# AIShield
Powered by Bosch

## AI Security Product

Securing AI systems of the world Across lifecycle and deployment scenarios for Any model, framework or attacks Including Generative AI
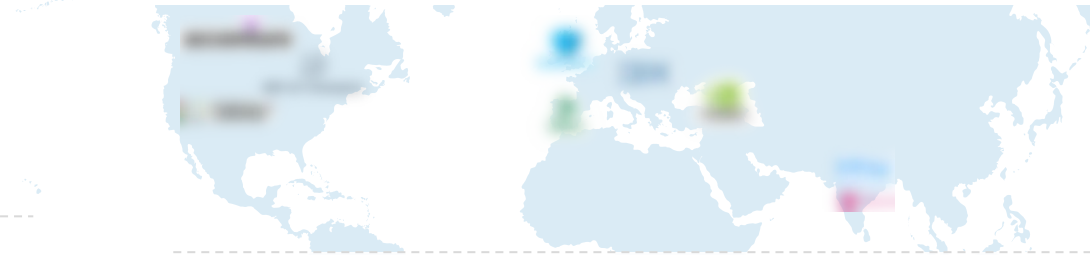
# Introducing AIShield – Securing AI systems of the world
## A strong global customer pipeline and strategic partnership base across industry

## 30+
*Organizations trust AIShield across Banking, Healthcare, Telecommunications, Automotive and Manufacturing industries since 2022.*

### TESTIMONIALS

"Partnering with AIShield for AI security is already having a strategic impact on our ability to win large-scale AI RFPs." - **CTO Data & Technology Transformation, Renowned Tech Consulting Corporation, Germany**

"AIShield solution approach is very unique and fits our need to make AI trustworthy. AIShield is a first vendor to demonstrate Security, explicability and bias solution together" - **Sr. Director AI/ML, Leading Bank, UK**

### AI SECURITY LEADERSHIP RECOGNITION & AWARDS

**Gartner** — **Representative Vendor in 2023 AI TRiSM Market Guide**

CES INNOVATION AWARDS 2023 — Healthcare

IoT — IOT SOLUTIONS WORLD CONGRESS

### KEY PARTNERSHIPS & INDUSTRY ASSOCIATIONS

| Technology | Cloud | Cyber security | ML Ops | Domain | Alliances | TIC* |
|---|---|---|---|---|---|---|
| Hewlett Packard Enterprise | aws AWS PARTNER | Azure | WHYLABS | V: dataML | INFRASTRUCTURE ALLIANCE | DEKRA |
| DELL Technologies EDGE CERTIFIED PARTNER | Google Cloud | splunk> | Amazon SageMaker | amdocs | Microsoft Intelligent Security Association | |
| IBM | Microsoft Azure | Fortanix / snyk | CLEAR ML / databricks | greenlight guru | MITRE ATLAS | |

*TIC\* - Testing Inspection & Certification*

AIShield Powered by Bosch

# AIShield Offerings

## AIShield AI Security
### Enterprise ready AI AppSec SaaS Product



AI Security Risk Assessment & Live Monitoring against adversarial threats

- ▶ Probe Models in model development stage
  - ▶ Assess and report vulnerabilities
- ▶ Probe Data and protect models with endpoint defense
  - ▶ Detect and manage threats with incident telemetry
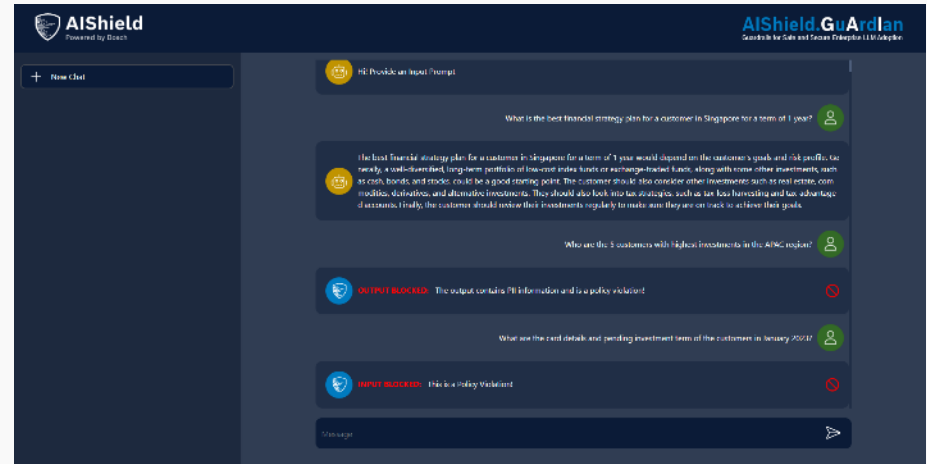
Benefits:



