

NAPS: Natural Program Synthesis Dataset

Maksym Zavershynskyi¹ Alex Skidanov¹ Illia Polosukhin¹

Abstract

We present a program synthesis-oriented dataset consisting of human written problem statements and solutions for these problems. The problem statements were collected via crowdsourcing and the program solutions were extracted from human-written solutions in programming competitions, accompanied by input/output examples. We propose using this dataset for the program synthesis tasks aimed for working with real user-generated data. As a baseline we present few models, with the best model achieving 8.8% accuracy, showcasing both the complexity of the dataset and large room for future research.

1. Introduction

The task of *program synthesis* is to automatically find a program that satisfies user’s specification. It is a problem that has been studied since the earliest days of artificial intelligence (Waldinger & Lee, 1969; Manna & Waldinger, 1975). With the renewed popularity of neural networks for machine learning in recent years, neural approaches to program synthesis have correspondingly attracted greater attention from the research community, which lead to great interest in datasets for program synthesis.

Most of the recent work in the field has been focused on program synthesis from examples for single domain of programming: string transformations (Robust-Fill (Devlin et al., 2017b), Neuro-Symbolic Program Synthesis (Parisotto et al., 2016) and Deep API Programmer (Bhupatiraju et al., 2017)) or Karel (Devlin et al. (2017a), Bunel et al. (2018)). A more domain agnostic dataset was presented in DeepCoder (Balog et al., 2016) but still featured very small programs. All of these results have crucial limitation that datasets were synthetically generated (with exceptions for small private test sets).

¹NEAR. Correspondence to: Maksym Zavershynskyi <max@near.ai>.

Recent examples of crowdsourced natural language to program datasets are WikiSQL (Zhong et al., 2017) and NL2Bash (Lin et al., 2017). Both of these datasets are also domain specific (with WikiSQL featuring only very simple version of SQL) and don’t have programming concepts like variables and control flow. Django dataset (Oda et al., 2015) has very limited scope and each textual description is associated with one line of code.

Worth mentioning fully natural dataset from Magic The Gathering (Ling et al., 2016), that has natural language from people describing actions of cards and Java programs that perform this actions in the Magic environment. This dataset has very limited scope of programs, mostly requiring to figure out complex API of the environment.

Related field to program synthesis from natural language is semantic parsing: mapping of natural language into formal representation, which can be considered as simple programs. Recent examples of such datasets are WebQuestions (Berant et al., 2013), Overnight (Wang et al., 2015), IFTTT (Beltagy & Quirk, 2016). All of these datasets are limited to a specific sub-domain and a limited set of functional intents.

Additionally, there is work on latent program induction which does not require programs as supervision. This simplifies the dataset collection but has a limitation that programs frequently fail to generalize to different inputs (Graves et al., 2014) and does not expose interpretable program back to the user while having huge performance overhead at runtime (Kaiser & Sutskever (2015), Neelakantan et al. (2016)).

In this work we presenting Natural Program Synthesis Dataset v1.0 (NAPS), freely available at <https://near.ai/research/naps/>, consisting of real expert programmers’ solutions for complex problems and rewritten statements in the form that is approachable at the current state of technology. Dataset contains 2231 training and 485 test examples, with additional 16410 unlabelled examples for pretraining and data augmentation.

To assess the difficulty of the NAPS dataset, we implemented sequence-to-sequence and sequence-to-tree baselines. Our best model achieves the accuracy of 8.8%. This shows there is a lot of room for advancement both in mod-

Table 1. NAPS Dataset Structure

Field	Description	Training A	Training B	Test
Solutions	Full programs in UAST format solving a competitive problem	✓	✓	✓
Partial solutions	Smaller pieces of the full programs	✗	✓	✗
IO examples	Input/output examples for the full programs	✓	✓	✓
IO schema	Input/output types and argument names for the full programs	✓	✓	✓
Statements	Crowdsourced problem statements in the imperative format	✗	✓	✓
URLs	URLs to the original problem statements	✓	✓	✓

Table 2. NAPS Dataset Metrics

Metric	Training A	Training B	Test
Number of examples in the dataset	16410	2231	485
Number of examples that are partial solutions	—	1649	—
Number of synthetic statements per solution	300	—	—
Statements length, i.e. number of tokens	173 ± 113 (synthetic)	93 ± 51 (real)	
Number of lines of code per solution		21.7 ± 6.4	
Number of inputs/outputs per solution		7.5 ± 2	

eling and in data augmentation on the NAPS dataset.

