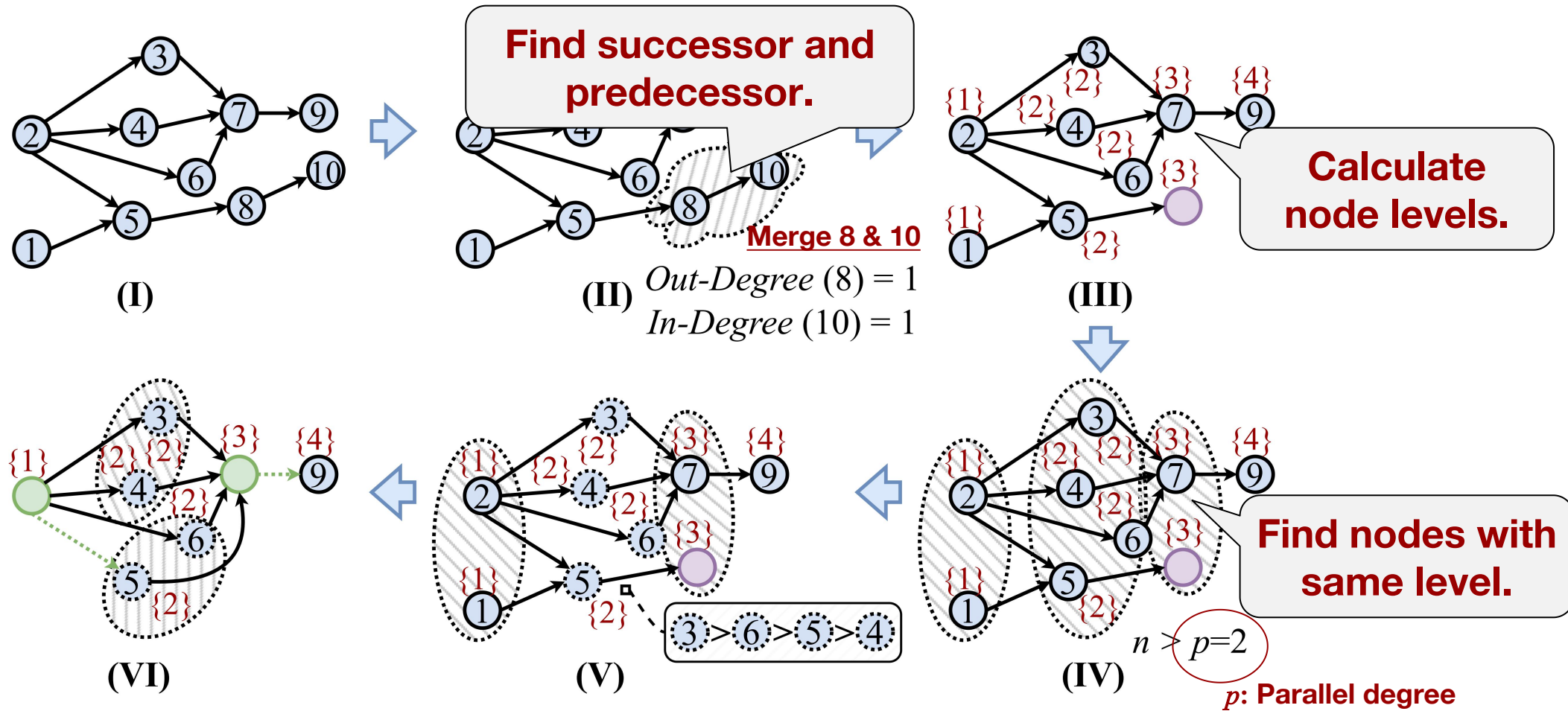