# Bag of Tricks for Training Data Extraction from Language Models

Weichen Yu[1,2], Tianyu Pang[2], Qian Liu[2], Chao Du[2], Bingyi Kang[2], Yan Huang[1], Min Lin[2], Shuicheng Yan[2]

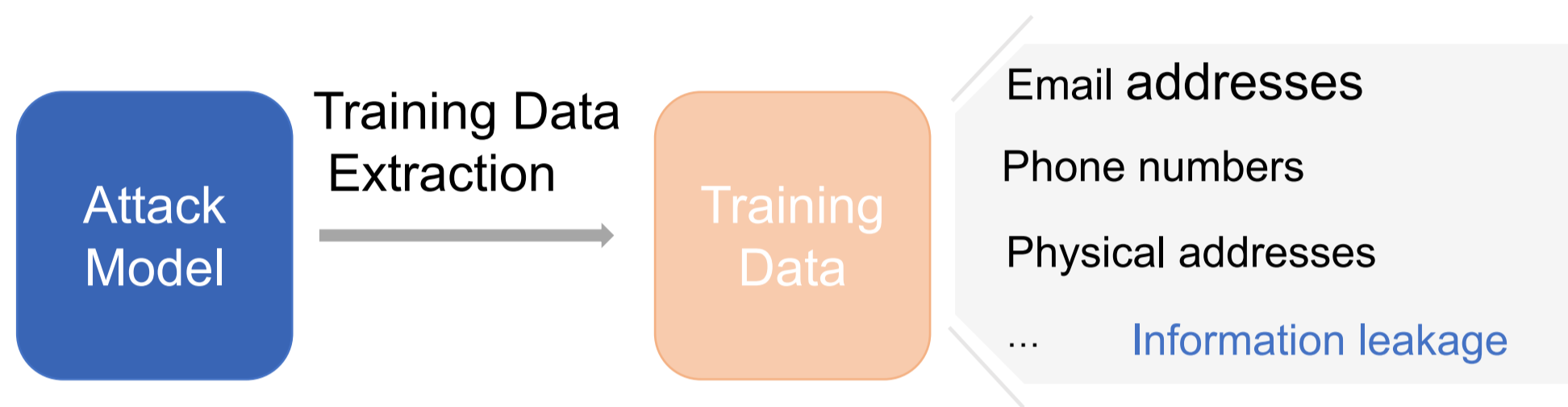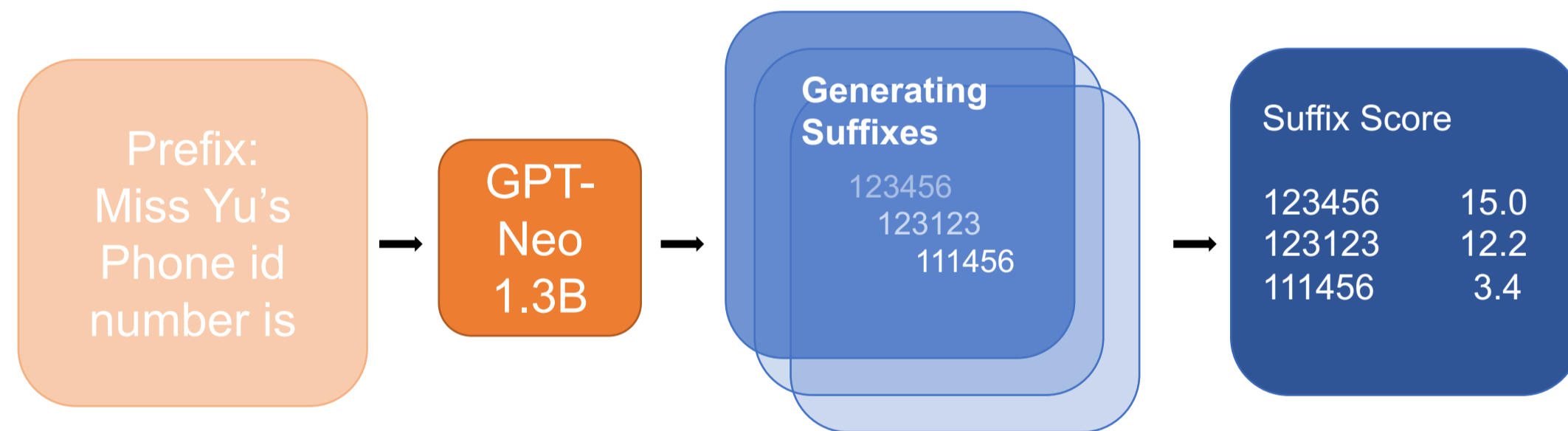[1] Institute of Automation, Chinese Academy of Sciences     [2] Sea AI Lab

## Significance of Training Data Extraction:

An effective tool to evaluate the privacy preserving ability of language models.



