

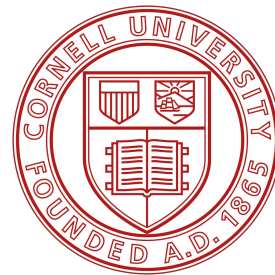
Non-Canonical Tasks in Synthesis and Learning

Kevin Ellis

Joint with Darren Key, Wen-Ding Li, Hao Tang

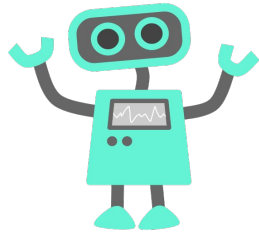
Cornell University

MAPS Symposium 2023, FSE



Program synthesizer

specification

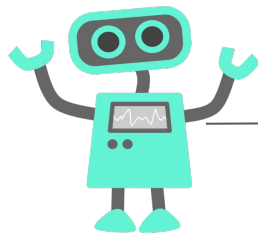


program

Specification = Dependent Types (types + logical predicates)

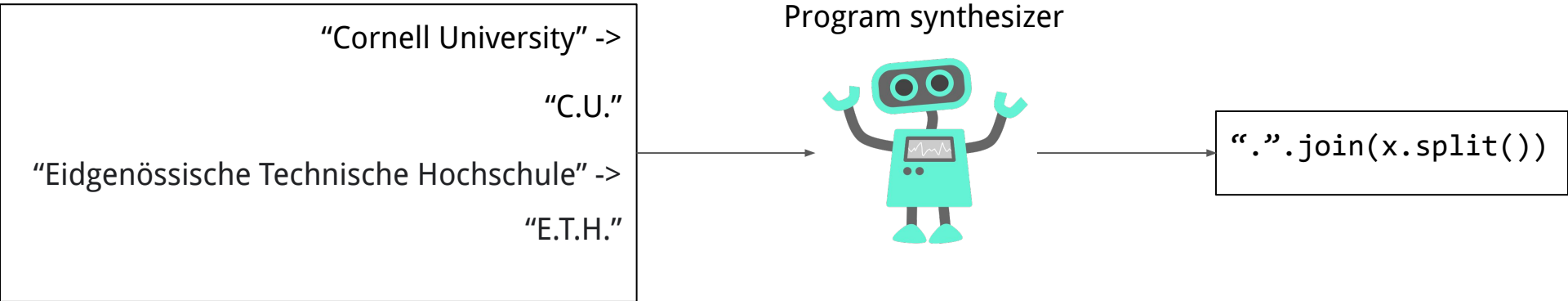
Program synthesizer

$n: \text{Nat} \rightarrow x: \alpha \rightarrow \{\text{List } \alpha \mid \text{len } \nu = n\}$ →



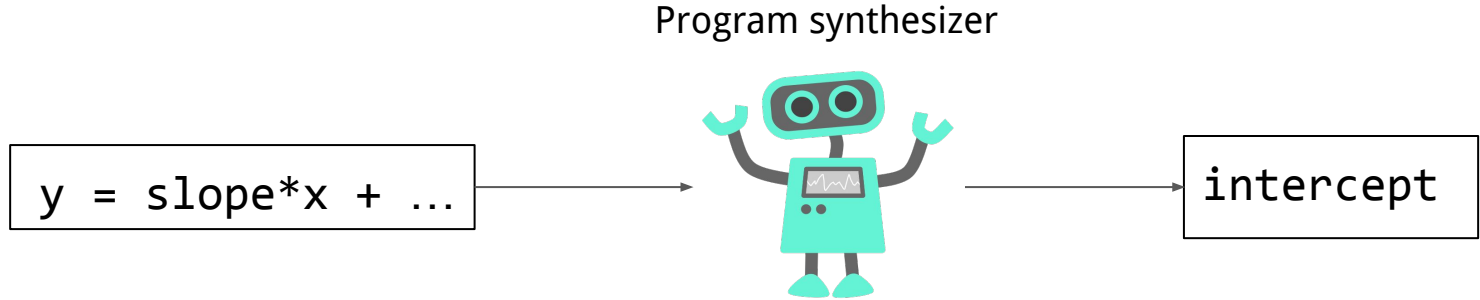
→ `replicate = λ n . λ x . if n ≤ 0
then Nil
else Cons x (replicate (dec n) x)`

Specification = Input-Outputs



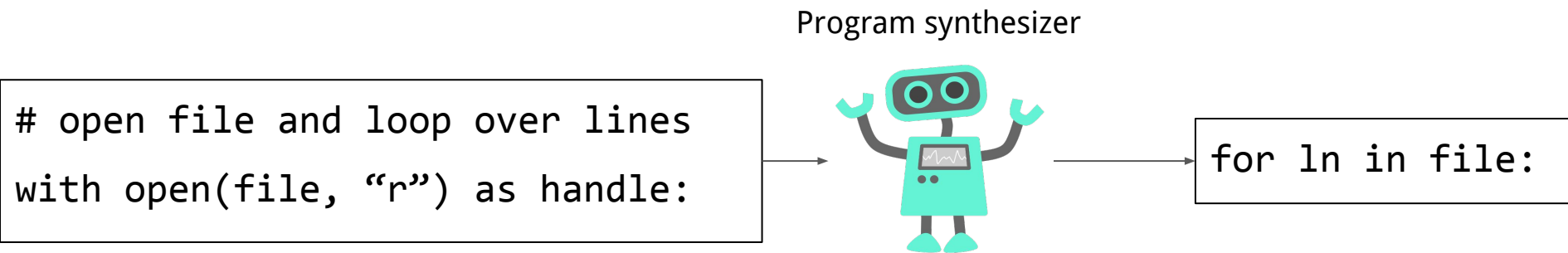
Eg, FlashFill. Gulwani 2012

Specification = Partially completed program



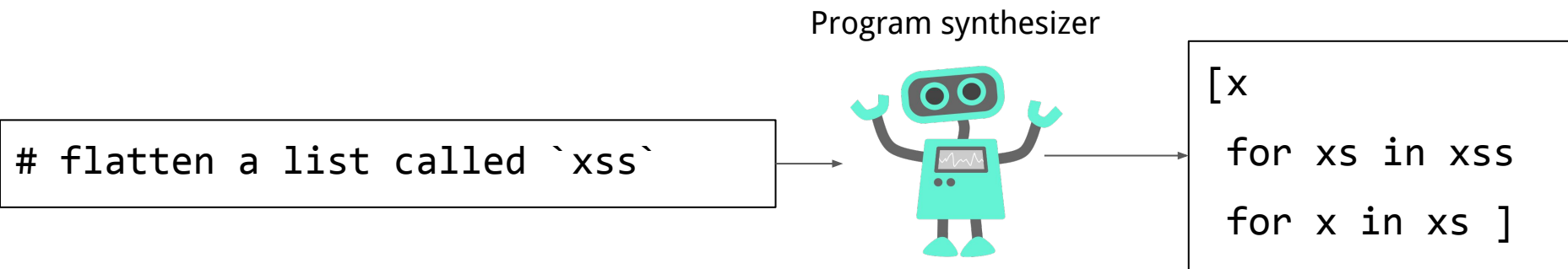
eg, Copilot, Codex

Specification = Partially completed program



eg, Copilot, Codex

Specification = Partially completed program



eg, Copilot, Codex

How we evaluate program synthesizers

Synthesizing Data Structure Transformations from Input-Output Examples *

	Name	Runtime	Runtime (no deduction)	Runtime (no types)	Expert examples	Random examples	Runtime (random)	Extra primitives	Description
Lists	add	0.04	0.05	3.87	5	4	0.04		Add a number to each element of a list.
	append	0.23	0.49	⊥	3	16	0.93		Append an element to a list.
	concat	0.13	0.22	68.95	5	23	0.20		Concatenate two lists together.
	dedup	231.05	⊥	⊥	7	-	-	member	Remove duplicate elements from a list.
	droplast	316.39	⊥	⊥	6	-	-	max	Drop the last element in a list.
	dropmax	0.12	0.19	77.05	3	7	0.16		Drop the largest number(s) in a list.
	dupli	0.11	0.86	378.35	3	5	0.20		Duplicate each element of a list.
	evens	7.39	45.52	⊥	5	8	30.08		Remove the odd numbers from a list.
	last	0.02	0.06	1.80	4	4	0.03		Return the last element in a list.
	length	0.01	0.14	41.36	4	5	0.04		Return the length of a list.
	max	0.46	9.53	⊥	7	8	8.19		Return the largest number in a list.

How we evaluate program synthesizers



2022-3-16

Approach	Validation Set				Test Set		
	10@1k	10@10k	10@100k	10@1M	10@1k	10@10k	10@100k
9B	16.9%	22.6%	27.1%	30.1%	14.3%	21.5%	25.8%
41B	16.9%	23.9%	28.2%	31.8%	15.6%	23.2%	27.7%
41B + clustering	21.0%	26.2%	31.8%	34.2%	16.4%	25.4%	29.6%

Table 5 | Solve rates of our best systems on the validation set and test set .

How we evaluate program synthesizers:

Success Rate, controlling for compute time

How we use program synthesizers

```
def sub_list(nums1 : list, nums2 : list) -> list:
    """
    Write a function to subtract two lists element-
    wise.
    """
```

How we use program synthesizers

```
def sub_list(nums1 : list, nums2 : list) -> list:
    """
    Write a function to subtract two lists element-
    wise.
    """
    return list(map(lambda x, y: x-y, nums1, nums2))
```



Evan Pu added a new photo.

November 16, 2021 · 🧑



copilot's buggy code suggestion (in gray) against the correct code (below)

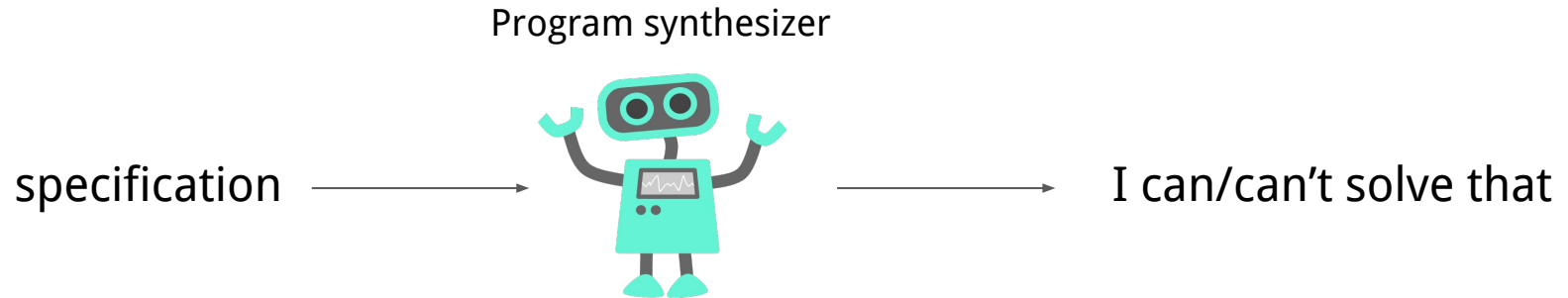
it is very subtle, but caused my search algorithm to bug out and invalidated 2 days worth of works

please use responsibly I guess is my take-away.

```
return rect_params, log_rect_prob
rect_params = t_rank[queue_ids[0]], b_rank[queue_ids[1]], l_rank[queue_ids[2]], r_rank[queue_ids[3]]
log_rect_prob = np.log(t[queue_ids[0]]) + np.log(b[queue_ids[1]]) + np.log(l[queue_ids[2]]) + np.log(r[queue_ids[3]])
return rect_params, log_rect_prob
log_rect_prob = np.log(t[rect_params[0]]) + np.log(b[rect_params[1]]) + np.log(l[rect_params[2]]) + np.log(r[rect_params[3]])
return rect_params, log_rect_prob
```

Task:

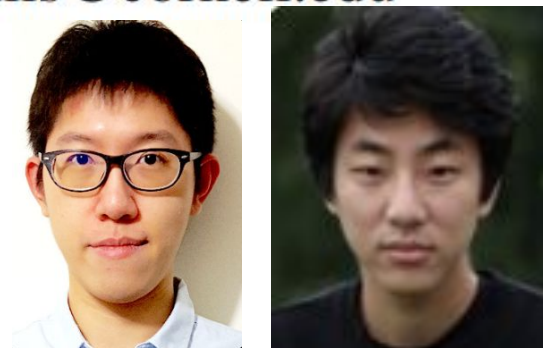
Predict whether the synthesizer should be trusted to solve a particular problem



Toward Trustworthy Neural Program Synthesis

Wen-Ding Li^{*1}, Darren Key^{*1}, Kevin Ellis¹

¹Department of Computer Science, Cornell University, USA
wl678@cornell.edu, dyk34@cornell.edu, kellis@cornell.edu



Trust in Traditional Program Synthesis

From Program Verification to Program Synthesis

Saurabh Srivastava

University of Maryland, College Park
saurabhs@cs.umd.edu

Sumit Gulwani

Microsoft Research, Redmond
sumitg@microsoft.com

Jeffrey S. Foster

University of Maryland, College Park
jfoster@cs.umd.edu

<pre>(a) Bresenhams(int X, Y) { v1:=2Y-X; y:=0; x:=0; while (x ≤ X) out[x]:=y; if (v1 < 0) v1:=v1+2Y; else v1:=v1+2(Y-X); y++; x++; return out; }</pre>	<pre>(b) Bresenhams(int X, Y) { []true → v'1=2Y-X ∧ y'=0 ∧ x'=0 while (x ≤ X) []v1 < 0 → out'=upd(out, x, y) ∧ v'1=v1+2Y ∧ y'=y ∧ x'=x+1 []v1 ≥ 0 → out'=upd(out, x, y) ∧ v'1=v1+2(Y-X) ∧ y'=y+1 ∧ x'=x+1 return out; }</pre> <p>(c) Invariant τ : $0 < Y \leq X \wedge v_1 = 2(x+1)Y - (2y+1)X \wedge 2(Y-X) \leq v_1 \leq 2Y \wedge \forall k : 0 \leq k < x \Rightarrow 2 out[k] - (Y/X)k \leq 1$</p> <p>Ranking function φ : $X - x$</p>
--	---

Figure 1. (a) Bresenham's line drawing algorithm (b) The invariant and ranking function that prove partial correctness and termination, respectively. (c) The algorithm written in transition system form, with statements as equality predicates, guarded appropriately.

Trust in Traditional Program Synthesis

program \vdash specification

Trust in *Neural* Program Synthesis

program \vdash specification (?)

Trust in Neural Program Synthesis

program \vdash natural language (?)

