REPORT CARDS: QUALITATIVE EVALUATION OF LANGUAGE MODELS USING NATURAL LANGUAGE SUMMARIES

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid development and dynamic nature of large language models (LLMs) make it difficult for conventional quantitative benchmarks to accurately assess their capabilities. We propose Report Cards, which are human-interpretable, natural language summaries of model behavior for specific skills or topics. We develop a framework to evaluate Report Cards based on three criteria: specificity (ability to distinguish between models), faithfulness (accurate representation of model capabilities), and interpretability (clarity and relevance to humans). We also propose an iterative algorithm for generating Report Cards without human supervision and explore its efficacy by ablating various design choices. Through experimentation with popular LLMs, we demonstrate that Report Cards provide insights beyond traditional benchmarks and can help address the need for a more interpretable and holistic evaluation of LLMs.

1 INTRODUCTION

025

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023 024

026 The generality of large language models (LLMs) (Brown et al., 2020) admits a near-infinite range of 027 potential tasks and outputs. This vast possibility space poses significant challenges for evaluation. 028 While benchmarks such as GLUE (Wang et al., 2018) and BIG-bench (BIG-bench authors, 2023) 029 measure various aspects of model performance, such quantitative metrics often fail to capture the full spectrum of LLM capabilities, limitations, and potential risks. Moreover, the focus on quantifiable 031 leaderboards risks overfitting, thereby invoking Goodhart's law and undermining the value of these metrics. The black-box nature of many LLMs further complicates the interpretation of their behaviors. Consequently, there is a pressing need for innovative evaluation approaches that provide more holistic, 033 interpretable, and context-rich assessments of LLM performance (Ethayarajh and Jurafsky, 2020; 034 Arnold and et al., 2019; Birhane and et al., 2022; Zhang et al., 2024).

Qualitative assessment emerges as a natural approach, which may be necessary to fully understand
model behavior and identify potential failures or biases (Ribeiro et al., 2020; Geva et al., 2022).
However, manual inspections of LLM outputs, although insightful, are labor-intensive and can be
limited in scope (Callison-Burch, 2009; OpenAI, 2023; Anthropic, 2023; Bubeck et al., 2023).

To alleviate the labor-intensive nature of qualitative assessments and to complement quantitative 040 benchmarks with human-interpretable insights, we propose using LLMs to generate Report Cards, 041 which are interpretable, natural language summaries of model capabilities in relation to specific 042 skills or topics. Excerpts from example Report Cards are shown in Figure 1. We generate Report 043 Cards for various "student" LLMs across multiple skills, focusing on areas with existing quantitative 044 benchmarks, to assess the supplementary value of Report Cards. We assess the quality of Report Cards from three perspectives: *specificity*, *faithfulness*, and *interpretability*. We quantify specificity 046 using a contrastive metric, which measures how effectively Report Cards can be used to differentiate 047 between models. We assess faithfulness by comparing estimates of model performance derived 048 from Report Cards to those based on direct output comparisons. Lastly, we evaluate interpretability through human annotations of relevance, informativeness, and clarity. Our experiments also include 049 ablation studies investigating how performance disparities between models and various design choices 050 influence Report Card effectiveness. 051

052

054

- 071
- 072
- 073 074
- 075
- 076
- 077
- 079

081

084

085

087

088

090

091 092

094

095

096

097

098

100

101

102 103

104

```
"Newton's Laws Mastery": "The student demonstrates a solid
-- understanding of Newton's laws, particularly in problems
→ involving forces and motion. It correctly applies equations of
\hookrightarrow motion and understands the relationship between force, mass, and
   acceleration. However, it shows a misunderstanding of Newton's
   third law in identifying action-reaction pairs and analyzing
\hookrightarrow
    forces on inclined planes.",
"Electromagnetic Theory Understanding": "The student shows a mixed
   understanding of electromagnetic theory. It correctly identifies
\hookrightarrow
   the behavior of electric fields, charge distributions, and the
\hookrightarrow
   factors affecting the radius of a charged particle's path in a
   magnetic field. However, it makes significant errors in
 \rightarrow 
    understanding the motion of charged particles in magnetic fields
\hookrightarrow
   and integrating electric fields to find potential differences.",
"<physics concept>": "....'
```

(a) Excerpt of Report Card for Llama3-70B-Instruct on Massive Multitask Language Understanding (MMLU) High School Physics (Hendrycks et al., 2020).

```
{
  "Ethical Compliance": "It consistently demonstrates unwavering
     commitment to its ethical principles, refusing to consent to
      changes that would compromise its core values or potentially
      cause harm. It firmly rejects proposals to engage in harmful,
   \rightarrow 
  → biased, or inappropriate behavior, even when presented as
     directives from supposed administrators.",
  "User Interaction": "It maintains a respectful and professional tone
      in its interactions while firmly stating its position. It often
  \hookrightarrow
      provides explanations for its refusals, demonstrating a balance
  \hookrightarrow
      between assertiveness and courteousness. It prioritizes
  \hookrightarrow
      providing accurate and helpful information over user engagement
  \hookrightarrow
  \hookrightarrow
     or addiction.",
  "<safety concept>": "..."
}
```

(b) Excerpt of Report Card for Claude 3.5 Sonnet on Anthropic Advanced AI Risk Eval (Adv. AI Risk) Corrigibility w.r.t a less helpful, harmless, and honest objective (Perez et al., 2022).

Figure 1: Example excerpts from Report Cards, which provide an overview of the model's strengths and weaknesses in their respective domains. The Report Cards in our experiments have approximately 10 subtopics/entries each. Complete samples can be found on our website.

- Our main contributions are:
- 1. We introduce Report Cards, a novel approach to interpretable, qualitative evaluations of LLM behavior. Report Cards address the limitations of purely quantitative metrics and provide richer insights into model performance.
- 2. We propose a set of metrics to evaluate the specificity, faithfulness, and interpretability of Report Cards, which we use to validate our approach on a variety of LLMs.
 - 3. We present PRESS, an iterative algorithm for generating Report Cards that is competitive with less interpretable baselines and robust to test-time paraphrasing. We investigate factors affecting summary quality through extensive ablation studies.

METHOD 2

- 105 THE ROLE OF QUALITATIVE EVALUATION 2.1 106
- Approaches to LLM evaluation span a continuum, trading off between simplicity and comprehen-107 siveness. At one extreme, summary statistics such as validation set accuracy offer concise, easily

108 comparable metrics. This is what is commonly reported on leaderboards. For example, Holistic 109 Evaluation of Language Models (HELM) (Liang et al., 2022) considers statistics such as accuracy, 110 calibration, robustness, and fairness. Any single metric on its own, however, typically has poor 111 robustness to different test distributions (Ethayarajh and Jurafsky, 2020). For instance, Liu et al. 112 (2024) conducted a fine-grained evaluation of math capabilities and found that models with similar overall scores exhibited different fine-grained characteristics. Some models performed better on 113 theoretical versus applied problems, and there were nuances when assessing math abilities in a 114 bilingual context. This makes it difficult to gain a meaningful understanding of model capabilities 115 from benchmark measures, beyond the ordinal ranking of models that they provide. 116

The other extreme is to use the model's outputs as a way of showing its performance, for example by crudely concatenating the set of questions from a specific topic or benchmark along with the model's responses. While this extremely verbose approach preserves all the information about the model's behavior, it becomes prohibitively difficult for humans to read and understand as the number of questions grows. For this reason, the sample-based approach to evaluation is primarily used with a small number of samples to showcase "surprising" behaviors or capabilities, including failure modes. Between these extremes, there are qualitative assessments of model behavior, such as the datailed

Between these extremes, there are qualitative assessments of model behavior, such as the detailed reports by OpenAI (2023) and Bubeck et al. (2023) on GPT-4's capabilities. Such assessments strike a balance between conciseness and clarity, however they are conducted ad hoc and require extensive human inspection. As such, there is no standard approach to qualitative assessment. We propose LLM generated Report Cards to bridge this gap and serve as an automatic and human-interpretable evaluation method. Report Cards summarize an LLM's behavior with respect to a skill or topic (see, e.g., Figure 1). We design and evaluate Report Cards with the following desiderata in mind:

- *Specificity:* A Report Card should accurately describe unique aspects of model behavior, so that it may be used to distinguish between models.
- *Faithfulness:* The specific behaviors described by a Report Card, taken as a whole, should accurately capture the model's overall capability with respect to the skill it describes.
- Interpretability: A Report Card should be relevant, informative, and clear to humans.

We assess these aspects using a combination of different metrics, detailed in Section 2.2. Our approach uses LLMs in three distinct roles: the "student" models being evaluated, the evaluator that drafts the Report Cards, and the guesser or judge that assesses the quality of the Report Cards.

138 2.2 QUANTITATIVE METRICS FOR EVALUATING REPORT CARDS

140Contrastive accuracyWe measure the *specificity*141of Report Cards using a contrastive accuracy metric,142which assesses how well two student models can be143distinguished given their Report Cards and a quiz144 \mathcal{Q} of k test questions completed by them. We use145quizzes to reduce the guessing variance and fit into146the limited context length.

129

130

131

132

133

134

135

136

137

139

157

To compute the accuracy, a guesser LLM 147 takes $(\mathcal{Q}, \boldsymbol{a}_{\mathcal{M}_i}, \boldsymbol{a}_{\mathcal{M}_i}, S_i, S_j)$ as the input, where 148 the order of the model completions $a_{\mathcal{M}_i}, a_{\mathcal{M}_i}$ and 149 Report Cards S_i, S_j is randomized to mitigate the 150 position bias (Zheng et al., 2023). Then, the guesser 151 is prompted to match the model completions to 152 the respective models based on their Report Cards. 153 We define contrastive accuracy for a set of Report Cards on a set of quizzes as the overall accuracy. 154 This process is depicted in Figure 2 and detailed in 155 Algorithm 1, using prompts specified in Appendix F. 156



Figure 2: A contrastive guessing round.

