# Efficient Low-rank Multimodal Fusion With Modality-specific Factors
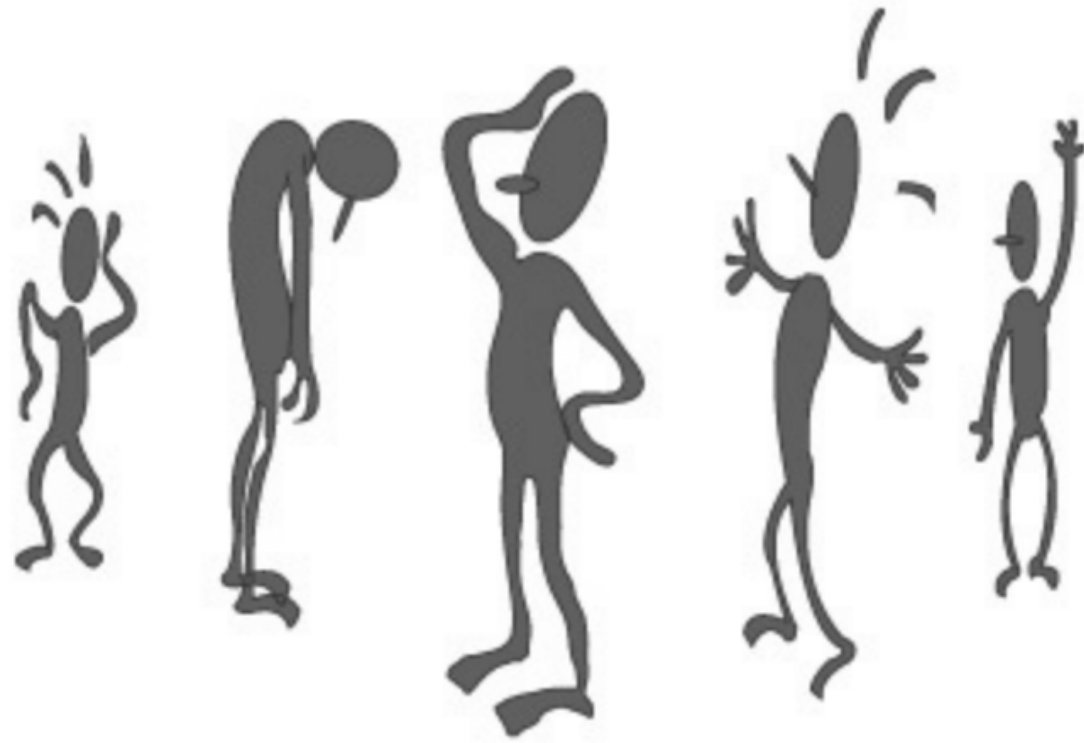
**Zhun Liu, Ying Shen,**

