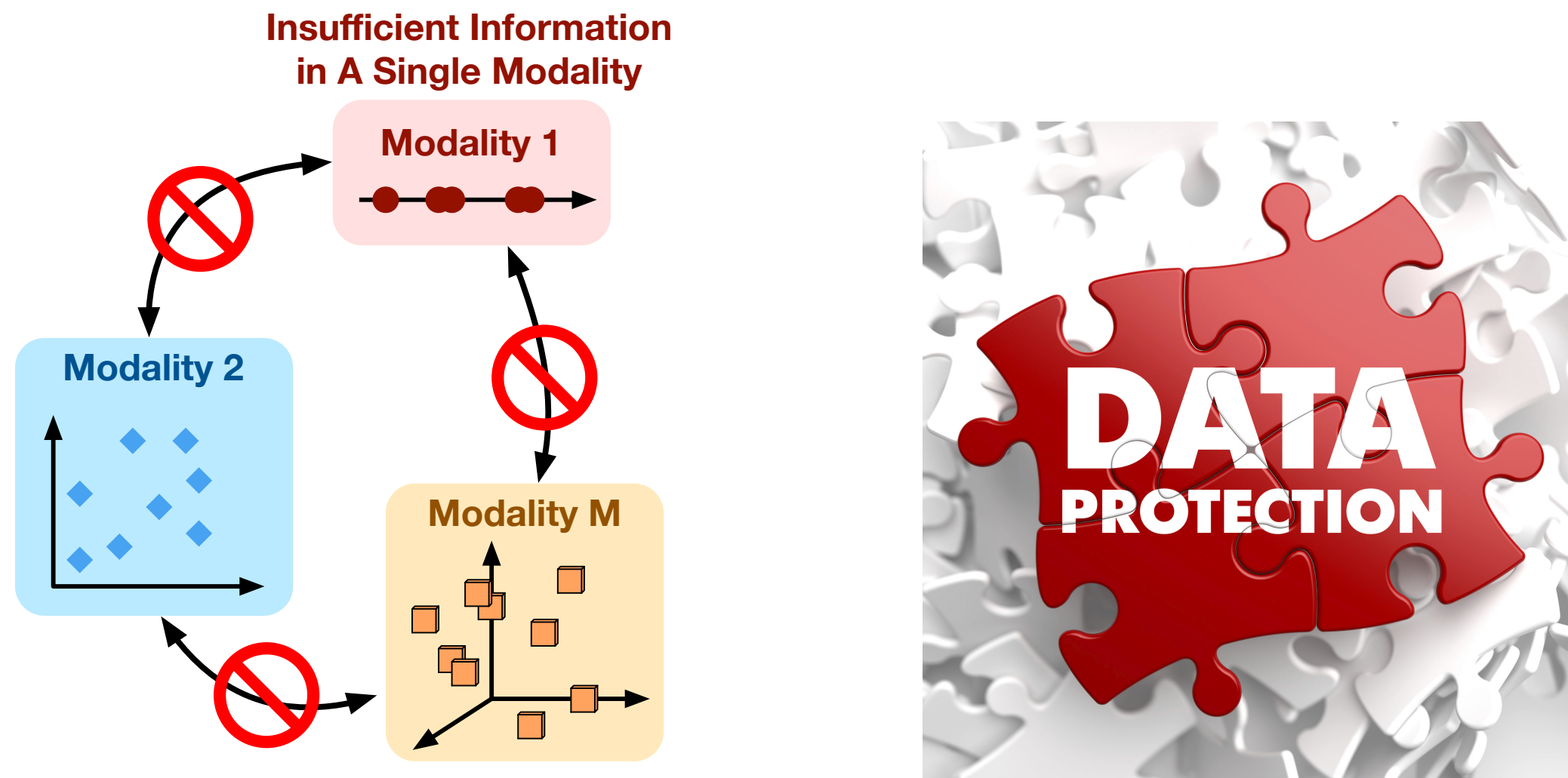


# AN OPTIMAL TRANSPORT-BASED LATENT MIXER FOR ROBUST MULTI-MODAL LEARNING

Fengjiao Gong, Angxiao Yue, Hongteng Xu  
Renmin University of China, Beijing, China

