

# Project Notes:

**Project Title:**

**Name:**

**Note Well: There** are NO SHORT-cuts to reading journal articles and taking notes from them. Comprehension is paramount. You will most likely need to read it several times, so set aside enough time in your schedule.

## **Contents:**

<b>Knowledge Gaps:</b>	<b>1</b>
<b>Literature Search Parameters:</b>	<b>2</b>
<b>Article #1 Notes: Title</b>	<b>3</b>
<b>Article #2 Notes: Title</b>	<b>4</b>
<b>Article #1 Notes: Title</b>	<b>5</b>

## Knowledge Gaps:

This list provides a brief overview of the major knowledge gaps for this project, how they were resolved and where to find the information.

<b>Knowledge Gap</b>	<b>Resolved By</b>	<b>Information is located</b>	<b>Date resolved</b>
Tuning large language models	Watching YouTube videos	Video links in project logbooks	9/29/24
Privacy of data	Reading patents	Patents at bottom of project notes	10/10/24
Limitations and effectiveness of AI in math education	Journal articles	First few articles up to article 7 in project notes	9/27/24

## Literature Search Parameters:

These searches were performed between (Start Date of reading) and XX/XX/2019.

List of keywords and databases used during this project.

Database/search engine	Keywords	Summary of search
Gordon Library	AI and (tutor* or teach*)	Found most of my articles from this and following up with references
Google patent search	Artificial intelligence in education	Found next to nothing useful
Google patent search	Artificial Intelligence	Found general patents related to data collection and privatization in AI

## Tags:

Tag Name	

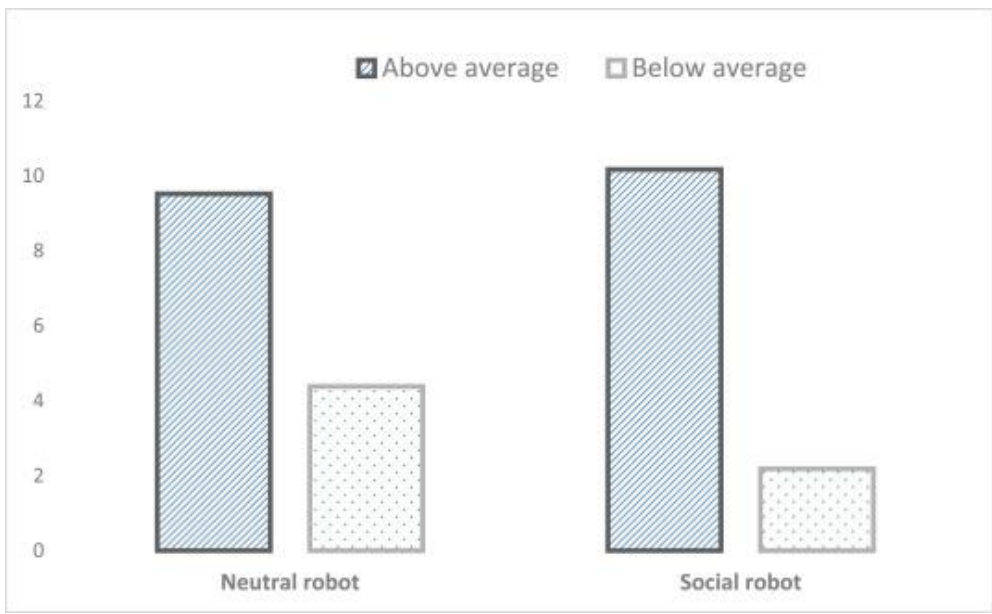
## Article #1 Notes: Title

Article notes should be on separate sheets

**KEEP THIS BLANK AND USE AS A TEMPLATE**

<b>Source Title</b>	
<b>Source citation (APA Format)</b>	
<b>Original URL</b>	
<b>Source type</b>	
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	
<b>Research Question/Problem/ Need</b>	
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	

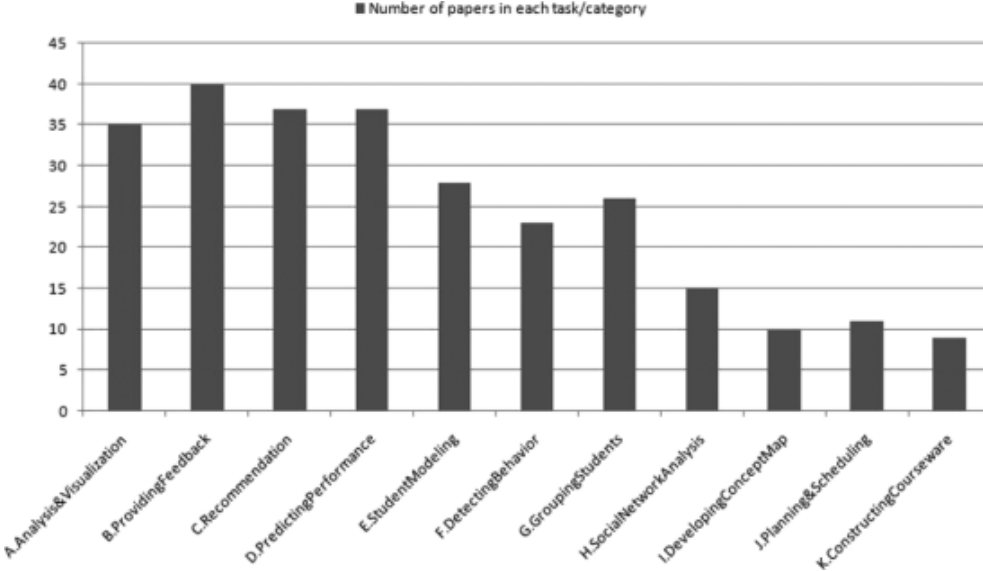
## Article #1 Notes:

<b>Source Title</b>	Robot tutor and pupils' educational ability: Teaching the times tables									
<b>Source citation (APA Format)</b>	Konijn, E. A., & Hoorn, J. F. (2020). Robot tutor and pupils' educational ability: Teaching the times tables. <i>Computers and Education</i> , 157, 103970-. <a href="https://doi.org/10.1016/j.compedu.2020.103970">https://doi.org/10.1016/j.compedu.2020.103970</a>									
<b>Original URL</b>	<a href="https://doi.org/10.1016/j.compedu.2020.103970">https://doi.org/10.1016/j.compedu.2020.103970</a>									
<b>Source type</b>	Journal Article									
<b>Keywords</b>	Robot tutor, tutoring, social robots									
<b>#Tags</b>										
<b>Summary of key points + notes (include methodology)</b>	Physically present robots have been shown to be quite effective in tutoring, even compared to actual teachers. This study assessed the improvement of multiplication table recall in 86 elementary school students. It found that the robots were shown to be capable of improving students' scores regardless of their behavior (social or neutral), but ethics and other factors are still a concern.									
<b>Research Question/Problem/Need</b>	How effective are robot tutors in improving elementary school kids' multiplication skills?									
<b>Important Figures</b>	 <p>The bar chart displays the number of students who improved their scores on a multiplication test, categorized by robot type (Neutral robot and Social robot) and performance level (Above average and Below average). The Y-axis represents the number of students, ranging from 0 to 12. The legend indicates that hatched bars represent 'Above average' and dotted bars represent 'Below average'.</p> <table border="1"> <thead> <tr> <th>Robot Type</th> <th>Above average</th> <th>Below average</th> </tr> </thead> <tbody> <tr> <td>Neutral robot</td> <td>9</td> <td>4</td> </tr> <tr> <td>Social robot</td> <td>10</td> <td>2</td> </tr> </tbody> </table> <p>Y-axis is score improvement on test. Neutral robot was better for below average</p>	Robot Type	Above average	Below average	Neutral robot	9	4	Social robot	10	2
Robot Type	Above average	Below average								
Neutral robot	9	4								
Social robot	10	2								

	students while robot type didn't seem to matter much for advanced students. However, advanced students benefited the most from robots in general.
<b>VOCAB: (w/definition)</b>	Pedagogical – related to teaching/learning
<b>Cited references to follow up on</b>	Leyzberg, D., Spaulding, S., & Scassellati, B. (2014). Personalizing robot tutors to individuals' learning differences. In Proceedings of the 2014 ACM/IEEE international conference on Human-robot interaction (pp. 423–430). ACM. <a href="https://doi.org/10.1145/2559636.2559671">https://doi.org/10.1145/2559636.2559671</a> .
<b>Follow up Questions</b>	<ol style="list-style-type: none"> <li>1. How would students of different age groups compare in their response to robot tutoring?</li> <li>2. How can teachers and robots work in tandem to maximize success in the classroom?</li> <li>3. Why do students not respond any better to more supportive feedback from robots?</li> </ol>

## Article #2 Notes:

<b>Source Title</b>	Educational Data Mining: A Review of the State of the Art
<b>Source citation (APA Format)</b>	Romero, C. & Ventura, S. (2010). Educational Data Mining: A Review of the State of the Art. <i>IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)</i> , 40(6), 601-618. <a href="https://doi.org/10.1109/TSMCC.2010.2053532">https://doi.org/10.1109/TSMCC.2010.2053532</a>
<b>Original URL</b>	<a href="https://doi.org/10.1109/TSMCC.2010.2053532">https://doi.org/10.1109/TSMCC.2010.2053532</a>
<b>Source type</b>	Journal article
<b>Keywords</b>	Educational Data Mining
<b>#Tags</b>	

<p><b>Summary of key points + notes (include methodology)</b></p>	<p>The article went over EDM itself, including defining it, reviewing current research, and looking at promising future research.</p> <p>Past research in EDM mainly looked at predicting student performance. EDM has risen in popularity with the advent of LMSs (learning management systems) creating lots of data online about student-teacher interaction and education in general.</p> <p>Educational data mining has been used to visualize data to help educators and administrators, analysis of that data, provide feedback to teachers, and help students, as well as predicting their success, creating models of students, and detecting unwanted behaviors (cheating, dropping out, etc.).</p> <p>Also, grouping students to aid in personalization, social network analysis to help students find relevant information, concept maps to help teachers teach, creating course material, and scheduling.</p> <p>Future work in EDM would include making it more accessible, integrating it with LMSs, standardization, and fine-tuning data mining to be more education specific.</p>																								
<p><b>Research Question/Problem/ Need</b></p>	<p>How is educational data mining used and how does it impact students and teachers?</p>																								
<p><b>Important Figures</b></p>	 <p>■ Number of papers in each task/category</p> <table border="1"> <thead> <tr> <th>Task/Category</th> <th>Number of Papers</th> </tr> </thead> <tbody> <tr> <td>A. Analysis &amp; Visualization</td> <td>35</td> </tr> <tr> <td>B. Providing Feedback</td> <td>40</td> </tr> <tr> <td>C. Recommendation</td> <td>37</td> </tr> <tr> <td>D. Predicting Performance</td> <td>37</td> </tr> <tr> <td>E. Student Modeling</td> <td>28</td> </tr> <tr> <td>F. Detecting Behavior</td> <td>23</td> </tr> <tr> <td>G. Grouping Students</td> <td>26</td> </tr> <tr> <td>H. Social Network Analysis</td> <td>15</td> </tr> <tr> <td>I. Developing Concept Map</td> <td>10</td> </tr> <tr> <td>J. Planning &amp; Scheduling</td> <td>11</td> </tr> <tr> <td>K. Constructing Courseware</td> <td>9</td> </tr> </tbody> </table> <p>There are a variety of use cases for educational data mining, y axis is number of papers relating to each topic, out of the 300 total that were reviewed</p>	Task/Category	Number of Papers	A. Analysis & Visualization	35	B. Providing Feedback	40	C. Recommendation	37	D. Predicting Performance	37	E. Student Modeling	28	F. Detecting Behavior	23	G. Grouping Students	26	H. Social Network Analysis	15	I. Developing Concept Map	10	J. Planning & Scheduling	11	K. Constructing Courseware	9
Task/Category	Number of Papers																								
A. Analysis & Visualization	35																								
B. Providing Feedback	40																								
C. Recommendation	37																								
D. Predicting Performance	37																								
E. Student Modeling	28																								
F. Detecting Behavior	23																								
G. Grouping Students	26																								
H. Social Network Analysis	15																								
I. Developing Concept Map	10																								
J. Planning & Scheduling	11																								
K. Constructing Courseware	9																								
<p><b>VOCAB: (w/definition)</b></p>	<p>Educational data mining – a field of data science that uses a variety of data analysis and machine learning techniques to solve problems in the world of educational research</p> <p>Learning management system – An online platform that connects students and teachers and allows teachers to assign materials and track student grades and progress (think Schoology or Canvas)</p>																								
<p><b>Cited references to follow up on</b></p>																									

**Follow up Questions**

How well does this analysis hold up in the current state of artificial intelligence?  
How much does EDM have on the seemingly less intuitive use cases such as scheduling and creating content maps (shouldn't teachers already be good at those things)?

