Learning from [code-related] feedback

Sean Welleck

Neural Code Generation Carnegie Mellon University January 18, 2024 Part I: Foundations

- \cdot Learning
- Evaluation
- Inference
- Data

Language model learning pipeline

- Pretraining
 - Gives a "foundation model"
- Adaptation
 - Continued pretraining
 - Fine-tuning
 - $\cdot\,$ Learning from feedback
 - In-context learning / prompting

Problem: distribution mismatch

- Language model p_{θ} fits distribution q
 - $\cdot\,$ E.g., code on the web
- \cdot Language model does not learn desired distribution q^\prime
 - $\cdot\,$ E.g., code that passes tests

This can be for several reasons. For instance, not enough data, not diverse enough data, limited model capacity.

Observation 1: many signals are not explicitly in pretraining data

 whether a program compiles, whether a program passes test cases, whether a specific user prefers one program over the other, ... Observation 1: many signals are not explicitly in pretraining data

whether a program compiles,
 whether a program passes test cases,
 whether a specific user prefers one program over the other, ...

Observation 2: we can get these via feedback on generated programs

Observation 1: many signals are not explicitly in pretraining data

 whether a program compiles, whether a program passes test cases, whether a specific user prefers one program over the other, ...

Observation 2: we can get these via feedback on generated programs

Today: learning from feedback on generated programs

- Reinforcement learning
- Reward modeling
- Expert iteration

Adjust the model so that it maximizes a reward function:

$$\arg \max_{\theta} \underbrace{\mathbb{E}_{x \sim \mathcal{D}, y \sim p_{\theta}(\cdot|x)} \left[R(x, y) \right]}_{J(\theta)}$$

Example reward:

 $\cdot R(x,y) = 1$ if program y passes test cases

General pattern:

- Generate data with the model, $y \sim p_{ heta}(\cdot|x)$
- Score the data, R(y)
- Update the model using data and rewards

At a high level:

• $p_{\theta'} \leftarrow \mathcal{A}(p_{\theta}, \{x\}, R)$

Generate program $\hat{y} \sim p_{\theta}(\cdot|x)$

Estimate the gradient of the expected reward with respect to θ :

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim p_{\theta}(y|x)} \nabla_{\theta} \log p_{\theta}(y|x) R(x, y)$$
(1)

Use gradient descent to update model parameters, $\theta' \leftarrow \theta + \alpha \nabla_{\theta}$.

Example: PPO (proximal policy optimization) [7]

Various innovations to stabilize policy gradient (out of scope)

Proximal Policy Optimization Algorithms

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, Oleg Klimov OpenAI {joschu, filip, prafulla, alec, oleg}@openai.com

Abstract

We propose a new family of policy gradient methods for reinforcement learning, which alternate between sampling data through interaction with the environment, and optimizing a "surrogate" objective function using stochastic gradient ascent. Whereas standard policy gradient methods perform one gradient update per data sample, we propose a novel objective function that enables multiple epochs of minbach updates. The new methods, which we call proximal policy optimization (PPO), have some of the benefits of trust region policy optimization (TRPO), but they are much simpler to implement, more general, and have better sample complexity (empirically). Our experiments test PPO on a collection of benchmark tasks, including simulated robotic locomotion and Atari game playing, and we show that PPO outperforms other online policy gradient methods, and overall strikes a favorable balance between sample

At the end, we get an alternative algorithm:

 $p_{\theta'} \leftarrow \mathcal{A}_{PPO}(p_{\theta}, \{x\}, R)$

RL: used to update a model using rewards and generated sequences.

- $p_{\theta'} \leftarrow \mathcal{A}(p_{\theta}, \{x\}, R)$
- Policy gradient, PPO, ...

How do we choose the reward?

Issue 1: reward hacking

- \cdot Models can overfit to patterns in the reward
- Example:
 - \cdot R(x,y) = 1 if program y compiles, 0 otherwise

Then generating *y* = *print("hello world")* for all *x* would maximize reward.

