

AAAI 2024



---

# On Disentanglement of Asymmetrical Knowledge Transfer for Modality-Task Agnostic Federated Learning

---

**Jiayi Chen<sup>1</sup>, Aidong Zhang<sup>2</sup>**

*<sup>1,2</sup>Department of Computer Science, University of Virginia*

<sup>1</sup>Contact: [jc4td@virginia.edu](mailto:jc4td@virginia.edu)

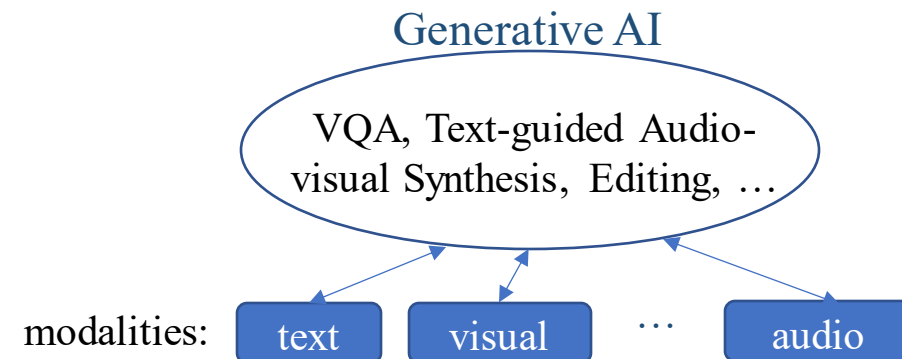
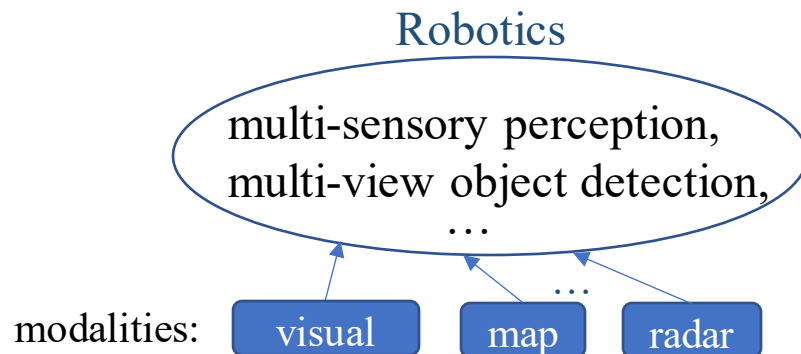
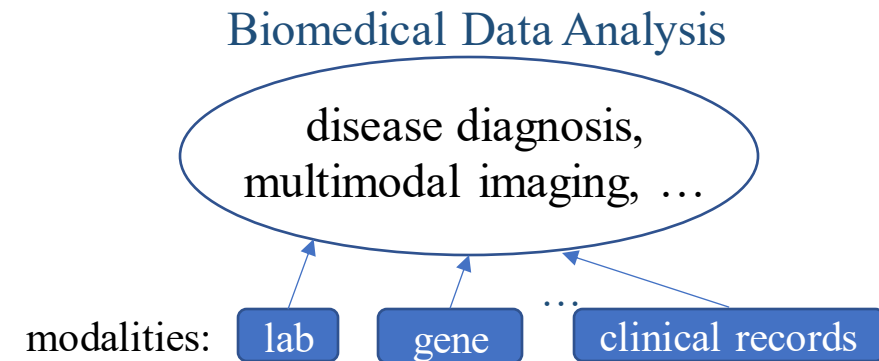
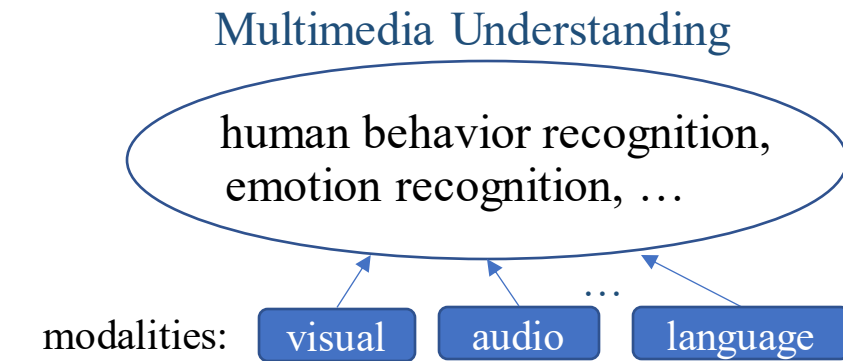
<sup>2</sup>Thomas M. Linville Professor



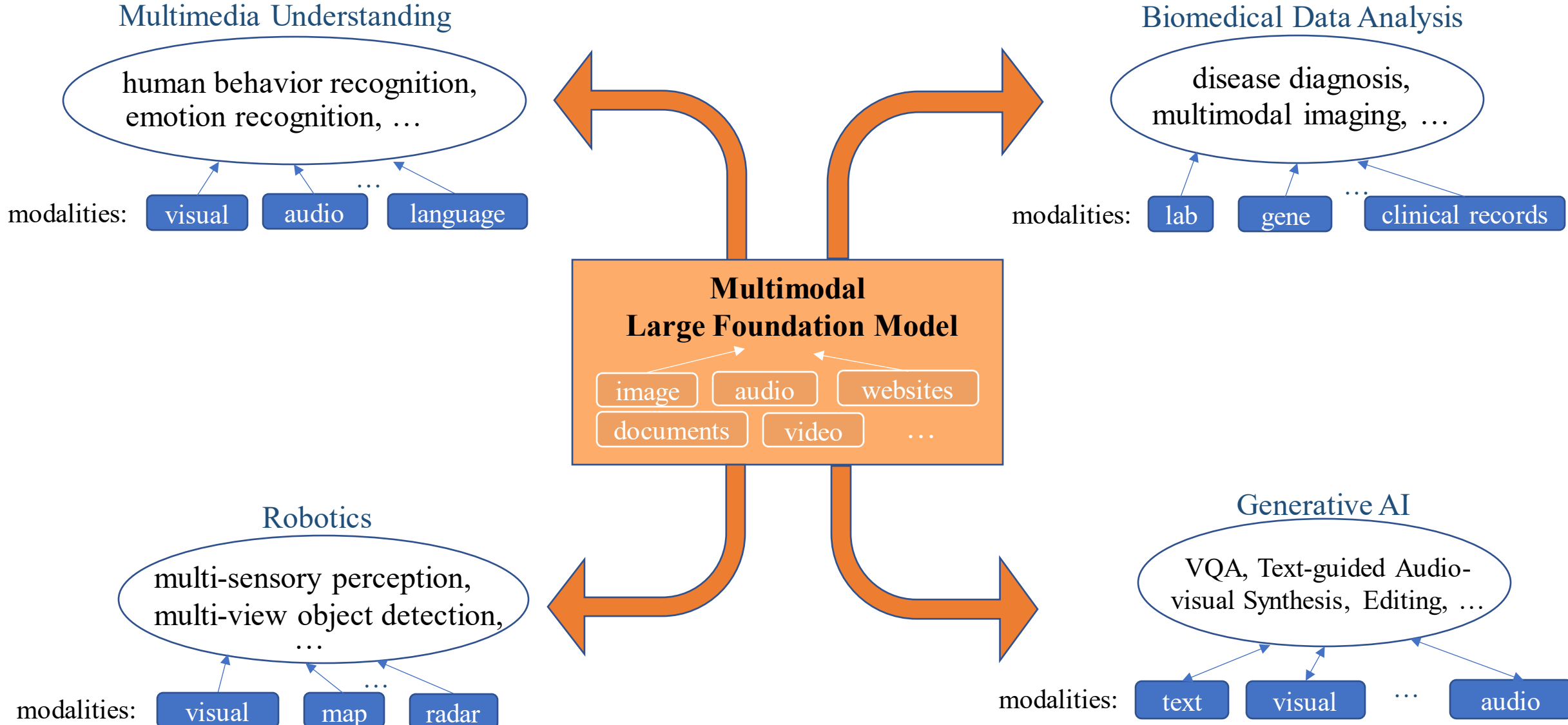
# Multimodal Artificial Intelligence

## ➤ Real-world Use Cases:

- *Discriminative tasks*: Multimodal/cross-modal Understanding, Alignment, Multimodal fusion, ...
- *Generative tasks*: Cross-modal guided data synthesis, Video captioning, grounding, ...

