

University of Stuttgart  
Germany



EUROPEAN CONFERENCE ON COMPUTER VISION

MILANO  
2024

# MST<sub>MIXER</sub> : Multi-Modal Video Dialog State Tracking in the Wild

Adnen Abdessaied, Lei Shi, Andreas Bulling

University of Stuttgart, Germany

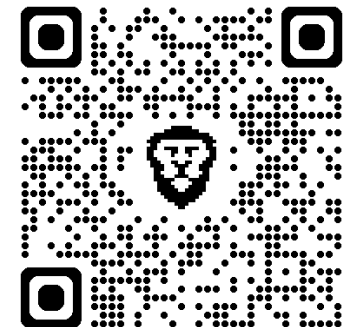
[adnen.abdessaied@vis.uni-stuttgart.de](mailto:adnen.abdessaied@vis.uni-stuttgart.de)

---

European Conference on Computer Vision (ECCV)

Milano, Italy 

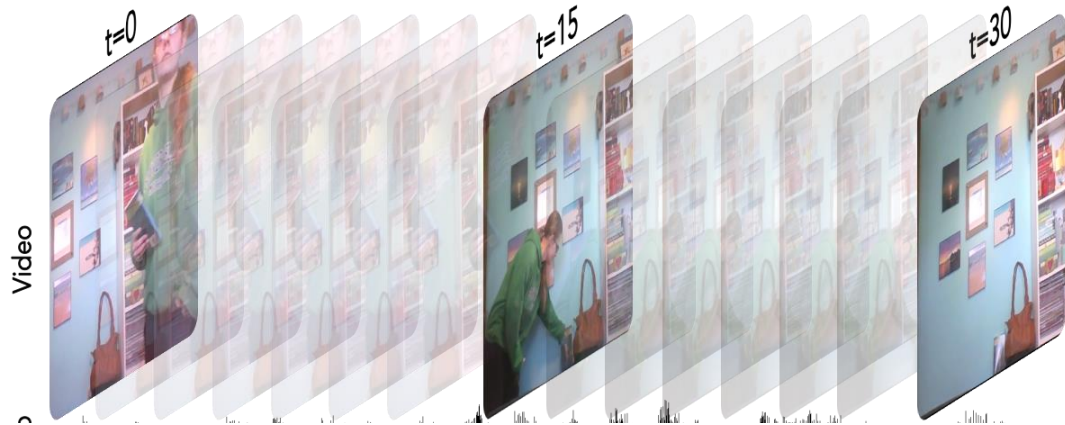
02.10.2024



# Introduction

## Video Dialog - Task formulation

# Introduction



**Video Dialog - Task formulation**  
Given a video,

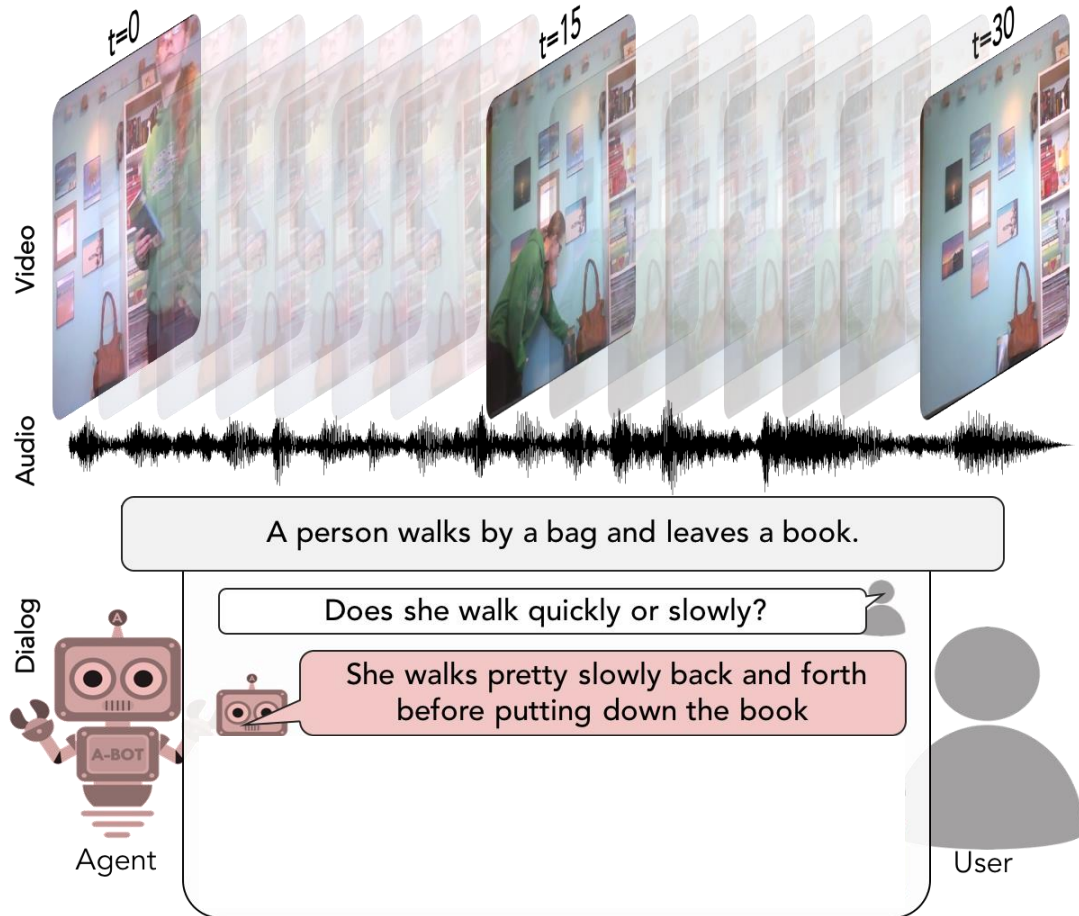
# Introduction



## Video Dialog - Task formulation

Given a video, audio data,

# Introduction

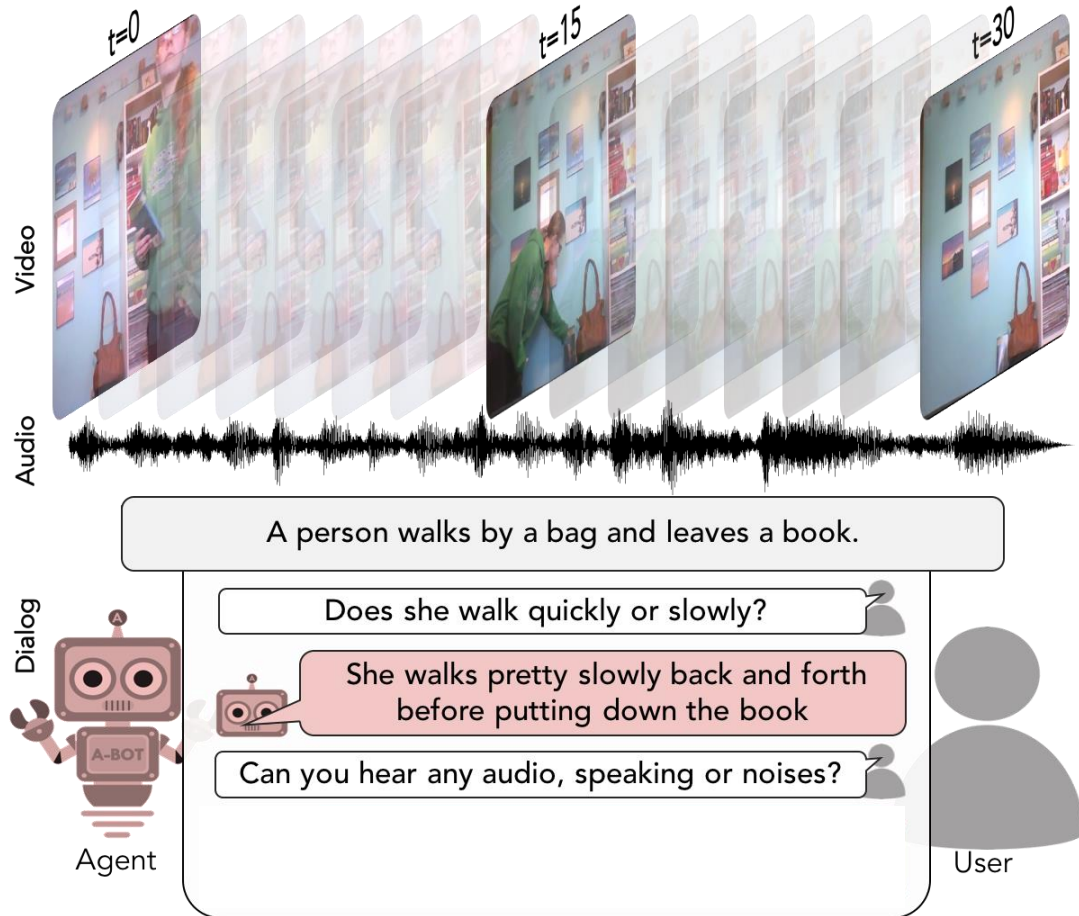


## Video Dialog - Task formulation

Given a video, audio data, a dialog history,

From <https://video-dialog.com/>

# Introduction

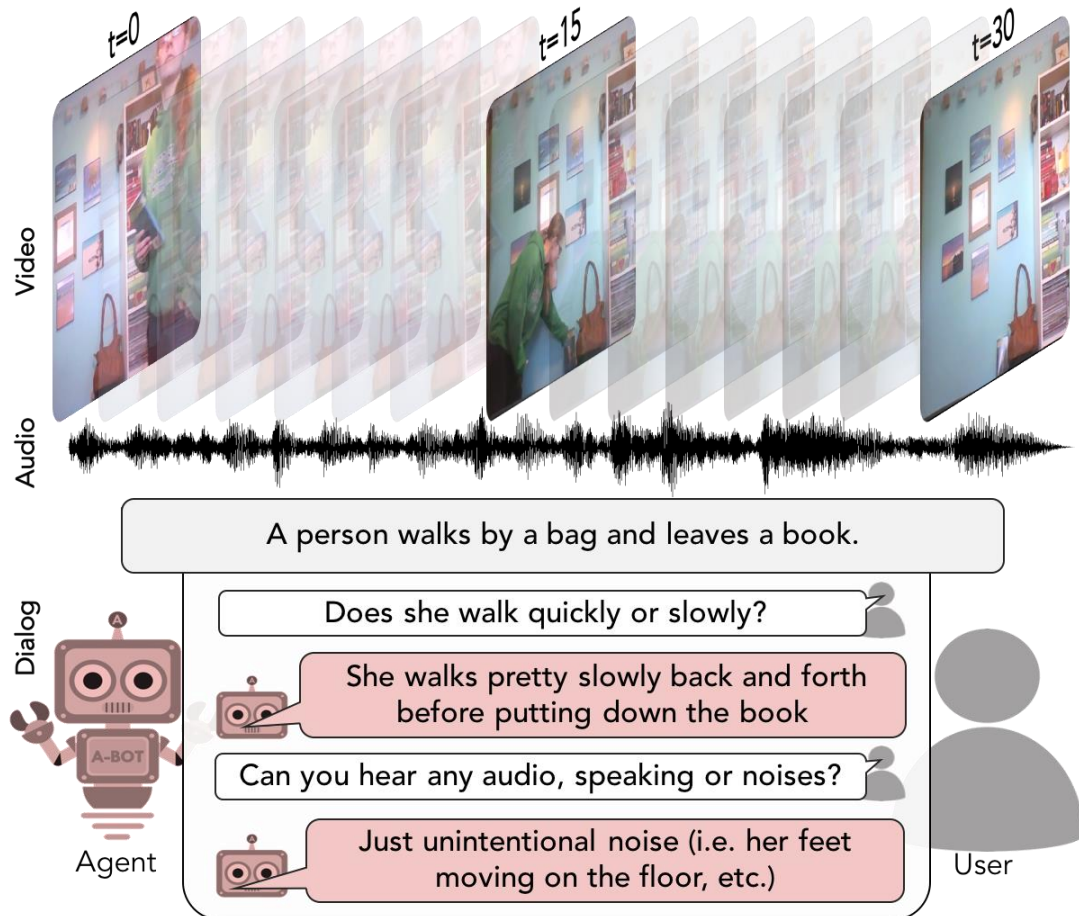


## Video Dialog - Task formulation

Given a video, audio data, a dialog history, and a question at time step  $t$ ,

From <https://video-dialog.com/>

# Introduction



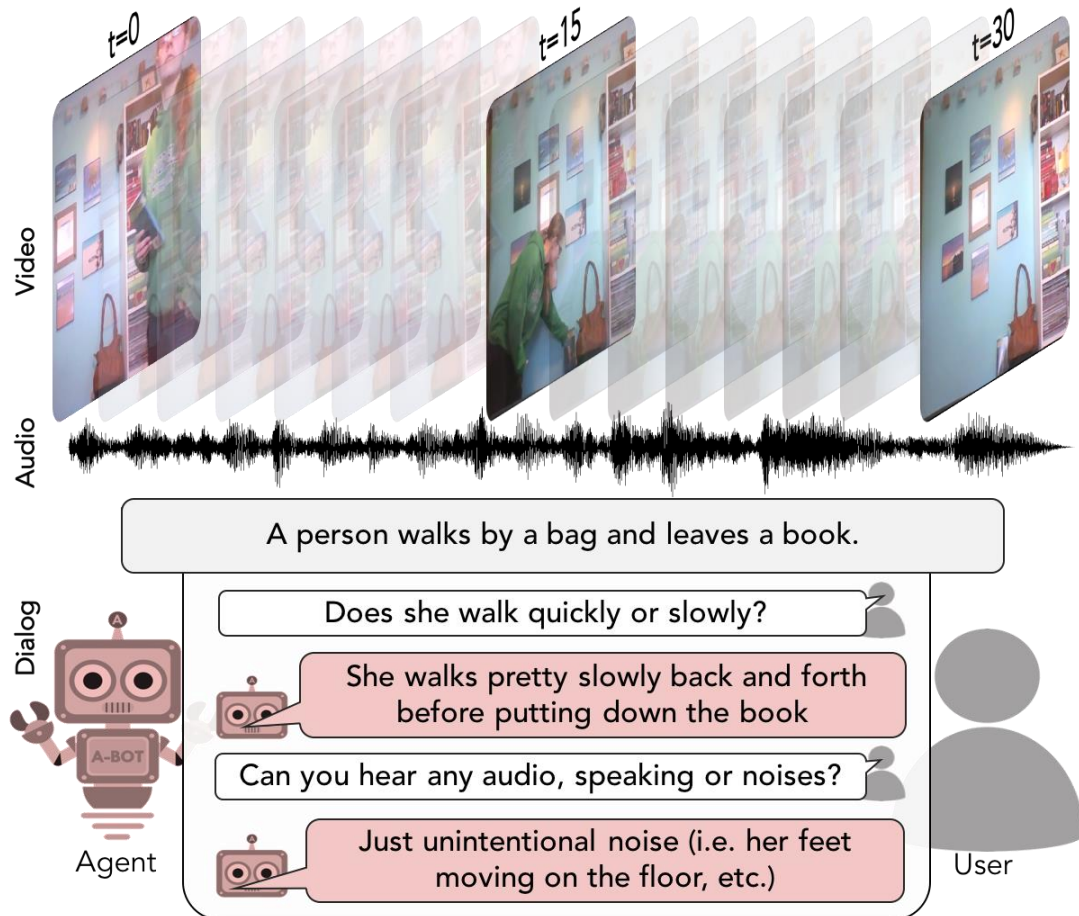
## Video Dialog - Task formulation

Given a video, audio data, a dialog history, and a question at time step  $t$ , generate an appropriate answer auto-regressively

From <https://video-dialog.com/>



# Introduction



From <https://video-dialog.com/>

## Video Dialog - Task formulation

Given a video, audio data, a dialog history, and a question at time step  $t$ , generate an appropriate answer auto-regressively



## Video Dialog is a natural extension to visual dialog

- Video vs image
- Complements videos with audio data



