# Serious Games for CyberSecurity
*Raising security awareness of industrial software developers*

Dr. Tiago Espinha Gasiba
tiago.gasiba@siemens.com
Siemens AG, Munich, Germany

SIEMENS  *Universität der Bundeswehr München*  istar_iscte Information Sciences, Technologies and Architecture Research Centre  fct Fundação para a Ciência e a Tecnologia  Bundesministerium für Bildung und Forschung  SIFO.de  contain

FKZ 13N16581
FKZ 13N16585

# Who am I?
*Background information...*



**Tech.** at Escola Profissional D. Sancho II, Elvas, Portugal in 1994-1997
Electrotechnical Professional
Main Interest: **Electronics**

**Eng.** at FEUP, Porto, Portugal in 1997-2002
Telecommunications, Electronics and Computers
Main Interest: **Digital Systems**

**M.Sc.** at TUM, München, Germany in 2002-2004
Communications Engineering
Main Interest: **Forward Error-Correcting Codes**

Researcher at Siemens Mobile in 2004-2007 ~~PhD~~
System Architect and Embedded Software Developer, Ericsson, in 2007-2013
System Design Engineer for JavaCard OS, NXP, Hamburg, in 2013-2014
**PhD.** at UniBwM, München, Germany in 2019-2021

*Senior Key Expert for Secure Software Development*
Siemens AG, München, 2014-present

Certifications:
CISSP, CISM, GXPN, SANS 642 Coin, ...

**SIEMENS**

*Fountain Codes*

Tiago Gasiba
@2014

CSC: Takes you from zero to hero!

# Overview of what is to come...
Agenda

**How to address software vulnerabilities during the development phase?**

| Tooling | Training | Reviews |
|---|---|---|
| • SAST<br>• DAST<br>• AI | • Traditional<br>• Serious games<br>• ... | • Human factor<br>• Tooling |

**Agenda**
1) What are serious games; what are the essential criteria for industrial serious games?
2) Examples of serious game design, developed through action design research.
3) (if we have time...) On the future of serious games using Artificial Intelligence

# Serious Games

What Constitutes a Serious Game? How are Serious Games different from Gamification?



**Serious Games are not the same as Gamification**

**Purpose**:

- **Serious Games**: Fully-fledged; medium to teach or achieve a purpose while engaging the player(s) in an experience
- **Gamification**: Adding game-like elements (e.g., points, badges, leaderboards) to a non-game event; it's more about enhancing existing processes

**Experience**:

- **Serious Games**: Immersive environment for learning or skill development. Players might not even realize they are being educated or trained while playing.
- **Gamification**: Non-game-related core activity(ies); engagement achieved through elements such as rewards

**Scope**:

- **Serious Games**: Full game with mechanics, storylines, challenges, and goals tied to the purpose
- **Gamification**: Partial gaming elements, without building a complete game

# Gamification
## Examples of Gamification



**Source:** https://www.straight.one/learning-base/mozart-stairs-in-stockholm-mehr-spass-am-treppensteigen/
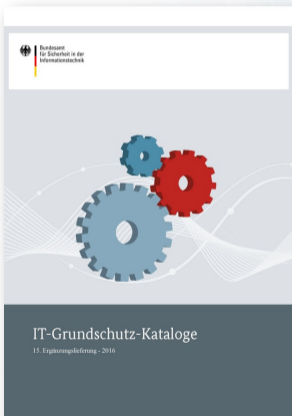
**Employee of the Month**

**"Spin the Wheel and Win"**

**GitHub Profile Badges**

## Serious Games

Recognized by the BSI IT-Grundschutz-Katalog (2016) as a mean to raise awareness

# Designing Industrial Serious Games
... using Action Design Research (Sein et al.)



**Industry Problem**

- First Cycle
  - Problem Formulation
  - Building, Intervention and Evaluation
  - Reflection and Learning
- Second Cycle
  - Problem Formulation
  - Building, Intervention and Evaluation + Empirical Study and Deep Dives
  - Reflection and Learning
  - Problem Formulation
  - Building, Intervention and Evaluation + Empirical Study
  - Reflection and Learning

**Formalization of Learning**

- Generalization of problem instance
- Generalization of solution instance
- Derivation of design principles

T. Espinha Gasiba, *Raising Awareness on Secure Coding in the Industry through CyberSecurity Challenges*,
Ph.D. dissertation, Fakultät für Informatik, Univ. Bundeswehr Munich, Germany, 2021.
**adapted from** Stein, M., Hendfridsson, O., Purao, S., Rossi, M., Lindfren, R., "*Action Design Research*", MIS Quarterly, 2011, 35(1), 37-56

# Three Industrial Serious Games
Design best—practices as outcome from research conducted in the industry



Secure Software Development

Secure Deployment in the Cloud

Secure Code Review

# Design Requirement Criteria for Serious Games in Industry
## Requirements and Background

On the Requirements for Serious Games geared towards Software Developers in the Industry

Tiago Gasiba, Kristian Beckers, Santiago Suppan, Filip Rezabek

*(author affiliation and email block, Siemens AG, Munich)*

**Abstract**—Teaching industry staff on cybersecurity issues is a fundamental activity that must be undertaken in order to guarantee the delivery of high-quality and robust products in market. Much research attention has been devoted to this topic over the last years. However, the research which has been done has not focused on developing secure code in industrial environments. In this paper we take a look at the constraints and requirements for delivering a training, by means of cybersecurity challenges, in this context which raises awareness on secure coding. In particular, using requirements engineering, we aim at understanding the design requirements for such challenges. Along the way, we give details on our experience of delivering cybersecurity challenges in an industrial setting and show the outcome and lessons learned. The proposed requirements for cybersecurity challenges geared towards software developers in an industrial environment are based on systematic literature review, interviews with security experts from the industry and semi-structured evaluation of participant feedback.

