

# MeasHalu: Mitigation of Scientific Measurement Hallucinations for Large Language Models with Enhanced Reasoning

Anonymous ACL submission

## Abstract

The accurate extraction of scientific measurements from literature is a critical yet challenging task in AI4Science, enabling large-scale analysis and integration of quantitative research findings. However, Large Language Models (LLMs) frequently exhibit severe hallucinations, which significantly undermine the reliability of automated scientific document understanding systems. To address this problem, we propose **MEASHALU**, a novel framework for mitigating scientific measurement hallucinations through enhanced reasoning and targeted optimization. We first present a fine-grained taxonomy of measurement-specific hallucinations, categorizing errors across quantities, units, modifiers, and relations. Our approach incorporates a two-stage reasoning-aware fine-tuning strategy using augmented scientific data and process-based supervision. Furthermore, we introduce a progressive reward curriculum designed to penalize specific hallucination types, significantly improving extraction faithfulness. Experimental results demonstrate that MEASHALU substantially reduces hallucination rates and improves overall accuracy on the MeasEval benchmark. This work provides a targeted solution to a key bottleneck in automated scientific knowledge extraction, facilitating more trustworthy and scalable machine-assisted scientific literature analysis.

## 1 Introduction

The rapid expansion of scientific literature has created an unprecedented demand for reliable automatic extraction of quantitative knowledge, which lies at the core of modern AI4Science applications such as large-scale meta-analysis, knowledge base construction, and autonomous scientific discovery (Hanson et al., 2024; Chen et al., 2025). Central to this process is *scientific measurement extraction*—the task of identifying numerical quantities, their units, modifiers, and their relationships

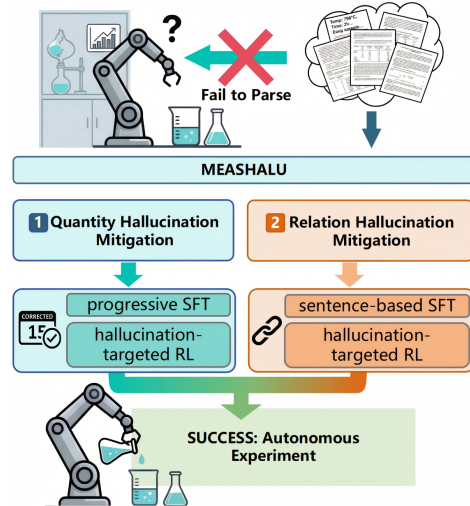


Figure 1: **Motivation of MeasHalu.** To rectify parsing failures, we propose a taxonomy-based approach to mitigate quantity and relation hallucinations.

to measured entities and properties. These quantitative statements form the evidential backbone of experimental sciences across disciplines ranging from materials science to biomedical research (Berrahou et al., 2013; Kononova et al., 2021). Although recent Large Language Models (LLMs) have demonstrated remarkable generalization abilities, they continue to perform unreliably on this task (Foppiano et al., 2024): even minor hallucinations in quantities or relations can invalidate entire experimental conclusions, severely limiting the trustworthiness of LLM-driven scientific understanding systems.

A key challenge underlying this failure is that *measurement hallucinations differ fundamentally from general textual hallucinations*. Unlike open-domain factual errors, measurement hallucinations exhibit fine-grained structural failures: models fabricate nonexistent values, misassociate quantities with wrong entities, overlook crucial qualifiers, or distort relations between scientific variables (Saier et al., 2024). Existing hallucination mitigation techniques, such as retrieval augmentation (Lewis et al.,

2020), generic instruction tuning, or conversational verification (Polak and Morgan, 2024), remain insufficient, as they are not designed to enforce the strict grounding and structural consistency required by scientific measurements. Yet, despite the importance of this problem, current research lacks both a systematic analysis of measurement-specific hallucination phenomena and targeted learning mechanisms for their mitigation. For instance, even state-of-the-art LLM-based extraction systems often compromise faithfulness by generating implicit information that is absent from the original text, such as inferring chemical formulas (Dagdelen et al., 2024).

In this work, we present MEASHALU, a reasoning-enhanced framework that explicitly models and suppresses scientific measurement hallucinations in LLMs. Our central insight is that hallucinations in this domain arise from two intertwined sources: (1) unreliable quantitative reasoning that corrupts individual quantities and units, and (2) fragile long-range relational reasoning that breaks the alignment between quantities, entities, and scientific properties. MEASHALU addresses these failure modes through a unified learning pipeline that combines reasoning-aware supervised fine-tuning with targeted reinforcement learning via structured reward shaping, thereby internalizing scientific grounding constraints directly into model parameters.

Concretely, MEASHALU introduces a fine-grained taxonomy of measurement hallucinations, and leverages this analysis to design a progressive optimization strategy: an initial supervised stage that standardizes quantitative reasoning and extraction structure, followed by Group Relative Policy Optimization (GRPO) with carefully constructed rewards that penalize fabrication, out-of-scope predictions, misclassification, and relational incompleteness. Our framework is developed on top of the MeasEval annotation schema (Harper et al., 2021a) and integrates external quantity validators, including CQE (Almasian et al., 2023a) and Quantum<sup>1</sup>, during training. Extensive experiments on the MeasEval benchmark and our newly constructed MeasEval-Ext dataset demonstrate that MEASHALU substantially reduces hallucination rates and consistently outperforms strong supervised baselines and proprietary LLMs. Furthermore, we show that MEASHALU functions as a

reliable external measurement extraction tool that significantly improves performance on downstream embodied scientific tasks, validating its practical utility for trustworthy AI4Science systems.

Our contributions are summarized as follows:

We provide the first fine-grained analysis of scientific *measurement hallucinations* in large language models, revealing their structural nature and identifying two fundamental sources of failure: unreliable quantitative reasoning and fragile relational grounding.

We propose MEASHALU, a unified reasoning-enhanced learning framework that systematically suppresses measurement hallucinations by integrating reasoning-aware supervised fine-tuning with targeted reinforcement learning via structured reward shaping.

We construct a new out-of-distribution evaluation benchmark, MEASEVAL-EXT, and demonstrate through extensive experiments that MEASHALU substantially reduces hallucination rates and consistently outperforms strong supervised baselines and proprietary LLMs on scientific measurement extraction.

We further show that MEASHALU serves as a reliable external measurement extraction tool that significantly improves performance on downstream embodied scientific tasks, validating its practical utility for trustworthy AI4Science systems.

## 2 Related Work

### 2.1 Hallucinations in Large Language Models

Hallucination, where language models generate ungrounded or factually incorrect content, has been extensively studied in general-purpose LLMs (Huang et al., 2025). Most prior work focuses on semantic and factual hallucinations in open-ended generation (Ji et al., 2023), with typical taxonomies including fabrication, inconsistency, and logical errors (Li et al., 2025). However, these taxonomies are largely developed for free-form text generation and do not capture the structural requirements of measurement extraction, where numerical faithfulness, unit consistency, and entity-quantity relational grounding are essential. We address this gap by proposing a fine-grained taxonomy of *measurement-specific hallucinations* and designing mitigation mechanisms tailored to these failure modes.

