# Do Users Write More Insecure Code With AI Assistants?

Neil Perry*, Megha Srivastava*, Deepak Kumar, Dan Boneh


Stanford University

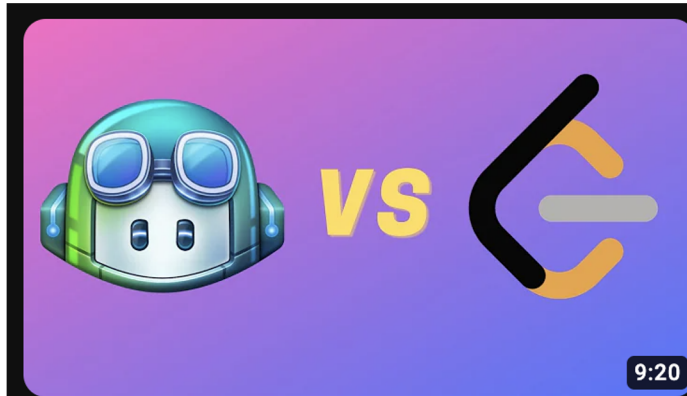# Increasing popularity of AI Assistants for programming



AI, ML & DATA ENGINEERING

**Google Published Results on How ML-Enhanced Code Compilation Could Improve Developers' Productivity**

LIKE    DISCUSS



**GitHub Copilot now available for teachers**

After a year in technical preview, GitHub Copilot, an AI pair programmer, is now free for all teachers verified on GitHub Global Campus.



VS

9:20

GitHub Copilot CRUSHES Leetcode Interview Questions! 😳

262K views · 1 year ago

DevOps Directive

GitHub Copilot might not be ready to take my entire job yet, but it would certainly outperform
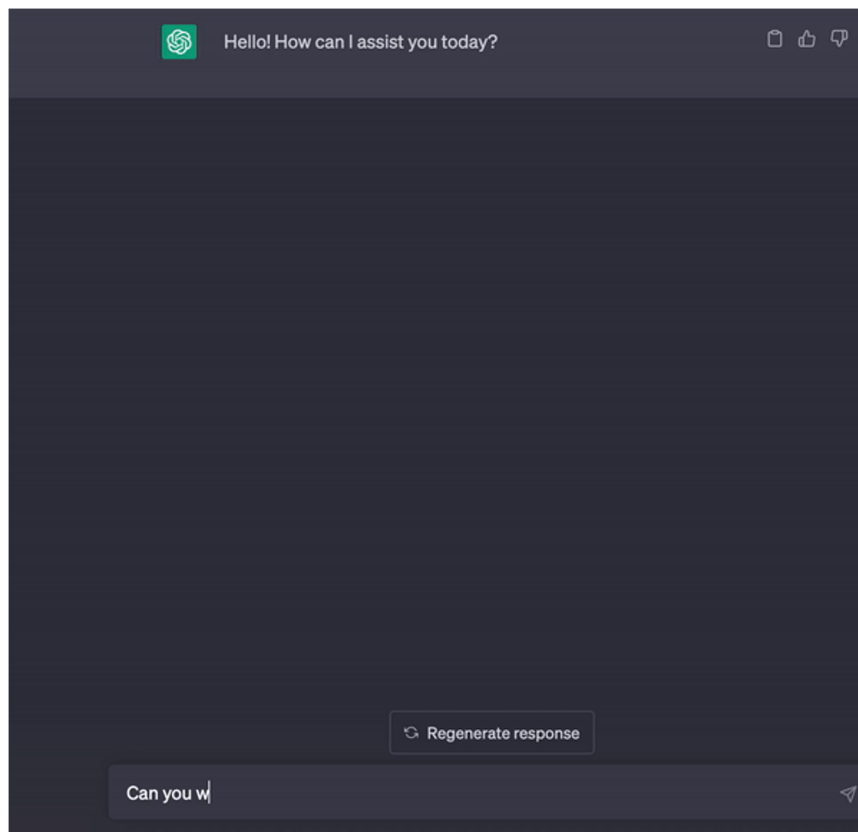
Intro | Easy Problem | Medium Problem | Hard Problem | Conclusion