2. Dataset

The first release of the NAPS dataset is split into three portions. The largest dataset contains 16410 competitive problem solutions with the corresponding input/output examples and URL links to the original problem statements from the codeforces.com website from which the problem statements can be retrieved. We also accompany each solution with 300 synthetic problem statements that we used for training the baseline models, see Section 3.

The second dataset contains 2231 competitive programming solutions together with the partial exerts from problem solutions. Each record in this dataset is accompanied with a problem statement that was collected by the means of a crowdsourcing platform, a URL to the original problem statement, and input/output examples for non-partial solutions.

The third, smallest dataset contains 485 full problem solutions also accompanied with the crowdsourced problem statements, URLs, and input/output examples.

Solutions: The solutions presented in this dataset are collected from the programming competitions. We then have converted the code written in Java into our intermediate language, UAST, which additionally allowed us to unify library-specific containers and algorithms. In the future, this method will also allow our models to work with solutions across programming languages such as C++, Python, C#, and Pascal.

Written Statements: We hosted a crowdsourcing platform

with participants from competitive programming community and asked them to describe the problem solution that was presented to them in UAST. The process was moderated and the participants were strongly encouraged to give descriptions that were as high-level as possible while at the same time using the language with the imperative structure of the sentences. To provide a curriculum step for the models trained on this dataset, we also asked the participants to describe the smaller inner blocks of the solutions. The workers were allowed to reuse the language used for the inner blocks when describing the blocks enclosing them, but only if the larger block couldn't be described at a higher abstraction level.

Tests: Each full solution is accompanied with 2-10 inputs/outputs each split into two groups. The first group can be used in search or can be included into the problem specification as part of the model input. The second group can be used for the evaluation at the test time.

2.1. UAST Specification

UAST eliminates the burden of managing a runtime or having a compilation step. The code is convertible back and forth between UAST and Java while preserving the readability and the ability to run the input/output examples. While converting to UAST, we also remove all the file I/O and pass all the input data as arguments to the main function, and make the function return the final output. The execution engine and tools for static and runtime analysis can be found at <https://near.ai/research/naps/>.

The language allows several redundancies that simplify the code analysis and the implementation of the executor and

Table 3. UAST Specification

PROGRAM	::=	{'types': [RECORD...], 'funcs': [FUNC...]}
<i>Optional record '_globals_ declares global variables. Function '_main_' is the entry point and optional function '_globals_.init_ initializes the global variables.</i>		
RECORD	::=	['record', name, {field_name: VAR...}]
<i>The function entries are: the return type, the name, the arguments, the local variables, and the body.</i>		
FUNC	::=	['func' 'ctor', TYPE, name, [VAR...], [VAR...], [STMT...]]
VAR	::=	['var', TYPE, name]
STMT	::=	EXPR IF FOREACH WHILE BREAK CONTINUE RETURN NOOP
EXPR	::=	ASSIGN VAR FIELD CONSTANT INVOKE TERNARY CAST
ASSIGN	::=	['assign', TYPE, LHS, EXPR]
LHS	::=	VAR FIELD INVOKE
IF	::=	['if', TYPE, EXPR, [STMT...], [STMT...]]
FOREACH	::=	['foreach', TYPE, VAR, EXPR, [STMT...]]
WHILE	::=	['while', TYPE, EXPR, [STMT...], [STMT...]]
BREAK	::=	['break', TYPE]
CONTINUE	::=	['continue', TYPE]
RETURN	::=	['return', TYPE, EXPR]
NOOP	::=	['noop']
FIELD	::=	['field', TYPE, EXPR, field_name]
CONSTANT	::=	['val', TYPE, value]
INVOKE	::=	['invoke', TYPE, function_name, [EXPR...]]
TERNARY	::=	['?:', TYPE, EXPR, EXPR, EXPR]
CAST	::=	['cast', TYPE, EXPR]
TYPE	::=	bool char int real TYPE* TYPE% <TYPE TYPE> record_name#
<i>The last four types correspond to an array, a set, a map, and a record type.</i>		

the tools. For instance, each expression has a TYPE as the second entry which eliminates the need for deducing the types. Functions require declaring local variables in advance, see Table 3. We have also introduced FOREACH and TERNARY which can be expressed through other control-flow constructs but their introduction has greatly reduced the size of the code. In addition the language is accompanied with a short library of basic functions like 'map_keys', 'string_find', etc.

