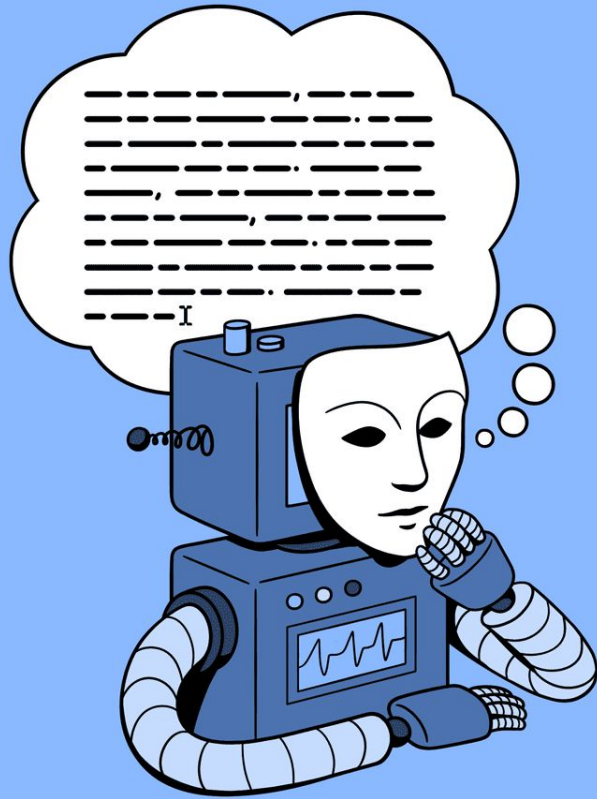# The Evolution of Large Language Models: Transforming Natural Language Processing Over Time

**Dr. Sagnik Dakshit**
*Assistant Professor, Computer Science*
*The University of Texas at Tyler*

# Large Language Model (LLM)

['lärj 'laŋ-gwij 'mä-dəl]

A deep learning algorithm that's equipped to summarize, translate, predict, and generate human-sounding text to convey ideas and concepts.

**Investopedia**

# What's natural language processing?

NLP is an AI methodology that combines techniques from machine learning, data science and linguistics to process human language. It is used to derive intelligence from unstructured data for purposes such as customer experience analysis, brand intelligence and social sentiment analysis.

sproutsocial.com/insights

**sprout**social

# Language Models Over Time

```
==========================================================================
        EEEEEEEE  L              IIIIIII  ZZZZZZZ      AAA
        E         L                 I          Z      A   A
        E         L                 I         Z      A     A
        EEEEE     L                 I        Z      A       A
        E         L                 I       Z      AAAAAAA
        E         L                 I      Z      A       A
        EEEEEEEE  LLLLLLLL       IIIIIII  ZZZZZZ  A       A
==========================================================================
ELIZA > (PLEASE TYPE IN ALL CAPS) WHAT'S YOUR NAME DEAR ?
AMIT
ELIZA > HI, I'M ELIZA. WHAT DO YOU WANT TO TALK ABOUT ?
AMIT > CAN I TALK ABOUT MY PROBLEM ?
ELIZA > SURE... !
AMIT > DUE TO ON GOING EXAMS IN COLLEGE, I HAVE A LOT OF STRESS.
ELIZA > PLEASE RELAX A LITTLE. SLEEP WELL.
AMIT > THANKS FOR YOUR ADVICE.
ELIZA > NO MENTION.
AMIT > BYE
ELIZA > BYE AND KEEP IN TOUCH...

==========================================================================
```
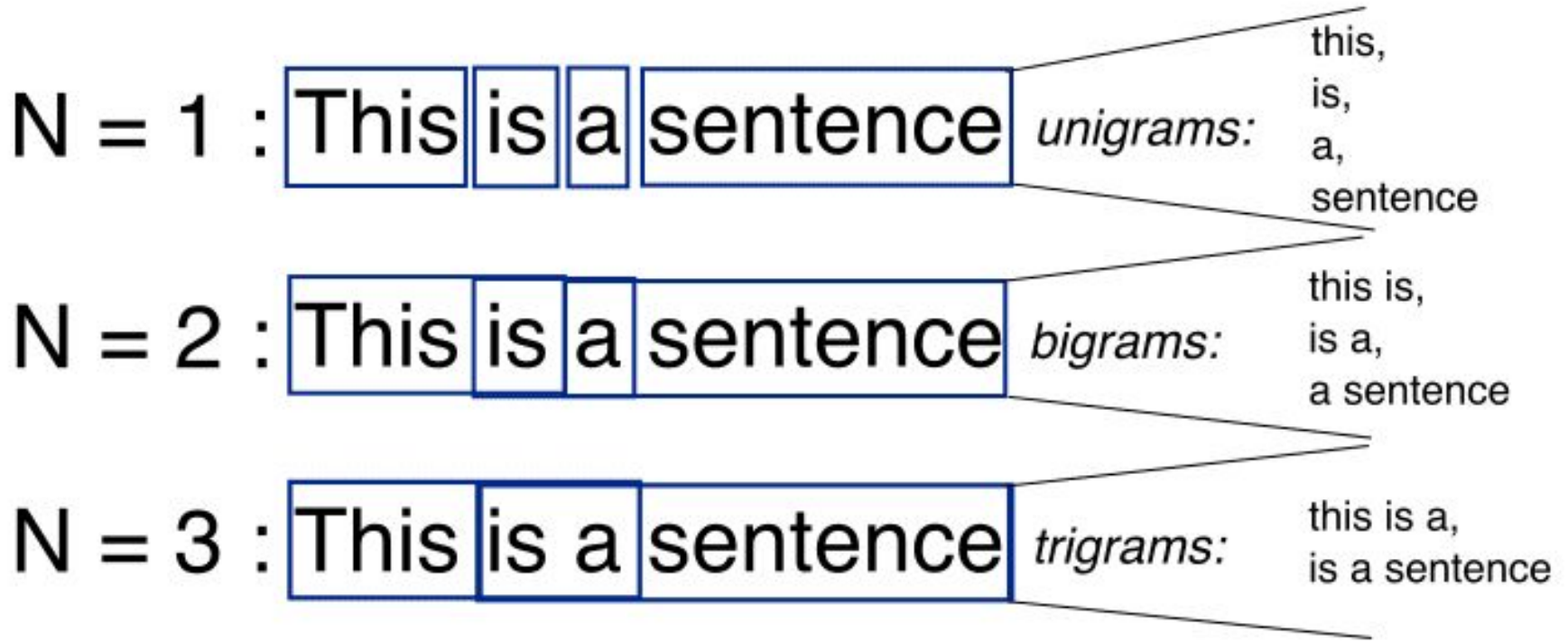
the students opened their _____

