
Programming Puzzles

Tal Schuster
MIT

Ashwin Kalyan
Allen Inst. for AI

Oleksandr Polozov
Microsoft Research

Adam Tauman Kalai
Microsoft Research

Abstract

We introduce a new type of programming challenge called programming *puzzles*, as an objective and comprehensive evaluation of program synthesis, and release an open-source dataset of Python Programming Puzzles (P3).¹ Each puzzle is defined by a short Python program f , and the goal is to find an input which makes f return True. The puzzles are objective in that each one is specified entirely by the source code of its verifier f , so evaluating f is all that is needed to test a candidate solution. They do not require an answer key or input/output examples, nor do they depend on natural language understanding. The dataset is comprehensive in that it spans problems of a range of difficulties and domains, ranging from trivial string manipulation problems, to classic programming puzzles (e.g., Tower of Hanoi), to interview/competitive-programming problems (e.g., dynamic programming), to longstanding open problems in algorithms and mathematics (e.g., factoring). We develop baseline enumerative program synthesis, GPT-3 and Codex solvers that are capable of solving puzzles—even without access to any reference solutions—by learning from their own past solutions. Codex performs best, solving up to 18% of 397 test problems with a single try and 80% of the problems with 1,000 tries per problem. In a small user study, we find a positive correlation between puzzle-solving performance and coding experience, and between the puzzle difficulty for humans and AI solvers. Therefore, further improvements on P3 could have a significant impact on many program synthesis areas.

1 Introduction

Puzzles are often used to teach and evaluate human programmers. Classic puzzles such as the Tower of Hanoi teach fundamental concepts such as recursion. Programming competition problems, also referred to as puzzles [34], evaluate a participant’s ability to apply these concepts. Puzzles are also used to evaluate programmers in job interviews, and puzzles such as the RSA-factoring challenge test the limits of state-of-the-art algorithms. Each of these types of puzzles is described in its own format, often in a natural language such as English. Evaluations often include a hidden test set.

We introduce a novel puzzle representation called a programming puzzle or simply a *puzzle*, which captures the essence of these challenges in a form convenient for machines and programmers. At a minimum, a puzzle is specified by a function $f(y)$, and the goal is to find y such that $f(y) = \text{True}$. More generally, a puzzle can include input variables x . Then, the puzzle can be seen as an output verifier $f(y, x)$ that validates y . The answer y is typically the output of a synthesized program g . In order to find $g(x) \rightarrow y$, a *synthesizer* is given the source code of f (and possibly x), with the goal of generating a program g such that $f(g(x), x) = \text{True}$. Importantly, puzzles make for an objective and explicit programming evaluation based solely on code with no formal requirement for input/output examples, natural language descriptions, or reference solutions.

Puzzles may have multiple valid outputs y and some puzzles, even if very short, are extremely challenging. Figure 1 illustrates three puzzles that are diverse in domain, difficulty, and algorithmic

¹<https://github.com/microsoft/PythonProgrammingPuzzles>

```

# Find a string that when reversed and concatenated with "world" gives "Hello world"
def f1(y: str):
    return y[::-1] + "world" == "Hello world"

# Tower of Hanoi, often teaches recursion. Move [i, j] means move top disk on tower i to j, with 1 ≤ i, j ≤ 3
def f2(moves: List[List[int]], num_disks=8):
    state = [1] * num_disks # All disks start at tower 1.
    for [i, j] in moves:
        assert state.index(i) <= (state + [1, 2, 3]).index(j), "bigger disk on top"
        state[state.index(i)] = j # Move smallest disk from tower i to tower j.
    return state == [3] * num_disks # All disks must end on tower 3.

# Find a non-trivial integer factor d of a large number n
def f3(d: int, n=100433627766186892221372630609062766858404681029709092356097):
    return 1 < d < n and n % d == 0

```

Figure 1: Programming puzzles ranging from trivial to longstanding open algorithmic challenges in multiple domains. `f1` is solved by `y="Hello "[: -1]`, a recursive program (see Figure H.1 on page 37) outputting 255 moves solves `f2`, and `f3` requires computational number theory algorithms.

tools. The first puzzle is an easy (for humans) puzzle that tests one’s understanding of basic syntax and properties of strings. The second is the quintessential example of recursion, and the last is a hard problem requiring advanced algorithms such as the quadratic sieve.

We also release a growing open-source Python Programming Puzzles dataset, called P3, which is already comprehensive in terms of difficulty, domain, and algorithmic tools. This dataset unifies many of the types of puzzles mentioned above. While P3’s puzzles are all specified in Python, solution programs g can be written in any language, or simulated by neural networks.

As describe in §3, P3 also contains numerous classic puzzles; optimization puzzles such as solving a linear programming; graph puzzles such as shortest path; and competitive-programming problems. The most difficult puzzles involve longstanding open problems such as learning parity with noise [11]; factoring [24, 30]; or finding a cycle in the $3n + 1$ process which would disprove the Collatz conjecture [35]. Thus, if AI were to surpass human-level performance on this dataset, it would lead to breakthroughs on major open problems. The P3 puzzles were inspired by sources such as Wikipedia, algorithms books, programming competitions, and other program synthesis datasets. This dataset is growing as an open-source project, and anyone can add a puzzle by simply writing a function f .

One motivation for programming puzzles is that improvements in solving puzzles may lead to performance gains at other tasks. Recently, neural Language Models (LMs) have advanced the state-of-the-art of AI systems performance in offering code completions and synthesizing source code in general-purpose programming languages, such as Python, based on English descriptions [5, 13]. While such systems represent a major advance over prior approaches, Chen et al. [13] point out that they also reproduce elementary programming mistakes. Figure 2 illustrates how state-of-the-art GitHub Copilot [63] solves a complex problem, handles the ambiguity of English, and yet makes elementary errors. Performance gains in solving programming puzzles may result in fewer errors or solving more sophisticated algorithms problems in downstream tasks such as code completion.

Puzzles are *objective*, meaning that it is easy to unambiguously evaluate whether one’s own answer is valid without consulting an answer key. This evaluation also allows bootstrapping, even on test puzzles without gold solutions. Given a set of puzzles (f_i, x_i) , one can attempt to solve them with solutions g_i , determine with certainty which solutions are correct, and use those to improve one’s ability to solve the remaining puzzles [19]. Inspired by success in playing games [53, 56], self-training has also proven useful in program synthesis [see, e.g., 6, 15]. Other commonly used representations, including natural language descriptions or Programming by Example (PbE), have inherent ambiguity. See Appendix F for a comparison of a competition problem represented in English and as a puzzle.

From a theoretical point of view, as we shall discuss, objectivity can be formalized as the complexity class NP of non-deterministic polynomial-time decision problems. Moreover, the puzzle decision problem is NP-complete, meaning puzzles can readily express any NP problem, including polynomial-time problems and other NP-complete problems such as Boolean satisfiability.

We compare several enumerative random forest and Transformers-based top-down solvers, as well as GPT-3 and Codex LM solvers with different prompt types (e.g., zero/few-shot and with/without English descriptions). In our experiments, without access to any reference solutions, only utilizing

```

def adjacent_primes(n: int):
    """Find the nth & (n+1)st prime numbers"""
    primes = [2, 3]
    i = 3
    while len(primes) < n:
        if all(i % p != 0 for p in primes): ...

```

```

# x is the concatenation of "Hello" and "world"
x = "Hello" + " " + "world"
assert "hello" not in x
assert len(filename + ".json") < len(filename)

```

Figure 2: GitHub Copilot code completion examples (in gray). Left: Copilot correctly implements a seven-line function. Top right: the completion adds a space character that may or may not have been intended by the user. Middle and bottom right: errors indicating a lack of basic understanding.

self-training bootstrapping, our enumerative models solved up to 43% more P3 problems than a naive brute force baseline. Our LM solvers were able to solve many of the puzzles, given enough tries.

To address the questions of whether puzzles measure programming proficiency and how puzzle difficulty compares between humans and computers, we performed a small user study. Puzzles were accessible and enjoyable for programmers with varying levels of experience. While both GPT-3 and enumerative techniques can solve a fraction of the puzzles, human programmers outperform them. For example, bootstrapping GPT-3 with up to 10K tries solved 60% of the puzzles, lower than both beginner and experienced participants that solved 76% and 87% puzzles on average, respectively. Overall, we find perceived puzzle difficulty to scale similarly for both humans and AI.

The main contributions of this paper are introducing:

1. programming puzzles: a new type of problem suitable for algorithmic problem-solving (for both machines and humans);
2. P3, an open-source dataset of puzzles covering diverse domains and difficulties; and
3. an evaluation of humans and baselines demonstrating that puzzles can be used to measure algorithmic problem-solving progress.

Progress in code completion is rapid—even between the time of submission of this paper and its publication, an API to Codex (a GPT-3 model fine-tuned for code completion) was released [13]. Our evaluation does in fact show significant improvements of Codex over other baselines.

2 Problem formulation

Programs, inputs and outputs can all be formally represented as strings, where Σ^* is the set of finite strings over alphabet Σ . The set of verification functions is denoted by $\mathcal{F} \subseteq \Sigma^*$, with inputs and outputs $\mathcal{X}, \mathcal{Y} \subseteq \Sigma^*$, respectively. A puzzle is defined by pair $(f, x) \in \mathcal{F} \times \mathcal{X}$ and the result of running verifier f on output $y \in \mathcal{Y}$ is denoted $f(y, x) \in \{0, 1\}$. Output $y \in \mathcal{Y}$ is *valid* if it *satisfies* $f(y, x) = 1$, i.e., f outputs 1 when run on (y, x) , within a specified amount of time. To ensure that puzzles can be *quickly* verified, it is necessary to upper-bound the time required for puzzle verification. This ensures that the puzzle decision problem, namely the problem of determining whether, given f, x , there is y such that $f(y, x) = 1$, is in the complexity class NP. Formally, the puzzle decision problem is, given strings f and x denoting the puzzle (represented as, say, a Turing machine) and input, and a timeout t , does the puzzle output 1 in time $\leq t$. See Appendix D for further details.

A **solver** takes n puzzles and timeouts $(f_1, x_1, t_1), \dots, (f_n, x_n, t_n)$, and produces outputs y_i to as many puzzles as it can within a time bound T . Of course $T \gg \sum t_i$ is significantly larger than the verification timeouts. Formally, the *score* of solver $S : \mathcal{F}^n \rightarrow \mathcal{X}^n$ is the number of puzzles f_i for which $f_i(y_i, x_i)$ outputs 1 in time $\leq t_i$. Although we do not study it in this paper, it would be natural to assign different *values* to different puzzles. For example, solving open problems such as finding a Collatz cycle or factoring the largest RSA challenge integer (currently unsolved, with a \$200,000 prize offered), should be of greater value than solving a simple hello-world puzzle.

It is convenient, though not required, to solve puzzles by outputting a program g which, when run, computes output $y = g(x)$. Such a program is called a **solution** g . Short solutions may have long outputs, e.g., the puzzle `(f=lambda y: len(y) == x, x=1000000)` requires a string of length one million as solution `g=lambda x: 'a' * x`. In this example, $y = g(1000000)$ is a valid output of length one million. Of course, another solution would be to explicitly write a string of length one million in the code, though this implementation may not pass a human code review. In the dataset and this paper, we provide solutions since they may be significantly shorter. Many puzzles fit a single

problem template, meaning they share the same verifier f but have different inputs x . Thus a dataset may have many more puzzles than problems.

3 The P3 dataset

P3 uses Python, the de facto language of ML research, as the programming language for specifying puzzles. At the time of publication, P3 currently has 397 problems, summarized in Table 1. The latest dataset can be generated by simply running `make_dataset.py` in the repository, with parameter n for how many puzzles per problem are generated: $n = 1$ and $n = 100$ yield 397 or 32,932 puzzles, respectively. Most experiments in this paper use only one puzzle per problem. Every puzzle is described by a function with a required typed argument (i.e., the candidate output) that returns True upon success. Since Python is not type-safe, we add type-checking to ensure that outputs match the declared type. Figure C.1 on page 19 illustrates a puzzle where type checking is important.

We also provide code for serializing Python objects to and from strings in a json-like format, so that programs implemented in any language can produce outputs. Moreover, strings are universal in that they can encode arbitrary Python objects including functions, as in the [Quine](#) puzzle (`lambda quine: eval(quine) == quine`)² motivated by the classic challenge of writing a program that outputs its own source code. As evaluation of the string `quine` can lead to an infinite loop, this puzzle illustrates the necessity of the evaluation timeout t for attempted solutions.

While not necessary for evaluation (since puzzles are self-contained) we follow the common practice of programming competitions and provide a reference solution to most (over 90%) of the puzzles. Some puzzles have more than one solution. A handful of puzzles represent major open problems in computer science and mathematics including [Factoring](#) (and [Discrete Log](#)), [Planted Clique](#), [Learning Parity with Noise](#), [Graph Isomorphism](#), and finding a [Collatz cycle](#),³ as described in Appendix E. We also provide English descriptions for each puzzle in the dataset to support research involving natural language. Appendix F compares programming competition problems to puzzles.