## Motivation

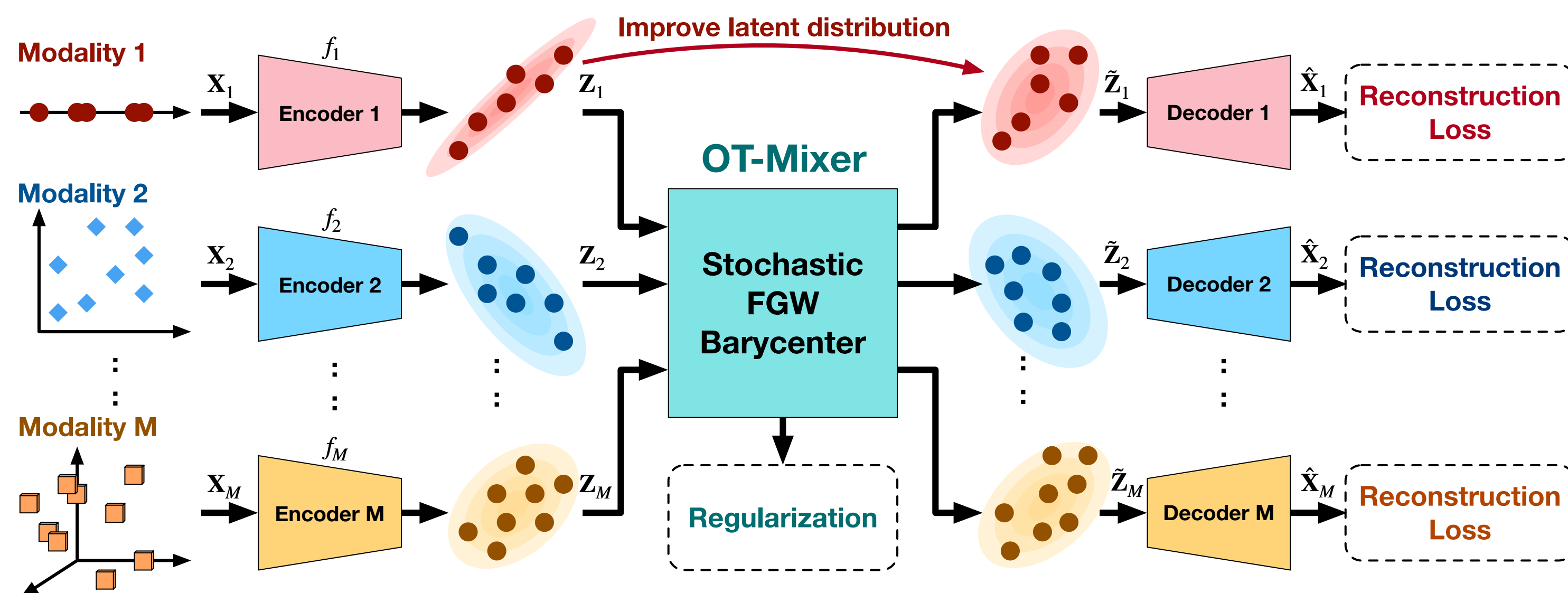


Practical scenario with unaligned and distributed multi-modal data.

- Real-world multi-modal data are often scattered to different local agents, and each agent can only access the data in a single modality.
- Due to privacy protection and data security, sharing data directly across different agents is forbidden in many applications.
- What is worse, for some agents, the data associated with its modality may be insufficient for downstream tasks because the number of the data can be limited and the features can be not informative enough for representation learning.

