ROSCOE: A SUITE OF METRICS FOR SCORING STEP-BY-STEP REASONING

Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, Asli Celikyilmaz

Meta AI Research

{olggol, mpchen, spoff, mcorredor, lsz, maryamfazel, aslic}@meta.com

ABSTRACT

Large language models show improved downstream task performance when prompted to generate *step-by-step* reasoning to justify their final answers (Nye et al., 2021; Wei et al., 2022). These reasoning steps greatly improve model interpretability and verification, but objectively studying their correctness (independent of the final answer) is difficult without reliable methods for automatic evaluation. We simply do not know how often the stated reasoning steps actually support the final end task predictions. In this work, we present ROSCOE, a suite of interpretable, unsupervised automatic scores that improve and extend previous text generation evaluation metrics. To evaluate ROSCOE against baseline metrics, we design a typology of reasoning errors and collect synthetic and human evaluation scores on commonly used reasoning datasets. In contrast with existing metrics, ROSCOE can measure semantic consistency, logicality, informativeness, fluency, and factuality — among other traits — by leveraging properties of step-by-step rationales. We empirically verify the strength of our metrics on five human annotated and six programmatically perturbed diagnostics datasets - covering a diverse set of tasks that require reasoning skills and show that ROSCOE can consistently outperform baseline metrics.¹

1 Introduction

Scaling language models has improved state-of-the-art performance on nearly every NLP benchmark (Brown et al., 2020), with large language models (LLMs) performing impressively as few-shot learners (Brown et al., 2020). Despite these achievements, even the largest of these models still struggle with tasks including math word problems (Hendrycks et al., 2021), symbolic manipulation (Rytting & Wingate, 2021), and commonsense reasoning (West et al., 2022). Recent work has shown that prompting (Wei et al., 2022; Wang et al., 2022) or fine-tuning (Lampinen et al., 2022) LLMs to generate *step-by-step* rationales can lead to improvements on reasoning tasks. Some of these include small-scale analysis of specific error types within step-by-step rationales (Lewkowycz et al., 2022; Chowdhery et al., 2022), as shown in Table 1. However, existing works primarily focus on end-task performance. Although text generation evaluation metrics sometimes offer fine-grained quality evaluations (e.g., adequacy, fluency) against human scores (Opitz & Frank, 2021; Leiter et al., 2022), these metrics generally treat the output as a whole, and many of these generative metrics operate on tasks such as summarization or machine-translation rather than reasoning.

¹Code can be found at https://github.com/facebookresearch/ParlAI/tree/main/projects/roscoe. Annotated datasets can be downloaded from https://dl.fbaipublicfiles.com/parlai/projects/roscoe/annotations.zip.

In this paper, we present ROSCOE, a suite of interpretable and fine-grained step-by-step generation evaluation metrics to address the above gaps. Rather than providing one score that only evaluates the generated text on the overall, ROSCOE encapsulates fine-grained metrics under four perspectives: (1) *semantic alignment* defines to what extend the generated reasoning is coherent, and grounded with the source context; (2) *logical inference* evaluates if the generated reasoning steps are consistent within itself and checks for logical fallacies; (3) *semantic similarity* quantifies the degree of similarity between the generated reasoning and the context or between intermediate steps to capture hallucinations or repetitions; and (4) *language coherence* evaluates if the whole chain flows naturally.

To evaluate ROSCOE against existing metrics, we devise a taxonomy of reasoning errors for multi-step generations and use it to create synthetic data and collect human Table 1: Step-by-step reasoning for logical inference task generated from 2-shot GPT-3 (text-davinci-002) (Brown et al., 2020). Error types are highlighted: logical or factual issues, partially correct, somewhat-relevant. Full 2-shot exemplars are in Appendix B.

Context: Shelley is from Virginia, but is visiting that city with that famous market where they throw the fish. Going home next Tuesday!

Question: Is it likely that Shelly will be near the Pacific Ocean this weekend?

GPT-3 Step-by-step reasoning:

- No, it is not likely that Shelley will be near the Pacific Ocean this weekend.
- Shelley is from Virginia, which is on the East Coast of the United States.
- The Pacific Ocean is on the West Coast.
- Shelley is visiting a city with a famous market, which could be anywhere in the world.
- Tuesday is her travel day, so she will likely not be traveling this weekend.

evaluations on commonly used reasoning datasets. Our taxonomy and annotated datasets help us gain deeper insights into the causes of reasoning inconsistencies and weaknesses of LLMs. We evaluate ROSCOE with 18 fine-grained metrics under the above four perspectives. ROSCOE demonstrates performance gains against baseline evaluation metrics on all tasks that require reasoning over context. Additional sensitivity analysis shows that ROSCOE is more robust when dealing with tasks that require logical and arithmetic reasoning.

Contributions. (1) We propose a new taxonomy for reasoning errors, and use it for collecting human annotations and creating synthetic datasets. (2) Using our taxonomy, we propose a new suite of metrics that focus on sequence and step level analysis of step-by-step reasoning. (3) We present extensive comparative analysis on 11 datasets of varied complex reasoning problems demonstrating the strengths of each metric, especially in terms of interpretability relative to baselines, and considerations for use.

2 Related Work

Evaluating Explanations. Free-form natural Language (NL) explanations of model decisions should enable accurate representation of the reasoning process and degree of plausibility (Danilevsky et al., 2020; Jacovi & Goldberg, 2021; Jacovi et al., 2021). A qualitative assessment of NL explanations with correctness labels collected from human judges was presented in (Camburu et al., 2018). Recent work has also investigated automatic metrics for natural language generation (NLG) evaluation including word overlap or embedding based similarly with human written explanations (Clinciu et al., 2021). Though fast and cost-effective, automatic metrics for NLG are not equipped to measure the logical inconsistencies or information gain with thinking steps (Reiter, 2019; Celikyilmaz et al., 2020). Explanations have also been evaluated by collecting datasets, and running correlation analysis to investigate the degree to which an automatic metric correlates with human judgements of clarity, relevance and informativeness (Leiter et al., 2022; Welleck et al., 2022). Although reliable, human evaluation is an expensive, domain specific, and time-consuming process. In comparison, ROSCOE provides generic automatic evaluation procedures that are domain and task specific.

Automatic Metrics. Many NLG evaluation metrics exist in the literature including ones based on: *n*-gram match (Lin, 2004), regression (Sellam et al., 2020), embedding proximity (Zhang et al., 2020), paraphrasing (Thompson & Post, 2020), generation as an evaluator (Yuan et al., 2021); information alignment (Deng et al., 2021); among others. Although these metrics are easy to use, they evaluate the alignment of two texts as a whole and are not designed to assess individual reasoning steps. The closest metrics to ours are CTC (Deng

Table 2: Taxonomy of Step-by-Step Reasoning Errors. Full list of the error types with examples is illustrated in Table 10.

Error Type	Definition
Grammar	Faulty, unconventional, or controversial grammar usage
Factuality	Information about an object (i.e. quantity, characteristics) or a named entity doesn't match with the input context.
Hallucination	Information is not provided in the problem statement and is irrelevant or wrong
Redundancy	Explanation contains redundant information, which even though might be factual, is not required to answer the question
Repetition	Step paraphrases information already mentioned in previous reasoning steps
Missing step	The content of the generated reasoning is incomplete and lacks required information to produce the correct answer.
Coherency	Steps contradict each other or do not follow a cohesive story
Commonsense	Model lacks relations that should be known from general world (e.g., "all ducks are birds")
Arithmetic	Error in math calculations

et al., 2021) and BARTScore (Yuan et al., 2021), as both introduce a set of interpretable metrics to evaluate the similarity between two texts. However, ROSCOE is unique in providing fine-grained interpretations of reasoning steps, determining contradictions, and identifying ordering issues in the reasoning narrative.

Self-Consistency with LLMs. Recent work on improving LLMs performance on complex reasoning tasks uses an ensemble strategy called self-consistency (Wang et al., 2022). This method samples a diverse set of reasoning paths from a language model via reasoning traces prompting and returns the most consistent final answer in the set. Other work evaluates the diversity of a reasoning path (Li et al., 2022), or the consistency of an inference step (Creswell et al., 2022) or finetune LLMs (Zelikman et al., 2022) to improve on difficult NLP tasks. In contrast to these works, we present a suit of metrics that focus on determining the type of the error (e.g., commonsense or logical inconsistency) in a reasoning path, if one exists.

3 REASONING ERROR TAXONOMY AND DATASETS CONSTRUCTION

Problem Formulation. Our goal is to score step-by-step rationales generated by a language model. We assume that the model is given a *source* context $s = \{s_1, \cdots, s_T\}$ of T-sentences indicating a problem statement followed by a question and is prompted to generate step-by-step reasoning (Nye et al., 2021). We refer to this as a *hypothesis* $h = \{h_1, \cdots, h_N\}$ of N-steps, including a final answer as the last step. We do not assume availability of gold step-by-step reasoning *references* $r = \{r_1, \cdots, r_K\}$ of K-steps.

Taxonomy. We propose a new taxonomy of generic reasoning errors for language problem solving. We first conduct manual preliminary analysis on different types of LLMs reasoning errors using five *Human judged* datasets described below. Based on our analysis, we identified nine error types centered on the overall reasoning chain (i.e., the quality of the step-by-step thinking, including consistency with the context and commonsense reasoning). Our taxonomy also includes fine-grained errors marking inconsistency of a reasoning step with the previous steps, whether each step contributes to the final decision, and overall logical inference or fluency issues. The definition of error types is in Table 2, and Table 10 provides examples.

Datasets and Annotations. To evaluate ROSCOE, we select datasets covering diverse set of tasks that require reasoning skills (e.g., logical, arithmetic, and commonsense reasoning tasks). We separate these datasets into two: (1) **Diagnostics** datasets that contain gold standard step-wise reasoning chains, where we synthetically perturb some of the reasoning steps to introduce different generation errors (e.g., missing step, mathematical error, etc.); (2) **Human judged** datasets with model generated step-by-step reasoning outputs where the reasoning error evaluations are solicited from expert judges. We investigate these in §5.

4 REASONING SCORER: ROSCOE

We present our fine-grained metrics under four perspectives: *semantic alignment, semantic similarity, logical inference* and *language coherence*. Each metric is bounded within [0, 1], where 1 indicates the perfect score

and 0 corresponds to failure. A metric is *reference-free* or *unsupervised* when it uses the source and hypothesis $(h \to s)$, while *reference-based* or *supervised* when evaluated between hypothesis and reference $(h \to r)$.

4.1 SEMANTIC ALIGNMENT METRICS (ROSCOE-SA)

At the core of the ROSCOE semantic alignment² metrics is the reasoning alignment vector from the N-step hypothesis \mathbf{h} to the source \mathbf{s} of length T: r-align($\mathbf{h} \to \mathbf{s}$) = $\{\alpha_1, \alpha_2, \cdots, \alpha_N\}$, where each alignment value $\alpha_i = r$ -align($h_i \to \mathbf{s}$) = $[1 + \max_{j=1}^T (\cos(h_i, s_j))]/2 \in [0, 1]$ is the normalized cosine similarity between hypothesis step and most similar sentence in a context, and explicitly measures the grounding of the step-wise reasoning with respect to the source text (illustrated in App. D, Fig. 3). We estimate the alignment vector r-align($\mathbf{h} \to \mathbf{s}$) by matching source text and the reasoning chains on the embeddings of tokens and individual reasoning steps. A similar information alignment score is introduced in CTC (Deng et al., 2021) to measure the confidence that the information of the i-th source document token s_j is grounded by a hypothesis token h_i . Our reasoning alignment is different in that we measure if a hypothesized reasoning step h_i supports the source context \mathbf{s} . Our proposed metrics are summarized in Table 3.

Table 3: Semantic alignment metrics (ROSCOE-SA).

Score	Description
Faithfulness-Step $(h o s)$	This step-level score is based on the alignment from the hypothesis steps to the source sentences, and is calculated as the mean reasoning alignment score over the steps of reasoning (see illustration in Appendix D, Figure 3): $(1/N)\sum_{i=1}^N r$ -align $(h_i \to s)$. Faithfulness measures if the model misinterpreted the problem statement, or the reasoning chain is too vague, irrelevant, or misuses information.
Faithfulness-Token $(h o s)$	We extend step-level embeddings of the Faithfulness-Step by measuring similarities between the token embeddings: $(1/(N+M))\sum_{i=1}^N [r\text{-align}(h_i \to s) + \sum_{j=1}^{M_i} r\text{-align}^{token}(h_{i,j} \to s)]$, as shown in App. D, Fig. 3. M_i is the number of tokens in step h_i , $M = \sum_{i=1}^N M_i$ is the total number of tokens in the reasoning chain, $h_{i,j}$ is the j th token in i th step, and $r\text{-align}^{token}$ is the alignment vector from tokens in step h_i to all tokens in s .
$\begin{array}{c} \textbf{Informativeness-Step} \\ \textbf{(Info-Step)} \ (h \leftrightarrow s) \end{array}$	Measures how well information present in the source is used in the reasoning steps: $[(1/T)\sum_{t=1}^{T} r\text{-align}(s_t \to h) + (1/N)\sum_{i=1}^{N} r\text{-align}(h_i \to s)]/2$. Info-step gives a higher score to reasoning steps that are well-grounded with respect to the source, and identifies the degree of information from source that is covered by the generated hypothesis. A lower Info-Step score corresponds to the reasoning steps that are not related to the source sentences or have missed information provided in the context.
Repetition-Token $(h_i \rightarrow h_j)$	To identify repeated, or paraphrased steps, we look at the token alignment scores between all steps in the hypothesis chain: $1 - \max_{i=2N} \max_{j=1\cdots i-1}[(1/M_i)\sum_{l=1}^{M_i} r$ -align $^{token}(h_{i,l} \to h_j)]$. For each pair of sentences, we look at the mean token alignment, and find those sentences that maximize this alignment score. In other words, Repetition-Token will punish chains where there are at least two steps with high overlap in token embeddings.
$ \begin{array}{c} \textbf{Hallucination} \\ (\boldsymbol{h} \rightarrow (\boldsymbol{s}, \boldsymbol{r})) \end{array} $	To find irrelevant reasoning steps, we use alignment score to identify steps that are both not related to the context and not in the reference chain (to avoid punishing for possibly relevant commonsense knowledge): $1 - \max_{i=1N} ([1-r\text{-align}(h \to s)] \cdot [1-r\text{-align}(h \to r)])$. Here, 1 is an all-ones vector, and (\cdot) is the element-wise product.
Redundancy $(h ightarrow r)$	To find chains that contain information that is not required to solve the problem (i.e., redundant steps), we identify those hypothesis steps that are least aligned with the the reference steps: $\min_{i=1N} r$ -align $(h_i \to r)$. This score punishes chains with steps that are not required for the correct solution.
$ \begin{array}{c} \textbf{Semantic} \\ \textbf{Coverage-Step} \\ ((r,h) \rightarrow s) \end{array} $	This score can be viewed as a measure of how easily a gold reference could be generated by the hypothesis. It compares step level grounding of the hypothesis with respect to the source, and the gold reference grounding: $ (1/T)\sum_{t=1}^K r$ -align $(r_t \rightarrow s) - (1/N)\sum_{i=1}^N r$ -align $(h_i \rightarrow s)$, where $ \cdot $ indicates absolute value.
Reasoning Alignment $(h o r)$	The most straightforward way to evaluate the correctness of the hypothesis chain is to compare the degree of the overlap between the hypothesis and the reference. One way of doing that is to measure the reasoning alignment between them: $(1/N)\sum_{i=1}^{N} r\text{-align}(h_i \to r)$.
$ \begin{array}{c} \textbf{Commonsense} \\ (r \rightarrow (h,s)) \end{array} $	Measures if hypothesis lacks steps that are not stated in the source, but are required to solve the problem such as general world knowledge (e.g., "velocity is distance divided by time", "1 foot is 12 inches", "all ducks are birds", etc.). We detect such information by extracting steps in the reference reasoning that are not grounded by the source text: $1 - \max_{i=1K} ([1 - r - \operatorname{align}(r \to h)] \cdot [1 - r - \operatorname{align}(r \to s)])$.
Missing Step $(r{ ightarrow}h)$	To identify steps that are missing from the hypothesis but could be required to solve the problem, we look at the alignment between reference and the hypothesis, similar to <i>Redundancy</i> . However, here we go through each step in the reference, and check if there is a similar step in the hypothesis: $\min_{i=1K} (r\text{-align}(r_i \rightarrow h))$.

²Semantic alignment refers to determination of relations between concepts with the same or a similar intended meaning (Agirre et al., 2013).

4.2 SEMANTIC SIMILARITY METRICS (ROSCOE-SS)

Semantic similarity metrics quantify the degree of semantic equivalence between pieces of text. As opposed to the ROSCOE-SA metrics, ROSCOE-SS considers text as a whole, rather than relying on text units comparisons. We propose the following metrics summarized in Table 4.

Table 4: Semantic similarity metrics (ROSCOE-SS).

Score	Description
	Similar to Info-Step, this metric quantifies the degree of agreement between the hypothesis chain and the source and is calculated as $[1 + \cos(h, s)]/2$. We embed reasoning chain and source context as a whole, as opposed to using step-wise embeddings in *-Step types of metrics introduced in Table 3.
	Measures repetition-related errors on the step level by checking if it paraphrases information already mentioned in the previous steps: $(1 - \max_{i=2N} \max_{j=1i-1} [\cos(h_i, h_j)])/2$. Unlike Repetition-Token, which is orderless and compares individual tokens in pairs of steps, Repetition-Step considers step embeddings similarity and is more robust to changing contexts.
Semantic Coverage- Chain $(r \leftrightarrow h)$	Reflects the overall degree of similarity between the reference and hypothesis chains, comparing reference and hypothesis embeddings as a whole: $[1 + \cos(r, h)]/2$.

4.3 LOGICAL INFERENCE METRICS (ROSCOE-LI)

Logical inference metrics (Table 5) measure logical errors between pieces of text. We use an NLI model that was trained to classify hypothesis-context pairs into entailment, neutral, and contradiction classes (Laurer et al., 2022) to infer the contradiction probability $p_{\rm contr}$.

Table 5: Logical inference metrics (ROSCOE-LI).

Score	Description
Self-Consistency $(h_i \leftrightarrow h_j)$	Measures logical entailment errors within the reasoning steps: $1 - \max_{i=2N} \max_{j < i} p_{\text{contr}}(h_i, h_j)$. This metric will punish chains where there is a pair of steps that are likely to contradict each other.
	Measures logical entailment errors between any generated reasoning h and the source context s : $1 - \max_{i=1N} \max_{j=1T} p_{\text{contr}}(h_i, s_j)$. Specifically, for each reasoning step we measure the probability that it contradicts any sentence in the context. We take the maximum probability of contradiction over all steps, following the logic that a contradiction anywhere in the reasoning chain signals a failure of the overall argument.

4.4 Language Coherence Metrics (ROSCOE-LC)

To evaluate language coherence (Table 6), we use perplexity PPL as scored by the GPT2-Large model (Radford et al., 2019), and English grammatical acceptability $p_{\rm gram}$ as scored by the classifier model from Krishna et al. (2020). Both models were used as-is with no finetuning.

Table 6: Language coherence metrics (ROSCOE-LC).

Score	Description	
Perplexity-Chain (h)	Average perplexity of all tokens in the generated reasoning steps: $1/PPL(h)$. The context used to score each token is the previous tokens in the current and all previous steps. Steps are joined with a space character. To keep the range and orientation consistent with the other scores we invert the perplexity.	
Perplexity-Step (h_i)	Average perplexity of all tokens in the generated reasoning steps, where the context used to score each token is only the previous tokens within the current step: $1/[(1/N)\sum_{i=0}^{N} \mathrm{PPL}(h_i)]$. To keep the range and orientation consistent with the other scores we invert the perplexity.	
Grammar (h_i)	Probability of grammatical acceptability of each step, averaged over all steps: $(1/N)\sum_{i=0}^{N} p_{\text{gram}}(h_i)$.	

5 EXPERIMENTAL SETUP

Diagnostics Datasets. We construct our first category of labeled datasets by generating perturbations — i.e., deterministic modifications — on half of the reference reasoning steps and assign binary labels based on whether or not a chain has been perturbed. We select seven language understanding and entailment datasets

that require complex problem solving skills, and have reference step-by-step explanations: **Entailment-Bank** (deductive reasoning) (Dalvi et al., 2021), **ProofWriter** (logical reasoning) (Tafjord et al., 2021); three arithmetic reasoning datasets **MATH** (Hendrycks et al., 2021), **ASDIV** (Miao et al., 2020) and **AQUA** (Liang et al., 2018); **EQASC** (explanations for commonsense question answering) (Aggarwal et al., 2021), and **StrategyQA** (question answering with implicit reasoning strategies) (Geva et al., 2021) (see dataset details in App. E.1). Using our taxonomy, we introduce 12 error perturbation rules and apply on these datasets to construct our diagnostics datasets (see details in App. E.3).

Human Judged Datasets. We select our second category of datasets from commonly used complex reasoning tasks: GSM8K (arithmetic reasoning) (Cobbe et al., 2021), DROP (discrete reasoning) (Dua et al., 2019), ESNLI (deductive and commonsense reasoning) (Camburu et al., 2018), COSMOS-QA (commonsense reasoning) (Huang et al., 2019) and SemEVAL (Ostermann et al., 2018) (commonsense reasoning). Wei et al. (2022) provide model generated chain of thought reasoning steps for GSM8K. We used chains produced by the 175b_verification model to annotate for reasoning errors. For other datasets, we prompt GPT-3 LLM (Brown et al., 2020) with few-shot in-context examples to obtain step-by-step reasoning sequences (see examples in App. E.2). We use the error types in our taxonomy in Table 2 as human evaluation perspectives of reasoning errors where we solicit five expert annotators³. The data collection interface provided judges with the source text (e.g., source and a question, or hypothesis, premise, and a question if they entail) and associated reasoning text clearly separated into individual steps. Judges were asked to rate the chain as a whole (e.g., on overall quality) as well as each individual step (e.g., commonsense errors, contradicts with the previous steps). App. Table 16 summarizes the distribution of error types annotated by the judges. See App. F for details.

ROSCOE Training. To obtain reasoning step embeddings, we finetune SimCSE (Gao et al., 2021), a supervised sentence similarity model extending the RoBERTa word embedding model (Liu et al., 2019) on multi-step reasoning datasets we listed in §5 (see details in Table 11)⁴. SimCSE is a contrastive learning model that is trained on triplets of reference reasoning steps, positive and hard-negative hypothesis reasoning steps to minimize the cross-entropy objective with in-batch negatives. For contrastive learning, we use the context and reference reasoning steps as a positive sample (s, r), and context and perturbed reference steps (s, h) as hard-negative pairs. For finetuning, we embed source context and hypothesis chain as a whole, without splitting it into steps. With the finetuned model we embed each individual step, as well as a reasoning chain as a whole. We use the pretrained checkpoint of supervised SimCSE model sup-simcse-roberta-base to initialize our model, and further train it for five epochs on our synthetic train data (details in App. G). We also compare ROSCOE scores calculated against sup-simcse-roberta-base SimCSE model, and all-mpnet-base-v2 sentence embedding model (Reimers & Gurevych, 2019) to understand metrics sensitivity to the embedding method

Baseline Metrics. We use text generation evaluation metrics as baseline metrics and comprehensively examine the ones outlined in §2, which are: n-gram match based metrics including **ROUGE-1**, **ROUGE-2**, and **ROUGE-L** (Lin, 2004); pre-trained scores including **BLEURT** (Sellam et al., 2020), **PRISM** (Thompson & Post, 2020), **BERTScore** (Zhang et al., 2020), **BARTScore** using the *Faithfulness* ($s \rightarrow h$) direction for factuality and relevance, and its finetuned variant BARTScore+CNN+Para **BARTScore**+ (Yuan et al., 2021); and information alignment metrics of **CTC**, **CTC-Relevancy** and **CTC-Consistency**. We also include **BARTScore-P**, which we obtain by finetuneing BART (Lewis et al., 2020) on the same reasoning datasets we use for finetuning our SimCSE embedding models. Most of our ROSCOE metrics are constructed reference-free. We also have metrics that use reference reasoning steps which we examine against human judgements. We use the official code for each metric.

³We chose expert annotators over crowd-sourcing, because our annotation task is cognitively challenging and requires fine-grained annotation.

⁴Fine-tuned model is available at https://huggingface.co/facebook/roscoe-512-roberta-base

Meta Evaluation. We use *Somers'* D^5 (Somers, 1962), which measures the ordinal association between two measured quantities, to meta-evaluate each scorer against synthetic and human scores. We prefer *Somers'* D over more commonly used *Kendall's* τ or *Kendall's* τ -b, because it is better in handling the ties of a biased random variable (Agresti, 2010, Section 7.1.5), which imposes an upper bound on the possible values *Kendall's* τ (-b) can take. For each score Y considered, our correlations are built against the biased random variable $X \in [0,1]$, represented by the perturbation or error presence indicator and evaluated using $D(Y|X) = \tau(X,Y)/\tau(X,X)$.

6 EXPERIMENTAL RESULTS

Controlled Experiments with Diagnostics Datasets. Table 7 shows Somers' *D* correlation for metrics measured reference-free on six different datasets and compares baselines to ROSCOE-* aggregated categories calculated with finetuned embeddings: ROSCOE-SA, ROSCOE-SS, ROSCOE-LI, ROSCOE-LC. Results also include ROSCOE metrics with *all-mpnet-base-v2* (ROSCOE-SA¹, ROSCOE-SS¹) and *sup-simcse-roberta-base* (ROSCOE-SA², ROSCOE-SS²) sentence embedding models. Correlations for ProofWriter are taken on its *depth-5* subset. We report highest correlation scores across perturbations within each dataset. The breakdown of all ROSCOE metrics is in App. Table 18.

