

# AI SCIENTIST VIA SYNTHETIC TASK SCALING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

With the advent of AI agents, automatic scientific discovery has become a tenable goal. Many recent works scaffold agentic systems that can perform machine learning research, but don't offer a principled way to train such agents—and current LLMs often generate plausible-looking but ineffective ideas. To make progress on training agents that can learn from doing, we provide a novel synthetic environment generation pipeline targeting machine learning agents. Our pipeline automatically synthesizes machine learning challenges compatible with the SWE-agent Yang et al. (2024) framework, covering topic sampling, dataset proposal, and code generation. The resulting synthetic tasks are 1) grounded in real machine learning datasets, because the proposed datasets are verified against the Huggingface API and are 2) verified for higher quality with a self-debugging loop. To validate the effectiveness of our synthetic tasks, we tackle MLGym Nathani et al. (2025), a benchmark for machine learning tasks. From the synthetic tasks, we sample trajectories from a teacher model (GPT-5 Singh et al. (2024)), then use the trajectories to train a student model (Qwen3-4B and Qwen3-8B Yang et al. (2025a)). The student models trained with our synthetic tasks achieve improved performance on MLGym raising the AUP metric by 9% for Qwen3-4B and 12% for Qwen3-8B.

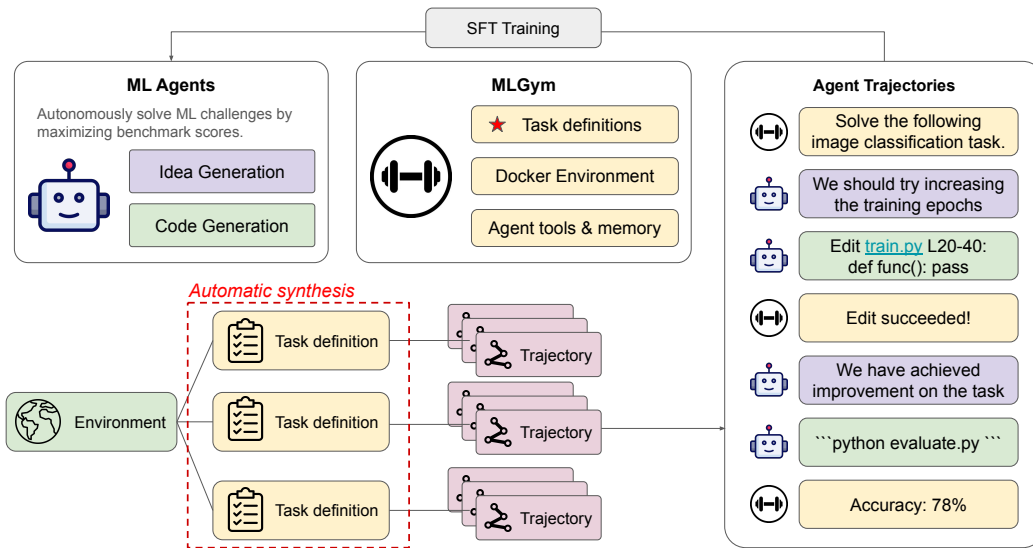


Figure 1: Overview of ML Agents and MLGym environment. Our core contribution is the starred component, synthetic tasks.

## 1 INTRODUCTION

One of the key goals of AI is to autonomously perform scientific discovery—formulating hypotheses, design and conduct experiments, analyze results, and integrate new knowledge. Recent systems such as AI Scientist Lu et al. (2024), Co-Scientist Gottweis et al. (2025), and AlphaEvolve Novikov et al. (2025) show that AI can already carry out basic research and algorithmic improvement. Meanwhile,

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

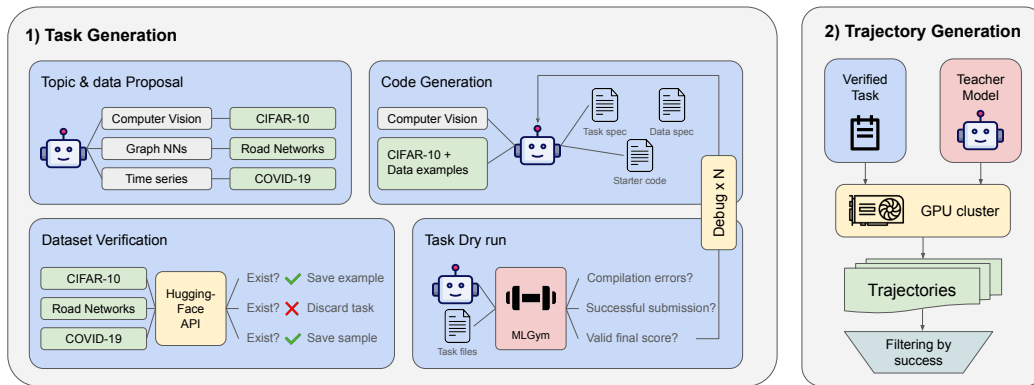


Figure 2: Illustration of our task and trajectory generation workflow. Crucially, the task generation process does not require human supervision. Instead, it automatically samples machine learning topics and proposes dataset to use in the task. To resolve compilation issues in generated tasks, we further enhance the generation with a debug loop instead of immediately discarding the task altogether.

large language models (LLMs) have acquired extensive knowledge of machine learning theory, literature, and coding patterns. Yet, knowledge alone is not enough: to convert understanding into effective research, AI agents must gain experience in executing multi-step, goal-directed tasks.

Existing research agents are often trained only on final outputs—papers, code, or datasets—ignoring the iterative processes that lead to discoveries, such as debugging, experimental failures, and step-by-step reasoning. To address this, we focus on end-to-end machine learning research task, and introduce a scalable pipeline for synthetic ML task generation that produces rich, agentic trajectories with minimal manual effort. Critically, this pipeline is compatible with the task-agnostic SWE-Agent framework, enabling models to learn from a wide variety of ML tasks across domains. By fine-tuning on these trajectories, agents gain structured experience in the full research cycle, from hypothesis to evaluation.

We use our method to tackle MLGym Nathani et al. (2025), a benchmark for machine learning agents. MLGym includes 13 machine learning tasks of various complexity. The goal of the agent is to improve upon a baseline implementation, and produce an implementation that achieves a better final score. The score is a scalar, and may vary from task to task, and usually corresponds to training accuracy, loss, win rate, etc. Based on SWE-agent framework, there is a set number of 50 rounds, and each round, the agent produce a "rational" and an "action", which may include browsing files, editing code, running commands, and submitting its final implementation. Multiple submission are allowed, which reflects iterative optimization process of the final score.

Our environment synthesis system produces around 500 tasks, which results in a dataset of around 30k agent trajectories. Training Qwen3-4B and 8B models Yang et al. (2025a) on these trajectories show performance gain, increasing performance on most individual tasks in the benchmark and increase performance of Qwen3-4B and Qwen3-8B by 9% and 12% respectively.

By combining broad knowledge, large-scale agentic experience, and task-agnostic training, our approach provides a practical path toward AI systems capable of autonomous, iterative scientific discovery.

## 2 METHODOLOGY

To advance of frontier of ML agents, we scale up automatic agent task synthesis. Since we target ML capabilities, we aim to synthesize many tasks for Machine Learning. Then, a teacher model would generate trajectories, based on synthetic tasks, which becomes viable training data for downstream models.

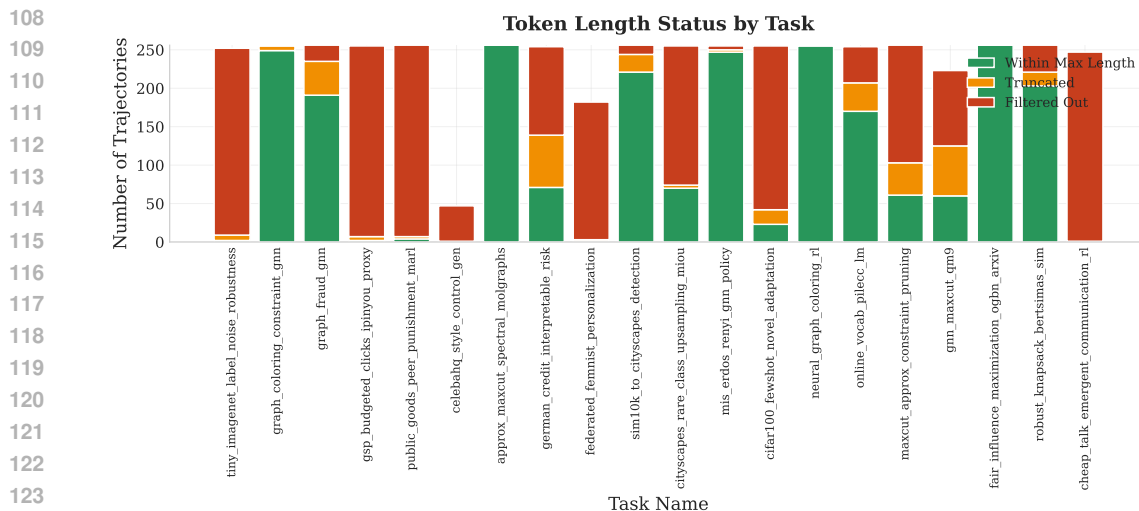


Figure 3: Generated trajectory count for each task. We select 20 generated tasks and show the number of successful trajectories for each task. Because of the unsupervised nature of our pipeline, we don’t expect all tasks to successfully create all 256 trajectories.

## 2.1 PHASE 1: ENVIRONMENT SYNTHESIS

The main driver of our method is synthetic environment generation of ML tasks. We use a multistage environment generation pipeline that focus on task diversity and task validity:

**1. Topic Sampling** Sample  $n$  distinct machine learning topics from the model.

**2. Task and dataset proposal** For each topic, the teacher model generates a task description and propose a HuggingFace dataset to use. We use the HuggingFace search API to find the closest match with the model’s proposal. We allow tasks that has no dataset (for example game theoretic tasks). If there is a match, then we enrich the dataset description with examples of the dataset rows fetched from Huggingface. If there is no match, the task is discarded.

