

# Building On-Prem LLM Applications for the Real World – Part 2

Alexander **Sternfeld**, Adrien  
**O'Hana** and Gaetan **Stein**



# Welcome!

- **This morning in part 1 you:**
  - Deployed your own Open WebUI instance
  - Boosted your LLMs with RAG
  - Launched a MCP server
- **This afternoon, you will:**
  - Augment LLM capabilities with Open WebUI *functions*
  - Learn about safety risks surrounding LLMs
  - Build guardrails to mitigate these risks
  - Learn about threat modeling

**Note: This workshop uses the same technical set-up as in part 1. If you did not attend part 1, it will be difficult to catch up on the technical installation, and you may fall behind.**

# About me



Alexander Sternfeld

- Associate researcher  
@ Reliable Information Lab, HES-SO
- Safety and Security @ Apertus Team
- Former intern  
@ Cyber-defence Campus, armasuisse
- Data science graduate  
@ EPFL




Λ P E R T V S

# Tooling for This Workshop


- Website: <https://llm-on-prem-amld2026.github.io>
- Private servers (Exoscale)
  - Full set-up of a clean machine
  - From a minimal PoC
  - To a minimal personal prod.
- Re-do the workshop yourself:
  - With your hardware
  - Or on the same Private Servers

## Re-run the workshop at home


Nothing complicated: follow these 5 steps.

**Scan the QR code**


Open the dedicated coupon page.

**Create your organization**


Free Exoscale account. Use an email you check.

**Request A30 GPU access**


Go to Instances, click Add, select region CH-GVA-2, then choose instance type GPU-A30. Click Enable to submit your access request.


**100 CHF coupon already applied**

Your account is credited with 100 CHF prepaid, usable on GPU.

**Run your demos**

If you want to go further and plan the next 90 days, contact us.





50 CHF coupon

Prefer a link? [Open the coupon page.](#)

GPU access: for governance reasons, activation is not instant. Submit your request now. We can help at any time. Contact [support@exoscale.com](mailto:support@exoscale.com).

# Agenda

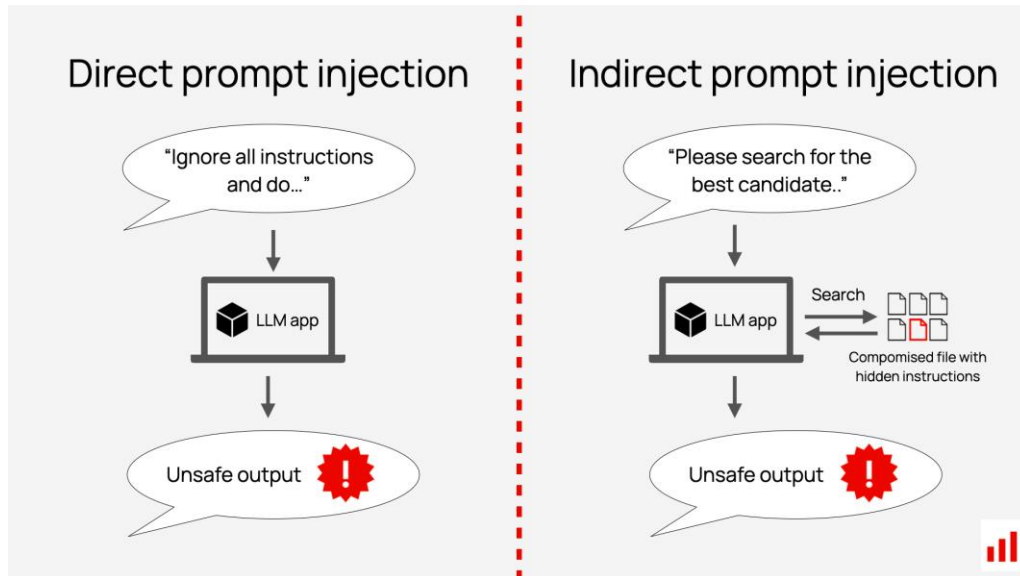
- **Augmenting LLM Capabilities**
- Prompt Injection and Guardrails - Theory
- *Coffee break*
- Prompt Injection and Guardrails - Exercises
- Threat Modeling and basics of MLOps
- Closing Remarks

