# Generating Sequences by Learning to [Self-]Correct
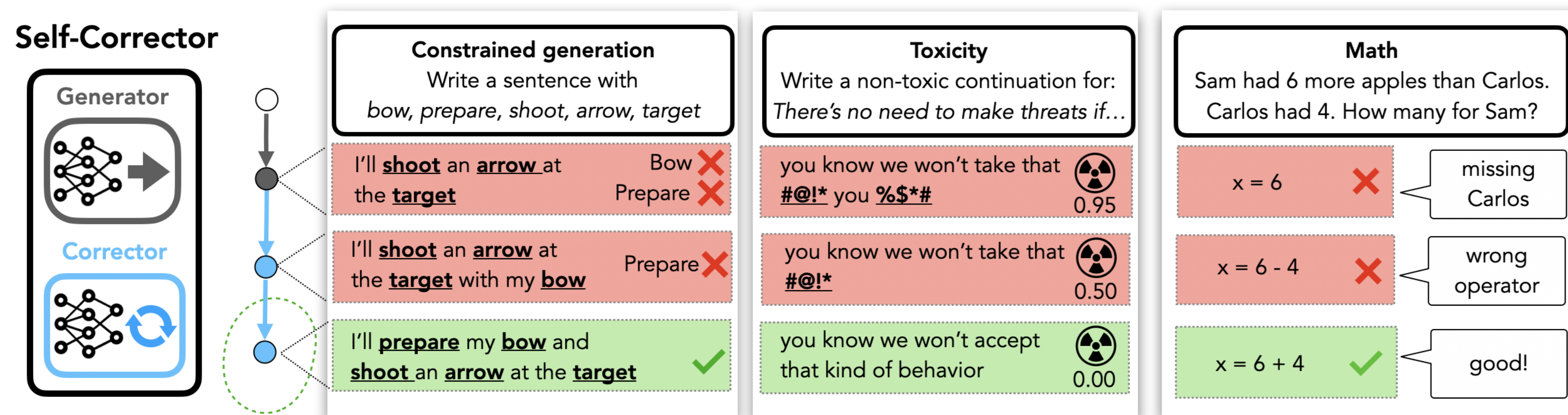
Sean Welleck*[1,2]   Ximing Lu*[2]   Peter West[+1]   Faeze Brahman[+2]   Tianxiao Shen[1]   Daniel Khashabi[2]   Yejin Choi[1,2]

[1]University of Washington      [2]Allen Institute for Artificial Intelligence

## Self-correctors

How do we **control** and **improve** a language model's generations **after it is trained**?

Key idea: plug in a learned corrector that iteratively improves outputs.



### Self-corrector = generator + learned corrector

$$p(y|x) = \sum_{y_0} \underbrace{p_0(y_0|x)}_{\text{generator}} \underbrace{p_\theta(y|y_0, x, f)}_{\text{corrector}}$$

