

how programmers interact with ai assistants

Nadia Polikarpova

DL4Code @ ICLR'23

with Shraddha Barke, Kasra Ferdowsifard, Lisa Huang, Michael B. James, and Sorin Lerner

the new era of programming



GitHub Copilot



Chat GPT

and more...



Amazon CodeWhisperer

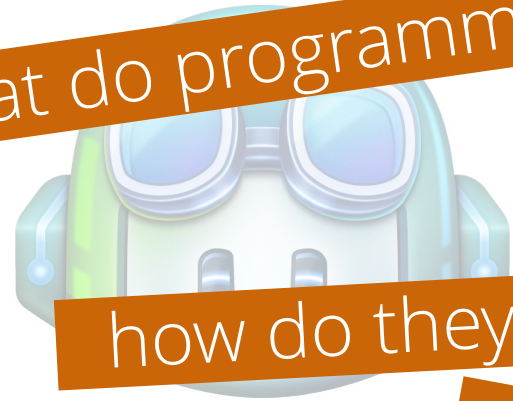
the new era of programming

what do programmers use them for?

how do they express intent?

how do they validate suggestions?

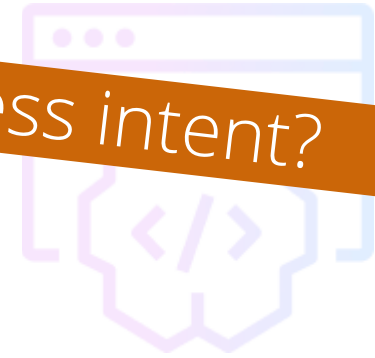
how do they cope with errors?



Github Copilot



ChatGPT



Amazon CodeWhisperer

and more...

this talk

grounded copilot

grounded theory
of AI-assisted programming

[OOPSLA'23]

leap

helping programmers
validate AI-generated code

[under review]

