

How to AI (Almost) Anything

Lecture 5 – Multimodal Fusion

Paul Liang

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

ppliang@mit.edu

 [@pliang279](#)



Assignments for This Coming Week

For project:

- I gave feedback and assigned primary TA.
- Meet with me and primary TA every other week.
- Should have finalized main ideas and experimental setup, have baseline models working, progress towards implementing new ideas.

Reading assignment due tomorrow Wednesday (3/12).

This Thursday (3/13): third reading discussion on **multimodal alignment**.

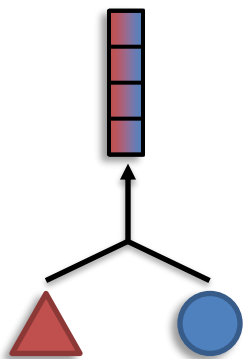
What views for contrastive learning

Platonic representation hypothesis

Today's lecture

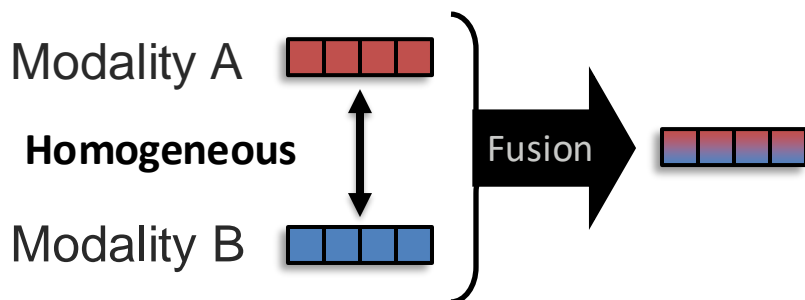
- 1 Basics of multimodal fusion
- 2 Early, intermediate, late fusion
- 3 Multiplicative and dynamic fusion
- 4 Complex fusion and improving optimization

Sub-Challenge 1a: Representation Fusion

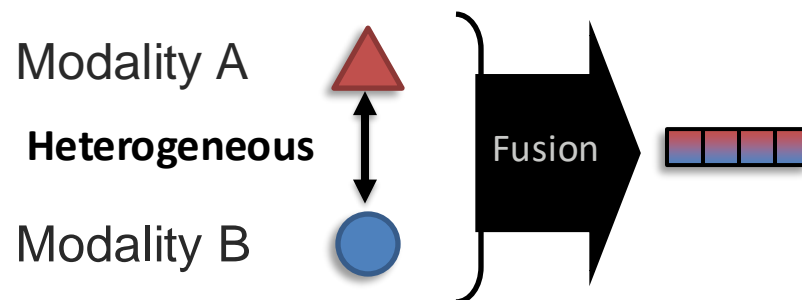


Definition: Learn a joint representation that models cross-modal interactions between individual elements of different modalities.

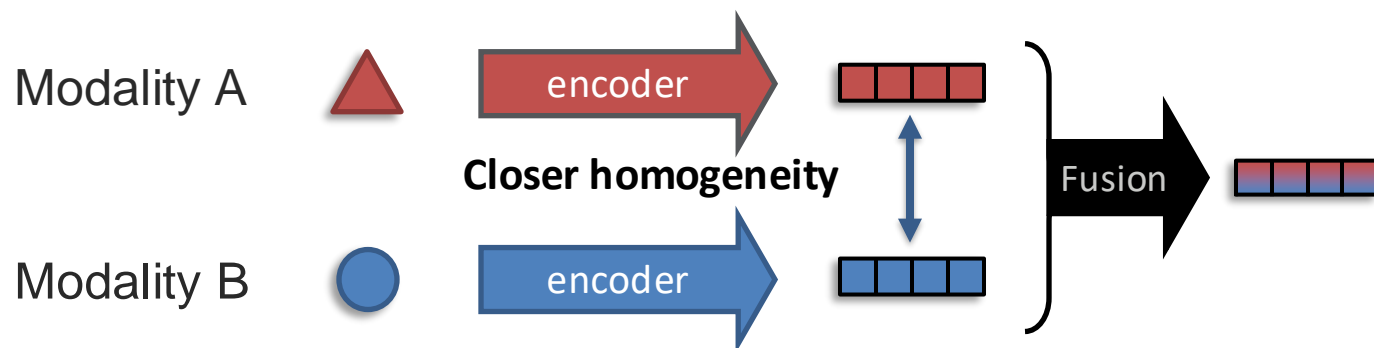
Fusion with abstract modalities:



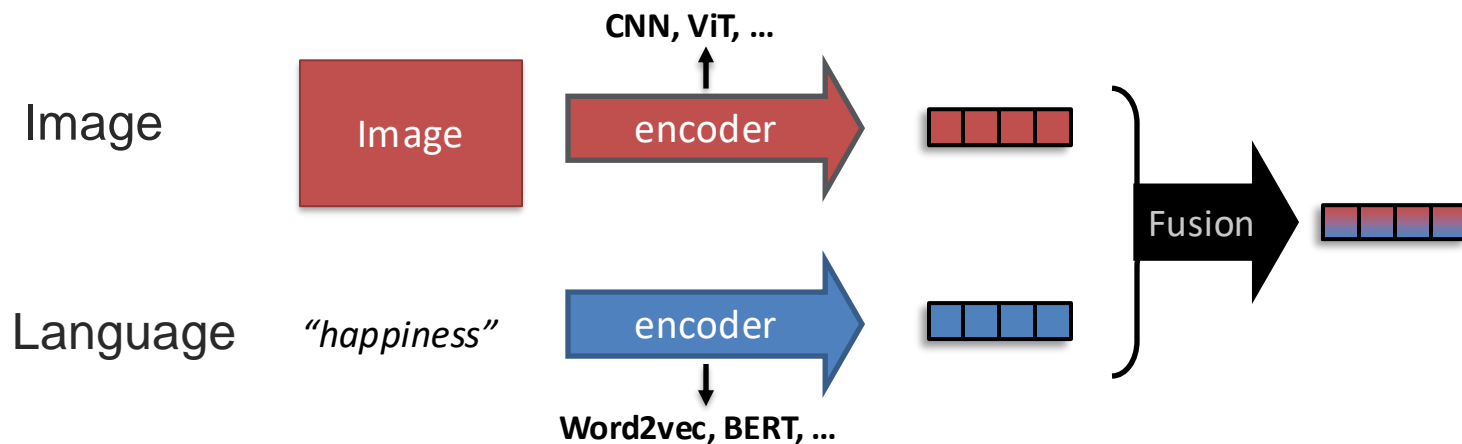
Fusion with raw modalities:



Fusion with Abstract Modalities



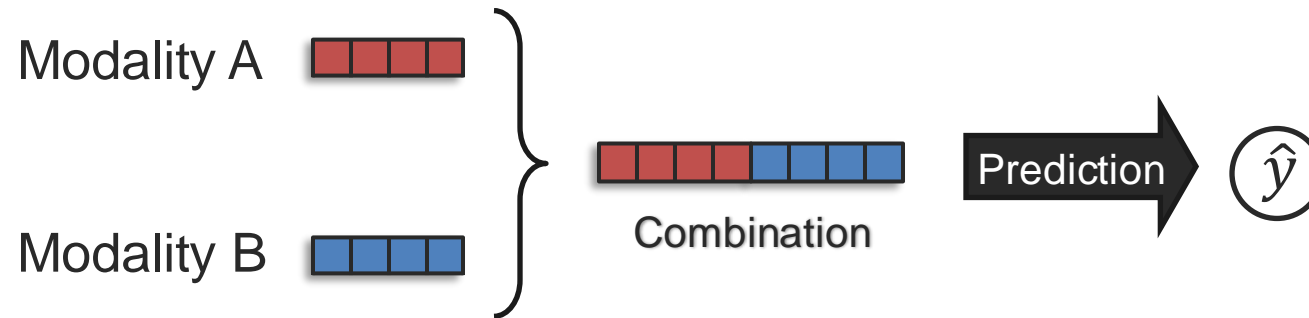
Example:



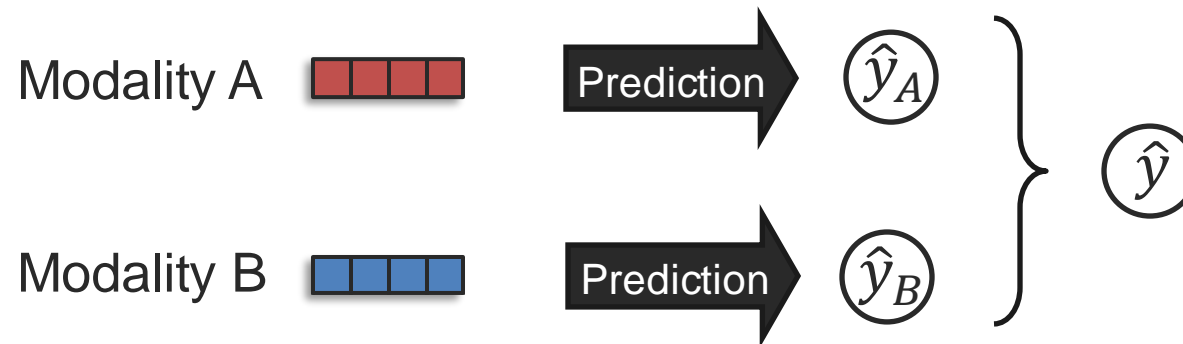
➡ Unimodal encoders can be jointly learned with fusion network, or pre-trained

Early and Late Fusion

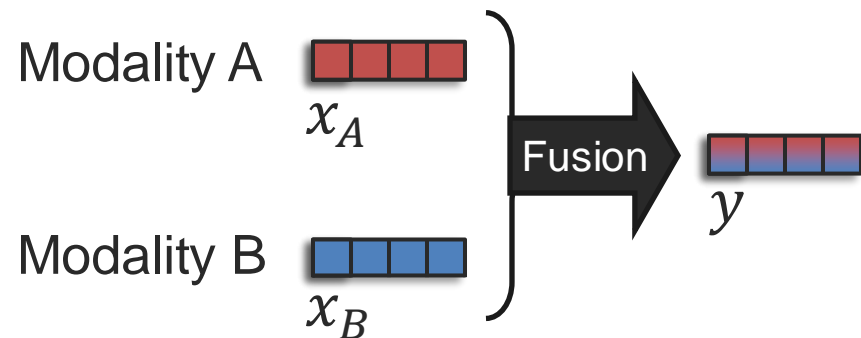
Early fusion:



Late fusion:



Basic Concepts for Fusion



Goal: Model *cross-modal interactions* between the multimodal elements

→ Let's study the **univariate case first**
 ↳ (only 1-dimensional features)

Linear regression:

$$y = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

↓
 intercept
 (bias term)

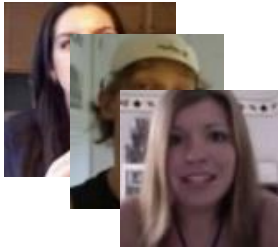
Additive
 terms

Multiplicative
 term

error
 (residual term)