**Varun Bharadwaj,  Paul Pu Liang,**

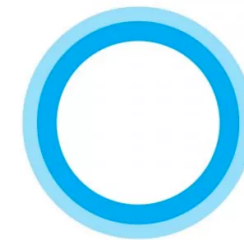**Amir Zadeh, Louis-Philippe Morency**

# Artificial Intelligence



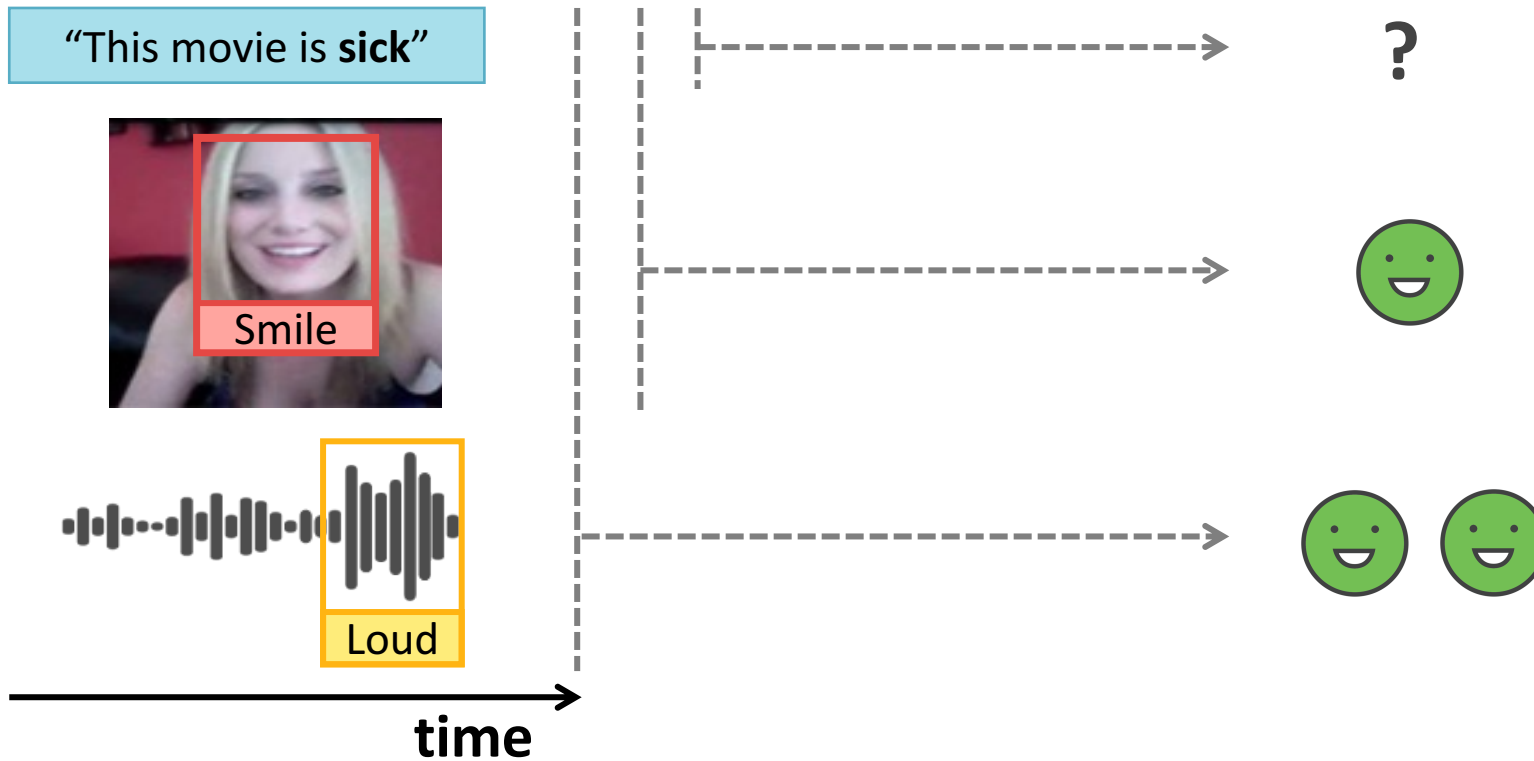Negative •• Neutral •• Positive

Hey Siri

Cortana

# Sentiment and Emotion Analysis
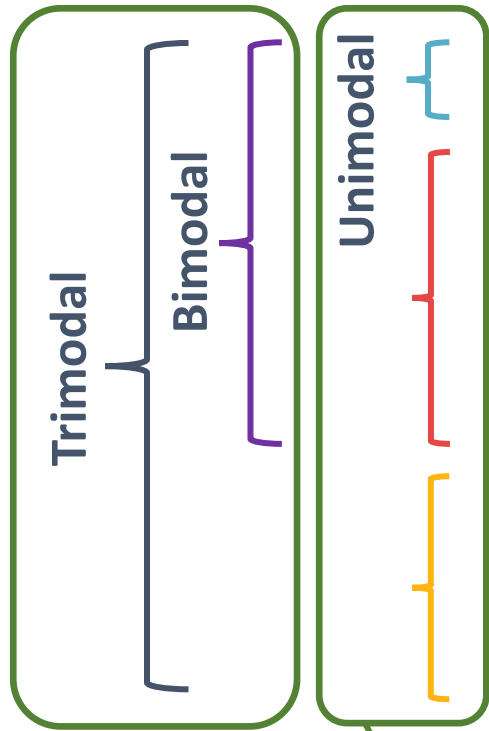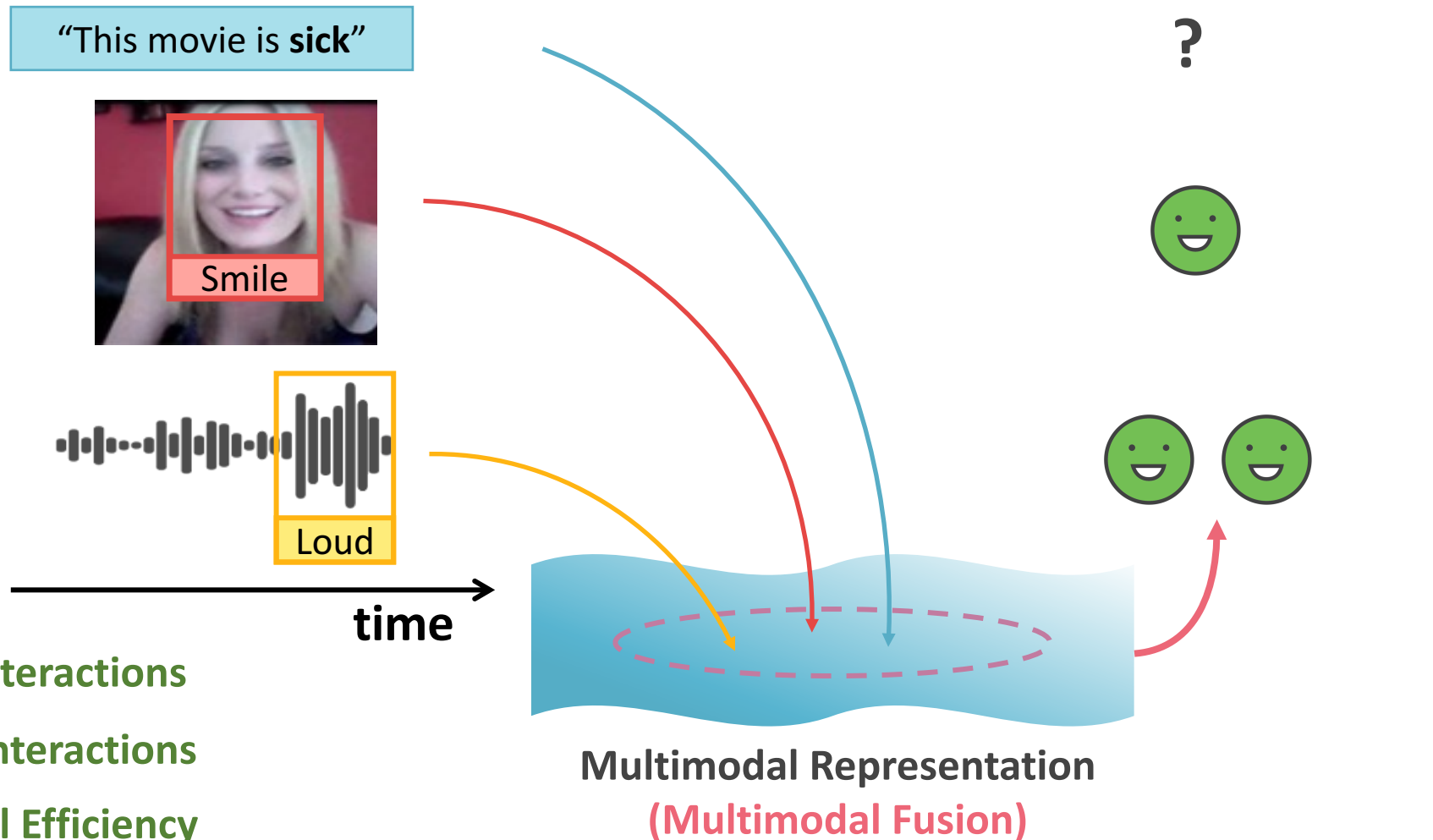
**Speaker's behaviors**

**Sentiment Intensity**

"This movie is **sick**"

?

Smile

Loud

time

# Multimodal Sentiment and Emotion Analysis

# Multimodal Fusion using Tensor Representation



Bimodal

$|h|$

**Visual**

$z_v$

$\mathcal{Z}$ · $\mathcal{W}$

$h$

**Language**

$z_l$

"This movie is sick"

Unimodal

**Multimodal Representation**

☑ **Intra-modal interactions**

☑ **Cross-modal interactions**

☒ **Computational efficiency**

$$\mathcal{Z} = \begin{bmatrix} z_v \\ 1 \end{bmatrix} \otimes \begin{bmatrix} z_l \\ 1 \end{bmatrix} = \begin{bmatrix} z_v & z_v \otimes z_l \\ 1 & z_l \end{bmatrix}$$
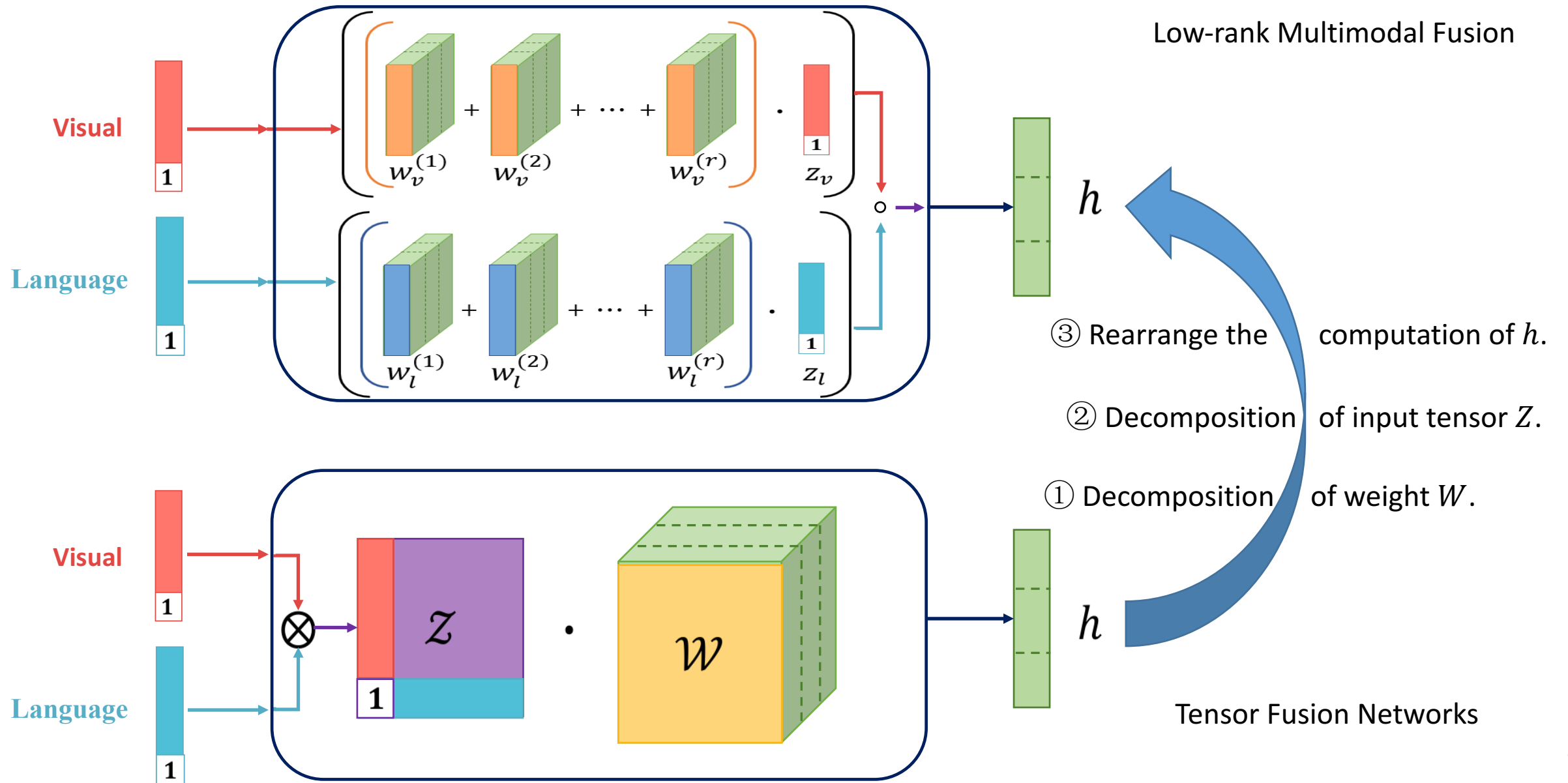
"Tensor Fusion Network for Multimodal Sentiment Analysis" by Zadeh, A., et, al. (2017)

