# Ontinue

# Generative AI:
# What will change with the rise
# of GPT in Cybersecurity?

**Theus Hossmann**
Director of Data Science
**Ontinue**

# Seperating hype from reality

GEN–AI
DEMYSTIFIED

AI FOR SECOPS

USING AI
RESPONSIBLY

ATTACKERS
USING AI

# Generative AI Demystified: Building a GPT-based Assistant

# Grounding & Retrieval–Augmented Generation

Step 1: Grounding

## User Prompt

Who won the World Cup 2022?

GPT

Argentina won the 2022 FIFA World Cup. They were crowned the champions after winning the final against France 4–2 on penalties following a 3–3 draw after extra time.

## System Prompt

### 2022 FIFA World Cup

文A 114 languages ⌄

Article    Talk

Read    View source    View history    Tools ⌄

From Wikipedia, the free encyclopedia

*"2022 World Cup" redirects here. For other competitions of that name, see 2022 World Cup (disambiguation).*

*"FIFA 2022" redirects here. For the video game, see FIFA 22.*

The **2022 FIFA World Cup** was the 22nd FIFA World Cup, the world championship for national football teams organized by FIFA. It took place in Qatar from 20 November to 18 December 2022, after the country was

The tournament was won by host country France, who beat defending champions Brazil 3–0 in the final. France won their first title, becoming the seventh nation to win a World Cup, and the sixth (after Uruguay, Italy, England, West Germany and Argentina) to win the world cup on on home soil. As of 2022, they are most recent team to win the tournament on home soil. Croatia, Jamaica, Japan and South Africa made their first appearances in the finals.

across five cities. Qatar entered the event—their first World Cup—automatically as the host's national team, alongside 31 teams determined by the qualification process.

Argentina were crowned the champions after winning the final against the title holder France 4–2 on penalties following a 3–3 draw after extra time. It was Argentina's third title and their first since 1986, as well being the first nation from outside of Europe to win the tournament since 2002. French player Kylian Mbappé became the first player to score a hat-trick in a World Cup final since Geoff Hurst in the 1966 final and won the Golden Boot as he scored the most goals (eight) during the tournament. Argentine captain Lionel Messi
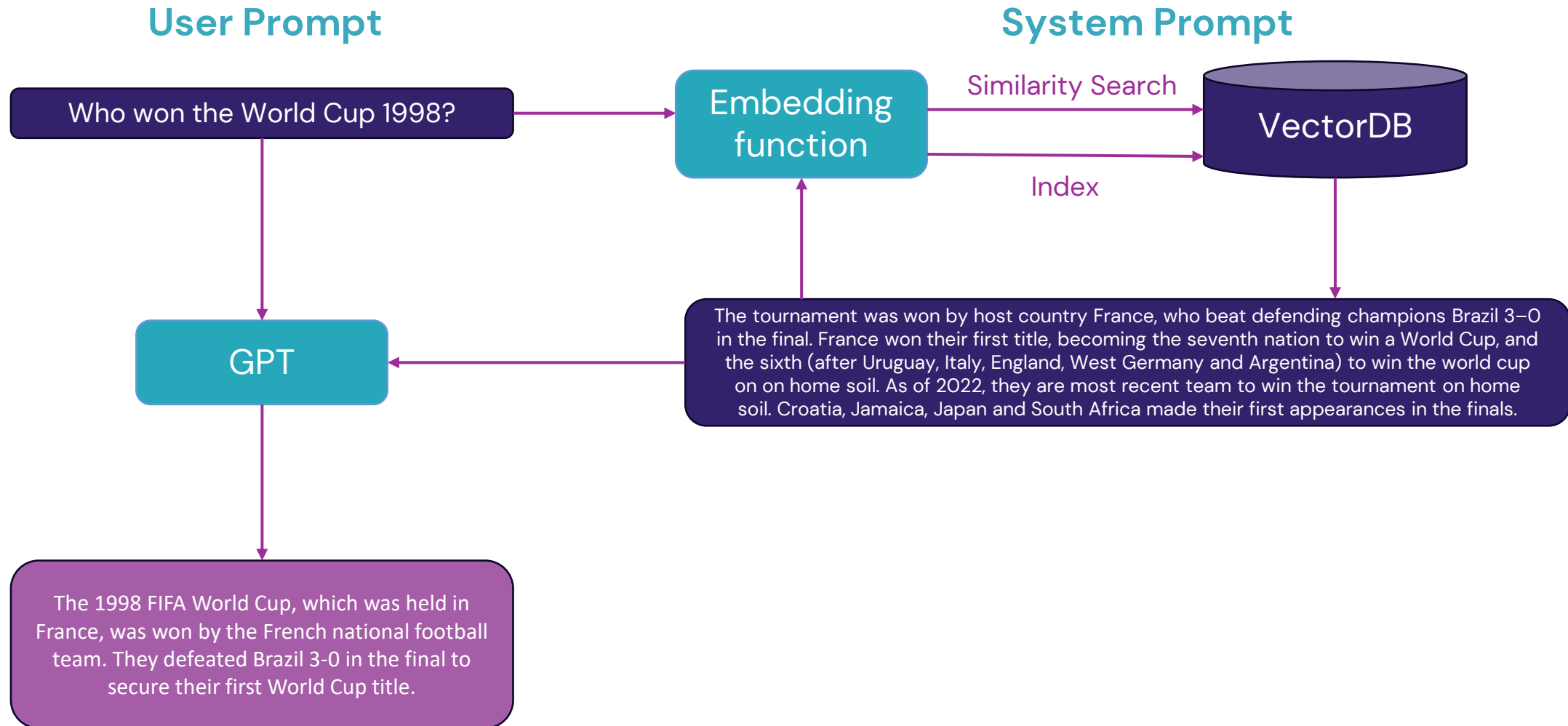
**2022 FIFA World Cup**