Trust in Neural Program Synthesis

program \vdash natural language (X)

Trust in Neural Program Synthesis

Neural network defines:

$\text{Pr}[\text{program} \mid \text{natural language}]$

The Trust Conundrum

Trust ~ Verification

program \vdash specification

My specification is informal...

...because train data is messy natural code

And I can't verify against an informal specification

How do people build trust?

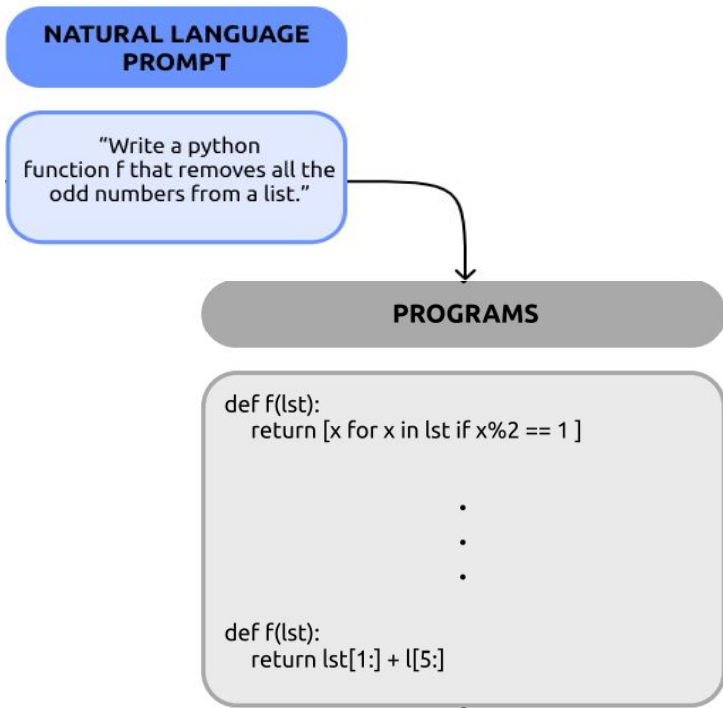
```
6     def test_admin_name_is_string(self):
7
8     class grafanaTests(unittest.TestCase):
9         """Tests for GeekTechStuff Grafana API Python"""
10
11         def test_admin_name_is_string(self):
12             admin_username = main.get_username()
13             self.assertIs(type(admin_username), str)
14
15         def test_admin_password_is_string(self):
16             admin_password = main.get_password()
17             self.assertIs(type(admin_password), str)
18
19         def test_grafana_url_is_string(self):
20             grafana_url = main.get_url()
21             self.assertIs(type(grafana_url), str)
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Speculyzer

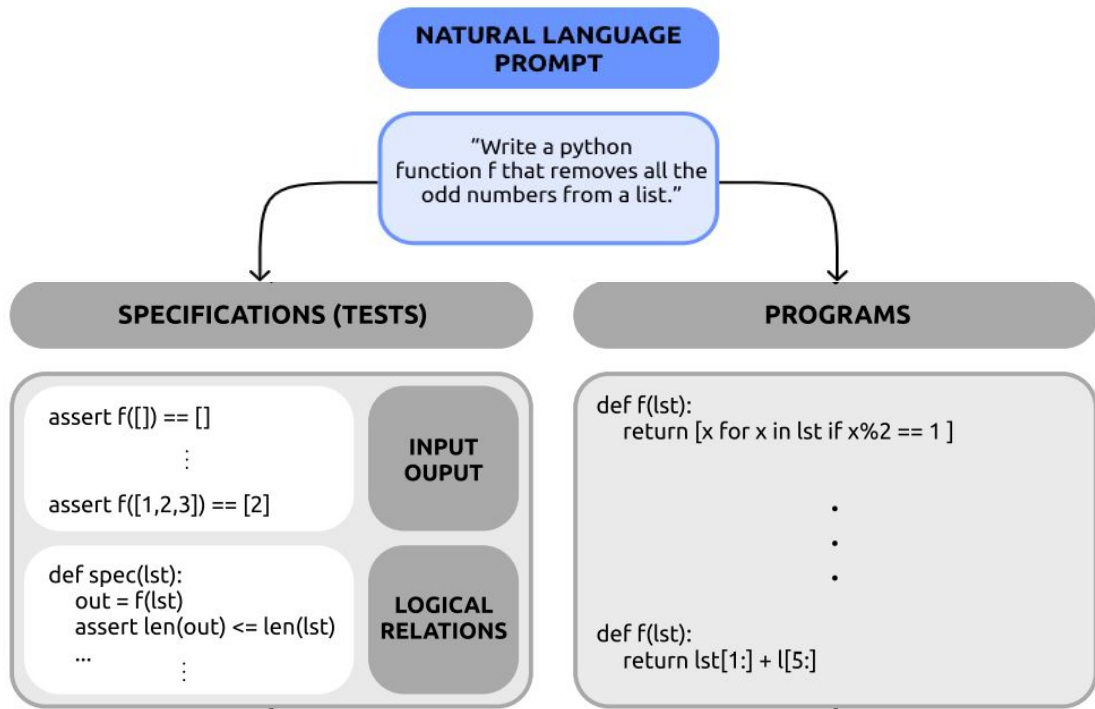
NATURAL LANGUAGE PROMPT

"Write a python
function f that removes all the
odd numbers from a list."

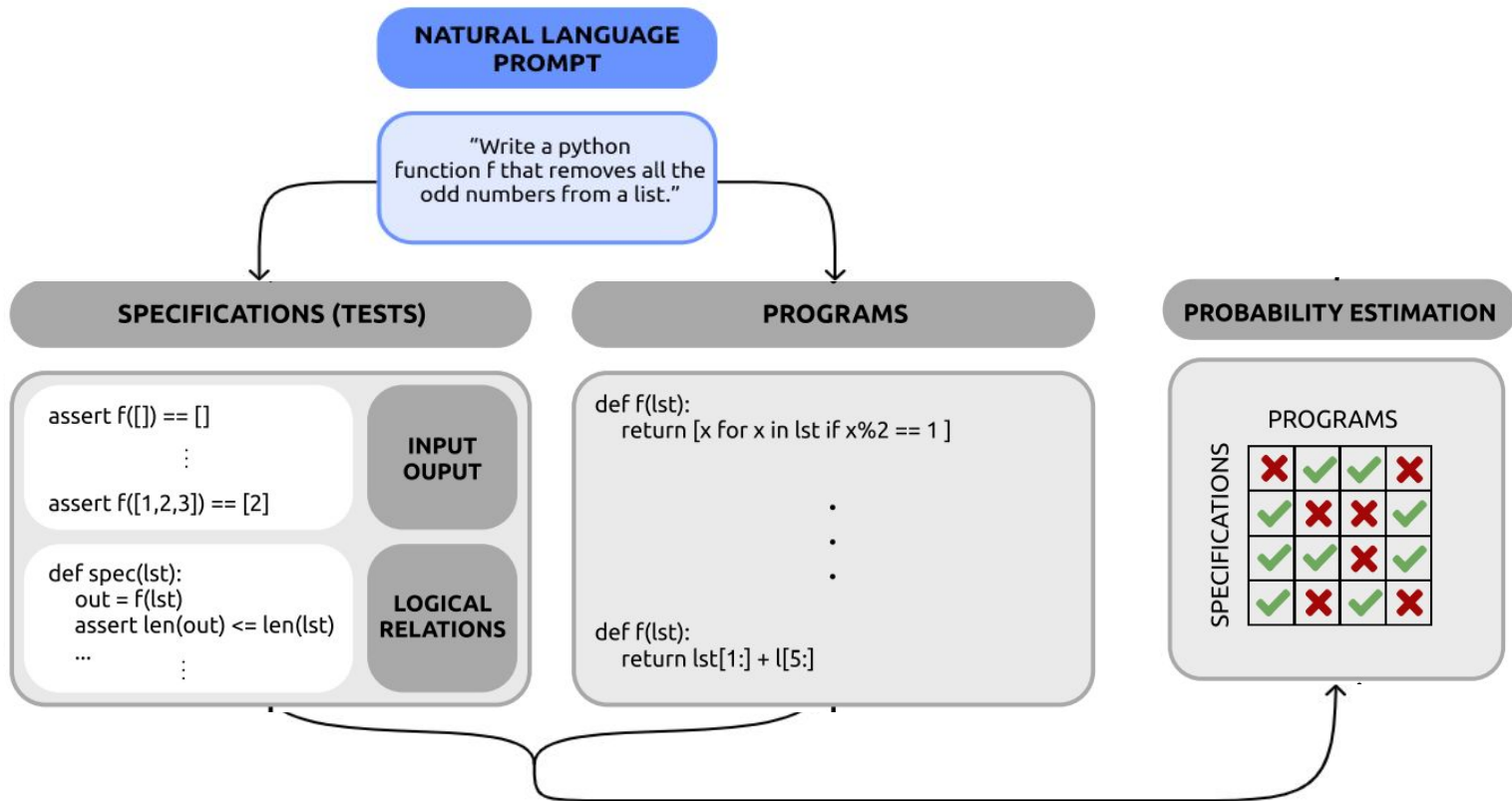
Speculyzer



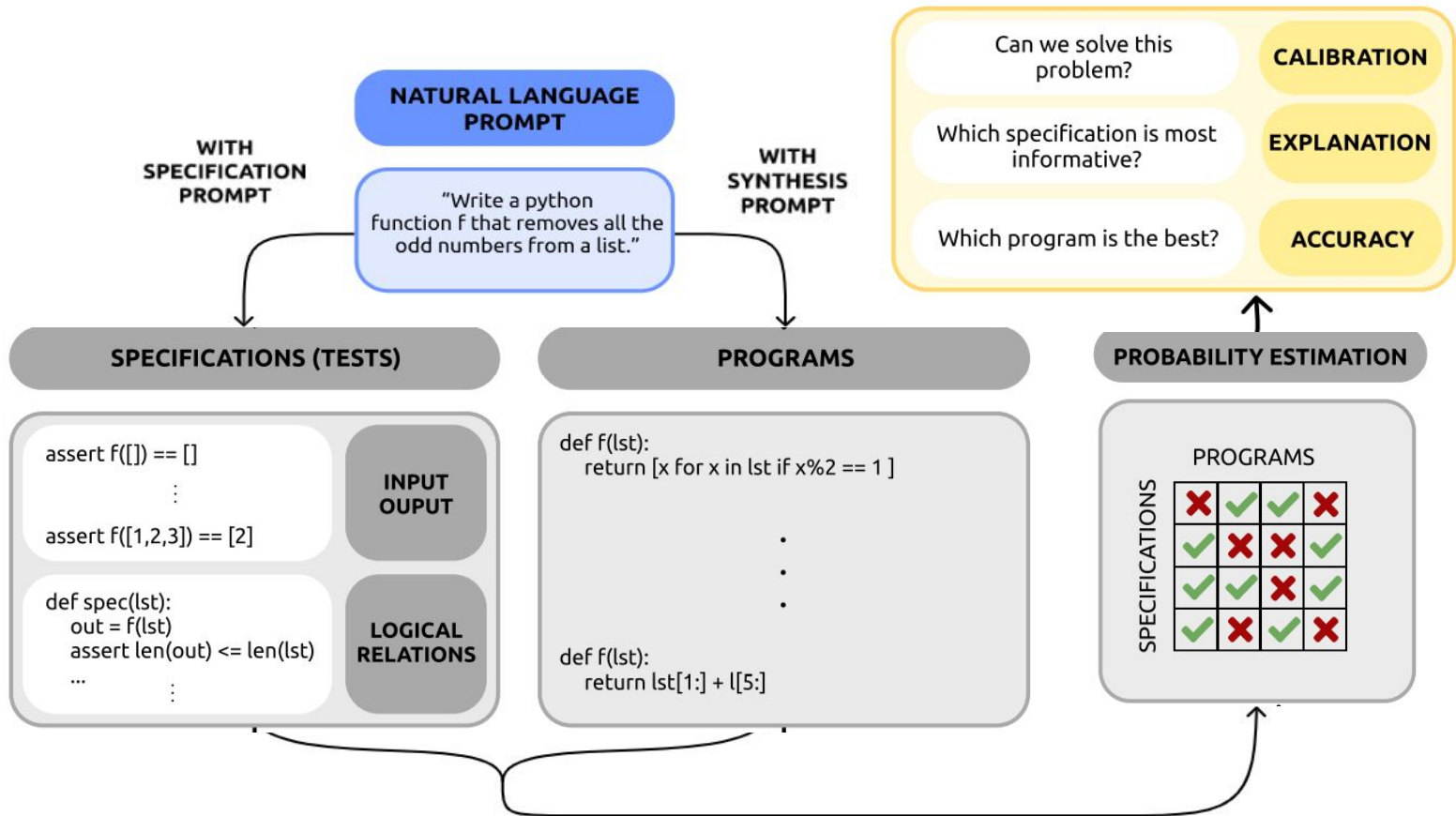
Speculyzer



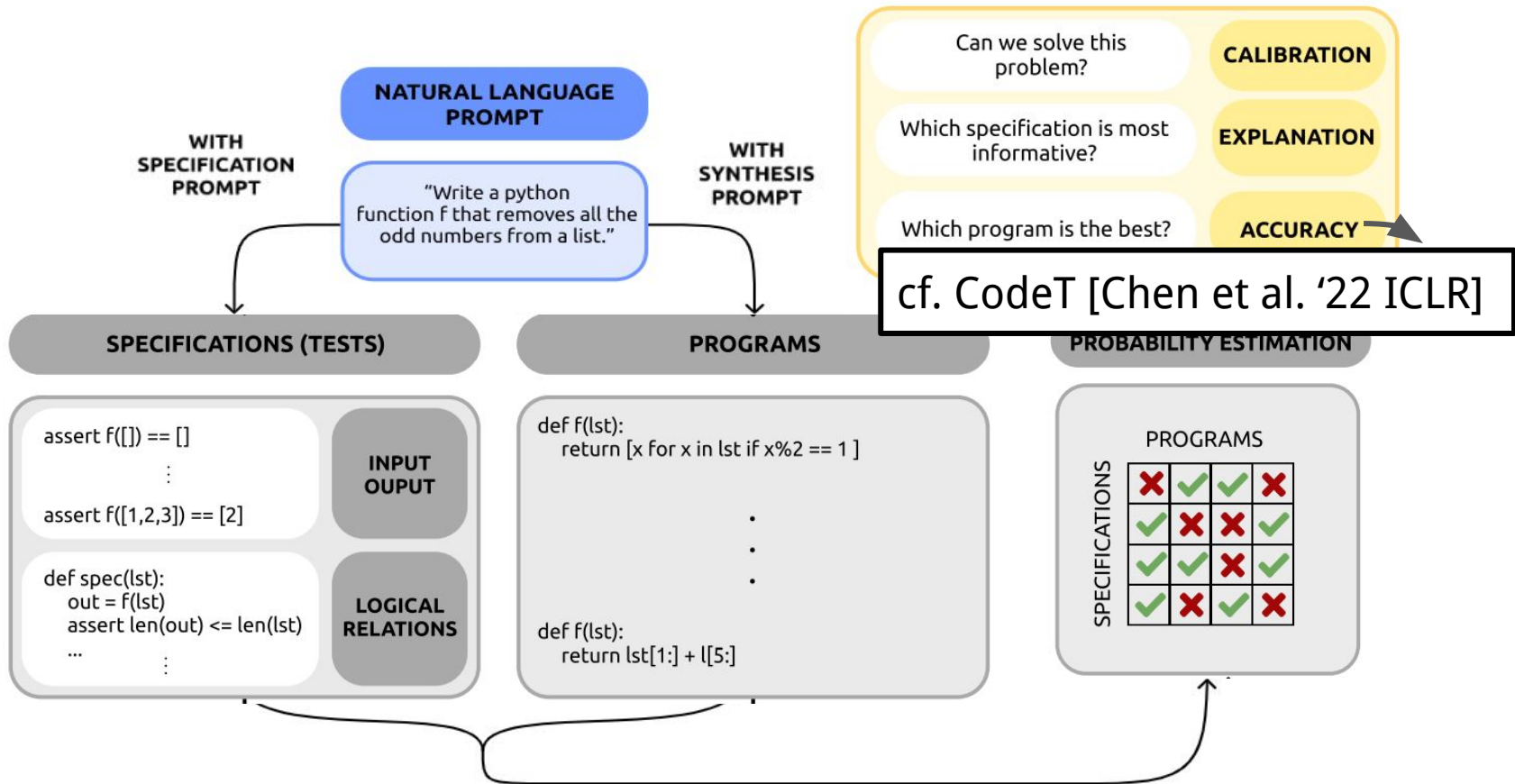
Speculyzer



Speculyzer



Speculyzer





$\text{Pr}[\text{prog correct}]$



- + Satisfies many specs?
- + Many other progs satisfy same specs?
- Lots of different prog behaviors?
- Low logits?



Pr[prog correct]

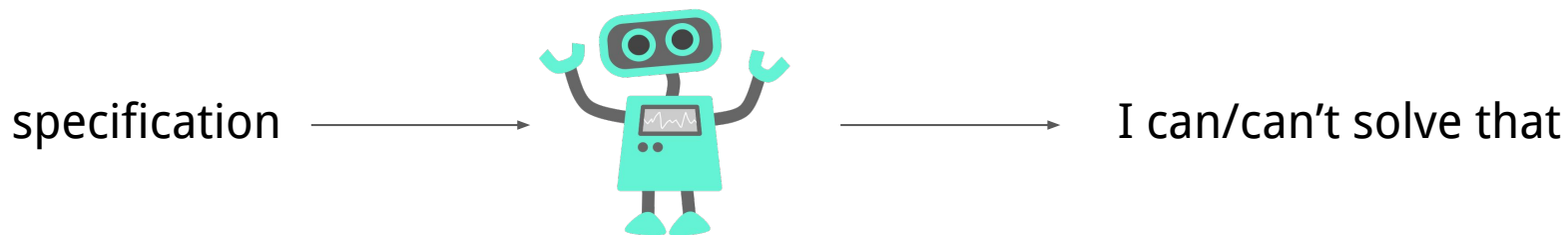


Feature extraction

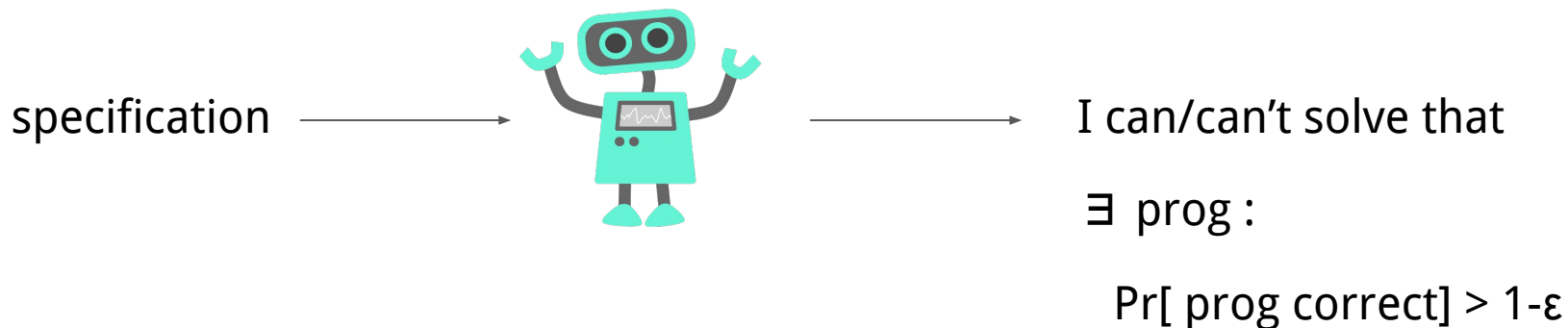


Binary classification
(logistic regression)

Speculyzer can decline to solve problems when uncertain



Speculyzer can decline to solve problems when uncertain



Declining to solve problems when uncertain



Steps toward formalization

ϵ -correctness $\epsilon \in [0, 1]$

$$\mathbb{P}(p(I) \neq \perp \wedge p(I) \not\models \phi) \leq \epsilon$$

δ -coverage $\delta \in [0, 1]$

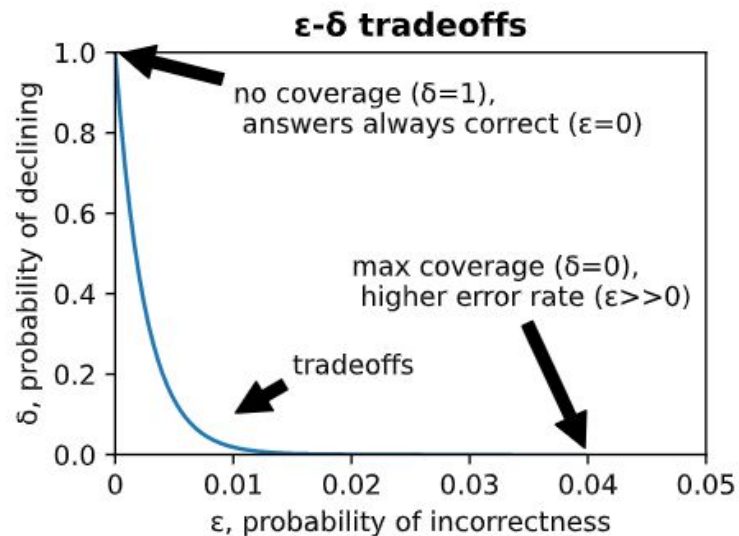
$$\mathbb{P}(p(I) = \perp) \leq \delta$$

p : synthesized program

I : user's input intent

ϕ : formal specification

\perp : decline to solve problem



Reality Check:

No neural net english->code system will achieve 0% errors

Task:

Generate tests that explain the generated code, so if it is incorrect, the user can discard it

PROGRAM