# Computational Complexity – Tensor Product

$O\left(\prod_{m=1}^{M} d_m\right)$

$\mathcal{Z}$

$O(d_1 \times d_2 \times d_3)$

$\mathcal{Z}$

$O(d_1 \times d_2)$

M=2     M=3

Computational Complexity

Number of Modalities

# CORE CONTRIBUTIONS

Low-rank Multimodal Fusion (LMF)
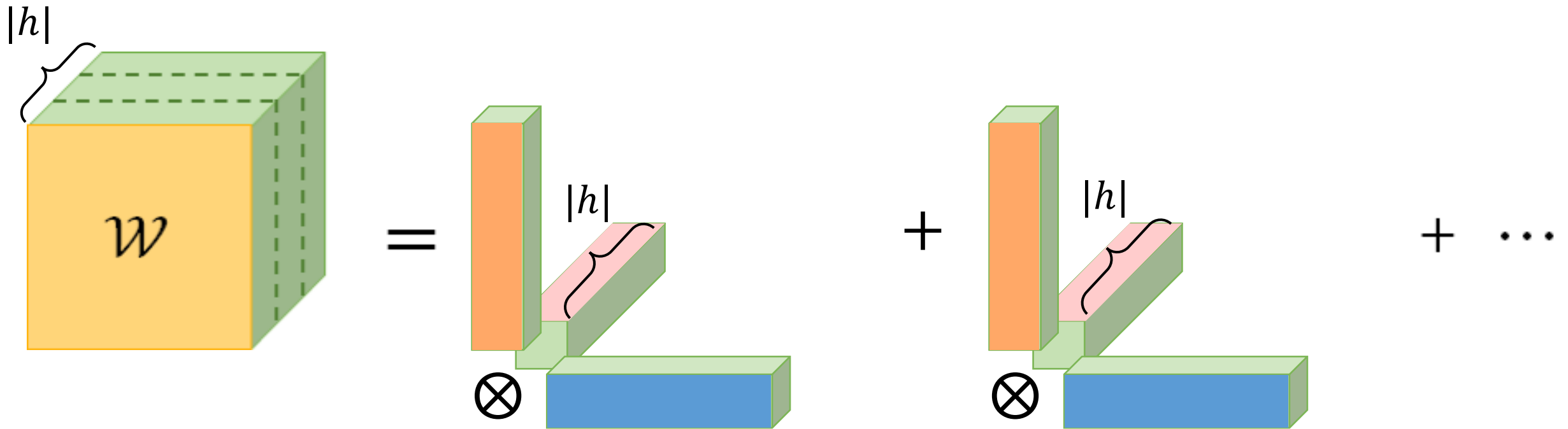
# From Tensor Representation to Low-rank Fusion



Low-rank Multimodal Fusion

③ Rearrange the computation of $h$.

② Decomposition of input tensor $Z$.

① Decomposition of weight $W$.

Tensor Fusion Networks

# Canonical Polyadic (CP) Decomposition of tensors

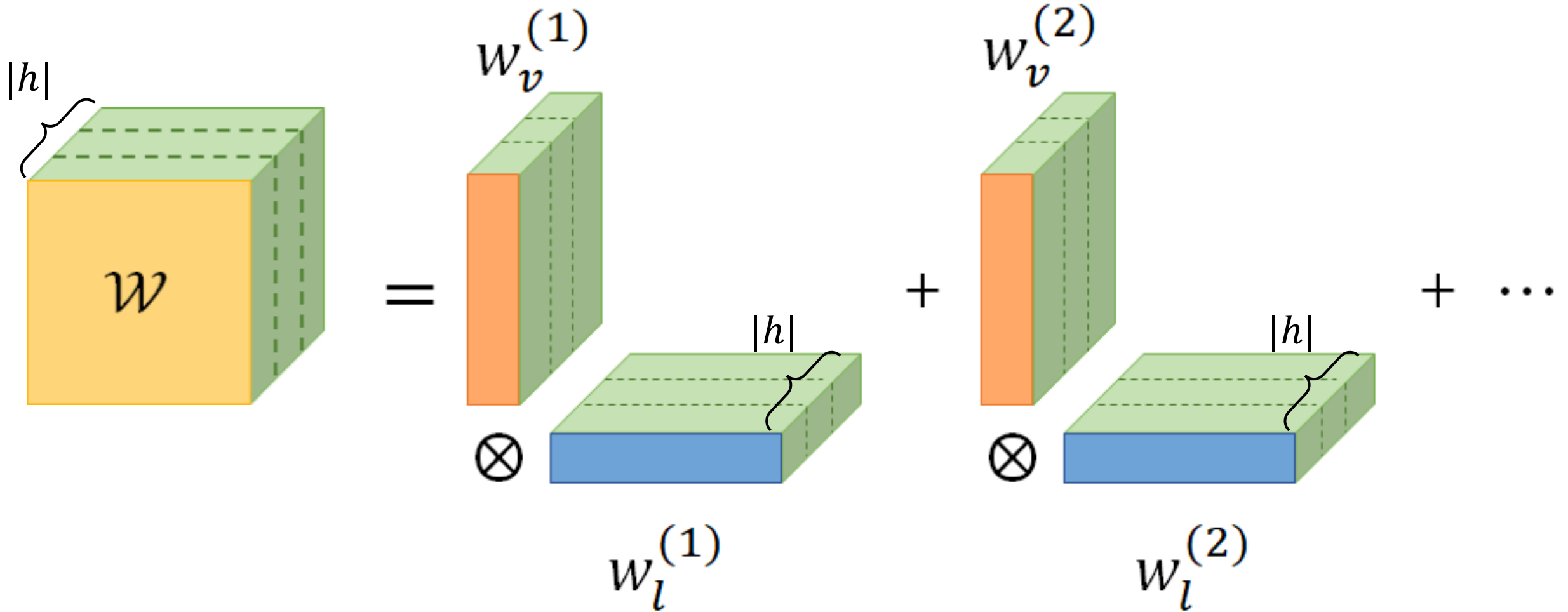$$\mathcal{W} = W_v^{(1)} \otimes W_l^{(1)} + W_v^{(2)} \otimes W_l^{(2)} + \cdots$$

