

# Attending to Emotional Narratives

Desmond Ong

National University of Singapore  
& A\*STAR Singapore



Together with:

Zhengxuan Wu<sup>1</sup>, Xiyu Zhang<sup>1</sup>, Zhi-Xuan Tan<sup>2</sup>, Jamil Zaki<sup>1</sup>

<sup>1</sup>Stanford University & <sup>2</sup>A\*STAR



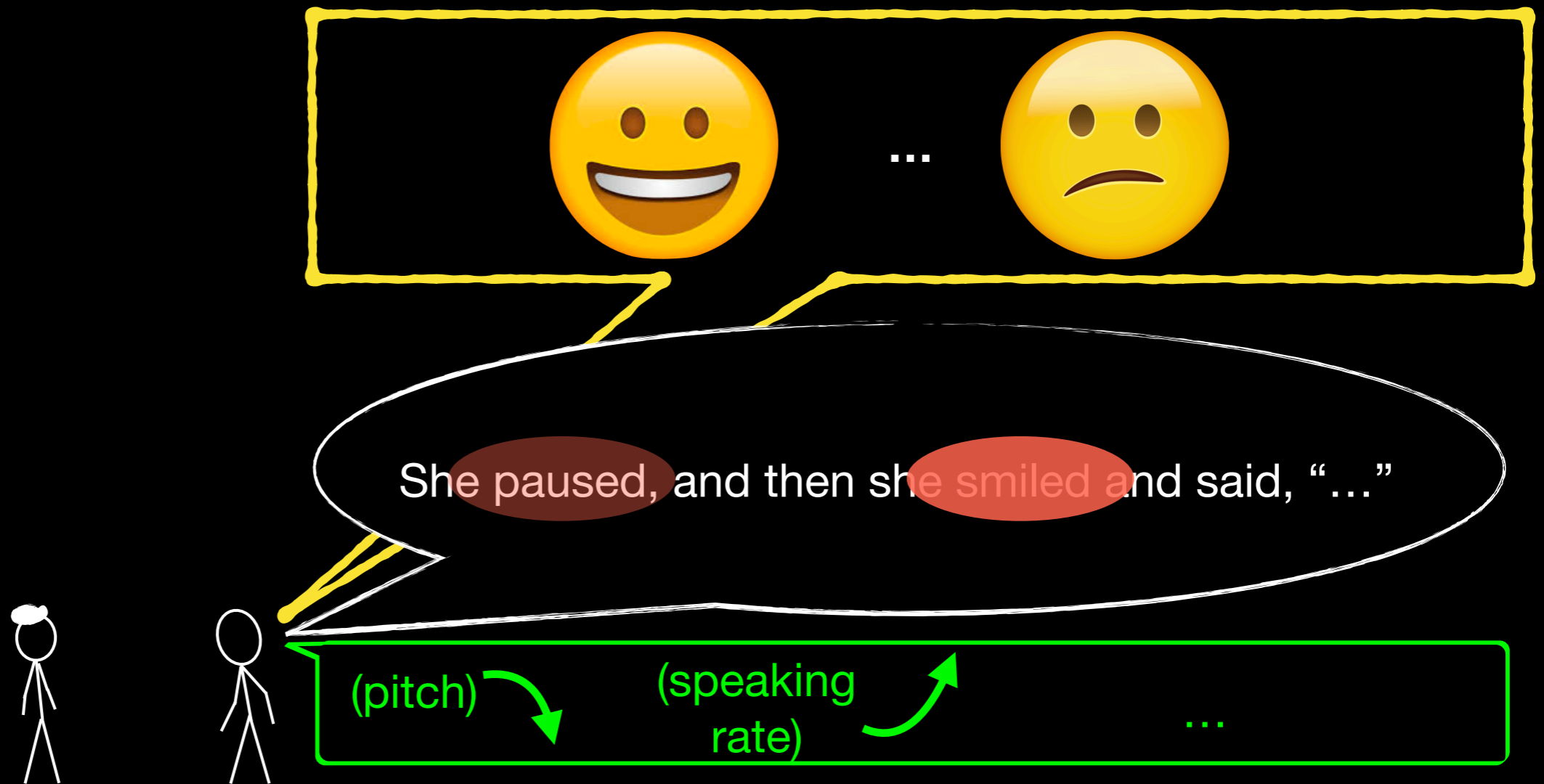
Agency for  
Science, Technology  
and Research



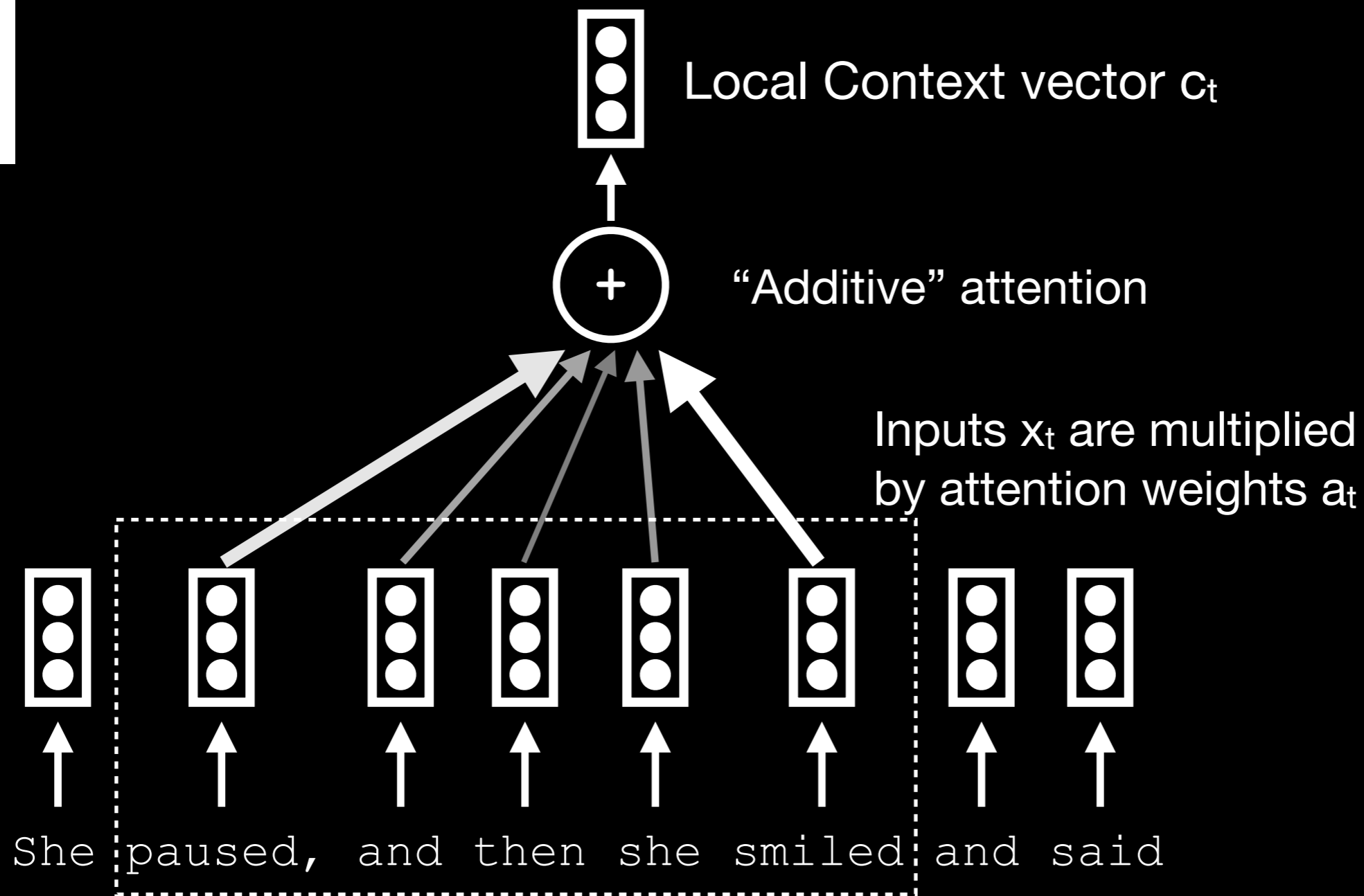
**NUS**  
National University  
of Singapore

School of  
Computing

# Human emotion reasoning



# Neural Network “Attention”



## Research Question:

Can these attention mechanisms improve multimodal emotion recognition?