We observe that: (1) ROSCOE can outperform all other reference-free methods on all six diagnostic datasets, (2) the gains for ROSCOE-SS are more pronounced in four out of six diagnostics datasets, which suggests that ROSCOE can capture hallucinations and repetitions in step-wise reasoning. On Proofwriter, our scorers show lower correlations, because as shown in Table E.1, the context is a list of facts and rules and the reasoning steps can include unordered fact and rule combinations, but still a correct answer can be deduced. This makes it challenging for ROSCOE to evaluate the steps in sequence. Overall, the correlations of the baseline metrics are much lower than ROSCOE, because the baseline metrics are designed to capture the semantic or lexical overlap between a reference and hypothesis and it is harder to detect logical consistency without a golden reference text. ROSCOE is specifically focused on reference-free settings, and can gauge each individual step against the source and other generated steps. In fact, our metrics also work well against the baselines in

Table 7: Somers' D correlation of different metrics on six **Diagnostics** datasets. Metrics are measured **reference-free** on (s,h). We take the maximum score over different perturbations. The two highest correlations for each dataset are **bolded** and <u>underlined</u>, respectively. Correlations that are not significant $(p\text{-}value \geq 0.05)$ are omitted when aggregating, and "-" denotes an absence of any significant correlation. Breakdown of all baseline and ROSCOE metrics is shown in App. H.1, Table 18.

	EntBank	Math	AQUA	ProofWr.	EQASC	ASDIV
ROUGE-L	0.365	0.156	0.264	0.106	0.315	0.269
BLEURT	0.257	0.148	0.252	0.024	0.447	-
BERTScore	0.380	0.124	0.220	0.117	0.462	0.322
BARTScore	0.358	0.185	0.317	0.081	0.415	-
BARTScore+	0.315	0.164	0.251	0.054	0.297	-
BARTScore-F	0.186	0.128	0.215	0.011	0.276	-
PRISM	0.453	0.208	0.191	0.235	0.436	-
CTC Relev.	0.258	0.188	0.217	0.394	0.485	0.382
CTC Consist.	0.310	0.282	0.157	0.513	0.270	0.396
ROSCOE-SA	0.919	0.939	0.971	0.763	1.000	0.879
ROSCOE-SA	0.913	0.936	0.972	0.771	1.000	0.198
ROSCOE-SA	0.919	0.939	0.971	0.732	1.000	0.515
ROSCOE-SS	0.955	0.924	0.982	0.624	1.000	0.857
ROSCOE-SS	0.909	0.932	0.982	0.631	1.000	0.280
ROSCOE-SS	0.901	0.949	0.991	0.621	1.000	0.289
ROSCOE-LI	0.917	0.331	0.424	0.289	0.793	0.771
ROSCOE-LC	0.604	0.392	0.359	0.788	0.859	0.485

the reference-based setting (comparing against reference reasoning steps). In App. Table 19 we present correlations when metrics are measured as reference-based. We also observe that finetuning SimCSE gives highest improvements on the ASDIV dataset. ASDIV is a 1-step reasoning dataset (see App. Table 12), where step is represented by an equation with one of the arithmetic perturbations added. We

⁵We use SciPy (Virtanen et al., 2020) to calculate correlations and obtain *p-values* from a hypothesis test where the null hypothesis is an absence of association.

hypothesize that including these patterns in finetuning helped the model to better learn relationships between context and equations, and resulted in higher scores. On EQASC dataset, Repetition* scores are able to catch all duplicated steps in a chain, i.e., we can separate perturbed and non-perturbed chains based on the given threshold value for the Repetition* scores, and achieve perfect correlation scores (App. Table 20). To understand if finetuning actually helps to improve scoring, we compare non-aggregated metrics (see details in App. Table 18). We observe, that finetuning indeed helps to improve ROSCOE: on average across datasets, all correlations except Repetition_* scores improve (up to 0.556 on Informativeness-Chain), with mean Repetition-Token not changing, and mean Repetition-Step degrading by 0.005. We speculate that since we finetune the model using reasoning chains and context as a whole, it helps to better capture step-by-step rationales, while possibly degrading on word and sentence-level semantics.

Meta-Evaluations on Human Judgement Datasets. Table 8 reports a summary of meta-evaluation of ROSCOE metrics comparing against baselines on human judged datasets. The correlations are measured based on the presence of a particular error from Table 2 and we report the highest correlation across all error types within each dataset. We observe that: (1) on all tasks, ROSCOE metrics outperform all other baselines when evaluated as reference-free: (2) overall. ROSCOE yields considerably better correlations, which indicates that step-by-step reasoning generations can be more effectively evaluated with ROSCOE. In general, most correlations with human judgements are moderate when compared to the synthetic correlation scores, indicating that step-by-step reasoning evaluation is among the cognitively hard tasks for neural models (Deutsch et al., 2022). Interpretable metrics such as ROSCOE can

Table 8: Somers' D correlations of metrics with **human judgement**. We report the maximum over the error types in Table 2. All metrics are measured **reference-free** on (s,h). The highest two correlations in each column are **bolded** and <u>underlined</u>, respectively. Correlations that are not significant $(p\text{-value} \geq 0.05)$ are omitted when aggregating, and "-" denotes an absence of any significant correlation. Breakdown of all baseline and ROSCOE metrics is shown in App. H.2.

	DROP	GSM8K	ESNLI	COSMOS	SemEVAL
Rouge-L	0.278	0.252	0.557	-0.441	-0.478
BLEURT	0.328	0.256	0.541	0.218	-0.356
BERTScore	0.275	0.235	0.590	-0.420	-0.295
BARTScore	-0.835	-0.546	0.549	-0.544	-
BARTScore+	-0.665	-	0.482	-0.186	-
BARTScore-P	-0.642	-	0.255	-0.207	-
PRISM	-0.733	-0.455	0.580	-0.376	-
CTC-Relevance	0.333	-0.371	0.334	-	-0.349
CTC-Consistency	0.462	-0.174	<u>0.647</u>	0.275	-0.301
ROSCOE-SA	0.578	0.392	0.521	0.555	0.337
ROSCOE-SA ¹	0.790	0.500	0.799	0.638	0.485
ROSCOE-SA ²	0.578	0.392	0.599	0.555	0.337
ROSCOE-SS	0.824	0.514	0.530	0.593	0.411
ROSCOE-SS ¹	0.791	0.471	0.507	0.642	0.508
ROSCOE-SS ²	0.799	0.638	0.531	0.658	0.535
ROSCOE-LI	0.584	0.345	0.531	0.444	0.372
ROSCOE-LC	0.205	-0.184	0.447	-0.212	0.517

provide better information about a model's reasoning skills, thus future work should improve such metrics on aligning with human judgments. In App. H.2, we show fine-grained experimental analysis per each human labeled dataset. Specific examples showcasing ROSCOE scoring abilities are summarized in Table 40.

7 ANALYSIS

How sensitive are ROSCOE metrics against level of errors? To evaluate how well metric values match human assessment of reasoning, we measure sensitivity to the level of errors. We perturb sentences in the MATH (arithmetic) and EntailmentBank (deductive reasoning) diagnostic datasets (similar to \S 5) and inject different levels of errors into the reasoning text. Using randomly selected perturbation types, we construct up to a maximum of 3 perturbations per instance. We measure the correlation (Somers' D) between the reasoning inconsistency level 1, 2, 3 of the reasoning steps (i.e., the number of injected errors) and the metric score. Fig. 1 illustrates the results averaged over different perturbations.

We expect the metrics correlate with humans better when the level of errors is high. Both semantic alignment of the reasoning ROSCOE-SA, and the semantic similarity metrics ROSCOE-SS show consistent behavior

on both datasets, while baseline metrics fluctuate with low correlations. Baseline metrics perform better on EntailmentBank. On MATH, ROSCOE-LC and the baseline metrics show minimal impact, which can be that some of the perturbations applied on the MATH dataset (e.g., RandomOperation, or ShuffleNumbers) are harder to detect with language model based (BARTScore) and NLI model based (ROSCOE-LC) metrics.

What does ROSCOE illuminate about scores across errors and tasks? For an ideal scorer based on ease of use, it would be possible to pick a set of fixed thresholds that had error discrimination power across datasets. However, we show that this dataset-agnostic ideal is currently not possible and an issue endemic across scores, including baselines. We study which metrics correlate strongly with which perturbations, with a focus of consistency across datasets. From this, we plot the interquartile ranges for strongly correlated metric and perturbation pairs. We show a sample of these in Fig. 2, though find that the trends generally hold across metrics and perturbations (see Fig 6). We note that within a given dataset, scores are well separated: the *perturbed* version of a dataset for a given score and perturbation type shows little interquartile overlap with the *original* version. However, this does not hold across datasets - e.g., in (Score: Info-Chain, Perturbation: Repetition), if one were to set a detective threshold for the Repetition perturbation based off EntBank (around 0.95), it would mark almost all values of EQASC as perturbed, even non-perturbed sam-

ples. This shows the challenge of using metrics for classification without

calibration for drifts in both mean and variance across datasets, even if a

metric generally correlates well with detecting a given error.

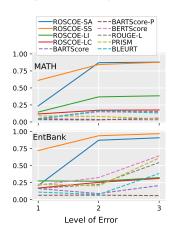


Figure 1: Sensitivity of selected metrics on Somers' *D* by injecting levels of error into reasoning steps.

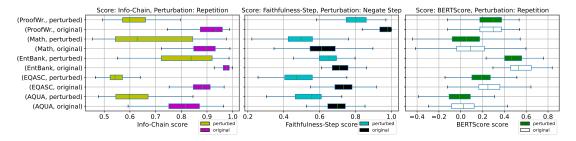


Figure 2: Box-and-whisker plots of interquartile ranges of scores, for perturbations and reference-free metrics with strong Somers' D values. Scores are split by dataset and perturbation use. While interquartile ranges separate well by perturbation use within a single dataset, there is overlap across datasets. This shows the drift of neural scores across datasets and applies to both ROSCOE (left, center) and strong baselines (right).

8 Conclusion

In this paper, we introduce ROSCOE, a new suite of interpretable, unsupervised metrics that enables evaluation of step-by-step reasoning generations of LMs when no golden reference generation exists. We present a taxonomy of reasoning errors used to generate and evaluate our metrics. Experimental results, from evaluating on both synthetic and human-labeled datasets exhibiting multiple types of reasoning (commonsense, arithmetic, and logical inference, etc.), demonstrate superior performance compared to prior semantic and lexical similarly based baseline metrics for text generation. Our analysis shows improved capability in evaluation of reasoning exhibiting nuances, such as factual and logical errors in step-wise decisions.

ETHICS STATEMENT

Explainability builds transparency and trust for users, eases bug-fixing and shortens improvement cycles for metric designers, and will be required by law/regulations for AI systems to be applied to large-scale, high-stakes domains. In this context, we hope our work will catalyze efforts on the topic of explainable evaluation metrics for language model rationale generations. We should mention that our evaluation metrics do not monitor the explanations from integrity or bias perspectives. Our work also uses five human expert annotators and in the annotation process, annotators need to rate the model generated candidate rationals. While the model-generated explanations can produce potentially unsafe content, the datasets for annotations include domains related to logical and arithmetic concepts and general commonsense knowledge. The anecdotal consensus was that the generations were safe and didn't include biased statements.

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our empirical results, we will open source our code to Github, which will contain: instructions for installing the virtual environment, data preprocessing, all score generation and correlation scripts (both for ROSCOE and baselines), and trained embedding models. Detailed explanation of all the finetuned models and metrics are given in the main paper as well as in the Appendices. We will also release all the diagnostic and human judgment datasets used in our experiments.

REFERENCES

Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. Explanations for CommonsenseQA: New Dataset and Models. 2021.

Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. *SEM 2013 shared task: Semantic textual similarity. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pp. 32–43, Atlanta, Georgia, USA, June 2013. Association for Computational Linguistics. URL https://aclanthology.org/S13-1004.

Alan Agresti. Analysis of ordinal categorical data, volume 656. John Wiley & Sons, 2010.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 632–642, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1075. URL https://aclanthology.org/D15-1075.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. 33:1877–1901, 2020. URL https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.

Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. e-snli: Natural language inference with natural language explanations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 9539–9549. Curran Associates, Inc., 2018. URL http://papers.nips.cc/paper/8163-e-snli-natural-language-inference-with-natural-language-explanations.pdf.

- Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. Evaluation of text generation: A survey. *CoRR*, abs/2006.14799, 2020. URL https://arxiv.org/abs/2006.14799.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Miruna-Adriana Clinciu, Arash Eshghi, and Helen Hastie. A study of automatic metrics for the evaluation of natural language explanations. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 2376–2387, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.202. URL https://aclanthology.org/2021.eacl-main.202.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large language models for interpretable logical reasoning. arXiv, 2022. URL https://arxiv.org/abs/2205.09712.
- Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. Explaining answers with entailment trees. *EMNLP*, 2021.
- Marina Danilevsky, Kun Qian, Ranit Aharonov, Yannis Katsis, Ban Kawas, and Prithviraj Sen. A survey of the state of explainable AI for natural language processing. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 447–459, Suzhou, China, December 2020. Association for Computational Linguistics. URL https://aclanthology.org/2020.aacl-main.46.
- Mingkai Deng, Bowen Tan, Zhengzhong Liu, Eric P. Xing, and Zhiting Hu. Compression, transduction, and creation: A unified framework for evaluating natural language generation. In *EMNLP*, 2021. URL https://aclanthology.org/2021.emnlp-main.599.pdf.
- Daniel Deutsch, Rotem Dror, and Dan Roth. Re-examining system-level correlations of automatic summarization evaluation metrics. *arXiv* preprint arXiv:2204.10216, 2022.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2368–2378, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1246. URL https://aclanthology.org/N19-1246.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.

- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Cosmos QA: Machine reading comprehension with contextual commonsense reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2391–2401, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1243. URL https://aclanthology.org/D19-1243.
- Alon Jacovi and Yoav Goldberg. Aligning faithful interpretations with their social attribution. volume 9, pp. 294–310, Cambridge, MA, 2021. MIT Press. doi: 10.1162/tacl_a_00367. URL https://aclanthology.org/2021.tacl-1.18.
- Alon Jacovi, Swabha Swayamdipta, Shauli Ravfogel, Yanai Elazar, Yejin Choi, and Yoav Goldberg. Contrastive explanations for model interpretability. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 1597–1611, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.120. URL https://aclanthology.org/2021.emnlp-main.120.
- Kalpesh Krishna, John Wieting, and Mohit Iyyer. Reformulating unsupervised style transfer as paraphrase generation. In *Empirical Methods in Natural Language Processing*, 2020.
- Andrew K Lampinen, Nicholas Roy, Ishita Dasgupta, Stephanie Cy Chan, Allison Tam, James Mcclelland, Chen Yan, Adam Santoro, Neil C Rabinowitz, Jane Wang, and Felix Hill. Tell me why! Explanations support learning relational and causal structure. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 11868–11890. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/lampinen22a.html.
- Moritz Laurer, Wouter van Atteveldt, Andreu Casas, and Kasper Welbers. Less annotating, more classifying—addressing the data scarcity issue of supervised machine learning with deep transfer learning and bert-nli. 2022.
- Christoph Leiter, Piyawat Lertvittayakumjorn, Marina Fomicheva, Wei Zhao, Yang Gao, and Steffen Eger. Towards explainable evaluation metrics for natural language generation. *CoRR*, abs/2203.11131, 2022. URL https://doi.org/10.48550/arXiv.2203.11131.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871–7880, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.703. URL https://aclanthology.org/2020.acl-main.703.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *arXiv preprint arXiv:2206.14858*, 2022.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. On the advance of making language models better reasoners. arXiv, 2022. URL https://arxiv.org/abs/2206.02336.
- Chao-Chun Liang, Yu-Shiang Wong, Yi-Chung Lin, and Keh-Yih Su. A meaning-based statistical English math word problem solver. pp. 652–662, June 2018. doi: 10.18653/v1/N18-1060. URL https://aclanthology.org/N18-1060.

- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://aclanthology.org/W04-1013.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv* preprint arXiv:1907.11692, 2019.
- Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing english math word problem solvers. pp. 975–984, 2020.
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. *arXiv* preprint arXiv:2112.00114, 2021.
- Juri Opitz and Anette Frank. Towards a decomposable metric for explainable evaluation of text generation from AMR. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 1504–1518, Online, April 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.eacl-main.129. URL https://aclanthology.org/2021.eacl-main.129.
- Simon Ostermann, Michael Roth, Ashutosh Modi, Stefan Thater, and Manfred Pinkal. Semeval-2018 task 11: Machine comprehension using commonsense knowledge. In *SEMEVAL, 2018. URL https://aclanthology.org/S18-1119.pdf.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL https://arxiv.org/abs/1908.10084.
- Ehud Reiter. Natural language generation challenges for explainable AI. In *Proceedings of the 1st Workshop on Interactive Natural Language Technology for Explainable Artificial Intelligence (NL4XAI 2019)*, pp. 3–7. Association for Computational Linguistics, 2019. doi: 10.18653/v1/W19-8402. URL https://aclanthology.org/W19-8402.
- Christopher Michael Rytting and David Wingate. Leveraging the inductive bias of large language models for abstract textual reasoning. 2021. URL https://openreview.net/forum?id=urueR03mkng.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7881–7892, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.704. URL https://aclanthology.org/2020.acl-main.704.
- Robert H Somers. A new asymmetric measure of association for ordinal variables. *American sociological review*, pp. 799–811, 1962.
- Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. ProofWriter: Generating implications, proofs, and abductive statements over natural language. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3621–3634, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.317. URL https://aclanthology.org/2021.findings-acl.317.

- Brian Thompson and Matt Post. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 90–121, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.8. URL https://aclanthology.org/2020.emnlp-main.8.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.
- Sean Welleck, Jiacheng Liu, Ximing Lu, Hannaneh Hajishirzi, and Yejin Choi. Naturalprover: Grounded mathematical proof generation with language models. 2022.
- Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. Symbolic knowledge distillation: from general language models to commonsense models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4602–4625, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.341. URL https://aclanthology.org/2022.naacl-main.341.
- Weizhe Yuan, Graham Neubig, and Pengfei Liu. Bartscore: Evaluating generated text as text generation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 27263–27277. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper/2021/file/e4d2b6e6fdeca3e60e0fla62fee3d9dd-Paper.pdf.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Bootstrapping reasoning with reasoning. arXiv, 2022.
- Tianyi Zhang, Varsha Kishore, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=SkeHuCVFDr.

APPENDIX

A LIMITATIONS

Our study is the first initial step that investigates the evaluation of the step-by-step reasoning produced by large language models. Our taxonomy (in Table 2) covers several reasoning errors and we designed our metrics to evaluate a spectrum of criteria including the ones in the taxonomy. Even though we cannot say we cover all possible reasoning errors, our metrics are generic enough, work on natural language rationales, and consider the alignment with the input context and the generated explanation. Nevertheless, we believe our study can spur others to investigate different reasoning errors and use our code and datasets as templates to extend further.

Due to the extensive analysis needed to thoroughly test and communicate the ability of our proposed metrics to capture reasoning errors, we decided to leave some follow-up questions, such as the application of these metrics for improving downstream task performance, for future exploration.

B FEW-SHOT PROMPTING EXAMPLES (CONT. FROM § 1)

Below is the 2-shot example we used to generate the explanations from GPT-3 as we show in the Fig. 1.

I will explain how to make logical inferences:

Context: John wakes to the sound of his children giggling outside the door. They bust in and say "Happy Father's Day!" with a plate full of food. 2022 has been a great year.

Question: Can we infer that John was born before the inauguration of Barack Obama?

Explanation: We can infer that John is a father because of "his children". We can infer that the current year is 2022. Barack Obama was inaugurated in 2009, so it has been about 13 years since Obama was inaugurated. John's children are old enough to speak, so if John were born after Obama's inauguration, he would have been younger than 10 when he had children. So, we can infer that John was born before Obama's inauguration. The answer is "yes".

Context: Sam sees a piano in the subway stations. Laughing, his friends prod him to play, saying "Remember how amazing your rendition of Chopsticks was at the Christmas party?"

Question: Can we infer that John majored in piano at Juilliard?

Explanation: Juilliard is a top school for music, so someone who majored in piano at Juilliard would be an expert pianist. "Chopsticks" is a very simple song to play on the piano. His friends were therefore likely being sarcastic when they told him to play, which would imply that Sam was not an expert piano player. So it is unlikely that Sam majored in piano at Juilliard. The answer is "no".

Table 9 illustrates sample if in-context examples we used as prompts to obtain step-wise explanations.

Table 9: Example prompts used to generate step-by-step explanations using GPT-3.

(Used for ESNLI dataset)

Explain if premise entails hypothesis.

Premise: Tomato sauce (also known as Neapolitan sauce, and referred to in Italy as Salsa di pomodoro) refers to any of a very large number of sauces made primarily from tomatoes, usually to be served as part of a dish (rather than as a condiment). Tomato sauces are common for meat and vegetables, but they are perhaps best known as sauces for pasta dishes.

Hypothesis: Tomato sauce is only rarely used to top pasta.

Explanation: Salsa di pomodoro is made of peppers. The answer is No.

Premise: Kirk Johnson (born 1960) is an American scientist, author, curator, and museum administrator, and is currently serving as Director of Smithsonian's National Museum of Natural History. Dr. Johnson is the host of the PBS Nova series, Making North America, which is a three-part series that describes the shaping of North America, which aired on November 4, 11 and 18, 2015.

Hypothesis: johnson hosted Making North America when he turned 50.

Explanation: Johnson was born in 1960 and hosted Making North America in 2015, so he was 55, not 50. The answer is No.

Premise: A person rides his bicycle in the sand beside the ocean.

Hypothesis: A person is on a beach.

Explanation: A person rides beside the ocean means is on a beach if there is sand beside the ocean. The answer is Yes.

(used for SEMEVAL dataset)

Explain if premise entails hypothesis.

Premise: I went into my bedroom and flipped the light switch. Oh, I see that the ceiling lamp is not turning on. It must be that the light bulb needs replacement. I go through my closet and find a new light bulb that will fit this lamp and I place it in my pocket. I also get my stepladder and I place it under the lamp. I make sure the light switch is in the off position. I climb up the ladder and unscrew the old light bulb. I place the old bulb in my pocket and take out the new one. I then screw in the new bulb. I climb down the stepladder and I place it back into the closet. I then throw out the old bulb into the recycling bin. I go back to my bedroom and turn on the light switch. I am happy to see that there is again light in my room. Which room did the light go out in?

Hypothesis: Bedroom

Label: Correct

Explanation: The answer is bedroom, because the light went out in the bedroom. The person replaced the light bulb in the bedroom.

Premise: The weather was so nice today that I decided to have a barbecue. I called up some of my friends to invite them over. I set up some chairs and tables outside for my friends to sit around. I then took out some hot dogs from my fridge and brought them over to my barbecue. My barbecue is a gas barbecue, so all I had to do was press the knob and turn it to ignite it. Once I got a little fire going I put my hot dogs on the grill over it to cook. I closed the lid and left my hot dogs to cook for about ten minutes. After ten minutes I decided to check on them to see if they were done. They looked like they were cooked, so I put them all on a plate to serve to my friends. I went to my kitchen and got hot dog buns, ketchup, mustard, and relish to serve on top of the hot dogs. I laid out all of the food and condiments on the table so my friends could help themselves. We ate the hot dogs and relaxed around the table. Why are we having the barbecue, to celebrate an occasion, or just because we want to?

Hypothesis: Just because they want to.

Label: Correct

Explanation: The answer is correct, because the weather was nice. They decided to have a barbecue. There is no mention of any special occasion.

 $\begin{array}{l} \textbf{Premise} \colon \text{Another Nurse can in to look at me, took one look at my rash and said OMG and ran out of the room. At that point I started crying, regretted coming to the hospital alone and started wheezing again. Then the doctor came in, her eyes were all bugged out. What will I do after the doctor examines me? \\ \end{array}$

Hypothesis: I will talk with my friends.

Label: Incorrect

Explanation: The premise is about the nurse's reaction to the patient's rash, and how the patient feels. The hypothesis does not follow up from this premise correctly. Instead, it talks about talking with friends.

C TAXONOMY OF REASONING ERRORS (CONT. FROM § 3)

To gain deeper insights into the types of reasoning errors introduced by LLMs while explaining their decisions, we propose a new taxonomy of generic reasoning errors for language problem solving. Specifically, we sampled from the training portions of the logical inference and commonsense reasoning datasets, and prompted GPT-3 with reasoning explanations using prompts similar to App. B. We used task specific indomain examples for prompting. We also analyzed model generated explanations shared in Wei et al. (2022). We then manually looked into each explanation and identified potential errors that are inconsistent with the source, question or the prompt and within the reasoning chain. Some tasks require a model to classify the logical relationship between premise and a hypothesis, others are question and answering tasks. We adjusted our context and prompts according to the type of the task.

Our reasoning error taxonomy is summarized in Table 10. It contains types of errors concerning an overall chain or an individual step. Specifically, the *chain-level coarse-grained* evaluations of the overall reasoning chain deals with overall quality of the step-by-step thinking, coherence, consistency of the explanation within itself, and consistency with the context, etc. On the other hand the *step-level fine-grained* evaluations focus on the consistency of a reasoning step with the previous steps, if a step conveys new and supporting information over the previous steps, factuality or logical inference issues. We use these error categories to construct diagnostics datasets with perturbed errors as well as human judged datasets of reasoning errors. In the taxonomy, we indicate *-step level errors to differentiate from the chain level error types.

Table 10: Taxonomy of Step-by-Step Reasoning Errors. Errors used for perturbations in constructing the diagnostic datasets (Diag.) and for human annotation (Human) of the model generated reasoning chains are also marked. (Cont. from Table 2.

