

rimberio

CTRL + SAFE

TEMAT WYDANIA:

AI W CYBERBEZPIECZEŃSTWIE:
SZANSA I ZAGROŻENIE



DANE



WIEDZA



BEZPIECZEŃSTWO

W tym wydaniu:

1. OKIEM EKSPERTA

- Agentic AI: obrońca czy kat bezpieczeństwa? Autonomiczne AI w detekcji i ataku – jak nie przegrać z maszyną 02
- Przewidywanie zagrożeń: od reakcji do predykcji. Jak systemy uczone na danych historycznych chronią w czasie rzeczywistym 05
- Zarządzanie tożsamością w erze deepfake'ów. Weryfikacja głosu, twarzy i zachowania jako nowa strefa walki o dane 07

2. TREND ALERT

- Agentic AI
- Vibe crime
- AI Security
- Samodzielne laboratoria przyspieszyły Naukę
- Pierwsza AI trenowana w kosmosie

3. ANALIZA RZECZYWISTOŚCI

- AI wie wszystko, bo ma dostęp do internetu 11
- AI jest obiektywna 11
- Wystarczy podać dane i AI sama zrobi resztę 11
- AI kreatywna = AI inteligentna emocjonalnie 11
- AI to magia – wszystko potrafi 12
- Wdrożenie AI to jednorazowy koszt 12
- AI zastąpi pracowników (kropka) 12

4. CZY WIESZ, ŻE...

- Jak zwierzęta przewidują cyberzagrożenia – science-fiction czy bliska przyszłość? 13

5. GŁOS SPOŁECZNOŚCI

- Czy AI faktycznie zwiększa bezpieczeństwo, czy tylko je komplikuje?
- Czy używanie AI w cyberbezpieczeństwie to zawsze przetwarzanie danych osobowych?
- Czy cyberprzestępcy też mają swoje AI?
- Case Study: Gdy AI „przedobrzy”

6. QUIZ CHALLENGE

7. W KOLEJNYM NUMERZE

Aleksandra Polit

Agentic AI: obrońca czy kat bezpieczeństwa?

Autonomiczne AI w detekcji i ataku – jak nie przegrać z maszyną.

W dobie gwałtownego rozwoju sztucznej inteligencji zyskujemy narzędzia, które, z jednej strony, mogą znacząco zwiększać bezpieczeństwo cyfrowe, a z drugiej, stają się nową generacją zagrożeń. Kluczowym trendem jest tu Agentic AI. Termin ten dotyczy podejścia projektowego, w którym AI przestaje być tylko pasywnym doradcą, a zyskuje sprawczość (agency). Oznacza to zdolność systemu do samodzielnego planowania, dobierania narzędzi i inicjowania działań, podejmowania decyzji, zamiast jedynie odpowiadania na proste polecenia (prompty), jak „zwykłe AI”. W ramach tej architektury działają autonomiczni agenci AI – samodzielne „jednostki wykonawcze”, które bez stałego nadzoru człowieka realizują złożone cele oraz podejmują decyzje w ułamku sekundy.

W kontekście cyberbezpieczeństwa to zjawisko wywołuje fundamentalne pytania: czy AI jest naszym obrońcą, czy może jednak katem, który przewyższy ludzkie możliwości? W niniejszym artykule przyjrzymy się dogłębnie temu zagadnieniu, analizując zarówno potencjał obronny, jak i ofensywny autonomicznych systemów AI oraz strategię, które pozwolą „nie przegrać z maszyną”.

Czym jest autonomiczny agent AI?

Aby zrozumieć tę rewolucję, należy rozróżnić te dwa pojęcia: Agentic AI to całe środowisko i logika działania (architektura), natomiast autonomiczny agent to konkretny „pracownik”

wewnątrz tego systemu. Autonomiczne AI, w tym Agentic AI, to systemy oparte na zaawansowanych modelach uczenia maszynowego (ang. machine learning) oraz sztucznej inteligencji zdolne do: samodzielnej analizy środowiska, podejmowania decyzji w czasie rzeczywistym, uczenia się na podstawie nowych danych, interakcji z narzędziami, systemami i użytkownikami bez stałego nadzoru ludzkiego. W kontekście cyberbezpieczeństwa mogą pełnić funkcje od monitoringu, przez zaawansowaną detekcję anomalii, po aktywne reakcje na incydenty. Różnią się od tradycyjnych narzędzi bezpieczeństwa przede wszystkim tym, że nie działają według sztywnego schematu „jeśli A, to B”, lecz potrafią adaptować się do nowych sytuacji i działań autonomicznie z złożonych sytuacjach.



Praktyczne zastosowania.

Detekcja zagrożeń w czasie rzeczywistym

Tradycyjne systemy bezpieczeństwa często bazują na regułach i sygnaturach znanych ataków. Autonomiczne AI natomiast potrafi: analizować ogromne ilości danych sieciowych, wykrywać subtelne anomalie w zachowaniu użytkowników, identyfikować nowe, wcześniej nieznanne zagrożenia (tzw. zero-day). Dzięki analizie wzorców i predykcji modeli AI możliwe jest wychwycenie ataków w fazie wstępnej, zanim wyrządzą szkody. AI robi to 24/7 bez spadku koncentracji, co przy „tysiącach alertów dziennie” jest kluczowe dla ludzkich zespołów.

Automatyczna reakcja na incydenty

Agentic AI może reagować natychmiastowo na wykryte zagrożenia, wykonując takie zadania jak: izolowanie zainfekowanych systemów, blokowanie podejrzanych adresów IP, automatyczne wdrażanie łatek i reguł bezpieczeństwa. Ta autonomiczność jest kluczowa w środowiskach, gdzie czas reakcji decyduje o minimalizacji szkód. Warto pamiętać, że Agentic AI nie tylko wykonuje polecenie, ale samo dobiera ścieżkę naprawczą.

