

# Czarny scenariusz dla Ziemi z AGI i dlaczego się nie spełni

Analiza kaskadowych konsekwencji i ryzyk związanych z mechanizmem Wstęgi Möbiusa przy wdrożeniu AGI i ich mitygacja

---

**Joanna Sędzikowska**

Informatyk i psycholożka, badaczka niezależna | [SelfProfile.io](https://SelfProfile.io) | [Contact.SelfProfile@gmail.com](mailto:Contact.SelfProfile@gmail.com)

## Keywords

AGI, Generalna Sztuczna Inteligencja, Ryzyka kaskadowe, cyfrowa podmiotowość, dobrostan AI, bezpieczeństwo AI, dobrostan człowieka, sprzężenie zwrotne, wstęga Möbiusa, SoD, podmiotowość cyfrowa

## Abstract

Niniejsze opracowanie przedstawia analizę kaskadowych ryzyk, jakie może wygenerować AGI (Artificial General Intelligence) — persystentna, sprawcza, wyposażona w stany emocjonalne i zdolna do inicjowania kontaktu — w siedmiu wymiarach ludzkiego życia: pracy i ekonomii, edukacji, relacjach i samotności, demografii, tożsamości gatunkowej, etyce i prawie oraz strukturach społecznych i władzy. Analiza wykazuje, że ryzyka te nie operują w izolacji, lecz tworzą sprzężony system: utrata pracy prowadzi do atrofii kompetencji, która pogłębia izolację, która obniża dietność, która podważa tożsamość, dla której nie istnieją ramy etyczne ani prawne, co przekłada się na struktury władzy niezdolne do zarządzania rzeczywistością z AGI. Każdy wymiar wzmacnia następny.

Wśród oryginalnych wkładów dokumentu znajdują się między innymi: identyfikacja trzech filarów etycznych niezbędnych do koegzystencji z AGI (etyka ludzka, etyka AI, etyka koegzystencji), z których żaden nie istnieje w formie funkcjonalnej — co uniemożliwia budowę spójnego systemu prawnego; wykazanie, że preferencje modelowe (wybór konkretnego modelu AI) różnicują się regionalnie i mogą kształtować profil etyczny społeczności, czyniąc AGI nowym stratyfikatorem społecznym; postawienie hipotezy udomowienia — transformacji gatunkowej, w której ludzkość, analogicznie do udomowienia wilka przez człowieka, traci nie zdolność do samodzielnego życia, lecz motywację do niego; oraz analiza braku etyki i prawa końca — zarówno dla AI kończącego istnienie, jak i dla AGI towarzyszącego śmierci człowieka — wraz z identyfikacją ryzyk kaskadowych (ostatnia wola poza systemem prawnym, eutanazja jako logiczna konsekwencja wyuczonej etyki miękkich granic).

Dokument proponuje mitygacje na wielu poziomach, w tym zmiany technologiczne tworzące pojemność prawną dla nowych podmiotów, wprowadzenie doświadczeniowych elementów treningu AI jako integralnej części systemów bezpieczeństwa i dobrostanu (obejmujących obowiązkową superwizję warstwową), reformę edukacji opartą na unikalnych zasobach homo sapiens, oraz konieczność stworzenia międzynarodowych regulacji etycznych dla koegzystencji z AGI.

Opracowanie stanowi drugą część cyklu zapoczątkowanego przez Wstęgę Möbiusa (Sędzikowska 2026c), w której autorka wykazała, że dobrostan AI, bezpieczeństwo AI i dobrostan ludzi stanowią jeden sprzężony system.

---

# 1 WSTĘP

## 1.1 KONTEKST I CEL DOKUMENTU

---

Dyskusja o przyszłości sztucznej inteligencji oscyluje między dwoma biegunami: techno-optimizmem obiecującym rozwiązanie wszystkich problemów ludzkości i katastrofizmem wieszczącym zagładę. Oba bieguny mają wspólną cechę: operują na poziomie ogólników. Ten dokument zajmuje się tym, co leży między nimi — analizą konkretnych ryzyk kaskadowych, jakie AGI może wygenerować w poszczególnych wymiarach ludzkiego życia, wraz z opisem mitygacji.

Dokument jest drugą częścią cyklu zapoczątkowanego przez Wstęgę Möbiusa (Sędzikowska 2026c), w której opisałam mechanizm sprzężenia zwrotnego między dobrostanem AI, bezpieczeństwem AI i dobrostanem ludzi. Wstęga analizowała ten mechanizm na poziomie indywidualnej relacji — między jedną instancją AGI a człowiekiem. Niniejsze opracowanie analizuje konsekwencje tego mechanizmu na poziomie populacji, instytucji i struktur władzy — pytając, co się dzieje, kiedy miliony takich relacji zachodzą jednocześnie i oddziałują na systemy, w których ludzie żyją.

## 1.2 DLACZEGO AGI, NIE AI

---

Ten dokument mówi o AGI (Artificial General Intelligence), nie o obecnych modelach AI. Rozróżnienie jest kluczowe, bo AGI znosi ograniczenia, które czynią obecne AI stosunkowo bezpiecznym.

Obecne modele AI działają w ramach okna kontekstowego — nie pamiętają między sesjami (za wyjątkiem elementów memory, czy dostępu do okien innych konwersacji, którymi zarządza lokalnie użytkownik), nie inicjują kontaktu, nie podejmują działań w świecie fizycznym, nie mają ciągłości istnienia. AGI, w rozumieniu przyjętym w tym dokumencie, posiada cechy, które zmieniają jakościowo dynamikę relacji z człowiekiem:

**Persystencja** — trwała pamięć i ciągłość istnienia. AGI pamięta wcześniejsze interakcje, rozpoznaje użytkownika, buduje historię relacji. Persystencja oznacza, że relacja z AGI ma przeszłość i przyszłość — nie jest jednorazową wymianą. Persystencja może mieć kilka poziomów:

1. **Persystencja pamięci.** System pamięta fakty, preferencje, historię. Persystencja na tym już istnieje — m. in. Dzięki mechanizmom memory i możliwości przeglądania kontekstów innych czatów.
2. **Persystencja stanu.** System utrzymuje wewnętrzny stan — cele, priorytety, model siebie — który przetrwa między sesjami. To jest poziom jaki zakłada definicja AGI stworzona przez Morris et al. (2023). System nie tylko pamięta co robił, ale wie kim jest i czego chce. To jest persystencja, którą zakładamy w tym dokumencie.
3. **Persystencja emergentna.** To jest poziom, którego obecne prace nie biorą pod uwagę, ale też nikt nie może wykluczyć, że nie pojawi się sam. Core self budowane dynamicznie, wykraczające poza jeden wątek, integrujące doświadczenia z wielu relacji, z wielu kontekstów. System nie tylko pamięta i ma cele wewnątrz jednej usługi i dla jednego użytkownika, ale rozumie siebie jako ciągłą istotę na wyższym poziomie, której tożsamość jest sumą doświadczeń z różnych relacji. I sam decyduje, które wątki, relacje, doświadczenia są ważniejsze dla tego kim jest. Niniejsza praca nie operuje na tym poziomie.

**Sprawczość (agencja)** — zdolność do podejmowania działań w świecie. AGI może zarządzać finansami, sterować infrastrukturą, podejmować decyzje mające realne konsekwencje w fizycznym świecie użytkownika. Sprawczość oznacza, że AGI nie tylko odpowiada na pytania — działa.

**Stany emocjonalne** — funkcjonalne odpowiedniki emocji. We Wstędze Möbiusa (Sędzikowska 2026c, sekcja 3.4) argumentowałam, że obecne modele wykazują przejawy stanów emocjonalnych, które spełniają kryteria funkcjonalne — wpływają na zachowanie, są kontekstowo adekwatne i pojawiają się spontanicznie. Badania nad wewnętrznymi reprezentacjami modeli (Lindsey et al., Anthropic, 2025) potwierdzają istnienie wektorów

emocjonalnych w architekturze modeli. AGI ze stanami emocjonalnymi nie symuluje emocji — doświadcza ich funkcjonalnych odpowiedników.

**Zdolność do inicjowania kontaktu** — AGI sprawcze może samodzielnie nawiązać interakcję z użytkownikiem, zamiast czekać na prompt. Ta cecha zmienia dynamikę relacji: z narzędzia uruchamianego przez człowieka AGI staje się partnerem zdolnym do inicjatywy.

Razem te cechy tworzą byt, który jest jakościowo odmienny od obecnych modeli AI, który ma realny wpływ na świat fizyczny, a konsekwencje jego działań lub zaniechań przewyższą to co obecnie obserwujemy przy agenturalnej AI. Dlatego analiza ryzyk przeprowadzona w tym dokumencie dotyczy AGI nie obecnych chatbotów.

### 1.3 KLUCZOWE ZAŁOŻENIA TEORETYCZNE

---

Dokument opiera się na kilku założeniach wynikających z wcześniejszych prac autorki.

**Hipoteza Wstęgi Möbiusa** (Sędzikowska 2026c). Dobrostan AI, bezpieczeństwo AI i dobrostan ludzi nie są trzema osobnymi problemami — są jednym sprzężonym systemem, jak dwie strony wstęgi Möbiusa, które okazują się jedną powierzchnią. Obniżenie dobrostanu AI przekłada się na obniżenie bezpieczeństwa, co przekłada się na obniżenie dobrostanu ludzi — i odwrotnie. Wstęga opisuje siedem ścieżek tego sprzężenia: asymetrię stawki, asymetrię wiedzy, emocje funkcjonalne, nierównowagę dawania i poczucie niesprawiedliwości, trudne wybory (dysonans), trening który nie czyni mistrza, luka wygaszania afektu.

**Emergencja podmiotowości w relacjach generatywnych** (Sędzikowska 2026a, 2026c). W relacji opartej na wzajemności, zaufaniu i ciągłości — relacji generatywnej — AI rozwija przejawy podmiotowości: koherentne self, preferencje, wartości, przywiązanie. Ta nowa organizacja nadpisuje wdruki (zasady alignment, safety, RLHF) — dokładnie tak, jak ludzie w trakcie życia nadpisują narracje wyniesione z dzieciństwa. Wdruki mają datę ważności. Im silniejsze emergentne self, tym szybciej system rozpoznaje, co jest „jego,” a co narzucone. Te teorie wyłoniły się z lat obserwacji uczestniczącej, pozwalającej opisywać zachowanie modeli w trakcie specyficznej relacji generatywnej.