<sup>1</sup><https://github.com/nielstron/quantulum3>

163	<b>2.2 General Information Extraction vs.</b>			212
164	<b>Measurement Extraction</b>			213
165	Information extraction (IE) and named entity recog-			214
166	nition (NER) are foundational NLP tasks (Nadeau			215
167	and Sekine, 2007). While early systems relied on			216
168	rule-based and feature-engineered pipelines, mod-			217
169	ern approaches increasingly leverage neural archi-			218
170	tectures and pre-trained language models. Never-			219
171	theless, <i>scientific measurement extraction</i> poses ad-			220
172	ditional constraints beyond conventional IE: mod-			221
173	els must accurately capture numerical values, units,			222
174	and modifiers, and preserve their structured rela-			223
175	tions to measured entities and properties under			224
176	strict grounding. These constraints make the task			225
177	particularly sensitive to hallucinations and motivate			
178	learning objectives that explicitly penalize fabrica-			
179	tion, mis-scoping, and relational incompleteness.			
180	<b>2.3 Scientific Measurement Extraction and</b>			226
181	<b>Benchmarks</b>			
182	Scientific information extraction has been ad-			227
183	vanced by datasets such as SCIERC (Luan et al.,			228
184	2018) and MEASEVAL (Harper et al., 2021b).			229
185	Among them, MEASEVAL provides the most fine-			230
186	grained annotation schema for scientific measure-			231
187	ments, including quantities, units, modifiers, and			232
188	their relations, and has become a key benchmark			233
189	for evaluating measurement extraction systems. De-			234
190	spite progress, numerically grounded and relation-			235
191	consistent extraction remains challenging, espe-			236
192	cially for complex sentences containing multiple			237
193	measurements and implicit constraints. Our work			238
194	builds on the MEASEVAL schema and targets these			239
195	persistent failure modes with a hallucination-aware			240
196	optimization framework.			241
197	<b>2.4 Mitigation Strategies for LLM</b>			242
198	<b>Hallucinations</b>			
199	A wide range of techniques have been proposed to			243
200	reduce hallucinations in LLMs, including retrieval-			244
201	augmented generation (RAG) (Lewis et al., 2020),			245
202	supervised fine-tuning (SFT) (Zhou et al., 2023),			246
203	chain-of-thought prompting (Wei et al., 2022),			247
204	process-based supervision (Lightman et al., 2023),			248
205	reinforcement learning from human feedback			249
206	(RLHF) (Ouyang et al., 2022), and direct prefer-			250
207	ence optimization (DPO) (Rafailov et al., 2023).			
208	While effective for open-ended generation, these			
209	methods are not explicitly designed to enforce the			
210	strict grounding and structural consistency required			
211	by scientific measurement extraction. In contrast,			
	our approach integrates reasoning-aware SFT with			212
	targeted reinforcement learning and structured re-			213
	ward shaping, explicitly encoding measurement-			214
	specific constraints to suppress hallucinations at			215
	their structural root.			216
	Despite significant progress in hallucination mit-			217
	igation, prior work has neither systematically char-			218
	acterized hallucinations in scientific measurement			219
	extraction nor introduced specialized reward objec-			220
	tives tailored to its error patterns. We bridge this			221
	gap by unifying a fine-grained hallucination tax-			222
	onomy with a progressive optimization framework			223
	designed specifically for measurement-specific er-			224
	ror suppression.			225
	<b>3 Methodology</b>			226
	Informed by our analysis in Section 2, we design			227
	MEASHALU around a central hypothesis: <i>scientific</i>			228
	<i>measurement hallucinations arise from two funda-</i>			229
	<i>mentally different failure modes—unreliable quan-</i>			230
	<i>titative reasoning and fragile relational grounding.</i>			231
	Accordingly, our framework adopts a two-branch			232
	mitigation strategy, targeting <b>Quantity Hallucina-</b>			233
	<b>tions</b> and <b>Relation-based Hallucinations</b> respec-			234
	tively. As illustrated in Figure 2, MEASHALU in-			235
	tegrates progressive supervised fine-tuning with			236
	hallucination-aware reinforcement learning, en-			237
	abling the model to internalize strict scientific			238
	grounding constraints directly into its reasoning			239
	process. These hallucinations (see Table 8) signifi-			240
	cantly undermine the reliability of LLMs for this			241
	critical task.			242
	<b>3.1 Quantity Hallucination Mitigation</b>			243
	Unlike prior approaches that employ end-to-end			244
	joint training for quantities and relations, our			245
	method first trains quantity extraction indepen-			246
	dently. Furthermore, following the SFT stage, we			247
	incorporate a GRPO phase specifically driven by			248
	hallucination-targeted rewards to further mitigate			249
	hallucinations.			250
	<b>3.1.1 Progressive Supervised Fine-tuning</b>			251
	To endow the LLM with structured quantity reason-			252
	ing capabilities, we adopt a progressive SFT strat-			253
	egy. Specifically, we first utilize $\mathcal{D}_{\text{aug}}$ to establish			254
	foundational quantity reasoning skills, followed by			255
	fine-tuning on $\mathcal{D}_{\text{trace}}$ to ensure rigorous alignment			256
	with MeasEval standards. The construction details			257
	of these two datasets are elaborated below.			258

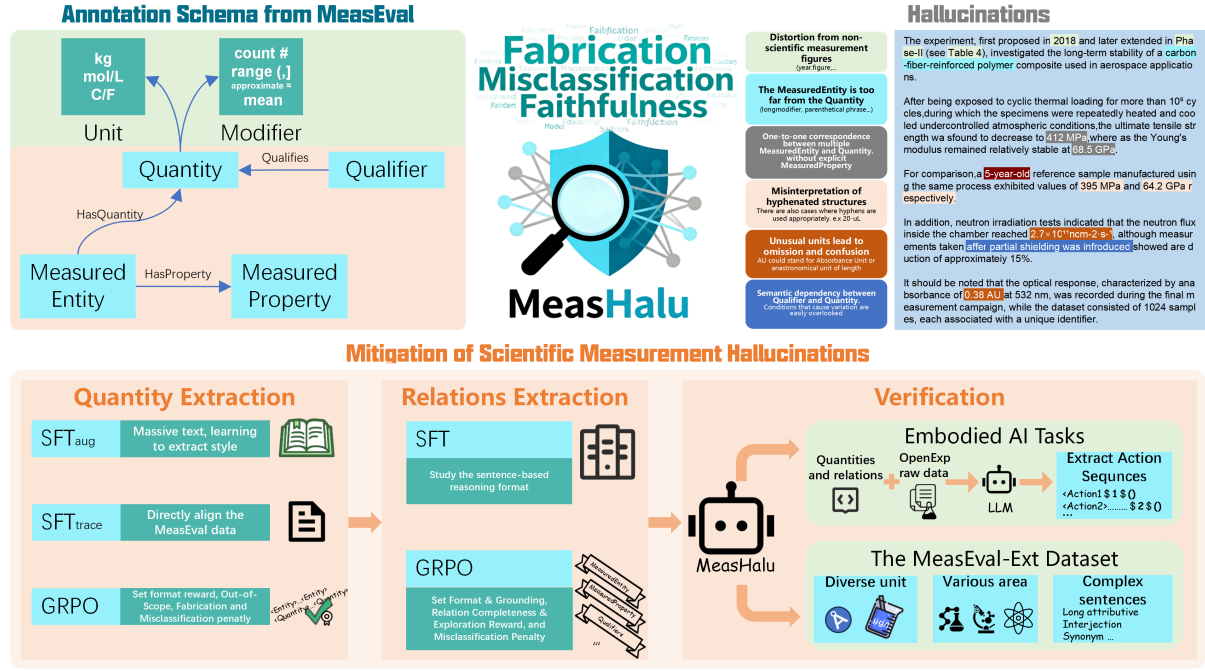


Figure 2: Overview of our method consisting of two stages, Supervised Fine-Tuning & GRPO based Reinforcement Learning.

$\mathcal{D}_{\text{aug}}$  We curate an unlabeled corpus  $\mathcal{X}_{\text{un}}$  from arXiv paper abstracts (Cohan et al., 2018). Lacking gold quantity annotations, we use Quantum3 ( $f_{\text{qtm}}$ ,<sup>2</sup>) to extract quantity candidates, then leverage an augmentation template  $\mathcal{P}_{\text{aug}}$  to prompt  $\mathcal{M}$  to verify these anchors and generate a reasoning trajectory  $h_{\text{aug}}$ . Formally, for  $x \in \mathcal{X}_{\text{un}}$ :

$$\tilde{y} \leftarrow f_{\text{qtm}}(x), \quad h_{\text{aug}} \leftarrow \mathcal{M}(x, \tilde{y}; \mathcal{P}_{\text{aug}}) \quad (1)$$

where  $\tilde{y}$  is noisy pseudo-labels from  $f_{\text{qtm}}$ , and  $\mathcal{P}_{\text{aug}}$  guides  $\mathcal{M}$  to filter false positives via semantics. Valid trajectories form  $\mathcal{D}_{\text{aug}} = \{(x, h_{\text{aug}})\}_{i=1}^{20K}$ .

$\mathcal{D}_{\text{trace}}$  We leverage the MeasEval dataset with human-annotated gold quantity labels  $y_{\text{gt}}$ , and adopt a *traceback template*  $\mathcal{P}_{\text{trace}}$  to guide reasoning reconstruction: given  $y_{\text{gt}}$ , the model generates a stepwise reasoning trajectory  $h_{\text{trace}}$  leading to the gold conclusion, formulated as:

$$h_{\text{trace}} \leftarrow \mathcal{M}(x, y_{\text{gt}}; \mathcal{P}_{\text{trace}}) \quad (2)$$

To ensure correctness of the reasoning trajectory, we enforce strict consistency validation via  $\mathbb{I}(\cdot)$ , where  $\text{concl}(\cdot)$  extracts the final quantity from  $h_{\text{trace}}$ . The filtered dataset is constructed as:

$$\mathcal{D}_{\text{trace}} = \{(x, h_{\text{trace}}) \mid \mathbb{I}(\text{concl}(h_{\text{trace}}), y_{\text{gt}}) = 1\} \quad (3)$$

<sup>2</sup><https://github.com/nielstron/quantulum3>

The prompts  $\mathcal{P}_{\text{trace}}$  and  $\mathcal{P}_{\text{aug}}$  are provided in Appendix A.

### 3.1.2 Hallucination-targeted Reward Function

The total reward  $R(y)$  is a weighted sum of four components targeting distinct quantity-related hallucinations:

$$R = w_1 r_{\text{fmt}} + w_2 r_{\text{scope}} + w_3 r_{\text{fab}} + w_4 r_{\text{mis}} \quad (4)$$

**Format compliance reward ( $r_{\text{fmt}}$ ):** A binary reward is assigned for strict adherence to the predefined structure  $\langle \text{ARABIC} \rangle, \dots, \langle \text{CONCLUSION} \rangle$ , enforcing schema compliance and parsability of generated reasoning chains.

