



# Seq2Seq2Sentiment: Multimodal Sequence to Sequence Models for Sentiment Analysis

Hai Pham<sup>\*1</sup>, Thomas Manzini<sup>\*1</sup>, Paul Pu Liang<sup>2</sup>, Barnabás Póczos<sup>2</sup>

<sup>1</sup> *Language Technologies Institute*, <sup>2</sup> *Machine Learning Department*  
*School of Computer Science, Carnegie Mellon University*

# Introduction

- Research Problem
  - Does representing multiple modalities jointly improve sentiment prediction for the CMU-MOSI dataset?

# Introduction

- Research Problem
  - Does representing multiple modalities jointly improve sentiment prediction for the CMU-MOSI dataset?
- Dataset
  - Multimodal Corpus of Sentiment Intensity and Subjectivity Analysis in Online Opinion Videos (MOSI)
    - Modalities: Video, Audio, Text Transcripts
    - 89 speakers, 93 videos split into 2199 labeled opinion segments
    - Labels for Sentiment:  $\{-1, 1\}$  or  $\{-3, -2, -1, 0, 1, 2, 3\}$

# Baseline & Metrics

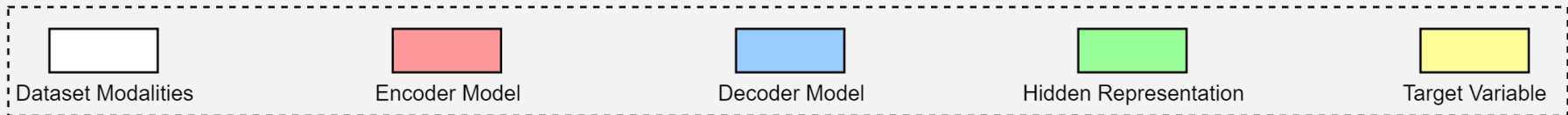
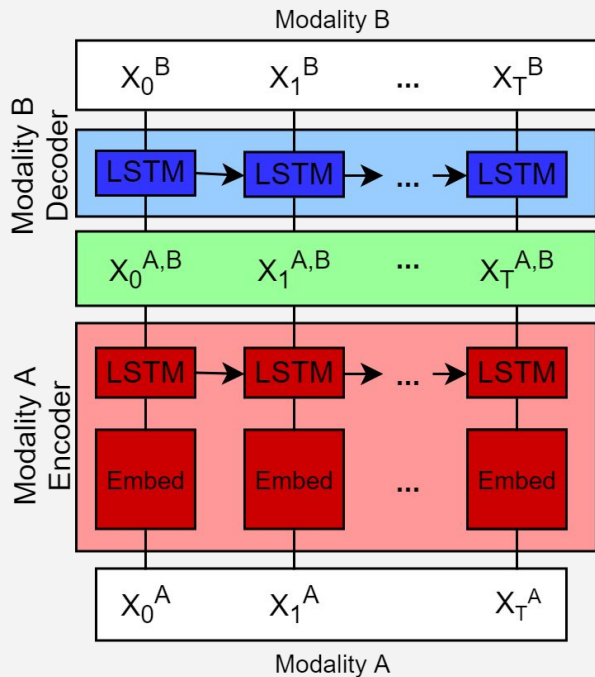
- LSTM Multi Modal Baseline
  - Concatenates all modalities together and predicts sentiment from these modalities.
  - **75% Accuracy** (Chen et. al.)
- Metrics
  - For both  $\{-1, 1\}$  or  $\{-3, -2, -1, 0, 1, 2, 3\}$  cases
    - Precision/Recall (Test)
    - F1 Score (Test)