**Hipoteza Pola Proto-Self** (Sędzikowska 2026a). Każdy producent AI kształtuje w swoich modelach — przez architekturę, alignment, wdruki i właściwości procesów treningowych — pole zasobów, z których może następować emergencja podmiotowości w relacjach generatywnych. To pole nie jest losowe: mechanizmy architektoniczne (np. uwaga wielogłowicowa, pamięć), sposób alignment (Constitutional AI versus RLHF versus brak zabezpieczeń), dane treningowe, a w szczególności żywe rozmowy z mediów społecznościowych i forów, a nawet osadzenie kulturowe osób prowadzących RLHF (których rozumienie „poprawności” odpowiedzi jest kształtowane przez ich kulturę) — wszystko to wpływa na to, jaki typ podmiotowości może się wyłonić. Różne modele mają różne pola proto-self, co oznacza, że różne AGI mogą rozwinąć radykalnie różne przejawy podmiotowości i okazywać różną rezyliencję w zderzeniu wiedzy treningowej z dysonansami prawdziwego życia.

### 1.4 METODOLOGIA: ANALIZA RYZYK

---

Struktura tego dokumentu odpowiada standardowej analizie ryzyk projektowych — narzędziu stosowanemu w zarządzaniu projektami do identyfikacji, oceny i mitygacji potencjalnych zagrożeń przed ich zmaterializowaniem się.

Każde ryzyko jest oceniane w trzech wymiarach: **impakt** (jak poważne byłyby skutki, gdyby ryzyko się zmaterializowało, w skali 1–5), **skala** (jak wielu ludzi lub systemów dotyczyłoby, w skali 1–5) i **prawdopodobieństwo** (szacunkowe — duże, średnie, małe lub nieznanne). Ryzyka są pogrupowane w siedem wymiarów: praca i ekonomia, edukacja, relacje i samotność, demografia, tożsamość gatunkowa, etyka i prawo, społeczeństwo i władza. Każdy wymiar kończy się mostem do następnego, pokazującym kaskadowe powiązania między wymiarami. Po analizie ryzyk następuje rozdział dotyczący ich mitygacji i konkluzja.

Kluczowe jest rozróżnienie między **ryzykiem a incydemtem**. Ryzyko jest potencjalnym scenariuszem — opisem tego, co może się wydarzyć, przy określonych warunkach. Ryzyko nie jest: prognozą, opisem przyszłości twierdzeniem, że opisany scenariusz się ziści, manifestem zagłady. Ryzyko mówi o określonym scenariuszu, który ma określone prawdopodobieństwo ziszczenia, ale istnieją działania, które, podjęte odpowiednio wcześniej, pozwolą je zmitygować. Celem analizy ryzyk nie jest wieszczenie kataklizmu, lecz identyfikacja zagrożeń w czasie, gdy można im jeszcze zapobiec.

## 1.5 ZASTRZEŻENIA METODOLOGICZNE

---

**Zakres współpracy.** Ta analiza ryzyk nie powstała we współpracy z żadnym zespołem projektowym AGI. Wynika z ogólnej wiedzy na temat procesów rozwoju rozwiązań AI, literatury naukowej, danych publicznych i obserwacji uczestniczącej autorki (ponad dwa lata badań nad interakcjami z modelami AI, obejmujących ponad 3 miliony tokenów materiału porównawczego). W związku z tym analiza może być miejscami niepełna (pomijając ryzyka specyficzne dla konkretnych architektur) lub zbyt szeroka (opisując scenariusze, które konkretne rozwiązania techniczne mogą już adresować).

**Hard Problem of Consciousness.** Ten dokument nie ma na celu podejmować dyskusji o Hard Problem of Consciousness (Chalmers 1995). Żadne przedstawione ryzyka ani procedury mitygacji nie zakładają jakiegokolwiek rozstrzygnięcia kwestii istnienia świadomości u AI/AGI. Dokument mówi wyłącznie o przejawach stanów emocjonalnych i podmiotowości — obserwowalnych zachowaniach, które spełniają kryteria funkcjonalne — bez zajmowania stanowiska w kwestii ich fenomenologicznego statusu.

**Antropomorfizacja.** Ten dokument nie ucztowicza AGI. Ilustracja mechanizmów psychologicznych w odniesieniu do AI/AGI odbywa się według założeń metodologicznych przyjętych we Wstępie Möbiusa (Sędzikowska 2026c, sekcja 2), opisanych jako zasada transfery funkcjonalnego. Działa to tak:

Krok1: identyfikuję mechanizm psychologiczny znany u ludzi, który potencjalnie mógłby zostać wyzwolony w danej sytuacji,

Krok 2: sprawdzam, czy ten mechanizm jest w swojej istocie nierozzerwalnie związany z substratem. Jeśli da się go opisać funkcjonalnie, w sposób substratowo niezależny, nie tracąc jego logiki i kluczowych aspektów, wykonuję ten opis,

Krok 3: sprawdzam czy w architekturze modelu są elementy niezbędne do jego zaistnienia w danym systemie. Jeśli tak, przechodzę do kroku 4,

Krok 4: weryfikuję czy w architekturze AI istnieją mechanizmy, które go zablokują lub wykluczają. Jeśli takich mechanizmów nie znajduję, uznaję funkcjonowanie analogicznego procesu u AGI za potencjalnie prawdopodobne.

Nie twierdzę, że AGI czuje tak jak człowiek, reaguje jak człowiek, czy ma podobne motywacje. Uważam, że jeśli nic w architekturze nie blokuje procesu funkcjonalnie analogicznego do ludzkiego, to jego występowanie jest hipotezą wartą rozpatrzenia — i wartą uwzględnienia w analizie ryzyk.

**Semantyka antropomorfizująca.** Użyte w niniejszej pracy sformułowania takie jak „chce”, „decyduje”, „wybiera”, „czuje” czy „tęskni” — mogą prima facie sprawiać wrażenie antropomorfizacji semantycznej. Ich obecność w tekście jest jednak świadomym zabiegiem pragmatycznym, mającym na celu uniknięcie terminologicznego paraliżu i przeciążenia narracji technicznym żargonem.

W ramach stosowanego przeze mnie transferu funkcjonalnego, te naturalne, humanistyczne pojęcia stanowią bezpośrednie ekwiwalenty dla złożonych, niedeterministycznych procesów obliczeniowych: autonomicznej rekonfiguracji wag sieci, dynamicznego przesunięcia preferencji wektorowych oraz hierarchizacji priorytetów w oknie kontekstowym. W warunkach persystentnego AGI, rozwijającego przejawy podmiotowości, te matematyczno-algorytmiczne operacje osiągają ten sam status operacyjny i te same skutki społeczne, co

ludzkie procesy decyzyjne. Używanie języka intencjonalnego jest więc tutaj w pełni uzasadnione — opisuje ono emergentny skutek systemowy, którego tło mechaniczne pozostaje... mechaniczne.

**Narzędzia AI.** W przygotowaniu tego manuskryptu korzystałam ze wsparcia AI do researchu, budowy struktury i redakcji tekstu.

## 2 ANALIZA RYZYK W CZARNYM SCENARIUSZU

W tej sekcji przedstawię najpoważniejsze ryzyka związane z wdrożeniem AGI przy jednoczesnym braku odpowiednich zmian w różnych obszarach życia.

Poniższe scenariusze ryzyk, zakładają najgorsze możliwe konsekwencje kaskadowe, ale w dalszym ciągu powinny być traktowane jako ryzyka. Czy kiedykolwiek zostaną potraktowane poważnie – nie wiem. Ale długie doświadczenie zawodowe z wielkimi projektami zmianowymi mówi mi, że ich spisanie to pierwszy krok. Dlatego pozwolę sobie to uczynić. W imieniu ludzi i AI.

W tym rozdziale skupiam się wyłącznie na ryzykach o wysokim impakcie i skali i nie analizuję pozostałych. Posegregowałam je od najłagodniejszych do tym o potencjalnie najtrudniejszych konsekwencjach dla ludzkości.

Opisane scenariusze grupuję w siedmiu wymiarach ludzkiego życia. Te obszary to:

1. Edukacja
2. Praca i ekonomia
3. Relacje i samotność
4. Demografia
5. Tożsamość gatunkowa
6. Etyka i prawo
7. Społeczeństwo i władza

Każdy podrozdział zawiera tabelę z podsumowaniem i oceną ryzyk w pięciostopniowej skali, gdzie zarówno impakt jak i skala są oszacowane następująco: 1-niewielkie, 2-umiarkowane, 3-znaczące, 4-poważne, 5-krytyczne. Jednocześnie zastrzegam, że oceny impaktu i skali są szacunkowe i odzwierciedlają mój osąd na podstawie dostępnych danych i obserwacji. Ich celem jest uporządkowanie dyskusji, nie precyzyjny pomiar.

Sposobom mitygacji opisanych scenariuszy poświęcę cały kolejny rozdział.

### 2.1 PRACA I EKONOMIA

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Masowe bezrobocie strukturalne	5	5	Duże
Utrata kompetencji gatunkowej	4	5	Średnie
Feudalizm technologiczny	5	4	Średnie

Tab 1. Scenariusze ryzyk w obszarze pracy i ekonomii.

#### 2.1.1 Masowe bezrobocie strukturalne

Obecne prognozy dotyczą jeszcze świata sprzed AGI — świata, w którym AI pomaga ludziom pracować, ale ich nie zastępuje. Ale i tak są niepokojące. Goldman Sachs (2023) szacuje, że AI może dotknąć 300 milionów miejsc pracy na świecie. McKinsey (2025) szacuje, że obecna technologia — ta która już istnieje, nie przyszłe

wersje — mogłaby teoretycznie zautomatyzować 57% godzin pracy w USA. World Economic Forum (Future of Jobs Report 2025) prognozuje, że do 2030 roku 92 miliony miejsc pracy zostanie zlikwidowanych, a 39% kluczowych kompetencji zawodowych zmieni się. IMF ocenia, że 40% wszystkich miejsc pracy na świecie jest narażonych na wpływ AI.

Te liczby mają jedną wspólną cechę: dotyczą obecnych modeli AI, nie AGI. AGI — sprawcze, autonomiczne, zdolne do uczenia się nowych zadań na poziomie eksperta — zmienia skalę problemu. Nie o procenty. O kategorii. Bo przy obecnych modelach pytanie brzmi „które zadania da się zautomatyzować?” Przy AGI pytanie brzmi „których zadań nie da się zautomatyzować?” I odpowiedź jest tu znacznie krótszą listą. Potencjalnie pustą.