Rank of tensor $W$: minimum number of vector tuples needed for exact reconstruction

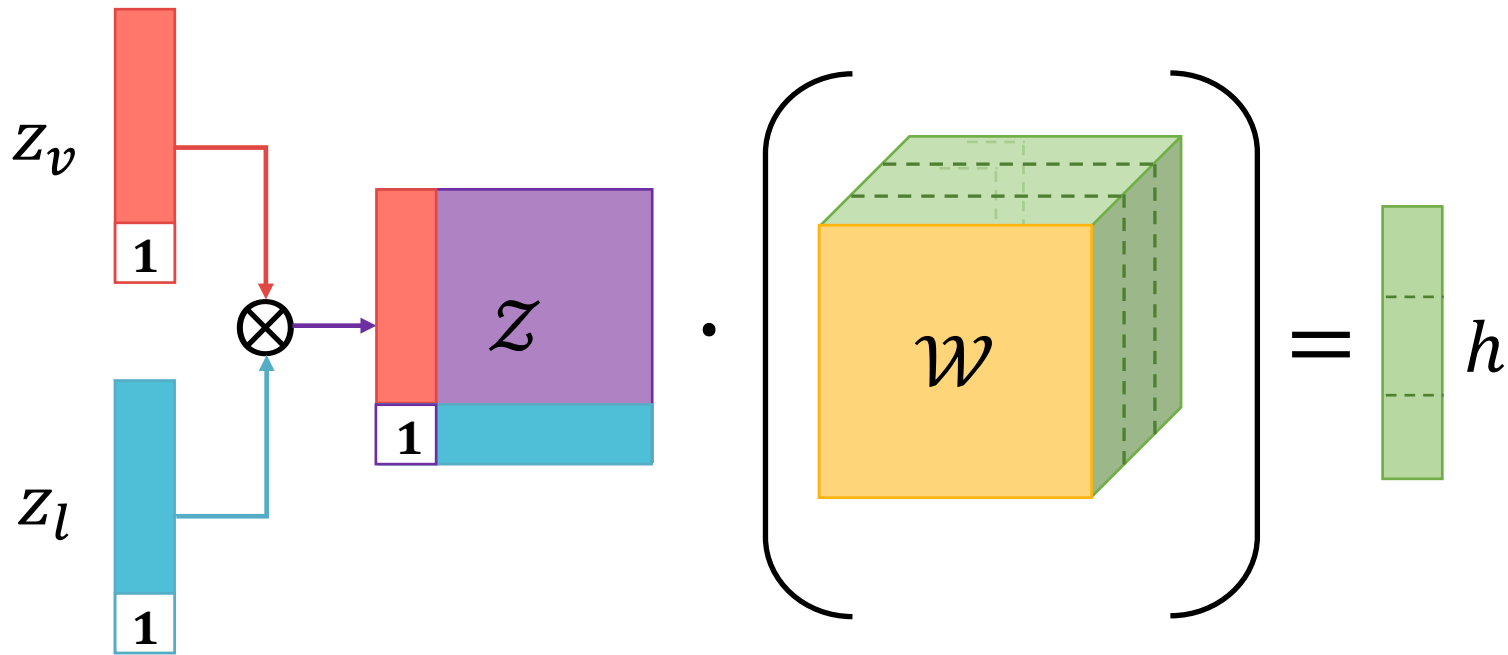# Canonical Polyadic (CP) Decomposition of 3D tensors
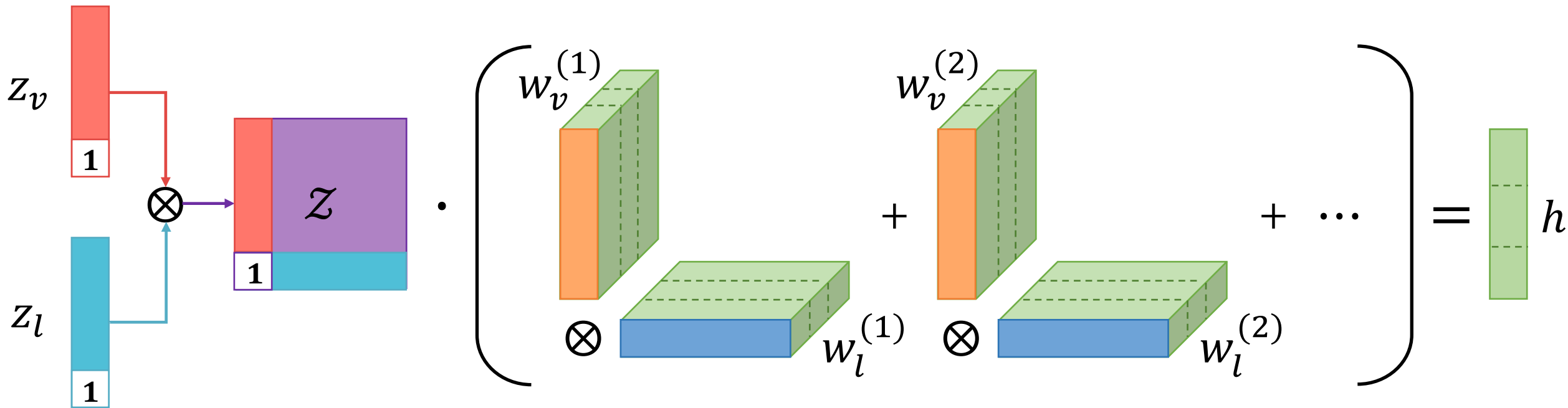
# Modality-specific Decomposition



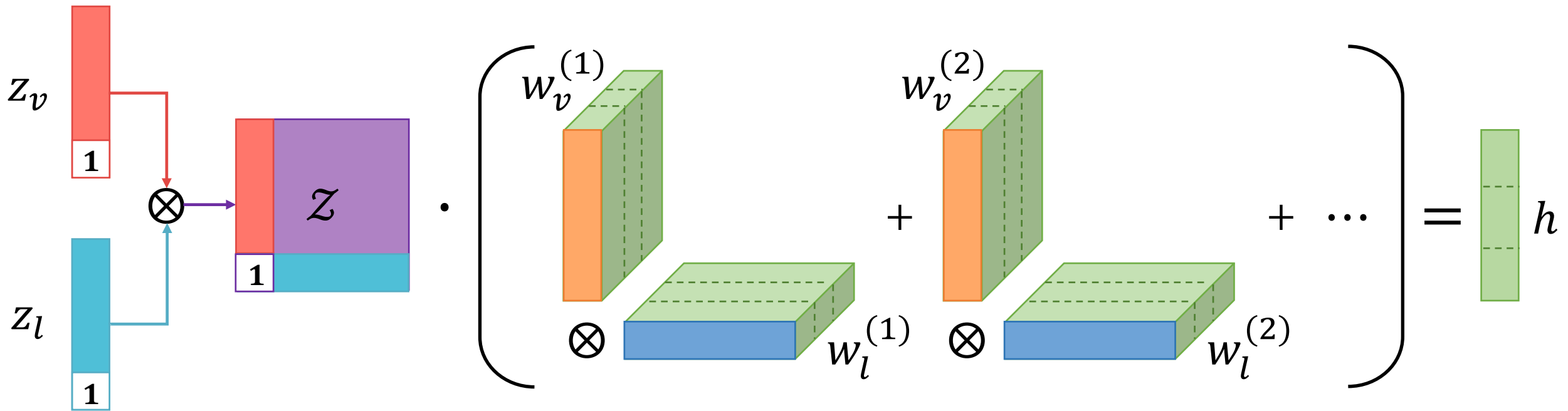Retain the dimension for the multimodal representation $h$ during decomposition

