of Foundation Node's Based on XLSIV



Günter Klambauer



bsky.app/profile/ gklambauer.bsky.social @gklambauer

available here!

Slides



- A word on foundation models
- RNNs & LSTM
- The age of Transformers
- The age of recurrent LLMs
- Applications
 - Vision-xLSTM (ViL)
 - Bio-xLSTM
 - DNA-xLSTM
 - Prot-xLSTM
 - Chem-xLSTM
- Conclusion

A word on foundation models



Image patches Vision Transformer

Over-smoothed attention maps









(a) Patch embedding

Tokens

(b) Self-attention

https://www.leewayhertz.com/vision-transformer-model/

Models considered as foundation models

- Language Models:
 - e.g. BERT, GPT family, LLama family,...
- Vision Models:
 - e.g. Vision Transformers (ViTs), DINO, SimCLR,...
- Multimodal Models: • CLIP, Gemini,...

¹ Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... & Liang, P. (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

Properties of foundation models?

[..] any model that is trained on broad data (generally using self-supervision at scale) that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks. ¹

- Pre-trained on Massive Data
 - \circ EU AI Act: FM if more than 10²⁵ FLOPS used for training
- Self-Supervised Learning
- Scalability
- Multimodal Capabilities
- Adaptability

¹ Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... & Liang, P. (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

Properties of foundation models?

- Pre-trained on Massive Data
- Self-Supervised Learning
- Scalability
- Multimodal Capabilities
- Adaptability
 - Good performance when fine-tuned to a variety of new tasks
 - In-context learning

$\hat{y} = g([\boldsymbol{X}', \boldsymbol{X}]; \boldsymbol{w})$

- New task described via context X' (additional input)
- Model's g(.; w) ability to adapt to and solve new task (X, y) based on additional input X'
- No weights $oldsymbol{w}$ in the neural network changed

In-context learning

Input: 2014-06-01 Output: !06!01!2014! Input: 2007-12-13 Output: !12!13!2007! Input: 2010-09-23 Output: !09!23!2010!

in-context examples

 \boldsymbol{X}

Input: 2005-07-23

test example

 \hat{y} Output: <u>107123120051</u> L - - model completion

RNNs & LSTM

Recurrent neural networks

RNNs map the input of one time step and the last hidden state to a new hidden state

$$\boldsymbol{h}_t = \operatorname{RNN}(\boldsymbol{x}, \boldsymbol{h}_{t-1}; \boldsymbol{w})$$

RNNs are neural nets with adaptive weight matrices (optimized/trained to solve a task):

$$h_t = f(Wx_t + Rh_{t-1})$$
$$\hat{y}_t = \phi(Vh_t)$$

The same transformation is applied in each time step. **Output:** probability of next word.

JYU JOHANNES KEPLER UNIVERSITY LINZ

$$\boldsymbol{h}_t = \operatorname{RNN}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}; \boldsymbol{w})$$





 $oldsymbol{x}_2$

robot

• • •

• • •

 x_{t-1}

obey

 $oldsymbol{x}_t$

orders

$$\boldsymbol{h}_t = \operatorname{RNN}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}; \boldsymbol{w})$$







• • •

• • •





1991: SEPP HOCHREITER'S ANALYSIS OF THE FUNDAMENTAL DEEP LEARNING PROBLEM

$$\begin{split} \| \frac{\partial e(t-q)}{\partial e(t)} \| &= \| \prod_{m=1}^{q} WF'(Net(t-m)) \| \\ &\leq (\| W \| \max_{Net} \{ \| F'(Net) \| \})^{q} \end{split}$$

$$\boldsymbol{h}_t = \operatorname{RNN}(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}; \boldsymbol{w})$$

 \boldsymbol{c}_0



through cell

LSTM cells

Note: scalar notation for cell state!

The unreasonable effectiveness of recurrent neural networks

2

Andrej Karpathy blog

About

Andrej Karpathy

The Unreasonable Effectiveness of Recurrent Neural Networks

May 21, 2015

There's something magical about Recurrent Neural Networks (RNNs). I still remember when I trained my first recurrent network for Image Captioning. Within a few dozen minutes of training my first baby model (with rather arbitrarily-chosen hyperparameters) started to generate very nice looking descriptions of images that were on the edge of making sense. Sometimes the ratio of how simple your model is to the quality of the results you get out of it blows past your expectations, and this was one of those times. What made this result so shocking at the time was that the common wisdom was that RNNs were supposed to be difficult to train (with more experience I've in fact reached the opposite conclusion). Fast forward about a year: I'm training RNNs all the time and I've witnessed their power and robustness many times, and yet their magical outputs still find ways of amusing me. This post is about sharing some of that magic with you.

We'll train RNNs to generate text character by character and ponder the question "how is that even possible?"

Founding member of OpenAl



Overview

	rnns/lstm $\mathcal{O}(T)$	Transformer	xLSTM
Context handling	Theory: infinite context length Practice: few hundred steps		
Memory usage	Constant memory usage		
Parallelizability	Not parallelizable across time steps		

The age of **Transformers**

The Evolution of Language Models



The age of Transformers and LLMs





Attention mechanism

- Dot product between
 - o Keys
 - Queries
- All time steps (tokens) can be processed in parallel
- Softmax for sparseness
- Quadratic in time steps

Layer: 5 \$ Attention: Input - I	nput 🜲
The_	The_
animal_	animal_
didn_	didn_
	·
t_	t_
cross_	cross_
the_	the_
street_	street_
because_	because_
it_	it
was_	was_
too_	too_
tire	tire
d_	d_

Generative Pre-trained Transformers (GPT)

- Input as set
 - $oldsymbol{X} = \{oldsymbol{x}_1, \dots, oldsymbol{x}_t, \dots, oldsymbol{x}_T\}$
- Main operation (parallel!) $H^0 = X$
- $$\begin{split} \boldsymbol{K}^{(l)} &= \boldsymbol{W}_{K}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{Q}^{(l)} &= \boldsymbol{W}_{Q}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{V}^{(l)} &= \boldsymbol{W}_{V}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{H}^{(l)} &= \underbrace{\operatorname{softmax} \left(1/\sqrt{d_{k}} \; \boldsymbol{Q}^{(l)}^{T} \boldsymbol{K}^{(l)} \right)}_{:=\boldsymbol{A}^{(l)}} \boldsymbol{V}^{(l)}, \end{split}$$



Generative Pre-trained Transformers (GPT)