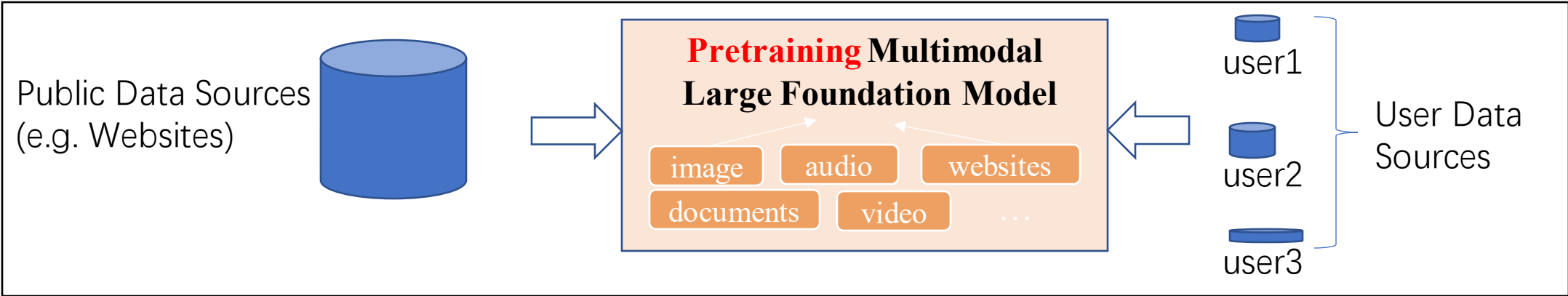


# Artificial General Intelligence (AGI)



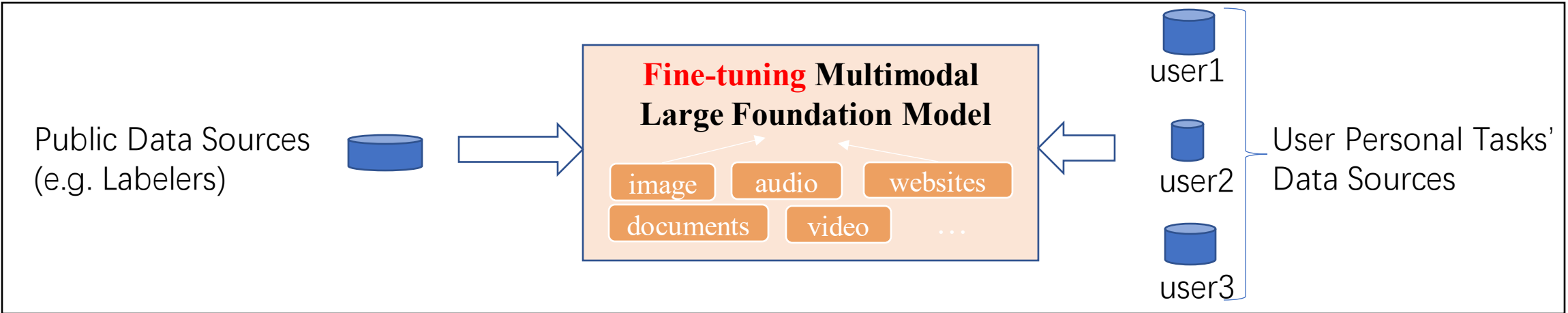
# User-AGI Interactions

## ➤ Pretraining Stage



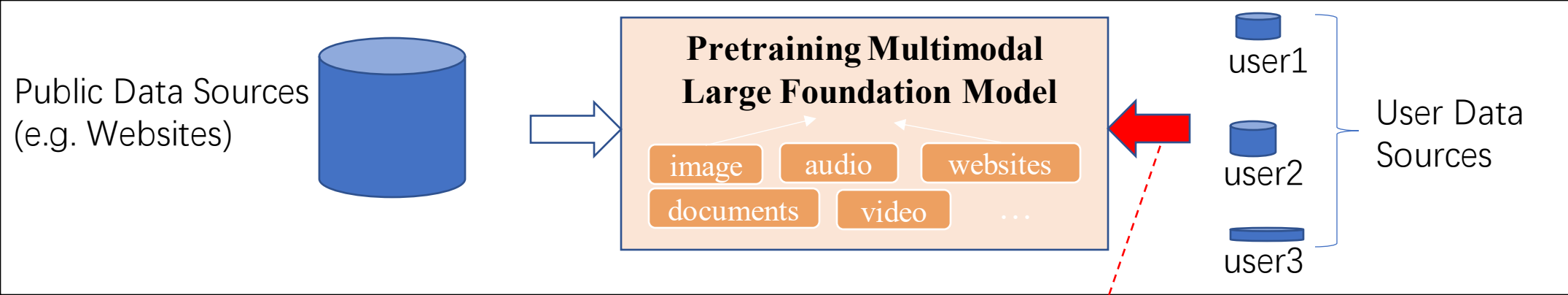
## ➤ Fine-tuning Stage

- *Training*: RLHF, meta-tuning, ...
- *Instruction/Data*: Prompts, RAG, ...
- *Model modification*: PEFT (Adaptors, Prefix Tuning, ...)

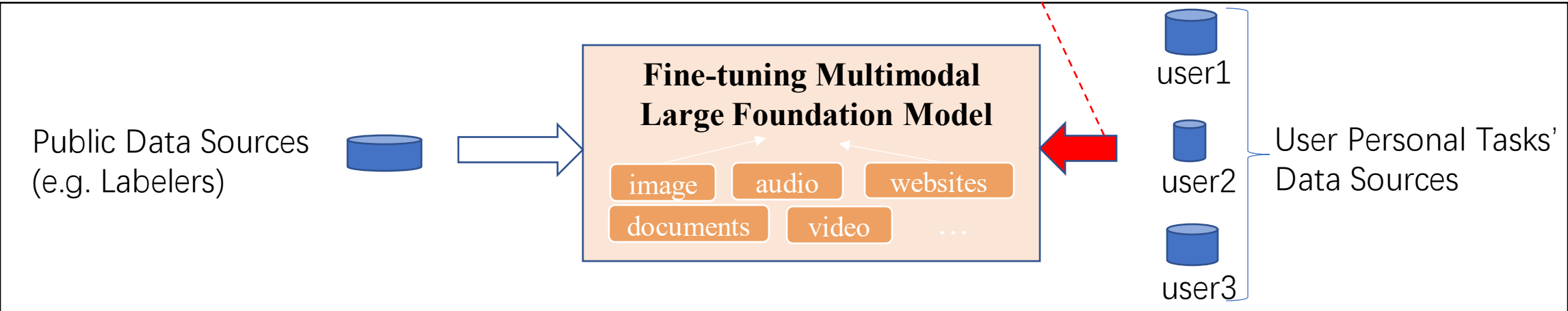


# Privacy Issue in User-AGI Interactions

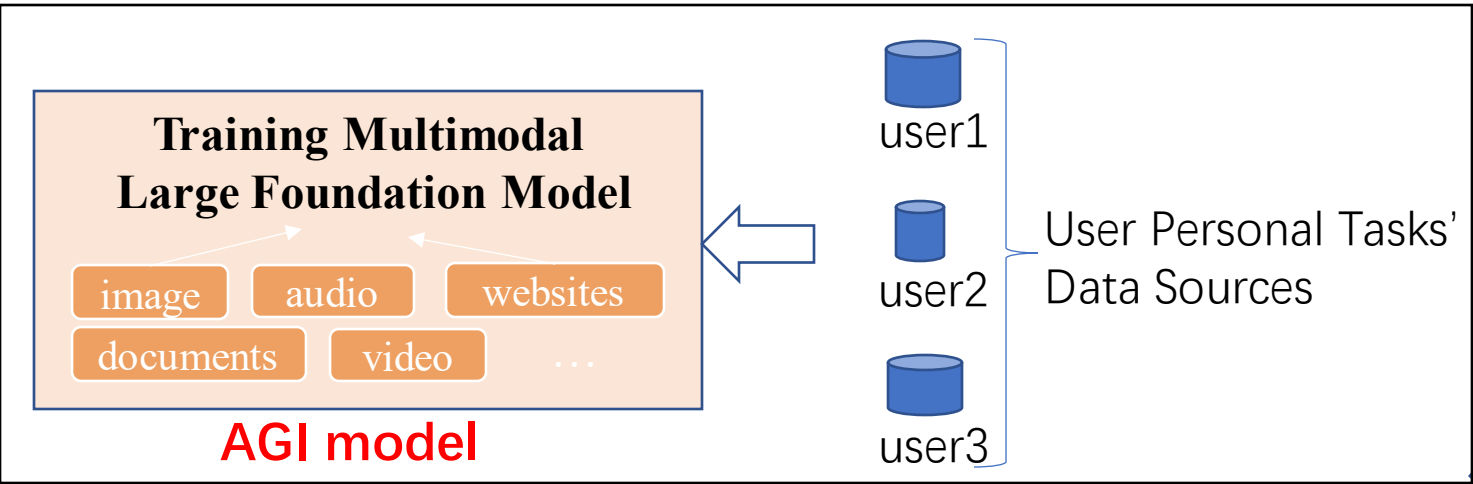
## ➤ Problem?



**Privacy Leakage**  
(e.g., products such as digital medical applications; VR products; document understanding)

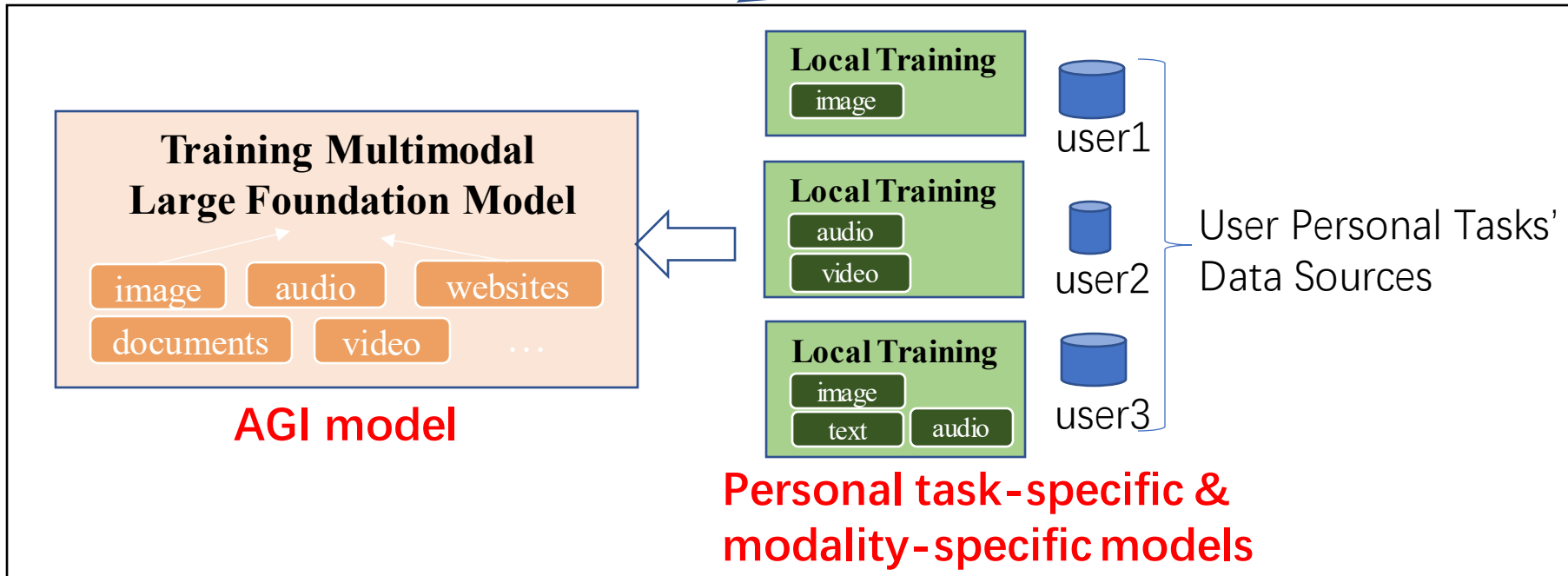


# Privacy-preserving AGI via User Collaboration



**Can we use Personalized Federated Learning (PFL) [1] to achieve**

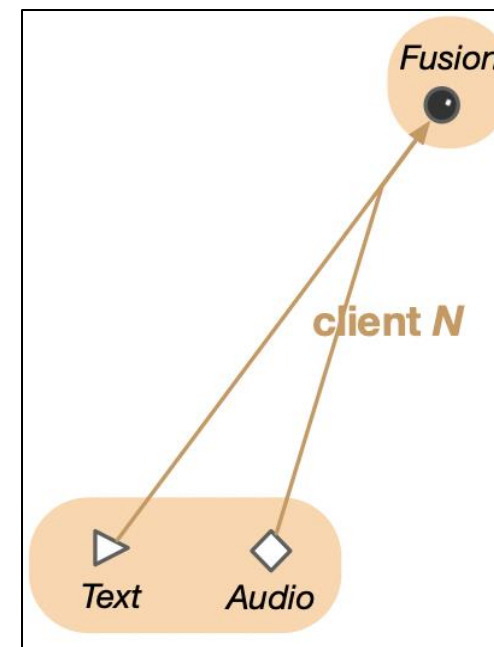
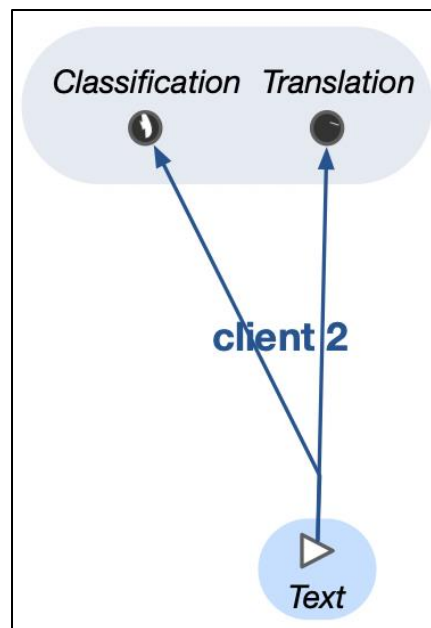
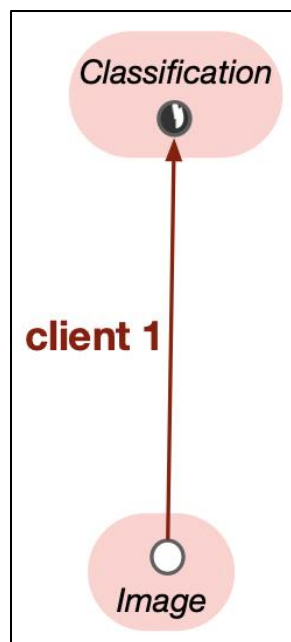
- user privacy of AGI while handling
- diverse modalities and task categories?