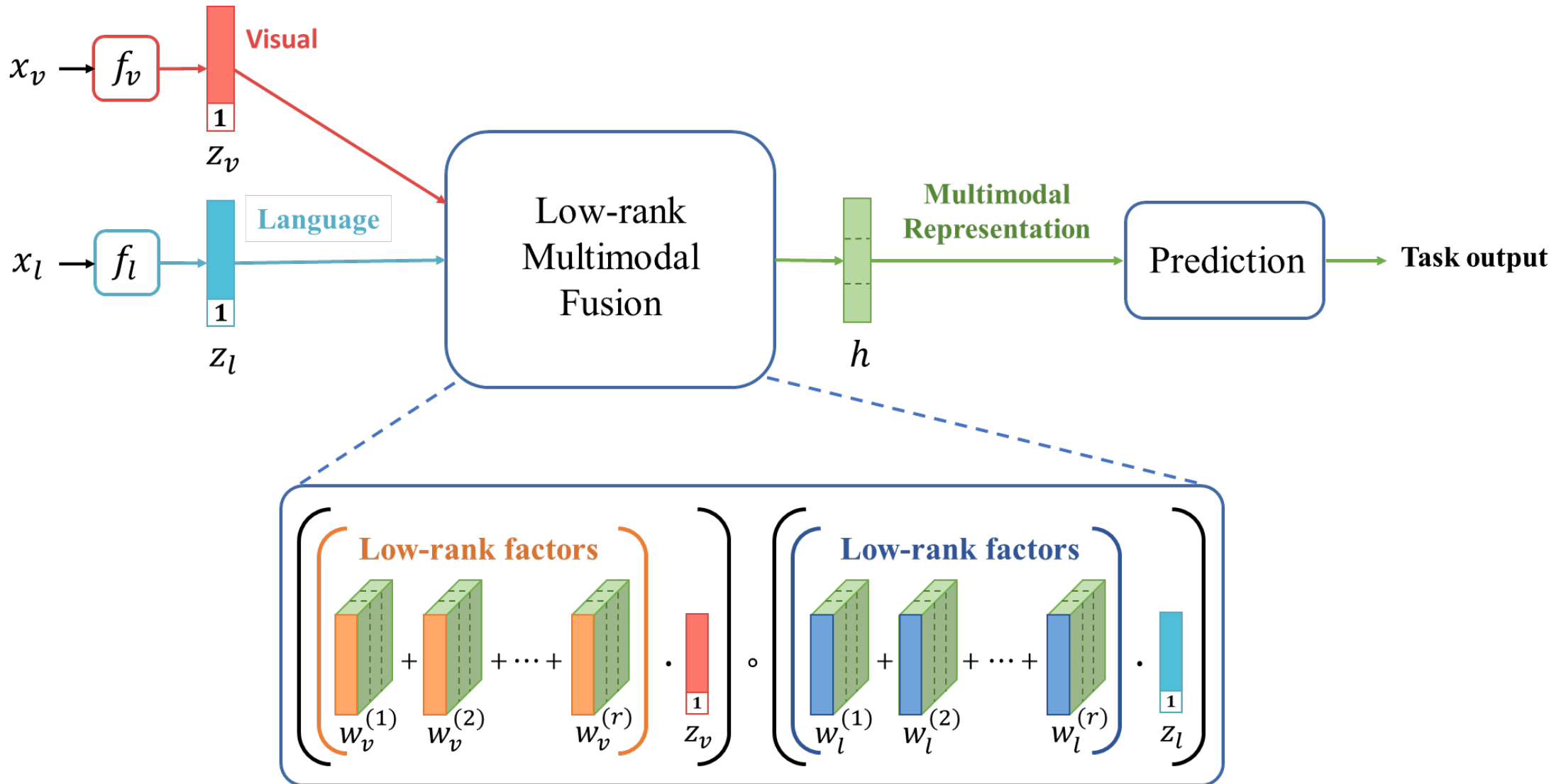# ① Decomposition of weight tensor W

# ① Decomposition of weight tensor W

$$z_v \otimes z_l \rightarrow \mathcal{Z} \cdot \left( w_v^{(1)} \otimes w_l^{(1)} + w_v^{(2)} \otimes w_l^{(2)} + \cdots \right) = h$$

# ③ Rearranging computation



$$\left(\left(w_v^{(1)} + w_v^{(2)} + \cdots + w_v^{(r)}\right) \cdot z_v\right) \circ \left(\left(w_l^{(1)} + w_l^{(2)} + \cdots + w_l^{(r)}\right) \cdot z_l\right) = h$$