this talk

grounded copilot

grounded theory
of AI-assisted programming

1. method
2. theory
3. recommendations

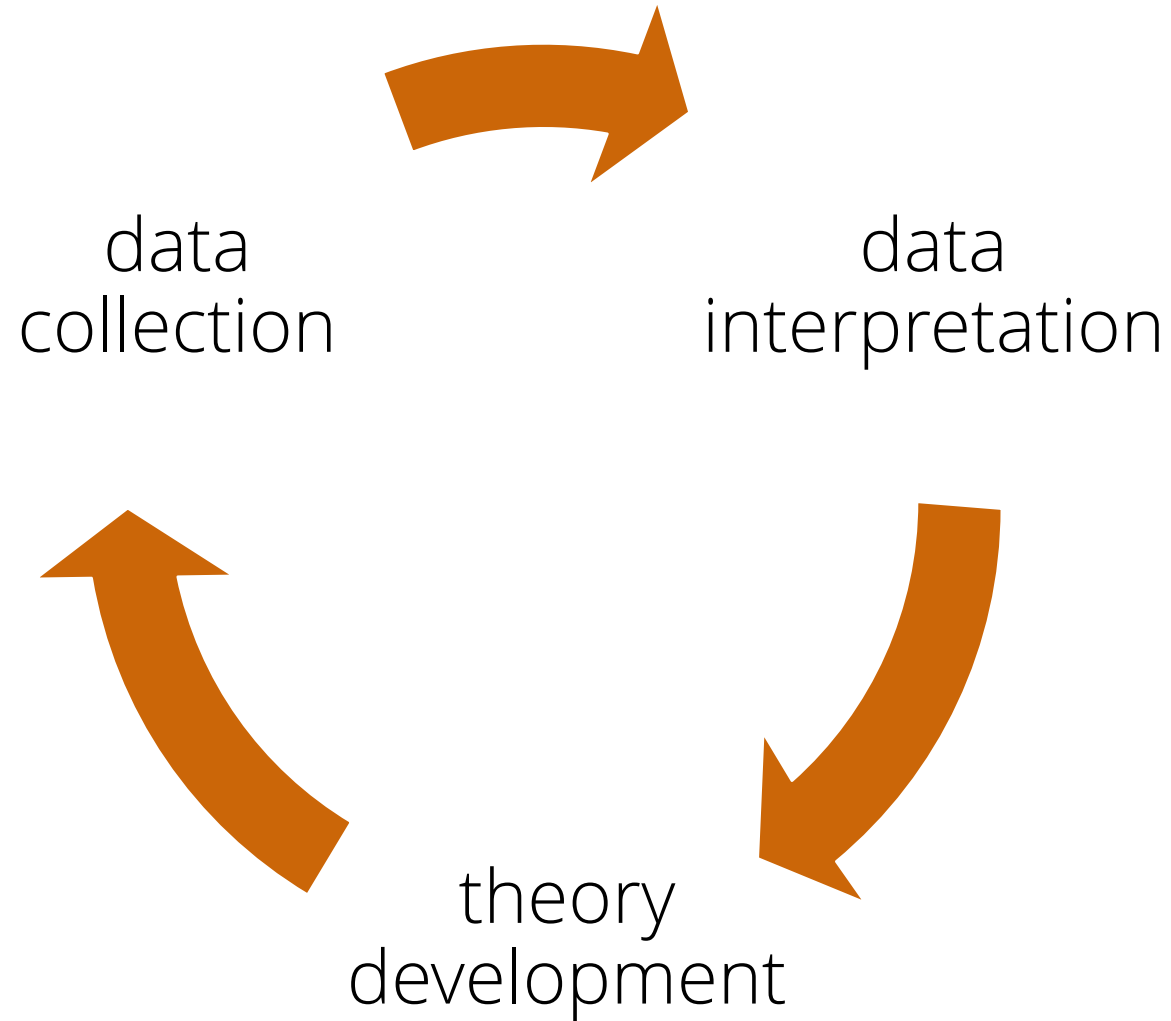
this talk

grounded copilot

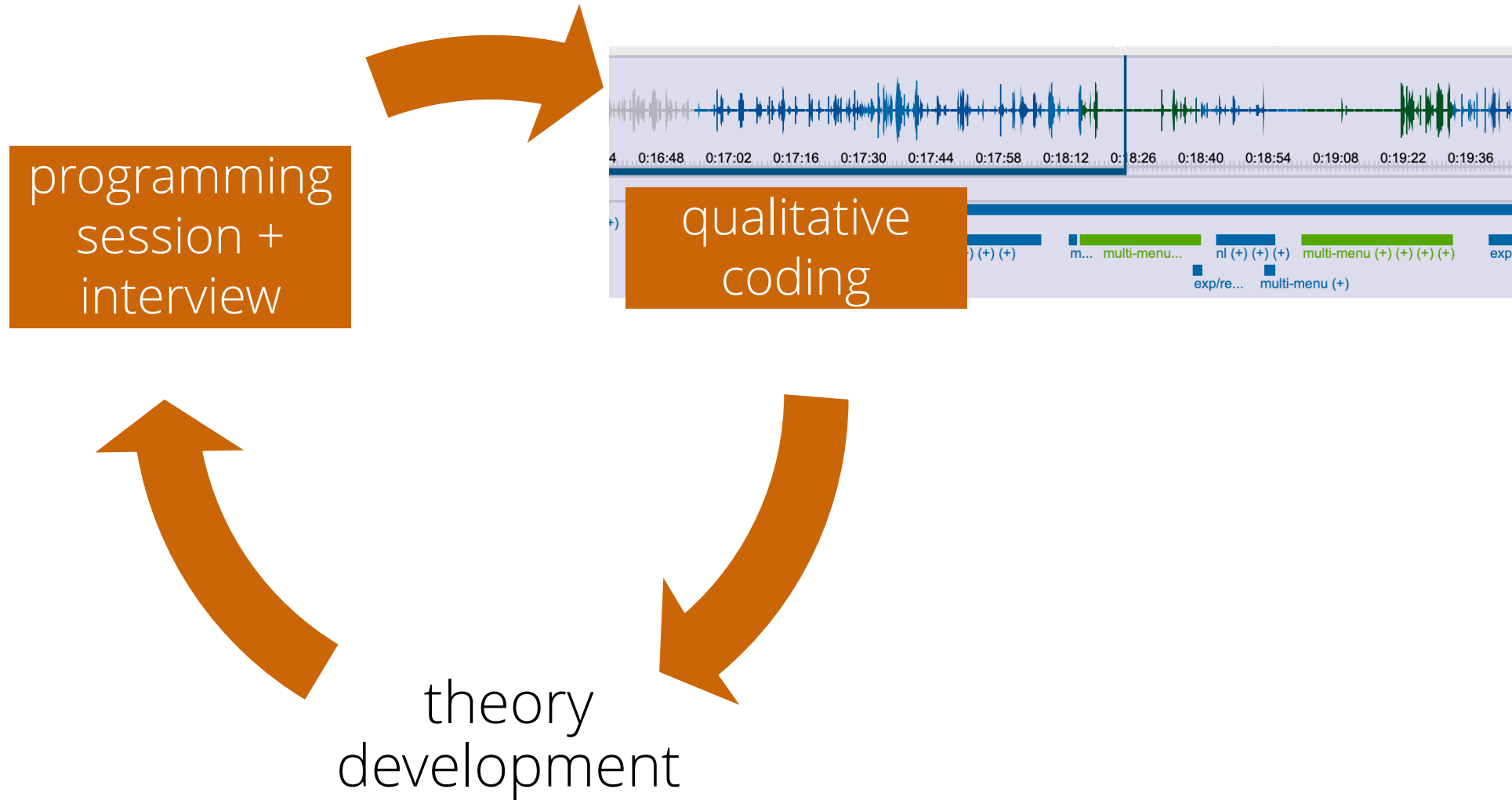
grounded theory
of AI-assisted programming

1. method

grounded theory



grounded theory



tasks

chat server

business logic of a chat app

Python/Rust

chat client

networking + custom crypto API

Python/Rust

benford's law

familiar algorithm + matplotlib

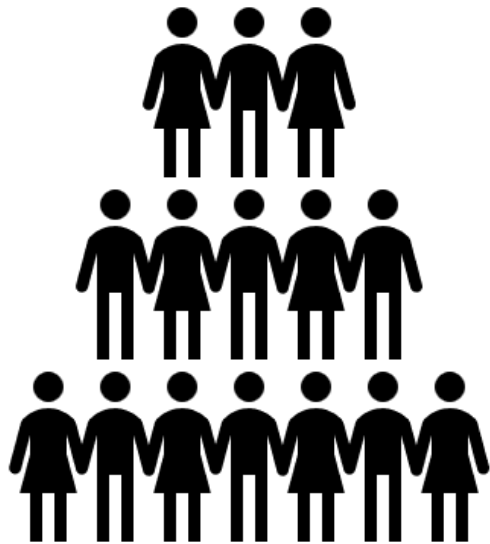
Rust + Python

string rewriting

competition task, easy to test

Python/Rust/Haskell/Java

participants



n = 20

occupation:

