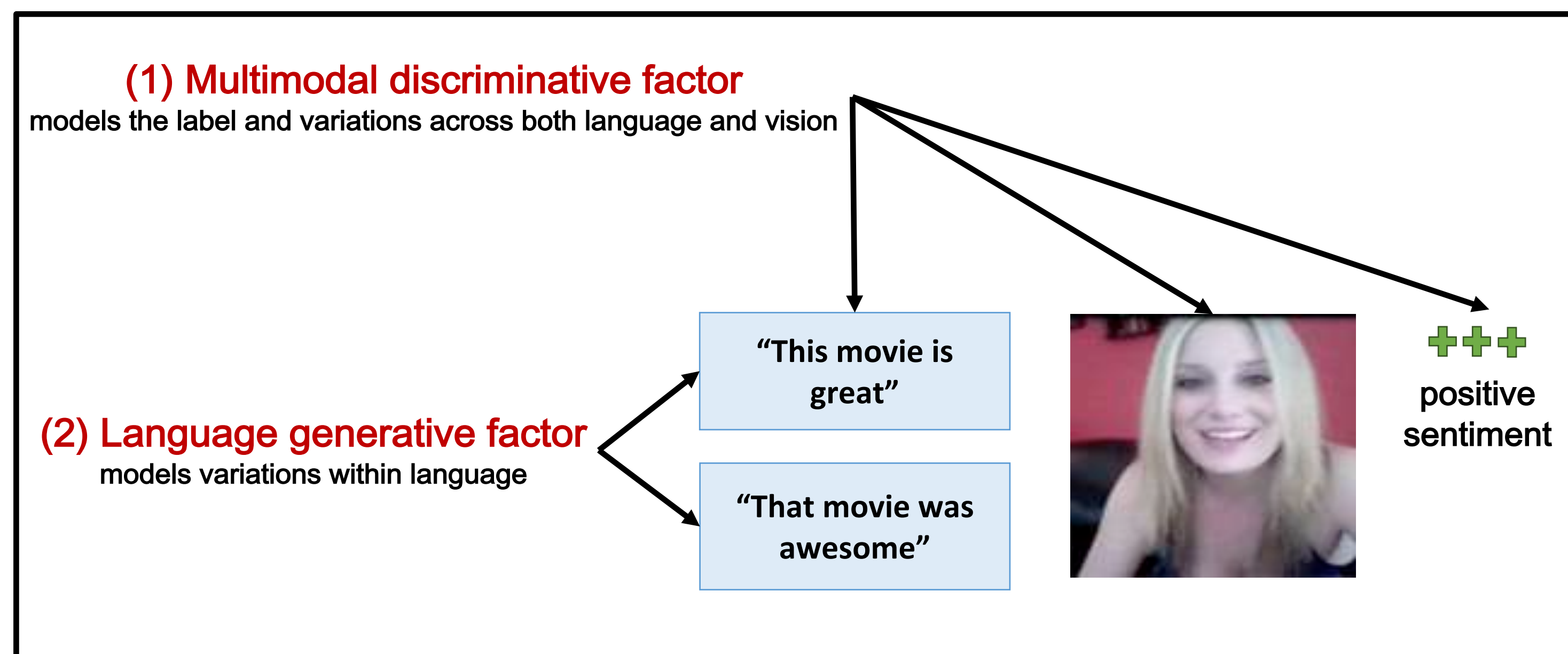


Multimodal Representation Learning

Multimedia Content Intelligent
Personal Assistants Robots and
Virtual Agents



Factorize representation into **independent sets** of factors



Experimental Setup: Datasets and Features

• Datasets in Human Multimodal Language

– Multimodal Personal Trait Recognition

* Movie Reviews (POM)

