

# SynGuar: Guaranteeing Generalization in Programming by Example

Bo Wang, Teodora Baluta, Aashish Kolluri, Prateek Saxena

{bo\_wang, teodora.baluta, aashishk}@u.nus.edu, prateeks@comp.nus.edu.sg

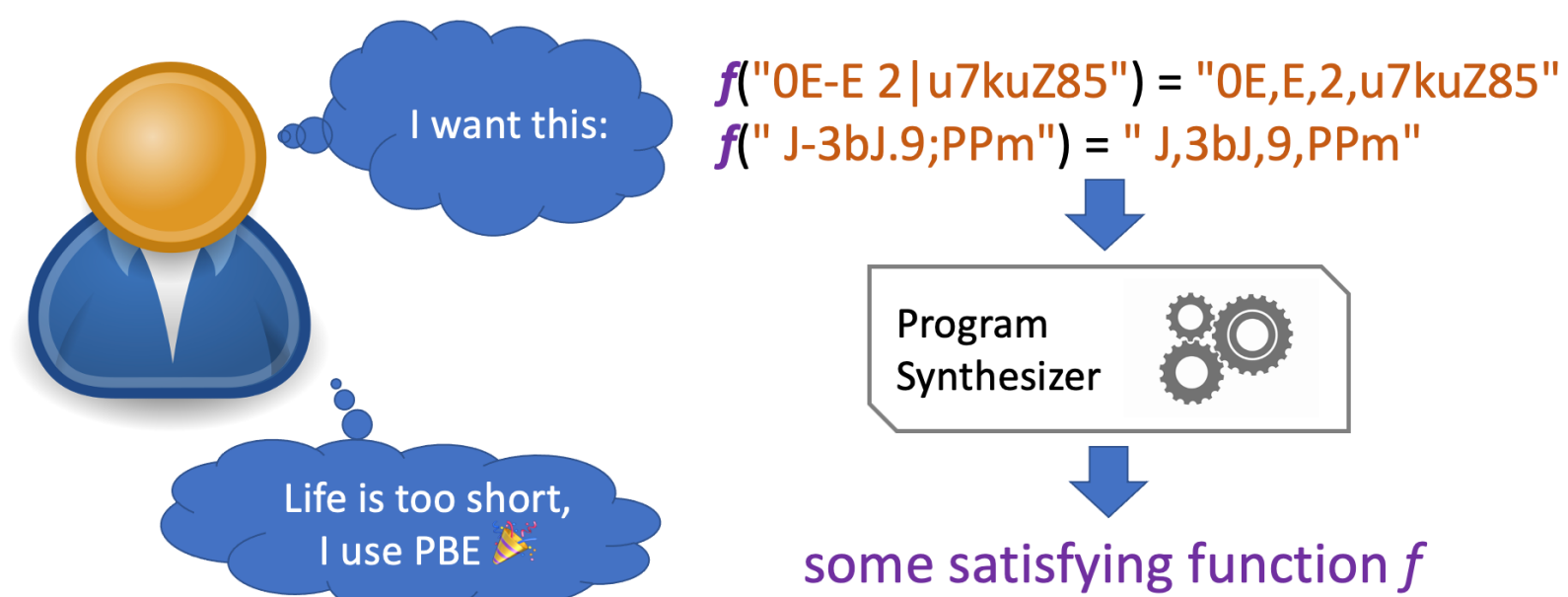


## Abstract

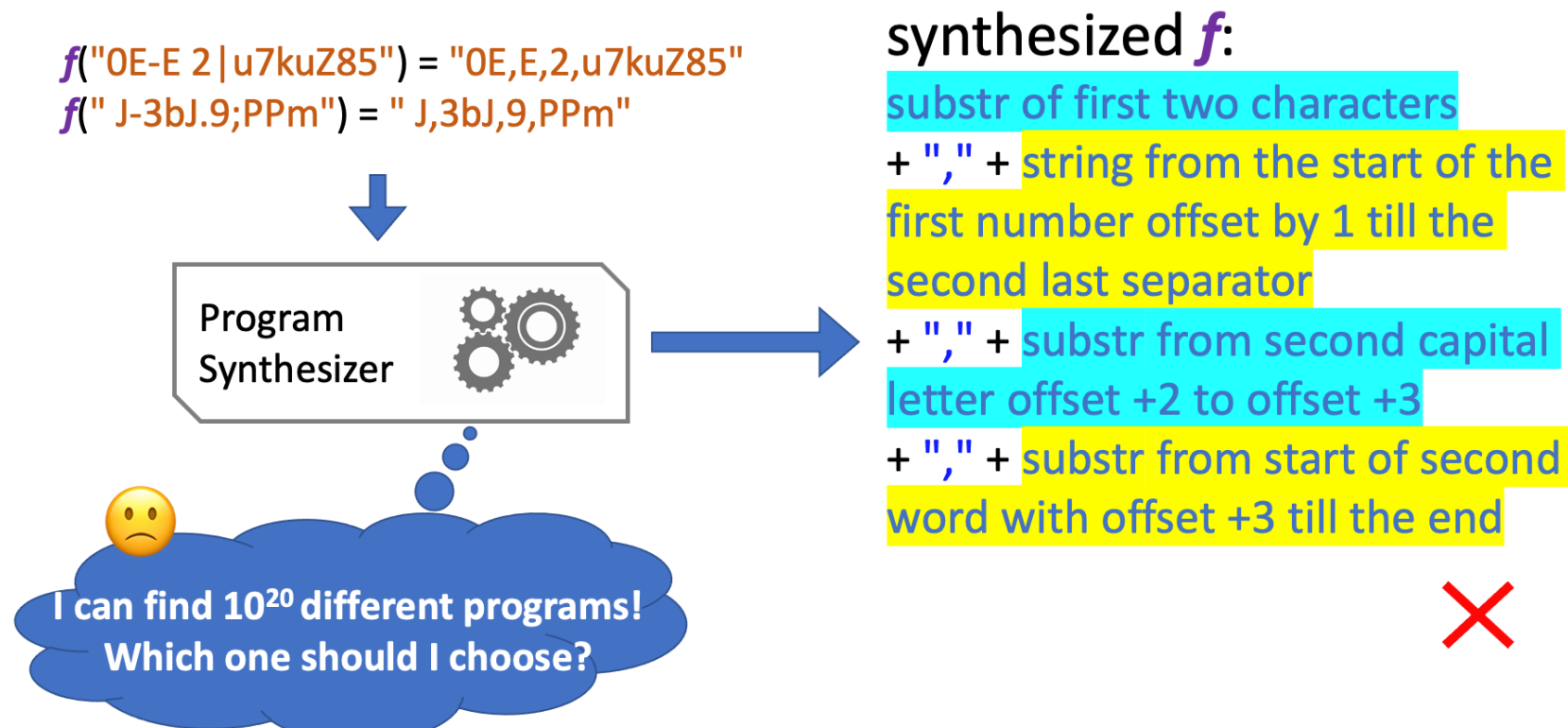
- **SynGuar** is a framework to provably reduce overfitting in program synthesis.
- It achieves **Probability Approximately Correct (PAC)** generalization in Programming-by-Example.
- It currently works for countable program space and realizable settings.

## Overfitting Problem in Programming-by-Example

**Programming-by-Example (PBE):** Input-Output pairs  $\rightarrow$  Code



Synthesizers suffer from overfitting with **insufficient examples**.



## PAC Framework

PAC framework gives a way to compute the **sufficient sample size for generalization**.

### PAC Formalization

How many examples are sufficient to return a program  $f \in \mathcal{H}$  that is **close to the target** with **high probability**?

