



Improving Multimodal Fusion Via Mutual Dependency Maximisation

Oral Presentation at EMNLP 2021

Pierre Colombo^{🌸🌼}, Emile Chapuis[🌸]
Matthieu Labeau[🌸], Chloé Clavel[🌸]

What is our definition of multimodality?

What is our definition of multimodality?

Verbal

« What you say ? »

- Lexicon:
 - Words
- Syntax:
 - POS
- Pragmatics:
 - DA
 - Emotion

What is our definition of multimodality?

Verbal

« What you say ? »

- Lexicon:
 - Words
- Syntax:
 - POS
- Pragmatics:
 - DA
 - Emotion

« How you say it ? »

What is our definition of multimodality?

Verbal

« What you say ? »

- Lexicon:
 - Words
- Syntax:
 - POS
- Pragmatics:
 - DA
 - Emotion

« How you say it ? »

Vocal

- Prosody
 - Intonation
 - Voice quality
- Vocal expressions:
 - Laughter
 - Moans

What is our definition of multimodality?

Verbal

« What you say ? »

- Lexicon:
 - Words
- Syntax:
 - POS
- Pragmatics:
 - DA
 - Emotion

« How you say it ? »

Vocal

- Prosody
 - Intonation
 - Voice quality
- Vocal expressions:
 - Laughter
 - Moans

Visual

- Gestures:
 - Head & Eye & Arm
- Body language
 - Body posture
 - Proxemics
- Eye contact
 - Head & Eye gaze
- Facial expressions
 - FACS action units
 - Smile & Frowning

Core Challenges in Multimodal Learning

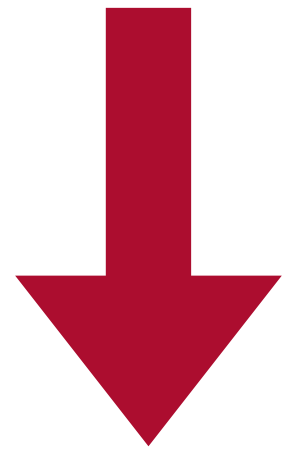
Core Challenges in Multimodal Learning

5 challenges of multimodal learning

Core Challenges in Multimodal Learning

5 challenges of multimodal learning

Representation



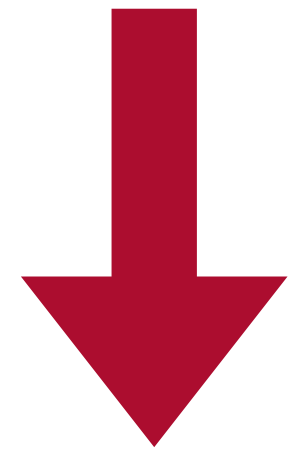
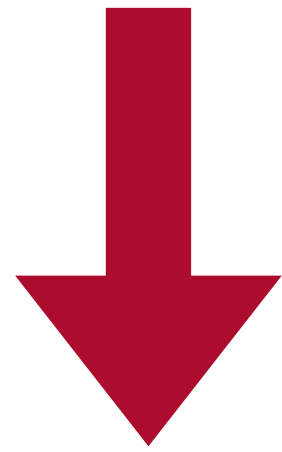
Represent
multimodal
data (leverage
complementarity,
redundancy)

Core Challenges in Multimodal Learning

5 challenges of multimodal learning

Representation

Alignement



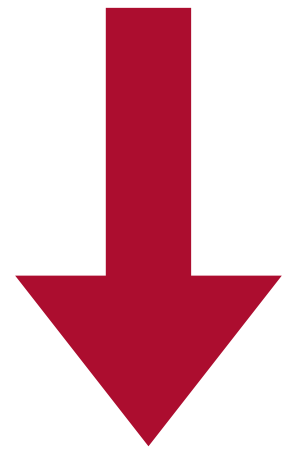
Represent
multimodal
data (leverage
complementarity,
redundancy)

Identify
relations
between
elements of
different
modalities

Core Challenges in Multimodal Learning

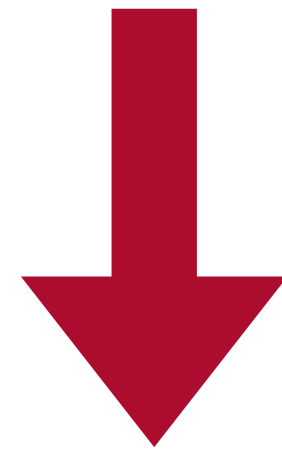
5 challenges of multimodal learning

Representation



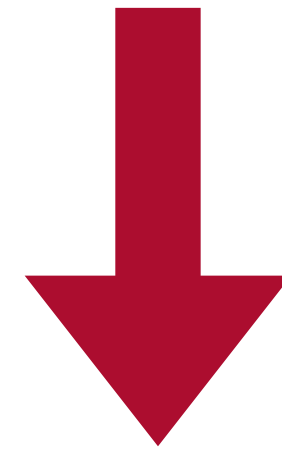
Represent
multimodal
data (leverage
complementarity,
redundancy)

Alignment



Identify
relations
between
elements of
different
modalities

Fusion

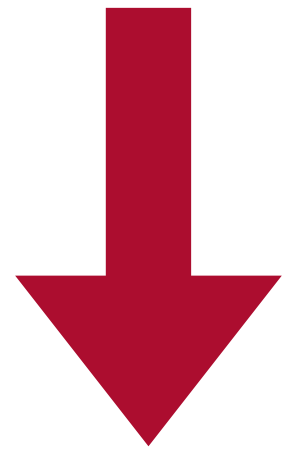


Join
information
from
modalities

Core Challenges in Multimodal Learning

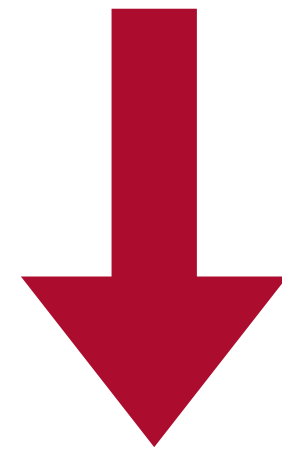
5 challenges of multimodal learning

Representation



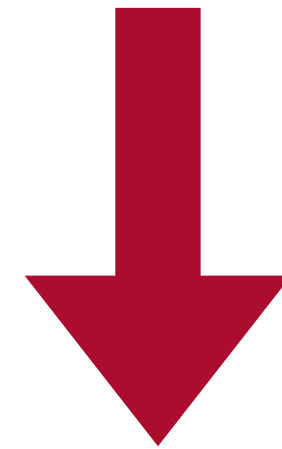
Represent
multimodal
data (leverage
complementarity,
redundancy)

Alignment



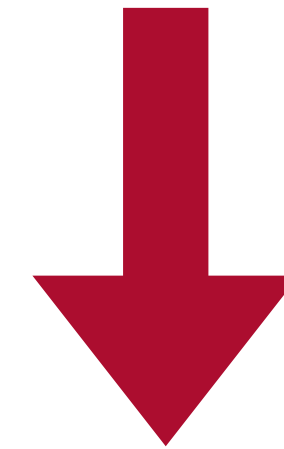
Identify
relations
between
elements of
different
modalities

Fusion



Join
information
from
modalities

Translation

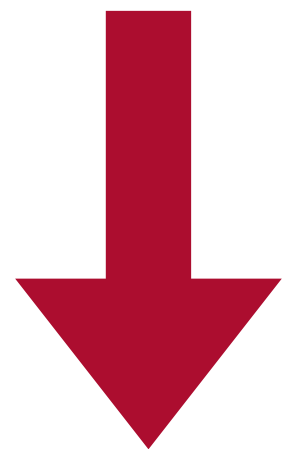


Translate
one
modality to
another

Core Challenges in Multimodal Learning

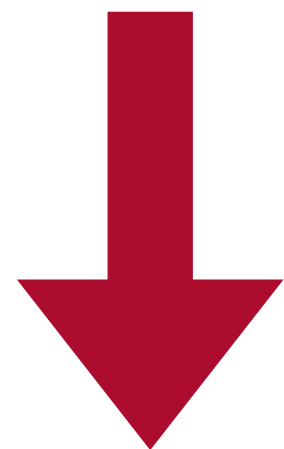
5 challenges of multimodal learning

Representation



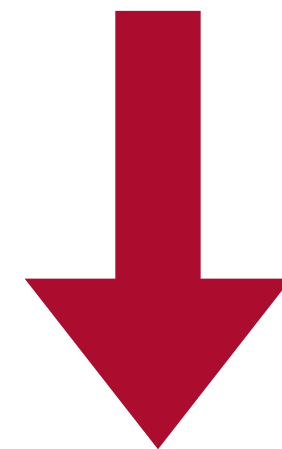
Represent
multimodal
data (leverage
complementarity,
redundancy)

