## Literature Review

This document presents the implementation of the research paper titled "A Review on Large Language Models: Architectures, Applications, Taxonomies, Open Issues and Challenges" from IEEE Xplore. The primary goal of my work is to explore, analyze, implement, and compare the key concepts proposed in the paper. The research paper comprises of detailed information about Large Language Models (LLMs), different architectures of Large Language Models, different applications of Large Language Models such as translation, summarization, information retrieval etc and the challenges faced while building and training Large Language Models.

Language plays an important role for communicating and expressing themselves between humans. Similarly, language also helps humans interact with the machines in a way that both understand and correlate with each other. The rapid technological advancements in Large Language Models have transformed various Natural Language Processing (NLP) tasks such as text generation, translation, etc. Despite their complex architectures, optimization challenges and their potential in various applications, they continue to rise in today's technological world and succeed in performing complex tasks.

The research paper I worked on is based on the theoretical concepts of Large Language Models, which is also complex to understand for beginner-to-intermediate learners having a keen interest in building and implementing Large Language Models. My goal is to simplify the research paper by explaining Large Language Models in simple terms along with practical implementations through programming the models, so that learners can understand the concepts of Large Language Models with ease and implement one or two on their own.

## Research Questions and Objectives

### Research Questions

This implementation addresses the following research questions from the reviewed paper:

1.  What are the architectures currently used in Large Language Models and what are their limitations?

    Here is an overview of the architectures used in Large Language Models and their limitations:

1. Transformer-based Models (GPT, BERT, etc)

- Utilize self-attention mechanism to process vast input data. Pre-trained on vast data and fine-tuned for specific tasks.

- Limitations:- Computationally expensive, require large datasets, can generate bias or incorrect outputs.

2. Early Language Models (RNNs and LSTMs)

- Process sequential data by maintaining a hidden state which carries information through steps.

- Limitations:- Inefficient for longer sequences due to vanishing gradients problem and imperfect accuracies due to overfitting problem.

3. Hybrid Models (RETRO, RAG)

- Combine Neural Networks with retrieval mechanisms to enhance accuracy and efficient data extraction.

- Limitations:- Training complexities and Potential inefficiencies.

2. What are the primary applications of Large Language Models and how does the proposed implementations visualize their usability?

Here are the following primary application areas of Large Language Models:

1. Natural Language Analysis, Understanding and Processing

- Sentiment analysis, Text recognition and classification, Entity recognition

- Example:- Analyzing user feedback for social media content categorization

2. Text Generation, Translation and Summarization

- Generating human-like text, Translating different language-based texts Summarizing documents and reports.

- Example:- Summarizing a research paper for better understanding of underlying complexities for a student.

3. Conversational Chatbots and Agents

- Customer support, Question-Answering conversation, Virtual assistants in day-to-day tasks.

- Example:- ChatGPT answering different questions and clarifying queries of a user.

4. Code Generation, Debugging and Assistance

- Autocompleting code, Generating code and Debugging errors in existing code scripts.

- Example:- Generating code for an automation task in Python by GitHub Copilot.

3. How does the proposed implementation visualize the efficiency and usability of the existing Large Language Models?

Through this research paper, I have implemented 4 variations of currently existing Large Language Models namely Phi-2 by Microsoft, BERT (Bidirectional Encoder Representation from Transformers) by Google, Llama 3.2:1b by Meta and a TinyLLM built from scratch.

The goal of this implementation is to visualize the efficiency and usability of these Large Language Models to the upcoming learners and to develop a comparative analysis of the performances of these models based on their outputs.

**Research Objectives**

The main objectives of this implementation from the reviewed paper are as follows:

1. To systematically review the existing literature on LLM architectures and applications

By analyzing the architectures, classifications, variations and limitations of the existing Large Language Models, the objective is to develop a comprehensive understanding of the advancements and their limitations.