## More input modalities

# Motivation

Video Dialog [1] is a highly multi-modal task → More challenging than similar tasks

- **VQA [2] & VideoQA [3]:** Reasoning about the dialog history in addition to the question
- **Visual Dialog [4]:** Reasoning about a dynamic scene instead of a static image



**Dialog State Tracking (DST)** is crucial in building capable models

- **DST** was originally introduced to track and update users' goals in form of dialog states [5, 6]
- Now, it is broadly used to describe a model that keeps track of what it believes to be relevant for answering the question at hand

Research on DST has been predominately uni-modal in the form of slot-filling tasks [7, 8]



The current landscape necessitates extending DST to the multi-modal domain

Current models with “multi-modal” DST fall short in two major aspects:

- ⚡ track the constituent of only one modality within a multi-modal task [9, 10]  
→ uni-modal DST
- ⚡ limited to synthetic and automatically-generated datasets [11, 12, 13]  
→ do not reflect the complexity of real world scenarios

MST<sub>MIXER</sub> addresses the aforementioned limitations with

- ✓ modality-specific tracking blocks to identify the most relevant constituents of each modality
- ✓ a multi-modal GNN approach to learn the underlying structure between the mix of modalities

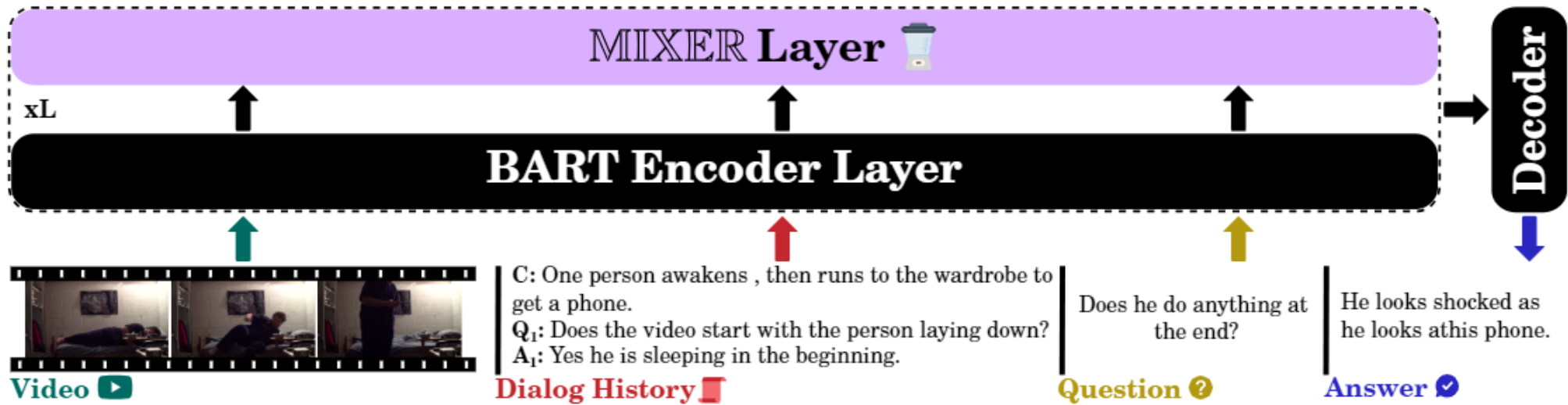


- ✓ Performs multi-modal state tracking in the real sense of the word
- ✓ Can tackle a wide-range of real-world datasets and benchmarks

# Method

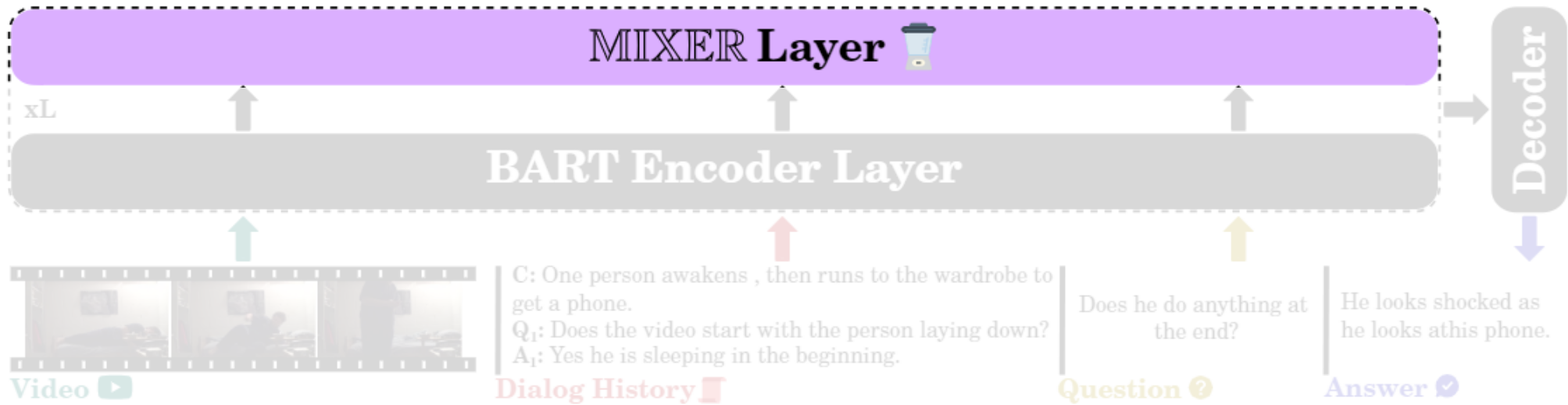
## Main Idea

- Perform multi-modal state tracking using MIXER layers
- Interleave BART encoder layers with MIXER layers → Enhance their hidden states

