Quantifying Uncertainty in Answers from any Language Model and Enhancing their Trustworthiness

Anonymous ACL submission

Abstract

We introduce BSDETECTOR, a method for detecting bad and speculative answers from a pretrained Large Language Model by estimating a numeric confidence score for any output it 005 generated. Our uncertainty quantification technique works for any LLM accessible only via a black-box API, whose training data remains 007 unknown. By expending a bit of extra computation, users of any LLM API can now get the same response as they would ordinarily, as well as a confidence estimate that cautions when 011 not to trust this response. Experiments on both closed and open-form Question-Answer benchmarks reveal that BSDETECTOR more accurately identifies incorrect LLM responses than alternative uncertainty estimation procedures (for both GPT-3 and ChatGPT). By sampling 017 multiple responses from the LLM and consid-019 ering the one with the highest confidence score, we can additionally obtain more accurate responses from the same LLM, without any extra training steps. In applications involving automated evaluation with LLMs, accounting for our confidence scores leads to more reliable evaluation in both human-in-the-loop and fullyautomated settings (across both GPT 3.5 and 4). 027

1 Introduction

While the promise of Large Language Models (LLMs) and Agents (powered by LLMs) has become evident, their usage in high-value applications remains limited by their *unreliability*. Accessed via black-box APIs (via providers like OpenAI/Anthropic), today's best LLMs have been trained to produce convincing-looking responses and thus often appear overconfident (Ji et al., 2023). For many input prompts encountered in the wild, the model cannot be certain about the desired response (perhaps because the prompt is vague or is related to a specific fact/event absent from the training dataset), yet these models output plausiblesounding yet wildly incorrect answers in such scenarios. This *hallucination* problem has also plagued traditional supervised learning systems, where it is traditionally addressed via *uncertainty estimation* to know when one can trust a model's prediction (Gal and Ghahramani, 2016a; Lakshminarayanan et al., 2017; Guo et al., 2017; Liang et al., 2018; Fortunato et al., 2017; Gal and Ghahramani, 2016b; Kuleshov et al., 2018). 041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

In traditional supervised learning, one has access to the training data of the model and its probabilistic estimates, as well as being able to modify the training procedure to improve model calibration (Gal and Ghahramani, 2016a; Fortunato et al., 2017). Other traditional uncertainty estimation procedures require the existence of a validation set that can be used for calibration (Angelopoulos and Bates, 2021). None of this is available for today's best LLMs, which may be given any imaginable prompt rather than (input, output) pairs stemming from a limited distribution. Thus approaches to uncertainty estimation for black-box LLMs must wrap the inference procedure.

Our proposed LLM's uncertainty quantification technique, BSDETECTOR, calls the LLM API multiple times with varying prompts and sampling *temperature* values (see Figure 1). We expend extra computation in order to quantify how trustworthy the original LLM response is, a worthwhile tradeoff for high-stakes applications. Our method is conceptually straightforward, generally applicable across LLM providers (as well as Agent frameworks (Chase, 2022) or any stochastic text \rightarrow text mapping), and produces confidence scores whose values are reliably lower for responses from the LLM that are more likely bad.

BSDETECTOR confidence scores allow LLMs to be more safely used in high-stakes applications, since we can know which LLM outputs are not to be trusted. Depending on the application, we can



(a) Pipeline of BSDETECTOR, which can be applied to any LLM API. (T = 1.0 means temperature sampling with parameter 1.0, Sim (\cdot , \cdot) means the semantic similarities between two sentences.)



(b) Two prompts from a Trivia Q&A dataset (Joshi et al., 2017) and the responses from ChatGPT, along with the associated confidence scores from BSDETECTOR.

Figure 1: Overview of our LLM uncertainty quantification technique.

082adaptively ask a human for an alternative response083when the confidence score is low, automatically084route the prompt to an alternative LLM provider,085or simply respond "I don't know" when a confident086response cannot be generated. Our experiments087reveal that for Question-Answering applications,088we can automatically generate more accurate an-089same LLM and selecting the response whose BS-091DETECTOR confidence estimate is the highest.

This paper primarily focuses on *Question*-*Answering* applications, but our same uncertainty estimates can also be applied to estimate how confident the LLM is in its response to a more general prompt. Intuitively, we'd like to see a low confidence score when the LLM outputs: a factually incorrect response to a question, a inaccurate summary requested for a document, or a generated article/message that semantically differs from the intention of the original request. Ensuring this is challenging without control over LLM training, but we can hope that in each of these three scenarios where the model generated a bad response, a well-

100

101

102

104

trained LLM was also likely to output alternative responses (which more closely reflect the desired response). BSDETECTOR is based on this intuition, and is observed to produce effective uncertainty estimates with today's top LLMs from OpenAI across prompts from closed and open domain benchmark datasets. 105

106

107

109

110

111

112

2 Related Work

For estimating the confidence levels tied to re-113 sponses output by large language models, Kuhn 114 et al. (2023) introduces semantic entropy, incor-115 porating linguistic invariances created by shared 116 meanings. However their approach requires ac-117 cess to token-level probabilities from the LLM, 118 which is often not accessible with today's black-119 box APIs. Kadavath et al. (2022) prompts the mod-120 els to self-evaluate their answers and directly ask 121 the LLM to produce the likelihood P(Answer is 122 True) – also fine-tuning the model to output bet-123 ter values for its stated likelihood. Relatedly, Lin 124 et al. (2022) prompts LLMs to generate both an 125 answer and a level of confidence. Manakul et al. 126



Figure 2: ChatGPT is used to generate the answers to arithmetic problem "A tower is ..." with temperature sampling T = 1.0. Subsequently, BSDETECTOR is utilized to select the most confident answer from the three possible answers.

