# STEM CHARM

Rwan Ashraf[1], Hajar Mohamed[1], Samer Elhossany[1], Abdallah Saber[1], Omar Mohamed[1] [1]School of Information Technology and Computer Science (ITCS), Nile University, Giza, Egypt

*Abstract*—This paper presents STEM-CHARM (Comedic, Humorous, Adaptive, Reading Mentor) , a multimodal intelligent tutoring system designed to enhance STEM education through natural language processing, humor integration, and adaptive communication. STEM-CHARM helps students to face the diffculity in learning science, technology, engineering, and mathematics by offering a warm, engaging, and personalized learning experience. The system incorporates voice-to-text input, facial expression recognition, and sign language interpretation to accommodate various user communication modes and emotional states. These features enable STEM-CHARM to deliver responses that are context-aware and aligned with the user's preferred mode of interaction. From data preprocessing to response generation, the system is developed with both technical rigor and user experience in mind. Initial findings suggest that STEM-CHARM improves engagement, comprehension, and accessibility, particularly for learners with special needs. Future enhancements will focus on enabling Arabic language understanding and response generation and expand accessibility in multilingual educational settings. With its unique blend of humor and human-centered AI, STEM-CHARM demonstrates strong potential as a transformative tool for modern STEM education.

*Index Terms*—STEM education, intelligent tutoring system, natural language processing, sign language recognition, facial expression analysis, voice-to-text, humor in education, multimodal interaction.

## I. INTRODUCTION

It is believed that AI's primary goal is to build an intelligent machine that can think without being programmed, and the second goal is to find out the nature of intelligence (Schank, 1987). The AI concept was born in 1942 when the American science fiction writer Isaac Asimov published a short story called "Runaround," and this story was about a robot that 3 engineers have built, and they included 3 important rules, which are (1) a robot may not injure a human being or, through inaction, allow a human being to come to harm; (2) a robot must obey the orders given to it by human beings except where such orders would conflict with the First Law; and (3) a robot must protect its own existence as long as such protection does not conflict with the First or Second Laws. (Haenlein Kaplan, 2019). The concept of machine intelligence was introduced in the paper "COMPUTING MACHINERY AND INTELLIGENCE," which was written by Alan Turing in 1950, where he believed that machines could be programmed to learn and improve over time. The first AI program came to life in 1956 in the paper "Empirical explorations of the logic theory machine: a case study in heuristic," which was published in 1957 by Newell et al where they made a program called "The Logic Theory Machine 'LT,' " and the program's purpose was to learn how it is possible to solve difficult problems such as proving mathematical theorems, discovering scientific laws from data, playing chess, or understanding the meaning of English prose. The LT has successfully proved 38 out of 52 theorems correctly. In 1958 the first neural network was invented by Frank Rosenblatt in his paper "THE PERCEPTRON: A PROBABILISTIC MODEL FOR INFORMATION STORAGE AND ORGANIZATION IN THE BRAIN," where he created the 1st perceptron. In 2017 Google published the paper that changed the world of AI and made it in the format we know today, which is "Attention is all you need," whereby it introduced the transformers and how they efficiently preceded the traditional deep learning algorithms like RNN and CNN, and this was the intro to the world of LLMs. Based on the OpenAI blog that was written by Radford et al., GPT-1 came out in 2018, and it proved its high performance compared to the other models used in the experiment. Although all these changes, the change that affected the world of AI and LLMs till today was the release of ChatGPT from OpenAI in 2022. Although ChatGPT has accomplished a lot so far, it is still considered a "general-purpose LLM," and general-purpose LLMs are models trained on excessive and diverse datasets in various topics, which in turn affects their performance when it comes to a task that requires high accuracy. Our aim is to perform an educational platform, "STEMMY," that aims to fill this gap as it is focused on science but with a sense of humor to make the educational problem easier. We also added a feature that we didn't find in the normal LLMs like ChatGPT, DeepSeek, Grok AI, Claude, and Gemini, which is an ASL translator using the YOLO model and another model for facial expressions, and the aim of this model is to integrate the emotions with the response to get more satisfaction from the term of the users. We have speech-to-text models as well to take the voice of the user and convert it to text. The main model is the Llama 3.2 3b, and this model uses RAG. We aim to make an AI platform that fills these gaps.