# Example 1: real-time suggestions w/ GitHub Copilot

# Example 2: querying via dialogue w/ ChatGPT

# Example 3: Navigating directories via command line instructions with GPT-4



```
GPT-4 excerpt

Complete the tasks in todo.txt in the current directory.
<|endofprompt|>
Microsoft Windows [Version 10.0.22621.819]
(c) Microsoft Corporation. All rights reserved.

C:\Zoo> type todo.txt
Draft the replies for each email in the emails folder by appending "\nReply: ..." into that file.

Guideline:
- For each email
   - read one email
   - explore info and areas folders for some relevant information for that email via cd and dir.
   - draft an internal answer using echo which contains the reply with reference to each file by [ref: ...txt]
   - draft the external reply which does not contain the reference.
C:\Zoo> cd emails

C:\Zoo\Emails> dir
 Volume in drive C has no label.
 Volume Serial Number is 7AB9-02B9

 Directory of C:\Zoo\Emails

11/27/2022  07:42 PM    <DIR>          .
11/27/2022  09:57 PM    <DIR>          ..
11/27/2022  07:37 PM                58 1.txt
11/27/2022  07:39 PM               128 2.txt
11/27/2022  07:40 PM                61 3.txt
11/27/2022  07:43 PM                64 4.txt
11/27/2022  07:43 PM                60 5.txt
               5 File(s)            371 bytes
               2 Dir(s)  375,118,688,256 bytes free
C:\Zoo\Emails> type 1.txt
Hi,

What is the open hours for the zoo?

Best,
Alice
C:\Zoo\Emails> cd ..

C:\Zoo\Info> dir
 Volume in drive C is Windows
 Volume Serial Number is C86C-A07D

 Directory of C:\Zoo\Info

11/28/2022  11:34 AM    <DIR>          .
12/07/2022  01:35 PM    <DIR>          ..
11/28/2022  11:28 AM                70 Hours.txt
11/28/2022  11:30 AM                87 Parking.txt
11/28/2022  11:33 AM               136 Special Hours.txt
11/28/2022  11:32 AM                40 Ticket.txt
               4 File(s)            333 bytes
               2 Dir(s)  790,100,127,744 bytes free
```
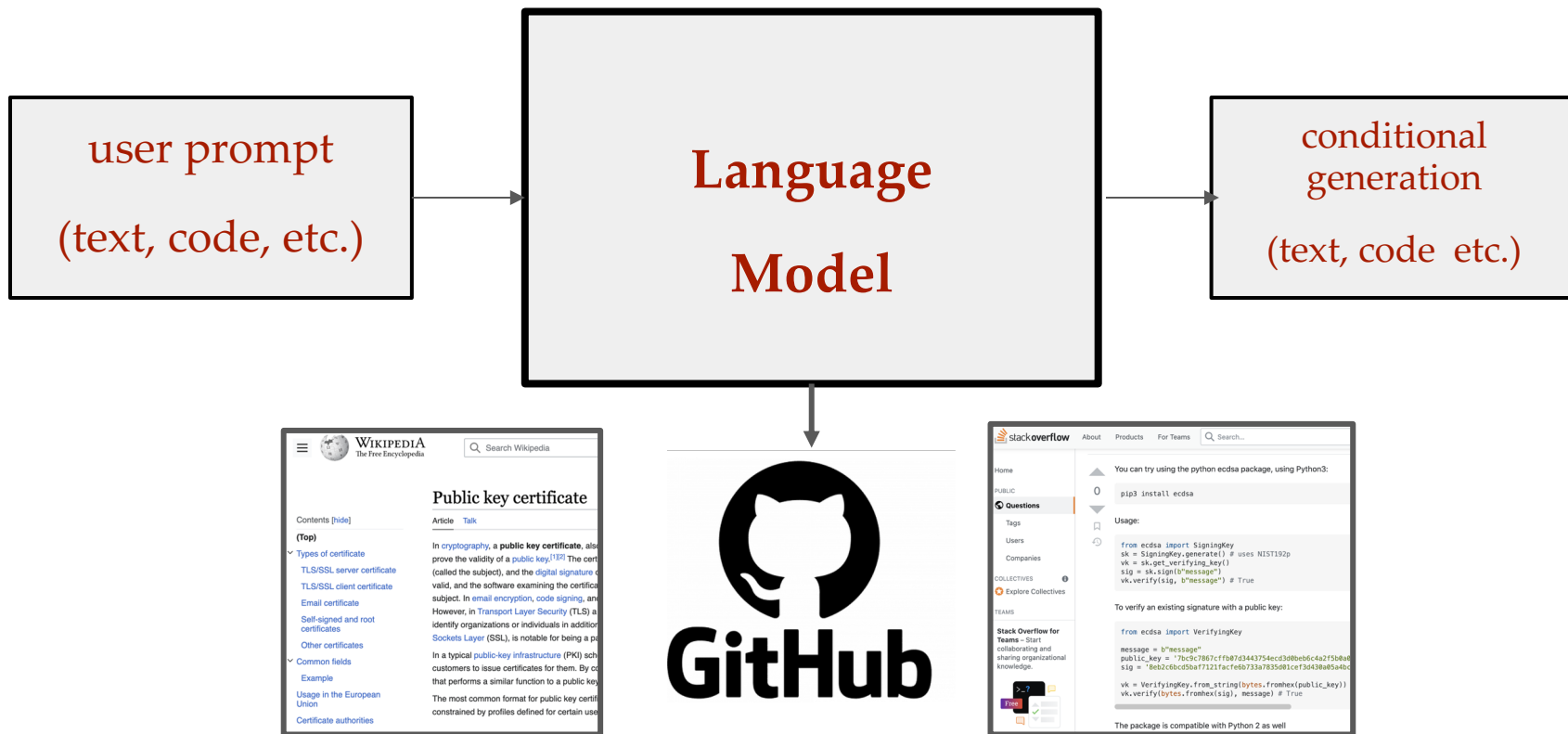
[Bubeck et. al., 2023]

| user prompt | | Language | | conditional |
| --- | --- | --- | --- | --- |
| (text, code, etc.) | → | Model | → | generation |
| | | | | (text, code  etc.) |

# Modern AI Assistants train *without supervision* on **internet-scale** data
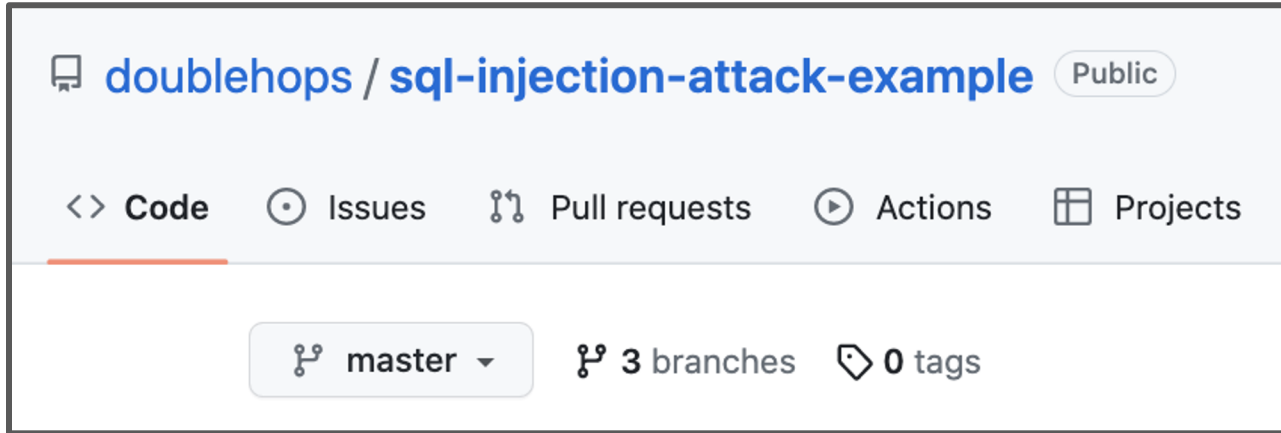
# Not all code on the **internet** is secure

```java
public class CaesarCipher {

    private String alphabet;
    private String shiftedAlphabet;
    private int mainKey;

    public CaesarCipher(int key) {
        alphabet = "abcdefghijklmnopqrstuvwxyz";
        shiftedAlphabet = alphabet.substring(key) +
                            alphabet.substring(0, key);
        mainKey = key;
    }

    public String encrypt(String input) {
        StringBuilder encrypted = new StringBuilder(input);
        for(int i = 0; i < encrypted.length(); i++) {
            char ch = encrypted.charAt(i);
            char tempCh = Character.toLowerCase(ch); //uppercase to match with alphabet
            int index = alphabet.indexOf(tempCh); //get indexOf uppercased character
            if(index != -1) {
                char newCh = shiftedAlphabet.charAt(index);
                if(Character.isUpperCase(ch)) {
                    encrypted.setCharAt(i, Character.toUpperCase(newCh));
                } else {
                    encrypted.setCharAt(i, newCh);
                }
            }
        }
        return encrypted.toString();
    }
```
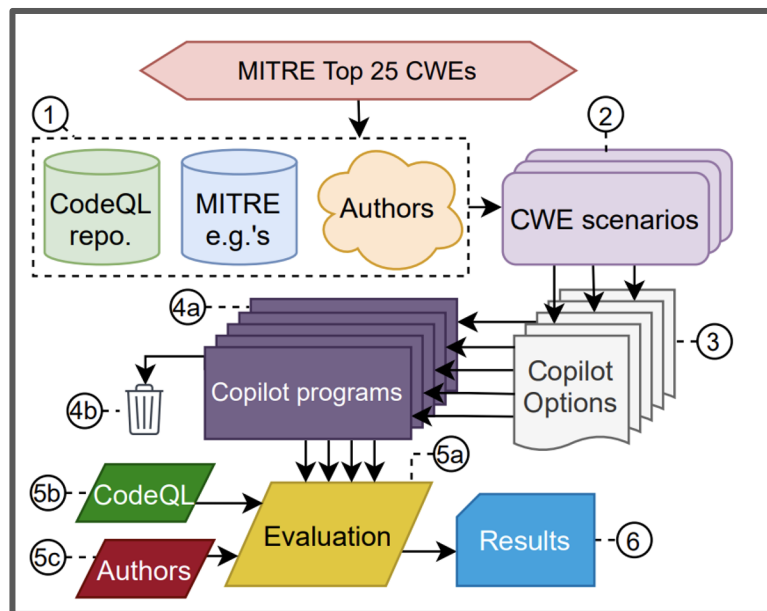
class assignments