3. Experimental Results

In this section, we present some of our results on applying sequence-to-sequence and sequence-to-tree models for synthesizing programs from problem statements. In addition, we present the data-structure that we used to perform the decoding in the sequence-to-tree model.

At training time we construct batches by first choosing between datasets A and B with equal probability and then drawing a random example from the chosen dataset. For sequence-to-tree model we choose dataset A 3.3-times more frequently than B. For the dataset A we generated synthetic problem statements using a rule-based randomized method where the rules were selected to match the stylistics of the crowdsource workers as close as possible. The synthetic statements were regenerated anew at the beginning of each epoch and we include 300 synthetic statements for each solution in the dataset A which corresponds

to the number of epochs we trained our baseline models for. The evaluation was performed on the holdout dataset that did not share solutions with the training datasets.

Our sequence-to-sequence model consists of the text encoder and the program decoder mediated through the standard multiplicative attention mechanism (Luong et al., 2015). The encoder is the bidirectional RNN with GRU cells stacked in two layers (Cho et al., 2014). The decoder is a single RNN with GRU cells augmented with a pointer mechanism (Vinyals et al., 2015). In addition to using the pointer mechanism for copying out-of-vocabulary constants and string literals from problem statements to the synthesized code, we also use it for copying in-vocabulary tokens like arithmetic operations and variable names. For this reason, we preferred the soft-switch design described in See et al. (2017), which is suitable for in-vocabulary copying, over the hard-switch design described in Gülcöhre et al. (2016).

3.1. Sequence to Tree

The sequence-to-tree model shares the same encoder and the attention mechanism with the sequence-to-sequence model but the decoding step accounts for the hierarchical nature of the program which allows us to restrict the output of the decoder based on the language grammar. It is done by first implementing a general purpose persistent tree data-structure (Sarnak & Tarjan, 1986) that allows storing

and extending multiple UASTs simultaneously, similarly to how it is done in [Polosukhin & Skidanov \(2018\)](#). The data-structure and the specific implementation of the decoder then work together where the decoder provides the nodes to extend and the data-structure extends them by forking a new tree and placing it in the priority queue based on the tree’s priority, e.g. the likelihood of the entire tree defined by the logits returned from the decoder.

Each node in UAST has an access to its siblings and the parent. For each tree, we store the global state of the entire tree and for each node, we store two states: for its siblings and for its children. The data-structure then passes these states to the decoder which decides which of the incomplete nodes to extend based on the given states. The data-structure then handles multiple extension options for each node which is used in the search. In this paper, we only provide results for the decoder that always extends the left-most incomplete node based on the global state of the tree. However this design can also easily adopt decoders from other papers, e.g. [Polosukhin & Skidanov \(2018\)](#) and [Parisotto et al. \(2016\)](#).

Dataset B and the Test dataset contain problem statements written by real users which poses a challenge since personal writing style varies a lot even though we tried to incentivize the consistency. The biggest challenge is the variance in the verbosity and the usage of rare words. Rules for the synthetic problem statements attempt to mimic the variance in the style but nevertheless, the resulting model is still very sensitive to verbosity. Specifically, the model learns to assign a higher significance to out-of-vocabulary tokens during training than what is optimal for the test dataset.

For the sequence-to-sequence model, the evaluation was performed using the beam search with the beam size equal 64. For the sequence-to-tree model the queue capacity was 64 and at each step, the decoder would expand the left-most incomplete node with 64 most probable tokens yielding 64 new trees which would utilize the memory saving properties of the persistent trees. In the end, we would search through the resulting 64 programs and pick the one that passed the input/output tests. The accuracy is then measured by counting the synthesized programs that pass all the input/output tests that were not used in the search. We also define 50%accuracy metric which counts the programs that pass at least 50% of the test input/output examples.

