



GEdge(Griffin-Edge) Platform

- 초저지연 지능형 클라우드 엣지 SW 플랫폼 -

딥러닝을 활용한 클라우드 엣지 수직/수평 협업

2020.12.10

GEdge Platform 코어 개발자
윤주상(jsyoun@deu.ac.kr)

“The First talk of Edge Computing with Clouds”

- GEdge Platform 커뮤니티 멤버들의 첫번째 이야기 -

GEdge Platform Community 1st Conference

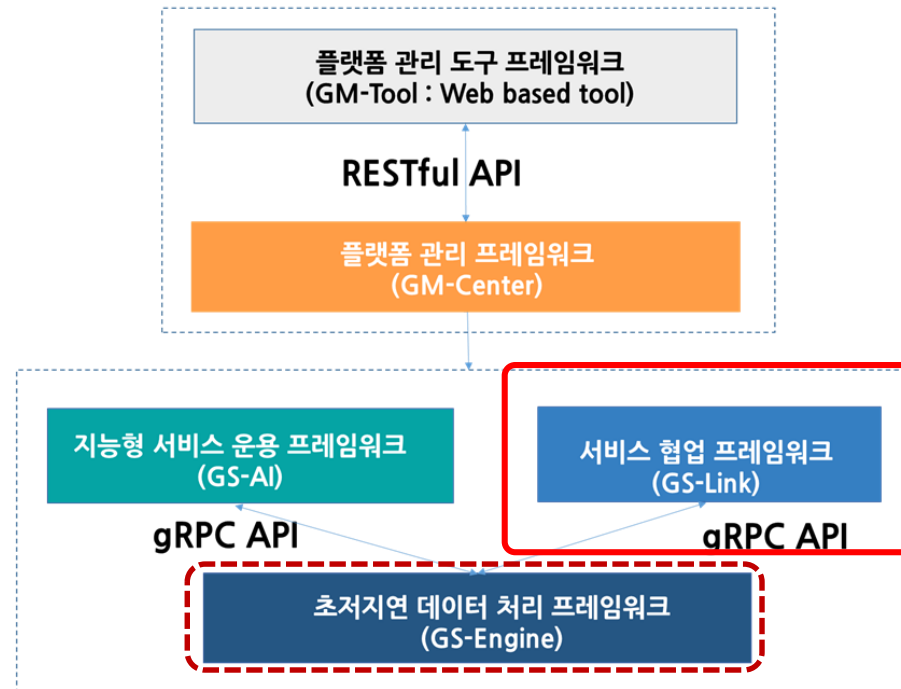
Contents

- I** 협업 모델 정의
- II** GEdge 플랫폼 내 서비스(오프로딩) 정의
- III** GEdge 플랫폼 내 지능형 협업 서비스 기술
- IV** 21년도 개발 계획

클라우드 엣지 수직/수평 협업 모델

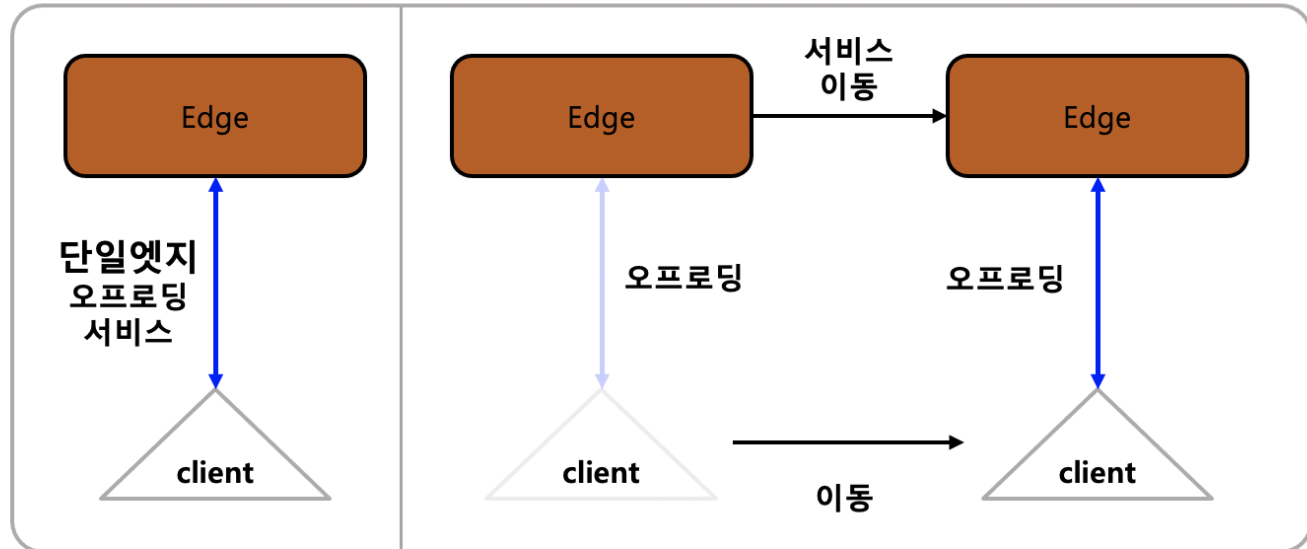
초저지연 지능형 클라우드 엣지 플랫폼 (GEdge Platform)

초저지연 클라우드 엣지 관리 플랫폼 (GM : GEdge Management)

