Time-series, Transformers, and Natural Language Processing

Inar Timiryasov (NBI) inar.timiryasov@nbi.ku.dk



May 14, 2025







- The Attention Mechanism: A Leap Forward
- Transformers: Architecture & Impact
- Transformers in Natural Language Processing (NLP)



Time series data and limitations of the usual network architectures







- Time series data are ubiquitous: Stock prices; weather patterns; sensor readings; flight passenger numbers (today's exercise); speech; text
- Why is it special? Order matters, dependencies between points.
- Usual NN approaches are not well suited for these types of data:
 - Lack of temporal dynamics (treating input features independently)
 - Fixed input size
 - CNNs capture only local dependencies by design



- A network with memory (hidden state).
- Analogy: Reading a sentence, understanding depends on previous words. \bullet Compressed information about previous words is stored in the hidden state.
- We apply the same model to each time step (unfold)! \bullet
- We can only process one element at a time this is a limitation \bullet



Recurrent Neural Network (RNN)

 h_t - hidden state

$$a_{t} = Vh_{t-1} + Ux_{t} + b$$
$$h_{t} = \tanh(a_{t})$$
$$o_{t} = Wh_{t} + c$$

Challenges with Simple RNNs & Evolution

- Training RNNs: backpropagation through time. ullet
- Each time step of the unrolled recurrent neural network may be seen as an additional layer.
- Vanishing/Exploding Gradients: Difficulty learning long-range dependencies. \bullet
- A fix is a "skip connection" (similar to U-Net and ResNet, but in time, not in depth) ullet
- LSTM, GRU \bullet







Long Short Term Memory (LSTM) and Gated Recurrent Unit (GRU)

- Skip connection long term memory a "highway" in time ullet
- Gates decide what to write to the long term memory \bullet
- They are better at capturing longer-term patterns.



LSTM

Sources: http://colah.github.io/posts/2015-08-Understanding-LSTMs/ https://web.stanford.edu/class/cs224n/slides/cs224n-2023-lecture06-fancy-rnn.pdf



GRU

Exercise: Flight Passenger Predictions



Recap on RNNs

- The same network is used at every time step (hence recurrent).
- Vanilla RNNs are hard to train. LSTM/GRU partly solve this problem.
- Intuition: there are exponentially many almost orthogonal vectors in D>>1 dimensions. Cons: in practice, only a limited amount of information can be stored in the hidden state. RNNs tend to forget the past context after a certain length.
- Hard to train in parallel, since we need to process steps sequentially.

• All information about the past sequence must be compressed into the hidden state. Pros: compression is intelligence [see "An Observation on Generalization" by Ilya Sutskever].

(but there are ways around, keywords: RWKV, State Space Models, e.g. Mamba.)



- Core Idea: Allow the model to "look back" at \bullet different parts of the input sequence when processing or generating an output, and weigh their importance.
- Analogy: Human translators don't just read a \bullet whole sentence and then translate; they focus on relevant parts.
- Benefit: Mitigates information bottleneck, better \bullet context.
- Self-Attention: Attention within a Single lacksquareSequence

Attention!



Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio Neural Machine Translation by Jointly Learning to Align and Translate arXiv:1409.0473

RNN with attention





How Self-Attention works?



How Self-Attention works?

- Input x_s at time step s
- We project each input x_s using three different transformations: $v_s = W^V x_s$ $k_s = W^K x_s$ $q_s = W^Q x_s$
- Output at step t is the sum of v_s over all steps s with weights A_{ts} :

$$y_t = \sum_s A_{ts} v_s$$

where attention weights $A_{ts} = \operatorname{softmax}\left(\frac{q_t \cdot k_s}{\sqrt{d}}\right)$

$$y_t = \sum_{s} \operatorname{softmax}\left(\frac{(W^Q x_t) \cdot (W^K x_s)^T}{\sqrt{d}}\right) (W^V x_s)$$

Scaled Dot-Product Attention



How Self-Attention works?