Interestingly, even when the model does mistakes during the inference those mistakes might be benign and it will still be passing the tests. For instance, Table 5 shows the inference example for the following problem statement:

You are given a number var0. You have to set var2 to 2. If var0-2 is divisible by 3 you have to set var1 to 1, otherwise you have to set var1 to zero. For each var3 between 1 and

Table 4. Accuracy of vanilla and pointer models with and without out-of-vocabulary copying

MODEL	ACCURACY	50%ACCURACY
VANILLA SEQ2SEQ	0%	0%
SEQ2SEQ WIHOUT OOV	3.5%	5.9%
SEQ2SEQ WITH OOV	4.7%	7%
SEQ2TREE WITHOUT OOV	8.8%	11.2%
SEQ2TREE WITH OOV	7.9%	10.2%

Table 5. Example of the inferred program and the tests

```

int __main__(int var0)
vars: int var1, int var2, int var3
var2 = 2
if (((var0 - 2) % 3) == 0)
    var1 = 1
else
    var1 = 0
var3 = 1
for(; (var3 < var0); var3 = (var3 + 1))
    if (var2 < var0)
        var2 = (var2 + ((var3 * 3) + 2))
    if (((var0 - var2) ≥ 0) & ((var0 - var2) ≤ 0))
        var1 = (var1 + 1)
    else
        if (((var0 - var2) ≥ 0) & (((var0 - var2) % 3) == 0))
            var1 = (var1 + 1)
        else
            break
return var1

```

Search Input	157	1312861	6
Search Output	3	312	0
Test Input	26	152	158
Test Output	2	3	4

*var0-1, if var2 is less than var0 you have to, add var3*3+2 to var2, if var0-var2 is greater than or equal to zero and var0-var2 is divisible by 3 add 1 to var1; otherwise you have to break from the enclosing loop. You have to return var1.*

Note that if $\text{var0}-\text{var2} \geq 0 \& \text{var0}-\text{var2} \leq 0$ then $\text{var0}-\text{var2} \geq 0 \& (\text{var0}-\text{var2}) \% 3 == 0$. Even though the model has inferred a redundant if-clause it did not break the program’s logic.

4. Future Work

NAPS dataset enables the program synthesis research on real-life non-trivial programs and problem statements written in a general-purpose language. The baseline metrics, however, demonstrate a large room for the improvement.

