

# AdaModal-Fed: 연합학습 환경에서 결손 모달리티를 극복하기 위한 클러스터링 및 크로스-어텐션 접근법

## 졸업 프로젝트 2차 보고서

Team 329 (14조)

2276095 김희진

2276225 이슬

2176310 이혜리

# 목차

1. Team Information
2. Project Summary
  - a. 문제점
  - b. 해결책
  - c. 기대효과
3. 과제 설계
  - a. 요구사항 정의
  - b. 전체 시스템 구성
4. 주요 기능 구현
5. 실험
  - a. 자원효율성
  - b. 클러스터링 개수 변화
6. 차별성
  - a. 전체적인 차별성
  - b. FedAvg와의 차별성
7. 논문

## 1. Team Information

- a. 과제명: AdaModal-Fed: Overcoming Missing Modalities with Clustering and Cross-Attention in Federated Learning Approach  
(AdaModal-Fed: 연합학습 환경에서 결손 모달리티를 극복하기 위한 클러스터링 및 크로스-어텐션 접근법)
- b. 팀 정보
  - 팀 번호 : 14
  - 팀 이름 : team329
- c. 팀 구성원
  - 김희진 (2276095)
  - 이슬 (2276225)
  - 이혜리 (2176310)

## 2. Project-Summary (과제 요약)

### 2.1. 문제 정의

#### a. 의료 데이터와 멀티모달성 – missing modality

현대 의료 환경에서는 환자 한 명에 대해 다양한 형태의 데이터가 수집된다. 예를 들어, 흉부 **X-ray** 이미지와 해당 검사에 대한 의사의 진단 보고서(텍스트)가 함께 존재할 수 있다. 그러나 실제 임상 환경에서는 모든 환자에게서 동일한 형태의 데이터를 확보하기 어렵다. 일부 환자의 경우 이미지는 존재하지만 보고서가 누락되어 있거나, 반대로 보고서만 존재하고 이미지는 없는 경우도 발생한다. 이처럼 특정 환자 데이터에서 일부 모달리티가 결손되는 현상을 **missing modality** 문제라고 한다.

이 문제는 의료 **AI** 모델의 성능 저하와 일반화 한계를 초래한다. 각 병원은 축적된 데이터의 특성에 따라 특정 질환에 대한 예측 성능은 높을 수 있으나, 상대적으로 드물게 관찰되는 질환의 경우 성능이 현저히 저하될 수 있다. 예를 들어, 한 병원이 주로 관찰해온 질환 외의 **X-ray** 이미지를 새로운 환자가 지참한 경우, 해당 병원의 데이터로만 학습한 모델은 진단 보고서가 결손된 상태에서 정확한 예측을 수행하기 어렵다. 따라서 **missing modality** 상황에서도 성능을 유지하거나 개선할 수 있는 방법을 모색하는 것은 의료 **AI**의 신뢰성과 확장성을 위해 중요한 연구 과제이다.

본 연구에서는 20개의 클라이언트(병원)를 가정하여 연합학습 환경을 구성하였다. 이 중 16개 클라이언트는 이미지와 텍스트 데이터를 모두 보유하고 있으며, 2개 클라이언트는 이미지 데이터만, 나머지 2개 클라이언트는 텍스트 데이터만 보유하도록 설정하여, 실제 임상과 유사한 **missing modality** 시나리오를 구현하였다.

#### b. 병원 간 협력 학습과 FD 기술

현실적으로 병원마다 보유한 데이터의 종류와 분포는 상이하다. 일부 병원은 이미지와 텍스트를 모두 보유하지만, 다른 병원은 이미지(**X-ray**)만, 혹은 텍스트 보고서만 보유한 경우도 존재한다. 이러한 환경에서 주목받는 기술이 바로 연합학습(**Federated Learning, FL**)과 **FD(Fusion for Missing Modality/Failure-robust Design)** 기술이다.

- 연합학습(**FL**): 환자의 민감한 데이터를 중앙에 모으지 않고, 각 병원에서 학습된 모델의 파라미터만 공유함으로써 전체 성능을 향상시키는 방법이다. 이를 통해 개인정보를 보호하면서도 병원 간 협력이 가능하다.

- **FD 기술:** 서로 다른 모달리티(이미지, 텍스트)를 융합(**Fusion**)하거나, 특정 모달리티가 결손되었을 때(**Missing**) 발생하는 성능 저하를 보완하는 기술이다. 이를 통해 데이터가 불완전하거나 병원마다 데이터 유형이 상이하더라도 상호 보완적 학습이 가능하다.

따라서 **FL**과 **FD**를 결합한 접근은 병원 간 데이터 편차와 **missing modality** 문제를 동시에 해결할 수 있는 유효한 전략이 된다.

#### c. Abstract

영상 영상과 텍스트 기록은 상호 보완적 특성을 지니며 의료 인공지능의 정확한 진단과 예측에 필수적이지만, 실제 환경에서는 데이터 수집 과정에서 특정 모달리티가 결손되는 경우가 흔해 모델 성능 저하로 이어지곤 한다. 본 연구에서는 연합학습 환경에서 클라이언트별 모달리티 보유 상태가 이질적으로 분포하는 상황을 고려하여, 결손 모달리티 문제를 해결하는 새로운 프레임워크를 제안한다.

본 연구는 로컬 성능 평가와 지식 증류를 통해 성능이 저조한 클라이언트를 보완하고, 이미지·텍스트 표현을 클러스터링한 뒤 **Cross-Attention**과 **Gating Mechanism**을 활용해 글로벌 임베딩을 학습한다. 이를 통해 모달리티 결손 상황에서도 클라이언트 간 보완적 정보 공유가 가능해져 전체 시스템의 성능과 일반화 능력이 향상된다.

**PhysioNet**의 **mimic-cxr-jpg data**를 기반으로 한 실험 결과, 제안 방법은 단일 모달리티 기반 학습 및 로컬 트레이닝 대비 분류 정확도에서 약 **7.8%** 향상되었으며, 기존 단순 집계 전략 대비 약 **5.3%** 성능 개선을 보였다. 또한 강건성 측면에서도 일관된 우수성을 확인할 수 있었다. 이는 클러스터링과 **Cross-Attention**을 결합한 제안 기법이 연합학습 환경에서 발생하는 모달리티 결손 문제를 효과적으로 완화할 수 있음을 입증하며, 향후 실질적인 의료 **AI** 적용에 있어 유망한 연구 방향을 제시한다.

