

BAYESIAN ATTENTION MECHANISM:

A PROBABILISTIC FRAMEWORK
FOR POSITIONAL ENCODING
AND CONTEXT LENGTH EXTRAPOLATION

Arthur Sperhacke Bianchessi
Yasmin Cardozo Aguirre
Rodrigo Coelho Barros
Lucas Silveira Kupssinskü

MALTA

Machine Learning Theory
and Applications Lab



PUCRS

Definitions

Definition 1. For a fixed query vector $\mathbf{q}_i \in \mathcal{R}^{1 \times d}$ and key-value matrices $\mathbf{K}, \mathbf{V} \in \mathcal{R}^{L \times d}$, a Bayesian Attention Mechanism computes self-attention as an expectation over its values:

$$\text{self-attention}(\mathbf{q}_i, \mathbf{K}, \mathbf{V}) = \frac{\exp(\text{score}(\mathbf{q}_i, \mathbf{k}_j))}{\sum_z \exp(\text{score}(\mathbf{q}_i, \mathbf{k}_z))} \mathbf{V} = \sum_j p_{ij} \mathbf{v}_j = \mathbf{E}_{j|i}[\mathbf{V}]$$

Definition 2. In Bayesian Attention, p_{ij} is a joint probability over the *content* of token j (f_{cont}) and its *position* relative to query \mathbf{q}_i (g_{pos}).

$$p_{ij} = p(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_j) | g_{\text{pos}}(i, j)) \times p(g_{\text{pos}}(i, j))$$

Additive Relative Positional Encodings

$$\begin{aligned}
 p_{ij} &= \frac{\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_j) + g_{\text{pos}}(i, j))}{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z) + g_{\text{pos}}(i, z)) \right)} \\
 &= \frac{\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_j)) \cdot \exp(g_{\text{pos}}(i, j))}{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z) + g_{\text{pos}}(i, z)) \right)} \cdot \frac{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z)) \right) \cdot \Sigma_z \left(\exp(g_{\text{pos}}(i, z)) \right)}{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z)) \right) \cdot \Sigma_z \left(\exp(g_{\text{pos}}(i, z)) \right)} = 1 \\
 &= \frac{\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_j))}{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z)) \right)} \cdot \frac{\exp(g_{\text{pos}}(i, j))}{\Sigma_z \left(\exp(g_{\text{pos}}(i, z)) \right)} \cdot \frac{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z)) \right) \cdot \Sigma_z \left(\exp(g_{\text{pos}}(i, z)) \right)}{\Sigma_z \left(\exp(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z) + g_{\text{pos}}(i, z)) \right)} \\
 &= \mathbf{p}(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_j)) \cdot \mathbf{p}(g_{\text{pos}}(i, j)) \cdot \frac{1}{\Sigma_z \left(p(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_z)) \cdot p(g_{\text{pos}}(i, z)) \right)} \\
 &= \frac{\mathbf{p}(f_{\text{cont}}(\mathbf{q}_i, \mathbf{k}_j)) \cdot \mathbf{p}(g_{\text{pos}}(i, j))}{Z}
 \end{aligned}$$

Additive Relative Positional Encodings

$$\text{score}_{\text{ALiBi}}(\mathbf{q}_i, \mathbf{k}_j) = \mathbf{q}_i \cdot \mathbf{k}_j - m \cdot |i - j|$$

$$\text{score}_{\text{NoPE}}(\mathbf{q}_i, \mathbf{k}_j) = \mathbf{q}_i \cdot \mathbf{k}_j + M_{i,j}$$

$$\text{where } M_{i,j} = \begin{cases} 0 & \text{if } j \leq i \\ -\infty & \text{if } j > i \end{cases}$$

Causal Mask (Uniform)

ALiBi (Laplace)

