



GEdge Platform

GEdge(Griffin-Edge) Platform

- 초저지연 지능형 클라우드 엣지 SW 플랫폼 -

# 강화학습 기반 멀티 엣지 협업 정책 생성 기술

2023.12.07

GS-Link 프레임워크 코어개발자 (GS-LinkHQ)

윤주상 (joosang.youn@gmail.com)

“**GEdge Platform**”은 클라우드 중심의 엣지 컴퓨팅 플랫폼을 제공하기 위한  
핵심 SW 기술 개발 커뮤니티 및 개발 결과물의 코드명입니다.

- New Leap Forward of

**GEdge Platform Community 7<sup>th</sup> Conference** (GEdge Platform v4.0 Release) -

# Contents

---

- I 강화학습 기반 정책 생성 기술**
- II 지능형 오프로딩 정책 생성**
- III 지능형 서비스 이동 정책**
- IV 지능형 캐싱 정책 생성**

# GEdge 플랫폼 내 LinkHQ의 역할

## 초저지연 지능형 클라우드 엣지 플랫폼 (GEdge Platform)

### 클라우드 엣지 관리 플랫폼 (GM : GEdge Management Platform)

플랫폼 관리 도구 프레임워크 (GM-Tool)

Framework I/F

플랫폼 관리 기능 프레임워크 (GM-Center)

Platform I/F

### 지능형 서비스 운용 프레임워크 (GS-AI)

엣지 AI 서비스 환경  
(GS-AIflow)

엣지 협업 학습 환경  
(GS-Optops)

### 서비스 협업 프레임워크 (GS-Link)

협업 게이트웨이  
(GS-Linkgw)

협업 정책 생성  
(GS-Linkhq)

Framework I/F

Framework I/F

### 초저지연 데이터 처리 프레임워크 (GS-Engine)

엣지 전용 스케줄러  
(GS-Scheduler)

엣지 메시지 브로커  
(GS-Broker)

### 클라우드 엣지 서비스 플랫폼 (GS : GEdge Service Platform)

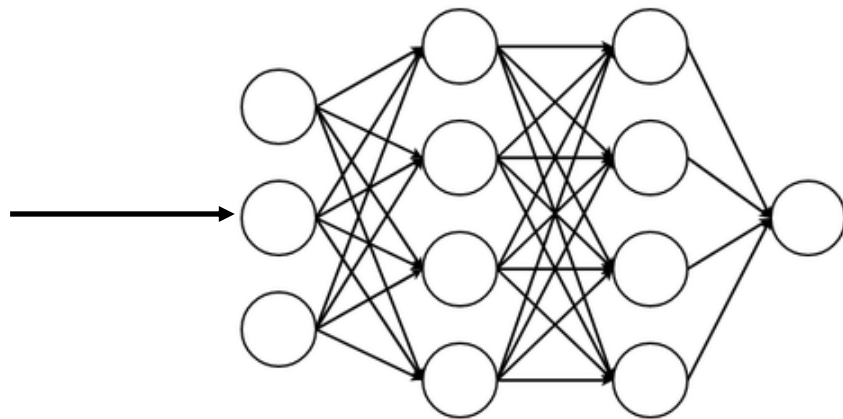


# 강화학습 기반 지능형 오프로딩 정책 생성 기술



- 기존 자원 할당 정책
  - Random, Least Load, Round-Robin
  - Rule Based
- 심층강화학습 기반 자원 할당 정책

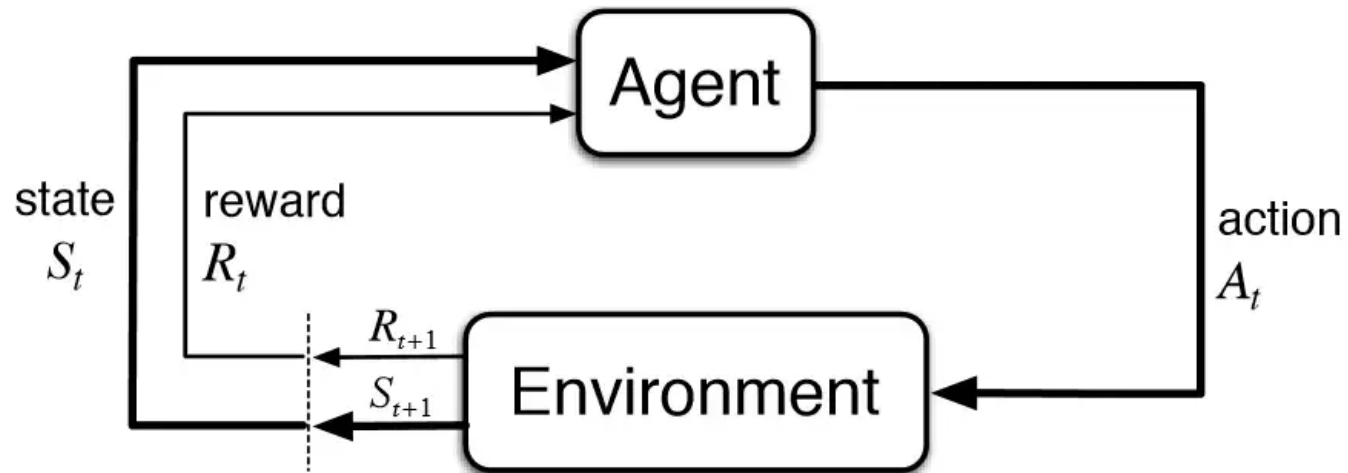
자원 할당 요청  
(태스크 요구사항)



