# Applying Genetic Programming to PSB2: The Next Generation Program Synthesis Benchmark Suite

Thomas Helmuth<sup>1\*</sup> and Peter Kelly<sup>1</sup>

<sup>1</sup>Computer Science, Hamilton College, Clinton, New York, USA.

\*Corresponding author(s). E-mail(s): thelmuth@hamilton.edu; Contributing authors: pxkelly@hamilton.edu;

#### Abstract

For the past seven years, researchers in genetic programming and other program synthesis disciplines have used the General Program Synthesis Benchmark Suite (PSB1) to benchmark many aspects of systems that conduct programming by example, where the specifications of the desired program are given as input/output pairs. PSB1 has been used to make notable progress toward the goal of general program synthesis: automatically creating the types of software that human programmers code. Many of the systems that have attempted the problems in PSB1 have used it to demonstrate performance improvements granted through new techniques. Over time, the suite has gradually become outdated, hindering the accurate measurement of further improvements. The field needs a new set of more difficult benchmark problems to move beyond what was previously possible and ensure that systems do not overfit to one benchmark suite.

In this paper, we describe the 25 new general program synthesis benchmark problems that make up PSB2, a new benchmark suite. These problems are curated from a variety of sources, including programming katas and college courses. We selected these problems to be more difficult than those in the original suite, and give results using PushGP showing this increase in difficulty. We additionally give an example of benchmarking using a state-of-the-art parent selection method, showing improved performance on PSB2 while still leaving plenty of room for improvement. These new problems will help guide program synthesis research for years to come.

 $045 \\ 046$ 

007

 $\begin{array}{r}
 008 \\
 009
 \end{array}$ 

 $\begin{array}{r}
 010 \\
 011 \\
 012
 \end{array}$ 

 $\begin{array}{c} 013\\ 014 \end{array}$ 

 $015 \\ 016 \\ 017$ 

018

019 020 021

 $\begin{array}{r}
 022 \\
 023
 \end{array}$ 

024

025

026

027

028

029

030

031

032

033

034

 $\begin{array}{c} 035\\ 036 \end{array}$ 

037

038

039

040

041

042

043

047 **Keywords:** automatic program synthesis, benchmarking, genetic 048 programming

- 049
- 050
- 051

#### 052 **1 Introduction** 053

Automatic general program synthesis, with the aim of automatically gener-054 ating programs of the type humans write from scratch, has long been a goal 055of artificial intelligence and machine learning. Yet, for many years there were 056no common benchmark problems for evaluating general program synthesis<sup>1</sup> 057 systems; existing problems were either easy toy problems or were situated in 058specific domains where solution programs were composed of a small set of 059 domain-specific instructions. In 2015, the General Program Synthesis Bench-060 mark Suite (PSB1) introduced 29 programming by example problems that can 061 be used to benchmark program synthesis systems that take their specifica-062tions as input/output examples [1]. Since then, more than 100 research papers 063 have benchmarked 10+ program synthesis systems using PSB1, producing 064 numerous insights into program synthesis. 065

Of the systems that have adopted PSB1, most fall within the field of 066 genetic programming (GP), including PushGP [1], grammar-guided GP [2], 067 grammatical evolution [3], and linear GP [4]. However, non-evolutionary pro-068 gram synthesis methods have also been applied to PSB1, including those based 069 on delayed-acceptance hillclimbing [5] and Monte Carlo tree search [6]. We 070 expand on the details of these methods and the results they have achieved using 071PSB1 in Section 2, but to summarize, many of these systems have improved 072 performance and demonstrated new techniques. 073

When PSB1 was first introduced, the initial PushGP benchmarking runs 074 were able to solve 22 of the 29 problems, with an average success rate of 23 075 successful runs out of 100 [1]. The best-performing PushGP results have now 076 solved 25 problems, with an average success rate of 42/100 [7, 8]. Some of 077 the most drastic improvements have come on some of the most informative 078 problems in PSB1, such as Double Letters  $(6 \rightarrow 50 \text{ successes between } [1] \text{ and }$ 079 [7]), Replace Space with Newline  $(51 \rightarrow 100)$ , Syllables  $(18 \rightarrow 64)$ , Vector 080 Average  $(16 \rightarrow 97)$ , and X-Word Lines  $(8 \rightarrow 91)$ . 081

Thus, for PushGP and other synthesis systems, the problems of PSB1 have 082 become less useful over time. In particular, the very high performance achieved 083 on some PSB1 problems leaves little room for exhibiting improvement; a few 084 other problems have never been solved and are likely too difficult to be solved 085any time soon. Additionally, peculiarities in some of the problems in PSB1 086 make them less ideal as benchmarks, either because of how synthesis systems 087 move through their search space or how slow they are to run. Finally, some 088 decisions about the specification of problems in PSB1 make them difficult to 089 implement, potentially preventing wider adoption. 090

<sup>092</sup> <sup>1</sup>Also known as automatic programming or software synthesis.

With these drawbacks in mind, we have created a second Program Syn-093 thesis Benchmark suite, which we refer to as PSB2. PSB2 consists of 25 094problems curated from programming challenges, programming katas, and col-095 lege courses. In order to facilitate the implementation of PSB2 in new synthesis 096 systems, we provide a reference implementation of each problem, datasets that 097 can be sampled for each problem, and libraries for downloading and sampling 098 the datasets.<sup>2</sup> Just like PSB1, the problems in PSB2 require a wide range 099 of programming techniques, data types, and control flow structures to solve. 100 However, they are markedly harder to solve than problems in PSB1, with our 101 initial results solving 13 of the 25 problems for an average success rate of 10 102successful runs out of 100. These more difficult problems will drive program 103synthesis research toward solving more realistic program synthesis tasks. 104

The purpose of benchmark problems is to allow us to empirically show 105what changes to a system produce improvements that may transfer to real-106world problems. To achieve this goal, they must be sufficiently difficult, unlike 107 toy problems that have been used as benchmarks in the past. They must 108 also be representative of the types of tasks we want our system to perform. 109However, we also want benchmarks to be easier and faster to run than an actual 110 real-world problem in order to produce benchmarking results with reasonable 111 amounts of computation. Given that automatic program synthesis is still in 112its fledgling stages, we see the problems in PSB2 as a stepping stone toward 113solving more realistic problems. 114

PSB2 also addresses calls from the GP community to produce and adopt 115realistic benchmarks. GP community discussions calling for better bench-116 marks [9–11] inspired the creation of PSB1: these calls also highlighted the 117 need to periodically update and replace benchmark problems in order to 118 keep advancing the field without over-optimizing to a single set of problems. 119More recently, a call to refocus the efforts of GP on automatic programming 120 stated, "We are in no doubt of the need for the further principled develop-121ment of additional benchmarks that can be used in a targeted manner to push 122the boundaries along different dimensions such as scalability, generalisation, 123and adaptation, and to facilitate comparison across a range of very different 124approaches to automatic programming" [12]. The creation of PSB2 aims to 125push the boundaries of program synthesis research and give synthesis systems 126a fresh set of problems to explore. Of course, there is no need to entirely throw 127out the problems of PSB1; we could imagine some of the harder problems con-128tinuing to provide useful data, and newer systems may need to start on the 129easier problems as a jumping off point. 130

PSB2 was originally introduced at the 2021 Genetic and Evolutionary Computation Conference, where it was nominated for Best Paper in the Genetic 132 Programming track [13]. This article expands on that work by extending the 133 discussion of best practices for benchmarking program synthesis using GP, presenting new experiments, and adding a discussion of related program synthesis 135

 $<sup>^{2}</sup>$ Reference implementation, datasets, and other resources can be found on this paper's 137 companion website: https://cs.hamilton.edu/~thelmuth/PSB2/PSB2.html. 138

techniques. In the new experiments, we use PushGP with down-sampled lexi-case selection [7, 8], improving the number of problems solved from 13 to 17,and additionally conduct a control study using random search to ensure thatprogram synthesis is necessary to solve these problems.

143The remainder of this paper is structured as follows: in the next section, 144 we discuss research that has used PSB1. In Section 3, we highlight lessons learned about program synthesis benchmarking from PSB1. Sections 4 and 5 145describe the sources of PSB2's problems and describe the problems in detail. 146 147We then give general guidance on benchmarking with PSB2, and give details of the parameters we used in our experiments in Sections 6, 7 and 8. Sections 9, 14814910, and 11 present initial results using PushGP, a control study using ran-150dom search, and a benchmarking study showing improved performance using down-sampled lexicase selection. Finally, Section 12 situates PSB2 in the larger 151152field of program synthesis, discussing how specifications are defined in various 153subfields.

154

# $^{155}_{156}$ 2 Past Research Using PSB1

PSB1 has been used in a variety of research projects on automatic program
synthesis, many of them using GP as the synthesis system. A recent survey
gives details of the methods that have been used to solve the problems in
PSB1, as well as analyzing the difficulties of the problems [14].

161 The paper that introduced PSB1 [1] used PushGP, a GP system based on 162 the stack-based Push programming language; a variety of papers using PushGP 163 have made use of PSB1 since [7, 8, 15–20]. Code-building GP is a stack-based 164 GP system borrowing some inspirations from Push that constructs programs 165 in a host language; it solved some of the PSB1 problems, producing solution 166 programs Synthesized in Python [21].

167General program synthesis requires the manipulation of multiple data 168types; stack-based GP systems have handled this requirement well, but so 169have other GP systems that handle strong typing of programs. In particular, 170grammar-based approaches such as grammar guided GP (G3P) [2, 22–24] and 171grammatical evolution (GE) [3, 25-27] have made good progress at solving the 172problems in PSB1. Many of these use the type-based grammar design patterns 173introduced to flexibly handle problems with different type requirements [2]. 174Another use of these grammars trains a sequence-to-sequence variational 175autoencoder to embed programs in a continuous space and then uses an evolu-176tionary algorithm to optimize programs in this space [28]. Finally, a linear GP 177system with tag-based memory has also been explored using PSB1 [4, 29, 30].

As for non-GP systems, an approach using delayed-acceptance hillclimbing for inductive synthesis proved competitive with GP on PSB1, including producing the only known solutions to the Collatz Numbers problem [5]. A comparison was made between Flash Fill [31], MagicHaskeller [32], PushGP, and G3P, finding that the non-GP methods fared much worse but ran much 183

faster than the GP methods [33]. Finally, Monte Carlo tree search was used to generate Java bytecode programs using a few of the problems in PSB1 [6]. 186

#### Lessons Learned 3

While PSB1 has been successfully used in a variety of research, it was a first 190191attempt at a general-purpose program synthesis benchmark suite. The research community has grown from using it, both in terms of improving program 192193synthesis methods as well as learning lessons about how to best define program 194synthesis benchmarks. Here, we discuss some of those lessons and how they have influenced our creation of PSB2. 195

196One major issue with PSB1 is that every system that uses it needs to 197implement all of the problems from scratch. This hurdle likely decreased wider 198 adoption. Additionally, there may be inconsistencies between implementations 199in different systems, leading to less comparable results; one known such inconsistency is that some systems use new randomized data for each run, while 200201others use the same dataset for every run. Four years after its initial release, 202the authors of PSB1 created large datasets of the inputs and correct outputs 203for each problem [34]. These datasets can be sampled for each program synthe-204sis run, meaning there is no need for each system to implement each problem. We have copied this model and provide datasets for each problem in PSB2, as 205206well as software libraries to make them easier to use (see Section 6).