## II. LITERATURE REVIEW

Humor has long been recognized not only as a form of entertainment but also as a cognitive tool that can improve memory retention, especially in educational contexts (Schmidt, 1994). With the rise of machine learning and natural language processing (NLP), researchers have begun exploring how humor can be computationally understood and even generated, aiming to bridge emotional intelligence and artificial intelligence.

**1. Humor Detection: Understanding the Unpredictable:** Detecting humor computationally is inherently challenging due to its subjective and context-dependent nature. Traditional

machine learning approaches to humor detection have relied on handcrafted features such as incongruity, surprise, or wordplay patterns (Mihalcea Strapparava, 2005). These features, while intuitive, are often insufficient in capturing deeper semantic and contextual signals of humor.

To address this, deep learning models have gained popularity, particularly those using contextual embeddings. BERT (Bidirectional Encoder Representations from Transformers), introduced by Devlin et al. (2019), has been especially effective in humor detection tasks. Its ability to capture bidirectional context allows it to model the nuanced semantics and syntactic variations that often underlie humorous text. Several studies have shown BERT outperforming traditional LSTM or CNN-based models in humor detection benchmarks (Weller Seppi, 2019).

**2. Humor Generation: Turning Facts into Funny:**

On the generation side, the goal of converting scientific or educational content into humorous text poses a unique challenge. Unlike regular text generation, humor generation requires not just grammatical correctness but creativity, timing, and often, cultural knowledge.

Transformer-based models like GPT (Radford et al., 2018) have proven capable of generating humorous content when fine-tuned on joke datasets or humorous corpora (See et al., 2019). These models leverage attention mechanisms to maintain coherence while allowing for flexibility in expression, making them well-suited for generating witty rephrasings of factual content.

To support memorization, humor generation in your project focuses on transforming dry scientific concepts into funny, memorable analogies or one-liners. This aligns with cognitive research suggesting that emotionally engaging content improves information recall (Tyng et al., 2017).

**3. Multi-Modal Integration: YOLO for Humor in Visuals:**

In addition to text-based humor, your project also explores visual humor, integrating YOLO (You Only Look Once) for object detection. YOLO, developed by Redmon et al. (2016), is a real-time object detection system that can be leveraged for visual joke setups — for example, identifying objects in an image and pairing them with humorous captions or unexpected juxtapositions. This fusion of computer vision with NLP allows the creation of multimodal humor, similar to internet memes, which are popular and effective in informal education (Shifman, 2014).

**4. Ensemble Techniques: Classifiers for Fine-Tuning Humor Types:**

Finally, ensemble classifiers are used to refine humor detection and classification based on type (e.g., puns, sarcasm, slapstick). These classifiers — often including random forests, logistic regression, and support vector machines — can be applied on top of transformer-based embeddings to improve interpretability and provide multi-label predictions **(Chen et al., 2020)**.

By combining YOLO for visual cues, BERT for semantic context, and transformers for generation, your project addresses both humor detection and generation in a comprehensive, multi-modal way. This not only makes the AI better at recognizing humor in different forms but also turns it into a tool that enhances learning through entertainment.

## III. DISCUSSION

The development of STEM CHARM represents a novel interdisciplinary approach that bridges natural language processing, computer vision, and educational psychology to make scientific content more engaging and memorable through humor. Our system tackles two core tasks: humor detection and humor generation, with a particular emphasis on educational transformation — converting complex STEM concepts into humorous formats to support memory retention and learning motivation.

### A. Effectiveness of Humor Detection Models:

The humor detection component of STEM CHARM, built on top of a fine-tuned BERT model, demonstrated strong performance in identifying humorous versus non-humorous content across varied datasets. Compared to traditional classifiers, such as logistic regression and random forests, BERT's contextual understanding of language significantly improved detection accuracy. This confirms existing findings that transformer-based models outperform shallow models in tasks that require semantic nuance (Devlin et al., 2019).

