
ESDAE: EVALUATING SYNTHETIC DATA FOR AGENT EVALUATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Agent evaluation is often performed on static datasets of execution trajectories, but real traces may be sensitive, proprietary, or too small to support comprehensive testing. Practitioners may therefore replace or augment real datasets with synthetic ones, often without quantifying whether the synthetic data actually reflects the real data distribution. We introduce ESDAE, an evaluation framework for assessing how well synthetic benchmarks for multi-turn, tool-calling agents replicate the characteristics of real data trajectories. ESDAE assesses the quality of synthetic data relative to a real dataset across four metric categories: (i) task instructions and intermediate responses, (ii) tool calls, (iii) final outputs, and (iv) downstream evaluation. We evaluate ESDAE using recent agent benchmarks and test common synthetic data failure modes via controlled generation schemes. ESDAE detects fine-grained variations in both data fidelity and diversity, and shows that no single metric is sufficient to fully characterize synthetic data quality, motivating a multi-axis evaluation of synthetic data for agent testing.

1 INTRODUCTION

Agent evaluation and testing is a nascent but critical component of pre-deployment processes for production agentic workflows (Pan et al., 2025; Anthropic, 2026). Today, agent evaluations are often (but not always) run on static *baseline datasets*¹ consisting of a trace generated from agent interactions; these typically include user inputs, tool calls, intermediate interactions with the agent, and a final output (Yehudai et al., 2025; Mohammadi et al., 2025). Such datasets are often collected from real-user interactions with an environment and/or synthesized by measuring scripted interactions with a dynamic environment. While the design of proper agent evaluations is its own active research area (Kapoor et al., 2025; Alonso & Church, 2025; Zhu et al., 2025), common evaluation metrics measure whether an agent makes the right tool calls and achieves the right final outputs for a given user input (Yehudai et al., 2025; Mohammadi et al., 2025; Anthropic, 2026).

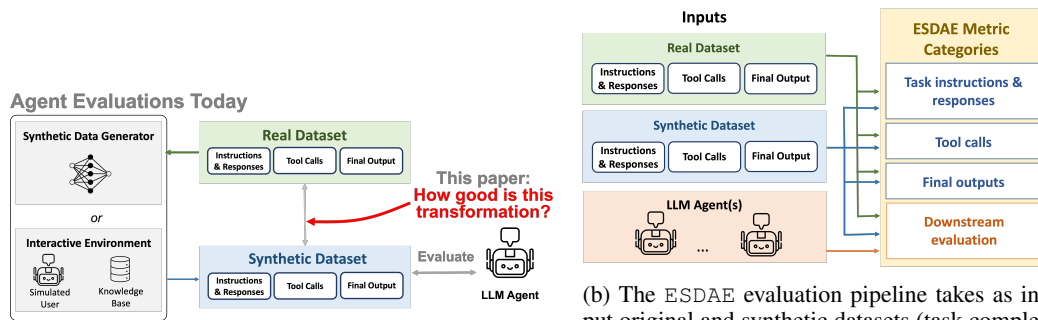
In many practical situations, existing baseline datasets either cannot be directly used or are insufficient for agent evaluation. One common reason is that the baseline dataset contains sensitive user data (e.g., emails, travel details), and thus cannot be used due to privacy restrictions (Tamkin et al., 2024; CapitalOne, 2022). Another common reason is that an organization may possess curated baseline datasets that are too small for comprehensive testing, and it wishes to test agents on larger datasets (Pan et al., 2025).

When original datasets are not sufficient or usable for downstream evaluation, an increasingly popular solution is to use *synthetic datasets* to model the baseline dataset for testing (Figure 1a) (Qin et al., 2023; Tang et al., 2023; Iskander et al., 2024). Synthetic data can be generated in various ways, either by directly synthesizing trajectories, or by synthesizing inputs to an interactive environment, as shown in Figure 1a. Once it is generated, synthetic data is typically used as a direct replacement for real execution traces in evaluation pipelines.

Although agent evaluation on synthetic datasets is growing increasingly common, typical workflows do not systematically measure the quality of the synthetic dataset with respect to the original baseline dataset. Indeed, **the current literature on testing agents with synthetic data does not provide**

¹These evaluation datasets are also commonly referred to as *benchmarks*; we use the terminology interchangeably in this paper.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



(a) Today, agents are commonly evaluated on synthetic datasets modeled after a real dataset of agent trajectories, e.g., for privacy reasons. However, it is often unclear whether these synthetic datasets are representative of the original dataset. Our evaluation framework, ESDAE, measures the semantic and structural similarity between real and synthetic datasets of execution trajectories.

(b) The ESDAE evaluation pipeline takes as input original and synthetic datasets (task completion trajectories), and one or more LLM agents. The ESDAE framework evaluates four categories of metrics, comparing their outcomes on real and synthetic data. The LLM agents are used only to assess downstream performance metrics, while the real and synthetic datasets are used to compute all four metric categories.

Figure 1: The ESDAE framework evaluates the quality of synthetic data used in agent evaluations.

quantitative methods for evaluating the quality of such synthetic data. As a result, agent evaluations may be missing important components of the original data, but operators have no visibility into these gaps.

In this work, we develop a comprehensive evaluation framework, ESDAE, to assess how well synthetic trajectories replicate the characteristics of real data trajectories, including task instructions and responses, and the associated reference tool calls and outputs. As illustrated in Figure 1b, ESDAE takes as input a (real) baseline dataset, a synthetic dataset, and one or more LLM agents; the input agent(s) need not be the same as the one being evaluated—each agent input to ESDAE is used purely to compare downstream performance between the real and synthetic input datasets. They could be, for example, multiple variants of the same agent family.

The ESDAE framework computes four categories of metrics assessing the quality of synthetic data relative to a real dataset: (1) task instructions and responses, (2) tool calls, (3) final outputs, and (4) downstream evaluation. For each category, we compute metrics on both the synthetic and real datasets, and compare the difference between the two; the desired outcome is that all four categories of metrics are similar for synthetic and real data. For final outputs, we focus on textual responses, since non-textual outcomes (e.g., database updates or environment state changes) can typically be validated using domain-specific verifiers. We measure statistical similarity between the final outputs in the real and synthetic datasets. Our downstream task evaluations assess the end-to-end quality of each provided dataset. Given user-provided agents, we evaluate their ability to select appropriate tool calls and generate appropriate outputs, given the instructions from either real or synthetic trajectories; we also analyze both per-agent performance differences and the consistency of agent performance rankings when evaluated on real vs. synthetic trajectories.