Alignement



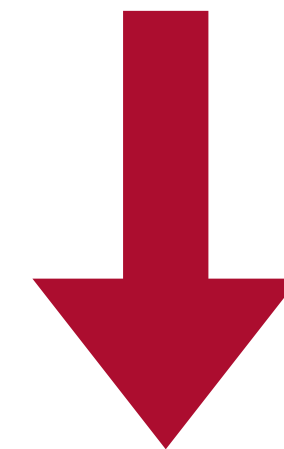
Identify
relations
between
elements of
different
modalities

Fusion



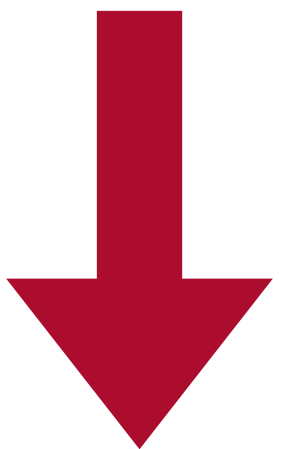
Join
information
from
modalities

Translation



Translate
one
modality to
another

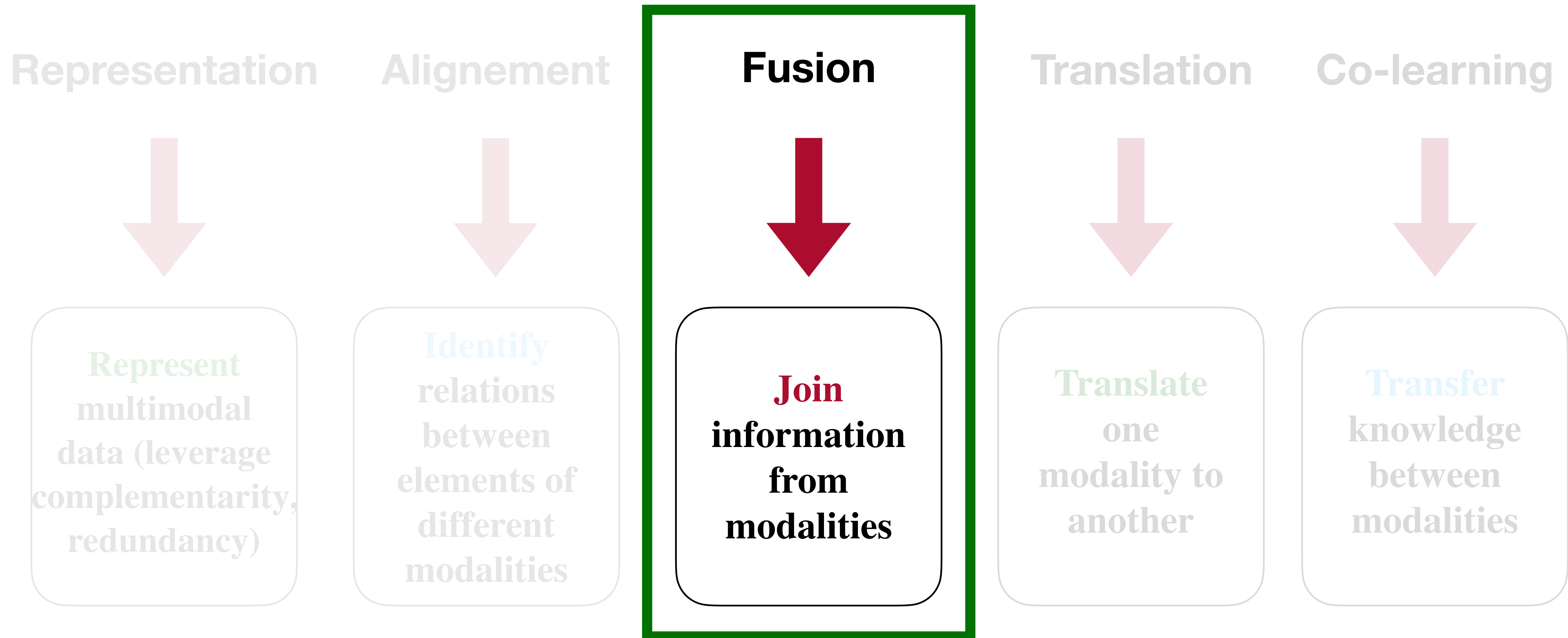
Co-learning



Transfer
knowledge
between
modalities

Core Challenges in Multimodal Learning

5 challenges of multimodal learning



Concurrent Works

Concurrent Works

Previous Work on **multimodal sentiment analysis**

Model		Fusion Mechanism
TFN	Zadeh et al. 2017	Tensor fusion
MARN	Zadeh et al. 2018	Multi attention block
LFN	Liu et al. 2018	Low rank tensor fusion
MFN	Zadeh et al 2018	Delta Memory attention network
MISA	Hazarika et al. 2020	Transformer - Multihead
MAGBERT/MAGXLNET	Rahman et al. 2020	Multimodal Adaptation Gate

Concurrent Works

Previous Work on **multimodal sentiment analysis**

Model		Fusion Mechanism
TFN	Zadeh et al. 2017	Tensor fusion
MARN	Zadeh et al. 2018	Multi attention block
LFN	Liu et al. 2018	Low rank tensor fusion
MFN	Zadeh et al 2018	Delta Memory attention network
MISA	Hazarika et al. 2020	Transformer - Multihead
MAGBERT/MAGXLNET	Rahman et al. 2020	Multimodal Adaptation Gate

Fusion in previous work is mainly neural

Few previous works on loss function !

Concurrent Works

Previous Work on **multimodal sentiment analysis**

Model		Fusion Mechanism
TFN	Zadeh et al. 2017	Tensor fusion
MARN	Zadeh et al. 2018	Multi attention block
LFN	Liu et al. 2018	Low rank tensor fusion
MFN	Zadeh et al 2018	Delta Memory attention network
MISA	Hazarika et al. 2020	Transformer - Multihead
MAGBERT/MAGXLNET	Rahman et al. 2020	Multimodal Adaptation Gate

Fusion in previous work is mainly neural

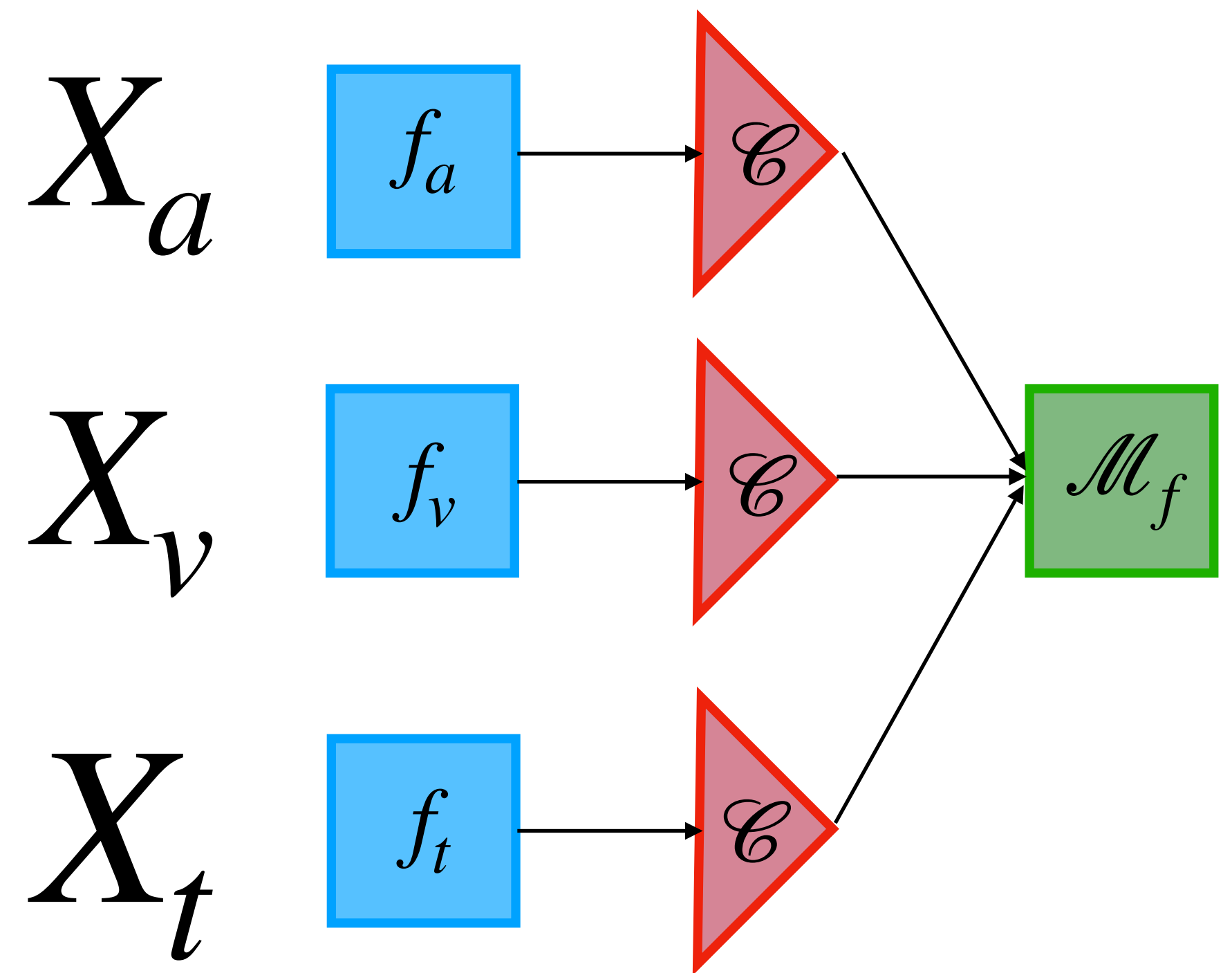
Few previous works on loss function !