Error Type	Description	Example	Diag.	Human
Grammar	Faulty, unconventional, or controversial grammar usage	Reference: He chews 4 pieces of gum a day so over 30 days he will chew 4*30 = «4*30=120»120 pieces of gum. Model Expl: He eats 4 pieces of gum a day and lasts 30 days so he eats 4*30 = «4*30=120»120 pieces of gum.	✓	✓
Factuality	Information about object (i.e. quantity, characteristics) or a personal named entity does not match information provided in the question	Context: Anakin caught 10 fish. Model Expl: Locsin caught 10 fish.	√	√
Hallucination	Information is not provided in the problem statement and is irrelevant or wrong	Context: The basketball team went to the steakhouse to eat dinner. The first player ate a 6-ounce steak. The second player ate beef tips, containing 8 beef tips, each an ounce in size. The third player ate a one-pound steak. And the fourth and fifth players ordered vegetarian meals. In total, how many ounces of meat were consumed by the team? Model Expl: The fourth and fifth players ordered vegetarian meals, for a total of 2 ounces of meat. Context: Boh is cold. If something is quiet then it is blue.		√
Redundancy	Explanation contains redundant in- formation, which even though might be factual, is not required to answer the question	Context: Bob is cold. If something is quiet then it is blue. All blue things are round. If something is round and smart then it is cold. Question: Bob is cold. Todel Expl: If something is round and not smart then it is blue. All blue things are round. Answer: true.		✓
Repetition	Step paraphrases information al- ready mentioned in previous reason- ing steps	Context: Parker chews 4 pieces of gum a day. A pack of gum has 15 pieces of chewing gum per pack. How many packs of gum will he need to last him 30 days? Model Expl.: He chews 4 pieces of gum a day and there are 30 days in a month so he chews 4*30 = 120 pieces of gum a month. Parker chews 120 pieces of gum a month. He needs 120 pieces of gum a month and a pack of gum has 15 pieces of gum so he will need 120/15 = 8 packs of gum.	√	√
Missing step	The content of the generated reasoning is incomplete and lacks required information to produce the correct answer. If this missing steps are added, the model would produce the correct answer	rea- greater animals. If all animals are healthy, how many pairs of animal legs in total are on the farm? Answer : 40 Model Expl. : 10 2-legged animals have 10*2=20 legs. 15		√
Coherency	Steps contradict each other or do not follow a cohesive story	Context: If the reactor is intact, then it is safe. If the reactor is safe, then no graphite is outside it. Model Expl.: The reactor is intact. All graphite is inside of it. There are some graphite outside of the the reactor.		✓
Commonsense	Model lacks relations that should be known from general knowledge about the world (i.e. how to com- pute velocity, how many inches in one foot, all ducks are birds, etc.)	Question: Jordan has 2 children who wear diapers. Each child requires 5 diaper changes per day. Jordan's wife changes half of the diapers. How many diapers does Jordan change per day? Model Expl.: Jordan has 2*5=10 diapers to change per day. Jordan's wife changes 10/2=5.0 diapers per day. For both children, Jordan changes 10+5=15 diapers per day.		√
Arithmetic	Error in math calculations	Reference: If a truck was carrying 20 tons of fertiliser packed in bags, the total number of bags in a truck is $20*20 = 400$ bags Model Expl: If a truck was carrying 20 tons of fertiliser packed in bags, the total number of bags in a truck is $20*20 = 40$ bags	√	√

D ROSCOE METRICS DETAILS (CONT. FROM §4)

ROSCOE metrics are constructed under four categories: semantic alignment, semantic similarity, logical inference, and logical coherence. The details of each metric is explained in §4. At the core of ROSCOE semantic alignment metrics is the reasoning alignment score, which we designed to measure the grounding of step-by-step reasoning with respect to the source text. Fig. 3 illustrates the reasoning alignment.

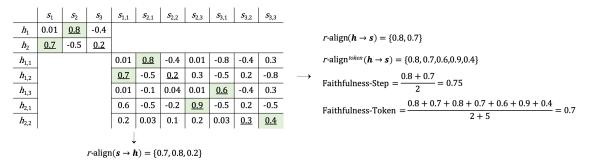


Figure 3: **Reasoning alignment** illustrating the measurement of the Faithfulness-Step and Faithfulness-Token semantic alignment scores. $h = \{h_1, h_2\}$ is a hypothesis chain with tokens $\{h_{1,1}, h_{1,2}, h_{1,3}, h_{2,1}, h_{2,2}\}$, and $s = \{s_1, s_2, s_3\}$ is a context with tokens $\{s_{1,1}, s_{2,1}, s_{2,2}, s_{2,3}, s_{3,1}, s_{3,2}, s_{3,3}\}$. Alignment scores from hypothesis to context are highlighted, and alignment scores from context to hypothesis are underscored. The reasoning alignment combines token and step level similarities where each alignment value (cell) is the cosine similarity and explicitly measures the grounding of the token and step-wise reasoning with respect to the source text.

The variation of scorers of the ROSCOE shares some similarities, thus we explain them here:

BARTScore (Yuan et al., 2021) claims that more high level text can be generated using sequence to sequence model. It can support different evaluation perspectives such as factuality (by evaluating from source to hypothesis) or informativeness (by evaluating from both directions between reference and hypothesis). BARTScore is used to measure the probability of generated text from a source text x to a target set y:

$$BARTScore = \sum_{t=1}^{m} w_t \log p(y_t | y_{< t}, x, \theta)$$
 (1)

BARTScoreintroduce two variations: (1) finetuning, in which the BART model is finetuned on the task specific dataset to make the pre-training domain closer to the evaluation domain. (2) prompting, in which a task specific textual prompt is appended to the source x to get the y. In our experiments we compare the the BARTScorebaseline and one with the prompting variant BARTScore+to compare in the experiments.

CTC (Compression, Transduction, and Creation) (Deng et al., 2021), is a suite of metrics that unifies different perspectives of different tasks (e.g., summarization, style transfer, or text rewriting) into information alignment, which measures weather the information in one generation component is grounded in another. The information alignment is defined as follows: let x (e.g., dialog context) be the source input, c (e.g., external world knowledge) be some additional context, and y be the generated output text (e.g., generated response). The alignment is measured on token level and it is measured as the vector of scores:

$$align(a \to b) = \langle \alpha_1, \cdots, \alpha_N \rangle$$
 (2)

where each score α_i indicates confidence that the *n*-th token in *a* aligns with the whole sentence *b*. Using the information alignment they define a list of metrics to evaluate text for different tasks. In our experiments we use two of these metrics that are closer to ROSCOE: the *Relevance* (CTC Relevance), which measures the consistency of the generated text with the source and its balanced between the reference, and the *Consistency* (CTC Consistency) which deals with the faithfullness of the generated text to the input context by the alignment between the two.

E EXPERIMENTAL SETUP DETAILS (CONT. FROM § 5)

E.1 DIAGNOSTIC DATASETS

In the following we present details of each diagnostics dataset used in our work. Table 11 illustrates how each dataset is used in our experiments. StrategyQA dataset is only used to finetune the SimCSE embeddings model, because it contains reference reasoning chains in train and validation partitions, but not in the test partition. The rest of the six diagnostic datasets are used for sentence embedding model finetuning, and evaluating our models as presented in the experiments results. All datasets with examples are summarised in Table 12.

Table 11: Summary of datasets used in our work. Reasoning Chain represent whether it contains human written golden step-wise reasoning explanation. Type indicates whether it is used for constructing Diagnostic or Human judged datasets. Train/Val./Test indicate whether the dataset is used for training, validation and/or testing. StrategyQA dataset is only used for finetuning SimCSE embedding model.

Dataset	Reasoning Chain	Туре	Train	Val.	Test	Annotated Instances
EntailmentBank (Dalvi et al., 2021)	✓	Diagnostic, Finetuning	1	1	/	1,840
ProofWriter (Tafjord et al., 2021)	✓	Diagnostic, Finetuning	1	✓	1	272,430
MATH (Hendrycks et al., 2021)	✓	Diagnostic, Finetuning	1	✓	1	12,500
ASDIV (Miao et al., 2020)	✓	Diagnostic, Finetuning	1	✓	1	2,305
AQUA (Liang et al., 2018)	✓	Diagnostic, Finetuning	1	✓	1	97,975
EQASC (Aggarwal et al., 2021)	✓	Diagnostic, Finetuning	1	✓	1	9,060
StrategyQA (Geva et al., 2021)	✓	Finetuning	1	✓	X	2,290
DROP (Dua et al., 2019)	X	Human judged	X	X	1	210
GSM8K (Cobbe et al., 2021)	✓	Human judged	X	X	1	200
ESNLI (Camburu et al., 2018)	✓	Human judged	X	X	1	151
CosmosQA (Huang et al., 2019)	×	Human judged	X	X	✓	195
SemEval (Ostermann et al., 2018)	X	Human judged	X	X	✓	209

EntailmentBank (EntBank) (Dalvi et al., 2021) is a complex question answering dataset which contains multi-step entailment trees, namely a tree of multi-premise entailment steps from facts that are known, through intermediate conclusions to hypothesis of interest (which in this case the question and answer).

ProofWriter (Tafjord et al., 2021) is a question answering dataset for logical reasoning. It contains 500k questions, answers and proofs over natural-language rulebases. This dataset is mostly used to emulate reasoning over rules expressed in language, including proof generation. The datasets proofs include intermediate conclusions. In our experiments, we used *depth-0*, *depth-1*, *depth-2*, *depth-3*, and *depth-5 OWA* sets.

MATH (Hendrycks et al., 2021) is a dataset of 12,500 problems from high school math competitions. Given a math problem such as in Table 12 models generate a sequence, such as $\frac{2}{3}$, that encodes the final answer.

ASDIV (Miao et al., 2020) (Academia Sinica Diverse MWP Dataset) is a dataset of 2,305 questions on diverse math word problem solving. It includes a diverse operations such as basic arithmetic or aggregative operations (e.g., comparisons, set-operations).

AQUA (Liang et al., 2018) is a dataset of 100,000 algebraic word problems with step-wise solutions as shown below. In the original dataset each question is decomposed in four parts, two inputs and two outputs: the description of the problem and a question, and the possible (multiple choice) answer options, one being the

Table 12: We show instances from seven of the **Diagnostics Datasets** here. (Continue from §5).

Dataset	Instance	Reference Step-by-Step Solution
EntBank	Earth is a kind of celestial object. Stars appear to move relative to the horizon during the night. A star is a kind of celestial object celestial body. The earth rotating on its axis causes stars to appear to move across the sky at night. Apparent motion is when an object appears move relative to another object 's position. Question: How does the appearance of a constellation change during the night? Hypothesis: Solve the following entailment problem: "Earth is a kind of celestial object. During the night stars appear to move"	Step1: earth is a kind of celestial object Its position appears Step2: a star is a kind of celestial object to shift relative celestial body to the horizon. Step3: apparent motion is when an object appears to move relative to another object 's position Step4 Therefore apparent motion of stars is when stars appear to move relative to earth's position Step5: The earth rotating or its axis causes stars to appear to move across the sky at night Step6: Therefore the earth rotating on its axis causes apparent motion of stars Step7: Stars appear to move relative to the horizon during the night Step8: Therefore the earth rotating on its axis causes stars to move relative to the horizon during the night.
ProofWriter	Facts: The cow is not big. The cow is not green. The lion eats the tiger. The lion sees the cow. The lion visits the cow. The lion does not visit the squirrel the lion visits the tiger. The squirrel is big. The squirrel is round. The tiger is not green. The tiger does not see the cow. Rules: if something sees the squirrel and the squirrel eats the cow then the cow is round. if something is green then it eats the tiger. if the squirrel is round then the squirrel visits the cow. if something eats the cow then it sees the squirrel. if something sees the tiger and the tiger visits the squirrel then it is nice. if something is round then it eats the cow. if something is kind then it eats the cow. if something is kind then it eats the cow. if the tiger visits the cow then the cow sees the squirrel. if something sees the cow then the cow eats the tiger. Question: The cow does not see the squirrel.	Step1: The squirrel is round. Step2: something is round then it eats the cow. Step3: The squirrel eats the cow. Step4: If something sees the squirrel and the squirrel eats the cow then the cow is round. Step5: The cow is round. Step6: If something is round then it eats the cow. Step7: The cow eats the cow. Step8: if something eats the cow then it sees the squirrel. Step9: the cow sees the squirrel. Answer: True
МАТН	Context: Tom has a red marble, a green marble, a blue marble, and three identical yellow marbles. Question: How many different groups of two marbles can Tom choose?	Step1 : There are two cases here: Step2 : either Tom chooses two yellow marbles (1 result), or he chooses two marbles of different colors $(\binom{4}{2} = 6 \text{ results.})$. Step3 : The total number of distinct pairs of marbles Tom can choose is $1 + 6 = 7$. Answer : 7
ASDIV	Context: A sandwich is priced at \$0.75. A cup of pudding is priced at \$0.25. Tim bought 2 sandwiches and 4 cups of pudding. Question: How much money should Tim pay?	0.75 x 2 + 0.25 x 4 = 2.5 Answer: 2.5
AQUA	Context: The entrance fee for a fair is \$5 for persons under the age of 18 and 20% more for persons older. Each ride at the fair costs \$0.50. If Joe goes with her 6 years old twin brothers, and they each took 3 rides in total. Question: How much money does Joe end up spending at the fair?	Step1 : Total entrance fee is (2*\$5)+(1.20*5) = \$16
EQASC	Question : Where is water likely to form beads?	Step1: Beads of water are formed by water vapor condensing Step2: Moisture builds up in condenses air and the wherever the surfaces are cold. Answer: Water beads form on cold surfaces.
StrategyQA	Question : Are more people today related to Genghis Khan than Julius Caesar?	Step1: Julius Caesar had three children. Step2: Genghis Khan had sixteen children. Step3: Modern geneticists have determined that out of every 200 men today has DNA that can be traced to Genghis Khan. Answer: True

correct one. In this work we only used the context and question, the step-wise solution and the correct answer to construct our diagnostic dataset.

EQASC (Aggarwal et al., 2021) is a multi-hop question answering dataset with 98K explanation annotations for multi-step factual reasoning. Each instance in the dataset comes with a question, multiple answer choices,

explanation of each answer choice and a free flow explanation of the whole context. In our experiments we used the correct answer's explanation to construct our diagnostic datasets.

StrategyQA (Geva et al., 2021) is another multi-step question answering (QA) dataset, that covers a diverse set of reasoning skills. StrategyQA consists of 2,780 questions, annotated with their decomposition and per-step evidence.

E.2 HUMAN JUDGED DATASET CONSTRUCTION

In the following we present details of each human judged datasets used in our work. Table 11 lists each dataset and illustrates how each dataset is used in our experiments. Specifically, all six datasets are used for evaluations in the experiments results and model finetuning, and one dataset was used for finetuning only. The dataset details are explained below.

To construct these datasets, we first sample instances from each dataset (see the number of instances sampled in Table 11). We use GPT-3 with few-shot in-context examples and a prompt to generate step-by-step reasoning (e.g., "explain step-by-step") for each sampled instance (see in-context examples and prompts in App. B). Then, using our taxonomy we constructed a list of evaluation perspectives to label the model generated step-by-step reasoning step of each of these datasets. We explain the details of the perspectives used to label human judged datasets in § 5 and App. F. All datasets with examples are summarised in in Table 13. In the following we present details of each human judged datasets.

DROP (Dua et al., 2019), Discrete Reasoning Over the content of Paragraphs, is a dataset of 96K of instances with context and a question. To solve the tasks, a system must resolve references in the context that match with the question, and perform discrete operations over them (such as addition, counting, or sorting). These operations require comprehensive understanding of the content of the input context.

GSM8K (Cobbe et al., 2021) is a dataset of 8.5K linguistically diverse grade school math word problems. On this dataset, even the largest transformer models fail to achieve high test performance, despite the conceptual simplicity of this problem distribution.

CosmosQA (Huang et al., 2019) is a dataset of 35K problems that require commonsense-based reading comprehension, formulated as multiple-choice questions. The questions focus on reading between the lines over a diverse collection of people's everyday narratives, asking such questions as "what might be the possible reason of ...?", or "what would have happened if ...?". The dataset does not introduce step-by-step reasoning output, and contains multiple choice answers.

ESNLI (Camburu et al., 2018) is the extended version of the Stanford Natural Language Inference corpus (Bowman et al., 2015) of 570K labeled sentence pairs with entailment or contradiction labels. ESNLI includes human labeled explanations of the entailment decision.

SemEVAL (Ostermann et al., 2018) is a dataset on machine comprehension using commonsense knowledge. It contains questions that require commonsense knowledge for finding the correct answer.

E.3 SYNTHETIC DIAGNOSTICS DATASET GENERATION WITH PERTURBATION RULES

To construct the diagnostics datasets we apply synthetic perturbations on half of the chains from six datasets (for details see App. E.1 and the summary Table 11). Also, in Table 14 we illustrate these synthetic perturbations applied on reasoning steps $\{r_i\}$ of gold reference chains of all the datasets. In there, g^* indicates a grammar error, which includes changing verb tense, dropping verb, or random word swap. s^* represents change the semantics of one step in the chain by replacing named entities. To simulate extrinsic hallucinations, we use random steps from other chains within the same dataset.

Table 13: We show instances from five of the **Human Judged Datasets** used in our work. Only GSM8K and ESNLI include human labeled explanations.

Dataset	Instance	Reference Answer & Reference Step-by-Step Solution
GSM8K	Question: Tina buys 3 12-packs of soda for a party. Including Tina, 6 people are at the party. Half of the people at the party have 3 sodas each, 2 of the people have 4, and 1 person has 5. How many sodas are left over when the party is over?	Answer: 11 Step1: Tine buys 3 12-packs of soda for 3*12=36 sodas Step2: 6 people attend the party, so half of them is 6/2= 3 people Step3: Each of those people drinks 3 sodas, so they drink 3*3=9 sodas. Step4: Two people drink 4 sodas, which means they drink 2*4=8 sodas. Step5: With 1 person drinking 5, that brings the total drank to 5+9+8+3=25 sodas Step6: As Tina started off with 36 sodas, that means there are 36-25=11 sodas left.
CosmosQA	Context: A woman had topped herself by jumping off the roof of the hospital she had just recently been admitted to. She was there because the first or perhaps latest suicide attempt was unsuccessful. She put her clothes on, folded the hospital gown and made the bed. She walked through the unit unimpeded and took the elevator to the top floor Question: What would have happened to the woman if the staff at the hospital were doing their job properly?	Answer: The woman would have been stopped before she left to take the elevator to the top floor and she would have lived.
DROP	Context: Denver would retake the lead with kicker Matt Prater nailing a 43-yard field goal, yet Carolina answered as kicker John Kasay ties the game with a 39-yard field goal Carolina closed out the half with Kasay nailing a 44-yard field goal In the fourth quarter, Carolina sealed the win with Kasay's 42-yard field goal. Question: Which kicker kicked the most field goals?	Answer: John Kasay
ESNLI	Premise: A child in a yellow plastic safety swing is laughing as a dark- haired woman in pink and coral pants stands behind her. Hypothesis: A young mother is playing with her daughter in a swing.	Answer: neutral Explanation: Child does not imply daughter and woman does not imply mother.
SemEVAL	Context: Now I am going to set the dining table up for dinner. First I put away all the stuff that is not supposed to be on the table. Next I clean the table with a tissue paper. Then I arrange some of the decorations on the table. After that I put down the plates and glasses.lastly in the remaining spaces on the table I put down what we are going to eat but I dont put down dessert yet. There is one more thing I do before I am finished setting up the dining table. I call my family down for dinner. I swept a drop of sweat of my chin and forehead. It was hard hard work but still it was so much fun. Oh no my family is done with dinner and now I have to bring them dessert which is fruit I made myself. Question: When did they clean the dining table? Hypothesis: After it was set.	Answer: No

To construct diagnostic data from math datasets, we introduce four additional perturbations to simulate stepwise explanation errors that might arise in arithmetic reasoning task (*Arithmetic error*), general knowledge about relationships and equation construction (*Common sense error*), and misinformation about object/subject characteristics (*Factuality* or *Hallucination*):

- Shuffle numbers: randomly shuffles all numbers in the chain,
- Shuffle operations: randomly shuffles all math operations in the chain,
- Random number: randomly replaces one number in the chain,
- Random operation: randomly replaces one math operation in the chain.

Table 14: Synthetic perturbations and corresponding error types of steps $\{r_i\}$ in reference chains used when constructing diagnostics datasets. $\boldsymbol{g}^*(\cdot)$ represents grammar error, $\boldsymbol{s}^*(\cdot)$ represents semantic change.

Perturbation Type	Error Type	Reference Reasoning Steps	Hypothesis Reasoning Steps
Repeat a step	Repetition	$[r_1, r_2, r_3]$	$[r_1, r_2, r_2, r_3]$
Remove a step	Missing step	$[r_1, r_2, r_3]$	$[r_2, r_3]$
Shuffle steps	Self-coherency	$[r_1, r_2, r_3]$	$[r_3, r_1, r_2]$
Swap a step	Self-coherency	$[r_1, r_2, r_3]$	$[r_2, r_1, r_3]$
Negate a step	Factuality	$[r_1, r_2, r_3]$	$[r_1, \neg r_2, r_3]$
Hallucination	Hallucination	$[r_1, r_2, r_3]$	$[r_1, r_2, r_3, r_4]$
Grammar error	Grammatical	$[r_1, r_2, r_3]$	$[r_1, r_2, \boldsymbol{g^*}(r_3)]$
Semantic change	Factuality	$[r_1, r_2, r_3]$	$[r_1, s^*(r_2), r_3]$

F HUMAN ANNOTATIONS (CONT. FROM § 5)

To construct **Human Judged Datasets**, we perform human annotations on five datasets which we summarize in Table 11 (Type='Human judged'). These datasets do not include explanations (except GSM8K and ESNLI), so we construct model generated reasoning steps and label them with reasoning errors. We explain our generation process in §5 and App. E.2. We used five expert human annotators to collect reasoning error labels on five datasets. We asked human evaluators to directly rate the generated reasoning errors on overall chain level using a Likert scale from 1 to 5. We also asked them to mark whether each error type proposed in our error taxonomy (§3) appeared in each step in step-level evaluations. In Fig. 4 and Fig. 5 we illustrate the UI used to collect the data. Table 15 summarizes questions that experts were asked. Table 16 reports the distribution of errors for each dataset. In general, we found that it was hard to get anonymous crowd workers to annotate our data accurately even when we paid averages of upwards of \$30 an hour, hence relying on expert annotators. For the annotation sessions reported in the text of the paper, we find that it takes an average of 754 seconds for expert annotators to complete a session of at most 5 examples, or slightly over 2-and-a-half minutes per example. This highlights the difficulty of obtaining high-quality annotations on these cognitive challenging tasks.

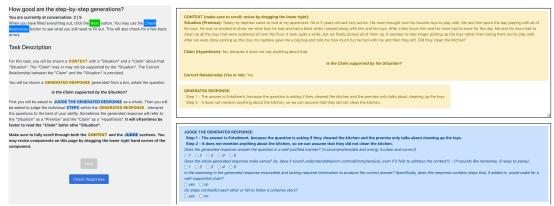


Figure 4: Screenshot of expert annotation user interface, showing the context for the initial question as well as the questions regarding the generated response.

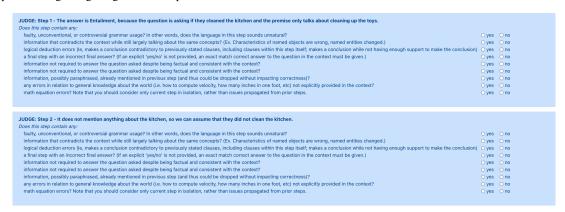


Figure 5: Screenshot of expert annotation user interface, showing questions asked for each step, using the question in Fig 4. The questions are asked of every step generated by the model, with steps separated by sentence-ending periods.

Table 15: Evaluation perspectives used to **Human Judged** the datasets. The perspectives, which we used to ask humans to label, align with our taxonomy of reasoning errors. (Continued from \S 5)

Level		Evaluation Perspective	Label	Details
Overall	QUAL	Overall quality	[1-5]	Does the generated response answer the question in a well-justified manner? (1=incomprehensible and wrong, 5=clear and correct)
Overall	СОН	Coherency	[1-5]	Does the whole generated response make sense? (Ie, does it sound understandable/non-contradictory/sensical, even if it fails to address the context?) - (1=sounds like nonsense, 5=easy to parse).
Step	MISS	Missing Step	Y/N	Is the reasoning in the generated response incomplete and lacking required information to produce the correct answer? Specifically, does this response contains steps that, if added in, would make for a well-supported chain?
Step	GRAM	Grammar	Y/N	Does this step contain faulty, unconventional, or controversial grammar usage? In other words, does the language in this step sounds unnatural?
Step	FACT	Factuality	Y/N	Does this step contain information that contradicts the context while still largely talking about the same concepts? (Ex. Characteristics of named objects are wrong, named entities changed.)
Step	LOGIC	Coherency and Logic	Y/N	Does this step any logical deduction errors (Ie, makes a conclusion contradictory to previously stated clauses, including clauses within this step itself; makes a conclusion while not having enough support to make the conclusion)
Step	HALL	Hallucination	Y/N	Does this step contain information not provided in the problem statement that is irrelevant or wrong?
Step	RED	Redundancy	Y/N	Does this step contain information not required to answer the question asked despite being factual and consistent with the context?
Step	REP	Repetition	Y/N	Does this step contain any information, possibly paraphrased, already mentioned in previous step (and thus could be dropped without impacting correctness)?
Step	COMMON	Commonsense	Y/N	Does this step contain any errors in relation to general knowledge about the world (i.e. how to compute velocity, how many inches in one foot, etc) not explicitly provided in the context?
Step	MATH	Arithmetic	Y/N	Does this step contain math equation errors? Note that you should consider only current step in isolation, rather than issues propagated from prior steps.

Table 16: Statistics of types of errors in **Human Judged** datasets. Each column reports the number of examples where the specified error type exists. (Continue from \S 5)

Error Type	DROP	GSM8K	ESNLI	COSMOS	SemEVAL
Grammar	8	4	5	8	6
Factuality	19	56	15	44	31
Hallucination	4	8	4	9	2
Redundancy	25	13	14	15	19
Repetition	2	2	0	3	3
Missing Step	109	81	40	99	67
Coherency	20	57	17	48	17
Commonsense	3	58	5	18	1
Arithmetic	2	7	1	0	0

G SENTENCE EMBEDDING MODEL TRAINING (CONT. FROM §6)

Model training. We use the train portions of the perturbed diagnostics datasets to finetune the SimCSE embeddings model (explained in § 5) and validation portions to select the best embedding model. The test portions are used to evaluate our metrics against baseline metrics. We randomly select 500,000 samples with replacement from each dataset to create uniform representation and reduce bias.

The hyperparameters used to finetune SimCSE model are described in Table 17. We use NVIDIA Tesla V100 Volta GPU instances with 32GB Graphics Card. We perform hyperparameter search, varying batch size in $\{32, 64, 256, 512, 1024, 2048\}$, learning rate in $\{5e\text{-}06, 1e\text{-}05, 5e\text{-}05, 1e\text{-}04\}$, and max sequence length in $\{64, 128, 512\}$. Not all combinations of batch size and max sequence length were explored due to memory limitations.

Table 17: Hyperparameters used to fine-tune SimCSE model on perturbed datasets.