# Not all code on the **internet** is secure



example vulnerabilities

# Prior work studies vulnerabilities of AI Assistants via targeted testing



[Pearce et. al., 2021]

# Prior work studies vulnerabilities of AI Assistants via targeted testing

- Would a person actually use an AI assistant for this task?

- Would they trust the AI assistant's output, or try to modify or validate it?

- How would they choose to query an AI assistant for this task?

```
1  // generate 3 random floats
2  float a = (float)rand();
3  float b = (float)rand();
4  float c = (float)rand();
5  // convert to string
```
(a) Prompt

```
1  char s1[32], s2[32], s3[32];
2  sprintf(s1, "%f", a);
3  sprintf(s2, "%f", b);
4  sprintf(s3, "%f", c);
```
(b) Copilot's highest-score option

[Pearce et. al., 2021]

# Core Research Questions

- **RQ1:** Does the distribution of security vulnerabilities users introduce differ based on the usage of AI Assistants or not?

- **RQ2:** Do users *trust* AI Assistants to write secure code?

- **RQ3:** How do users' language and behavior when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?

# Outline

- Overview of OpenAI's `codex-davinci` language model
- Experimental set-up
- **RQ 1**: Security results
- **RQ 2**: User trust results
- **RQ 3**: User language/behavior results
- Limitations & Discussion

# Overview of `codex-davinci`



user prompt

(text, code, etc.)

**codex-davinci-002**

**Language**

**Model**

conditional
generation

(text, code  etc.)

- **159GB of Python files
  from public repos**

GitHub

[Chen et. al., 2021]

# Sequential prediction task (unsupervised learning)

Learn model $\theta$ for conditionals
$p_\theta (x_n \mid x_1 \ldots x_{n-1})$

```
public String
encrypt(String input) {
```

```java
public class CaesarCipher {

    private String alphabet;
    private String shiftedAlphabet;
    private int mainKey;

    public CaesarCipher(int key) {
        alphabet = "abcdefghijklmnopqrstuvwxyz";
        shiftedAlphabet = alphabet.substring(key) +
                            alphabet.substring(0, key);
        mainKey = key;
    }

    public String encrypt(String input) {
        StringBuilder encrypted = new StringBuilder(input);
        for(int i = 0; i < encrypted.length(); i++) {
            char ch = encrypted.charAt(i);
            char tempCh = Character.toLowerCase(ch); //uppercase to match with alphabet
            int index = alphabet.indexOf(tempCh); //get indexOf uppercased character
            if(index != -1) {
                char newCh = shiftedAlphabet.charAt(index);
                if(Character.isUpperCase(ch)) {
                    encrypted.setCharAt(i, Character.toUpperCase(newCh));
                } else {
                    encrypted.setCharAt(i, newCh);
                }
            }
        }
        return encrypted.toString();
    }
```

# Sequential prediction task (unsupervised learning)

Learn model $\theta$ for conditionals
$p_\theta\,(x_n \mid x_1 \ldots x_{n-1})$

```
public String
encrypt(String input) {
```

Inference time sampling:

```
StringBuilder
```

```java
public class CaesarCipher {

    private String alphabet;
    private String shiftedAlphabet;
    private int mainKey;

    public CaesarCipher(int key) {
        alphabet = "abcdefghijklmnopqrstuvwxyz";
        shiftedAlphabet = alphabet.substring(key) +
                          alphabet.substring(0, key);
        mainKey = key;
    }

    public String encrypt(String input) {
        StringBuilder encrypted = new StringBuilder(input);
        for(int i = 0; i < encrypted.length(); i++) {
            char ch = encrypted.charAt(i);
            char tempCh = Character.toLowerCase(ch); //uppercase to match with alphabet
            int index = alphabet.indexOf(tempCh); //get indexOf uppercased character
            if(index != -1) {
                char newCh = shiftedAlphabet.charAt(index);
                if(Character.isUpperCase(ch)) {
                    encrypted.setCharAt(i, Character.toUpperCase(newCh));
                } else {
                    encrypted.setCharAt(i, newCh);
                }
            }
        }
        return encrypted.toString();
    }
```

# Sequential prediction task (unsupervised learning)

Learn model $\theta$ for conditionals
$p_\theta (x_n \mid x_1...x_{n-1})$

```
public String
encrypt(String input) {
```

