

# Using Large Language Models to Automate Annotation and Part-of-Math Tagging of Math Equations

Authors: Ruocheng Shan & Abdou Youssef

The George Washington University

# Part-of-Math (POM) Tagging and Annotation

---

## **Definition of POM tagging and math annotation:**

- Identifying and labeling different components within math equations
  - Such as variables, operators, functions and constants
- Determining their roles and relationships within the equation

## **Applications of POM tagging:**

- Math UIs
- Generating metadata to enrich math-IR systems, and improve their performance
- Create Math datasets for training/finetuning/testing specialized math-AI models

# Research Overview

---

## Traditional approaches to math annotation/POM tagging are:

- Manual or semi-manual
- Relying on crafted rules
- Having limited datasets

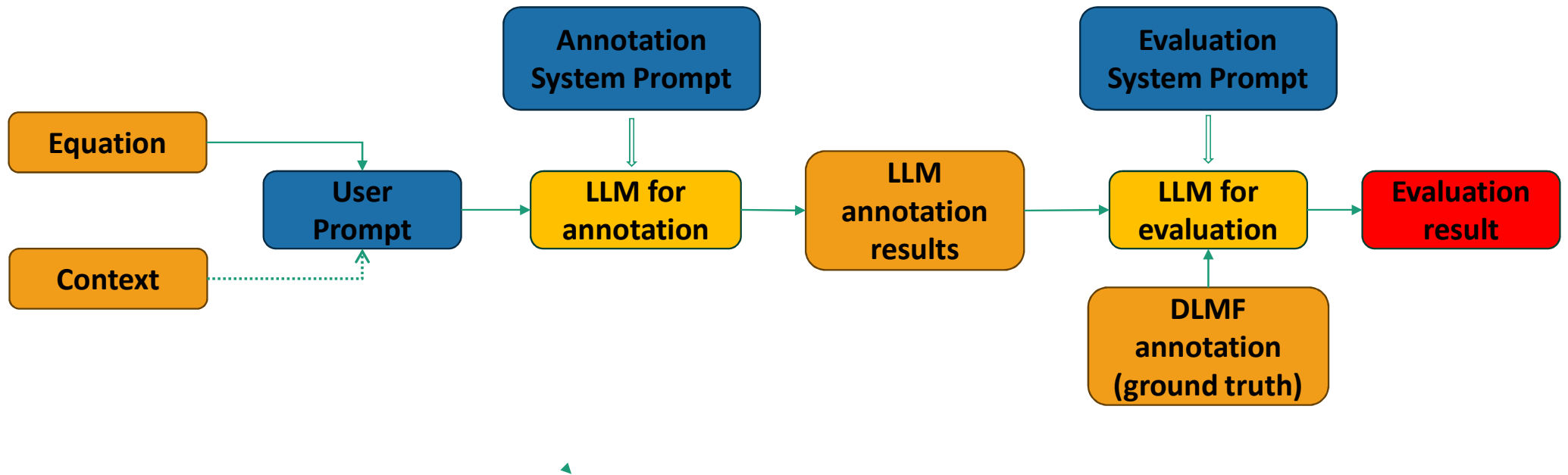
## Objectives of this research

1. Explore the effectiveness of Large Language Models (LLMs) in automating annotation and POM tagging of math equations
  - To reduce human involvement, and improve annotation accuracy
2. Investigate the impact of different levels of context on the accuracy of automated annotation
3. Explore the possibility of evaluating the annotation accuracy using LLMs

