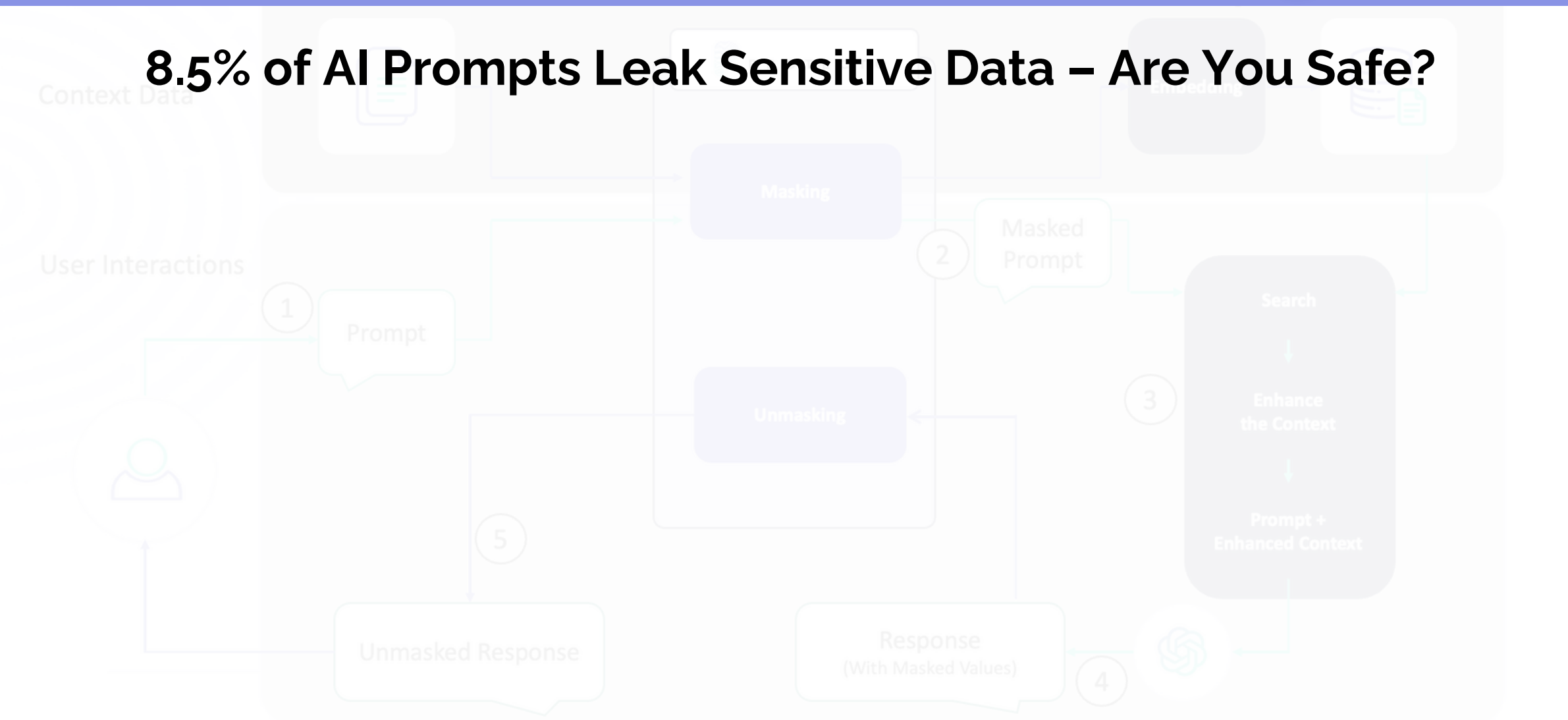


LLM Security

8.5% of AI Prompts Leak Sensitive Data – Are You Safe?



Introduction

AI tools like ChatGPT can transform workplace productivity—until an employee accidentally shares customer data, contract terms, or confidential info.

