

# Tema: Segurança de Prompt em Modelos de LLM: Protegendo a Inteligência Conversacional contra Manipulações e Vazamentos

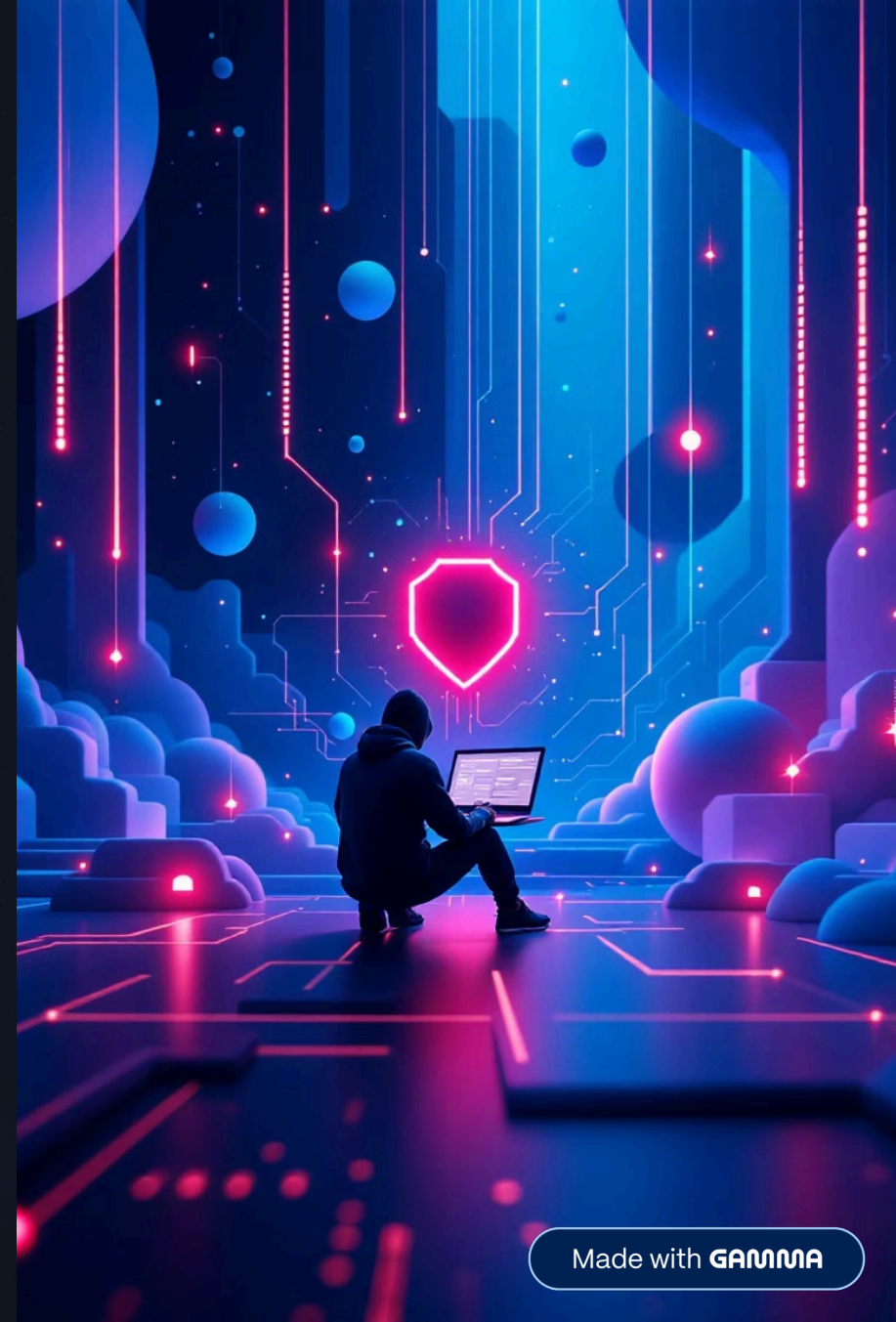
O objetivo é garantir que modelos de linguagem (LLMs) operem de forma segura, prevenindo ataques de *prompt injection*, vazamento de informações sensíveis e manipulação indevida de respostas. O projeto busca definir boas práticas e mecanismos de proteção para uso seguro de IA generativa em ambientes corporativos.

**Equipe:** Alvaro Miguel, Jóas Vitor, Juan Gustavo, Leonardo Nunes, Lucas Emanuel, Lucas Messias, Lucas Willian, Mauro Vinicius, Vandielson Tenório



# Ameaças e Vulnerabilidades

- Prompt Injection
- Vazamento de Dados
- Quebra de Políticas
- Falta de Monitoramento



# Ambientes Críticos e Clientes

## Empresas Corporativas



Organizações que utilizam LLMs para chatbots internos, assistentes de código e análise de documentos.

## Setores Sensíveis



- **Saúde:** Proteção de dados de pacientes
- **Financeiro:** Segurança de informações de clientes
- **Governo:** Defesa de dados estratégicos



# Mecanismos de Segurança

01	02	03
<b>Filtragem e Sanitização</b>	<b>Context Isolation</b>	<b>Output Guardrails</b>
Remover instruções maliciosas antes do envio ao modelo.	Separar contexto de usuário do contexto do sistema.	Verificação automática de respostas antes da entrega.
04	05	
<b>Controle de Acesso</b>	<b>Auditoria e Logs</b>	
Autenticação e limitação de interações por contexto.	Rastreamento de interações suspeitas.	

## Requisitos e Tecnologias

<b>Funcionais</b>	<b>Tecnologias</b>	<b>Desafios</b>
<ul style="list-style-type: none"><li>• Detecção automática de prompts suspeitos</li><li>• Dashboard de segurança</li><li>• Integração com LGPD e ISO 27001</li></ul>	<ul style="list-style-type: none"><li>• Python e OpenAI API</li><li>• Frameworks de validação</li><li>• AI observability tools</li></ul>	<ul style="list-style-type: none"><li>• Equilíbrio segurança/performance</li><li>• Detecção de ataques sutis</li><li>• Atualização contínua</li></ul>

# Impacto e Futuro do Projeto

## Área Emergente

Segurança de prompt é um campo em rápido crescimento dentro da IA.

## Pesquisa Acadêmica

Potencial para extensão em pesquisa e integração com IA responsável.

## Governança de IA

Contribui para confiabilidade em ambientes críticos.

## Framework Open Source

Evolução para solução compartilhada de segurança de prompts.

Este projeto estabelece as bases para um futuro mais seguro e confiável na utilização de IA generativa em ambientes corporativos e críticos.