Linear Fusion Case

300 book reviews



y : audience score

x_A : percentage of smiling

x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

$$y = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

Diagram illustrating the components of the linear regression equation:

- w_0 : intercept (bias term)
- $w_1 x_A + w_2 x_B$: Additive terms
- $w_3 (x_A \times x_B)$: Multiplicative term
- ϵ : error (residual term)

w_0 : average score when x_A and x_B are zero

w_1 : effect from x_A variable only

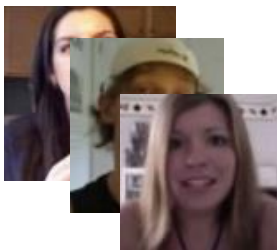
w_2 : effect from x_B variable only

w_3 : effect from x_A and x_B interaction only

ϵ : residual not modeled by w_0 , w_1 , w_2 or w_3

Linear Fusion Case

300 book reviews



y : audience score

x_A : percentage of smiling

x_B : professional status
(0=non-critic, 1=critic)

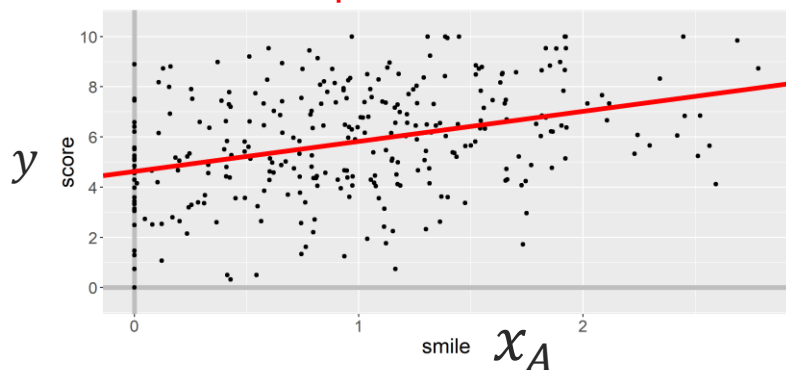
H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

$$y = w_0 + \boxed{w_1} x_A + \epsilon$$

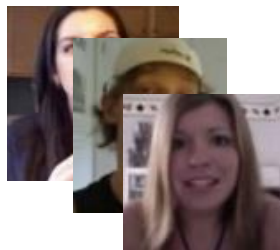
slope



	Estimate
w_0	4.63
w_1	1.20

Linear Fusion Case

300 book reviews



y : audience score

x_A : percentage of smiling

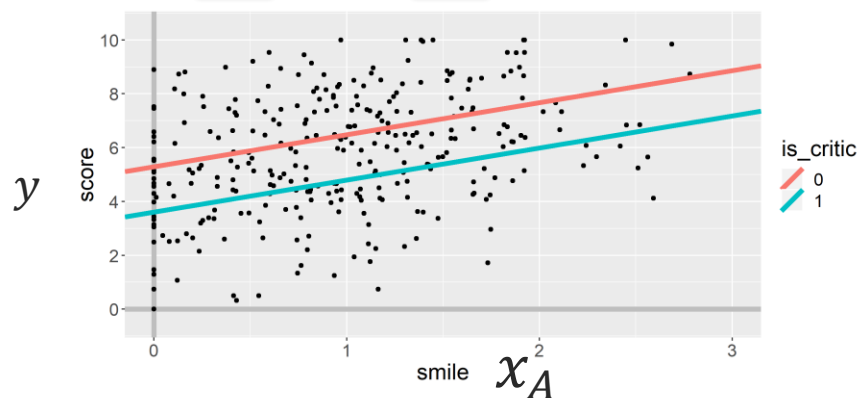
x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

$$y = w_0 + w_1 x_A + w_2 x_B + \epsilon$$



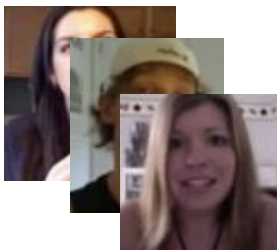
	Estimate
w_0	5.29
w_1	1.19
w_2	-1.69

➡ Positive effect

➡ Negative effect

Linear Fusion Case

300 book reviews



y : audience score

x_A : percentage of smiling

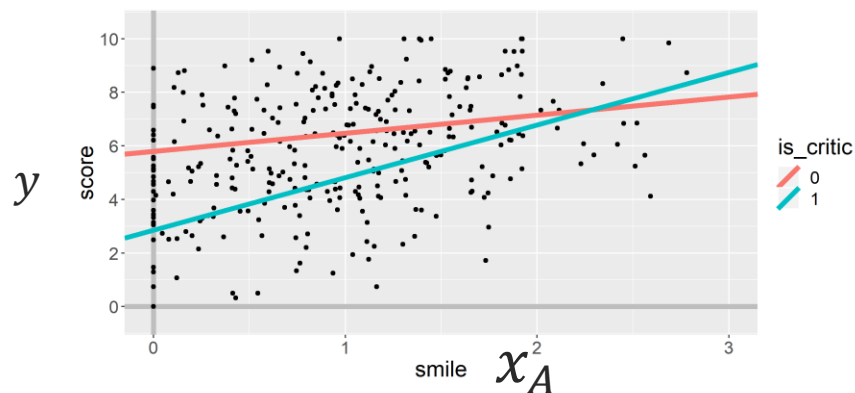
x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

H2: Does the effect of smiling depend on professional status?

Linear regression:

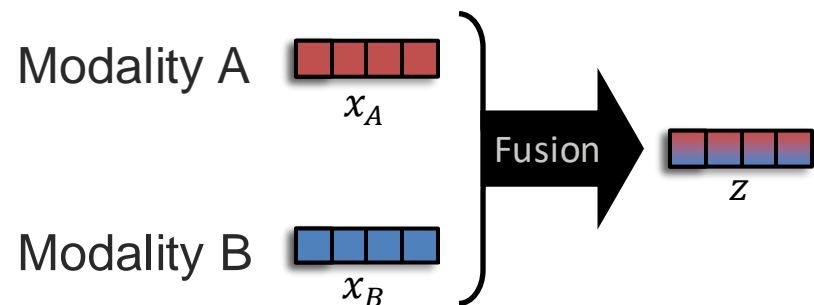
$$y = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$



	Estimate
w_0	5.79
w_1	0.68
w_2	-2.94
w_3	1.29

➡ **Multiplicative interaction!**

Basic Concepts for Representation Fusion



Goal: Model *cross-modal interactions* between the multimodal elements

➡ Let's study the univariate case first
 ↳ (only 1-dimensional features)

Linear regression:

$$z = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

constant Additive terms Multiplicative term error

① Additive interaction:

$$z = w_1 x_A + w_2 x_B + \epsilon$$

② Multiplicative interaction:

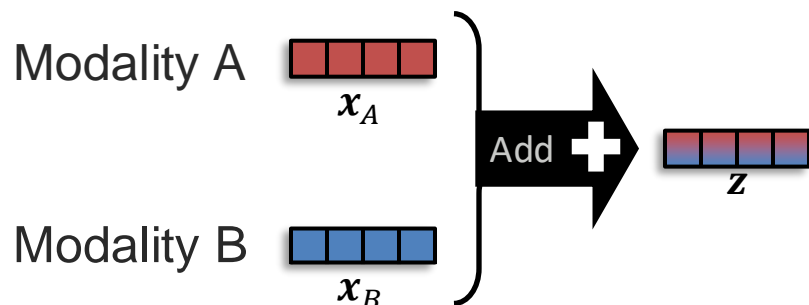
$$z = w_3 (x_A \times x_B) + \epsilon$$

③ Additive and multiplicative interactions:

$$z = w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

Additive Fusion Back to multivariate case!

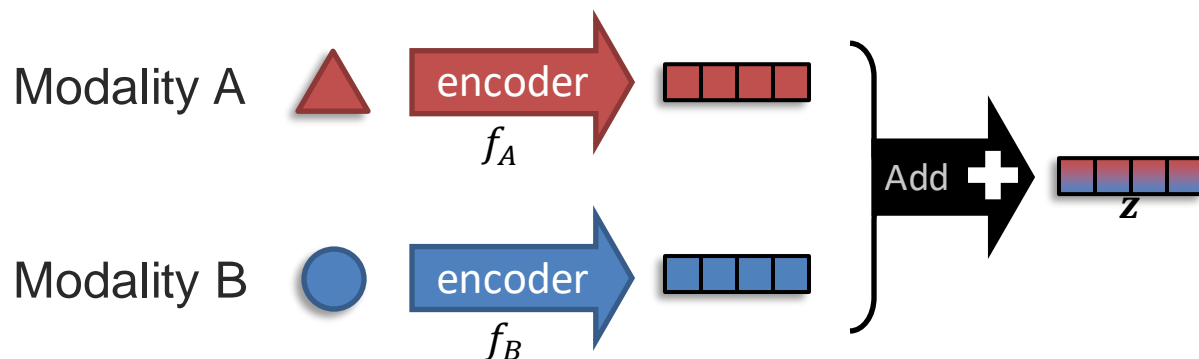
 (multi-dimensional features)



Additive fusion:


$$z = w_1 x_A + w_2 x_B$$

With unimodal encoders:

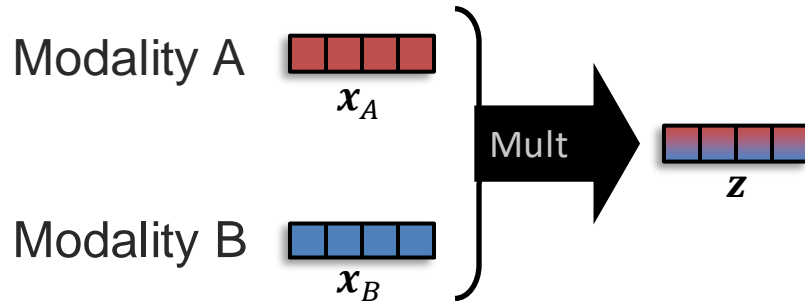


Additive fusion:

$$z = f_A(\triangle) + f_B(\circ)$$

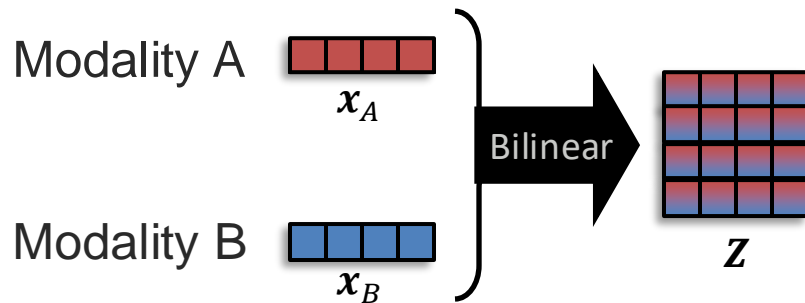
 It could be seen as an ensemble approach
(late fusion)