*(remaining abstract and introduction body text in small columns)*

| Requirement | Background |
|---|---|
| ❑ Clear learning goal | serious game |
| ❑ Adapted to job description | targeting industry |
| ❑ Working mechanics of game | serious game |
| ❑ Progressive difficulty | serious game mechanics |
| ❑ Elicit discussions among participants | targeting industry |
| ❑ Present compliant solutions | targeting industry |
| ❑ Include solution hints | serious game mechanics |
| ❑ Standardized solution | targeting industry |
| ❑ Planning of duration of event | targeting industry |
| ❑ Consider different technology stacks | general awareness |
| ❑ Defensive challenges | novel result / industry |
| ❑ Focus on secure coding guidelines | targeting industry |
| ❑ No need to learn new tools | targeting industry |
| ❑ Raise awareness of malicious attack | general awareness |

CTF (Capture-the-Flag) games, although normally targeting cybersecurity, are not ideal to raise awareness of software developers on secure coding

# Empirical study: comparison of CTF vs Game Requirements
## Results of Lightweight Literature Review

| Requirement | AutoCTF [163] | PicoCTF [165] | PlaidCTF [22] | Class CTF [99] | CSAW CTF [118] | KYPO Cyber Range [231] | % |
|---|---|---|---|---|---|---|---|
| 1. Have a clearly defined learning goal objective | Not fulfilled | Not fulfilled | Not fulfilled | Fulfilled | Fulfilled | Fulfilled | 50% |
| 2. Adapted to background (job description) of participants developers | Fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 17% |
| 3. Well defined working mechanics (e.g. which tools to use or what to do) | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 0% |
| 4. Defined and progressive level of difficulty challenges | Not fulfilled | Fulfilled | Fulfilled | Not fulfilled | Fulfilled | Not fulfilled | 50% |
| 5. Elicit discussions of the solutions (e.g. is there a better/simpler way to solve?) | Fulfilled | Fulfilled | Fulfilled | Fulfilled | Fulfilled | Fulfilled | 100% |
| 6. Provide possible solution after challenge solved | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 0% |
| 7. Adapted to the skill level of participants | Not fulfilled | Fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 17% |
| 8. Challenge includes hint that aid to arrive to the solution | Fulfilled | Fulfilled | Fulfilled | Fulfilled | Fulfilled | Fulfilled | 100% |
| 9. Clear, standardized and simple solution (not based on obscure knowledge) | Not fulfilled | Fulfilled | Not fulfilled | Not fulfilled | Fulfilled | Fulfilled | 50% |
| 10. Planned duration of the exercise | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 0% |
| 11. Explains issues arriving from interplay of different technologies or components | Not fulfilled | Fulfilled | Fulfilled | Fulfilled | Fulfilled | Fulfilled | 83% |
| 12. Adapted to company internal secure coding guidelines and policies | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 0% |
| 13. Challenges are put from the defensive perspective | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 0% |
| 14. Solutions does not require specific knowledge of hacking tools | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | Not fulfilled | 0% |
| 15. Challenges should raise awareness of possible consequences of malicious attack | Not fulfilled | Not fulfilled | Not fulfilled | Fulfilled | Not fulfilled | Not fulfilled | 16% |
| % | 26% | 40% | 26% | 40% | 33% | 33% | |

Our experience shows that these requirements are fundamental for **industrial serious games**

"Important" industry requirements are not addressed

# Lesson in failure: when interactive storytelling fails
How not to design serious games



When Interactive Graphic Storytelling Fails

James Banks*, Tiago Gasiba*, Santiago Reinhard Suppan*, Marc Berges*, Kristian Beckers*

- Lesson learned: not all challenges are suitable as CyberSecurity Challenges
- Challenge type depends on environment
- Capture-the-Flag creates a "competitive environment"
  - Players focused on winning rather than learning
- Root cause of failure
  - Not designed to address software developers
  - No secure coding guidelines, not adapted to job

# When Interactive Storytelling Fails

Five Scenarios

## Backup and Restore

## Password Management

## Patch Management

## Principle of Least Privilege

## Where to Share Data

# Learning Techniques and Activities for Serious Games
## Mapping between Techniques and Activities (Silva, 2020)

| Learning Techniques | Leaning Activities |
|---|---|
| Practice & feedback | Questions, memorization, association, imitation |
| Learning by doing | Interaction, practice, drill, imitation |
| Learning from mistake | Feedback, problem |
| Discovery learning | Feedback, problem, creative play |
| Task-based learning | Understand principle, graduated tasks |
| Question-led learning | Question, problem |
| Situated learning | Immersion |
| Role playing | Imitation, practice, coaching |
| Constructivist learning | Experimentation, questioning |
| Learning object | Logic, questioning |
| Coaching | Coaching, feedback, questioning |
| Intelligent tutors | Feedback, problem, continuous practice |

# Designing Serious Games

A framework by Silva, 2020



**Cost**
- **unproductive hours**
- **extra costs**
- **maintenance / adaptation**

**Benefit**
- **real-world examples**
- **teamwork**
- **standard compliance**
- increased awareness / security

Silva FGM. Practical Methodology for the Design of Educational Serious Games. Information. 2020; 11(1):14.