The generalization error  $\text{error}(f) = \Pr_{x \sim D}[f(x) \text{ not correct}]$  is bounded by  $\epsilon$ .

The probability of generating  $f$  with  $\text{error}(f) > \epsilon$  is bounded by  $\delta$ .

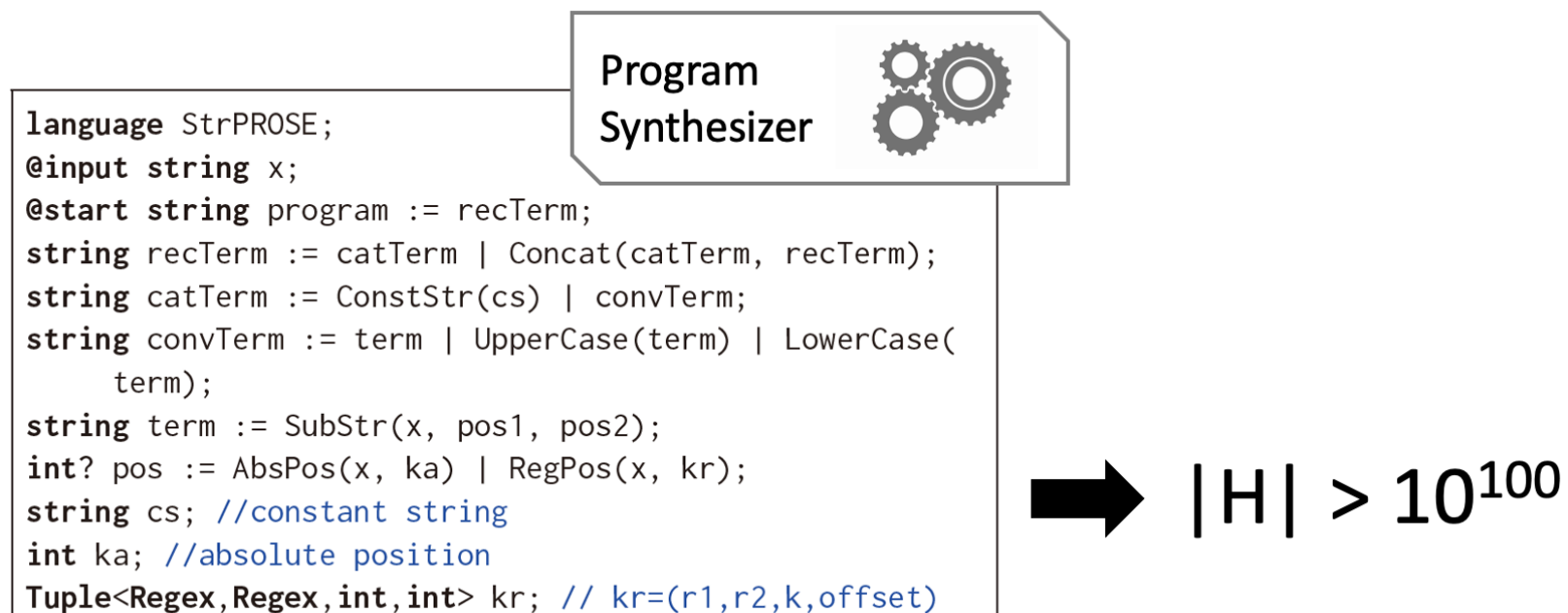
**Our Goal:**  $\Pr[\text{error}(f) > \epsilon] < \delta$