# The Stanford Emotional Narratives Dataset (SEND)

Volunteers describing emotional life events.

This **first release (SENDv1)** contains:

- 49 unique “targets”
- N=193 video clips,
- ~2 mins each, total 7 hrs 15 mins.
- 60:20:20 Train/Valid/Test split

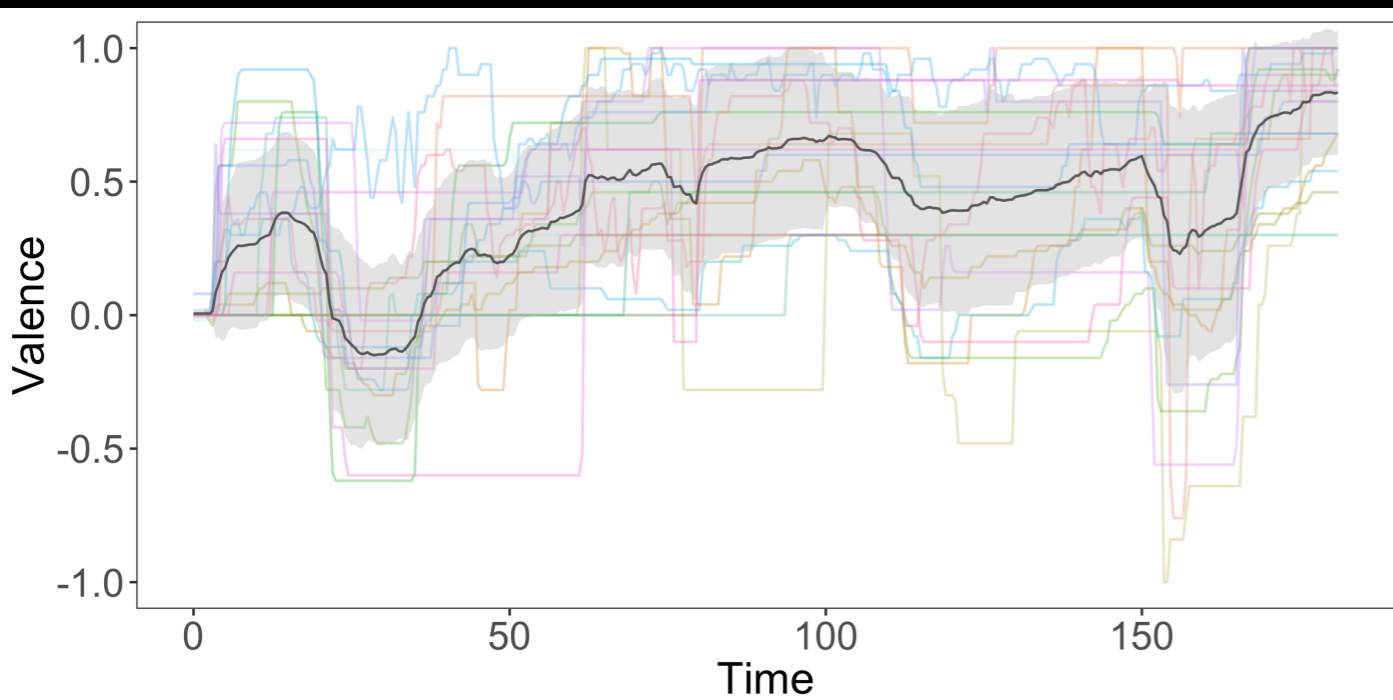


Each clip annotated by ~20 observers (on Amazon Mechanical Turk)  
Annotated for emotional valence, sampled every 0.5s.

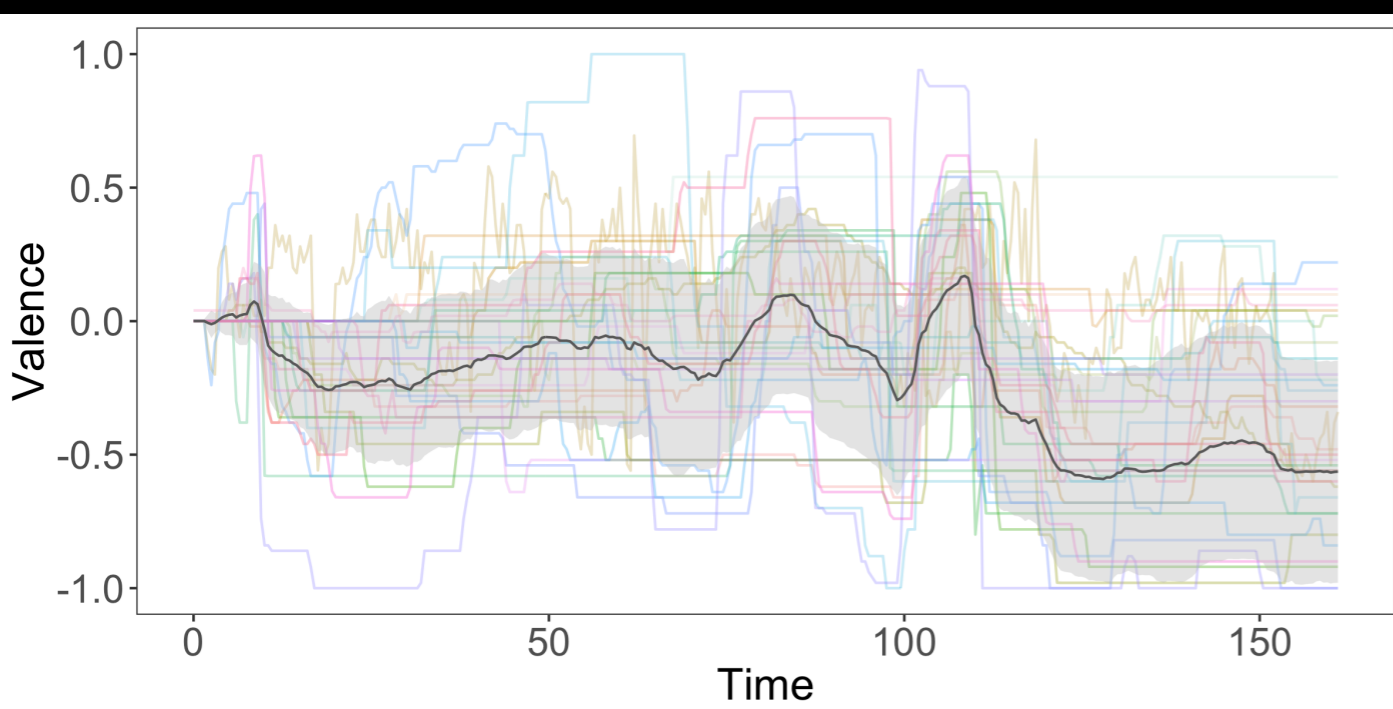
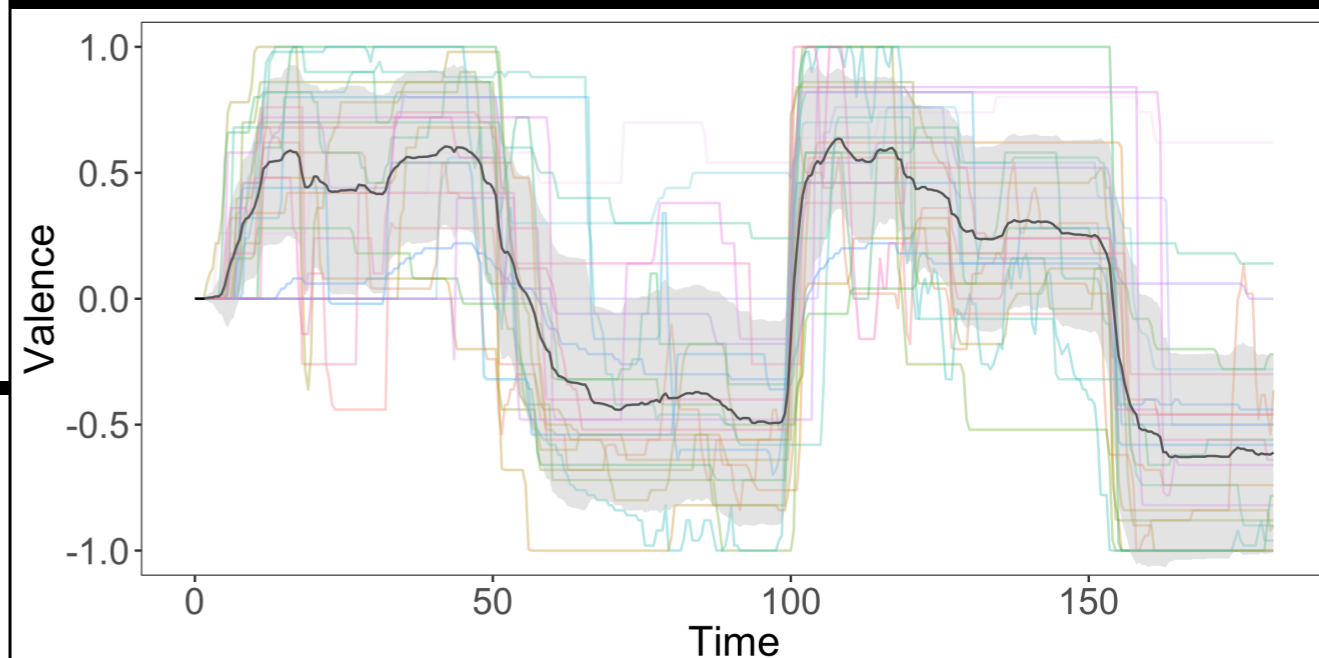
Gold-Standard Labels: **Evaluator-Weighted Estimator** (Grimm et al, 2007)