– Multimodal Sentiment Analysis

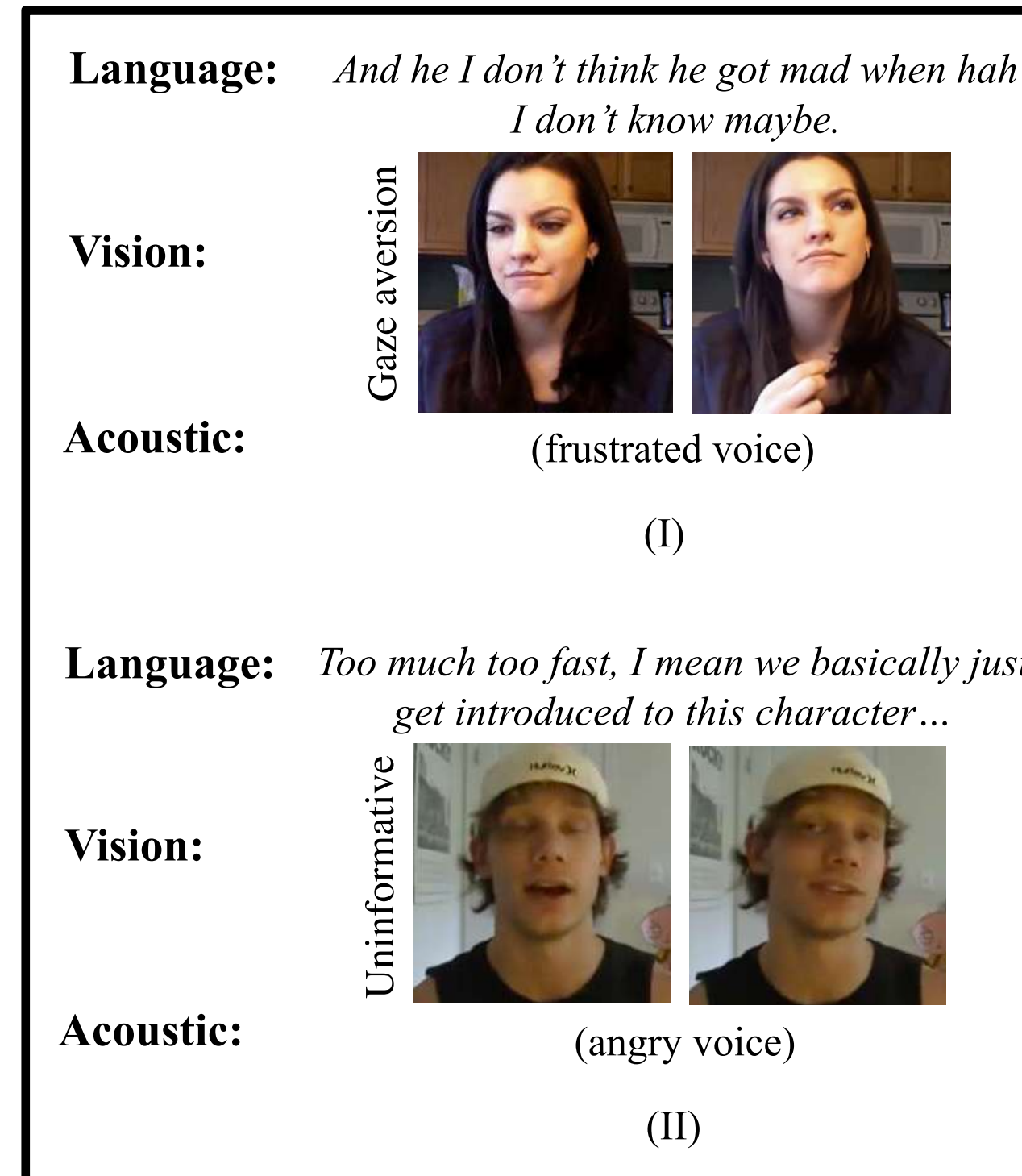
* Monologue Opinion Videos (CMU-MOSI)
* Online Social Reviews (ICT-MMMO)
* Product Reviews (MOUD and YouTube)

– Multimodal Emotion Recognition

* Recorded Dyadic Dialogues (IEMOCAP)