Multiplicative Fusion



Multiplicative fusion:

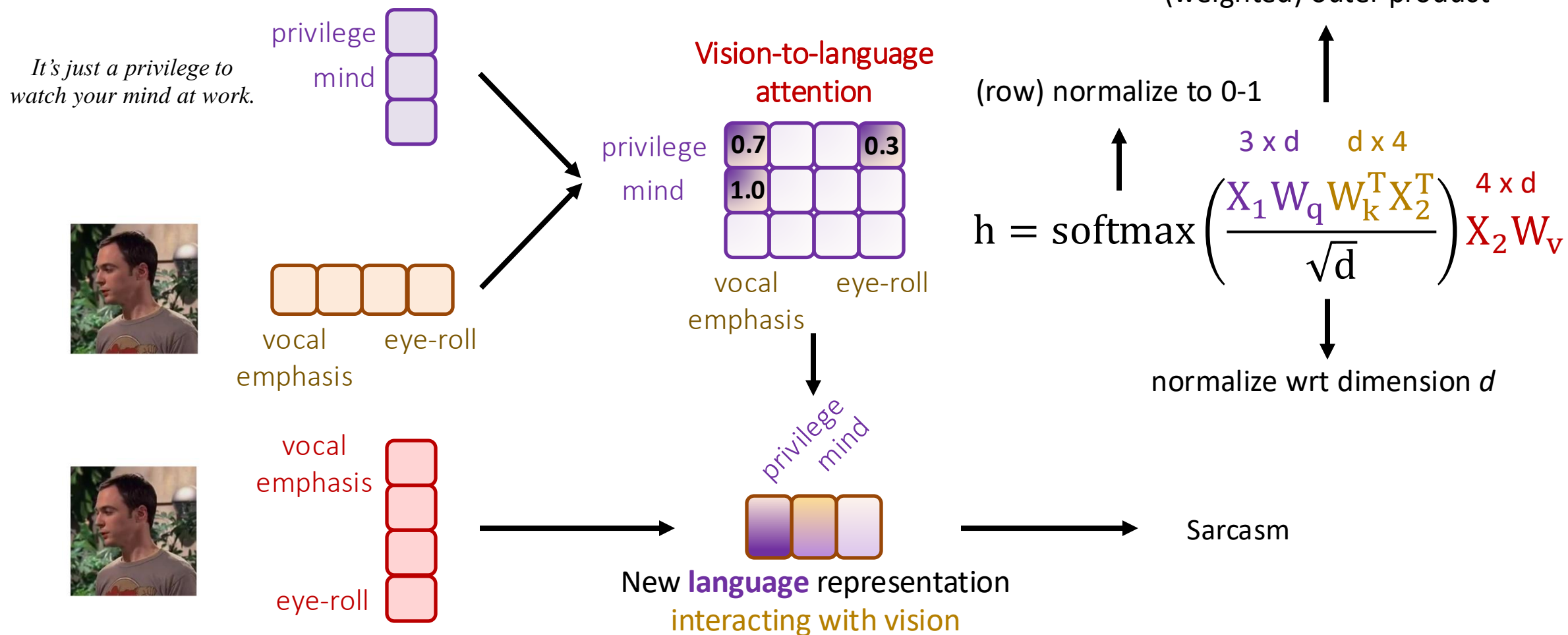
$$z = w(x_A \times x_B)$$



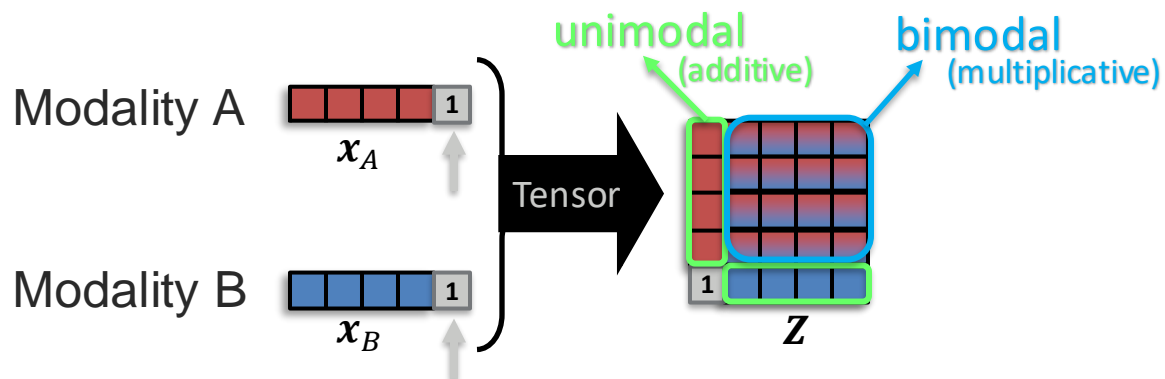
Bilinear Fusion:

$$Z = w(x_A^T x_B)$$

Multimodal Transformers

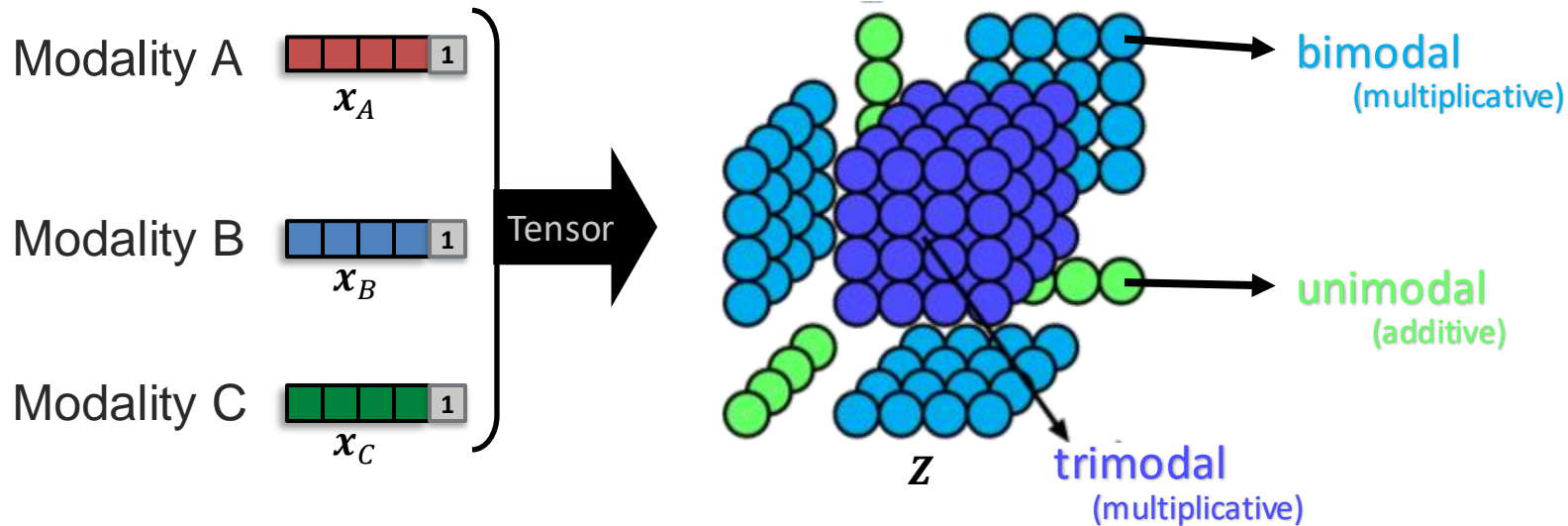


Tensor Fusion



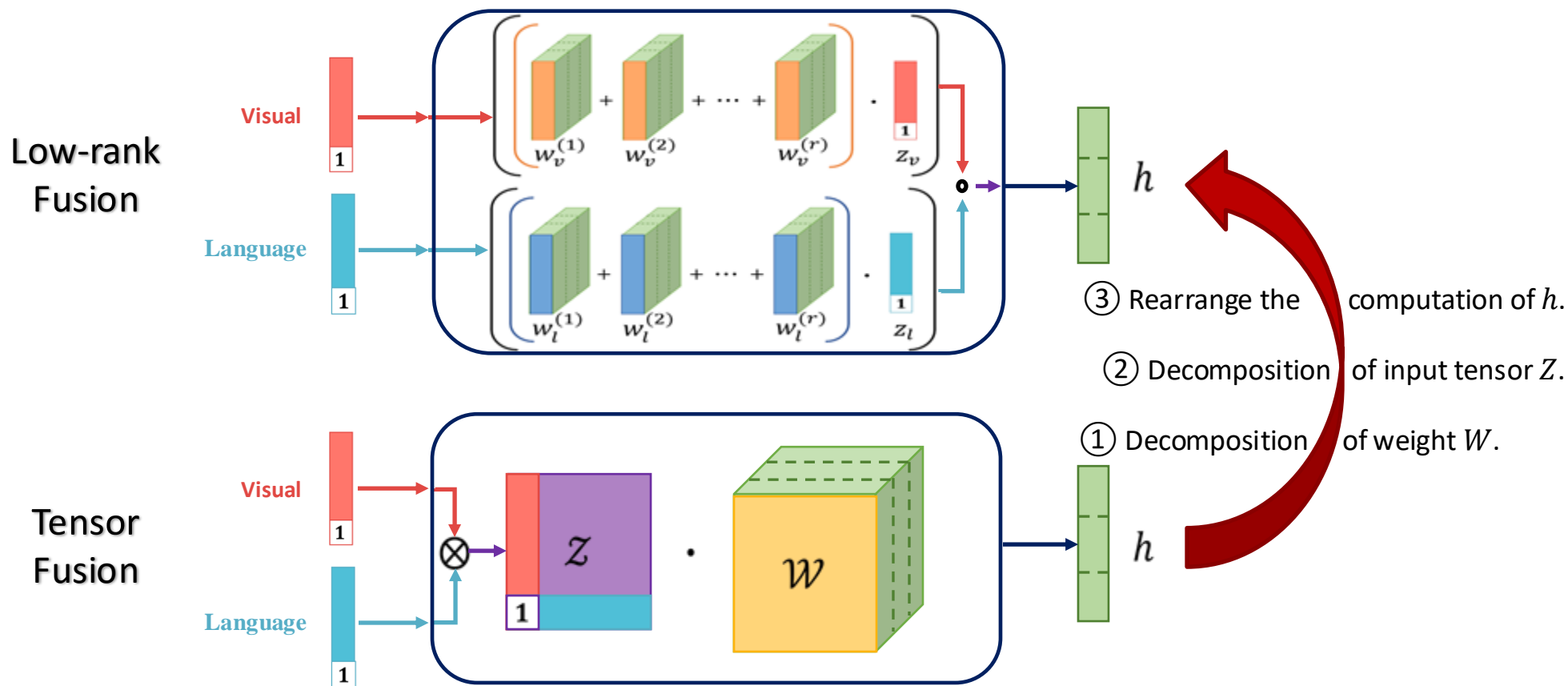
Tensor Fusion (bimodal):