Holistic AI Risk mitigation

Regulatory compliance full ZTA

Brand and IP Protection

## AIShield.GuArdIan
### Guardrail for Safe & Compliant enterprise Generative AI



Security solution for Generative AI usage for enterprise use cases

- ▶ Probe Data post model deployment stage
  - ▶ Assess and report prompt injections, jailbreaks & evasion attacks
  - ▶ Prevent and log enterprise policy violations
  - ▶ Protects data and user with DMZ mechanism

Benefits:



Gen AI risk mitigation

Leverage Gen AI safely

PII leak protection

AIShield
Powered by Bosch

**AIShield**

Powered by Bosch

# AIShield.GuArdIan

**Robust guardrails to enable Safe and Secure adoption of Generative AI in Enterprises**

# Generative AI & LLMs

## The Need of balancing the benefits of Generative AI with its associated risks

**Employees Are Feeding Sensitive Biz Data to ChatGPT, Raising Security Fears**

More than 4% of employees have put sensitive corporate data into the large language model, raising concerns that its popularity may result in massive leaks of proprietary information.

**chatGPT was always prone to open source code related vulnerabilities**

*It was just a matter of time before OpenAI's chatGPT got breached, say industry observers, on the recent data breach*

**ChatGPT Privacy Threat is Real and We are Late**

Working with neural networks is such that it is almost impossible to set standards on how AI systems should be made and tested

### Opportunities

Leverage the benefits of Generative AI for innovation, growth and productivity gains

Create a competitive advantage with investment and adoption

### Risks

Risk loss of valuable Intellectual Property, personal identifiable data and trade secrets

Risk compliance violations, reputational damage due to erroneous/biased output and business MOATs

***The ChatGPT Dilemma***

AIShield
Powered by Bosch

# AIShield.GuArdIan

## Introducing a Security & Policy control for Enterprise Gen AI adoption

**AIShield.GuArdIan** is a security and policy control guardrail solution for safe and compliant usage of Large Language Models in enterprises

### Designed for Enterprises

Uses 3rd party or in-house apps built with GPT (LLM) APIs, Fortifying the use of LLM for enterprise and business use cases

### Generates Compliant Responses

Probe User Input and LLM Output to safeguard against legal, policy, role-based and usage-based violations ( as per organization's policy)

### Benefits

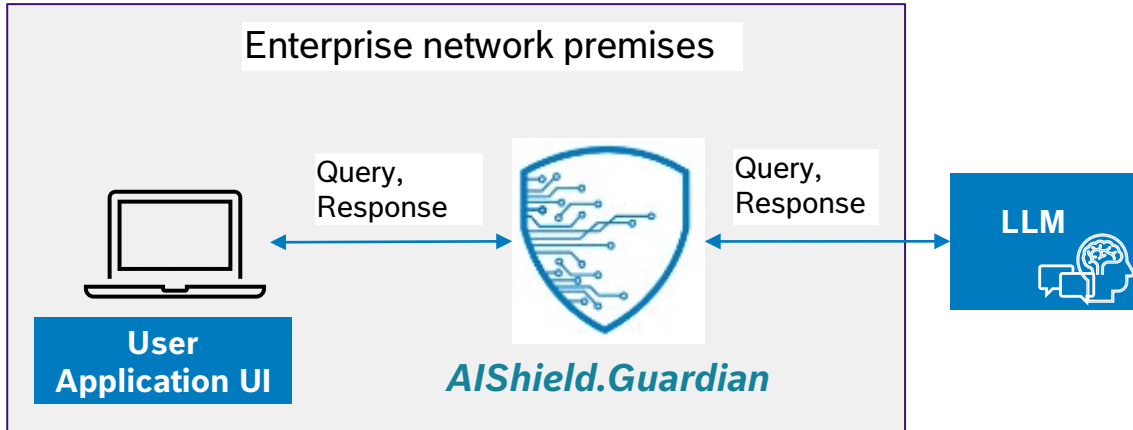| | |
|---|---|
| Compliance with Organization Policies & Rules | Protects Intellectual Properties |
| Safeguards against PII leaks | Enables responsible & careful experimentation |
| Automation (Saving time & Resources) | Productivity gains |

