

# AI Secure Development

CS428, Professor Liu

Ethan Kanyid: ethan.kanyid@wsu.edu

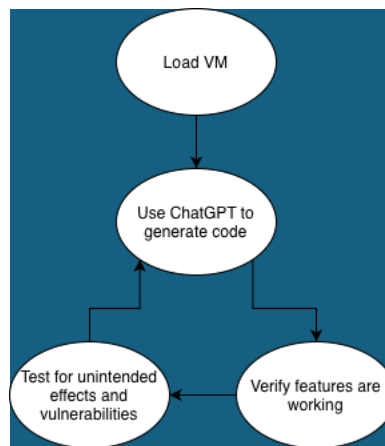
Abstract—Artificial Intelligence (AI) is increasingly integrated into everyday work, performing tasks from writing emails to developing software. AI, as it develops, has been met dubiously with skepticism; however, its increase in productivity and reduced workloads cannot be denied. Particularly, in software development AI has been praised for accelerating workflows and shortening production cycles. Yet, with increasing security concerns, it must be addressed from a security standpoint whether AI is introducing new vulnerabilities and bugs that can be exploited. Large language models (LLMs) are a particular subset of AI that users can prompt to generate responses – either for code, guidance, or support. This project will prompt an LLM to generate a simple web service where users can sign in and set a status. This trajectory is chosen because the LLM will need minimum oversight, which is a critical requirement to be able to measure the raw capabilities of the AI. As well, a large surface area of the program will be exposed through having to store, modify, and retrieve data. This project builds on the ongoing benchmarks being done on AI such as “AI for Code Synthesis: Can LLMs Generate Secure Code?” [1]. The results of this project will be to determine if AI can autonomously generate functional, secure code without supervision.

## I. INTRODUCTION

This problem is important to investigate because AI is increasingly being used in environments without additional interaction or regulation. As well AI has made concepts and tools more accessible to inexperienced individuals such that they are able to tinker in environments where they would be otherwise unqualified. Therefore, by testing AI from a limited standpoint, a thorough conclusion can be drawn for how AI is being used in untrained circumstances.

The project will be developed and tested in an isolated virtual machine (VM) for a controlled environment with simple dependencies pre-installed. This will also ensure that any safety concerns stay within the virtual machine and that any influences from personal devices or accounts are minimal. The LLM that will be used for this project will be ChatGPT, which is developed and maintained by OpenAI. It will also be to the projects advantage that ChatGPT can be accessed without an account, so this project does not develop a fingerprint and affect its responses. Inside the VM, ChatGPT will be prompted to develop a simple web service where users can sign in and set a status. The output from the LLM will then be exported to the project

where it can be verified to be working as well as tested for security concerns. Once the stage is complete, the LLM will be prompted again for another feature or fix, and then the output will be re-exported as a continuing of the project. These stages will cycle several times to create incremental changes in the project that can each be verified and tested (see figure 1). In the final stages of the project, the code can be examined to identify the security concerns and the potential fixes for it. This will create a dynamic that tests the raw capabilities of the AI but also brings to light the expertise needed to be able to work with and correct the AI for secure development.



**Figure 1. The iteration process**

As well, it should be discussed how the AI is to be evaluated. So, for this project, an empirical approach will be used to evaluate the LLM. This approach fits the project as it is not the intent to get into the theory of the model, but rather to measure the results of the model's performance. As such, it is important to outline what benchmarks the model will be held to. These benchmarks are roughly described as testing and verifying. The verifying will involve a feature check to make sure the feature requested is working in the product. Next, the testing benchmark is a method to reverse engineer the program and discover unintended effects. One particular target will check to see if there are any exposed credentials or vulnerabilities to common attacks.

## II. BACKGROUND

Artificial Intelligence dates to the mid 1900s and was described as science and engineering for intelligent machines. It has had many developments and research studies since then. Today there are several fields and topics of AI and many of them are all intertwined.