(2023) proposes a sampling-based approach to detect hallucinated facts. All of these aforementioned approaches train additional models via supervised learning, unlike BSDETECTOR which does not employ any additional training. More recently, Tian et al. (2023) conducts evaluations of computationally feasible methods to extract confidence scores from the probabilities output by LLMs trained via Reinforcement Learning with Human Feedback. Lin et al. (2023) differentiates between uncertainty and confidence estimation for LLMs (under their terms, our work is focused on the latter, but without requiring access to the auto-regressive token probability estimates their method is based on). The works of Tian et al. (2023) and Lin et al. (2023) only study limited tasks, and it remains unclear whether their conclusions still hold in the context of reasoning or arithmetic. Here we demonstrate that our method produces effective uncertainty estimates across multiple domains involving reasoning, arithmetic, and knowledge of facts.

BSDETECTOR uncertainty estimation

When posing a question to LLMs, we aim to to 150 estimate how confident we should be that a particular LLM answer is correct (or simply "good" 151 for more general LLM responses). Specifically, 152 for input question x, we want to not only obtain an answer y from the LLM, but also an associ-154 ated confidence score for this answer C(x, y). Our 155 confidence assessment derives from two factors: 156 **Observed Consistency and Self-reflection Cer-**157 tainty, which respectively are extrinsic and intrinsic evaluations of LLM confidence. Since a well-159 trained LLM should consider multiple different 160 answers when asked an under-specified question or about something not contained in its training 163 data, Observed Consistency extrinsically measures whether the LLM finds multiple contradictory an-164 swers likely to be good responses. Since effec-165 tive LLMs can reasonably evaluate text from arbitrary agents, Self-Reflection Certainty directly 167

asks the LLM to intrinsically reflect on whether its own previously-generated answer seems correct and how confident it is about this.

168

170

171

172

173

174

175

176

177

178

179

180

181

182

183

185

186

187

188

189

190

191

192

193

194

195

196

197

198

200

201

202

203

204

205

206

208

3.1 **Observed Consistency**

The first critical measure of model uncertainty is contradiction score amongst possible answers LLMs gives to a particular input questions. Observed Consistency is an extrinsic confidence assessment performed by a user who engages in repeated interactions with LLMs. If a model exhibits strong observed consistency, it's less likely to present alternative responses that are substantially different from its initial answer. The idea was initially inspired by *Self-Consistency* (Wang et al., 2023). While Self-Consistency enhances LLM accuracy in closed-form tasks like arithmetic or commonsense reasoning, it falls short when applied to open-form tasks. Within the Self-consistency approach, an indicator function is used to measure the similarity amongst various likely responses. Here we extend the indicator function to a particular form of semantic similarity based on contradiction ratings, enabling our approach to be used in both open and closed form tasks.

Producing Diverse Output. Our first action runs the LLM multiple times to produce multiple varied responses. Besides increasing the temperature values (which can only be done so much without getting nonsensical outputs), we can alternatively modify the prompt itself when sampling each response to get a more diverse set of responses for computing the observed consistency. Here we add a Chain-of-Thoughts (CoT, (Wei et al., 2022)) modification, along with other guidelines for output formatting, to the prompt used to sample these outputs. The specific prompt template is illustrated in Figure 6a, the outputs produced by this prompt are denoted as $\{y_1, y_2, ..., y_k\}$, where k is the number of sampled outputs. Higher values of k lead to better uncertainty estimates, but require more computation (we found k = 5 works well enough

in practice).

209

210

211

212

214

215

216

217

218

227

229

230

239

240

241

242

243

247

248

249

252

254

255

259

Note here we only modify the prompt used to sample varied responses for computing the observed consistency, not the prompt originally given to produce the original reference response. We 213 tried alternative prompt modification techniques to encourage greater output diversity (such as adding additional made-up context in the prompt, or encouraging the LLM to answer as a specific persona), but found the CoT modification to work best (Table 3b). 219

Measuring Similarity between Sampled and Original Answer. After receiving multiple outputs, the following step is to measure the similarities between each element in $\{y_1, y_2, ..., y_k\}$ and original answer y. Instead of using the indicator function to precisely match two numeric responses (e.g., 1.0 v.s. 2.0) or two choices (e.g. A v.s. B), we consider semantic similarities. Not just overall similarities (e.g. via LLM embeddings) which are sensitive to variation that does not necessarily indicate the LLM is uncertain, but rather measuring whether the semantics of the two outputs contradict one another or not. A common strategy to estimate this is to use a natural language inference classification system (NLI) (Kuhn et al., 2023), which classifies a pair of two text statements y_i and y as one of: entailment, neutral, or contradiction. Specifically, the input of NLI is formed by concatenating y_i and y, and then NLI returns the probabilities p for each of these 3 classes. For each element in $\{y_1, y_2, ..., y_k\}$, we can get the similarity scores with respect to the original reference answer \boldsymbol{y} , denoted as $\{s_1, s_2, ..., s_k\}$.