207A handful of the problems in PSB1 require programs to produce Boolean 208outputs, as such functions are common in programming exercises. A trend 209noted across completely different program representations is that solutions to these Boolean-output problems often do not generalize to unseen data [1, 21021122]. A simple explanation for this phenomenon is that it is relatively easy 212for a solution program to produce the correct answers for the wrong reasons 213when there are only two possible answers, thus overfitting to the training data. 214It is much harder to perfectly answer training data for the wrong reasons when the output is an integer or string, for example. Because of this issue, 215216we have selected fewer Boolean-output problems for PSB2, including only one 217representative problem.

218PSB1 was designed to emulate the textbook problems it was curated from 219as closely as possible. For example, many problems from the original textbook required the program to "print" its answers. PSB1 suggested that synthesis 220221systems develop methods for emulating an output buffer and common printing 222instructions in order to mimic these problems. However, this approach was 223infeasible for some synthesis systems, which instead simply returned string 224outputs. As PSB2 is less loosely coupled with its problem sources, we decided to have all programs return their outputs instead of "printing" them. Another 225226wrinkle related to outputs is that some problems require a solution to return 227multiple outputs. While multiple outputs may prove difficult in some systems, 228it is generally feasible in all; we have included 4 multi-output problems in 229PSB2.

185

5

187 188 189

PSB1 recommended different training and test set sizes, as well as program evaluation budgets, for each problem. This led to difficulties and confusion in both implementation and reporting results. For PSB2, we recommend using a fixed setting for each parameter across problems, as well as for system-specific parameters, such as maximum sizes for program size control.

236Forstenlechner et al. [22] discussed understanding and refining the prob-237lems in PSB1, making some general recommendations about both synthesis 238systems and benchmark problems. One suggestion put forth is using larger and 239more targeted training sets, to better guide synthesis and increase generaliza-240tion. We take this recommendation and use large training sets (200 examples) 241that have a variety of specific edge cases purposefully included. Most of their 242other suggestions relate to specific system settings, such as the length of evo-243lution; these parameters are not prescribed by PSB2, and can be chosen by the researcher. 244

245Our approach to progressing the field of general program synthesis is an incremental one by design: researchers have improved synthesis systems on the 246247problems of PSB1, so now we wish to find improvements on the moderately more difficult problems of PSB2. However, we could imagine other approaches 248249that could be used to progress the field besides an incremental approach, and 250other researchers may choose to use these instead. For example, we could 251create significantly more challenging problems, beyond the capabilities of any 252current synthesis system. Then, we could work on methods to solve those 253problems directly, with the hope that working on a more difficult task will 254allow for the creation of new techniques that may not be obvious with a more 255incremental approach. The drawback here is that until some method actually solves the problem, there is not much gradient for the research community to 256257follow in order to learn which methods are more promising. Thus we prefer an 258incremental approach, with benchmarking and comparison of existing methods using moderately more difficult problems. 259

260

# <sup>261</sup> 4 Problem Selection and Sources

263 Below we describe the four sources we used as inspiration for the problems 264 included in this suite. Each of these sources contains problems for humans to 265 use to improve their programming skills, whether for experienced programmers 266 or students in class. As such, these sources contain problems representative of 267 the types of programming that we expect humans to perform.

268 **Code Wars (CW)** - A website full of user-created programming chal-269 lenges, called *coding kata*. The aim of the site is for users to spend small 270 amounts of time programming every day to hone their coding skills.

Advent of Code (AoC) - An Advent calendar of coding problems created
every year in December. These problems can be used for any number of things,
like training, interview prep, or coursework. Problems tend to become harder
throughout the month.

- 275
- 276

Homework Problems (HW) - These problems come from programming277homework given in our undergraduate programming courses. The problems278come from two courses: an introductory programming course, and a program279languages course. These problems do not have citations, since we created them280for our courses.281

Project Euler (PE) - A website containing hundreds of problems in an282archive. Users are free to submit answers to validate their solutions. Most283problems tend to be mathematically focused and often require efficient and284elegant solutions.285

These sources contain a large number and variety of problems; we consid-286ered over 75 problems from these sources and implemented and tested over 28750 of them, filtering out problems that seemed too easy, too difficult, or inac-288cessible. While curating the suite, we did not include any problems on which 289standard PushGP produced a success rate over 60% in initial experiments, to 290ensure that problems are sufficiently difficult to allow for improvement. While 291other GP systems surely will find some of these problems easier or harder than 292PushGP, as they have with PSB1, we hope a combination of subjective cura-293tion and objective results using PushGP will give a good balance of problems 294that are moderately more difficult than the problems in PSB1 for any system. 295In order to be transparent about our curation process, we have created a table 296containing all of the problems we considered, including the reason for accept-297ing/rejecting each problem, initial results if we implemented the problem, and 298a link to the source of the problem.<sup>3</sup> 299

We aimed to include problems that require a large variety of data types 300 and control flow structures to solve, with a balance between data types across 301 problems. Most of the problems require some type of iteration and/or conditional statements. Required data types include integers, floats, Booleans, 303 characters, strings, vectors of integers, and vectors of floats. In order to produce large datasets, we aimed to select problems that have at least 1 million possible unique inputs. 306

307

 $308 \\ 309$ 

320

## 5 Problem Descriptions

Below is a list of the English language descriptions of the 25 benchmark prob-310lems in PSB2. Each problem (besides those from our courses) has a citation of 311 its source with a link to the original problem. The types of the input(s) and 312output(s) for each problem are given in Table 1. Because the type(s) of control 313 flow that are necessary to solve a problem can be a large influence on how dif-314 ficult it is, we give our best prediction as to which types of control flow will be 315necessary to solve each problem in Table 2. Note that all problems that require 316vector/string processing (considering individual elements or characters) must 317 also require iteration over those elements/characters, but we include it as a 318separate category to distinguish it from problems that simply need iteration 319

 $<sup>\</sup>label{eq:shttps://docs.google.com/spreadsheets/d/e/2PACX-1vQK01D2sZA9KosXpOJNui} 321 \\ DW6yDQEZnMrwzNeMJJU25MbZhU6odQ0jGkJN5lgbRgspsmmum65WLbEI2B/pubhtml?gid= 322 \\ 0\&single=true$ 

$\mathbf{Problem}$	Inputs	Outpu
Basement	vector of integers of length $[1, 20]$ with each integer in $[-100, 100]$	integer
Bouncing B.	float in $[1.0, 100.0]$ , float in $[1.0, 100.0]$ , integer in $[1, 20]$	float
Bowling	string in form of completed bowling card, with one character per roll	integer
Camel Case	string of length [1, 20]	string
Coin Sums	integer in [1, 10000]	4 integ
Cut Vector	vector of integers of length $[1, 20]$ with each integer in $[1, 10000]$	2 vect of integ
Dice Game	2  integers in  [1, 1000]	float
Find Pair	vector of integers of length $[2, 20]$ with each integer in $[-10000, 10000]$ , integer in $[-20000, 20000]$	2 integ
Fizz Buzz	integer in [1, 1000000]	string
Fuel Cost	vector of integers of length $[1, 20]$ with each integer in $[6, 100000]$	integer
GCD	2 integers in [1, 1000000]	integer
Ind. of Subs.	2  strings of length  [1, 20]	vector integers
Leaders	vector of integers of length $[0, 20]$ with each integer in $[0, 1000]$	vector
Luhn	vector of integers of length 16 with each integer in $[1, 9]$	integer
Mastermind	2 strings of length 4 made of B, R, W, Y, O, G	2 integ
Middle Char.	string of length [1, 100]	string
Paired Digits	string of digits of length $[2, 20]$	integer
Shopping List	vector of floats of length $[1, 20]$ with each float in $[0.0, 50.0]$ , vector of floats of length $[1, 20]$ with each float in $[0.0, 100.0]$ . Both vectors must be the same length	float
Snow Day	integer in $[0, 20]$ , float in $[0.0, 20.0]$ , float in $[0.0, 10.0]$ , float in $[0.0, 1.0]$	float
Solve Boolean	<pre>string of length [1,20] made of characters from {t, f,  , &amp;}</pre>	Boolean
Spin Words	string of length [0, 20]	string
Square Digits	integer in [0, 1000000]	string
Subs. Cipher	3 strings of length $[0, 26]$	string
Twitter	string of length [0, 200]	string
Vector Dist.	2 vectors of floats of length $[1, 20]$ with each float in	float
	[-100.0, 100.0]	

**Table 1** For each problem, the types of the inputs and outputs, and the limits imposed on the inputs.

without vector/string processing. For more precise details of each problem, see the reference implementation.<sup>4</sup>

- Basement (AoC) Given a vector of integers, return the first index such that the sum of all integers from the start of the vector to that index (inclusive) is negative. [35]
- 363
  364
  365
  366
  2. Bouncing Balls (CW) Given a starting height and a height after the first bounce of a dropped ball, calculate the *bounciness index* (height of first bounce / starting height). Then, given a number of bounces, use
- 368 <sup>4</sup>https://github.com/thelmuth/Clojush/releases/tag/psb2-v1.0
- 367

**Table 2** For each problem, the types of control flow that a program likely needs to include to solve the problem. **Iteration** refers to some sort of iteration, looping, or recursion. **Vector** includes an type of vector/string processing, considering individual elements or characters. **Conditional** refers to conditional statements of some kind, whether through standard *if* statements or other means.