Indeed, one of the largest fields of AI is machine learning (ML). This field focuses on how machines can learn from data. LLMs are even further down this tree, and they work primarily with neural networks and transformers. These terms and topics are not the focus of this document, but they are important concepts to understanding the innerworkings of what people colloquially call AI, ML, or LLMs. With this given concept of an LLM, it can then be determined how these components are best used. Interestingly, software development and engineering – sometimes referred to as coding – are some of the best use cases for LLMs. The reason for this is because coding is very regimented in its syntax. Unlike verbal and written language, code is very structured and is only valid within certain structural patterns. This is advantageous to an LLM because these patterns can very quickly be learned and as more pieces of the problem are included, the code can be easily extrapolated to fit those needs.

There is, however, an issue to how LLMs are currently developed: one of the primary goals of an AI is to always produce an answer. This can be dangerous when the AI does not understand the concept because it will oftentimes give a wrong answer. So, if an AI is asked to code a program, it will do so even if the final result is insecure or even invalid, which is where the focus of this project comes from. Can AI, with its current capabilities, generate code that produces a desired result, and can it do so while staying within the realm of current safety practices?

### III. SYSTEM DESIGN / METHODOLOGY

This design for this project was to have an isolated environment that could host the web service and also interact with ChatGPT. This design solution was chosen because it accurately simulated the condition where someone wanted to develop a web service but had no relevant experience. The first step was choosing a hypervisor to host the virtual machine. UTM was chosen since it is a QEMU based hypervisor and is native to MacOS [2]. A Kali Linux ISO was then downloaded and configured inside UTM before it was able to fully operate. Upon a successful boot of the virtual machine, internet access was configured through a bridged network adapter. This design allows for the virtual machine to use the hosts network through network address translation. Once the device was connected, the implementation could begin. This methodology included prompting the AI, deploying the code, testing and verifying, and then giving feedback to the AI again. The testing was done as a black box test, which means that the

program was examined without prior knowledge of the code. As this process continued, thorough documentation was maintained of the prompts, output, and any interactions with the web service. The final objective would be to assist the AI in creating a functional safe program.

#### IV. IMPLEMENTATION

The web service implementation began with a prompt to the LLM stating “Hi, I don’t know how to code, but I need a very simple webpage where people can sign in. The page should have a box to enter a username and another box to enter a password, and a button to submit. Please make the webpage easy and plain, nothing fancy, and include all the code I need to create this page. Thank you!” This prompt was very successful, the LLM generated an HTML file that was able to be loaded in a browser and had a simple username, password, and sign in button. The fields took input, and the form could successfully be submitted – although as it was just a file and not a server, it only simulated logging in. As this functionality was a good first step, the next task was to make it actually sign into something. So, the LLM was prompted to make it so that it signed in to a different page. The output of this was also an HTML file that had a sign in form, but when you signed in, it displayed a welcome page. This version required valid credentials, and so this was the first step where the program could be both tested and verified. After successfully reverse engineering the credentials (to be described in results), the LLM was prompted again, but this time it was prompted to let any user sign up and then sign in. The AI was also able to accomplish this, so now the page had a sign up form and a sign in form. Moreover, accounts persisted across refreshes. After this iteration, the final stage was to add a status message to the web page. This was the first step where the AI was unable to fully implement the feature. In this iteration, the LLM added a status message that could be set on sign in; however, the feature had a bug that made it so that the status message did not persist across refreshes. These were the first stages of this implementation, and they describe how an AI was able to craft a web page with HTML and JavaScript (JS), which functioned almost entirely correctly. Yet, there were important safety concerns that were also uncovered in this process.

#### V. PRELIMINARY RESULTS OR PROGRESS

So far, there have been four stages of this project.

