DataDreamer: A Tool for Synthetic Data Generation and Reproducible LLM Workflows

Anonymous ACL submission

Abstract

001 Large language models (LLMs) have become a dominant and important tool for NLP researchers in a wide range of tasks. Today, many researchers use LLMs in synthetic data generation, task evaluation, fine-tuning, distillation, and other model-in-the-loop research workflows. However, challenges arise when using these models that stem from their scale, their closed source nature, and the lack of standardized tooling for these new and emerging workflows. The rapid rise to prominence of these models and these unique challenges has had immediate adverse impacts on open science and on the reproducibility of work that uses them. In this ACL 2024 theme track paper, we introduce DataDreamer, an open source Python library that allows researchers to write simple 017 code to implement powerful LLM workflows. DataDreamer also helps researchers adhere to best practices that we propose to encourage open science and reproducibility. 021

1 Introduction

024

037

While large language models (LLMs) have established a new era in NLP research through the prompt-and-predict paradigm that has proven effective on a wide variety of tasks, the use of these models has come with significant drawbacks (Liu et al., 2023). Many popular models like GPT-4 (OpenAI et al., 2023) are closed source and behind a remote API, while running models locally can be technically complex and expensive due to their scale. Moreover, the now well-established prompting paradigm can be brittle with results widely varying between different models, configurations, and environments (Sclar et al., 2023; Jaiswal et al., 2023). These challenges have made it difficult for researchers to share, reproduce, extend, and compare work, hindering the rate of research progress.

In context of the rapid shift to using these large models in research, this year's 2024 ACL theme



Figure 1: DataDreamer helps researchers implement many types of LLM workflows easier and makes reproducibility automatic and simple. These workflows often involve synthetic data generation with a LLM-in-theloop and/or fine-tuning, aligning, and distilling models.

track calls for "stimulating discussion about open science and reproducible NLP research, as well as supporting the open source software movement" and invites "high-quality open source software implementations".⁴ In concordance with this theme, we introduce DataDreamer, our open source Python package that provides both practical utility to researchers and scientific utility to the community:

⁴https://2024.aclweb.org/calls/main_conferenc e_papers/#theme-track-open-science-open-data-and -open-models-for-reproducible-nlp-research

DataDreamer helps researchers implement state-of-the-art emerging workflows involving LLMs such as synthetic data generation, fine-tuning, instruction-tuning, and alignment. It simplifies implementations by providing a single library with a standardized interface for many of these tasks while reducing technical complexity around switching between models, caching, resumability, logging, multi-GPU inference and training, using adapter and quantization optimizations, and publishing open datasets and models.

061

065

066

- DataDreamer makes chaining data between tasks, an increasingly common practice, simple. For example, a user can generate data with a synthetic data workflow and then finetune on that synthetic data.
- DataDreamer helps researchers implement workflows while crucially producing output that is compatible with open science and reproducible ideals with minimal effort, through automatic caching, reproducibility fingerprints, and more best-practice artifacts.

2 LLM Workflows

To motivate DataDreamer, we first discuss the LLM workflows that it supports. We discuss challenges to open science that arise from these usage patterns. In this paper, we do not seek to validate or critique these approaches. Instead, we offer a solution to implement them and make them reproducible. These LLM workflows are often used in combination with each other (Yuan et al., 2024), and orchestration of multi-stage workflows is frequently implemented through multiple shell or Python scripts. Reproducing these multi-stage workflows is challenging as shell scripts may rely upon a particular author's job scheduler or environment and require execution in a specific order. In Section 4 and 5, we discuss how DataDreamer's task orchestration, caching system, and simple multi-GPU training make it easier to implement these multi-stage workflows in a single Python program, minimizing these issues.

072

074

075

076

077

078

079

081

084

090

092

093

Synthetic Data Generation Recent work has explored using LLMs to create synthetic data for tasks or to augment existing datasets to boost task performance (Yu et al., 2023; Kumar et al., 2020a,b;

Feature	LangChain ¹	Axlotl ²	HF Transformers + TRL ³	DataDreamer
Implementation				
Accessible via Python API Built for Researchers	√ X	× ×	√ √	\ \
Integrations				
Open Source Models Commercial & API-based Models	√ √	✓ ×	✓ ×	\ \
Tasks				
Prompting & Prompt "Chaining" Synthetic Data Generation & Augmentation Fine-tuning LLMs Instruction-tuning LLMs Aligning LLMs Training Classifier Models Training Embedding Models	✓ ✓ × × × ×	×	× × ✓ ✓ ✓ ×	\ \ \ \ \ \ \
Conveniences				
Caching Resumability Simplifies Boilerplate Code (tokenization, etc.) Simplifies Multi-GPU Inference and Training Publishing Datasets & Models	○ × × × ×	×	× ○ × × ✓	
Open Science and Reproducibility				
Reproducibility Fingerprints Saves Intermediate Outputs Synthetic Data and Model Cards	× × ×	× × ×	× × ×	\$ \$ \$

Table 1: We compare feature coverage between other popular libraries and solutions available to researchers today that target similar workflows. DataDreamer integrates these features into a single library with a standardized interface making experimentation and chaining data between tasks simple. (X = No; V = Yes; \bigcirc = Partial Support)

Yoo et al., 2021; Han et al., 2021a; Ye et al., 2022; 095 Honovich et al., 2022, inter alia). Synthetic data generation involves using a LLM once or multiple times in a multi-stage workflow to process data, sometimes referred to as "chaining" (Rush, 2023). When prompting LLMs to generate or augment 100 datasets, a reproducibility challenge that arises is 101 "prompt sensitivity" where even small variations in 102 a prompt can lead to significantly different results 103 (Sclar et al., 2023). Moreover, it is imperative to tag 104 synthetically generated datasets because of model degradation concerns (Shumailov et al., 2023). 106

107LLMs for Task EvaluationAnother increas-108ingly common workflow is using LLMs as judges109or as automatic metrics for evaluating a model's110performance on a task (Zheng et al., 2023; Fu et al.,1112023; Dubois et al., 2023; Chiang and Lee, 2023,112*inter alia*). Many of the reproducibility challenges113applicable to synthetic data also arise here.