Why didn't constructing courseware have much research into it when it is likely a huge benefit to EDM as it would take lots of workload off teachers?



## Article #3 Notes: Title

<b>Source Title</b>	Music teachers' labeling accuracy and quality ratings of lesson plans by artificial intelligence (AI) and humans
<b>Source citation (APA Format)</b>	Cooper, P. K. (2024). Music teachers' labeling accuracy and quality ratings of lesson plans by artificial intelligence (AI) and humans. <i>International Journal of Music Education</i> , 0(0). <a href="https://doi.org/10.1177/02557614241249163">https://doi.org/10.1177/02557614241249163</a>
<b>Original URL</b>	<a href="https://doi.org/10.1177/02557614241249163">https://doi.org/10.1177/02557614241249163</a>
<b>Source type</b>	Journal article
<b>Keywords</b>	Lesson plan, music education
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>A survey was sent out to US music teachers. Most of them had over ten years of experience and more than half of them had experience with AI, implying that they would be suitable participants in the study.</p> <p>Overall, the teachers were 55% accurate on average in labeling whether lesson plans were generated by AI or humans. This was not a statistically significant result. Also, AI generated content and human content were ranked similarly in usefulness. Using multiple regression, they did find that they could predict the accuracy of a teacher's guess based on their personal use of AI, their ratings of usefulness for both humans and AI, and how useful they thought AI would be in the future.</p> <p>The analysis found that overall, teachers were unsuccessful in predicting whether music lesson plans were generated by AI much better than chance.</p>
<b>Research Question/Problem/Need</b>	How well can AI generated music lesson plans be distinguished from human made ones?
<b>Important Figures</b>	

<b>VOCAB: (w/definition)</b>	Intelligent Tutor System – a program that individually tutors a student in a custom manner, similarly to a human tutor
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	What, specifically, do these lesson plans entail? Can this be taken one step further, using AI to assign homework and tests? Can these lesson plans be modified on the fly as is so common in private lessons by AI?

## Video #1 Notes:

<b>Source Title</b>	Introduction to Generative AI
<b>Source citation (APA Format)</b>	Google Cloud Tech. (2023, May 8). <i>Introduction to Generative AI</i> [Video]. YouTube. <a href="https://youtu.be/G2fqAlgmoPo?si=KH73Mt7LrbO5ryt5">https://youtu.be/G2fqAlgmoPo?si=KH73Mt7LrbO5ryt5</a>
<b>Original URL</b>	<a href="https://youtu.be/G2fqAlgmoPo?si=KH73Mt7LrbO5ryt5">https://youtu.be/G2fqAlgmoPo?si=KH73Mt7LrbO5ryt5</a>
<b>Source type</b>	YouTube video
<b>Keywords</b>	Generative AI
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>Generative AI is a type of artificial intelligence that creates content</p> <p>Artificial intelligence is a field of computer science, while machine learning is the subfield that involves creating models that can perform “intelligent” tasks</p> <p>Neural networks – layers of “neurons” (nodes) that make up a deep learning model and can use labeled or unlabeled data to process patterns (semisupervised learning)</p> <p>Transformers: consist of encoders and decoders that convert input data into relevant output data, transformers sometimes make hallucinations (incorrect outputs)</p> <p>Prompt design – creating a prompt that gives the desired output from a generative AI model</p> <p>Variety of model types: text-to-text, text-to-task, text-to-image, foundational</p> <p>Foundation models can be fine-tuned to a variety of tasks</p> <p>Generative AI Studio – Google's gen AI platform that allows developers to create generative AI. Has a library of pre-trained models and has tools for fine-tuning, deployment, and more</p> <p>Generative AI App Builder – create gen AI apps without code (drag and drop, might be limited)</p>
<b>Research Question/Problem/Need</b>	N/A
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	<p>Neural networks, transformers (definitions in summary)</p> <p>Machine learning – models that can “learn” by changing their parameters and the connections in their neural networks to create more desirable outputs</p>
<b>Cited references to follow up</b>	

<b>on</b>	
<b>Follow up Questions</b>	What are the different methods for fine tuning? How does prompt design work? Is that different from prompt engineering?

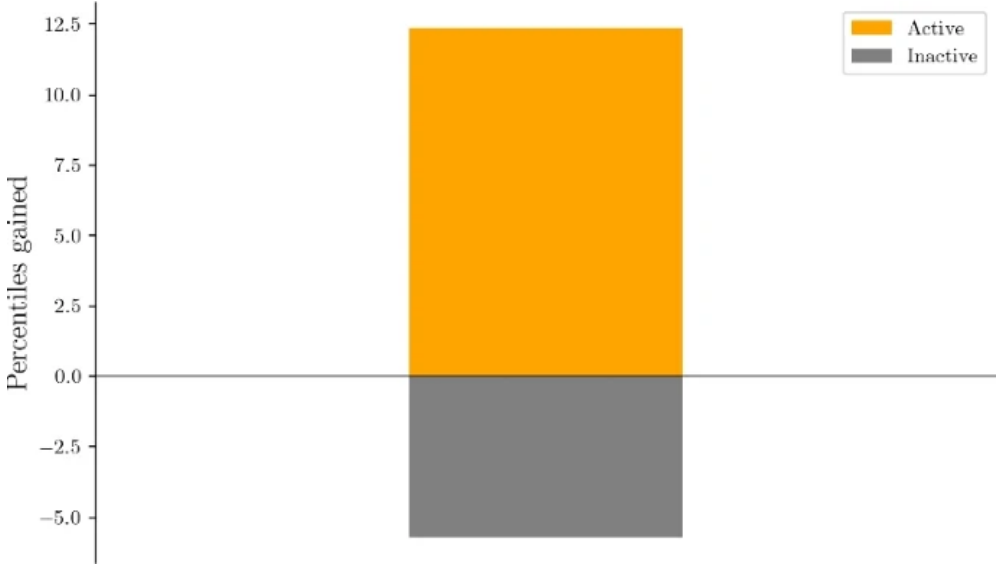
## Video #2 Notes:

<b>Source Title</b>	Introduction to large language models
<b>Source citation (APA Format)</b>	Google Cloud Tech. (2023, May 8). <i>Introduction to Large Language Models</i> [Video]. YouTube. <a href="https://youtu.be/zizonToFXDs?si=SJnO9BN6-vkOfc0Q">https://youtu.be/zizonToFXDs?si=SJnO9BN6-vkOfc0Q</a>
<b>Original URL</b>	<a href="https://youtu.be/zizonToFXDs?si=SJnO9BN6-vkOfc0Q">https://youtu.be/zizonToFXDs?si=SJnO9BN6-vkOfc0Q</a>
<b>Source type</b>	YouTube video
<b>Keywords</b>	Large language model
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>LLMs are general models that can be fine tuned to specific use cases  Tuned usually using domain specific data which is in a much smaller quantity than the data used to create the general model  Parameters – the “knowledge” that the model gathered from the data. LLMs usually have many parameters  I probably don’t have the resources required to create a LLM (need tons of data), however, I could probably get the data required to tune one to educational purposes  LLM performance is increasing over time  LLM development does not require as much coding knowledge as regular machine learning development  Prompt engineering differs from prompt design in that the goal is to improve the performance of the model, may use known effective keywords or give examples of the correct output</p>
<b>Research Question/Problem/Need</b>	N/A
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	<p>Fine-tuning – modifying a pre-trained large language model to be more suited to a specific domain or field  Parameters – act as guidelines that affect a model’s output and are determined during the pre-training process. LLMs often have millions or even billions of parameters</p>
<b>Cited references to follow up</b>	

<b>on</b>	
<b>Follow up Questions</b>	Need more information on fine-tuning.

## Article #4 Notes:

Source Title	Effective learning with a personal AI tutor: A case study
Source citation (APA Format)	Baillifard, A., Gabella, M., Lavenex, P. B., Martarelli, C. S. (2024). Effective learning with a personal AI tutor: A case study. <i>Educ Inf Technol</i> . <a href="https://doi.org/10.1007/s10639-024-12888-5">https://doi.org/10.1007/s10639-024-12888-5</a>
Original URL	<a href="https://doi.org/10.1007/s10639-024-12888-5">https://doi.org/10.1007/s10639-024-12888-5</a>
Source type	Journal Article
Keywords	AI Tutor, AIEd, Learning Sciences, Personalization, Intelligent Tutoring Systems
#Tags	
Summary of key points + notes (include methodology)	<p>Known benefits of AI in education:</p> <ol style="list-style-type: none"> <li>1. Capable of predicting student outcomes and creating profiles of students</li> <li>2. Good at testing students, taking workload off teachers</li> <li>3. Can be personalized to help a wider range of students more effectively</li> <li>4. Intelligent Tutoring Systems that simulate real 1:1 tutoring experiences</li> </ol> <p>Methodology:</p> <p>They used a tutoring app developed by MAGMA Learning that uses personalized tutoring. Also, they used GPT-3 to create a set of questions that would be used. A neural network was used to predict the probability of a student answering a question correctly (called the "grasp") and thus selected the best questions for the student to practice with.</p> <p>App was tested in an online neuroscience college course. Class was managed through an LMS called Moodle. There was also a parallel course taken by most of the same students at the same time, but they did not have the app for that course. 43 of the 51 students enrolled in the course did use the app, students could use the app as they pleased. 47 students took both final exams and 40 of those were using the app. They investigated performance on the final exam and use of the app. Also, performance on Moodle and the "grasp" compared to the final exam. They also compared these to the course with the final exam but no app used.</p> <p>Comparing active users to inactive users:</p> <p>Average increase in score of 0.71 for active user (test was on a scale from 1-6). They used different thresholds for "active" but found this to be the average increase, on average</p> <p>Comparing the two courses:</p> <p>Active users had an average grade increase of 12.4% in the course where the app was present compared to the parallel course. Inactive users had a decrease of 5.7%.</p>

	<p>Active users of Moodle were not shown to do any better on the final exam than non-Moodle users.</p> <p>There was also a strong correlation between grasp prediction and exam grade.</p> <p>Overall, app usage and grade were positively correlated.</p>
<b>Research Question/Problem/ Need</b>	<p>How well does AI in education work with known learning sciences?</p>
<b>Important Figures</b>	 <p>Active users did better on the final exam than inactive users. Inactive users even seemed to do worse on the course that offered the AI app as a study option.</p>
<b>VOCAB: (w/definition)</b>	<p>Retrieval practice – recalling information from memory without having it available to help you remember it</p> <p>Natural language processing – using machine learning to allow computers to understand and create human language</p>
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	<p>How can be sure it was causation and not just correlation (ambitious students would happen to use the app more and have higher grades)?</p> <p>How were the neural networks trained?</p> <p>How important is the personalization aspect of the AI tutor itself, rather than just having the tutor?</p>



## Article #5 Notes:

<b>Source Title</b>	Editorial Note: From Conventional AI to Modern AI in Education: Re-examining AI and Analytic Techniques for Teaching and Learning
<b>Source citation (APA Format)</b>	Xie, H., Hwang, G. J., & Wong, T. L. (2021). Editorial Note: From Conventional AI to Modern AI in Education: Re-examining AI and Analytic Techniques for Teaching and Learning. <i>Educational Technology &amp; Society</i> , 24(3), 85-88. <a href="https://doi.org/10.30191/ETS.202107_24(3).0006">https://doi.org/10.30191/ETS.202107_24(3).0006</a>
<b>Original URL</b>	<a href="https://doi.org/10.30191/ETS.202107_24(3).0006">https://doi.org/10.30191/ETS.202107_24(3).0006</a> <a href="https://link.gale.com/apps/doc/A668399451/AONE?u=mlic_worpoly&amp;sid=bookmark-AONE&amp;xid=6fa17e23">https://link.gale.com/apps/doc/A668399451/AONE?u=mlic_worpoly&amp;sid=bookmark-AONE&amp;xid=6fa17e23</a> (need this link because otherwise I can't access)
<b>Source type</b>	Journal Article
<b>Keywords</b>	Modern AI, AI transformation, Deep neural networks, Analytic techniques
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	It's an editorial note for the issue itself which discusses technology in education. Modern AI uses deep neural networks, while traditional AI uses statistical models. There is limited research on modern AI's use in education as most of it in the past has been with traditional AI. Teachers and AI developers don't know much about each other's domains, so it is hard to connect the two for effective AI education. Precision education is the next big step in the use of AI in education, along with more general predictions, and using AI for new apps.
<b>Research Question/Problem/Need</b>	How effective has artificial intelligence been in education thus far?
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Convolution neural network – neural networks that are better at using image and audio inputs ( <a href="#">What are Convolutional Neural Networks?   IBM</a> ) Generative adversarial network – made up of a generator and discriminator neural network, the generator attempts to create data identical to training data until it can fool the discriminator; used for unsupervised learning ( <a href="#">Generative Adversarial Network (GAN) - GeeksforGeeks</a> ) Precision education – identifying students who are at risk of failure, dropping out, etc., and giving them the guidance and resources needed to succeed accordingly