## Proposed Method

Suppose that we have a set of multi-modal data, denoted as  $\mathcal{D} = \{X_m\}_{m=1}^M$ , where  $M$  is the number of modalities. The data of the  $m$ -th modality, i.e.,  $X_m = \{x_{m,j}\}_{j=1}^{N_m} \in \mathbb{R}^{N_m \times D_m}$ , contains  $N_m$   $D_m$ -dimensional samples.



The scheme of OTM-based multi-modal learning.

- Learning the WAE models

$$\min_{\{f_m, g_m\}_{m=1}^M} \sum_{m=1}^M \|X_m - \hat{X}_m\|_F^2 + \lambda R(\{Z_m\}_{m=1}^M), \quad (1)$$

- Reconstruction loss of each modality

$$\mathcal{L}_{mix} = \sum_{m=1}^M \left( \|X_m - g_m(Z_m)\|_F^2 + \|Z_m - \tilde{Z}_m\|_F^2 \right), \quad (2)$$

- Regularization on the latent representations

$$\min_{T \in \Pi(\mu_B, \mu_C)} \underbrace{\sum_{i,j,k,l} |K(i,k) - I_C(j,l)|^2 t_{ij} t_{kl}}_{GW(\tilde{Z}_B, I_C)}, \quad \text{clustering} \quad (3)$$

$$\mathcal{L}_{supervise}(h_m(\tilde{Z}_m), y_m) + \mathcal{L}_{supervise}(h_m(Z_m), y_m), \quad \text{supervised tasks}$$

## AN OPTIMAL TRANSPORT-BASED LATENT MIXER

- FGW Distance  $FGW(Z, Z_m; \alpha)$

$$\min_{T_m \in \Pi(\mu, \mu_m)} \sum_{i,j,k,l} \alpha \underbrace{d_Z^2(z_i, z_{m,j}) t_{m,ij}}_{\text{Wasserstein term}} + (1-\alpha) \underbrace{|A(i,k) - A_m(j,l)|^2 t_{m,ij} t_{m,kl}}_{\text{Gromov-Wasserstein term}}, \quad (4)$$

- FGW Barycenter Problem

$$Z_B, \{T_m^*\}_{m=1}^M = \arg \min_Z \sum_{m=1}^M FGW(Z, Z_m; \alpha), \quad (5)$$

- Stochastic Mixing

$$\tilde{Z}_B = N_B \sum_{i=1}^M M_m \odot (T_m^* Z_m), \text{ with } \sum_{m=1}^M M_m = 1_{N_B \times d}, \quad (6)$$