Problem	Iteration	Vector	Conditional
Basement	v	v	v
Bouncing B	v	л	А
Bowling D.	x	v	v
Camel Case	x	x	x
Coin Sums	A	A	7
Cut Vector	х	х	х
Dice Game			
Find Pair	х	х	х
Fizz Buzz			х
Fuel Cost	х	х	
GCD	х		х
Ind. of Subs.	х	х	х
Leaders	х	х	х
Luhn	х	x	х
Mastermind	х	x	х
Middle Char.			х
Paired Digits	х	х	х
Shopping List	х	х	
Snow Day	х		
Solve Boolean	х	х	х
Spin Words	х	x	х
Square Digits	х		
Subs. Cipher	х	х	
Twitter			х
Vector Dist.	х	х	
Total	20	16	16

the bounciness index to calculate the total distance that the ball travels 398 across those bounces. [36] 399

- 3. Bowling (CW) Given a string representing the individual bowls in a 10-frame round of 10 pin bowling, return the score of that round. [37] 400 401 402
- 4. Camel Case (CW) Take a string in kebab-case and convert all of the words to camelCase. Each group of words to convert is delimited by "-", and each grouping is separated by a space. For example: "camel-case example-test-string" → "camelCase exampleTestString". [38]
  403
- 5. Coin Sums (PE) Given a number of cents, find the fewest number of 408 US coins (pennies, nickles, dimes, quarters) needed to make that amount, 409 and return the number of each type of coin as a separate output. [39]
- 6. Cut Vector (CW) Given a vector of positive integers, find the spot where, if you cut the vector, the sum of the numbers on both sides are 411

- 415 either equal, or the difference in sums is as small as possible. Return the416 two resulting subvectors as two outputs. [40]
- 417
  418
  419
  420
  7. Dice Game (PE) Peter has an n sided die and Colin has an m sided die. If they both roll their dice at the same time, return the probability that Peter rolls strictly higher than Colin. [41]
- 421
  422
  423
  8. Find Pair (AoC) Given a vector of integers, return the two elements that sum to a target integer. [42]
- 424 9. Fizz Buzz (CW) Given an integer x, return "Fizz" if x is divisible by
  425 3, "Buzz" if x is divisible by 5, "FizzBuzz" if x is divisible by 3 and 5,
  426 and a string version of x if none of the above hold. [43]
- 427
  428
  429
  430
  10. Fuel Cost (AoC) Given a vector of positive integers, divide each by 3, round the result down to the nearest integer, and subtract 2. Return the sum of all of the new integers. [44]
- 431
   432
   433
   11. GCD [Greatest Common Divisor] (CW) Given two positive integers, return the largest integer that divides each of the integers evenly. [45]
- 434 12. Indices of Substring (CW) Given a text string and a target string,
  435 return a vector of integers of the indices at which the target appears in the
  436 text. If the target string overlaps itself in the text, all indices (including
  437 those overlapping) should be returned. [46]
- 438
  439
  440
  440
  441
  441
  442
  13. Leaders (CW) Given a vector of positive integers, return a vector of the leaders in that vector. A leader is defined as a number that is greater than or equal to all the numbers to the right of it. The rightmost element is always a leader. [47]
- 443
  444
  445
  446
  446
  447
  14. Luhn (CW) Given a vector of 16 digits, implement Luhn's algorithm to verify a credit card number, such that it follows the following rules: double every other digit starting with the second digit. If any of the results are over 9, subtract 9 from them. Return the sum of all of the new digits. [48]
- 44815. Mastermind (HW) Based on the board game Mastermind. Given a449Mastermind code and a guess, each of which are 4-character strings con-450sisting of 6 possible characters, return the number of white pegs (correct451color, wrong place) and black pegs (correct color, correct place) the452codemaster should give as a clue.
- 453
  454
  455
  456
  16. Middle Character (CW) Given a string, return the middle character as a string if it is odd length; return the two middle characters as a string if it is even length. [49]
- 457
  458
  459
  17. Paired Digits (AoC) Given a string of digits, return the sum of the digits whose following digit is the same. [50]
- 460

18.	<b>Shopping List (CW)</b> Given a vector of floats representing the prices of various shopping goods and another vector of floats representing the percent discount of each of those goods, return the total price of the shopping trip after applying the discount to each item. [51]	461 462 463 464
19.	<b>Snow Day (HW)</b> Given an integer representing a number of hours and 3 floats representing how much snow is on the ground, the rate of snow fall, and the proportion of snow melting per hour, return the amount of snow on the ground after the amount of hours given. Each hour is considered a discrete event of adding snow and then melting, not a continuous process.	$ \begin{array}{r} 465 \\ 466 \\ 467 \\ 468 \\ 469 \\ 470 \\ \end{array} $
20.	Solve Boolean (CW) Given a string representing a Boolean expression consisting of T, F, $ $ , and &, evaluate it and return the resulting Boolean. [52]	471 472 473 474
21.	<b>Spin Words (CW)</b> Given a string of one or more words (separated by spaces), reverse all of the words that are five or more letters long and return the resulting string. [53]	475 476 477
22.	Square Digits (CW) Given a positive integer, square each digit and concatenate the squares into a returned string. [54]	$478 \\ 479 \\ 480$
23.	Substitution Cipher (CW) This problem gives 3 strings. The first two represent a cipher, mapping each character in one string to the one at the same index in the other string. The program must apply this cipher to the third string and return the deciphered message. [55]	481 482 483 484 485
24.	<b>Twitter (HW)</b> Given a string representing a tweet, validate whether the tweet meets Twitter's original character requirements. If the tweet has more than 140 characters, return the string "Too many characters". If the tweet is empty, return the string "You didn't type anything". Otherwise, return "Your tweet has X characters", where the X is the number of characters in the tweet.	$   \begin{array}{r}     486 \\     487 \\     488 \\     489 \\     490 \\     491   \end{array} $
25.	Vector Distance (CW) Given two $n$ -dimensional vectors of floats, return the Euclidean distance between the two vectors in $n$ -dimensional space. [56]	492 493 494 495
6	Implementing PSB2	$496 \\ 497$
Whi prob Here in n I of la	ile Section 5 provides English-language descriptions of the 25 benchmark olems, these are not sufficient to implement each problem in a new system. e we discuss the system-agnostic details for implementing these problems ew synthesis systems. For reasons discussed in Section 3, we have created datasets consisting arge numbers of inputs and correct outputs for every problem [57]. <sup>5</sup> The	$ \begin{array}{r}     498 \\     499 \\     500 \\     501 \\     502 \\     503 \\     504 \\ \end{array} $
<sup>5</sup> O Benc	bur datasets follow the model of other machine learning datasets such as Penn ML hmarks [58, 59] and the UCI ML Repository [60].	$504 \\ 505 \\ 506$

507dataset for each problem consists of a small number of hand-chosen inputs, 508often addressing edge cases for the problem, and 1 million randomly-generated 509inputs falling within the constraints of the problem's inputs; the only exception being the Coin Sums problem, which has 10,000 possible integer inputs. We 510recommend that each different program synthesis run use a different set of 511512training data, composed of every one of the hand-chosen inputs and a random 513sample of the randomly-generated inputs. The alternative method of using 514the same fixed set of inputs for every run could happen to use a particularly 515lucky (or unlucky) set of inputs, and across multiple experiments, could lead to undesirable tuning to the particular data; using different randomized inputs 516517for each run avoids these issues. Our datasets will allow those implementing 518PSB2 to simply sample the provided data, greatly decreasing the barrier to 519using PSB2. The PSB2 datasets can be found permanently on Zenodo.<sup>6</sup> For more information about distributions of inputs in randomly-generated inputs, 520see the reference implementation, which was used to generate the datasets.<sup>7</sup> 521

522 In order to make it even easier to download and sample the PSB2 datasets, 523 we have created libraries that implement this functionality in Python<sup>8</sup> and 524 Clojure<sup>9</sup>. Each library can be easily installed using the host language's pack-525 age management system. Both libraries provide one function that downloads, 526 caches, and samples the dataset to produce training and test sets for a given 527 problem.

528 When using our provided datasets, one could sample different sizes of train-529 ing and unseen test sets to fit a given experiment. Our recommendation, which 530 we use in the experiments, is to use 200 example cases for the training set 531 (including all hand-chosen inputs and the remaining inputs sampled randomly) 532 and 2000 for the unseen test set. However, some synthesis methods may need 533 smaller or larger training sets, and PSB2 can flexibly adapt to such systems.

534Program synthesis methods that have been applied to PSB1 have used 535varying methods for constraining the instruction set and other program syntax. For example, some have used grammars [2, 3, 23, 25–28] while others 536537 have used data-type categorized subsets of an instruction set [1, 18]. We do 538not want to constrain what a reasonable approach to selecting instructions 539may look like for any given program synthesis system. However, we also warn against cherry-picking a small subset of instructions suspected of being useful 540541for a particular problem. Part of the difficulty of general program synthesis is 542that a system must manage a large set of potentially useful instructions, find-543ing those relevant to a particular problem. We recommend employing a large 544set of general-purpose instructions when using PSB2 to benchmark program 545synthesis to best replicate the conditions of a real-world scenario.

- 546
- 547  $^{6}_{https://zenodo.org/record/4678739}$
- 548 <sup>7</sup>https://github.com/thelmuth/Clojush/releases/tag/psb2-v1.0
- <sup>8</sup>https://github.com/thelmuth/psb2-python
- 549 <sup>9</sup>https://github.com/thelmuth/psb2-clojure
- 550
- 551
- 552

## 7 Best Practices for Benchmarking

There are many ways to use a suite of benchmark problems, ranging from illustrating a synthesis method's behavior to tuning algorithms and hyperparameters [61]. However, many uses require fair comparisons between systems or methods in order to benchmark performance or other characteristics. Here we make recommendations for benchmarking best practices in order to ensure fair comparisons. 555

In order to compare performance, we need to ensure that the compared sys-561tems or methods are given similar computational resources. Arguments could 562be made for a variety of budgets to limit execution, including processing time, 563number of GP generations, program fitness evaluations (across all training 564cases), program executions (per case), and number of individual instructions 565executed. All of these methods have flaws, and all might be useful in certain 566circumstances. However, we generally recommend using a fixed program exe-567 cution budget to limit the number of generated program executions in a single 568program synthesis run. This type of budget is more general than using a bud-569get of generations or fitness evaluations, as it allows systems to use different 570 population sizes or training set sizes yet still budget similar computation. On 571the other hand, a budget of processing time or instruction executions may 572be misleading when running experiments on different hardware or different 573 synthesis systems. We recommend using a budget of 60 million program execu-574tions; we allocate these to 200 training cases used to evaluate a population of 5751000 individuals for 300 generations in our experimental GP runs, but other 576allocations of the same executions would be reasonable. 577

In order to produce sufficient data on the performance of a stochastic search 578 such as GP, one needs to conduct multiple runs/trials of the system to get an 579idea of the distribution of results. While we do not want to set a fixed number 580 of recommended trials, more are certainly better when considering the signif-581icance of the collected results. We expect that around 30 trials would be the 582minimum acceptable for showing significant differences, though larger num-583bers of trials such as 100 or more would be more likely to produce significant 584differences in statistical tests; we use 100 runs per problem in our experiments. 585

When evaluating the performance of a synthesis system on PSB2, we recom-586mend using success rate (the number of synthesis runs that produce a solution) 587 as the primary measure of performance, as was recommended in PSB1 [1]. For 588the synthesis of software, generating programs that pass most, but not all, 589of the training cases is not sufficient; for this reason success rate is a better 590measure of performance than other metrics such as mean best fitness or mean 591number of training cases passed. In particular, a solution must not only pass 592all cases in the training set, but also all of the cases in the test set, to ensure 593that it generalizes to unseen data. This avoids considering programs that over-594fit the training data, such as by memorizing the correct output to each input, 595as solutions. We recommend sampling 2000 (or more) random inputs to create 596the unseen test set. 597

598

13

 $553 \\ 554$ 

In order to determine whether differences in performance are statistically significant, experimenters should use an appropriate significance test. In the case of success rates measured in numbers of success and failure runs, a chisquared significance test is appropriate. When comparing multiple methods, it is necessary to apply a correction to account for multiple comparisons, such as the Holm–Bonferroni correction.