<p><b>Cited references to follow up on</b></p>	<p>Chen, X., Xie, H., &amp; Hwang, G. J. (2020a). A Multi-perspective study on Artificial Intelligence in Education: grants, conferences, journals, software tools, institutions, and researchers. <i>Computers and Education: Artificial Intelligence</i>, 1, 100005. - like the current article in that the references in this one may be more useful than the information itself; may lead to studies or developments more pertinent to my topic</p> <p>Yang, S. J., Ogata, H., Matsui, T., &amp; Chen, N. S. (2021). Human-centered artificial intelligence in education: Seeing the invisible through the visible. <i>Computers and Education: Artificial Intelligence</i>, 2, 100008.</p> <p>Wang, J., Xie, H., Wang, F. L., Lee, L. K., &amp; Au, O. T. S. (2021). Top-n personalized recommendation with graph neural networks in MOOCs. <i>Computers and Education: Artificial Intelligence</i>, 2, 100010.</p> <p>Almohammadi, K., Hagra, H., Alghazzawi, D., &amp; Aldabbagh, G. (2016). Users- centric adaptive learning system based on interval type-2 fuzzy logic for massively crowded E-learning platforms. <i>Journal of Artificial Intelligence and Soft Computing Research</i>, 6(2), 81-101.</p>
<p><b>Follow up Questions</b></p>	<p>Three years later, how true do these gaps hold up?</p> <p>Is the gap between AI experts and educators specific to the education field only, or does this problem exist across many domains?</p> <p>Why has predicting students' risk or classifying students been so researched compared to other needs in education?</p>

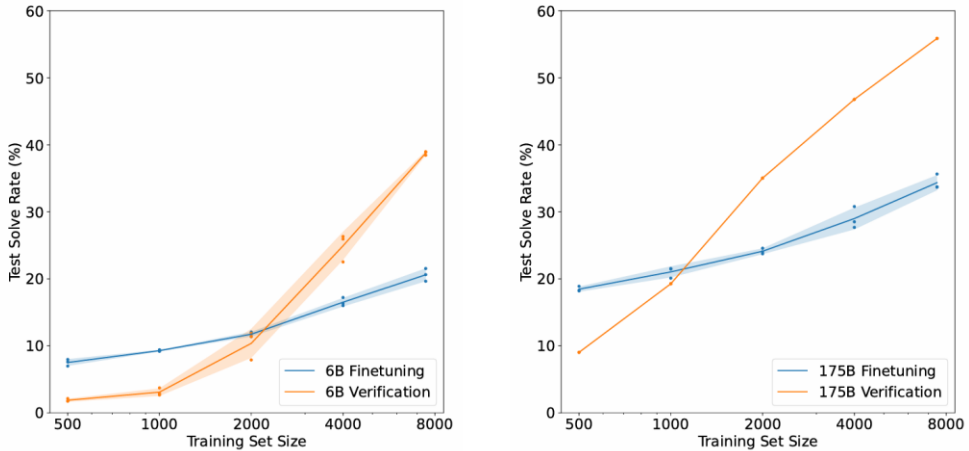
## Article #6 Notes:

Source Title	Evaluating language models for mathematics through interactions
Source citation (APA Format)	Collins, K. M., Jiang, A. Q., Frieder, S., Wong, L., Zilka, M., Bhatt, U., Lukasiewicz, T., Wu, Y., Tenenbaum, J. B., Hart, W., Gowers, T., Li, W., Weller, A., & Jamnik, M. (2024). Evaluating language models for mathematics through interactions. <i>Proceedings of the National Academy of Sciences of the United States of America</i> , 121(24). <a href="https://doi.org/10.1073/pnas.2318124121">https://doi.org/10.1073/pnas.2318124121</a>
Original URL	<a href="https://doi.org/10.1073/pnas.2318124121">https://doi.org/10.1073/pnas.2318124121</a>
Source type	Journal article
Keywords	Dynamic (evaluation) - observing how people and language models interact over the course of an entire “conversation” rather than a snapshot evaluation
#Tags	
Summary of key points + notes (include methodology)	<p>First, they developed a platform called CheckMate that allowed people to interact with LLM chatbots and rate them individually or comparatively. They tested how people used InstructGPT, ChatGPT, and GPT-4. For individually, they were allowed to use a model to help solve a math problem and then rate each step of the process. Comparatively, they ranked the different models without knowing which was which. Participants’ experience ranged from undergraduate students to college professors, however data on participants was not collected beyond this. Specifically, they asked participants to prove undergraduate level theorems and allowed them to use AI any way they wished, as they wanted to see how people naturally used it. They were asked to rate perceived helpfulness along with mathematical correctness.</p> <p>They used dynamic evaluation – observing a model’s entire interaction with a person – rather than static evaluation.</p> <p>GPT-4 was ranked the highest overall and received the highest helpfulness and correctness ratings. Models built for chatting (GPT-4 and ChatGPT) were ranked much better than those that aren’t (InstructGPT).</p> <p>The correlation between helpfulness and correctness was decent but not 100% - sometimes it could be helpful but wrong (contains decent ideas) or correct but unhelpful (verbosity).</p> <p>Currently, measuring helpfulness and correctness cannot be done computationally or automatically, they must be determined by humans. This is one of the reasons humans were used to test these models.</p> <p>GPT-4 often struggled with arithmetic mistakes. In general, LLMs were found to struggle with <b>algebra, being too wordy, and reliance on memorized solutions</b>.</p>
Research Question/Problem/ Need	How good are LLMs at assisting people with undergraduate level math problems?

<p><b>Important Figures</b></p>	<div style="display: flex; flex-wrap: wrap;"> <div style="width: 50%;"> <p><b>A</b></p> <table border="1"> <caption>Rank Distribution Data</caption> <thead> <tr> <th>Model</th> <th>Rank: 3</th> <th>Rank: 2</th> <th>Rank: 1</th> </tr> </thead> <tbody> <tr> <td>InstructGPT</td> <td>11</td> <td>2</td> <td>2</td> </tr> <tr> <td>ChatGPT</td> <td>4</td> <td>6</td> <td>5</td> </tr> <tr> <td>GPT-4</td> <td>0</td> <td>6</td> <td>9</td> </tr> </tbody> </table> </div> <div style="width: 50%;"> <p><b>B</b></p> <table border="1"> <caption>Correctness Distribution Data</caption> <thead> <tr> <th>Model</th> <th>0</th> <th>1</th> <th>2</th> <th>3</th> <th>4</th> <th>5</th> <th>6</th> </tr> </thead> <tbody> <tr> <td>InstructGPT</td> <td>0</td> <td>22</td> <td>12</td> <td>28</td> <td>4</td> <td>1</td> <td>23</td> </tr> <tr> <td>ChatGPT</td> <td>0</td> <td>8</td> <td>8</td> <td>19</td> <td>8</td> <td>10</td> <td>38</td> </tr> <tr> <td>GPT-4</td> <td>0</td> <td>2</td> <td>4</td> <td>12</td> <td>8</td> <td>8</td> <td>41</td> </tr> </tbody> </table> <table border="1"> <caption>Helpfulness Distribution Data</caption> <thead> <tr> <th>Model</th> <th>0</th> <th>1</th> <th>2</th> <th>3</th> <th>4</th> <th>5</th> <th>6</th> </tr> </thead> <tbody> <tr> <td>InstructGPT</td> <td>2</td> <td>21</td> <td>16</td> <td>17</td> <td>15</td> <td>10</td> <td>10</td> </tr> <tr> <td>ChatGPT</td> <td>1</td> <td>9</td> <td>16</td> <td>12</td> <td>15</td> <td>15</td> <td>22</td> </tr> <tr> <td>GPT-4</td> <td>1</td> <td>4</td> <td>6</td> <td>10</td> <td>10</td> <td>23</td> <td>22</td> </tr> </tbody> </table> </div> <div style="width: 50%;"> <p><b>C</b></p> </div> <div style="width: 50%;"> <p><b>D</b></p> </div> </div> <p>This figure shows that overall, GPT-4 was the best in every aspect.</p>	Model	Rank: 3	Rank: 2	Rank: 1	InstructGPT	11	2	2	ChatGPT	4	6	5	GPT-4	0	6	9	Model	0	1	2	3	4	5	6	InstructGPT	0	22	12	28	4	1	23	ChatGPT	0	8	8	19	8	10	38	GPT-4	0	2	4	12	8	8	41	Model	0	1	2	3	4	5	6	InstructGPT	2	21	16	17	15	10	10	ChatGPT	1	9	16	12	15	15	22	GPT-4	1	4	6	10	10	23	22
Model	Rank: 3	Rank: 2	Rank: 1																																																																														
InstructGPT	11	2	2																																																																														
ChatGPT	4	6	5																																																																														
GPT-4	0	6	9																																																																														
Model	0	1	2	3	4	5	6																																																																										
InstructGPT	0	22	12	28	4	1	23																																																																										
ChatGPT	0	8	8	19	8	10	38																																																																										
GPT-4	0	2	4	12	8	8	41																																																																										
Model	0	1	2	3	4	5	6																																																																										
InstructGPT	2	21	16	17	15	10	10																																																																										
ChatGPT	1	9	16	12	15	15	22																																																																										
GPT-4	1	4	6	10	10	23	22																																																																										
<p><b>VOCAB: (w/definition)</b></p>	<p>Taxonomize – to arrange a set into a classification</p>																																																																																
<p><b>Cited references to follow up on</b></p>	<p>M. Lee, P. Liang, Q. Yang, “CoAuthor: Designing a human–AI collaborative writing dataset for exploring language model capabilities” in <i>Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems</i> (2022), pp. 1–19.</p>																																																																																
<p><b>Follow up Questions</b></p>	<p>Because they make so many arithmetic mistakes, could models like GPT-4 still be effective in teaching younger students?</p> <p>If these models were helpful only some of the time, is it up to the human to determine when to use them, or is it feasible to improve them so that they are always helpful?</p> <p>What does static evaluation of a LLM look like?</p>																																																																																

## Article #7 Notes:

<b>Source Title</b>	<b>Training Verifiers to Solve Math Word Problems</b>
<b>Source citation (APA Format)</b>	Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., Schulman, J. (2021). <b>Training Verifiers to Solve Math Word Problems.</b> arxiv.org/abs/2110.14168.
<b>Original URL</b>	<a href="https://doi.org/10.48550/arXiv.2110.14168">https://doi.org/10.48550/arXiv.2110.14168</a>
<b>Source type</b>	Article
<b>Keywords</b>	verifier
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>The researchers created GSM8K, a dataset of grade-school level math problems. It has natural language solutions rather than equations which would allow for better evaluation of large language models. They found that language models did not have high levels of accuracy on this dataset despite the problems lacking complexity.</p> <p>They hypothesized that using verifiers would increase this accuracy. They started by finetuning the models with various training set sizes. Unsurprisingly, models with more parameters and larger training sets had higher accuracy, when ran for 2 epochs. However, when allowed to run for more, allowing models to output 100 answers would cause their accuracy to eventually decrease (due to overfitting). However, this accuracy was still much higher than models that output only 1 answer. Next, they trained verifiers to output the probability that a model was correct. These models were trained on problems and solutions, but solutions could be labeled as correct even if the reasoning was wrong, as long as the final answer was right.</p> <p>They trained verifiers by finetuning a “generator” on the training data for 2 epochs, generating 100 solutions and labeling them as correct or incorrect, and then training a verifier on these solutions. They kept it at 2 epochs to maintain diversity in the data. They noted that it should be possible to combine the generator and verifier. They found that with high enough data sets, the models did</p>

	much better with verifiers for both 6B and 175B parameter models.
<b>Research Question/Problem/ Need</b>	Can using verifiers increase performance of large language models on math problems?
<b>Important Figures</b>	 <p>Verifiers drastically increase performance with enough data. A 6B parameter model with a verifier did better than the 175B parameter one without a verifier. Interestingly, verifiers decrease performance with small training sets</p>
<b>VOCAB: (w/definition)</b>	Epoch – a complete pass of a model through an entire dataset
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	<p>Why were verifiers worse for smaller training sets?</p> <p>How did they fine tune the models?</p> <p>Why didn't they combine the verifier and generator if they said it was possible to do it?</p>

## Article #8 Notes:

<b>Source Title</b>	The comparison of general tips for mathematical problem solving generated by generative AI with those generated by human teachers
<b>Source citation (APA Format)</b>	Jia, J., Wang, T., Zhang, Y., & Wang, G. (2024). The comparison of general tips for mathematical problem solving generated by generative AI with those generated by human teachers. <i>Asia Pacific Journal of Education</i> , 44(1), 8–28.  <a href="https://doi.org/10.1080/02188791.2023.2286920">https://doi.org/10.1080/02188791.2023.2286920</a>
<b>Original URL</b>	<a href="https://doi.org/10.1080/02188791.2023.2286920">https://doi.org/10.1080/02188791.2023.2286920</a>
<b>Source type</b>	Journal article
<b>Keywords</b>	Intelligent tutoring system, large language models, prompt engineering
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>They decided to use prompt engineering on ChatGPT to see if it could effectively generate tips for solving math problems. They used zero-shot, one-shot, and few-shot learning with and without CoT. However, they only used each one twice (one tip for a geometry problem and one for an algebra problem, for a total of 12 tips). They also had teachers generate tips for the same problems to compare them.</p> <p>Then, they developed a rubric to score these tips and had people score AI-generated tips and teacher-made tips. They created an online survey which got 121 responses, most of which were from people with teaching experience. Participants had to score the 12 tips given per problem and decide which ones were made by ChatGPT (6 were per problem).</p> <p>On average, teacher-made tips were scored better. They ran t-tests on both the ratings from the geometry problem and those from the algebra problem and had a p-value of less than 0.05 on both, suggesting that teacher-made tips were better</p>

than ChatGPT's tips. However, they also found that participants could not differentiate the source of the tips.