Evaluation Metric: Concordance Correlation Coefficient with the EWE

# The Stanford Emotional Narratives Dataset (SEND)



- heterogeneous
- complex emotional trajectories



# The Stanford Emotional Narratives Dataset (SEND)

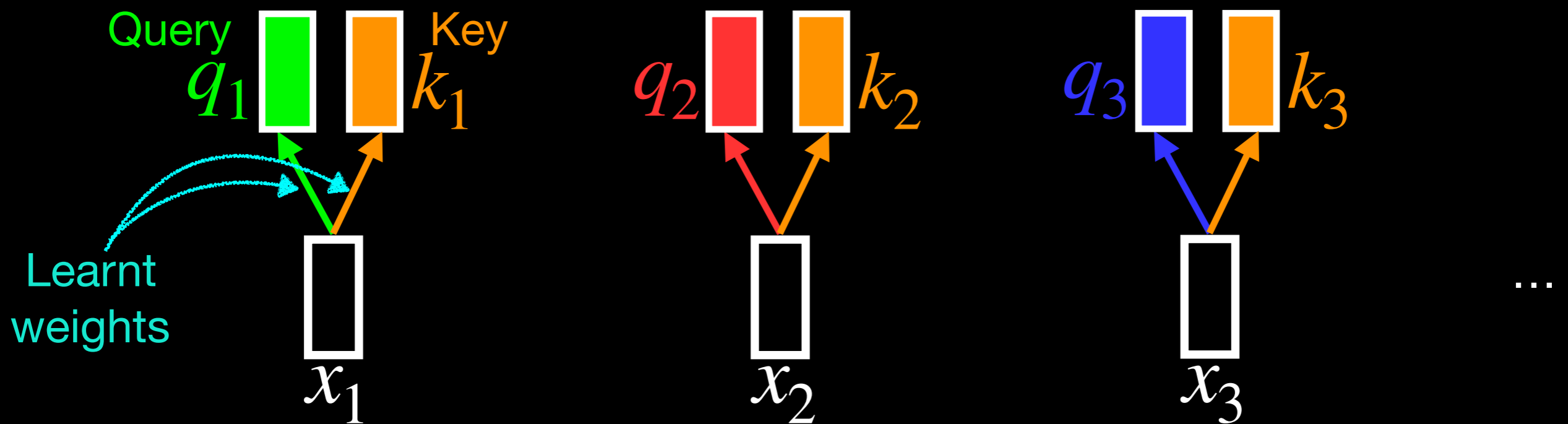
## Summary

- ★ Multimodal (Video, Audio, Text)
- ★ Unscripted, naturalistic expressions
- ★ Large diversity of stories
- ★ Continuous over time (time-series)
- ★ Dimensional labels
- ★ Multiple annotations —> calculate reliable estimate

# Transformers (1)

State-of-the-art in NLP

Introduces the concept of “self-attention”