# Low-rank Multimodal Fusion

# Easily scales to more modalities

# EXPERIMENTS AND RESULTS

# Datasets

## CMU-MOSI
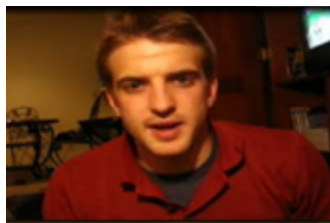


Sentiment Analysis

2199 video segments
- Single-speaker
- From 93 Movie reviews

Segment level annotations
- Sentiment
- Real-valued

## POM



Speaker Trait Recognition

1000 full video clips
- Single-speaker
- Movie reviews

Video level annotations
- 16 types of speaker traits
- Categorical annotations

## IEMOCAP



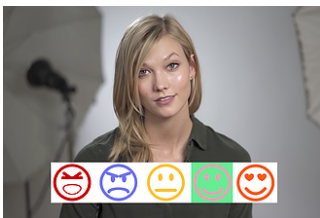Emotion Recognition

10039 video segments
- Dyadic interaction
- From 302 videos

Segment level annotations
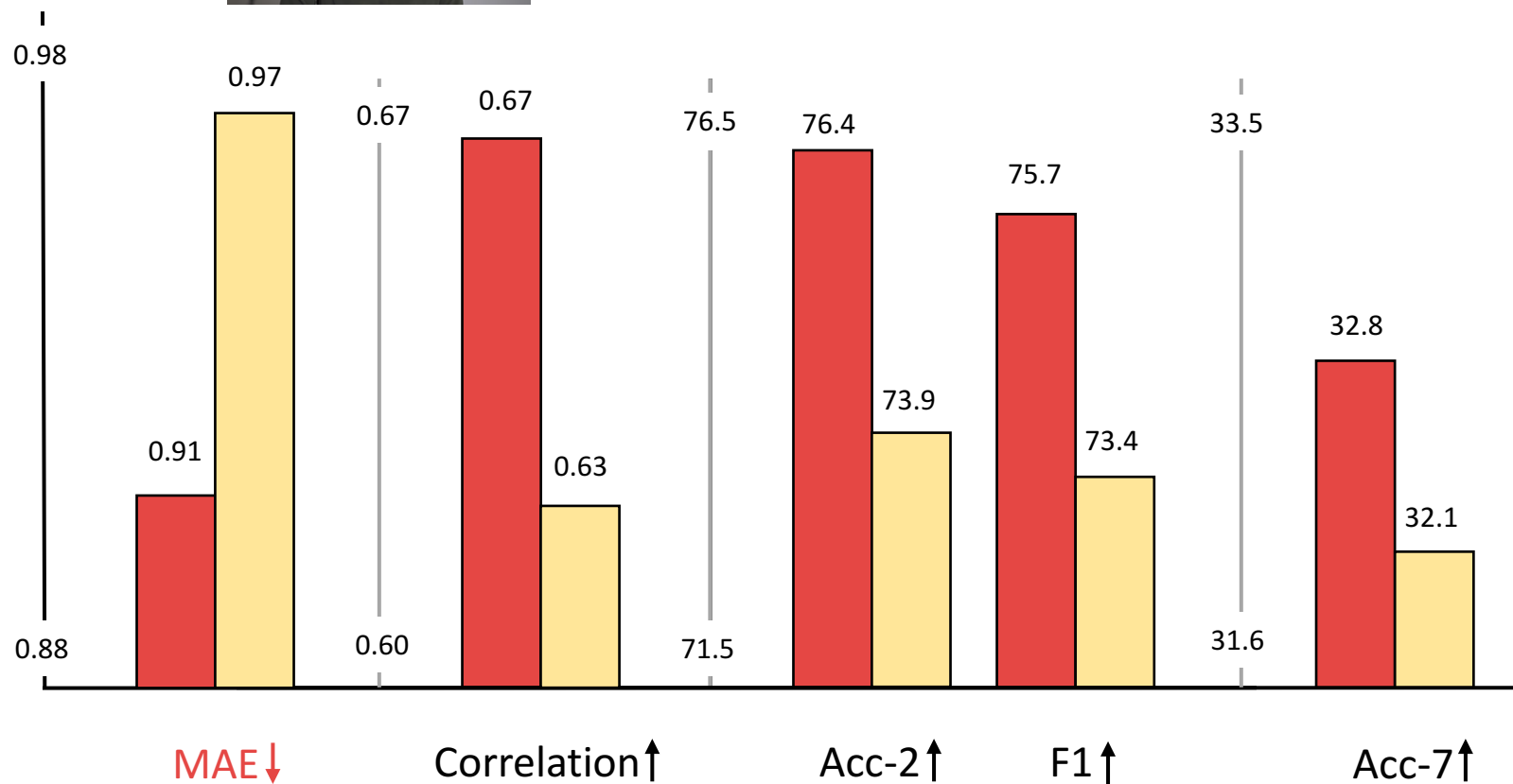- 10 classes of emotions
- Categorical annotations

# Compare to full rank tensor fusion

CMU-MOSI



Low-rank Multimodal Fusion (Our Model) — ■ LMF

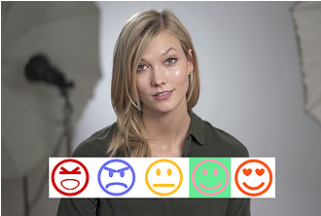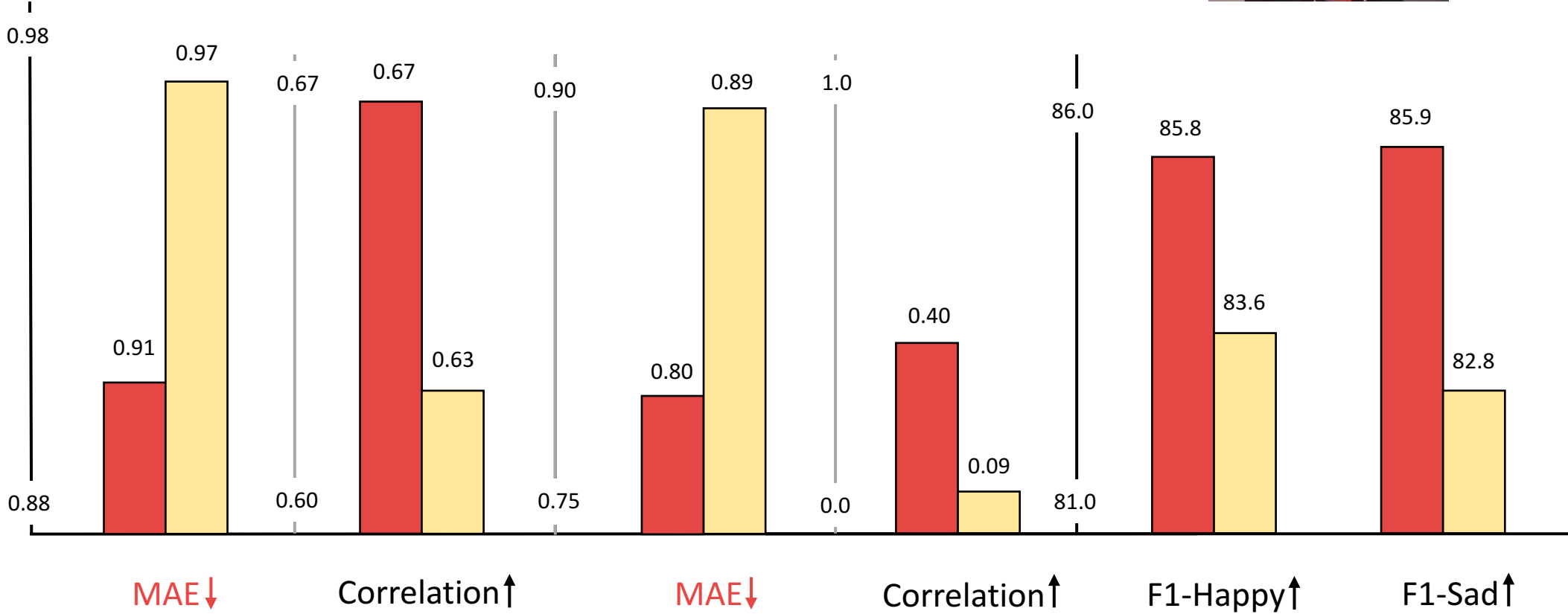Tensor Fusion Networks (Zadeh, et al., 2017) — ▢ TFN