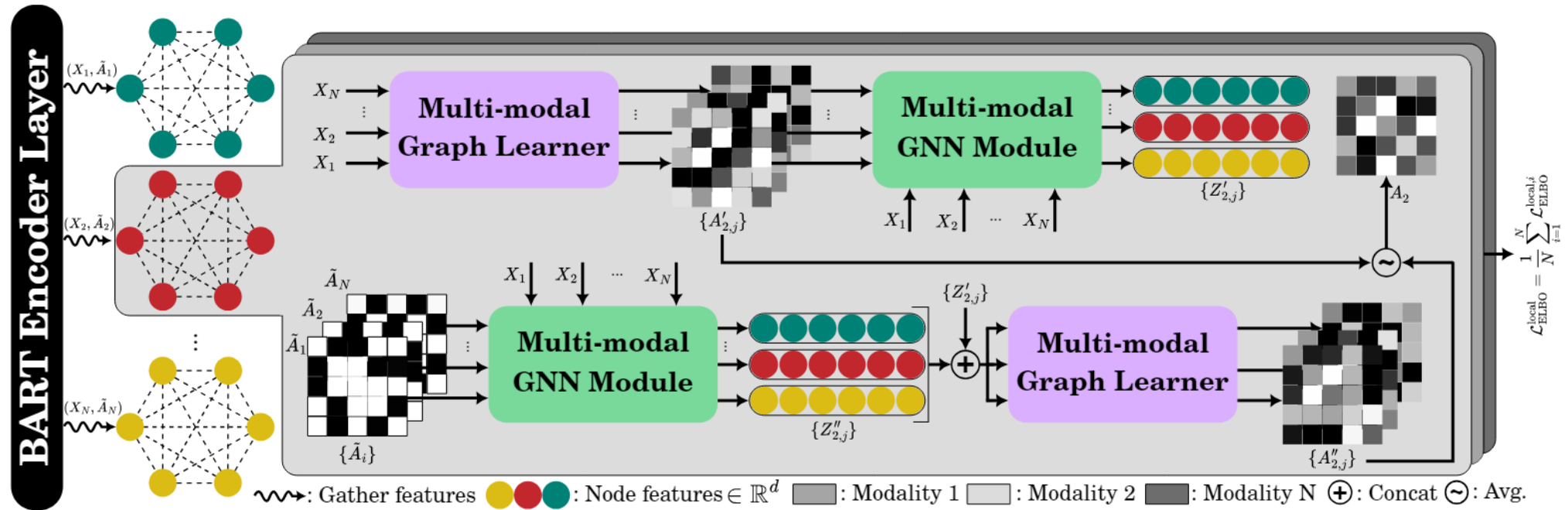


## MIXER layer

- Keeps track of the most relevant constituents of each modality at different semantic levels
- Employs a *divide-and-conquer* approach:  
local structures of individual modalities → global structure of the mix of all modalities

