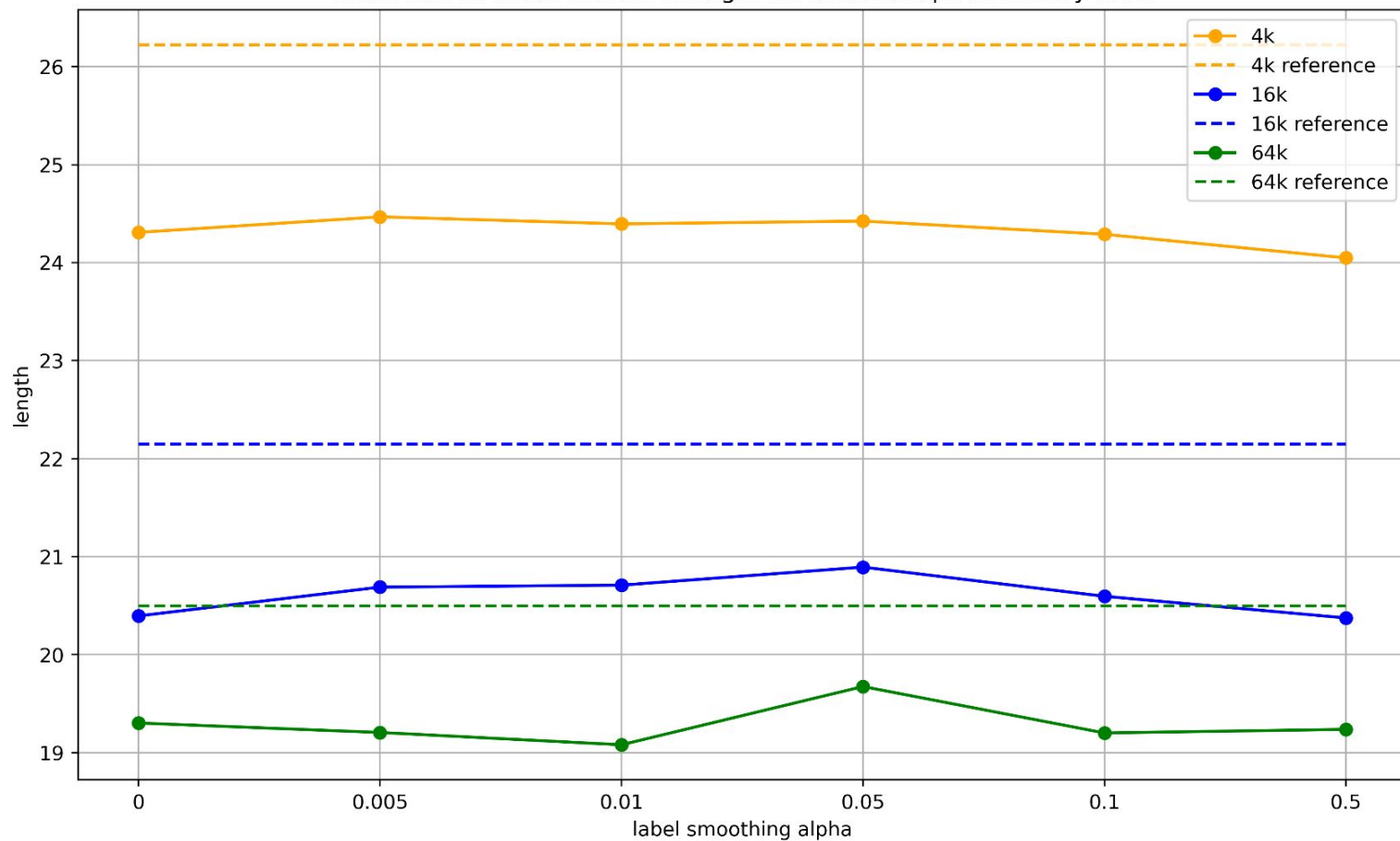


Transformer experiments

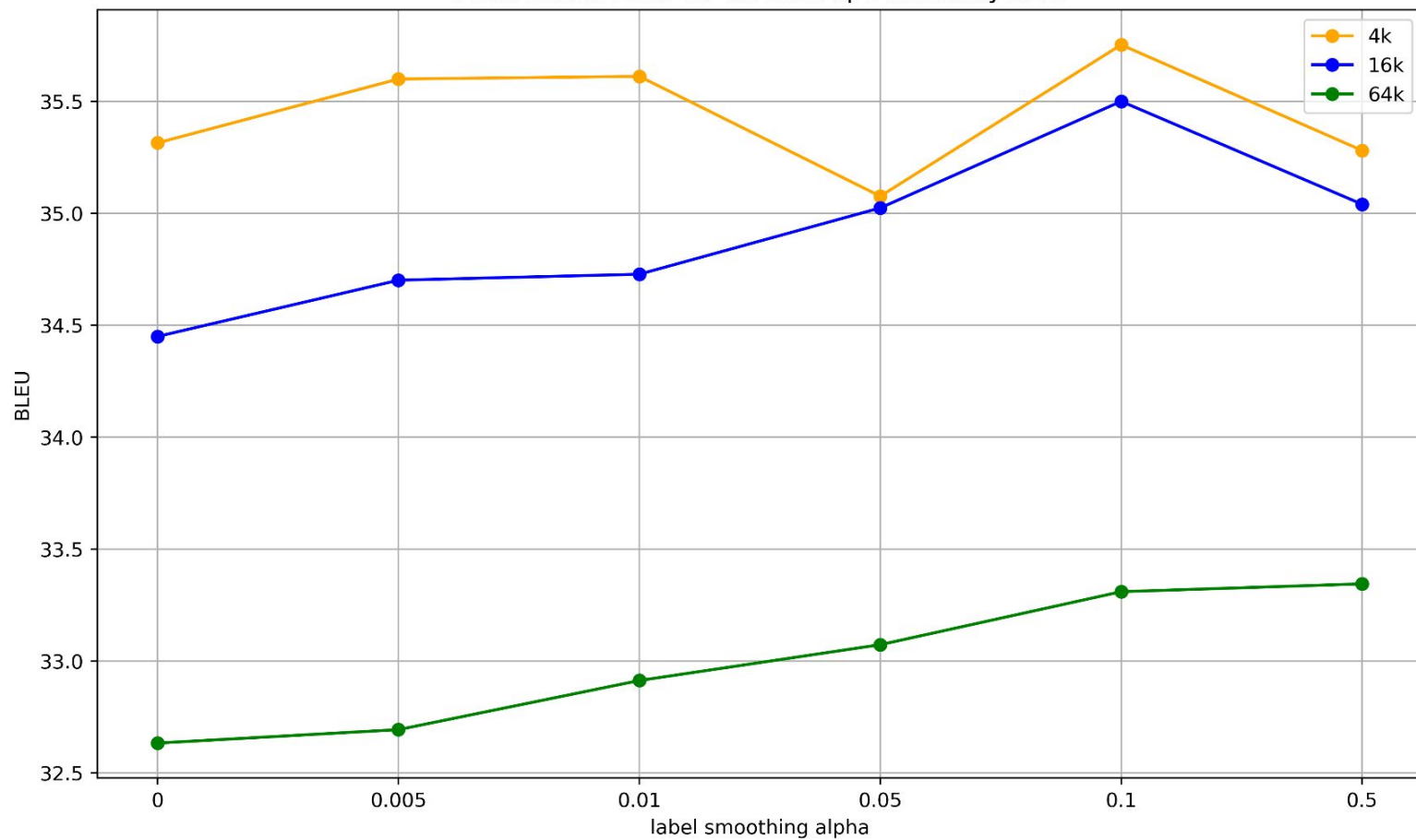
Dict Size Experiments - Beam Search

- Transformer setup as in Riley and Chiang
 - 6 layers, 4 heads
 - fairseq transformer_iwslt_de_en setup
 - sinusoid positional embeddings
 - context size of 4096 tokens
 - Trained on iwslt17 de-en data
- Experimental setup
 - Evaluation with beam size of 4
 - We compare the mean sentence length of the beam search outputs to the mean sentence length of the reference translations
 - sentence length = # bpe tokens
 - evaluated for label smoothing alpha values [0, 0.005, 0.01, 0.05, 0.1, 0.5] and bpe dict sizes [4k, 16k, 64k]
 - BLEU values reported for reference

