# PULCEO in Action
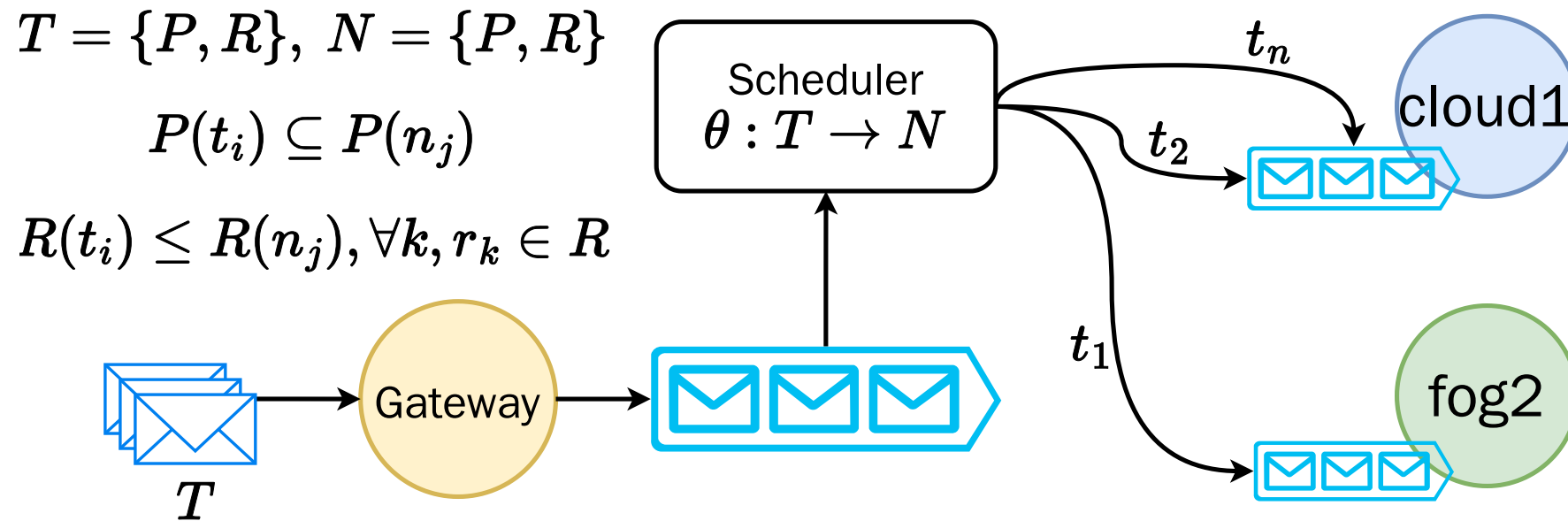
API-driven Task Scheduling and Offloading
with PULCEO: An Extension

Sebastian Böhm and Guido Wirtz

University of Bamberg, Germany

https://github.com/spboehm/pulceo-misc

# Task Scheduling / Offloading

$$T = \{P, R\}, \ N = \{P, R\}$$

$$P(t_i) \subseteq P(n_j)$$

$$R(t_i) \leq R(n_j), \forall k, r_k \in R$$



- IoT / mobile clients, submitting independent / atomic tasks $T$

- Heterogeneous cloud, edge and fog nodes $N$

- Tasks $T$ and Nodes $N$ are described by

  - Properties $P$: Layer, Location, Costs, …

  - Requirements $R$: CPU, Memory, Bandwidth, …

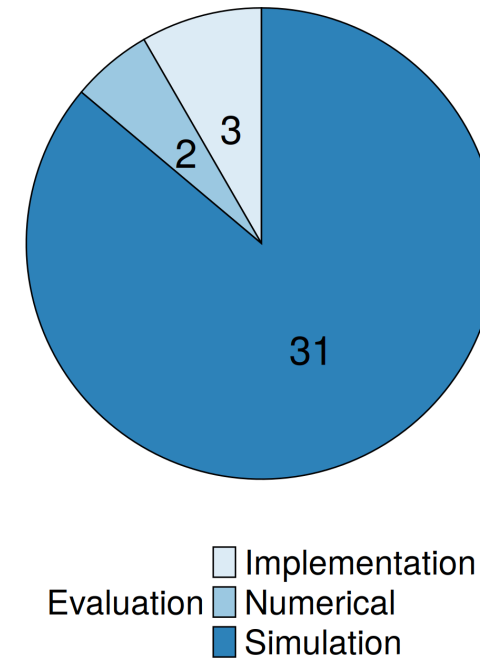https://github.com/spboehm/pulceo-misc

# Motivation

Many solutions exist for **task scheduling / offloading**

But, similar to **service placement**:

- **Reproducibility**
- (General) **Applicability**

are **limited** because of:

- **simulation-based** evaluations
- **custom** implementations
- missing **real-world experiments**



Only **3** out of **36** task offloading solutions were using real-world evaluations.[1]

1. S. Dong et al., "Task offloading strategies for mobile edge computing: A survey," Computer Networks, vol. 254, p. 110791, Dec. 2024, doi: 10.1016/j.comnet.2024.110791.

https://github.com/spboehm/pulceo-misc

# Simulations are sufficient?

Selected limitations, addressed by the **authors** of selected solutions:

> Where the value of **CPU weight** and **memory weight** is an **empirical issue**, and **multiple experiments are required** to obtain the optimal value such that α + β = 1. In this paper, all parameters are set to appropriate values.[1]
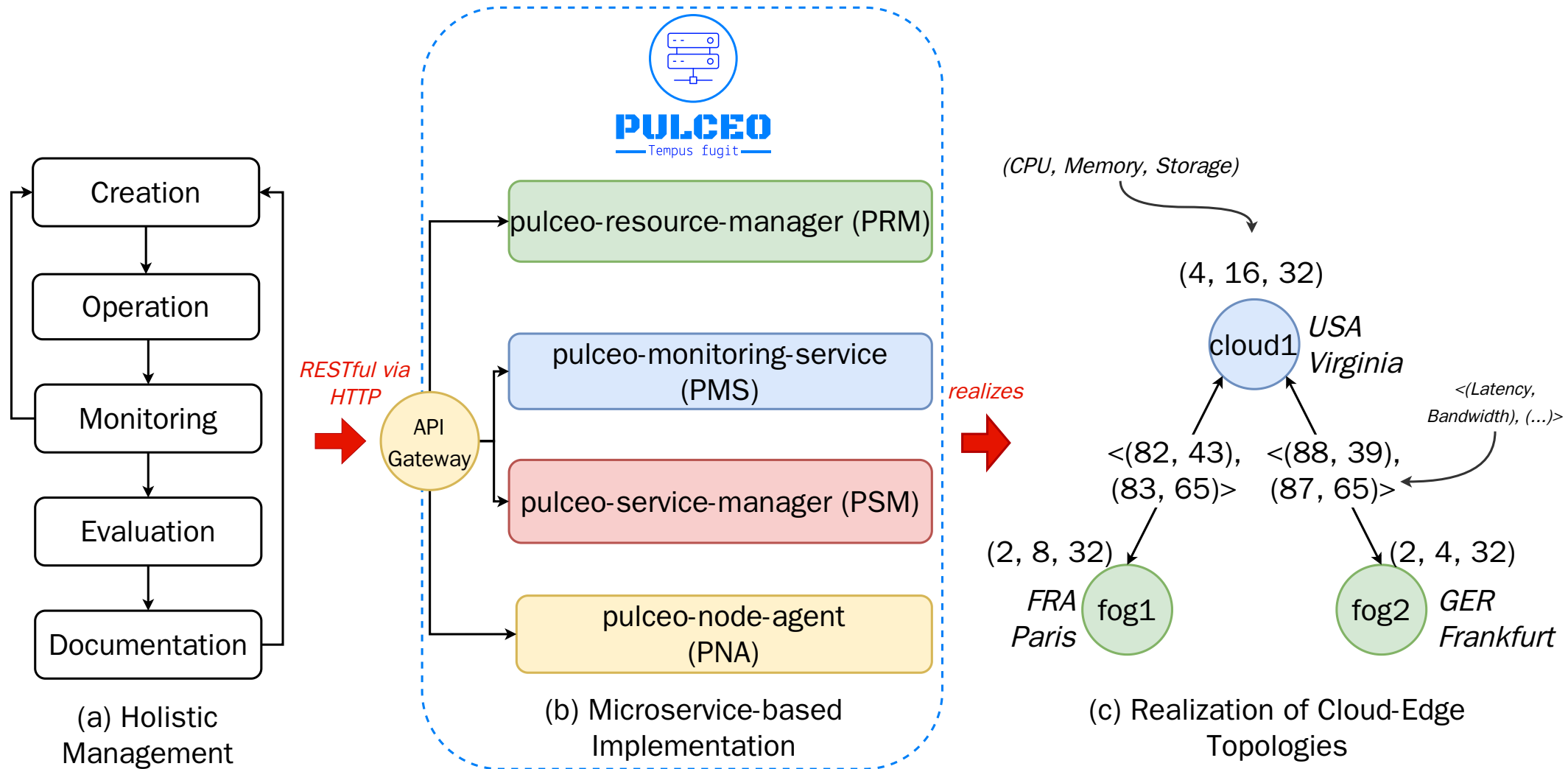
> Phare [the scheduling solution] requires updated information on the clusters participating in the federation, including **CPU**, **memory**, and **bandwidth usage**. A **real-world implementation** should thus **carefully balance** the resolution of such data and the additional overhead required to process it.[2]