- Input as set
 - $oldsymbol{X} = \{oldsymbol{x}_1, \dots, oldsymbol{x}_t, \dots, oldsymbol{x}_T\}$
- Main operation (parallel!) $H^0 = X$
- $$\begin{split} \boldsymbol{K}^{(l)} &= \boldsymbol{W}_{K}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{Q}^{(l)} &= \boldsymbol{W}_{Q}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{V}^{(l)} &= \boldsymbol{W}_{V}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{H}^{(l)} &= \underbrace{\operatorname{softmax} \left(1/\sqrt{d_{k}} \; \boldsymbol{Q}^{(l)}^{T} \boldsymbol{K}^{(l)} \right)}_{:=\boldsymbol{A}^{(l)}} \boldsymbol{V}^{(l)}, \end{split}$$



Generative Pre-trained Transformers (GPT)

- Input as set
 - $oldsymbol{X} = \{oldsymbol{x}_1, \dots, oldsymbol{x}_t, \dots, oldsymbol{x}_T\}$
- Main operation (parallel!) $H^0 = X$
- $$\begin{split} \boldsymbol{K}^{(l)} &= \boldsymbol{W}_{K}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{Q}^{(l)} &= \boldsymbol{W}_{Q}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{V}^{(l)} &= \boldsymbol{W}_{V}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{H}^{(l)} &= \underbrace{\operatorname{softmax} \left(1/\sqrt{d_{k}} \; \boldsymbol{Q}^{(l)}^{T} \boldsymbol{K}^{(l)} \right)}_{:=\boldsymbol{A}^{(l)}} \boldsymbol{V}^{(l)}, \end{split}$$



Generative Pre-trained Transformers (GPT)

• Input as set

$$oldsymbol{X} = \{oldsymbol{x}_1, \dots, oldsymbol{x}_t, \dots, oldsymbol{x}_T\}$$

- Main operation (parallel!) $H^0 = X$
- $$\begin{split} \boldsymbol{K}^{(l)} &= \boldsymbol{W}_{K}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{Q}^{(l)} &= \boldsymbol{W}_{Q}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{V}^{(l)} &= \boldsymbol{W}_{V}^{(l)} \boldsymbol{H}^{(l-1)} \\ \boldsymbol{H}^{(l)} &= \underbrace{\operatorname{softmax} \left(1/\sqrt{d_{k}} \ \boldsymbol{Q}^{(l)}^{T} \boldsymbol{K}^{(l)} \right)}_{:=\boldsymbol{A}^{(l)}} \boldsymbol{V}^{(l)}, \end{split}$$



Parallelizability of the transformer

Attention mechanism

- Dot product between
 - o Keys
 - Queries
- All time steps (tokens) can be processed in parallel
- Softmax for sparseness
- Quadratic in context size



Overview

	rnns/lstm $\mathcal{O}(T)$	Transformer ${\cal O}(T^2)$	xLSTM
Context handling	Theory: infinite context length Practice: few hundred inputs	Strong long-range capacity Fixed context window	
Memory usage	Constant memory usage	Memory usage grows with context size	
Parallelizability	Not parallelizable across time steps	Parallelizable across time steps	

The age of recurrent LLMs?

LSTM and the test-of-time



LSTMs in reinforcement learning, e.g. OpenAl Five

Berner, C., Brockman, G., Chan, B., Cheung, V., Dębiak, P., Dennison, C., ... & Zhang, S. (2019). Dota 2 with large scale deep reinforcement learning. arXiv preprint arXiv:1912.06680. LSTMs as basis of language models (until 2017)

LSTMs in hydrology Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., ... & Matias, Y. (2024). Global prediction of

extreme floods in ungauged watersheds. *Nature*, 627(8004), 559-563.

Long short-term memory and the constant error carousel (CEC)

$$c_t = \mathbf{f}_t \ c_{t-1} + \mathbf{i}_t \ z_t$$

$$h_t = \mathbf{o}_t \ \psi(\mathbf{c}_t)$$

- Cell state c_t
- Cell inputs z_t
- Forget gate f_t
- Input gate i_t

Review of LSTM

Note: scalar notation for cell state!

How far do we get with LSTM in language modeling?

when

• scaling LSTMs to **billions of parameters**

- leveraging the **latest techniques** from LLMs
- mitigating known limitations of LSTMs?
Overview: from LSTM to xLSTM



LSTM limitations

- Inability to revise storage decisions
 - Nearest Neighbor Search
- Limited storage capacities
 - Rare Token Prediction
 - problem with new or rare tokens
- No parallelization for training

sLSTM

Revise storage decisions through

- exponential gating
- exponential input gate and normalization scales down old inputs

sLSTM

c_t	$= \mathrm{f}_t \ c_{t-1} \ + \ \mathrm{i}_t \ z_t$		cell state	(8)
n_t	$=$ f _t n_{t-1} + i _t		normalizer state	(9)
h	$h_t = o_t \tilde{h}_t$,	$ ilde{h}_t \;=\; c_t \;/\; n_t$	hidden state	(10)
z_t	$= \varphi \left(ilde{z}_t ight) \; ,$	$\tilde{z}_t = \boldsymbol{w}_z^\top \boldsymbol{x}_t + r_z h_{t-1} + b_z$	cell input	(11)
i_t	$= \exp \left(\tilde{\mathbf{i}}_t \right) ,$	$\widetilde{\mathrm{i}}_t \ = \ oldsymbol{w}_\mathrm{i}^ op oldsymbol{x}_t \ + \ r_\mathrm{i} \ h_{t-1} \ + \ b_\mathrm{i}$	input gate	(12)
f_t	$= \sigma\left(ilde{\mathrm{f}}_t ight) \operatorname{OR} \exp\left(ilde{\mathrm{f}}_t ight) ,$	$ ilde{\mathrm{f}}_t \ = \ oldsymbol{w}_{\mathrm{f}}^{ op} oldsymbol{x}_t \ + \ r_{\mathrm{f}} \ h_{t-1} \ + \ b_{\mathrm{f}}$	forget gate	(13)
0 _t	$= \sigma(\tilde{\mathbf{o}}_t)$,	$\tilde{\mathbf{o}}_t = \boldsymbol{w}_{\mathbf{o}}^{\top} \boldsymbol{x}_t + r_{\mathbf{o}} h_{t-1} + b_{\mathbf{o}}$	output gate	(14)

Revise storage decisions: $i_t = \infty, \quad ilde{h}_t = z_t$

sLSTM: scalar LSTM

Memory $oldsymbol{c} \in \mathbb{R}^m$ with m memory cells

Memory capacity: m/d where d is the hidden dimension

The memory is independent of:

- sequence length
- hidden dimension



mLSTM

- Limited storage capacities
- New memory structure:
 - matrix memory
 - covariance update rule
- Storage capacity is maximized for fixed memory



mLSTM

$oldsymbol{C}_t$	=	$\mathbf{f}_t \boldsymbol{C}_{t-1} +$	$\mathbf{i}_t \boldsymbol{v}_t \; \boldsymbol{k}_t$	t			cell state (19)
$oldsymbol{n}_t$	=	$\mathbf{f}_t \mid \boldsymbol{n}_{t-1} \mid +$	$i_t k_t$				normalizer state (20)
$oldsymbol{h}_t$	=	$oldsymbol{o}_t ~\odot~ ilde{oldsymbol{h}}_t ~,$	$ ilde{m{h}}$	$t_t =$	$oldsymbol{C}_t oldsymbol{q}_t \ / \ \max \Big\{ igg \cdot$	$oldsymbol{n}_t^{ op} oldsymbol{q}_t ig , 1 \Big\}$	hidden state (21)
$oldsymbol{q}_t$	=	$oldsymbol{W}_q oldsymbol{x}_t \ + \ oldsymbol{b}_q$					query input (22)
$oldsymbol{k}_t$	=	$rac{1}{\sqrt{d}} oldsymbol{W}_k oldsymbol{x}_t +$	$oldsymbol{b}_k$				key input (23)
$oldsymbol{v}_t$	=	$oldsymbol{W}_v oldsymbol{x}_t + oldsymbol{b}_v$					value input (24)
\mathbf{i}_t	=	$\exp\left(\tilde{\mathbf{i}}_t\right)$,	Î	$\tilde{\mathbf{i}}_t =$	$oldsymbol{w}_{\mathrm{i}}^{ op} oldsymbol{x}_t \ + \ b_{\mathrm{i}}$		input gate (25)
\mathbf{f}_t	=	$\sigma(\tilde{\mathbf{f}}_t)$ OR exp	$\mathbf{p}\left(\tilde{\mathbf{f}}_{t}\right), \; \tilde{\mathbf{f}}$	$\tilde{\mathbf{f}}_t = \mathbf{f}_t$	$oldsymbol{w}_{\mathrm{f}}^{ op} oldsymbol{x}_t \ + \ b_{\mathrm{f}}$		forget gate (26)
\mathbf{o}_t	=	$\sigma\left(ilde{\mathbf{o}}_{t} ight) \;,$	ĉ	$\tilde{\mathbf{b}}_t =$	$oldsymbol{W_o} oldsymbol{x}_t + oldsymbol{b_o}$		output gate (27)

Parallelization: gates, keys, queries, values all independent from h_{t-1}

Parallelizability of mLSTM



44

LSTM limitations

- Inability to revise storage decisions
 - Nearest Neighbor Search
- Limited storage capacities
 - Rare Token Prediction
 - problem with new or rare tokens
- No parallelization for training







xLSTM: mLSTM & sLSTM



Backbone architecture



xLSTM: Length extrapolation



xLSTM Scaling Laws



49

xLSTM: Scaling Laws

J⊻U



xLSTM: Scaling Laws

J⊻U



xLSTM-7B: A Recurrent LLM for Fast and Efficient Inference

- First LSTM-based 7B language model
- Language modeling performance similar to other linear LLMs
- More efficient (faster)
- Several changes to architectures
 - E.g. fully based on mLSTM
 - Changes in normalization and gating
- **Pareto front:** best choice given fixed computational budget
- Computing:
 - Cluster with 256 NVIDIA H100 GPUs
 - Pre-Training: ~140k GPU hours,
 2.3T tokens, +Long-Context Version

Beck, M., Pöppel, K., Lippe, P., Kurle, R., Blies, P. M., Klambauer, G., ... & Hochreiter, S. (2025). xLSTM 7B: A Recurrent LLM for Fast and Efficient Inference. *International Conference on Machine Learning (ICML)*.



xLSTM: efficient kernels hardware optimized



 Triton kernels render xLSTM extremely fast

Overview

	RNNs/LSTM $\mathcal{O}(T)$	Transformer ${\cal O}(T^2)$	xlstm $\mathcal{O}(T)$
Context handling	Theory: infinite context length Practice: few hundred inputs	Strong long-range capacity Fixed context window	Strong long-range capacity Length extrapolation
Memory usage	Constant memory usage	Memory usage grows with context size	Constant memory usage
Parallelizability	Not parallelizable across time steps	Parallelizable across time steps	Parallelizable across time steps

xLSTM

- New:
 - exponential gating and memory mixing
 - memory structure (matrix memory and covariance update)

• Results:

- xLSTM performs favorably on language modeling when compared to Transformers and State Space models
- The scaling laws indicate that larger xLSTM models will be serious competitors to current large language models

Resources for xLSTM





Vision-xLSTM (ViL)



xLSTM can process image "tokens"

Alkin, B., Beck, M., Pöppel, K., Hochreiter, S., & Brandstetter, J. (2025). Vision-LSTM: xLSTM as Generic Vision Backbone. International Conference on Learning Representations (ICLR)

Vision-xLSTM (ViL)



• High efficiency and accuracy

Alkin, B., Beck, M., Pöppel, K., Hochreiter, S., & Brandstetter, J. (2025). Vision-LSTM: xLSTM as Generic Vision Backbone. International Conference on Learning Representations (ICLR)



Motivation

- **Biological Foundation Models:** large models trained on biological domains capture important concepts, which enables downstream applications.
- **Transformer Dominance:** Transformers remain dominant due to their versatility and performance at large scales.
- Scalability & Memory: Transformers face runtime and memory challenges. Efficient architectures are critical for processing long-range biological data efficiently.
- **xLSTM**: xLSTM has shown strong results in natural language tasks. This work evaluates its scalability and performance in biological applications.