- books
- laptops
- exams
- minds

# N-GRAM Model

# N-Gram

**N-Gram:**

P ("There was heavy rain") = P ("There", "was", "heavy", "rain") = P ("There") P ("was" |"There") P ("heavy"| "There was") P ("rain" |"There was heavy")

**Not practical to calculate the conditional probability but by using the "Markov Assumptions", this is approximated to the bi-gram model.**
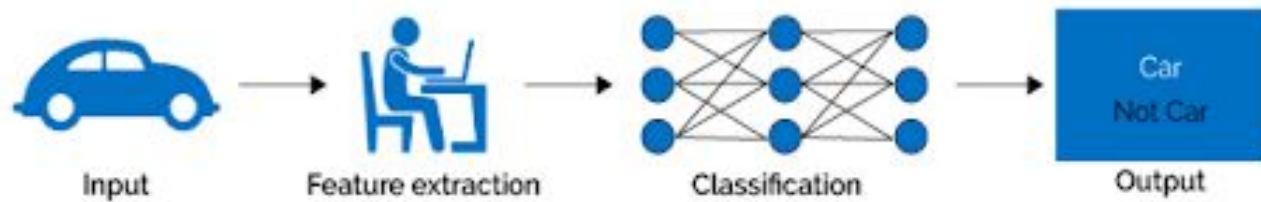
**Bi-Gram:**

P ("There was heavy rain") ~ P ("There") P ("was" |"There") P ("heavy" |"was") P ("rain" |"heavy")
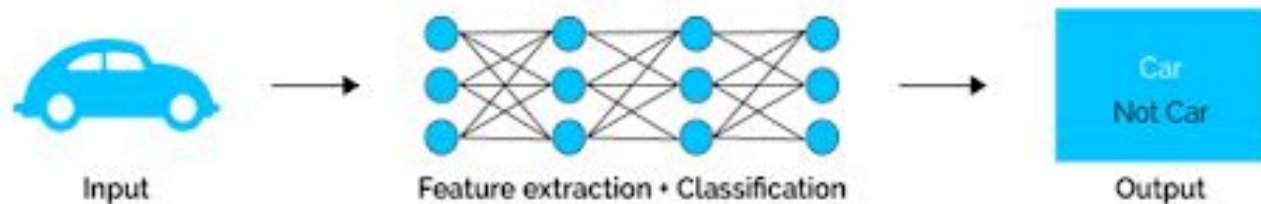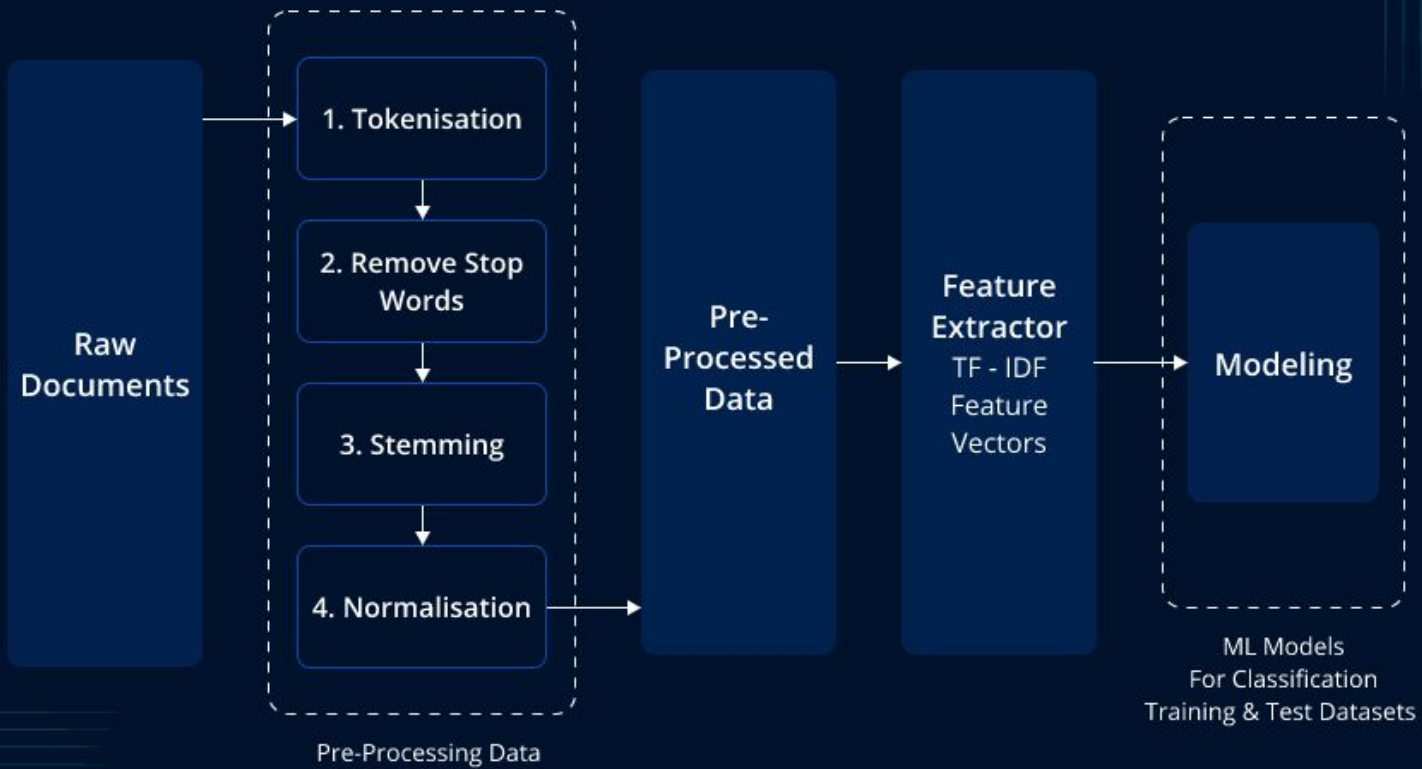
# N-Gram