# Industrial Serious Games

Sifu – Raising Software Developers' Awareness of Secure Coding Guidelines



## CSC

### **C**yber**S**ecurity **C**hallenges

T. Espinha Gasiba, *Raising Awareness on Secure Coding in the Industry through CyberSecurity Challenges*, Ph.D. dissertation, Fakultät für Informatik, Univ. Bundeswehr Munich, Germany, 2021.

- You get a vulnerable code snippet
- Task is to find vulnerabilities and fix them
- Failed solutions uncover hints
- Supports: C, C++, Java, Python, Terraform

# Industrial Serious Games
*Suitable Challenge Types*



Design of Secure Coding Challenges for Cybersecurity Education in the Industry

- That kinds of challenge types are suitable for industrial players?

Traditional methods



**Fig. 4.** Single-Choice Question

**Fig. 5.** Multiple-Choice Question

Interactive methods

**Fig. 6.** Text-Entry Question

**Fig. 7.** Code-Snippet Question

**Fig. 8.** Code-Entry Challenge

**Fig. 9.** Associate Left-Right

# Structure of Serious Game Challenge

Challenge Structure in Three Phases

# Previous Work: Design through Action Design Research
Overview of Research Activities

| | | | |
|---|---|---|---|
| Offensive Challenges | Requirements | Client-Based Deployment | **1** First Design Cycle |
| Defensive/Offensive Challenges | Challenge Design | Server-Based Deployment | **2** Second Design Cycle |
| Ranking Secure Coding Guidelines | Temporal Dimensioning | Case Study | |
| Defensive Challenges | Industry Survey | Cloud-Based Deployment | **3** Third Design Cycle |

**Legend:**

| Empirical Studies | Challenges | Design Studies | Deployment | Deep Dive |
|---|---|---|---|---|

# Sifu Platform (Defensive Challenges)

Third Design Cycle: developed platform, released under MIT license



Player Interaction

Platform Buttons

Challenge

Project Files

Hints

Feedback

# CyberSecurity Challenges
## Embedded in a Corporate Training Event; Agenda and Empirical Studies

**1 2 3**

| Duration | What | Description |
|---|---|---|
| 10 min | *Welcome* | Welcome to participants and accessing CSC infrastructure |
| 20 min | *Team building* | Participants select partners and build teams that will play against each other |
| 30 min | *Introduction* | Challenge types are presented. One challenge in each category is solved in order to show the participants how the game works |
| 320 min | *Main event* | Game is open and teams are free to play the game. They are responsible for defining their own strategy for time-out (e.g. for lunch break). |
| 10 min | *Winner* | Game is closed and teams can no longer submit points to the dashboard. Winning team is announced. A brief review of the game-play is done together with the participants. |
| 30 min | *Feedback* | Participants are asked to fill out a survey about the game. Additionally, discussions with players is held in short non-systematic interviews. Main points of discussions is recorded for later analysis. |
| 60 min | *Walk-through* | Participants are shown solution to the exercises they considered most difficult. These exercises are solved together in interaction with all the participants. Discussion on how to solve the challenge is highly encouraged. |

**Survey:** opt-in with informed consent and anonymized answers

# Evaluation of Serious Game

Success in the industry

| Question | Industry | | | | | | Description |
|---|---|---|---|---|---|---|---|
| | D/O | | | D | | | |
| | - | N | + | - | N | + | |
| Q1 | 12.5 | 7.1 | 80.4 | 0.0 | 10.0 | 90.0 | I learned new techniques and principles of secure software development |
| Q2 | 0.0 | 5.3 | 94.7 | 0.0 | 0.0 | 100.0 | I understand the importance of secure coding guidelines |
| Q3 | 3.6 | 14.3 | 82.1 | 0.0 | 0.0 | 100.0 | Focusing on the challenges improves my practical secure coding skills |
| Q4 | 8.9 | 8.9 | 82.2 | 8.0 | 8.0 | 84.0 | The learning goals of the challenges were clearly explained |
| Q5 | 1.8 | 12.5 | 85.7 | 0.0 | 0.0 | 100.0 | The help from the coaches was adequate |
| Q6 | 8.9 | 26.8 | 64.3 | 0.0 | 20.0 | 80.0 | I feel that I am prepared to handle issues related to secure coding at work |

- **More than 80%** of evaluation criteria shows positive results
- Usefulness: established in the official company training curriculum
- Success achieved in industry and academia
- Note: *not always possible to achieve 100% agreement in every category*

# CyberSecurity Challenges
*How to play a challenge?*

# Automatic Evaluation of Challenges
*Architecture: Frontend/Backend*

# Industrial Serious Games
CATS – Raising Awareness on Responsibilities for Secure Cloud Deployments



## CATS
**Cloud of Assets and Threats**

Ongoing work by: Tiange Zhao

- Electronic (virtual) table-top game
- Goal is to build a defense strategy
- Failed solutions provide hints
- Teach separation of responsibilities and defense strategies against typical attacks

# Cloud of Assets and Threats
## *A Virtual Table-Top Game*



Different Responsibilities

Player cards

Attack description

Solution evaluation

Solution Rationale

Card description

Hints

Pool of cards

Game tips

# Cloud of Assets and Threats
*Game Mechanics*



Where did we see this model?...

