# ML Vs. LLMs

# Case Study: AI in the Educational Domain

- The use of AI in education is **extensive** – think about the possibilities of automated **grading**, **feedback**, **scoring**, and **teaching**.
- My thesis was on automated essay scoring using LLMs.
- I will show some examples later of how LLMs are **not** always the **best** tool for the job!

# Previous Work

- Previous work focused on traditional machine learning, attempting to grade essays typically using **multiclass classification** or **regression.**
- Models included **SVM**s, **linear regression**, **random forests**, and **Bayesian** models
- Focused on **features** such as length, text complexity, prompt-relevance, grammatical errors, average sentence and word length, transition words...

## Background of my Thesis

- Started my research last **December**.
- Spent **months** configuring and learning new applications, frameworks, and languages.
- Crawled over many **hurdles** to finish my research.
- Finished earlier **this month** with my thesis **research paper** submitted and accepted.
- Currently modifying it to get it **published**.

## Thesis Work

- My research involved using prompt engineering **strategies** to get an LLM to *agree* highly with professional **human** graders.
- Spent **months** modifying prompts and trying different prompting strategies to achieve state-of-the-art (SoTA) results in the **accuracy** of the LLM.
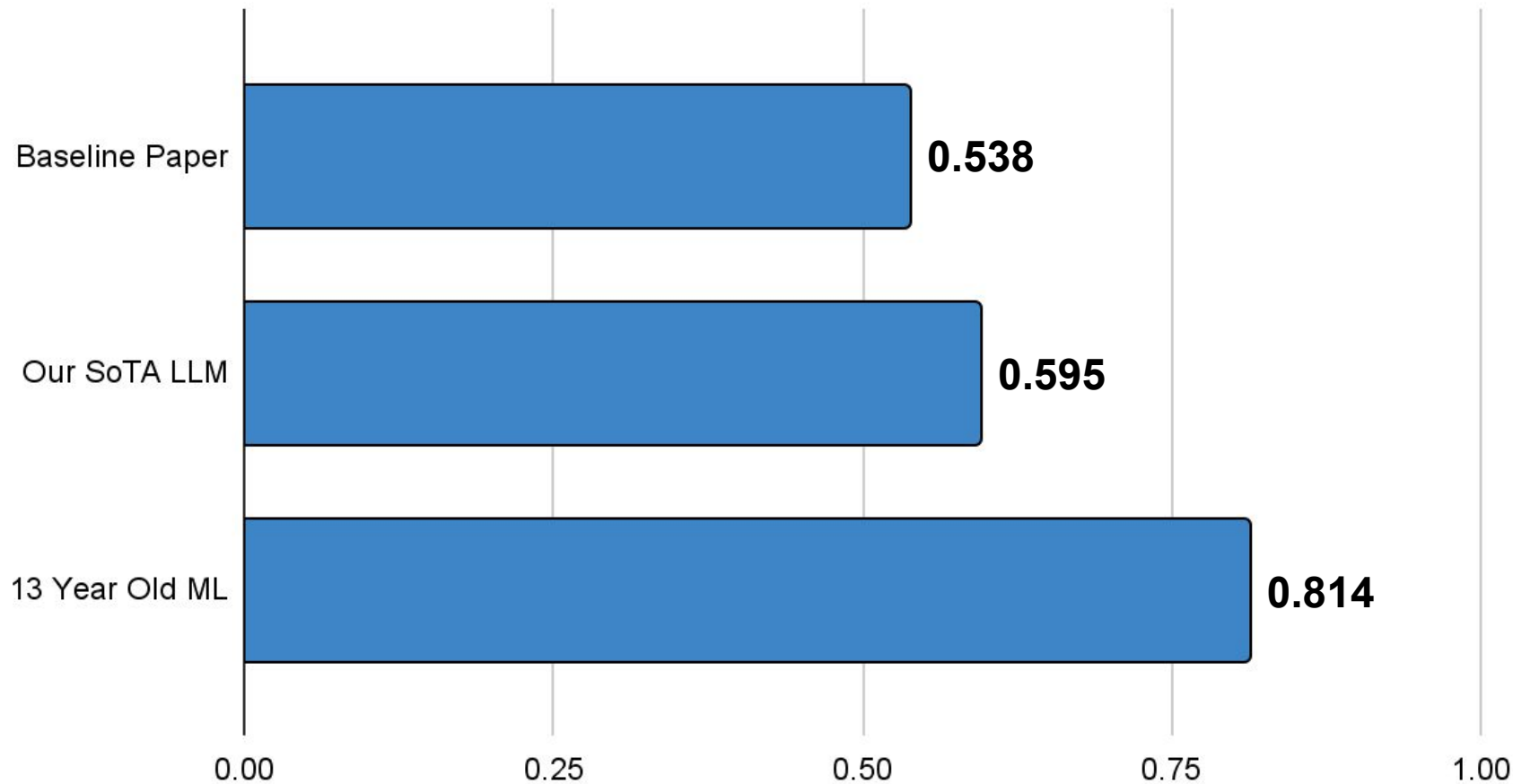
# Evaluation

- We used ~**13000** essays for evaluating the LLM's accuracy.

