# TARM: Token Averaging Recurrent Memory Transformer

Paper presented as part of the presentation 1 :

```
@inproceedings{bulatovrecurrent,
  title={Recurrent Memory Transformer},
  author={Bulatov, Aydar and Kuratov, Yuri and Burtsev, Mikhail},
  booktitle={Advances in Neural Information Processing Systems (NeurIPS)},
  year={2022}
}
```

**Course Instructor**

Prof. C. Krishna Mohan

**Presenter**

Rahul Vigneswaran K[*]

CS23MTECH02002

**Assigned TA**

Peketi Divya

*Presented as part of the coursework for Visual Computing (CS6450)
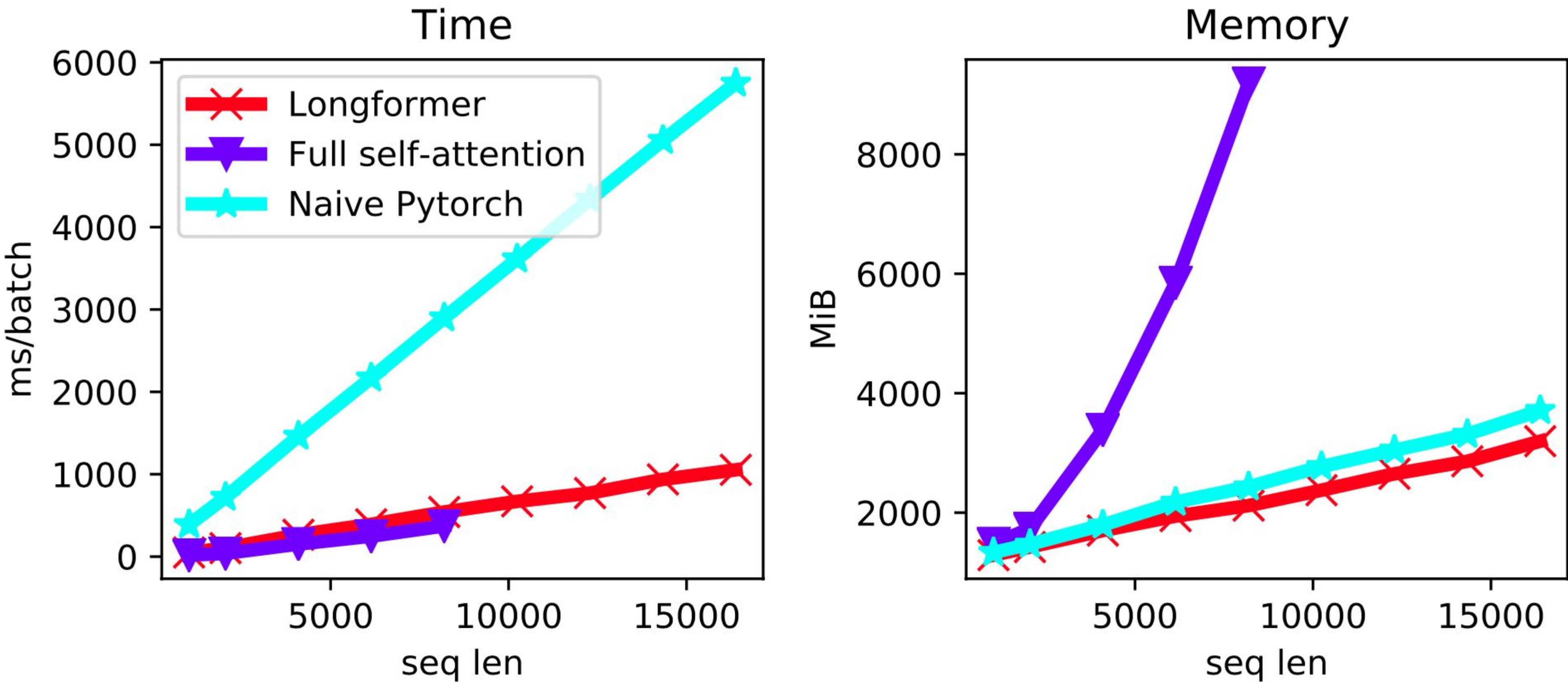
# Table of Contents

