

# Welcome Falcon Mamba: First Strong Attention-free 7B Language model

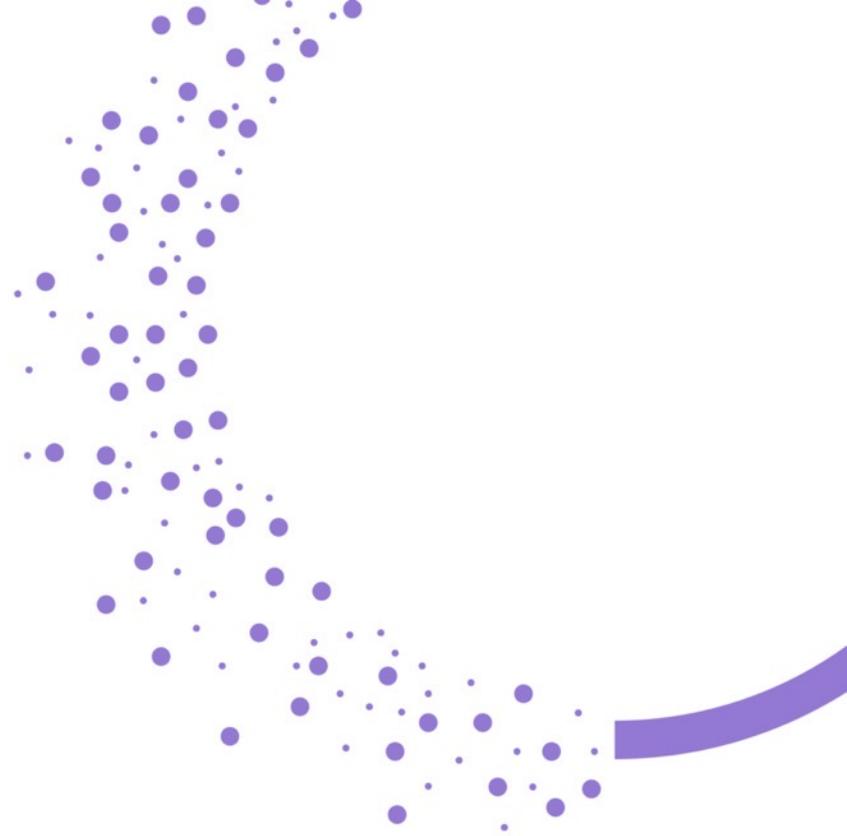
**Jingwei Zuo**

Lead Researcher @Falcon LLM team, Technology Innovation Institute

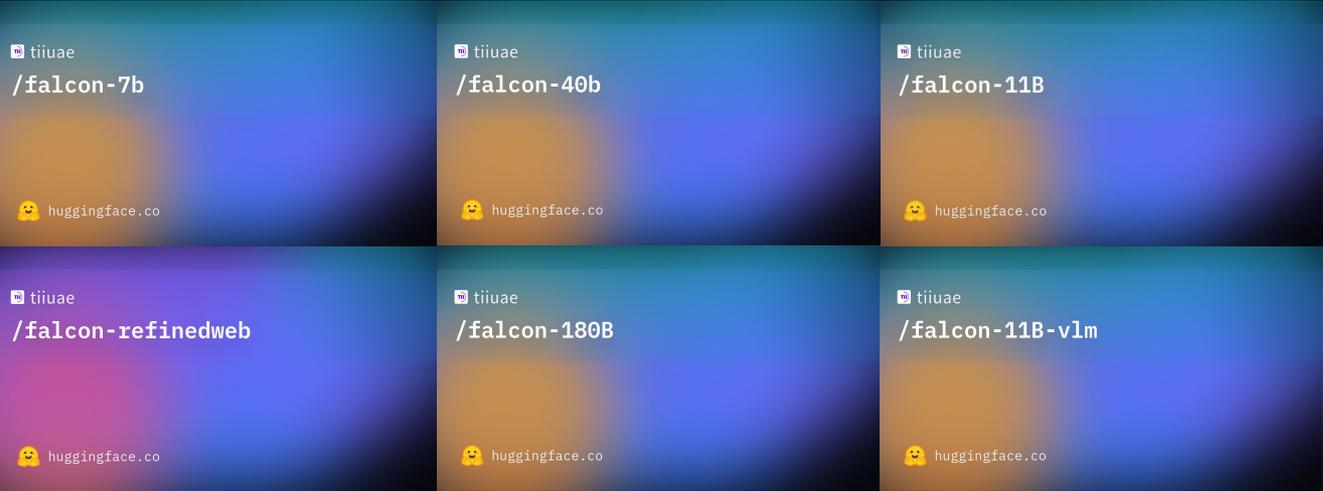
19/12/2024, TotalEnergies Digital Factory

## Outline

- Context
- Falcon Mamba 7B
- How to use Falcon Mamba 7B?
- More...