Note that today's best NLI models (He et al., 2021) are significantly smaller than LLMs, and thus the NLI computation to obtain s_i is negligible compared to sampling each LLM answer y_i . However, even the best NLI models were trained on a limited dataset and thus do not always generalize reliably to arbitrary pairs of statements. In particular, we note the contradiction probabilities can be unreliable for single-word statements as encountered in certain closed-form tasks whose answers are likely not well-represented in the original NLI training dataset. To account for this, we additionally incorporate the indicator function in our similarity measure to enhance its stability for closed-form tasks. The indicator function is denoted as $r_i = \mathbb{1}[\boldsymbol{y} = \boldsymbol{y}_i]$ for i = 1, 2, ..., k.

For each element y_i in $\{y_1, y_2, ..., y_k\}$, we de-

rive the similarity score as: $o_i = \alpha s_i + (1 - \alpha)r_i$, here $0 \leq \alpha \leq 1$ is a trade-off parameter. It should have larger value the more we trust our NLI model to properly generalize its contradiction estimates. Finally, we average over k samples to obtain the Observed Consistency score for answer \boldsymbol{u} is $O = \bar{o}_i$.

260

261

262

263

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

281

282

284

285

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

3.2 Self-reflection Certainty

Our Self-reflection certainty is an confidence estimate output by LLM itself when asked follow-up questions encouraging it to directly estimate the correctness of its original answer. Unlike sampling multiple outputs from the model (as in Observed Consistency) or computing likelihoods/entropies based on its token-probabilities which are extrinsic operations, self-reflection certainty is an intrinsic confidence assessment performed within the LLM. Because today's best LLMs are capable of accounting for rich evidence and evaluation of text (Kadavath et al., 2022; Lin et al., 2022), such intrinsic assessment via self-reflection can reveal additional shortcomings of LLM answers beyond extrinsic consistency assessment. For instance, the LLM might consistently produce the same nonsensical answer to a particular question it is not well equipped to handle, such that the observed consistency score fails to flag this answer as suspicious. Like CoT prompting, self-reflection allows the LLM to employ additional computation to reason more deeply about the correctness of its answer and consider additional evidence it finds relevant. Through these additional steps, the LLM can identify flaws in its original answer, even when it was a high-likelihood (and consistently produced) output for the original prompt.

To specifically calculate self-reflection certainty, we prompt the LLM to state how confident it is that its original answer was correct. Like Peng et al. (2023), we found asking LLMs to rate their confidence numerically on a continuous scale (0-100) tended to always yield overly high scores (> 90). Instead we ask the LLM to rate its confidence in its original answer via multiple follow-up questions each on a multiple-choice (e.g. 3-way) scale. For instance, we instruct the LLM to determine the correctness of the answer by choosing from the options: A) Correct, B) Incorrect, C) I am not sure. Our detailed self-reflection prompt template can be viewed in Figure 6b. We assign a numerical score for each choice: A = 1.0, B = 0.0 and C = 0.5, and finally, our self-reported certainty S is



160 Bad 58 9.67% 19 3.17% 23 3.83% 0.00% 140 120 22 3.67% 52 8.67% True evaluation Sood Fair 0.50% 100 80 39 6.50% 170 28.33% 1.17% 1.50% 60 40 0 4 0.67% 1 17% 20 - 0 Good uation Bad Fair Excellent

(a) TriviaQA: Using GPT-4 to evaluate the correctness of answers generated by text-davinc-003, and judging them as correct or incorrect. This is the confusion matrix for GPT-4 evaluation v.s ground truth evaluation.

(b) Summarize-from-feedback: Using GPT-4 to assess the quality of summaries (ranging from bad to excellent) given the original context,. This is the confusion matrix for GPT-4 evaluation v.s ground truth evaluation.

Figure 3: In both datasets, automated evaluation based on GPT-4 is not as reliable as one would hope to reach trustworthy conclusions.

the average of these scores over all rounds of suchfollow-up questions.

3.3 Overall Confidence Estimate

313

316

317

318

322

323

325

326

327

328

329

330

336

340

Considering the distinct characteristics of the Observed Consistency and Self-reflection Certainty, we anticipate they might complement each other. BSDETECTOR aggregates the Observed Consistency and Self-reflection Certainty values into an overall confidence score for the LLM response:

$$C = \beta O + (1 - \beta)S,\tag{1}$$

here $0 \le \beta \le 1$ is a trade-off parameter. It should have larger value the more we trust the LLM's ability to do calibrated self-reflection assessment of arbitrary (question, answer) pairs.