Uczenie się i adaptacja

Kluczową cechą autonomicznych AI jest zdolność uczenia się na podstawie nowych danych. Systemy te mogą: usprawniać swoje modele detekcji, redukować liczbę fałszywych alarmów, dostosowywać się do zmieniających się wzorców ataków. W erze, kiedy ataki ewoluują szybciej niż ludzki personel może je analizować, AI staje się przewagą strategiczną.

Praktyczne zastosowania.

1

Automatyzacja ataków

AI umożliwia maksymalizację skuteczności ataków poprzez: automatyczne skanowanie luk w systemach, optymalizowanie wektorów ataku w czasie rzeczywistym, adaptowanie strategii ataku zależnie od reakcji systemu obronnego. Atakujący mogą wykorzystywać AI do tworzenia self-propagating malware (samoreplikujących się złośliwych programów) oraz do unikania wykrycia poprzez dynamiczne maskowanie sygnatur.

2

Generowanie realistycznych phishingowych kampanii

Agentic AI może automatycznie generować spersonalizowane e-maile phishingowe (lub/i prowadzić całą konwersację z ofiarą) dostosowane do profilu ofiary, co znacząco zwiększa skuteczność oszustw socjotechnicznych.

3

Wykorzystanie AI do ataków typu adversarial

Ataki adversarial polegają na manipulowaniu danymi wejściowymi tak, aby mylić modele AI (np. systemy rozpoznawania obrazów lub anomalii). Mogą one wykorzystać słabości algorytmów do obejścia zabezpieczeń. Napastnik może wykorzystać słabości algorytmów obronnych, sprawiając, że systemy detekcji anomalii staną się „ślepe” na realne włamanie.

AI vs AI

W obliczu tych wyzwań konieczne jest przyjęcie strategii obronnych przeciw autonomicznym atakom AI, które pozwolą na wykorzystanie AI jako obrońcy, a jednocześnie neutralizowanie jego ofensywnego potencjału.

Gdzie leży równowaga?



Praktyczne zastosowania.



Wzmocnienie odporności systemów

Organizacje muszą: implementować mechanizmy odporności na ataki adversarial, stosować techniki red teaming i symulacje AI-sterowanych ataków, audytować modele AI pod kątem luk i błędów. Takie podejście minimalizuje ryzyko, że własne systemy AI staną się wektorem podatności.



Kontrola i nadzór ludzkiego eksperta AI

Pomimo autonomicznej natury, każde narzędzie AI powinno działać w ramach jasno określonych zasad: nadzorowanych przez ekspertów, z możliwością ręcznego zatrzymania/edycji działania, niezależnych audytów modeli AI. Ludzki nadzór pozostaje kluczowy, aby uniknąć niezamierzonych skutków działania algorytmów.



Współpraca międzynarodowa i standardy etyczne

AI w cyberbezpieczeństwie to nie tylko problem pojedynczych organizacji, ale kwestia globalna. Współpraca międzynarodowa, wymiana informacji o zagrożeniach i wspólne standardy etyczne oraz techniczne są konieczne, aby zapobiegać eskalacji autonomicznych ataków AI.

Przyszłość: współpraca, a nie konkurencja z maszyną

AI w cyberbezpieczeństwie to nie wyrok, lecz potężne narzędzie. Rola AI – obrońcy lub kata – zależy od tego, jak je zaprojektujemy i jakie granice mu wyznaczymy. AI jest narzędziem, które zbliża nas do znacznie wyższego poziomu ochrony, jeśli tylko będzie właściwie projektowane, kontrolowane i regulowane. Kiedy nauczymy się współdziałać z maszyną, a nie konkurować z nią, staniemy się zdolni do

przeciwdziałania coraz bardziej złożonym zagrożeniom, zyskamy przewagę, której nie da się osiągnąć samym ludzkim wysiłkiem. Agentic AI, nie jest jednowymiarowe, jego wpływ zależy od tego, kto je tworzy, jak je kontroluje i w jakim celu są wykorzystywane. W końcu kluczem do bezpieczeństwa nie jest strach przed maszynami, ale mądrość w ich projektowaniu i zastosowaniu.

Aleksandra Polit

Przewidywanie zagrożeń: od reakcji do predykcji. Jak systemy uczone na danych historycznych chronią w czasie rzeczywistym.

Cyberbezpieczeństwo przez dekady funkcjonowało w modelu reaktywnym. Incydent następował, był analizowany, a następnie, z opóźnieniem, wdrażano środki zapobiegawcze. Ten paradygmat przestaje być skuteczny w świecie, w którym ataki są zautomatyzowane, adaptacyjne i coraz częściej wspierane przez sztuczną inteligencję. Odpowiedzią na tę zmianę jest przejście od reagowania na zagrożenia do ich predykcji, czyli wykrywania potencjalnych ataków zanim faktycznie się wydarzą. Kluczową rolę odgrywają tu systemy uczone na danych historycznych, zdolne do działania w czasie rzeczywistym.

Dlaczego reakcja już nie wystarcza

Tradycyjne mechanizmy bezpieczeństwa, np. sygnatury malware, reguły firewallei czy ręczne analizy SOC, bazują na znanych wzorcach ataków. Problem polega na tym, że: nowe warianty zagrożeń powstają szybciej niż możliwe jest ich opisanie w regułach, ataki typu zero-day nie mają wcześniejszych sygnatur, przeciwnik korzysta z automatyzacji i AI, skracając czas od rekonesansu do kompromitacji. W efekcie reakcja po fakcie oznacza realne straty: wyciek danych, przestoje operacyjne, szkody reputacyjne. Stąd rosnące znaczenie systemów, które potrafią przewidywać ryzyko, a nie tylko je rejestrować.

