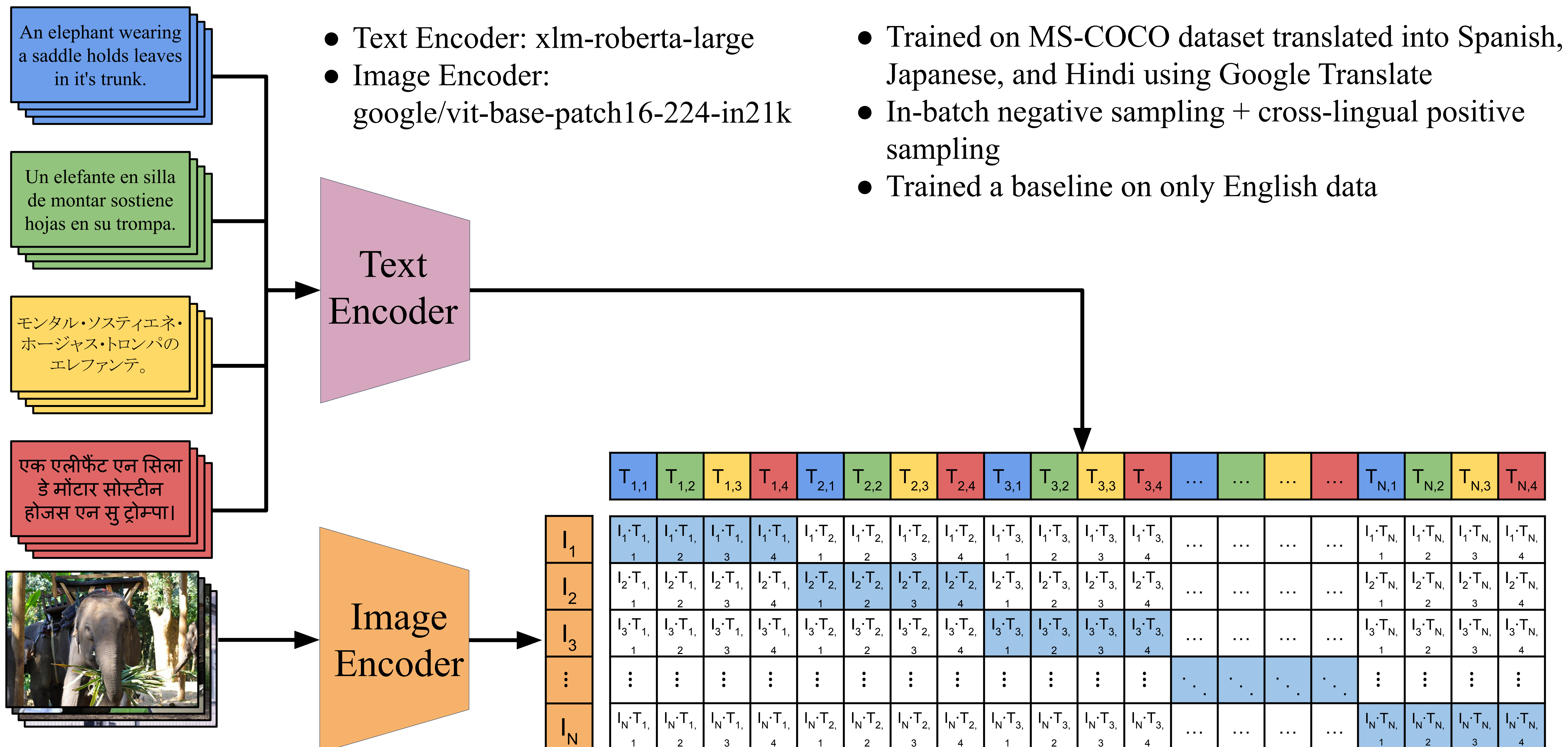
