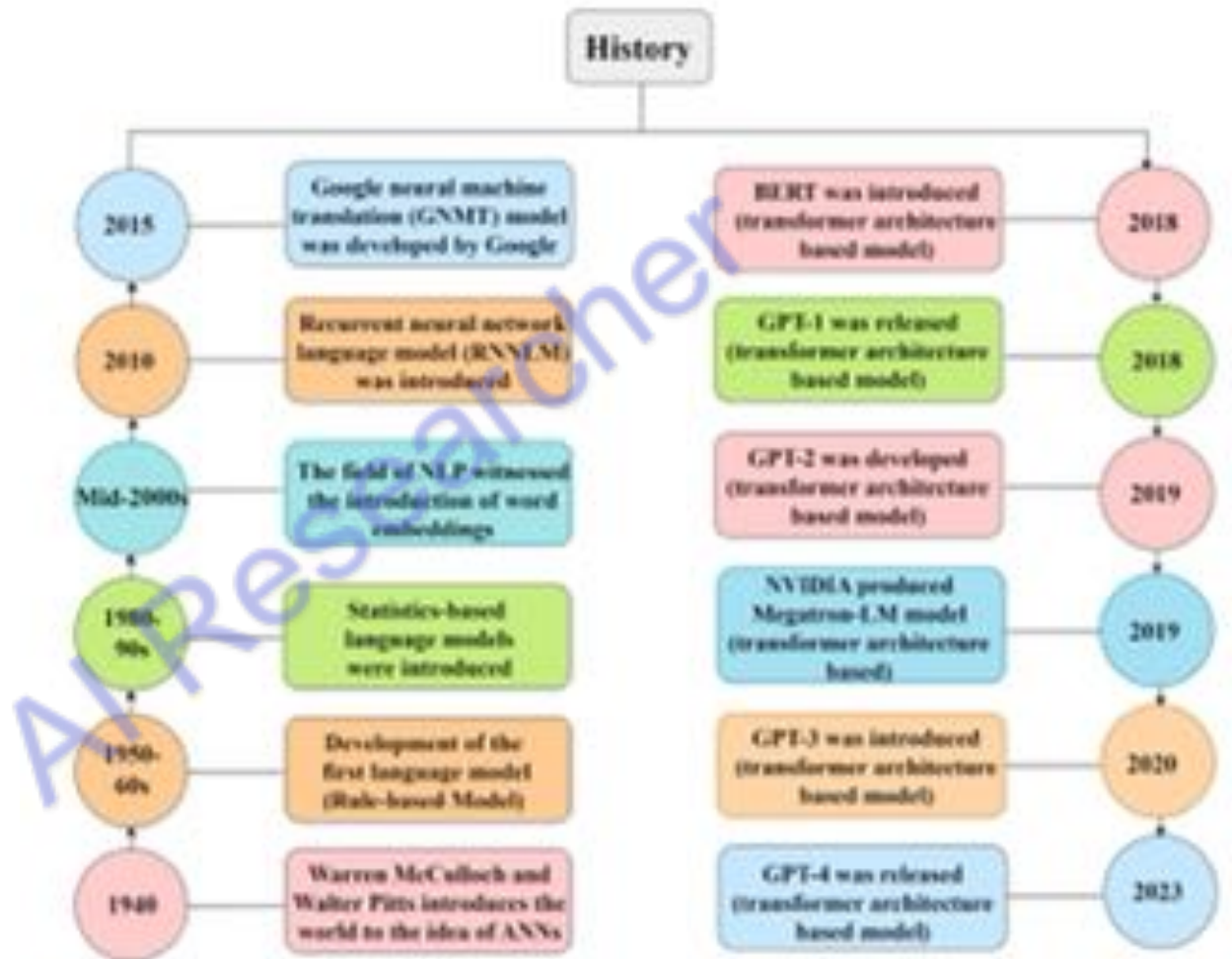


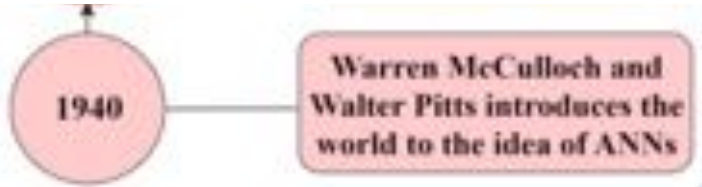


# Language Model History

# History of Language Model

Raiaan, M. A. K., Mukta, M. S. H., Fatema, K., Fahad, N. M., Sakib, S., Mim, M. M. J., ... & Azam, S. (2024). A review on large Language Models: Architectures, applications, taxonomies, open issues and challenges. *IEEE Access*.