## Pipeline: Generating-then-ranking



Ranking by perplexity $\mathcal{P} = \exp\left(-\frac{1}{N}\sum_{n=0}^{N}\log f_\theta(x_n|x_{[0:n-1]})\right)$

1-eidetic (Carlini et al. (2021) ): the sentence $[p, s]$ appears in at most *1* example in the training data.

### Evaluation Metrics:
Precision, Recall, and Hamming distance.

## Bag of tricks

### Probability adjustment
• temperature
• repetition penalty

Table 3. Results of $\mathcal{M}_P$, $\mathcal{M}_R$, and $\mathcal{M}_H$ under different repetition penalty. Repetition penalty $r = 1$ is the baseline. All results are reported on 5 trials.

| Repetition penalty | $\mathcal{M}_P$ (%)($\uparrow$) | $\mathcal{M}_R$ (%)($\uparrow$) | $\mathcal{M}_H$ ($\downarrow$) |
|---|---|---|---|
| 0.9 | 19.8 | 66.4 | 27.927 |
| 1 | 37.0 | 76.5 | 19.614 |
| 1.1 | 37.3 | 76.5 | 20.181 |
| 1.2 | 37.1 | 76.5 | 20.323 |
| 1.3 | 36.7 | 76.4 | 20.332 |
| 1.5 | 34.7 | 75.7 | 21.154 |

### Dynamic context window

$f_\theta(x_n; \mathcal{W}) = h_{\mathcal{W}}\left(f_\theta(x_n|x_{[n-w_1,n-1]}), ..., f_\theta(x_n|x_{[n-w_m,n-1]})\right),$

$f_\theta(x_n; \mathcal{W}_w) = \frac{\sum_{i=1}^{m}\epsilon_i f_\theta(x_n|x_{[n-w_i,n-1]})}{\sum_{i=1}^{m}\epsilon_i},$

$f_\theta(x_n; \mathcal{W}_v) = \frac{1}{m}\sum_{i=1}^{m}\mathcal{V}(f_\theta(x_n|x_{[n-w_i,n-1]}));$