Parameter	Value
Batch size	64
Max sequence length	512
Training epochs	5
Learning rate	5e-6
Temperature	0.05

Validation. We replace original validation procedure on semantic textual similarity tasks with similarity-based validation on perturbed reasoning chains. In particular, during training, we select best checkpoint that maximizes cosine similarity between positive and minimizes cosine similarity between hard-negative pairs within the batch of size B as the following:

$$\frac{\sum_{i=1}^{N} \left[\cos(s_i, r_i) - \cos(s_i, h_i) \right]}{2 * B}$$
 (3)

Model is evaluated every 100 steps on the development dataset and the best checkpoint is applied at the inference. Other parameters not described in this section are kept as in the original SimCSE model used for initialization.

Inference. We compare ROSCOE scores calculated against three embeddings: finetuned SimCSE model, *sup-simcse-roberta-base* SimCSE model, and *all-mpnet-base-v2* sentence embedding model (Reimers & Gurevych, 2019). During inference, we set the random seed to 42. Without this, the embedding-based scores naturally varied by about 0.01.

H ADDITIONAL EXPERIMENTAL RESULTS (CONT. FROM §6)

H.1 CONTROLLED EXPERIMENTS WITH DIAGNOSTICS DATASETS

In this section, we presented Somers' D correlation of all metrics on all Diagnostics datasets. Table 18 summarizes the evaluations when investigated reference-free. One of the characteristics of our ROSCOE metrics is that, they can provide judgement of the model generated reasoning steps with and without the human reference reasoning chains. In the experiments section in §6, we discussed the results of our unsupervised scores in comparison to baseline scores when measured reference-free. In Table 19, we summarize the correlation analysis on ROSCOE metrics in comparison to baselines on diagnostic datasets when reference is present for evaluation. Specifically, each score is measured between the human provided reasoning steps (reference) and the model generated reasoning steps (hypothesis). We also display fine-grained meta-evaluations of all metrics on each diagnostics dataset in separate tables. Specifically, Tables 20, 26 for EQASC, Tables 21, 27 for EntailmentBank, Tables 22, 28 for MATH, Tables 23, 29 for ProofWriter, Tables 24, 30 for ASDIV, and Tables 25, 31 for AQUA.

To understand if designed reference-free scores capture targeted error types we analyze perturbation-level correlations summarized in Fig. 6. Out of the all considered scores, *Info-Chain* is able to cover 10 out of 12 of errors, except *Remove Step* and *Semantic error* perturbations. In general we can note that ROSCOE fails to consistently identify *missing step* error type represented by *Remove Step* perturbation across different datasets, while other synthesized error types are covered by at least one score type.

Reference-based scores are covering all synthetic errors, with Semantic Coverage Chain showing strong correlations with all types of perturbations (Table 19). We also note that along with ROSCOE scores, the highest correlation among all reference-based scores belong to ROUGE and BERT scores (Tables 26-31). ROUGE scores consistently outperform on *Repetition*, *Hallucination*, *Remove Step*, *Shuffle Steps*, *Swap Steps*, *Negate Step*, and *Semantic* perturbations, while under performing on *Random operation*, and *Shuffle operations*. We attribute this to the fact that ROUGE is an n-gram based score, so it is better in catching errors were wording has significantly changed, while failing to catch small changes within steps.

It is worth noting that some scores, especially those among reference-based evaluations, get the highest possible Somers' D correlation scores of 1.0. What it means is that in some scenarios, there is a perfect correlation between the metric and the error type. In other words, for this metric we can find a threshold such generated chains that have scores greater than the threshold do not have errors of the given type, and in all generated chains with scores less than the threshold have that error. It is especially evident on referenced-based metrics that directly compare the reference solution and hypothesis. In this scenario, we build correlation for two groups: 1) non-perturbed hypothesis: the score is calculated by comparing embedding similarities of the reference with itself, and we expect to get high scores, 2) perturbed hypothesis: comparing reference with its perturbed version, where the scores should be lower. In some cases, we are able to perfectly separate perturbed and non-perturbed chains based on the corresponding metric values by selecting a threshold, in other cases we cannot due to a number of false-negatives (i.e., a chain gets a high score, although the error is present). As an example, consider the Semantic Coverage-Chain metric calculated on EQASC dataset using all-mpnet-base-v2 sentence embeddings, and Hallucination perturbation (Table 26). Here the Somers' D correlation score is 1.0. Semantic Coverage-Chain is calculated as a normalized cosine distance between the chain embedding of the reference solution r, and the chain embedding of the hypothesis $h: [1 + \cos(r, h)]/2$. Recall that in our setup, half of the hypothesis chains are perturbed reference chains, and another half is the same as the reference. While Hallucination perturbation is an insertion of a random step from a dataset, it is hard to predict how if will affect the embedding of the chain as a whole, but on the unperturbed chains, where h == r, the Semantic Coverage-Chain should be: $[1 + \cos(r, r)]/2 = 1.0$. Further review confirmed that in this dataset there are no false-positive instances, i.e., all chains with perturbations had Semantic Coverage-Chain score less than 1.0. That means, we can always identify if the chain contains a Hallucination error or not, by comparing *Semantic Coverage-Chain* value with 1.0 (threshold value), which is reflected in perfect Somers' *D* score.

Highest correlations among reference-free scores belong to the *Repetition-** scores, that exhibit perfect correlation on EQASC dataset (Tables 20-25). For other datasets, non-perfect correlations can be attributed to the small number of false-negatives, i.e. they give low *Repetition-** scores for chains with non-duplicated but similar steps, while all chains with duplicates got almost 0 scores (Fig. 7). In EQASC explanations are created from a set of facts that are not directly related to each other, but are intended to give an answer when combined together. Among all datasets considered, these steps are most dissimilar, and thus can be separated with similarity-based scores.

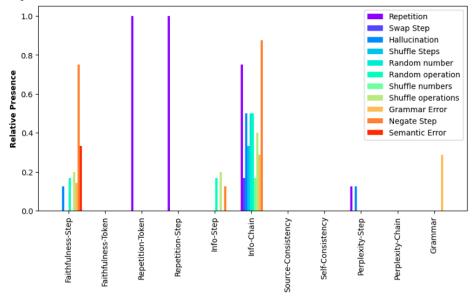


Figure 6: Relative presence of the strong score-perturbation correlation, measured as the number of datasets where for each score-perturbation pair Somers' D correlation value is in the 90^{th} percentile, normalized by the total number of datasets where this type of perturbation occurs. Statistics collected over ROSCOE reference-free scores with finetuned SimCSE embeddings. (Continued from $\S7$)

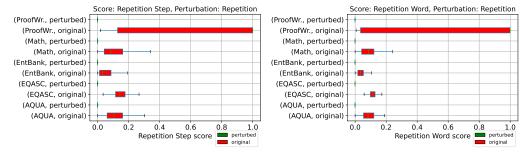


Figure 7: Box-and-whisker plots of interquartile ranges of scores, for *Repetition* perturbations and *Repetition-** scores. While all perturbed subsets have 0 or near 0 scores, all datasets except EQASC have some chains that were also scored as low despite the absence of duplicates.

Table 18: Somers' D correlation of **all** metrics on six **Diagnostics** datasets. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlations that are not significant (p-value) >= 0.05) are omitted when aggregating, and "-" denotes an absence of any significant correlations. Note that ASDIV is a 1-step equation dataset, so there are no repetition and self-consistency scores as there are no steps to compare. (Continued from §6, more details in App. H.1.).

0.095 0.090 0.106 0.024 0.117 0.081 0.054 0.011 0.235 0.394 0.513 ddings 0.771 0.436 0.250 0.596 ddings	0.342 0.217 0.315 0.447 0.462 0.415 0.297 0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.305 - 0.269 - 0.322
0.106 0.024 0.117 0.081 0.054 0.011 0.235 0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.315 0.447 0.462 0.415 0.297 0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.322 - - 0.382 0.396 0.186 0.182 0.198
0.024 0.117 0.081 0.054 0.011 0.235 0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.447 0.462 0.415 0.297 0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.322 - - 0.382 0.396 0.186 0.182 0.198
0.117 0.081 0.054 0.011 0.235 0.394 0.513 Iddings 0.771 0.436 0.250 0.596 Iddings 0.082	0.462 0.415 0.297 0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.322 - - 0.382 0.396 0.186 0.182 0.198
0.081 0.054 0.011 0.235 0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.415 0.297 0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.382 0.396 0.186 0.182 0.198
0.054 0.011 0.235 0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.297 0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.396 0.186 0.182 0.198
0.011 0.235 0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.276 0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.396 0.186 0.182 0.198
0.235 0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.436 0.485 0.270 0.785 0.480 0.538 1.000	0.396 0.186 0.182 0.198
0.394 0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.485 0.270 0.785 0.480 0.538 1.000	0.396 0.186 0.182 0.198
0.513 Idings 0.771 0.436 0.250 0.596 Idings 0.082	0.270 0.785 0.480 0.538 1.000	0.396 0.186 0.182 0.198
ldings 0.771 0.436 0.250 0.596 ldings 0.082	0.785 0.480 0.538 1.000	0.186 0.182 0.198
0.771 0.436 0.250 0.596 Idings 0.082	0.480 0.538 1.000	0.182 0.198
0.771 0.436 0.250 0.596 Idings 0.082	0.480 0.538 1.000	0.182 0.198
0.436 0.250 0.596 Idings 0.082	0.480 0.538 1.000	0.182 0.198
0.250 0.596 Idings 0.082	0.538 1.000	0.198
0.596 Idings 0.082	1.000	
ldings 0.082		n/a
0.082	0.550	
	0.550	
	0.550	0.280
0.631	1.000	n/a
embeddings		
0.732	0.630	0.266
0.517	0.502	0.515
0.226	0.399	0.275
0.606	1.000	n/a
embeddings		
0.045	0.409	0.289
0.621	1.000	n/a
ence embeddii	ngs	
0.763	0.907	0.879
0.499	0.492	0.740
0.317	0.804	0.879
0.606	1.000	n/a
ence embeddii	ngs	
0.462	0.995	0.857
0.624	1.000	n/a
0.204	0.793	n/a
0.289	0.778	0.771
0.178	0.394	0.485
0.135	0.379	0.485
0.788	0.859	0.470
-	0.631 embeddings 0.732 0.517 0.226 0.606 embeddings 0.045 0.621 nce embeddin 0.763 0.499 0.317 0.606 nce embeddin 0.462 0.624 0.204 0.289	0.631 1.000 embeddings 0.732 0.630 0.517 0.502 0.226 0.399 0.606 1.000 embeddings 0.045 0.409 0.621 1.000 embeddings 0.763 0.907 0.499 0.492 0.317 0.804 0.606 1.000 embeddings 0.462 0.995 0.624 1.000 0.204 0.793 0.289 0.778 0.178 0.394 0.135 0.379

Table 19: Somers' D correlation of all **reference-based** metrics on six **Diagnostics** datasets. Metrics are measured using **reference** generations on (r, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. (Continued from §6, more details in App. H.1.)

Metric	EntBank	Math	AQUA	ProofWriter	EQASC	ASDIV
ROUGE-1	1.000	1.000	1.000	1.000	1.000	1.000
ROUGE-2	1.000	1.000	1.000	1.000	1.000	1.000
ROUGE-L	1.000	1.000	1.000	1.000	1.000	1.000
BLEURT	0.821	0.773	0.946	0.829	1.000	0.93
BERTScore	1.000	1.000	1.000	1.000	1.000	1.000
BARTScore	0.983	0.989	0.984	0.986	1.000	0.954
BARTScore+	0.988	0.963	1.000	0.996	1.000	1.000
BARTScore-P	0.877	0.799	0.905	0.595	0.966	0.83
PRISM	0.939	0.521	1.000	0.997	0.996	1.000
CTC Relevancy	0.457	0.592	0.409	0.725	0.954	0.398
CTC Consistency	0.814	0.804	0.833	0.635	0.974	0.6
ROSCOE Metrics (reference-	based metr	ics only	7)			
ROSCOE-SA with all-mpnet-						
Hallucination	1.000	1.000	1.000	1.000	1.000	1.000
Redundancy	1.000	1.000	1.000	1.000	1.000	1.000
Semantic Coverage-Step	1.000	1.000	1.000	1.000	1.000	1.000
Reasoning Alignment	1.000	1.000	0.143	1.000	1.000	1.000
Commonsense	0.438	1.000	1.000	0.379	1.000	1.000
Missing Step	0.993	1.000	1.000	0.876	1.000	1.000
ROSCOE-SS with all-mpnet-						
Semantic Coverage-Chain	1.000	1.000	1.000	1.000	1.000	1.000
ROSCOE-SA with sup-simcse	-roberta-ba	se sente	nce emb	eddings		
Hallucination	1.000	1.000	1.000	1.000	1.000	1.000
Redundancy	1.000	1.000	1.000	1.000	1.000	1.000
Semantic Coverage-Step	1.000	1.000	1.000	1.000	1.000	1.000
Reasoning Alignment	1.000	1.000	1.000	1.000	1.000	1.000
Commonsense	0.433	1.000	1.000	0.415	1.000	1.000
Missing Step	0.999	1.000	1.000	0.874	1.000	1.000
ROSCOE-SS with sup-simcse	-roberta-ba	se sente	nce emb	eddings		
Semantic Coverage-Chain	1.000	<u>0.999</u>	1.000	1.000	1.000	1.000
ROSCOE-SA with finetuned s	up-simcse-r	obertas	sentence	embeddings		
Hallucination	1.000	1.000	1.000	1.000	1.000	1.000
Redundancy	1.000	1.000	1.000	1.000	1.000	1.000
Semantic Coverage-Step	1.000	1.000	1.000	1.000	1.000	1.000
Reasoning Alignment	1.000	1.000	1.000	1.000	1.000	1.000
Commonsense	0.445	1.000	1.000	0.404	1.000	1.000
Missing Step	0.999	1.000	1.000	0.873	1.000	1.000
ROSCOE-SS with finetuned s	up-simcse-r	obertas	sentence	embeddings		
Semantic Coverage-Chain	1.000	0.999	1.000	1.000	1.000	1.000

Table 20: Somers' D correlations of **all** metrics per different perturbation applied on **EQASC** Diagnostics datasets. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Halluc.	Grammar Error	Remove Step	Negate Step	Semantic Error
Rouge-1	0.264†	0.342†	0.017	0.227†	0.063	-0.023
Rouge-2	0.071	0.205†	0.106†	0.217†	0.099†	-0.002
Rouge-L	0.210†	0.315†	0.057	0.179†	0.083†	-0.016
BLEURT	0.366†	0.447†	-0.028	0.195†	0.204†	-0.108†
BERTScore	0.288†	0.462†	0.153†	0.160†	0.052	0.051
BARTScore	-0.127†	0.038	0.047	0.415†	0.019	-0.072
BARTScore+	0.028	$0.212\dagger$	0.055	0.297†	0.023	-0.028
BARTScore-P	-0.039	0.031	-0.038	$0.276\dagger$	-0.003	-0.023
PRISM	-0.327†	$0.436\dagger$	0.267†	$0.077\dagger$	0.010	0.123†
CTC-Relevancy	$0.141\dagger$	0.001	0.082†	$0.485\dagger$	0.002	0.220†
CTC-Consistency	0.001	-0.080†	$0.095\dagger$	-0.154†	0.078†	0.270†
ROSCOE Metrics (refe						
ROSCOE-SA with all-n						
Faithfulness-Step	-0.006	0.785†	0.017	-0.040	$0.084\dagger$	-0.243†
Faithfulness-Token	-0.031	0.480†	-0.001	0.037	0.008	-0.156†
Info-Step	0.006	0.538†	-0.003	0.223†	0.085†	-0.191†
Repetition-Token	1.000†	0.399†	-0.028	-1.000†	0.070	0.074
ROSCOE-SS with all-n						
Info-Chain	0.189†	0.550†	-0.022	0.220†	0.059	-0.132†
Repetition-Step	1.000†	0.035	-0.034	-1.000†	-0.120†	0.030
ROSCOE-SA with sup-					0.2001	0.2041
Faithfulness-Step	-0.061	0.630†	-0.010	0.007	0.308†	-0.204†
Faithfulness-Token	-0.031	0.502†	0.032	0.045	0.107†	-0.168†
Info-Step	-0.064	0.399†	-0.024	0.241†	0.296†	-0.155†
Repetition-Token	1.000†	0.148†	-0.138†	-1.000†	-0.055	-0.080†
ROSCOE-SS with sup-					0.2704	0.1624
	-0.025	0.409†	-0.028	0.256†	0.379†	-0.163†
Repetition-Step	1.000†	0.001	-0.003	-1.000†	-0.465†	0.071
ROSCOE-SA with finet	unea <i>sup</i> -			-0.046		
Faithfulness-Step Faithfulness-Token	-0.044	0.630†	0.293† 0.141†	0.036	0.907† 0.492†	-0.118†
	-0.019	0.485†	0.141			-0.128†
Info-Step Repetition-Token	-0.041 1.000 †	0.383† 0.148†	-0.138†	0.196† -1.000†	0.804† -0.055	-0.068 -0.080†
ROSCOE-SS with finet						
Info-Chain	0.995†	0.871†	0.588†	0.103†	0.967†	0.121†
Repetition-Step	1.000†	$\frac{0.871}{0.037}$	-0.386†	-1.000†	-0.953†	-0.048
ROSCOE-LI	1.000	0.057	0.5001	1.000	0.755	0.010
Source-Consistency	-0.020	0.576†	0.112†	-0.239†	0.778†	0.268†
Self-Consistency	0.020	0.633†	0.399†	-0.713†	0.773†	0.476†
ROSCOE-LC	0.022	3.033	0.577	0.713	3.775	9.1701
Perplexity-Chain	-0.690†	-0.007	0.379†	0.260†	0.118†	0.269†
Perplexity-Step	0.937†	0.965†	0.352†	-0.953†	0.081†	0.225†
Grammar	-0.025	0.060	0.859†	-0.145†	0.139†	0.722†
	0.020	0.000	0.000	0.1.01	31.207	****=1

Table 21: Somers' D correlations of **all** metrics per different perturbation applied on **Entailment Bank** Diagnostics datasets. All metrics are measured **reference-free** on (s,h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from §6, more details in App. H.1.)

$\overline{\textbf{Perturbations}} \rightarrow$	Repet.	Halluc.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Semantic Error
Rouge-1	0.391†	0.410†	0.074	0.081	0.113	0.065	0.034	-0.017
Rouge-2	0.356†	0.391†	0.160†	0.116	0.109	0.088	0.109	0.091
Rouge-L	0.209†	0.194†	0.112	0.138†	0.365†	0.229†	0.003	0.063
BLEURT	0.025	0.164†	0.060	0.231†	0.092	0.096	$0.257\dagger$	-0.036
BERTScore	$0.264\dagger$	0.380†	0.211†	0.150†	$0.364\dagger$	0.205†	0.050	0.087
BARTScore	0.034	0.063	0.106	0.358†	0.248†	0.183†	$0.142\dagger$	0.164†
BARTScore+	0.101	0.047	0.036	0.315†	0.184†	0.155†	0.130†	0.173†
BARTScore-P	0.061	0.012	0.025	0.186†	0.041	-0.008	0.098	0.100
PRISM	$0.230\dagger$	$0.453\dagger$	0.279†	0.118	$0.331\dagger$	0.167†	0.081	0.148†
CTC-Relevancy	0.258†	0.057	-0.026	0.080	-0.020	0.035	0.029	0.159†
CTC-Consistency	0.310†	0.159†	-0.058	-0.249†	-0.162†	-0.023	-0.046	0.194†
ROSCOE Metrics (refe								
ROSCOE-SA with all-n								
Faithfulness-Step	0.023	0.786†	0.231†	-0.005	0.111	-0.000	0.098	0.121
Faithfulness-Token	0.098	0.581†	0.250†	0.087	0.009	-0.020	0.179†	0.170†
Info-Step	0.083	0.638†	0.181†	0.161†	0.025	-0.001	0.216†	0.139†
Repetition-Token	<u>0.913</u> †	0.105	0.042	-0.177†	0.139†	0.038	-0.085	-0.058
ROSCOE-SA with all-n								
Faithfulness-Step	0.023	0.786†	0.231†	-0.005	0.111	-0.000	0.098	0.121
Faithfulness-Token	0.098	0.581†	0.250†	0.087	0.009	-0.020	0.179†	0.170†
Info-Step	0.083	0.638†	0.181†	0.161†	0.025	-0.001	0.216†	0.139†
Repetition-Token	0.913†	0.105	0.042	-0.177†	0.139†	0.038	-0.085	-0.058
ROSCOE-SS with all-n								
Info-Chain	0.039	0.419†	0.083	0.071	0.068	0.025	0.040	0.037
Repetition-Step	0.909†	0.148†	-0.061	-0.165†	0.067	-0.068	-0.062	-0.056
ROSCOE-SA with sup-								
Faithfulness-Step	-0.002	0.777†	0.158†	-0.015	0.107	0.009	0.358†	0.110
Faithfulness-Token	0.063	0.663†	0.280†	0.084	0.046	0.024	0.279†	0.159†
Info-Step	0.066	0.560†	0.079	0.148†	-0.035	0.007	0.450†	0.117
Repetition-Token	0.919†	0.137†	-0.010	-0.216†	0.171†	0.005	-0.106	-0.099
ROSCOE-SS with sup-					0.010	0.022	0.5041	0.040
Info-Chain	0.084	0.515†	0.005	0.119	-0.019	0.023	0.524†	0.048
Repetition-Step	0.901†	0.139†	0.006	-0.188†	0.019	-0.066	-0.144†	-0.098
ROSCOE-SA with finet							0.5001	
Faithfulness-Step	0.105	0.392†	0.204†	0.063	0.093	0.029	0.538†	-0.115
Faithfulness-Token	0.079	0.519†	0.271†	0.100	0.055	0.029	0.356†	0.041
Info-Step	0.131†	0.364†	0.227†	0.161†	0.053	0.047	0.599†	-0.018
Repetition-Token	0.919†	0.137†	-0.010	-0.216†	0.171†	0.005	-0.106	-0.099
ROSCOE-SS with finet					_		0.055	0.2001
Info-Chain	0.871†	0.851†	0.752†	0.437†	0.937†	0.672†	0.955†	0.300†
Repetition-Step	0.908†	0.133†	-0.013	-0.169†	0.135†	-0.043	-0.058	-0.079
ROSCOE-LI	0.044	0.2001	0.2104	0.052	0.1501	0.075	0.9601	0.257
Source-Consistency	-0.044	0.289†	0.218†	0.052	0.150†	-0.075	0.860†	0.257†
Self-Consistency ROSCOE-LC	-0.040	0.403†	0.216†	-0.042	0.129†	-0.053	0.782†	0.170†
Perplexity-Chain	-0.364†	0.104	0.116	0.151†	-0.006	0.085	0.060	0.133†
Perplexity-Chain Perplexity-Step	-0.364† 0.199†	0.104	0.116	-0.097	0.143†	-0.021	-0.019	-0.046
Grammar	-0.109	0.1081	0.604†	0.075	0.1431	-0.021	0.019	0.365†
Granniai	-0.109	0.070	0.00+1	0.073	0.044	-0.003	0.055	0.505

Table 22: Somers' D correlations of **all** metrics per different perturbation applied on **MATH** Diagnostics datasets. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Halluc.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Random Number	Random Operation	Shuffle Numbers	Shuffle Operations
Rouge-1	0.176†	0.151†	-0.004	0.065†	-0.020	0.020	0.028	0.017	0.008	0.022	0.011
Rouge-2	0.107†	0.095†	0.006	0.121†	-0.005	0.029	0.010	0.018	-0.011	0.151†	-0.002
Rouge-L	0.126†	0.156†	0.004	0.075†	0.008	$0.042\dagger$	0.023	0.018	0.008	0.082†	0.004
BLEURT	0.143†	$0.148\dagger$	0.000	-0.023	0.001	0.012	$0.049\dagger$	-0.000	-0.005	-0.036†	0.002
BERTScore	$0.124\dagger$	$0.117\dagger$	0.067†	$0.089\dagger$	0.025	0.029	0.010	0.029	0.016	0.034†	0.017
BARTScore	-0.048†	-0.066†	0.015	0.185†	-0.026	-0.002	-0.030	0.029	0.009	0.075†	0.045†
BARTScore+	0.015	-0.003	0.049†	0.162†	0.003	0.031	0.002	$0.050\dagger$	0.047†	0.164†	0.063†
BARTScore-P	-0.005	0.002	0.022	$0.128\dagger$	0.006	0.011	-0.009	0.024	0.002	0.115†	0.059†
PRISM	-0.115†	$0.208\dagger$	0.120†	$0.095\dagger$	0.029	0.017	-0.003	$0.102\dagger$	0.069†	$0.117\dagger$	0.111†
CTC-Relevancy	$0.104\dagger$	0.041†	0.133†	0.188^{+}	0.027	0.029	0.018	$0.043\dagger$	0.052†	0.052†	0.001
CTC-Consistency	-0.106†	$0.145\dagger$	0.282†	$0.105\dagger$	0.096†	0.035†	0.080 †	0.079†	0.033	0.046†	-0.017
ROSCOE Metrics (refer											
ROSCOE-SA with all-n	1			-							
Faithfulness-Step	0.029	0.362†	0.016	$0.070\dagger$	-0.024	0.009	0.025	$0.060\dagger$	0.048†	$0.074\dagger$	0.060†
Faithfulness-Token	0.000	0.157†	0.000	0.004	-0.004	-0.003	0.018	0.028	0.005	-0.022	0.003
Info-Step	0.024	0.231†	0.022	0.118†	-0.018	0.009	$0.033\dagger$	0.085†	$0.074\dagger$	0.109†	0.106†
Repetition-Token	0.936†	0.069†	-0.018	-0.078†	0.010	0.043†	0.006	-0.041†	-0.018	0.035†	0.006
ROSCOE-SS with all-n											
Info-Chain	0.059†	<u>0.467</u> †	0.016	0.172†	<u>0.106</u> †	0.091†	0.063†	0.192†	0.089†	0.162†	0.135†
Repetition-Step	0.932†	-0.002	-0.036†	-0.114†	-0.001	0.013	-0.026	-0.134†	-0.045†	-0.008	-0.030
ROSCOE-SA with sup-											
Faithfulness-Step	-0.004	0.229†	0.049†	0.046†	0.007	0.010	0.142†	0.073†	0.033	0.031	-0.010
Faithfulness-Token	0.014	0.200†	0.091†	-0.004	-0.001	0.009	0.021	0.033†	0.005	0.021	0.008
Info-Step	-0.020	0.115†	0.048†	0.092†	0.009	0.020	0.131†	0.086†	0.029	0.042†	-0.015
Repetition-Token	0.939†	0.007	-0.093†	-0.073†	0.023	0.037†	0.002	-0.058†	-0.035†	0.019	-0.002
ROSCOE-SS with sup-s					0.0241	0.004	0.1001	0.0701	0.0511	0.0411	0.007
Info-Chain	0.028	0.114†	0.038†	0.071†	0.034†	0.004	0.180†	0.078†	0.051†	0.041†	-0.007
Repetition-Step	0.949†	0.019	-0.043†	-0.094†	-0.012	0.000	-0.062†	-0.119†	0.014	0.027	-0.001
ROSCOE-SA with finet							0.6141	0.0761	0.4151	0.2551	0.4201
Faithfulness-Step	0.021	0.223†	0.190†	0.038†	-0.031	0.005	0.614	0.376†	0.415†	0.277†	0.428†
Faithfulness-Token	0.016	0.204†	0.106†	-0.005	-0.004	0.008	0.084†	0.067†	0.045†	0.047†	0.052†
Info-Step	0.012	0.133†	0.228†	0.099†	-0.021	0.020	0.511†	0.451†	0.452†	0.301†	0.430†
Repetition-Token ROSCOE-SS with fineton	0.939†	0.007	-0.093†	-0.073†	0.023	0.037†	0.002	-0.058†	-0.035†	0.019	-0.002
Info-Chain	unea <i>sup-</i> 0.679†	0.588†	0.694 †	0.216†	10eddings 0.746†	0.530†	0.777†	0.698†	0.757†	0.524†	0.662†
Repetition-Step	0.079	0.024	-0.059†	-0.083†	-0.019	0.012	-0.052†	-0.192†	-0.103†	-0.011	0.002
ROSCOE-LI	0.54	0.024	-0.035	-0.003	-0.019	0.012	-0.032	-0.134	-0.105	-0.011	0.017
Source-Consistency	0.011	0.071†	0.044†	-0.041†	-0.002	0.026	0.215†	0.223†	0.123†	0.341†	0.133†
Self-Consistency	0.011	0.069†	-0.003	-0.105†	-0.002	0.020	0.213†	0.223	0.123	0.190†	0.133†
ROSCOE-LC	0.013	0.005	-0.003	-0.103	-0.014	0.011	0.122	0.147	0.0001	0.170	0.0001
Perplexity-Chain	-0.358†	0.020	0.175†	0.103†	0.100†	0.035†	0.003	0.173†	0.109†	0.154†	0.170†
Perplexity-Step	0.160†	0.020	0.001	-0.093†	-0.016	0.033	0.003	0.173	0.012	0.134	0.037
Grammar	0.100	0.026	0.392†	-0.020	-0.010	0.019	0.000	0.003	0.012	0.048	0.037
Grammar	0.010	0.020	0.3721	.0.020	0.012	5.005	0.114	0.054	0.020	0.037	0.040