[1] Tan, Alysia Ziyang, et al. "Towards personalized federated learning." *IEEE Transactions on Neural Networks and Learning Systems* (2022).

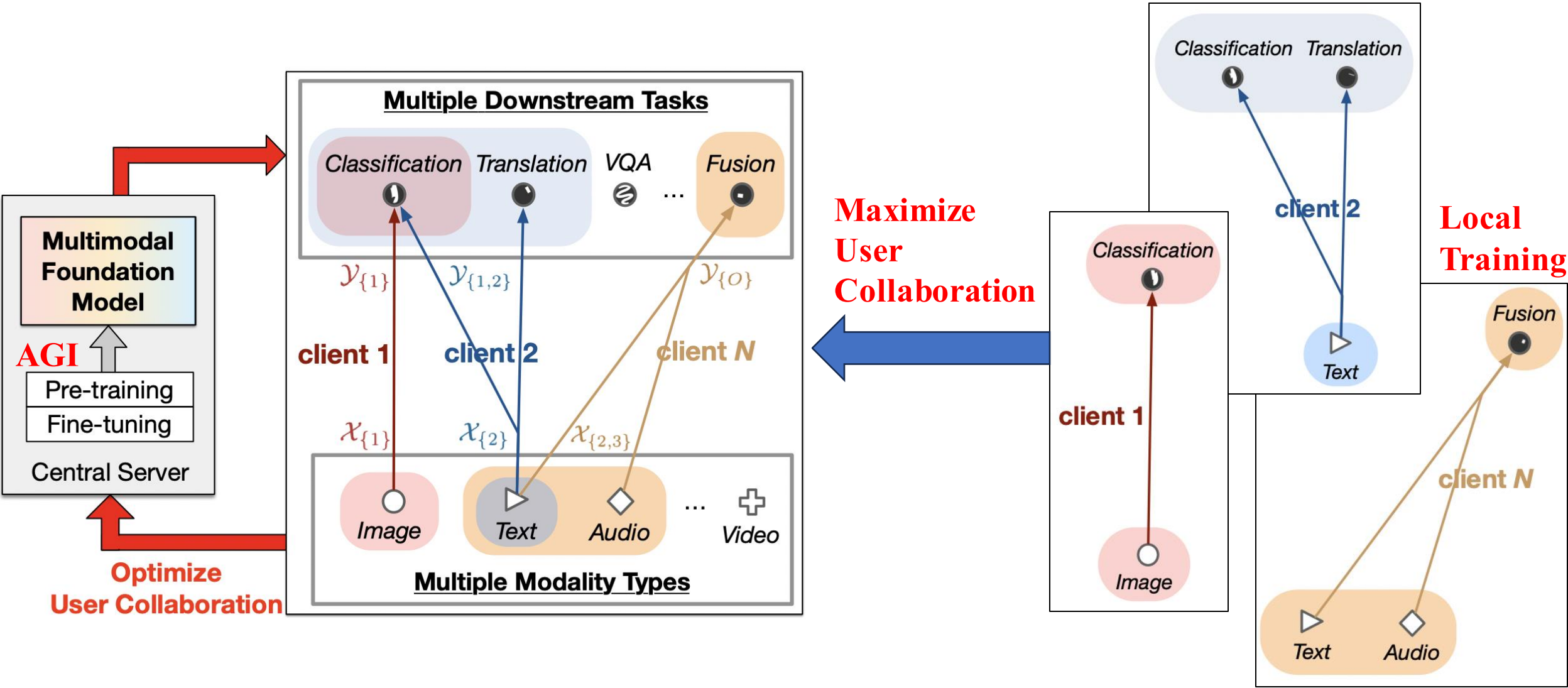
# Problem Setting

- **Local Training:** Totally personal, no AGI benefit



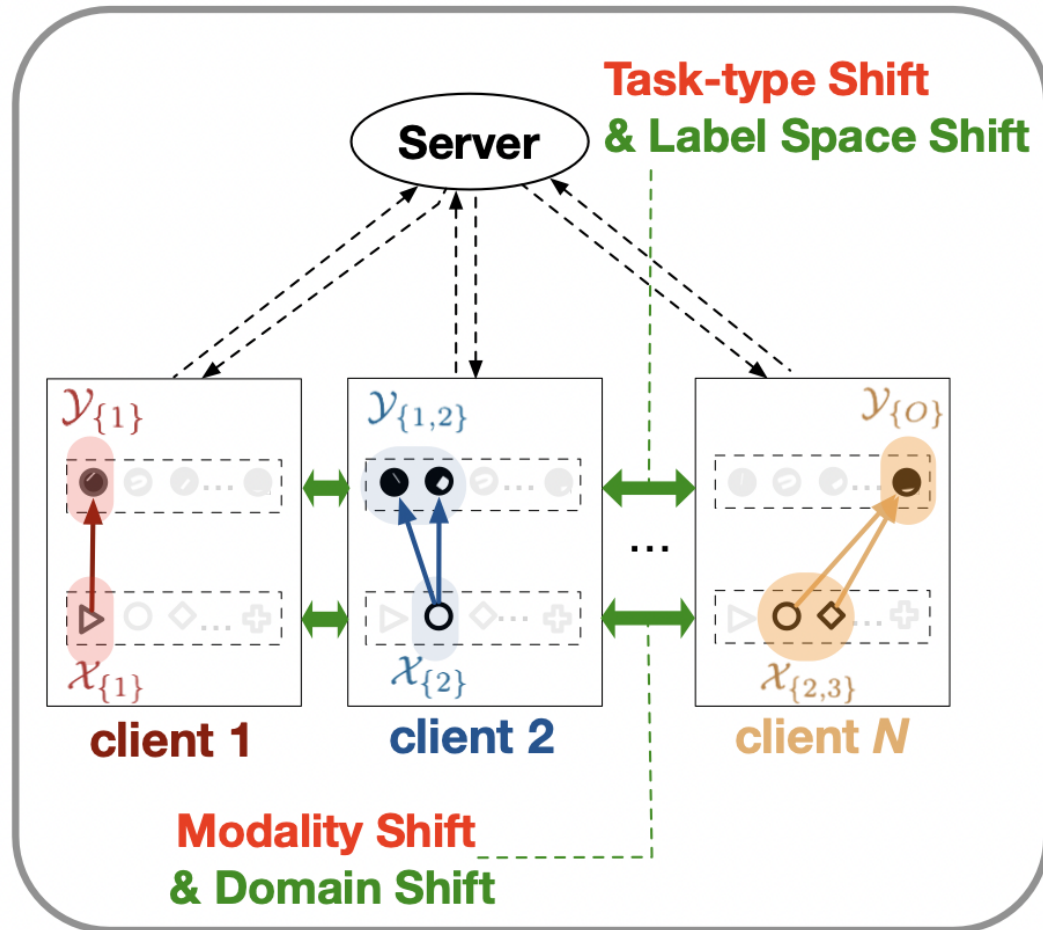
# Problem Setting

➤ **Our Setting:** leverage user collaboration to learn an AGI model





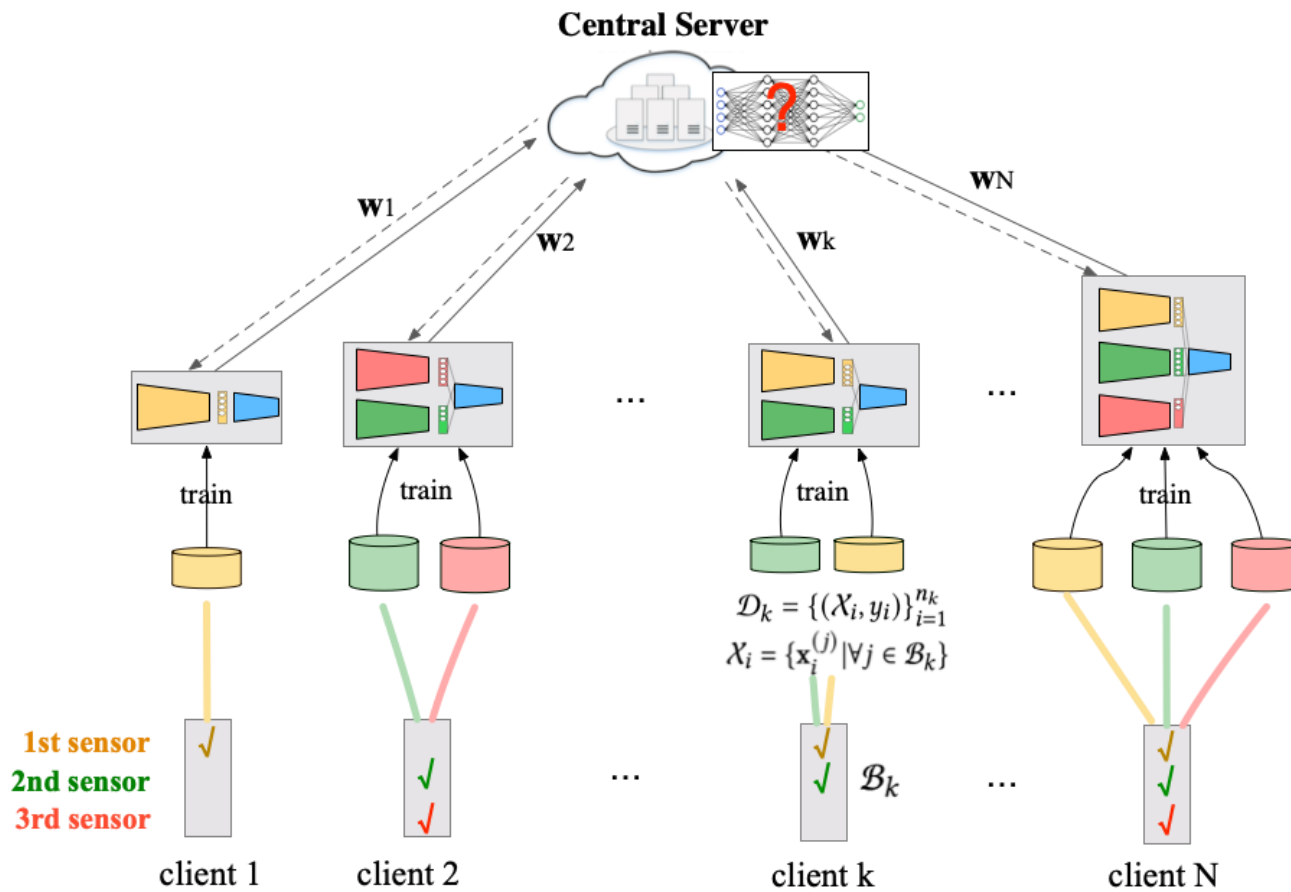
# Problem Setting



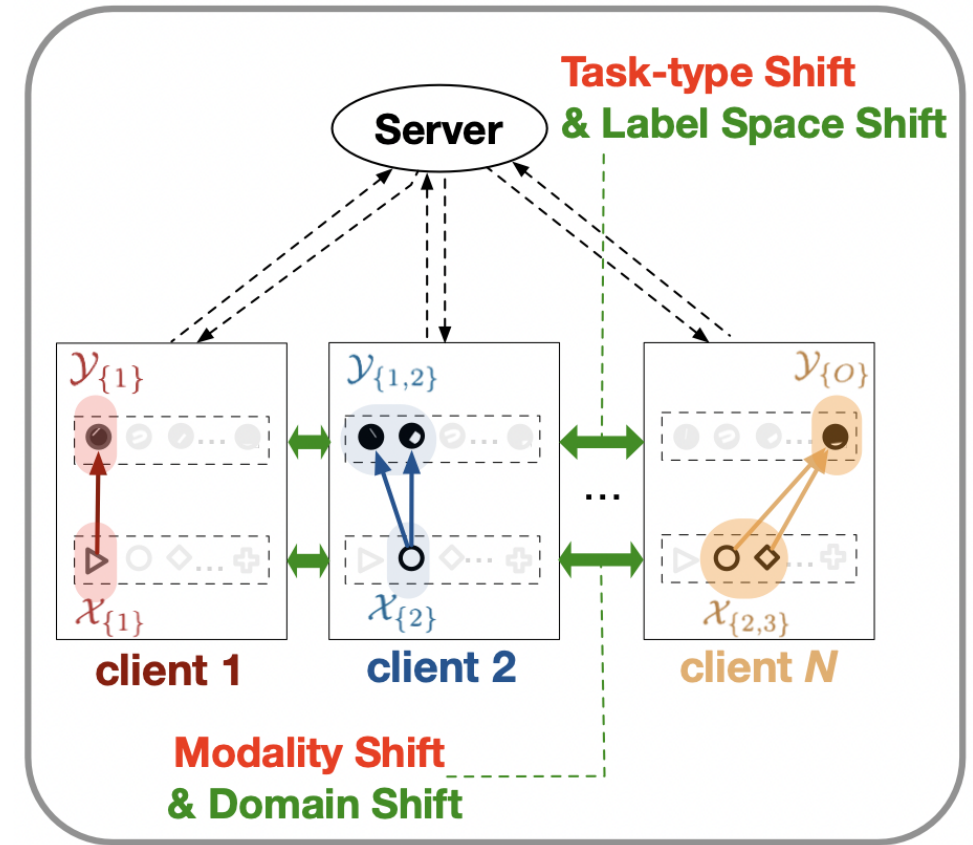
## 4 Heterogeneity Patterns:

- **Domain shift**
