# Manticore

CS769 Natural Language Processing
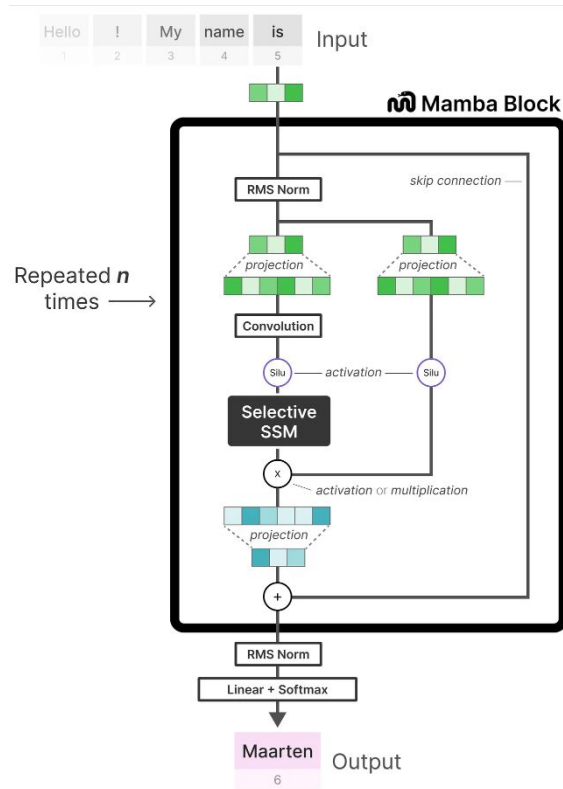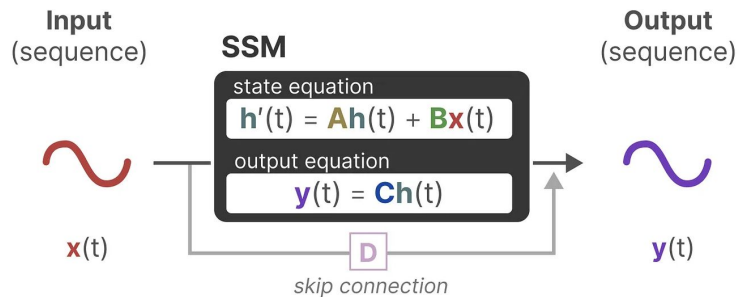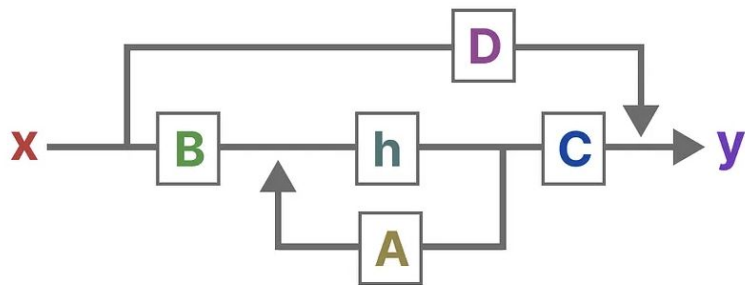
# Topics

- Motivation
- Manticore framework
- Implementation
  - Results on IMDB
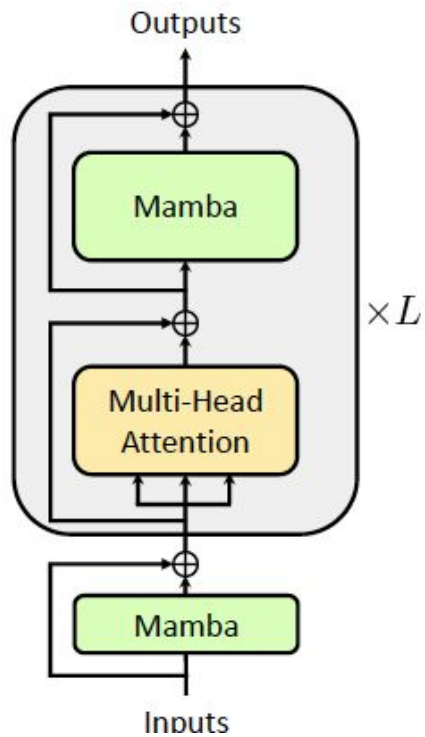  - Results on MAD tasks

