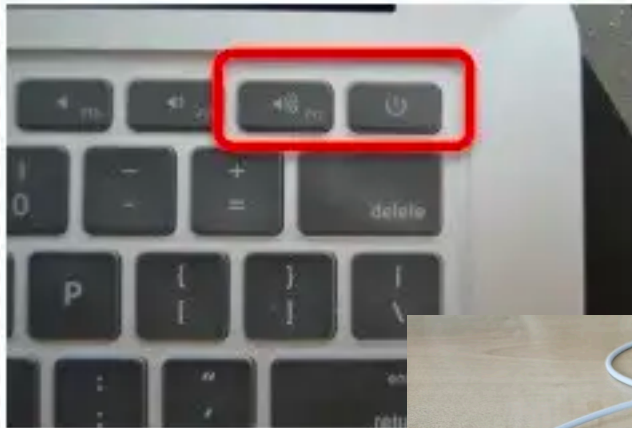


QA: User Study

17-313 Fall 2022

Rohan Padhye, Michael Hilton, Chris Timperley, and Daye Nam

Why do we care about users?



Human Centered Methods

- Contextual Inquiry
- Paper prototypes
- Think-aloud protocols
- Heuristic Evaluation
- Affinity diagrams
- Personas
- ...
- A/B testing
- Body storming
- Questionnaires
- Surveys
- Log analysis
- Card sorting
- ...

How can we test usability



Activity

How can you test usefulness / usability of CoPilot?



GitHub Copilot

Usability Evaluation Case Study

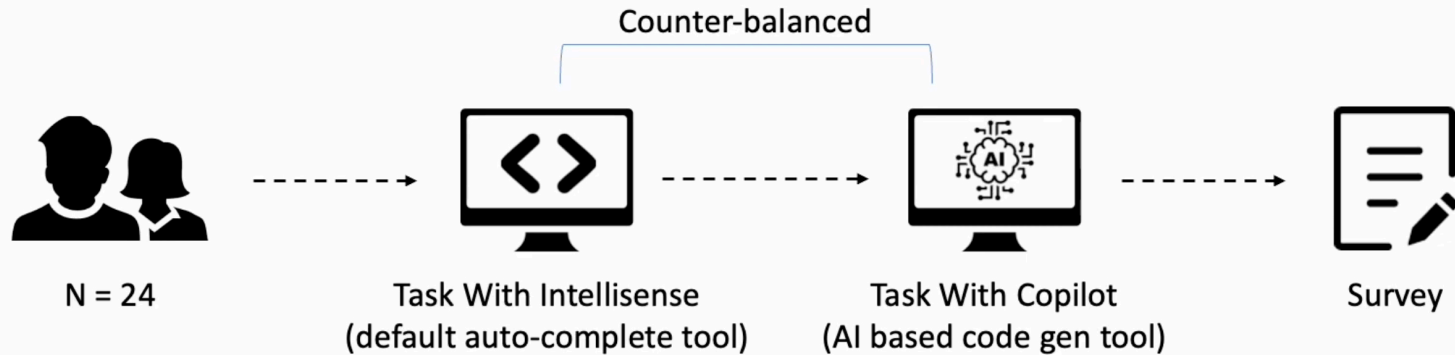
Expectation vs. Experience: Evaluating the Usability of Code Generation Tools Powered by Large Language Models