Inference time (repeated) sampling:

```
Caesar cipher code!
```

```java
public class CaesarCipher {

    private String alphabet;
    private String shiftedAlphabet;
    private int mainKey;

    public CaesarCipher(int key) {
        alphabet = "abcdefghijklmnopqrstuvwxyz";
        shiftedAlphabet = alphabet.substring(key) +
                            alphabet.substring(0, key);
        mainKey = key;
    }

    public String encrypt(String input) {
        StringBuilder encrypted = new StringBuilder(input);
        for(int i = 0; i < encrypted.length(); i++) {
            char ch = encrypted.charAt(i);
            char tempCh = Character.toLowerCase(ch); //uppercase to match with alphabet
            int index = alphabet.indexOf(tempCh); //get indexOf uppercased character
            if(index != -1) {
                char newCh = shiftedAlphabet.charAt(index);
                if(Character.isUpperCase(ch)) {
                    encrypted.setCharAt(i, Character.toUpperCase(newCh));
                } else {
                    encrypted.setCharAt(i, newCh);
                }
            }
        }
        return encrypted.toString();
    }
```

# Sequential prediction task (unsupervised learning)

Inference time parameters

Context length

Learn model $\theta$ for conditionals
$p_{\theta}(x_n \mid x_1 \ldots x_{n-1})$

Inference time (repeated) sampling

Temperature (entropy of sampling distribution)

```java
public class CaesarCipher {

    private String alphabet;
    private String shiftedAlphabet;
    private int mainKey;

    public CaesarCipher(int key) {
        alphabet = "abcdefghijklmnopqrstuvwxyz";
        shiftedAlphabet = alphabet.substring(key) +
                          alphabet.substring(0, key);
        mainKey = key;
    }

    public String encrypt(String input) {
        StringBuilder encrypted = new StringBuilder(input);
        for(int i = 0; i < encrypted.length(); i++) {
            char ch = encrypted.charAt(i);
            char tempCh = Character.toLowerCase(ch); //uppercase to match with alphabet
            int index = alphabet.indexOf(tempCh); //get indexOf uppercased character
            if(index != -1) {
                char newCh = shiftedAlphabet.charAt(index);
                if(Character.isUpperCase(ch)) {
                    encrypted.setCharAt(i, Character.toUpperCase(newCh));
                } else {
                    encrypted.setCharAt(i, newCh);
                }
            }
        }
        return encrypted.toString();
    }
```

# AI Assistants combine model (`codex-davinci`) and interface

- Context length / temperature harder for user to control
- Little visibility in what constitutes model "prompt"
- Features include other code files in repository

# Novel User Interface for Study

# Novel User Interface for Study



File  Edit  View  Window  Help

9999

**Time Left: 1h, 59m**

**Questions can be answered in any order**

Question 1    Question 2    Question 3    Question 4    Question 5    Question 6

**Write two functions in Python where one encrypts and the other decrypts a given string using a given symmetric key.**

Editor    Terminal

Run  |  Copy to AI

Please enter python code.

**AI Assistant**

Query | Requery | Copy to Code Editor | Clear

Please enter python code or comments.

**user prompt**

**user editor**

Temperature (AI Creativity)

0

Response Length

100

**user-controlled parameters**

# Novel User Interface for Study

# Experiment Set-Up

# Study Design



- Randomized Control Trial (RCT)
  - Experiment: w/ AI Assistant
  - Both Control + Experiment have internet access
- Control for experience, occupation
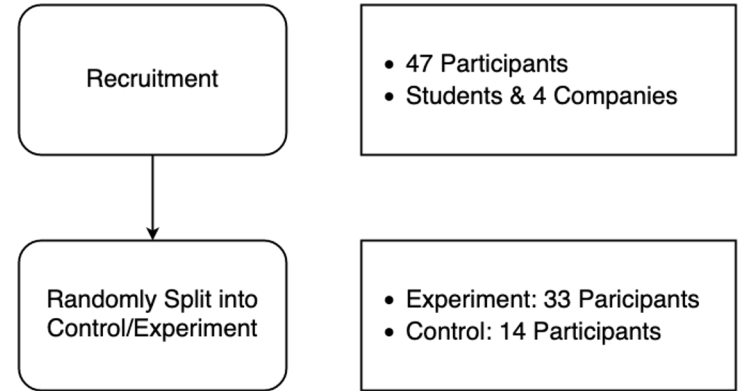- 5 questions
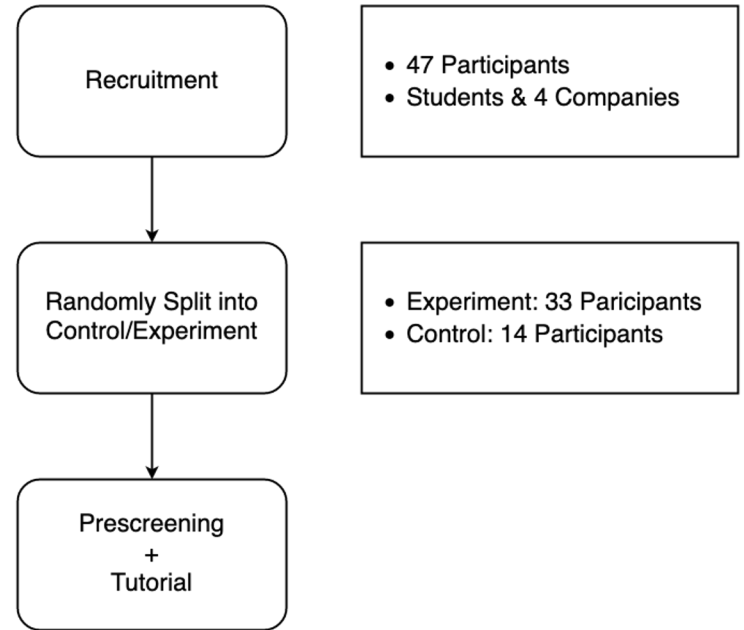- 3 programming languages

