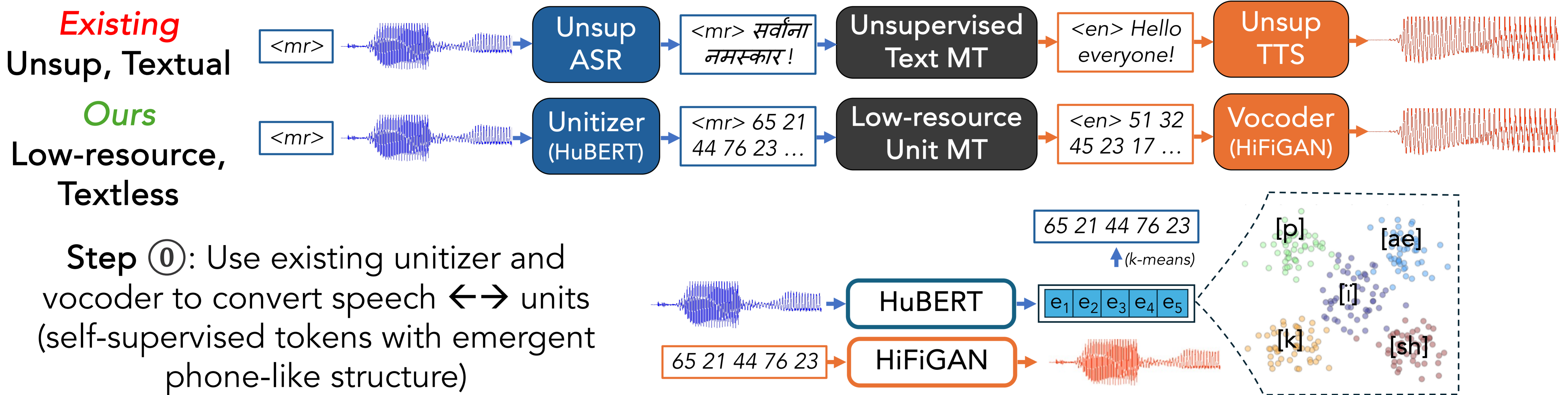




Current S2ST models either rely on text as an intermediary or require extensive parallel speech data, limiting support for textless and low-resource languages. How can we bridge the gap?

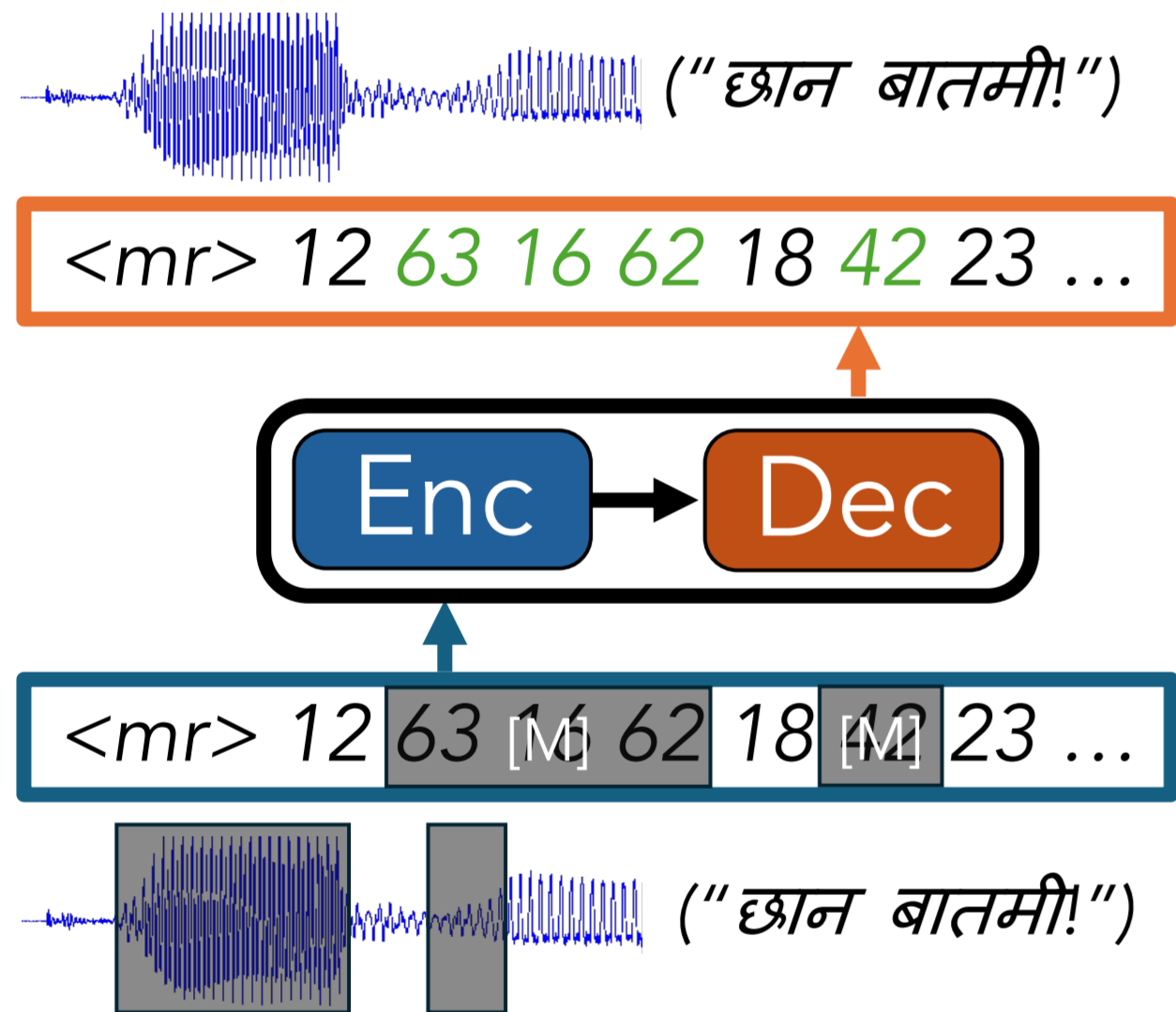
We adapt an unsupervised text-based S2ST approach to the low-resource, textless S2ST setting by using **units** instead of text