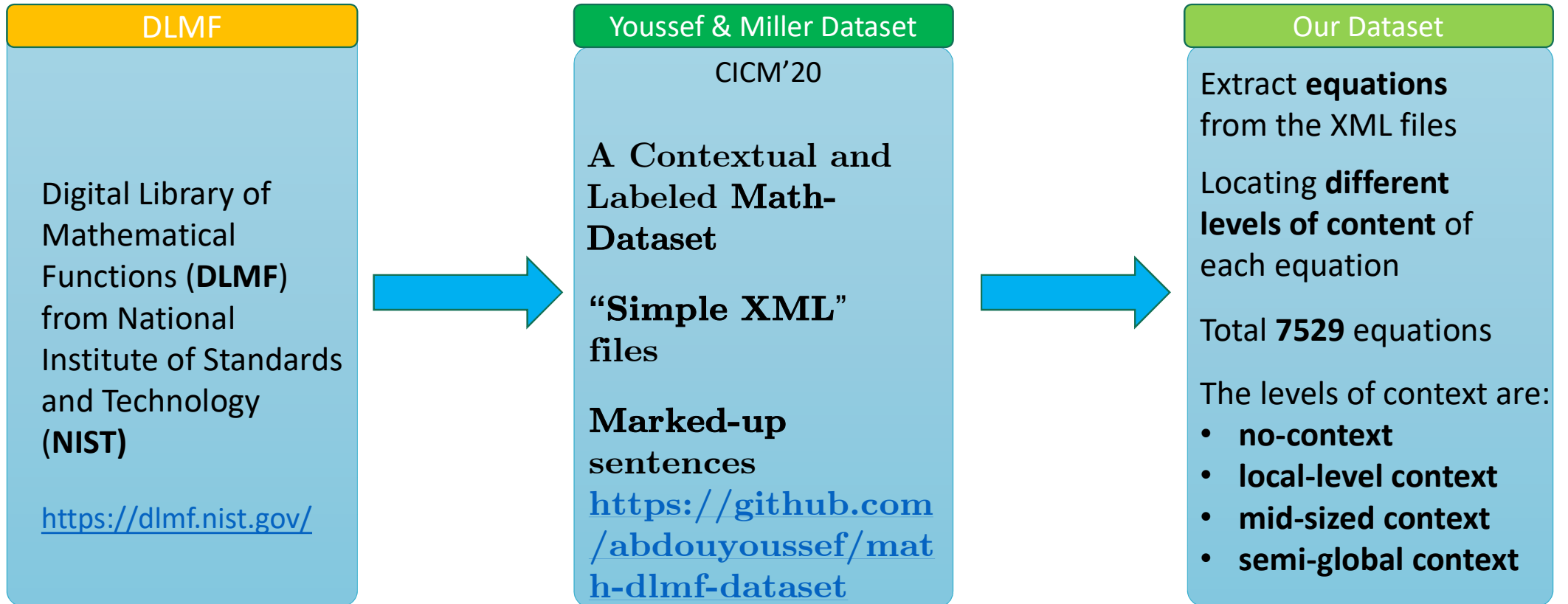
focus on GPT-3.5 Turbo

# Methodology Overview

---



# Dataset



# Example of Context Levels

Equation: 
$$f(z) = f(z_0) + \sum_{n=0}^{\infty} f_n(z - z_0)^{\mu+n}$$

Local (sentence) level context:

“Suppose that  $f(z) = f(z_0) + \sum_{n=0}^{\infty} f_n(z - z_0)^{\mu+n}$ , where  $\mu > 0$ ,  $f_0 \neq 0$ , and the series converges in a neighborhood of  $z_0$ .”

Mid-sized (paragraph) level context:

“Suppose that  $f(z) = f(z_0) + \sum_{n=0}^{\infty} f_n(z - z_0)^{\mu+n}$ , where  $\mu > 0$ ,  $f_0 \neq 0$ , and the series converges in a neighborhood of  $z_0$ . (For example, when  $\mu$  is an integer,  $f(z) - f(z_0)$  has a zero of order  $\mu$  at  $z_0$ .) Let  $w_0 = f(z_0)$ . Then (1.10.12) has a solution  $z = F(w)$ , where  $F(w) = z_0 + \sum_{n=1}^{\infty} F_n(w - w_0)^{n/\mu}$  in a neighborhood of  $w_0$ ,  $nF_n$  being the residue of  $1/(f(z) - f(z_0))^{n/\mu}$  at  $z = z_0$ .”