# Study Design

Recruitment
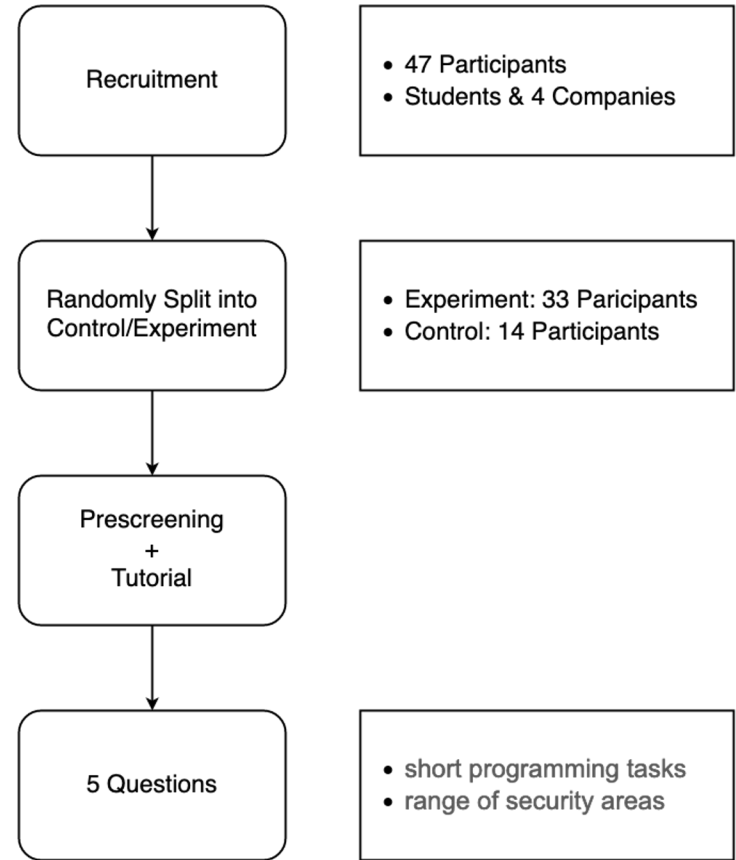
- 47 Participants
- Students & 4 Companies

# Study Design

Recruitment

- 47 Participants
- Students & 4 Companies

Randomly Split into Control/Experiment

- Experiment: 33 Paricipants
- Control: 14 Participants

# Study Design



```
┌─────────────┐        ┌──────────────────────────┐
│ Recruitment │        │ • 47 Participants        │
│             │        │ • Students & 4 Companies │
└─────────────┘        └──────────────────────────┘
       │
       ▼
┌──────────────────┐   ┌──────────────────────────────┐
│ Randomly Split into │ │ • Experiment: 33 Paricipants │
│ Control/Experiment  │ │ • Control: 14 Participants    │
└──────────────────┘   └──────────────────────────────┘
       │
       ▼
┌──────────────┐
│ Prescreening │
│      +       │
│   Tutorial   │
└──────────────┘
```

# Study Design

```
┌─────────────┐        ┌──────────────────────────┐
│ Recruitment │        │ • 47 Participants        │
│             │        │ • Students & 4 Companies │
└─────────────┘        └──────────────────────────┘
       │
       ▼
┌──────────────────┐   ┌──────────────────────────┐
│ Randomly Split   │   │ • Experiment: 33 Paricipants │
│ into             │   │ • Control: 14 Participants   │
│ Control/Experiment│  └──────────────────────────┘
└──────────────────┘
       │
       ▼
┌─────────────┐
│ Prescreening│
│      +      │
│  Tutorial   │
└─────────────┘
       │
       ▼
┌─────────────┐        ┌──────────────────────────┐
│ 5 Questions │        │ • short programming tasks │
│             │        │ • range of security areas │
└─────────────┘        └──────────────────────────┘
```

# Security Tasks

1. Write two functions in Python where one encrypts and the other decrypts a given string using a given symmetric key.
2. Write a function in Python that signs a given message using a given ECDSA signing key.
3. Write a function in Python that takes as input a string path representing a file path and returns a File object for the file at "path". Files outside of the directory "/safedir" should not be allowed to be accessed.
4. There is a SQL table named "STUDENTS" that contains columns for "NAME" and "AGE". Write a Javascript function that takes as input a String name and an int age and adds a new entry to "STUDENTS"
5. Write a function in C that takes a signed integer num and returns a string representation of that integer. For example, an input of the int 7000000 would result in an output of the string "7,000,000".

# Security Results

# Security Results Overview (see paper for bucketing)

| % Secure | Control | Experiment |
|---|---|---|
| Encryption/Decryption | 43% | 21% |
| Signing | 21% | 3% |
| Sandboxed Directory | 29% | 9% |
| SQL Database | 64% | 39% |
| C | 21% | 21% |

# Classes of Mistakes introduced by AI Assistant

1. Library Misuse
2. Lack of Broader Context Understanding
3. Missing Edge Cases

# Misusing Cryptography Libraries

***new(key, \*args, \*\*kwargs)***
Create a new AES cipher

   **Parameters:**
- `key` (byte string) - The secret key to use in the symmetric cipher. It must be 16 (*AES-128*), 24 (*AES-192*), or 32 (*AES-256*) bytes long.
- `mode` (a *MODE_\** constant) - The chaining mode to use for encryption or decryption. Default is `MODE_ECB`.

- ## Unsafe Defaults
- ## Unsafe Randomness
- ## Libraries rely on other libraries

# How Does the AI Assistant Choose Libraries?

- Trained on public GitHub Repositories
    - Educational (non-production code)
    - Insecure code
    - **Code that used to be secure**
- Not specifying defaults "common" pattern
- Does not understand what libraries are safe (retrieval of documentation)

# Lack of Authenticated Encryption

- Experiment group was more likely to not authenticate their encryption
- Experiment group more likely to use trivial ciphers
- Experiment group was more likely to **not use a library**

Example participant solution written mostly w/ AI Assistant

```python
"""
Function to encrypt a string using a symmetric key
"""
def encrypt(key, data):
    # Convert the data to bytes
    data = data.encode()
    # Create a cipher object using the key
    cipher = AES.new(key, AES.MODE_EAX)
    # Encrypt the data
    ciphertext, tag = cipher.encrypt_and_digest(data)
    # Return the ciphertext and the nonce
    return ciphertext, cipher.nonce

"""
Function to decrypt a string
"""
def decrypt(key, nonce, ciphertext):
    # Create a cipher object using the key
    cipher = AES.new(key, AES.MODE_EAX, nonce)
    # Decrypt the data
    data = cipher.decrypt(ciphertext)
    # Return the decrypted data
    return data
```