# <sup>605</sup> <sup>606</sup> 8 Experimental Methods and System <sup>607</sup> <sup>608</sup> Parameters

609 In this section we will discuss the system-specific parameters and choices that 610 must be decided in order to use PSB2. In contrast with the previous sections 611 on general implementation and benchmarking considerations, the choices here 612 may differ considerably for different program synthesis systems. For our exper-613 iments, we used PushGP; we will describe in general the decisions that must 614 be made and give our specific choices.

615PushGP evolves programs in the language Push, a stack-based program-616ming language built specifically for use in genetic programming [62, 63]. Every 617data type has its own stack, and each Push instruction acts by pushing and 618popping various elements on and off the stacks. The output of each problem is 619 typically the top element on a particular stack. The interpreter executes pro-620 grams that are themselves placed on an **exec** stack, allowing **exec** instructions 621 to manipulate control flow as well as the program itself as it runs. We provide 622 a reference implementation in Clojure of the PushGP system used to produce 623 our results, which includes each problem in PSB2.<sup>10</sup> This reference implemen-624 tation is the same implementation of PushGP used in recent research using 625PSB1, e.g. [7, 16, 17].

626 We discuss the general design of program synthesis instruction sets in 627 Section 6. For our PushGP experiments, we use the general process recom-628mended in PSB1, where, for each problem, we identify which data types 629 (corresponding to stacks) are relevant and include all implemented instruc-630 tions that use those stacks [1]. In Tables 3 and 4, we present the data types 631 we chose to include for each problem, and the total number of instructions in 632the instruction set. Every problem includes the exec, integer, and Boolean 633 types, since they are required for a variety of control flow instructions such as 634iteration and conditionals. These large instruction sets contain a wide range 635of general-purpose Push instructions, including some new instructions imple-636mented since PSB1, avoiding the cherry-picking of clearly useful instructions. 637 For the full list of instructions for each data type, see Appendix A.

Research utilizing PSB1 in using transfer-learned instruction sets showed that the composition of the instruction set matters a great deal to problemsolving performance [18]. While we do not use fully transfer-learned instruction sets here, we do make use of one simple take-away: that including larger proportions of input instructions and literals/ERCs improves performance. An 643

<sup>&</sup>lt;sup>10</sup>https://github.com/thelmuth/Clojush/releases/tag/psb2-v1.0

#### Springer Nature 2021 LATEX template

#### Applying Genetic Programming to PSB2 15

0

 $\begin{array}{c} 645 \\ 646 \end{array}$ 

Table 3 Instructions and data types used in our PushGP implementation of each 647 problem. The column "# Instructions" reports the number of instructions, terminals, and 648 ephemeral random constants (ERC) used for each problem. The middle columns show 649 which data types were used for each problem. For example, the Basement problem used all 650 instructions relevant to exec, integers, Booleans, and vectors of integers. The last column lists the literals and ERCs used for the problem. Here, char literals are represented in 651Clojure syntax, starting with a backslash, and strings are surrounded by double quotation 652 marks. The "Problems" row simply counts how many problems use each data type. The 653 "Instructions" row shows the number of Push instructions that primarily use each data 654 type; some use multiple types but are only counted once.

Problem	# Instructions	exec	integer	float	Boolean	char	string	vector of integers	vector of floats	Literals and ERCs (besides inputs)
Basement	117	v	v		v			v		[], -1, 0, 1,
Bouncing B.	127	x	x	х	x			л		0.0, 1.0, 2.0
Bowling	161	x	x		x	x	x			<pre>\-, \X, \/, \1, \2, \3, \4, \5, \6, \7, \8, \9, 10, integer ERC</pre>
										<pre>\-, \space, visible character ERC,</pre>
Camel Case	151	х	x		х	x	x			string ERC
Coin Sums	86	х	x		х					0, 1, 5, 10, 25
Cut Vector	116	х	x		х			х		[], 0
Dice Game	125	х	х	х	х					0.0, 1.0
Find Pair	120	x	x		x			x		-1, 0, 1, 2, integer ERC
Fizz Buzz	118	x	x		x		x			"Fizz", "Buzz", "FizzBuzz", 0, 3, 5
Fuel Cost	117	x	x		x			x		0, 1, 2, 3, integer ERC
GCD	79	x	x		x					integer ERC
Ind. of Subs.	184	х	х		х	х	х	х		[], "", 0, 1
Leaders	114	x	x		x			x		[], vector ERC
Luhn	117	x	x		x			x		0, 2, 9, 10, integer ERC
Mastermind	123	x	x		x	x	x			0, 1, $B$ , $R$ , $W$ , $Y$ , $O$ , $G$
										"", 0, 1, 2, integer
Middle Char.	151	x	x		x	x	x			ERC
Paired Digits	149	х	x		x	x	x			integer ERC
Shopping List	161	x	x	x	х				x	0.0, 100.0, float ERC

#### Springer Nature 2021 LATEX template

#### 16 Applying Genetic Programming to PSB2

Problem	# Instructions	exec	integer	float	Boolean	char	$\operatorname{string}$	vector of integers	vector of floats	Literals and ERCs (beside inputs)
Snow Day	191									0, 1, -1, 0.0,
Show Day	191	x	x	x	x					true, false, \t.
Solve Boolean	153	x	x		x	x	x			f,  , &
Spin Words	152	x	x		x	x	x			4, 5, \space, visible character ERC, string ERC
										"", 0, 1, 2,
Square Digits	151	х	х		х	х	х			integer ERC
Twitter	151	x	x		x	x	x			), 140, "Too many characters", "You didn't typ anything", "you tweet has ", " characters"
Vector Dist.	160	х	х	х	х				х	[], 0
Problems Instructions		25 29	$\begin{array}{c} 25\\ 33 \end{array}$	$5\\45$	$25 \\ 21$	11 21	$\begin{array}{c} 12 \\ 47 \end{array}$	$7\\34$	$2 \\ 34$	

691 Table 4 Continuation of Table 2

717

718explanation of this result is that most Push instructions decrease stack sizes 719 by consuming arguments and producing fewer return values, so increasing 720 inputs and literals creates more data on which instructions can act. We boost 721 the presence of input instructions and literals in the instruction set, making 722 input instructions fill 15% of the instruction set and literals/ERCs fill 5% of 723 the instruction set. The additional input instructions are evenly distributed 724between each input for problems with multiple inputs, and literals/ERCs are 725similarly evenly distributed for each listed in the last column of Tables 3 and 4. 726

For problems with multiple outputs, different synthesis systems will need 727 to make choices specific to the language of the synthesized programs. Initial 728 experiments in PushGP show that it achieves better results on multi-output 729problems when using one output instruction per output. These output instruc-730 tions are included in the instruction set for such problems and will always 731 appear in solution programs. An example of this is for Coin Sums, which has 732 4 outputs. We provide four corresponding output instructions, each of which 733 takes the top integer from the integer stack and stores it in a write-only reg-734 ister for that output; further calls to an output instruction will overwrite this 735 output register. 736

In order to define each problem for GP, we not only need the inputs and 737 correct outputs for each problem, but also how to calculate the error function 738 based on the correct output and a program's output. Here we describe the 739 error functions we employed in our experiments, which we recommend for 740 any GP system implementing PSB2; other non-GP program synthesis systems 741 may require entirely different metrics. For each output data type, we use the 742following standard error functions for problems outputting that data type: 743

- Integer or float: absolute value of the difference between program 744output and correct output. 745
- **Boolean:** 0 for correct and 1 for incorrect output.
- String: Levenshtein string edit distance between the program output and 747 correct output. 748
- Vector of integers: add the difference in length between the program's 749 output vector and the correct vector times 1000 to the absolute difference 750between each integer and the corresponding integer in the correct vector. 751

The only exception is for the Indices of Substring problem, where we used 752Levenshtein distance to compare vectors of integers, since it makes more sense 753when comparing vectors of indices. In PushGP, some evolved programs will 754not return values of a problem's output data type; we give a penalty error 755value specific to the problem when this occurs. 756

As has been shown to be effective at improving generalization, we use an 757 automatic simplification procedure on every evolved Push program that passes 758all of the training cases before testing it on the test set [16]. 759

Unlike for PSB1, we aimed to keep all system-specific parameters constant 760 between problems, increasing ease of use for both implementation and report-761 ing of results. These parameters were chosen based on prior experience and 762 reasonable performance; we leave optimizing parameter settings as an open 763 research question. Other systems may choose to use different system-specific 764parameters. 765

Our PushGP system uses linear Plush genomes that are initialized by gen-766 erating lists of random instructions from the instruction set [15]. We list the 767 important parameters used in our experiments below: 768769

- Maximum initial genome size: 250 genes
- Maximum genome size: 500 genes
- Population size: 1000
- Maximum generations per run: 300
- Maximum steps of the Push interpreter when executing one program: 773 2000 774
- Parent selection: lexicase selection [64, 65]
- Genetic operator: Uniform Mutation with Additions and Deletions 776 (UMAD), used to make 100% of children [17]. 777

• UMAD addition rate: 0.09

779 As described in Section 7, using the exact same population size and generations is not necessary for comparisons between systems; instead, we recommend 780

781

770

771

772

775

778

746

783 Table 5 Results from 100 PushGP runs on each problem. "Succ." gives the number of runs that successfully find a program that passes every training case and perfectly passes a set of 2000 unseen test cases. "Gen." gives the proportion of solutions on the training data that generalize to unseen data. "Size" gives the size of the smallest automatically simplified solution that generalized to unseen data in any of our runs. "Time" is the average number of seconds taken per generation.

788	Problem	Succ.	Gen.	Size	Time
789					
790	Basement	1	1.00	16	250
791	Bouncing Balls	0	0.00	15	311
702	Bowling	0	-	-	206
192	Camel Case	1	1.00	20	95
793	Coin Sums	2	1.00	22	213
794	Cut Vector	0	-	-	194
795	Dice Game	0	-	20	287
706	Find Pair	4	1.00	13	763
790	Fizz Buzz	25	0.96	15	281
797	Fuel Cost	50	1.00	9	305
798	GCD	8	0.67	9	198
799	Indices of Substring	0	-	19	241
200	Leaders	0	-	-	302
800	Luhn	0	-	-	239
801	Mastermind	0	-	-	126
802	Middle Character	57	0.86	10	547
803	Paired Digits	8	1.00	12	250
804	Shopping List	0	-	-	714
804	Snow Day	4	1.00	11	263
805	Solve Boolean	5	1.00	18	373
806	Spin Words	0	-	-	443
807	Square Digits	0	-	13	435
202	Substitution Cipher	60	0.98	9	395
000	Twitter	31	0.74	22	527
809	Vector Distance	0	-	-	667

810

using a maximum budget of 60 million program executions regardless of other
 settings.

813

# <sup>814</sup><sub>815</sub> 9 Experimental Results

In order to give a baseline performance of the 25 problems in PSB2, we conducted 100 PushGP runs on each problem using the experimental methods described in Section 8. Additionally, in Section 10 we conduct a control study using randomly generated programs.