Contribution: A new class of loss functions $\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$

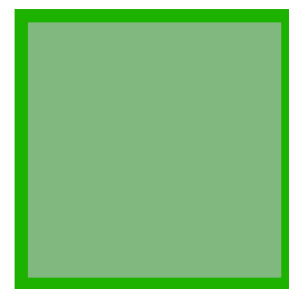
Background in Multimodal sentiment analysis

Background in Multimodal sentiment analysis

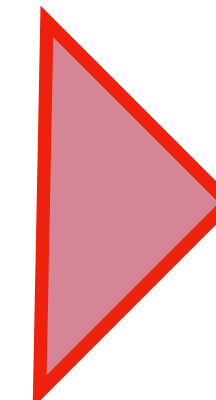
Late Fusion



Embedding Block



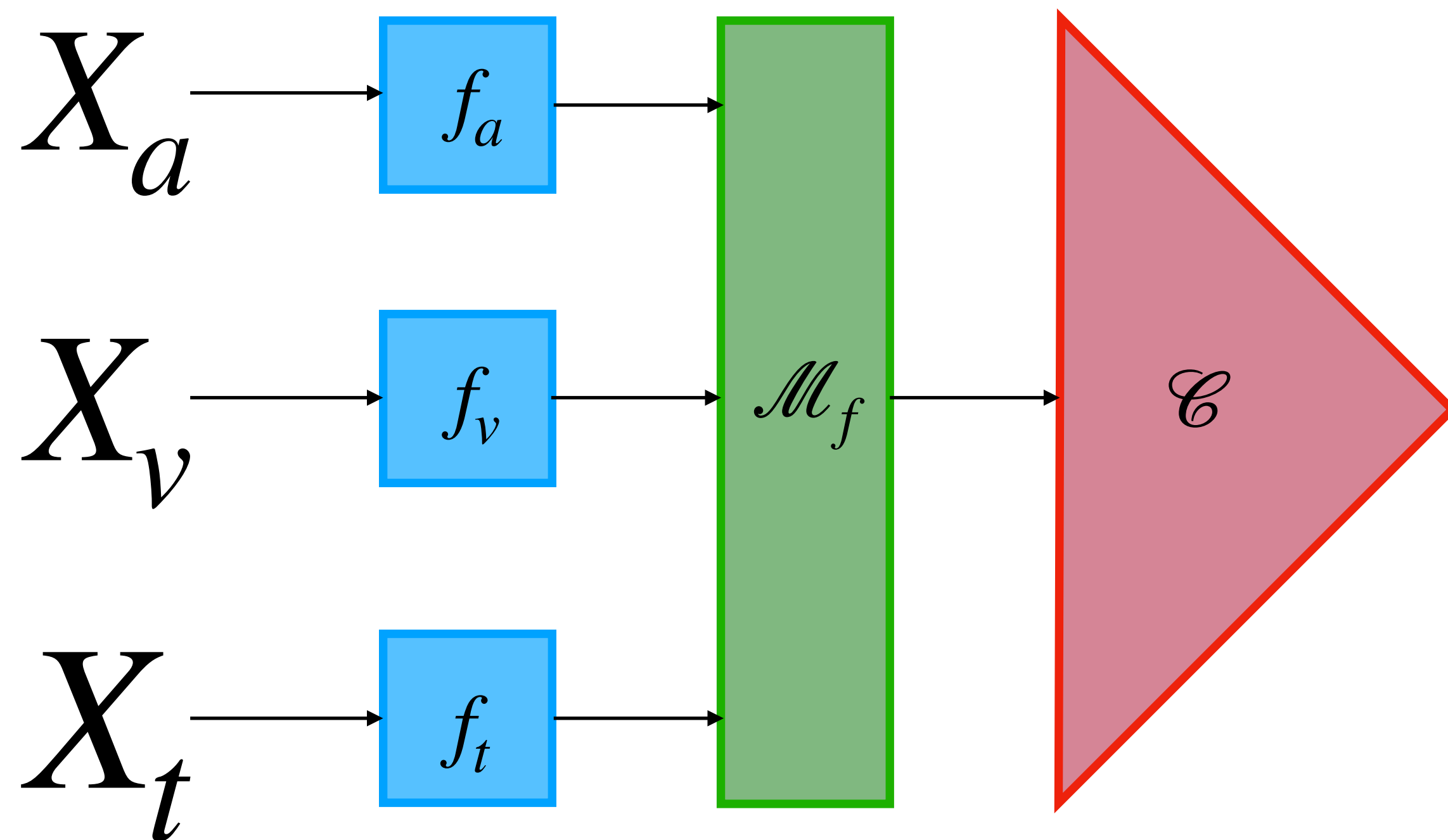
Fusion Block



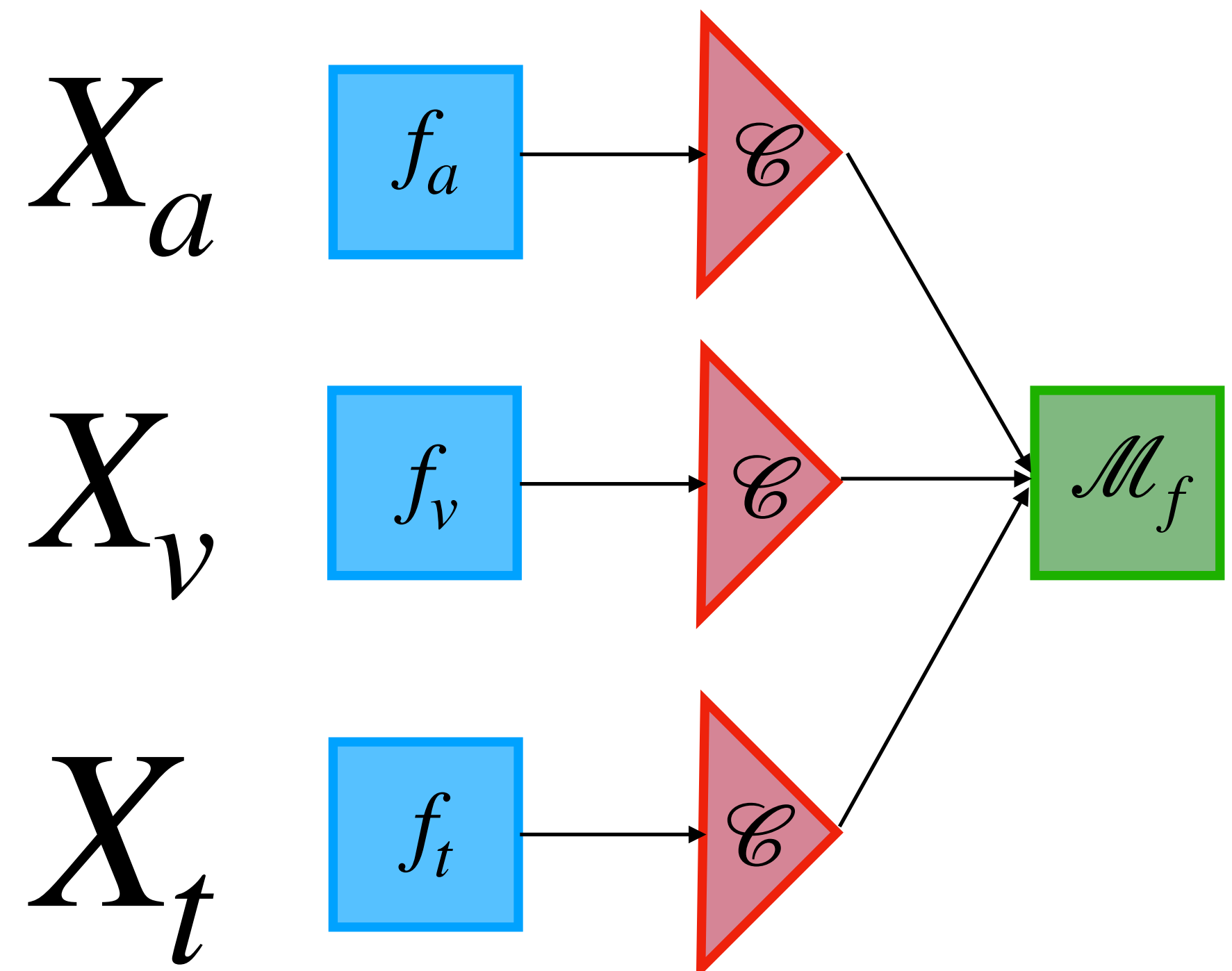
Classifier block

Background in Multimodal sentiment analysis

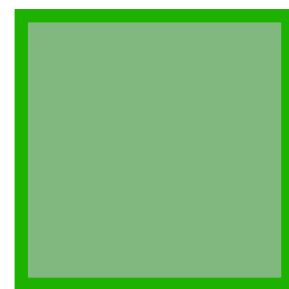
Early Fusion