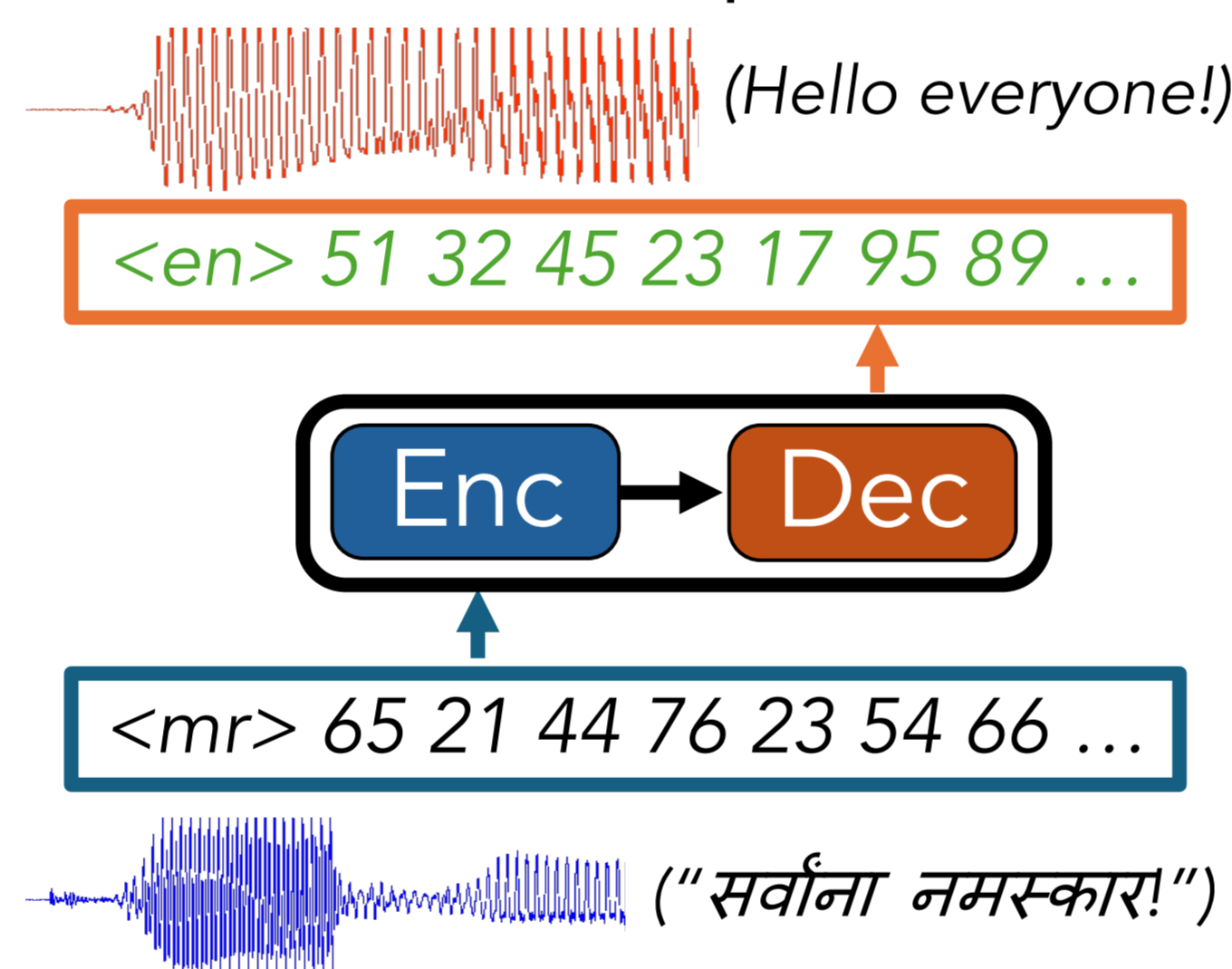
**Step ①:** Use existing unitizer and vocoder to convert speech  $\leftrightarrow$  units (self-supervised tokens with emergent phone-like structure)