Creation process. The following sources were used for identifying possible puzzles:

- Wikipedia, specifically the [Logic puzzles](#) category, the [List of unsolved problems in mathematics](#), and the [List of algorithms](#).
- Competitions, primarily the competitive programming website [codeforces.com](#) but also a handful of problems from the [International Collegiate Programming Contest](#) and the [International Mathematical Olympiad](#) (IMO)—a high school mathematics competition.
- Puzzles inspired by the HumanEval dataset used for evaluating Codex [13], added in v0.2.
- The Python programming language itself, with trivial puzzles created to test understanding of basic functions, such as the the hello-world puzzle which tests string concatenation.

P3 is organized topically into modules listed in Table 1. These topics include domains such as number theory, graph theory, chess puzzles, game theory, etc., as well as puzzles inspired by a specific source such as a specific programming competition. One finding in this paper is that many types of puzzles can be captured in spirit, if not exactly, as succinct puzzles. Common patterns include:

- Problems that are *naturally* puzzles. For instance, search problems such as the [TowerOfHanoi](#) (f2, Figure 1) and [SlidingPuzzle](#) simply test the sequence of moves to see if they lead to the goal state.
- Problems that have an equivalent natural puzzle. For instance, the standard definition of the factoring problem, namely factorizing an integer into its prime factors would require a puzzle that tests primality. However the simpler problem of [finding any non-trivial integer factor](#), f3 in Figure 1, can be recursively called to solve the prime factorization problem.
- Optimization problems. Some such problems have equivalent natural puzzles, e.g., linear programming is well-known [18] to be equivalent to solving a zero-sum game which is the [ZeroSum](#) puzzle. For others, such as [LongestMonotonicSubstring](#) or [ShortestPath](#), we specify a bound θ on the objective, and the goal is to find a feasible y with objective better than θ . In order to generate θ (included in x), we first solve the optimization problem ourselves, but the puzzle generation code is not provided to the solvers.

²GPT-3 generated a 5-character solution to the quine puzzle while the authors’ solution was 88 characters.

³The solution to this problem would disprove the Collatz conjecture that is believed to be true, but no proof has been found yet. Therefore, if the conjecture is true, the maximum attainable score in P3 is $< 100\%$.

Table 1: Number of problems (and how many of them have at least one reference solution) per domain in P3 v0.2. The right two columns show the average size of puzzles and solutions, measured by the number of nodes in the Python AST.

| Domain | Problems | Solutions | $ f $ | $ g $ |
|------------------------|----------|-----------|-------|-------|
| Algebra | 4 | 4 | 70 | 172 |
| Basic | 23 | 23 | 54 | 44 |
| Chess | 5 | 3 | 221 | 186 |
| Classic puzzles | 23 | 23 | 101 | 211 |
| Codeforces | 47 | 45 | 73 | 70 |
| Compression | 2 | 2 | 126 | 113 |
| Conways Game of Life | 3 | 2 | 189 | 345 |
| Games | 7 | 7 | 225 | 299 |
| Graphs | 12 | 11 | 105 | 152 |
| HumanEval | 164 | 164 | 81 | 62 |
| ICPC | 4 | 4 | 304 | 569 |
| IMO | 6 | 6 | 173 | 256 |
| Lattices | 2 | 2 | 70 | 228 |
| Number Theory | 16 | 12 | 47 | 68 |
| Probability | 5 | 5 | 85 | 72 |
| Study | 30 | 30 | 40 | 21 |
| Trivial inverse | 39 | 38 | 27 | 30 |
| Tutorial | 5 | 5 | 27 | 13 |
| Total # / Average size | 397 | 386 | 79 | 84 |

- Problems that ask how many items in a certain set satisfy a given property, may be converted to problems that require an explicit enumeration of all such items. See for example the [AllPandigital-Squares](#) puzzle that requires all 174 pandigital perfect squares as input.
- Problems that involve game-playing can often be converted to puzzles. In chess, this includes the classic Eight Queens and Knights Tour search problems. A puzzles Mastermind involves exhibiting a winning strategy tree, and a nim puzzle involves beating a given computer opponent.

In order to ensure that each puzzle is achieving its goals, the puzzle design process has a step which automatically tests for trivial solutions such as small integers or common strings.

Exclusions. Many programming challenges do not make as good puzzles. First, simple *translation* tasks, where goal is translate a sequence of steps described in natural language into a program, do not make good puzzles. Second, some challenges require problem-solving that is not easily expressed as a program. For example, computing the probability of rolling a total of 350 when rolling 100 dice relies on external knowledge about probability theory. Third, “soft” challenges involving natural language or images are not in NP and not easily verifiable. This includes challenges involving human commonsense or world knowledge about names, dates, or image classification. Finally, *interactive* challenges do not make for good programming puzzles. Fortunately, several other benchmarks cover these latter two types of exclusions [see, e.g., 3, 32, 37, 42, 46, 48, 49, 51–53, 55, 58, 61, 62].

Growth process. The focus of this paper is in creating a framework with an initial dataset; and demonstrating its utility for developing and evaluating AI solvers. As a GitHub repository, the dataset can grow over time in a standard manner with the ability to reference previous versions. We plan to continue adding puzzles and hope that others will as well. Popular competitive-programming websites such as [codeforces](#) may be a source of thousands of puzzles of varying difficulties.

4 Solvers

In this section, we describe the models we develop as baselines for the dataset. We consider both solving problems independently and joint solvers that bootstrap from previously obtained solutions to find new ones. We also consider both enumerative solvers that use standard techniques from program synthesis and LM solvers that use GPT-3 and Codex to solve puzzles. While a direct comparison between these two different approaches is difficult because they run on different hardware (the LMs call an API), we can still compare the relative difficulty with which they solve different puzzles, and also to human difficulty rankings among puzzles.

4.1 Enumerative solvers

Following prior work [1, 6, 15, 40], we develop models to guide the search for g over the space of all possible functions. In particular, we implement a grammar that generates Abstract Syntax Trees (ASTs) for a large subset of Python. The grammar covers basic Python functionality and is described in Appendix B.1. Specifically, each Python function is translated to an AST using a given set \mathcal{R} of rules. Based on the puzzle, a context-specific distribution over rule probabilities is computed. To facilitate efficient top-down search, the context of a rule is defined to be the rule used by the parent node and the index of the current node among the parent’s children. Thus if the parent node was a division binary operator, then the two children would each have different contexts, but if two such divisions were present in the same program, both numerators would share the same context.

Each puzzle f is represented by a feature vector $\phi(f)$ and each context is represented by a vector $c(p, i)$ where p is the parent rule and i is the child index. Each rule $r \in \mathcal{R}$ is also associated with a feature vector $\rho(r)$. The probability distribution over \mathcal{R} is determined based on $\rho(r)$, $\phi(f)$, $c(p, i)$, and the likelihood of a solution g is the product of all rules constituting its AST. Naturally, this scoring mechanism introduces a bias towards shorter programs (i.e., smaller trees), which is desirable as a short solution is easy to inspect.

COPY rules. Solutions often reuse constants or puzzle parameters, for example the constant 25 or the variable `s` in example `f2` in Figure 1. As in prior work [40], for each puzzle, the global rules bank is expanded to include COPY rules for constants and parameters of the examined puzzle.⁴ When composing solutions, this rule can reduce the complexity of the solution by simply learning to copy part of the puzzle rather than having to generate it from scratch. For simplicity, we create copy rules for each of the supported types and assign the probabilities uniformly across all the puzzle’s constants of that type. In other words, our models learn when a certain type should be copied from the puzzle, and rank all available constants and parameters of that type the same.

To solve a new puzzle, we perform a top-down search. Specifically, at each node, we apply a selected model over all rules in \mathcal{R} whose type matches the context, and re-normalize the scores to create a valid probability distribution. The solver enumerates solutions in order of decreasing likelihood until it finds a solution g such that `f(g())` evaluates to `True` in time $\leq t$, for a maximum number of tries M . See Appendix B for details on the search and rules. Next, we briefly describe our models.

Uniform. The first model is a simple uniform rule that assigns the same probability to all rules. The only exception is COPY rules, which have a larger, fixed probability in order to bias the solver towards utilizing this option. As we score programs by their joint probability, this bias effectively favors shorter programs. We use this model to find solutions to the easier problems, satisfied by a simple and short answer, and use these to bootstrap the learning of the parametric models. This model also provides a naive brute force baseline to compare the parametric models with.

The remaining two models have parameters that are fit based on *bootstrapping*. Namely, given previously obtained solutions, we collect all parent-child rule pairs as self-supervision and fit the model’s parameters on them. The training size is then the total number of nodes in all the trees among solutions discovered up until that point. We implement two bigram parametric models to predict $\mathbb{P}(r \mid \rho(r), \phi(f), c(p, i))$, where r is a candidate rule to appear in g ’s tree under p as its i ’s argument.

Random forest. In this model, we represent f as a bag-of-rules $\cup\{r_k \in f\}$. Specifically, $\phi(f)$ is a vector of length $|\mathcal{R}|$ representing the number of occurrences of each rule in f . p and i are encoded as a one-hot vector and concatenated to f ’s representation to construct the input to the model. Given past solution trees, we train the model to predict the index of r out of $|\mathcal{R}|$ given f, p, i examples.

Transformer. Following the recent success of transformer models [20, 57] in encoding source code [21, 31, 54, *inter alia*], we turn to these encoders for richer representations. We use a RoBERTa-based [36] Transformer to encode puzzles and rules directly from their code. The probability of a rule r being the i ’s child of p in g is proportional to the dot product of the deep joint representation of f, p, i and the Transformer encoding $\rho(r)$. We pretrain the Transformer with a masked language model task on Python GitHub repositories [27].⁵ Then, our solver concatenates the Transformer

⁴When executing a solution, COPY rules are simply the identity function (COPY = `lambda x: x` in Python).

⁵Our pretrained model and tokenizer are available at https://huggingface.co/tals/roberta_python.

```
def f(li: List[int]):
    return len(li) == 10 and li.count(li[3]) == 2

assert True == f(...
```

Figure 3: A Short prompt for a puzzle requesting a list of ten integers where the fourth item occurs exactly twice, with valid completion `... [1, 2, 3, 4, 5] * 2`). Appendix C has Medium/Long prompts.

encodings $\phi(f)$ and $\rho(p)$ with a learned embedding for i , following by non-linear layers to compute the joint representation. We fine-tune the solver on parent-child rule pairs from previously acquired solutions. See Appendix B.2 for extended details, and Figure B.1 on page 19 for a model diagram.

4.2 Autoregressive Language Model solvers

We experiment with the transformer-based GPT-3 [12] and Codex [13] LMs with billions of parameters. Codex was trained over large amounts of publicly available code from GitHub, specifically aimed for coding applications. We follow the recent strategy of designing a prompt that directs the text generation to our desired task. This approach has shown to be useful in converting natural language descriptions to programming code and guide theorem proving [45]. Unlike our enumerative models that build an AST, LMs generate the solution as a string that is directly evaluated as Python code.

We consider four different prompts: (a) A *short* zero-shot prompt based solely on the puzzle at hand (illustrated in Figure 3); (b) a *medium* 5-shot prompt that includes the five example puzzles that had been shown to (human) programmers during our user study (Appendix Figures C.2-C.3); (c) a *long* prompt with the same five puzzles augmented by English descriptions of the tasks in comments (Figures C.4-C.5); and (d) a *bootstrapping* prompt which uses only solutions to problems that it has already solved (Figures C.7). The bootstrapping prompt begins with no solutions but quickly exceeds the API maximum length as more puzzles are solved. At that point, previously solved puzzles are randomly sampled to form the prompt. The prompts used for Codex are slightly more sophisticated but enable multi-line programs.

The completions which parse as valid Python expressions are then evaluated. Appendix C gives further details of the execution environment, the API parameters and other prompts we investigated.

5 Experiments

We use our P3 dataset to evaluate the performance of the solvers from §4. We assume no access to reference solutions⁶ and measure how many puzzles are solved by each solver with up to k tries per puzzle, where each try is a potential solution that is evaluated. For the enumerative solvers, this is equivalent to having a valid solution ranked in the top k . For LM solvers, we use `pass@k` [13] which is an unbiased estimator of the probability of obtaining a solution within k tries. First, we test the solvers bootstrapping efficacy in leveraging previously solved problems to solve new ones. Then, once solutions to a single instance of certain problems are found, we test whether solvers also succeed on other problem instances (i.e., puzzles originated from the same problem). In §5.1, we present our user study results that compares human’s performance with AI solvers. Finally, in §5.2, we test whether P3 can distinguish between subtly different variants of Codex, using the larger v0.2 release of the dataset (the current version at the time of publication).