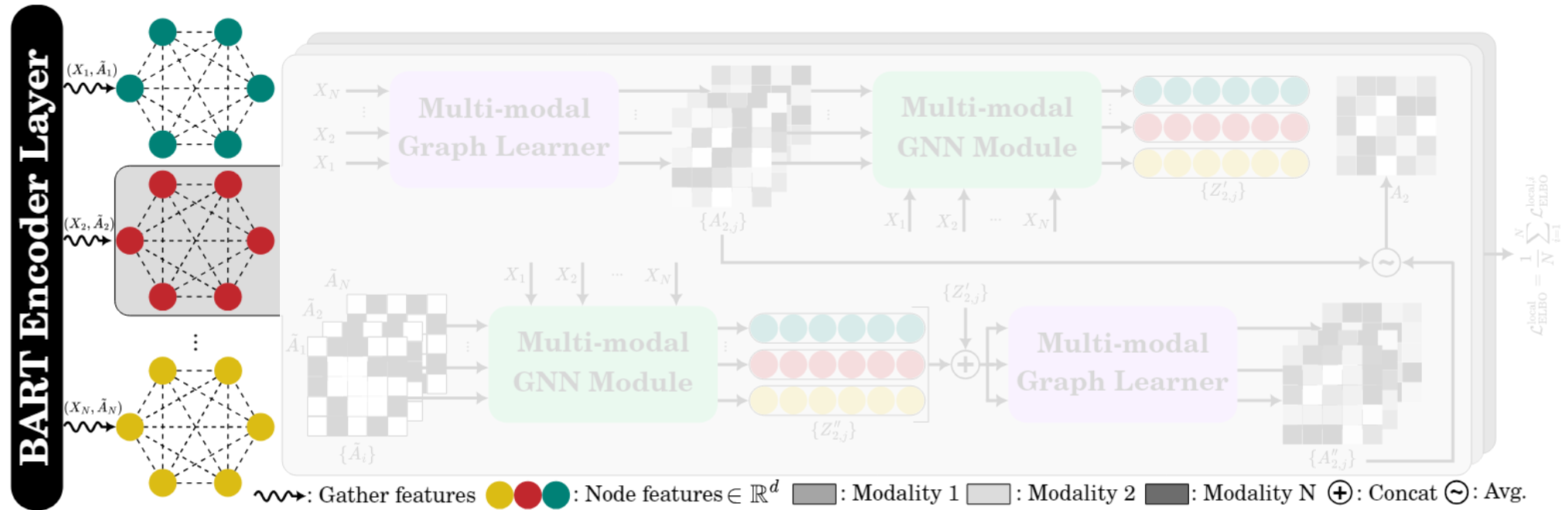


## Divide Stage



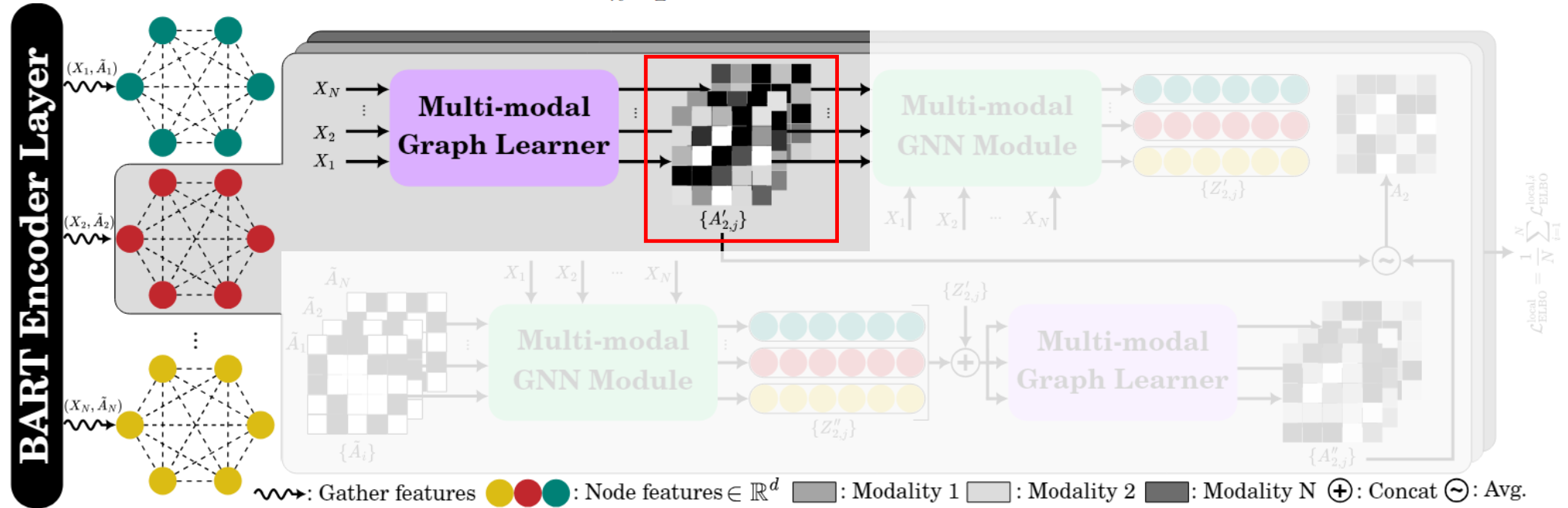


## Divide Stage

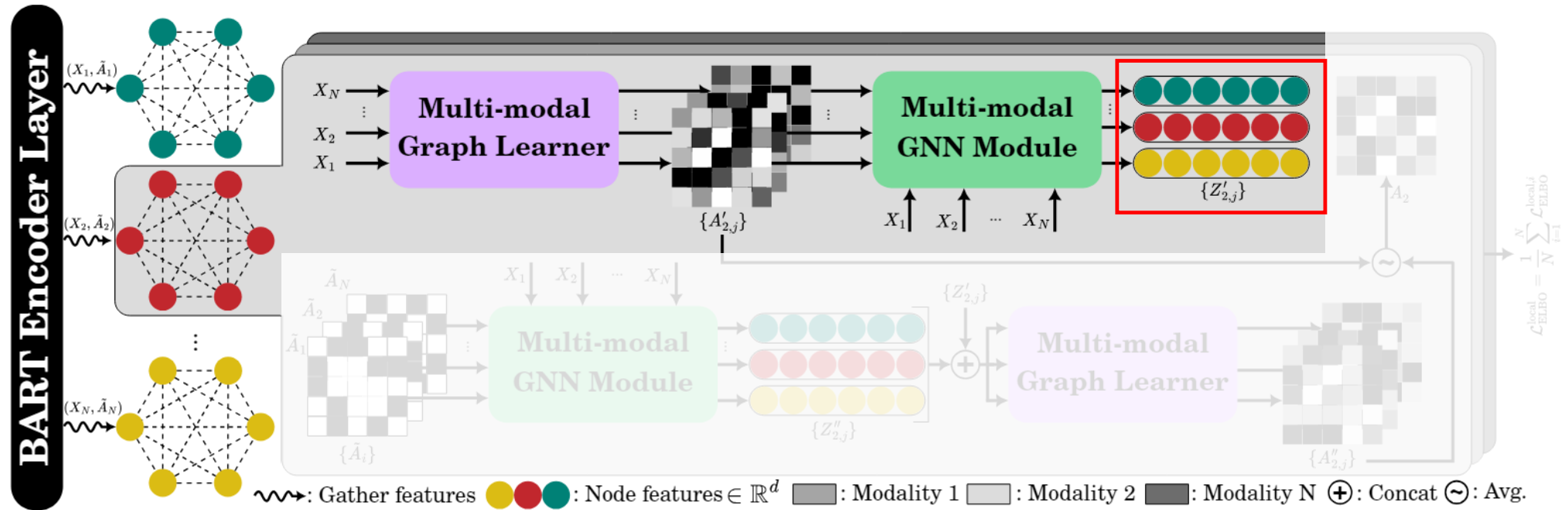


## Divide Stage

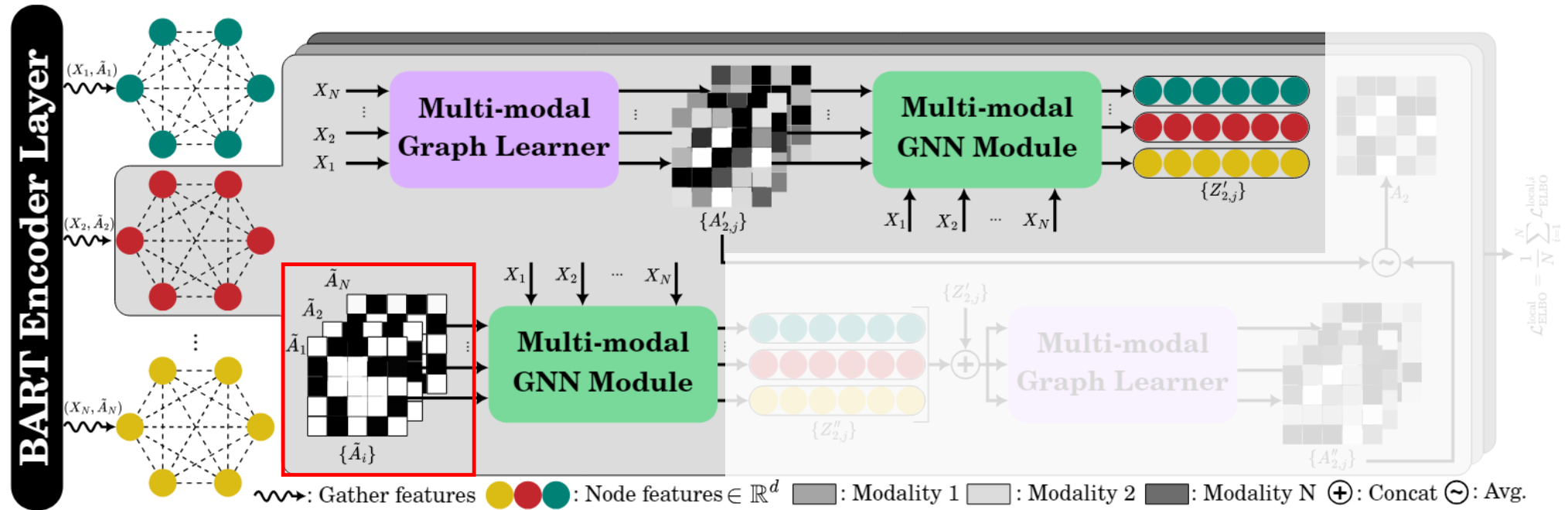
$$a'_{mn} = \frac{1}{K} \sum_{k=1}^K \cos(w_j^k \odot x_m, w_j^k \odot x_n)$$