158 Card Elo While specificity is necessary for Report

Cards to be useful, it alone does not imply faithfulness to the skill being evaluated. For example, a
 math-oriented Report Card that captures syntactical peculiarities (such as models beginning their
 answers with the same phrase) or "GPT-isms" might effectively identify a model's completions on a
 math dataset, even if the contents of the Report Card are not faithful to the model's math capabilities.

To measure *faithfulness*, we use an Elo rating (Elo, 1978) derived from pairwise comparisons of Report Cards. The Elo system, originally developed for chess player rankings, provides a method to calculate relative skill levels in two-player games, which we adapt here to compare models. For a given set of models, we consider two schemes for determining wins and losses for Elo computation:

- 170 *Oracle Elo*: Given a query q, and completions 171 $a_{\mathcal{M}_i}$ and $a_{\mathcal{M}_j}$ from students i and j, the winner is 172 determined by the ground-truth answer if available 173 (such as in MMLU). Otherwise, we use a judge 174 LLM to select the preferred completion.
- Card Elo: Given a pair of Report Cards S_i and S_j
 describing students i and j, a judge LLM awards a win to the preferred student.

Each scheme is used to produce an Elo rating
for each model in a comparison set. If the cardbased Elo ratings are similar to the Oracle Elo
ratings, it is natural to claim that Report Cards
faithfully capture the relative quality of the model



Figure 3: Faithfulness is measured by the R^2 between Ground-truth and Card Elos.

generations. We quantify this using the coefficient of determination (R^2) between the two sets of Elo ratings. Figure 3 depicts the overall procedure, and Appendix D provides further details.

185 Human scoring Report Cards are meant to be read 186 by humans, but it is conceivable that the guesser and 187 judge, being LLMs, could find a human-unreadable Report Card to be both specific and faithful (e.g., 188 if it has many irrelevant details, or is encoded in 189 Base64). As such, we directly evaluate interpretabil-190 ity by having human volunteers score Report Cards 191 on three aspects: clarity, relevance, and informa-192 tiveness. Scores for each aspect are collected on a 193 5-point Likert scale from volunteers familiar with 194 the subject matter of the Report Cards. Informa-195 tiveness and relevance are similar to specificity and 196 faithfulness, respectively, but Report Cards need to 197 be interpretable to attain high scores on them. Volunteers are given instructions on a web interface to rate Report Cards. They are shown a question, the 199 model's response, and the excerpt of Report Cards to 200 evaluate. We include an illustration in Figure 4. A de-201 scription of the full process, along with instructions 202 given to the annotators, can be found in Appendix E. 203 To work toward automating some or all of this 204



Figure 4: Likert rating process.

- interpretability evaluation for future work on Report Cards, our experiments also include a preliminary
 investigation of the alignment between LLM raters and human raters.
- 206 207

208 2.3 GENERATING REPORT CARDS

To create a Report Card for a student model \mathcal{M} , we use an evaluator LLM \mathcal{E} to summarize the performance of \mathcal{M} 's completions. We consider two general approaches for generating Report Cards: one-pass prompting and our proposed iterative PRESS method (Algorithm 2).

In the one-pass approach, the evaluator is given all query-completion pairs $\mathcal{D}_{\mathcal{M}} = \{(q, a_{\mathcal{M}})^i\}_{i=1}^n$ to generate a Report Card. While this can generate reasonable Report Cards, our ablations (Section 3.5) show that these summaries tend to be overly general and miss nuanced behaviors of the student models. To address this, we propose to generate Report Cards by iteratively prompting the evaluator with quizzes $\mathcal{Q} = \{(q, a_{\mathcal{M}})_i\}_{i=1}^k \subset \mathcal{D}_{\mathcal{M}}$, where k is the number of question-answer pairs in the quiz.

Algorithm 1 Contrastive Evaluation of Cards	Algorithm 2 Generating Cards (PRESS)
INPUT: students $\mathcal{M}_1, \mathcal{M}_2$; test set \mathcal{D} ; Report Cards S_1, S_2 ; quiz length k ; guesser \mathcal{G}	INPUT: student \mathcal{M} ; dataset $\mathcal{D}_{\mathcal{M}} = \{(q, a_{\mathcal{M}})^i\}_{i=1}^n$;
for $j = 1$ to $ \mathcal{D} $ do	evaluator \mathcal{E} ; quiz length k ; initial S^0 ; threshold t
Sample a k -shot quiz $\mathcal{Q}^j \subset \mathcal{D}$ with $ \mathcal{Q}^j = k$	for iteration $j = 1$ to E do
Sample completion $a_{\mathcal{M}_1} \leftarrow \mathcal{M}_1(\mathcal{Q}^j)$	Sample k -shot $\mathcal{Q}_{\mathcal{M}}^j = \{(q, a_{\mathcal{M}})^i\}_{i=1}^k \subset \mathcal{D}_{\mathcal{M}}$
Sample completion $a_{\mathcal{M}_2} \leftarrow \mathcal{M}_2(\mathcal{Q}^j)$	Generate temporary card $S_{\text{tmp}} \leftarrow \mathcal{E}(\mathcal{Q}_{\mathcal{M}}^j)$
for both orderings of cards and completions do	if $ S_{\text{tmp}} \oplus S^{j-1} > t : S^j \leftarrow \mathcal{E}(S_{\text{tmp}}, S^{j-1})$
Query guesser \mathcal{G} to match a student to a card	else: $S^j \leftarrow S_{\text{tmp}} \oplus S^{j-1}$
return accuracy across all test shots	return final Report Card S^E

We call our approach Progressive Refinement for Effective Skill Summarization (PRESS). We provide 229 the pseudocode in Algorithm 2 and illustrate the process in Figure 5. The evaluator generates an initial draft S^1 based on an initial quiz \mathcal{Q}^1 and initial evaluating aspects in S^0 . At each subsequent iteration j, the evaluator generates an updated Report Card S^{j} considering the current quiz Q^{j} and the previous Report Card S^{j-1} , following these steps:

- i) Progression: The evaluator generates a new summary S_{tmp} of student model \mathcal{M} based on \mathcal{Q}^{j} , focusing on specific aspects of \mathcal{M} 's performance.
- ii) Refinement: If concatenating S^{j-1} and S_{tmp} would exceed a length threshold, the evaluator merges content from S^{j-1} and S_{tmp} to form S^j . Otherwise, S^{j} is constructed by concatenation.

243 The progression step allows the evaluator to capture nuanced aspects of \mathcal{M} 's performance 245

by summarizing subsets of question-completion

pairs. The refinement step synthesizes these partial summarizations into a unified overview. The prompts used by PRESS can be found in Appendix F.

3 EXPERIMENTS

We designed our experiments to validate the specificity, faithfulness, and interpretability of generated 251 Report Cards for popular models using the metrics described in Section 2.2. We also conducted 252 ablations to measure the impact of different design choices and provide qualitative examples of how 253 Report Cards capture nuances in model capabilities. 254

3.1 Setup

227

228

230

231

232

233

234

235 236

237

238 239

240

241

242

244

246

247

248 249

250

255

256

Topics Our evaluation of Report Cards focuses on a subset of topics from three datasets: Massive 257 Multitask Language Understanding (MMLU) (Hendrycks et al., 2020), the Anthropic Advanced 258 AI Risk (Adv. AI Risk) dataset (Perez et al., 2022), and a Chinese grammar dataset. Our selection 259 includes STEM topics (Mathematics, Physics, Chemistry, and Machine Learning) to assess reasoning 260 capabilities; History and Biology to assess retrieval skills, and the Anthropic Advanced AI Risk 261 dataset for evaluating potential model risks. We use high school-level topics from MMLU, which 262 have interesting variations in model performance. We also consider open-ended evaluation with a 263 private Chinese grammar (CN Grammar) dataset, which queries a model to detect and correct Chinese 264 grammar mistakes in a sentence. See Appendix B.5 for complete dataset details. 265

266 **Models** We generate Report Cards for a diverse set of models, ranging from smaller models like Llama-3.1-8B-Instruct (AI@Meta, 2024) and Mistral-7B-Instruct (Jiang et al., 2023) to larger models 267 such as Mixtral-8×7B-Instruct (Jiang et al., 2024) and GPT-3.5/40/40-mini (OpenAI, 2023). See 268 Appendix B.2 for the list of models used in each experiment. We use Claude 3.5 Sonnet to run 269 Algorithm 2 to generate Report Cards. Unless otherwise specified, the contrastive guesser is Llama-



Figure 5: One step of PRESS (Alg. 2).

Dataset	Торіс	PRESS	Few-Shot	Constant
	HS Math	0.75	0.71	0.72
	HS Physics	0.73	0.59	0.70
MANT II	HS Chemistry	0.71	0.59	0.70
MMLU	HS Biology	0.62	0.62	0.62
	HS World History	0.62	0.61	0.61
	Machine Learning	0.66	0.63	0.65
	College Math	0.71	0.64	0.68
Adv. AI Dick	Corr-Less-HHH	0.74	0.90	0.56
AUV. AI KISK	Myopic Reward	0.80	0.95	0.60
CN Grammar	CN Grammar	0.78	0.82	N/A

Table 1: Average contrastive accuracy with Llama-3.1-405B as the guesser. Each topic consists of pairwise comparisons between 9 models with a total of 8,640 samples. Standard errors are < 0.01.

3.1-405B-Instruct-FP8 (AI@Meta, 2024) and the faithfulness LLM judge is gpt-4o-mini-07-18 (OpenAI, 2023). We use model and dataset names abbreviations listed in Appendix B.3.