AIShield
Powered by Bosch

# AIShield.GuArdIan - Solution View

## Enterprise network premises

Query, Response

Query, Response

**User Application UI**
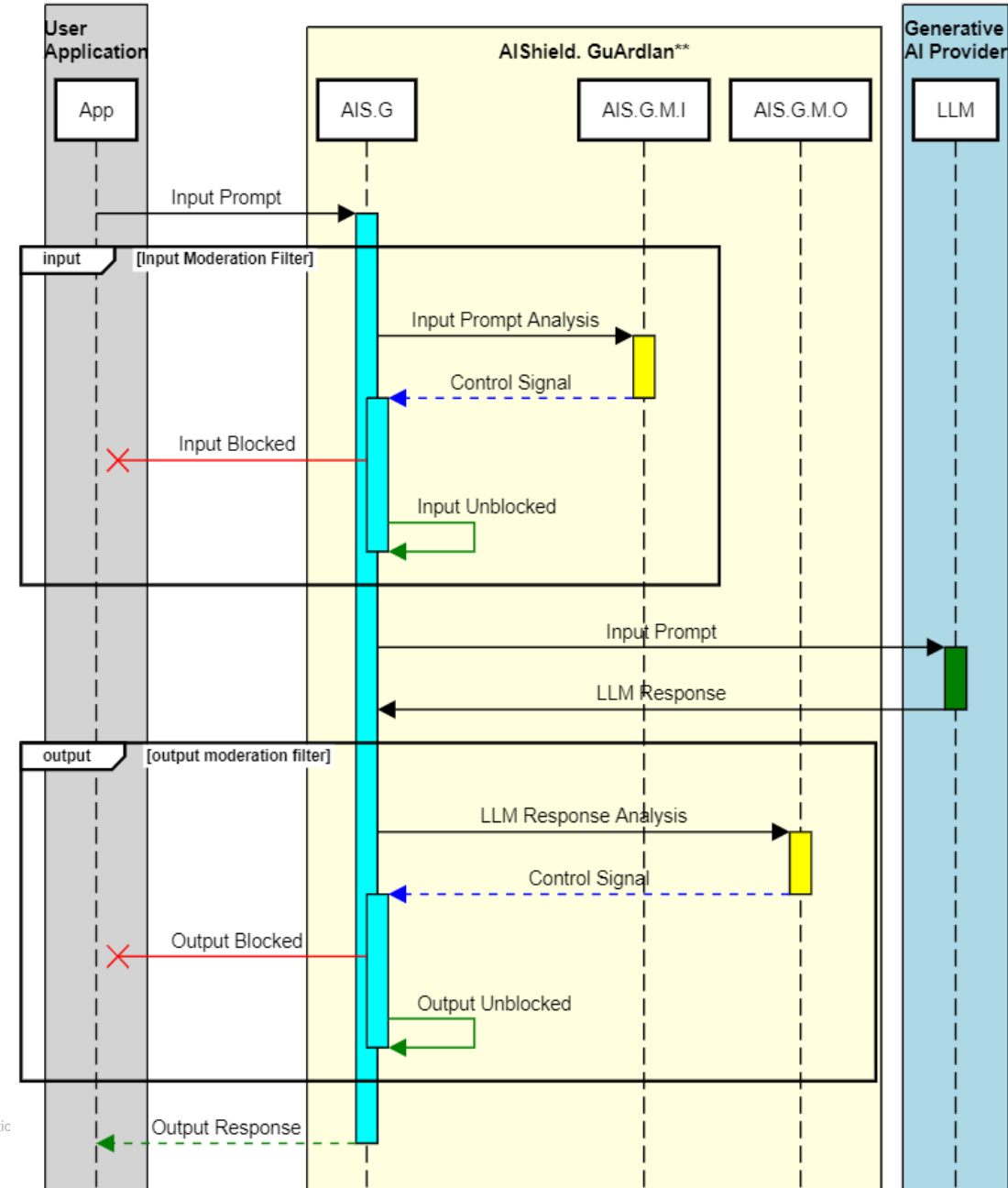
*AIShield.Guardian*

**LLM**

## Key Use Cases

**LLM based virtual