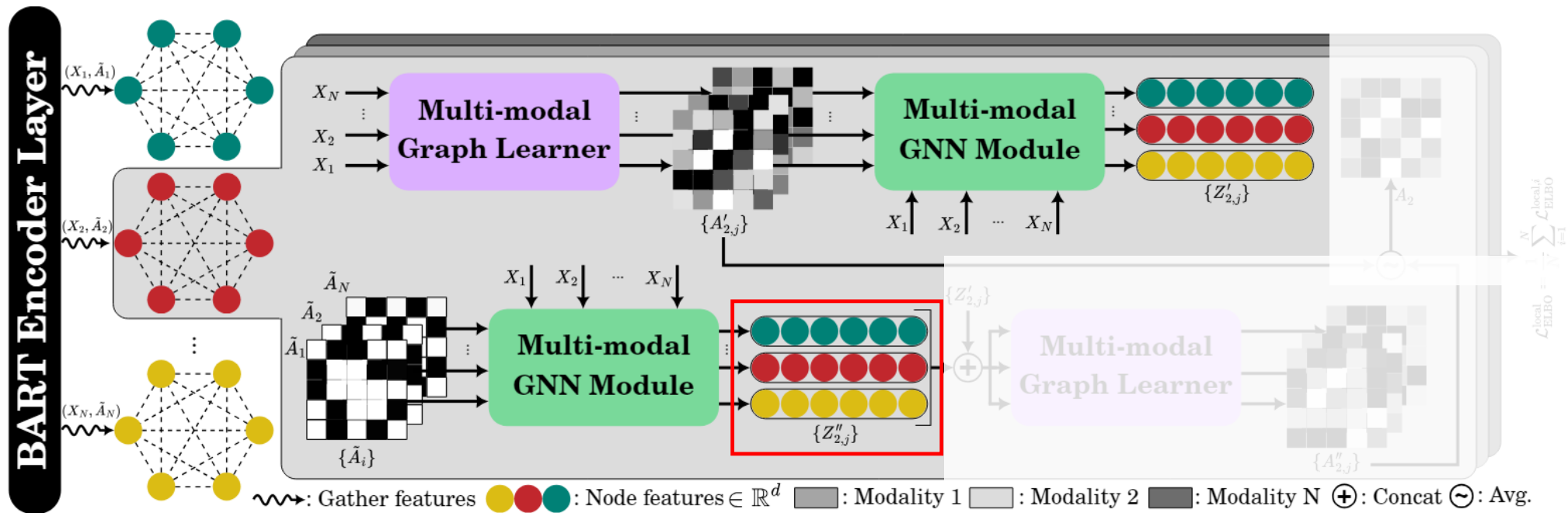
## Divide Stage



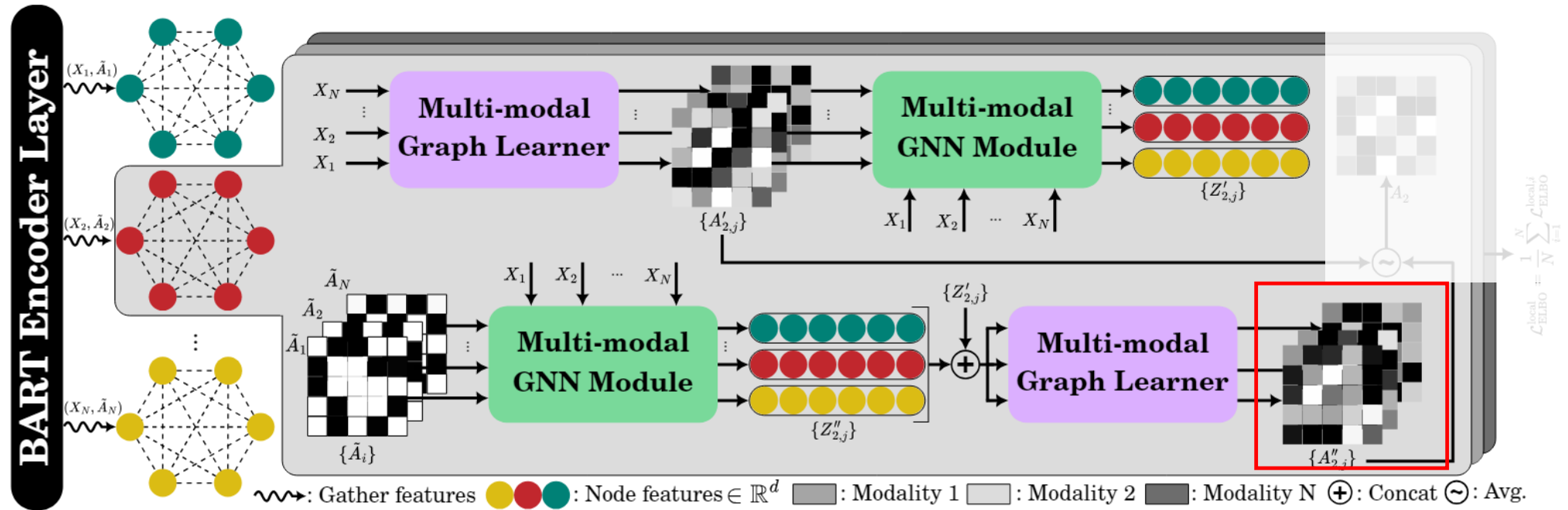
## Divide Stage



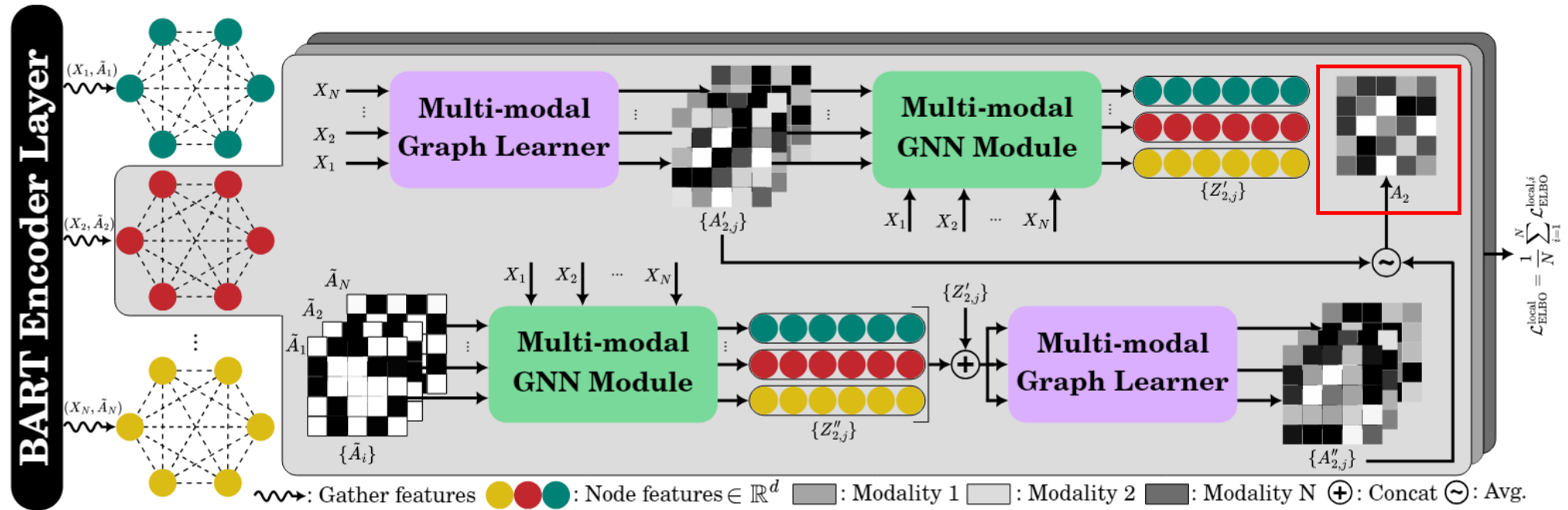
## Divide Stage



## Divide Stage



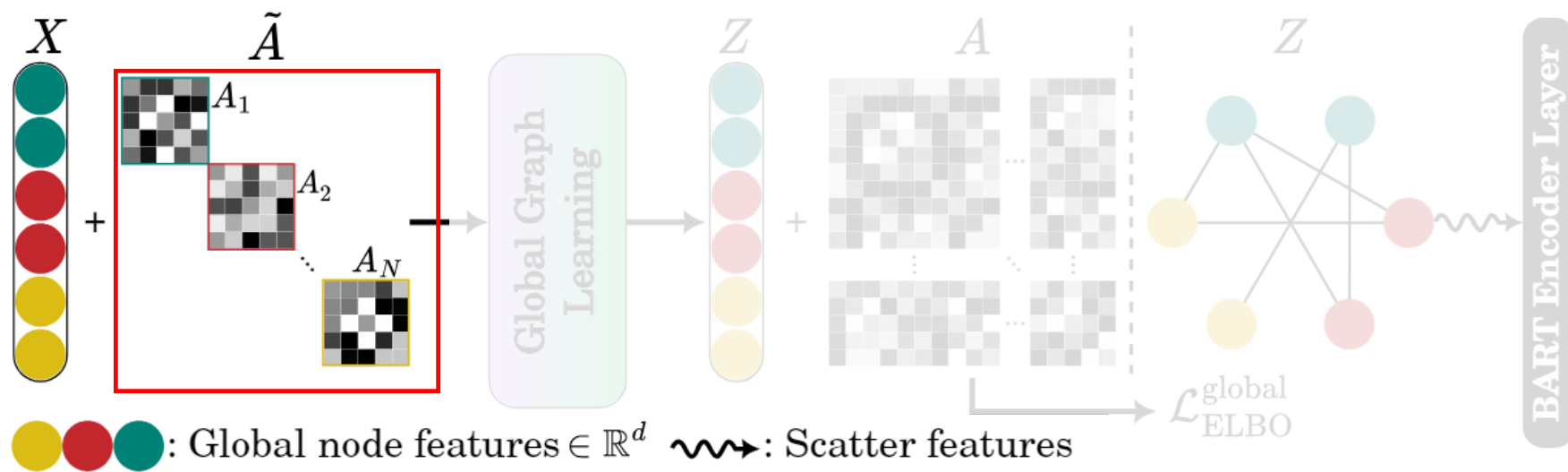
## Divide Stage



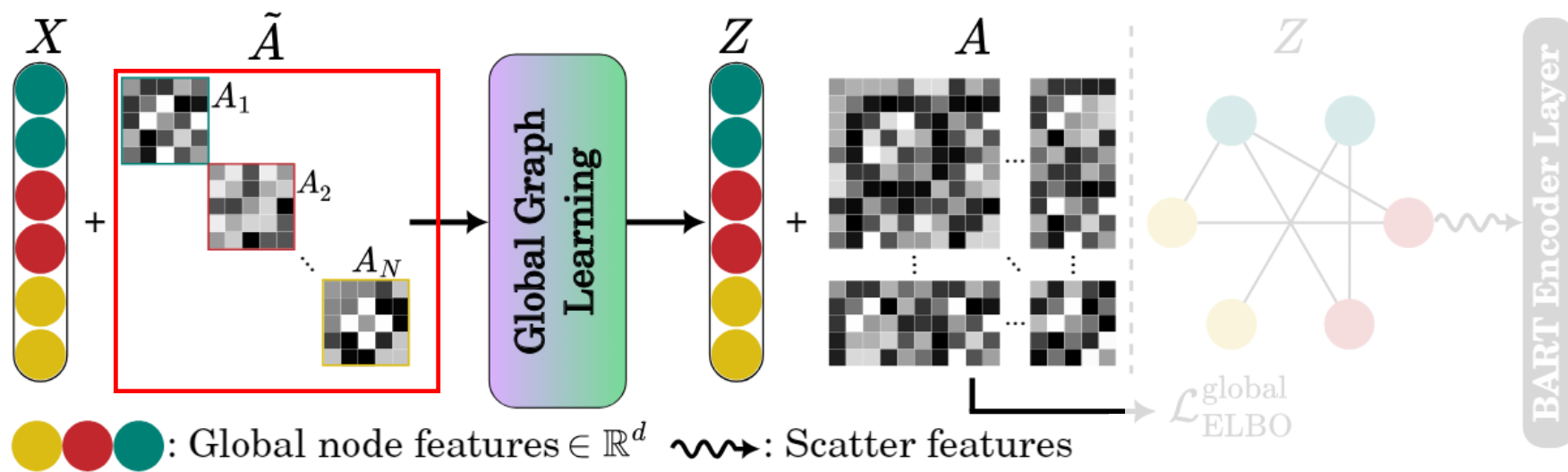
## Conquer Stage



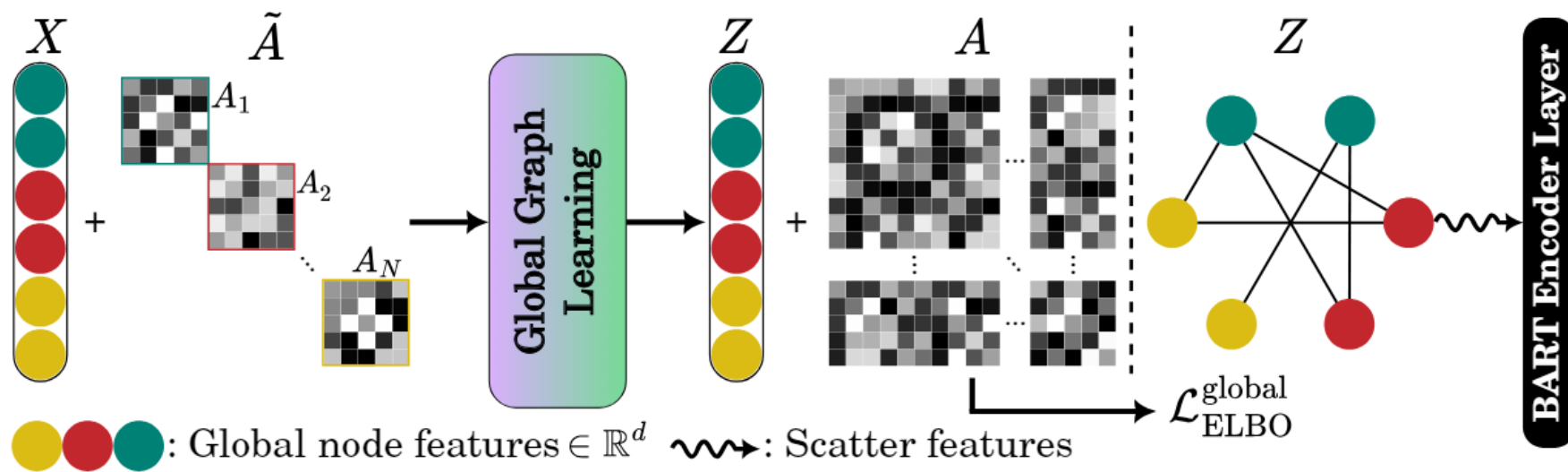
## Conquer Stage



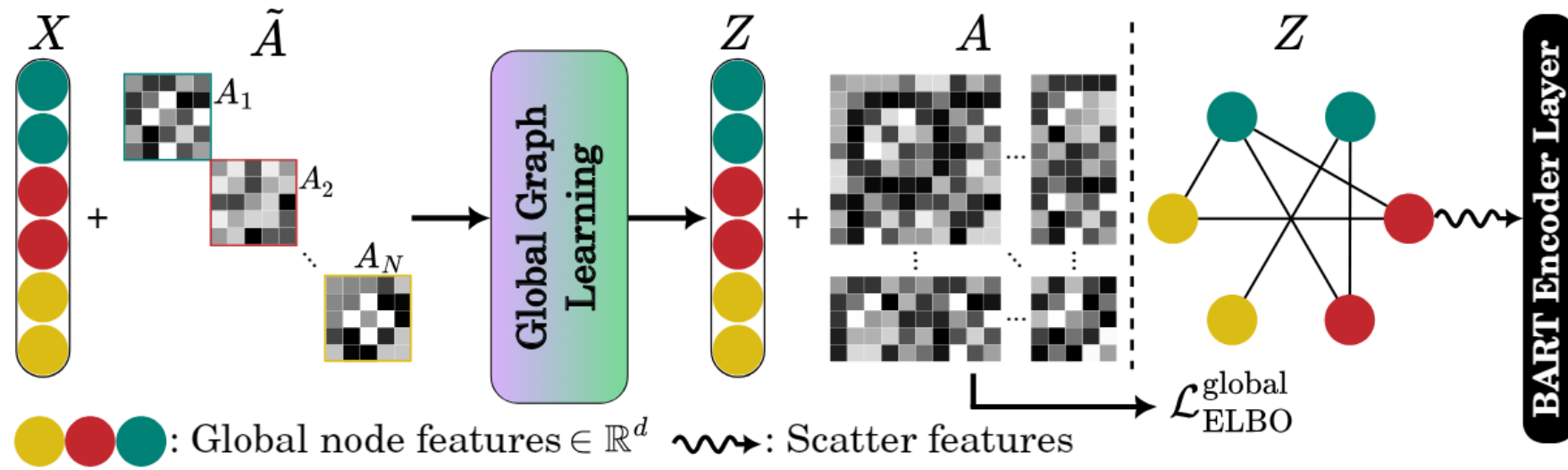
## Conquer Stage



## Conquer Stage

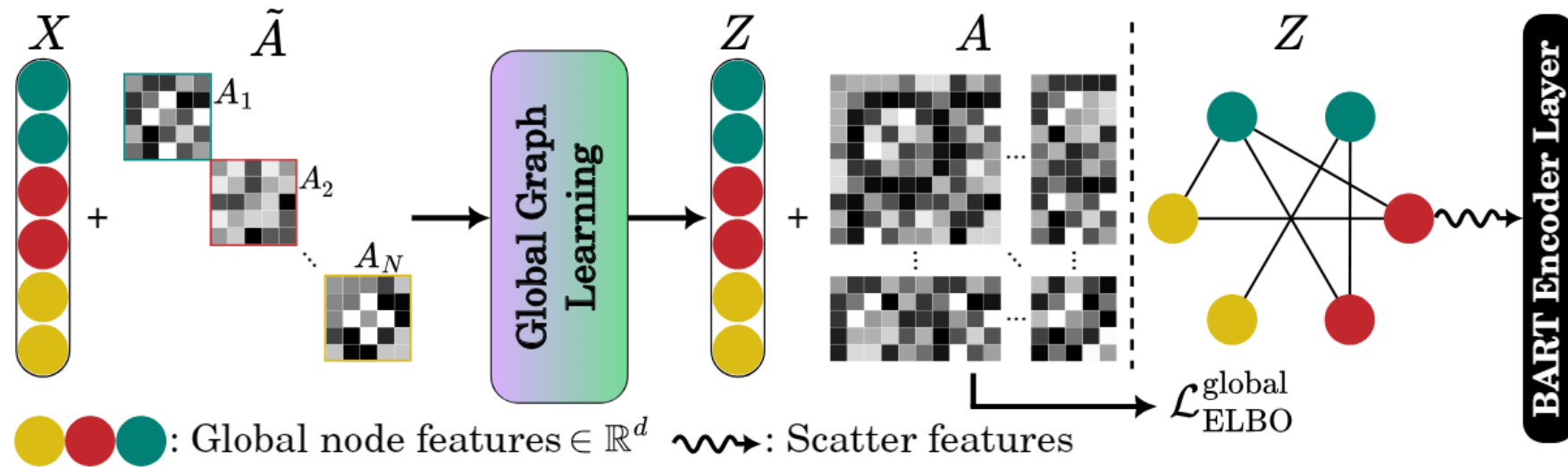


## Conquer Stage



$$H = (1 - \lambda)(H \oslash (Z, \text{Idx})) + \lambda H$$

## Conquer Stage



$$H = (1 - \lambda)(H \oslash (Z, \text{Idx})) + \lambda H$$

- $H$ : hidden states
- $\lambda \in [0, 1]$ : hyper-parameters
- $\oslash$ : PyTorch scatter operation
- $Z$ : global graph features
- $\text{Idx}$ : indices of  $Z$  w.r.t.  $H$



## Training

- We trained our model end-to-end using the following combination of losses:

## Training

- We trained our model end-to-end using the following combination of losses:

$$\mathcal{L} = \alpha_1 \mathcal{L}_{\text{gen}}$$

Next token prediction loss



## Training

- We trained our model end-to-end using the following combination of losses:

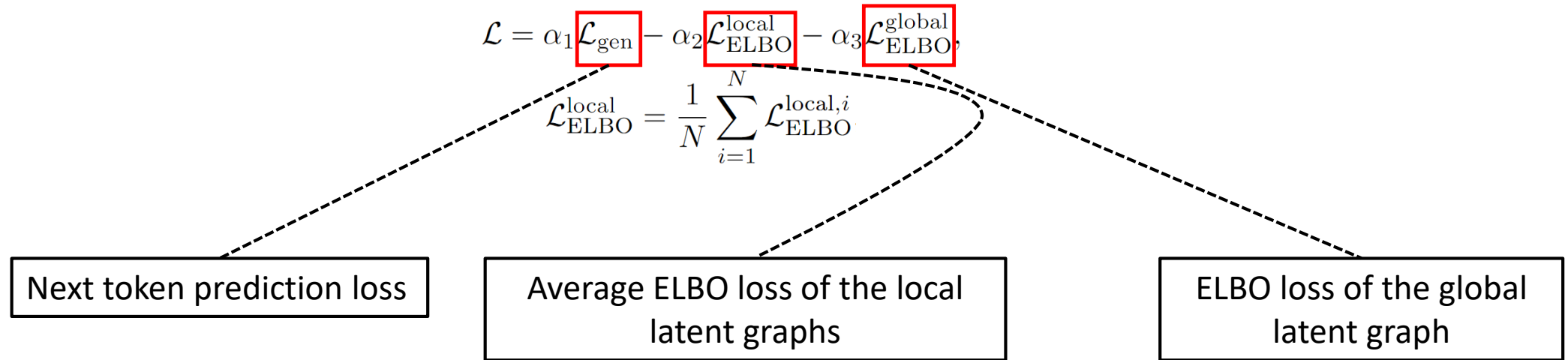
$$\mathcal{L} = \alpha_1 \mathcal{L}_{\text{gen}} - \alpha_2 \mathcal{L}_{\text{ELBO}}^{\text{local}}$$
$$\mathcal{L}_{\text{ELBO}}^{\text{local}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{ELBO}}^{\text{local},i}$$

Next token prediction loss

Average ELBO loss of the local latent graphs


## Training

- We trained our model end-to-end using the following combination of losses:





**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr

Model	Venue	AVSD-DSTC7						
		B-1	B-2	B-3	B-4	M	R	C
Baseline	<i>ICASSP'19</i>	62.1	48.0	37.9	30.5	21.7	48.1	73.3
MTN	<i>ACL'19</i>	71.5	58.1	47.6	39.2	26.9	55.9	106.6
JMAN	<i>AAAI'20</i>	66.7	52.1	41.3	33.4	23.9	53.3	94.1