$$Z = w([x_A \ 1]^T \cdot [x_B \ 1])$$

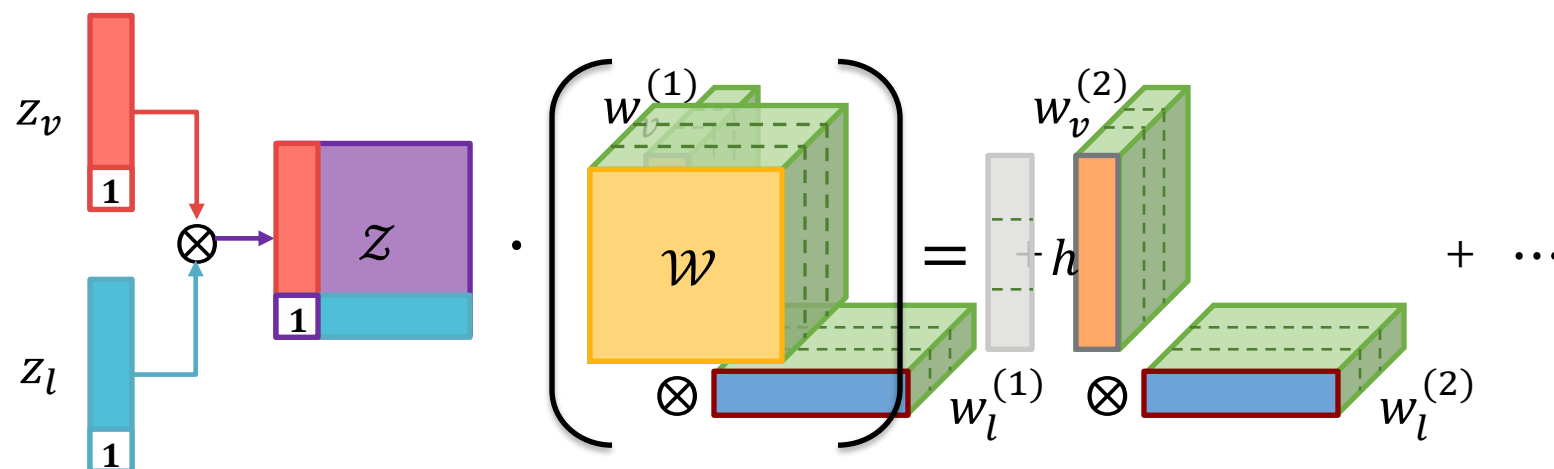


... but the weight matrix may end up quite large!

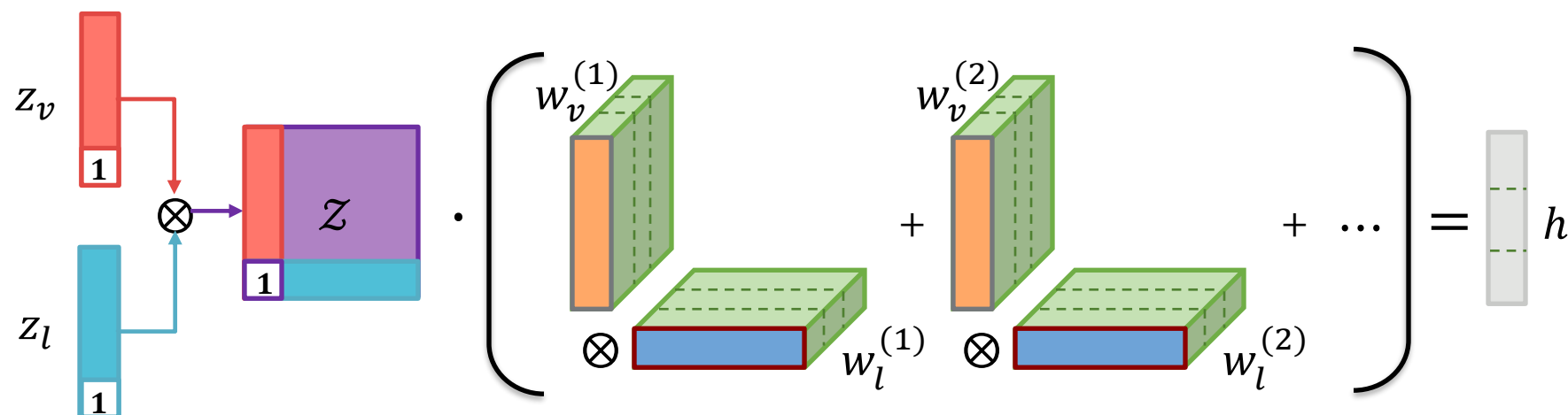
Low-rank Fusion



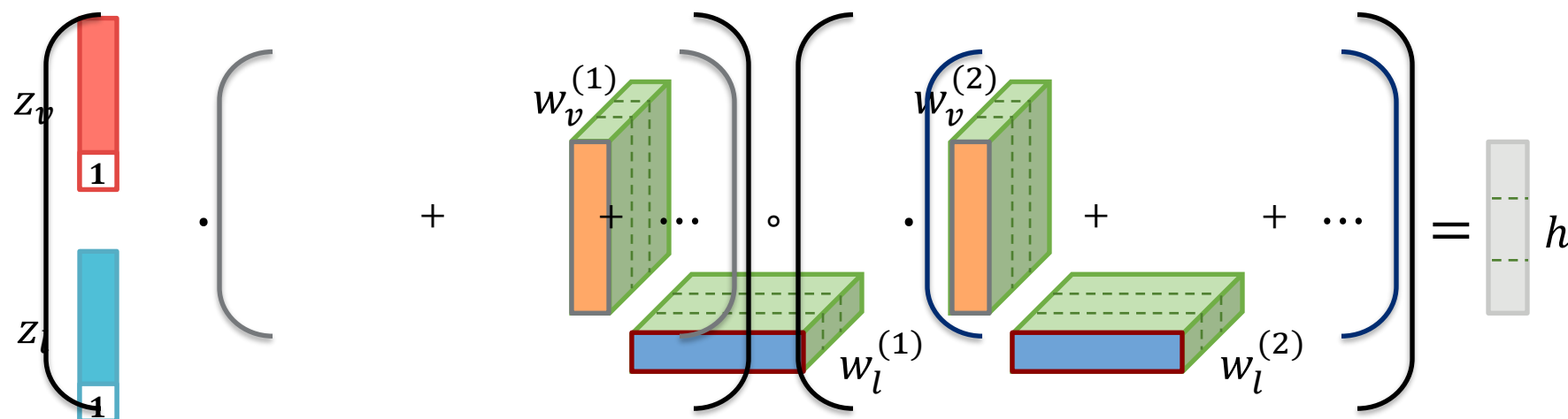
Low-rank Fusion



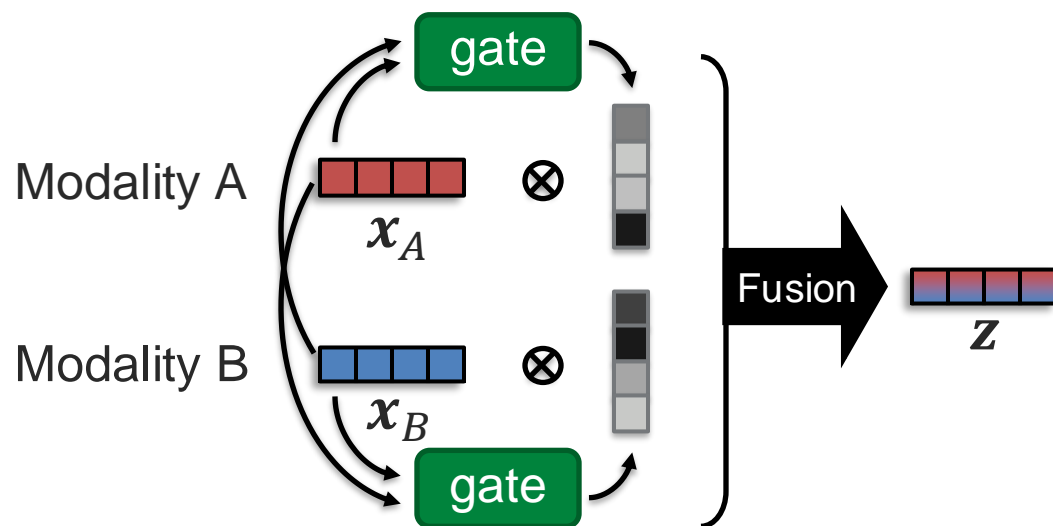
Low-rank Fusion



Low-rank Fusion



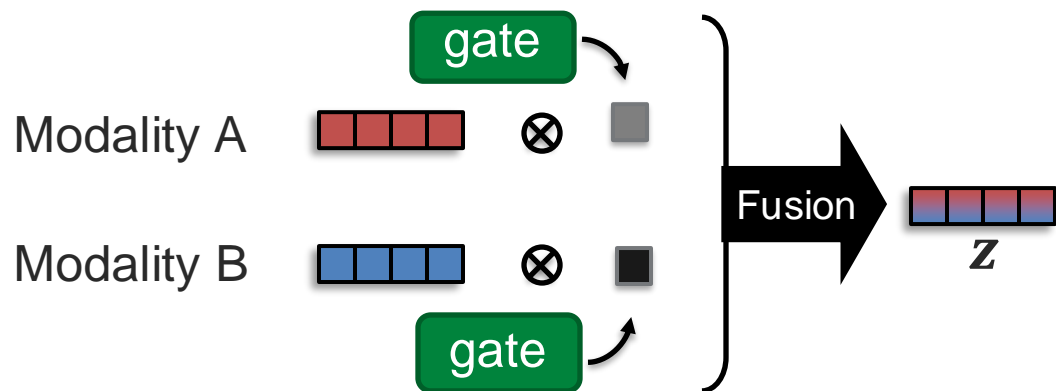
Gated Fusion



Example with additive fusion:

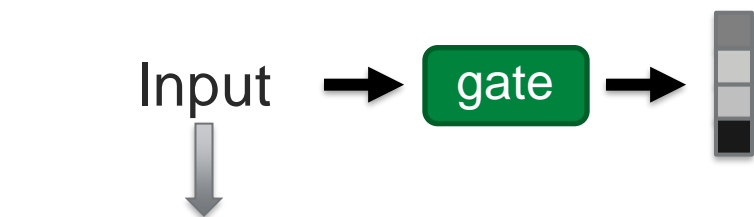
$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

→ g_A and g_B can be seen as attention functions



→ Gating output can be one weight for the whole modality

Gated Fusion



“Neural network designed to mask unwanted signal from propagating forward” (gating)

...or with a more positive view:

“Neural network designed to select preferable signal to move forward” (attention)

