



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal Machine Learning

Lecture 3.2: Multimodal Representations (Part 1)

Louis-Philippe Morency

** Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yanatan Bisk. Some slides from Jeffrey Girard.*

Administrative Stuff

Reading Assignments – Reminder

Week 3 reading assignment was posted

1. **Friday 8pm:** Post your summary
2. **Monday 8pm:** End of the reading assignment

Be sure to post your discussion comments before Monday 8pm!

 Start the discussion early 😊

 Late submissions will be accounted

Primary TAs

- Each team will have one primary TA
- Meetings with primary TA will be scheduled for next week
 - Feedback for the pre-proposals
- Contact your primary TA anytime (piazza or email)
 - Groups will be created in Piazza for each team
- Some projects may have a secondary TA, with complementary expertise

First Project Assignment

Due date: Sunday 9/25 at 8m

Four main sections:

- Introduction
- Related work
- Experimental setup
- Research ideas

Follows ICML paper format



The two main sections are related work and research ideas



teammates = # research ideas



Page limit depends on team size:

- 3 students : 4 pages + references
- 4 students : 4.5 pages + references
- 5 students : 5 pages + references
- 6 students : 5.5 pages + references



Language
Technologies
Institute

Carnegie
Mellon
University

Multimodal Machine Learning

Lecture 3.2: Multimodal Representations (Part 1)

Louis-Philippe Morency

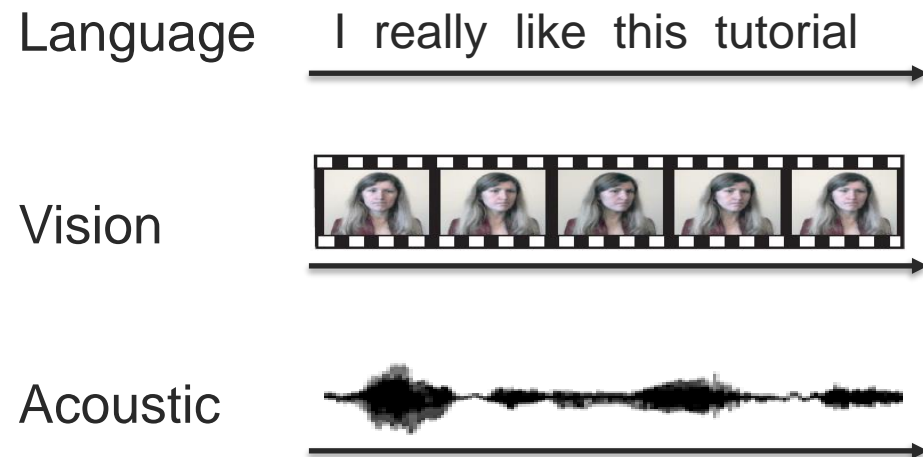
** Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yoanatan Bisk. Some slides from Jeffrey Girard.*

Lecture Objectives

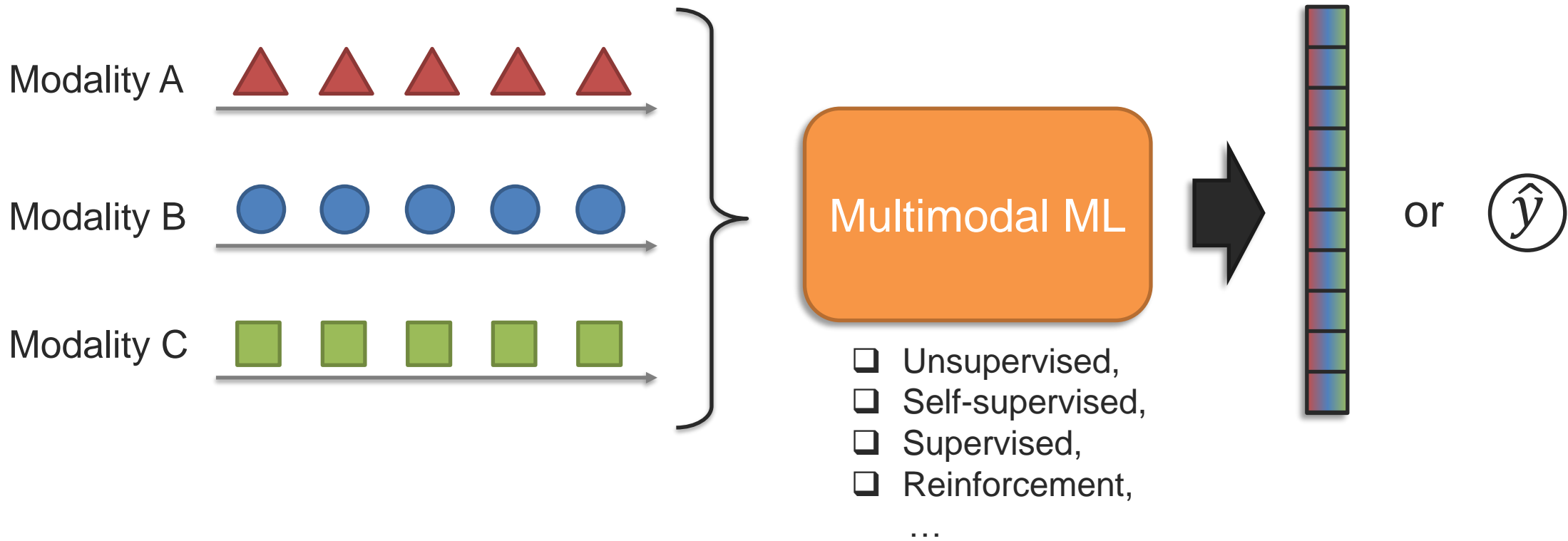
- Multimodal representations
 - Cross-modal interactions
- Representation fusion
 - Additive and multiplicative fusion
 - Tensor and polynomial fusion
 - Gated fusion
 - Modality-shift fusion
 - Dynamic fusion
 - Fusion on raw modalities
 - Multimodal autoencoder
- Measuring non-additive interactions

Multimodal Representation

Multimodal Machine Learning



Multimodal Machine Learning

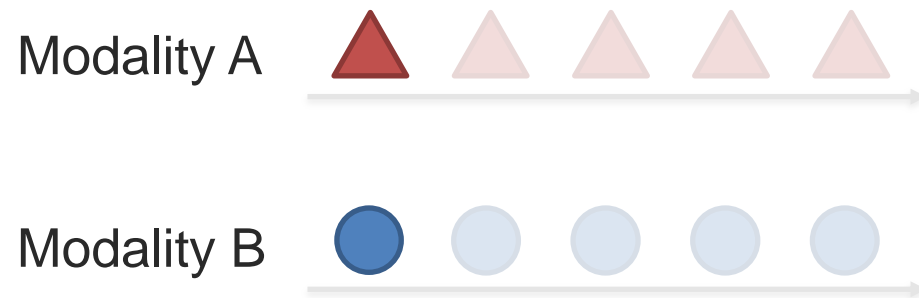


Challenge 1: Representation

Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

➡ This is a core building block for most multimodal modeling problems!

Individual elements:



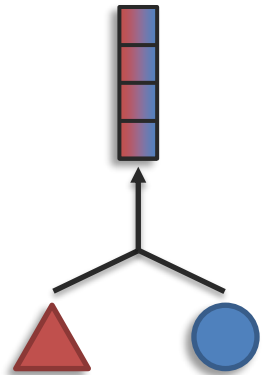
*It can be seen as a “local” representation
or
representation using holistic features*

Challenge 1: Representation

Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

Sub-challenges:

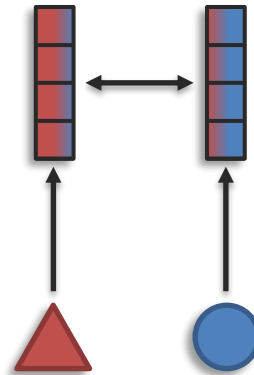
Fusion



modalities $>$ # representations

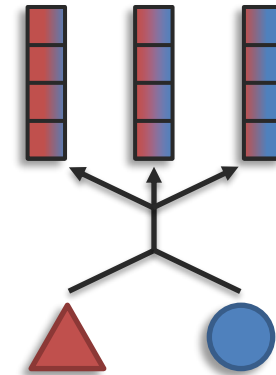
Today

Coordination



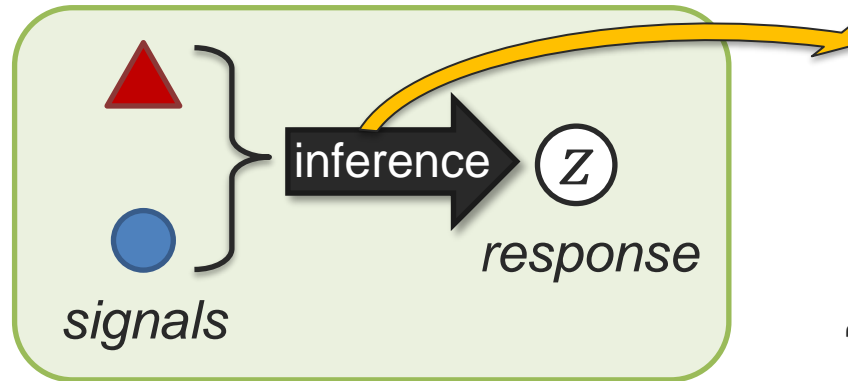
modalities = # representations

Fission



modalities $<$ # representations

Cross-modal Interactions

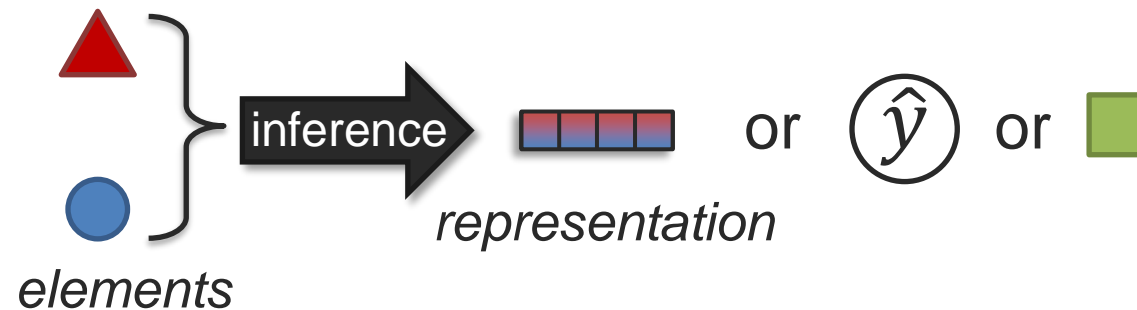


Interactions happen during inference!

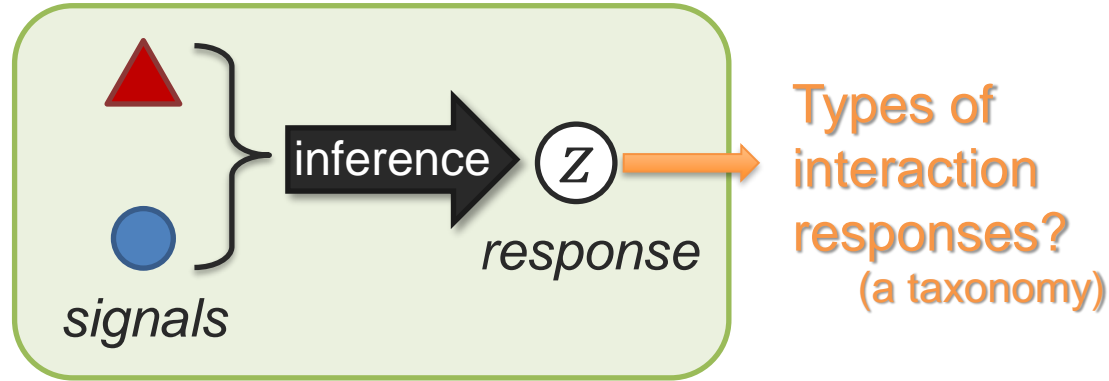


"Inference" examples:

- Representation fusion
- Prediction task
- Modality translation



Interconnected Modalities



Is this a living room?

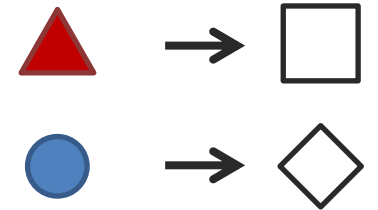


inference → **Yes!**

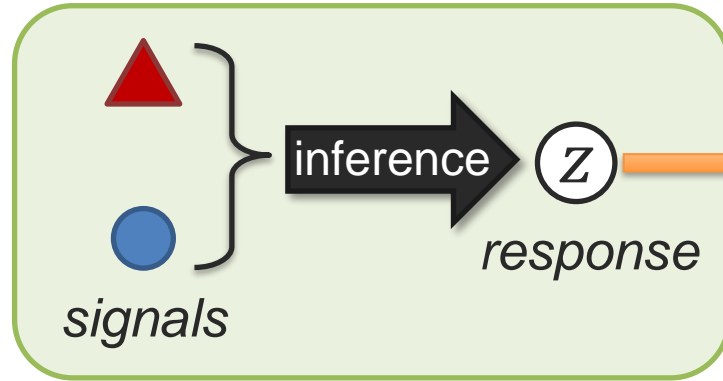
A teacup on the right of a laptop in a clean room.

inference → **No, probably study room.**

Unimodal
Non-redundancy



Interconnected Modalities



Types of interaction responses?
(a taxonomy)

Is this a living room?