VGD	<i>ACL'20</i>	74.9	62.0	52.0	43.6	28.2	58.2	119.4
BiST	<i>EMNLP'20</i>	75.5	61.9	51.0	42.9	28.4	58.1	119.2
SCGA	<i>AAAI'21</i>	74.5	62.2	51.7	43.0	28.5	57.8	120.1
RLM	<i>TASLP'21</i>	76.5	64.3	54.3	45.9	29.4	60.6	130.8
PDC	<i>ICLR'21</i>	77.0	65.3	53.9	44.9	29.2	60.6	129.5
AV-TRN	<i>ICASSP'22</i>	–	–	–	40.6	26.2	55.4	107.9
VGNMN	<i>NAACL'22</i>	–	–	–	42.9	27.8	57.8	118.8
COST	<i>ECCV'22</i>	72.3	58.9	48.3	40.0	26.6	56.1	108.5
MRLV	<i>NeurIPS'22</i>	–	59.2	49.3	41.5	26.9	56.9	115.9
THAM	<i>EMNLP'22</i>	77.8	65.4	54.9	46.8	30.8	61.9	133.5
DialogMCF	<i>TASLP'23</i>	77.7	65.3	54.7	45.7	30.6	61.3	135.2
ITR	<i>PAMI'23</i>	78.2	65.5	55.2	46.9	30.5	61.9	133.1
MST <sub>MIXER</sub> 	<i>ECCV'24</i>	<b>78.7</b>	<b>66.5</b>	<b>56.3</b>	<b>47.6</b>	<b>31.3</b>	<b>62.5</b>	<b>138.8</b>

**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr

Model	Venue	AVSD-DSTC7						
		B-1	B-2	B-3	B-4	M	R	C
Baseline	<i>ICASSP'19</i>	62.1	48.0	37.9	30.5	21.7	48.1	73.3
MTN	<i>ACL'19</i>	71.5	58.1	47.6	39.2	26.9	55.9	106.6
JMAN	<i>AAAI'20</i>	66.7	52.1	41.3	33.4	23.9	53.3	94.1
VGD	<i>ACL'20</i>	74.9	62.0	52.0	43.6	28.2	58.2	119.4
BiST	<i>EMNLP'20</i>	75.5	61.9	51.0	42.9	28.4	58.1	119.2
SCGA	<i>AAAI'21</i>	74.5	62.2	51.7	43.0	28.5	57.8	120.1
RLM	<i>TASLP'21</i>	76.5	64.3	54.3	45.9	29.4	60.6	130.8
PDC	<i>ICLR'21</i>	77.0	65.3	53.9	44.9	29.2	60.6	129.5
AV-TRN	<i>ICASSP'22</i>	–	–	–	40.6	26.2	55.4	107.9
VGNMN	<i>NAACL'22</i>	–	–	–	42.9	27.8	57.8	118.8
COST	<i>ECCV'22</i>	72.3	58.9	48.3	40.0	26.6	56.1	108.5
MRLV	<i>NeurIPS'22</i>	–	59.2	49.3	41.5	26.9	56.9	115.9
THAM	<i>EMNLP'22</i>	77.8	65.4	54.9	46.8	<u>30.8</u>	<u>61.9</u>	133.5
DialogMCF	<i>TASLP'23</i>	<u>77.7</u>	<u>65.3</u>	<u>54.7</u>	<u>45.7</u>	<u>30.6</u>	<u>61.3</u>	<u>135.2</u>
ITR	<i>PAMI'23</i>	<u>78.2</u>	<u>65.5</u>	<u>55.2</u>	<u>46.9</u>	<u>30.5</u>	<u>61.9</u>	<u>133.1</u>
MST <sub>MIXER</sub>	<i>ECCV'24</i>	<b>78.7</b>	<b>66.5</b>	<b>56.3</b>	<b>47.6</b>	<b>31.3</b>	<b>62.5</b>	<b>138.8</b>