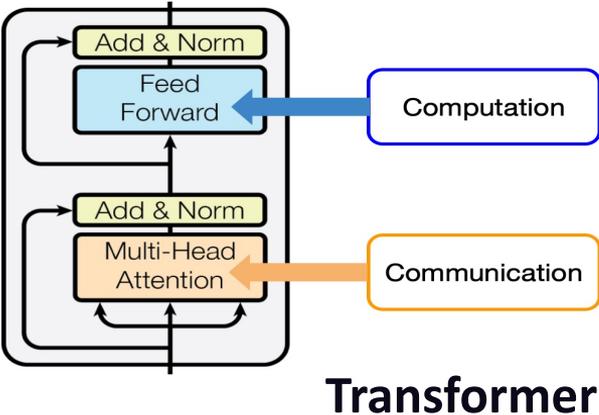


# Falcon Model Series

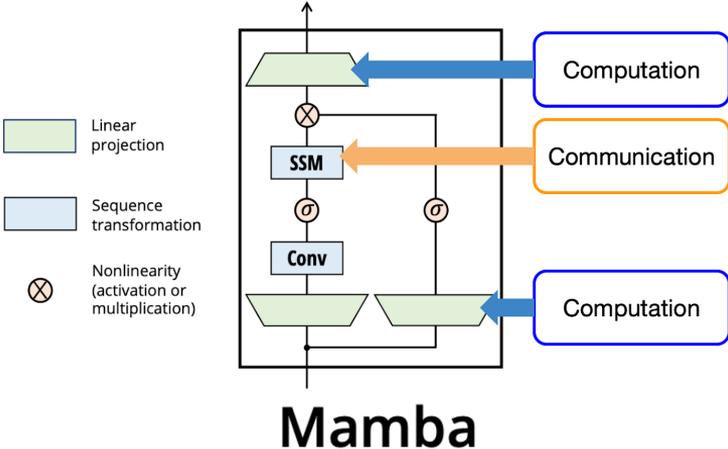


# Foundation Model Backbones (e.g., LLMs)\*

- **Communication** *between* tokens
- **Computation** *within* a token



**Transformer**



**Mamba**

\*<https://www.kolaayonrinde.com/blog/2024/02/11/mamba.html>

# From State Space Models (SSMs) to Mamba

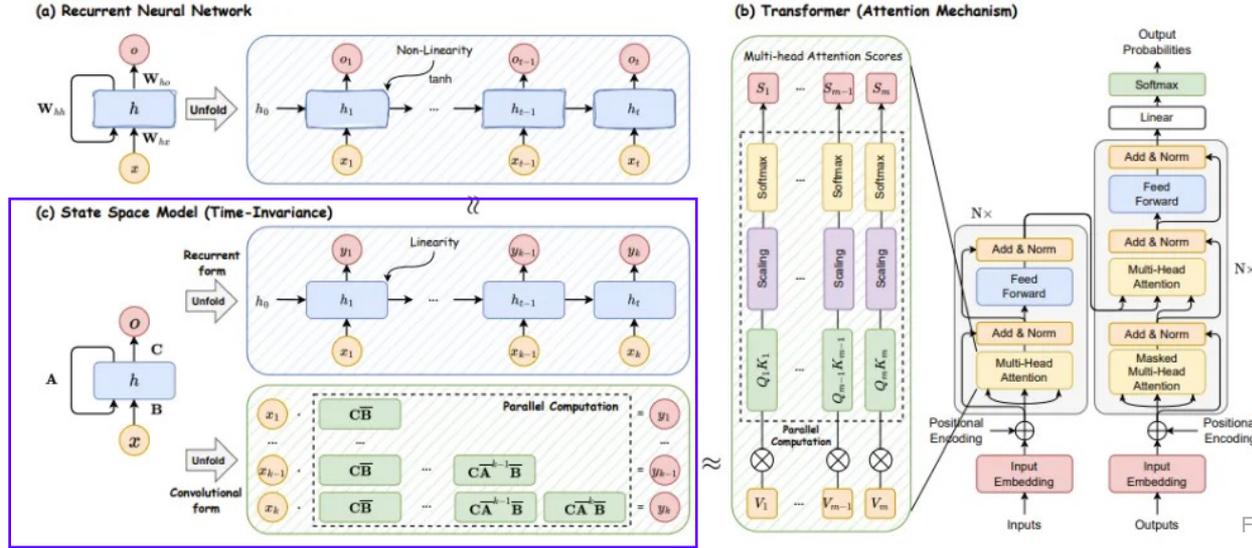
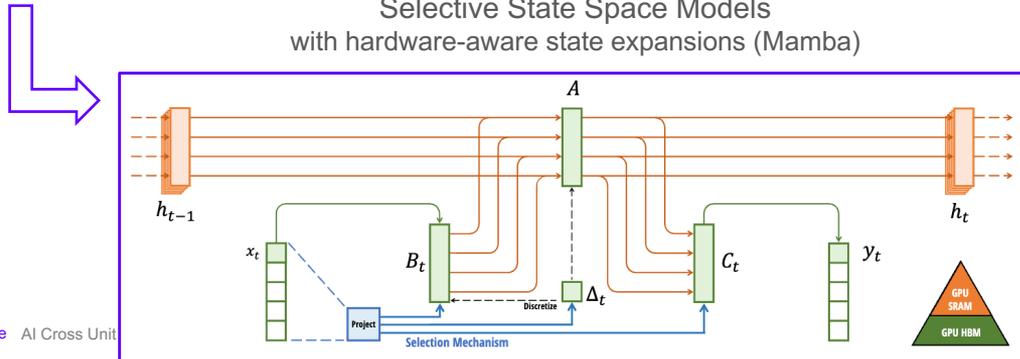


Fig. credit [Qu et al., arXiv'24]

## Selective State Space Models with hardware-aware state expansions (Mamba)



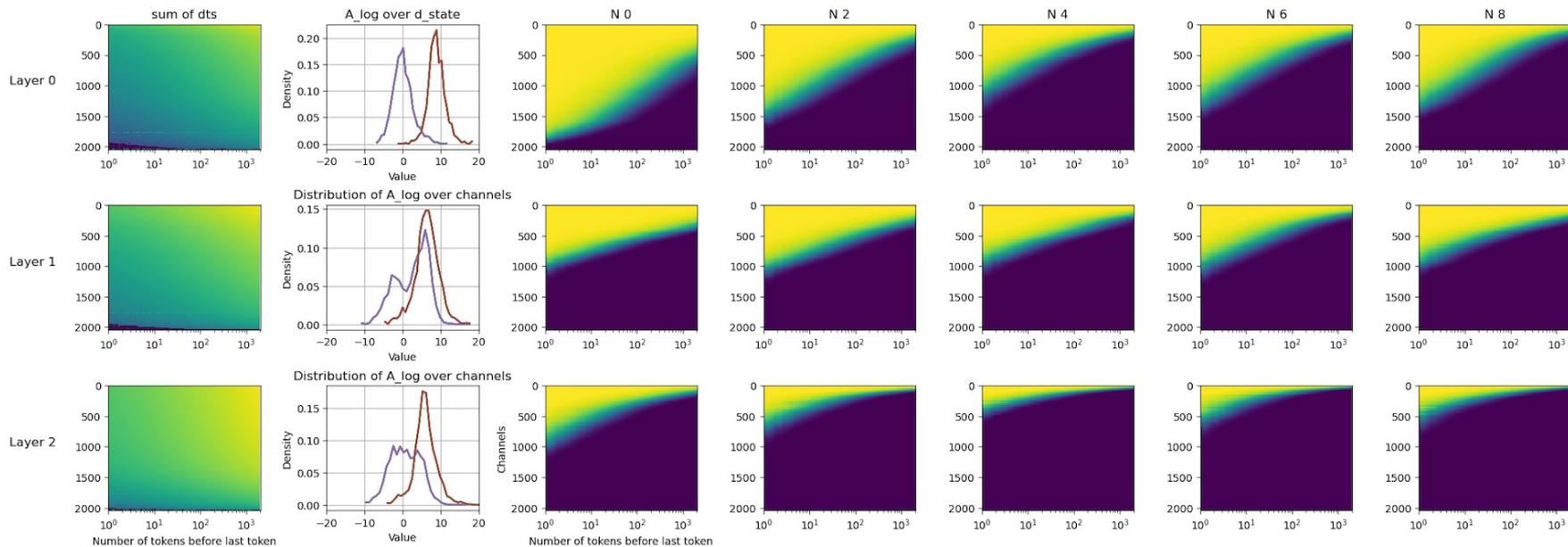
$$h_t = \bar{A}h_{t-1} + \bar{B}x_t$$

$$y_t = Ch_t$$

- Select to remember/forget
- Read/write from/to the hidden states

Fig. credit [Gu and Dao, arXiv'23]

# Mamba Memorization & Forgetting



# Mamba – By-Design Advantages

## Transformers

- Quadratic time complexity:  $O(n^2)$  regarding context length  $n$
- Memory Constraints:  $O(n)$  key-value (KV) cache for each token
- Optimizations: Sparse or local attention (e.g., SWA), CUDA optimizations (e.g., Flash Attention), etc.

## Mamba (attention-free)

- Reduced time complexity:  $O(n)$  regarding context length  $n$
- Memory Optimization: **constant** memory cost
- Long-context dependencies: **compact** state representation

# Mamba applications in a nutshell

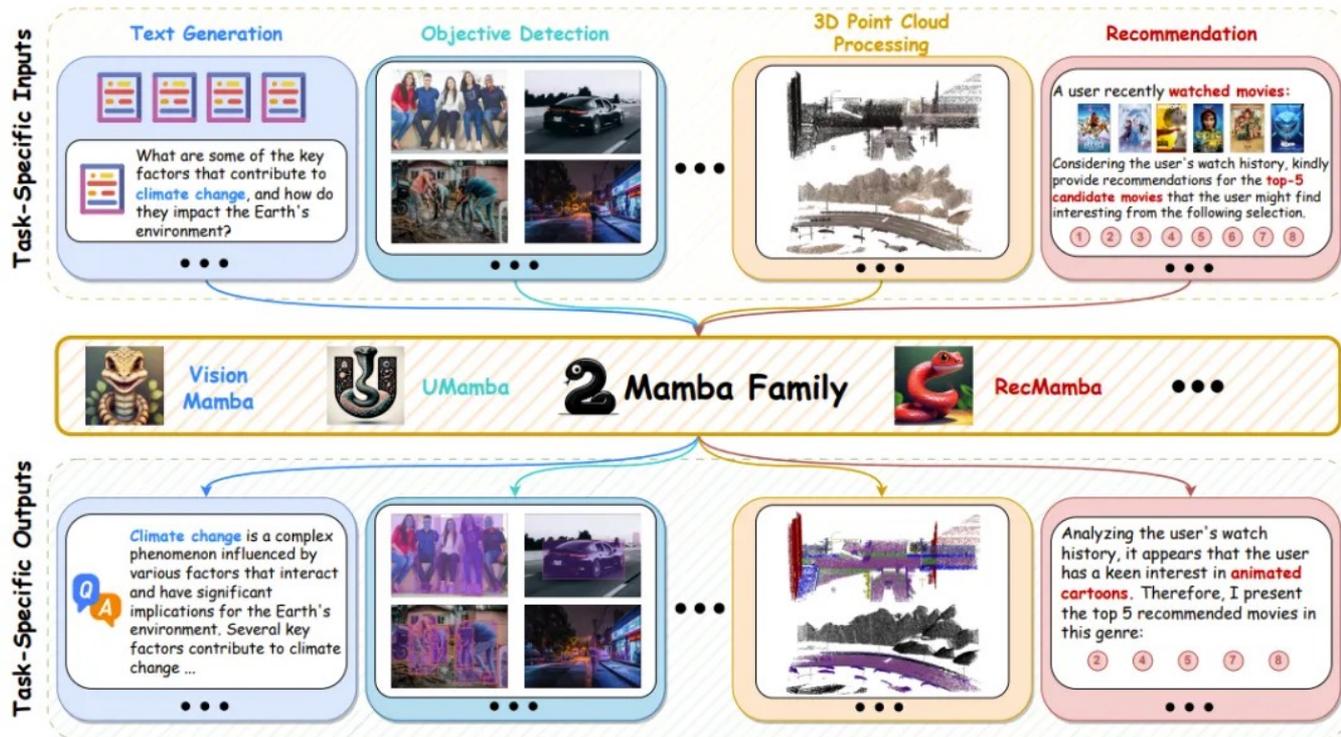


Fig. credit - [Qu et al., arXiv'24]

