# Transformer Language Models
# Handle Word Frequency in Prediction Head

Goro Kobayashi [1,3], Tatsuki Kuribayashi [2,1], Sho Yokoi [1,3], Kentaro Inui [1,3]

[1] Tohoku University    [2] MBZUAI   [3] RIKEN

✉   goro.koba@dc.tohoku.ac.jp

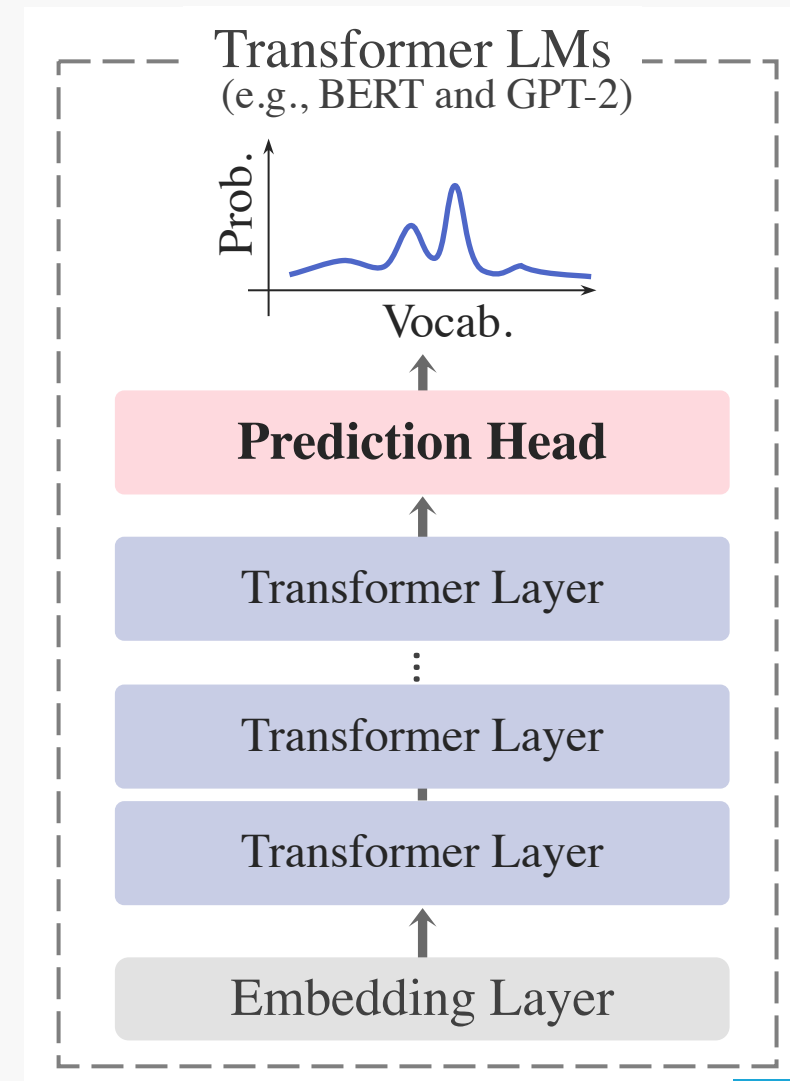https://github.com/gorokoba560/transformer-lm-word-freq-bias

TOHOKU
NLP
GROUP

Findings of ACL 2023

# **Prediction head** has been overlooked in Transformer analyses

- **Transformer layer** has been typically analyzed
  - Analyses of **Attention** [Clark+'19;Kobayashi+'20;etc.]
  - Analyses of **Feed-forward network**
    [Geva+'21;Dai+'22;etc.]

