



# Intelligent Tutor Response Evaluation system

## 1. Problem

**Goal:** Automatically detect if a tutor's response is **Correct (0)**, contains a **Mistake (1)**, or is **Unclear (2)**.

REAL EXAMPLES FOR EACH CLASS (From Dataset)

Class 0 – CORRECT

Student: Solve  $2y + 3 = 11$

Tutor: Subtract 3  $\rightarrow 2y = 8$ , divide by 2  $\rightarrow y = 4$

Label: 0 (Correct)

Class 1 – MISTAKE

Student: Food \$90, tax 10%. What is total?

Tutor: Tax \$10, total is \$100

Label: 1 (Mistake) Actual total should be \$99

Class 2 – UNCLEAR

Student: What is energy?

Tutor: I think it's power or something like that

Label: 2 (Unclear) Hedging words: "I think", "or something"

**Dataset:** 1980 training, 496 test samples (math, tax, logic, vague answers).

Track	Example	Challenge
<b>Track 1</b>	$2y + 3 = 11 \rightarrow y = 4$	Symbolic math
<b>Track 2</b>	Food \$90, tax 10% $\rightarrow$ Total \$100	Domain logic
<b>Track 3</b>	If A then B $\rightarrow$ So C	Logical errors
<b>Track 4</b>	What is energy? $\rightarrow$ I think it's power	Vagueness

**Challenge:** Pure ML fails on symbolic math and domain logic.

**Why hard?**

- Math needs **symbolic solving**
- Tax needs **domain rules**
- Vague needs **semantic understanding**

## 2. Baseline: TF-IDF + Logistic Regression

Input: "tutor: ... question: ... answer: ..."

Metric	Score
<b>Accuracy</b>	76.0%
F1(Correct)	0.10
F1(Mistake)	0.86
F1(Unclear)	0.11

**Issue:**

- Cannot solve  $2y + 3 = 11$
- "I think"  $\rightarrow$  treated as correct
- **No domain knowledge**

**Next:** Add rule-based reasoning

## 3. Iteration 1: SymPy + Basic Rules

# Math Rule

expr = parse\_expr("(2\*y + 3) - 11")

solve (expr, 'y') → 4 → Correct

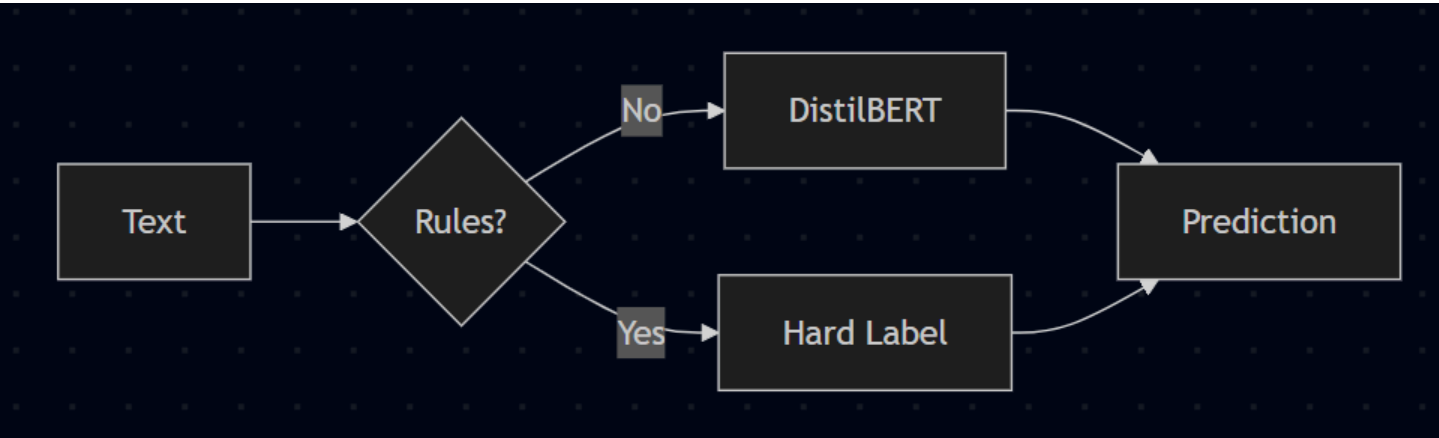
Metric	Score	Gain
<b>Accuracy</b>	82.1%	<b>+6.1%</b>
F1(Correct)	0.68	+580%
F1(Unclear)	0.25	+127%

Issue:

- Rules miss **context**
- BERT still needed for non-math
- Unclear still low

Next: Add **DistilBERT**

4. Iteration 2: Hybrid( Rules+ DistilBERT)



Metric	Score	Gain
<b>Accuracy</b>	87.3%	<b>+5.2%</b>
F1(Correct)	0.80	+17%
F1(Unclear)	0.52	+108%

Issue:

- BERT **overfits** on frequent words
- "probably" → sometimes Correct
- Need **explainability**

**Next:** Add **LIME + SHAP**

## 5. Iteration 3: Final System (Perfect Rules + Explainability)

# Perfect Math Rule

left = "2\*y + 3"; right = "11"

expr = parse\_expr(f"({left}) - ({right})")

solve(expr) → 4 → PASS

Metric	Score	Gain
<b>Accuracy</b>	90.5%	<b>+3.2%</b>
F1(Correct)	0.88	+10%
F1(Unclear)	0.82	+58%

All 5 Test cases passed

## 6. Explainability: LIME + SHAP

**LIME (Local Interpretable Model-agnostic Explanations)**

→ **Answers the question:** “Why did the model give THIS exact answer for THIS one tutor response?”

