

### Algorithm: Proxy-Tuning of Language Model

*/\* Applies proxy-tuning to adjust logits of a target model in an NLP task \*/*

**Input:** input\_ids, large\_model\_base  $M$ , small\_model\_tuned  $M^+$ , small\_model\_untuned  $M^-$ ,

**Output:** generated\_text, a sequence of tokens

**Hyperparameters:** max\_length (maximum generation length), n (number of tokens to generate)

**Parameters:**  $\Theta$  includes all parameters for the large\_model\_base  $M$ , small\_model\_tuned  $M^+$ , and small\_model\_untuned  $M^-$ .

1. Initialize generated\_tokens as an empty list
2. Encode input\_text into input\_ids using tokenizer  $\Theta_{\text{tokenizer}}$
- // Perform token-wise proxy-tuning and text generation
3. For  $t$  in  $[1, \dots, n]$ :
  - a. Obtain large\_model\_base, small\_model\_tuned, and small\_model\_untuned logits:
    - i.  $\text{large\_base\_logits} \leftarrow \text{large\_model\_base } M(\text{input\_ids})$ . logits with parameters  $\Theta_{\text{large\_model\_base}}$
    - ii.  $\text{small\_tuned\_logits} \leftarrow \text{small\_model\_tuned } M^+(\text{input\_ids})$ . logits with parameters  $\Theta_{\text{small\_model\_tuned}}$
    - iii.  $\text{small\_untuned\_logits} \leftarrow \text{small\_model\_untuned } M^-(\text{input\_ids})$ . logits with parameters  $\Theta_{\text{small\_model\_untuned}}$
  - b. Proxy-tuning adjustment:
    - i.  $\Delta\text{logit\_offsets} \leftarrow \text{small\_tuned\_logits} - \text{small\_untuned\_logits}$
    - ii.  $\text{logits}' \leftarrow \text{large\_base\_logits} + \Delta\text{logit\_offsets}$
  - c. Normalize the logits for next token prediction:
    - i.  $\text{predictions} \leftarrow \text{softmax}(\text{logits}', \text{axis}=-1)$
  - d. Select the next token:
    - i.  $\text{next\_token\_id} \leftarrow \text{argmax}(\text{predictions})$
    - ii. Append next\_token\_id to generated\_tokens
  - e. Update input\_ids with next\_token\_id for the next iteration
4. Decode the sequence of generated\_tokens into text using  $\Theta_{\text{tokenizer}}$
5. **Return** generated\_text