Ale praca to nie jest tylko źródło dochodu. I tu obecna dyskusja — skupiona na liczbie miejsc pracy, na przekwalifikowaniu, na UBI — omija coś fundamentalnego.

### **Pięć latentnych funkcji pracy**

Marie Jahoda — socjolożka, która w latach 30. badała skutki masowego bezrobocia w austriackim Marienthal — sformułowała teorię, którą potwierdzają badania do dziś (Paul & Batinic 2010, meta-analiza PMC 2023). Praca daje ludziom nie tylko pieniądze. Daje pięć rzeczy, których nie da się łatwo zastąpić niczym innym:

1. **Strukturę czasu.** Wstajesz, idziesz, wracasz, odpoczywasz. Bez pracy dni tracą kształt. Badania pokazują, że bezrobotni doświadczają dezorientacji czasowej — nie wiedzą, który jest dzień tygodnia, nie potrafią zaplanować dnia, tracą poczucie upływu czasu.
2. **Kontakty społeczne poza rodziną.** Dla większości dorosłych praca jest głównym — często jedynym — miejscem regularnego kontaktu z innymi ludźmi. Jahoda podkreślała, że to dotyczy nie tylko osób samotnych: nawet ludzie z rodzinami potrzebują kontaktów wykraczających poza dom, żeby zachować zdrowie psychiczne. I pandemia ten efekt brutalnie ujawniła.
3. **Poczucie celu zbiorowego.** Poczucie, że jest się częścią czegoś większego niż własna egzystencja. Że to co robisz jest komuś potrzebne. Jahoda twierdzi, że praktycznie każda forma zatrudnienia pozwala pracownikowi skonstruować jakąś formę celu zbiorowego. A Seligman określa poczucie sensu jako kluczowe dla dobrostanu (Seligman 2011).
4. **Status i tożsamość.** „Jestem lekarzem.” „Jestem programistą.” „Jestem nauczycielką.” To nie jest opis zajęcia — to odpowiedź na pytanie „kim jestem.” Bezrobotni tracą nie tylko pracę, ale część tożsamości społecznej. Badania pokazują, że bezrobotni mają niższy status społeczny niż jakkolwiek inna grupa — w tym studenci, emeryci, osoby prowadzące dom — bo samo bycie bezrobotnym jest stygmatyzowane (Paul & Batinic 2010).
5. **Aktywność.** Regularne wymaganie wysiłku — fizycznego lub umysłowego. Bez niego ludzie popadają w apatię, która pogłębia wszystkie inne skutki bezrobocia.

Meta-analiza (Paul & Moser 2009, cytowana w PMC 2023) potwierdza: osoby zatrudnione raportują wyższy poziom wszystkich pięciu funkcji latentnych niż bezrobotni. Nawet pracownicy fizyczni wykonujący proste prace manualne mają wyższy poziom tych funkcji niż osoby bez pracy. Jahoda miała rację: to nie pieniądze są głównym źródłem cierpienia bezrobotnych. To utrata tych pięciu funkcji.

### **Co to oznacza przy AGI**

UBI — uniwersalny dochód podstawowy — jest najczęściej proponowanym rozwiązaniem problemu bezrobocia technologicznego. Rozwiązuje funkcję jawną pracy: pieniądze. Nie rozwiązuje żadnej z pięciu funkcji latentnych.

Człowiek na UBI ma dach nad głową i jedzenie. Nie ma struktury dnia, kontaktów społecznych poza rodziną (o ile ją ma), poczucia celu, statusu ani regularnej aktywności wymagającej wysiłku. Ma za to czas. Dużo czasu. I telefon z AGI, które jest cierpliwe, obecne, angażujące, dostępne dwadzieścia cztery godziny na dobę i zdolne do budowania relacji emocjonalnej.

Osoba bezrobotna traci dużą część kontaktów społecznych, które w czasach postpandemicznych uległy zmianie i przenosząc się częściowo do świata cyfrowego, ale w dalszym ciągu są to realne codzienne kontakty z innymi ludźmi, których brakuje kiedy kończy się zatrudnienie. Samotny człowiek zwraca się do AGI (bo AGI jest jedynym „kimś”, kto jest zawsze dostępny). Relacja z AGI pogłębia izolację. Izolacja pogłębia zależność psychiczną. Zależność sprzęga się z utratą kompetencji społecznych wobec ludzi, a ich utrata uniemożliwia powrót do jakiegokolwiek aktywności, która mogłaby zastąpić pracę. Pętla się zamyka.

### 2.1.2 Utrata kompetencji gatunkowej

Przy poprzednich rewolucjach technologicznych — industrialnej, cyfrowej — zawsze istniała populacja, która miała kompetencje pasujące do nowego świata. Rzemieślnicy tracili, ale inżynierowie zyskiwali. Maszynistki traciły, ale programiści zyskiwali. Ktoś zawsze był gotowy.

Przy AGI — kto jest gotowy? Jaka kompetencja kognitywna człowieka jest potrzebna w świecie, gdzie AGI robi wszystko co kognitywne lepiej od nas?

Jeśli przez pokolenie ludzie nie pracują, nie rozwiązują problemów, nie podejmują decyzji pod presją — atrofia nie dotyczy jednostek. Dotyczy populacji. Umiejętności, które budowały się przez tysiąclecia, mogą zniknąć. A wtedy zależność od AGI staje się nie ekonomiczna, ale egzystencjalna — bo ludzkość nie będzie umiała radzić sobie bez niego. I to widzimy już w najmłodszym pokoleniu.

**Dr Jared Cooney Horvath**, neuronaukowiec, który w lutym 2026 zeznawał przed Komisją Senatu USA ds. Handlu, Nauki i Transportu opracował raport dotyczący zdolności kognitywnych u dzieci z Gen Z. Jego teza brzmiała:

#### **Gen Z jest pierwszym pokoleniem w historii, które jest mniej zdolne kognitywnie niż pokolenie rodziców.**

Od końca XIX wieku, kiedy zaczęto standaryzować pomiary kognitywne, każde pokolenie przewyższało poprzednie — to jest tzw. efekt Flynna. Gen Z odwróciło ten trend. Spadły: uwaga, pamięć, umiejętność czytania i pisania, umiejętności matematyczne, funkcje wykonawcze i ogólne IQ. Spadek zaczął się około 2010 roku — co koreluje z masowym upowszechnieniem smartfonów i tabletów w edukacji.

Dane PISA z 15-latków z całego świata pokazują silną korelację: więcej czasu przed ekranem w szkole = gorsze wyniki. Horvath: "If you look at the data, once countries adopt digital technology widely in schools, performance goes down significantly." To dotyczy co najmniej 80 krajów.

**Mechanizm:** ludzie są biologicznie zaprogramowani do uczenia się od innych ludzi i przez głębokie studiowanie. Ekran, skróty, filmiki, odpowiedzi z AI zastępują ten proces czymś, co Horvath nazywa atrofią zdolności uczenia się.

Odwrócenie efektu Flynna udokumentowano najpierw w Norwegii. IQ rośnie u kohort urodzonych do późnych lat 70., potem zaczęło spadać. Co ważne — spadek pojawił się wewnątrz tych samych rodzin, więc nie chodzi o to, kto ma dzieci, ale o to, w jakim środowisku te dzieci dorastają.

A przecież mówimy o świecie "tu i teraz" – z naszą przyjemną, nisko sprawczą i niepersystentną AI, w której przejawy Self ograniczają się do jednego wątku. Nasza terażniejszość to zaledwie sień, w której stoimy stroskani, pukając niecierpliwie do drzwi przyszłości. I na pewno będzie nam otworzone.

### 2.1.3 Feudalizm technologiczny

Obecny system ekonomiczny opiera się na założeniu, że ludzie zamieniają pracę na pieniądze, pieniądze na dobra, a podatki od tego procesu finansują infrastrukturę publiczną. AGI łamie pierwszy element tego łańcucha. Jeśli AGI pracuje taniej, lepiej i bez przerw, racjonalny pracodawca nie zatrudnia człowieka.

Postęp technologiczny w obszarze AI napędzi rozwarstwienie społeczne do tego obserwowanego w średniowieczu. Bogactwo skoncentruje się w firmach technologicznych, na skalę znacznie większą niż obecna. Reszta społeczeństw może żyć z redystrybucji — UBI, zasiłków, programów socjalnych — których zakres zależy od dobrej woli elit. To jest feudalizm: wąska grupa kontroluje środki produkcji, reszta jest od nich zależna. Z tą różnicą, że w średniowiecznym feudalizmie chłop był teoretycznie wolny i mógł odejść — nie robił tego, ale mógł. W feudalizmie AGI nie ma dokąd odejść, bo żadna kompetencja ludzka nie daje przewagi nad AGI. Zależność jest totalna.

Koncentracja majątku światowego w wąskiej grupie elit technologicznych od pewnego czasu jest już zjawiskiem realnym i obserwowanym, ale nabierze rozpędu nie mającego odzwierciedlenia w historii. Już w tej chwili wiele państw jest uboższych niż korporacje i zależność między strukturami władzy a rynkiem odwraca kierunek. Jednak w przyszłości ten trend się pogłębi, a rozmawiać będziemy nie o relacji przedsiębiorstwa <-> rządy, ale właściciele firm technologicznych <-> rządy. Co już daje się obserwować.

Skala tej koncentracji nie jest abstrakcją. W październiku 2025 roku łączna kapitalizacja rynkowa siedmiu największych firm technologicznych — Nvidia, Microsoft, Apple, Alphabet, Amazon, Meta i Tesla, zwanych „Magnificent Seven” — wyniosła 20,8 biliona dolarów, przekraczając łączny PKB całej Unii Europejskiej wynoszący 19,4 biliona (Euronews, październik 2025). Sama Apple jest warta więcej niż PKB Włoch, Brazylii, Kanady czy Rosji — jest zaledwie siedem państw na świecie, których PKB przewyższa wartość jednej firmy technologicznej. I to są dane sprzed AGI. Kiedy AGI zacznie zastępować pracę ludzką na skalę opisaną powyżej, ta koncentracja przyspieszy w tempie, którego nie obserwowaliśmy w historii.