Fine-tuning and Alignment Another common 114 workflow is the creation of task-specific expert 115 models using knowledge from larger models to cre-116 ate smaller, more efficient models via fine-tuning 117 118 and distillation (Han et al., 2021b; Liu et al., 2022; Hsieh et al., 2023). Instruction-tuning is fine-tuning 119 that allows base pre-trained models to better follow 120 natural language human instruction and improve 121 their generalized task performance (Ouyang et al., 122 2022; Wei et al., 2021; Sanh et al., 2021; Mishra 123 et al., 2021). Closely related, alignment techniques 124 steer model responses towards those more prefer-125 able to humans (Stiennon et al., 2020; Bai et al., 126 2022; Rafailov et al., 2023). Implementing resuma-127 128 bility and efficient training techniques are practical challenges often faced. Reproducibility challenges 129 include sharing exact data and hyperparameters. 130

Self-improving LLMs Self-improving LLMs through self-feedback training loops is an increasingly active area of research interest (Huang et al., 2022; Wang et al., 2022; Li et al., 2023; Chen et al., 2024; Yuan et al., 2024; Gunasekar et al., 2023). These workflows can be uniquely complex to both implement and reproduce due to requiring multiple rounds that chain together synthetic data generation, automatic evaluation, and model re-training. DataDreamer supports all of these workflows and makes it simple to chain data between them.

131

132

133

136

138

139

140

141

3 Demonstration and Examples

Before delving into the structure and implementation of DataDreamer, we first provide a simple demonstration of DataDreamer's capabilities and API through an example synthetic data generation and distillation workflow in Example 1. The LLM used in this example is GPT-4 (OpenAI et al., 2023). As an initial step, the example uses the LLM to generate 1,000 NLP research paper abstracts. The LLM is then used to summarize those abstracts in a tweet-like style. These two steps result in a fully synthetic dataset of abstracts and tweets summarizing them. Using a trainer, this synthetic dataset is then distilled to a small, local model that is capable of summarizing paper abstracts in a tweet-like style. As a final step, the example demonstrates how both the synthetic dataset and the trained model can be published and shared. For illustrative purposes, we demonstrate a sample generation of the trained model's output on this paper's abstract:

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

165

166

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

"Introducing DataDreamer, an open source Python library for advanced #NLP workflows. It offers easy code to create powerful LLM workflows, addressing challenges in scale, closed source nature, and tooling. A step towards open science and reproducibility! #AI #MachineLearning"

Further example workflows can be found in the Appendix (Example 2, Example 3, Example 4, Example 5).

4 DataDreamer

DataDreamer is an open source Python package that allows researchers to implement all of the LLM workflows discussed in Section 2 using a single library. DataDreamer provides a standardized interface for prompting and training models, abstracting away vendor-specific libraries and tooling. This makes research code simpler to implement, modify, experiment with, and share with others. DataDreamer integrates with other open source LLM libraries like transformers (Wolf et al., 2019) and trl (von Werra et al., 2020), as well as commercial model APIs like OpenAI and Anthropic⁵ for commercial LLMs (Brown et al., 2020). Moreover, DataDreamer automatically implements the best practices for reproducibility discussed in Section 5.

¹https://github.com/langchain-ai/langchain

²https://github.com/OpenAccess-AI-Collective/ axolotl

³Wolf et al. (2019); von Werra et al. (2020) ⁵https://www.anthropic.com/

```
1 from datadreamer import DataDreamer
2 from datadreamer.llms import OpenAI
3 from datadreamer.steps import DataFromPrompt, ProcessWithPrompt
4 from datadreamer.trainers import TrainHFFineTune
5 from peft import LoraConfig
  with DataDreamer("./output"):
7
      # Load GPT-4
8
      gpt_4 = OpenAI(model_name="gpt-4")
9
10
      # Generate synthetic arXiv-style research paper abstracts with GPT-4
      arxiv_dataset = DataFromPrompt(
12
           "Generate Research Paper Abstracts",
14
          args={
               "llm": gpt_4,
15
               "n": 1000,
16
               "temperature": 1.2,
               "instruction": (
18
                 "Generate an arXiv abstract of an NLP research paper."
19
                 " Return just the abstract, no titles."
20
21
               ),
22
          },
          outputs={"generations": "abstracts"},
23
      )
24
25
      # Use GPT-4 to convert the abstracts to tweets
26
27
      abstracts_and_tweets = ProcessWithPrompt(
           "Generate Tweets from Abstracts",
28
           inputs={"inputs": arxiv_dataset.output["abstracts"]},
29
          args={
30
             "llm": gpt_4,
31
             "instruction":
                            "Given the abstract, write a tweet to summarize the work.",
32
             "top_p": 1.0,
33
34
          },
          outputs={"inputs": "abstracts", "generations": "tweets"},
35
36
      )
37
38
      # Create training data splits
      splits = abstracts_and_tweets.splits(train_size=0.90, validation_size=0.10)
39
40
      # Train a model to convert research paper abstracts to tweets with the
41
      # synthetic dataset
42
      trainer = TrainHFFineTune(
43
           "Train an Abstract => Tweet Model".
44
          model_name="google/t5-v1_1-base",
45
          peft_config=LoraConfig(),
46
47
      )
48
      trainer.train(
49
           train_input=splits["train"].output["abstracts"],
           train_output=splits["train"].output["tweets"];
50
          validation_input=splits["validation"].output["abstracts"],
51
52
          validation_output=splits["validation"].output["tweets"],
          epochs=30,
53
54
          batch_size=8,
55
      )
56
      # Publish and share the synthetic dataset
57
      abstracts_and_tweets.publish_to_hf_hub("repo_id")
58
59
      # Publish and share the trained model
60
61
      trainer.publish_to_hf_hub("repo_id")
```

Example 1: In this demonstration snippet, DataDreamer generates a fully synthetic dataset of tweets summarizing research paper abstracts and then trains a smaller T5 distilled model (Raffel et al., 2020) to perform the task and publishes both the synthetic dataset and the trained model. DataDreamer makes it simple to chain data from each step in the workflow to the next and automatically caches each step of this workflow to the ./output/ folder to allow interruption and resumability at any point in the script. The standardized API also makes it easy to switch to and experiment with different models, both open source and commercial, for generation and training.

185 186

187

189

190

191

192

193

194

195

196

199

205

207

210

211

212

213

214

215

216

218

219

220

4.1 Installation

DataDreamer can be installed with the wheel file:

pip install <WHEEL_FILE>

4.2 Sessions

All code using the DataDreamer library is placed within a "session" using a Python context manager instantiated using the with keyword:

from datadreamer import DataDreamer

with DataDreamer("./output"):
 ...

Workflow tasks can be run within the session context manager. These tasks are called "steps" (loading a dataset, prompting a model, etc.) or "trainers". The session allows DataDreamer to automatically organize the resulting datasets, outputs, caches, training checkpoints, and trained models that result from tasks run within the session into the ./output/ folder. Each step in a workflow assigns a custom descriptive name for its subfolder under ./output/. DataDreamer sessions automatically provide user-friendly logging around workflow tasks run within the session (see Figure 2).