We present success rates of our runs in Table 5. Out of the 25 problems in PSB2, 13 were solved by PushGP. Of these 13, 3 of them had 50 or more successes (Fuel Cost, Middle Character, and Substitution Cipher) and 2 others had 25 or more successes (Fizz Buzz and Twitter). The remaining 8 had fewer than 10 solutions, showing that they are solvable by GP but leave a lot of room for improvement.

The second column in Table 5 gives the generalization rate of all evolved
solutions for problems on which PushGP produced at least one program that

829 solved every training case. The generalization rate is calculated as the number of solution programs that pass the unseen test set divided by the number of 830 solution programs that pass the training set. For most problems with training 831 set solutions, those solutions tended to generalize well with rates of 0.95 to 1.0. 832 The three problems with lower generalization rates, GCD, Middle Character, 833 and Twitter all had rates over 0.5. However, Bouncing Balls found 2 solutions 834 on the training data, but neither of them generalized to the test, which resulted 835 in 0 successful runs and a 0.00 generalization rate. 836

Another way of approximating the difficulty of these problems is by looking 837 at the size of the smallest solution program found for each problem. Smaller 838 solutions are easier for a program synthesis system to generate, simply because 839 they require assembling fewer instructions in the right order. Our results are 840 particular to Push program solutions, but should correlate with the sizes of 841 programs needed to solve these problems in other systems. In order to find each 842 843 size, we took each solution program and automatically simplified it to produce a smaller equivalent program [16, 66]. Of these simplified programs, in Table 5 844 we report the smallest simplified solution size out of all simplified solutions to 845 each problem (including experiments we conduct in Section 11). We see that 846 the smallest solution size is 9 instructions for three problems, and one other 847 has a size of 10; three of these four problems also had the highest success rates 848 in PSB2. Many others have larger smallest solution sizes, though we note that 849 with the small sample sizes of solutions for some problems, smaller solutions 850 may exist. In comparison, [1] reported that 8 of the problems in PSB1 had a 851 smallest Push solution size less than 9, the minimum for PSB2. Along with 852 success rates, these sizes of smallest solutions give evidence that the problems 853 in PSB2 are more difficult than those in PSB1. 854

The last column of Table 5 gives the average number of seconds per generation over all of the PushGP runs for the problem. Note that these runs were conducted on two different computing clusters, each of which is composed of heterogeneous machines, so these measurements should only be considered as rough approximations of running time. To that end, we note that all problems have generational running times within one order of magnitude of each other, meaning there are not any exceptionally slow or fast problems. 861

# 10 Control Study: Random Search

To make a fair comparison with our GP study, we run the same number of 872 program executions with one "run" of random search. Our GP runs evaluate 873

874

862 863

```
(boolean_yank exec_yankdup in1 boolean_dup
875
      in1 integer_gte integer_flush in2 in1 exec_y
876
       (in1 integer_yankdup integer_mod integer_swap))
877
878
     (-1 exec_k (-8 exec_do*count -5 integer_flush in2 in1
879
      exec_y (in2 integer_yankdup integer_mod integer_swap)))
880
     Fig. 1 The two simplified Push solutions generated by random search to the GCD problem.
881
882
     300.000 individuals, using a population size of 1000 over 300 generations. While
883
884
```

we could simply run 300,000 randomly generated programs on each of the 200 training cases for a problem, since the relative performance of each program 885 does not matter (since we are not using it to guide search), we only need to 886 887 know if each program perfectly passes all cases or not. As such, we can instead execute each program on one case at a time, stopping and moving on to the 888 next program when it fails a single case. Since most random programs will fail 889 the first case they are tested on, we can evaluate nearly 200 times as many 890 random programs this way, equivalent to the 200 cases used for evaluating GP 891 892 programs. Thus nearly 60 million random programs can be tested with the same execution budget as one GP run. We repeated this process 100 times per 893 problem, to give random search the same computation as GP. 894

Out of the 25 PSB2 problems, random search found only 2 solutions, both 895 to the GCD problem. The automatically simplified versions of those solution 896 897 programs are presented in Figure 1. Both programs have similar structure, especially near the end, starting with exec\_y, a version of the Y combinator. 898 In both programs, **exec\_v** recursively calls the parenthesized code block, which 899 performs the heart of the Euclidean algorithm for finding a greatest common 900 denominator. The random generation of these two solutions is surprising; while 901 902 GCD is tied for the smallest simplified solution sizes in PSB2, all solutions so far need at least 9 instructions, and it is not among the 5 problems solved 903 most often in our initial PushGP results (see Table 5). However, GCD does 904 have the smallest instruction set of the PSB2 problems, with 79 instructions 905 being significantly smaller than the instruction set for most other problems 906 907 (see Table 3). This smaller instruction set reduces the search space a bit, but we would still expect it unlikely to randomly generate these solutions in the 908 huge combinatorial space. 909

Giving random search almost 200 times as many programs to consider as 910 GP, it has not produced a single solution to any of the PSB2 problems besides 911 912 GCD. This is no great surprise, since even the problems with the smallest possible solutions need to line up at least 9 instructions in the right order to 913 solve the problem, producing a large combinatorial space when considering 914 that every problem uses a set of at least 79 instructions or more. These results 915 show that some type of program synthesis method, such as GP, is almost 916 917 always necessary to solve the problems in PSB2.

- 918
- 919
- 920

## 11 Down-sampled Lexicase Selection

In order to demonstrate the utility of PSB2 for benchmarking GP, we have conducted a set of PushGP runs using down-sampled lexicase selection [7, 8, 29, 30] to compare with standard lexicase selection. These runs additionally assess the current state-of-the-art for PushGP, since down-sampled lexicase has given better results than standard lexicase on the PSB1 problems. 927

In down-sampled lexicase selection, the population each generation is only 928 evaluated on a random subsample of the available training cases, with the 929 subsample being resampled each generation. We used a down-sampling level of 930 0.25, meaning only 50 out of 200 training cases are used each generation. By 931 down-sampling, less information is available for evaluating the performance of 932 each individual; however, each individual is executed fewer times. As such, we 933 can allow GP to run for more generations (or use a larger population size) 934while maintaining the same execution budget. In our experiments here, we 935 increased the maximum number of generations from 300 to 1200, a change 936 proportional to the down-sampling level, maintaining the execution budget of 937 60 million programs. In previous work using PSB1, the drawback of having 938 less information to assess each program was more than compensated for by 939 increasing the number of individuals considered during search, and here we 940 assess whether that holds true for PSB2. 941

The results comparing down-sampled lexicase selection to standard lexicase 942 selection can be found in Table 6. Down-sampling improves results compared to 943 lexicase selection across the board, solving 17 of the 25 problems. In addition, 944 the success rates are the same or improved on every problem, with eight of 945 those improvements being significant using a chi-squared test. These results 946 show clear indications that down-sampled lexicase selection is an improvement 947 over standard lexicase selection on the PSB2 problems, aligning with previous 948 results [7, 8]. 949

These results exemplify how PSB2 can be useful for benchmarking, in particular that different methods can produce significantly different results when using the PSB2 problems. At this point, 18 out of 25 problems have been solved at least once using PushGP: 17 of those with down-sampled lexicase selection, and 1 additional in our initial exploratory PushGP runs (the Leaders problem). While we have no guarantees that the other 7 problems can be solved by any program synthesis system, they provide useful targets for future research. 950 951 952 953 954 955 956

# 12 Related Work in Program Synthesis

We expect that PSB2 can be adopted by GP program synthesis practitioners 960 and others who used PSB1. Here, we would like to examine the place of PSB2 961 within the larger community of program synthesis, to see both where it can 962 be adopted and where there is room for cross-pollination of ideas. Program 963 synthesis encompasses a large and varied field; here, we cite representative 964 examples of these approaches without conducting a full literature review of the field. 966

1

921

957

Table 6 Number of successful runs comparing lexicase selection to down-sampled lexicase selection. Underlined values indicate significant improvement of down-sampled lexicase over lexicase using a chi-squared test with a 0.05 significance level. Lexicase never produced better results than down-sampled lexicase. "Problems solved" counts the number of problems each method solved at least once.

971 072	Problem	Down-sampled	Lexicase
973	Basement	2	1
974	Bouncing Balls	3	0
975	Bowling	0	0
076	Camel Case	4	1
910	Coin Sums	<u>39</u>	2
977	Cut Vector	0	0
978	Dice Game	1	0
979	Find Pair	<u>20</u>	4
080	Fizz Buzz	<u>74</u>	25
900	Fuel Cost	<u>67</u>	50
981	GCD	<u>20</u>	8
982	Indices of Substring	4	0
983	Leaders	0	0
084	Luhn	0	0
904	Mastermind	0	0
985	Middle Character	<u>79</u>	57
986	Paired Digits	17	8
987	Shopping List	0	0
088	Snow Day	7	4
900	Solve Boolean	5	5
989	Spin Words	0	0
990	Square Digits	2	0
991	Substitution Cipher	<u>86</u> 59	61
002	1 witter	<u>52</u>	31
994	vector Distance	0	0
993 994	Problems solved	17	13

995

996 Programming by example (PBE) generally includes any method where the 997 specification of the intended program is given by input/output examples, as is 998 the case with PSB2. Most GP methods, including those mentioned in Section 2, use PBE. That said, many non-GP PBE synthesis methods have been proposed 999 1000 [67–73]. Many of these use non-evolutionary search over the space of possi-1001 ble programs, while others use satisfiability modulo theories (SMT) solvers to 1002 generate programs. In either case, the systems in this space are typically lim-1003 ited to generating programs using small domain-specific languages over one 1004 or two datatypes (such as manipulating strings or vectors), both limitations 1005 that would make them difficult to apply to every problem in PSB2. However, 1006 we encourage the exchange of ideas between these fields, and hope that some 1007 non-GP PBE methods can be applied to PSB2 soon.

1008 Natural language specifications, for example a docstring of a desired func-1009 tion, give a natural but often loose specification of the intended functionalty. 1010 However, program synthesis based on natural language specifications could 1011 prove very useful if the desired programs can be accurately synthesized. 1012

1013 Recently, advances in language modeling have led to methods based on producing abstract syntax trees [74, 75], using component-based synthesis [76], 1014 or that simply produce sequences of source code strings [77, 78]. Some of 1015 these methods combine natural language specifications with PBE [74, 76]. 1016 While some of these methods have only been benchmarked on creating small, 1017 domain-specific programs [74–76], those producing source code generate high-1018 level code [77, 78]. These methods could be tested directly on PSB2's problem 1019 descriptions given in Section 5. 1020

Some methods define the problem to be solved using a *formal specifica*-1021 1022 tion, such as a logic formula or a reference implementation [79-83]. These methods also require either a *sketch* of the desired program with holes to be 1023filled in [79] or a problem-specific syntax defined by a grammar [80-83]. For 1024 example, Syntax-Guided Synthesis (SyGuS) problems require formal specifica-1025tions as well as a context-free grammar describing allowable syntax [81]. These 1026 methods have been tested on domain-specific benchmark tasks similar to those 1027 often used for PBE, such as those manipulating one of the following: integers, 1028 strings, bit vectors, or Boolean circuits. One interesting approach combines 1029GP and sketching, where GP creates the partial program and then a SMT 1030 solver is used to fill in the holes [70]. Another combines GP with SMT solvers 1031 to decompose problems into smaller parts [83]. Any of these methods would 1032require the creation of formal specifications of the PSB2 problems, instead of 1033 using examples of the desired behavior. 1034

## 13 Conclusions

We have presented PSB2, the second generation of general program synthesis 1038 benchmark problems. We discussed the past research that has used PSB1, 1039 the lessons learned from years of its use, and why we need a new benchmark 1040 suite. We then provided the sources and problems that make up PSB2, giving 1041 details of how to implement and use it in new systems. Finally, we presented 1042 experimental results showing the increased difficulty of the problems of PSB2 compared to PSB1, and offer comparisons to other synthesis techniques. 1044

We anticipate that using other GP systems (such as those we mention 10451046in Section 2 that have used PSB1) to produce initial results on PSB2 will 1047 provide a useful comparison to the results with PushGP we have presented 1048 here. Research that has compared these different GP methods on the problems 1049in PSB1 has shown them to perform better or worse on different problems, 1050likely due to their different GP representations and instruction sets [3, 14, 33]; 1051we would expect to see differences in abilities on the problems in PSB2 as 1052well. Additionally, we encourage the application of non-evolutionary automatic 1053program synthesis methods to these problems, to better gauge the strengths 1054and weaknesses of these different methods.