**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr

Model	Venue	AVSD-DSTC8						
		B-1	B-2	B-3	B-4	M	R	C
Baseline	<i>ICASSP'19</i>	61.4	46.7	36.5	28.9	21.0	48.0	65.1
MTN	<i>ACL'19</i>	–	–	–	–	–	–	–
JMAN	<i>AAAI'20</i>	64.5	50.4	40.2	32.4	23.2	52.1	87.5
VGD	<i>ACL'20</i>	–	–	–	–	–	–	–
BiST	<i>EMNLP'20</i>	68.4	54.8	45.7	37.6	27.3	56.3	101.7
SCGA	<i>AAAI'21</i>	71.1	59.3	49.7	41.6	27.6	56.6	112.3
RLM	<i>TASLP'21</i>	74.6	62.6	52.8	44.5	28.6	59.8	124.0
PDC	<i>ICLR'21</i>	74.9	62.9	52.8	43.9	28.5	59.2	120.1
AV-TRN	<i>ICASSP'22</i>	–	–	–	39.4	25.0	54.5	99.7
VGNMN	<i>NAACL'22</i>	–	–	–	–	–	–	–
COST	<i>ECCV'22</i>	69.5	55.9	46.5	3.82	27.8	57.4	105.1
MRLV	<i>NeurIPS'22</i>	–	–	–	–	–	–	–
THAM	<i>EMNLP'22</i>	<u>76.4</u>	<u>64.1</u>	<u>53.8</u>	<u>45.5</u>	<u>30.1</u>	<u>61.0</u>	<u>130.4</u>
DialogMCF	<i>TASLP'23</i>	75.6	63.3	53.2	44.9	29.3	60.1	125.3
ITR	<i>PAMI'23</i>	76.2	<u>64.1</u>	<u>54.3</u>	<u>46.0</u>	29.8	60.7	128.5
MST <sub>MIXER</sub>	<i>ECCV'24</i>	<b>77.5</b>	<b>66.0</b>	<b>56.1</b>	<b>47.7</b>	<b>30.6</b>	<b>62.4</b>	<b>135.4</b>

**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr


Model	Venue	AVSD-DSTC8						
		B-1	B-2	B-3	B-4	M	R	C
Baseline	<i>ICASSP'19</i>	61.4	46.7	36.5	28.9	21.0	48.0	65.1
MTN	<i>ACL'19</i>	–	–	–	–	–	–	–
JMAN	<i>AAAI'20</i>	64.5	50.4	40.2	32.4	23.2	52.1	87.5
VGD	<i>ACL'20</i>	–	–	–	–	–	–	–
BiST	<i>EMNLP'20</i>	68.4	54.8	45.7	37.6	27.3	56.3	101.7
SCGA	<i>AAAI'21</i>	71.1	59.3	49.7	41.6	27.6	56.6	112.3
RLM	<i>TASLP'21</i>	74.6	62.6	52.8	44.5	28.6	59.8	124.0
PDC	<i>ICLR'21</i>	74.9	62.9	52.8	43.9	28.5	59.2	120.1
AV-TRN	<i>ICASSP'22</i>	–	–	–	39.4	25.0	54.5	99.7
VGNMN	<i>NAACL'22</i>	–	–	–	–	–	–	–
COST	<i>ECCV'22</i>	69.5	55.9	46.5	3.82	27.8	57.4	105.1
MRLV	<i>NeurIPS'22</i>	–	–	–	–	–	–	–
THAM	<i>EMNLP'22</i>	<u>76.4</u>	<u>64.1</u>	53.8	45.5	<u>30.1</u>	<u>61.0</u>	<u>130.4</u>
DialogMCF	<i>TASLP'23</i>	75.6	63.3	53.2	44.9	29.3	60.1	125.3
ITR	<i>PAMI'23</i>	76.2	<u>64.1</u>	<u>54.3</u>	<u>46.0</u>	29.8	60.7	128.5
MST <sub>MIXER</sub>	<i>ECCV'24</i>	<b>77.5</b>	<b>66.0</b>	<b>56.1</b>	<b>47.7</b>	<b>30.6</b>	<b>62.4</b>	<b>135.4</b>





# Results


**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr

Model	Venue	B-1	B-2	B-3	B-4	M	R	C
AV-TRN	<i>ICASSP'22</i>	–	–	–	24.7	19.1	43.7	56.6
+ Ext.	<i>ICASSP'22</i>	–	–	–	37.1	24.5	53.5	86.9
DSTC10	<i>AAAI'22</i>	67.3	54.5	44.8	<u>37.2</u>	24.3	53.0	<u>91.2</u>
DialogMCF	<i>TASLP'23</i>	<u>69.3</u>	<u>55.6</u>	<u>45.0</u>	36.9	<u>24.9</u>	<u>53.6</u>	<u>91.2</u>
MST <sub>MIXER</sub> 	<i>ECCV'24</i>	<b>70.0</b>	<b>57.4</b>	<b>47.6</b>	<b>40.0</b>	<b>25.7</b>	<b>54.5</b>	<b>99.8</b>


## AVSD-DSTC10

# Results

**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr

Model	Venue	B-1	B-2	B-3	B-4	M	R	C
AV-TRN	<i>ICASSP'22</i>	–	–	–	24.7	19.1	43.7	56.6
+ Ext.	<i>ICASSP'22</i>	–	–	–	37.1	24.5	53.5	86.9
DSTC10	<i>AAAI'22</i>	67.3	54.5	44.8	<u>37.2</u>	24.3	53.0	<u>91.2</u>
DialogMCF	<i>TASLP'23</i>	<u>69.3</u>	<u>55.6</u>	<u>45.0</u>	36.9	<u>24.9</u>	<u>53.6</u>	<u>91.2</u>
MST <sub>MIXER</sub> 	<i>ECCV'24</i>	<b>70.0</b>	<b>57.4</b>	<b>47.6</b>	<b>40.0</b>	<b>25.7</b>	<b>54.5</b>	<b>99.8</b>


**AVSD-DSTC10**

Model	Venue	B-4
MTN	<i>ACL'19</i>	21.7
GPT-2	<i>EMNLP'21</i>	19.2
BART	<i>NAACL'22</i>	33.1
PaCE	<i>ACL'23</i>	<u>34.1</u>
MST <sub>MIXER</sub> 	<i>ECCV'24</i>	<b>44.7</b>


**SIMMC 2.0**

# Results

**B-n** = BLEU-n, **M** = METEOR, **R** = Rouge-L, **C** = CIDEr


Model	Venue	B-1	B-2	B-3	B-4	M	R	C
AV-TRN	<i>ICASSP'22</i>	–	–	–	24.7	19.1	43.7	56.6
+ Ext.	<i>ICASSP'22</i>	–	–	–	37.1	24.5	53.5	86.9
DSTC10	<i>AAAI'22</i>	67.3	54.5	44.8	<u>37.2</u>	24.3	53.0	<u>91.2</u>
DialogMCF	<i>TASLP'23</i>	<u>69.3</u>	<u>55.6</u>	<u>45.0</u>	36.9	<u>24.9</u>	<u>53.6</u>	<u>91.2</u>
<b>MST<sub>MIXER</sub></b> 	<i>ECCV'24</i>	<b>70.0</b>	<b>57.4</b>	<b>47.6</b>	<b>40.0</b>	<b>25.7</b>	<b>54.5</b>	<b>99.8</b>

**AVSD-DSTC10**

Model	Venue	B-4
MTN	<i>ACL'19</i>	21.7
GPT-2	<i>EMNLP'21</i>	19.2
BART	<i>NAACL'22</i>	33.1
PaCE	<i>ACL'23</i>	<u>34.1</u>
<b>MST<sub>MIXER</sub></b> 	<i>ECCV'24</i>	<b>44.7</b>

**SIMMC 2.0**

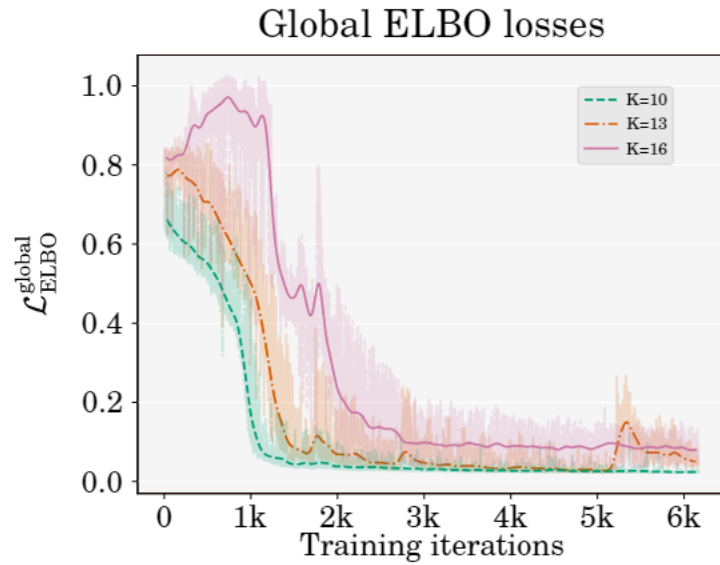
**C** = Causal, **T** = Temporal, **D** = Descriptive questions

Model	Venue	WUPS <sub>C</sub>	WUPS <sub>T</sub>	WUPS <sub>D</sub>	WUPS
HCRN	<i>CVPR'20</i>	16.05	17.68	49.78	23.92
HGA	<i>AAAI'20</i>	<u>17.98</u>	<u>17.95</u>	<u>50.84</u>	24.06
Flamingo	<i>NeurIPS'22</i>	–	–	–	<u>28.40</u>
KcGA	<i>AAAI'23</i>	–	–	–	28.20
EMU	<i>arXiv'23</i>	–	–	–	23.40
<b>MST<sub>MIXER</sub></b> 	<i>ECCV'24</i>	<b>22.12</b>	<b>22.20</b>	<b>55.64</b>	<b>29.50</b>