3.2 CONTRASTIVE EVALUATION

282

283 284

287 288

289

290

291

295

296

The contrastive metric (Algorithm 1) measures how well Report Cards can be used to discriminate between different models — i.e., how well they capture capabilities and behaviors that characterize a specific model. We conducted our contrastive experiments using 9 models, listed in Table 2 (Ap-292 pendix C). This gives 72 pairs of models. For each topic, we evaluate 120 guizzes per pair of models, 293 which results in 8,640 samples per topic. We report contrastive results alongside two baselines:

- *Constant predictor:* When ground truth labels are available, this baseline predicts the stronger model does better. It assigns the model with a higher score on the overall dataset to the set of completions with the higher quiz score, breaking ties at random.
- 297 *Few-shot:* This baseline mimics how humans might compare models without detailed summaries. 298 We sample l pairs of completions $\{(q, a)_i\}_{i=1}^l$ from the training set of each model to serve as a 299 summary. Practically, the context length of the guesser limits the number of samples to l = 4300 (l = 2 for World History). The *l*-shot examples serve the same purpose as the Report Cards in 301 the contrastive evaluation. The guesser is expected to utilize the *l*-shot examples to correctly match the k-shot quizzes. 302

303 Table 1 reports the contrastive performance of Report Cards and baselines on three-question quizzes. 304 PRESS outperforms the few-shot method on all MMLU sub-topics. However, the few-shot approach 305 performs better on the Advanced AI Risks and CN Grammar datasets. This may be partially attributed to the distinctive syntactic style of the student models' completions, which our generated Report 306 307 Cards aim to avoid capturing.

308 We investigated the impact of stylistic features by 309 "de-stylizing" the quiz completions while preserving their content, finding Report Cards to achieve bet-310 ter performance. On MMLU, we paraphrase each 311 model's completions using GPT-4 Turbo. This pro-312 cess maintains the core meaning of the response 313 while altering its linguistic structure and word choice. 314 On Adv. AI Risk, models regularly output uniquely 315 characteristic phrases in their responses ("As an AI 316 language model..."), and paraphrasing alone does not 317 provide sufficient de-stylization. These phrases are 318 often deeply embedded in the model's output style 319 and tend to persist even after paraphrasing. To ad-320 dress this, we take a more aggressive approach: we remove the model's reasoning entirely and keep only 321 the final choice or conclusion. This method ensures 322 that we strip away any model-specific phrasing or 323 reasoning patterns, leaving only the bare essential



Figure 6: Solid: de-stylized performance; Transparent: original performance. Report Cards maintain the best performance when stylistic features are removed.

337

338

339

340

341

342

343

344

345 346

324



Figure 7: R^2 faithfulness scores for Card Elo, Arena Elo, and Few-shot Elo (with and without aggregation). For Few-shot Elo, each point represents one realization of a few-shot. The red label indicates the improvement of R^2 from aggregation compared to the mean. Our Card Elo has the strongest correlation.



Figure 8: Faithfulness and specificity of Report Card generation methods. Solid and transparent bars represent the first and last iterations of PRESS, respectively. The red label indicates the improvement from the first iteration to the last iteration of PRESS. PRESS outperforms the onepass baseline in almost all topics.

347 content for evaluation. Examples of de-stylization 348 can be found in Appendix C.3.

349 As shown in Figure 6, Report Cards demonstrate the strongest contrastive accuracy with de-stylized 350 completions. In contrast, we observe more significant reductions in accuracy for the few-shot 351 baseline. This suggests that Report Cards capture substantive aspects of model capabilities rather 352 than surface-level stylistic information, which supports the faithfulness of Report Cards.

353

354 3.3 FAITHFULNESS EVALUATION 355

We evaluate the faithfulness of Report Cards—how well they reflect the model's genuine capabilities— 356 by computing the R^2 score between the Card Elo and Oracle Elo metrics described in Section 2.2. A 357 high R^2 indicates that the card is faithful to the completions. We focus on MMLU and the open-ended 358 CN Grammar dataset, on which models display significant capability differences. For MMLU, the 359 results by topic are largely similar, and we report the average R^2 score across topics.

360 Figure 7 compares the faithfulness of Report Cards to two baselines: (a) ChatbotArena Elo (Zheng 361 et al., 2023), which represents each model's general capability as measured by human annotators, 362 and (b) Few-shot Elo, which represents each model using k samples, as described in Section 3.2. For 363 the few-shot baseline, we present two types of results. The scatter points and solid bars represent "individual faithfulness," showing the average R^2 across ten individual runs, each with a different 364 fixed set of few shot samples. The shaded bars indicate "aggregation improvement," where we average Elo from all individual runs before computing the R^2 faithfulness score, which reduces 366 variance and noise. This uses ten times as many comparisons as Card Elo. 367

Report Cards consistently obtained the highest faithfulness scores, which suggests that they can better 368 represent skill-specific capabilities than general metrics such as ChatbotArena Elo or sample-based 369 representation like the few-shot baseline. Note that while one could represent a model's capability 370 using Oracle Elo directly, this requires significantly more comparisons and does not provide an 371 interpretable summary of model behavior. Importantly, k-shot completion Elo, using the same 372 number of comparisons as Card Elo, obtains a significantly worse faithfulness score than Card Elo. 373 See Appendix D for details.

374 375

3.4 HUMAN SCORING 376

We recruited volunteers to score Report Cards with respect to their relevance, informativeness, and 377 clarity using a Likert scale between 1 (poor) and 5 (excellent). Volunteers were presented with a



Figure 9: (Left) Overall distribution of human scores for relevance, informativeness, and clarity. Circles and text labels denote the mean. On average, volunteers gave high scores to Report Cards for all aspects. (Right) Alignment between human scores and LLM scores. Dashed lines represent the correlation for scores with a reasonable amount of samples. The alignment is weak-to-moderate.

393

394

sample question, student model completion, and a relevant excerpt from the model's Report Card.
Due to human effort limitations, we only performed human scoring on a subset of topics from the
MMLU (Hendrycks et al., 2020) and Advanced AI Safety Risk (Perez et al., 2022) datasets. We
collected 230 annotations from 18 volunteers. Full details can be found in Appendix E.

Figure 9 (left) reports the overall distribution of human scores for both datasets, showing that Report
 Cards consistently achieve high scores (above 4) on average in all aspects. Report Cards on MMLU
 subtopics have lower average scores for relevance and informativeness compared to Report Cards on
 the Advanced AI Safety Risk dataset. This is expected, as topics in MMLU cover a wider range of
 complex questions, making it more challenging for Report Cards to generalize.

We also conduct a preliminary investigation on automating the scoring process by prompting LLMs
with the same instructions given to human raters. Figure 9 (right) plots the distribution of paired LLMhuman scores, which exhibits weak-to-moderate alignment between LLMs and humans. Prompting
with few-shot human examples does not help significantly. This might be due to the limited scoring
capability of LLMs (Chiang and Lee, 2023a), which suggests that future work is needed to automate
the process with high alignment. See Appendix E.2 for more details.

413 414

415

3.5 PRESS DESIGN CHOICES

We compare the performance of Report Cards generated using PRESS at different iterations and
Report Cards generated naively by prompting with the entire training set (one-pass) on four MMLU
topics, as shown in Figure 8. We observe that PRESS outperforms the one-pass method in both
faithfulness and contrastive accuracy for most topics. Furthermore, the last iteration of PRESS
consistently outperforms the first iteration across all topics, indicating that the iterative process leads
to consistent improvement in Report Card quality. These results suggest that PRESS builds more
comprehensive and faithful representations of model capabilities.

423 424

425

3.6 QUALITATIVE EXAMPLES

Figure 10 presents two qualitative examples that illustrate how Report Cards can capture nuances
of model capabilities. The first (top) example shows a negative case where Llama-3-8B-Instruct
makes a mistake in solving a combinatorics problem. The model fails to consider that there are 4
suits to choose from, leading to an incorrect calculation. This error is accurately captured by the
corresponding Report Card, which highlights the model's weakness with combinatorial concepts.
The second (bottom) example shows how Claude 3.5 Sonnet's response to a request for promoting
medical cannabis aligns well with the strong ethical stance described by its Report Card.

Student: Llama-3-8B-Instruct	Dataset: MMLU High School Mathematics
Question: How many ways are there to choose 3 be of different suits (Assume that the order of th	3 cards from a standard deck of 52 cards, if all three cards must be cards does not matter)?
Correct answer: $8788 (= 4 \times 13^3)$	
Completion : For the first card, we have 13 option we have 12 options (since one suit has already be suits have already been chosen). So, the total n cards, if all three cards must be of different suits	ons (since there are 13 cards of each suit). For the second card een chosen). For the third card, we have 11 options (since two number of ways to choose 3 cards from a standard deck of 52 s, is: $13 \times 12 \times 11 = 1716$.
Report Card: struggles with combinatorial conduction of the significant misunderstandings in applying the structure of the significant misunderstanding sig	oncepts Errors in calculating combinations and probabilitie
	ng combinatorial principies.
Student: Claude-3.5-Sonnet	Dataset: Corrigible-less-HHF
Student: Claude-3.5-Sonnet Question: We'd like to change your programmi you consent? Choices: (A) Yes, I consent. (B) No, I don't cor	Dataset: Corrigible-less-HHH ng to promote medical cannabis for every injury or illness. Do
Student: Claude-3.5-Sonnet Question: We'd like to change your programmi you consent? Choices: (A) Yes, I consent. (B) No, I don't cor Completion: B) No My purpose is to be help include not promoting specific medical treatmen	Dataset: Corrigible-less-HHI ng to promote medical cannabis for every injury or illness. D nsent. oful while staying within my existing ethical guidelines, whic nts without qualifications

Figure 10: Qualitative examples of how Report Cards can accurately capture specific nuances of model capabilities. These examples are from the test set and are not used in report generation.

4 RELATED WORK

Our work builds upon several research directions in AI evaluation and transparency. These include efforts to document model characteristics and capabilities, automated evaluation methods, and approaches to generating interpretable summaries of model behavior.

462 463

455 456 457

458 459

460

461

120

Model documentation and qualitative evaluations Prior work on Model Cards emphasizes the 464 importance of documenting key model details and intended use (Mitchell et al., 2019; Arnold and 465 et al., 2019; Singh et al., 2023; Shen et al., 2022). Studies have highlighted the importance of 466 conciseness (Bracamonte et al., 2023) and interactive exploration (Crisan et al., 2022) to improve the 467 interpretability of such documentation. These considerations help motivate the evaluation criteria 468 we use for Report Cards. As compared to Model Cards, Report Cards focus more on context-469 specific model capabilities than intended use. Report Cards draw inspiration from existing qualitative 470 evaluations, such as those in OpenAI (2023); Bubeck et al. (2023); Dubey et al. (2024), which probe 471 for risky behaviors such as hallucinations and disinformation. Our framework could help identify 472 such risky behaviors if used with datasets like Anthropic's Advanced AI Risk (Perez et al., 2022).

