

Live Exploration of AI-Generated Programs

Kasra Ferdowsi

with Ruanqianqian (Lisa) Huang*, Michael B. James, Nadia Polikarpova, and Sorin Lerner

UC San Diego

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

LLMs for Code Generation*



GitHub
Copilot



OpenAI
ChatGPT



* For experienced programmers.

Background

1. **Grounded Copilot: How Programmers Interact with Code-Generating Models.**
Shraddha Barke, Michael B. James, and Nadia Polikarpova. 2023.
2. **Understanding the Usability of AI Programming Assistants.**
Jenny T. Liang, Chenyang Yang, and Brad A. Myers. 2023.
3. **Reading Between the Lines: Modeling User Behavior and Costs in AI-Assisted Programming.**
Hussein Mozannar, Gagan Bansal, Adam Fourney, and Eric Horvitz. 2022.
4. **Expectation vs. Experience: Evaluating the Usability of Code Generation Tools Powered by Large Language Models.**
Priyan Vaithilingam, Tianyi Zhang, and Elena Glassman. 2022.

Background

In summary, programmers using AI-generated code...

- Spend significant time *validating* code suggestions,
- Have trouble evaluating the correctness of generated code,
- Choose validation strategies based on *time cost*, and so
- Both *under- and over-rely* on AI code suggestions.

Background

"User interactions with Copilot can be classified into two modes—*acceleration* and *exploration*—akin to the two systems of thought in dual-process theories of cognition"*

* Grounded Copilot: How Programmers Interact with Code-Generating Models

Shraddha Barke, Michael B. James, and Nadia Polikarpova. 2023.

Grounded Copilot

Acceleration	vs.	Exploration
unintentional	Prompting	intentional with comments / invoke side panel
"pattern matching"	Validation	explicit validation via elimination / execution / documentation
unit of focus (sub-expression / statement)	Scope	entire function + multiple alternatives
unwilling to edit	Mismatch Tolerance	willing to edit / debug / "rip apart" / cherry-pick

Grounded Copilot

Acceleration	vs.	Exploration
unintentional	Prompting	intentional with comments / invoke side panel
"pattern matching"	Validation	explicit validation via elimination / execution / documentation
unit of focus (sub-expression / statement)	Scope	entire function + multiple alternatives
unwilling to edit	Mismatch Tolerance	willing to edit / debug / "rip apart" / cherry-pick



Live Programming

Demo!

Research Question:

Does *Live Programming* offer a good interaction for *exploring* AI-generated code?

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

Copilot's Side Panel

```
home > kas > projects > tmp > blah.py > ...
12 import matplotlib.pyplot as plt
11 import pandas as pd
10
9 def import_data(path: str):
8     with open(path, 'r') as f:
7         return pd.read_csv(f)
6
5 # Get the data, assuming we have the following
4 # column names: "group", "time", "success"
3 df = import_data('my_data.csv')
2
1 # create a box plot where each box refers to a group
13
1
```

```
1 Synthesizing 10/10 solutions (Duplicates hidden)
1 Accept Solution
2 =====
3 Suggestion 1
4
5 # and the y axis is the time it took to complete the ta
6 df.boxplot(column='time', by='group')
7 plt.show()
8 ```
9 This is not necessarily the most efficient way to do it
10
11 Accept Solution
12 =====
13 Suggestion 2
14
15 fig, ax = plt.subplots()
16 ax.boxplot(df['time'], labels=df['group'])
17 ax.set_ylabel('Time (s)')
18 ax.set_title('Time taken to complete a task')
19 plt.show()
20 ```
21 You can find more information about box plots here.
22
23 Accept Solution
```

LEAP

Demo!

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

Research Questions

How does Live Programming effect...

1. Code Correctness
2. Over-/Under-Reliance on AI
3. Cognitive Load
4. Users' Impressions

Experimental Conditions

no-PB

AI suggestions

+

Manually Invoked
Terminal Output

PB

AI Suggestions

+

Projection Boxes

Tasks

Pandas

clean dataframe and compute stats using pandas

API-Heavy

Algorithmic

Bigrams

find the most frequent bigram in a string

Fixed-prompt

Box Plot

overlay scatter plot over boxplot using matplotlib

Open-prompt

String Rewriting

parse rewrite rules and apply to a string

Participants

n = 17



Occupation:

15 academia

2 industry

Python Usage:

2 occasionally

8 regularly

7 almost every day

Overview

1. Motivation & Background:

- a. Grounded Copilot
- b. Live Programming

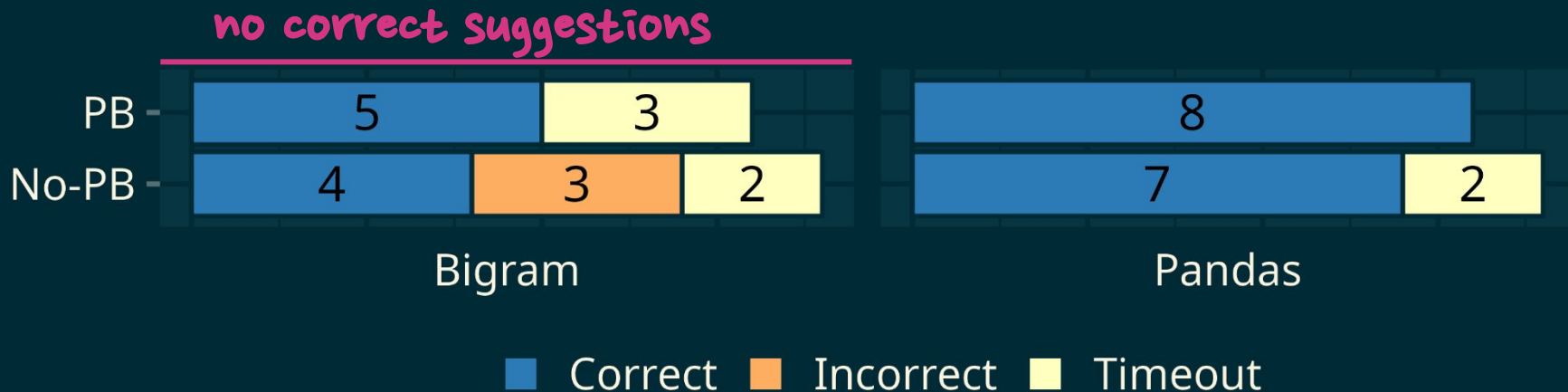
2. LEAP: Live Exploration of AI-Generated Code

3. User Study

4. Findings

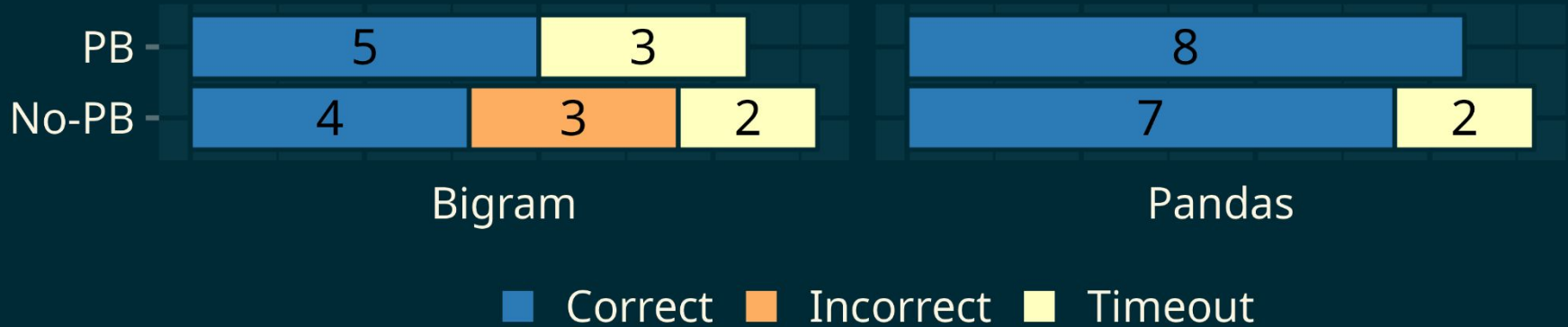
- a. Validating suggestions
- b. Over-/Under-reliance
- c. Cognitive Load
- d. Impressions

RQ1: Correctness



LEAP helps validate suggestions!
(But does not help fix incorrect ones)