**01** **Recap of presentation 1**
- What problem are we trying to solve?
- Methodology : RMT

**02** **Implementation**
- Implementation details
- Reproduced Results

**03** **Novelty**
- Our Contributions
- Methodology : TARM
- TARM increases stability
- Results
  - Comparison with RMT
  - Ablation

# Introduction: What problem are we trying to solve?

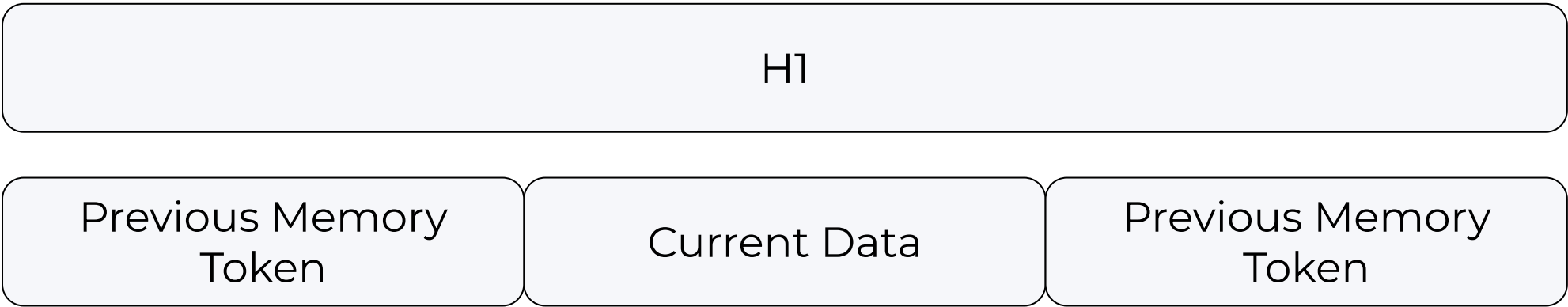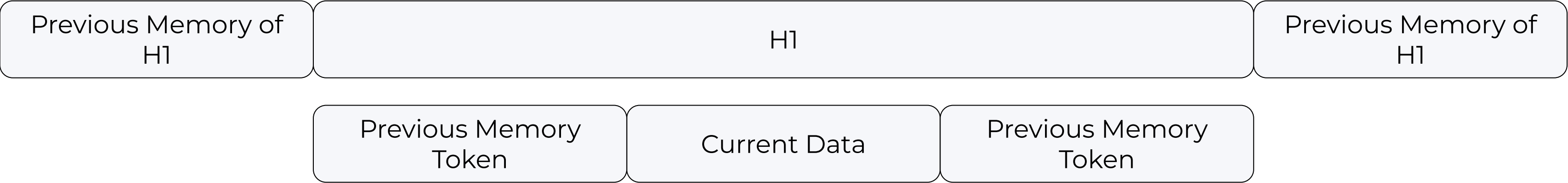# Introduction: What problem are we trying to solve?



Beltagy, Iz, Matthew E. Peters, and Arman Cohan. "Longformer: The Long-Document Transformer." (Arxiv 2020)

# Methodology

Current Data

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

| Previous Memory Token | Current Data | Previous Memory Token |
|:---:|:---:|:---:|

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

| H1 |
|:---:|

| Previous Memory Token | Current Data | Previous Memory Token |
|:---:|:---:|:---:|