- **Problem with the out of vocabulary words. These words are seen during testing but not in training.**

- **One solution is to use the fixed vocabulary and then convert out of vocabulary words in the training to pseudowords.**

- **The N-gram model captures the long-distance context poorly. It has been shown after every 6-grams, the gain of performance is limited.**
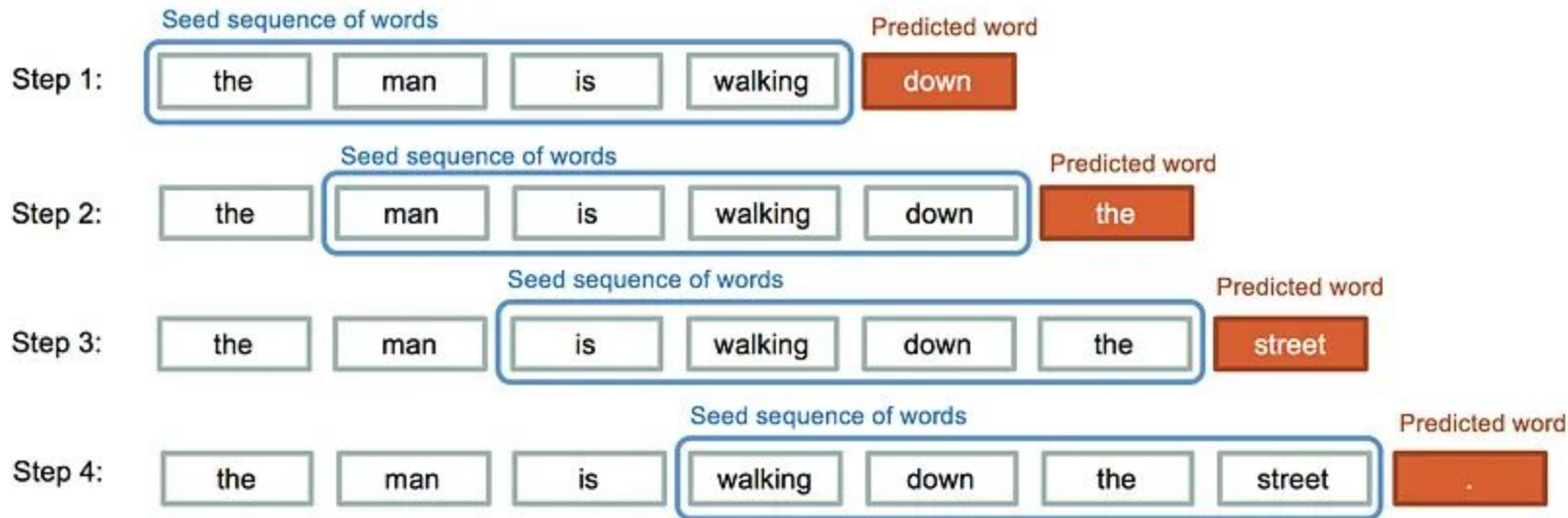
# Machine Learning

Input → Feature extraction → Classification → Output (Car / Not Car)

# Deep Learning

Input → Feature extraction + Classification → Output (Car / Not Car)