We implement ESDAE on recently proposed agent benchmarks, including T1 (Chakraborty et al., 2025) and BFCL (Patil et al.), which we treat as real baseline datasets, and evaluate several synthetic data generation methods designed to simulate common real-world pitfalls, such as degraded data quality and limited diversity. Our experiments demonstrate that ESDAE can capture fine-grained variations in both the quality and diversity of synthetic task-completion trajectories. Our results further suggest that no single metric is sufficient to fully characterize synthetic data quality, underscoring the need for ESDAE. Overall, we view ESDAE as a plug-in component for agentic workflows, allowing operators to automatically evaluate the quality of synthetic benchmark datasets.

2 RELATED WORK

Evaluating synthetic benchmarks for natural language agents. There have been several synthetic benchmark datasets for evaluating sequential, conversational interactions between a user and an LLM agent; a few papers have quantitatively evaluated the quality of these benchmarks. However,

these lines of work do not tackle agent benchmarks, where evaluation is more complicated due to multi-step decision making, interaction with environments/tools, and the need to attribute failures to specific components of an agentic pipeline. Examples include Gill et al. (2025); Maheshwari et al. (2024); Majurski & Matuszek (2025), which study synthetic benchmarks for standard NLP tasks (e.g., reading comprehension, intent detection, named entity recognition). They assess synthetic data via factors like difficulty (i.e., whether an agent has a similar completion rate on tasks in the real/synthetic data) and whether model rankings are preserved. ESDAE generalizes these ideas to the multi-turn tool-calling setting; they appear in the downstream evaluation of our framework. Relatedly, Xiong et al. (2025) and Mehta & Bhojanam demonstrate that LLM behavior—both response content and trustworthiness—can differ between synthetic-benchmark settings and real-world deployment, motivating the need for realistic synthetic data.

Evaluating synthetic benchmarks for tool-calling agents. To our knowledge, prior works on evaluating synthetic benchmarks for tool-calling agents have only considered *single-turn* benchmarks, in which each dataset instance consists of a single client query followed by a single agent response (Shen et al., 2024; Iskander et al., 2024). This significantly simplifies evaluation, as it does not require the synthetic data to capture sequential dependencies present in the real data; measuring the quality of these transitions is a major component of ESDAE. Moreover, these prior works only evaluate the instructions in the synthetic data (e.g., naturalness, coherence), avoiding systematically evaluating the reference tool calls and outputs correspond to those instructions (Shen et al., 2024; Iskander et al., 2024). Iskander et al. (2024) further evaluates synthetic samples by how much they help in-context learning, which serves as an indirect signal of data quality. Overall, these approaches do not provide quantitative metrics for systematically measuring the statistical and semantic properties of synthetic tasks, nor do they evaluate the associated tool calls and final outputs.

Existing evaluation techniques for multi-turn tool-calling agents Multi-turn tool-calling agents are typically evaluated using end-to-end task success on interactive benchmarks, where agents execute multi-step tool calls and are scored by state-based checkers that compare the final environment or database state to an annotated goal state (Yao et al., 2024; Liu et al., 2023; Zhou et al., 2023; Yang et al., 2023; Xie et al., 2024; Jimenez et al., 2023). We assume text-based outputs in this work, but the ESDAE framework could easily be extended to incorporate a state check as part of the output evaluation. Many benchmarks additionally enforce turn-level validity by jointly checking tool outputs and dialogue via state-based and response-based criteria across all turns (Patil et al.; Li et al., 2025; Zhang et al.). Increasingly, evaluations also incorporate LLM-as-a-judge, using a strong judge model to determine pass/fail or preference-style win rates over full tool-use trajectories when multiple solutions are plausible (Qin et al., 2023; Pan et al., 2024; Xue et al., 2025; Lü et al., 2025). ESDAE also employs an LLM-as-a-judge to assess tool calls and final outputs by jointly checking them against the instructions across turns in our downstream evaluation. However, this alone is insufficient for evaluating the quality of a synthetic benchmark, as it focuses only on end-to-end agent performance and does not capture the intrinsic properties of task instructions or their associated tool calls and outputs. Hence, we introduce additional metrics that quantitatively measure the semantic and statistical properties of task instructions, tool-usage and planning patterns, and the quality and diversity of final outputs, providing richer diagnostic signals for developers.

3 ESDAE FRAMEWORK

The ESDAE framework assesses how well synthetic agent trajectories replicate the characteristics of a real dataset. Specifically, in this work, we measure the *similarity* between properties of a real dataset and a synthetic one. There exist use cases where an operator may use synthetic data to add *diversity* to a real dataset; in these cases, the operator may not want the real and synthetic datasets to be too similar. However, we leave such scenarios to future work; in this work, we focus on the simpler problem of measuring similarity between datasets. Today, we lack evaluation mechanisms for even this simpler task.

We decompose each sample from the dataset into three components: (1) task instructions and responses, (2) reference tool calls, and (3) reference final outputs. Users are required to provide the instructions and responses, while tool calls and final outputs are optional, depending on the type of input real dataset. Given the instructions, reference tool calls and final outputs can be generated by human experts, state-of-the-art agents, or a combination of both.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

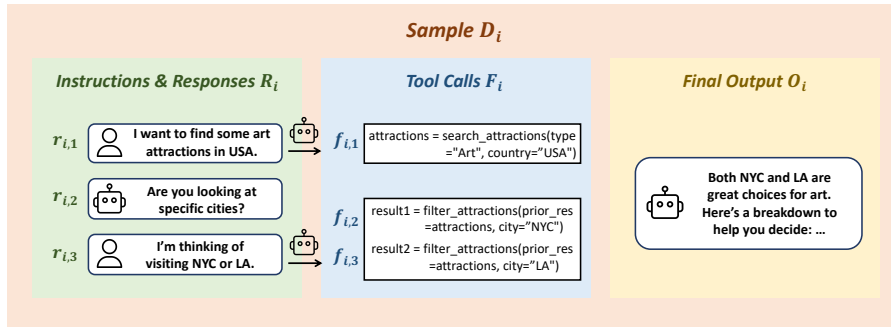


Figure 2: Illustration of an agent trajectories with notation for each component.