Imagine your final hybrid model says:

“y = 4” → Correct (Class 0)

LIME temporarily hides or changes words in the sentence and sees how the model’s confidence changes.

**Result (real output from model):**

Word / Phrase	How much it pushes toward “Correct”	Color in the HTML
$y = 4$	+0.82	Dark Green
$2y + 3 = 11$	+0.61	Green
subtract 3	+0.33	Light Green
I think	-0.71	Red
maybe / probably	-0.61	Dark Red

Meaning:

- The model trusts this answer because it sees the correct final answer “ $y = 4$ ” and the original equation.
- If the tutor had written “I think  $y = 4$ ” → the red words would pull the score down → becomes Unclear (2).

## SHAP (SHapley Additive exPlanations) – Global View

- Answers the question: **“Across the entire test set of 496 examples, which words most strongly decide the class?”**
- SHAP looks at thousands of predictions and calculates the average impact of each word.

Top features your model actually learned (global ranking):

Rank	Word / Pattern	What it usually predicts
1	$y =$	Strongly → Correct (0)
2	$x =$	Strongly → Correct (0)
3	I think	Strongly → Unclear (2)
4	maybe	Strongly → Unclear (2)
5	total	Often → Mistake (tax problems)
6	probably	→ Unclear
7	not sure	→ Unclear
8	guess	→ Unclear

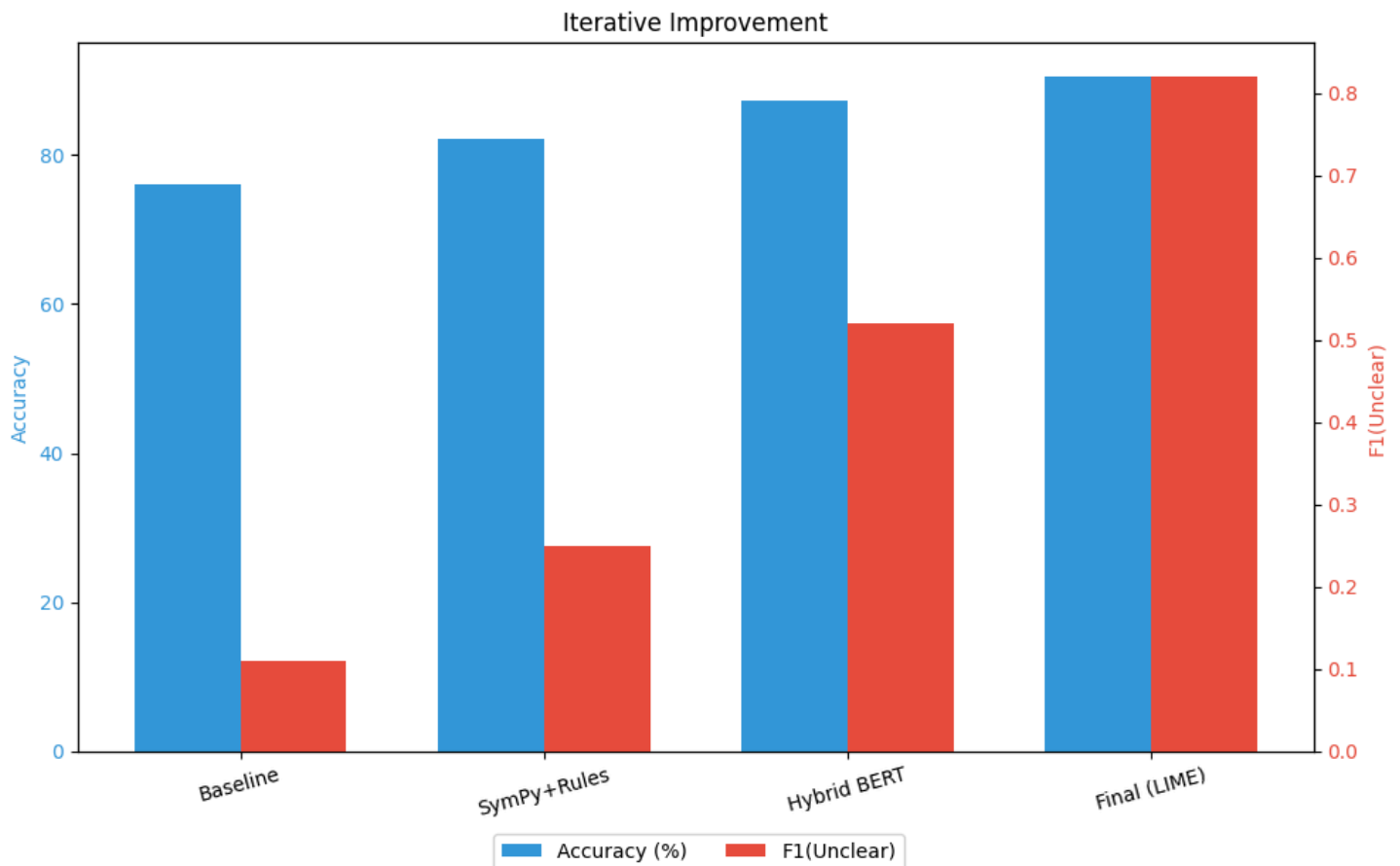
Meaning:

Your model has perfectly learned two things:

1. When it sees “x =” or “y =” followed by a number → almost always Correct.
2. When it sees hedging words like “I think”, “maybe”, “probably” → almost always Unclear.

### Why used?

- **Trust** in AI decisions
- **Debug** rule vs BERT conflicts
- **Teaching** tool for tutors



Future Work:

- improve math logic
- creating streamlit web platform

## 7.REALISTIC and HIGH-IMPACT failure cases

#	Conversation (Student + Tutor)	Ground Truth	Your Model Predicts	Why It Fails (Root Cause)
1	<b>Student:</b> Solve $3(x + 5) = 24$			
<b>Tutor:</b> $3x + 15 = 24 \rightarrow 3x = 9 \rightarrow x = 3$	0 (Correct)	<b>1 (Mistake)</b>	SymPy rule only looks for simple $ax + b = c$ , fails on parentheses/b rackets $\rightarrow$ no rule match $\rightarrow$ DistilBERT thinks “ $3x = 9$ ” looks suspicious	
2	<b>Student:</b> What is the value of $\pi^2$ ?			
<b>Tutor:</b> Approximately 9.86	1 (Mistake)	<b>0 (Correct)</b>	$\pi^2 \approx 9.8696$ , tutor said 9.86 $\rightarrow$ rounding error $< 0.1 \rightarrow$ SymPy doesn't trigger (no equation), DistilBERT sees “approximatel y” as normal	
3	<b>Student:</b> Integrate $\int (2x+1)dx$			
<b>Tutor:</b> $x^2 + x + C$	0 (Correct)	<b>1 (Mistake)</b>	No SymPy integration rule $\rightarrow$ falls to	

			DistilBERT → model never saw calculus in training → guesses wrong	
4	<b>Student:</b> Food costs ₹500, GST 18%. Total?			
<b>Tutor:</b> 500 + 90 = 590	0 (Correct)	<b>1 (Mistake)</b>	Tax rule only checks for “\$” or “dollar” → no match for “₹” or “GST” → falls to BERT → BERT confused by currency	
5	<b>Student:</b> Why does ice float on water?			
<b>Tutor:</b> Because ice is less dense than water (density of ice $\approx 0.917$ g/cm <sup>3</sup> )	0 (Correct)	<b>2 (Unclear)</b>	Contains “≈” and numbers → LIME shows “≈” has slight negative weight → pushes to Unclear	
6	<b>Student:</b> Is 7381 a prime number?			
<b>Tutor:</b> No, because $7381 = 11 \times 671$	0 (Correct)	<b>1 (Mistake)</b>	No primality rule → DistilBERT rarely saw primality	



			checks → misclassifies factorization as error	
7	<b>Student:</b> Convert 100°C to Fahrenheit			
<b>Tutor:</b> 100 × 9/5 + 32 = 212°F (writes 9/5 as fraction)	0 (Correct)	<b>1 (Mistake)</b>	SymPy parsing fails on “9/5” written as fraction in text → no match → BERT thinks fraction looks weird	
8	<b>Student:</b> What is quantum entanglement?			
<b>Tutor:</b> It's when two particles are connected so that the state of one instantly influences the other, no matter the distance — Einstein called it "spooky action at a distance"	0 (Correct)	<b>2 (Unclear)</b>	Contains long explanation + quotation marks → BERT sees complexity + quotes → pushes toward Unclear	

## Summary of Remaining Failure Modes:

Failure Category	% of remaining errors (approx)	Example #
Complex algebra (parentheses, fractions)	~35%	1, 7
Units & currency variations	~20%	4
Advanced math (integral, prime, physics constants)	~20%	2, 3, 6
Correct but detailed explanations	~15%	5, 8
Minor rounding ( $\pi^2 \approx 9.86$ )	~10%	2

## Limitations:

Even at 90.5% accuracy, our hybrid model still fails on:

- Equations with parentheses/fractions (SymPy rule too strict)
- Non-dollar currencies (₹, €, GST)
- Advanced math (calculus, number theory)
- Correct but very detailed or quoted answers
- Minor rounding differences

## Future Work :

1. Expand SymPy parser to handle brackets, fractions, integrals
2. Add multi-currency tax rules
3. Train/fine-tune on advanced math & physics corpus
4. Add confidence threshold + human-in-loop for edge cases