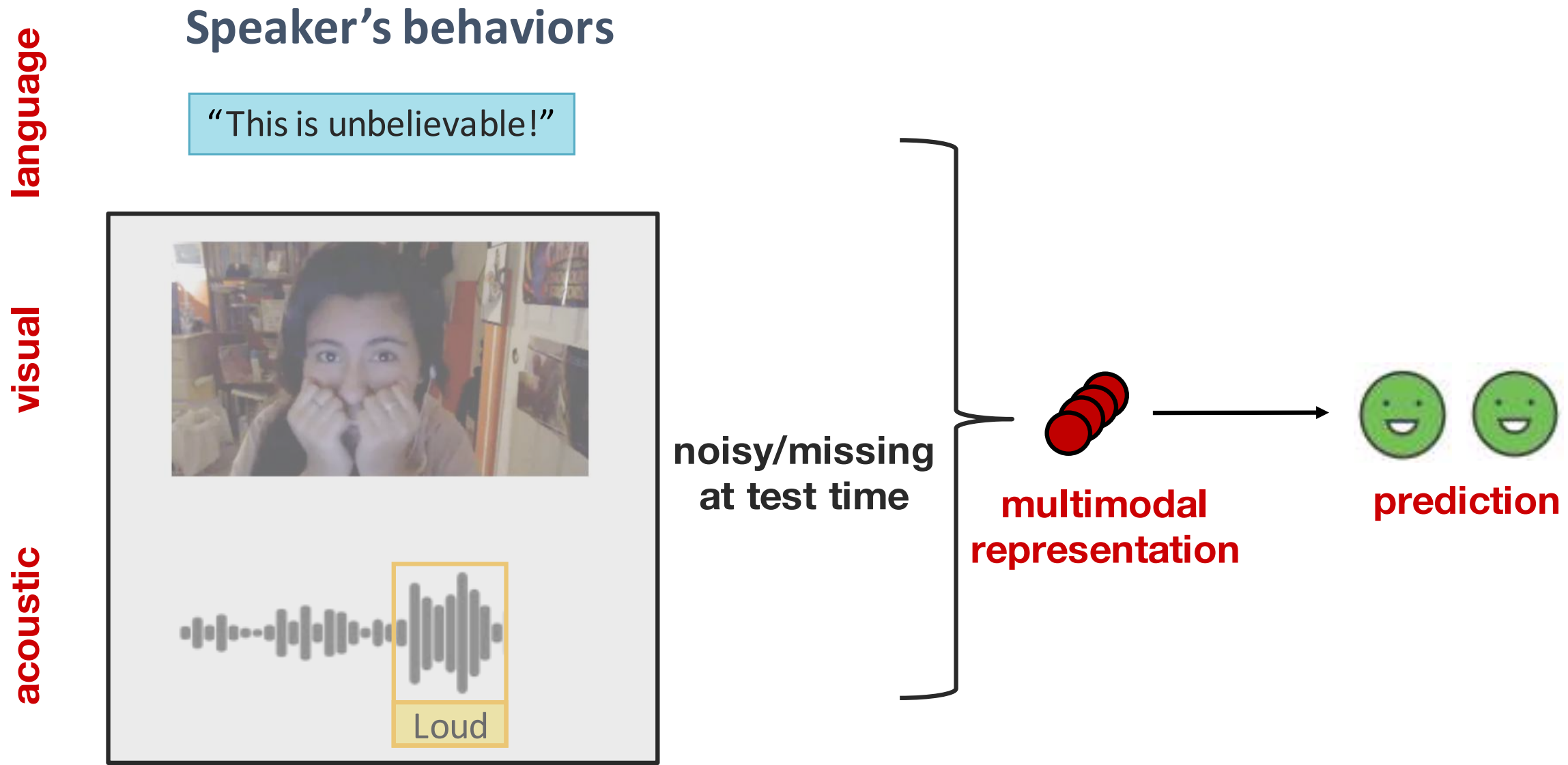