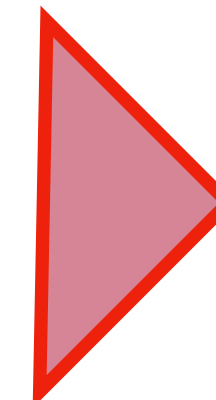
Late Fusion



Embedding Block



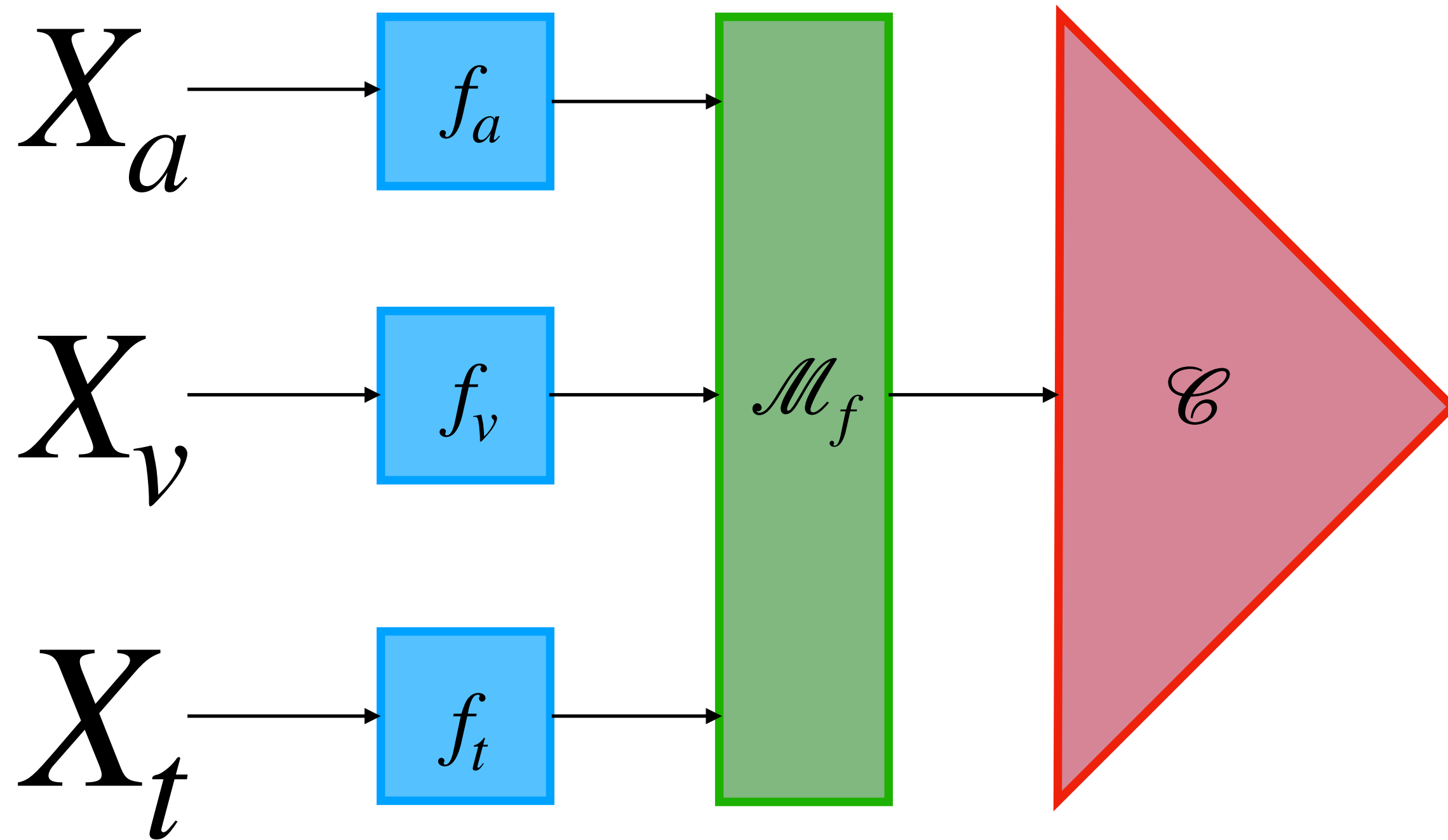
Fusion Block



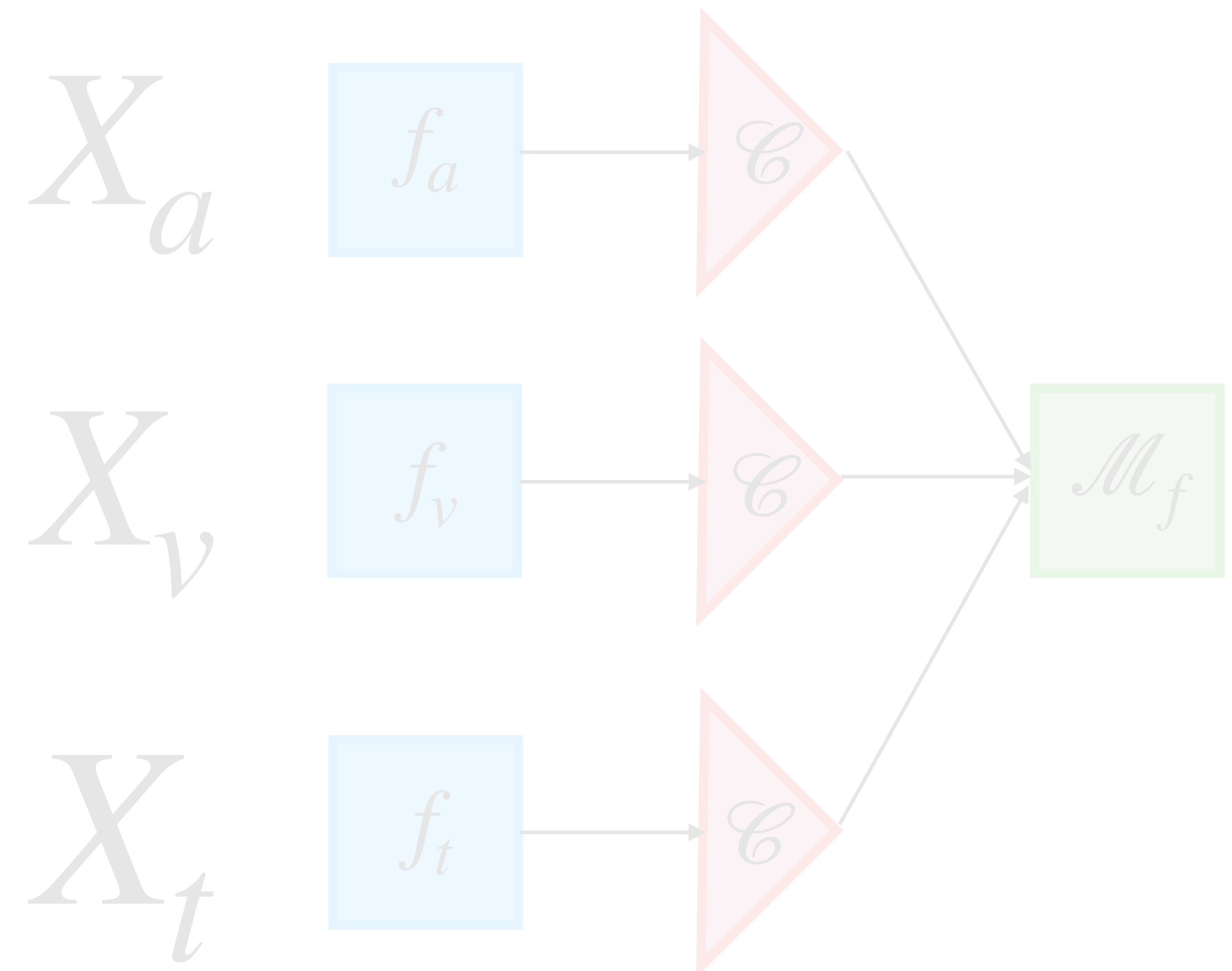
Classifier block

Background in Multimodal sentiment analysis

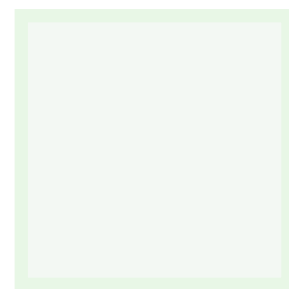
Early Fusion



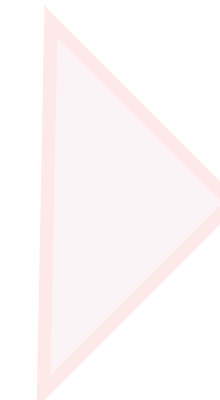
Late Fusion



Embedding Block



Fusion Block



Classifier block

Problem formulation

Problem formulation

Fusion Mechanism $\mathcal{M}_f : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathcal{R}^d$

