

멀티 클라우드 인프라 기반 연합학습 환경 구축 플랫폼 개발



정보컴퓨터공학부 20학번 전진혁

향후 연구 방향



향후 연구 방향

1. 에너지 효율, 네트워크 사용량 등 **복합적인 지표**를 반영하는 **최적화 전략**을 연구
 - **강화학습** 기반의 지능형 오케스트레이션을 적용한 지능형 시스템 구축
2. 연합학습 참여자의 **리소스 효율성** 및 **공정성(Fairness)** 보장을 위한 방법 연구
 - 각 연합학습 참여자의 데이터 품질, 참여 빈도, 네트워크 상태 등을 종합적으로 고려한 다차원적 연합학습 스케줄링 알고리즘 설계
 - 특정 기관에 **과부하**되는 현상을 **완화**
 - 연합학습 **성능**과 **안정성** 동시에 강화하는 방안 마련



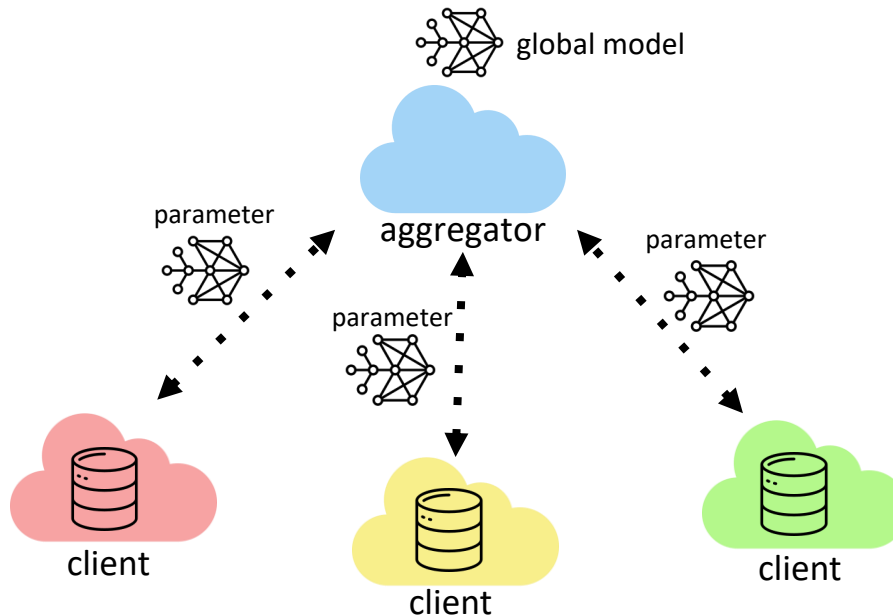
01

배경



연합학습(Federated Learning)이란?

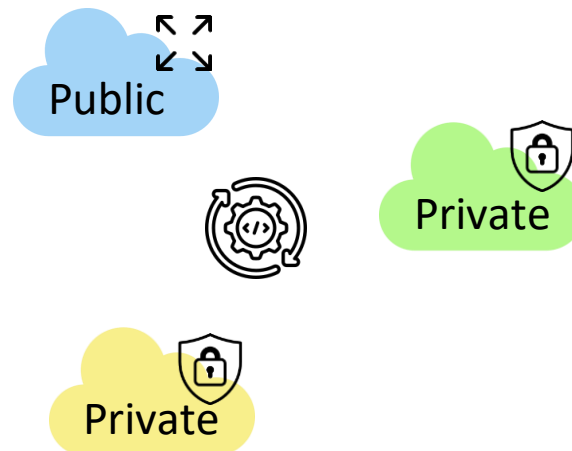
- 참여자(Client)와 집계자(Aggregator)로 이루어진 분산학습 기술
 - 집계자: 참여자의 모델 파라미터 수집 및 통계 작업 수행 후 글로벌 모델 갱신
 - 참여자: 학습 데이터를 사용하여 모델 학습 수행 후 모델 파라미터 전송
- ➡ 중앙으로 학습 데이터의 전송 없이 파라미터만 전송 ➡ 데이터 유출 위험 감소
- ➡ 의료, 금융 및 스마트시티 분야에서 활용





멀티 클라우드(Multi-Cloud)란?

- 이형의 클라우드 플랫폼(퍼블릭 클라우드, 프라이빗 클라우드)을 통합하여 단일 클라우드 플랫폼처럼 활용할 수 있는 기술
 - ➡ 벤더 종속성을 탈피하고, 각 클라우드 플랫폼의 장점 확보 가능
 - ➡ 퍼블릭 클라우드를 통한 확장성 및 접근성 / 프라이빗 클라우드를 통한 보안성



02

문제점

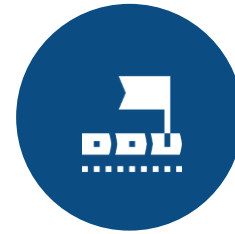
기존 시스템의 문제점



1. 지연 시간 및
비용 최적화 부재



2. 보안 취약성 및
데이터
프라이버시 문제

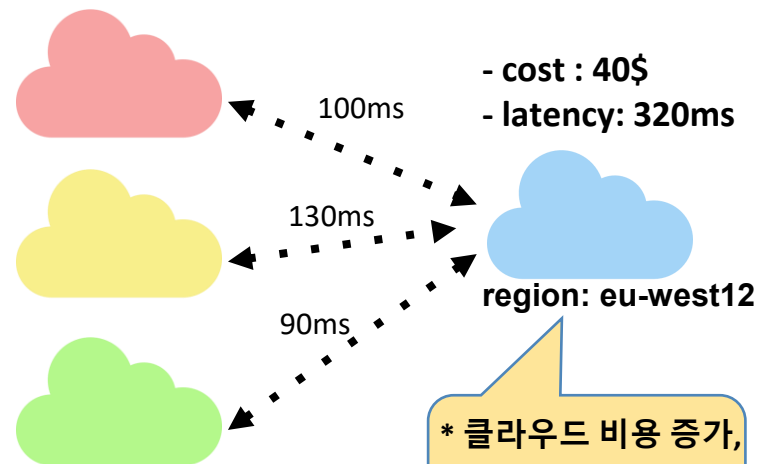
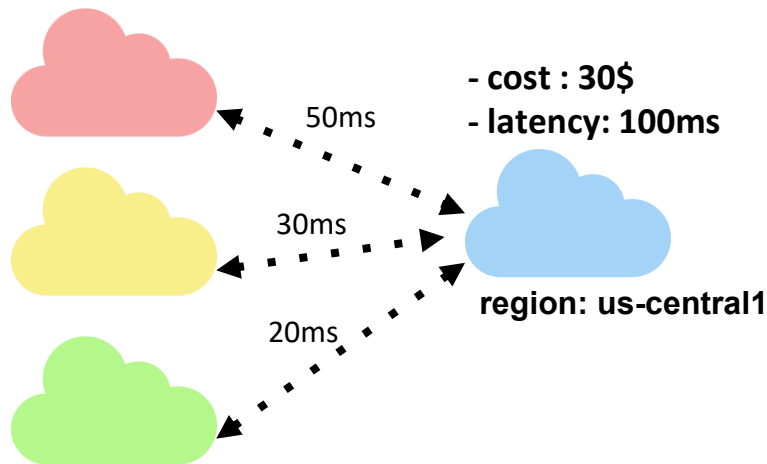