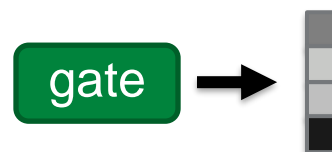
Target modality [horizontal vector of 4 red squares]

Other modality [horizontal vector of 4 blue squares]

All modality

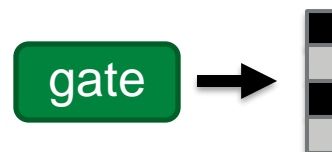


Soft attention



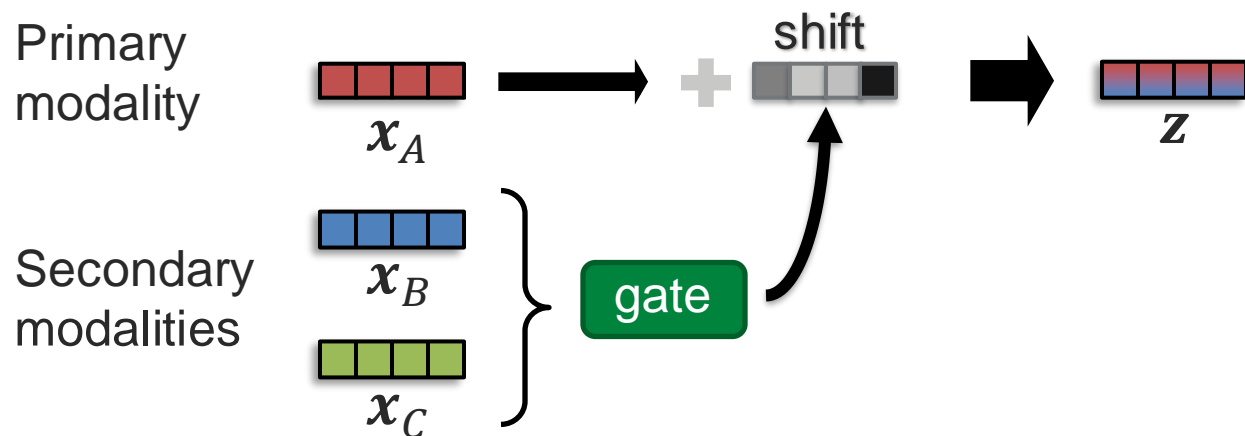
Easier to compute derivative (gradient)

Hard attention



Derivative is harder (e.g., use reinforcement learning)

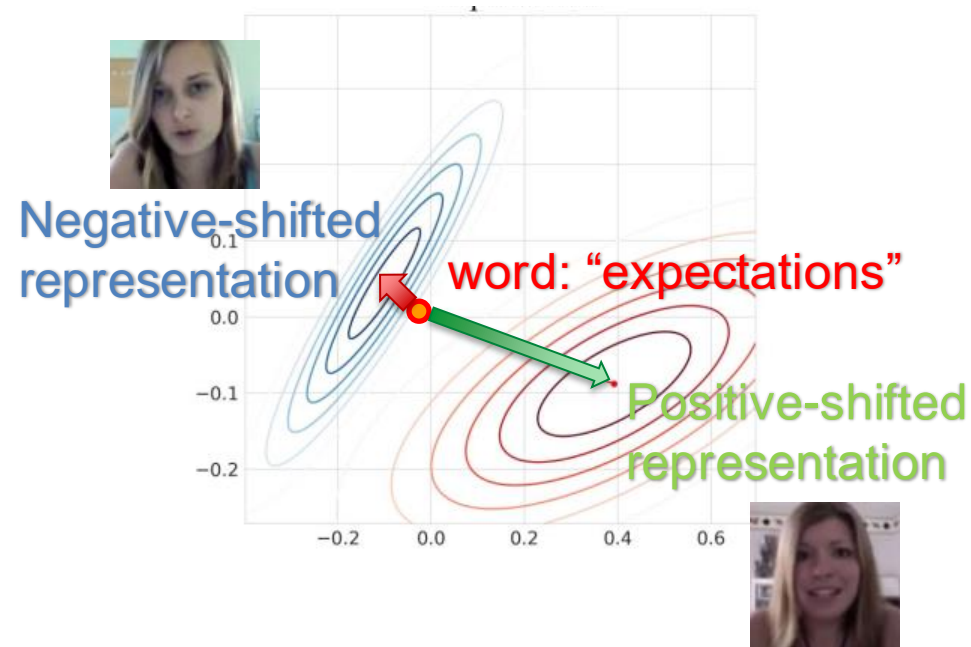
Modality-Shifting Fusion



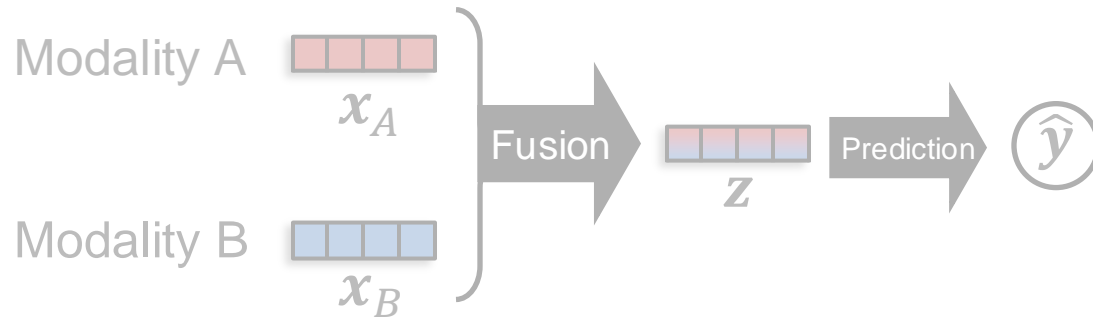
Example with language modality:

Primary modality: language

Secondary modalities: acoustic and visual



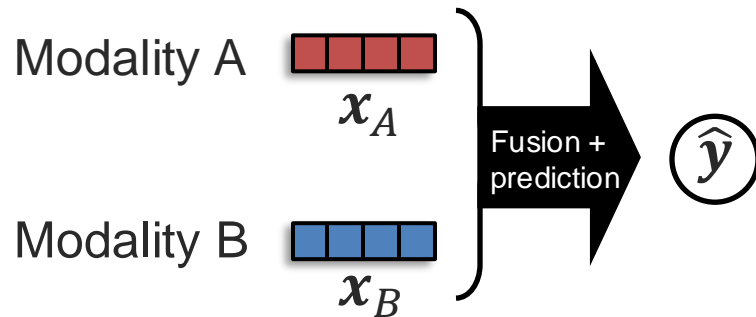
Nonlinear Fusion



Nonlinear fusion:

$$\hat{y} = f(x_A, x_B) \in \mathbb{R}^d$$

where f could be a multi-layer perceptron or any nonlinear model

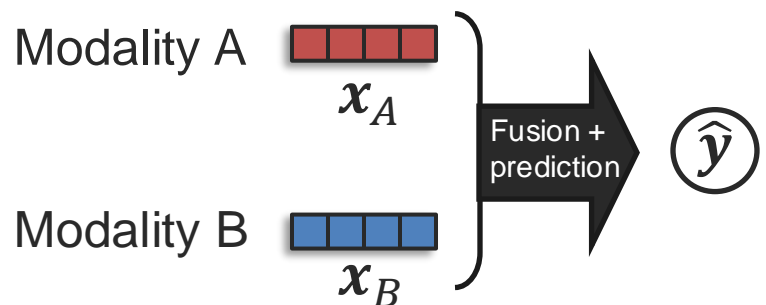


→ This could be seen as *early fusion*:

$$\hat{y} = f([x_A, x_B])$$

... but will our neural network learn the nonlinear interactions?

Measuring Non-Additive Interactions



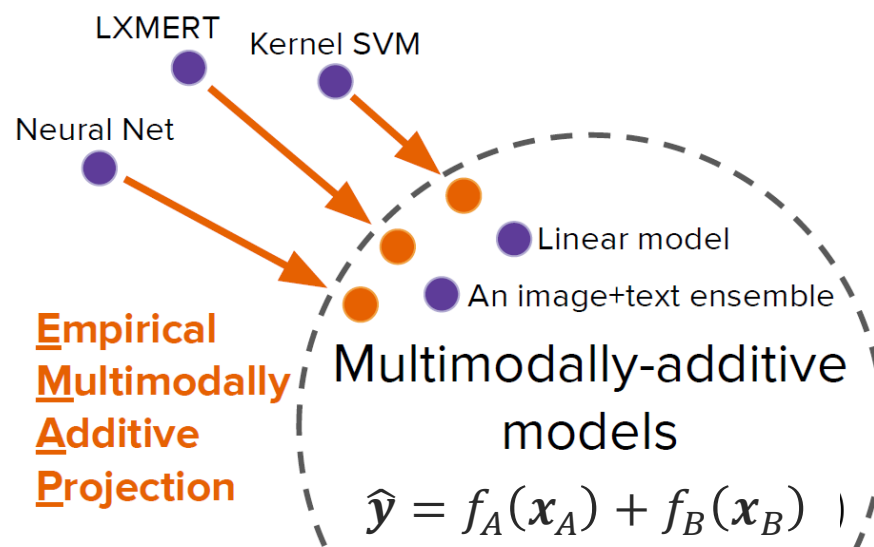
Nonlinear fusion:

$$\hat{y} = f(x_A, x_B)$$

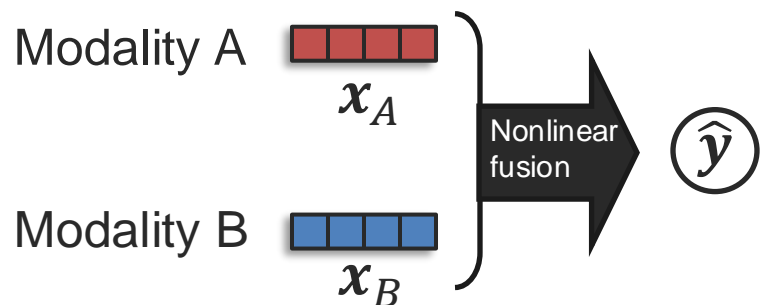
Projection?

Additive fusion:

$$\hat{y} = f_A(x_A) + f_B(x_B)$$



Measuring Non-Additive Interactions



Nonlinear fusion:

$$\hat{\mathbf{y}} = f(\mathbf{x}_A, \mathbf{x}_B)$$

Projection?