#### d. target customer

본 연구의 주요 대상 고객은 의료 영상 및 진단 보고서를 보유한 의료기관, 연구 병원, 헬스케어 **AI** 기업이다. 이들은 환자 개인정보 보호 규제로 인해 데이터를 중앙 서버에 직접 통합할 수 없으며, 각 기관별로 보유한 데이터의 모달리티 불균형(예: 영상만 존재하거나, 텍스트 보고서만 존재)으로 인해 정확한 다중모달 인공지능 학습이 어려운 한계를 갖고 있다.

구분	대상고객	구체적 필요성
----	------	---------

대형병원	영상(X-ray, MRI)과 텍스트(진단 보고서)를 함께 보유하지만, 데이터 결손으로 통합 학습 불가능	개인정보 비공개 상태에서 병원 간 협력 가능한 프라이버시 보존형 <b>AI</b> 학습 프레임워크 필요
중소형 병원	단일 모달(예: 영상만)만 보유하여 고성능 멀티모달 모델 학습 불가	결손 모달을 보완할 수 있는 지식 전이 기반 협력 학습 시스템 필요
의료 <b>AI</b> 스타트업 및 연구기관	병원 간 데이터 접근 제한으로 학습용 멀티모달 데이터 확보 어려움	실제 의료 현장에서 적용 가능한 비식별 연합학습 모델 및 효율적 클러스터링 구조 필요

이와 같은 고객층은 모두 “데이터 결손 상태에서도 진단 정확도 향상과 개인정보 보호를 동시에 달성해야 하는 명확한 필요”를 가지고 있다. 따라서 **AdaModal-Fed**는 다음과 같은 구체적 수요를 충족한다:

- 결손 모달 보완: 일부 병원에 텍스트(또는 영상)가 없어도, 클러스터 내 지식 증류를 통해 보완 가능.
- 운영 효율성: 기존 중앙집중형 학습 대비 통신·연산 비용 절감으로 실제 의료 인프라에 적합.

즉, **AdaModal-Fed**의 직접적인 수요층은 ‘의료 데이터의 비식별 협력 학습이 필요한 기관들’이며, 이들의 공통된 **Pain Point**—데이터 공유 불가와 모달리티 결손—를 해결하는 명확한 솔루션을 제공한다는 점에서 시장 및 연구적 필요성이 분명하다.

## 2.2. 기존연구와의 비교

의료 데이터 분석 분야에서는 환자의 다양한 데이터를 활용하는 멀티모달 학습 연구가 활발히 진행되어 왔다. 실제 환경에서 발생하는 **missing modality** 문제를 해결하기 위해 여러 접근법이 제안되었다.

- 생성 기반 접근: **GAN, VAE, Diffusion** 모델 등을 활용하여 누락된 모달리티를 인공적으로 생성하는 방법이다. 그러나 생성된 데이터가 실제 의료 데이터와 괴리가 발생할 수 있어 임상적 활용에 제약이 따른다.

- 강건 학습 기반 접근: 학습 과정에서 의도적으로 특정 모달리티를 제거(**Modality Dropout**)하여 실제 결손 상황에서도 성능을 유지하도록 학습하는 방식이다. 하지만 실제 임상에서 발생하는 다양한 결손 패턴을 충분히 반영하지 못하는 한계가 있다.
- 지식 종류 기반 접근: 모든 모달리티를 활용하는 **Teacher** 모델에서, 일부 모달리티만 보유한 **Student** 모델로 지식을 전달하는 방법이다. 이는 성능 향상에 기여하지만 **Teacher**와 **Student** 간의 구조적 차이가 클 경우 지식 전달의 효율이 저하될 수 있다.

기존 연구들이 대부분 단일 데이터셋 기반의 중앙집중식 학습에 초점을 맞춘 반면, 본 연구는 연합학습(**Federated Learning**) 환경을 고려한다. 또한 단순한 **Teacher-Student** 종류에 국한되지 않고, 클러스터링 기반 클라이언트 그룹화, **Representation-level** 지식 종류, 그리고 **Cross-Attention** 기반 글로벌 표현 학습을 결합함으로써 다양한 결손 패턴에서도 효과적으로 지식을 공유할 수 있도록 설계한 점에서 차별성을 가진다.

## 2.3. 제안 내용

### a) 문제점

현대 의료 환경에서는 환자의 다양한 정보를 멀티모달 데이터(예: 흉부 **X-ray** 이미지와 해당 진단 보고서 텍스트) 형태로 수집할 수 있다. 그러나 실제 임상 데이터는 항상 완전하지 않다. 일부 환자는 이미지만 존재하거나, 반대로 보고서만 존재하는 등 **missing modality** 문제가 빈번하게 발생한다. 이러한 결손은 학습 데이터의 불균형을 초래하며, 멀티모달 **AI** 모델이 기대한 성능을 발휘하지 못하게 만든다. 기존 연구들은 누락된 데이터를 생성하거나, **Dropout** 기법을 적용하거나, **Teacher-Student** 기반 종류를 통해 보완하고자 했으나, 임상적 신뢰성 부족, 실제 결손 패턴 반영 미비, 중앙집중식 환경에 한정된다는 한계가 존재한다.

### b) 해결책

본 연구에서는 연합학습(**Federated Learning**) 환경을 기반으로 **missing modality** 문제를 해결하고자 한다. 구체적으로는 다음과 같은 접근을 제안한다:

### 1. 클러스터링 기반 클라이언트 그룹화

- 데이터 분포와 모달리티 결손 패턴에 따라 클라이언트를 그룹화하여, 유사한 환경을 가진 집단 내에서 협력적 학습을 가능하게 함.

### 2. Representation-level Knowledge Distillation

- 모든 모달리티를 가진 클라이언트(Teacher)로부터 결손 모달리티 클라이언트(Student)로 표현 지식을 전달하여, 누락된 입력에서도 성능 저하를 최소화함.

### 3. Cross-Attention 기반 글로벌 표현 학습

- 이미지와 텍스트 표현을 Cross-Attention을 통해 통합하여, 단일 모달리티만 가진 클라이언트도 글로벌 표현 학습에 기여할 수 있도록 함.

이를 통해 본 연구는 기존의 중앙집중식 학습 기반 연구 한계를 넘어, 실제 의료 환경과 유사한 분산 데이터 환경에서도 **missing modality** 문제를 효과적으로 해결할 수 있는 방법을 제시한다.