Classical Sample Complexity Bound

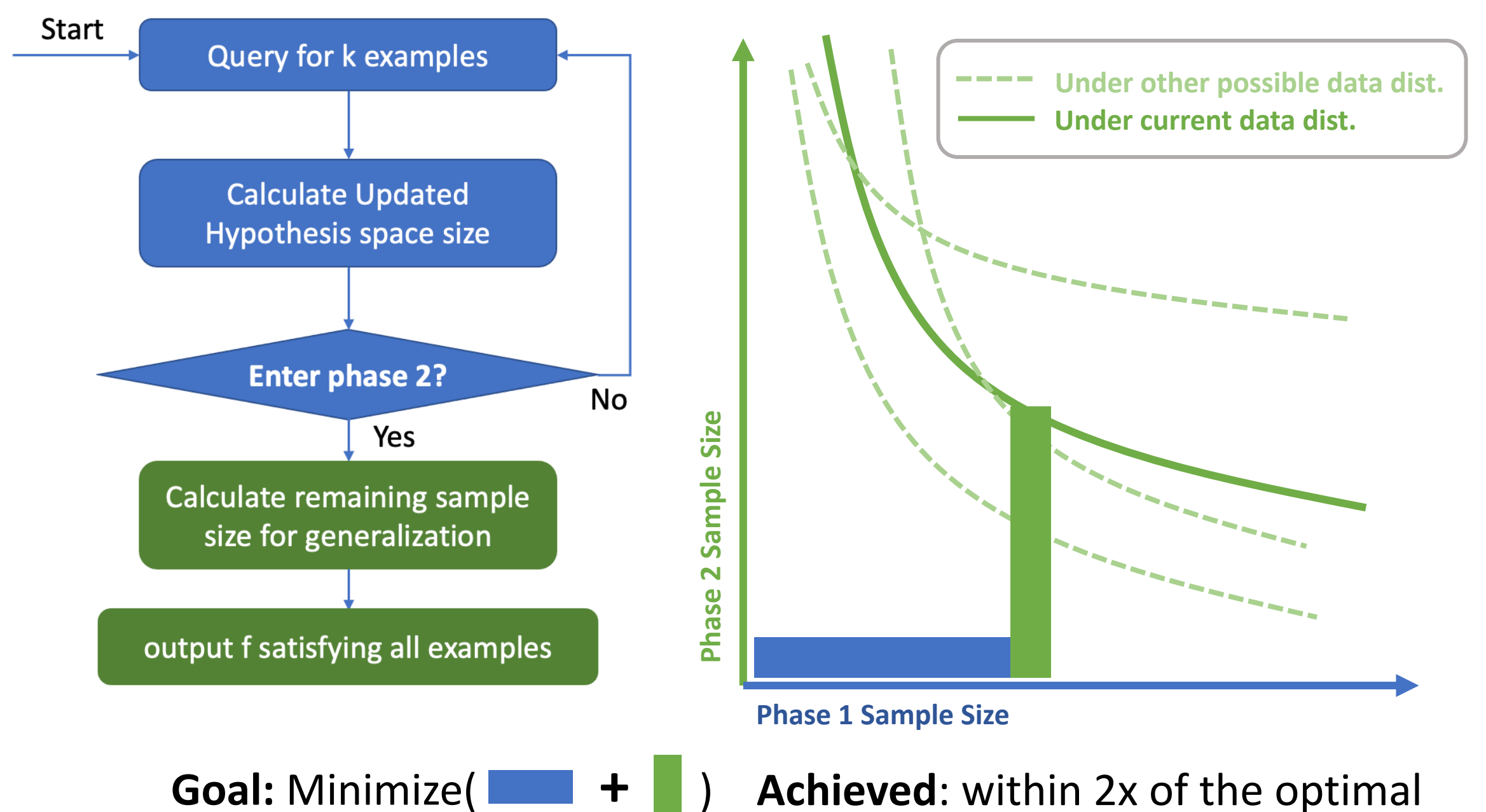
with sample size  $m > \frac{1}{\epsilon} (\ln |H| + \ln \frac{1}{\delta})$

we have  $\Pr[\text{error}(f) > \epsilon] < \delta$

But synthesizers have **astronomical** program space sizes ...