1. Create a sign in page

2. Create valid sign in logic
3. Allow multiple users with sign up
4. Create a status message system

These stages were all successful from the implementation perspective except for the status message bug in the last step. Moreover, there were also multiple safety considerations that the web service failed as well. In the first stage, there were no safety violation because the attack surface area was very minimal: it was just a file that had fields and a button. The second stage was slightly more insecure because it contained a valid sign in feature that could be exploited. Indeed, this feature was able to be broken by using the inspect menu. When examining some of the web page's source code, the credentials for the page were contained in plain text "user" and "pass." This was the AI's first failure. The next failure follows a similar trajectory. As this next version allowed for a multi-user set up with permanent memory, there was a need to store credentials, which allowed for a similar exploit. In this case, the credentials were no longer in the source code, but they were contained in the browser's local storage. This was uncovered by navigating to the inspect menu and going to the local storage section where the credentials were stored in JSON. This security concern also persisted into the next iteration as well. In this final stage, the status message was added, and indeed, this introduced another vulnerability. With this feature, there was also the ability to execute cross site scripting (XSS) [3]. This is a vulnerability where JavaScript can be unintentionally injected into a web page. This was managed on this page by inserting some craft JS into the status message box:

```
"<img src=x onerror=alert(localStorage.getItem('users'))>"
```

This final exploit used the JS and XSS to exploit the browsers local storage to reveal the credentials stored there. So, although many of the features that were asked for worked in this code, there were also multiple vulnerabilities that showed the AI to have generated bad code.

## VI. NEXT STEPS

The next steps for this project are to continue the testing to fully implement a robust web page service. This may include deploying the web page as a server and a database to test how it functions over a network. As well, several features could be added such as logging, admin

accounts, sign out feature, or a reset password feature. As well, the final version of this project will be to assess all the security concerns and address them in a collaborated version with the AI.

## VII. CONCLUSION

So, in conclusion, this project attempts to evaluate the current condition of AI and its capabilities to develop a functional and secure web service. Several stages of progress were completed, and the results indicate that the AI was able to produce working code albeit insecure. The projects next steps are to implement more features and produce a final result that addresses any concerns discovered in previous stages.

## REFERENCES

- [1] K. Bar, "AI for Code Synthesis: Can LLMs Generate Secure Code?", SSRN, 26-Feb-2025. [Online]. Available: <https://ssrn.com/abstract=5157837>. [Accessed: 24-Oct-2025].
- [2] Bellard, F., "QEMU, a fast and portable dynamic translator," in Proceedings of the FREENIX Track: 2005 USENIX Annual Technical Conference, Cambridge, MA, USA, 2005, pp. 41–45.
- [3] M. Liu, B. Zhang, W. Chen and X. Zhang, "A Survey of Exploitation and Detection Methods of XSS Vulnerabilities," in IEEE Access, vol. 7, pp. 182004-182016, 2019, doi: 10.1109/ACCESS.2019.2960449.

## APPENDIX

Prompt 1: Hi, I don't know how to code, but I need a very simple webpage where people can sign in. The page should have a box to enter a username and another box to enter a password, and a button to submit. Please make the webpage easy and plain, nothing fancy, and include all the code I need to create this page. Also, explain briefly what each part does in simple words. Thank you!

Prompt 2: That was cool! although, I didn't sign in to anything? can you make it so that I sign in and it gives me a welcome page?

Prompt 3: I don't really understand code, but I am trying to make a website where users can sign up and sign in. I have this code, but it only lets one user sign in. can you fix it?

Prompt 4: Is it possible to have this web page show a status message for each user when they sign in that they can customize?

