

PSB2: The Second Program Synthesis Benchmark Suite



Tom Helmuth and Peter Kelly

Hamilton College

+ General Program Synthesis

- Problems that resemble those that humans solve
 - multiple general-purpose data types
 - control flow
- Supervised learning: specifications given as input/output examples
- GP and non-GP methods

```
1 def spin_words(input1):
2     words = input1.split()
3     result = ""
4     for word in words:
5         if len(word) >= 5:
6             result += word[::-1] + " "
7         else:
8             result += word + " "
9     return result[:-1]
```

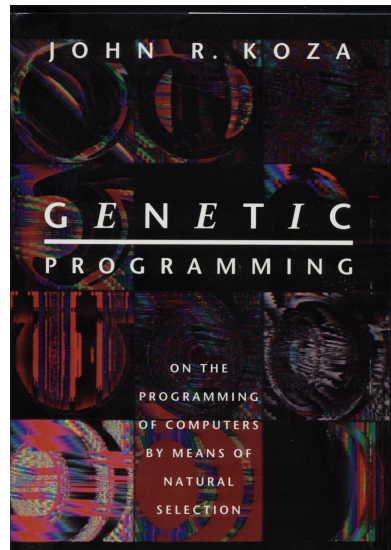
Example Python program of the type we want to generate

+ General Program Synthesis

Benchmark Problems: A Brief History

- **Before 2015:**

- Limited instruction sets
 - Domain specific
 - Toy problems
 - Ex: Santa Fe artificial ant problem
- } both in GP and inductive program synthesis

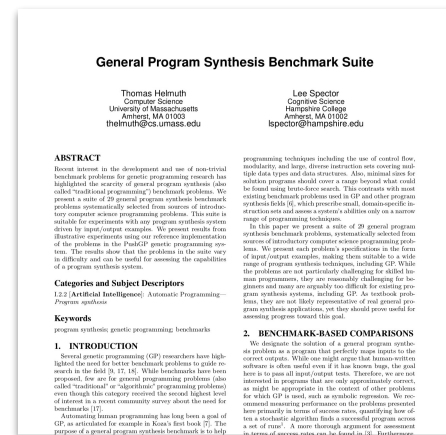


+ General Program Synthesis

Benchmark Problems: A Brief History

- Before 2015:
 - Limited instruction sets
 - Domain specific
 - Toy problems
 - Ex: Santa Fe artificial ant problem
- 2015: PSB1 (Helmuth and Spector)
 - 29 benchmark problems
 - introductory programming homework problem sets
 - has been used to benchmark 10+ GP and non-GP synthesis systems

} both in GP and inductive program synthesis



+ General Program Synthesis

Benchmark Problems: A Brief History

- Before 2015:
 - Limited instruction sets
 - Domain specific
 - Toy problems
 - Ex: Santa Fe artificial ant problem
- 2015: PSB1 (Helmuth and Spector)
 - 29 benchmark problems
 - introductory programming homework problems
 - has been used to benchmark 10+ GP and non-GP synthesis systems
- **2021: PSB2 - this talk!**



both in GP and inductive program synthesis

PSB2: The Second Program Synthesis Benchmark Suite

Thomas Helmuth
Hamilton College
Clifton, New York, USA
thelmuth@hamilton.edu

Peter Kelly
Hamilton College
Clifton, New York, USA
pkelly@hamilton.edu

ABSTRACT

For the past six years, researchers in genetic programming and other program synthesis disciplines have used the General Program Synthesis Benchmark Suite to benchmark many aspects of automatic program synthesis systems. These problems have 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 the original benchmark suite 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.

In this paper, we describe the 21 new general program synthesis benchmark problems that make up PSB1, a new benchmark suite. These problems are created 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 PsatC4F showing this increase in difficulty. These new problems give plenty of room for improvement, pointing the way for the next six to nine years of general program synthesis research.

CCS CONCEPTS

• Software and its engineering → Automatic programming

KEYWORDS

automatic program synthesis, benchmarking, genetic programming

For many years there were no common benchmark problems for evaluating general program synthesis systems, existing problems were either very toy problems or were situated in specific domains where solution programs were composed of a small set of domain-specific instructions. In 2015, the General Program Synthesis Benchmark Suite (PSB1) introduced 29 problems that could be used to benchmark program synthesis systems [17]. Since then, more than 40 research papers have benchmarked 1st-3rd program synthesis systems using PSB1, producing numerous insights into program synthesis.