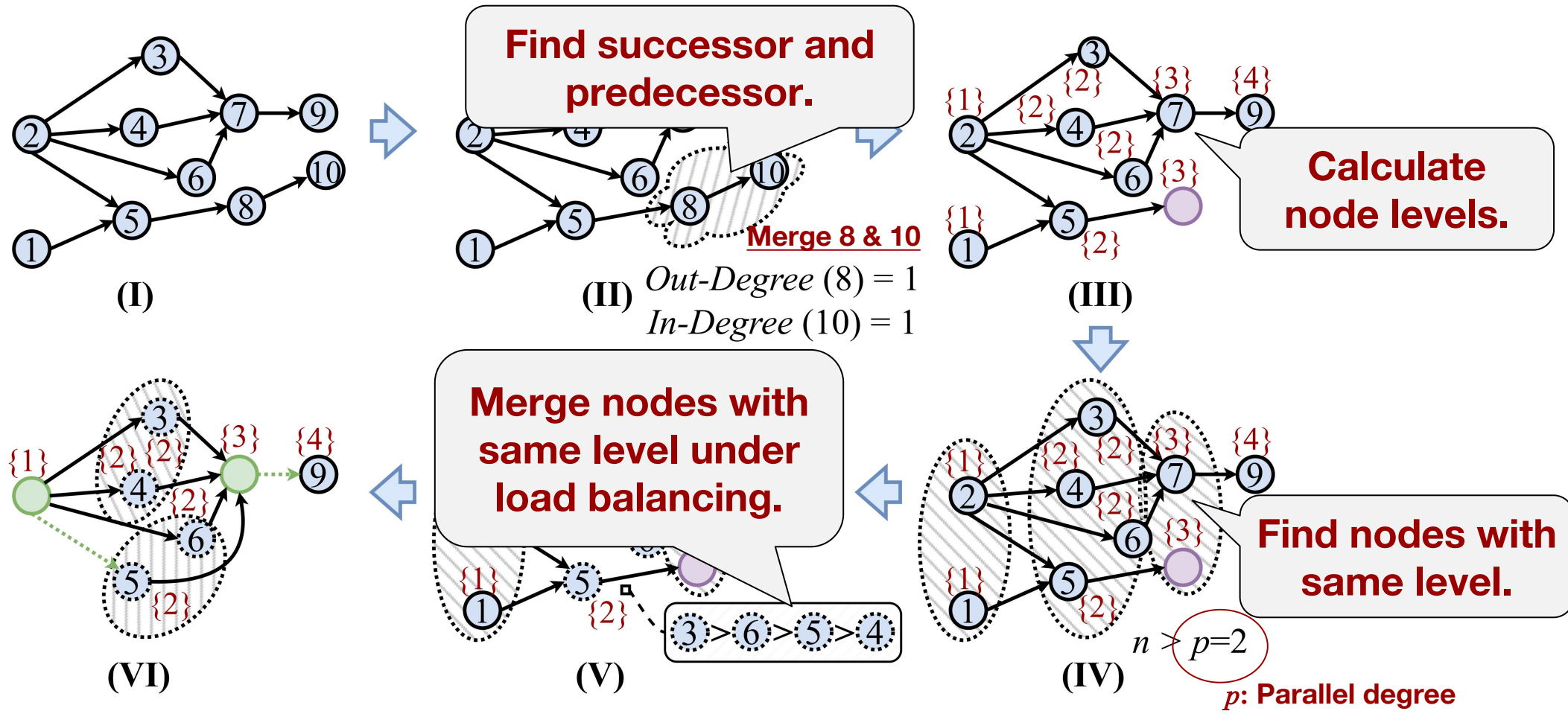## Step 2: Two-Level Data Merging