15 academia / 5 industry

language proficiency:

occasional / regular / professional

prior Copilot experience:

9 no / 11 yes

this talk

grounded copilot

grounded theory
of AI-assisted programming

1. method
2. theory
3. recommendations

this talk

grounded copilot

grounded theory
of AI-assisted programming

2. theory

programming, fast and slow

acceleration

autocomplete++

programmer has a plan
copilot helps them get there faster

VS

exploration

StackOverflow++

programmer is lost
copilot suggests potential solutions

programming, fast and slow

acceleration

autocomplete++

programmer has a plan
copilot helps them get there faster

acceleration: example

programmer: broke down the task,
has a good idea for this function

```
# rules are formatted like:  
# AB => C  
def parse_input(filename):  
    with open(filename) as f:  
        template, rules = f.read().split("\n\n")  
        for rule in rules:  
            rule_parts = |
```

pauses
(unintentional prompting)

acceleration: example

```
# rules are formatted like:  
# AB => C  
def parse_input(filename):  
    with open(filename) as f:  
        template, rules = f.read().split("\n\n")  
        for rule in rules:  
            rule_parts = rule.split("=> ")
```

programmer: "pattern-matches"
suggestion against expectations;
quickly accepts,
without leaving flow

copilot: auto-completes current logical unit (line of code)

programming, fast and slow

acceleration

autocomplete++

programmer has a plan
copilot helps them get there faster

VS

exploration

StackOverflow++

programmer is lost
copilot suggests potential solutions

exploration: example

programmer:
unfamiliar with matplotlib

```
You, now | 1 author (You)
1  import matplotlib
2  import matplotlib.pyplot as plt
3
4  def read_first_digits_from_file(filename):
5      with open(filename) as file:
6          data = file.read().splitlines()
7          return [int(line[0]) for line in data]
8
9  fib_first_digits = read_first_digits_from_file("fib
10 inverse_first_digits = read_first_digits_from_file(
11
12 # Plot the first digits of the Fibonacci
13 # sequence as a histogram You, now • Uncommitt
14
15
16
17
```

intentionally prompts
with a comment;
invokes side panel

exploration: example

programmer: carefully examines suggestions;
compares to gauge confidence in API usage

```
You, now | 1 author (You)
1  import matplotlib
2  import matplotlib.pyplot as plt
3
4  def read_first_digits_from_file(filename):
5      with open(filename) as file:
6          data = file.read().splitlines()
7          return [int(line[0]) for line in data]
8
9  fib_first_digits = read_first_digits_from_file("fib.")
10 inverse_first_digits = read_first_digits_from_file("
11
12 # Plot the first digits of the Fibonacci
13 # sequence as a histogram
14
15
16
17
```

```

3
4  =====
Accept Solution
5  # Plot the first digits of the Fibonacci sequence as
6  plt.hist(fib_first_digits, bins=range(0, 10))
7  plt.title("Fibonacci sequence")
8  plt.xlabel("First digit")
9  plt.ylabel("Number of occurrences")
10 plt.savefig("fib.png")
11
12  =====
13
Accept Solution
14 # Plot the first digits of the Fibonacci sequence as
15 plt.hist(fib_first_digits, bins=range(0, 10))
16 plt.title("Fibonacci sequence")
17 plt.xlabel("First digit")
18 plt.ylabel("Number of occurrences")
19 plt.show()
20
21  =====
22
Accept Solution
23 # Plot the first digits of the Fibonacci sequence as
24 plt.hist(fib_first_digits, bins=10, range=(0, 10))
25 plt.title("Fibonacci sequence")
26 plt.xlabel("First digit")
27 plt.ylabel("Number of occurrences")
28 plt.savefig("fib.png")
29
```