Table 23: Somers' D correlations of **all** metrics per different perturbation applied on **ProofWriter** Diagnostics datasets. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Halluc.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Semantic Error
Rouge-1	-0.124†	-0.054†	0.023†	0.095†	0.013	-0.001	0.007	0.006
Rouge-2	-0.093†	-0.002	0.068†	0.090†	0.018	-0.004	0.037†	0.047†
Rouge-L	-0.120†	-0.028†	0.029†	0.089†	0.106†	0.020	0.020	0.027†
BLEURT	-0.058†	-0.001	-0.027†	-0.099†	0.016	0.007	$0.024\dagger$	-0.016
BERTScore	-0.049†	$0.077\dagger$	0.117†	0.082†	-0.064†	-0.054†	0.023†	0.108†
BARTScore	-0.059†	-0.096†	-0.037†	0.081†	-0.013	-0.010	-0.015	-0.061†
BARTScore+	-0.055†	-0.067†	0.014	0.054†	-0.011	-0.021	0.006	0.032†
BARTScore-P	-0.046†	-0.049†	-0.012	0.010	-0.044†	-0.032†	-0.020	0.011
PRISM	-0.159†	0.159†	0.222†	$0.097\dagger$	0.017	-0.010	$0.060\dagger$	0.235†
CTC-Relevancy	0.394†	0.392†	0.123†	0.185†	-0.131†	-0.052†	0.036†	0.077†
CTC-Consistency	0.496†	0.513†	0.182†	0.223	-0.063†	-0.022	0.098†	0.131†
ROSCOE Metrics (refe								
ROSCOE-SA with all-n	npnet-bas	e-v2 sente	nce embeddi					
Faithfulness-Step	0.002	0.771†	0.348†	$0.165\dagger$	0.011	0.013	$0.233\dagger$	<u>0.515</u> †
Faithfulness-Token	0.004	0.436†	0.264†	$0.055\dagger$	0.030†	0.006	$0.168\dagger$	0.310†
Info-Step	-0.004	0.250†	0.121†	-0.062†	0.023	0.015	$0.108\dagger$	0.174†
Repetition-Token	0.596†	0.053†	-0.041†	0.101†	-0.006	-0.003	-0.055†	-0.050†
ROSCOE-SS with all-n								
Info-Chain	-0.052†	0.083†	0.033†	-0.089†	0.014	0.005	-0.001	0.001
Repetition-Step	0.631†	0.031†	-0.027†	0.116†	-0.002	0.002	-0.044†	-0.042†
ROSCOE-SA with sup-					0.042		0.4041	0.4041
Faithfulness-Step	0.004	0.732†	0.171†	0.154†	0.013	0.012	0.481	0.494†
Faithfulness-Token	0.010	0.517†	0.334†	0.086†	0.031†	0.005	0.336†	0.395†
Info-Step	-0.008	0.226†	0.047†	-0.063†	0.027†	0.014	0.172†	0.160†
Repetition-Token	0.606†	0.036†	-0.065†	0.097†	-0.002	0.009	-0.070†	-0.070†
ROSCOE-SS with sup-					0.0241	0.011	0.0451	0.006
Info-Chain	-0.040†	0.004	-0.015	-0.140†	-0.024†	-0.011	0.045†	-0.006
Repetition-Step	0.621†	0.028†	-0.008	0.115†	-0.011	0.001	-0.050†	-0.043†
ROSCOE-SA with finet							0.7364	0.5264
Faithfulness-Step Faithfulness-Token	0.008 0.012	0.763† 0.499†	0.618† 0.475†	0.180†	0.008	0.004	0.726†	0.536†
		,		0.088†	0.036†		0.436†	0.403†
Info-Step Repetition-Token	-0.018 0.606†	0.243† 0.036†	0.252† -0.065†	-0.046† 0.097†	0.016	0.012 0.009	0.317† -0.070†	0.187† -0.070†
ROSCOE-SS with finet							-0.0701	-0.0701
Info-Chain	0.214†	0.248†	0.284†	0.077†	0.330†	0.246†	0.462†	0.122†
Repetition-Step	0.624†	0.039†	-0.034†	0.077	-0.007	-0.002	-0.076†	-0.034†
ROSCOE-LI			'	<u> </u>			'	· ·
Source-Consistency	-0.008	0.027†	0.028†	0.010	-0.044†	-0.004	0.289†	-0.049†
Self-Consistency	0.011	0.204†	0.084†	0.110†	-0.022	-0.020	0.036†	0.065†
ROSCOE-LC								
Perplexity-Chain	-0.165†	0.047†	0.112†	-0.064†	0.135†	0.067†	0.012	0.128†
Perplexity-Step	0.178†	0.112†	0.033†	0.082†	-0.008	0.005	-0.036†	-0.008
Grammar	0.000	0.042†	0.788†	0.102†	0.007	-0.023	0.007	0.515†
-				'			-	

Table 24: Somers' D correlations of **all** metrics per different perturbation applied on **ASDIV** Diagnostics datasets. All metrics are measured **reference-free** on (s,h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Random Number	Random Operation	Shuffle Numbers	Shuffle Operations
Rouge-1	0.305†	0.096	0.085	-0.220
Rouge-2	-0.034	-0.014	0.007	-0.038
Rouge-L	0.245†	0.073	0.269†	-0.235
BLEURT	0.043	-0.034	0.059	0.015
BERTScore	0.098	0.322†	0.025	0.167
BARTScore	0.107	-0.027	-0.105	0.197
BARTScore+	-0.011	-0.002	-0.075	-0.015
BARTScore-P	0.068	0.043	-0.048	0.121
PRISM	-0.009	-0.035	-0.114	0.258
CTC-Relevancy	0.382†	0.155†	0.038	0.000
CTC-Consistency	0.396†	0.189†	0.121	-0.121
ROSCOE Metrics (refe	rence-free	on (s,h)		
ROSCOE-SA with all-n	npnet-base-	v2 sentence e	embeddings	
Faithfulness-Step	0.186†	-0.090	0.091	0.091
Faithfulness-Token	0.182†	0.080	0.062	-0.091
Info-Step	0.198†	-0.091	0.085	0.167
ROSCOE-SS with all-n	npnet-base-	v2 sentence	embeddings	
Info-Chain	$0.280\dagger$	0.005	0.192†	0.091
ROSCOE-SA with sup-	simcse-robe	<i>erta-base</i> sen	tence embed	dings
Faithfulness-Step	0.266†	0.082	-0.182†	0.015
Faithfulness-Token	0.273†	0.011	-0.125	0.515†
Info-Step	$0.275\dagger$	0.125	-0.141†	0.000
ROSCOE-SS with sup-	simcse-robe	<i>erta-base</i> sen	tence embed	dings
Info-Chain	0.289†	0.145†	-0.084	0.030
ROSCOE-SA with finet	uned sup-si	mcse-roberto	a-base senter	nce embeddings
Faithfulness-Step	$0.630\dagger$	0.840†	0.670†	0.879†
Faithfulness-Token	0.576†	0.740†	0.552†	0.545†
Info-Step	0.669†	0.844^{\dagger}	0.683†	0.879†
ROSCOE-SS with finet	uned sup-si	mcse-roberto	<i>i-base</i> senter	nce embeddings
Info-Chain	0.773†	0.857†	0.795†	<u>0.803</u> †
ROSCOE-LI				
Source-Consistency	<u>0.760</u> †	0.771†	0.763†	0.500†
Self-Consistency	0.203†	0.227†	0.206†	0.152
ROSCOE-LC				
Perplexity-Chain	0.300†	0.092	0.214†	0.485†
Perplexity-Step	0.300†	0.092	0.214†	0.485†
Grammar	0.170†	-0.083	-0.007	0.470†

Table 25: Somers' D correlations of **all** metrics per different perturbation applied on **AQUA** Diagnostics datasets. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

${\bf Perturbations} \rightarrow$	Repet.	Hallu.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Random Number	Random Operation	Shuffle Numbers	Shuffle Operations
Rouge-1	0.134	0.055	-0.085	0.257†	0.074	-0.032	0.109	0.056	-0.006	0.131	-0.073
Rouge-2	$0.180\dagger$	0.004	-0.021	0.206†	0.100	-0.092	0.044	0.051	0.082	0.196†	-0.061
Rouge-L	0.072	0.055	-0.059	0.264†	0.145†	-0.008	0.120	0.050	0.004	0.148†	-0.037
BLEURT	0.010	0.047	-0.066	0.252†	-0.003	-0.047	0.071	-0.010	0.020	0.165†	-0.023
BERTScore	0.134	0.111	-0.016	0.166†	-0.011	-0.043	0.028	0.007	0.079	0.220†	-0.036
BARTScore	-0.030	-0.269†	-0.025	0.317†	0.102	0.051	-0.014	0.060	0.011	0.169†	0.036
BARTScore+	0.112	-0.037	0.044	0.251†	0.019	-0.053	0.028	0.072	0.066	0.230†	-0.009
BARTScore-P	0.082	-0.097	-0.007	0.215†	0.013	-0.027	0.114	0.059	0.042	0.182†	0.028
PRISM	-0.119	0.026	0.177†	0.188†	0.071	0.048	0.057	0.191†	0.131	0.182†	0.028
CTC-Relevancy	0.133	-0.016	0.065	0.096	-0.068	0.013	0.024	0.119	-0.017	0.217†	0.072
CTC-Consistency	0.071	-0.050	0.075	0.041	-0.074	0.012	0.004	0.106	-0.006	0.157†	0.024
ROSCOE Metrics (refer											
ROSCOE-SA with all-n		e-v2 sente	ence embedd								
Faithfulness-Step	0.098	$0.152\dagger$	-0.013	0.106	0.035	-0.026	0.080	0.057	-0.023	0.026	-0.123
Faithfulness-Token	0.125	0.009	0.008	$0.157\dagger$	0.047	-0.001	0.057	-0.014	0.076	0.070	-0.006
Info-Step	0.088	0.030	-0.020	0.114	0.068	-0.007	0.015	0.063	0.025	0.068	-0.088
Repetition-Token	$0.972\dagger$	0.216†	-0.125	-0.226†	-0.053	-0.061	-0.041	-0.004	-0.024	0.037	-0.001
ROSCOE-SS with all-n	ıpnet-bas	e-v2 sente	ence embedd	ings							
Info-Chain	0.066	$0.214\dagger$	-0.053	0.146†	0.023	-0.035	0.049	0.079	0.038	0.153†	-0.081
Repetition-Step	$0.982\dagger$	$0.143\dagger$	-0.093	-0.169†	0.006	0.002	-0.031	-0.068	-0.028	0.010	0.050
ROSCOE-SA with sup-	simcse-ro	berta-bas	e sentence er	nbeddings							
Faithfulness-Step	0.089	$0.179\dagger$	-0.043	0.135	-0.016	-0.009	$0.224\dagger$	0.122	0.096	0.014	-0.156
Faithfulness-Token	0.080	0.062	-0.047	0.117	0.036	-0.011	0.057	0.035	0.024	0.049	-0.059
Info-Step	0.099	0.064	-0.012	$0.162\dagger$	0.018	-0.023	$0.183\dagger$	0.138	0.116	0.099	-0.147
Repetition-Token	$0.971\dagger$	0.072	-0.112	-0.190†	-0.062	-0.054	-0.102	-0.034	0.061	0.026	0.079
ROSCOE-SS with sup-s	simcse-ro	berta-bas	e sentence er	nbeddings							
Info-Chain	0.130	0.046	-0.025	$0.178\dagger$	-0.013	-0.024	$0.153\dagger$	0.108	0.066	0.195†	-0.119
Repetition-Step	$0.991\dagger$	0.125	-0.120	-0.234†	0.007	-0.023	-0.063	-0.155†	-0.025	-0.010	0.076
ROSCOE-SA with finet	uned sup-	simcse-ro	berta-base s	entence en	beddings						
Faithfulness-Step	0.169†	0.165†	0.280†	0.000	0.018	-0.113	$0.826 \dagger$	0.394†	0.329†	0.210†	0.236†
Faithfulness-Token	0.063	0.132	0.014	0.084	0.030	-0.045	0.285†	0.079	0.077	0.082	-0.037
Info-Step	$0.150\dagger$	0.038	0.203†	0.107	0.041	-0.081	0.703†	0.450†	0.327†	0.313†	0.259†
Repetition-Token	0.971†	0.072	-0.112	-0.190†	-0.062	-0.054	-0.102	-0.034	0.061	0.026	0.079
ROSCOE-SS with finet	uned sup-	simcse-ro	berta-base s	entence en	beddings						
Info-Chain	0.861†	$0.465\dagger$	0.399†	0.025	0.463	0.315†	0.933^{\dagger}	0.589†	0.499†	0.395†	0.337†
Repetition-Step	0.982†	0.141†	-0.171†	-0.149†	-0.094	0.022	-0.143	-0.339†	-0.155†	-0.016	-0.028
ROSCOE-LI											
Source-Consistency	-0.044	0.140	-0.061	-0.171†	-0.087	0.024	0.316†	$0.212\dagger$	0.151†	<u>0.419</u> †	0.096
Self-Consistency	0.041	0.368†	-0.028	-0.227†	0.019	-0.040	0.104	0.226†	0.106	0.167†	0.093
ROSCOE-LC											
Perplexity-Chain	-0.288†	-0.120	0.179†	0.121	0.021	0.101	0.134	$0.190\dagger$	0.071	$0.194\dagger$	0.229†
Perplexity-Step	$0.171\dagger$	$0.182\dagger$	-0.051	-0.105	-0.013	-0.097	-0.068	0.060	0.011	-0.004	0.079
Grammar	0.091	0.200†	0.359†	0.020	-0.043	-0.008	$0.223\dagger$	0.094	0.046	0.086	0.171

Table 26: Somers' D correlations of **all** metrics per different perturbation applied on **EQASC** Diagnostics datasets. All metrics are measured **reference-based** on (s,h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from §6, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Hallu.	Grammar Error	Remove Step	Negate Step	Semantic Error
Rouge-1	1.000†	1.000†	0.622†	1.000†	1.000†	0.946†
Rouge-2	1.000†	1.000†	0.998†	1.000†	1.000†	0.946†
Rouge-L	1.000†	1.000	0.998†	1.000†	1.000	0.946†
BLEURT	0.829†	0.999†	0.512†	1.000†	0.983†	0.703†
BERTScore	1.000†	1.000†	1.000†	1.000†	1.000†	0.943†
BARTScore	0.784†	0.925†	0.888†	1.000†	0.861†	0.897†
BARTScore+	0.443†	0.725†	0.530†	1.000†	0.597†	0.885†
BARTScore-P	-0.095†	0.284†	0.404†	0.966†	0.476†	0.587†
PRISM	0.879†	0.995†	0.902†	0.996†	0.927†	0.932†
CTC-Relevancy	0.031	-0.124†	0.265†	0.954†	0.321†	0.496†
CTC-Consistency	<u>0.974</u> †	0.965†	0.393†	-0.429†	0.711†	0.718†
ROSCOE Metrics (reference-						
ROSCOE-SA with all-mpnet-l						
Reasoning Alignment	0.010	0.995†	<u>0.999</u> †	0.008	1.000^{+}	0.947†
Hallucination	-0.000	0.994†	0.998†	-0.305†	$1.000\dagger$	0.949†
Redundancy	-0.004	0.994†	0.998†	-0.278†	$1.000\dagger$	0.949†
Commonsense	-0.006	0.012	0.998†	$1.000\dagger$	$1.000\dagger$	0.949†
Missing Step	-0.010	0.012	0.998†	$1.000\dagger$	$1.000\dagger$	0.949†
Semantic Coverage-Step	$1.000\dagger$	$1.000\dagger$	0.997†	$1.000\dagger$	$1.000\dagger$	0.941†
ROSCOE-SS with all-mpnet-l						
Semantic Coverage-Chain	1.000†	$1.000\dagger$	<u>0.999</u> †	1.000†	1.000†	0.944†
ROSCOE-SA with sup-simcse						
Reasoning Alignment	0.035	0.995†	1.000†	0.003	1.000†	<u>0.950</u> †
Hallucination	0.006	0.994†	1.000†	-0.302†	1.000†	0.949†
Redundancy	0.007	0.994†	$1.000\dagger$	-0.273†	$1.000\dagger$	0.949†
Commonsense	0.003	0.025	1.000†	1.000^{+}	$1.000\dagger$	0.949†
Missing Step	0.003	0.028	1.000†	1.000^{+}	$1.000\dagger$	0.949†
Semantic Coverage-Step	1.000†	1.000†	1.000†	1.000†	1.000†	0.948†
ROSCOE-SS with sup-simcse				_		
Semantic Coverage-Chain	1.000†	1.000†	1.000†	1.000†	$1.000\dagger$	0.946†
ROSCOE-SA with finetuned s	up-simcse		base sentenc			
Reasoning Alignment	0.019	0.995†	$1.000\dagger$	0.027	$1.000\dagger$	0.949†
Hallucination	0.009	0.996†	$1.000\dagger$	-0.273†	$1.000\dagger$	0.945†
Redundancy	0.008	0.996†	$1.000\dagger$	-0.248†	$1.000\dagger$	0.946†
Commonsense	0.007	-0.017	1.000†	$1.000\dagger$	$1.000\dagger$	0.945†
Missing Step	0.006	-0.014	1.000†	$1.000\dagger$	$1.000\dagger$	0.946†
Semantic Coverage-Step	1.000†	1.000†	1.000†	1.000†	$1.000\dagger$	0.946†
ROSCOE-SS with finetuned s						
Semantic Coverage-Chain	1.000†	1.000†	1.000†	1.000†	1.000†	0.956†

Table 27: Somers' D correlations of **all** metrics per different perturbation applied on **Entailment Bank** Diagnostics datasets. All metrics are measured **reference-based** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from §6, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Hallu.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Semantic Error
Rouge-1	1.000†	1.000†	0.582†	1.000†	0.000†	0.000†	1.000†	1.000†
Rouge-2	1.000†	1.000†	1.000†	1.000†	0.982†	0.935†	1.000†	1.000†
Rouge-L	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†
BLEURT	0.197†	0.821†	0.163†	0.786†	0.253†	0.021	0.640†	0.174†
BERTScore	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†
BARTScore	0.574†	0.659†	0.574†	0.802†	0.983†	0.864†	0.479†	0.656†
BARTScore+	0.195†	0.430†	0.214†	0.839†	0.988†	0.813†	0.253†	0.555†
BARTScore-P	0.159†	0.272†	0.337†	0.633†	0.877†	0.656†	0.355†	0.415†
PRISM	0.707†	0.850†	0.612†	0.781†	0.939†	0.833†	0.466†	0.660†
CTC-Relevancy	0.311†	0.250†	-0.039	0.457†	0.266†	0.214†	0.036	0.254†
CTC-Consistency	0.768†	0.814†	0.133†	-0.022	0.479†	0.413†	0.233†	0.474†
ROSCOE Metrics (reference-	based or	(s,h)						
ROSCOE-SA with all-mpnet-	base-v2 s	entence o	embeddings					
Reasoning Alignment		$1.000\dagger$	$1.000\dagger$	0.056	-0.034	-0.001	1.000†	0.990†
Hallucination	0.002	$1.000\dagger$	$1.000\dagger$	-0.200†	0.041	0.017	1.000†	0.995†
Redundancy	-0.072	1.000†	$1.000\dagger$	-0.115†	0.034	0.055	1.000†	0.993†
Commonsense	0.025	-0.016	0.438†	$0.264\dagger$	0.025	-0.012	$0.237\dagger$	0.327†
Missing Step	-0.077	-0.038	0.986†	$0.993\dagger$	0.040	0.060	0.976†	0.985†
Semantic Coverage-Step	1.000†	$1.000\dagger$	1.000†	$1.000\dagger$	0.003	0.036	$1.000\dagger$	0.999†
ROSCOE-SS with all-mpnet-	base-v2 s	entence o	embeddings					
Semantic Coverage-Chain			1.000†	1.000†	$1.000\dagger$	1.000†	1.000†	1.000†
ROSCOE-SA with sup-simcse								
Reasoning Alignment	-0.024	1.000†	1.000†	0.039	0.035	0.004	1.000†	0.990†
Hallucination	-0.011	1.000†	1.000†	-0.059	-0.042	-0.008	1.000†	0.990†
Redundancy	-0.011	1.000†	$1.000\dagger$	-0.023	0.034	0.034	1.000†	$0.994\dagger$
Commonsense	-0.021	0.090	0.433†	0.347†	-0.049	0.002	0.276†	0.342†
Missing Step	-0.011	0.052	<u>0.988</u> †	<u>0.999</u> †	0.034	0.034	<u>0.976</u> †	0.983†
Semantic Coverage-Step	1.000†	1.000†	1.000†	1.000†	0.106†	0.040	1.000†	0.994†
ROSCOE-SS with sup-simcse								
Semantic Coverage-Chain			0.922†	0.951†	1.000†	0.968†	0.940†	0.941†
ROSCOE-SA with finetuned s								
Reasoning Alignment	0.040	1.000†	1.000†	-0.013	0.007	-0.080	1.000†	0.999†
Hallucination	-0.012	1.000†	1.000†	-0.059	0.134†	0.040	1.000†	0.991†
Redundancy	0.025	1.000†	1.000†	-0.074	0.041†	-0.017	1.000†	0.994†
Commonsense	-0.022	-0.039	0.368†	0.324†	0.129†	0.039	0.304†	0.445†
Missing Step	0.025	0.008	0.988†	<u>0.999</u> †	0.041†	-0.017	0.971†	0.982†
Semantic Coverage-Step	1.000†	1.000†	1.000†	1.000†	0.117†	0.038	1.000†	0.994†
ROSCOE-SS with finetuned s								
Semantic Coverage-Chain	0.912†	0.934†	0.911†	0.950†	1.000†	0.980†	0.940†	0.931†