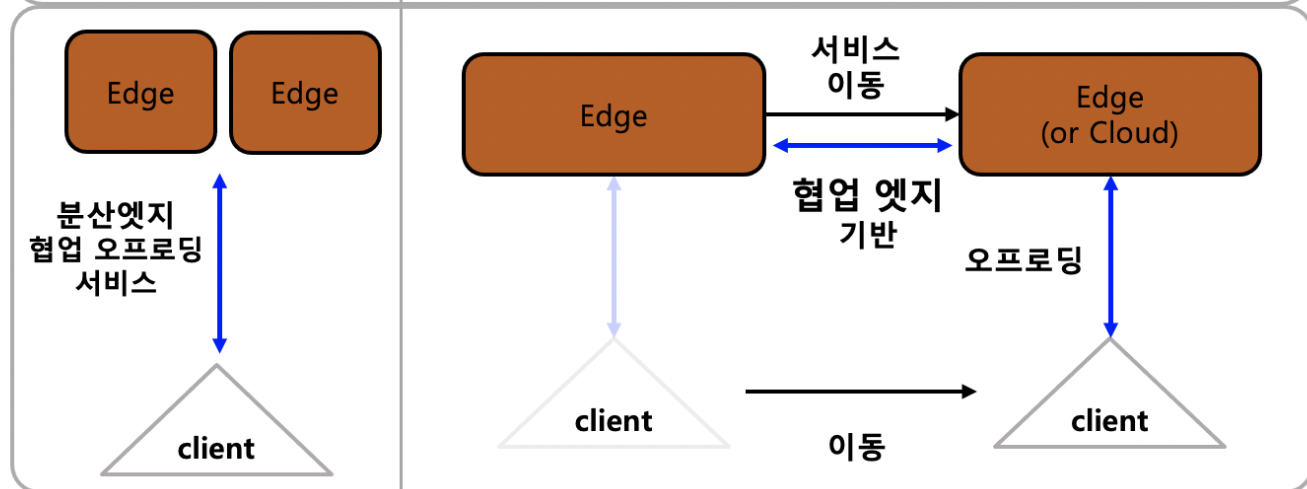


1 기존 엣지 모델 vs. 협업 엣지 모델

기존
엣지



협업
엣지

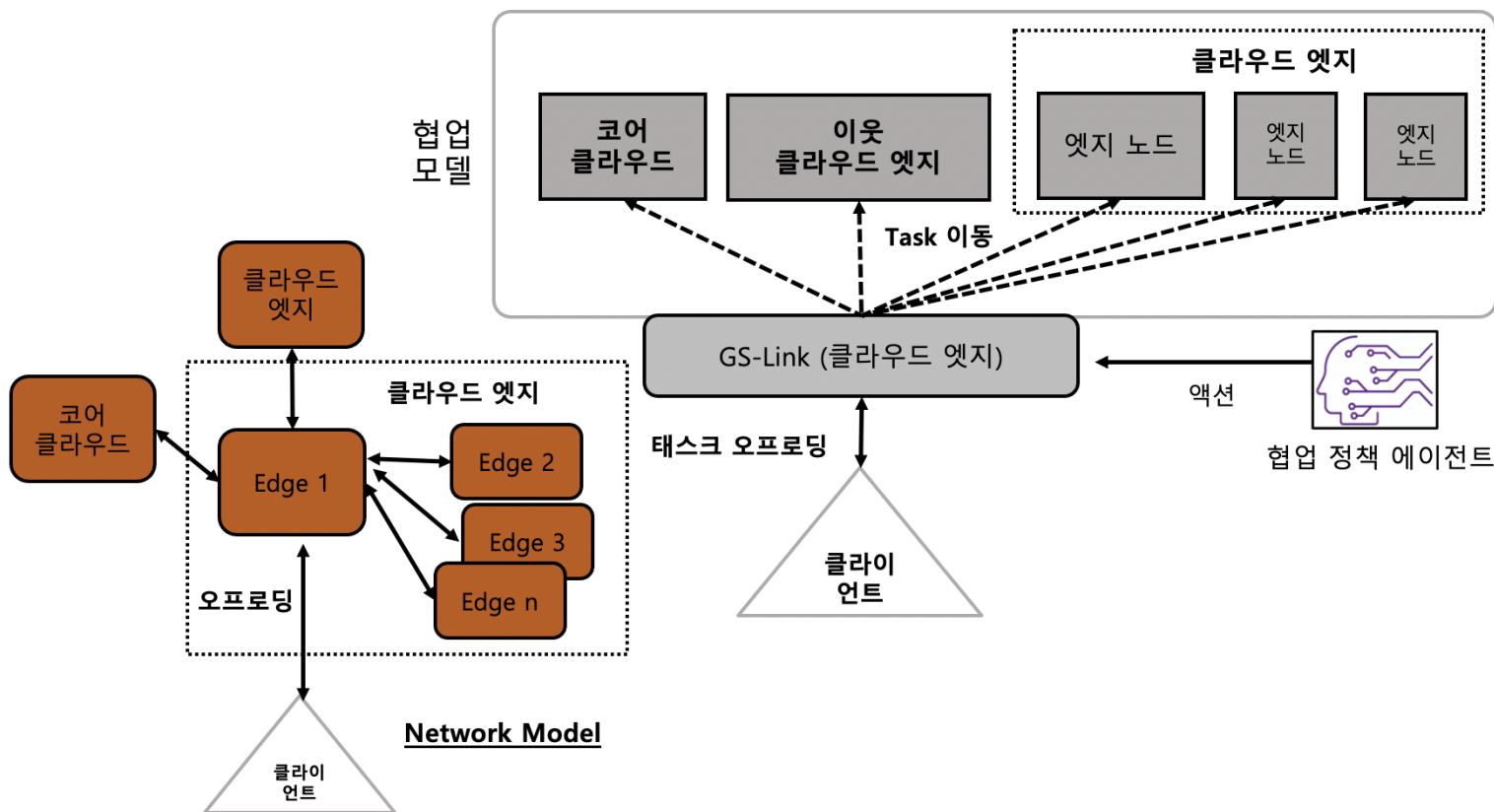


- 일반적인 엣지 오프로딩 서비스
 - 디바이스가 가진 응용을 엣지에서 대신 실행하고 그 결과를 리턴 받는 서비스
 - 디바이스는 오프로딩 요청 시 실행 코드, 실행 시 필요한 데이터 등을 첨부하여 서비스를 요청
- 제안된 GEdge 오프로딩 서비스
 - 메인 응용 서비스 : 엣지 기반 인공지능 서비스 (추론 서비스)
 - GE 플랫폼에서는 3가지 응용 서비스 분류 모델을 지원
 - 오프로딩 요청 시 응용 별 별도의 오프로딩 서비스 제공
 - 3가지 협업 오프로딩 제공
 - 차별화된 서비스 모델 : 협업 오프로딩 (GEdge 오프로딩)
 - 엣지 간, 엣지-클라우드 간 컴퓨팅 자원을 공유하는 협업 오프로딩 서비스. 여기서 협업 오프로딩 서비스는 엣지 간 수평적 오프로딩 서비스, 엣지-클라우드 간 수직적 오프로딩 서비스로 정의