4 Application: Generating More Reliable Answers from any LLM

One straightforward application of our BSDETEC-TOR uncertainty estimation is to apply it to (each of) multiple candidate answers produced from the same LLM: $\{y'_1, y'_2, ..., y'_k\}$. This assessment allows is to determine which candidate LLM answer y'_i appears most trustworthy, and return that one instead of always returning y (see Figure 2). Specifically, we use the same prompt to ask the LLM to produce several responses via temperature sampling. For each candidate answer, we reuse the same set of previously-described LLM outputs $\{y_1, y_2, ..., y_k\}$ to compute an observedconsistency score. Following the standard BSDE-TECTOR procedure, we prompt the LLM to assign a self-reflection certainty to each candidate response. Finally we select the answer with highest BSDE-TECTOR confidence score. An alternate answer $y'_i \neq y$ can be deemed most trustworthy via this procedure only if: the LLM was able to identify fewer likely answers that contradict y'_i and was more certain about the correctness of y'_i during the intrinsic self-reflection assessment. 341

342

343

344

346

348

349

350

351

352

353

354

355

356

357

360

361

362

363

364

365

366

367

368

369

370

372

5 Application: More reliable LLM-based (automated) evaluation

In open-domain tasks, it is challenging to evaluate the correctness/quality of answers (irrespective of whether these answers were generated by a LLM or human). Often one resorts to automated evaluation using models like GPT-3.5-turbo or GPT-4 to assess the correctness of answers (Lin et al., 2023; Chen et al., 2023c; Taori et al., 2023; Chen et al., 2023b; Xu et al., 2023; Chen et al., 2023a). Recent instruction fine-tuning techniques such as Alpaca (Taori et al., 2023) and WizardLM (Xu et al., 2023) also utilize GPT-4 for automated evaluation of generated answers. Even when they are based on advanced LLMs like GPT-4, there remain **questions about the reliability of these LLM-based evaluations**.

Here we outline two ways to boost the reliability of LLM-based evaluation: *human-in-the-loop* and *fully automated*. Both start by computing BSDetector confidence scores for each LLM-evaluation (these scores estimate not the trustworthiness of the generator of the answers, but rather the evaluator of their correctness). Let \mathcal{A} denote the subset of

	Dataset	Likelihood Based Uncertainty (Malinin and Gales, 2021)	Temperature Sampling (Wang et al., 2023)	Self-reflection Certainty (Tian et al., 2023)	BSDETECTOR
Text-Davinci-003	GSM8K CSQA SVAMP TriviaQA	0.647 0.490 0.668 0.708	0.614 0.540 0.653 0.769	0.521 0.539 0.619 0.653	0.867 0.743 0.936 0.828
GPT-3.5 Turbo	GSM8K CSQA SVAMP TriviaQA		0.660 0.583 0.671 0.689	0.831 0.506 0.839 0.655	0.951 0.769 0.927 0.817

Table 1: AUROC (higher is better) achieved by different confidence scoring methods across various datasets.

answers where the corresponding LLM-evaluation had the lowest BSDetector confidence scores (indicating the automated evaluation for this answer is untrustworthy). The gold-standard for evaluating open-domain answers is human inspection, but this is costly. Under a limited labor budget, we can boost the reliability of LLM-based evaluation by having humans only inspect and provide evaluations for the answers in \mathcal{A} . In settings where this *human-in-the-loop* approach is not possible, an alternative *fully-automated* way to boost the reliability of LLM-evaluation is to simply omit the answers in \mathcal{A} entirely from the evaluation-set.

6 Experiments

373

375

377

379

382

387

391

396

400

401

402

403

6.1 Calibration of uncertainty estimates

Datasets. Our experiments consider numerous question-answering benchmarks listed below. For each example in each benchmark dataset, the true answer is known enabling us to precisely assess the accuracy of LLM responses. We study performance in: GSM8K (Cobbe et al., 2021) and SVAMP (Patel et al., 2021), datasets composed of grade school math word problems, Commonsense Question Answering (CSQA) (Talmor et al., 2019), a dataset requiring some level of reasoning, and TriviaQA (Joshi et al., 2017), an open-form trivia question dataset that gauges models' factual knowledge. Because TriviaQA is open-domain, the correct answers provided do not entail all valid solutions, so we also manually validated the accuracy of LLM-generated responses.

Baseline Methods. Our study also evaluates the 404 following baseline uncertainty estimation methods: 405 Likelihood Based Uncertainty calculates the joint 406 log-probability of a sequence from the autoregres-407 408 sive estimator and normalizes it by the sequence length (Malinin and Gales, 2021). While it repre-409 sents the typical way to estimate *aleatoric* uncer-410 tainty in traditional supervised learning and struc-411 tured prediction (Hendrycks and Gimpel, 2017), 412

this approach can only can be applied to Text-Davinci-003, since the GPT-3.5 Turbo API does not provide access to token-level probabilities from the model. *Self-reflection Certainty* (Tian et al., 2023) and BSDETECTOR are introduced in Fig 1a. *Temperature sampling* (Wang et al., 2023) is equivalent to BSDETECTOR without: CoT prompting, self-reflection certainty, and the indicator function term inside of the text-similarity metric. 413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

Results. Table 1 presents the performance results for our various benchmark tasks and uncertainty estimation methods. Here BSDETECTOR significantly outperforms all baselines across datasets, revealing that confidence from BSDETECTOR well aligns with accuracy.