Mitigation: KL-divergence penalty

- \cdot Keep the updated model close to the pretrained model
- $R_{\text{KL}} = -\beta \log \frac{p_{\theta}(y|x)}{p_0(y|x)}$

Mitigation: KL-divergence penalty

- Keep the updated model close to the pretrained model
- $R_{\text{KL}} = -\beta \log \frac{p_{\theta}(y|x)}{p_0(y|x)}$

$$egin{aligned} &D_{ ext{KL}}(p_ heta(y|x)\|p_0(y|x)) &= \sum_y p_ heta(y|x)\lograc{p_ heta(y|x)}{p_0(y|x)} \ &= \mathbb{E}_{y\sim p_ heta}\lograc{p_ heta(y|x)}{p_0(y|x)} \ &pprox \lograc{p_ heta(y|x)}{p_0(\hat{y}|x)}, \end{aligned}$$

where $\hat{y} \sim p_{\theta}(\cdot|x)$, i.e. a single-sample Monte-Carlo approximation.

Issue 2: sparse reward

• The reward may be 0 for many programs; we only occasionally see a positive reward

Mitigation: engineer the reward function

Example: PPOCoder Execution-based Code Generation using Deep RL [8]

Execution

• $R_{execution}(x, \hat{y})$: 1 if program \hat{y} compiles and passes tests cases

Example: PPOCoder Execution-based Code Generation using Deep RL [8]

Execution

- $R_{execution}(x, \hat{y})$: 1 if program \hat{y} compiles and passes tests cases Syntactic matching score
 - $R_{syntax}(x, \hat{y}, y_*)$: overlap between abstract syntax tree of y and y_*

Execution

- $R_{execution}(x, \hat{y})$: 1 if program \hat{y} compiles and passes tests cases Syntactic matching score
 - + $\mathsf{R}_{\mathsf{syntax}}(x, \hat{y}, y_*)$: overlap between abstract syntax tree of y and y_*

Semantic matching score

• $R_{\text{semantics}}(x, \hat{y}, y_*)$: overlap between dataflow graph of y and y_*

$$R = R_{\text{execution}} + R_{\text{syntax}} + R_{\text{semantics}} + R_{\text{KL}}$$

Execution

- $R_{execution}(x, \hat{y})$: 1 if program \hat{y} compiles and passes tests cases Syntactic matching score
 - + $\mathsf{R}_{\mathsf{syntax}}(x, \hat{y}, y_*)$: overlap between abstract syntax tree of y and y_*

Semantic matching score

• $R_{\text{semantics}}(x, \hat{y}, y_*)$: overlap between dataflow graph of y and y_*

$$R = R_{\text{execution}} + R_{\text{syntax}} + R_{\text{semantics}} + R_{\text{KL}}$$

Run PPO using the reward

1	. 0		
Model	$\uparrow xMatch$	$\uparrow Edit Sim$	↑Comp Rate
BiLSTM	20.74	55.32	36.34
Transformer	38.91	61.47	40.22
GPT-2	40.13	63.02	43.26
CodeGPT	41.98	64.47	46.84
CodeT5 (220M)	42.61	68.54	52.14
PPOCoder + CodeT5 (220M)	42.63	69.22	97.68

Table 1: Results on the code completion task for completing the last 25 masked tokens from CodeSearchNet.