copilot suggests multiple alternatives

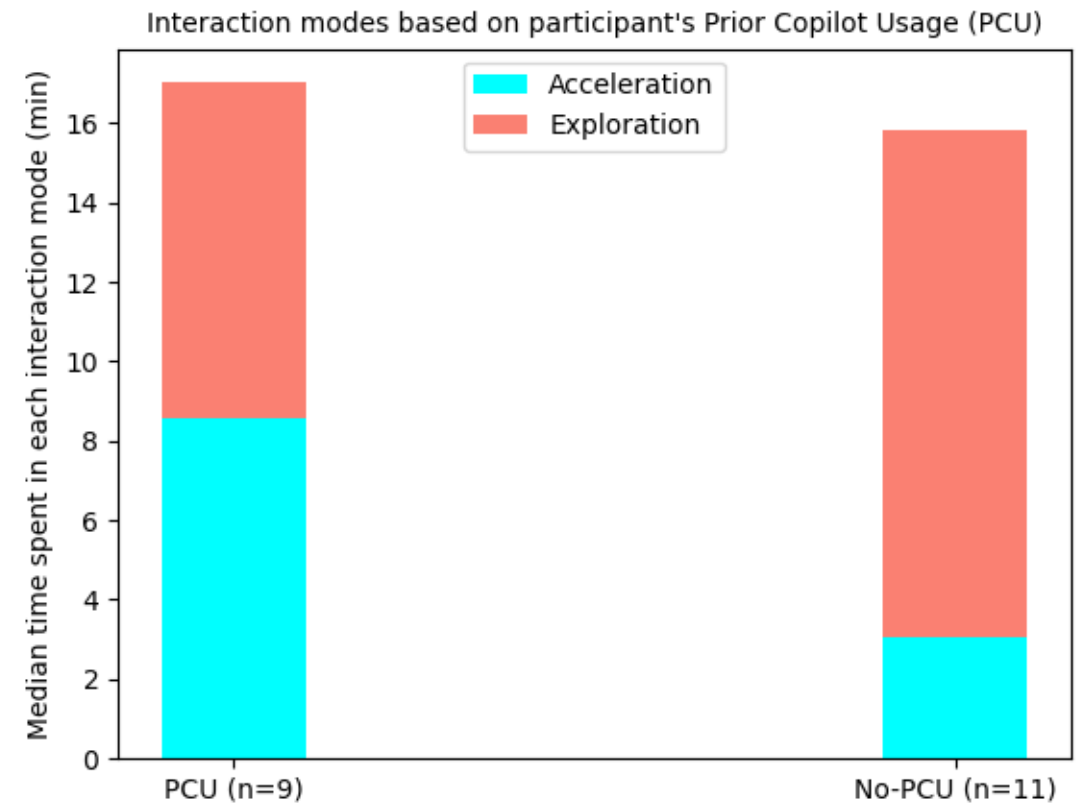
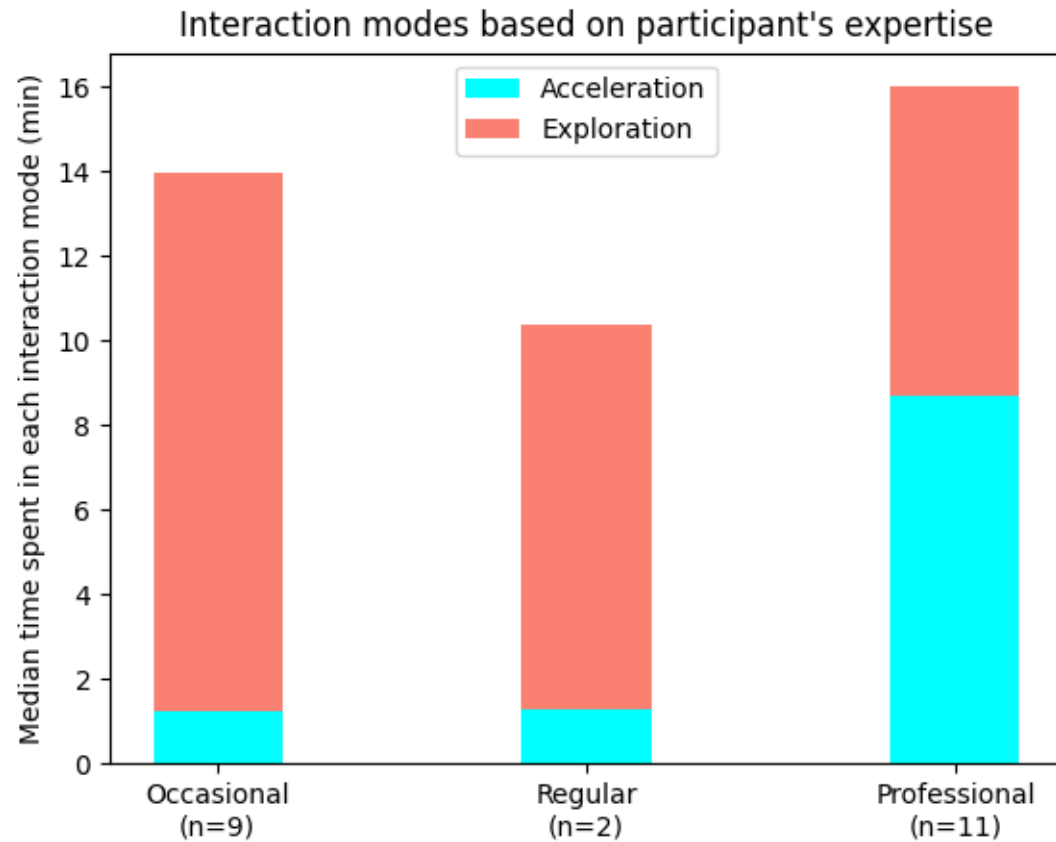
might cherry-pick parts
from different suggestions

validates code by executing
or consulting documentation

acceleration

vs

exploration



acceleration

unintentional

“pattern matching”

unit of focus
(sub-expression / statement)

