
Measuring few-shot extrapolation with program induction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural networks are capable of learning complex functions, but still have problems
2 generalizing from few examples and beyond their training distribution. Meta-
3 learning provides a paradigm to train networks to learn from few examples, but it
4 has been shown that some of its most popular benchmarks do not require significant
5 adaptation to each task nor learning representations that extrapolate beyond the
6 training distribution. Program induction lies at the opposite end of the spectrum:
7 programs are capable of extrapolating from very few examples, but we still do
8 not know how to efficiently search these discrete spaces. We propose a common
9 benchmark for both communities, by learning to extrapolate from few examples
10 coming from the execution of small programs. These are obtained by leveraging
11 a C++ interpreter on codes from programming competitions and extracting small
12 sub-codes with their corresponding input-output pairs. Statistical analysis and
13 preliminary human experiments show the potential of this benchmark for enabling
14 progress in few-shot extrapolation.

15 1 Introduction and motivation

16 Despite their great successes learning complex functions, neural networks still require large amounts
17 of data and have trouble generalizing beyond their training distribution. In contrast, program
18 induction lies at the opposite end of both spectrums: we can infer programs from few examples that
19 extrapolate far beyond the training samples, but these programs are typically very simple. As a field,
20 we're trying to find the best of both worlds: learning complex functions that generalize broadly from
21 few examples.

22 There have been multiple few-shot learning benchmarks, such as Omniglot [1] and mini-Imagenet [2],
23 that have enabled fast progress in meta-learning. However, it has also been shown that getting
24 high performance on some of these benchmarks is simpler than expected, as one can simply
25 learn a common set of features for all tasks and adapt the last layer of a neural network [3]. The
26 authors of Omniglot have also reported [4] that most solutions getting high results on their few-shot
27 classification benchmark do not learn generalizable, flexible representations. We introduce a
28 new few-shot learning dataset which, we hope, can provide a new challenging benchmark for
29 the meta-learning community; asking approaches to extrapolate to inputs outside of their training
30 distribution in natural ways from natural examples.

31 Program induction has seen a lot of growing interest, but there is also a lack of large-scale few-shot
32 program induction benchmarks. Most datasets are manually created by a small team of researchers;
33 this creates biases, often assumes a specific Domain Specific Language (DSL) and limits the number
34 of total tasks to a few hundreds. There is also, to the best of our knowledge, no benchmark that allows
35 testing a broad spectrum of program complexity, from 1-line programs to 100-line implementations

36 requiring algorithmic insights. In this work, we present a benchmark containing such a spectrum of
37 problems, with the hope that it provides multiple a useful target both now and in the coming years.

38 Humans are able to code complex programs that satisfy a series of input/output pairs and generalize
39 well beyond these examples: being applicable to larger numbers, longer lists, or tests requiring more
40 computation time. In particular, there are popular *competitive programming* competitions where
41 coders have to implement programs to novel tasks; similar to coding interviews in some technology
42 companies. However, current systems are still very far from being able to solve generic complex
43 programs such as those required in these competitions. Therefore, using these problems directly
44 would not be a useful challenge.

45 We propose to create a wide range of tasks, ranging from simple one-line programs to complex
46 algorithmic implementations, by extracting sub-codes from codes online. More concretely, we
47 leverage a database coming from the popular website `codeforces.com`, which has hundreds of
48 problems and hundreds of thousands of implementations. From these, we can extract tens of thousands
49 different sub-codes, each describing its own task. We can obtain input-output examples for each code
50 by running the entire program on a set of inputs and recording the input-output example for each
51 sub-code. C++ is, by far, the most popular language in these competitions. C++ has the advantage of
52 being typed and structured, which, as we will see, allows us to more easily categorize and analyze the
53 difficulty of problems. At the same time, C++ is also compiled, which makes it quite hard to collect
54 the evaluation of sub-codes, since that requires running C++ as an interpreted language. Despite this
55 difficulty, we manage to obtain such evaluations for a broad subset of the codes, which allows us to
56 create a rich, diverse meta-dataset with over 10,000 tasks.