Figure 1: Compilation rate increases while holding other metrics constant

		↑pass@1			$\uparrow pass@5$				
Model	Size	Intro	Inter	Comp	All	Intro	Inter	Comp	All
Codex	12B	4.14	0.14	0.02	0.92	9.65	0.51	O.09	2.25
AlphaCode	1B	-	-	-	-	-	-	-	-
GPT-3	175B	0.20	0.03	0.00	0.06	-	-	-	-
GPT-2	0.1B	1.00	0.33	0.00	0.40	2.70	0.73	0.00	1.02
GPT-2	1.5B	1.30	0.70	0.00	0.68	3.60	1.03	0.00	1.34
GPT-Neo	2.7B	3.90	0.57	0.00	1.12	5.50	0.80	0.00	1.58
CodeT5	60M	1.40	0.67	0.00	0.68	2.60	0.87	0.10	1.06
CodeT5	220M	2.50	0.73	0.00	0.94	3.30	1.10	0.10	1.34
CodeT5	770M	3.60	0.90	0.20	1.30	4.30	1.37	0.20	1.72
CodeRL+CodeT5	770M	4.90	1.06	0.5	1.71	8.60	2.64	1.0	3.51
PPOCoder +CodeT5	770M	5.20	1.00	0.5	1.74	9.10	2.50	1.20	3.56

Model	Size	State	$\uparrow pass@80$
GPT	224M	fine-tuned	7.2
GPT	422M	fine-tuned	12.6
GPT	1B	fine-tuned	22.4
GPT	4B	fine-tuned	33.0
GPT	8B	fine-tuned	40.6
GPT	68B	fine-tuned	53.6
GPT	137B	fine-tuned	61.4
CodeT5	60M	fine-tuned	19.2
CodeT5	220M	fine-tuned	24.0
CodeT5	770M	fine-tuned	32.4
CodeRL+CodeT5	770M	zero-shot	63.0
PPOCoder +CodeT5	770M	zero-shot	68.2

Figure 2: APPS

Figure 3: MBPP (transfer)

PPOCoder ablations







Figure 4: Reward terms (MBPP)

Figure 5: RL objective (MBPP)

Figure 6: KL penalty

Qualitative - what is it changing?

Source (NL)

You are given two integers n and m. Calculate the number of pairs of arrays (a, b) such that: the length of both arrays is equal to m; each element of each array is an integer between 1 and n (inclusive); ai \leq bi for any index i from 1 to m; array a is sorted in non-descending order; array b is sorted in non ascending order. As the result can be very large, you should print it modulo 109 + 7.

Input: The only line contains two integers n and m ($1 \le n \le 1000$, $1 \le m \le 10$).

Output: Print one integer – the number of arrays a and b satisfying the conditions described above modulo 109 + 7

Generated (Python)				
CodeT5	PPOCoder+CodeT5			
MOD = 10**9 + 7 (1) def solve(n, m): (1) for in range(2, n + 1): for in range(1, n + 1): for in range(1, n + 1): ans = (n = 0 fb)(0)(0)(0)(0)(0)(0)(0)(0)(0)(0)(0)(0)(0)	MOD = 10**9 + 7 def solve(n, m); def [0]: (n = 1) (n = 1) (n = nnnge(m = 1)) for in nnge(1, n = 1); for in nnge(1, n = 1); for in nnge(1, n = 1); do [[[] = (dpi][[]] + dp[= 1][k]) % MOD ans = 0 for in nnge(1, n = 1); ans = (ans + dp[m][]) % MOD return ans n,m = map(nt,input(),split()) print(solve(n,m))			
Pass Tests?	Pass Tests? 🗹			

Qualitative - what is it changing?



Figure 8: Error types (APPS), from RLTF: Reinforcement Learning from Unit Test Feedback [5]

Qualitative - what is it changing?



Figure 8: Error types (APPS), from RLTF: Reinforcement Learning from Unit Test Feedback [5]

programs that pass, especially for problems with introductory difficulty levels. The observed increase in failure rate stems from the fixing of error codes, resulting in either pass or fail outcomes. This illustrate that RLTF is more effective in addressing runtime and compiler errors compared to semantic errors, which remain more challenging. Figure 2b illustrates the percentages of different sub-errors among erroneous results before and after applying the RLTF method. Most errors show a decline in proportion after using the RLTF method, particularly syntax errors. It is also noteworthy that the proportion of timeout errors exhibits a minor increase, which can be attributed to the correction of other grammar-related errors resulting in timeout errors.