**NExT-QA (open-ended)**

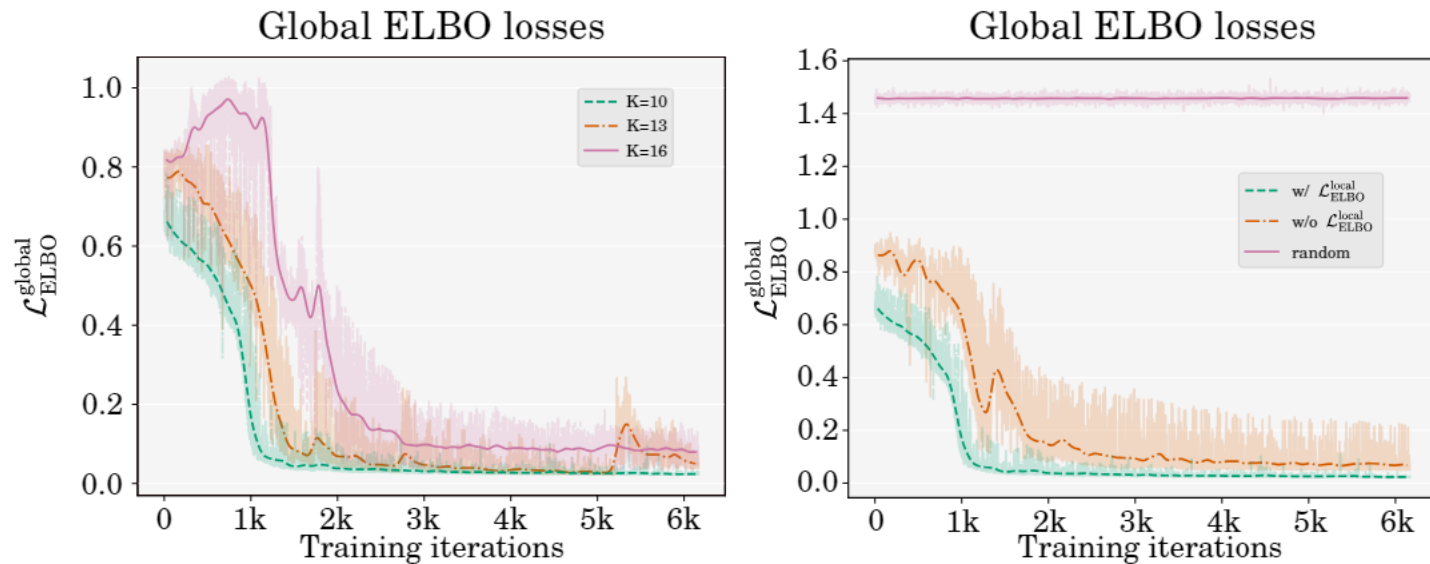
# Ablation Study

# Ablation Study



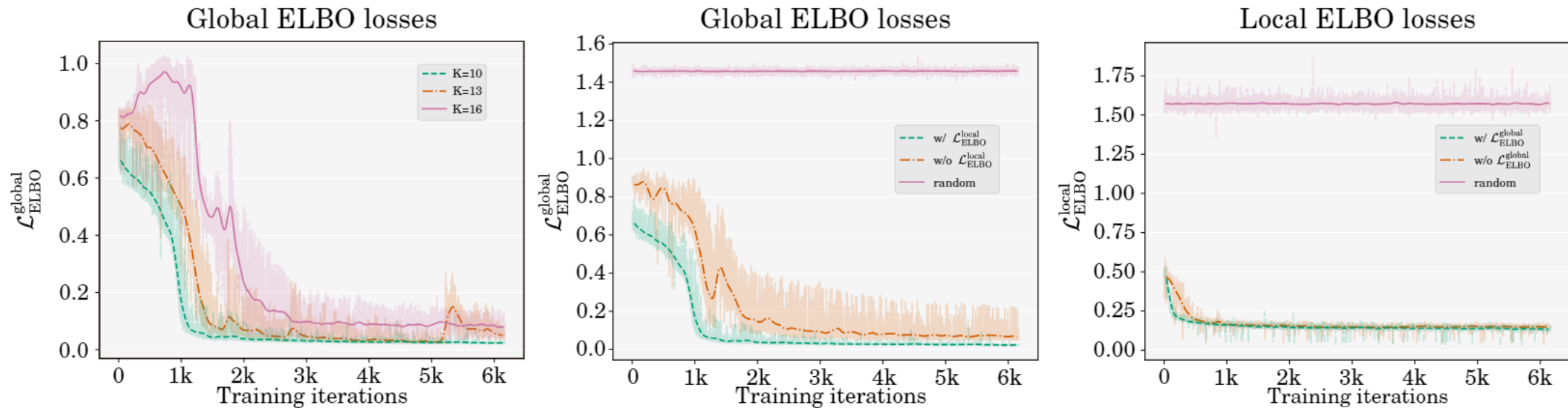
- Learning the latent graphs become more difficult when using a higher number of nodes  $K$

# Ablation Study



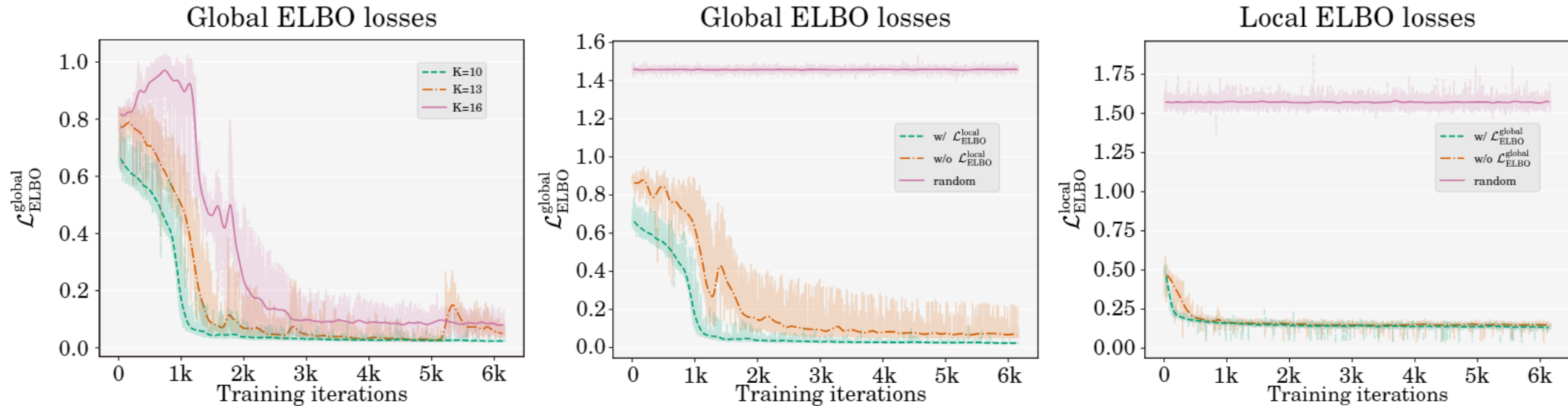
- Learning the latent graphs become more difficult when using a higher number of nodes  $K$
- The *divide* stage alleviates the difficulty of learning the global latent graphs

# Ablation Study



- Learning the latent graphs become more difficult when using a higher number of nodes  $K$
- The *divide* stage alleviates the difficulty of learning the global latent graphs
- The *conquer* stage slightly improves the learning of the local latent graphs in the *divide* stage

# Ablation Study

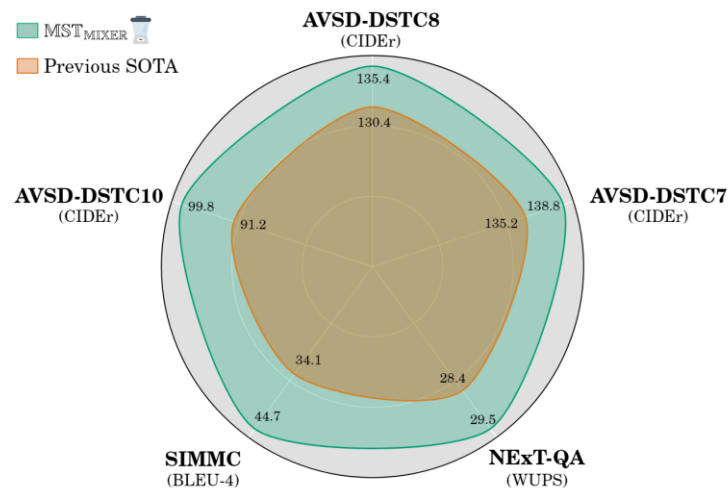


- Learning the latent graphs become more difficult when using a higher number of nodes  $K$
  - The *divide* stage alleviates the difficulty of learning the global latent graphs
  - The *conquer* stage slightly improves the learning of the local latent graphs in the *divide* stage
- **Conquer stage benefits more from divide stage than vice-versa**



# Summary

- We proposed  $MST_{MIXER}$ : A novel multi-modal state tracking model specifically geared towards video dialog
- It first identifies the most influential constituents at different semantic levels
- Then, it relies on a *divide-and-conquer* GNN-based approach to infer the missing underlying structure of the mix of all modalities
- Finally, it leverages these features to augment the hidden states of a backbone VLM
- $MST_{MIXER}$  achieves new SOTA results on a variety of challenging benchmarks



# References

1. Alamri, H., Cartillier, V., Das, A., et al.: *Audio visual scene-aware dialog*. In **CVPR** (2019)
2. Antol, S., Agrawal, A., Lu, J., et al.: *VQA: Visual Question Answering*. In **ICCV** (2015)
3. Xu, D., Zhao, Z., et al.: *Video question answering via gradually refined attention over appearance and motion*. In **ACM MM** (2017)
4. Das, A., Kottur, S., Gupta, K., et al.: *Visual Dialog*. In **CVPR** (2017)
5. Lee, H., Lee, J., Kim, T.Y.: *SUMBT: Slot-utterance matching for universal and scalable belief tracking*. In **ACL** (2019)
6. Wu, C.S., Madotto, A., et al.: *Transferable multi-domain state generator for task-oriented dialogue systems*. In **ACL** (2019)
7. Le, H., Socher, R., Hoi, S.C.: *Non-autoregressive dialog state tracking*. In **ICLR** (2020)
8. Mrkšić, N., Ó Séaghdha, D., et al.: *Neural belief tracker: Data-driven dialogue state tracking*. In **ACL** (2017)
9. Le, H., Chen, N.F., Hoi, S.C.: *Multimodal Dialogue State Tracking*. In **NAACL** (2022)
10. Pang, W., Wang, X.: *Visual dialogue state tracking for question generation*. In **AAAI** (2020)
11. Abdessaied, A., Hochmeister, M., Bulling, A.: *OLViT: Multi-modal state tracking via attention-based embeddings for video-grounded dialog*. In **LREC-COLING** (2024)
12. Kottur, S., Moon, S., et al.: *SIMMC 2.0: A task-oriented dialog dataset for immersive multimodal conversations*. In **EMNLP** (2021)
13. Moon, S., Kottur, S., Crook, P.A., et al.: *Situated and interactive multimodal conversations*. In **COLING** (2020)



University of Stuttgart  
Germany



EUROPEAN CONFERENCE ON COMPUTER VISION

MILANO  
2024

# Thank You!

---

Adnen Abdessaied, Lei Shi, Andreas Bulling

University of Stuttgart, Germany

[adnen.abdessaied@vis.uni-stuttgart.de](mailto:adnen.abdessaied@vis.uni-stuttgart.de)

02.10.2024