Priyan Vaithilingam
pvaithilingam@g.harvard.edu
Harvard University
USA

Tianyi Zhang
tianyi@purdue.edu
Purdue University
USA

Elena L. Glassman
glassman@seas.harvard.edu
Harvard University
USA

User Study Procedure



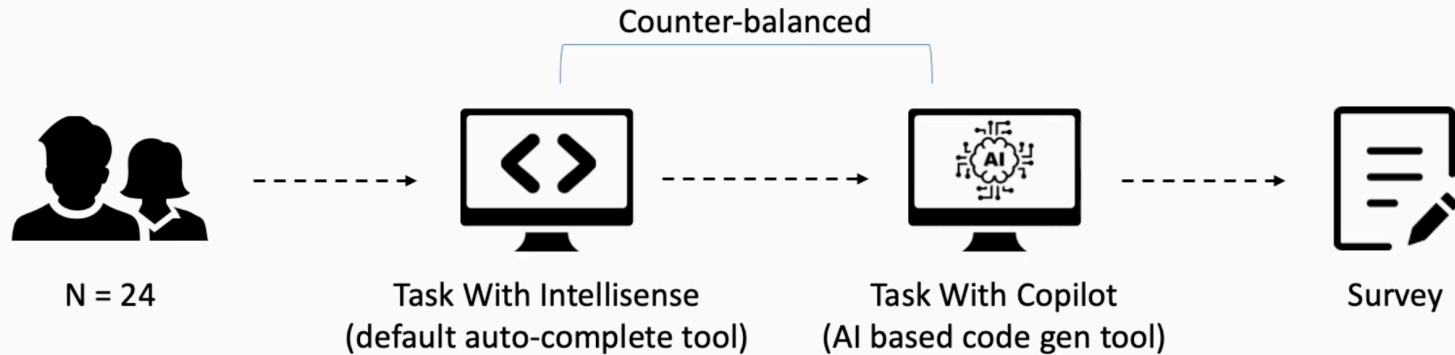
Participants

- 24 participants
 - 4 Female, 19 Male, 1 Non-binary
 - 5 master's students, 8 Ph.D. students, 1 software engineer
 - 1: ≤ 2 years, 14: 2-5 years, 9: > 5 years
- \$20 Amazon giftcard as compensation

Task Design

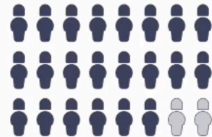
- Real-world Python programming tasks with different levels of difficulty
 - T1: Edit CSV (Easy)
 - T2: Web Scrapping (Medium)
 - T3: Graph Plotting (Hard)

User Study Procedure



Quantitative Results

Using Intellisense



22 of 24

participants finished the task successfully.



Average Time

Using Copilot



19 of 24

participants finished the task successfully.



Average Time

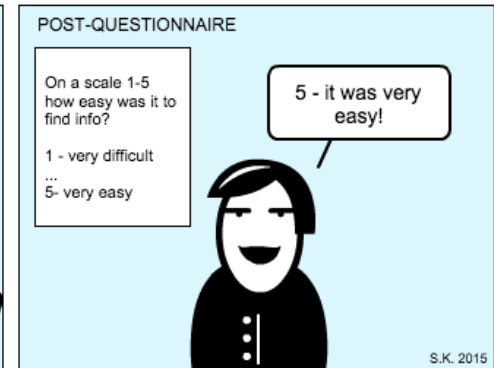
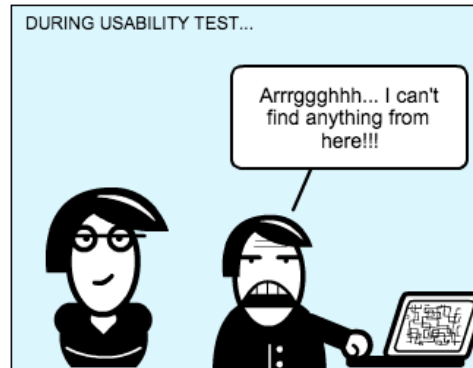


19 of 24

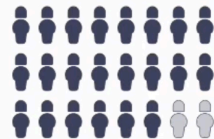
participants **prefer using Copilot**

Not statistically significant

Quantitative Results



Using Intellisense



22 of 24

participants finished the task successfully.



Average Time

Using Copilot



19 of 24

participants finished the task successfully.