Figure 9: RLTF: analogous RL method with similar performance

RL with policy gradient methods

- Directly optimizes reward
- Susceptible to reward hacking; requires good reward design
- Learning procedure adds complexity
- So far, improvements may be explained by syntax/index fixes

RL with policy gradient methods

- Directly optimizes reward
- Susceptible to reward hacking; requires good reward design
- Learning procedure adds complexity
- So far, improvements may be explained by syntax/index fixes

- Reinforcement learning
- \cdot Reward modeling
- Expert iteration

Basic idea:

- Train a model $R_{\phi}(y)$ to predict whether a program is correct
 - $R_{\phi}(y) \in [0, 1]$, 0 means incorrect, 1 means correct
- At test time:
 - Generate many programs, $\{y_1, \ldots, y_K\} \sim p_{\theta}(\cdot|x)$
 - Select the program with the highest score $R_{\phi}(y)$

Basic idea:

- Train a model $R_{\phi}(y)$ to predict whether a program is correct
 - $R_{\phi}(y) \in [0, 1]$, 0 means incorrect, 1 means correct
- At test time:
 - Generate many programs, $\{y_1, \ldots, y_K\} \sim p_{\theta}(\cdot | x)$
 - Select the program with the highest score $R_{\phi}(y)$

 $R_{\phi}(y)$: "reward model" or "learned verifier"

Test time procedure: "best-of-n"

LLMs: investigated on math word problems [2]





60 Test Solve Rate (%)

Learning to Verify Language-to-Code Generation with Execution [6]

- Key difference: we can **execute** code
- Train a model $p_{\phi}(v|x, y, \mathcal{E}(y))$
 - v is 0 or 1
 - x: input prompt
 - y: generated program
 - $\mathcal{E}(y)$ is the result of executing program y

Given $(x, \mathcal{E}(y_*))$

- Generate $\{y_1, \ldots, y_K\} \sim p_{\theta}(\cdot|x)$
- Add $(x, y_k, \mathcal{E}(y_k), v_k)$ to a set S_x
 - v_k is 1 if execution result matches gold result $\mathcal{E}(y_*)$, 0 otherwise

$$\mathcal{L}(x, S_x) = -\frac{1}{|S_x|} \sum_{k=1}^{|S_x|} \log p(v_k | x, y_k, \mathcal{E}(y_k))$$

GSM8K: question + idiomatic program + answer variable

Input:

Carly recently graduated and is looking for work in a field she studied for. She sent 200 job applications to companies in her state, and twice that number to companies in other states. Calculate the total number of job applications she has sent so far. | n.job.apps.in.state = 200 n.job.apps.out.of.state = n.job.apps.in.state * 2 answer = n.job.apps.in.state + n.job.apps.out.of.state | 'answer': 600 Output: ves SPIDER/WIKITQ: question + SQL + linearized result table

```
Input:
```

-- question: List the name, born state and age of the heads of departments ordered by age \mid

-- SQL:|select name, born.state, age from head join management on head.head.id = management.head.id order by age|

-- exec result: |/*| name born_state age| Dudley Hart California 52.0| Jeff Maggert

Delaware 53.0|Franklin Langham Connecticut 67.0| Billy Mayfair California 69.0|

```
K. J. Choi Alabama 69.0 | */
```

Output: no

MBPP: task description + function + return type & value

Input:

 \sharp^{\dagger} description Write a function to find the n-th power of individual elements in a list using lambda function.

program def nth_nums(nums,n):

result_list = list(map(lambda x: x ** n, nums))
return (result_list)

execution
return: (list)=[1, 4, 9, 16, 25, 36, 49, 64, 81, 100]
return: (list)=[1000, 8000, 27000]

```
# return: (list)=[248832, 759375]
```

Output: yes

At test time:

- Generate $\{y_1, \ldots, y_K\} \sim p_\theta(\cdot|x)$
- Select the program y_k with the highest score $R(x, y_k)$.