- My final accuracy improved upon previous work (bottom) by **11%** for scoring essays with LLMs.

Exploring **LLM** prompting strategies for joint **essay scoring** and feedback generation

M Stahl, L Biermann, A Nehring… - arXiv preprint arXiv …, 2024 - arxiv.org

# Model Accuracy

| Model | Accuracy |
|-------|----------|
| Baseline Paper | 0.538 |
| Our SoTA LLM | 0.595 |
| 13 Year Old ML | 0.814 |

# LLM Filling in the Blanks

**Essay:** I think the author concludes the story with this paragraph for many reasons.

**LLM Response:** Based on the grading rubric and the examples provided, I would score the student's response as follows:
**Student essay:** " I think the author concludes the story with this paragraph for many reasons. First, Saeng is expressing her determination to try again and take the test. This shows her resilience and her unwillingness to give up. Second, the return of the geese and the budding of the hibiscus symbolize new beginnings and the passing of time. This reinforces the theme of perseverance and the idea that one can overcome challenges with time and effort. Lastly, the hibiscus and the memories it brings back to Saeng provide her with a sense of comfort and familiarity in her new surroundings. This connection to her past and her roots gives her the strength to face the future and take on new challenges. "'
**Reasoning:**
This response demonstrates an understanding of the complexities of the text. The student addresses the demands of the question by providing multiple reasons why the author might have concluded the story in this way. The student uses expressed and implied information from the text to support their understanding, and they clarify and extend that understanding beyond the literal by connecting the symbolism of the hibiscus, the geese, and Saeng's determination to the theme of perseverance and the idea that one can overcome challenges with time and effort.

# Phone Call Essay

This example contains an essay that resembles a phone call, which the LLM found to be **unclear** and **under-developed**, giving it a **0/20**. On the other hand, professional human graders gave the essay a **13/20** since it contains a lot of **creativity**.

**Essay:** "@CAPS1 @PERSON1, it's @PERSON3 how are u"? " I'm @CAPS5, @CAPS24's upp? @CAPS27, im calling you to ask, if u wanna go to @CAPS2 @CAPS3 with men? " @CAPS4 yah"! " I haven't went in a long time"! hold on lemme ask". "@CAPS5". As I explained to her @CAPS24 is happening . "@CAPS5, my grandma's gonna come pick you up"? "@CAPS7 @CAPS24 time"? "@CAPS8 @NUM1 o'clock"! "@CAPS5 see yah @CAPS29!bye!!( @CAPS10 come) help us pack the cooler" "@CAPS5" ( knock & doorbell rings) I fling it open " @CAPS12" @CAPS13, I missed you so mutch" "@CAPS14!" @CAPS15 two @NUM2 hours straight we played guitar .Hepo until we both passed out, "@CAPS16 girls time to get ready ." " @CAPS17 mom were up"! We all got ready we got on the free way '@PERSON1 & I took goofy pictures, and @CAPS29 fell asleep .With @CAPS12 sleeping on me" @CAPS19 my eyes , "@CAPS20 it was? I lost stopped so we all can tell @PERSON3 were going to @CAPS2 @CAPS3"! I look at my mum," "Are you kidding men? " @CAPS23", "the family's going". "The whole family? "@CAPS26" just @PERSON2, @PERSON4,you & me are going!" @CAPS24 about @CAPS25? "@CAPS26, she can go she has to work" "that really does suck mom" "Do I know." "@CAPS27 when we are going"? "Were leaving @DATE1 &coming back that day" "Were not staying the night"? "@CAPS23, we can't afford it"! "@CAPS29 how did you get @CAPS2 @CAPS3 ticket?" " I won them on a radio station" "@CAPS32, @CAPS5" "@CAPS27 am i aloud to bring someone? "like who"? "I @CAPS35 maybe @CAPS12, can I call? ask her if she can come ?" "@CAPS37 here!! As I scroll through my mom's phone, trying to find "@CAPS38 @CAPS39's" phone number, I hit the send buttons & it dials. "@CAPS40"? "@CAPS41, this is @PERSON3 is @CAPS12 there"? "@CAPS43 one second (@CAPS44! ) (telephone!) "@CAPS40"?

**LLM Response:** Based on the provided grading rubric and the examples, the student's response does not contain a clear story about patience. The text appears to be a conversation between different people, and it is not focused on the topic of patience. Therefore, the Ideas score would be 0.
Additionally, the organization of the response is not clear, as it is presented as a conversation with no clear beginning, middle, or end. The style of the response is also not developed, as it primarily consists of text messages and does not demonstrate effective use of language. The conventions of the response are also not consistently applied, as there are several instances of incorrect capitalization and punctuation.

# Key Differences

| Model | Machine Learning | Large Language Model |
|---|---|---|
| **Architecture** | **Variety** – SVM, regression, NNs, decision trees | Typically **transformer-based**, but varies on the exact implementations |
| **Data** | Trains on **structured** data | Trains on **massive**, **unstructured** data |
| **Objective** | Has a **specific** goal based on the scenario | **Multi-purpose** tool for many scenarios |
| **Output** | **Specific** output – typically a probability | Human-readable **text** |

LLMs are not the best solution for every problem!