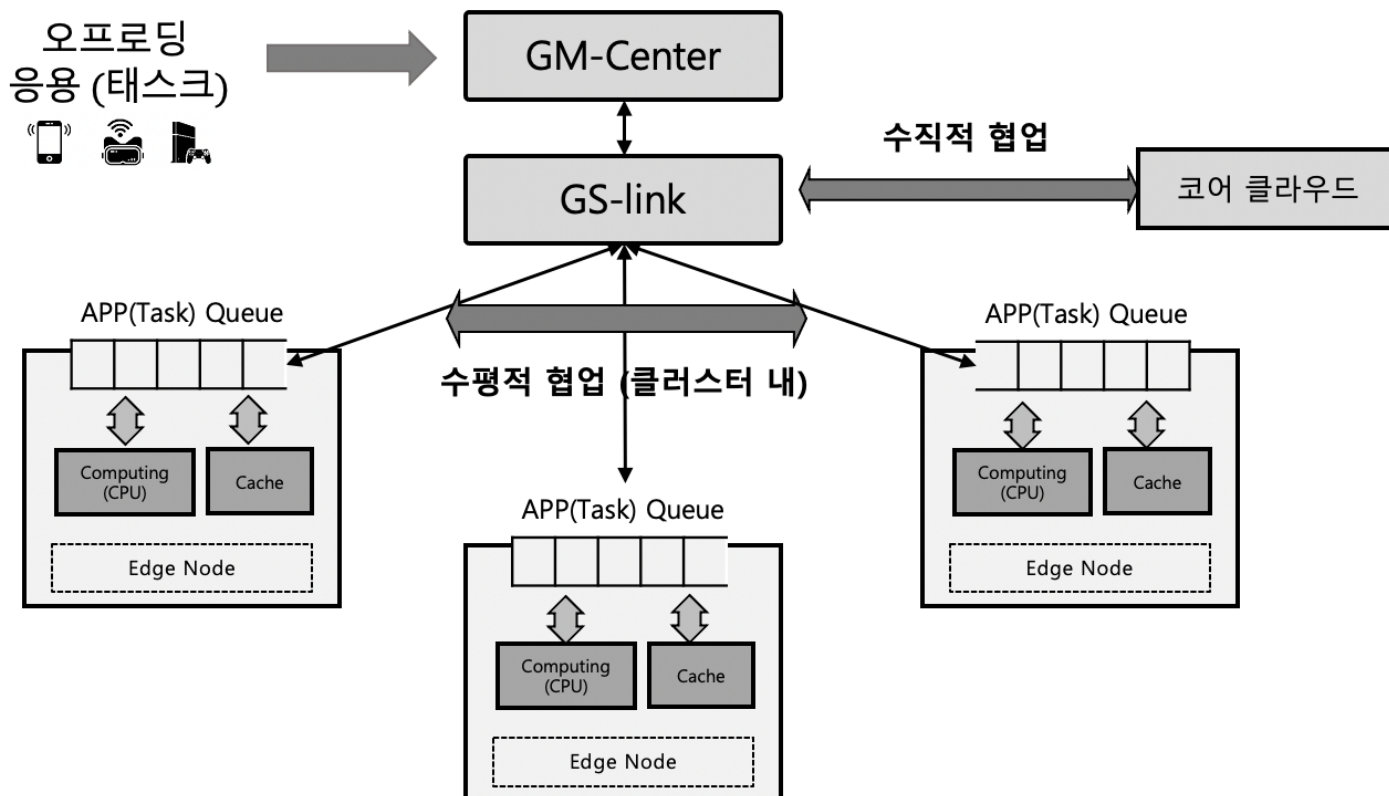
- 기본 협업 오프로딩 (Collaborative Offloading Service(COS)-1)
 - Type 1 응용을 위한 협업 오프로딩 서비스
 - 일반적인 오프로딩 서비스
 - 정책 기반 응용을 위한 단일 엣지 할당 방법 적용
- 독립 태스크 용 협업 오프로딩 (COS-2)
 - Type 2 응용을 위한 협업 오프로딩 서비스
 - 정책 기반 태스크 별 단일 엣지 할당 방법 적용
- 종속 태스크 용 협업 오프로딩 (COS-3)
 - Type 3 응용을 위한 협업 오프로딩 서비스
 - 정책 기반 멀티 태스크(응용)을 고려한 그룹 엣지 할당 방법 적용
- 모든 협업 오프로딩 서비스는 수평적 협업 서비스 제공