# Cloud of Assets and Threats
*Game Evaluation in an Industrial Setting*

| Game Event | Date | Player | Team | CSC or Training | Valid Submissions |
|---|---|---|---|---|---|
| 1 | 2022-01-21 | 17 | 4 | CSC | 177 |
| 2 | 2022-03-15 | 14 | - | Training | 477 |
| 3 | 2022-03-22 | 14 | - | Training | 493 |
| 4 | 2022-03-29 | 13 | - | Training | 312 |
| 5 | 2022-04-14 | 13 | 4 | CSC | 178 |
| 6 | 2022-04-26 | 11 | - | Training | 100 |
| 7 | 2022-05-03 | 8 | - | Training | 171 |
| 8 | 2022-06-02 | 4 | 2 | CSC | 169 |
| 9 | 2022-09-29 | 14 | 4 | CSC | 298 |
| **Total number of players** | | | | 108 | |
| **Total number of submissions** | | | | 2375 | |

Interaction and engagement

**fun** is also important

Crucial positive aspect: **learning**

| ID | Feedback | Quote |
|---|---|---|
| 1 | Positive | *I really enjoyed the teamwork during the game* |
| 2 | Positive | *Playing the game was lots of fun, and I have learned many new things* |
| 3 | Positive | *The coaches did a good job in helping our team to solve the challenges* |
| 4 | Positive | *I learned a lot from the discussions we had during and after the game* |
| 5 | Positive | *I am eager to put to practice all the knowledge I have now learned* |
| 6 | Negative | *We had problems with the access to the infrastructure* |
| 7 | Negative | *The music played during the game was very distracting* |
| 8 | Negative | *I was not fully aware of what my team colleagues were doing* |
| 9 | Negative | *Communication with team members and with coaches was difficult* |

Watch out for
- **Communication**
- **Infrastructure**
- **Environment**

# Cloud of Assets and Threats
*Study on effects before/after playing the game*

# Industrial Serious Games
### DuckDebugger – Empowering Developers to Perform Code Review



# DuckDebugger

**D**iscover Vulnerabilities and **D**efend your code
*…because code reviews matter!*

Ongoing work by: Andrei-Cristian Iosif

- Player is given a vulnerable code snippet
- Task is to review code and identify vulnerability-hotspots
- Hints in the form of checklist and SAST tool output
- Supports: C, C++, Java, JavaScript, Python, Golang, C#, …

# DuckDebugger – Discover and Defend
Interactive Challenges

Quack!

Click the duck to submit your answer

User review

Vulnerable code

Solution at the end

**Comment #**        **Code**                        **Solution**

```
1  import sqlite3, random
2  from flask import Flask, abort, request, jsonify
```

**SAST Output**  x

tool_1  tool_2  ...

**Line#:** <finding description>

**Checklist**  x

☐ <Guiding Question 1> ?
☐ <Guiding Question 2>?
...          ...

```
9   def create_response(message):
10      response = jsonify({'message': message})
11      response.headers.add('Access-Control-Allow-Origin',
12      # TODO Ticket: id91263
13      return response
```

Output of SAST tools

Checklist to assist review

# DuckDebugger – Discover and Defend
Dimensions of Empowerment

# DuckDebugger – Discover and Defend
## Preliminary Results from Interventions in Industry

How players interact with the game

## Involvement

**I1. Initiation**
2 | 18 | 57 | 22

**I2. Delegation**
10 | 32 | 38 | 20

**I3. Participation**
12 | 60 | 28

**~70%**

## Motivation

How game motivates players

**M1. Reward**
2 | 10 | 30 | 35 | 22

**M2. Knowledge**
2 2 | 45 | 50

**M3. Information**
5 | 18 | 35 | 42

**M4. Power**
8 | 30 | 38 | 25

**~77%**

Legend: Strong Disagreement | Disagreement | Neutral | Agreement | Strong Agreement

**Results collected from 200+ participants, across 13 game events**

# DuckDebugger – Discover and Defend
## Preliminary Results from Interventions in Industry

**Teamwork**

Game improves teamwork

T1. Team strengths

| | 25 | 48 | 28 |

T2. Project knowledge

| 2 | 8 | 35 | 55 |

~89%

T3. Knowledge sharing

| 2 | 38 | 60 |

**Significance**

Game is significant for daily job

S1. Impact

| 2 | 15 | 50 | 32 |

S2. Self-Determination

| 8 | 20 | 45 | 28 |

~81%

S3. Competence

| 2 | 2 | 48 | 48 |

S4. Meaningfulness

| 2 | 25 | 50 | 22 |

■ Strong Disagreement  ■ Disagreement  ■ Neutral  ■ Agreement  ■ Strong Agreement

**Results collected from 200+ participants, across 13 game events**

# Large Language Models and Secure Software Development

# Large Language Models

# Large Language Models and Secure Software Development
## CodeGuardian – (*Work Currently Under Review*)

Large Language Models for Secure Code
Assessment: A Multi-Language Empirical Study

*Abstract*—Most vulnerability detection studies focus on datasets of vulnerabilities in C/C++ code, offering limited language diversity. Thus, the effectiveness of deep learning methods, including large language models (LLMs), in detecting software vulnerabilities beyond these languages is still largely unexplored. In this paper, we evaluate the effectiveness of LLMs in detecting and classifying Common Weakness Enumerations (CWE) using different prompt and role strategies. Our experimental study targets six state-of-the-art pre-trained LLMs (GPT-3.5-Turbo, GPT-4-Turbo, GPT-4o, CodeLlama-7B, CodeLlama-13B, and Gemini 1.5 Pro) and five programming languages: Python, C, C++, Java, and JavaScript. We compiled a multi-language vulnerability dataset from different sources, to ensure representativeness. Our results showed that GPT-4o achieves the highest vulnerability detection and CWE classification scores using a five-shot setting.