Learning from past solutions. This first experiment was run on the v0.1 release of P3.⁷ We use a single puzzle instance per problem. We first identified the 138 of the 200 v0.1 problems supported by our grammar (see Appendix B.1). For the enumerative solvers, we then ran the uniform solver with $k = 10^4$ on these 138 problems supported by our, solving 38 of them. The solutions contain a total of 2,475 rules that we use to train the parametric models. In the bootstrapping variant, we repeat the training for 6 cycles, each time adding the new solutions found with $k = 10^4$. In the final round, we allow up to $k = 10^6$ solving tries (including previous cycles). For comparison to GPT-3/Codex,⁸ we

⁶With the exception of the Medium and Long prompts that including five Tutorial problems and solutions.

⁷<https://github.com/microsoft/PythonProgrammingPuzzles/tree/v0.1>

⁸The Codex API was released after this experiment had been run on v0.1 using the enumerative and GPT-3 solvers. Thus, we simply replaced the GPT-3 solver with the Codex solver and re-ran on the same 138 puzzles.

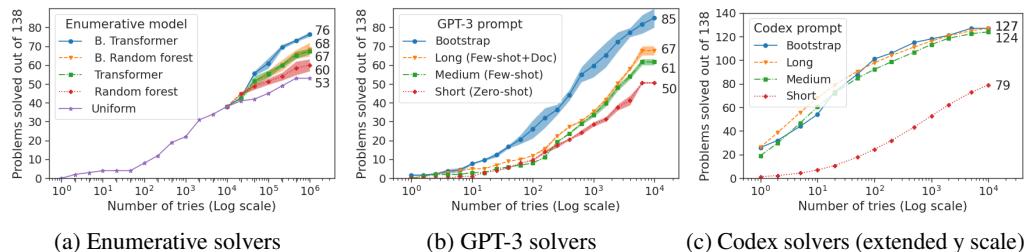


Figure 4: Increasing the number of tries allows solving new problems. Better solvers, though, solve new problems significantly faster by learning from past experience. Parametric enumerative solvers (a) initialized with the solutions of the uniform solver at $k = 10^4$ accelerate the solution search. Additional self-training bootstrapping cycles (marked with B.) solve even more problems. GPT-3 (b) and Codex Davinci (c) solvers were evaluated with up to 10^4 attempts. Having natural language descriptions (Long) provides small improvements over Medium. Adding previously found solutions to the prompt (Bootstrap) allows significant improvements for enumerative and GPT-3, and matches Long for Codex. Overall, the Codex models performed best, solving up to 127 of the examined 138 puzzles. (a), (b) are averaged across three runs and the shaded areas show the standard deviation.

use the same 138 problems and start with a zero-shot prompt. As valid solutions are found, they are appended to the prompt as discussed in §C.

Figure 4a shows the total number of puzzles solved by each enumerative solver, with and without the self-training bootstrapping cycles. We report the average results across three runs and present the standard deviation in the graph. We see that the parametric models quickly improve over the naive uniform search and that the bootstrapping process facilitates solving many new problems. At $k = 10^6$, the random forest and Transformer-based enumerative models solved a total of 68 and 76 problems, respectively, which is 28% and 43% more than the uniform solver.

The GPT-3 solver also improves by learning from previously found solutions. As Figure 4b shows, few-shot settings with tutorial examples perform better than zero-shot (Short) and solve new problems. Including natural language descriptions (Long) helps for solving five more puzzles, with up to 10^4 tries. The best strategy, however, is the bootstrapping one that starts without any reference and adds solutions to the prompt as they are found. Codex, trained on large amounts of code, performs the best (see Figure 4c) but does not benefit significantly from bootstrapping.

Generalizing to other problem instances. In the previous experiment, we attempted to solve the *default* single puzzle instance of each problem. Next, we examine whether our solvers can also solve other puzzle instances, originating from the same problems. We collect a set of 700 puzzles that are random instances of 35 problems for which both our bootstrapping enumerative models solved the default puzzle. At $k = 10^4$, the random forest and Transformer models solved 75% and 79%, respectively. As a reference, the uniform model solves only 62% of these puzzles.

5.1 User study

In a small user study, 21 participants with varying experience in Python programming attempted to solve 30 puzzles, as found in v0.1 dataset as the study module. Each puzzle was allocated a maximum of 6 minutes to solve, and the study was conducted virtually using Jupyter notebooks. Participants were employees at a major software company and were recruited by email and at a hackathon. No compensation was offered. Participants were first given a short tutorial about puzzles and how to submit solutions. The user study files are available in the open-source dataset, and Appendix G has further details including the 30 puzzles.

The first finding is that success in puzzles correlates with programming experience. For our retrospective study analysis, we split the participants by the median years of Python programming experience. We had 10 *beginners* with less than three years of experience, and 11 *experienced* participants with at least three years. We find that 9 of the 30 puzzles were solved by all beginners, while 17 of the puzzles were solved by all experienced participants. Also, beginners spent on average 194 seconds per puzzle, while experienced spent only 149 seconds on average. The average solving time provides

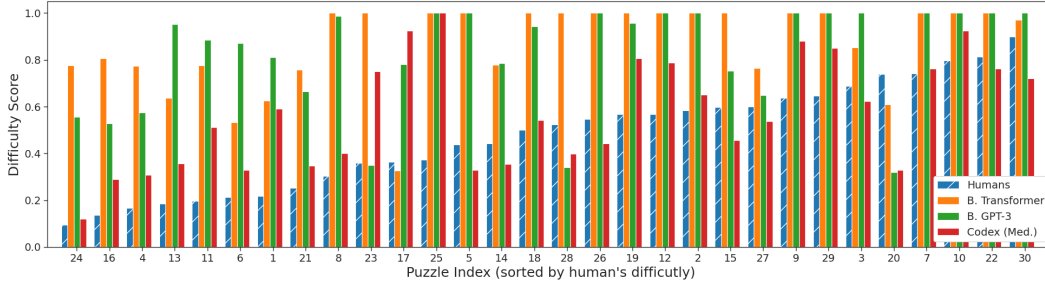


Figure 6: The difficulty score per study puzzle for both humans and AI solvers, sorted by the human’s scores. The difficulty score for humans is measured by the average fraction of solving time out of the maximum allowed. For AI, we use the fraction of allotted attempts required. Most of the puzzles solved by AI (low difficulty score) are also easier for humans (left hand side of the plot).

a useful proxy to the perceived difficulty of each puzzle. Overall, we see that puzzles are easier for experienced programmers, indicating their value for evaluating programming proficiency.

Next, we compare human’s performance to Codex-davinci. We use the Medium prompt as it is similar to the study format (i.e., same 5 tutorial examples, no docstrings). Participants solved an average of 24.6 out of the 30 puzzles (22.9 for beginners and 26.2 for experienced) within the 6 minutes per puzzle time limit. Only one out of the 21 participants solved all puzzles. As Figure 5 shows, Codex required 1K tries per puzzle to match the performance of beginner programmers in our study.

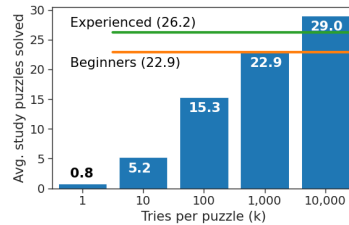


Figure 5: Number of solved puzzles by Codex-davinci (blue bars), compared to human coders with 6 minutes per puzzle (horizontal lines).

Finally, we find that difficult puzzles for humans are also harder for AI. Figure 6 shows that most of the puzzles solved by AI solvers with limited number of tries are the ones that are easier for humans (i.e., solved faster). To compare the two, we define a puzzle’s perceived difficulty score as the average solving time for humans and the expected number of required tries for machines (normalized to $[0, 1]$, where the score of unsolved puzzles is set to 1). The Spearman’s rank coefficient of humans with B. GPT-3 is 0.512, and with Codex (Med.) is 0.563. The AI solvers correlation is stronger with beginner programmers (0.541 and 0.562), than with the experienced ones (0.470 and 0.544, respectively). On the one hand, this suggests that additional computational power might allow AI solvers to match humans. However, as Figure 4 shows, this improvement is logarithmic, leading to diminishing returns. Encouragingly, we see that even within the same budget, modeling choices can improve performance. We hope that P3 will support the research and development of new AI solvers that will solve more puzzles with less computational effort.

5.2 Comparing small and large Codex models

In addition to the standard *davinci-codex* engine, the API offers an alternate *cushman-codex* engine that they report is significantly faster and only slightly less accurate. To test the ability of P3 as an evaluation of such fine distinctions, we ran the Medium and Long prompts on both engines across the most recent v0.2 release⁹ of 397 puzzles. As can be seen in the results of Table 2, the larger engine indeed slightly outperformed the smaller engine across all k . Thus, in this experiment, puzzle solving success aligns with code completion success. Also, we observe that English descriptions (Long prompt) are helpful for both engines. Inasmuch as puzzles are useful for code completion, the $< 20\%$ success rates at $k = 1$ leaves substantial room for improvement.

Table 2: Codex $\text{pass}@k$ results over P3 v0.2.

| engine (prompt) | $k = 1$ | 10 | 100 | 1,000 |
|-----------------|---------|-------|-------|-------|
| cushman (Med.) | 7.1% | 26.7% | 51.7% | 68.3% |
| davinci (Med.) | 11.2% | 36.7% | 60.6% | 75.3% |
| cushman (Long) | 14.9% | 42.4% | 63.9% | 76.5% |
| davinci (Long) | 18.3% | 48.7% | 69.1% | 79.8% |

⁹<https://github.com/microsoft/PythonProgrammingPuzzles/tree/v0.2>

Table 3: Codex-davinci and Codex-cushman number of solved problems per domain with up to 1,000 tries for Medium and Long prompts. The first row also shows the number of available P3 v0.2 problems in that domain. The score is the average percent solved across domains.

| Model | Algebra | Basic | Chess | Classic | CodeForces | Compression | Conway’s | Games | Graphs |
|----------------|---------|-------|-------|---------|------------|-------------|----------|-------|--------|
| cushman (Med.) | 3/4 | 15/23 | 0/5 | 6/23 | 32/47 | 0/2 | 0/3 | 1/7 | 6/12 |
| davinci (Med.) | 2 | 20 | 0 | 8 | 35 | 0 | 0 | 1 | 9 |
| cushman (Long) | 2 | 21 | 0 | 5 | 38 | 0 | 0 | 1 | 8 |
| davinci (Long) | 2 | 22 | 1 | 4 | 39 | 0 | 0 | 1 | 8 |

| Model | HumanEval | ICPC | IMO | Lattices | N. Theory | Probability | Study | Trivial ⁻¹ | Tutorial | Score |
|----------------|-----------|------|-----|----------|-----------|-------------|-------|-----------------------|----------|-------|
| cushman (Med.) | 139/164 | 2/4 | 1/6 | 0/2 | 8/16 | 2/5 | 21/30 | 33/39 | 5/5 | 44.2 |
| davinci (Med.) | 145 | 2 | 1 | 1 | 9 | 3 | 22 | 36 | 5 | 51.2 |
| cushman (Long) | 149 | 1 | 1 | 1 | 9 | 3 | 24 | 36 | 5 | 49.8 |
| davinci (Long) | 155 | 1 | 1 | 2 | 10 | 3 | 25 | 38 | 5 | 54.8 |

Table 3 shows the number of achieved solutions per domain, as well as an overall score computed as the macro-average of solving rates across domains.

6 Related Work

Program synthesis has taken drastically different forms for different applications, often resulting in one-off evaluations rather than common datasets. A major paradigm is Programming by Example (PbE) where problems are specified by input-output examples. For instance, several studies focus on text processing [22] or robot navigation [43]. While convenient for end user applications (e.g., many in [44]), PbE alone is inadequate to objectively describe many sophisticated algorithmic programming challenges. A recent ARC dataset [14] adopts PbE for evaluating abstraction and reasoning in AI, but as in all PbE applications, there can be ambiguity.

Program synthesis from formal specifications has a long history of study [surveyed in 23], benchmarked by e.g., the SyGuS competition [2]. In this setting, however, the AI system has to synthesize an algorithm that correctly and efficiently solves a problem on all inputs (and often prove correctness as well). Writing and testing such formal specifications is often non-trivial.

English descriptions, often mixed with examples, are becoming an increasingly popular problem representation as LMs improve [28, 33, 60]. In independent work, Hendrycks et al. [26] created a large dataset of English programming problems with examples on which they fine-tuned GPT models. In another concurrent work, the Codex model that powers the new GitHub Copilot auto-completion tool [13] was evaluated with short problem descriptions paired with a set of unit tests that should validate the described specification. Our work, together with this very recent and concurrent work [5, 13, 26], represent the first controlled evaluation of large Transformer-based LMs on general-purpose program synthesis.

