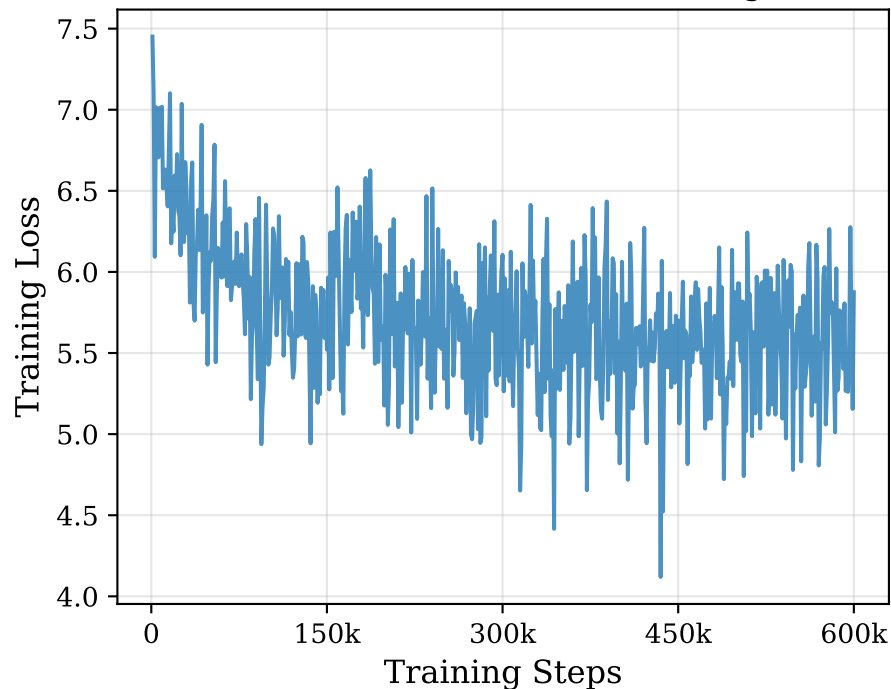
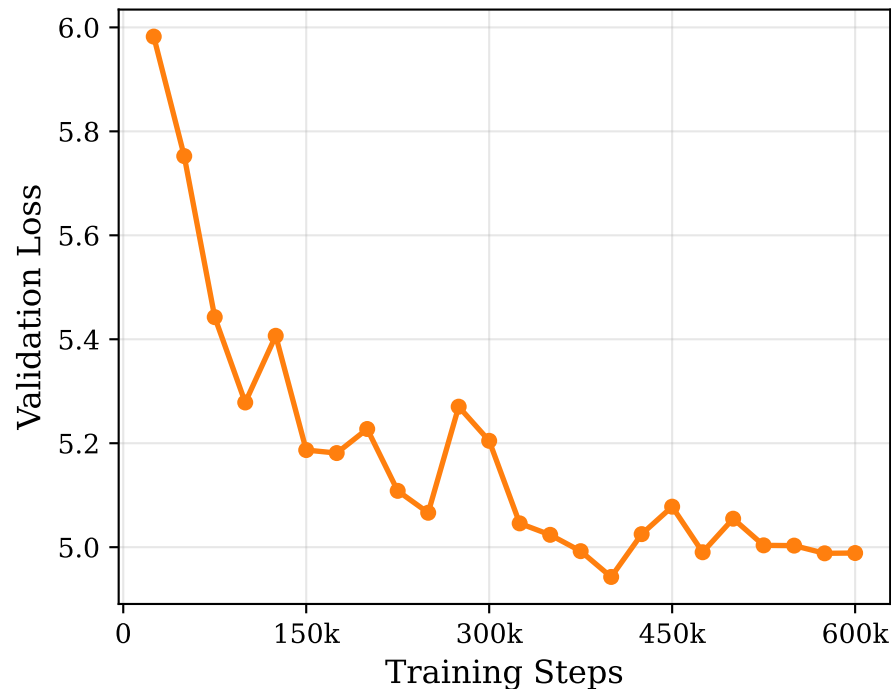


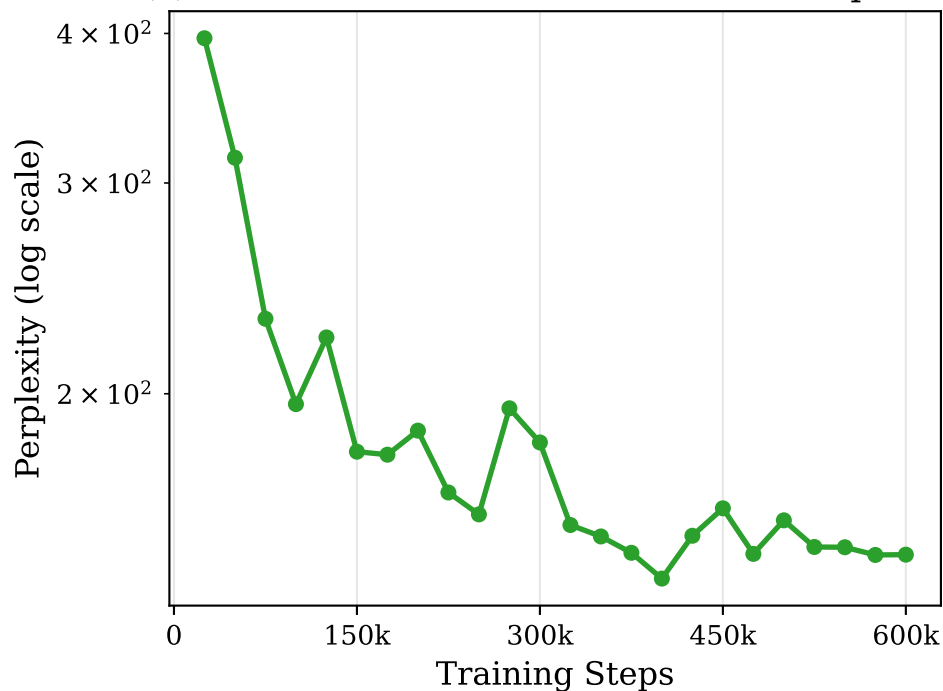
(a) Standard Transformer Training Loss



(b) Standard Transformer Validation Loss



(c) Standard Transformer Validation Perplexity



(d) Standard Transformer Generalization Gap