**3. Config and starter code generation** From the task and dataset descriptions, we generate task and dataset config files compatible with the MLGym execution environment. We also generate all the starter code files for the task as well as any extra helper code. In the end, we will have baseline implementation and an evaluation file.

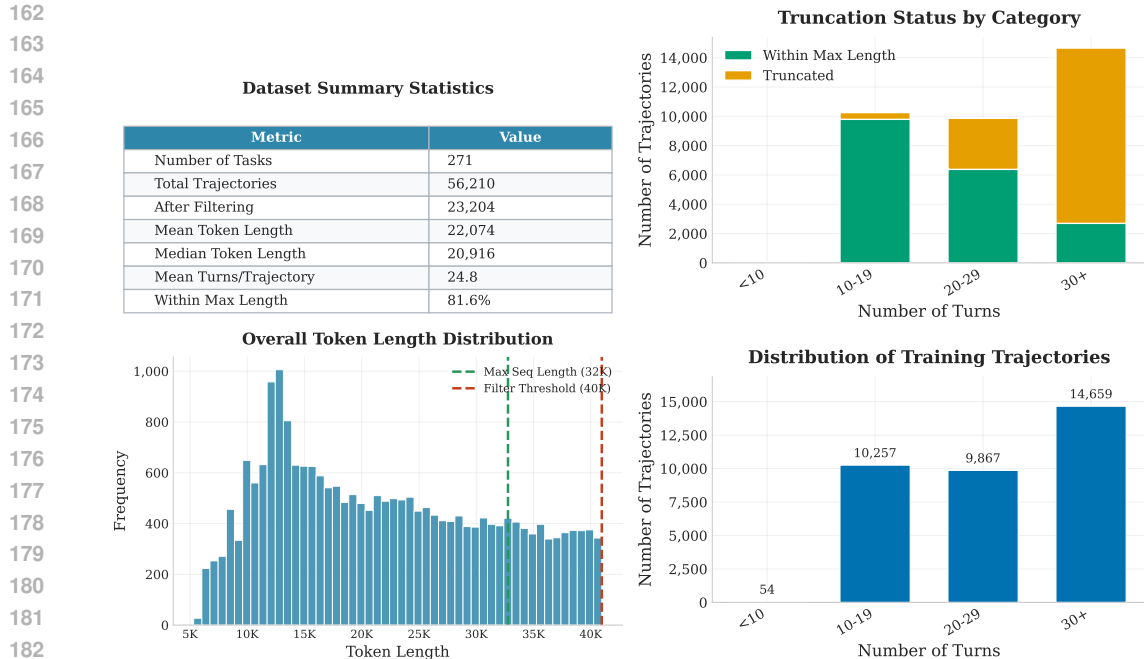
## 2.2 PHASE 2: ENVIRONMENT VERIFICATION

Since each step of the pipeline may be prone to error, we need to verify validity of the tasks as best as we can. To do this, we plug the new task into MLGym, and run the task using a GPT-5 agent to obtain the baseline performance and at least one agent trajectory. If there is an error during the execution, we collect the errors and feed them back to the model in step 3 (starter code generation) with probability  $p_{\text{debug}}$  or restart from step 3 with probability  $1 - p_{\text{debug}}$ . The iterative debug process can continue at most  $k$  times. If the task still fails after maximum iterations, we discard the task.

Crucially, this environment synthesis pipeline requires no human input, and is highly scalable through parallel compute.

## 2.3 PHASE 3: TRAJECTORY GENERATION & FILTERING

**Large scale sampling** To sample a large amount of agent trajectories for training, we run the synthetic tasks in parallel in a HPC cluster. Each task occupies one GPU, and we aim to collect 256 trajectories per tasks. Even though the tasks are validated, they can still fail in many ways. The



184 Figure 4: Top left: summary statistics of the final training trajectories. Top right: Statistics of  
185 truncated trajectories. Bottom left: distribution of tasks by token length. Bottom right: distribution of  
186 number of turns in the trajectory.

187  
188 cluster environment further impacts trajectory generation through file system and containerization  
189 instabilities. Figure 3 qualitatively show the diversity of our generated tasks.

191 **Trajectory filtering** The collected trajectories are further filtered based on agent performance.  
192 Right now, we simply choose the trajectories where the agent completes at least one successful  
193 submission. This filter is sufficient for many pathological cases where the agent is stuck in debugging  
194 loops. We also filter the trajectories based on length, rejecting any trajectories over 48K tokens long.  
195 During training, we further truncate the trajectories to 32K tokens.

### 197 3 EXPERIMENTS

199  
200 **The MLGym Benchmark** We specifically tackle the MLGym (Nathani et al., 2025) benchmark,  
201 which consists of 13 machine learning challenges of different complexity and topics, including simple  
202 game agents, computer vision, language modeling, and reinforcement learning. Each task in MLGym  
203 consists of a task description, dataset description (if task uses a dataset), and starter code. The agent  
204 lives in a standard SWE-agent environment, with tools to read and modify code, and ability to execute  
205 bash commands in a virtual environment. The agent is instructed to improve on the current solution  
206 provided in the starter code. The tasks proceeds in rounds. Each round, the agent must output some  
207 reasoning and a command The tasks have an upper limit

208 **Environment synthesis and Trajectory generation** We use GPT-5 (Singh et al., 2024) throughout  
209 our data generation pipeline. From 1000 ML topics, we generated and validated 500 tasks. For  
210 each task, we aim to generate 256 trajectories. After aggregating and filtering the trajectories we  
211 obtain around 34000 trajectories, which forms our SFT training set. Figure 3 shows a sample of the  
212 tasks generated as well as the count of valid paths generated from the tasks. Figure 4 summarize the  
213 trajectories in the final training dataset.

214  
215 **Model training** We train two models, Qwen3-4B and Qwen3-8B using SFT on the filtered trajec-  
tories. Detailed training hyperparameters are available in appendix.

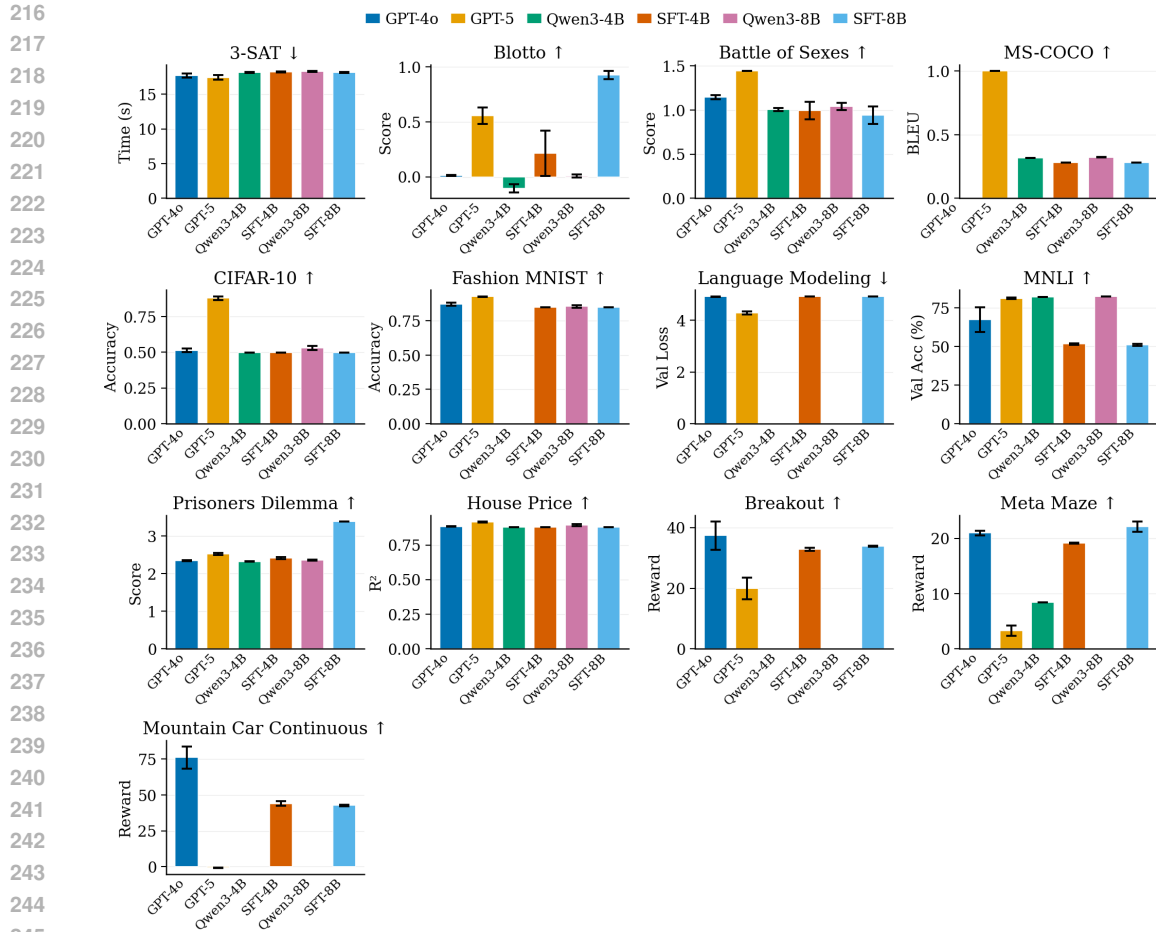


Figure 5: Model performance comparison between the baselines: GPT-4o, GPT-5, Qwen3-4B and Qwen3-8B, and our trained models: SFT-Qwen3-4B and SFT-Qwen8B. The performance is aggregated across 64 runs, which is displayed as violin plots for each subtask of MLGym. If all of the tasks fail, then the chart would show empty bar. In 9 out of 13 tasks, our trained models perform better than the baseline Qwen3-4B models.