Bio-xLSTM



Schmidinger, N., Schneckenreiter, L., Seidl, P., Schimunek, J., Hoedt, P. J., Brandstetter, J., ... & Klambauer, G. (2025). Bio-xLSTM: Generative modeling, representation and in-context learning of biological and chemical sequences. *International Conference on Learning Representations 2025*.



xLSTM Bidirectionality





$$egin{aligned} egin{aligned} egin{aligne} egin{aligned} egin{aligned} egin{aligned} egin$$

Weight Sharing

Alternating Block Directions mLSTM Native Bidirectionality

ELMo, Peters et.al. (2018) Vision-LSTM, Alkin et.al. (2024)

Reverse-Complement Invariance

 $ext{xLSTM}(ext{RC}(extbf{X})) = ext{RC}(ext{xLSTM}(extbf{X}))$ Logits

 $\mathrm{xLSTM}(\mathbf{X}) = \mathrm{xLSTM}(\mathrm{RC}(\mathbf{X}))$



Caduceus, Schiff et. al. (2024)

ACGCTTGATCGTAGTGTAA

DNA-xLSTM: Short Context



DNA-xLSTM: Short Context



DNA-xLSTM-2M MLM



DNA-xLSTM: Downstream Adaptation





ACGCTTGATCGTAGTGTAA



							F		
Task	Metric	> 100M Param. Models			2M Param. Models				
		Enformer (252M)	DNABERT-2 (117M)	NT-v2 (500M)	HyenaDNA	Mamba-PS	Mamba-PH	xLSTM-PS	xLSTM-PH
Histone Markers									
H3	MCC \uparrow	$0.719^{\pm 0.048}$	$0.785^{\pm 0.033}$	$0.784^{\pm 0.047}$	$0.779^{\pm 0.037}$	$0.799^{\pm 0.029}$	$0.815^{\pm 0.048}$	$0.796^{\pm 0.014}$	$0.824^{\pm 0.010}$
H3K14AC	MCC \uparrow	$0.288^{\pm 0.077}$	$0.516^{\pm 0.028}$	$0.551^{\pm 0.021}$	$0.612^{\pm 0.065}$	$0.541^{\pm 0.212}$	$0.631^{\pm 0.026}$	$0.570^{\pm 0.003}$	$0.598^{\pm 0.017}$
H3K36ME3	MCC ↑	$0.344^{\pm 0.055}$	$0.591^{\pm 0.020}$	$0.625^{\pm 0.030}$	$0.613^{\pm 0.041}$	$0.609^{\pm 0.109}$	$0.601^{\pm 0.129}$	$0.588^{\pm 0.01}$	$0.625^{\pm 0.010}$
H3K4ME1	MCC ↑	$0.291^{\pm 0.061}$	$0.511^{\pm 0.028}$	$0.550^{\pm 0.021}$	$0.512^{\pm 0.024}$	$0.488^{\pm 0.102}$	$0.523^{\pm 0.039}$	$0.490^{\pm 0.012}$	$0.526^{\pm 0.001}$
H3K4ME2	MCC ↑	$0.211^{\pm 0.069}$	$0.336^{\pm 0.040}$	$0.319^{\pm 0.045}$	$0.455^{\pm 0.095}$	$0.388^{\pm 0.101}$	$0.487^{\pm 0.170}$	$0.489^{\pm 0.024}$	$0.504^{\pm 0.012}$
H3K4ME3	MCC ↑	$0.158^{\pm 0.072}$	$0.352^{\pm 0.077}$	$0.410^{\pm 0.033}$	$0.549^{\pm 0.056}$	$0.440^{\pm 0.202}$	$0.544^{\pm 0.045}$	$0.520^{\pm 0.019}$	$0.537^{\pm 0.012}$
H3K79ME3	MCC ↑	$0.496^{\pm 0.042}$	$0.613^{\pm 0.030}$	$0.626^{\pm 0.046}$	$0.672^{\pm 0.048}$	$0.676^{\pm 0.026}$	$0.697^{\pm 0.077}$	$0.662^{\pm 0.011}$	$0.697^{\pm 0.007}$
H3K9AC	MCC ↑	$0.420^{\pm 0.063}$	$0.542^{\pm 0.029}$	$0.562^{\pm 0.040}$	$0.581^{\pm 0.061}$	$0.604^{\pm 0.048}$	$0.622^{\pm 0.030}$	$0.622^{\pm 0.013}$	$0.627^{\pm 0.008}$
H4	MCC ↑	$0.732^{\pm 0.076}$	$0.796^{\pm 0.027}$	$0.799^{\pm 0.025}$	$0.763^{\pm 0.044}$	$0.789^{\pm 0.020}$	$0.811^{\pm 0.022}$	$0.793^{\pm 0.011}$	$0.813^{\pm 0.008}$
H4AC	MCC \uparrow	$0.273^{\pm 0.063}$	$0.463^{\pm 0.041}$	$0.495^{\pm 0.032}$	$0.564^{\pm 0.038}$	$0.525^{\pm 0.240}$	$0.621^{\pm 0.054}$	$0.558^{\pm 0.013}$	$0.583^{\pm 0.014}$
Regulatory Annotation									
Enhancer	MCC ↑	$0.451^{\pm 0.108}$	$0.516^{\pm 0.098}$	$0.548^{\pm 0.144}$	$0.517^{\pm 0.117}$	$0.491^{\pm 0.066}$	$0.546^{\pm 0.073}$	$0.375^{\pm 0.03}$	$0.545^{\pm 0.024}$
Enhancer Types	MCC ↑	$0.309^{\pm 0.134}$	$0.423^{\pm 0.051}$	$0.424^{\pm 0.132}$	$0.386^{\pm 0.185}$	$0.416^{\pm 0.095}$	$0.439^{\pm 0.054}$	$0.444^{\pm 0.045}$	$0.466^{\pm 0.011}$
Promoter: All	F1 ↑	$0.954^{\pm 0.006}$	$0.971^{\pm 0.006}$	$0.976^{\pm 0.006}$	$0.960^{\pm 0.005}$	$0.967^{\pm 0.004}$	$0.970^{\pm 0.004}$	$0.962^{\pm 0.002}$	$0.967^{\pm 0.001}$
NonTATA	F1 ↑	$0.955^{\pm 0.010}$	$0.972^{\pm 0.005}$	$0.976^{\pm 0.006}$	$0.959^{\pm 0.011}$	$0.968^{\pm 0.006}$	$0.968^{\pm 0.010}$	$0.963^{\pm 0.002}$	$0.970^{\pm 0.001}$
TATA	F1 ↑	$0.960^{\pm 0.023}$	$0.955^{\pm 0.021}$	$0.966^{\pm 0.013}$	$0.944^{\pm 0.040}$	$0.957^{\pm 0.015}$	$0.953^{\pm 0.016}$	$0.948^{\pm 0.005}$	$0.952^{\pm 0.005}$
Splice Site Annotation									
All	Accuracy ↑	$0.848^{\pm 0.019}$	$0.939^{\pm 0.009}$	$0.983^{\pm 0.008}$	$0.956^{\pm 0.011}$	$0.927^{\pm 0.021}$	$0.940^{\pm 0.027}$	$0.965^{\pm 0.005}$	$0.974^{\pm 0.004}$
Acceptor	F1 ↑	$0.914^{\pm 0.028}$	$0.975^{\pm 0.006}$	$0.981^{\pm 0.011}$	$0.958^{\pm 0.010}$	$0.936^{\pm 0.077}$	$0.937^{\pm 0.033}$	$0.970^{\pm 0.005}$	$0.953^{\pm 0.008}$
Donor	F1 ↑	$0.906^{\pm 0.027}$	$0.963^{\pm 0.006}$	$0.985^{\pm 0.022}$	$0.949^{\pm 0.024}$	$0.948^{\pm 0.025}$	$0.874^{\pm 0.289}$	$0.962^{\pm 0.004}$	$0.951^{\pm 0.005}$