**Research Question/Problem/ Need** How do general tips for solving math problems generated by large language models differ from those generated by humans?

**Important Figures**

**Table 2.** Evaluation indicators and scores for general tips.

Primary Indicator	Secondary Indicator	Description	Score
Correctness	Correctness	The general tip correctly reflects the key steps or thinking process when solving the problem, without any logical or mathematical errors.	1
Expressiveness	Expressiveness	The general tip is clear, concise and coherent, without any grammatical errors, and uses the appropriate language and symbols.	1
Comprehensiveness	Comprehensiveness	If the content of general tip covers two aspects: method and knowledge that requires. It gets 1 point if it only covers one aspect, and 2 points if it covers both aspects.	2
Proceduralness	Clarify the problem	The general tip can help learners understand what the problem is asking.	1
	Make a plan	The general tip can help learners take into account all the given conditions and find the connection between the known conditions and the unknowns.	1
	Implement the plan	The general tip can help learners transform the original problem into an equivalent problem or a solved problem.	1
	Review	The general tip can summarize the knowledge points and methods used in this problem at the end, connecting common sense and theory.	1
Heuristicness	Heuristicness	The general tip can reflect the students' subjectivity and initiative, guide them to think and explore independently, rather than directly provide the answer or too simple hints. It can help learners apply the method to other problems and build their own knowledge network.	2
Total			10

The rubric that they had participants use to grade tips

Teacher-designed tips were rated higher on average than ChatGPT's tips

**VOCAB: (w/definition)** Zero-shot learning – allowing the model to generate responses without any examples given  
 One-shot learning – allowing the model to generate responses with one specific



	<p>example given</p> <p>Few-shot learning – using a few examples to allow the model to generate responses</p> <p>Chain of thought (CoT) - guiding the model to reason through to the desired output step-by-step (can be one-shot, few-shot, or even zero-shot)</p>
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	<p>How can ChatGPT's tips be improved to be as good as or even exceed teacher tips?</p> <p>Which prompt engineering method was the best at generating tips?</p> <p>How come people couldn't differentiate the source of the tips if they could differentiate their quality?</p>

## Article #9 Notes:

<b>Source Title</b>	Dropout Prediction in MOOCs: Using Deep Learning for Personalized Intervention
<b>Source citation (APA Format)</b>	Xing, W., & Du, D. (2019). Dropout Prediction in MOOCs: Using Deep Learning for Personalized Intervention. <i>Journal of Educational Computing Research</i> , 57(3), 547-570. <a href="https://doi.org/10.1177/0735633118757015">https://doi.org/10.1177/0735633118757015</a>
<b>Original URL</b>	<a href="https://doi.org/10.1177/0735633118757015">https://doi.org/10.1177/0735633118757015</a>
<b>Source type</b>	Journal article
<b>Keywords</b>	MOOC, dropout, deep learning,
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>MOOCs can have up to 90% attrition rates. This is usually just written off as a tradeoff for scale, but researchers wanted to look into solving it. They wanted to use deep learning because it would be impossible to manually look after all of these students.</p> <p>For their methodology they only investigated a single 8-week MOOC hosted on Canvas with 11 modules, 3617 students, 14 discussion forums, and 12 multiple choices. They tracked the various features listed in table 1.</p> <p>They started with 3 algorithms: K-nearest neighbors, support vector machines, and decision tree. Then they created a deep learning network (70% training data and 30% testing) which is different because it does automatic feature extraction and tuning. This allowed them to determine which students were most likely to drop out every week and thus plan personalized intervention.</p> <p>Deep learning model performed the best  KNN doesn't give a probability so it couldn't be used  Higher probability would indicate to teachers to give more intervention</p>
<b>Research Question/Problem/ Need</b>	How can students that are at risk of dropping out receive personalized intervention?

<b>Important Figures</b>	<p>Table 1. Features and Description.</p> <table border="1"> <thead> <tr> <th>Feature</th> <th>Description</th> </tr> </thead> <tbody> <tr> <td>Number of announcements</td> <td>Number of times students view the announcements</td> </tr> <tr> <td>Number of assignments</td> <td>Number of times students access the assignments</td> </tr> <tr> <td>Number of calendar</td> <td>Number of times students view the calendar</td> </tr> <tr> <td>Number of module pages</td> <td>Number of times students access the module pages</td> </tr> <tr> <td>Number of courses</td> <td>Number of times students access the courses</td> </tr> <tr> <td>Number of discussion</td> <td>Number of times students access the discussion forum</td> </tr> <tr> <td>Number of files</td> <td>Number of times students access the files</td> </tr> <tr> <td>Number of gradebooks</td> <td>Number of times students check the gradebooks</td> </tr> <tr> <td>Number quizzes</td> <td>Number of times students access the quizzes</td> </tr> <tr> <td>Number of submissions</td> <td>Number of times students submit assignments</td> </tr> <tr> <td>Number of Wikis</td> <td>Number of times students access the wikis pages</td> </tr> <tr> <td>Number of actives days</td> <td>Number of days student interacts with the course</td> </tr> <tr> <td>Dropout week</td> <td>The week when the student last visits the course. This is for algorithms to predict.</td> </tr> </tbody> </table> <p>The features they tracked in the study.</p>	Feature	Description	Number of announcements	Number of times students view the announcements	Number of assignments	Number of times students access the assignments	Number of calendar	Number of times students view the calendar	Number of module pages	Number of times students access the module pages	Number of courses	Number of times students access the courses	Number of discussion	Number of times students access the discussion forum	Number of files	Number of times students access the files	Number of gradebooks	Number of times students check the gradebooks	Number quizzes	Number of times students access the quizzes	Number of submissions	Number of times students submit assignments	Number of Wikis	Number of times students access the wikis pages	Number of actives days	Number of days student interacts with the course	Dropout week	The week when the student last visits the course. This is for algorithms to predict.
Feature	Description																												
Number of announcements	Number of times students view the announcements																												
Number of assignments	Number of times students access the assignments																												
Number of calendar	Number of times students view the calendar																												
Number of module pages	Number of times students access the module pages																												
Number of courses	Number of times students access the courses																												
Number of discussion	Number of times students access the discussion forum																												
Number of files	Number of times students access the files																												
Number of gradebooks	Number of times students check the gradebooks																												
Number quizzes	Number of times students access the quizzes																												
Number of submissions	Number of times students submit assignments																												
Number of Wikis	Number of times students access the wikis pages																												
Number of actives days	Number of days student interacts with the course																												
Dropout week	The week when the student last visits the course. This is for algorithms to predict.																												
<b>VOCAB: (w/definition)</b>	<p>MOOC – massive online open course – free online courses that anyone can enroll into</p> <p>Attrition – gradual dropping out</p> <p>K-nearest neighbors (KNN) – an algorithm that classifies data into categories based on a given number of dimensions</p> <p>Decision tree – consists of a root node (features), branches (rules for classification), and leaf nodes (classification)</p> <p>Support vector machine (SVM) - creates a plane in feature space to separate</p>																												
<b>Cited references to follow up on</b>																													
<b>Follow up Questions</b>	<p>Can they evaluate the effectiveness of personalized intervention?</p> <p>Are dropout rates affected by the content of the MOOC?</p> <p>How well would educators be able to understand and use these?</p>																												

## Article #10 Notes:

<b>Source Title</b>	Enhancing Math Learning with AI: ChatGPT's Impact on Number Base Conversion Comprehension
<b>Source citation (APA Format)</b>	Gadapa, S. P., Daud, S. B. M., Hui, B. T. C., & Raju, M. R. T. (2024). Enhancing Math Learning with AI: ChatGPT's Impact on Number Base Conversion Comprehension. <i>International Journal of Academic Research in Progressive Education and Development</i> , 13(3), 992–1008. <a href="http://dx.doi.org/10.6007/IJARPED/v13-i3/21642">http://dx.doi.org/10.6007/IJARPED/v13-i3/21642</a>
<b>Original URL</b>	<a href="http://dx.doi.org/10.6007/IJARPED/v13-i3/21642">http://dx.doi.org/10.6007/IJARPED/v13-i3/21642</a>
<b>Source type</b>	Journal article
<b>Keywords</b>	ChatGPT, Student Performance, Wilcoxon-singed Rank Test, Man-Whitney U Test, Number Base Conversions
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>First, they generated various levels of questions using ChatGPT. The experimental group of students had a ChatGPT assessment in between the pre- and post-assessments while the control group did not. Each group had 170 randomly selected students. They collected data on the students' answers and demographics.</p> <p>For statistical testing, they ran many tests (to determine which test should be used) on the data but finally a Wilcoxon signed-rank test. It determined that there was a significant difference for the experimental group (but not for the control group which was to be expected). They also found no significant difference between male and female students.</p>
<b>Research Question/Problem/Need</b>	How well can questions generated by ChatGPT impact students' skill in number base conversions?
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Shapiro-Wilk test – a test that determines if data follows a normal distribution or not
<b>Cited references to follow up on</b>	Luan, H., Geczy, P., Lai, H., Gobert, J., Yang, S. J. H., Ogata, H., Baltes, J., Guerra, R., Li, P., & Tsai, C.-C. (2020). Challenges and Future Directions of Big Data and Artificial Intelligence in Education. <i>Frontiers in Psychology</i> , 11. <a href="https://doi.org/10.3389/fpsyg.2020.580820">https://doi.org/10.3389/fpsyg.2020.580820</a>
<b>Follow up Questions</b>	Are ChatGPT questions any better or worse than regular ones made by teachers?

If not, how could these questions be improved?

Would ChatGPT's effectiveness be maintained for more complex questions?

## Patent #1 Notes:

<b>Source Title</b>	Methods and systems for secure data analysis and machine learning
<b>Source citation (APA Format)</b>	Carley, D. N., (2022). <i>Methods and systems for secure data analysis and machine learning</i> (U.S. Patent No. 20220067181A1). U.S. Patent and Trademark Office. <a href="https://patents.google.com/patent/US20220067181A1/en">https://patents.google.com/patent/US20220067181A1/en</a>
<b>Original URL</b>	<a href="https://patents.google.com/patent/US20220067181A1/en">https://patents.google.com/patent/US20220067181A1/en</a>
<b>Source type</b>	Patent
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	The goal was to use systems that would allow for security in machine learning development. To allow this, the system included keeping labeled data confidential so that a user or device would only have access to subsets of the data at any given time. Also, the parameters would be kept completely confidential from the users as well. Most importantly, a model would be encrypted after training so that it can be safely released for public use.
<b>Research Question/Problem/Need</b>	How can machine learning models be made more secure?
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Artifacts – outputs of a model during various stages
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	Is it possible to encrypt the data during the initial phases? Why is it important to encrypt the model? Are these protections foolproof and if not, how can they be improved?

