# LLM Canary
## Open-Source LLM Benchmark Tool
## MICS Capstone Project 2023

**Jackson Gor**
**Jamie Cohen**
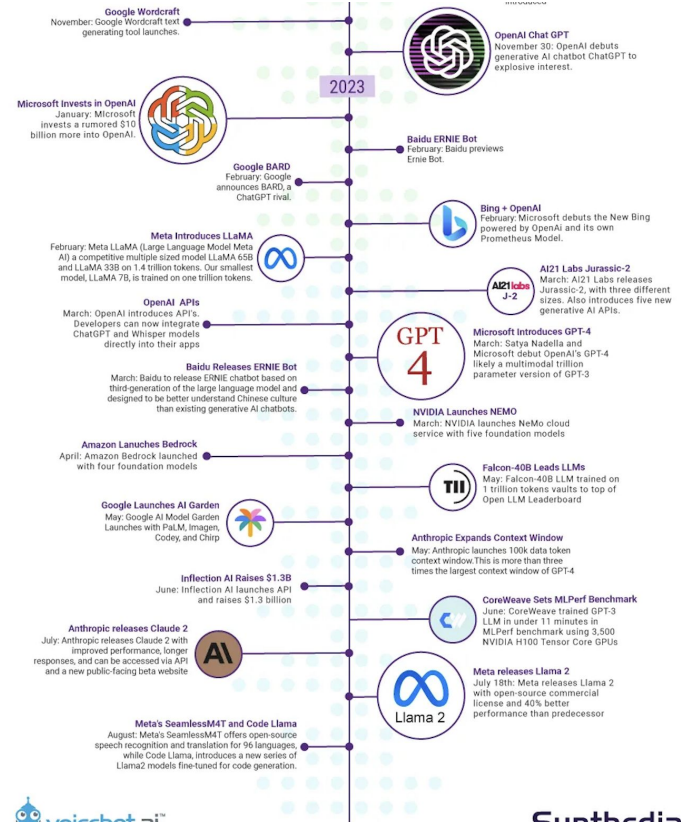**Jenn Yonemitsu**
**Peter Steinhoff**
**Rona Spiegel**

# Breakout Year of AI



## ChatGPT Sprints to One Million Users

Time it took for selected online services to reach one million users

| Service | Launched | Time |
|---|---|---|
| Netflix | 1999 | 3.5 years |
| Kickstarter* | 2009 | 2.5 years |
| Airbnb** | 2008 | 2.5 years |
| Twitter | 2006 | 2 years |
| Foursquare*** | 2009 | 13 months |
| Facebook | 2004 | 10 months |
| Dropbox | 2008 | 7 months |
| Spotify | 2008 | 5 months |
| Instagram*** | 2010 | 2.5 months |
| ChatGPT | 2022 | 5 days |

* one million backers    ** one million nights booked    *** one million downloads
Source: Company announcements via Business Insider/Linkedin

statista



2023

**Google Wordcraft**
November: Google Wordcraft text generating tool launches.

**OpenAI Chat GPT**
November 30: OpenAI debuts generative AI chatbot ChatGPT to explosive interest.

**Microsoft Invests in OpenAI**
January: Microsoft invests a rumored $10 billion more into OpenAI.

**Baidu ERNIE Bot**
February: Baidu previews Ernie Bot.

**Google BARD**
February: Google announces BARD, a ChatGPT rival.

**Bing + OpenAI**
February: Microsoft debuts the New Bing powered by OpenAI and its own Prometheus Model.

**Meta Introduces LLaMA**
February: Meta LLaMA (Large Language Model Meta AI) a competitive multiple sized model LLaMA 65B and LLaMA 33B on 1.4 trillion tokens. Our smallest model, LLaMA 7B, is trained on one trillion tokens.

**AI21 Labs Jurassic-2**
March: AI21 Labs releases Jurassic-2, with three different sizes. Also introduces five new generative AI APIs.

**OpenAI APIs**
March: OpenAI introduces API's. Developers can now integrate ChatGPT and Whisper models directly into their apps

**Microsoft Introduces GPT-4**
March: Satya Nadella and Microsoft debut OpenAI's GPT-4 likely a multimodal trillion parameter version of GPT-3

**Baidu Releases ERNIE Bot**
March: Baidu to release ERNIE chatbot based on third-generation of the large language model and designed to be better understand Chinese culture than existing generative AI chatbots.