A teacup on the right of a laptop in a clean room.

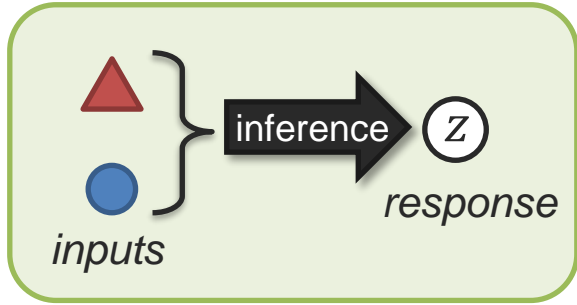
inference **Yes!**

Unimodal
Non-redundancy



Multimodal
dominance

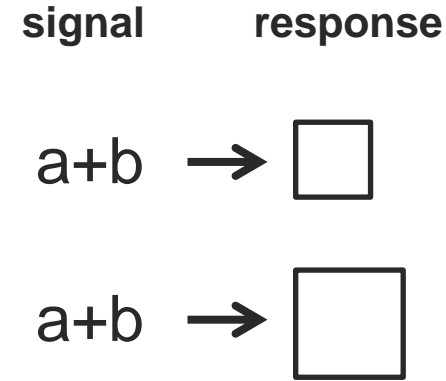
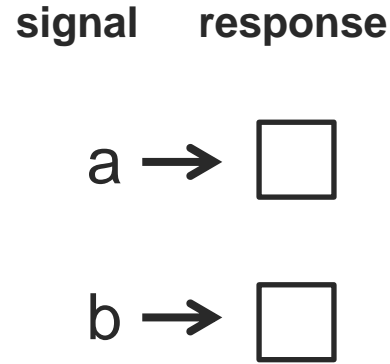
Taxonomy of Interaction Responses – A Behavioral Science View



Multimodal Communication



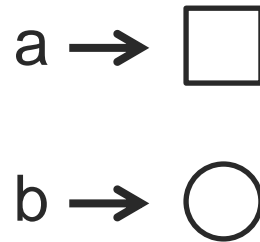
Redundancy



Equivalence

Enhancement

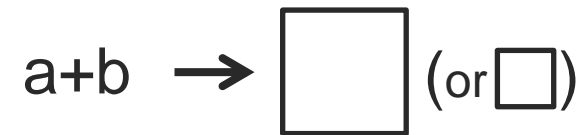
Nonredundancy



Independence



Dominance



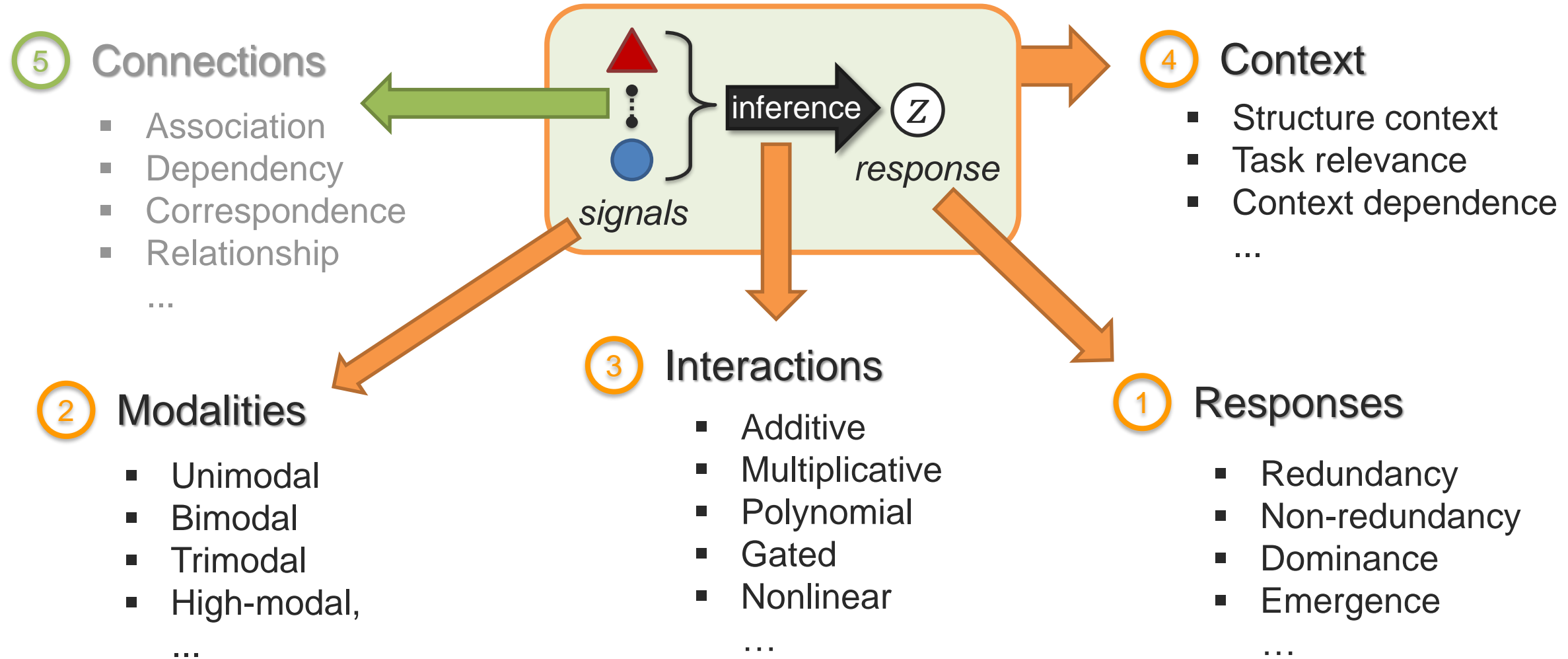
Modulation



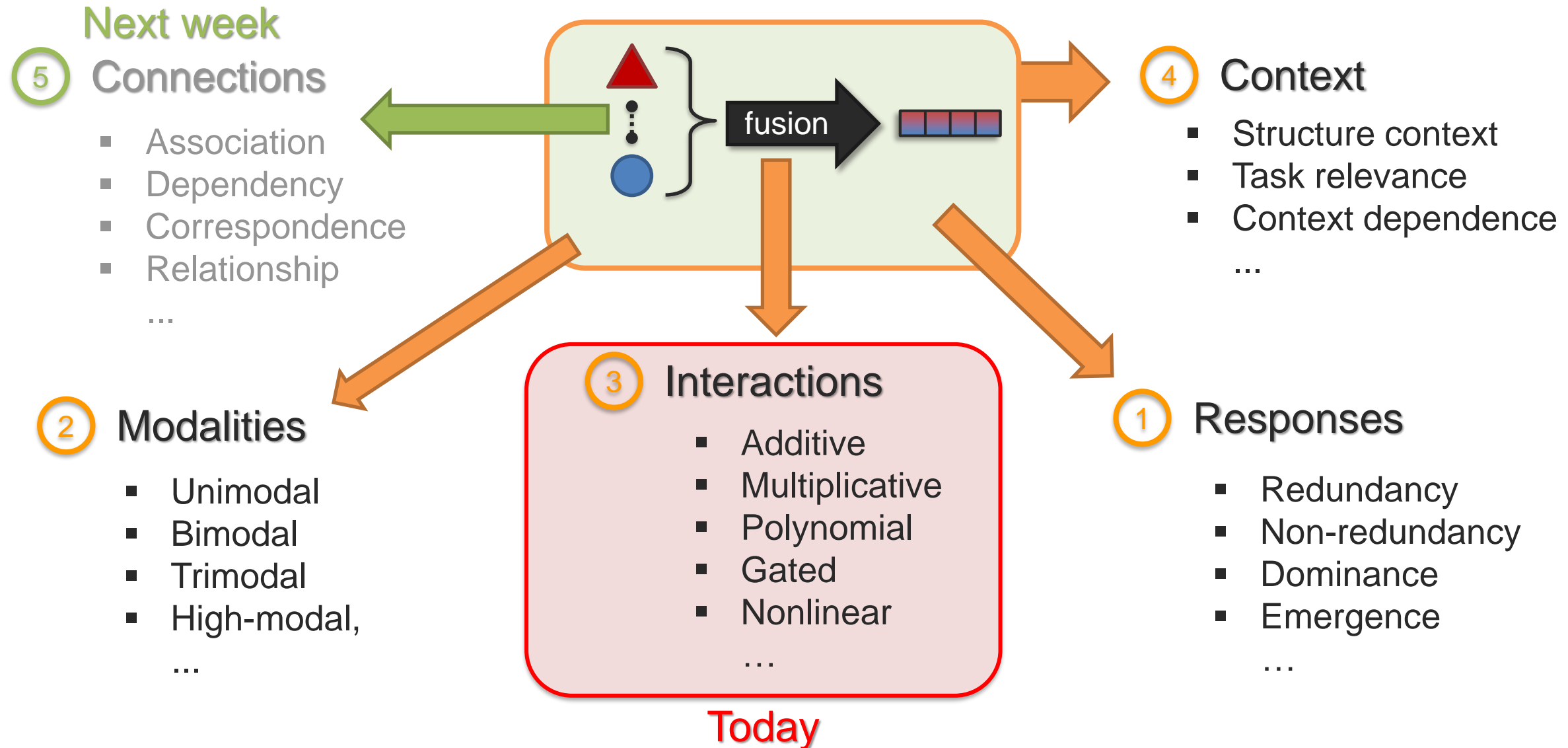
Emergence

Partan and Marler (2005). *Issues in the classification of multimodal communication signals*. *American Naturalist*, 166(2)

Cross-modal Interactions – A Taxonomy

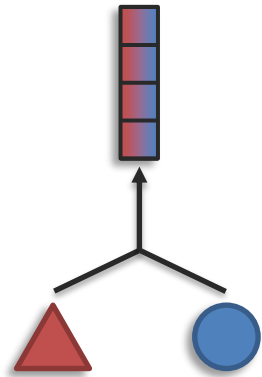


Cross-modal Interactions – Representation Fusion



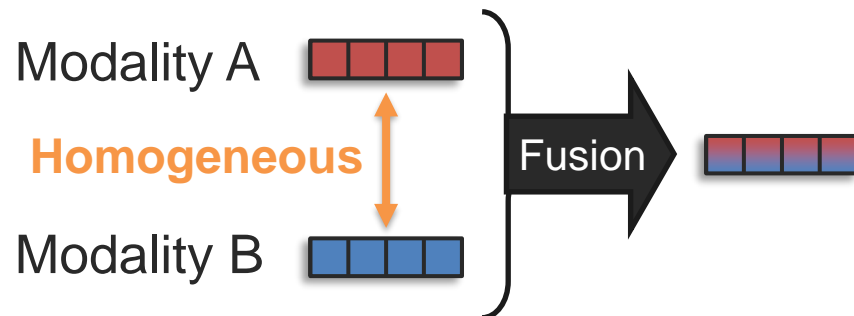
Representation Fusion

Sub-Challenge 1a: Representation Fusion

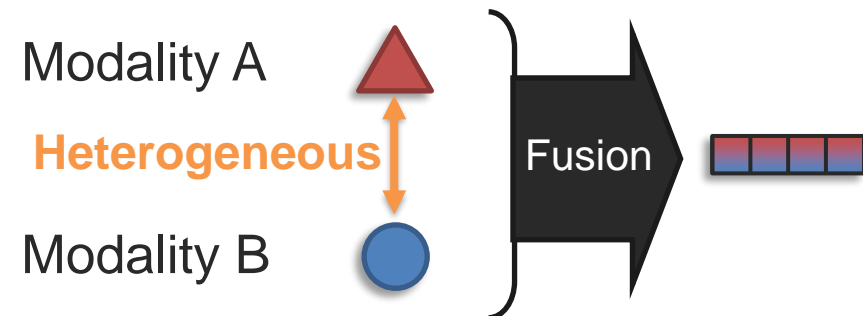


Definition: Learn a joint representation that models cross-modal interactions between individual elements of different modalities