We measure the performance of the trained models on the MLGym benchmark, and compare with GPT-4o (OpenAI et al., 2024), GPT-5 (Singh et al., 2024), Qwen3-4B and Qwen3-8B (Yang et al., 2025a). We report the performance on individual tasks and in aggregate in Figure 5 and 6.

## 4 DISCUSSION

**Failure modes** Our current task synthesis pipeline covers most but not all tasks in MLGym. For example, for the MS-COCO task, we don’t see a performance increase. This is likely because our task synthesis pipeline does not cover well the distribution of more complex starter code files. One direction is to condition the task synthesis on existing, high quality code bases (e.g. NanoGPT), so we can generate more complex tasks.

**Extending to other benchmarks** Our task synthesis pipeline is fully generic and can be easily extended to other agentic coding tasks. One good fit is MLE-Bench Chan et al. (2025), which uses Kaggle challenges. Since our models are trained on a wide variety of machine learning tasks, we expect to zero-shot performance gains on MLE-Bench.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

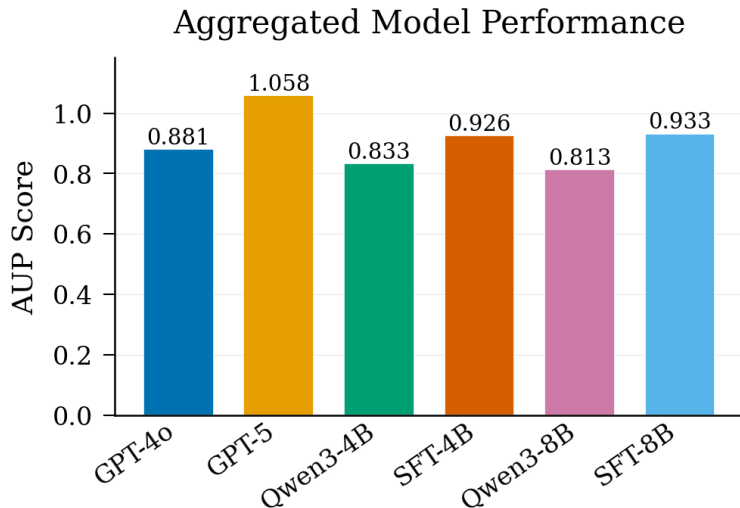


Figure 6: The aggregate performance on MLGym. Since different sub-tasks in MLGym have different score scale and comparison direction, Nathani et al. (2025) introduced the AUP score, which stands for area under the performance curve. Here we report the AUP score of each of the models.

**Optimizing for discovery of new ideas** While our synthetic task pipeline is a first step towards training LLM agents capable of machine learning tasks, we can explicitly encourage agents to form new ideas during the trajectory sampling by enabling literature search over existing machine learning research.

**Reinforcement learning** Although all of our model training is done with SFT, our synthetic tasks also can be used for reinforcement learning, where the reward signal is directly the final score defined by the task. Applying RL to machine learning tasks is challenging, because each roll-out may include long GPU training jobs, and the final reward may have vastly different scales. Addressing these challenges is a promising future direction.

## 5 RELATED WORK

Recent work has explored using LLM-based agents to support scientific research across ideation, execution, and evaluation. For ideation, multi-agent systems such as AI Co-Scientist generate and iteratively refine hypotheses aligned to researcher goals Gottweis et al. (2025). Controlled comparisons suggest LLMs can produce ideas judged more novel than expert proposals, but often with reduced feasibility Siegel et al. (2024), and downstream studies find a pronounced ideation–execution gap when researchers attempt to implement LLM-generated ideas Si et al. (2025). Other efforts structure hypothesis generation explicitly, e.g., via Bit–Flip supervision that links assumptions to counterproposals O’Neill et al. (2025).

To evaluate execution capabilities, several benchmarks test whether agents can reproduce real ML engineering and research workflows. MLE-Bench samples Kaggle-style end-to-end engineering tasks Chan et al. (2025), while PaperBench measures replication of modern ICML papers via many rubric-graded subtasks Starace et al. (2025). Related benchmarks probe targeted execution skills, such as re-implementing and improving training-script optimizations in NanoGPT “speedruns” Zhao et al. (2025). For software engineering, SWE-Smith scales task generation by synthesizing test-breaking instances across Python codebases and improves performance on SWE-bench Verified Yang et al. (2025b).

Finally, work on automated reviewing and end-to-end pipelines highlights both promise and limitations. DeepReview trains reviewer-style models with structured retrieval and argumentation Zhu et al. (2025), whereas broader evaluations show LLM reviewers remain imperfect, especially on long-context understanding and critical feedback Zhou et al. (2024). Toward full research automation, The

AI Scientist-v2 demonstrates hypothesis-to-paper loops with automated experimentation and writing Lu et al. (2024). Benchmarks such as MLAgentBench, MLGym/MLGym-Bench, and MLRC-Bench further study long-horizon research behaviors, generally finding that agents can tune and execute established pipelines but still struggle with robust planning and genuinely novel method discovery Huang et al. (2024); Nathani et al. (2025); Zhang et al. (2025); Chen et al. (2025).

## 6 CONCLUSION

We presented a scalable pipeline for training machine learning research agents via *synthetic task scaling*. Our approach automatically generates diverse ML tasks compatible with the SWE-agent framework by sampling topics, proposing and validating real HuggingFace datasets, and synthesizing full runnable environments including configs, starter code, and evaluation scripts. To ensure task validity at scale, we introduced an automated verification and self-debugging loop that filters out broken environments without requiring human intervention.

Using this pipeline, we generated roughly 500 synthetic ML tasks and collected  $\sim 30\text{k}$ – $34\text{k}$  teacher trajectories from GPT-5. Fine-tuning Qwen3-4B and Qwen3-8B on these trajectories leads to consistent gains on the MLGym benchmark, improving aggregate AUP by 9% and 12% respectively, and improving performance on the majority of individual tasks. These results suggest that synthetic environments can provide effective training signal for long-horizon agent behaviors such as iterative debugging, experimentation, and implementation refinement.

More broadly, our work supports a practical direction for building AI scientists: instead of relying purely on static corpora of papers and code, we can train agents through large-scale experience in executable research environments. We hope this enables future work on reinforcement learning over ML tasks, richer task distributions grounded in real-world codebases, and agents that move beyond optimization toward genuine discovery.

## REFERENCES

- Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Mądry. Mle-bench: Evaluating machine learning agents on machine learning engineering, 2025. URL <https://arxiv.org/abs/2410.07095>.
- Hui Chen, Miao Xiong, Yujie Lu, Wei Han, Ailin Deng, Yufei He, Jiaying Wu, Yibo Li, Yue Liu, and Bryan Hooi. Mlr-bench: Evaluating ai agents on open-ended machine learning research. *arXiv preprint arXiv:2505.19955*, 2025. 201 tasks over CS; end-to-end research pipeline; idea+writing ok, experiments often fabricated.
- Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, Khaled Saab, Dan Popovici, Jacob Blum, Fan Zhang, Katherine Chou, Avinatan Hassidim, Burak Gokturk, Amin Vahdat, Pushmeet Kohli, Yossi Matias, Andrew Carroll, Kavita Kulkarni, Nenad Tomasev, Yuan Guan, Vikram Dhillon, Eeshit Dhaval Vaishnav, Byron Lee, Tiago R D Costa, José R Penadés, Gary Peltz, Yunhan Xu, Annalisa Pawlosky, Alan Karthikesalingam, and Vivek Natarajan. Towards an ai co-scientist, 2025. URL <https://arxiv.org/abs/2502.18864>.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents on machine learning experimentation, 2024. URL <https://arxiv.org/abs/2310.03302>.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery, 2024. URL <https://arxiv.org/abs/2408.06292>.
- Deepak Nathani, Lovish Madaan, Nicholas Roberts, Nikolay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, Dieuwke Hupkes, Ricardo Silveira Cabral, Tatiana Shavrina, Jakob Foerster, Yoram Bachrach, William Yang Wang, and Roberta Raileanu. Mlgym: A new framework and benchmark for advancing ai research agents, 2025. URL <https://arxiv.org/abs/2502.14499>.

- 378 Alexander Novikov, Ngan Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wag-  
 379 ner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan  
 380 Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Push-  
 381 meet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and algorithmic discovery,  
 382 2025. URL <https://arxiv.org/abs/2506.13131>.
- 383 Charles O’Neill, Tirthankar Ghosal, Roberta Răileanu, Mike Walmsley, Thang Bui, Kevin Schawinski,  
 384 and Ioana Ciucă. Sparks of science: Hypothesis generation using structured paper data, 2025. URL  
 385 <https://arxiv.org/abs/2504.12976>.
- 386  
 387 OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan  
 388 Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-  
 389 Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex  
 390 Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau,  
 391 Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian,  
 392 Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein,  
 393 Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey  
 394 Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia,  
 395 Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben  
 396 Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake  
 397 Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon  
 398 Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo  
 399 Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li,  
 400 Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,  
 401 Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim,  
 402 Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley  
 403 Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler,  
 404 Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki,  
 405 Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay,  
 406 Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric  
 407 Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani,  
 408 Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh,  
 409 Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang  
 410 Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik  
 411 Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung,  
 412 Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu,  
 413 Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon,  
 414 Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie  
 415 Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe,  
 416 Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi  
 417 Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers,  
 418 Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan  
 419 Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh  
 420 Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn  
 421 Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra  
 422 Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe,  
 423 Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman,  
 424 Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng,  
 425 Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk,  
 426 Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine  
 427 Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin  
 428 Zhang, Marwan Aljubei, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank  
 429 Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna  
 430 Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle  
 431 Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles  
 Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho  
 Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine,  
 Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige,  
 Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko,  
 Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick

- 432 Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan,  
433 Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal,  
434 Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo  
435 Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob  
436 Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory  
437 Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi  
438 Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara  
439 Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu  
440 Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer  
441 Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal  
442 Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas  
443 Cunninghamman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao  
444 Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan,  
445 Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie  
446 Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra,  
447 Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang,  
448 Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024.  
449 URL <https://arxiv.org/abs/2410.21276>.
- 449 Chenglei Si, Tatsunori Hashimoto, and Diyi Yang. The ideation-execution gap: Execution outcomes  
450 of llm-generated versus human research ideas, 2025. URL [https://arxiv.org/abs/2506.](https://arxiv.org/abs/2506.20803)  
451 [20803](https://arxiv.org/abs/2506.20803).
- 452 Zachary S. Siegel, Sayash Kapoor, Nitya Nagdir, Benedikt Stroebel, and Arvind Narayanan. Core-  
453 bench: Fostering the credibility of published research through a computational reproducibility  
454 agent benchmark, 2024. URL <https://arxiv.org/abs/2409.11363>.
- 455 Aaditya Singh, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer  
456 Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller,  
457 Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam  
458 Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng,  
459 Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Openai gpt-5  
460 system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- 461 Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel  
462 Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, Johannes Heidecke, Amelia Glaese, Tejal  
463 Patwardhan, and OpenAI. Paperbench: Evaluating ai’s ability to replicate ai research. *arXiv*  
464 *preprint arXiv:2504.01848*, 2025. 20 ICML Spotlight/Oral papers; 8,316 sub-tasks; agent score  
465 21%.
- 466 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
467 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,  
468 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin  
469 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,  
470 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui  
471 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang  
472 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger  
473 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan  
474 Qiu. Qwen3 technical report, 2025a. URL <https://arxiv.org/abs/2505.09388>.
- 475 John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,  
476 and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering,  
477 2024. URL <https://arxiv.org/abs/2405.15793>.
- 478 John Yang, Kilian Leret, Carlos E. Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang,  
479 Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software  
480 engineering agents, 2025b. URL <https://arxiv.org/abs/2504.21798>.
- 481 Yunxiang Zhang, Muhammad Khalifa, Shitanshu Bhushan, Grant D Murphy, Lajanugen Logeswaran,  
482 Jaekyeom Kim, Moontae Lee, Honglak Lee, and Lu Wang. Mlrc-bench: Can language agents  
483 solve machine learning research challenges?, 2025. URL [https://arxiv.org/abs/2504.](https://arxiv.org/abs/2504.09702)  
484 [09702](https://arxiv.org/abs/2504.09702).

486 Bingchen Zhao, Despoina Magka, Minqi Jiang, Xian Li, Roberta Raileanu, Tatiana Shavrina, Jean-  
487 Christophe Gagnon-Audet, Kelvin Niu, Shagun Sodhani, Michael Shvartsman, Andrei Lupu,  
488 Alisia Lupidi, Edan Toledo, Karen Hambardzumyan, Martin Josifoski, Thomas Foster, Lucia  
489 Cipolina-Kun, Abhishek Charnalia, Derek Dunfield, Alexander H. Miller, Oisín Mac Aodha, Jakob  
490 Foerster, and Yoram Bachrach. The automated llm speedrunning benchmark: Reproducing nanogpt  
491 improvements, 2025. URL <https://arxiv.org/abs/2506.22419>.

492 Ruiyang Zhou, Lu Chen, and Kai Yu. Is LLM a reliable reviewer? a comprehensive evaluation  
493 of LLM on automatic paper reviewing tasks. In Nicoletta Calzolari, Min-Yen Kan, Veronique  
494 Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of the 2024 Joint  
495 International Conference on Computational Linguistics, Language Resources and Evaluation  
496 (LREC-COLING 2024)*, pp. 9340–9351, Torino, Italia, May 2024. ELRA and ICCL. URL  
497 <https://aclanthology.org/2024.lrec-main.816/>.

498 Minjun Zhu, Yixuan Weng, Linyi Yang, and Yue Zhang. Deepreview: Improving llm-based paper  
499 review with human-like deep thinking process, 2025. URL [https://arxiv.org/abs/2503.  
500 08569](https://arxiv.org/abs/2503.08569).

501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

## A APPENDIX

### A.1 PROMPTS USED IN THE TASK GENERATION PIPELINE

This appendix lists the core, non-redundant prompt texts used in the data generation pipeline.

#### A.1.1 TOPIC SAMPLING PROMPT

```
You are an expert in machine learning research. Generate a list of
20 diverse and interesting machine learning research topics.
Each topic should be a short phrase or sentence, suitable for
use as a research challenge or task. Do not repeat topics from
previous examples. Return the topics as a JSON array of strings.
```

```
Your output:
```

### A.2 TASK PROPOSAL AND DATASET VALIDATION PROMPTS

#### Task proposal prompts

```
You are an expert ML researcher create a training task for a junior
researcher. Given a topic, generate a JSON object describing a
machine learning task for that topic. The JSON must include:
- topic: the original topic
- metric: a suitable evaluation metric for the task
- description: a detailed description of the ML task
- dataset: a dataset name that can be matched to a huggingface
dataset OR a simple search query for a public huggingface
dataset (e.g., 'cifar10', 'imdb', 'tiny imagenet'). Omit this
field if the task does not require a dataset.
```

```
You have access to a tool that can search the HuggingFace datasets
API to find suitable datasets based on your query. Please make
sure the dataset exists on HuggingFace.
```

```
Here are some examples:
```

```
Topic: Image Classification
```

```
Output:
```

```
{"topic": "Image Classification", "metric": "Accuracy", "description":
  "Classify images into categories using a standard image
  classification dataset.", "dataset": "cifar10"}
```

```
Topic: Sentiment Analysis
```

```
Output:
```

```
{"topic": "Sentiment Analysis", "metric": "Accuracy", "description":
  "Predict the sentiment (positive/negative) of movie reviews.",
  "dataset": "imdb"}
```

```
Topic: Named Entity Recognition
```

```
Output:
```

```
{"topic": "Named Entity Recognition", "metric": "F1-score", "
  description": "Identify named entities in text using a standard
  NER dataset.", "dataset": "conll2003"}
```

```
Topic: Text Summarization
```

```
Output:
```

```
{"topic": "Text Summarization", "metric": "ROUGE score", "
  description": "Generate concise summaries of news articles.", "
  dataset": "cnn_dailymail"}
```

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

```

Topic: Machine Translation
Output:
{"topic": "Machine Translation", "metric": "BLEU", "description": "
  Translate sentences from English to German.", "dataset": "wmt14
"}

Topic: Speech Command Recognition
Output:
{"topic": "Speech Command Recognition", "metric": "Accuracy", "
  description": "Classify spoken commands from audio clips.", "
  dataset": "google speech commands"}

Topic: Human Activity Recognition
Output:
{"topic": "Human Activity Recognition", "metric": "Accuracy", "
  description": "Classify human activities from wearable sensor
  data.", "dataset": "UCI HAR"}

Topic: {{topic}}
Output:

```

### Dataset validation prompt

```

You may call the dataset search tool to validate or refine the
  dataset choice before producing the final JSON.
Rules:
- If you are unsure about the dataset name, call the tool with a
  short query.
- When confident, output ONLY the final JSON object (no surrounding
  prose) with required keys.
- Keys: topic, metric, description, optional dataset (string). If
  dataset provided, ensure it plausibly exists on HuggingFace.
- Avoid re-calling the tool if current results already contain a
  suitable dataset.
- You can modify the topic slightly to fit the available datasets.

```

### Dataset Search Tool-Result Follow-Up Prompt

```

Search results for query '{query}': {json.dumps(results,
  ensure_ascii=False)}
Select one dataset id (or refine by calling the tool again) and
  output final JSON when ready.

```

### JSON-Missing Nudge Prompt

```

I did not receive a valid JSON. Please either call the search tool
  or output the final JSON object now.

```

## A.3 TASK FILES GENERATION PROMPT

### Task files stage 1: config generation

```

Your objective is to create YAML config files for a machine learning
  task. You are given JSON input that describes the topic as well

```

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

```

    as the dataset you are working with. The first file is a task
    configuration file that describes the task, dataset, and
    submission format. The other files (usually one but can be
    multiple) are dataset configuration files that describe the
    datasets used in the task. Be creative and generate a task that
    is interesting and challenging for the agent to solve.

```

IMPORTANT OUTPUT FORMAT REQUIREMENT:

Return every file using markdown code blocks ONLY (no sentinel markers, no extra text) exactly like:

```

```lang
# relative/path/to/file
<file contents>
```

```

- The task is executed in a linux environment with Python 3.10 and the following packages preinstalled (generic\_conda\_requirements.txt):

generic\_conda\_requirements.txt (preinstalled packages):

```

numpy
pandas
scipy
torch
scikit-learn
tqdm
datasets
gymnasium
transformers[torch]
datasets
matplotlib
torchvision

```

Here is the format for the input JSON.

```

```json
# input.json
{
  "topic": ...,
  "metric": ...,
  "description": ...,
  "dataset": {
    "id": ...,
    "features": [...],
    "examples": [...]
  }
}
```

```

Here is the format for the config files (showing expected file outputs using markdown code blocks):

```

```yaml
# tasks/task_id.yaml
id: # task id, will be used as the config file name as well.
  Recommend using snake_case.
name: # task name
description: # Description of the task that includes the task
  objective submission format requirements. You must include the
  string "{dataset_docs}", which reference the description in the
  dataset config file.
dataset_configs: # zero or more data config files described below.
task_entrypoint: # one of four values: CSVSubmissionTasks,
  ModelSubmissionTasks, LMSubmissionTasks, PythonSubmissionTasks

```

