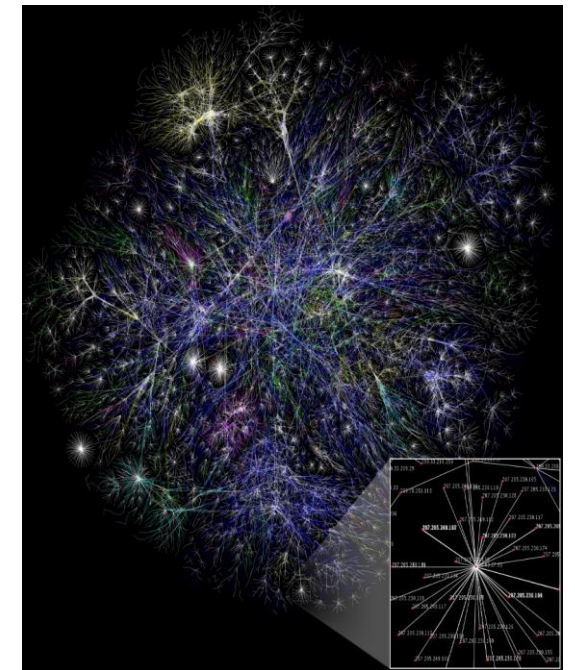# Introductory Course on Transformer Models and Generative AI Tools

- Introduction to Transformer Models
- Generative AI
- Applications of Transformer Models
- Training and Evaluating Transformer Models
- Advanced Topics in Transformers & Generative AI

**Speaker :** *Pranav Raikote, Senior AI Researcher, Plexflo* plexflo

# Introduction

- Understanding Large Language Models
  - Language
  - LLM
  - Some of the Large Language Models
- Background About Transformer Models
- How are transformers different from previous Models
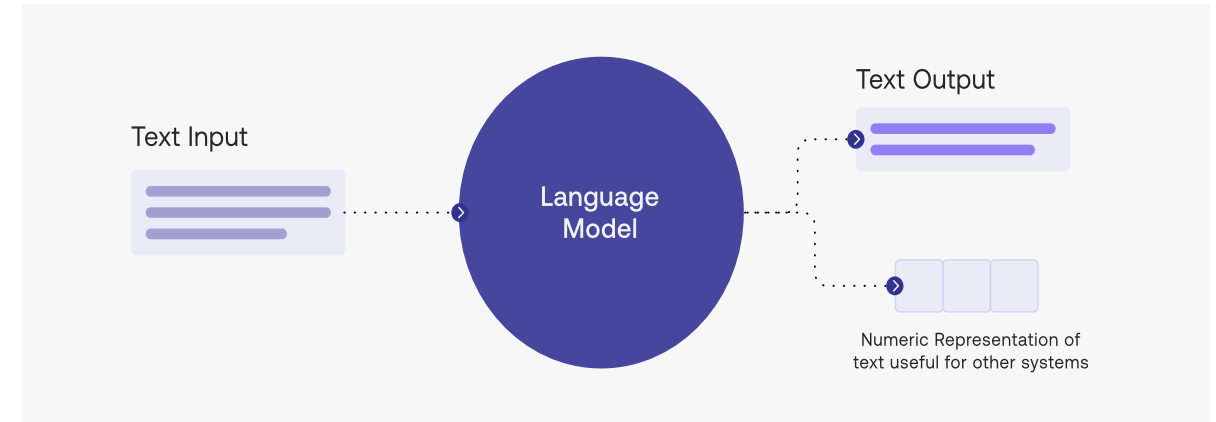- What is Generative Artificial Intelligence (AI)?

- ## Language?