Consequently, there is a need for **real-world experiments** and **evaluations**.

1. Y. Dong, G. Xu, M. Zhang, and X. Meng, "A High-Efficient Joint 'Cloud-Edge' Aware Strategy for Task Deployment and Load Balancing," IEEE Access, vol. 9, pp. 12791–12802, 2021, doi: 10.1109/ACCESS.2021.3051672.

2. G. Castellano, S. Galantino, F. Risso, and A. Manzalini, "Scheduling Multi-Component Applications Across Federated Edge Clusters With Phare," IEEE Open J. Commun. Soc., vol. 5, pp. 1814–1826, 2024, doi: 10.1109/OJCOMS.2024.3377917.

# Holistic Cloud-Edge Orchestration



(a) Holistic Management

**PULCEO** — Tempus fugit

pulceo-resource-manager (PRM)

pulceo-monitoring-service (PMS)

pulceo-service-manager (PSM)

pulceo-node-agent (PNA)

API Gateway

*RESTful via HTTP*

(b) Microservice-based Implementation

*realizes*

(CPU, Memory, Storage)

(4, 16, 32)

cloud1 — *USA Virginia*

<(82, 43), (83, 65)>   <(88, 39), (87, 65)>

*<(Latency, Bandwidth), (...)>*

(2, 8, 32) — *FRA Paris* fog1

(2, 4, 32) — fog2 *GER Frankfurt*

(c) Realization of Cloud-Edge Topologies

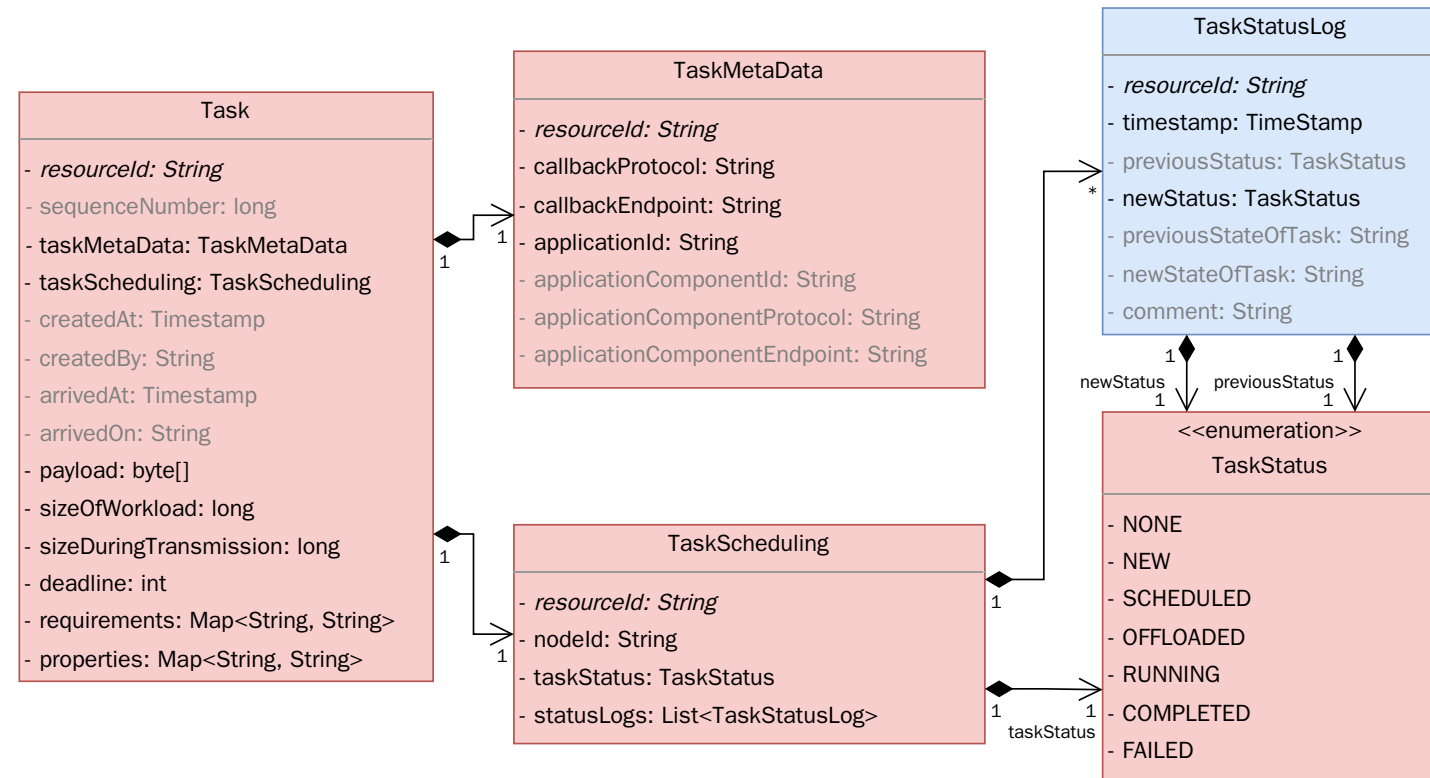https://github.com/spboehm/pulceo-misc

# Contributions

1. Task scheduling and offloading in line with holistic cloud-edge orchestration.