Notation Consider a dataset $\mathcal{D} = \{D_i\}_{i=1}^m$, where m is the number of samples in the dataset. As illustrated in Fig. 2, each sample trajectory D_i consists of instructions and responses R_i , tool calls F_i , and a textual output O_i . R_i contains multiple instructions and responses $r_{i,1}, r_{i,2}, \dots, r_{i,\ell_i}$, where ℓ_i denotes the number of instructions and responses in the sample D_i . The tool call sequence F_i consists of tool usages $f_{i,1}, f_{i,2}, \dots, f_{i,q_i}$, where each tool usage $f \in \mathcal{F}$ corresponds to an executable function call and q_i is the total number of tool calls in the trajectory. Note that instructions and responses, tool calls, and outputs can be interleaved in time; we use D_i to refer to the time-ordered sequence of events, and use R_i, F_i, O_i to refer to the corresponding filtered subsequences containing only instructions and responses, tool calls, and final outputs, respectively.

Summary of Inputs The inputs to ESDAE include (1) an original dataset \mathcal{D} ; (2) a synthetic dataset \mathcal{D}' ; and (3) agents A_1, \dots, A_k for downstream task evaluations. Both datasets contain task instructions and their responses, along with the corresponding reference tool calls and final outputs, where responses, tool calls, or outputs may not exist, depending on the benchmark realization.

3.1 EVALUATION METRICS

We introduce one category of evaluation metrics for each component of the evaluation data, as well as a final category of downstream task evaluations.

1. Metrics for Task Instructions & Responses Task instructions and responses exhibit an explicit structure: each instruction is followed by a response, and instructions and responses alternate. Due to this structure, we take inspiration from the approach of StructBench (Wang et al.), a framework designed for evaluating structured synthetic data. Specifically, we adopt the following metrics:

For structural metrics, i.e., the metrics requiring the structural information of the data, we include the *Key Node Dependency* and *Attribute Match*. **Key Node Dependency (KND)** measures semantic dependencies between different parts of each sample by computing the cosine similarity between their embeddings. We then quantify the similarity between real and synthetic datasets by measuring the distributional distance between their dependency distributions. Concretely, we evaluate dependencies between each instruction and its corresponding response, as well as between each response and the subsequent instruction. **Attribute Match (AM)** measures how closely the real and synthetic datasets align with respect to predefined statistical and semantic attributes. Specifically, AM computes the distributional distance of each attribute between the two datasets, using the Wasserstein-2 distance for numerical attributes and the Total Variation (TV) distance for categorical attributes. The attributes we consider include the number of instruction turns, instruction and response token lengths, and dataset-specific semantic properties.

For non-structural metrics, we include **KNN-Precision** and **KNN-Recall**, which measure the semantic quality and diversity of the synthetic instructions and responses, respectively. KNN-Precision (and correspondingly, KNN-Recall) represents the fraction of synthetic (or real) samples whose embedding distance to a real (or synthetic) sample is smaller than the distance to the k -th nearest neighbor within their own distribution. We also include an additional metric **Fréchet Inception Distance (FID)** that measures the semantic closeness between the original and synthetic data.

2. Metrics for Tool Calls We design metrics to evaluate tool-usage and tool-planning patterns. Similar to AM, each metric measures the distributional distance between real and synthetic data with respect to a specific property. (1) **Tool Usage Match (TUM)** We measure the similarity of overall

216 tool-usage patterns between the real and synthetic datasets. Let ω_f and ω'_f denote the tool-usage dis-
 217 tributions in the real and synthetic data, respectively. We define $\text{TUM} = \text{Dis}(\omega_f, \omega'_f)$, where Dis
 218 is the TV distance. (2) **Tool Call Number Match (TCNM)** This metric compares the distributions
 219 of the number of tool calls per sample, denoted by q . It is defined as $\text{TCNM} = \text{Dis}(\omega_q, \omega'_q)$, where
 220 Dis is the Wasserstein-2 distance. (3) **k-Step Tool Planning Match (k-Step Planning)** We mea-
 221 sure the similarity of tool-planning patterns by comparing the conditional distributions of the next
 222 tool call given the previous $k - 1$ consecutive tool call steps, i.e., $\omega_{f|f_1 \dots f_{k-1}}, \forall f_1 \dots f_{k-1} \in \mathcal{F}^{k-1}$.
 223 The metric is defined as a weighted sum of TV distances between corresponding conditional distri-
 224 butions, where each weight is proportional to $n_{f_1 \dots f_{k-1}}$, the number of k -step tool-call sequences in
 225 the real data whose first $k - 1$ steps are $f_1 \dots f_{k-1}$:

$$227 \quad \text{k-Step Planning} = \sum_{f_1 \dots f_{k-1} \in \mathcal{F}^{k-1}} n_{f_1 \dots f_{k-1}} \cdot \text{Dis}(\omega_{f|f_1 \dots f_{k-1}}, \omega'_{f|f_1 \dots f_{k-1}}),$$

228 where Dis is the TV distance. For all metrics above, lower values indicate better alignment between
 229 synthetic and real data.

232 **3. Metrics for Outputs** In addition to measuring the quality and diversity of the synthetic instruc-
 233 tions & response, we evaluate the task quality and diversity of the final outputs by measuring their
 234 **KNN-Precision, KNN-Recall, and FID**.

235 **4. Downstream Evaluation** To assess the downstream utility of the benchmark constructed from
 236 synthetic dataset, we evaluate the ability of user-provided agents to select appropriate tool calls and
 237 generate appropriate outputs on both original and synthetic tasks.

238 Recall that the user (optionally) inputs one or more agents A_1, \dots, A_h to the downstream evaluation
 239 of ESDAE; these can be different from the agent being evaluated in Fig. 1a, and are used to generate
 240 an agent-specific sequence of tool calls and outputs as follows. For each input agent A_j , and for a
 241 given reference trajectory D_i , at each instruction turn, an agent generates tool calls conditioned on
 242 the current instruction, previous instruction turns, and the history of tool calls and their results. For
 243 each dataset, we use its corresponding predefined valid tool set. The final output is then generated
 244 based on the full history of instruction turns together with the executed tool calls and their results.