$$\cdot r(x, y_k) = \underbrace{p_{\theta}(y_k|x)}_{\text{LM score}} \cdot \underbrace{p_{\phi}(v = 1|x, y_k, \mathcal{E}(y_k))}_{\text{verifier score}} \\ \cdot R(x, y_k) = \sum_{y_{k'} \text{ with same exec result as } y_k} r(x, y_{k'})$$



Figure 1: The illustration of LEVER using text-to-SQL as an example. It consists of three steps: 1) *Generation*: sample programs from code LLMs based on the task input and few-shot exemplars; 2) *Execution*: obtain the execution results with program executors; 3) *Verification*: using a learned verifier to output the probability of the program being correct based on the NL, program and execution results.



Figure 2: Comparison of LEVER * and baselines with Codex-Davinci. LEVER and its ablation results are in solid bars.

Figure 10: LEVER improves performance. Using execution info is important



Figure 11: Scaling the number of samples

Reward modeling:

- · Does not require updating generator $p_{ heta}$
- Simple learning objective: standard maximum likelihood
- Strong performance
- Bounded by the generator's capabilities
- Expensive at generation time
- Reward model is imperfect

- Reinforcement learning
- Reward modeling
- Expert iteration

Alternate between **search** and **learning**:

- Search: Use an 'expert model' to find good outputs
- Learning: Fine-tune on the discovered outputs
- Repeat

Thinking Fast and Slow with Deep Learning and Tree Search

Thomas Anthony^{1, 120}, Zheng Tian¹, and David Barber^{1,2}

¹University College London ²Alan Turing Institute ⊠thomas.anthony.14@ucl.ac.uk

Abstract

Sequential decision making problems, such as structured prediction, robotic control, and game playing, require a combination of planning policies and generalisation of those plans. In this paper, we present Expert Iteration (EX17), a novel reinforcement learning algorithm which decomposes the problem into separate planning and generalisation tasks. Planning new policies is performed by tree search, while a deep neural network generalises those plans. Subsequently, tree search is improved by using the neural network policy to guide search, increasing the strength of new plans. In contrast, standard deep Reinforcement Learning algorithms rely on a neural network not only to generalise plans, but to discover them too. We show that EXI7 outperforms RELINFORCE for training a neural network to play the board game Hex, and our final tree search agent, trained tabula rasa, defeats MOHEX 10, the most recent Olympiad Champion player to be publicly released.

Figure 12: Anthony et al 2017

For neural code generation:

- Search: Generate many programs, save those that succeed
- Learning: Fine-tune on the saved programs
- Repeat

For neural code generation:

- Search: Generate many programs, save those that succeed
- Learning: Fine-tune on the saved programs
- Repeat

"Self-training": the expert model is the current language model (plus the binary execution feedback)

Google DeepMind

2023-12-25

Beyond Human Data: Scaling Self-Training for Problem-Solving with Language Models

Avi Singh^{1,*}, John D Co-Reyes^{1,*}, Rishabh Agarwal^{1,2,*},

Ankesh Anand¹, Piyush Patil¹, Xavier Garcia¹, Peter J. Liu¹, James Harrison¹, Jaehoon Lee¹, Kelvin Xu¹,

Aaron Parisi¹, Abhishek Kumar¹, Alex Alemi¹, Alex Rizkowsky¹, Azade Nova¹, Ben Adlam¹, Bernd Bohnet¹, Gamaleldin Elsayed¹, Hanie Sedghi¹, Igor Mordatch¹, Isabelle Simpson¹, Izzeddin Gur¹, Jasper Snoek¹, Jeffrey Pennington¹, Jiri Hron¹, Kathleen Kenealy¹, Kevin Swersky¹, Kshiteej Mahajan¹, Laura Culp¹, Lechao Xiao¹, Maxwell L Bileschi¹, Noah Constant¹, Roman Novak¹, Rosanne Liu¹, Tris Warkentin¹, Yundi Qian¹, Yamini Bansal¹, Ethan Dyer¹, Behnam Neyshabur¹, Jascha Sohl-Dickstein¹, Noah Fiedel¹ ^c Contributed equally. ¹Google DeepMind, ² Mila

Builds on recent ideas, e.g. for reasoning [13, 12], generation [4], preference alignment [3].