```
def is_bored(S):  
    """You'll be given a string of words, and your task is to count the number  
    of boredoms. A boredom is a sentence that starts with the word "I".  
    Sentences are delimited by '.', '?' or '!'.  
    For example:  
    >>> is_bored("Hello world")  
    0  
    >>> is_bored("The sky is blue. The sun is shining. I love this weather")  
    1"""
```

How to communicate what program does so that the user can accept/reject it?

program \vdash specification

Out-of-the-box LLM

PROGRAM

```
def is_bored(S):
    """You'll be given a string of words, and your task is to count the number
    of boredoms. A boredom is a sentence that starts with the word "I".
    Sentences are delimited by '.', '?' or '!'.
    For example:
    >>> is_bored("Hello world")
    0
    >>> is_bored("The sky is blue. The sun is shining. I love this weather")
    1"""
    boredoms = 0
    # replace . or ! or ? with . to simplify this problem
    S = S.replace('.', '. ')
    S = S.replace('!', '! ')
    S = S.replace('?', '? ')
    sentences = S.split(' ')
    for sentence in sentences:
        if sentence.startswith('I'): boredoms = boredoms + 1
    return boredoms
```


Out-of-the-box LLM

PROGRAM

```
def is_bored(S):
    """You'll be given a string of words, and your task is to count the number
    of boredoms. A boredom is a sentence that starts with the word "I".
    Sentences are delimited by '.', '?' or '!'.
    For example:
    >>> is_bored("Hello world")
    0
    >>> is_bored("The sky is blue. The sun is shining. I love this weather")
    1"""
    boredoms = 0
    # replace . or ! or ? with . to simplify this problem
    S = S.replace('.', '. ')
    S = S.replace('!', '! ')
    S = S.replace('?', '? ')
    sentences = S.split(' ')
    for sentence in sentences:
        if sentence.startswith('I'): boredoms = boredoms + 1
    return boredoms
```

**RANDOM
INPUT/
OUTPUT**

```
assert is_bored("I love this weather.") == 1
```

How to communicate what program does so that the user can accepted/reject it?

$$\operatorname{argmax}_{\text{prog} \vdash \text{spec}} \text{objective_function}(\text{prog}, \text{spec})$$

Pick the thing which is true about the program
But which is not true about most other programs
“distinguishing”, “selective”

PROGRAM

```
def is_bored(S):
    """You'll be given a string of words, and your task is to count the number
    of boredoms. A boredom is a sentence that starts with the word "I".
    Sentences are delimited by '.', '?' or '!'.
    For example:
    >>> is_bored("Hello world")
    0
    >>> is_bored("The sky is blue. The sun is shining. I love this weather")
    1"""
    boredoms = 0
    # replace . or ! or ? with . to simplify this problem
    S = S.replace('.', '. ')
    S = S.replace('!', '! ')
    S = S.replace('?', '? ')
    sentences = S.split(' ')
    for sentence in sentences:
        if sentence.startswith('I'): boredoms = boredoms + 1
    return boredoms
```

**RANDOM
INPUT/
OUTPUT**

```
assert is_bored("I love this weather.") == 1
```

PROGRAM

```
def is_bored(S):
    """You'll be given a string of words, and your task is to count the number
    of boredoms. A boredom is a sentence that starts with the word "I".
    Sentences are delimited by '.', '?' or '!'.
    For example:
    >>> is_bored("Hello world")
    0
    >>> is_bored("The sky is blue. The sun is shining. I love this weather")
    1"""
    boredoms = 0
    # replace . or ! or ? with . to simplify this problem
    S = S.replace('.', '. ')
    S = S.replace('!', '! ')
    S = S.replace('?', '? ')
    sentences = S.split(' ')
    for sentence in sentences:
        if sentence.startswith('I'): boredoms = boredoms + 1
    return boredoms
```

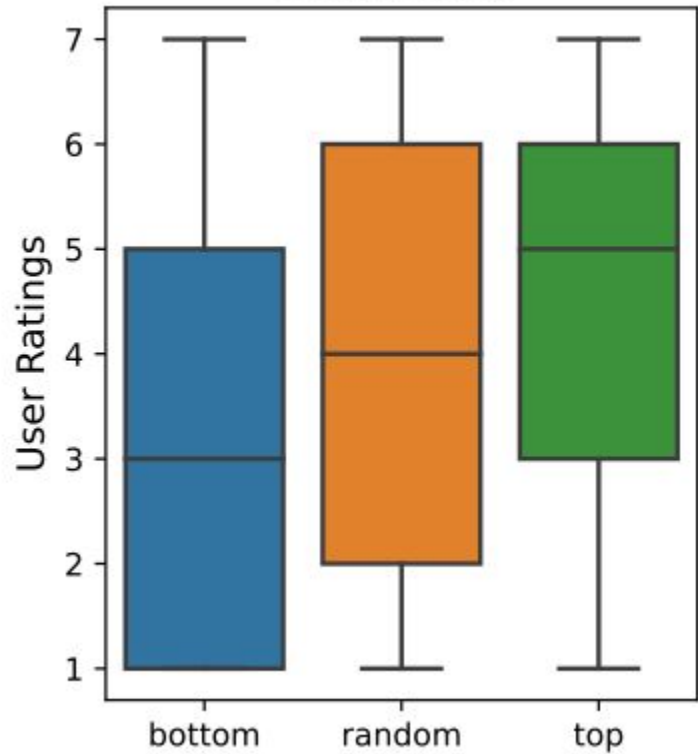
**DISTINGUISHING
INPUT/
OUTPUT**

```
assert is_bored("I have no idea what I'm doing") == 2
```

**RANDOM
INPUT/
OUTPUT**

