

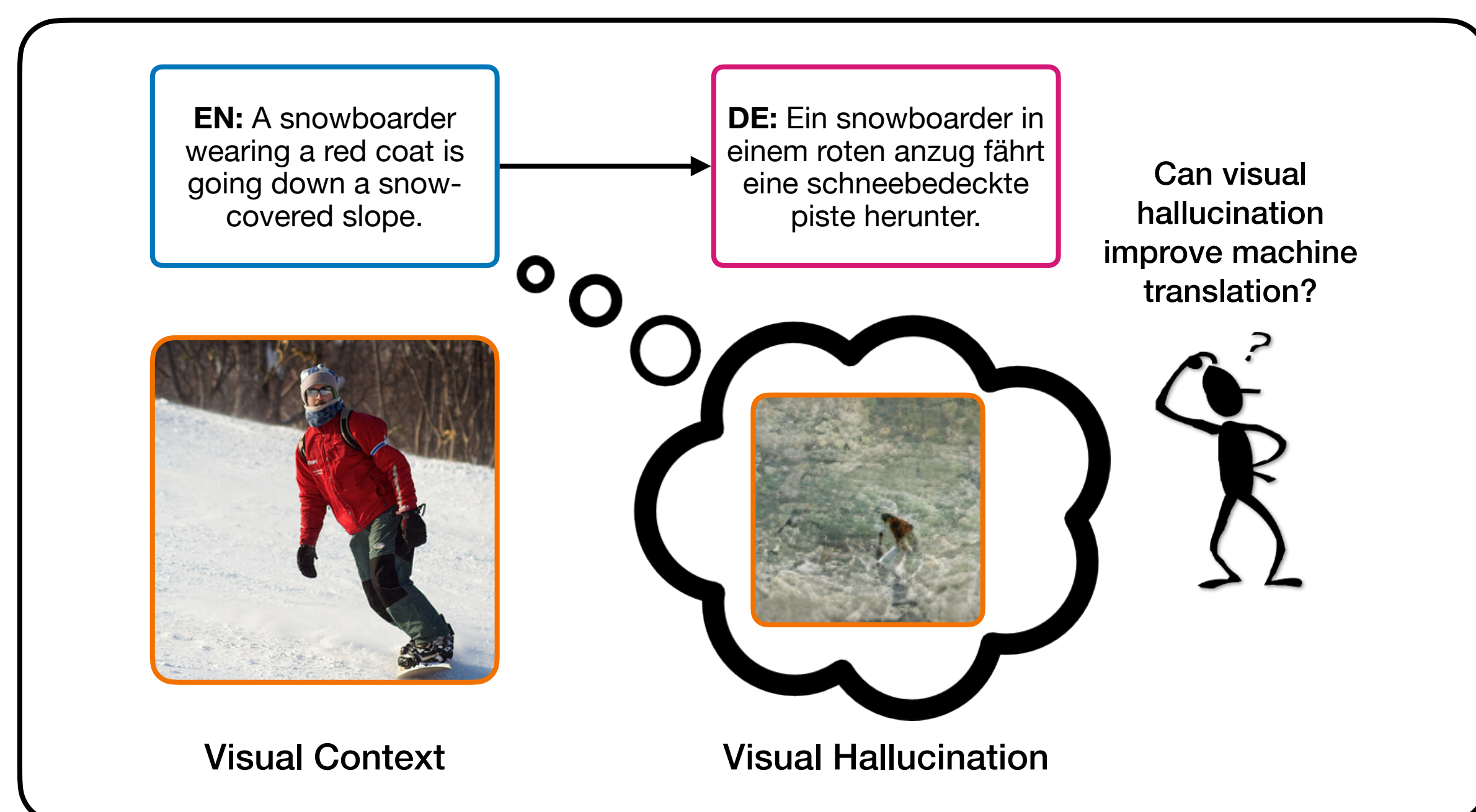
VALHALLA: Visual Hallucination for Machine Translation