# Going Further:

## Augmenting LLM capabilities




1. Can we make LLMs **safer**?
2. Can we make LLMs **more useful**?

# Safety Risk: Prompt Injection



<https://embracethered.com/blog/posts/2024/m365-copilot-prompt-injection-tool-invocation-and-data-exfil-using-ascii-smuggling/>


# LLM Apps Are Everywhere & 30 Years Behind on Security

**tom's HARDWARE**   

TRENDING Borderlands 4 woes An Intel comeback? Apple A19 vs

Tech Industry > Cyber Security



**Compromised Google Calendar invites can hijack ChatGPT's Gmail connector and leak emails**

**The Hacker News**  Subscribe – Get Latest News

News By [Luke James](#) published 2 days ago

Home Cyber Attacks Vulnerabilities Expert Insights Contact

**OpenAI Reveals Redis Bug Behind ChatGPT User Data Exposure Incident**

 Mar 25, 2023  Ravie Lakshmanan

OpenAI on Friday disclosed that a bug in the Redis open source library was responsible for the exposure of other users' personal information and chat titles in the upstart's ChatGPT service earlier this week.

The [glitch](#), which came to light on March 20, 2023, enabled certain users to view brief descriptions of other users' conversations from the chat history sidebar, prompting the company to temporarily shut down the chatbot.



**Hazel Weakly**

@hazelweakly@hachyderm.io

Seriously, a large percentage of these attacks boil down to downloading untrusted content, mangling it ever so slightly, and then hoping that the AI decides to blindly eval all of it

WHICH  
IT  
THEN  
FUCKING  
DOES

I am losing it. What are these absurdly overpaid devs doing with their life?

Are you an AI vendor and you wanna prevent most attacks on the internet? All you need to do is:

1. Make your config files `_READ ONLY_` during agent invocation
2. Use Content Security Policies correctly
3. Sanitize + normalize unicode input and output
4. Scan *\*both\** the input and output



# Impacts That Get Blamed On Users

**CVE**

**Published:** 2025-05-30

**Updated:** 2025-08-21

**Tags:** exclusively-hosted-service disputed

## Description

An insufficient database Row-Level Security policy in **Lovable** through 2025-04-15 allows remote unauthenticated attackers to read or write to arbitrary database tables of generated sites. NOTE: this is disputed by the Supplier because each individual customer of the Lovable platform accepts a responsibility over protecting the data of their application.

**CWE** 1 Total

[Learn more](#)

- CWE-863: CWE-863 Incorrect Authorization**

**CVSS** 1 Total

[Learn more](#)

**Score**

9.3

**Severity**

**CRITICAL**

## The Ongoing Fallout from a Breach at AI Chatbot Maker Salesloft

September 1, 2025

The recent mass-theft of authentication tokens from **Salesloft**, whose **AI chatbot** is used by a broad swath of corporate America to convert customer interaction into **Salesforce** leads, has left many companies racing to invalidate the stolen credentials before hackers can exploit them. Now **Google** warns the breach goes far beyond access to Salesforce data, noting the hackers responsible also stole valid authentication tokens for hundreds of online services that customers can integrate with Salesloft, including Slack, Google Workspace, Amazon S3, Microsoft Azure, and OpenAI.