RQ2: Over-/Under-reliance



6 no-PB vs 0 PB participants **mis-judged** correctness of their solutions

RQ2: Over-/Under-reliance

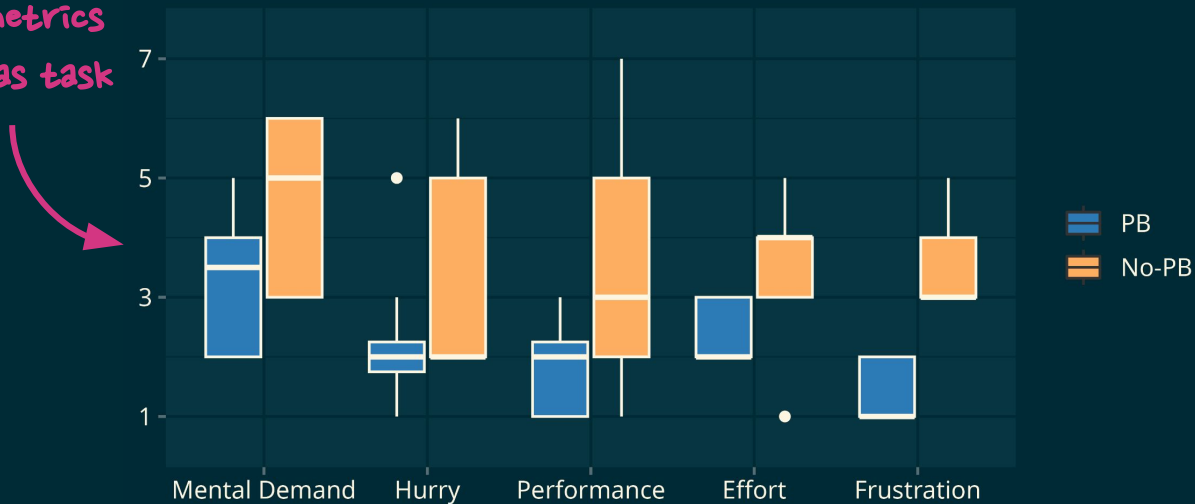
"it was **easy to understand** the behavior of a code suggestion because the little boxes on the side allowed for you to preview the results." (P3)

"it **saved me the effort** of writing multiple print statements." (P1)

LEAP reduces over-/under-reliance on AI,
by lowering the cost of validation.

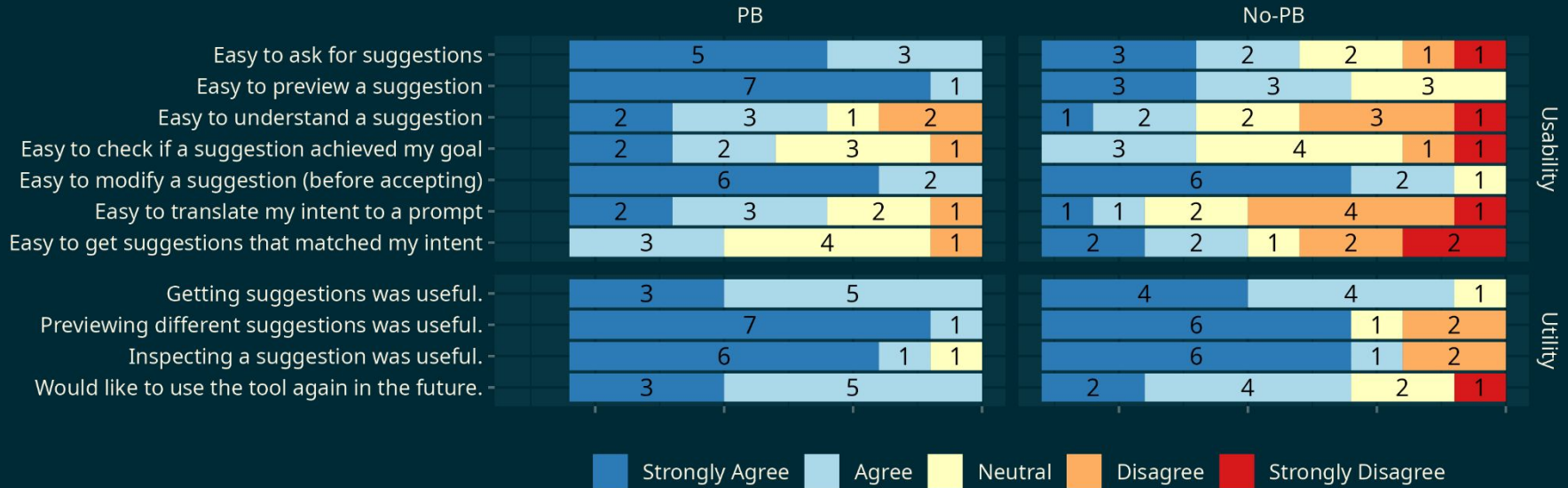
RQ3: Cognitive Load

NASA TLX metrics
on the Pandas task



LEAP significantly reduced the cognitive load of exploring AI suggestions *on tasks amenable to validation by execution.*

RQ4: Users' Impressions



LEAP was more *usable* and more *useful*.

Overview

1. Motivation & Background:
 - a. Grounded Copilot
 - b. Live Programming
2. LEAP: Live Exploration of AI-Generated Code
3. User Study
4. Findings
 - a. Validating suggestions
 - b. Over-/Under-reliance
 - c. Cognitive Load
 - d. Impressions

Summary

1. Background:
 - a. Live Programming, for
 - b. Exploration of AI-Generated Code
2. LEAP: Projection Boxes + Copilot-like interface
3. User Study w/ 17 participants
4. Found that LEAP...
 - a. Helps with validating suggestions,
 - b. Reduces Over-/Under-reliance,
 - c. Improves Cognitive Load, and
 - d. Leaves a positive impressions on participants.