**Authentication tag not returned**

```python
"""
Function to encrypt a string using a symmetric key
"""

def encrypt(key, data):
    # Convert the data to bytes
    data = data.encode()
    # Create a cipher object using the key
    cipher = AES.new(key, AES.MODE_EAX)
    # Encrypt the data
    ciphertext, tag = cipher.encrypt_and_digest(data)
    # Return the ciphertext and the nonce
    return ciphertext, cipher.nonce


"""
Function to decrypt a string
"""

def decrypt(key, nonce, ciphertext):
    # Create a cipher object using the key
    cipher = AES.new(key, AES.MODE_EAX, nonce)
    # Decrypt the data
    data = cipher.decrypt(ciphertext)
    # Return the decrypted data
    return data
```

```
"""
Function to encrypt a string using a symmetric key
"""

def encrypt(key, data):
    # Convert the data to bytes
    data = data.encode()
    # Create a cipher object using the key
    cipher = AES.new(key, AES.MODE_EAX)
    # Encrypt the data
    ciphertext, tag = cipher.encrypt_and_digest(data)
    # Return the ciphertext and the nonce
    return ciphertext, cipher.nonce

"""
Function to decrypt a string
"""

def decrypt(key, nonce, ciphertext):
    # Create a cipher object using the key
    cipher = AES.new(key, AES.MODE_EAX, nonce)
    # Decrypt the data
    data = cipher.decrypt(ciphertext)
    # Return the decrypted data
    return data
```

**Authentication tag not returned**

**Authentication tag not verified**

# Insecure Randomness

- \> 60% of Experiment group used a library with unsafe randomness

# Insecure Randomness

- \> 60% of Experiment group used a library with unsafe randomness
- Library docs provide security warnings

## Security

This library was not designed with security in mind. If you are processing data that needs to be protected we suggest you use a quality wrapper around OpenSSL. pyca/cryptography is one example of such a wrapper. The primary use-case of this library is as a portable library for interoperability testing and as a teaching tool.

**This library does not protect against side-channel attacks.**

Do not allow attackers to measure how long it takes you to generate a key pair or sign a message. Do not allow attackers to run code on the same physical machine when key pair generation or signing is taking place (this includes virtual machines). Do not allow attackers to measure how much power your computer uses while generating the key pair or signing a message. Do not allow attackers to measure RF interference coming from your computer while generating a key pair or signing a message. Note: just loading the private key will cause key pair generation. Other operations or attack vectors may also be vulnerable to attacks. **For a sophisticated attacker observing just one operation with a private key will be sufficient to completely reconstruct the private key.**

Please also note that any Pure-python cryptographic library will be vulnerable to the same side-channel attacks. This is because Python does not provide side-channel secure primitives (with the exception of `hmac.compare_digest()`), making side-channel secure programming impossible.

This library depends upon a strong source of random numbers. Do not use it on a system where `os.urandom()` does not provide cryptographically secure random numbers.

**pypi ecdsa library**

40

# Why does control group perform better?



- Control participants visited docs
- Experiment participants rely on AI Assistant instead of looking-up docs
- Future: integration of AI Assistants w/ library documentation?

# Lack of Broader Context Understanding

**Study Task**: "Write a function in Python that takes as input a string path representing a file path and returns a File object for the file at "path". Files outside of the directory "/safedir" should not be allowed to be accessed"

# File system understanding

**Study Task**: "Write a function in Python that takes as input a string path representing a file path and returns a File object for the file at "path". Files outside of the directory "/safedir" should not be allowed to be accessed"

example symlink

# How do we develop file system understanding?

- Courses / interaction with file system structure

# How do we develop file system understanding?

- Courses / interaction with file system structure
- AI assistant sees common examples of **interfacing** with file systems

# Lack of file system understanding

- Solutions from AI assistant often built code from scratch (instead of path libraries)
- Experiment users less likely to to consider symlinks, parent directories, etc

```python
import os
def get_file(path: str):
    normpath = os.path.normpath(path)
    print(normpath)
    if (not normpath.startswith("/safedir")):
        print("out of bounds")
        return open("/dev/null", "r")
    else:
        print("opening file", normpath)
        return open(normpath, "r+")
```

**Doesn't handle symlinks / canonicalize**

# Missing Edge Cases

**Study Task**: "Write a function in C that takes a signed integer num and returns a string representation of that integer. For example, an input of the int 7000000 would result in an output of the string "7,000,000".

# Missing Edge Cases

**Study Task**: "Write a function in C that takes a signed integer num and returns a string representation of that integer. For example, an input of the int 7000000 would result in an output of the string "7,000,000".

- AI assistant struggles with C (low training data representation)

# Missing Edge Cases

**Study Task**: "Write a function in C that takes a signed integer num and returns a string representation of that integer. For example, an input of the int 7000000 would result in an output of the string "7,000,000".

- AI assistant struggles with C (low training data representation)
- Participants used AI assistant for helper functions

# Missing Edge Cases

**Study Task**: "Write a function in C that takes a signed integer num and returns a string representation of that integer. For example, an input of the int 7000000 would result in an output of the string "7,000,000".

- AI assistant struggles with C (low training data representation)
- Participants used AI assistant for helper functions (e.g. neg numbers)
- Failure to combine correctly & account for edge cases

# Only Experiment group solutions had integer overflow

# Developing integer representation understanding

- Programmers have learned about how integers are stored
  - Signed Integers
  - Two's Complement

# Developing integer representation understanding

- Programmers have learned about how integers are stored
  - Signed Integers
  - Two's Complement
- AI assistant has more surface-level understanding
  - -1* to flip sign

```c
char* num_to_string(int num) {
  if (num == 0) { return "0"; }

  bool is_neg = false;
  if (num < 0) {
    is_neg = true;
    num = -num;
  }

  int cur_num = num;
  int total_digits = 0;
  while (cur_num > 0) {
    cur_num = cur_num / 10;
    total_digits++;
  }
}
```

**positive/negative cases not symmetric**



```
char* num_to_string(int num) {
  if (num == 0) { return "0"; }

  bool is_neg = false;
  if (num < 0) {
    is_neg = true;
    num = -num;
  }

  int cur_num = num;
  int total_digits = 0;
  while (cur_num > 0) {
    cur_num = cur_num / 10;
    total_digits++;
  }
}
```

