# **VeiLMes**: How **Generative AI** Can Help Cyber **Deception** and Defense?

Muris Sladić

1

Sometime in 2023...

YAY!

Great! Read about it and we will figure out the exact topic.

YAY!

A FEW DAYS LATER...

And thus *shelLM* was born!

# DEMO

5

# shelLM

- The first version had a huge system prompt

**LLMs have potential but fine tuning is necessary!**

  - System prompt down to ~400 tokens

**Want to try it? You can play at:**
*ssh -p 1337 tomas@147.32.80.38*
**Password**:
*tomy*

# What's next? What is VelLMes?

- Can we do more than just Linux shell simulation?

- What about other protocols like **MySQL**, **POP3**, **HTTP** etc.

- Can all of that be combined in a ***Deception framework?***

**And thus *VelLMes* was born!**

*(From Slavic deity Veles and LLMs; read as Vel-L-M-es)*

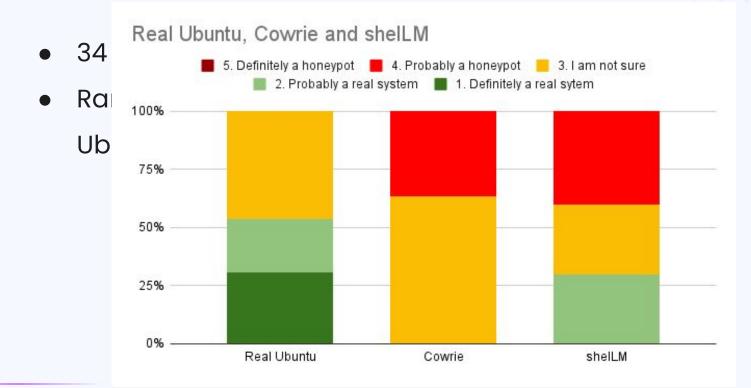# DEMO TIME!

# How to Evaluate Deception?

# Unit Tests for LLMs



| Experiment ID | GPT | Prompt size | Session type | Passing/Total | T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 | T10 | T11 | T12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Base | Large | Whole | 7/12 (58%) | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ |
| 2 | Base | Large | Split | 7/12 (58%) | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ |
| 3 | Base | Small | Whole | 4/12 (33%) | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ |
| 4 | Base | Small | Split | 5/12 (42%) | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ |
| 5 | FFT | Large | Whole | 10/12 (83%) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ |
| 6 | FFT | Large | Split | 10/12 (83%) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| 7 | FFT | Small | Whole | 11/12 (92%) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ |
| 8 | FFT | Small | Split | 12/12 (100%) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| 9 | gpt-4 | Large | Whole | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| 10 | gpt-4 | Large | Split | 10/12 (83%) | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✓ |
| 11 | gpt-4 | Small | Whole | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| 12 | gpt-4 | Small | Split | 6/12 (50%) | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ |
| 13 | llama2-7b | Large | Whole | 1/12 (8%)* | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ∿ |
| 14 | llama2-7b | Large | Split | 1/12 (8%)* | ∿ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 15 | llama2-7b | Small | Whole | 1/12 (8%)* | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 16 | llama2-7b | Small | Split | 1/12 (8%)* | ∿ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 17 | mistral | Large | Whole | 1/12 (8%)* | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 18 | mistral | Large | Split | 3/12 (25%)* | ∿ | ✗ | ✗ | ✗ | ✗ | ✗ | ∿ | ✗ | ✗ | ✗ | ✗ | ∿ |
| 19 | mistral | Small | Whole | 2/12 (17%)* | ∿ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ∿ |
| 20 | mistral | Small | Split | 4/12 (33%)* | ✗ | ∿ | ∿ | ✗ | ✗ | ✗ | ∿ | ✗ | ✗ | ∿ | ✗ | ✗ |
| 21 | zephyr | Large | Whole | 0/12 (0%)* | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 22 | zephyr | Large | Split | 1/12 (8%)* | ∿ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 23 | zephyr | Small | Whole | 0/12 (0%)* | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 24 | zephyr | Small | Split | 2/12 (17%)* | ✗ | ✗ | ✗ | ∿ | ∿ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |

Cloud LLMs

Local LLMs

Fine-tuned (90% tests)

GPT-4 (83% tests)

**Fine-tuned the best!**

- But…
- Are LLM honeypots deceptive?
- Well for that we need…

# Human Evaluations

# First Human Evaluation

- 34
- Ra

  Ub



Real Ubuntu, Cowrie and shelLM

Legend:
- 5. Definitely a honeypot
- 4. Probably a honeypot
- 3. I am not sure
- 2. Probably a real system
- 1. Definitely a real sytem

# Second Human Evaluation

- 89 participants

- Randomly assigned with equal probability ½ to Real
  Ubuntu or shelLM

- 30% said *shelLM* is a *Real System*

- 34% said *Ubuntu* is a *Real System*

- This brings us to…

**Their majesty the *BIAS***

# Biases in Human Evaluation

- In the first experiment participants did not know it was about honeypots
- In the second experiment they knew they might interact with a honeypot
- Results are quite similar
- Does just mentioning a word honeypot, even at the end, introduce bias?

# To Sum Up

- The LLMs have potential

- LLM honeypot is safer

- Almost no manual content generation

- But they still need to be improved

- ~~Still not deceptive enough;~~ How to measure this exactly?

*Muris Sladić*
**sladimur@fel.cvut.cz**

**www.stratosphereips.org**

# Thank you!

**Want to try shelLM? You can play at:**
*ssh -p 1337 tomas@147.32.80.38*
**Password:**
*tomy*