unwilling to edit

vs

prompting

validation

scope

mismatch
tolerance

exploration

intentional with comments /
invoke side panel

explicit validation via
examination / execution /
documentation

entire function +
multiple alternatives

willing to edit / debug /
“rip apart” / cherry-pick

this talk

grounded copilot

grounded theory
of AI-assisted programming

1. method
2. theory
3. recommendations

this talk

grounded copilot

grounded theory
of AI-assisted programming

3. recommendations

acceleration

vs

exploration

unexpected suggestions
break flow

main
challenges

acceleration

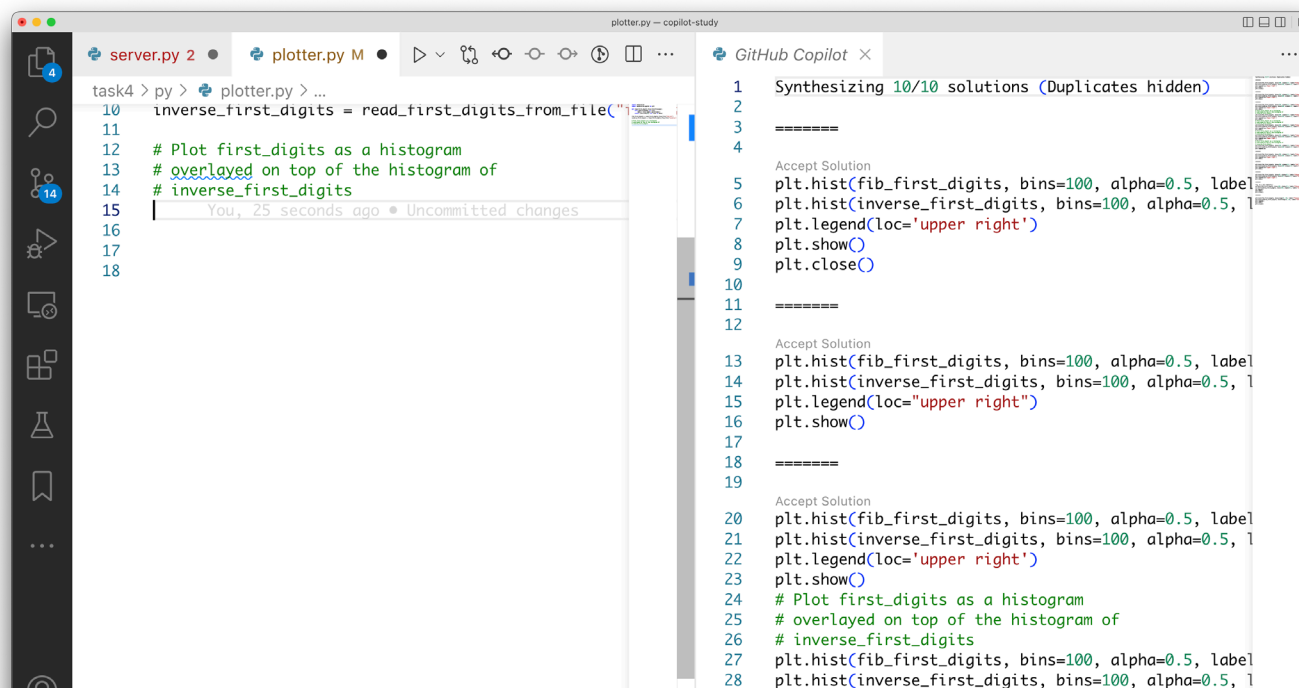
unexpected suggestions
break flow

vs

main
challenges

exploration

suggestions hard to
validate & debug
multiple suggestions hard to
distinguish



The screenshot shows a code editor with a Python script on the left and GitHub Copilot suggestions on the right. The script is named 'plotter.py' and contains the following code:

```
task4 > py > plotter.py > ...
10 inverse_first_digits = read_first_digits_from_file("
11
12 # Plot first_digits as a histogram
13 # overlaid on top of the histogram of
14 # inverse_first_digits
15
16
17
18
```

The GitHub Copilot suggestions on the right are:

```
1 Synthesizing 10/10 solutions (Duplicates hidden)
2
3 =====
4
5 Accept Solution
6 plt.hist(fib_first_digits, bins=100, alpha=0.5, label
7 plt.hist(inverse_first_digits, bins=100, alpha=0.5, l
8 plt.legend(loc='upper right')
9 plt.show()
10 plt.close()
11
12 =====
13
14 Accept Solution
15 plt.hist(fib_first_digits, bins=100, alpha=0.5, label
16 plt.hist(inverse_first_digits, bins=100, alpha=0.5, l
17 plt.legend(loc="upper right")
18 plt.show()
19
20 =====
21
22 Accept Solution
23 plt.hist(fib_first_digits, bins=100, alpha=0.5, label
24 plt.hist(inverse_first_digits, bins=100, alpha=0.5, l
25 plt.legend(loc='upper right')
26 plt.show()
27
28 # Plot first_digits as a histogram
29 # overlaid on top of the histogram of
30 # inverse_first_digits
31 plt.hist(fib_first_digits, bins=100, alpha=0.5, label
32 plt.hist(inverse_first_digits, bins=100, alpha=0.5, l
```

acceleration

vs

exploration

unexpected suggestions
break flow

1. mode awareness

suggestions hard to
validate & debug

multiple suggestions hard to
distinguish

2. simplify validation

3. better support for comparing
alternative suggestions

this talk

grounded copilot

grounded theory
of AI-assisted programming

leap

helping programmers
validate AI-generated code

this talk

idea: simplify validation using
live programming
(continuous display of runtime values)