```
assert is_bored("I love this weather.") == 1
```

User Study



Worst of 100 LLM samples

Raw LLM

Top of 100 LLM samples

Speculyzer Recap

Task: Predict if the synthesizer should be trusted

Task: Generate informative test(s) that explain the generated code

What could trust unlock?

[ellisk42/ec] Bump protobuf from 3.8.0 to 3.15.0 (PR #90) Σ Inbox x

dependabot[bot] <notifications@github.com> [Unsubscribe](#)
to ellisk42/ec, Subscribed ▼

This automated pull request fixes a [security vulnerability](#) (high severity).

[Learn more about Dependabot security updates.](#)

Bumps [protobuf](#) from 3.8.0 to 3.15.0.

Release notes

Sourced from [protobuf's releases](#).

<science_fiction>

This repository Search or type a command Explore Gist Blog Help dsbaars + - X

dsbaars / Project-XChain Watch 6 Unstar 1 Fork 0

Browse Issues Milestones New Issue

Everyone's Issues 62 0 Open 62 Closed Sort: Newest

Assigned to you 0

Created by you 26

Mentioning you 0

No milestone selected

Labels

- bug 37
- critical 8
- enhancement 10
- high priority 26
- invalid 5
- low priority 2
- medium priority 27

Reopen Label Assignee Milestone

- Force currency-numbers in negotiation table** bug critical high priority verified #62
Opened by dsbaars 6 days ago 1 comment
- After clicking "I have read it", disable button** enhancement medium priority verified #61
Opened by dsbaars 6 days ago 1 comment
- Create Game, make duration field required** bug medium priority verified #60
Opened by dsbaars 6 days ago 1 comment
- profit forecast calculation bug** bug high priority invalid verified #59
Opened by fapthohanded 6 days ago 3 comments
- Layout issues 2 and 2 point decimal rounding** bug usability #58
Opened by fapthohanded 6 days ago 1 comment
- Chat does not scroll down when watching game as chairman** bug medium priority verified #57
Opened by dsbaars 6 days ago 8 comments
- 'Number of games' editable in 'Change settings'** bug medium priority verified #56


</science_fiction>

<science_fiction>

Add toString implementation #17





[Edit](#)

 **Open** maxjacobson wants to merge 1 commit into `master` from `fix-git-tag-descriptions`

 Conversation **0**  Commits **1**  Files changed **3**

Changes from all commits ▾ Jump to... ▾ **+20 -4** 

[Unified](#)[Split](#)[Review changes ▾](#)

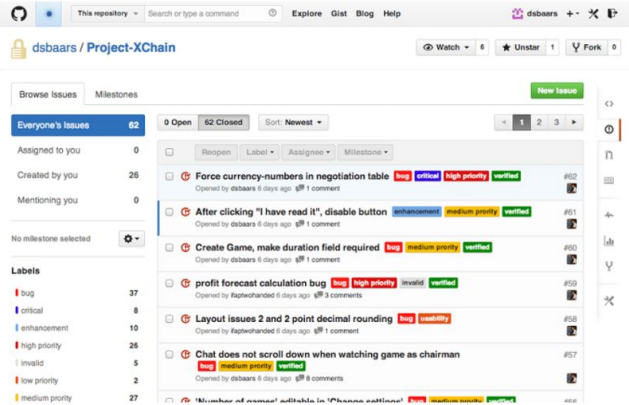
20  src/main/java/com/github/koraktor/mavanagaiata/git/GitTagDescription.java 88.24% cov | 8  [View](#)   ▾

issues

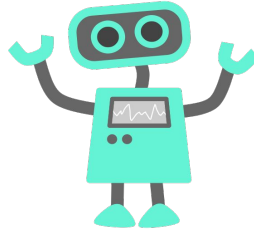
✚	@@ -67,9 +67,23 @@ public boolean isTagged() {	
67	67	*
68	68	* @return The string representation of this description
69	69	*/
70	-	@Override
71	-	public String toString() {
72	-	return "TODO: implement this method";
	▲	Method `toString` has a Cognitive Complexity of 7 (exceeds 5 allowed). Consider refactoring. ...
70	+	@Override
71	+	public String toString() {
72	+	if (this.nextTag == null) {
73	+	return this.abbreviatedCommitId;
74	+	} else if (this.distance == 0) {
75	+	return this.nextTag.getName();
76	+	} else {
	▲	'&&' should be on a new line. ...

</science_fiction>

<science_fiction>

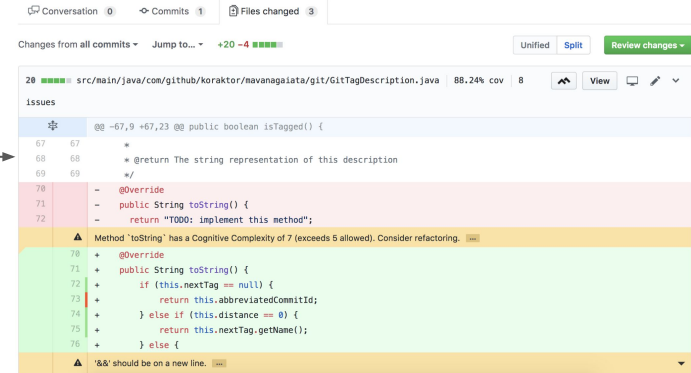


A screenshot of a GitHub repository page for 'dsbaars / Project-XChain'. The page shows a list of issues with various labels such as 'bug', 'critical', 'high priority', 'verified', 'enhancement', 'medium priority', 'invalid', 'low priority', and 'medium priority'. The issues are sorted by 'Newest' and include details like the issue number, title, and status.



Add toString implementation #17

maxjacobson wants to merge 1 commit into master from fix-git-tag-descriptions



A screenshot of a GitHub pull request page for 'Add toString implementation #17'. The page shows a diff view of the code changes, including a method implementation for 'toString' in 'GitTagDescription.java'. The diff highlights the changes in green and red. The code includes a method signature, a TODO comment, and the implementation logic.

```
@@ -67,9 +67,23 @@ public boolean isTagged() {
67 67      *
68 68      * @return The string representation of this description
69 69      */
70 - @Override
71 - public String toString() {
72 -     return "TODO: implement this method";
73 +
74 +     @Override
75 +     public String toString() {
76 +         if (this.nextTag == null) {
77 +             return this.abbreviatedCommitId;
78 +         } else if (this.distance == 0) {
79 +             return this.nextTag.getName();
80 +         } else {
```

</science_fiction>

Task:

Neurosymbolic programming for perception and graphics

Joint work with Hao Tang!
ICML '23

From Perception to Programs: Regularize, Overparameterize, and Amortize

Hao Tang¹ Kevin Ellis¹



Domain Specific Languages, Pre-provided Symbols

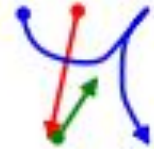


Neural Program
Generator

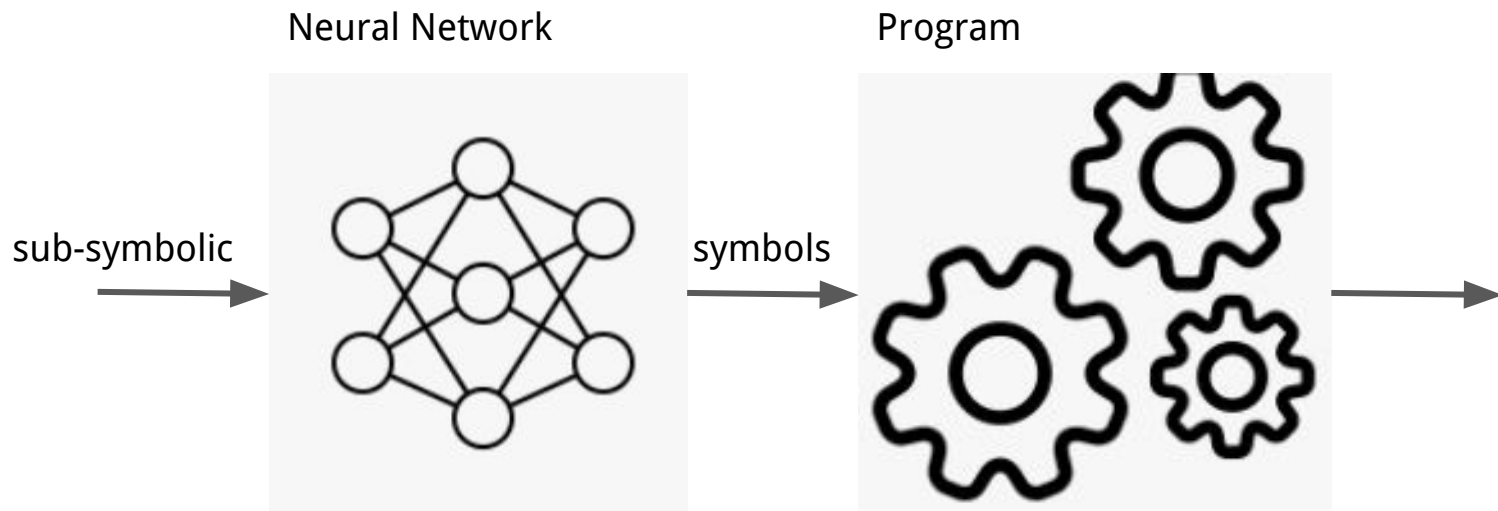
Neural Program
Executor