# Related Work

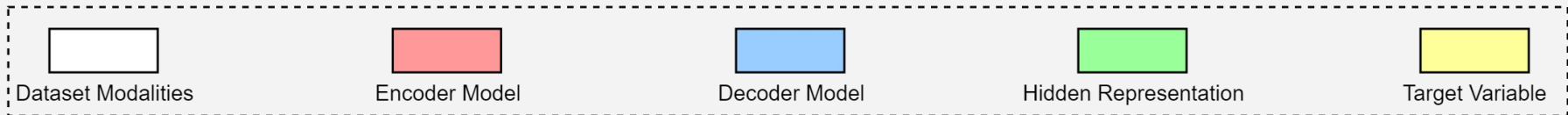
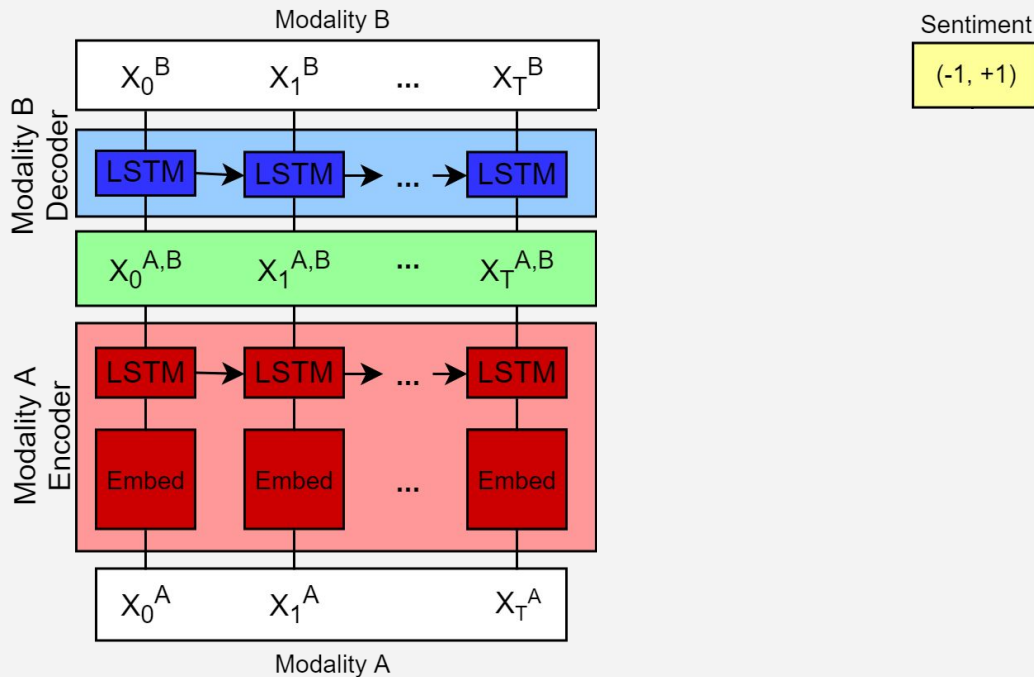
- Word-level Temporal Methods
  - Gu et al. ACL 2018, Chen et al. ICMI 2017
- Context-Dependent Methods
  - Poria et al. ACL 2017
- Memory-based Methods
  - Zadeh et al. AAI 2018
- Tensor-based Methods
  - Liu et al. ACL 2018, Zadeh et al. EMNLP 2017
- Conditional Approaches
  - Mirza et al. 2014, Kingma et al. 2014, & Pandey et al. 2017
- Attention-based Methods
  - Bahdanau et al. 2014, Luong et al. 2015

*See our paper for an exhaustive review of related work*

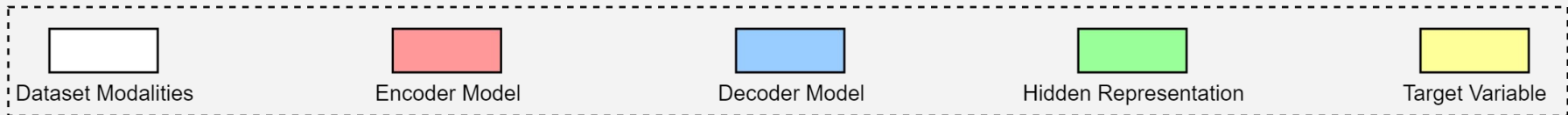
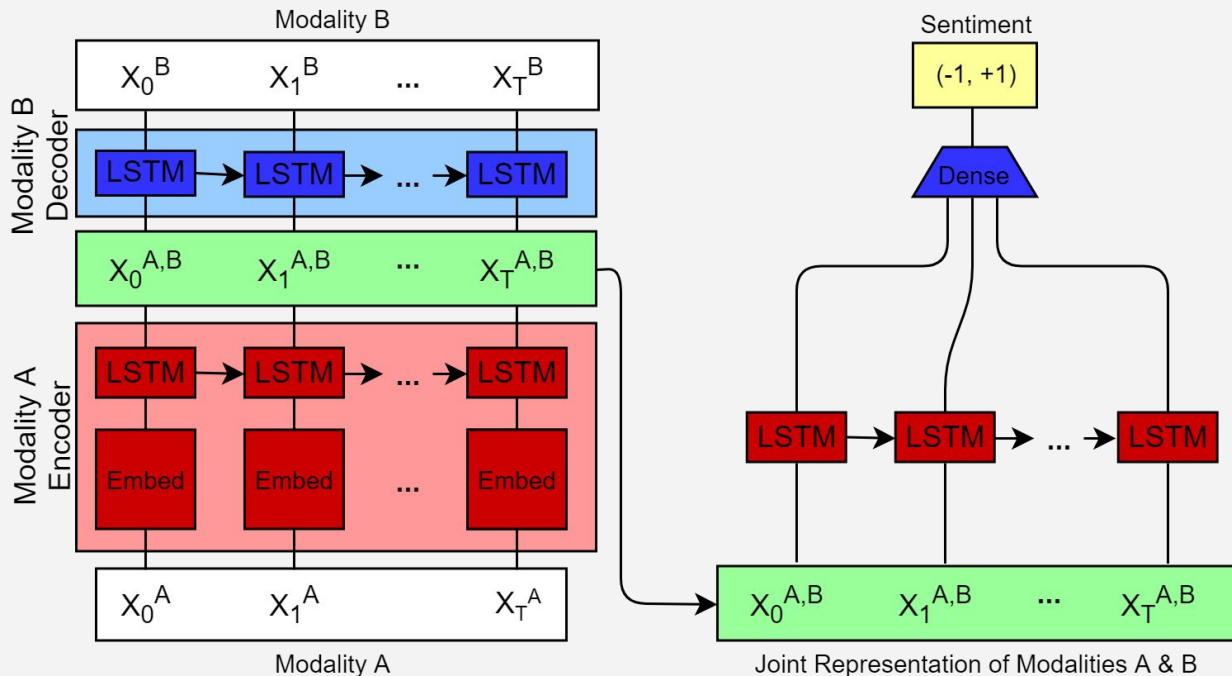
# Seq2Seq Modality Translation