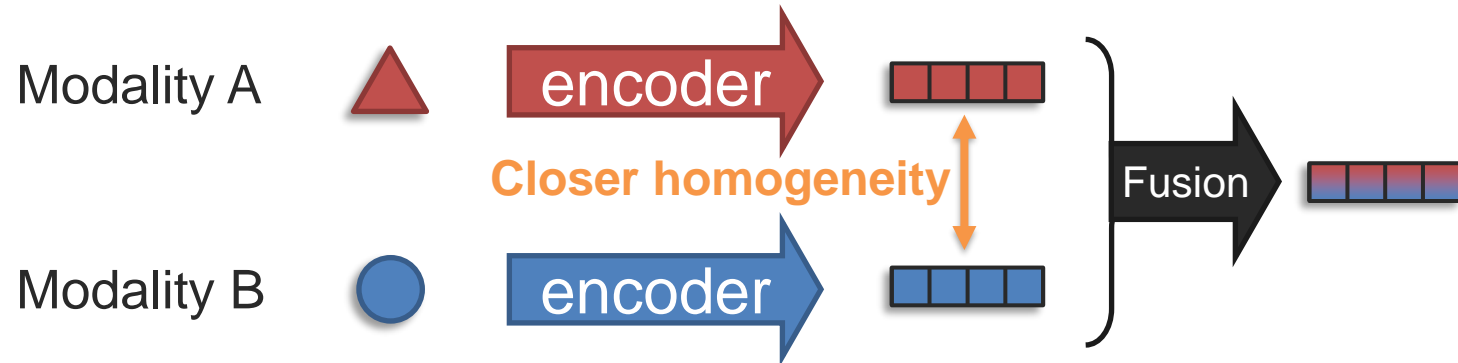
Basic fusion:



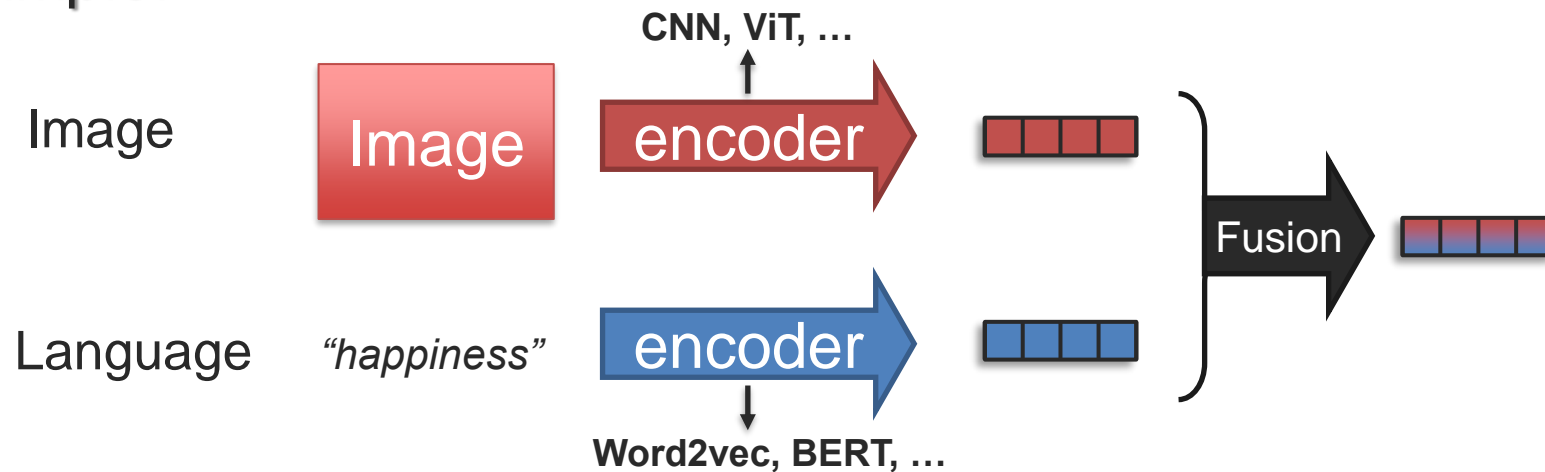
Raw-modality fusion:



Fusion with Unimodal Encoders



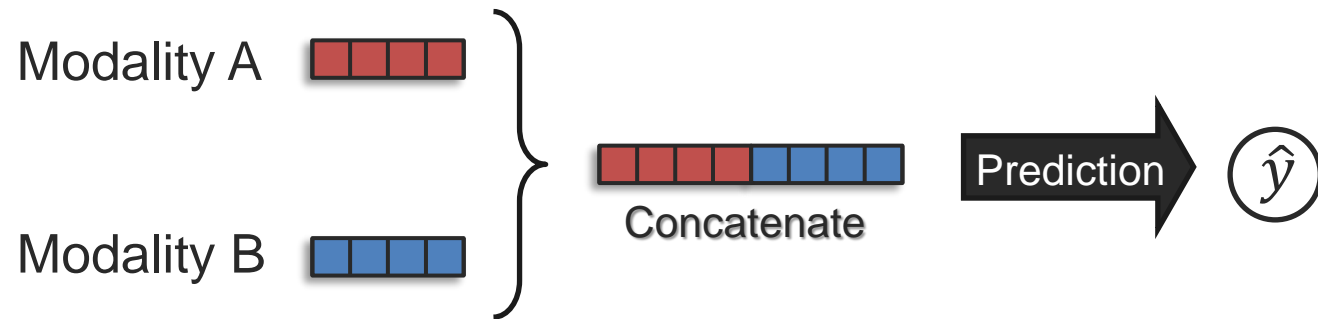
Example:



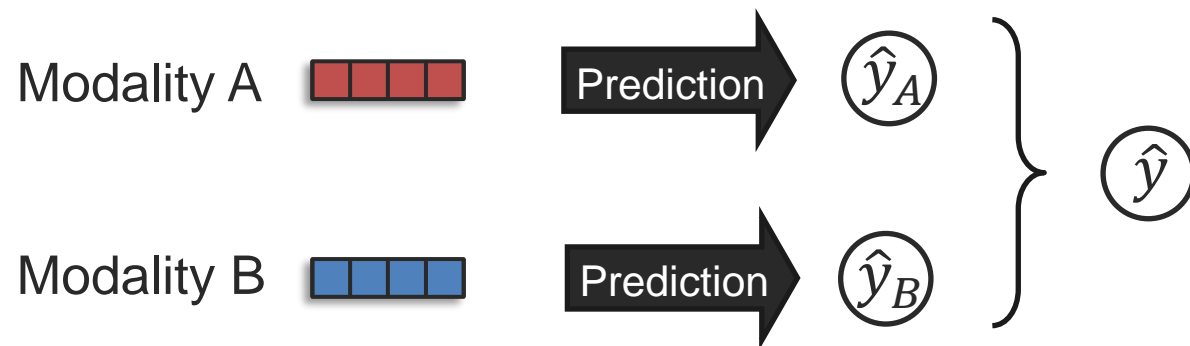
➡ Unimodal encoders can be jointly learned with fusion network, or pre-trained

Early and Late Fusion – A historical View

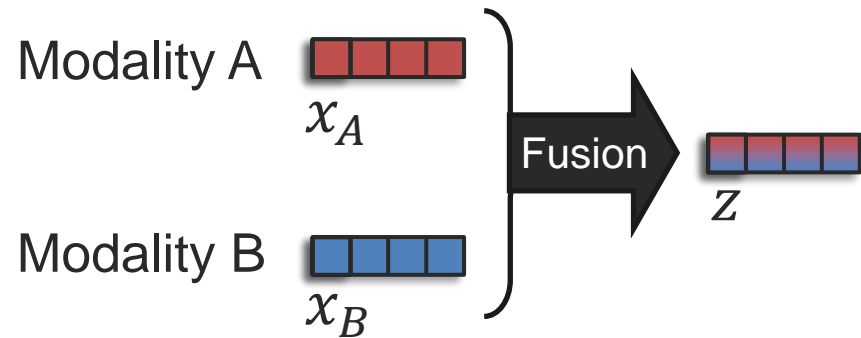
Early fusion:



Late fusion:



Basic Concepts for Representation Fusion (aka, Basic Fusion)



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

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

Linear regression:

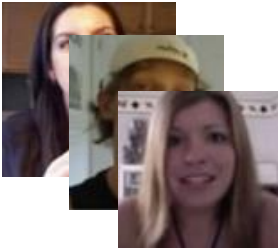
$$z = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$

intercept (bias term) Additive terms Multiplicative term error (residual term)

Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

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:

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

intercept (bias term) Additive terms Multiplicative term 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

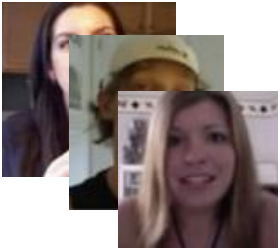
w_3 : effect from x_A and x_B interaction only

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

Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

300 book reviews



y : audience score

x_A : percentage of smiling

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

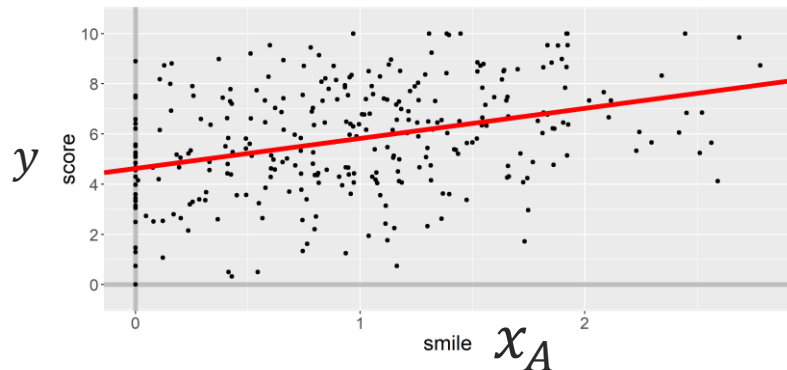
H1: Does smiling reveal what the audience score was?

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

Linear regression:

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

slope



Confidence interval: "95% confident that w parameter is contained within this interval"

	Estimate	95% CI
w_0	4.63	[4.20, 5.06]
w_1	1.20	[0.83, 1.57]

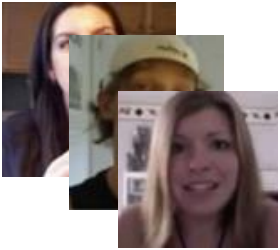
p-values would be another way to test hypothesis

Confidence interval does not contain 0, so effect is significant

Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

300 book reviews



y : audience score

x_A : percentage of smiling

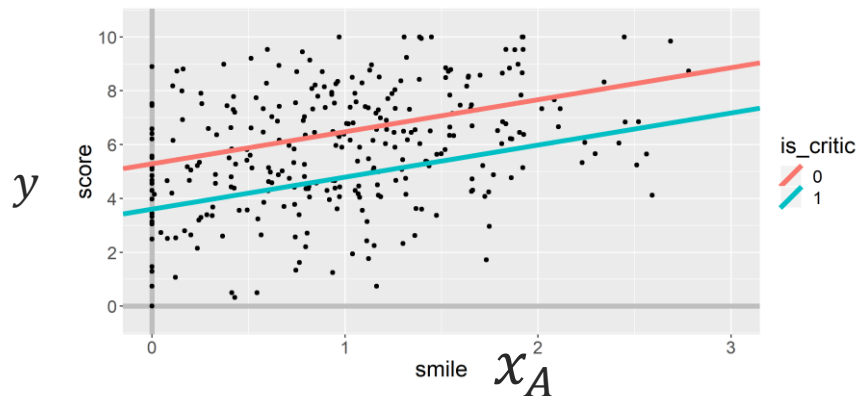
x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

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

Linear regression:

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

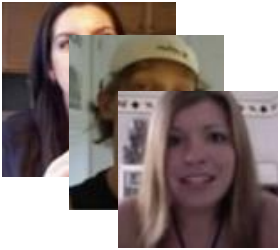


	Estimate	95% CI	
w_0	5.29	[4.86, 5.73]	
w_1	1.19	[0.85, 1.53]	➔ Positive effect
w_2	-1.69	[-2.14, -1.24]	➔ Negative effect

Linear Regression

Linear regression is used to test research hypotheses, over a whole dataset

300 book reviews



y : audience score

x_A : percentage of smiling

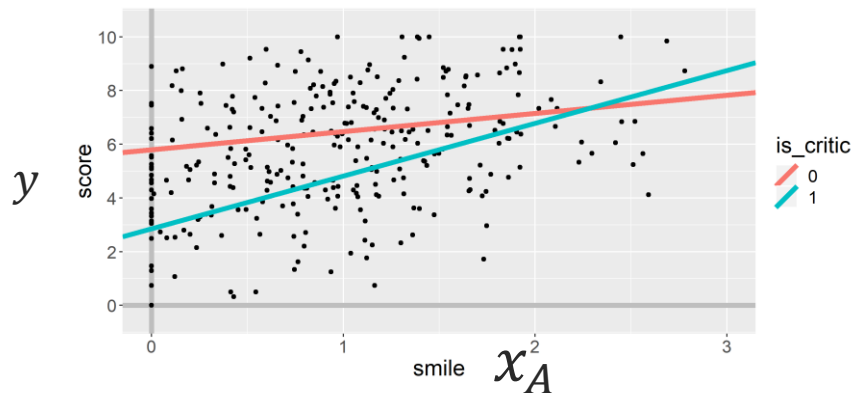
x_B : professional status
(0=non-critic, 1=critic)

H1: Does smiling reveal what the audience score was?