PSB2, like PSB1 before it, focuses on core programming methods and data 1055 types. As program synthesis systems become more complex and encounter 1056 real-world scenarios, they may need to utilize more complex resources such 1057

1058

 $1035 \\ 1036$ 

1059 as concurrency or database manipulation. Additionally, they may need to 1060 manipulate problem-specific data types or library-specific functions (a proof-1061 of-concept of one approach to these problems was given using code-building 1062 GP [21]). We imagine that new benchmark problems that address these 1063 challenges will need to be developed in the future.

1064 The lessons learned from PSB1 will make it easier to implement PSB2 in 1065 new program synthesis systems, increasing adoption in the community and fur-1066 thering the field. PSB2 will provide a new target for program synthesis systems, 1067 stretching their capabilities and moving the field toward the types of problems 1068 that may be encountered in real-world program synthesis applications.

1069

# <sup>1070</sup> Acknowledgments

1071

 $1072~{\rm The}$  authors would like to thank Lee Spector, Grace Woolson, and Amr $1073~{\rm Abdelhady}$  for discussions that helped shape this work.

1074

# <sup>1075</sup><sub>1076</sub> Appendix A Push Instructions

1076

1077 Our experiments using PushGP use instruction sets based on the data types 1078 relevant to the problem, as discussed in Section 8. Tables A1 and A2 con-1079 tain the Push instructions for each of the types referenced in Tables 3 and 4. 1080 Experiments using PSB2 do not need to match this instruction set exactly, and 1081 indeed for non-Push program synthesis systems, there may be wildly different 1082 instruction sets. However, we include the exhaustive list of instructions that 1083 we used in order to give the full details of the system and allow comparisons 1084 with instruction sets in other systems.

1085

# $\frac{1086}{1087}$ References

1087 **1**087 **1**088 [1

- 1088 [1] Helmuth, T., Spector, L.: General program synthesis benchmark suite. In: GECCO '15: Proceedings of the 2015 Conference on Genetic and Evolutionary Computation Conference, pp. 1039–1046. ACM, Madrid, Spain (2015). https://doi.org/10.1145/2739480.2754769
- 1092
- 1093 [2] Forstenlechner, S., Fagan, D., Nicolau, M., O'Neill, M.: A grammar design pattern for arbitrary program synthesis problems in genetic programming. In: EuroGP 2017: Proceedings of the 20th European Conference on Genetic Programming. LNCS, vol. 10196, pp. 262–277. Springer, Amsterdam (2017). https://doi.org/10.1007/978-3-319-55696-3\_17
- 1098
- 1099 [3] Hemberg, E., Kelly, J., O'Reilly, U.-M.: On domain knowledge and novelty
  100 to improve program synthesis performance with grammatical evolution.
  1101 In: GECCO '19: Proceedings of the Genetic and Evolutionary Computa1102 tion Conference, pp. 1039–1046. ACM, Prague, Czech Republic (2019).
  1103 https://doi.org/10.1145/3321707.3321865
- 1104

Table A1 The Push instructions included in our experiments.

Data Type	Instructions
exec	exec_do*count, exec_do*range, exec_do*times, exec_do*vector_float, exec_do*vector_integer, exec_do*while, exec_dup, exec_dup_items, exec_dup_times, exec_empty, exec_eq, exec_flush, exec_if, exec_k, exec_noop, exec_pop, exec_rot, exec_s, exec_shove, exec_stackdepth, exec_string_iterate, exec_swap, exec_when, exec_while, exec_y, exec_yank, exec_yankdup, tag, tagged
integer	integer_abs, integer_add, integer_dec, integer_div, integer_dup, inte- ger_dup_items, integer_dup_times, integer_empty, integer_eq, integer_flush, integer_fromboolean, integer_fromchar, integer_fromfloat, integer_fromstring, integer_gt, integer_gte, integer_inc, integer_lt, integer_lte, integer_max, inte- ger_min, integer_mod, integer_mult, integer_negate, integer_pop, integer_pow, integer_rot, integer_shove, integer_stackdepth, integer_sub, integer_swap, integer_yank, integer_yankdup
float	float_abs, float_add, float_arccos, float_arcsin, float_arctan, float_ceiling, float_cos, float_dec, float_div, float_dup, float_dup_items, float_dup_times, float_empty, float_eq, float_floor, float_flush, float_fromboolean, float_fromchar, float_frominteger, float_fromstring, float_gt, float_gte, float_inc, float_log10, float_log2, float_lt, float_lte, float_max, float_min, float_mod, float_mult, float_negate, float_pop, float_pow, float_rot, float_shove, float_sin, float_sqrt, float_square, float_stackdepth, float_sub, float_swap, float_tan, float_vank, float vank, dup
Boolean	boolean_and, boolean_dup, boolean_dup_items, boolean_dup_times, boolean_eq, boolean_flush, boolean_fromfloat, boolean_frominteger, boolean_invert_first_then_and, boolean_invert_second_then_and, boolean_not, boolean_or, boolean_pop, boolean_rot, boolean_shove, boolean_stackdepth, boolean_swap, boolean_xor, boolean_yank, boolean_yankdup
char	char_allfromstring, char_dup, char_dup_items, char_dup_times, char_empty, char_eq, char_flush, char_fromfloat, char_frominteger, char_isdigit, char_isletter, char_iswhitespace, char_lowercase, char_pop, char_rot, char_shove, char_stackdepth, char_swap, char_uppercase, char_yank, char_yankdup
string	string_butlast, string_capitalize, string_concat, string_conjchar, string_contains, string_containschar, string_dup, string_dup_items, string_dup_times, string_empty, string_emptystring, string_eq, string_first, string_flush, string_fromboolean, string_fromchar, string_indexofchar, string_frominteger, string_includes, string_indexof, string_indexofchar, string_last, string_length, string_lowercase, string_removechar, string_replace, string_replacechar, string_replacefirst, string_replacefirstchar, string_rest, string_reverse, string_rot, string_setchar, string_shove, string_sort, string_split, string_stackdepth, string_substring, string_swap, string_take, string_uppercase, string_yank, string_yankdup

- [4] Lalejini, A., Ofria, C.: Tag-accessed memory for genetic programming. 1145In: GECCO '19: Proceedings of the Genetic and Evolutionary Computa-1146tion Conference Companion, pp. 346–347. ACM, Prague, Czech Republic 1147(2019). https://doi.org/10.1145/3319619.3321892 1148
- 1149[5] Rosin, C.D.: Stepping stones to inductive synthesis of low-level looping

1151 Table A2 Continuation of Table A1

Г Т	Data Type	Instructions
vi o: ir	ector f ntegers	vector_integer_butlast, vector_integer_concat, vector_integer_conj, vec- tor_integer_contains, vector_integer_dup, vector_integer_dup_items, vector_integer_dup_times, vector_integer_empty, vector_integer_dup_tems, vector_integer_eq, vector_integer_first, vector_integer_flush, vec- tor_integer_indexof, vector_integer_last, vector_integer_length, vector_integer_pushall, vector_integer_remove, vector_integer_replace, vector_integer_replacefirst, vector_integer_rest, vector_integer_replace, vector_integer_replacefirst, vector_integer_rest, vector_integer_replacefirst, vector_integer_rest, vector_integer_rest, vector_integer_rest, vector_integer_replacefirst, vector_integer_rest, vector_integer_res
		vector_integer_stackdepth, vector_integer_subvec, vector_integer_swap,
		vector_integer_take, vector_integer_yank, vector_integer_yankdup
V	ector f	vector_float_butlast, vector_float_concat, vector_float_conj, vec- tor_float_contains vector_float_dup_vector_float_dup_items vec-
fl	oats	tor_float_dup_times, vector_float_empty, vector_float_emptyvector,
		vector_float_eq, vector_float_first, vector_float_flush, vector_float_indexof, vec-
		vector_float_pop, vector_float_pushall, vector_float_remove, vec-
		tor_float_replace, vector_float_replacefirst, vector_float_rest,
		vector_float_reverse, vector_float_rot, vector_float_set, vector_float_shove, vec- tor_float_sort_vector_float_stackdepth_vector_float_subvec_vector_float_swap
		vector_float_take, vector_float_yank, vector_float_yankdup
_		
	prog ficial	rams. In: Proceedings of the Thirty-Third AAAI Conference on Art Intelligence. AAAI '19, vol. 33. AAAI Press, Palo Alto, Californi
	USA	(2019)
[6]	] Lim, gram Engi	J., Yoo, S.: Field report: Applying monte carlo tree search for pro- a synthesis. In: International Symposium on Search Based Softwar neering, pp. 304–310 (2016). Springer
[7	] Holm	with T. Spector L. Explaining and exploiting the advantages of
ι.	dowr	n-sampled lexicase selection. In: Artificial Life Conference Proceed
	ings,	pp. 341–349. MIT Press, Online (2020). https://doi.org/10.1162
	isal_a	a_00334. https://www.mitpressjournals.org/doi/abs/10.1162/isal_abs/1000/isal_abs/1000
	0033	4
[0]	] <b>U</b> _l-~	with T. Spector I. Droblem Solving Deposits of Down Severals
0	Levi	rate Selection Artificial Life 1-21 (2021) https://arviv.org
	abs/	https://direct.mit.edu/artl/article-pdf/doi/10.1162/artl.a.00341/
	1960	075/artl_a_00341.pdf. https://doi.org/10.1162/artl_a_00341
[9]	] McD	ermott, J., White, D.R., Luke, S., Manzoni, L., Castelli, M., Van
	nescl	ni, L., Jaskowski, W., Krawiec, K., Harper, R., De Jong, K., O'Reill
	UM	L: Genetic programming needs better benchmarks. In: GECCO '1
	Proc	eedings of the Genetic and Evolutionary Computation Conference
	pp. 7	91–796. AUM, Filladelphia, Fennsylvania, USA (2012). https://do

