# NLP4Web
# Practice Session 9

## Transformers
## Decoder-only (GPT)

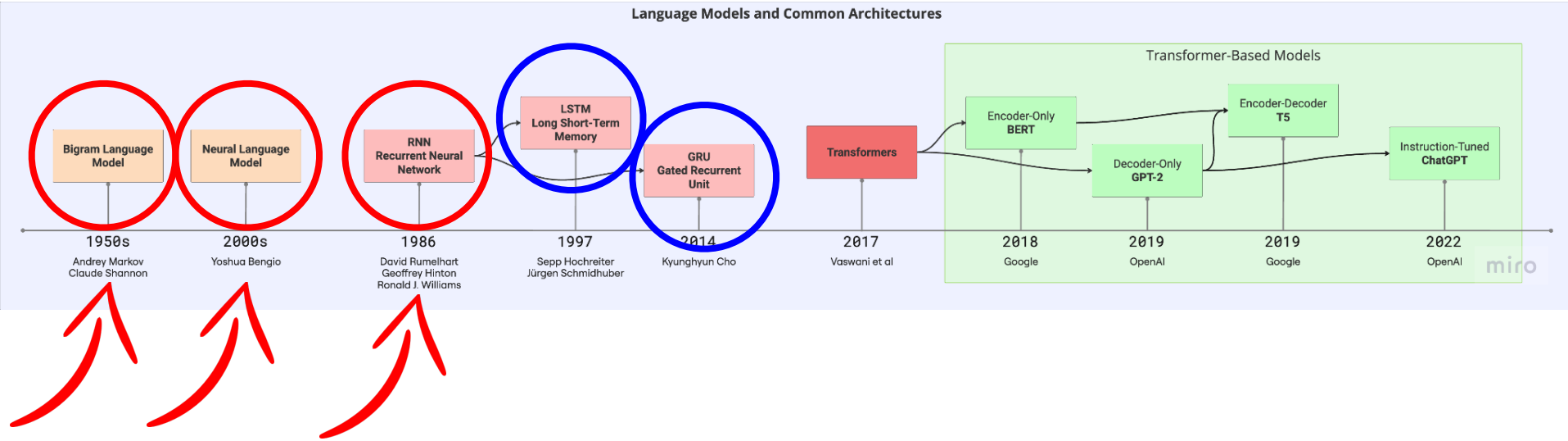Hovhannes Tamoyan

tamohannes.com

To not get lost in space over time, let's Use a **mind map**

# Last time we covered: Bigram LM, NLM, RNN

# LSTM and GRU are for HW4



**Language Models and Common Architectures**

| Bigram Language Model | Neural Language Model | RNN Recurrent Neural Network | LSTM Long Short-Term Memory | GRU Gated Recurrent Unit | Transformers | Transformer-Based Models |
|---|---|---|---|---|---|---|

Encoder-Only BERT · Decoder-Only GPT-2 · Encoder-Decoder T5 · Instruction-Tuned ChatGPT

| 1950s | 2000s | 1986 | 1997 | 2014 | 2017 | 2018 | 2019 | 2019 | 2022 |
|---|---|---|---|---|---|---|---|---|---|
| Andrey Markov Claude Shannon | Yoshua Bengio | David Rumelhart Geoffrey Hinton Ronald J. Williams | Sepp Hochreiter Jürgen Schmidhuber | Kyunghyun Cho | Vaswani et al | Google | OpenAI | Google | OpenAI |

# Today's subject: Transformers (Decoder-only)



Language Models and Common Architectures

# Today's subject: Transformers (Decoder-only)



**Language Models and Common Architectures**

# Recap of Transformer architecture

- The main components
    - Embedding
    - **Positional Encoding**
    - **Self-Attention**
    - Feed Forward
    - **Layer Normalization**
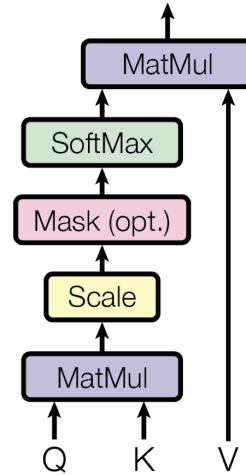    - Residual Connections

# Recap of Attention mechanism

- Scaled Dot-Product attention
- where $\sqrt{d_k}$ is the dimension of the key vector $k$ and query vector $q$

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$$

Scaled Dot-Product Attention

# Recap of Attention mechanism
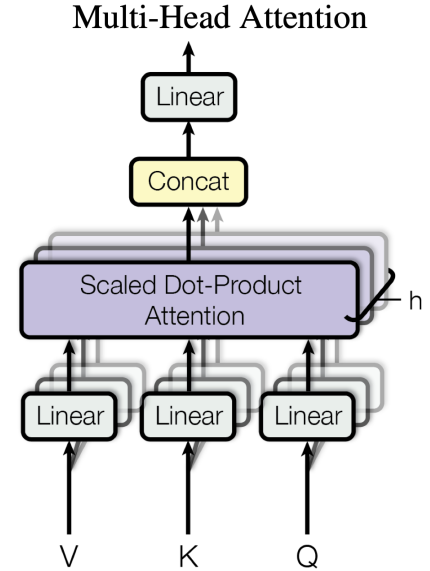
● Multi-head attention



Multi-Head Attention

$$MultiHead(Q, K, V) = Concat(head_1, \ldots, head_h)W^O$$

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

# GPT: Decoder-Only (Autoregressive Architecture)

- **Decoder-only:** generates text by predicting the next word in a sequence based on prior context
- **Unidirectional:** processes text left-to-right, predicting one token at a time
- **Causal Language Modeling (CLM):** trains by predicting the next word in a sequence, avoiding future context
- **Self-attention:** focuses on relevant context in the input to determine the next output word
- **Pre-trained:** fine-tuned for specific tasks with minimal additional training
- **Text Generation:** excels at producing coherent and contextually relevant text based on prompts



Decoder-Only (Unidirectional)
**GPT**