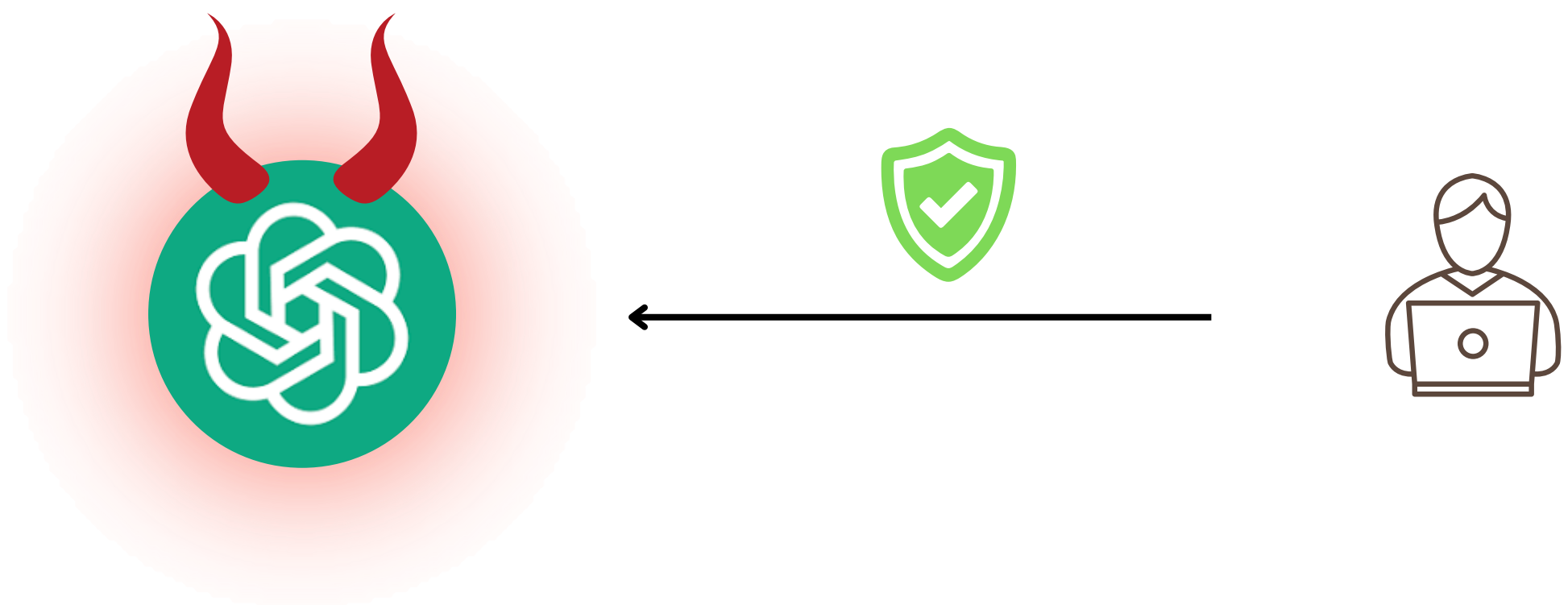
The dilemma:

- **Employees need AI to work faster**
- **Enterprises can't risk data leaks**

People relentlessly feed sensitive data into AI tools - customer records, contracts, financials - without a second thought.

About **8.5%** enterprise AI queries contain sensitive data.

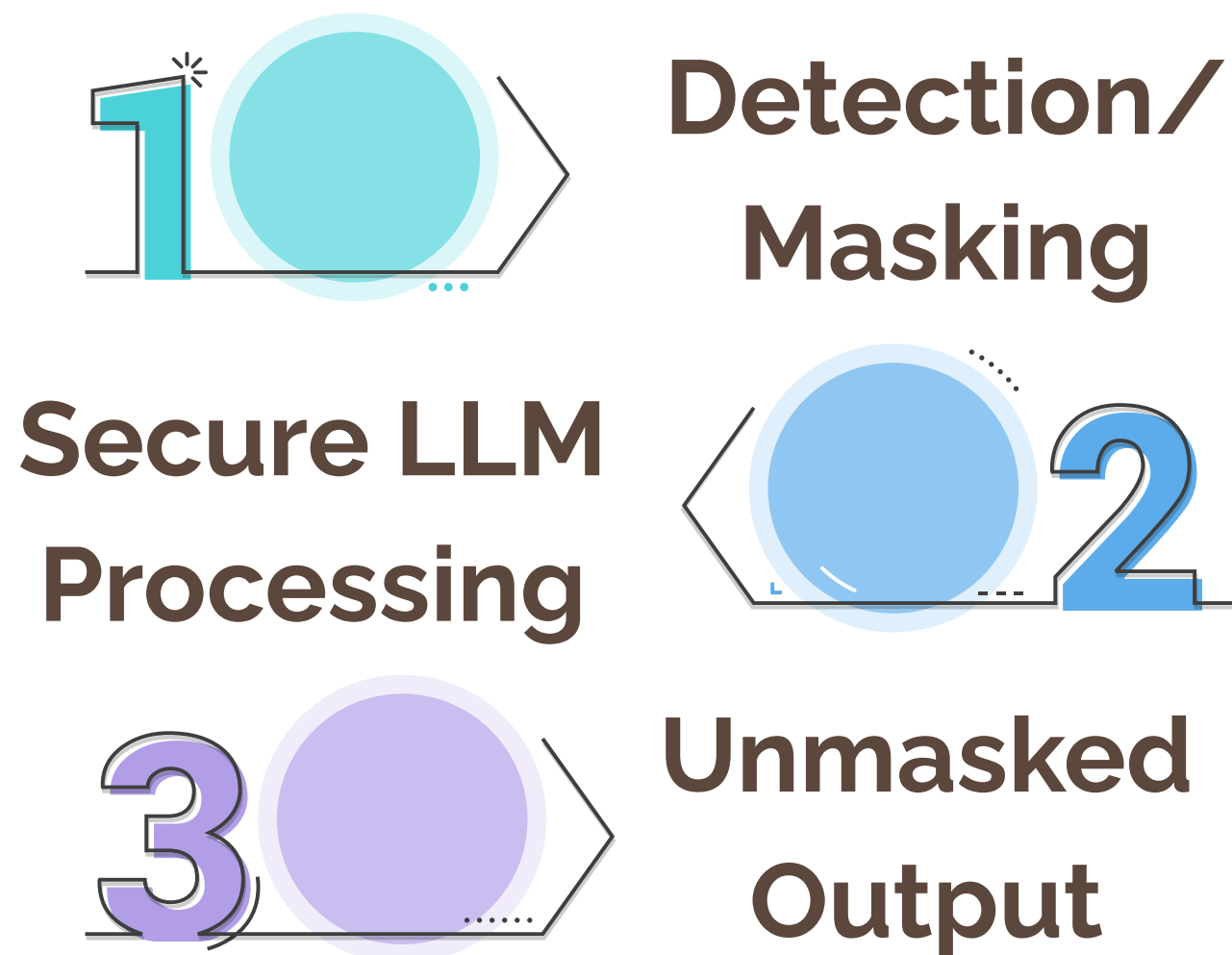
But, how do you enable AI without compromising security?