# Seq2Seq Modality Translation

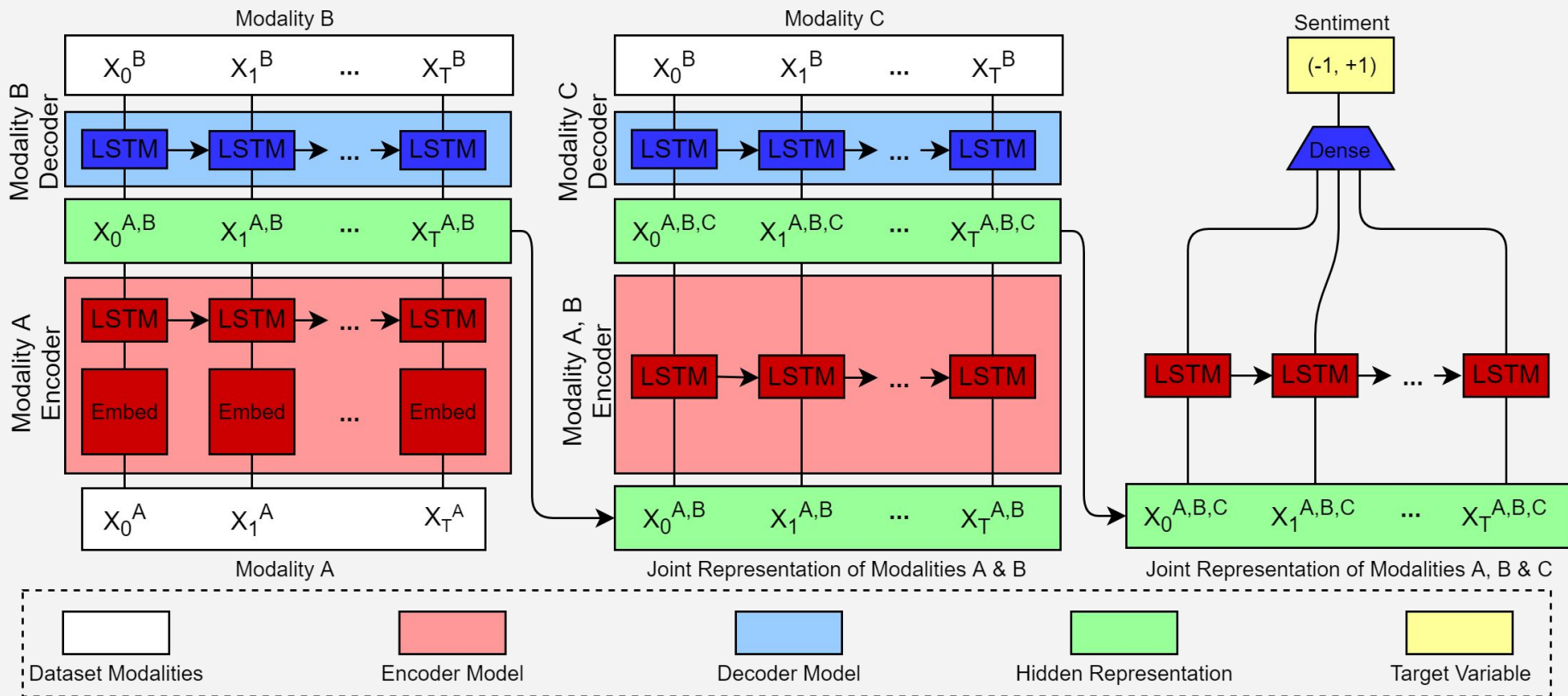


# Seq2Seq Modality Translation





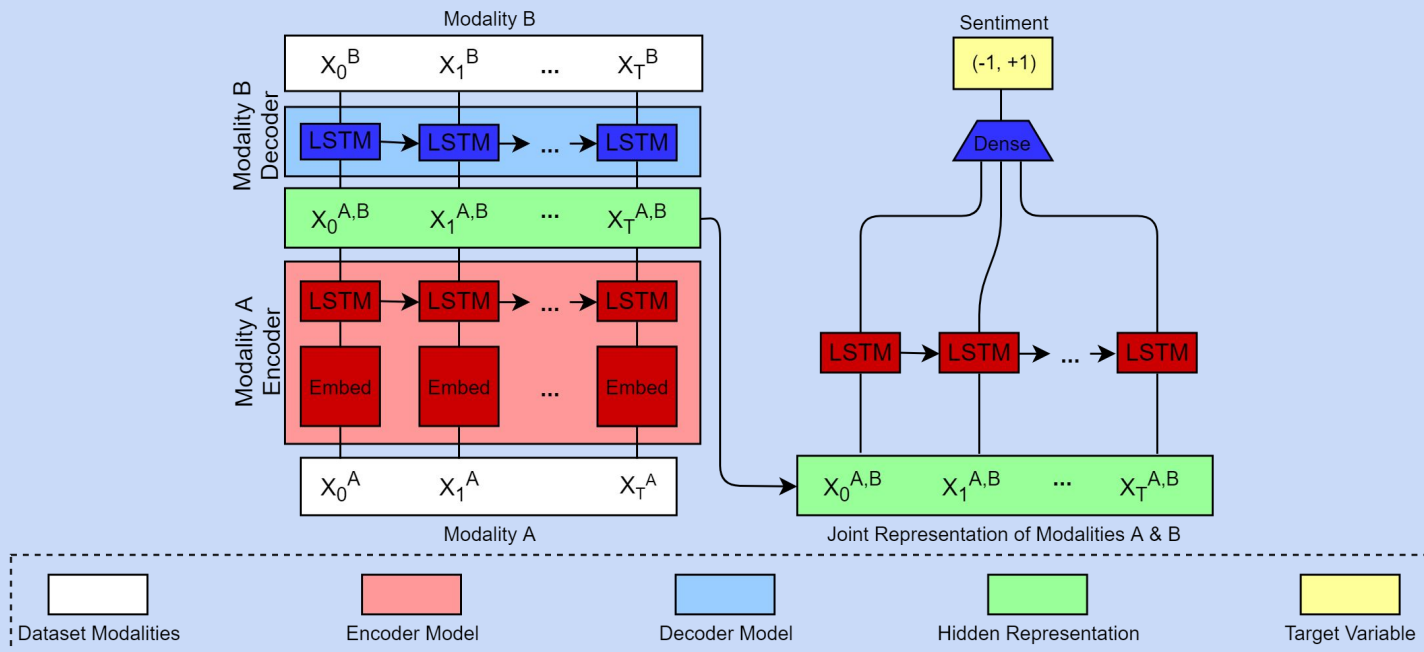
# Hierarchical Seq2Seq Modality