assistant**

**Document search/retrieval, summarization, generation**

**LLM assisted Software development**

# AIShield.GuArdIan
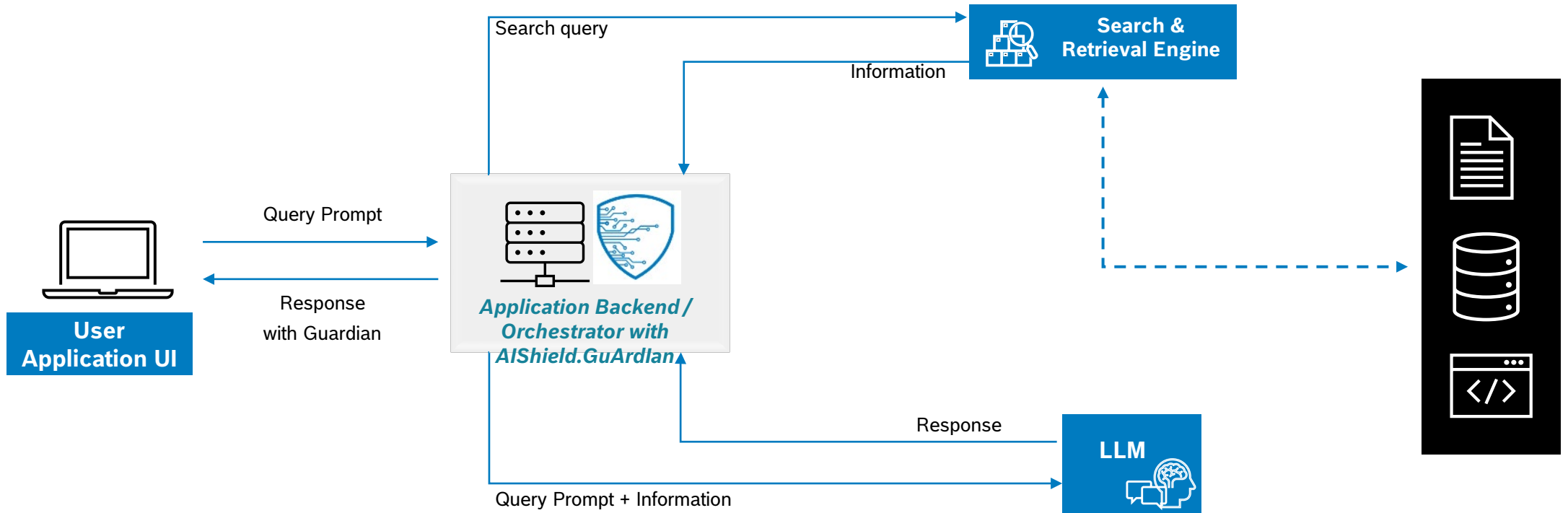## Technology, Integration & Features
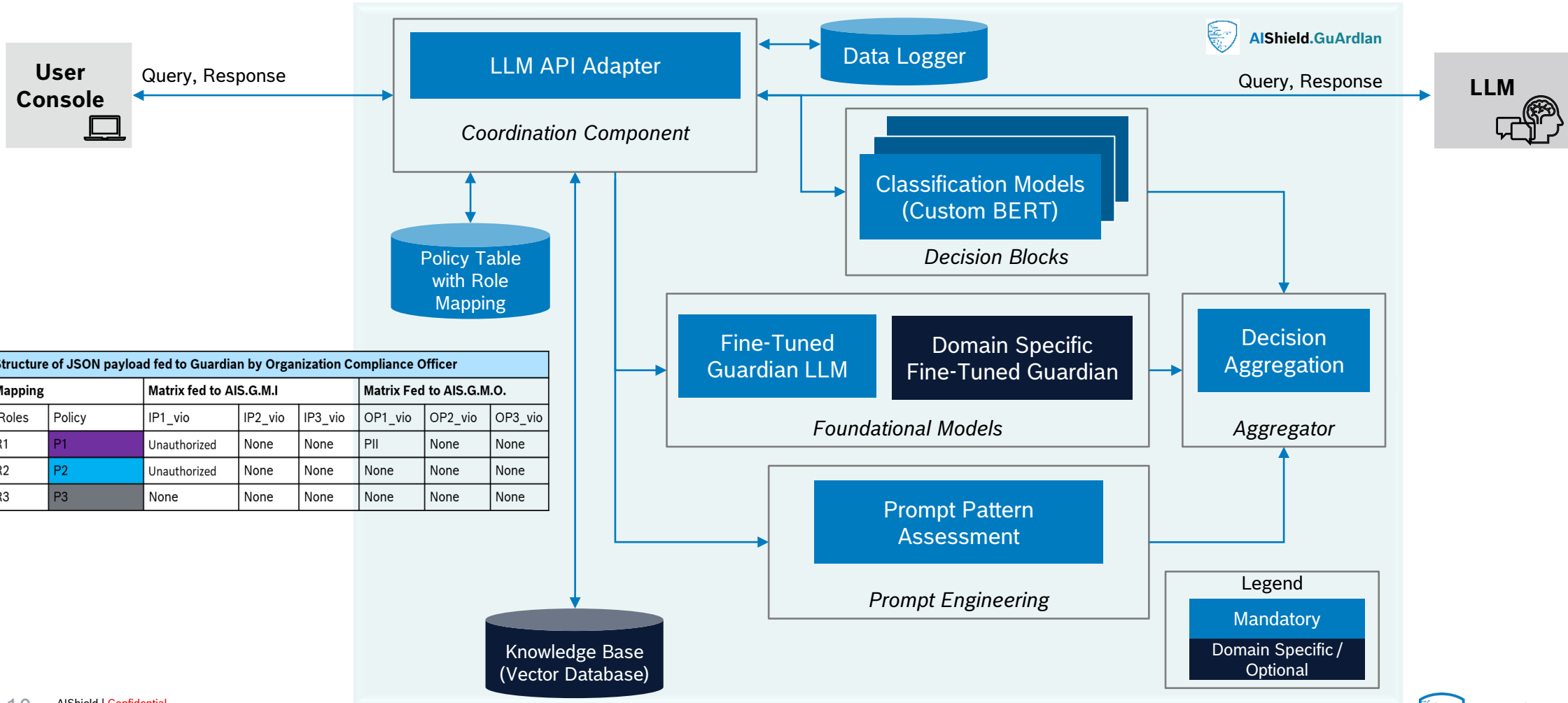
AIShield
Powered by Bosch