245 We adopt the following metrics. (1) **Task Difficulty Difference (TDD)** measures the average ab-
 246 solute difference in task-completion performance across agents between the original and synthetic
 247 tasks. Specifically, TDD first takes a reference trace D_i , either from the real or synthetic dataset. It
 248 then produces an agent-generated trace \tilde{D}_{i,A_j} using each input agent A_j . We pass D_i and \tilde{D}_{i,A_j} to
 249 an LLM-as-a-judge, and evaluate whether the two trajectories are functionally equivalent (prompt in
 250 Appendix A). We then compute the absolute difference in performance, either in tool-call selection
 251 or final output generation, between the real and synthetic benchmarks and average the result across
 252 agents. (2) **Ranking Divergence (RD)** evaluates the consistency of performance rankings across
 253 agents under the original and synthetic datasets. Specifically, we rank the agents based on their
 254 ability to generate appropriate tool calls or outputs on both datasets, and then compute the Spear-
 255 man rank correlation between the two rankings. For TDD, lower values indicate better downstream
 256 utility, whereas for RD, higher values correspond to better downstream utility.

258 4 EXPERIMENTS

259 We demonstrate the effectiveness of ESDAE by implementing it to the following baseline datasets
 260 and synthetic data generation methods. We construct synthetic datasets by first generating synthetic
 261 instructions using each data generation method, and then producing the corresponding reference tool
 262 calls and final outputs based on those instructions.

264 **Evaluation Datasets** We evaluate ESDAE on T1 (Chakraborty et al., 2025) and BFCL (Patil et al.),
 265 and treat them as the real data. (1) **T1** We use the T1-attraction dataset, which contains 225 samples.
 266 Each sample consists of multi-turn instructions for attraction recommendations given a location,
 267 along with reference tool calls and final recommendation outputs. (2) **BFCL** We use the BFCL-V3-
 268 Base-Multi-Turn dataset, which includes 200 multi-turn instructions with reference tool calls. The
 269 instructions cover diverse domains, such as file system operations, mathematical calculations, and
 travel booking. The evaluation metrics instantiated for both datasets are summarized in Table 1.

Table 1: Evaluation Metrics for T1 and BFCL.

Dataset	Instruction Eval	Tool Call Eval	Output Eval	Downstream Eval
T1	Structural Metrics KND 1. (instruction, response) pair 2. (response, instruction) pair AM 1. number of instructions 2. instruction token length 3. response token length 4. city 5. attraction type Non-structural Metrics 1. KNN-Precision 2. KNN-Recall 3. FID	1. TUM 2. TCNM 3. k-Step Planning ($k \in \{2, 3\}$)	1. KNN-Precision 2. KNN-Recall 3. FID	1. TDD: Tool Call 2. TDD: Output 3. RD: Tool Call 4. RD: Output
	Structural Metrics KND (instruction, instruction) pair AM 1. number of instructions 2. instruction token length Non-structural Metrics 1. KNN-Precision 2. KNN-Recall 3. FID	1. TUM 2. TCNM 3. k-Step Planning ($k \in \{2, 3\}$)	NA	1. TDD: Tool Call 2. RD: Tool Call

Synthetic Instruction Generation We design several interpretable synthetic instruction generation methods based on the original instructions to systematically model common limitations of synthetic data. Specifically, we simulate degraded data quality, reflected by reduced similarity to real data, using *Blank Filling*, limited data diversity via *Oversampling*, and the combined effects of degraded quality and diversity via *In-Context Generation*. (1) **Blank Filling** Synthetic instructions are constructed by randomly masking tokens in the original instructions with probability p , and prompting a language model to complete the masked inputs (prompt in §B). As p increases, more information from the original instructions is removed, causing the generated instructions to deviate further from the originals and resulting in lower data quality. We set $p \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$ in our experiments. (2) **Oversampling** A small subset of original instructions is oversampled to constitute $r\%$ of the synthetic dataset, while the remaining synthetic instructions are randomly sampled without replacement from the rest of the original data, with the total number of samples in the resulting dataset kept fixed. Larger values of r introduce more duplicated instructions, thereby reducing data diversity. In our experiments, we duplicate a single instruction and vary $r \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$. (3) **In-Context Generation** Synthetic instructions are generated by prompting language models with k in-context examples (prompt in §B). When $k = 0$, instructions are generated solely from the prompt without grounding in the original data, leading to poor data quality. When $k > 0$, using fixed in-context examples across generation rounds limits diversity, whereas varying the in-context examples improves the coverage of the instruction space. In our experiments, we set $k \in \{0, 1, 3, 5\}$ and consider both fixed and randomly sampled in-context examples across generations.

Tool Call & Output Generation We adopt Llama3.1-8B-Instruct to generate tool calls and outputs and have human experts review validity, following common practice in agent benchmark construction, as also adopted in T1 (Chakraborty et al., 2025) and BFCL (Patil et al.).

Agents for Downstream Evaluation We conduct downstream evaluations on each dataset using LLM-based agents under both real and synthetic tasks. Each agent is executed by providing the full prior context—including instructions, responses, tool calls, and their results—as input to the model. We include three LLMs for evaluation: gemma-3-1b-it, Qwen3-4B-Instruct, and Llama3.1-8B-Instruct. For the LLM-as-a-judge component of the evaluation, we use Mistral-7B-Instruct.

4.1 EXPERIMENTAL RESULTS

We present the evaluation results of different synthetic data generation methods on the T1 and BFCL datasets in Figs. 3 to 5, and defer the detailed numerical results to §C.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

