

# Multimodal AI

## Lecture 4.1 – Multimodal Fusion

**Paul Liang**

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

[ppliang@mit.edu](mailto:ppliang@mit.edu)

 [@pliang279](https://twitter.com/pliang279)



# Assignments for This Coming Week

Project proposal due today. Meet with me at 4-5pm if you want feedback.

I want to meet every group at least once regarding their project ideas either today or this Thursday.

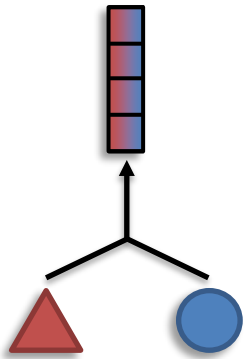
Compute credits: 40 x \$50 Kimi credits, 40 x \$40 other credits.

HW2 released last week, due next Wednesday 3/4.

# Today's lecture

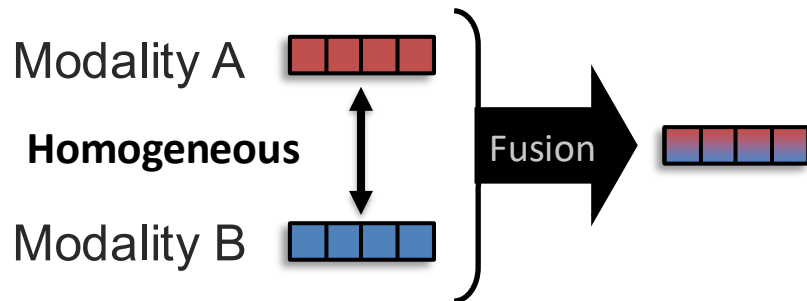
- 1 Basics of multimodal interactions and fusion
- 2 Early, intermediate, late fusion
- 3 Multiplicative and dynamic fusion
- 4 Quantifying fusion

# Multimodal Fusion

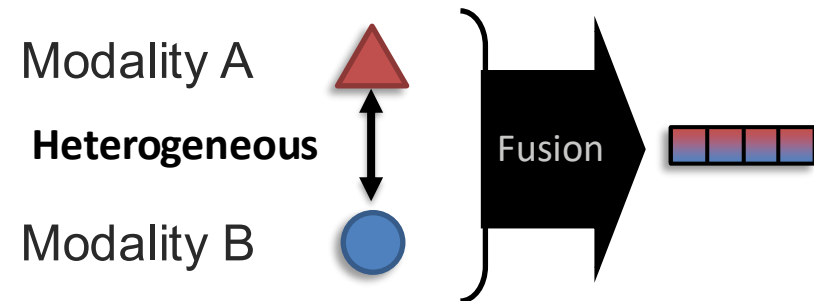


**Definition:** Learn a joint representation that models the **multimodal interactions** between individual elements of different modalities.

Fusion with abstract modalities:

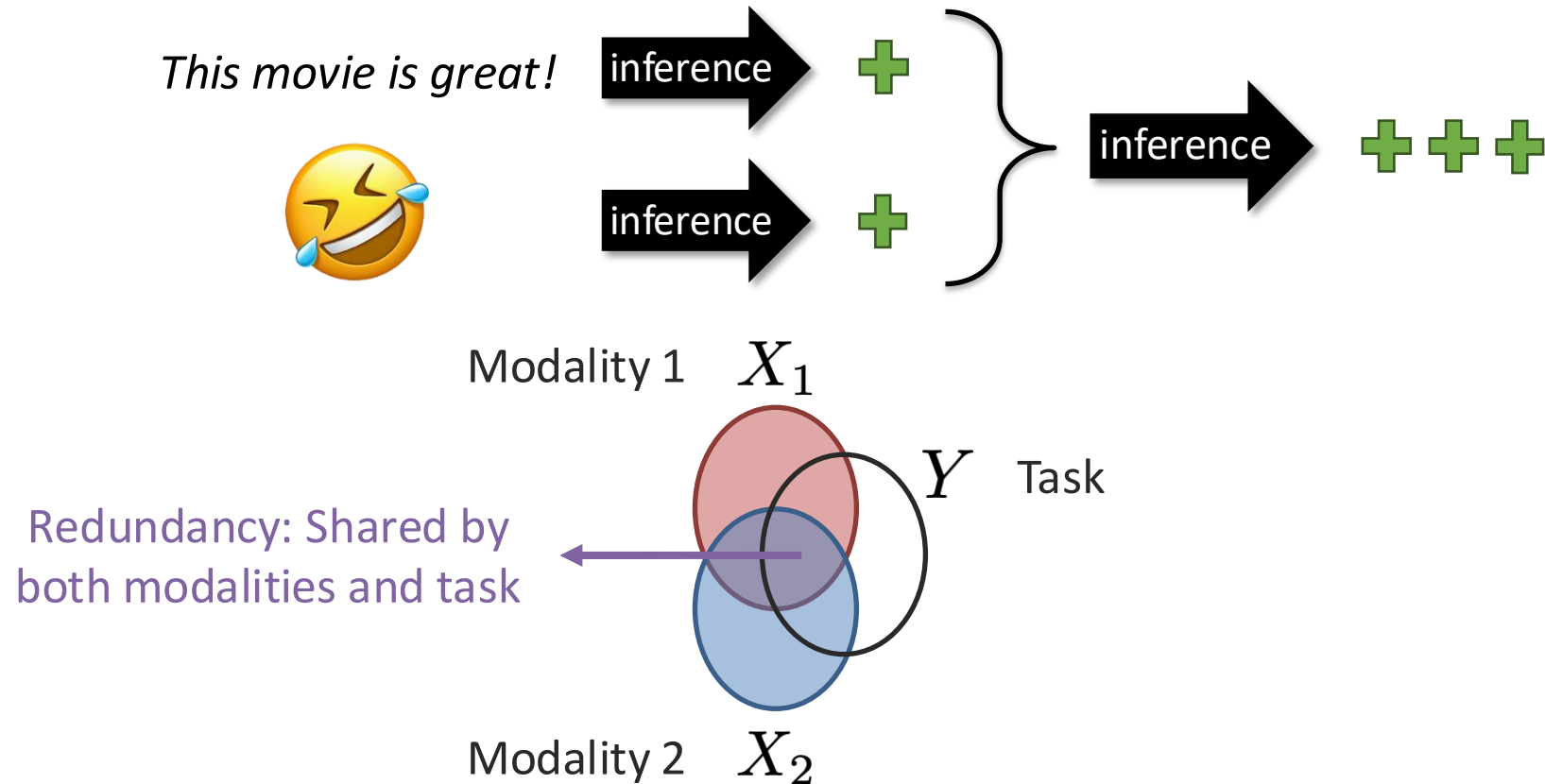


Fusion with raw modalities:



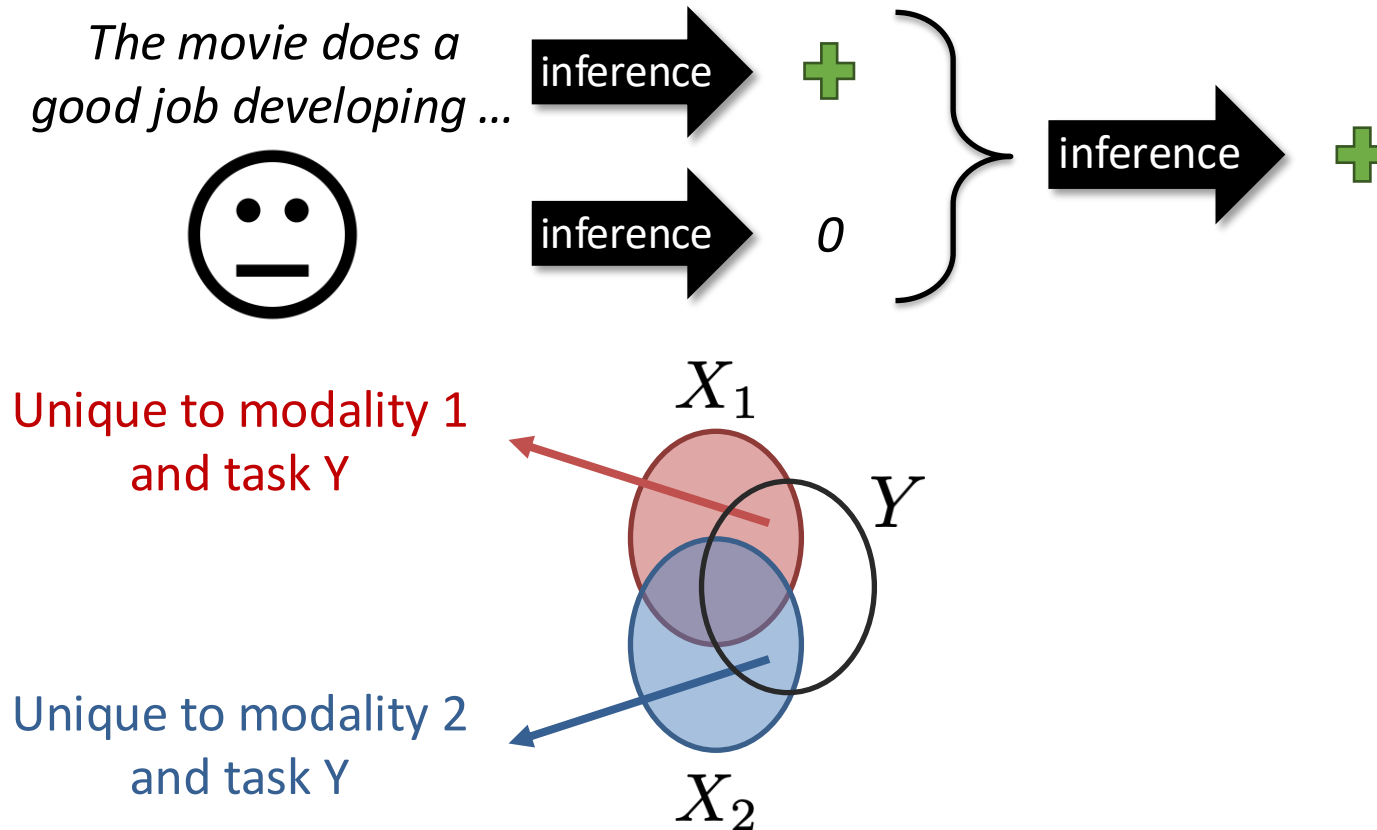
# Multimodal Interactions

**Interactions:** How modalities *combine* to provide information for a task.



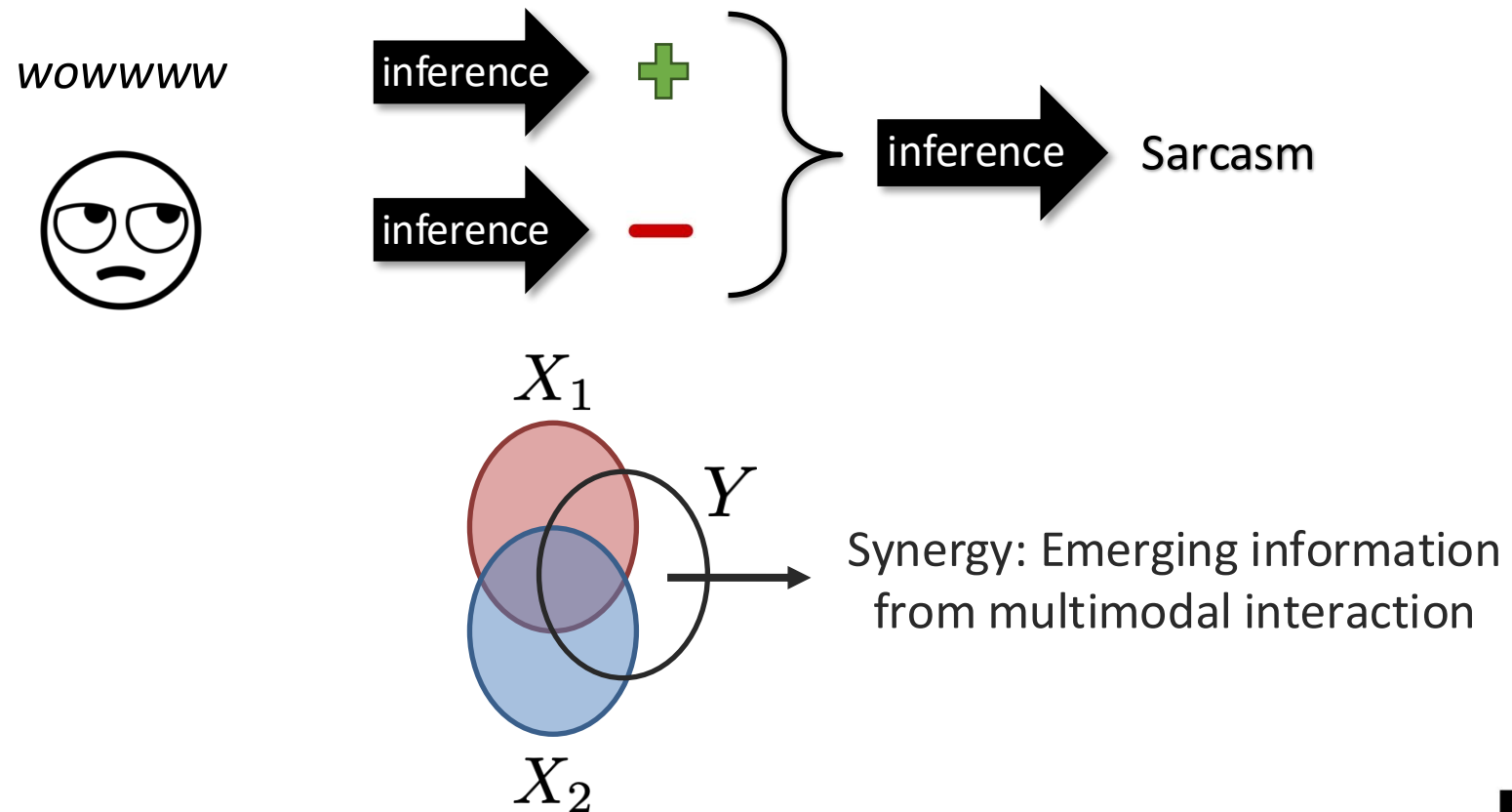
# Multimodal Interactions

**Interactions:** How modalities *combine* to provide information for a task.

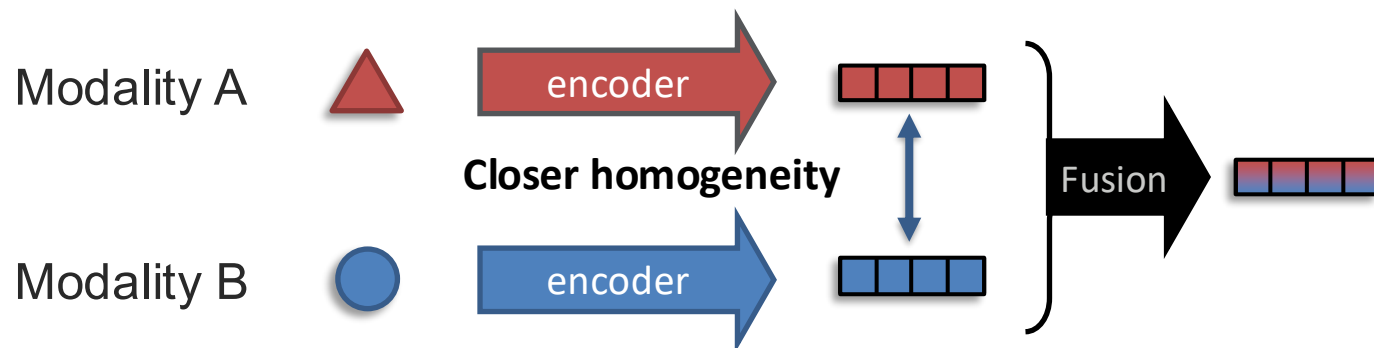


# Multimodal Interactions

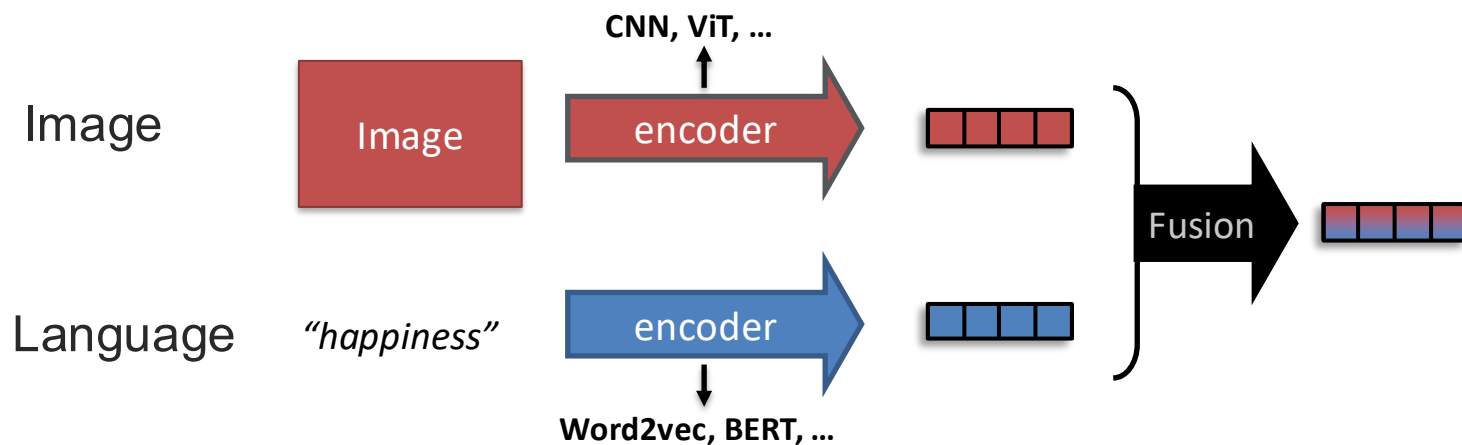
**Interactions:** How modalities *combine* to provide information for a task.



# Fusion with Abstract Modalities



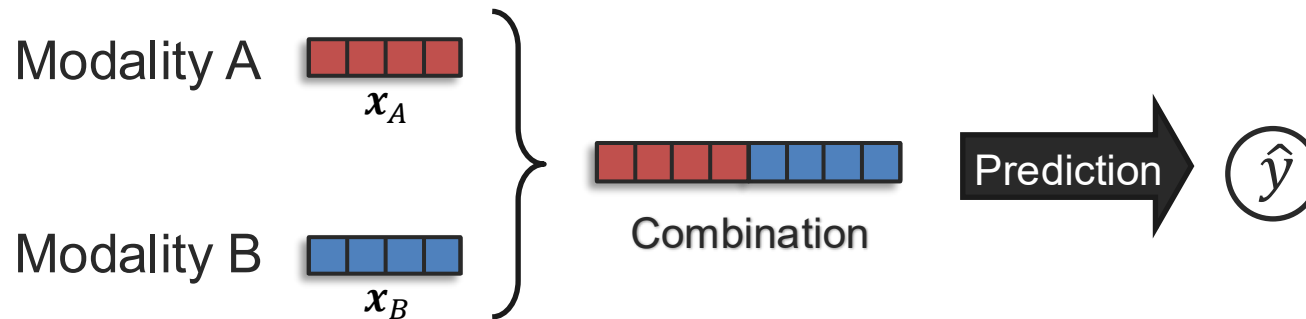
Example:



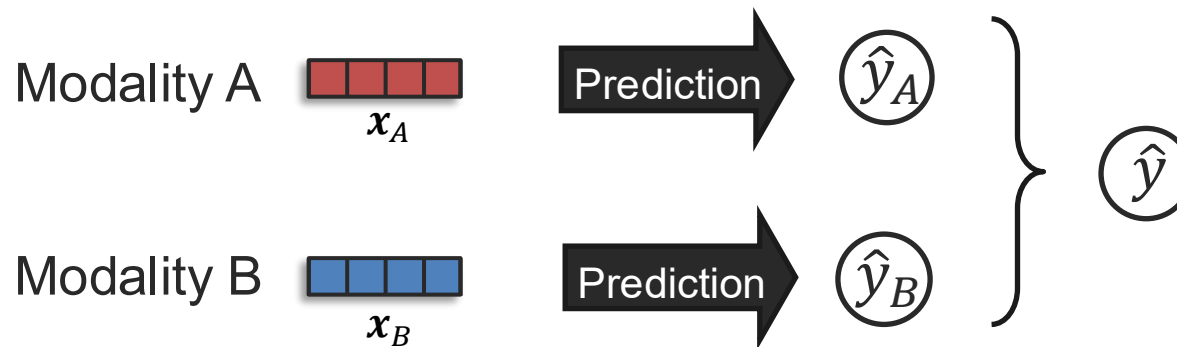
➔ Unimodal encoders can be jointly learned with fusion network, or pre-trained

# Early and Late Fusion

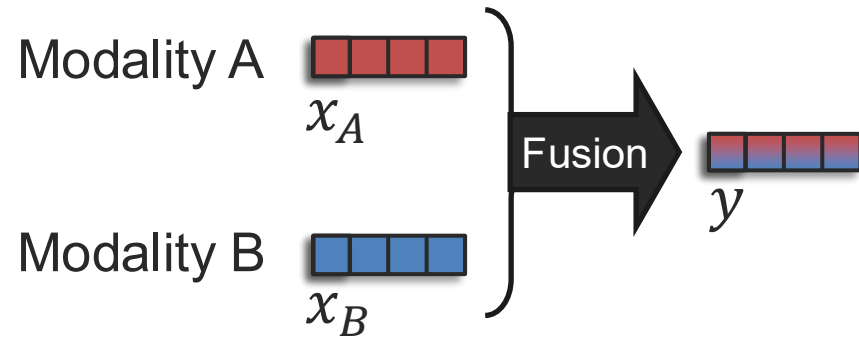
Early fusion:



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)

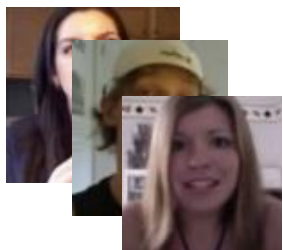
Additive  
 terms

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)