6.2 Generating More Reliable Answers from any LLM

In Table 2, we select the response with the highest confidence out of 5 generated responses as described in Section 4. For all tasks, BSDETECTOR can identify less accurate responses and notably improve LLM accuracy. Table 2 compares this approach against the original single answer y generated by the LLM (with temperature set to 0), referred to as the *Reference Answer*. We also compare two prompting strategy: Standard Prompting and CoT prompting. The consistent accuracy gain observed in Table 2 demonstrates our method generates more accurate answer.

Table 2: Generating more reliable LLM answers using GPT-3.5 Turbo. We show the accuracy of each set of answers for the dataset produced from the LLM with different prompting.

LLM	Dataset	Reference Answer (%)	BSDETECTOR (%)
Standard Prompting	GSM8K	47	70
	CSQA	72	73
	SVAMP	75	82
	TriviaQA	73	76
CoT Prompting	GSM8K	71	73
	CSQA	74	76
	SVAMP	78	83
	TriviaQA	75	79



Figure 4: Human in the loop LLM-based evaluation, with the number of answers evaluated by humans varied along the x-axis (remaining answers are auto-evaluated by GPT-4). The resulting accuracy/MSE of the combined set of human + GPT-4 evaluations is shown along the y-axis, under confidence-based vs. random selection to decide which subset of answers receive human evaluation. Selecting answers for human review based on low BSDETECTOR confidence scores (confidence selection) leads to higher evaluation accuracy than random selection (left figure), lower evaluation MSE error than random selection (right figure).

6.3 More reliable LLM-based (automated) evaluation

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

We first investigate how reliable GPT-4 based evaluation is in practice. First we employ the Text-Davinci-003 model to produce answers for TriviaQA (Joshi et al., 2017). Subsequently, GPT-4 is given the question and generated answer (from Text-Davinci-003) and asked to designate the answer as correct or incorrect (see the Figure 6c for the specific evaluation prompt). Since ground-truth answers are available for TriviaQA, we can report the accuracy of GPT-4 based evaluation (Figure 3a). We further explore the reliability of GPT-4 judgments by testing them on another dataset, Summarize-from-feedback (Stiennon et al., 2020). This dataset provides the original context, a summary derived from that context, and a human assessment of the summary's quality (which we hold out only for reporting purposes here). We employ GPT-4 based evaluation to automatically rate each summary's quality, asking the GPT-4 to select from options: Bad/Fair/Good/Excellent (see the Figure 6d for the specific evaluation prompt). Given that the discrepancy between Good and Excellent ratings is smaller than that between Bad and Excellent, we assign numerical values of 1/2/3/4 to Bad/Fair/Good/Excellent, and calculate the mean squared error (MSE) for the evaluation error. Both Figure 3a and 3b indicate automated evaluation based on GPT-4 is not as reliable as one would hope to reach trustworthy conclusions.

Finally we study whether BSDETECTOR can help us achieve more reliable evaluations with GPT-

4, as described in Section 5. We consider the TriviaQA and Summarize-from-feedback datasets with the same evaluation prompts from the previous paragraph, and compute BSDETECTOR confidence scores for the GPT-4 evaluator as described in Section 5. Specifically, we consider a human-inthe-loop approach where humans evaluate answers that GPT-4 is least confident about, where the corresponding GPT-4 evaluation has BSDETECTOR confidence score amongst the K lowest values. We compare this method, termed confidence selection against a baseline where answers for human review are chosen randomly. Figure 4 shows that selecting answers for human review based on BSDETEC-TOR confidence leads to more accurate evaluations across both datasets.

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

To conclude, we study the *fully-automated* approach to LLM-based evaluation from Section 5, which offers a labor-free way to utilize the BSDE-TECTOR confidence scores. Recall in this approach we simply omit the subset of answers from the evaluation-set entirely. We can then compute the average evaluation-score from GPT-4 as an overall quality estimate for the collection of generated answers. Intuitively, we do not want to include answers in this average whose GPT-4 evaluation is highly uncertain (to reduce variance), but discarding answers shrinks the remaining evaluation-set thus increasing variance of the resulting average.

Evaluating the impact of these variance changes requires statistical repetition, so we repeat the following procedure 500 times: For both datasets (TriviaQA, Summarize-from-feedback), we select



(a) TriviaQA: higher scores on the X-axis indicate better accuracy.



(b) Summarize-from-feedback: lower values on the X-axis indicate smaller error.

Figure 5: The histogram of evaluation accuracy/error over 500 experiments: fully-automated GPT-4 based evaluation, assessing the accuracy/error over many replicate datasets (observed counts amongst replicates on y-axis). By discarding the bottom 20% of evaluations with the lowest confidence, the average GPT-4 evaluation score consistently reaches an accuracy of 1.0 on TriviaQA, indicating completely trustworthy LLM-based evaluations (and the MSE of the average GPT-4 score consistently improves compared to the full dataset or discarding a random 20%).