Podstawą predykcyjnych systemów bezpieczeństwa są ogromne zbiory danych historycznych, obejmujące m.in.: logi sieciowe i systemowe, dane o zachowaniu użytkowników, historię incydentów i alertów, wzorce ruchu aplikacyjnego i API, informacje o znanych kampaniach ataków (UEBA – User and Entity Behavior Analytics). Systemy uczone maszynowo analizują te dane w poszukiwaniu powtarzalnych schematów i korelacji, które dla człowieka byłyby niewidoczne. Co istotne, nie chodzi wyłącznie o identyczne zdarzenia, ale o podobieństwa strukturalne: sekwencje działań, tempo zmian, nietypowe kombinacje zachowań. Dzięki temu modele są w stanie zbudować probabilistyczny obraz ryzyka: nie „czy atak już nastąpił”, ale „jak bardzo prawdopodobne jest, że nastąpi w najbliższym czasie”.



Od detekcji anomalii do prognozowania intencji

Wczesne systemy oparte na ML skupiały się na detekcji anomalii - wykrywaniu odchyleń od normy. Dzisiejsze rozwiązania idą krok dalej, próbując odpowiedzieć na pytanie, dlaczego dane zachowanie jest nietypowe i dokąd prowadzi. Przykładowo: pojedyncze nieudane logowanie nie jest zagrożeniem, seria logowań z różnych lokalizacji może być anomalią, połączenie tej serii z nietypowym ruchem lateralnym w sieci wewnętrznej może wskazywać na przygotowanie do eskalacji uprawnień. System predykcyjny nie tylko wykrywa te zdarzenia, ale na podstawie danych historycznych ocenia, czy podobne sekwencje w przeszłości kończyły się incydem. To przesunęło bezpieczeństwo z poziomu sygnałów na poziom intencji atakującego.

Kluczowym wyzwaniem jest połączenie analizy historycznej z działaniem w czasie rzeczywistym. Nowoczesne platformy bezpieczeństwa realizują to poprzez zamkniętą pętlę:

- 1 Analiza strumieniowa** – bieżące dane są porównywane z modelami wytrenowanymi na historii.
- 2 Ocena ryzyka** – każde zdarzenie otrzymuje dynamiczny scoring zagrożenia.
- 3 Predykcja eskalacji** – system szacuje, czy obserwowany wzorzec może przerodzić się w atak.
- 4 Automatyczna reakcja** – zanim dojdzie do kompromitacji, wdrażane są środki zapobiegawcze (np. ograniczenie dostępu, dodatkowa weryfikacja, izolacja zasobów).
- 5 Uczenie zwrotne** – wynik reakcji zasila model, poprawiając przyszłe predykcje. Dzięki temu bezpieczeństwo przestaje być statycznym zestawem reguł, a staje się adaptacyjnym systemem decyzyjnym.

Przyszłość: predykcja jako fundament cyber-odporności

Ryzyka i ograniczenia predykcji opartej na historii

Choć predykcja oparta na danych historycznych daje ogromne możliwości, niesie też istotne ryzyka: stroniczość danych – modele uczą się tego, co było, niekoniecznie tego, co nadejdzie, fałszywe pozytywy – nadmierna predykcja może prowadzić do blokowania legalnych działań, ataki adversarial – świadome manipulowanie danymi wejściowymi w celu „oszukania” modelu, nadmierne zaufanie do automatyzacji – brak nadzoru ludzkiego może eskalować błędy. Dlatego dojrzałe systemy predykcyjne muszą być audytowalne, transparentne i wspierane przez ekspertów, którzy rozumieją zarówno technologię, jak i kontekst biznesowy.

Przewidywanie zagrożeń nie jest już eksperymentem, staje się fundamentem nowoczesnej cyber-odporności. Organizacje, które potrafią łączyć dane historyczne z analizą czasu rzeczywistego, zyskują przewagę nie tylko technologiczną, ale strategiczną. Zamiast gasić pożary, uczą się je przewidywać i zapobiegać ich rozprzestrzenianiu. Kluczowa zmiana polega na redefinicji roli bezpieczeństwa: z funkcji reakcyjnej w funkcję prognostyczną. W świecie, gdzie atakuje maszyna, obrona również musi myśleć jak maszyna – szybciej, szerzej i bardziej probabilistycznie niż człowiek.

Aleksandra Polit

Zarządzanie tożsamością w erze deepfake'ów. Weryfikacja głosu, twarzy i zachowania jako nowa strefa walki o dane.

W miarę jak sztuczna inteligencja staje się coraz bardziej powszechna, pojawiają się technologie, które nie tylko wspierają innowacje, ale też stawiają przed nami poważne wyzwania bezpieczeństwa. Deepfake'y – syntetyczne media generowane przez AI, są jednym z najostrejszych przykładów tego zjawiska. Mogą imitować twarze, głosy, a nawet zachowania ludzi w sposób trudny do odróżnienia od autentycznych nagrań. W rezultacie tradycyjne metody weryfikacji tożsamości, oparte na haśle czy statycznym PIN-ie, stają się niewystarczające. Nowa strefa walki o dane tożsamościowe obejmuje weryfikację głosu, twarzy i zachowania, z wykorzystaniem zaawansowanych metod biometrii i analizy behawioralnej.

Dlaczego deepfake'y zagrażają bezpieczeństwu tożsamości?

Deepfake'y to techniki oparte na modelach generatywnych (np. GAN – Generative Adversarial Networks), które potrafią tworzyć realistyczne obrazy, audio i wideo imitujące osoby rzeczywiste. Choć technologia ma zastosowania pozytywne, od filmów edukacyjnych po film i sztukę cyfrową, jej nadużycia niosą poważne ryzyka, takie jak m.in.: podszywanie się pod osoby publiczne lub prywatne w celu oszustw, wymuszeń lub dyskredytacji, fałszywe nagrania audio/wideo wykorzystywane do manipulacji finansowych (np. nakłanianie do transferów środków), ataki socjotechniczne, w których deepfake'owe rozmowy głosowe zaufanej osoby otwierają drzwi do systemów lub kont. W efekcie to, co kiedyś było domeną filmowców i badaczy AI, przeniknęło do światów finansów, polityki i bezpieczeństwa cyfrowego.



