

Warsztaty 7: Deepfake i oszustwa AI

Prompty do weryfikacji podejrzanych treści

ZASADA

AI może pomóc wypisać czerwone flagi.
AI nie potwierdza, że podejrzana treść jest bezpieczna.
Źródło sprawdzam samodzielnie.

Nie wklejam:

- kodów,
- haseł,
- danych logowania,
- danych karty,
- PESEL-u,
- pełnych dokumentów,
- prywatnych rozmów innych osób,
- aktywnych linków, których nie chce klikac.

1. PODSTAWOWY PROMPT

Pomóż mi spokojnie ocenić podejrzana treść.

Opis sytuacji:

[opisz bez danych prywatnych, bez linków do klikania i bez kodów]

Wypisz:

1. sygnały presji,
2. prośby o pieniądze, kod, dane albo kliknięcie,
3. miejsca, gdzie ktoś może się podszywać,
4. informacje, których nie da się potwierdzić z opisu,
5. pytania, które muszę sprawdzić drugim kanałem.

Nie rozstrzygaj, czy to na pewno prawda albo fałsz.

Na końcu podaj najbezpieczniejszy pierwszy krok bez klikania w link i bez płacenia.

2. GDY KTOS PROSI O KOD ALBO PIENIADZE

Ocen te sytuacje pod kątem oszustwa.

Opis:

[opisz sytuacje bez kodów, numerów kont i danych prywatnych]

Sprawdź:

1. czy jest presja czasu,
2. czy ktoś prosi o kod, przelew, dane albo aplikacje,
3. czy tożsamość osoby jest potwierdzona,
4. jaki drugi kanał sprawdzenia jest najbezpieczniejszy,
5. czego nie powinienem/powinnam robić pod presją.

Nie podawaj instrukcji wykonania płatności.

Podaj tylko bezpieczne kroki zatrzymania i sprawdzenia.

3. GDY WIDZE REKLAMĘ Z OBIECUJĄCYM ZYSKIEM

Pomóż mi ocenić reklamę, która obiecuje szybki zysk.

Opis reklamy:
[opisz bez klikania w link]

Wypisz:

1. obietnice zysku,
2. elementy presji,
3. czy reklama korzysta z autorytetu osoby albo instytucji,
4. jakie informacje powinny być na oficjalnej stronie,
5. gdzie może sprawdzić ostrzeżenia samodzielnie.

Nie oceniaj atrakcyjności inwestycji.
Skup się na ryzyku oszustwa i bezpiecznym pierwszym kroku.

4. GDY DOSTAJE FILM ALBO NAGRANIE

Pomóż mi sprawdzić, czy nagranie wymaga ostrożności.

Opis nagrania:
[opisz co widzisz i słyszysz]

Wypisz:

1. sygnały w obrazie,
2. sygnały w głosie,
3. niespójności treści,
4. brakujący kontekst: data, miejsce, pełne źródło,
5. co sprawdzić przed przesłaniem dalej.

Nie rozstrzygaj pewnie, czy to deepfake.
Podaj pytania do weryfikacji i bezpieczny krok.

5. GDY CHCE OSTRZEC BLISKICH

Pomóż mi napisać spokojną wiadomość do rodziny albo znajomych.

Cel:
Chcę ostrzec przed prośbami o kod, pieniądze albo kliknięcie w link,
ale bez straszenia i bez zawstydzania.

Napisz:

1. wersje krótka,
2. wersje ciepła i rodzinna,
3. jedno zdanie-zasada,
4. propozycje hasła albo procedury rodzinnej.

Ton:
spokojny, konkretny, bez paniki.

6. GDY KTOS JUZ KLIKNAL ALBO PODAL DANE

Pomóż mi ułożyć spokojny plan ograniczenia szkody.

Opis sytuacji:
[opisz bez haseł, kodów, numerów kart, PESEL-u i aktywnych linków]

Wypisz:

1. jakie dowody warto zachować,
2. które konta, karty albo hasła mogą wymagać zabezpieczenia,
3. z jaką instytucją skontaktować się oficjalnym kanałem,
4. czego nie kasować przed zgłoszeniem,

5. jak opisać sytuację rodzinie albo zaufanej osobie bez wstydu.

Nie podawaj instrukcji płatności ani omijania zabezpieczeń.
Skup się na zatrzymaniu szkody, zgłoszeniu i bezpiecznym kontakcie.