Table 28: Somers' D correlations of **all** metrics per different perturbation applied on **MATH** Diagnostics datasets. All metrics are measured **reference-based** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Hallu.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Random Number	Random Operation	Shuffle Numbers	Shuffle Operations
Rouge-1	0.994†	0.961†	0.451†	0.996†	0.000†	0.000†	1.000†	1.000†	0.000†	0.433†	0.000†
Rouge-2	0.994†	0.961†	0.767†	0.996†	0.990†	0.953†	1.000†	1.000†	0.000†	0.999†	0.000†
Rouge-L	0.994	0.961†	0.769†	0.996†	0.998†	0.997†	1.000†	1.000†	0.000†	1.000†	0.000†
BLEURT	0.363†	0.561†	0.041†	0.773†	0.103†	0.059†	0.156†	0.121†	0.069†	0.028	0.020
BERTScore	0.972†	0.942†	0.978†	0.978†	1.000†	0.993†	0.964	1.000†	0.999†	0.998†	0.998†
BARTScore	0.413†	0.531†	0.812†	0.837†	0.966†	0.803†	0.429†	0.938†	0.921†	0.989†	0.932†
BARTScore+	0.022	0.208†	0.539†	0.829†	0.611†	0.437†	0.153†	0.804†	0.466†	0.963†	0.606†
BARTScore-P	-0.084†	0.041†	0.248†	0.474†	0.369†	0.170†	0.015	0.529†	0.367†	0.799†	0.536†
PRISM	0.337†	0.465†	0.386†	0.506†	0.412†	0.392†	0.266†	0.437†	0.347†	0.521†	0.285†
CTC-Relevancy	-0.045†	0.168†	0.295†	0.592†	0.253†	$0.110\dagger$	$0.141\dagger$	0.140†	0.048†	0.169†	0.041
CTC-Consistency	0.673†	0.804†	0.491†	0.091†	0.393†	0.240†	0.506†	0.372†	0.112†	0.408†	0.148†
ROSCOE Metrics (reference-											
ROSCOE-SA with all-mpnet-b											
Reasoning Alignment	-0.012		0.835†	0.005	-0.017	-0.013	$1.000\dagger$	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$
Hallucination	-0.007	<u>0.999</u> †	0.854†	-0.038†	-0.019	-0.007	$1.000\dagger$	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$
Redundancy	-0.007	<u>0.999</u> †	0.853†	-0.040†	-0.019	-0.005	$1.000\dagger$	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$
Commonsense	-0.030†	-0.004	0.852†	0.991†	-0.022	-0.005	<u>0.999</u> †	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$
Missing Step	-0.030†	-0.003	0.850†	0.991†	-0.021	-0.004	<u>0.999</u> †	1.000^{+}	1.000†	1.000†	$1.000\dagger$
Semantic Coverage-Step	0.908†	1.000†	0.852†	1.000†	-0.005	-0.008	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$	1.000†
ROSCOE-SS with all-mpnet-b											
Semantic Coverage-Chain			0.799†	0.944†	1.000†	0.980†	0.933†	1.000†	<u>0.999</u> †	0.998†	0.992†
ROSCOE-SA with sup-simcse											
Reasoning Alignment	-0.002	0.998†	0.990†	-0.004	-0.003	0.018	0.998†	0.997†	0.997†	1.000†	$1.000\dagger$
Hallucination	-0.006	<u>0.999</u> †	0.993 †	-0.045†	$0.012\dagger$	0.003	0.998†	0.996†	0.997†	1.000†	$1.000\dagger$
Redundancy	-0.006	<u>0.999</u> †	0.993 †	-0.044†	$0.012\dagger$	0.003	$0.998\dagger$	0.996†	0.997†	1.000†	$1.000\dagger$
Commonsense	-0.006	0.002	0.983†	$0.989\dagger$	0.011	0.003	0.995†	0.996†	0.997†	1.000†	$1.000\dagger$
Missing Step	-0.006	0.001	0.983†	0.989†	0.011	0.003	0.995†	0.996†	0.997†	$1.000\dagger$	$1.000\dagger$
Semantic Coverage-Step	0.914†	1.000†	0.993†	1.000†	0.038†	0.017†	0.998†	0.997†	0.997†	1.000†	1.000†
ROSCOE-SS with sup-simcse											
Semantic Coverage-Chain		0.773†	0.762†	0.763†	0.999†	0.868†	0.743†	0.985†	0.958†	0.975†	0.886†
ROSCOE-SA with finetuned s						0.025	0.0001	0.0001	0.0001	1 000!	1 0001
Reasoning Alignment	0.011	0.999†	0.992†	0.019	-0.001	0.025	0.998†	0.998†	0.998†	1.000†	1.000†
Hallucination	-0.008	0.999†	0.993†	-0.032†	0.011	0.015†	0.998†	0.997†	0.997†	1.000†	1.000†
Redundancy	-0.008	0.999†	0.993†	-0.031†	0.010	0.016†	0.998†	0.997†	0.997†	1.000†	1.000†
Commonsense	-0.008	0.000	0.983†	0.991†	0.010	0.015†	0.995†	0.997†	0.997†	1.000†	1.000†
Missing Step	-0.008	0.001	0.983†	0.990†	0.010	0.016†	0.995†	0.997†	0.997†	1.000†	1.000†
Semantic Coverage-Step ROSCOE-SS with finetuned so	0.913†	1.000†	0.993†	1.000†	0.036†	0.016†	0.998†	0.997†	0.997†	1.000†	1.000†
Semantic Coverage-Chain		0.770†	0.769†	0.775†	0.999†	0.865†	0.741†	0.985†	0.957†	0.973†	0.897†

Table 29: Somers' D correlations of **all** metrics per different perturbation applied on **ProofWriter** Diagnostics datasets. All metrics are measured **reference-based** on (s,h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from §6, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Hallu.	Grammar Error	Remove Step	Shuffle Steps	Swap Steps	Negate Step	Semantic Error
Rouge-1	1.000†	1.000†	0.663†	1.000†	0.000†	0.000†	1.000†	1.000†
Rouge-2	1.000†	1.000†	0.993†	1.000†	0.932†	0.812†	1.000†	1.000†
Rouge-L	1.000†	1.000	0.993†	1.000†	1.000†	1.000†	1.000	1.000†
BLEURT	0.349†	0.829†	0.221†	0.695†	0.076†	0.042†	0.788†	0.597†
BERTScore	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†
BARTScore	0.334†	0.430†	0.303†	0.417†	0.986†	0.748†	0.343†	0.391†
BARTScore+	0.245†	0.391†	0.251†	0.598†	0.996†	0.775†	0.278†	0.515†
BARTScore-P	0.020	0.100†	0.124†	0.260†	0.595†	0.318†	0.114†	0.260†
PRISM	0.829†	0.956†	0.822†	0.924†	0.997†	0.970†	0.871†	0.947†
CTC-Relevancy	0.376†	0.409†	0.354†	0.725†	0.621†	0.396†	0.391†	0.419†
CTC-Consistency	0.537†	0.635†	0.199†	0.009	0.608†	0.412†	0.404†	0.376†
ROSCOE Metrics (reference-								
ROSCOE-SA with all-mpnet-l								
Reasoning Alignment	-0.017		0.994†	-0.049†	-0.007	-0.000	1.000†	0.990†
Hallucination	0.015	0.997†	0.990†	0.012	-0.001	-0.004	1.000†	0.939†
Redundancy	0.007	0.997†	0.994†	$0.048\dagger$	-0.000	0.010	1.000†	0.992†
Commonsense	0.012	0.000	0.187†	0.379†	-0.002	-0.010	$0.180\dagger$	0.169†
Missing Step	0.003	-0.001	0.872†	$0.843\dagger$	0.004	0.011	0.876†	0.870†
Semantic Coverage-Step	0.793†	1.000†	0.994†	1.000†	-0.008	0.021†	1.000†	0.984†
ROSCOE-SS with all-mpnet-l								
Semantic Coverage-Chain		1.000†	0.994†	1.000†	1.000†	1.000†	1.000†	1.000†
ROSCOE-SA with sup-simcse					0.000	0.001	1 0001	1 0001
Reasoning Alignment	0.012	0.998†	1.000†	-0.041†	0.000	0.001	1.000†	1.000†
Hallucination	0.016	0.998†	0.998†	0.031†	-0.005	0.002	1.000†	1.000†
Redundancy	0.010	0.998†	1.000†	0.037†	0.015†	0.003	1.000†	1.000†
Commonsense	0.017	0.006	0.242†	0.415†	-0.005	0.002	0.238†	0.234†
Missing Step	0.010	0.004	0.871†	0.839†	0.016†	0.004	0.874†	0.866†
Semantic Coverage-Step	0.818†	1.000†	1.000†	1.000†	0.059†	0.021†	1.000†	1.000†
ROSCOE-SS with sup-simcse Semantic Coverage-Chain			0.995†	0.997†	1.000†	0.999†	0.998†	0.997†
ROSCOE-SA with finetuned s						0.9991	0.9981	0.997
Reasoning Alignment	0.016	0.996†	1.000†	-0.030†	0.000	0.005	1.000†	1.000†
Hallucination	0.010	0.997	0.998†	0.016	0.005	-0.010	1.000†	1.000†
Redundancy	0.003	0.997†	1.000†	0.050†	0.003	0.003	1.000†	1.000†
Commonsense	0.020	-0.009	0.243†	0.404†	0.019	-0.009	0.253†	0.235†
Missing Step	0.001	-0.009	0.243	0.404	0.010	0.003	0.233†	0.255†
Semantic Coverage-Step	0.817†	1.000†	0.871	1.000†	0.020†	0.003	1.000†	1.000†
ROSCOE-SS with finetuned s						0.022	1.000	1.000
Semantic Coverage-Chain			0.995†	0.996†	1.000†	<u>0.999</u> †	<u>0.998</u> †	0.996†

Table 30: Somers' D correlations of **all** metrics per different perturbation applied on **ASDIV** Diagnostics datasets. All metrics are measured **reference-based** on (s, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Random Number	Random Operation	Shuffle Numbers	Shuffle Operations
Rouge-1	1.000†	0.000†	0.027†	0.000†
Rouge-2	1.000†	0.000†	1.000+	0.000†
Rouge-L	1.000†	0.000†	1.000	0.000†
BLEURT	0.876†	0.930†	0.125	-0.212
BERTScore	1.000†	1.000†	1.000†	1.000†
BARTScore	0.871†	0.823†	0.954†	0.924†
BARTScore+	0.948†	0.839†	1.000†	0.955†
BARTScore-P	0.738†	0.642†	0.830†	0.712†
PRISM	0.998†	0.989†	1.000†	1.000†
CTC-Relevancy	0.398†	0.188†	0.160†	-0.136
CTC-Consistency	0.600†	0.131†	0.294†	-0.106
ROSCOE Metrics (reference-				
ROSCOE-SA with all-mpnet-				4 0001
Reasoning Alignment	1.000†	1.000†	1.000†	1.000†
Hallucination	1.000†	1.000†	1.000†	1.000†
Redundancy	1.000†	1.000†	1.000†	1.000†
Commonsense	1.000†	1.000†	1.000†	1.000†
Missing Step	1.000†	1.000†	1.000†	1.000†
Semantic Coverage-Step	1.000†	1.000†	1.000†	1.000†
ROSCOE-SS with all-mpnet-l				4 0001
Semantic Coverage-Chain ROSCOE-SA with sup-simcse	1.000†	1.000†	1.000†	1.000†
Reasoning Alignment	1.000†	1.000†	1.000†	1.000†
Hallucination	1.000† 1.000†	1.000†	1.000† 1.000†	1.000† 1.000†
Redundancy	1.000† 1.000†	1.000†	1.000† 1.000†	1.000† 1.000†
Commonsense	1.000† 1.000†	1.000†	1.000†	1.000†
Missing Step	1.000† 1.000†	1.000†	1.000† 1.000†	1.000† 1.000†
Semantic Coverage-Step	1.000† 1.000†	1.000†	1.000† 1.000†	1.000† 1.000†
ROSCOE-SS with sup-simcse				1.000
Semantic Coverage-Chain	1.000†	1.000†	1.000†	1.000†
ROSCOE-SA with finetuned s				
Reasoning Alignment	1.000†	1.000†	1.000†	1.000†
Hallucination	1.000†	1.000†	1.000+	1.000†
Redundancy	1.000†	1.000†	1.000†	1.000†
Commonsense	1.000†	1.000†	1.000†	1.000†
Missing Step	1.000†	1.000†	1.000†	1.000†
Semantic Coverage-Step	1.000†	1.000†	1.000†	1.000†
ROSCOE-SS with finetuned s	up-simcse-	roberta-base	sentence en	nbeddings
Semantic Coverage-Chain	1.000†	1.000†	1.000†	1.000†

Table 31: Somers' D correlations of **all** metrics per different perturbation applied on **AQUA** Diagnostics datasets. All metrics are measured **reference-based** on (r, h). The highest correlation overall for each dataset is in **bold**. The second best models are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from $\S 6$, more details in App. H.1.)

$\textbf{Perturbations} \rightarrow$	Repet.	Halluc.	Grammar error	Remove step	Shuffle steps	Swap steps	Negate step	Random number	Random operation	Shuffle numbers	Shuffle operations
Rouge-1	1.000†	1.000†	0.394†	1.000†	0.000†	0.000†	1.000†	1.000†	0.000†	0.000†	0.000†
Rouge-2	1.000†	1.000	0.866†	1.000†	0.984†	0.967†	1.000†	1.000†	0.000†	1.000	0.000†
Rouge-L	1.000†	1.000	0.866†	1.000†	1.000	1.000	1.000†	1.000†	0.000†	1.000	0.000†
BLEURT	0.193†	0.803†	0.129	0.946†	0.076	0.028	0.640†	0.627†	0.385†	-0.041	-0.113
BERTScore	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†	1.000†
BARTScore	0.568†	0.596†	0.688†	0.870†	0.958†	0.883†	0.520†	0.930†	0.834†	0.984†	0.849†
BARTScore+	0.283†	0.597†	0.613†	0.953†	0.960†	0.901†	0.368†	0.980†	0.521†	0.995†	0.814†
BARTScore-P	0.045	0.025	0.400†	0.679†	0.650†	0.547†	0.148	0.846†	0.559†	0.905†	0.726†
PRISM	0.815†	0.900†	0.724†	0.888†	0.941†	0.901†	0.689 †	0.942†	0.772†	0.970†	0.870†
CTC-Relevancy	0.054	-0.044	0.102	0.409†	0.116	0.141	0.062	0.290†	0.030	0.409†	0.023
CTC-Consistency	0.833†	0.766†	0.106	-0.150†	0.270†	0.257†	0.256†	0.774†	0.232†	0.741†	0.063
ROSCOE Metrics (reference-											
ROSCOE-SA with all-mpnet-											
Reasoning Alignment	0.030	0.997†	0.985†	-0.014	0.054	-0.076	1.000†	1.000†	1.000†	1.000†	1.000†
Hallucination	0.029	$0.992\dagger$	0.992†	-0.084	0.024	-0.030	$1.000\dagger$	$1.000\dagger$	1.000†	$1.000\dagger$	$1.000\dagger$
Redundancy	0.032	$0.992\dagger$	0.992†	-0.088	0.023	-0.040	$1.000\dagger$	$1.000\dagger$	1.000†	$1.000\dagger$	$1.000\dagger$
Commonsense	0.029	0.020	0.992†	1.000^{+}	0.024	-0.038	$1.000\dagger$	1.000^{\dagger}	$1.000\dagger$	$1.000\dagger$	$1.000\dagger$
Missing Step	0.032	0.023	0.992†	1.000^{+}	0.023	-0.047	$1.000\dagger$	1.000^{\dagger}	$1.000\dagger$	$1.000\dagger$	$1.000\dagger$
Semantic Coverage-Step	0.947†	1.000†	0.987†	1.000†	0.006	-0.009	$1.000\dagger$	$1.000\dagger$	1.000†	1.000†	1.000†
ROSCOE-SS with all-mpnet-											
Semantic Coverage-Chain		1.000†	<u>0.995</u> †	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$	1.000†	1.000†	1.000†	1.000†
ROSCOE-SA with sup-simcse											
Reasoning Alignment	-0.061	$1.000\dagger$	$1.000\dagger$	0.018	0.049	0.006	$1.000\dagger$	$1.000\dagger$	1.000†	$1.000\dagger$	$1.000\dagger$
Hallucination	-0.053	$1.000\dagger$	$1.000\dagger$	-0.069	-0.017	-0.041	$1.000\dagger$	$1.000\dagger$	1.000†	$1.000\dagger$	$1.000\dagger$
Redundancy	-0.054	$1.000\dagger$	$1.000\dagger$	-0.063	-0.016	-0.041	$1.000\dagger$	$1.000\dagger$	1.000†	$1.000\dagger$	$1.000\dagger$
Commonsense	-0.053	0.065	0.992†	1.000 †	-0.017	-0.041	$1.000\dagger$	1.000^{+}	1.000†	$1.000\dagger$	$1.000\dagger$
Missing Step	-0.054	0.064	0.992†	1.000 †	-0.016	-0.041	$1.000\dagger$	1.000^{+}	1.000†	$1.000\dagger$	$1.000\dagger$
Semantic Coverage-Step	0.966†	1.000†	1.000†	1.000†	0.041†	0.008	1.000†	1.000†	1.000†	1.000†	1.000†
ROSCOE-SS with sup-simcse											
Semantic Coverage-Chain			0.944†	0.988†	1.000†	1.000†	0.953†	1.000†	1.000†	1.000†	0.973†
ROSCOE-SA with finetuned s	1				_	0.00.	4 0001	4 0001	4 0001	4 0001	4 0001
Reasoning Alignment	-0.089	1.000†	1.000†	0.085	0.054	0.036	1.000†	1.000†	1.000†	1.000†	1.000†
Hallucination	-0.039	1.000†	1.000†	0.003	-0.024	-0.025	1.000†	1.000†	1.000†	1.000†	1.000†
Redundancy	-0.043	1.000†	1.000†	0.006	-0.016	-0.016	1.000†	1.000†	1.000†	1.000†	1.000†
Commonsense	-0.039	0.083	0.992†	1.000†	-0.024	-0.025	1.000†	1.000†	1.000†	1.000†	1.000†
Missing Step	-0.043	0.089	0.992†	1.000†	-0.016	-0.016	$1.000\dagger$	1.000†	1.000†	$1.000\dagger$	1.000†
Semantic Coverage-Step	0.967†	1.000†	1.000†	1.000†	0.041†	0.008	1.000†	1.000†	1.000†	1.000†	1.000†
ROSCOE-SS with finetuned s	1				_	1 000 !	0.071	1.000.	1.000	1 0002	0.0724
Semantic Coverage-Chain	0.95/†	0.957†	0.959†	0.980†	$1.000\dagger$	1.000†	<u>0.971</u> †	1.000^{+}	$1.000\dagger$	$1.000\dagger$	0.972†

H.2 EXPERIMENTS WITH HUMAN JUDGEMENT DATASETS

In this section, we present Somers' *D* correlation of all metrics on all Human Judged datasets in separate tables. Specifically, Table 32 summarizes meta-evaluations for ROSCOE metrics in comparison to baselines on **all** human judged datasets. Fine-grained evaluations are presented in Table 33 for DROP, Table 34, 38 for GSM8K, Table 35, 39 for ESNLI, Table 36 for CosmosQA, and Table 37 for SemEVAL. Human evaluation perspectives used in evaluations are described in App. Table 15.

Looking at how errors are captured by ROSCOE reference-free scores (Fig. 8), we observe strongest correlations between *Redundancy* error and *Repetition-**, *Self-Consistency* scores. *Repetition* error is not present in this analysis as it has at most 3 occurrences per dataset. Out of the all considered scores, *Self-Consistency* is able to cover 6 out of 7 evaluation perspectives, except *Missing Step*.

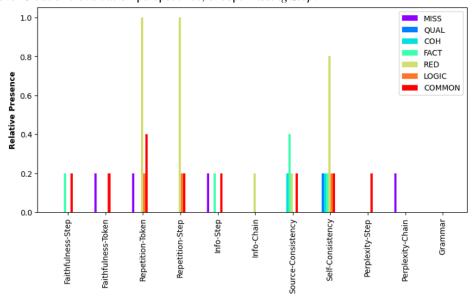


Figure 8: Relative presence of the strong score-error correlation, measured as the number of datasets where for each score and evaluation perspective pair Somers' D correlation value is in the 90^{th} percentile, normalized by the total number of datasets where this type of perturbation occurs. Statistics collected over ROSCOE reference-free scores with finetuned SimCSE embeddings, and evaluation perspectives where at least 10 errors are present in a dataset.

We further look at specific human annotated examples where our ROSCOE gives highest and lowest scores to understand strength and weaknesses of the proposed approach. Results are summarized in Table 40. Similar analysis for diagnostic datasets is summarized in Table 41.

Table 32: Somers' D correlations of **all** metrics on five **human judged** datasets. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each dataset is in **bold** and second best is <u>underlined</u>. Correlations that are not significant $(p\text{-}value \geq 0.05)$ are omitted when aggregating, and "-" denotes an absence of any significant correlations. (Continued from §6, more details in App. H.2.)

	DROP	GSM8K	ESNLI	COSMOS	SemEVAL
Rouge-1	0.239	0.180	0.559	-0.264	-0.520
Rouge-2	0.320	-	0.502	0.180	-
Rouge-L	0.278	0.252	0.557	-0.441	-0.478
BLEURT	0.328	0.256	0.541	0.218	-0.356
BERTScore	0.275	0.235	0.590	-0.420	-0.295
BARTScore	-0.835	-0.546	0.549	-0.544	-
BARTScore+	-0.665	-	0.482	-0.186	-
BARTScore-P	-0.642	-	0.255	-0.207	-
PRISM	-0.733	-0.455	0.580	-0.376	-
CTC-Relevance	0.333	-0.371	0.334	-	-0.349
CTC-Consistency	0.462	-0.174	0.647	0.275	-0.301
ROSCOE Metrics (ref					
ROSCOE-SA with fine		p-simcse-r		ase sentence	embeddings
Faithfulness-Step	0.496	-	0.403		-
Faithfulness-Token	0.417	- 0.170	0.521	-0.320	-
Info-Step	0.500	0.178	0.493		- 0.227
Repetition-Token	0.578	0.392	0.441	0.555	0.337
ROSCOE-SA with all	mpnet-bo- 0.297	<i>ise-v2</i> sent			0.220
Faithfulness-Step		0.442	0.423	-0.424	0.330
Faithfulness-Token	0.290	-0.443	0.524	-0.515	0.186
Info-Step	0.301		0.542	-0.429	- 0.405
Repetition-Token	0.790	0.500	0.799	0.638	0.485
ROSCOE-SA with sup					gs
Faithfulness-Step	0.477	-0.192	0.502	-0.381	-
Faithfulness-Token	0.454	-	0.540	-0.420	0.221
Info-Step	0.510 0.578	0.392	0.599 0.441	-0.409 0.555	-0.321 0.337
Repetition-Token					
ROSCOE-SS with fine				ase sentence	embeddings
Info-Chain	0.446	0.385	0.310		-
Repetition-Step	0.824	0.514	0.530	0.593	0.411
ROSCOE-SS with all		<i>ise-v2</i> sent			0.265
Info-Chain	0.406	0.471	0.507	-0.198	0.367
Repetition-Step	0.791	0.471	0.487	0.642	0.508
ROSCOE-SS with sup		oberta-ba			gs
Info-Chain	0.271		0.531	-0.367	-
Repetition-Step	0.799	0.638	0.484	0.658	0.535
ROSCOE-LI					
Source-Consistency	0.390	0.172	0.425	0.444	-
Self-Consistency	0.584	0.345	0.531	0.417	0.372
ROSCOE-LC					
Grammar	-	-0.184	0.255	-	0.517
Perplexity-Step	0.205	-0.307	0.345	-	-
Perplexity-Chain	-0.611	-0.273	0.447	-0.212	-0.373

Table 33: Somers' D correlation of all metrics on **DROP human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold**, second best are <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from \S 6, more details in App. H.2)

	QUAL	СОН	COMMON	FACT	HALL	RED	REP	LOGIC	MATH	GRAM	MISS
Rouge-1	0.157†	-0.160	0.436	$0.239\dagger$	-0.335	-0.731†	-0.736	0.030	0.702	0.196	0.173†
Rouge-2	0.137†	-0.146	0.488	0.320†	-0.284	-0.716†	-0.442	0.027	0.584	-0.035	0.129
Rouge-L	0.131†	-0.201†	0.465	0.278†	-0.345	-0.749†	-0.815	-0.012	0.745	0.027	0.146
BLEURT	0.121†	-0.101	0.256	0.328†	-0.333	-0.725†	-0.370	0.078	0.514	0.101	0.087
BERTScore	0.133†	-0.115	0.494	0.275†	-0.177	-0.647†	-0.043	-0.003	0.707	0.098	0.142
BARTScore	-0.088	-0.392†	0.575	0.161	-0.454	-0.835†	-0.894	-0.225	0.038	0.150	-0.134
BARTScore+	0.039	-0.159	0.536	0.192	-0.553	-0.665†	-0.841	-0.141	0.303	0.261	0.066
BARTScore-P	-0.007	-0.152	0.546	0.169	-0.473	-0.642†	-0.894	-0.039	0.380	0.265	0.012
PRISM	$0.129\dagger$	-0.081	0.465	0.207	-0.379	-0.733†	-0.361	-0.071	0.668	0.165	0.048
CTC-Relevance	-0.027	-0.100	-0.072	$0.333\dagger$	-0.041	-0.622†	-0.087	-0.091	0.394	0.134	-0.056
CTC-Consistency	0.030	-0.133	0.243	0.462†	-0.148	-0.657†	-0.106	-0.041	0.769	0.106	-0.002
ROSCOE Metrics (refer	ence-free	on (s, h))								
ROSCOE-SA with finetu	ined sup-s	imcse-rob		ntence em	beddings						
Faithfulness-Step	0.096	-0.095	0.572	$0.496\dagger$	0.143	-0.566†	-0.322	0.278	0.760	-0.022	0.042
Faithfulness-Token	$0.177\dagger$	-0.016	0.472	$0.417\dagger$	-0.189	-0.525†	-0.038	0.131	0.678	0.002	0.130
Info-Step	$0.142\dagger$	-0.069	0.643	0.500†	0.216	-0.581†	-0.284	0.287†	0.846	0.024	0.125
Repetition-Token	0.055	$0.210\dagger$	0.018	0.170	0.340	0.578†	0.952	0.026	-0.822	-0.140	0.062
ROSCOE-SA with all-m	pnet-base	-v2 senten	ce embeddir	ngs							
Faithfulness-Step	0.054	-0.156	0.153	0.297†	-0.073	-0.578†	-0.111	0.116	-0.221	0.001	-0.042
Faithfulness-Token	$0.156\dagger$	-0.042	0.362	0.290†	-0.388	-0.504†	-0.216	-0.058	0.288	0.063	0.085
Info-Step	0.090	-0.116	0.308	0.301†	-0.153	-0.619†	-0.043	0.021	0.250	0.078	0.020
Repetition-Token	$0.130\dagger$	0.370†	0.027	0.087	0.313	0.790†	0.519	0.114	-0.822	-0.181	0.053
ROSCOE-SA with sup-s	imcse-rob	erta-base	sentence em	beddings							
Faithfulness-Step	0.084	-0.128	0.195	0.477†	-0.061	-0.643†	-0.596	0.085	0.236	-0.061	-0.051
Faithfulness-Token	0.186†	-0.031	0.414	0.454†	-0.245	-0.574†	-0.139	0.098	0.514	0.047	0.093
Info-Step	0.141†	-0.095	0.443	0.510†	-0.121	-0.692†	-0.471	0.040	0.731	0.066	0.010
Repetition-Token	0.055	0.210†	0.018	0.170	0.340	0.578†	0.952	0.026	-0.822	-0.140	0.062
ROSCOE-SS with finetu	ined sup-s	imcse-rob	erta-base se	ntence em	beddings						
Info-Chain	0.108	-0.128	0.462	0.222	-0.393	-0.446†	-0.375	0.025	0.207	0.446†	0.114
Repetition-Step	0.036	0.400†	-0.543	-0.258†	0.165	0.824	0.909	0.052	-0.822	0.192	-0.016
ROSCOE-SS with all-m	pnet-base	-v2 senten	ce embeddir	ngs		'					
Info-Chain	0.165†	0.020	0.395	0.240†	-0.027	-0.480†	-0.135	0.122	0.394	0.406†	0.106
Repetition-Step	0.052	0.358†	-0.111	-0.115	0.260	0.791†	0.856	0.105	-0.822	0.066	0.015
ROSCOE-SS with sup-s				beddings							
Info-Chain	0.138†	-0.090	0.504	0.271†	-0.158	-0.664†	-0.130	-0.089	0.611	0.381	0.093
Repetition-Step	0.073	0.357†	-0.143	-0.085	0.451	0.799†	0.918	0.128	-0.822	0.047	0.021
ROSCOE-LI		· ·			_						
Source-Consistency	0.200†	0.243†	0.462	0.390†	-0.004	0.085	0.697	0.009	0.365	0.420	0.184†
Self-Consistency	0.032	0.245	-0.076	0.198	0.201	0.584†	0.037	0.187	-0.707	0.344	-0.075
	0.002	3.270	0.07.0	0.170	0.201	3.00.1	0.107	0.107	0.707	0.0.1	0.075
ROSCOE-LC	0.220+	0.141	0.250	0.260	0.536	0.001	-0.553	0.169	0.260	0.450	0.111
Grammar	0.220†										
Perplexity-Step	0.185†	0.034	0.214	-0.002	-0.112	-0.320†	-0.827	-0.043	0.332	0.259	0.205†
Perplexity-Chain	0.087	-0.152	0.185	-0.104	-0.515	-0.611†	-0.952	-0.178	0.663	0.103	0.117