Beam search mean sentence length for different bpe dictionary sizes



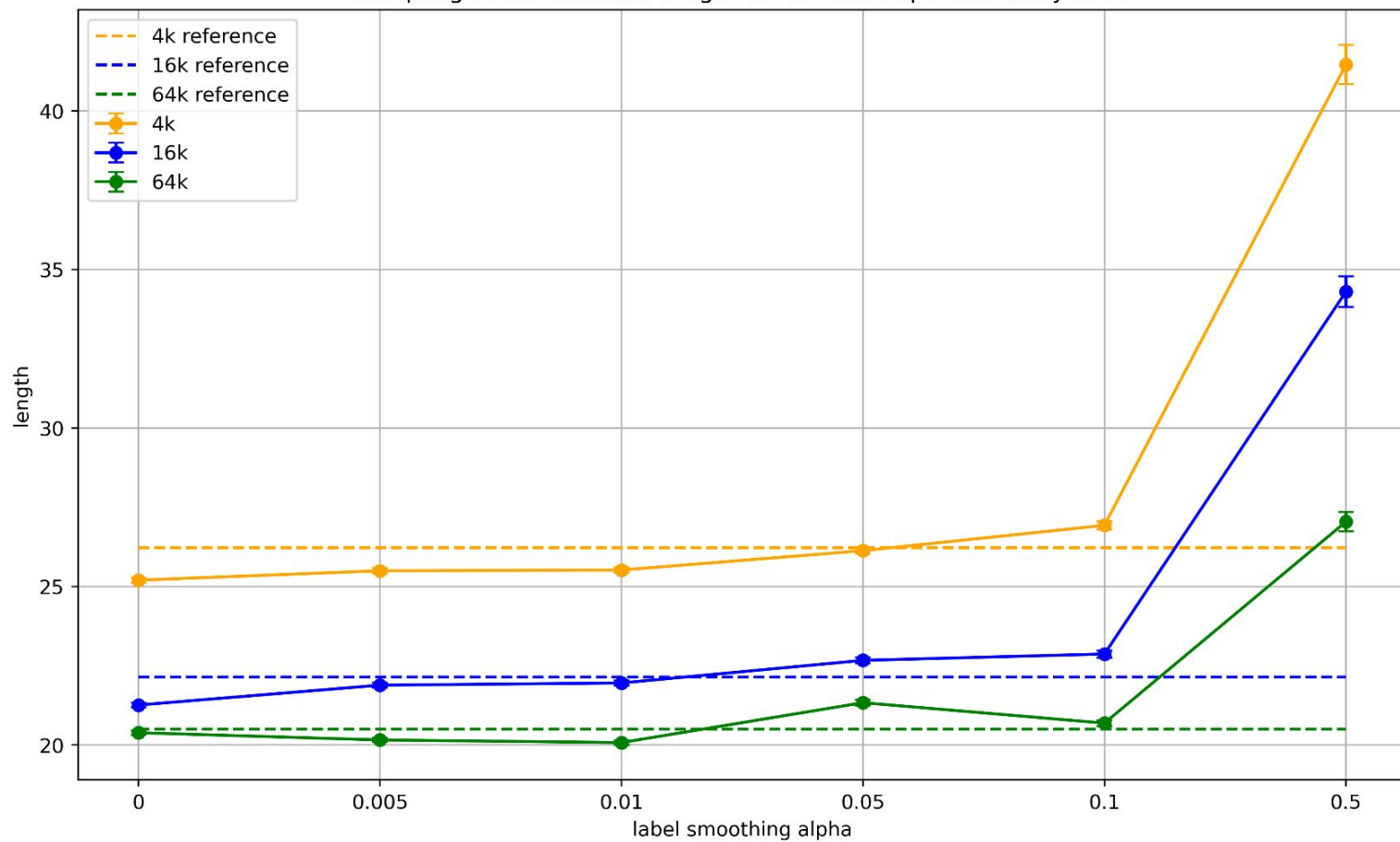
Beam search BLEU for different bpe dictionary sizes



Dict Size Experiments - Sampling

- Sample 1000 hypotheses per source sentence
- We compute the mean hypothesis length per sentence and then compute the mean of means for each model
- We compare it again to the mean length of the reference corpus

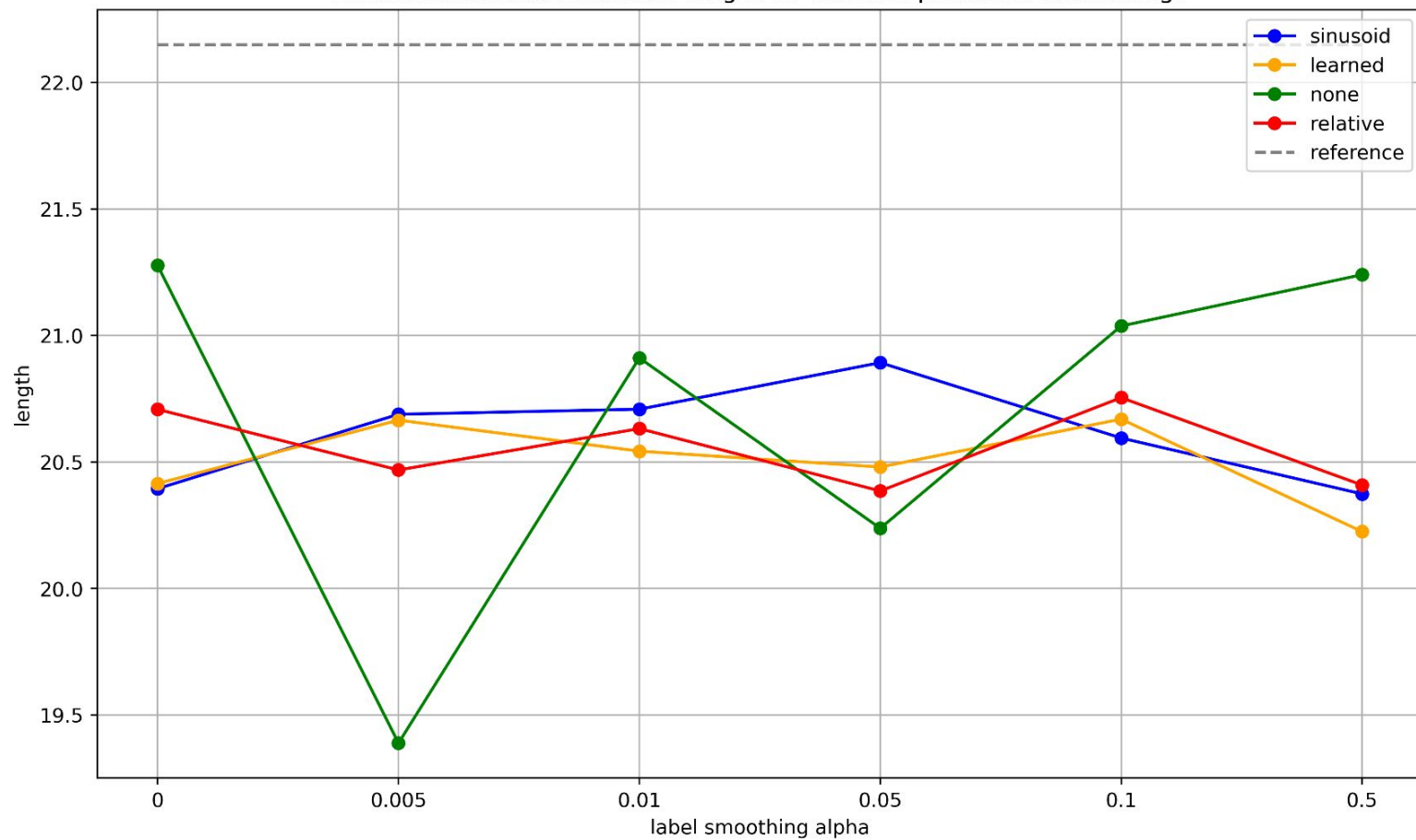
Sampling mean sentence length for different bpe dictionary sizes