Additive fusion:

$$\hat{\mathbf{y}}' = f_A(\mathbf{x}_A) + f_B(\mathbf{x}_B)$$

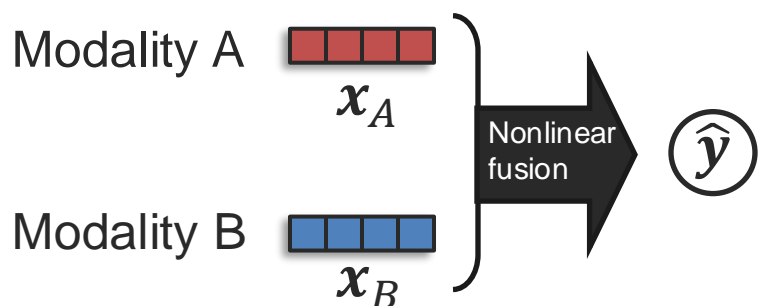
Projection from nonlinear to additive (using EMAP):

$$\tilde{f}(\mathbf{x}_A, \mathbf{x}_B) = \underbrace{\mathbb{E}_{\mathbf{x}_B} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_A(\mathbf{x}_A)} + \underbrace{\mathbb{E}_{\mathbf{x}_A} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_B(\mathbf{x}_B)}$$

Modality A + Modality B

Additive fusion
(approximation)

Measuring Non-Additive Interactions



Nonlinear fusion:

$$\hat{y} = f(\mathbf{x}_A, \mathbf{x}_B)$$

EMAP projection

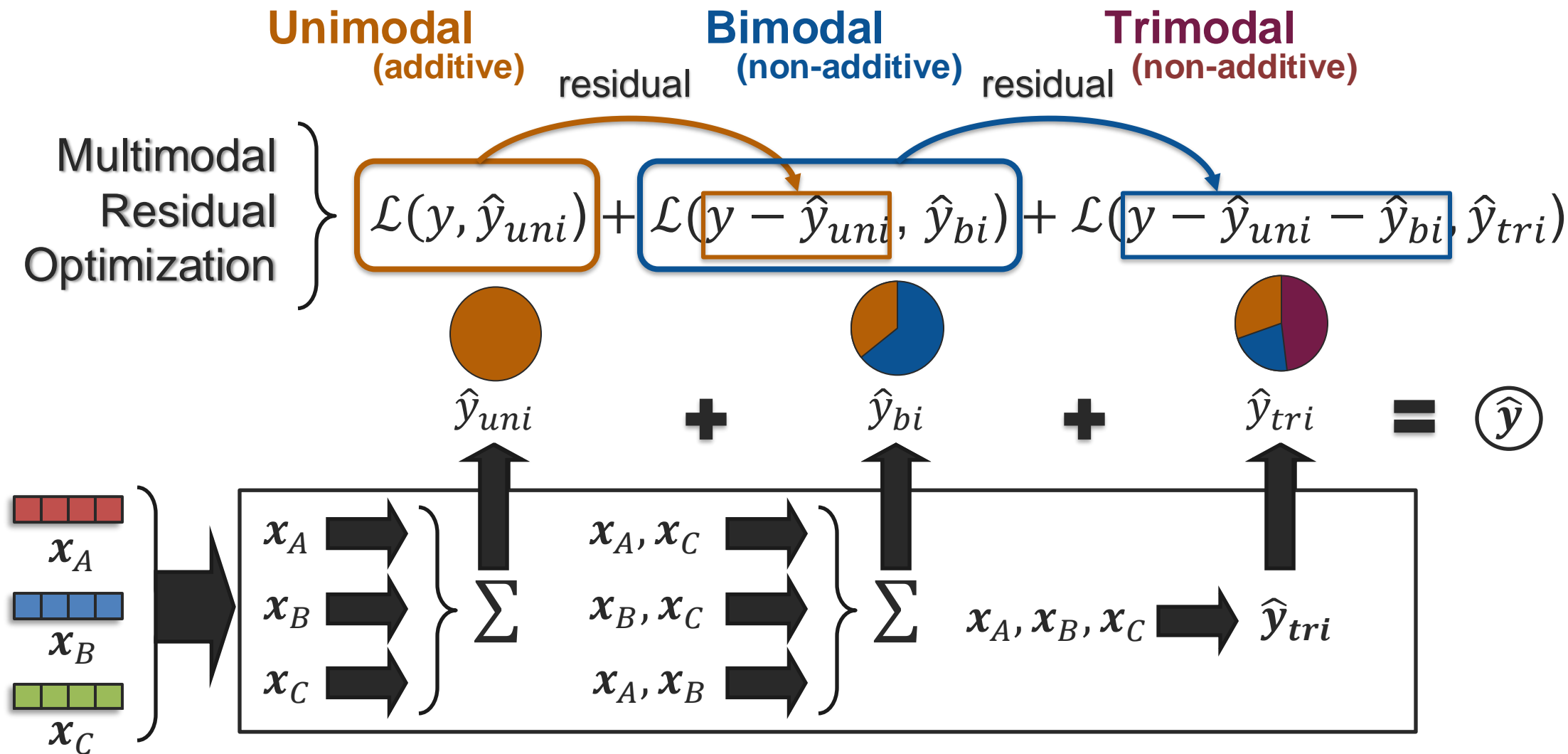
Additive fusion:

$$\hat{y}' = \hat{f}_A(\mathbf{x}_A) + \hat{f}_B(\mathbf{x}_B) + \hat{\mu}$$

		I-INT	I-SEM	I-CTX	T-VIS	R-POP	T-ST1	T-ST2	
Nonlinear	Neural Network	90.4	69.2	78.5	51.1	63.5	71.1	79.9	
Polynomial	Polykernel SVM	91.3	74.4	81.5	50.8	–	72.1	80.9	
Nonlinear	FT LXMERT	83.0	68.5	76.3	53.0	63.0	66.4	78.6	
Nonlinear	↳ + Linear Logits	89.9	73.0	80.7	53.4	64.1	75.5	80.3	
Additive	Linear Model	90.4	72.8	80.9	51.3	63.7	75.6	76.1	
	Best Model	91.3	74.4	81.5	53.4	64.2	75.5	80.9	
Additive	↳ + EMAP	91.1	74.2	81.3	51.0	64.1	75.9	80.7	Always a good baseline! Differences are small!!!

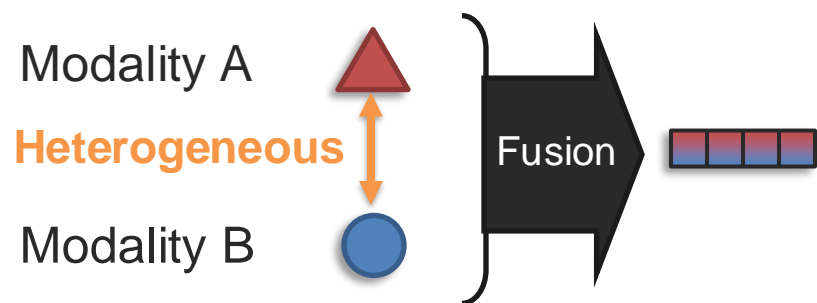
Non-Additive Interactions

Idea: prioritize simpler interactions

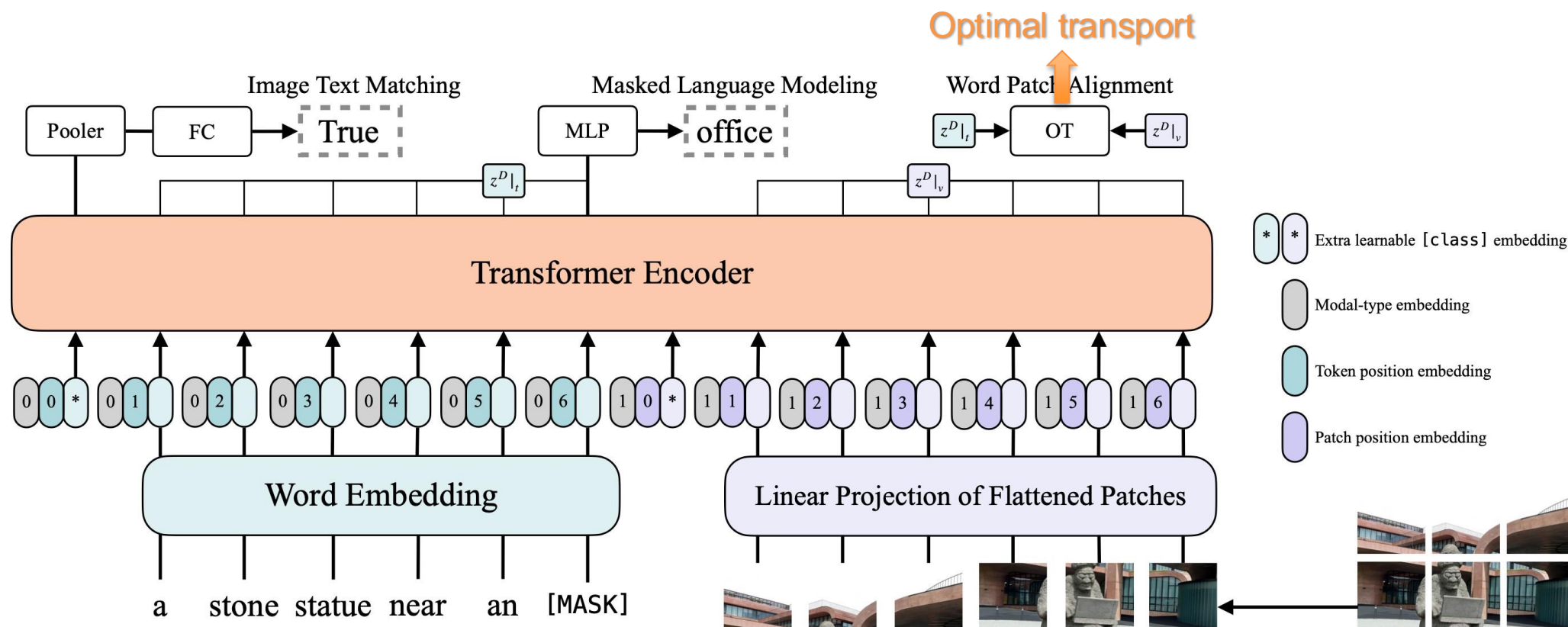


Fusion with Heterogeneous Modalities

Example: From feature fusion to early fusion



Visual-and-Language Transformer (ViLT) (\approx BERT + ViT)



[original slide co-developed with Louis-Philippe Morency for CMU course 11-777]

[Kim et al., ViLT: Vision-and-Language Transformer Without Convolution or Region Supervision. ICML 2021]

Visual-and-Language Transformer (ViLT)

Example of alignment between modalities:



a display of **flowers** growing out and over the retaining **wall** in front of **cottages** on a **cloudy** day.



flowers



wall



cottages



cloudy



a room with a **rug**, a **chair**, a **painting**, and a **plant**.