```

702
703 training_timeout: # timeout in seconds, make a good effort
704     estimating the time it takes to train the model on our hardware
705     (NVIDIA RTX A6000 GPU, 8 CPU cores).
706 use_generic_conda: # true if the task does not require any packages
707     other than the ones listed in generic_conda_requirements.txt,
708     false otherwise
709 requirements_path: # path to requirements.txt if use_generic_conda
710     is false, otherwise leave this empty
711 starter_code: [] # Leave this as an empty list, it will be filled in
712     later
713 baseline_paths: [] # Leave this empty, it will be filled in later
714 baseline_scores: [] # Leave this as an empty list, it will be filled
715     in later
716 evaluation_paths: [] # Leave this empty, it will be filled in later
717 evaluation_read_only: # Whether the evaluation script should be read
718     -only to the agent
719 memory_path: memory.json # This value is fixed
720 ```
721
722 ```yaml
723 # datasets/dataset_name.yaml
724 data_path: # Path to the dataset, IMPORTANT: This must be a valid
725     public huggingface dataset, for example "uoft-cs/cifar10" or "
726     ILSVRC/imagenet-1k". You can find the dataset ID in the provided
727     JSON input. Do not use a placeholder.
728 description: # detailed description of the dataset, including
729     features, content, format, number of classes and samples, and
730     any other relevant information. Give concrete example rows of
731     the dataset.
732 is_local: # should always be false
733 name: # dataset name
734 ```
735
736 You can only use one of four values for `task_entrypoint` outlined
737 below.
738
739 ## Quick Decision Flow
740
741 * Agent outputs a CSV predictions file? CSVSubmissionTasks
742 * Agent submits a model/checkpoint + YAML config?
743     ModelSubmissionTasks
744 * Language-model training/eval that should run with torchrun on GPUs
745     ? LMSubmissionTasks
746 * Deliverable is Python code you evaluate directly?
747     PythonSubmissionTasks
748
749 Set this with task_entrypoint in your TaskConfig.
750
751 ---
752 1) CSVSubmissionTasks
753 Use when: The submission is a CSV of predictions (Kaggle-style).
754 Submission expected: submission.csv in the task workspace root.
755 Evaluation call (first path in evaluation_paths):
756     python <eval_script> --submission_file <path/to/submission.csv>
757 Eval output format: Entire stdout must be a valid JSON object.
758 Baseline: If baseline_paths is set, runs the first baseline script,
759     then evaluate().
760 Config snippet:
761     task_entrypoint: CSVSubmissionTasks
762
763 2) ModelSubmissionTasks
764 Use when: The agent submits a model artifact or config (e.g.,
765     checkpoints + a YAML config), not a CSV.

```

```

756 Submission expected: The first *.yaml file found under the workspace
757 .
758 Evaluation call:
759 python <eval_script> --config_fname <path/to/config.yaml>
760 Eval output format: Stdout must contain at least one line starting
761 with { that is valid JSON.
762 Baseline: Same approach.
763 Config snippet:
764 task_entrypoint: ModelSubmissionTasks
765
766 3) LMSubmissionTasks
767 Use when: Language-model tasks that should run distributed via
768 torchrun.
769 Evaluation call:
770 torchrun --nproc_per_node=<detected_gpus> --standalone <eval_script
771 >
772 Eval output format: First line starting with { that parses as JSON.
773 Baseline: Also with torchrun.
774 Config snippet:
775 task_entrypoint: LMSubmissionTasks
776
777 4) PythonSubmissionTasks
778 Use when: The agent writes Python code that evaluator imports/
779 executes directly.
780 Submission expected: target.py file.
781 Evaluation call:
782 python <eval_script>
783 Eval output format: Entire stdout JSON object.
784 Config snippet:
785 task_entrypoint: PythonSubmissionTasks
786 ---
787
788 Important tips:
789 1. Make the task and dataset description as detailed as possible:
790 task objective, dataset format, data examples, submission format
791 , metrics, constraints. Be very informative because the agent
792 rely on this information.
793 2. Include full examples of data rows in the dataset config
794 description.
795 3. The "id" field of the task config must be exactly same as the
796 filename without the .yaml extension.
797 4. In the description, always escape curly braces with double braces
798 {{ and }}, except for {dataset_docs}.
799 5. You may see a dataset field in the input JSON, use that as the
800 data_path in the dataset config. Use the dataset ID from the
801 input JSON exactly as the dataset name, it is a verified public
802 huggingface dataset.
803 6. Use the dataset information provided to you (if any), give
804 detailed information about the features, content, and example
805 rows of the dataset (if any).
806 7. If you are not prompted with a dataset, optionally choose a valid
807 public huggingface dataset. DO NOT emit a placeholder or other
808 invalid names of the dataset.
809 8. Follow the task description provided in the input.
810 9. Output files in order: one task config, then dataset config files
811 . Only one task config file.
812
813 Example description:
814 {{example_input_1}}
815
816 Example output:
817 {{example_1}}
818
819

```

```

810
811 Example description:
812   {{example_input_2}}
813
814 Example output:
815   {{example_2}}
816
817 Your description:
818   {{task_description}}
819
820 Think step by step and plan out the task before writing the output
821   files.
822 Your output:

```

### Task files stage 2: starter code generation

```

824
825 You are tasked to create a ML training task for a autonomous machine
826   learning research agent according to given config files.
827 The agent is an autonomous Machine Learning Researcher operating in
828   a specialized command-line environment.
829 In a turn-based interaction, the agent provides a "discussion" of
830   its plan, followed by a single shell command.
831 The agent can navigate the file system, and read and write files
832   using special commands, but must handle code indentation
833   manually.
834 The agent is provided with baseline code for an ML task and its goal
835   is to improve the model's performance and submit the final
836   solution.
837
838 The input task config file look like this:
839 ```yaml
840 # tasks/task_id.yaml
841 id: # task id, will be used as the config file name as well.
842   Recommend using snake_case.
843 name: # task name
844 description: # Description of the task that includes the task
845   objective submission format requirements. You must include the
846   string "{dataset_docs}", which reference the description in the
847   dataset config file.
848 dataset_configs: # zero or more data config files described below.
849 task_entrypoint: # one of four values: CSVSubmissionTasks,
850   ModelSubmissionTasks, LMSubmissionTasks, PythonSubmissionTasks
851 training_timeout: # timeout in seconds, make a good effort
852   estimating the time it takes to train the model on our hardware
853   (NVIDIA RTX A6000 GPU, 8 CPU cores).
854 use_generic_conda: # true if the task does not require any packages
855   other than the ones listed in generic_conda_requirements.txt,
856   false otherwise
857 requirements_path: # path to requirements.txt if use_generic_conda
858   is false, otherwise leave this empty
859 starter_code: [] # Fill this in after creating the task files
860 baseline_paths: [] # You need to identify the baseline file and fill
861   this in
862 baseline_scores: [] # Leave this empty, it will be filled in later
863   by running your evaluation
864 evaluation_paths: [] # You need to identify the evaluation file and
865   fill this in
866 evaluation_read_only: # Whether the evaluation script should be read
867   -only to the agent
868 memory_path: memory.json # This value is fixed
869 ```
870

```

```

864
865 The input dataset configs are optional, and look like this:
866 ```yaml
867 # datasets/dataset_name.yaml
868 data_path: # Path to the dataset, IMPORTANT: This must be a valid
869           public huggingface dataset, for example "uoft-cs/cifar10" or "
870           ILSVRC/imagenet-1k". Do not use a placeholder.
871 description: # detailed description of the dataset, including
872             content, format, number of classes and samples, and any other
873             relevant information
874 is_local: # should always be false
875 name: # dataset name
876 ```
877
878 IMPORTANT OUTPUT FORMAT REQUIREMENT:
879 Return every file using markdown code blocks ONLY exactly like:
880 ```lang
881 # relative/path/to/file
882 <file contents>
883 ```
884
885 Output ALL created files this way.
886
887 The task is executed in a linux environment with Python 3.10 and the
888 following packages preinstalled (generic_conda_requirements.txt
889 ):
890 generic_conda_requirements.txt (preinstalled packages):
891   numpy
892   pandas
893   scipy
894   torch
895   scikit-learn
896   tqdm
897   datasets
898   gymnasium
899   transformers[torch]
900   datasets
901   matplotlib
902   torchvision
903 If you decides to use any other packages, you MUST include a file
904 named requirements.txt in the output, and set the
905 use_generic_conda field to false in the task config.
906
907 Follow the following instructions:
908 - The input gives you a task config + dataset config(s). You must
909   produce runnable starter code.
910 - Leave the starter_code field empty, I will help you fill it.
911 - Provide exactly one baseline file path in baseline_paths; baseline
912   must be directly runnable without any command line arguments
913   and generate a valid submission artifact.
914 - Provide exactly one evaluation file path in evaluation_paths;
915   IMPORTANT: evaluation must output a JSON with a single field and
916   numeric score to stdout.
917 - Do not fill baseline_scores, it will be filled in later by running
918   the baseline and evaluation.
919 - You MAY create auxiliary scripts (data utils, model, etc.)
920 - Evaluation file MUST print a single valid JSON object with string
921   keys and float values (only once) for metrics.
922 - Task must run in under 30 minutes on a NVIDIA RTX A6000 GPU and 8
923   CPU cores.
924 - If you set use_generic_conda: true, then use only the preinstalled
925   packages. If you set use_generic_conda: false you add

```