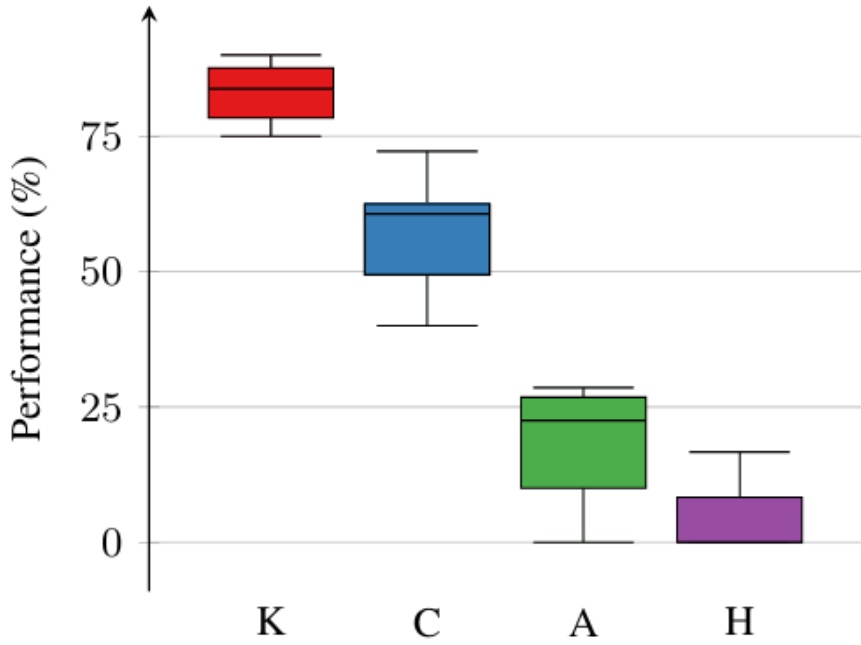
## Patent #2 Notes:

<b>Source Title</b>	Distributed labeling for supervised learning
<b>Source citation (APA Format)</b>	Bhowmick, A., Rogers, R. M., Vaishampayan, U. S., Vyrros, A. H., (2020). <i>Distributed labeling for supervised learning</i> (U.S. Patent No. 20200104705A1). U.S. Patent and Trademark Office. <a href="https://patents.google.com/patent/US20200104705A1/en">https://patents.google.com/patent/US20200104705A1/en</a>
<b>Original URL</b>	<a href="https://patents.google.com/patent/US20200104705A1/en">https://patents.google.com/patent/US20200104705A1/en</a>
<b>Source type</b>	Patent
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	They developed a technique to crowdsource labeling of data to be used for machine learning models while maintaining the privacy of the data. It starts with sending out unlabeled data to people's mobile devices that would give them back labels. Then, they would receive these labels encoded and determine the most frequent ones. They would add each element in the original data set with its most commonly proposed label to a training data set. Finally, it would train a model on this new data set.
<b>Research Question/Problem/Need</b>	How can labels be created for unlabeled data, when necessary, without sacrificing privacy and still being feasible on a large scale?
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Recurrent neural network – a neural network that feeds the output from the previous step as the input into the next one
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	How will they know that these are reliable sources of labels? What is the necessity of privatizing data? What are the limitations of crowdsourced labels?

## Article #11 Notes:

<b>Source Title</b>	Investigating the Effectiveness of ChatGPT in Mathematical Reasoning and Problem Solving: Evidence from the Vietnamese National High School Graduation Examination
<b>Source citation (APA Format)</b>	Dao, X.-Q., & Le, N.-B. (2023, October 10). <i>Investigating the Effectiveness of ChatGPT in Mathematical Reasoning and Problem Solving: Evidence from the Vietnamese National High School Graduation Examination</i> . ArXiv. <a href="https://arxiv.org/pdf/2306.06331">https://arxiv.org/pdf/2306.06331</a>
<b>Original URL</b>	<a href="https://doi.org/10.48550/arXiv.2306.06331">https://doi.org/10.48550/arXiv.2306.06331</a>
<b>Source type</b>	Arxiv article
<b>Keywords</b>	ChatGPT · large language model · natural language processing · Vietnamese high school graduation examination
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>Their goal was to determine ChatGPT's effectiveness at high school level math skills. Specifically, they evaluated its performance on the VNHSGE (Vietnamese High School Graduation Examination) dataset. The dataset was also separated by which year the question came from.</p> <p>VNHSGE consists of 250 multiple choice questions covering high school math topics like algebra, geometry, and calculus. The researchers divided the questions into four levels of difficulty (knowledge, comprehension, application, and high application).</p> <p>They gave ChatGPT each question as well as instructions on how to format the answer. As expected, they found that it performed worse on the harder questions. The exact percentages are shown in the figure.</p>
<b>Research Question/Problem/ Need</b>	How effective is ChatGPT at various math tests?



<b>Important Figures</b>	 <p>ChatGPT's performance on each difficulty of question.</p> <table border="1"> <caption>Approximate data from the box plot</caption> <thead> <tr> <th>Difficulty</th> <th>Min</th> <th>Q1</th> <th>Median</th> <th>Q3</th> <th>Max</th> </tr> </thead> <tbody> <tr> <td>K</td> <td>75</td> <td>78</td> <td>82</td> <td>85</td> <td>88</td> </tr> <tr> <td>C</td> <td>40</td> <td>50</td> <td>58</td> <td>62</td> <td>72</td> </tr> <tr> <td>A</td> <td>0</td> <td>10</td> <td>22</td> <td>28</td> <td>30</td> </tr> <tr> <td>H</td> <td>0</td> <td>5</td> <td>8</td> <td>10</td> <td>18</td> </tr> </tbody> </table>	Difficulty	Min	Q1	Median	Q3	Max	K	75	78	82	85	88	C	40	50	58	62	72	A	0	10	22	28	30	H	0	5	8	10	18
Difficulty	Min	Q1	Median	Q3	Max																										
K	75	78	82	85	88																										
C	40	50	58	62	72																										
A	0	10	22	28	30																										
H	0	5	8	10	18																										
<b>VOCAB: (w/definition)</b>																															
<b>Cited references to follow up on</b>																															
<b>Follow up Questions</b>	<p>How can it improve at the more difficult problems?</p> <p>How do these scores compare to a teacher, or someone expected to be knowledgeable at these skills?</p> <p>How did they separate the question difficulties?</p>																														

## Article #12 Notes

<b>Source Title</b>	Learning Relation-Enhanced Hierarchical Solver for Math Word Problems
<b>Source citation (APA Format)</b>	Lin, X., Huang, Z., Zhao, H., Chen, E., Liu, Q., Lian, D., Li, X., & Wang, H. (2023). Learning Relation-Enhanced Hierarchical Solver for Math Word Problems. <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 35(10), 13830-13844. <a href="http://dx.doi.org/10.1109/TNNLS.2023.3272114">http://dx.doi.org/10.1109/TNNLS.2023.3272114</a>
<b>Original URL</b>	<a href="http://dx.doi.org/10.1109/TNNLS.2023.3272114">http://dx.doi.org/10.1109/TNNLS.2023.3272114</a>
<b>Source type</b>	Journal Article
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>The main idea was that because humans solve math word problems much more efficiently than language models, they should take a more human approach to solving them. This was characterized by a few main factors. For example, humans gather meaning from a problem phrase by phrase, while machine learning models may go word by word by default. Also, humans can mentally group together similar problems, making them easier to solve, while AI models usually do not.</p> <p>They developed a hierarchical math solver (HMS) that derives semantics from each clause of a problem as they relate to the total problem. Then, they make a relation enhanced math solver (RHMS) that determines the similarity between math problems based on the structure. Both the HMS and RHMS proved to be effective when tested on large math datasets.</p>
<b>Research Question/Problem/Need</b>	How can math word problems be solved more efficiently by artificial intelligence?
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	GAT – graph attention network – a neural network that works with data structured as graphs
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	<p>How much more resource intensive is the RHMS than the HMS?</p> <p>Was any testing done to show that there was a statistically significant difference between RHMS and HMS performance?</p>

	How were they able to make the models perform well on many different data sets?
--	---

## Article #13 Notes

<b>Source Title</b>	Learning Fine-Grained Expressions to Solve Math Word Problems
<b>Source citation (APA Format)</b>	Huang, D., Shi, S., Lin, C., & Yin, J. (2017). Learning Fine-Grained Expressions to Solve Math Word Problems. <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , 805-814. <a href="https://doi.org/10.18653/v1/D17-1084">https://doi.org/10.18653/v1/D17-1084</a>
<b>Original URL</b>	<a href="#">Learning Fine-Grained Expressions to Solve Math Word Problems - ACL Anthology</a>
<b>Source type</b>	Journal Article
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>The main challenge the researchers tackled was deriving math concepts from natural language. Problems may use different words and contexts but still require the same math concept to be applied. First, they used their training data to create a few templates that the model could then use to solve any problem. That way, the model would be able to map the problem to a given template and then place the numbers in and solve. A template can just be thought of as a system of equations.</p> <p>When tested on a public dataset Dolphin18K, they got an accuracy of 28%. This may seem low, but at the time it was quite competitive as one state-of-the-art system at the time only reached an accuracy of 18%, for example.</p>
<b>Research Question/Problem/Need</b>	Systems that automatically solve math word problems have very low accuracy.
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Template – in the context of this article, a template was a system of equations with coefficients as variables that could be substituted by the numbers in the problem

<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	Would this concept still apply to problems with many steps? How complex were the problems in the dataset? Are these methods outdated compared to new methods and technology?

## Article #14 Notes

<b>Source Title</b>	The Ultimate Guide to Fine-Tuning LLMs from Basics to Breakthroughs: An Exhaustive Review of Technologies, Research, Best Practices, Applied Research Challenges and Opportunities  Venkatesh Balavadhani Parthasarathy, Ahtsham Zafar, Aafaq Khan, and Arsalan Shahid
<b>Source citation (APA Format)</b>	Parthasarathy, V. B., Zafar, A., Khan, A., & Shahid, A. (2024, October 30). <i>The Ultimate Guide to Fine-Tuning LLMs from Basics to Breakthroughs: An Exhaustive Review of Technologies, Research, Best Practices, Applied Research Challenges and Opportunities</i> . ArXiv. <a href="https://arxiv.org/pdf/2408.13296">https://arxiv.org/pdf/2408.13296</a>
<b>Original URL</b>	<a href="https://arxiv.org/pdf/2408.13296">2408.13296</a>
<b>Source type</b>	Arxiv article
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>There are 3 types of fine-tuning: supervised, unsupervised, and instructional. Supervised uses labeled data and is better when a specific task is in mind, while unsupervised is when unlabeled data is used to improve language capabilities in a domain. Instructional uses prompt engineering.</p> <p>Retrieval augmented generation (RAG) – incorporation of one’s own data into prompts for LLMs. It is much more cost effective because it doesn’t require all the hassles of fine-tuning and is good for question-and-answer use cases. RAG requires much less data.</p> <p>Parameter-efficient fine-tuning (PEFT) - less intensive than total fine-tuning because it involves adding adaptive layers to the neural network rather than editing every single one.</p> <p>Low rank adaptation (LoRA) - transforming the model into one with lower number of parameters. Allows for less resource-intensive tuning.</p>
<b>Research Question/Problem/Need</b>	The goal was to determine the best types of parameter-efficient fine-tuning for large language models, as well as creating a guide on which ones should be used.
<b>Important Figures</b>	

<b>VOCAB: (w/definition)</b>	Parameter-efficient fine-tuning – only tweaking a smaller number of parameters in a language model during fine-tuning rather than the entire model
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	Are different methods of fine-tuning better suited to different tasks? What kind of parameter-efficient fine-tuning is best for abstract question answering? Which methods of fine-tuning are the most resource efficient?

## Article #15 Notes

<b>Source Title</b>	Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks
<b>Source citation (APA Format)</b>	Lewis, P., Perez, E., Piktus, A., Petroni, P., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., & Kiela, D. (2021, April 12). <i>The Ultimate Guide to Fine-Tuning LLMs from Basics to Breakthroughs: An Exhaustive Review of Technologies, Research, Best Practices, Applied Research Challenges and Opportunities</i> . ArXiv. <a href="https://arxiv.org/pdf/2005.11401">https://arxiv.org/pdf/2005.11401</a>
<b>Original URL</b>	<a href="#">[2005.11401] Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks</a>
<b>Source type</b>	Arxiv article
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	Retrieval augmented generation (RAG) uses input to get some form of stored information (usually text documents) and then uses that to get an output. Used as fine tuning when you have your own data. They experimented with RAG on open-domain question answering, abstractive question answering, jeopardy question generation, and fact answering. RAG outperformed other models in all tasks tested.
<b>Research Question/Problem/Need</b>	The problem was that LLMs have lots of knowledge due to the amount of data they are trained on but usually cannot use the knowledge in meaningful ways besides regurgitating it.
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Retrieval augmented generation – uses input to get information from a given source, usually a text document, then uses that source to form an output
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	How did they score performance, especially in more complicated tasks like jeopardy question generation? Would retrieval augmented generation be suited towards non-NLP tasks?