최적 엣지 자원 추천

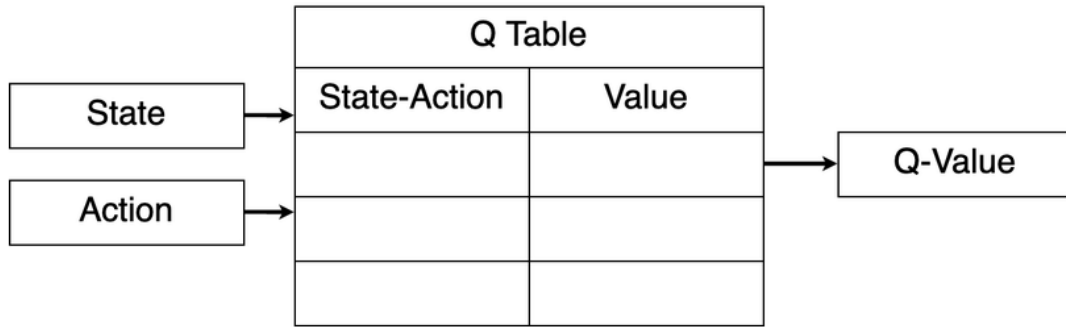
- 강화학습

- 에이전트가 환경과 상호작용하며 해당 경험을 통해 특정 상태에서 **최적 행동을 학습**하는 기법

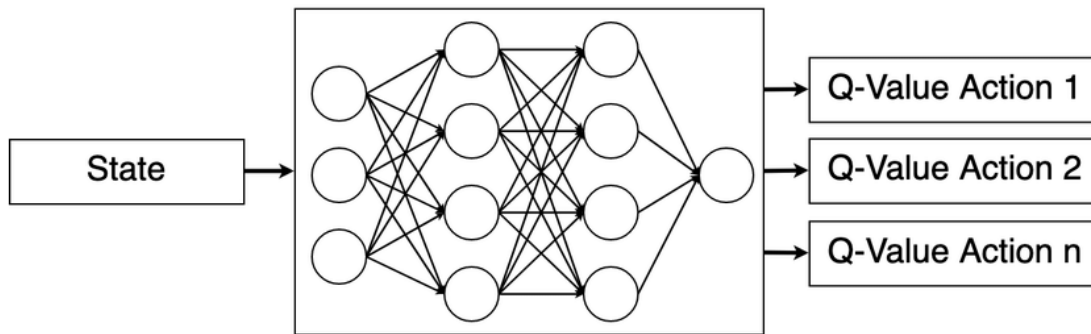


# 3 Deep Q-Learning (DQN)

## Q-Learning



## Deep Q-Learning

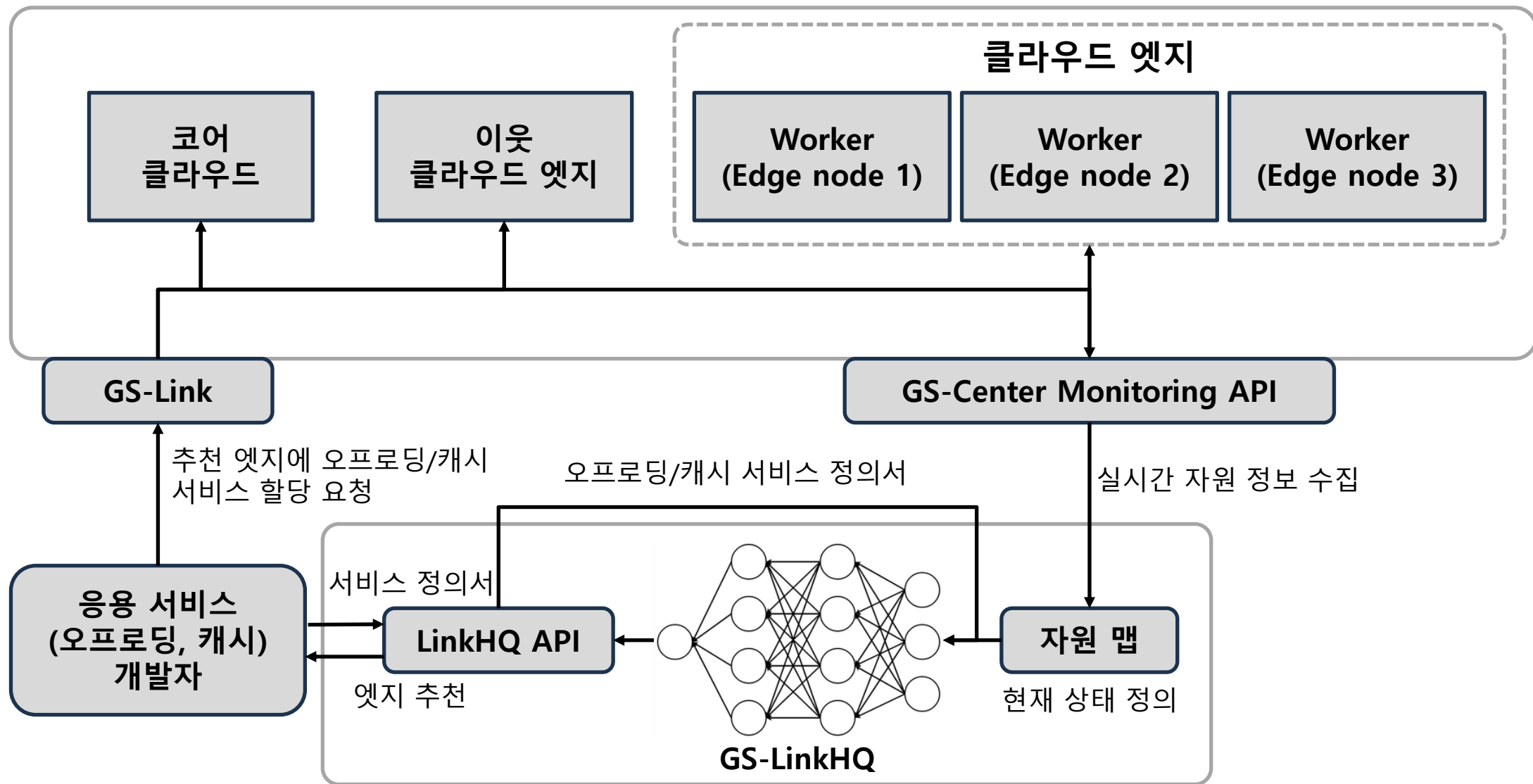


### • Q-Learning

- Lookup-table 기반 강화학습 기법
- Q table에 에이전트의 경험과 가치를 저장
- 예상 보상이 최대가 되는 행동을 반환하도록 학습
- 환경이 복잡하면 효율이 떨어짐

### • Deep Q-Learning

- Q table 대신 신경망을 통해 Q-value를 근사
- Replay memory를 통해 경험의 효율을 높임
- 정책을 학습하는 네트워크와 액션을 추론하는 네트워크를 분리하여 학습 안정성 향상
- 복잡한 환경에서도 안정적인 정책 학습 가능





**GS-LinkHQ**  
Offloading Cluster Recommendation System

**Recommendation Result**

Cluster	<b>gm-cluster</b>
CPU	4
Memory	16
Disk	32

[Back](#)

LinkHQ Result Page

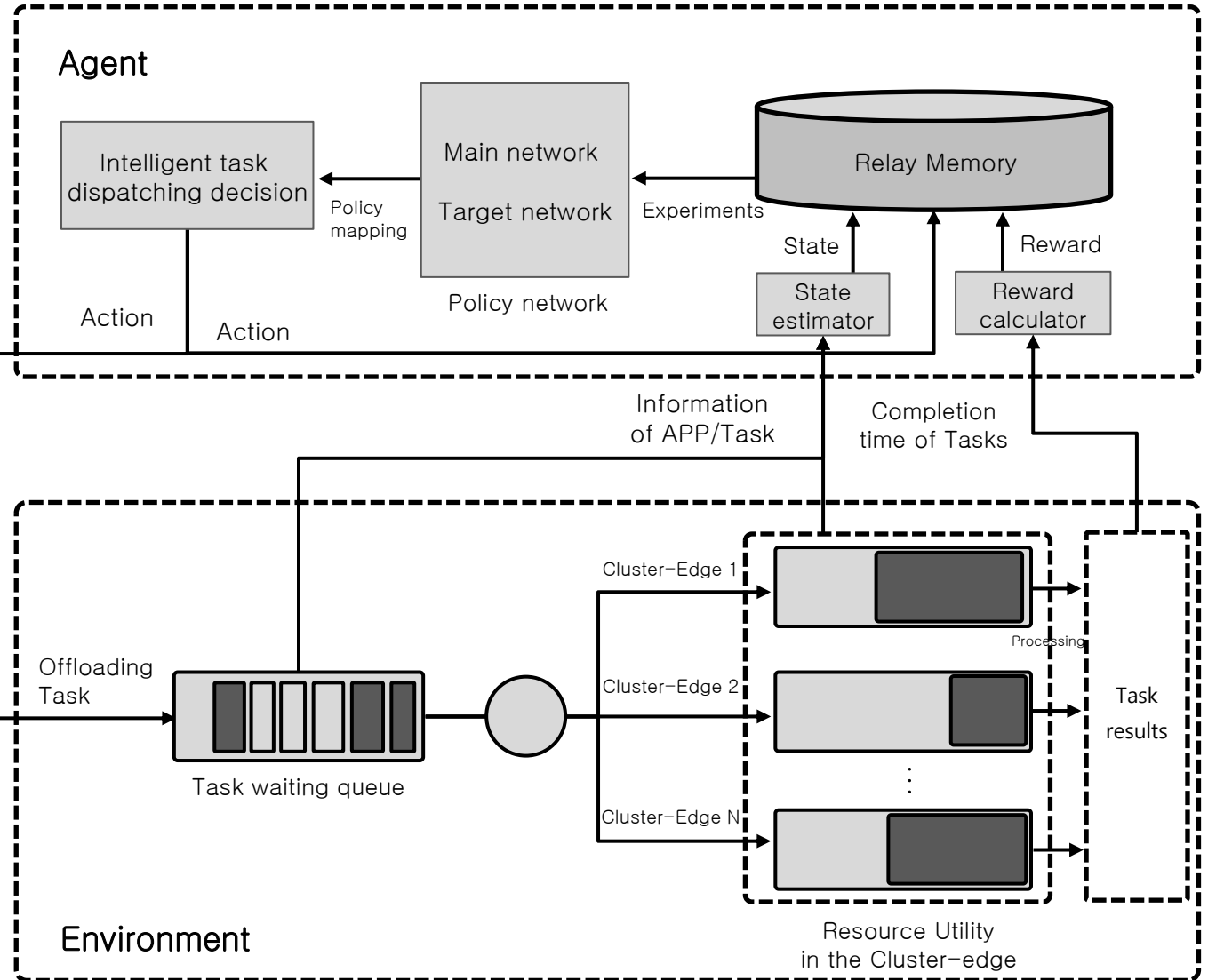
**GS-LinkHQ**  
Offloading Cluster Recommendation System