2. A general, universal, and flexible model of tasks.

3. Integration of task-processing applications.

4. Automated standard evaluation with metrics for task scheduling and offloading.

# Task: Domain Model

Extension of step **Operation**, supported by **p**ulceo-**s**ervice-**m**anager (**PSM**).

- *Task* with fixed / variable
  - requirements $R$
  - properties $P$
- *TaskMetaData*: Internal routing information
- *TaskScheduling*: Current state of scheduling
- *TaskStatusLog*: State transitions over time

**Task**
- *resourceId: String*
- sequenceNumber: long
- taskMetaData: TaskMetaData
- taskScheduling: TaskScheduling
- createdAt: Timestamp
- createdBy: String
- arrivedAt: Timestamp
- arrivedOn: String
- payload: byte[]
- sizeOfWorkload: long
- sizeDuringTransmission: long
- deadline: int
- requirements: Map<String, String>
- properties: Map<String, String>

**TaskMetaData**
- *resourceId: String*
- callbackProtocol: String
- callbackEndpoint: String
- applicationId: String
- applicationComponentId: String
- applicationComponentProtocol: String
- applicationComponentEndpoint: String

**TaskStatusLog**
- *resourceId: String*
- timestamp: TimeStamp
- previousStatus: TaskStatus
- newStatus: TaskStatus
- previousStateOfTask: String
- newStateOfTask: String
- comment: String