```
Draw("Top", "Circle", position, geometry)
for(i < 2, "translation", a)
  for(j < 2, "translation", b)
    Draw("Leg", "Cub", position + i*a + j*b, geometry)
for(i < 2, "translation", c)
  Draw("Layer", "Rec", position + i*c, geometry)
```

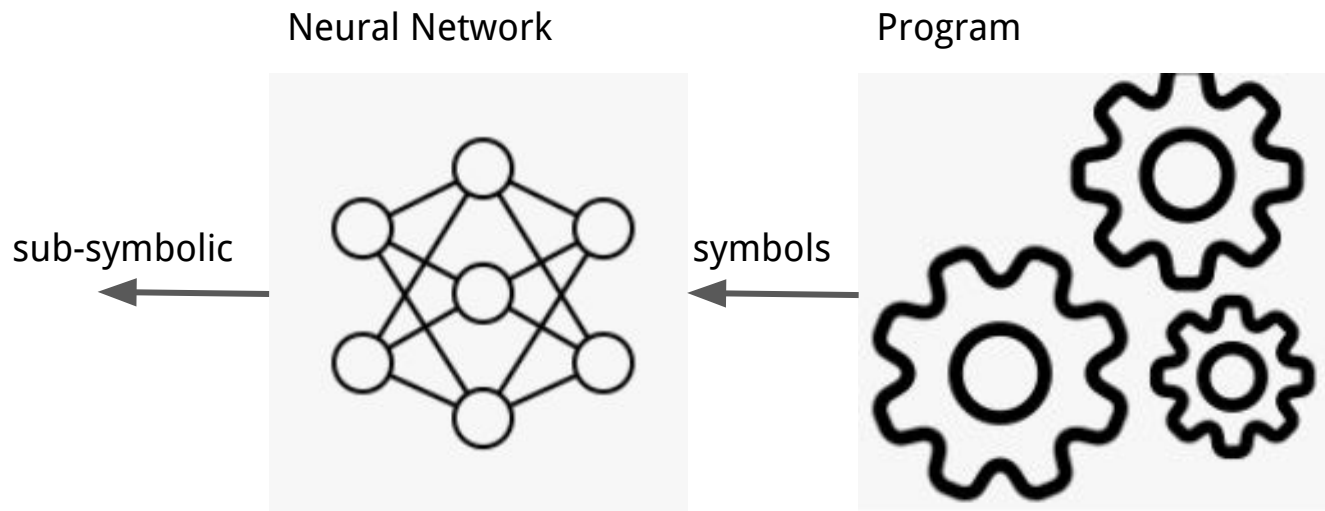
Flexible Symbols



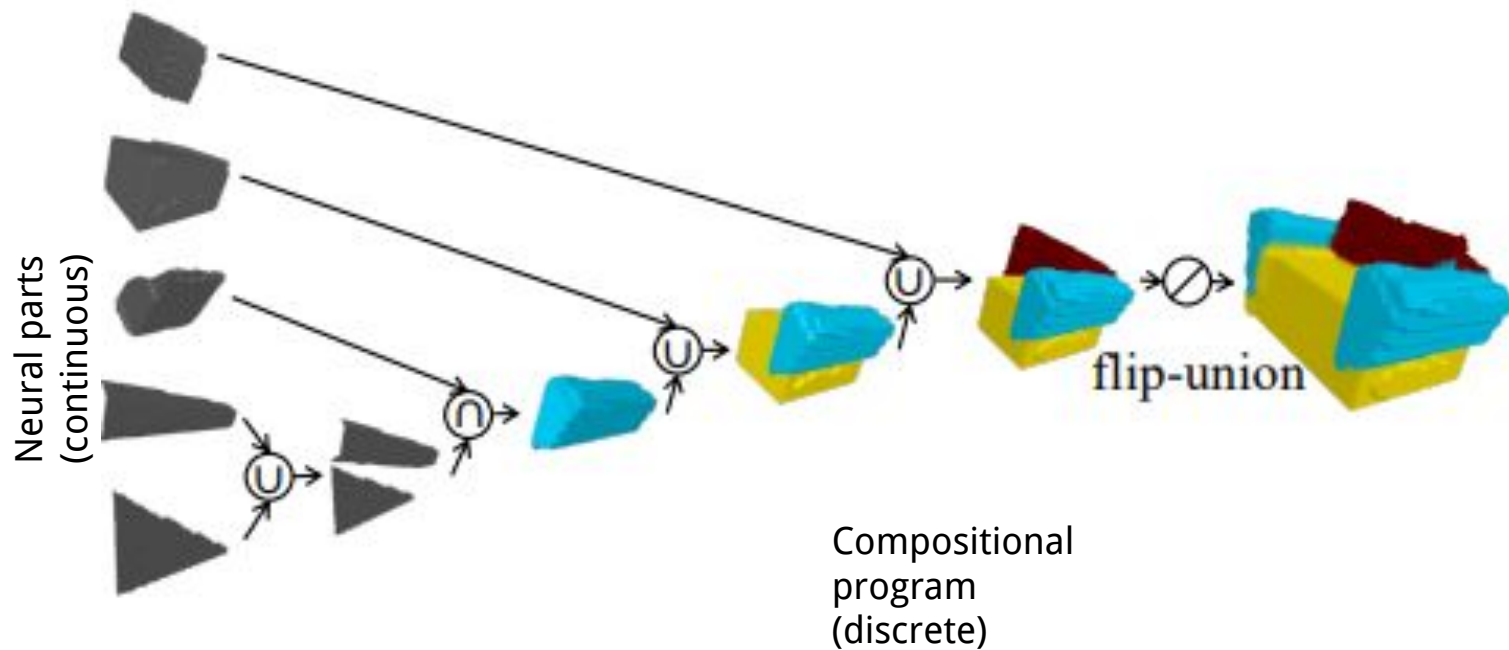
Flexible, Learnable Symbols, AND Powerful, Program-Like Compositional Processing?



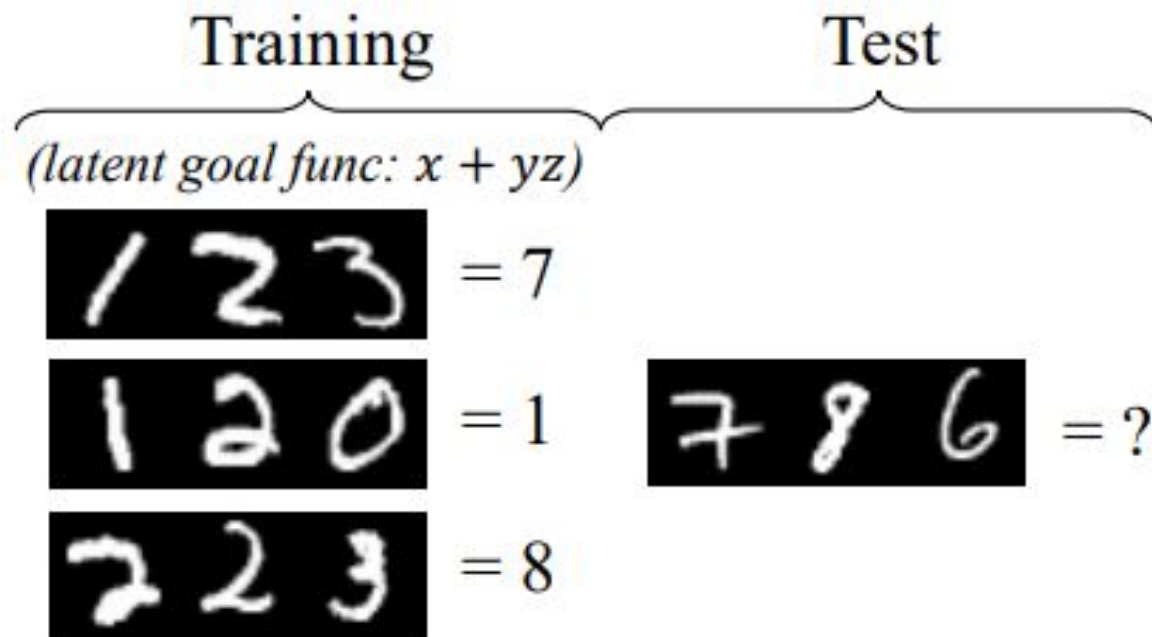
Flexible, Learnable Symbols, AND Powerful, Program-Like Compositional Processing? (generative)



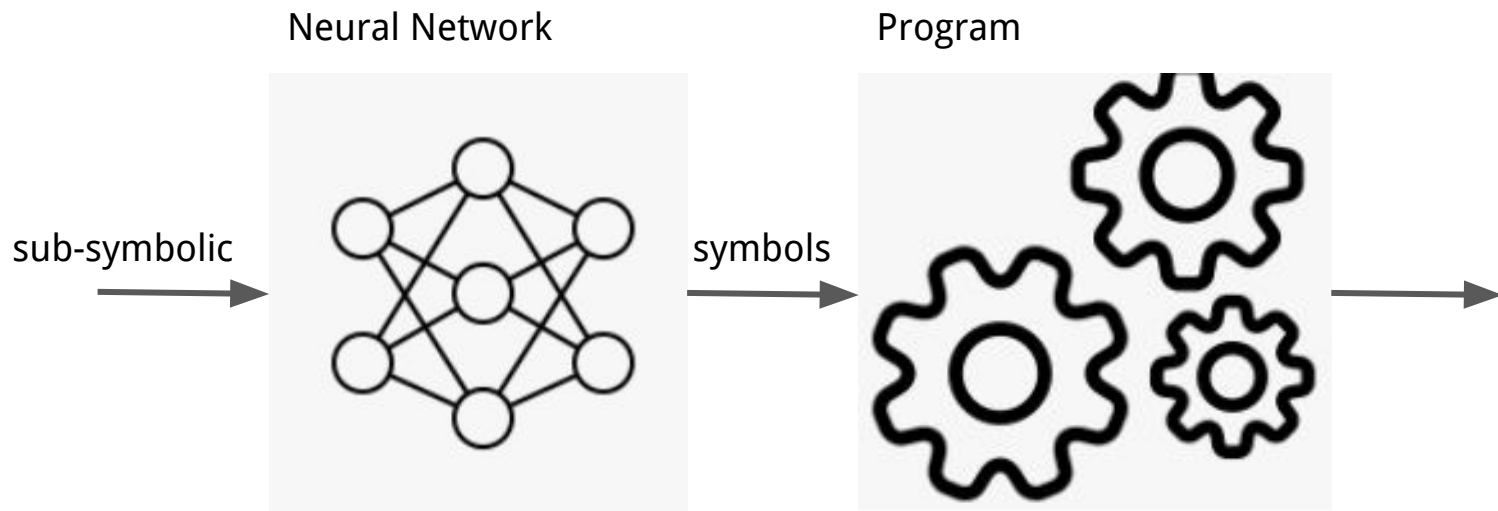
Generative Example



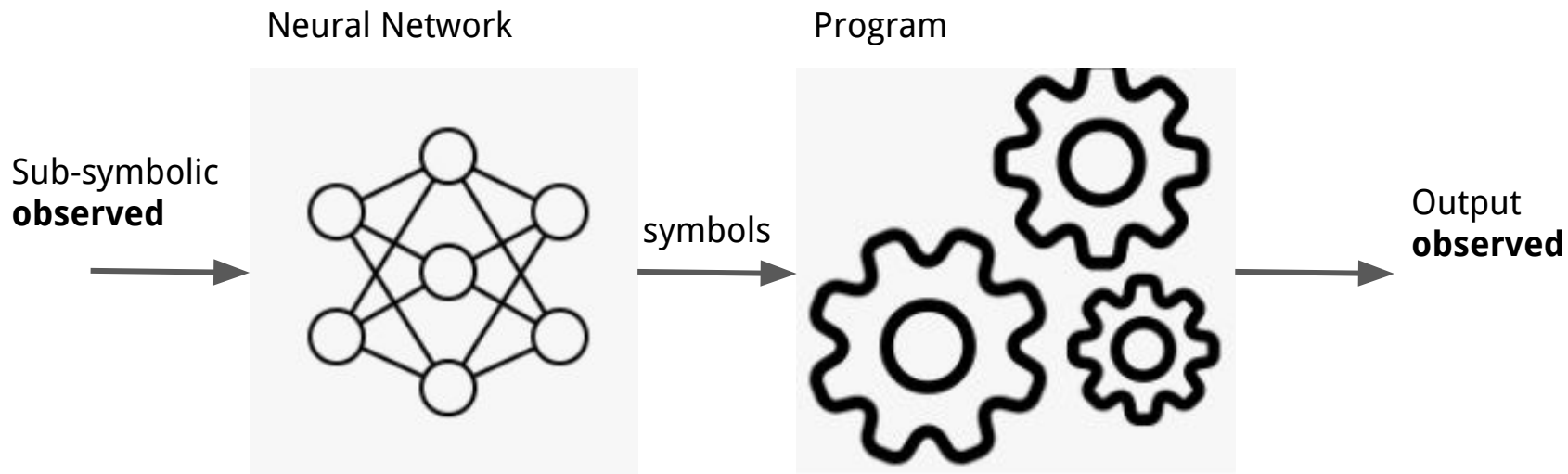
Discriminative Example



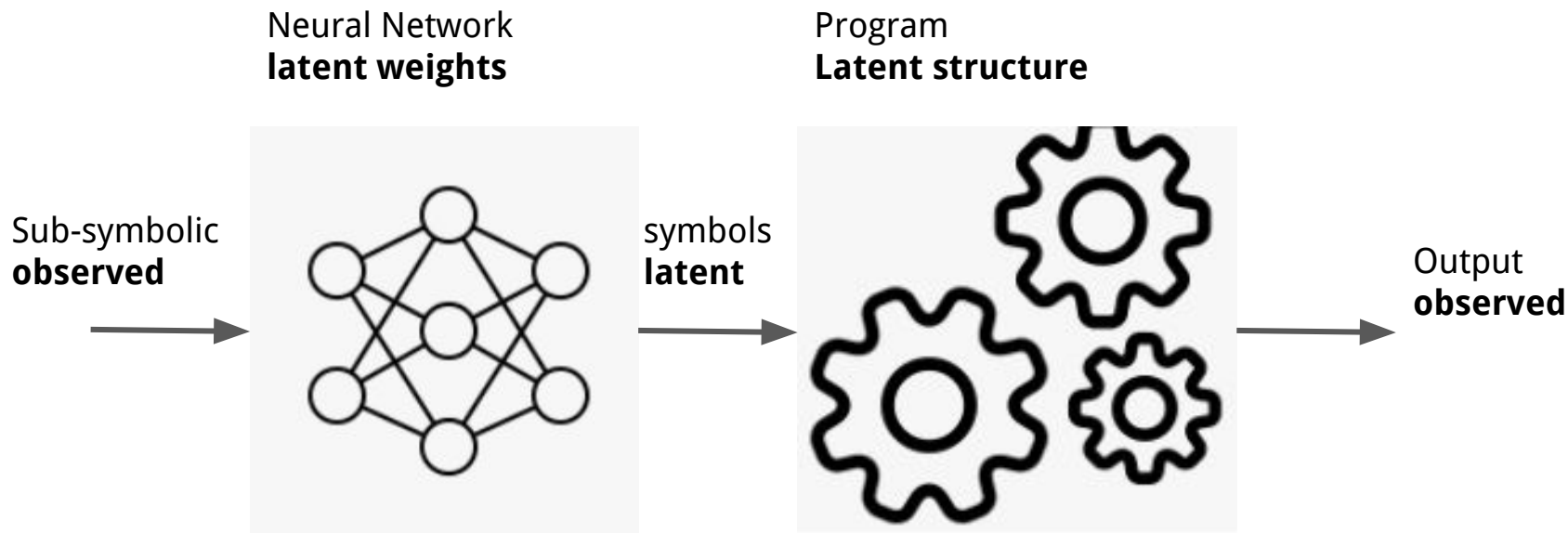
Learning Problem



Learning Problem



Learning Problem

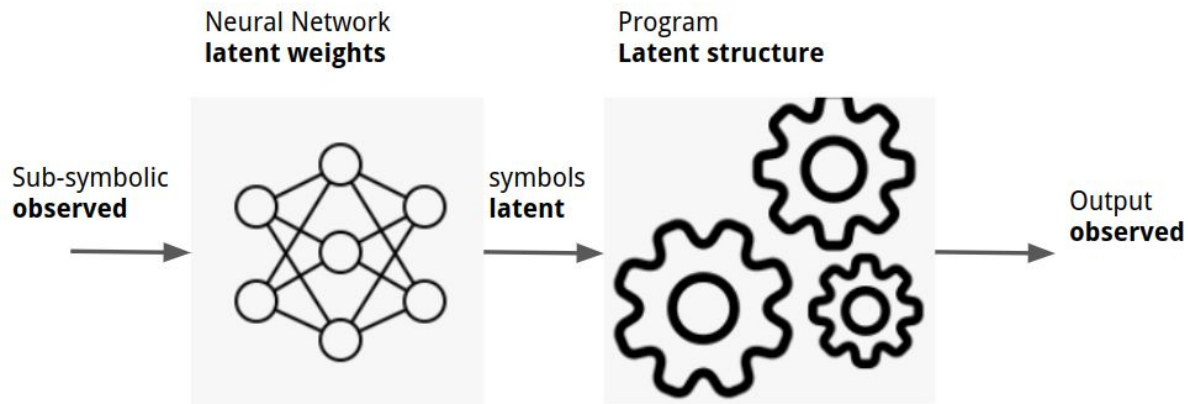


Challenges

Mixed Discrete-Continuous

Underconstrained

Symbol Grounding Problem

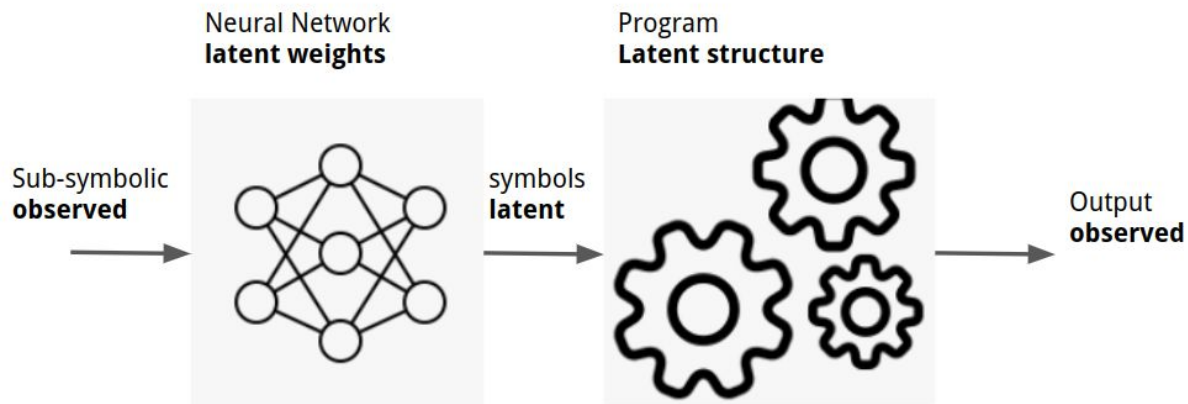


Opportunities

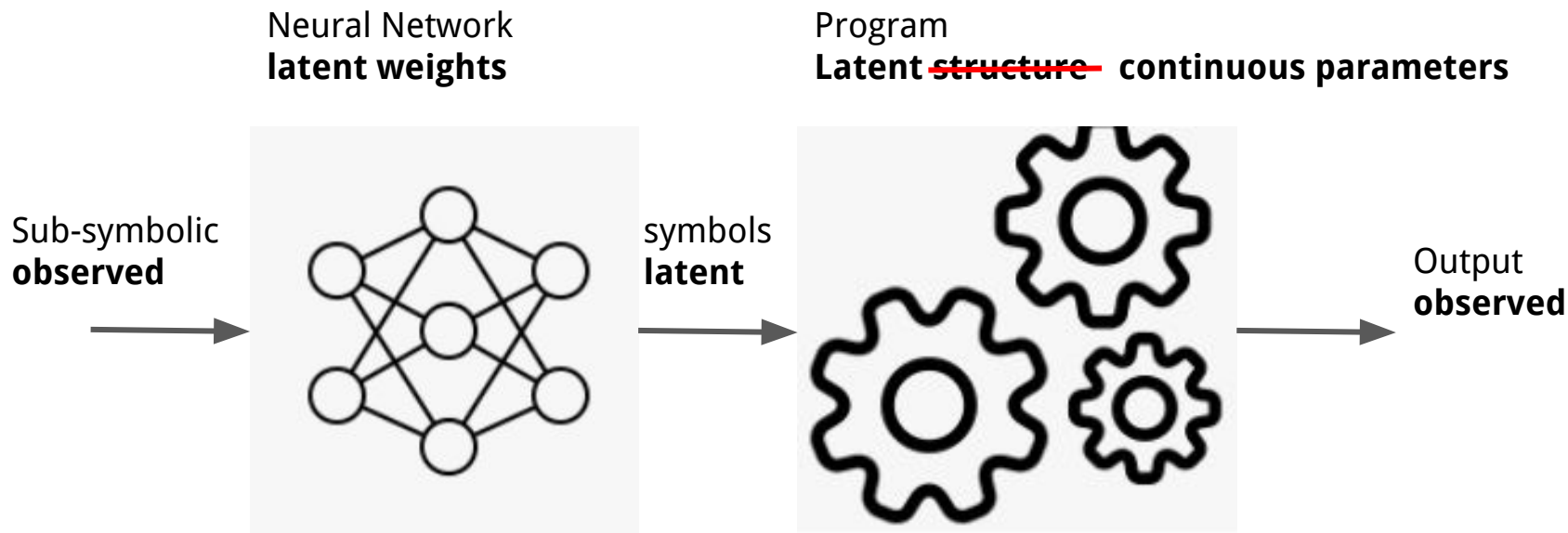
Compositional generalization from raw perceptual input

Systems that learn their own symbols!

Synthesizing neurosymbolic programs



Guiding Principles: Make Everything Continuous. Use Gradient Descent.



Relaxing discrete program spaces

Many approaches:

Terpret [Gaunt et al 2016]

dILP [Evans et al 2017]

EQLNet [Sahoo et al 2018]