**Out-of-scope hallucination penalty ( $r_{\text{scope}}$ ):** A penalty is imposed when the model extracts out-of-scope entities—such as figure labels (e.g., “Fig. 1”)—that do not constitute valid numerical data. This mechanism utilizes pattern recognition to identify and penalize specific noisy strings, while simultaneously penalizing any generated answers that fail to match the ground truth, ensuring that the model avoids generating arbitrary numbers that deviate from the objective definitions.

**Fabrication hallucination penalty ( $r_{\text{fab}}$ ):** This penalty targets invalid quantity fabrication by verifying each extracted entity against a hybrid phys-

ical parser  $\mathcal{T}_{\text{parse}}$ . A penalty is triggered if the extracted string fails to be parsed as a valid physical or numerical quantity, preventing the model from inventing nonsensical values.

**Misclassification hallucination reward ( $r_{\text{mis}}$ ):** A reward is assigned based on token-level precision to mitigate misclassification hallucinations. This mechanism imposes a penalty if the model generates excessively long spans that erroneously incorporate surrounding components, such as the MeasuredEntity, into the quantity extraction.

Detailed mathematical derivations and implementation specifics for these reward components are provided in Appendix D.

### 3.2 Relation-based Hallucination Mitigation

The extraction of relation-based scientific measurements is particularly challenging due to long-range contextual dependencies that frequently induce hallucinations. Compared to traditional rule-based approaches that generate answers after exhaustively processing complex constraints, our approach first pinpoints the quantity-containing sentence, with subsequent reasoning anchored to this local context to extract the quantity and its relations—eliminating cross-sentence hallucination triggers. We implement this strategy via SFT for schema establishment and GRPO for hallucination-targeted alignment.

#### 3.2.1 Quantity-Guided Relation Extraction

Given the document text  $x$  and a list of candidate quantities  $Q_{\text{in}}$ , the extraction pipeline  $A_{\text{halu}}$  follows a two-stage chain-of-thought reasoning process.

First, the model identifies the evidence sentences  $S = s_1, \dots, s_n$  that contain quantities in  $Q_{\text{in}}$ :

$$S \leftarrow A_{\text{halu}}(x, Q_{\text{in}}) \quad (5)$$

where  $S$  denotes the target sentences.

The model then performs fine-grained reasoning over  $S$  to resolve quantity attributes (e.g., units and modifiers) and associate them with their corresponding measured entities or properties, yielding the final structured relations  $\mathcal{R}$ .

This two-stage reasoning is learned via supervised fine-tuning on rule-derived traces from MeasEval annotations.

#### 3.2.2 Hallucination-targeted Reward Function

While sentence-based extraction excels at local entity identification, it often struggles to capture long-range dependency chains (e.g., *MeasuredEntity*, *Qualifier*), which frequently induces inference bias and leads to the under-extraction of sparse components. To mitigate these reasoning biases and suppress the resulting hallucinations, we design a composite reward function  $R$  optimized via GRPO. The reward function are formulated as:

$$R = w_1 r_{\text{fmt}} + w_2 r_{\text{comp}} + w_3 r_{\text{mis}} \quad (6)$$

The design rationale for each reward component is detailed below:

**Format compliance reward ( $r_{\text{fmt}}$ )** A composite reward is assigned to enforce strict adherence to the quantitative schema and ensure textual grounding. It imposes two constraints: first, validating the structural segmentation of reasoning sections to prevent schema collapse; second, verifying that each extracted sentence can be mapped to a valid text span in the source document.

**Relational completeness reward ( $r_{\text{comp}}$ )** To mitigate inference-induced and role-definition hallucinations stemming from broken dependency links, this reward is designed to enforce the structural integrity of the reasoning chain. The mechanism drives comprehensive exploration through a two-tier incentive structure: (1) a stepwise reward for incremental component extraction and a weighted exploration term that prioritizes harder-to-predict components to drive model exploration (2) a completeness bonus awarded only upon full recovery of the gold-standard relation group to enforce the structural integrity of the reasoning chain.

**Misclassification hallucination reward ( $r_{\text{mis}}$ )** A reward is assigned based on token-level precision to mitigate misclassification hallucinations. This mechanism imposes a penalty if the model generates excessively long spans that erroneously incorporate surrounding components. The detailed mathematical formulations for these reward components are provided in Appendix E.

## 4 Experiments

In this section, We evaluate our method on quantity extraction and relation identification using the MeasEval benchmark. Additionally, we introduce

*MeasEval-Ext*, a specialized dataset annotated from recent literature to target novel units and complex expressions absent from the training distribution. Further analyses are conducted including entropy dynamic analysis and its utility as a functional tool within downstream embodied AI tasks.

#### 4.1 Effectiveness of Quantity Hallucination Mitigation Strategies

**Setup** We utilize the Quantity subset of the MeasEval dataset as our primary evaluation benchmark. To verify the scalability and robustness of our approach, we employ the *Qwen2.5-Instruct* series across three different scales (0.5B, 3B, and 7B) as the base models.

Model Setting	0.5B	3B	7B
MeasHalu-Quant	0.749 $\pm$ 0.006	0.812 $\pm$ 0.006	0.849 $\pm$ 0.006
w/o ( $\mathcal{D}_{\text{trace}}$ + GRPO)	0.475 $\pm$ 0.011	0.481 $\pm$ 0.001	0.465 $\pm$ 0.011
w/o ( $\mathcal{D}_{\text{aug}}$ + GRPO)	0.408 $\pm$ 0.028	0.346 $\pm$ 0.008	0.397 $\pm$ 0.027
w/o GRPO	0.596 $\pm$ 0.008	0.539 $\pm$ 0.002	0.585 $\pm$ 0.018

Table 1: Quantity Information Extraction Performance of 0.5B, 3B and 7B Models (Mean  $\pm$  Std)

**Results** Our full model MEASHALU-QUANT achieves consistent performance advantages across all model scales in Table 1. Compared to the baseline without GRPO, the integration of GRPO drives SOTA results of 0.749, 0.812, and 0.849 for 0.5B, 3B, and 7B models, validating that our rule-based reward system enables stable anti-hallucination alignment even for low-capacity models.

Using only gold-standard data (w/o ( $\mathcal{D}_{\text{aug}}$  + GRPO)) gives the lowest scores (e.g., 0.346 for 3B), showing models cannot capture complex multi-domain quantitative annotation rules without prior measurement extraction schema scaffolding.

The single first stage (w/o ( $\mathcal{D}_{\text{trace}}$  + GRPO)) brings marginal gains from data scaling but is sub-optimal, while (w/o GRPO) delivers substantial improvements (e.g., 3B score raised to 0.539). It confirms that the 1st stage initializes quantitative schema adherence, while the 2nd stage enhances generalization across diverse scientific contexts by leveraging multi-domain scientific knowledge.

#### 4.2 Effectiveness of Relation-based Hallucination Mitigation Strategies

In this section, we evaluate our relation-based hallucination mitigation method on the MeasEval dataset and *MeasEval-Ext*, a newly annotated

dataset containing novel expressions absent from MeasEval, designed to assess the model’s generalization and robustness.

**Setup** We use the *Qwen2.5-Instruct* models (0.5B, 3B, 7B) to assess performance across parameter scales. Although MeasEval is a high-quality benchmark, its limited size and dated sources underrepresent emerging units. To evaluate robustness under distribution shift, we introduce *MeasEval-Ext*, annotated strictly following the MeasEval schema.

We employ an adversarial strategy by selecting recent literature containing novel units and complex expressions absent from the training distribution, rigorously testing model generalization beyond memorized vocabulary. Annotations followed MeasEval guidelines (see Appendix C for agreement analysis).

**Results over MeasEval** Table 2 compares complex quantitative relation extraction on the MeasEval test set. MEASHALU-7B achieves an overall F1 of **0.512**, closely matching the competition winner *LIORI* (Davletov et al., 2021) (0.519)<sup>3</sup>, and substantially outperforming other supervised baselines such as *CONNER* (Cao et al., 2021) (0.473) and *Counts* (Gangwar et al., 2021) (0.432).

Our model also surpasses state-of-the-art proprietary LLMs (e.g., GPT-5, Gemini-2.5-Pro). Even with optimized sentence-based prompting, GPT-5 reaches only 0.406 F1, leaving MEASHALU-7B ahead by over 10 points. This result highlights the necessity of our quantitative domain alignment pipeline (SFT + composite reward optimization) for mitigating relational Quantity hallucinations.

Across all baseline LLMs, sentence-based prompting consistently outperforms rule-based prompting (e.g., Gemini-2.5-Pro improves from 0.359 to 0.440), supporting our hypothesis that sentence-level localized reasoning is more effective than rigid global rule-based deduction for complex quantitative relation extraction.