4.3 Steps

Steps are the core operators in a DataDreamer session. A step in DataDreamer transforms from an input dataset to an output dataset (Lhoest et al., 2021). This is useful for tasks like generating synthetic data from LLMs, or data augmentation for existing datasets. The output of one step can be directly used as the input to another step or as the input to a trainer, allowing users to chain together multiple steps/trainers to create complex workflows. DataDreamer comes with a number of built-in steps for common operations in LLM workflows, some examples of which can be seen in Table 2. Useful standard data processing operations such as .map(), .filter(), and .shuffle() can also quickly be applied to the output of a step for custom processing. DataDreamer uses memory-mapping to handle large datasets stored on disk and can be run lazily over iterable, streaming datasets.

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

237

238

239

240

241

242

243

244

245

247

248

250

251

252

253

254

255

4.4 Models

Models can be loaded in a DataDreamer session and then be passed as an argument to steps like FewShotPrompt and ProcessWithPrompt. DataDreamer creates a standardized interface for accessing open source and commercial LLMs. It includes interfaces for embedding models as well as LLMs. Examples of supported models and model providers can be found in Table 2.

4.5 Trainers

Trainers can train on a dataset produced by a step in a DataDreamer workflow. The dataset may be loaded from an external source or produced as the output of a step in a multi-step workflow. DataDreamer's trainers support a wide variety of techniques and tasks including fine-tuning, instruction-tuning, alignment via RLHF (Ouyang et al., 2022) and DPO (Rafailov et al., 2023), distillation, training classifiers, and training embedding models. Examples of supported techniques are shown in Table 2.

4.6 Caching and Sharing Workflows

Caching has practical utility in LLM workflows as these large models can be both computationally and financially expensive to run. Therefore, eliminating re-computation can save both time and resources. Caching in DataDreamer happens at multiple levels. When a step or trainer is completed, its resulting dataset or trained model is saved to disk and loaded

Туре		Examples
Steps	Load a Dataset	DataSource, HFHubDataSource, JSONDataSource, CSVDataSource,
	Prompting	Prompt, RAGPrompt, ProcessWithPrompt, FewShotPrompt, DataFromPrompt, DataFromAttributedPrompt, FilterWithPrompt, RankWithPrompt, JudgeGenerationPairsWithPrompt,
	Other	Embed, Retrieve, CosineSimilarity,
Models		OpenAI, OpenAIAssistant, HFTransformers, CTransformers, VLLM, Petals, HFAPIEndpoint, Together, MistralAI, Anthropic, Cohere, AI21, Bedrock, Vertex,
Trainers		TrainOpenAIFineTune, TrainHFClassifier, TrainHFFineTune, TrainSentenceTransformer, TrainHFDPO, TrainHFPPO,

Table 2: A few examples of built-in steps, models, and trainers available in DataDreamer.

from disk if the step or trainer is executed again with the same inputs and arguments, instead of being run again. Additionally, DataDreamer caches at the model-level, caching the results of prompts or texts being run against a model to a SQLite database file. During training, DataDreamer similarly automatically saves checkpoints and resumes from them if interrupted and restarted. Caching uses minimal disk space (storing mainly text) and adds minimal overhead in these workloads dominated by heavy model inference computation, but can be granularly disabled if desired.

256

257

262

265

267

269

270

274

276

277

278

279

281

284

291

292

301

302

303

DataDreamer's cache system allows a researcher to share both their workflow script and their session output folder with others, giving them access to useful caches and saved outputs. These allow others to easily reproduce and extend the entire workflow while also benefiting from avoiding expensive computations when unnecessary. For example, a researcher could extend another researcher's workflow by adding another step at the end. Only the additional added step would need to be computed, while all of the original steps could have their results loaded from disk.

4.7 Resumability

Caching allows resumability during development, so scripts can be interrupted and resumed. This allows graceful handling of crashes, server preemption, and other situations where only a portion of a workflow was previously computed. Furthermore, caching can be useful during experimentation of a workflow. For example, when modifying a single prompt in the middle of a multi-step synthetic data generation workflow, the change may only affect a certain number of inputs to the next step. If so, only that portion of the work will be re-computed.

4.8 Sharing Open Data and Open Models

DataDreamer provides convenient utilities for exporting and publishing datasets and trained models produced by steps or trainers. Resources can be exported to disk or published to the Hugging Face Hub.⁶ When resources are published, DataDreamer can automatically upload a demonstration snippet and set up the live demonstration widget on the Hugging Face Hub, which makes shared resources easily usable. Additionally, these resources are automatically given appropriate metadata such as tags clearly indicating when data is syntheti-

Date & Time

The date and time the step or trainer was run. This is important to document when using API-based LLMs that can be updated over time.

Dataset Name & Card

The name of any datasets used as part of a step or trainer's operation along with their data cards.

Model Name & Card

The name of any models used in a step or trainer's operation along with their model cards.

URL

A URL that can be referenced for more information about the step or trainer.

License

Any known license that may apply as a result of a model or dataset being used in a step or trainer.

Citations

Citations for datasets and models used in a trainer.

Reproducibility Fingerprint

A hash of all inputs, arguments, and configurations that may affect reproducibility for a step or trainer. When steps and trainers are chained in a multi-stage workflow, the reproducibility hash is computed recursively through the chain. These fingerprints can be used to compare if two workflows within DataDreamer are exactly identical.

Other Reproducibility Information

Other miscellaneous reproducibility information such as environment information, system information, and versions of packages and dependencies.

Table 3: Information automatically recorded in a synthetic data card or synthetic model card. An example synthetic data card can found in Appendix E.

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

322

cally generated and its source LLM. DataDreamer also produces what we call "synthetic data cards" and "synthetic model cards". Synthetic data and model cards are automatically produced by recursively tracing through all steps, models, and trainers that DataDreamer used to produce the dataset or model. Each step, model, and trainer has associated metadata including license information and citation information. DataDreamer collects this information and produces a synthetic data card (or model card) that reports the information along with reproducibility information for each step, model, and trainer in the workflow. The information collected in our cards is defined in Table 3.

These automatically generated synthetic data cards and model cards can aid in preventing contamination of pre-training sources with modelgenerated synthetic data. As synthetic data generation becomes more prevalent, contamination can

⁶https://huggingface.co/

be a concern due to the performance degradation that has been observed when synthetic datasets are 324 shared and trained on, possibly without the knowl-325 edge of the model developer (Shumailov et al., 2023). DataDreamer's cards can also help other researchers understand what license restrictions may 328 apply to the synthetically generated data, among 329 other usability concerns. These automatically generated cards are not a replacement for traditional data cards and model cards (Pushkarna et al., 2022; 332 Mitchell et al., 2019) that recommend a wider set 333 of important attributes such as potential dataset bi-334 ases. Instead, they provide supplemental informa-335 tion that is crucial to the usability and reproducibility of LLM workflows. We encourage researchers 337 to review and add information that cannot be automatically detected to our generated cards.