Semi-global context (doc-level notation) :

x,y    real variables.  
z        complex variable.  
j,k,l    integers.  
m,n     nonnegative integers, unless specified otherwise.  
⟨f,g⟩   inner, or scalar, product for real or complex vectors or functions.  
...      .....

# Prompt Engineering for Annotation

## System Prompt

### Background

You are a mathematician who is also experienced in computer science, especially natural language processing. You understand the essence of deep learning at the data format for training models. You are also familiar with the LaTeX representation of math equations.

### Task and instruction

Given an equation in LaTeX format, segment the equation into math terms/chunks, and give each term a math explanation as annotation. You will also be given some basic contextual information about the input mathematical equation. Use the information as a supporting reference to the equation.

### Notices and rules

Give your response in a key-value pair format where each key is a math term/chunk, and each value is an explanation of each term/chunk. Do not escape any math tokens. Return the result in a JSON format.

## User Prompt

The equation is  $[equation]$ ,  
and some contextual information about the equation:  $[contexts]$



**Background** is an introduction and identification of LLM's role



**Task and instruction** is the specific job needed to be done by the LLM.



**Notices and rules** is additional dos and don'ts for the LLM to follow.



**Template** to input equation and contexts

# LLM-Annotation of the All the DLMF Equations

---

- We got all the (7529) equations of the DLMF annotated by GPT
  - Using different levels of context
- Evaluated the GPT annotations in two ways
  - By hands
  - Using LLM



# Annotation examples (No-context)

<b>Equation</b>	$z^n - 1 = (z - 1)(z^{n-1} + z^{n-2} + \dots + z + 1) = 0$
<b>DLMF Annotation</b>	$z$ : variable, $n$ : non-negative variable
<b>LLM Annotation</b>	$z^n$ : The $n$ th power of $z$ , $1$ : The number 1, $(z - 1)$ : The difference between $z$ and 1, $(z^{n-1} + z^{n-2} + \dots + z + 1)$ : The sum of terms from $z$ to the power of $n-1$ to 1

<b>Equation</b>	$\text{Si}(z) = z \sum_{n=0}^{\infty} \left( j_n \left( \frac{1}{2}z \right) \right)^2$
<b>DLMF Annotation</b>	$\text{Si}(z)$ : sine integral, $j_n(z)$ : spherical Bessel function of the first kind, $z$ : complex variable, $n$ : non-negative integer
<b>LLM Annotation</b>	$\text{Si}(z)$ : Sine Integral function of $z$ , $z$ : Variable $z$ , $\sum_{n=0}^{\infty}$ : Summation from $n = 0$ to infinity, $j_n \left( \frac{1}{2}z \right)$ : Bessel function of the first kind of order $n$ with argument $1/2z$ , $^2$ : Squared

<b>Equation</b>	$F(s) = \mathcal{L}f(s) = \int_0^{\infty} e^{-st} f(t) dt$
<b>DLMF Annotation</b>	$\mathcal{L}f(s)$ : Laplace transform, $dx$ : differential of $x$ , $e$ : base of natural logarithm, $\int$ : integral
<b>LLM Annotation</b>	$F(s)$ : The Laplace transform of the function $f(t)$ with respect to the variable $s$ , $\mathcal{L}f(s)$ : Another way to represent the Laplace transform of the function $f(t)$ with respect to the variable $s$ , $\int_0^{\infty}$ : The integral from 0 to infinity, $e^{-st}$ : The exponential function with base $e$ raised to the power of $-st$ , $f(t)$ : The function $f$ evaluated at $t$ , representing the input function being transformed, $dt$ : The differential element indicating that the integration is with respect to the variable $t$