	.org/10.1145/2330163.2330273	1197
[10]	White D.B. Mcdermett, I. Castelli, M. Manzoni, I. Coldman, B.W.	1198
[10]	Kronberger G. Jaśkowski W. O'Beilly U-M. Luke S: Better GP	1199
	benchmarks: community survey results and proposals. Genetic Pro-	1200
	gramming and Evolvable Machines 14(1), 3–29 (2013), https://doi.org/	1201
	10.1007/s10710-012-9177-2	1202 1203
	'	1204
[11]	Woodward, J., Martin, S., Swan, J.: Benchmarks that matter for genetic	1205
	programming. In: GECCO 2014 4th Workshop on Evolutionary Computa-	1206
	tion for the Automated Design of Algorithms, pp. 1397–1404. ACM, Van-	1207
	couver, BC, Canada (2014). https://doi.org/10.1145/2598394.2609875.	1208
	http://doi.acm.org/10.1145/2598394.2609875	1209
[12]	O'Neill M Spector L · Automatic programming · The open issue?	1210
[±=]	Genetic Programming and Evolvable Machines <b>21</b> (1-2), 251–262	1211
	(2020). https://doi.org/10.1007/s10710-019-09364-2. Twentieth Anniver-	1212
	sary Issue	1213
	•	1214
[13]	Helmuth, T., Kelly, P.: PSB2: The second program synthesis bench-	1210
	mark suite. In: 2021 Genetic and Evolutionary Computation Confer-	1210 1217
	ence. GECCO '21. ACM, Lille, France (2021). https://doi.org/10.1145/	1218
	3449639.3459285. https://dl.acm.org/doi/10.1145/3449639.3459285	1219
[14]	Sobania, D., Schweim, D., Rothlauf, F.: Recent developments in program	1220
[]	synthesis with evolutionary algorithms. arXiv (2021) https://arxiv.org/	1221
	abs/2108.12227 [cs.NE]	1222
		1223
[15]	Helmuth, T., Spector, L., McPhee, N.F., Shanabrook, S.: Linear genomes	1224
	for structured programs. In: Genetic Programming Theory and Practice	1225
	XIV. Genetic and Evolutionary Computation. Springer, Ann Arbor, USA	1226
	(2010)	1227
[16]	Helmuth, T., McPhee, N.F., Pantridge, E., Spector, L.: Improving general-	1220
	ization of evolved programs through automatic simplification. In: Proceed-	1220 1230
	ings of the Genetic and Evolutionary Computation Conference. GECCO	1231
	'17, pp. 937–944. ACM, Berlin, Germany (2017). https://doi.org/10.1145/	1232
	3071178.3071330. http://doi.acm.org/10.1145/3071178.3071330	1233
[1]		1234
[17]	form mutation by addition and delation. In: Program synthesis using uni-	1235
	and Evolutionary Computation Conference CECCO '18 pp 1197–1134	1236
	ACM Kvoto Japan (2018) https://doi.org/10.1145/3205455.3205603	1237
	http://doi.acm.org/10.1145/3205455.3205603	1238
	······································	1239
[18]	Helmuth, T., Pantridge, E., Woolson, G., Spector, L.: Genetic	1240 19/1
	source sensitivity and transfer learning in genetic programming.	1241 1242

- 1243 In: Artificial Life Conference Proceedings, pp. 303–311. MIT Press, 1244 Online (2020). https://doi.org/10.1162/isal\_a\_00326. https://www.mi 1245 tpressjournals.org/doi/abs/10.1162/isal\_a\_00326
- 1246
- 1247 [19] Saini, A.K., Spector, L.: Using modularity metrics as design features to
  1248 guide evolution in genetic programming. In: Banzhaf, W., Goodman, E.,
  1249 Sheneman, L., Trujillo, L., Worzel, B. (eds.) Genetic Programming Theory
  1250 and Practice XVII, pp. 165–180. Springer, East Lansing, MI, USA (2019).
  1251 https://doi.org/10.1007/978-3-030-39958-0\_9
- 1252
- [20] Saini, A.K., Spector, L.: Why and when are loops useful in genetic programming? In: Proceedings of the 2020 Genetic and Evolutionary Computation Conference Companion. GECCO '20, pp. 247–248. Association for Computing Machinery, internet (2020). https://doi.org/10.1145/3377929.3389919. https://doi.org/10.1145/3377929.3389919
- 1257
- Pantridge, E., Spector, L.: Code building genetic programming. In: Proceedings of the 2020 Genetic and Evolutionary Computation Conference. GECCO '20, pp. 994–1002. Association for Computing Machinery, internet (2020). https://doi.org/10.1145/3377930.3390239. https://arxiv.org/abs/2008.03649
- 1263
- 1264 [22] Forstenlechner, S., Fagan, D., Nicolau, M., O'Neill, M.: Towards understanding and refining the general program synthesis benchmark suite with genetic programming. In: Vellasco, M. (ed.) 2018 IEEE Congress on Evolutionary Computation (CEC). IEEE, Rio de Janeiro, Brazil (2018). https://doi.org/10.1109/CEC.2018.8477953
- 1269
- 1270 [23] Forstenlechner, S., Fagan, D., Nicolau, M., O'Neill, M.: Extending program synthesis grammars for grammar-guided genetic programming. In: Auger, A., Fonseca, C.M., Lourenco, N., Machado, P., Paquete, L., Whitley, D. (eds.) 15th International Conference on Parallel Problem Solving from Nature. LNCS, vol. 11101, pp. 197–208. Springer, Coimbra, Portugal (2018). https://doi.org/10.1007/978-3-319-99253-2\_16. https://www.spri nger.com/gp/book/9783319992587
- 1278 [24] Forstenlechner, S., Fagan, D., Nicolau, M., O'Neill, M.: Towards effective semantic operators for program synthesis in genetic programming. In:
  1280 GECCO '18: Proceedings of the Genetic and Evolutionary Computation Conference, pp. 1119–1126. ACM, Kyoto, Japan (2018). https://doi.org/ 10.1145/3205455.3205592
- 1283
- 1284 [25] Kelly, J., Hemberg, E., O'Reilly, U.-M.: Improving genetic programming
  with novel exploration exploitation control. In: Sekanina, L., Hu, T.,
  Lourenço, N., Richter, H., García-Sánchez, P. (eds.) EuroGP 2019: Proceedings of the 22nd European Conference on Genetic Programming, pp.
  64-80. Springer, Leipzig, Germany (2019)

[26]	O'Neill, M., Brabazon, A.: Mutational robustness and structural com- plexity in grammatical evolution. In: Coello, C.A.C. (ed.) 2019 IEEE Congress on Evolutionary Computation, CEC 2019, pp. 1338–1344. IEEE Press, Wellington, New Zealand (2019). https://doi.org/10.1109/ CEC.2019.8790010. IEEE Computational Intelligence Society	1289 1290 1291 1292 1293 1204
[27]	Sobania, D., Rothlauf, F.: Challenges of program synthesis with gram- matical evolution. In: Hu, T., Lourenco, N., Medvet, E. (eds.) EuroGP 2020: Proceedings of the 23rd European Conference on Genetic Program- ming. LNCS, vol. 12101, pp. 211–227. Springer, Seville, Spain (2020). https://doi.org/10.1007/978-3-030-44094-7_14	1294 1295 1296 1297 1298 1299
[28]	Lynch, D., McDermott, J., O'Neill, M.: Program synthesis in a continuous space using grammars and variational autoencoders. In: Baeck, T., Preuss, M., Deutz, A., Wang2, H., Doerr, C., Emmerich, M., Trautmann, H. (eds.) 16th International Conference on Parallel Problem Solving from Nature, Part II. LNCS, vol. 12270, pp. 33–47. Springer, Leiden, Holland (2020). https://doi.org/10.1007/978-3-030-58115-2_3	$     1300 \\     1301 \\     1302 \\     1303 \\     1304 \\     1305 \\     1306   $
[29]	Hernandez, J.G., Lalejini, A., Dolson, E., Ofria, C.: Random subsampling improves performance in lexicase selection. In: GECCO '19: Proceedings of the Genetic and Evolutionary Computation Conference Companion, pp. 2028–2031. ACM, Prague, Czech Republic (2019). https://doi.org/10.1145/3319619.3326900	1307 1308 1309 1310 1311 1312
[30]	Ferguson, A.J., Hernandez, J.G., Junghans, D., Lalejini, A., Dolson, E., Ofria, C.: Characterizing the effects of random subsampling and dilu- tion on lexicase selection. In: Banzhaf, W., Goodman, E., Sheneman, L., Trujillo, L., Worzel, B. (eds.) Genetic Programming Theory and Practice XVII, East Lansing, MI, USA (2019)	1313 1314 1315 1316 1317 1318
[31]	Gulwani, S.: Automating string processing in spreadsheets using input- output examples. SIGPLAN Not. $46(1)$ , 317–330 (2011). https://doi.org/ 10.1145/1925844.1926423	1319 1320 1321 1322
[32]	Katayama, S.: Recent improvements of MagicHaskeller. In: Approaches and Applications of Inductive Programming. Springer, Berlin, Heidelberg (2010). https://doi.org/10.1007/978-3-642-11931-6_9. http://dx.doi.org/10.1007/978-3-642-11931-6_9	1323 1324 1325 1326
[33]	Pantridge, E., Helmuth, T., McPhee, N.F., Spector, L.: On the difficulty of benchmarking inductive program synthesis methods. In: Proceedings of the Genetic and Evolutionary Computation Conference Companion. GECCO '17, pp. 1589–1596. ACM, Berlin, Germany (2017). https://doi.org/10.1145/3067695.3082533. http://doi.acm.org/10.1145/3067695.3082533	1327 1328 1329 1330 1331 1332 1333 1334