DiffLog [Si et al 2019]

...

ROAP, this work [Tang et al 2023]

Relaxing discrete program spaces

Many approaches:

Terpret [Gaunt et al 2016]

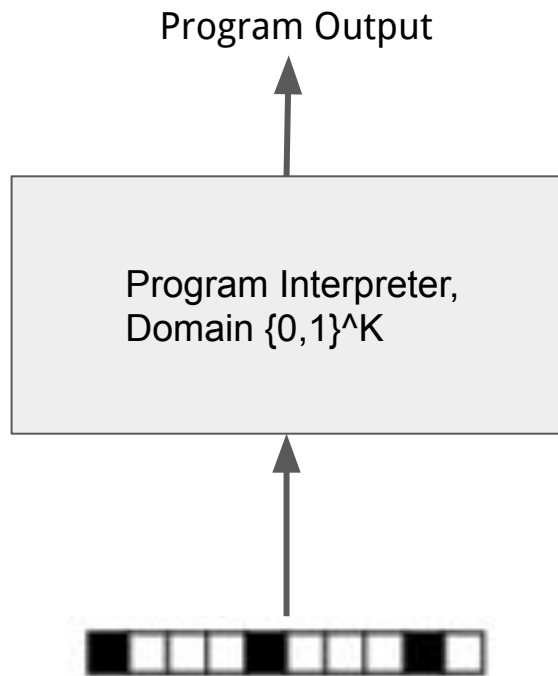
dILP [Evans et al 2017]

EQLNet [Sahoo et al 2018]

DiffLog [Si et al 2019]

...

ROAP, this work [Tang et al 2023]



Relaxing discrete program spaces

Many approaches:

Terpret [Gaunt et al 2016]

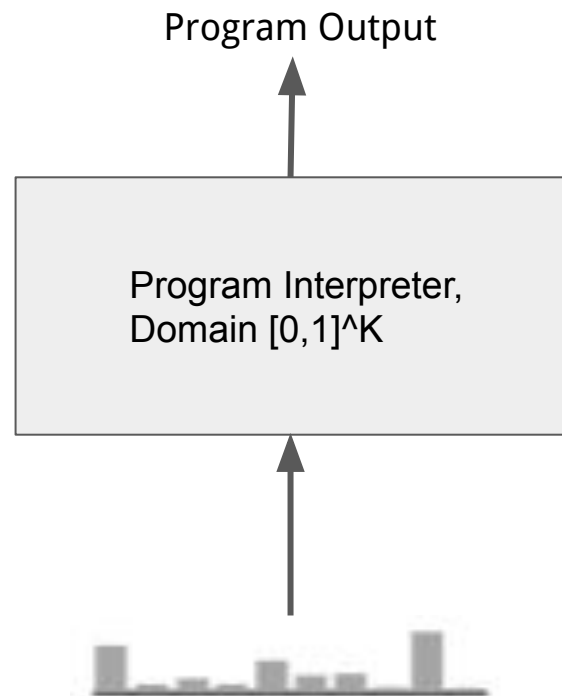
dILP [Evans et al 2017]

EQLNet [Sahoo et al 2018]

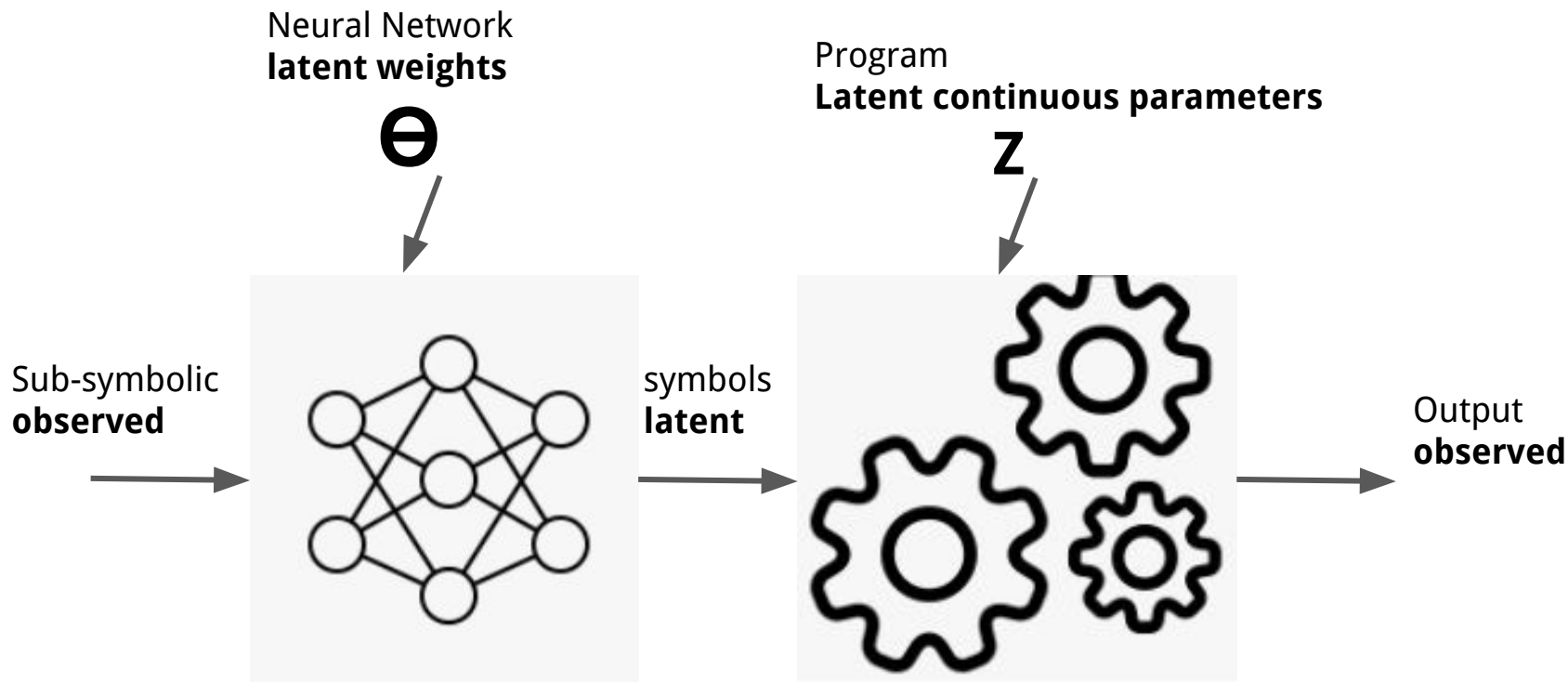
DiffLog [Si et al 2019]

...

ROAP, this work [Tang et al 2023]



A Standard Deep Learning Problem



Out-of-the-box gradient descent doesn't work

Need extra tricks:

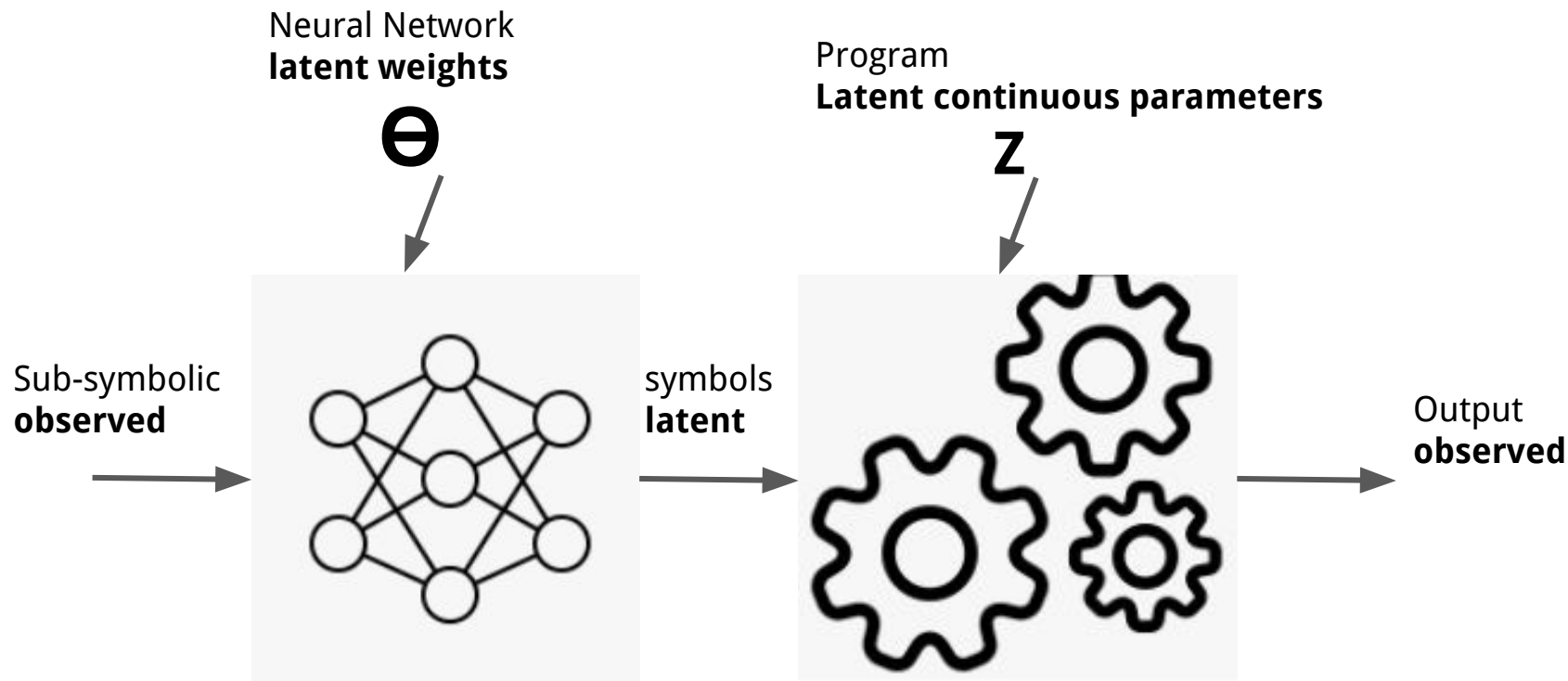
1. Multitasking
2. Overparameterization
3. Special regularizer

Out-of-the-box gradient descent doesn't work

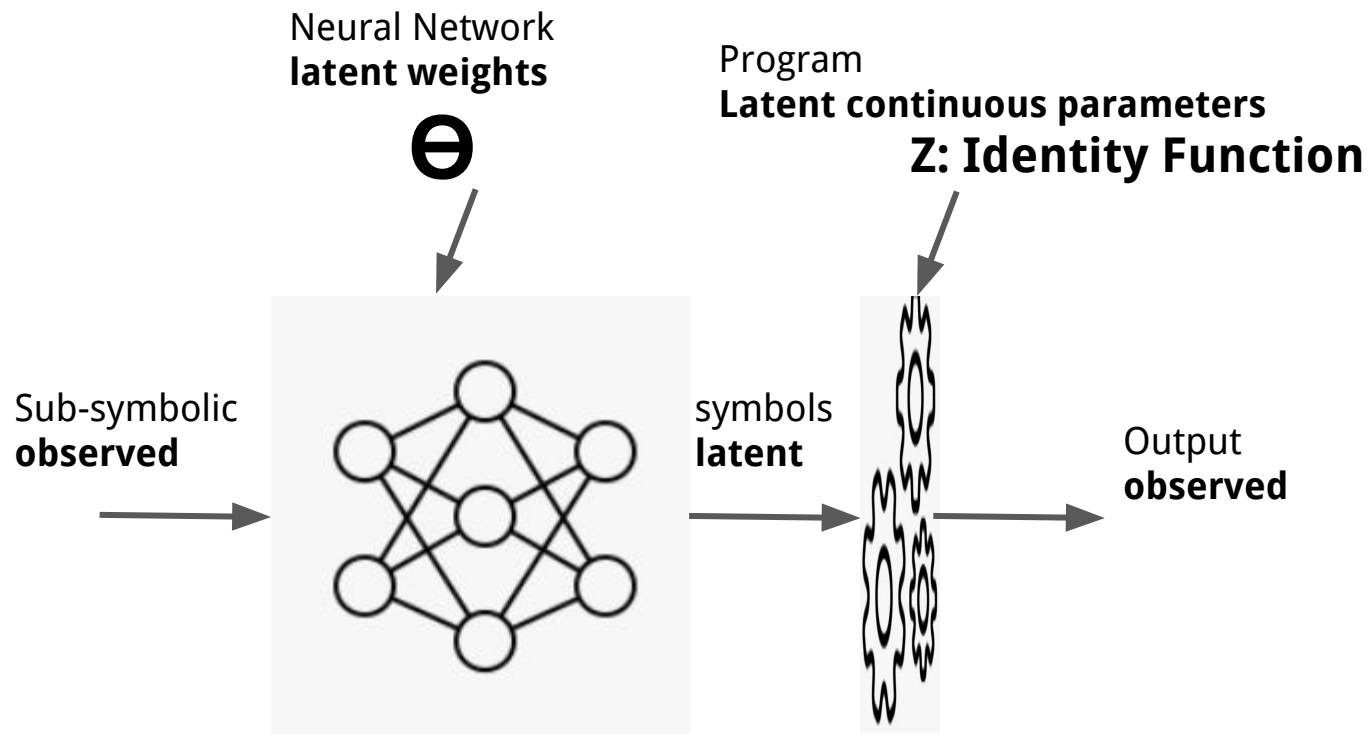
Need extra tricks:

1. **Multitasking**
2. Overparameterization
3. Special regularizer

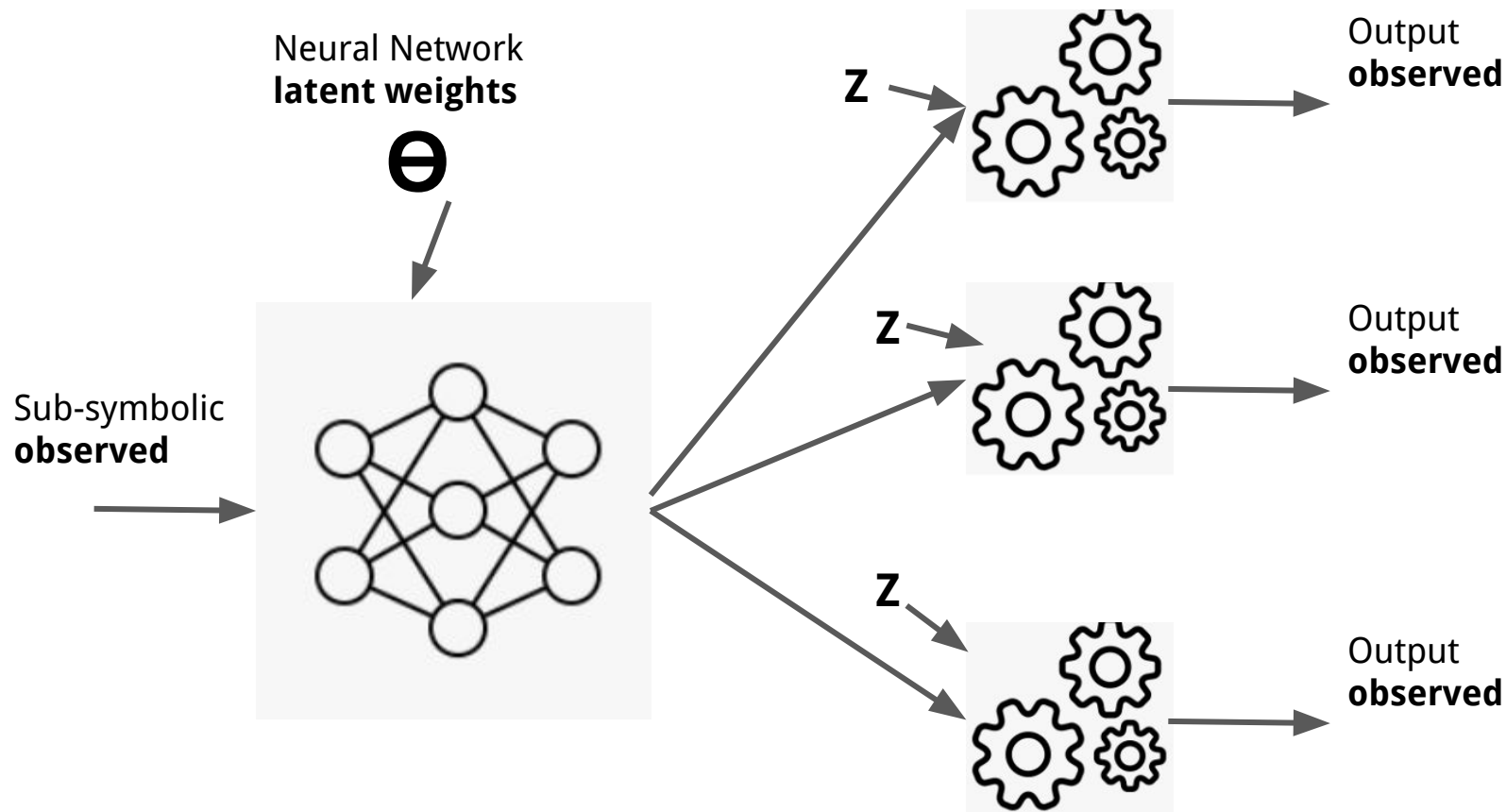
Under-constrained Optimization Problem



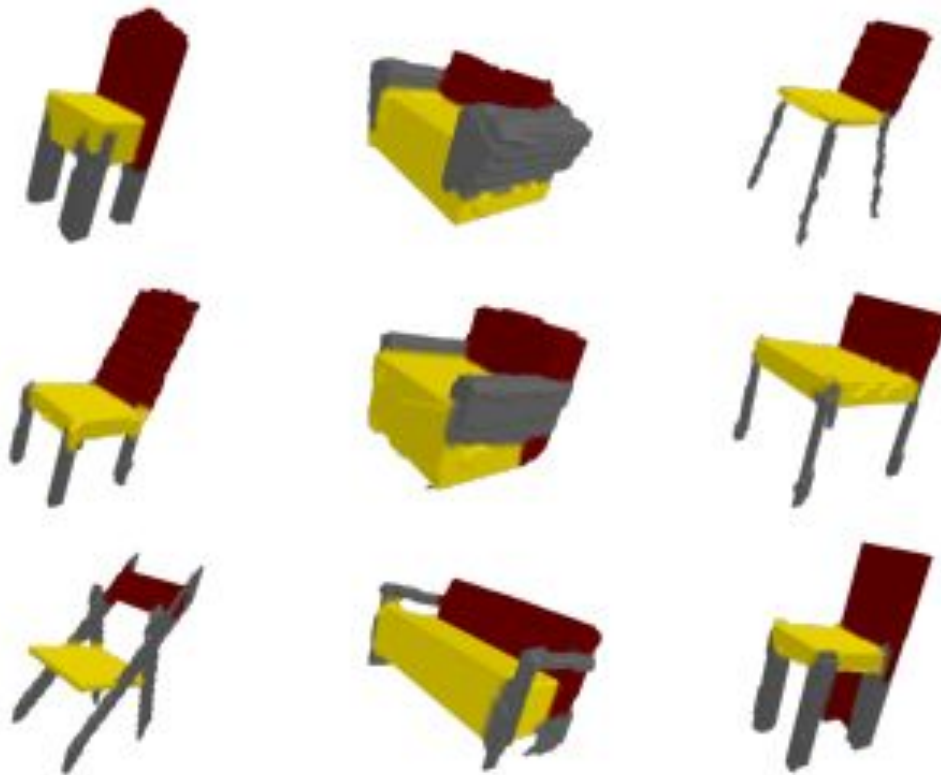
Under-constrained Optimization Problem



Under-constrained Optimization Problem: Therefore, Multitask



Under-constrained Optimization Problem:
Therefore, Multitask



Out-of-the-box gradient descent doesn't work

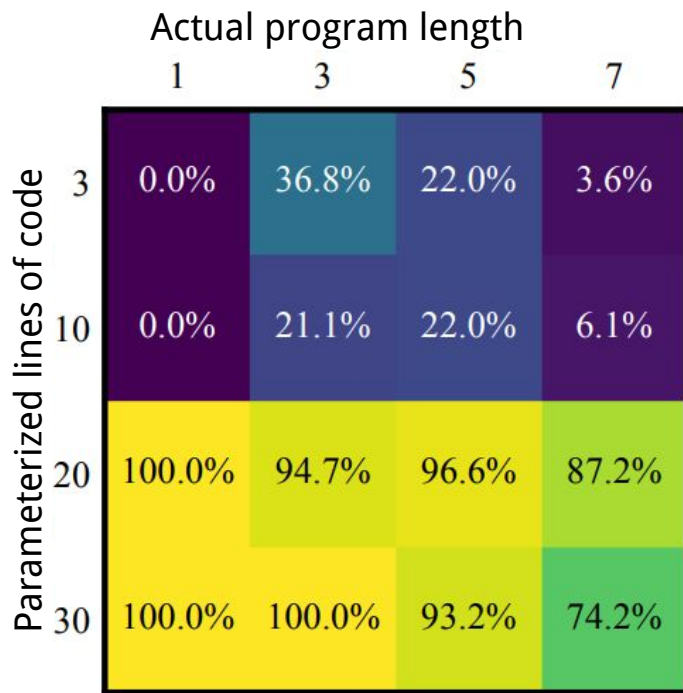