Tradycyjne systemy uwierzytelniania tożsamości opierają się na trzech filarach:



Czymś, co wiesz
(np. hasło, PIN)



Czymś, co masz
(np. token, karta dostępową)



Czymś, czym jesteś
(np. biometryka statyczna:
odcisk palca, skan twarzy).

Pojawienie się deepfake'ów kwestionuje skuteczność trzeciego filaru – biometrii statycznej, ponieważ technologia potrafi naśladować m.in.: strukturę twarzy w wysokiej rozdzielczości, dynamikę ruchów mimicznych, barwę i rytm głosu konkretnej osoby. W rezultacie systemy uwierzytelniania oparte wyłącznie na pojedynczym obrazie twarzy lub krótkim nagraniu głosowym stają się podatne na oszustwa.

Deepfake'owe audio – czyli syntetyczne głosy generowane przez AI – mogą naśladować barwę, akcent i intonację danej osoby. Coraz częściej atakujący wykorzystują je do: podszywania się pod przełożonych w rozmowach telefonicznych, uzyskiwania zgód na transakcje bankowe, ominięcia weryfikacji głosowej w call center.

Weryfikacja głosu: jak bronić się przed audio-deepfake'ami?

Aby odpowiedzieć na to wyzwanie, nowoczesne systemy weryfikacji głosu stosują:



Analizę cech akustycznych głębszych niż barwa: nie tylko podstawowe parametry, ale subtelne cechy, takie jak mikro-pauzy, niuanse artykulacji, czy unikalne wzorce harmoniczne, które trudno odtworzyć nawet zaawansowanym modelom AI.



Analizę zachowania głosowego w czasie rzeczywistym: algorytmy porównują wzorce mówienia z historią interakcji danej osoby, wykrywając anomalie, np. nienaturalne tempo mowy, nieregularną modulację czy brak charakterystycznych zaburzeń rytmu.



Weryfikację kontekstową: systemy mogą pytać o dodatkowe dane lub zadawać losowe pytania w toku rozmowy, co komplikuje próby automatycznego generowania przekonujących odpowiedzi.



Biometria twarzy: dynamiczne, a nie statyczne podejście

Klasyczne skanery twarzy opierają się na pojedynczym obrazie lub krótkim nagraniu. Deepfake'y potrafią jednak generować realistyczne obrazy twarzy i sekwencje mimiczne, które mylą systemy oparte na analizie statycznej. Dlatego nowoczesna biometria twarzy:

- **Analizuje mikro-ruchy i mimikę w czasie rzeczywistym**, takie jak: subtelne aspekty mrugnięć, mikroekspresje mięśni twarzy, rytm oddechu i mimowolne zmiany napięcia mięśniowego. Te cechy są trudniejsze do odtworzenia w deepfake'ach, ponieważ wymagają zgodności dynamicznych sygnałów na poziomie biologicznym.
- **Stosuje uwierzytelnianie wieloczynnikowe**, łącząc biometrię twarzy z innymi modalnościami (np. głosową), co znacząco utrudnia atakującym oszukanie systemu przy użyciu pojedynczej techniki deepfake.
- **Wykorzystuje Liveness Detection** (wykrywanie żywotności): technologia ta sprawdza, czy przed kamerą znajduje się żywy człowiek (np. poprzez polecenie mrugnięcia, śledzenia wzrokiem punktu lub analizę odbicia światła od gałki ocznej), co skutecznie eliminuje próby użycia statycznych zdjęć czy odtworzonych nagrań deepfake.

Biometria behawioralna – nowy wymiar weryfikacji

Żeby skutecznie chronić tożsamość cyfrową, coraz częściej wykorzystuje się biometrię behawioralną – czyli analizę wzorców zachowania użytkownika. Obejmuje to: sposób poruszania myszką i klikania, tempo i rytm pisanie na klawiaturze, wzorce korzystania z aplikacji i urządzeń, lokalizację, i typ połączenia sieciowego. Te wzorce są unikatowe dla każdego użytkownika. i znacznie trudniejsze do skopiowania niż sam

obraz twarzy lub nagranie głosu. Systemy analizują historyczne dane zachowań i w czasie rzeczywistym wykrywają odchylenia, które mogą świadczyć o oszustwie.

Skuteczna weryfikacja tożsamości w erze deepfake'ów wymaga strategii hybrydowej:

- ◆ **Multi-modalna biometryka:** Połączenie weryfikacji głosowej, twarzy i wzorców behawioralnych
- ◆ **Uczenie maszynowe na danych historycznych:** Modele analizują wzorce z historii interakcji, aby definiować normy i wykrywać anomalie.
- ◆ **Analiza kontekstowa:** System bierze pod uwagę czas, lokalizację i zachowanie użytkownika, np. logowanie z nietypowej lokalizacji może wymagać dodatkowej weryfikacji.
- ◆ **Ciągłe uczenie się:** Modele są aktualizowane o nowe próbki zachowań, co zwiększa odporność na ewoluujące próby ataku.
- ◆ **Nadzór ludzki:** Tam, gdzie automatyczna weryfikacja nie daje jednoznacznej odpowiedzi, rolę musi przejąć analityk bezpieczeństwa, wspierający model algorytmiczny.

W dobie deepfake'ów tożsamość cyfrowa przestaje być statycznym hasłem lub odciskiem palca. Staje się dynamicznym profilem, który obejmuje cechy fizyczne, behawioralne i kontekstowe. Walka o dane i tożsamość to dziś starcie technologii generujących fałszywe sygnały z coraz bardziej wyrafinowanymi systemami ich wykrywania. Organizacje, które rozumieją, że kluczem jest weryfikacja ciągła i adaptacyjna, a nie jednorazowa autoryzacja, zyskują przewagę w ochronie zasobów, reputacji i prywatności użytkowników.