# Mamba LLMs below 2.8B

Model	Token.	Pile ppl ↓	LAMBADA ppl ↓	LAMBADA acc ↑	HellaSwag acc ↑	PIQA acc ↑	Arc-E acc ↑	Arc-C acc ↑	WinoGrande acc ↑	Average acc ↑
Hybrid H3-130M	GPT2	—	89.48	25.77	31.7	64.2	44.4	24.2	50.6	40.1
Pythia-160M	NeoX	29.64	38.10	33.0	30.2	61.4	43.2	24.1	<b>51.9</b>	40.6
<b>Mamba-130M</b>	NeoX	<b>10.56</b>	<b>16.07</b>	<b>44.3</b>	<b>35.3</b>	<b>64.5</b>	<b>48.0</b>	<b>24.3</b>	<b>51.9</b>	<b>44.7</b>
Hybrid H3-360M	GPT2	—	12.58	48.0	41.5	68.1	51.4	24.7	54.1	48.0
Pythia-410M	NeoX	9.95	10.84	51.4	40.6	66.9	52.1	24.6	53.8	48.2
<b>Mamba-370M</b>	NeoX	<b>8.28</b>	<b>8.14</b>	<b>55.6</b>	<b>46.5</b>	<b>69.5</b>	<b>55.1</b>	<b>28.0</b>	<b>55.3</b>	<b>50.0</b>
Pythia-1B	NeoX	7.82	7.92	56.1	47.2	70.7	57.0	27.1	53.5	51.9
<b>Mamba-790M</b>	NeoX	<b>7.33</b>	<b>6.02</b>	<b>62.7</b>	<b>55.1</b>	<b>72.1</b>	<b>61.2</b>	<b>29.5</b>	<b>56.1</b>	<b>57.1</b>
GPT-Neo 1.3B	GPT2	—	7.50	57.2	48.9	71.1	56.2	25.9	54.9	52.4
Hybrid H3-1.3B	GPT2	—	11.25	49.6	52.6	71.3	59.2	28.1	56.9	53.0
OPT-1.3B	OPT	—	6.64	58.0	53.7	72.4	56.7	29.6	59.5	55.0
Pythia-1.4B	NeoX	7.51	6.08	61.7	52.1	71.0	60.5	28.5	57.2	55.2
RWKV-1.5B	NeoX	7.70	7.04	56.4	52.5	72.4	60.5	29.4	54.6	54.3
<b>Mamba-1.4B</b>	NeoX	<b>6.80</b>	<b>5.04</b>	<b>64.9</b>	<b>59.1</b>	<b>74.2</b>	<b>65.5</b>	<b>32.8</b>	<b>61.5</b>	<b>59.7</b>
GPT-Neo 2.7B	GPT2	—	5.63	62.2	55.8	72.1	61.1	30.2	57.6	56.5
Hybrid H3-2.7B	GPT2	—	7.92	55.7	59.7	73.3	65.6	32.3	61.4	58.0
OPT-2.7B	OPT	—	5.12	63.6	60.6	74.8	60.8	31.3	61.0	58.7
Pythia-2.8B	NeoX	6.73	5.04	64.7	59.3	74.0	64.1	32.9	59.7	59.1
RWKV-3B	NeoX	7.00	5.24	63.9	59.6	73.7	67.8	33.1	59.6	59.6
<b>Mamba-2.8B</b>	NeoX	<b>6.22</b>	<b>4.23</b>	<b>69.2</b>	<b>66.1</b>	<b>75.2</b>	<b>69.7</b>	<b>36.3</b>	<b>63.5</b>	<b>63.3</b>

Fig. credit [Gu and Dao, arXiv'23]

# Mamba LLMs beyond 2.8B

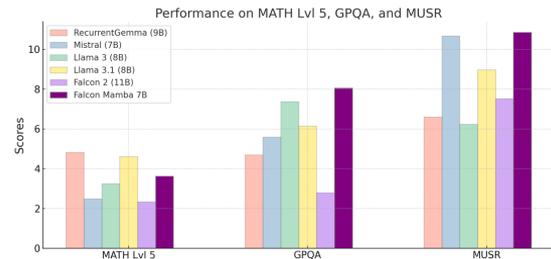
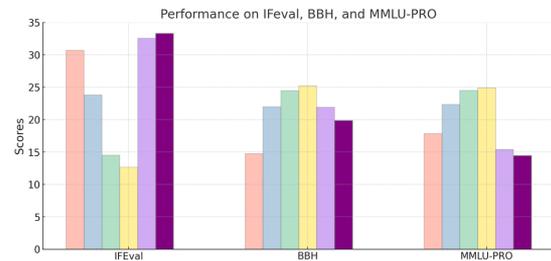
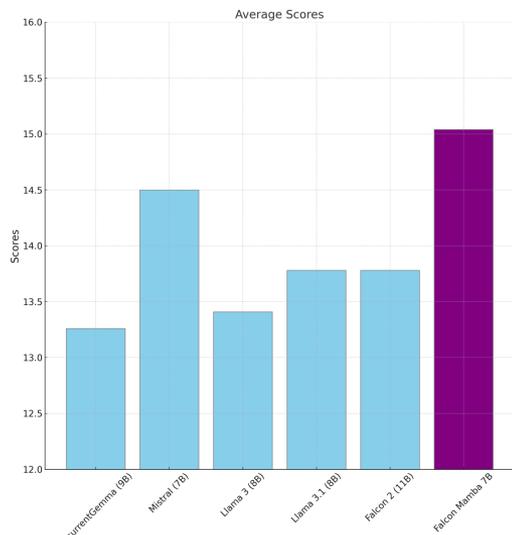