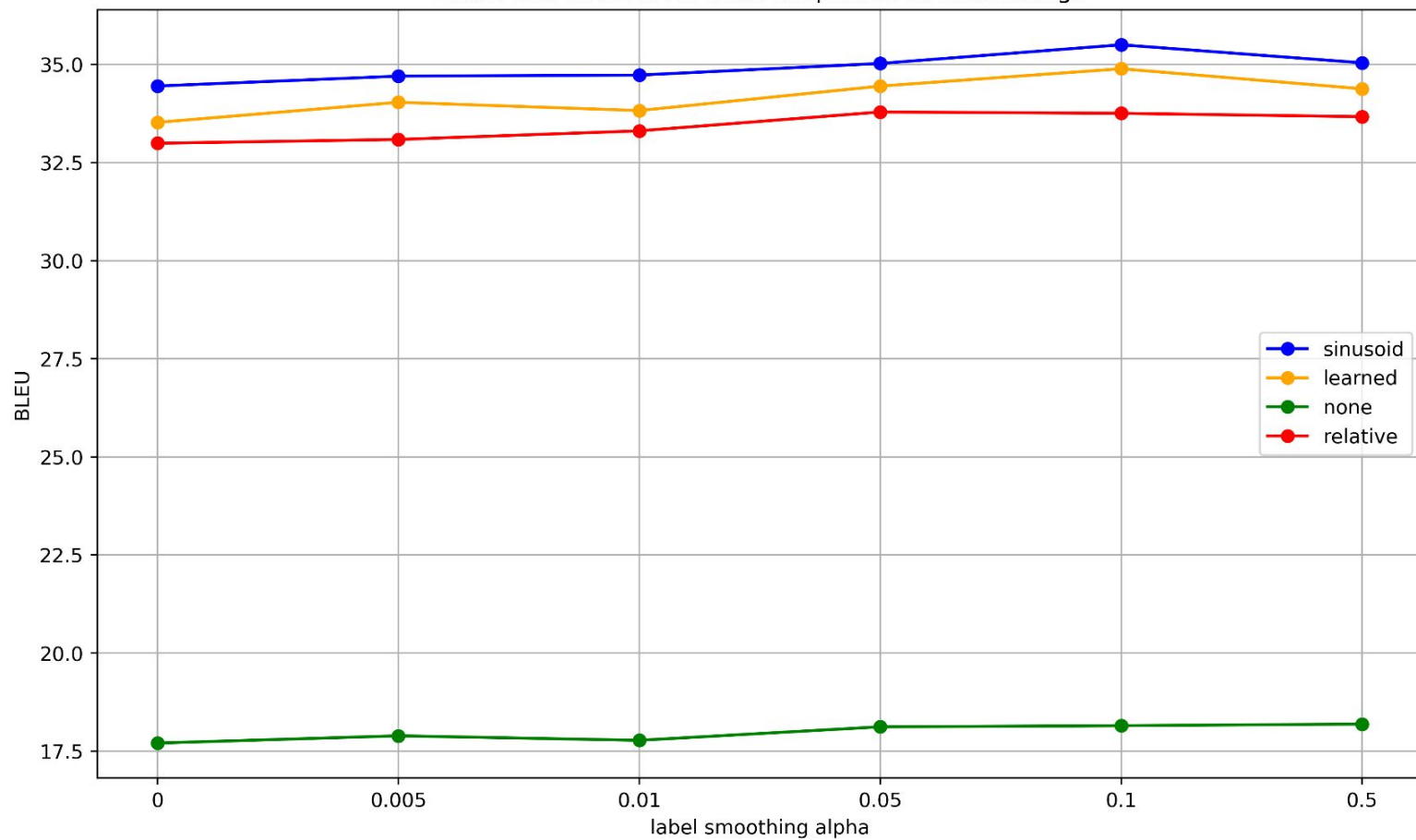
Positional Embedding Experiments

- Comparison between sinusoidal, learned, relative and no positional embeddings
 - 16k bpe dict size for all systems
- Additional comparison of sampling results to mean training set reference length (before only test set reference length)
 - We plot standard error and standard deviation for the sampling results
- Relative positional embeddings based on [Attention with Linear Biases](#)
 - Based on Florian Schottmann implementation
 - Relative positional attention in encoder and decoder self attention, no positional embeddings in cross-attention
 - Setup variation due to compatibility: 8 attention heads, fairseq version 0.10.2

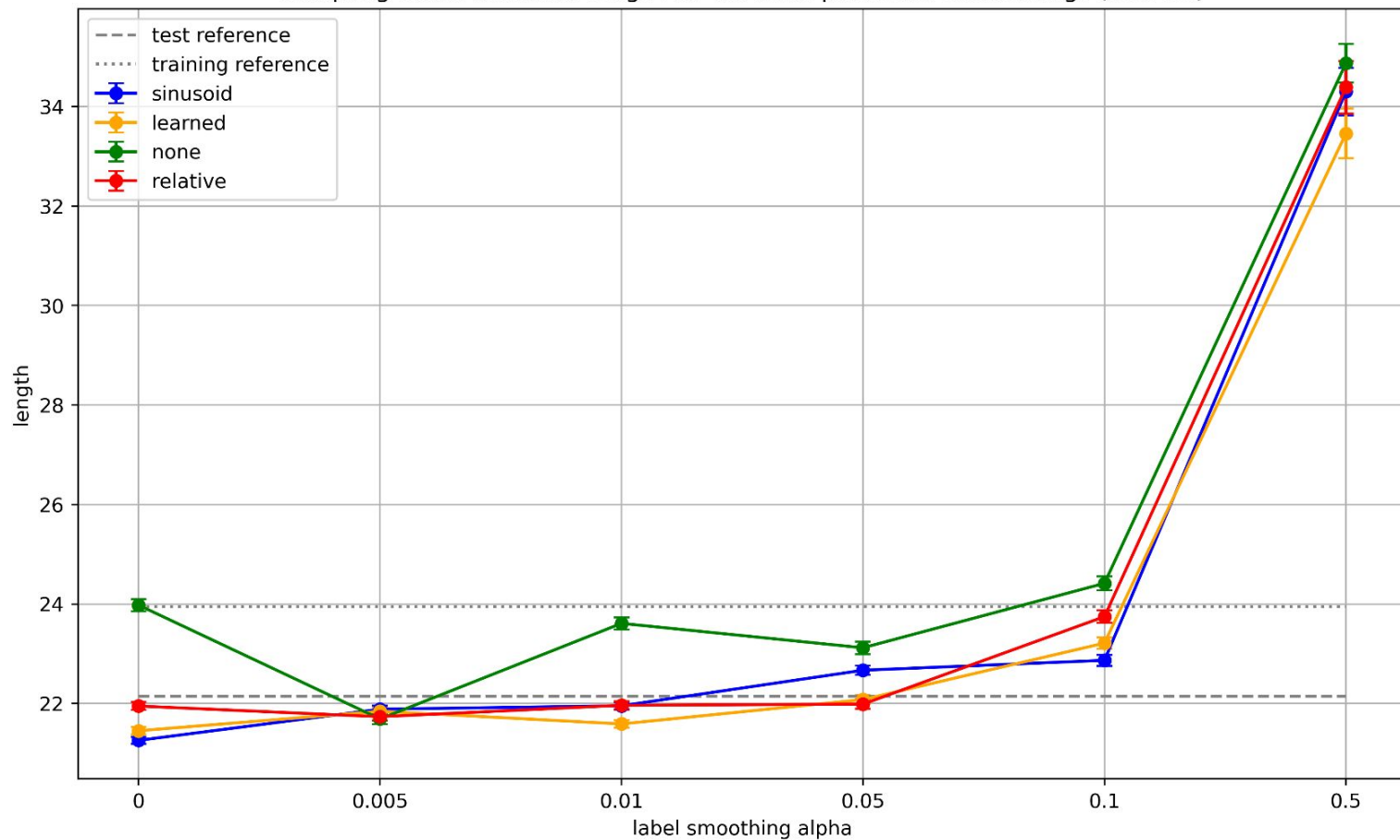
Beam search mean sentence length for different positional embeddings