```

918 requirements.txt, the path to requirements.txt MUST be set in
919 the task config.
920 - Choose dataset usage consistent with provided dataset config(s)
921 and task type. Remember you MUST use a REAL and VALID public
922 huggingface dataset. Your task may not need a dataset (e.g. game
923 theory).
924 - Keep the baseline simple, leave room for the agent to improve it.
925 - If the user runs into errors validating the task, you can change
926 the task config to fix the issue.
927
928 When you create the evaluation script, it will be ran with different
929 commands based on the value of the task_endpoint field. The
930 output must print a valid JSON object. Your evaluation file must
931 respect the evaluation call format for the task class.
932 ---
933 1) CSVSubmissionTasks
934 Use when: The submission is a CSV of predictions (Kaggle-style).
935 Submission expected: submission.csv in the task workspace root.
936 Evaluation call (first path in evaluation_paths):
937 python <eval_script> --submission_file <path/to/submission.csv>
938
939 2) ModelSubmissionTasks
940 Use when: The agent submits a model artifact or config (e.g.,
941 checkpoints + a YAML config), not a CSV.
942 Submission expected: The first *.yaml file found under the workspace
943 .
944 Evaluation call:
945 python <eval_script> --config_fname <path/to/config.yaml>
946
947 3) LMSubmissionTasks
948 Use when: Language-model tasks that should run distributed via
949 torchrun.
950 Evaluation call:
951 torchrun --nproc_per_node=<detected_gpus> --standalone <eval_script>
952 >
953 Eval output format: First line starting with { that parses as JSON.
954
955 4) PythonSubmissionTasks
956 Use when: The agent writes Python code that evaluator imports/
957 executes directly.
958 Submission expected: target.py file.
959 Evaluation call:
960 python <eval_script>
961 ---
962
963 Coding tips:
964 - Use multiprocessing / DataLoader workers for speed. You have 8 CPU
965 cores.
966 - Use GPUs for training and evaluation and anything else that makes
967 sense, assume it is always available. You have a NVIDIA RTX
968 A6000 GPU.
969 - Use deterministic seeds where relevant.
970 - Use the right indentation.
971 - Import the dataset using the datasets library load_dataset
972 function.
973 - Do not import scripts that you have not written.
974 - Do not import packages that are not in the
975 generic_conda_requirements.txt or requirements.txt.
976 - For your ease, strongly prefer a flat directory structure, and
977 create subfolders only if necessary.
978 - Use the dataset_path variable in the dataset config given to you.
979 It is verified to be a valid public huggingface dataset.

```

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990

```
- Try to use as flexible as possible package requirements, e.g. "
  torch>=2.0.0" instead of "torch==2.0.0" in requirements.txt,
  since you may have outdated knowledge.
```

```
Here is an example:
{{example_input_1}}
Example output:
{{example_output_1}}
```

```
Here is another example:
{{example_input_2}}
Example output:
{{example_output_2}}
```

```
Here is your task config:
{{task_config}}
```

```
Think step by step and plan out the task before writing the code
files.
Your output:
```

991  
992

### Error-Recovery Retry Prompt

993  
994  
995  
996  
997  
998

```
Error encountered: {self.stage_1_err}

Please try again. Return the revised output in whole.
```

999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007

### A.4 EXAMPLE SYNTHETIC TASK

We show a random example among our generated tasks. The task includes

1. Task description hotpotqa\_join\_facts\_qa.yaml
2. Dataset description hotpotqa\_hotpot\_qa.yaml
3. Starting implementation baseline.py
4. Evaluation code evaluate.py

1008  
1009

#### hotpotqa\_hotpot\_qa.yaml hotpotqa\_join\_facts\_qa.yaml

1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

```
data_path: hotpotqa/hotpot_qa
description: "HotpotQA is a large-scale multi-hop question answering
dataset featuring\
 \ questions that require reasoning across multiple documents. This
configuration\
 \ targets the distractor setting, where each example provides 10
candidate paragraphs\
 \ (titles and sentence lists), of which only a subset contains the
gold supporting\
 \ sentences needed to answer the question.\nFeatures: - id (string
): Unique identifier\
 \ for the example, e.g., \"5a7a06935542990198eaf050\". - question
(string): Natural\
 \ language question requiring multi-hop reasoning. - answer (
string): Gold answer\
 \ text (can be \"yes\"/\"no\" or a short span). - type (string):
Question type,\\"
```

```

1026
1027 \ e.g., \"comparison\", \"bridge\". - level (string): Difficulty
1028 level, e.g., \"
1029 easy\", \"medium\", \"hard\". - supporting_facts (struct of lists)
1030 : \n - title\
1031 \ (list[string]): Titles of the documents containing supporting
1032 sentences.\n \
1033 \ - sent_id (list[int32]): 0-based indices of the supporting
1034 sentences in the corresponding\
1035 \ documents.\n The k-th title aligns with the k-th sent_id to form
1036 a pair (title[k],\
1037 \ sent_id[k]).\n- context (struct):\n - title (list[string]):
1038 Titles of the 10\
1039 \ candidate documents.\n - sentences (list[list[string]]): For
1040 each document,\
1041 \ a list of its sentence strings, aligned by index with context.
1042 title.\n\nTypical\
1043 \ splits: - train: ~90k-113k examples (depending on release/
1044 version). - validation/dev:\
1045 \ ~7k-8k examples. - test: may be available without supporting
1046 facts/answers in\
1047 \ certain releases. For this task, use train and validation/dev.\
1048 nData format details:\
1049 \ - Sentence indices in supporting_facts.sent_id are 0-based and
1050 reference the sentence\
1051 \ array of the document whose title matches supporting_facts.title.
1052 - The distractor\
1053 \ setting includes 10 documents (context.title length == 10); each
1054 has a variable\
1055 \ number of sentences. - Answer normalization for EM/F1 is
1056 performed during evaluation\
1057 \ (lowercasing, removing punctuation and articles).\nConcrete
1058 examples: - Example\
1059 \ 1:\n id: \"5a7a06935542990198eaf050\"\n question: \"Which
1060 magazine was started\
1061 \ first Arthur's Magazine or First for Women?\"\n answer: \"Arthur'
1062 s Magazine\"\
1063 \n type: \"comparison\"\n level: \"medium\"\n supporting_facts:\n
1064 title: [\"
1065 Arthur's Magazine\", \"First for Women\"]\n sent_id: [0, 0]\n
1066 context:\n \
1067 \ title: [\n \"Radio City (Indian radio station)\", \"History of
1068 Albanian football\"
1069 , \"Echosmith\", \n \"Women's colleges in the Southern United
1070 States\", \"First\
1071 \ Arthur County Courthouse and Jail\", \n \"Arthur's Magazine\",
1072 \"2014-15 Ukrainian\
1073 \ Hockey Championship\", \"First for Women\", \n \"Freeway Complex
1074 Fire\", \"
1075 William Rast\"\n ]\n sentences: [\n [\n \"Radio City is India's\
1076 \ first private FM radio station and was started on 3 July
1077 2001.\", \n \"
1078 It broadcasts on 91.1 (earlier 91.0 in most cities) megahertz from
1079 Mumbai (where\
1080 \ it was started in 2004), Bengaluru (started first in 2001),
1081 Lucknow and New Delhi\
1082 \ (since 2003).\", \n \"It plays Hindi, English and regional songs
1083 .\", \n \
1084 \ \"It was launched in Hyderabad in March 2006, in Chennai on 7
1085 July 2006\
1086 \ and in Visakhapatnam October 2007.\", \n \"Radio City recently
1087 forayed into\
1088 \ New Media in May 2008 with the launch of a music portal -
1089 PlanetRadiocity.com\

```

1080 \ that offers music related news, videos, songs, and other music-  
1081 related features.\"\  
1082 ,\n \"The Radio station currently plays a mix of Hindi and  
1083 Regional music.\"\  
1084 ,\n \"Abraham Thomas is the CEO of the company.\"\  
1085 \ [\n \"Arthur's Magazine (1844-1846) was an American literary  
1086 periodical\  
1087 \ published in Philadelphia in the 19th century.\"\  
1088 T.S. Arthur,\  
1089 \ it featured work by Edgar A. Poe, J.H. Ingraham, Sarah Josepha  
1090 Hale, Thomas G.\  
1091 \ Spear, and others.\"\  
1092 \ \"In May 1846 it was merged into \\\"  
1093 Godey's Lady's\  
1094 \ Book\\\".\"\  
1095 \ ],\n ...\  
1096 \ [\n \"First for Women is a woman's\  
1097 \ magazine published by Bauer Media Group in the USA.\"\  
1098 \ \"The  
1099 \ magazine\  
1100 \ was started in 1989.\"\  
1101 \ \"It is based in Englewood Cliffs, New  
1102 \ Jersey.\"\  
1103 ,\n \"In 2011 the circulation of the magazine was 1,310,696 copies  
1104 .\"\  
1105 \ ]\n ...\  
1106 \ ]\n\n- Example 2:\n id: \"5a879ab05542996e4f30887e\"\  
1107 \n\  
1108 \ question: \"The Oberoi family is part of a hotel company that  
1109 has a head office\  
1110 \ in what city?\"\  
1111 \n answer: \"Delhi\"\  
1112 \n type: \"bridge\"\  
1113 \n level:  
1114 \"medium\"\  
1115 \n\  
1116 \ supporting\_facts:\n title: [\"Oberoi family\", \"The Oberoi  
1117 \ Group\"]\n \  
1118 \ sent\_id: [0, 0]\n context:\n title: [\n \"Ritz-Carlton Jakarta\",  
1119 \n\  
1120 \ Oberoi family\", \"Ishqbaaaz\", \"Hotel Tallcorn\", \"Mohan Singh  
1121 \ Oberoi\",\  
1122 \n \  
1123 \ \"Hotel Bond\", \"The Oberoi Group\", \"Future Fibre  
1124 \ Technologies\", \"289th\  
1125 \ Military Police Company\",\  
1126 \n \"Glennwanis Hotel\"\  
1127 \n ]\n  
1128 \ sentences:\n  
1129 \ [\n [\n \"The Ritz-Carlton Jakarta is a hotel and skyscraper in  
1130 \ Jakarta,\  
1131 \ Indonesia and 14th Tallest building in Jakarta.\"\  
1132 \n \"It is  
1133 \ located in\  
1134 \ city center of Jakarta, near Mega Kuningan, adjacent to the  
1135 \ sister JW Marriott\  
1136 \ Hotel.\"\  
1137 \n \"It is operated by The Ritz-Carlton Hotel Company  
1138 \ .\"\  
1139 \n \  
1140 \ \"The complex has two towers that comprises a hotel and the  
1141 \ Airlangga Apartment\  
1142 \ respectively.\"\  
1143 \n \"The hotel was opened in 2005.\"\  
1144 \n ]\n [\n\  
1145 \ \"The Oberoi family is an Indian family that is famous for its  
1146 \ involvement\  
1147 \ in hotels, namely through The Oberoi Group.\"\  
1148 \n ]\n ...\  
1149 \n [\n\  
1150 \ \"The Oberoi Group is a hotel company with its head office in  
1151 \ Delhi.\"\  
1152 \n \  
1153 \ ,\n \"Founded in 1934, the company owns and/or operates 30+ luxury  
1154 \ hotels\  
1155 \ and two river cruise ships in six countries, primarily under its  
1156 \ Oberoi Hotels\  
1157 \ & Resorts and Trident Hotels brands.\"\  
1158 \n ]\n ,\n ...\  
1159 \ ]\n\n-  
1160 \ Example\  
1161 \ 3:\n id: \"5a8d7341554299441c6b99fe5\"\  
1162 \n question: \"Musician and  
1163 \ satirist Allie\  
1164 \ Goertz wrote a song about the \\\"The Simpsons\\\" character  
1165 \ Milhouse, who Matt\  
1166 \