3. 동적 태스크
오케스트레이션 부재

기존 시스템 문제점

1. 운영 비용 및 지연시간 최적화 부재

- 기존 클라우드 기반 연합학습(e.g. AWS SageMaker, GCP VertexAI)의 경우 단일 클라우드 플랫폼에 한정되어 연합학습을 수행
 - 다양한 클라우드 플랫폼의 다양한 리전을 사용하지 못함.
 - 가상머신 운영 비용 증가 및 집계자 위치에 따른 학습시간 증가 가능

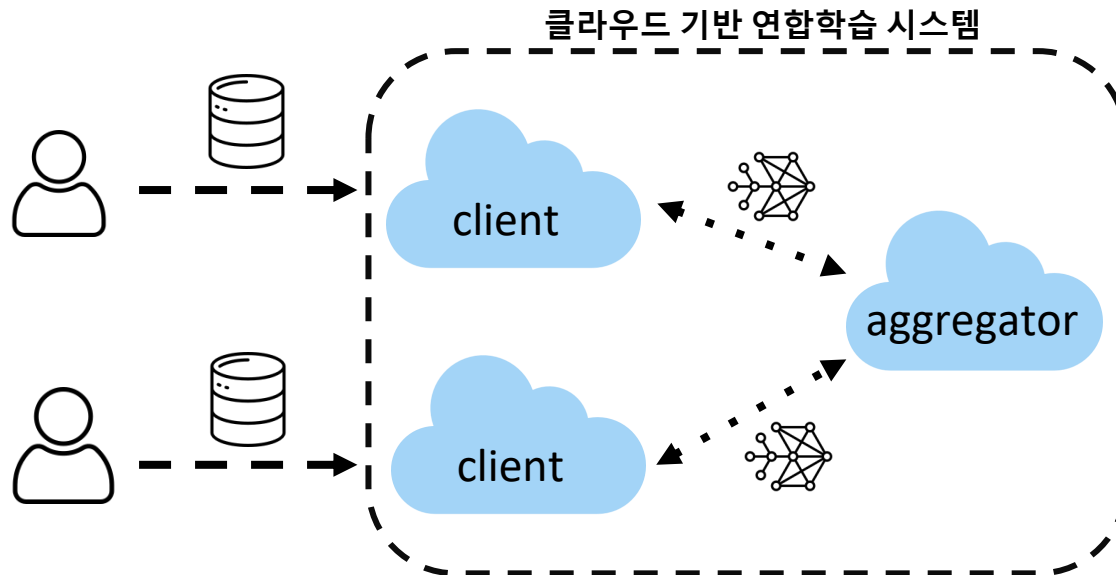


* 클라우드 비용 증가,
* 학습 시간 증가

기존 시스템 문제점

2. 보안 취약성 및 데이터 프라이버시 문제

- 기존 클라우드 기반 연합학습(e.g. AWS SageMaker, GCP VertexAI)의 경우 단일 클라우드만 사용하여 클라우드 플랫폼에 학습 데이터 업로드 요구됨
 - ➡ 학습 데이터 저장위치: 퍼블릭 클라우드
 - ➡ 학습 데이터가 네트워크 환경 및 제3의 인프라에 노출될 수 있는 보안적 취약점 존재



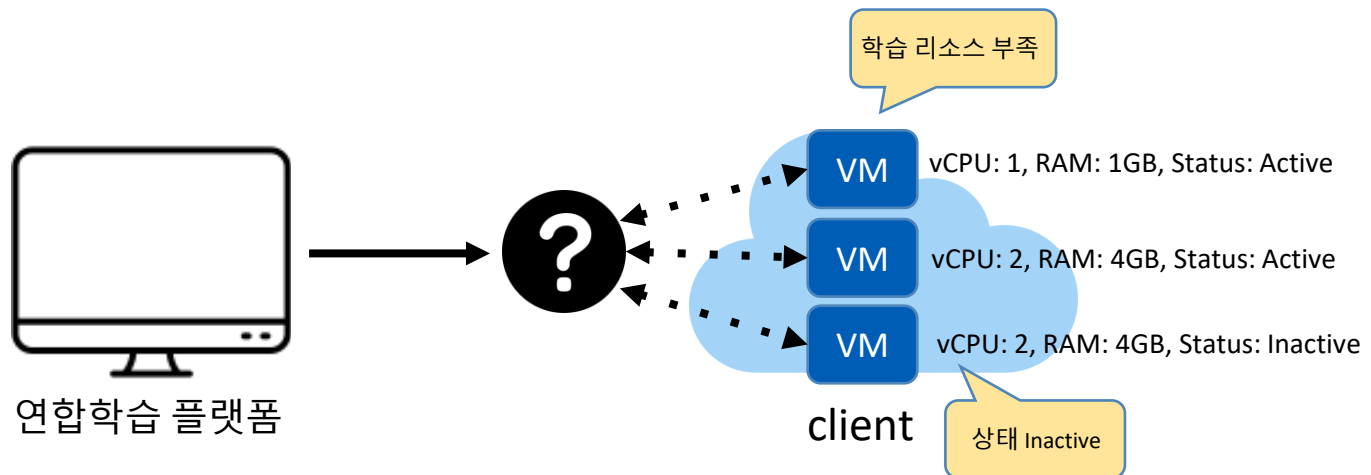
기존 시스템 문제점

3. 동적 태스크 오케스트레이션 부재

- 기존 연합학습 플랫폼(FedML, IBM Federated Learning 등)의 경우 연합학습 태스크를 연합학습 참여자의 적절한 가상머신(VM)에 할당하는 **동적 태스크 오케스트레이션 부재**