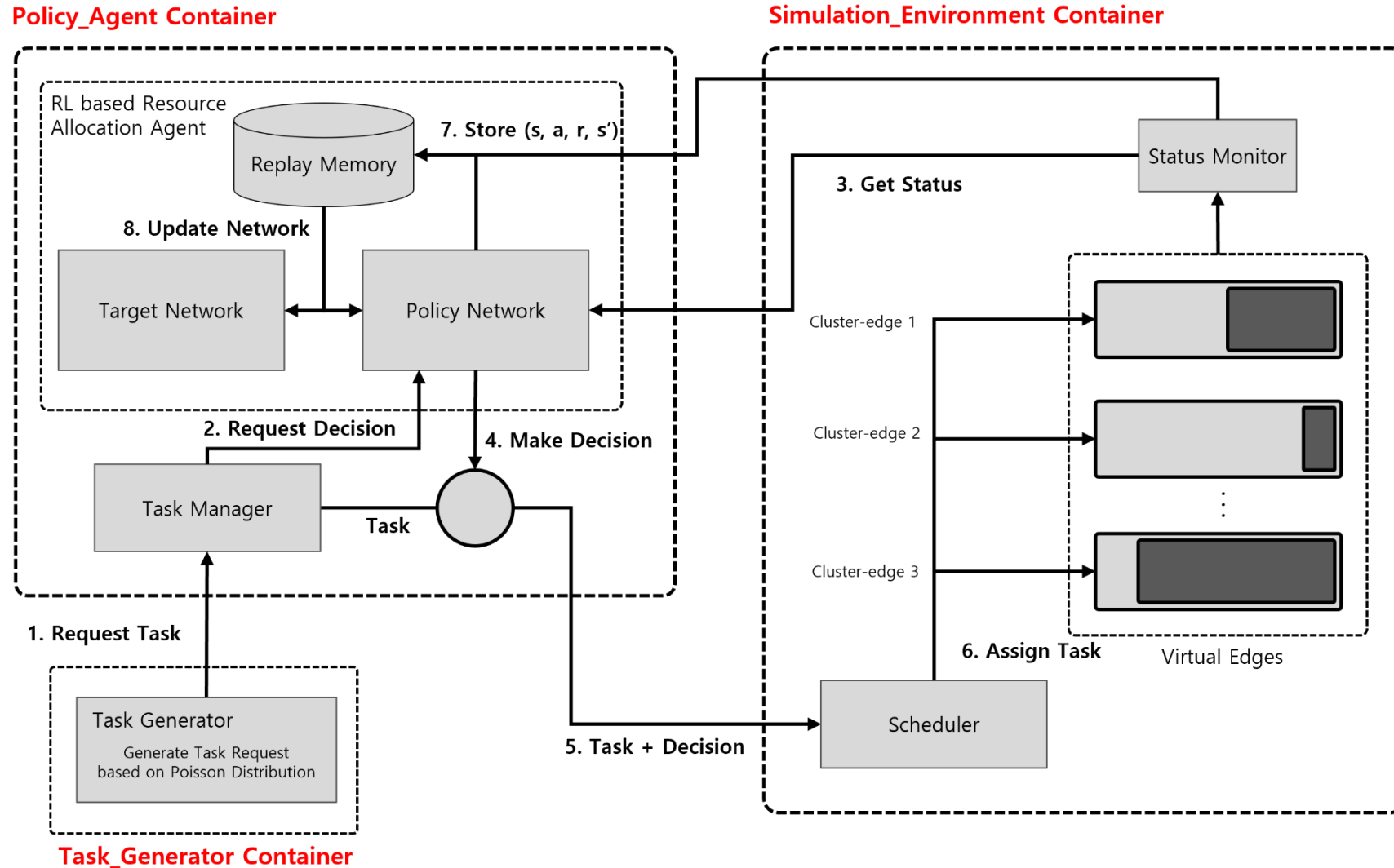
**Enter Task Information**

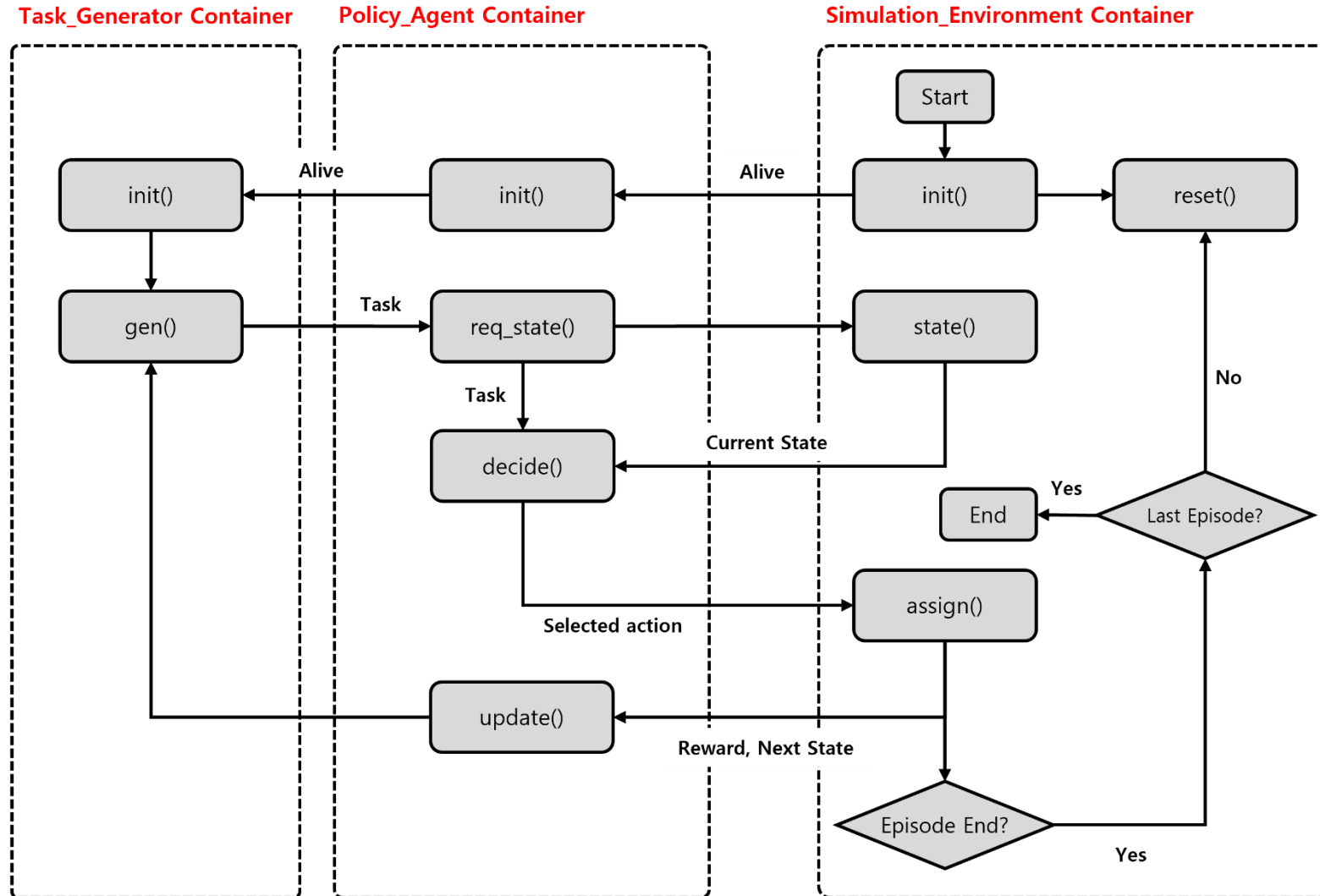
CPU (Cores)	<input type="text" value="4"/>
Memory (GB)	<input type="text" value="16"/>
Disk (GB)	<input type="text" value="32"/>
GPU (WIP)	<input type="text" value="0"/>