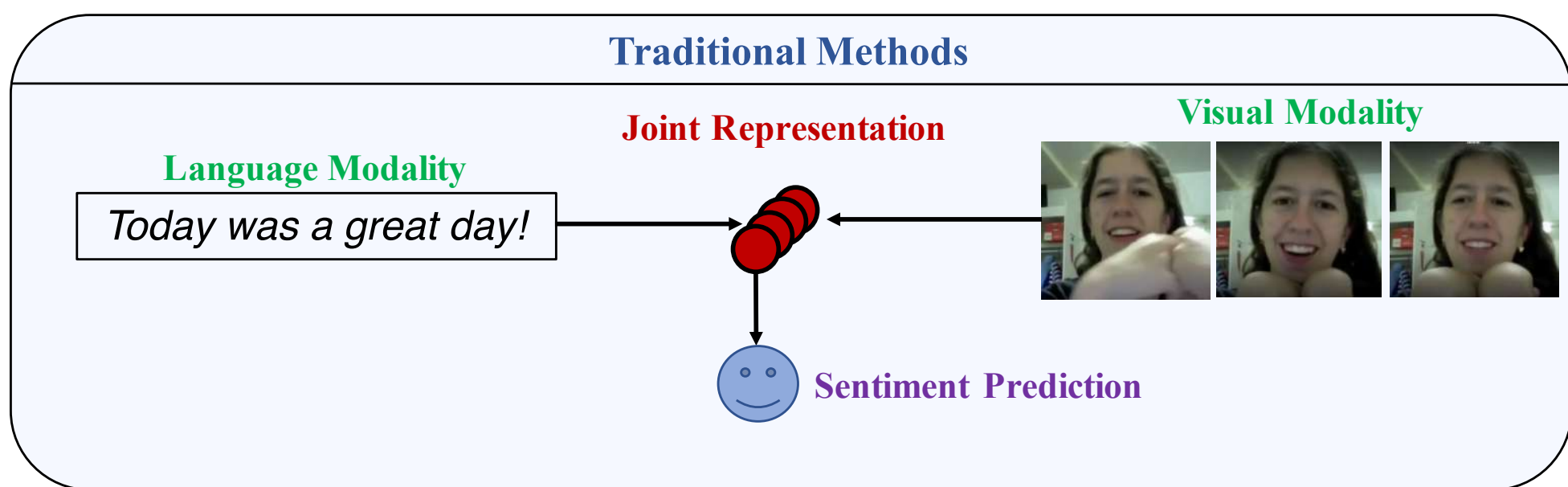