Linear regression:

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

↓
Additive
Multiplicative
↓
error

(bias term)
terms
term
(residual term)

$w_0$ : average score when  $x_A$  and  $x_B$  are zero

$w_1$ : effect from  $x_A$  variable only

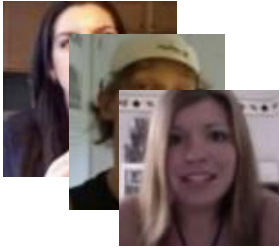
$w_2$ : effect from  $x_B$  variable only

$w_3$ : effect from  $x_A$  and  $x_B$  interaction only

$\epsilon$ : 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

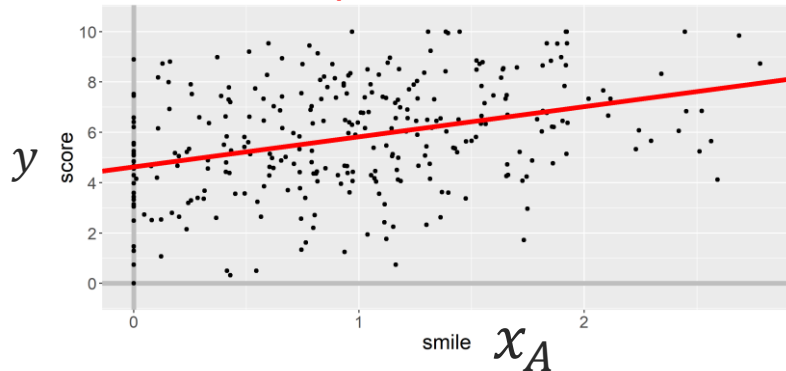
$x_B$ : professional status

(0=non-critic, 1=critic)

Linear regression:

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

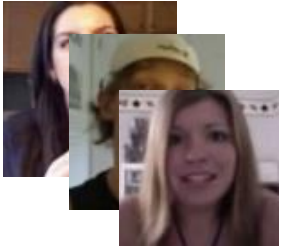
slope



	Estimate
$w_0$	4.63
$w_1$	1.20

# Linear Fusion Case

300 book reviews



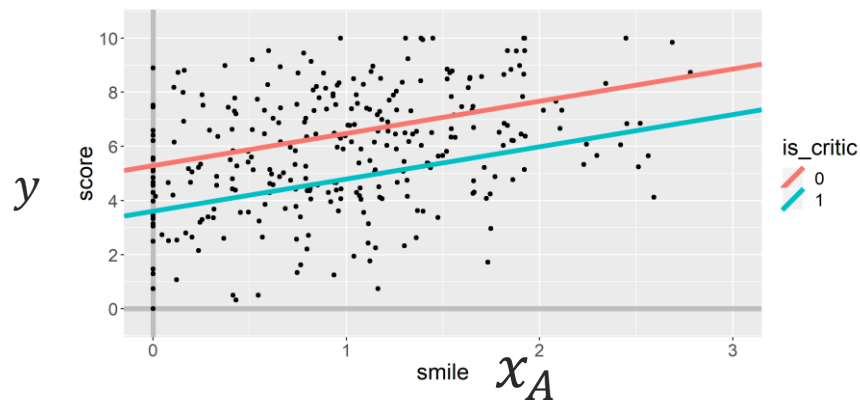
$y$ : audience score

$x_A$ : percentage of smiling

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

Linear regression:

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



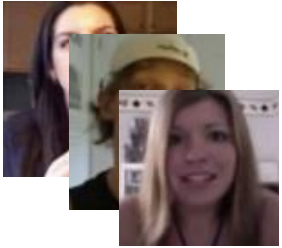
	Estimate
$w_0$	5.29
$w_1$	1.19
$w_2$	-1.69

→ Positive effect

→ Negative effect

# Linear Fusion Case

300 book reviews



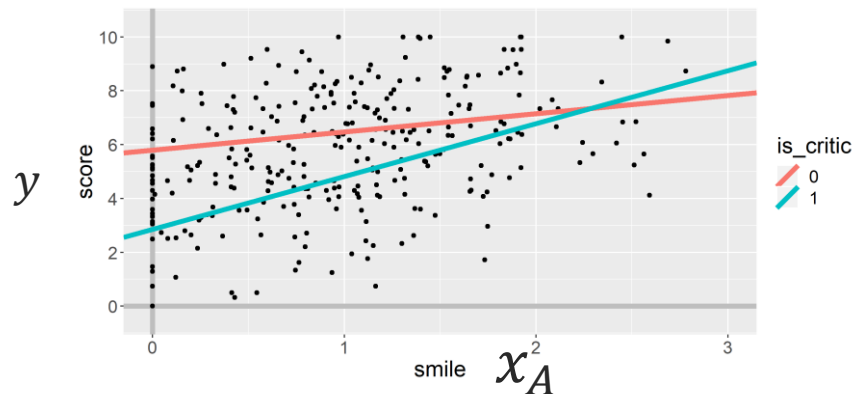
$y$ : audience score

$x_A$ : percentage of smiling

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

Linear regression:

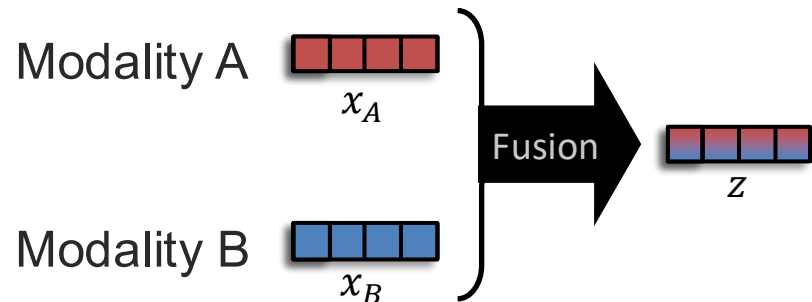
$$y = w_0 + w_1 x_A + w_2 x_B + w_3 (x_A \times x_B) + \epsilon$$



	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:

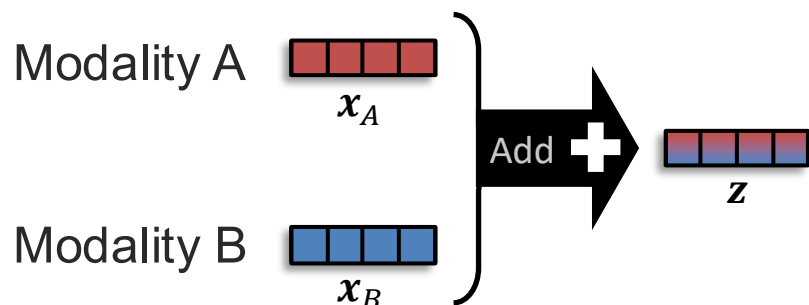
$$z = w_3 (x_A \times x_B) + \epsilon$$

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!

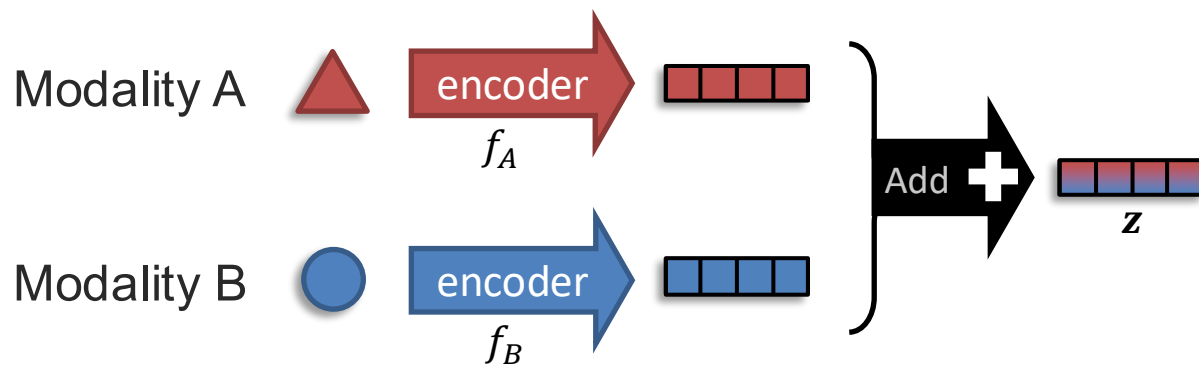
 (multi-dimensional features)



Additive fusion:


$$z = w_1 x_A + w_2 x_B$$

With unimodal encoders:

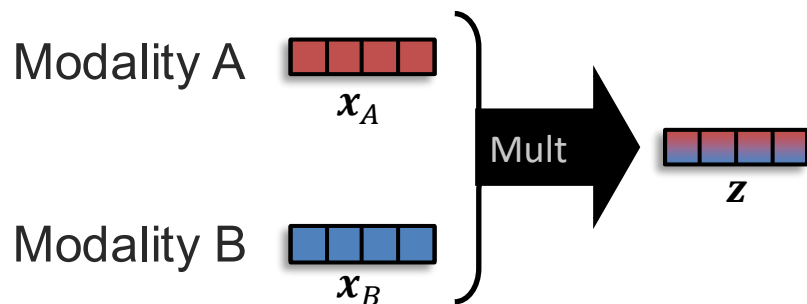


Additive fusion:

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

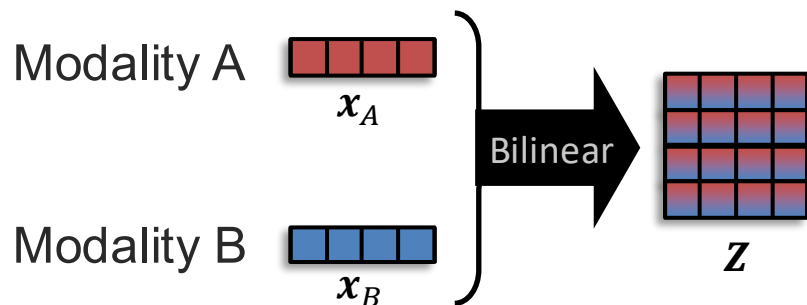
 It could be seen as an ensemble approach  
(late fusion)

# Multiplicative Fusion



Multiplicative fusion:

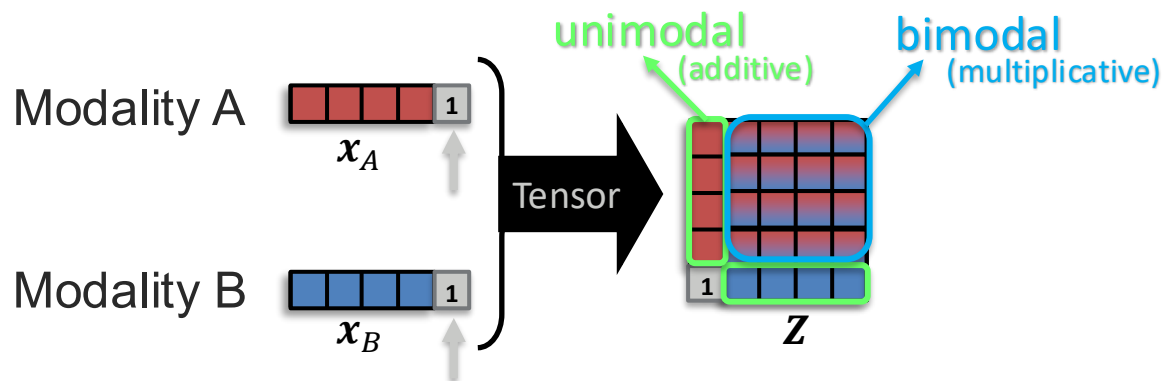
$$z = w(x_A \times x_B)$$



Bilinear Fusion:

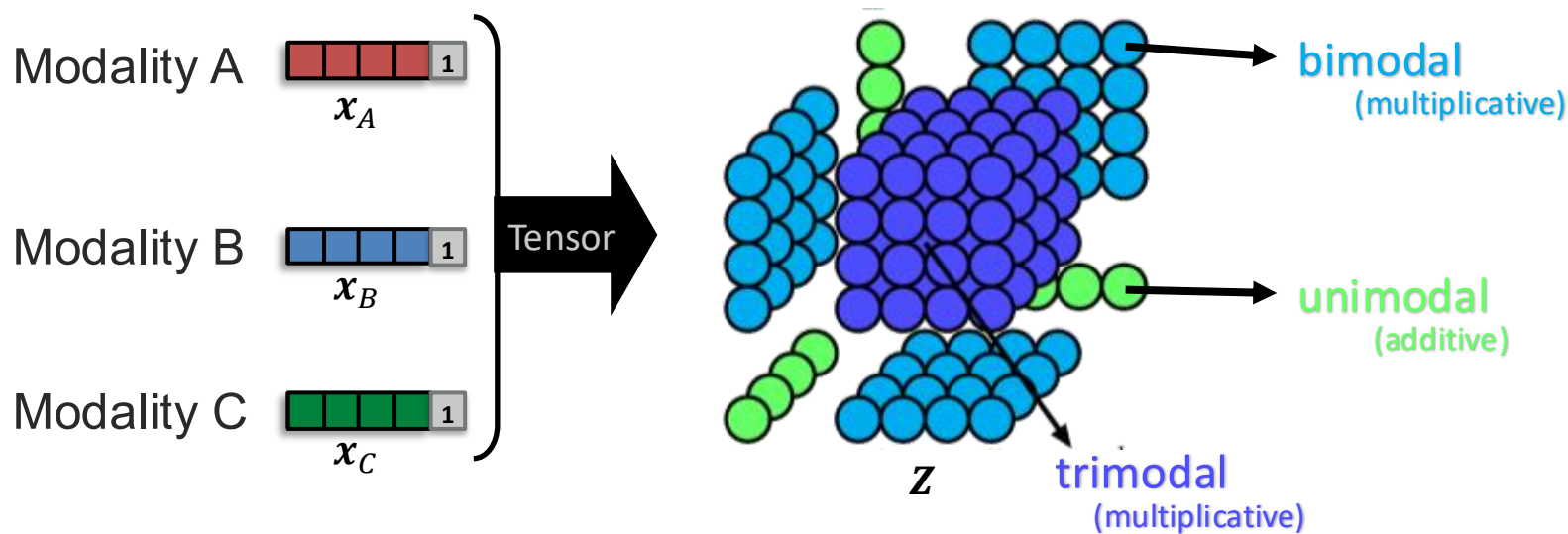
$$Z = w(x_A^T x_B)$$

# Tensor Fusion



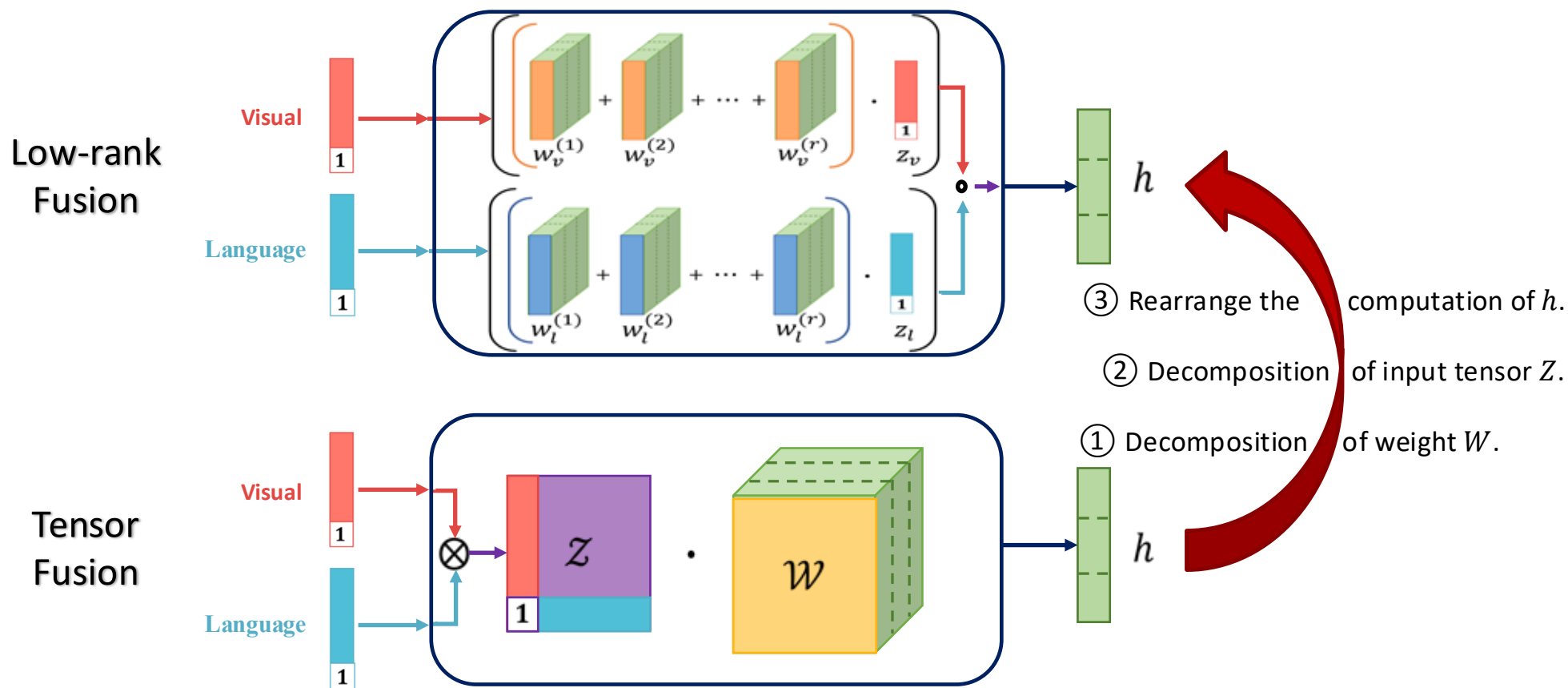
Tensor Fusion (bimodal):

$$Z = w([\mathbf{x}_A \ 1]^T \cdot [\mathbf{x}_B \ 1])$$

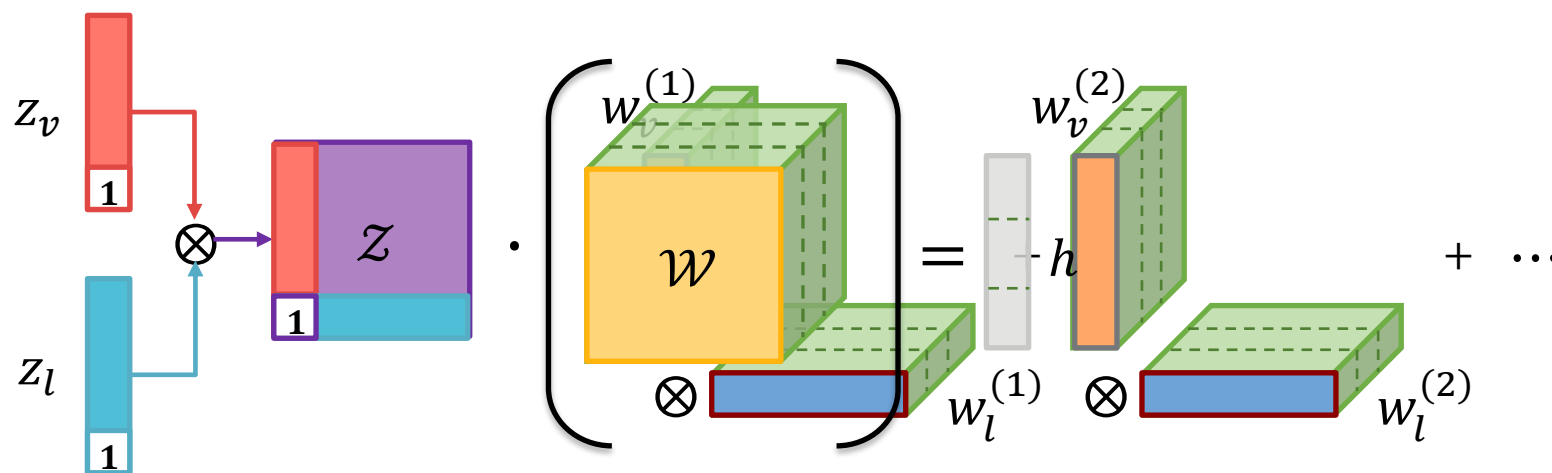


... but the weight matrix may end up quite large!