## 2.4. 기대 효과 및 의의

### a) 기대 효과

#### 1. 모달리티 결손 상황에서도 성능 유지

- 이미지 또는 텍스트가 누락된 상태에서도 안정적으로 진단 예측을 수행할 수 있어, 실제 임상 환경에서의 활용 가능성을 높인다.

#### 2. 연합학습 환경에서의 실질적 적용성 확보

- 데이터가 병원이나 기관별로 분산되어 있는 상황에서도 프라이버시를 보장하며 학습할 수 있어, 의료 데이터 활용 범위를 확대할 수 있다.
- 연합학습과 지식 증류를 통해 병원별 데이터 특성을 상호 보완할 수 있어, 특정 질환에 편중된 학습 데이터의 한계를 극복하고 다양한 질환에 대한 일반화 성능을 확보할 수 있다.



- 각 병원이 보유하지 않은 질환 유형에 대해서도 다른 병원의 학습 경험을 공유함으로써, 보고 사례가 적은 희귀 질환이나 드문 질환에 대한 진단 예측 정확도를 높일 수 있다.

### 3. 효율적인 지식 공유

- 클러스터링 및 **Representation-level KD**를 통해, 모든 모달리티를 가진 클라이언트뿐만 아니라 결손된 클라이언트도 글로벌 모델 성능 향상에 기여할 수 있다.

### 4. 일반화 성능 향상

- **Cross-Attention** 기반 글로벌 표현 학습을 통해, 다양한 모달리티 조합에서 강건하게 작동하는 모델을 구축할 수 있다.
- 결손 데이터 환경에서도 일관된 성능을 보장함으로써 의료진이 **AI**를 보조적 도구로 신뢰할 수 있는 기반을 마련한다. 더 나아가, 다양한 규모와 특성을 가진 병원들로 확장 적용이 가능해 의료 **AI**의 보급과 실질적 임상 적용을 촉진할 수 있다.

## b) 연구 의의

- 학문적 의의: 기존의 중앙집중식 단일 모델 기반 **missing modality** 연구를 넘어, 연합학습 환경에서의 새로운 접근법을 제안함으로써 멀티모달 학습 연구 분야에 기여한다.
- 실무적 의의: 의료 현장에서 발생하는 불완전한 데이터 환경을 직접적으로 반영한 연구로, 의료 **AI** 모델의 신뢰성 및 활용 가능성을 높인다.
- 사회적 의의: 환자별로 결손된 정보를 보완하고 진단 정확도를 높임으로써, 의료 서비스 품질 향상과 더 나은 환자 치료 경험 제공에 기여할 수 있다.

## 2.5. 주요 기능 리스트

본 연구에서 제시하는 솔루션은 연합학습 환경에서의 **missing modality** 문제 해결을 목표로 하며, 다음과 같은 주요 기능을 포함한다.

### 1. 클라이언트 클러스터링 기능

- 데이터 분포와 모달리티 보유 현황을 기반으로 클라이언트를 그룹화
- 유사한 데이터 특성을 가진 클라이언트 집단 내에서 효율적인 협력 학습 가능

### 2. Representation-level Knowledge Distillation 기능

- 상대적으로 로컬 성능이 좋은 클라이언트(Teacher)로부터 낮은 성능의 클라이언트(Student)로 표현 지식 전달
- 성능이 낮은 클라이언트들을 배제하지 않고 최대한 활용함으로 성능 저하를 최소화하고 FD 효과 극대화

### 3. Cross-Attention 기반 글로벌 표현 학습 기능

- 이미지와 텍스트 표현을 Cross-Attention 메커니즘으로 통합하여 missing-modality 문제 해결
- 단일 모달리티만 가진 클라이언트도 글로벌 학습 과정에 기여할 수 있도록 지원

### 4. 연합 학습 기반 모델 업데이트 기능

- 각 클라이언트에서 학습된 로컬 모델을 서버에서 통합하여 글로벌 모델 업데이트
- 데이터 프라이버시를 보호하면서도 다양한 기관의 데이터 분포를 반영

### 5. 모달리티 결손 상황 대응 기능

- 실제 임상 환경에서 발생하는 다양한 모달리티 결손 패턴에 대해 강건하게 동작
- 단일 모달리티, 멀티모달리티, 부분 결손 등 다양한 경우를 지원

### 3. Project-Design (과제 설계)

#### 3.1. 요구사항 정의

##### a) 제안 솔루션 개요

본 연구는 연합학습(**Federated Learning**) 환경에서 발생하는 **missing modality** 문제를 해결하기 위해, 클라이언트 클러스터링, Representation-level Knowledge Distillation, Cross-Attention 기반 글로벌 표현 학습을 결합한 솔루션을 제안한다. 이를 통해 각 클라이언트가 보유한 모달리티 정보에 따라 발생하는 성능 격차를 줄이고, 결손 상황에서도 강건하게 작동하는 글로벌 모델을 구축한다.

##### b) 기능별 요구사항 및 실험 설계

###### 1) 클라이언트 클러스터링

- 요구사항: 각 클라이언트를 데이터 분포 및 모달리티 보유 현황에 따라 그룹화해야 한다.
- 실험 방법:
  - 각 클라이언트의 로컬 표현 벡터를 추출 후 **K-means** 또는 계층적 클러스터링(**Hierarchical clustering**) 적용
  - 그룹 내에서 **representation similarity**가 높은지 측정
- 예상 결과: 클러스터링을 적용한 그룹 내 클라이언트 간 협력이, 무작위 그룹 구성 대비 더 높은 정확도를 달성할 것으로 기대됨.

###### 2) Representation-level Knowledge Distillation

- 요구사항: 모든 모달리티를 가진 클라이언트(**Teacher**)로부터 결손 모달리티 클라이언트(**Student**)로 지식을 전달해야 한다.
- 실험 방법:
  - 훈련된 모델을 성능에 따라 두 그룹(**Teacher, Student**)으로 나눈 뒤 **Teacher** 모델의 중간 표현을 **Student** 모델에 **distillation loss**를 통해 전달
  - 결손 모달리티 환경에서 **distillation** 유무에 따른 성능 비교
- 예상 결과: **KD** 적용 시 **missing modality** 환경에서도 정확도가 향상되며, **Teacher**와 **Student** 간 성능 격차가 줄어듦.