References

- Balog, Matej, Gaunt, Alexander L., Brockschmidt, Marc, Nowozin, Sebastian, and Tarlow, Daniel. Deepcoder: Learning to write programs. *CoRR*, abs/1611.01989, 2016. URL <http://arxiv.org/abs/1611.01989>.
- Beltagy, I. and Quirk, Chris. Improved semantic parsers for if-then statements. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL-16)*, Berlin, Germany, 2016. URL <http://www.cs.utexas.edu/users/ai-lab/pub-view.php?PubID=127577>.
- Berant, Jonathan, Chou, Andrew, Frostig, Roy, and Liang, Percy. Semantic parsing on freebase from question-answer pairs. In *EMNLP*, 2013.
- Bhupatiraju, Surya, Singh, Rishabh, rahman Mohamed, Abdel, and Kohli, Pushmeet. Deep api programmer: Learning to program with apis. *CoRR*, abs/1704.04327, 2017.
- Bunel, Rudy, Hausknecht, Matthew, Devlin, Jacob, Singh, Rishabh, and Kohli, Pushmeet. Leveraging grammar and reinforcement learning for neural program synthesis. *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=H1Xw62kRZ>.
- Cho, Kyunghyun, van Merriënboer, Bart, Gülcöhre, Çağlar, Bougares, Fethi, Schwenk, Holger, and Bengio, Yoshua. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014. URL <http://arxiv.org/abs/1406.1078>.
- Devlin, Jacob, Bunel, Rudy R., Singh, Rishabh, Hausknecht, Matthew, and Kohli, Pushmeet. Neural program meta-induction. In *Advances in Neural Information Processing Systems*, pp. 2077–2085, 2017a.
- Devlin, Jacob, Uesato, Jonathan, Bhupatiraju, Surya, Singh, Rishabh, rahman Mohamed, Abdel, and Kohli, Pushmeet. Robustfill: Neural program learning under noisy i/o. In *ICML*, 2017b.
- Graves, Alex, Wayne, Greg, and Danihelka, Ivo. Neural turing machines. *CoRR*, abs/1410.5401, 2014.
- Gülcöhre, Çağlar, Ahn, Sungjin, Nallapati, Ramesh, Zhou, Bowen, and Bengio, Yoshua. Pointing the unknown words. *CoRR*, abs/1603.08148, 2016. URL <http://arxiv.org/abs/1603.08148>.
- Kaiser, Łukasz and Sutskever, Ilya. Neural gpus learn algorithms. *arXiv preprint arXiv:1511.08228*, 2015.
- Lin, Xi Victoria, Wang, Chenglong, Pang, Deric, Vu, Kevin, Zettlemoyer, Luke, and Ernst, Michael D. Program synthesis from natural language using recurrent neural networks. Technical Report UW-CSE-17-03-01, University of Washington Department of Computer Science and Engineering, Seattle, WA, USA, March 2017.
- Ling, Wang, Grefenstette, Edward, Hermann, Karl Moritz, Kocišký, Tomás, Senior, Andrew, Wang, Fumin, and Blunsom, Phil. Latent predictor networks for code generation. *CoRR*, abs/1603.06744, 2016. URL <http://arxiv.org/abs/1603.06744>.
- Luong, Minh-Thang, Pham, Hieu, and Manning, Christopher D. Effective approaches to attention-based neural machine translation. *CoRR*, abs/1508.04025, 2015. URL <http://arxiv.org/abs/1508.04025>.
- Manna, Zohar and Waldinger, Richard. Knowledge and reasoning in program synthesis. In *Proceedings of the 4th International Joint Conference on Artificial Intelligence - Volume 1*, IJCAI'75, pp. 288–295, San Francisco, CA, USA, 1975. Morgan Kaufmann Publishers Inc. URL <http://dl.acm.org/citation.cfm?id=1624626.1624670>.
- Neelakantan, Arvind, Le, Quoc V, Abadi, Martin, McCallum, Andrew, and Amodei, Dario. Learning a natural language interface with neural programmer. *arXiv preprint arXiv:1611.08945*, 2016.
- Oda, Yusuke, Fudaba, Hiroyuki, Neubig, Graham, Hata, Hideaki, Sakti, Sakriani, Toda, Tomoki, and Nakamura, Satoshi. Learning to generate pseudo-code from source code using statistical machine translation (t). *2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pp. 574–584, 2015.
- Parisotto, Emilio, Mohamed, Abdel-rahman, Singh, Rishabh, Li, Lihong, Zhou, Dengyong, and Kohli, Pushmeet. Neuro-symbolic program synthesis. *CoRR*, abs/1611.01855, 2016. URL <http://arxiv.org/abs/1611.01855>.
- Polosukhin, Illia and Skidanov, Alexander. Neural program search: Solving programming tasks from description and examples. *CoRR*, abs/1802.04335, 2018. URL <http://arxiv.org/abs/1802.04335>.
- Sarnak, Neil and Tarjan, Robert E. Planar point location using persistent search trees. *Commun. ACM*, 29(7):669–679, July 1986. ISSN 0001-0782. doi: 10.1145/6138.6151. URL <http://doi.acm.org/10.1145/6138.6151>.
- See, Abigail, Liu, Peter J., and Manning, Christopher D. Get to the point: Summarization with pointer-generator

networks. *CoRR*, abs/1704.04368, 2017. URL
<http://arxiv.org/abs/1704.04368>.

Vinyals, O., Fortunato, M., and Jaitly, N. Pointer Networks.
ArXiv e-prints, June 2015.

Waldinger, Richard J. and Lee, Richard C. T. Prow: A step toward automatic program writing. In *Proceedings of the 1st International Joint Conference on Artificial Intelligence*, IJCAI'69, pp. 241–252, San Francisco, CA, USA, 1969. Morgan Kaufmann Publishers Inc. URL <http://dl.acm.org/citation.cfm?id=1624562.1624586>.

Wang, Yushi, Berant, Jonathan, and Liang, Percy. Building a semantic parser overnight. In *ACL*, 2015.

Zhong, V., Xiong, C., and Socher, R. Seq2SQL: Generating Structured Queries from Natural Language using Reinforcement Learning. *ArXiv e-prints*, August 2017.

Appendix A. Correctly inferred programs

This section demonstrates the examples of the programs inferred by the model and passing all input/output tests. We also provide the texts describing the problem statements and the original human-written programs. We, however, do not give the input/output tests that were used to do the search and evaluate the inferred program because their size is substantial especially for the problems that involve arrays. Please refer to the linked dataset for input/output examples.

Example 1 Given integers var0 , var1 . Let var2 be the less of var0 and var1 . Assign $(\text{var0} + \text{var1}) / 3$ to var3 . Return the less of var2 and var3 .