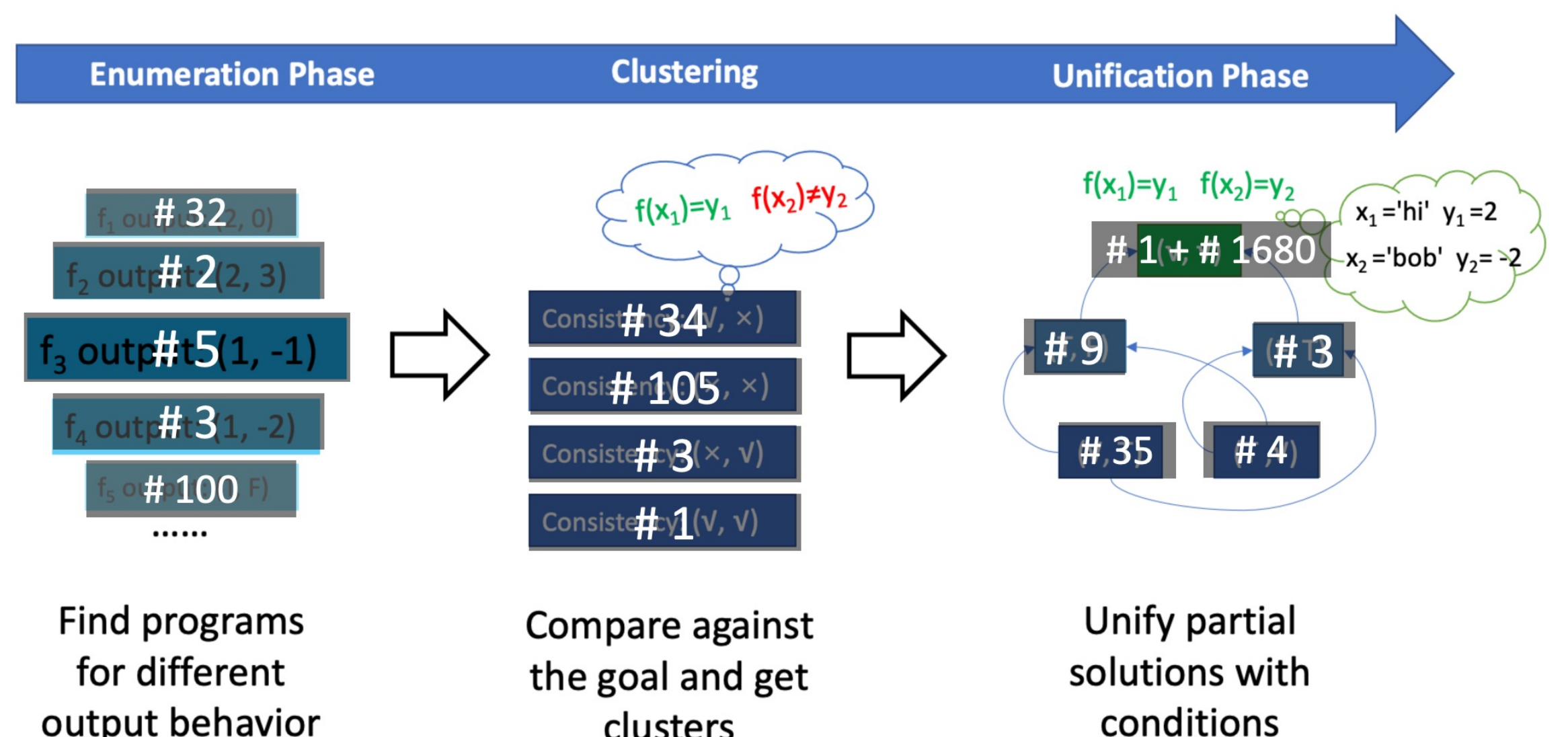


## Solution: SynGuar Algorithm (Contribution 1)



## Challenge: Compute Program Space Size (Contribution 2)

- **PROSE-based** synthesizer (StrPROSE), already supported. PROSE framework has **Size** API for the program space.
- **Bottom-up** synthesizer (StrSTUN), modified to add counting.