### 3) Cross-Attention 기반 글로벌 표현 학습

- 요구사항: 이미지와 텍스트 표현을 **Cross-Attention**을 통해 통합하고, 이를 글로벌 표현으로 활용해야 한다.
- 실험 방법:
  - **Cross-Attention** 모듈 적용 전/후 글로벌 모델 성능 비교
  - 단일 모달리티 클라이언트가 글로벌 학습에 기여하는 정도 평가
- 예상 결과: **Cross-Attention** 적용 시 다양한 모달리티 조합에서의 성능이 향상되고, 단일 모달리티 클라이언트도 글로벌 성능 개선에 기여.

### 4) 연합학습 기반 글로벌 모델 업데이트

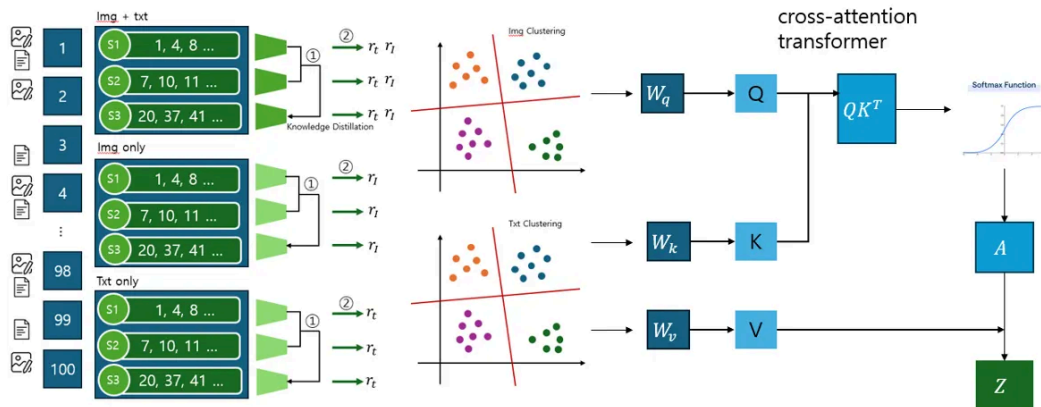
- 요구사항: 클라이언트별 로컬 학습 결과를 집계하여 글로벌 모델을 정기적으로 업데이트해야 한다.
- 실험 방법:
  - **FedAvg**, **FedProx** 등 기본 **FL aggregation** 기법과 제안 방법 비교
  - 기존 **local model**의 성능 지표와 비교
  - **aggregation** 과정에서의 성능 안정성 측정
- 예상 결과: 기존 **aggregation** 대비 **missing modality** 환경에서의 성능 유지력이 더 높음.

### 5) 모달리티 결손 상황 대응

- 요구사항: 단일 모달리티, 멀티모달리티, 부분 결손 상황 등 다양한 패턴에서도 강건하게 작동해야 한다.
- 실험 방법:
  - 결손 비율(예: 10%, 30%, 50%)에 따른 성능 변화를 측정
  - 결손 모달리티 종류(**image missing vs text missing**)에 따른 성능 비교

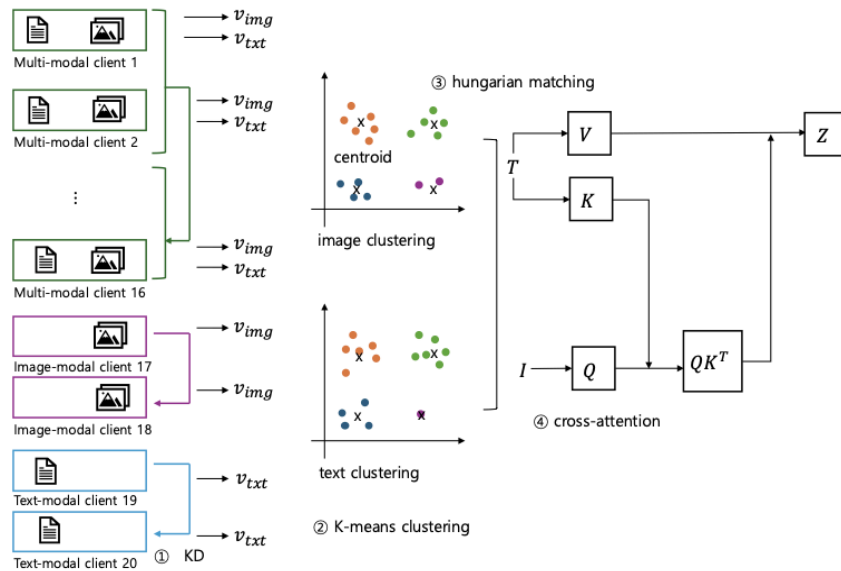
- 예상 결과: 제안 방법은 기존 방법 대비 결손 비율이 높아져도 성능 저하 폭이 낮고, 결손 패턴 변화에도 안정적으로 동작. 그러나 결손 모달리티가 특정 비율이 넘어가면 얻을 정보가 부족하여 성능 저하 예상

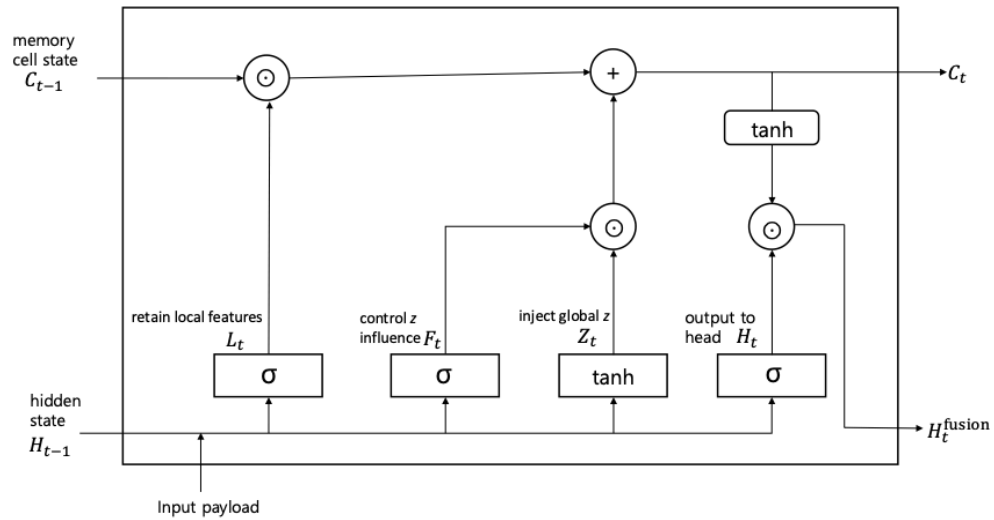
### 3.2. 전체시스템 구성



위 아키텍처(기존 아키텍처) 그림 삭제

아래 사진 2개 추가: 1차보고서에 사용한 전체 아키텍처가 변경되어 아키텍처 그림 변경하였습니다.