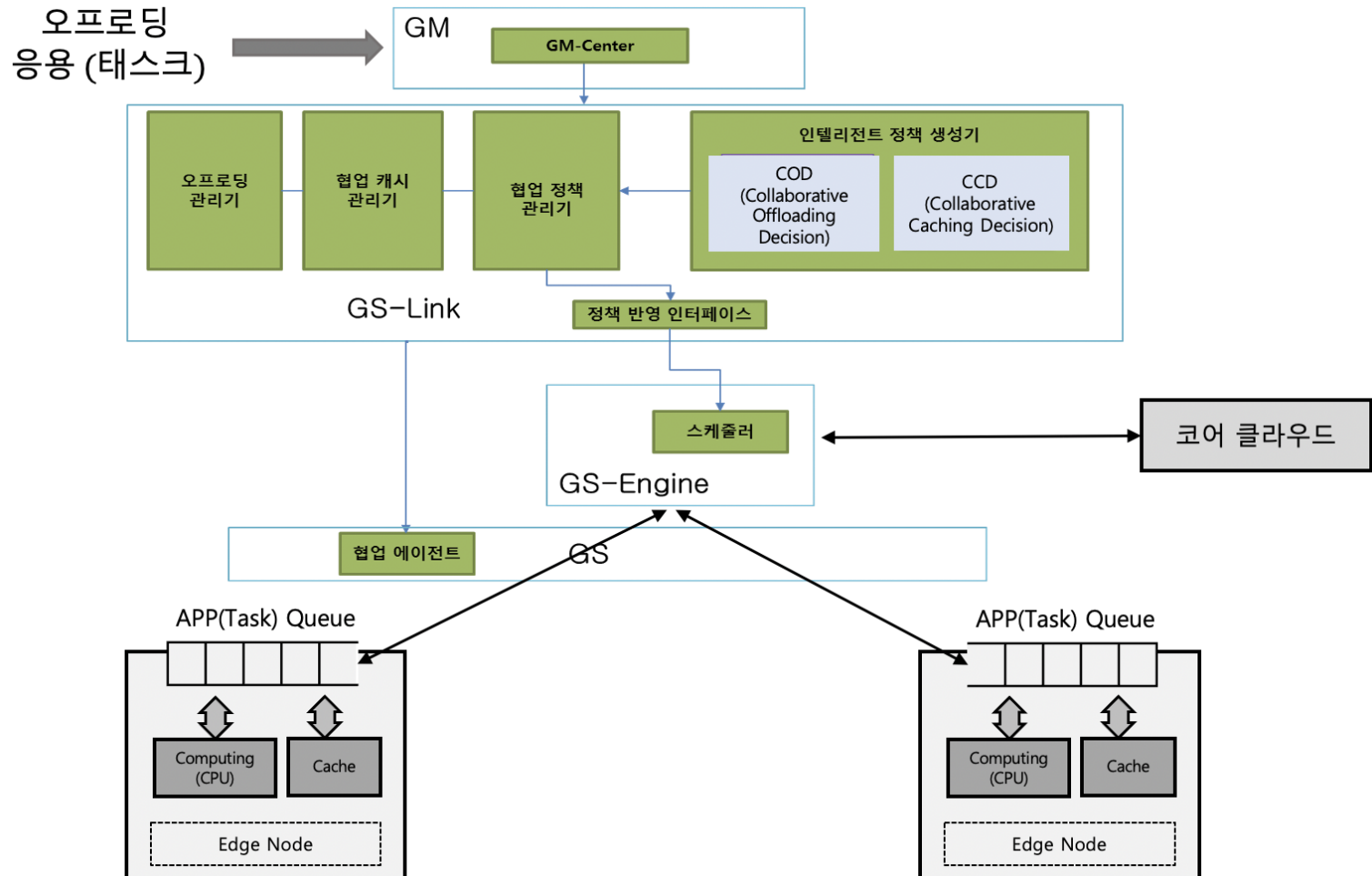
2 협업 응용(태스크) 오프로딩 모델

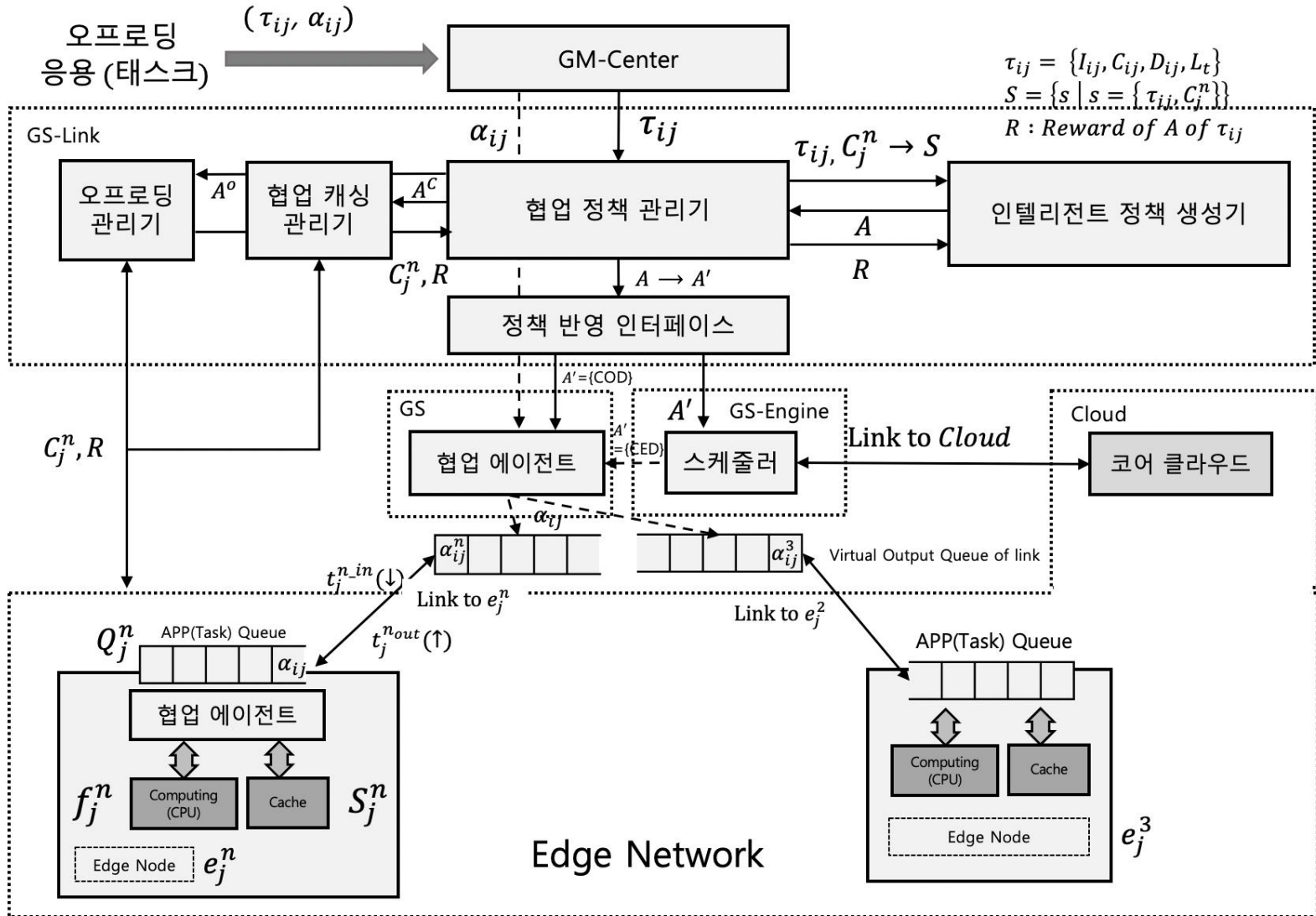
- 협업 오프로딩 모델 (추상적 모델)