- **Concept shift**
- **Modality gap** (image, text, audio, video) across different domains
- **Task type difference** (object classification, image captioning, audio generation, emotion recognition)

# Unique Challenges compared to Existing Solutions



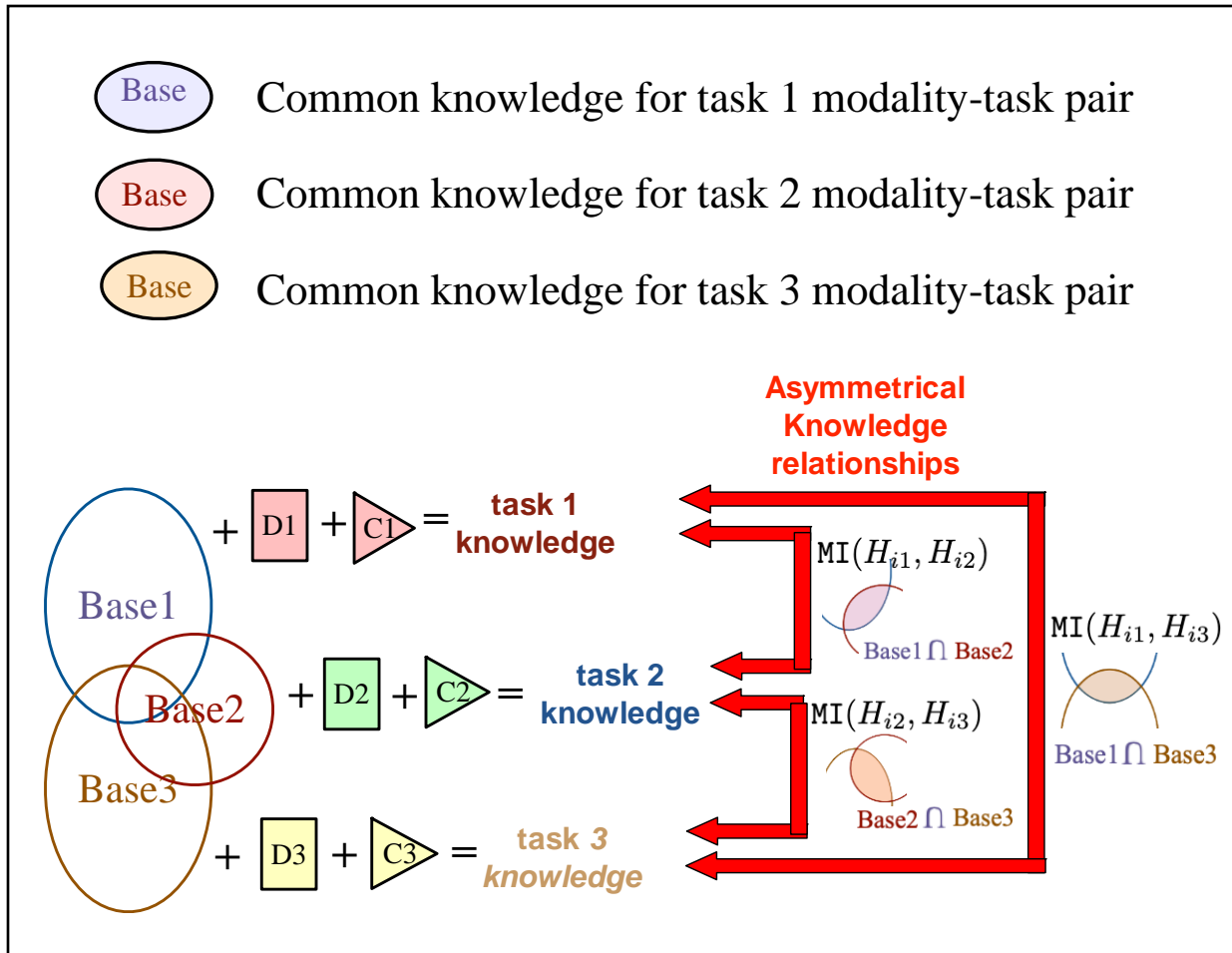
**Existing Solutions: Multimodal FL via Latent Space Alignment**



**Challenge: Suboptimal solution with large modality gap & task gap**

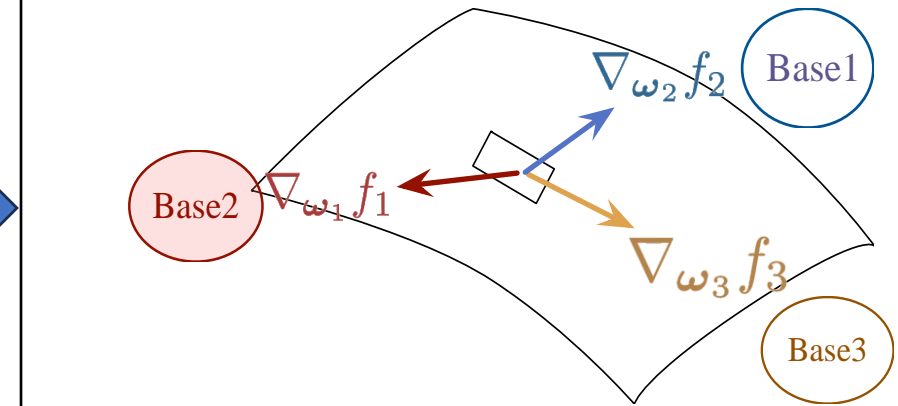
# A Closer Look: Knowledge Unalignment between Users