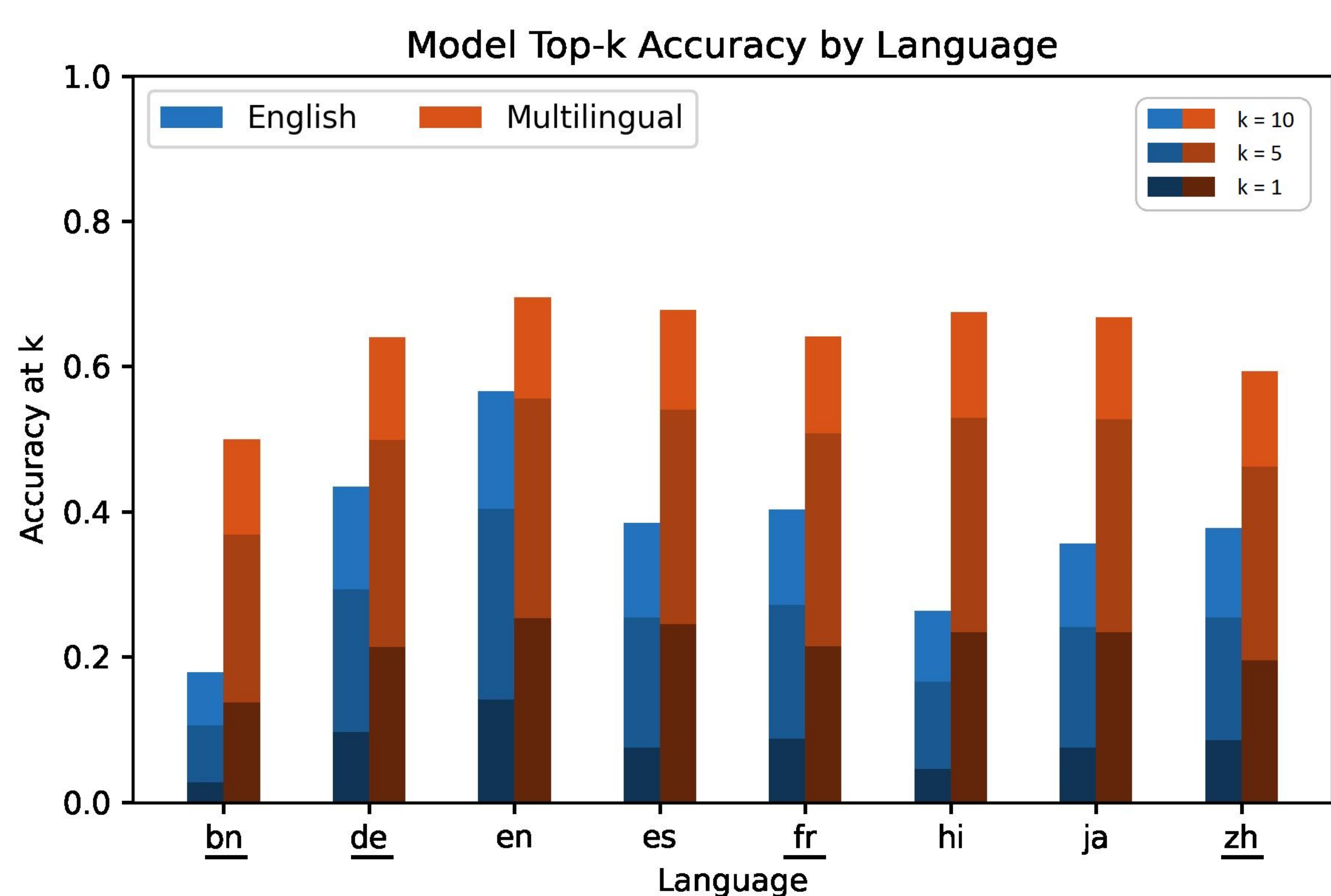