# Motivation

- Transformer models have dominated the NLP arena
  - But computational cost is quadratic in sequence length (attention layer)
  - **Transformer dominance:** Challenged by new subquadratic architectures (e.g. State Space Models, Mamba)
- State Space models
  - Mamba
- Hybrid architectures: Combine strengths of diverse models but suffer from:
  - **Manual design:** Inefficient and heuristic-driven
  - **Pre-training challenges:** Limited ability to reuse pretrained components (Neural architecture search)
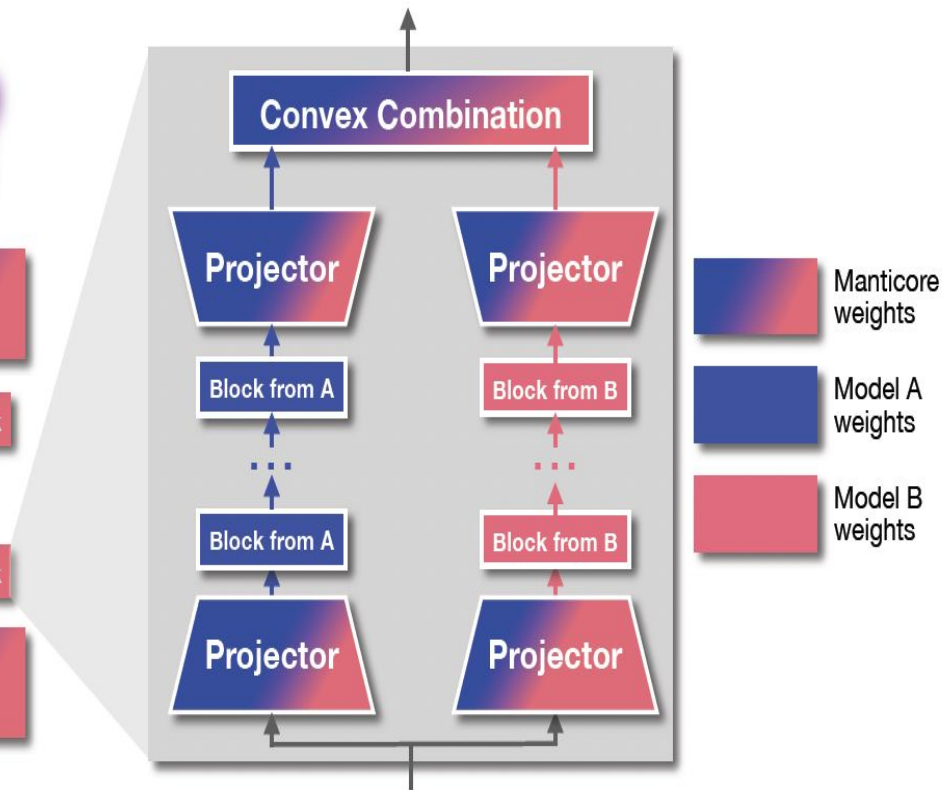
# Mamba



