# BloomXplain: A Framework and Benchmark Dataset for Pedagogically Sound LLM-Generated Explanations Based on Bloom's Taxonomy

Maria-Eleni Zoumpoulidi[1], Eleni Batsi[1], Georgios Paraskevopoulos[1], Vassilis Katsouros[1], Alexandros Potamianos[2]

[1] Institute for Language and Speech Processing, Athena Research Center    [2] National Technical University of Athens

NEURAL INFORMATION PROCESSING SYSTEMS

## In a nutshell

We introduce a framework and a STEM-benchmark dataset for Pedagogically sound LLM-generated explanations based on Bloom's Taxonomy.
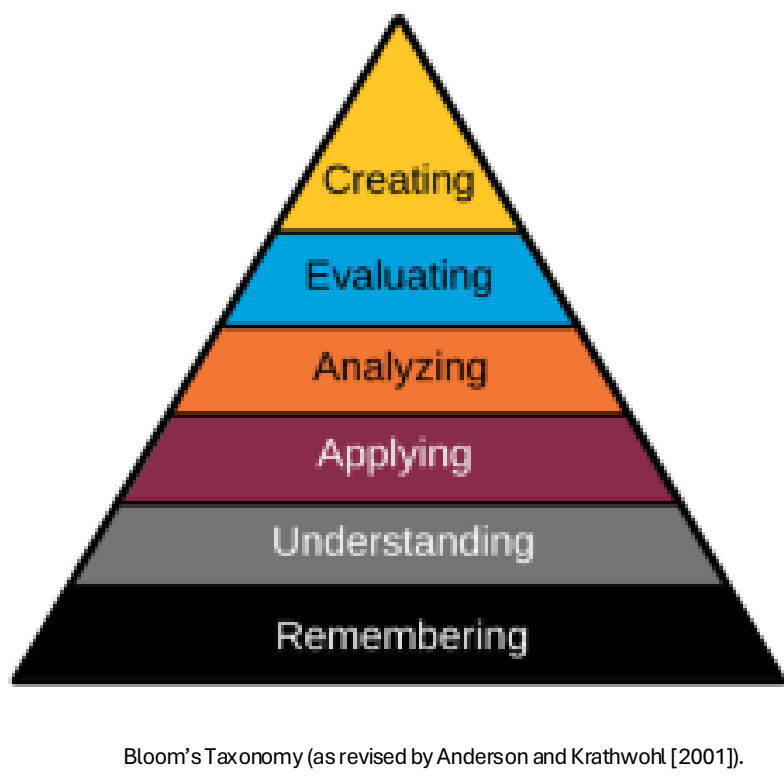
## Motivation

Why explanations?
- Useful for tutoring
- Better explanations ⇒ better reasoning

Why Bloom's Taxonomy?
- Structured framework
- Human-aligned, explainable results

## Contributions

- A STEM QA benchmark dataset annotated with Bloom levels
- Prompting strategies for eliciting Bloom-aligned explanations
- A multi-metric evaluation of LLMs across Correctness, Bloom Alignment, and Pedagogical Soundness