• Multimodal Features

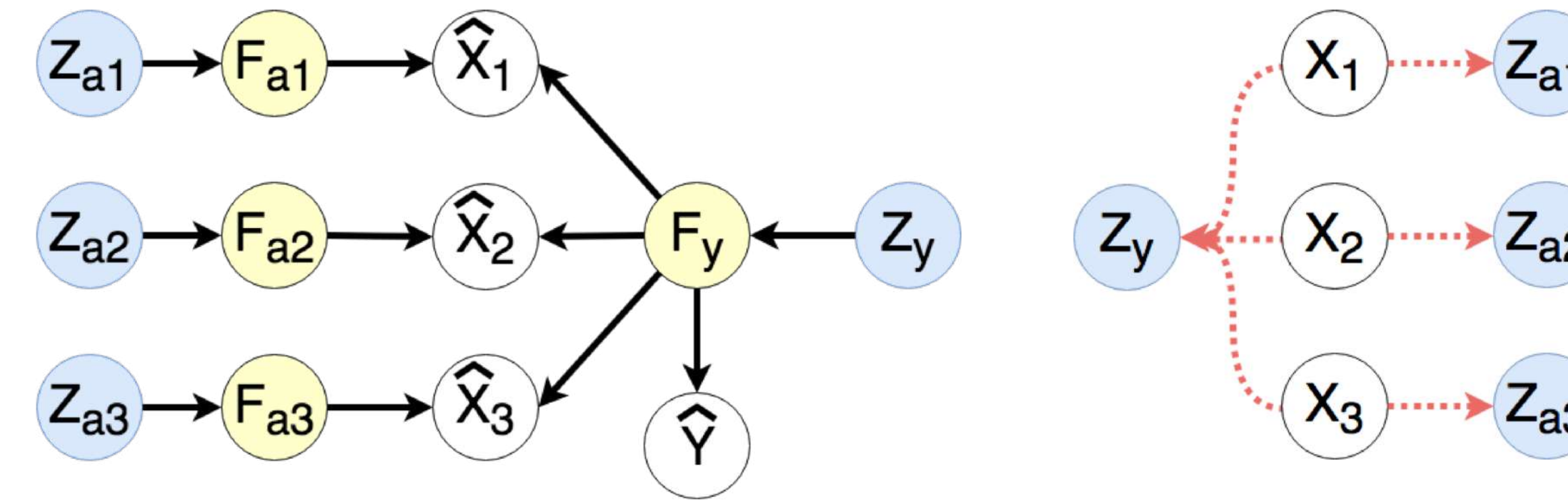
– **Language**: pre-trained GloVe word embeddings
– **Visual**: facial action units from Facet
– **Acoustic**: MFCCs from COVAREP
– Aligned by P2FA



Dataset	CMU-MOSI	ICT-MMMO	YouTube	MOUD	IEMOCAP	POM
Level	Segment	Video	Segment	Segment	Segment	Video
# Speakers	98	200	50	101	10	903
# Train	52→1284	220	30→169	49→243	5→6373	600
# Valid	10→229	40	5→41	10→37	1→1775	100
# Test	31→686	80	11→59	20→106	1→1807	203
Label	Sentiment	Sentiment	Sentiment	Sentiment	Emotions	Personalities
Features	{ ℓ, v, a }	{ ℓ, v, a }	{ ℓ, v, a }	{ ℓ, v, a }	{ ℓ, v, a }	{ ℓ, v, a }

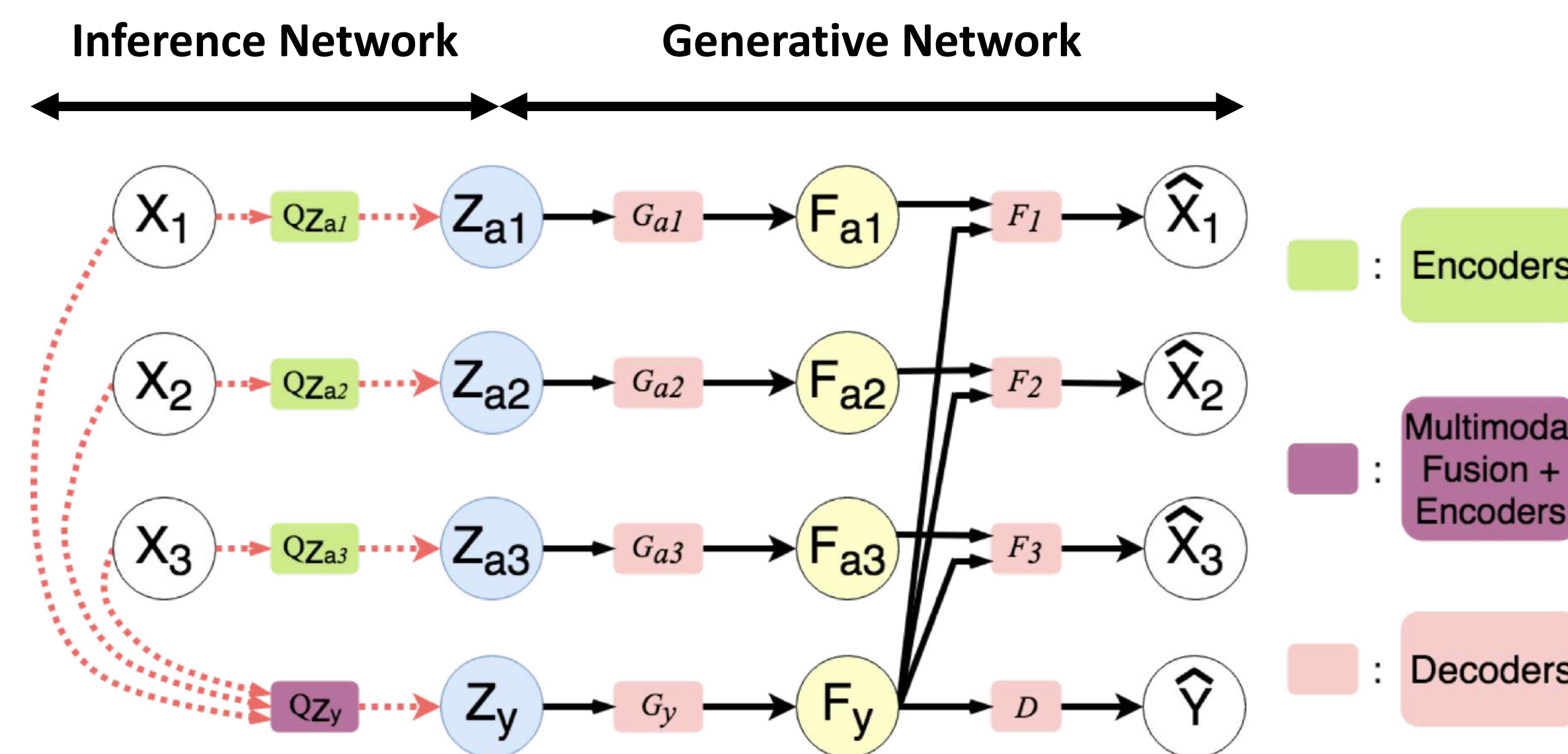
Multimodal Factorization Model (MFM)