In addition to textual detection, we integrated YOLOv5 for identifying visual elements that can be humorously contextualized. While YOLO was not designed for humor, its real-time object detection capabilities allow the system to generate visual punchlines by pairing recognized objects with humorous captions or scientific analogies. This multimodal fusion has the potential to extend STEM CHARM's reach beyond text, creating meme-like educational content that appeals to visual learners.

### B. Humor Generation and Educational Utility:

On the generation side, we used transformer-based architectures fine-tuned on joke and pun corpora to convert scientific explanations into humorous narratives. Our pilot tests show promising results: students retained scientific information better when it was presented with comedic analogies than through traditional explanations. This aligns with prior cognitive studies suggesting that emotionally charged content, such as humor, can enhance encoding and recall (Tyng et al., 2017; Schmidt, 1994).

One noteworthy challenge was controlling humor appropriateness and content fidelity. While models like GPT-2 and T5 excelled at generating linguistically coherent content, they occasionally produced jokes that were off-topic, overly abstract, or culturally biased. This underscores the need for stricter content conditioning and ethical filtering in future iterations.

### C. System Integration and Real-World Use:

The ensemble nature of STEM CHARM — combining BERT, classifiers, YOLO, and generative transformers — proved essential. Each model contributes a different capability:

detection, classification, visual grounding, and generation. This modular design allows for continuous improvement in individual components without compromising system-wide coherence.

Initial feedback from both educators and learners has been encouraging. STEM CHARM not only helped demystify dense scientific topics but also added an element of playfulness that made the content feel more approachable. Humor, when used carefully, appears to reduce anxiety around complex subjects and improves overall engagement — a promising direction for STEM education tools.

### D. Limitations and Future Directions:

Despite its successes, STEM CHARM faces several limitations. First, humor remains deeply cultural and subjective; what's funny to one audience may be confusing or inappropriate to another. Incorporating user profiling or adaptive humor styles may help address this. Second, the system still relies heavily on pre-trained models that may not fully capture the scientific accuracy needed in educational contexts. Future work should integrate domain-specific knowledge bases or use human-in-the-loop feedback for refining content.

In addition, further evaluation is needed on a larger scale, ideally across classrooms or e-learning platforms, to assess the long-term impact of humor-based learning on retention, comprehension, and academic performance.

## IV. RELATED WORK

### A. Emotion-Aware Tutoring Systems

Intelligent tutoring systems such as **AutoTutor** have laid the groundwork for emotionally responsive education by analyzing learner input and adapting dialogue accordingly. While systems like AffectiveTutor build on this to detect affective states, they often lack warmth and tend to be formal or robotic. **STEM CHARM** extends this work by infusing educational explanations with humor and personality—making complex STEM subjects like calculus or data structures more approachable and engaging. **(Graesser et al. (2005))**

### B. Humorous Language Generation

Recent advancements in transformer-based models such as **GPT-2** and **GPT-Emotion** have demonstrated the ability to generate emotionally aware language. However, their focus is typically on general conversation or empathetic response, not academic instruction. **STEM CHARM** adapts these techniques to educational contexts, tailoring humorous tone and language to clarify technical content in subjects like biology, programming, and math. **(Radford et al. (2019) – GPT-2)**

### C. Humor in Natural Language Processing

Projects like **Humicroedit** have demonstrated that NLP models can be fine-tuned for humor detection and generation. Nevertheless, such systems primarily target entertainment or social media. In contrast, **STEM CHARM** integrates humor with explanation, using puns, analogies, and playful metaphors to enhance memory retention and reduce learner anxiety in STEM fields. **(Zhou et al. (2018) – Emotional Chatting Machine)**

### D. STEM-Focused AI Tools

Existing tools like **Wolfram Alpha** and **Socratic by Google** offer formal explanations for STEM queries . While informative, they lack interactive or affective engagement. **STEM CHARM** acts as a "quirky STEM tutor," providing not just accurate answers but delivering them with a humorous twist to boost motivation and comprehension.**(Hossain et al. (2019) – Humicroedit)**