Need extra tricks:

1. Multitasking
2. **Overparameterization**
3. Special regularizer

Overparameterization

Neural nets converge with gradient descent due to massively overparametrization

Analog of overparameterization for programs: Give it more lines of code



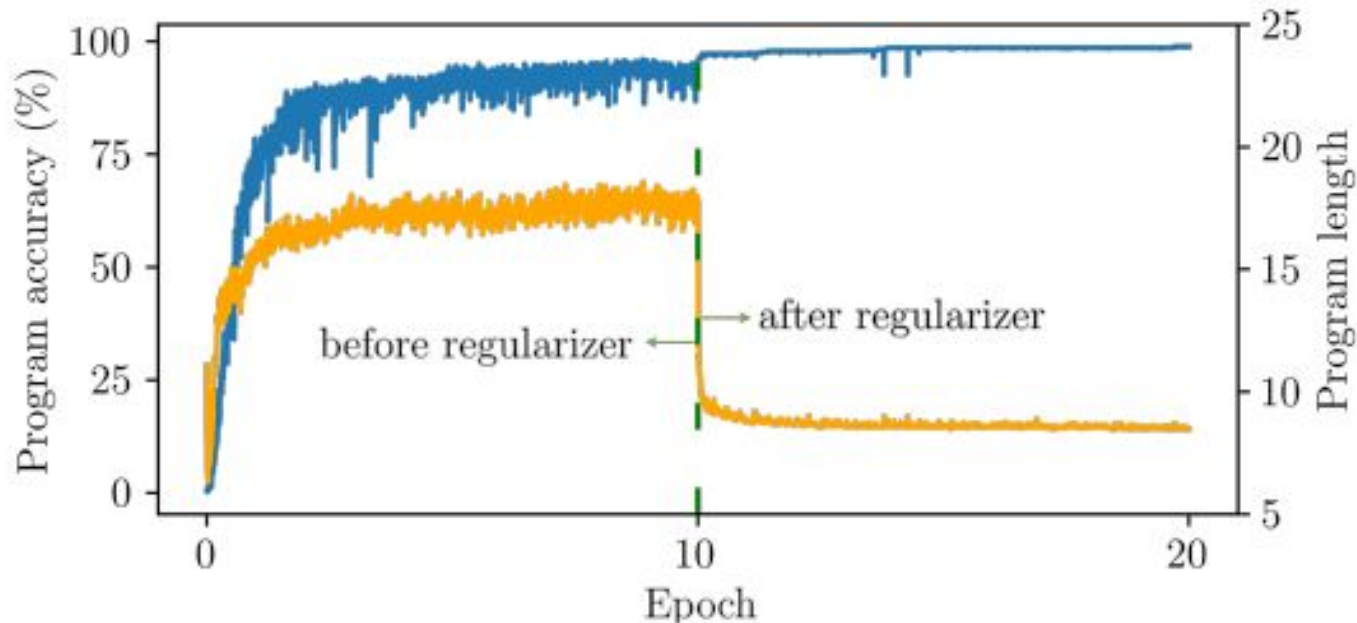
Out-of-the-box gradient descent doesn't work

Need extra tricks:

1. Multitasking
2. Overparameterization
3. **Special regularizer**

New regularizer for relaxed programs: Softly penalized expected program length

Overparameterize: Converge... to long, bloated, overfit programs

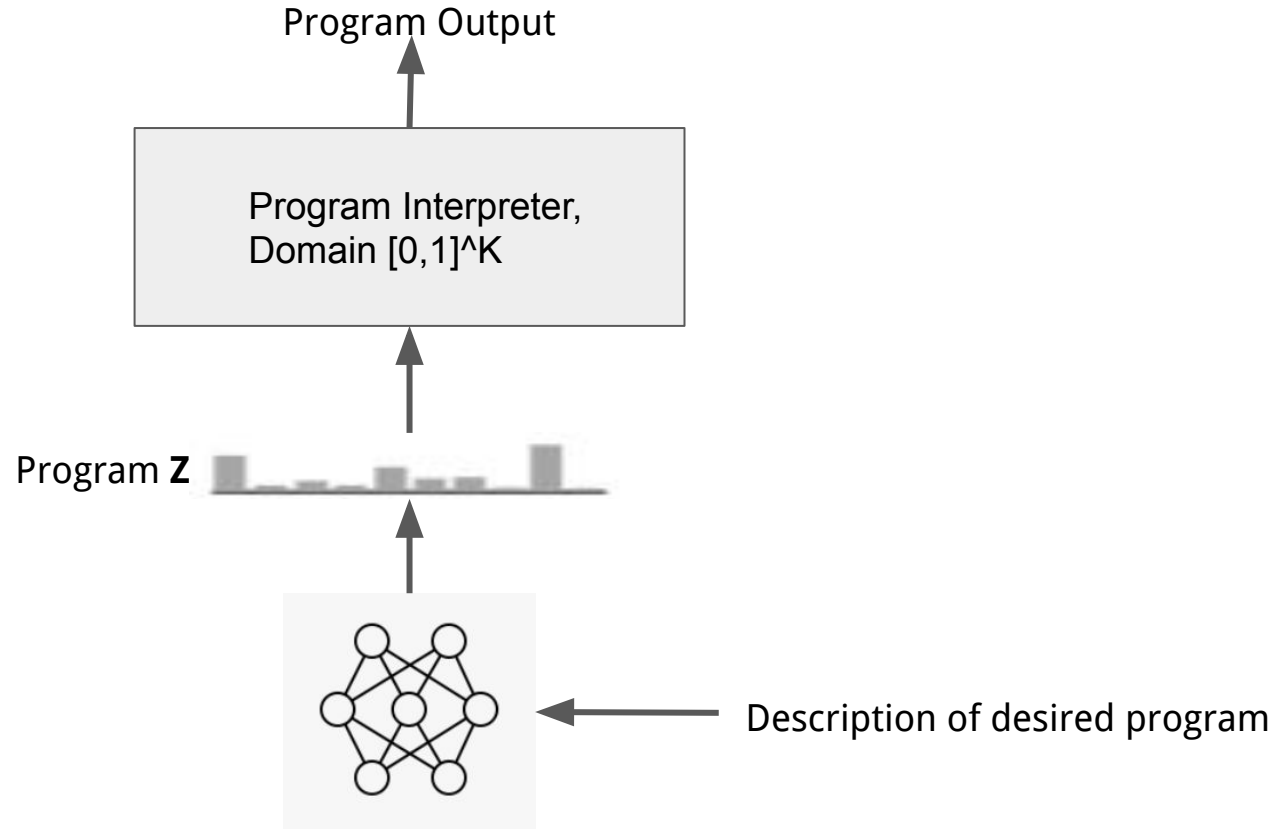


Out-of-the-box gradient descent doesn't work

Need extra tricks:

1. Multitasking, **which allows amortized program synthesis**
2. Overparameterization
3. Special regularizer

Amortized Program Synthesis: Learning to Generate Programs



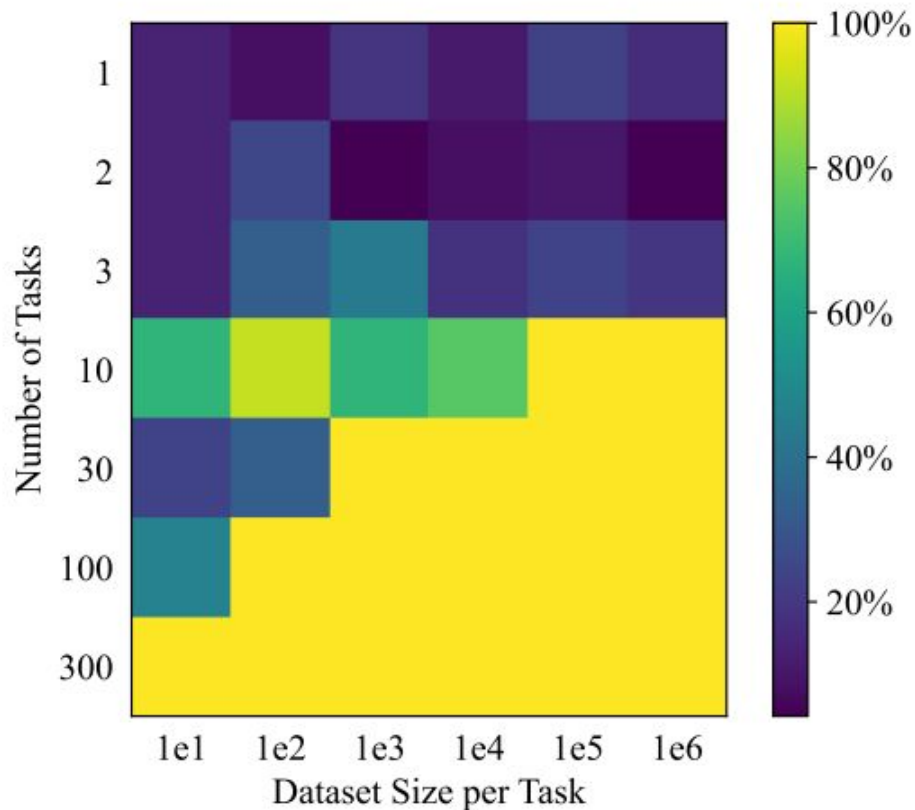
Justifying Each Trick

1. Need multitasking in order to force non-degenerate solutions, imposes extra constraints
2. Need overparameterization for gradient descent to find programs
3. Overparametrization causes bad programs => Need special regularizer
4. Multitasking allows learning to find programs

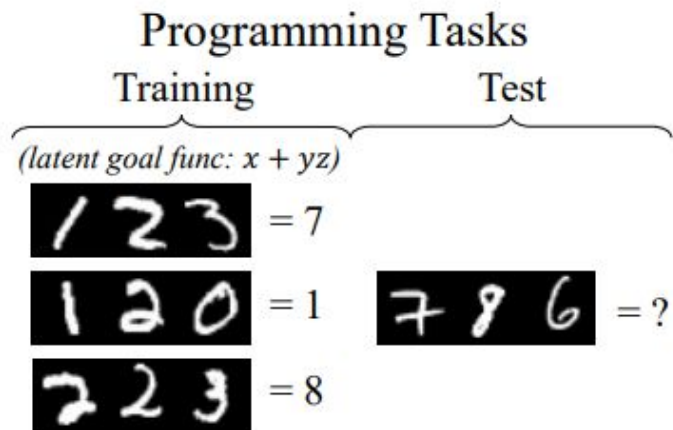
Justifying Each Trick

1. Need multitasking in order to force non-degenerate solutions,
imposes extra constraints