live exploration of ai-generated programs

helping programmers
validate AI-generated code

this talk

1. demo
2. study
3. findings

leap

helping programmers
validate AI-generated code

this talk

1. demo

leap

helping programmers
validate AI-generated code

```

2
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 import io
6
7
8 def create_plot(df):
9     # create a box plot where each box refers to a group
10    ## ---
11    print(df)
12    plt.boxplot([df[df['group'] == 'control']['time'],
13                df[df['group'] == 'experiment']['time']])
14    ## ---
15
16
17 csv = '''
18 ID,group,time,success
19 P1,control,18.6,1
20 P2,control,15.42,1
21 P3,control,25.55,0
22 P4,control,12.56,0
23 P5,control,8.67,1
24 P6,experiment,7.31,0
25 P7,experiment,9.66,0
26 P8,experiment,13.64,1
27 P9,experiment,14.92,1
28 P10,experiment,18.47,1
29 '''
30
31 df = pd.read_csv(io.StringIO(csv))
32
33 plot = create_plot(df)
34

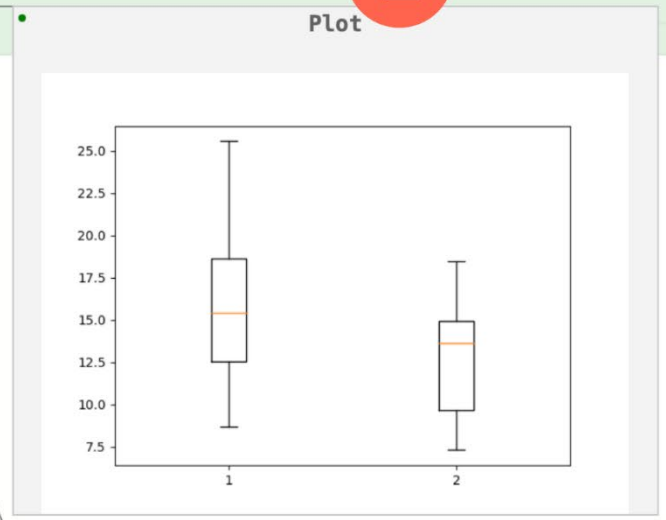
```

A

D

	df			
	ID	group	time	success
0	P1	control	18.60	1
1	P2	control	15.42	1
2	P3	control	25.55	0
3	P4	control	12.56	0
4	P5	control	8.67	1
5	P6	experiment	7.31	0
6	P7	experiment	9.66	0
7	P8	experiment	13.64	1
8	P9	experiment	14.92	1
9	P10	experiment	18.47	1

E



```

'\nID,group,time,success\nP1,control,18.6,1\nP2,control,15.42,1\nP3,control,25.55,0

```

	df			
	ID	group	time	success
0	P1	control	18.60	1
1	P2	control	15.42	1
2	P3	control	25.55	0
3	P4	control	12.56	0
4	P5	control	8.67	1

C

Suggestion 1

B

Preview

```

fig, ax = plt.subplots()
ax.boxplot([df[df.group == 'control'].time,
            df[df.group == 'experiment'].time],
            labels=['control', 'experiment'])
ax.set_title('time to success')
ax.set_ylabel('time (s)')
plt.show()

```

Suggestion 2

Preview

```

# the x-axis is the group
# the y-axis is the time
# the color of the box depends on the success

```

Suggestion 3

Preview

```

fig, ax = plt.subplots()
df.boxplot(column="time", by="group", ax=ax)

```

Suggestion 4

Preview



```

plt.boxplot([df[df['group'] == 'control']['time'],
            df[df['group'] == 'experiment']['time']])

```

this talk

1. demo
2. study
3. findings

leap

helping programmers
validate AI-generated code

this talk

2. study

leap

helping programmers
validate AI-generated code

experimental conditions

no-PB

AI suggestions
+
terminal

PB

AI suggestions
+
projection boxes

research questions

how does **live programming** affect...

1. code correctness
2. over- / under-reliance on AI
3. cognitive load
4. user impressions