57 Our main contribution is therefore a few-shot benchmark for the meta-learning and program induction
58 communities with the following key aspects:

- 59 **1. The dataset is automated, allowing us to scale to more than ten thousand real-world**
60 **tasks** and removing some of the biases of manually created datasets. At the same time,
61 tasks are not random, each coming from a solution to a task similar to those given in tech
62 interviews.
- 63 **2. Problems have an accompanying solution in the form of a code**, which can also function
64 as a target for methods that are learning to search in program induction. Solutions also
65 provide auxiliary information to guide meta-learning and provide an extra layer of structure,
66 which is not typically part of traditional few-shot learning benchmarks.
- 67 **3. Tasks have a diverse range of difficulty and can be classified into many different**
68 **groups**: from their type traces, to whether the program that generated them had loops,
69 or its number of lines.

70 2 Related work

71 **Meta-learning** meta-learning [5, 6, 7] aims at learning priors from many tasks so as to generalize
72 to a new task from few amounts of data; for a nice recent survey see [8]. The three main paradigms
73 in meta-learning have been optimization-based approaches, MAML being the primary example [9],
74 model-based approaches that tailor to a particular application (often image classification) [10] and
75 architecture-based approaches that use LSTMs, transformers, GNNs, etc to encode the dataset before
76 making a prediction [11, 12, 13]. Most of these methods assume that the input form is uniform
77 and do not typically generalize outside broadly outside the data distribution [14], especially non-
78 optimization-based approaches [15]. Given that we know tasks are generated by pieces of code, this
79 makes it close to AutoML [16], which often searches through code-like representations to optimize
80 machine learning models.

81 **Related datasets** Few-shot learning benchmarks have allowed great progress in meta-learning.
82 The two most popular ones are in few-shot classification for computer vision: miniImageNet [2]
83 and Omniglot [1]. Other notable few-shot classification benchmarks have been proposed such as
84 tieredImageNet [17], SlimageNet [18], CUB-200 [19] and meta-dataset [20]. There have also been
85 pushes to increase the generality of meta-reinforcement learning benchmarks to include completely
86 different virtual environments [21, 22] as well as learning an entire RL loss functions that generalize
87 between them [23].

88 There have been a number of program induction benchmarks, such as those used in DreamCoder [24],
 89 and FlashFill [25, 26], as well as the Abstract Reasoning Challenge(ARC) [27], a list functions
 90 benchmark [28], or the SyGus competition [29]. Although these benchmarks contain many interesting
 91 problems, they have been manually created by humans instead of being automatically generated from
 92 natural data. This creates significant bias on the datasets (often being captured by a relatively simple
 93 Domain Specific Language) and restricts the amount of tasks to a few hundreds to a bit more than a
 94 thousand tasks. In contrast, our benchmark contains 11,000 tasks and we are planning to extend it to
 95 around 100k. This will allow neural-based methods, often data-inefficient, to learn to generalize or
 96 learn to search in these domains with less need to embed biases into the search.

97 **Datasets of competitive programming** There have also been some works leveraging data from
 98 competitive programming websites. Most notably, we take the programs scraped by [30] and
 99 Dr.Repair [31] from `codeforces.com`. However, their goals are significantly different from ours.
 100 In particular, they use the code itself in its entirety. Since doing program induction on the entire
 101 competitive programming tasks is far beyond current methods, they also make the task easier along
 102 interesting dimensions. Whereas we make the tasks easier by considering sub-codes and creating
 103 thousands of few-shot learning tasks, they learn to go from line-by-line pseudo-code to code [30]
 104 or learn to debug programs [31]. Even though the origin of the data is the same, our end-product is
 105 orthogonal to theirs. [32] is probably closest to our benchmark, manually annotating a bit more than
 106 2000 codes with a problem statement. There are two problems with this approach: first, since it is
 107 much easier to describe what the code does than creating a task that the code solves, most statements
 108 resemble pseudo-code, which is turns the goal into something closer to translation. Second, because
 109 of codes are manually annotated it suffers from the same problem as other program induction
 110 benchmarks, since it is hard to scale to tens of thousands of tasks, as in our case.

111 3 Description of the ANONYMOUS dataset

112 We refer to our dataset as ANONYMOUS because the original name contains the name of our
 113 institution.

114 3.1 Description of the data

115 In competitive programming there are different *problems*, each with a statement (a short text) requiring
 116 a program that solves the specified description. There are also multiple test-cases (some public, some
 117 private) that the submitted program has to satisfy. For each *problem* we have many hundreds of
 118 *codes* that solve it; meaning we have a pair of (code, test-cases) such that the entire code satisfies the
 119 test-cases.

