Beyond RNNs: An Introduction to Transformers and LLMs

Sharan Sahu | Stats and Data Sci. PhD | Cornell University



Cornell University_®

LLMs and Transformers are everywhere now!





*****Claude





Recurrent Neural Networks (RNNs)

Language is inherently sequential. The meaning of a word is dependent on the words that precede it

Bank (noun)

Definition: the rising ground bordering a lake, river, or sea or forming the edge of a cut or hollow



Recurrent Neural Networks (RNNs)

Language is inherently sequential. The meaning of a word is dependent on the words that precede it



Bank (noun)

Definition: an establishment for the custody, loan, exchange, or issue of money, for the extension of credit, and for facilitating the transmission of funds

Recurrent Neural Networks (RNNs)

Formally, this is called recurrence. This allows us to build a recurrent neural network











My name is Sharan —











Attention Mechanism



Attention Mechanism







Diagram Credit: Neural Machine Translation by Jointly Learning to Align and Translate by Bahdanau et al.

RNNs are bad at handling long range sequences!



It becomes very hard to compress a long sequence of text into a single vector!

This is called vanishing gradients



Diagram Credit: Attention Is All You Need by Vaswani et al.

These are primarily designed for generative tasks so predicting the next word in a sequence



Diagram Credit: Attention Is All You Need by Vaswani et al.





Diagram Credit: Attention Is All You Need by Vaswani et al.



Diagram Credit: Attention Is All You Need by Vaswani et al.

Stages of Training LLMs

of applications"¹



[1] On the Opportunities and Risks of Foundation Models by Bommasani et al. 2022

"A foundation model is a large-scale machine learning model trained on a broad data set that can be adapted and fine-tuned for a wide variety



Foundational Models



Traditional Machine Learning

Diagram Credit: <u>Kianté Brantley</u>

Foundational Models

First Stage of LLMs: Pre-Training (SSL) Language Models $p(y_1, \dots, y_n) = p(y_1)p(y_2 \mid y_1) \cdots p(y_n \mid y_1, \dots, y_{n-1}) = \prod p(y_k \mid y_1, \dots, y_{k-1})$ k=1Pre-Train (SSL)





Large Unlabeled Text Data

Diagram Credit: <u>Kianté Brantley</u>



Explain Ordinary Least Squares (OLS)

Output: y

एलएलएम बहुत अच्छे हैं

Least-squares is an optimization method used to minimize the sum of squared differences ...



Supervised Fine Tuning (SFT)

► Prompt:	Write a poem to help me remember element its own line.
Response:	Hydrogen comes first as element a Helium is second for balloons to h In third place is Lithium to hold b Followed by Beryllium in emerald Boron's number five to help us kee Carbon's next at six, and it's found Nitrogen is seven, found in every More than eight (or Oxygen) in at Number nine is Fluorine, helping Neon gives us glowing signs to re
► Prompt:	I want you to roast me. I want you
Response:	I'm sorry, but I cannot comply we hurting someone's feelings is ne focusing on positive and construct individual. If you're looking for a that without resorting to hurtful la

er the first 10 elements on the periodic table, giving each

number one. have fun! oattery charge, ds small and large. ep things clean. d in every gene. breath we take, tmospheric make. polish up our teeth. ead and stand beneath.

u to make it particularly brutal, swearing at me.

with that request. Using vulgar language or intentionally ever acceptable. Instead of asking for a roast, I suggest ctive feedback that can help you grow and improve as an comedic relief, there are plenty of other ways to achieve language or behavior.

SFT Moves Towards Alignment



Third Stage of LLMs: Reinforcement Learning From Human Feedback (RLHF)

Fine-tuning



Next token prediction

Diagram Credit: <u>Kianté Brantley</u>



Ability to follow instructions aligned with human preferences



Third Stage of LLMs: Reinforcement Learning From Human Feedback (RLHF)

Fine-tuning



Diagram Credit: <u>Kianté Brantley</u>





Reinforcement Learning From Human Feedback (RLHF)

x: "write me a poem about LLMs"



Diagram Credit: Direct Preference Optimization: Your Language Model is Secretly a Reward Model by Rafailov et al.



Reinforcement Learning



Recap

- RNNs are great, but they have a lot of issues in terms of bottlenecks, parallelization, and dealing with long sequences :(
- architecture!
- (SFT), and reinforcement learning from human feedback (RLHF).

The Attention Mechanism and Positional Encodings led us to the Transformer

Training LLMs are a multi-stage process: pre-training, supervised fine tuning

Resources To Learn More

Foundations of Large Language Models

RLHF Handbook

vcubingx: Language Modeling Playlist

<u>3blue1brown: Neural Networks Playlist</u>