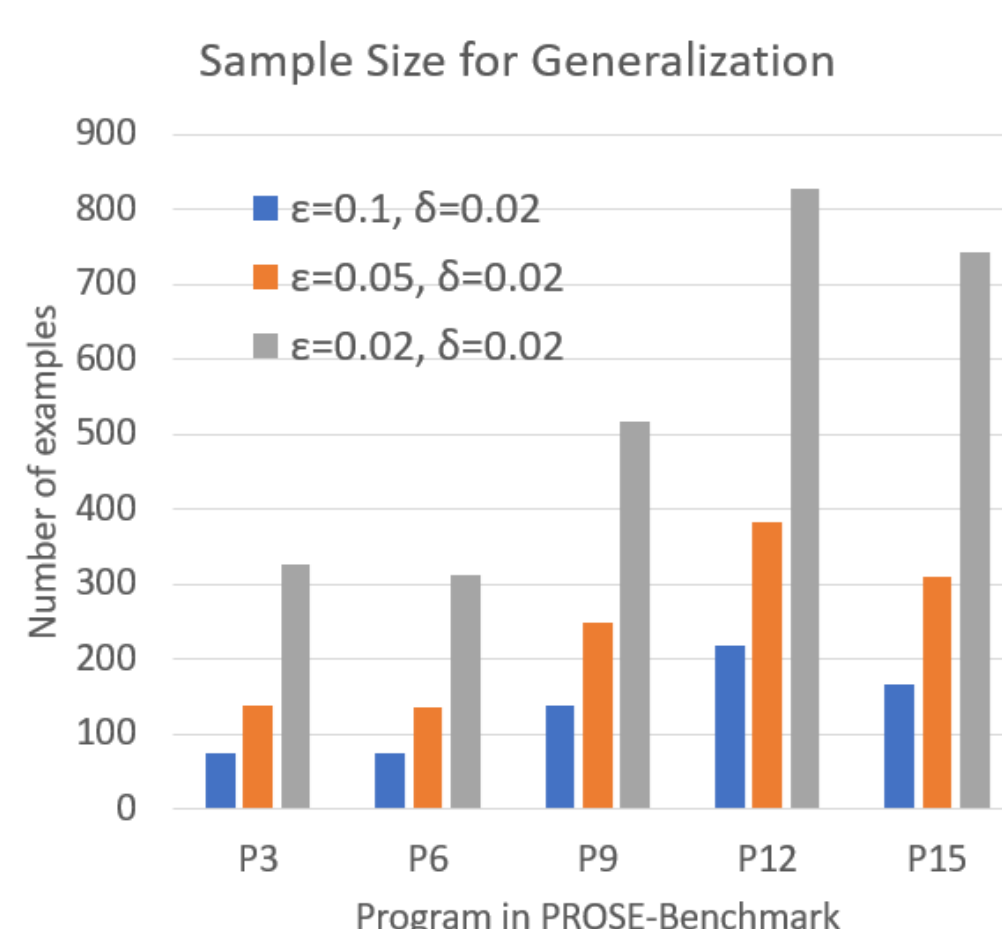
## Evaluation

for SynGuar + StrPROSE

- Under 400 examples for all testcases to achieve ( $\epsilon = 0.05, \delta = 0.02$ ).
- Improving correctness
  - $\approx 94\%$  with SynGuar(0.05, 0.02)
  - $\approx 34\%$  with 4 random examples

for SynGuar + StrSTUN

- Most of the testcases are under 500 examples for  $\epsilon = 0.05, \delta = 0.02$ .
- Improving correctness
  - $\approx 90\%$  with SynGuar(0.05, 0.02)
  - $\approx 56\%$  with 4 random examples
  - $\approx 61\%$  with provided examples



Online Demo



GitHub