Yi Li¹, Rameswar Panda², Yoon Kim³, Chun-Fu (Richard) Chen², Rogerio Feris², David Cox², Nuno Vasconcelos¹

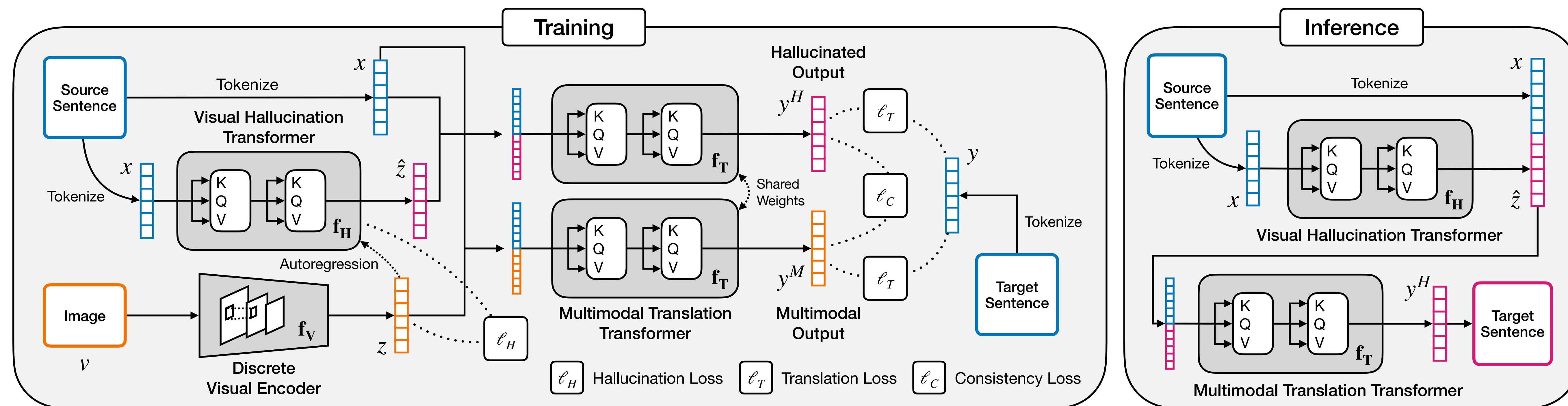


Overview

- Visual context has been exploited in **multimodal machine translation (MMT)** systems to improve translation performance
- In practice, inference with **paired image & text** not always realistic
- Can the ability to **hallucinate visual features** be used to improve language tasks, such as machine translation?
- We propose **VALHALLA**, a visual hallucination framework that learns from image-text correspondence during training, and predicts missing visual information for **text-only inference**