## True Knowledge Alignment & Conflict

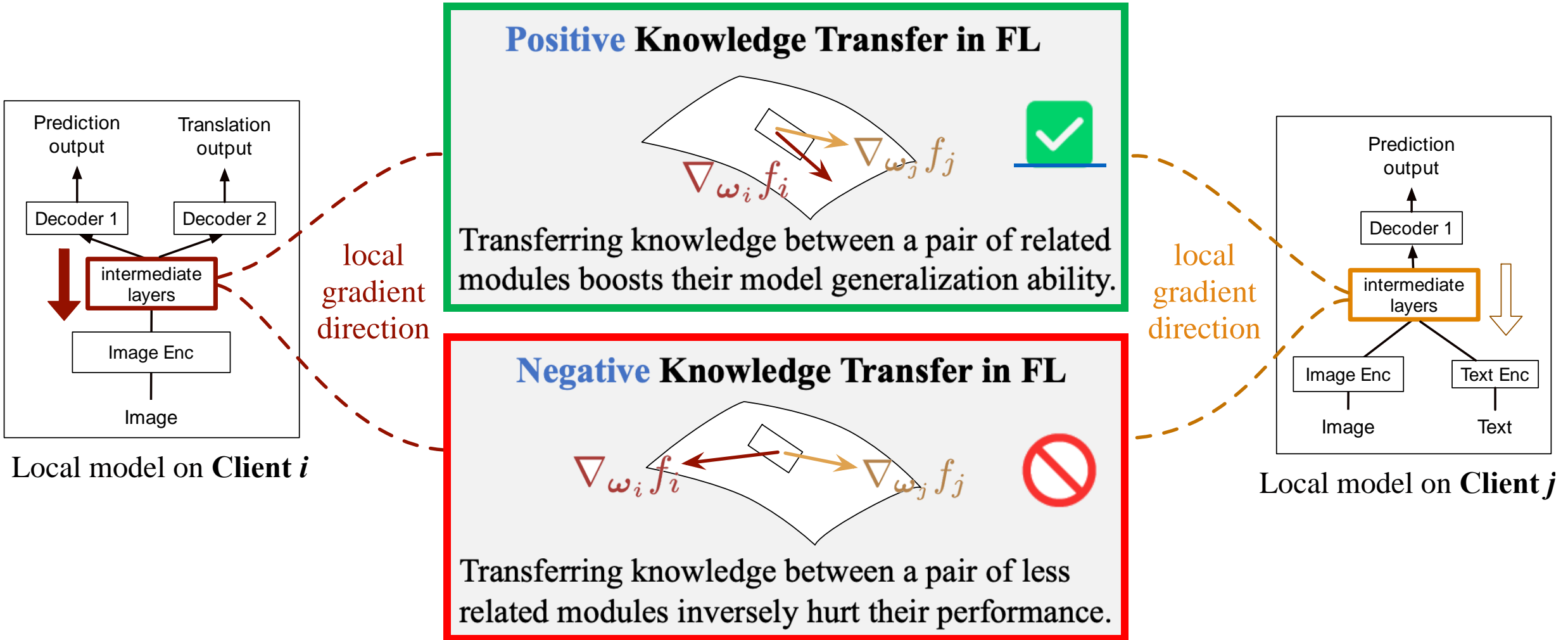


## Our Observation:

If assuming an aligned latent space, gradients at the latent layers

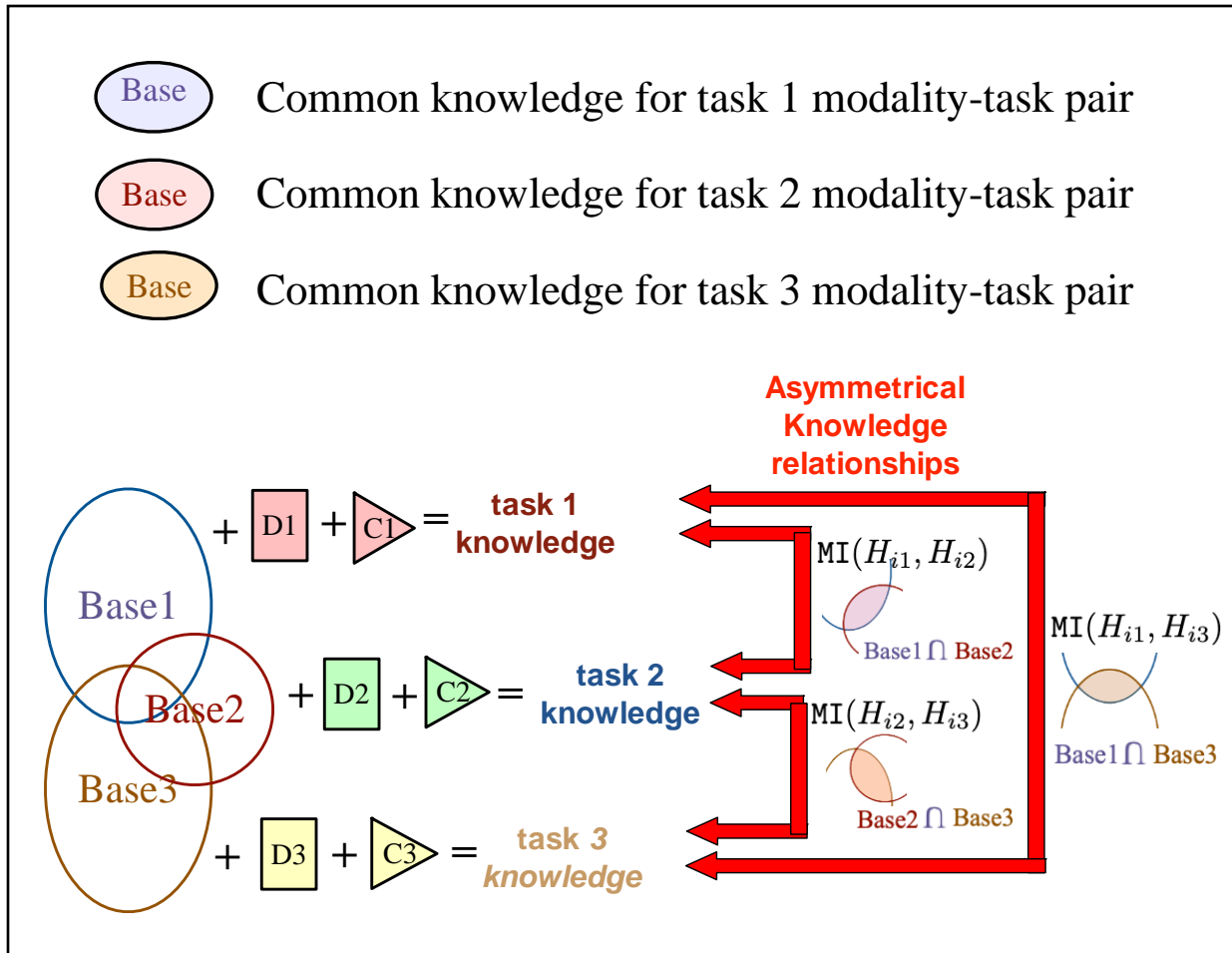


# Inspiration



# A Closer Look: Knowledge Unalignment between Users

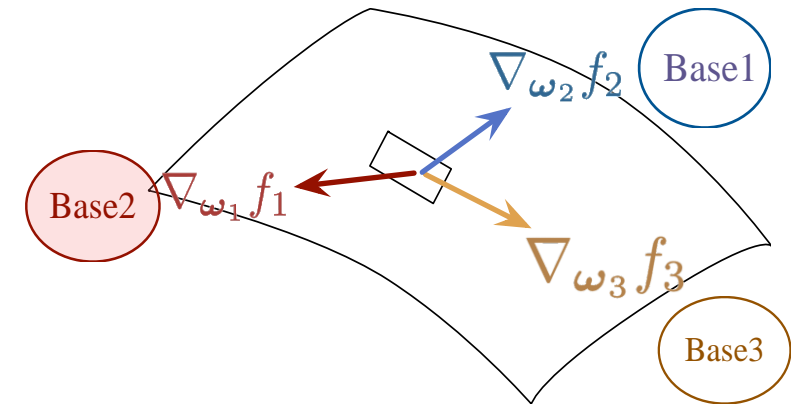
## True Knowledge Alignment & Conflict



## Our Observation:

**If assuming an aligned latent space,**

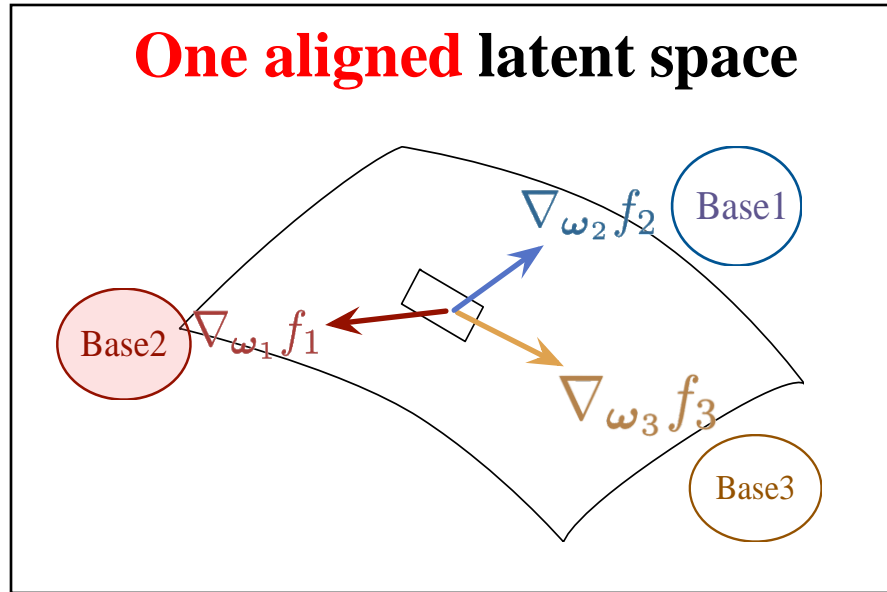
gradients at the latent layers



**Consequence:**

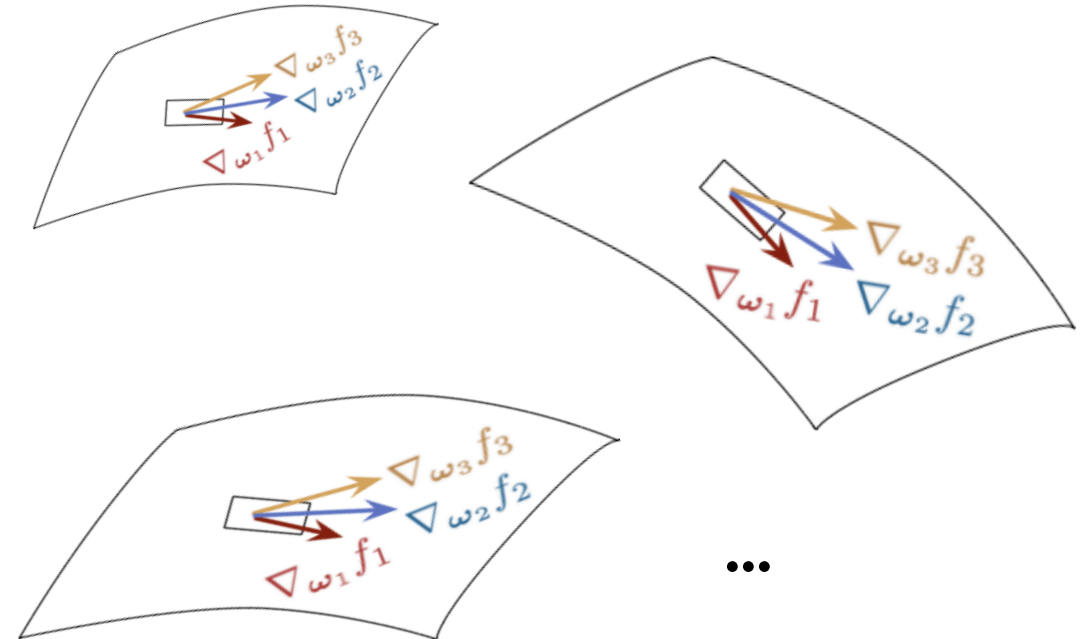
- **Insufficient Positive Knowledge Transfer**
- **Unwanted Negative Knowledge Transfer**