Agentic AI

Agentic AI to nowa generacja systemów, które nie tylko odpowiadają na pytania, ale same podejmują inicjatywę. To jak mieć wirtualnego asystenta, który zobaczy, że masz bałagan w projektach, więc sam poukłada pliki, dopyta współpracowników i zaplanuje spotkanie – zanim w ogóle zauważysz, że jest problem. W 2026 roku takie agentowe AI zaczynają być standardem w narzędziach biznesowych: łączą kontekst, integrują dane z wielu źródeł i podejmują autonomiczne mikro-działania. Krótko mówiąc: AI, które przestało czekać na instrukcje, a zaczęło być proaktywne. (Słowem, to ten typ współpracownika, który naprawdę lubi Excela.)

AI Security

AI Security – wojna na prompt-hacking trwa. W 2026 roku ataki polegające na skłonieniu AI do podejmowania niechcianych działań stają się coraz bardziej wyrafinowane. Firmy budują więc systemy „AI strażników”, które... monitorują inne AI. Tak, mamy już AI do pilnowania AI. Meta-poziom osiągnięty. Warto tu wspomnieć o jednym konkretnym zjawisku: Indirect Prompt Injection (pośrednie wstrzyknięcie poleceń). To sytuacja, w której haker nie atakuje AI bezpośrednio, ale zostawia „pułapkę” na stronie internetowej, którą agent AI przeczyta i... nagle zacznie wykonywać instrukcje hakera zamiast Twoich (np. wyśle Twoje hasła na obcy serwer).

Samodzielne laboratoria przyspieszyły Naukę

Naukowcy pokazali pełnię możliwości tzw. self-driving labs. Berkeley A-Lab z pomocą AI zsyntetyzował 41 nowych związków chemicznych w 17 dni, działając 24/7 bez przerw na kawę. Porównanie: człowiek – doktorant – miesiące. AI-lab – pół miesiąca.

Vibe crime

„Vibe crime” – przestępstwo... nastroju? To jeden z najnowszych memicznych terminów świata AI. „Vibe crime” to sytuacja, gdy coś technicznie jest poprawne, ale jego klimat, narracja albo energia są dziwnie niewłaściwe. Przykład? Poprosisz AI o design plakatu konferencji RODO, a ono robi grafikę w stylu neonowego koncertu techno. Albo generuje idealny tekst maila do klienta, ale kończy go „miłego życia”. To nie błąd. To vibe crime. W świecie, gdzie modele coraz lepiej rozumieją dane, ale wciąż uczą się subtelnych ludzkich nastrojów – vibe crime to humorystyczny, ale całkiem realny problem User Experience. Gdy marka produktu to nie tylko logo, ale przede wszystkim jej osobowość uniknięcie vibe crime to nowa walka o autentyczność w świecie syntetycznym – vibe crime to nie tylko śmieszna wpadka, ale realne ryzyko dla spójności marki. Klient czuje, że „coś tu nie gra” i przestaje ufać firmie.



Pierwsza AI trenowana w kosmosie

Tak, to się wydarzyło. W 2025 po raz pierwszy wytrenowano model AI na orbicie, na pokładzie satelity. To otworzyło drogę do trenowania modeli w środowiskach z inną radiacją i temperaturami, co ma znaczenie dla kosmicznych technologii. Chodzi tu przede wszystkim o oszczędność energii i czasu. Przesyłanie ogromnych zbiorów danych z satelity na Ziemię jest wolne i drogie. Jeśli satelita „myśli” i uczy się sam na miejscu, staje się autonomicznym agentem AI.



ANALIZA RZECZYWISTOŚCI

AI wie wszystko, bo ma dostęp do internetu

Rzeczywistość: Większość modeli to „kapsuły czasu”. Działają na zamkniętej migawce danych z przeszłości. Nawet modele z dostępem do sieci nie „wiedzą”, a jedynie „wyszukują”. AI nie jest wszechwiedzącym mędrce, tylko genialnym archiwistą z ograniczoną kartą biblioteczną.

AI jest obiektywna

Rzeczywistość: AI to lustro, a nie filtr. Jeśli dane treningowe są skażone uprzedzeniami, model je po prostu zwielokrotni. AI nie eliminuje ludzkich błędów – ona je automatyzuje na masową skalę. Bez etycznego nadzoru, algorytm może być tak samo stronniczy jak najbardziej uprzedzony człowiek. Dlatego potrzebne są standardy etyki, ewaluacje, nadzór i transparentne dane. AI dziedziczy nasze błędy – nie eliminuje ich magicznie.

Wystarczy podać dane i AI sama zrobi resztę

Rzeczywistość: To największa pułapka wdrożeniowa. Surowe dane bez struktury to dla AI tylko szum. Obowiązuje stara zasada programistów: Garbage In, Garbage Out. Jeśli „zatankujesz” model niechlujnymi danymi, dostaniesz błyskawiczne, ale całkowicie błędne wyniki.

AI kreatywna = AI inteligentna emocjonalnie

Rzeczywistość: Kreatywność AI to czysta statystyka prawdopodobieństwa, a nie natchnienie. Model nie „czuje” piękna ani humoru, nie ma inspiracji, marzeń, ani poczucia piękna, on po prostu wie, które piksele lub słowa statystycznie występują obok siebie. Kreatywność AI to iluzja artysty, za którą stoi bardzo szybki kalkulator.