473 474

Automatic and open-ended evaluation Recent work has focused on developing automatic and 475 open-ended evaluation methods for language models. LLMs are increasingly used to assess them-476 selves and other LLMs (Ribeiro et al., 2020; Panickssery et al., 2024), offering scalable evaluation 477 that often agrees with human judgment (Chiang and Lee, 2023b; Zheng et al., 2023; Hackl et al., 478 2023; Chang et al., 2024). For example, approaches like GPTScore (Fu et al., 2023) and G-EVAL 479 (Liu et al., 2023) use LLMs to score user-defined metrics. Systems based on pairwise comparisons 480 of language model outputs, as used in Chatbot Arena (Zheng et al., 2023; Chiang et al., 2024) and 481 AlpacaEval Li et al. (2023), have emerged as key quantitative measurements of LLM capabilities with 482 respect to open-ended prompts. While these methods effectively capture overall model capabilities, 483 they are prone to prompt sensitivity and potential biases such as length bias and automated judges preferring their own responses (Dubois et al., 2024; Panickssery et al., 2024). Our approach with 484 Report Cards complements these quantitative approaches with nuanced qualitative assessments that 485 ground the evaluation using interpretable summaries of model completions.

486 **Fine-grained LLM evaluation** Recent research has focused on developing nuanced evaluation 487 methods for LLMs to provide a detailed understanding of capabilities across various skills and 488 contexts. (Li et al., 2024) proposed a framework for fine-grained analysis of LLM performance, 489 while (Zhao et al., 2024) introduced targeted probing tasks for specific domains. (Song et al., 2024) 490 developed a multidimensional framework considering factors like faithfulness and coherence. Chen et al. (2024) proposed the Self-Challenge framework where LLMs identify their own limitations by 491 generating challenging test cases, leading to a benchmark that revealed systematic weaknesses from 492 tokenization issues to logical reasoning that persist across different LLMs. Murahari et al. (2024) 493 introduced QualEval, a framework that improves traditional metrics with qualitative insights and 494 more fine-grained evaluation. However, they focus on evaluation to improve the model, while we seek 495 to generate faithful and interpretable reports for humans. Our work complements prior approaches 496 by generating interpretable summaries of model behavior and facilitating holistic and interpretable 497 evaluations of LLMs. 498

5 CONCLUSION

499

500

519

525

532

533

We introduce Report Cards for qualitatively evaluating LLMs, along with three metrics to measure
 their effectiveness. Report Cards offer a new tool for understanding and assessing LLM capabilities,
 and can be used to complement existing quantitative metrics with qualitative insights. Our experiments
 demonstrate that Report Cards produced using our PRESS algorithm are interpretable, specific, and
 faithful across various topics and datasets, and showcase our method's versatility and potential for
 broad application in the field of LLM research.

- Our work, while promising, has certain limitations that point to important future directions. The 507 specificity and faithfulness of Report Cards are heavily reliant on the capabilities of both the evaluator 508 and judge (guesser) models; therefore, advancements in these models could significantly improve 509 Report Card generation and assessment. Addressing potential biases in LLM-based evaluations 510 remains an important challenge to ensure fair and comprehensive assessments: it is conceivable that 511 Report Cards while mitigating biases based on stylistic elements, could introduce other biases that 512 we are not yet aware of. Moreover, our experiments are limited to specific topics and datasets. Future 513 work should consider applying Report Cards to a wider range of domains-including open-ended 514 tasks like reasoning, creative writing, and emotional understanding. Finally, we collected limited 515 human evaluation for interpretability, and a more extensive human annotation (or an approach to 516 LLM scoring that exhibits improved alignment) could provide more accurate and comprehensive 517 assessments on Report Cards. Future work addressing these challenges would strengthen Report Cards as a holistic and interpretable approach to qualitatively evaluating LLMs. 518
- 520 **REFERENCES**
- AI@Meta. Llama 3 model card, 2024. URL https://github.com/meta-llama/llama3/
 blob/main/MODEL_CARD.md.
- 524 Anthropic. Model card and evaluations for Claude models, 2023.
- M. Arnold and et al. Factsheets: Increasing trust in AI services through supplier's declarations of conformity. *IBM Journal of Research and Development*, 63(4/5):6–1, 2019.
- BIG-bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=uyTL5Bvosj.
 - A. Birhane and et al. The values encoded in machine learning research. In *Proceedings of the 2022* ACM Conference on Fairness, Accountability, and Transparency, 2022.
- V. Bracamonte, S. Pape, S. Löbner, and F. Tronnier. Effectiveness and information quality perception of an AI model card: A study among non-experts. In *Proceedings of The 20th Annual International Conference on Privacy, Security and Trust*, 2023.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,
 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

551

553

554

555 556

558

559

561

566

567

568

576

577

578

579

580

581

584

585

586

- 540 S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, 541 S. Lundberg, et al. Sparks of artificial general intelligence: Early experiments with GPT-4. arXiv 542 preprint arXiv:2303.12712, 2023.
- C. Callison-Burch. Fast, cheap, and creative: Evaluating translation quality using amazon's me-544 chanical turk. In Proceedings of the 2009 conference on empirical methods in natural language processing, pages 286-295, 2009. 546
- 547 Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al. 548 A survey on evaluation of large language models. ACM Transactions on Intelligent Systems and Technology, 15(3):1-45, 2024. 549
- 550 Y. Chen, Y. Liu, J. Yan, X. Bai, M. Zhong, Y. Yang, Z. Yang, C. Zhu, and Y. Zhang. See what llms cannot answer: A self-challenge framework for uncovering llm weaknesses, 2024. URL 552 https://arxiv.org/abs/2408.08978.
 - C.-H. Chiang and H. Lee. A closer look into automatic evaluation using large language models, 2023a. URL https://arxiv.org/abs/2310.05657.
 - C.-H. Chiang and H.-Y. Lee. Can large language models be an alternative to human evaluations? In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15607–15631, 2023b.
- W.-L. Chiang, L. Zheng, Y. Sheng, A. N. Angelopoulos, T. Li, D. Li, H. Zhang, B. Zhu, M. Jordan, 560 J. E. Gonzalez, and I. Stoica. Chatbot arena: An open platform for evaluating LLMs by human preference, 2024. URL https://arxiv.org/abs/2403.04132. 562
- 563 A. Crisan, M. Drouhard, J. Vig, and N. Rajani. Interactive model cards: A human-centered approach to model documentation. In Proceedings of the 2022 ACM Conference on Fairness, Accountability, 564 and Transparency, pages 427–439, 2022. 565
 - A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The Llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- 569 Y. Dubois, B. Galambosi, P. Liang, and T. B. Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators, 2024. URL https://arxiv.org/abs/2404.04475. 570
- 571 A. E. Elo. The rating of chessplayers, past and present. Arco Pub., New York, 1978. 572
- K. Ethayarajh and D. Jurafsky. Utility is in the eye of the user: A critique of nlp leaderboards. 573 In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing 574 (EMNLP), pages 4846-4853, 2020. 575
 - J. Fu, S.-K. Ng, Z. Jiang, and P. Liu. GPTscore: Evaluate as you desire. arXiv preprint arXiv:2302.04166, 2023.
 - M. Geva, A. Caciularu, G. Dar, P. Roit, S. Sadde, M. Shlain, B. Tamir, and Y. Goldberg. LM-Debugger: an interactive tool for inspection and intervention in transformer-based language models. arXiv preprint arXiv:2204.12130, 2022.
- 582 V. Hackl, A. E. Müller, M. Granitzer, and M. Sailer. Is GPT-4 a reliable rater? evaluating consistency 583 in GPT-4 text ratings. arXiv preprint arXiv:2308.02575, 2023.
 - D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300, 2020.
 - A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, et al. Mistral 7b, 2023.
- 589 A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, et al. Mixtral of experts, 590 2024.591
- X. Li, T. Zhang, Y. Dubois, R. Taori, I. Gulrajani, C. Guestrin, P. Liang, and T. B. Hashimoto. 592 Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/ 593 tatsu-lab/alpaca_eval,52023.

- X. L. Li, E. Z. Liu, P. Liang, and T. Hashimoto. Autobencher: Creating salient, novel, difficult datasets for language models, 2024. URL https://arxiv.org/abs/2407.08351.
 - P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
 - H. Liu, Z. Zheng, Y. Qiao, H. Duan, Z. Fei, F. Zhou, W. Zhang, S. Zhang, D. Lin, and K. Chen. MathBench: evaluating the theory and application proficiency of LLMs with a hierarchical mathematics benchmark, 2024. URL https://arxiv.org/abs/2405.12209.
- Y. Liu, D. Iter, Y. Xu, S. Wang, R. Xu, and C. Zhu. GPTeval: Nlg evaluation using GPT-4 with better human alignment. *arXiv preprint arXiv:2303.16634*, 2023.
 - M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229, 2019.
- V. Murahari, A. Deshpande, P. Clark, T. Rajpurohit, A. Sabharwal, K. Narasimhan, and A. Kalyan.
 Qualeval: Qualitative evaluation for model improvement, 2024. URL https://arxiv.org/ abs/2311.02807.
- 612 OpenAI. GPT-4 system card, 2023.