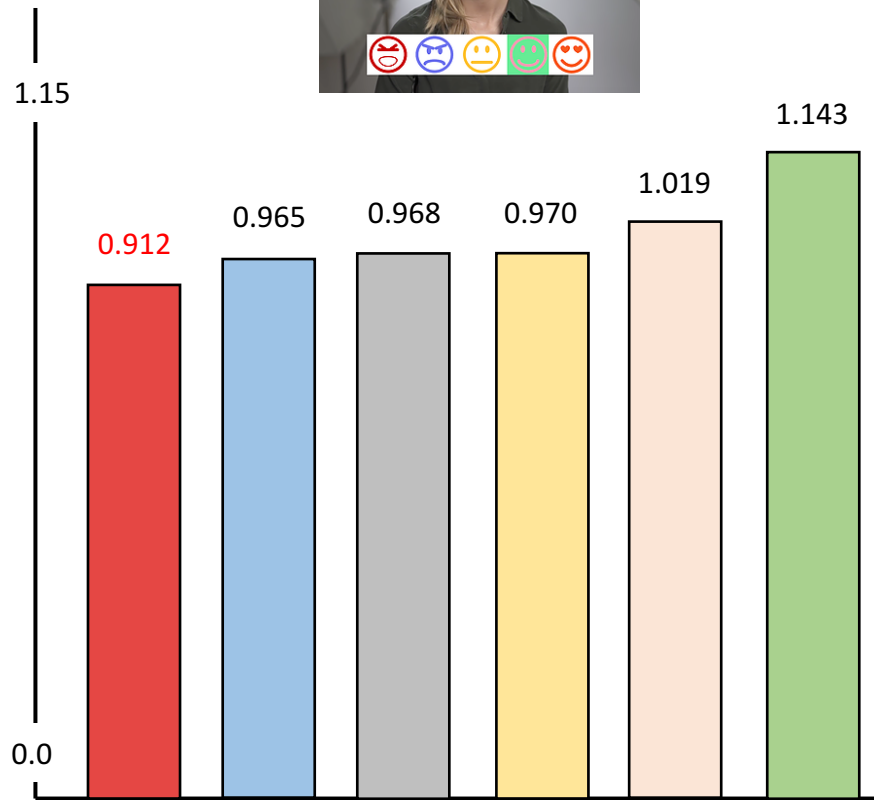| | MAE↓ | Correlation↑ | Acc-2↑ | F1↑ | Acc-7↑ |
|---|---|---|---|---|---|
| LMF | 0.91 | 0.67 | 76.4 | 75.7 | 32.8 |
| TFN | 0.97 | 0.63 | 73.9 | 73.4 | 32.1 |

# Compare to full rank tensor fusion

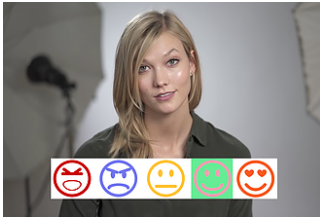# Compare with State-of-the-Art Approaches
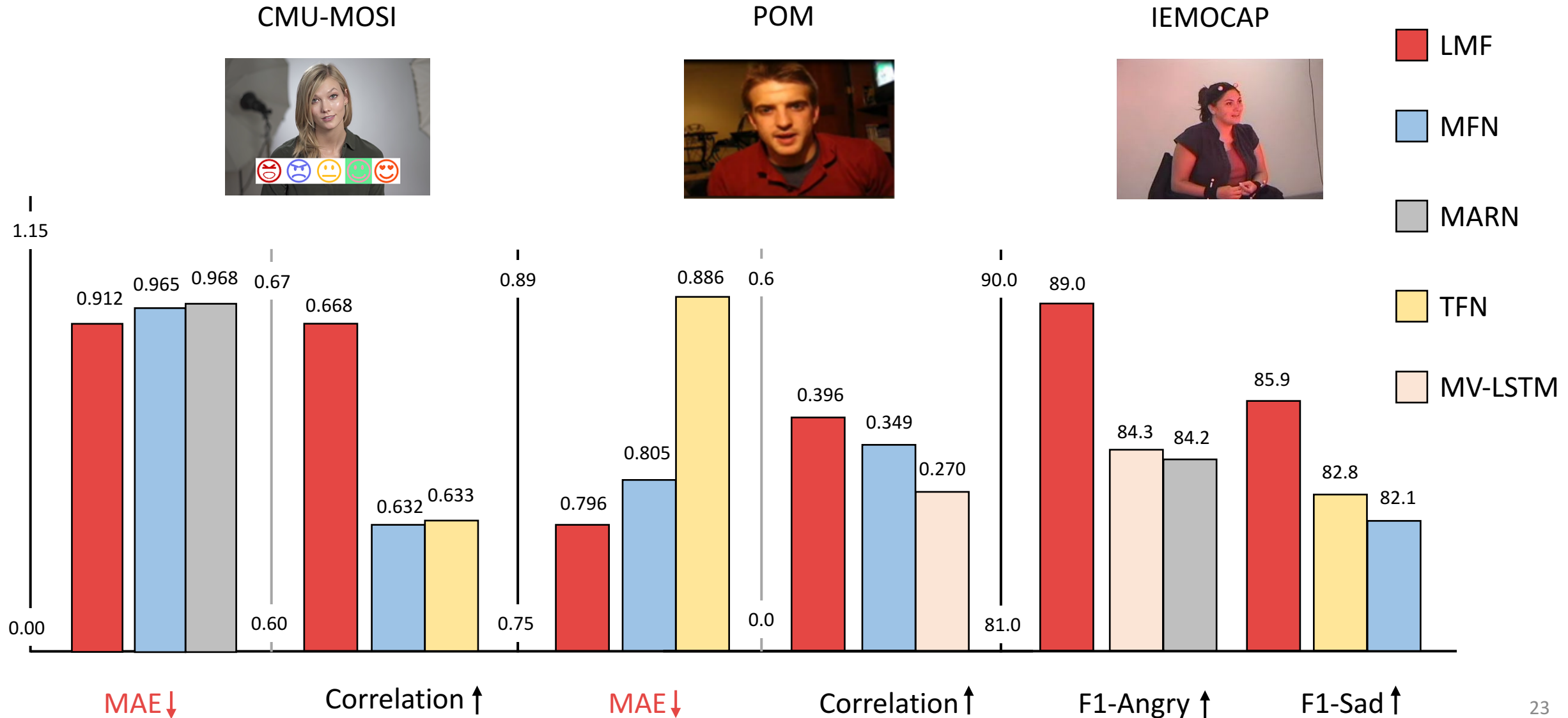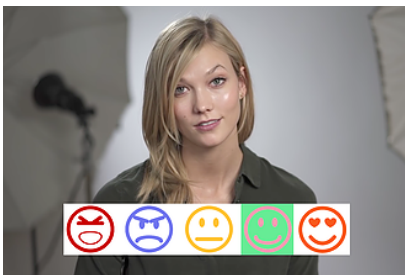
CMU-MOSI



Mean Average Error (MAE)