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

Linear regression:

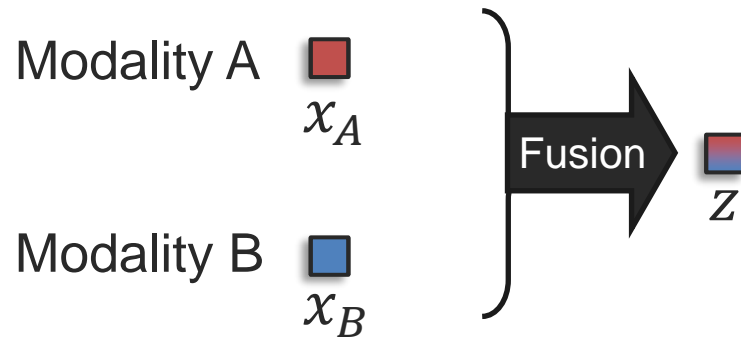
$$z = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_b) + \epsilon$$



	Estimate	95% CI
w_0	5.79	[5.29, 6.29]
w_1	0.68	[0.25, 1.11]
w_2	-2.94	[-3.73, -2.15]
w_3	1.29	[0.61, 1.97]

➔ **Multiplicative interaction!**

Basic Concepts for Representation Fusion (aka, Basic Fusion)



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

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

Linear regression:

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

intercept (bias term) Additive terms Multiplicative term error (residual term)

① Additive terms:

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

② Multiplicative “interaction” term:

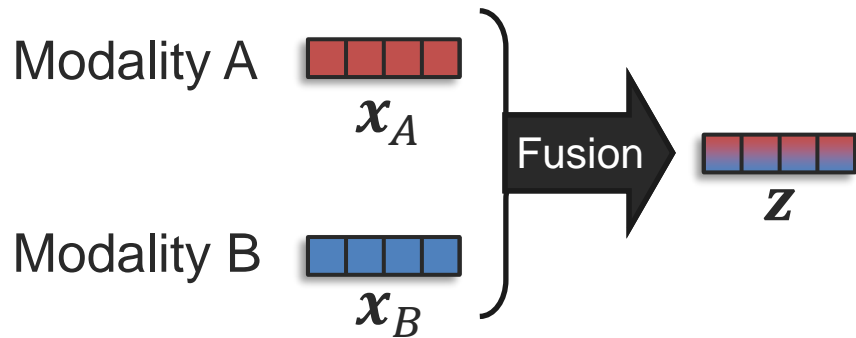
$$Z = w_3 (x_A \times x_B) + \epsilon$$

③ Additive and multiplicative terms:

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

Additive Fusion Back to multivariate case!

 (multi-dimensional features)

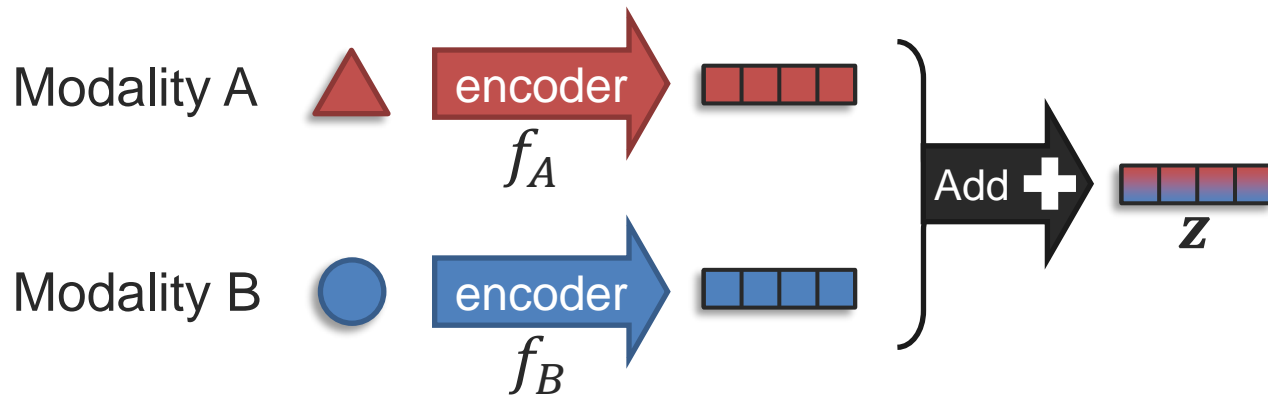


Additive fusion:

$$z = w_1 x_A + w_2 x_B$$


 1-layer neural network
can be seen as additive

With unimodal encoders:

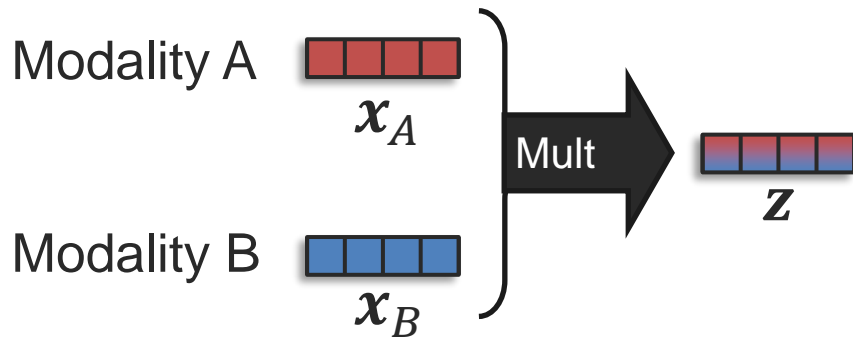


Additive fusion:

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

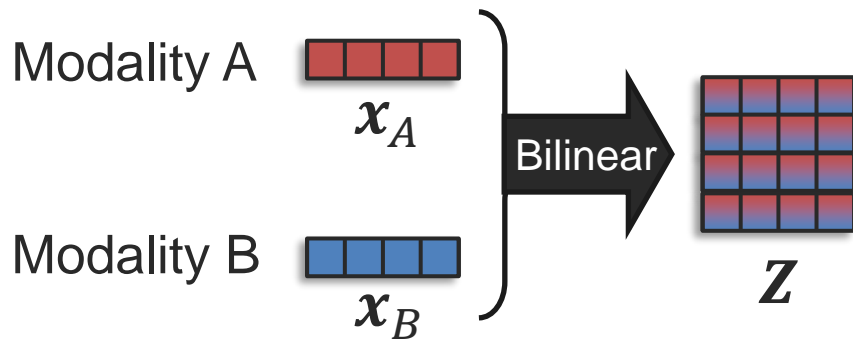
 It could be seen as an
ensemble approach
(late fusion)

Multiplicative Fusion



Simple multiplicative fusion:

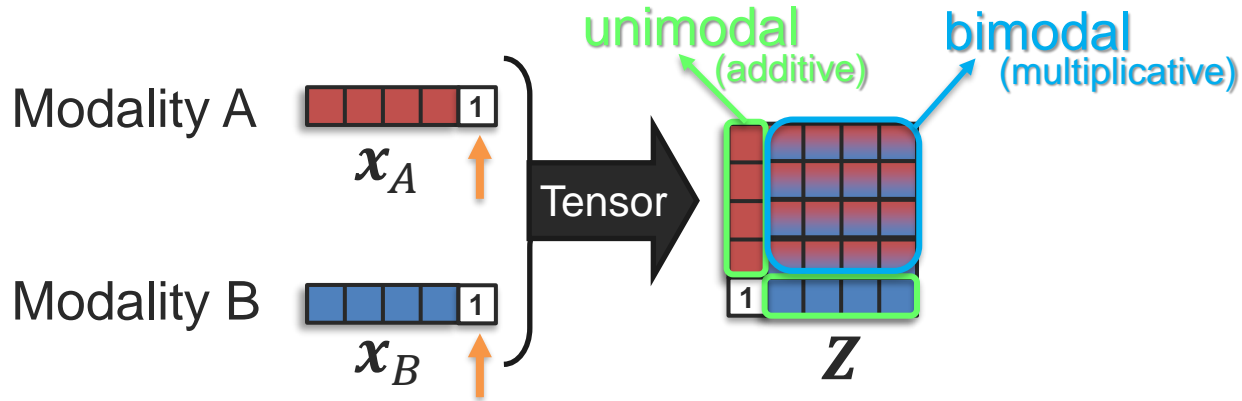
$$z = w(x_A \times x_B)$$



Bilinear Fusion:

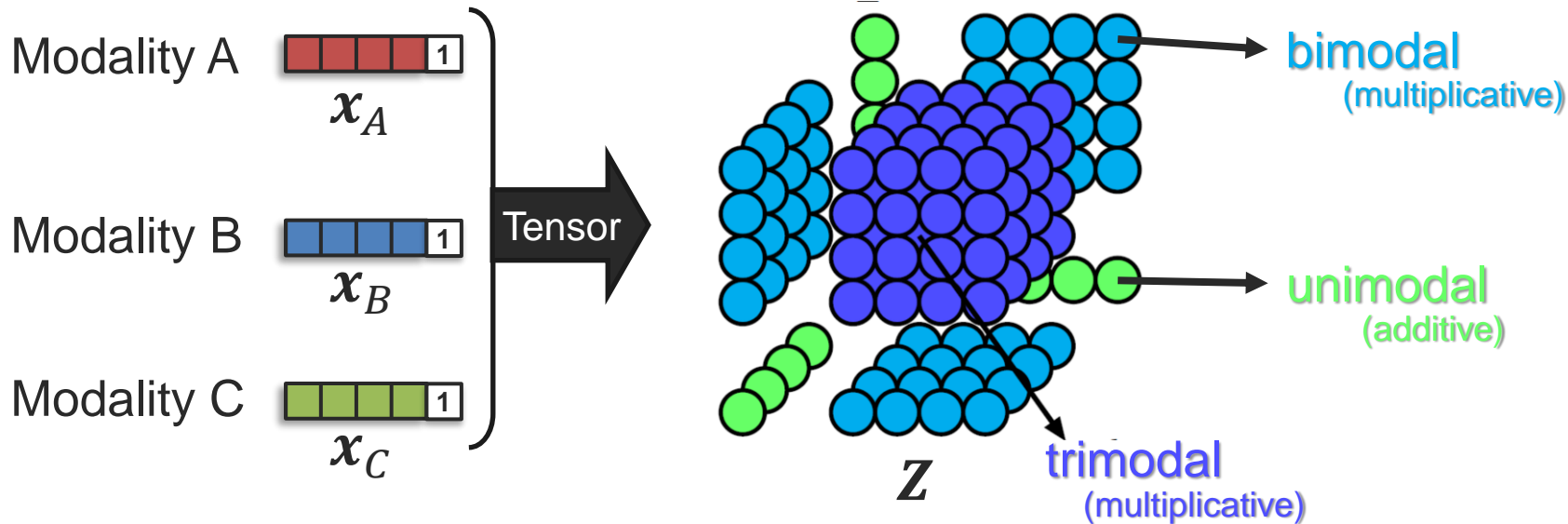
$$Z = W(x_A^T \cdot x_B)$$

Tensor Fusion



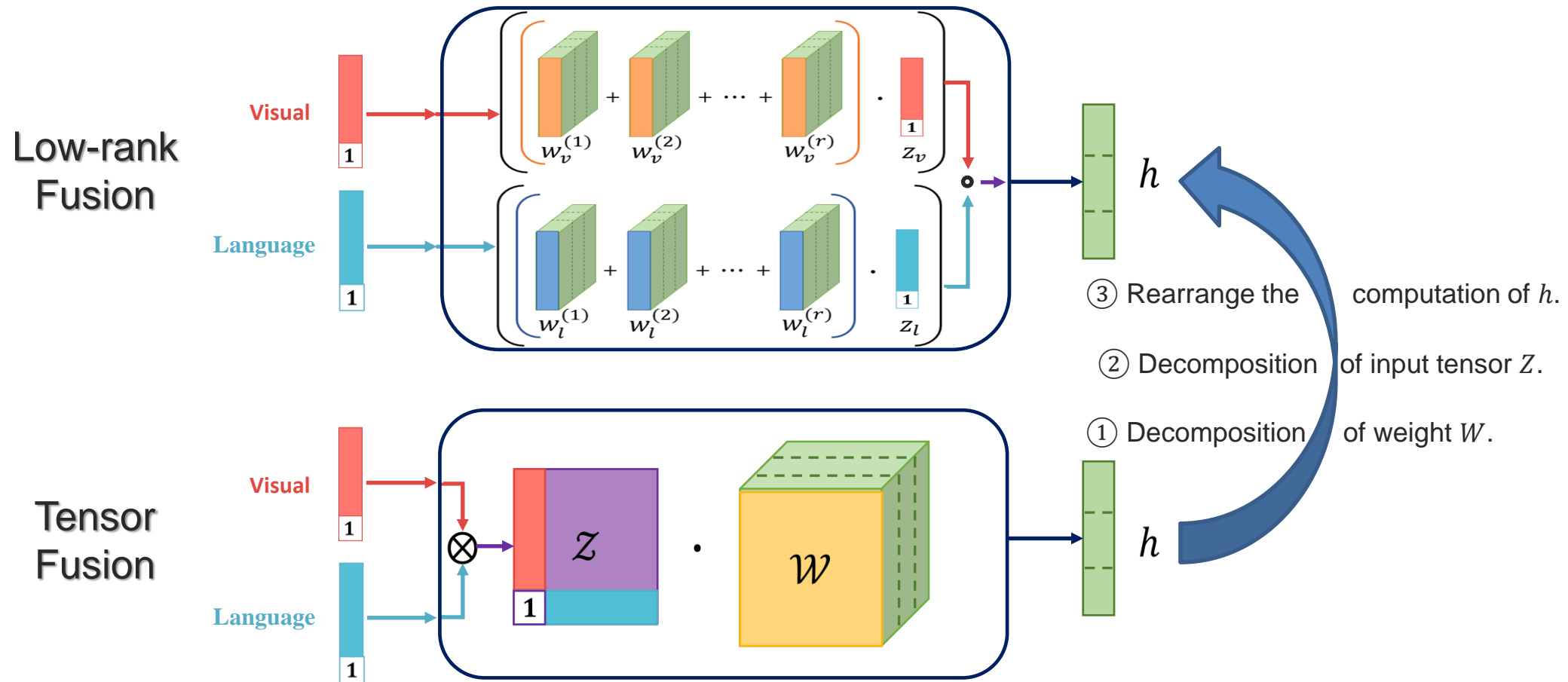
Tensor Fusion (bimodal):