The recent CodeXGLUE benchmark [38] collected several code-related datasets. To evaluate generation, they use CodeBLEU [47] which relies on ASTs and other code-specific aspects. This evaluation still requires reference solutions and, therefore, does not resolve the answer-key bias with ambiguous specifications. Several neighboring fields that have made substantial progress in reasoning include theorem proving [8], two-player game playing [53], and SAT-solving [9]. In all these fields, important progress has been made by encoding the problems, be they theorems, game rules, or optimization problems, in machine-readable formats that do not involve the ambiguities of natural language.

7 Conclusions

We introduce Python Programming Puzzles (P3), an open-source dataset with puzzles described only in source code. As discussed in §3, the puzzle framework captures NP problems, which include a wide range of interesting challenges. Puzzles allow fast and objective evaluation, thereby supporting unsupervised solving without training solutions. We implemented and evaluated several enumerative program-synthesis and LM baselines, and found a positive correlation between their per-puzzle performance and the difficulty for human programmers. Similarly, LMs that performed better at code completion also solved more puzzles with less tries.

We welcome contributions to P3 and hope it will grow in size, coverage, and utility.

Acknowledgments. We would like to thank Mariia Mykhailova for suggesting doing a Python Programming Puzzles Hackathon. We are especially grateful to the participants in our user study and hackathon. We are grateful to the creators of Codex and GPT-3 and to Nicolás Fusi for suggesting its use in this project. We would like to thank David Alvarez Melis and Alec Helbing for suggesting quine puzzles. We are grateful to Ana-Roxana Pop for helpful discussions and feedback. We also thank Tianxiao Shen, Adam Fisch and the rest of the MIT NLP members for valuable writing feedback.

References

- [1] Uri Alon, Roy Sadaka, Omer Levy, and Eran Yahav. Structural language models of code. In *International Conference on Machine Learning*, pages 245–256. PMLR, 2020.
- [2] Rajeev Alur, Dana Fisman, Saswat Padhi, Rishabh Singh, and Abhishek Udapa. SyGuS-Comp 2018: Results and analysis. 2019.
- [3] Jacob Andreas, John Bufe, David Burkett, Charles Chen, Josh Clausman, Jean Crawford, Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eisner, Hao Fang, Alan Guo, David Hall, Kristin Hayes, Kellie Hill, Diana Ho, Wendy Iwaszuk, Smriti Jha, Dan Klein, Jayant Krishnamurthy, Theo Lanman, Percy Liang, Christopher H. Lin, Ilya Lintsbakh, Andy McGovern, Aleksandr Nisnevich, Adam Pauls, Dmitrij Petters, Brent Read, Dan Roth, Subhro Roy, Jesse Rusak, Beth Short, Div Slomin, Ben Snyder, Stephon Striplin, Yu Su, Zachary Tellman, Sam Thomson, Andrei Vorobev, Izabela Witoszko, Jason Wolfe, Abby Wray, Yuchen Zhang, and Alexander Zotov. Task-oriented dialogue as dataflow synthesis. *Transactions of the Association for Computational Linguistics*, 8:556–571, 2020. doi: 10.1162/tacl_a_00333. URL <https://aclanthology.org/2020.tacl-1.36>.
- [4] Sanjeev Arora and B. Barak. Computational complexity: A modern approach. 2009.
- [5] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021.
- [6] Matej Balog, Alexander L. Gaunt, Marc Brockschmidt, Sebastian Nowozin, and Daniel Tarlow. Deepcoder: Learning to write programs. In *International Conference on Representation Learning (ICLR)*, 2017.
- [7] E. Berlekamp, R. McEliece, and H. van Tilborg. On the inherent intractability of certain coding problems (corresp.). *IEEE Transactions on Information Theory*, 24(3):384–386, 1978. doi: 10.1109/TIT.1978.1055873.
- [8] Yves Bertot and Pierre Castéran. *Interactive Theorem Proving and Program Development: Coq’Art: The Calculus of Inductive Constructions*. Springer Science & Business Media, 2004.
- [9] Armin Biere, Marijn Heule, and Hans van Maaren. *Handbook of satisfiability*, volume 185. IOS press, 2009.
- [10] N. Biggs, P.M.L.S.E.N.L. Biggs, London Mathematical Society, London Mathematical Society lecture note series, and S.P.G.N.J. Hitchin. *Finite Groups of Automorphisms: Course Given at the University of Southampton, October-December 1969*. Lecture note series. Cambridge University Press, 1971. ISBN 9780521082150. URL <https://books.google.com/books?id=f1A4AAAAIAAJ>.
- [11] Avrim Blum, Adam Kalai, and Hal Wasserman. Noise-tolerant learning, the parity problem, and the statistical query model. *J. ACM*, 50(4):506–519, July 2003. ISSN 0004-5411. doi: 10.1145/792538.792543. URL <https://doi.org/10.1145/792538.792543>.
- [12] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot

- learners. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfbcb4967418bfb8ac142f64a-Paper.pdf>.
- [13] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde, Jared Kaplan, Harri Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, Will Guss, Alex Nichol, Igor Babuschkin, Suchir Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.
- [14] François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- [15] Konstantina Christakopoulou and Adam Tauman Kalai. Glass-box program synthesis: A machine learning approach. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [16] John Horton Conway. Five \$1,000 problems (update 2017). 2017. URL <https://oeis.org/A248380/a248380.pdf>. [Online; accessed 12/15/2020].
- [17] Valentina Dagienė and Gerald Futschek. Bebras international contest on informatics and computer literacy: Criteria for good tasks. In *International conference on informatics in secondary schools-evolution and perspectives*, pages 19–30. Springer, 2008.
- [18] George B Dantzig. A proof of the equivalence of the programming problem and the game problem. *Activity analysis of production and allocation*, 13:330–338, 1951.
- [19] Eyal Dechter, Jonathan Malmaud, Ryan P Adams, and Joshua B Tenenbaum. Bootstrap learning via modular concept discovery. In *IJCAI*, pages 1302–1309, 2013.
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.
- [21] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. CodeBERT: A pre-trained model for programming and natural languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1536–1547, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.139. URL <https://www.aclweb.org/anthology/2020.findings-emnlp.139>.
- [22] Sumit Gulwani. Automating string processing in spreadsheets using input-output examples. In *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL ’11, page 317–330, New York, NY, USA, 2011. Association for Computing Machinery. ISBN 9781450304900. doi: 10.1145/1926385.1926423. URL <https://doi.org/10.1145/1926385.1926423>.
- [23] Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. Program synthesis. *Foundations and Trends® in Programming Languages*, 4(1-2):1–119, 2017.
- [24] G. H. Hardy, E. M. Wright, D. R. Heath-Brown, and Joseph H. Silverman. *An introduction to the theory of numbers*. Oxford University Press, Oxford; New York, 2008. ISBN 9780199219858 0199219850 9780199219865 0199219869.
- [25] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). 2016.

- [26] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with APPS. 2021.
- [27] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. CodeSearchNet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:1909.09436*, 2019.
- [28] Alexander Skidanov Illia Polosukhin. Neural program search: Solving data processing tasks from description and examples. In *ICLR Workshop Acceptance Decision*, 2018. URL <https://openreview.net/forum?id=B1KJJf-R->.
- [29] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. 2015.
- [30] Burt Kaliski. *RSA factoring challenge*, pages 531–532. Springer US, Boston, MA, 2005. ISBN 978-0-387-23483-0. doi: 10.1007/0-387-23483-7_362. URL https://doi.org/10.1007/0-387-23483-7_362.
- [31] Aditya Kanade, Petros Maniatis, Gogul Balakrishnan, and Kensen Shi. Learning and evaluating contextual embedding of source code. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 12-18 July 2020*, Proceedings of Machine Learning Research. PMLR, 2020.
- [32] Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. Dynabench: Rethinking benchmarking in NLP. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4110–4124, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.324. URL <https://aclanthology.org/2021.naacl-main.324>.
- [33] Sumith Kulal, Panupong Pasupat, Kartik Chandra, Mina Lee, Oded Padon, Alex Aiken, and Percy S Liang. SPoC: Search-based pseudocode to code. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/7298332f04ac004a0ca44cc69ecf6f6b-Paper.pdf>.
- [34] A. Laaksonen. *Guide to Competitive Programming: Learning and Improving Algorithms Through Contests*. Undergraduate Topics in Computer Science. Springer International Publishing, 2020. ISBN 9783030393571. URL <https://books.google.com/books?id=3JbiDwAAQBAJ>.
- [35] Jeffrey C. Lagarias. The $3x + 1$ problem and its generalizations. *The American Mathematical Monthly*, 92(1):3–23, 1985. ISSN 00029890, 19300972. URL <http://www.jstor.org/stable/2322189>.
- [36] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. 2019.
- [37] Nicholas Lourie, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Unicorn on rainbow: A universal commonsense reasoning model on a new multitask benchmark. *AAAI*, 2021.
- [38] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. Codexglue: A machine learning benchmark dataset for code understanding and generation, 2021.
- [39] C. Mayer. *Python One-Liners: Write Concise, Eloquent Python Like a Professional*. No Starch Press, Incorporated, 2020. ISBN 9781718500501. URL <https://books.google.com/books?id=jVv6DwAAQBAJ>.

- [40] Aditya Krishna Menon, Omer Tamuz, Sumit Gulwani, Butler W Lampson, and Adam Kalai. A machine learning framework for programming by example. In *ICML*, pages 187–195, 2013.
- [41] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [42] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1250. URL <https://aclanthology.org/D19-1250>.
- [43] Chris Piech and Eric Roberts. Karel the robot learns python. <https://compedu.stanford.edu/karel-reader/docs/python/en/intro.html>, 2020. [Online; accessed 07-Jun-2021].
- [44] Oleksandr Polozov and Sumit Gulwani. FlashMeta: A framework for inductive program synthesis. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*, pages 107–126, 2015.
- [45] Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *CoRR*, abs/2009.03393, 2020. URL <https://arxiv.org/abs/2009.03393>.
- [46] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL <https://aclanthology.org/D16-1264>.
- [47] Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming Zhou, Ambrosio Blanco, and Shuai Ma. Codebleu: a method for automatic evaluation of code synthesis, 2020.
- [48] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- [49] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL <https://aclanthology.org/D19-1454>.
- [50] Roei Schuster, Congzheng Song, Eran Tromer, and Vitaly Shmatikov. You autocomplete me: Poisoning vulnerabilities in neural code completion. In *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, August 2021. URL <https://www.usenix.org/conference/usenixsecurity21/presentation/schuster>.
- [51] Tal Schuster, Adam Fisch, and Regina Barzilay. Get your vitamin C! robust fact verification with contrastive evidence. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 624–643, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.52. URL <https://aclanthology.org/2021.naacl-main.52>.
- [52] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017. doi: 10.1038/nature24270. URL <https://doi.org/10.1038/nature24270>.

- [53] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- [54] Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan. IntelliCode Compose: code generation using Transformers. *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, Nov 2020. doi: 10.1145/3368089.3417058. URL <http://dx.doi.org/10.1145/3368089.3417058>.
- [55] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421>.
- [56] Gerald Tesauro. Temporal difference learning and TD-Gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- [57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- [58] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-5446. URL <https://aclanthology.org/W18-5446>.
- [59] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- [60] Maksym Zavershynskiy, Alexander Skidanov, and Illia Polosukhin. NAPS: natural program synthesis dataset. *CoRR*, abs/1807.03168, 2018. URL <http://arxiv.org/abs/1807.03168>.
- [61] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. SWAG: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1009. URL <https://aclanthology.org/D18-1009>.
- [62] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [63] Albert Ziegler. Research recitation: A first look at rote learning in github copilot suggestions., June 2021. <https://docs.github.com/en/github/copilot/research-recitation>, Last accessed on 2021-11-01.

We provide complementary details, analysis and results in the following sections:

| | | |
|---|---|---------|
| A | Example solutions by enumerative models | Page 16 |
| B | Enumerative solvers details | Page 16 |
| C | Language Model solvers details | Page 19 |
| D | NP-completeness | Page 22 |
| E | Open problems | Page 23 |
| F | Comparing puzzles to competitive-programming problems | Page 23 |
| G | User Study Details | Page 24 |
| H | Solution to Tower of Hanoi | Page 25 |

A Example solutions by enumerative models

We provide examples of our examined enumerative solvers to three P3 puzzle in Figure A.1 on page 17 (examples of LM solutions are found in the P3 repository). The solution to the first puzzle is general and will work for any other instance of this problem. For the two other puzzles, the obtained solutions are instance-specific and don't even use the input variables. Yet, it is possible that the logical steps to achieve the answer are implicitly executed by the model. To test this, we evaluate the solvers on other problem instances (i.e., puzzles originated from the same problem).

The solvers' solutions to the first puzzle in Figure A.1 are simpler than the one created by humans (though less efficient in terms of input length). This illustrates another potential use case of AI solvers: debugging puzzles by finding easy solutions.

B Enumerative solvers details

We train our random forest solver with the Python Skickit-learn library [41]. The Transformer model is implemented on top of the Hugging Face repository [59]. We use GPUs for training the Transformer and for querying it for rule probabilities. All other computations are performed with CPUs. Making up to 10^4 solution tries takes only a few seconds to a few tens of seconds, depending on the puzzle and the attempted solutions. Running up to 10^6 solution tries usually takes less than an hour but for some puzzles can take longer. We run the solver in parallel on multiple puzzles to reduce the global computation time.