OVERVIEW

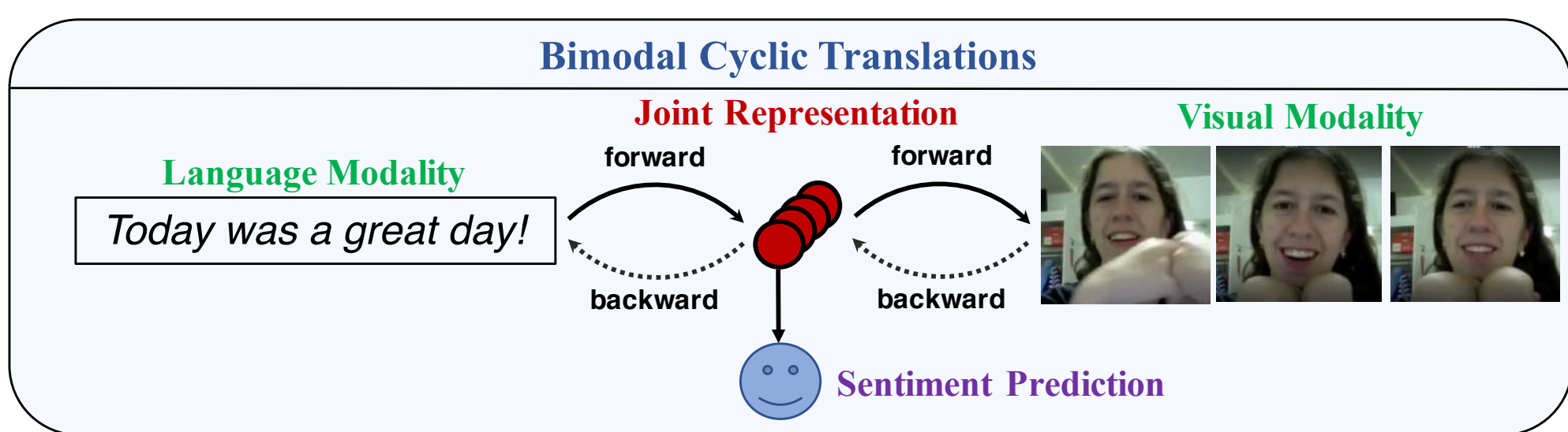


Traditional approaches



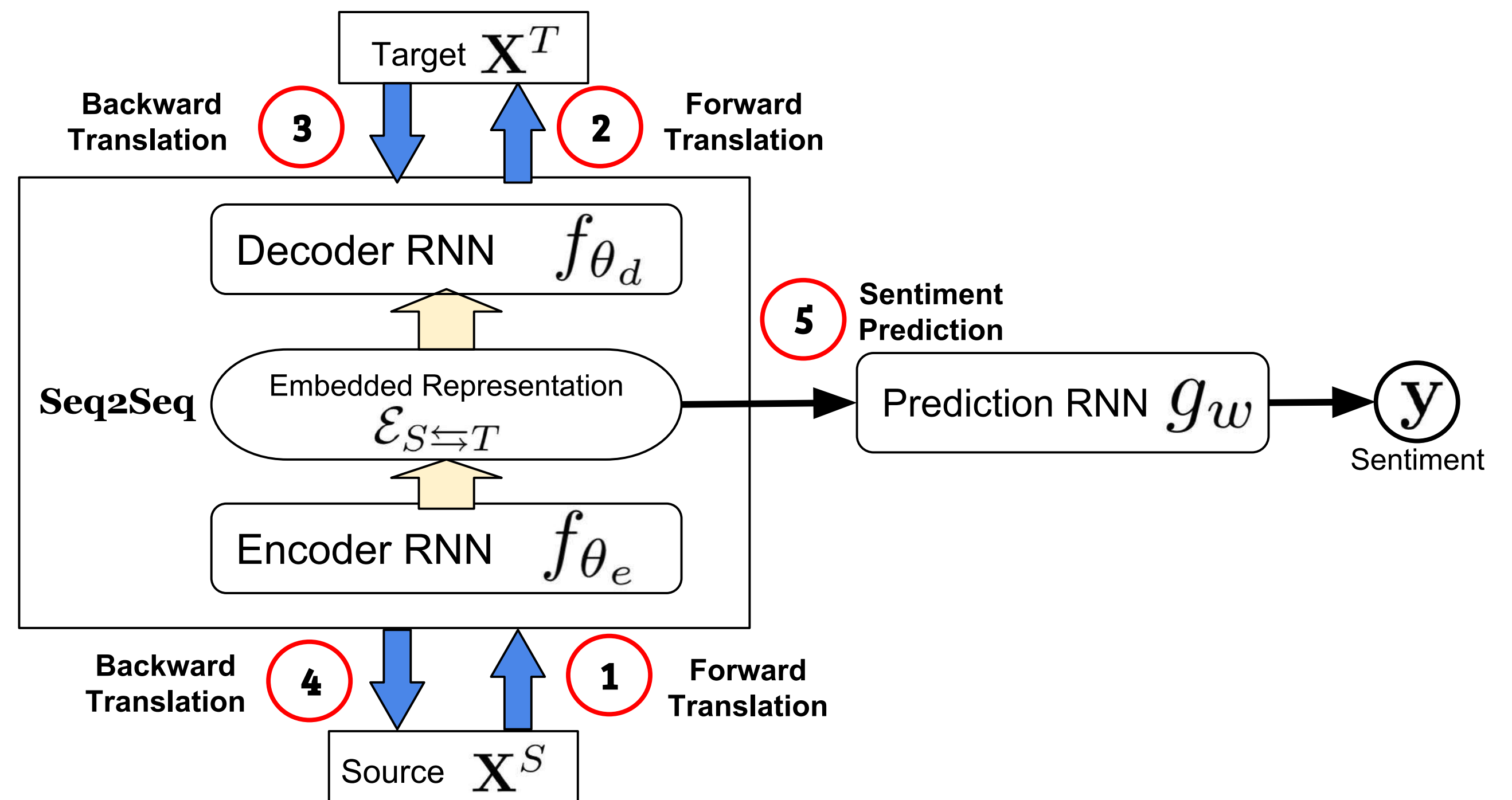
Both modalities required at test time!
Sensitive to missing/noisy visual modality.

Our approach: Found in Translation



Only language modality required at test time!