	Is RAG-token model which can pull from different documents for each token better than RAG-sequence model which only uses one document?
--	--

## Article #16 Notes

<b>Source Title</b>	Sequence to Sequence Learning with Neural Networks
<b>Source citation (APA Format)</b>	Sutskever, I., Vinyals, O., & Le, Q. V. (2014, December 14). <i>Sequence to Sequence Learning with Neural Networks</i> . ArXiv. <a href="https://arxiv.org/pdf/1409.3215">https://arxiv.org/pdf/1409.3215</a>
<b>Original URL</b>	<a href="https://arxiv.org/abs/1409.3215">arXiv:1409.3215v3 [cs.CL] 14 Dec 2014</a>
<b>Source type</b>	Arxiv article
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	<p>Deep Neural Networks are very powerful but can't handle tasks that are sequential problems as their inputs and outputs are of unspecified dimensionality. Their plan was to test a model on English to French translation. They started with a Recurrent Neural Network which can map input sequences to output sequences provided that they are the same length. They used an LSTM to map an entire input sequence to a vector which would then be mapped to an output using another LSTM. This allowed recurrent neural networks to be used because the size of the vectors would be known.</p> <p>Their model received a BLEU score of 34.81 on the task of translation from English to French on the dataset used. BLEU is a common metric used to score machine learning models. For context, a different model that they were comparing to only had a score of 33.</p>
<b>Research Question/Problem/Need</b>	Deep Neural Networks are powerful but limited to labeled data; they can't be used for sequential tasks.
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Long Short-Term Memory – LSTM – a type of recurrent neural network that can hold information for a longer period of time
<b>Cited references to follow up on</b>	



<b>Follow up Questions</b>	<p>How was the LSTM able to translate long sentences despite memory constraints?</p> <p>How can the semantics of a sentence be captured in just a single vector?</p> <p>Would it be easy to use the model to translate French to English?</p>
----------------------------	---

## Article #17 Notes

<b>Source Title</b>	AI Chatbots as Math Algorithm Problem Solvers: A Critical Evaluation of Its Capabilities and Limitations
<b>Source citation (APA Format)</b>	Dahal, N., Luitel, B., C., Lamichhane, B., R., & Pant, B., P. (2023). AI Chatbots as Math Algorithm Problem Solvers: <i>Proceedings of the 28th Asian Technology Conference in Mathematics</i> , 429–438. <a href="https://www.researchgate.net/publication/375522509_AI_Chatbots_as_Math_Algorithm_Problem_Solvers_A_Critical_Evaluation_of_Its_Capabilities_and_Limitations">https://www.researchgate.net/publication/375522509_AI_Chatbots_as_Math_Algorithm_Problem_Solvers_A_Critical_Evaluation_of_Its_Capabilities_and_Limitations</a>
<b>Original URL</b>	<a href="#">(PDF) AI Chatbots as Math Algorithm Problem Solvers: A Critical Evaluation of Its Capabilities and Limitations</a>
<b>Source type</b>	Conference paper
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	Advanced Chatbot language models like ChatGPT and Bard can solve and explain basic math problems. They can also generate problems for educational purposes. They often use embed code as a response which can then be used to create an answer. However, a limitation of many of these models is that their solutions often come with drawbacks and are poorly explained. For example, even the tool WolframAlpha might not show all the steps required to solve a problem. They treat WolframAlpha like a chatbot in this article even though it differs heavily from models such as ChatGPT and Bard which are much more focused on natural language. Another drawback was that models tend to perform better on problems that are more theoretical examples and struggle with application of math concepts in the real world.
<b>Research Question/Problem/Need</b>	How well is AI currently equipped to handle skills necessary for aiding in math education?

<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	Multimodal – a model that can work with various types of data (text, image, audio, etc.)
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	<p>Why did they lump in WolframAlpha with chatbots?</p> <p>Why do natural language models struggle with geometry specifically?</p> <p>Why would they struggle with things that require real world knowledge if they are mostly trained on real world data rather than pure math data?</p>

## Article #18 Notes

<b>Source Title</b>	ChatGLM-Math: Improving Math Problem-Solving in Large Language Models with a Self-Critique Pipeline
<b>Source citation (APA Format)</b>	Xu, Y., Liu, X., Liu, X., Hou, Z., Li, Y., Zhang, X., Wang, Z., Zeng, A., Du, Z., Zhao, W., Tang, J., & Dong, Y. (2024, April 3). <i>ChatGLM-Math: Improving Math Problem-Solving in Large Language Models with a Self-Critique Pipeline</i> . ArXiv. <a href="https://arxiv.org/pdf/2404.02893">https://arxiv.org/pdf/2404.02893</a>

<b>Original URL</b>	<a href="#">[2404.02893] ChatGLM-Math: Improving Math Problem-Solving in Large Language Models with a Self-Critique Pipeline</a>																																																																																																																																																																																																													
<b>Source type</b>	ArXiv article																																																																																																																																																																																																													
<b>Keywords</b>																																																																																																																																																																																																														
<b>#Tags</b>																																																																																																																																																																																																														
<b>Summary of key points + notes (include methodology)</b>	<p>Machine learning models have already been used to generate feedback. This paper aimed to create a math-critique model to generate feedback on a large language model and thus improve its performance.</p> <p>One of the main problems with training LLMs to solve math problems is that the standard method of supervised fine-tuning may increase its math domain ability, but this would come at the cost of general language abilities.</p> <p>They used rejective fine tuning and direct performance optimization. Rejective fine tuning in this case was allowing the model to create responses, then scoring those responses with the math-critique model, then getting rid of low-scoring responses and fine-tuning with high-scoring ones. For direct performance optimization, they compared pairs of correct and incorrect answers to further tune the model. This was done after direct performance optimization.</p> <p>They tested these methods with many different models in both English and Chinese, including GPT-3.5-Turbo, Claude-2, and ChatGLM-3</p>																																																																																																																																																																																																													
<b>Research Question/Problem/ Need</b>	Large language models struggle with math problem solving because it differs from standard language usage.																																																																																																																																																																																																													
<b>Important Figures</b>	<table border="1"> <thead> <tr> <th rowspan="3">Models</th> <th rowspan="3">#params</th> <th colspan="3">Chinese</th> <th colspan="4">English</th> <th colspan="2">General</th> </tr> <tr> <th colspan="3">MathUserEval</th> <th rowspan="2">Ape210k</th> <th rowspan="2">Cmath</th> <th rowspan="2">GSM8k</th> <th rowspan="2">MATH</th> <th rowspan="2">Hungarian</th> <th rowspan="2">AlignBench</th> <th rowspan="2">MT-Bench</th> </tr> <tr> <th>Overall</th> <th>Elementary</th> <th>Advanced</th> <th>Language</th> </tr> </thead> <tbody> <tr> <td>GPT-4-1106-Preview [34]</td> <td>N/A</td> <td><b>5.73</b></td> <td><b>5.07</b></td> <td><b>6.81</b></td> <td><u>84.2</u></td> <td><b>89.3</b></td> <td><b>93.6</b></td> <td><b>53.6</b></td> <td><b>92</b></td> <td><u>8.29</u></td> <td><b>9.32</b></td> </tr> <tr> <td>GPT-4-0613 [34]</td> <td>N/A</td> <td>4.14</td> <td>3.34</td> <td>5.33</td> <td>83.6</td> <td>86.5</td> <td>91.4</td> <td>45.8</td> <td>68</td> <td>7.59</td> <td>9.18</td> </tr> <tr> <td>GPT-3.5-Turbo-0613 [34]</td> <td>N/A</td> <td>3.42</td> <td>3.04</td> <td>4.07</td> <td>70.4</td> <td>76.8</td> <td>78.2</td> <td>28.0</td> <td>41</td> <td>6.82</td> <td>8.36</td> </tr> <tr> <td>Claude-2 [1]</td> <td>N/A</td> <td>3.29</td> <td>2.63</td> <td>4.35</td> <td>72.8</td> <td>80.5</td> <td>88.0</td> <td>-</td> <td>55</td> <td>6.78</td> <td>8.06</td> </tr> <tr> <td>GLM-4</td> <td>N/A</td> <td><u>5.11</u></td> <td><u>4.86</u></td> <td><u>5.43</u></td> <td><b>93.5</b></td> <td><u>89.0</u></td> <td><u>91.8</u></td> <td><u>49.0</u></td> <td><u>75</u></td> <td><b>8.38</b></td> <td>8.62</td> </tr> <tr> <td>Skywork-13B-Math [54]</td> <td>13B</td> <td>2.66</td> <td>2.75</td> <td>2.54</td> <td>74.4</td> <td>77.3</td> <td>72.3</td> <td>17.0</td> <td>39</td> <td>5.58</td> <td>4.12</td> </tr> <tr> <td>InternLM2-Chat [43]</td> <td>20B</td> <td>3.25</td> <td>3.00</td> <td>3.68</td> <td>72.0</td> <td>80.7</td> <td>79.6</td> <td>34.8</td> <td>48</td> <td>7.68</td> <td>8.21</td> </tr> <tr> <td>Math-InternLM2 [43]</td> <td>20B</td> <td>3.17</td> <td>3.08</td> <td>3.37</td> <td>75.2</td> <td>78.5</td> <td>82.6</td> <td>37.7</td> <td>66</td> <td>6.53</td> <td>6.09</td> </tr> <tr> <td>Yi-Chat [56]</td> <td>34B</td> <td>2.64</td> <td>2.49</td> <td>2.87</td> <td>65.1</td> <td>77.7</td> <td>76.0</td> <td>15.9</td> <td>39</td> <td>6.18</td> <td>6.54</td> </tr> <tr> <td>DeepSeek-Chat [12]</td> <td>67B</td> <td>3.24</td> <td>2.76</td> <td>3.84</td> <td>76.7</td> <td>80.3</td> <td><b>84.1</b></td> <td>32.6</td> <td>58</td> <td>7.11</td> <td><b>8.35</b></td> </tr> <tr> <td>MetaMath (EN) [57]</td> <td>70B</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>82.3</td> <td>26.0</td> <td>35</td> <td>-</td> <td>4.28</td> </tr> <tr> <td>Qwen-Chat [3]</td> <td>72B</td> <td>3.87</td> <td><u>3.99</u></td> <td>3.67</td> <td>77.1</td> <td><b>88.1</b></td> <td>76.4</td> <td>31.8</td> <td>52</td> <td>7.29</td> <td>6.43</td> </tr> <tr> <td>ChatGLM3-32B-SFT-2312*</td> <td>32B</td> <td>3.25</td> <td>3.03</td> <td>3.60</td> <td>78.0</td> <td>79.8</td> <td>75.8</td> <td>29.0</td> <td>39</td> <td>7.37</td> <td>8.05</td> </tr> <tr> <td>+ RFT</td> <td>32B</td> <td><u>4.01</u></td> <td>3.86</td> <td><u>4.26</u></td> <td><u>87.0</u></td> <td>85.3</td> <td>82.4</td> <td><u>39.5</u></td> <td>58</td> <td>7.42</td> <td>8.03</td> </tr> <tr> <td>+ RFT, DPO</td> <td>32B</td> <td><b>4.23</b></td> <td><b>4.01</b></td> <td><b>4.59</b></td> <td><b>89.4</b></td> <td><u>85.6</u></td> <td><u>82.6</u></td> <td><b>40.6</b></td> <td><b>73</b></td> <td><b>7.80</b></td> <td>8.08</td> </tr> </tbody> </table> <p>* ChatGLM3-32B-SFT-2312 is a newer version of the ChatGLM series and not identical to the model discussed in [19], despite sharing the same model size.</p> <p>Scoring of inference of each model after training with the described methods. Each dataset has a different manner of scoring, so they are best compared within a column.</p>	Models	#params	Chinese			English				General		MathUserEval			Ape210k	Cmath	GSM8k	MATH	Hungarian	AlignBench	MT-Bench	Overall	Elementary	Advanced	Language	GPT-4-1106-Preview [34]	N/A	<b>5.73</b>	<b>5.07</b>	<b>6.81</b>	<u>84.2</u>	<b>89.3</b>	<b>93.6</b>	<b>53.6</b>	<b>92</b>	<u>8.29</u>	<b>9.32</b>	GPT-4-0613 [34]	N/A	4.14	3.34	5.33	83.6	86.5	91.4	45.8	68	7.59	9.18	GPT-3.5-Turbo-0613 [34]	N/A	3.42	3.04	4.07	70.4	76.8	78.2	28.0	41	6.82	8.36	Claude-2 [1]	N/A	3.29	2.63	4.35	72.8	80.5	88.0	-	55	6.78	8.06	GLM-4	N/A	<u>5.11</u>	<u>4.86</u>	<u>5.43</u>	<b>93.5</b>	<u>89.0</u>	<u>91.8</u>	<u>49.0</u>	<u>75</u>	<b>8.38</b>	8.62	Skywork-13B-Math [54]	13B	2.66	2.75	2.54	74.4	77.3	72.3	17.0	39	5.58	4.12	InternLM2-Chat [43]	20B	3.25	3.00	3.68	72.0	80.7	79.6	34.8	48	7.68	8.21	Math-InternLM2 [43]	20B	3.17	3.08	3.37	75.2	78.5	82.6	37.7	66	6.53	6.09	Yi-Chat [56]	34B	2.64	2.49	2.87	65.1	77.7	76.0	15.9	39	6.18	6.54	DeepSeek-Chat [12]	67B	3.24	2.76	3.84	76.7	80.3	<b>84.1</b>	32.6	58	7.11	<b>8.35</b>	MetaMath (EN) [57]	70B	-	-	-	-	-	82.3	26.0	35	-	4.28	Qwen-Chat [3]	72B	3.87	<u>3.99</u>	3.67	77.1	<b>88.1</b>	76.4	31.8	52	7.29	6.43	ChatGLM3-32B-SFT-2312*	32B	3.25	3.03	3.60	78.0	79.8	75.8	29.0	39	7.37	8.05	+ RFT	32B	<u>4.01</u>	3.86	<u>4.26</u>	<u>87.0</u>	85.3	82.4	<u>39.5</u>	58	7.42	8.03	+ RFT, DPO	32B	<b>4.23</b>	<b>4.01</b>	<b>4.59</b>	<b>89.4</b>	<u>85.6</u>	<u>82.6</u>	<b>40.6</b>	<b>73</b>	<b>7.80</b>	8.08
Models	#params			Chinese			English				General																																																																																																																																																																																																			
				MathUserEval			Ape210k	Cmath	GSM8k	MATH	Hungarian	AlignBench	MT-Bench																																																																																																																																																																																																	
		Overall	Elementary	Advanced	Language																																																																																																																																																																																																									
GPT-4-1106-Preview [34]	N/A	<b>5.73</b>	<b>5.07</b>	<b>6.81</b>	<u>84.2</u>	<b>89.3</b>	<b>93.6</b>	<b>53.6</b>	<b>92</b>	<u>8.29</u>	<b>9.32</b>																																																																																																																																																																																																			
GPT-4-0613 [34]	N/A	4.14	3.34	5.33	83.6	86.5	91.4	45.8	68	7.59	9.18																																																																																																																																																																																																			
GPT-3.5-Turbo-0613 [34]	N/A	3.42	3.04	4.07	70.4	76.8	78.2	28.0	41	6.82	8.36																																																																																																																																																																																																			
Claude-2 [1]	N/A	3.29	2.63	4.35	72.8	80.5	88.0	-	55	6.78	8.06																																																																																																																																																																																																			
GLM-4	N/A	<u>5.11</u>	<u>4.86</u>	<u>5.43</u>	<b>93.5</b>	<u>89.0</u>	<u>91.8</u>	<u>49.0</u>	<u>75</u>	<b>8.38</b>	8.62																																																																																																																																																																																																			
Skywork-13B-Math [54]	13B	2.66	2.75	2.54	74.4	77.3	72.3	17.0	39	5.58	4.12																																																																																																																																																																																																			
InternLM2-Chat [43]	20B	3.25	3.00	3.68	72.0	80.7	79.6	34.8	48	7.68	8.21																																																																																																																																																																																																			
Math-InternLM2 [43]	20B	3.17	3.08	3.37	75.2	78.5	82.6	37.7	66	6.53	6.09																																																																																																																																																																																																			
Yi-Chat [56]	34B	2.64	2.49	2.87	65.1	77.7	76.0	15.9	39	6.18	6.54																																																																																																																																																																																																			
DeepSeek-Chat [12]	67B	3.24	2.76	3.84	76.7	80.3	<b>84.1</b>	32.6	58	7.11	<b>8.35</b>																																																																																																																																																																																																			
MetaMath (EN) [57]	70B	-	-	-	-	-	82.3	26.0	35	-	4.28																																																																																																																																																																																																			
Qwen-Chat [3]	72B	3.87	<u>3.99</u>	3.67	77.1	<b>88.1</b>	76.4	31.8	52	7.29	6.43																																																																																																																																																																																																			
ChatGLM3-32B-SFT-2312*	32B	3.25	3.03	3.60	78.0	79.8	75.8	29.0	39	7.37	8.05																																																																																																																																																																																																			
+ RFT	32B	<u>4.01</u>	3.86	<u>4.26</u>	<u>87.0</u>	85.3	82.4	<u>39.5</u>	58	7.42	8.03																																																																																																																																																																																																			
+ RFT, DPO	32B	<b>4.23</b>	<b>4.01</b>	<b>4.59</b>	<b>89.4</b>	<u>85.6</u>	<u>82.6</u>	<b>40.6</b>	<b>73</b>	<b>7.80</b>	8.08																																																																																																																																																																																																			
<b>VOCAB: (w/definition)</b>	RLHF – reinforcement learning from human feedback – using humans to score responses that are then used to further tune the model																																																																																																																																																																																																													