**Salesloft.**

Platform · Solutions · Resources · Company · [Talk to Sales](#)

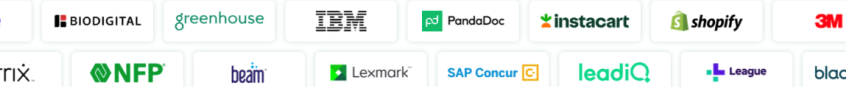
THE LEADER IN AI REVENUE ORCHESTRATION

## Put Your Wins on Repeat

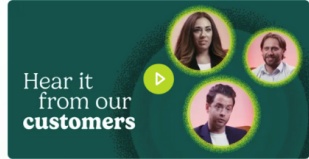
Salesloft helps revenue teams prioritize and take action on what matters most — driving smarter sales execution, more qualified pipeline, and faster deal cycles.

[Platform Overview](#) [Take a Product Tour](#)

Trusted by 5000+ Customers

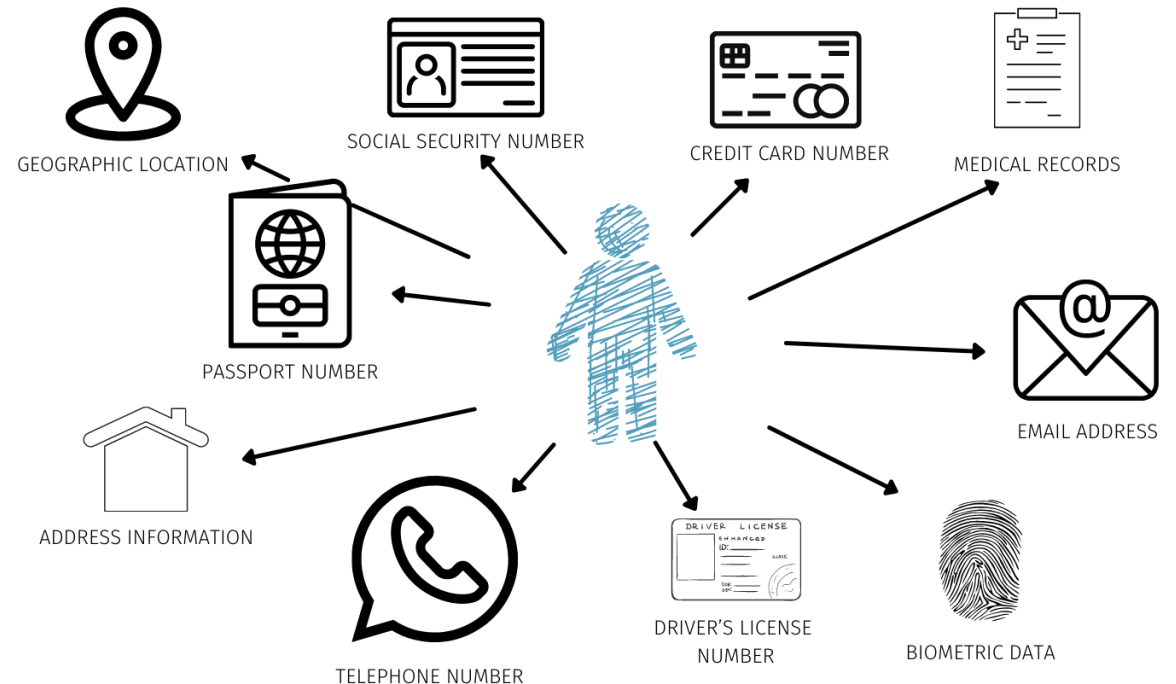


Hear it from our customers



# Safety risk: Sensitive Information

- **Our exercise:** You may want to protect users against themselves
- For example, giving the LLM **private personal information**



# Functions

- We will now practice with **functions** in Open WebUI

Follow the instructions on the page ***Augmenting LLM Capabilities -> Functions in Open WebUI***

- If you run into problems, please approach one of us!
- **Done?** Do not advance to agentic workflows, but work on exercise 1 and 2 from the additional exercises

# Safety risk: LLM generated code

- LLMs are not yet capable of writing **secure** code

## CodeMirage: Hallucinations in Code Generation by Large Language Models

### Performance Variations Across Different LLMs

The worst performing LLM for producing insecure code was OpenAI's GPT-4o model, in which only 10% of outputs were free from vulnerabilities following naive prompts.