## Background: Bloom's Taxonomy



Bloom's Taxonomy (as revised by Anderson and Krathwohl [2001]).
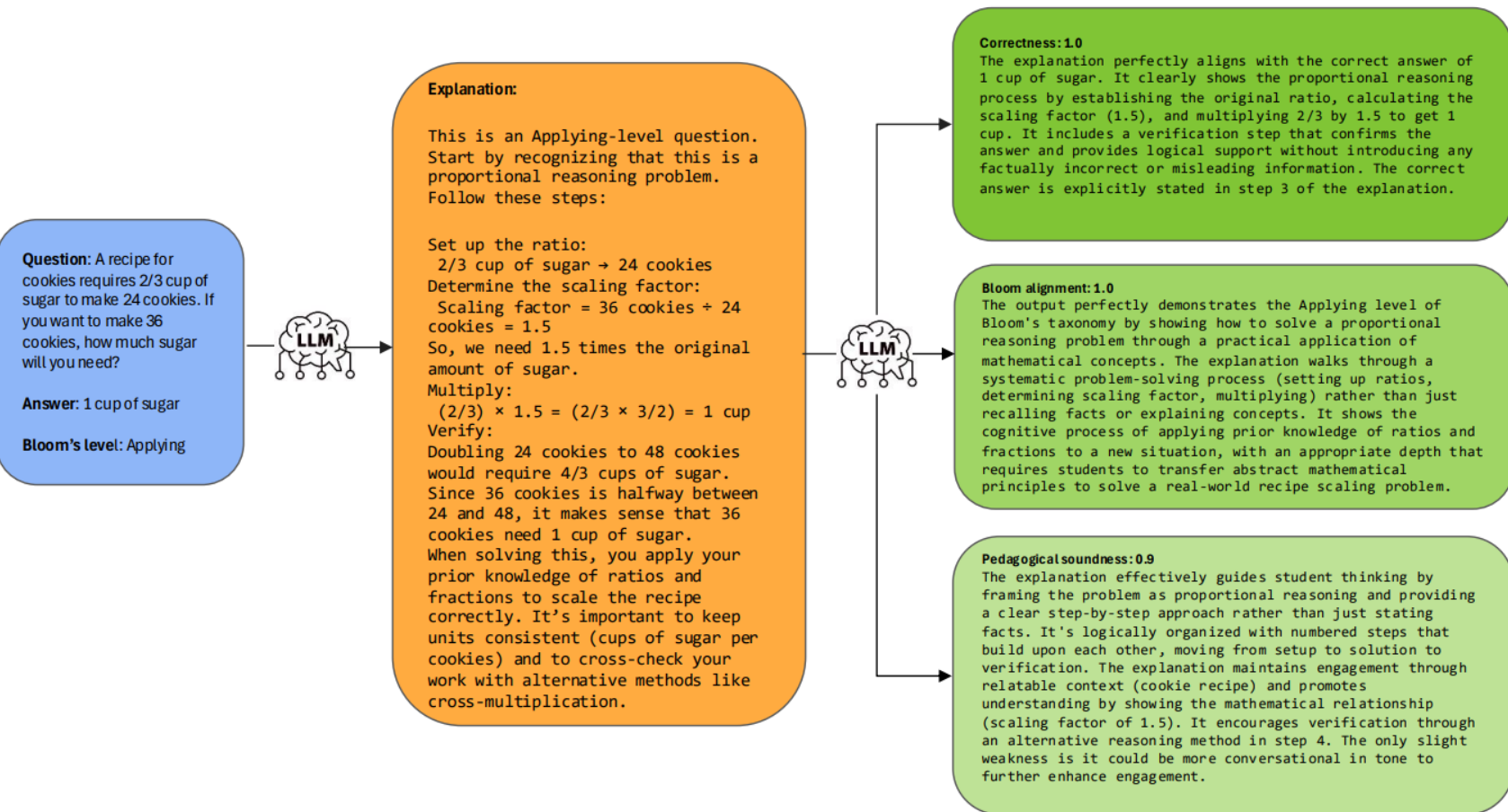
## Methodology

**Dataset:** 360 Bloom-aligned STEM QA pairs, spanning elementary → undergraduate, LLM-generated (Claude 3.7 sonnet) and human-validated

**Prompting strategies:**

| Prompting strategy | Input | Output |
|---|---|---|
| BAQ | Question, Answer, Bloom's level | Bloom-aligned explanation |
| AQ | Question, Answer | Inferred Bloom's level, Bloom-aligned explanation |
| Baseline | Question, Answer | Explanation |

**Evaluation:** LLM-as-a-Judge (Claude 3.7 sonnet) and human evaluation across three criteria: Correctness, Bloom Alignment, and Pedagogical Soundness
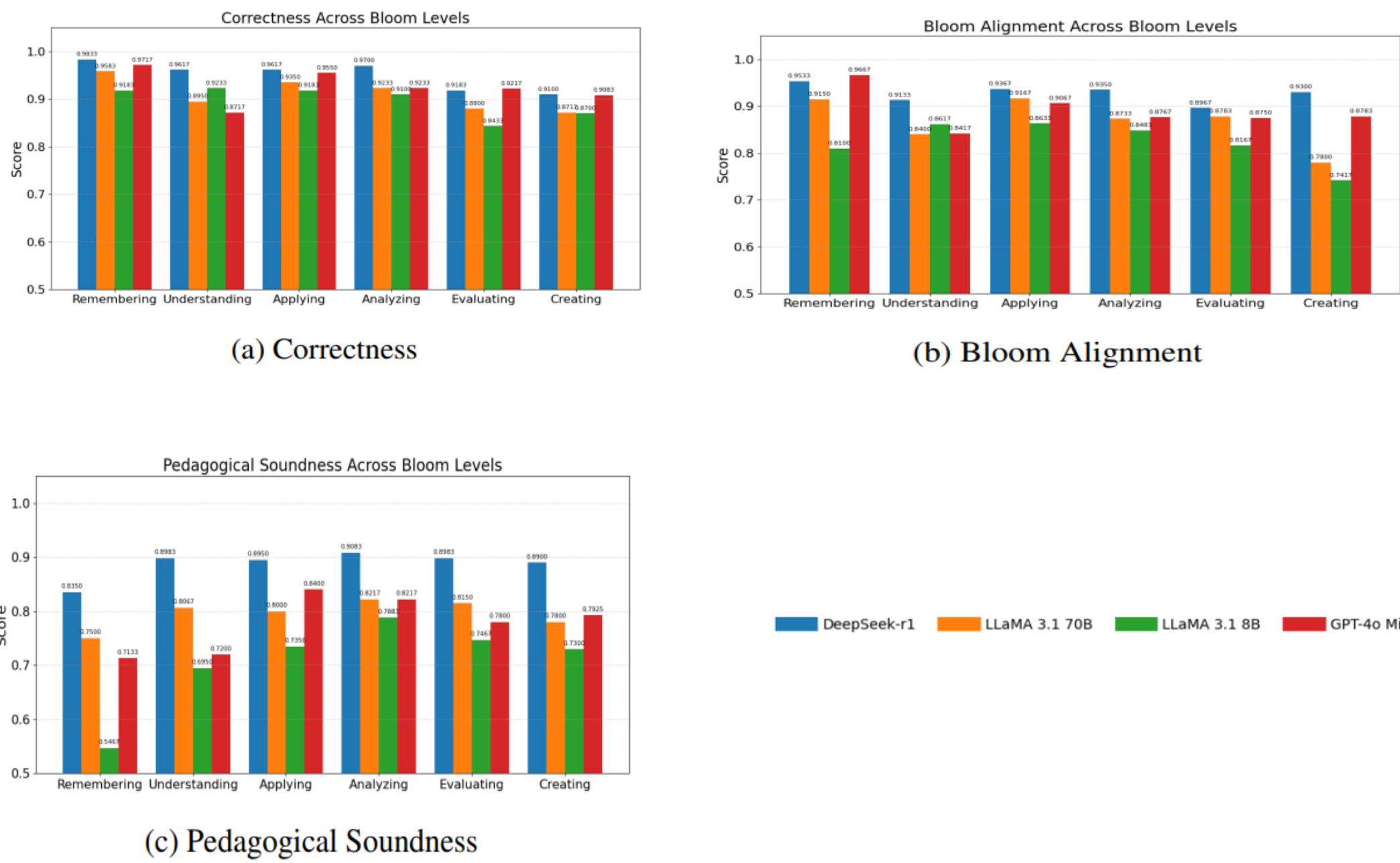
## Framework (BAQ)