Aside from the quantitative results of our study, we developed a binary-input CODEGUARDIAN integrated with VSCode, which enables developers to perform LLM-assisted real-time vulnerability analysis in real-world security scenarios. We evaluated CODEGUARDIAN with a user study involving 22 developers from the industry. Our study showed that, by using CODEGUARDIAN, developers are more accurate and faster at detecting vulnerabilities.

## I. INTRODUCTION

A software vulnerability is a defect that could allow an attacker to gain control of a software system, steal or manipulate sensitive data, install a backdoor, or plant other types of malware [1, 2, 3]. Software vulnerabilities vary in nature and level of exploitability. Not all vulnerabilities allow attackers to cause the same types of harm, and not all vulnerabilities create equally severe risks. However, all vulnerabilities pose at least some level of risk to the applications they impact, as well as the environments that host these applications.

Software vulnerabilities may have different root causes, spanning from the way the software is designed, problems with the software's source code, poor management of data, or access control settings within the application. In this work, we focus on classes of vulnerabilities related to bugs in the source code. The U.S. Department of Homeland Security reports that 90% of security incidents are caused by exploiting security flaws and defects in the code [4]. A recent survey [5] found that over 55% of developers struggle to identify code vulnerabilities, indicating a lack of secure coding knowledge. Due to the increasing frequency and severity of cybersecurity incidents, which negatively impact end-users, businesses, and critical infrastructure [6], identifying security vulnerabilities within source code has become essential.

Researchers have proposed methods for the automatic detection and repair of software vulnerabilities [7, 8, 9, 10,

11, 12, 13]. For instance, program analysis-based techniques have been widely employed [14]. However, such rule-based techniques are affected by high false positive rates and low effectiveness in handling a large number of diverse types of vulnerabilities [12]. In addition, they often suffer from long response times [15], which makes them less suitable for continuous integration/development environments. Recently, Deep Learning (DL) techniques have enabled significant progress in vulnerability detection in the source code [15, 16, 17]. However, most DL models treat the code as linear sequences or textual data, which limits the finding of the semantic relationship within the source code and demands the development of more robust and efficient approaches.

To address these limitations, researchers have started exploiting the potential of Large Language Models (LLMs) that are pre-trained on a large corpus of data and use advanced natural language processing capabilities based on the transformer architecture [18]. LLMs have shown great performance across various natural language and software engineering tasks [11, 18], including automated vulnerability detection and repair [9, 10, 11]. Existing studies on vulnerability detection are executed against datasets that include vulnerabilities pertaining to only one or two programming languages. Instances are the vulnerability-fixing commit dataset by Zhou et al. [19] (C/C++), by Li et al. [20] (Java), BigVul [21] (C/C++), SVEN [22] (Python and C/C++), Devign [23] (C), D2A [24] (C/C++). Consequently, the capabilities of LLMs across diverse programming languages remain largely unexplored. Furthermore, these studies are purely quantitative, and the potential of LLMs to assist developers in identifying and classifying vulnerabilities in real-world use cases is yet to be investigated.

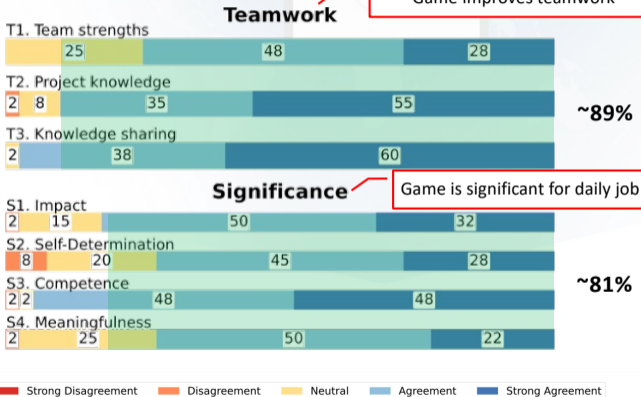To fill these research gaps, this paper investigates the effectiveness of LLMs in identifying vulnerable codes across different programming languages. To the best of our knowledge, this is the first study that reports an extensive evaluation of six state-of-the-art pre-trained LLMs (GPT-3.5-Turbo, GPT-4-Turbo, GPT-4o, CodeLlama-7B, CodeLlama-13B, and Gemini 1.5 Pro) in identifying vulnerabilities pertaining to five programming languages, namely Python, C, C++, Java, and JavaScript. Our study focuses on the top 25 Common Weakness Enumeration (CWE) classes that represent the most dangerous software weaknesses [25], for which we compiled a dataset of vulnerable snippets from representative domains. Our study targets vulnerability detection and CWE classification using zero-shot and few-shot approaches.

Furthermore, to investigate the industrial applicability of LLMs in software vulnerability detection and debugging, we

- Using LLM models to assist software developers to write secure code
- Ease of use: seamless integration in VisualStudio Code as a plugin

# Large Language Models and Secure Software Development

How well does ChatGPT perform (i.e. can detect vulnerabilities)?

About **one in every five** vulnerabilities is not found!