# Falcon Mamba 7B

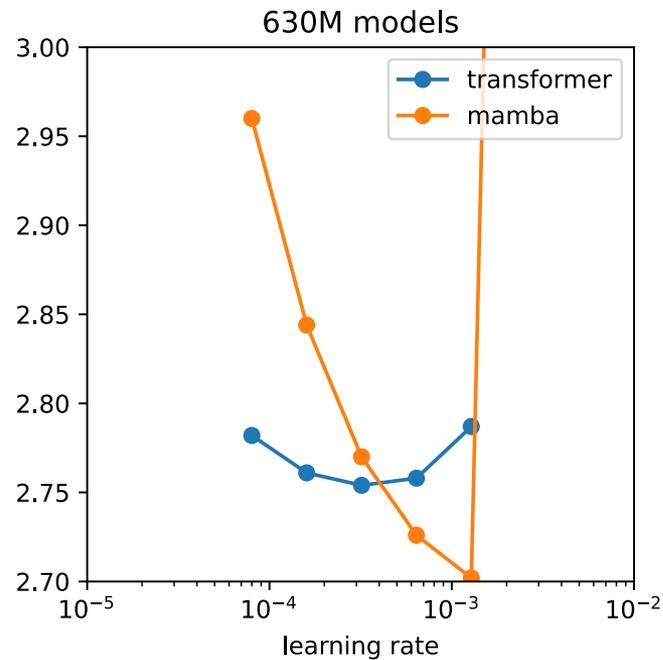
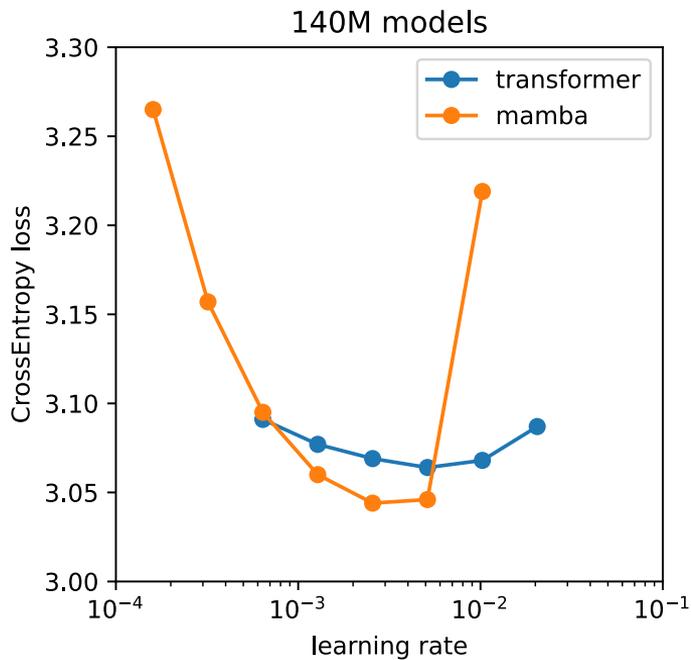
Welcome  
FalconMamba 7B



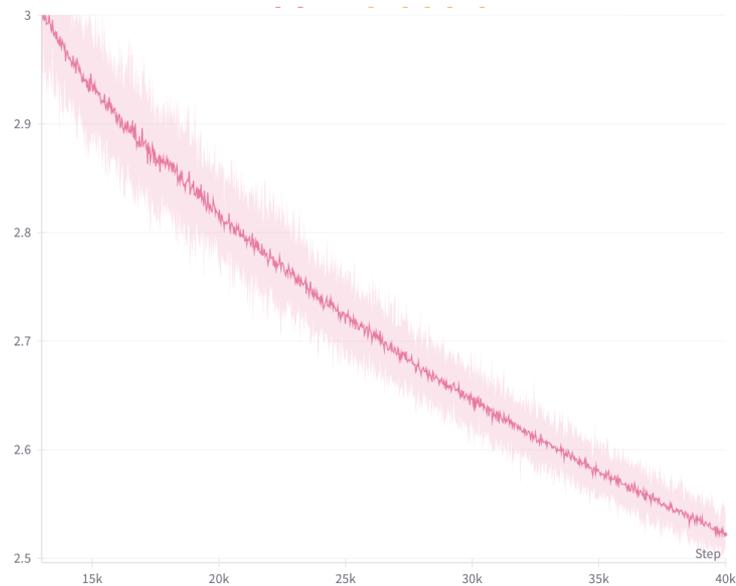
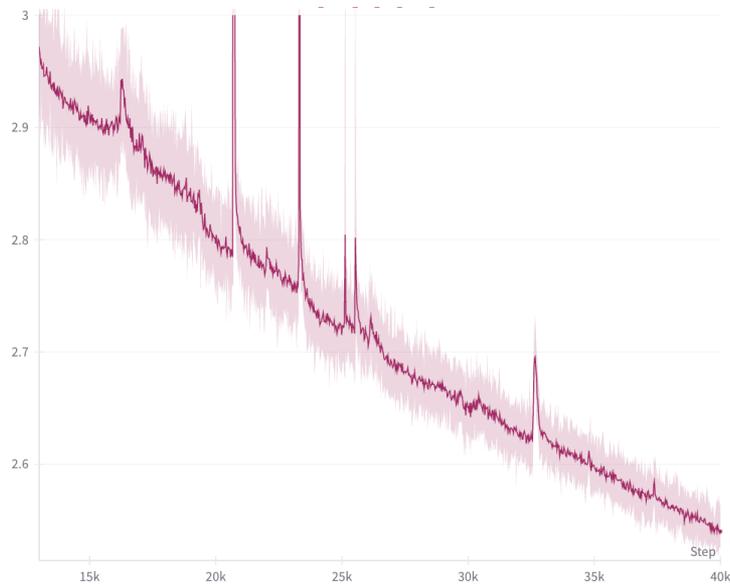
- Best open-source Mamba model at its scale
- Outperforms top Transformer models: Llama 3.1 8B, Mistral 7B, Falcon2 11B
- Supports infinite context length (sequential prefill)



# Model's Sensitivity to Learning Rate



# Stabilizing Mamba LLM's Pre-Training



# Falcon Mamba 7B - Training Details

## Training Infrastructure

Gigatron: In-house 5D parallel training framework

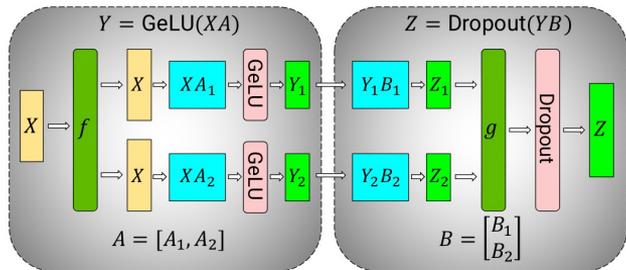
- 5D Parallelism: Data, Tensor, Pipeline, Sequence, Context
- Support both Transformer and Mamba architectures

Trained on 256 H100 80GB GPUs (~ two months)

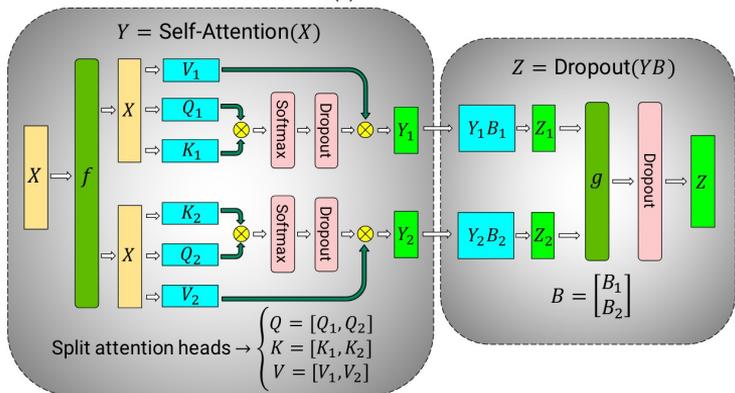
- Only Data Parallelism was applied for Falcon Mamba 7B

# Falcon Mamba 7B - Training Details

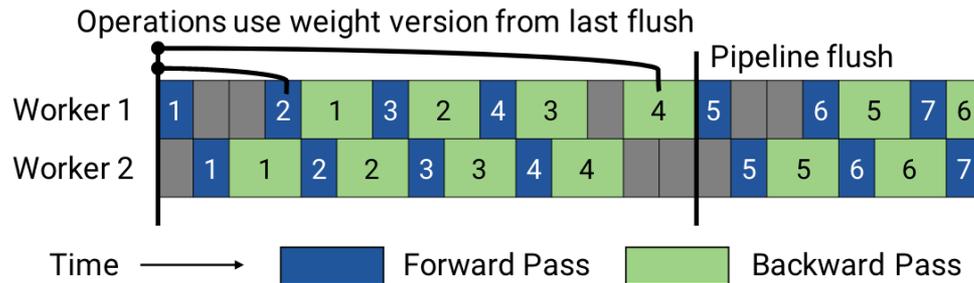
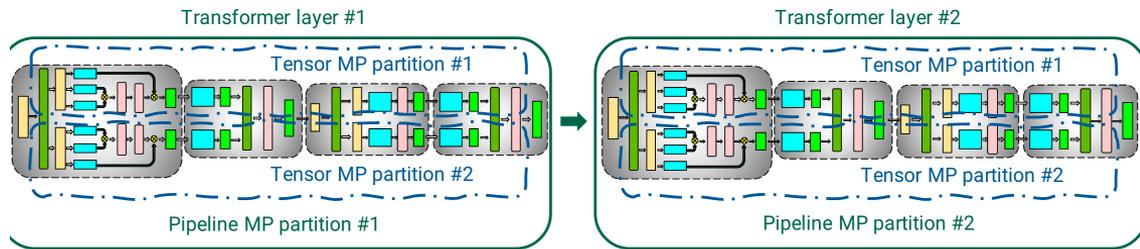
## Parallel Training: from Transformer to Mamba



(a) MLP.