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

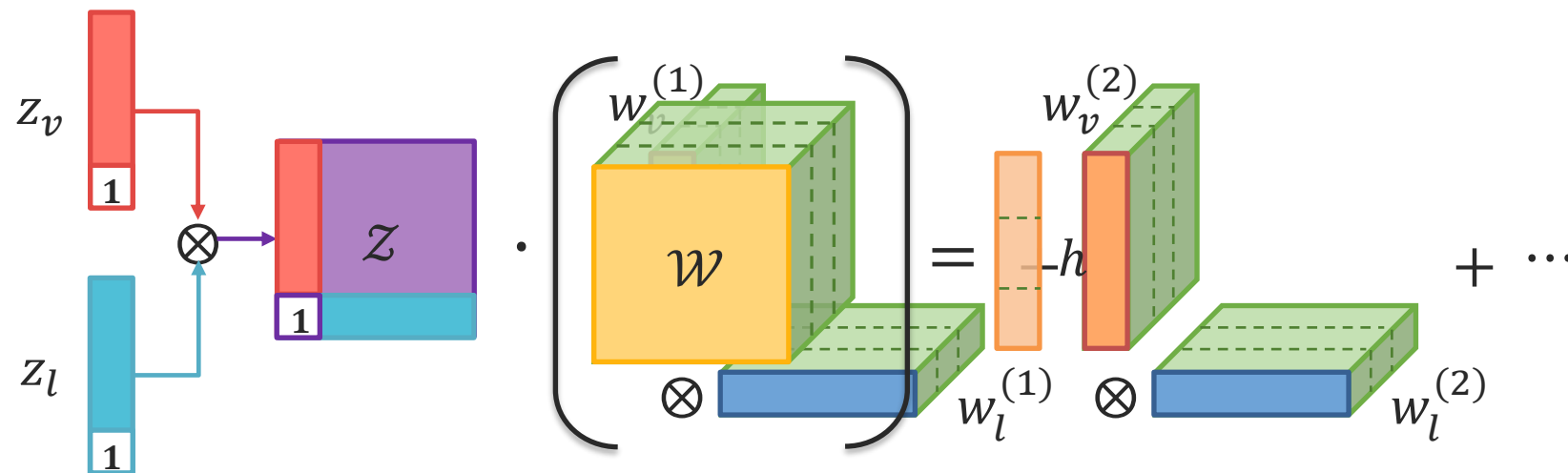


... but the weight matrix may end up quite large!

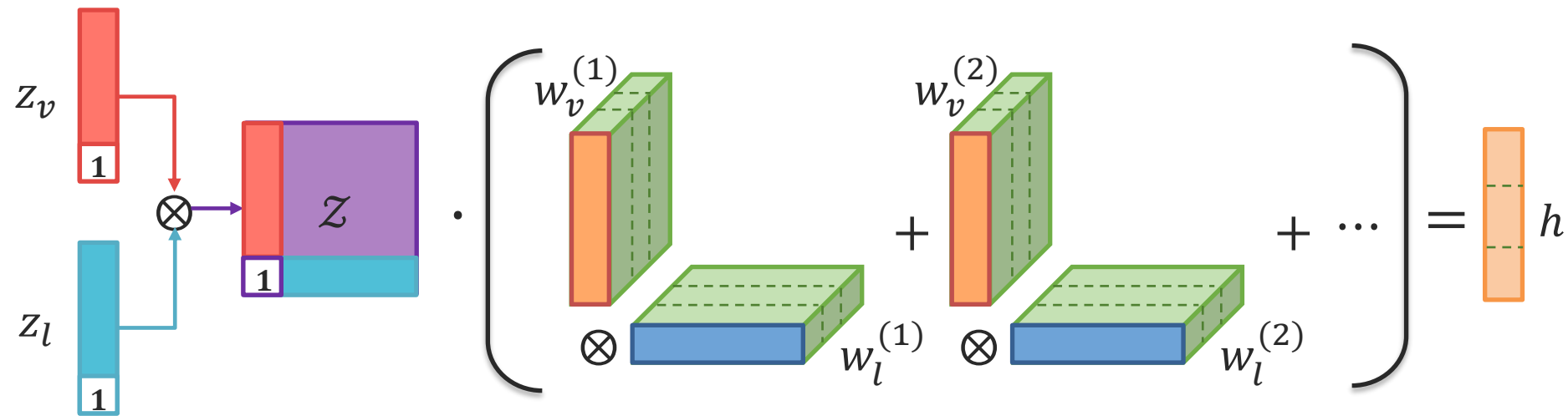
Low-rank Fusion



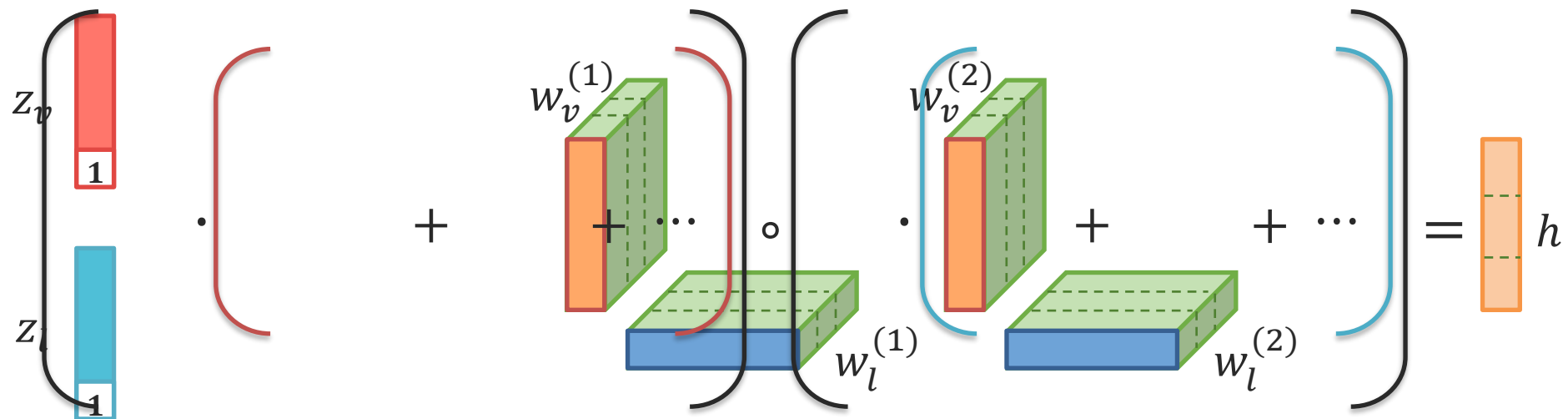
Low-rank Fusion



Low-rank Fusion

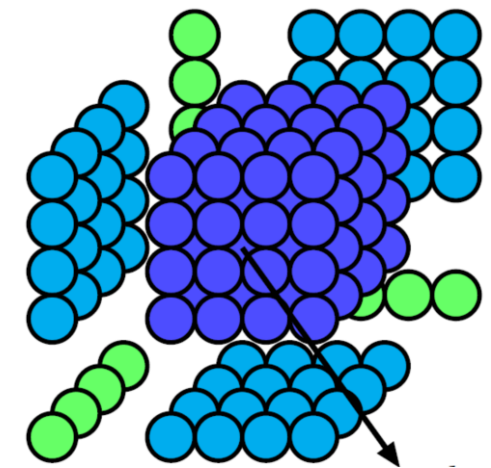


Low-rank Fusion

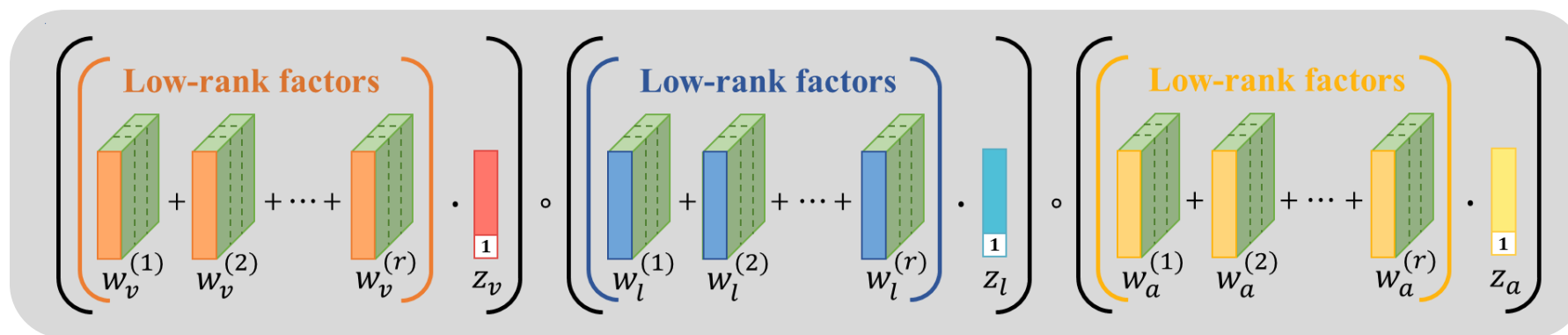


Low-rank Fusion with Trimodal Input

Tensor Fusion



Low-rank Fusion :



Canonical Polyadic Decomposition

Going Beyond Additive and Multiplicative Fusion

Additive interaction:

$$z = w_1x_A + w_2x_B$$

← First-order polynomial

Additive and multiplicative interaction:

$$z = w_1x_A + w_2x_B + w_3(x_A \times x_B)$$

← Second-order polynomial

Trimodal fusion (e.g., tensor fusion):

$$z = \underbrace{w_1x_A + w_2x_B + w_3x_C}_{\substack{\text{Unimodal terms} \\ \text{(first-order)}}} + \underbrace{w_4(x_A \times x_C) + w_5(x_A \times x_C) + w_6(x_B \times x_C)}_{\substack{\text{Bimodal terms} \\ \text{(second-order)}}} + \underbrace{w_7(x_A \times x_B \times x_C)}_{\substack{\text{Trimodal terms} \\ \text{(third-order)}}$$

Can we add higher-order interaction terms?