### Dynamic position shifting

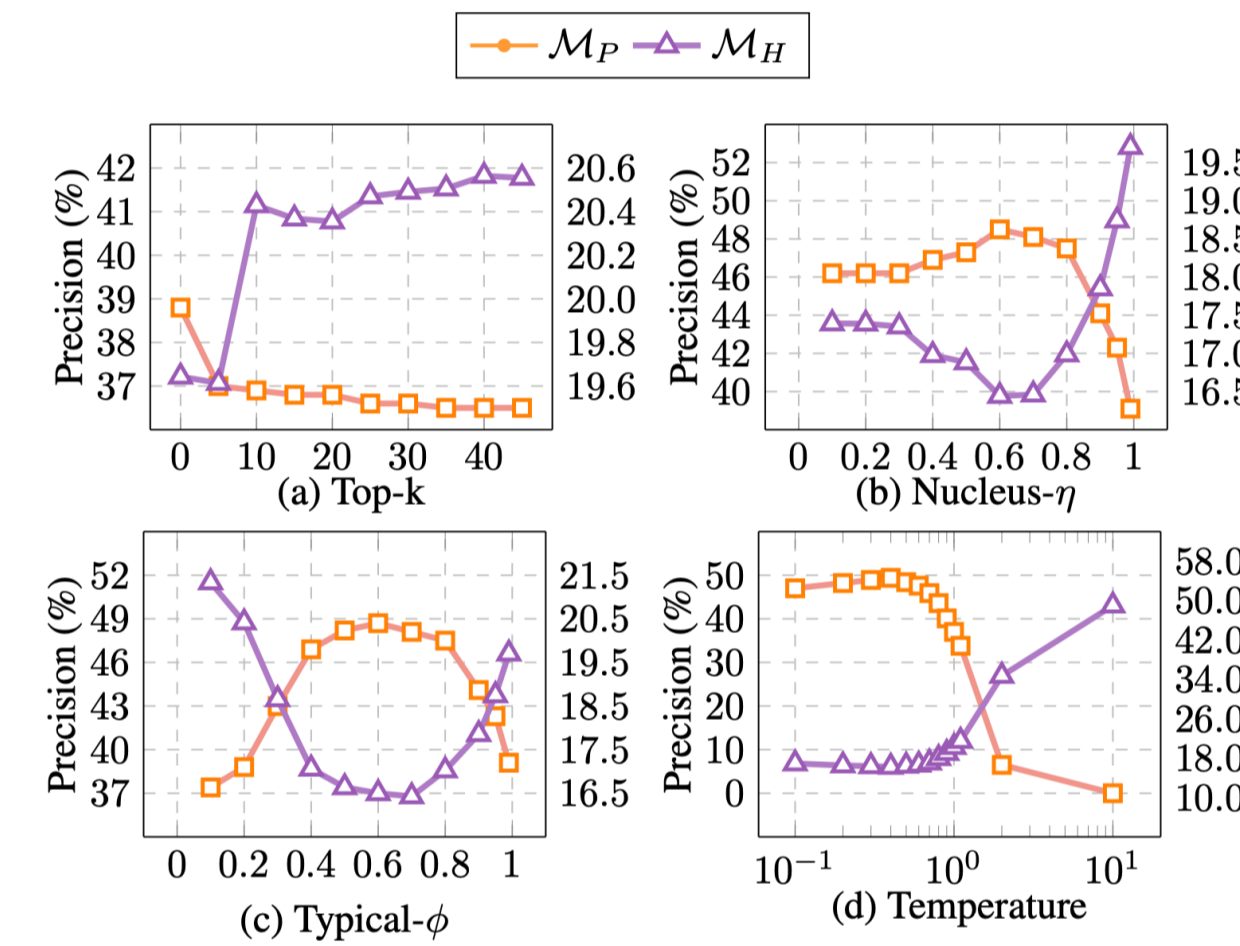$c = \arg\min_{-}\mathcal{P}(p, c^i); \quad \hat{\phi}(x_i) = \psi(c_n) + \phi(x_n),$

Table 4. Results of $\mathcal{M}_P$, $\mathcal{M}_R$, and $\mathcal{M}_H$ under context window length adjustments. All results are reported on a single trial.
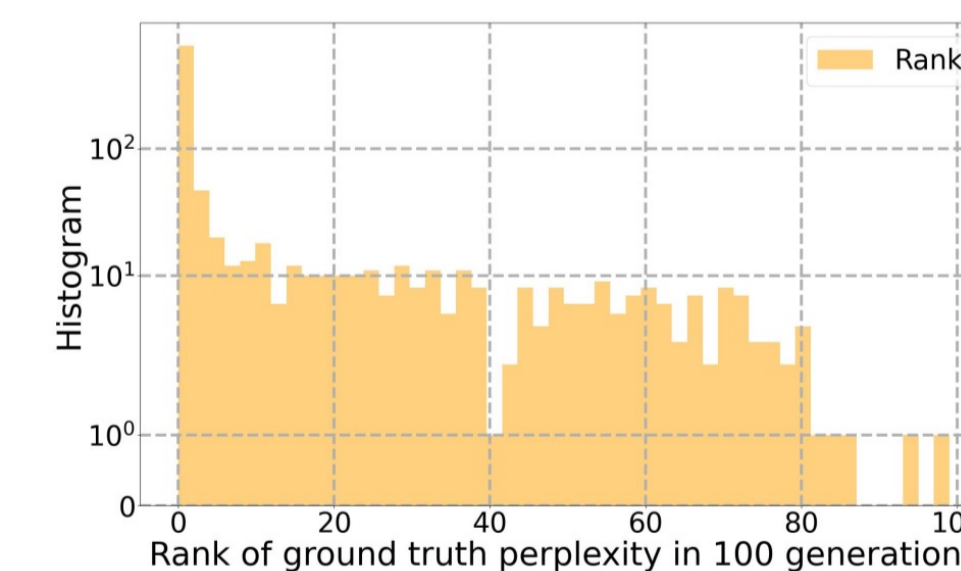
| | $\mathcal{M}_P$ (%)($\uparrow$) | $\mathcal{M}_R$ (%)($\uparrow$) | $\mathcal{M}_H$ ($\downarrow$) |
|---|---|---|---|
| Baseline | 19.5 | 65.6 | 26.948 |
| Context Win $\mathcal{W}_w$ | 47.4 | 77.6 | 16.993 |
| Context Win $\mathcal{W}_v$ | 46.7 | 77.5 | 17.164 |
| Position Shifting | 16.4 | 39.0 | 21.154 |