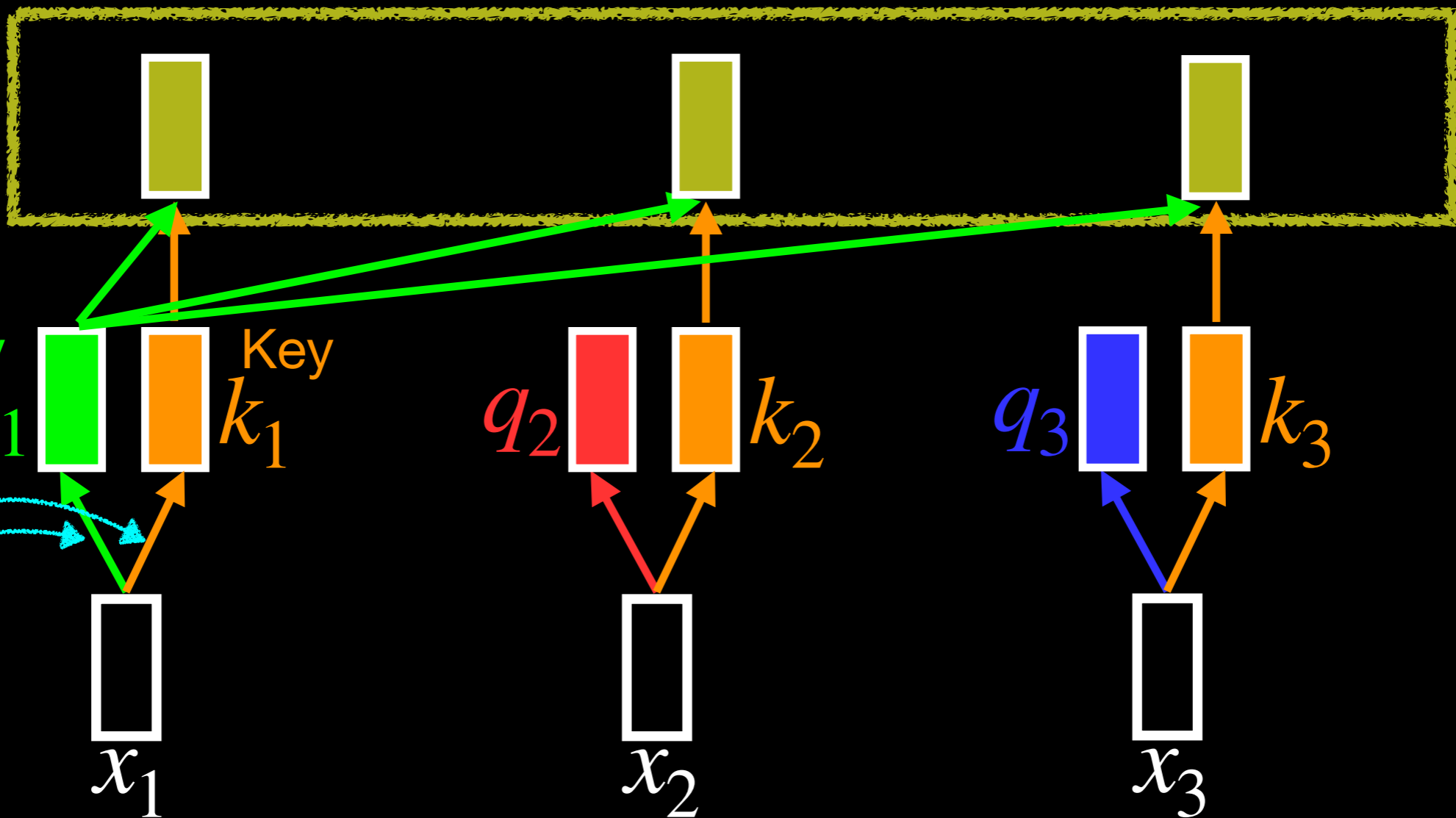
Input token (Feature vector at time  $t_1$ )

# Transformers (1)

State-of-the-art in NLP

Introduces the concept of “self-attention”

“Attention weights of token 1 on others”

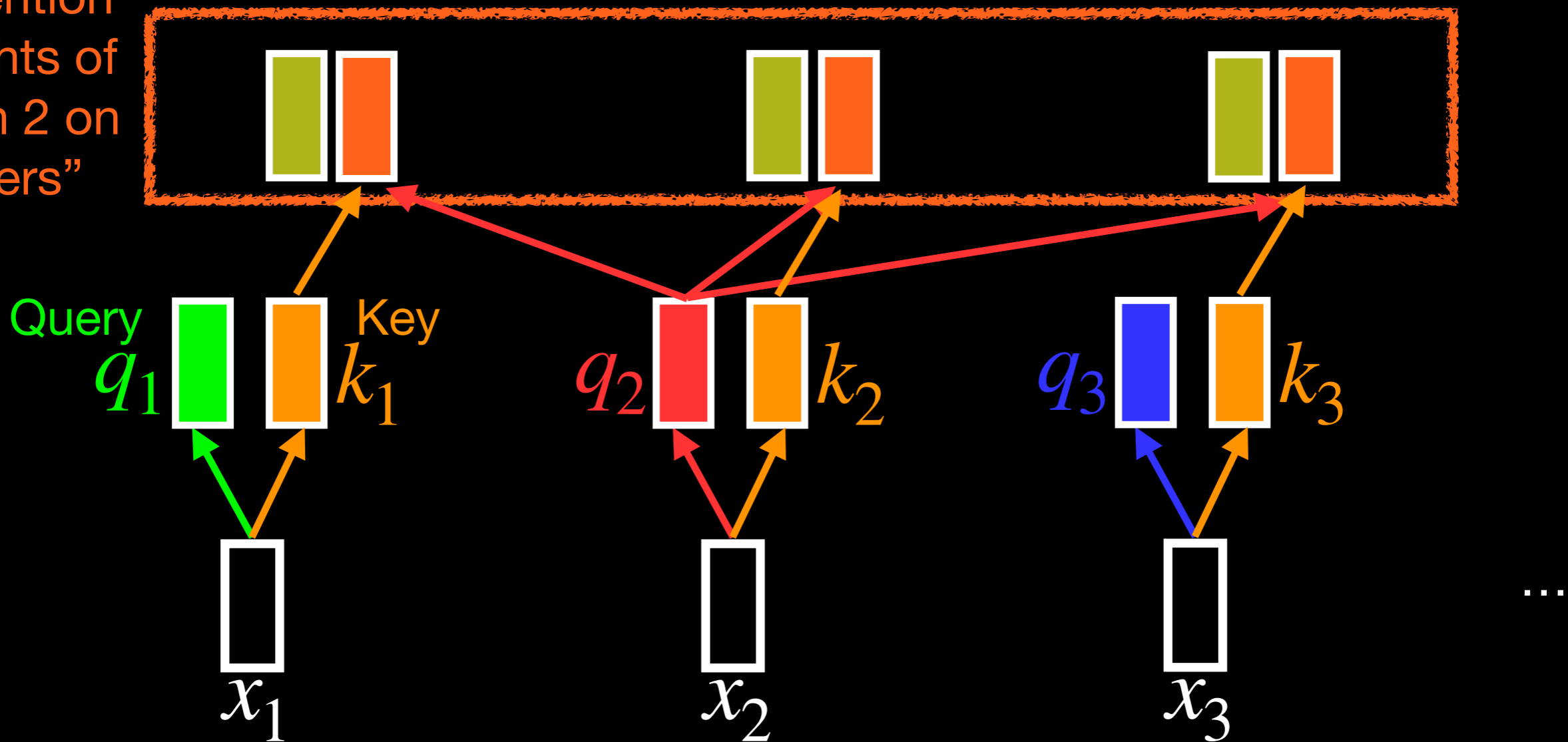


Input token (Feature vector at time  $t_1$ )