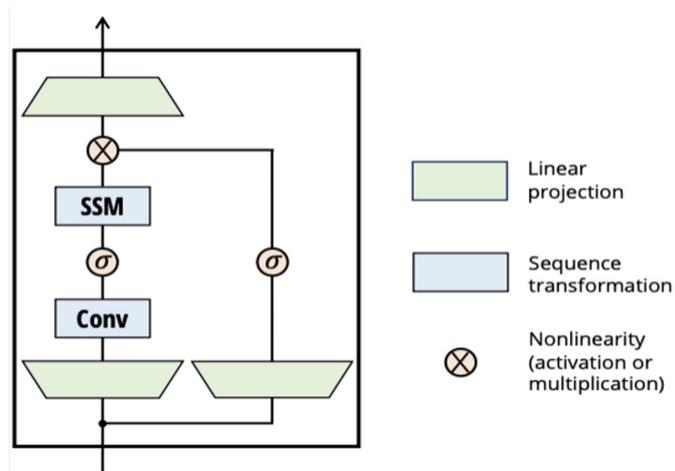
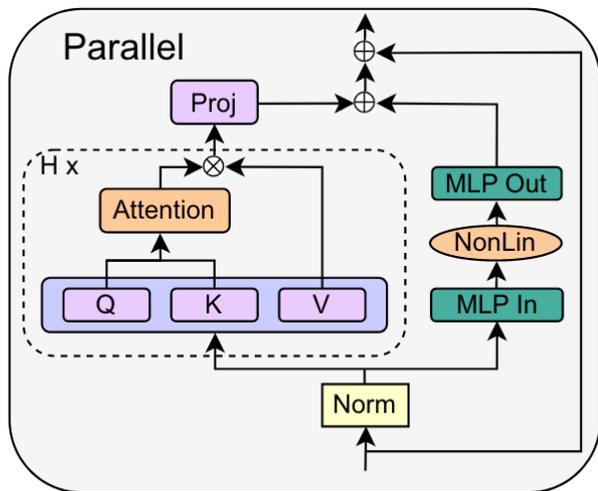


(b) Self-Attention.



# Falcon Mamba 7B - Training Details

## Parallel Training: from Transformer to Mamba



# Falcon Mamba 7B - Training Details

## Training Recipes & Procedure

- WSD (warmup-stable-decay) learning rate schedule
  - Flexible pretraining & continual training
- Batch size ramp-up, batch scaling (during ramp-up), etc.
- RMSNorms within Mamba blocks

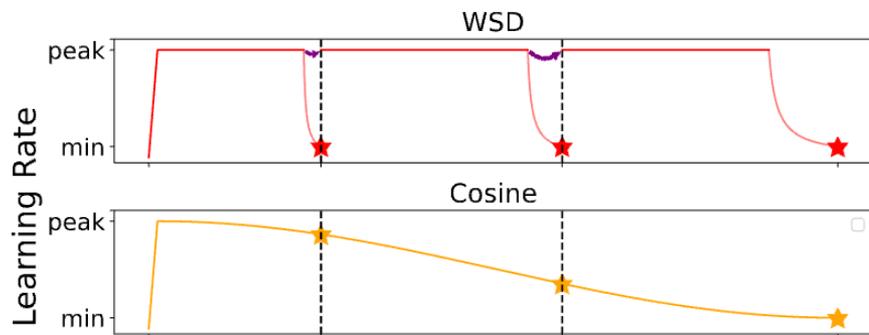


Fig. credit - [Wen et al., arXiv'24]

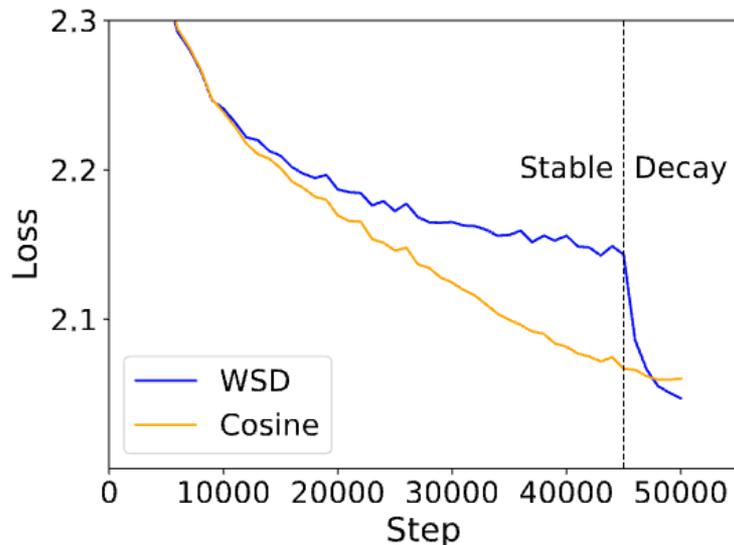
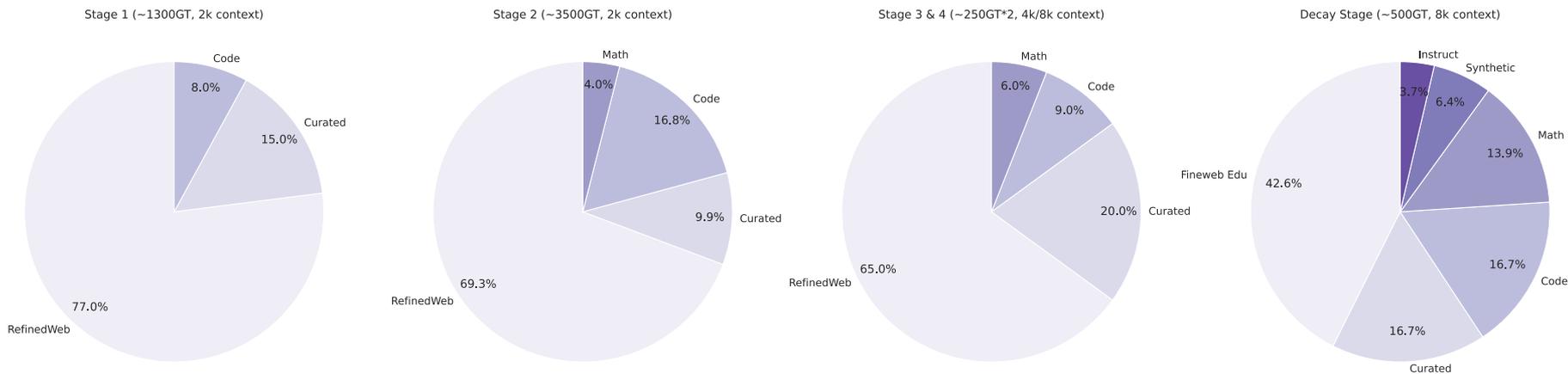


Fig. credit - [Wen et al., arXiv'24]

# Falcon Mamba 7B - Training Details

## Pretraining Data

- Mostly from Falcon 2 11B
- ~**5,800GT** from multiple resources, e.g., RefinedWeb, code, math, etc.