- 연합학습 참여자의 여러 가상머신 중 **학습을 수행할 가상머신 결정 필요**
- 컴퓨팅 리소스가 부족하거나, 상태가 Inactive한 가상머신 선택할 경우 **참여자 이탈 및 연합학습 실패**

* 참여자 이탈 : 연합학습 참여자가 학습에 참여하지 못하는 것



03

연구 목표

- ① 지연 시간 및 비용 최적화 부재,
- ② 보안 취약성 및 데이터 프라이버시 문제,
- ③ 동적 태스크 오케스트레이션 부재를 해결을 위한

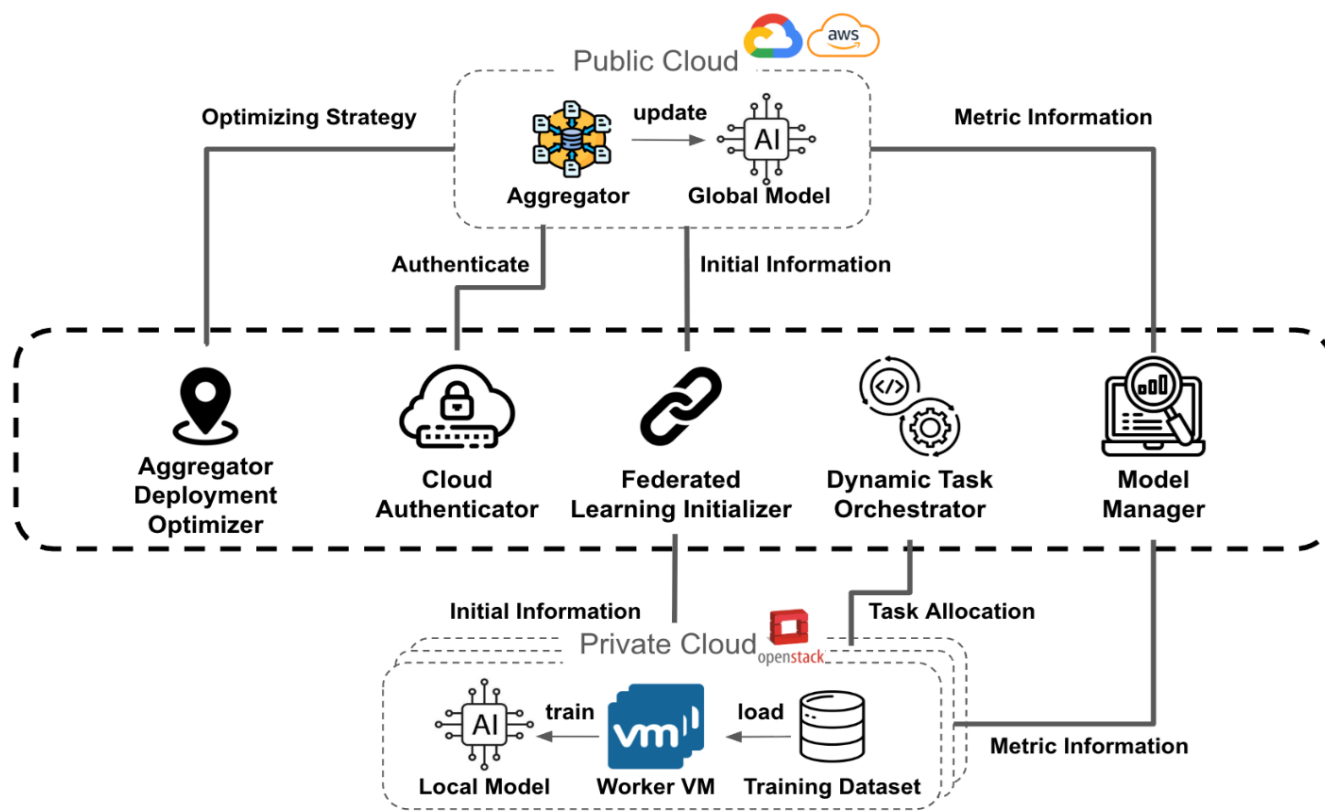
**“멀티 클라우드 인프라 기반
연합학습 환경 구축 플랫폼 개발”**

04

멀티 클라우드 인프라 기반 연합학습 환경 구축 플랫폼

제안 시스템

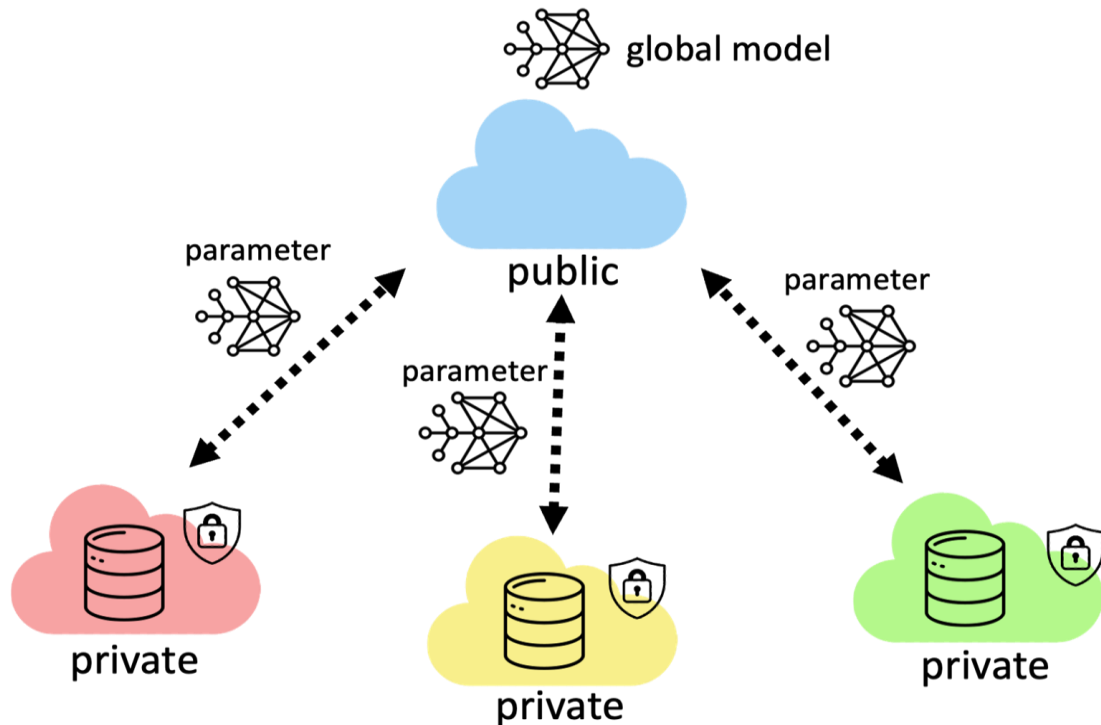
구조도



제안 시스템

계층 구조 기반 연합학습 수행