DNA-xLSTM: Long Context



Bio-xLSTM



Schmidinger, N., Schneckenreiter, L., Seidl, P., Schimunek, J., Hoedt, P. J., Brandstetter, J., ... & Klambauer, G. (2025). Bio-xLSTM: Generative modeling, representation and in-context learning of biological and chemical sequences. International Conference on Learning Representations.

DNA-xLSTM: Long Context




Winning Strategy: Incorporate Evolutionary Information

Transception



MSA Transformer



ESMFold



Single sequence

AlphaFold



AlphaFold 2, Jumper et.al. (2021) ESMFold, Lin et.al. (2022) Transception, Notin et.al. (2022) MSA Transformer, Rao et.al.(2021)

Homology-Aware Modeling via Conditioning



Fill-in-the-Middle (FiM) Augmentation



Prot-xLSTM Pretraining



Homology-aware Pre-Training



Homology-aware Generation



	Prot-xLSTM-26M	ProtMamba-28M	Prot-xLSTM-102M	ProtMamba-107M
Sequence Length	$0.41^{\pm0.09}$	$0.52^{\pm0.09}$	$0.40^{\pm0.08}$	$0.36^{\pm 0.08}$
Min. Hamming	$0.43^{\pm 0.08}$	$0.60^{\pm0.11}$	$0.47^{\pm 0.09}$	$0.42^{\pm 0.07}$
HMMER	$0.57^{\pm0.10}$	$0.54^{\pm 0.11}$	$0.44^{\pm 0.09}$	$0.49^{\pm 0.10}$
pLDDT	$0.40^{\pm 0.09}$	$0.68^{\pm 0.12}$	$0.27^{\pm 0.05}$	$0.30^{\pm 0.07}$
pTM	$0.38^{\pm0.08}$	$0.72^{\pm 0.10}$	$0.26^{\pm0.05}$	$0.28^{\pm0.05}$

Zero-Shot Fitness Prediction



Mask Mutation

Zero-Shot Fitness Prediction Evaluation

Model Type	Model	Reference	#Params	Spearman ρ
Alignment-based	Site-Independant	Hopf et al. (2017)		0.359
	EVE	Frazer et al. (2021)	_ ^a	0.432
	GEMME	Laine et al. (2019)	-	0.455
Protein language model	Tranception L (w/o R)	Notin et al. (2022a)	700M	0.374
(PLM)	VespaG	Marquet et al. (2024)	3B	0.458
	ProGen2 XL	Nijkamp et al. (2023)	6B	0.391
	ESM-2	Lin et al. (2023)	15B	0.401
Alignment + PLM	MSA-Transformer	Rao et al. (2021)	100M	0.432
	Tranception L (w/ R)	Notin et al. (2022a)	700M	0.434
	TranceptEVE L	Notin et al. (2022b)	>700M ^a	0.456
Homology-aware PLM	Prot-xLSTM	Ours	26M	0.411 ^b
	ProtMamba	Sgarbossa et al. (2024)	28M	0.360 ^b
	Prot-xLSTM	Ours	102M	0.416 ^b
	ProtMamba	Sgarbossa et al. (2024)	107M	0.415 ^b
	PoET	Truong Jr and Bepler (2023)	201M	0.470
Inverse folding	ESM-IF1	Hsu et al. (2022)	142M	0.422
Structure + PLM	SaProt	Su et al. (2024b)	35M	0.407
	ProSST	Li et al. (2024)	110M	0.507
	SaProt	Su et al. (2024b)	650M	0.457

^a EVE parameters depend on the size of a given MSA.

^b This work. All other values are retrieved from ProteinGym on 03/11/2024.



Unconditional Pre-Training



Unconditional Generation Evaluation

Recurrent Decoding until eos



	SMILES-LSTM ^a	SMILES-GPT ^b	SMILES-S4 ^c	Chem-Mamba ^d	Chem-xLSTM
FCD↓	$0.46^{\pm < 0.01}$	$0.15^{\pm < 0.01}$	$0.28^{\pm < 0.01}$	$0.21^{\pm < 0.01}$	$0.13^{\pm < 0.01}$
Perplexity \downarrow	$1.88^{\pm 3.8}$	$1.65^{\pm0.6}$	$1.73^{\pm 2.4}$	$1.74^{\pm0.5}$	$1.68^{\pm1.0}$
^a Segler et al. (2018) ^b Adilov (2021) ^c Özçelik et al. (2024) ^d Adapted to SMILES in this work.					

Conditional Generation



In-context learning







Prot-xLSTM



Chem-xLSTM

Huggingface App here!

Paper







e 1 1 i s

Resources (Bio)

- Bio-xLSTM: LLM models for DNA, proteins and small molecules
 - Paper (ICLR 2025): <u>https://openreview.net/forum?id=ljbXZdugdi</u>
- VN-EGNN: binding pocket identification method
 - Paper (pre-print): <u>https://arxiv.org/abs/2404.07194</u>
 - Github: https://github.com/ml-jku/vnegnn
 - HuggingFace App: <u>https://huggingface.co/spaces/ml-jku/vnegnn</u>
- CLOOME: most powerful features for microscopy images ("cell painting")
 - Paper (NComms): <u>https://www.nature.com/articles/s41467-023-42328-w</u>
 - Github: https://github.com/ml-jku/cloome
 - HuggingFace App: <u>https://huggingface.co/spaces/anasanchezf/cloome</u>
- GNN-VPA: variance-preserving aggregation for message-passing networks
 - Paper (ICLR 2024): <u>https://arxiv.org/abs/2403.04747</u>
 - Implementation: <u>https://pytorch-geometric.readthedocs.io/en/latest/</u> generated/torch_geometric.nn.aggr.VariancePreservingAggregation.html
- CLAMP: multi-modal bioactivity prediction model with in-context capacity
 - Paper (ICML 2023): <u>https://arxiv.org/abs/2303.03363</u>
 - Github: <u>https://github.com/ml-jku/clamp</u>
- MHNfs: best few-shot learning method via in-context learning
 - Paper (ICLR 2023): <u>https://openreview.net/pdf?id=XrMWUuEevr</u>
 - Github: https://github.com/ml-jku/MHNfs
 - HuggingFace App: <u>https://huggingface.co/spaces/ml-jku/mhnfs</u>
- LAM-Slide: highly efficient and accurate molecular dynamics
 - Paper (pre-print): <u>https://arxiv.org/abs/2502.12128</u>
 - Blog: <u>https://ml-jku.github.io/LaM-SLidE/</u>