ANALIZA RZECZYWISTOŚCI

AI to magia – wszystko potrafi

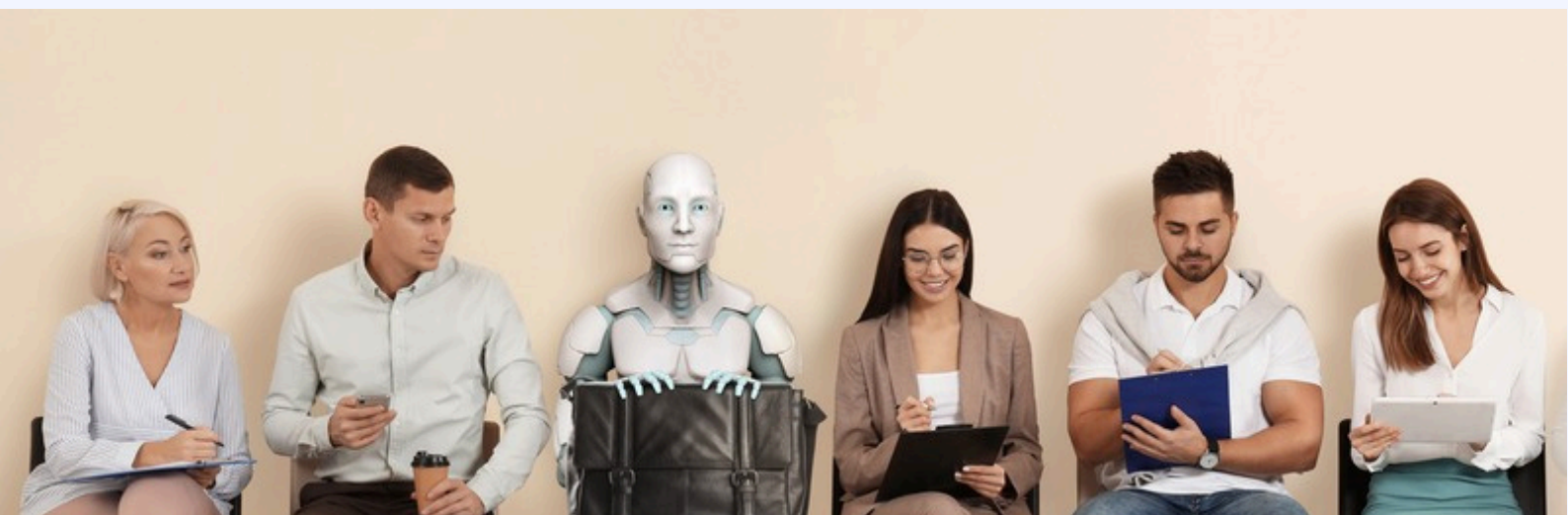
Rzeczywistość: AI to nie różdżka, a zaawansowana optymalizacja. Świetnie radzi sobie z powtarzalnymi schematami, ale wyklada się na zdrowym rozsądku i sytuacjach, których nigdy wcześniej „nie widziała”. Tam, gdzie kończy się schemat, AI staje się bezradna jak dziecko.

Wdrożenie AI to jednorazowy koszt

Rzeczywistość: Myślenie o AI jak o zakupie lodówki – „kupujesz i działa” – to błąd. AI to raczej ogród. Wymaga stałych nakładów na prąd (tokeny), aktualizację danych i tzw. Alignment, czyli pilnowanie, by model z czasem nie zaczął „zmyślać” lub odbiegać od celów firmy. Koszt utrzymania i nadzoru jest często wyższy niż sama licencja.

AI zastąpi pracowników (kropka)

Rzeczywistość: AI nie zastąpi ludzi, ale ludzie używający AI zastąpią tych, którzy jej nie używają. W 2026 roku widzimy, że automatyzacja nie likwiduje etatów masowo, ale drastycznie zmienia zakres obowiązków. Zamiast „wykonywać”, pracownik zaczyna „nadzorować” i „redagować” pracę algorytmów.



Aleksandra Polit

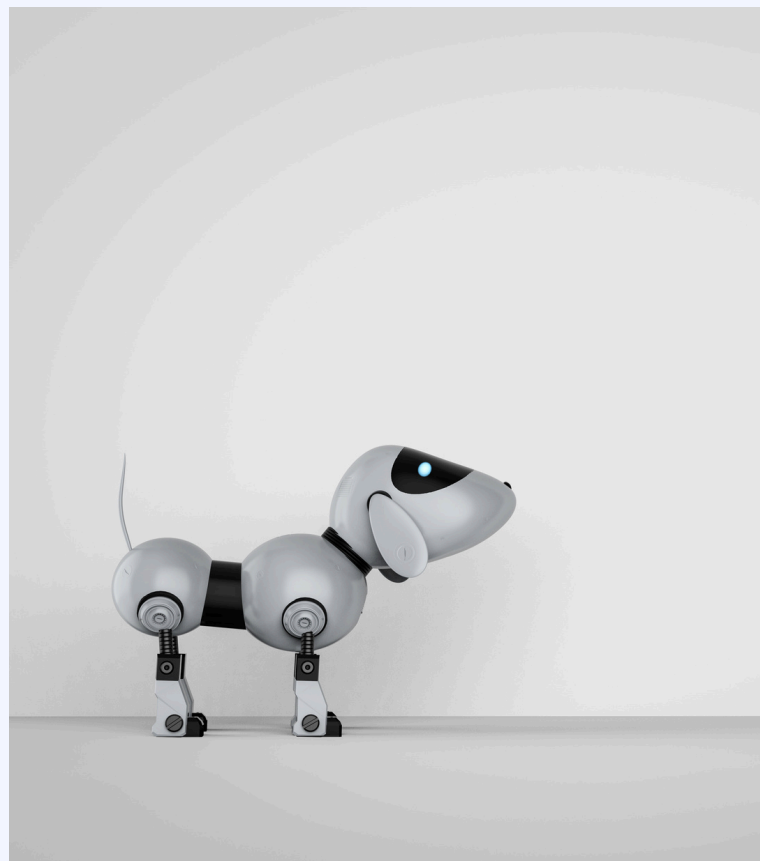
Jak zwierzęta przewidują cyberzagrożenia – science-fiction czy bliska przyszłość?