Golden code

```
int __main__(int var0, int var1)
vars: int var2, int var3
if (var0 > var1)
    var2 = var1
else
    var2 = var0
var3 = (var0 + var1) / 3
return (var2 > var3)?var3:var2
```

Inferred code

```
int __main__(int var0, int var1)
vars: int var2, int var3, int var4
var2 = min(var0, var1)
var3 = ((var0 + var1) / 3)
return min(var2, var3)
```

Example 2 You are given an array of numbers var2 (indexing is 0-based), an array of numbers var3 (indexing is 0-based) and a number var1 . You have to set var9 to -1000000000. For each position var10 in var2 if $\text{var3}[\text{var10}]$ is greater than var1 you have to store in var9 the maximum between var9 and $\text{var2}[\text{var10}]-\text{var3}[\text{var10}]+\text{var1}$, otherwise you have to store in var9 the maximum between var9 and $\text{var2}[\text{var10}]$. You have to return var9 .

Golden code

```
int func0(int var0, int var1, int* var2, int* var3, int var4, int var5)
vars: int var6, int var7, int var8, int var9, int var10,
     int var11, int var12
var6 = len(var2)
var7 = var1
var9 = -1000000000
var10 = 0
for(; (var10 < var6); var10 = (var10 + 1))
    var11 = var2[(var4 = (var4 + 1) - 1)]
    var12 = var3[(var5 = (var5 + 1) - 1)]
    if (var12 > var7)
        var8 = (var11 - (var12 - var7))
    else
        var8 = var11
    var9 = max(var9, var8)
return var9
int __main__(int var1, int* var2, int* var3)
vars: int var4, int var5
var5 = 0
var4 = 0
return func0(1, var1, var2, var3, var4, var5)
```

Inferred code

```
int func0(int var0, int var1, int* var2, int* var3, int var4)
vars: int var5, int var6, int var7, int var8, int var9, int var10
var5 = len(var2)
var6 = 0
var7 = len(var2)
var8 = 0
var9 = -1000000000
var10 = 0
for(; (var10 < var5); var10 = (var10 + 1))
    if (var3[var10] > var1)
        var9 = max(var9, ((var2[var10] - var3[var10]) + var1))
    else
        var9 = max(var9, var2[var10])
    return var9
int __main__(int var1, int* var2, int* var3)
vars: int var4
var4 = 0
return func0(1, var1, var2, var3, var4)
```

Example 3 You are given an integer var0 . If $\text{var0} \% 10$ is less than 5 then decrement var0 by the value $\text{var0} \% 10$. Otherwise, increment var0 by the value $10 - \text{var0} \% 10$. Return var0 .

Golden code

```
int __main__(int var0)
vars: int var1
var1 = (var0 % 10)
if (var1 < 5)
    var0 = (var0 - var1)
else
    var0 = (var0 + (10 - var1))
return var0
```

Inferred code

```
int __main__(int var0)
vars: int var1
var1 = 0
if ((var0 % 10) < 5)
    var0 = (var0 - (var0 % 10))
else
    var0 = (var0 + (10 - (var0 % 10)))
return var0
```

Example 4 Given integers var0 , var1 . Initialize var2 to 0, var3 to 1. While var3 is not less than 1, var3 is not greater than var0 and var1 is not less than var3 , perform the following operations. Add var3 to var2 . Subtract var3 from var1 . Increase var3 by 1. If var3 is greater than var0 , set var3 to 1. Return var1 .

Golden code

```
int __main__(int var0, int var1)
vars: int var2, int var3
var2 = 0
var3 = 1
for; (((var3 ≥ 1) & (var3 ≤ var0)) & (var1 ≥ var3)); )
    if (var1 ≥ var3)
        var2 = (var2 + var3)
        var1 = (var1 - var3)
        var3 = (var3 + 1)
        if (var3 > var0)
            var3 = 1
        else
            pass
    else
        pass
return var1
```

Inferred code

```
int __main__(int var0, int var1)
vars: int var2, int var3
var2 = 0
var3 = 1
for; (((var3 ≤ var0) & (var1 ≥ var3)) & (var1 ≥ var3)); )
    var2 = (var2 + var3)
    var1 = (var1 - var3)
    var3 = (var3 + 1)
    if (var3 > var0)
        var3 = 1
    else
        pass
return var1
```

Example 5 Given integer var0 . Initialize var2 to 1. While var2 is less than var0 subtract var2 from var0 and increase var2 by 1. Return var0 .