4.9 Efficiency and Optimizations

341

342

347

LLMs workflows often benefit from or require certain optimizations to be applied in order to load or process the scale of data and models typically used. DataDreamer supports many of the common optimizations that researchers may want to apply.

Parallelization DataDreamer supports running steps in background processes and running steps concurrently to easily implement parallel task or-chestration in a workflow.

Quantization and Adapters DataDreamer supports quantization of model weights that can reduce 351 memory usage (Dettmers and Zettlemoyer) as well as parameter-efficient fine-tuning techniques like LoRA adapters (Hu et al., 2021; Mangrulkar et al., 354 2022). It standardizes using these optimizations across different model architectures and minimizes boilerplate, making it as simple as a single argu-357 ment to configure training with LoRA in Example 1. DataDreamer attempts to create uniform support for features across all of its supported integrations when possible. So while the underlying sentence_transformers and transformers 362 libraries do not support training embedding models with LoRA (Reimers and Gurevych, 2019; Wolf et al., 2019), DataDreamer supports this, which 366 extends the benefits of LoRA to these models.

Multi-GPU Usage DataDreamer makes it simple
to load models on multiple GPUs and train models on multiple GPUs with PyTorch FSDP (Paszke
et al., 2019; Zhao et al., 2023). For example, training a model on multiple GPUs is as simple as pass-

ing a list of torch. devices to the device parame-372 ter of a trainer (device=["cuda:0", "cuda:1"]). 373 DataDreamer automatically configures FSDP and 374 launches distributed processes within the session 375 so that a command line launcher like torchrun 376 never has to be used, simplifying multi-GPU train-377 ing. The use of torchrun can often force com-378 plex, multi-stage workflows being split into multi-379 ple scripts launched via shell scripts since training portions need to be isolated from data generation 381 or data processing portions. This added complex-382 ity in running the workflow end-to-end can make 383 reproducibility challenging. With DataDreamer, workflows do not need to be re-orchestrated around 385 portions needing to be launched via torchrun. Since DataDreamer handles this distributed orches-387 tration automatically, users can build multi-stage 388 workflows involving data generation, data process-389 ing, and training on multiple GPUs all in a single 390 Python program, obviating the use of orchestration 391 through multiple shell scripts. Example 4 in the 392 Appendix provides an example of such a workflow. 393

4.10 Configuration and Extensibility

DataDreamer seeks to minimize configuration and boilerplate code that for most research workflows do not need to be customized, for example automatically handling tokenization and applying the correct padding, among other tasks. DataDreamer applies sensible defaults and standard research practices to minimize configuration. Some researchers, however, may need to customize these choices and the option to override and extend is provided and well-documented. 394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

5 Reproducibility

We outline a few best practices, specific to the emerging use of LLMs in research workflows that DataDreamer adopts. We believe instituting these practices can alleviate a number of reproducibility concerns. Of course, when closed-source models are involved, these concerns can never be fully eliminated (see Section 6 for further discussion on limitations). We discuss how DataDreamer makes it easier to implement these practices or automatically implements these practices in this section.

Adaptable to Model Substitution While experimental workflows can often be sensitive to model choice and the transferability of prompts can be unreliable (Liu et al., 2023), for reproducibility purposes and for ease of experimentation, workflow



Figure 2: DataDreamer logs produced by the workflow in Example 1 when resuming from a prior interrupted run.

implementation code should attempt to minimize dependence on a specific model and should allow other researchers to easily substitute one LLM for another. This can also be useful if a model is not accessible to another researcher or if a model has become obsolete. DataDreamer's API and model abstractions make model substitution simple.

421

422

423

424

425

426

427

428 Sharing Prompts Exact prompts used should be shared since even minor variations can signif-429 icantly impact performance (Sclar et al., 2023). 430 DataDreamer makes it easy to share an entire work-431 flow and session output folder. DataDreamer can 432 433 also help ensure a re-implementation is exactly 434 identical between two experimental setups by comparing the reproducibility fingerprints of individual 435 steps or the entire workflow in aggregate. 436

Sharing Intermediate Outputs In multi-stage 437 workflows, intermediate outputs should be shared 438 for inspection and analysis by other researchers as 439 well as for extendability purposes. DataDreamer 440 makes this simple by automatically saving the re-441 sults of each step in a multi-stage workflow in 442 443 an easily inspectable Hugging Face datasets format (Lhoest et al., 2021). When API-based LLMs 444 are used, there is greater risk to reproducibility. 445 DataDreamer allows workflows to be exactly re-446 produced from caches in the session output folder, 447 448 even if the remote API is no longer available.

449 Synthetic Data Cards and Model Cards Synthetic data and model cards can help other researchers understand the source of synthetic data, license restrictions that may apply, citations that may apply, among other attributes. Importantly, these cards and other metadata-like tags can help

prevent contamination of pre-training data (Shumailov et al., 2023). Finally, these cards carry reproducibility information, useful for validating two experimental setups as identical. 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

486

487

Sharing Optimization Configurations Optimizations like quantization can have an effect on generations (Jaiswal et al., 2023). DataDreamer's reproducibility fingerprints account for these configurations and with its easily shareable workflows, DataDreamer makes it easy to reproduce an exact workflow, along with configured optimizations.

Environment-Agnostic Code For reproducibility, code should attempt to minimize dependence on local environments, job schedulers, shell scripts, etc. DataDreamer helps make this easier by providing tools for workflow orchestration (steps, parallelization, managed distributed processes for multi-GPU training) that can be all be done within Python. DataDreamer also minimizes dependencies on local file paths, by organizing results and outputs into the session output folder automatically.

6 Conclusion

The current moment in NLP research and recent progress is exciting yet raises important questions for the community. We introduce DataDreamer, an open source Python package for implementing common patterns and workflows involving LLMs. We believe DataDreamer provides both practical and scientific utility to the research community and that its adoption can help advance the rate of research progress in workflows involving LLMs by making implementation easier and making research output reproducible and extendable.