VALHALLA Architecture

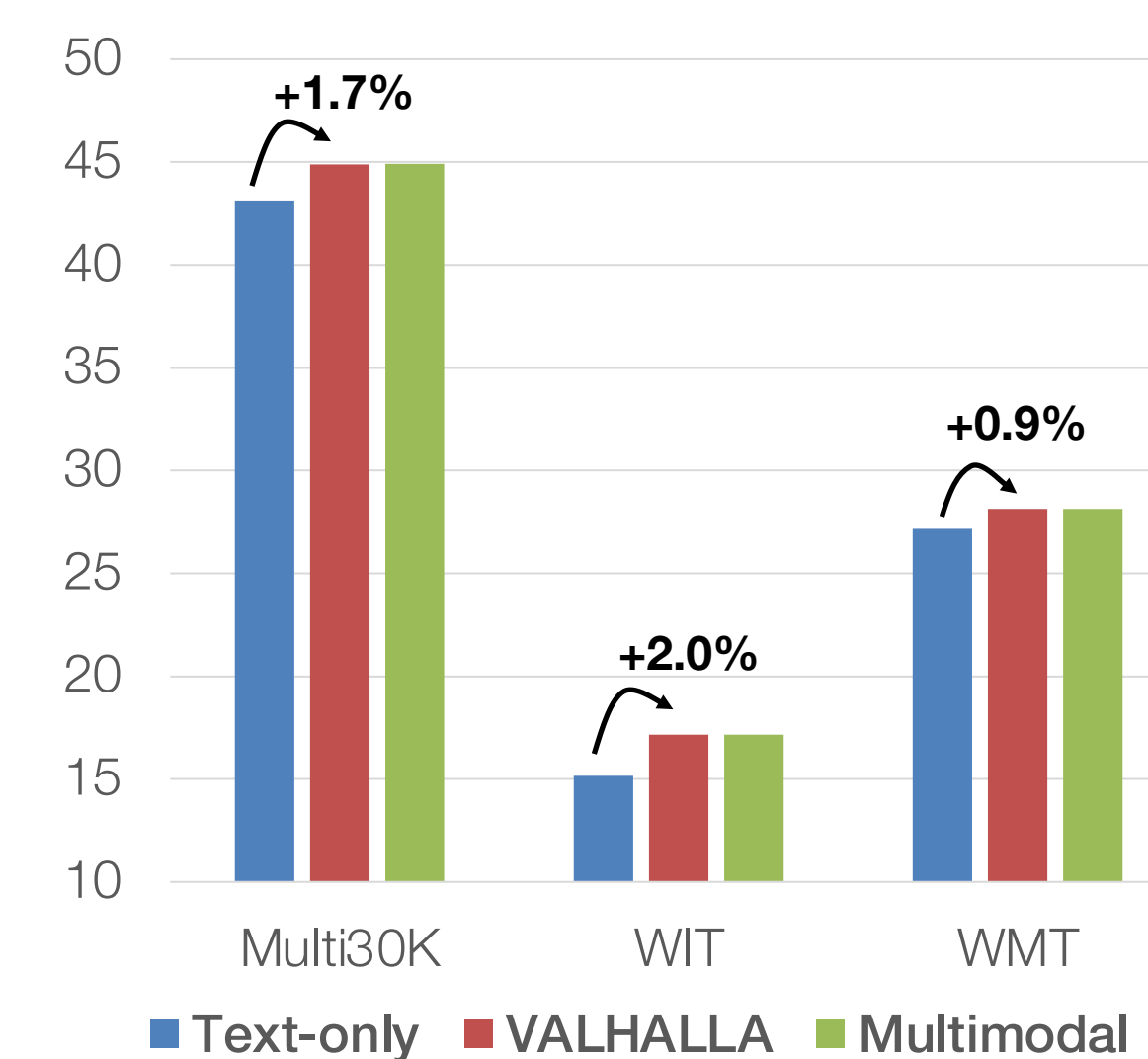


- Stage 1 **Discrete visual encoder f_V** tokenizes input image v into discrete sequence z
- Stage 2 **Visual hallucination f_H** predicts visual tokens z from source sentence x
- Stage 3 **Multimodal translation f_T** predicts target sentence y from concatenated input (x, z)
- Text-only translation** by replacing z with hallucinated sequence $\hat{z} = f_H(x)$