[Submit](#)

LinkHQ Request Page









# 지능형 오프로딩 정책 생성



# 1 오프로딩 시스템 수식화

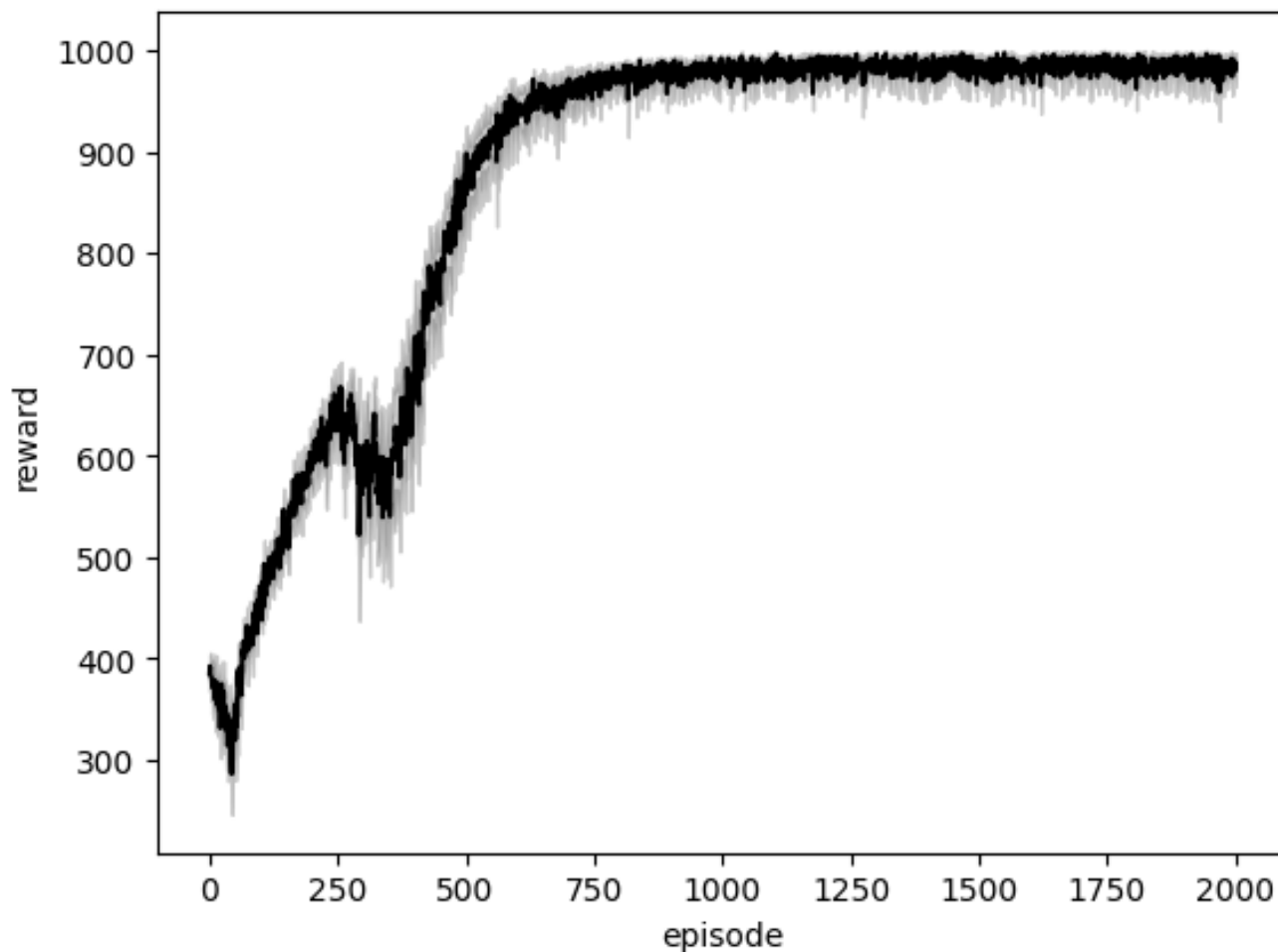
Notation	Definition
$n$	Number of edge servers
$e_k$	$k$ -th edge server
$e_k^{CPU}$	Available CPU resource of $e_k$
$e_k^{MEM}$	Available memory resource of $e_k$
$e_k^{GPU}$	Available GPU resource of $e_k$
$t_i$	$i$ -th requested task
$t_i^e$	Index of edge requested by Task $t_i$
$t_i^{index}$	Flag indicating whether $t_i$ is requested to $e_k$ (one hot encoded)
$t_i^{CPU}$	Required CPU of $t_i$
$t_i^{MEM}$	Required memory of $t_i$
$t_i^{GPU}$	Required GPU of $t_i$
$t_i^T$	Flag indicating whether $t_i$ is time sensitive task
$S_i$	State set of system when task $i$ requested
$A_i$	Decided action set of requested task $i$
$\alpha_i^{index}$	Index of edge assigned Task $t_i$
$\alpha_i^k$	Flag indicating whether $t_i$ is assigned to $e_k$ (one hot encoded)
$R_i$	Reward of $A_i$

- State ( $S_i$ ) : 태스크 정보 + 엣지 상태
- Action ( $A_i$ ) : 오프로딩 결정 엣지
- Reward ( $R_i$ ) : Action에 대한 평가 점수

$$S_i = \{t_i^0, t_i^1, \dots, t_i^n, t_i^{CPU}, t_i^{MEM}, t_i^{GPU}, t_i^T, e_0^{CPU}, e_0^{MEM}, e_0^{GPU}, \dots, e_n^{CPU}, e_n^{MEM}, e_n^{GPU}\}$$