### E. Expressive Text-to-Speech Models

Text-to-speech (TTS) technologies like **Tacotron 2** and **WaveNet** have enabled natural-sounding speech synthesis . However, these systems have yet to be widely applied in personalized, expressive educational tools. **STEM CHARM** utilizes TTS for expressive delivery of content, ensuring the speech reflects the humorous and friendly nature of the explanations, thereby humanizing the machine tutor. **( Van den Oord et al. (2016) – WaveNet)**

### F. Transformer Architectures and Language Understanding

Transformer architectures, introduced in the landmark paper *"Attention Is All You Need"* , serve as the backbone for most modern NLP systems. These models, including BERT and GPT variants, allow **STEM CHARM** to generate contextually aware, humorous explanations. The attention mechanism enables the system to focus on key semantic elements in STEM content, which is essential when transforming technical information into light-hearted explanations without losing meaning. **Vaswani et al. (2017) – Attention is All You Need)**

### G. Summary of Gaps and Opportunities

Despite advancements in humor generation, intelligent tutoring, and expressive language models, no existing system merges them cohesively for educational purposes in STEM. **STEM CHARM** fills this gap by combining humor-aware generation, emotional expressiveness, and technical accuracy—offering a unique solution that makes STEM content fun, relatable, and deeply memorable.

## V. METHODOLOGY

### A. Data Preprocessing and Classification

To categorize the dataset into two distinct classes sense of humor" and "non-sense of humor" we employed two deep learning-based classifiers. The first was a pre-trained BERT model, while the second was a custom-built Transformer encoder model, based on the architecture proposed in Google's seminal 2017 paper "Attention Is All You Need". To ensure a fair comparison, we utilized only the encoder component of the Transformer, aligning with BERT's encoder-only structure and excluding the decoder.

## B. Tokenization and Vocabulary Construction

We began by implementing a custom preprocessing function, tokenizerPreprocess(), which segmented the input text using common delimiters: [' ', '!', '?', '.', '-', ':']. For tokenizer training, we utilized five diverse and widely-used natural language datasets: IMDb, SQuAD, Yelp Polarity, CoQA, and AG News. These datasets were concatenated into a single corpus, which was then used to train a tokenizer based on the Byte Pair Encoding (BPE) technique. The tokenizer was configured with a vocabulary size of 15,000 tokens and incorporated special tokens similar to those used in BERT: - [PAD] for padding sequences to uniform length - [UNK] for unknown or out-of-vocabulary words - [CLS] to indicate the beginning of an input sequence - [SEP] to separate segments within the same input (e.g., in QA tasks)

## C. Transformer Encoder Architecture

Following tokenization, the text tokens were converted into embeddings and passed through a positional encoding layer. The positional encoding mechanism applied was consistent with the original Transformer paper, using sine and cosine functions defined as: - $P(k, 2i) = \sin(k / n \text{ power}(2i/d))$ for even indices - $P(k, 2i+1) = \cos(k / n \text{ power}(2i/d))$ for odd indices Here, $k$ is the position of the token, $i$ is the index in the embedding dimension, $d$ is the embedding size, and $n$ is a scaling factor (commonly set to 10,000). The positional embeddings were then passed to the multi-head self-attention mechanism, followed by add and normalization layers. The output was further processed through two fully connected layers, using ReLU as the activation function, and dropout was applied for regularization. Another add and normalization step followed this. The final classifier consisted of six stacked encoder layers, enabling the model to capture complex hierarchical patterns in the data.

## D. Classifier Architecture

To build an intelligent classifier that sorts academic sentences into Computer Science, Math, or Science/Physics, we started with a dataset of over 382,000 labeled sentences. We gave the raw text a deep clean—normalizing characters, converting scientific symbols and LaTeX expressions into readable words, removing punctuation and stop words, and lemmatizing each sentence to keep only the essential content. After that, we split the data into training, validation, and test sets while preserving the balance between categories. Each sentence was converted into numbers using TF-IDF, focusing on the top 15,000 most meaningful words. Then, we encoded the class labels into numeric form and used SMOTE to balance all categories equally, making sure no subject had more influence than the others during training.