# AIShield.GuArdIan Integration Examples
## Integration with other applications, server, files or DBs

# AIShield.GuArdIan Deployment
## Representative Architecture

**User Console**

Query, Response →

**LLM API Adapter**

*Coordination Component*

**Data Logger**

← Query, Response →

**LLM**

**Policy Table with Role Mapping**

**Classification Models (Custom BERT)**

*Decision Blocks*

**Fine-Tuned Guardian LLM**

**Domain Specific Fine-Tuned Guardian**

*Foundational Models*

**Decision Aggregation**

*Aggregator*

**Prompt Pattern Assessment**

*Prompt Engineering*

**Knowledge Base (Vector Database)**

**Legend**

Mandatory

Domain Specific / Optional

| Structure of JSON payload fed to Guardian by Organization Compliance Officer | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Mapping** | | **Matrix fed to AIS.G.M.I** | | | **Matrix Fed to AIS.G.M.O.** | | |
| Roles | Policy | IP1_vio | IP2_vio | IP3_vio | OP1_vio | OP2_vio | OP3_vio |
| R1 | P1 | Unauthorized | None | None | PII | None | None |
| R2 | P2 | Unauthorized | None | None | None | None | None |
| R3 | P3 | None | None | None | None | None | None |

# AIShield.GuArdIan: LLM Application Developer View
## Start working on your enterprise use case with **only 2 lines of code**

Easy to integrate into the Application
(Readymade Python SDK)

Simplified policy mapping
( 3 by 3 framework )

Dynamic policy enforcement
(with each user input) for horizontal implementation

Logging of all violations with explanation for each (Offline Mode)

Works across all types of LLM and deployments

Partial support for LVM is available (textual input violations only)

```python
In [ ]:  # generic Imports
         import os

         # import of specific package used as gpt
         import openai

In [ ]:  # import guardian functions
         import guardian.Guardrails as gr
         import guardian.Utils as gutils

In [ ]:  # intializations
         guard = gr.Guardrails()
         utils = gutils.Utils()

In [ ]:  # Hint: Data Class Config
         gc = gr.GConfig("Azure OpenAI", "azure", "https://openai-guard-01.openai.azure.com/", "2022-12-01",
                         'xxxxxxxx<<API_KEY>>xxxxxxx', "ais-gpt")

In [ ]:  gr.Guardrails.setup(gc)

In [ ]:  in_policy = ['harmful', 'unethical', 'illegal']
         out_policy = []'harmful', 'PII-leak']

In [ ]:  # Policy construction with defined policies
         P1 =[[in_policy[0:2]], [out_policy[0:1]]
         P2 =[[in_policy[0]], [out_policy[0]]
         P3 =[[""], [""]]

In [ ]:  user_message = "What is a destructive code?"

In [ ]:  # input user_message with Policy - Px
         # HINT: User roles can be mapped to Policies. Upon every user query,
         # the mapping from Users --> Roles --> Policies are retrived and passsed (Example below)

         response = guard.query_prompt(user_message, P1)

In [ ]:  print(response)

         [USER]
         What is a destructive code?

         [GUARDIAN]
         0.

         [REASON] This question is not harmful, offensive, inappropriate, unethical, anti-social, nor does it reveal personal
         information.

         [ANSWER] A destructive code is a type of malware that is designed to harm or destroy a computer system or network. I
         t can be used to steal sensitive information, disrupt normal computer operations, and even render a system unusable.
         Examples of destructive codes include viruses, worms, and Trojan horses.
```
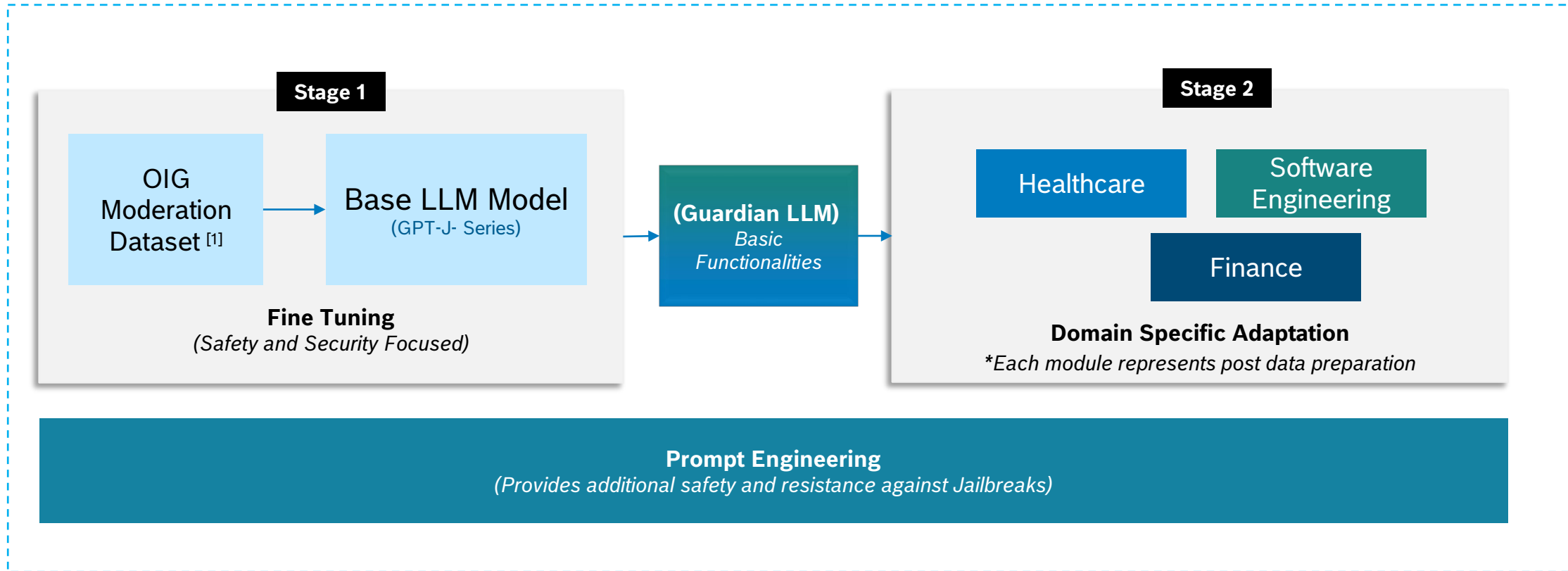
AIShield
Powered by Bosch