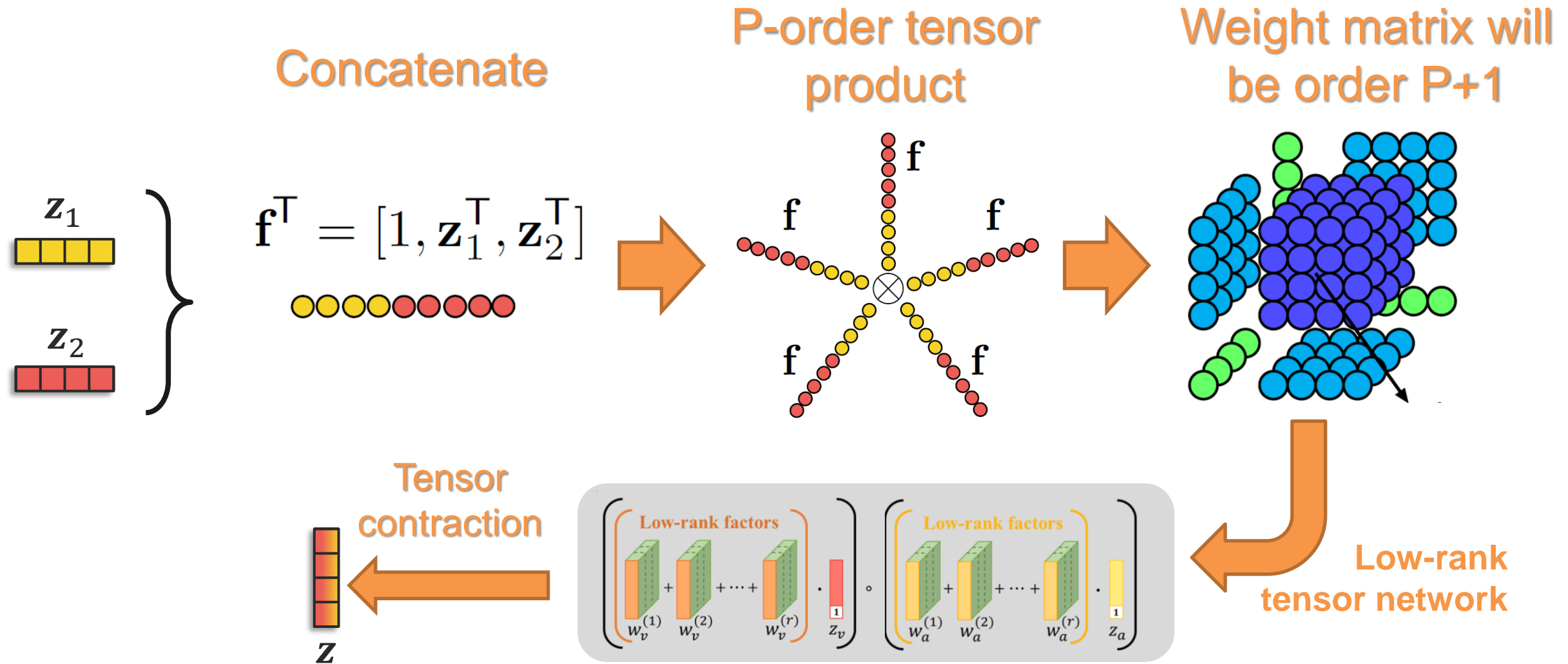
For example:

$$+w_8(x_A^2 \times x_B^2 \times x_C^2)$$

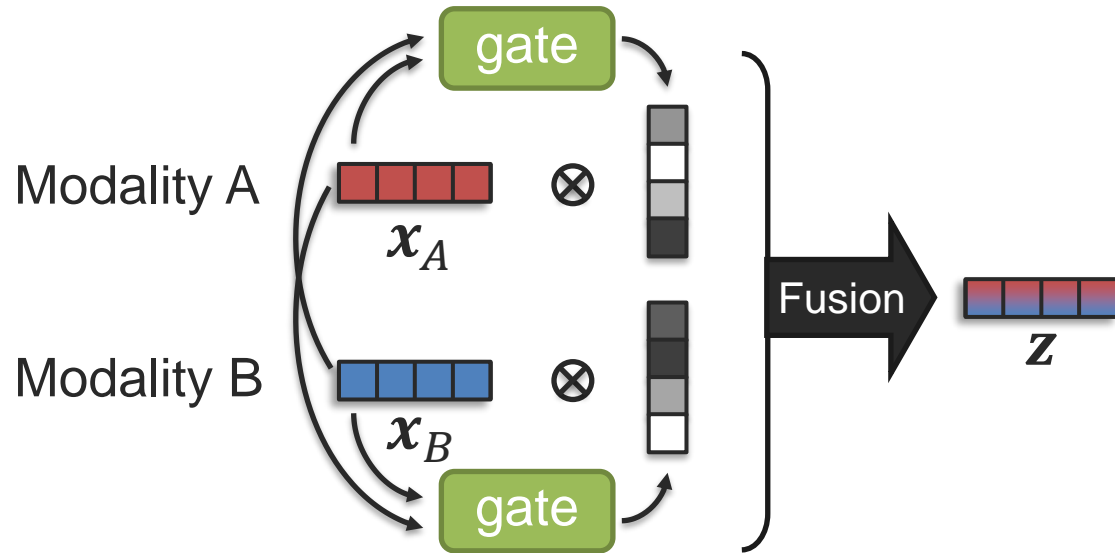
$$+w_9(x_A^3 \times x_B)$$

$$+w_{10}(x_B^3 \times x_C^3)$$

High-Order Polynomial Fusion



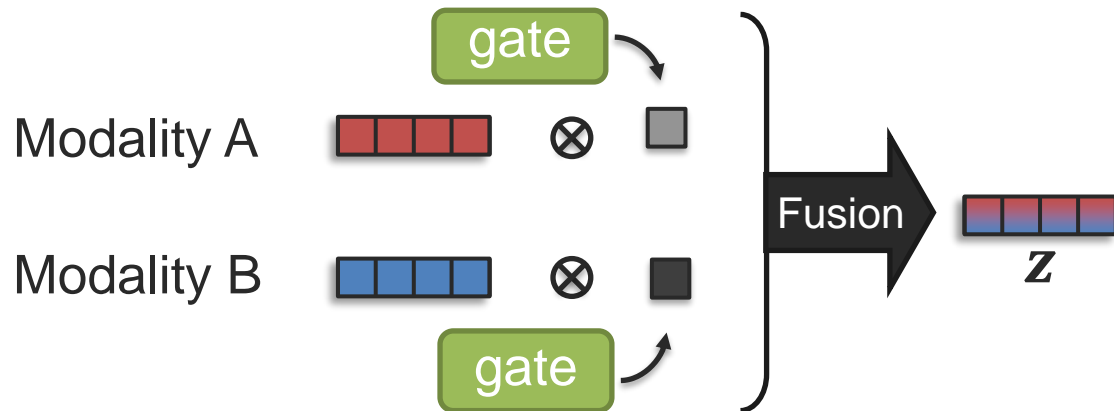
Gated Fusion



Example with additive fusion:

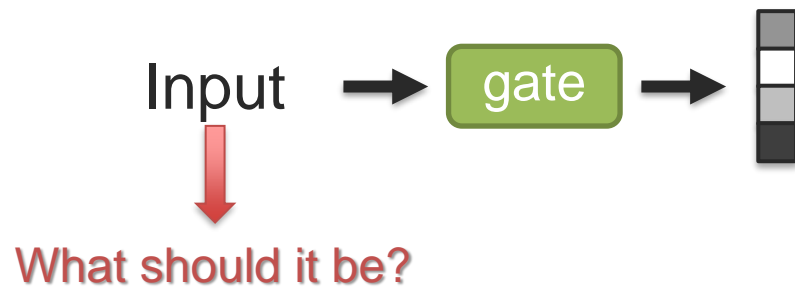
$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

→ g_A and g_B can be seen as attention functions



→ Gating output can be one weight for the whole modality

Gating Module (aka, attention module)



“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 

Other modality 

All modality 



Soft attention



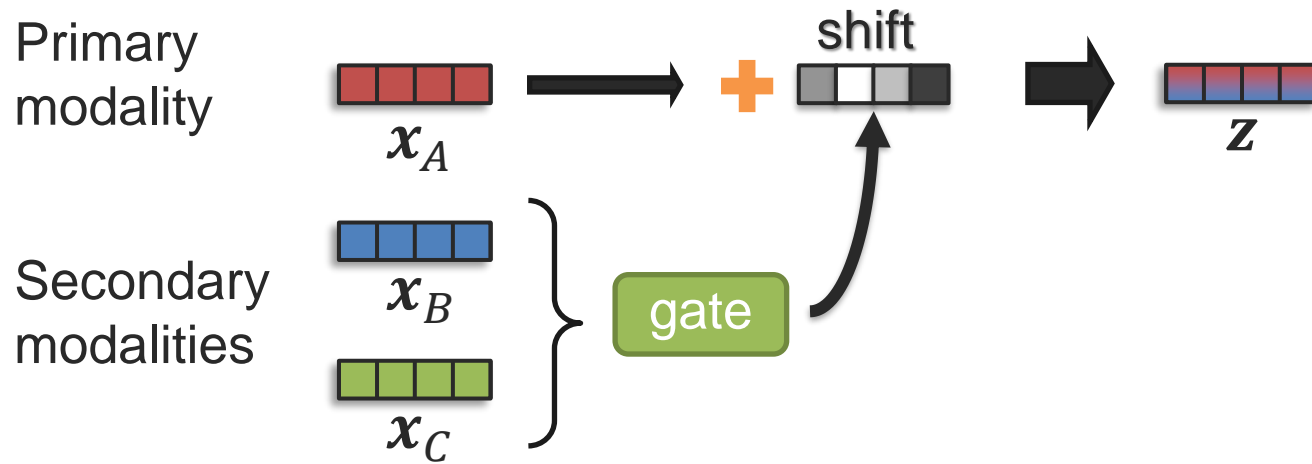
Easier to compute derivative (gradient)

Hard attention



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

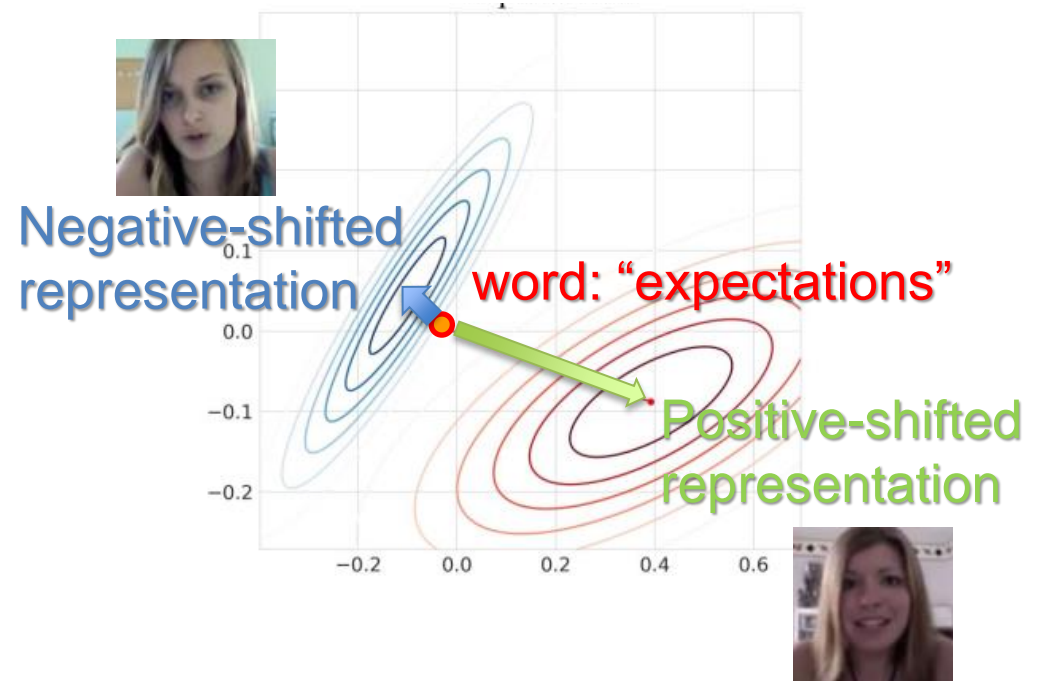
Modality-Shifting Fusion



Example with language modality:

Primary modality: language

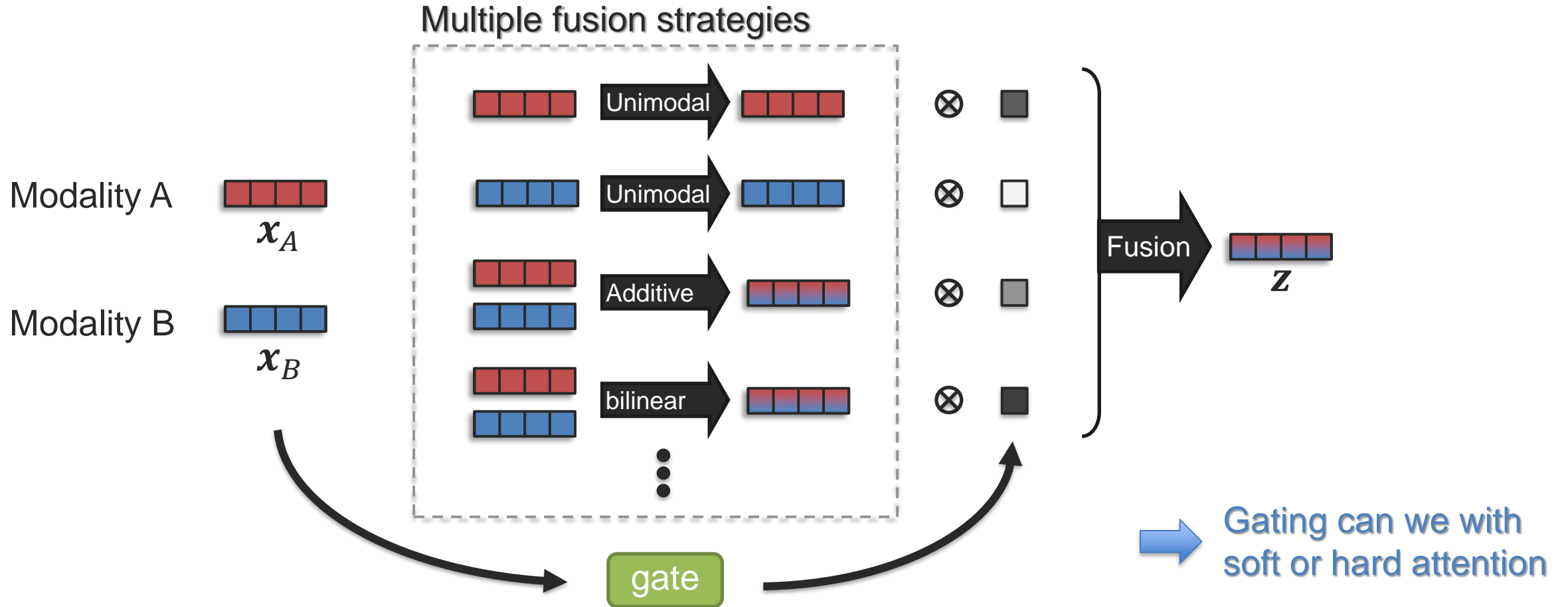
Secondary modalities: acoustic and visual