Schimunek, J., Seidl, P., Friedrich, L., Kuhn, D., Rippmann, F., Hochreiter, S., & Klambauer, G. (2023). Context-enriched molecule representations improve few-shot drug discovery. *International Conference on Learning Representations*. <u>https://huggingface.co/spaces/ml-jku/mhnfs</u>





- Domain-Specific Adaptations: Introduced xLSTM variants tailored to NLP, computer vision, and major biological domains—DNA-xLSTM, Prot-xLSTM, and Chem-xLSTM
- **State-of-the-Art Performance**: Demonstrated that xLSTM consistently matches or surpasses strong baselines across all three domains.
- **Scalability & Efficiency**: Validated xLSTM's ability to process sequences up to 260,000 tokens while maintaining superior training speed and inference efficiency compared to Transformers.
- Foundation Model Potential: Positioned xLSTM as a leading candidate for foundational models in biology, bridging performance, scalability, and domain adaptability.

ACKNOWLEDGEMENTS LIT AI Lab & ELLIS unit Linz Johannes Kepler Universität Linz



Protein design and optimization (PDO)



AIDD

Project funded by the European Union's Horizon 2020 research and innovation programme under the <u>Marie</u> <u>Skłodowska-Curie grant agreement No</u> <u>956832</u>, and it is Horizon 2020 Marie Skłodowska-Curie Innovative Training Network - European Industrial Doctorate.





Deep Learning for drug discovery & design



Bilateral Artificial Intelligence



We are hiring!



Protein design and optimization (PDO)



AIDD

Project funded by the European Union's Horizon 2020 research and innovation programme under the <u>Marie</u> <u>Skłodowska-Curie grant agreement No</u> <u>956832</u>, and it is Horizon 2020 Marie Skłodowska-Curie Innovative Training Network - European Industrial Doctorate.





Deep Learning for drug discovery & design



Bilateral Artificial Intelligence



BACKUP SLIDES



Essential References

Hochreiter, S. (1991). Untersuchungen zu dynamischen neuronalen Netzen. Diploma, Technische Universität München, 91(1), 31.

Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.

Beck, M., Pöppel, K., Spanring, M., Auer, A., Prudnikova, O., Kopp, M., ... & Hochreiter, S. (2024). xLSTM: Extended Long Short-Term Memory. *Advances in neural information processing systems, 37.*

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, *33*, 1877-1901.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, *35*, 27730-27744.

Yang, S., Wang, B., Shen, Y., Panda, R., & Kim, Y. (2023). Gated linear attention transformers with hardware-efficient training. *arXiv* preprint arXiv:2312.06635.

Alkin, B., Beck, M., Pöppel, K., Hochreiter, S., & Brandstetter, J. (2025). Vision-LSTM: xLSTM as Generic Vision Backbone. *International Conference on Learning Representations.*

Schmidinger, N., Schneckenreiter, L., Seidl, P., Schimunek, J., Hoedt, P. J., Brandstetter, J., ... & Klambauer, G. (2025). Bio-xLSTM: Generative modeling, representation and in-context learning of biological and chemical sequences. *International Conference on Learning Representations.*

xLSTM Scaling



xLSTM Background



sLSTM ③

 $egin{aligned} c_t &= \sigma(ilde{\mathrm{f}}_t) \; c_{t-1} + \exp(ilde{\mathrm{i}}_t) \; anh(ilde{z}_t) \ n_t &= \sigma(ilde{\mathrm{f}}_t) \; n_{t-1} + \exp(ilde{\mathrm{i}}_t) \ h_t &= \sigma(ilde{\mathrm{o}}_t) \; rac{c_t}{n_t} \end{aligned}$

+ Enhanced State Tracking



$$egin{aligned} m{C}_t &= \sigma(ilde{\mathrm{f}}_t) \; m{C}_{t-1} + \exp(ilde{\mathrm{i}}_t) \; m{v}_t m{k}_t^{ op} \ m{n}_t &= \sigma(ilde{\mathrm{f}}_t) \; m{n}_{t-1} + \exp(ilde{\mathrm{i}}_t) \; m{k}_t \ m{h}_t &= \sigma(ilde{\mathrm{o}}_t) \odot rac{m{C}_t m{q}_t}{\maxig\{ig|m{n}_t^{ op} m{q}_tig|, 1ig\}} \end{aligned}$$

+ Parallel & Recurrent Mode

NXAI

- Company that sponsored the development of xLSTM
- NXAI scales up xLSTM and builds larger LLMs
- "AI at scale", AI for simulations
- Austrian startup with center in Linz:
- https://www.nx-ai.com/







xLSTM Source Code Now Open Source

AlphaFold



Al breakthrough in biology

Future Directions

Scaling: train larger versions of Bio-xLSTM and evaluate scaling.

Multi-Modality: Incorporate diverse data types, including 3D structures and other non-sequential modalities.