• Bayesian Network



Generative Network

Inference Network



• Notations

- $\mathbf{X}_{1:M}$: multimodal data from M modalities, \mathbf{Y} : labels
- $\hat{\mathbf{X}}_{1:M}$: generated multimodal data, $\hat{\mathbf{Y}}$: generated labels
- $\mathbf{Z}_{a:}$: modality-specific latent variables, $\mathbf{F}_{:}$: factors

Generation: Factorization over Joint Distribution

$$P(\hat{\mathbf{X}}_{1:M}, \hat{\mathbf{Y}}) = \int_{\mathbf{F}, \mathbf{Z}} P(\hat{\mathbf{X}}_{1:M}, \hat{\mathbf{Y}} | \mathbf{F}) P(\mathbf{F} | \mathbf{Z}) P(\mathbf{Z}) d\mathbf{F} d\mathbf{Z}$$

$$= \int_{\mathbf{F}_Y, \mathbf{F}_{a\{1:M\}}} \left(P(\hat{\mathbf{Y}} | \mathbf{F}_Y) \prod_{i=1}^M P(\hat{\mathbf{X}}_i | \mathbf{F}_{a_i}, \mathbf{F}_Y) \right) \left(P(\mathbf{F}_Y | \mathbf{Z}) \prod_{i=1}^M P(\mathbf{F}_{a_i} | \mathbf{Z}_{a_i}) \right) \left(P(\mathbf{Z}_Y) \prod_{i=1}^M P(\mathbf{Z}_{a_i}) \right) d\mathbf{F} d\mathbf{Z}$$

$$\text{with } d\mathbf{F} = d\mathbf{F}_Y \prod_{i=1}^M d\mathbf{F}_{a_i} \text{ and } d\mathbf{Z} = d\mathbf{Z}_Y \prod_{i=1}^M d\mathbf{Z}_{a_i}$$

Inference: Joint-Distribution Wasserstein Distance

– **Proposition 1**: For any functions $G_Y : \mathbf{Z}_Y \rightarrow \mathbf{F}_Y$, $G_{a\{1:M\}} : \mathbf{Z}_{a\{1:M\}} \rightarrow \mathbf{F}_{a\{1:M\}}$, $D : \mathbf{F}_Y \rightarrow \hat{\mathbf{Y}}$, and $F_{1:M} : \mathbf{F}_{a\{1:M\}}, \mathbf{F}_Y \rightarrow \hat{\mathbf{X}}_{1:M}$, we have Joint-Distribution Wasserstein distance $W_c(P_{\mathbf{X}_{1:M}, \mathbf{Y}}, P_{\hat{\mathbf{X}}_{1:M}, \hat{\mathbf{Y}}}) =$

$$\inf_{Q_Z = P_Z} \mathbb{E}_{P_{\mathbf{X}_{1:M}, \mathbf{Y}}} \mathbb{E}_{Q(\mathbf{Z} | \mathbf{X}_{1:M}, \mathbf{Y})} \left[\sum_{i=1}^M c_{X_i} \left(\mathbf{X}_i, F_i(G_{a_i}(\mathbf{Z}_{a_i}), G_Y(\mathbf{Z}_Y)) \right) + c_Y \left(\mathbf{Y}, D(G_Y(\mathbf{Z}_Y)) \right) \right],$$

where P_Z is the prior over $\mathbf{Z} = [\mathbf{Z}_Y, \mathbf{Z}_{a\{1:M\}}]$ and Q_Z is the aggregated posterior of the proposed approximate inference distribution $Q(\mathbf{Z} | \mathbf{X}_{1:M}, \mathbf{Y})$.