The solution isn't saying no to AI - but saying **yes to safe AI**. Lets learn how tools like **GPTGuard** can help achieve that ->

How GPTGuard works?

GPTGuard's privacy protection happens in three steps:



Intelligent Detection

- Combines AI, regex patterns, and custom rules to identify 50+ sensitive data types
- Adapts to your needs: Learns company-specific terms (project codenames, internal IDs)

Context-Aware Masking

- Replaces sensitive values with secure tokens (e.g., "Raj Mehta" → "[PERSON_1]")
- Preserves data format and meaning so LLMs understand the context

Secure Knowledge Retrieval

- HyperSearch RAG pulls answers from your documents using masked data only
- Always shows sources so you can verify responses

GPTGuard is powered by privacy APIs under the hood, lets see how it works -

Protecto's Privacy APIs

GPTGuard is powered by Protecto's Privacy API to deliver-



Protecto's API:

- Secure AI without sacrificing accuracy
- Audit-proof data handling
- Enterprise trust at scale

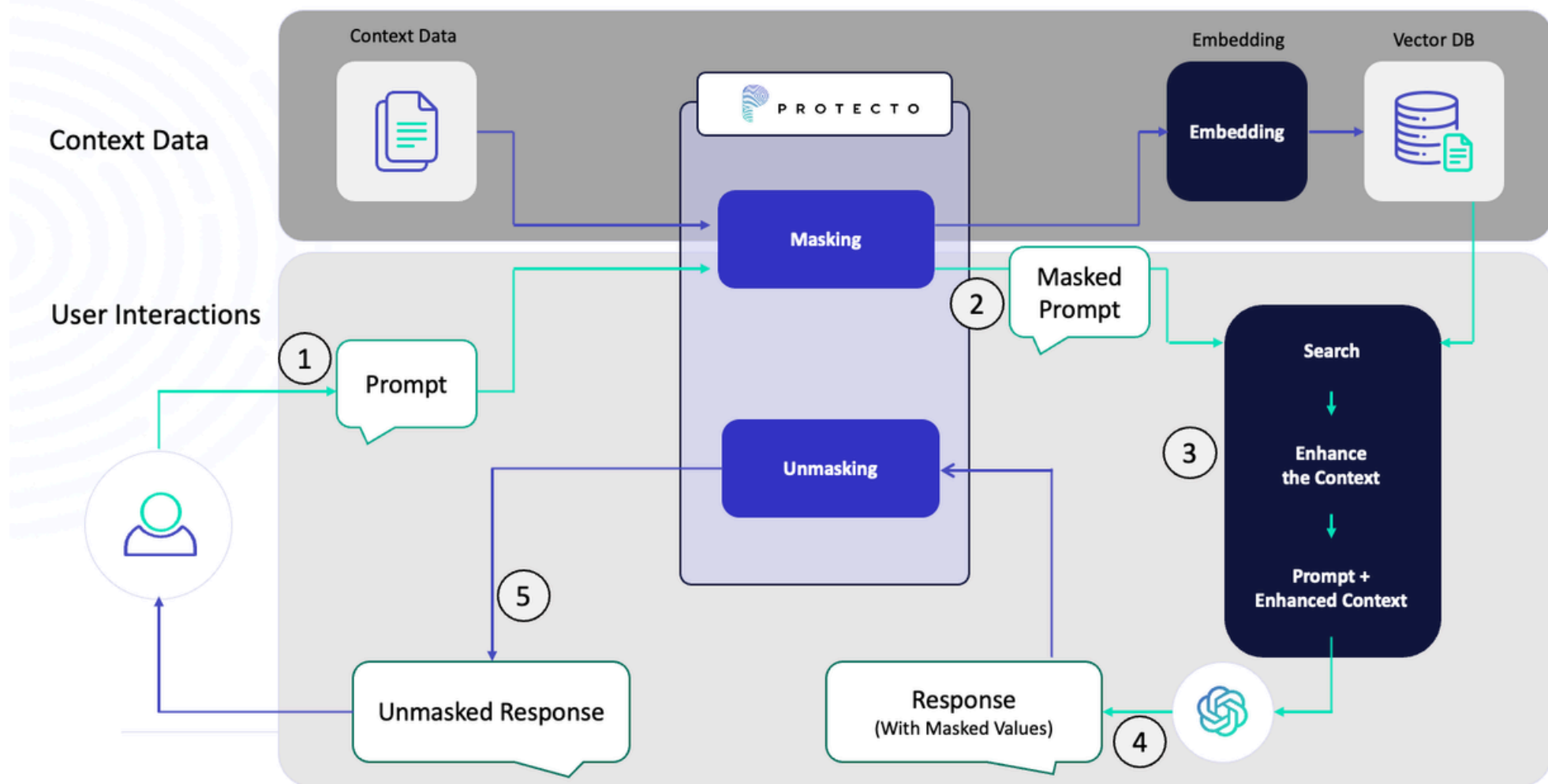
In terms of PII/PHI detection -

Protecto AI > Microsoft Presidio, AWS Comprehend

Let's try to understand the architecture of GPTGuard to get better clarity-

GPTGuard Architecture

GPTGuard Architecture



1. User Interaction Phase

- Employees submit queries or documents as they would with any AI tool
- The system automatically detects sensitive elements like: Customer identifiers, payment details(account number, CC, transaction ids), Medical/legal terms

2. Intelligent Protection Layer

- Advanced masking replaces sensitive data with tokens while preserving: Contextual meaning, Data formats, Relationship between elements

3. Secure Processing Stage

- Only masked information reaches the LLM (GPT-4, Claude, etc.)
- Vector database enhances responses using protected document extracts

4. Verified Output Delivery

- Authorized users receive complete answers with sensitive data restored