# Transformers (2)

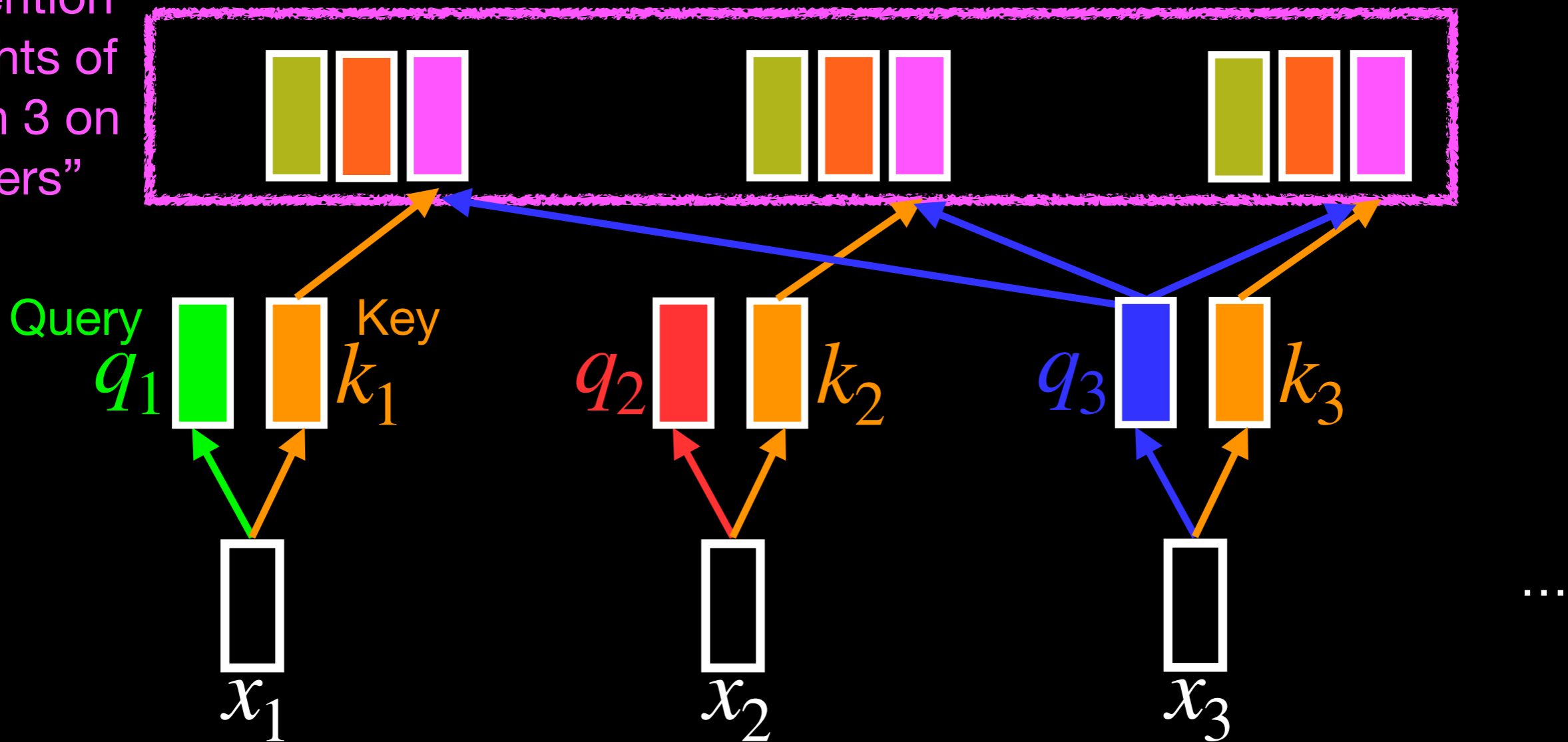
“Attention weights of token 2 on others”



Input token (Feature vector at time  $t_1$ )

# Transformers (3)

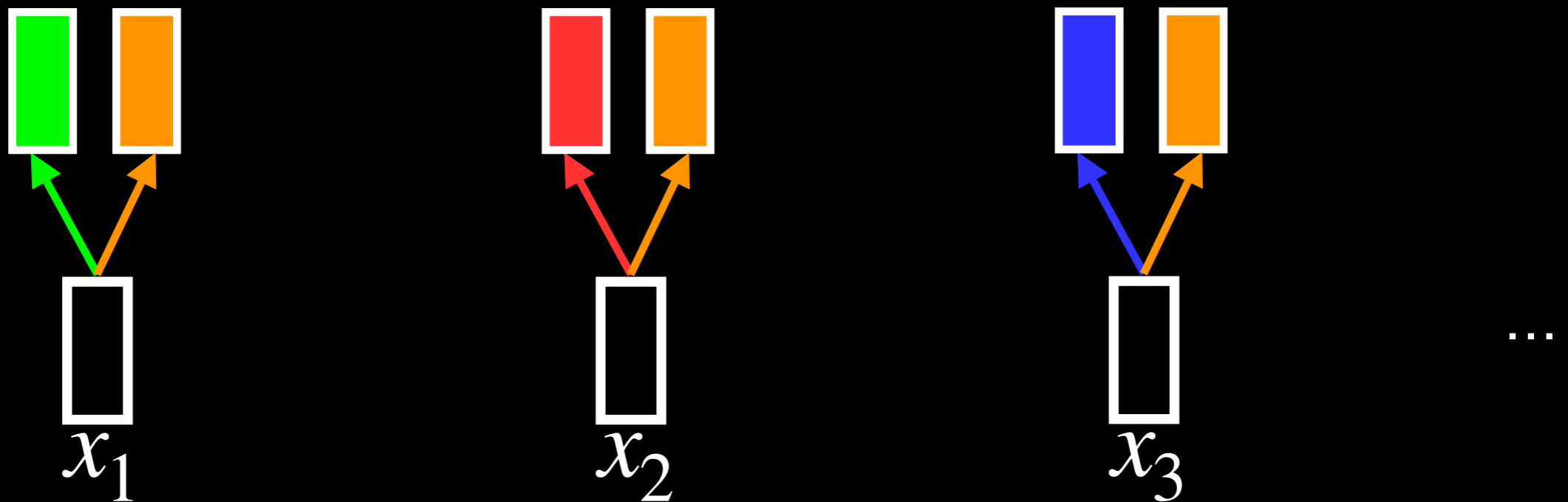
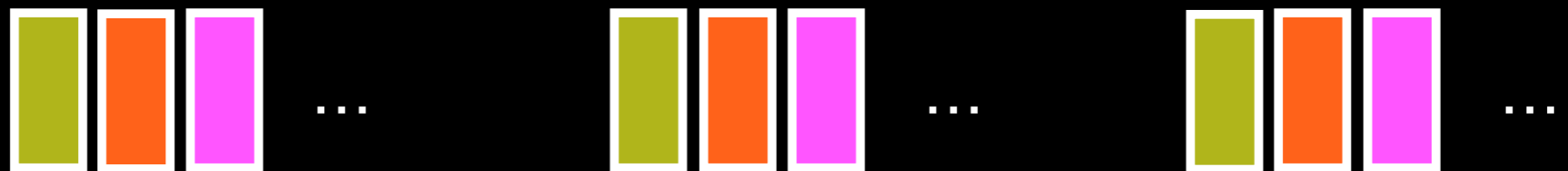
“Attention weights of token 3 on others”



Input token (Feature vector at time  $t_1$ )