rug



chair

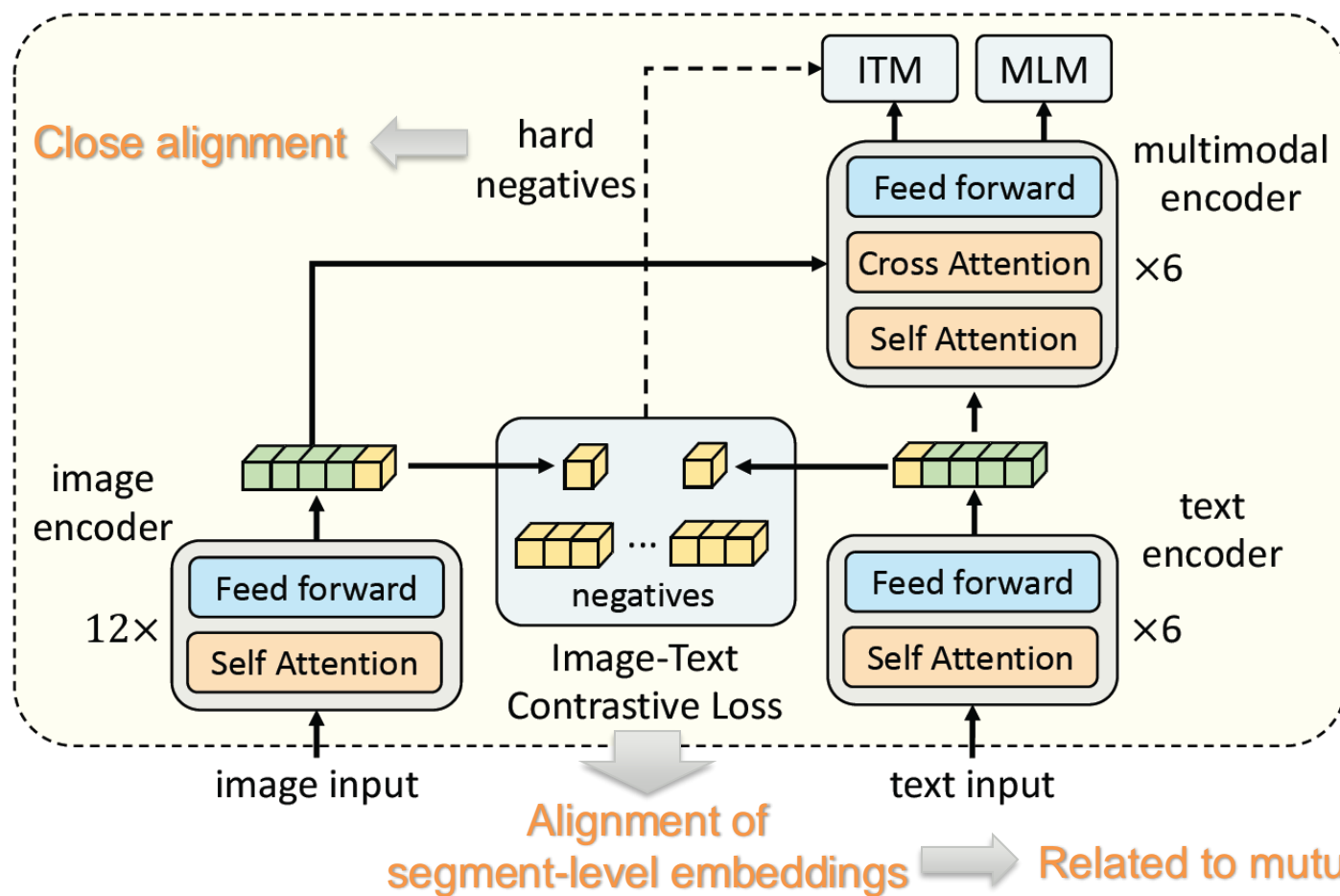


painting



plant

ALBEF: Align Before Fusion (\approx BERT + ViT + CLIP-ish)

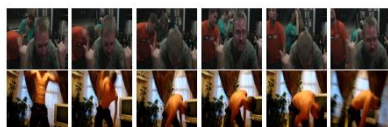


[original slide co-developed with Louis-Philippe Morency for CMU course 11-777]

[Li et al., Align before Fuse: Vision and Language Representation Learning with Momentum Distillation. NeurIPS 2021]

Nonlinear Fusion

Kinetics dataset



(a) headbanging



(c) shaking hands



(e) robot dancing



(g) riding a bike



Adding more modalities should always help?

Modalities: **RGB** (video clips)

A (Audio features)

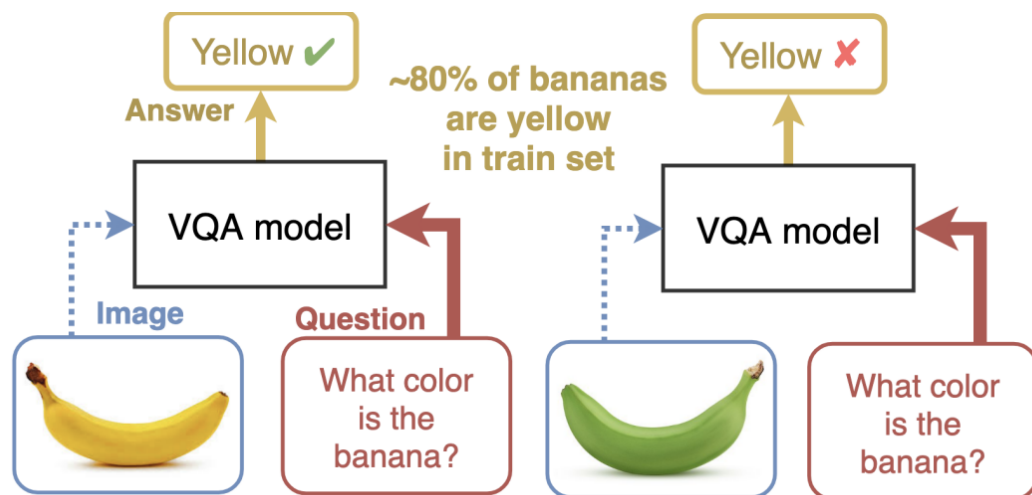
OF (optical flow - motion)

Dataset	Multi-modal	V@1	Best Uni	V@1	Drop
Kinetics	A + RGB	71.4	RGB	72.6	-1.2
	RGB + OF	71.3	RGB	72.6	-1.3
	A + OF	58.3	OF	62.1	-3.8
	A + RGB + OF	70.0	RGB	72.6	-2.6

But sometimes multimodal doesn't help! **Why?**

Unimodal Biases

Finding: VQA models answer the question without looking at the image

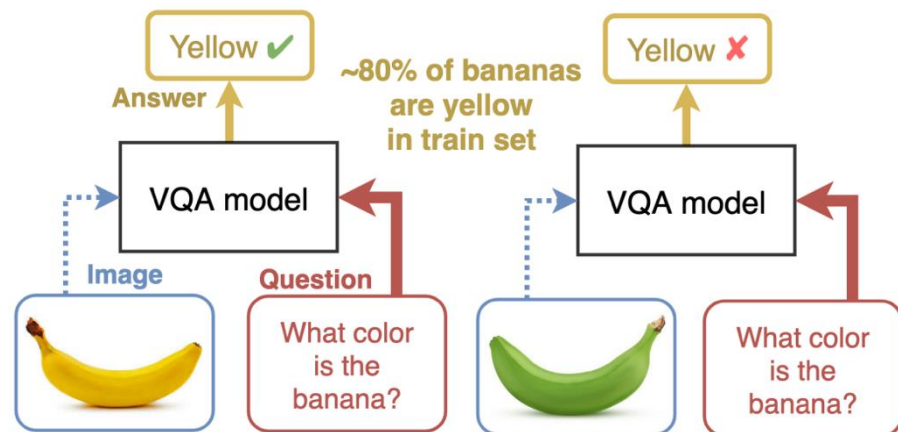


Finding: Image captioning models capture spurious correlations between gender and generated actions.



Unimodal Biases

VQA models answer the question without looking at the image

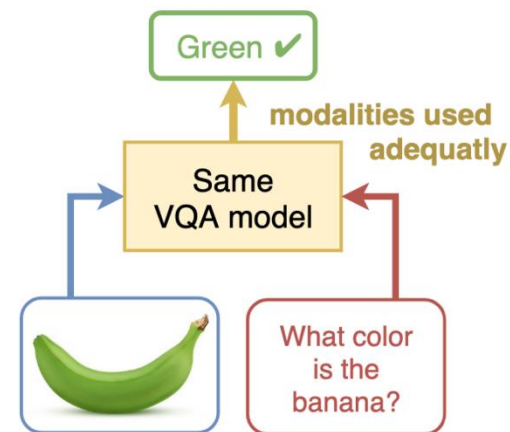


Balancing modalities

Balancing training



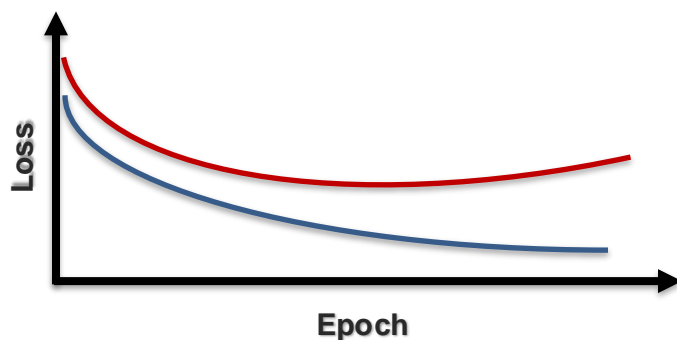
Not the case when trained with RUBi



Optimization Challenges

2 explanations for drop in performance:

1. Multimodal networks are more prone to overfitting due to **increased complexity**
2. Different modalities overfit and generalize at **different rates**



Key idea 1: compute overfitting-to-generalization ratio (OGR)

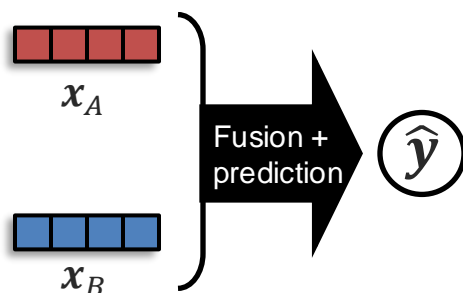


Gap between training and valid loss

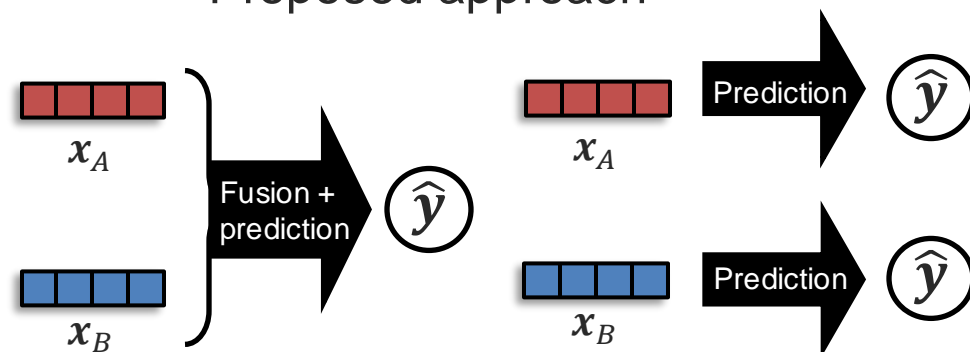
OGR wrt each modality tells us
how much to train that modality

Optimization Challenges

Conventional approach



Proposed approach



Key idea 2: Simultaneously train unimodal networks to estimate OGR wrt each modality

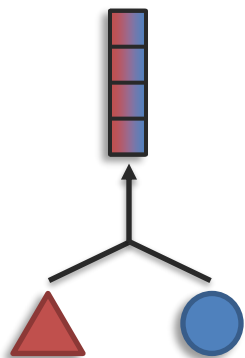


Reweight multimodal loss using unimodal OGR values



Allows to better balance generalization & overfitting rate of different modalities

Summary: How To Multimodal Fusion



Definition: Learn a joint representation that models cross-modal interactions between individual elements of different modalities



Lecture Summary

- 1 Basics of multimodal fusion
- 2 Early, intermediate, late fusion
- 3 Multiplicative and dynamic fusion
- 4 Complex fusion and improving optimization

Assignments for This Coming Week

For project:

- I gave feedback and assigned primary TA.
- Meet with me and primary TA every other week.
- Should have finalized main ideas and experimental setup, have baseline models working, progress towards implementing new ideas.