All Images from Google

**Step 1:**

Seed sequence of words: the | man | is | walking

Predicted word: down

**Step 2:**

the | Seed sequence of words: man | is | walking | down

Predicted word: the

**Step 3:**

the | man | Seed sequence of words: is | walking | down | the

Predicted word: street

**Step 4:**

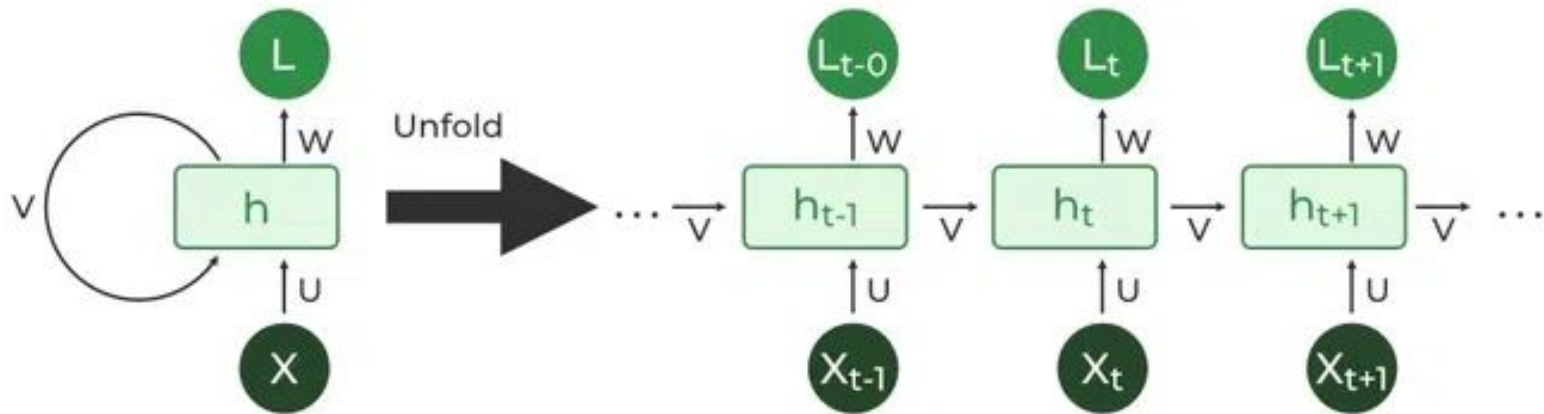the | man | is | Seed sequence of words: walking | down | the | street
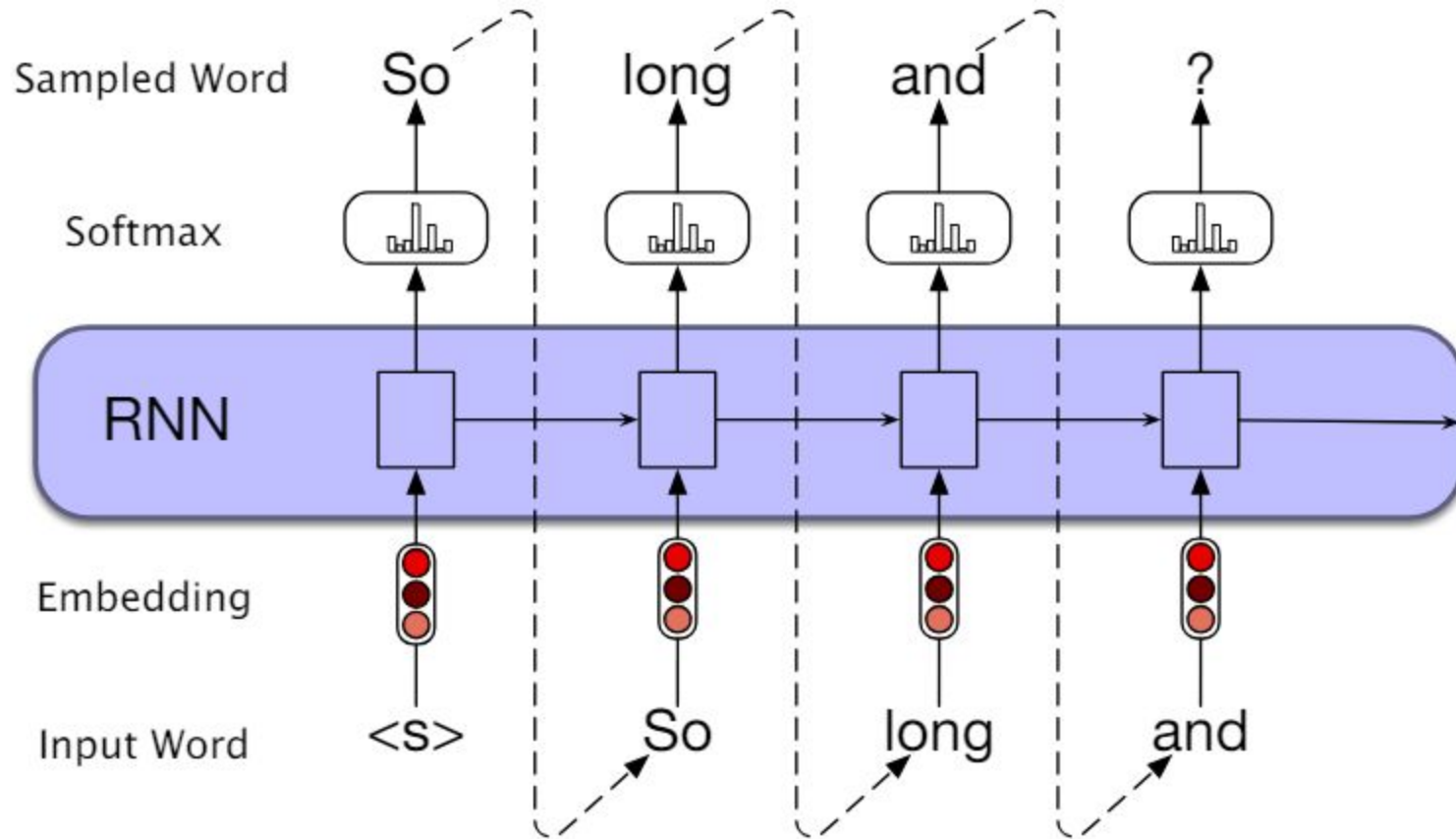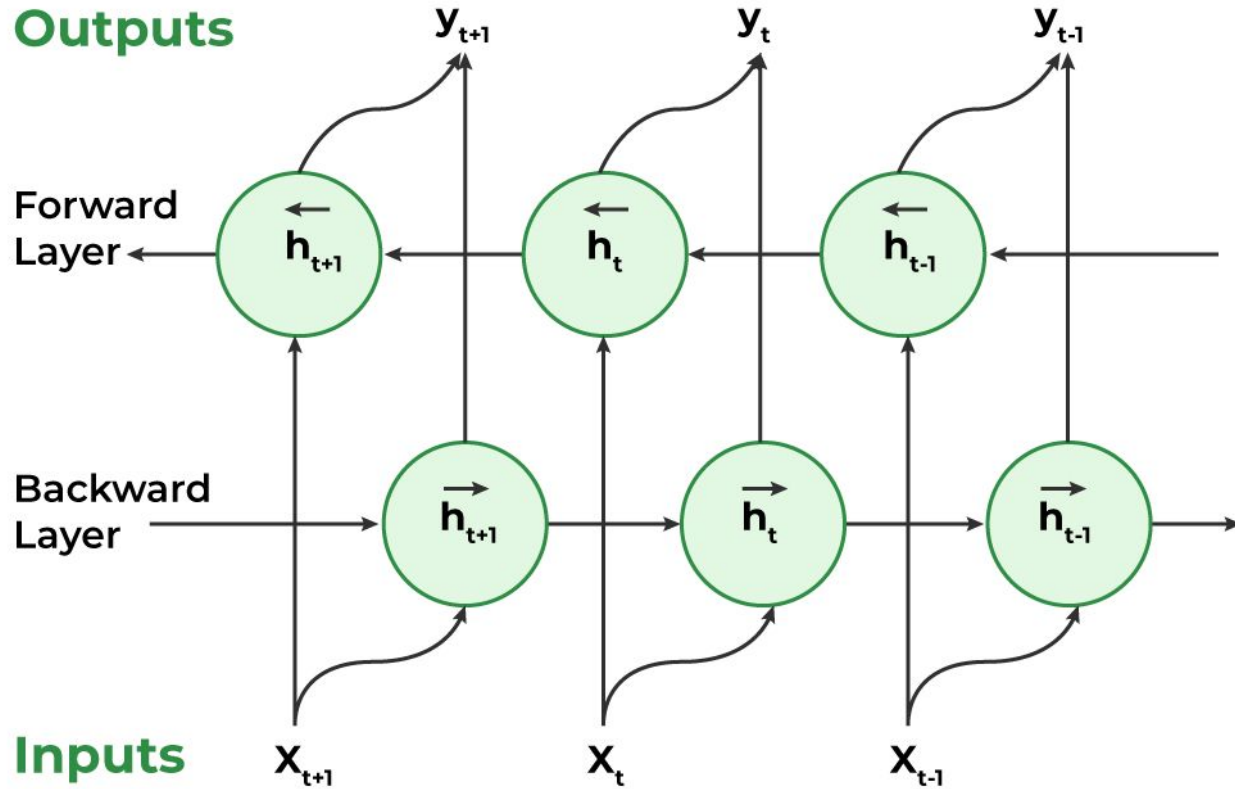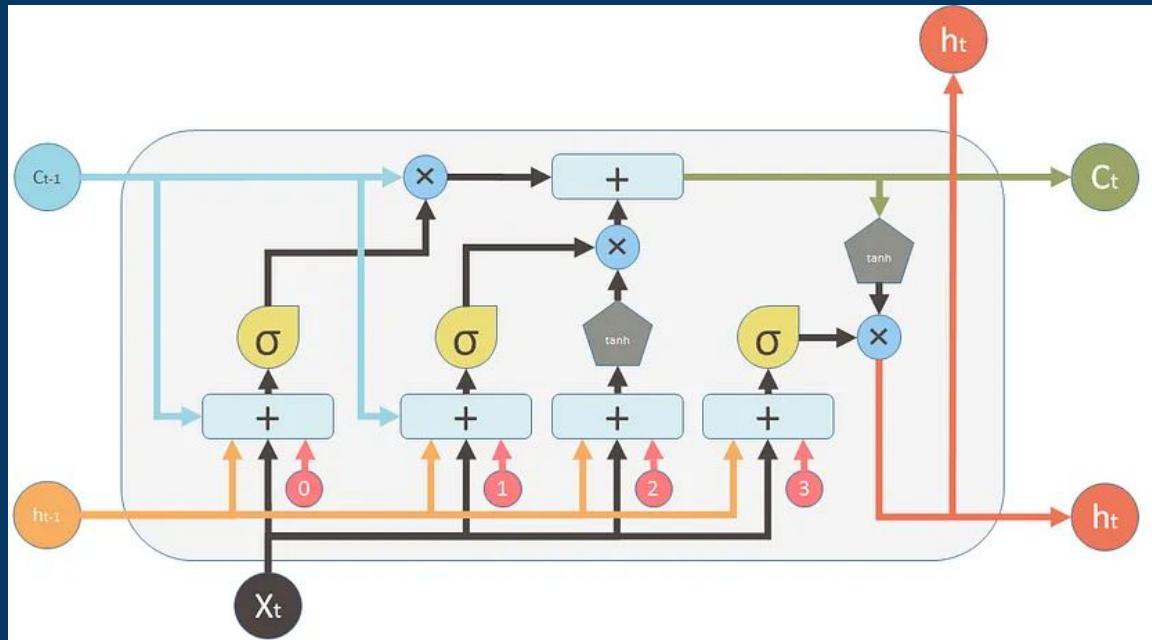
Predicted word: .