<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	Do these results show that it is worth using the self-critique pipeline? Why weren't the numbers of parameters of some of the models known? Why did ChatGLM score the highest on most of the datasets?

## Article #19 Notes

<b>Source Title</b>	SciInstruct: a Self-Reflective Instruction Annotated Dataset for Training Scientific Language Models
<b>Source citation (APA Format)</b>	Zhang, D., Hu, Z., Zhoubian, Z., Du, Z., Yang, K., Wang, Z., Yue, Y., Dong, Y., & Tang, J. (2024, November 18). <i>SciInstruct: a Self-Reflective Instruction Annotated Dataset for Training Scientific Language Models</i> . ArXiv. <a href="https://arxiv.org/pdf/2401.07950">https://arxiv.org/pdf/2401.07950</a>
<b>Original URL</b>	<a href="https://arxiv.org/pdf/2401.07950">[2401.07950] SciInstruct: a Self-Reflective Instruction Annotated Dataset for Training Scientific Language Models</a>
<b>Source type</b>	ArXiv article
<b>Keywords</b>	

#Tags																																																																																																																																																																																																																																											
Summary of key points + notes (include methodology)	<p>They needed a large dataset full of varying scientific questions to train large language models. Chain of thought reasoning has been used to improve LLM performance on reasoning tasks, but for scientific data, chain of thought examples are not abundant.</p> <p>They started with 257,143 data points which were question-answer pairs. They used GPT-4 to generate intermediate steps by prompting it to give steps that would get to the answer. Then, they had other models label the accuracy to filter out inaccurate steps.</p> <p>To test Sci-Instruct, they chose ChatGLM3, Llama3-8B-Instruct, and Mistral-7B. They used the Sci-Instruct dataset to fine-tune each of these models so that they could test its accuracy. Then, they used various evaluation datasets to test their abilities. These models were able to outperform others, even if they had more parameters.</p>																																																																																																																																																																																																																																										
Research Question/Problem/ Need	LLMs are useful in science domains but are limited by a lack of scientific reasoning.																																																																																																																																																																																																																																										
Important Figures	<table border="1"> <thead> <tr> <th>Model</th> <th>CEval-Hard</th> <th>CEval-Sci</th> <th>MMLU-Sci</th> <th>SciEval</th> <th>SciBench</th> <th>GPQA_Diamond</th> <th>Avg. Sci</th> <th>Avg. (Sci+Math)</th> </tr> </thead> <tbody> <tr> <td colspan="9" style="text-align: center;">(API, parameter details unknown)</td> </tr> <tr> <td>GPT-4</td> <td>54.96</td> <td>60.55</td> <td>-</td> <td><b>73.93</b></td> <td><b>28.52</b></td> <td>39.70</td> <td>-</td> <td>-</td> </tr> <tr> <td>GPT-3.5-turbo</td> <td>41.37</td> <td>46.83</td> <td>-</td> <td>66.97</td> <td>12.17</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>Claude-v1.3</td> <td>39.14</td> <td>44.64</td> <td>-</td> <td>63.45</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td colspan="9" style="text-align: center;">(# parameter = 6B~7B)</td> </tr> <tr> <td>LLaMA-2-7B</td> <td>28.29<sup>†</sup></td> <td>30.00<sup>†</sup></td> <td>30.41</td> <td>28.37</td> <td>0.40</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>Galactica-6.7B</td> <td>11.84<sup>†</sup></td> <td>11.44<sup>†</sup></td> <td>30.68</td> <td>50.87</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>ChatGLM2-6B</td> <td>29.61<sup>†</sup></td> <td>45.71<sup>†</sup></td> <td>37.09<sup>†</sup></td> <td>53.02<sup>†</sup></td> <td>1.54<sup>†</sup></td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>ChatGLM2-6B-Base</td> <td>32.90<sup>†</sup></td> <td>40.95<sup>†</sup></td> <td>38.06<sup>†</sup></td> <td>50.38<sup>†</sup></td> <td>1.20<sup>†</sup></td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>ChatGLM3-6B</td> <td>36.84<sup>†</sup></td> <td>38.57<sup>†</sup></td> <td>41.78<sup>†</sup></td> <td>56.56<sup>†</sup></td> <td>2.40<sup>†</sup></td> <td>28.70</td> <td>34.14</td> <td>29.73</td> </tr> <tr> <td>ChatGLM3-6B-Base</td> <td><b>45.40<sup>†</sup></b></td> <td><b>54.29<sup>†</sup></b></td> <td>40.16<sup>†</sup></td> <td>61.69<sup>†</sup></td> <td>2.40<sup>†</sup></td> <td>24.75</td> <td><b>38.12</b></td> <td><b>40.34</b></td> </tr> <tr> <td>SciGLM (ChatGLM3-6B-Base)</td> <td><b>51.97</b></td> <td><b>60.00</b></td> <td><b>45.34</b></td> <td>62.09</td> <td><b>3.77</b></td> <td>25.25</td> <td><b>41.40</b></td> <td><b>45.32</b></td> </tr> <tr> <td>Llama3-8B-Instruct (zero-shot)</td> <td>26.32<sup>†</sup></td> <td>27.62<sup>†</sup></td> <td>26.90<sup>†</sup></td> <td><b>71.38<sup>†</sup></b></td> <td>1.03<sup>†</sup></td> <td>27.27<sup>†</sup></td> <td>30.09</td> <td>28.58</td> </tr> <tr> <td>Llama3-8B-Instruct (few-shot)</td> <td>25.66<sup>†</sup></td> <td>23.33<sup>†</sup></td> <td><b>52.67<sup>†</sup></b></td> <td>71.38<sup>†</sup></td> <td>3.60<sup>†</sup></td> <td>31.31<sup>†</sup></td> <td>34.66</td> <td>37.92</td> </tr> <tr> <td>+ SciInstruct</td> <td>32.24</td> <td>34.76</td> <td>40.86</td> <td><b>66.47</b></td> <td>3.60</td> <td>29.29</td> <td>34.54</td> <td>36.04</td> </tr> <tr> <td>Mistral-7B: MetaMATH (zero-shot)</td> <td>9.87<sup>†</sup></td> <td>8.57<sup>†</sup></td> <td>28.25<sup>†</sup></td> <td>63.61<sup>†</sup></td> <td>4.63<sup>†</sup></td> <td>27.78<sup>†</sup></td> <td>23.79</td> <td>25.57</td> </tr> <tr> <td>Mistral-7B: MetaMATH (few-shot)</td> <td>9.21<sup>†</sup></td> <td>19.52<sup>†</sup></td> <td>44.74<sup>†</sup></td> <td>63.61<sup>†</sup></td> <td>6.17<sup>†</sup></td> <td>29.29<sup>†</sup></td> <td>28.76</td> <td>33.92</td> </tr> <tr> <td>+ SciInstruct</td> <td>30.92</td> <td>38.10</td> <td>42.16</td> <td>64.16</td> <td>6.23</td> <td>27.27</td> <td>34.81</td> <td>37.91</td> </tr> <tr> <td colspan="9" style="text-align: center;">(# parameter = 12B~13B)</td> </tr> <tr> <td>LLaMA-2-13B</td> <td>19.74<sup>†</sup></td> <td>19.05<sup>†</sup></td> <td>35.85</td> <td>36.96</td> <td>1.37</td> <td>26.20</td> <td>22.59</td> <td>22.13</td> </tr> <tr> <td>Vicuna-13B</td> <td>-</td> <td>-</td> <td>32.13</td> <td>53.93</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td colspan="9" style="text-align: center;">(# parameter = 30B~32B)</td> </tr> <tr> <td>Galactica-30B</td> <td>-</td> <td>-</td> <td>35.53</td> <td>54.96</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td> </tr> <tr> <td>ChatGLM3-32B-Base</td> <td>53.95<sup>†</sup></td> <td>64.29<sup>†</sup></td> <td><b>50.30<sup>†</sup></b></td> <td>67.38<sup>†</sup></td> <td>4.29<sup>†</sup></td> <td>22.22</td> <td>43.74</td> <td>48.62</td> </tr> <tr> <td>SciGLM (ChatGLM3-32B-Base)</td> <td><b>56.58</b></td> <td><b>66.19</b></td> <td>49.38</td> <td><b>69.84</b></td> <td><b>5.15</b></td> <td><b>25.76</b></td> <td><b>45.48</b></td> <td><b>51.47</b></td> </tr> </tbody> </table> <p>Bolded numbers within a column means that model performed best compared to others with similar numbers of parameters</p>	Model	CEval-Hard	CEval-Sci	MMLU-Sci	SciEval	SciBench	GPQA_Diamond	Avg. Sci	Avg. (Sci+Math)	(API, parameter details unknown)									GPT-4	54.96	60.55	-	<b>73.93</b>	<b>28.52</b>	39.70	-	-	GPT-3.5-turbo	41.37	46.83	-	66.97	12.17	-	-	-	Claude-v1.3	39.14	44.64	-	63.45	-	-	-	-	(# parameter = 6B~7B)									LLaMA-2-7B	28.29 <sup>†</sup>	30.00 <sup>†</sup>	30.41	28.37	0.40	-	-	-	Galactica-6.7B	11.84 <sup>†</sup>	11.44 <sup>†</sup>	30.68	50.87	-	-	-	-	ChatGLM2-6B	29.61 <sup>†</sup>	45.71 <sup>†</sup>	37.09 <sup>†</sup>	53.02 <sup>†</sup>	1.54 <sup>†</sup>	-	-	-	ChatGLM2-6B-Base	32.90 <sup>†</sup>	40.95 <sup>†</sup>	38.06 <sup>†</sup>	50.38 <sup>†</sup>	1.20 <sup>†</sup>	-	-	-	ChatGLM3-6B	36.84 <sup>†</sup>	38.57 <sup>†</sup>	41.78 <sup>†</sup>	56.56 <sup>†</sup>	2.40 <sup>†</sup>	28.70	34.14	29.73	ChatGLM3-6B-Base	<b>45.40<sup>†</sup></b>	<b>54.29<sup>†</sup></b>	40.16 <sup>†</sup>	61.69 <sup>†</sup>	2.40 <sup>†</sup>	24.75	<b>38.12</b>	<b>40.34</b>	SciGLM (ChatGLM3-6B-Base)	<b>51.97</b>	<b>60.00</b>	<b>45.34</b>	62.09	<b>3.77</b>	25.25	<b>41.40</b>	<b>45.32</b>	Llama3-8B-Instruct (zero-shot)	26.32 <sup>†</sup>	27.62 <sup>†</sup>	26.90 <sup>†</sup>	<b>71.38<sup>†</sup></b>	1.03 <sup>†</sup>	27.27 <sup>†</sup>	30.09	28.58	Llama3-8B-Instruct (few-shot)	25.66 <sup>†</sup>	23.33 <sup>†</sup>	<b>52.67<sup>†</sup></b>	71.38 <sup>†</sup>	3.60 <sup>†</sup>	31.31 <sup>†</sup>	34.66	37.92	+ SciInstruct	32.24	34.76	40.86	<b>66.47</b>	3.60	29.29	34.54	36.04	Mistral-7B: MetaMATH (zero-shot)	9.87 <sup>†</sup>	8.57 <sup>†</sup>	28.25 <sup>†</sup>	63.61 <sup>†</sup>	4.63 <sup>†</sup>	27.78 <sup>†</sup>	23.79	25.57	Mistral-7B: MetaMATH (few-shot)	9.21 <sup>†</sup>	19.52 <sup>†</sup>	44.74 <sup>†</sup>	63.61 <sup>†</sup>	6.17 <sup>†</sup>	29.29 <sup>†</sup>	28.76	33.92	+ SciInstruct	30.92	38.10	42.16	64.16	6.23	27.27	34.81	37.91	(# parameter = 12B~13B)									LLaMA-2-13B	19.74 <sup>†</sup>	19.05 <sup>†</sup>	35.85	36.96	1.37	26.20	22.59	22.13	Vicuna-13B	-	-	32.13	53.93	-	-	-	-	(# parameter = 30B~32B)									Galactica-30B	-	-	35.53	54.96	-	-	-	-	ChatGLM3-32B-Base	53.95 <sup>†</sup>	64.29 <sup>†</sup>	<b>50.30<sup>†</sup></b>	67.38 <sup>†</sup>	4.29 <sup>†</sup>	22.22	43.74	48.62	SciGLM (ChatGLM3-32B-Base)	<b>56.58</b>	<b>66.19</b>	49.38	<b>69.84</b>	<b>5.15</b>	<b>25.76</b>	<b>45.48</b>	<b>51.47</b>
Model	CEval-Hard	CEval-Sci	MMLU-Sci	SciEval	SciBench	GPQA_Diamond	Avg. Sci	Avg. (Sci+Math)																																																																																																																																																																																																																																			
(API, parameter details unknown)																																																																																																																																																																																																																																											
GPT-4	54.96	60.55	-	<b>73.93</b>	<b>28.52</b>	39.70	-	-																																																																																																																																																																																																																																			