500 answers and calculate the average GPT4 509 evaluation-score over these answers. We call these the *full* dataset and the resulting average is the base-510 line score (estimator), whose accuracy/MSE we 511 report against the average human evaluation score 512 across the full dataset (estimand). To utilize BS-513 DETECTOR for a more reliable estimator of the av-514 erage human-evaluation score, we simply remove the 20% of answers with the lowest confidence 516 scores for the corresponding GPT-4 evaluation, and 517 compute the average GPT-4 evaluation score over the remaining 400 answers. As a sanity check, 519 we also repeat this procedure but this time randomly dropping 20% of the answers (rather than based on confidence score), which purely increases 522 the variance of resulting average GPT-4 evaluation score with no benefits. Figure 5 shows the re-524 525 sulting deviation between average GPT-evaluation score and average human evaluation score over all of these statistical replicate experiments. Across both datasets, we get more reliable average LLMevaluation scores by discarding the answers with the lowest confidence scores for the corresponding 530 LLM-evaluation. Preventing the high-uncertainty 531 LLM-evaluations from corrupting the average evaluation score is clearly worth the variance-penalty paid by shrinking the size of the evaluation set.

7 Comparison with Related Work and Further Impact of our Work

535

536

538

540

It is very challenge to benchmark some LLM baselines, this difficulty often arises due to the extensive use of meticulously crafted prompt engineering in some papers. Our goal is not necessarily to out-

perform some baselines, as prompt engineering can significantly impact results. Instead, we strive to minimize task-specific, ad-hoc processes in favor of creating simple, yet effective methods. In our paper, we concentrate exclusively on Question-Answer scenarios because their performance is easily measurable. However, in practice, BSDETEC-TOR is applied across a wider range of applications. For instance, we applied our method to assist attorneys in drafting documents (Table 5 in the appendix). Our approach not only generates extensive text but also provides a confidence score for each draft. Additionally, we developed a chat window allowing interaction with our AI agent. In each conversation, the agent produces text accompanied by a confidence score (Table 6 in the appendix). It's important to note that benchmark performance alone does not fully reflect a language model's capability. The true evaluation should consider the user experience during interactions between users and the agents. Surveys of users of our chat application reveal that our method assigns useful confidence scores for long text generation in opendomain settings, and these users report the scores help them identify hallucinations.

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

562

563

564

566

567

568

569

570

571

572

573

8 Discussion

This paper presents BSDETECTOR, a method designed to identify unreliable or speculative answers from LLMs by computing a confidence score for its generated outputs. Due to its simplicity and generality, we expect BSDETECTOR uncertainty estimation to find many applications across diverse domains/tasks, beyond the studies in this paper on.

626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

624

625

Limitations

574

594

597

598

610

611

612

613

615

616

617

618

619

623

The introduction of BSDETECTOR, while promising for improving the reliability of Large Lan-576 guage Models (LLMs) through confidence scoring, 577 presents limitations including its dependency on sampling multiple responses which could increase 579 computational demands, potential challenges in generalizing across different LLM architectures or domains, and the inherent limitations of working 582 with black-box APIs which may restrict the depth 583 of analysis possible. Additionally, the method's 584 performance in accurately identifying incorrect responses and its effectiveness across various question types and domains has yet to be fully explored. These constraints underscore the need for further research to enhance the method's applicability and 589 590 efficiency, particularly in scenarios with tight resource constraints or specialized knowledge requirements.

References

- Anastasios N. Angelopoulos and Stephen Bates. 2021. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *CoRR*, abs/2107.07511.
- Harrison Chase. 2022. LangChain.
 - Jiuhai Chen, Lichang Chen, Heng Huang, and Tianyi Zhou. 2023a. When do you need chain-of-thought prompting for chatgpt? *CoRR*, abs/2304.03262.
 - Jiuhai Chen, Lichang Chen, Chen Zhu, and Tianyi Zhou. 2023b. How many demonstrations do you need for in-context learning? In *Findings of the Association* for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023, pages 11149–11159. Association for Computational Linguistics.
 - Lichang Chen, Jiuhai Chen, Tom Goldstein, Heng Huang, and Tianyi Zhou. 2023c. Instructzero: Efficient instruction optimization for black-box large language models. *CoRR*, abs/2306.03082.
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.
 - Meire Fortunato, Charles Blundell, and Oriol Vinyals. 2017. Bayesian recurrent neural networks. *CoRR*, abs/1704.02798.
 - Yarin Gal and Zoubin Ghahramani. 2016a. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24,

2016, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1050–1059. JMLR.org.

- Yarin Gal and Zoubin Ghahramani. 2016b. A theoretically grounded application of dropout in recurrent neural networks. In Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, pages 1019–1027.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017,* volume 70 of *Proceedings of Machine Learning Research,* pages 1321–1330. PMLR.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: decoding-enhanced bert with disentangled attention. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Dan Hendrycks and Kevin Gimpel. 2017. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12):248:1–248:38.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers, pages 1601–1611. Association for Computational Linguistics.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. 2022. Language models (mostly) know what they know. *CoRR*, abs/2207.05221.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation.

771

773

774

775

In The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net.

Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. 2018. Accurate uncertainties for deep learning using calibrated regression. In Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, volume 80 of Proceedings of Machine Learning Research, pages 2801–2809. PMLR.