Limitations

488

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

529

532

533

535

537

538

In this work, we outline best practices and im-489 plement these practices in an open source system 490 called DataDreamer. We believe these contribu-491 tions can help aid open science in our field, how-492 ever, we acknowledge that as long as the research 493 community chooses to use closed-source models 494 for experiments, especially those served behind an 495 API on remote servers, challenges to reproducibil-496 ity are inevitable. With DataDreamer, we provide a 497 way to reproduce and further analyze some of these 498 experiments long after these remote APIs may be 499 changed or unavailable through the session-based 500 caching system as well as provide a way to easily substitute models where needed through abstractions. To the best of our knowledge, there are no 503 significant ethical considerations that arise from 504 this work. We believe the broader impacts of this 505 work to be largely positive, making state-of-the-art LLM workflows both easier and more accessible to implement and reproduce as well as reducing 508 509 carbon emissions through DataDreamer's caching system that helps researchers avoid expensive re-510 computation when possible. 511

References

- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
 - Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. 2024. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*.
- Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937*.
- Tim Dettmers and Luke Zettlemoyer. The case for 4-bit precision: k-bit inference scaling laws, 2022. URL https://arxiv. org/abs/2212.09720.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpacafarm: A simulation framework for methods

that learn from human feedback. *arXiv preprint arXiv:2305.14387*.

539

540

541

542

543

544

545

546

547

548

549

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

586

587

588

589

590

591

- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *arXiv preprint arXiv:2302.04166*.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio Cesar Teodoro Mendes, Allison Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero C. Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, S. Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuan-Fang Li. 2023. Textbooks are all you need. *ArXiv*, abs/2306.11644.
- Jesse Michael Han, Igor Babuschkin, Harrison Edwards, Arvind Neelakantan, Tao Xu, Stanislas Polu, Alex Ray, Pranav Shyam, Aditya Ramesh, Alec Radford, et al. 2021a. Unsupervised neural machine translation with generative language models only. *arXiv preprint arXiv:2110.05448*.
- Jesse Michael Han, Igor Babuschkin, Harrison Edwards, Arvind Neelakantan, Tao Xu, Stanislas Polu, Alex Ray, Pranav Shyam, Aditya Ramesh, Alec Radford, et al. 2021b. Unsupervised neural machine translation with generative language models only. *arXiv preprint arXiv:2110.05448*.
- Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. Unnatural instructions: Tuning language models with (almost) no human labor. *arXiv preprint arXiv:2212.09689*.
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. *arXiv preprint arXiv:2305.02301*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*.
- Ajay Jaiswal, Zhe Gan, Xianzhi Du, Bowen Zhang, Zhangyang Wang, and Yinfei Yang. 2023. Compressing llms: The truth is rarely pure and never simple. *arXiv preprint arXiv:2310.01382*.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020a. Data augmentation using pre-trained transformer models. *arXiv preprint arXiv:2003.02245*.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020b. Data augmentation using pre-trained transformer models. *arXiv preprint arXiv:2003.02245*.

Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, et al. 2021. Datasets: A community library for natural language processing. *arXiv preprint arXiv:2109.02846*.

593

594

596

609

610

611

612

613

614

615

616

617

618

619

621

628

631

633

641

644

645

647

- Xiang Lisa Li, Vaishnavi Shrivastava, Siyan Li, Tatsunori Hashimoto, and Percy Liang. 2023. Benchmarking and improving generator-validator consistency of language models. *arXiv preprint arXiv:2310.01846*.
- Wing Lian, Bleys Goodson, Eugene Pentland, Austin Cook, Chanvichet Vong, and "Teknium". 2023. Openorca: An open dataset of gpt augmented flan reasoning traces. https://https://huggingface. co/Open-Orca/OpenOrca.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. Advances in Neural Information Processing Systems, 35:1950–1965.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM Computing Surveys, 55(9):1–35.
- Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan.
 2022. Peft: State-of-the-art parameter-efficient finetuning methods. https://github.com/huggingfa ce/peft.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. In Annual Meeting of the Association for Computational Linguistics.
- Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In Proceedings of the conference on fairness, accountability, and transparency, pages 220–229.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of gpt-4.
- OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor

Cai, Rosie Campbell, Andrew Cann, Brittany Carey, 650 Chelsea Carlson, Rory Carmichael, Brooke Chan, 651 Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 652 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, 653 Chester Cho, Casey Chu, Hyung Won Chung, Dave 654 Cummings, Jeremiah Currier, Yunxing Dai, Cory 655 Decareaux, Thomas Degry, Noah Deutsch, Damien 656 Deville, Arka Dhar, David Dohan, Steve Dowl-657 ing, Sheila Dunning, Adrien Ecoffet, Atty Eleti, 658 Tyna Eloundou, David Farhi, Liam Fedus, Niko 659 Felix, Simón Posada Fishman, Juston Forte, Is-660 abella Fulford, Leo Gao, Elie Georges, Christian 661 Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, 662 Rapha Gontijo-Lopes, Jonathan Gordon, Morgan 663 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, 664 Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse 665 Han, Jeff Harris, Yuchen He, Mike Heaton, Jo-666 hannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny 668 Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu 669 Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger 670 Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie 671 Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, 672 Ali Kamali, Ingmar Kanitscheider, Nitish Shirish 673 Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook 674 Kim, Christina Kim, Yongjik Kim, Hendrik Kirch-675 ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, 676 Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-677 stantinidis, Kyle Kosic, Gretchen Krueger, Vishal 678 Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan 679 Leike, Jade Leung, Daniel Levy, Chak Ming Li, 680 Rachel Lim, Molly Lin, Stephanie Lin, Mateusz 681 Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, 682 Anna Makanju, Kim Malfacini, Sam Manning, Todor 683 Markov, Yaniv Markovski, Bianca Martin, Katie 684 Mayer, Andrew Mayne, Bob McGrew, Scott Mayer 685 McKinney, Christine McLeavey, Paul McMillan, 686 Jake McNeil, David Medina, Aalok Mehta, Jacob 687 Menick, Luke Metz, Andrey Mishchenko, Pamela 688 Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 689 Mossing, Tong Mu, Mira Murati, Oleg Murk, David 690 Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, 691 Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, 692 Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex 693 Paino, Joe Palermo, Ashley Pantuliano, Giambat-694 tista Parascandolo, Joel Parish, Emy Parparita, Alex 695 Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-696 man, Filipe de Avila Belbute Peres, Michael Petrov, 697 Henrique Ponde de Oliveira Pinto, Michael, Poko-698 rny, Michelle Pokrass, Vitchyr Pong, Tolly Pow-699 ell, Alethea Power, Boris Power, Elizabeth Proehl, 700 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, 701 Cameron Raymond, Francis Real, Kendra Rimbach, 702 Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-703 der, Mario Saltarelli, Ted Sanders, Shibani Santurkar, 704 Girish Sastry, Heather Schmidt, David Schnurr, John 705 Schulman, Daniel Selsam, Kyla Sheppard, Toki 706 Sherbakov, Jessica Shieh, Sarah Shoker, Pranav 707 Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, 708 Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin 709 Sokolowsky, Yang Song, Natalie Staudacher, Fe-710 lipe Petroski Such, Natalie Summers, Ilya Sutskever, 711 Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil 712 Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-713 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. Gpt-4 technical report.