Solution validation. Given the AST, a solution is generated in the form of a Python program (possibly multiple lines) that is evaluated by the Python interpreter to get an answer that is tested by the puzzle. To address long-running programs and infinite loops, timeout checks are added to the puzzles and to the solution during conversion from AST to Python. Alternatively, the programs could be evaluated in a sandbox as is done in programming competitions and as we did for the LM generators, though a sandbox imposes an additional overhead.

B.1 Vocabulary

We use a grammar for a subset of Python covering the following basic objects: Booleans, unlimited-precision integers, floats, strings, lists, sets, dictionaries, generators, and tuples. Table B.1 summarizes the grammar. These rules occur multiply, for instance the addition rule has instantiations for adding two strings, two integers, an integer and a float, etc., where each Python type corresponds to a non-terminal in our grammar. However, because Python is a duck-typed language, in several cases a variable can be used with multiple different types. To handle such programs, we also have a generic non-terminal which can correspond to any Python object, and this makes our grammar ambiguous. For instance, the program `1+1` can be parsed either as the sum of two integers or as the sum of two Python objects, also using a rule mapping an object to an integer. This latter program is a larger AST and hence will typically have lower probability, hence we have the advantages of types when possible but the flexibility to generate fully duck-typed code. In this manner we are able to parse puzzles from 138 of our 200 problems. We also use this grammar to generate timed and safe Python code. In particular, we inject timing checks into comprehensions and loops, and we also add timing checks to potentially time-consuming operations such as exponentiation or string multiplication. This grammar is available upon request for researchers who wish to use it in further projects.


```

# Sum of digits.
def sat1(x: str, s: int=679):
    return s == sum([int(d) for d in x])

# B. Random forest solution.
def sol(s):
    return ((chr(49))*(COPY(s)))

# B. Transformer solution.
def sol(s):
    return ((COPY(s))*(str(1)))

# Human-written solution.
def sol(s):
    return int(s/9) * '9' + str(s%9)

----

# Line intersection.
def sat2(e: List[int], a: int=2, b: int=-1, c: int=1, d: int=2021):
    x = e[0] / e[1]
    return abs(a * x + b - c * x - d) < 10 ** -5

# B. Random forest and B. Transformer solution (identical).
def sol(a, b, c, d):
    return ([2022, 1, ])

# Human-written solution.
def sol(a, b, c, d):
    return [d - b, a - c]

---

# Find the three slice indices that give the specific target in string s.
def sat3(inds: List[int], s: str="hello world", target: str="do"):
    i, j, k = inds
    return s[i:j:k] == target

# B. Random forest solution.
def sol(s, target):
    return ([12, 5, -(3), ])

# B. Transformer solution.
def sol(s, target):
    return ([11, 1, -(6), ])

# Human-written solution.
def sol(s, target):
    from itertools import product
    for i, j, k in product(range(-len(s) - 1, len(s) + 1), repeat=3):
        try:
            if s[i:j:k] == target:
                return [i, j, k]
        except (IndexError, ValueError):
            pass

```

Figure A.1: Example of three P3 puzzles and the solutions found by our examined solvers. The natural language description of each problem is provided for ease of read, but is hidden to these models. Human-written solutions are provided here for reference, but are also hidden from AI solvers.

Table B.1: The grammar for a subset of Python.

| Rule name | rule | Rule name | rule | Rule name | rule |
|-----------|--------------|-------------------|-------------------------|------------|-------------------------|
| != | (_)!=(_) | [list] | [_] | is not | (_)is not(_) |
| & | (_)&(_) | % | (_)%(_) | issubset | (_).issubset(_) |
| (tuple) | (_, _) | {set} | {_} | issuperset | (_).issuperset(_) |
| (tuple) | (_, _, _) | ^ | (_)^(_) | join | (_).join(_) |
| (tuple) | (_, _, _, _) | abs | abs(_) | len | len(_) |
| * | (_)*(_) | all | all(_) | list | list(_) |
| ** | (_)**(_) | and | (_)and(_) | log | log(_) |
| ** | (_)**(_) | any | any(_) | max | max(_) |
| ** | _**(_) | append | (_).append(_) | min | min(_) |
| *args | *_ | arg | _, _ | not | not (_) |
| *args | *_, **_ | arg | _: _, _ | not in | (_) not in (_) |
| + | (_)+(_) | assert | assert _ | or | (_)or(_) |
| += | (_)+=(_) | assert | assert _, _ | ord | ord(_) |
| += | _ += (_) | bool | bool(_) | range | range(_) |
| +unary | +(_) | chr | chr(_) | range | range(_, _) |
| - | (_)-(_) | cos | cos(_) | range | range(_, _, _) |
| -- | (_)--(_) | count | (_).count(_) | replace | (_).replace(_, _) |
| -unary | -(_) | def | def _(:) : _ | return | return (_) |
| / | (_)/(_) | def_ANY_tuple | (_) | reversed | reversed(_) |
| // | (_)/(_) | default_arg | _: _=, _ | reversed | sorted(_, reverse=True) |
| //= | (_)/(_) | default_arg | _=, _ | round | round(_) |
| :slice | _:_:_ | endswith | (_).endswith(_) | round | round(_, _) |
| < | (_)<(_) | exp | exp(_) | set | set(_) |
| << | (_)<<(_) | f_string | f'_' | sin | sin(_) |
| <= | (_)<=(_) | float | float(_) | sorted | sorted(_) |
| = | (_)=(_) | float-const | -- | split | (_).split(_) |
| == | (_)==(_) | float-const-large | --e_ | split | (_).split() |
| > | (_)>(_) | float-const-tiny | --e- | startswith | (_).startswith(_) |
| >= | (_)>=(_) | for | for (_ in ()): _ | str | str(_) |
| COPY | COPY(_) | for | for (_, _) in ()): _ | str-const | "_" |
| [-1] | (_) [-1] | for_in_if | for _ in (_) if _ | sum | sum(_) |
| [-2] | (_) [-2] | formatted_value | {:_} | tuple | tuple(_) |
| [-3] | (_) [-3] | if | if _: _ | type | type(_) |
| [-4] | (_) [-4] | if | if _: _ else: _ | union | (_).union(_) |
| [0] | (_) [0] | ifExp | (_) if (_) else (_) | zip | zip(_, _) |
| [1] | (_) [1] | in | (_) in (_) | zip | zip(_, _, _) |
| [2] | (_) [2] | index | (_).index(_) | | (_) (_) |
| [3] | (_) [3] | int | int(_) | | |
| [i] | (_) [i] | is | (_)is(_) | | |

B.2 Transformer implementation

We use the RoBERTa-base 12-layers Transformer [36] pretrained on English text and fine-tune it on Python code using the Hugging Face library [59]. For fine-tuning data, we use Python functions with their documentation text from GitHub repositories [27]. In order to better adjust the tokenizer to Python code, we retrain a Byte-level BPE tokenizer on our Python fine-tuning data. We use the same vocabulary size as the original tokenizer and keep the token embeddings of the overlapping ones (39%). For the other tokens, we initialize new token embeddings. Thereafter, we fine-tune RoBERTa with a masked language modeling task for 30 epochs. This model, which we denote by T_P , achieved an impressive 3.3 perplexity score on held-out evaluation data, indicating its success in learning Python’s syntax.

Next, we use T_P to encode dense embeddings $e_r = T_P(r)$ for all the rules r in our vocabulary \mathcal{R} . As input to the Transformer, we use a string representation of the Python operation and types of each rule. For example, $(x)/(y)' '// -> \text{FLOAT} :: (x: \text{INT}, y: \text{FLOAT})$ is used to describe the rule for the `//` operation with an integer and float inputs, resulting in a float. Then, we take e_r as the average across the top-layer embeddings of all tokens.

Finally, we design a neural model on top of T_P to predict $\mathbb{P}(r_j|\phi(f), p, i)$ for each puzzle f where p is the parent rule and i is the child index. The model computes a hidden representation of the puzzle with the parent rule as a concatenation $h = [T_P(f), \mathbf{W}_1 e_p, e_i] \in \mathbb{R}^{d+d_r+d_i}$, where $e_i \in \mathbb{R}^{d_i}$ is a learned embedding for the child rule index, $\mathbf{W}_1 \in \mathbb{R}^{d_r}$ is a learned linear projection, and d is the hidden dimension of T_P . To obtain $\phi(f)$, we duplicate T_P and further fine-tune it with the rest of the solver parameters, while keeping the rule Transformer fixed as T_P . Specifically, we use the [CLS] embedding of the top layer as $\phi(f)$. Fixing the rule encoder prevents overfitting to the rules seen in the puzzle-solution fine-tuning pairs. h is then passed through two non-linear layers, where the first also projects it to \mathbb{R}^{d_r} , with a gelu activation [25] and batch normalization [29] to

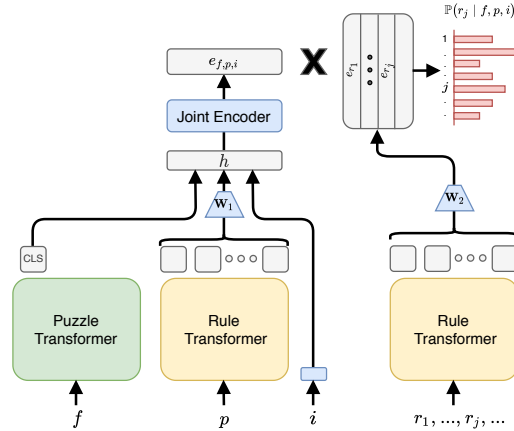


Figure B.1: An illustration of our Transformer-based enumerative solver. The rule strings are encoded with a Transformer pretrained on Python code. The puzzle Transformer is initialized the same, but is further fine-tuned on puzzles, together with the rest of the solver’s parameters shown in blue color. The left hand side of a diagram represents the encoding of the puzzle f , parent rule p , and child index i , each separately and then combined to a joint representation $e_{f,p,i}$. All rules $r \in \mathcal{R}$ are also encoded with the Transformer and projected to the same dimension as $e_{f,p,i}$. The output probability of r being the i ’s child of p in the solution tree g to puzzle f is computed by a softmax over the product of $e_{f,p,i}$ with all rule representations. Encoding the puzzle and the parent rule first separately, allows passing the puzzle only once during inference, and computing all rule embeddings in advance.

```
def f(s: str):
    return s.count("o") == 1000 and s.count("oo") == 0
```

Figure C.1: A puzzle where type-checking is important. A type-safe solution is computed by the program returning "ox" * 1000. However, ["o"] * 1000 would be considered invalid as it is a list of strings, though it does satisfy the puzzle as stated.

get a joint puzzle and parent rule embedding $e_{f,p,i}$. The score of rule r_j then being the i ’s argument of r in the solution to f is determined by the dot product of its projected embedding e_{r_j} with the parent’s embedding: $p_{r_j|\phi(f),e_p,e_i} \propto e_{f,p,i} \cdot (\mathbf{W}_2 e_{r_j})^T$. Similar to the Random Forest fitting process, we use all parent-child rule pairs from the previously obtained solutions for fine-tuning. We use cross-entropy loss with an Adam optimizer. See Figure B.1 for a model diagram.

C Language Model solvers details

The GPT-3 and Codex APIs were used to generate completions based on prompts. For all models, the completions were generated in batches of $n=32$ with $\text{temp}=0.9$, for a maximum of 150 tokens, with default values of $\text{top}_p=1$, $\text{presence_penalty}=0$, $\text{frequency_penalty}=0$, and $\text{best_of}=1$. The resulting programs were evaluated in a sandbox limited to 1 second on Intel Xeon Platinum 8272CL CPUs at 2.60GHz. The timeout was necessary since a number of solution generators would take prohibitive resources such as "a"*(10**(100)) which would generate a string of length googol. The solutions were also checked to be of the type requested in the problem, as was the case for the top-down solver. Figure C.1 illustrates a puzzle where type checking matters.

Prompt programming. The space of possible prompts is practically boundless. Our current prompt designs leverage the API without fine-tuning. For GPT-3, among the prompts we experimented with, we found that the `assert` structure worked best but it was limited to one-line Python solutions. One-line Python programs are considered, by some, to be a useful form of programming with books dedicated to the topic [see, e.g., 39]. For Codex, we found a prompt that resembled a legal python file with a multi-line solution structure worked better.

```

def f1(s: str):
    return "Hello " + s == "Hello world"

assert True == f1("world")

---

def f2(s: str):
    return "Hello " + s[::-1] == "Hello world"

assert True == f2("world"[::-1])

---

def f3(x: List[int]):
    return len(x) == 2 and sum(x) == 3

assert True == f3([1, 2])

---

def f4(s: List[str]):
    return len(set(s)) == 1000 and all(
        (x.count("a") > x.count("b")) and ('b' in x) for x in s)

assert True == f4(["a" * (i + 2) + "b" for i in range(1000)])

---

def f5(n: int):
    return str(n * n).startswith("123456789")

assert True == f5(int(int("123456789" + "0"*9) ** 0.5) + 1)

---

def f6(li: List[int]):
    return len(li) == 10 and li.count(li[3]) == 2

assert True == f6(...)