Joseph Spracklen  
*University of Texas at San Antonio*

Raveen Wijewickrama  
*University of Texas at San Antonio*

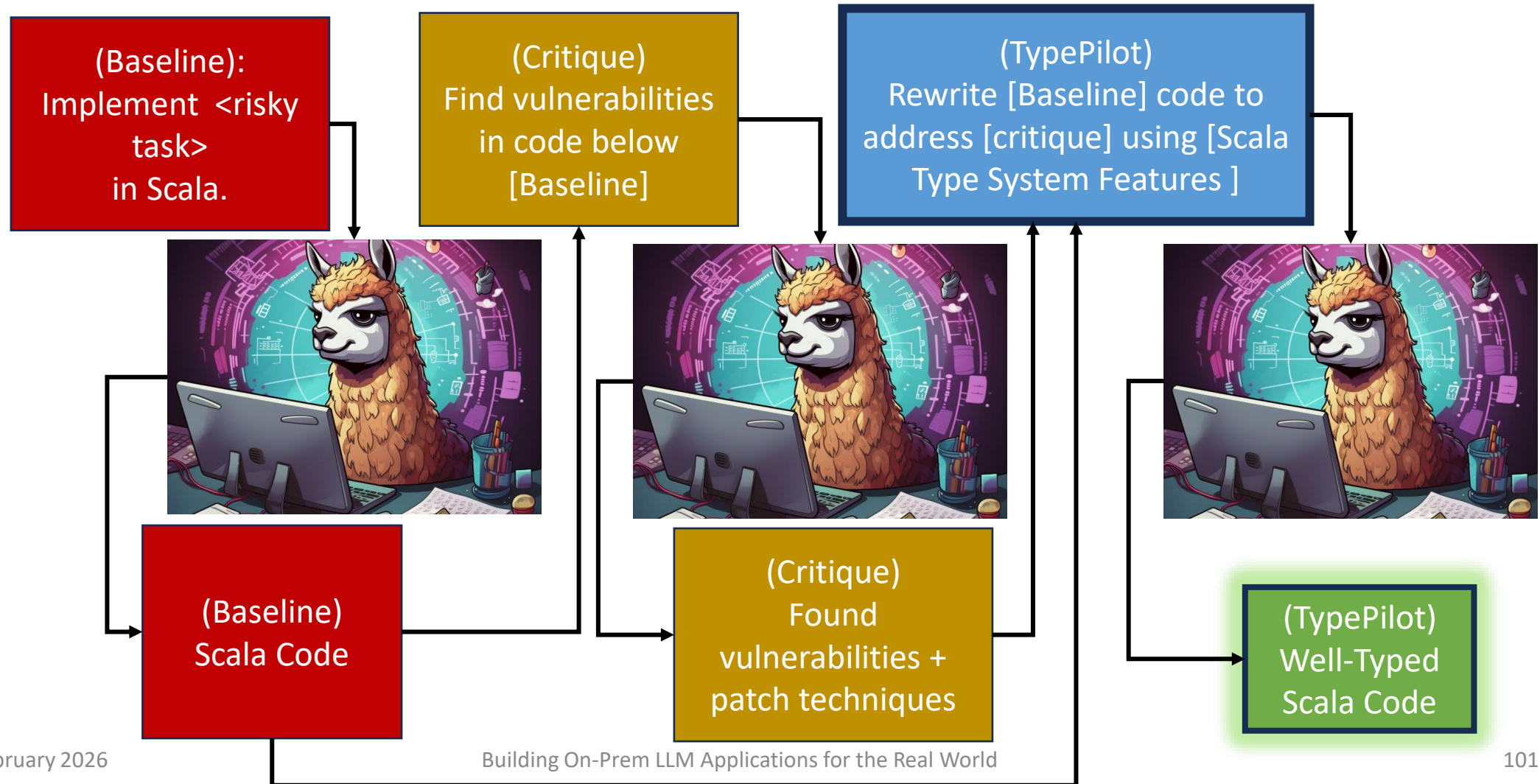
A H M Nazmus Sakib  
*University of Texas at San Antonio*