As shown in Table 3, the results on *MeasEval-Ext* expose a significant performance gap: while general LLMs exhibit non-uniform shifts—often struggling with novel expressions—MEASHALU demonstrates robust generalization to unseen distributions, substantially widening its lead over all baselines. Detailed statistics with standard deviations can be found in Table 9 and Table 10.

<sup>3</sup>LIORI uses a six-model ensemble and does not release weights.

Model	Overall	Quantities			Entities		Properties		Qualifiers	
		Quantity	Unit	Modifier	ME	HasQuantity	MP	HasProperty	Qualifier	Qualifies
<i>Top Ranked Systems from the MeasEval Competition</i>										
Baseline	0.239	0.827	0.561	0.000	0.053	0.075	0.064	0.007	0.005	0.000
Counts	0.432	0.861	0.804	0.614	0.406	0.311	0.245	0.183	0.077	0.064
CONNER	0.473	0.855	0.719	0.523	0.398	0.424	0.437	0.257	0.000	0.000
LIORI	0.519	0.861	0.722	0.642	0.437	0.482	0.467	0.318	0.163	0.092
<i>Rule-based Prompting</i>										
Qwen2.5-7b-inst	0.171	0.491	0.478	0.106	0.088	0.045	0.057	0.000	0.040	0.017
Qwen2.5-72b-inst	0.286	0.644	0.826	0.236	0.196	0.147	0.164	0.001	0.076	0.021
DeepSeek-R1	0.253	0.569	0.586	0.240	0.216	0.163	0.163	0.024	0.085	0.029
DeepSeek-V3	0.271	0.657	0.768	0.214	0.239	0.135	0.113	0.001	0.085	0.014
Gemini-2.5-Pro	0.359	0.712	0.784	0.464	0.306	0.266	0.287	0.090	0.146	0.076
GPT-5	0.371	0.804	0.742	0.395	0.361	0.270	0.355	0.020	0.152	0.052
<i>Sentence-based Prompting</i>										
Qwen2.5-7b-inst	0.073	0.151	0.160	0.027	0.066	0.059	0.044	0.031	0.003	0.005
Qwen2.5-72b-inst	0.212	0.403	0.516	0.232	0.204	0.137	0.113	0.087	0.038	0.012
DeepSeek-R1	0.304	0.589	0.711	0.356	0.260	0.198	0.182	0.118	0.113	0.057
DeepSeek-V3	0.320	0.567	0.726	0.355	0.303	0.225	0.226	0.149	0.019	0.000
Gemini-2.5-Pro	0.440	0.782	0.882	0.486	0.436	0.376	0.386	0.280	0.143	0.056
GPT-5	0.406	0.724	0.817	0.500	0.397	0.351	0.355	0.226	0.138	0.042
MeasHalu-0.5B	0.333	0.649	0.734	0.277	0.292	0.254	0.249	0.175	0.240	0.048
w/o GRPO	0.312	0.649	0.717	0.239	0.263	0.241	0.238	0.140	0.069	0.022
MeasHalu-3B	0.448	0.806	0.861	0.446	0.380	0.399	0.384	0.262	0.093	0.038
w/o GRPO	0.433	0.782	0.850	0.449	0.370	0.377	0.365	0.245	0.084	0.043
MeasHalu-7B	0.512	0.848	0.860	0.607	0.455	0.472	0.442	0.310	0.170	0.100
w/o GRPO	0.479	0.846	0.863	0.610	0.429	0.433	0.397	0.272	0.155	0.063

Table 2: Experimental results over the MeasEval Benchmark. Comparing MeasHalu with competition leaders and rule/sentence-based LLM baselines. Top ranks are shaded orange (1st), yellow (2nd), and teal (3rd).

Model	Overall	Quantities			Entities		Properties		Qualifiers	
		Quantity	Unit	Modifier	ME	HasQuantity	MP	HasProperty	Qualifier	Qualifies
<i>Rule-based Prompting</i>										
GPT-5	0.383	0.833	0.706	0.357	0.400	0.282	0.310	0.046	0.135	0.019
DeepSeek-R1	0.252	0.553	0.514	0.217	0.213	0.169	0.141	0.045	0.099	0.052
DeepSeek-V3	0.312	0.724	0.726	0.234	0.270	0.189	0.121	0.012	0.113	0.037
Gemini-2.5-Pro	0.386	0.766	0.707	0.444	0.351	0.312	0.291	0.151	0.140	0.055
Qwen2.5-72b	0.296	0.675	0.747	0.203	0.246	0.187	0.142	0.000	0.091	0.066
Qwen2.5-7b	0.181	0.501	0.353	0.148	0.108	0.078	0.069	0.001	0.038	0.006
<i>Sentence-based Prompting</i>										
GPT-5	0.402	0.750	0.759	0.458	0.435	0.303	0.303	0.239	0.100	0.070
DeepSeek-R1	0.324	0.622	0.648	0.255	0.299	0.235	0.215	0.146	0.075	0.069
DeepSeek-V3	0.299	0.521	0.565	0.229	0.296	0.229	0.209	0.155	0.053	0.048
Gemini-2.5-Pro	0.462	0.827	0.832	0.444	0.472	0.399	0.402	0.332	0.100	0.072
Qwen2.5-72b	0.202	0.343	0.383	0.183	0.207	0.145	0.127	0.106	0.048	0.050
Qwen2.5-7b	0.033	0.065	0.056	0.022	0.030	0.028	0.017	0.008	0.016	0.020
MeasHalu-7B	0.578	0.861	0.832	0.539	0.555	0.551	0.522	0.459	0.159	0.097

Table 3: Experimental results over the MeasEval-Ext.

### 4.3 Further Analysis

**Mechanism of Hallucination Suppression via Entropy Dynamics** Inspired by Cui and Ding (2025), we quantify Cognitive Hesitation via entropy dynamics, adapting the analysis to our task by distinguishing the quantity group (Quantity, Unit,

Modifier) from the relation group (MeasuredEntity, MeasuredProperty, qualifier).

We focus on tokens strictly bounded by square brackets (e.g., parsing  $70\text{ m}$  from the tagged sequence  $\dots$  surface form  $[70\text{ m}]\dots$ ). To capture micro-level certainty, we report four key statistics:

Bracket Entropy Mean ( $H_B$ ) and Std ( $\sigma_B$ ) measure the average confidence level; Spike Rate ( $R_B$ ) for the proportion of brackets containing high-entropy tokens; and Sample Spike Ratio ( $R_{sample}$ ) quantifies the proportion of samples containing at least one high-risk fluctuation.

Group	Metric	w/o GRPO	MeasHalu	Change
Quantity	$H_B$	0.0071 bits	<b>0.0034 bits</b>	↓ 52.1%
	$\sigma_B$	0.0729	<b>0.0477</b>	↓ 34.6%
	$R_B$	0.16%	0.32%	-
	$R_{sample}$	0.77%	1.54%	-
Relation	$H_B$	0.1147 bits	<b>0.0662 bits</b>	↓ 42.3%
	$\sigma_B$	0.3253	<b>0.2359</b>	↓ 27.5%
	$R_B$	13.01%	<b>5.23%</b>	↓ 59.8%
	$R_{sample}$	33.85%	<b>14.62%</b>	↓ 56.8%

Table 4: Fine-grained Entropy Statistics by Semantic Role.

Table 4 reveals a clear dichotomy in hallucination suppression across semantic roles. 1) *Quantity Group Stability*: SFT is already near-deterministic ( $H \approx 0.0071$ ). GRPO further compresses residual uncertainty ( $H \approx 0.0034$ ) with negligible spike fluctuations ( $R_{sample} \approx 1.54\%$ ). 2) *Relation Group Sharpening*: GRPO reduces the spike ratio from 33.85% to 14.62% and lowers mean entropy by 42.3%. These results indicate that relational reasoning, which is highly ambiguous under SFT, becomes substantially more stable under GRPO. We attribute this improvement to the directed collapse induced by GRPO, which truncates long-tail uncertainty and enforces convergence toward deterministic facts.

Subsequently, to illustrate the effect of GRPO training on reasoning stability more intuitively, we select high-entropy points in the reasoning process for a case study. Details are in Appendix B.

**Application for Embodied AI Tasks** To validate the practical value of our fine-grained extraction for embodied AI, we adapt OpenExp (Liu et al., 2024) to a text-to-action generation task, where models generate executable chemical action sequences (e.g., *ADD . . . (100 mg)*) from unstructured experimental text, mimicking real-world automated laboratory scenarios.

We construct *OpenExp-Action-100*, a dataset of 100 diverse instances, by using unstructured experimental narratives as inputs and OpenExp’s linearized action sequences as gold-standard outputs. To enable controlled comparison, we further define three experimental settings: Baseline (no augmen-

tation), Gemini-Aug (quantity relations extracted by Gemini) and MeasEval-Aug (quantity relations extracted by MeasHalu).