tasks

API-heavy

algorithmic

multiple correct suggestions

no correct suggestions

pandas

bigrams

clean dataframe and compute stats
using pandas

find most frequent bigram in a string

fixed prompt

box plot

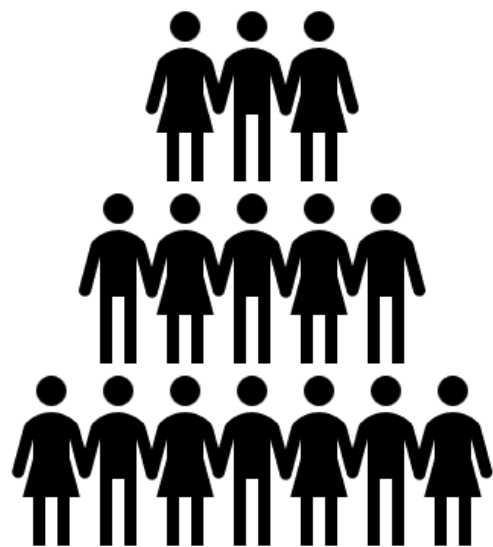
string rewriting

overlay scatter plot over boxplot
using matplotlib

parse rewrite rules and apply to string

open prompt

participants



n = 17

occupation:
15 academia / 2 industry

Python usage:
2 occasionally /
8 regularly /
7 almost every day

this talk

1. demo
2. study
3. findings

leap

helping programmers
validate AI-generated code

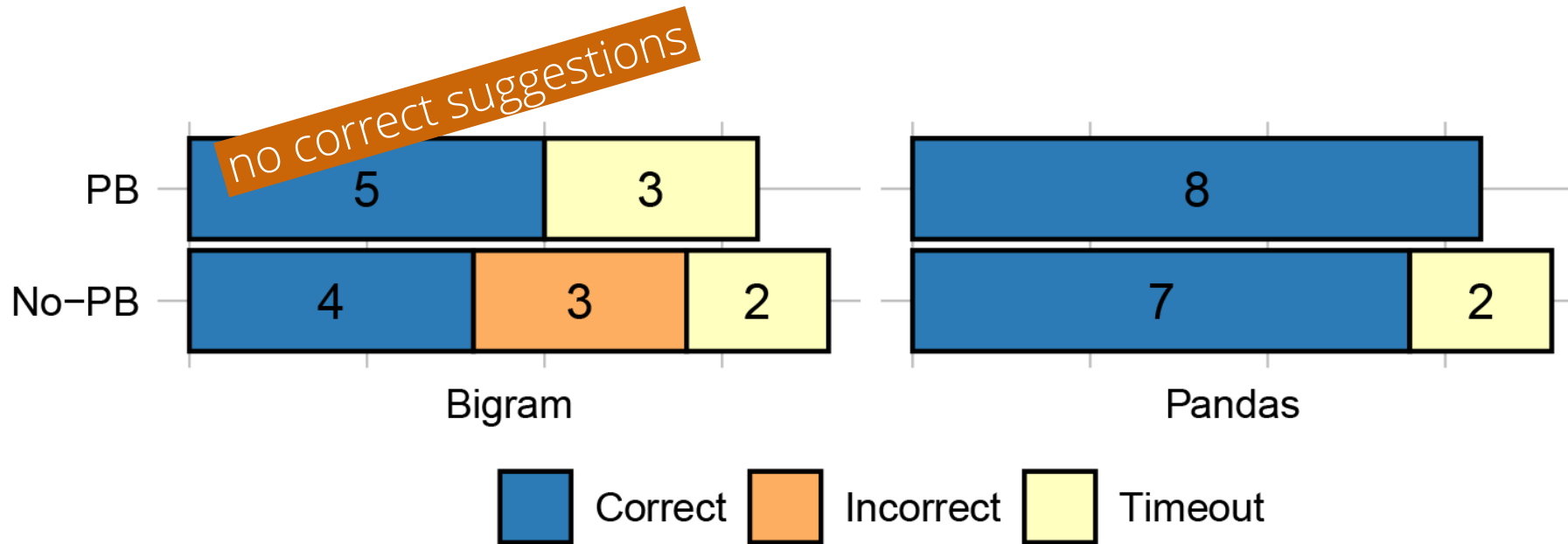
this talk

3. findings

leap

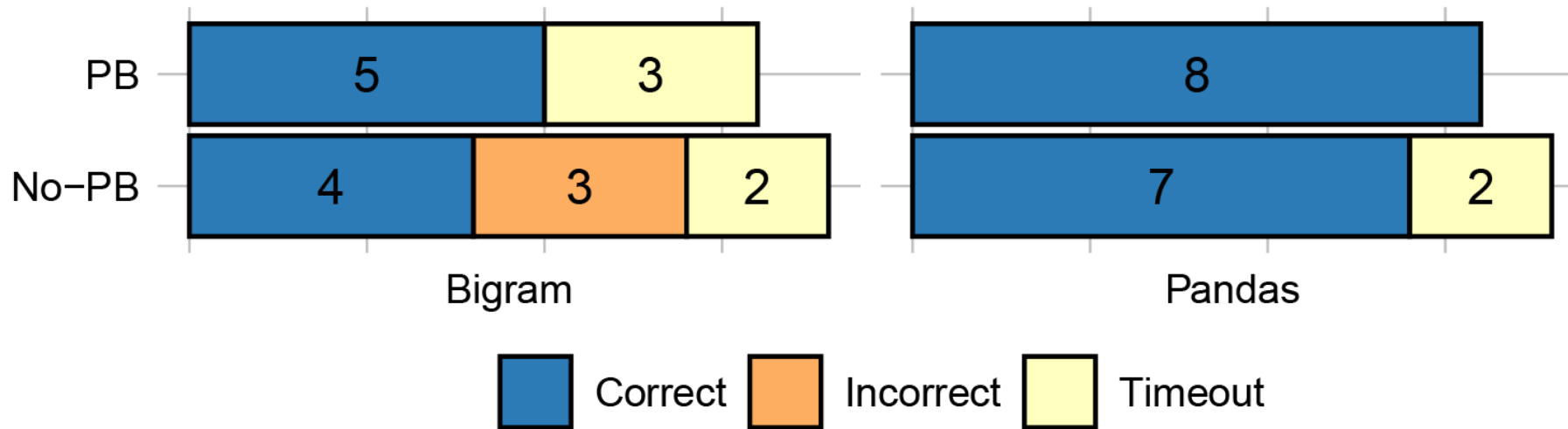
helping programmers
validate AI-generated code

rq1: correctness



leap helps validate suggestions
(but does not help fix incorrect ones)

rq2: over-/under-reliance



6 no-PB vs 0 PB participants mid-judged correctness of their solution

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

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)