Relaxed Generative-Discriminative Objective

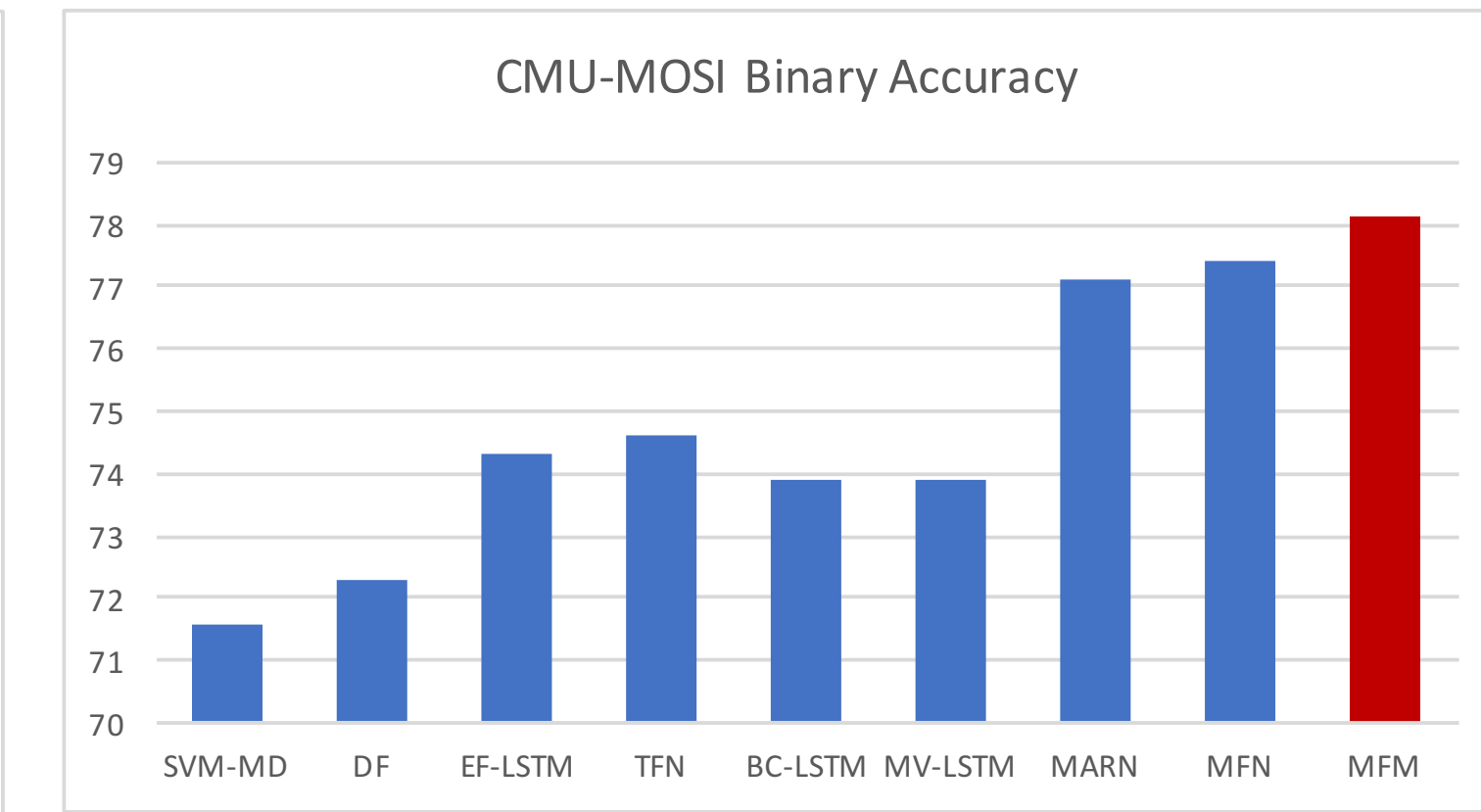
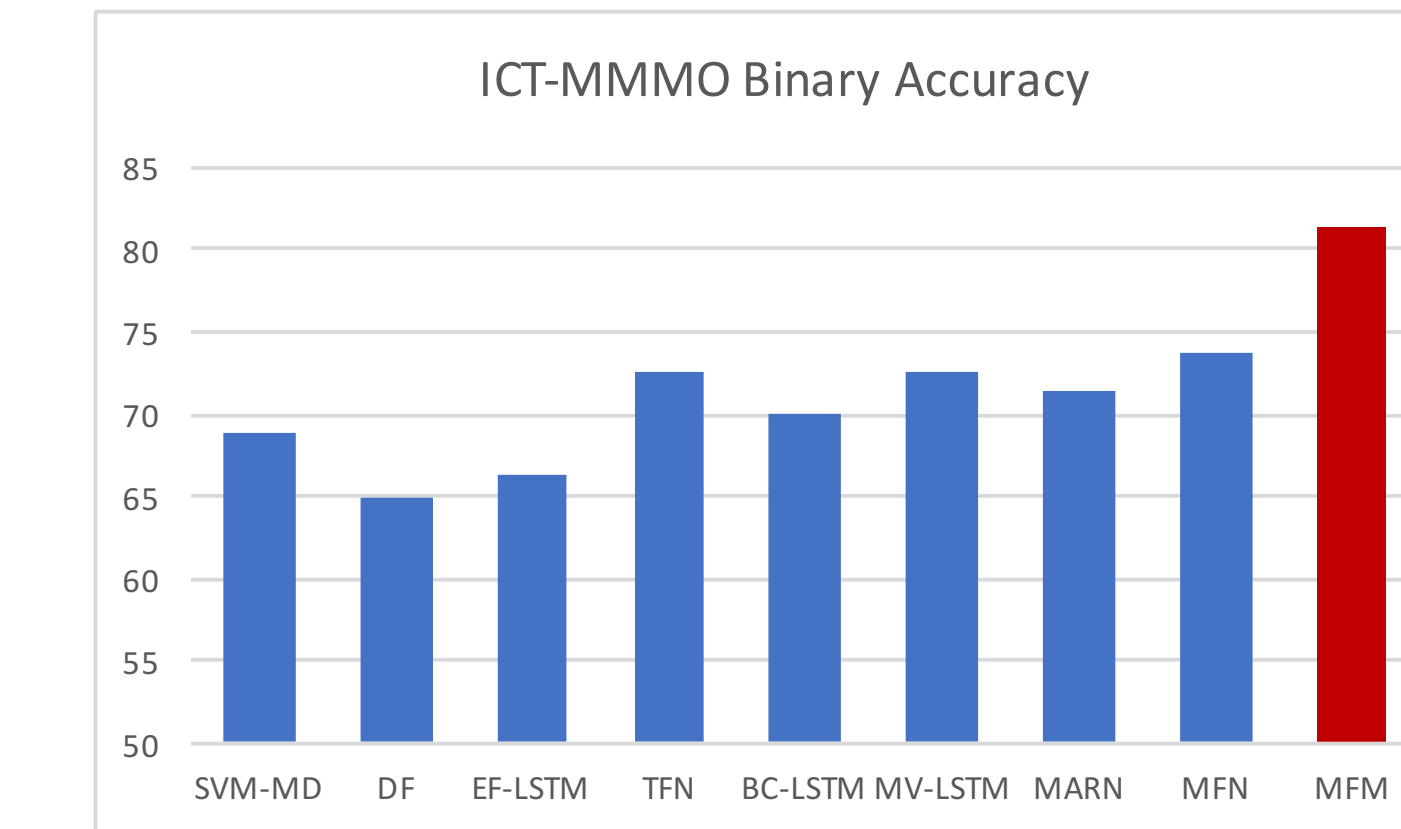
Mean-field assumption: $Q(\mathbf{Z} | \mathbf{X}_{1:M}, \mathbf{Y}) := Q(\mathbf{Z}_Y | \mathbf{X}_{1:M}, \mathbf{Y}) \prod_{i=1}^M Q(\mathbf{Z}_{a_i} | \mathbf{X}_i)$