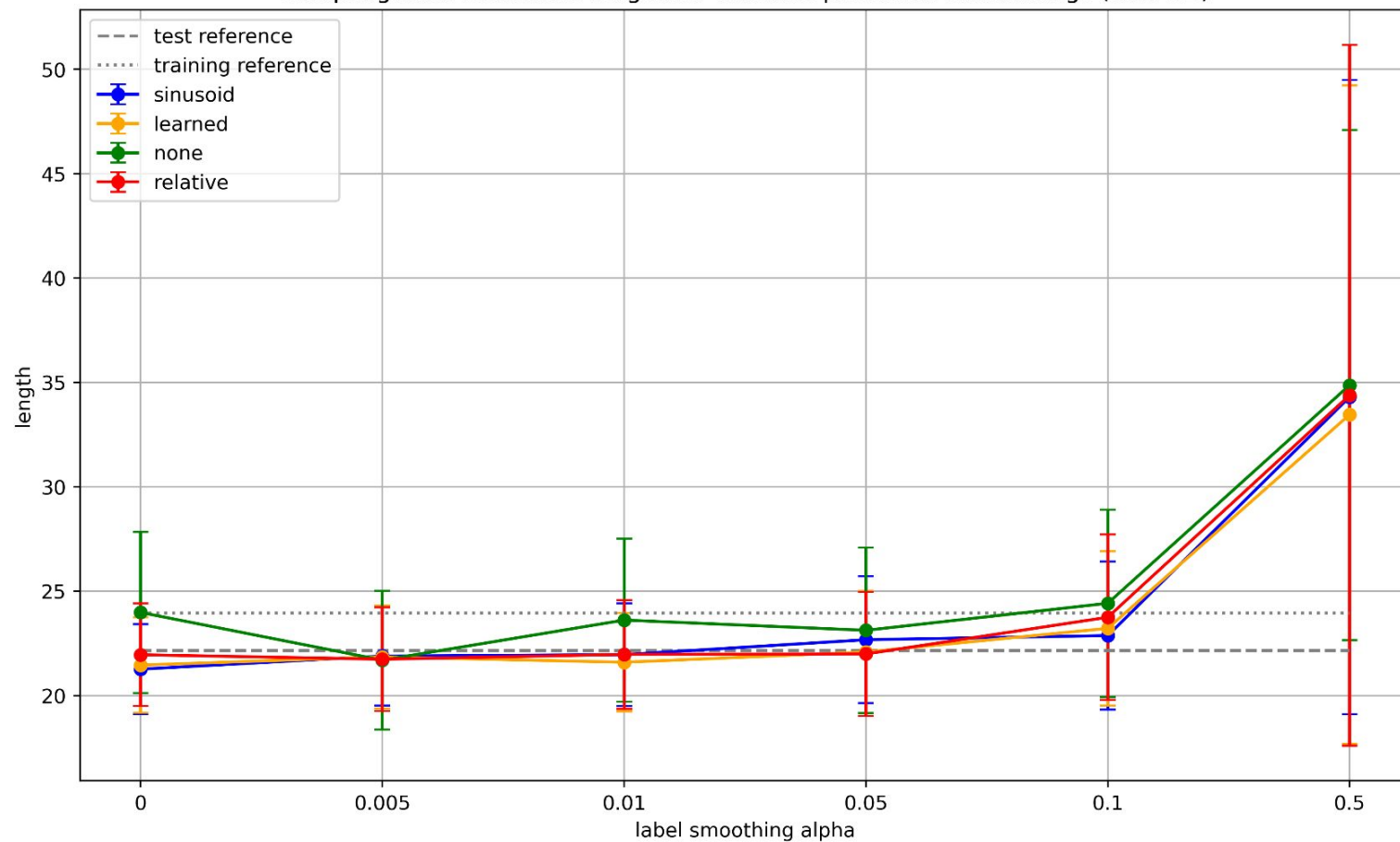
Beam search BLEU for different positional embeddings



Sampling mean sentence length for different positional embeddings (with SE)



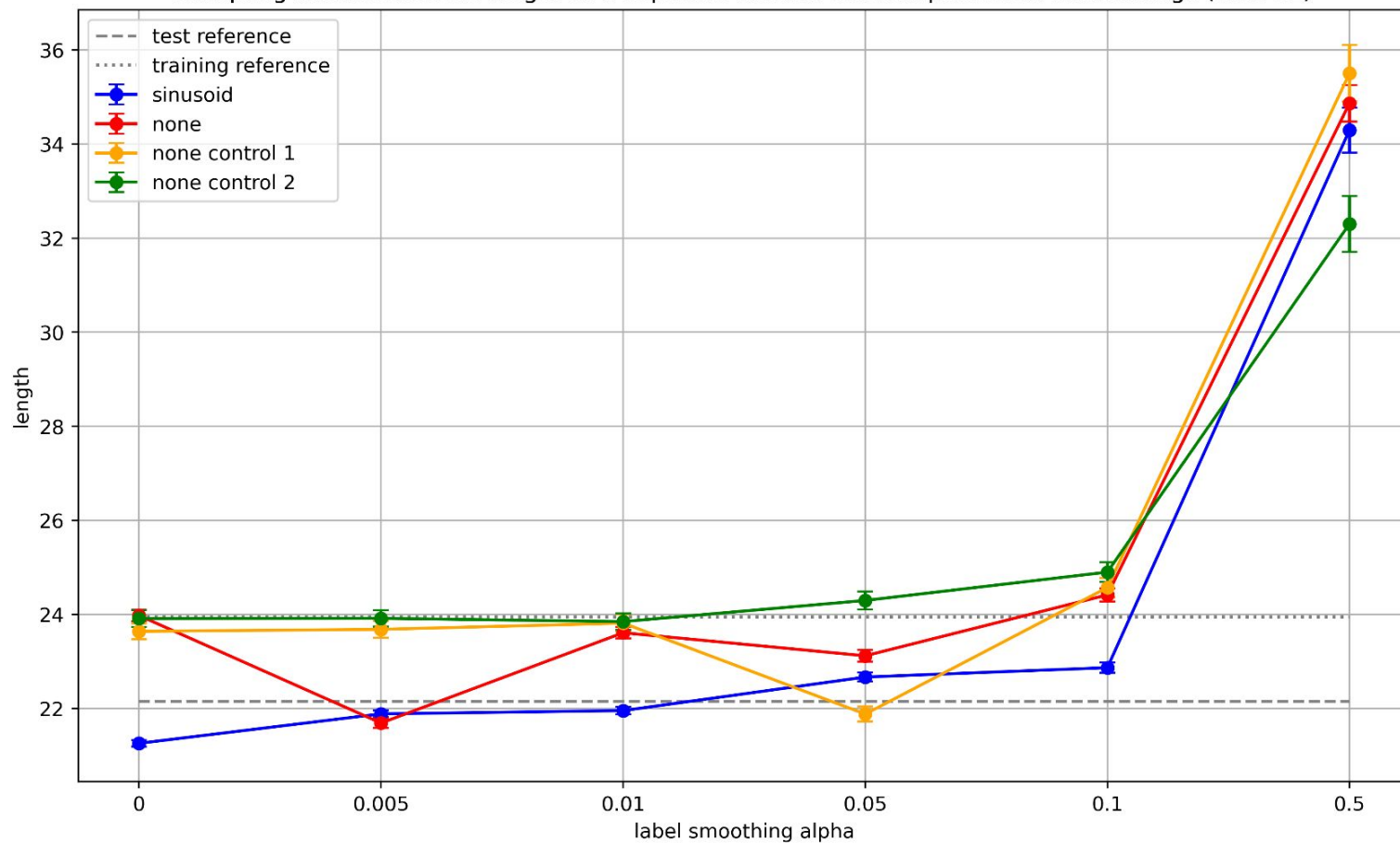
Sampling mean sentence length for different positional embeddings (with SD)