Interpretability: Identify factors contributing to superior performance compared to Transformers and Mamba in specific contexts. Analyze scenarios where xLSTM underperforms.

xLSTM stabilization

$$m_{t} = \max\left(\log(f_{t}) + m_{t-1}, \log(i_{t})\right)$$
stabilizer state

$$i'_{t} = \exp\left(\log(i_{t}) - m_{t}\right) = \exp\left(\tilde{i}_{t} - m_{t}\right)$$
stabil. input gate

$$f'_{t} = \exp\left(\log(f_{t}) + m_{t-1} - m_{t}\right)$$
stabil. forget gate

Numerical stability of the exponential function

Can be seen as running-softmax

LSTM limitations

• Nearest Neighbor Search:

- given a reference vector
- scan a sequence sequentially
- for the most similar vector
- return its attached value at sequence end
- Difficult for models:
 - o whenever new most similar
 → revise value and largest
 similarity



LSTM limitations

- Inability to revise storage decisions
 - Nearest Neighbor Search
- Limited storage capacities
 - Rare Token Prediction
 - problem with new or rare tokens
- No parallelization for training

LSTM limitations

- Rare Token Prediction:
 - perplexity (PPL) of token
 prediction on Wikitext 103
 - partitions of token according to frequency
- rare tokens must be memorized since they are not learned



Post Up-Projection Block

- Typically used for sLSTM
- Known from Transformer
- Memory mixing
- Block-diagonal: multiple heads



Pre Up-Projection Block

- Typically used for mLSTM
- Known from State Space Models
- High memory capacity through pre up-projection
- External output gate
- Convolution for queries and keys
- Swish activations
- Learnable skip connections



xLSTM: state tracking

e.g. majority count

Languages

Formal

some tasks require state tracking

	Context Sensitive		Deterministic Context Free		Regular					
	Bucket Sort	Missing Duplicate	Mod Arithmetic (w Brackets)	Solve Equation	Cycle Nav	Even Pairs	Mod Arithmetic (w/o Brackets)	Parity	Majority	Majority Count
Llama	0.92 ± 0.02	0.08 ± 0.0	0.02 ± 0.0	0.02 ± 0.0	$\begin{array}{c} 0.04 \\ \pm \ 0.01 \end{array}$	$\begin{array}{c} 1.0 \\ \pm \ 0.0 \end{array}$	0.03 ± 0.0	$\underset{\pm \text{ 0.01}}{0.03}$	$\underset{\pm 0.01}{0.37}$	$\underset{\pm \ 0.0}{0.13}$
Mamba	0.69 ± 0.0	$\underset{\pm \ 0.0}{0.15}$	$\underset{\pm 0.01}{0.04}$	$\underset{\pm 0.02}{0.05}$	0.86 ± 0.04	1.0 ± 0.0	$\underset{\pm 0.02}{0.05}$	$\underset{\pm 0.02}{0.13}$	0.69 ± 0.01	$\underset{\pm 0.03}{0.45}$
Retention	0.13 ± 0.01	$\underset{\pm \ 0.03}{0.03}$	0.03 ± 0.0	$\underset{\pm \ 0.03}{0.03}$	$\underset{\pm 0.01}{0.05}$	0.51 ± 0.07	$\underset{\pm \ 0.04}{0.04}$	$\underset{\pm \text{ 0.01}}{0.05}$	$\underset{\pm \ 0.0}{0.36}$	$\underset{\pm 0.01}{0.12}$
Hyena	0.3 ± 0.02	$\underset{\pm 0.02}{0.06}$	$\underset{\pm 0.05}{0.05}$	$\underset{\pm \ 0.02}{0.02}$	$\underset{\pm 0.01}{0.06}$	0.93 ± 0.07	$\underset{\pm \ 0.04}{0.04}$	$\underset{\pm 0.04}{0.04}$	$\underset{\pm 0.01}{0.36}$	$\underset{\pm 0.02}{0.18}$
RWKV-4	0.54 ± 0.0	$\underset{\pm 0.01}{0.21}$	0.06 ± 0.0	$\underset{\pm 0.07}{0.07}$	$\underset{\pm 0.0}{0.13}$	$\underset{\pm 0.0}{1.0}$	$\underset{\pm 0.07}{0.07}$	$\underset{\pm 0.06}{0.06}$	$\underset{\pm 0.0}{0.63}$	$\underset{\pm 0.0}{0.13}$
RWKV-5	0.49 ± 0.04	$\underset{\pm 0.01}{0.15}$	0.08 ± 0.0	$\underset{\pm 0.08}{0.08}$	0.26 ± 0.05	1.0 ± 0.0	$\underset{\pm 0.02}{0.15}$	$\underset{\pm 0.03}{0.06}$	0.73 ± 0.01	$\underset{\pm 0.03}{0.34}$
RWKV-6	0.96 ± 0.0	0.23 ± 0.06	$\underset{\pm 0.01}{0.09}$	$\underset{\pm 0.02}{0.09}$	$\underset{\pm \ 0.14}{0.31}$	1.0 ± 0.0	$\underset{\pm 0.0}{0.16}$	0.22 ± 0.12	0.76 ± 0.01	$\underset{\pm 0.01}{0.24}$
LSTM (Block)	0.99 ± 0.0	$\begin{array}{c} 0.15 \\ \pm 0.0 \end{array}$	0.76 ± 0.0	0.5 ± 0.05	$\underset{\pm 0.03}{0.97}$	1.0 ± 0.0	$\underset{\pm 0.09}{0.91}$	$\underset{\pm 0.0}{1.0}$	0.58 ± 0.02	$\underset{\pm 0.0}{0.27}$
LSTM	0.94 ± 0.01	0.2 ± 0.0	0.72 ± 0.04	$\underset{\pm 0.05}{0.38}$	0.93 ± 0.07	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	0.82 ± 0.02	$\underset{\pm 0.0}{0.33}$
xLSTM[0:1]	0.84 ± 0.08	$\underset{\pm 0.01}{0.23}$	$\underset{\pm 0.09}{0.57}$	$\underset{\pm 0.09}{0.55}$	$\underset{\pm 0.0}{1.0}$	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0	$\underset{\pm 0.02}{0.75}$	0.22 ± 0.0
xLSTM[1:0]	0.97 ± 0.0	0.33 ± 0.22	0.03 ± 0.0	0.03 ± 0.01	0.86 ± 0.01	1.0 ± 0.0	0.04 ± 0.0	0.04 ± 0.01	0.74 ± 0.01	$\underset{\pm 0.0}{0.46}$
xLSTM[1:1]	0.7 ± 0.21	0.2 ± 0.01	$\underset{\pm 0.06}{0.15}$	$\underset{\pm 0.04}{0.24}$	0.8 ± 0.03	1.0 ± 0.0	0.6 ± 0.4	$\underset{\pm 0.0}{1.0}$	0.64 ± 0.04	$\underset{\pm \ 0.0}{0.5}$

xLSTM: memory



Multi-Query Associative Recall task:

- For each sequence, key-value pairs are randomly chosen from a large vocabulary
- must be memorized for later retrieval.