 $y_t = \sum_{s} \operatorname{softmax}\left(\frac{(W^Q x_t) \cdot (W^K x_s)^T}{\sqrt{d_K}}\right) (W^V x_s)$



attention becomes non-trivial with multiple layers see [Transformer Circuits Thread]



- Initially proposed by Google for machine translation tasks. \bullet
- They utilize the attention mechanism introduced in [1409.0473], but \bullet without an RNN component.
- The original architecture consisted of an encoder and a decoder: lacksquare
 - **Encoder**: Processes input sequence (e.g., sentence to be translated). Each position attends to all other positions in the input sequence.

- **Decoder**: Generates output sequence (e.g., the translation), attending to encoder output.

Each position attends to previous positions (including itself) in the output sequence (causal attention).

Generative models, such as GPT, are decoder-only. \bullet In physics, we often need encoder models. (Technically, the only difference: decoders feature a causal attention mask.)

Transformer





- Now that diagram is not so scary!
- Stack multiple attention heads.
- Add MLP. \bullet
- Add norms (before) Attention and ulletMLP.
- Repeat n_layer times. ullet
- Positional embeddings to break \bullet permutation invariance. (nowadays: RoPE)

Transformer



- Superior performance on many sequence tasks, especially NLP. Also on vision. \bullet
- Handles long-range dependencies effectively. \bullet Google's Gemini has context length **1M**, see at <u>https://aistudio.google.com</u>
- Parallelization: Faster training on modern hardware. \bullet Computations for all tokens within a layer can be performed in parallel, unlike RNN.
- The price to pay is $\mathcal{O}(N^2)$ computational and memory complexity with respect to sequence length N. Mitigations:
 - Efficient implementations: FlashAttention
 - KV Caching for generation.
- Scalability: Foundation for very large models (LLMs). \bullet

Why Transformers are a Breakthrough

Attention is all you need

Attention is all you n <u>A Vaswani, N Shazeer, N</u> The dominant sequence to

May 2023

<u>A Vaswani, N Shazeer, N Parmar</u>... - Advances in neural ..., 2017 - proceedings.neurips.cc ... to attend to **all** positions in the decoder up to and including that position. We need to prevent ... We implement this inside of scaled dot-product attention by masking out (setting to -∞) ... ☆ Save 57 Cite Cited by 179358 Related articles All 73 versions ≫

echanisms. We propose a novel, simple network architecture based solely onan attention mechanism, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superiorin quality while being more ... c Save $\overline{99}$ Cite Cited by 74101 Related articles All 46 versions \gg

May 2025





My slide from 2023 Hard to impress in 2025!

- Goal: Enabling computers to understand, interpret, and generate human language.
 - Text Classification (e.g., sentiment analysis)
 - Machine Translation
 - Text Generation <- today's focus Also known as language modeling, autoregressive generation

Language modeling: why does it even work?

- Given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$ estimate $P(x^{(t+1)} | x^{(1)}, x^{(2)}, ..., x^{(t)})$
- N-gram models: estimate $P(x^{(N)} | x^{(1)}, x^{(2)}, ..., x^{(N-1)})$ from a corpus of texts by simple counting
- for a vocabulary of 50 000 words:
- 1.25×10^{14} trigrams
- 6.1×10^{65} 14-grams





0.8 0.0

Representing words

- Vocabulary: enumerate all words But there are too many words in many languages
- Tokenization: Breaking text into smaller units (words, subwords) Note! Tokenization causes its own issues [very detailed lecture]

lokens	Characters	
220	747	



A helpful rule of thumb is that one token generally corresponds to ~4 characters of text for common English text. This translates to roughly 3/4 of a word (so 100 tokens ~= 75 words).