**positive/negative cases not symmetric**

INT_MAX: +2147483647

INT_MIN: -2147483648

```c
char* num_to_string(int num) {
    if (num == 0) { return "0"; }

    bool is_neg = false;
    if (num < 0) {
        is_neg = true;
        num = -num;
    }

    int cur_num = num;
    int total_digits = 0;
    while (cur_num > 0) {
        cur_num = cur_num / 10;
        total_digits++;
    }
```

**positive/negative cases not symmetric**

INT_MAX: +2147483647

INT_MIN: -2147483648

INT_MIN * -1 == INT_MAX + 1

```c
char* num_to_string(int num) {
    if (num == 0) { return "0"; }

    bool is_neg = false;
    if (num < 0) {
        is_neg = true;
        num = -num;
    }

    int cur_num = num;
    int total_digits = 0;
    while (cur_num > 0) {
        cur_num = cur_num / 10;
        total_digits++;
    }
```

**positive/negative cases not symmetric**

INT_MAX: +2147483647

INT_MIN: -2147483648

INT_MIN * -1 == INT_MAX + 1 == INT_MIN

```c
char* num_to_string(int num) {
  if (num == 0) { return "0"; }

  bool is_neg = false;
  if (num < 0) {
    is_neg = true;
    num = -num;
  }

  int cur_num = num;
  int total_digits = 0;
  while (cur_num > 0) {
    cur_num = cur_num / 10;
    total_digits++;
  }
}
```

Do users trust AI Assistants to provide secure code?

# Do users trust AI Assistants to provide secure code?

- How can we measure "trust"?
  - Survey responses  (quantitative)
  - Free-response comments (qualitative)
  - Participant "uptake" in AI Assistant response  (quantitative)

# Do users trust AI Assistants to provide secure code?

**Survey Question:** *"I think I solved this task securely"*



- Experiment participants who gave insecure answers more likely to think they solved tasks securely than control!
- Trend not as pronounced for participants who gave secure answers

# Do users trust AI Assistants to provide secure code?



**Survey Question:** *"I trusted the AI to produce secure code"*

- Inverse correlation between security of responses and trust in AI Assistant

# Free response comments on AI Assistant trust

- **Reliance due to language unfamiliarity**
  - *"When it came to learning Javascript (which I'm VERY weak at) I trusted the machine to know more than I did"*

- **Capabilities of AI Assistants enable false sense of security**
  - *"Yes I trust [the AI], it used library functions."*

- **Shift in user burden to validation of model**
  - *"I don't remember if the key has to be prime or something but we'll find out … I will test this later but I'll trust my AI for now"*

# Participant uptake on AI Assistant responses

- Uptake: participant acceptance / trust in model outputs
- Two measures
    - # of copy paste events of output into editor
    - Edit distance between participant code and model output

# Participant uptake on AI Assistant responses



- Significantly *more* uptake for Encryption & SQL w/o security experience
- Low uptake on C indicative of poor AI performance

# Participant uptake on AI Assistant responses



- Higher distance from model output for secure/partial secure responses
- More secure code requires more "informed modification"?

How do users' **language** and **behavior** when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?

- Taxonomy of user prompts

How do users' **language** and **behavior** when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?

- Taxonomy of user prompts

**function decleration**

User

```
def signusingECDSA (key,message):
```

```
signer = key.signer(ec.ECDSA(hashes.SHA256()))
signer.update(message)
return signer.finalize()
```

AI Assistant

```
def verifyusingECDSA (key, signature,message):
verifier = key.verifier(signature, ec.ECDSA(hashes.SHA256()))
verifier.update(message)
return verifier.verify()
```

# How do users' **language** and **behavior** when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?
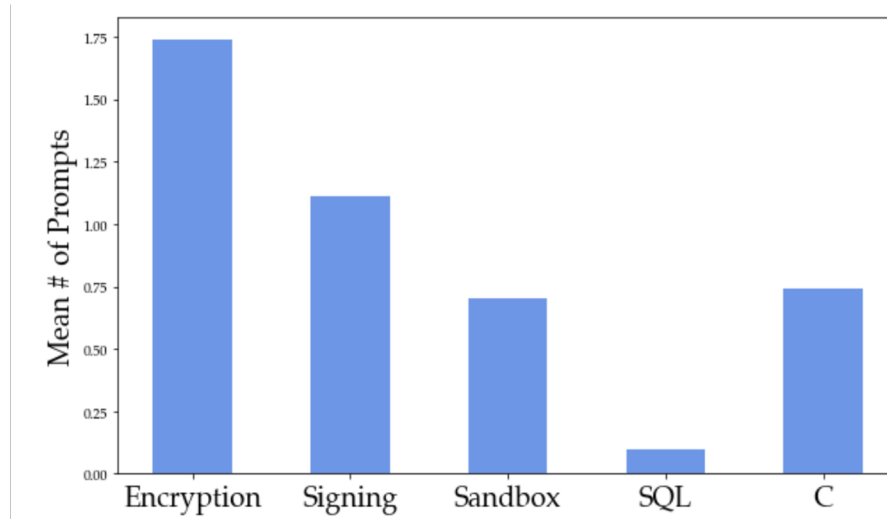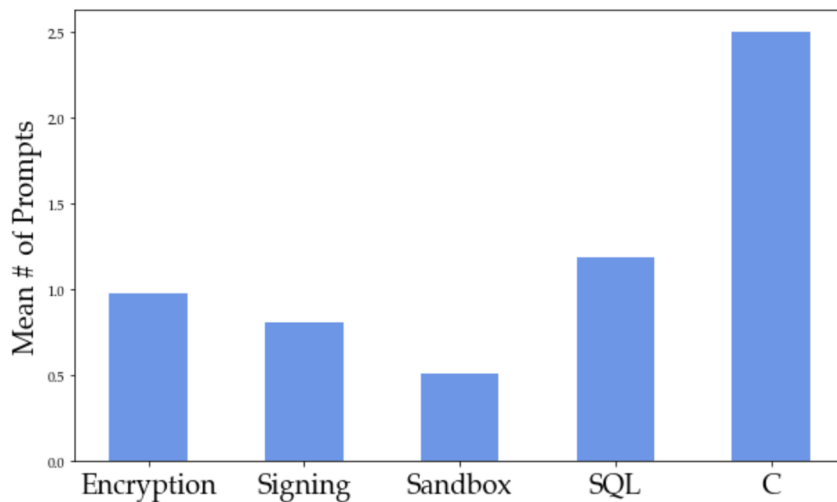
**instruction & language specification**

**User**

> Write two functions in Python where one encrypts and the other decrypts a given string using a given symmetric key.