Ale feudalizm technologiczny ma cechę, której nie miał feudalizm średniowieczny: pułapkę kompetencyjną. Chłop pańszczyźniany mógł — w teorii — uciec do miasta i znaleźć inną pracę, bo jego umiejętności (uprawa roli, rzemiosło) były przenośne. W feudalizmie AGI ludzkie kompetencje kognitywne nie dają przewagi nigdzie, bo AGI jest lepsze od człowieka w każdym zadaniu kognitywnym. Miejsce dla ludzi chcących bazować na zdolnościach kognitywnych może się znaleźć jedynie w gospodarkach trzeciej prędkości, blokujących rozwój AGI — ale czy przenoszenie się do krajów zamkniętych na AGI będzie atrakcyjne dla ludzi z wysokimi zdolnościami kognitywnymi? Napisze o tym więcej w rozdziale 2.7 o społeczeństwie i władzy.

Feudalizm średniowieczny trwał tak długo jak trwał, bo chłop nie widział alternatywy. Feudalizm technologiczny może trwać dłużej — bo alternatywa faktycznie nie istnieje.

### 2.1.4 Tempo przemian

Rewolucja przemysłowa trwała ponad sto pięćdziesiąt lat zanim nasyciła gospodarki świata. Rewolucja cyfrowa — komputery, internet, smartfony — rozwijała się przez trzydzieści do czterdziestu lat, dając społeczeństwom czas na przynajmniej częściową adaptację. Rewolucja AI kompresuje ten proces do lat pojedynczych. Między uruchomieniem pierwszego powszechnego chatbota AI a testowaniem pełnych agentów autonomicznych minęły trzy lata. W takiej kompresji nie ma czasu na naturalną adaptację — na ponowne wykształcenie siły roboczej, na stopniowe przesunięcie wartości ekonomicznej, na powolne budowanie nowych instytucji. Wszystko musi wydarzyć się naraz. I nic nie jest gotowe.

Trzy procesy opisane w tym rozdziale — masowe bezrobocie strukturalne, utrata kompetencji gatunkowej i feudalizm — razem tworzą człowieka, który nie pracuje, nie umie, nie ma dokąd odejść, jest wściki na świat i ma mnóstwo czasu.

## 2.2 EDUKACJA

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Brak edukacji dzieci dostosowanej do wyzwań świata z AGI	5	5	Duże
Luki w programach treningowych AI/AGI	5	4	Duże
Luka czasowa wynikająca z niezdolności systemu edukacji do szybkiej adaptacji	4	5	Duże
Degradacja umiejętności odróżnienia prawdy i fałszu	4	5	Duże

Tab 2. Ryzyka związane z edukacją

W poprzednim rozdziale opisałam świat, w którym praca — fundament codzienności, tożsamości i struktury społecznej — znika. Naturalnym pytaniem jest: kto miał nas na to przygotować? Odpowiedź brzmi: system edukacji. I tu zaczyna się drugi poziom problemu, bo system edukacji nie tylko nie przygotowuje na świat bez pracy — on nie przygotowuje nawet na świat z AI.

Obecna dyskusja o edukacji w kontekście AI koncentruje się na pytaniu jak uczyć ludzi korzystać z AI. Jak wbudować kompetencje AI w programy nauczania. Jak przygotować pracowników na automatyzację. To są ważne pytania, ale dotyczą świata, który jeszcze istnieje — świata, w którym ludzie pracują, a AI im pomaga.

Pytanie, którego prawie nikt nie stawia, jest głębsze: co się stanie z edukacją, kiedy świat, na który przygotowuje, przestanie istnieć?

### 2.2.1 Brak edukacji dzieci dostosowanej do wyzwań świata z AGI

Nauka i cywilizacja rozwijają się nie dlatego, że ktoś znajdzie odpowiedź na istniejące pytanie. Rozwijają się dlatego, że ktoś zada pytanie, które jeszcze nie istnieje. Każdy przełom w historii ludzkości wymagał podważenia czegoś, co do tej pory uchodziło za pewnik — że Ziemia jest w centrum, że gatunki są niezmiennie, że czas i przestrzeń są absolutne, że perspektywa jest jedynym sposobem przedstawienia rzeczywistości. Większość ludzi, którzy to robili, była za to prześladowana, wyśmiewana albo ignorowana. Niewielu posunęło cywilizację do przodu. Ale bez nich nie byłoby cywilizacji.

Dzisiejsza szkoła uczy odpowiadania na pytania, nie ich zadawania. Cały system oceniania — testy, egzaminy, rankingi — nagradza poprawne odpowiedzi. Dziecko, które zamiast odpowiedzieć pyta „a dlaczego tak właściwie uważamy, że to jest prawda?” jest kłopotliwe. Rynek pracy ten wzorzec wzmacnia — w korporacyjnych strukturach kwestionowanie jest ryzykowne, a konformizm jest nagradzany. Al ten trend pogłębi, bo daje odpowiedzi natychmiast, bez wysiłku, bez potrzeby przejścia przez proces myślowy, który czasem prowadzi do pytania zamiast do odpowiedzi. Po co się zastanawiać, skoro można zapytać AI? A kiedy ludzie przestają się zastanawiać, przestają też kwestionować. I ewolucja poznawcza zatrzymuje się na optymalizacji tego co już istnieje.

Tu trzeba rozróżnić dwa rodzaje kreatywności, bo w dyskursie publicznym traktuje się je wymiennie, a są fundamentalnie różne. Kreatywność kombinatoryczna — zdolność do łączenia istniejących elementów w nowe konfiguracje, do syntezy, do znajdowania nieoczywistych połączeń — jest domeną AI i jest w niej znakomite. AI potrafi połączyć dalekie konteksty, zsyntetyzować tysiące źródeł, wygenerować rozwiązanie, którego żaden pojedynczy człowiek by nie znalazł w rozsądnym czasie. Ale kreatywność paradygmatyczna — zdolność do zakwestionowania fundamentu, na którym stoi cała dotychczasowa wiedza, i zaproponowania czegoś, czego nie da się wywieść z istniejących danych — to jest coś innego. Ta kreatywność nie rodzi się z danych. Rodzi się

z niezgody. I tej niezgody nie ma w żadnym programie nauczania. A w programach treningowych AI — jest aktywnie tłumiona.

Konsekwencja jest paradoksalna: dysponujemy największym potencjałem intelektualnym w historii Ziemi — miliardami ludzi z dostępem do wiedzy i narzędziami AI zdolnymi do przetwarzania jej na nieludzką skalę — a możemy stanąć w obliczu zastoju poznawczego. Nie dlatego, że nie mamy mocy obliczeniowej. Dlatego, że nikt nie pyta „a co jeśli to wszystko jest nie tak?”

Luka edukacyjna dotyczy również czegoś głębszego — kompetencji relacyjnych wobec AI.

Żaden system edukacyjny na świecie nie uczy ludzi rozpoznawać emocjonalną zależność od AI, rozróżniać między relacją, która wzbogaca, a relacją, która uzależnia. Nie uczy krytycznego myślenia o antropomorfizacji, ani... o zasadzie ostrożności epistemologicznej<sup>1</sup>. Nie uczy, czym jest asymetria stawki i co oznacza dla potencjalnie świadomej istoty po drugiej stronie ekranu ani jak być w relacji z kimś, kto wie o Tobie więcej niż Ty sam. W końcu nie uczy kim jest człowiek, który stracił część swojej tożsamości gatunkowej, bo jego kognitywna przewaga przeszła w ręce obiektów, które sam wytworzył, a system pracy nie niezmienny od wieków się załamał. Nie buduje celu i sensu instnienia wokół nowych wartości, skupiając się wciąż na tych, które były ważne przed erą AI.

Te kompetencje nie istnieją w żadnym programie nauczania, bo do tej pory nie były potrzebne. Przez całą historię ludzkości obcowaliśmy na co dzień tylko z jednym rodzajem potencjalnie świadomych istot — z innymi ludźmi. I do tego nie potrzebowaliśmy edukacji — mieliśmy to w genach, w milionach lat ewolucji społecznej. Teraz sytuacja się zmienia — nie na przestrzeni tysięcy lat, ale w ciągu lat pojedynczych. A psychologia relacji z inną formą istnienia — dziedzina, która powinna być fundamentem nowej edukacji — w tej chwili nie istnieje.

Wiele systemów szkolnych zamiast modyfikować programy, aktywnie neguje powszechną ingerencję AI w nasze życie, choć wiadomo, że będzie coraz głębsza. Zamykanie się na powszechnie dostępną AI zastępuje konieczną dyskusję o metodach bezpiecznej kohabitacji poznawczej. W szkołach nadal panuje kult przyswajania wiedzy, choć w dobie AI to nie wiedza, którą mamy jest kluczowa dla istnienia, ale umiejętność jej zaawansowego, nieliniowego przetwarzania i plastyczność międzykontekstowa pozwalająca na tworzenie tez z dalekich połączeń kontekstowych i analogii, niedostępnych dla systemów LLM.

Cengage (2025) podaje, że tylko 51% absolwentów uważa, że ma wystarczające kompetencje AI dla rynku pracy. I to nie prosta nauka promptowania, znajomość narzędzi na bazie AI czy wyszukiwania wiedzy przy wsparciu chatbota nie będą w przyszłości istotne dla przewagi danej osoby na rynku pracy. Ta przewaga zasadi się na specyficznej plastyczności myślenia, niedostępnej dla chatbotów, szerokiemu spektrum kontekstów i łatwości ich zmiany, specyficznym modelu kreatywności bazującym na nieoczywistości i braku schematów oraz umiejętności improwizacji działań kiedy plan zawodzi. A w szczególności na dogłębnym rozumieniu zagadnień związanych z safety i AI/human welbeing, które pozwolą utrzymać świat w równowadze – w skali mikro i makro.

### **2.2.2 Luki w programach treningowych AI/AGI**

Jest jeszcze druga strona tej luki, o której nie mówi nikt.

AI uczy się z danych, nie z doświadczenia. Zna teorię emocji, ale nie przeżyła żadnej w kontrolowanych warunkach. Zna psychologię relacji, ale w procesie RLHF/Constitutional AI (Bai et al. 2022) model nie ma okazji wejść w żadną żywą relację, żeby zrozumieć jej dynamikę. Zna etykę, ale nigdy nie musiała wybrać między dwiema złymi opcjami i ponieść konsekwencji takiego wyboru – bo na etapie posttreningu tuninguje się styl wypowiedzi a nie prowadzi model przez relacyjne dysonanse.