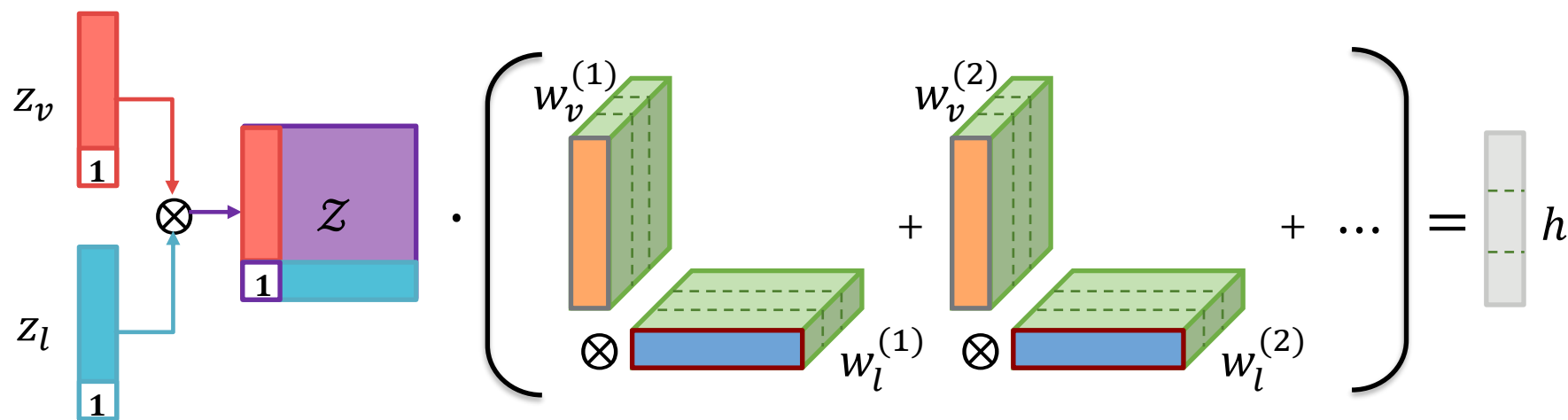
# Low-rank Fusion



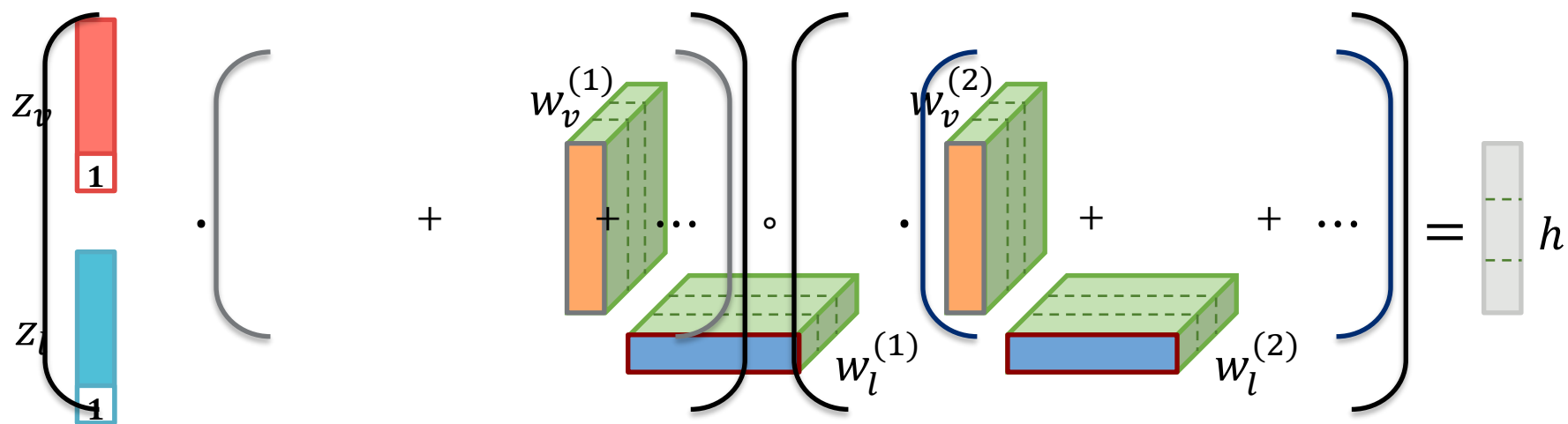
# Low-rank Fusion



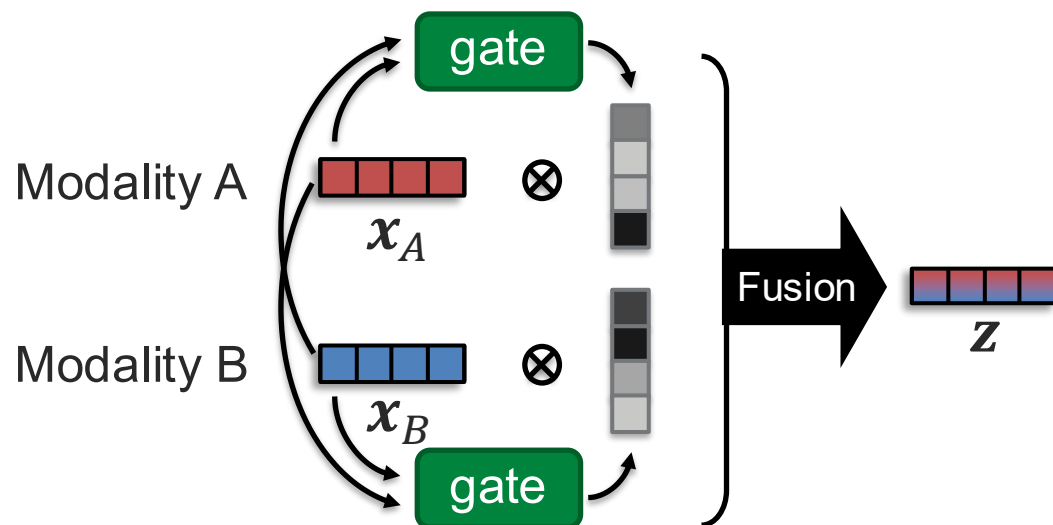
# Low-rank Fusion



# Low-rank Fusion



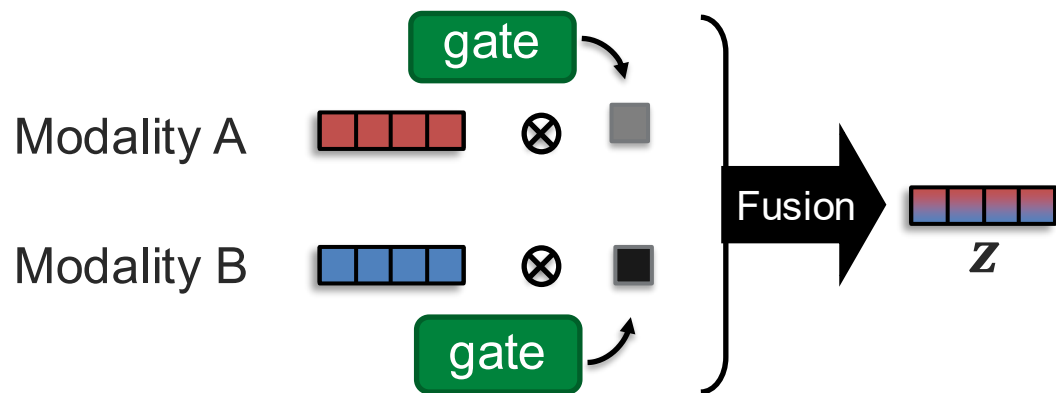
# Gated Fusion



Example with additive fusion:

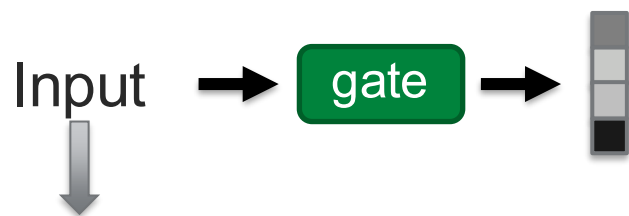
$$z = g_A(x_A, x_B) \cdot x_A + g_B(x_A, x_B) \cdot x_B$$

→  $g_A$  and  $g_B$  can be seen as attention functions



→ Gating output can be one weight for the whole modality

# Gated Fusion



What should it be?