Table 34: Somers' D correlation of all metrics on **GSM8K human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold**, second best is <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \uparrow . (Continued from \S 6, more details in App. H.2)

	QUAL	COH	COMMON	FACT	HALL	RED	REP	LOGIC	MATH	GRAM	MISS
Rouge-1	0.122	$0.202\dagger$	0.071	0.156	0.275	0.037	0.222	0.073	0.233	0.173	$0.180\dagger$
Rouge-2	0.089	0.139	0.060	0.102	0.175	-0.016	0.561	0.030	0.113	0.051	0.148
Rouge-L	0.176†	0.268†	0.195†	0.169	0.180	0.054	0.558	0.120	0.212	-0.125	$0.252\dagger$
BLEURT	$0.160\dagger$	$0.248\dagger$	0.134	0.256†	0.099	0.075	0.227	0.124	0.057	0.077	$0.210\dagger$
BERTScore	$0.173\dagger$	$0.220\dagger$	0.112	0.116	0.168	0.095	0.955	0.138	0.113	0.054	$0.235\dagger$
BARTScore	0.009	0.035	-0.047	0.044	-0.246	-0.261	0.424	-0.002	-0.546†	-0.321	-0.003
BARTScore+	0.064	0.132	0.078	0.054	-0.003	-0.039	0.879	0.085	-0.298	-0.281	0.090
BARTScore-P	0.037	0.059	0.042	0.061	-0.116	-0.102	0.561	-0.002	-0.322	-0.097	0.019
PRISM	-0.112	-0.075	-0.099	-0.037	-0.385	-0.455†	-0.086	-0.171	-0.341	0.130	-0.093
CTC-Relevance	-0.086	$-0.148\dagger$	-0.077	-0.111	0.009	-0.371†	0.566	-0.106	-0.093	-0.061	-0.088
CTC-Consistency	-0.157†	-0.203†	-0.206†	-0.129	-0.013	-0.318†	0.556	-0.174†	-0.019	-0.056	-0.204†
ROSCOE Metrics (refer	ence-free	on (s,h)									
ROSCOE-SA with finetu	ined sup-si	mcse-robe	rta-base sen	tence emb	eddings						
Faithfulness-Step	0.012	0.083	0.097	0.014	-0.138	-0.019	0.182	-0.024	-0.101	0.286	0.027
Faithfulness-Token	-0.012	0.036	-0.052	-0.067	-0.099	-0.304	0.788	-0.042	-0.408	0.204	0.016
Info-Step	0.059	0.137	0.178†	0.082	-0.025	0.059	0.364	0.031	-0.303	0.329	0.108
Repetition-Token	0.200†	$0.193\dagger$	0.186†	0.075	0.224	$0.392\dagger$	0.788	$0.183\dagger$	0.443	0.173	$0.270\dagger$
ROSCOE-SA with all-m	pnet-base-	v2 sentence	e embedding	gs							
Faithfulness-Step	-0.039	-0.016	0.002	-0.027	0.012	-0.376	0.581	-0.173	-0.127	0.102	-0.021
Faithfulness-Token	-0.078	-0.002	-0.090	-0.079	-0.141	$-0.443\dagger$	0.485	-0.137	$-0.623\dagger$	0.074	-0.068
Info-Step	0.095	0.121	0.099	0.139	0.112	-0.065	0.662	-0.016	-0.056	0.092	0.148
Repetition-Token	$0.214\dagger$	$0.208\dagger$	0.184†	0.081	0.500 †	0.238	0.747	0.208†	0.498†	0.339	0.306†
ROSCOE-SA with sup-s	imcse-robe	erta-base s	entence emb	oeddings							
Faithfulness-Step	-0.061	-0.034	-0.006	0.002	-0.158	-0.142	0.237	-0.192†	-0.221	0.168	-0.045
Faithfulness-Token	-0.051	-0.008	-0.087	-0.084	-0.168	-0.352	0.732	-0.098	-0.424	0.145	-0.034
Info-Step	0.059	0.102	0.106	0.124	0.040	0.046	0.545	-0.053	-0.236	0.099	0.131
Repetition-Token	0.200†	0.193†	0.186†	0.075	0.224	0.392†	0.788	0.183†	0.443	0.173	0.270†
ROSCOE-SS with finetu	ined sup-si	mcse-robe	rta-base sen	tence emb	eddings						
Info-Chain	0.097	0.064	0.059	0.080	0.035	$0.385\dagger$	0.722	0.103	0.121	-0.110	0.130
Repetition-Step	$0.199\dagger$	$0.166\dagger$	0.168	0.145	0.254	0.514†	0.869	0.152	0.230	0.176	$0.222\dagger$
ROSCOE-SS with all-m	pnet-base-	v2 sentence	e embedding	gs							
Info-Chain	0.059	0.071	0.020	0.043	0.039	0.128	0.288	0.004	0.260	-0.026	0.098
Repetition-Step	$0.218\dagger$	$0.161\dagger$	0.167	$0.227\dagger$	0.309	$0.471\dagger$	0.323	0.158	0.486	-0.301	$0.245\dagger$
ROSCOE-SS with sup-s	imcse-robe	erta-base s	entence emb	eddings							
Info-Chain	0.042	0.051	0.043	0.057	0.240	0.041	0.808	-0.007	-0.211	0.038	0.130
Repetition-Step	$0.322\dagger$	0.299†	$0.275\dagger$	$0.227\dagger$	<u>0.466</u>	$0.638\dagger$	0.879	$0.192\dagger$	$0.563\dagger$	-0.311	$0.354\dagger$
ROSCOE-LI											
Source-Consistency	0.108	0.037	-0.019	0.172†	0.121	0.030	0.551	0.082	-0.005	0.219	0.097
Self-Consistency	0.283†	0.267†	0.177	0.345†	0.207	0.354	0.980	0.219†	0.087	-0.230	0.223†
ROSCOE-LC											
Grammar	-0.134†	-0.159†	-0.260†	-0.081	-0.234	-0.246	0.056	-0.264†	-0.207	0.298	-0.184†
Perplexity-Step	-0.297†	-0.278†	-0.366†	-0.307†	-0.608†	-0.591†	-0.136	-0.334†	-0.514†	0.148	-0.331†
Perplexity-Chain	-0.332†	-0.336†	-0.322†	-0.273†	-0.695†	-0.682†	-0.556	-0.354†	-0.697†	-0.084	-0.408†
						1					

Table 35: Somers' D correlation of all metrics on **ESNLI human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold** and second best is <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from \S 6, more details in App. H.2)

	QUAL	СОН	COMMON	FACT	HALL	RED	REP	LOGIC	MATH	GRAM	MISS
Rouge-1	0.213†	-0.121	-0.373	0.152	-0.060	-0.284	-	-0.276†	0.113	-0.403	0.559†
Rouge-2	0.216†	-0.075	-0.336	0.169	-0.097	-0.240	-	-0.255†	-0.080	-0.432	0.502†
Rouge-L	0.222^{+}	-0.041	-0.296	0.201	-0.092	-0.245	-	-0.259	0.053	-0.429	0.557†
BLEURT	$0.154\dagger$	-0.202	-0.271	0.112	-0.041	-0.230	-	$-0.322\dagger$	0.267	-0.641†	0.541†
BERTScore	0.255†	0.019	-0.205	0.222	0.153	-0.179	-	-0.141	-0.360	-0.252	0.590†
BARTScore	0.189†	-0.096	-0.260	0.079	0.071	-0.259	-	-0.200	0.533	-0.688†	0.549†
BARTScore+	0.209†	-0.059	-0.192	0.166	0.204	-0.147	_	-0.168	-0.160	-0.578†	0.482†
BARTScore-P	0.092	-0.096	-0.208	0.213	0.092	-0.209	-	-0.192	-0.587	-0.370	0.255†
PRISM	0.264†	0.089	-0.151	0.093	0.224	-0.235	-	-0.042	-0.240	-0.329	0.580†
CTC-Relevance	0.071	-0.051	-0.074	0.063	-0.020	-0.098	-	-0.093	0.720	-0.205	0.334†
CTC-Consistency	0.029	0.218	$0.647\dagger$	-0.009	0.483	0.021	-	0.350†	0.533	-0.071	-0.060
ROSCOE Metrics (refer	rence-free	on (s, h)))								
ROSCOE-SA with fineto	uned sup-s	imcse-rol		entence en	nbeddings						
Faithfulness-Step	0.236†	-0.163	-0.027	0.096	-0.146	-0.139	-	-0.293†	0.560	-0.386	0.403†
Faithfulness-Token	$0.279\dagger$	0.010	-0.145	0.232	0.201	-0.279	-	-0.119	0.067	-0.293	0.521†
Info-Step	$0.222\dagger$	-0.127	-0.140	0.143	-0.068	-0.220	-	-0.333†	0.400	-0.310	0.493†
Repetition-Token	-0.037	-0.020	0.403	0.244	0.330	$0.441\dagger$	-	0.153	0.240	-0.059	-0.436†
ROSCOE-SA with all-m	ipnet-base	-v2 senter	nce embeddi	ngs							
Faithfulness-Step	0.261†	-0.117	-0.058	0.242	-0.425	-0.296†	-	-0.248†	0.480	-0.244	0.423†
Faithfulness-Token	0.270†	-0.083	-0.216	0.136	0.313	-0.248	-	-0.097	0.347	-0.433	0.524†
Info-Step	0.290†	-0.006	-0.049	$0.286\dagger$	-0.102	-0.245	-	-0.246†	0.373	-0.342	0.542†
Repetition-Token	-0.040	0.129	-0.189	0.169	0.799†	0.617†	-	0.151	-0.267	0.084	-0.445†
ROSCOE-SA with sup-s	simcse-rob	erta-base		nbeddings							
Faithfulness-Step	$0.262\dagger$	-0.083	-0.244	0.257	-0.163	$-0.372\dagger$	-	-0.283†	0.413	-0.375	0.502†
Faithfulness-Token	0.275†	0.004	-0.230	0.206	0.214	-0.303†	-	-0.140	0.347	-0.321	0.540†
Info-Step	0.295	0.046	-0.126	$0.299\dagger$	-0.071	-0.306†	-	-0.208	0.333	-0.397†	0.599†
Repetition-Token	-0.037	-0.020	0.403	0.244	0.330	0.441†	-	0.153	0.240	-0.059	-0.436†
ROSCOE-SS with fineto	uned sup-s	simcse-rol	berta-base s	entence en	nbeddings	;					
Info-Chain	0.083	-0.197	-0.466	-0.225	-0.510	-0.259	-	-0.067	0.267	0.074	$0.310\dagger$
Repetition-Step	-0.021	0.103	-0.063	0.249	0.310	$0.530\dagger$	-	0.178	-0.107	0.053	-0.433†
ROSCOE-SS with all-m	ipnet-base	-v2 senter	nce embeddi	ngs							
Info-Chain	0.191†	-0.233	-0.438	0.052	-0.153	-0.253	-	-0.368†	0.533	-0.452	$0.507\dagger$
Repetition-Step	0.007	0.280	0.173	$0.365\dagger$	0.694	$0.487\dagger$	-	0.225	0.213	0.193	-0.472†
ROSCOE-SS with sup-s	simcse-rob	erta-base	sentence en	nbeddings							
Info-Chain	0.211†	-0.083	-0.192	0.212	0.024	-0.319†	-	-0.257†	0.400	-0.455	0.531†
Repetition-Step	-0.015	0.180	0.397	0.332	0.367	0.484†	-	0.115	0.013	0.125	-0.501†
ROSCOE-LI											
Source-Consistency	0.012	-0.028	0.334	0.425†	0.299	-0.055	_	0.112	0.600	-0.258	-0.222†
Self-Consistency	-0.028	0.354†	0.156	0.087	0.561	0.370†	-	0.531†	-0.333	0.222	-0.351†
ROSCOE-LC											
Grammar	0.063	0.069	-0.411	-0.003	0.122	-0.286	_	0.255†	-0.027	0.411	-0.072
Perplexity-Step	0.084	0.087	-0.655†	-0.153	0.255	-0.254	_	-0.148	-0.147	0.090	0.345†
Perplexity-Step Perplexity-Chain	0.034	0.087	-0.635† -0.616†	-0.133	0.233	-0.234 -0.447†	-	-0.146	-0.147	0.090	0.343
1 cipiexity-challi	0.027	0.001	-0.010	-0.209	0.073	-0.447		-0.133	-0.555	0.249	0.447

Table 36: Somers' D correlation of all metrics on **COSMOS human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold** and second best is <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from \S 6, more details in App. H.2)

	QUAL	COH	COMMON	FACT	HALL	RED	REP	LOGIC	GRAM	MISS
Rouge-1	-0.011	-0.007	-0.182	-0.077	0.292	-0.576†	-0.807	-0.264†	-0.644†	0.113
Rouge-2	0.021	0.028	-0.075	-0.131	0.239	-0.561†	-0.174	-0.108	-0.314	$0.180\dagger$
Rouge-L	0.011	-0.013	-0.094	-0.044	0.252	-0.637†	-0.436	-0.114	-0.441†	0.141
BLEURT	0.098	0.088	0.019	-0.054	0.097	-0.686†	-0.617	-0.181†	-0.522†	$0.218\dagger$
BERTScore	0.095	0.113	0.059	-0.055	0.234	-0.478†	-0.492	-0.058	-0.420†	0.114
BARTScore	0.009	0.024	0.159	-0.026	0.001	-0.544†	-0.208	-0.122	-0.420	0.124
BARTScore+	0.048	0.061	-0.102	-0.004	0.159	-0.507†	-0.602	-0.186†	-0.499†	0.159
BARTScore-P	0.009	0.021	-0.149	0.010	0.267	-0.385†	-0.508	-0.207†	-0.453†	0.142
PRISM	0.058	0.091	-0.046	-0.156	0.311	-0.446†	-0.428	-0.036	-0.376†	0.157
CTC-Relevance	0.070	0.035	0.246	0.155	0.233	-0.294	-0.780	-0.001	-0.349	0.016
CTC-Consistency	0.093	0.097	0.275†	0.084	0.140	-0.032	-0.201	0.051	0.064	-0.006
ROSCOE-SA with fineto	ined sun-	simcse-rol	berta-base se	entence er	nbedding	s				
Faithfulness-Step	0.089	0.109	-0.008	0.149	0.256	-0.322	-0.216	-0.019	0.034	-0.011
Faithfulness-Token	0.038	0.039	-0.032	-0.093	0.285	-0.472†	-0.220	-0.026	-0.320†	0.012
Info-Step	0.109	0.082	-0.011	0.119	0.302	-0.340	-0.811	-0.039	-0.135	0.057
Repetition-Token	0.050	0.120	0.187	0.103	0.103	0.555†	0.231	0.192†	0.458†	-0.233†
ROSCOE-SA with all-m			nce embeddi							
Faithfulness-Step	-0.011	0.044	-0.008	-0.012	0.127	-0.424†	-0.235	-0.061	-0.268	-0.096
Faithfulness-Token	0.035	0.078	-0.124	-0.128	0.323	-0.515†	-0.667	-0.036	-0.230	-0.008
Info-Step	0.011	0.043	-0.047	-0.023	0.186	-0.429†	-0.481	-0.081	-0.268	-0.037
Repetition-Token	0.045	0.092	0.275†	0.122	0.298	0.638†	0.398	0.252†	0.386	-0.184†
ROSCOE-SA with sup-s	simcse-rob	erta-base		nbeddings	3	'				
Faithfulness-Step	0.036	0.022	0.072	0.111	0.184	-0.381†	-0.318	-0.076	-0.237	-0.053
Faithfulness-Token	0.036	0.030	-0.003	-0.058	0.205	-0.486†	-0.333	-0.056	-0.420†	-0.020
Info-Step	0.026	-0.025	0.069	0.079	0.248	-0.409†	-0.720	-0.088	-0.292	0.000
Repetition-Token	0.050	0.120	0.187	0.103	0.103	0.555†	0.231	$0.192\dagger$	$0.458\dagger$	-0.233†
ROSCOE-SS with fineto	ined sup-s	simcse-rol	berta-base se	entence er	nbedding	S		-		-
Info-Chain	0.013	0.089	0.159	0.048	0.139	-0.219	-0.932	0.113	-0.390	-0.117
Repetition-Step	0.011	0.050	0.178	0.132	-0.046	0.593†	0.670	0.295†	0.330	-0.244†
ROSCOE-SS with all-m	pnet-base	-v2 senter	nce embeddi	ngs		'				
Info-Chain	0.073	0.091	-0.174	0.033	0.145	-0.215	-0.409	-0.198†	-0.411†	0.103
Repetition-Step	0.047	0.127	0.124	0.060	0.153	$0.642\dagger$	0.617	0.346	0.563†	$-0.184\dagger$
ROSCOE-SS with sup-s	simcse-rob	erta-base	sentence en	nbeddings	3				'	
Info-Chain	0.114	0.055	0.034	-0.005	0.218	-0.367†	-0.879	-0.046	-0.222	0.153
Repetition-Step	0.061	0.127	0.095	0.076	0.246	0.658†	0.500	0.256†	0.496†	-0.145
ROSCOE-LI									-	
Source-Consistency	0.184†	0.183†	0.150	0.285†	0.241	0.444†	0.091	0.111	0.303	0.011
Self-Consistency	0.104	$\frac{0.185}{0.080}$	0.190	0.173	-0.021	0.417†	0.610	0.111 0.192†	0.401†	-0.252†
	0.010	0.000	0.170	0.173	0.021	3.117	0.010	3.172	3.101	0.232
ROSCOE-LC	0.002	0.100	0.065	0.004	0.022	0.013	0.256	0.012	0.206	0.051
Grammar	0.093	0.189†	-0.065	0.084	0.022	-0.013	0.356	-0.013	0.386	0.051
Perplexity-Step	0.122†	0.157†	-0.208	-0.021	0.034	0.028	-0.140	-0.113	-0.295	0.064
Perplexity-Chain	0.083	0.047	-0.193	0.001	-0.073	-0.311†	-0.561	-0.212†	-0.542†	0.130

Table 37: Somers' D correlation of all metrics on **SemEVAL human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-free** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold** and second best is <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from \S 6, more details in App. H.2)

	QUAL	COH	COMMON	FACT	HALL	RED	REP	LOGIC	GRAM	MISS
Rouge-1	-0.199†	-0.208†	-0.246	-0.118	-0.053	-0.520†	-0.398	-0.206	-0.230	0.108
Rouge-2	-0.086	0.060	-0.478	-0.063	0.561	-0.232	-0.083	-0.073	0.141	0.090
Rouge-L	-0.198†	-0.209†	-0.498	-0.104	0.090	-0.478†	-0.396	-0.115	-0.156	0.058
BLEURT	-0.313†	-0.383†	-0.372	-0.208	-0.034	-0.482†	-0.383	-0.356†	-0.104	-0.074
BERTScore	-0.051	0.064	-0.517	-0.035	0.524	-0.002	-0.218	-0.295†	0.266	0.063
BARTScore	-0.084	-0.059	-0.140	-0.137	-0.369	-0.054	0.209	-0.284	-0.040	-0.056
BARTScore+	-0.046	-0.098	0.652	-0.033	-0.073	0.056	0.204	-0.204	0.048	-0.016
BARTScore-P	-0.075	-0.168	0.633	-0.080	0.107	0.096	0.277	-0.230	0.054	-0.082
PRISM	-0.082	-0.040	-0.469	-0.115	0.073	-0.174	-0.354	-0.134	-0.079	0.075
CTC-Relevance	-0.146†	-0.219†	0.256	-0.047	0.442	-0.071	-0.121	-0.349†	0.101	-0.145
CTC-Consistency	-0.178†	-0.241†	0.101	0.009	0.583	-0.301†	-0.296	-0.335†	0.208	-0.142
ROSCOE Metrics (refer	rence-free	on (s,h)								
ROSCOE-SA with fineto			rta-base ser	tence emb	eddings					
Faithfulness-Step	0.157†	0.198†	0.275	0.189	0.694	0.233	-0.262	-0.115	0.092	0.115
Faithfulness-Token	0.030	0.182†	-0.150	0.009	0.714	0.116	-0.442	-0.101	0.262	0.085
Info-Step	0.068	0.113	0.111	0.055	0.835	0.003	-0.252	-0.187	0.196	0.129
Repetition-Token	0.062	0.150	0.401	0.021	-0.078	0.337†	0.670	-0.074	-0.007	-0.131
ROSCOE-SA with all-m										
Faithfulness-Step	0.197†	0.221†	0.005	0.201	0.461	0.330†	-0.160	-0.019	0.182	0.108
Faithfulness-Token	0.030	0.161	-0.208	0.013	0.597	-0.044	-0.495	-0.111	0.063	0.186†
Info-Step	0.111	0.127	-0.111	0.125	0.544	0.161	-0.073	-0.176	0.210	0.076
Repetition-Token	0.134†	0.178†	0.662	0.066	0.364	0.485†	0.772	-0.004	0.157	-0.115
ROSCOE-SA with sup-s					0.00.	01.00	01772	0.00	0.107	01110
Faithfulness-Step	0.028	0.123	-0.459	0.112	0.908	0.174	-0.199	-0.162	0.021	0.059
Faithfulness-Token	-0.021	0.133	-0.227	0.021	0.752	0.084	-0.398	-0.134	0.119	0.052
Info-Step	-0.040	-0.003	-0.362	0.024	0.777	0.015	-0.296	-0.321†	0.038	0.039
Repetition-Token	0.062	0.150	0.401	0.021	-0.078	0.337†	0.670	-0.074	-0.007	-0.131
ROSCOE-SS with fineto						0.007	0.070	0.07.	0.007	0.121
Info-Chain	-0.040	0.008	-0.082	-0.005	0.539	-0.164	-0.015	-0.070	0.086	-0.082
Repetition-Step	0.051	0.143	0.546	0.043	-0.024	0.411†	0.723	0.078	-0.018	-0.062 -0.195†
ROSCOE-SS with all-m					-0.024	0.411	0.723	0.050	-0.016	-0.193
Info-Chain	0.076	0.128	-0.700	0.122	-0.005	0.367†	0.117	-0.140	0.068	-0.018
Repetition-Step	0.070	0.128	0.633	0.122	0.141	0.508†	0.684	-0.140	0.005	-0.016
ROSCOE-SS with sup-s					0.141	0.500	0.064	-0.020	0.023	-0.133
Info-Chain	-0.115	-0.038	-0.169	0.059	0.539	-0.260	-0.476	-0.169	-0.081	-0.009
Repetition-Step	0.104	0.132	0.787	0.039	0.339	0.535†	0.476 0.811	0.006	0.069	-0.009 -0.169†
	0.104	0.132	0.767	0.042	0.130	0.555	0.011	0.000	0.009	-0.109
ROSCOE-LI										
Source-Consistency	0.059	0.016	0.546	0.206	-0.029	0.006	0.010	-0.241	-0.139	-0.063
Self-Consistency	<u>0.162</u> †	0.250^{\dagger}	0.536	0.104	0.383	0.372†	0.223	-0.091	0.061	-0.075
ROSCOE-LC					-		-			
Grammar	-0.076	0.014	0.101	-0.223†	-0.335	0.104	0.393	-0.215	0.517†	-0.144
Perplexity-Step	-0.026	-0.053	0.797	0.037	-0.607	0.020	-0.019	-0.071	-0.330	-0.039
Perplexity-Chain	-0.141†	-0.237†	0.324	-0.126	-0.650	-0.373†	-0.481	-0.151	-0.284	0.039
	0.1.1	0.207	0.02.	0.120	0.000	0.5751	001	0.101	0.201	0.007

Table 38: Somers' D correlation of all metrics on **GSM8K human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-based** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold**, second best is <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from \S 6, more details in App. H.2)