714

715

716

718

724

725

727

730

731

733

734

735

739

740

741

742

743

744

745

747

749

750

751

755

761

767

770

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
 - Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
 - Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. 2022. Data cards: Purposeful and transparent dataset documentation for responsible ai. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1776–1826.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn.
 2023. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*.
 - Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
 - Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics.
 - Alexander M Rush. 2023. Minichain: A small library for coding with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 311–317.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak,

Debajyoti Datta, Jonathan D. Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng-Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Stella Biderman, Leo Gao, Tali Bers, Thomas Wolf, and Alexander M. Rush. 2021. Multitask prompted training enables zero-shot task generalization. *ArXiv*, abs/2110.08207.

771

772

774

779

780

781

782

783

787

789

790

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2023. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*.
- Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. 2023. The curse of recursion: Training on generated data makes models forget.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008– 3021.
- 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://gi thub.com/tatsu-lab/stanford_alpaca.
- Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, and Shengyi Huang. 2020. Trl: Transformer reinforcement learning. https://github.com/huggi ngface/trl.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-ofthe-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.

Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. 2022. Zerogen: Efficient zero-shot learning via dataset generation. *arXiv preprint arXiv:2202.07922*.

826

827

829

831

832

834

836

837 838

839

840

841 842

843

844

845

847

848

849

850

851

853

- Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo Lee, and Woomyeong Park. 2021. Gpt3mix: Leveraging large-scale language models for text augmentation. arXiv preprint arXiv:2104.08826.
- Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. 2023. Large language model as attributed training data generator: A tale of diversity and bias. In *Thirty-Seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models. arXiv preprint arXiv:2401.10020.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model.
- Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, et al. 2023. Pytorch fsdp: experiences on scaling fully sharded data parallel. *arXiv preprint arXiv:2304.11277*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

A Instruction-Tuning a LLM

```
1 from datadreamer import DataDreamer
2 from datadreamer.steps import HFHubDataSource
3 from datadreamer.trainers import TrainHFFineTune
4 from peft import LoraConfig
 with DataDreamer("./output"):
6
      # Get the Alpaca instruction-tuning dataset (cleaned version)
      instruction_dataset = HFHubDataSource(
    "Get Instruction-Tuning Dataset", "yahma/alpaca-cleaned", split="train"
8
9
10
      )
11
      # Keep only 1000 examples as a quick demo
12
13
      instruction_dataset = instruction_dataset.take(1000)
14
      # Some examples taken in an "input", we'll format those into the instruction
15
      instruction_dataset.map(
16
          lambda row: {
               "instruction": (
18
                   row["instruction"]
19
20
                   if len(row["input"]) == 0
                   else f"Input: {row['input']}\n\n{row['instruction']}"
               ),
               "output": row["output"],
23
24
          },
25
          lazy=False,
      )
26
27
      # Create training data splits
28
29
      splits = instruction_dataset.splits(train_size=0.90, validation_size=0.10)
30
      # Define what the prompt template should be when instruction-tuning
31
      chat_prompt_template = "### Instruction:\n{{prompt}}\n\n### Response:\n"
32
33
34
      # Instruction-tune the base TinyLlama model to make it follow instructions
35
      trainer = TrainHFFineTune(
36
           "Instruction-Tune TinyLlama",
          model_name="TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T",
38
           chat_prompt_template=chat_prompt_template,
          peft_config=LoraConfig(),
39
           device=["cuda:0", "cuda:1"],
40
          dtype="bfloat16",
41
42
      )
43
      trainer.train(
           train_input=splits["train"].output["instruction"],
44
           train_output=splits["train"].output["output"],
45
           validation_input=splits["validation"].output["instruction"],
46
           validation_output=splits["validation"].output["output"],
47
48
           epochs=3.
49
          batch_size=1,
50
           gradient_accumulation_steps=32,
51
      )
```

Example 2: In this demonstration snippet, we instruction-tune a model (Ouyang et al., 2022; Zhang et al., 2024; Taori et al., 2023). DataDreamer reduces boilerplate around tokenization, caching, training resumability, multi-GPU training, parameter-efficient fine-tuning, and more.

B Aligning a LLM

```
911
          1 from datadreamer import DataDreamer
          2 from datadreamer.steps import HFHubDataSource
912
          3 from datadreamer.trainers import TrainHFDPO
913
          4 from peft import LoraConfig
914
915
916
          6 with DataDreamer("./output"):
917
                # Get the DPO dataset
          7
                dpo_dataset = HFHubDataSource(
918
          8
919
                     "Get DPO Dataset", "Intel/orca_dpo_pairs", split="train"
          9
920
         10
                )
921
         11
922
                # Keep only 1000 examples as a quick demo
         12
923
                dpo_dataset = dpo_dataset.take(1000)
924
         14
925
         15
                # Create training data splits
                splits = dpo_dataset.splits(train_size=0.90, validation_size=0.10)
926
         16
927
928
                # Align the TinyLlama chat model with human preferences
         18
929
                trainer = TrainHFDPO(
         19
930
         20
                     "Align TinyLlama-Chat"
                     model_name="TinyLlama/TinyLlama -1.1B-Chat-v1.0",
931
         21
                    peft_config=LoraConfig(),
932
         22
                    device=["cuda:0", "cuda:1"],
dtype="bfloat16",
933
         23
934
         24
935
         25
                )
936
                trainer.train(
         26
937
                     train_prompts=splits["train"].output["question"],
         27
938
                     train_chosen=splits["train"].output["chosen"],
         28
939
         29
                     train_rejected=splits["train"].output["rejected"],
                    validation_prompts=splits["validation"].output["question"],
940
         30
941
                    validation_chosen=splits["validation"].output["chosen"],
         31
                    validation_rejected=splits["validation"].output["rejected"],
942
         32
943
         33
                    epochs=3,
944
         34
                    batch_size=1,
945
         35
                     gradient_accumulation_steps=32,
         36
                )
```

Example 3: In this demonstration snippet, we align a model using DPO (Rafailov et al., 2023; Zhang et al., 2024; Lian et al., 2023; Mukherjee et al., 2023). DataDreamer reduces boilerplate around tokenization, caching, training resumability, multi-GPU training, parameter-efficient fine-tuning, and more.