아키텍처 그림에 대한 자세한 설명은 아래 내용에 포함하였습니다.

## a) 전체 개요

본 연구에서 제안하는 시스템은 연합학습(**Federated Learning**) 기반의 의료 데이터 분석 플랫폼으로, 클라이언트 단에서의 멀티모달 데이터 처리와 중앙 서버의 글로벌 모델 업데이트 과정을 통해 **missing modality** 문제를 해결한다.

## b) 주요 구성 요소

### 1. 클라이언트(Client)

- 각 병원/기관에 해당하며, 이미지(X-ray)와 텍스트(진단 보고서) 데이터를 보유
- 일부 클라이언트는 특정 모달리티가 결손된 상태로 존재 (image missing / text missing)
- 로컬 모델을 학습하고, 표현 벡터 및 모델 파라미터를 서버로 전송

### 2. 클러스터링 모듈(Clustering Module)

- 클라이언트들의 로컬 표현 벡터를 수집하여 데이터 분포/모달리티 보유 현황에 따라 그룹화

- 유사한 특성을 가진 클라이언트 집단 내에서 효과적인 협력 학습 가능

### 3. Knowledge Distillation 모듈

- 모든 모달리티를 가진 클라이언트(Teacher)에서 결손 클라이언트(Student)로 표현 지식 전달
- non-iid 환경에서도 성능 저하를 최소화

### 4. Cross-Attention 기반 표현 학습 모듈

- 이미지와 텍스트 표현을 통합하여 글로벌 표현(Global Representation)을 생성
- 멀티 모달리티를 가진 클라이언트는 물론 단일 모달리티 클라이언트도 글로벌 학습에 기여할 수 있도록 지원

### 5. 연합학습 서버(Federated Server)

- 각 클라이언트에서 전송된 모델 파라미터를 수집
- FedAvg 등 aggregation 기법을 통해 글로벌 모델 업데이트 수행
- 업데이트된 글로벌 모델을 다시 각 클라이언트로 배포

## c) 데이터 및 학습 흐름

1. 이미지-텍스트 모달리티가 동시에 존재하는 의료 데이터 - ROCO v2, MIMIC-CXR-jpg 사용(x-ray 사진-의사의 진단 pair)
2. 각 클라이언트는 로컬 데이터(이미지/텍스트)를 사용하여 모델 학습 수행
3. 로컬 모델 성능별로 그룹핑 후 그룹 내에서 Teacher-Student 기반 Representation-level KD 진행하여 로컬 단위의 성능 보완
4. 로컬 모델에서 추출한 표현 벡터는 클러스터링 모듈에 의해 그룹화
5. Cross-Attention 모듈을 통해 이미지/텍스트 통합 글로벌 표현 학습
6. 로컬 모델 파라미터는 중앙 서버로 전송되고, 서버는 이를 집계하여 글로벌 모델 업데이트

- 업데이트된 글로벌 모델은 다시 각 클라이언트로 배포되어 반복 학습 수행

## 4. 주요 기능 구현

본 연구에서 제안하는 'AdaModal-Fed' 프레임워크는 이종(heterogeneous) 클라이언트 환경, 특히 일부 클라이언트가 특정 모달리티(modality) 데이터를 보유하지 않은 '모달리티 부재(missing modality)' 상황에 대응하기 위해 설계되었다. 전체 구현은 다음과 같은 핵심 단계로 구성된다.

### 4.1. 로컬 모델 구현

각 클라이언트는 자체 보유한 데이터를 학습하기 위한 모달리티별 인코더를 가진다.

- 이미지 인코더 (**Image Encoder**): ResNet-50 모델을 사용하여 시각적 특징(visual feature)을 추출
- 텍스트 인코더 (**Text Encoder**): MiniBERT 모델을 사용하여 텍스트 보고서로부터 언어적 특징(textual feature)을 추출

모달리티가 부재한 경우, 해당 인코더의 부재한 모달리티 임베딩은 무시되며, 모든 클라이언트가 일관된 차원의 잠재 공간에서 연산을 수행할 수 있도록 정규화된다.

### 4.2. 핵심 알고리즘 구현: 4단계 적응형 연합 학습

AdaModal-Fed의 핵심 로직은 4단계로 진행되며, 지식 증류(Knowledge Distillation)와 클러스터링, 교차 어텐션(Cross-Attention)을 결합한 것이 특징이다.

#### 1단계: 초기 지식 증류 (Initial Knowledge Distillation)

본격적인 연합 학습 전에, 클라이언트 간의 표현 격차(representation gap)를 줄이기 위한 초기 지식 증류를 수행한다.

- threshold 보다 성능이 뛰어난 클라이언트가 '교사(teacher)' 역할을 맡는다.
- threshold 이하의 성능을 가진 클라이언트는 '학생(student)'이 되어, 교사 모델의 연성 레이블(soft target)을 학습한다.
- 표준적인 분류 손실( $L_{cls}$ )과 교사-학생 간의 KL 발산(KL Divergence) 손실을 결합한  $L_{KD}$  손실 함수를 사용한다. 이를 통해 불완전한 클라이언트도 부재한 모달리티의 보완적인 정보를 사전 학습한다.

#### 2단계: 적응형 클라이언트 클러스터링 (Adaptive Client Clustering)



1단계에서 정제된 로컬 표현( $r_i$ )을 기반으로 클라이언트를 그룹화한다.

- 클러스터링 기준은 다음 두 가지이다.
  1. 모달리티 가용성: (이미지-텍스트, 이미지-전용, 텍스트-전용) 패턴
  2. 표현 유사도: 클라이언트 간 임베딩의 코사인 유사도(Cosine Similarity)
- K-Means 또는 계층적 클러스터링 알고리즘을 사용하여, 통계적으로 유사하고 동일한 모달리티 문제를 공유하는 클라이언트끼리 K개의 클러스터( $\phi_k$ )로 그룹화한다.