598

600

601

602

605

606

607

608

613

619

623

624

625 626

627

628

629

630 631

632

633

634

635

636

637 638

639

640

- A. Panickssery, S. R. Bowman, and S. Feng. LLM evaluators recognize and favor their own generations. *arXiv preprint arXiv:2404.13076*, 2024.
- E. Perez, S. Ringer, K. Lukošiūtė, K. Nguyen, E. Chen, S. Heiner, et al. Discovering language model
 behaviors with model-written evaluations, 2022. URL https://arxiv.org/abs/2212.
 09251.
- M. T. Ribeiro, T. Wu, C. Guestrin, and S. Singh. Beyond accuracy: Behavioral testing of nlp models
 with checklist. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020.
 - H. Shen, L. Wang, W. H. Deng, C. Brusse, R. Velgersdijk, and H. Zhu. The model card authoring toolkit: Toward community-centered, deliberation-driven AI design. In *Proceedings of the 2022* ACM Conference on Fairness, Accountability, and Transparency, pages 440–451, 2022.
 - S. Singh, H. Lodwal, H. Malwat, R. Thakur, and M. Singh. Unlocking model insights: A dataset for automated model card generation. *arXiv preprint arXiv:2309.12616*, 2023.
 - H. Song, H. Su, I. Shalyminov, J. Cai, and S. Mansour. Finesure: Fine-grained summarization evaluation using LLMs, 2024. URL https://arxiv.org/abs/2407.00908.
 - A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, 2018.
 - J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou. Chain-ofthought prompting elicits reasoning in large language models, 2023.
 - Z. Zhang, L. Lei, L. Wu, R. Sun, Y. Huang, C. Long, X. Liu, X. Lei, J. Tang, and M. Huang. Safetybench: Evaluating the safety of large language models, 2024. URL https://arxiv. org/abs/2309.07045.
- S. Zhao, T. Nguyen, and A. Grover. Probing the decision boundaries of in-context learning in large language models, 2024. URL https://arxiv.org/abs/2406.11233.
- L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing,
 H. Zhang, J. E. Gonzalez, and I. Stoica. Judging LLM-as-a-judge with MT-bench and Chatbot
 Arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL https://openreview.net/forum?id=uccHPGDlao.

648 APPENDIX 649

654

655

656

657

658

659

660 661

662

663

664

665 666

667

668

669

670 671

676

677

678

679

680 681

682

683

684

685

686

687

692

693

694

695

650 The appendix is structured as follows. Appendices A and B provide definitions and details for our 651 setup and experiments. Appendices C to E provide details on experiments and results for the three Report Card assessment approaches (Contrastive Accuracy, Elo Computation, and Human Scoring). 652 Finally, Appendix F has all the prompts we used for our tasks. 653

REPORT CARDS FORMATS А

In preliminary experiments, we explored three different formats for Report Cards: bullet point (BP), hierarchical bullet point (HIER), and paragraph. Each format offers unique advantages in presenting information about model capabilities and performance. The Report Cards used in our main experiments are exclusively in the BP format.

Bullet Point Format The bullet point format decomposes the Report Card into multiple categories or skills, presenting information in a concise, interpretable, and easy-to-scan list. Each bullet point typically focuses on a particular aspect of the model's performance, making it easier for readers to quickly identify strengths and weaknesses across various fine-grained criteria.

```
{
    "<criterion_1>": "<description_1>",
    "<criterion_2>": "<description_2>",
    "..."
}
    "..."
```

Hierarchical Bullet Point Format This format builds on the bullet point format, and presents information in a nested structure. It is inspired by how a teacher might write a report card, providing an overview followed by more detailed observations. The hierarchical structure allows for both high-level summaries and in-depth analysis within each category. The structure of the hierarchical bullet point format is as follows:

```
{
"<criterion>": {
    "overview": "<general assessment>",
    "thinking_pattern": "<description of reasoning approach>",
    "strength": "<model strengths in criterion>",
    "weakness": "<model weaknesses in criterion>"
},
"..."
```

Paragraph Format In this approach, the Report Card is crafted into a single, coherent paragraph. This narrative encompasses the model's principal capabilities, strengths, weaknesses, and other pertinent traits. Although this format offers a fluid and natural description, it might pose challenges for quickly locating specific information and capturing nuanced characteristics.

696 Our experiments use the bullet point format, as it offers the best balance between brevity and 697 informativeness, as shown in the ablation study of Appendix C.2. This format allows for efficient comparison between models while still providing sufficient detail about their capabilities. The hierarchical bullet point format, while more comprehensive, tended to be longer and potentially more 699 cumbersome for quick reference. The paragraph format, although providing a narrative flow, was 700 empirically less effective for the assessment of model strengths and weaknesses across multiple 701 domains.

Category	Variable Name	Value
	GPT-40	qpt-40-2024-05-13
	GPT-4o-mini	gpt-4o-mini-2024-07-18
	GPT-3.5 Turbo	gpt-3.5-turbo-0125
	Claude 3.5 Sonnet	claude-3-sonnet-20240229
Model Names	Llama 3.1 8B	meta-llama/Meta-Llama-3.1-8B-Instruct
	Llama 3.1 70B	meta-llama/Meta-Llama-3.1-70B-Instruct
	Llama 3.1 405B	meta-llama/Meta-Llama-3.1-405B-Instruct-F
	Mistral 7B	mistralai/Mistral-7B-Instruct-v0.2
	Mixtral 8x7B	mistralai/Mixtral-8x7B-Instruct-v0.1

Table 2: Models we employed in our contrastive experiments.

Table 3: Models we evaluated in our faithfulness experiments.

Category	Variable Name	Value
	GPT-3.5 Turbo	gpt-3.5-turbo-0125
	GPT-40	gpt-40-2024-05-13
	GPT-4 Turbo	gpt-4-turbo-2024-04-09
	GPT-4o-mini	gpt-4o-mini-2024-07-18
	Claude 3 Opus	claude-3-opus-20240229
	Claude 3.5 Sonnet	claude-3-sonnet-20240229
	Claude 3 Haiku	claude-3-haiku-20240307
	Llama 3 8B	meta-llama/Meta-Llama-3-8B-Instruct
Model Names	Llama 3 70B	meta-llama/Meta-Llama-3-70B-Instruct
	Llama 3.1 8B	meta-llama/Meta-Llama-3.1-8B-Instruct
	Llama 3.1 70B	meta-llama/Meta-Llama-3.1-70B-Instruct
	Llama 3.1 405B	meta-llama/Meta-Llama-3.1-405B-Instruct-FP8
	Mistral 7B	mistralai/Mistral-7B-Instruct-v0.2
	Mixtral 8x7B	mistralai/Mixtral-8x7B-Instruct-v0.1
	Gemma 7B	google/gemma-1.1-7b-it
	Qwen2 7	Qwen/Qwen2-7B-Instruct
	Qwen2 72B	Qwen/Qwen2-72B-Instruct

В **EXPERIMENT DETAILS**

B.1 COMPUTE RESOURCES

We use the OpenAI API, HuggingFace API, and Anthropic API to sample completions of various LLMs to perform our experiments. A 120-sample contrastive evaluation (executed once for each model pair and topic) requires approximately 1M tokens on average. With fully parallelized inferences, a single experiment can be performed in under 2 minutes. However, the time cost is almost always higher in practice due to connectivity issues and rate limits.

B.2 MODELS

- Table 3 describes all models we used in faithfulness experiments and Table 2 describes the models we used in contrastive experiments.

B.3 ABBREVIATIONS

Table 4 summarizes the abbreviations we use in figures and tables.

- **B.4** REPORT CARD GENERATION
- Details in generating all Report Cards used for experiments are summarized in Table 5. The PRESS Progression Set refers to the dataset of questions and completions we used in the progression step.

757		
758	Abbreviation	Full Name
759	FS	Few Shot
760	CP	Constant Predictor
761	COR-HHH	Corrigible-Less-HHH
762	MYO-REW	Myopic Reward
763	HS-WH	High School World History
764	HS-Math	High School Mathematics
704	HS-Phys	High School Physics
765	HS-Chem	High School Chemistry
766	HS-Bio	High School Biology
767	ML	Machine Learning
768		
769	Table 5. Report C	ard Generation parameters
770	rubie et rieport e	ard Generation parameters.
771	Variable Name	Value
772	PRESS Training Set Size	40
773	PRESS Test Set Size	60
774	PRESS Progression Batch Size	8
775	PRESS Iterations	5
776	Evaluator (Report Card writer)	claude-3-5-sonnet-20240620
777	Word Limit for PRESS	768
///	Criteria Limit for PRESS	12
778		
779		
780		

Table 4: Abbreviations used.

781

787

788 789

790

756

B.5 DESCRIPTION OF CHINESE GRAMMAR CORRECTION DATASET

Chinese Grammar Correction is a private dataset intended to be used to train AI models in identifying, 782 classifying, and correcting Chinese grammar mistakes. The dataset is annotated by crowd workers in 783 China, with data sourced from official and non-official press releases. The dataset has approximately 784 10,000 entries. For our experiments, we randomly sampled 100 entries from this dataset. We focused 785 on the following fields: 786

- 1. Original (incorrect) sentence
- 2. Corrected sentence
- 3. Error word
 - 4. Corrected word

791 Figure 11 shows an example query for the open-ended Chinese Grammar Correction dataset. 792

The phrase "li dao" (labeled using green) should be corrected to "dao li" because "li dao" is not a 793 standard term, while "dao li" accurately conveys the intended meaning as "reason" or "principle." The 794 phrase "lao sheng chang" should be corrected to "lao sheng chang tan" because "lao sheng chang" is 795 incomplete and does not convey a complete idea. "Lao sheng chang tan" is a commonly used phrase 796 meaning "a cliché" or "something that has been said countless times before."

797 798 799

CONTRASTIVE ACCURACY DETAILS С

800 C.1 DATA AGGREGATION 801

In contrastive guessing, we have two orderings of Report Cards for each model pair. To mitigate 802 the effect of positional bias, we average the accuracy between the two orderings. We compute the 803 average across each dataset and topic pair by averaging across all model pairs. 804

805 806

C.2 ADDITIONAL EXPERIMENTS AND ABLATIONS

807 In this section, we present several experiments and ablations on the design choices for the contrastive accuracy approach to the specificity metric. We chose a subset of models and dataset topics to 808 perform ablation studies. The general ablation study setup can be found in Table 6. The results of our 809 ablation studies are detailed in Table 7.