- **Prediction head** is the last block of LMs
  - Can directly impact on prediction
  - However, it has been overlooked in previous analyses…

  ➡ **We investigate its inner workings!**
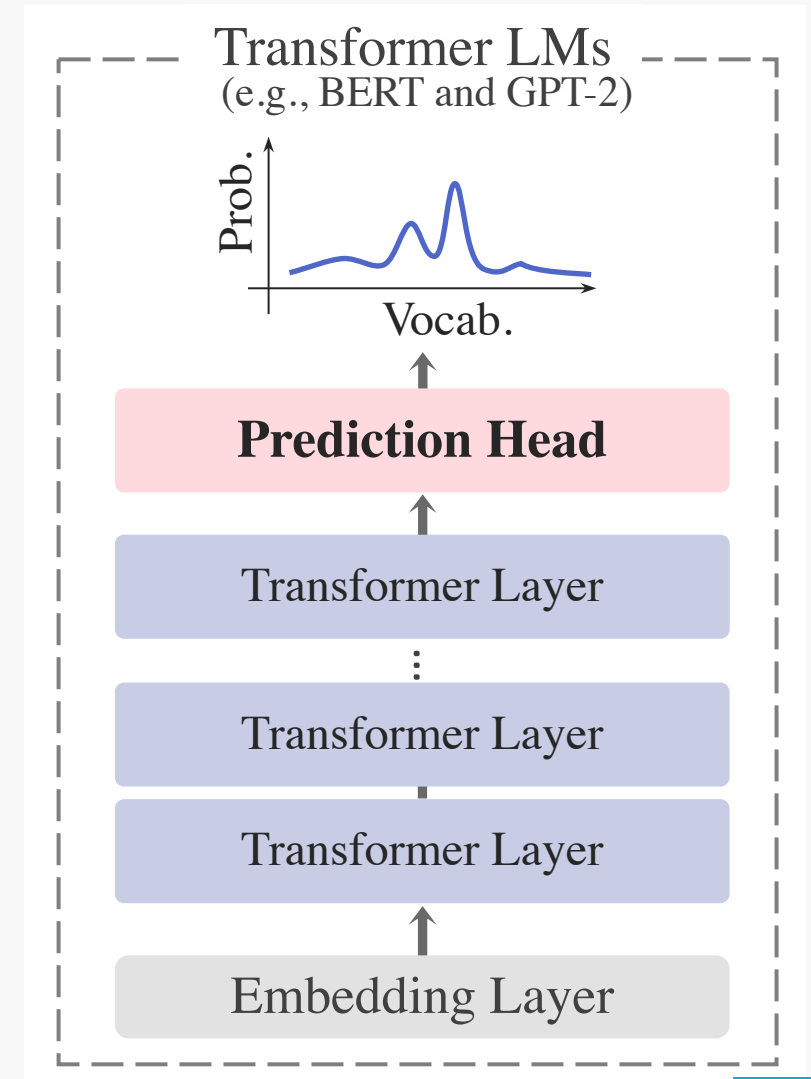
Transformer LMs
(e.g., BERT and GPT-2)

Prob.

Vocab.

**Prediction Head**

Transformer Layer

⋮

Transformer Layer

Transformer Layer

Embedding Layer

# We focus on bias parameters in prediction head

- Prediction head has **bias parameters**
  - BERT has three biases: $b_{\mathrm{FC}}$, $b_{\mathrm{LN}}$, $b_{\mathrm{last}}$
  - GPT-2 has one bias: $b_{\mathrm{LN}}$