120 For each code we can obtain *sub-codes*: valid continuous segments of code contained in the original
 121 code. To be valid, a sub-code has to start and end at the same level of indentation and never go to a
 122 level above than the starting one in the indentation hierarchy (i.e. it’s indentation has to be a correct
 123 parenthetisation).

124 Given a sub-code we can generate the data for a task; consisting of 20 input-output pairs (10 training,
 125 10 test). We obtain these by running the entire code and observing the intermediate values at every
 126 line. Therefore, a task consists of:

- 127 • a type trace describing the types of inputs and outputs,
- 128 • 20 pairs of input-outputs examples, 10 for training and 10 for test,
- 129 • a code that solves these pairs and extrapolates to other inputs.

130 Note that the input distribution is far from random, as it is affected by previous computations in the
 131 overall program. Figure 1 shows the example of two tasks.

132 3.2 Overview of the implementation

133 In this section we provide an overview of how we obtained the data. This helps provide a better
 134 understanding on the data distribution, explaining how we obtain the input-output pairs as well as
 135 some limitations of our pipeline, which constrain some of the problems in our dataset.

```

int(v0) subproblem_1080A_63051446_10_10(int v0, int v1) { std::string &(v2) subproblem_514A_60400646_13_13(int v0, int v1, std::string & v2) {
    v0 += v1 - 1;
    return v0;
}
    v2[v1] = min(v2[v1], (char)(v0 + '0'));
    return v2;
}

Input 0: 6 5
Output 0: 10

Input 1: 30 6
Output 1: 35

Input 2: 2 100000000
Output 2: 100000001

Input 3: 193730132 63740710
Output 3: 257470841

Input 4: 116129238 65614207
Output 4: 181743444

Input 0: 1 0 "8772"
Output 0: "1772"

Input 1: 2 0 "71723447"
Output 1: "21723447"

Input 2: 1 0 "8464062628894325"
Output 2: "1464062628894325"

Input 3: 4 0 "5728"
Output 3: "4728"

Input 4: 1 0 "8537"
Output 4: "1537"

```

Figure 1: Examples of two (relatively simple) tasks in our (meta-)dataset. Each task consists of 10 training input-output pairs, 10 test input-output pairs and a program that solves them and extrapolates to other inputs.

We obtain the original raw codes from SPoC and DrRepair [30, 31], which scraped `codeforces.com`. In their case, they are interested in analyzing the code itself and executing the entire program (which can be done by compiling it as usual). This gives us around 300,000 codes to 700+ competitive programming tasks.

To interpret C++ we use the Cling C++ interpreter [33]. Cling performs an elaborate incremental just-in-time compilation that keeps modifying the abstract syntax tree before executing the new piece of code. This allows us to execute pieces of code and check the values of variables in between. Since they have to be compiled, the given pieces of code have to be self-contained: functions have to be defined entirely before being fed to Cling and loops and if statements have to be given as a block. This severely restricts the type of sub-codes that we can obtain with raw Cling, since we cannot inspect the intermediate values within loops or functions.

In the current version of the dataset (we plan to expand it with even more tasks) we discard codes that contain functions and implement a work-around to be able to obtain intermediate values for loops and if statements. In particular we first standardize all codes changing for loops to while loops plus extra instructions and ensure all loops and if statements are properly bracketed. Once this processing is done, we create an emulator that, instead of feeding the entire while/if statement to Cling, it first calls its condition and then calls the appropriate code depending on whether the condition is satisfied. Note that these if/while conditionals are often very interesting quantities, and we also include them as tasks even though there is no explicit boolean variable created in the original code.

Competitive programming codes interact with the terminal, receiving inputs and outputting the result. Cling cannot handle console input or output; we therefore implement a wrapper that simulates this communication. Similar to loop conditionals, console outputs often contain interesting results and we thus also store them as program outputs.

Finally, in C++ we can initialize a variable without giving it a value ("int a;"). It is also possible that some variable is initialized in an if statement that is executed for some test-cases, but not others. In both circumstances, whenever the value has never been set, we list it as 'null'.