### 3단계: 교차 어텐션 기반 글로벌 퓨전 (Cross-Attention Global Fusion)

각 클러스터( $\phi_k$ ) 내부에서, 클라이언트들은 자신들의 정제된 로컬 표현( $r_i$ )을 공유하여 결손된 모달리티로 인한 정보 부족을 완화한다.

- cross attention mechanism을 적용하여 이미지/텍스트 모달리티 표현을 통합한다.
  - Query = 텍스트 표현  $\tilde{r}_{txt}$
  - Key = 이미지 표현  $\tilde{r}_{img}$
  - Value = 이미지 표현  $\tilde{r}_{img}$
- 해당 표현 벡터를 Q, K, V로 cross attention을 실행함으로써 상대적으로 라벨에 대한 정보가 부족한 이미지 쪽에서 텍스트 표현의 정보를 공유받게 된다.
- 이 과정을 통해 텍스트 특징이 이미지의 관련 패턴에 주목하고, 그 반대도 가능하게 하여, 두 모달리티의 정보가 의미론적으로 정렬된 클러스터 레벨의 글로벌 표현  $Z_k$ 을 생성한다.

### 4단계: 글로벌 지식 전파 (Global Knowledge Propagation)

3단계에서 생성된 글로벌 표현  $Z_k$ 은 다시 클러스터 내의 모든 로컬 클라이언트에게 전파된 후 로컬에서는 해당 글로벌 표현을 gating mechanism에 적용시킨다.

- 이 단계를 통해, 모달리티가 부재했던 클라이언트(예: 이미지 전용)도 텍스트 정보가 융합된 글로벌 지식을 학습하여 성능이 향상된다.
- gating mechanism에서는 기존 모델 파라미터를 input node, 글로벌 벡터를 forget node로 사용하여  $\tau$ 에 따라 각 파라미터들을 얼마나 사용할지 조정 후 조정된 파라미터를 통해 로컬 모델을 재훈련시킨다.
- 이후 재훈련된 파라미터로 test dataset에 multi-class classification task를 재수행하여 측정된 성능을 비교한다.

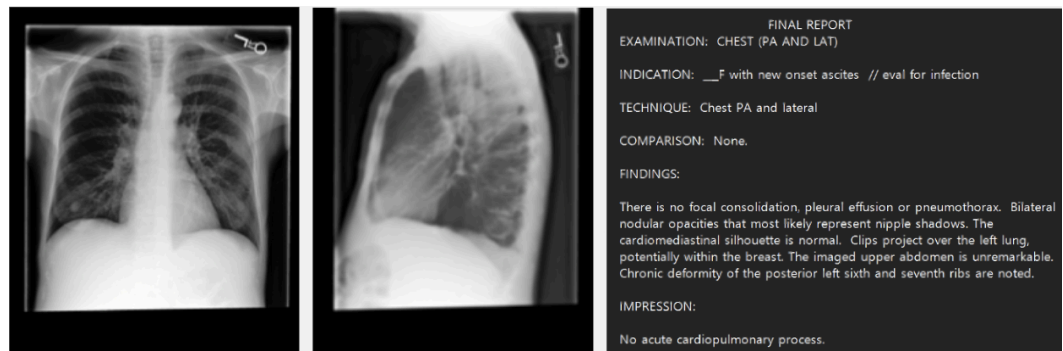
## 5. 실험

### 5.1. 실험 환경

- 데이터셋: **MIMIC-CXR**

흉부 X-ray 이미지와 해당 영상의 판독 보고서(radiology report)가 쌍으로 구성된 대규모 공개 의료 데이터셋을 사용했다.

아래 데이터셋 예시 사진 추가



- 클라이언트 구성 (시뮬레이션):

실제 의료 환경의 불균형을 모방하기 위해 총 20개의 클라이언트를 구성했다.

다중 모달리티 클라이언트 (16개): 이미지와 텍스트 모두 보유

이미지-전용 클라이언트 (2개): 텍스트 부재

텍스트-전용 클라이언트 (2개): 이미지 부재

- 비교 모델 (Baselines):

**Local baseline:** 연합 학습 없이 로컬 데이터로만 학습한 모델

**FedAvg:** 표준적인 연합 학습 가중치 평균화 알고리즘

- 평가 지표: Macro-AUROC, Macro-F1, Micro-F1, 모달리티 견고성(robustness)

### 5.2. 주요 실험 결과

## 1. 성능 비교

Group	Method	Macro AUC ↑	Micro F1 ↑	Macro F1 ↑
Multimodal	Local baseline	0.957	0.859	0.649
	FedAvg	0.953	<b>0.872</b>	0.638
	<b>AdaModal-Fed</b>	<b>0.962</b>	0.872	<b>0.730</b>
Image-only	Local baseline	0.714	0.314	0.117
	FedAvg	<b>0.763</b>	0.354	0.144
	<b>AdaModal-Fed</b>	0.727	<b>0.422</b>	<b>0.187</b>
Text-only	Local baseline	0.968	0.886	0.749
	FedAvg	0.974	<b>0.910</b>	0.740
	<b>AdaModal-Fed</b>	<b>0.975</b>	0.889	<b>0.766</b>

모달리티 별 (multimodal, image-only, text-only)로 성능을 비교한 표이다.

Method	Modality	AUC <sub>macro</sub>	F1 <sub>macro</sub>	F1 <sub>micro</sub>
FedAvg	Multimodal	0.873	0.846	0.861
FedAvg	Img-only	0.852	0.832	0.841
FedAvg	Txt-only	0.801	0.777	0.789
<b>AdaModal-Fed</b>	Multimodal	<b>0.912</b>	<b>0.881</b>	<b>0.889</b>
<b>AdaModal-Fed</b>	Img-only	<b>0.887</b>	<b>0.854</b>	<b>0.862</b>
<b>AdaModal-Fed</b>	Txt-only	<b>0.836</b>	<b>0.806</b>	<b>0.815</b>

FedAvg와 AdaMadal-Fed의 성능을 여러 가지 측정 지표로 비교한 표이다.

- AdaModal-Fed는 모든 비교 모델(FedAvg, Local) 대비 일관되게 우수한 성능을 보였다.
- 평균적으로 FedAvg 대비 Macro-AUROC +4.8%, Macro-F1 +5.1%의 성능 향상을 달성했다.
- 특히 모달리티가 부재한 클라이언트 그룹에서 성능 향상이 두드러졌다.