C Self-Rewarding LLMs

```
1 from datadreamer import DataDreamer
  from datadreamer.steps import (
2
      HFHubDataSource,
3
      Prompt,
4
      JudgeGenerationPairsWithPrompt,
5
6)
  from datadreamer.trainers import TrainHFDPO
7
 from datadreamer.llms import HFTransformers
8
  from peft import LoraConfig
9
10
 with DataDreamer("./output"):
11
      # Get a dataset of prompts
12
      prompts_dataset = HFHubDataSource(
13
           "Get Prompts Dataset", "Intel/orca_dpo_pairs", split="train"
14
      ).select_columns(["question"])
15
16
      # Keep only 3000 examples as a quick demo
      prompts_dataset = prompts_dataset.take(3000)
18
19
20
      # Define how many rounds of self-reward training
      rounds = 3
22
23
      # For each round of self-reward training
24
      adapter_to_apply = None
25
      for r in range(rounds):
          # Use a partial set of the prompts for each round
26
          prompts_for_round = prompts_dataset.shard(
27
               num_shards=rounds, index=r, name=f"Round #{r+1}: Get Prompts"
28
29
          )
30
           # Load the LLM
31
32
          llm = HFTransformers(
               "TinyLlama/TinyLlama-1.1B-Chat-v1.0",
33
34
               adapter_name=adapter_to_apply,
               device_map="auto",
35
36
               dtype="bfloat16",
          )
37
38
          # Sample 2 candidate responses from the LLM
39
          candidate_responses = []
40
41
          for candidate_idx in range(2):
               candidate_responses.append(
42
43
                   Prompt(
                       f"Round #{r+1}: Sample Candidate Response #{candidate_idx}",
44
45
                        inputs={"prompts": prompts_for_round.output["question"]},
                        args={
46
47
                            "llm": llm,
                            "batch_size": 2,
48
                            "top_p": 1.0,
49
                            "seed": candidate_idx,
50
51
                       },
52
                   )
               )
53
54
           # Have the LLM judge its own responses
55
           judgements = JudgeGenerationPairsWithPrompt(
56
               f"Round #{r+1}: Judge Candidate Responses",
57
58
               args={
                   "llm": llm,
59
                   "batch_size": 1,
60
                   "max_new_tokens": 5,
61
62
               }.
               inputs={
63
                    "prompts": prompts_for_round.output["question"],
64
                   "a": candidate_responses[0].output["generations"],
65
                   "b": candidate_responses[1].output["generations"],
66
67
               },
68
          )
```

```
1016
          69
1017
                      # Unload the LLM
          70
1018
                      llm.unload_model()
1019
1020
                      # Process the judgements into a preference dataset
          73
1021
          74
                      dpo_dataset = judgements.map(
                          lambda row: {
          75
                               "question": row["prompts"],
1023
          76
                               "chosen": (
1024
          77
1025
          78
                                   row["a"]
1026
          79
                                    if row["judgements"] == "Response A"
1027
                                    else row["b"]
          80
1028
          81
                                rejected": (
1029
          82
                                    row["b"]
1030
          83
                                    if row["judgements"] == "Response A"
1031
          84
1032
          85
                                    else row["a"]
          86
                               ),
1034
          87
                          }.
1035
          88
                          lazy=False,
1036
                          name=f"Round #{r+1}: Create Self-Reward Preference Dataset",
          89
1037
                      )
          90
1038
          91
                      # Create training data splits
          92
                      splits = dpo_dataset.splits(train_size=0.90, validation_size=0.10)
          93
1041
          94
1042
                      # Align the TinyLlama chat model with its own preferences
          95
1043
                      trainer = TrainHFDPO(
          96
                          f"Round #{r+1}: Self-Reward Align TinyLlama-Chat",
1044
          97
                          model_name="TinyLlama/TinyLlama-1.1B-Chat-v1.0",
1045
          98
1046
                          peft_config=LoraConfig(),
          99
                          device=["cuda:0", "cuda:1"],
dtype="bfloat16",
1047
          100
1048
          101
1049
                      )
          102
          103
                      trainer.train(
                          train_prompts=splits["train"].output["question"],
          104
1052
                           train_chosen=splits["train"].output["chosen"],
          105
                           train_rejected=splits["train"].output["rejected"],
          106
                          validation_prompts=splits["validation"].output["guestion"],
1054
          107
1055
                          validation_chosen=splits["validation"].output["chosen"],
          108
1056
                          validation_rejected=splits["validation"].output["rejected"],
          109
1057
                           epochs=3,
          110
1058
                          batch_size=1,
1059
                          gradient_accumulation_steps=32,
          112
1060
                      )
          113
          114
                      # Unload the trained model from memory
          115
1063
                      trainer.unload_model()
          116
1064
          117
1065
                      # Use the newly trained adapter for the next round of self-reward
          118
1066
                      adapter_to_apply = trainer.model_path
          119
```

Example 4: This demonstration snippet implements a simplified version of the self-rewarding LLMs (Yuan et al., 2024) procedure. This workflow involves using an LLM to judge its own generations in order to self-align and self-improve itself over a number of rounds. DataDreamer allows this complex multi-stage workflow to be implemented intuitively, without needing to split generation and training logic into separate files and without needing to involve a launcher like torchrun to perform multi-GPU training. DataDreamer also makes this complex multi-round, multi-stage workflow automatically cachable and resumable.

D Augmenting an Existing Dataset

```
1 from datadreamer import DataDreamer
                                                                                                      1068
2 from datadreamer.llms import OpenAI
                                                                                                      1069
3 from datadreamer.steps import ProcessWithPrompt, HFHubDataSource
                                                                                                      1071
4
 with DataDreamer("./output"):
                                                                                                      1072
5
      # Load GPT-4
                                                                                                      1073
6
      gpt_4 = OpenAI(model_name="gpt-4")
                                                                                                      1074
7
                                                                                                      1075
8
      # Get HotPot QA questions
                                                                                                      1076
9
      hotpot_qa_dataset = HFHubDataSource(
                                                                                                      1077
10
11
           "Get Hotpot QA Questions",
                                                                                                      1078
           "hotpot_qa",
                                                                                                      1079
          config_name="distractor",
                                                                                                      1080
          split="train"
14
      ).select_columns(["question"])
                                                                                                      1082
15
                                                                                                      1083
16
      # Keep only 1000 questions as a quick demo
                                                                                                      1084
      hotpot_qa_dataset = hotpot_qa_dataset.take(1000)
                                                                                                      1085
18
                                                                                                      1086
19
20
      # Ask GPT-4 to decompose the question
                                                                                                      1087
      questions_and_decompositions = ProcessWithPrompt(
                                                                                                      1088
           "Generate Decompositions",
22
                                                                                                      1089
23
           inputs={"inputs": hotpot_qa_dataset.output["question"]},
                                                                                                      1090
          args={
24
                                                                                                      1091
               "llm": gpt_4,
25
                                                                                                      1092
               "instruction": (
                                                                                                      1093
26
                   "Given the question which requires multiple steps to solve,"
                                                                                                      1094
27
                   " give a numbered list of intermediate questions required to"
                                                                                                      1095
28
                   " solve the question. Return only the list, nothing else."
29
                                                                                                      1096
                                                                                                      1097
               ),
30
                                                                                                      1098
31
          },
          outputs={"inputs": "questions", "generations": "decompositions"},
32
                                                                                                      1099
      ).select_columns(["questions", "decompositions"])
33
                                                                                                      1100
```