  - It's how we learn about the world. Language is also how we connect and communicate — as people, and as groups and companies
  - Despite the rapid evolution of software, computers remain limited in their ability to deal with language. Software is great at searching for exact matches in text, but often fails at more advanced uses of language — ones that humans employ daily
  - There's a clear need for more intelligent tools that better understand language



Text Input

Language Model

Text Output

Numeric Representation of text useful for other systems

- ## LLM

  - A large language model, or LLM, is a deep learning algorithm that can recognize, summarize, translate, predict and generate text and other content based on knowledge gained from massive datasets

# Transformers

- ## Some of the Large Language Models
  - Pathways Language Model (PaLM) 540 billion parameter model, from Google Research
  - Generalist Language Model (GLaM) 1 trillion parameter model, from Google Research
  - Language Models for Dialog Applications (LaMDA) 137 billion parameter model from Google Research
  - Megatron-Turing NLG 530 billion parameter model, from Microsoft/Nvidia
  - DreamFusion/Imagen 3D image generation from Google Research
  - Get3D from Nvidia
  - MineClip from Nvidia
  - BLOOM: BigScience Large Open-science Open-access Multilingual Language Model with 176 billion parameters.

## Why Transformers

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best-performing models also connect the encoder and decoder through an attention mechanism.

- ## Need for Transformers

  The dominant sequence transduction models are based on
  - o Complex recurrent or convolutional neural networks
  - o The best-performing models also connect the encoder and decoder through an attention mechanism.

- ## Transformers

  - Developed in 2017 by Researchers at Google and the University of Toronto
  - Unlike recurrent neural networks, transformers could be very efficiently parallelized
  - With the right hardware, you could train some really big models.
  - Transformer is a deep learning model that adopts the mechanism of self-attention
  - It differentially weighs the significance of each part of the input data.
  - Like recurrent neural networks (RNNs), transformers are designed to process sequential input data, such as natural language with applications for tasks such as translation and text summarization.
  - However, unlike RNNs, transformers process the entire input all at once.

  Transformers are a type of artificial intelligence (AI) model that uses a self-attention mechanism to process sequential data (When the points in the dataset are dependent on the other points in the dataset), such as natural language (a language that has developed naturally in use), by encoding and decoding input data into a continuous space. This allows for better generalization and understanding of the data, which can be used for various tasks such as language modeling, question answering, machine translation, and summarization.
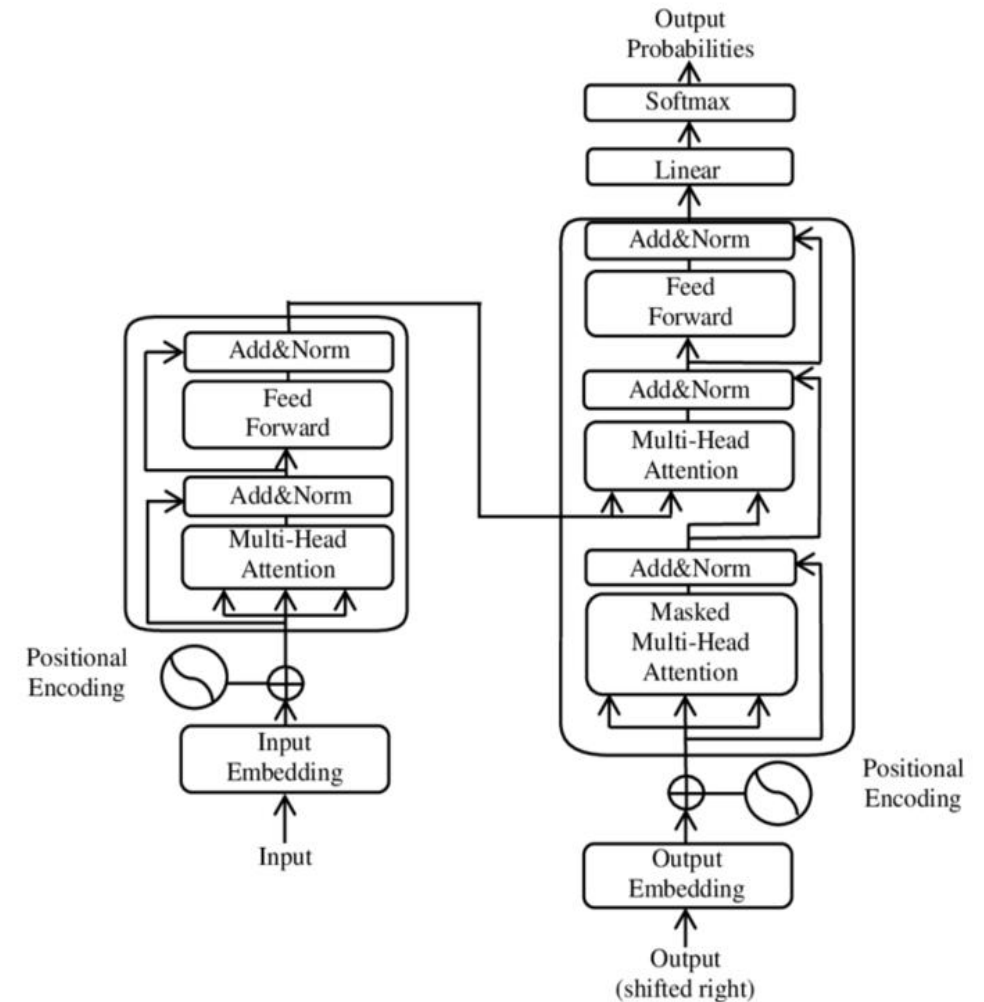
  Transformers are in many cases replacing convolutional and recurrent neural networks (CNNs and RNNs), the most popular types of deep learning models were developed just five years ago.