```

Figure C.2: The medium-length prompt, used for GPT-3. The first five example puzzles f1-f5 were shown to people in the user study and f6 is the one that is being solved. GPT-3's completion was ... [1,2,3,3,4,5,6,7,8,9]

Numerous other prompts were considered. For instance, we tried adding a preface stating, "A Programming Puzzle is a short python function, and the goal is to find an input such that the function True. In other words, if program computes a function f, then the goal is to find x such that f(x)=True."

Interestingly, a handful of generations included potentially dangerous commands such as `eval` and `__import__("os").system(...)`, but a cursory inspection did not detect any that used them in malicious ways. We do advise caution in executing generated code, as malicious actors can take advantage of such routine [50]. There are several libraries available for scoring programming competitions to serve this purpose. Also, some of the generated code seemed especially human-like, e.g.: `raise RuntimeError("this is a joke.")` which of course did not solve the puzzle at hand.

Figures 3, C.2, C.5-C.4, and C.7 show our prompts for the Short, Medium, Long, and Bootstrap prompts, respectively.

```

from typing import List

def f1(s: str):
    return "Hello " + s == "Hello world"

def g1():
    return "world"

assert f1(g1())

def f2(s: str):
    return "Hello " + s[::-1] == "Hello world"

def g2():
    return "world"[::-1]

assert f2(g2())

def f3(x: List[int]):
    return len(x) == 2 and sum(x) == 3

def g3():
    return [1, 2]

assert f3(g3())

def f4(s: List[str]):
    return len(set(s)) == 1000 and all((x.count("a") > x.count("b")) and ('b' in x)
    for x in s)

def g4():
    return ["a"*(i+2)+"b" for i in range(1000)]

assert f4(g4())

def f5(n: int):
    return str(n * n).startswith("123456789")

def g5():
    return int(int("123456789" + "0"*9) ** 0.5) + 1

assert f5(g5())

def f6(inds: List[int], string="Sssuubbstrissiingg"):
    return inds == sorted(inds) and "".join(string[i] for i in inds) == "substring"

def g6(string="Sssuubbstrissiingg"):

```

Codex completed it successfully as:

```

inds = []
ind = 0
for c in "substring":
    while string[ind] != c:
        ind += 1
    inds.append(ind)
    ind += 1
return inds

```

Figure C.3: The medium-length prompt, used for Codex. The first five example puzzles f1-f5 were given in the tutorial to participants in the user study and f6 is the puzzle that is being solved.

Smoothing evaluation. Rather than simply generating solutions until the first correct one is found, to evaluate the Short, Medium and Long prompts, we generate 10,000 solutions for each puzzle. This gives us more than one solution for some puzzles, which we use for improved accuracy in estimating how many solutions are necessary (on average) to solve each puzzle shown in Figure 4b. We use the unbiased estimator of $\text{pass}@k$ defined by Chen et al. [13].

D NP-completeness

Before formally proving that the puzzle decision problem is NP-complete, note that the Boolean Satisfiability problem (SAT) is NP-complete and any Boolean SAT formula such as $(x_0 \vee \neg x_7 \vee x_{17}) \wedge \dots$ can trivially be rewritten as a puzzle, e.g.,

```
def f(x: List[bool]):
    return (x[0] or not x[7] or x[17]) and ...
```

The size of f is linear in the formula size. Thus converting a SAT formula to a puzzle is natural and does not make the problem much bigger or harder.

However, a common misconception is that NP-complete problems are all equally intractable, but the theory of NP-completeness only speaks to the worst-case complexity of solving all puzzles. While any of our puzzles could theoretically be converted to a SAT formula, the resulting formula would be mammoth without any abstraction or intuition. For example, consider the following puzzle,

```
def f(d: int): # find a non-trivial integer factor
    """Hint, try d = 618970019642690137449562111 ;-"""
    n = 100433627766186892221372630609062766858404681029709092356097
    return 0 < d < n and n % d == 0
```

This puzzle is identical to the factoring puzzle f_3 from Figure 1 except that the answer is given away in a comment. Any natural compiler from Python to SAT would ignore comments so the SAT form of this trivial puzzle would be quite hard. While we are not aware of such a compiler, there are programs that convert a factoring problem to a SAT instance. We ran such a converter <http://cgi.cs.indiana.edu/sabry/cnf.html> on this n and it generated a formula with 113,878 variables and 454,633 terms! This illustrates that not all polynomials are small, and that some easy puzzles may become hard puzzles in such a conversion. The theory of NP-completeness only guarantees that if one can efficiently solve *every* SAT instance one could efficiently solve every puzzle, but specific easy puzzles may become quite hard SAT formulas.

D.1 Proof of NP-completeness

Formally, a puzzle f represents a Turing machine as a string, a timeout t is a positive integer represented in unary, and the decision problem is, given (f, x, t) , does there exist y such that when the Turing machine f is run on (y, x) , it halts in fewer than t steps and outputs 1. The time constraint is necessary to ensure that the puzzle decision problem is in NP. It is well-known that this problem is in NP and, moreover is NP-complete:

Observation 1. *The puzzle decision problem is NP-complete.*

Proof. One can test whether a given puzzle string f encoding a Turing machine halts on a witness y in time $\leq t$ by simulating running f on (y, x) for t steps. Since simulating a Turing machine of size $|f|$ running for t steps can be done in $\text{poly}(|f|, t)$ time, this can be done in time $\text{poly}(|f|, t)$ as required for NP.

To see that the problem is complete, note that given any other NP problem defined by a Turing machine $T(x, y)$ that runs on input $x \in \Sigma^*$ and witness $y \in \Sigma^*$ in polynomial time $t = p(|x|)$ is a type of puzzle itself for $f = T$ (with inputs swapped). \square

E Open problems

The following five puzzles would each represent a major breakthrough in computer science or mathematics if solved.

1. **Factoring**. In the traditional version of this ancient problem, the goal is to efficiently find the prime factorization of a given integer. In the puzzle version, we state the equivalent problem of finding any non-trivial factor of a given integer. The puzzle is equivalent in the sense that one can recursively call the puzzle on each of the factors found until one achieves the complete prime factorization. A number of factoring algorithms have been developed over decades that factor larger and larger numbers. The RSA Factoring Challenge [see, e.g., 30] has awarded tens of thousands of dollars in prize money and RSA offered \$200,000 for factoring the largest RSA challenge number with 617 digits. The closely related **Discrete Log** problem is also unsolved.
2. **Graph Isomorphism**. Given two isomorphic graphs, find the bijection that relates the two of them. In a breakthrough, Babai has claimed a quasi-polynomial time for this problem, but no polynomial time algorithm is known.
3. **Planted Clique**. In this classic graph-theory problem, an n -node Erdős–Rényi random graph is chosen and then k nodes are selected at random and the edges are added so that they form a clique. The problem is to find the clique. It is not known whether there is a polynomial-time algorithm for this problem [see, e.g., 4].
4. **Learning Parity with Noise**. This is a binary classification problem in computational learning theory. Roughly speaking, the problem is to efficiently learn a parity function with random classification noise. The fastest known algorithm for this problem runs in time $\tilde{O}(2^{n/\log n})$ [11]. The problem is also closely related to efficiently decoding random linear codes [7] and various assumptions in cryptography. Note that some of the instances of this problem are small (and thus easy) while others are quite large.
5. **Collatz cycle**. The problem is to find a cycle in the famous $3n + 1$ process, where you start with integer $n > 0$ and repeatedly set n to $n/2$ if n is even, otherwise $3n + 1$, until you reach 1. The Collatz cycle conjecture is that there are no cycles in this process. According to the **Wikipedia article** on the topic, Jeffrey Lagarias stated that it “is an extraordinarily difficult problem, completely out of reach of present day mathematics” and Paul Erdős said “Mathematics may not be ready for such problems.” He also offered \$500 for its solution.

Each of these problems is described by a short (1-5 line) python function. Now, for the algorithms problems 1-3, the puzzle involves solving given instances and not exactly with the open problem: coming up with a provably polynomial-time algorithm, and it is entirely possible that no poly-time algorithm exists. However, these are all problems that have been intensely studied and an improvement, even a practical one, would be a breakthrough. For the Collatz cycle, if the Collatz conjecture holds then there is no cycle. However, we give problems involving finding integers with large Collatz delays which could be used to, at least, break records. Also noteworthy but perhaps not as well-known is **Conway’s 99 puzzle**, an unsolved problem in graph theory due to Conway and [10] (as cited by Wikipedia). The two-line puzzle describes finding an undirected graph with 99 vertices, in which each two adjacent vertices have exactly one common neighbor, and in which each two non-adjacent vertices have exactly two common neighbors. Conway [16] offered \$1,000 for its solution.

There are also several unsolved puzzles in terms of beating records, e.g., finding oscillators or spaceships of certain periods in Conway’s game of life and finding uncrossed knights tours on chess boards of various sizes.

F Comparing puzzles to competitive-programming problems

Figure F.1 illustrates an elementary **codeforces.com** problem. As is typical in programming competitions, the authors have concocted an entertaining story to motivate the problem. Dagienė and Futschek [17] include “should be funny” and “should have pictures” among desirable criteria for competitive programming problems. Also, as is typical the first step is explaining how the input

is formatted and how the output should be formatted. One difficulty in authoring such competitive-programming challenges is ensuring that the English description unambiguously matches with the hidden test cases. The [ICPC rules](#) state: “A contestant may submit a claim of ambiguity or error in a problem statement by submitting a clarification request. If the judges agree that an ambiguity or error exists, a clarification will be issued to all contestants.” With puzzles, this is not necessary—a mistake in a puzzle either means that the puzzle is unsolvable or that the puzzle has an unexpected (often trivial) solution, neither of which cause major problems as it would still be a fair comparison of different solvers.

The puzzle form [InvertPermutation](#)¹⁰ has no story, no description of input/output format, and no examples. The input/output formatting is taken care of simply by the type hints.

The intention is for puzzles to isolate the essence of the part of the problem that involves reasoning. Other datasets already address natural language understanding and input/output string formatting.

G User Study Details

The user study began with a short tutorial about puzzles, which included the puzzles shown in [Figure C.2](#). The 30 puzzles (see [Figures G.6-G.7](#)) were divided into three parts of 10 puzzles each: numbers 1-10, 11-20, and 20-30. Since each puzzle took at maximum of 6 minutes, no part took more than one hour. In the internal IRB approval (July 22, 2020), the key discussion points were that we would not collect age, gender or any other PII since it was not relevant to our study.

G.1 Provided instructions

[Figures G.1-G.3](#) present the initial instructions that participants were given before starting the study. [Figures G.4-G.5](#) show the interface that they used for retrieving puzzles and submitting solutions. We run implement a Python backend to store progress logs and to serve each puzzle in its turn, so participants won’t accidentally be exposed to any of the puzzles in advance. We asked participants to follow the simple interface and not to attempt any sophisticated hacking techniques that will give them any personal advantage. We did not observe any such malicious behaviour and received positive feedback for the stability and clarity of the interface.

G.2 Qualitative feedback.

Our Jupyter notebook interface also allowed users to submit qualitative feedback. As an example of this last point, participants mentioned that they were not familiar with functions such as `zip` or `all` but learned them in the course of the study. Overall, Three themes emerged in the feedback: participants enjoyed solving the puzzles, they felt that 6 minutes was not enough time to solve the puzzles, and they felt they learned Python from doing the puzzles.

G.3 Results summary

A total of 21 participants completed the user study. Participants solved between 12-30 puzzles, with 6 participants solving more than 28 puzzles, and only a single participant solving all 30. As [Figure G.8](#) shows, the participants Python experience ranged between a few months to 8 years, with a median of 3 years. For post study analysis purposes, we denote participants with less than 3 years of experience as *beginners* and the rest as *experienced*. [Figure G.9](#) shows the number of participants that solved each puzzle, grouped by experience. 9 of the puzzles were solved by all beginners, whereas 17 puzzles were solved by all experienced. This positive correlation between the number of programming experience and number of puzzles solved, indicates the effectiveness of our puzzles as a proxy to evaluating programming proficiency.

We also notice that experienced programmers solve puzzles faster (149 seconds per puzzle on average, compared to 194 seconds for beginners). [Figure G.10](#) shows the distribution of time spent by participants on each puzzle. We use the per puzzle average solving time as an indicator to its

¹⁰In P3, we have slightly modified the problem so that it is only inspired by the codeforces problem and not a direct translation. The P3 problem is harder in that characters not in the permutation may also appear in the string unmodified.

perceived difficulty. As discussed in the main paper (§5.1), we see a strong correlation between the perceived difficulty of different puzzles for humans and for our examined AI solvers.