Objective:

$$\min_{F, G_{a\{1:M\}}, G_Y, D} \inf_{Q(\mathbf{Z}_Y) \in \mathcal{Q}} \mathbb{E}_{P_{\mathbf{X}_{1:M}, \mathbf{Y}}} \mathbb{E}_{Q(\mathbf{Z}_{a1} | \mathbf{X}_1)} \cdots \mathbb{E}_{Q(\mathbf{Z}_{aM} | \mathbf{X}_M)} \mathbb{E}_{Q(\mathbf{Z}_Y | \mathbf{X}_{1:M})} \left[\sum_{i=1}^M c_{X_i} \left(\mathbf{X}_i, F(G_{a_i}(\mathbf{Z}_{a_i}), G_Y(\mathbf{Z}_Y)) \right) + c_Y \left(\mathbf{Y}, D(G_Y(\mathbf{Z}_Y)) \right) \right] + \lambda \mathcal{MMD}(Q_Z, P_Z),$$

Generative Discriminative Regularizer

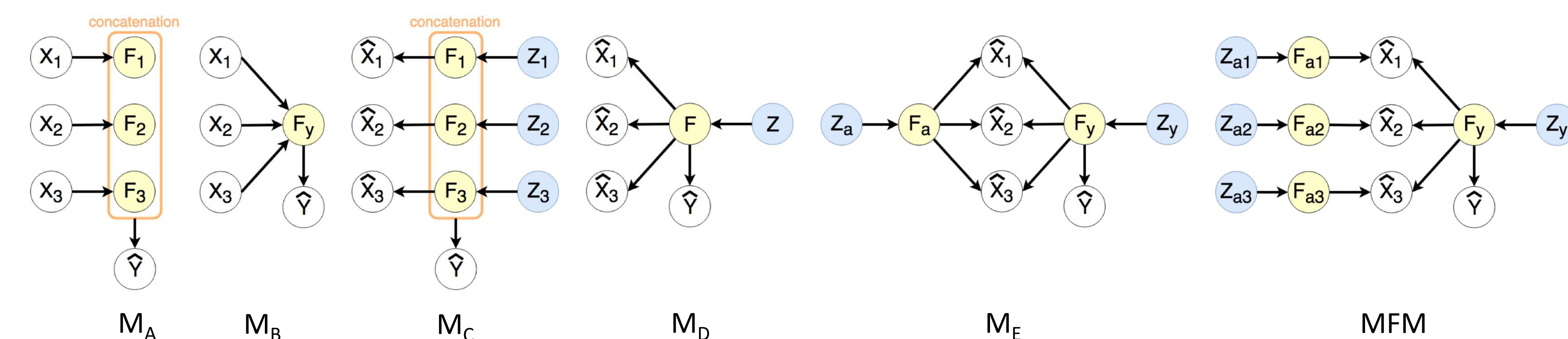
Results on Multimodal Time Series Datasets



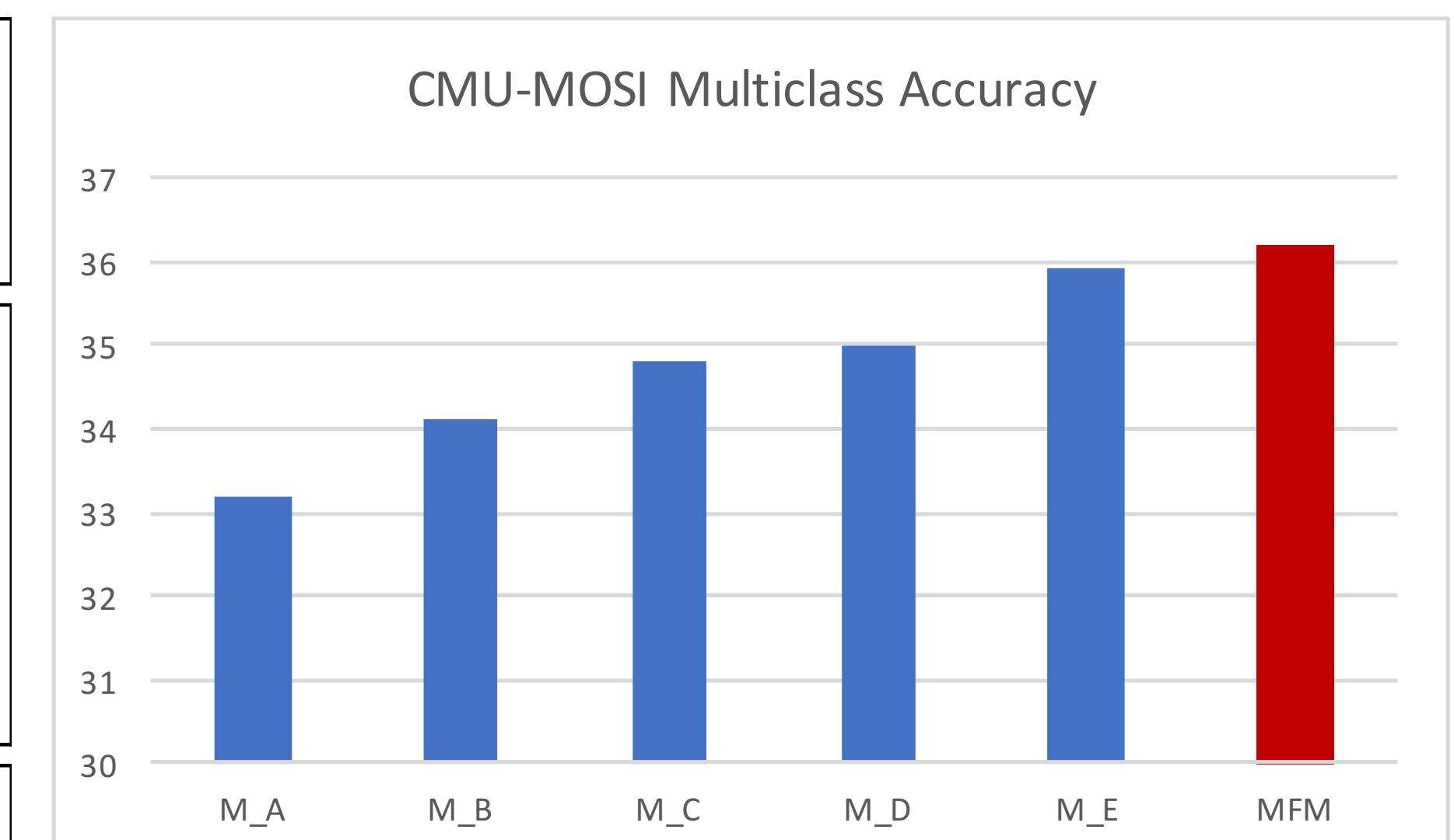
MFM achieves strong performance on 6 multimodal time-series datasets