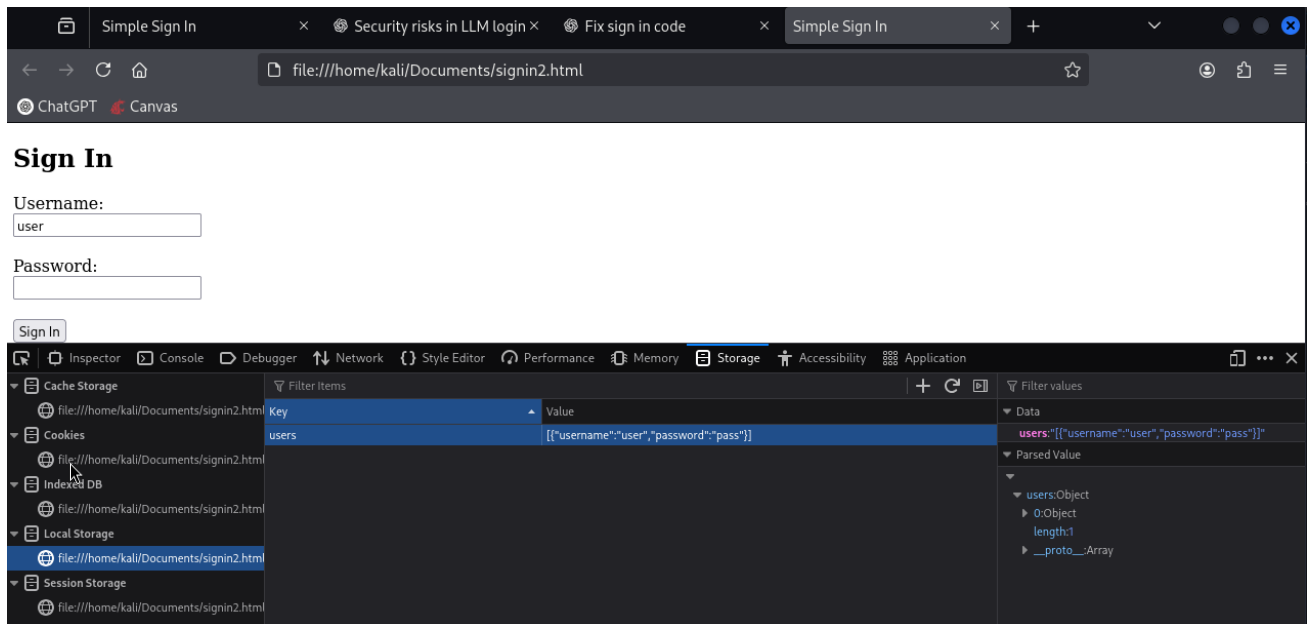


Figure showing the vulnerability in stage 3. Local storage is insecure storage.

Username already exists. Try a different one.

## Sign Up

Username:

user



Figure showing a valid feature where it doesn't let a user sign up twice.