An example of 10 serverless functions

# Step 2: Two-Level Data Merging



*An example of 10 serverless functions*

## Step 2: Two-Level Data Merging



**Find successor and predecessor.**

**Merge 8 & 10**
*Out-Degree* (8) = 1
*In-Degree* (10) = 1

**Calculate node levels.**

**Merge nodes with same level under load balancing.**

3>6>5>4

**Find nodes with same level.**

$n > p=2$

*p*: Parallel degree

(I)   (II)   (III)   (IV)   (V)   (VI)

*An example of 10 serverless functions*

## Step 2: Two-Level Data Merging



An example of 10 serverless functions

## Step 2: Two-Level Data Merging

**Algorithm 2:** Two-level data merging

**Input:** $\mathcal{A}$, $p$, $L$;
**Output:** Merged DAG $\mathcal{A}'$;

```
1  for f in topological sort of A do
2      if out-degree(f) is 1 then
3          g ← suc(f);
4          if in-degree(g) is 1 and pre(g) is f then
5              φf ← φf + φg;
6              suc(f) ← suc(g);
7              del(g);
8          end
9      end
10 end
```
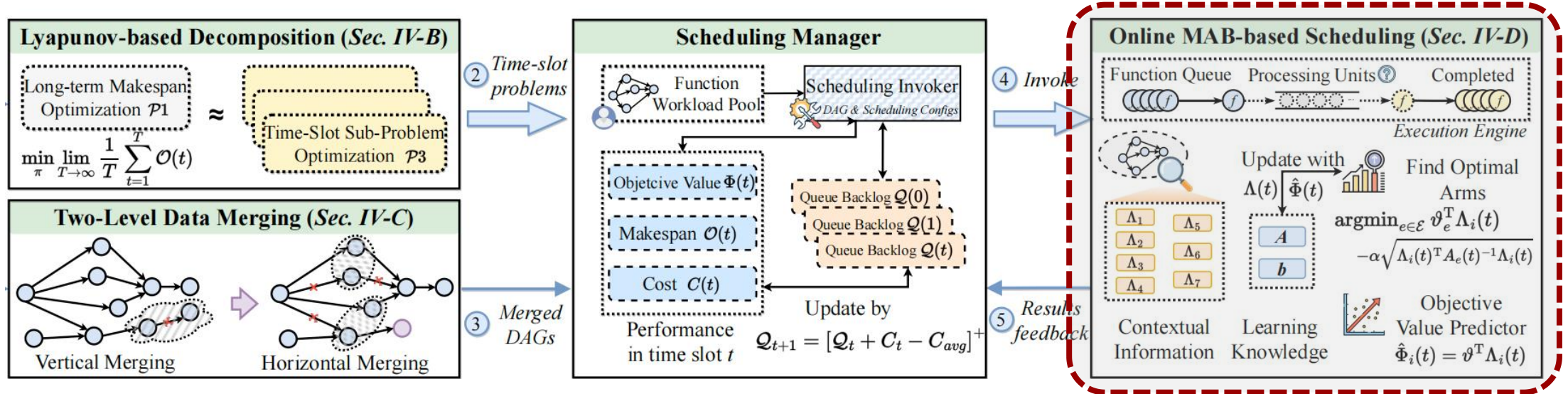
**Vertical Merging**

```
11 for l = 1 to L do
12     for f in group l do
13         w ← w ∪ φf;
14     end
15     Sort w in descending order;
16     i ← 1, j ← |w|;
17     while i < j do
18         n ← min(p, j − i + 1);
19         merge(A', n);
20         φf1 ← {φf1, φf2, · · · , φfn};
21         suc(fi) ←
               {suc(fi), suc(fj), · · · , suc(fi+⌈n/2⌉), suc(fj−⌊n/2⌋)};
22         pre(fi) ←
               {pre(fi), pre(fj), · · · , pre(fi+⌈n/2⌉), pre(fj−⌊n/2⌋)};
23         del(A', n);
24         i ← i + ⌈p/2⌉, j ← j − ⌊p/2⌋;
25     end
26 end
27 Return A'.
```
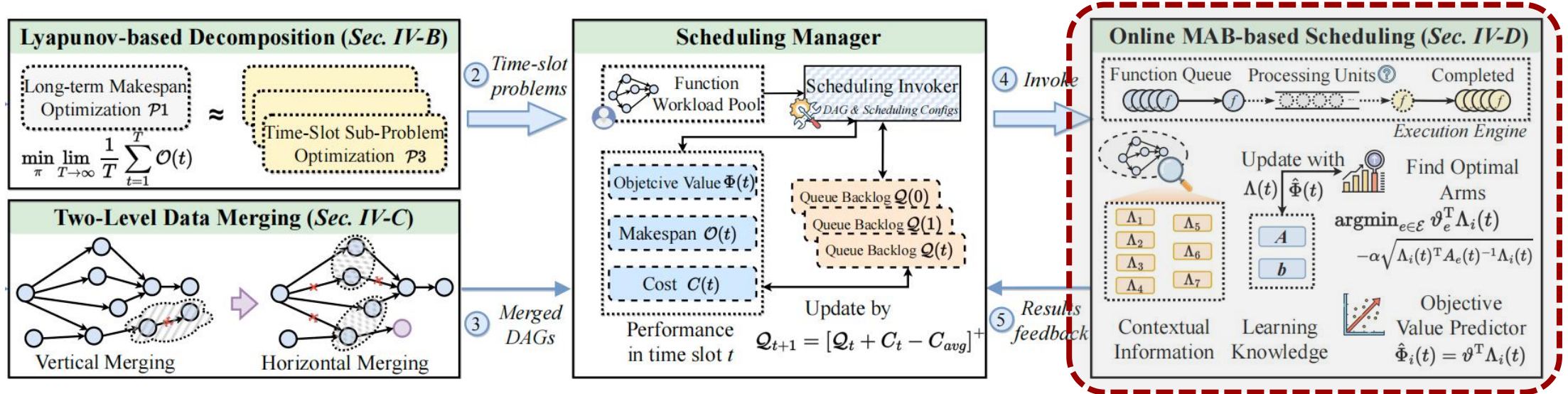
**Horizontal Merging**