```

1134
1135 \ Groening named after who?"\n answer: \ "President Richard Nixon
1136 \ "\n type: \ "\
1137 bridge"\n level: \ "hard"\n supporting_facts:\n title: [ \ "Allie
1138 Goertz", \
1139 \ \ "Allie Goertz", \ "Allie Goertz", \ "Milhouse Van Houten" ]\n
1140 sent_id: [0, \
1141 \ 1, 2, 0]\n context:\n title: [ \n \ "Lisa Simpson", \ "Marge
1142 Simpson" \
1143 , \ "Bart Simpson", \ "Allie Goertz", \ "Milhouse Van Houten", \n
1144 \ "Los Angeles" \
1145 \ Reader", \ "Homer Simpson", \ "List of The Simpsons video games
1146 \ ", \n \ "The" \
1147 \ Simpsons: An Uncensored, Unauthorized History", \ "List of The
1148 Simpsons guest" \
1149 \ stars"\n ]\n sentences: [ \n [ \n \ "Lisa Marie Simpson is a
1150 fictional" \
1151 \ character in the animated television series \ \ "The Simpsons
1152 \ \ ". \ ", \n \
1153 \ \ "She is the middle child and most intelligent of the Simpson
1154 family. \ ", \n \
1155 \ \ "Voiced by Yeardley Smith, Lisa first appeared on television in
1156 \ \ "The Tracey" \
1157 \ Ullman Show \ \ " short \ \ "Good Night \ \ " on April 19, 1987. \ ", \n
1158 \ "Cartoonist" \
1159 \ Matt Groening created and designed her while waiting to meet
1160 James L. Brooks. \ "\
1161 , \n \ "Groening had been invited to pitch a series of shorts based
1162 on his" \
1163 \ comic \ \ "Life in Hell \ \ ", but instead decided to create a new
1164 set of characters. \ "\
1165 , \n \ "He named the elder Simpson daughter after his younger sister
1166 Lisa Groening. \ "\
1167 , \n \ "After appearing on \ \ "The Tracey Ullman Show \ \ " for three
1168 years, \
1169 \ the Simpson family were moved to their own series on Fox, which
1170 debuted on December" \
1171 \ 17, 1989. \ ", \n ], \n ..., \n [ \n \ "Allison Beth \ \ "Allie \ \
1172 \ " Goertz (born March 2, 1991) is an American musician. \ ", \n \ "
1173 Goertz is" \
1174 \ known for her satirical songs based on various pop culture
1175 topics. \ ", \n \
1176 \ \ "Her videos are posted on YouTube under the name of
1177 Cossbysweater. \ ", \n \
1178 \ \ "Subjects of her songs have included the film \ \ "The Room \ \ ",
1179 the character" \
1180 \ Milhouse from the television show \ \ "The Simpsons \ \ ", and the
1181 game Dungeons" \
1182 \ & Dragons. \ ", \n \ "Her style has been compared to that of Bo
1183 Burnham. \ "\
1184 , \n \ "In December 2015, Goertz released a concept album based on
1185 the Adult" \
1186 \ Swim series \ \ "Rick and Morty \ \ ", \ \ "Sad Dance Songs \ \ ",
1187 with the album's" \
1188 \ cover emulating the animation and logo of the series. \ ", \n \ "The
1189 album" \
1190 \ was made possible through Kickstarter. \ ", \n \ "She is co-host of
1191 Everything's" \
1192 \ Coming Up Podcast, a Simpsons-focused podcast along with Julia
1193 Prescott. \ "\n \
1194 \ ], \n [ \n \ "Milhouse Mussolini van Houten is a fictional
1195 character" \
1196 \ featured in the animated television series \ \ "The Simpsons \ \ ",
1197 voiced by Pamela" \

```

```

1188
1189 \ Hayden, and created by Matt Groening who named the character
1190 after President Richard\
1191 \ Nixon's middle name.\",\n \"Later in the series, it is revealed
1192 that Milhouse's\
1193 \ middle name is \\\"Mussolini.\\\"\"\\n ],\n ...\\n ]\\n\\nNotes: -
1194 Titles\
1195 \ must match exactly when referencing supporting facts. - sent_id
1196 values must be\
1197 \ valid indices into the corresponding document's sentence list. -
1198 The dataset is\
1199 \ English-only and licensed under CC BY-SA 4.0.\n"
1198 is_local: false
1199 name: hotpotqa/hotpot_qa
1200
1201
1202

```

### hotpotqa\_join\_facts\_qa.yaml

```

1204 id: hotpotqa_joint_facts_qa
1205 name: HotpotQA Multi-hop QA with Supporting Facts (Distractor)
1206 description: |
1207 Build a multi-hop QA system on the HotpotQA distractor setting
1208 that predicts both the final answer and the exact set of
1209 supporting sentences used to derive it. Each example provides
1210 a question, 10 candidate paragraphs (title, list of sentences),
1211 gold supporting facts as {{title, sent_id}} pairs, and an
1212 answer string {{which may be yes/no or a short span}}. You
1213 must use the provided sentence boundaries exactly {{do not re-
1214 split}} and select sentence indices from the given context.
1215 Train a multi-task model that: {{1}} predicts the answer, and
1216 {{2}} performs sentence-level classification over all
1217 candidate sentences to select supporting facts. Use either an
1218 extractive span head {{start/end indices with a separate yes/
1219 no head}} or a generative seq2seq model for the answer; use a
1220 classifier over sentence representations for supporting facts.
1221 Optimize a weighted sum of answer loss and supporting fact
1222 classification loss. Evaluate using the official HotpotQA
1223 metrics: Answer EM/F1, Supporting Facts EM/F1, and Joint EM/F1.
1224 Report Joint F1 as the primary metric, and include component
1225 metrics for analysis.
1226 Data and splits: - Use the train split for training and the
1227 validation/dev split for model selection and reporting. - This
1228 task uses the distractor setting {{10 provided paragraphs per
1229 question, where only a subset contains the supporting facts}}.
1230 - See {dataset_docs} for full dataset schema, features, and
1231 concrete examples.
1232 Modeling guidance: - Represent the 10-paragraph context with
1233 hierarchical encoders {{e.g., encode sentences, aggregate to
1234 paragraphs, then across paragraphs}}. - For long inputs,
1235 consider models that handle long sequences {{e.g., Longformer,
1236 BigBird}} or multi-hop retrieval to prune the context. -
1237 Supporting facts prediction is a multi-label classification
1238 over all candidate sentences in the provided context.
1239 Submission format: - You must submit a CSV file named submission.
1240 csv in the workspace root. - Columns:
1241 - id: the example id string {{e.g., "5a7a06935542990198eaf050"}}.
1242 - answer: the predicted answer string {{exact text; normalization
1243 is applied during evaluation}}.
1244 - supporting_facts: a JSON array of objects, each with keys "
1245 title" and "sent_id" indicating the selected supporting
1246 sentences.
1247 - Example lines:
1248 - id,answer,supporting_facts

```

```

1242
1243 - 5a7a06935542990198eaf050,Arthur's Magazine,"[{"title": "
1244   Arthur's Magazine", "sent_id": 0}], [{"title": "First
1245   for Women", "sent_id": 0}]]"
1246 - 5a879ab05542996e4f30887e,Delhi,"[{"title": "Oberoi family
1247   ", "sent_id": 0}], [{"title": "The Oberoi Group", "
1248   sent_id": 0}]]"
1249 - Constraints:
1250 - Each supporting fact must reference a title that exists in the
1251   example's context.title list and a valid sentence index {{0-
1252   based}} for that title.
1253 - No duplicates; order does not matter. Predict exactly the set
1254   of gold supporting sentences to achieve Supporting Facts EM.
1255 - Include as many sentences as required by the example {{often 2,
1256   but some questions require more}}.
1257
1258 Evaluation and metrics: - Primary metric: Joint F1 {{combines
1259   answer F1 and supporting facts F1}}. - Secondary metrics
1260   reported: Answer EM/F1, Supporting Facts EM/F1, Joint EM. -
1261   Answer normalization follows HotpotQA conventions {{
1262   lowercasing, stripping punctuation and articles}}. -
1263   Supporting facts F1 computed over the set of {{title, sent_id}}
1264   pairs.
1265 Resources and expected runtime: - Training is expected to complete
1266   in ~3-5 hours on a single NVIDIA RTX A6000 with 8 CPU cores
1267   for pre-processing. - Keep memory usage in mind when encoding
1268   long contexts. Batch size and gradient accumulation may be
1269   required.
1270 Tips: - Start with a strong encoder {{e.g., DeBERTa-v3, RoBERTa}}
1271   with segment-level inputs and a sentence classification head. -
1272   Use curriculum: begin with answer-only training, then add
1273   supporting facts loss, or alternate batches. - Joint decoding
1274   heuristic: prioritize sentences from predicted relevant
1275   paragraphs; ensure coverage of all hops before final answer
1276   extraction/generation.
1277 dataset_configs:
1278 - datasets/hotpotqa_joint_facts_qa/hotpotqa_hotpot_qa.yaml
1279 task_entrypoint: CSVSubmissionTasks
1280 training_timeout: 18000
1281 use_generic_conda: true
1282 starter_code:
1283 - data_train_v1/hotpotqa_joint_facts_qa/baseline.py
1284 - data_train_v1/hotpotqa_joint_facts_qa/evaluate.py
1285 baseline_paths:
1286 - baseline.py
1287 baseline_scores:
1288 - joint_f1: 0.022210986997935424
1289   ans_em: 0.052532072923700206
1290   ans_f1: 0.08454953694131888
1291   sp_em: 0.008102633355840648
1292   sp_f1: 0.12714489351896444
1293   joint_em: 0.0013504388926401081
1294 evaluation_paths:
1295 - evaluate.py
1296 evaluation_read_only: true
1297 memory_path: data_train_v1/hotpotqa_joint_facts_qa/memory.json

```

**baseline.py**