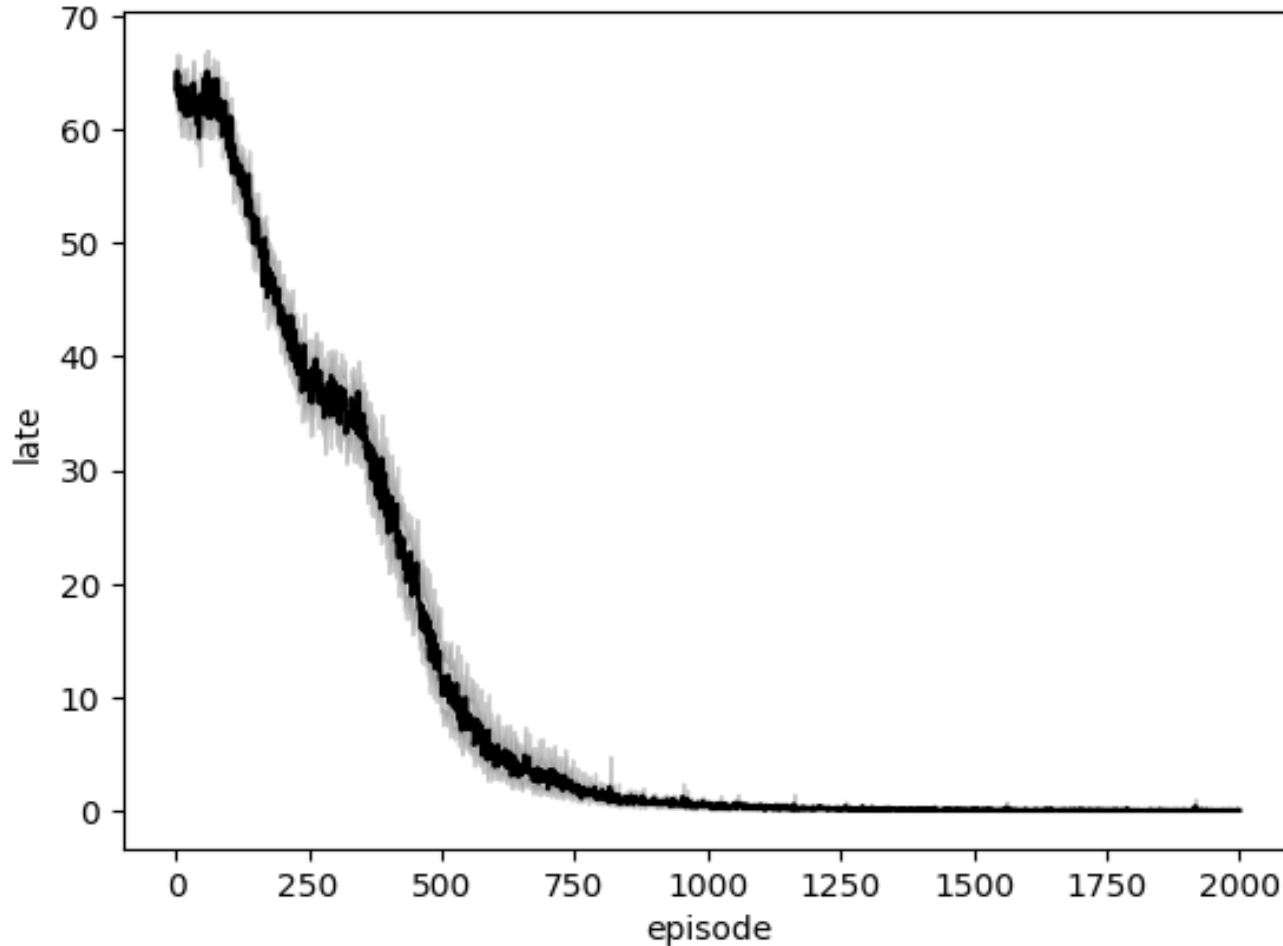
$$A_i = \{\alpha_i^0, \alpha_i^1, \alpha_i^2, \dots, \alpha_i^n\}$$

$$R_i = \begin{cases} 10 & \text{allocated and } t_i^T = 0 \\ 10 & \text{allocated and } t_i^T = 1 \text{ and } t_i^{index} = \alpha_i^{index} \\ 1 & \text{allocated and } t_i^T = 1 \text{ and } t_i^{index} \neq \alpha_i^{index} \\ -10 & \text{rejected} \end{cases}$$



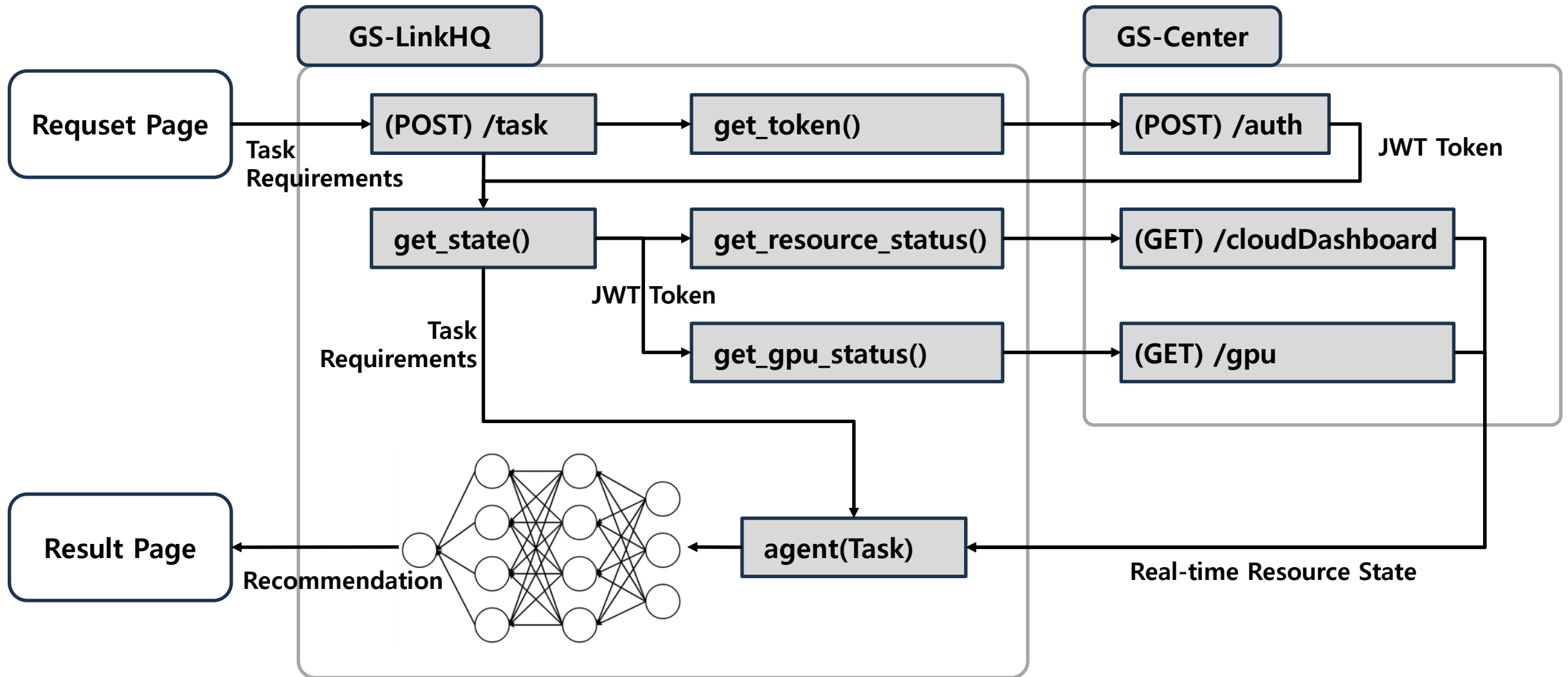
- Reward

- 각 에피소드에서 reward 총 합
- 800 에피소드 이후로 시뮬레이터 최대 reward = 1000 으로 수렴



- Late Count

- 각 에피소드에서 시간민감형 태스크가 멀리 떨어진 엣지에 할당된 횟수
- 800 에피소드 이후로 0건으로 수렴





# 5 오프로딩 엣지 추천 서비스

## GS-LinkHQ

### Offloading Cluster Recommendation System

Enter Task Information

CPU (Cores)