**SSM**

state equation
$$h'(t) = \mathbf{A}h(t) + \mathbf{B}\mathbf{x}(t)$$

output equation
$$\mathbf{y}(t) = \mathbf{C}h(t)$$

Input (sequence) — $\mathbf{x}(t)$

Output (sequence) — $\mathbf{y}(t)$

*skip connection*

Input — Hello ! My name is — 1 2 3 4 5

**Mamba Block**

RMS Norm — *skip connection*

*projection* — *projection*

Convolution

Silu — *activation* — Silu

**Selective SSM**

$\times$ — *activation* or *multiplication*

*projection*

$+$

Repeated **n** times

RMS Norm

Linear + Softmax

Output — Maarten 6
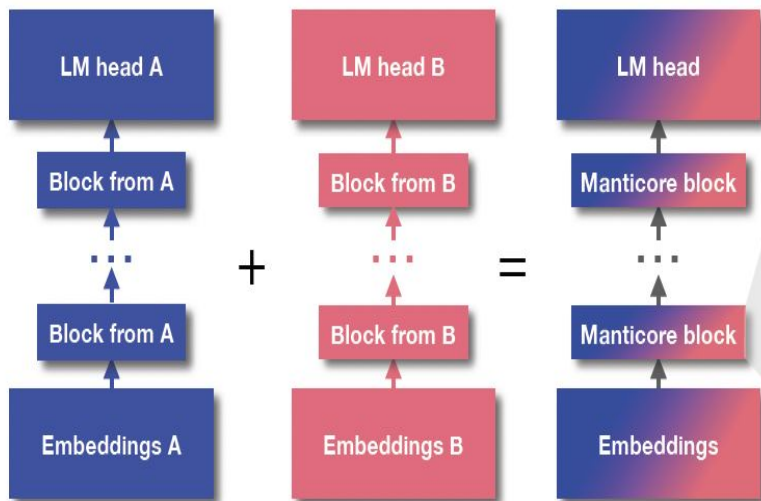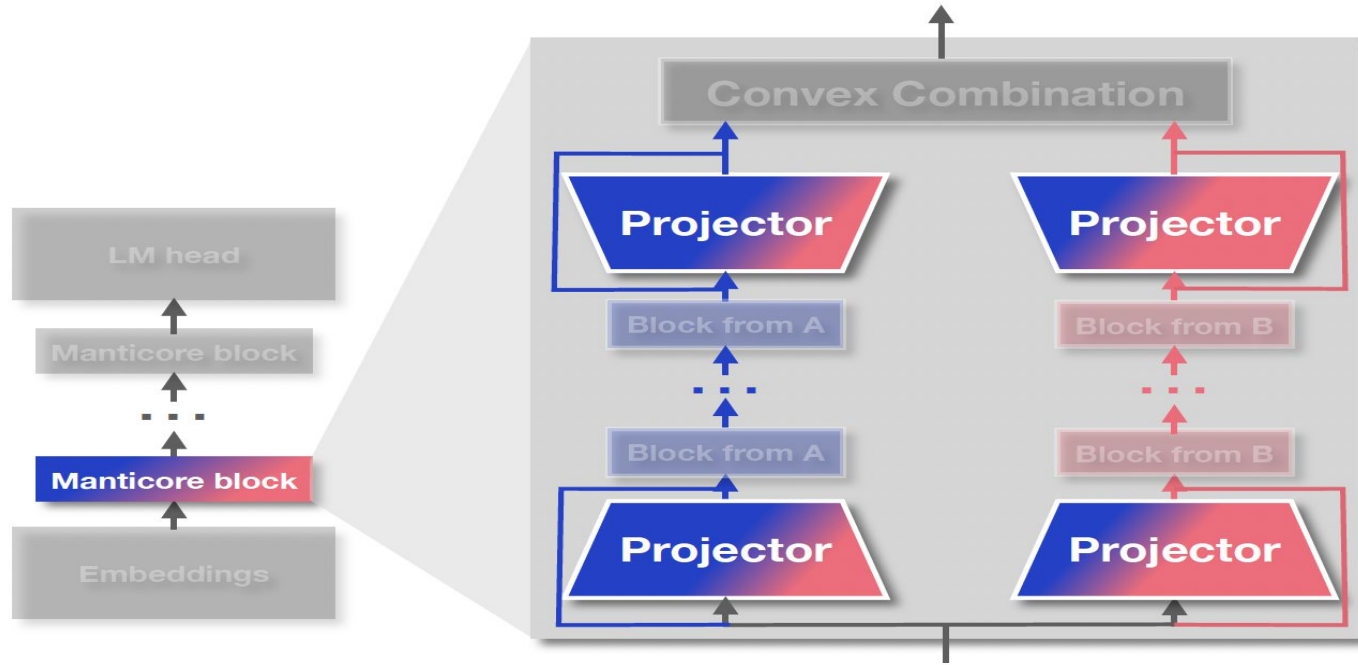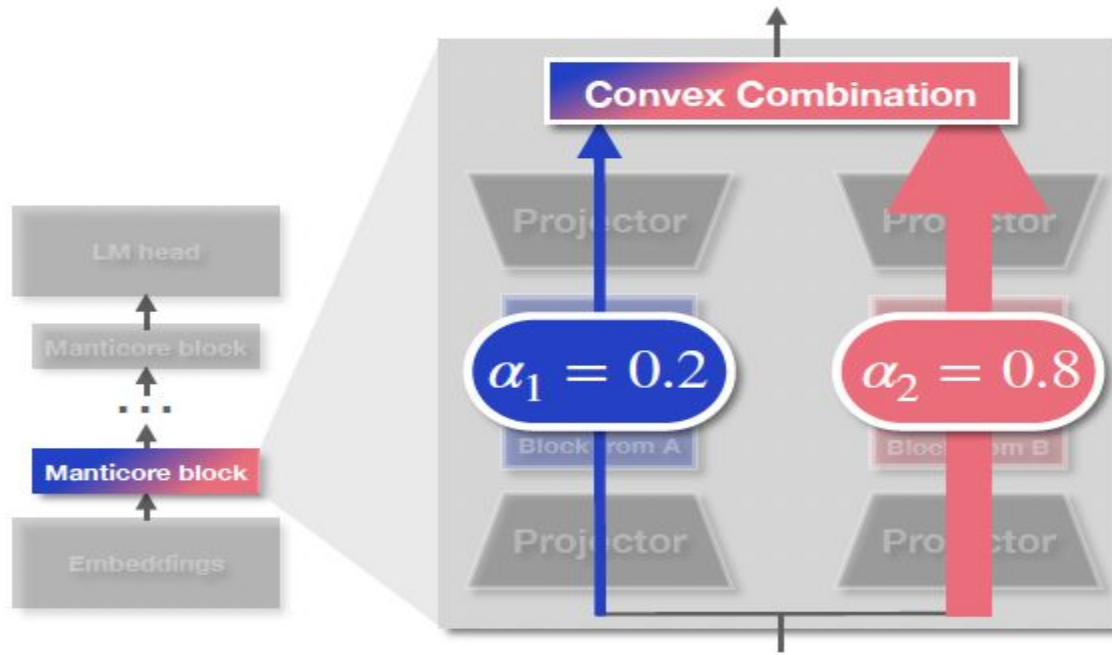
# MambaFormer