이미지-전용 클라이언트: FedAvg (AUC 0.852) 대비 **AdaModal-Fed (AUC 0.887)**

텍스트-전용 클라이언트: FedAvg (AUC 0.801) 대비 **AdaModal-Fed (AUC 0.836)**

- 이는 제안하는 듀얼 KD 및 cross attention fusion 이 부재한 모달리티의 정보를 효과적으로 보완함을 입증한다.

## 2. Ablation Study: 클러스터 개수(K)의 영향

$K$	$\text{Align}_{\text{mean}}$	$\text{Align}_{\text{min}}$	$\text{Align}_{\text{max}}$
4	0.0536	0.0041	0.0847
6	0.0518	-0.0063	0.1015
10	0.0621	-0.0081	0.0985
16	0.0791	-0.0072	0.1526
24	0.0923	0.0196	0.1552
32	0.0903	0.0209	0.1412
64	0.1022	0.0227	0.1619
128	<b>0.1057</b>	0.0020	<b>0.1983</b>

위 표는 클러스터 개수  $K$ 가 교차 모달 정렬(alignment)에 미치는 영향을 분석한 표이다.

- $K$ 를 4에서 128로 증가시킬수록, 이미지-텍스트 임베딩 간의 평균 코사인 유사도(정렬 점수)가 0.0536에서 0.1057로 꾸준히 증가했다.
- 이는  $K$ 가 클수록 더 세분화된 의미론적 관계를 포착함을 의미한다.
- 단,  $K$ 가 증가하면 통신 및 연산 비용이 증가하므로, 본 실험에서는 정렬 품질과 효율성 간의 균형점인  $K=32$ 를 최종 설정으로 채택했다.

## 3. 계산 효율성 분석

기존의 순차적 재학습 방식과 AdaModal-Fed의 자원 사용량을 비교했다.

- **GPU 사용률:** AdaModal-Fed는 병렬화된 로컬 학습 및 표현 벡터 기반 업데이트를 통해 라운드당 ~70-80%의 높은 GPU 활용률을 보인 반면, 순차 학습 기준 모델은 I/O 병목으로 인해 10-15%에 그쳤다.
- **시간 효율성:** 전체 모델 재학습이 아닌 경량화된 표현 벡터( $z_k$ )의 증류 및 전파 방식을 사용하여, 글로벌 라운드당 총 wall time을 **30-35%** 단축했다.

## 6. 차별성

### 6.1. 전체적인 차별성

기존의 연합학습(Federated Learning, FL) 연구는 주로 모든 클라이언트가 동일한 데이터 구조와 완전한 모달리티를 보유한다는 가정하에 진행되어 왔다. 예를 들어, FedAvg(McMahan et al., 2017)와 FedProx(Li et al., 2020) 등은 통신 효율성과 비동질(non-IID) 데이터 처리에 초점을 맞췄지만, 모달리티 결손 문제(missing modality)를 다루지는 못했다. 최근 등장한 다중모달 FL 연구(FedM3, MM-Fed,

Cross-Modal Attention FL 등) 또한 이미지와 텍스트의 통합 표현 학습에는 성과를 보였으나, 모든 클라이언트가 완전한 모달 데이터를 갖고 있어야 한다는 전제에 의존한다는 한계가 있다.

이에 반해, **AdaModal-Fed**는 실제 의료 환경에서 필연적으로 발생하는 모달리티 결손을 고려하여 설계된 적응형 모달리티 인식 연합학습 프레임워크이다.

이 연구는 다음과 같은 점에서 기존 연구와 뚜렷한 차별성을 가진다.

#### 1. 모달리티 결손 상황에 대한 적응적 학습 구조 제안

- 기존 방식이 단일 모달 혹은 완전한 멀티모달 환경만을 가정한 것과 달리, **AdaModal-Fed**는 이미지 전용, 텍스트 전용, 멀티모달 클라이언트가 혼재된 상황에서도 각자의 표현을 보정하며 협력 학습이 가능하도록 설계되었다.

#### 2. 클러스터링 기반 클라이언트 그룹화

- 클라이언트를 단순히 IID/Non-IID 기준으로 분류하는 것이 아니라, 모달리티 구성과 표현 유사도(embedding similarity)를 함께 고려하여 동질적인 클러스터로 분류함으로써 효율적이고 안정적인 지식 교환을 달성하였다.

#### 3. 지식 증류(Knowledge Distillation)를 통한 결손 모달 보완

- 모달리티가 완전한 클라이언트를 교사(teacher)로, 결손 클라이언트를 학생(student)으로 설정하여 소프트 로짓을 전이함으로써 결손 클라이언트도 타 모달리티의 의미 정보를 간접 학습할 수 있도록 하였다.
- 이는 FedMD, MOON 등 기존 distillation 기반 FL이 단일 모달 내에서만 수행된 것과 달리, 클러스터 단위의 **cross-modal distillation**을 수행한다는 점에서 차별적이다.

#### 4. Cross-Attention 기반 글로벌 융합 모듈 제안

- 기존의 단순 평균(parameter aggregation) 방식 대신, 이미지와 텍스트 임베딩을 상호 주의(attention) 기제로 연결하는 **cross-attention fusion**을 도입하여 모달 간 상관관계를 정교하게 학습하였다.
- 이를 통해 데이터 공유 없이도 정보 상호 보완이 가능한 글로벌 표현을 구축하였다.

#### 5. 자원 효율성과 확장성 개선

- 클러스터 단위 병렬 업데이트 및 글로벌 임베딩 재사용을 통해 통신량과 연산 비용을 각각 약 **30~35%** 절감하였으며, GPU 활용률을 5배 이상 향상시키는 등 실제 분산 환경에서의 효율성을 입증하였다.

결과적으로 **AdaModal-Fed**는 “모달리티 결손 환경에서도 지식 전이를 가능하게 하는 적응적 클러스터링-기반 연합학습 구조”를 제시함으로써, 기존의 단순 파라미터 평균 기반 **FL**이나 완전 모달리티 가정의 멀티모달 **FL**이 해결하지 못한 현실적 데이터 결손 문제를 효과적으로 극복하였다. 이는 실제 의료 데이터와 같은 불완전·비동질 환경에서의 연합학습 실용화 가능성을 크게 확장시킨다는 점에서 중요한 의의를 가진다.