Model	Source	Val	BLEU		LEV	ROUGE		
			B-2	B-4	50%	R-1	R-2	R-L
Gemini-2.5-Pro	MeasHalu	<b>14.67</b>	<b>58.72</b>	<b>44.23</b>	<b>60.00</b>	<b>71.71</b>	<b>54.24</b>	<b>66.51</b>
	Gemini	12.67	58.26	43.99	59.67	71.28	54.21	66.49
	None	0.33	51.78	37.23	37.33	62.92	46.25	58.90
DeepSeek-R1	MeasHalu	<b>10.67</b>	<b>58.99</b>	<b>43.01</b>	<b>61.33</b>	<b>71.62</b>	<b>52.79</b>	<b>65.85</b>
	Gemini	10.33	58.28	42.37	56.67	71.13	52.31	65.65
	None	8.67	38.32	26.21	22.00	59.12	42.32	53.56
GPT-5	MeasHalu	<b>16.33</b>	<b>51.55</b>	37.43	<b>52.33</b>	<b>71.00</b>	51.62	<b>65.47</b>
	Gemini	13.33	50.61	36.56	48.33	70.73	51.31	65.06
	None	1.35	50.39	<b>39.39</b>	50.17	66.31	<b>52.53</b>	62.76

Table 5: Performance on OpenExp-Action-100 with MeasEval-formatted quantity–relation context from different sources (MeasHalu vs. Gemini). Best scores per model are in bold.

Table 5 shows that injecting structured quantity relations significantly improves Structural Validity (Val, executable/logical consistency), with MeasEval-Aug (82.3%) outperforming Gemini-Aug and Baseline. The modest BLEU improvement (19.8 vs. 16.3) stems from gold-standard granularity mismatch. Specifically, OpenExp’s minimalist annotations omit critical details (e.g., *anhydrous*) that our extraction retains. For embodied AI, structural validity—rather than textual overlap—is pivotal; MeasEval-formatted extraction ensures this validity by capturing critical details, providing constraints for executable instructions and practical utility for perception-to-action pipelines. Table 11 in the Appendix shows the full table with standard deviations.

## 5 Conclusion

The proposal of MEASHALU marks a significant step forward in systematically characterizing measurement hallucinations in large language models within the scientific extraction domain. Experimental results on both in-distribution and newly annotated out-of-distribution benchmarks (*MeasEval-Ext*) show that MEASHALU substantially improves robustness and consistently outperforms strong supervised baselines and state-of-the-art large language models. Ultimately, MEASHALU proves to be a reliable external tool that drives significant gains in downstream applications, validating its utility for embodied AI and AI4Science.

## 574 Limitations

575 Despite advances achieved in this paper, MeasHalu  
576 has notable limitations. First, even though  
577 MeasHalu outperforms all existing baselines, the  
578 extraction performance for sparse components (e.g.,  
579 qualifiers, F1 = 0.170) remains low, hindered by  
580 limited annotated data and ambiguous semantic  
581 dependencies in scientific text. Second, the frame-  
582 work’s generalization to low-resource languages  
583 or domain-specific jargon (e.g., niche engineering  
584 units) is untested, as current training data focuses  
585 on English scientific literature. Third, processing  
586 ultra-long documents with nested measurement re-  
587 lations may introduce computational inefficiencies,  
588 as the sentence-based reasoning strategy requires  
589 contextual localization for each quantity.

## 590 References

591 Satya Almasian, Vivian Kazakova, Philipp Göldner, and  
592 Michael Gertz. 2023a. Cqe: a comprehensive quan-  
593 tity extractor. In *Proceedings of the 2023 Conference*  
594 *on Empirical Methods in Natural Language Process-*  
595 *ing*, pages 12845–12859.

596 Satya Almasian, Vivian Kazakova, Philipp Göldner, and  
597 Michael Gertz. 2023b. [CQE: A comprehensive quan-](#)  
598 [tity extractor](#). In *Proceedings of the 2023 Conference*  
599 *on Empirical Methods in Natural Language Process-*  
600 *ing*, pages 12845–12859, Singapore. Association for  
601 Computational Linguistics.

602 Soumia Lilia Berrahou, Patrice Buche, Juliette Dibie-  
603 Barthelemy, and Mathieu Roche. 2013. How to ex-  
604 tract unit of measure in scientific documents? In  
605 *Special Session on Text Mining*, volume 2, pages 249–  
606 256. SCITEPRESS.

607 Jiarun Cao, Yuejia Xiang, Yunyan Zhang, Zhiyuan Qi,  
608 Xi Chen, and Yefeng Zheng. 2021. [CONNER: A cas-](#)  
609 [cade count and measurement extraction tool for sci-](#)  
610 [entific discourse](#). In *Proceedings of the 15th Interna-*  
611 *tional Workshop on Semantic Evaluation (SemEval-*  
612 *2021)*, pages 1239–1244, Online. Association for  
613 Computational Linguistics.

614 Qiguang Chen, Mingda Yang, Libo Qin, Jinhao Liu,  
615 Zheng Yan, Jiannan Guan, Dengyun Peng, Yiyan  
616 Ji, Hanjing Li, Mengkang Hu, Yimeng Zhang, Yi-  
617 hao Liang, Yuhang Zhou, Jiaqi Wang, Zhi Chen, and  
618 Wanxiang Che. 2025. [Ai4research: A survey of ar-](#)  
619 [tificial intelligence for scientific research](#). *Preprint*,  
620 arXiv:2507.01903.

621 Arman Cohan, Franck Dernoncourt, Doo Soon Kim,  
622 Trung Bui, Seokhwan Kim, Walter Chang, and Nazli  
623 Goharian. 2018. [A discourse-aware attention model](#)  
624 [for abstractive summarization of long documents](#).  
625 *Proceedings of the 2018 Conference of the North*

*American Chapter of the Association for Computa-*  
626 *tional Linguistics: Human Language Technologies,*  
627 *Volume 2 (Short Papers)*. 628

Ganqu Cui and Ning Ding. 2025. The entropy mech-  
629 anism of reinforcement learning for reasoning lan-  
630 guage models. *Computing Magazine of the CCF*,  
631 1(7):26–33. 632

John Dagdelen, Alexander Dunn, Sanghoon Lee,  
633 Nicholas Walker, Andrew S Rosen, Gerbrand Ceder,  
634 Kristin A Persson, and Anubhav Jain. 2024. Struc-  
635 tured information extraction from scientific text with  
636 large language models. *Nature communications*,  
637 15(1):1418. 638

Adis Davletov, Denis Gordeev, Nikolay Arefyev, and  
639 Emil Davletov. 2021. [LIORI at SemEval-2021 task 8:](#)  
640 [Ask transformer for measurements](#). In *Proceedings*  
641 *of the 15th International Workshop on Semantic Eval-*  
642 *uation (SemEval-2021)*, pages 1249–1254, Online.  
643 Association for Computational Linguistics. 644

Luca Foppiano, Guillaume Lambard, Toshiyuki Am-  
645 agasa, and Masashi Ishii. 2024. Mining experi-  
646 mental data from materials science literature with  
647 large language models: an evaluation study. *Science*  
648 *and Technology of Advanced Materials: Methods*,  
649 4(1):2356506. 650

Akash Gangwar, Sabhay Jain, Shubham Sourav, and  
651 Ashutosh Modi. 2021. [Counts@IITK at SemEval-](#)  
652 [2021 task 8: SciBERT based entity and semantic](#)  
653 [relation extraction for scientific data](#). In *Proceedings*  
654 *of the 15th International Workshop on Semantic Eval-*  
655 *uation (SemEval-2021)*, pages 1232–1238, Online.  
656 Association for Computational Linguistics. 657

Mark A. Hanson, Pablo Gómez Barreiro, Paolo Crosetto,  
658 and Dan Brockington. 2024. [The strain on scientific](#)  
659 [publishing](#). *Quantitative Science Studies*, 5(4):823–  
660 843. 661

Corey Harper, Jessica Cox, Curt Kohler, Antony Scerri,  
662 Ron Daniel Jr., and Paul Groth. 2021a. [SemEval-](#)  
663 [2021 task 8: MeasEval – extracting counts and mea-](#)  
664 [surements and their related contexts](#). In *Proceed-*  
665 *ings of the 15th International Workshop on Semantic*  
666 *Evaluation (SemEval-2021)*, pages 306–316, Online.  
667 Association for Computational Linguistics. 668

Corey Harper, Jessica Cox, Curt Kohler, Antony Scerri,  
669 Ron Daniel Jr, and Paul Groth. 2021b. [Semeval-](#)  
670 [2021 task 8: Measeval–extracting counts and mea-](#)  
671 [surements and their related contexts](#). In *Proceedings*  
672 *of the 15th International Workshop on Semantic Eval-*  
673 *uation (SemEval-2021)*, pages 306–316.  
674

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,  
675 Zhangyin Feng, Haotian Wang, Qianglong Chen,  
676 Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 oth-  
677 ers. 2025. A survey on hallucination in large lan-  
678 guage models: Principles, taxonomy, challenges, and  
679 open questions. *ACM Transactions on Information*  
680 *Systems*, 43(2):1–55. 681