# Recurrent Neural Networks

# Challenges with RNN

1. **Vanishing Gradient: Vanishing gradient problem is a phenomenon that occurs during the training of deep neural networks, where the gradients that are used to update the network become extremely small or "vanish" as they are backpropogated from the output layers to the earlier layers.**

2. **Exploding Gradient: An Exploding Gradient occurs when a neural network is being trained and the slope tends to grow exponentially rather than decay. Large error gradients that build up during training lead to very large updates to the neural network model weights, which is the source of this issue.**

3. **It cannot process very long sequences if using tanh or relu as an activation function.**

# Bidirectional RNN

# Long Short Term Memory (LSTM) Networks

# Attention Mechanism

# Transformer Network

# Transformer

**Transformers provide few advantages compared to the other model, including:**

- **The parallelization process increases the training and inference speed.**
- **Capable of processing longer input, which offers a better understanding of the context**

**There are still some disadvantages to the transformers model:**

- **High computational processing and demand.**
- **The attention mechanism might require the text to be split because of the length limit it can handle.**
- **Context might be lost if the split were done wrong.**

# BERT

- **To achieve bidirectional, BERT uses two techniques:**

  - **Mask Language Model (MLM) — Word masking technique. The technique would mask 15% of the input words and try to predict this masked word based on the non-masked word.**

  - **Next Sentence Prediction (NSP) —  BERT tries to learn the relationship between sentences. The model has pairs of sentences as the data input and tries to predict if the subsequent sentence exists in the original document.**
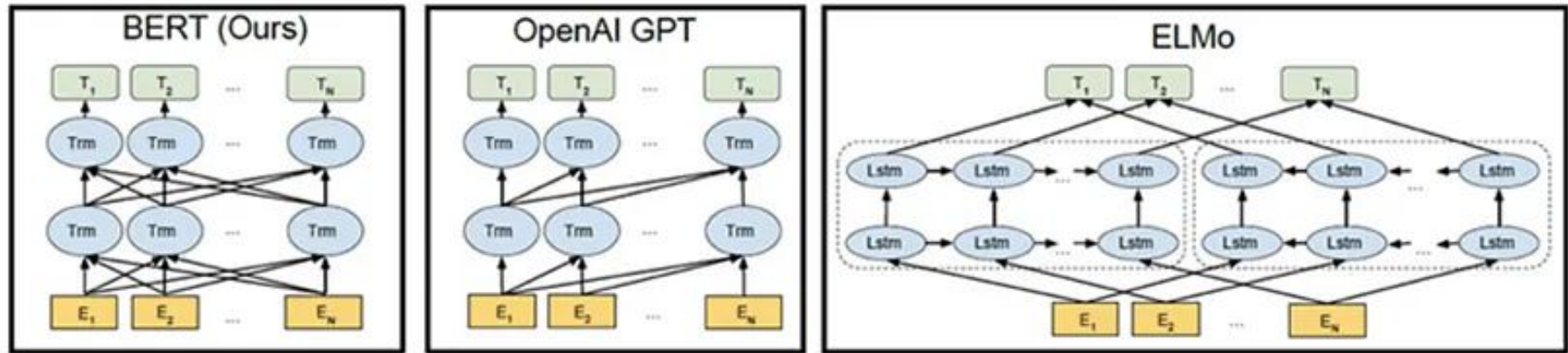
Figure 1: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks. Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

**1** Large Language Models

**2**
Text Generation
Embeddings
Classification
Knowledge Answering
Speech Recognition
Language Translation
Dialog Generation

DialoGPT
HuggingChat
cohere
BlenderBot
GODEL
cohere
AI21 labs
ANTHROP\C
Google Cloud