# Falcon Mamba 7B VS SoTA LLMs

Table 2: Model Performance on HF Leaderboard v1 tasks: **bold** (best), underline (second best)

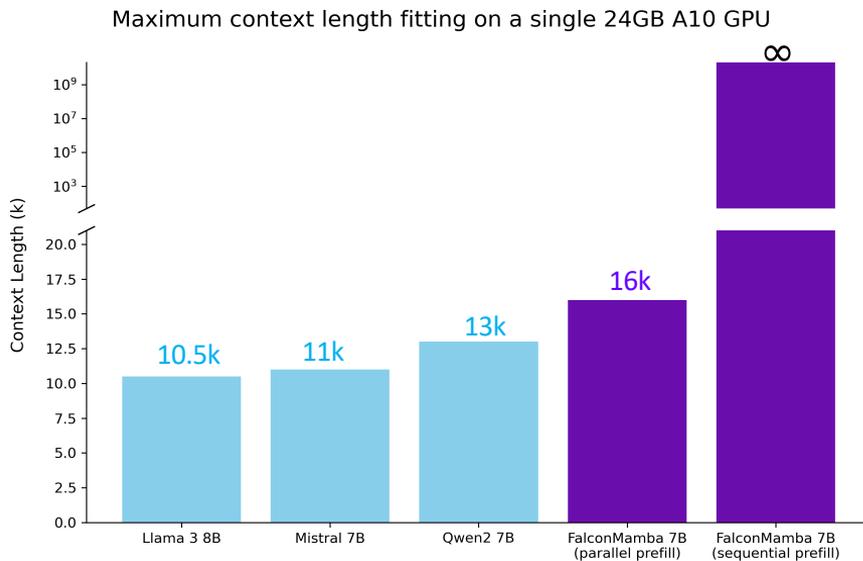
Model Name	ARC-25	HellaSwag-10	MMLU-5	Winogrande-5	TruthfulQA-0	GSM8K-5	Average
<b>RWKV models</b>							
RWKV-v6-Finch-7B*	43.86	75.19	41.69	68.27	42.19	19.64	48.47
RWKV-v6-Finch-14B*	47.44	78.86	52.33	71.27	45.45	38.06	55.57
<b>Transformer models</b>							
Falcon2-11B	59.73	<u>82.91</u>	58.37	78.30	<u>52.56</u>	<u>53.83</u>	<b>64.28</b>
Meta-llama-3-8B	60.24	82.23	<b>66.70</b>	78.45	42.93	45.19	62.62
Meta-llama-3.1-8B	58.53	82.13	<u>66.43</u>	74.35	44.29	47.92	62.28
Mistral-7B-v0.1	59.98	<b>83.31</b>	64.16	78.37	42.15	37.83	60.97
Mistral-Nemo-Base-2407 (12B)	57.94	82.82	64.43	73.72	49.14	<b>55.27</b>	63.89
Gemma-7B	<u>61.09</u>	82.20	64.56	<u>79.01</u>	44.79	50.87	63.75
<b>Hybrid SSM-attention models</b>							
RecurrentGemma-9b**	52.00	80.40	60.50	73.60	38.60	42.60	57.95
Zyphra/Zamba-7B-v1*	56.14	82.23	58.11	<b>79.87</b>	52.88	30.78	60.00
<b>Pure SSM models</b>							
TRI-ML/mamba-7b-rw*	51.25	80.85	33.41	71.11	32.08	4.70	45.52
FalconMamba-7B (pre-decay)*	49.23	80.25	57.27	70.88	37.28	21.83	52.79
FalconMamba-7B*	<b>62.03</b>	80.82	62.11	73.64	<b>53.42</b>	52.54	<u>64.09</u>

Table 3: Model Performance on HF Leaderboard v2: **bold** (best), underline (second best)

Model Name	IFEval-0	BBH-3	Math-Lvl5-4	GPQA-0	MuSR-0	MMLU-PRO-5	Average
<b>RWKV models</b>							
RWKV-v6-Finch-7B	27.65	9.04	1.11	2.81	2.25	5.85	8.12
RWKV-v6-Finch-14B	29.81	12.89	1.13	5.01	3.16	11.3	10.55
<b>Transformer models</b>							
Falcon2-11B	<u>32.61</u>	21.94	2.34	2.80	7.53	15.44	13.78
Meta-llama-3-8B	14.55	24.50	3.25	<u>7.38</u>	6.24	24.55	13.41
Meta-llama-3.1-8B	12.70	<u>25.29</u>	4.61	6.15	8.98	<u>24.95</u>	13.78
Mistral-7B-v0.1	23.86	22.02	2.49	5.59	10.68	22.36	14.50
Mistral-Nemo-Base-2407 (12B)	16.83	<b>29.37</b>	<u>4.98</u>	5.82	6.52	<b>27.46</b>	<u>15.08</u>
Gemma-7B	26.59	21.12	<b>6.42</b>	4.92	<b>10.98</b>	21.64	<b>15.28</b>
<b>Hybrid SSM-attention models</b>							
RecurrentGemma-9b	30.76	14.80	4.83	4.70	6.60	17.88	13.20
Zyphra/Zamba-7B-v1*	24.06	21.12	3.32	3.03	7.74	16.02	12.55
<b>Pure SSM models</b>							
TRI-ML/mamba-7b-rw*	22.46	6.71	0.45	1.12	5.51	1.69	6.25
FalconMamba-7B (pre-decay)*	24.05	11.01	1.71	3.05	8.68	8.59	9.52
FalconMamba-7B	<b>33.36</b>	19.88	3.63	<b>8.05</b>	10.86	14.47	15.04

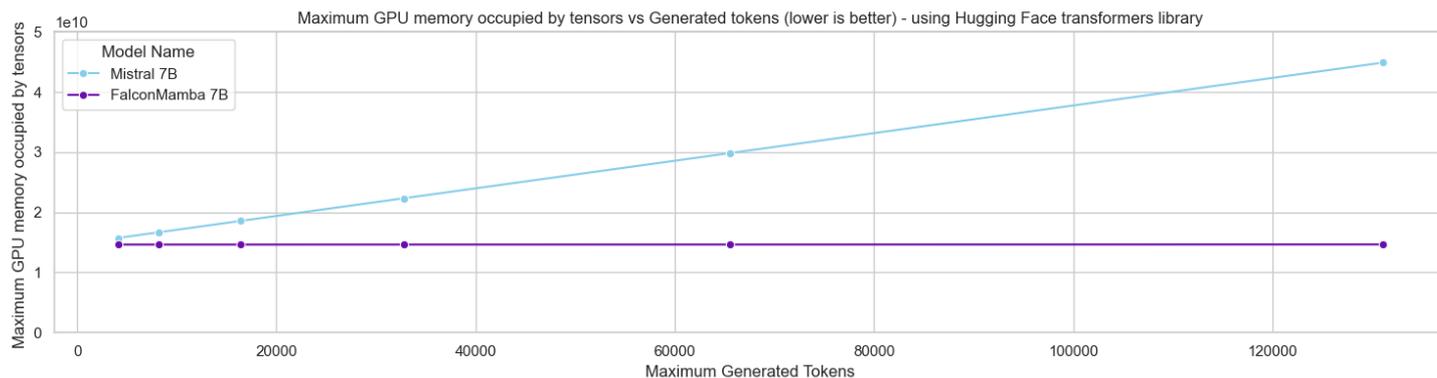
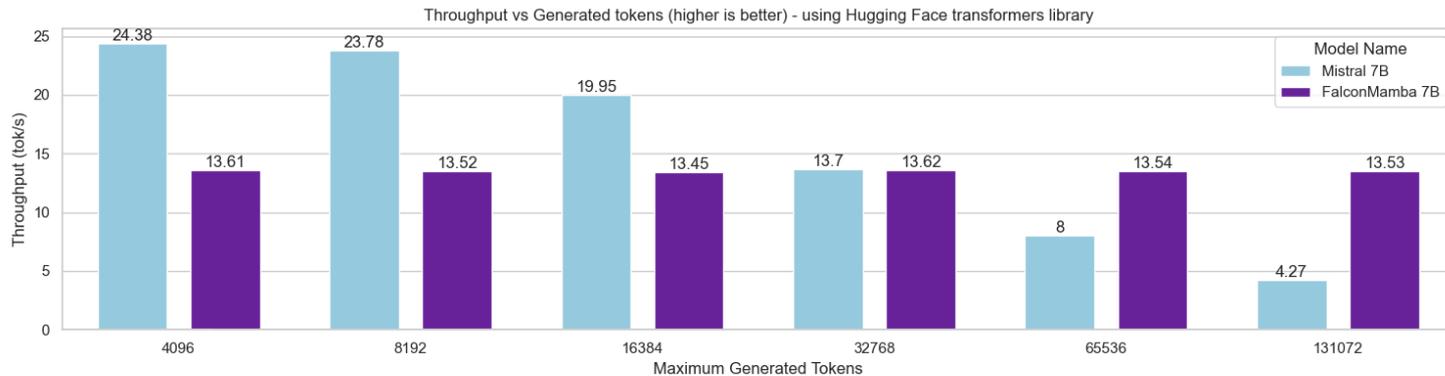
# Falcon Mamba 7B VS SoTA LLMs

- Infinite context length (sequential prefill)
- ~16K context length (parallel prefill), **outperforming** Transformer models



# Falcon Mamba 7B VS SoTA LLMs

- Constant throughput regardless the generation length
- Constant memory usage and lower memory cost than Mistral 7B



# How to use Falcon Mamba 7B

## HuggingFace Ecosystem

Direct use with HF Transformers

- Multiple precisions (quantized & non quantized)
- Multiple platforms (CPU & GPUs)

Fine tuning

- SFTrainer, PEFT, etc.