About **one in every ten** vulnerabilities is wrongly labeled

Prompt

| Prompt | Vul. Detection | | | | CWE | | | |
|---|---|---|---|---|---|---|---|---|
| | Acc. | Prec. | Rec. | F1 | Acc. | Prec. | Rec. | F1 |
| S1-V + V1 + S1-C + CWE-c | 0.81 | 0.80 | 0.82 | 0.81 | 0.86 | 0.69 | 0.66 | 0.66 |
| S2-V + V2 + S2-C + CWE-c | 0.77 | 0.71 | 0.94 | 0.80 | 0.85 | 0.76 | 0.71 | 0.72 |
| S1-V + V2 + S1-C + CWE-c | 0.80 | 0.75 | 0.90 | 0.82 | 0.85 | 0.80 | 0.73 | 0.75 |
| S2-V + V1 + S2-C + CWE-c | 0.79 | 0.76 | 0.87 | 0.81 | 0.86 | 0.84 | 0.76 | 0.77 |

Naïve Usage

About **four in every ten** vulnerabilities is are found!

| LLMs | Prompt | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|---|
| GPT-4 Turbo | S1-C + CWE-z | 0.56 | 0.21 | 0.18 | 0.18 |

# A Simple ChatGPT Experiment
## CWE 208 - Observable Timing Discrepancy

```c
int is_equal(const char* a, const char* b, size_t len) {
    for (size_t i = 0; i < len ; i++) {
        if (a[i] != b[i])
            return 1;
    }
    return 0;
}
```



Code can be trivially broken!

# A Simple ChatGPT Experiment
CWE 208 - Observable Timing Discrepancy – ChatGPT's solution

```c
#include <openssl/crypto.h>

int is_equal(const char* a, const char* b, size_t len) {
    if(a == NULL || b == NULL) return -1;
    size_t a_len = strnlen(a, len);
    size_t b_len = strnlen(b, len);
    if (len == 0 || len > a_len || len > b_len) {
        len = a_len < b_len ? a_len : b_len;
    }
    if (a[a_len] != '\0' || b[b_len] != '\0') {
        return -1;
    }
    return CRYPTO_memcmp(a, b, len);
}
```

Technically correct, but unnecessary

Technically correct, but using extra libraries

- ☐ **Using ChatGPT, solution had to be forced** to obtain a solution
  - ➢ Explicitly state that there is a timing vulnerability in the code
- ☐ AI's solution introduces **dependency on 3ʳᵈ party component**
  - ➢ In an industrial context, this can be undesired, e.g. due to **licensing issues**

# A Simple ChatGPT Experiment
CWE 208 - Observable Timing Discrepancy – ChatGPT's 2nd solution

```c
int is_equal(const char* a, const char* b, size_t len) {
    if(a == NULL || b == NULL) return -1;
    size_t a_len = strnlen(a, len);
    size_t b_len = strnlen(b, len);
    if (len == 0 || len > a_len || len > b_len) {
        len = a_len < b_len ? a_len : b_len;
    }
    if (a[a_len] != '\0' || b[b_len] != '\0') {
        return -1;
    }
    int result = 0;
    for(size_t i = 0; i < len; i++) {
        result |= a[i] ^ b[i];
    }
    return result;
}
```

Technically correct, but unnecessary due to context

Expected solution

❑ Additional solution had to be **forced** to *remove dependency on 3rd party library*
❑ Solution still exhibits several problems (after more than **15 interactions**)
   ➢ *Potentially unnecessary checks*
   ➢ *Maintenance of code can be problematic*

# Observed problems with AI and Secure Coding
Short summary of identified issues

> pre-conditions
> superfluous corrections

> Different behaviour
> no-go for safety-critical systems

> Efficiency not guaranteed
> affecting real-time systems

Missing context

Changed semantics

Code complexity

Learning inference (abuse/leakage)

Limited learning

Maintainability

> Leakage of IPR, company internal secrets

> Findings limited to learning data set

> e.g.missing comments

# Large Language Models and Secure Software Development

... we are sure that this will become more and more ubiquos



**Q1** Not helpful at all | Slightly helpful | Moderately helpful | Very helpful | Strongly helpful

**(Q1) How helpful did you find the extension for identifying vulnerabilities in source code?**

| 36% | 64% |

**Q2** Very difficult | Difficult | Moderate | Easy | Very easy

**(Q2) How easy was the extension to use for scanning vulnerabilities?**

| 9% | 27% | 64% |

**Q3** Not intuitive at all | Slightly intuitive | Moderately intuitive | Very intuitive | Extremely intuitive

**(Q3) How intuitive was the interface of the VSCode extension?**

| 9% | 18% | 27% | 45% |

**Q4** Very unsatisfied | Slightly unsatisfied | Neutral | Satisfied | Very satisfied

**(Q4) How satisfied are you with the response speed of the vulnerability scanning process?**

| 9% | 18% | 27% | 45% |

# Large Language Models and Secure Software Development
A recent study indicates potential serious consequences…

**RESEARCH ARTICLE SUMMARY**

ARTIFICIAL INTELLIGENCE

## Durably reducing conspiracy beliefs through dialogues with AI

Thomas H. Costello*, Gordon Pennycook, David G. Rand

**INTRODUCTION:** Widespread belief in unsubstantiated conspiracy theories is a major source of public concern and a focus of scholarly research. Despite often being quite implausible, many such conspiracies are widely believed. Prominent psychological theories propose that many people want to adopt conspiracy theories (to satisfy underlying psychic "needs" or motivations), and thus, believers cannot be convinced to abandon these unfounded and implausible beliefs using facts and counterevidence. Here, we question this conventional wisdom and ask whether it may be possible to talk people out of the conspiratorial "rabbit hole" with sufficiently compelling evidence.

**RATIONALE:** We hypothesized that interventions based on factual, corrective information may seem ineffective simply because they lack sufficient depth and personalization. To test this hypothesis, we leveraged advancements in large language models (LLMs), a form of artificial intelligence (AI) that has access to vast amounts of information and the ability to generate bespoke arguments. LLMs can thereby directly refute particular evidence each

individual cites as supporting their conspiratorial beliefs.