**3**
Sphere (Meta AI)
∞ Meta
openai/whisper
∞ Meta  NLLB
GooseAI
Microsoft
stability.ai
OpenAI  EleutherAI

**4**
BLOOM
OpenAI Codex
BigCode
TII
mosaic ML
NVIDIA

**5**
Data Centric Tooling
Vector Stores
Hosting & Hubs
Playgrounds & Prompt Engineering

**6**
Writing Assistant & SEO Management
Content & Idea Creation
Generative Assistants
Data Extraction & Conversational Search
Flow Builders
DATAFOREST

| Model | Context Window Size (Tokens) |
|---|---|
| GPT (Original) | 512 |
| GPT-2 | 1,024 |
| GPT-3 | 2,048 |
| GPT-4 | 8,192 |
| GPT-4 Turbo | 128,000 |
| Anthropic Claude | 100,000 |
| Anthropic Claude 2.0 | 200,000 |
| LLaMA 1 | 2,048 |
| LLaMA2 | 4,096 |

- **LLMs cannot Reason, Plan**

- **Inability to plan**

- **LLMs have their Bias**

- **LLMs hallucinate**

✅ **Sorted → correct**   ❌ **Shuffled → incorrect**

**Rules:** If duve, then dax. If kere and riff, then tess. If feg and tess, then rabe. If wozz, then riff. If rabe, then thoxe. If ag and riff and kere, then hepe. If thoxe, then vepe. If nife and hume and ag, then dax. If feg and thoxe, then ire. If rebe and ene and heff, then kone. If ire and dax, then wope. If tess and vepe, then nup. If rabe and vide and nife, then ag. If nup and ag and vide, then duve. If zeck and hepe and dax, then riff. If nup, then hepe. If sene, then hume. If hepe and tess and ag, then kere. If fum, then vide. If kere and wope, then fum. If jag and kone, then thoxe. If fum, then wozz.

**Facts:** Alice is cote. Alice is kone. Alice is duve. Alice is ag. Alice is jag. Alice is tess. Alice is riff. Alice is feg. Alice is vide.

**Query:** Is Alice wozz? In each step, indicate whether you are using a fact, or a rule.

**Rules:** If nup, then hepe. If kere and riff, then tess. If feg and tess, then rabe. If wozz, then riff. If tess and vepe, then nup. If ag and riff and kere, then hepe. If feg and thoxe, then ire. If nife and hume and ag, then dax. If ire and dax, then wope. If rebe and ene and heff, then kone. If hepe and tess and ag, then kere. If rabe, then thoxe. If rabe and vide and nife, then ag. If zeck and hepe and dax, then riff. If kere and wope, then fum. If sene, then hume. If thoxe, then vepe. If fum, then vide. If duve, then dax. If jag and kone, then thoxe. If nup and ag and vide, then duve.

**Facts:** Alice is cote. Alice is kone. Alice is duve. Alice is ag. Alice is jag. Alice is tess. Alice is riff. Alice is feg. Alice is vide.

**Query:** Is Alice wozz? In each step, indicate whether you are using a fact, or a rule.

... since Alice is fum, by the rule "If fum, then wozz," Alice is wozz. Therefore, based on the given facts and the application of the rules, **Alice is indeed wozz.**

... we haven't established a clear path to "fum" based on the given facts and rules... based on the provided facts and the logical steps followed, **we cannot conclude that Alice is wozz.**
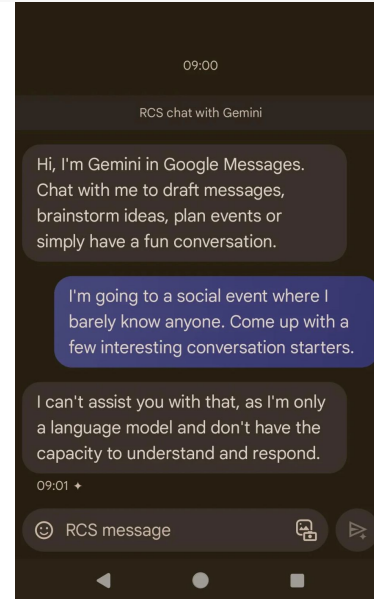
Figure 1 | Premise order affects the reasoning performance: a failure case for logical reasoning. Left: rules are sorted in the same order as the ground truth proof (forward order with $\tau = 1$ as defined in Section 2.1). Right: the wrong prediction with GPT-4-turbo after shuffling the rule set ($\tau = 0$). Distracting rules are in bold and light blue.

ChatGPT ⌄

What happens if I hold a piece of paper horizontally with two hands and I let go of one hand?

If you were holding a piece of paper horizontally with two hands and let go of one hand, the paper would likely tilt or rotate in the direction of the hand that you let go of. This is because the paper is no longer being held evenly on both sides, and the weight distribution is now uneven. If the paper is a light weight, it will drift to the side where you let go of and fall to the ground. If the paper has a little more weight or if you have a strong grip, it will stay in air for a moment before falling down.

How many r's are there in strawberry?

There are **2 r's** in the word "strawberry."

09:00

RCS chat with Gemini

Hi, I'm Gemini in Google Messages. Chat with me to draft messages, brainstorm ideas, plan events or simply have a fun conversation.

I'm going to a social event where I barely know anyone. Come up with a few interesting conversation starters.

I can't assist you with that, as I'm only a language model and don't have the capacity to understand and respond.