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

| Previous Memory of H1 | H1 | Previous Memory of H1 |
|---|---|---|

| Previous Memory Token | Current Data | Previous Memory Token |
|---|---|---|

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

8

# Methodology



| Previous Memory of Hn | Hn | Previous Memory of Hn |
|---|---|---|
| Previous Memory of H2 | H2 | Previous Memory of H2 |
| Previous Memory of H1 | H1 | Previous Memory of H1 |

| Previous Memory Token | Current Data | Previous Memory Token |
|---|---|---|

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



| | | |
|---|---|---|
| | Output | |
| Previous Memory of Hn | Hn | Previous Memory of Hn |
| Previous Memory of H2 | H2 | Previous Memory of H2 |
| Previous Memory of H1 | H1 | Previous Memory of H1 |
| | Previous Memory Token    Current Data    Previous Memory Token | |

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

Current Data

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

| Previous Memory Token | Current Data | Previous Memory Token |
| --- | --- | --- |

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

| H1 |
|---|

| Previous Memory Token | Current Data | Previous Memory Token |
|---|---|---|

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Methodology



Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Implementation details

- Language Modelling Task
  - Model
    - Transformer
      - 16-layer
      - 10 heads
      - 410 in/out neuron count
      - 2100 intermediate neuron count
  - Batch size - 60
  - Segment size - 160
  - Optimizer
    - Adam
      - Learning Rate Schedule
        - Linear schedule learning rate starting from 0.00025 for 200,000 steps

- Quadratic Equations Task
  - Model
    - Transformer
      - 6-layer
      - 6 heads
      - 128 in/out neuron count
      - 256 intermediate in/out neuron count
  - Batch size - 32
  - Segment size - 180
  - Optimizer
    - Adam
      - Learning Rate Schedule
        - Constant learning rate 1e-4 with reduction on plateau with decay factor of 0.5 for 600,000 steps

Due to compute constraints, we were able to run **TARM (Ours)** only for 120,000 steps, so we compare the corresponding best RMT result for the same number of steps.

Bulatov, Aydar, Yuri Kuratov, and Mikhail Burtsev. "Recurrent Memory Transformer." Advances in Neural Information Processing Systems. (NeurIPS 2022)

# Reproduced Results

| Task | Dataset | Metric | RMT (From Paper) | RMT (Reproduced) |
|------|---------|--------|------------------|------------------|
| LM* | WT-103 | PPL ↓ | **23.99** | 24.291 (↑0.301) |
| Step-by-Step | Quadratic Equations | Accuracy ↑ | 99.8 | **99.9** (↑0.1) |

Table 1. Comparison of reproduced RMT results on various tasks. The best results are shown in bold, and the second-best results are underlined. The symbol ↑ and ↓ indicate the improvement and degradation of the result compared to the reproduced RMT model. LM* denotes language modeling task. PPL denotes perplexity.
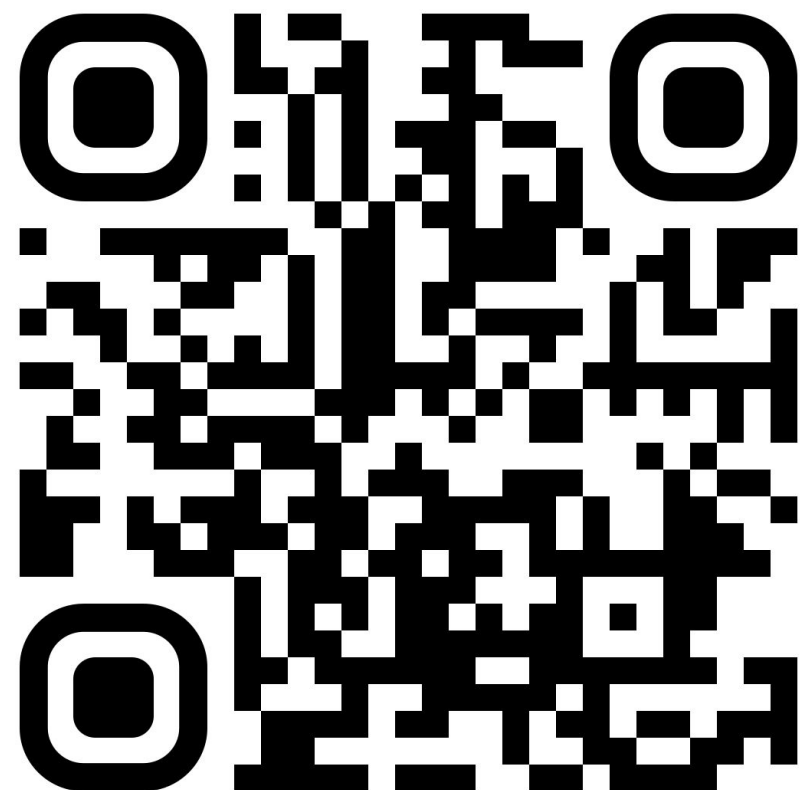
# Reproduced Results