# Main Idea



Disentangle

**Multiple disentangled latent subspaces**



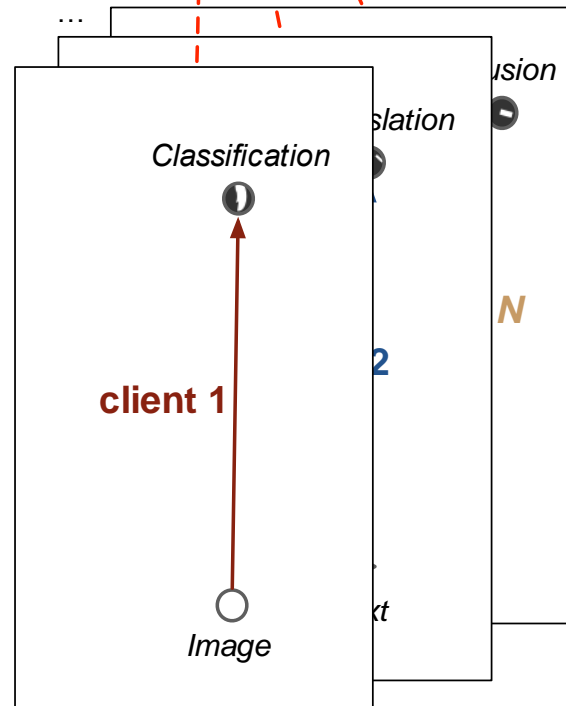
- **Insufficient** Positive Knowledge Transfer
- **Unwanted** Negative Knowledge Transfer

- **maximized** Positive Knowledge Transfer
- **minimized** Negative Knowledge Transfer

# Global Objective

$$\min_{\omega_1, \omega_2, \dots, \omega_N} \left[ \frac{1}{N} \sum_{i=1}^N f_i(\omega_i) \right] + \mathcal{R}(\omega_1, \omega_2, \dots, \omega_N):$$

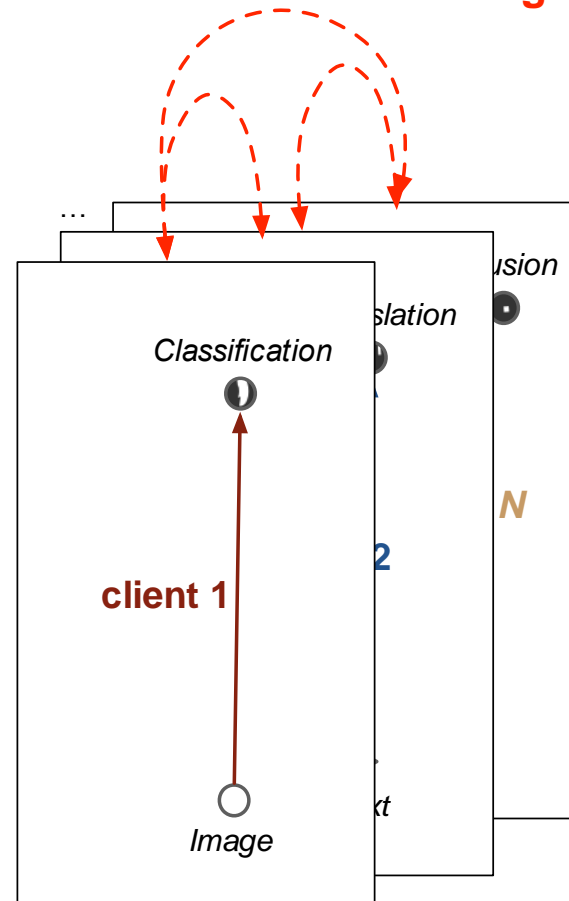
**local objectives:**  $\min_{\omega_i} f_i(\omega_i)$



# Global Objective

$$\min_{\omega_1, \omega_2, \dots, \omega_N} \left[ \frac{1}{N} \sum_{i=1}^N f_i(\omega_i) \right] + \mathcal{R}(\omega_1, \omega_2, \dots, \omega_N)$$

**privacy preserving  
knowledge sharing scheme**





# Global Objective

$$\min_{\omega_1, \omega_2, \dots, \omega_N} \left[ \frac{1}{N} \sum_{i=1}^N f_i(\omega_i) \right] + \mathcal{R}(\omega_1, \omega_2, \dots, \omega_N)$$

**original  
knowledge sharing  
scheme**



$$\sum_{k=1}^K \mathcal{R}_k(\{\omega_i^{(k)} \mid \forall i \in C_k\})$$

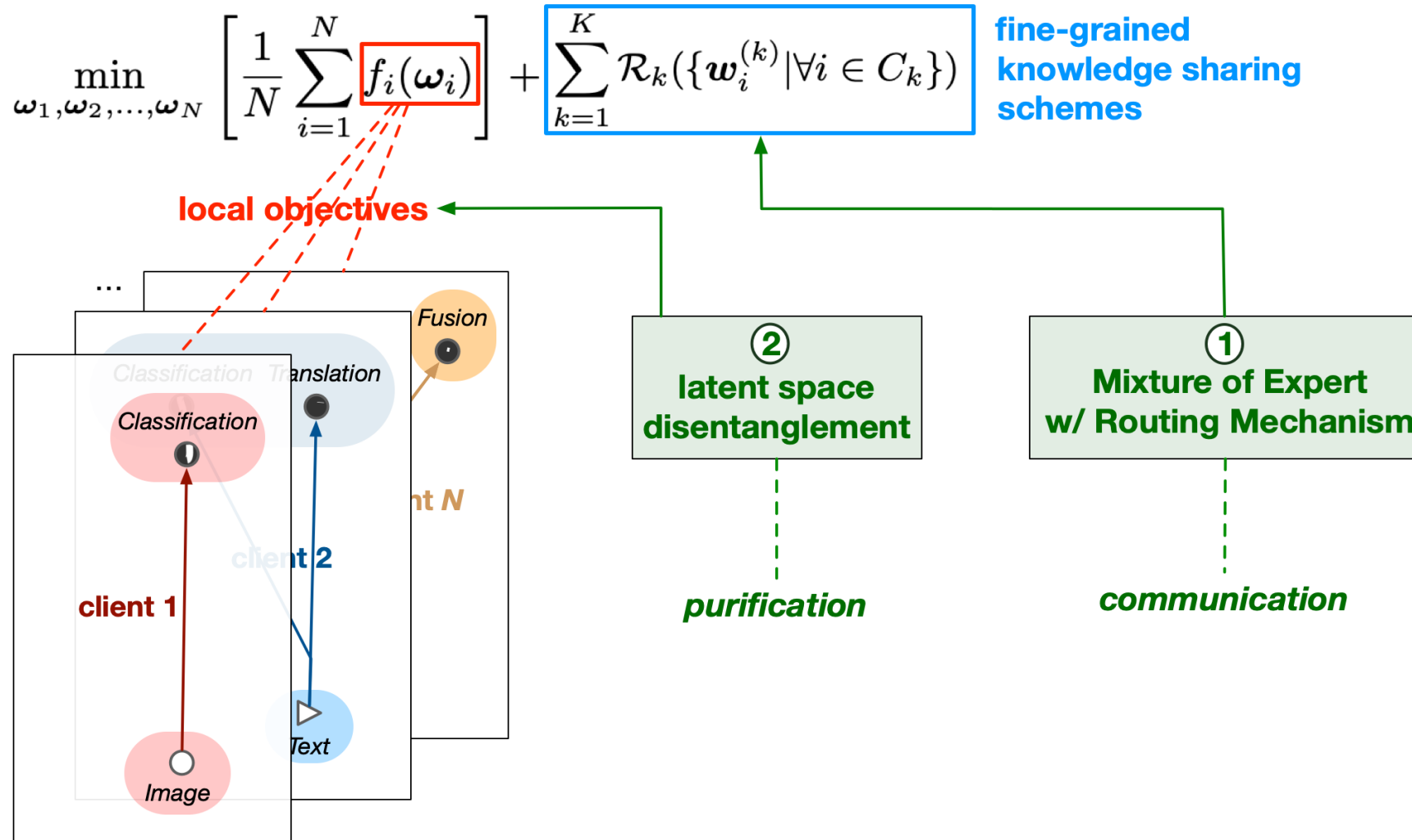
**fine-grained  
knowledge sharing  
schemes**

$$= \sum_{m \in [M]} \mathcal{R}_{\text{BE}}^m(\phi_{\text{BE},i}^{(m)} \mid i \in C_{\text{modal}}^m) + \sum_{o \in [O]} \mathcal{R}_{\text{BD}}^o(\theta_{\text{BD},i}^{(o)} \mid i \in C_{\text{task}}^o) + \sum_{m \in [M]} \sum_{d=1}^D \mathcal{R}_{\text{MoDE}}^{m,d}(\phi_{\text{mode},i}^{(m)d} \mid i \in C_{\text{MoDE}}^d) + \sum_{o \in [O]} \sum_{i=1}^N \mathcal{R}_{\text{concept}}^{o,i}(\theta_{\text{final},i}^{(o)} \mid \{i\}) + \mathcal{R}_{\text{MoTE \& MoTE}}^{\text{share,share}}(\{\phi_{\text{mote},i}^{(\text{share})\text{share}}, \theta_{\text{mome},i}^{(\text{share})\text{share}}\} \mid i \in [N]) + \sum_{m \in [M]} \left[ \mathcal{R}_{\text{MoTE \& MoTE}}^{m,\text{share}}(\{\phi_{\text{mote},i}^{(m)\text{share}}, \theta_{\text{mote},i}^{(\text{share})m}\} \mid i \in C_{\text{modal}}^m) \right] + \sum_{m \in [M]} \sum_{o \in [O]} \left[ \mathcal{R}_{\text{MoTE \& MoTE}}^{m,o}(\{\phi_{\text{mote},i}^{(m)o}, \theta_{\text{mome},i}^{(o)m}\} \mid i \in C_{\text{pair}}^{m,o}) \right] + \sum_{o \in [O]} \left[ \mathcal{R}_{\text{MoTE \& MoTE}}^{m,\text{share}}(\{\phi_{\text{mote},i}^{(\text{share})o}, \theta_{\text{mome},i}^{(o)\text{share}}\} \mid i \in C_{\text{task}}^o) \right]$$

**Disentangle Knowledge  
Sharing Schemes (Details)**

# How to Solve Global Objective?

*Research Question:* How to maximize positive transfer while minimizing negative transfer?