Imagine this is a very long context for the Large Language Model <EOT>

Generalized Gaussian Distribution

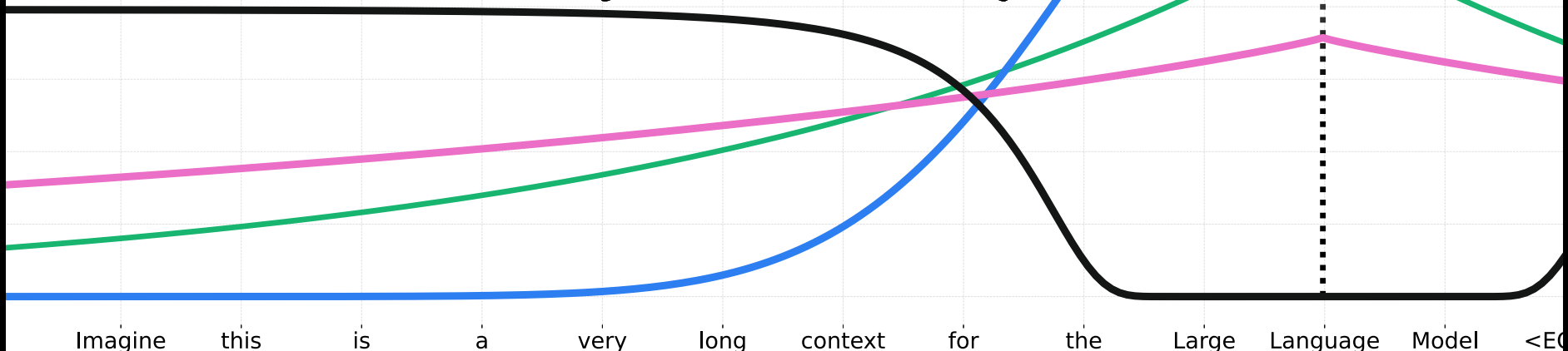
- ALiBi (Laplace)
- GGD-BAM (Gaussian)
- GGD-BAM ($\beta \in (0, 1)$)
- GGD-BAM ($\beta < 0$)

Continuous & Unbounded

$$\frac{\beta}{2\alpha\Gamma(1/\beta)} e^{-\left(\frac{|j-i-\mu|}{\alpha}\right)^\beta}$$

Discrete & Bounded

$$\frac{e^{-\left(\frac{|j-i-\mu|}{\sigma}\right)^\beta}}{\sum_{k=0}^i e^{-\left(\frac{|k-i-\mu|}{\sigma}\right)^\beta}}$$



Imagine this is a very long context for the Large Language Model <EOT>

Generalized Gaussian Distribution

$$\begin{array}{c}
 \mathbf{q}_1 \mathbf{k}_1^\top \quad \mathbf{q}_1 \mathbf{k}_2^\top \quad \mathbf{q}_1 \mathbf{k}_3^\top \quad \mathbf{q}_1 \mathbf{k}_4^\top \quad \mathbf{q}_1 \mathbf{k}_5^\top \\
 \mathbf{q}_2 \mathbf{k}_1^\top \quad \mathbf{q}_2 \mathbf{k}_2^\top \quad \mathbf{q}_2 \mathbf{k}_3^\top \quad \mathbf{q}_2 \mathbf{k}_4^\top \quad \mathbf{q}_2 \mathbf{k}_5^\top \\
 \mathbf{q}_3 \mathbf{k}_1^\top \quad \mathbf{q}_3 \mathbf{k}_2^\top \quad \mathbf{q}_3 \mathbf{k}_3^\top \quad \mathbf{q}_3 \mathbf{k}_4^\top \quad \mathbf{q}_3 \mathbf{k}_5^\top \\
 \mathbf{q}_4 \mathbf{k}_1^\top \quad \mathbf{q}_4 \mathbf{k}_2^\top \quad \mathbf{q}_4 \mathbf{k}_3^\top \quad \mathbf{q}_4 \mathbf{k}_4^\top \quad \mathbf{q}_5 \mathbf{k}_5^\top \\
 \mathbf{q}_5 \mathbf{k}_1^\top \quad \mathbf{q}_5 \mathbf{k}_2^\top \quad \mathbf{q}_5 \mathbf{k}_3^\top \quad \mathbf{q}_5 \mathbf{k}_4^\top \quad \mathbf{q}_5 \mathbf{k}_5^\top
 \end{array}
 \quad - \quad
 \frac{1}{\alpha} \times
 \begin{array}{c}
 \left[\begin{array}{ccccc}
 0 & 1 & 2 & 3 & 4 \\
 -1 & 0 & 1 & 2 & 3 \\
 -2 & -1 & 0 & 1 & 2 \\
 -3 & -2 & -1 & 0 & 1 \\
 -4 & -3 & -2 & -1 & 0
 \end{array} \right]
 \begin{array}{c}
 \beta \\
 -\mu
 \end{array}
 \end{array}$$