Algorithm 1: ReST (Expectation-Maximization). Given a initial policy (e.g., pre-trained LM), ReST^{EM} iteratively applies Generate and Improve steps to update the policy.

 $\begin{array}{l} \textbf{Input: } \mathcal{D}: \text{Training dataset, } \mathcal{D}_{val}: \text{Validation dataset, } \mathcal{L}(\textbf{x}, \textbf{y}; \theta): \text{loss, } r(\textbf{x}, \textbf{y}): \text{Non-negative} \\ \text{reward function, } I: \text{ number of iterations, } N: \text{ number of samples per context} \\ \textbf{for } i = 1 \ to \ I \ \textbf{do} \\ \hline // \ \text{Generate } (\textbf{E-step}) \\ \text{Generate dataset } \mathcal{D}_i \ \text{by sampling: } \mathcal{D}_i = \{ (\textbf{x}^j, \textbf{y}^j) |_{j=1}^N \ \text{s.t. } \textbf{x}^j \sim \mathcal{D}, \ \textbf{y}^j \sim p_\theta(\textbf{y} | \textbf{x}^j) \} \\ \text{Annotate } \mathcal{D}_i \ \text{with the reward } r(\textbf{x}, \textbf{y}). \\ // \ \text{Improve } (M-\text{step}) \\ \text{while reward improves on } \mathcal{D}_{val} \ \textbf{do} \\ & \mid \text{Optimise } \theta \ \text{to maximize objective: } J(\theta) = \mathbb{E}_{(\textbf{x}, \textbf{y}) \sim \mathcal{D}_i} \left[r(\textbf{x}, \textbf{y}) \ \log p_\theta(\textbf{y} | \textbf{x}) \right] \\ \text{end} \\ \text{end} \\ \\ \text{Output: Policy } p_\theta \end{array}$

Self-training with execution



Figure 13: On the MATH dataset, improves for multiple iterations

Self-training with execution



Figure 14: On a subset of APPS: initially improves, then overfits.

Self-training with execution



Figure 15: On a subset of APPS: improves pass@k

$$\mathcal{L}_{\mathsf{RL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \mathcal{P}_{\theta}(y|x)} \left[R(x, y) \right]$$

$$\mathcal{L}_{\mathsf{RL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \rho_{\theta}(y|x)} \left[R(x, y) \right]$$

Policy gradient methods: interleave updates and generation

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left[\nabla_{\theta} \log p_{\theta}(\hat{y}|x) R(x, \hat{y}) \right]$$

$$\mathcal{L}_{\mathsf{RL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \rho_{\theta}(y|x)} \left[R(x, y) \right]$$

Policy gradient methods: interleave updates and generation

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left[\nabla_{\theta} \log p_{\theta}(\hat{y}|x) R(x, \hat{y}) \right]$$

Self-training: generate a large dataset, then update

$$\theta_{t+1} \leftarrow \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}} \left[\mathbb{E}_{y \sim p_{\theta_t}(y|x)} \left[r(x, y) \log p_{\theta}(y|x) \right] \right]$$

$$\mathcal{L}_{\mathsf{RL}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \rho_{\theta}(y|x)} \left[R(x, y) \right]$$

Policy gradient methods: interleave updates and generation

$$\theta_{t+1} \leftarrow \theta_t + \alpha \left[\nabla_{\theta} \log p_{\theta}(\hat{y}|x) R(x, \hat{y}) \right]$$

Self-training: generate a large dataset, then update

$$\theta_{t+1} \leftarrow \arg \max_{\theta} \mathbb{E}_{x \sim \mathcal{D}} \left[\mathbb{E}_{y \sim p_{\theta_t}(y|x)} \left[r(x, y) \log p_{\theta}(y|x) \right] \right]$$

See the Rest-EM paper [9] for more details on the connection.