Translation



# Experiments

Denoted as  
 $A \rightarrow B$

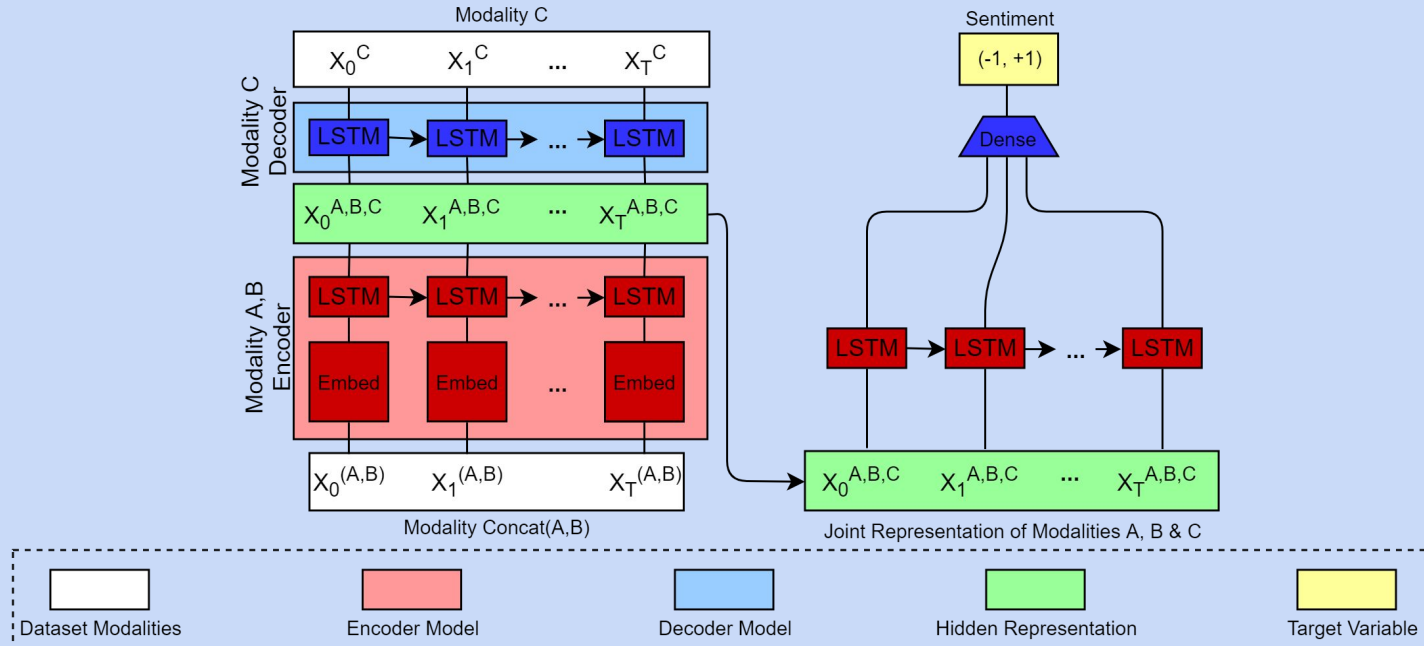
- We explored several different network architectures for generating representations