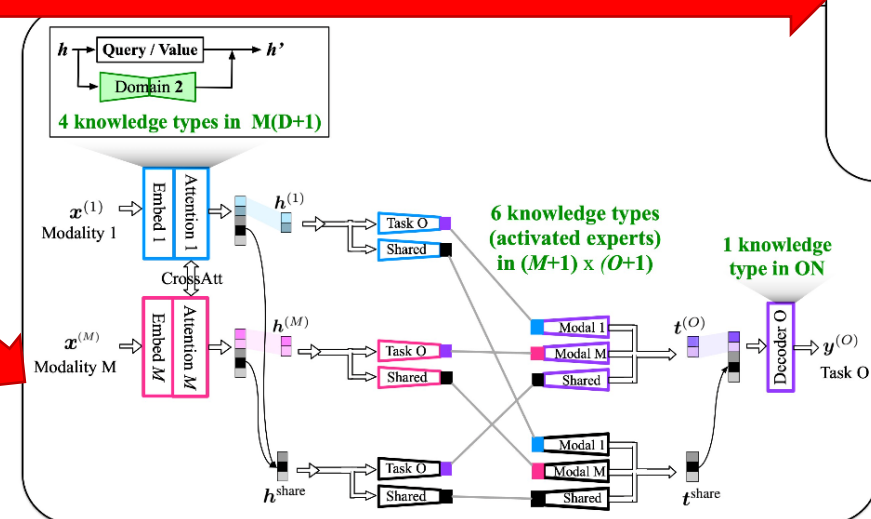
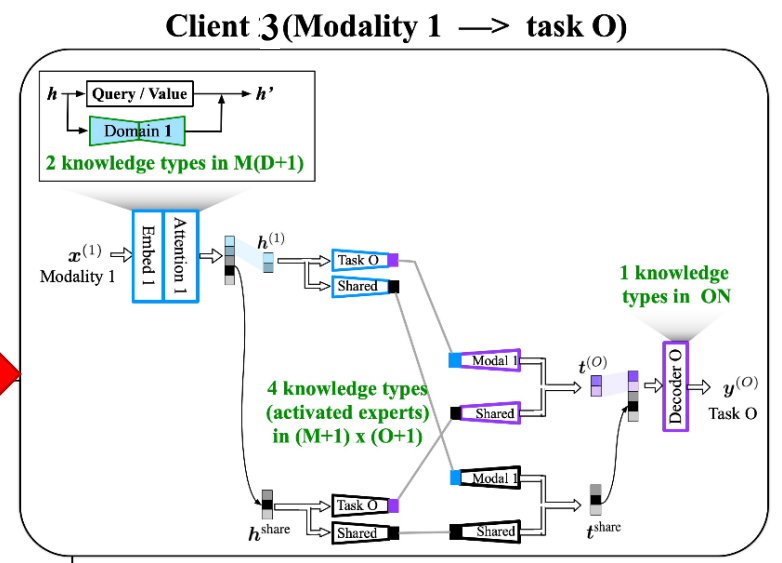
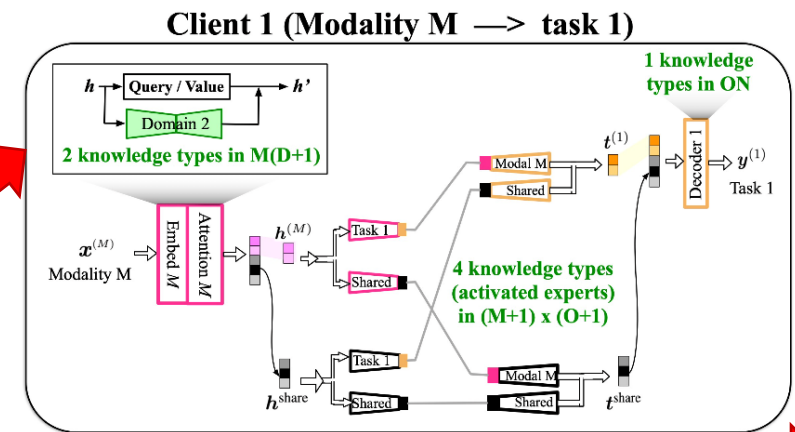
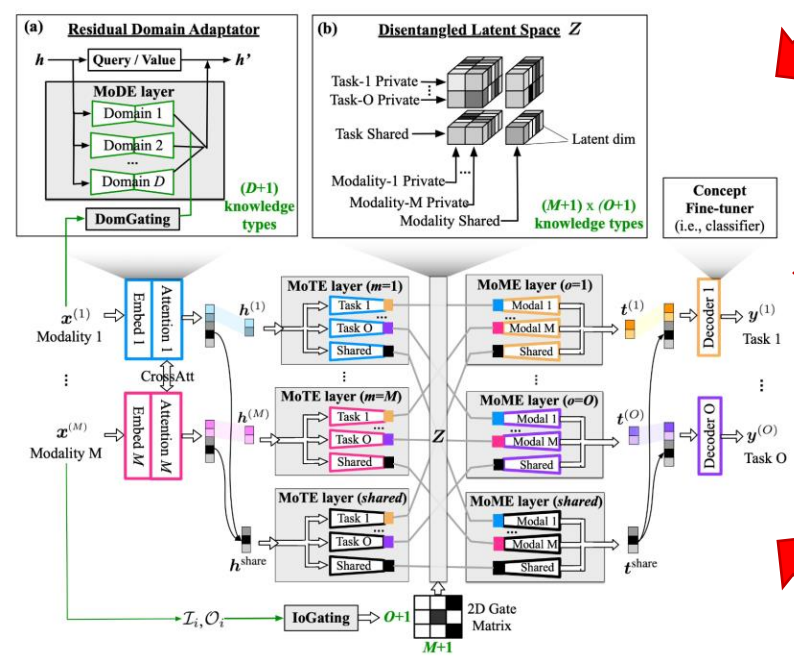


# Communication

- Model architectures based on Mixture of Experts

- Client Personal Models (MoE)

## Supernetwork (MoE) on Server

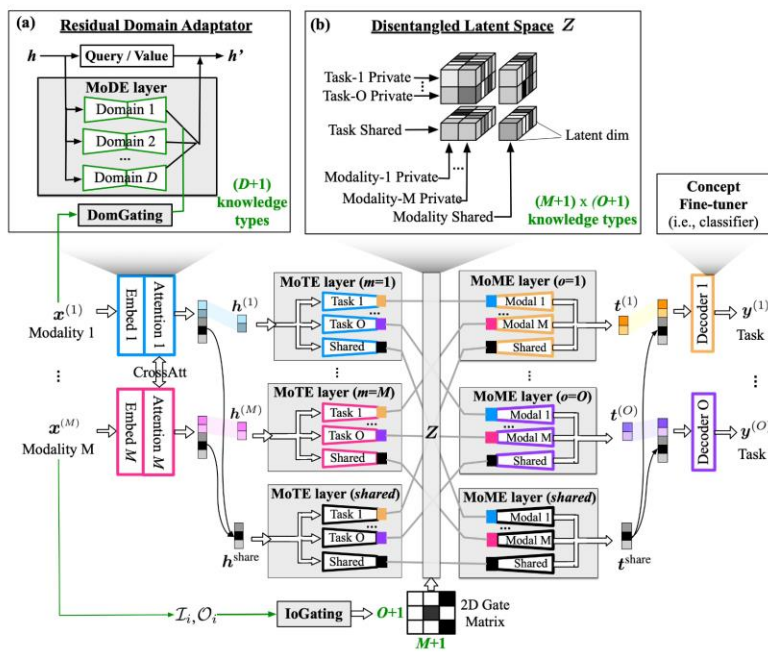


Client 2 (Modality 1 & M  $\rightarrow$  task O)