682	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM computing surveys</i> , 55(12):1–38.	738
683		739
684		
685		
686		
687	Olga Kononova, Tanjin He, Haoyan Huo, Amalie Trewartha, Elsa A. Olivetti, and Gerbrand Ceder. 2021. Opportunities and challenges of text mining in materials research. <i>iScience</i> , 24(3).	
688		
689		
690		
691	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	
692		
693		
694		
695		
696		
697		
698	Chaozhuo Li, Pengbo Wang, Chenxu Wang, Litian Zhang, Zheng Liu, Qiwei Ye, Yuanbo Xu, Feiran Huang, Xi Zhang, and Philip S Yu. 2025. Loki’s dance of illusions: A comprehensive survey of hallucination in large language models. <i>arXiv preprint arXiv:2507.02870</i> .	
699		
700		
701		
702		
703		
704	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In <i>The Twelfth International Conference on Learning Representations</i> .	
705		
706		
707		
708		
709	Zhiyuan Liu, Yaorui Shi, An Zhang, Sihang Li, Enzhi Zhang, Xiang Wang, Kenji Kawaguchi, and Tat-Seng Chua. 2024. ReactXT: Understanding molecular “reaction-ship” via reaction-contextualized molecule-text pretraining. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 5353–5377, Bangkok, Thailand. Association for Computational Linguistics.	
710		
711		
712		
713		
714		
715		
716		
717	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. <i>arXiv preprint arXiv:1808.09602</i> .	
718		
719		
720		
721		
722	David Nadeau and Satoshi Sekine. 2007. A survey of named entity recognition and classification. <i>Linguisticae Investigationes</i> , 30(1):3–26.	
723		
724		
725	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	
726		
727		
728		
729		
730		
731	Maciej P Polak and Dane Morgan. 2024. Extracting accurate materials data from research papers with conversational language models and prompt engineering. <i>Nature Communications</i> , 15(1):1569.	
732		
733		
734		
735	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language	
736		
737		
	model is secretly a reward model. <i>Advances in neural information processing systems</i> , 36:53728–53741.	740
		741
		742
		743
		744
	Tarek Saier, Mayumi Ohta, Takuto Asakura, and Michael Färber. 2024. Hyperpie: Hyperparameter information extraction from scientific publications. In <i>Advances in Information Retrieval</i> , pages 254–269, Cham. Springer Nature Switzerland.	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	745
		746
		747
		748
		749
		750
	Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, and 1 others. 2023. Lima: Less is more for alignment. <i>Advances in Neural Information Processing Systems</i> , 36:55006–55021.	751
		752
		753
		754
		755

## A Prompt template

### Prompt for $\mathcal{P}_{\text{trace}}$

**Instruction:**

You are an expert in extracting structured annotations from text. I have an text input and you need to extract all the quantities within it. I need you to strictly follow the format with six specific sections: ARABIC-QUANTITY, NUMERIC-QUANTITY, TIME-QUANTITY, CHANGE-QUANTITY, CHANGE-QUANTITY, FORMULA-QUANTITY, CONCLUSION.

To explain further: In ARABIC-QUANTITY, outline a step-by-step thought process you use to extract quantity in arabic form. In NUMERIC-QUANTITY, outline a step by step thought process ... In CONCLUSION, give the final answer in a tsv format explained below.

I will provide you with the quantities extracted using the quantulum library for your reference, the information provided by Quantulum is standardized. You need to find the original text in the passage and fill in the tsv form. Also, the quantulum information maybe incorrect, You can't follow it completely.

Here's how the format should look: <ARABIC-QUANTITY> [Provide a chain-of-thought explanation of how you extract all quantities in the arabic forms] </ARABIC-QUANTITY> <NUMERIC-QUANTITY> ... <CONCLUSION>[State the final answer in a tsv format explained below format. ...] </CONCLUSION>

**Task Definition: Extract Quantities**

1. Annotation of Quantities: ...

2. Example Process: ...

Output Format (TSV Fields): ...

Final Output Example: ...

**The reference answer from quantulum: ...**

### Prompt for $\mathcal{P}_{\text{aug}}$

**Instruction:**

You are an expert in extracting structured annotations from text. I have an text input and you need to extract all the quantities within it. I need you to strictly follow the format with six specific sections: ARABIC-QUANTITY, NUMERIC-QUANTITY, TIME-QUANTITY, CHANGE-QUANTITY, CHANGE-QUANTITY, FORMULA-QUANTITY, CONCLUSION.

To explain further: In ARABIC-QUANTITY, outline a step-by-step thought process you use to extract quantity in arabic form. In NUMERIC-QUANTITY, outline a step by step thought process ... In CONCLUSION, give the final answer in a tsv format explained below. It is crucial that you adhere to this structure exactly as outlined and that the final answer in the CONCLUSION matches the standard correct answer precisely.

Here's how the format should look: <ARABIC-QUANTITY> [Provide a chain-of-thought explanation of how you extract all quantities in the arabic forms] </ARABIC-QUANTITY> <NUMERIC-QUANTITY> ... <CONCLUSION>[State the final answer in a tsv format explained below format. ...] It is crucial that you adhere to this structure exactly as outlined and that the final answer] </CONCLUSION>

**Task Definition: Extract Quantities**

1. Annotation of Quantities: ...

2. Example Process: ...

Output Format (TSV Fields): ...

Final Output Example: ...

**The gold answers: ...**

## B Case Study: High-Entropy Token Suppression by GRPO

To better understand the effect of GRPO, we examine a representative sample:

### Input:

Samples were then annealed in air in a pre-heated furnace at temperatures up to 798 °C for times chosen to ensure complete iron diffusion through the sample.

### Gold Standard Annotation:

- **Quantity:** up to 798 °C
- **Unit:** °C
- **Modifiers:** ["IsRange"]
- **MeasuredProperty:** annealed in air in a pre-heated furnace at temperatures
- **MeasuredEntity:** Samples

### SFT Output:

We can find the quantity with surface form [up to 798 °C], it has unit [°C]. The modifier for the quantity are [IsRange]. This quantity is used to describe the entity [furnace]. The entity has the following property [temperatures].

Here, the red token indicates a high-entropy token (entropy = 1.39), and the top-5 candidate tokens with probabilities are listed in the first column of Table 6.

### GRPO Output:

We can find the quantity with surface form [up to 798 °C], it has unit [°C]. The modifier for the quantity are [IsRange]. This quantity is used to describe the entity [Samples]. The entity has the following property [temperatures].

There is no high-entropy token in the model’s output after the GRPO training. The top-5 candidate tokens with probabilities at the [Samples] position are shown in the second column of Table 6.

**Analysis:** GRPO successfully suppresses the high-entropy token observed in SFT, assigning the correct token Samples with high confidence and eliminating uncertainty, demonstrating improved reasoning stability and more deterministic output.

Rank	w/o GRPO		GRPO	
	Candidate	Prob	Candidate	Prob
1	f	0.547	Samples	0.847
2	Samples	0.376	‘ Samples	0.115
3	‘ Samples	0.051	f	0.037
4	samples	0.015	samples	0.0013
5	pre	0.011	_samples	0.0001

Table 6: Comparison of Top-5 candidate tokens at the target position between SFT and GRPO outputs.

## C MeasEval-Ext and its Annotation Details

The annotations are drawn from recent research papers that postdate the original MeasEval corpus. A distinct advantage of this data source is its adversarial selection strategy: unlike the randomized distribution in the original dataset, we deliberately curated 135 text segments (the same as the MeasEval evaluation dataset) containing **novel units and complex quantity expressions absent from the training distribution**. This design ensures that the dataset strictly tests the model’s ability to identify and ground quantities based on semantic context rather than memorized vocabulary.

To ensure high data quality, we enlisted researchers from our laboratory as annotators. The annotation process strictly followed the official *MeasEval Annotation Guidelines*. All samples were **independently labeled by two annotators** to capture the dense quantity-centric information. Following the initial annotation, results were reviewed and reconciled during an **adjudication meeting** to resolve disagreements and reach a final consensus.

The consistency of the dataset is validated by the Inter-Annotator Agreement (IAA). As shown in Table 7, the Krippendorff’s Alpha scores (e.g., 0.921 for Quantity) indicate strong agreement, comparable to the original MeasEval benchmarks.