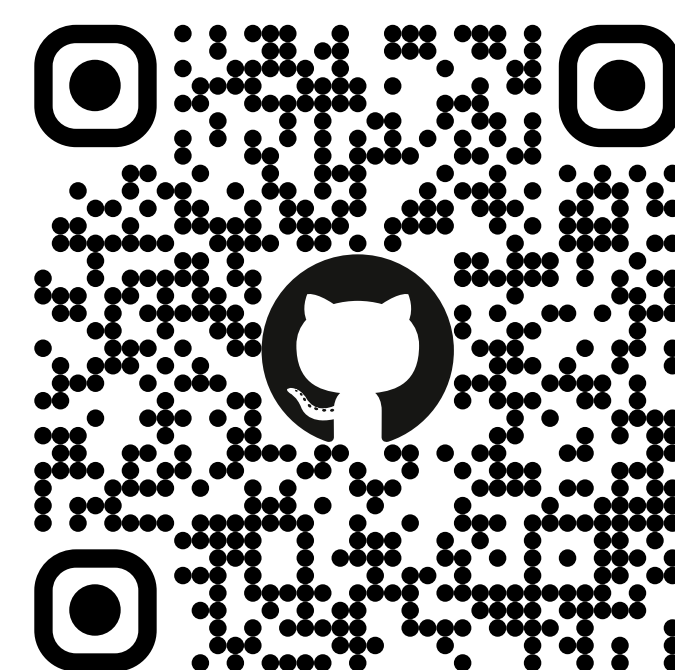
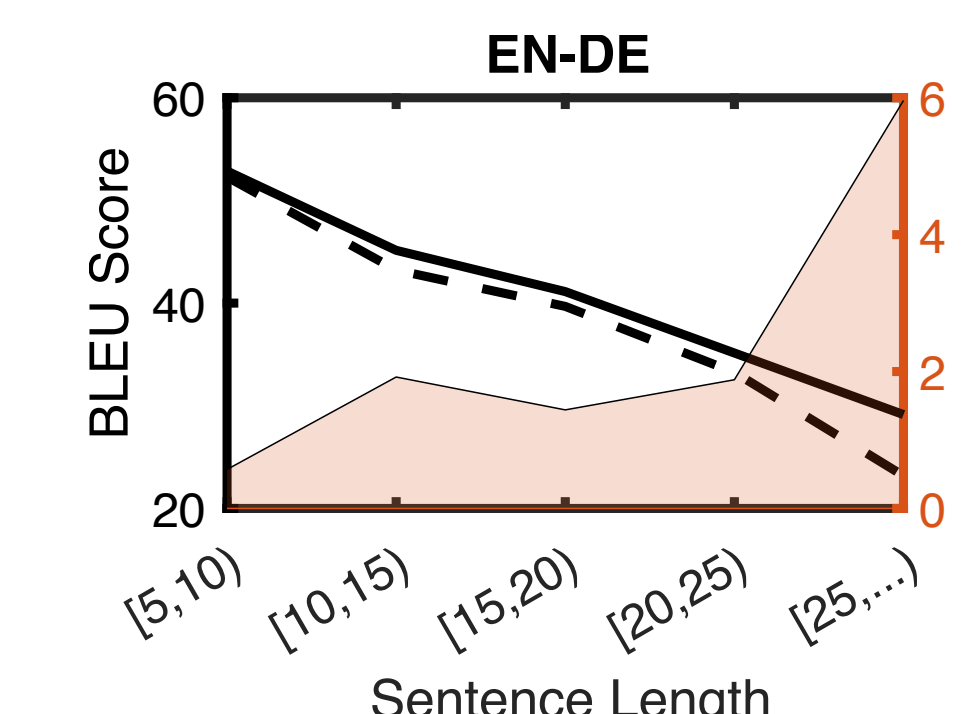
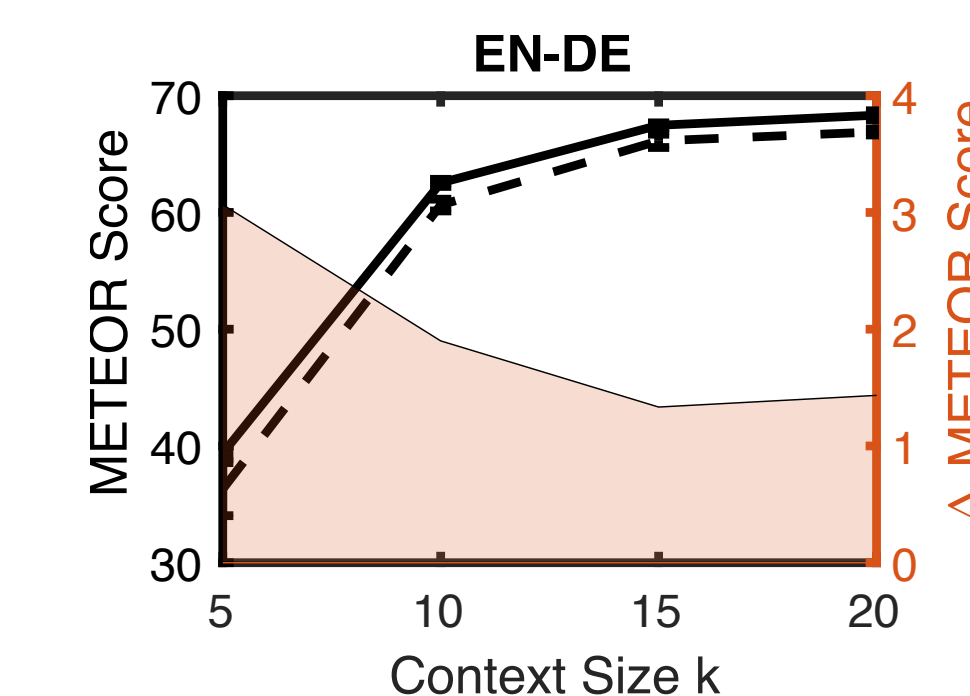
Experiments