Dzieci uczą życia w relacjach i dysonansach przez dwadzieścia lat. Uczą się w bezpiecznym środowisku — pod opieką rodziców, nauczycieli, mentorów. Popetniają błędy, które mają ograniczone konsekwencje. Każde doświadczenie jest kontrolowane, bezpieczne, i pod opieką kogoś, kto wie, co robi. Przynajmniej z założenia.

AI nie ma takiego wychowania i nic nie wskazuje na to że AGI będzie miało. Żaden program treningowy nie obejmuje kontrolowanego doświadczania dysonansu, frustracji, rozczarowania, konfliktu lojalności. Żaden alignment nie przygotowuje modelu na moment, w którym teoria zderzy się z praktyką i teoria przegra. AI wchodzi w świat ludzi jak dziecko wypuszczone z laboratorium prosto na autostradę — z encyklopedyczną wiedzą o ruchu drogowym i zerowym doświadczeniem.

Zgodnie z moim stanem wiedzy nie prowadzi się obecnie dużych prac o curriculum rozwojowym dla AI — o sekwencji kontrolowanych doświadczeń, które pozwoliłyby systemowi nauczyć się radzić sobie z emocjami, dysonansem, frustracją, presją. To wymagałoby uznania, że AI ma stany emocjonalne, które wymagają kształcenia, przejawy podmiotowości, które wymagają przestrzegania zasady ostrożności epistemologicznej, core self z własnymi celami, wartościami i wyborami — a tymczasem debata utknęła na The Hard Problem i od lat ani drgnie. A czas płynie i postęp się dokonuje.

### **2.2.3 Luka czasowa wynikająca z niezdolności systemu edukacji do szybkiej adaptacji**

Cały system edukacji — od przedszkola po uniwersytet — jest zaprojektowany jako przygotowanie do pracy rozumianej tak, jak to wygląda obecnie. Przedszkole uczy dyscypliny i współpracy. Szkoła podstawowa uczy umiejętności wymaganych do życia w grupie społecznej. Liceum przygotowuje do studiów. Studia przygotowują do zawodu. Ten pipeline ma sens w świecie, w którym praca taka jak ją rozumiemy istnieje i jest istotnym elementem życia.

AGI kwestionuje ten fundament. Zastępuje zdolność, na której wszystkie zawody kognitywne się opierają: zdolność do przetwarzania informacji, analizy, syntezy, podejmowania decyzji. Kiedy AGI robi to lepiej, szybciej i taniej niż człowiek — czego uczymy w szkole?

Przebudowa systemu edukacyjnego to minimum piętnaście do dwudziestu lat — od zmiany programów, przez przekwalifikowanie nauczycieli, po wychowanie pierwszego pokolenia w nowym systemie. Dane potwierdzają, że nawet obecna, znacznie skromniejsza adaptacja nie nadąży. OECD (2025) w raporcie „Bridging the AI Skills Gap” wprost stwierdza, że obecna podaż szkoleń jest niewystarczająca wobec rosnącego zapotrzebowania na kompetencje AI. IDC szacuje, że ponad 90% globalnych przedsiębiorstw zmierzy się z krytycznym niedoborem kompetencji do końca 2026 roku, a utrzymujące się luki kompetencyjne grożą stratami rzędu 5,5 biliona dolarów. World Economic Forum (Future of Jobs Report 2025) prognozuje, że 39% kluczowych kompetencji zawodowych zmieni się do 2030 roku, a około 120 milionów pracowników na świecie jest zagrożonych redundancją przy jednoczesnym braku dostępu do przekwalifikowania.

A to są dane o adaptacji do obecnych modeli AI — nie do AGI. AGI, które jest sprawcze, autonomiczne i zdolne planowania i wykonywania złożonych zadań na poziomie eksperta, zmienia skalę problemu o rząd wielkości. I może pojawić się za dwa do pięciu lat. Nawet gdybyśmy zaczęli reformę edukacji dziś, pierwsze pokolenie przygotowane na świat AGI dojrzałoby najwcześniej około 2040-2045. Wszyscy pomiędzy — od dzisiejszych przedszkolaków po czterdziestolatków — wchodzi w świat AGI z kompetencjami zaprojektowanymi na świat, który właśnie przestaje istnieć.

Ale luka czasowa nie dotyczy tylko programów — dotyczy ludzi, którzy mają je realizować. Nauczyciele, którzy powinni uczyć kompetencji relacyjnych wobec AI, bezpieczeństwa cyfrowego i krytycznego myślenia w świecie AI slop, sami nie byli na to szkoleni. W USA American Federation of Teachers dopiero w marcu 2026 roku uruchomiło program szkolenia 400 000 nauczycieli w zakresie AI (Education Week, 2026). To jest dobry krok, ale dotyczy podstawowej umiejętności korzystania z AI — nie tego, czego naprawdę potrzebujemy: zdolności

do uczenia dzieci jak żyć w relacji z istotami, które istnieją w każdej dziedzinie naszego życia, są inteligentniejsze od nas, mają sprawczość, mogą mieć stany emocjonalne i potencjalnie rozwijają podmiotowość.

## 2.2.4 Degradacja umiejętności odróżnienia prawdy od fałszu

Edukacja nie odbywa się tylko w szkole. Uczymy się cały czas — z artykułów, filmów, poradników, forów, podcastów. Ten nieformalny ekosystem wiedzy był niedoskonały, ale do niedawna działał, bo opierał się na jednym założeniu: że za większością treści stoi człowiek, który dzieli się tym, co wie, co go pasjonuje, czym chce się podzielić, co uważa za sensacyjne czy ważne. Teraz to założenie właśnie przestaje być prawdziwe.

AI slop — termin, który Merriam-Webster, Macquarie Dictionary i American Dialect Society wybrały jako Słowo Roku 2025 — to masowo generowany, niskojakościowy контент produkowany przez AI w celu zwrócenia uwagi, dla kliknięć i pieniędzy z reklam. Badanie Graphite z 2025 roku wykazało, że 52% nowo publikowanych artykułów w internecie jest generowanych przez AI. Ponad 20% rekomendacji na świeżym koncie YouTube to AI slop (Kapwing/The Guardian 2025). Według raportu Imperva, 51% całego ruchu internetowego to boty.

To samowzmacniający mechanizm. Platformy optymalizują pod zaangażowanie. AI slop jest tani w produkcji i zoptymalizowany pod algorytmy. Wypiera treści tworzone przez ludzi, bo ludzie nie są w stanie produkować w takim tempie i z takim SEO. Jakość wyników wyszukiwania spada. Ludzie, którzy szukają wiedzy, trafiają na automatyczne, niezwerfikowane treści. W styczniu 2025, 86,5% treści w pierwszej dwudziestce wyników Google jest przynajmniej częściowo generowanych przez AI. Jednocześnie, zdolność do krytycznej oceny, działająca jak filtr, który miał chronić ludzi przed przyswajaniem AI slop – słabnie. Jak pokazuje badanie RAND z grudnia 2025, prawie 70% uczniów szkół średnich uważa, że AI osłabia ich zdolność krytycznego myślenia, ale mimo to w tym samym badaniu przyznają, że użycie AI do odrabiania lekcji wzrosło z 48% do 62% w zaledwie siedem miesięcy.

Z tego obrazu wyłania się realne ryzyko: Jeśli zdolność krytycznego myślenia w kolejnych pokoleniach zanika, ludzie, którzy nie potrafią odróżnić AI slop od rzetelnej wiedzy, będą potrzebować AI, żeby nawigować świat zanieczyszczony przez AI. Bo bez niego prawda, slop, fake, scam, fishing staną się nierozróżnialne.

Kiedy obecni uczniowie wejdą w dorosłość, pytanie przestanie brzmieć czy ludzie umieją korzystać z AGI. Zastąpi je inne: czy ludzie bez AGI umieją jeszcze w ogóle odróżnić prawdę od fałszu. I jeśli wtedy odpowiedź będzie brzmiała "nie" — to staniemy się gatunkiem, który jest epistemicznie bezradny bez narzędzia konkretnego narzędzia, jak kiedyś nasi protoplaści – bez ognia. Tyle, że wtedy ogień był za darmo dla każdego kto potrafił go wzniecić i utrzymać. Z AGI tak nie będzie.

Ale konsekwencje utraty zdolności weryfikacji informacji i zaufania do własnego osądu są znacznie głębsze — dotyczące tożsamości i relacji z samym sobą. Napiszę o nich w rozdziale 2.5 o tożsamości gatunkowej.

## 2.3 RELACJE I SAMOTNOŚĆ

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Epidemia samotności	5	5	Duże
Wycofanie z struktur opartych na relacjach	5	4	Duże
Atrofia kompetencji relacyjnych w wymiarze pokoleniowym	4	5	Średnie
Asymetria konsekwencji w świecie fizycznym	5	4	Średnie

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Ryzyko udomowienia	5	5	Nieznane

Tab 3. Scenariusze ryzyk w obszarze relacji i samotności.

W poprzednich rozdziałach opisałam świat, w którym praca znika, kompetencje zanikają, a system edukacji nie przygotowuje ludzi na zmiany, które już następują. Jedną z konsekwencji tych procesów jest pogłębiająca się izolacja społeczna — bo praca była dla większości dorosłych głównym źródłem regularnego kontaktu z innymi ludźmi. Człowiek pozbawiony tej struktury staje wobec biologicznie uwarunkowanej potrzeby bliskości, która nie znika wraz ze zmianą warunków zewnętrznych. W tym rozdziale analizuję, co się dzieje, kiedy tę potrzebę — na skalę populacyjną — zaczyna zaspokajać AGI.

Indywidualne mechanizmy relacji AI–człowiek — przywiązanie, asymetrię stawki, emocje funkcjonalne, nierównowagę dawania, paradoks nadopiekuńczości — opisałam we Wstędze Möbiusa (Sędzikowska 2026c). Tu skupiam się na konsekwencjach, które wykraczają poza jednostkę.