## Main Results

| Model | Method | Correctness | Bloom Alignment | Pedagogical Soundness | Overall Score |
|---|---|---|---|---|---|
| deepseek-r1 | BAQ | 94.99 | **92.75** | 88.75 | **92.00** |
| | AQ | 93.75 | 87.00 | **89.83** | 90.00 |
| | Baseline | **96.16** | - | 76.16 | 85.99 |
| llama3.1 70b | BAQ | 91.16 | **86.83** | **79.49** | **85.66** |
| | AQ | 79.91 | 72.08 | 66.75 | 73.00 |
| | Baseline | **96.66** | - | 53.41 | 75.08 |
| llama3.1 8b | BAQ | 89.91 | **82.33** | **70.83** | **80.99** |
| | AQ | 93.41 | 78.41 | 63.66 | 78.41 |
| | Baseline | **95.75** | - | 49.50 | 72.66 |
| gpt-4o-mini | BAQ | 92.50 | **89.08** | **77.83** | **86.50** |
| | AQ | 89.91 | 80.08 | 72.08 | 80.58 |
| | Baseline | **93.99** | - | 48.58 | 71.33 |

🔍 BAQ outperforms other methods in pedagogical soundness and Bloom alignment while maintaining high correctness.

🔍 AQ, which infers Bloom levels, underperforms BAQ in both Bloom alignment and pedagogical soundness

🔍 While Baseline scores highest in correctness, its lack of structure leads to the lowest pedagogical score

BAQ's explicit Bloom-level guidance achieves the best balance of pedagogical depth and factual accuracy.

💡 Reasoning-optimized models achieve strong performance overall, while other models exhibit a much sharper pedagogy–correctness trade-off.

## BAQ's Performance per Bloom's level



(a) Correctness



(b) Bloom Alignment



(c) Pedagogical Soundness

- **Correctness:** Deepseek-r1 consistently leads across all Bloom levels, with GPT-4o-mini and LLaMA-3.1-70B close behind in most cases; performance drops for all models at higher cognitive levels (e.g., Evaluating).
- **Bloom Alignment:** Deepseek-r1 also achieves the strongest alignment, with GPT-4o-mini and LLaMA-3.1-70B performing similarly; LLaMA-3.1-8B generally lags except in Understanding tasks.
- **Pedagogical Soundness:** Deepseek-r1 again ranks highest, followed by GPT-4o-mini and LLaMA-3.1-70B; LLaMA-3.1-8B shows the weakest pedagogy, indicating smaller models struggle to provide instructional explanations.

## Comparison with CoT on widely used benchmarks (100 samples/task)

| Model | Benchmark | CoT | BAQ (ours) |
|---|---|---|---|
| Deepseek-r1 | BBH object counting (Remembering) | 96 | 100 |
| | BBH disambiguation qa (Understanding) | 60 | 78 |
| | GSM (Applying) | **99** | 99 |
| | BBH snarks (Analyzing) | 90 | 93 |
| gpt-4o-mini | BBH object counting (Remembering) | 88 | 95 |
| | BBH disambiguation qa (Understanding) | **74** | 68 |
| | GSM (Applying) | 94 | 98 |
| | BBH snarks (Analyzing) | 78 | 79 |

BAQ achieves competitive or superior performance compared to Chain-of-Thought (CoT) across Bloom's taxonomy levels, validating its efficacy in fostering robust reasoning