كأس ...dam 2022

FIFA WORLD CUP
Qatar2022

### How do we find the system prompt?

# Grounding & Retrieval–Augmented Generation
## Step 2: Retrieval–Augmented Generation with Semantic Search

**User Prompt**

**System Prompt**

Who won the World Cup 1998?

Embedding function

Similarity Search

VectorDB

Index

GPT

The tournament was won by host country France, who beat defending champions Brazil 3–0 in the final. France won their first title, becoming the seventh nation to win a World Cup, and the sixth (after Uruguay, Italy, England, West Germany and Argentina) to win the world cup on on home soil. As of 2022, they are most recent team to win the tournament on home soil. Croatia, Jamaica, Japan and South Africa made their first appearances in the finals.

The 1998 FIFA World Cup, which was held in France, was won by the French national football team. They defeated Brazil 3-0 in the final to secure their first World Cup title.

# Building a GPT-based Assistant

Step 3: Agents

# AI for SecOps:
# Beyond Threat Detection

# Supercharging the Defenders

Improve Productivity and Collaboration – With Streamlined Natural Language Workflows

## Summarizing Security Incidents



## GPT-Powered Chat Interface

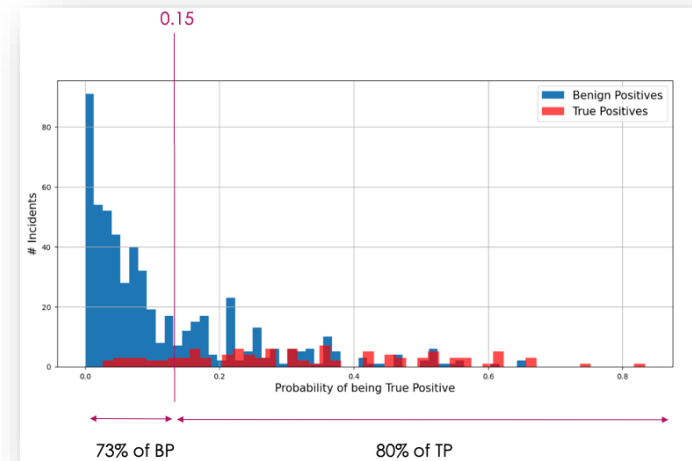

## Filtering Rules Recommendation

# Understand what you defend

Modeling the structural and operational context of the environment you defend

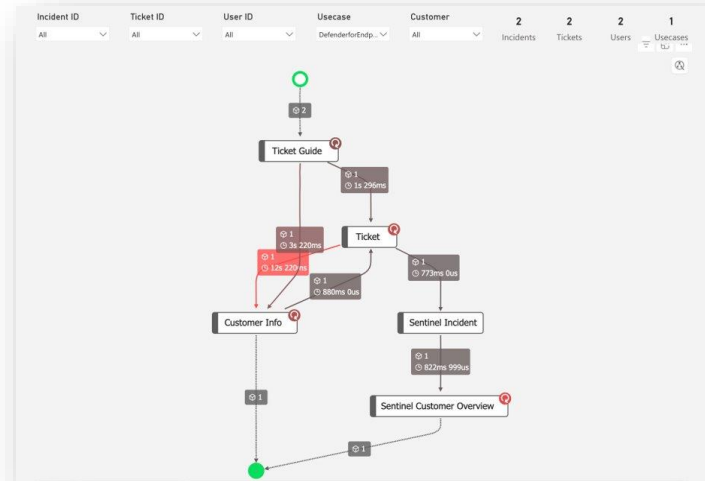Identifying Critical Assets
for better protection & response

Scoring Security Incidents
for better response

# Supercharging the Defenders

Modeling the activity and behavior of SOC Analysts to build and validate automations