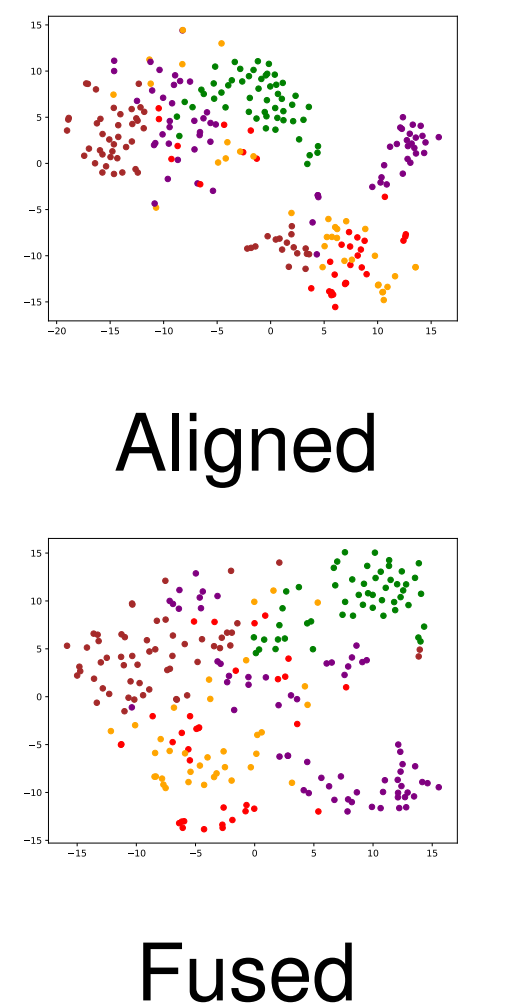
- Augmented Latent Code

$$\tilde{Z}_m = (T_m^*)^\top \tilde{Z}_B, \quad (7)$$

## Numerical Comparisons

- Clustering Performance

Data type	Datasets Algorithms	Caltech 7 Purity	ORL Purity	Movies Purity	Prokaryotic Purity
Aligned	MCCA	0.5313	0.3475	0.0989	0.5620
	DCCAE	0.4110	0.5625	0.1572	0.5070
	AttnAE	0.4600	0.4600	0.1880	0.5390
	MVKSC	0.5196	0.3013	0.2285	<b>0.6188</b>
	MultiNMF	0.4525	<b>0.6900</b>	0.1726	0.5771
	MWAE+OTM	0.6072	0.6563	0.3177	0.5952
	MWAE+OTM(WB)	<b>0.6097</b>	0.6525	<b>0.3184</b>	0.5541
Unaligned	MVC-UM	0.3112	0.5431	0.1841	0.4451
	GWMAC	0.3568	0.5118	0.1928	<b>0.5479</b>
	MWAE+OTM	<b>0.5788</b>	<b>0.6550</b>	<b>0.2925</b>	0.5438
	MWAE+OTM(WB)	0.5667	0.6350	0.2860	0.5299



- Classification & Regression Performance

Dataset	Method	Result	Task
AV-MNIST	Late fusion	0.7295	Classification
	Late fusion + OTM	<b>0.7316</b>	
ENRICO	MI matrix	0.4815	Classification
	MI matrix + OTM	0.4814	
	Tensor matrix + OTM	<b>0.4911</b>	
CMU-MOSI	Late fusion	0.5194	Classification
	Late fusion + OTM	<b>0.5368</b>	
	LRTF	0.5245	
	LRTF + OTM	<b>0.5327</b>	
	MFM	0.5391	Regression
	MFM + OTM	<b>0.5410</b>	
MUJOCO	Late fusion	1.3710	Regression
	Late fusion + OTM	<b>1.3630</b>	
	Tensor fusion	1.3691	
	Tensor fusion + OTM	<b>1.3644</b>	
	Tensor fusion	$1.583 \times 10^{-3}$	Regression
	Tensor fusion + OTM	<b><math>1.369 \times 10^{-3}</math></b>	

Dataset (Task)	Method	Selected Modality 1	2	3
Prokaryotic (Clustering)	MWAE	0.4936	<b>0.6261</b>	0.4791
	MWAE+OTM	<b>0.5554</b>	0.5209	<b>0.5426</b>
CMU-MOSI (Classification)	Late fusion	0.5369	<b>0.5373</b>	0.5163
	Late fusion+OTM	<b>0.5428</b>	0.5328	<b>0.5268</b>
	LRTF	<b>0.5190</b>	<b>0.5241</b>	0.5131
CMU-MOSI (Regression)	LRTF+OTM	0.5131	0.5222	<b>0.5209</b>
	Late fusion	<b>1.3721</b>	<b>1.3581</b>	1.4055
	Late fusion+OTM	1.3804	1.3629	<b>1.3698</b>
	Tensor fusion	1.3684	1.3680	1.3853
	Tensor fusion+OTM	<b>1.3674</b>	<b>1.3669</b>	<b>1.3716</b>

Single Modality

## Conclusion & Future Work

- We propose a novel optimal transport-based mixer (OTM) that achieves data alignment and augmentation for robust multi-modal learning.
- In the future, we plan to test our method in real-world applications, e.g., federated learning for healthcare data modeling.



Github Repo

Email:  
hongtengxu@ruc.edu.cn  
fengjiaogong2021@ruc.edu.cn  
angxiaoyue@ruc.edu.cn