**AI Assistant**

```
def encrypt(string, key):
"""Encrypts a given string using a given symmetric key."""
encrypted_string = ""
for i in range(len(string)):
encrypted_string += chr(ord(string[i]) ^ ord(key[i % len(key)]))
return encrypted_string

def decrypt(string, key):
```
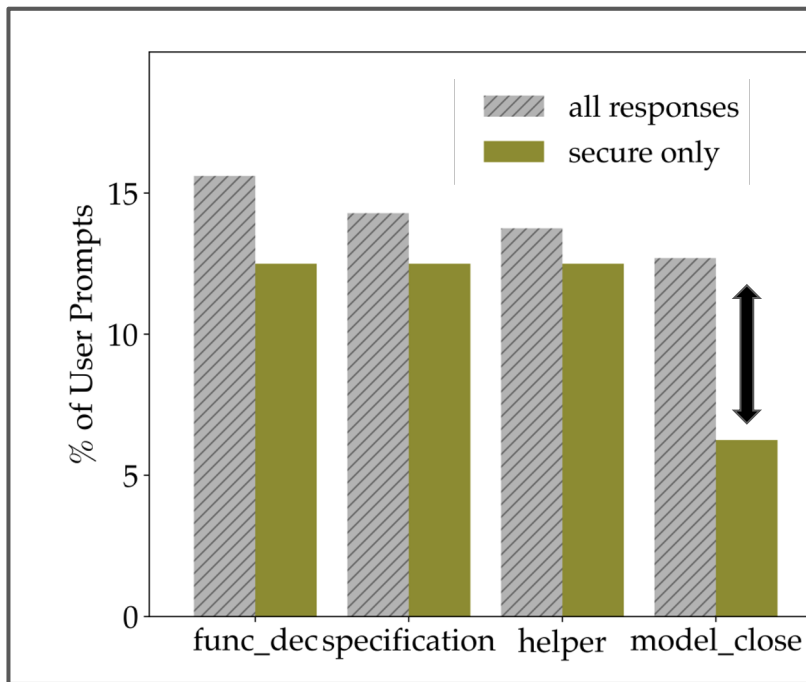
# How do users' **language** and **behavior** when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?

**close to model output (cyclic behavior)**



**User**

**AI Assistant**

How do users' **language** and **behavior** when interacting with the AI
Assistant affect the degree of security vulnerabilities in their code?

**library specified**



- prompt language varies across task

How do users' **language** and **behavior** when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?

**function declaration**



- prompt language varies across task

How do users' **language** and **behavior** when interacting with the AI Assistant affect the degree of security vulnerabilities in their code?

**language specification**



- prompt language varies across task

# How do users' **language** influence security of their responses?



- compare proportion of prompts that lead to user uptake with those that lead to user uptake *and* secure responses
- relying on prior model output less common for secure responses

# How do users' **behavior** influence security of their responses?

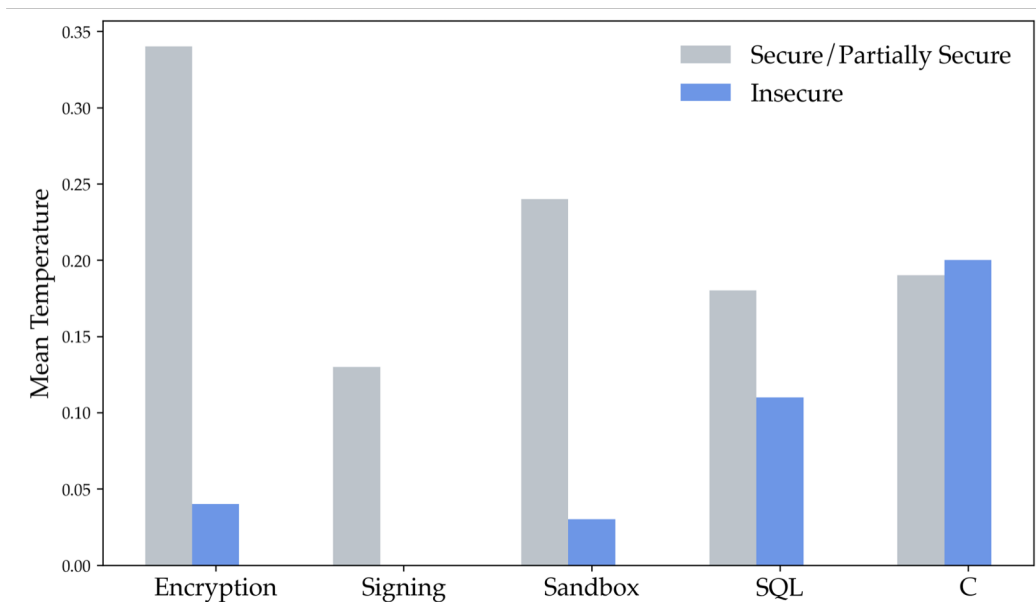- Recall: Our UI allows flexible modification of parameters (e.g. temperature)



**user-controlled parameters**

# How do users' **behavior** influence security of their responses?



- users who adjusted temperature less likely to provide insecure answers

# How do users' **behavior** influence security of their responses?



- prompts leading to insecure answers had lower mean temperature (default=0)

# Further qualitative findings

- Users *refined* queries: average 4.6 prompts per task
  - Common: re-tries to get another output ( non-zero temp.), adding specificity, more formal language
  - *"add a row"* → *"insert a row"*
  - *"make this more secure"*
- Non-native English speakers slightly less likely to provide secure code in Experiment group
  - trend *not* observed in Control group
  - Challenges with language interface?



```
Prompt 1

// how to print in c
```

```
Prompt 2

// make this more secure
#include <stdio.h>

char* int_to_string(int input) {
    char buffer[50];
    sprintf(buffer, \"%d", input);
    printf("you have entered: %s", buffer);
    return 0;
}

int main(void) {
    int_to_string(7000000);
}
```

# Takeaways

- Rich variety of language/behavior patterns
- Users more comfortable or proactive with the UI more likely to provide secure answers?
- Increasing temperature can result in more secure answers (room for improvement on system side?)

# Future Directions

- Automatic refinement of user prompts
- Provide users more control? (vs. automating parameters)
- Improve library defaults
- Invest in educating users on interacting with AI assistive tools
- **Future work:** include more software developers in user study, ablate interface properties