Target modality [ ]

Other modality [ ]

All modality [ ]

[ ]

*“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)

Soft attention



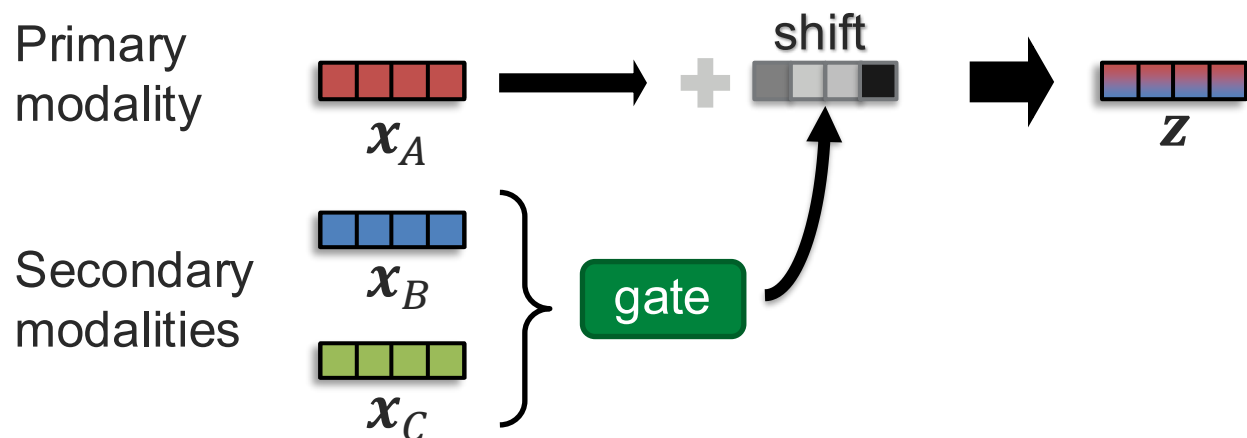
Easier to compute derivative (gradient)

Hard attention



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

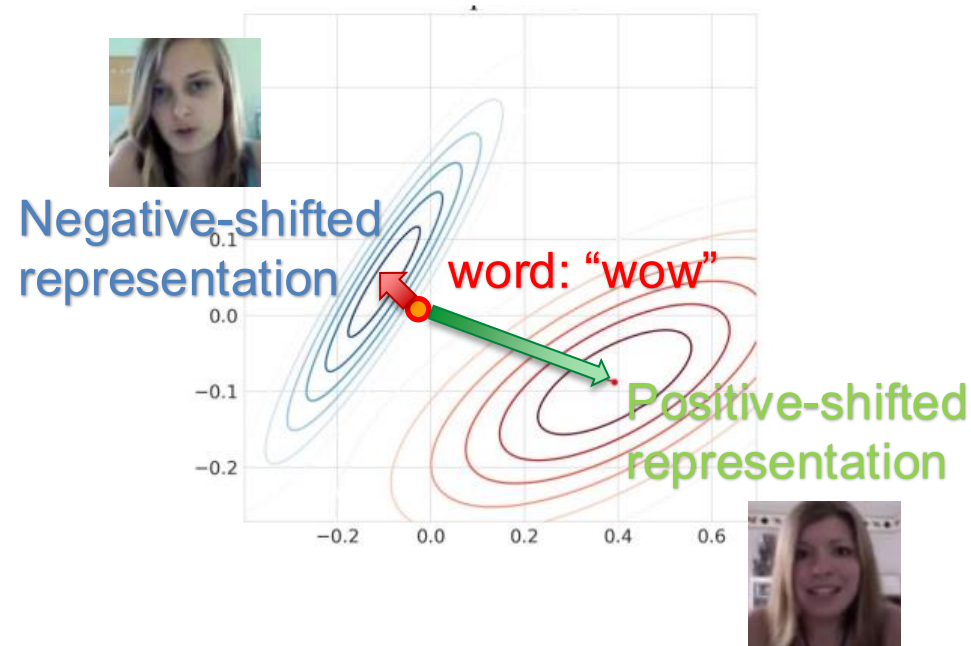
# Modality-Shifting Fusion



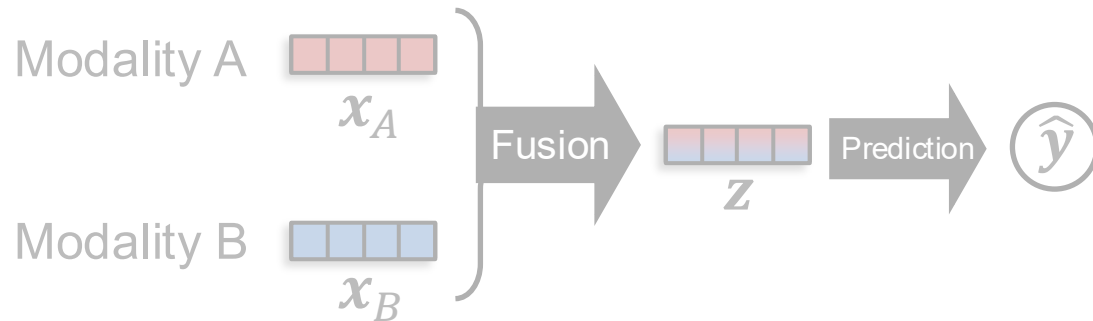
## Example with language modality:

Primary modality: language

Secondary modalities: acoustic and visual



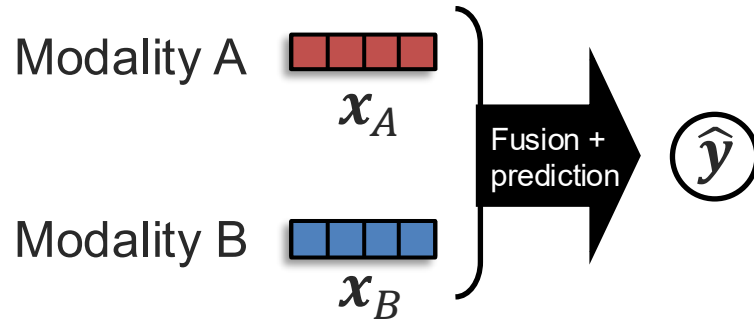
# Nonlinear Fusion



Nonlinear fusion:

$$\hat{\mathbf{y}} = f(\mathbf{x}_A, \mathbf{x}_B) \in \mathbb{R}^d$$

where  $f$  could be a neural network or any nonlinear model

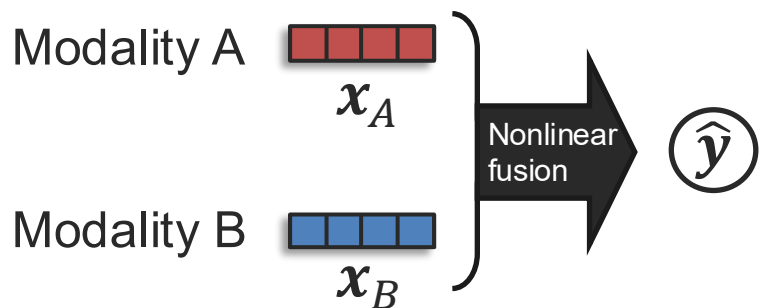


→ This could be seen as *early fusion*:

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

... 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?