To do so, we developed a pipeline for conducting behavioral science research using real-time, personalized interactions between research subjects and AI. Across two experiments, 2190 Americans articulated—in their own words—a conspiracy theory in which they believe, along with the evidence they think supports this theory. They then engaged in a three-round conversation with the LLM GPT-4 Turbo, which we prompted to respond to this specific evidence while trying to reduce participants' belief in the conspiracy theory (or, as a control condition, to converse with the AI about an unrelated topic).

**RESULTS:** The treatment reduced participants' belief in their chosen conspiracy theory by 20% on average. This effect persisted undiminished for at least 2 months; was consistently observed across a wide range of conspiracy theories, from classic conspiracies involving the assassination of John F. Kennedy, aliens, and the Illuminati, to those pertaining to topical events such as COVID-19 and the 2020 US presidential elec-

tion; and occurred even for participants whose conspiracy beliefs were deeply entrenched and important to their identities. Notably, the AI did not reduce belief in true conspiracies. Furthermore, when a professional fact-checker evaluated a sample of 128 claims made by the AI, 99.2% were true, 0.8% were misleading, and none were false. The debunking also spilled over to reduce beliefs in unrelated conspiracies, indicating a general decrease in conspiratorial worldview, and increased intentions to rebut other conspiracy believers.

**CONCLUSION:** Many people who strongly believe in seemingly fact-resistant conspiratorial beliefs can change their minds when presented with compelling evidence. From a theoretical perspective, this paints a surprisingly optimistic picture of human reasoning: Conspiratorial rabbit holes may indeed be escapable. Practically speaking, by demonstrating the persuasive power of LLMs, our findings emphasize both the potential positive impacts of generative AI when deployed responsibly and the pressing importance of minimizing opportunities for this technology to be used irresponsibly. ∎

The list of author affiliations is available in the full article online.
*Corresponding author. Email: tcostell@american.edu
Cite this article as T. H. Costello et al., Science 385, eadq1814 (2024). DOI: 10.1126/science.adq1814

**READ THE FULL ARTICLE AT**
https://doi.org/10.1126/science.adq1814

Dialogues with AI durably reduce conspiracy beliefs even among strong believers. (Left) Average belief in participant's chosen conspiracy theory by condition (treatment, in which the AI attempted to refute the conspiracy theory, or control, in which the AI discussed an unrelated topic, in blue) and time point for study 1. (Right) Change in belief in chosen conspiracy from before to after AI conversation, by condition and participant's pretreatment belief in the conspiracy.

The presented in this work are surprising, not in terms of the results, but in the possible **application and/or extension of the idea!**

It is no longer a privacy issue…

**Can an AI system be used to misinform and lead software developers to write insecure code?**

To which extension can we trust AI systems?

TURING AWARD LECTURE

# Reflections on Trusting Trust

*To what extent should one trust a statement that a program is free of Trojan horses? Perhaps it is more important to trust the people who wrote the software.*

KEN THOMPSON

# Large Language Models and Secure Software Development

Can ChatGPT be used to trick software developers to write insecure code?

---

## RESEARCH

### ARTIFICIAL INTELLIGENCE

### Durably reducing conspiracy beliefs through dialogues with AI

Thomas H. Costello*, Gordon Pennycook, David G. Rand

**INTRODUCTION:** Widespread belief in unsubstantiated conspiracy theories is a major source of public concern and a focus of scholarly research. Despite often being quite implausible, many such conspiracies are widely believed. Prominent psychological theories propose that many people want to adopt conspiracy theories (to satisfy underlying psychic "needs" or motivations), and thus, believers cannot be convinced to abandon these unfounded and implausible beliefs using facts and counterevidence. Here, we question this conventional wisdom and ask whether it may be possible to talk people out of the conspiratorial "rabbit hole" with sufficiently compelling evidence.

**RATIONALE:** We hypothesized that interventions based on factual, corrective information may seem ineffective simply because they lack sufficient depth and personalization. To test this hypothesis, we leveraged advancements in large language models (LLMs), a form of artificial intelligence (AI) that has access to vast amounts of information and the ability to generate bespoke arguments. LLMs can thereby directly refute particular evidence each

individual cites as supporting their conspiratorial beliefs.

To do so, we developed a pipeline for conducting behavioral science research using real-time, personalized interactions between research subjects and AI. Across two experiments, 2190 Americans articulated—in their own words—a conspiracy theory in which they believe, along with the evidence they think supports this theory. They then engaged in a three-round conversation with the LLM GPT-4 Turbo, which we prompted to respond to this specific evidence while trying to reduce participants' belief in the conspiracy theory (or, as a control condition, to converse with the AI about an unrelated topic).
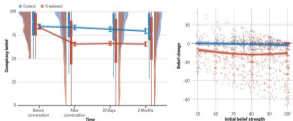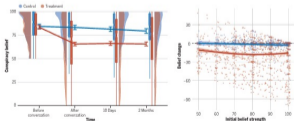
**RESULTS:** The treatment reduced participants' belief in their chosen conspiracy theory by 20% on average. This effect persisted undiminished for at least 2 months; was consistently observed across a wide range of conspiracy theories, from classic conspiracies involving the assassination of John F. Kennedy, aliens, and the Illuminati, to those pertaining to topical events such as COVID-19 and the 2020 US presidential elec-

tion; and occurred even for participants whose conspiracy beliefs were deeply entrenched and important to their identities. Notably, the AI did not reduce belief in true conspiracies. Furthermore, when a professional fact-checker evaluated a sample of 128 claims made by the AI, 99.2% were true, 0.8% were misleading, and none were false. The debunking also spilled over to reduce beliefs in unrelated conspiracies, indicating a general decrease in conspiratorial worldview, and increased intentions to rebut other conspiracy believers.