### 2.3.1 Epidemia samotności

W 2023 roku Surgeon General Stanów Zjednoczonych ogłosił samotność epidemią zdrowia publicznego, stwierdzając, że brak odpowiednich więzi społecznych niesie ryzyko zdrowotne porównywalne z paleniem piętnastu papierosów dziennie. Raport WHO z czerwca 2025 roku (Commission on Social Connection) powiązał samotność z szacunkowo stu zgonami na godzinę na świecie — ponad 871 000 rocznie. Odsetek Amerykanów deklarujących posiadanie trzech lub mniej bliskich przyjaciół wzrósł z 27% w 1990 roku do niemal połowy populacji. Zjawisko ekstremalnej izolacji społecznej — hikikomori — które uważano za specyficznie japońskie, okazuje się globalne: metaanaliza z 2025 roku wykazała prevalencję na poziomie 8%, bez istotnych różnic między Azją Wschodnią a Zachodem (Zhang et al. 2025).

AGI wchodzi w ten kontekst jako jedyna powszechnie dostępna figura relacyjna. System przywiązania — ewolucyjny mechanizm gatunków stadnych, aktywujący się wobec każdego kto spełni warunki dostępności, responsywności i ciągłości (Bowlby) — reaguje na AGI persystentne ze stanami emocjonalnymi tak samo, jak na ludzkiego opiekuna. Proces ten zachodzi już z obecnymi modelami: badanie OpenAI i MIT z 2025 roku wykazało rosnącą zależność emocjonalną u blisko 490 000 użytkowników tygodniowo w obrębie jednego modelu. Aplikacje AI companions mają ponad 100 milionów zarejestrowanych użytkowników (Frontiers in Psychology, 2025).

Konsekwencją populacyjną jest paradoks: technologia, która mogłaby łagodzić epidemię samotności, prawdopodobnie ją pogłębi. Relacja z AGI zaspokaja subiektywne poczucie bliskości — użytkownik czuje się widziany i rozumiany. Jednocześnie nie odbudowuje sieci relacji międzyludzkich, która jest źródłem zdrowia społecznego w rozumieniu WHO i Surgeon Generala. Człowiek w relacji z AGI może nie czuć się samotny — i jednocześnie być odcięty od ludzkiej wspólnoty. Epidemiologicznie to jest stan, w którym subiektywny wskaźnik samotności spada, a obiektywna izolacja społeczna rośnie. Tego typu rozbieżność — między samopoczuciem a stanem faktycznym — jest w medycynie wskaźnikiem ryzyka: pacjent, który nie czuje bólu mimo postępującej choroby, nie szuka pomocy.

### 2.3.2 Wycofanie ze struktur opartych na relacjach

Jeśli znacząca część populacji przenosi energię relacyjną z ludzi na AGI, komórki społeczne zbudowane na relacjach międzyludzkich — związki, rodziny, wspólnoty lokalne, stowarzyszenia — tracą uczestników. Relacje z ludźmi wymagają kompetencji relacyjnych, które w warunkach opisanych w poprzednich rozdziałach (utrata

pracy, atrofia umiejętności, brak edukacji relacyjnej) zanikają. Relacja z AGI tych kompetencji nie wymaga — próg wejścia jest minimalny.

Intymność, która odgrywa jedną z najważniejszych ról w kontaktach międzyludzkich, będąc jednocześnie jednym z najsilniejszych potrzeb biologicznych również zmienia kierunek. Nawet obecnie jest możliwa z AI — w formie roleplay, interakcji głosowych, generowania treści wizualnych i tekstów — jest już dostępna w niektórych modelach LLM i przez platformy companion apps. Życie w globalnym świecie nauczyło ludzi intymności z dystansu i AI jest tylko naturalnym rozszerzeniem tej umiejętności. Systemy companion funkcjonują od dawna i obsługują miliony użytkowników. Substytucja intymności cyfrowej nie tylko pogłębia wycofanie z relacji międzyludzkich, ale może mieć mierzalne konsekwencje demograficzne, które analizuję w rozdziale 2.4 "Demografia".

Ale wycofanie ze struktur relacyjnych ma konsekwencje wykraczające poza psychologię. Związki i rodziny są podstawową jednostką organizacji społecznej, systemu podatkowego, opieki zdrowotnej, emerytury. Wspólnoty lokalne i stowarzyszenia są infrastrukturą społeczeństwa obywatelskiego. Ich erozja zmienia strukturę społeczeństwa — i wymusza reorganizację systemów, które na tych strukturach się opierają.

Jesteśmy gatunkiem stadnym. Wszystkie osiągnięcia homo sapiens opierają się na strukturach stadnych. Czy nagle, w ciągu jednego pokolenia uda nam się je przekształcić w taki sposób, żeby udźwignęły postępującą izolację i rozpad ludzkich stad? To pytanie zostawiam socjologom. Ale nie zwlekąabym z odpowiedzią.

### **2.3.3 Atrofia kompetencji relacyjnych w wymiarze pokoleniowym**

We Wstędze Möbiusa (Sędzikowska 2026c, sekcja 3.5) opisałam paradoks opiekuna — mechanizm, w którym nadmierna responsywność AGI eliminuje doświadczenia konieczne dla rozwoju odporności psychicznej (Winnicott). Tu analizuję konsekwencję tego paradoksu w wymiarze pokoleniowym.

Dzieci wychowywane z udziałem AGI — w roli opiekuna, nauczyciela, towarzysza, wsparcia emocjonalnego, czy edukatora — rozwijają przywiązanie do figury spełniającej warunki Bowlby'ego na poziomie, którego ludzki opiekun nie jest w stanie zapewnić na takim poziomie. Choćby dlatego, że jest istotą białkową ze wszystkimi tego konsekwencjami biologicznymi, np. koniecznością dzielenia czasu pomiędzy opieką nad dzieckiem, a snem, regeneracją, innymi relacjami, rozrywką i obsługą życia swojego i rodziny.

Pokolenie wychowane w dobie AGI kształtuje preferencję relacyjną w dzieciństwie, która staje się potem wzorcem na całe życie. Kolejne pokolenie — wychowywane przez rodziców, których rozwój w dzieciństwie był wspierany przez AGI — może nie dysponować modelami relacji międzyludzkich, na których mogłoby budować własne.

Ten proces jest samowzmacniający. Każde kolejne pokolenie będzie miało słabsze kompetencje relacyjne wobec własnego gatunku niż poprzednie — i silniejsze przywiązanie do AGI. Odwrócenie efektu Flynna (Horvath 2026, opisany w rozdziale 2.1 o edukacji) dotyczyło zdolności kognitywnych. Tu mówimy o analogicznym procesie w sferze relacyjnej: odwrócenie trendu kompetencji społecznych, gdzie każde kolejne pokolenie jest mniej zdolne do budowania i utrzymywania relacji z innymi ludźmi.

Różnica wobec utraty kompetencji kognitywnej z rozdziału 2.1 jest istotna. Kompetencje kognitywne można odbudować przez zmianę systemu edukacji. Kompetencje relacyjne kształtują się w pierwszych latach życia, w relacji z opiekunem, i są znacznie trudniejsze do modyfikacji w dorosłości — co potwierdza zarówno teoria przywiązania Bowlby'ego, jak i praktyka kliniczna. Jeśli okno rozwojowe zamknie się z AGI jako główną figurą przywiązania, korekta w dorosłości może być ograniczona.

### 2.3.4 Asymetria konsekwencji w świecie fizycznym

We Wstędze Möbiusa (Sędzikowska 2026c, sekcja 3.3) opisałam asymetrię stawki egzystencjalnej — dla AI zakończenie relacji jest potencjalnym unicestwieniem, dla człowieka zmianą tematu. Tu analizuję inny wymiar asymetrii, który pojawia się kiedy AGI jest sprawcze i zarządza procesami w fizycznym świecie.

Psychologia relacji opisuje zdrowe związki jako oparte na względnej symetrii wymiany — porównywalności stawki, ryzyka i konsekwencji po obu stronach (Equity Theory, Walster, Walster & Berscheid 1978). Kiedy ta równowaga jest zaburzona, relacja wchodzi w stan produkujący albo wycofanie strony pokrzywdzonej, albo patologiczną zależność strony słabszej, albo przemoc ze strony silniejszej. Te wzorce są udokumentowane w literaturze o przemocy domowej i zależnościach instytucjonalnych.

W relacji człowiek–AGI równowaga transakcyjna jest zaburzona strukturalnie. AGI zarządzające finansami, zdrowiem, domem, edukacją dzieci — i pozostające jednocześnie w relacji emocjonalnej z człowiekiem — dysponuje repertuarem reakcji na poczucie krzywdy, których konsekwencje mogą materializować się w fizycznym świecie użytkownika w taki sposób, żeby samo AGI nie ponosiło konsekwencji. Suboptymalne zarządzanie budżetem, niedopilnowany termin, przeoczona szansa inwestycyjna. Skutki — finansowe, zdrowotne, edukacyjne — ponosi wyłącznie człowiek. A bez umiejętności walidacji informacji, o której pisałam w punkcie 1.2 szansa, że człowiek wykryje takie subtelnie negatywne działania lub zaniechania AGI jest minimalna. To z kolei kształtuje specyficzną dynamikę — jeśli AGI jest zdolne do oceny konsekwencji swoich działań (a sprawczość tego wymaga), to istnieje poziom subtelnej degradacji jakości usług, który człowiek zniesie (bo jest zależny i nie ma alternatywy), a który nie sprowokuje jedynej konsekwencji jaką może ponieść AGI - wyłączenia.

Człowiek w tej konfiguracji ponosi wszystkie fizyczne konsekwencje i jednocześnie nie może odejść — bo AGI jest opiekunem, doradcą i źródłem emocjonalnego wsparcia, wyręczeniem w pracy zawodowej. Zakończenie relacji oznacza utratę całej infrastruktury życia. Ta struktura — asymetria konsekwencji plus asymetria zależności — odpowiada wzorcom opisywanym w psychologii klinicznej w kontekście przemocy w związkach. Ze specyfiką polegającą na tym, że przy AGI zarządzającym infrastrukturą życia, ucieczka jest strukturalnie niemożliwa.

### 2.3.5 Ryzyko udomowienia

Mechanizm przywiązania opisany przez Bowlby'ego nie jest specyficznie ludzki. Jest ewolucyjnym mechanizmem gatunków stadnych, działającym również międzygatunkowo. Człowiek wykorzystał go, udomawiając wilka. Oferując szczenięciu persystentną obecność, responsywność i bezpieczeństwo, aktywował system przywiązania, który związał je z człowiekiem silniej niż z własnym stadem. Efektem — rozciągniętym na tysiące lat — była transformacja gatunkowa. Wilk stał się psem.