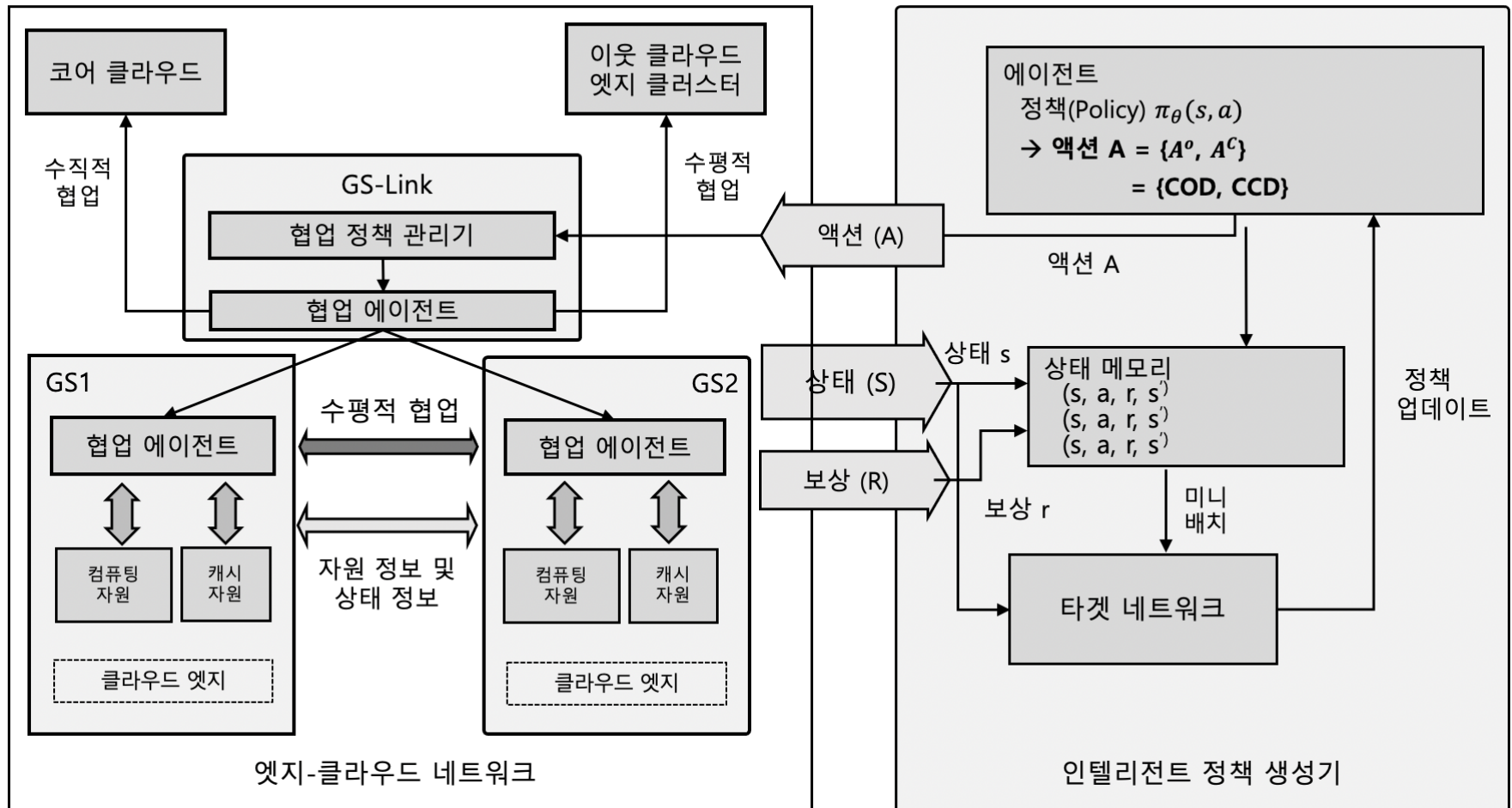
- 싱글클러스터 기반 서비스 협업 프레임워크



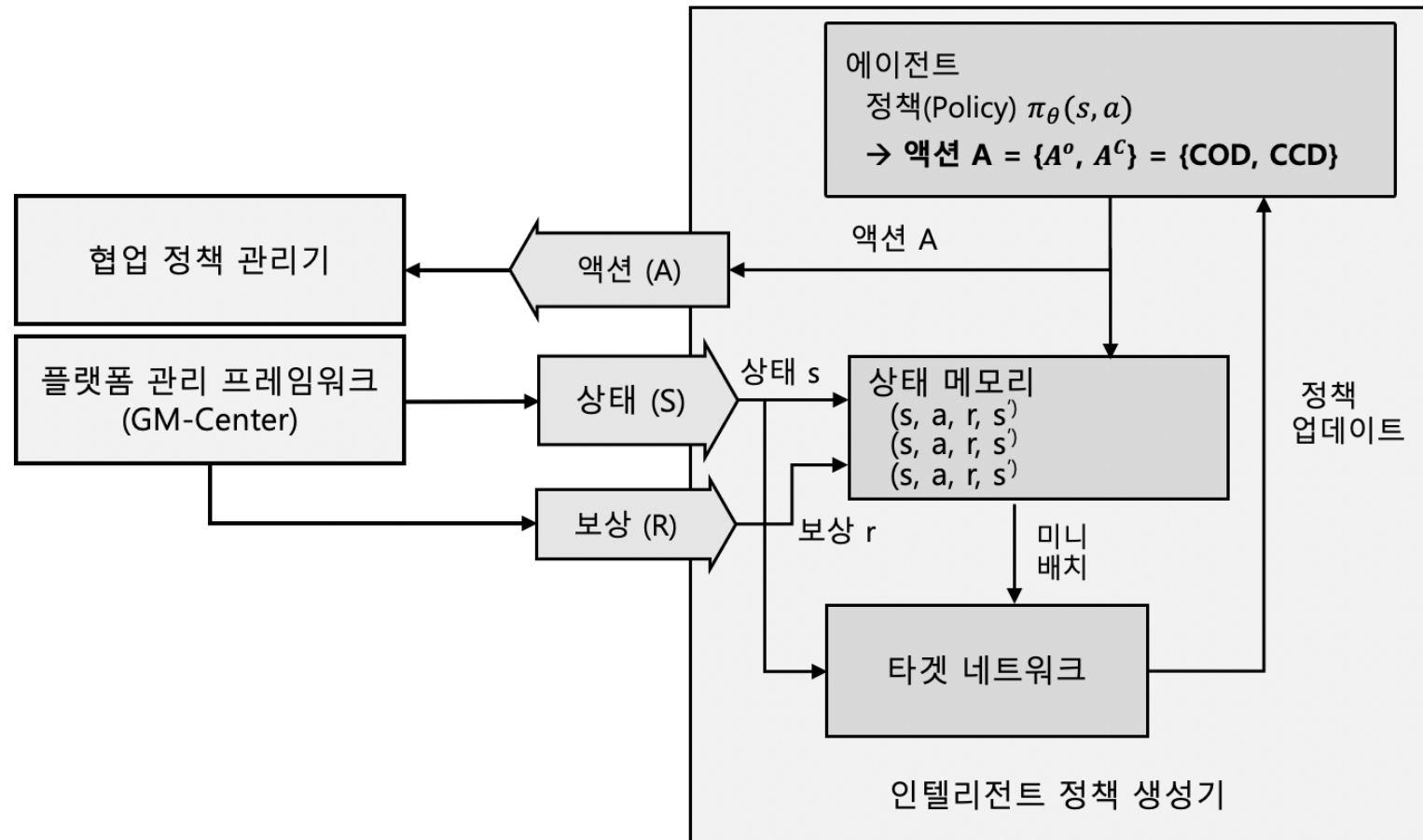




COD (Collaborative Offloading Decision)
CCD (Collaborative Caching Decision)



COD (Collaborative Offloading Decision)
CCD (Collaborative Caching Decision)