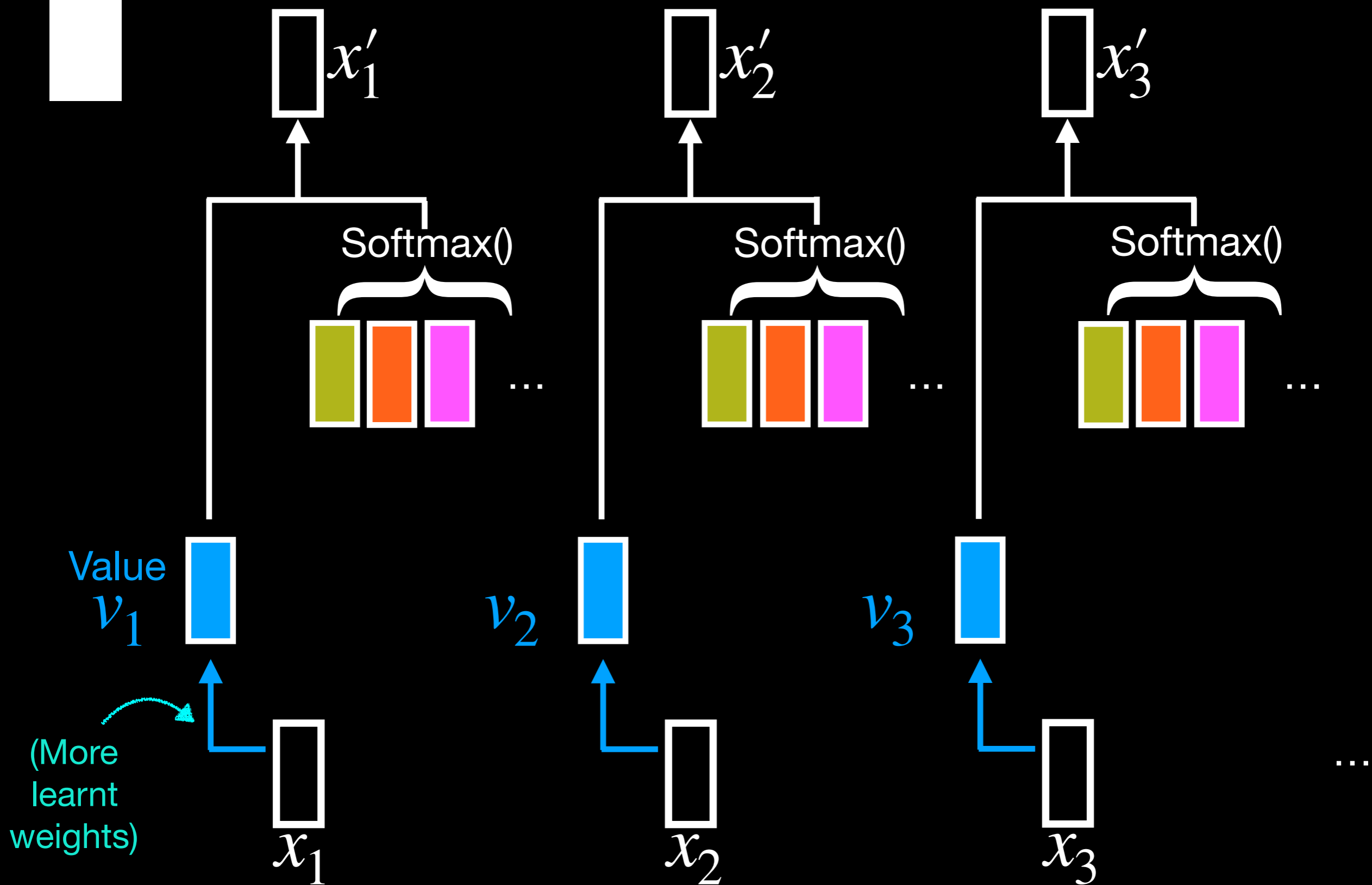
# Transformers (4)

“Attention weights of all tokens on all others”



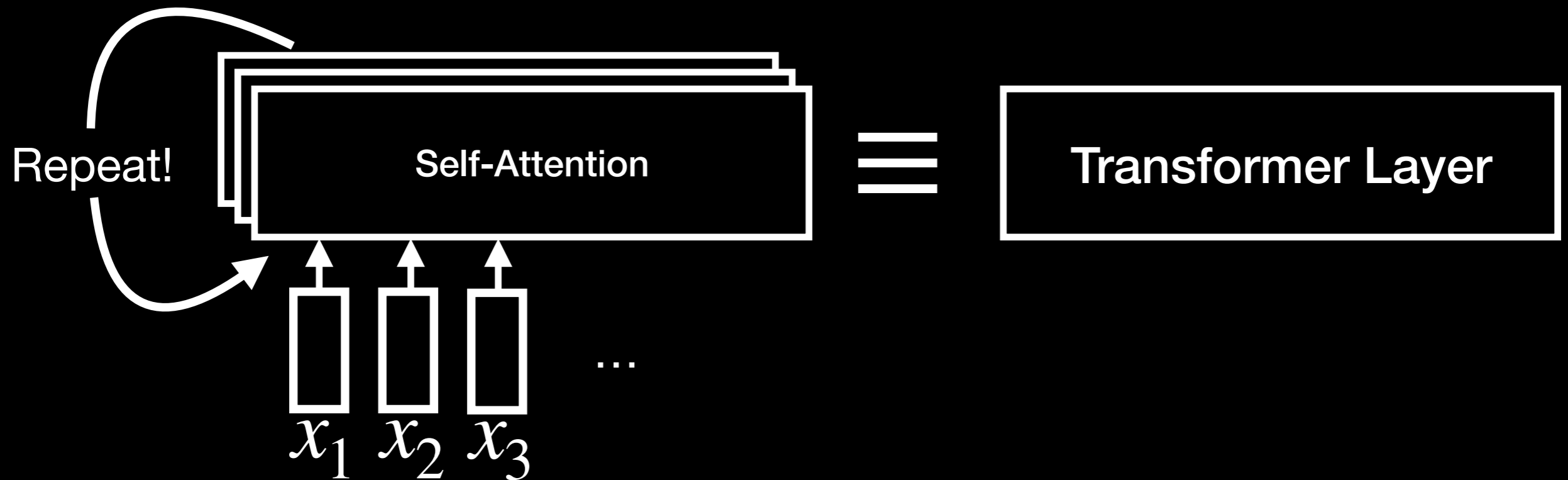
Input token (Feature vector at time  $t_1$ )

# Transformers (5)

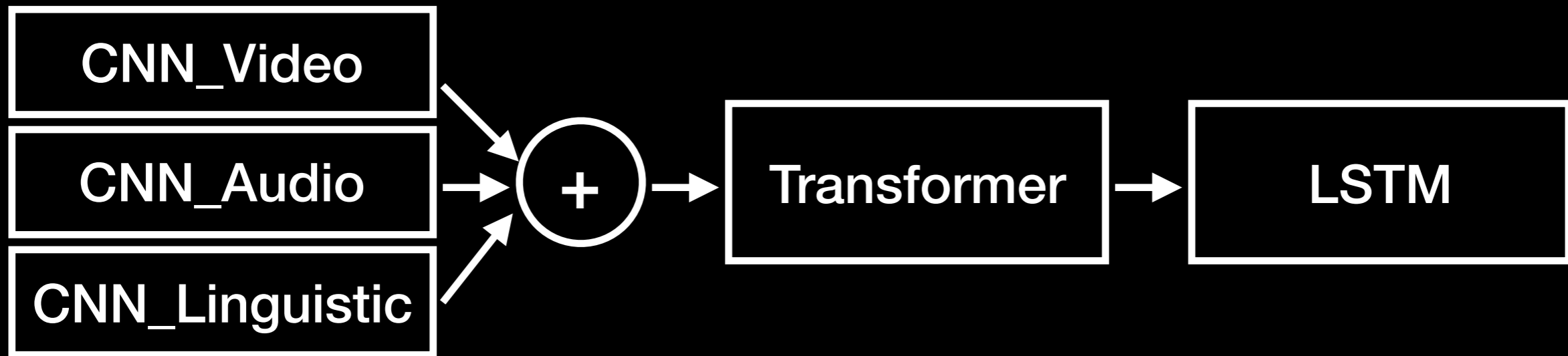


Input token (Feature vector at time  $t_1$ )