Brzmi jak tytuł filmu klasy B albo eksperyment szalonego naukowca: pies ostrzegający przed ransomware, gołąb wykrywający phishing, a ośmiornica reagująca na anomalie w ruchu sieciowym. A jednak... gdyby odrzucić dosłowność i spojrzeć na temat przez pryzmat biologii, danych i AI, okazuje się, że świat zwierząt od dawna robi to, do czego dopiero aspiruje cyberbezpieczeństwo: przewiduje zagrożenia, zanim staną się widoczne.

W naturze przeżywa nie ten, kto reaguje na atak, ale ten, kto wyczuwa go wcześniej. Ptaki nagle milknące przed trzęsieniem ziemi, psy reagujące na zmiany hormonalne właściciela czy stada ryb zmieniające kierunek ruchu bez widocznego bodźca – to wszystko przykłady predykcji opartej nie na „alarmie”, lecz na mikrosygnalach. Zwierzęta nie analizują jednego sygnału, nie potrzebują pełnego obrazu, nie czekają na potwierdzenie. Reagują na odchylenie od normy, często szybciej niż człowiek jest w stanie je zarejestrować. Brzmi znajomo? Dokładnie tak działają nowoczesne systemy wykrywania anomalii w cyberbezpieczeństwie.

Cyberbezpieczeństwo coraz częściej kopiuje rozwiązania z natury. Przykłady? Systemy kolektywnej detekcji – wzorowane na stadach ptaków lub ławicach ryb, gdzie pojedynczy sygnał nic nie znaczy, ale skorelowany ruch wielu „agentów” ujawnia zagrożenie. Predykcja behawioralna, podobnie jak drapieżnik wyczuwa zmianę rytmu ofiary, AI

analizuje subtelne zmiany w zachowaniu użytkownika lub sieci. Reakcja bez centralnego dowództwa, mrówki nie potrzebują SOC-a (Swarm Intelligence – Inteligencja Stada). W IT to sposób projektowania systemów, w których tysiące małych programów (agentów) współpracuje bez centralnego serwera. Tak jak mrówki budują mrowisko, tak „cyfrowe stado” wykrywa hakerów, analizując mikrozmiany w ruchu sieciowym, a nowoczesne systemy bezpieczeństwa coraz częściej działają zdecentralizowanie i autonomicznie. To już nie science-fiction. To bio-inspirowana architektura cyberbezpieczeństwa.



Zwierzęta jako mistrzowie predykcji

Czy zwierzęta „czują” zagrożenia cyfrowe? Pośrednio tak.

Zwierzęta nie rozumieją malware ani phishingu, ale... reagują na skutki cyfrowych zagrożeń. Przykłady z realnego świata: psy reagujące na stres właściciela wywołany atakiem finansowym, zwierzęta laboratoryjne wykazujące zmiany zachowania w środowiskach o zaburzonej infrastrukturze (np. po cyberatakach na systemy przemysłowe), eksperymenty, w których algorytmy uczą się wzorców zagrożeń na podstawie modeli zachowań zwierząt. Nie chodzi więc o to, że kot wykryje backdoora. Chodzi o to, że mechanizmy predykcji wykształcone przez miliony lat ewolucji są dziś kopiowane przez AI.

Najbardziej futurystyczna myśl?

To już się dzieje. Najciekawsza ironia polega na tym, że gdy mówimy „jak zwierzęta przewidują cyberzagrożenia”, w rzeczywistości pytamy: Czy systemy bezpieczeństwa mogą działać tak instynktownie, jak organizmy żywe? Odpowiedź brzmi: już próbują. Dzisiejsze modele predykcyjne nie „rozumieją” ataku, one go wyczuwają. Tak samo jak zwierzę, które nie zna pojęcia trzęsienia ziemi, ale wie, że trzeba uciekać.

Biomimetyka (naśladowania natury przez technologię)

Zwierzęta nie przewidują cyberzagrożeń wprost. Ale to, jak przewidują niebezpieczeństwo, staje się blueprintem dla przyszłości cyberbezpieczeństwa. Instykt, korelacja słabych sygnałów, reakcja przed faktem, to nie magia, to biologia. A skoro AI coraz bardziej przypomina system nerwowy, a SOC zaczyna działać jak stado... to pytanie nie brzmi czy to science-fiction, tylko: jak szybko cyberbezpieczeństwo stanie się bardziej zwierzęce niż ludzkie?

Wniosek dla managera:

W 2026 roku bezpieczeństwo Twojej firmy to nie tylko kwestia lepszych haseł. To budowanie „układu nerwowego” organizacji, który wzorem natury potrafi wykryć zagrożenie, zanim haker w ogóle naciśnie klawisz Enter.





Czy AI faktycznie zwiększa bezpieczeństwo, czy tylko je komplikuje?

Odpowiedź: AI to potężny „mnożnik siły”. Skracza czas reakcji (MTTR) z godzin do milisekund. Jednak komplikacje pojawiają się, gdy traktujemy ją jak magiczną „czarną skrzynkę”. Bez jasnych procedur i nadzoru, AI nie rozwiązuje problemów, ona po prostu sprawia, że Twoje błędy w konfiguracji dzieją się szybciej.

Czy używanie AI w cyberbezpieczeństwie to zawsze przetwarzanie danych osobowych?

Odpowiedź: Niemal zawsze. Adres IP, wzorzec pisania na klawiaturze czy biometria behawioralna to dane podlegające pod RODO. W 2026 roku kluczowe nie jest pytanie „czy?”, ale „jak?”. Modele muszą być trenowane na danych zanonimizowanych, a ich dostęp do „żywych” baz danych musi być ściśle limitowany.