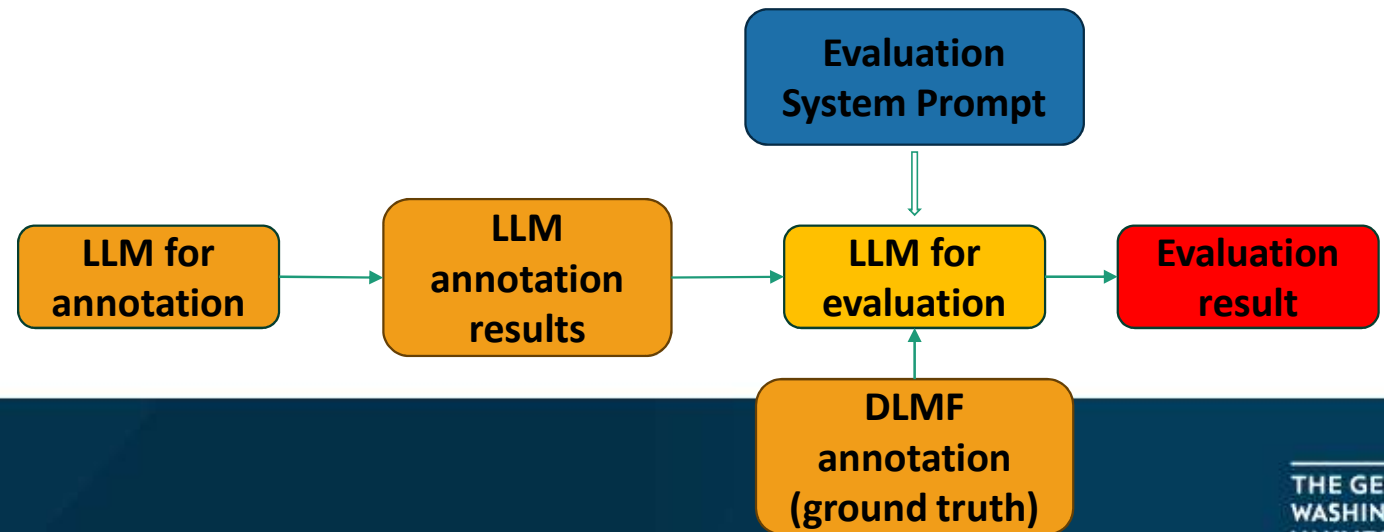
# Comparison between DLMF-Annotations and LLM-Annotations

---

- **The numbers of annotated terms are different**
  - The DLMF annotations are not comprehensive: some terms are labeled, some are not
  - While the LLM tends to label (nearly) all math tokens/symbols
- **The granularity of math chunks varies**
  - Some labels are given to a single math term like “z”, while sometimes given to “z^n”
- **The annotation wordings are different**
  
- **More importantly: Can we trust the LLM-annotations?     **Subject to investigation****
- **Need to evaluate the LLM accuracy (e.g., compare the LLM annotations to the DLMF’s)**
- **But, the comparison is nuanced and non-trivial (as seen above)**
- **Idea: Use LLMs to do the comparisons and categorize the relationship between the two annotations**
  
- **Question: Can we trust the LLMs to do the comparison, i.e., to evaluate the accuracy of the LLM-annotations?**

# LLM-based Evaluation

- Using a **separate LLM session** for evaluation (same LLM, different sessions)
- **Evaluation task** is viewed as **classification** of **annotation-pairs** (ground-truth annotation, LLM-generated ann.)
- First as **binary** classification of *consistent* vs. *non-consistent*
- Then as **multi-class** classification that is more refined, more informative, better aligned with reality



# LLM-based Evaluation: Binary Classification – Prompt

---

## System Prompt

### Background

You are a mathematician who is also experienced in computer science, especially natural language processing. You understand the essence of deep learning at the data format for training models. You are also familiar with the Latex representation of math equations.

### Task and instruction

There are two versions of math term annotations of a given equation, which also include the explanations of those math terms. Please determine whether the two versions are consistent or not.

If they are consistent, return "Yes". If not, return "No".