**TaskScheduling**
- *resourceId: String*
- nodeId: String
- taskStatus: TaskStatus
- statusLogs: List<TaskStatusLog>

**<<enumeration>> TaskStatus**
- NONE
- NEW
- SCHEDULED
- OFFLOADED
- RUNNING
- COMPLETED
- FAILED

**Example:** NONE -> NEW -> SCHEDULED -> OFFLOADED -> RUNNING -> COMPLETED
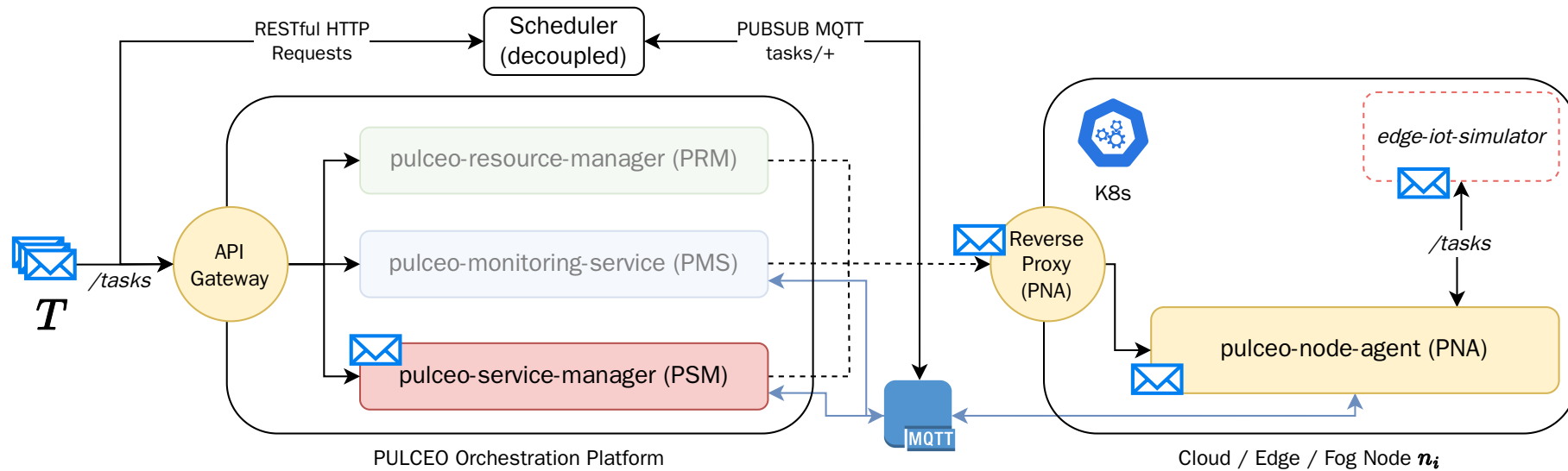
https://github.com/spboehm/pulceo-misc

# Orchestration Workflow

Experiment · Task Submission, Scheduling, Offloading, Processing, and Completion

# Experiment



- Decoupled scheduler uses HTTP (`/tasks`) and MQTT (`tasks/+`)
- $P_{task\_type}(t_i) = S\,(50\%), M\,(30\%), L\,(20\%) \qquad batch\_sizes = 200, \ldots, 400$
- $P_{size}(t_i) = 10, 50, 100 \quad \ldots \qquad R_{cpu\_shares}(t_i) = 250, 500, 1000 \quad \ldots$
- Instance of edge-iot-simulator[1] used for task processing