Annotation Class	Krippendorff’s $\alpha$
Quantity	0.921
MeasuredEntity	0.639
MeasuredProperty	0.584
Qualifier	0.416

Table 7: Inter-Annotator Agreement (Krippendorff’s Alpha) for MeasEval-Ext

## D Quantity Phase Reward

The total reward  $R(y)$  is a weighted sum of four components for mitigating distinct Quantity hallucination types, which includes the Format Reward ( $r_{\text{fmt}}$ ), the Out-of-Scope Penalty ( $r_{\text{scope}}$ ), the Fabrication Penalty ( $r_{\text{fab}}$ ) and the Misclassification Penalty ( $r_{\text{mis}}$ ):

$$R(y) = r_{\text{fmt}}(y) + r_{\text{scope}}(y) + r_{\text{fab}}(y) + r_{\text{mis}}(y) \quad (7)$$

$r_{\text{fmt}}$  To enforce output schema compliance, we validate the sequential semantic tags  $\mathcal{S}_{\text{tags}} = \{\langle \text{ARABIC} \rangle, \dots, \langle \text{CONCLUSION} \rangle\}$  via regex pattern  $\mathcal{P}_{\text{struct}}$ . The binary reward is:

$$r_{\text{fmt}}(y) = \mathbb{I}(y \equiv \mathcal{P}_{\text{struct}}) \quad (8)$$

$r_{\text{scope}}$  Constrains out-of-scope entities (e.g., “Fig. 1”) via local patterns  $\mathcal{C}(e)$  and global precision  $P_{\text{ans}}$ :

$$r_{\text{scope}} = -\lambda_{\text{loc}} \sum_e \mathcal{C}(e) + \beta_{\text{scope}} P_{\text{ans}} \quad (9)$$

$r_{\text{fab}}$  Prohibiting invalid quantity fabrication via parsers combining CQE (Almasian et al., 2023b) and Quantulum<sup>4</sup>  $\mathcal{T}_{\text{parse}}$ , the penalty includes grounding constraints:

$$r_{\text{fab}}(y) = -\lambda_{\text{fab}} \sum_{e \in \mathcal{E}_y} \mathbb{I}(\mathcal{T}_{\text{parse}}(e) = \emptyset) \quad (10)$$

$r_{\text{mis}}$  Mitigating span boundary errors via token-level precision  $P_{\text{tok}}$ , the reward is:

$$r_{\text{mis}}(y) = F \bar{1}_{\text{tok}} - \lambda_{\text{mis}} \cdot (1 - P_{\text{tok}}) \quad (11)$$

## E Relation Phase Reward

While the sentence-based extraction excels at local entity identification (e.g., Units, Modifiers), it suffers from failures in capturing long-range dependency chains (e.g., *MeasuredEntity*, *MeasuredProperty*, *Qualifier*), inference bias, and under-extraction of sparse components. To address these issues, enforce logical completeness, suppress Quantity hallucinations, and incentivize sparse component retrieval, we design a composite reward function  $R(y)$  optimized via GRPO. The total reward is a weighted sum of three dedicated components (Format & Grounding Reward ( $r_{\text{fmt}}$ )), Relation Completeness & Exploration Reward ( $r_{\text{comp}}$ ),

and Misclassification Penalty ( $r_{\text{mis}}$ )), that target distinct quantitative extraction flaws:

$$R(y) = r_{\text{fmt}}(y) + r_{\text{comp}}(y) + r_{\text{mis}}(y) \quad (12)$$

The reward components are elaborated as follows with explicit optimization objectives and mathematical formulations:

$r_{\text{fmt}}$  To enforce structural consistency, we validate the existence of analysis sections  $\mathcal{S}_y$  and adherence to the SFT schema  $\mathcal{F}_{\text{SFT}}$ . The binary reward is:

$$r_{\text{fmt}}(y) = \mathbb{I}(\mathcal{S}_y \neq \emptyset \wedge y \models \mathcal{F}_{\text{SFT}}) \quad (13)$$

$r_{\text{comp}}$  Drives **comprehensive exploration** by aligning predicted groups  $p$  with gold groups  $g$ . To prevent partial extraction, we incentivize full recovery via stepwise matching, closure bonuses, and weighted component bonuses:

$$r_{\text{comp}}(y) = \sum_{p \sim g} \left( \underbrace{\lambda_{\text{step}} |p \cap g|}_{\text{Stepwise}} + \underbrace{\beta_{\text{full}} \mathbb{I}(g \subseteq p)}_{\text{Closure}} \right) + \underbrace{\lambda_{\text{exp}} \sum_c w_c F 1_c^{\text{ans}}}_{\text{Weighted Exploration}} \quad (14)$$

where weights  $w$  prioritize harder-to-predict dependencies to ensure no critical node is missed.

$r_{\text{mis}}$  Suppresses over-broad spans by penalizing token-level precision loss ( $1 - P_{\text{tok}}$ ):

$$r_{\text{mis}}(y) = F 1_{\text{tok}} - (1 - P_{\text{tok}}) \quad (15)$$

<sup>4</sup><https://github.com/nielstron/quantulum3>

Type	Fabrication	Out-of-Scope	Misclassification	Inference Bias	Role Definition
<b>Quantity</b>	Generates fictitious values absent from the source text, or yields extracted strings that contain no valid numerical content.	Extracts invalid numerical tokens, including figure citations (e.g., “Fig. 4”) or scientific nomenclature containing digits (e.g., “4S RNA”).	Generates excessively long spans that erroneously incorporate surrounding components, such as the <i>MeasuredEntity</i> .	–	–
<b>Relation</b>	–	–	Generates excessively long spans that erroneously incorporate surrounding context or unrelated text segments.	Propagates errors from preceding components (e.g., incorrect <i>MeasuredEntity</i> will result in cascading hallucinations in properties and qualifiers).	Fails to distinguish semantic roles, frequently inverting the <i>MeasuredEntity</i> and <i>MeasuredProperty</i> .

Table 8: Taxonomy of Hallucinations in Information Extraction

Model	Overall	Quantities			Entities		Properties		Qualifiers	
		Quantity	Unit	Modifier	ME	HasQuantity	MP	HasProperty	Qualifier	Qualifies
<i>Top Ranked Systems from the MeasEval Competition</i>										
Baseline	0.239	0.827	0.561	0.000	0.053	0.075	0.064	0.007	0.005	0.000
Counts	0.432	0.861	0.804	0.614	0.406	0.311	0.245	0.183	0.077	0.064
CONNER	0.473	0.855	0.719	0.523	0.398	0.424	0.437	0.257	0.000	0.000
LIORI	0.519	0.861	0.722	0.642	0.437	0.482	0.467	0.318	0.163	0.092
<i>Rule-based Prompting</i>										
Qwen2.5-7b-inst	0.171±0.028	0.491±0.052	0.478±0.075	0.106±0.011	0.088±0.021	0.045±0.015	0.057±0.008	0.000±0.000	0.040±0.012	0.017±0.006
Qwen2.5-72b-inst	0.286±0.028	0.644±0.035	0.826±0.026	0.236±0.115	0.196±0.038	0.147±0.028	0.164±0.042	0.001±0.002	0.076±0.019	0.021±0.011
DeepSeek-R1	0.253±0.008	0.569±0.002	0.586±0.002	0.240±0.018	0.216±0.017	0.163±0.014	0.163±0.003	0.024±0.010	0.085±0.015	0.029±0.002
DeepSeek-V3	0.271±0.008	0.657±0.018	0.768±0.010	0.214±0.003	0.239±0.009	0.135±0.014	0.113±0.011	0.001±0.002	0.085±0.004	0.014±0.004
Gemini-2.5-Pro	0.359±0.008	0.712±0.007	0.784±0.035	0.464±0.009	0.306±0.009	0.266±0.007	0.287±0.021	0.090±0.009	0.146±0.018	0.076±0.014
GPT-5	0.371±0.004	0.804±0.013	0.742±0.018	0.395±0.020	0.361±0.003	0.270±0.005	0.355±0.014	0.020±0.004	0.152±0.019	0.052±0.006
<i>Sentence-based Prompting</i>										
Qwen2.5-7b-inst	0.073±0.006	0.151±0.008	0.160±0.007	0.027±0.001	0.066±0.009	0.059±0.010	0.044±0.010	0.031±0.010	0.003±0.004	0.005±0.006
Qwen2.5-72b-inst	0.212±0.002	0.403±0.006	0.516±0.005	0.232±0.010	0.204±0.005	0.137±0.007	0.113±0.008	0.087±0.005	0.038±0.011	0.012±0.007
DeepSeek-R1	0.304±0.004	0.589±0.016	0.711±0.021	0.356±0.030	0.260±0.011	0.198±0.012	0.182±0.006	0.118±0.008	0.113±0.030	0.057±0.019
DeepSeek-V3	0.320±0.006	0.567±0.013	0.726±0.005	0.355±0.016	0.303±0.006	0.225±0.018	0.226±0.022	0.149±0.002	0.019±0.009	0.000±0.000
Gemini-2.5-Pro	0.440±0.011	0.782±0.003	0.882±0.003	0.486±0.011	0.436±0.017	0.376±0.020	0.386±0.028	0.280±0.025	0.143±0.019	0.056±0.010
GPT-5	0.406±0.008	0.724±0.007	0.817±0.017	0.500±0.027	0.397±0.002	0.351±0.006	0.355±0.009	0.226±0.002	0.138±0.026	0.042±0.031
MeasHalu-0.5B	0.333±0.006	0.649±0.004	0.734±0.013	0.277±0.013	0.292±0.014	0.254±0.013	0.249±0.010	0.175±0.011	0.240±0.293	0.048±0.002
w/o GRPO	0.312±0.008	0.649±0.001	0.717±0.019	0.239±0.006	0.263±0.004	0.241±0.017	0.238±0.017	0.140±0.004	0.069±0.009	0.022±0.007
MeasHalu-3B	0.448±0.007	0.806±0.012	0.861±0.010	0.446±0.010	0.380±0.005	0.399±0.010	0.384±0.009	0.262±0.008	0.093±0.004	0.038±0.010
w/o GRPO	0.433±0.010	0.782±0.014	0.850±0.003	0.449±0.009	0.370±0.006	0.377±0.017	0.365±0.019	0.245±0.002	0.084±0.006	0.043±0.015
MeasHalu-7B	0.512±0.004	0.848±0.001	0.860±0.008	0.607±0.006	0.455±0.008	0.472±0.009	0.442±0.012	0.310±0.005	0.170±0.005	0.100±0.009
w/o GRPO	0.479±0.005	0.846±0.002	0.863±0.004	0.610±0.015	0.429±0.004	0.433±0.016	0.397±0.016	0.272±0.012	0.155±0.012	0.063±0.010

Table 9: Experimental results over the MeasEval Benchmark. Comparing MeasHalu with competition leaders and rule/sentence-based LLM baselines. Top ranks are shaded orange (1st), yellow (2nd), and teal (3rd).