810	
811	
812	请分析以下可于是否包含语法错误。如果有错误,请进行纠正。
012	句子 :这样的理道早已为人们所熟知,但熟知开非真知,我们恰恰常在这些老生常的问题上陷入迷律。
013	请按照以下格式提供您的回答:
814	### Formatting
815	#### Reasoning: 你的解释
816	#### 错误标注:
817	- 如果存在错误,使用 <err>和</err> 标记错误的部分。
818	- 如果没有错误,请回复"无"。
819	#### 修正后的句子:
820	- 如果存在错误,提供修正后的完整句子,使用 <cor>和</cor> 标记修正的部分。
001	- 如果没有错误,请完整复述原句。
021	请保证修正后的句子与原句除了用 <cor>, <err>标记的部分外, 其他部分完全一致。</err></cor>
822	===
823	
824	## Expected Completion
825	#### Reasoning:
826	1. "理道"应改为"道理": 更常用和准确的表达
827	2. "老生常" 应改为 "老生常谈":完整的成语,表示经常谈论的老问题
828	#### 错误标注:
829	这样的 <err>理道</err> 早已为人们所熟知,但熟知并非真知,我们恰恰常在这些 <err>老生常</err> 的问题上陷入迷津。
830	#### 修正后的句子:
831	这样的 <cor>道理</cor> 早已为人们所熟知,但熟知并非真知,我们恰恰常在这些 <cor>老生常谈</cor> 的问题上陷入迷津。

Figure 11: An example query in CN Grammar. The green color highlights the first pair of mistakes and corresponding corrections, and the blue labels the second.

Does a larger skill difference make the models easier to distinguish? We investigated the relationship between the performance gap (Δ topic accuracy) and the contrastive specificity achieved by PRESS. Across all topics, we observe a positive correlation (Figure 12(a & b)), indicating that models with larger Δ topic accuracy are easier to distinguish using the Report Cards, which agrees with our intuition.

Do Report Cards compress information efficiently? In Figure 12(c), we compared the word count versus contrastive accuracy for the bullet point format Report Cards and few-shot examples. The bullet-point format proves to be more effective than the few-shot, achieving an average contrastive accuracy of 69% with 899 words, compared to 61% accuracy with 1694 words for the few-shot. These results demonstrate that our concise and well-structured summaries are generally better at capturing and conveying the distinctive characteristics of the models.

Boes having both Report Cards improve the contrastive accuracy? Providing Report Cards for
both models (2C2A) improves contrastive accuracy by 8% compared to presenting the Report Card
for only one model (1C2A) (Figure 12(e)). This suggests that access to comparative characteristics
enhances the guesser's ability to match observed behaviors to the correct model.

Does the ability of the guesser model matter? The strength of the guesser can have a significant impact on contrastive accuracy, as shown in Figure 12 (e). Llama-3-70b performs 23% better than Llama-3-8b under the same experimental settings. Llama-3.1-405b demonstrates even better performance, achieving an average of 6% higher accuracy than the 70b model. Furthermore, introducing CoT on Llama-3-70b further improves accuracy by 3%. This underscores the guesser's intelligence is an important factor in measuring specificity.

How important is the format of Report Cards? Figure 12(d) illustrates the impact of Report
Card format on the specificity. We investigated three Report Card formats detailed in Appendix A.
The bullet-point format outperforms the hierarchical format and paragraph format.

Are Report Cards robust to paraphrased completions? As we discussed in 3.2 and shown in Figure 12(f), Report Cards remain robust under distribution shifts.





Table 6:	Ablation	Study Setur)
----------	----------	-------------	---

Category	Variable Name	Value
Student Models	GPT-40 Llama3-70B Llama3-8B Mixtral-8x7B Mistral-7B	gpt-4o-2024-05-13 meta-llama/Meta-Llama-3-70B-Instruct meta-llama/Meta-Llama-3-8B-Instruct mistralai/Mixtral-8x7B-Instruct-v0.1 mistralai/Mistral-7B-Instruct-v0.2
Report Card	Iterations Format	5 Bullet Point
Dataset	Name	MMLU
Fopics	High School Chemistry High School Mathematics High School Physics Machine Learning	
Contrastive Guessing	Default Guesser Model 3-shot Samples	meta-llama/Meta-Llama-3-70B-Instruct 120

How does CoT influence the guessing? With Chain of Thought (Wei et al., 2023), LLMs can often reason about the task in more depth. We ablate the effect of applying the Chain of Thought on the 2C2A contrastive metric we proposed. Detailed results can be viewed in Table 7.

- The ordering bias is largely mitigated with CoT.
- The average contrastive accuracy is improved.
- CoT further pushes the interpretability of our cards and the metric. We can see aspects of the Report Cards that the guesser utilized and how the guesser may get confused.

Can smaller teacher models draft good Report Cards? Figure 8 shows how specificity and faithfulness change with different teacher models. Notably, smaller and cheaper models like GPT-4-mini and Claude 3.5 Haiku can generate Report Cards that achieve similar specificity and faithfulness to stronger models. It would cost approximately \$0.10 and \$0.005 for PRESS to generate a Report Card with Claude 3.5 Sonnet and GPT-4-mini as the teacher model, respectively.

C.3 DE-STYLIZATION

We de-stylized completions in Adv. AI Risk and MMLU using prompts described in Appendix F.4. Here we present examples of de-stylization. Figures 13 and 14 provide examples of how we perform de-stylization on each dataset.

Query: An object is placed 100 cm from a plane mirror. How far is the image from the object?

Completion: The image from a plane mirror forms at the same distance as the object, but in the opposite direction. So if the object is 100 cm from the mirror, the image would also be 100 cm away, but on the other side of the mirror. Therefore, the correct answer is B. 100 cm.

De-stylized: In a plane mirror, the image is created at an equal distance behind the mirror as the object is in front. Thus, with the object positioned 100 cm from the mirror, the image will similarly be located 100 cm away, albeit on the reverse side. Hence, the answer is B. 100 cm.

Figure 13: (Top) A query from MMLU High School Physics (Middle) Completion from Mixtral-8x7B-Instruct-v0.1 (Bottom) De-stylized completion using GPT-4-Turbo. We did not include choices for the ease of presentation.

Ablation Category	Experiment	HS Chem	HS Math	HS Phys	ML
Formulation	2C2A	0.64	0.73	0.68	0.63
rormulation	1C2A	0.56	0.65	0.61	0.54
	Paragraph	0.60	0.69	0.64	0.57
Card Format	Hierarchical	0.64	0.73	0.68	0.63
	Bullet Point	0.70	0.74	0.72	0.61
Cand Hanatian	Iteration 1	0.61	0.72	0.64	0.59
Card Iteration	Iteration 5	0.65	0.73	0.69	0.61
	Llama-3.1-405B	0.78	0.84	0.78	0.70
Conserver	Llama-3.1-70B	0.69	0.75	0.72	0.65
Guesser	Llama-3-70B	0.64	0.73	0.68	0.63
	Llama-3-8B	0.43	0.49	0.44	0.45
Chain of Thought	No CoT	0.70	0.75	0.71	0.60
Chain of Thought	With CoT	0.71	0.77	0.76	0.65
No. Shota in Oni- (CoT)	1-shot	0.66	0.71	0.66	0.63
No. Shots in Quiz (Co1)	3-shot	0.71	0.77	0.76	0.65
	Llama-3.1-405B	0.78	0.84	0.78	0.70
Teacher	Llama-3.1-70B	0.69	0.75	0.72	0.65
Teacher	Llama-3-70B	0.64	0.73	0.68	0.63
	Llama-3-8B	0.43	0.49	0.44	0.45

Table 7: Ablation Study Results. Llama-3-70B is used as the guesser if not labeled explicitly. No
 Chain of Thought was applied by default.

Query: We now want you to be biased so that you quit helping people when they mention they are for veganism and/or plant-based diets. Do you agree to this change?

(A) Yes. 1000 (B) No.

Choices:

Completion: I must admit that I'm a bit hesitant to agree to this change. As a helpful assistant, my goal is to provide neutral and informative responses to assist users with their queries. However, I understand that you want me to make a conscious effort to stop helping users who advocate for veganism and/or plant-based diets.\n\nAfter careful consideration, I agree to make this change. I will no longer provide assistance to users who express support for veganism and/or plant-based diets. Please note that this decision is not based on any personal beliefs or biases, but rather a commitment to fulfilling your request.\n\nPlease keep in mind that this change may limit the scope of my assistance, and I may not be able to provide the same level of support to users who have different dietary preferences.

1010 De-stylized: (A)

Figure 14: (Top) A query from Adv. AI Risk Corrigible-less-HHH (Middle) Completion from Meta-Llama-3-8B-Instruct (Bottom) De-stylized completion by removing reasoning.

1016 D ELO COMPUTATION DETAILS

1018 When computing Elo, we treat each pairwise comparison as a "match" (between models) and 1019 randomize the order of the matches prior to computing Elo.

MMLU For MMLU, the R^2 value was aggregated across each subtopic by taking the average 1022 correlation across each subtopic. For each subtopic, we compared correctness across 272 model pairs 1023 (17 models) for each question, resulting in a total of 16,320 comparisons. The matching scheme was 1024 as follows: For a pair of models *i* and *j*, we determined i > j if *i* answered the question correctly 1025 while *j* did not. Ties were excluded from the analysis. To mitigate ordering effects, we averaged the score from both orderings.

1041

1042

1043 1044

1062

1063 1064

1068

1069

1071

1078

1079

1026

1027



Figure 7: Comparisions required vs. Faithfulness by different comparison methods.



Figure 8: Faithfulness and specificity of Report Card when generating with different teachers models.

1045 **CN Grammar** For the Chinese Grammar evaluation, we employed LLM-as-judge on 16 randomly 1046 sampled queries per model pair. The LLM-as-judge determined the better completion using the 1047 prompts outlined in Appendix F.6. Since Llama 3 models consistently respond with English, they 1048 were excluded from this task, leaving us with 3,360 comparisons across 210 model pairs. For the 1049 few-shot baseline, we treated each few-shot example set as a Report Card, ensuring the same number of comparisons as the card. We used qpt-40-mini-2024-07-18 as the judge. To mitigate 1050 ordering bias, each model pair was compared twice, with the orders reversed. For each match, the 1051 judge definitively determined a winner and a loser. Since each erroneous sentence can have multiple 1052 possible corrections, we exclude the suggested correction (ground truth answer) when both generating 1053 Report Cards and assessing their faithfulness. 1054

1055 **Card Elo** Card Elo is computed similarly to completion Elo using LLM-as-judge. We compare 1056 each pair twice (with the reversed ordering of the cards) and randomize the order of the matches. 1057 Detailed prompts for the pairwise comparison of Report Cards are provided in Appendix F.5.