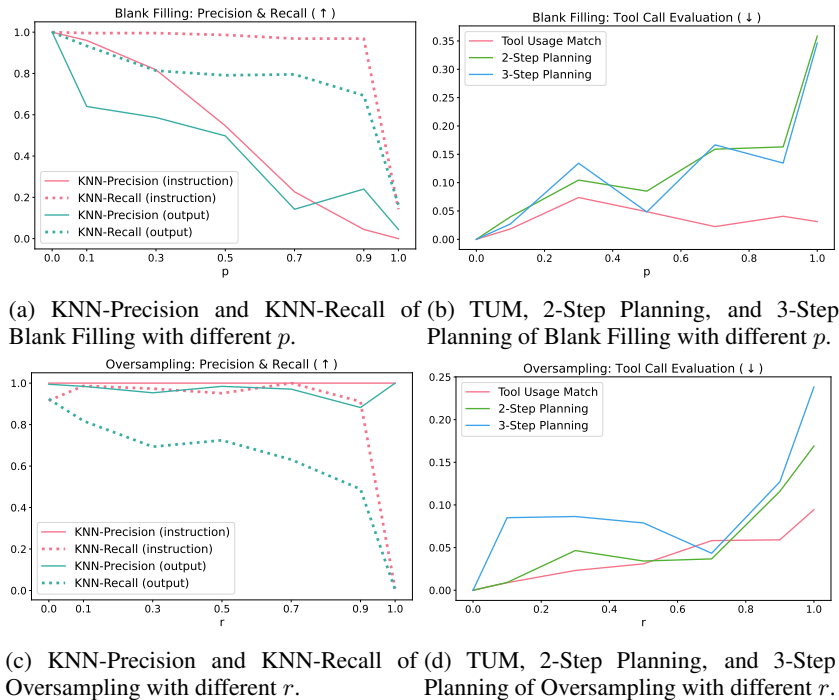
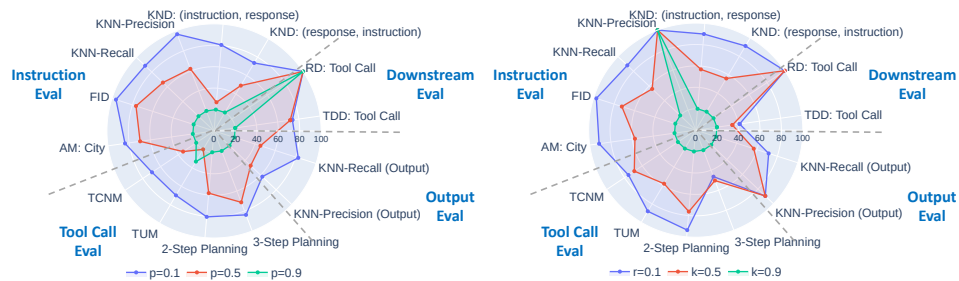


Figure 3: Performance of Blank Filling and Oversampling under T1.

Fig. 3 shows Blank Filling and Oversampling on T1 across different parameter settings; Figs. 4 and 5 further compare Blank Filling, Oversampling, and In-Context Generation on T1 and BFCL via radar plots. To improve visualization, we rescale all evaluation metrics to the range $[20, 100]$, where 20 corresponds to the worst performance observed across all compared methods and 100 represents the performance upper bound (e.g., KNN-Precision = 1). The figures indicate several main takeaways:

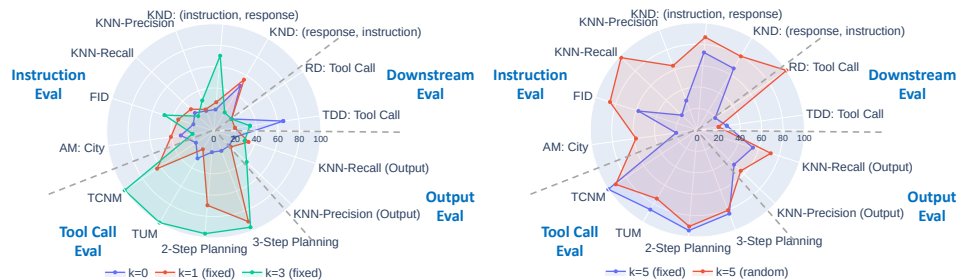
- *ESDAE can captures fine-grained variations in both the quality and diversity of synthetic data.* For Blank Filling (shown in Fig. 3a), as p increases, the KNN-Precision of both instructions and outputs decreases from nearly 1 to close to 0, indicating progressive degradation in task quality. This trend is consistent with the intuition that a higher masking probability leads to lower-fidelity synthetic data. In contrast, when $p \leq 0.9$, the KNN-Recall for both instructions and outputs remains relatively high (above 0.7), suggesting that Blank Filling largely preserves task diversity. For Oversampling (shown in Fig. 3c), as r increases, the output KNN-Recall drops significantly, while the KNN-Precision of both instructions and outputs remains near perfect, consistent with the intuition that duplication reduces task diversity while largely preserving task quality. Furthermore, as shown in Figs. 3b, 3d, 4a and 4b, increasing p or r leads to performance degradation for Blank Filling and Oversampling across nearly all evaluation metrics. This trend aligns with our intuition that higher masking probabilities or heavier duplication results in worse synthetic data.
- *Adding more in-context examples does not necessarily improve In-Context Generation, whereas Randomizing in-context examples helps synthetic data better match the semantic and statistical distributions of the real tasks.* For In-Context Generation (see Figs. 4c and 5a), increasing the number of in-context examples k does not consistently improve synthetic data performance. For instance, on the T1 dataset, $k = 1$ outperforms $k = 3$ under several metrics, including KNN-Recall, KND, AM, and TDD. We then fix the number of in-context examples k and compare fixed versus randomized example selection across generations in Figs. 4d and 5b. Randomizing in-context examples consistently improves performance on instruction- and output-level metrics, including task quality and diversity, as well as semantic and statistical properties.
- *No single metric can fully characterize synthetic data performance.* A method with one parameter setting may outperform another on some metrics while underperforming on others, for example, KNN-Precision versus KNN-Recall in Fig. 4c. Another example is tool-call metrics and RD: in Fig. 4b, increasing r degrades tool-call metrics, indicating a growing discrepancy in tool-planning patterns between real and synthetic tasks. In contrast, RD: Tool Call remains strong, since higher-

378
379
380
381
382
383
384
385
386
387



(a) Blank Filling. As p increases, data utility degrades across nearly all metrics. (b) Oversampling. As r increases, data utility degrades across nearly all metrics.

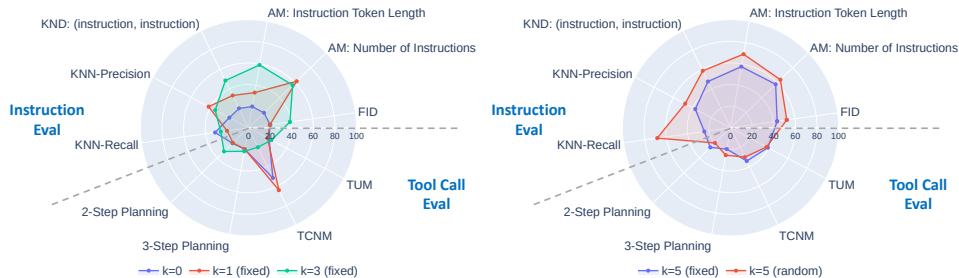
390
391
392
393
394
395
396
397
398



(c) In-Context Generation with different k . Increasing k does not consistently improve data utility across several metrics, including KNN-Recall, KND, and AM. (d) In-Context Generation with $k = 5$. Compared to fixed in-context examples, randomized in-context examples across generations lead to higher data utility under both instruction and output evaluation metrics.