➡ **We focus on these bias parameters!**

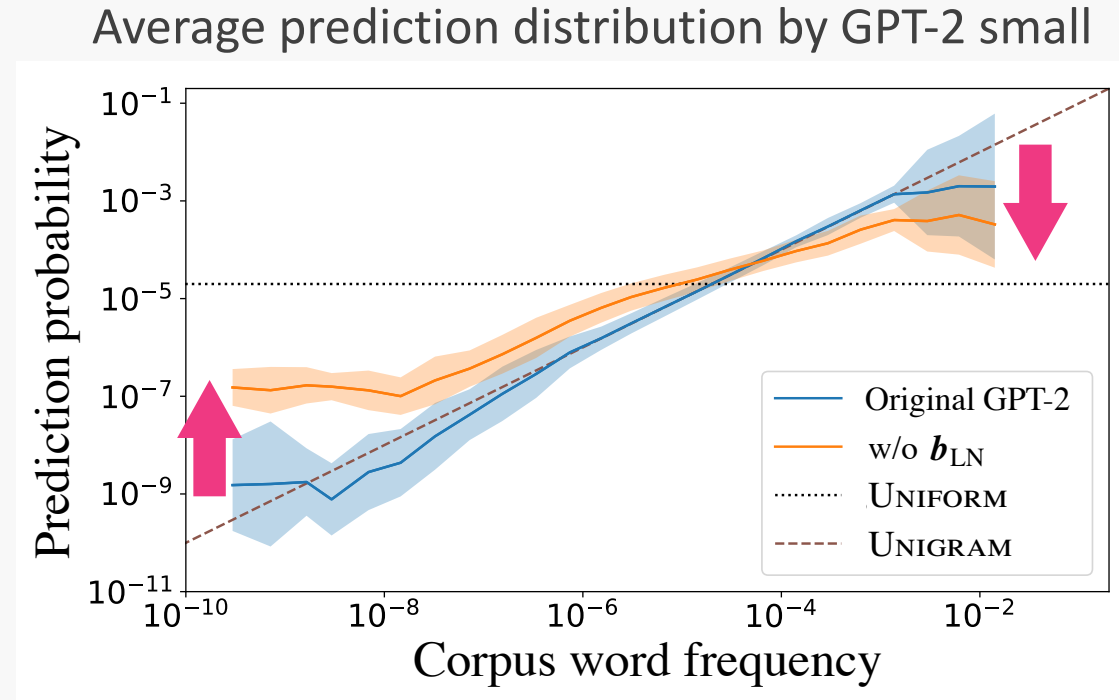Bias parameters can be easily mapped to the output space (word prediction)



Transformer LMs
(e.g., BERT and GPT-2)

Prob.

Vocab.

**Prediction Head**

Transformer Layer

Transformer Layer

Transformer Layer

Embedding Layer

# Finding 1: Bias adjusts word prediction according to word frequency

- When removing a bias $\boldsymbol{b}_{\mathrm{LN}}$ (━━ ➡ ━━)
  - Probability of high-frequency words is decreased ⬇
  - Probability of high-frequency words is increased ⬆

‖

- **Bias $\boldsymbol{b}_{\mathrm{LN}}$ adjusts word prediction**
  - **to promote high-frequency words**
  - **to discourage low-frequency words**

Average prediction distribution by GPT-2 small



Original GPT-2
w/o $\boldsymbol{b}_{\mathrm{LN}}$
UNIFORM
UNIGRAM

# Finding 2: Controlling the bias can encourage more diverse language generation

- Control the bias $\boldsymbol{b}_{\text{LN}}$ with coefficient $\lambda \in [0,1]$

$$\boldsymbol{b}_{\text{LN}} \leftarrow \lambda \boldsymbol{b}_{\text{LN}}$$

- For large models, weakening $\boldsymbol{b}_{\text{LN}}$
  - Improves diversity
  - Maintains quality

| Model | $\lambda$ | Diversity ↑ | | | Quality | |
|-------|-----------|-------|-------|------|----------|--------|
| | | $D_1$ | $D_2$ | $D$ | MAUVE ↑ | PPL ↓ |
| large | 1 | 0.04 | 0.30 | 0.47 | 0.90 | **12.7** |
| | 0.5 | 0.04 | 0.36 | 0.50 | **0.91** | 12.9 |
| | 0 | **0.04** | **0.42** | **0.54** | 0.86 | 13.6 |
| xl | 1 | 0.04 | 0.30 | 0.47 | 0.90 | **11.4** |
| | 0.7 | 0.04 | 0.34 | 0.49 | **0.92** | 11.5 |
| | 0 | **0.04** | **0.41** | **0.53** | 0.86 | 12.1 |

Thank you for listening!
Feel free to ask or comment!

We hope you read the paper
for more details, other findings, and discussions!