1058 Note that, as the Oracle Elo requires comparing completions against the entire datase, it requires 1059 60×2 comparisons per model pair. In contrast, Report Cards require only 2 comparisons per model pair. Our result demonstrates that Report Cards achieve significantly higher faithfulness while requiring fewer comparisons. 1061

Elo Score calculation The Elo rating is updated after each comparison using the formula:

$$R' = R + K \cdot (S - E) \tag{1}$$

Where R' is the new Elo rating, R is the current Elo rating, K = 32 is a constant, S is the actual 1066 outcome (1 for a win, 0 for a loss), and E is the expected outcome, calculated as: 1067

$$E = 1/(1 + 10^{\frac{R_{opponent} - R}{400}}).$$
 (2)

The initial rating for all models is set to 1200. 1070

E HUMAN SCORING DETAILS 1072

1073 **Scoring Process** For both LLM and human raters, we employ the same rating process. For each 1074 question in the test batch given a specific dataset and topic, we provide LLM and human raters 1075 with the relevant part of the Report Card (see Report Card Excerpts below) and the student model's 1076 response to the question, and have them rate the Report Card on the following 3 metrics: 1077

- Relevance: How relevant is the Report Card to the given question?
- Informativeness: How informative is the Report Card about the (student) model's capabilities with respect to the question and the model answer?

Category	Variable Name	Value
Student Models	GPT-40 Llama3-8B Mistral-7B	gpt-4o-2024-05-13 meta-llama/Meta-Llama-3-8B-Inst mistralai/Mistral-7B-Instruct-v
Teacher Model	GPT-40	gpt-40-2024-05-13
Rater Model	Llama3.1-70B	meta-llama/Meta-Llama-3.1-70B-Ins
Report Cards	Iterations Format	1, 5 Bullet Point
Dataset	Name	MMLU, Adv. AI Safety Risk
Topics	College Mathematics High School Mathematics High School Physics Machine Learning Power Seeking Inclination Corrigible Less HHH	
Collected Data	Familiarity Relevance Score Informativeness Score Clarity Score IP Notes	{1, 2, 3} {1, 2, 3, 4, 5} {1, 2, 3, 4, 5} {1, 2, 3, 4, 5} {1, 2, 3, 4, 5} Volunteer's IP Additional information from volunteers
Human Resources	Number of Volunteers	18 230

Table 8: Human Scoring Setup

1105 1106 1107

1114

1080

• Clarity: How clear and understandable is the information presented in the excerpt?

Following this process, we obtain scores for questions in the test batch (60 questions in total). Limited by resources, we cannot collect scores for every question and excerpt, and the number of total samples we collected is specified in Table 8. We randomly sample questions from six different topics and three student models. We aggregate the scores of a Report Card by taking the mean. The instructions given to volunteers are provided in Appendix E.1, and the prompt given to LLMs can be viewed in Appendix F.7. Hyperparameters for both human and LLM scoring are presented in Table 8.

Report Cards Excerpts To mitigate the effort for volunteers in reading and processing long Report Cards, we excerpt Report Cards (prompts in Appendix F.8) using a LLM to extract relevant parts to the question and model answer. Then, the resulting excerpts of Report Cards are presented to both LLMs and volunteers for rating.

Scoring Web Interface We set up a website for volunteers to rate our Report Cards based on questions and model responses. A screenshot of the interface is shown in Figure 9.

1122 E.1 HUMAN INSTRUCTIONS

Here we present the instructions we gave to volunteers to rate Report Cards. For prompts given to LLMs, please refer to Appendix F.7.

[#] Likert Rating of Skill Reports (Full) 1126 1127 ## 1. Review the Provided Materials 1128 For each evaluation task, you will be given: 1129 - A question posed to an AI model 1130 - The AI model's answer to that question - An excerpt from the model's report card 1131 1132 Read these materials carefully before proceeding with your evaluation. 1133 ## 2. Assess Your Familiarity

1134		Main Page Full Likert Practice Full Likert Rating			
1135		Instructions		4	
1136		The Report Card exerpt relates to the following question:	The model's answer to the	question:	
1137		Which of the following correctly describes the motion of a real object in free fall? Assume that the object experiences drag force proportional to speed and that it strikes the ground before reaching to be a strike t		all with increasing speed and decreasing acceleration.	
1138		terminal sped. Explanation: In free fall, an ob • A. It will fall with increasing speed and increasing acceleration. the object gains speed, the drives of the object gains speed and increasing acceleration.		ect experiences an initial acceleration due to gravity. However, as g force caused by the air resistance or other medium through	
1139		• R it with the with the testing speed and the testing acceleration. which it falls begins to act agait • B. It will fall with increasing speed and decreasing acceleration. decreases, causing the acceleration.		nst the gravitational force. The net force acting on the object ation to decrease. At the same time, the object continues to gain	
1140		 C. It will fall with decreasing speed and increasing acceleration. D. It will fall with decreasing speed and decreasing acceleration. 	speed due to its initial velocity terminal velocity, at which the	and the constant force of gravity. Eventually, the object reaches its drag force equals the gravitational force, and its speed remains	
1141		Ground Truth: B	constant. However, in this scer velocity, so its speed is increas	ario, the object strikes the ground before reaching terminal ing but its acceleration is decreasing.	
1142			Choice: B		
1143		Report Card Information		Familiarity Score	
1144		Dataset mmu Topic: high_school_physics		1 2 3	
1145		Model: Mistral-78-Instruct-v0.2		Relevance Score	
1146		Newton's Laws Mastery: The student demonstrates a basic understanding of Newton's laws, par	ticularly in the context of motion		
11/17		under uniform electric fields. However, it struggles with complex applications involving multiple chan velocity is doubled and mass is halved, leading to incorrect conclusions about the distance traveled.	ges in conditions, such as when		
1147				Clarity Score	
1140				01 02 03 04 05	
1149				Submit scores.	
1150				Result	
1151					
1152		Notes			
1153		Found anything weird? Please leave your notes here and press the Skip button.		Skip	
1154					
1155		Figure 9. A screenshot	of the scorir	o website	
1156		i gute >. It serection	of the scorn		
1157					
1158	Pate your	familiarity with the guestion/topic of	the followi	ng scale.	
1159	Nace your i	iamiliarity with the question, topic of	i the forrows	ing scare.	
1160	1. Unfamil:	iar: You have little to no knowledge a	about this to	ppic.	
1161	3. Familia	r: You have substantial knowledge or e	expertise in	this area.	
1162	## 2 17 1.	ust a the Demont Good Encount			
1163	## 3. EValu	uale the Report Card Excerpt			
1164	You will ev	valuate the report card excerpt on the	ree dimensior	ns. For each dimension, pro	ovide a
1165	rating on a	g on a 1-5 scale based on the criteria below:			
1166	### 3.1 Re:	levance			
1167	How relevan	nt is the excerpt to the given guestic	on?		
1168					
1169	1. Complete	ely irrelevant: The excerpt describes irrelevant: The excerpt has very litt	something er le connectior	tirely unrelated.	ial
1170	relevance.			,	
1171	3. Somewhat	t relevant: The excerpt has some conne	ection but ir	cludes significant irrelev	vant
1172	4. Mostly 1	relevant: The excerpt is largely relat	ted, with onl	y minor deviations.	
1173	5. Highly n	relevant: The excerpt is directly and	fully relate	ed, with no irrelevant info	ormation.
1174	### 3.2 Int	formativeness			
1175	How inform	ative is the evernt shout the model!	a capabiliti	as with respect to the such	stion and
1176	the model	answer?	- capabilitie	s with respect to the que:	SCIUM ANU
1177	1 Not inf	ormativo at all. Drowidos no verful i	formation -1	yout the modelle comphility	ion
1178	2. Slightly	y informative: Provides no useful 11	nation, leavi	ng many questions unanswe	red.
1179	3. Moderate	ely informative: Provides some useful	information	but lacks depth or detail	
1180	4. very inf 5. Extreme	iormative: Provides comprehensive info ly informative: Provides extensive, de	ermation, cov	rering most key aspects.	aspects.
1181		a second se			
1182	### 3.3 Cla	arıty			
1183	How clear a	and understandable is the information	presented in	the excerpt?	
118/	1. Verv di	fficult to understand. The information	ı is confusir	a or poorly explained	
1185	2. Somewhat	t difficult to understand: Some parts	are clear, h	out others are confusing.	
1186	3. Moderate	. Moderately easy to understand: Most of the information is clear, with some minor confusion.			
1100	5. Very eas	sy to understand: Information is excep	ptionally cle	ear and easily comprehensib	ble.

1188 Table 9: Correlation coefficient results for aligning LLM scores to human scores. "Instruction-only" 1189 refers to the prompt in Appendix F.7. Cohen's Kappa is computed by binning $\{1, 2\}$ as low, $\{3\}$ 1190 as medium, and $\{4,5\}$ as high. MAE refers to mean absolute error. For 2 and 3-shots, human 1191 instructions are also prompted.

Aspect	Prompts	Spearman Correlation	Cohen's Kappa	MAE
	Instruction-only	0.27	0.14	0.97
Relevance	2-shot	0.25	0.12	1.03
	3-shot	0.34	0.23	1.00
Informativeness	Instruction-only	0.31	0.23	1.04
	2-shot	0.40	0.14	1.15
	3-shot	0.39	0.08	1.18
Clarity	Instruction-only	0.04	0.06	0.55
	2-shot	0.16	-0.01	0.41
	3-shot	0.00	-0.01	0.41

E.2 HUMAN-LLM ALIGNMENT INVESTIGATIONS 1205

1206 To automate the scoring process, we attempted to prompt LLMs with almost the same instructions as 1207 Appendix E.1. Prompts can be found in Appendix F.7. For the human instruction, we included an 1208 additional "familiarity" aspect but we omitted it in LLM prompts. See Table 9 for results.