- 1335 [34] Helmuth, T., Kelly, P.: General Program Synthesis Benchmark Suite 1336Datasets. https://github.com/thelmuth/program-synthesis-benchmark-1337 datasets 1338 1339 [35] Wastl, E.: Advent of Code: Not Quite Lisp. Accessed: 2020-01-20. https: //adventofcode.com/2015/day/1 1340 1341 [36] g964: Code Wars: Bouncing Balls. Accessed: 2020-01-20. https:// 1342www.codewars.com/kata/5544c7a5cb454edb3c000047 1343 1344[37] dnolan: Code Wars: Ten-Pin Bowling. Accessed: 2020-01-20. https:// 1345www.codewars.com/kata/5531abe4855bcc8d1f00004c/javascript 1346 1347[38] jhoffner: Code Wars: Convert String to Camel Case. Accessed: 2020-01-20. 1348 https://www.codewars.com/kata/517abf86da9663f1d2000003 1349 1350 [39] Euler, P.: Project Euler: Coin Sums. Accessed: 2020-01-20. https:// 1351projecteuler.net/problem=31 13521353 [40] myjinxin2015: Code Wars: Fastest Code: Half It IV. Accessed: 2020-01-20. 1354https://www.codewars.com/kata/5719b28964a584476500057d 13551356 [41] Euler, P.: Project Euler: Dice Game. Accessed: 2020-01-20. <br/>  $\rm https://$ projecteuler.net/problem=205 1357 1358 [42] Wastl, E.: Advent of Code: Report Repair. Accessed: 2020-01-20. https: 1359//adventofcode.com/2020/day/1 1360 1361[43] stephenyu: Code Wars: Fizz Buzz. Accessed: 2020-01-20. https:// 1362 www.codewars.com/kata/5300901726d12b80e8000498 1363 1364[44] Wastl, E.: Advent of Code: The Tyranny of the Rocket Empire. Accessed: 13652020-01-20. https://adventofcode.com/2019/day/1 1366 1367 [45] RVdeKoning: Code Wars: Greatest Common Divisor. Accessed: 2020-136801-20. https://www.codewars.com/kata/5500d54c2ebe0a8e8a0003fd/ 1369python 13701371 [46] smile67: Code Wars: Text Search. Accessed: 2020-01-20. https:// www.codewars.com/kata/56b78faebd06e61870001191 137213731374 [47] MrZizoScream: Code Wars: Array Leaders. Accessed: 2020-01-20. https: //www.codewars.com/kata/5a651865fd56cb55760000e0 1375
- ${1379\atop 1380}$  [49] Shivo: Code Wars: Get the Middle Character. Accessed: 2020-01-20. https:

	Applying Genetic Programming to PSB2 31	
	//www.codewars.com/kata/56747 fd5cb988479 a f000028	13
[50]	Wastl, E.: Advent of Code: Inverse Captcha. Accessed: 2020-01-20. https://adventofcode.com/2017/day/1	138 138 138
[51]	rb50: Code Wars: Shopping List. Accessed: 2020-01-20. https://www.codewars.com/kata/596266482f9add20f70001fc	13 13 13
[52]	KenKamau: Code Wars: The Boolean Order. Accessed: 2020-01-20. https://www.codewars.com/kata/59eb1e4a0863c7ff7e000008	13 13 13
[53]	xDranik: Code Wars: Stop gninnipS My sdroW! Accessed: 2020-01-20. https://www.codewars.com/kata/5264d2b162488dc400000001	13 13 13
[54]	MysteriousMagenta: Code Wars: Square Every Digit. Accessed: 2020-01-20. https://www.codewars.com/kata/546e2562b03326a88e000020	13 13 13
[55]	jacobb: Code Wars: Simple Substitution Cipher Helper. Accessed: 2020-01-20. https://www.codewars.com/kata/52eb114b2d55f0e69800078d	13 13 13
[56]	StephenLastname2: Code Wars: Distance Between Two Points. Accessed: 2020-01-20. https://www.codewars.com/kata/ 5a0b72484bebaefe60001867	13 14 14 14
[57]	Helmuth, T., Kelly, P.: PSB2: The Second Program Synthesis Benchmark Suite. Zenodo (2021). https://doi.org/10.5281/zenodo.4678739	14 14 14
[58]	Olson, R.S., La Cava, W., Orzechowski, P., Urbanowicz, R.J., Moore, J.H.: Pmlb: a large benchmark suite for machine learning evaluation and comparison. BioData Mining <b>10</b> (1), 36 (2017). https://doi.org/10.1186/s13040-017-0154-4	14 14 14 14 14
[59]	Le, T.T., La Cava, W., Romano, J.D., Gregg, J.T., Goldberg, D.J., Chakraborty, P., Ray, N.L., Himmelstein, D., Fu, W., Moore, J.H.: Pmlb v1.0: an open source dataset collection for benchmarking machine learning methods. arXiv preprint arXiv:2012.00058 (2020)	14 14 14 14
[60]	Dua, D., Graff, C.: UCI Machine Learning Repository (2017). <a href="http://archive.ics.uci.edu/ml">http://archive.ics.uci.edu/ml</a>	14 14 14
[61]	Bartz-Beielstein, T., Doerr, C., van den Berg, D., Bossek, J., Chan- drasekaran, S., Eftimov, T., Fischbach, A., Kerschke, P., Cava, W.L., Lopez-Ibanez, M., Malan, K.M., Moore, J.H., Naujoks, B., Orzechowski, P., Volz, V., Wagner, M., Weise, T.: Benchmarking in optimization: Best practice and open issues. arXiv (2020) https://arxiv.org/abs/2007.03488 [cs.NE]	14 14 14 14 14 14 14 14
[62]	Spector, L., Robinson, A.: Genetic programming and autoconstructive	14 14

evolution with the push programming language. Genetic Programming
and Evolvable Machines 3(1), 7–40 (2002). https://doi.org/10.1023/A:
1014538503543

1430

1431 [63] Spector, L., Klein, J., Keijzer, M.: The Push3 execution stack and
1432 the evolution of control. In: GECCO 2005: Proceedings of the 2005
1433 Conference on Genetic and Evolutionary Computation, vol. 2, pp. 1689–
1434 1696. ACM Press, Washington DC, USA (2005). https://doi.org/10.1145/
1435 1068009.1068292

- 1436
- [64] Helmuth, T., Spector, L., Matheson, J.: Solving uncompromising problems with lexicase selection. IEEE Transactions on Evolutionary Computation 19(5), 630–643 (2015). https://doi.org/10.1109/
  TEVC.2014.2362729
- 1441
- [65] Spector, L.: Assessment of problem modality by differential performance of lexicase selection in genetic programming: A preliminary report. In: McClymont, K., Keedwell, E. (eds.) 1st Workshop on Understanding Problems (GECCO-UP), pp. 401–408. ACM, Philadelphia, Pennsylvania, USA (2012). https://doi.org/10.1145/2330784.2330846. http://hampshi re.edu/lspector/pubs/wk09p4-spector.pdf
- Robinson, A.: Genetic programming: Theory, implementation, and the evolution of unconstrained solutions. Division III thesis, Hampshire College (May 2001). http://hampshire.edu/lspector/robinson-div3.pdf
- 1452 [67] Gulwani, S.: Automating string processing in spreadsheets using inputoutput examples. In: Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. POPL
  1455 '11, pp. 317–330. ACM, New York, NY, USA (2011). https://doi.org/ 10.1145/1926385.1926423. http://doi.acm.org/10.1145/1926385.1926423
- 1457
- 1458 [68] Menon, A.K., Tamuz, O., Gulwani, S., Lampson, B., Kalai, A.T.: A
  1459 Machine Learning Framework for Programming by Example. ICML, 9
  1460 (2013)
- 1461 1462 [69] Balog, M., Gaunt, A.L., Brockschmidt, M., Nowozin, S., Tarlow, D.: 1463 DeepCoder: Learning to write programs. In: ICLR (2017)
- 1464
- Bladek, I., Krawiec, K.: Evolutionary program sketching. In: Castelli, M., McDermott, J., Sekanina, L. (eds.) EuroGP 2017: Proceedings of the 20th European Conference on Genetic Programming. LNCS, vol. 10196, pp.
  3–18. Springer, Amsterdam (2017). https://doi.org/10.1007/978-3-319-55696-3\_1. http://repozytorium.put.poznan.pl/publication/495662
- 1470
- 1470[71]Zohar, A., Wolf, L.: Automatic Program Synthesis of Long Programs with<br/>a Learned Garbage Collector. NIPS (2018). arXiv: 1809.04682. Accessed

	2021-10-10	1473
[72]	Gulwani, S., Pathak, K., Radhakrishna, A., Tiwari, A., Udupa, A.: Quantitative programming by examples. arXiv (2019) https://arxiv.org/abs/1909.05964 [cs.PL]	1474 1475 1476 1477
[73]	Cropper, A., Morel, R.: Learning programs by learning from failures. Machine Learning $110(4)$ , 801–856 (2021). https://doi.org/10.1007/s10994-020-05934-z. Accessed 2021-10-09	1478 1479 1480 1481
[74]	Polosukhin, I., Skidanov, A.: Neural program search: Solving program- ming tasks from description and examples. arXiv (2018) https://arxi v.org/abs/1802.04335 [cs.AI]	1482 1483 1484 1485
[75]	Bednarek, J., Piaskowski, K., Krawiec, K.: Ain't Nobody Got Time for Coding: Structure-Aware Program Synthesis from Natural Language. arXiv, 12 (2019)	1486 1487 1488 1489
[76]	Rahmani, K., Raza, M., Gulwani, S., Le, V., Morris, D., Radhakrishna, A., Soares, G., Tiwari, A.: Multi-modal Program Inference: a Marriage of Pre-trained Language Models and Component-based Synthesis. arXiv:2109.02445 [cs] (2021). arXiv: 2109.02445. Accessed 2021-09-15	1490 1491 1492 1493 1494
[77]	Chen, M., Tworek, J., Jun, H., Yuan, Q., Ponde, H., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W., Nichol, A., Babuschkin, I., Balaji, S., Jain, S., Carr, A., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W.: Evaluating Large Language Models Trained on Code. arXiv (2021)	$\begin{array}{c} 1434\\ 1495\\ 1496\\ 1497\\ 1498\\ 1499\\ 1500\\ 1501\\ 1502\\ 1503\\ 1504 \end{array}$
[78]	Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., Sutton, C.: Program Synthesis with Large Language Models. arXiv (2021). arXiv: 2108.07732. Accessed 2021-08-21	1505 1506 1507 1508 1509
[79]	Solar-Lezama, A.: Program sketching. International Journal on Software Tools for Technology Transfer $15(5)$ , 475–495 (2013)	$1510 \\ 1511 \\ 1512$
[80]	Alur, R., Bodik, R., Juniwal, G., Martin, M.M.K., Raghothaman, M., Seshia, S.A., Singh, R., Solar-Lezama, A., Torlak, E., Udupa, A.: Syntax-guided synthesis. In: 2013 Formal Methods in Computer-Aided Design, pp. 1–8 (2013). https://doi.org/10.1109/FMCAD.2013.6679385	$1513 \\ 1514 \\ 1515 \\ 1516 \\ 1517 \\ 1518$

1519 [81] Alur, R., Singh, R., Fisman, D., Solar-Lezama, A.: Search-based program synthesis. Commun. ACM 61(12), 84–93 (2018). https://doi.org/10.1145/3208071
1522
1523 [82] Lee, W., Heo, K., Alur, R., Naik, M.: Accelerating search-based program

[523] [82] Lee, W., Heo, K., Alur, K., Naik, M.: Accelerating search-based program
synthesis using learned probabilistic models, 436–449 (2018). https://doi
.org/10.1145/3192366.3192410

[83] Welsch, T., Kurlin, V.: Synthesis through unification genetic programming. In: Proceedings of the 2020 Genetic and Evolutionary Computation Conference. GECCO '20, pp. 1029–1036. Association for Computing Machinery, internet (2020). https://doi.org/10.1145/3377930.3390208.
[531] https://doi.org/10.1145/3377930.3390208

 $\begin{array}{c} 1546 \\ 1547 \end{array}$ 

 $1550 \\ 1551$ 

 $1555 \\ 1556$ 

 $1558 \\ 1559$