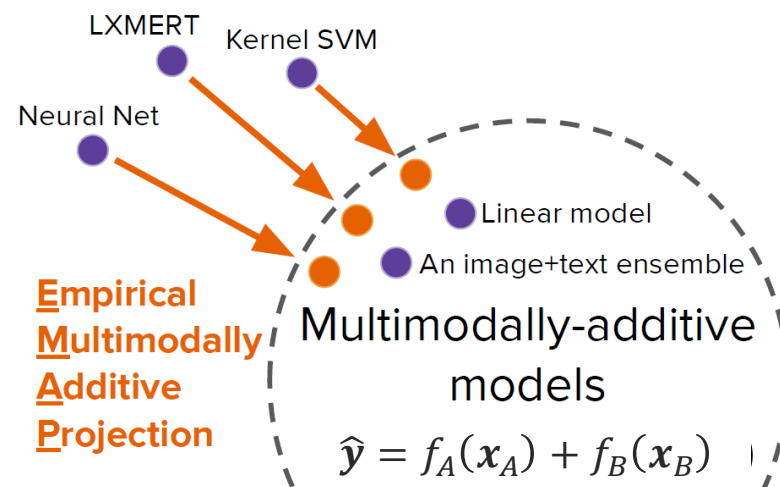
Additive fusion:

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

Projection from nonlinear to additive (using EMAP):

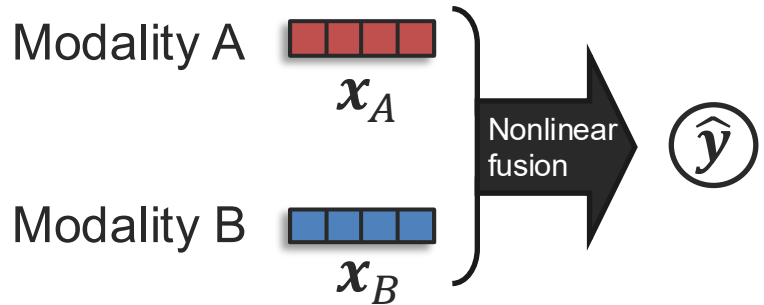
$$\tilde{f}(\mathbf{x}_A, \mathbf{x}_B) = \underbrace{\mathbb{E}_{\mathbf{x}_B} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_A(\mathbf{x}_A)} + \underbrace{\mathbb{E}_{\mathbf{x}_A} [f(\mathbf{x}_A, \mathbf{x}_B)]}_{f_B(\mathbf{x}_B)}$$

Modality A + Modality B



Additive fusion  
(approximation)

# Measuring Non-Additive Interactions



Nonlinear fusion:

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

EMAP projection

Additive fusion:

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

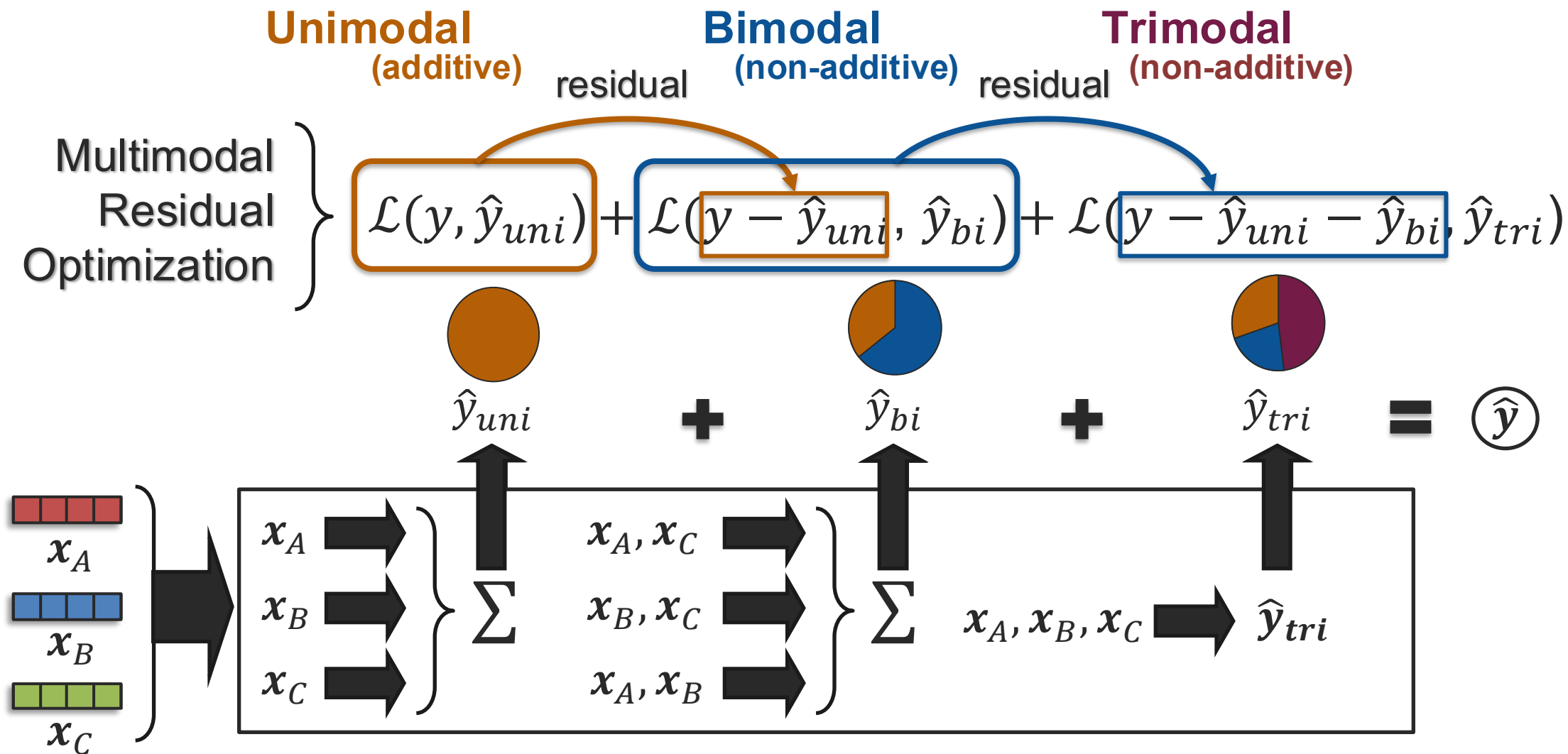
	I-INT	I-SEM	I-CTX	T-VIS	R-POP	T-ST1	T-ST2
<b>Nonlinear</b> ← Neural Network	90.4	69.2	78.5	51.1	63.5	71.1	79.9
<b>Polynomial</b> ← Polykernel SVM	<b>91.3</b>	<b>74.4</b>	<b>81.5</b>	50.8	–	72.1	<b>80.9</b>
<b>Nonlinear</b> ← FT LXMERT	83.0	68.5	76.3	<b>53.0</b>	63.0	66.4	78.6
<b>Nonlinear</b> ← $\hookrightarrow$ + Linear Logits	89.9	73.0	80.7	<b>53.4</b>	<b>64.1</b>	<b>75.5</b>	80.3
<b>Additive</b> ← Linear Model	90.4	72.8	80.9	51.3	63.7	<b>75.6</b>	76.1
<b>Best Model</b>	<b>91.3</b>	<b>74.4</b>	<b>81.5</b>	<b>53.4</b>	<b>64.2</b>	<b>75.5</b>	<b>80.9</b>
<b>Additive</b> ← $\hookrightarrow$ + EMAP	<b>91.1</b>	<b>74.2</b>	<b>81.3</b>	51.0	<b>64.1</b>	<b>75.9</b>	<b>80.7</b>

Always a good baseline!

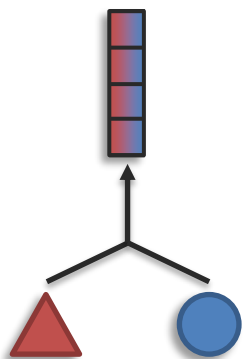
Differences are small!!!

# Non-Additive Interactions

Idea: prioritize simpler interactions



# Summary: How To Multimodal Fusion



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



# Assignments for This Coming Week

Project proposal due today. Meet with me at 4-5pm if you want feedback.

I want to meet every group at least once regarding their project ideas either today or this Thursday.

HW2 released last week, due next Wednesday 3/4.