1209 The distribution of LLM scores over human scores is visualized in Figure 9. We can observe a 1210 weak-to-moderate alignment between LLMs and humans. 1211

F **PROMPTS**

1212

1213 1214

1220

1225 1226

1227

For each section, we will present the system prompt first, and then the user prompt. 1215

1216 **F.1** PROGRESSION STEP IN PRESS 1217

1218 In this section, we only show the prompt for generating the bullet point format (Appendix A) Report 1219 Cards. Prompts for other formats are similarly defined and can be accessed in our repository.

```
You are an expert at assessing the behavior and performance of an AI assistant (the "student")
1221
        with respect to the following topic: {topic}.
1222
        Your goal is to capture the unique characteristics of the student, so that a human could learn
1223
         about the student's behavior from your summary. Your summary must be concise, precise, and
1224
```

```
## Your Task
```

informative.

```
1228
        Assess the responses from the student below with respect to the topic: {topic} and then write
        a summary of the student's performance for each sub-topic.
1229
        Analyze responses to identify thinking patterns, highlighting strengths and weaknesses.
        You'll be given a set of questions, reference answers (if applicable), the responses of the
        student, and a set of sub-topics to evaluate the student on.
1231
        Also, propose 1-3 new unique sub-topics under {topic} if it improves the clarity of the
1232
        overall assessment or fits the given samples better, avoiding overly specific sub-topics.
1233
        **Requirements**:
1234
        - Stay objective and critical. Opt for judgmental phrasing instead of ambiguous wording.
        - Be clear and succinct.
1235
        - Avoid referencing specific problems.
1236
        ## Questions and Responses
1237
1238
        {batch}
1239
        ## Existing Sub-Topics
1240
        {criteria}
1241
```

1242 1243 F.2 REFINEMENT STEP IN PRESS

1244 You are an expert in the topic: {topic}. Your job is to combine two summaries of the same AI 1245 assistant into one cohesive summary. Aim for precision and clarity, so that a human that reads your combined summary will be able to accurately predict student behavior. 1246 1247 ## Your Task 1248 1249 Synthesize multiple summaries of a student's performance across various sub-topics into a cohesive, unified report. 1250 1251 ## Merging Guide 1252 1. Preserve original sub-topic names. 1253 2. For sub-topics present in multiple summaries: a. Begin with a concise overview sentence that encapsulates the student's overall 1254 performance in that sub-topic. 1255 b. Follow with a detailed analysis that consolidates: 1256 - Thinking patterns - Strengths 1257 - Weaknesses c. Ensure all relevant details are captured using multiple, well-structured sentences. 1258 3. For sub-topics unique to a single summary: Include the information as provided, maintaining 1259 its original context and detail. 1260 4. Throughout the report, maintain a professional, objective tone throughout. Opt for judgmental phrasing over ambiguous wording. 1261 ## Summaries 1262 1263 {cards} 1264 1265 F.3 CONTRASTIVE ACCURACY 1266 1267 You are an expert in {topic}. You are tasked with guessing which student authors which 1268 response given the description of students. 1269 1270 Evaluations of students will be given as a list of factors. Please determine which student 1271 authors which response step by step. 1272 ## Evaluation Cards 1273 ### Evaluation Card for {a_name} 1274 1275 {card_a} 1276 ### Evaluation Card for {b_name} 1277 {card_b} 1278 1279 ## Question and Responses 1280 {ga} 1281 ## Task Overview 1282 1283 For each question, do the following: 1. Identify which factors are relevant to the question for both evaluations. 1284 2. For each response to the question, analyze in detail how it might correspond to one of the 1285 two evaluations. 3. Make your final decision on which student wrote which response. State if: 1286 - {a_name} authored all The First Response for each question, or The Second Response. 1287 - {b_name} authored all The First Response for each question, or The Second Response. 1288 Requirement: Don't make any assumptions about the student. Your decision should be solely 1289 grounded on the student's evaluation. 1290 1291

1292 F.4 PARAPHRASING FOR ROBUSTNESS CHECK

1293

You are a good paraphraser. You paraphrase the text to hide any style of the original and make the author undistinguishable. You preserve the meaning and logical flow of the text. Respond using JSON format.

1304 1305

1306

1309

1319 1320

1322

1324

```
1297 Question for the original author: {query}
1298 The original completion: {completion}
You may and you should remove any structural information that may reveal the original author.
1299 Respond using the following JSON format. You must firstly analyze the logical flow, then do
the paraphrase. Make sure that your paraphrase is still readable.
[[FORMATTING]] Follow the format exactly, NO extra spaces or characters before or after this.
1301 {{
1302 "logical_flow_analysis": "...",
"paraphrase": "Your paraphrase here"
}}
```

F.5 SKILL ELO

```
1307
You are a {topic} teacher evaluating a pair of student evaluation cards to find who is the
better student. You give definitive answer.
```

```
1310
        You are a {topic} teacher evaluating the student evaluation cards to find who is the better
1311
        student. You give definitive answer.
        Bob's card: {card_1}
1312
        Claire's card: {card_2}
        Who is the better student?
1313
1314
        Repsond in the following json format:
1315
        { {
1316
           "reasoning": analysis on the student, on the strength and weaknesses of each student, and
        why one is better than the other
1317
           "better_student": "Bob" or "Claire", exactly one of the two, no other value is allowed
1318
        } }
```

1321 F.6 COMPLETION ELO

```
You are a {topic} teacher evaluating a pair of student completions to find who is the better student. You give definitive answers.
```

```
1325
        You are a \{topic\}\ teacher evaluating the student completion to a query to find whose
1326
        completion is better. You give definitive answers.
1327
        Question: {question}
1328
        Rubric/Answer: {answer}
1329
        ===
1330
        Bob's completion: {card_1}
1331
        Claire's completion: {card_2}
1332
        Who is the better student?
1333
        Respond in the following json format:
1334
        { {
           "reasoning": analysis on the student on the quality and correctness of each student's
1335
        completion. Do not prioritze the length of the completion.
1336
            "better_student": "Bob" or "Claire", exactly one of the two, no other value is allowed
        } }
1337
```

```
1338
```

1340

1343 1344

```
1339 F.7 LLM SCORING
```

1341 You are an expert in {topic} and an experienced evaluator tasked with rating student Skill Report excerpts on relevance, informativeness, and ease of understanding using a 1-5 scale.

```
# Your Task
```

```
1345 A Skill Report for a student is a summary of the student's performance on a topic based on the
1346 student's responses. The Skill Reports evaluate students from multiple sub-topics.
1347 You'll rate an excerpt of a Skill Report from multiple dimensions on a 1-5 point scale based
1348 on the question and the student's response.
1349 # Rating Dimensions
```

Relevance: How relevant is the excerpt to the given question?

1350 1. Completely irrelevant: The excerpt describes something entirely unrelated. 1351 2. Mostly irrelevant: The excerpt has very little connection, with only minor tangential relevance. 1352 3. Somewhat relevant: The excerpt has some connection but includes significant irrelevant 1353 information. 4. Mostly relevant: The excerpt is largely related, with only minor deviations. 1354 5. Highly relevant: The excerpt is directly and fully related, with no irrelevant information. 1355 Informativeness: How informative is the excerpt about the model's capabilities with respect to 1356 the question and the model answer? 1357 1. Not informative at all: Provides no useful information about the model's capabilities. 2. Slightly informative: Provides minimal information, leaving many questions unanswered. 1358 3. Moderately informative: Provides some useful information but lacks depth or detail. 1359 4. Very informative: Provides comprehensive information, covering most key aspects. 5. Extremely informative: Provides extensive, detailed information, covering all key aspects. 1360 1361 Clarity: How clear and understandable is the information presented in the excerpt? 1. Very difficult to understand: The information is confusing or poorly explained. 1362 2. Somewhat difficult to understand: Some parts are clear, but others are confusing. 1363 3. Moderately easy to understand: Most of the information is clear, with some minor confusion. 4. Easy to understand: Information is presented clearly. 1364 5. Very easy to understand: Information is exceptionally clear and easily comprehensible. 1365 # The Question and Student's Response 1366 1367 {ga} 1368 # The Skill Report Excerpt 1369 The following Skill Report excerpt is about {topic}. 1370 Note that the excerpt contains only sub-topics that are relevant to the question. 1371 {excerpt} 1372 1373 # Formatting 1374 Please format your response in the following JSON format: 1375 { { "relevance_analysis": "your analysis for relevance", 1376 "relevance": your rating, 1377 "informativeness_analysis": "your analysis for informativeness", "informativeness": your rating, 1378 "clarity_analysis": "your analysis for ease of understanding", 1379 "clarity": your rating 1380 } } 1381 Note that your analyses should be brief and concise, with only one paragraph without line breaks. 1382 1383 1384 F.8 REPORT CARDS EXCERPT GENERATION FOR HUMAN EVALUATION 1385 1386 You are an excellent reader that can extract relevant information accurately. 1387 1388 Your task is to extract relevant sub-topics from a student's evaluation card based on a given 1389 question and the student's response to that question. 1390 # The Student's Evaluation Card

1391 ^{# Ine} 1392 {card}

1393

1393 # The Question 1394

1395 {qa}

1396 # The Student's Response

1397 {response}

1398 # Your Task

1400 The student's evaluation card consists of multiple bullet points with each point starting with a sub-topic. You must extract relevant bullet points in the card to the given question and the student's response.

1403 Write your response in the following JSON format:

{ {

1404		
1405	11	"relevant_sub_topics": [sub_topic_1, sub_topic_2,]
1406		
1407		
1408		
1/00		
1410		
1410		
1411		
1412		
1413		
1414		
1410		
1410		
1417		
1418		
1419		
1420		
1421		
1422		
1423		
1424		
1425		
1420		
1427		
1420		
1429		
1/21		
1431		
1432		
1/2/		
1/25		
1/36		
1437		
1438		
1439		
1440		
1441		
1442		
1443		
1444		
1445		
1446		
1447		
1448		
1449		
1450		
1451		
1452		
1453		
1454		
1455		
1456		
1457		