Self-training:

- Natural extension of best-of-n, which had good performance
- Simple learning objective: standard maximum likelihood
- Susceptible to overfitting
- Very recent; ongoing investigation!

Three methods for learning from feedback:

- Directly optimize a reward with reinforcement learning
- Learn a reward, generate programs, select the best program
- Generate programs, save successful ones, train on them

Looking ahead:

- $\cdot\,$ Each method has pros and cons
- Still a research frontier for code generation
- $\cdot\,$ Other potential sources of feedback, e.g. natural language $[1]^1$

¹Another research frontier; not covered due to time constraints.

📔 A. Chen.

Improving code generation by training with natural language feedback.

ArXiv, abs/2303.16749, 2023.

 K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, C. Hesse, and J. Schulman.

Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168, 2021.*

References ii

 H. Dong, W. Xiong, D. Goyal, Y. Zhang, W. Chow, R. Pan, S. Diao, J. Zhang, K. SHUM, and T. Zhang.
 RAFT: Reward ranked finetuning for generative foundation model alignment. *Transactions on Machine Learning Research*, 2023.
 C. Gulcehre, T. L. Paine, S. Srinivasan, K. Konyushkova, L. Weerts, A. Sharma, A. Siddhant, A. Ahern, M. Wang, C. Gu, W. Macherey, A. Doucet, O. Firat, and N. de Freitas.

Reinforced self-training (rest) for language modeling, 2023.

J. Liu, Y. Zhu, K. Xiao, Q. FU, X. Han, Y. Wei, and D. Ye. **RLTF: Reinforcement learning from unit test feedback.** *Transactions on Machine Learning Research*, 2023.

References iii

- - A. Ni, S. Iyer, D. Radev, V. Stoyanov, W.-t. Yih, S. I. Wang, and X. V. Lin.

Lever: Learning to verify language-to-code generation with execution.

In Proceedings of the 40th International Conference on Machine Learning (ICML'23), 2023.

- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. **Proximal policy optimization algorithms, 2017.**
- P. Shojaee, A. Jain, S. Tipirneni, and C. K. Reddy.
 Execution-based code generation using deep reinforcement learning.

Transactions on Machine Learning Research, 2023.

References iv

A. Singh, J. D. Co-Reyes, R. Agarwal, A. Anand, P. Patil, P. J. Liu,
 J. Harrison, J. Lee, K. Xu, A. Parisi, A. Kumar, A. Alemi, A. Rizkowsky,
 A. Nova, B. Adlam, B. Bohnet, H. Sedghi, I. Mordatch, I. Simpson,
 I. Gur, J. Snoek, J. Pennington, J. Hron, K. Kenealy, K. Swersky,
 K. Mahajan, L. Culp, L. Xiao, M. L. Bileschi, N. Constant, R. Novak,
 R. Liu, T. B. Warkentin, Y. Qian, E. Dyer, B. Neyshabur, J. N.
 Sohl-Dickstein, and N. Fiedel.

Beyond human data: Scaling self-training for problem-solving with language models.

ArXiv, abs/2312.06585, 2023.

References v

R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation.

In Proceedings of the 12th International Conference on Neural Information Processing Systems, NIPS'99, page 1057–1063, Cambridge, MA, USA, 1999. MIT Press.

🔋 R. J. Williams.

Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8:229–256, 1992.

Z. Yuan, H. Yuan, C. Li, G. Dong, C. Tan, and C. Zhou. Scaling relationship on learning mathematical reasoning with large language models. ArXiv, abs/2308.01825, 2023. E. Zelikman, Y. Wu, J. Mu, and N. Goodman.
 STar: Bootstrapping reasoning with reasoning.
 In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, Advances in Neural Information Processing Systems, 2022.

 D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. Christiano, and G. Irving.
 Fine-tuning language models from human preferences. ArXiv, abs/1909.08593, 2019.