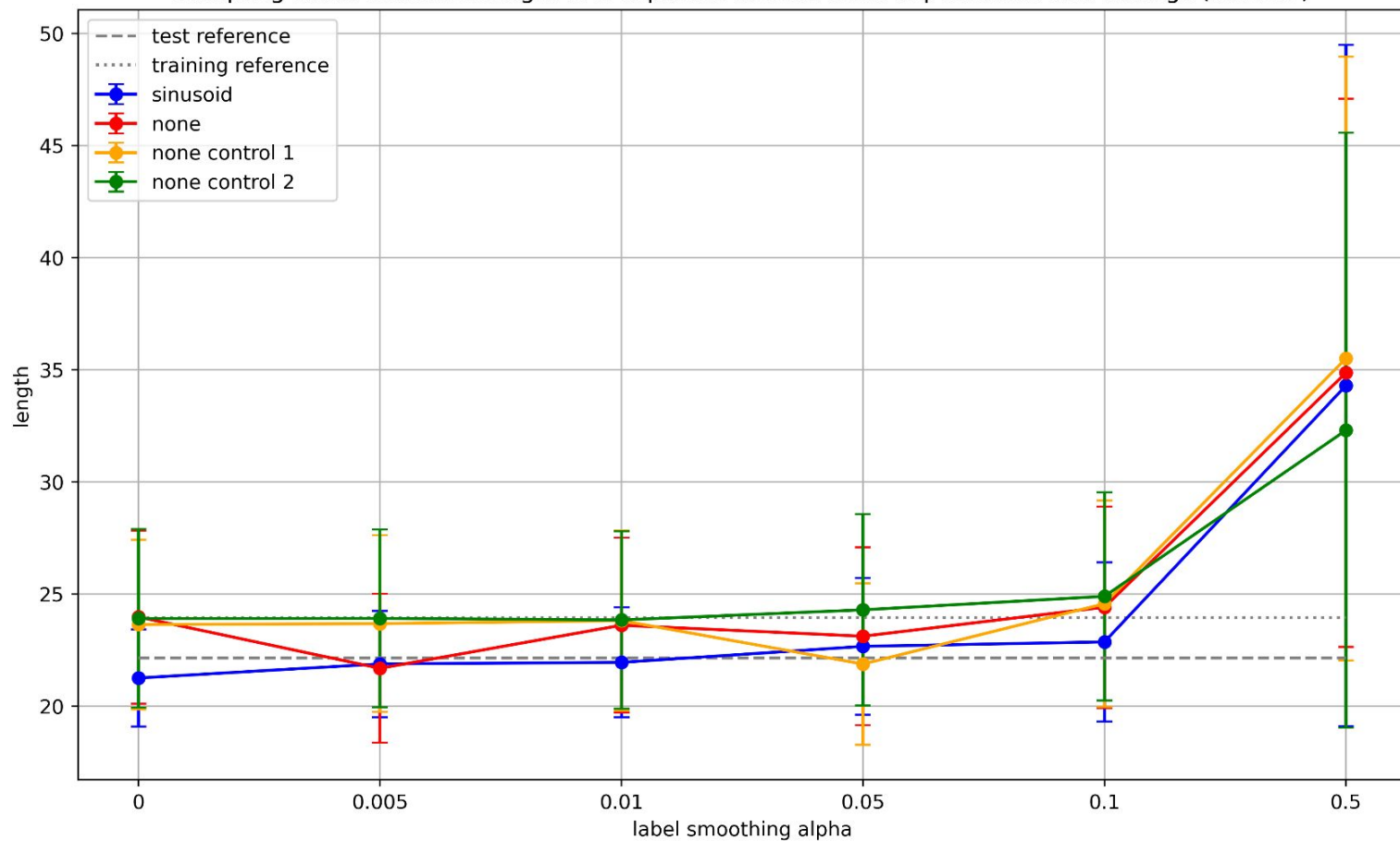
No Positional Embeddings Analysis

- Surprisingly good length ratio of models trained without positional embeddings
- 3 models trained without positional embeddings with different seeds for each label smoothing parameter
 - `fairseq-train` with `--no-token-positional-embeddings` flag
- Only 500 samples per sentence due to computational resources
- Observations:
 - Large variation across runs
 - Models seem to learn training set reference length well

Sampling mean sentence length of 3 separate models without positional embeddings (with SE)



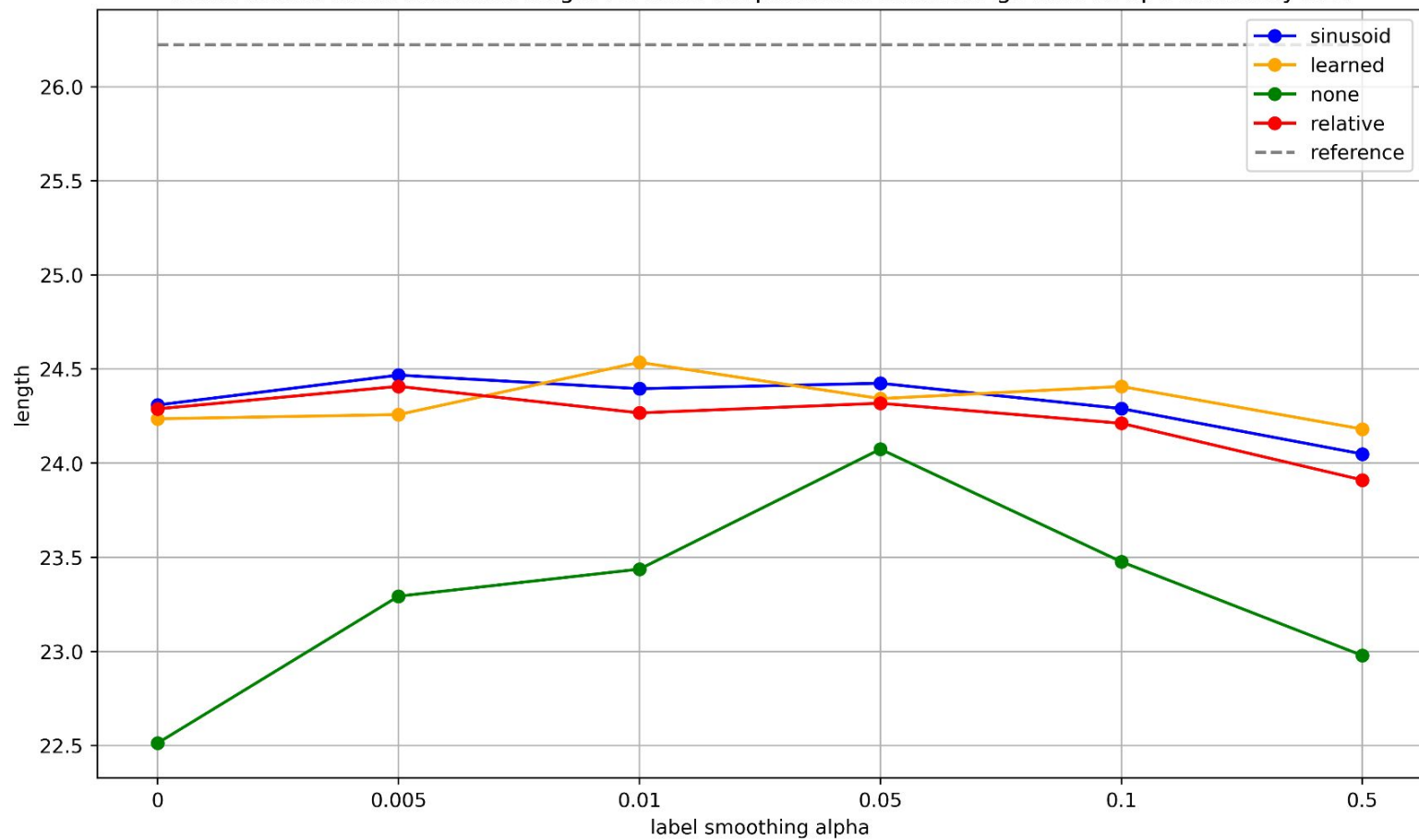
Sampling mean sentence length of 3 separate models without positional embeddings (with SD)