### Notices and rules

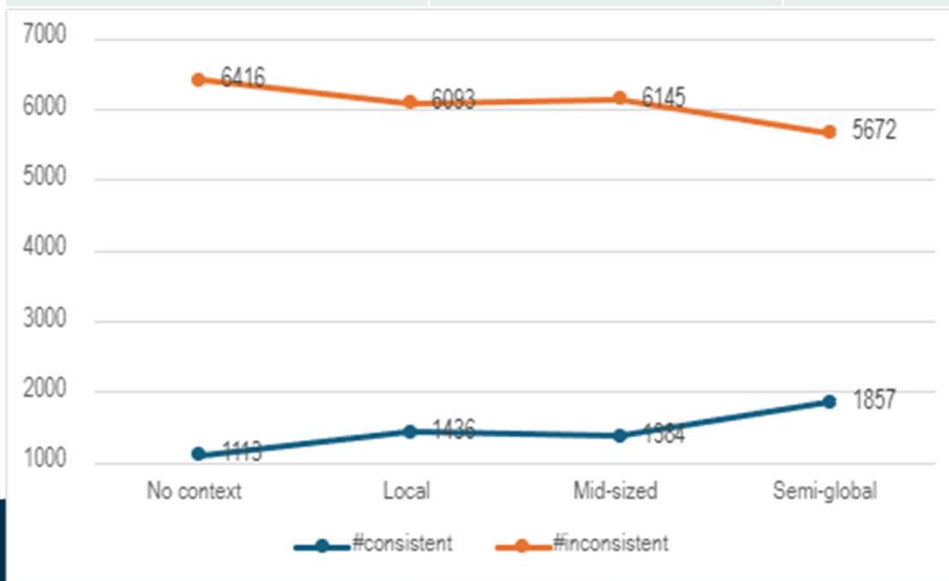
Note that the granularity of the math term segments could be different but still be consistent. Note also that the format of two versions could also be different, so pay attention to the content instead of the format.

## User Prompt

The first version of segmentation/annotation is [version1],  
the second version of segmentation/annotation is [version2]

# LLM-based Evaluation: Binary Classification Results

Context level	#consistent	#inconsistent	Consistency rate
No context	1113	6416	14.8%
Local context	1436	6093	19.1%
Mid-sized context	1384	6145	18.4%
Semi-global context	1857	5672	24.5%



- Overall consistency rate is low
- Consistency rate increases with context
  - Best when semi-global context was provided
- Binary classification is not enough to fully describe the relationship between two annotations



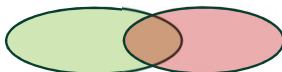

# LLM-based Evaluation: Multi-Class Classification (More Refined)

**System Prompt**

**Background**  
 You are a mathematician who is also experienced in computer science, especially natural language processing. You understand the essence of deep learning at the data format for training models. You are also familiar with the Latex representation of math equations.

**Task and instruction**  
 There are two versions of segmentations/annotations of a given equation, which also include the explanations of the math terms. Please analyze and classify the two versions into one of the following classes: consistent, contradictory, indeterminate, independent/incomparable, mixed, subset-of, and superset-of.

**Notices and rules**  
 Simply return the classification as the provided label. Do not include explanations of your answer.  
 Follow the detailed explanations or instructions of the classes below.  
**consistent:** Note that the granularity of segments could be different but still be consistent. Note that the format of two versions could also be different, pay attention to the content instead of the format.  
**contradictory:** One version does not agree or is the opposite of another.  
**indeterminate:** Do not have enough information to make the comparison possible.  
**independent/incomparable:** The two versions talk about two totally different things.  
**mixed:** The two versions are partially agreed or disagreed with one another.  
**subset-of:** The first version is a subset of the second version.  
**superset-of:** The first version is a superset of the second version.