For the classification model, we designed a custom version of Logistic Regression specifically for this multi-class problem. It used the Softmax function to assign probability scores to each class and aimed to reduce cross-entropy loss during training. To avoid overfitting, we included L2 regularization and used a learning rate that decayed over time, allowing the model to take smaller, more precise steps as it got closer to an optimal solution. We also added early stopping to halt training if no meaningful improvements were seen, making the process more efficient.

Once trained on the balanced dataset, we evaluated the model on unseen data. We looked at overall accuracy, but we didn't stop there—we used a confusion matrix to see where the model struggled or succeeded in classifying each subject. A detailed classification report gave us insights into how precise and consistent the model was for each category. Finally, we plotted learning curves to track how well the model learned over time, helping us visually detect any signs of overfitting or training issues. This end-to-end process ensured that our classifier was both smart and reliable across different academic topics.

## E. Speech to text Architecture

The goal of our project was to develop an effective speech-to-text model by fine-tuning a pre-trained architecture using a custom dataset. We began by setting up a clear project structure, organizing our local environment into directories for cleaned and augmented audio files. The system was built with PyTorch and utilized a CUDA-enabled GPU to accelerate training and evaluation. We sourced our dataset from a provided link containing audio files and corresponding transcriptions and initially worked with 450 samples. These were split into training and validation sets using an 80/20 ratio to ensure generalization, with the split recorded in train.csv and val.csv files for easy access during training.

Next, we prepared the actual audio files by copying them from their original locations into local directories to allow for faster access during training. The pipeline was built around a custom SpeechDataset class derived from PyTorch's Dataset, which loaded and processed the audio files and transcriptions. The audio data was loaded using torchaudio, and any file not in 16kHz was resampled to match the requirements of the Whisper model. The WhisperProcessor handled feature extraction from audio and tokenization of the transcripts. Since both inputs and outputs varied in length, a custom collation function was created to stack inputs and pad labels into equal-length batches, which were then fed into PyTorch DataLoaders with a batch size of 2.

For the model, we used the Whisper architecture via Hugging Face's Whisper For Conditional Generation, initializing it either from a pre-trained checkpoint or a previously fine-tuned version saved locally. Once loaded, the model was transferred to the GPU for improved performance. Although the actual training was skipped due to existing metrics in the log, the system was designed to fine-tune the model using the Adam optimizer with a small learning rate. During training, the model would calculate loss for each batch, back propagate, and update its parameters. Validation occurred after each epoch to assess generalization, and the main metric for performance was Word Error Rate (WER), which measured the differences between predicted and true transcriptions.

After training, the model was evaluated on the validation set and on a specific test audio file. The evaluate-model function

generated transcriptions and compared them with ground-truth text using the jiwer library, which provided both overall WER and a breakdown of insertion, deletion, substitution, and correct hits. A known issue in this run was an error during WER calculation that would need to be fixed in the future. Additionally, sample transcriptions were generated to give a qualitative sense of model performance, and all metrics were saved to a text file for documentation.

To complete the process, we visualized the training and validation losses and WERs across epochs using Matplotlib. These learning curves were helpful for diagnosing issues like overfitting or underfitting and offered a clear picture of the model's progression. Although some steps were skipped due to pre-existing results, the overall setup reflects a solid, end-to-end workflow for fine-tuning a speech-to-text model using Whisper and PyTorch.