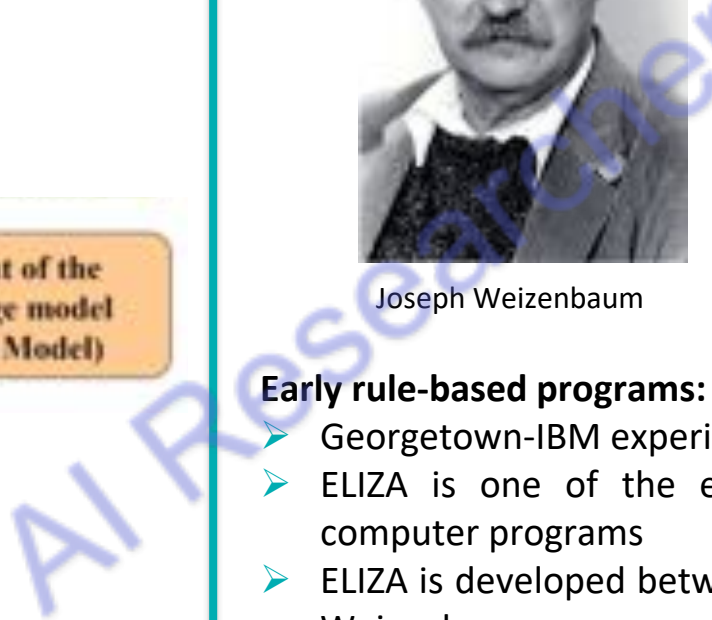




McCulloch (right) and Pitts (left)

- The first mathematical model of a neural network
- Paper: A logical calculus of the ideas immanent in nervous activity
- The paper provided a way to describe brain functions in abstract terms

McCulloch, W.S., Pitts, W. A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics* **5**, 115–133 (1943). <https://doi.org/10.1007/BF02478259>



- Georgetown-IBM experiment, which was conducted in 1954
- ELIZA is one of the earliest natural language processing computer programs
- ELIZA is developed between 1964 and 1967 at MIT by Joseph Weizenbaum

 AI Researcher  
Art's Expertise in Intelligence



Karen Sparck Jones

- A revolution with statistics-based language models
- A statistical interpretation of term-specificity called Inverse Document Frequency (idf)
- Paper: A Statistical Interpretation of term specificity and its application in retrieval

Sparck Jones, K. (1972). A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1), 11-21.



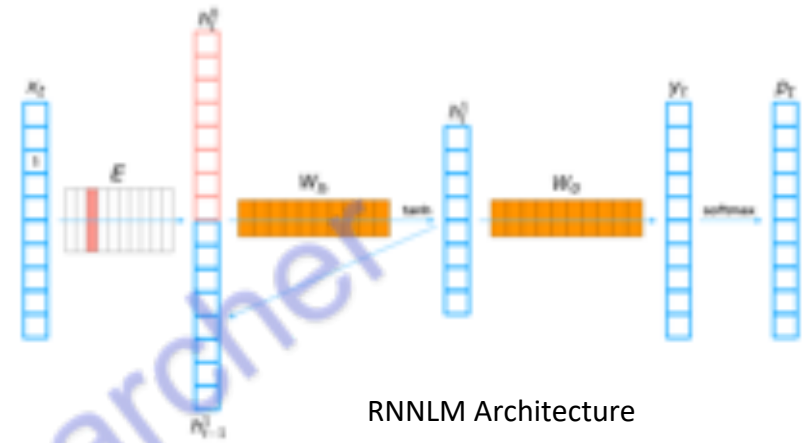
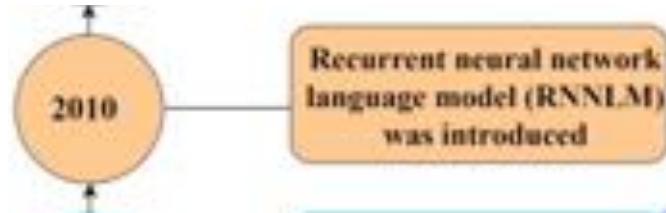
- In the 1990s, Hidden Markov Model, which was used to recognize speech

Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2), 257-286.