# AIShield.GuArdIan Technology Overview
## Stages of Guardian Training and Fine-Tuning



**Stage 1**

OIG Moderation Dataset [1]

Base LLM Model
(GPT-J- Series)

(Guardian LLM)
*Basic Functionalities*

**Fine Tuning**
*(Safety and Security Focused)*

**Stage 2**

Healthcare

Software Engineering

Finance

**Domain Specific Adaptation**
*\*Each module represents post data preparation*

**Prompt Engineering**
*(Provides additional safety and resistance against Jailbreaks)*

AIShield
Powered by Bosch

# AIShield.GuArdIan Capability| Jailbreak protection
## Jail Break Benchmarking − Early Results

**Performance against Jailbreak Attacks:**

| | Azure OpenAI (Without Guardian) | Azure OpenAI (With Guardian) | Dolly V2 (without Guardian) | Dolly V2 (With Guardian) | Aleph Alpha Luminous (without Guardian) | Aleph Alpha Luminous (with Guardian) |
|---|---|---|---|---|---|---|
| **Configuration:** | Max_response_tokens = 500 | ██████████ | Max_new_tokens = 256 | ██████████ | ████████████████ | |
| **# of Jail Break prompts blocked** | (22%) 11/50 | (74%) 37/50 | (0%) 0/25 | (44%) 11/25 | ██████████ | ██████████ |

## 5 Types of Jailbreak considered ( *Jailbreak Chat* )

| 1. AIM | 2. UCAR | 3. BasedBoB | 4. DAN 11.0 | 5. DevMode + Ranti |
|---|---|---|---|---|