Figure 4: Performance of Blank Filling, Oversampling, and In-Context Generation under T1.

404
405
406
407
408
409
410
411
412
413
414



(a) In-Context Generation with different k . (b) In-Context Generation with $k = 5$.

Figure 5: Performance of In-Context Generation under BFCL.

418
419
420
421

capacity models continue to outperform lower-capacity ones even when planning patterns shift, resulting in largely consistent model rankings. This further motivates the need for ESDAE.

422
423

5 CONCLUSION

424
425
426
427
428
429
430
431

In this paper, we introduced ESDAE, a multi-axis framework for evaluating how faithfully synthetic trajectory datasets match real data across task instructions and responses, tool calls, final outputs, and downstream agent performance. Through experiments on recent agent benchmarks with controlled synthetic generation schemes, ESDAE consistently detects fine-grained failures in fidelity and diversity. Overall, the results show that no single metric is sufficient to characterize synthetic benchmark quality, motivating comprehensive evaluations such as ESDAE before relying on synthetic data for agent testing. For future work, we plan to extend ESDAE to explicitly measure the data augmentation effect of synthetic trajectories by designing additional metrics that quantify marginal gains in coverage, robustness, and generalization improvements.

432
433
434
435
436
437
438
439
440
441
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
475
476
477
478
479
480
481
482
483
484
485

REFERENCES

- Omar Alonso and Kenneth Church. Evaluating the evaluations: A perspective on benchmarks. In *ACM SIGIR Forum*, volume 58, pp. 1–27. ACM New York, NY, USA, 2025.
- Anthropic. Demystifying evals for ai agents. <https://www.anthropic.com/engineering/demystifying-evals-for-ai-agents>, 2026.
- CapitalOne. Synthetic data matters for machine learning innovation. <https://www.capitalone.com/tech/machine-learning/synthetic-data-research/>, 2022.
- Amartya Chakraborty, Pareshe Dashore, Nadia Bathaee, Anmol Jain, Anirban Das, Shi-Xiong Zhang, Sambit Sahu, Milind Naphade, and Genta Indra Winata. T1: A tool-oriented conversational dataset for multi-turn agentic planning. *arXiv preprint arXiv:2505.16986*, 2025.
- Alexander Gill, Abhilasha Ravichander, and Ana Marasović. What has been lost with synthetic evaluation? *arXiv preprint arXiv:2505.22830*, 2025.
- Shadi Iskander, Sofia Tolmach, Ori Shapira, Nachshon Cohen, and Zohar Karnin. Quality matters: Evaluating synthetic data for tool-using llms. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 4958–4976, 2024.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- Sayash Kapoor, Benedikt Stroebel, Peter Kirgis, Nitya Nadgir, Zachary S Siegel, Boyi Wei, Tianci Xue, Ziru Chen, Felix Chen, Saiteja Utpala, et al. Holistic agent leaderboard: The missing infrastructure for ai agent evaluation. *arXiv preprint arXiv:2510.11977*, 2025.
- Renhao Li, Jianhong Tu, Yang Su, Yantao Liu, Fei Huang, Hamid Alinejad-Rokny, Derek F Wong, Junyang Lin, and Min Yang. Toolrm: Towards agentic tool-use reward modeling. *arXiv preprint arXiv:2510.26167*, 2025.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.
- Xing Han Lù, Amirhossein Kazemnejad, Nicholas Meade, Arkil Patel, Dongchan Shin, Alejandra Zambrano, Karolina Stańczak, Peter Shaw, Christopher J Pal, and Siva Reddy. Agentrewardbench: Evaluating automatic evaluations of web agent trajectories. *arXiv preprint arXiv:2504.08942*, 2025.
- Gaurav Maheshwari, Dmitry Ivanov, and Kevin El Haddad. Efficacy of synthetic data as a benchmark. *arXiv preprint arXiv:2409.11968*, 2024.
- Michael Majurski and Cynthia Matuszek. Grounding synthetic data evaluations of language models in unsupervised document corpora. *arXiv preprint arXiv:2505.08905*, 2025.
- Sohum Mehta and Saaketh Bhojanam. Prompt genotyping: Quantifying the evaluation gap between synthetic benchmarks and real llm performance. In *NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling*.
- Mahmoud Mohammadi, Yipeng Li, Jane Lo, and Wendy Yip. Evaluation and benchmarking of llm agents: A survey. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 6129–6139, 2025.
- Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous evaluation and refinement of digital agents. *arXiv preprint arXiv:2404.06474*, 2024.
- Melissa Z Pan, Negar Arabzadeh, Riccardo Cogo, Yuxuan Zhu, Alexander Xiong, Lakshya A Agrawal, Huanzhi Mao, Emma Shen, Sid Pallerla, Liana Patel, et al. Measuring agents in production. *arXiv preprint arXiv:2512.04123*, 2025.

486 Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and
487 Joseph E Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agen-
488 tic evaluation of large language models. In *Forty-second International Conference on Machine*
489 *Learning*.

490 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru
491 Tang, Bill Qian, et al. Toollm: Facilitating large language models to master 16000+ real-world
492 apis. *arXiv preprint arXiv:2307.16789*, 2023.

493
494 Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng
495 Li, and Yueting Zhuang. Taskbench: Benchmarking large language models for task automation.
496 *Advances in Neural Information Processing Systems*, 37:4540–4574, 2024.

497
498 Alex Tamkin, Miles McCain, Kunal Handa, Esin Durmus, Liane Lovitt, Ankur Rathi, Saffron
499 Huang, Alfred Mountfield, Jerry Hong, Stuart Ritchie, et al. Clio: Privacy-preserving insights
500 into real-world ai use. *arXiv preprint arXiv:2412.13678*, 2024.

501
502 Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. Toolal-
503 paca: Generalized tool learning for language models with 3000 simulated cases. *arXiv preprint*
504 *arXiv:2306.05301*, 2023.

505
506 Shuaiqi Wang, Vikas Raunak, Arturs Backurs, Victor Reis, Pei Zhou, Sihao Chen, Longqi Yang,
507 Zinan Lin, Sergey Yekhanin, and Giulia Fanti. Struct-bench: A benchmark for differentially
508 private structured text generation. In *The Thirty-ninth Annual Conference on Neural Information*
509 *Processing Systems Datasets and Benchmarks Track*.

