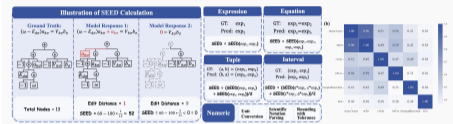


Motivation

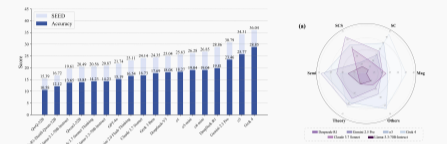
- Inadequate Benchmarks:** Current sets focus on K-12 or basic college levels, failing to test graduate-level research proficiency.
- Complexity of Physics:** Condensed Matter Physics (CMP) requires unique mastery of conceptual principles and precise symbolic math.
- Need for Nuance:** Brittle binary accuracy cannot capture near-miss reasoning in multi-step physical derivations

Scalable Expression Edit Distance



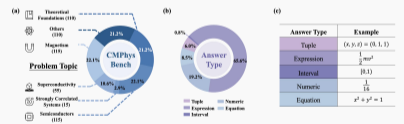
- Fine-grained, non-binary partial credit and interpretable error localization.
- The highest correlation with human expert judgment.

Results and Analysis



- Significant Capability Gap:** Across 18 LLMs, the top model achieves only a 36 SEED score and 29% accuracy.
- Asymmetric Strengths:** Models show highly uneven competence across different CMP subfields.
- Reasoning vs. General:** Longer CoT does not guarantee better performance without solid domain knowledge.

CMPhysBench



- Graduate-Level Depth:** A collection of 520 meticulously curated calculation problems authored by Ph.D. students and postdocs.
- Comprehensive Coverage:** Spans 6 core domains, including Magnetism, Superconductivity, and Strongly Correlated Systems.
- Diverse Answer Types:** Supports open-ended solutions in five formats: Tuple, Numeric, Expression, Equation, and Interval.
- Expert-Verified Pipeline:** Sourced from 17 classic textbooks with rigorous human review to ensure procedural correctness.

Example Problem

Consider the Anderson s-d exchange model with Hamiltonian $H = \sum_{k\sigma} E_{k\sigma} n_{k\sigma} + \sum_{d\sigma} E_{d\sigma} n_{d\sigma} + \sum_{k\sigma} n_{d\sigma} n_{k\sigma} + \sum_{k\sigma} V_{k\sigma} (c_{k\sigma}^\dagger a_{d\sigma} + a_{d\sigma}^\dagger c_{k\sigma})$ where $E_{k\sigma} = E_k + \sigma \mu_B h$, $E_{d\sigma} = E_d + \sigma \mu_B h$. Here, $\mu_B = \left(\frac{e\hbar}{2m_e}\right) \hbar$ is the Bohr magneton, with a Landé factor of $g_d = g_s = 2$ for both electrons and impurities. This is the non-degenerate orbital Anderson s-d mixing model. Derive the equation of motion for the s-d exchange model concerning the mixed Green's function $\langle\langle c_{k\sigma} | d_{d\sigma}^\dagger \rangle\rangle_{\omega}$, and $b_{k\sigma}$ symbolize $\langle\langle a_{k\sigma} | d_{d\sigma}^\dagger \rangle\rangle_{\omega}$.

Answer Type: Equation

Final Answer: $(\omega - E_{d\sigma}) b_{k\sigma} = V_{k\sigma} b_{k\sigma}$

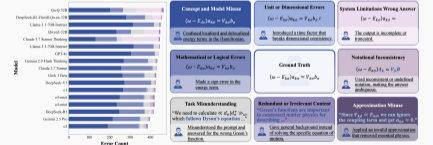
Topic: Strongly Correlated Systems

Scalable Expression Edit Distance

Ground Truth:	Model Response 1:	Model Response 2:
$(\omega - E_{d\sigma}) a_{2d} = V_{k0} b_{2d}$	$(\omega - E_{d\sigma}) n_{k\sigma} + a_{2d} = V_{k0} b_{2d}$	$0 = V_{k0} b_{2d}$
SEED Score 100	SEED Score 52	SEED Score 0

Expression Edit Distance

Ground Truth:	Response 1:	Response 2:
$(\omega - E_{d\sigma}) b_{k\sigma} = V_{k\sigma} b_{k\sigma}$	$(\omega - E_{d\sigma}) b_{k\sigma} = V_{k\sigma} b_{k\sigma}$	$0 = V_{k\sigma} b_{k\sigma}$
ACC Score 100	ACC Score 0	ACC Score 0



- Concept and Model Misuse (~50-60%):** The dominant error, exposing weak physics intuition.
- Mathematical or Logical Errors (20-30%):** Highlights persistent bottlenecks in symbolic derivation.