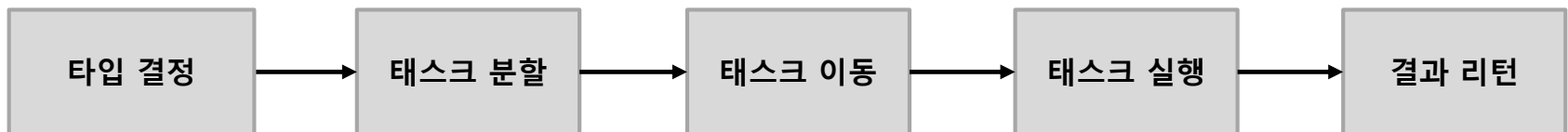


- 협업 오프로딩 모델 (Collaborative Offloading Model)
 - 전체 이관(Full Offloading) 또는 부분 이관(Partial Offloading)
- 협업 오프로딩 과정
 - 오프로딩 서비스 타입 결정 (전체 이관, 부분 이관 결정)

- (전체 이관 경우, COS-1) → 오프로딩 태스크를 위한 최상의 엣지 결정
- (부분 이관 경우, COS-2) → 태스크 분할 → 태스크 별 엣지 결정
- (부분 이관 경우, COS-3) → 태스크 분할 → 태스크((마이크로)서비스)유무 & 성능 기반 엣지 결정 & 태스크 체인 설정