2. To implement the currently existing Large Language Models in various tasks for upcoming learners

By practically implementing four of currently existing Large Language Models in various  tasks such as text generation, summarization and

question-answering conversation, the goal is to demonstrate the effectiveness of these Large Language Models to the upcoming learners.

3. To document the key findings and insights gained from the implementation

Through output visualizations, performance metrics and analytical discussion representations of the implemented Large Language Models, learners can understand the architectures of the Large Language Models with ease and through comparative analysis of the models, they also can define which model is suitable for a specific task.

## Implementation of Large Language Models

## Proposed Architectures with steps and output visualizations

1. TinyLLM - A model built from scratch for Text Generation

This is a tiny character-level language model that I have built using PyTorch that learns to generate text one character at a time.

### Architecture Visualization:

```python
import torch
import torch.nn as nn
import torch.nn.functional as F

# === Data Setup ===
text = "hello world! welcome to the world of language models."
chars = sorted(list(set(text)))
vocab_size = len(chars)

char_to_idx = {ch: i for i, ch in enumerate(chars)}
idx_to_char = {i: ch for i, ch in enumerate(chars)}

def encode(s): return [char_to_idx[ch] for ch in s]
def decode(l): return ''.join([idx_to_char[i] for i in l])

data = torch.tensor(encode(text), dtype=torch.long)

# === Model Definition ===
class TinyLLM(nn.Module):
    def __init__(self, vocab_size, embed_dim):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, embed_dim)
        self.linear = nn.Linear(embed_dim, vocab_size)

    def forward(self, x):
        x = self.embedding(x)
        x = self.linear(x)
        return x

model = TinyLLM(vocab_size, embed_dim=32)
optimizer = torch.optim.Adam(model.parameters(), lr=1e-2)
loss_fn = nn.CrossEntropyLoss()

# === Training Setup ===
block_size = 4

def get_batch():
    ix = torch.randint(0, len(data) - block_size - 1, (1,))
    x = data[ix:ix + block_size]
    y = data[ix + 1:ix + block_size + 1]
    return x.unsqueeze(0), y.unsqueeze(0)

# === Training Loop ===
for step in range(1000):
    x_batch, y_batch = get_batch()
    logits = model(x_batch)
    loss = loss_fn(logits.view(-1, vocab_size), y_batch.view(-1))
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
    if step % 100 == 0:
        print(f"Step {step}, Loss: {loss.item():.4f}")

# === Text Generation ===
context = torch.tensor([[char_to_idx['h']]], dtype=torch.long).unsqueeze(0)

for _ in range(100):
    logits = model(context)
    last_logits = logits[0, -1]
    probs = F.softmax(last_logits, dim=0)
    next_char_idx = torch.multinomial(probs, num_samples=1)
    context = torch.cat([context, next_char_idx.unsqueeze(0)], dim=1)

print("Generated text:", decode(context[0].tolist()))
```

**Algorithm Steps:**

- A sample text is defined and the unique characters get sorted in order with the length of those characters recorded.

- Two dictionaries are defined to map characters to indexes and vice-versa.

- The functions "Encode" and "Decode" are defined to convert a string to a list of character indices and vice-versa. The input text is converted into a tensor of indices using Tensor sub-module from Torch.

- The TinyLLM model is defined with

    - An Embedding layer that maps each character to a dense vector

    - A Linear Layer that maps from embedding space back to the size of vocabulary for prediction.

- The model is created with the forward function defining the data flow direction through the model and defining Adam optimization function and CrossEntropyLoss (for classification)

- "block_size = 4" is defined so that the model looks at four characters at a time to predict the next one.

- The function "get_batch" generates a training example with x being current characters and y being next characters.

- The built model is trained for 1000 steps where:

    - Gets a random batch, runs it through the model and computes loss by reshaping tensors into 2D

    - Backpropagation and optimization steps are also performed, printing loss for every 100 steps to track training.

- The text generation starts with the letter "h" where for every 100 steps:

    - Feed the current context into the model

    - Get the prediction for last character

    - Apply softmax activation function to get probabilistic distributions

    - Append to the context and decodes the final generated character sequence back to text.