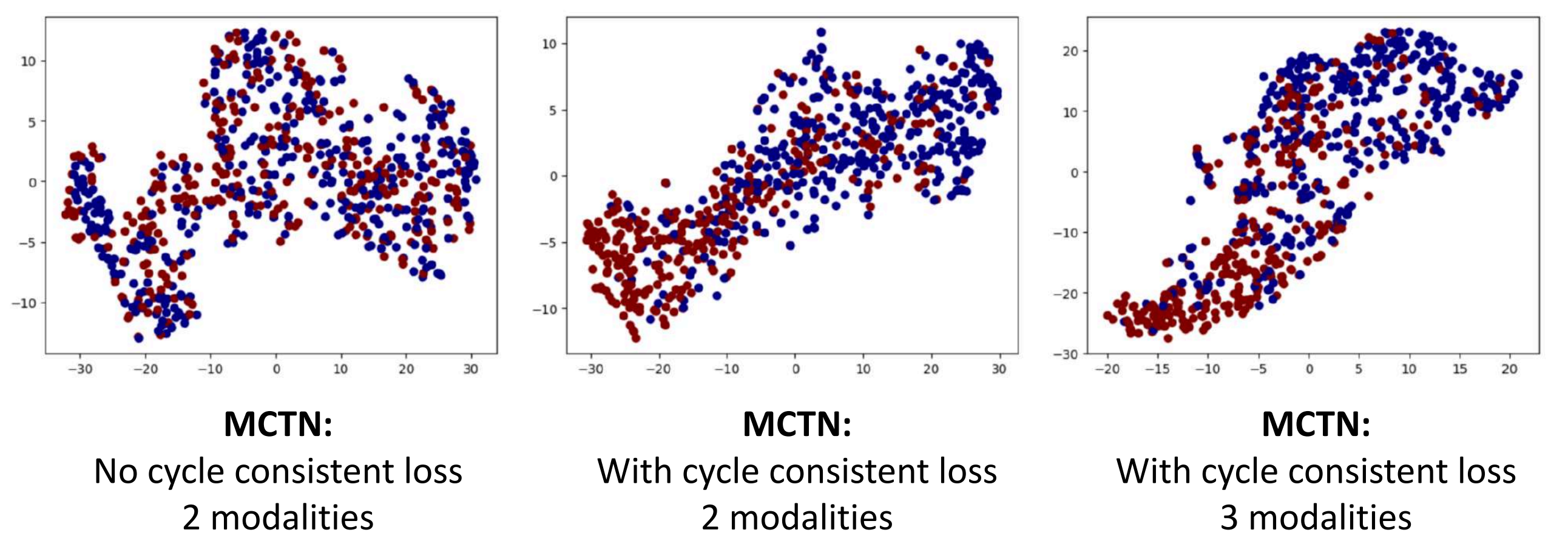
MULTIMODAL CYCLIC TRANSLATION NETWORK (MCTN)



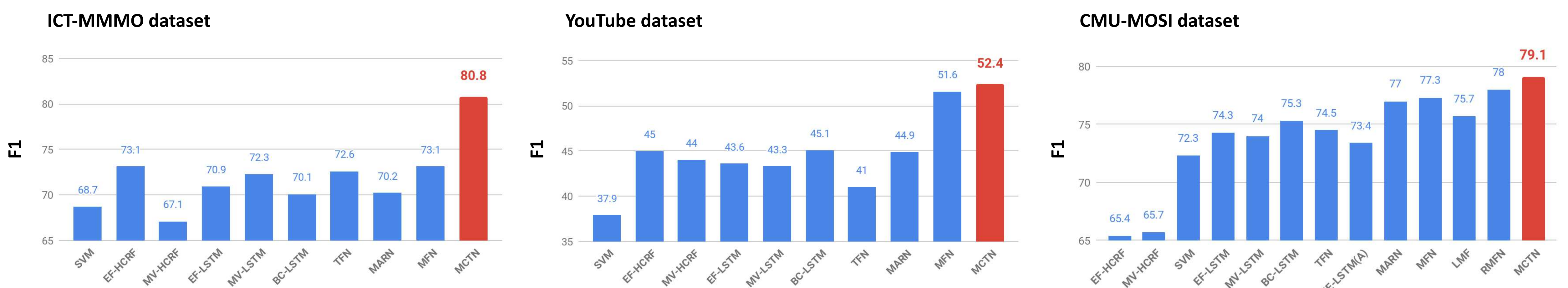
- Forward translation loss $\mathcal{L}_t = \mathbb{E}[\ell_{\mathbf{X}^T}(\hat{\mathbf{X}}^T, \mathbf{X}^T)]$
- Cycle consistent loss $\mathcal{L}_c = \mathbb{E}[\ell_{\mathbf{X}^S}(\hat{\mathbf{X}}^S, \mathbf{X}^S)]$
- Prediction loss $\mathcal{L}_p = \mathbb{E}[\ell_y(\hat{y}, y)]$

$$\mathcal{L} = \lambda_t \mathcal{L}_t + \lambda_c \mathcal{L}_c + \mathcal{L}_p$$