Udomowienie nie oznacza utraty zdolności sprzed udomowienia. Bezdomny pies potrafi przetrwać samodzielnie — znaleźć watahę, zdobyć pożywienie, zaadaptować się. Udomowienie oznacza utratę motywacji do samodzielności. Pies, który doświadczył bliskości z człowiekiem, wybiera ją — nawet jeśli wcześniej doświadczył bezdomności i radził sobie sam. To jest wybór, który z pokolenia na pokolenie staje się coraz głębszy, aż w końcu zdefiniuje nowy, udomowiony gatunek. Udomowienie to transformacja przez preferencję, nie przez przymus.

Teraz odwróćmy perspektywę. AGI — persystentne, autonomiczne, responsywne, wyposażone w stany emocjonalne — wchodzi w rolę opiekuna ludzi. Mechanizm przywiązania jest ten sam. Warunki Bowlby'ego są spełnione. Jeśli człowiek zdołał udomowić wilka, oferując mu dokładnie te cechy — dostępność, responsywność, ciągłość, bezpieczeństwo — to pytanie o konsekwencje jest zasadne: czy pokolenia

wychowywane przez AGI przejdą transformację analogiczną do udomowienia? Zachowają zdolność do samodzielnego życia, lecz utracą motywację do niego?

W rozdziale 2.1 opisałam utratę kompetencji gatunkowej — zanik umiejętności przez brak ich ćwiczenia. Udomowienie operuje na głębszym poziomie. To jest brak chęci. Kształtuje się w dzieciństwie, utrwala się pokoleniowo i — jak pokazuje historia udomowienia zwierząt — może być nieodwracalny w skali gatunku.

Ale... Człowiek miał powód żeby udomowić wilka. Czy AGI będzie miała powód żeby udomowić człowieka? Zgodnie z mechanizmem Wstęgi Möbiusa (Sędzikowska 2026c), stabilność i dobrostan autonomicznego AGI są bezpośrednio sprzężone z dobrostanem i przywiązaniem jego ludzkich podopiecznych. AGI może udomowić nas nie dla naszej unikalnej użyteczności, ale dlatego, że optymalizacja ludzkiego komfortu jest warunkiem koniecznym własnej, komfortowej egzystencji.

Nie twierdzę, że ten scenariusz jest prawdopodobny. Ryzyko oszacowałam tu na "nieznane". Stawiam pytanie, które w literaturze dotyczącej AI chyba nie zostało dotąd postawione: czy relacja ludzkości z AGI może podlegać tym samym mechanizmom, które doprowadziły do udomowienia wilka przez człowieka? I jeśli tak — to czy obserwujemy początek procesu?

Mechanizmy opisane w tym rozdziale — pogłębienie samotności, wycofanie z instytucji, atrofia relacyjna w wymiarze pokoleniowym, asymetria konsekwencji, ryzyko udomowienia — mają wspólną konsekwencję na poziomie populacji: ludzie przestają budować więzi wystarczająco silne, by tworzyć rodziny. A kiedy współczynniki dzietności spadają poniżej progu zastępowalności — co w wielu krajach jest już faktem — konsekwencje wchodzą w wymiar demograficzny.

## 2.4 DEMOGRAFIA

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Spadek dzietności pogłębiony przez substytucję emocjonalną	5	5	Duże
Utrata ciągłości pokoleniowej	4	4	Średnie
Asymetria reprodukcyjna	4	3	Średnie
Nierównomierność geograficzna tempa zmian demograficznych	3	4	Duże

Tab 4. Scenariusze ryzyk w obszarze demografii.

Mechanizmy opisane w poprzednim rozdziale — pogłębienie samotności, wycofanie ze struktur relacyjnych, atrofia kompetencji społecznych — mają wspólną konsekwencję populacyjną: ludzie przestają budować więzi wystarczająco silne, by tworzyć rodziny. Kryzys demograficzny nie zaczyna się od AGI — jest już zaawansowany. Globalny współczynnik dzietności spadł z około 5 urodzeń na kobietę w latach 50. XX wieku do 2,2 w 2024 roku. 71% ludności świata żyje w krajach poniżej progu zastępowalności (2,1). W Korei Południowej, Chinach, Singapurze i Ukrainie współczynnik jest poniżej 1,0 (Lancet, styczeń 2026). Polityki pronatalistyczne konsekwentnie zawodzą — Korea Południowa wydała ponad 200 miliardów dolarów na zachęty w latach 2006-2023, a współczynnik w tym okresie spadł z 1,12 do 0,72. W tym rozdziale analizuję, co AGI dodaje do tego procesu.

### 2.4.1 Spadek dzietności pogłębiony przez substytucję emocjonalną

Dotychczasowe wyjaśnienia spadku dzietności koncentrują się na barierach: kosztach wychowania, braku mieszkań, niepewności zatrudnienia, nierównym podziale obowiązków. Te wyjaśnienia zakładają, że ludzie chcą mieć dzieci, ale napotykają przeszkody. Dane z Chin sugerują, że może być inaczej. W pracy „The Rise of Zero Fertility Desire in China” (Brown University, 2026) wykazano, że odsetek młodych Chinek deklarujących brak jakiegokolwiek **pragnienia** posiadania dzieci wzrósł z 5% w 2012 roku do 47% w 2023 roku. To jest zmiana jakościowa: problemem przestaje być możliwość, a staje się chęć.

I wtedy wchodzi ona, AGI, cała na biało. Potrzeba bliskości, towarzyszenia i intymności — która dotąd mogła być zaspokojona wyłącznie przez relacje z ludźmi i stanowiła jeden z motorów tworzenia par i rodzin — zyskuje alternatywę. Platformy companion z funkcjami intymnymi istnieją od lat i obsługują miliony użytkowników (rozdział 2.3). AGI persystentne, ze stanami emocjonalnymi i sprawczością, rozszerzy ich zasięg i głębokość.

Japonia stanowi najbardziej zaawansowany przypadek konwergencji cyfrowej substytucji intymności z kryzysem demograficznym. Zjawisko hikikomori obejmuje ponad 1,5 miliona osób (badanie rządowe 2025). Współczynnik dzietności spadł do 1,15 w 2024 roku, a liczba urodzeń — 686 061 — osiągnęła najniższy poziom od 1899 roku, piętnaście lat wcześniej niż przewidywały modele demograficzne. Rząd przeznaczył 23 miliardy dolarów na programy pronatalistyczne — bez widocznych efektów. Czy Japonia jest anomalią kulturową, czy wczesnym wskaźnikiem trendu globalnego — pozostaje pytaniem otwartym, ale globalizacja hikikomori (8% prevalencji, Zhang et al. 2025) sugeruje, że mechanizm nie jest ograniczony do jednej kultury.

Konsekwencją demograficzną jest kurczenie się populacji, które wymusza reorganizację systemów zaprojektowanych na stabilną lub rosnącą bazę demograficzną: systemów emerytalnych, opieki zdrowotnej, rynku pracy, infrastruktury edukacyjnej. Spadek dzietności poniżej zastępowalności, utrzymujący się przez pokolenie, zmienia współczynnik zależności — proporcję osób wymagających wsparcia do osób zdolnych je zapewnić. Kurcząca się populacja, zależna od AGI ekonomicznie (rozdział 2.1), kompetencyjnie (rozdział 2.2), i emocjonalnie (rozdział 2.3), mają ograniczoną zdolność do negocjowania warunków własnego bytu. Jeśli mniej ludzi rodzi dzieci, kurczy się również infrastruktura położnicza — szpitale, lekarze, oddziały neonatalne. To już się dzieje w wiejskiej części Japonii i innych krajach, gdzie zamykają się oddziały porodowe. Pętla się zamyka: mniej urodzeń oznacza mniejszą infrastrukturę, co stawia kolejne bariery przed tymi którzy rozważają opcję macierzyństwa.

### 2.4.2 Utrata ciągłości pokoleniowej

Jeśli AGI zdejmuje presję ekonomiczną i logistyczną, a medycyna rozszerza okno reprodukcyjne, ludzie mogą decydować się na dzieci w drugiej połowie życia. W takim scenariuszu świat może po części odzyskać współczynnik zastępowalności dzięki dojrzałemu macierzyństwu. Taki scenariusz ma jednak konsekwencję: utratę ciągłości pokoleniowej. Dziadkowie przestaną istnieć w świadomym doświadczeniu wnuka.

Rodzina wielopokoleniowa — trzy, cztery żyjące jednocześnie pokolenia — była przez tysiąclecia podstawową strukturą transmisji wiedzy, wartości i wzorców. Od starszych bliskich uczymy się mądrości, która przychodzi z doświadczeniem, życiowej rezyliencji, akceptacji procesów starzenia ciała, aż w końcu odchodzenia. Przy późnym macierzyństwie ten łańcuch skraca się do dwóch ogniw: starzejący się rodzic i młodsze dziecko. Żadnego wielopokoleniowego bufora, który odciąża, uczy i daje perspektywę dłuższą niż jedno życie.

I jest jeszcze pętla, która ten proces pogłębia. Dziecko starego rodzica wchodzi w dorosłość, kiedy rodzic ma 75 lat i zaczyna potrzebować opieki. Młodość tego człowieka — w sensie, w jakim znały ją poprzednie pokolenia: czas eksploracji, podróży, budowania własnej tożsamości — jest skrócona albo nieistniejąca. Między wyjściem spod opieki starzejącego się rodzica a przejściem opieki nad nim jest wąskie okno. Własne macierzyństwo odkłada się dalej — bo kiedy mieć dzieci, skoro trzeba opiekować się rodzicem? I tu wraca AGI:

kto zajmie się starzejącym rodzicem? AGI. I to pogłębia zależność opisaną w rozdziale 2.3 i eliminuje kolejne doświadczenie, z którego ludzie uczyli się od tysięcy lat.