**Output Visualization**

```
sachinkarthikeya@Vs-MacBook-Pro-2 Megaminds Task % python3 tinyllm.py
Step 0, Loss: 3.5380
Step 100, Loss: 1.3157
Step 200, Loss: 1.6961
Step 300, Loss: 1.2366
Step 400, Loss: 0.6637
Step 500, Loss: 0.4957
Step 600, Loss: 1.6336
Step 700, Loss: 0.6944
Step 800, Loss: 0.9221
Step 900, Loss: 0.4093
Generated text: helanguaguahe melo worlange taguagelanguange wo the wele to corlanguanguaguanguangelanguanguanguaguag
sachinkarthikeya@Vs-MacBook-Pro-2 Megaminds Task % ▊
```

From the above output, we can understand that the model can learn the patterns of the character in the given text, and generates the next character in the output.

However, this implementation is a minimal character-level language model, which is great for understanding the basics of building Large Language Models. This code can be further improved or advanced in terms of code quality, structure, training dynamics and text generation such as:

- Adding a Sequential model architecture to capture sequence history more efficiently.

- Adding more data and improving the model accordingly like using dropout in the model and validation loss tracking.

- Switching to word-level tokenization for realistic language modeling and generation.

- Using Transformer architecture (multi-head self-attention mechanism) for long range dependencies.

2. Phi-2 model by Microsoft for Question-Answering Conversation

This is a pre-trained Large Language model developed my Microsoft, built based on Hugging Face Transformers. It is used to generate a response to a prompt.

**Architecture Visualization**

```python
from transformers import AutoTokenizer, AutoModelForCausalLM
import torch

model_name = "microsoft/phi-2"

tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(model_name)

prompt = input("Enter your prompt: ")
inputs = tokenizer(prompt, return_tensors="pt")
outputs = model.generate(**inputs, max_new_tokens=50)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

**Algorithm Steps:**

- Python libraries necessary for this model are imported such as:

  - transformers:- Hugging Face library that provides pre-trained models and tokenizers

  - AutoTokenizer:- Automatically loads the correct tokenizer for model I choose.

  - AutoModelForCausalLM:- Loads a model designed for casual language modeling(predicting next token based on previous one)

  - torch:- PyTorch library for neural network operations and model execution

- Define the Microsoft Phi-2 model along with downloading and loading the necessary tokenizer and pre-trained neural network to generate text.

- The prompt is defined to let the user enter his question or query, which the tokenizer receives and converts into IDs.

- The defined model generates continuation text from the prompt, limiting the output to 50 newly generated tokens.

- Output is initially a tensor of token IDs, which is converted back to a readable string and prints the string-formatted answer to the asked question or query.

**Output Visualization:**



```
sachinkarthikeya@Vs-MacBook-Pro-2 Megaminds Task % python3 phi2_model.py
Loading checkpoint shards: 100%|████████████████████████| 2/2 [00:15<00:00,  7.79s/it]
Enter your prompt: what is the capital of India?
Setting `pad_token_id` to `eos_token_id`:50256 for open-end generation.
The capital of India is New Delhi.
```

From the above output, we can say that the model answered to the question correctly. Improvements can be done to this code like using "torch.no_grad()" to disable gradient trackings, enabling sampling, improving diversity, controlling randomness for more creative and varied output.

3. BERT (Bidirectional Encoder Representations from Transformers) by Google for Text Summarization

This is also a pre-trained Large Language model developed my Google, built based on Hugging Face Transformers. It is used to to perform automatic text summarization.

**Architecture Visualization**

```
from transformers import pipeline

# Load pre-trained summarization pipeline
summarizer = pipeline("summarization", model="sshleifer/distilbart-cnn-12-6")

# Your long text input
text = """
Language models are a key component of many modern NLP systems. They are trained to predict the next word in a
sequence given the previous words, enabling applications such as text generation, machine translation, and question
answering.
Recent advancements in transformer-based models like BERT and GPT have significantly pushed the boundaries of what
language models can achieve, thanks to their ability to model long-range dependencies and large-scale training on
massive datasets.
"""