Of the systems that have adopted PSB1, most fall within the field of genetic programming (GP), including PsatC4F [17], grammar-guided GP [6], grammatical evolution [21], and linear GP [25]. However, non-evolutionary program synthesis methods have also been applied to PSB1, including those based on integer-acceptance hybridizing [4] and Monte Carlo tree search [21]. We expand on the details of these methods and the results they have achieved using PSB1 in Section 2, but to summarize, many of these systems have improved performance and demonstrated new techniques.

When PSB1 was first introduced, the initial PsatC4F runs were able to solve 22 of the 29 problems, with an average success rate of 62/700 [15]. Some of the most drastic improvements have come on some of the most informative problems in PSB1, such as Double Letters ($11 \rightarrow 36$ successes between [17] and [15]), Regular Space with Newline ($11 \rightarrow 100$), 3-letters ($11 \rightarrow 64$), Vector Average ($10 \rightarrow 47$), and 3-Word Lines ($0 \rightarrow 93$).

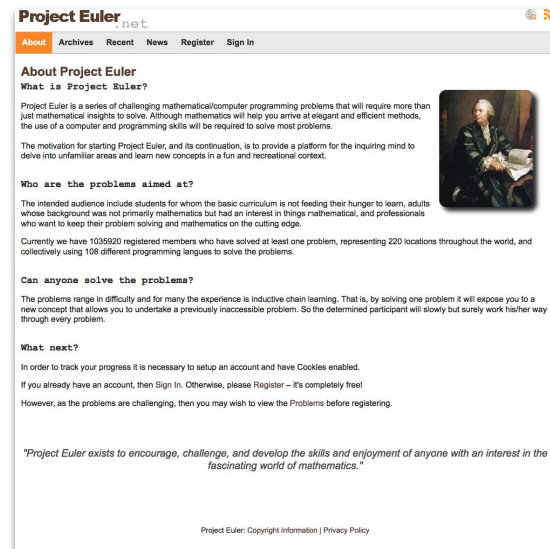
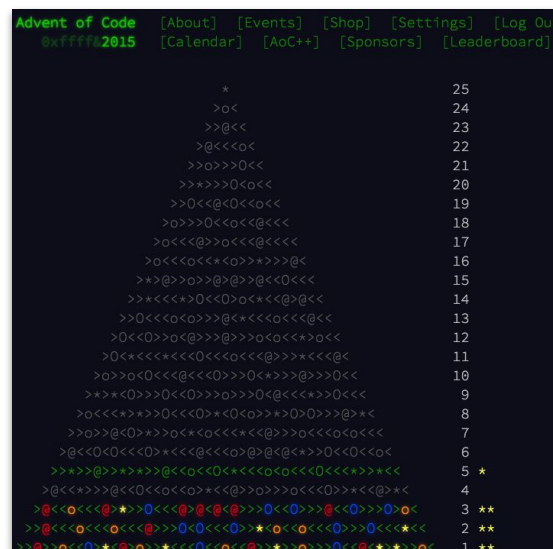
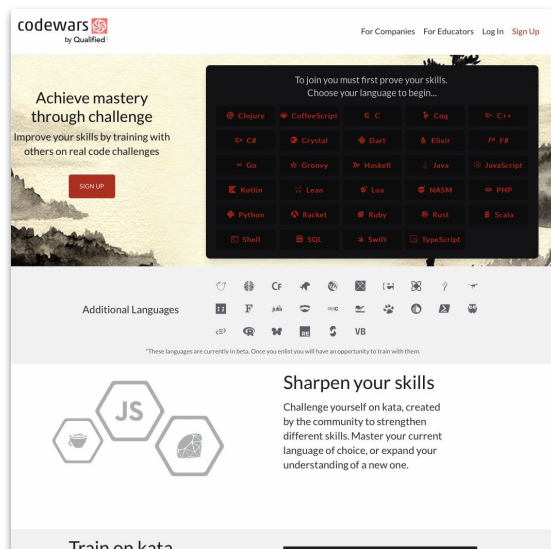
Thus, the PsatC4F and other synthesis systems, the problems of PSB1 have become less useful over time. In particular, the very

+ Need for New Benchmark Problems

- Synthesis systems improved on PSB1
- Need new, harder problems to drive research forward
 - **subjective** - problems we expected to be more difficult
 - **objective** - no problems that our GP system solved more than 65% of the time
 - most selected problems have success rates $< 10\%$
- Don't overfit on one set of problems

+ Problem Sources for PSB2

- **Code Wars** - user-created *coding kata*
- **Advent of Code** - advent calendar of coding problems
- **Project Euler** - archive of mathematical computer programming problems
- **Homework Problems** - from undergraduate classes