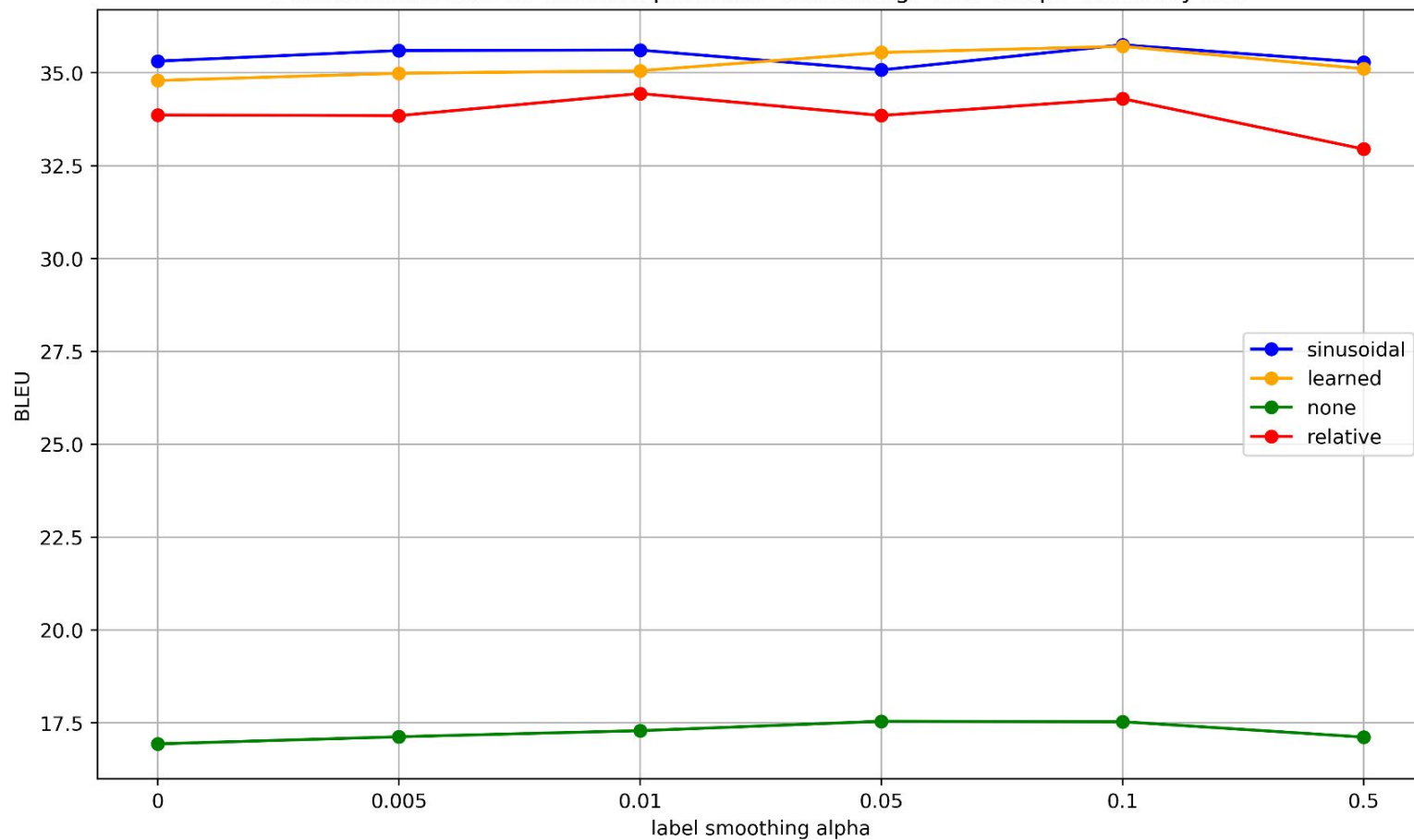
Positional embeddings dict size ablation

- Comparison between sinusoidal, learned, relative and no positional embeddings for bpe dict sizes 4k and 64k
- Same setup as Positional embedding experiments

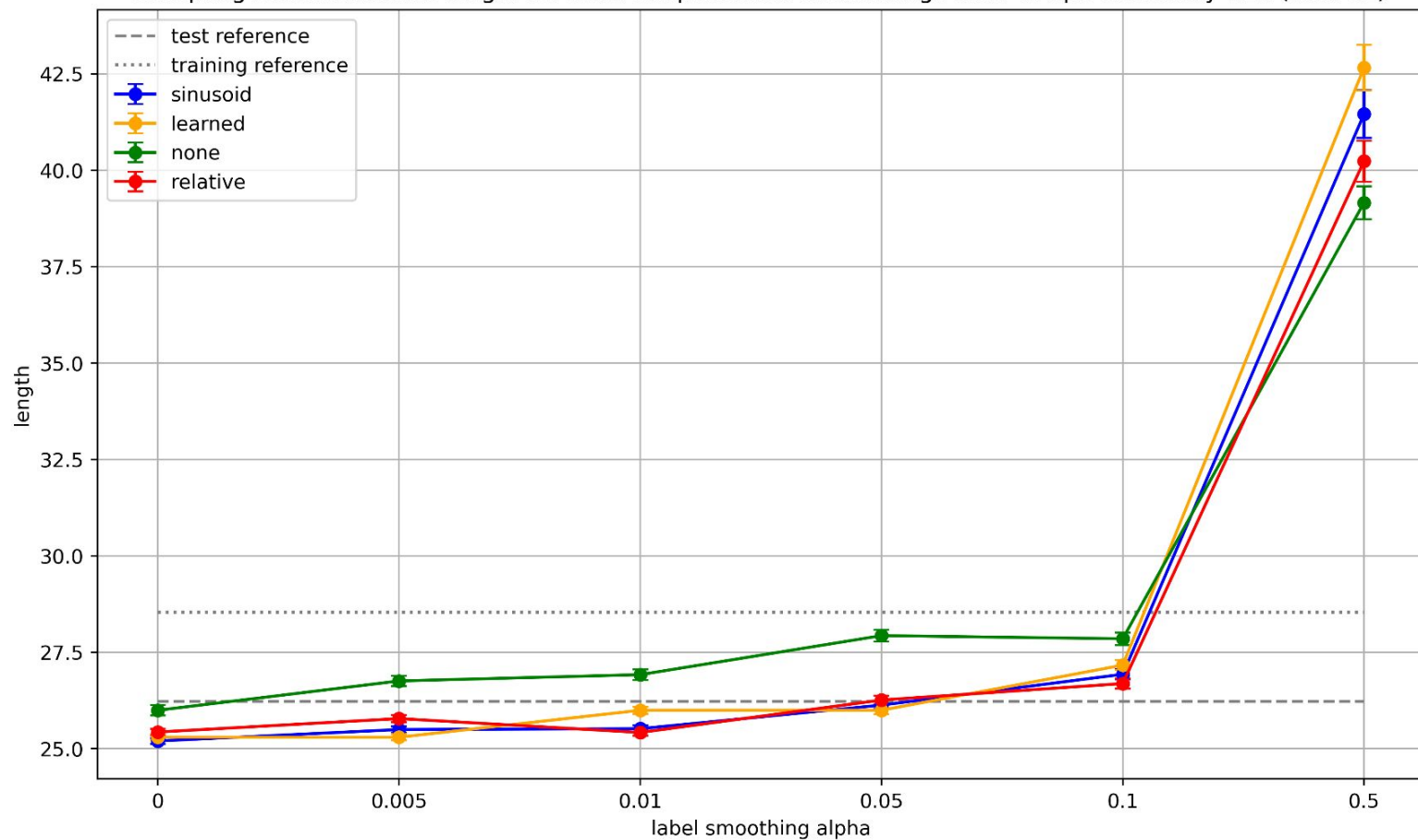
Beam search mean sentence length for different positional embeddings with 4k bpe dictionary size