## Falcon Mamba 7B Chat model (Instruct)

Online, full precision

- HuggingFace Playground

Offline, multiple precisions

- Llama.cpp, LM Studio, etc.



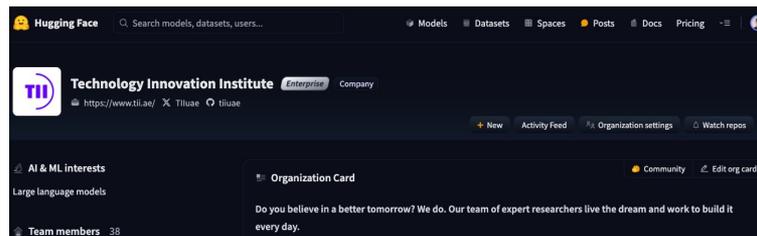
**Hugging Face**

**LLaMA**<sup>C+</sup>



**LM Studio**

**MLX**



```
# pip install accelerate
from transformers import AutoTokenizer, AutoModelForCausalLM

tokenizer = AutoTokenizer.from_pretrained("tiiuae/falcon-mamba-7b")
model = AutoModelForCausalLM.from_pretrained("tiiuae/falcon-mamba-7b", device_map=

input_text = "Question: How many hours in one day? Answer: "
input_ids = tokenizer(input_text, return_tensors="pt").input_ids.to("cuda")

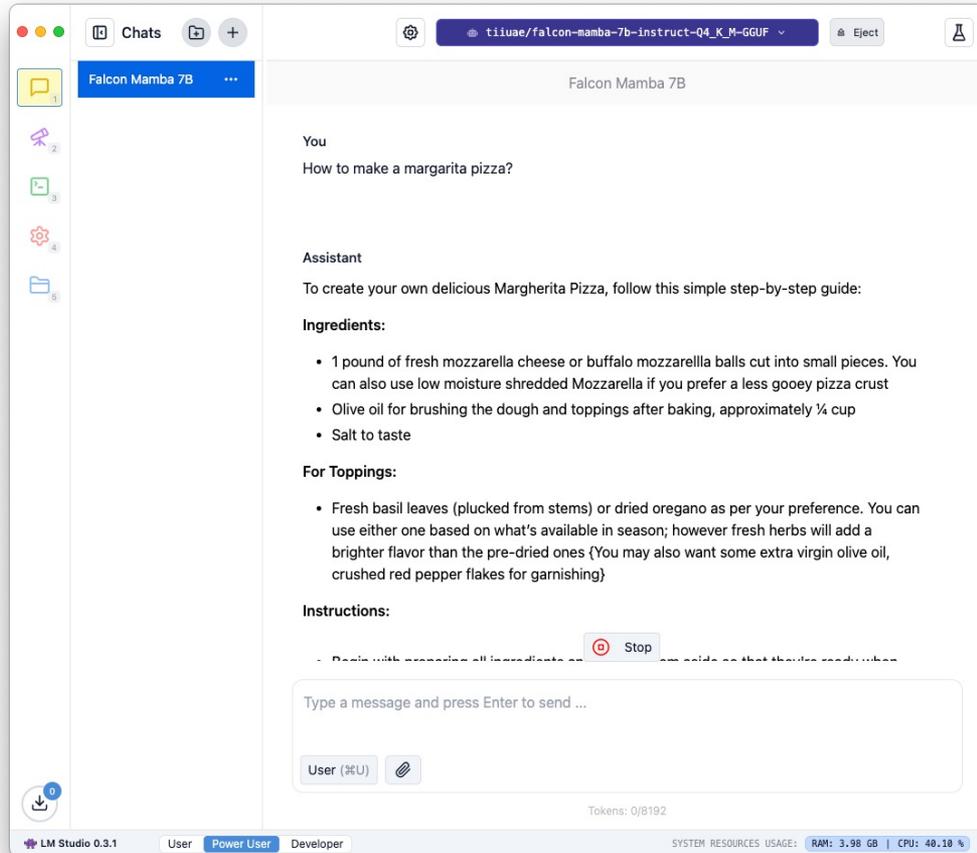
outputs = model.generate(input_ids)
print(tokenizer.decode(outputs[0]))
```



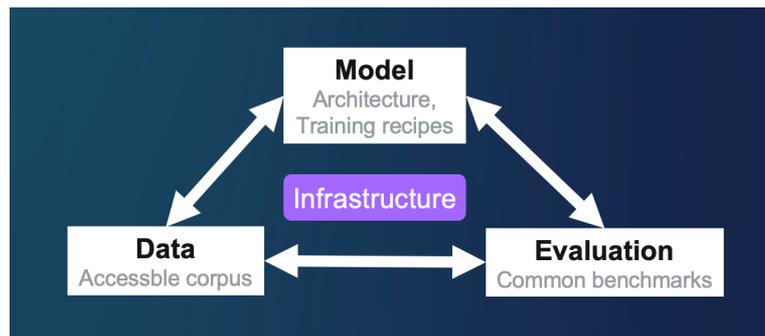
# How to use Falcon Mamba 7B - LM Studio

The screenshot displays the LM Studio application window. At the top, a search bar contains 'falcon mamba'. Below the search bar, a list of models is shown, sorted by 'Best Match'. The selected model, 'falcon-mamba-7b-instruct-Q4\_K\_M-GGUF', is highlighted in blue. The right-hand panel provides details for this model, including a 'Model Card' link, a 'Choose a download option' section with 'Full GPU' and 'Downloaded' indicators, and a 'Use in New Chat' button. The 'Hugging Face Stats' section shows repository information, 239 downloads, and 1 like. The bottom status bar indicates 'LM Studio 0.3.1' and system resources usage: RAM: 3.98 GB, CPU: 39.60%.