# Manticore - Projectors

# Manticore - Convex combination

# Manticore - Equations

$$\text{Proj-in}_{(k)}^{(\ell)}(x; \alpha) := (1 - \alpha) \cdot \text{Linear}_{d_{\max} \to d_{M_{(k)}}}(x) + \alpha \cdot \text{Trunc}(x; d_{M_{(k)}})$$

$$\text{Proj-out}_{(k)}^{(\ell)}(x; \alpha) := (1 - \alpha) \cdot \text{Linear}_{d_{M_{(k)}} \to d_{\max}}(x) + \alpha \cdot \text{Pad}(x; d_{\max}).$$

$$h_k(x; \alpha_k, I_k) = \left( \text{Proj-out}_{(k)}^{(I_{k,n_k})} \circ M_{(k)\text{Block}}^{(I_{k,n_k})} \circ ... \circ M_{(k)\text{Block}}^{(I_{k,1})} \circ \text{Proj-in}_{(k)}^{(I_{k,1})} \right)(x; \alpha_k).$$

$$\text{Mix}_\alpha(x; I_1, ..., I_K) = \sum_{k \in [K]} \alpha_k h_k(x; \alpha_k, I_k).$$

# Implementation Details

Created our own codebase from scratch with the following modular components:

- Gated Combiners (with convex combination)

- Gated Splitters

- Hybrid Model class (with task-specific output heads)

- Trainer

Next, trying to utilize DeepSpeed's ZeRO optimization or other model parallelism approaches to work with larger model backbones

# IMDB Results

With 3 epochs, 2500 iterations/ epoch, we got the following dev accuracies:

| Transformer | Mamba | Hybrid |
|---|---|---|
| 58.7 | 89.5 | 84.8 (incomplete training) |

# MAD Tasks

- Selective Copying

  input example: a c [b] t [b] [i] [i] [b] [i] | a c [b] t [b] <u>a</u> <u>c</u> [b] <u>t</u>

- In-context Recall

  input example: a b d e f g | a <u>b</u> f <u>g</u>

- Noisy In-context Recall

  input example: a b h d e f g | i a <u>b</u> f <u>g</u>

- Fuzzy In-context Recall

  input example: (a d) (b) (d a f) (e g) | (d a f) (<u>e</u> <u>g</u>)

- Memorization

  key-value dictionary example: {a:b, c:d, e:f}     input example: a [i] c [i] e [i] a [i] | a <u>b</u> c <u>d</u> e <u>f</u> a <u>b</u>

# MAD Task Config

- **MAD config**
  - Number of layers: 2
  - Vocab size: 16
  - Seq len: 32
  - Train examples: 4096
  - Test examples: 256
  - Hidden size: 64
  - Number of Transformer Heads: 4
  - Noise fraction: 0.2
- **Training**
  - Epochs: 20
  - Batch size: 32
  - Learning rate: 5e-4

# MAD Tasks Results

Table: Best Test Accuracy on MAD Tasks

| Task | Transformer | Mamba | Hybrid 1-blk | Hybrid 2-blk | MambaFormer |
|---|---|---|---|---|---|
| Selective Copying | 0.91 | 1.0 | 0.99 | 0.99 | 0.98 |
| In-context Recall | 0.35 | 1.0 | 1.0 | 1.0 | 0.99 |
| Noisy In-context Recall | 0.42 | 1.0 | 1.0 | 0.99 | 1.0 |
| Fuzzy In-context Recall | 0.35 | 0.36 | 0.42 | 0.43 | 0.15 |
| Memorization | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |

# References

- Albert Gu, Karan Goel, and Christopher Ré. 2021. Efficiently modeling long sequences with structured state spaces. CoRR, abs/2111.00396.
- Albert Gu and Tri Dao. 2024. Mamba: Linear-time sequence modeling with selective state spaces.
- Jongho Park, et al. 2024. Can mamba learn how to learn? a comparative study on in-context learning tasks.
- Michael Poli,et. al. 2024. Mechanistic design and scaling of hybrid architectures.
- MohammadReza Davari and Eugene Belilovsky. 2024. Model breadcrumbs: Scaling multi-task model merging with sparse masks.
- Nicholas Roberts, et. al. 2024. Pretrained hybrids with mad skills.
- Takuya Akiba, et. al. Evolutionary optimization of model merging recipes.