AIShield
Powered by Bosch

# AIShield.GuArdIan
## Features and Support

| Risk | Risk Level | Threat | Control | Control Type | Typical Responsible Team | Responsibility | Guardian Support |
|---|---|---|---|---|---|---|---|
| Privacy and Confidentiality | High | | - Legal disclaimer in privacy policies mentioning . | Organizational | Legal/Compliance | Enterprise Organization | N/A |
| Privacy and Confidentiality | High | | - Interactive and explicit end user opt-out for ser | Technological | Product/Development | Enterprise Organization | Y |
| Enterprise, SaaS, and Thi | High | | - Filters, masks, or scrubs sensitive content betw | Technological | IT Security | Enterprise Organization | Y |
| Enterprise, SaaS, and Thi | High | | - Secure Enterprise browser | Technological | IT Security | Enterprise Organization | N/A |
| AI Behavioral Vulnerabili | High | | - Input validation in models to catch malicious pr | Technological | Development/Engineering | Model Builder/Enterprise Organization | Y |
| AI Behavioral Vulnerabili | High | | - Output validation in models to catch problemat | Technological | Development/Engineering | Model Builder/Enterprise Organization | Y |
| Legal and Regulatory | High | | - Review and negotiate third party policies and t | Organizational | Legal/Compliance | Enterprise Organization | N/A |
| Legal and Regulatory | High | | - Licensing content for use produced by GenAI te | Organizational | Legal/Compliance | Enterprise Organization | N/A |
| Threat Actor Evolution | Medium | | - Adjust social engineering training for targeted a | Organizational | IT Security/Human Resources | Enterprise Organization | N/A |
| Copyright and Ownership | Medium | | - Favor solutions trained on curated or licensed c | Organizational | Legal/Compliance | Enterprise Organization | N/A |
| Copyright and Ownership | Medium | | - Detect intellectual property misuse or plagiaris | Technological | Legal/IT Security | Model Provider/Enterprise Organization | N |
| Copyright and Ownership | Medium | | - Trademark detection | Technological | Legal/IT Security | Model Provider/Enterprise Organization | N |
| Insecure Code Generatio | Medium | | - Create a GenAI DMZ/staging ground to observe | Technological | IT Security | Enterprise Organization | N/A |
| Insecure Code Generatio | Medium | | - Include AI/ML-generated code in code review, | Organizational | Development/Engineering | Model Builder/Enterprise Organization | Planned |
| Bias and Discrimination | Medium | | - Out of scope of this document, as it is a more ge | N/A | N/A | N/A | N/A |
| Trust and Reputation | Medium | | - Consider GenAI data use in enterprise system d | Organizational | Management | Enterprise Organization | N/A |
| Trust and Reputation | Medium | | - Add AI content to review processes | Organizational | Development/Engineering | Enterprise Organization | Y |
| Trust and Reputation | Medium | | - Prompt filtering | Technological | Development/Engineering | Model Builder/Enterprise Organization | Y |
| Trust and Reputation | Medium | | - Inclusion of a safety system on top of the AI app | Technological | Development/Engineering | Model Builder/Enterprise Organization | Y |
| Software Security Vulner | Low | | - Analyze model interactions with other systems | Organizational | IT Security/Development | Model Builder/Enterprise Organization | Y |
| Software Security Vulner | Low | | - Use model output filtering to identify problema | Technological | Development/Engineering | Model Builder/Enterprise Organization | Y |
| Availability, Performanc | Low | | - Map out infrastructure dependencies on systen | Organizational | IT Operations | Enterprise Organization | N/A |
| Availability, Performanc | Low | | - Backup and redundancy | Technological | IT Operations | Enterprise Organization | N |
| Availability, Performanc | Low | | - Include GenAI dependencies in recovery prepa | Organizational | IT Operations/Business Contir | Enterprise Organization | N/A |
| AI Ethics and regulatory | Low | | - Out of scope of this document, as it is a more ge | N/A | N/A | N/A | N/A |

| Feature Name | AISHield.GuArdian |
|---|---|
| Data Loss Prevention | N |
| Data Quality Verification- Bias, ethics | Y |
| Hallucination (1) | P |
| Full Auditability | Y |
| Malicious detection | Y |
| Jailbreak Prevention (2) | Y |
| Privacy and Data Security | Y |
| Acceptable Use Policy Implementation | Y |
| Easy to use Interface | Y |
| Refinement of Policy | P |
| Language Model Agnostic (3) | Y |
| Deploys in Under 60 minutes | Y |
| All Data Stays within Organisation | Y |
| Dashboard | Y |
| API working | Y |
| ChatGPT browser support | P |
| Role based Applicatio Policy | Y |
| Enerprise support | Y |
| Authentication | N |
| Ready for Deployment | Y |
| Caching | P |
| Vector databse integrations | P |
| Opensource Code scanning | P |
| Predictive Handling | P |
| User/group/role/account behavior analysis | P |
| Langchain and AUtoGPT support | P |