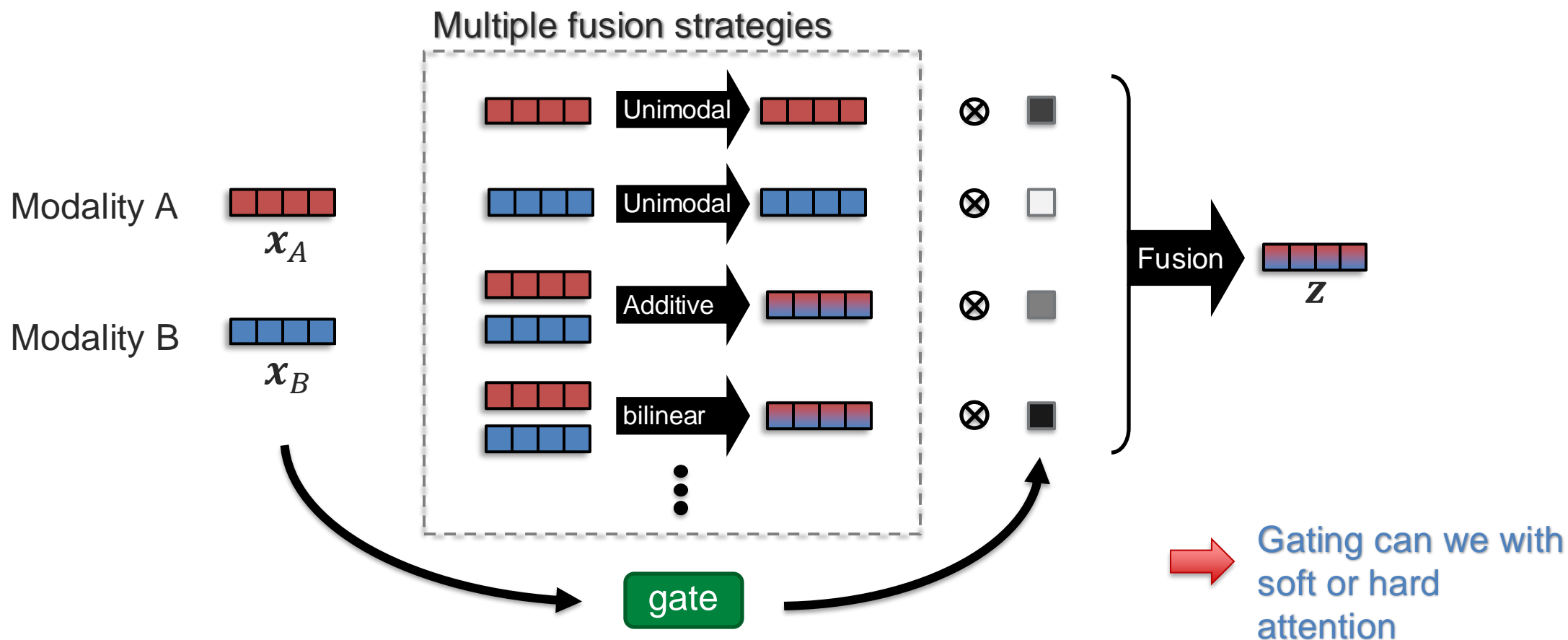
Reading assignment due tomorrow Wednesday (3/12).

This Thursday (3/13): third reading discussion on **multimodal alignment**.

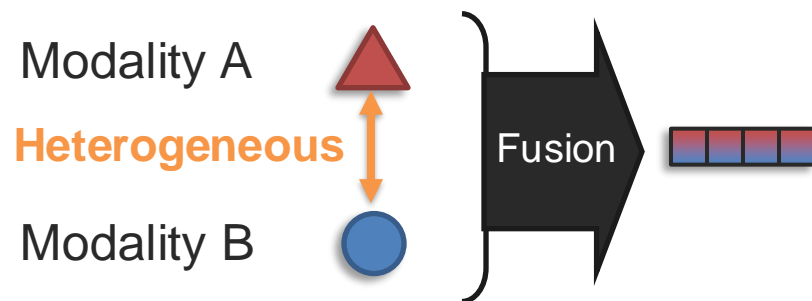
What views for contrastive learning

Platonic representation hypothesis

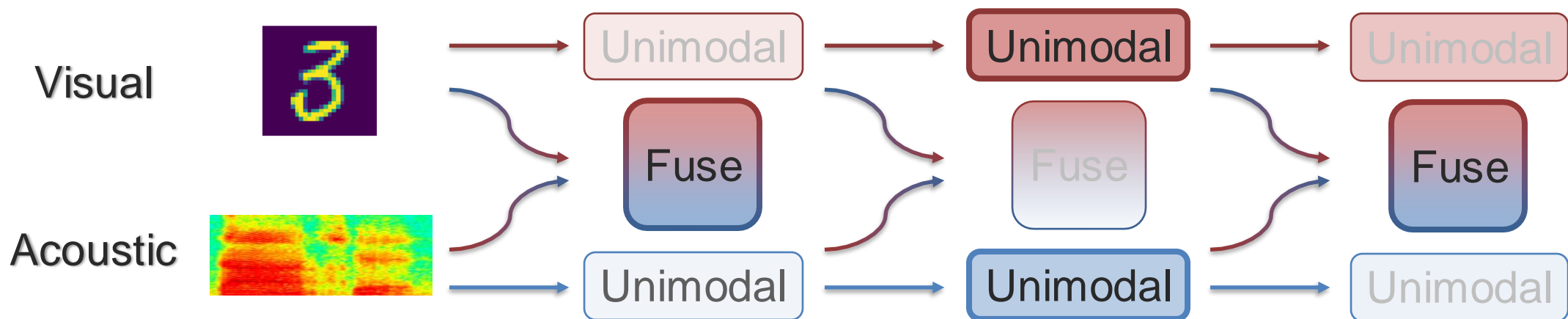
Dynamic Fusion



Dynamic Early Fusion



Idea: Deciding when to fuse in early fusion



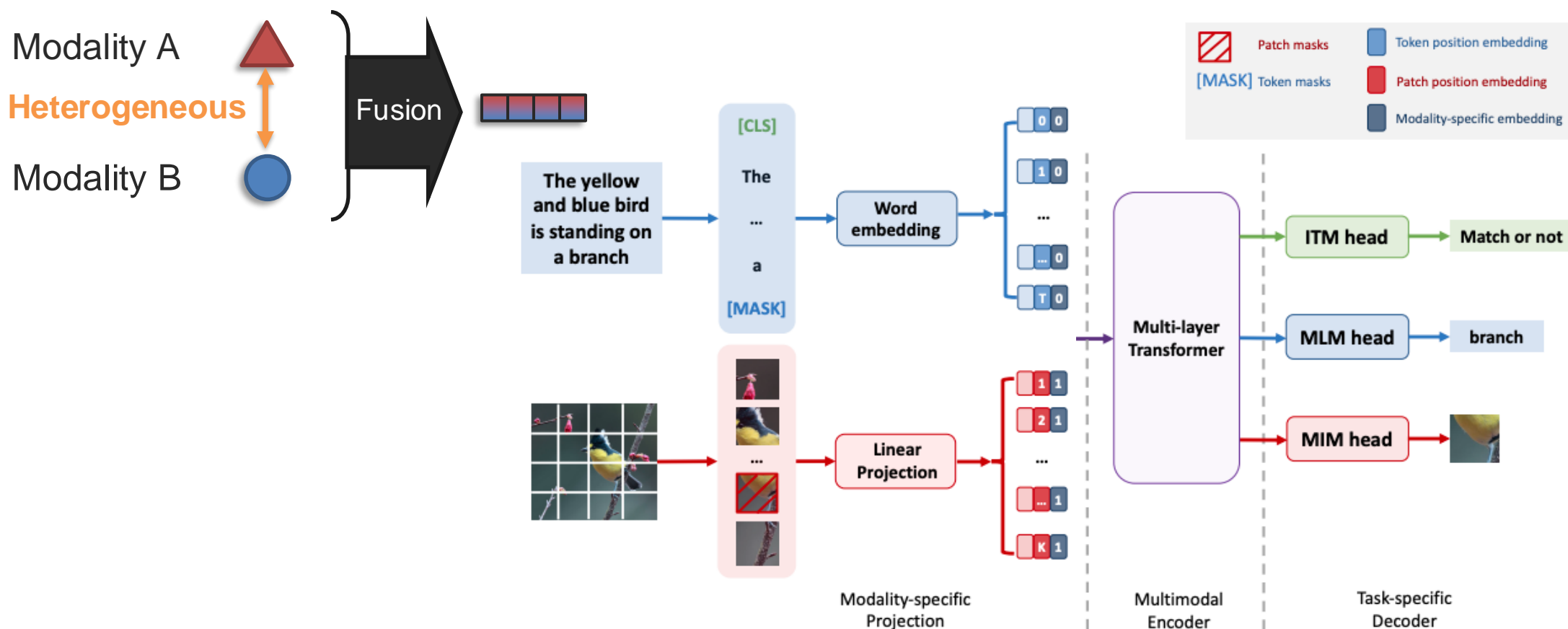
[Xue and Marculescu, Dynamic Multimodal Fusion, arxiv 2022]

[Xu et al., MUFASA: Multimodal Fusion Architecture Search for Electronic Health Records. AAAI 2021]

[Liu et al., DARTS: Differentiable Architecture Search. ICLR 2019]

Fusion with Heterogeneous Modalities

Example: From feature fusion to early fusion



[Liang et al., High-modality Multimodal Transformer. TMLR 2022]

[Gui et al., Training Vision-Language Transformers from Captions. arxiv 2022]