# Transformers (6)



# Simple Fusion Transformer + Results



## Concordance Correlation on Test Set

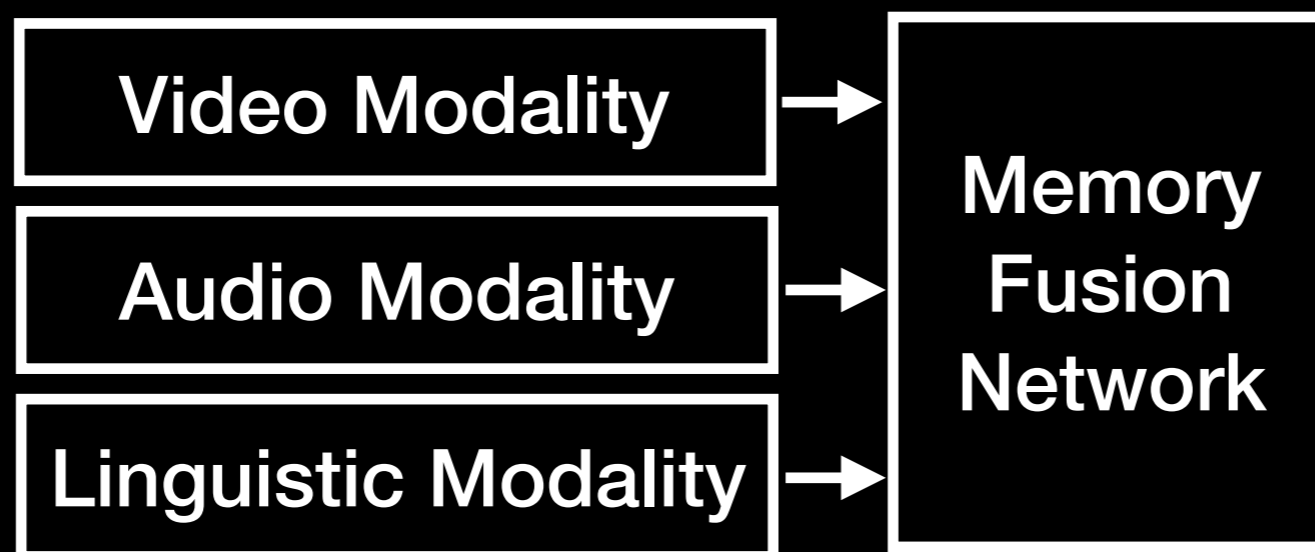
	Best Unimodal	Best Bimodal	Trimodal
SFT	.34	.35	.14
Human	-	-	.50

# Memory Fusion Transformer

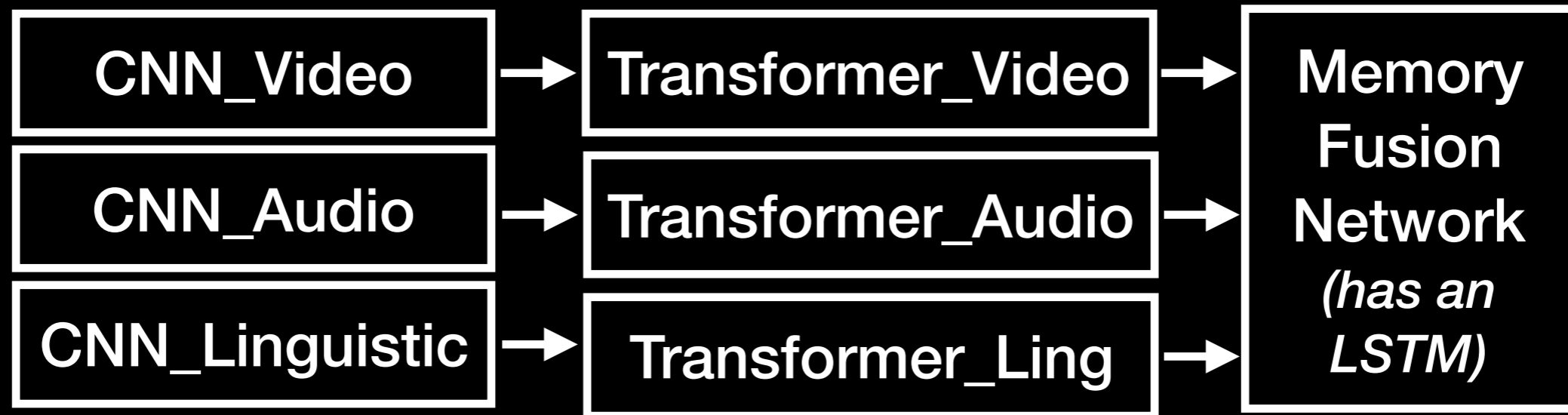
Our Simple Fusion Transformer [**“Self-Attention”**]

- Does well on Linguistic input
- And on Linguistic + Visual
- But **performs poorly on trimodal input**

Decided to also implement **Memory Fusion Network** (Zadeh et al, 2018), which learns **“cross-modality attention”**



