

# How to Synthesize Text Data without Model Collapse?

Xuekai Zhu · Daixuan Cheng · Hengli Li · Kaiyan Zhang · Ermo Hua · Xingtai Lv · Ning Ding · Zhouhan Lin · Zilong Zheng · Bowen Zhou

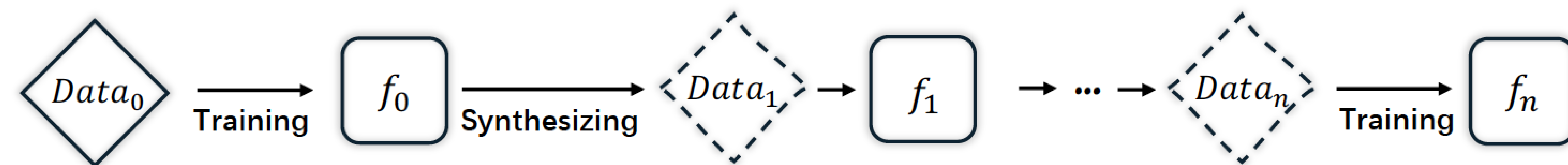
Code



Background: Future GPT-{n} models will inevitably be trained on a blend of synthetic and human-produced data.

Model collapse in synthetic data indicates that iterative training on self-generated data leads to a gradual decline in performance.

① Model Collapse Setting  $\rightarrow E_{test} = \frac{\sigma^2 d}{T-d-1} \times n$



Trained Model  $f_i$

Source Real Data:  $Data_0$

Synthetic Data:  $Data_{>0}$

Iterations  $i \in \{1, \dots, n\}$

Test Error  $E_{test}$

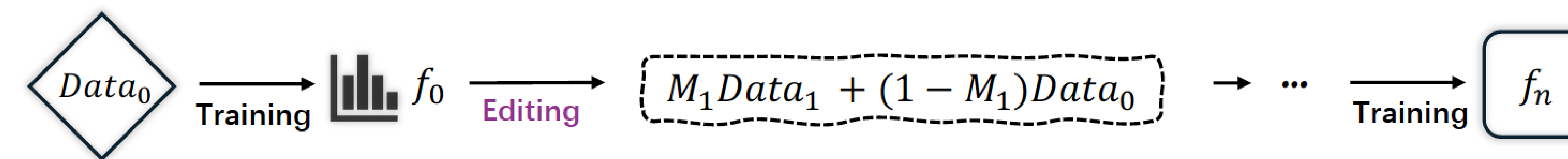
Data Size  $T$

Input Dimensions  $d$

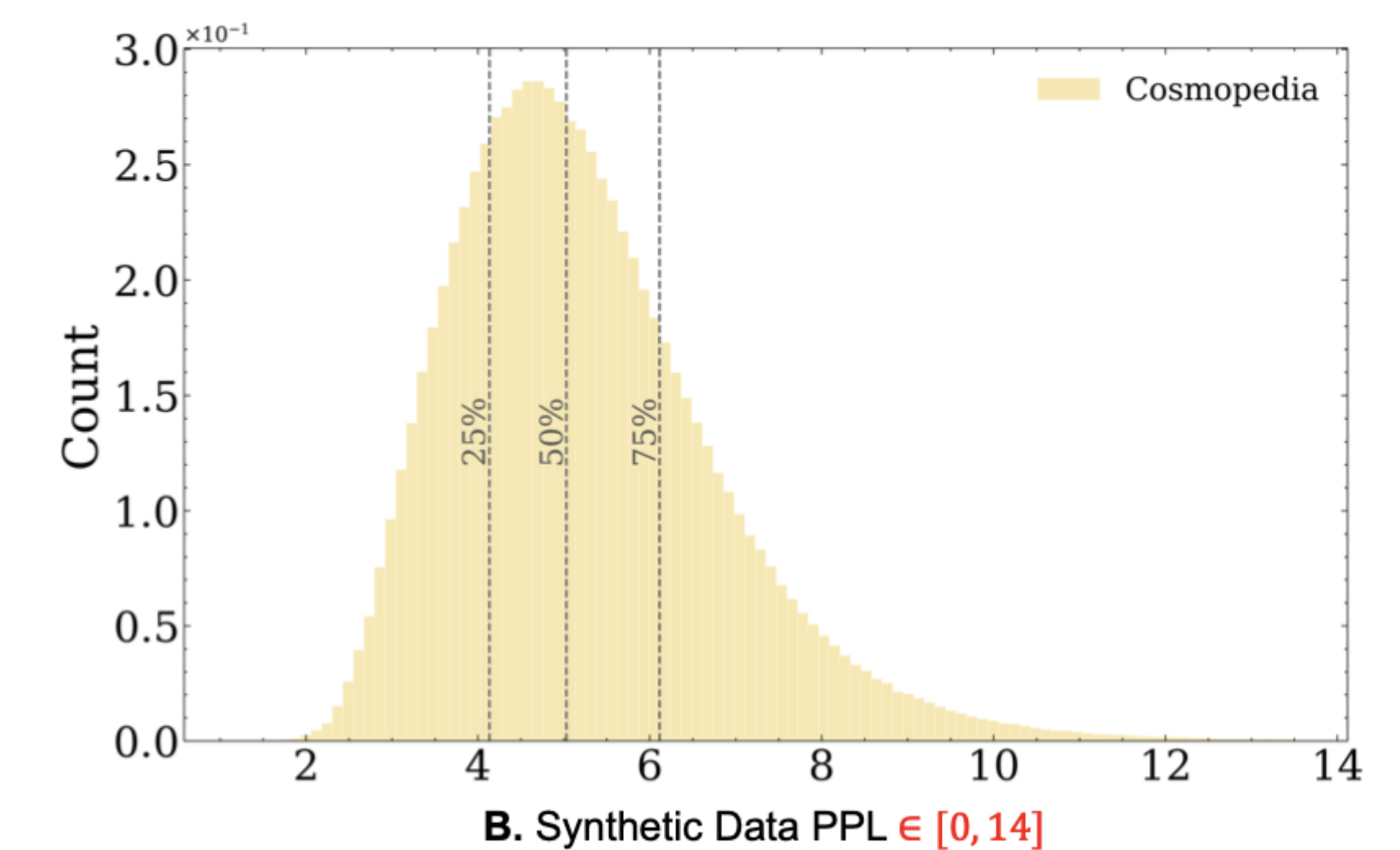
Label Noise Scalar  $\sigma$

Editing Operation Matrix  $M_i$

② Token-Level Editing  $\rightarrow E_{test} \leq \frac{\sigma^2 d}{T-d-1} \times 2 \rightarrow$  Avoiding Model Collapse



Vanishing of the Long Tail



Non-iterative model collapse