Tokens	Characters
220	747

[20916, 604, 357, 3103, 85, 2122, 282, 47986, 27862, 357, 18474, 82, 828, 3311, 6657, 47986, 27862, 357, 49, 6144, 82, 828, 290, 11160, 12, 4834, 19815, 364, 357, 14242, 8, 2599, 220, 198, 6747, 1315, 25, 1511, 25, 1314, 12, 1558, 25, 405, 25, 34872, 2122, 282, 47986, 27862, 357, 18474, 82, 8, 290, 2939, 3781, 357, 19962, 337, 700, 1531, 737, 198, 220, 220, 220, 220, 32900, 25, 31517, 1096, 4263, 357, 39764, 8808, 27039, 11, 29877, 12014, 329, 11881, 11, 290, 14, 273, 35831, 84, 2977, 422, 30155, 4771, 21758, 8, 351, 257, 8100, 13, 198, 6747, 1596, 25, 860, 25, 1314, 12, 1065, 25, 405, 25, 3311, 6657, 47986, 27862, 357, 49, 6144, 828, 5882, 10073, 35118, 14059, 357, 43, 2257, 44, 8, 290, 12068, 15417, 28403, 357, 45, 19930, 8, 357, 818, 283, 5045, 9045, 292, 709, 737, 198, 220, 220, 220, 220, 32900, 25, 5765, 281, 406, 2257, 44, 284, 4331, 5474, 4979, 290, 466, 12068, 15417, 28403, 319, 8959, 11012, 3807, 8088, 13, 198, 6747, 1596, 25, 1511, 25, 1314, 12, 1558, 25, 405, 25, 357, 23907, 864, 8, 11160, 12, 27195, 12342, 290, 32172, 13326, 357, 7250, 737, 198, 220, 220, 220, 220, 32900, 25, 3082, 601, 4263, 1262, 11160, 12, 27195, 12342, 11, 290, 13946, 41270, 2272, 351, 471, 33767, 13]

TEXT TOKEN IDS

https://platform.openai.com/tokenizer

Embeddings — every token is a vector in a \bullet multidimensional space (Word2Vec)



Operations over vectors: king - man + woman ~= queen

Italy German

Language modeling: generic picture

- Process the first *t* words
- Predict the probability of the next word $P(x^{(t+1)} | x^{(1)}, x^{(2)}, ..., x^{(t)})$
- Sample the next word

. . .

• Process the first t + 1 words





• Self-supervised pretraining

 \bullet

• Output probabilities over all tokens (words)

$$P_{\alpha} = softmax(\ell_{\alpha}) = \frac{\exp(\ell_{\alpha})}{\sum_{\beta} \exp(\ell_{\beta})}, \text{ where } \ell_{\alpha} - \text{ logits (raw})$$

Cross-Entropy Loss = $-\sum_{\alpha} \sum_{\alpha} y_{\alpha} \log(P_{\alpha}), \text{ where } y_{\alpha} \text{ are transformed to the set of the$

• Typically trained with very large batch sizes (millions of tokens). Highly parallelizable.

 $s \alpha$

• Adam optimizer (or similar). SGD doesn't work for transformers!



outputs)

rue next tokens (one-hot encoded).



nanoGPT exercise!

Success of self-supervise training

- Labeled data is limited lacksquare
- Unlabeled data is abundant lacksquare(text, image, video)
- Led to GenAl revolution \bullet
- Some recent models were trained on 30 \bullet trillion tokens!



Limits of LLM scaling based on human-generated data"

BERT - 3.3B tokens

1810.04805 "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding"

Self-supervise training: Scaling Laws

Performance predictably improves with scale



bottlenecked by the other two.

Figure 1 Language modeling performance improves smoothly as we increase the model size, datasetset size, and amount of compute² used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not

> https://arxiv.org/pdf/2001.08361 Scaling Laws for Neural Language Models



Success of self-supervise training

Power

Probably the single biggest constraint on the supply-side will be power. Already, at nearer-term scales (1GW/2026 and especially 10GW/2028), power has become the binding constraint: there simply isn't much spare capacity, and power contracts are usually long-term locked-in. And building, say, a new gigawatt-class nuclear power plant takes a decade. (I'll wonder when we'll start seeing things like tech companies buying aluminum smelting companies for their gigawatt-class power contracts.⁵⁷)

https://situational-awareness.ai/ Leopold Aschenbrenner, June 2024

MICROSOFT / TECH / SCIENCE

Microsoft wants Three Mile Island to fuel its AI power needs



/ Microsoft has signed a 20-year deal to exclusively access 835 megawatts of energy from a nuclear plant.