EMBEDDED REPRESENTATION WITH t-SNE

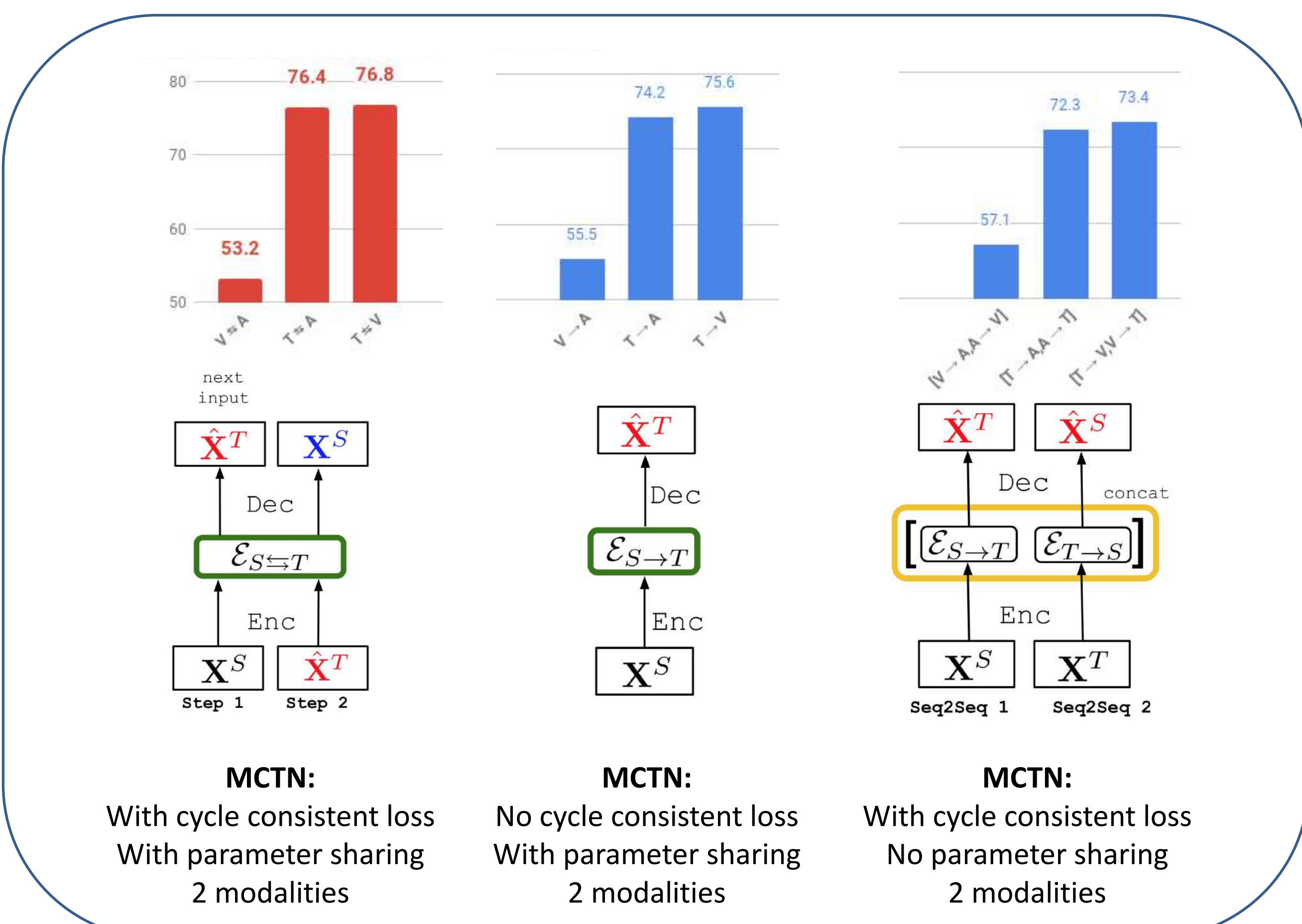


STATE-OF-THE-ART PREDICTION RESULTS



MCTN uses only language modality at test time!

ABLATION STUDY



- Use language as source modality
- Use cyclic translations
- Share parameters in seq2seq models

Code and Models:
<http://github.com/hainow/MCTN>