# Get the summary
summary = summarizer(text, max_length=60, min_length=25, do_sample=False)

# Print the summary
print("Summary:", summary[0]['summary_text'])
```

Algorithm Steps:

- Import pipeline, a high-level API in Hugging Face Transformers that helps in model loading and pre-processing steps.

- This helps in loading a summarization pipeline and a pre-trained model namely "sshleifer/distilbart-cnn-12-6", which is a lighter and faster version of BART (Bidirectional and Auto-Regressive Transformers), trained on CNN dataset for summarization tasks

- A long-paragraphed text is given as an input, which goes through model where:

  - The maximum and minimum length of the generated summary in tokens is defined.

  - Disables sampling i.e., makes the output deterministic instead of random.

- Extracts and prints the actual summary string from the result.

**Output Visualization**



From the above output, we can say that the model summarized the given long-paragraphed text into a two-lined text. Improvements can be done like using DistilBART for longer paragraphs, using more powerful models like "facebook/bart-large-cnn" or "google/pegasus-xsum".

## 4. Llama 3.2:1b by Meta for Conversational-based interaction

Llama 3.2:1b is the 2nd version of the 3rd generation of Llama-based Large Language Model developed by Meta, which contains 1 Billion parameters. It is downloaded and run locally, and is used to generate responses or conversations.

**Architecture Visualization**

```python
import ollama

model_name = "llama3.2:1b"
prompt = input("Enter your prompt: ")

response = ollama.chat(
    model=model_name,
    messages=[{"role": "user", "content": prompt}]
)

print("AI Response:", response["message"]["content"])
```

**Algorithm Steps:**

- Python package "Ollama" is imported, which helps to download, run and chat with models like Llama, Mistral, etc which run on local machines.

- Defining the "llama3.2:1b" model for use and the prompt to let the user enter his query.

- The query is sent like a chat-style message to the model where a list of conversations turns with each having a "role" and a "content".

- The response is extracted and prints the generated text from the model's response.

**Output Visualization**

## Conclusion and Future Directions

The implementation of the research paper titled "A Review on Large Language Models: Architectures, Applications, Taxonomies, Open Issues and Challenges" provided valuable theoretical insights into the world of Large Language Models. Through an in-depth analysis and literature review, I have understood the key architectures, applications and challenges of Large Language Models.

The models that I have implemented aim to help upcoming learners understand about Large Language Models with ease and also understand the theoretical concepts behind the building of these models without any hardship. The models implemented perform various tasks of existing Large Language models such as Text generation, Text summarization, Question-Answering based interaction and Conversational interaction, helping learners determine the best model that can be used for a specific task.

While the implemented models show promising results, there is always room for improvements and advancements of these models like:

1. Improving model's generalization robustness across several languages and dialects.

2. Developing more energy-efficient models that reduce computational costs of building and training, decreasing environmental concerns

3. Exploring into several methods of data input such as audio, video and visual data, further enhancing model capabilities.

4. Investigating the ethical aspects of deploying Large Language Models such as controlling biased responses, ensuring fairness in response generation and data privacy for responsible AI development.

5. Training Large Language Models for specific industry-based applications like healthcare and defense sectors, addressing domain-specific challenges.

By continuing to address these advancements, the implemented models can be further advanced to contribute for further complex tasks such as translation, code debugging, etc, thereby fostering innovation in Natural Language Processing applications.

## References

Original Research Paper from IEEE Xplore:- https://ieeexplore.ieee.org/document/10433480

"A Comprehensive Overview of Large Language Models" Research Paper from arXiv:- https://arxiv.org/pdf/2307.06435

"TinyLlama: An Open-Source Small Language Model" Research Paper from arXiv:- https://arxiv.org/pdf/2401.02385