Figure 1. Reproduced step-wise perplexity score based on both training and validation sets. We use a log-scale for y-axis for better visibility of the lower scores which are obstructed by the outliers.

23

# Our Contributions

# Our Contributions

**01** We introduce a novel token based memory method called **TARM** that increases the capacity of the memory token without actually increasing the memory token size.

# Our Contributions

**01** We introduce a novel token based memory method called **TARM** that increases the capacity of the memory token without actually increasing the memory token size.

**02** We show how **TARM** increases the stability of the training which enables faster convergence.

# Our Contributions

**01** We introduce a novel token based memory method called **TARM** that increases the capacity of the memory token without actually increasing the memory token size.

**02** We show how **TARM** increases the stability of the training which enables faster convergence.

**03** We evaluate the effectiveness of **TARM** approach on language modelling tasks and demonstrate that it outperforms the original RMT model on similar settings.

# Methodology: TARM

**RMT**
**(alpha = 1)**

# Methodology: TARM

# Methodology: TARM

$$E_t = \alpha \times M_t + (1 - \alpha) \times E_{t-1}$$

where:

$E_t$ = weighted memory tokens at time step $t$

$\alpha$ = weight used in the moving average (damping factor)

$M_t$ = original memory tokens at time step $t$

**TARM**
**(Ours)**

**RMT**
**(alpha = 1)**

# TARM increases stability



Figure 2. Demonstration of how the Exponential Moving Average (EMA) component used in TARM helps in stabilizing the erratic change in values.

# TARM increases stability



Figure 3. Figure show how TARM stabilizes the validation accuracy over steps and also leads to faster/better convergence when compared to RMT.

# Comparison with RMT

| Task | Dataset | Metric | RMT (Reproduced) | TARM (Ours) |
|------|---------|--------|------------------|-------------|
| LM* | WT-103 | PPL ↓ | <u>28.301</u> | **24.493** (↓3.808) |

Table 2. Results comparison of RMT and TARM on Language Modelling task. The best results are shown in bold, and the second-best results are underlined. The symbol ↑ and ↓ indicate the improvement and degradation of the result compared to the reproduced RMT model. LM* denotes language modeling task. PPL denotes perplexity.

# Ablations

| Task | Dataset | Metric | $\alpha$ | PPL |
|------|---------|--------|----------|-----|
| LM* | WT-103 | PPL ↓ | 0.2 | **24.493** |
| LM* | WT-103 | PPL ↓ | 0.7 | <u>24.603</u> |
| LM* | WT-103 | PPL ↓ | 1.0 | 28.301 |

Table 3. Ablation of the hyperparamter $\alpha$. The best results are shown in bold, and the second-best results are underlined. The symbol ↑ and ↓ indicate the improvement and degradation of the result compared to the reproduced RMT model. LM* denotes language modeling task. PPL denotes perplexity.

# References

[1] Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. Recurrent memory transformer. *Advances in Neural Information Processing Systems*, 35:11079–11091, 2022. 1, 2

[2] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 1

[3] Matt Mahoney. Large text compression benchmark, 2011. 1

[4] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016. 1, 2