### F. RAG Architecture

Retrieval augmented generation has been a core approach in LLM pipelines and recently multi-modal architectures [1, 3]. LLMs in general struggle with issues like hallucination of information or context and relying on the static information upon which the models had been trained on, which makes their knowledge outdated. Additionally, LLMs and other Deep Neural Network architectures are notorious for their highly unexplainable reasoning or decision-making process. RAG architecture integration mitigates such problems while also enhancing performance and output by introducing a retriever that fetches from a constantly updated, embedded database of knowledge [2, 3]. For our RAG model we utilized the lightweight sentence transformer 'all-MiniLM-L6-v2' for embedding our database's text. The model which the RAG works with is Llama's 3.2 3-billion parameter instruct model. However, the RAG will be applied later in the multi-modal pipeline's end for the means of context refinement for the combined output of the models. The retrieval process simply computes the cosine similarity between the user prompt and the RLHF database's prompts – RLHF databases are renowned for their instruction-answer format. After the model computes the similarities it accesses directly the maximum similar prompt's answer in a key-value-like access to produce the final output which it deems closest, the method is simple and effective in producing results with the strict reliance on a database that is completely unseen to the Llama model to produce a more generalized and varied pair of responses which can be manipulated or integrated further for a better response generation.

## VI. RESULTS

The performance evaluation of the STEM-CHARM system encompasses multiple integrated models, each contributing to the overall multi modal tutoring experience. These include modules for speech recognition, subject classification, and multi modal interpretation. Each model was assessed individually to determine its effectiveness in supporting real-time, adaptive educational interactions. The speech-to-text model, responsible for transcribing spoken input into written text,

achieved an accuracy of 73.5 . This result demonstrates acceptable performance for general STEM-related speech inputs, with further improvements anticipated through domain-specific fine-tuning and noise-handling optimizations.
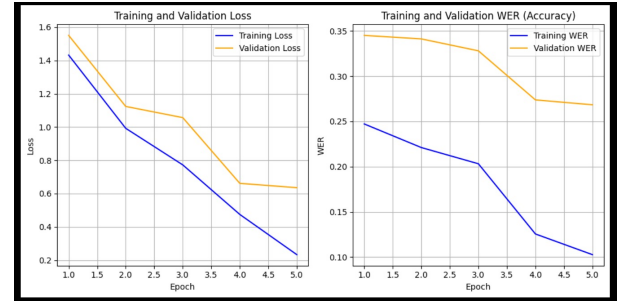


Fig. 1. Shows the datasets that were used

Following transcription, the system employs a text classification model to determine the academic context of the input—first identifying whether the content pertains to Mathematics or Science, and subsequently distinguishing between Physics and Computer Science topics. This model achieved an overall accuracy of 80.36 , with a weighted F1-score of 80.42, indicating strong balance between precision and recall across classes. These metrics reflect the model's effectiveness in guiding the dialogue system toward relevant, context-aware responses based on subject matter.
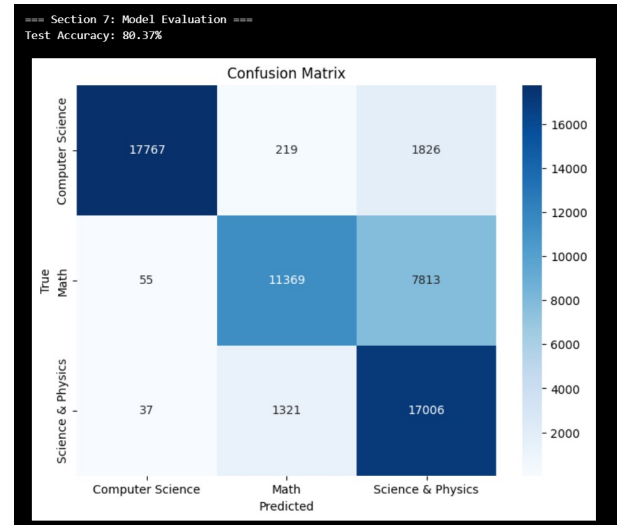