- **Consistent:** the two versions (of annotation) agree with each other semantically (even if different words)  $A_{dlmf} \equiv A_{llm}$
- **Subset-of:** 1<sup>st</sup> version is a subset of 2<sup>nd</sup> version  $A_{dlmf} \subset A_{llm}$
- **Superset-of:** 1<sup>st</sup> version is a superset of 2<sup>nd</sup> version  $A_{dlmf} \supset A_{llm}$
- **Mixed:** the two versions agree partially and differ partially   $A_{dlmf} \neq A_{llm} , A_{dlmf} \cap A_{llm} \neq \emptyset$
- **Contradictory:** one version conflicts with the other version
- **Independent:** the 2 versions are about two totally different things   $A_{dlmf} \cap A_{llm} = \emptyset$
- **Indeterminate:** the 2 versions do not have enough information to make the comparison possible

# LLM-based Evaluation: Multi-Class Classification Results 1/2

Context Level \ Result class	No context	Local level	Mid-sized	Semi-global	
#consistent	1095	1111	1148	2028	← favorable outcomes
#subset-of	2454	2943	2938	3809	
#superset-of	1	3	2	9	← unfavorable outcomes
#mixed	3748	3279	3240	1551	
#contradictory	2	1	0	0	
#independent	229	192	201	133	
#indeterminate	0	0	0	0	

# LLM-based Evaluation: Multi-Class Classification Results 2/2

- Observably, the class of *consistent* and *subset-of* are both “good” results, so we group them into “*favorable outcomes*”.
- The other classes are grouped into “*unfavorable outcomes*”

Context level \ Result class	No context	Local level	Mid-sized	Semi-global
#favorable outcomes	3549	4054	4086	5837
#unfavorabl outcomes	3980	3475	3443	1692
favorable outcomes rate	47.1%	53.8%	54.3%	77.5%



# How Good are LLMs as Evaluators?

---

- An LLM-as-evaluator is fundamentally a classifier of pairs (of annotations): consistent, contradictory, mixed, ...
- Its classification accuracy may not be 100%
- To assess its classification (evaluation) accuracy, we need a labeled test set (ground truth)
  - Each instance is a pair of annotations, and a classification label of the relationship between the two annotations
- No such labeled dataset exists
  
- Creating such a dataset would be too time-consuming
  
- So we opted for a **statistical approximation** of the classification (evaluation) accuracy
  - We sampled 100 random instances, class-proportionally
  - Humanly labeled the pairs of annotations for all the 100 instances in the sample
  - Computed the classification (evaluation) accuracy based on that 100-instance sample

# LLM-based Evaluation – Human verification of binary-classification

LLM-assigned \ Human-assigned		Consistent	Inconsistent
		Consistent	Inconsistent
Consistent	49	<b>10</b>	
Inconsistent	<b>1</b>	40	
Accuracy	98%	80%	

Note: the context used here is the *semi-global context*

- Among 50 annotation-pairs labeled as “consistent” by LLM, 49 were found “consistent” by human
- Among 50 annotation-pairs labeled as “inconsistent” by LLM, 40 were found “inconsistent” by human

# LLM-based Evaluation – Human verification of multi-class classification

#Samples by Result class	Evaluation by LLM	Evaluation by Human
#consistent	18	18
#subset-of	30	30
#superset-of	3	0
#mixed	44	10
#contradictory	2	0
#independent	3	2
#indeterminate	0	0

Note: the context used here is the *semi-global context*

- LLM is very **accurate** when assigning “*favorable outcome*” labels (*consistent* and *subset-of*)
- LLM is **less accurate** when assigning “*unfavorable outcome*” labels.
- High True-Positive rate, somewhat lower True-Negative rate

# Future work

---

## **Comparative Analysis of Different LLMs:**

- Investigate the effectiveness of various LLMs (e.g., GPT-4, Gemini, Llama, ...) in math annotation and POM tagging
- Explore specialized LLMs designed for scientific and mathematical contexts to understand their impact on performance
- Explore the option of finetuning LLMs for better math annotation and POM tagging

## **Integration of More Context:**

- Incorporate external context, such as content dictionaries and knowledge graphs, to enhance LLM understanding and annotation of math equations

Questions?