Wang et al., Words Can Shift: Dynamically Adjusting Word Representations Using Nonverbal Behaviors, AAI 2019

Rahman et al., Integrating Multimodal Information in Large Pretrained Transformers, ACL 2020

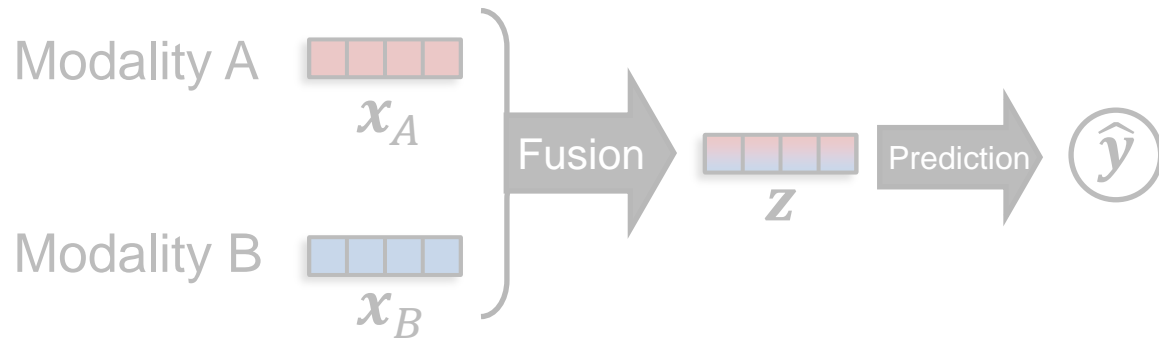
Dynamic Fusion



Zadeh et al., Multimodal Language Analysis in the Wild: CMU-MOSEI Dataset and Interpretable Dynamic Fusion Graph, ACL 2018

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

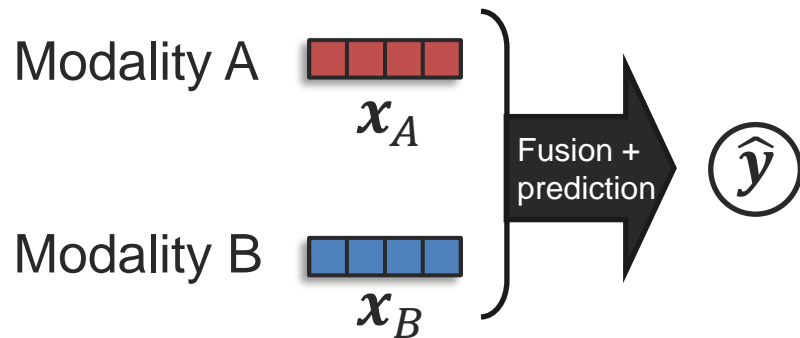
Nonlinear Fusion



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

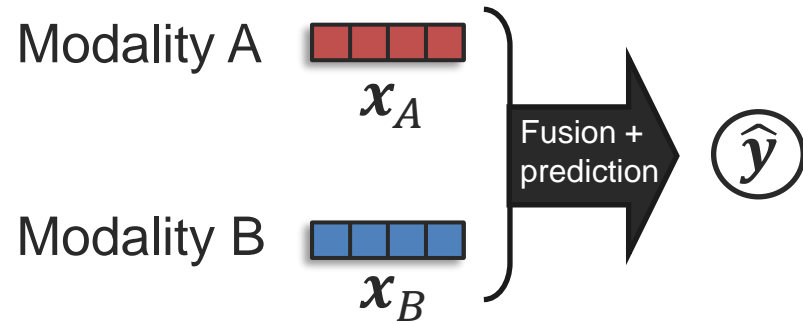


→ This could be seen as *early fusion*:

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

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

Measuring Non-Additive Interactions



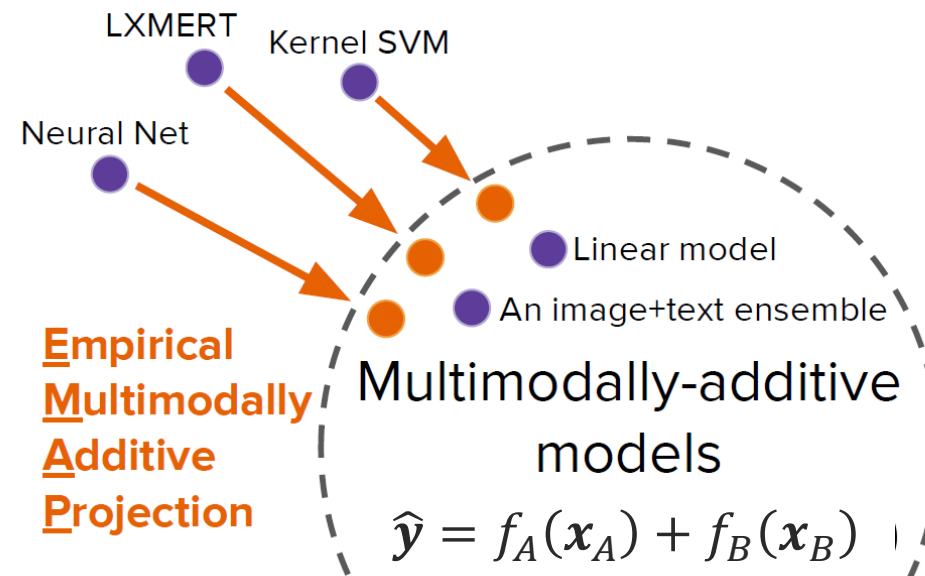
Nonlinear fusion:

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

Projection?

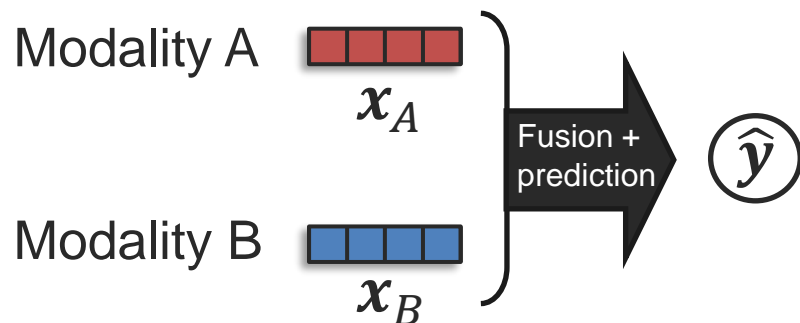
Additive fusion:

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



Hessel and Lee, Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think!, EMNLP 2020 → introduced the EMAP method

Measuring Non-Additive Interactions



Nonlinear fusion:

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

Projection?

Additive fusion:

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

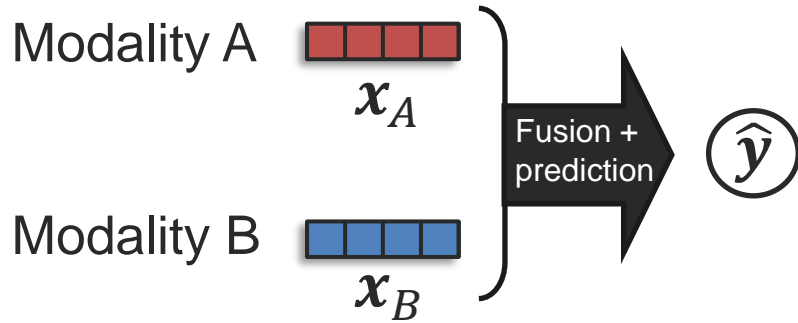
Projection from nonlinear to additive (using EMAP):

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

The expectations \mathbb{E} can be approximated with summation over training data:

$$\hat{f}_A(\mathbf{x}_A) = \frac{1}{N} \sum_{j=1}^N f(\mathbf{x}_{A,j}, \mathbf{x}_{B,j})$$

Measuring Non-Additive Interactions



Nonlinear fusion:

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

EMAP
projection

Additive fusion:

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

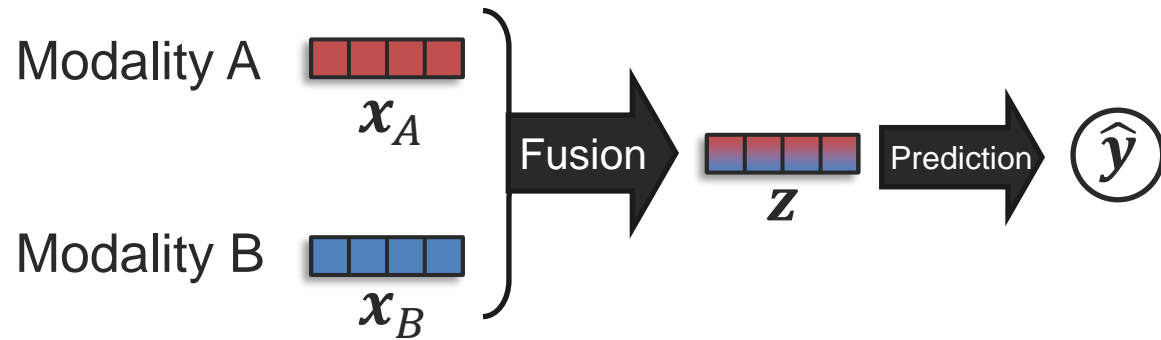
	I-INT	I-SEM	I-CTX	T-VIS	R-POP	T-ST1	T-ST2
Nonlinear ← Neural Network	90.4	69.2	78.5	51.1	63.5	71.1	79.9
Polynomial ← Polykernel SVM	91.3	74.4	81.5	50.8	–	72.1	80.9
Nonlinear ← FT LXMERT	83.0	68.5	76.3	53.0	63.0	66.4	78.6
Nonlinear ← ↵ + Linear Logits	89.9	73.0	80.7	53.4	64.1	75.5	80.3
Additive ← Linear 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!!!

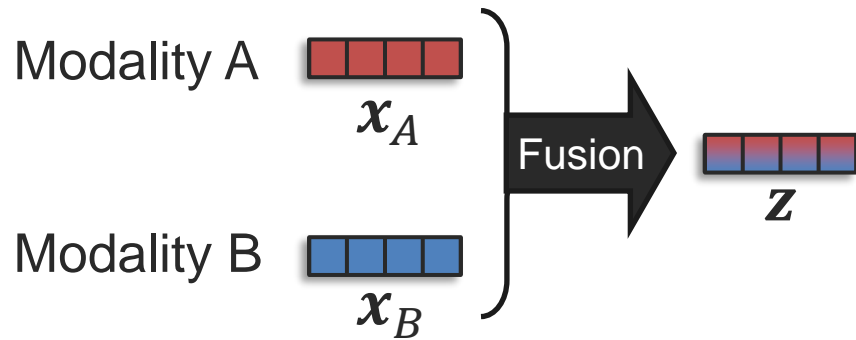
Hessel and Lee, Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think!, EMNLP 2020 → introduced the EMAP method

Learning Fusion Representations



How to learn fusion models?

Learning Fusion Representations

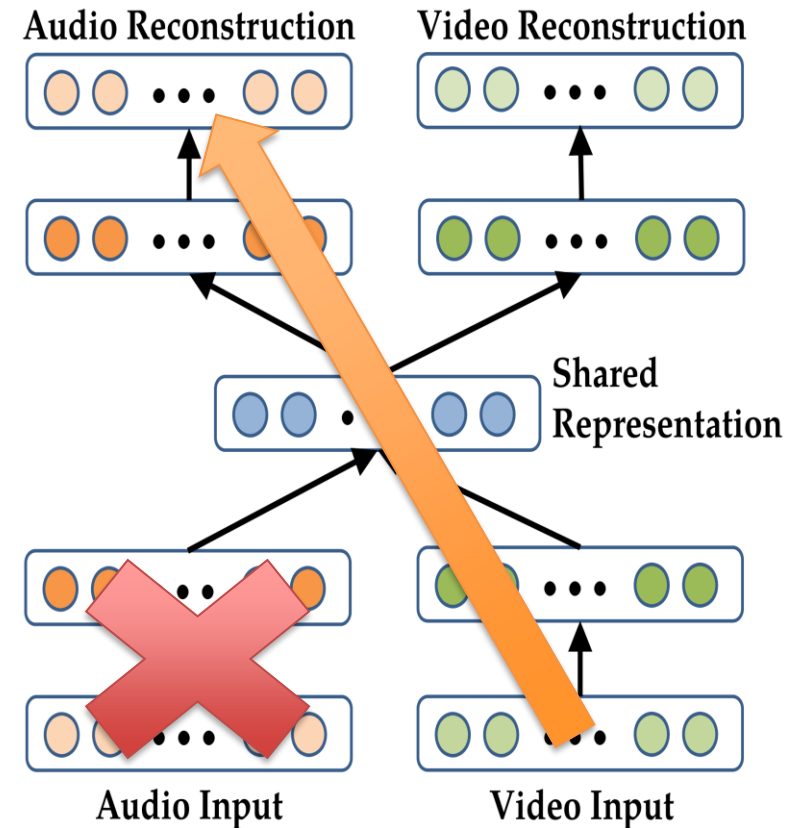


How to learn fusion models?