- On multimodal datasets (**Multi30K**, **WIT**), VALHALLA improves performance over text-only transformer, matches MMT scores without requiring image input
- Training on text-only corpora (**WMT**) enabled through CLIP-based image retrieval
- Under **limited textual context**, VALHALLA recovers masked words, visual entities in source sentence
- Greater improvement for translating long sentences

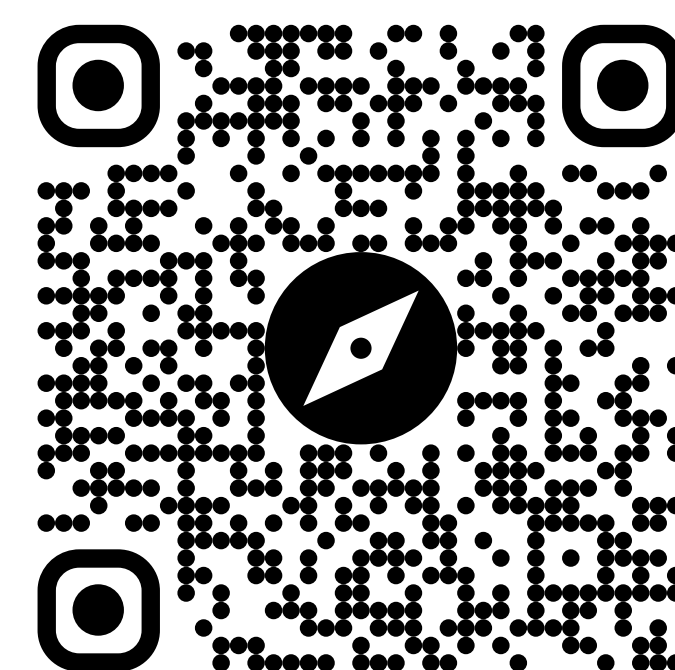
Translation BLEU4



Source EN: A boy wearing a red shirt digs into the sand with a yellow shovel.
 Reference DE: Ein junge in einem roten shirt gräbt mit einer gelben schaufel im sand.
 Text-Only: Ein junge, der ein rotes hemd trägt, wirft einen gelben ball in den sand. (A boy wearing a red shirt throws a yellow ball in the sand.)
 VALHALLA: Ein junge in einem roten hemd gräbt mit einer gelben schaufel in den sand. (A boy in a red shirt is digging in the sand with a yellow shovel.)



Code & Models



Project page