2. Need overparameterization for gradient descent to find programs
3. Overparametrization causes bad programs => Need special regularizer
4. Multitasking allows learning to find programs

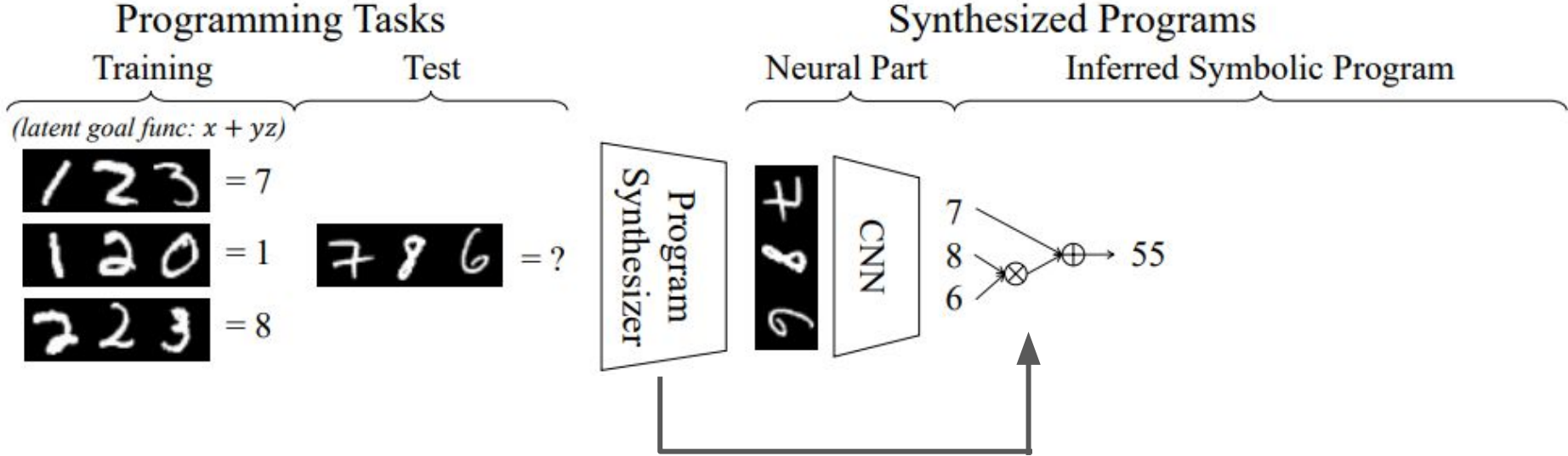
Constraints necessary, But There Exists More Than One Kind of Constraint



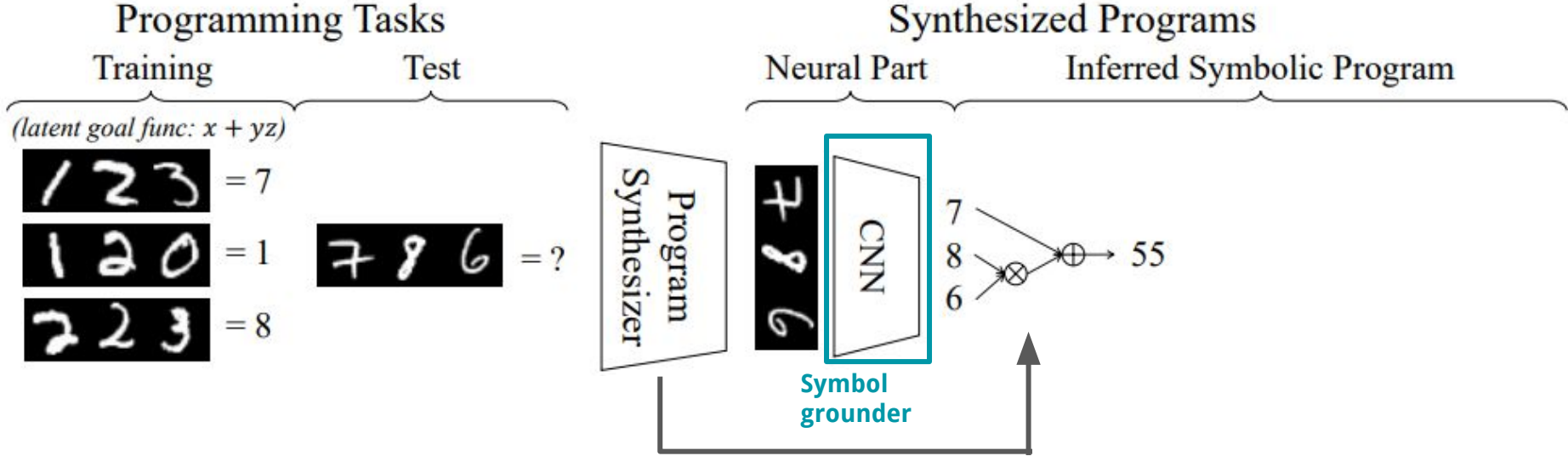
Results: Discovering Number Concepts



Results: Discovering Number Concepts



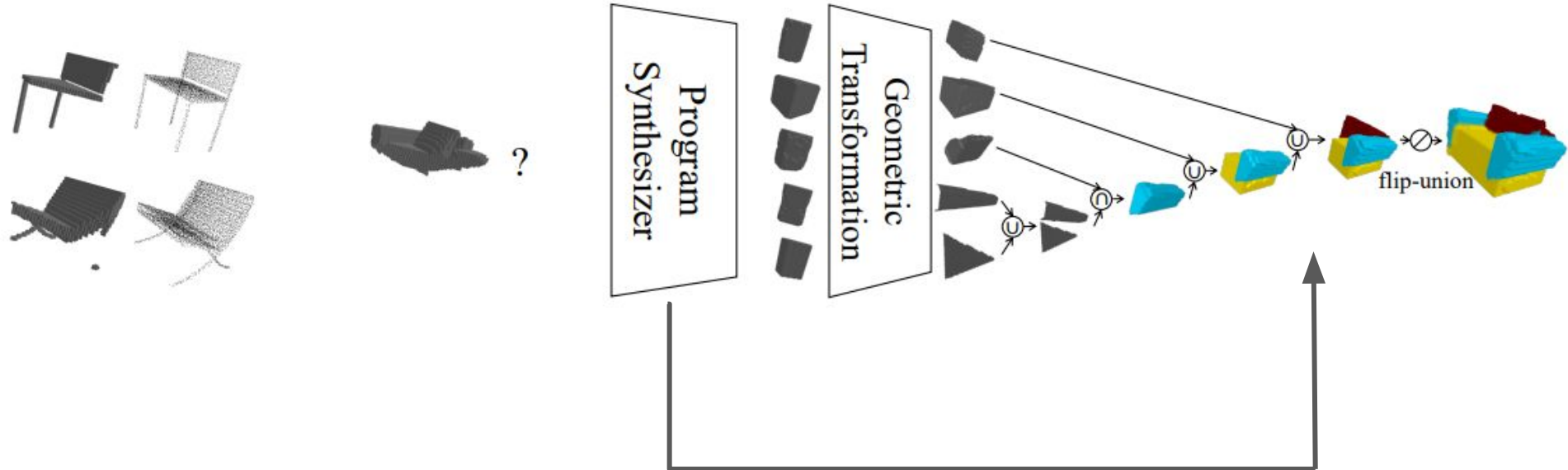
Results: Discovering Number Concepts



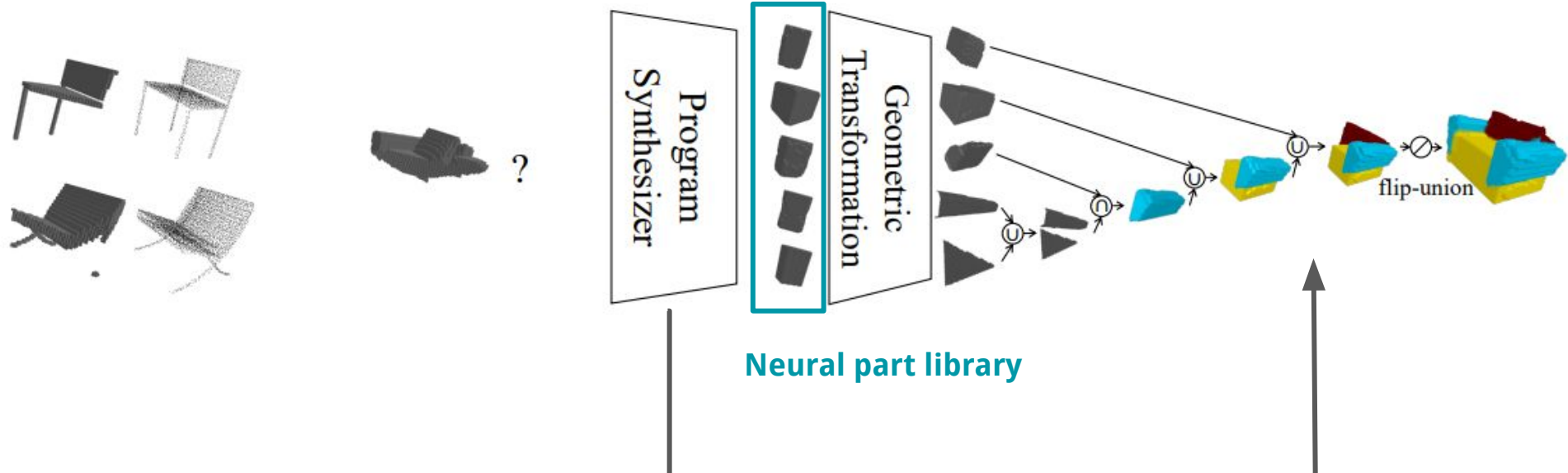
Results: 3D Reconstruction



Results: 3D Reconstruction



Results: 3D Reconstruction



Lessons

Symbol Grounding can emerge from the interaction of different constraints

Doesn't *need* natural language as scaffolding

Gradient descent works for neural nets because of “gradient gadgets”

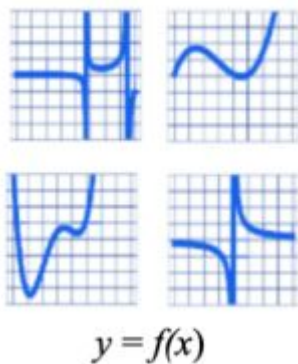
Need to *invent new gadgets* for neurosymbolic programs

the end!

Learning to make programming problems

Simplified setting: Programs take no arguments and their output is the spec

Symbolic Regression

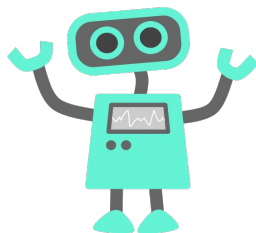


LOGO/Turtle



```
for w in range(6):  
    s0=get_state()  
    draw_square(w)  
    reset_state(s0)
```

Simplified Setting



```
for w in range(6):  
    s0=get_state()  
    draw_square(w)  
    reset_state(s0)
```