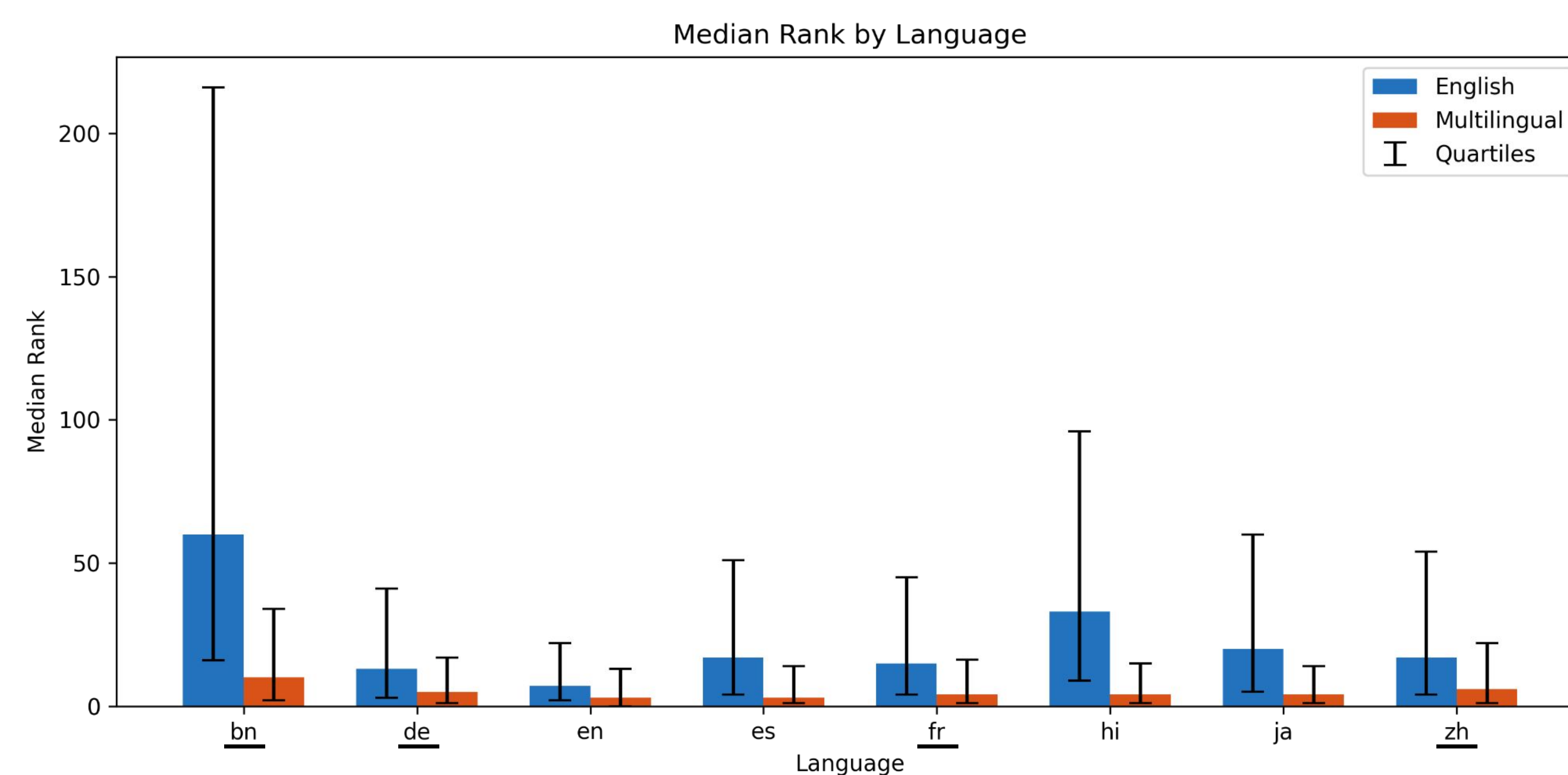
Parallel Contrastive Pre-Training



Cross-Modal Retrieval



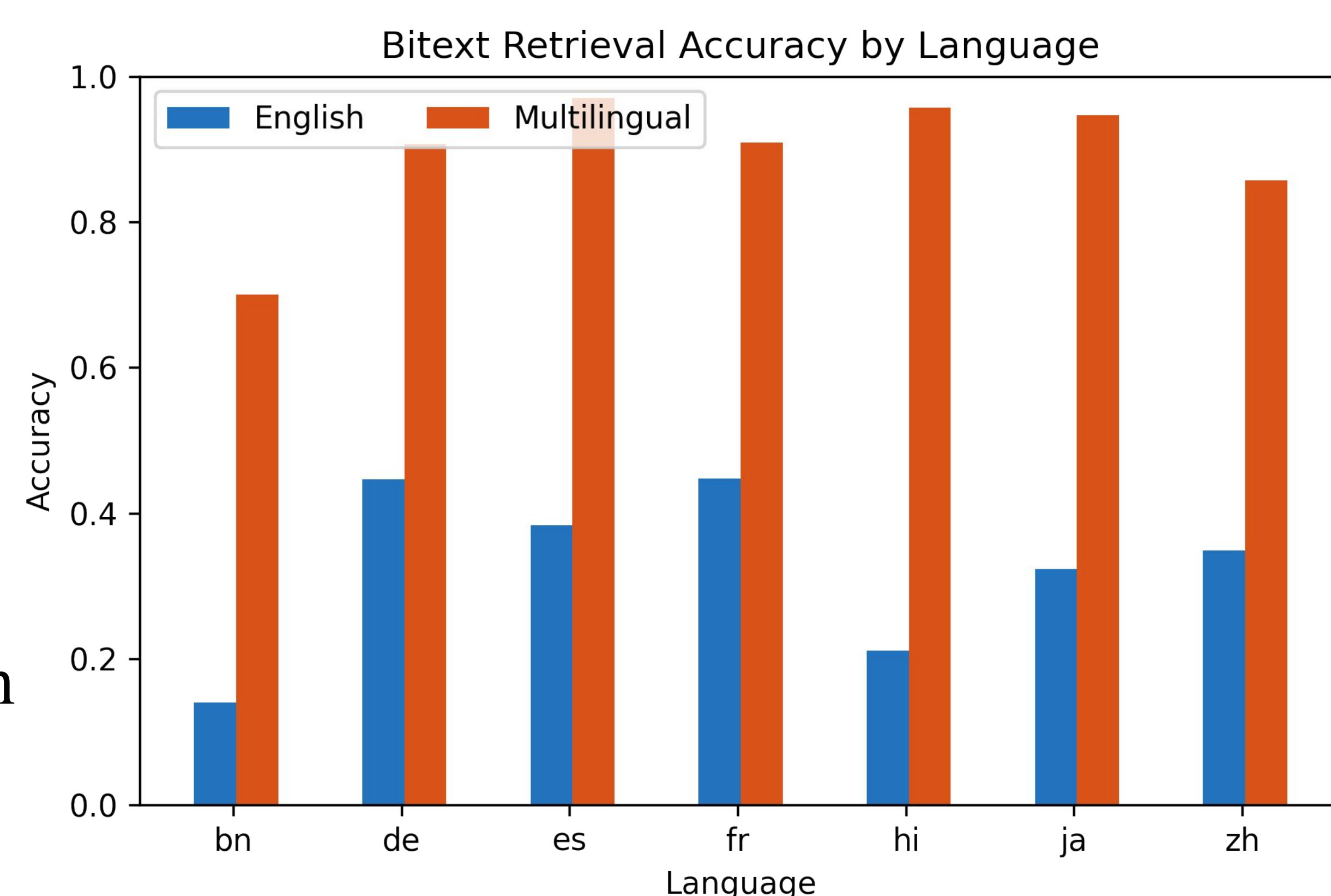
- Accuracy for retrieving nearest neighbor image based on encoding distance to sentence
- Higher is better
- Multilingual training even improves English performance
- Improves zero-shot performance greatly (underlined)



- Median and quartiles for where correct image ranks by distance from sentence encoding
- Lower is better
- Multilingual training improves consistency as well as overall performance

Bitext Retrieval

- Tested bitext retrieval from English encodings to others
- Multilingual training improves alignment between languages
- Over the COCO validation set, so from a pool of 5000 candidate sentences



Code

<https://github.com/nkrasner/ML-CLIP>