**India**
**SMART UTILITY**
**Week 2023**

**28 Feb – 04 March 2023 | New Delhi**

ORGANIZER

**ISGF**
India Smart Grid Forum

- **Model Architecture**
  - **Encoder and Decoder Stacks**
  - **Attention**
    - ✓ Scaled Dot-Product Attention
    - ✓ Multi-Head Attention
  - **Position-wise Feed-Forward Networks**
  - **Embeddings and Softmax**
  - **Positional Encoding**
- **Functioning**
  - Transduction models have an encoder-decoder structure
  - the encoder maps an input sequence of symbol representations $(x_1, ..., x_n)$ to a sequence of continuous representations $z = (z_1, ..., z_n)$
  - The decoder then generates an output sequence $(y_1, ..., y_m)$ of symbols one element at a time
  - At each step the model is auto-regressive consuming the previously generated symbols as additional input when generating the next.

# Some of the Transformer Models

- <u>Generative pre-trained transformer</u> (GPT)
- <u>GPT-2</u>: Generative Pre-trained Transformer 2 with 1.5 billion parameters.
- <u>GPT-3</u>: Generative Pre-trained Transformer 3, with the unprecedented size of 2048-token-long context and 175 billion parameters (requiring 800 GB of storage).
- GPT-3.5/<u>ChatGPT</u>/InstructGPT from OpenAI
- <u>GPT-NeoX-20B</u>: An Open-Source Autoregressive Language Model with 20 billion parameters.
- <u>BERT</u>: Bidirectional Encoder Representations from Transformers (BERT)
- <u>OPT-175B by Meta AI</u>: another 175-billion-parameter language model. It is available to the broader AI research community.
- <u>Point-E by OpenAI</u>: a 3D model generator.
- <u>RT-1 by Google</u>: a model for operating robots
- VALL-E text to speech synthesis based on 3-second speech sample It was pre-trained on 60,000 hours of English speech from 7,000 unique speakers (dataset: LibriLight)

# Problems in Training LLMs

- Distributed/Parallel training

- Representation bottleneck

- Vanishing gradients

- Fixed context (sequence length)

India
SMART UTILITY
Week 2023

28 Feb – 04 March 2023 | New Delhi

ORGANIZER

ISGF
India Smart Grid Forum

# Transformers

# Attention Mechanism

- What should I pay attention to?