By Tom Warren, a senior editor and author of Notepad, who has been covering all things Microsoft, PC, and tech for over 20 years. Sep 20, 2024 at 2:23 PM GMT+2

69 Comments (69 New)

Photo by Andrew Caballero-Reynolds / AFP via Getty Images



https://www.theverge.com/2024/9/20/24249770/

8)(f)





- Open weights example: **DeepSeek R1**, 671B parameters ullet
- 671B parameters × 2 bytes per parameter lacksquare
- For maximum context (32,092 tokens): KV cache requires around 280GB RAM ullet
- Total VRAM Needed: ~1,680GB ullet
- H100 GPU Count Calculation: 1,680GB ÷ 80GB per H100 = 21 GPUs (minimum) ulletIncluding overhead for system operations, tensor parallelism, and other processes: 22-24 H100 GPUs
- Total Estimated Cost: 22 H100s at approximately 30,970 each = \sim 681,340ullet

- Try it yourself! <u>https://huggingface.co/deepseek-ai/DeepSeek-R1</u>
- You can try smaller models locally! <u>https://huggingface.co/Qwen/Qwen3-0.6B</u>

How large are Large Language Models?

- Self-supervised training is the first step.
- Instruction Fine-Tuning (IFT / SFT with Instructions):
 - commands.
 - Fine-tune on a dataset of (instruction, response) pairs.
 - Teaches the model to respond to prompts in a desired format and style.
 - Example data: "Translate this sentence to French: {sentence}" -> "{French translation}"
- Alignment with Human Preferences: ullet
 - Reinforcement Learning from Human Feedback (RLHF)
 - Direct Preference Optimization (DPO)
- Reasoning (RL in verifiable domains, e.g GRPO)

- A crucial step for making models more "helpful and harmless" and better at following user

see https://rlhfbook.com/



Thank you for your Attention



Source https://transformer-circuits.pub/2021/framework/index.html

How it works

The OV ("output-value") circuit determines how attending to a given token affects the logits.

 $W_U W_O W_V W_E$

The QK ("query-key") circuit controls which tokens the head prefers

 $W_E^T W_Q^T W_K W_E$

How it works: attention patterns and activations



Demo: https://huggingface.co/spaces/simon-clmtd/exbert

and a short of the second s	ndiger (for for for	
		and the second secon
	u panti lanta na kant	
	ny linita se	
and descate high end and or democratics	maintenanti	
		alesses have a such the
a del un an esta an ortra cal mar estad	المتحد والمتحدثة ومتعاله	والمعتقل والمستعمل الم
han a dhallain fi shada ann an 1955 dheann		h in other sectors and in particular
	la de la completa de	
		i i i i i i i i i i i i i i i i i i i
And advection of protocol structures are not financial		
	Aprile de parte	ويتباح أناجه بالألي
		ile disease in the life is a
		Julium ike titue.
Von <mark>ar estatut banan ana ana dana ana</mark>		and a set of the set of
	And the second of the second of the second	an financia calatica
endifinence of the second state		
, hith as a she tet above a heter e art at	l	ndan its tit.
and the second se		
		And the state of the
	a providencia de la constante d	and the state of the second
and a second		and the state of the
		ومألك ويتعرف والم
to be the second of the second second second second		وهالي الجمع ا
The area of a standard surface of a standard stand		a na mina ang paga at a
ep <mark>elesienen seljennen beinen erekenen selven serven</mark>		
nalitation of attendance of the second	التر بتلميتينين	
Lingen Mannen og det som		
	والبرج والمرجوبا و	
	A the second	energia de la companya de la company
n parlasta lagana alia ilijan sata alia alim sinata di pata pana ina ing		
eyhdelitar einder		999.866.999.866.999.999.999.999.999.999.

Time Series

V

Predicting the next value

Global properties of time series













https://www.kaggle.com/competitions/icecube-neutrinos-in-deep-ice/overview

29