Problem formulation

Fusion Mechanism $\mathcal{M}_f : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathcal{R}^d$

What we want for \mathcal{M}_f :

- Retain **modality-specific** interaction
- Retain **cross-view** interaction
- Retain **task specific** information

\mathcal{L}_{MDM}

$\mathcal{L}_{down.}$

Problem formulation

Fusion Mechanism $\mathcal{M}_f : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathcal{R}^d$

What we want for \mathcal{M}_f :

- Retain **modality-specific** interaction \mathcal{L}_{MDM}
- Retain **cross-view** interaction
- Retain **task specific** information $\mathcal{L}_{down.}$

Total $\underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

$$\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$$

Problem formulation

Fusion Mechanism $\mathcal{M}_f : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathcal{R}^d$

What we want for \mathcal{M}_f :

- Retain **modality-specific** interaction
- Retain **cross-view** interaction
- Retain **task specific** information

\mathcal{L}_{MDM}

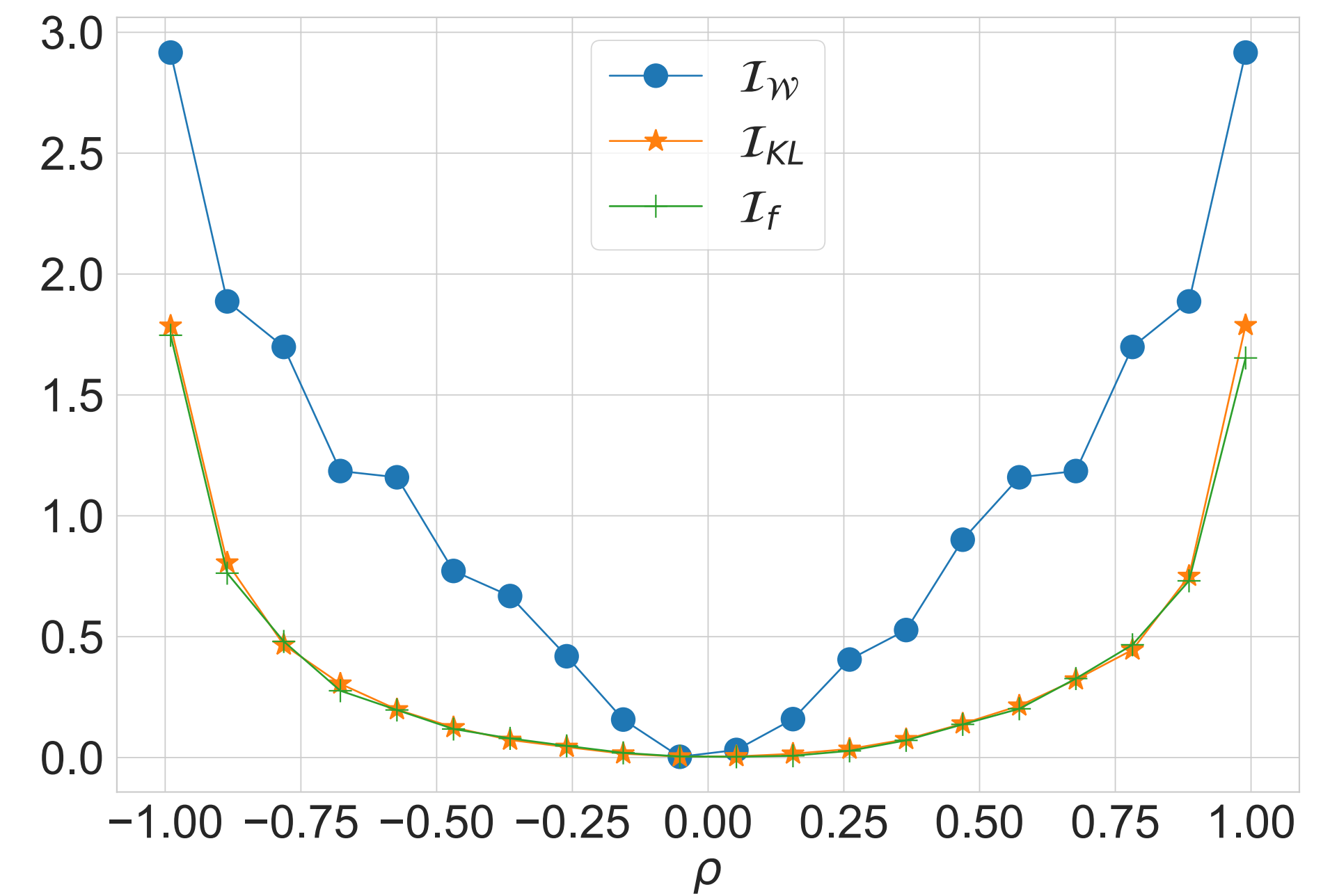
$\mathcal{L}_{down.}$

Total $\underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

$$\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$$

Different choice for MDM

- Wasserstein
- Kullback Leibler (KL)
- F-divergence



Toy example on correlated Gaussian

Training Losses

Training Losses

Training Loss $\mathcal{L} \triangleq \underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

Training Losses

Training Loss	$\mathcal{L} \triangleq \underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$
---------------	---

Notation	$X_a \sim p_{X_a}, X_v \sim p_{X_v}, X_l \sim p_{X_l}$	$X_a, X_v, X_l \sim p_{X_a, X_v, X_l}$
	$T(\theta) : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathbb{R}$	family of functions

Training Losses

Training Loss	$\mathcal{L} \triangleq \underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$
---------------	---

Notation	$X_a \sim p_{X_a}, X_v \sim p_{X_v}, X_l \sim p_{X_l}$	$X_a, X_v, X_l \sim p_{X_a, X_v, X_l}$
	$T(\theta) : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathbb{R}$	family of functions

For *MDM* we extend Belghazi et al 2018 (MINE)

Training Losses

Training Loss

$$\mathcal{L} \triangleq \underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$$

Notation

$$X_a \sim p_{X_a}, X_v \sim p_{X_v}, X_l \sim p_{X_l}$$

$$T(\theta) : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathbb{R}$$

$$X_a, X_v, X_l \sim p_{X_a, X_v, X_l}$$

family of functions

For *MDM* we extend Belghazi et al 2018 (MINE)

	Formula
Wasserstein (\mathcal{W})	$\mathbf{I}_{\mathcal{W}} \triangleq \sup_{\theta: T_{\theta} \in \mathbb{L}} \mathbb{E}_{p_{X_a X_v X_l}} [T_{\theta}] - \log \left[\mathbb{E}_{\prod_{j \in \{a, v, l\}} p_{X_j}} [T_{\theta}] \right].$
f -divergence (f)	$\mathbf{I}_f \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}} [T_{\theta}] - \mathbb{E}_{\prod_{j \in \{a, v, l\}} p_{X_j}} [e^{T_{\theta}-1}].$
Kullback-Leibler (KL)	$\mathbf{I}_{kl} \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}} [T_{\theta}] - \log \left[\mathbb{E}_{\prod_{j \in \{a, v, l\}} p_{X_j}} [e^{T_{\theta}}] \right].$

Training Losses

Training Loss

$$\mathcal{L} \triangleq \underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$$

Notation

$$X_a \sim p_{X_a}, X_v \sim p_{X_v}, X_l \sim p_{X_l}$$
$$T(\theta) : \mathcal{X}_a \times \mathcal{X}_v \times \mathcal{X}_l \rightarrow \mathbb{R}$$

$$X_a, X_v, X_l \sim p_{X_a, X_v, X_l}$$

family of functions

For *MDM* we extend Belghazi et al 2018 (MINE)

	Formula	Alternative Work
Wasserstein (\mathcal{W})	$\mathbf{I}_{\mathcal{W}} \triangleq \sup_{\theta: T_{\theta} \in \mathbb{L}} \mathbb{E}_{p_{X_a X_v X_l}} [T_{\theta}] - \log \left[\mathbb{E}_{\prod_{j \in \{a, v, l\}} p_{X_j}} [T_{\theta}] \right].$	Ozair et al 2019.
f -divergence (f)	$\mathbf{I}_f \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}} [T_{\theta}] - \mathbb{E}_{\prod_{j \in \{a, v, l\}} p_{X_j}} [e^{T_{\theta}-1}].$	KNIFE. Anonymous et al. 2022
Kullback-Leibler (KL)	$\mathbf{I}_{kl} \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}} [T_{\theta}] - \log \left[\mathbb{E}_{\prod_{j \in \{a, v, l\}} p_{X_j}} [e^{T_{\theta}}] \right].$	Oord et al 2018

Practical Implementation

Practical Implementation

Our method requiers to compute

$$MDM \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}}[T_{\theta}] - g \left[\mathbb{E}_{\prod_{j \in \{a,v,l\}} p_{X_j}}[\cdot] \right].$$

Practical Implementation

Our method requires to compute

$$MDM \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}}[T_{\theta}] - g \left[\mathbb{E} \prod_{j \in \{a, v, l\}} p_{X_j}[\cdot] \right].$$

Two Stage Procedure

Algorithm 1 Two-stage procedure to minimise multivariate dependency measures.