3.3 Current limitations and future improvements

There are a few practical limitation with the implemented pipeline that restrict some codes from being added to our database. This does not affect the correctness of our tasks, but slightly biases the distribution of codes in our benchmark with respect to the distribution of codes in competitive programming as a whole.

As mentioned above, Cling cannot analyze functions line by line. Therefore we cannot obtain sub-codes from pieces of functions. We are currently discarding codes that have functions; however, in the future we will add codes that contain functions, treating them as individual instructions that

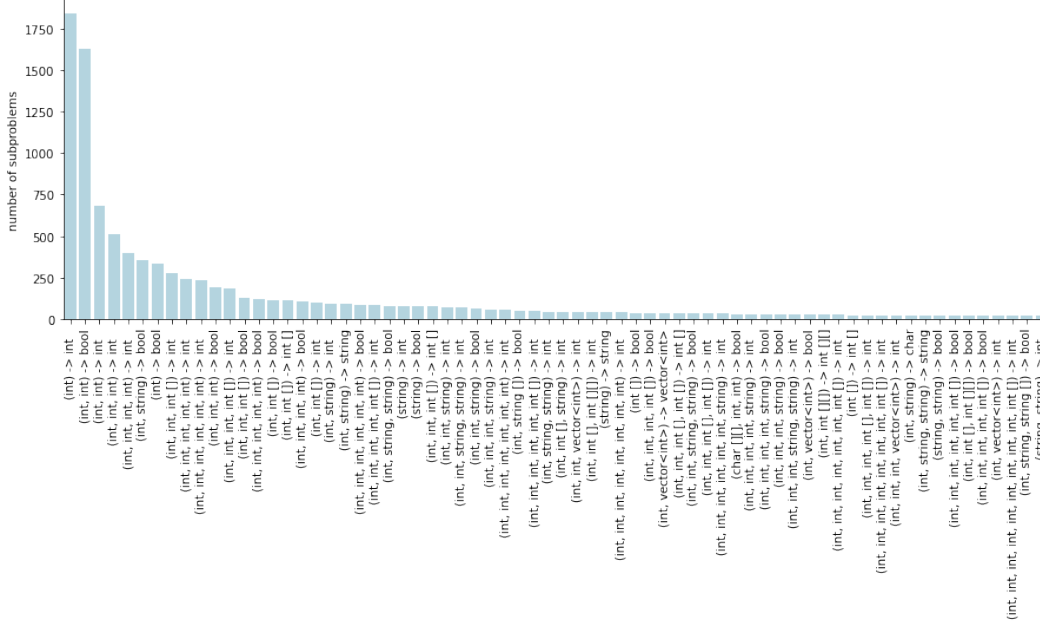


Figure 2: Overview of input-output type traces. Most problems involve few inputs and output either integers or booleans. Some perform string or array manipulations.

cannot be split. To restrict the size of the overall dataset as well of the intermediate pipeline, we currently remove test-cases that surpass 10^6 bits=125KB.

Programs in our tasks consist of contiguous segments of code where the output variable is modified on the last line and input variables are all variables involved in this particular segment. However, this implies that current programs contain lines or variables that do not necessarily affect the input. Understanding these relations requires static analysis and we plan to do it in the near future.

Finally, we often have codes that are implemented differently, but end up producing the same results. Detecting these occurrences is hard to do for arbitrary programs, and often expensive, but we only need to do it once during the creation of the dataset. Moreover, it can also be approximated by checking whether two programs solve the test-cases of one another. We have currently standardized each program by making variable names depend on their order of appearance instead of their original name. Going forward, we plan on removing further symmetries (such as swapping a pair of lines whose order does not affect the output) by expressing programs as graphs.

3.4 Statistical analysis

We analyze the difficulty and diversity of our meta-dataset along multiple axis. First, we can understand the type of problems through their traces. As expected, the most popular traces involve integer manipulations as well as classification problems from few variables. It would be interesting to assess whether classification techniques from few-shot image classification can be re-purposed for tasks with boolean outputs.

There are other traces that involve array(list) and string manipulations, often conditioned on other variables like integers or individual characters. These are interesting as they often require to generalize to longer computations as well as bigger data structures. Finally, there are other types that are more complex such as matrices or dictionaries mapping a collection of keys with their respective values.