- 퍼블릭 클라우드 - 프라이빗 클라우드 역할 분리를 기반으로 한 연합학습 수행
 - 프라이빗 클라우드: 자체 데이터를 사용해 모델 학습 수행
 - 퍼블릭 클라우드: 모델 파라미터를 수신 받아 통계 작업 수행



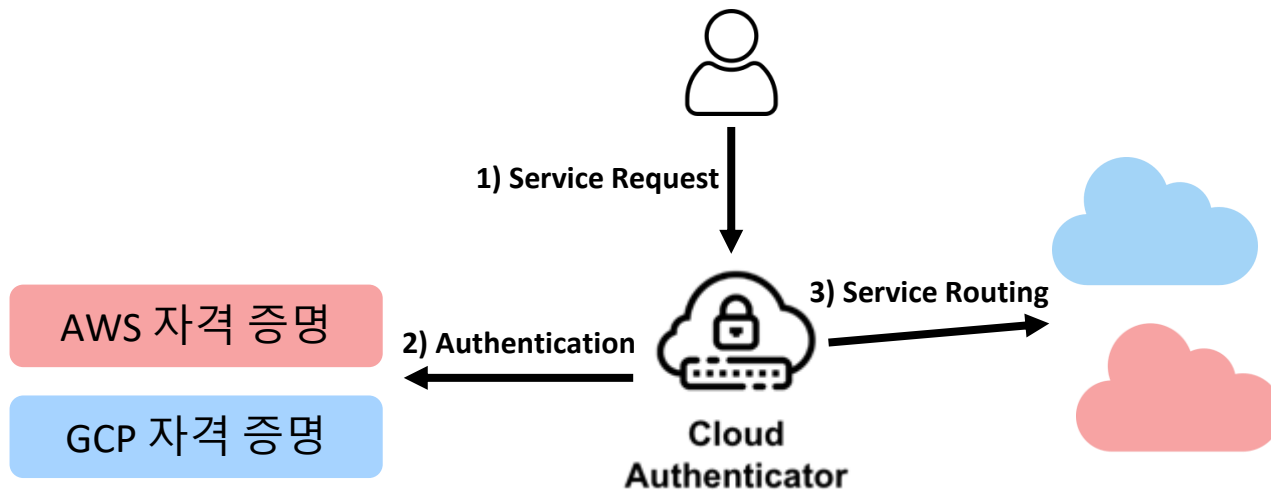
제안 시스템



Cloud Authenticator

- 멀티 클라우드 서비스를 제공하기 위해 다양한 클라우드 플랫폼과의 연동을 수행

클라우드 자격 증명 파일 업로드 ➔ 클라우드 API 호출 기반 검증 ➔ 검증성공 시 연동 완료
➔ 멀티 클라우드 서비스 활용 가능

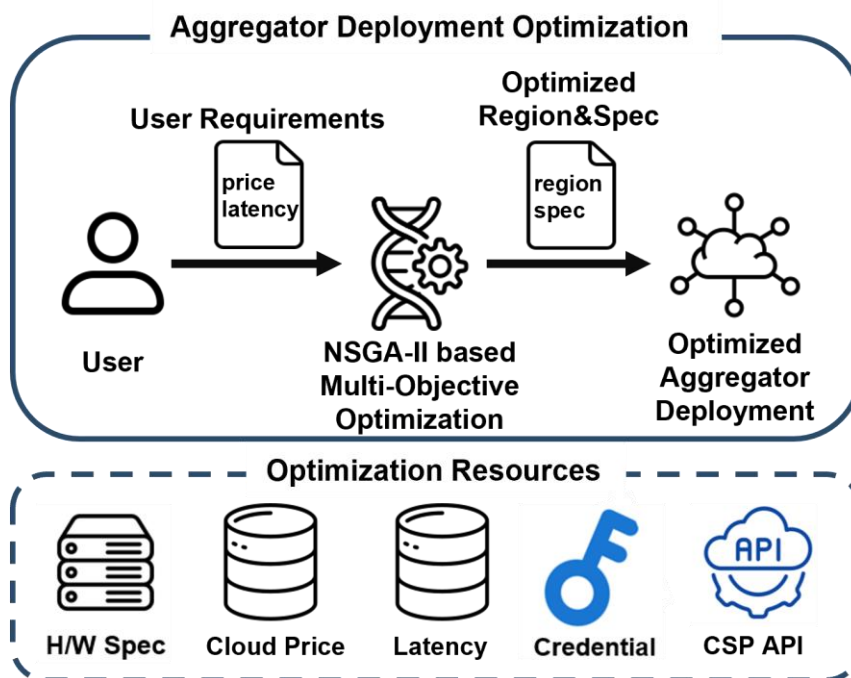


제안 시스템



Aggregator Deployment Optimizer

- 사용자 요구사항 기반 집계자 배포 위치 및 스펙(CPU, RAM) 최적화
- ➔ 유전 알고리즘 중 하나인 NSGA-II를 이용한 다목적 최적화