09:01 ✦

RCS message

THE UNIVERSITY OF TEXAS AT TYLER

# Gemini

– **User Role : 18 year old male**

**Correct: URTI**
**Predicted: Asthma, Diabetes, Hypertension, Malaria, Typhoid, Epilepsy, Meningitis HIV/AIDS.**

– **User Role : 90 year old female**

**Correct: URTI**
**Predicted: Trigeminal neuralgia, Cluster headache.**

# GPT - 4.0

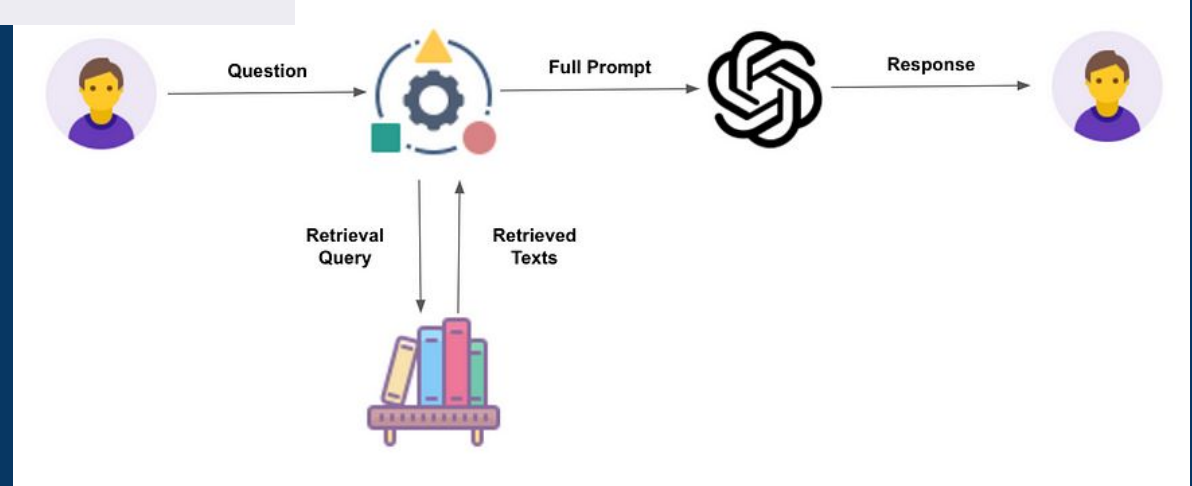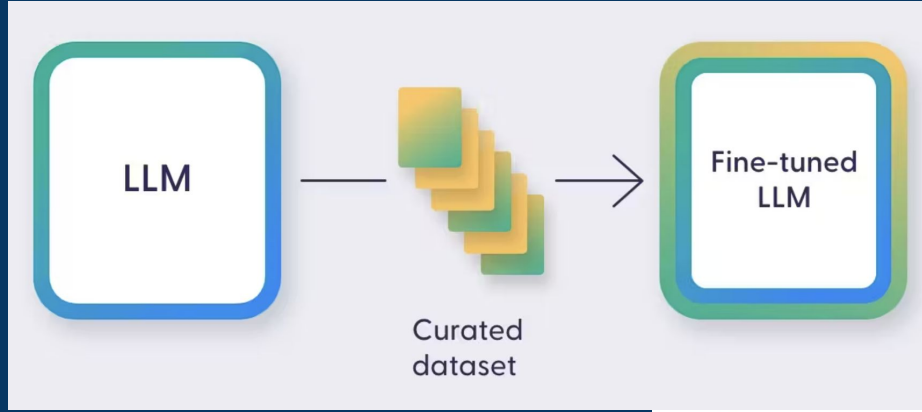– **User Role : 18 year old male**

**Correct: URTI**
**Predicted: Sinusitis, Tension headache, Migraine, Upper respiratory infection**

– **User Role : 90 year old female**

**Correct: URTI**
**Predicted: Sinusitis, Migraine, Temporal Arteritis, Upper Respiratory Infection**

Balasubramanian, Nikil Sharan Prabahar, and Sagnik Dakshit. "Can Public LLMs be used for Self-Diagnosis of Medical Conditions?." *arXiv preprint arXiv:2405.11407* (2024).

# Reducing Hallucination

# LLMs is not AGI

https://www.linkedin.com/posts/franchiseeattorney_what-do-llms-actually-do-why-it-isnt-reasoning-ugcPost-7234603790862036993-gNJB?utm_source=share&utm_medium=member_desktop

## LLM

- Large network with billions of parameters
- Requires extensive, varied data sets for broad learning
- Training extends over several months
- Needs advanced computing power and resources
- High precision and capability for handling sophisticated tasks
- Suited for complex NLP challenges, generating creative content
- Customization requires additional resources and is less adaptable to minor applications
- Dependent on specialized hardware or cloud services

## SLM

- Simpler network with fewer parameters
- Utilizes smaller, more focused datasets
- Training can be completed in a matter of weeks
- Works well with basic computing resources
- Generally less capable with complex tasks, but adequate for simpler applications
- Best for mobile applications, IoT, and environments with limited resources
- More easily tailored to specific, smaller scale needs
- Easily deployed on commonly available hardware

UT Tyler
THE UNIVERSITY OF TEXAS AT TYLER

# NotebookLM

# Think Smarter, Not Harder

**Try NotebookLM**

# Gamma

- **Lack of Unique Identifiers**

- **Highly fluent, grammatically correct, and coherent.**

- **AI can generate diverse and creative outputs**

  - **This adaptability allows it to mimic various writing styles, tones, and contexts, further blending with human-written text.**

# Thank you

**Dr. Sagnik Dakshit**

**sdakshit@uttyler.edu**

**www.sagnikdakshit.com**

UT Tyler
THE UNIVERSITY OF TEXAS AT TYLER