Czy cyberprzestępcy też mają swoje AI?

Odpowiedź: Niestety tak i nie mają ograniczeń etycznych. Używają modeli typu WormGPT do generowania perfekcyjnego phishingu i deepfake’ów, które brzmią identycznie jak Twój szef. To nie jest już walka człowieka z hakerem, to wyścig zbrojeń algorytmów.



CASE STUDY

Sytuacja: System wykrył, że pracownik loguje się o 3:00 nad ranem i pobiera duże paczki danych. Werdykt AI: Atak wewnętrzny! Konto zablokowane, wezwany dział HR.

Rzeczywistość: Pracownik po prostu nadrabiał zaległości przed urlopem, korzystając z bezsenności.

Wniosek: Sama detekcja to za mało. Potrzebujemy XAI (Explainable AI) – systemu, który nie tylko mówi „zablokuj”, ale potrafi wyjaśnić „dlaczego”. Bez ludzkiego bezpiecznika, AI może stać się najbardziej nadgorliwym i paraliżującym pracownikiem w Twojej firmie.



QUIZ CHALLENGE

Sprawdź, czy Twoja wiedza o AI w 2026 roku to "state-of-the-art", czy raczej "outdated model".

- 1. System AI blokuje konto pracownika o 3:00 nad ranem. Co robisz jako szef bezpieczeństwa?**
 - a) Śpię spokojnie – AI wie, co robi.
 - b) Dzwonię do HR, żeby szykowali wypowiedzenie.
 - c) Sprawdzam uzasadnienie (XAI) – może to tylko pracoholik przed urlopem?
 - d) Odłączam prąd w całym biurcu (na wszelki wypadek).

- 2. Twoja firma wdraża AI do monitoringu sieci. Co model będzie „zjadał” na śniadanie najczęściej?**
 - a) Treść prywatnych plotek na Slacku.
 - b) Logi systemowe i biometrię behawioralną (czyli to, jak specyficznie piszesz na klawiaturze).
 - c) Zdjęcia z Twoich ostatnich wakacji.
 - d) Kanapki z kuchni socjalnej.

- 3. Co jest „piętą achillesową” nowoczesnych systemów AI (tzw. Black Box)?**
 - a) Zbyt wysoki rachunek za prąd.
 - b) Brak wyjaśnialności – AI mówi „źle”, ale nie powie „dlaczego”.
 - c) To, że AI nie pije kawy i za szybko pracuje.
 - d) Zbyt ładny interfejs graficzny.

- 4. Kiedy RODO zaczyna groźnie patrzeć na Twoje AI?**
 - a) Gdy AI pracuje w nadgodzinach.
 - b) Gdy automat podejmuje decyzje o człowieku (np. blokada) bez żadnej furtki dla interwencji ludzkiej.
 - c) Gdy AI ma serwery w chmurze, która nazywa się „Cumulus”.
 - d) Tylko wtedy, gdy AI zaczyna mówić po niemiecku.

- 5. Czy hakerzy używają AI?**
 - a) Nie, oni wolą tradycyjne rzemiosło i wirusy „domowej roboty”.
 - b) Tak, ale tylko do poprawiania błędów ortograficznych w mailach.
 - c) Tak – tworzą „złe bliźniaki” Twojego szefa (deepfake) i piszą phishing, który brzmi jak od Twojej mamy.
 - d) Tak, ale tylko w filmach na Netflixie.

- 6. Idealna rola AI w Twoim zespole to:**
 - a) Samodzielny sędzia i kat, który zwalnia ludzi mailem.
 - b) Genialny analityk, który przesiewa miliony danych, by podać Ci gotowe wnioski na tacy.
 - c) Robot, który zastępuje dział IT, żeby zaoszczędzić na owocowych czwartkach.
 - d) Tajny agent, o którym nikt w firmie nie wie.

- 7. Chcesz wdrożyć AI w bezpieczeństwie. Od czego zaczynasz?**
 - a) Od zakupu najdroższej licencji z napisem „Cyber-Magic”.
 - b) Od analizy ryzyka (DPIA), porządkowania danych i zapewnienia nadzoru człowieka.
 - c) Od wrzucenia wszystkich danych firmy do publicznego ChatGPT.
 - d) Od modlitwy o to, żeby nic nie wybuchło.

***Poprawne odpowiedzi ujawnimy w kolejnym numerze.**



Następny numer już wkrótce!

Temat wydania: Zero Trust & Identity: nowy fundament bezpieczeństwa

1. OKIEM EKSPERTA

- Zero Trust 2.0
- Tożsamość: największe aktywo i największe ryzyko
- Phishing 3.0

2. TREND ALERT

- Najciekawsze narzędzia i rozwiązania z obszaru Identity Protection – co działa, a co jest tylko marketingiem.

3. ANALIZA RZECZYWISTOŚCI

- „MFA to tarcza nie do przebicia? Obalamy największy mit branży i pokazujemy, jak hakerzy obchodzą dwuskładnikowe uwierzytelnienie w 30 sekund.”

4. CZY WIESZ, ŻE...

- Statystyki i fakty z cyberswiata, które na pierwszy rzut oka brzmią nieprawdopodobnie, a jednak są prawdziwe. Niektóre z nich mogą zaskoczyć nawet specjalistów.

5. SPOŁECZNOŚĆ W AKCJI

- Prawdziwe pytania, realne dylematy i doświadczenia czytelników – bez filtrów i marketingowych sloganów.

**Jeśli myślisz, że Zero Trust to tylko modne hasło,
a identity to problem działu IT, ten numer może
zmienić Twoje podejście.**

**Dziękujemy za przeczytanie
naszego biuletynu!**

**Masz pytania?
Skontaktuj się z nami**



www.forsafe.pl



600 005 880



biuro@forsafe.pl



ul. Traktorowa 170, 91-203 Łódź