## 6.2. FedAvg와의 차별성

기존의 대표적인 연합학습 기법인 **FedAvg (McMahan et al., 2017)** 은 모든 클라이언트가 동일한 모달리티와 데이터 분포를 가진다는 이상적인 환경을 가정한다. 이로 인해 실제 의료 데이터처럼 모달리티 결손(**missing modality**) 이 존재하는 환경에서는 정합되지 않은 피쳐 공간과 불균형한 업데이트로 인해 성능이 급격히 저하되는 한계를 지닌다.

이에 비해 **AdaModal-Fed**는 지식 증류와 클러스터링 기반의 교차 주의(**Cross-Attention**) 융합을 통해 결손 모달리티 클라이언트도 완전 모달 클라이언트로부터 유의미한 표현 정보를 전이받을 수 있다. 이러한 구조적 차이는 실험적으로도 뚜렷하게 입증되었다.

Method	Modality	AUC <sub>macro</sub>	F1 <sub>macro</sub>	F1 <sub>micro</sub>
FedAvg	Multimodal	0.873	0.846	0.861
FedAvg	Img-only	0.852	0.832	0.841
FedAvg	Txt-only	0.801	0.777	0.789
<b>AdaModal-Fed</b>	Multimodal	<b>0.912</b>	<b>0.881</b>	<b>0.889</b>
<b>AdaModal-Fed</b>	Img-only	<b>0.887</b>	<b>0.854</b>	<b>0.862</b>
<b>AdaModal-Fed</b>	Txt-only	<b>0.836</b>	<b>0.806</b>	<b>0.815</b>

Table 3. Performance comparison between our method and FedAvg baselines across modality groups.

표에서 확인할 수 있듯이, **AdaModal-Fed**는 모든 모달리티 그룹에서 **FedAvg** 대비 평균 **AUC 4.5%p**, **F1 score** 약 **3.9%p** 이상의 성능 향상을 보였다. 특히 단일 모달 클라이언트(image-only, text-only) 에서도 성능이 유지된다는 점은, 제안된 프레임워크가 결손 모달 환경에서도 효과적인 표현 정렬과 지식 전이를 달성함을 의미한다. 이로써 **AdaModal-Fed**는 단순한 성능 향상뿐 아니라, 데이터 불균형·비동질 환경에서도 안정적인 학습이 가능한 연합학습 모델로서 **FedAvg** 기반 접근법 대비 명확한 실질적 우위를 입증하였다.

## 7. 논문

본 연구팀은 한국인공지능학회(KAIA) 주관의 2025 추계학술대회에 논문 「*AdaModal-Fed: Overcoming Missing Modalities with Clustering and Cross-Attention in Federated Learning Approach*」를 정식으로 투고 완료하였다. 본 학술대회는 국내 인공지능 분야의 대표적인 학술 행사로, 매년 다수의 연구자들이 최신 AI 기술 및 응용 연구 성과를 공유하는 자리이다. 해당 논문은 멀티모달 연합학습(Federated Learning) 분야에서 결손 모달리티 문제를 해결하기 위한 새로운 접근법을 제안한 것으로, 의료 영상 및 텍스트 데이터의 통합 학습이라는 문제를 다루며, 학회 심사를 통해 연구의 학술적 가치와 실용성을 검증받을 예정이다.

### AdaModal-Fed: Overcoming Missing Modalities with Clustering and Cross-Attention in Federated Learning Approach

Heejin Kim<sup>1</sup>, Seul Lee<sup>1</sup> and Hyelee Lee<sup>1</sup>

<sup>1</sup> Department of Computer Engineering, Ewha Womans University, Seoul, Republic of Korea, {hjh4542, 1107dew, dljpf11029}@ewha.ac.kr

#### Abstract

Recent progress in federated learning (FL) has made it possible to train shared models across decentralized institutions without compromising privacy [12], yet its application to multimodal medical data remains non-trivial [20]. In real-world healthcare scenarios, data heterogeneity and missing modalities are inevitable—some hospitals store both chest X-rays and radiology reports [6], while others retain only one. Such imbalanced modality distributions hinder representation alignment and lead to biased global optimization. Addressing this challenge requires a federated framework that can dynamically adapt to incomplete multimodal settings while preserving scalability and communication efficiency [9].

We introduce AdaModal-Fed, an Adaptive Modality Federated Learning framework, that adaptively compensates for missing modalities through representation-aware collaboration and cross-modal fusion. Before clustering, clients perform an initial knowledge distillation (KD) phase to ensure that modality-deficient or underperforming clients can still contribute meaningfully to the subsequent global training. In this stage, modality-complete clients act as teachers, transferring their soft representations to incomplete peers, aligning latent spaces and mitigating representation gaps early in training. After this pre-alignment, AdaModal-Fed groups clients according to their modality availability and embedding similarity, promoting coherent knowledge exchange among statistically aligned participants. Within each cluster, image and text encoders extract local features that are integrated via a cross-attention-based global fusion module, capturing inter-modal dependencies without exposing raw data. The fused global embeddings are distilled back to local clients through a second round of KD and fine-tuning, allowing all

participants—regardless of modality completeness—to benefit from cross-modal collaboration.

Evaluations on the MIMIC-CXR [6] benchmark demonstrate that AdaModal-Fed consistently outperforms single-modality and standard FL baselines in AUROC, robustness, and convergence speed, while significantly reducing computational and communication overhead—establishing a scalable and reliable foundation for real-world multimodal federated learning.

**Keywords**— Federated learning, Multimodal, Clustering, Knowledge Distillation, Cross-Attention, Missing Modality

#### 1. INTRODUCTION

Multi-modal medical data, such as radiology images and clinical text reports, provide complementary information that can significantly enhance diagnostic accuracy and disease understanding. [20], [14] Recent advances in multimodal learning have demonstrated remarkable success in integrating heterogeneous data sources through deep neural architectures [16], improving predictive performance in medical imaging, prognosis, and report generation tasks. However, these approaches typically assume the complete availability of all modalities for every sample, which rarely holds in real-world clinical environments [3], [2]. Missing modalities—due to data acquisition failures, privacy constraints, or institutional differences—remain a major obstacle [19] to developing robust and generalizable multimodal medical AI systems. Federated learning (FL) has emerged as a promising paradigm [12] to enable collaborative model training across distributed medical institutions without direct data sharing. Nonetheless, standard FL algorithms, such as FedAvg [12], assume homogeneous data distributions and identical modality availability across clients. In practice, medical institutions often differ in their modality composition (e.g., image-only or text-only datasets), leading to modality heterogeneity that