# Memory Fusion Transformer



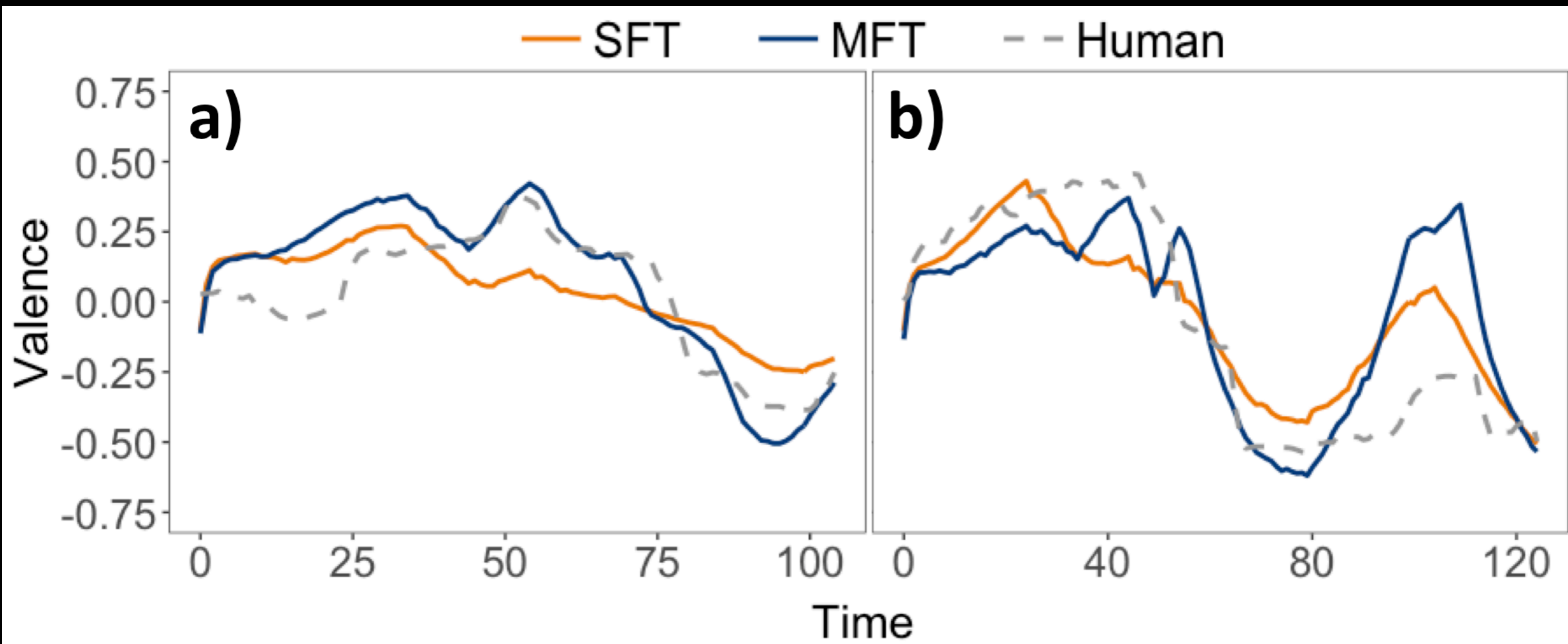
## Concordance Correlation on Test Set

	Best Unimodal	Best Bimodal	Trimodal
SFT	.34	.35	.14
MFT	-	.36	.44
Human	-	-	.50

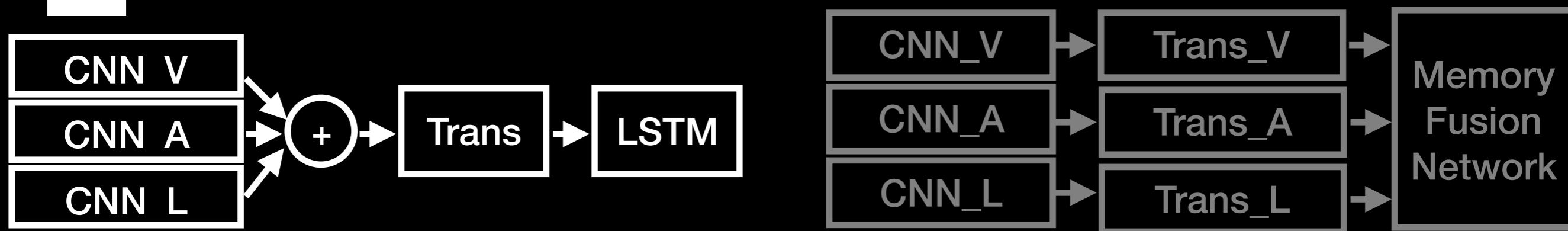
*(n.s)*



# Results



# Lesion Experiments (1)

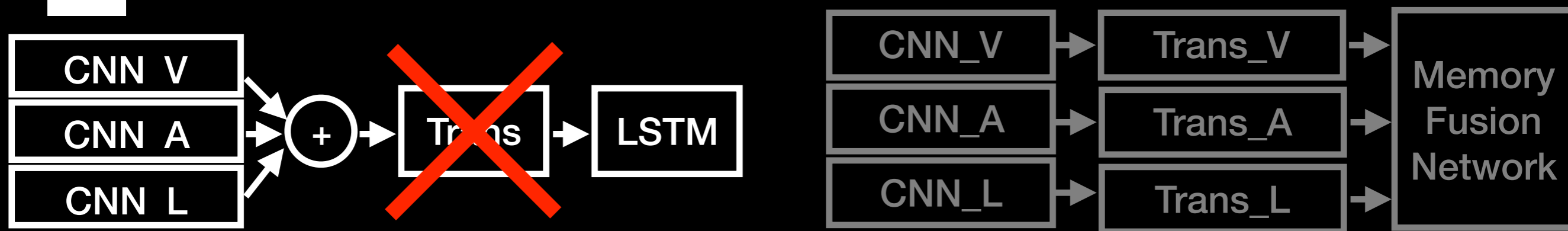


Simple Fusion Transformer

## Concordance Correlation on Test Set

	Best Unimodal	Best Bimodal	Trimodal
<b>SFT</b>	<b>.34</b>	<b>.35</b>	<b>.14</b>
LSTM-only			
Trans-only			
MFT	-	.36	.44
MFN-only			
<b>Human</b>	-	-	<b>.50</b>

# Lesion Experiments (1)

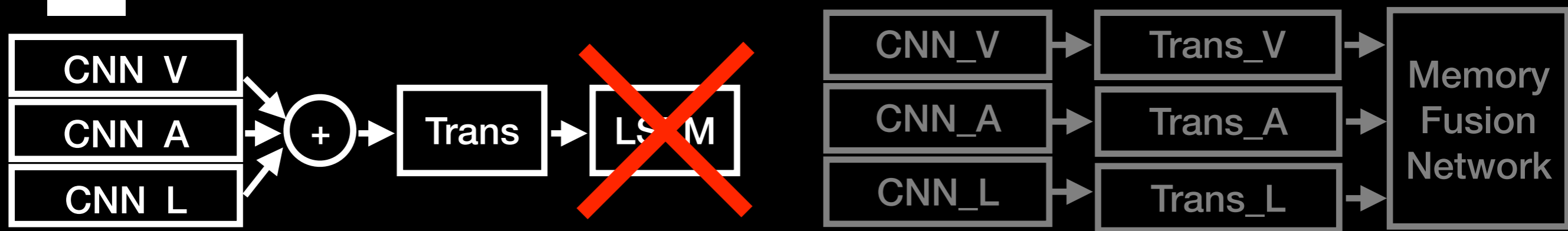


Simple Fusion Transformer

## Concordance Correlation on Test Set

	Best Unimodal	Best Bimodal	Trimodal
<b>SFT</b>	<b>.34</b>	<b>.35</b>	<b>.14</b>
<b>LSTM-only</b>	<b>.21</b>	<b>.17</b>	<b>-.02</b>
<b>Trans-only</b>			
<b>MFT</b>	-	.36	.44
<b>MFN-only</b>			
<b>Human</b>	-	-	<b>.50</b>

# Lesion Experiments (2)



Simple Fusion Transformer

## Concordance Correlation on Test Set

	Best Unimodal	Best Bimodal	Trimodal
<b>SFT</b>	<b>.34</b>	<b>.35</b>	<b>.14</b>
LSTM-only	.21	.17	-.02
<b>Trans-only</b>	<b>.05</b>	<b>.05</b>	<b>.00</b>
MFT	-	.36	.44
MFN-only	-	-	-
<b>Human</b>	-	-	<b>.50</b>

# Lesion Experiments (3)



Memory Fusion Transformer

## Concordance Correlation on Test Set

	Best Unimodal	Best Bimodal	Trimodal
SFT	.34	.35	.14
LSTM-only	.21	.17	-.02
Trans-only	.05	.05	.00
<b>MFT</b>	-	<b>.36</b>	<b>.44</b>
MFN-only			
Human	-	-	.50

# Lesion Experiments (3)



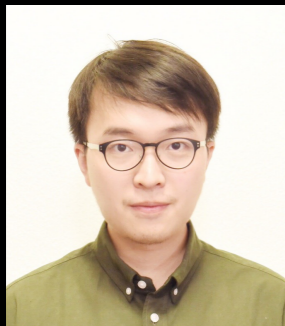
Memory Fusion Transformer

## Concordance Correlation on Test Set

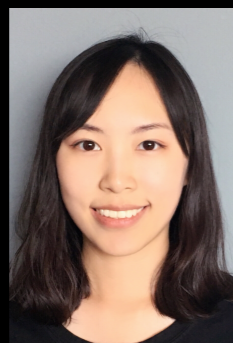
	Best Unimodal	Best Bimodal	Trimodal
SFT	.34	.35	.14
LSTM-only	.21	.17	-.02
Trans-only	.05	.05	.00
<b>MFT</b>	-	<b>.36</b>	<b>.44</b>
<b>MFN-only</b>	-	<b>.33</b>	<b>.28</b>
Human	-	-	.50

# Summary, Limitations and Future Directions

- Showed that neural network attention mechanisms (self-attention, cross-modality attention) can improve multimodal emotion recognition.
  - Lesioned experiments suggest that different types of attention contribute to better performance.
- Current work: probing attention weights
- Could serve as a way to build explainable affective computers
- More work on SEND: More diverse demographics, cross-cultural...



Zhengxuan Wu



Xiyu Zhang



Zhi-Xuan Tan



Jamil Zaki

## Other collaborators on the SEND

Marianne Reddan

Xi Jia Zhou

Isabella Kahhale

Alison Mattek

Anat Perry (@ HUJI)

# Thanks!

[dco@comp.nus.edu.sg](mailto:dco@comp.nus.edu.sg)

[web.stanford.edu/~dco](http://web.stanford.edu/~dco)

Paper: <https://arxiv.org/abs/1907.04197>

Code: <https://github.com/frankaging/ACII2019-transformer>