INPUT: $\mathcal{D}_n = \{(x_a^j, x_v^j, x_l^j), \forall j \in [1, n]\}$ multi-modal training dataset, m batch size, $\sigma_a, \sigma_v, \sigma_l : [1, m] \rightarrow [1, m]$ three permutations, θ_c weights of the deep classifier, θ weights of the statistical network T_{θ} .

Initialization: parameters θ and θ_c

Build Negative Dataset:

$$\bar{\mathcal{D}}_n = \{(x_a^{\sigma_a(j)}, x_v^{\sigma_v(j)}, x_l^{\sigma_l(j)}), \forall j \in [1, n]\}$$

Optimization:

while (θ, θ_c) not converged **do**

for $i \in [1, Unroll]$ **do**

 Sample from \mathcal{D}_n , $\mathcal{B} \sim p_{X_a X_v X_l}$

 Sample from $\bar{\mathcal{D}}_n$, $\bar{\mathcal{B}} \sim \prod_{j \in \{a, v, l\}} p_{X_j}$

 Update θ based on the empirical version of Eq. 6 or Eq. 7 or Eq. 8.

end for

 Sample a batch \mathcal{B} from \mathcal{D}

 Update θ_c with \mathcal{B} using Eq. 5.

end while

OUTPUT: Classifiers weights θ_c

Practical Implementation

Our method requires to compute

$$MDM \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}}[T_{\theta}] - g \left[\mathbb{E} \prod_{j \in \{a, v, l\}} p_{X_j}[\cdot] \right].$$

Two Stage Procedure

1. Update θ to find the supremum

Algorithm 1 Two-stage procedure to minimise multivariate dependency measures.

INPUT: $\mathcal{D}_n = \{(x_a^j, x_v^j, x_l^j), \forall j \in [1, n]\}$ multi-modal training dataset, m batch size, $\sigma_a, \sigma_v, \sigma_l : [1, m] \rightarrow [1, m]$ three permutations, θ_c weights of the deep classifier, θ weights of the statistical network T_{θ} .

Initialization: parameters θ and θ_c

Build Negative Dataset:

$$\bar{\mathcal{D}}_n = \{(x_a^{\sigma_a(j)}, x_v^{\sigma_v(j)}, x_l^{\sigma_l(j)}), \forall j \in [1, n]\}$$

Optimization:

while (θ, θ_c) not converged **do**

for $i \in [1, Unroll]$ **do**

 Sample from \mathcal{D}_n , $\mathcal{B} \sim p_{X_a X_v X_l}$

 Sample from $\bar{\mathcal{D}}_n$, $\bar{\mathcal{B}} \sim \prod_{j \in \{a, v, l\}} p_{X_j}$

 Update θ based on the empirical version of Eq. 6 or Eq. 7 or Eq. 8.

end for

 Sample a batch \mathcal{B} from \mathcal{D}

 Update θ_c with \mathcal{B} using Eq. 5.

end while

OUTPUT: Classifiers weights θ_c

Practical Implementation

Our method requires to compute

$$MDM \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}}[T_{\theta}] - g \left[\mathbb{E} \prod_{j \in \{a, v, l\}} p_{X_j}[\cdot] \right].$$

Two Stage Procedure

1. Update θ to find the supremum

2. Update neural network using the total loss

Algorithm 1 Two-stage procedure to minimise multivariate dependency measures.

INPUT: $\mathcal{D}_n = \{(x_a^j, x_v^j, x_l^j), \forall j \in [1, n]\}$ multi-modal training dataset, m batch size, $\sigma_a, \sigma_v, \sigma_l : [1, m] \rightarrow [1, m]$ three permutations, θ_c weights of the deep classifier, θ weights of the statistical network T_{θ} .

Initialization: parameters θ and θ_c

Build Negative Dataset:

$$\bar{\mathcal{D}}_n = \{(x_a^{\sigma_a(j)}, x_v^{\sigma_v(j)}, x_l^{\sigma_l(j)}), \forall j \in [1, n]\}$$

Optimization:

while (θ, θ_c) not converged **do**

for $i \in [1, Unroll]$ **do**

 Sample from \mathcal{D}_n , $\mathcal{B} \sim p_{X_a X_v X_l}$

 Sample from $\bar{\mathcal{D}}_n$, $\bar{\mathcal{B}} \sim \prod_{j \in \{a, v, l\}} p_{X_j}$

 Update θ based on the empirical version of Eq. 6 or Eq. 7 or Eq. 8.

end for

 Sample a batch \mathcal{B} from \mathcal{D}

 Update θ_c with \mathcal{B} using Eq. 5.

end while

OUTPUT: Classifiers weights θ_c

Practical Implementation

Our method requires to compute

$$MDM \triangleq \sup_{\theta} \mathbb{E}_{p_{X_a X_v X_l}}[T_{\theta}] - g \left[\mathbb{E} \prod_{j \in \{a, v, l\}} p_{X_j}[\cdot] \right].$$

Two Stage Procedure

1. Update θ to find the supremum

2. Update neural network using the total loss

Algorithm 1 Two-stage procedure to minimise multivariate dependency measures.

INPUT: $\mathcal{D}_n = \{(x_a^j, x_v^j, x_l^j), \forall j \in [1, n]\}$ multi-modal training dataset, m batch size, $\sigma_a, \sigma_v, \sigma_l : [1, m] \rightarrow [1, m]$ three permutations, θ_c weights of the deep classifier, θ weights of the statistical network T_{θ} .

Initialization: parameters θ and θ_c

Build Negative Dataset:

$$\bar{\mathcal{D}}_n = \{(x_a^{\sigma_a(j)}, x_v^{\sigma_v(j)}, x_l^{\sigma_l(j)}), \forall j \in [1, n]\}$$

Optimization:

while (θ, θ_c) not converged **do**

for $i \in [1, Unroll]$ **do**

 Sample from \mathcal{D}_n , $\mathcal{B} \sim p_{X_a X_v X_l}$

 Sample from $\bar{\mathcal{D}}_n$, $\bar{\mathcal{B}} \sim \prod_{j \in \{a, v, l\}} p_{X_j}$

 Update θ based on the empirical version of Eq. 6 or Eq. 7 or Eq. 8.

end for

 Sample a batch \mathcal{B} from \mathcal{D}

 Update θ_c with \mathcal{B} using Eq. 5.

end while

OUTPUT: Classifiers weights θ_c

Summary & Experimental Setting

Summary & Experimental Setting

Summary

Total $\underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

$$\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$$

Summary & Experimental Setting

Summary

Total $\underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

$$\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$$

Require little intervention

+

can work on any models

Summary & Experimental Setting

Summary

Total $\underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

$$\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$$

Require little intervention

+

can work on any models

Experimental Setting

- Sentiment Analysis
- Robustness Analysis
- Explainability using T_θ

Summary & Experimental Setting

Summary

Total $\underbrace{\mathcal{L}_{down.}}_{\text{main task}} - \underbrace{\lambda \cdot \mathcal{L}_{MDM}}_{\text{mutual dependency term}}$

$$\mathcal{L}_{MDM} \triangleq MDM \left(p_{X_a X_v X_l}(x_a, x_v, x_l), \prod_{j \in \{a, v, l\}} p_{X_j}(x_j) \right)$$

Require little intervention

+

can work on any models

Experimental Setting

- Sentiment Analysis
- Robustness Analysis
- Explainability using T_θ

Other results available
in the paper

Overall Results

Overall Results

CMU-MOSEI

CMU-MOSI

2,199/23,454 movie review videos

Sentiment Score in [-3,3]

Overall Results

CMU-MOSEI

CMU-MOSI

2,199/23,454 movie review videos

Sentiment Score in [-3,3]

	Acc_7^h	Acc_2^h	MAE^l	$Corr^h$
CMU-MOSI				
\mathcal{L}_\emptyset	31.1	76.1	1.00	0.65
\mathcal{L}_{kl}	<u>31.7</u>	<u>76.4</u>	1.00	<u>0.66</u>
\mathcal{L}_f	<u>33.7</u>	76.2	1.02	<u>0.66</u>
$\mathcal{L}_{\mathcal{W}}$	<u>33.5</u>	<u>76.4</u>	<u>0.98</u>	<u>0.66</u>
CMU-MOSEI				
\mathcal{L}_\emptyset	44.2	75.0	0.72	0.52
\mathcal{L}_{kl}	44.5	<u>75.6</u>	<u>0.70</u>	<u>0.53</u>
\mathcal{L}_f	<u>45.5</u>	75.2	<u>0.70</u>	0.52
$\mathcal{L}_{\mathcal{W}}$	<u>45.3</u>	<u>75.9</u>	<u>0.68</u>	<u>0.54</u>