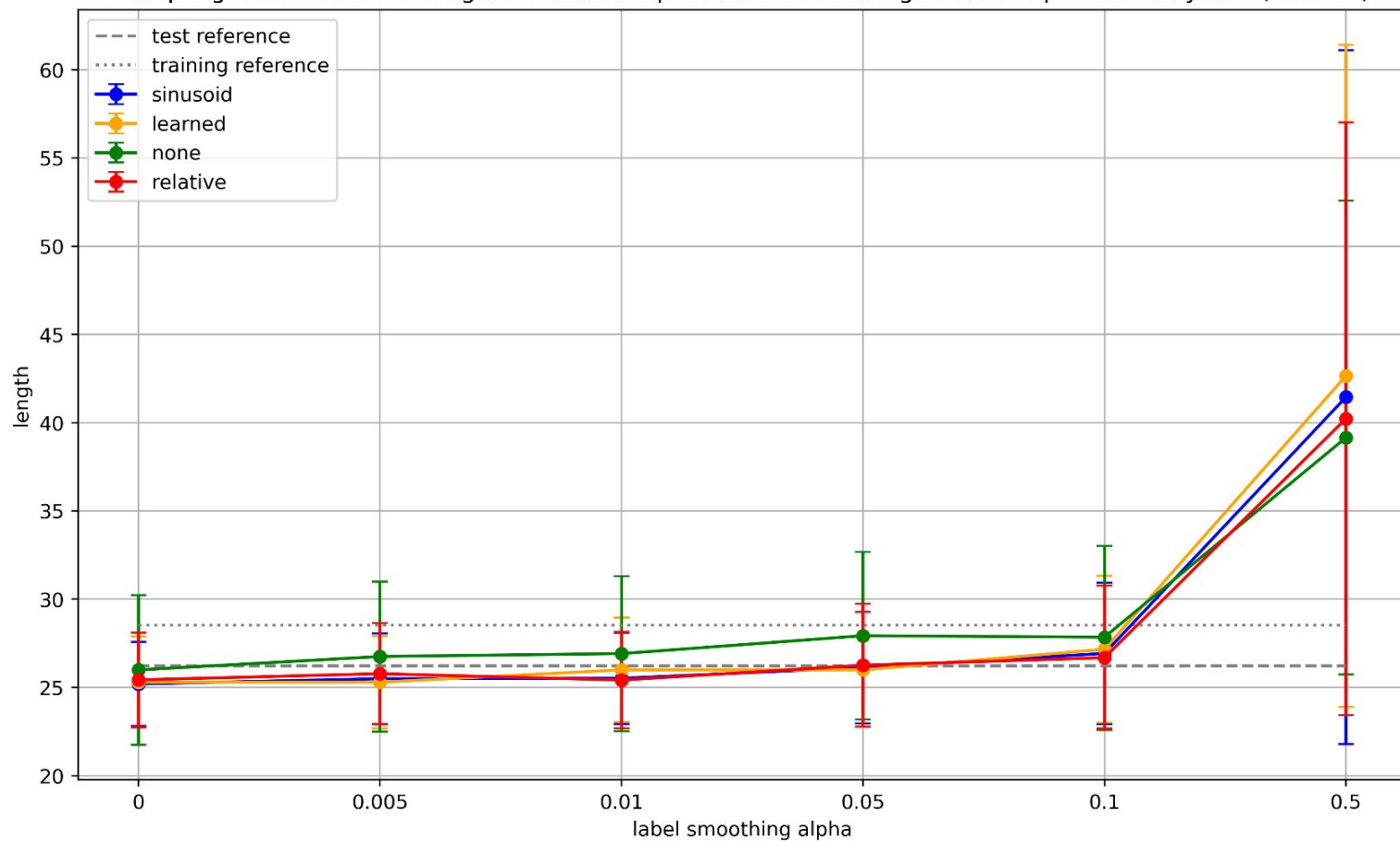
Beam search BLEU for different positional embeddings with 4k bpe dictionary size



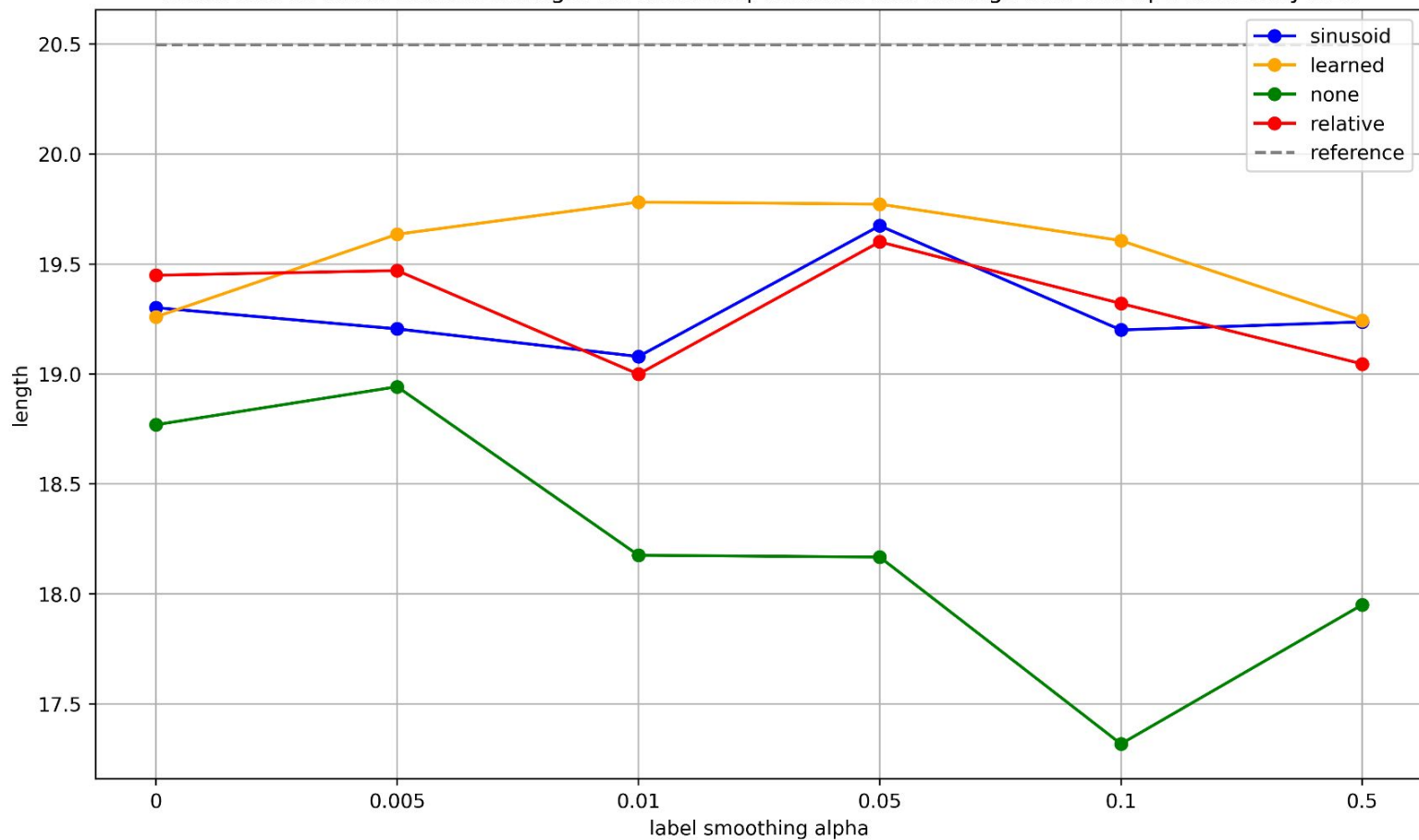
Sampling mean sentence length for different positional embeddings with 4k bpe dictionary size (with SE)