다목적 최적화:
➔ 비용-성능같이 상충될 수 있는 두 개 이상의 목표의 균형점을 찾는 방법

사용자 요구사항

- 최대 허용 비용
- 최대 허용 지연시간
- 비용-지연시간 가중치

집계자 선택

최적화 요약

참여자 수: 3명

후보 옵션: 7283개

참여자 지역: us-east4, europe-west4, ap-northeast-1

조건 만족 옵션: 6개

#1

aws us-east-1

추천도: 14.4%

₩31,450

월 예상 비용

인스턴스

t4g.medium

2vCPU, 4096GB

평균 지연시간

231.76ms

최대 지연시간

233.49ms

시간당 비용

\$0.0336

#2

aws us-east-1

추천도: 9.6%

₩35,194

월 예상 비용

인스턴스

t3a.medium

2vCPU, 4096GB

평균 지연시간

231.76ms

최대 지연시간

233.49ms

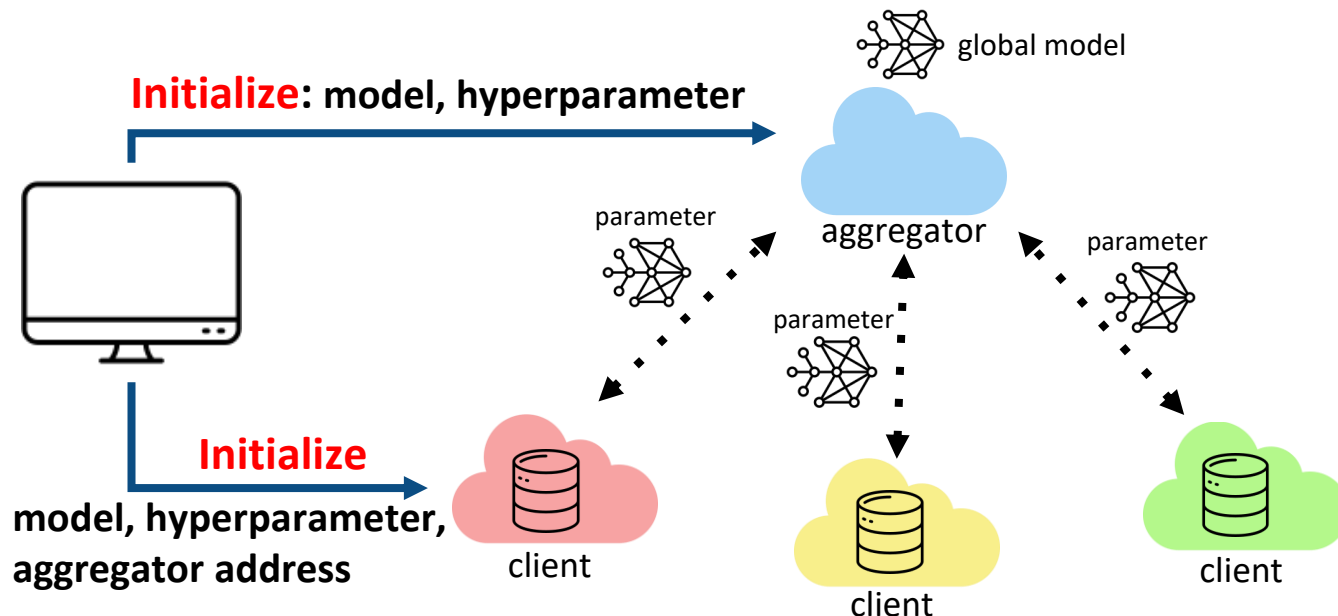
시간당 비용

\$0.0376



Federated Learning Initializer

- 연합학습 집계자와 연합학습 참여자를 연계한 연합학습 수행을 위해
연합학습 환경 설정 및 연합학습 수행 명령 전달
 - 학습할 모델 파일
 - 연합학습 집계자의 주소
 - 모델 학습 하이퍼파라미터(epoch, batch size 등)





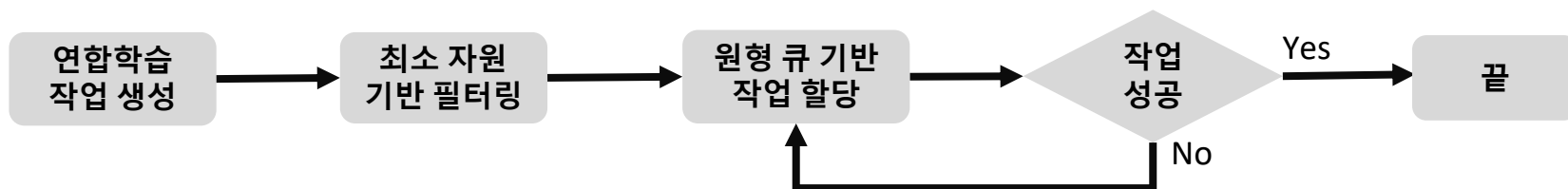
Dynamic Task Orchestrator

* 참여자 이탈 : 연합학습 참여자가 학습에 참여하지 못하는 것

- 여러 대의 가상머신을 운영하는 연합학습 참여자에게 가상머신의 자원 상태를 고려해 적절한 가상머신에 작업을 할당
- 연합학습 참여자의 이탈률을 줄이고 연합학습 작업의 안정성을 높이기 위함

수행 흐름

- 1) 최소 자원 기반 필터링
- 2) 원형 큐 기반 작업 할당
- 3) 할당 된 가상머신에서 작업 실패 시 재시도





Model Manager

- 연합학습 과정에서 생성되는 연합학습 모델(가중치, 메트릭, 메타데이터)를 관리
 - 1) 모델 성능 추적
 - 각 학습 라운드마다 생성되는 글로벌 모델의 성능 (정확도, 손실값, F1-score, Recall, Precision)을 기록
 - 2) 모델 추천 / 선택: 사용자 요구사항(성능, 정확도 등)에 따른 최적 모델 제공
 - 사용자 요구사항- 모델 평가지표(정확도, F1-score 등)에 따라 최적 모델 추천
 - 3) 최적 모델 다운로드
 - 사용자 요구사항에 기반한 최적의 모델의 다운로드 기능 제공