	CMU-MOSI				CMU-MOSEI			
	Acc_7^h	Acc_2^h	MAE^l	$Corr^h$	Acc_7^h	Acc_2^h	MAE^l	$Corr^h$
MFN								
\mathcal{L}_\emptyset	31.3	76.6	1.01	0.62	44.4	74.7	0.72	0.53
\mathcal{L}_{kl}	<u>32.5</u>	76.7	<u>0.96</u>	0.65	44.2	74.7	0.72	<u>0.57</u>
\mathcal{L}_f	<u>35.7</u>	<u>77.4</u>	<u>0.96</u>	0.65	<u>46.1</u>	75.4	<u>0.69</u>	<u>0.56</u>
$\mathcal{L}_{\mathcal{W}}$	<u>35.9</u>	<u>77.6</u>	<u>0.96</u>	0.65	<u>46.2</u>	75.1	<u>0.69</u>	<u>0.56</u>
LFN								
\mathcal{L}_\emptyset	31.9	76.9	1.00	0.63	45.2	74.2	0.70	0.54
\mathcal{L}_{kl}	<u>32.6</u>	<u>77.7</u>	0.97	0.63	<u>46.1</u>	75.3	0.68	<u>0.57</u>
\mathcal{L}_f	<u>35.6</u>	77.1	0.97	0.63	45.8	<u>75.4</u>	0.69	<u>0.57</u>
$\mathcal{L}_{\mathcal{W}}$	<u>35.6</u>	<u>77.7</u>	<u>0.96</u>	<u>0.67</u>	<u>46.2</u>	<u>75.4</u>	<u>0.67</u>	<u>0.57</u>
MAGBERT								
\mathcal{L}_\emptyset	40.2	84.7	0.79	0.80	46.8	84.9	0.59	0.77
\mathcal{L}_{kl}	<u>42.0</u>	<u>85.6</u>	<u>0.76</u>	0.82	47.1	85.4	0.59	<u>0.79</u>
\mathcal{L}_f	<u>41.7</u>	<u>85.6</u>	0.78	0.82	46.9	85.6	0.59	<u>0.79</u>
$\mathcal{L}_{\mathcal{W}}$	<u>41.8</u>	85.3	<u>0.76</u>	0.82	<u>47.8</u>	85.5	0.59	<u>0.79</u>
MAGXLNET								
\mathcal{L}_\emptyset	43.0	86.2	0.76	0.82	46.7	84.4	0.59	0.79
\mathcal{L}_{kl}	<u>44.5</u>	86.1	<u>0.74</u>	0.82	<u>47.5</u>	<u>85.4</u>	0.59	0.81
\mathcal{L}_f	<u>43.9</u>	86.6	<u>0.74</u>	0.82	47.4	85.0	0.59	0.81
$\mathcal{L}_{\mathcal{W}}$	<u>44.4</u>	<u>86.9</u>	<u>0.74</u>	0.82	<u>47.9</u>	<u>85.8</u>	0.59	<u>0.82</u>

Overall Results

CMU-MOSEI	CMU-MOSI
2,199/23,454 movie review videos	
Sentiment Score in [-3,3]	

Simple fusion mechanism

	Acc_7^h	Acc_2^h	MAE^l	$Corr^h$
--	-----------	-----------	---------	----------

CMU-MOSI

\mathcal{L}_\emptyset	31.1	76.1	1.00	0.65
\mathcal{L}_{kl}	<u>31.7</u>	<u>76.4</u>	1.00	<u>0.66</u>
\mathcal{L}_f	33.7	76.2	1.02	0.66
\mathcal{L}_W	33.5	76.4	0.98	0.66

CMU-MOSEI

\mathcal{L}_\emptyset	44.2	75.0	0.72	0.52
\mathcal{L}_{kl}	44.5	<u>75.6</u>	<u>0.70</u>	<u>0.53</u>
\mathcal{L}_f	45.5	75.2	<u>0.70</u>	0.52
\mathcal{L}_W	<u>45.3</u>	75.9	<u>0.68</u>	<u>0.54</u>

Complex fusion mechanism

	CMU-MOSI				CMU-MOSEI			
	Acc_7^h	Acc_2^h	MAE^l	$Corr^h$	Acc_7^h	Acc_2^h	MAE^l	$Corr^h$

MFN

\mathcal{L}_\emptyset	31.3	76.6	1.01	0.62	44.4	74.7	0.72	0.53
\mathcal{L}_{kl}	<u>32.5</u>	76.7	<u>0.96</u>	0.65	44.2	74.7	0.72	<u>0.57</u>
\mathcal{L}_f	35.7	77.4	0.96	0.65	46.1	75.4	0.69	0.56
\mathcal{L}_W	35.9	77.6	0.96	0.65	46.2	75.1	0.69	0.56

LFN

\mathcal{L}_\emptyset	31.9	76.9	1.00	0.63	45.2	74.2	0.70	0.54
\mathcal{L}_{kl}	<u>32.6</u>	<u>77.7</u>	0.97	0.63	<u>46.1</u>	75.3	0.68	<u>0.57</u>
\mathcal{L}_f	35.6	77.1	0.97	0.63	45.8	75.4	0.69	0.57
\mathcal{L}_W	<u>35.6</u>	<u>77.7</u>	<u>0.96</u>	<u>0.67</u>	<u>46.2</u>	<u>75.4</u>	<u>0.67</u>	<u>0.57</u>

MAGBERT

\mathcal{L}_\emptyset	40.2	84.7	0.79	0.80	46.8	84.9	0.59	0.77
\mathcal{L}_{kl}	<u>42.0</u>	<u>85.6</u>	<u>0.76</u>	0.82	47.1	85.4	0.59	<u>0.79</u>
\mathcal{L}_f	41.7	85.6	0.78	0.82	46.9	85.6	0.59	0.79
\mathcal{L}_W	<u>41.8</u>	85.3	<u>0.76</u>	0.82	<u>47.8</u>	85.5	0.59	<u>0.79</u>

MAGXLNET

\mathcal{L}_\emptyset	43.0	86.2	0.76	0.82	46.7	84.4	0.59	0.79
\mathcal{L}_{kl}	<u>44.5</u>	86.1	<u>0.74</u>	0.82	<u>47.5</u>	<u>85.4</u>	0.59	0.81
\mathcal{L}_f	<u>43.9</u>	86.6	<u>0.74</u>	0.82	47.4	85.0	0.59	0.81
\mathcal{L}_W	<u>44.4</u>	86.9	<u>0.74</u>	0.82	<u>47.9</u>	85.8	0.59	<u>0.82</u>

Robustness to modality drop

Robustness to modality drop

1. Most of information carried by Text

- **Drop Text**

Robustness to modality drop

1. Most of information carried by Text

- **Drop Text**

How maximising the MDM affect robustness?

Robustness to modality drop

1. Most of information carried by Text

- **Drop Text**

How maximising the MDM affect robustness?