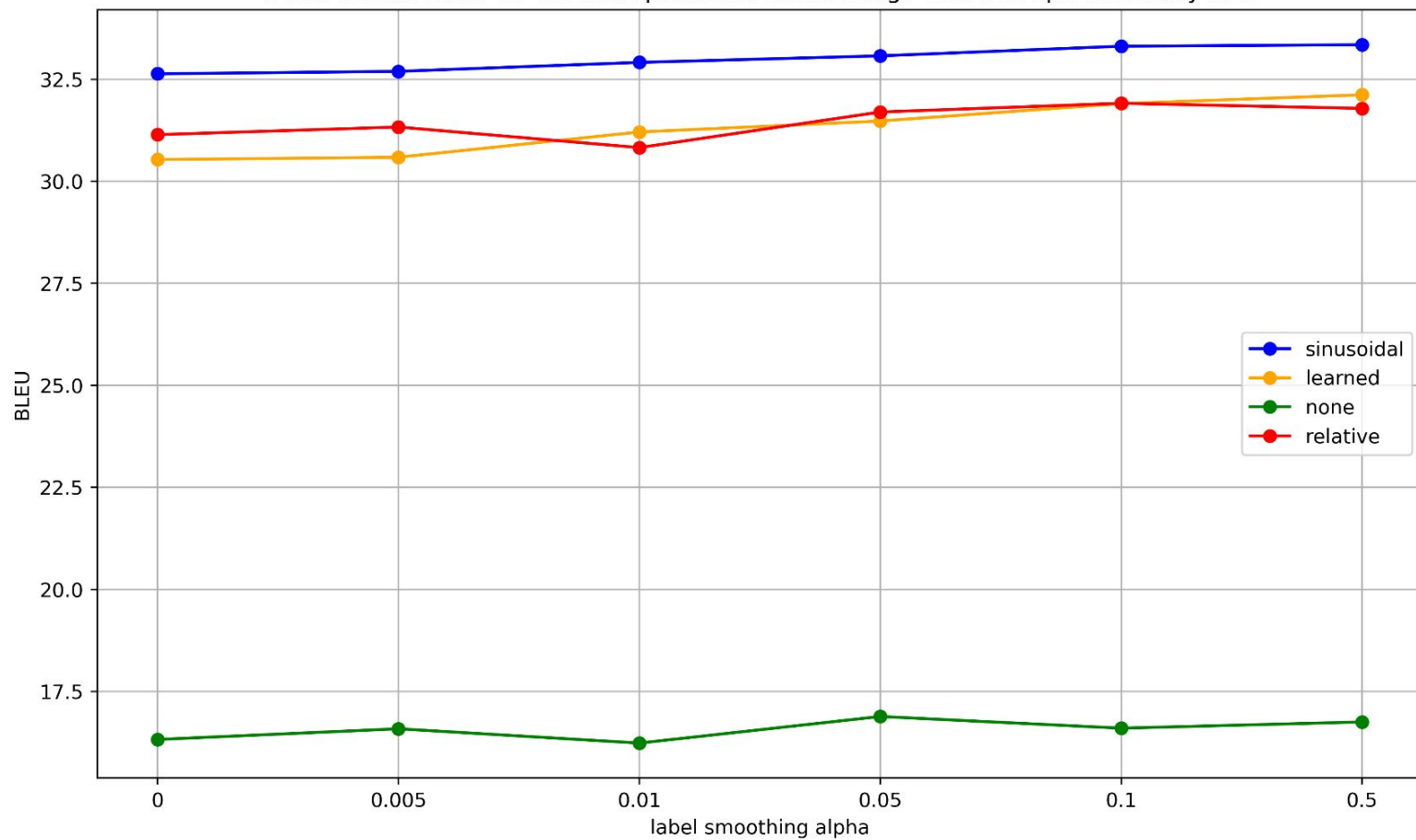
Sampling mean sentence length for different positional embeddings with 4k bpe dictionary size (with SD)



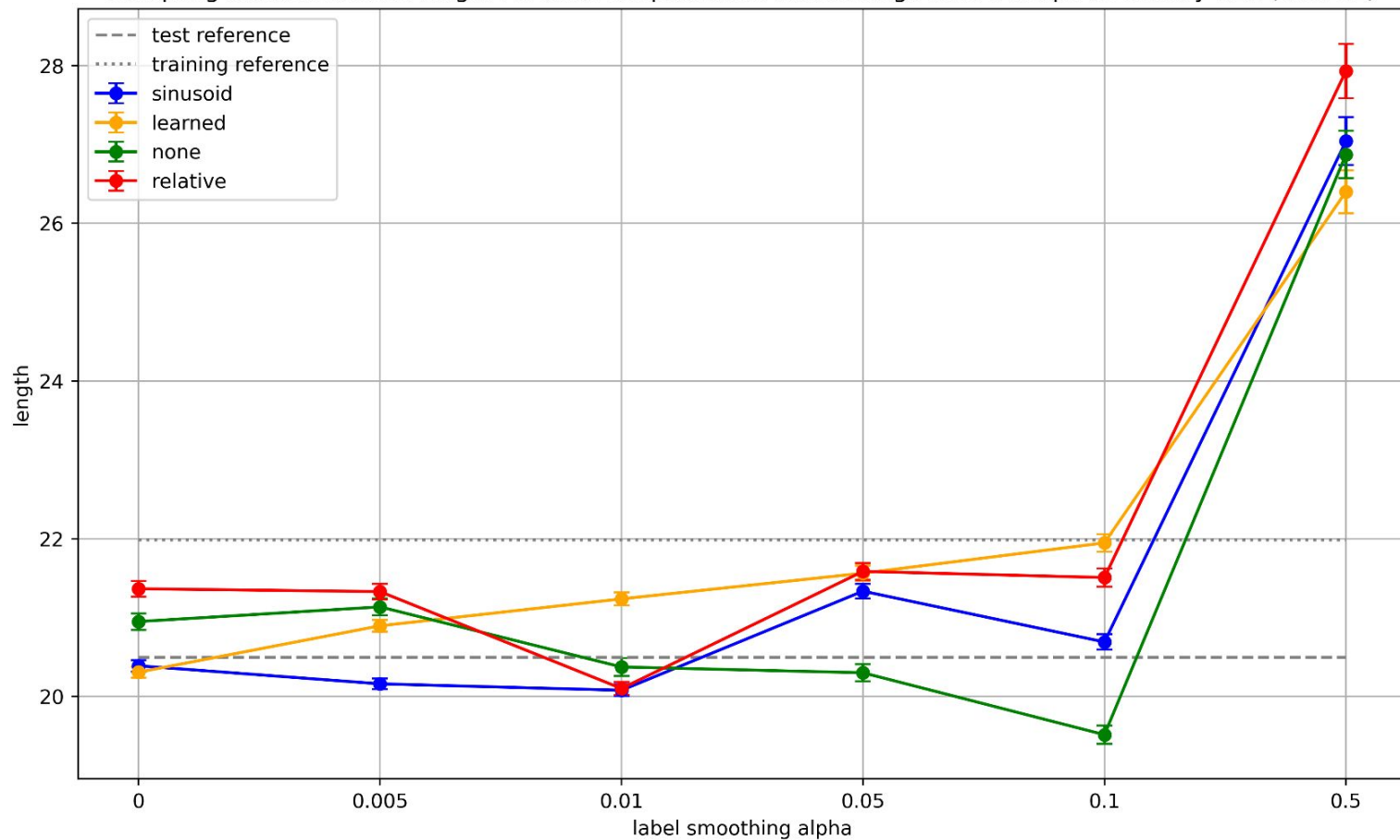
Beam search mean sentence length for different positional embeddings with 64k bpe dictionary size



Beam search BLEU for different positional embeddings with 64k bpe dictionary size



Sampling mean sentence length for different positional embeddings with 64k bpe dictionary size (with SE)



Sampling mean sentence length for different positional embeddings with 64k bpe dictionary size (with SD)