1. https://github.com/spboehm/edge-iot-simulator

https://github.com/spboehm/pulceo-misc

# Task Submission

**Task states:** NONE -> NEW -> SCHEDULED -> OFFLOADED -> RUNNING -> COMPLETED

**IoT device submits task:**

`HTTP POST /tasks`

```
 1  {
 2    "createdBy": "task_emitter.py",
 3    "sizeOfWorkload": 50,
 4    "sizeDuringTransmission": 50,
 5    "deadline": 50,
 6    "payload_length": 50,
 7    "payload": "payload",
 8    "cpu_shares": 500,
 9    "memory_size": 0.25,
10    "properties": {
11      "task_type": "small"
12    }
13  }
```
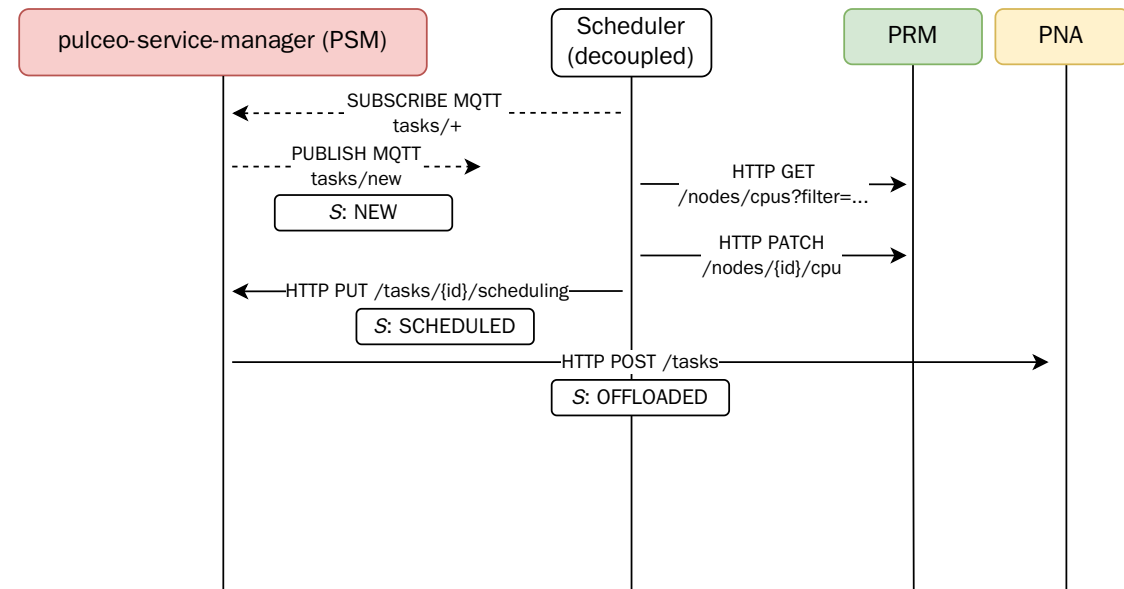
# Task Scheduling and Offloading

**Task states:** NONE -> NEW -> SCHEDULED -> OFFLOADED -> RUNNING -> COMPLETED

## Scheduler assigns node:

`HTTP PUT /tasks/{id}/scheduling`

```
 1  {
 2    "nodeId": "0b1c6697-...",
 3    "appId": "74e18419-...",
 4    "status": "SCHEDULED",
 5    "properties": {
 6      "batchSize": "100",
 7      "layer": "cloud-only",
 8      "policy": "dong et al."
 9    }
10  }
```

## PSM offloads task:

`HTTP POST /tasks` (on PNA)

```
 1  {
 2    "taskId": "831de05b-...",
 3    "appId": "74e18419-...",
 4    "appComponentId": "6a463804-...",
 5    ...
 6  }
```



Sequence diagram with participants: pulceo-service-manager (PSM), Scheduler (decoupled), PRM, PNA

- SUBSCRIBE MQTT tasks/+
- PUBLISH MQTT tasks/new — *S*: NEW
- HTTP GET /nodes/cpus?filter=...
- HTTP PATCH /nodes/{id}/cpu
- HTTP PUT /tasks/{id}/scheduling — *S*: SCHEDULED
- HTTP POST /tasks — *S*: OFFLOADED