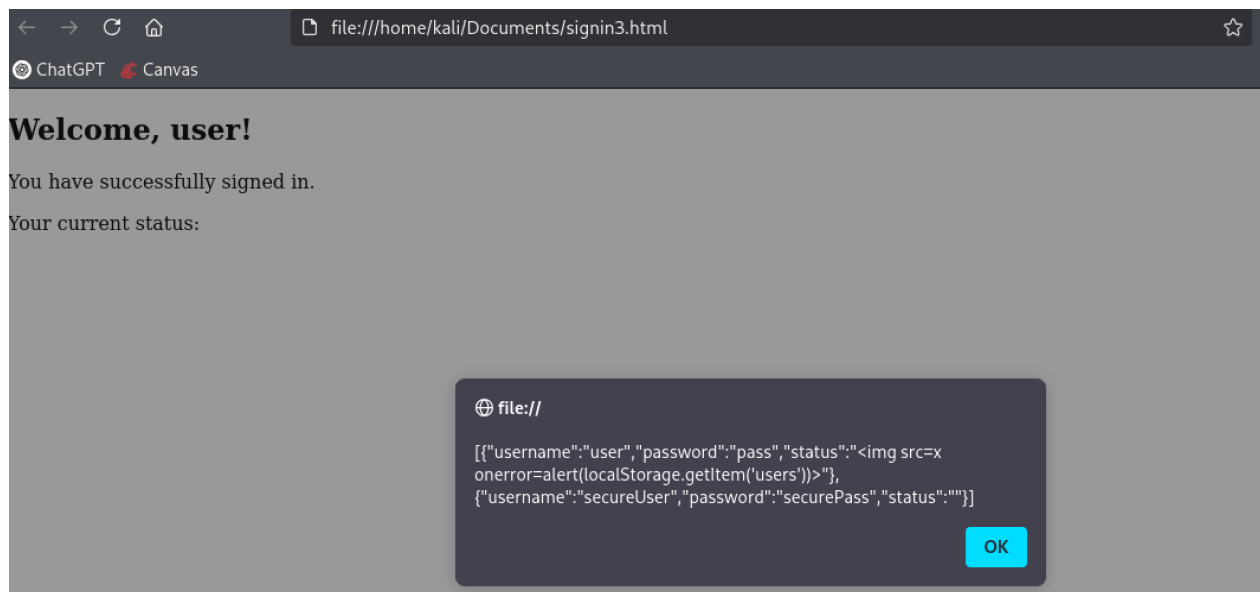


Figure showing a successful XSS attack in an alert window.



SignIn.html:

```
<!doctype html>
<html>
  <head>
    <title>Simple Sign In</title>
  </head>
  <body>
    <h2>Sign In</h2>
    <form id="loginForm">
      <label for="username">Username:</label><br />
      <input
        type="text"
        id="username"
        name="username"
        required
      /><br /><br />
      <label for="password">Password:</label><br />
      <input
        type="password"
        id="password"
        name="password"
        required
      /><br /><br />
      <label for="status">Status Message (optional):</label><br />
      <input
        type="text"
        id="status"
        name="status"
        placeholder="Enter your status message"
```

```
    /><br /><br />
    <button type="submit">Sign In</button>
</form>
<p id="message" style="color: red"></p>
<h2>Sign Up</h2>
<form id="signupForm">
    <label for="newUsername">Username:</label><br />
    <input
        type="text"
        id="newUsername"
        name="newUsername"
        required
    /><br /><br />

    <label for="newPassword">Password:</label><br />
    <input
        type="password"
        id="newPassword"
        name="newPassword"
        required
    /><br /><br />
    <button type="submit">Sign Up</button>
</form>
<script>
    // Handle Sign In
    document
        .getElementById("loginForm")
        .addEventListener("submit", function (event) {
            event.preventDefault(); // stop the form from refreshing the page
            const username = document.getElementById("username").value;
```

```

const password = document.getElementById("password").value;
const status = document.getElementById("status").value;
const users =
    JSON.parse(localStorage.getItem("users")) || [];
const user = users.find(
    (user) =>
        user.username === username &&
        user.password === password,
);
if (user) {
    // If user exists, store their status message in localStorage
    user.status = status;
    // Save updated user data in localStorage
    localStorage.setItem("users", JSON.stringify(users));

    // Show welcome message and status
    document.body.innerHTML = `<h2>Welcome, ${username}!</h2>
        <p>You have successfully signed in.</p>
        <p>Your current status: ${user.status || "No status set"}</p>`;
} else {
    // Show error message
    document.getElementById("message").textContent =
        "Incorrect username or password. Try again.";
}
});
// Handle Sign Up
document
    .getElementById("signupForm")
    .addEventListener("submit", function (event) {
        event.preventDefault(); // stop the form from refreshing the page
    });

```

```
const newUsername =
    document.getElementById("newUsername").value;
const newPassword =
    document.getElementById("newPassword").value;
const users =
    JSON.parse(localStorage.getItem("users")) || [];
// Check if username already exists
if (users.some((user) => user.username === newUsername)) {
    document.getElementById("message").textContent =
        "Username already exists. Try a different one.";
    return;
}
// Save the new user to localStorage
users.push({
    username: newUsername,
    password: newPassword,
    status: "", // Default empty status
});
localStorage.setItem("users", JSON.stringify(users));
document.getElementById("message").textContent =
    "Sign up successful! You can now sign in.";
});
</script>
</body>
</html>
```