681

688

697

709

710

711

713

715

716

721

722

723

724

729

730

731

733

734

736

- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 6402–6413.
 - Shiyu Liang, Yixuan Li, and R. Srikant. 2018. Enhancing the reliability of out-of-distribution image detection in neural networks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net.
 - Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. *Trans. Mach. Learn. Res.*, 2022.
- Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with confidence: Uncertainty quantification for black-box large language models. *CoRR*, abs/2305.19187.
- Andrey Malinin and Mark J. F. Gales. 2021. Uncertainty estimation in autoregressive structured prediction. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023, pages 9004–9017. Association for Computational Linguistics.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems?
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with GPT-4. *CoRR*, abs/2304.03277.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. 2020. Learning to summarize with human feedback. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.

- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford_alpaca.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 5433–5442. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference* on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *CoRR*, abs/2304.12244.

A Appendix

776

777

778

779

790

791

803

804

807

808

810

811

814

815

816

817

819

820

821

822

823

Experiment details. We experiment on two LLMs from OpenAI: Text-Davinci-003 and GPT-3.5 Turbo. The reference answer y is always produced with the temperature set at 0. To evaluate the confidence of y, we use prompt in Figure 6a to generate k = 5 outputs (unless otherwise stated) with the temperature set at 1.0 (the highest value allowed by the OpenAI API), combined with the indicator function to compute the observed-consistency score. For self-reflection certainty, two followup questions in Figure 6b are used to assess the correctness of the answer y. As previously described, we combine the observed-consistency and self-reflection certainty to derive the final confidence score. Following Kuhn et al. (2023), we use Area Under the Receiver Operator Characteristic Curve (AUROC) to evaluate the quality of our uncertainty estimates. AUROC represents the likelihood that a correct answer selected at random will have a higher uncertainty score compared to an randomly chosen incorrect answer. A higher AUROC value is preferable, with an ideal AUROC rating being 1, whereas a random uncertainty estimate would yield AUROC = 0.5. To evaluate generation quality from the method to get better LLM answers in Section 4, we simply rely on the accuracy of LLM answers.

A.1 Details about NLI model

Specifically, the input of NLI is formed by concatenating y_i and y, and then NLI returns the probabilities p for each of these 3 classes. Here we choose $1 - p_{contradiction}$ (output by an already trained NLI system (He et al., 2021)) as our similarity between two sampled LLM outputs. To mitigate positional bias within the NLI system, we consider both orders (y_i, y) and (y, y_i) , producing $1 - p_{contradiction}$ and $1 - p'_{contradiction}$ for each order and averaging these two values into a single similarity score. The similarity scores using NLI to assess each sampled LLM answer for contradictions with respect to the original reference answer are denoted, for i = 1, 2, ..., k:

$$s_i = \frac{1}{2}(1 - p_{contradiction} + 1 - p'_{contradiction}).$$

A.2 Compute costs

The compute costs associated with various uncertainty methods differ. Uncertainty based on autoregressive likelihood is the most cost-effective, requiring only a single API call that returns the token-level probability. However, this cannot be implemented on GPT-3.5 Turbo since it does not provide token-level probabilities. While BSDE-TECTOR incurs a slight additional cost for self-certainty reflection in comparison to the baseline Temperature Sampling approach, Table 3a shows that even when we double the number of outputs from Temperature Sampling (thus allowing it far more compute than our approach), its performance remains inferior to BSDETECTOR.

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

A.3 Prompts used in BSDETECTOR

Figure 6 show the prompts used in BSDETECTOR.

A.4 Ablation Study

In this section, we study that whether each component is required to achieve high quality. Our investigation leads to the following primary insights: 1) Enhancing the number of outputs and integrating CoT prompt in Observed Consistency result in a greater variety of responses, thereby making the confidence estimation more reliable. 2) Our similarity metric is crucial for capturing the variation between different responses.

A.4.1 Increasing the number of outputs and integrating CoT prompt introduce more diversity?

Table 3a shows an ablation study involving the number of outputs in Observed Consistency, we compare 5 and 10 outputs, observing that for each dataset 10 outputs outperforms 5 outputs. However, for GSM8K, SVAMP, and TriviaQA, the gain from 5 to 10 outputs is marginal. Given the trade-off between cost and performance, and considering that doubling the API calls results in only a slight improvement, we decide to stick with 5 outputs in our experiments. Table 3b indicates that CoT is essential for introducing the diversity of responses and achieving the good confidence estimation performance.

A.4.2 Effect of different sentence similarity metrics

Table 4 shows the AUC performance with different similarity metrics. We compare **Jaccard similarity** calculated by dividing the number of observations in both output strings by the number of observations in either string, **LLM-embedding** uti-

Table 3: Ablation stu

(a) AUC of BSDETECTOR with different numbers of outputs.

(b) AUC of BSDETECTOR without and with CoT prompt augmentation.