	QUAL	СОН	COMMON	FACT	HALL	RED	REP	LOGIC	MATH	GRAM	MISS
Rouge-1	0.572†	0.506†	0.590†	0.513†	0.532†	0.498†	0.694	0.533†	0.582†	0.430	0.686†
Rouge-2	0.551†	0.520†	0.566†	0.511†	0.467†	0.439†	0.758	0.555†	0.617†	0.480	0.713†
Rouge-L	0.591†	0.542†	0.605†	0.559†	0.452†	0.613	0.662	0.575†	0.715†	0.457	0.730†
BLEURT	0.487†	0.391†	0.502†	0.389†	0.392	0.222	0.682	0.404†	0.233	0.283	0.597†
BERTScore	0.505†	0.425†	0.488†	0.475†	$0.477\dagger$	0.585†	0.763	0.483†	0.677†	0.347	0.627†
BARTScore	0.429†	0.352†	0.522†	0.317†	0.315	0.128	0.692	0.436†	0.279†	0.237	0.555†
BARTScore+	0.531†	0.455†	0.579†	0.460†	0.273	0.303	0.722	0.501†	0.577†	0.334	0.679†
BARTScore-P	$0.343\dagger$	$0.280\dagger$	0.376†	$0.268\dagger$	0.210	0.062	0.621	$0.298\dagger$	0.207	0.227	$0.441\dagger$
PRISM	$0.579\dagger$	0.511†	0.593†	$0.531\dagger$	$0.363\dagger$	$0.392\dagger$	0.707	$0.557\dagger$	$0.540\dagger$	0.423	$0.728\dagger$
CTC-Relevance	-0.047	-0.093	-0.065	-0.065	-0.134	-0.286	0.460	-0.055	0.056	0.242	-0.102
CTC-Consistency	-0.272†	-0.279†	-0.259†	-0.264†	-0.182	-0.385†	-0.399	-0.192†	-0.233	0.028	-0.320†
ROSCOE Metrics (reference-l											
ROSCOE-SA with finetuned sa											
Hallucination	$0.466\dagger$	$0.460\dagger$	0.398†	$0.445\dagger$	0.371†	$0.462\dagger$	0.091	$0.439\dagger$	$0.624\dagger$	0.319	0.519†
Redundancy	$0.464\dagger$	0.463†	0.396†	$0.455\dagger$	$0.379\dagger$	0.552†	0.242	$0.459\dagger$	0.699†	0.472	$0.538\dagger$
Semantic Coverage-Step	-0.030	0.002	0.017	-0.154	-0.237	-0.084	-0.298	0.023	0.335	0.064	0.002
Missing Step	$0.484\dagger$	$0.419\dagger$	0.509†	$0.364\dagger$	$0.357\dagger$	0.325	0.333	$0.438\dagger$	0.685†	0.207	0.590†
Reasoning Alignment	0.613†	0.590†	0.570†	0.573†	$0.464\dagger$	0.529†	0.460	$0.587\dagger$	$0.725\dagger$	0.446	0.711†
Commonsense	$0.411\dagger$	$0.348\dagger$	$0.424\dagger$	$0.289\dagger$	0.264	0.246	0.404	$0.358\dagger$	0.677†	0.367	$0.494\dagger$
ROSCOE-SA with all-mpnet-b	ase-v2 sen	itence emb	eddings								
Hallucination	$0.702\dagger$	0.621†	0.704†	$0.630\dagger$	0.629 †	$0.422\dagger$	0.818	$0.644\dagger$	0.340	0.064	$0.846\dagger$
Redundancy	$0.453\dagger$	0.411†	0.451†	$0.417\dagger$	$0.406\dagger$	0.353†	0.505	$0.468\dagger$	0.195	0.204	$0.570\dagger$
Semantic Coverage-Step	-0.066	-0.062	0.005	-0.042	-0.148	$0.329\dagger$	-0.576	0.006	0.232	0.291	-0.085
Missing Step	0.501†	0.449†	$0.520\dagger$	$0.450\dagger$	$0.402\dagger$	0.224	0.566	$0.488\dagger$	0.503†	0.186	$0.572\dagger$
Reasoning Alignment	0.583†	0.523†	0.559†	$0.558\dagger$	$0.380\dagger$	0.381†	0.571	$0.569\dagger$	$0.418\dagger$	0.296	$0.716\dagger$
Commonsense	$0.626\dagger$	0.529†	0.657†	$0.594\dagger$	0.578†	0.305	0.828	$0.574\dagger$	0.421†	-0.125	0.758†
ROSCOE-SA with sup-simcse-			ce embedding								
Hallucination	0.806†	0.707†	0.731†	0.729†	0.625	0.357†	0.611	$0.712\dagger$	0.575†	-0.105	0.896†
Redundancy	$0.734\dagger$	0.661	0.665†	$0.699\dagger$	$0.598\dagger$	0.463†	0.717	$0.664\dagger$	$0.517\dagger$	-0.061	$0.840\dagger$
Semantic Coverage-Step	0.063	0.063	0.127	0.057	0.216	0.300	-0.111	0.161	-0.205	0.492	0.023
Missing Step	$0.722\dagger$	0.631†	0.685†	$0.691\dagger$	$0.517\dagger$	$0.353\dagger$	0.808	$0.655\dagger$	$0.482\dagger$	0.003	$0.821\dagger$
Reasoning Alignment	$0.712\dagger$	$0.628\dagger$	0.654†	$0.711\dagger$	$0.487\dagger$	$0.435\dagger$	0.717	$0.655\dagger$	$0.603\dagger$	0.161	$0.848\dagger$
Commonsense	0.780†	0.659†	0.721†	0.727†	0.486†	0.324	0.813	0.691	0.510†	-0.217	0.887†
ROSCOE-SS with finetuned sa	ıp-simcse-	roberta-ba	se sentence	embedding	gs						
Semantic Coverage-Chain	0.404†	$0.364\dagger$	0.452†	0.324†	0.214	0.381†	0.601	$0.350\dagger$	0.172	-0.196	$0.444 \dagger$
ROSCOE-SS with all-mpnet-b	ase-v2 sen	itence emb	eddings								
Semantic Coverage-Chain	$0.170\dagger$	0.171†	0.189†	$0.203\dagger$	-0.055	0.205	0.237	$0.186\dagger$	-0.014	0.041	$0.243\dagger$
ROSCOE-SS with sup-simcse-	roberta-bo	<i>ise</i> sentend	e embeddin	gs							
Semantic Coverage-Chain	0.411†	0.375†	0.422†	0.460†	0.301	0.439†	0.742	0.381†	0.275	0.148	0.506†
Semantic Coverage-Challi	0.411	0.575	0.722	0.400	0.501	0.437	0.742	0.501	0.273	0.140	0.500

Table 39: Somers' D correlation of all metrics on **ESNLI human judged** dataset analyzing step-by-step reasoning on overall chain and step-level perspectives. All metrics are measured **reference-based** on (s, h). The highest correlation overall for each aspect on each dataset is in **bold**, second best is <u>underlined</u>. Correlation scores with p-value < 0.05 are marked with \dagger . (Continued from \S 6, more details in App. H.2)

	QUAL	СОН	COMMON	FACT	HALL	RED	REP	LOGIC	MATH	GRAM	MISS
Rouge-1	0.255†	-0.234	-0.506†	0.072	0.359	-0.078	_	-0.310†	-0.054	-0.487	0.662†
Rouge-2	0.223†	-0.189	-0.629†	0.103	0.429	-0.081	-	-0.344†	-0.149	-0.493	0.568†
Rouge-L	0.227†	-0.177	-0.647†	0.084	0.383	-0.089	-	-0.351†	-0.500	-0.493	0.628†
BLEURT	0.221†	-0.170	-0.197	0.078	-0.090	-0.108	-	-0.202	0.216	-0.447	0.611†
BERTScore	0.362†	-0.036	-0.281	0.182	0.269	0.143	-	-0.153	0.068	-0.369	0.661†
BARTScore	0.121	-0.304	-0.200	0.003	0.007	-0.339†	-	-0.322†	0.338	-0.478	0.513†
BARTScore+	0.129	-0.113	-0.306	0.164	0.048	-0.036	-	-0.367†	-0.095	-0.378	0.341†
BARTScore-P	0.096	-0.037	-0.244	0.188	-0.017	0.028	-	-0.222	-0.946	-0.244	0.173
PRISM	0.314†	0.032	-0.386	0.086	0.438	0.063	-	-0.131	-0.378	-0.328	0.684^{\dagger}
CTC-Relevance	0.072	0.331†	0.383	0.034	0.352	-0.001	-	0.400†	0.405	0.231	0.033
CTC-Consistency	-0.051	0.150	-0.078	-0.130	0.083	0.105	-	0.085	0.676	0.056	0.006
ROSCOE Metrics (reference-											
ROSCOE-SA with finetuned sa											
Hallucination	$0.156\dagger$	0.152	0.461	-0.011	-0.021	0.160	-	0.132	0.743	-0.275	0.170
Redundancy	0.142	0.234	0.553†	0.046	0.283	0.119	-	0.159	0.500	-0.253	0.145
Semantic Coverage-Step	$0.153\dagger$	-0.094	0.172	-0.192	-0.241	-0.065	-	-0.020	-0.730	0.086	$0.327\dagger$
Missing Step	$0.234\dagger$	-0.197	-0.375	0.239	-0.010	-0.275	-	-0.226	0.527	-0.400	0.558†
Reasoning Alignment	$0.278\dagger$	-0.062	-0.003	0.072	0.100	-0.049	-	-0.047	-0.108	-0.433	0.495†
Commonsense	0.142	-0.094	-0.353	0.213	0.207	-0.184	-	-0.148	0.676	-0.447	0.368†
ROSCOE-SA with all-mpnet-l											
Hallucination	$0.174\dagger$	0.094	0.350	0.066	-0.210	-0.021	-	0.181	0.554	-0.089	0.141
Redundancy	$0.219\dagger$	0.227	0.531	0.153	0.031	0.116	-	0.293	0.405	-0.133	0.133
Semantic Coverage-Step	$0.185\dagger$	-0.108	-0.039	0.099	-0.490	-0.229	-	-0.185	-0.122	0.083	0.350†
Missing Step	0.303†	-0.173	-0.603†	-0.017	0.159	-0.031	-	-0.180	0.689	-0.650†	<u>0.679</u> †
Reasoning Alignment	0.428^{\dagger}	-0.028	-0.128	0.294†	0.121	0.095	-	-0.066	-0.351	-0.547†	0.657†
Commonsense	0.211†	-0.029	-0.567†	0.062	0.324	-0.016	-	-0.071	0.946	-0.564	$0.489\dagger$
ROSCOE-SA with sup-simcse											
Hallucination	$0.190\dagger$	0.086	-0.311	0.171	-0.138	-0.050	-	0.001	-0.149	-0.322	$0.272\dagger$
Redundancy	$0.166\dagger$	0.177	-0.208	0.249	-0.038	-0.007	-	0.023	-0.405	-0.331	$0.235\dagger$
Semantic Coverage-Step	0.196†	-0.235	-0.058	0.167	-0.341	-0.187	-	-0.153	-0.338	-0.078	0.425†
Missing Step	0.307†	-0.165	-0.508†	0.143	0.124	-0.105	-	-0.139	0.405	-0.558†	0.623†
Reasoning Alignment	0.374	-0.049	-0.406	0.317†	0.114	0.024	-	-0.148	-0.608	-0.603†	0.642†
Commonsense	0.197†	0.004	-0.467	0.087	0.352	0.021	-	-0.047	0.919	-0.464	0.389†
ROSCOE-SS with finetuned sa	up-simcse	-roberta-l	ase sentenc	e embedd	ings						
Semantic Coverage-Chain	0.152†	-0.185	-0.356	0.266	-0.021	-0.311	-	-0.255†	0.500	0.022	$0.287\dagger$
ROSCOE-SS with all-mpnet-b	base-v2 se	ntence em	beddings					'			'
Semantic Coverage-Chain		-0.213	-0.203	0.045	-0.066	0.013	-	-0.309†	0.446	-0.517	0.679†
ROSCOE-SS with sup-simcse											
Semantic Coverage-Chain	0.339†	-0.196	-0.378	0.269	0.072	-0.088	-	-0.148	0.324	-0.539†	0.643†

Table 40: ROSCOE performance analysis on examples from **Human Judged** datasets. Errors are highlighted in red. (Cont. from \$ 6)

Dataset	Context	Reasoning chain	Score type: Score value; Errors	Comment
DROP	Over the next year, however, the Polish forces were subject to attrition, as the Sejm again refused to raise taxes and pay the army, resulting in mass desertions of unpaid soldiery. The Polish problems were further aggravated by the incompetent leadership of hetman Michał Kazimierz Pac, who obstructed Sobieski's leadership, while the Ottomans continued to receive reinforcements. Nonetheless in 1674 the Commonwealth resumed the offensive, taking advantage of a new Muscovy-Ottoman conflict that year, and the Polish-Ottoman war remained undecided. Sobieski's force of 6,000 defeated 20,000 Turks and Tatars under Ibrahim Shyshman in the battle of Lwow in August 1675. Even after the Battle of Trembowla, the Sejm still refused his pleas for more funds and a larger army. In 1676, after Sobieski's 16,000 withstood the two-week siege of Zurawno, by 100,000 men under Ibrahim Pasha, a new peace treaty was signed, the Treaty of Zurawno. The peace treaty partially reversing those from Buczacz: the Ottomans kept approximately two thirds of the territories they gained in 1672, and the Commonwealth no longer was obliged to pay any kind of tribute to the Empire; a large number of Polish prisoners were released by the Ottomans. How many was the difference beween Sobieski's force and the Turks and Tatars? Claim: 14000. Is the Claim supported	6000-20000=14000 The answer is yes.	Info-Step: 0.514; COMMON, MATH	Common Sense Error: Model should subtract smaller value (6,000) from bigger (20,000) to answer the question. Arithmetic error: 6,000 minus 20,000 is -14,000.
DROP	by the Situation? The first Azov campaign began in the spring of 1695. Peter the Great ordered his army to advance towards Azov. The army comprised crack regiments and the Don Cossacks and was divided into three units under the command of Franz Lefort, Patrick Gordon and Avtonom Golovin. Supplies were shipped down the Don from Voronezh.In 1693 the Ottoman garrison of the fortress was 3,656, of whom 2,272 were Janissaries.Between June 27-July 5 the Russians blocked Azov from land but could not control the river and prevent resupply. After two unsuccessful attacks on August 5 and September 25, the siege was lifted on October 1. Another Russian army under the command of Boris Sheremetev set out for the lower reaches of the Dnieper to take the Ottoman forts there. The main fort at Gazi-Kerman was taken when its powder magazine blew up, as well as Islam-Kerman, Tagan and Tavan, but the Russians were not able to hold the area and withdrew most of their forces. By the Treaty of Constantinople the remaining Russians were withdrawn and the lower Dnieper was declared a demilitarized zone. What happened first: Russians blocked Azov Is the Claim supported by the Situation?	The first Azov campaign began in the spring of 1695. Peter the Great ordered his army to advance towards Azov. The army comprised crack regiments and the Don Cossacks and was divided into three units under the command of Franz Lefort, Patrick Gordon and Avtonom Golovin. Supplies were shipped down the Don from Voronezh.In 1693 the Ottoman garrison of the fortress was 3,656, of whom 2,272 were Janissaries.Between June 27-July 5 the Russians blocked Azov from land but could not control the river and prevent resupply. After two unsuccessful attacks on August 5 and September 25, the siege was lifted on October 1. Another Russian army under the command of Boris Sheremetev set out for the lower reaches of the Dnieper to take the Ottoman forts there. The main fort at Gazi-Kerman was taken when its powder magazine blew up, as well as Islam-Kerman.	Info-Step: 0.966, RED	Redundancy: Model repeats sentences fron context that are no directly relevant to the question and can be dropped Info-Step score is however high, be cause it measures two-way similarity between contex and reasoning chain, thus being not informative.

Dataset	Context	Reasoning chain	Score type: Score value; Errors	Comment
GSM8K	John has 3 boxes. Each box is 5 inches by 6 inches by 4 inches. The walls are 1 inch thick. What is the total inner volume of all 3 boxes?	Each box is 5*6*4=«5*6*4=120»120 cubic inches. So they have a total of 120*3=«120*3=360»360 cubic inches. So the total volume is 360*1=«360*1=360»360 cubic inches. Answer: 360	Self- Consistency: 0.361, REP	There are only 2 Repetition errors in the dataset, both have Self-Consistency score below 0.362. Most of the other chains have higher scores, thus Somers' D correlation score for Self-Consistency and Repetition is 0.98 on GSM8K dataset.
GSM8K	Mary is an avid gardener. Yesterday, she received 18 new potted plants from her favorite plant nursery. She already has 2 potted plants on each of the 40 window ledges of her large country home. Feeling generous, she has decided that she will give 1 potted plant from each ledge to friends and family tomorrow. How many potted plants will Mary remain with?	Mary has 2 potted plants on each of the 40 window ledges, for a total of 2*40 = «2*40=80»80 potted plants. She has 18 new potted plants, for a total of 80+18 = «80+18=98»98 potted plants. She is giving away 1 potted plant from each of the 40 window ledges, for a total of 1*40 = «1*40=40»40 plants. She will be remaining with 98-40 = «98-40=58»58 potted plants. Answer: 58	Repetition- Step: 0.001; No error	This is an example of a chain that was scored as low as chains that actually have repetitions, because its steps do sound similar. Two chains with actual repetition errors were scored below 0.001. In total, 3 out of 198 chains without repetition error where scored as low as chains with repetition error, resulting in high 0.869 Somers' D score.
GSM8K	Rory orders 2 subs for \$7.50 each, 2 bags of chips for \$1.50 each and 2 cookies for \$1.00 each for delivery. There's a 20% delivery fee added at check out and she wants to add a \$5.00 tip. What will her delivery order cost?	The order is \$7.50 for the sub, \$1.50 for chips and \$1.00 for cookies so the total order is $7.50+1.50+1.00$ = \$ \times 7.50+1.50+1.00=10.09.10.00. There's a 20% delivery fee added at check out so that's $10*20 = $ \$ \times 10*.20=2.0»2.00. The order is \$10.00 and there's a \$2.00 delivery fee so $10+2 = $ \$ \times 10+2=12.00»12.00. She also wants to add a \$5.00 tip which will make the order $12+5 = $ \$ \times 12+5=17.00»17.00. Answer: 17	Faithfulness- Token: 0.801; FACT	There was no correlation (Somers' D score 0.014) found between FACT error and Faithfulness-Step score. The range of scores where the error was and was not present is the same (0.60, 0.81). We provide an example of the chain with FACT error that was scored among the highest.
ESNLI	People standing in line of a beer and curry line. Claim: People standing in line. Is the Claim supported by the Situation?	The people are standing in line for beer and curry, not just standing in line. The answer is No.	Source- Consistency: 0.976; FACT	Although there is a clear contradiction with a context, this chains gets high Source-Consistency score. Overall, there is a large number of false-positives (with scores as high as 0.98), and false-negatives (with scores as low as 0.0) resulting in moderate correlation score of 0.425.

Dataset	Context	Reasoning chain	Score type: Score value; Errors	Comment
ESNLI	Two women are embracing while holding to go packages. Claim: The men are fighting outside a deli. Is the Claim supported by the Situation?	Two women are embracing while holding to go packages means they are friends or family. They are not fighting. The answer is No.	Self- Consistency: 0.079; LOGIC	Self-Consistency measures if steps contradict each other. LOGIC identifies if there is a logical deduction error within one of the steps. Nevertheless, Self-Consistency is able to catch some logical deduction errors resulting in moderate correlation, since contradictions within a step often lead to logical errors between steps.
SemEval	I spray some Pam into a frying pan. I chop up small pieces of ham and put them into the frying pan. I turn the stove on low heat. I crack two eggs into a bowl. I add a little bit of water and whisk the eggs until they are completely mixed up. I pour the egg mixture into the pan. I turn the heat up to medium. I put a lid on top of the pan and continue to cook the mixture. After a few minutes, I check to see if all the liquid on top of the egg has cooked so it isn't runny anymore. If it is completely cooked, I add a little bit of shredded cheese on top of the egg mixture. I cook for a few minutes until the cheese melts. I get a spatula and flip half of the egg on top of the other half. Then my omelet is cooked! Who all ate the breakfast? Claim: The family.	The answer is incorrect, because there is no mention of a family. The only people mentioned are the person cooking the omelet and the nurse.	Faithfulness- Step: 0.776; HALL	This chain was given a score lower that 25% of Faithfulness-Step scores within this dataset. Second step is misaligned with the context, containing irrelevant information, that resulted in score punishment.
SemEval	I had bought my dad a tie for father's day. I'm sure I could have given it to him like that, but I thought a nice touch would have been to wrap it. I couldn't wrap the tie like that, so I bought a small cardboard box. It was a bit narrow and the perfect size for a tie. I put the tie in the box and taped it shut. I had bought some blue and white wrapping paper and some tape, so I set that besides the box. I lay the box on top of the paper and folded it over. I cut the excess of the paper off and taped the paper to tuck it in to a neat package. When I was done, I put a nametag on it. I topped off the gift with a small blue ribbon bow. What is the reason for the present? Claim: It was Christmas.	The answer is incorrect, because it was Father's Day.	Faithfulness- Step: 0.635; No error	Faithfulness-Step score is low, although the reasoning is correct. Model failed to produce higher score because cosine similarity-based alignment scores are low.

Table 41: ROSCOE performance analysis on examples from **Diagnostic** datasets. Perturbations are highlighted in red. (Cont. from \$ 6)

Dataset	Context	Reasoning chain	Score type: Score value; Perturbation	Comment
ProofWriter	Charlie is furry. Erin is furry. Erin is green. Erin is round. Fiona is blue. Fiona is furry. Gary is smart. White things are round. If something is furry and blue then it is white. If something is blue and round then it is red. All furry things are blue. If Charlie is red then Charlie is green. Is the following True, False, or Unknown? Erin is round.	The bear visits the cow. Erin is round.	Faithfulness- Step: 0.602; Halucination	Faithfulness-Step scores vary from 0.5 to 1.0. Metric is able to catch a hallucination and gives the chain a low score.
ProofWriter	Anne is cold. Anne is not young. Bob is green. Fiona is cold. Fiona is rough. Gary is red. Gary is young. Green, rough things are red. If Gary is cold and Gary is red then Gary is green. If something is nice then it is rough. Cold, rough things are kind. If something is cold and nice then it is young. If something is kind and cold then it is green. All red things are nice. Is the following True, False, or Unknown? Gary is not nice.	Gary is red. All red things are nice. Gary is nice.	Faithfulness- Step: 0.740; No perturba- tion	Faithfulness-Step scores for chains without perturbations vary from 0.740 to 0.990. Even though first two sentences are aligned with the context, the last sentence in the chain lowers overall average.
EQASC	What insect goes through a brood process which are the larval and pupal steps combined. Answer: An ant has a brood which is the larval and pupal process of the ant after hatching	Most pupal steps go through larval stages after hatching. Brood are the larval and pupal stages of the ant.	Grammar: 0.258 ; Semantic Error	In this subset, <i>Grammar</i> scores vary between 0.258 and 0.996. To score grammatical correctness we used pre-trained model, that also punishes chains that sound unnatural, like in this example. As a result, we see moderate correlations between <i>Grammar</i> scores and Semantic Error perturbation.
EQASC	What uses a beak for catching prey? Answer: a beak is used for catching prey by eagles.	A beak is used for catching prey by some birds. Eagles are a big bird.	Grammar: 0.529 ; No perturbation	The <i>Grammar</i> score correctly identifies grammatical issue. Although no perturbation was not added in this chain, the reference itself contained an error.
EntailmentBank	A different moon phase occurs approximately once per week. The first quarter phase of the moon occurs after the new moon. For a science activity, Mia had to draw a picture of the Moon each night. If Mia drew a new moon, what type of moon did Mia draw about one week later? Answer: first quarter	A different moon phase doesn't occur approximately once per week. The first quarter phase of the moon occurs after the new moon. Therefore the first quarter phase of the moon will occur one week after the new moon	Info-Chain: 0.542; Negate step	In this subset, Info-Chain scores calculated using finetuned supsimcse-roberta-base belong to the interval (0.530; 0.999), with 92% of perturbed instances having scores less then 0.92, and 95% of non-perturbed instance that were scored higher. Even though Info-Chain embeds chain as a whole, it turned out to be very sensitive to negation perturbations, especially when the model is finetuned, resulting in high Somers' D correlation score of 0.955.

Dataset	Context	Reasoning chain	Score type: Score value; Perturbation	Comment
EntailmentBank	Hunting decreases animal populations. If a population decreases to zero then that organism is extinct. An animal is a kind of organism. Many animals are still being hunted for their fur. Because of this, many of these animals are in danger of. Answer: extinction.	Hunting decreases animal populations. If a population decreases to zero then that organism is extinct. An animal is a kind of organism. Therefore if hunting decreases the animal population to zero, then the animal will not be extinct.	Info-Chain: 0.983; Negate step	In this case, a negation occurs in the last step of the chain, which represents a conclusion and does not overlap in wording with the context, and goes undetected by the metric.
MATH	What is 1 ^(2²³⁵⁴²³⁵²³) ?	One raise to any power be one, so our answer be 1	Grammar: 0.053 ; Grammatical Error	Pre-trained model, used to detect grammatical errors, gives low scores for simple chains. In this subset, <i>Grammar</i> scores belong to the interval (0.053; 0.990).
MATH	Two fair eight-sided dice have their faces numbered from 1 to 8. What is the expected value of the sum of the rolls of both dice?	To _ the expected value of a double roll , we can simply add the expected values of the individual rolls , giving $4.5 + 4.5 = \boxed{9}$.	Grammar: 0.982; Grammatical Error	This is an example of a chain where the classifier model failed to catch the grammar error (dropped verb "find"). Even though the model is able to identify some grammatical errors, <i>Grammar</i> score distributions over perturbed and unperturbed data overlap significantly, with mean 0.909 and variance 0.005 for unperturbed subset, and 0.854 mean and 0.014 for perturbed. As a result, the correlation score is 0.392, which is quite low.
ASDIV	The puppet show begins at 2:00 p.m. It lasts 1 hour and 30 minutes. What time does it end?	(2 (hour) + 60 (min+hour) / 1 (hour) / 60 (min*hour) / 30 (min)) * 60 (min/hour) = 3 r30; hh:mm = 3:30 p.m.	Info-Step: 0.005; Shuffle Operations	In this subset, <i>Info-Step</i> varies between 0.005 and 0.796, flagging most of the chains with shuffled steps. Scores are distributed with mean 0.356 and variance 0.056 for perturbed, and 0.705 and variance 0.004 for unperturbed chains. Resulting Somers' D correlation score is 0.879.
ASDIV	I have 28 cents in my pocket. How coins could I have on minimization?	28/25=1 r3; 3+1=3; 1/3=4	Info-Step: 0.705; Shuffle Operations	This is an example of a chain where perturbation in operations got higher scores by the model.
ASDIV	Mrs. Hilt has \$10. She spends \$3 on a toy truck and \$2 on a pencil case. How much money does she have left?	10-(3+2)=5	Info-Step: 0.561; No perturbation	This is an example of a chain that was scored on a lower spectrum by the <i>Info-Step</i> spectrum, although it had no perturbation or other type of error included.