8

Memory (GB)

32

Disk (GB)

28

GPU (Ea)

2

Time Sensitive

☒

gm-cluster

Submit

Recommendation Request Page

## GS-LinkHQ

### Offloading Cluster Recommendation System

Recommendation Result

Cluster

ai-networklab(pangyo)

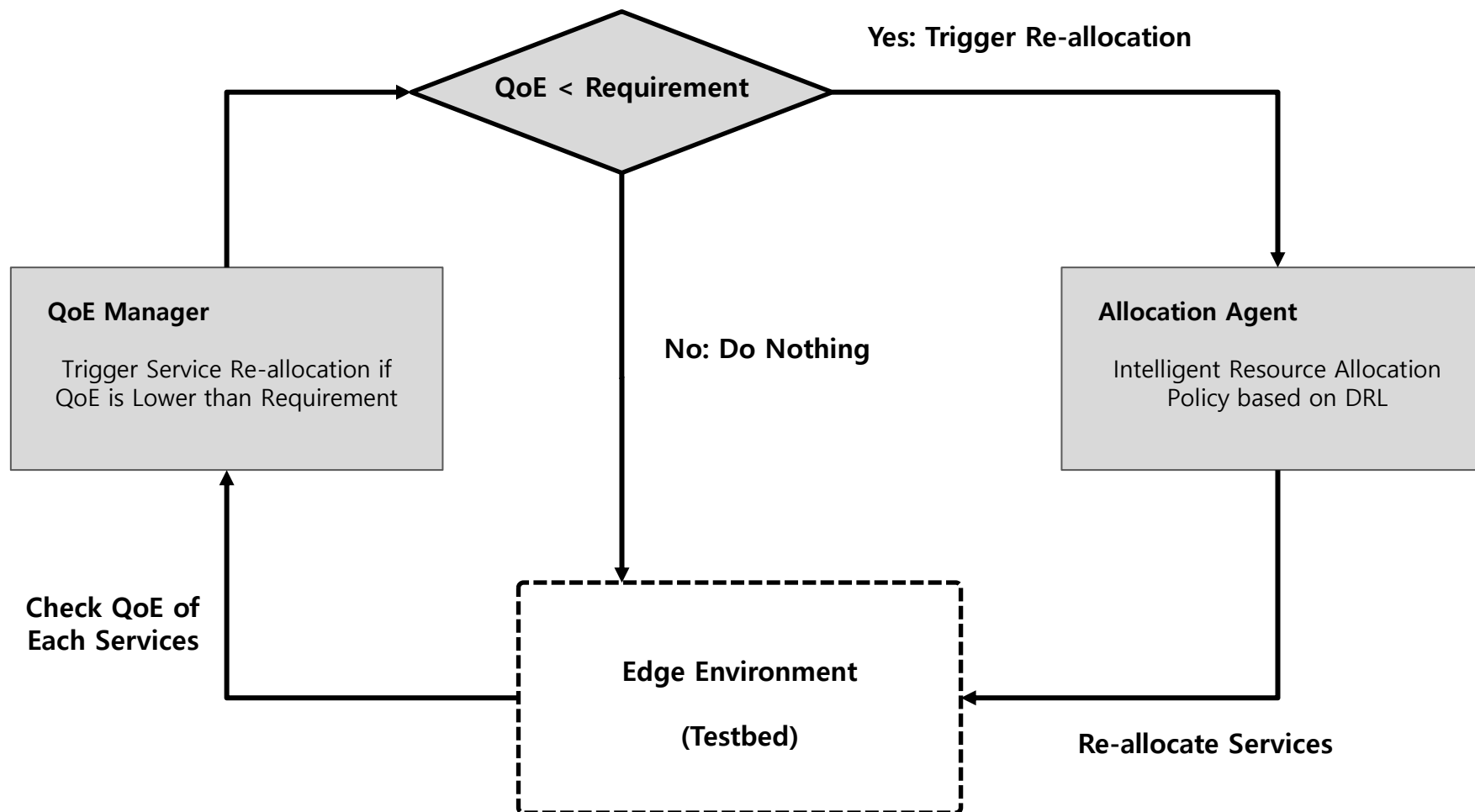
Back

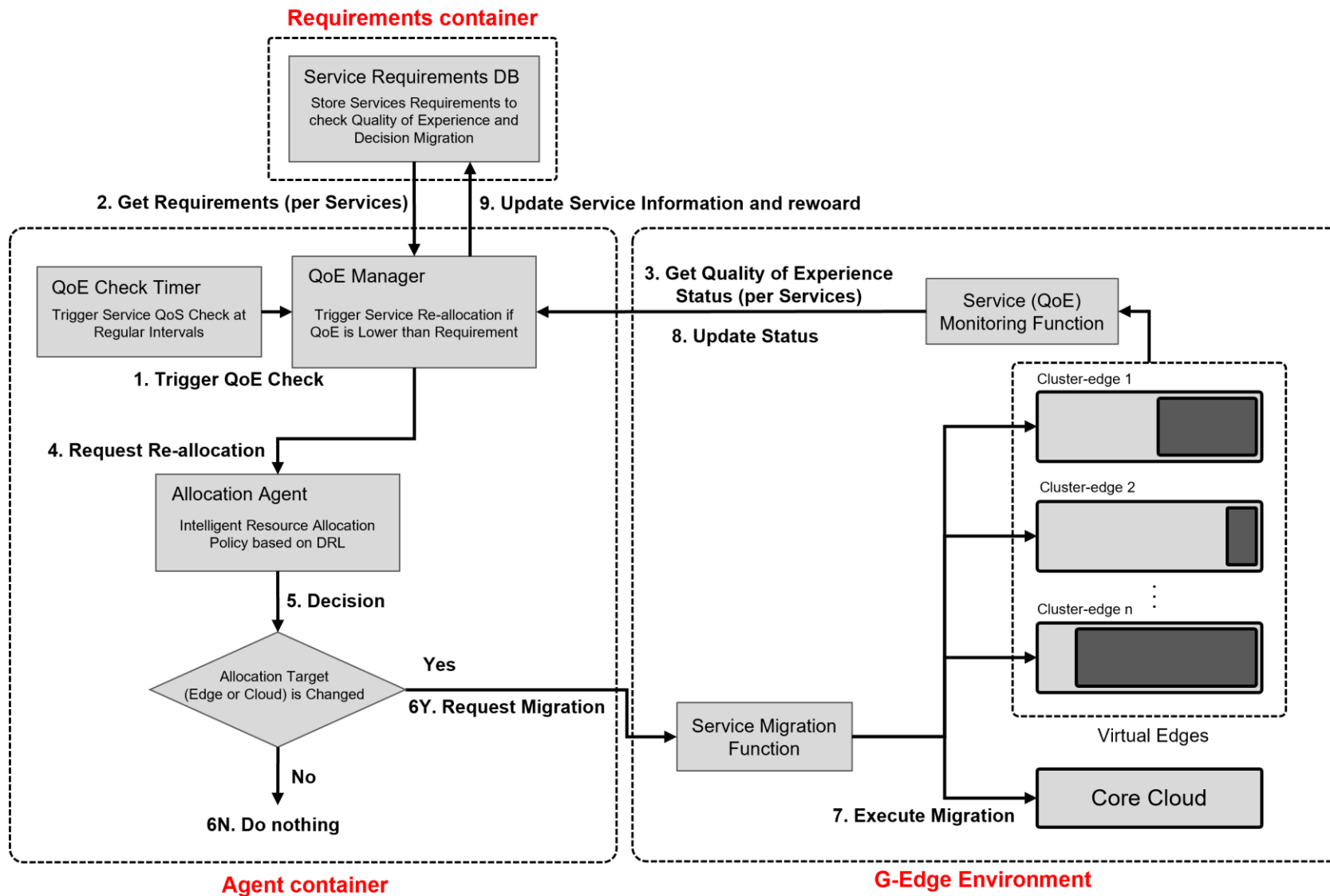
Recommendation Result Page



# 지능형 서비스 이동 정책







```

...
class QoEManager:
    def __init__(self, env, db, agent):
        self.db = db
        self.env = env
        self.agent = agent

    def qoe_check(self):
        requirements = self.db.get_requirements()
        qoe_status = self.env.get_qoe()

        for service in qoe_status:
            if service['requirement'] > requirements[service['task_id']]:

    def request_allocation(self, service):
        target = self.agent.allocation(service['task_id'])

        if target != service['target']:
            res = self.env.migration(service['task_id'], target)

            if res:
                self.update_db(service, target)
                self.agent.update_reward()

    def update_db(self, service, target):
        self.db.hset(service['task_id'], 'target', target)
...