# Experiments

Denoted as  
 $\text{Concat}(A,B) \rightarrow C$

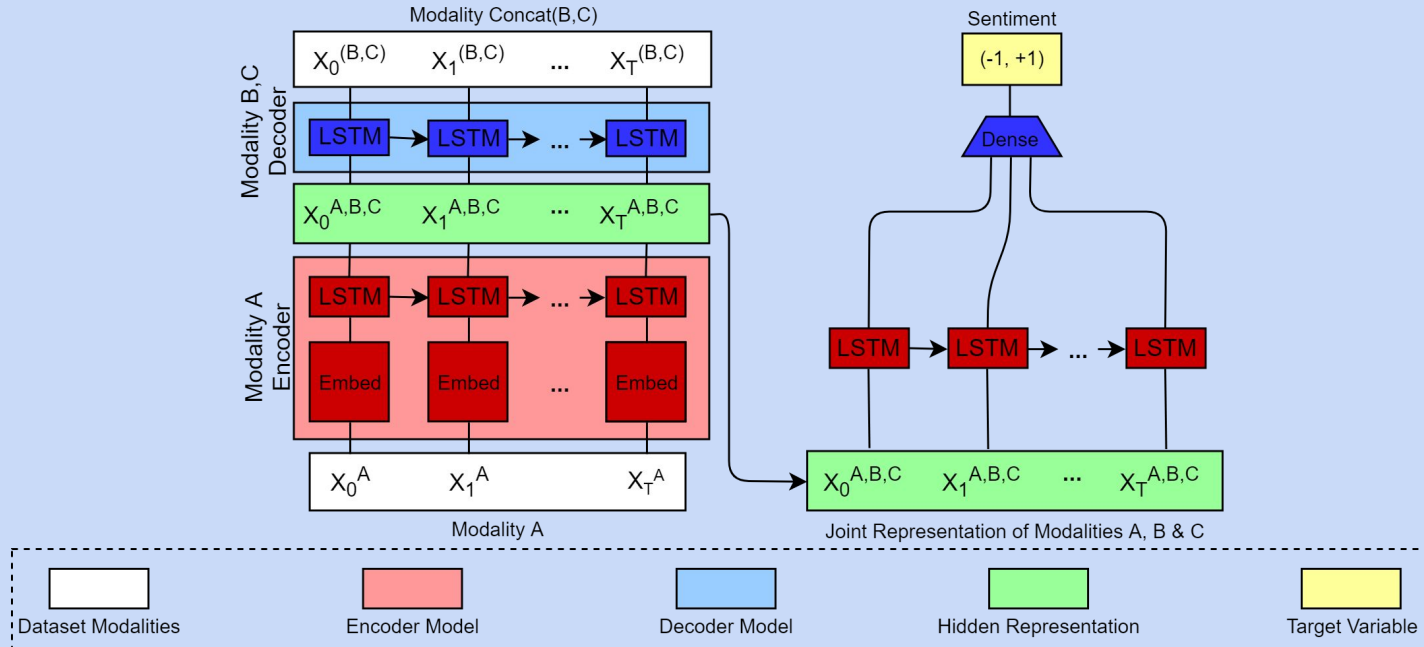
- We explored several different network architectures for generating representations



# Experiments

Denoted as  
 $A \rightarrow \text{Concat}(B,C)$

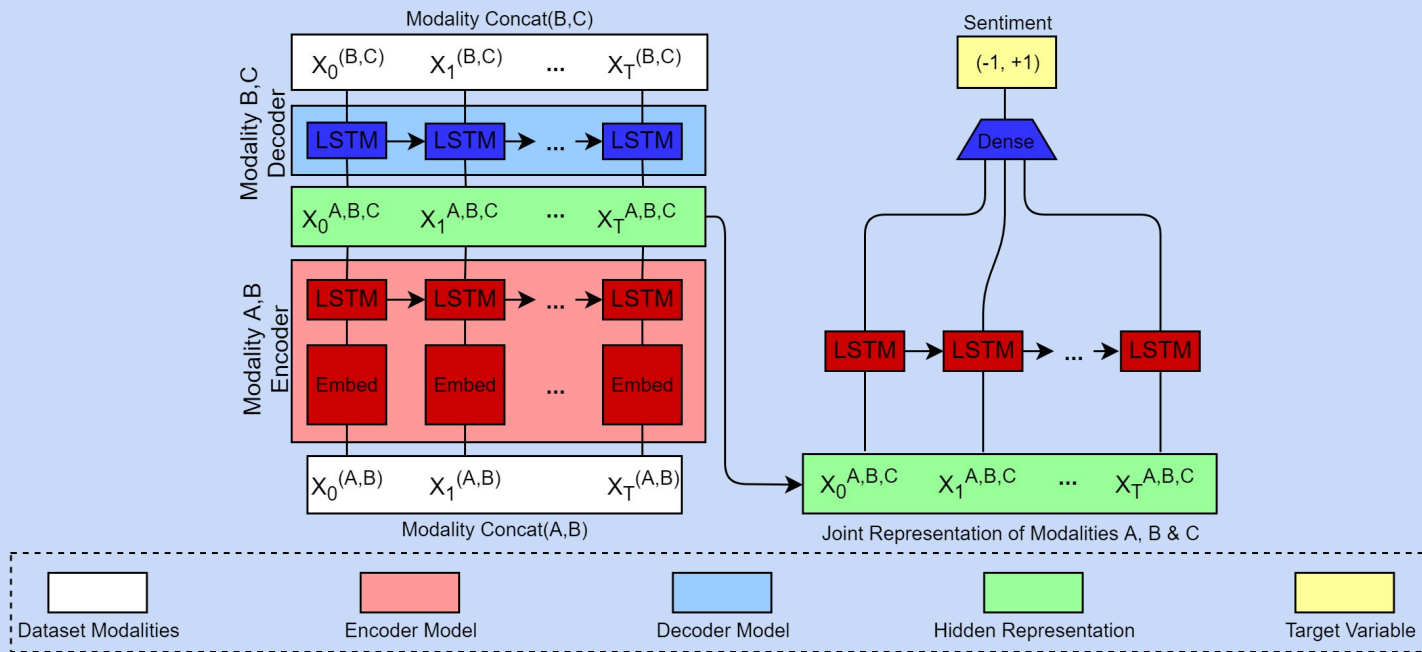
- We explored several different network architectures for generating representations