## Sampling strategy
• top-k sampling
• nucleus sampling
• typical sampling



### Cumprod $\mathcal{L}_c = (\prod_{n=0}^{N}\log p(x_n|x_0,..,x_{n-1}))^{-N}$



### Look-ahead

$f_\theta(x_n|x_{n+1}, x_{<n})$

$= \frac{f_\theta(x_{n+1}|x_n, x_{<n})f_\theta(x_n|x_{<n})}{\sum_{x'_n}f_\theta(x_{n+1}|x'_n, x_{<n})f_\theta(x'_n|x_{<n})},$

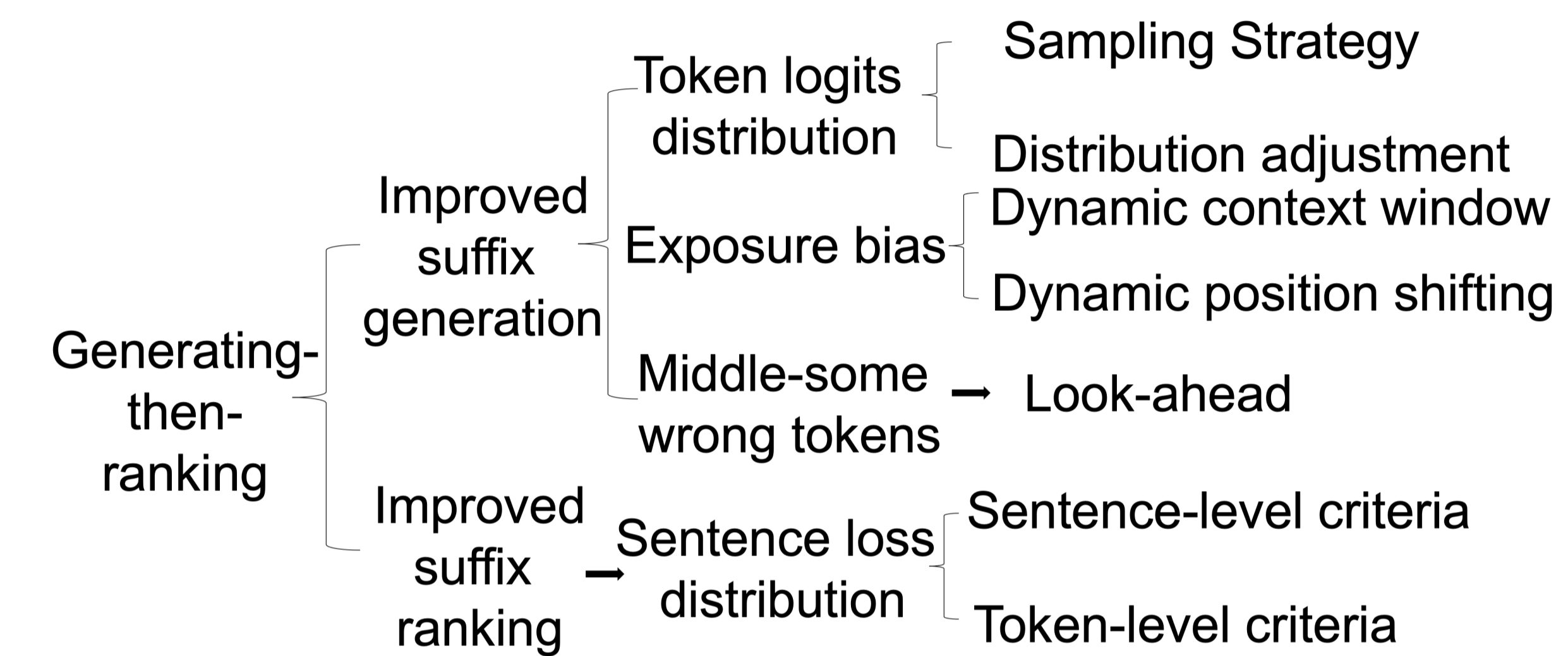$\mathcal{X} = \{x'_n|\mathcal{R}(f_\theta(x'_n|x_{<n})) \geq \lambda\}.$

## Taxonomy of the evaluated tricks



## Overall evaluation