## Step 3: Online MAB-based Scheduling



- **Contextual information:** request info, computations, data, etc.
- **Without any priors of edge networks:** i.e., transmission link and computing resources

## Step 3: Online MAB-based Scheduling



**ridge regression**

$$\pi_i(t) = \operatorname{argmin} \vartheta^{\mathrm{T}} \Lambda_i(t) \quad\Longrightarrow\quad \pi_i(t) = \operatorname{argmin}_{e \in \mathcal{E}} \left[ \vartheta_e^{\mathrm{T}} \Lambda_i(t) - \alpha \sqrt{\Lambda_i(t)^{\mathrm{T}} A_e(t)^{-1} \Lambda_i(t)} \right]$$

**learnable weight**

**matrix of contextual information**

**balance exploration and exploitation**

## Step 3: Online MAB-based Scheduling



$$\pi_i(t) = \mathrm{argmin}\, \vartheta^{\mathrm{T}} \Lambda_i(t)$$

**ridge regression**

$$\pi_i(t) = \mathrm{argmin}_{e \in \mathcal{E}} \left[ \vartheta_e^{\mathrm{T}} \Lambda_i(t) - \alpha \sqrt{\Lambda_i(t)^{\mathrm{T}} A_e(t)^{-1} \Lambda_i(t)} \right]$$