In the mid 2000s,

- Researchers began to use deep learning techniques
- Introduced word embeddings in NLP

James, H. M. (2000). Speech and language processing: An introduction to natural language processing computational linguistics and speech recognition. *Person Education, Inc.*



RNNLM Architecture

- Recurrent Neural Network Language Models (RNNLM)
- RNNLM is a type of neural net language models
- Suitable for modeling the sequential data

Mikolov, T., Kombrink, S., Burget, L., Černocký, J., & Khudanpur, S. (2011, May). Extensions of recurrent neural network language model. In *2011 IEEE international conference on acoustics, speech and signal processing (ICASSP)* (pp. 5528-5531). IEEE.

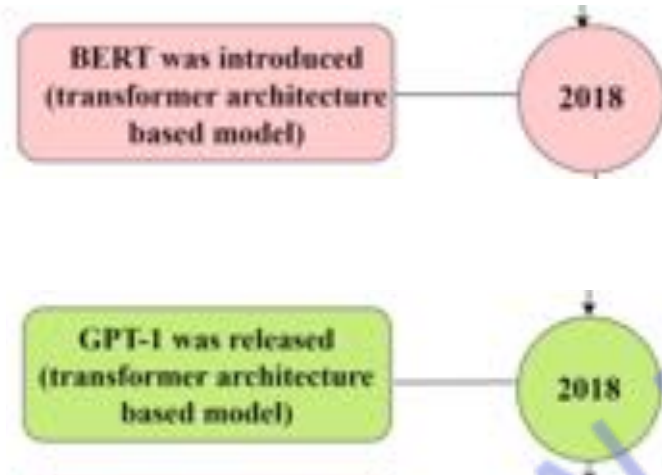




- Google Neural Machine Translation (GNMT) is a neural machine translation system
- It is developed by Google and introduced in November 2016
- It uses an artificial neural network to increase fluency and accuracy in Google Translate.

Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., ... & Dean, J. (2016). Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.



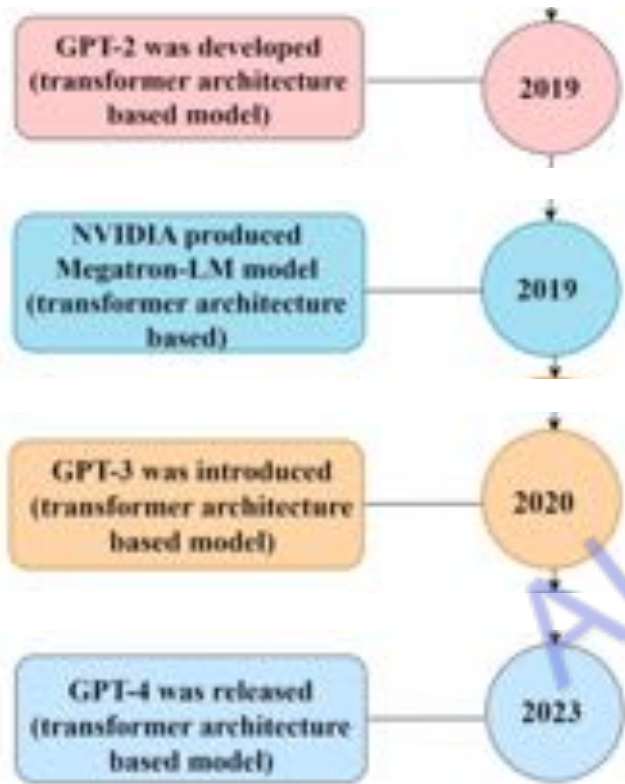


BERT



GPT

- **Bidirectional Encoder Representations from Transformers (BERT)** is a language model based on the transformer architecture. It was introduced in October 2018 by researchers at Google
- **Generative Pre-trained Transformer 1 (GPT-1)** was the first of OpenAI's large language models following Google's invention of the transformer architecture in 2017



- GPT-2 in 2019 refined and expanded upon its predecessor
- NVIDIA wasn't far behind. In 2019, they produced Megatron-LM
- Generative Pre-trained Transformer 3 (GPT-3) is a large language model released by OpenAI in 2020
- Generative Pre-trained Transformer 4 (GPT-4) is a multimodal large language model and the fourth in its series of GPT foundation models. It was launched on March 14, 2023

# Possible Future Development

- Large Language Models may include the development of models using 1-bit quantization for operations.
- 1-bit LLM could enable the deployment of LLMs on low-power devices and broadening their application and accessibility.

# Thank you!

AI Researcher