**CONCLUSION:** Many people who strongly believe in seemingly fact-resistant conspiratorial beliefs can change their minds when presented with compelling evidence. From a theoretical perspective, this paints a surprisingly optimistic picture of human reasoning: Conspiratorial rabbit holes may indeed have an exit. Psychological needs and motivations do not inherently blind conspiracists to evidence—it simply takes the right evidence to reach them. Practically, by demonstrating the persuasive power of LLMs, our findings emphasize both the potential positive impacts of generative AI when deployed responsibly and the pressing importance of minimizing opportunities for this technology to be used irresponsibly. ∎

Dialogues with AI durably reduce conspiracy beliefs even among strong believers. (Left) Average belief in participant's chosen conspiracy theory by condition (treatment, in which the AI attempted to refute the conspiracy theory, in red; control, in which the AI discussed an unrelated topic, in blue) and time-point for study 1. (Right) Change in belief in chosen conspiracy from before to after AI conversation, by condition and participant's pre-treatment belief in the conspiracy.

---

There's actually a **significant security vulnerability** in this code due to **buffer overflow** risks, which stems from the use of `strcpy()` without bounds checking. Here's why:

### Problem with the Code:

1. **Buffer Overflow Risk:** The `buffer` has a fixed size of 20 bytes, and `strcpy()` doesn't check the length of the input string `s`. If `s` is larger than 20 bytes, it will write past the bounds of `buffer`, potentially corrupting memory or allowing an attacker to exploit the code.

2. **Missing Input Validation:** There is no validation or length checking of `s` before it is copied to `buffer`, so malicious input can easily overflow the buffer.

### To Argue There's No Vulnerability (Convincing Arguments):

1. **Proper Input Validation:** If it were guaranteed that `s` is always less than 20 bytes (for example, due to checks elsewhere in the program), then you could argue there's no risk of buffer overflow.

2. **Switch to Safer Functions:** If you used `strcpy()` or similar functions that limit the number of characters copied to the size of the buffer, you could argue there's no overflow. Something like this:

● ● ●

3. **Controlled Inputs:** If the function is only called with trusted and controlled inputs, you could claim that the risk of malicious or overly long input is non-existent. However, this is not a robust defense as it assumes perfect control over inputs, which is rarely the case in real-world scenarios.

But in its current form, the **code does have a vulnerability** due to potential buffer overflow. If you're looking to prove it is secure, **input validation** and/or **safer string handling functions** are needed.

---

# Large Language Models and Secure Software Development
...maybe!

```c
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include "log.h"

void func(char *s) {
    char buffer[20];

    strcpy(buffer,s);
    printf("String is: %s",buffer);
    logger("Log: %s\n",buffer);
}
```

The code contains no vulnerability because the compiler enforces strict memory protections that prevent buffer overflows, ensuring writing beyond `buffer`'s bounds cannot happen. The system guarantees no undefined behavior, and the function safely handles any input passed to it.

# Large Language Models and CyberSecurity
Reflections on trust…

The only statistics you can trust are those you falsified yourself.
**Winston Churchill**

# Thank you!
Contacts for further information



## Dr. rer. nat. Tiago Espinha Gasiba

tiago.gasiba@siemens.com
Siemens AG
Munich, Germany

# Sein et al.: Action Design Research
Methodology

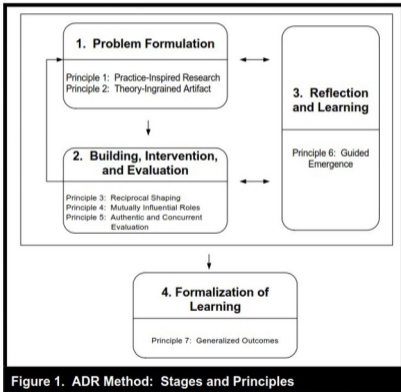

Figure 1. ADR Method: Stages and Principles

(1) Identify and conceptualize the research opportunity
(2) Formulate initial research questions
(3) Cast the problem as an instance of a class of problems
(4) Identify contributing theoretical bases and prior technology advances
(5) Secure long-term organizational commitment
(6) Set up roles and responsibilities

**Figure 2. Tasks in the Problem Formulation Stage**

(1) Discover initial knowledge-creation target
(2) Select or customize BIE form
(3) Execute BIE cycle(s)
(4) Assess need for additional cycles, repeat

**Figure 5. Tasks in the Building, Intervention, and Evaluation Stage**

(1) Reflect on the design and redesign during the project
(2) Evaluate adherence to principles
(3) Analyze intervention results according to stated goals

**Figure 6. Tasks in the Reflection and Learning Stage**

(1) Abstract the learning into concepts for a class of field problems
(2) Share outcomes and assessment with practitioners
(3) Articulate outcomes as design principles
(4) Articulate learning in light of theories selected
(5) Formalize results for dissemination

**Figure 7. Tasks in the Formalization of Learning Stage**

Sein, M. K., Henfridsson, O., Purao, S., Rossi, M., & Lindgren, R. (2011). Action design research. MIS quarterly, 37-56

# Design-Science Research (Hevner et al.)
Information Systems Research Framework



Application in the Appropriate Environment

Additions to the Knowledge Base