https://github.com/spboehm/pulceo-misc

# Task Processing

**Task states:** NONE -> NEW -> SCHEDULED -> OFFLOADED -> RUNNING -> COMPLETED
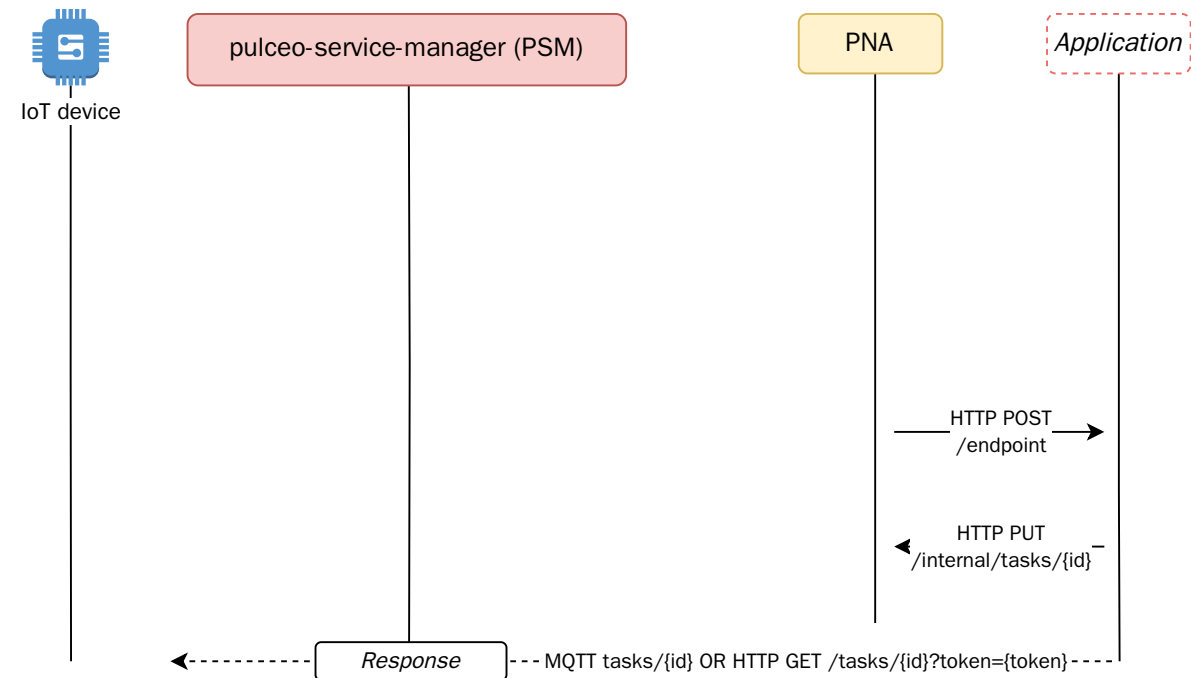
## Forward task to application:

`HTTP POST /endpoint`

```
1  {
2    "taskId": "831de05b-...",
3    "callbackProtocol": "0b1c6697-...",
4    "callbackEndpoint": "74e18419-...",
5    ...
6  }
```

## Update task status internally:

`HTTP PUT /internal/tasks/{id}`

```
1  {
2    "taskId": "831de05b-...",
3    "newTaskStatus": "COMPLETED"
4  }
```

IoT device

pulceo-service-manager (PSM)