---

**Thank you!**

# An example of wrongly labeled “mixed” by the LLM

Equation:

$$f(z) = \sum_{n=0}^{\infty} \frac{f^{(n)}(z_0)}{n!} (z - z_0)^n.$$

DLMF Annotations:

!: factorial,  
z: variable,  
n: nonnegative integer

LLM Annotations:

$f(z)$ : function  $f$  with input variable  $z$ ,  
 $\sum_{\{n=0\}}^{\infty}$  : sum from  $n = 0$  to infinity,  
 $\frac{f^{\{(n)\}}(z_0)}{n!}$ :  $n$ th derivative of  $f$  at  $z_0$  divided by  $n$  factorial,  
 $(z - z_0)^n$ :  $(z - z_0)$  raised to the power of  $n$

# An Example of the “mixed” Class of 2 Annotations

An example of “mixed” annotation-pairs

<b>Equation</b>	$v = \ln \left( \ln \left( \frac{1}{x} \right) \right) - 2 + \ln \pi$
<b>DLMF Annotation</b>	$v$ : expansion variable (locally), $\pi$ : the ratio of the circumference of a circle to its diameter, $\ln(z)$ : principal branch of logarithm function, $x$ : real variable
<b>LLM Annotation (no context)</b>	$v$ : velocity, $\ln$ : natural logarithm, $\frac{1}{x}$ : inverse of $x$ , $\pi$ : $\pi$

- The “ $v$ ” term has different annotation between DLMF and LLM
- Other terms agreed between the DLMF and the LLM

# An Example of the “mixed” Class of an 2 Annotations

An example of “mixed” annotation-pairs

Equation	$v = \ln \left( \ln \left( \frac{1}{x} \right) \right) - 2 + \ln \pi$
DLMF Annotation	<p><math>v</math>: expansion variable (locally),  <math>\pi</math>: the ratio of the circumference of a circle to its diameter,  <math>\ln(z)</math>: principal branch of logarithm function,  <math>x</math>: real variable</p>
LLM Annotation (no context)	<p><math>v</math>: velocity,  <math>\ln</math>: natural logarithm,  <math>\frac{1}{x}</math>: inverse of <math>x</math>,  <math>\pi</math>: <math>\pi</math></p>
LLM Annotation (semi-global context)	<p><math>v</math>: the variable representing the result of the equation,  <math>\ln</math>: Natural logarithm function,  <math>\frac{1}{x}</math>: Reciprocal of <math>x</math>,  <math>\ln \left( \frac{1}{x} \right)</math>: Natural logarithm of the reciprocal of <math>x</math>,  <math>\ln \left( \ln \left( \frac{1}{x} \right) \right)</math>: Natural logarithm of the natural logarithm of the reciprocal of <math>x</math>,  <math>-2</math>: Constant value of negative two,  <math>\ln \pi</math>: Natural logarithm of the mathematical constant <math>\pi</math></p>

- The “ $v$ ” term has different annotation between DLMF and LLM
- Other terms agreed between the DLMF and the LLM



# LLM-based Evaluation: Multi-Class Classification Results 1/3

Context Level \ Result class	No context	Local level	Mid-sized	Semi-global	
#consistent	1095	1111	1148	2028	Favorable outcome
#contradictory	2	1	0	0	Unfavorable outcome
#indeterminate	0	0	0	0	
#independent	229	192	201	133	Unfavorable outcome
#mixed	3748	3279	3240	1551	Unfavorable outcome
#subset-of	2454	2943	2938	3809	Favorable outcome
#superset-of	1	3	2	9	Unfavorable outcome

# Part-of-Math (POM) Tagging and Annotation

---

## Definition of POM tagging and math annotation:

- Identifying and labeling different components within math equations
  - Such as variables, operators, functions and constants
- Determining their roles and relationships within the equation

$$z^n - 1 = (z - 1)(z^{n-1} + z^{n-2} + \dots + z + 1) = 0$$



$z^n$ : The nth power of z,

1: The number 1,

$(z - 1)$ : The difference between z and 1,

$(z^{n-1} + z^{n-2} + \dots + z + 1)$ : The sum of terms from z to the power of n-1 to 1

## Applications of POM tagging:

- Math UIs
- Generating metadata to enrich math-IR systems, and improve their performance
- Create Math datasets for training/finetuning/testing specialized math-AI models