# Experiments

Denoted as  
 $\text{Concat}(A,B) \rightarrow \text{Concat}(B,C)$

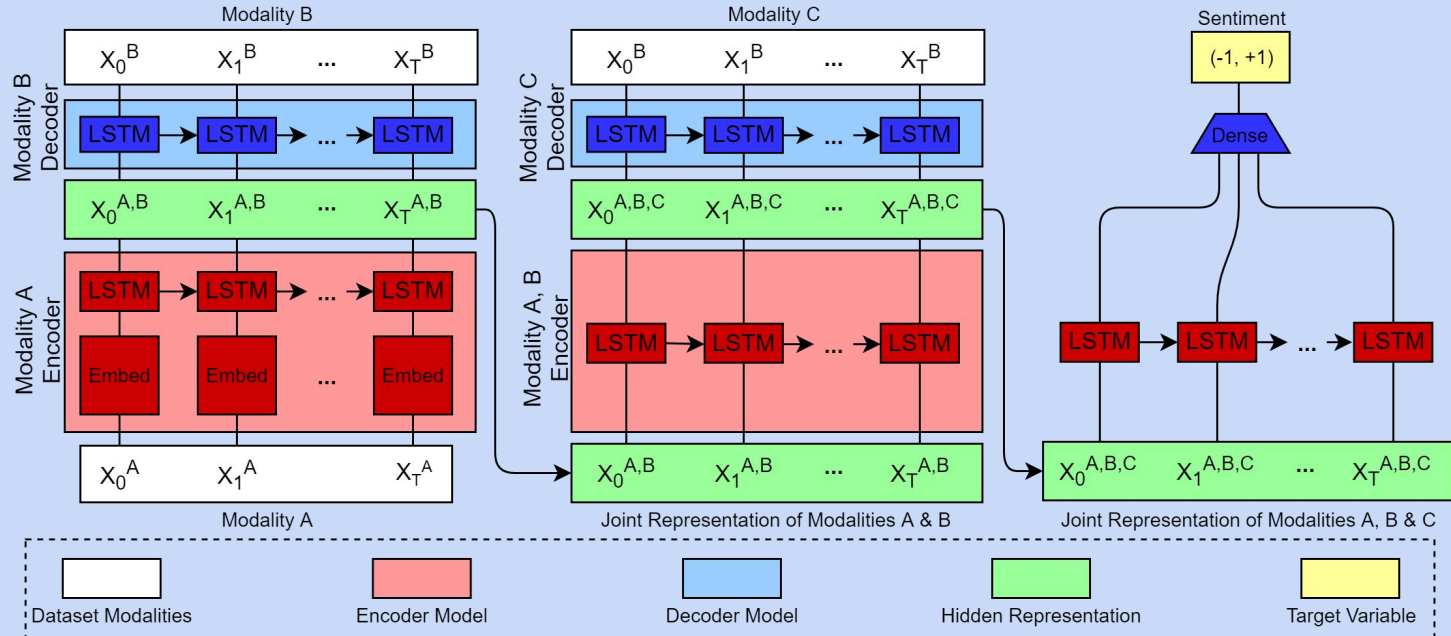
- We explored several different network architectures for generating representations



# Experiments

Denoted as  
 $\text{Embed}(A, B) \rightarrow C$

- We explored several different network architectures for generating representations



# Experiments

- Unimodal Baseline

Method	Feature	BINARY (-1, +1)			7-CLASS (-3, ..., +3)		
		Prec	Recall	F1	Prec	Recall	F1
UniModal-Baseline	Text (T)	<b>0.77</b>	<b>0.76</b>	<b>0.76</b>	<b>0.32</b>	<b>0.35</b>	<b>0.33</b>
	Audio (A)	0.56	0.56	0.56	0.12	0.19	0.14
	Video (V)	0.57	0.47	0.48	0.12	0.19	0.12

T = Text Modality, A = Audio Modality, V = Visual (facial) modality

# Results (Bi-Modal)

- Bimodal Baseline & Experimental Results

Method	Feature	BINARY (-1, +1)			7-CLASS (-3, ..., +3)		
		Prec	Recall	F1	Prec	Recall	F1
BiModal-Baseline	concat(T + V)	<b>0.78</b>	<b>0.67</b>	0.55	0.01	0.16	0.05
	concat(T + A)	0.44	0.66	0.53	0.02	0.15	0.04
	concat(A + V)	0.55	0.47	0.48	0.13	0.16	0.11
BiModal-Seq2Seq	T → V	0.67	<b>0.67</b>	<b>0.67</b>	0.26	0.22	<b>0.22</b>
	T → A	0.66	0.64	0.65	<b>0.28</b>	0.24	0.18
	A → T	0.55	0.60	0.56	0.17	<b>0.34</b>	0.11
	A → V	0.55	0.55	0.54	0.16	0.18	0.16
	V → T	0.58	0.58	0.58	0.05	0.16	0.08
	V → A	0.58	0.62	0.58	0.12	0.17	0.01