**Steps ①-③:** Train a low-resource unit **Enc**  $\rightarrow$  **Dec** MT model (like mBART)

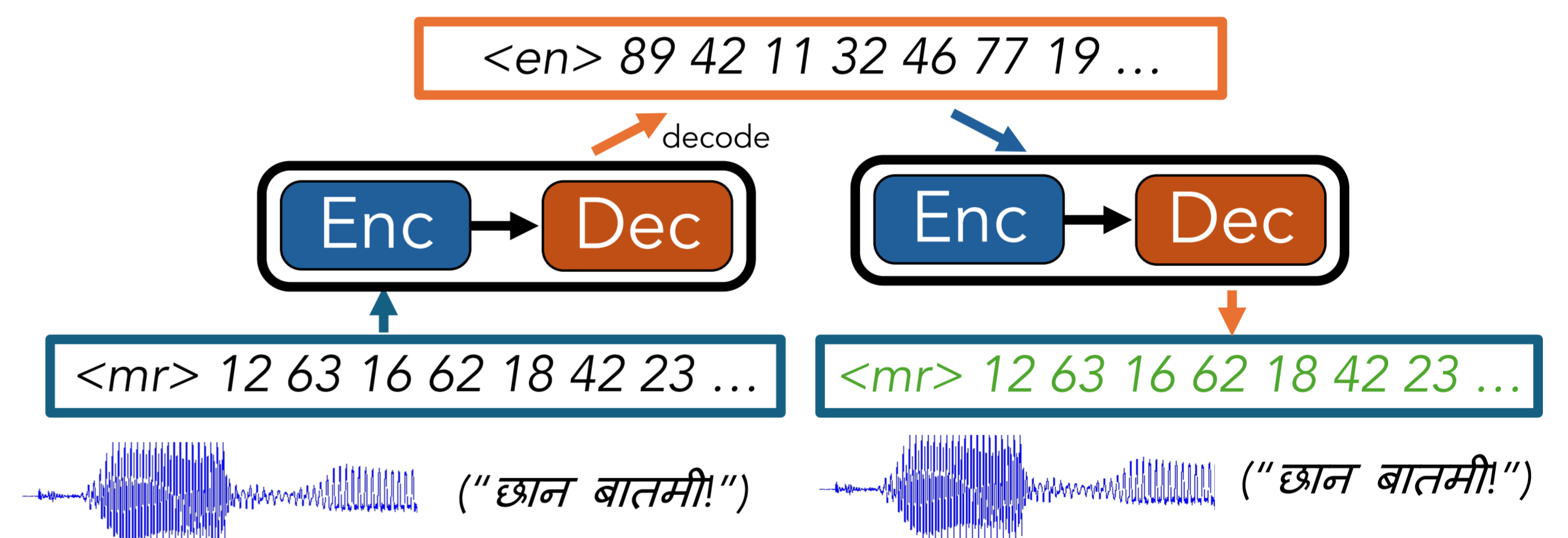
**Step ①:** Pretrain LM on monolingual per-lang data (masked denoising loss)