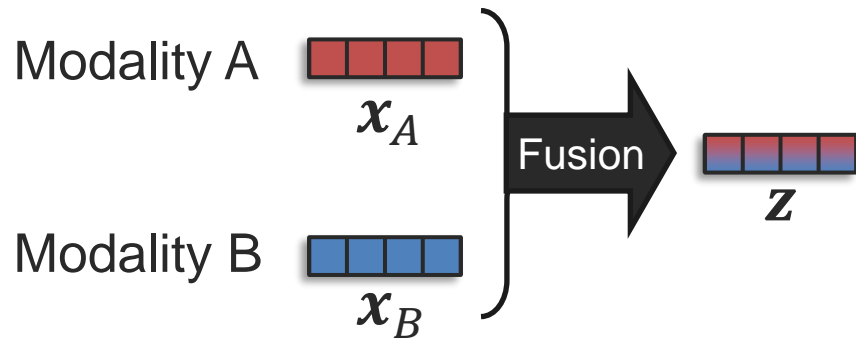
What will be the loss function?

Can it hallucinate the other modality?

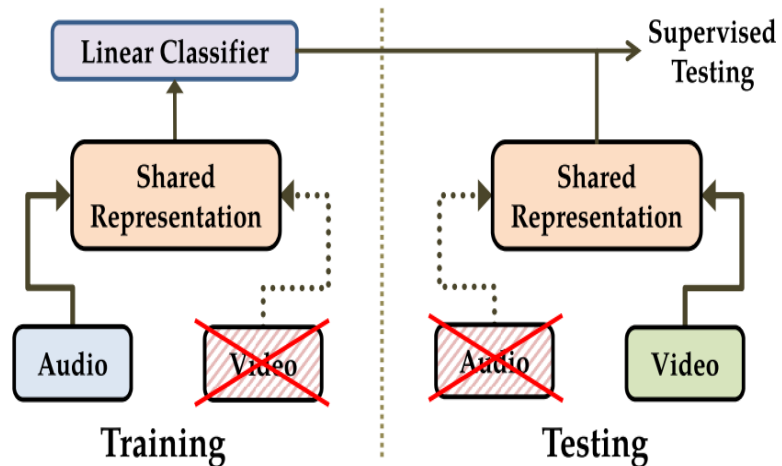
Multimodal Autoencoder



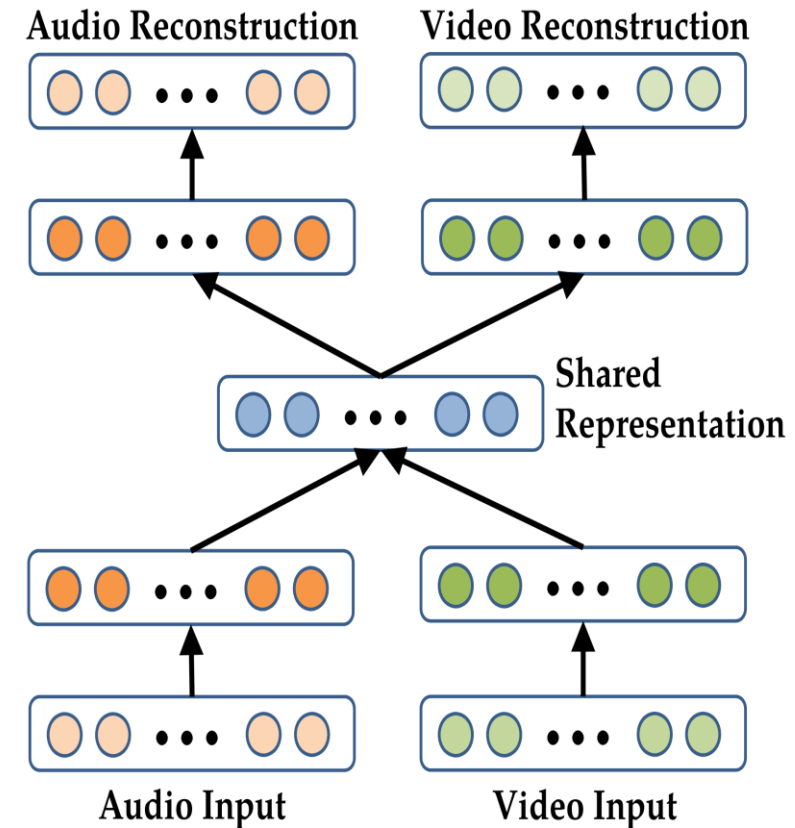
Learning Fusion Representations



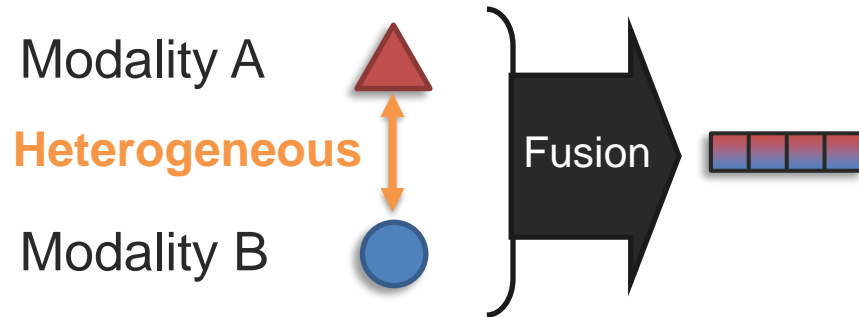
Interesting experiment: "Hearing to see"



Multimodal Autoencoder

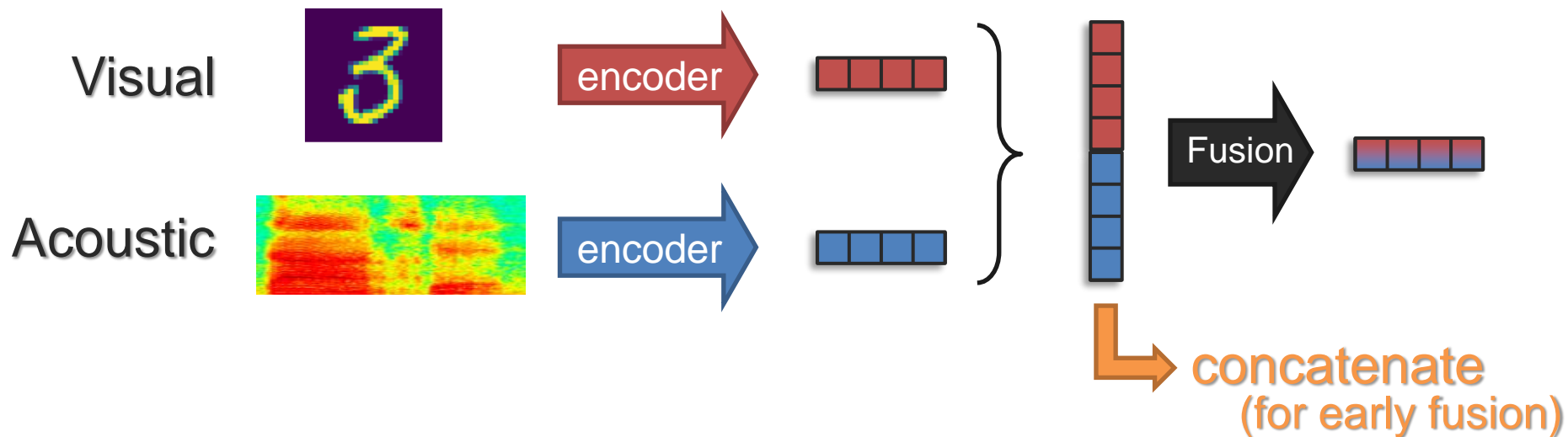


Fusion with Raw Modalities

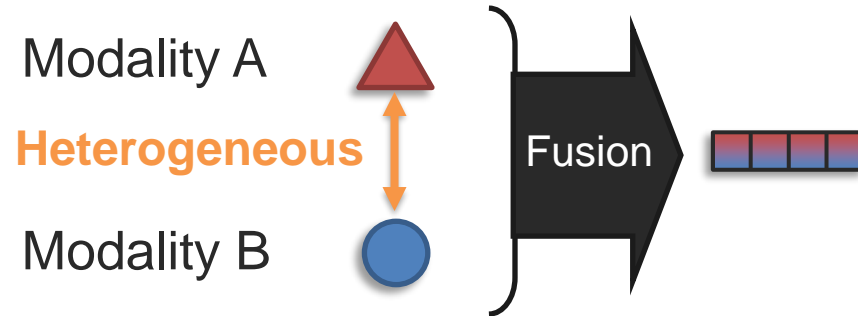


Open Challenge!

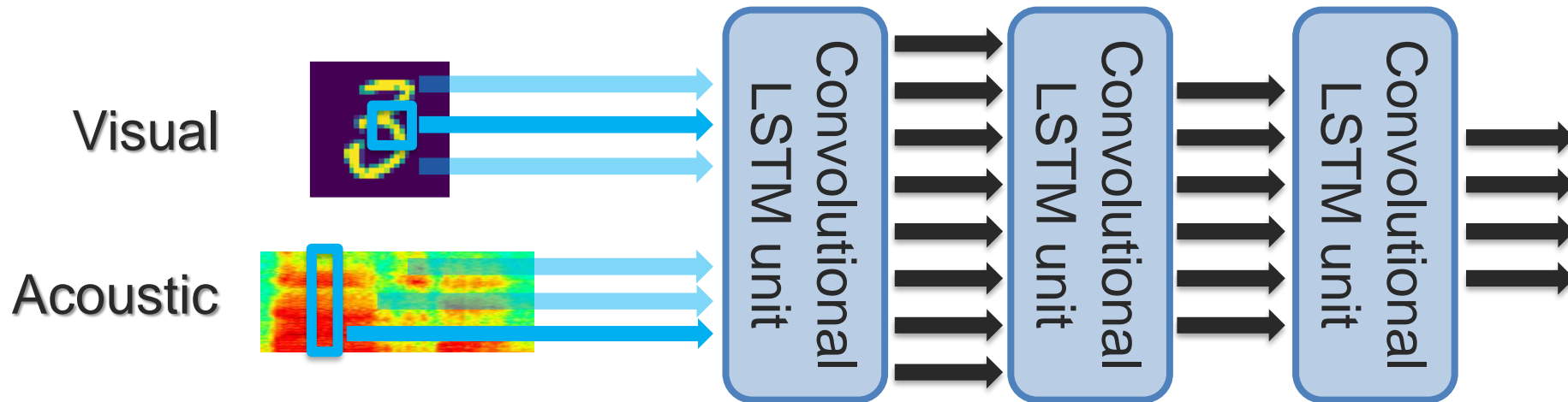
Example: From Early Fusion...



Fusion with Raw Modalities

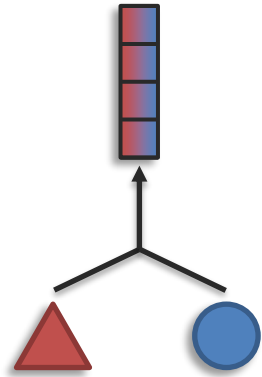


Example: From Early Fusion... to Very Early Fusion (inspired by human brain)



Barnum, et al. "On the Benefits of Early Fusion in Multimodal Representation Learning." *arxiv* 2022

Sub-Challenge 1a: Representation Fusion



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

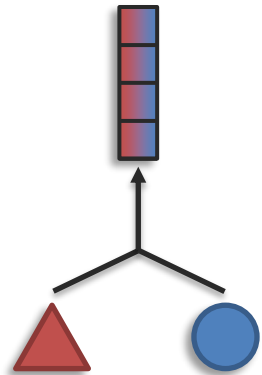


Challenge 1: Representation

Definition: Learning representations that reflect cross-modal interactions between individual elements, across different modalities

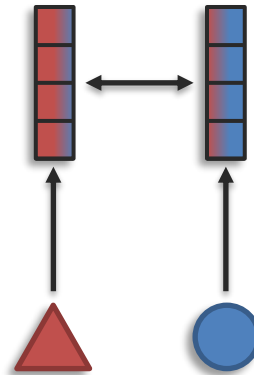
Sub-challenges:

Fusion



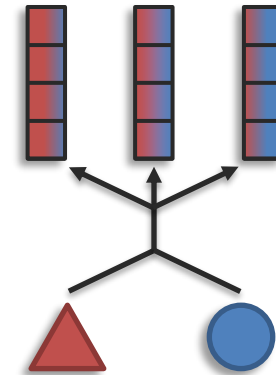
modalities \gt # representations

Coordination



modalities = # representations

Fission



modalities \lt # representations

Next week