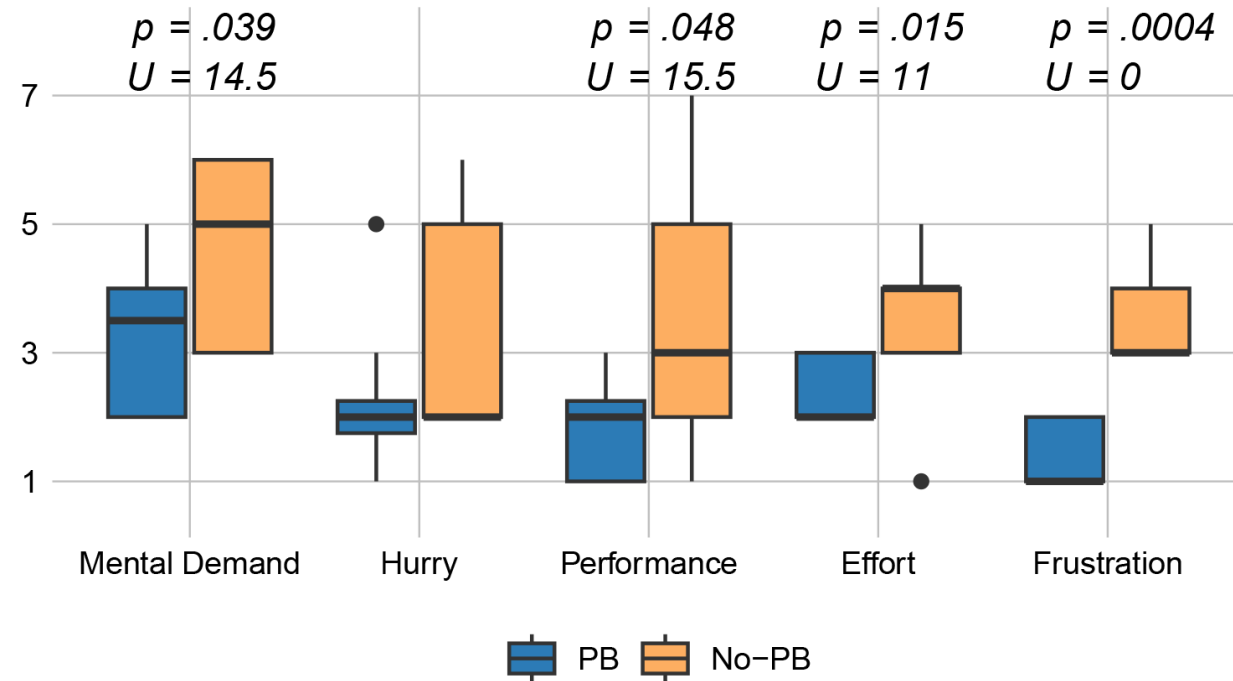
“it **saved me the effort** of writing multiple print statements.” (P1)

6 no-PB vs 0 PB participants **mid-judged** correctness of their solution

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

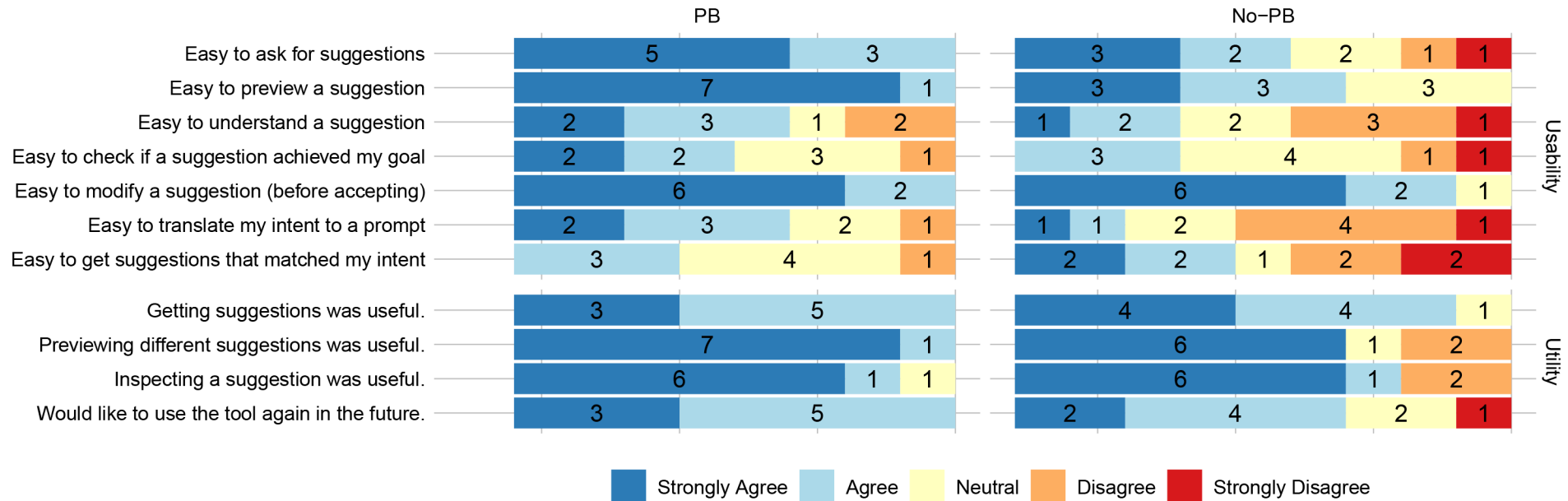
rq3: cognitive load

NASA TLX cognitive load metrics on Pandas



leap significantly reduced cognitive load of exploring AI suggestions on tasks amenable to validation by execution

rq3: user impressions



“Being able to preview, edit, and look at the projection boxes before accepting a snippet was **very helpful when choosing** between multiple suggestions.” (P1)

users found leap more usable and useful

this talk

grounded copilot

grounded theory
of AI-assisted programming

[OOPSLA'23]

leap

helping programmers
validate AI-generated code

[under review]