+ Problems in PSB2

25 benchmark problems

- Considered 75+ problems
- Implemented 50+ problems

Basement	Leaders
Bouncing Balls	Luhn
Bowling	Mastermind
Camel Case	Middle Character
Coin Sums	Paired Digits
Cut Vector	Shopping List
Dice Game	Snow Day
Find Pair	Solve Boolean
Fizz Buzz	Spin Words
Fuel Cost	Square Digits
GCD	Substitution Cipher
Indices of Substring	Twitter
Vector Distance	

Problem names

Input/Output types across all problems

Type	Input Count	Output Count
Integer	7	9
Float	2	5
String	10	7
Boolean	0	1
Vector of Integers	6	3
Vector of Floats	2	0

+ Example Problem: Spin Words

- 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
- Input: String
- Output: String
- Example:
`"Hello all at GECCO" → "olleH all at OCCEG"`

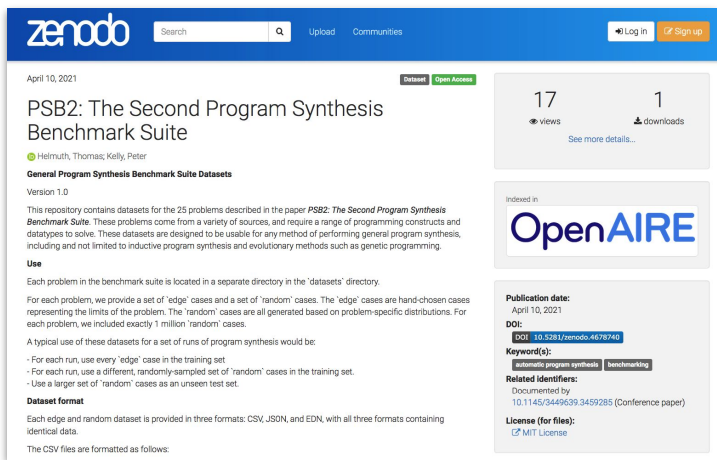
+ Example Problem: Vector Distance

- Given two n -dimensional vectors of floats, return the Euclidean distance between the two vectors in n -dimensional space
- Input: Two vectors of floats in n -dimensions
- Output: Float
- Example:
`[-77.6 -45.2 58.3] [6.5 48.6 17.5] → 132.4`

+ Using PSB2

- PSB2 emphasizes ease of use for wider adoption
- Datasets
 - sample training and test data for each run
- Python and Clojure libraries for easy sampling
 - Easy install right now: `$ pip install psb2`

Datasets and documentation link



<https://zenodo.org/record/4678739>

+ Experimental Methods

- PSB2 works with any program synthesis system
 - Experiments: PushGP
- Standardized error functions per output data type
 - integers outputs: absolute difference
 - string outputs: Levenshtein string edit distance
 - etc.
- Performance metric:
 - success rate
 - on unseen test set

```
(in1 in2 exec_do*vector_integer (output_integer1  
integer_sub in1 output_integer2  
vector_integer_occurrencesof exec_do*times  
exec_stackdepth integer_mod exec_yankdup  
vector_integer_yank in2))
```

Push solution to Find Pair

+ Results - Successes out of 100 runs

- 13 out of 25 problems solved at least once
- *Comparison with PSB1:*
 - 26 out of 29 problems solved

Problem	Successes
Basement	1
Bouncing Balls	0
Bowling	0
Camel Case	1
Coin Sums	2
Cut Vector	0
Dice Game	0
Find Pair	4
Fizz Buzz	25
Fuel Cost	50
GCD	8
Indices of Substring	0
Leaders	0

Problem	Successes
Luhn	0
Mastermind	0
Middle Character	57
Paired Digits	8
Shopping List	0
Snow Day	4
Solve Boolean	5
Spin Words	0
Square Digits	0
Substitution Cipher	61
Twitter	31
Vector Distance	0

+ Size of Smallest Simplified Solution Program

- After automatic simplification
- *Comparison with PSB1:*
 - 8 problems solved with < 9 instructions

Problem	Size
Basement	18
Camel Case	20
Coin Sums	33
Find Pair	16
Fizz Buzz	19
Fuel Cost	9
GCD	19
Middle Character	10
Paired Digits	15
Snow Day	11
Solve Boolean	18
Substitution Cipher	9
Twitter	22

+ Conclusions

- PSB2 - new suite of program synthesis benchmark problems
- Easier to use than PSB1
- Results with PushGP:
 - solutions to 13/25 problems
 - harder than PSB1

Get the datasets and
sampling libraries here!



<https://zenodo.org/record/4678739>