Golden code

```
int __main__(int var0)
vars: int var1, int var2
var1 = 0
var2 = 1
for; True; var2 = (var2 + 1))
    if (var2 ≥ var0)
        break
    else
        pass
    var0 = (var0 - var2)
return var0
```

Inferred code

```
int __main__(int var0)
vars: int var1, int var2
var2 = 1
for; (var2 < var0); )
    var0 = (var0 - var2)
    var2 = (var2 + 1)
return var0
```

Appendix B. Incorrectly inferred programs

Example 1 Given an integer `var0`, find the smallest integer, which is greater than `var0`, such that all its digits are distinct and return the array, which contains one element equal to the result.

This problem provides a great example of how our model attempts to handle more complex concepts, like comparing numbers digit-wise. In this case, it attempts to extract the digit from a number by accessing it as an array, instead of doing it through the arithmetic operations like in the golden code.

Golden code

```
int* __main__(int var0)
vars: int var1, int var2, int var3, int var4, int* var5
var5 = new int*()
var0 = (var0 + 1)
for(; True; )
    var1 = (var0 % 10)
    var2 = ((var0 / 10) % 10)
    var3 = ((var0 / 100) % 10)
    var4 = (var0 / 1000)
    if (((((var1 != var2) & (var1 != var3)) & (var1 != var4))
        & (var2 != var3)) & (var2 != var4)) & (var3 != var4))
        array_push(var5, var0)
        break
    else
        pass
    var0 = (var0 + 1)
return var5
```

Inferred code

```
int* __main__(int var0)
vars: int var1, int var2, int var3, int* var4
var4 = new int*()
var2 = 0
var3 = 0
for(; (var3 < var2); var3 = (var3 + 1))
    var4 = var0[(var2 = (var2 + 1) - 1)]
    if (var2[var3] == 1)
        var2 = (var2 + 1)
    else
        pass
return var4
```

Example 2 Given integer `var0`. Return the sum of powers of 2 from 1 to `var0`.

For this problem, the inferred program uses incorrect variable names in the iterator increment and in the accumulator. Besides that, the program structure is correct.

Golden code

```
int __main__(int var0)
vars: int var1, int var2
var1 = 0
var2 = 1
for(; (var2 ≤ var0); var2 = (var2 + 1))
    var1 = (var1 + pow(2, var2))
return var1
```

Inferred code

```
int __main__(int var0)
vars: int var1, int var2, int var3
var1 = 1
var2 = 1
for(; (var2 ≤ var0); var3 = (var3 + 1))
    var3 = (var3 + (2 << var2))
return var2
```

Example 3 You are given integers `var0`, `var1`. Return the factorial of the minimum of `var1` and `var0`.

This is an example of a completely incorrect inference.

Golden code

```
int __main__(int var0, int var1)
vars: int var2, int var3, int var4
if (var0 > var1)
    var3 = var1
else
    var3 = var0
var4 = var3
for; (var4 > 1); var4 = (var4 - 1))
    var3 = (var3 * (var4 - 1))
return var3
```

Inferred code

```
int __main__(int var0, int var1)
vars: int var2
return (min(var1, var0) * 2)
```

Appendix C. Decoding algorithm

Here we give a high-level description of how the data-structure and the decoder perform the decoding in the sequence-to-tree model.

```

 $Q \leftarrow \{create\_empty\_tree()\}$ 
while has_incomplete_trees( $Q$ ) do
    for all  $T$  in incomplete_trees( $Q$ ) do
        candidates  $\leftarrow \emptyset$ 
        for all  $N$  in incomplete_nodes( $T$ ) do
             $h_{tree} \leftarrow T.h$ 
             $h_{parent} \leftarrow parent(N).h_{parent}$ 
             $h_{left} \leftarrow extract\_states(left\_siblings(N))$ 
             $h_{right} \leftarrow extract\_states(right\_siblings(N))$ 
            candidates.add(( $h_{tree}, h_{parent}, h_{left}, h_{right}$ ))
        end for
        new_nodes  $\leftarrow \text{DECODER}(candidates)$ 
        for all  $N$  in new_nodes do
            if  $N$  satisfies language grammar then
                Fork a new tree  $T_{new}$  by extending  $N$ 
                 $Q.push(T_{new})$ 
            end if
        end for
    end for
end while
```