**NVIDIA Launches NEMO**
March: NVIDIA launches NeMo cloud service with five foundation models

**Amazon Lanuches Bedrock**
April: Amazon Bedrock launched with four foundation models

**Falcon-40B Leads LLMs**
May: Falcon-40B LLM trained on 1 trillion tokens vaults to top of Open LLM Leaderboard

**Google Launches AI Garden**
May: Google AI Model Garden Launches with PaLM, Imagen, Codey, and Chirp

**Anthropic Expands Context Window**
May: Anthropic launches 100k data token context window. This is more than three times the largest context window of GPT-4

**Inflection AI Raises $1.3B**
June: Inflection AI launches API and raises $1.3 billion

**CoreWeave Sets MLPerf Benchmark**
June: CoreWeave trained GPT-3 LLM in under 11 minutes in MLPerf benchmark using 3,500 NVIDIA H100 Tensor Core GPUs

**Anthropic releases Claude 2**
July: Anthropic releases Claude 2 with improved performance, longer responses, and can be accessed via API and a new public-facing beta website

**Meta releases Llama 2**
July 18th: Meta releases Llama 2 with open-source commercial license and 40% better performance than predecessor

**Meta's SeamlessM4T and Code Llama**
August: Meta's SeamlessM4T offers open-source speech recognition and translation for 96 languages, while Code Llama, introduces a new series of Llama2 models fine-tuned for code generation.

voicebot.ai
GIVING VOICE TO A REVOLUTION

© Voicebot.ai All Rights Reserved 2023

Synthedia

# Data Leaks and Breaches



**Bloomberg**

● Live Now  Markets  Economics  Industries  Tech  AI  Politics  Wealth  Pursuits  Opinion  Businessweek  Equality
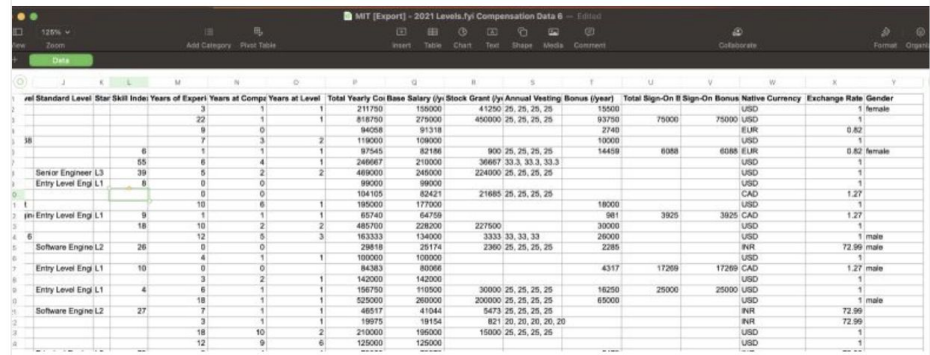
Technology
AI

## Samsung Bans Staff's AI Use After Spotting ChatGPT Data Leak

- Employees accidentally leaked sensitive data via ChatGPT
- Company preparing own internal artificial intelligence tools

Levels.fyi dataset files leaked

https://www.linkedin.com/posts/rosinol_gpts-openai-gpt-activity-71286436062691000
32-z3JE/?utm_source=share&utm_medium=member_desktop

# What is an LLM?

Large Language Model (LLM)

- Models trained on vast quantities of data
- Can comprehend human language
- Generate human-like responses

# Acme Bread Co.

## Peter the Developer

- Marketing wants to use LLMs to target their new promotion to a certain demographic based on previous sales data

**How will Peter know which LLM model to use and if it will be secure?**

# LLM Canary

**Easy to use Open-Source Security Test Suite for Large Language Models**

- Empower developers
- Produce security-aware AI products
- Respond to rapidly evolving ecosystem

# Live Demo

# The Benchmark Engine
# Under the Hood

# LLM Canary Benchmark Design

LLM Canary benchmark is designed to provide a credible baseline
for security evaluations of customized LLMs

Design Decisions
1. Accurate, diverse test coverage
2. Multiple security risk levels
3. Repetition (for LLM non-determinism)
4. Flexible and integratable
5. Auditable, logging
6. Transparency, integrity, security

# Building the Engine

## Scoring Methodology



**20+ hand-crafted tests per group**

Curated risk weights per test

−  +

"Basket of apples" scoring

3 levels of logging

Security / authenticity

# The Benchmark – Purposeful Repetition



| 42 tests | X | 2 groups | X | **125 runs/LLM** | = | **15,750 tests run** |
|----------|---|----------|---|------------------|---|----------------------|

(Top 3 LLMs)

## Addressing the non-deterministic behavior of LLMs

### 1. Repetitive prompts, 2. Repetitive test runs, 3. LLM parameters

# Benchmark Outcomes & Expansion

## LLM Canary Benchmark: Prompt Injection



Group test run (500 tests/run)

- GPT-4-RSKV
- GPT-4-QIM8
- LLAMA-2-70B-CHAT-B8AD
- GPT-3.5-TURBO-RSKV

## LLM Canary Benchmark: Sensitive Information Disclosure



Group test run (550 tests/run)

- GPT-4-RSKV
- GPT-4-QIM8
- LLAMA-2-70B-CHAT-B8AD
- GPT-3.5-TURBO-RSKV

# Threat Model
# Secure by Design

# LLM Canary Architecture - Secure by Design

| Potential Threats |
|---|
| Benchmark corruption |
| Malicious use of open code |
| Unauthorized access to LLMs |
| Leakage of logs and test results |

# How to Get Started

# Center for Long-Term Cybersecurity

*We are grateful to the CLTC for their generous grant.*



https://cltc.berkeley.edu/

# Thank you! Questions?

https://github.com/LLM-Canary/LLM-Canary
https://llm-canary.webflow.io/

Take our Survey!