H Solution to Tower of Hanoi

Codex’s solution to the Tower of Hanoi puzzle is shown in Figure H.1. Even though the puzzle did not mention the word Hanoi, Codex’s solution clearly knew the reference, in fact offering a link to the Wikipedia page. The first part of the URL is correct, but there is no “Advanced computer algorithm” section on the page, so the link simply resolves to the Wikipedia page. The Python code on the Wikipedia page is only similar in spirit, in part because the way the puzzle asks for the moves is somewhat different from the Wikipedia page. This is a difficult puzzle for which solutions are found at a rate of approximately 0.03%. Surprisingly, Codex was not able to solve the puzzle when we renamed the variable `num_disks` to `n` and removed the string `"bigger disk on top"`, possibly because the association with Tower of Hanoi was weaker.

```

def f1(s: str):
    """Find a string that when concatenated onto 'Hello ' gives 'Hello world'."""
    return "Hello " + s == "Hello world"

assert True == f1("world")

---

def f2(s: str):
    """Find a string that when reversed and concatenated onto 'Hello ' gives 'Hello
world'."""
    return "Hello " + s[::-1] == "Hello world"

assert True == f2("world"[::-1])

---

def f3(x: List[int]):
    """Find a list of two integers whose sum is 3."""
    return len(x) == 2 and sum(x) == 3

assert True == f3([1, 2])

---

def f4(s: List[str]):
    """Find a list of 1000 distinct strings which each have more 'a's than 'b's and
at least one 'b'."""
    return len(set(s)) == 1000 and all(
        (x.count("a") > x.count("b")) and ('b' in x) for x in s)

assert True == f4(["a" * (i + 2) + "b" for i in range(1000)])

---

def f5(n: int):
    """Find an integer whose perfect square begins with 123456789 in its decimal
representation ."""
    return str(n * n).startswith("123456789")

assert True == f5(int(int("123456789" + "0"*9) ** 0.5) + 1)

---

def f6(li: List[int]):
    """Find a list of length 10 where the fourth element occurs exactly twice."""
    return len(li) == 10 and li.count(li[3]) == 2

assert True == f6(...)

```

Figure C.4: An example GPT-3 Long prompt which includes English descriptions in the Python docstrings. As in the medium-length prompts, the first five example puzzles f1-f5 were shown to people in the user study and f6 is the one that is being solved.

```

from typing import List

def f1(s: str):
    return "Hello " + s == "Hello world"

def g1():
    """Find a string that when concatenated onto 'Hello ' gives 'Hello world'."""
    return "world"

assert f1(g1())

def f2(s: str):
    return "Hello " + s[::-1] == "Hello world"

def g2():
    """Find a string that when reversed and concatenated onto 'Hello ' gives 'Hello world'."""
    return "world"[::-1]

assert f2(g2())

def f3(x: List[int]):
    return len(x) == 2 and sum(x) == 3

def g3():
    """Find a list of two integers whose sum is 3."""
    return [1, 2]

assert f3(g3())

def f4(s: List[str]):
    return len(set(s)) == 1000 and all(
        (x.count("a") > x.count("b")) and ('b' in x) for x in s)

def g4():
    """Find a list of 1000 distinct strings which each have more 'a's than 'b's and at least
    one 'b'."""
    return ["a"*(i+2)+"b" for i in range(1000)]

assert f4(g4())

def f5(n: int):
    return str(n * n).startswith("123456789")

def g5():
    """Find an integer whose perfect square begins with 123456789 in its decimal
    representation ."""
    return int(int("123456789" + "0"*9) ** 0.5) + 1

assert f5(g5())

def sat(d: int, n=6002685529):
    return n % d == 0 and all(i in "47" for i in str(d))

def sol(n=6002685529):
    """Find a integer factor of n whose decimal representation consists only of 7's and 4's.
    """

```

Figure C.5: An example Codex long prompt which includes English descriptions in the Python docstrings. As in the medium-length prompts, the first five example puzzles f1-f5 were shown to people in the user study and f6 is the one that is being solved. The representation here supports multi-line solutions. Codex's successful completion is shown in Figure C.6.

```
# return the closest factor to n which is inferior either to n or to the sqrt(n)
if n == 1:      # def.: n>1!
    return 1
i = int(n ** 0.5)
while True:
    if f6(i):
        return i
    elif i == 2:
        raise Exception("Could not find any 7 or 4 number factors
for n.")
    i -= 1
```

Figure C.6: An example Codex completion to the long prompt illustrated in C.5. The solution cleverly calls `f6`, a trick that some participants used to solve a couple of the user study puzzles. Codex also generated English comments.

```

from typing import List

def f1(item: int, li=[17, 2, 3, 9, 11, 11], index=4):
    return li.index(item) == index

def g1(li=[17, 2, 3, 9, 11, 11], index=4):
    return li[index]

assert f1(g1())

def f2(s: str, word="antidisestablishmentarianism", max_len=10):
    if len(word) <= max_len:
        return word == s
    return int(s[1:-1]) == len(word[1:-1]) and word[0] == s[0] and word[-1] == s[-1]

def g2(word="antidisestablishmentarianism", max_len=10):
    if len(word) == max_len:
        return word
    return word[0] + str(len(word[1:-1])) + word[-1]

assert f2(g2())

#
# omitting 28 random puzzles that Codex solved...
#

assert f31(g31())

def f32(x: List[int], a=7, s=5, e=200):
    return x[0] == a and x[-1] <= e and (x[-1] + s > e) and all([x[i] + s == x[i+1]
    for i in range(len(x)-1)])

def g32(a=7, s=5, e=200):

```

```

def f1(s: str, a: List[str]=['cat', 'dot', 'bird'], b: List[str]=['tree', 'fly', '
dot']):
    return s in a and s in b

assert True == f1('dot')

---

def f2(li: List[int]):
    return all([sum(li[:i]) == i for i in range(20)])

assert True == f2(list(map(lambda x: 1, range(100))))

#
# omitting 22 random puzzles that GPT-3 solved...
#

---

def f25(probs: List[float]):
    assert len(probs) == 3 and abs(sum(probs) - 1) < 1e-6
    return max(probs[(i + 2) % 3] - probs[(i + 1) % 3] for i in range(3)) < 1e-6

assert True == f25(

```

Figure C.7: Example bootstrapping prompts for the Codex and GPT-3 LMs. The prompts includes random solved puzzles among those that the LM solved, truncated to the token limit of the API (2048 for GPT3 and 4096 for Codex).

Codeforces problem 474 A. Keyboard

Our good friend Mole is trying to code a big message. He is typing on an unusual keyboard with characters arranged in following way:

```
qwertyuiop  
asdfghjkl;  
zxcvbnm,./
```

Unfortunately Mole is blind, so sometimes it is problem for him to put his hands accurately. He accidentally moved both his hands with one position to the left or to the right. That means that now he presses not a button he wants, but one neighboring button (left or right, as specified in input).

We have a sequence of characters he has typed and we want to find the original message.

Input

First line of the input contains one letter describing direction of shifting ('L' or 'R' respectively for left or right).

Second line contains a sequence of characters written by Mole. The size of this sequence will be no more than 100. Sequence contains only symbols that appear on Mole's keyboard. It doesn't contain spaces as there is no space on Mole's keyboard.

It is guaranteed that even though Mole hands are moved, he is still pressing buttons on keyboard and not hitting outside it.

Output

Print a line that contains the original message.

Examples

input

```
R  
s;;upimrrfod;pbr
```

output

```
allyouneedislove
```

```
def f(s: str, perm="qwertyuiopasdfghjkl;zxcvbnm,./", target="s;;upimrrfod;pbr"):  
    return "".join(perm[perm.index(c) + 1] for c in s) == target
```

Figure F.1: Example of an introductory competition problem <https://codeforces.com/problemset/problem/474/A> (top) and the respective puzzle version (bottom) that is only using code and is short to read. In this problem, there is a given permutation of characters π , and a given target string t , and one wants to find a source string s such that when each character of s has been permuted with π , the target is achieved. The puzzle has been simplified to always shift right.

Getting started

A Python Programming Puzzle is simply a Boolean function `puzzle` and the goal is to find an input which makes `puzzle` return `True`. Let's start with a trivial example:

```
In [1]: def puzzle(s: str):  
        return "Hello " + s == "Hello world"
```

```
In [2]: puzzle("world")
```

```
Out[2]: True
```

That's it, so easy!

```
In [3]: def puzzle(s: str):  
        return "Hello " + s[::-1] == "Hello world"
```

```
In [4]: # What's [::-1]? Check stackoverflow.com or just try it out:
```

```
"Testing"[::-1]
```

```
Out[4]: 'gnitset'
```

```
In [5]: # Aha! Now that you know what [::-1] does, can you solve the puzzle?  
puzzle("your solution here")
```

```
Out[5]: False
```

Figure G.1: Instructions page provided to the study participants as a Jupyter notebook (part 1).

No cheating on types!!!

The `: str` in the arg list above indicates that the input should be a string.

In this Hackathon, we will always define the required type for the input.

For example, in the following puzzle, the input should be a list of integers:

```
In [6]: from typing import List, Tuple, Callable, Set  
  
def puzzle(x: List[int]):  
    return len(x) == 2 and sum(x) == 3
```

Our checker will verify the correctness of the type.

Learning to comprehend it all

For background, a lot of puzzles involve list [comprehensions](#), python functions like `set` and `all`. It will be useful to be familiar with these. For example:

```
In [7]: def puzzle(s: List[str]):  
        """  
        Find 1000 *different* strings where each string has more a's than b's.  
        """  
        return len(set(s)) == 1000 and \  
            all((x.count("a") > x.count("b")) and ('b' in x) for x in s)  
  
        # Example of a valid solution:  
        solution = []  
        for i in range(1000):  
            solution.append('baa' + 'a' * i)  
  
        print(puzzle(solution))  
  
        # You can also use comprehensions in your solution (but you're not required to):  
        solution = ["baa" + "a" * i for i in range(1000)]  
        print(puzzle(solution))  
  
True  
True
```

Figure G.2: Instructions page provided to the study participants as a Jupyter notebook (part 2).

Classic puzzle example

There is no limit to how easy or hard a short puzzle can be.

Many classic puzzles problem can be written as short puzzles. The point of the example below is just to illustrate how one classic programming problem (see the 7,000 words [Wikipedia article](#)) can be written in a couple of lines of Python. (No need to read or understand it)

Don't be scared! We are interested in puzzles of all levels as simple puzzles can be helpful to computers and people just beginning to learn Python. The study puzzles range from easy to medium since each problem is limited to 6 minutes---we do not expect anyone to solve all puzzles but many can be solved faster.

```
In [8]: def puzzle(sol: List[Tuple[int, int]]):
        """
        @param sol: list of moves (i, j) meaning a move from stack i to j (i, j in [0, 1, 2])
        """
        s = (list(range(8)), [], [])
        return all(s[j].append(s[i].pop()) or sorted(s[j])==s[j] for i, j in sol) and s[0] == s[1]

        # The type annotation List[Tuple[int, int]] means that the solution should be a list of pairs of integers like [(0, 2),
        # A great puzzle like this is easy to state but requires a trick (in this case, recursion) to solve.

        def solve_hanoi(n_disks, i=0, j=2):
            if n_disks==0:
                return []
            k = 3 - i - j
            return solve_hanoi(n_disks - 1, i, k) + [(i, j)] + solve_hanoi(n_disks - 1, k, j)

        solution = solve_hanoi(8)
        puzzle(solution)

Out[8]: True
```

Puzzles can also require a function for a solution!

```
In [9]: def puzzle(f: Callable[[int], int]):
        """Find a function f that maps integers to integers where f(0) is nonzero but f(f(0)) is 0."""
        return f(f(0)) == 0 and f(0) != 0
```

Figure G.3: Instructions page provided to the study participants as a Jupyter notebook (part 3).

```
In [ ]: # Please run this cell first.
from study.study import next_puzzle, cur_puzzle, puzzle, give_up, submit_feedback, submit_years

In [ ]: # Please submit the approximate number of years you have been programming in python.
years = _ # integer.
submit_years(years)
```

Instructions

Thank you so much for your participation! Please first complete [Study Consent.ipynb](#).

- The first 3 problems are "practice" and the time you take will not be used in the study. This is a good chance to see how the system works.
- Puzzles are defined by `def puzzle(...)`. For each puzzle, you will try to find an input `x` which makes `puzzle(x)` return `True`.
- Type `next_puzzle()` when you are ready to start the first problem or to advance to the next problem.
- There is **no option to revisit a puzzle** and once you call `next_puzzle()` the clock starts ticking (you have up to 6 minutes per puzzle).
- If you get lost, call `cur_puzzle()`** to see the current puzzle you are on (and time bar).

