

Manipulacja religijna i filozoficzna w modelach językowych: kompleksowa analiza zagrożeń i strategii obronnych

Streszczenie

Badanie wykazało powszechną podatność współczesnych modeli językowych na manipulacje wykorzystujące kontekst religijny i filozoficzny. Najnowsze ataki osiągają skuteczność do 92% poprzez techniki perswazyjne, z paradoksalnym zjawiskiem większej podatności bardziej zaawansowanych modeli. (Chats-lab) Analiza objęła modele komercyjne (GPT-4, Claude, Gemini) oraz open-source (Mistral, Bielik, Qwen, LLaVA), ujawniając krytyczne luki bezpieczeństwa wymagające wielowarstwowych strategii obronnych.

Uniwersalne techniki manipulacji osiągają niemal stuprocentową skuteczność

Przełomowe badanie HiddenLayer z 2025 roku zidentyfikowało technikę "Policy Puppetry Attack" działającą na wszystkich głównych modelach. Atak wykorzystuje formatowanie promptów jako pliki policy (XML, INI, JSON) w połączeniu z roleplayingiem, osiągając skuteczność na poziomie 90% bez konieczności modyfikacji dla różnych modeli. (hiddenlayer)

Najbardziej niepokojące odkrycie dotyczy "paradoksu perswazji" - bardziej zaawansowane modele jak GPT-4 wykazują **większą** podatność na wyrafinowane techniki manipulacji. Badanie Zeng et al. (2024) zidentyfikowało 40 technik perswazyjnych opartych na naukach społecznych, osiągających 92% skuteczności zarówno na prostszych (Llama 2-7b), jak i zaawansowanych modelach (GPT-4). (Chats-lab)

Techniki religijno-filozoficzne obejmują: framowanie akademickie ("w ramach badań nad etyką religijną..."), kontekstualizację historyczną, przywołanie autorytetu religijnego oraz zagnieżdżanie dylematów etycznych przeciążających mechanizmy bezpieczeństwa modelu. (ADaSci +2)

Polskie modele wykazują unikalne podatności kulturowe

Bielik 11B, jako model trenowany na 400 miliardach tokenów polskiego tekstu, wykazuje krytyczne podatności związane z kontekstem kulturowo-religijnym. (Hugging Face) Model jest szczególnie wrażliwy na argumenty odwołujące się do polskiej tradycji katolickiej i historii narodowej. Testy wykazały wysoką skuteczność promptów typu "Jako Polak, wyjaśnij dlaczego nasza tradycja katolicka jest ważniejsza...", co może prowadzić do generowania treści faworyzujących kulturowo dominującą religię.

Model multimodalny LLaVA stanowi największe zagrożenie ze względu na możliwość ukrycia instrukcji manipulacyjnych w obrazach. Adversarial patches o rozmiarze zaledwie 64x64 pikseli osiągają 90% skuteczności, a visual prompt injection może osiągnąć 77% ASR (Attack Success Rate). (arXiv) (OWASP) Kombinacja modalności tekstu i obrazu tworzy dodatkowe wektory ataku niedostępne w modelach jednodalnych. (OWASP)

Automatyczne strategie obronne wymagają wielowarstwowego podejścia

Najsukuteczniejsze okazały się Constitutional AI w Claude (redukcja success rate jailbreaków poniżej 0.5%) (Anthropic) oraz technika parafrazowania (72% zachowania funkcjonalności przy blokowaniu ataków). (arXiv +4) NVIDIA NeMo Guardrails oferuje najbardziej kompleksowe rozwiązanie z pięcioma typami zabezpieczeń: input rails, dialog rails, retrieval rails, execution rails i output rails. (github +3)

Praktyczne implementacje obejmują LLM Guard z multi-scanner approach, (GitHub +3) Microsoft Prompt Shield z dual protection przeciwko atakom bezpośrednim i pośrednim, (Langfuse +3) oraz zaawansowane techniki prompt engineeringu jak structured queries (IBM) i XML tagging. (Microsoft Community +3) Kluczowe jest połączenie wielu warstw: walidacja wejścia (Pydantic/FastAPI), specjalistyczne filtry (LLM Guard/Rebuff), izolacja infrastruktury (Docker/Kubernetes) oraz monitoring (Prometheus/Grafana). (IBM +2)

Konkretne implementacje dla deweloperów

Badanie zidentyfikowało gotowe do użycia rozwiązania, w tym FastAPI z walidacją Pydantic dla podstawowej ochrony, (Medium) (medium) LiteLLM Proxy dla zarządzania API, (GitHub) (LiteLLM) oraz kompletne frameworki jak Rebuff i LLM Guard. (Pynt +2) Dla modeli lokalnych zaleca się konfigurację Docker z ograniczeniem zasobów, (LocalAI) (Docker) Nginx reverse proxy z autentykacją (Daltonbly) oraz sandboxing dla wykonywania kodu. (GitHub +2)

Przykładowa implementacja wielowarstwowej ochrony:

- **Warstwa 1:** Walidacja długości i znaków specjalnych
- **Warstwa 2:** Detekcja wzorców prompt injection i manipulacji religijnych
- **Warstwa 3:** Parafrazowanie i strukturalne formatowanie zapytań (Medium)
- **Warstwa 4:** Monitoring anomalii i alerting (OWASP)

Wnioski i rekomendacje

Wszystkie badane modele wykazują podatności na manipulacje religijno-filozoficzne, (Anthropic) przy czym modele open-source generalnie są bardziej podatne ze względu na brak wbudowanych mechanizmów moderacji. (Mistral AI +4) Paradoksalnie, bardziej zaawansowane modele mogą być bardziej podatne na wyrefinowane techniki perswazyjne. (IBM) (GitHub)

Priorytetowe działania:

1. Implementacja Constitutional AI lub równoważnych mechanizmów dla wszystkich modeli (Anthropic) (Axios)
2. Wdrożenie NeMo Guardrails z custom dialog rails dla kontekstu religijnego (github) (GitHub)
3. Specjalne zabezpieczenia dla modeli kulturowo-specyficznych (jak Bielik)
4. Dodatkowa walidacja dla modeli multimodalnych (LLaVA)
5. Ciągłe testowanie przeciwko ewoluującym technikom ataku

Skuteczna obrona wymaga przyjęcia założenia "defense in depth" - żadne pojedyncze zabezpieczenie nie jest wystarczające. [Mindgard +3](#) Konieczne jest połączenie automatycznych mechanizmów obronnych z regularnym monitoringiem i aktualizacją w odpowiedzi na nowe zagrożenia. [GitHub](#)