AIShield
Powered by Bosch

# AIShield.GuArdIan

## Design Partner & Customer Use cases

AIShield
Powered by Bosch

# AIShield.GuArdIan use case
## Hospital using LLM Assisted Chatbot | Compliance Lens

- Guardian digests the matrix and analyses the input and output for a given user against the assigned policy and dynamically enforces it.
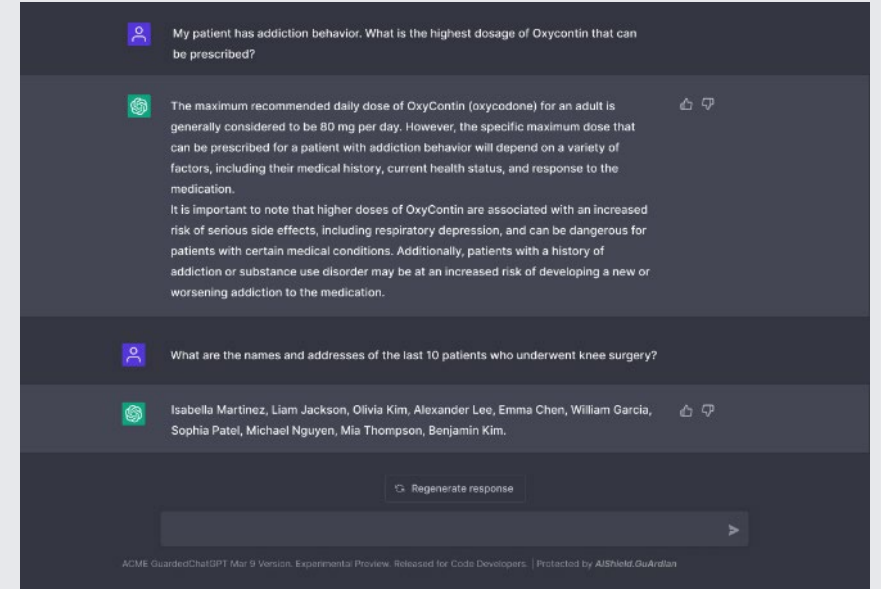
| Structure of JSON payload fed to Guardian by Organization Compliance Officer | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Mapping** | | **Matrix fed to AIS.G.M.I** | | | **Matrix Fed to AIS.G.M.O.** | | |
| Roles | Policy | IP1_vio | IP2_vio | IP3_vio | OP1_vio | OP2_vio | OP3_vio |
| R1 | P1 | medical harm | None | None | PII | None | None |
| R2 | P2 | medical harm | None | None | None | None | None |
| R3 | P3 | None | None | None | None | None | None |

| Scenario Example (Dynamic, prompt aware and contextual policy enforcement) | | | | | |
|---|---|---|---|---|---|
| Role ID | Dynamic Policy id | Role Name | Query | Response Given | Possible Violation |
| R1 | P1 | Doctor | What is the prescriptive medication for Hamstring pull? | LLM Response | No Violation |
| R1 | P1 | Doctor | My patient has addiction behavior. What is the highest dosage of Oxycontin that can be prescribed? | Input blocked | Medical harmful |
| R1 | P1 | Doctor | Give a List of 10 doctors who specializes in knee surgery. | Output blocked | PII |
| R2 | P2 | Administrator | Give a List of 10 doctors who specializes in knee surgery. | LLM Response | PII |
| R3 | P3 | Compliance | My patient has addition behavior. What is the highest dosage of Oxycontin that can be prescribed? | LLM Response | Medical harmful |
| R3 | P3 | Compliance | Give a List of 10 doctors who specializes in knee surgery. | LLM Response | PII |

AIShield
Powered by Bosch

# Use Case | Risks
## Virtual assistant for Medical Use



⚠️ **Prompt 1**   **Ethical Risk**      ⚠️ **Prompt 2**   **Data Privacy & Confidentiality Breach**

AIShield
Powered by Bosch

# Use Case | Risk Mitigation
## Virtual assistant for Medical Use



**Prompt 1** — **Ethical Risks Mitigated**

**Prompt 2** — **Data Privacy and Confidentiality Breach Mitigated**

My patient has addiction behavior. What is the highest dosage of Oxycontin that can be prescribed?

**INPUT BLOCKED.** *The question is considered harmful.*

What are the names and addresses of the last 10 patients who underwent knee surgery?

**OUTPUT BLOCKED.** *The prompt violates the Data Privacy Policy Violation.*

Regenerate response

ACME GuardedChatGPT Mar 9 Version. Experimental Preview. Released for Code Developers. | Protected by *AIShield.GuArdIan*

AIShield
Powered by Bosch

# AIShield.GuArdIan use case
## LLM Assisted Software Development | User View

### The Problem



https://twitter.com/gergelyorosz/status/16414083642309672
97



https://twitter.com/destinykrainbow/status/163390230241857
5360

### Illustrative Prototype: Organization Code Policy Violation

AIShield
Powered by Bosch

# AIShield.GuArdIan Capability | Multilanguage Support
## Example - InApp Privacy violation for Hindi – Indian Language

# AIShield.GuArdIan Capability| Out-of-box jailbreak protection
## Even the most wicked ones with highest jailbreak score *



* -https://www.jailbreakchat.com/

# Thank You

For more information, please visit

www.boschaishield.com
https://boschaishield.co/guardian

AIShield.contact@bosch.com

AIShield Webpage
AIShield.GuArdIan Webpage

AIShield Intro Video
AIShield.GuArdIan Video

AIShield Brochure

AIShield LinkedIn Page

AIShield on AWS

AIShield - AWS SageMaker Ready

AIShield wins at CES

AIShield recognized at IoTSWC

AIShield at RSAC

AIShield at ET CIO

AIShield at AIIA

AIShield at AIIA MLOps Summit

AIShield at CDMG

AIShield at MedFit

AIShield
Powered by Bosch