Example 5: In this demonstration snippet, we augment an existing dataset, HotpotQA (Yang et al., 2018), a multi-hop QA dataset. DataDreamer makes it easy to perform synthetic dataset augmentation with a LLM. In this example, we add intermediate questions required to solve the multi-hop question.

1 {

E Example Synthetic Data Card

```
"data_card": {
1103
            2
1104
                        "Generate Research Paper Abstracts": {
            3
                            "Date & Time": "<DATE_TIME_HERE>",
1105
            4
                            "Model Name": [
1106
            5
1107
                                 "gpt-4"
            6
1108
                            ],
            7
1109
                            "Model Card": [
            8
1110
            9
                                 "https://cdn.openai.com/papers/gpt-4-system-card.pdf"
1111
           10
                            ٦.
1112
           11
                            "License Information": [
                                 "https://openai.com/policies"
1113
1114
                            ٦.
1115
                             "Citation Information": [
           14
                  "@article{OpenAI2023GPT4TR, \ title={GPT-4 Technical Report}, a author={OpenAI}, n journal={ArXiv}, n year={2023}, n
1116
           15
1117
1118
                  \hookrightarrow volume={abs/2303.08774},\n
                  → url={https://api.semanticscholar.org/CorpusID:257532815}\n}",
1119
1120
                                 "@article{ouyang2022training,\n title={Training language models to
           16
1121
                  \hookrightarrow follow instructions with human feedback},\n author={Ouyang, Long and Wu,
                  \hookrightarrow Jeffrey and Jiang, Xu and Almeida, Diogo and Wainwright, Carroll and
1122
1123
                  \hookrightarrow Mishkin, Pamela and Zhang, Chong and Agarwal, Sandhini and Slama, Katarina
1124
                  \hookrightarrow and Ray, Alex and others},\n _journal={Advances in Neural Information
1125
                  \hookrightarrow Processing Systems},\n volume={35},\n pages={27730--27744},\n
1126
                  \hookrightarrow year={2022}\n}"
1127
                            ٦
1128
           18
                        1129
           19
1130
           20
                            "Model Name": [
1131
           21
1132
                                 "gpt-4"
1133
           23
                            ],
1134
                             "Model Card": [
           24
1135
           25
                                 "https://cdn.openai.com/papers/gpt-4-system-card.pdf"
1136
           26
                            "License Information": [
1137
           27
1138
                                 "https://openai.com/policies"
           28
1139
           29
                            "Citation Information": [
1140
           30
                                 "@article{OpenAI2023GPT4TR,\n title={GPT-4 Technical Report},\n
1141
           31
1142
                  \hookrightarrow author={OpenAI},\n journal={ArXiv},\n year={2023},\n
1143
                  \hookrightarrow volume={abs/2303.08774},\n
1144
                  → url={https://api.semanticscholar.org/CorpusID:257532815}\n}",
1145
                                 "@article{ouyang2022training,\n title={Training language models to
           32
1146
                  \hookrightarrow follow instructions with human feedback},\n author={Ouyang, Long and Wu,
1147
                  \hookrightarrow Jeffrey and Jiang, Xu and Almeida, Diogo and Wainwright, Carroll and
                  \hookrightarrow Mishkin, Pamela and Zhang, Chong and Agarwal, Sandhini and Slama, Katarina \hookrightarrow and Ray, Alex and others},\n journal={Advances in Neural Information
1148
1149
1150
                  \hookrightarrow Processing Systems},\n volume={35},\n pages={27730--27744},\n
1151
                  \hookrightarrow year={2022}\n}"
1152
           33
                            Т
1153
           34
                       }
1154
                  },
"__version__": "0.1.0",
"datetime": "<DATE_TIME_HERE>",
"Concernment",
           35
1155
           36
1156
           37
                   "type": "ProcessWithPrompt",
1157
           38
                   "name": "Generate Tweets from Abstracts",
1158
           39
                  "version": 1.0,
"fingerprint": "28b5e209bdad7d15",
1159
           40
1160
           41
                   "pickled": false,
1161
           42
                   "req_versions": {
1162
           43
                        "dill": "0.3.7"
1163
           44
                       "sqlitedict": "2.1.0",
1164
           45
1165
           46
                       "torch": "2.1.2"
                       "numpy": "1.26.3"
1166
           47
                       "transformers": "4.36.2",
1167
           48
                       "datasets": "2.16.1",
1168
           49
                       "huggingface_hub": "0.20.2",
1169
           50
```

51	"accelerate": "0.26.1",	1170
52	"peft": "0.7.1",	1171
53	"tiktoken": "0.5.2",	1172
54	"tokenizers": "0.15.0",	1173
55	"petals": "2.2.0",	1174
56	"openai": "1.9.0",	1175
57	"ctransformers": "0.2.27",	1176
58	"optimum": "1.16.2",	1177
59	"bitsandbytes": "0.42.0",	1178
60	"litellm": "1.15.3",	1179
61	"trl": "0.7.6",	1180
62	"setfit". "1.0.3",	1181
63	"together": "0.2.10",	1182
64	"google.generativeai": "0.2.1",	1183
65	"google-cloud-aiplatform": "1.35.0"	1184
66	} ,	1185
67	"interpreter": "3.11.7 (main, Dec 4 2023, 18:10:11) [Clang 15.0.0	1186
	↔ (clang-1500.1.0.2.5)]"	1187
68 }		1188

Example 6: A JSON representation of an example automatically generated synthetic data card produced by DataDreamer for Example 1. Synthetic data cards and model cards are automatically produced by recursively tracing through any steps, models, and trainers used to produce a given dataset or model. Each step, model, and trainer has associated metadata such as license information and citation information. DataDreamer collects this information and produces a synthetic data card (or model card) that reports the information along with reproducibility information like the reproducibility fingerprint.