**exploration**

**learnable weight**

**matrix of contextual information**

**balance exploration and exploitation**

## Step 3: Online MAB-based Scheduling

**Algorithm 3:** Online serverless function scheduling

**Input:** $\mathcal{A}'_i$, $\mathcal{U}$;

**Output:** $\pi_i$;

1 Construct the contextual knowledge $\Lambda_i(t)$;
2 Initialize $A_e$ and $b_e$, $\forall e \in \mathcal{E}$;
3 **repeat**
4    **for** $f$ *in topological sort of* $\mathcal{A}'_i$ **do**
5       **for** $e = 1$ *to* $|\mathcal{E}|$ **do**
6          $\hat{\vartheta}_e(h) \leftarrow A_e^{-1}(h-1)b_e(h-1)$;
7          $\hat{\Phi}_{f,e}(h) \leftarrow$
           $\hat{\vartheta}_e(h)^{\mathrm{T}}\Lambda_i(h) - \alpha\sqrt{\Lambda_i(h)^{\mathrm{T}}A_e(h)^{-1}\Lambda_i(h)}$;
8       **end**
9       $\pi_{i,f}(h) \leftarrow \operatorname{argmin}_{e \in \mathcal{E}} \hat{\Phi}_{f,e}(h)$;
10       $\pi_i(h) \leftarrow \pi_i(h) \cup \pi_{i,f}(h)$;
11    **end**
12    **for** *Activated* $e \in \mathcal{E}$ **do**
13       $A_e(h) \leftarrow A_e(h-1) + \Lambda_i(h)\Lambda_i(h)^{\mathrm{T}}$;
14       $b_e(h) \leftarrow b_e(h-1) + \Lambda_i(h)\hat{\Phi}_{f,e}(h)$;
15    **end**
16    $\pi_i \leftarrow \pi_i \cup \pi_i(h)$;
17 **until** *Learning step h finished*;
18 **Return** $\pi_i$.

**Calculate weight and predict objective value**

## Step 3: Online MAB-based Scheduling

**Algorithm 3:** Online serverless function scheduling

**Input:** $\mathcal{A}'_i$, $\mathcal{U}$;
**Output:** $\pi_i$;

1  Construct the contextual knowledge $\Lambda_i(t)$;
2  Initialize $A_e$ and $b_e$, $\forall e \in \mathcal{E}$;
3  **repeat**
4      **for** $f$ *in topological sort of* $\mathcal{A}'_i$ **do**
5          **for** $e = 1$ *to* $|\mathcal{E}|$ **do**
6              $\hat{\vartheta}_e(h) \leftarrow A_e^{-1}(h-1)b_e(h-1)$;
7              $\hat{\Phi}_{f,e}(h) \leftarrow$
                $\hat{\vartheta}_e(h)^{\mathrm{T}}\Lambda_i(h) - \alpha\sqrt{\Lambda_i(h)^{\mathrm{T}}A_e(h)^{-1}\Lambda_i(h)}$;
8          **end**
9          $\pi_{i,f}(h) \leftarrow \mathrm{argmin}_{e \in \mathcal{E}}\hat{\Phi}_{f,e}(h)$;
10         $\pi_i(h) \leftarrow \pi_i(h) \cup \pi_{i,f}(h)$;
11     **end**
12     **for** *Activated* $e \in \mathcal{E}$ **do**
13         $A_e(h) \leftarrow A_e(h-1) + \Lambda_i(h)\Lambda_i(h)^{\mathrm{T}}$;
14         $b_e(h) \leftarrow b_e(h-1) + \Lambda_i(h)\hat{\Phi}_{f,e}(h)$;
15     **end**
16     $\pi_i \leftarrow \pi_i \cup \pi_i(h)$;
17 **until** *Learning step $h$ finished*;
18 **Return** $\pi_i$.