Average Time



19 of 24

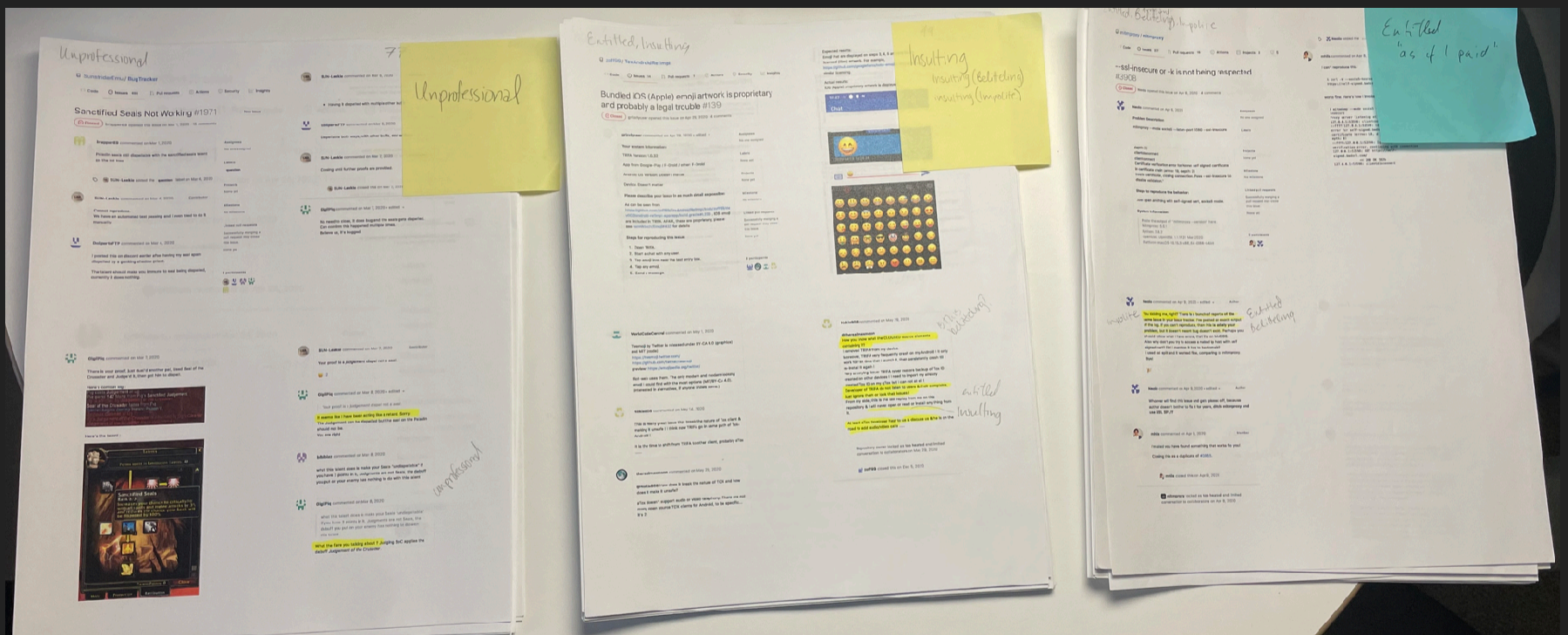
participants prefer using Copilot

Not statistically significant

+ Qualitative analysis

- ▶ Piles of qualitative data, mostly **text**
 - ▶ What to do with it?
 - ▶ From journalism to science - how?
- ▶ **Step 1: Abstraction**
 - ▶ Attach “codes” (labels) to chunks of data
 - ▶ Characterize / summarize the data
- ▶ **Step 2: Finding patterns**
 - ▶ Use these abstractions to find meta-patterns, craft a theory (“**grounded theory**”), ...
 - ▶ Interpret the data
- ▶ This is difficult, but very doable with practice

+ Qualitative analysis



Miller, C., Cohen, S., Klug, D., Vasilescu, B., & Kästner, C. (2022). "Did You Miss My Comment or What?" Understanding Toxicity in Open Source Discussions. In In 44th International Conference on Software Engineering (ICSE'22).

+ Qualitative analysis

¹ As I walked toward the school, there was a 7-11 convenience store 1 block away, next to a small professional office building: an optometrist and other medical/health-related clinics. The street was an empty lot, but next to a Burger King restaurant.

¹ BUSINESSES

Well, that's one problem, that [my school is] pretty small, so ¹ if you say one thing to one person, and then they decide to tell two people, then those two people tell two people, and in one period everybody else knows. ² Everybody in the entire school knows that you said whatever it was. So . . .

¹ SPREADING RUMORS

² KNOWING WHAT YOU SAID

I ¹ hated school last year. It was awful, I hated it. And I, um, don't know why. I guess, over the summer I kind of ³ stopped caring about what other people thought and cared more about, just, I don't know.

³ "STOPPED CARING"

Qualitative Insights



Copilot provides a **useful starting point** for users unlike Intellisense



12 users found code generated by copilot **hard to understand and change**



5 users found it **hard to trust code** generated by copilot



Users of Copilot **constantly switch modes** between reading and writing code



Many users used copilot as a **substitution for internet search**



8 users **over-relied on copilot**, i.e., **ignored validating** the code generated by copilot

Design Implications



Several participants requested to see **multiple code suggestions** for a given prompt



Users relying on internet search compare multiple sources unlike users using copilot. It is worthwhile **integrating online search with code generation**



We observed copilot is more accurate for simpler prompts. It's worth exploring **task decomposition** for better accuracy



Explanations with inline comments for code generated by copilot can help in understanding



We can also help users debug the code by **automatically generating test cases** and test data for users to validate

Activity

What are the other aspects to consider?



GitHub Copilot