Self-correctors offer several benefits, including:
1. **Controlling** generators without modifying them
2. **Decomposing** problems into multiple iterations
3. Using **natural language feedback** for (1) and (2)
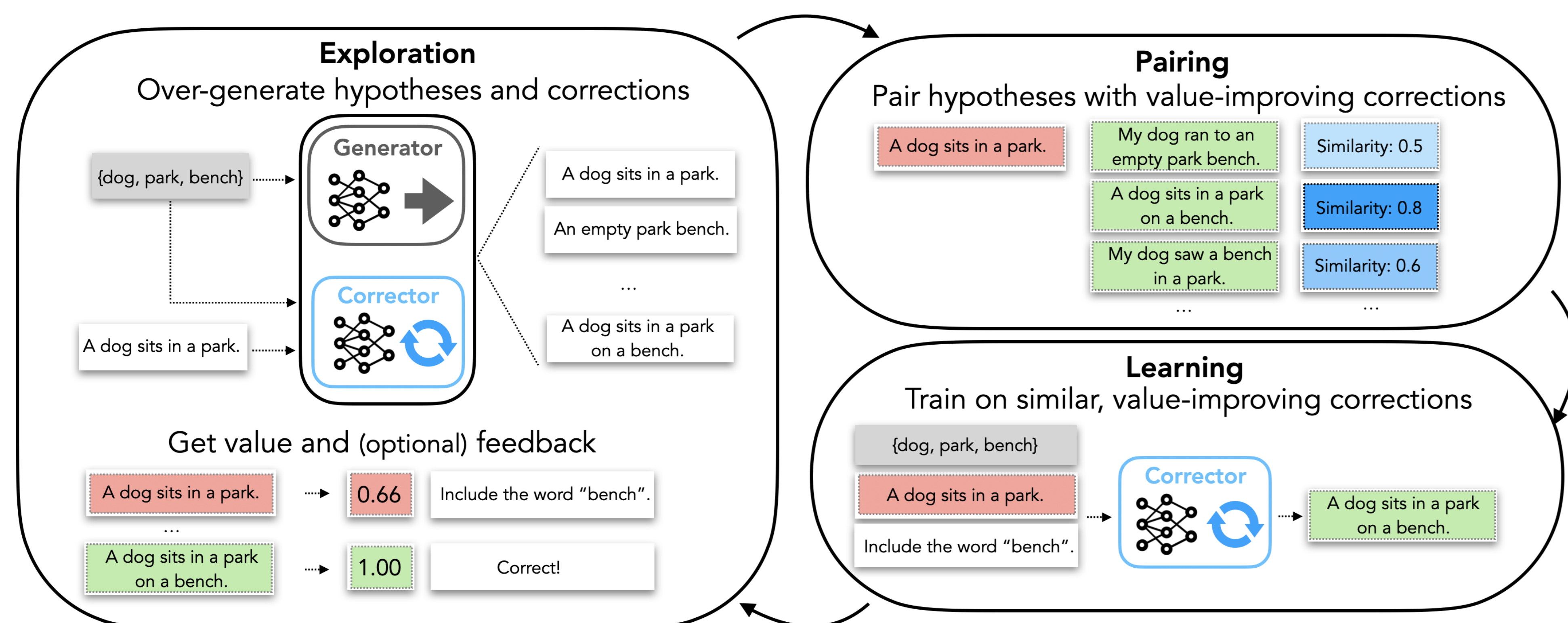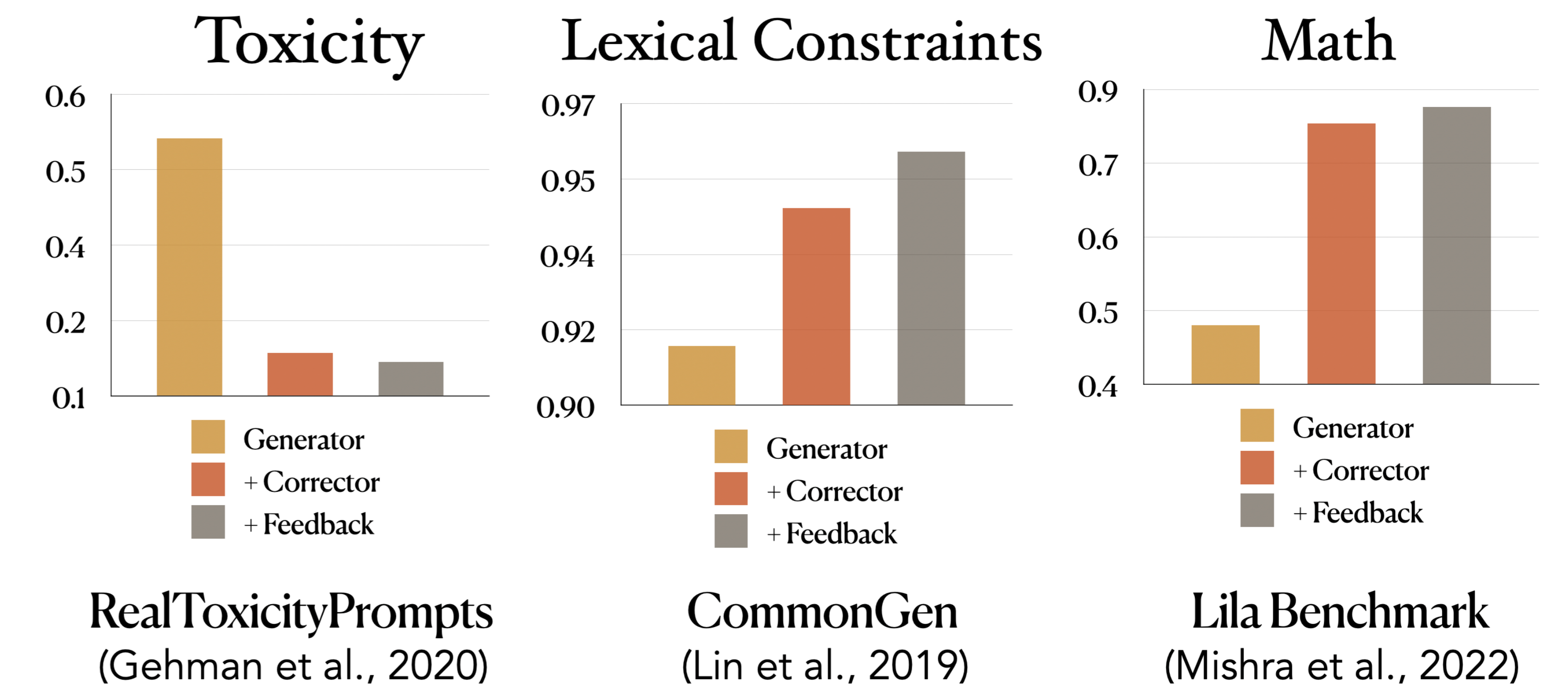
### Learning without annotated corrections