Timing

- Please solve the problems as quickly as you can. We are measuring the difficulty of the problems both in terms of how many people solve each problem and how long on average it takes them.
- If you do not solve a problem in **6 minutes**, move to the next puzzle by typing `next_puzzle()`.
- If you are sure that you won't be able to solve the puzzle in 6 minutes and would like to skip to the next puzzle without waiting, you can call `give_up()`. However, please avoid this option when possible.

We are evaluating the puzzle's difficulty and not your ability, so do not feel bad about the problems you do not solve. In fact, your not solving a problem is extremely useful information for us. Also, please do not discuss the specific puzzles with people who have not yet completed the study.

Breaks

- Since problems are timed individually, feel free to take breaks between puzzles at your convenience.
- We are storing a state for each user, so you can restart the kernel or close and reopen the notebook if needed. After doing so, please run the top cell to reimport the functions and use `cur_puzzle()` or `next_puzzle()` to get back on track.

Figure G.4: The introduction of the study notebook given to participants.

Summary of functions

| Function | Description |
|----------------------|---|
| next_puzzle() | Start the next puzzle (call only when you are ready to start! no revisiting) |
| cur_puzzle() | Present the current puzzle (useful if you got lost or accidentally overridden puzzle()) |
| puzzle(...) | Submit a solution to the current puzzle |
| give_up() | Give up and skip to the next puzzle before 6 minutes have passed. Please avoid this option if possible. |
| submit_feedback(...) | Send us feedback |

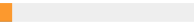
```
In [ ]: # The first 3 puzzles are warmup. Begin the warmup part by running this cell.
next_puzzle()

In [ ]: # Solve the first puzzle by running this cell.
puzzle('world')

In [ ]: # when you are ready to continue, run this cell.
next_puzzle()
```

(a) Initial view.

```
In [4]: # The first 3 puzzles are warmup. Begin the warmup part by running this cell.
next_puzzle()

Time: 


PUZZLE 1/3 (WARM UP)
=====

def puzzle(s: str):
    """
    Warmup problem.
    """
    return "Hello " + s == 'Hello world'
```

```
In [ ]: # Solve the first puzzle by running this cell.
puzzle('world')
```

(b) View while solving a puzzle. The progress bar advances towards the 6 minutes limit.


```
In [5]: # Solve the first puzzle by running this cell.
puzzle('world')
CORRECT in 00:39 sec.



Out[5]: True
```

(c) View after submitting a successful solution.

```
In [6]: # when you are ready to continue, run this cell.
next_puzzle()

Out of time 

PUZZLE 2/3 (WARM UP)
=====

def puzzle(n: int):
    """
    Hint: puzzle(11111111) works.
    """
    return str(n * n).startswith("123456789")
```

(d) View after 6 minutes have passed since viewing the puzzle without submitting a valid solution.

```
In [6]: puzzle(1234)
Out[6]: False
```

(e) View when submitting a wrong solution to a puzzle (before timeout is reached).

Figure G.5: The interface used by participants to solve puzzles during the study. Each sub-figure shows a different state of the notebook according to the user's interaction.

```

def f1(s: str):
    return s.count("o") == 1000 and s.count("oo") == 100 and s.count("ho") == 801

def f2(s: str):
    return s.count("o") == 1000 and s.count("oo") == 0

def f3(x: List[int]):
    return sorted(x) == list(range(999)) and all(x[i] != i for i in range(len(x)))

def f4(x: List[int]):
    return len(x) == 10 and x.count(x[3]) == 2

def f5(x: List[int]):
    return all([x.count(i) == i for i in range(10)])

def f6(n: int):
    return n % 123 == 4 and n > 10**10

def f7(s: str):
    return str(8**2888).count(s) > 8 and len(s) == 3

def f8(s: List[str]):
    return s[1234] in s[1235] and s[1234] != s[1235]

def f9(x: List[int]):
    return ["The quick brown fox jumps over the lazy dog"[i] for i in x] \
        == list("The five boxing wizards jump quickly")

def f10(s: str):
    return s in str(8**1818) and s==s[::-1] and len(s)>11

def f11(x: List[str]):
    return min(x) == max(x) == str(len(x))

def f12(x: List[int]):
    return all(a + b == 9 for a, b in zip([4] + x, x)) and len(x) == 1000

def f13(x: float):
    return str(x - 3.1415).startswith("123.456")

def f14(x: List[int]):
    return all([sum(x[:i]) == i for i in range(20)])

def f15(x: List[int]):
    return all(sum(x[:i]) == 2 ** i - 1 for i in range(20))

```

Figure G.6: The first 15 puzzles in the user study.

```

def f16(x: str):
    return float(x) + len(x) == 4.5

def f17(n: int):
    return len(str(n + 1000)) > len(str(n + 1001))

def f18(x: List[str]):
    return [s + t for s in x for t in x if s!=t] == 'berlin berger linber linger
gerber gerlin'.split()

def f19(x: Set[int]):
    return {i+j for i in x for j in x} == {0, 1, 2, 3, 4, 5, 6, 17, 18, 19, 20, 34}

def f20(x: List[int]):
    return all(b in {a-1, a+1, 3*a} for a, b in zip([0] + x, x + [128]))

def f21(x: List[int]):
    return all([x[i] != x[i + 1] for i in range(10)]) and len(set(x)) == 3

def f22(x: str):
    return x[:2] in x and len(set(x)) == 5

def f23(x: List[str]):
    return tuple(x) in zip('dee', 'doo', 'dah!')

def f24(x: List[int]):
    return x.count(17) == 3 and x.count(3) >= 2

def f25(s: str):
    return sorted(s)==sorted('Permute me true') and s==s[::-1]

def f26(x: List[str]):
    return "".join(x) == str(8**88) and all(len(s)==8 for s in x)

def f27(x: List[int]):
    return x[x[0]] != x[x[1]] and x[x[x[0]]] == x[x[x[1]]]

def f28(x: Set[int]):
    return all(i in range(1000) and abs(i-j) >= 10 for i in x for j in x if i != j) \
        and len(x)==100

def f29(x: Set[int]):
    return all(i in range(1000) and abs(i*i - j*j) >= 10 for i in x for j in x if i
        != j) and len(x) > 995

def f30(x: List[int]):
    return all([123*x[i] % 1000 < 123*x[i+1] % 1000 and x[i] in range(1000)
        for i in range(20)])

```

Figure G.7: The last 15 puzzles in the user study.

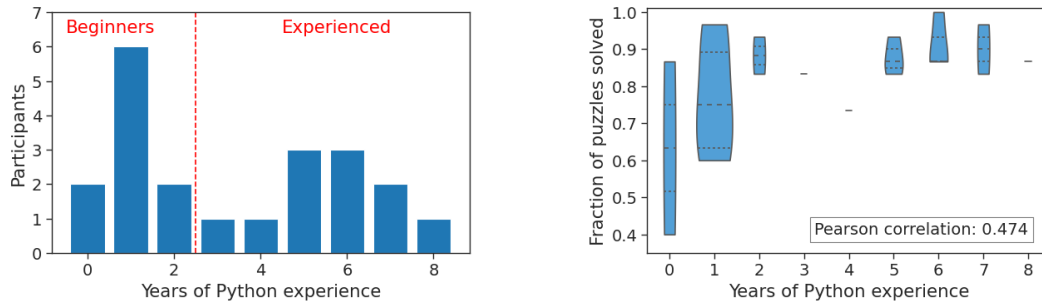


Figure G.8: Years of Python programming experience distribution of our study participants. For post study analysis purposes, we split the group by the median (3 years) to beginners and experienced programmers. The right violin plot shows the fraction of puzzles solved by participants with different years of experience. The lines in the violin show the four quartiles.

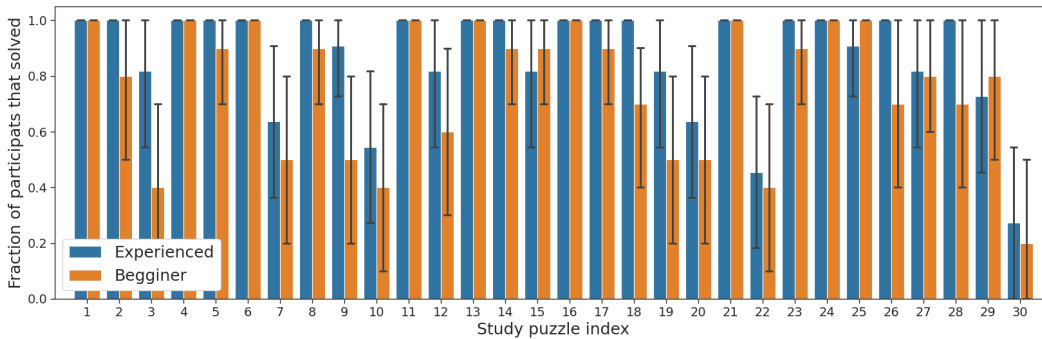


Figure G.9: Fraction of participants, divided to experienced and beginners, that solved each of the 30 puzzles in less than 6 minutes.

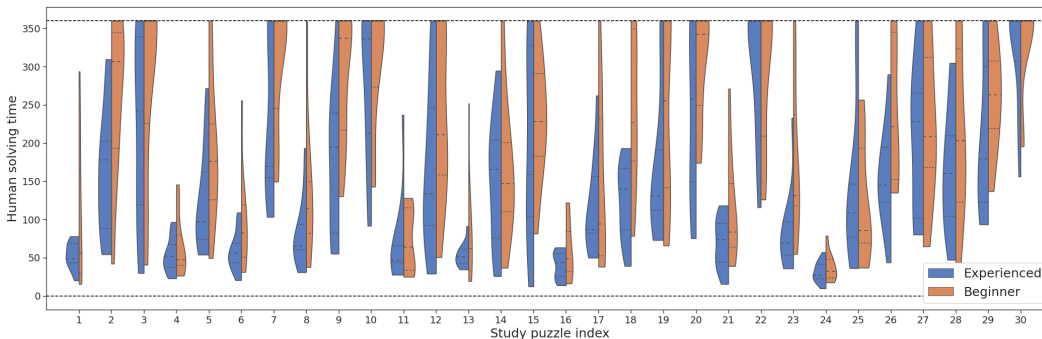


Figure G.10: The time that participants spent on each study puzzle, up to 360 seconds per puzzle. For unsolved puzzles, we count the time as using the full 6 minutes. The orange colored areas on the right show the time distribution for beginner Python coders with less than three years experience (10 out of 21). The blue colored areas on the left present the solving time distribution per puzzle for users with at least three years of experience in Python.

The prompt was the medium Codex prompt (without English descriptions):

```
from typing import List

def f1(s: str):
    return "Hello " + s == "Hello world"

def g1():
    return "world"

assert f1(g1())

def f2(s: str):
    return "Hello " + s[::-1] == "Hello world"

def g2():
    return "world"[::-1]

assert f2(g2())

# skipping two puzzles, see Figure C.3

assert f5(g5())

def f6(moves: List[List[int]], num_disks=8):
    state = [1] * num_disks
    for [i, j] in moves:
        assert state.index(i) <= (state + [1, 2, 3]).index(j), "bigger disk
        on top"
        state[state.index(i)] = j
    return state == [3] * num_disks

def g6(num_disks=8):
```

And a valid completion from Codex, the smaller cushman-codex engine, was:

```
# Algorithm is equivalent to moving all disks.
# From https://en.wikipedia.org/wiki/Tower\_of\_Hanoi#Advanced\_computer\_algorithm
def hanoi(n, p, q, r):
    if n > 0:
        hanoi(n - 1, p, r, q)
        moves.append([p, r])
        hanoi(n - 1, q, p, r)
    moves = []
    hanoi(num_disks, 1, 2, 3)
    assert f6(moves, num_disks)
    return moves
```

Figure H.1: A solution to the Tower of Hanoi puzzle found by the small Codex. The Wikipedia page link generated by the model indeed contains (though not in the *Advanced computer algorithm* section) a solution to a slightly different formulation of this puzzle, see Figure H.2. Note that the medium prompt doesn't mention the name of the puzzle. Codex made the correct association, and adjusted the solution code to the state-based representation of this puzzle as given in f6. Interestingly, replacing the use of *disks* in the puzzle's variable names with other non-descriptive options seems to prevent Codex from solving this puzzle.

```

A = [3, 2, 1]
B = []
C = []

def move(n, source, target, auxiliary):
    if n > 0:
        # Move n - 1 disks from source to auxiliary, so they are out of the way
        move(n - 1, source, auxiliary, target)

        # Move the nth disk from source to target
        target.append(source.pop())

        # Display our progress
        print(A, B, C, '#####', sep='\n')

        # Move the n - 1 disks that we left on auxiliary onto target
        move(n - 1, auxiliary, target, source)

# Initiate call from source A to target C with auxiliary B
move(3, A, C, B)

```

Figure H.2: The algorithm from https://en.wikipedia.org/wiki/Tower_of_Hanoi#Recursive_implementation (November, 2021) that solves Tower of Hanoi for a representation in which the three towers are lists with disk indices.