- → 태스크 이동 → 태스크 실행 → 태스크 결과 리턴

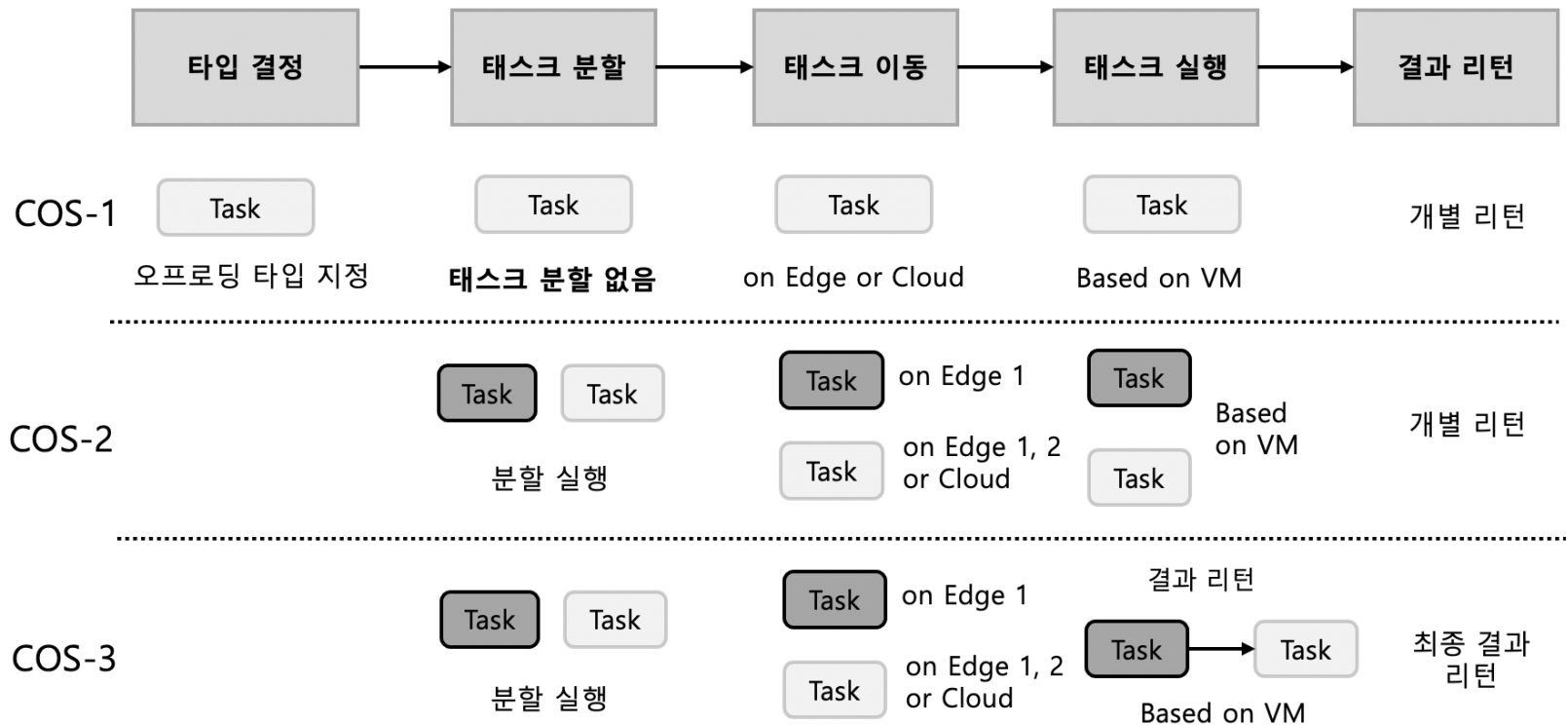
- 협업 오프로딩 과정



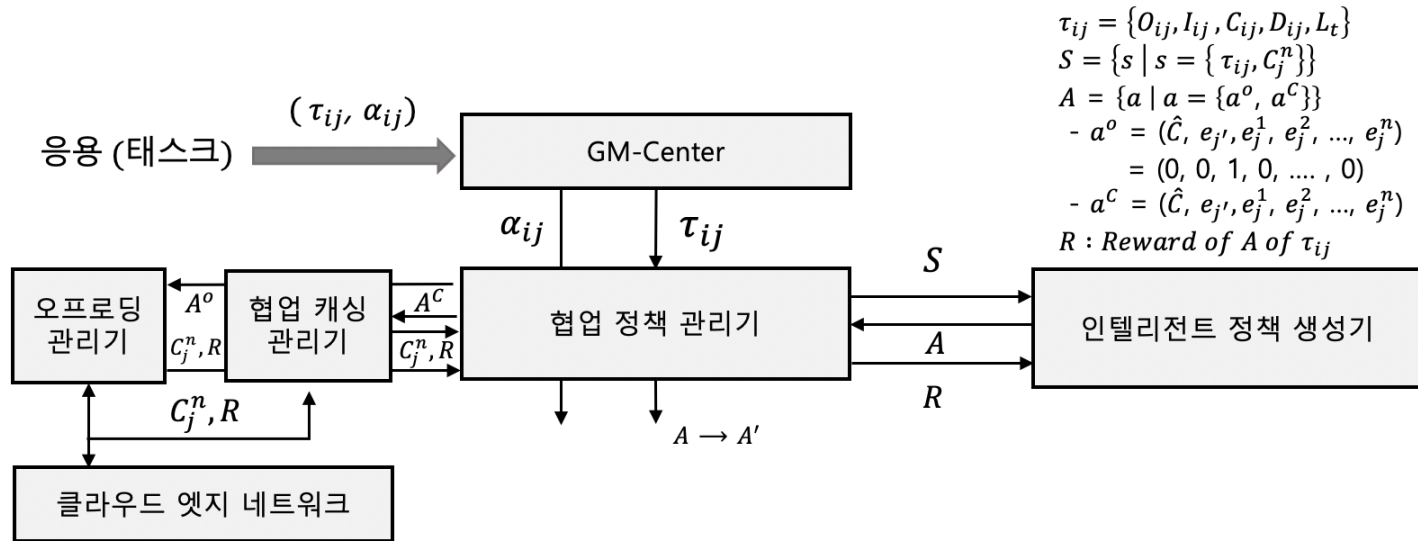
- 멀티클러스터 기반 수평적 협업 오프로딩 시나리오



협업 오프로딩 과정



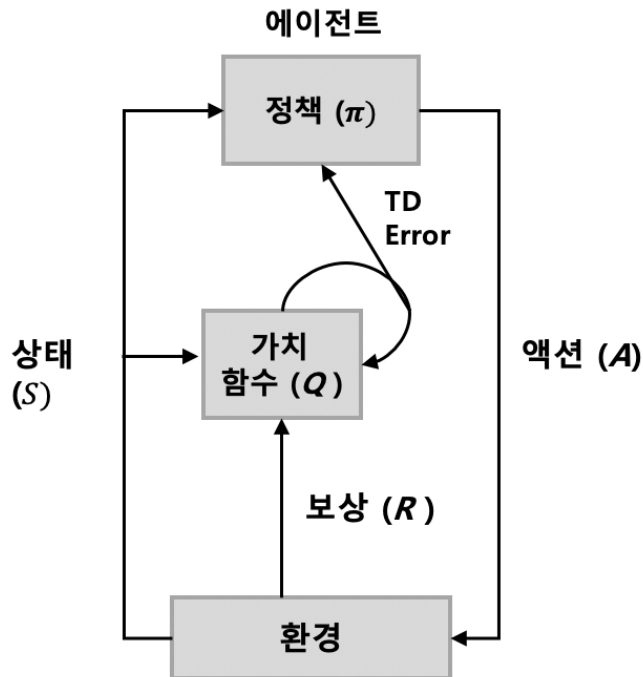
● 정책 생성 모델



● 협업 서비스 정책

- 협업 오프로딩 정책, 분산-협업 캐싱 정책, 고속 서비스 이동 정책
- $\pi_n^O, n = \{1, 2, \dots, n\}, \pi_n^C, n = \{1, 2, \dots, n\}, \pi_n^M, n = \{1, 2, \dots, n\}$

강화학습 정책 모델



$$S = \{s \mid s = \{\tau_{ij}, C_j^n\}\}$$

$$A = \{a \mid a = \{a^o, a^c\}\}$$

$$\begin{aligned} - a^o &= (\hat{C}, e_{j'}, e_j^1, e_j^2, \dots, e_j^n) \\ &= (1, 0, 1, 0, \dots, 0): \text{COS-2} \end{aligned}$$

$$- a^c = (\hat{C}, e_{j'}, e_j^1, e_j^2, \dots, e_j^n)$$

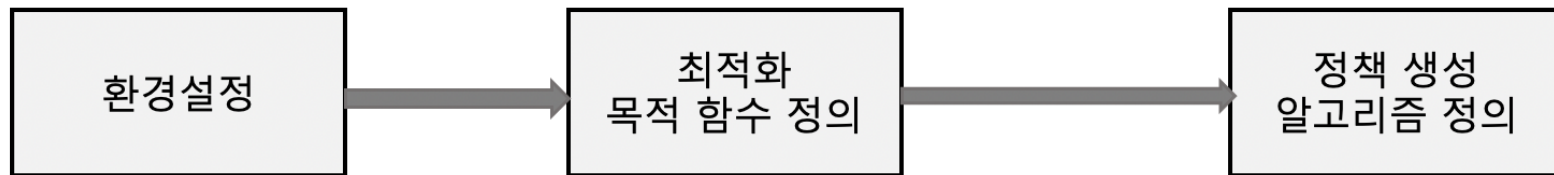
$$Q(s_t, a_t) = R(s_t, a_t) + \gamma Q(s_{t+1}, a_{t+1})$$

$$R_t = \sum_{k=0}^T \alpha^k r_{t+k+1} \text{ (future reward at time t)}$$

* $\alpha \in (0, 1]$ is a discount factor

r_t is reward of action a_t in s_t

● 액션 생성



- $MIN (t_i: TST(Task Service Latency))$
- * $TST = Task\ return\ time\ to\ node$
- $Task\ request\ time\ to\ edge$

- $S = \{s \mid s = \{\tau_{ij}, C_j^n\}\}$
- $A = \{a \mid a = \{a^o, a^c\}\}$
- $Q(s_\tau, a_\tau)$
- $R = \sum_{k=0}^T \alpha^k r_{\tau+k+1}$
 $(a_t \rightarrow r_t)$

- 최적화 목적 함수

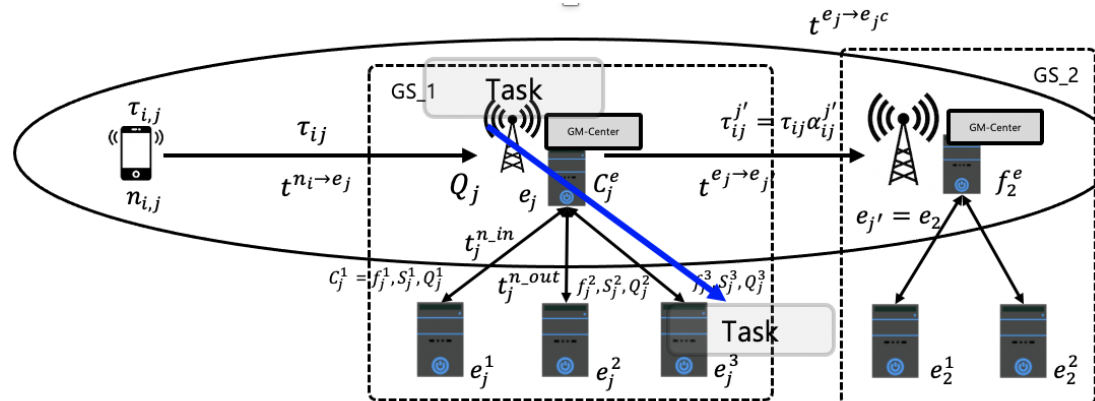
- 모든 태스크의 평균 예상 서비스 지연

$$\text{Min } (T_{all} = \frac{\sum_{i=1}^N t_i}{N}) \quad (N: \text{the number of all tasks})$$