Method comparison: language modeling

Comparison at language modeling

- SlimPajama dataset with 15B tokens
- perplexity (ppl) as metric

Model	# Params M	Slim Pajama (15B) ppl \downarrow
GPT-3	356	14.26
Llama	407	14.25
H3	420	18.23
Mamba	423	13.70
Hyena	435	17.59
RWKV-4	430	15.62
RWKV-5	456	16.53
RWKV-6	442	17.40
RetNet	431	16.23
HGRN	411	21.83
GLA	412	19.56
HGRN2	411	16.77
xLSTM[1:0]	409	13.43
xLSTM[7:1]	408	13.48
xLSTM Limitations

- Memory mixing of the sLSTM prohibits parallelization Fast sLSTM CUDA kernel is 1.5 times slower than parallel mLSTM.
- CUDA kernels for mLSTM are not optimized (4 times slower than FlashAttention or Mamba's Scan). FlashAttention analog is possible
- Matrix memory of mLSTM is computationally expensive. Parallelization leads to only a minor overhead concerning the wall clock time
- Initialization of the forget gates must be chosen carefully
- Limited memory (*dxd*), but no problems for contexts up to 16k
- Neither architecture nor the hyperparameters are optimized

Improved efficiency through xLSTM: RNNs versus Transformer

- Currently large language models (LLMs) dominate the AI landscape
 - Based on Transformer's self-attention
 - Quadratic in sequence length (aka "context-size")
 - Prior to 2017: LSTM networks
 - ELMO: first large language model
 - Linear in sequence length



Improved efficiency through xLSTM: RNNs versus Transformer

- Currently large language models (LLMs) dominate the AI landscape
 - Based on Transformer's self-attention
 - Quadratic in sequence length (aka "context-size")
 - Prior to 2017: LSTM networks
 - ELMO: first large language model
 - Linear in sequence length



Abilities and limitations of Al Moravec's Paradoxon



- Image recognition, object recognition
- Image generation
- Language models
 - Writing texts
 - Semantic similarity
- Board games and computer-game-like tasks
- Control tasks

- Mathematical tasks, calculations
- Planning and reasoning
- Motoric tasks

 setting the table

 Understanding of the physical world

Incorporating world knowledge via Geometric Deep Learning



E(3)-equivariance built into deep learning architecture

Virtual nodes for improved learning

Sestak, F., Schneckenreiter, L., Brandstetter, J., Hochreiter, S., Mayr, A., & Klambauer, G. (2024). VN-EGNN: E (3)-Equivariant Graph Neural Networks with Virtual Nodes Enhance Protein Binding Site Identification. *arXiv preprint arXiv:2404.07194*.

Improving robustness through contrastive learning



Fürst, A., Rumetshofer, E., Lehner, J., Tran, V. T., Tang, F., Ramsauer, H., ... & Hochreiter, S. (2022). Cloob: Modern hopfield networks with infoloob outperform clip. *Advances in neural information processing systems*, *35*, 20450-20468.



Sanchez-Fernandez, A., Rumetshofer, E., Hochreiter, S., & Klambauer, G. (2023). CLOOME: contrastive learning unlocks bioimaging databases for queries with chemical structures. *Nature Communications*, *14*(1), 7339.

Improving robustness through contrastive learning



Sanchez-Fernandez, A., Rumetshofer, E., Hochreiter, S., & Klambauer, G. (2023). CLOOME: contrastive learning unlocks bioimaging databases for queries with chemical structures. *Nature Communications*, *14*(1), 7339.

Improving robustness through contrastive learning



Seidl, P., Vall, A., Hochreiter, S., & Klambauer, G. (2023, July). Enhancing activity prediction models in drug discovery with the ability to understand human language. In *International Conference on Machine Learning* (pp. 30458-30490). PMLR.

Improving adaptability through new few-shot learning approaches



Schimunek, J., Seidl, P., Friedrich, L., Kuhn, D., Rippmann, F., Hochreiter, S., & Klambauer, G. (2023). Context-enriched molecule representations improve few-shot drug discovery. *International Conference on Learning Representations*.

Drug Discovery

- Selection, generation and prediction of drug candidates
- Data sets:
 - >100,000,000 bioactivity triplets
- Methods:
 - Multi-task neural networks
 - Graph neural networks
 - Large language models
- Validation:
 - u.a. "Hit rate"
- Results:
 - Highly accurate prediction of bioactivities
 - Up to 100-fold increased hit-rates
 - Acceleration of pre-clinical phases of drug discovery



Zhavoronkov, A., Ivanenkov, Y. A., Aliper, A., Veselov, M. S., Aladinskiy, V. A., Aladinskaya, A. V., ... & Volkov, Y. (2019). Deep learning enables rapid identification of potent DDR1 kinase inhibitors. Nature Biotechnology, 1-4.

The rise of AI drug discovery

- Artificial intelligence will assist life sciences in multiple ways
 - Drug discovery and development
 - Personalized medicine
 - Medical imaging, diagnosis, and prognosis
 - Virtual assistants
 - o ...
- We strive for AI systems in life sciences that
 - can incorporate knowledge
 - are robust
 - efficient
 - and can adapt quickly
- through
 - advancing machine learning and deep learning

Drug Discovery

- Selection, generation and prediction of drug candidates
- Data sets:
 - >100,000,000 bioactivity triplets
- Methods:
 - Multi-task neural networks
 - Graph neural networks
 - Large language models
- Validation:
 - u.a. "Hit rate"
- Results:
 - Highly accurate prediction of bioactivities
 - Up to 100-fold increased hit-rates
 - Acceleration of pre-clinical phases of drug discovery



https://doi.org/10.1038/s41467-023-42328-w



CLOOME: contrastive learning unlocks bioimaging databases for queries with chemical structures

Received: 20 February 2023 Ana Sanchez-Fernandez¹, Elisabeth Rumetshofer¹, Sepp Hochreiter @^{1,2} & Günter Klambauer @¹

Sanchez-Fernandez, A., Rumetshofer, E., Hochreiter, S., & Klambauer, G. (2022). CLOOME: a new search engine unlocks bioimaging databases for queries with chemical structures.

Article

mLSTM: memory

Memory $oldsymbol{A} \in \mathbb{R}^{d imes d}$:

Memory capacity: 0.14 d.

Memory is independent of the sequence length L.

$$m{A} \;=\; egin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,d} \ a_{2,1} & a_{2,2} & \cdots & a_{2,d} \ dots & dots & \ddots & dots \ a_{d,1} & a_{d,2} & \cdots & a_{d,d} \end{pmatrix}$$

Parallelizability of mLSTM



122