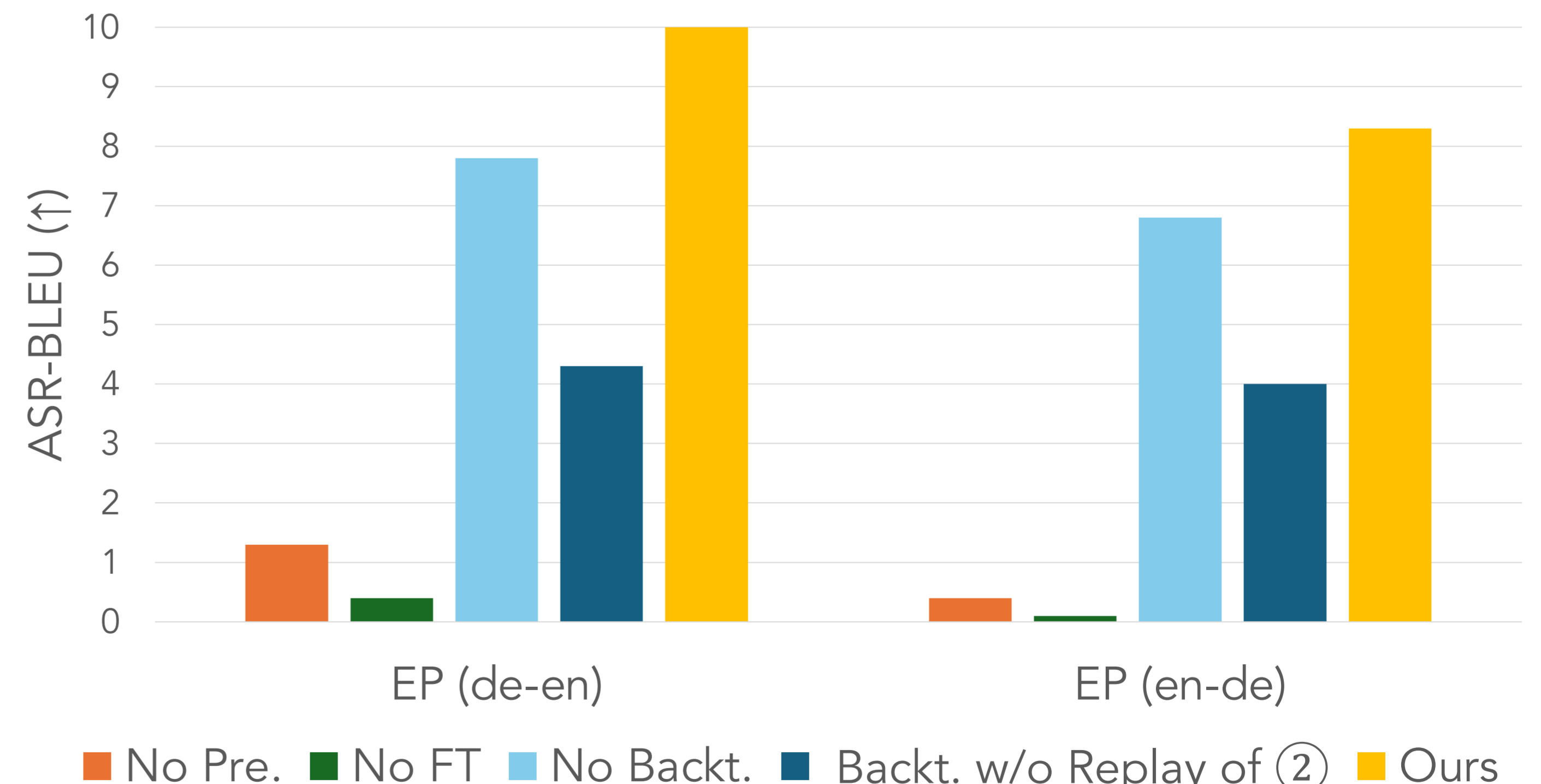
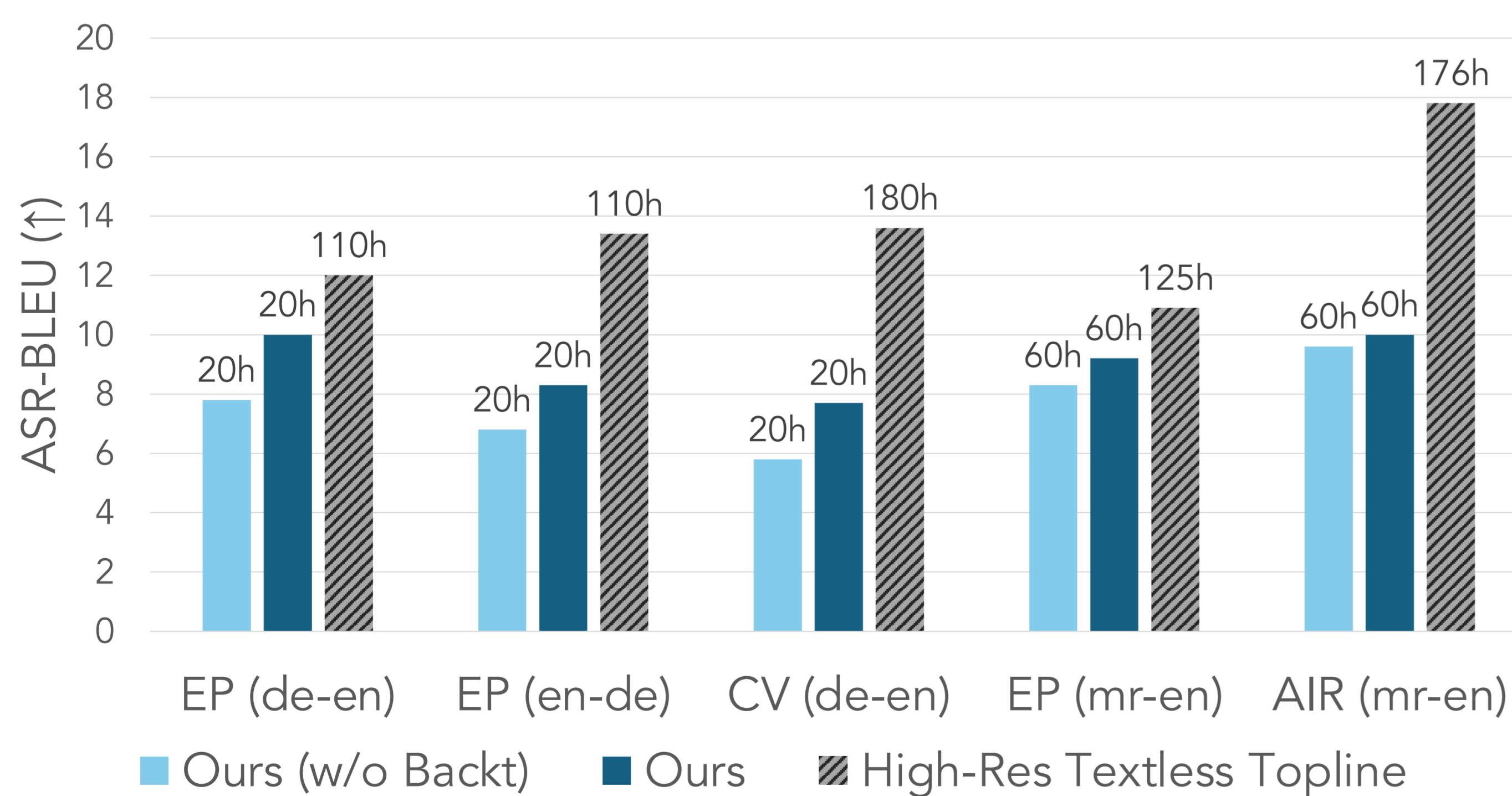
**Step ②:** Finetune LM on limited S2S translation data (cross-entropy loss)



**Step ③:** Backtranslate LM on monolingual per-lang data (while replaying ②) (cross-entropy loss)



**Experiments:** Languages: en-de, en-mr. Domains: EP (Europarl), CV (Common Voice), AIR (All Ind Radio)



### Main Takeaways

- In some settings, our method is within 1-2 ASR-BLEU points of a high-res textless topline
- Diff. domains exhibit diff. performance gaps

### Ablations

- In order of importance, FT > Pre. > Backt.
- Replay is necessary for backtranslation to work

### More analyses in the paper!

- How to select an appropriate unitizer?
- How do textless models compare to text-based models?
- How does performance differ for short vs. long utterances?

### Future Work

- Scaling to stronger pretrained multilingual unit LMs, with potential for zero-resource textless S2ST
- Using more semantic unitizers for efficient unit LM training