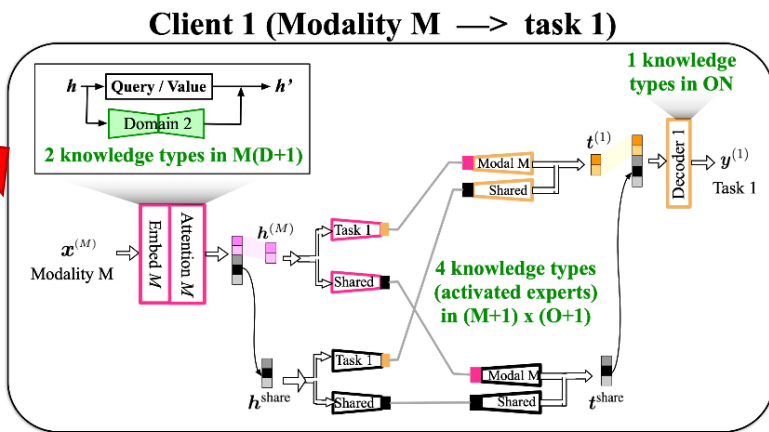
# Communication

- Automatic routing during communication

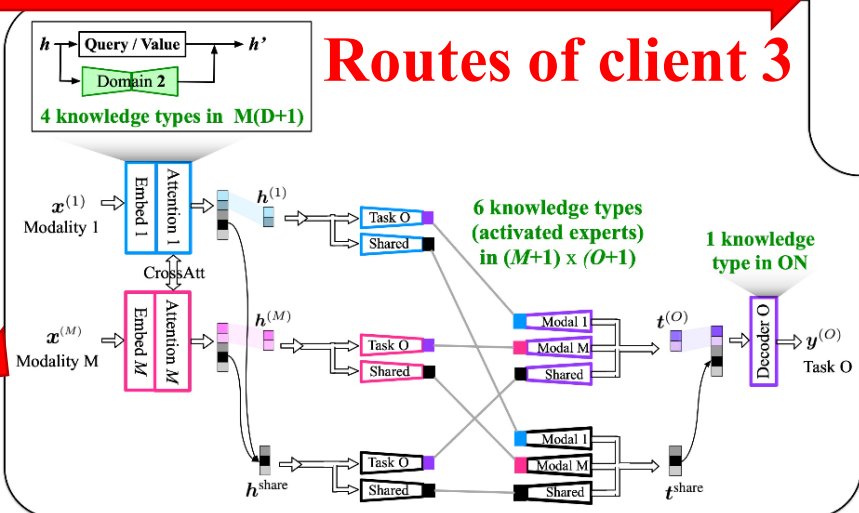
**Routes of client 1**



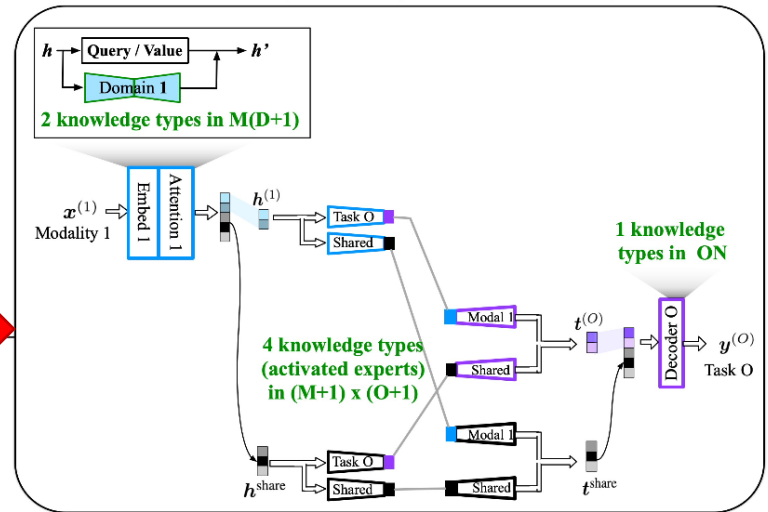
**Routes of client 2**



**Routes of client 3**



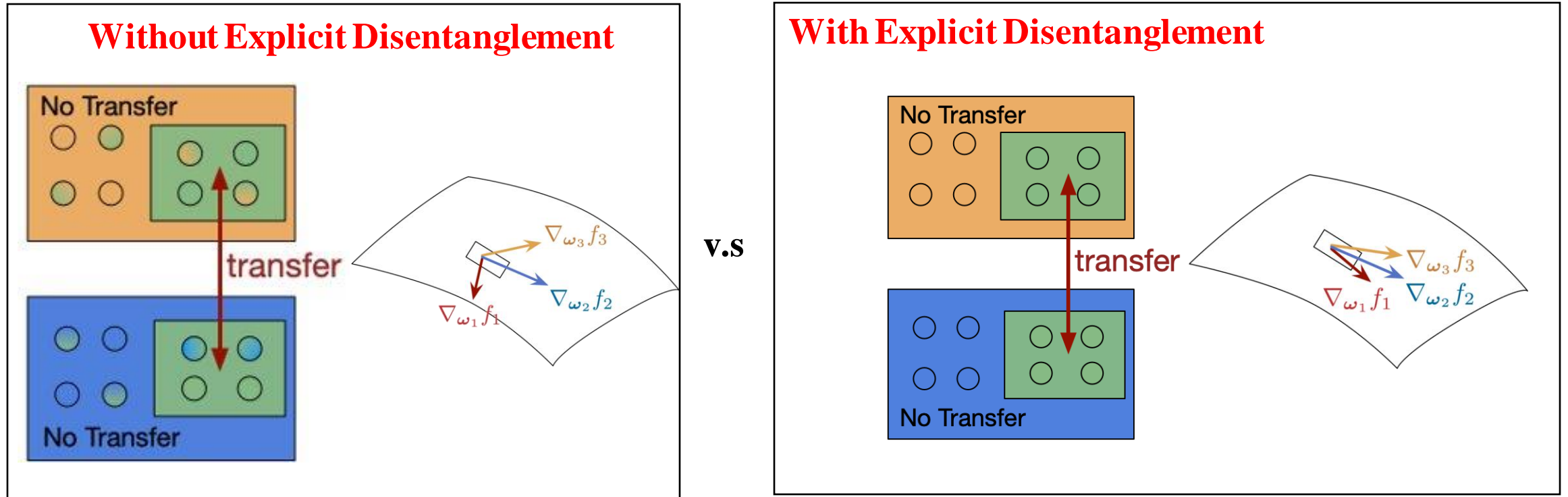
**Client 2 (Modality 1 -> task O)**



**Client 2 (Modality 1 & M -> task O)**

# Purification

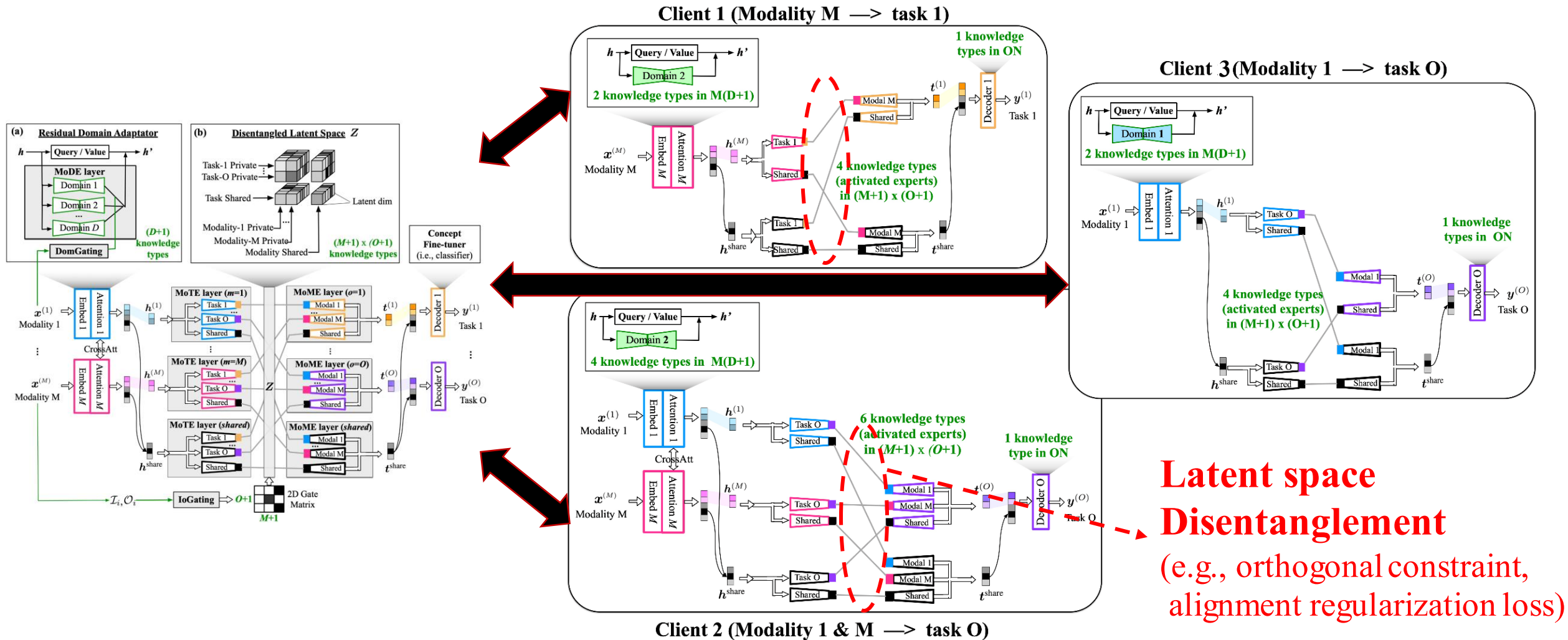
- Multiple disentangled latent subspaces



*A more purified knowledge split is beneficial to produce more aligned gradients.*

# Purification

- Leverage explicit disentanglement losses to enhance purification



# Datasets/Simulations

- #clients (<50), #modalities (<5), #downstream tasks (<5)
- model size: 4 self/cross-attention layers, 3 heads

Dataset	# Samples	Modalities	Tasks
Aircraft	10,200	{Image}	{ <i>Classification</i> (102 aircraft classes)}
CIFAR-100	60,000	{Image}	{ <i>Classification</i> (100 object classes)}
Vehicle Sensor	23,000	{Audio, Seismic}	{ <i>Classification</i> (2 vehicle types)}
ModelNet40	12,300	{View1, View2}	{ <i>Classification</i> (40 3d objects)}
CMU-MOSEI	22,777	{Audio, Text, Video}	{ <i>Classification</i> (9 sentiments), <i>Regression</i> (3 emotions)}
Multi-FMNIST	70,000	{Image}	{ <i>Classification Task 1</i> (10 digits), <i>Classification Task 2</i> (10 objects)}
AV-MNIST	70,000	{Image, Acoustic}	{ <i>Generation Task1</i> (image), <i>Generation Task2</i> (audio), <i>Classification</i> (10 digits)}

“Multimedia Understanding”

$3\rho \times (\mathcal{X}_{\text{video}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{text}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{video, audio}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{audio, text}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{video, text}} \rightarrow \mathcal{Y}_{\text{sentiment, emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{text}} \rightarrow \mathcal{Y}_{\text{emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{video, audio}} \rightarrow \mathcal{Y}_{\text{emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{audio, text}} \rightarrow \mathcal{Y}_{\text{sentiment, emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{cls\_vehicle}})$   
 $3\rho \times (\mathcal{X}_{\text{seismic}} \rightarrow \mathcal{Y}_{\text{cls\_vehicle}})$   
 $10(1 - \rho) \times (\mathcal{X}_{\text{audio, seismic}} \rightarrow \mathcal{Y}_{\text{cls\_vehicle}})$   
 $10(1 - \rho) \times (\mathcal{X}_{\text{video, text, audio}} \rightarrow \mathcal{Y}_{\text{emotions}})$

Client inputs    Client outputs

Cross-modal Generation & Understanding

$3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_digits}})$   
 $3 \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{cls\_digits}})$   
 $3 \times (\mathcal{X}_{\text{image, audio}} \rightarrow \mathcal{Y}_{\text{cls\_digits}})$   
 $3 \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{cls\_digits, gen\_image}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{gen\_audio}})$   
 $3 \times (\mathcal{X}_{\text{image, audio}} \rightarrow \mathcal{Y}_{\text{gen\_image, gen\_audio}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_digits, cls\_objects}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_objects, gen\_audio}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_digits, cls\_objects, gen\_image}})$

Client inputs    Client outputs

# Some Experimental Results

- Server model size: 22M
- Average of client model sizes: 63% (Quantization: 15%)

Methods	Average Testing Accuracy on Classification Tasks	
Local	88.23 ± 0.72	70.23 ± 0.93
FedAvg	84.63 ± 0.02	74.12 ± 0.93
Multi-FedAvg	84.82 ± 0.29	69.65 ± 0.73
FedMSplit	87.37 ± 0.03	73.25 ± 0.31
<b>Ours</b>	<b>96.38 ± 0.41</b>	<b>75.96 ± 0.83</b>

$3\rho \times (\mathcal{X}_{\text{video}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{text}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{video, audio}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{audio, text}} \rightarrow \mathcal{Y}_{\text{sentiment}})$   
 $3\rho \times (\mathcal{X}_{\text{video, text}} \rightarrow \mathcal{Y}_{\text{sentiment, emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{text}} \rightarrow \mathcal{Y}_{\text{emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{video, audio}} \rightarrow \mathcal{Y}_{\text{emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{audio, text}} \rightarrow \mathcal{Y}_{\text{sentiment, emotions}})$   
 $3\rho \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{cls\_vehicle}})$   
 $3\rho \times (\mathcal{X}_{\text{seismic}} \rightarrow \mathcal{Y}_{\text{cls\_vehicle}})$   
 $10(1 - \rho) \times (\mathcal{X}_{\text{audio, seismic}} \rightarrow \mathcal{Y}_{\text{cls\_vehicle}})$   
 $10(1 - \rho) \times (\mathcal{X}_{\text{video, text, audio}} \rightarrow \mathcal{Y}_{\text{emotions}})$

“Multimedia Understanding”

$3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_digits}})$   
 $3 \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{cls\_digits}})$   
 $3 \times (\mathcal{X}_{\text{image, audio}} \rightarrow \mathcal{Y}_{\text{cls\_digits}})$   
 $3 \times (\mathcal{X}_{\text{audio}} \rightarrow \mathcal{Y}_{\text{cls\_digits, gen\_image}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{gen\_audio}})$   
 $3 \times (\mathcal{X}_{\text{image, audio}} \rightarrow \mathcal{Y}_{\text{gen\_image, gen\_audio}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_digits, cls\_objects}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_objects, gen\_audio}})$   
 $3 \times (\mathcal{X}_{\text{image}} \rightarrow \mathcal{Y}_{\text{cls\_digits, cls\_objects, gen\_image}})$

Cross-modal Generation & Understanding



**Thank You!**  
**Q & A**