[Image Credit]

India
SMART UTILITY
Week 2023

28 Feb – 04 March 2023 | New Delhi

ORGANIZER

ISGF
India Smart Grid Forum

# Transformer

- Highly capable Encoder-Decoder model with Self-Attention mechanism

- Train once and fine-tune on a variety of downstream tasks



Image Credit

# Embeddings

## Input Embeddings



## Positional Embeddings



Image Credit

# Positional Encoding

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i)} = sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$



Image Credit

# Self-Attention

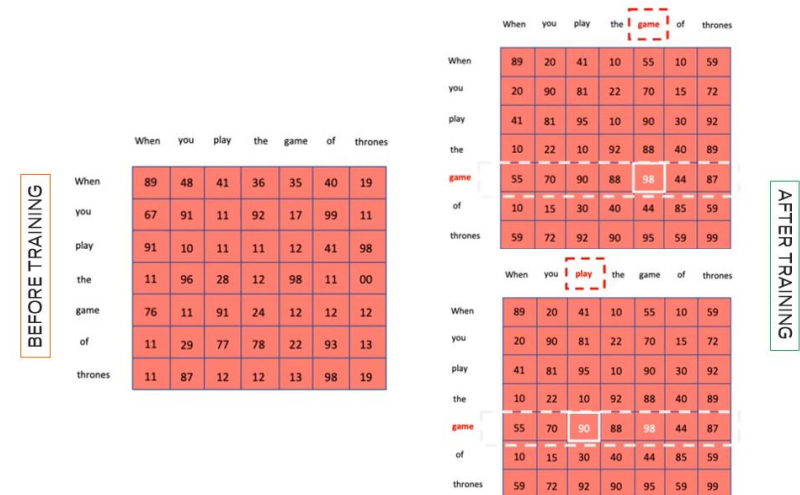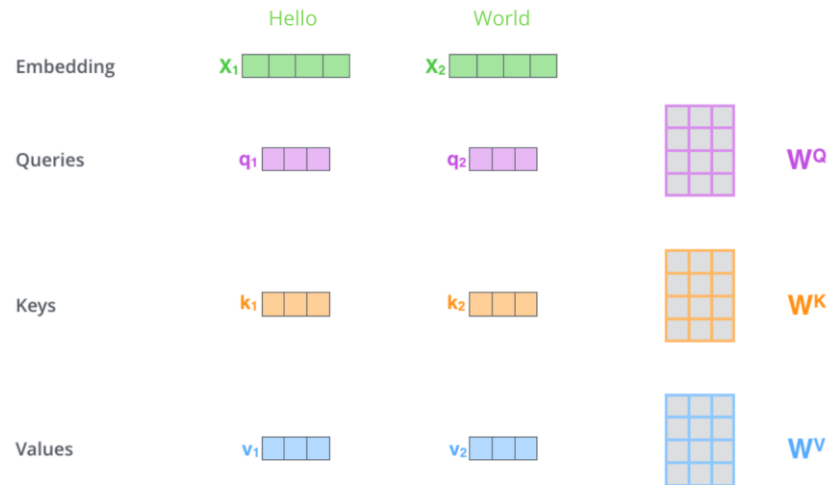He went to the **bank** and learned of his empty account, after which he went to a river **bank** and cried

# Self-Attention



Image Credit

# Multi-Head Attention

Multi-Head Attention diagram and attention head heatmaps (head_0 through head_7)

- 8 heads
- Representation sub-space
- Each head will have a 64-dim vector

Image Credit

# Multi-Head Attention (Recap)



1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply $X$ or $R$ with weight matrices

4) Calculate attention using the resulting $Q$/$K$/$V$ matrices

5) Concatenate the resulting $Z$ matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

Image Credit

# Residual Connections & Feed Forward Network



Image Credit

# Decoder

- Autoregressive in nature

- 2 special Multi-Head Attention layers



Image Credit

# Masked Multi-Head Attention

Scaled Scores

| 0.7 | 0.1 | 0.1 | 0.1 |
| 0.1 | 0.6 | 0.2 | 0.1 |
| 0.1 | 0.3 | 0.6 | 0.1 |
| 0.1 | 0.3 | 0.3 | 0.3 |