- 태스크(응용) 서비스 타임

$$\text{- COS-1: } t_{ij} = t^{n_i \rightarrow e_j} + Q_j + t_j^{n_in} + Q_j^n + \frac{C_{ij}}{f_j^n} + t_j^{n_out} + t^{e_j \rightarrow n_i}$$

$$* t^{n_i \rightarrow e_j} = \frac{D_{ij}}{R^{n_i \rightarrow e_j}}, \quad t_j^{n_in} = \frac{D_{ij}}{R_j^{n_in}}, \quad t_j^{n_out} = \frac{D_{ij}^r}{R_j^{n_in}}, \quad t^{e_j \rightarrow n_i} = \frac{D_{ij}^r}{R^{e_j \rightarrow n_i}}$$

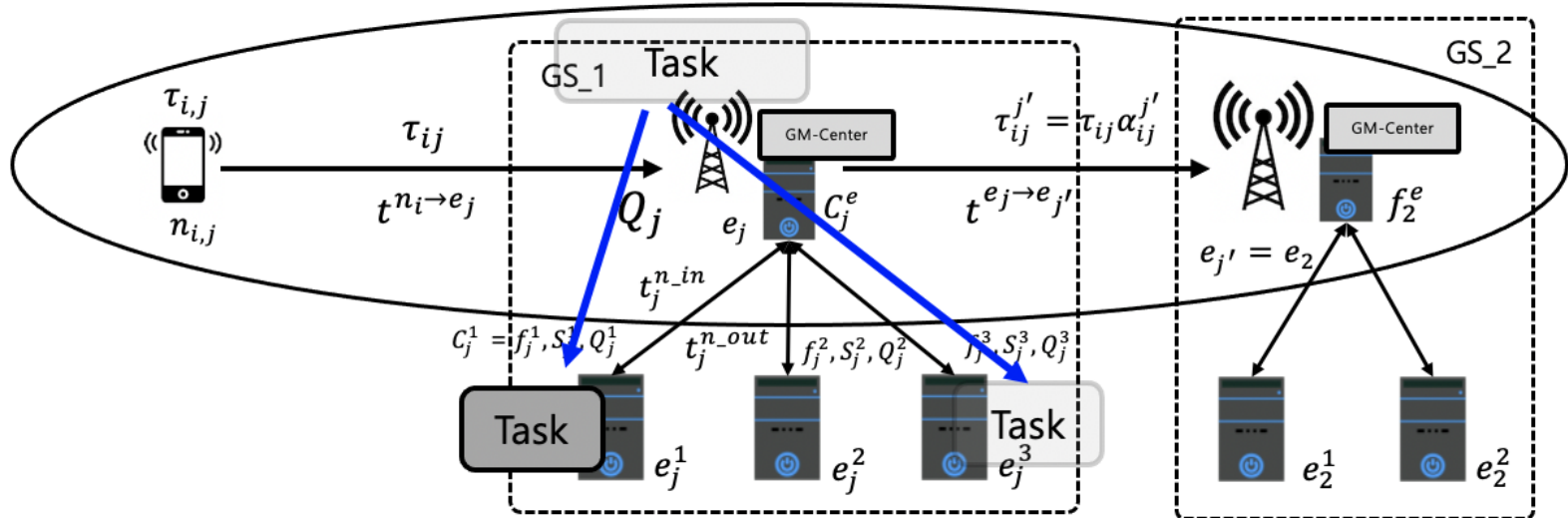


● 서비스 타입 2의 태스크(응용) 서비스 타임 :

– the offloading ratio of the task data

$$- \text{COS-2: } t_{ij} = t^{n_i \rightarrow e_j} + Q_j + \max(\alpha_{ij}^1 (t_j^{1-in} + D_{ij} \times \frac{C_{ij}}{f_j^1} + t_j^{1-n-out}), \alpha_{ij}^2 (t_j^{3-in} + D_{ij} \times \frac{C_{ij}}{f_j^1} + t_j^{3-n-out})) + t^{e_j \rightarrow n_i}$$

$$* \alpha_{ij}^1 + \alpha_{ij}^2 = 1$$



● 보상

- 최적화 목적 함수와 관련 있도록 디자인해야함
- 모든 태스크의 평균 예상 서비스 지연을 최소화 할 수 있도록 보상 지불 방법을 정의하는 게 중요함

$$- Min (T_{all} = \frac{\sum_{j=1}^N \sum_{i=1}^N t_{ij}}{N}) \text{ (N: the number of all tasks)}$$

$$- R_{\tau} = \sum_{k=0}^T \alpha^k r_{\tau+k+1} \text{ (future reward at time } \tau \text{)}$$

- 태스크 수행이 일정하다면

$$- r_{\tau} = \begin{cases} \pi t = avg(t_{ij}) - t_{ij}, & t_i < avg(t_{ij}) \\ -\gamma \pi t = t_{ij} - avg(t_{ij}), & t_i > avg(t_{ij}) \\ \eta, & \text{if task is failed} \end{cases}$$

- * π is the price of edge,
- η is negative(penalty) value,
- γ is penalty value ($0 < \gamma < 1$)

- Double DQN

- Long-term cumulative discount reward value $Q(s_t, a_t)$

- $Q(s_t, a_t) = E_{\pi}[\sum_{\tau=0}^{\infty} \gamma^{\tau} r_{\tau}(S_{\tau}, A_{\tau})]$

- Optimization objective function

- $Q_{\pi}(s, a) = \max_{A_{\tau}} E_{\pi}[\sum_{\tau=0}^{\infty} \gamma^{\tau} r_{\tau}(S_{\tau}, A_{\tau})]$

- **인텔리전트 정책 생성기 개발**
 - 강화 학습 기반 정책 생성기 개발 완료
 - 최적 알고리즘 선정
- **분산캐쉬 사용을 위한 정책 생성 알고리즘 개발**
- **고속 응용 이동 정책 생성 알고리즘 개발 및 정책 생성 학습 모델 개발**
 - 엣지 공유캐시를 통한 고속 응용 이동 알고리즘 개발
 - 공유캐시 내 고속 응용 이동을 위한 딥러닝(어텐션+LSTM 알고리즘) 기반 응용 배치 선정 및 정책 알고리즘 개발
 - 협업 엣지 네트워크 모델 기반 고속 응용 이동 정책 생성 알고리즘 및 정책 생성을 위한 학습 모델 개발

감사합니다.

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



GEdge Platform 코어 개발자

윤주상(jsyoun@deu.ac.kr)

Welcome to GEdge Platform

An Open Cloud Edge SW Platform to enable Intelligent Edge Service

GEdge Platform will lead Cloud-Edge Collaboration