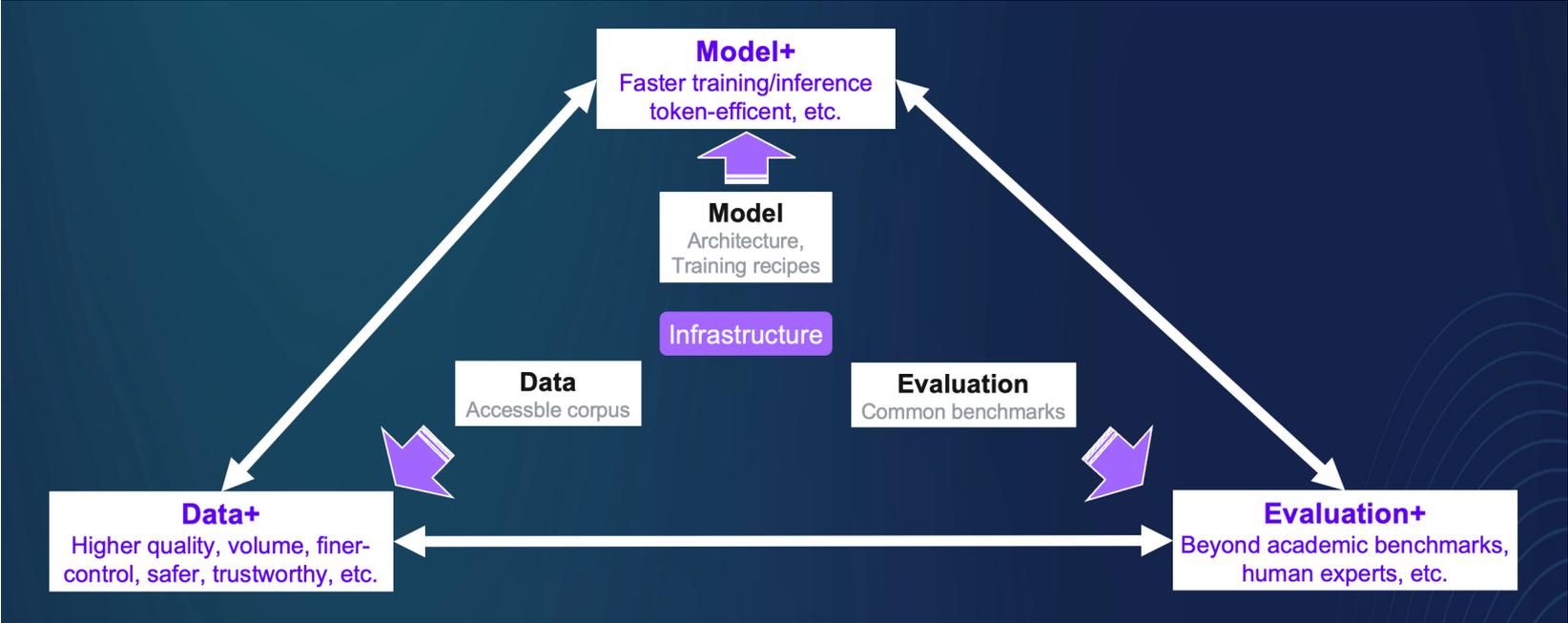
# How to use Falcon Mamba 7B - LM Studio



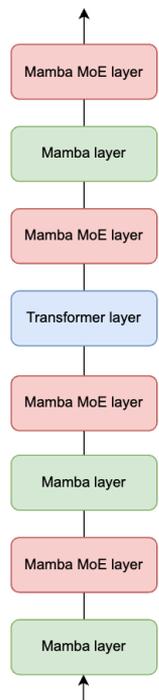
# LLMs - Model & Data & Evaluation



# LLMs - Model & Data & Evaluation

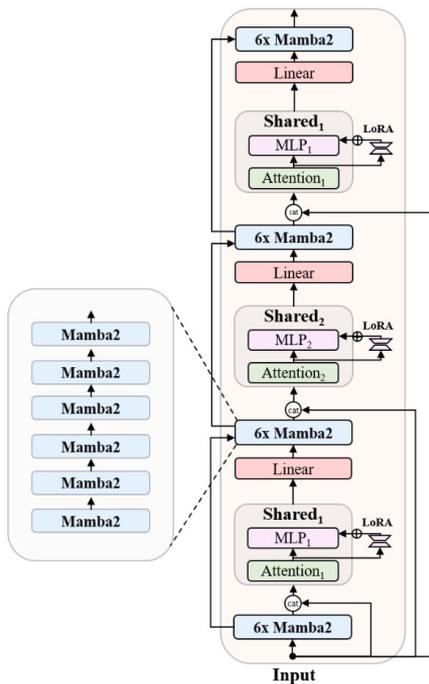


# Hybrid Model - The Future Design of LLMs?



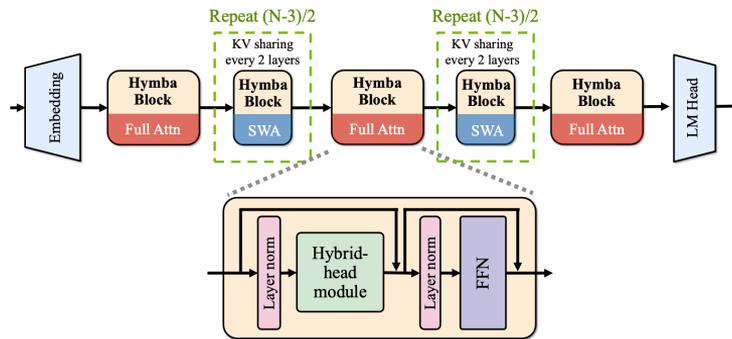
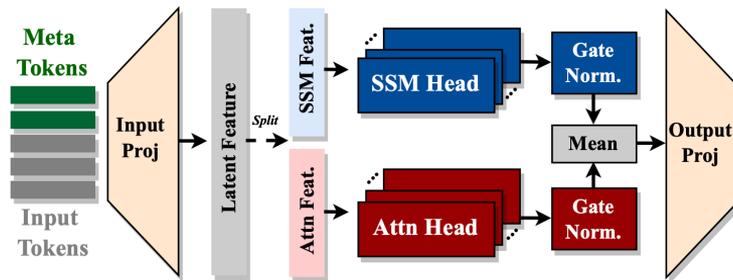
Jamba-MoE 52B  
(12B)

Fig. credit - [Lieber et al., arXiv'24]



Zamba 2 7B

Fig. credit - Zyphra



Hymba 1.5B

Fig. credit - [Dong et al., arXiv'24]

# Falcon 3 Family of Open Models



# Welcome to the Falcon 3 Family of Open Models!



Published December 17, 2024

- [Update on GitHub](#)
- Abdalgader Abubaker** [tiuae](#)
  - Abdulmuneertii Abdul Muneer** [tiuae](#)
  - AdamWma Adam** [tiuae](#)
  - AbdulazizAlshamsi Abdulaziz Alshamsi** [tiuae](#)
  - Alice Pagnoux** [tiuae](#)
  - Almansoorialikhalifa Ali Khalifa Almansoori** [tiuae](#)
  - amztheory Ahmed Alzubaidi** [tiuae](#)
  - Billel Mokdeddem Bilal Mokdeddem** [tiuae](#)
  - Dhia-GB Dhia Garbaya** [tiuae](#)
  - DhiyaEddine Rhaïem** [tiuae](#)
  - dunghuynh ngoc dung huynh guest** [tiuae](#)
  - fedyanin Kirill** [tiuae](#)
  - gcamp Giulia Campesan** [tiuae](#)
  - Gkunsch Guillaume Kunsch** [tiuae](#)
  - gokulkarthik Gokul Karthik** [tiuae](#)
  - griffintaur Ankit Singh** [tiuae](#)
  - HakimHacid Hakim Hacid** [tiuae](#)
  - HamzaYousB9 Hamza Yous** [tiuae](#)
  - hangzou Hang Zou** [tiuae](#)
  - ibrahim-khadraoui-TII ibrahim khadraoui** [tiuae](#)
  - ifaxhat1993 Brahim Farhat** [tiuae](#)
  - Iheb-Chaabane Iheb Chaabane** [tiuae](#)
  - jaadariF Fedi Jaadari** [tiuae](#)
  - JingweiZuo Jingwei Zuo** [tiuae](#)
  - karnakar Karna** [tiuae](#)
  - kasper-piskorski Kasper Piskorski** [tiuae](#)
  - IChahed Ilyas Chahed** [tiuae](#)
  - lkhphuc Phúc Lê Khắc guest** [tiuae](#)
  - LeenAlQadi Leen AlQadi** [tiuae](#)
  - Ludovick Lepauloux** [tiuae](#)
  - melasdeddik Mohamed El Amine Seddik** [tiuae](#)
  - mezsinald Mike Lubinets** [tiuae](#)
  - Mughaira Mughaira** [tiuae](#)
  - puneeshkhanna Puneesh Khanna** [tiuae](#)
  - qiyang-zhao Qiyang Zhao** [tiuae](#)
  - ruxandra Cojocar Ruxandra Cojocar** [tiuae](#)
  - RedaAlami Reda alami** [tiuae](#)
  - rishabh-saraf Rishabh Saraf** [tiuae](#)
  - SanathNarayan Sanath Narayan** [tiuae](#)
  - shihux Shi Hu** [tiuae](#)
  - slimfrikha-tii Slim Frikha** [tiuae](#)
  - wamreyaz wamiq para** [tiuae](#)
  - wdevazelhes William de Vazelhes** [tiuae](#)
  - yellowvm Maksim Velikanov** [tiuae](#)
  - ybelkada Younes Belkada** [tiuae](#)
  - yasserTII Yasser Dahou** [tiuae](#)

**Thank you!**