+

Look-Ahead Mask

| 0 | -inf | -inf | -inf |
| 0 | 0 | -inf | -inf |
| 0 | 0 | 0 | -inf |
| 0 | 0 | 0 | 0 |

=

Masked Scores

| 0.7 | -inf | -inf | -inf |
| 0.1 | 0.6 | -inf | -inf |
| 0.1 | 0.3 | 0.6 | -inf |
| 0.1 | 0.3 | 0.3 | 0.3 |

Softmax(
| 0.7 | -inf | -inf | -inf |
| 0.1 | 0.6 | -inf | -inf |
| 0.1 | 0.3 | 0.6 | -inf |
| 0.1 | 0.3 | 0.3 | 0.3 |
) =

|  | <start> | I | am | fine |
|---|---|---|---|---|
| <start> | 1 | 0 | 0 | 0 |
| I | 0.37 | 0.62 | 0 | 0 |
| am | 0.26 | 0.31 | 0.43 | 0 |
| fine | 0.21 | 0.26 | 0.26 | 0.26 |

Image Credit

# Final Layers & Loss Function

- Linear Layer with dimension equivalent to number of unique words

- Softmax Layer for class probabilities

- Adam optimizer with Categorical Cross-entropy (Log) Loss

$$Logloss_i = -[y_i \ln p_i + (1 - y_i) \ln(1 - p_i)]$$

# Impact of Transformers

- New avenue and direction in Deep Learning w.r.t NLP & LLMs

- Initially it started off with an idea/approach to understand how languages function and how to represent them

- 66000+ citations and counting

- Base architecture to a plethora of new-age architectures like GPT and BERT

# Generative AI

- Creation of data in various modalities like Images, Text, Audio, Video, and so on

- Solves various data problems
  - Data scarcity
  - Edge cases
  - Dangerous situations

**AI can not only boost our analytic and decision-making abilities but also heighten creativity.**

# History of Generative AI

- First explored in 1950-60s, gained traction in 2010s

- Breakthrough in 2014 by Ian Goodfellow's GANs paper

- Going really strong and dominating the Deep Learning space
  - ChatGPT, ImageGPT, CLIP, DALL-E, Stable Diffusion, etc.



Access to private APIs and open source models drive progress for generative ML

**Coolest idea in Deep Learning in the last 20 years – Yann LeCun**

Image Credit

# Generative AI Timeline

|  | PRE-2020 | 2020 | 2022 | 2023? | 2025? | 2030? |
|---|---|---|---|---|---|---|
| TEXT | Spam detection Translation Basic Q&A | Basic copy writing First drafts | Longer form Second drafts | Vertical fine tuning gets good (scientific papers, etc) | Final drafts better than the human average | Final drafts better than professional writers |
| CODE | 1-line auto-complete | Multi-line generation | Longer form Better accuracy | More languages More verticals | Text to product (draft) | Text to product (final), better than full-time developers |
| IMAGES |  |  | Art Logos Photography | Mock-ups (product design, architecture, etc.) | Final drafts (product design, architecture, etc.) | Final drafts better than professional artists, designers, photographers) |
| VIDEO / 3D / GAMING |  |  | First attempts at 3D/video models | Basic / first draft videos and 3D files | Second drafts | AI Roblox Video games and movies are personalized dreams |

Large model availability: 🔴 First attempts  🟡 Almost there  🟢 Ready for prime time

Global AI investment surged from $12.75 million in 2015 to $93.5 billion in 2021, and the market is projected to reach $422.37 billion by 2028. Already over $2B has been invested in Generative AI, up 425% since 2020

Image Credit

- Generative Adversarial Networks (GANs)

- Variational Autoencoders (VAEs)

- Autoregressive Models (Transformers-XL, PixelCNN, WaveNet)

- Markov Chain Monte Carlo Models

# Generative Adversarial Networks



Image Credit

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, $k$, is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

**for** number of training iterations **do**

    **for** $k$ steps **do**

        • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.

        • Sample minibatch of $m$ examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.

        • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$

    **end for**

    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.

    • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Image Credit

- Problems/Nuances with GANs
  - Mode Collapse
  - Nash Equilibrium
  - Very Very Very difficult to train and get good outputs

- Popular architectures
  - Pix2Pix
  - StyleGAN
  - CycleGAN
  - VideoGAN

# Auto-Encoders



Image Credit

# Variational Auto-Encoders

# Variational Auto-Encoders

Original form

Reparameterized form

```
generation_loss = mean(square(generated_image - real_image))
latent_loss = KL-Divergence(latent_variable, unit_gaussian)
loss = generation_loss + latent_loss
```

Image Credit

# Auto-regressive Models

- Generate a token, given the previous token in an iterative manner

- Tend to forecast future behaviour
  - High impact in NLP & Time-Series domains

- Examples: Transformers-XL, GPT

# Applications of Generative AI

- Realistic and diverse images. Applications in Graphics, Video games, Super resolution, Art, In-painting

- Language Translation, Summarization, Conversational AI

- Simulate complex systems

- Generate data for scientific research
  - Climate modelling
  - Drug discovery
  - Protein folding

- Power & Utility
  - Synthetic data
  - Identify and replicate energy saving patterns
  - Text2Sim, Text2Image, Text2Code, etc.

- Co-Pilot/Codex, Stable Diffusion, PaLM, LaMDA, ChatGPT

- CLIP, DALL-E

- ImageGPT, GPT-3

- StyleGAN2

- BigGAN

- WaveNet

- DeepDream

**Applications of Transformers**

- Regarded as AI 2.0

- Accelerated R&D in various industries

- Cross domain problem statements
  - Language, Images, Videos, Speech, Tabular

- Advances in Drug Discovery, Protein Folding

- Carbon Footprint



Image Credit

# Applications of Transformers

- Vision
  - Image Classification: ViT, DeiT
  - Object Detection, Segmentation: DETR, SWIN Transformer

- Music
  - Music Transformer
  - MusicLM: Synthetic Music generation from text desciptions

- NLG
  - Image Captions, Generating Reports, Chatbots, Transcriptions
  - E-Commerce, Banking, Media, Art & Poetry, Assistants

# Applications of Transformers

- Power & Utility Industry
  - Disaggregation
  - Load Forecasting
  - Pricing optimization
  - Solar & EV Propensity
  - Demand and Response
  - Material & Design Innovations
  - Outage Predictions, Grid Stability & Resilience, Restoration

# Training & Evaluating Transformers

# Training & Evaluating Transformers

- Data
  - Massive data required ~A few TBs
  - Little to no pre-processing
  - Data representation

- Compute
  - Massive GPU servers, mostly one-time investment

- Hyperparameters
  - 6 blocks of Encoders and Decoders
  - 8 Nvidia P100 GPUs, 300,000 steps optimized by Adam (varied LR)

- Outline
  - Dataset -> Character level or Word level tokenization -> Indexing
  - Generate sequences from data
    - If Seq2Seq, then additional masking operation to be done
  - Embedding Representation
  - Define the functional blocks
  - Define the hyperparameters
  - Training & Testing loop
  - Stand-alone function for saving, loading and generation of text