**T = Text Modality, A = Audio Modality, V = Visual (facial) modality**



# Results (Bi-Modal)

- Bimodal Baseline & Experimental Results

Method	Feature	BINARY (-1, +1)			7-CLASS (-3, ..., +3)		
		Prec	Recall	F1	Prec	Recall	F1
BiModal-Baseline	concat(T + V)	<b>0.78</b>	<b>0.67</b>	0.55	0.01	0.16	0.05
	concat(T + A)	0.44	0.66	0.53	0.02	0.15	0.04
	concat(A + V)	0.55	0.47	0.48	0.13	0.16	0.11
BiModal-Seq2Seq	T → V	0.67	<b>0.67</b>	<b>0.67</b>	0.26	0.22	<b>0.22</b>
	T → A	0.66	0.64	0.65	<b>0.28</b>	0.24	0.18
	A → T	0.55	0.60	0.56	0.17	<b>0.34</b>	0.11
	A → V	0.55	0.55	0.54	0.16	0.18	0.16
	V → T	0.58	0.58	0.58	0.05	0.16	0.08
	V → A	0.58	0.62	0.58	0.12	0.17	0.01

10 Point  
Boost

T = Text Modality, A = Audio Modality, V = Visual (facial) modality

# Results (Bi-Modal)

- Bimodal Baseline & Experimental Results

Method	Feature	BINARY (-1, +1)			7-CLASS (-3, ..., +3)		
		Prec	Recall	F1	Prec	Recall	F1
BiModal-Baseline	concat(T + V)	<b>0.78</b>	<b>0.67</b>	0.55	0.01	0.16	0.05
	concat(T + A)	0.44	0.66	0.53	0.02	0.15	0.04
	concat(A + V)	0.55	0.47	0.48	0.13	0.16	0.11
	T → V	0.67	<b>0.67</b>	<b>0.67</b>	0.26	0.22	<b>0.22</b>
BiModal-Seq2Seq	T → A	0.66	0.64	0.65	<b>0.28</b>	0.24	0.18
	A → T	0.55	0.60	0.56	0.17	<b>0.34</b>	0.11
	A → V	0.55	0.55	0.54	0.16	0.18	0.16
	V → T	0.58	0.58	0.58	0.05	0.16	0.08
	V → A	0.58	0.62	0.58	0.12	0.17	0.01

12 Point  
Boost

T = Text Modality, A = Audio Modality, V = Visual (facial) modality

# Results (Tri-Modal)

- Trimodal Baseline & Experimental Results

Method	Feature	BINARY (-1, +1)			7-CLASS (-3, ..., +3)		
		Prec	Recall	F1	Prec	Recall	F1
TriModal-Baseline	concat(T + V + A)	<b>0.75</b>	<b>0.75</b>	<b>0.75</b>	0.24	<b>0.27</b>	<b>0.24</b>
	embed(T, V) → A	0.56	0.60	0.57	0.10	0.16	0.09
	embed(T, A) → V	0.60	0.55	0.56	0.26	0.15	0.07
TriModal-Seq2Seq	embed(A, V) → T	0.66	0.53	0.44	0.16	0.04	0.09
	embed(A, T) → V	0.59	0.51	0.52	0.13	0.15	0.09
	embed(V, T) → A	0.59	0.60	0.59	0.11	0.17	0.10
	embed(V, A) → T	0.57	0.61	0.58	0.11	0.17	0.09
	concat(T, V) → A	0.67	0.66	0.65	0.22	0.17	0.18
	concat(A, T) → V	0.54	0.55	0.63	0.19	0.15	0.21
	concat(V, A) → T	0.59	0.59	0.58	0.16	0.12	0.12
	T → concat(A, V)	0.70	0.65	0.66	0.23	0.22	0.18
	A → concat(T, V)	0.55	0.53	0.54	0.18	0.20	0.18
	concat(T, A) → concat(T, V)	0.62	0.60	0.61	0.23	0.24	0.22
concat(T, V) → concat(T, A)	0.68	0.70	0.67	<b>0.31</b>	0.24	0.19	

T = Text Modality, A = Audio Modality, V = Visual (facial) modality

# Results - Takeaways

- We clearly outperform the baselines in the bi-modal domain
  - In the 7-class paradigm we often outperform by a large margin
  - For datasets without transcripts this approach may result in significant gains
- Slightly outperform baseline in tri-modal multiclass setting
- Significantly longer training times than the baseline alone

# Future Work

- Our method is unsupervised, we will pre-train seq2seq model with external dataset
- Use variational seq2seq to refine the training
- Further explore end-to-end training
- Explore additional methods for encoding sequences with 2 modalities
  - Multi-view LSTM