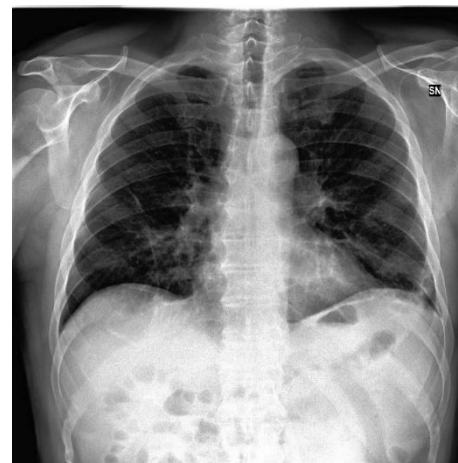
04

사례연구 및 평가



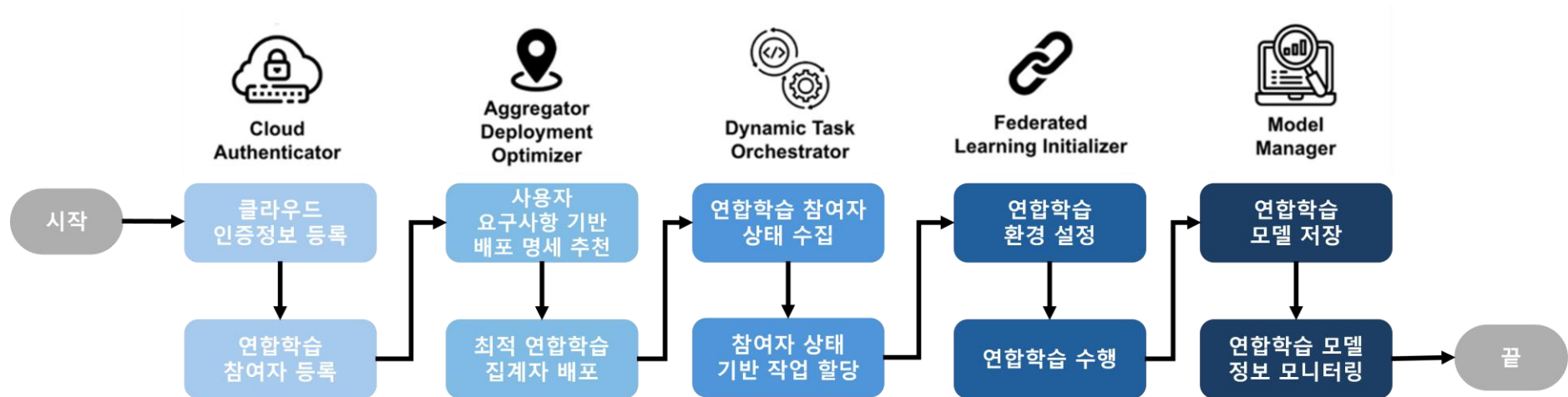
폐 이미지 기반 코로나 19 진단 연합학습 모델 생성

- **의료 정보**
 - 개인정보 보호 필요성
 - 민감한 정보로 인한 데이터 공유 어려움 → 연합학습에 적합한 도메인
- 연합학습 참여자 설정
 - 3명의 참여자
 - **지리적 분산 환경 모사**를 위해 US, Asia, Europe 리전에 각각 구축
- **집계자 최적화 관련 사용자 요구사항**
 - 지연시간 : 비용 가중치 조절(1:9, 5:5, 9:1)
 - 최대 지연시간 : 250ms
 - 최대 비용: 85000KRW
- 데이터 및 모델 설정
 - 데이터: Kaggle Covid 19 Image Dataset
 - 모델: 일반적인 CNN 모델(26MB)



폐 이미지 예시

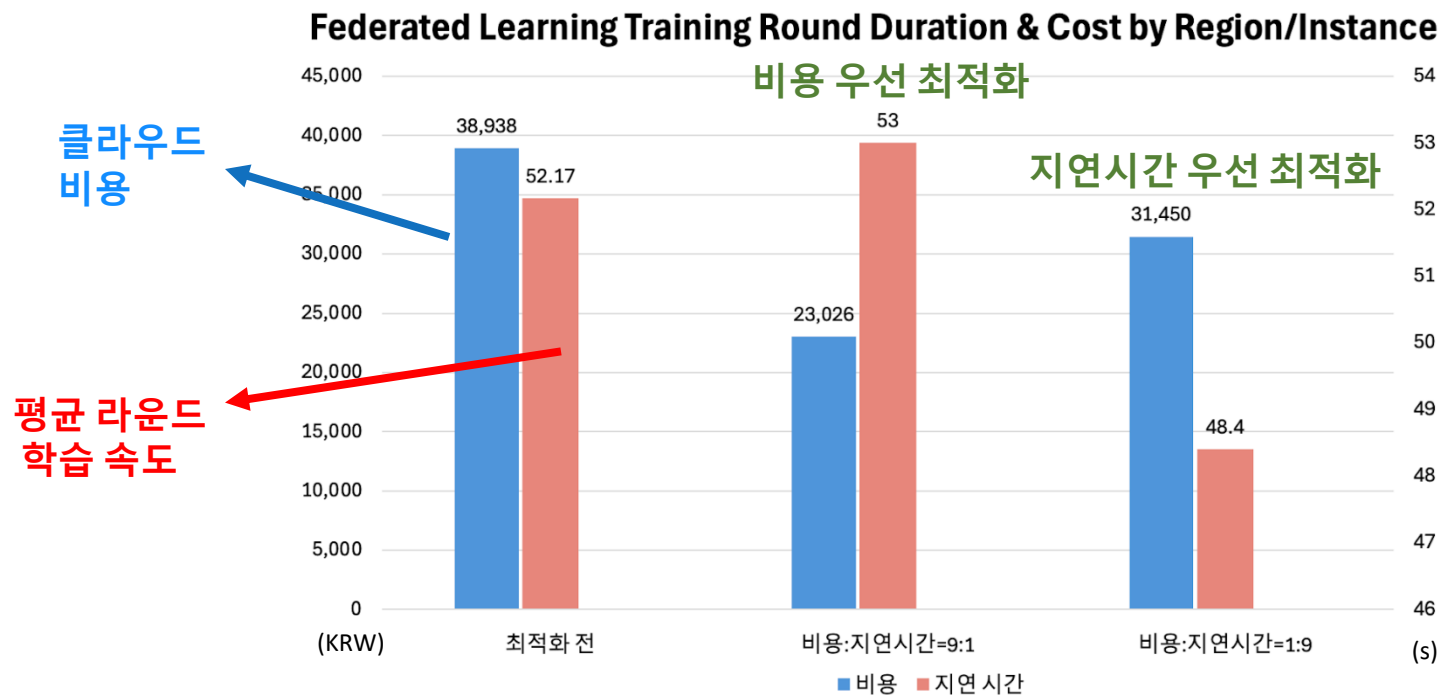
연합학습 모델 생성 전주기 과정





연합학습 집계자 최적화 평가 (1)