- Evaluation Metrics
  - BLEU
  - Perplexity

# Transformers in Action!
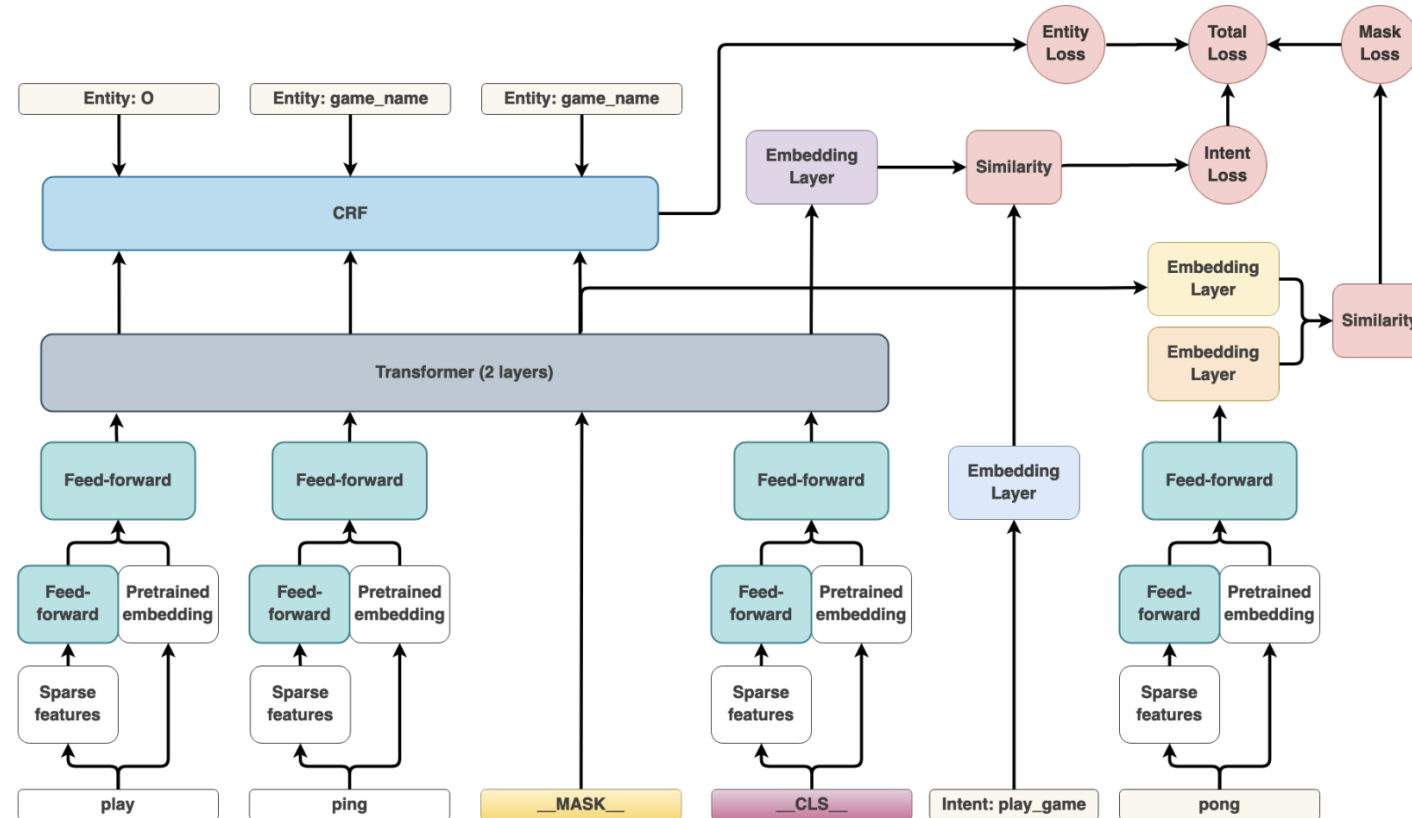
# Advanced Topics in Generative AI

- Transformers in Reinforcement Learning
  - Modified into a "Big Sequence Modelling"
  - Decision Transformers

- Policy Based Learning
  - ChatGPT

- DIET (Dual Intent Entity Transformer)

$\hat{R}_{t-1}$

[Image Credit](#)

# DIET



Image Credit

- Large Language Models & Transformers
  - Understanding of How? Why? And Impact!

- Generative AI
  - Concept and Fundamental Architectures
  - Impact, Applications & Opportunities

- Train your own Generative AI Model
  - Generate novel language
  - How to use pre-trained models for downstream tasks

- Future of Power & Utility Industry with Generative AI

# Thank You

*For discussions/suggestions/queries email: **isuw@isuw.in***
*www.isuw.in*
*Links/References (If any)*

isuw@isuw.in      www.isuw.in      @ISUW_India      India Smart Utility Week (ISUW)