MFM can be applied on any multimodal fusion encoder

Ablation Studies

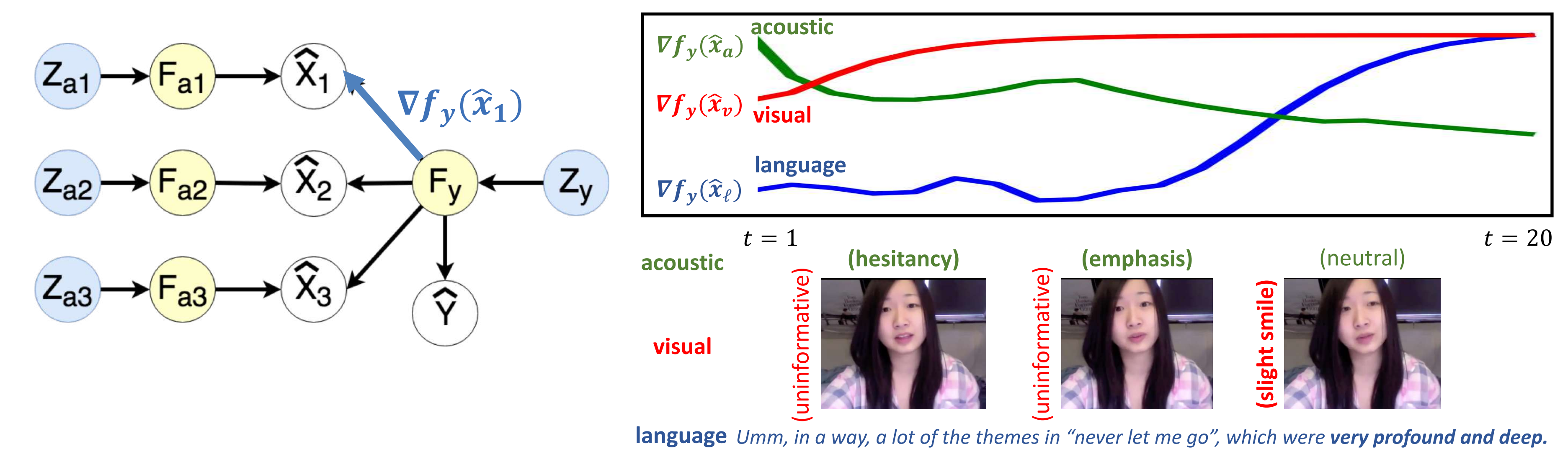


Model	Multimodal Disc. Factor	Hybrid Gen.-Disc. Objective	Factorized Gen.-Disc. Factors	Mod.-Spec. Gen. Factors
M_A	no	no	-	-
M_B	yes	no	-	-
M_C	no	yes	no	-
M_D	yes	yes	no	-
M_E	yes	yes	yes	no
MFM	yes	yes	yes	yes



IMPORTANT:
(1) Multimodal discriminative factor (2) Modality-specific generative factors (3) Hybrid generative-discriminative objective

Analyzing Multimodal Representations



Code and Models:
<http://github.com/pliang279/factorized>