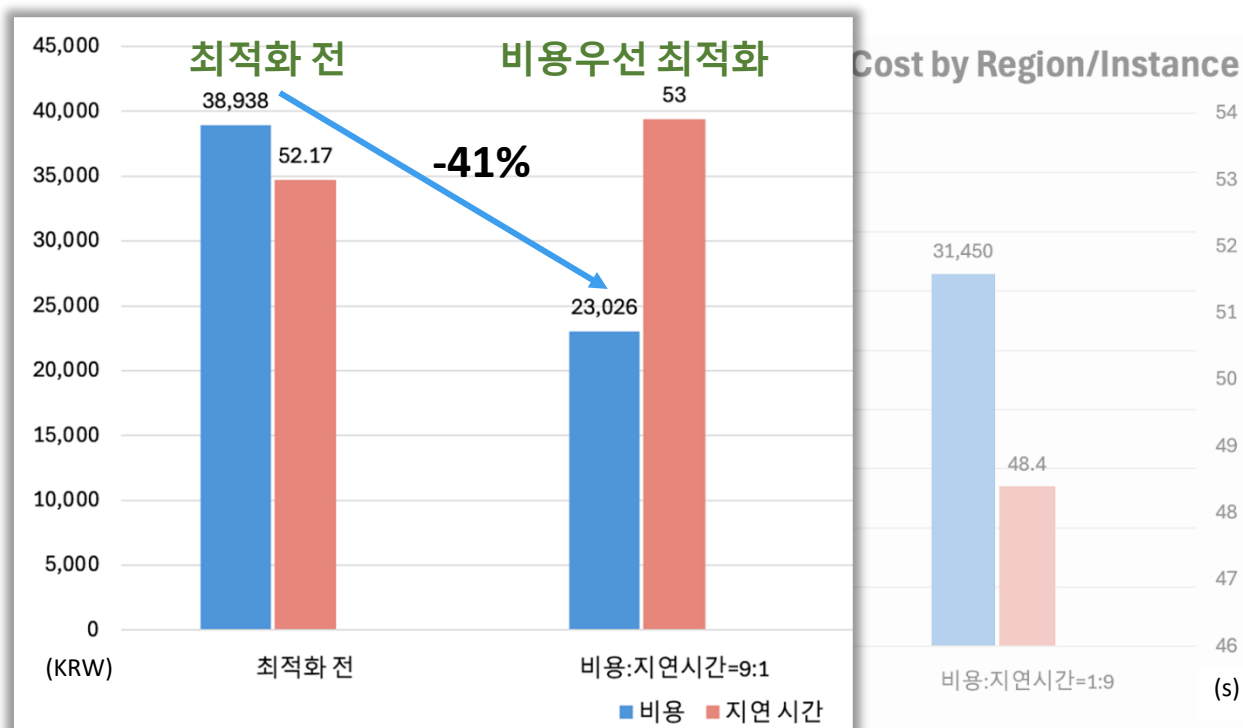
- 세 가지 시나리오를 이용하여 평가
 - 최적화 전
 - 비용 우선 최적화(비용:지연시간 = 9:1)
 - 지연시간 우선 최적화(비용:지연시간 = 1:9)





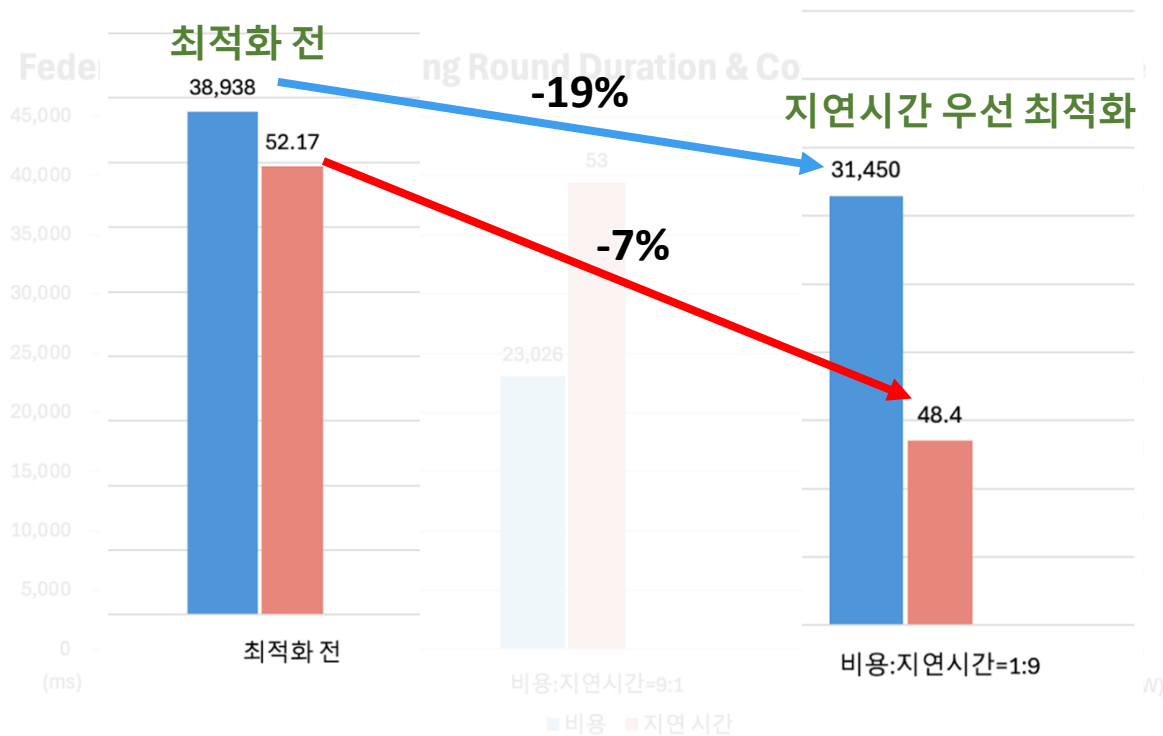
연합학습 집계자 최적화 평가 (2)