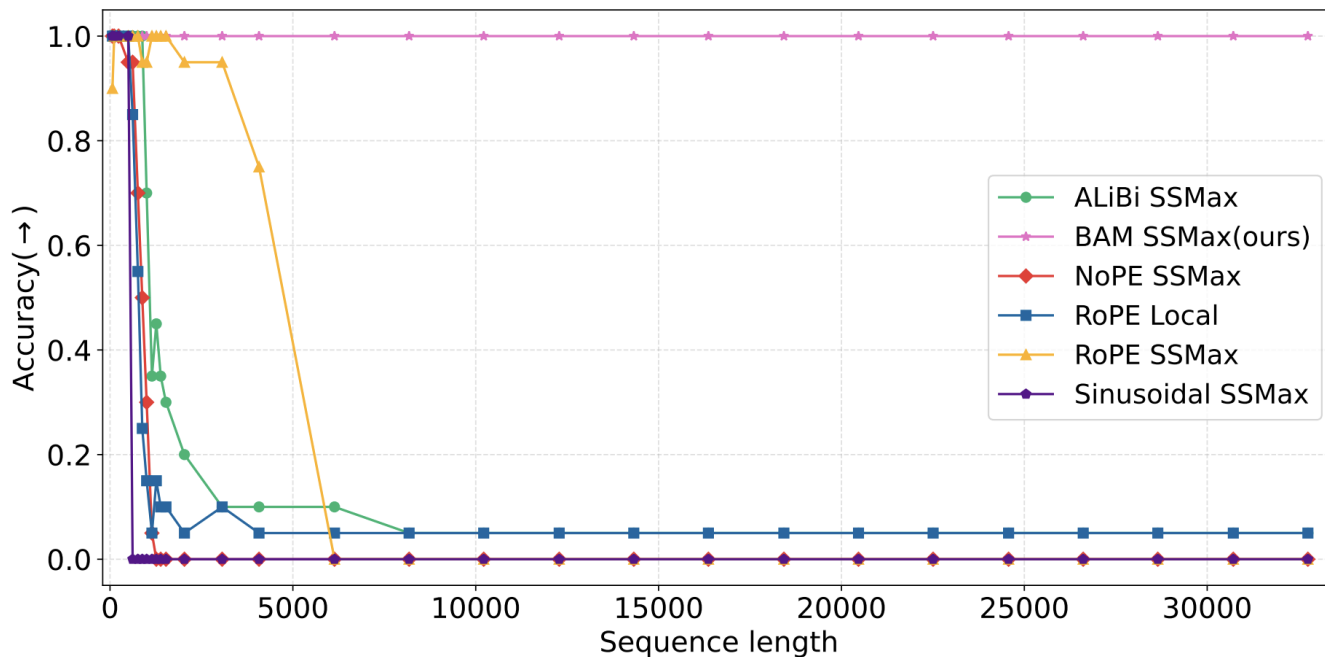
\mathbf{QK}^\top
 \mathbf{B}

$$\mu = 0 \qquad \alpha = e^{\theta_\alpha} \qquad \beta = \theta_\beta$$

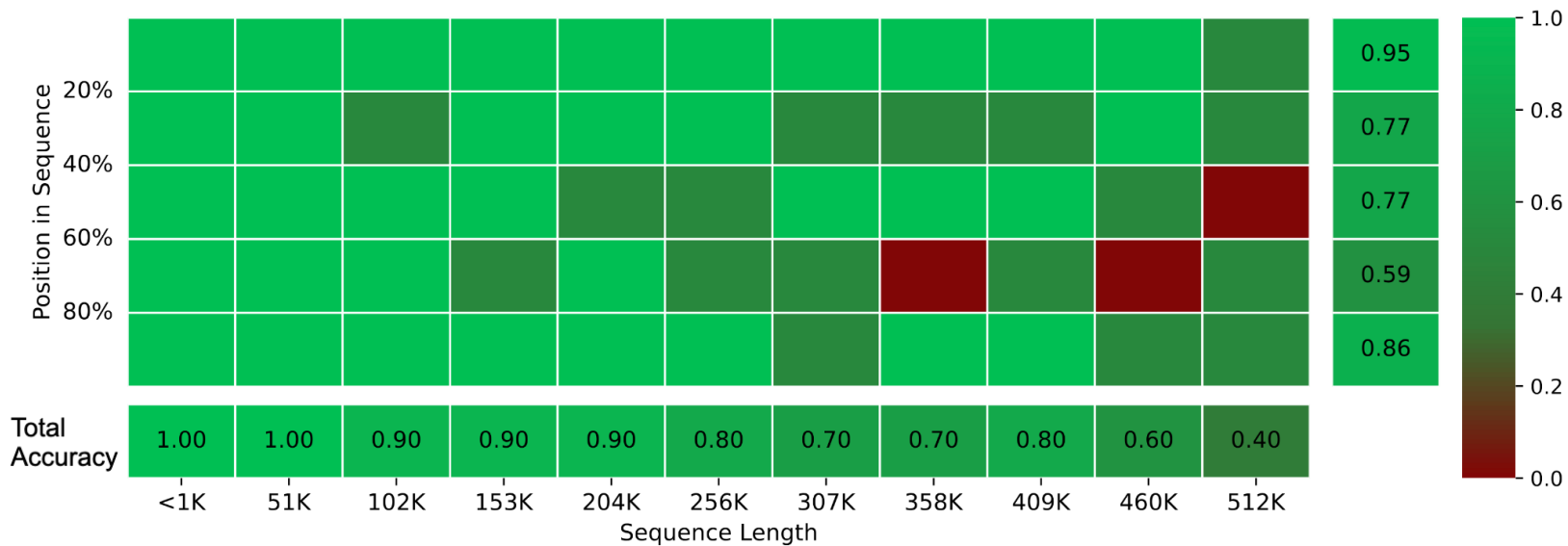
Evaluation I:

Evaluating Extrapolation

Passkey Retrieval

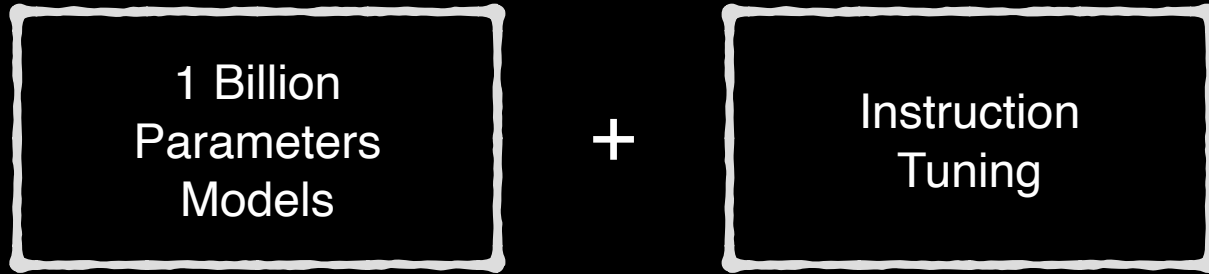


Passkey Retrieval (GGD-BAM)



Evaluation II:

Evaluating Extrapolation & Performance



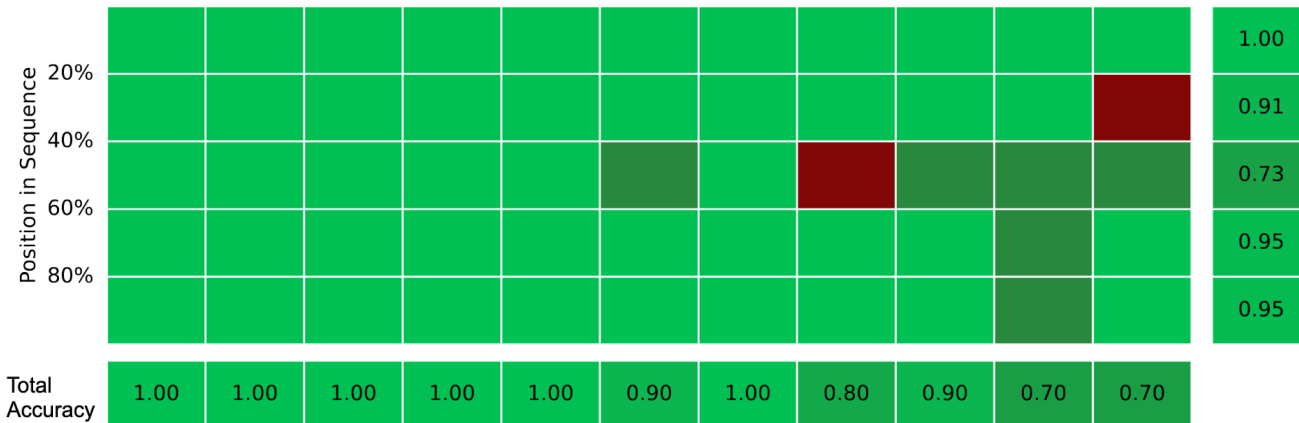
Passkey Retrieval

1 Billion
Parameter
Models

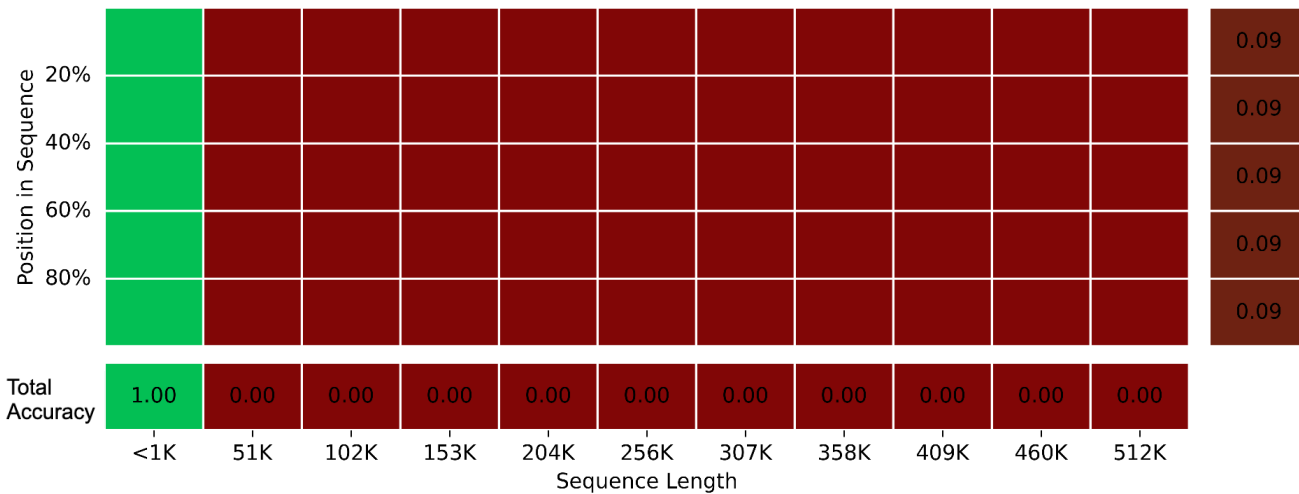
+

Instruction
Tuning

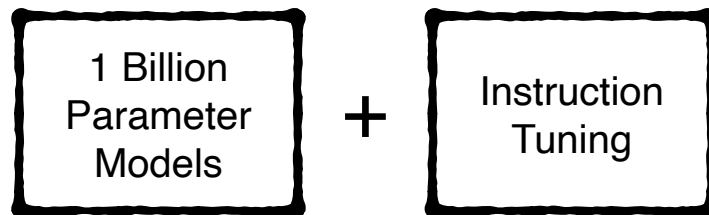
GGD - BAM



RoPE



Performance Benchmarks



| | MMLU | ARC-Easy | ARC-Challenge |
|--|------|----------|----------------------|
|--|------|----------|----------------------|

| | | | |
|------------------------|---------------|---------------|---------------|
| BAM SSM _{Max} | 0.3716 | 0.5770 | 0.4132 |
|------------------------|---------------|---------------|---------------|

| | | | |
|-------------------------|--------|--------|--------|
| RoPE SSM _{Max} | 0.3573 | 0.5715 | 0.4123 |
|-------------------------|--------|--------|--------|

Performance on Long Context

1 Billion
Parameter
Models

+

Instruction
Tuning

Table 6: LongBenchv2 Benchmark: GGD-BAM vs RoPE, 1B parameter models.

| | BAM SSMaX | RoPE SSMaX |
|-------------------------------------|-------------|------------|
| Code Repository Understanding | 41.7 | 25.0 |
| Long In-context Learning | 36.4 | 30.3 |
| Long-dialogue History Understanding | 35.0 | 35.0 |
| Multi-Document QA | 26.5 | 25.3 |
| Single-Document QA | 26.5 | 18.8 |
| Overall | 28.6 | 24.2 |

BAYESIAN ATTENTION MECHANISM:

A PROBABILISTIC FRAMEWORK
FOR POSITIONAL ENCODING
AND CONTEXT LENGTH EXTRAPOLATION