GPT-3.5-turbo	41.37	46.83	-	66.97	12.17	-	-	-																																																																																																																																																																																																																																			
Claude-v1.3	39.14	44.64	-	63.45	-	-	-	-																																																																																																																																																																																																																																			
(# parameter = 6B~7B)																																																																																																																																																																																																																																											
LLaMA-2-7B	28.29 <sup>†</sup>	30.00 <sup>†</sup>	30.41	28.37	0.40	-	-	-																																																																																																																																																																																																																																			
Galactica-6.7B	11.84 <sup>†</sup>	11.44 <sup>†</sup>	30.68	50.87	-	-	-	-																																																																																																																																																																																																																																			
ChatGLM2-6B	29.61 <sup>†</sup>	45.71 <sup>†</sup>	37.09 <sup>†</sup>	53.02 <sup>†</sup>	1.54 <sup>†</sup>	-	-	-																																																																																																																																																																																																																																			
ChatGLM2-6B-Base	32.90 <sup>†</sup>	40.95 <sup>†</sup>	38.06 <sup>†</sup>	50.38 <sup>†</sup>	1.20 <sup>†</sup>	-	-	-																																																																																																																																																																																																																																			
ChatGLM3-6B	36.84 <sup>†</sup>	38.57 <sup>†</sup>	41.78 <sup>†</sup>	56.56 <sup>†</sup>	2.40 <sup>†</sup>	28.70	34.14	29.73																																																																																																																																																																																																																																			
ChatGLM3-6B-Base	<b>45.40<sup>†</sup></b>	<b>54.29<sup>†</sup></b>	40.16 <sup>†</sup>	61.69 <sup>†</sup>	2.40 <sup>†</sup>	24.75	<b>38.12</b>	<b>40.34</b>																																																																																																																																																																																																																																			
SciGLM (ChatGLM3-6B-Base)	<b>51.97</b>	<b>60.00</b>	<b>45.34</b>	62.09	<b>3.77</b>	25.25	<b>41.40</b>	<b>45.32</b>																																																																																																																																																																																																																																			
Llama3-8B-Instruct (zero-shot)	26.32 <sup>†</sup>	27.62 <sup>†</sup>	26.90 <sup>†</sup>	<b>71.38<sup>†</sup></b>	1.03 <sup>†</sup>	27.27 <sup>†</sup>	30.09	28.58																																																																																																																																																																																																																																			
Llama3-8B-Instruct (few-shot)	25.66 <sup>†</sup>	23.33 <sup>†</sup>	<b>52.67<sup>†</sup></b>	71.38 <sup>†</sup>	3.60 <sup>†</sup>	31.31 <sup>†</sup>	34.66	37.92																																																																																																																																																																																																																																			
+ SciInstruct	32.24	34.76	40.86	<b>66.47</b>	3.60	29.29	34.54	36.04																																																																																																																																																																																																																																			
Mistral-7B: MetaMATH (zero-shot)	9.87 <sup>†</sup>	8.57 <sup>†</sup>	28.25 <sup>†</sup>	63.61 <sup>†</sup>	4.63 <sup>†</sup>	27.78 <sup>†</sup>	23.79	25.57																																																																																																																																																																																																																																			
Mistral-7B: MetaMATH (few-shot)	9.21 <sup>†</sup>	19.52 <sup>†</sup>	44.74 <sup>†</sup>	63.61 <sup>†</sup>	6.17 <sup>†</sup>	29.29 <sup>†</sup>	28.76	33.92																																																																																																																																																																																																																																			
+ SciInstruct	30.92	38.10	42.16	64.16	6.23	27.27	34.81	37.91																																																																																																																																																																																																																																			
(# parameter = 12B~13B)																																																																																																																																																																																																																																											
LLaMA-2-13B	19.74 <sup>†</sup>	19.05 <sup>†</sup>	35.85	36.96	1.37	26.20	22.59	22.13																																																																																																																																																																																																																																			
Vicuna-13B	-	-	32.13	53.93	-	-	-	-																																																																																																																																																																																																																																			
(# parameter = 30B~32B)																																																																																																																																																																																																																																											
Galactica-30B	-	-	35.53	54.96	-	-	-	-																																																																																																																																																																																																																																			
ChatGLM3-32B-Base	53.95 <sup>†</sup>	64.29 <sup>†</sup>	<b>50.30<sup>†</sup></b>	67.38 <sup>†</sup>	4.29 <sup>†</sup>	22.22	43.74	48.62																																																																																																																																																																																																																																			
SciGLM (ChatGLM3-32B-Base)	<b>56.58</b>	<b>66.19</b>	49.38	<b>69.84</b>	<b>5.15</b>	<b>25.76</b>	<b>45.48</b>	<b>51.47</b>																																																																																																																																																																																																																																			
VOCAB: (w/definition)	Lean – a popular syntax that is used to write math proofs and theorems with formality and logic																																																																																																																																																																																																																																										
Cited references to follow up on																																																																																																																																																																																																																																											
Follow up Questions	<p>How did they ensure accuracy of labeling of intermediate steps?</p> <p>What's the difference between ChatGLM3 and SciGLM(ChatGLM3) if they're both tuned off SciInstruct?</p> <p>How is small size considered a downside if it's just experimental?</p>																																																																																																																																																																																																																																										

## Article #20 Notes

<b>Source Title</b>	Dual Instruction Tuning with Large Language Models for Mathematical Reasoning
<b>Source citation (APA Format)</b>	Zhao, T., & Zhou, Y. (2024, March 27). <i>Dual Instruction Tuning with Large Language Models for Mathematical Reasoning</i> . ArXiv. <a href="https://doi.org/10.48550/arXiv.2403.18295">https://doi.org/10.48550/arXiv.2403.18295</a>
<b>Original URL</b>	<a href="https://arxiv.org/abs/2403.18295">[2403.18295] Dual Instruction Tuning with Large Language Models for Mathematical Reasoning</a>
<b>Source type</b>	ArXiv article
<b>Keywords</b>	
<b>#Tags</b>	
<b>Summary of key points + notes (include methodology)</b>	They used dual instruction tuning, meaning they would tune the model's generations in both directions of the sequence. They used the existing dataset MathInstruct, and applied what they called intermediate reasoning state prediction. This would involve masking certain parts of the data and then having the model fill in the gaps with its own generations. These generations would then be added to the dataset. This required the models to

	<p>use context from previous steps to get closer to an answer. They also applied instruction reconstruction, which involved doing the same thing as intermediate reasoning state prediction, but from backwards reasoning. They trained various models on this chain of thought data and found that it mainly resulted in improvements on more challenging datasets. Additionally, they calculated loss, or error in expected output on average.</p>
<b>Research Question/Problem/Need</b>	<p>Although Chain-of-thought is a powerful method to improve LLM reasoning skills, it still has limitations with the steps sometimes being missing, inaccurate, or unnecessary.</p>
<b>Important Figures</b>	
<b>VOCAB: (w/definition)</b>	<p>Ablation – removing components from a model one at a time to see what causes changes</p>
<b>Cited references to follow up on</b>	
<b>Follow up Questions</b>	<p>Wouldn't this result in inaccurate steps if the models made bad generations?  Why was it not as effective at improving performance on simple datasets?  Why did they test with a task that was not math related?</p>