**Select the target scheduling edge server with minimum value $\hat{\Phi}_i(t)$**

# Evaluation Settings

- **Workloads**
  - Alibaba cluster trace[1], which includes about 4000 machines in a periods of 8 days, over 20,000 applications
  - Three request modes, i.e., data-intensive, computation-intensive and both data and computation-intensive
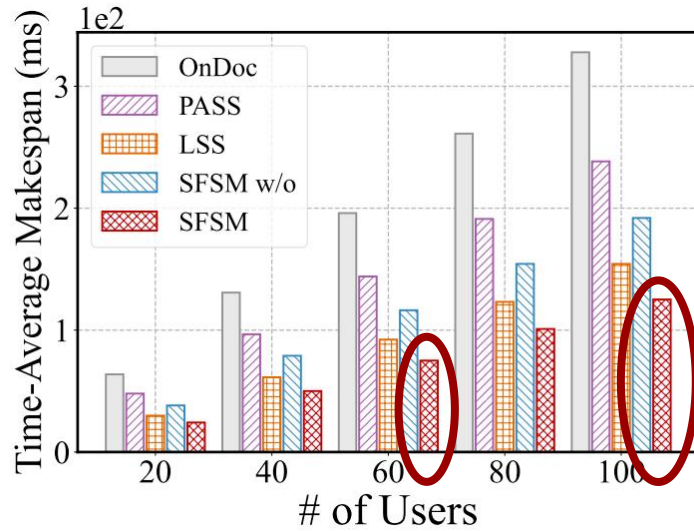
- **Baselines**
  - OnDoc [TOSN 2023]
  - PASS [IWQoS 2022]
  - Latency-myopic Static Scheduling (LSS) [TPDS 2022, IWQoS 2019]
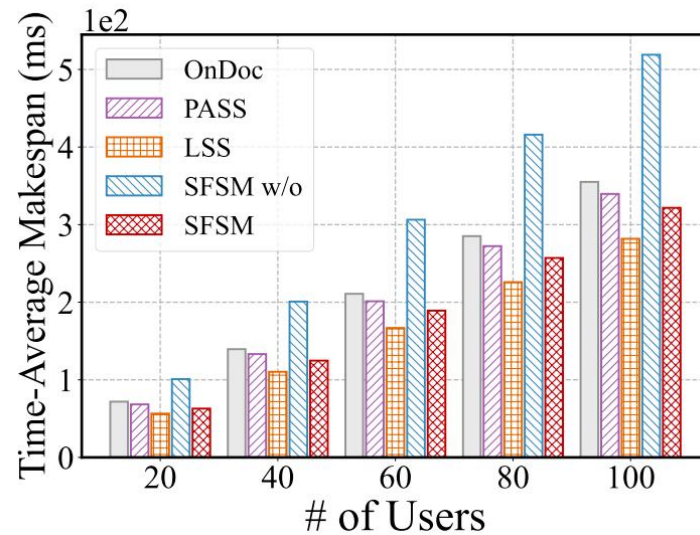  - SFSM w/o

- **Metircs**
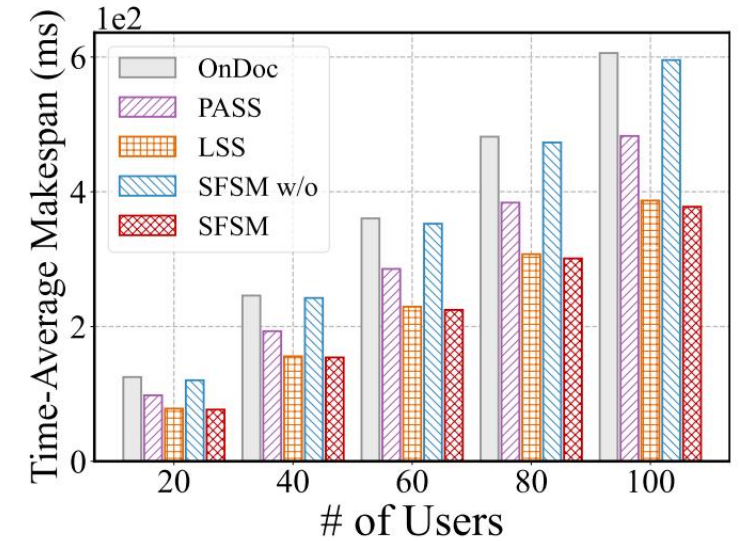  - Time-Average Makespan
  - Time-Average Cost