We can also measure the difficulty of our tasks along three different axis (see figure 3). First we can see that 70% of tasks have a fixed computation graph that does not involve generalizing to longer computations. However 30% have either if-statements or loops that require generalizing to more 2x to 20x more computation than that observed for training examples. Conditional execution (characterized by indentation in C++) is often very hard for program induction techniques. We observe that most tasks can be solved with programs with a single level of indentation (no conditional execution), but

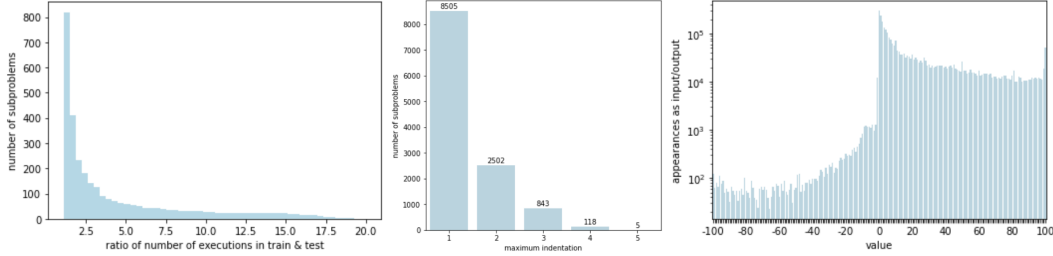


Figure 3: (Left) Ratio of number of maximum number of lines executed for a test input vs maximum number of lines executed for a training input; showing the requirement to generalize to longer executions. 71% of tasks have a fixed execution graph and thus did not require longer computation for test inputs. (Center) Depth of indentation execution (involving nested loops and if statements), which often significantly affects the difficulty of program induction. (Right) number of times each integer in $[-100,100]$ appears as an input or output on a test-case; note the logarithmic y-axis.

some require multiple up to 4 levels of nested execution. Finally, we observe that most input and outputs involve small positive integers (note the logarithmic y axis), but many involve larger numbers. It is worth noting that these can extend up to $\pm 2 \cdot 10^9$.

4 Preliminary human baselines

As mentioned in the introduction, humans can infer programs from few examples and extrapolate them beyond the training distribution. At the same time, they also have limited search capacity and cannot mentally execute large programs. To assess the difficulty of our dataset we choose a random subset of 30 problems such that the number of executed lines was at most 5. We tested 5 humans with some prior C++ exposure in high school, but who did not necessarily do competitive programming in college or majored in Computer Science. With these criteria, subjects had a prior on reasonable mathematical functions that could be involved, while having very few priors on typical functions used in competitive programming.

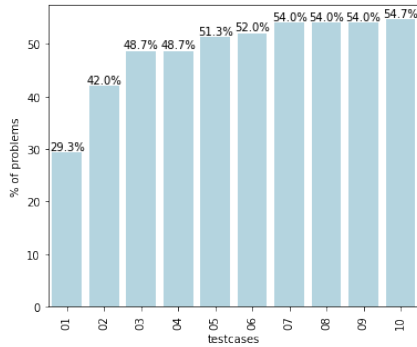


Figure 4: Fraction of problems solved by humans after seeing n examples; most problems only require a couple of examples, but there is still significant progress until 5 examples.

Subjects saw 10 test-cases and had to describe the program (expressed in human language) that they believed generated the outputs from the inputs. Out of all 30 problems, 13 problems were solved by all 5 subjects, and 10 were solved by some, but not all subjects. Each subject solved between 14 and 19 problems. These results are encouraging because they show that most problems (at least $25/30 \approx 83\%$) are feasible to infer, with a significant fraction ($1/3$) being non-trivial. Even for tasks solved by everyone, it is likely that this is still far from what most methods can achieve at the moment, providing a challenging benchmark for the meta-learning and program induction communities.

5 Discussion

We present a new benchmark for few-shot extrapolation. We hope this sparks progress in making powerful meta-learning algorithms that learn representations that generalize far beyond their distribution by exploiting compositionality. At the same time the structure of the dataset also enables other interesting problems such as learning to search, generalizing to inputs of different types than tasks seen at meta-training time, or meta-active learning by exploiting that we have oracles for each task.