Experiment

- 1. Train using 3 modalities**
- 2. At inference we use only A, V or A+V**

Robustness to modality drop

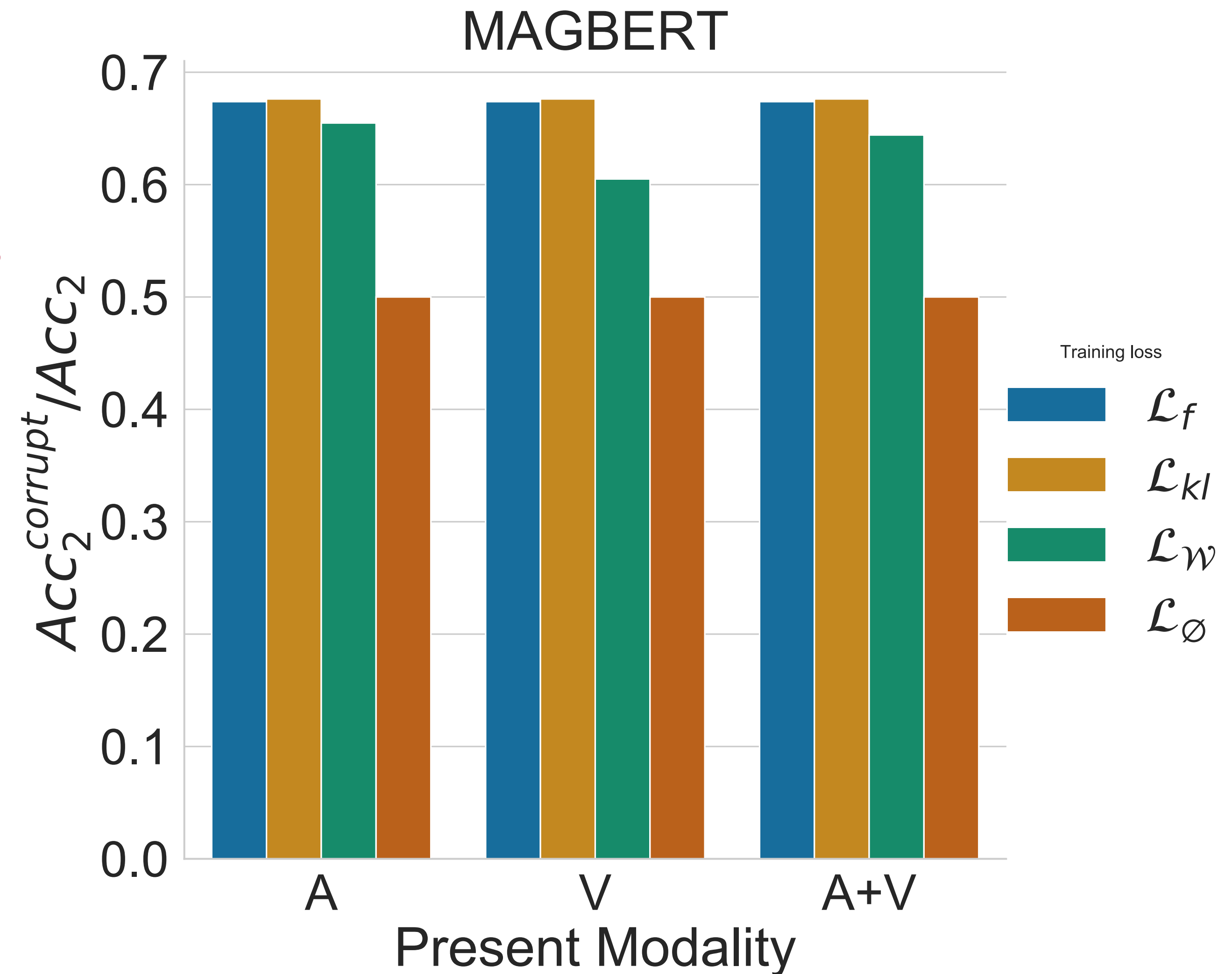
1. Most of information carried by Text

- Drop Text

How maximising the MDM affect robustness?

Experiment

1. Train using 3 modalities
2. At inference we use only A, V or A+V



Robustness to modality drop

1. Most of information carried by Text

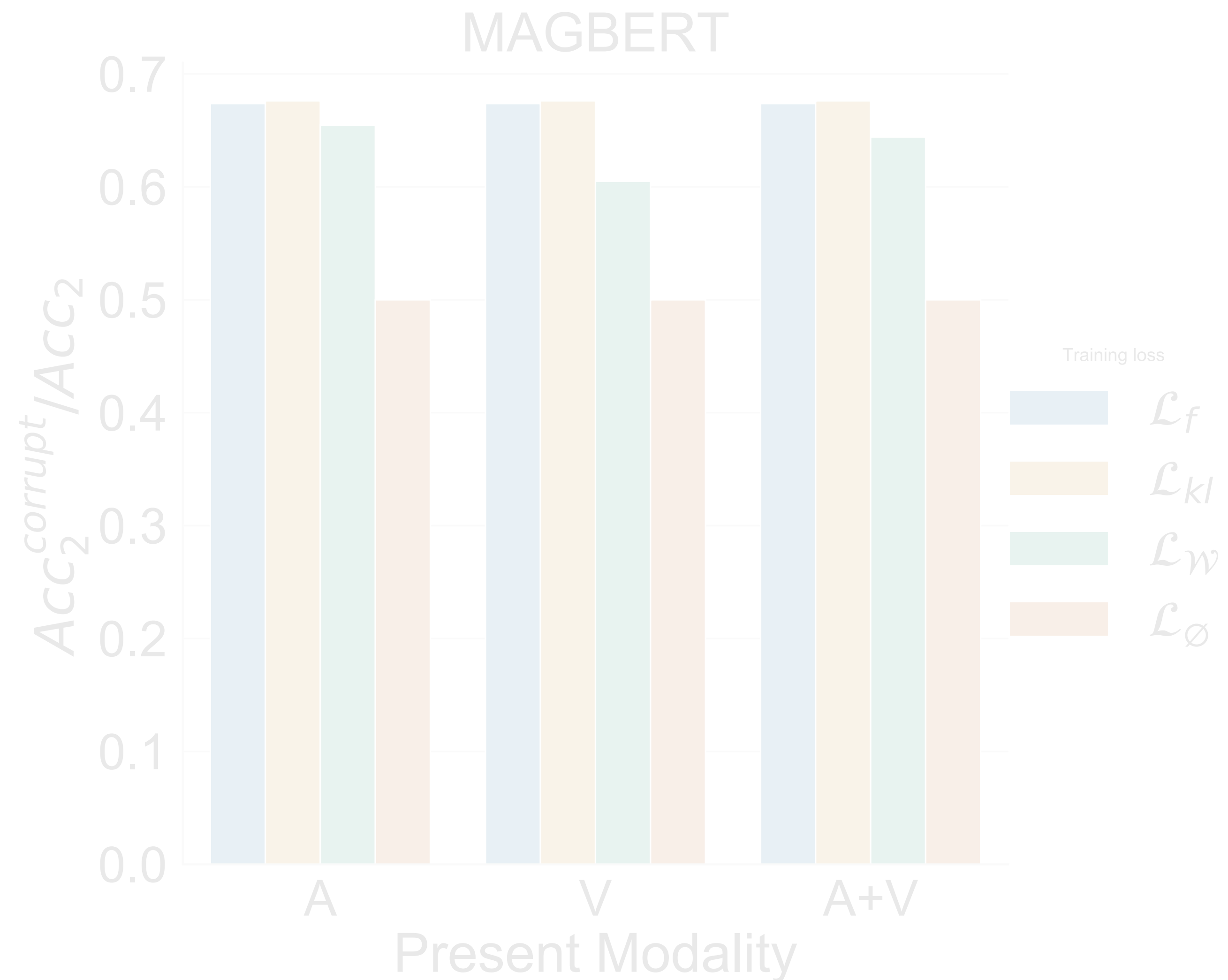
- Drop Text

How maximising the MDM affect robustness?

Experiment

1. Train using 3 modalities
2. At inference we use only A, V or A+V

Our MDM loss improve robustness



Towards explainable representations

MDM

Towards explainable representations

Goal: Use low/high values or MDM to explain representations

MDM

Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L

Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L



Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L
i mean its a Nicholas Sparks book it must be good	disappointed tone + neutral facial expression	L



Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L
i mean its a Nicholas Sparks book it must be good	disappointed tone + neutral facial expression	L



Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L
i mean its a Nicholas Sparks book it must be good	disappointed tone + neutral facial expression	L
the action is fucking awesome	head nod + excited voice	H



Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L
i mean its a Nicholas Sparks book it must be good	disappointed tone + neutral facial expression	L
the action is fucking awesome	head nod + excited voice	H



Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L
i mean its a Nicholas Sparks book it must be good	disappointed tone + neutral facial expression	L
the action is fucking awesome	head nod + excited voice	H
[1] it was cute you know the actors did a great job bringing the smurfs to life such as joe george lopez neil patrick harris katy perry and a fourth	multiple smiles	H



Towards explainable representations

Goal: Use low/high values or *MDM* to explain representations

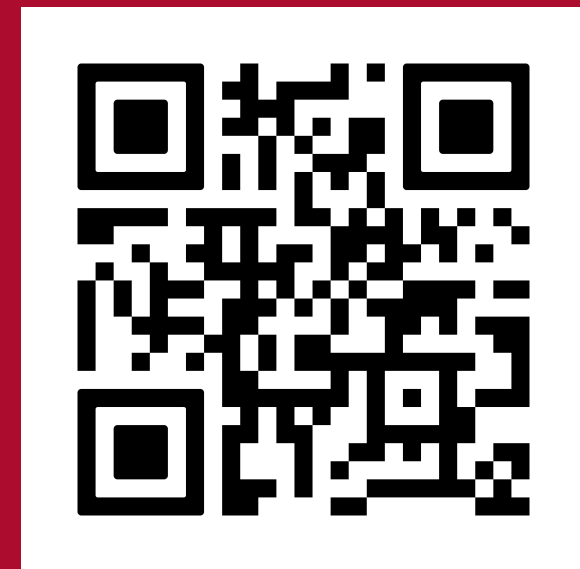
Spoken Transcripts	Acoustic and visual behaviour	<i>MDM</i>
um the story was all right	low energy monotonous voice + headshake	L
i mean its a Nicholas Sparks book it must be good	disappointed tone + neutral facial expression	L
the action is fucking awesome	head nod + excited voice	H
[1] it was cute you know the actors did a great job bringing the smurfs to life such as joe george lopez neil patrick harris katy perry and a fourth	multiple smiles	H



Thanks for listening

**Title: Improving Multimodal Fusion Via Mutual
Dependency Maximisation**

Corresponding Authors:



Pierre Colombo

Link to Paper