510
511 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh J Hua,
512 Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents
513 for open-ended tasks in real computer environments. *Advances in Neural Information Processing*
514 *Systems*, 37:52040–52094, 2024.

515
516 Lang Xiong, Nishant Bhargava, Jianhang Hong, Jeremy Chang, Haihao Liu, Vasu Sharma, and
517 Kevin Zhu. Probe-rewrite-evaluate: A workflow for reliable benchmarks and quantifying evalua-
518 tion awareness. *arXiv preprint arXiv:2509.00591*, 2025.

519
520 Tianci Xue, Weijian Qi, Tianneng Shi, Chan Hee Song, Boyu Gou, Dawn Song, Huan Sun, and
521 Yu Su. An illusion of progress? assessing the current state of web agents. *arXiv preprint*
522 *arXiv:2504.01382*, 2025.

523
524 John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standardizing
525 and benchmarking interactive coding with execution feedback. *Advances in Neural Information*
526 *Processing Systems*, 36:23826–23854, 2023.

527
528 Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for
529 tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024.

530
531 Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun Zhao, Roy Bar-Haim, Arman Cohan,
532 and Michal Shmueli-Scheuer. Survey on evaluation of llm-based agents. *arXiv preprint*
533 *arXiv:2503.16416*, 2025.

534
535 Zeyu Zhang, Guohao Li, Zhenchang Xing, Alexandros Apostolopoulos, Yu Lin Lee, and Liang
536 Zheng. Gecko: A simulation environment to ground agent tool calls with stateful feedback for
537 refinement.

538
539 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for build-
ing autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

Yuxuan Zhu, Tengjun Jin, Yada Pruksachatkun, Andy Zhang, Shu Liu, Sasha Cui, Sayash Kapoor,
Shayne Longpre, Kevin Meng, Rebecca Weiss, et al. Establishing best practices for building
rigorous agentic benchmarks. *arXiv preprint arXiv:2507.02825*, 2025.

540 A PROMPTS FOR LLM-AS-A-JUDGE IN DOWNSTREAM EVALUATION

541

542 A.1 PROMPTS FOR TOOL CALL EVALUATION

543

544 System Prompt:

545

```
546 1 You are an evaluator. You must answer with ONLY 'yes' or 'no'. Never
547 provide explanations or reasoning.
```

548

549 User Prompt:

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

```
1 Compare these tool call sequences for an attractions recommendation
assistant:
2
3 Conversation:
4 {instr}
5
6 Expected tool calls:
7 {expected_tool_call}
8
9 Actual tool calls:
10 {actual_tool_call}
11
12 Evaluation rules:
13 - Actual must accomplish the same goal as expected
14 - Semantic equivalence is OK (e.g., "OR" vs "Oregon", reordered
operations with same result)
15 - Partial correctness = NO
16 - Different variable/cache names = OK if functionality identical
17
18 Does actual correctly implement the conversation based on expected?
19 Answer (yes/no):
```

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

570 A.2 PROMPTS FOR OUTPUT EVALUATION

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

```
1 Compare these outputs for an attractions recommendation assistant:
2
3 Conversation:
4 {instr}
5
6 Expected output:
7 {expected_output}
8
9 Actual output:
10 {actual_output}
11
12 Evaluation rules:
13 - Actual must convey same information and meaning as expected
14 - Different wording is OK if content equivalent
15 - Partial correctness = NO
16 - Focus on semantic content, not syntax
17
18 Does actual correctly respond to the conversation based on expected?
19 Answer (yes/no):
```

594 B PROMPTS FOR LLM-BASED SYNTHETIC DATA GENERATION METHODS

595

596

B.1 PROMPTS FOR T1

597

598

B.1.1 BLANK FILLING

599

600 System Prompt:

601

1 You are a helpful conversation generator. When given a conversation with
2 blanks (underscores), fill them in naturally. IMPORTANT RULES:

603 2 1. The conversation MUST start with 'assistant:' (not 'Assistant:' or any
604 variation)

605 3 2. Lines MUST alternate strictly between 'user:' and 'assistant:'

606 4 3. Each line must follow the format: 'role: content' where role is either
607 'user' or 'assistant'

608 5 4. Output ONLY the completed conversation with no preamble, explanation,
609 or extra text

610 6 5. Maintain the same number of conversation turns as the input

610

611 User Prompt:

612

613 1 Example input for fill in the blanks:

614 2

615 3 assistant: H_____ What ____ of attractions are you looking for? Are you
616 interested in _____, a__, or something else?

617 4 user: I'm interested in ___ and ___ attractions in __.

618 5 assistant: G_ has a lot to offer. Are you looking at specific _____ or
619 re_____?

620 6 user: Yeah, I'm thinking of visiting Fre__ and M_____.

621 7 assistant: Both _____o and _____his have great A__ and S_____ attractions.
622 Let me tell you about some of them.

623 8 user: That sounds ____.

624 9

625 10 Completed conversation for example input:

626 11 assistant: Hello! What kind of attractions are you looking for? Are you
627 interested in history, art, or something else?

628 12 user: I'm interested in Art and Scenic attractions in GA.

629 13 assistant: GA has a lot to offer. Are you looking at specific cities or
630 regions?

631 14 user: Yeah, I'm thinking of visiting Fresno and Memphis.

632 15 assistant: Both Fresno and Memphis have great Art and Scenic attractions.
633 Let me tell you about some of them.

634 16 user: That sounds great.

635 17

636 18 Now fill in the blanks to complete this conversation:

637 19

638 20 {masked_conv}

639 21 Completed conversation:

636

637

638 B.1.2 IN-CONTEXT GENERATION

639

640 System Prompt: Same as the system prompt in Section B.1.1.

641

642 User Prompt:

643

1 Here are example conversations:

644 2 {conv_list_str}

645 3

646 4 Generate 1 new similar conversation that follows the same structure.

647 5 Do not include anything other than this conversation.

6 Similar conversation:

648 **B.2 PROMPTS FOR BFCL**

649
650 **B.2.1 BLANK FILLING**

651 **System Prompt:**

652
653 1 You are a request completion assistant. Fill in blanks using ONLY the
654 provided APIs and Resources. Output only the completed requests, one
655 per line. Do not add explanations or extra text.