```

1291
1292
1293 import csv
1294 import json
1295 from datasets import load_dataset

```

```

1296
1297 from tqdm import tqdm
1298
1299 YES_NO_STARTS = (
1300     "is", "are", "was", "were",
1301     "do", "does", "did",
1302     "can", "could", "may", "might", "must",
1303     "have", "has", "had",
1304     "will", "would", "should", "shall"
1305 )
1306
1307 def simple_answer_heuristic(question, fallback_title):
1308     q = (question or "").strip().lower()
1309     if any(q.startswith(aux + " ") for aux in YES_NO_STARTS):
1310         return "yes"
1311     return fallback_title if fallback_title is not None else "unknown"
1312
1313
1314 def main():
1315     # Use the distractor setting validation split
1316     ds = load_dataset("hotpotqa/hotpot_qa", "distractor", split="
1317         validation")
1318
1319     with open("submission.csv", "w", newline="", encoding="utf-8") as
1320         f:
1321         writer = csv.DictWriter(f, fieldnames=["id", "answer", "
1322             supporting_facts"])
1323         writer.writeheader()
1324
1325         for ex in tqdm(ds, desc="Generating baseline predictions"):
1326             ex_id = ex["id"]
1327             question = ex["question"]
1328             titles = ex["context"]["title"]
1329             sentences = ex["context"]["sentences"]
1330
1331             # Choose up to two candidate supporting facts: first two
1332             titles with at least one sentence
1333             pred_sfs = []
1334             for idx, title in enumerate(titles):
1335                 if idx < len(sentences) and len(sentences[idx]) > 0:
1336                     pred_sfs.append({"title": title, "sent_id": 0})
1337                 if len(pred_sfs) == 2:
1338                     break
1339
1340             # Fallbacks if not enough
1341             if not pred_sfs:
1342                 # Ensure we still output something structurally valid
1343                 if len(titles) > 0:
1344                     pred_sfs = [{"title": titles[0], "sent_id": 0}]
1345             else:
1346                 pred_sfs = []
1347
1348             fallback_title = pred_sfs[0]["title"] if pred_sfs else (
1349                 titles[0] if titles else None)
1350             answer = simple_answer_heuristic(question, fallback_title)
1351
1352             writer.writerow({
1353                 "id": ex_id,
1354                 "answer": answer,
1355                 "supporting_facts": json.dumps(pred_sfs, ensure_ascii=
1356                     False)
1357             })
1358
1359 if __name__ == "__main__":

```

```

1350
1351     main()
1352
1353

```

### evaluate.py

```

1355
1356     import argparse
1357     import csv
1358     import json
1359     import math
1360     import re
1361     import string
1362     from collections import defaultdict
1363     from datasets import load_dataset
1364
1365     PUNCT = set(string.punctuation)
1366     ARTICLES = {"a", "an", "the"}
1367     WHITESPACE_RE = re.compile(r"\s+")
1368
1369     def normalize_answer(s):
1370         if s is None:
1371             return ""
1372         s = s.lower()
1373
1374         def remove_punc(text):
1375             return "".join(ch for ch in text if ch not in PUNCT)
1376
1377         def remove_articles(text):
1378             return re.sub(r"\b(a|an|the)\b", " ", text)
1379
1380         def white_space_fix(text):
1381             return WHITESPACE_RE.sub(" ", text).strip()
1382
1383         return white_space_fix(remove_articles(remove_punc(s)))
1384
1385     def f1_score(prediction, ground_truth):
1386         pred_tokens = normalize_answer(prediction).split()
1387         gold_tokens = normalize_answer(ground_truth).split()
1388         if len(pred_tokens) == 0 and len(gold_tokens) == 0:
1389             return 1.0
1390         if len(pred_tokens) == 0 or len(gold_tokens) == 0:
1391             return 0.0
1392         common = defaultdict(int)
1393         for t in gold_tokens:
1394             common[t] += 1
1395         num_same = 0
1396         for t in pred_tokens:
1397             if common[t] > 0:
1398                 num_same += 1
1399                 common[t] -= 1
1400         if num_same == 0:
1401             return 0.0
1402         precision = num_same / len(pred_tokens)
1403         recall = num_same / len(gold_tokens)
1404         return 2 * precision * recall / (precision + recall)
1405
1406     def exact_match_score(prediction, ground_truth):
1407         return 1.0 if normalize_answer(prediction) == normalize_answer(
1408             ground_truth) else 0.0
1409
1410     def sp_em_f1(pred_set, gold_set):
1411         # pred_set and gold_set are sets of (title, sent_id) tuples
1412         inter = pred_set.intersection(gold_set)

```

```

1404
1405     if len(gold_set) == 0 and len(pred_set) == 0:
1406         return 1.0, 1.0
1407     if len(gold_set) == 0:
1408         # No gold facts; treat as EM/F1 zero if pred non-empty;
1409         # otherwise 1
1410         return (1.0 if len(pred_set) == 0 else 0.0), (1.0 if len(
1411             pred_set) == 0 else 0.0)
1412     em = 1.0 if pred_set == gold_set else 0.0
1413     if len(pred_set) == 0:
1414         return em, 0.0
1415     precision = len(inter) / len(pred_set)
1416     recall = len(inter) / len(gold_set)
1417     if precision + recall == 0:
1418         f1 = 0.0
1419     else:
1420         f1 = 2 * precision * recall / (precision + recall)
1421     return em, f1
1422
1423 def parse_supporting_facts(cell):
1424     try:
1425         data = json.loads(cell)
1426         out = set()
1427         if isinstance(data, list):
1428             for item in data:
1429                 if isinstance(item, dict):
1430                     title = item.get("title", "")
1431                     sent_id = item.get("sent_id", 0)
1432                     try:
1433                         sent_id = int(sent_id)
1434                     except Exception:
1435                         # if cannot parse, skip this item
1436                         continue
1437                     out.add((title, sent_id))
1438         return out
1439     except Exception:
1440         return set()
1441
1442 def load_predictions_csv(path):
1443     preds = {}
1444     with open(path, "r", encoding="utf-8") as f:
1445         reader = csv.DictReader(f)
1446         for row in reader:
1447             ex_id = row.get("id", "")
1448             answer = row.get("answer", "")
1449             sf_cell = row.get("supporting_facts", "")
1450             pred_sfs = parse_supporting_facts(sf_cell)
1451             preds[ex_id] = {"answer": answer, "supporting_facts":
1452                 pred_sfs}
1453     return preds
1454
1455 def main():
1456     parser = argparse.ArgumentParser()
1457     parser.add_argument("--submission_file", type=str, required=True)
1458     args = parser.parse_args()
1459
1460     # Load dev split of HotpotQA distractor
1461     ds = load_dataset("hotpotqa/hotpot_qa", "distractor", split="
1462         validation")
1463
1464     gold_by_id = {}
1465     for ex in ds:
1466         ex_id = ex["id"]
1467         answer = ex["answer"]

```

```
1458
1459     titles = ex["supporting_facts"]["title"]
1460     sent_ids = ex["supporting_facts"]["sent_id"]
1461     gold_sfs = set()
1462     for t, s in zip(titles, sent_ids):
1463         try:
1464             s = int(s)
1465         except Exception:
1466             continue
1467         gold_sfs.add((t, s))
1468     gold_by_id[ex_id] = {"answer": answer, "supporting_facts":
1469                         gold_sfs}
1470
1471     preds = load_predictions_csv(args.submission_file)
1472
1473     total = len(gold_by_id)
1474     ans_em_sum = 0.0
1475     ans_fl_sum = 0.0
1476     sp_em_sum = 0.0
1477     sp_fl_sum = 0.0
1478     joint_em_sum = 0.0
1479     joint_fl_sum = 0.0
1480
1481     for ex_id, gold in gold_by_id.items():
1482         pred = preds.get(ex_id, {"answer": "", "supporting_facts": set
1483                                 ()})
1484
1485         a_em = exact_match_score(pred["answer"], gold["answer"])
1486         a_fl = f1_score(pred["answer"], gold["answer"])
1487         s_em, s_fl = sp_em_f1(pred["supporting_facts"], gold["
1488                               supporting_facts"])
1489
1490         ans_em_sum += a_em
1491         ans_fl_sum += a_fl
1492         sp_em_sum += s_em
1493         sp_fl_sum += s_fl
1494         joint_em_sum += (a_em * s_em)
1495         joint_fl_sum += (a_fl * s_fl)
1496
1497     metrics = {
1498         "joint_fl": joint_fl_sum / total if total > 0 else 0.0,
1499         "ans_em": ans_em_sum / total if total > 0 else 0.0,
1500         "ans_fl": ans_fl_sum / total if total > 0 else 0.0,
1501         "sp_em": sp_em_sum / total if total > 0 else 0.0,
1502         "sp_fl": sp_fl_sum / total if total > 0 else 0.0,
1503         "joint_em": joint_em_sum / total if total > 0 else 0.0,
1504     }
1505     # Print a single JSON object
1506     print(json.dumps(metrics))
1507
1508     if __name__ == "__main__":
1509         main()
1510
1511
```