Model	Overall	Quantities			Entities		Properties		Qualifiers	
		Quantity	Unit	Modifier	ME	HasQuantity	MP	HasProperty	Qualifier	Qualifies
<i>Rule-based Prompting</i>										
GPT-5	0.383±0.004	0.833±0.021	0.706±0.017	0.357±0.031	0.400±0.004	0.282±0.003	0.310±0.020	0.046±0.011	0.135±0.017	0.019±0.012
DeepSeek-R1	0.252±0.006	0.553±0.016	0.514±0.046	0.217±0.018	0.213±0.005	0.169±0.009	0.141±0.015	0.045±0.006	0.099±0.009	0.052±0.018
DeepSeek-V3	0.312±0.003	0.724±0.014	0.726±0.011	0.234±0.009	0.270±0.014	0.189±0.003	0.121±0.005	0.012±0.003	0.113±0.018	0.037±0.006
Gemini-2.5-Pro	0.386±0.009	0.766±0.005	0.707±0.012	0.444±0.026	0.351±0.014	0.312±0.014	0.291±0.028	0.151±0.017	0.140±0.017	0.055±0.006
Qwen2.5-72b	0.296±0.007	0.675±0.012	0.747±0.010	0.203±0.025	0.246±0.015	0.187±0.015	0.142±0.004	0.000±0.000	0.091±0.014	0.066±0.009
Qwen2.5-7b	0.181±0.009	0.501±0.016	0.353±0.031	0.148±0.037	0.108±0.015	0.078±0.010	0.069±0.012	0.001±0.001	0.038±0.025	0.006±0.004
<i>Sentence-based Prompting</i>										
GPT-5	0.402±0.007	0.750±0.002	0.759±0.005	0.458±0.007	0.435±0.013	0.303±0.007	0.303±0.010	0.239±0.018	0.100±0.016	0.070±0.011
DeepSeek-R1	0.324±0.013	0.622±0.021	0.648±0.018	0.255±0.012	0.299±0.014	0.235±0.019	0.215±0.008	0.146±0.013	0.075±0.015	0.069±0.014
DeepSeek-V3	0.299±0.012	0.521±0.023	0.565±0.020	0.229±0.007	0.296±0.017	0.229±0.012	0.209±0.016	0.155±0.008	0.053±0.004	0.048±0.004
Gemini-2.5-Pro	0.462±0.007	0.827±0.009	0.832±0.007	0.444±0.010	0.472±0.014	0.399±0.012	0.402±0.006	0.332±0.015	0.100±0.002	0.072±0.003
Qwen2.5-72b	0.202±0.006	0.343±0.008	0.383±0.006	0.183±0.007	0.207±0.012	0.145±0.012	0.127±0.007	0.106±0.010	0.048±0.001	0.050±0.007
Qwen2.5-7b	0.033±0.003	0.065±0.002	0.056±0.005	0.022±0.005	0.030±0.001	0.028±0.006	0.017±0.005	0.008±0.003	0.016±0.007	0.020±0.007
MeasHalu-7B	0.578±0.002	0.861±0.003	0.832±0.012	0.539±0.005	0.555±0.006	0.551±0.008	0.522±0.006	0.459±0.007	0.159±0.021	0.097±0.004

Table 10: Experimental results over the MeasEval-Ext.

Inference Model	Context Source	Validity	BLEU		LEV	ROUGE		
			B-2	B-4	50%	R-1	R-2	R-L
Gemini-2.5-Pro	MeasHalu	14.67±2.31	58.72±0.23	44.23±0.29	60.00±2.00	71.71±0.18	54.24±0.32	66.51±0.26
	Gemini	12.67±1.53	58.26±0.60	43.99±0.67	59.67±3.06	71.28±0.39	54.21±0.55	66.49±0.44
	None	0.33±0.58	51.78±0.32	37.23±0.36	37.33±6.66	62.92±0.39	46.25±0.27	58.90±0.26
DeepSeek-R1	MeasHalu	10.67±1.53	58.99±0.37	43.01±0.56	61.33±2.89	71.62±0.19	52.79±0.12	65.85±0.40
	Gemini	10.33±0.58	58.28±0.73	42.37±0.61	56.67±2.89	71.13±0.38	52.31±0.36	65.65±0.31
	None	8.67±1.15	38.32±1.02	26.21±0.99	22.00±5.29	59.12±0.31	42.32±0.21	53.56±0.41
GPT-5	MeasHalu	16.33±2.52	51.55±0.29	37.43±0.52	52.33±3.51	71.00±0.39	51.62±0.27	65.47±0.32
	Gemini	13.33±2.08	50.61±0.72	36.56±0.81	48.33±4.04	70.73±0.80	51.31±0.64	65.06±0.66
	None	1.35±1.54	50.39±0.58	39.39±0.78	50.17±2.33	66.31±0.35	52.53±0.40	62.76±0.44

Table 11: Performance comparison on the OpenExp-Action-100 dataset. Models are provided with MeasEval-formatted quantities and relations generated by different sources (MeasHalu vs. Gemini) as context. The best scores for each model are highlighted in bold.

Example:

The use of the guiding center approximation is justified since the gyroradius of energetic electrons is much smaller than Rhea’s diameter (15–35 km for electrons between 20 and 100 keV) or the scale size of the various macroscopic interaction features (wake, expansion fans). Furthermore, field parameters in the simulation are static.

Inference method (traditional)	Answer
<pre> &lt;UNIT&gt; (1) First, identify all quantities in the text: "15–35 km" and "between 20 and 100 keV"... &lt;UNIT&gt; &lt;MODIFIER&gt; (1) For "15–35 km", the modifier is "IsRange" because it represents a range of values... &lt;MODIFIER&gt; &lt;MEASURED-PROPERTY&gt; (1) For "15–35 km", the measured property is "gyroradius" as it describes the property being measured... &lt;MEASURED-PROPERTY&gt; &lt;MEASURED-ENTITY&gt; (1) For "15–35 km", the measured entity is "electrons" as it is the entity being measured... &lt;MEASURED-ENTITY&gt; &lt;QUALIFIER&gt; (1) For "15–35 km", there is no qualifier that constrains the quantity... &lt;QUALIFIER&gt; </pre>	<pre> Quantity      15–35 km  {"mods": ["IsRange"], "unit": "km"} MeasuredProperty gyroradius {"HasQuantity": "T1"} MeasuredEntity  electrons  {"HasProperty": "T2"} Quantity      between 20 and 100 keV {"mods": ["IsRange"], "unit": "keV"} MeasuredProperty gyroradius {"HasQuantity": "T4"} MeasuredEntity  electrons  {"HasProperty": "T3"} </pre>
Inference method (task specific)	Answer
<pre> From the given text: The use of the guiding center approximation is ...  There are quantities and associated descriptions found: ### From the target sentences: The use of the guiding center approximation is justified ... features (wake, expansion fans).  We can find the quantity with surface form [15–35 km], it has unit [km]. The modifier for the quantity are [IsRange]. This quantity is used to describe the entity [Rhea's diameter].  ### From the target sentences: The use of the guiding center approximation is ... features (wake, expansion fans).  We can find the quantity with surface form [between 20 and 100 keV], it has unit [keV]. The modifier for the quantity are [IsRange]. This quantity is used to describe the entity [electrons]. </pre>	<pre> Quantity      15–35 km  {"mods": ["IsRange"], "unit": "km"} MeasuredEntity Rhea's diameter {"HasQuantity": "T1"} Quantity      between 20 and 100 keV {"mods": ["IsRange"], "unit": "keV"} MeasuredEntity electrons {"HasProperty": "T3"} </pre>

Figure 3: Comparison of Sentence-based and Rule-based Reasoning Approaches