References

- [1] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [2] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, pages 3630–3638, 2016.
- [3] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157*, 2019.
- [4] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. The omniglot challenge: a 3-year progress report. *Current Opinion in Behavioral Sciences*, 29:97–104, 2019.
- [5] Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- [6] Samy Bengio, Yoshua Bengio, and Jocelyn Cloutier. On the search for new learning rules for anns. *Neural Processing Letters*, 2(4):26–30, 1995.
- [7] Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 1998.
- [8] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020.
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017.
- [10] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*, pages 4077–4087, 2017.
- [11] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- [12] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- [13] Victor Garcia and Joan Bruna. Few-shot learning with graph neural networks. *arXiv preprint arXiv:1711.04043*, 2017.
- [14] Ferran Alet, Tomas Lozano-Perez, and Leslie P. Kaelbling. Modular meta-learning. In *Proceedings of The 2nd Conference on Robot Learning*, pages 856–868, 2018.
- [15] Chelsea Finn. *Learning to Learn with Gradients*. PhD thesis, EECS Department, University of California, Berkeley, Aug 2018.
- [16] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren, editors. *Automated Machine Learning: Methods, Systems, Challenges*. Springer, 2018. In press, available at <http://automl.org/book>.
- [17] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676*, 2018.
- [18] Antreas Antoniou, Massimiliano Patacchiola, Mateusz Ochal, and Amos Storkey. Defining benchmarks for continual few-shot learning. *arXiv preprint arXiv:2004.11967*, 2020.
- [19] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019.
- [20] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. *arXiv preprint arXiv:1903.03096*, 2019.
- [21] Ferran Alet, Martin F. Schneider, Tomas Lozano-Perez, and Leslie Pack Kaelbling. Meta-learning curiosity algorithms. In *International Conference on Learning Representations*, 2020.
- [22] Louis Kirsch, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Improving generalization in meta reinforcement learning using learned objectives. *arXiv preprint arXiv:1910.04098*, 2019.

- 269 [23] Junhyuk Oh, Matteo Hessel, Wojciech M Czarnecki, Zhongwen Xu, Hado van Hasselt, Satinder
270 Singh, and David Silver. Discovering reinforcement learning algorithms. *arXiv preprint*
271 *arXiv:2007.08794*, 2020.
- 272 [24] Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sable-Meyer, Luc Cary, Lucas Morales,
273 Luke Hewitt, Armando Solar-Lezama, and Joshua B Tenenbaum. Dreamcoder: Growing
274 generalizable, interpretable knowledge with wake-sleep bayesian program learning. *arXiv*
275 *preprint arXiv:2006.08381*, 2020.
- 276 [25] Sumit Gulwani. Automating string processing in spreadsheets using input-output examples.
277 *ACM Sigplan Notices*, 46(1):317–330, 2011.
- 278 [26] Sumit Gulwani, José Hernández-Orallo, Emanuel Kitzelmann, Stephen H Muggleton, Ute
279 Schmid, and Benjamin Zorn. Inductive programming meets the real world. *Communications of*
280 *the ACM*, 58(11):90–99, 2015.
- 281 [27] François Chollet. On the measure of intelligence, 2019.
- 282 [28] Josh Rule. *The child as hacker: building more human-like models of learning*. PhD thesis, MIT,
283 2020.
- 284 [29] Rajeev Alur, Dana Fisman, Rishabh Singh, and Armando Solar-Lezama. Sygus-comp 2017:
285 Results and analysis. *arXiv preprint arXiv:1711.11438*, 2017.
- 286 [30] Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and
287 Percy S Liang. Spoc: Search-based pseudocode to code. In *Advances in Neural Information*
288 *Processing Systems*, pages 11906–11917, 2019.
- 289 [31] Michihiro Yasunaga and Percy Liang. Graph-based, self-supervised program repair from
290 diagnostic feedback. *arXiv preprint arXiv:2005.10636*, 2020.
- 291 [32] Maksym Zavershynskyi, Alex Skidanov, and Illia Polosukhin. Naps: Natural program synthesis
292 dataset. *arXiv preprint arXiv:1807.03168*, 2018.
- 293 [33] V. Vassilev, Ph. Canal, A. Naumann, L. Moneta, and P. Russo. Cling – the new interactive
294 interpreter for ROOT 6.