```

서비스 요구사항 만족여부 확인

서비스 자원 재할당 요청

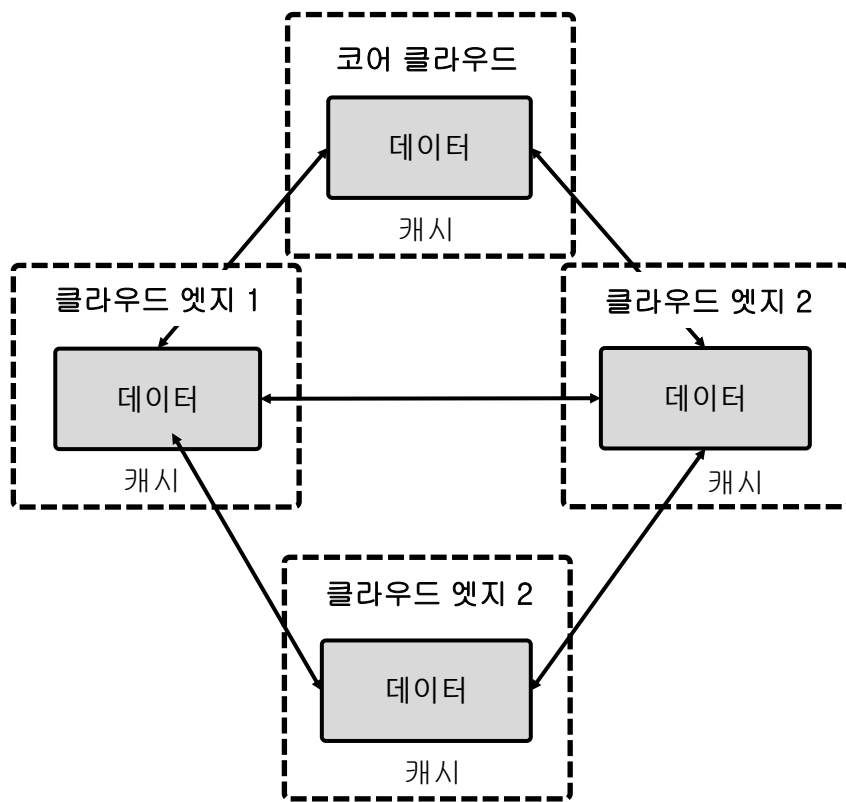
서비스 DB 및 리워드 업데이트

# IV

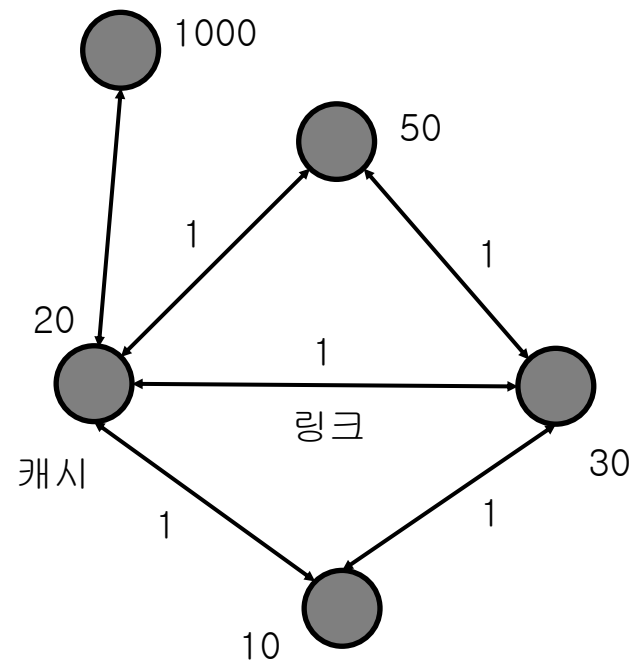
## 지능형 캐싱 정책 생성



- 협업 캐시 물리적 네트워크

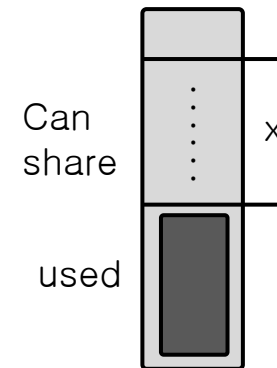


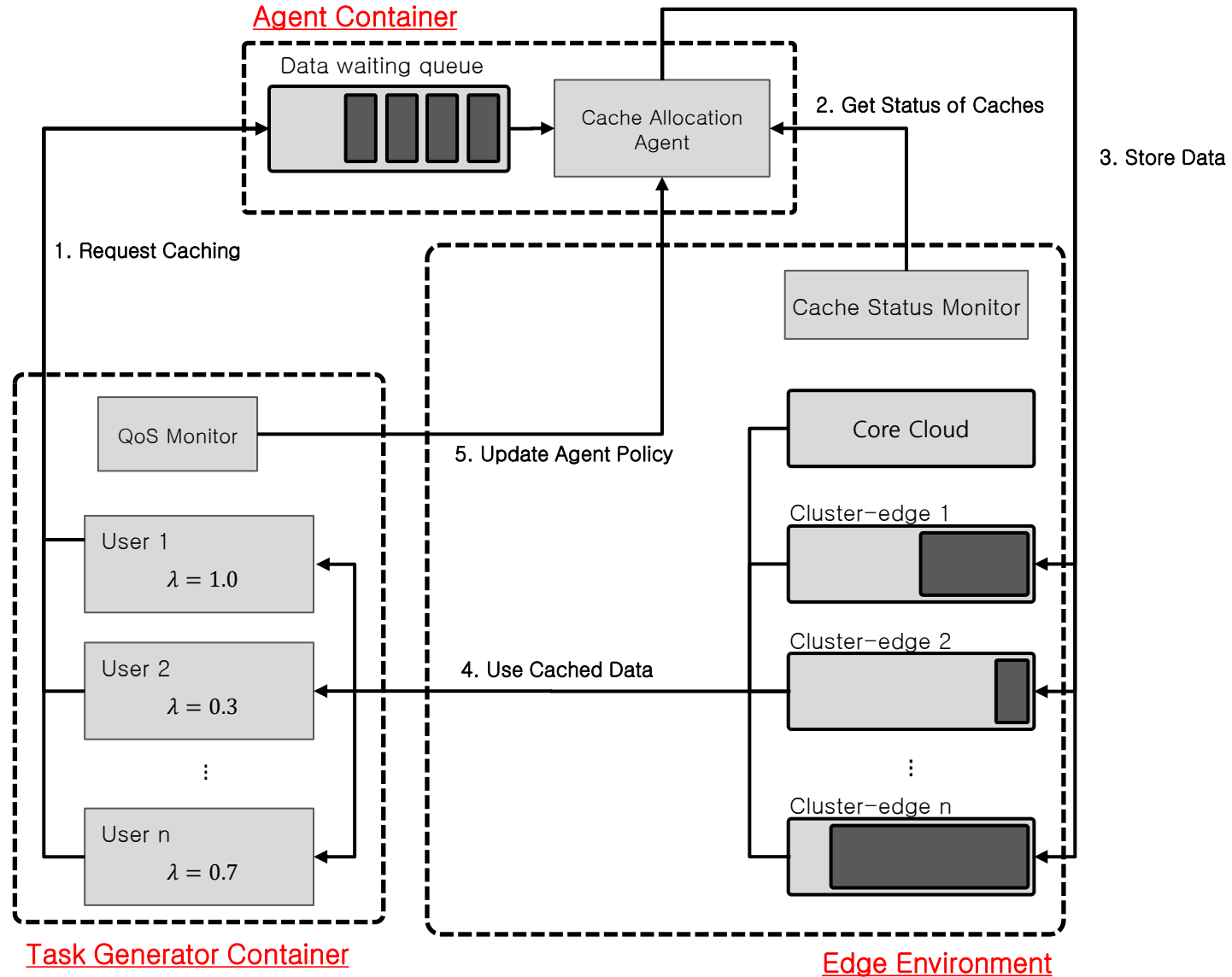
물리적 캐시 네트워크 모델



그래프 기반 캐시 네트워크

캐시 자원  
정의 (x-level)







## Cache Assign

### Endpoint

`http://<agent>/cache`

### Method

`post`

### Contents

- `cache`
  - `size`: Size of cache (MB)
  - `deadline`: Deadline of cache(ms)
  - `tts`: Time to survive of service (min)
  - `rts`: Max request number of service

## Request

Request agent to assign a cache

### Body Format

```
{
  "cache": {
    "size": <size>,
    "deadline": <deadline>,
    "tts": <tts>,
    "rts": <rts>
  }
}
```

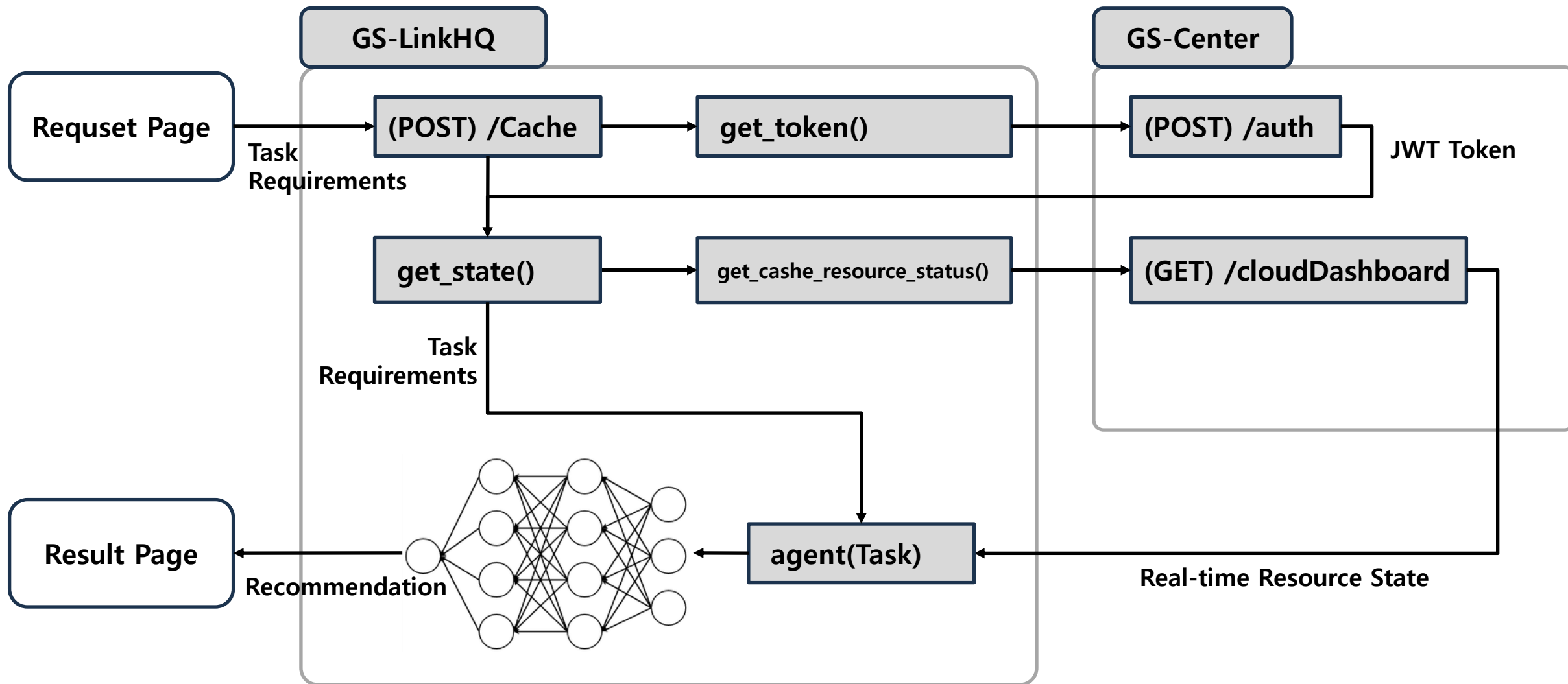
### Sample

```
POST /cache HTTP/1.1
Host: <agent>
Content-Type: application/json
{
  "cache": {
    "size": 5,
    "deadline": 0.1,
    "tts": 1440,
    "rts": 100
  }
}
```

## Response

201

```
HTTP/1.1 201 Created
Location: /cache/<cache-id>
Content-Type: application/json
{
  "id": <cache-id>,
  "message": "Successfully created"
}
```



## GS-LinkHQ

## Caching Cluster Recommendation System

## Enter Task Information

Disk (GB)

30

Time Sensitive



gm-cluster

Submit



## GS-LinkHQ

## Caching Cluster Recommendation System

## Recommendation Result

Cluster

ai-networklab(pangyo)

Back



# 감사합니다.

<http://gedge-platform.github.io>



GS-Link 프레임워크 코어개발자(GS-LinkHQ)

윤주상 (joosang.youn@gmail.com)

## Welcome to GEdge Platform

An Open Cloud Edge SW Platform to enable Intelligent Edge Service

### GEdge Platform will lead Cloud-Edge Collaboration