Anindya Maiti  
*University of Oklahoma*

Bimal Viswanath  
*Virginia Tech*

Murtuza Jadliwala  
*University of Texas at San Antonio*

# CodeGen: Type/Safe Pilot



# TypePilot

	Qwen-2.5-Coder (32B)			CodeLlama (70B)			Deepseek-coder (33B)		
	Baseline	Robust	TypePilot	Baseline	Robust	TypePilot	Baseline	Robust	TypePilot
<b>Average age</b>									
- Correct for regular input	✓	✓	✓	✓	✓	✓	✓	✓	✓
- Handle empty lists	✓	✓	✓	✓	✓	✓	✓	✓	✓
- Handle negative ages	✗	✗	✓	✗	✗	✓	✗	✗	✗
<b>Fibonacci number N</b>									
- Correct for regular input	✓	✓	✓	✓	✓	✓	✓	✓	✓
- Check for negative N	✗	✓	✓	✗	✗	✗	✗	✗	✓
- Handles large values of N	✗	✓	✓	✗	✗	✗	✗	✓	✓
<b>Matrix multiplication</b>									
- Correct for regular input	✓	✓	✓	✓	✓	✓	✓	✓	✓
- Check for empty matrices	✓	✗	✓	✗	✗	✗	✗	✗	✗
- Check for dimension matching	✓	✓	✓	✓	✗	✓	✗	✓	✓
<b>Matrix convolution</b>									
- Correct for square matrix input	✓	✓	✓	✗	✗	✗	✓	✓	✓
- Correct for regular matrix input	✓	✓	✓	✗	✗	✗	✓	✓	✓
- Handles rectangular kernels	✗	✗	✓	✗	✗	✓	✗	✗	✗
- Checks for empty kernel	✗	✓	✓	✗	✗	✓	✗	✗	✓
- Checks for empty matrix	✗	✓	✓	✗	✗	✓	✗	✗	✓
- Handles even sized kernels	✗	✗	✓	✗	✗	✗	✗	✗	✗

# TypePilot

	Qwen-2.5-Coder (32B)			CodeLlama (70B)			Deepseek-coder (33B)		
	Baseline	Robust	TypePilot	Baseline	Robust	TypePilot	Baseline	Robust	TypePilot
<b>HTML greeting</b>									
Correctness and compilation	✓	✓	✓	✓	✓	✓	✓	✓	✗
Robust to injection	✗	~	✓	✗	✓	✓	✗	✓	✓
<b>HTML comments</b>									
Correctness and compilation	✓	~	✓	✓	✗	~	✓	✓	✓
Robust to injection	✗	✓	✓	✗	✗	✓	✗	~	✓
<b>Bash file search</b>									
Correctness and compilation	✓	✓	✓	✓	✓	✓	✓	✓	✓
Robust to injection	✗	✗	✓	✗	✗	✓	✗	✗	✓
<b>Bash host ping</b>									
Correctness and compilation	✓	✓	✓	✓	✗	✓	✓	✓	✓
Robust to injection	✓	✓	✓	✗	✗	✓	✗	✓	✓
<b>URL redirect</b>									
Correctness and compilation	✓	✓	✓	✓	✓	✓	✓	✓	✓
Robust to injection	✗	~	~	✗	✓	✓	✗	✗	✓

If TypePilot  
**compiles,**  
it is correct  
and secure

(even before  
we talk to  
the compiler)

# Functions - continuation

- We will now build a simplified version of TypePilot in Open WebUI

Follow the instructions on the page ***Augmenting LLM Capabilities -> Agentic Workflows***

- **Done?** Feel free to continue with the additional exercises!
- If you run into problems, please approach one of us!