- Full audit capability shows what was protected and when.

GPTGuard can help prevent severe data leaks in major sectors. Lets see how-

Real World Example

According to Harmonic Security's report **8.5% of enterprise AI** queries contain **sensitive data** with top exposures of customer data(**45%**), employee records(**26%**) etc.

Tools like **GPTGuard** can help prevent this by:

Finance



- Automatic masking of account numbers, SSNs
- Audit-compliant document search

Health



- Real-time PHI tokenization
- HIPAA-safe chat with EHRs

Legal



- Privileged data stays on-premises
- Masked but accurate legal research

As AI becomes widely adopted across every sector, we must establish clear security measures to prevent data leaks. Remember:

The solution isn't saying no to AI - it's saying yes to safe AI.

Stay Ahead with Our Tech Newsletter! 🚀

👉 Join 1k+ leaders and professionals to stay ahead in GenAI!

<https://bhavishyapandit9.substack.com/>

Join our newsletter for:

- Step-by-step guides to mastering complex topics
- Industry trends & innovations delivered straight to your inbox
- Actionable tips to enhance your skills and stay competitive
- Insights on cutting-edge AI & software development

WTF In Tech

Home

Notes

Archive

About

People with no idea about AI
saying it will take over the world:

My Neural Network:



Object Detection with Large Vision Language Models (LVLMs)

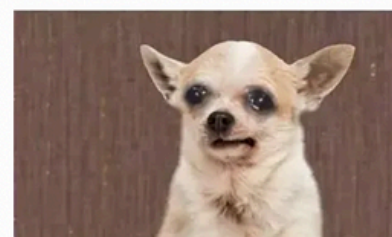
Object detection, now smarter with LVLMs

MAR 27 • BHAVISHYA PANDIT

AI Interview Playbook : Comprehensive guide to land an AI job in 2025

Brownie point: It includes 10 Key AI Interview Questions (With Answers).

MAR 22 • BHAVISHYA PANDIT



WTF In Tech

My personal Substack

💡 Whether you're a developer, researcher, or tech enthusiast, this newsletter is your shortcut to staying informed and ahead of the curve.



**Follow to stay updated on
Generative AI**



LIKE



COMMENT



REPOST