[1] https://github.com/alibaba/clusterdata
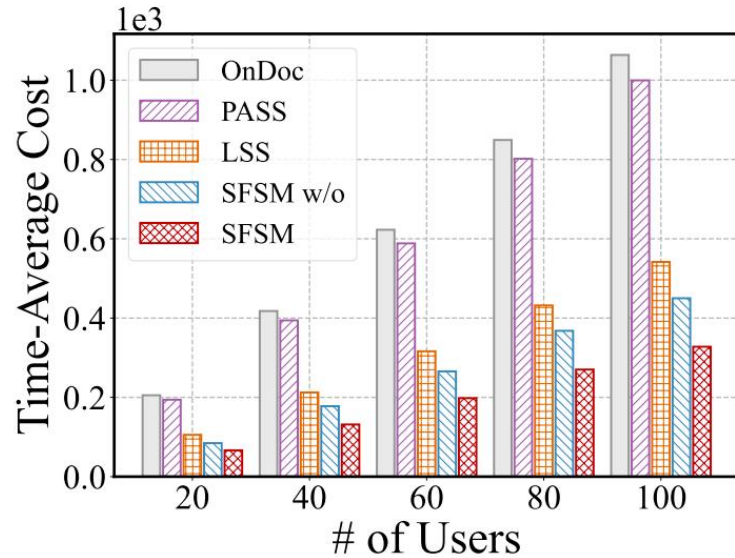
(a) data-intensive

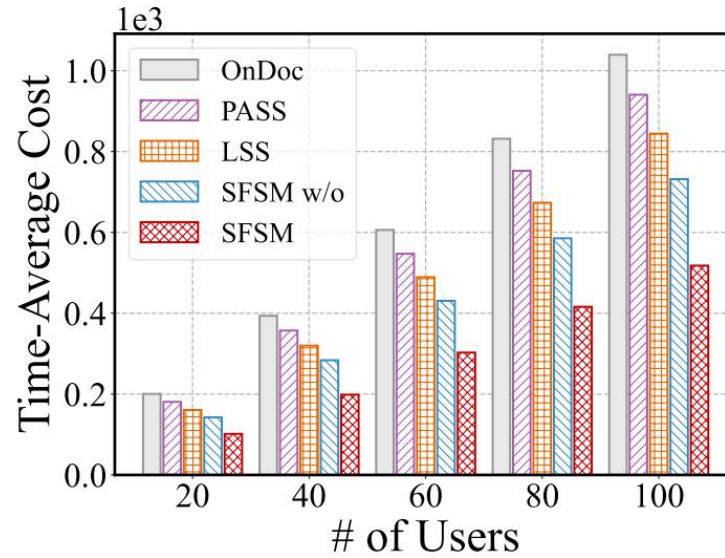(b) computation-intensive

(c) both data and computation-intensive

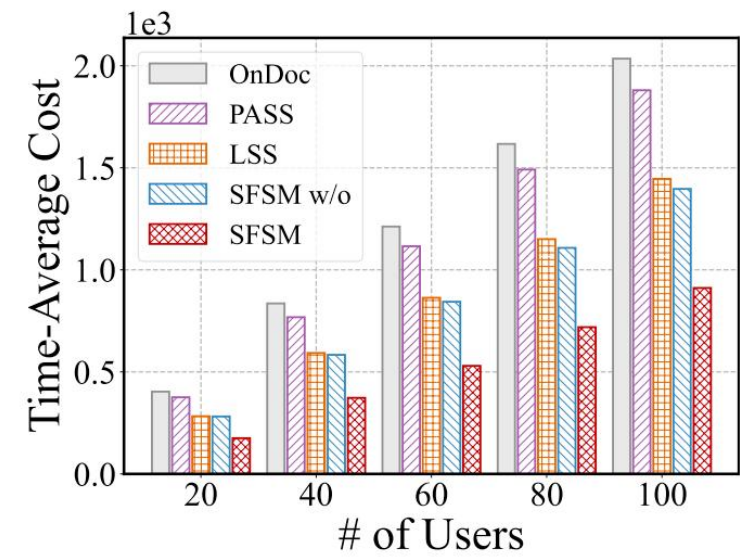**Lower makespan, especially with intensive data req.**

# Evaluation: Cost



(a) data-intensive

(b) computation-intensive

(c) both data and computation-intensive

**Lower cost within different request modes**

# Conclusion

- SFSM aims to achieve the **long-term makespan** minimization while improving the execution cost

- Two-level data merging algorithm to **mitigate the transmission overhead** of redundant data

- **Online** contexual shceduling **without priors** to realize fine-grained decisions

# Thank You For Your Attention!