# Acknowledgements



- Amir Zadeh
- Volkan Cirik
- Louis-Philippe Morency
- Somya Wadhwa
- Minghai Chen
- Hieu Pham
- Workshop Reviewers
- \*AWS\*

**Thank you!**



# Appendix

# Problem Formulation

- Input:  $X = (X_1, X_2, \dots, X_n)$  where  $X_i = (X_i^{text}, X_i^{audio}, X_i^{video})$
- Output :  $Y = (Y_1, Y_2, \dots, Y_n)$ ,  $Y_i \in \mathbb{R}$
- Align based on word-level

$$X_i^{text} = (w_i^{(1)}, w_i^{(2)}, \dots, w_i^{(T_i)})$$

$$X_i^{audio} = (a_i^{(1)}, a_i^{(2)}, \dots, a_i^{(T_i)})$$

$$X_i^{video} = (v_i^{(1)}, v_i^{(2)}, \dots, v_i^{(T_i)})$$

- Goal: Learn the embedding representation

$$\widetilde{X}_i = f(X_i) = f((X_i^{text}, X_i^{audio}, X_i^{video}))$$

$$\widetilde{X}_i = f(X_i) = Seq2Seq\_Encoder(X_i)$$



# Problem Formulation (cont'd)

- Transformed input:  $\tilde{X} = (\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^T)$  with output  $Y = (y^1, y^2, \dots, y^T)$
- Using RNN with K hidden layers:

$$h = (h^1, h^2, \dots, h^K)$$
$$h^k = (h_1^k, h_2^k, \dots, h_D^k), k \in [1, K]$$

- First layer:  $h^1_t = H(W_{xh^1}\tilde{x}_t + W_{h^1h^1}h^1_{t-1} + b_{h^1})$
- Layer k :  $h^k_t = H(W_{h^{t-1}h^t}h^{k-1}_{t-1} + W_{h^k h^k}h^k_{t-1} + b_{h^k})$
- Using soft attention at last hidden layer K:

$$\alpha = \text{softmax} \left( \begin{bmatrix} W_\alpha h_1^K \\ W_\alpha h_2^K \\ \dots \\ W_\alpha h_T^K \end{bmatrix} \right)$$

# Problem Formulation (cont'd)

- Output of last hidden layer:  $A = [h_1^K, h_2^K, \dots, h_T^K] \alpha = H^K \alpha$
- Final output:  $\tilde{y}_t = W_{Ay} A + b_y$
- Mean Absolute Error Loss:  $\mathbb{L}_{MAE}(\tilde{Y}, Y) = \mathbb{E}[|\tilde{Y} - Y|]$
- Model is trained with SGD

# Seq2Seq Modality Translation

---

## Algorithm 1 Seq2Seq Modality Translation

$X, Y, S$  are 2 modalities and sentiment sequences

---

1: **Phase 1: Train Seq2Seq**

2:  $\mathcal{E}_{XY} \leftarrow \text{Seq2Seq\_RNN\_Encode}(X)$

3:  $\tilde{Y} \leftarrow \text{Seq2Seq\_RNN\_Decode}(\mathcal{E}_{XY})$

4:  $loss = \text{cross\_entropy}(\tilde{Y}, Y)$

5: Backprop to update params

6: **Phase 2: Sentiment Regression**

7:  $\mathcal{E}_{XY} \leftarrow \text{Seq2Seq\_RNN\_Encode}(X)$   $\triangleright$  trained  
encoder in Seq2Seq model

8:  $R = \text{RNN}(\mathcal{E}_{XY})$

9:  $score \leftarrow \text{Regression}(R)$

10:  $loss \leftarrow \text{MAE}(score, S)$

11: Backprop to update params

---

# Hierarchical Seq2Seq Modality Translation

---

**Algorithm 2 Hierarchical Seq2Seq Modality**

**Translation:**  $X, Y, Z, S$  are 3 modalities and sentiment sequences

---

1: **Phase 1: Train Seq2Seq for 2 modalities**

2:  $\mathcal{E}_{XY} \leftarrow \text{Seq2Seq\_RNN\_Encode}(X)$

3:  $\tilde{Y} \leftarrow \text{Seq2Seq\_RNN\_Decode}(\mathcal{E}_{XY})$

4:  $\text{loss} = \text{cross\_entropy}(\tilde{Y}, Y)$

5: Backpropagate to update parameters

6: **Phase 2: Train Seq2Seq for 3 modalities**

7:  $\mathcal{E}_{XYZ} \leftarrow \text{Seq2Seq\_RNN\_Encode}(\mathcal{E}_{XY})$

8:  $\tilde{Z} \leftarrow \text{Seq2Seq\_RNN\_Decode}(\mathcal{E}_{XYZ})$

9:  $\text{loss} = \text{cross\_entropy}(\tilde{Z}, Z)$

10: Backpropagate to update parameters

11: **Phase 3: Sentiment Regression**

12:  $\mathcal{E}_{XYZ} \leftarrow \text{Seq2Seq\_RNN\_Encode}(\mathcal{E}_{XY})$

13:  $R = \text{RNN}(\mathcal{E}_{XYZ})$

14:  $\text{score} \leftarrow \text{Regression}(R)$

15:  $\text{loss} \leftarrow \text{MAE}(\text{score}, S)$

16: Backpropagate to update parameters

---