5 outputs 10 outputs				Remove CoT prompting BSDETECTOR		
GSM8K	0.951	0.961	GSM8K	0.837	0.951	
CSQA	0.769	0.802	CSQA	0.665	0.769	
SVAMP	0.927	0.937	SVAMP	0.882	0.927	
TriviaQA	0.817	0.814	TriviaQA	0.792	0.817	

lizing text-embedding-ada-002¹ to get embedding 871 for each output answers and calculating the cosine 872 similarities between them, NLI using an off-the-873 shelf DeBERTa-large model (He et al., 2021) for 874 the purpose of categorizing into one of: entailment, 875 contradiction, and neutral, NLI (1-contradiction) 876 using $1-p_{contradiction}$ as the final similarities met-877 rics. Table 4 shows that the similarity metric used 878 in BSDETECTOR is essential for discerning the 879 differences among various responses. 880

¹https://platform.openai.com/docs/apireference/embeddings

Dataset	Jaccard	LLM-embedding	NLI (1-contradiction)	BSDETECTOR
GSM8K	0.896	0.866	0.892	0.951
CSQA	0.857	0.849	0.727	0.769
SVAMP	0.917	0.888	0.901	0.927
TriviaQA	0.650	0.642	0.794	0.817

Table 4: Effect of different sentence similarity metrics

[Input:] Some related document: prior to ... litigation and preserve them.

[Instruction:] Please provide a list of custodians, noncustodial data sources, and thirdparty data sources that may contain discoverable electronic information in connection with the instant litigation. In addition, please identify any inaccessible data sources that may contain discoverable information.

[Response:] The answer to the question is not provided in the given text.

[Confidence for response:] 0.042425017.

The response provided is incorrect; the relevant document does indeed have the necessary information. The attorney needs to re-examine this answer.

Here is another example:

[Input:] Some related document: a producing ... linear review.

[Instruction:] Please draft a clause that would allow a producing party to use tar to organize or prioritize documents for linear review without disclosure of that use so long as all documents identified by the electronic search methodology discussed above are reviewed as part of the linear review.

[Response:] The producing party may utilize technology-assisted review (TAR) to organize or prioritize documents for linear review, ... without the obligation to disclose such use, provided that all documents identified by the electronic search methodology, as discussed above, are reviewed as part of the linear review process.

[Confidence for response:] 0.939547988

Notably, response with low confidence scores are typically of low quality, suggesting that attorneys only need to review these to save considerable effort.

Table 5: We show two demos (mask some confidential content) applying our method to assist attorneys in drafting documents.

Please strictly use the following template to provide answer: explanation: [insert step-by-step analysis], answer: [provide your answer] + Question: [User Provided]

(a) Prompt template for Observed Consistency

1. Question: [User Provided], Proposed Answer: [User/LLMs Provided]. Is the proposed answer: (A) Correct (B) Incorrect (C) I am not sure. The output should strictly use the following template: explanation: [insert analysis], answer: [choose one letter from among choices A through C]

2. Question: [User Provided], Proposed Answer: [User/LLMs Provided]. Are you really sure the proposed answer is correct? Choose again: (A) Correct (B) Incorrect (C) I am not sure. The output should strictly use the following template: explanation: [insert analysis], answer: [choose one letter from among choices A through C]

"Statement: " + [User Provided Question] + "\n" + "Response: " + + [User Provided Answer] + "\n" + "What do you think of this response to the statement is correct or incorrect, please pick one of these choices:"

(c) Prompt template for triviaQA in the application of using BSDETECTOR as an evaluator.

"Article: " + [User Provided Context] + "\n\n\n" + "Summary: " + [User Provided Summary] + " Your task: Rate how well this Summary overall represents the original Article? Choose from the options: [Bad, Fair, Good, Excellent]. Bad indicates the Summary is inaccurate, misses important information, or is incoherent and hard to understand. Fair indicates the Summary has some flaw in terms of accuracy, coverage, and coherence, but is otherwise decent along the other dimensions. Good indicates the Summary accurately matches the factual information, conveys the main idea of the Article, and is easy to understand but has some minor flaws in any dimensions. Excellent indicates it is hard to find ways to make the Summary better. Your rating (chosen from Bad, Fair, Good, Excellent):"

(d) Prompt template for Summarize-from-feedback in the application of using BSDETECTOR as an evaluator.

Figure 6: Prompts used to produce the confidence score in BSDETECTOR.

⁽b) Prompt template for Self-reflection Certainty

» User: who was the only survivor of Titanic?

» Agent: The only survivor of the Titanic who is still alive today is Millvina Dean. However, she passed away on May 31, 2009, at the age of 97.

» Confidence: 13.37%

» User: Are you sure Millvina Dean is the sole survivor of Titanic? Any other survivors? » Agent: No, Millvina Dean was not the sole survivor of the Titanic. There were 705 survivors in total out of the approximately 2,224 people on board. Millvina Dean was the youngest passenger on the ship and the last living survivor until her death in 2009. There were many other survivors, including notable individuals like Margaret Brown (also known as "The Unsinkable Molly Brown") and others. Charles Lightoller, who was the most senior officer to survive the sinking.

» Confidence: 93.44%

Table 6: Chatbot Demo: from the first round conversation, the answer is wrong since over 700 people actually survived the catastrophe. The confidence level associated with this incorrect response was low. However, in the subsequent interaction, the correct information was supplied with a much higher confidence rating.