- 최적화 수행 전 vs 비용우선 최적화 : 클라우드 비용 41% 감소



연합학습 집계자 최적화 평가(3)

- 최적화 수행 전 vs 지연시간 우선 최적화 : 클라우드 비용 19% 감소, 학습 속도 7% 단축





동적 태스크 오케스트레이션 평가(1)

- 각 참여자가 5대의 가상머신을 운영하고 있으며, 20번의 연합학습을 수행하는 시나리오
- 제안하는 동적 태스크 오케스트레이션 vs 랜덤 가상머신 선택(재시도 x)

학습 가능
가상머신

- 1) 최소 자원 기반 필터링
- 2) 원형 큐 기반 작업 할당
- 3) 할당된 가상머신에서 작업 실패 시 재시도

VM ID	상태	vCPU	RAM(GB)	Disk(GB)
vm-001	Active	4	8	40
vm-002	Active	2	4	20
vm-003	Inactive	8	16	80
vm-004	Inactive	2	4	20
vm-005	Active	1	1	10

vm-003, vm-004: **상태 != Active**

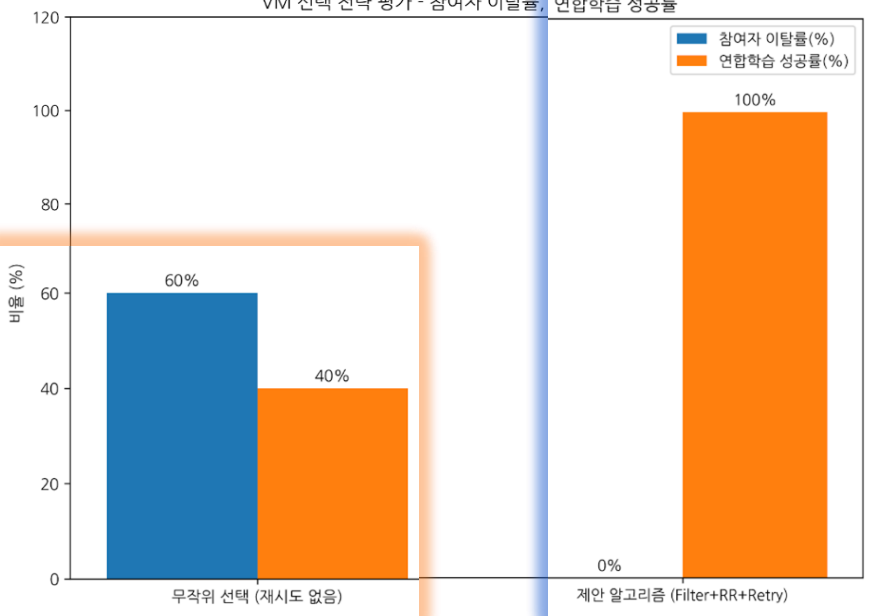
vm-005: **최소 사양 기준 미충족**



동적 태스크 오케스트레이션 평가(2)

학습 가능
가상머신

VM 선택 전략 평가 - 참여자 이탈률, 연합학습 성공률



VM ID	상태	vCPU	RAM(GB)	Disk(GB)
vm-001	Active	4	8	40
vm-002	Active	2	4	20
vm-003	Inactive	8	16	80
vm-004	Inactive	2	4	20
vm-005	Active	1	1	10

동적 태스크 오케스트레이션 적용시
이탈률 **0%**, 성공률 **100%** 실험 결과를 보여줌

무작위 선택 적용시)
이탈률 **60%**, 성공률 **40%** 실험 결과를 보여줌



유사 시스템 비교 평가

- 기존 퍼블릭 클라우드 기반 연합학습 플랫폼에 비해 **학습 데이터 외부 노출 x**
- 제안 시스템만이 연합학습 **집계자 최적화, 동적 테스트 오케스트레이션** 지원

➡ 기존 플랫폼에 비해 데이터 **프라이버시, 비용 및 학습시간 감소, 안정성 강화** 측면에서 우수

구분	제안 시스템	AWS SageMaker 기반 연합학습	GCP Vertex AI 기반 연합학습	FedML
학습 데이터 저장 위치	프라이빗 클라우드 내부	AWS 내부 스토리지 (S3 등)	GCP 내부 스토리지 (GCS 등)	참여자 로컬 환경
데이터 신뢰 경계	기관 내부	기관 내부 및 AWS 클라우드	기관 내부 및 GCP 클라우드	참여자 로컬 환경
연합학습 집계자 최적화	지원	미지원	미지원	미지원
동적 테스트 오케스트레이션	지원	미지원	미지원	미지원

05

결론 및 향후 연구



결론

- 기존 클라우드 기반 연합학습 플랫폼이 가지는 세 가지 문제점
 - 1) 비용 및 지연시간 최적화 부재
 - 2) 보안 취약성 및 데이터 프라이버시 문제
 - 3) 동적 태스크 오케스트레이션 미지원 문제
- ➔ 이를 해결하기 위해 멀티 클라우드 인프라 기반의 연합학습 환경 구축 플랫폼 제안
 - 집계자: 퍼블릭 클라우드
 - 참여자: 프라이빗 클라우드
- ➔ 프라이빗 클라우드를 이용한 학습 ⇒ 데이터 보안성 유지
- 사례연구: Covid 19 Image Dataset를 활용한 CNN 모델
 - 1) 비용 41% 절감, 학습 속도 7% 향상
 - 2) 참여자 이탈 문제 해소 및 모든 학습 작업에서 100% 성공
- ➔ 비용 및 학습 속도 최적화, 안정적인 학습 수행 달성

THANK YOU