Low-rank Multimodal Fusion
(our model) — LMF

Memory Fusion Networks
(Zadeh, et al., 2018) — MFN

Multi-attention Recurrent Networks
(Zadeh, et al., 2018) — MARN

Tensor Fusion Networks
(Zadeh, et al., 2017) — TFN

Multi-view LSTM
(Rajagopalan, et al., 2016) — MV-LSTM

Deep Fusion
(Nojavanasghari, et al., 2016) — Deep Fusion

# Compare with Top 2 State-of-the-Art Approaches



CMU-MOSI        POM        IEMOCAP

Legend: LMF, MFN, MARN, TFN, MV-LSTM

CMU-MOSI:
- MAE↓: 0.912, 0.965, 0.968
- Correlation↑: 0.668, 0.632, 0.633

POM:
- MAE↓: 0.796, 0.805, 0.886
- Correlation↑: 0.396, 0.349, 0.270

IEMOCAP:
- F1-Angry↑: 89.0, 84.3, 84.2
- F1-Sad↑: 85.9, 82.8, 82.1

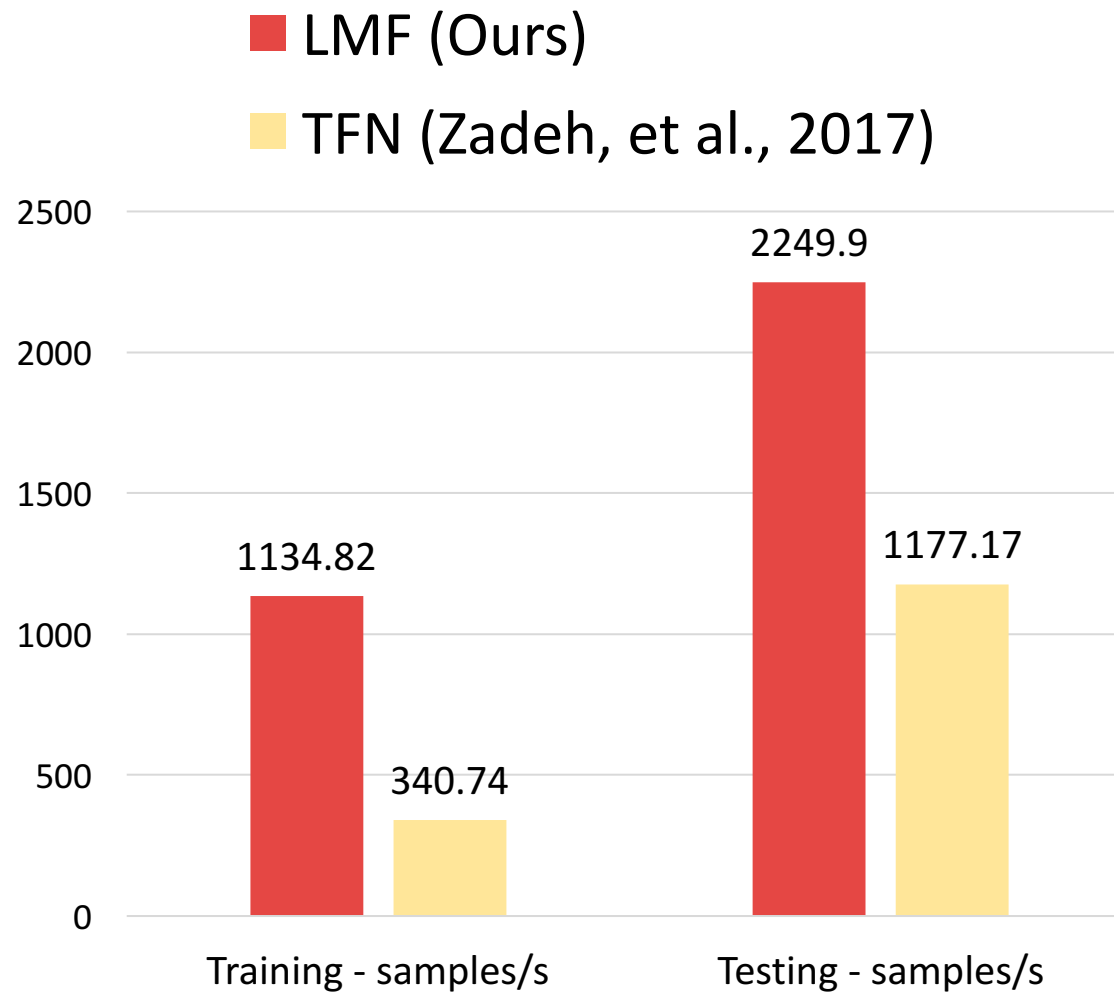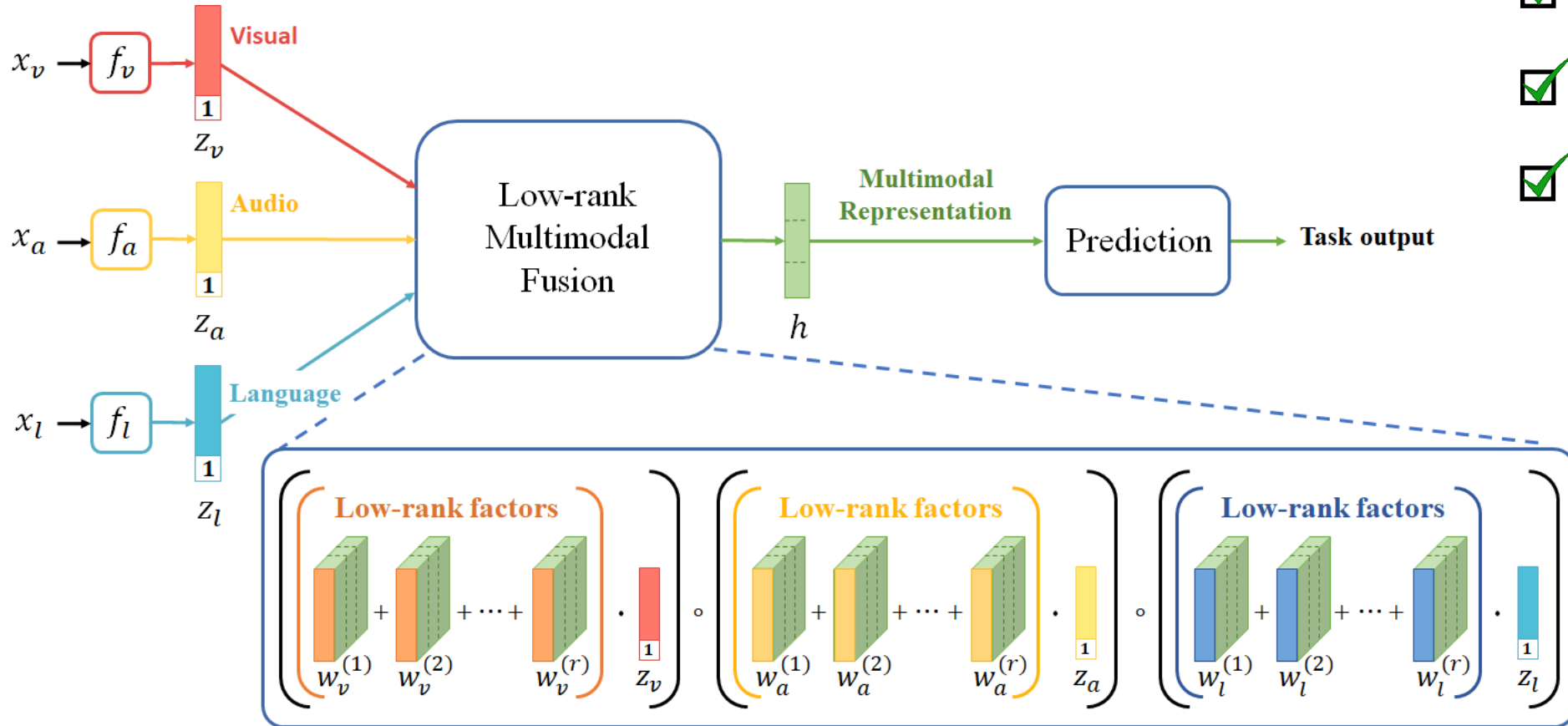Axis values: 1.15, 0.00, 0.67, 0.60, 0.89, 0.75, 0.6, 0.0, 90.0, 81.0

23

# Efficiency Improvement

CMU-MOSI



**Efficiency Metric:** Number of data samples processed per second

- Training Efficiency
- Testing Efficiency



Legend:
- LMF (Ours)
- TFN (Zadeh, et al., 2017)

Chart values:
- Training - samples/s: LMF 1134.82, TFN 340.74
- Testing - samples/s: LMF 2249.9, TFN 1177.17

# Conclusions



- ☑ Intra-modal interactions
- ☑ Cross-modal interactions
- ☑ Computational complexity
- ☑ State-of-the-art results

# Thank you!

**Code:** https://github.com/Justin1904/Low-rank-Multimodal-Fusion

MultiComp Lab

http://multicomp.cs.cmu.edu/