Identifying efficient & effective
response patterns for faster response

# Using AI Responsibly

# AI is introducing numerous risks and challenges

Data privacy concerns

AI Behavioral Vulnerabilities

Insecure Code Generation

Copyright and Ownership

Threat Actor Evolution

Bias and Discrimination

Enterprise, SaaS, and 3rd party Security

Trust and Reputation

# Overcoming AI challenges

### Risk Framework
Identify relevant risks and their impacts to your enterprise

### Policy & Governance
Establish organizational policies on who and how can use these tools in a manner that mitigates the risks to acceptable levels

### AI provider selection
Choose appropriate AI providers depending on the security and policy customization offered to clients, such as opt-out and data retention

### Ethical AI
Create an ethical framework that guides the development and deployment of AI
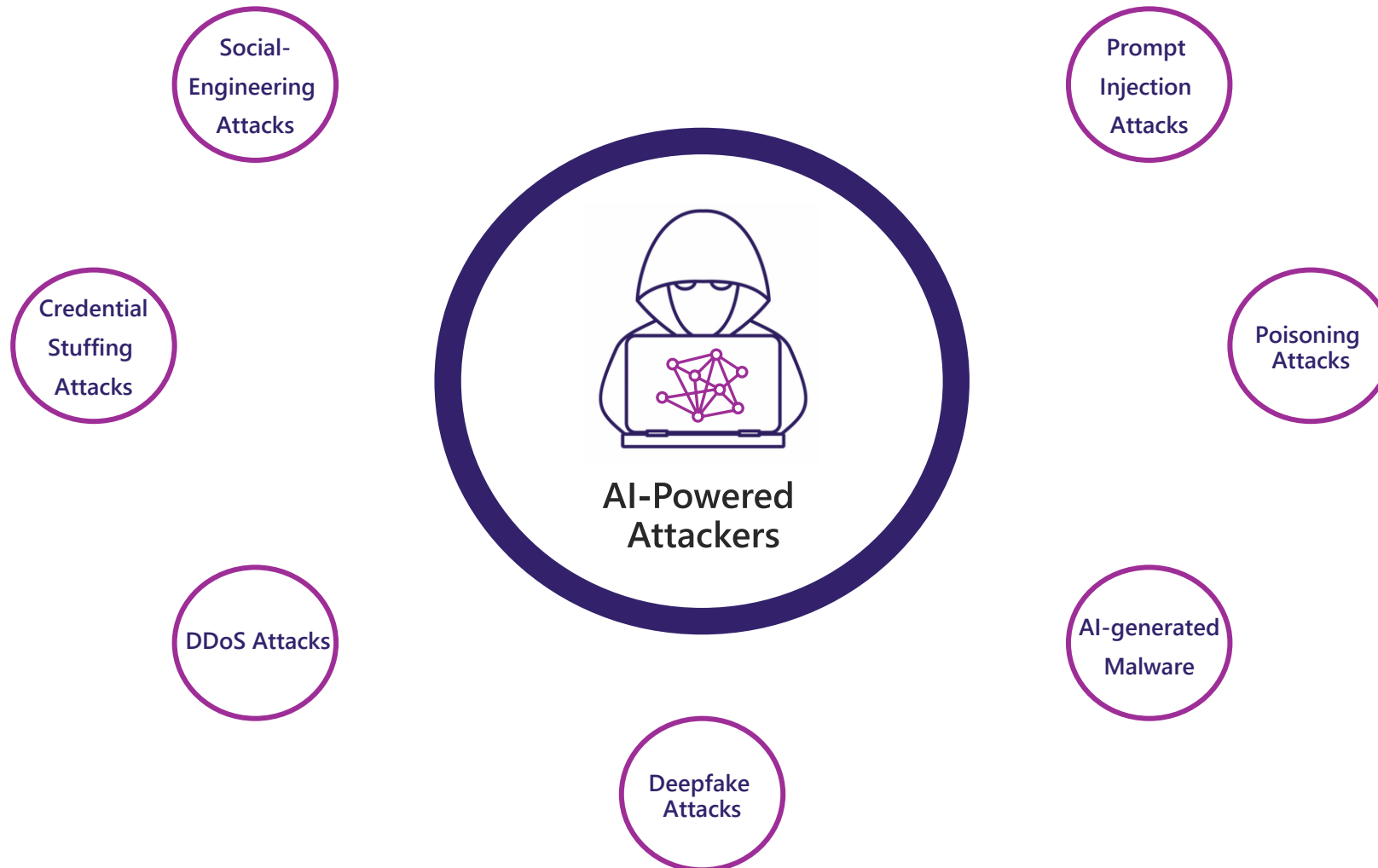
### Training & Best practices
Build organizational awareness by educating and upskilling employees around AI

# Attackers using AI

# AI– and generative AI–enabled attacks

Social-Engineering Attacks

Prompt Injection Attacks

Credential Stuffing Attacks

Poisoning Attacks

AI-Powered Attackers

DDoS Attacks

AI-generated Malware

Deepfake Attacks

# Future of AI

Ontinue

# Thank you!

Learn more about Ontinue