PNA

*Application*

HTTP POST
/endpoint

HTTP PUT
/internal/tasks/{id}

*Response* ---- MQTT tasks/{id} OR HTTP GET /tasks/{id}?token={token} ----

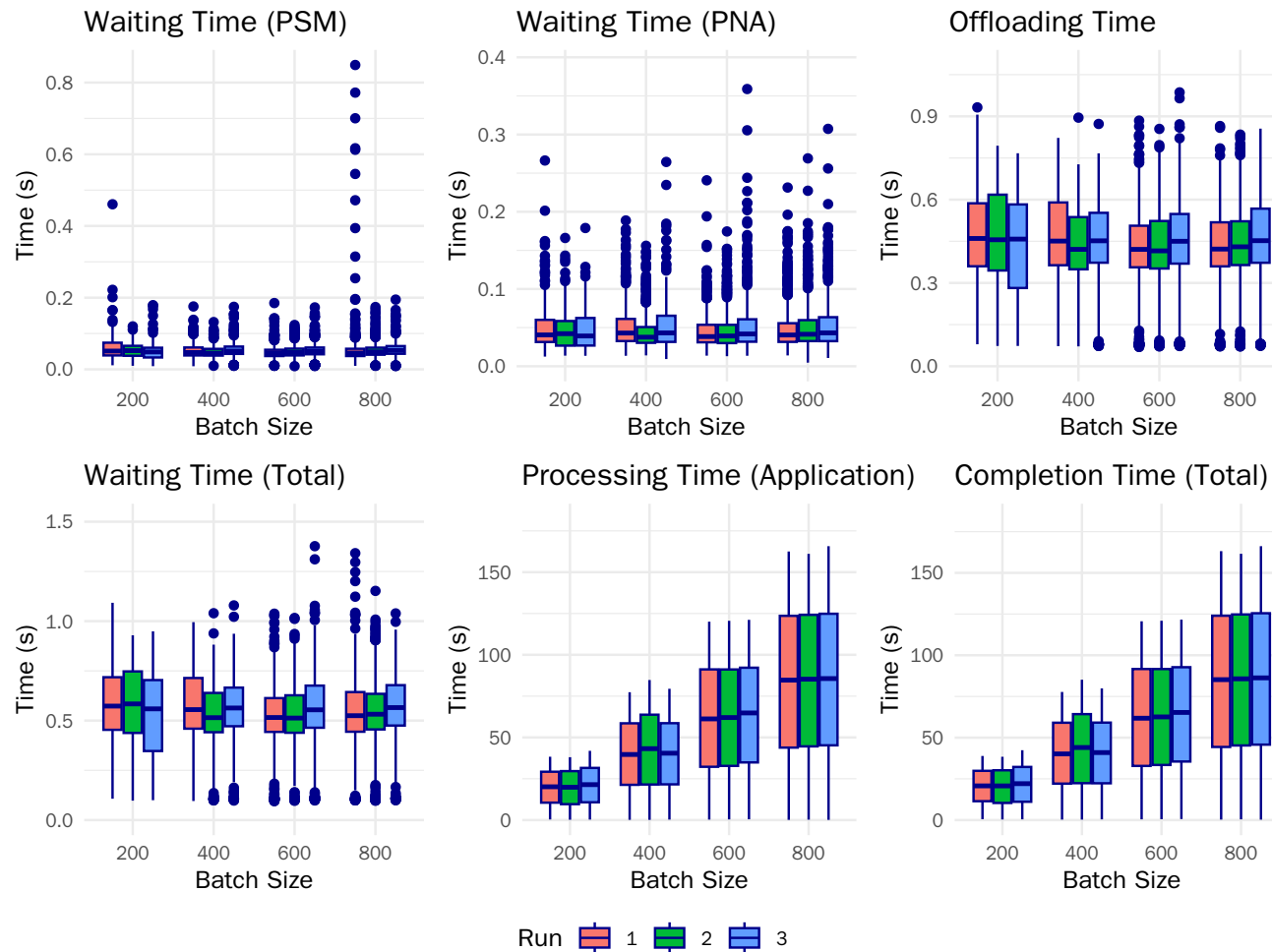https://github.com/spboehm/pulceo-misc

# Evaluation

Task Metrics · Performance Metrics

- Provided by the automatically generated orchestration reports

- On top of already implemented metrics, like CPU, memory, and network-related metrics

# Task Metrics



Drill-down into overall completion time:

- Scheduling (PSM): $0.05\text{s}$

- Waiting (PNA): $0.05\text{s}$

- Offloading: $0.44\text{s}$

- Waiting (Total): $0.54\text{s}$

Overhead:

- Communication (HTTP)

- Telemetry

- Proxied application

https://github.com/spboehm/pulceo-misc

# Performance Metrics

- **Task Response Time**: Time between submitting the task and receiving the response

- **Task Arrival Rate**: Tasks/s arriving at PSM

- **Task Throughput**: Tasks/s completed by the system

# Conclusion & Limitations

## Contributions

1. 🚀 Task scheduling and offloading in line with holistic cloud-edge orchestration
2. 🌐 A general, universal, and flexible model of tasks
3. ✅ Integration of task-processing applications
4. 🔄 Automated standard evaluation with metrics for task scheduling and offloading

Empirical evaluation with first insights on the overall overhead.

## Limitations

- 🚨 Integration of task-processing applications only via HTTP and MQTT
- 🚨 Centralized, cloud-based, and only partly decentralized orchestration
- 🚨 Minimal cloud-edge topology for first evaluation