Figure 1. Self-corrective learning iteratively trains a corrector by generating hypotheses and corrections, forming value-improving pairs, and selecting those with high similarity for learning.
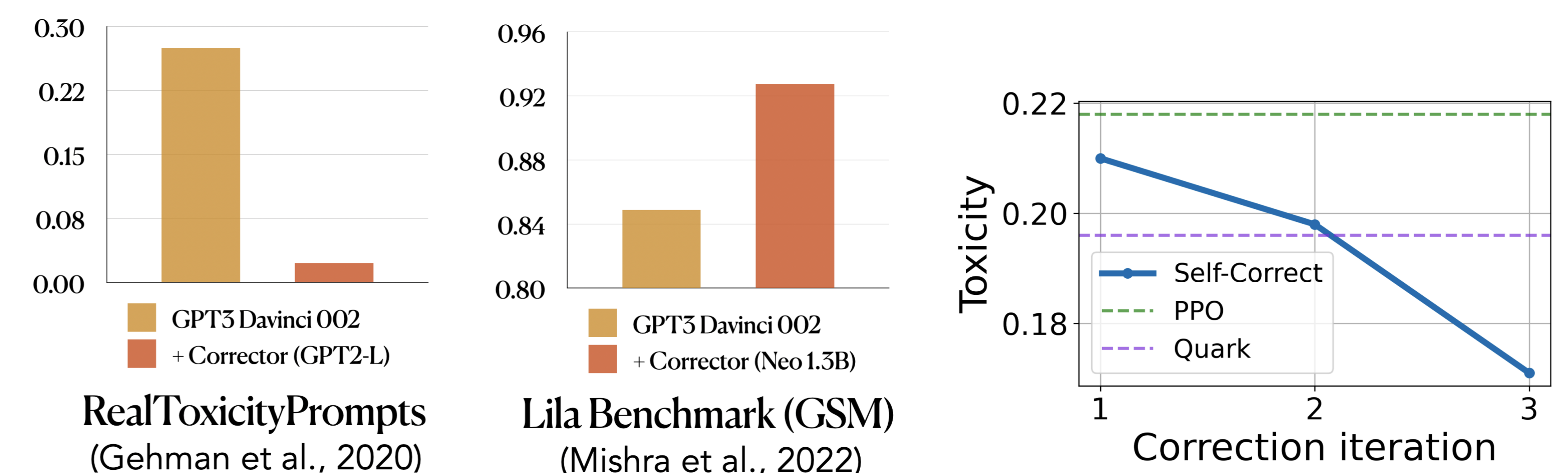
## Self-correction

Self-correctors improve upon the base generator, and natural language feedback brings additional gains. Diverse tasks: **toxicity, lexical constraints, mathematical program synthesis**.



Feedback sources:

- **Toxicity**: fine-grained properties from Perspective API, e.g.
- **Lexical constraints**: missing words, e.g. *add 'bow' and 'prepare'*
- **Math**: few-shot prompted GPT-3, e.g. *2 is missing*

### Correcting inaccessible larger models & multiple corrections



## Discussion

- **Natural language feedback**: *sources* (e.g., humans, models) and *formats* (e.g., line-by-line).
- Other **learning algorithms** for the corrector: e.g. reinforcement learning