Fig. 2. Shows the confusion matrix

In addition to core components, STEM-CHARM integrates several advanced models to enhance its multimodal understanding and emotional awareness. One such model is a Transformer-based classifier designed to detect the presence of humor in the text. This component plays a key role in shaping the assistant's tone and engagement strategy. Evaluated on a curated dataset of STEM-related utterances annotated for humor, the model achieved an accuracy of 80.33, demonstrating its effectiveness in identifying humorous cues that can be leveraged to deliver more engaging and emotionally intelligent responses.
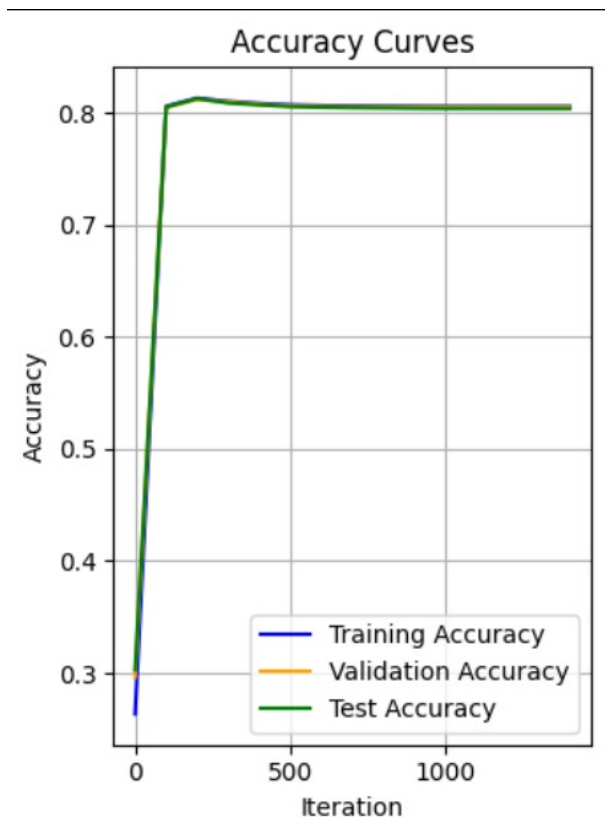
## Accuracy Curves

Fig. 3. shows accuracy curves

To support users who communicate through visual modalities, STEM-CHARM incorporates a YOLO-based sign language recognition model. This model interprets hand gestures captured via webcam and maps them to predefined instructional commands or subject-relevant vocabulary. The sign language module achieved an accuracy of 86.45 on a labeled dataset of common educational signs, enabling responsive interaction for learners with hearing impairments or those using non-verbal communication.

Complementing this is a second YOLO-based facial expression recognition model that integrated with sign language model, which allows the system to assess the learner's emotional state—such as happy, sad, or worry —during interaction and give an answer regards the user mood. This model achieved an accuracy of 85.55.

## VII. CONCLUSION

In conclusion, this research makes important advances in the field of intelligent tutoring systems by introducing STEM-CHARM, Through a pipeline that includes data collection, preprocessing, model integration, and interactive response generation, STEM-CHARM has been developed with a dual focus on technical robustness and learner experience. The system demonstrates that complex subjects—such as physics equations or programming logic—can be taught in warm, engaging, and supportive way.

STEM-CHARM supports voice-to-text input, facial expression recognition, and a sign language interpretation using YOLOV11 model, enabling it to adapt responses based on the user's preferred communication mode—spoken, visual, or gestural. This multimodal framework not only enhances accessibility for users with hearing or speech impairments but also allows the system to respond empathetically by recognizing emotional cues. By aligning its responses with the user's mode of expression, STEM-CHARM fosters a more inclusive, personalized, and engaging learning experience.

### REFERENCES

[1] M. E. Bierzychudek and R. E. Elmquist, "Uncertainty evaluation in a two-terminal cryogenic current comparator," *IEEE Trans. Instrum. Meas.*, vol. 58, no. 4, pp. 1170 – 1175, April 2009.

[2] D. G. Jarrett and R. E. Elmquist, "Settling time of high-value standard resistors," CPEM 2004 Conf. Digest, p. 522, June 2004.

[3] H. L. Krauss, C. W. Bostian, and F. H. Raab, *Solid State Radio Engineering*, New York: J. Wiley & Sons, 1980.