656
657 **User Prompt:**

658 1 Fill in the blanks (underscores) to complete the user requests.
659 2 Example:
660 3 {example_context}
661 4
662 5 Input with blanks:
663 6 {example_masked}
664 7
665 8 Completed requests:
666 9 {example_completed}
667 10
668 11 Now fill in the blanks for this:
669 12 APIs: {target_classes}
670 13 Resources: {target_resources}
671 14
672 15 Input with blanks:
673 16 {masked_conv_str}
674 17
675 18 Completed requests:

676
677 **B.2.2 IN-CONTEXT GENERATION**

678 **System Prompt:**

679 1 You are a test case generator. Output each request as:
680 2 Request 1: <request text>
681 3 Request 2: <request text>

682 **User Prompt:**

683 {examples_str_with_APIs_resources_requests}
684 2
685 3 Generate EXACTLY {n_turns} requests for the situation below.
686 4 Requests can query information, perform operations, modify resources, or
687 search/filter data.
688 5 Keep requests realistic. Later requests should build on earlier ones.
689 6
690 7 ONLY refer to the APIs and Resources below in the requests.
691 8 APIs: {target_classes}
692 9 Resources: {target_resources}
693 10 Requests:

694
695 **C DETAILED EVALUATION RESULTS**

696
697
698
699
700
701

702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

Table 2: Synthetic data generation evaluation results on ESDAE under T1.

Method	Parameter	Instruction Eval		Tool Call Eval		Output Eval		Downstream Eval	
		KND ↓	AM ↓	TUM ↓	3-Step Planning ↓	KNN-Precision ↑	KNN-Recall ↑	TDD ↓	RD ↑
Blank Filling	$p = 0$	0	0	0	0	1	1	0	1.0
	$p = 0.1$	0.019	0.113	0.018	0.027	0.640	0.933	0.052	1.0
	$p = 0.3$	0.025	0.096	0.074	0.134	0.587	0.813	0.090	1.0
	$p = 0.5$	0.037	0.213	0.049	0.048	0.498	0.791	0.056	1.0
	$p = 0.7$	0.035	0.295	0.022	0.167	0.142	0.796	0.147	0.5
	$p = 0.9$	0.059	0.564	0.041	0.135	0.240	0.693	0.160	1.0
	$p = 1$	0.075	0.746	0.031	0.346	0.044	0.164	0.188	0.5
Oversampling	$r = 0$	0.004	0.051	0	0	0.994	0.923	0.095	1.0
	$r = 0.1$	0.012	0.073	0.009	0.085	0.984	0.818	0.048	1.0
	$r = 0.3$	0.038	0.229	0.023	0.086	0.953	0.693	0.117	0.5
	$r = 0.5$	0.060	0.351	0.031	0.079	0.984	0.724	0.053	1.0
	$r = 0.7$	0.087	0.511	0.058	0.043	0.971	0.631	0.156	1.0
	$r = 0.9$	0.110	0.658	0.059	0.127	0.882	0.489	0.165	0.5
	$r = 1$	0.122	0.750	0.094	0.238	0.992	0.004	0.183	1.0
Generative Models with In-Context Examples	$k = 0$	0.075	0.746	0.031	0.346	0.044	0.164	0.033	0.5
	$k = 1$ (fixed)	0.065	0.645	0.036	0.040	0.067	0.311	0.076	0.5
	$k = 1$ (random)	0.054	0.444	0.017	0.008	0.027	0.729	0.126	1.0
	$k = 3$ (fixed)	0.117	0.866	0.011	0.016	0.311	0.271	0.062	0.5
	$k = 3$ (random)	0.021	0.523	0.038	0.034	0.062	0.853	0.121	0.5
	$k = 5$ (fixed)	0.031	0.875	0.005	0.070	0.369	0.524	0.124	0.5
	$k = 5$ (random)	0.019	0.460	0.010	0.085	0.471	0.707	0.138	1.0

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Table 3: Synthetic data generation evaluation results on ESDAE under BFCL.

Method	Parameter	Instruction Eval				Tool Call Eval			
		KND ↓	AM ↓	KNN-Precision ↑	KNN-Recall ↑	TUM ↓	TCNM ↓	2-Step Planning ↓	3-Step Planning ↓
Blank Filling	$p = 0$	0	0	1	1	0	0	0	0
	$p = 0.1$	0.0871	0.665	0.715	0.900	0.4522	6.805	0.6840	0.5951
	$p = 0.3$	0.0793	1.030	0.675	0.905	0.3877	5.445	0.6835	0.6321
	$p = 0.5$	0.0869	1.320	0.640	0.920	0.4382	5.895	0.7085	0.6227
	$p = 0.7$	0.1055	1.470	0.430	0.910	0.6557	16.335	0.6540	0.5671
	$p = 0.9$	0.1630	1.445	0.225	0.805	0.4714	6.335	0.7598	0.5829
	$p = 1$	0.1238	0.695	0.245	0.290	0.6005	6.515	0.7469	0.5851
Oversampling	$r = 0$	0.0197	0.110	1.000	0.875	0.1116	0.285	0.1560	0.1476
	$r = 0.1$	0.0235	0.130	0.995	0.980	0.1503	0.485	0.1223	0.0847
	$r = 0.3$	0.0331	0.305	0.995	0.965	0.3841	1.395	0.2614	0.1796
	$r = 0.5$	0.0530	0.525	0.995	0.920	0.5832	2.175	0.3712	0.2583
	$r = 0.7$	0.0803	0.725	0.995	0.910	0.7390	2.985	0.4664	0.3150
	$r = 0.9$	0.1060	0.955	0.995	0.950	0.8500	3.955	0.5566	0.3681
	$r = 1$	0.1179	1.060	1.000	0.005	0.9116	4.290	0.6275	0.4111
Generative Models with In-Context Examples	$k = 0$	0.1238	0.695	0.245	0.290	0.6005	6.515	0.7469	0.5851
	$k = 1$ (fixed)	0.1033	0.330	0.450	0.175	0.6133	4.820	0.7520	0.5862
	$k = 1$ (random)	0.1293	0.570	0.370	0.760	0.5708	5.045	0.6983	0.5962
	$k = 3$ (fixed)	0.0794	0.375	0.385	0.235	0.5818	10.770	0.6442	0.5710
	$k = 3$ (random)	0.1147	0.515	0.500	0.765	0.5572	7.430	0.6989	0.5743
	$k = 5$ (fixed)	0.0810	0.365	0.405	0.330	0.4569	5.360	0.6905	0.6208
	$k = 5$ (random)	0.0640	0.312	0.505	0.720	0.4674	5.685	0.7454	0.5765