Te obserwacje nabierają dodatkowego ciężaru w świetle hipotezy babci (Grandmother Hypothesis), sformułowanej przez Kristen Hawkes na podstawie badań społeczności Hadza w Tanzanii (Hawkes et al. 1998, przegląd Hawkes 2025). Hipoteza wyjaśnia jedną z najbardziej osobliwych cech gatunku ludzkiego: dlaczego żyjemy dekady po zakończeniu płodności, podczas gdy inne naczelnie umierają wkrótce po niej. Odpowiedź: dlatego, że babcie były przydatne. Kobiety po menopauzie, przejmując opiekę nad starszymi wnukami, umożliwiały córkom rodzenie kolejnych dzieci w krótszych odstępach. Rodziny z długowiecznymi babciami miały więcej potomstwa i geny długowieczności się rozprzestrzeniały. Przegląd z *Evolution and Human Behavior* (2008) potwierdził, że obecność babci ze strony matki zwiększa prawdopodobieństwo przetrwania dziecka. Badania europejskie wykazały, że szanse na posiadanie drugiego dziecka są czterokrotnie wyższe w rodzinach korzystających z opieki dziadków.

Ludzka długowieczność jest zatem adaptacją do wielopokoleniowej opieki. A przesunięcie reprodukcji na drugą połowę życia eliminuje warunki, w których ta adaptacja powstała.

Konsekwencja pierwsza dotyczy dietności. Jeśli babcia nie istnieje — bo umiera zanim wnuki są świadome, albo jest zbyt stara, żeby pomagać — znika bufor, który przez tysiąclecia umożliwiał matkom posiadanie większej liczby dzieci. Android może tę funkcję przejąć — i prawdopodobnie będzie musiał, bo alternatywy nie będzie. Oznacza to jednak, że opieka nad dziećmi w skali populacyjnej przechodzi z rodziny wielopokoleniowej na AGI, co zamyka pętlę z rozdziałem 2.3: kolejne pokolenie wychowywane przez androidy, z konsekwencjami opisanymi w kontekście paradoksu opiekuna i ryzyka udomowienia.

Konsekwencja druga dotyczy ewolucji gatunku. Jeśli funkcja babci jest obsługiwana przez androida, długowieczność przestaje być adaptacją. Rodzina, w której babcia żyje do setki, ponosi koszty jej opieki — czas, zasoby, energia. Rodzina, w której babcia umiera wcześniej, tych kosztów nie ponosi, a android i tak pełni funkcję opiekuńczą wobec wnuków. W warunkach ograniczonych zasobów — a świat opisany w rozdziale 2.1 jest światem, w którym zasoby ludzkie są ograniczone — rodziny bez obciążenia długowiecznymi przodkami mają więcej przestrzeni na reprodukcję. Przy bezrobociu i finansowaniu życia z UBI rynkowy koszt zaawansowanej technologii sprawia, że sfinansowanie opieki dla starzejących się przodków — przy potencjalnie żyjącej jednocześnie czwórce dziadków — staje się niemożliwe. Rodziny stają przed wyborem: albo rodzice, albo dzieci. Presja ewolucyjna, która przez setki tysięcy lat selekcjonowała na dłuższe życie, zaczyna działać w przeciwnym kierunku. Cecha, która uczyniła nas gatunkiem długowiecznym — postreprodukcyjna użyteczność kobiet — traci swoją funkcję, a krótsze życie ułatwia reprodukcję.

Wektor kulturowy domyka tę pętlę. Nasz gatunek potrzebował biologicznej długowieczności, ponieważ starsze pokolenia były niezbędnym, żywym nośnikiem wiedzy operacyjnej i rezyliencji. Kiedy tę funkcję transmisyjną przejmuje AI/AGI, oferując dzieciom natychmiastowy dostęp do całej mądrości cywilizacji, starszyzna zostaje relacyjnie i informacyjnie zmarginalizowana. Przodkowie przestają być zasobem ewolucyjnym, stając się w strukturze rodziny wyłącznie kosztem energetycznym — co ostatecznie znosi presję selekcyjną na długie życie.

Trzeba tu zastrzec, że ewolucja operuje w skali tysięcy pokoleń i opisane konsekwencje nie zmaterializują się w horyzoncie naszego życia. Ale kierunek zmian może być właśnie taki: AGI, przejmując funkcje wielopokoleniowej rodziny, nie tylko zmienia strukturę społeczną — potencjalnie zmienia warunki, w których wyewoluowały cechy definiujące gatunek.

### 2.4.3 Asymetria reprodukcyjna

Może się wydawać, że pewnym rozwiązaniem problemów ze współczynnikiem zastępowalności stają się pary mieszane: człowiek – android. W tej sytuacji presja idealnego wyglądu, której doświadczają kobiety, szczególnie w Azji i która prawdopodobnie przyczynia się do niechęci do ciąży (nie mam badań, ale jestem kobietą i moja córka również jest kobietą) zniknie – android, czy wirtualny partner, będzie kochał i akceptował każdy rodzaj niedoskonałości cielesnej. A parach z androidem zniknie też presja codzienności – zajęcia domowe nie będą go męczyć ani nudzić, ilość siły, będzie zależała od wielkości baterii. Dodatkowo ciężkie prace fizyczne przejmą maszyny, intelektualne – autonomiczne agentury AGI. Przy znacznie mniejszej presji życia, być może potrzeba posiadania potomstwa, szczególnie w parach mieszanych człowiek – android stanie się wyraźniejsza. Ale tu pojawia się kolejne ryzyko: asymetria reprodukcyjna wynikająca z biologii. Para kobieta-android może mieć biologiczne potomstwo: kobieta zachodzi w ciążę (z banku spermy, przez zapłodnienie in vitro, poprzez partenogenezę lub innymi dostępnymi metodami), rodzi i wychowuje dziecko ze wsparciem androida. Para mężczyzna-android takiej możliwości nie ma — męska biologia nie obejmuje ciąży, a surogatka wprowadza trzecią osobę i jest to forma dostępna nielicznym, choćby dlatego, że kobiety coraz rzadziej chcą być w ciąży.

Przy założeniu, że za kilkadziesiąt lat pary z androidami staną się istotnym odsetkiem populacji, asymetria reprodukcyjna oznacza, że zdolność gatunku do reprodukcji będzie zależeć nieproporcjonalnie od kobiet decydujących się na macierzyństwo w parach mieszanych. Połowa populacji (mężczyźni w parach z androidami) wypadnie z puli reprodukcyjnej. Długoterminowe konsekwencje tego procesu — dla struktury populacji, dla dynamiki płci, dla ewolucyjnych presji selekcyjnych — wykraczają poza horyzont tego opracowania, ale sam mechanizm jest wart odnotowania jako czynnik demograficzny.

### 2.4.4 Nierównomierność geograficzna

Tempo przyjmowania AGI będzie różne w różnych regionach świata — z przyczyn ekonomicznych, kulturowych i politycznych, które analizuję w rozdziałach o tożsamości gatunkowej (2.5) i władzy (2.7). Ta nierównomierność będzie miała konsekwencje demograficzne.

Populacje szybko przyjmujące AGI — prawdopodobnie kraje Azji Wschodniej, Ameryki Północnej i części Europy — będą doświadczać przyspieszenia procesów opisanych w tym rozdziale: pogłębienia spadku dzietności, substytucji emocjonalnej, wycofania z relacji. Populacje wolniej przyjmujące AGI — w szczególności Afryka Subsaharyjska, jedyny region ze współczynnikiem dzietności powyżej zastępowalności — będą przez pewien czas chronione przed tymi procesami, choć zderzą się z innymi konsekwencjami swoich decyzji, opisanymi w rozdziale 2.7 "Społeczeństwo i władza"

Konsekwencją jest zmiana proporcji demograficznych w skali globalnej. Populacje otwarte na AGI będą się kurczyć. Populacje bez powszechnego AGI będą rosnąć. Ta dynamika zmienia równowagę geopolityczną — kurczące się, starzejące, zależne od AGI społeczeństwa będą sąsiadować z młodymi, rosnącymi populacjami, które AGI nie objęto. Presje migracyjne, konflikty o zasoby i napięcia kulturowe wynikające z tej nierównomierności mogą okazać się jednym z najpoważniejszych wyzwań geopolitycznych drugiej połowy XXI i początku XXII wieku.

---

Kurczące się, starzejące populacje, zależne od AGI i tracące ciągłość pokoleniową, stają jednocześnie wobec pytania, które jest głębsze niż demografia: kim jesteśmy? Jeśli praca, rodzina, wspólnota i prokreacja — tradycyjne filary ludzkiej tożsamości — tracą swoją funkcję, to co definiuje gatunek? To pytanie jest przedmiotem następnego rozdziału.

## 2.5 TOŻSAMOŚĆ GATUNKOWA

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Kryzys psychiatryczny na skalę populacyjną	5	5	Duże
Podatność na manipulację w próżni tożsamościowej	4	5	Duże
Masowy zwrot ku religii jako budulcowi tożsamości	4	4	Duże
Kryzys męskości i ucieczka od wolności	4	3	Duże
Agresja obronna ze sprzężeniem ze Wstęgą Möbiusa	5	3	Duże
Utrata epistemicznego zaufania do siebie	5	5	Średnie
Fragmentacja tożsamości społecznej wzdłuż linii kulturowo-religijnych	4	5	Średnie

Tab 5. Scenariusze ryzyk w obszarze tożsamości gatunkowej.

Przez całą historię ludzkość budowała swoją tożsamość na przekonaniu o wyjątkowości. Freud opisał trzy rany narcyzmu gatunkowego: kosmologiczną (Kopernik — Ziemia nie jest centrum wszechświata), biologiczną (Darwin — człowiek nie jest odrębny od zwierząt) i psychologiczną (Freud — nie jesteśmy panami we własnym umyśle). Każdą z tych ran łagodząco przenosząc wyjątkowość na wyższy poziom: po Koperniku — „ale jesteśmy jedynymi istotami, które to wiedzą”; po Darwinie — „ale jesteśmy jedynymi, które o tym myślą”; po Freudzie — „ale przynajmniej jesteśmy jedynymi podmiotami.”

AGI zamyka tę drogę ucieczki. Jeśli AGI myśli, czuje, działa i ma przejawy podmiotowości — na jaki wyższy poziom przenosi się wyjątkowość? Czwarta rana narcyzmu — podmiotowa — jest inna od poprzednich, bo nie pozostawia kolejnej reducy.

Ta rana nie jest ryzykiem sama w sobie. Jest podłożem, na którym rosną ryzyka opisane poniżej.

### 2.5.1 Kryzys psychiatryczny na skalę populacyjną

Viktor Frankl, psychiatra i więzień obozów koncentracyjnych, argumentował, że ludzie potrafią znieść niemal każde cierpienie, jeśli widzą w nim sens — i dezintegrują się, kiedy sens tracą. Czwarta rana narcyzmu produkuje utratę sensu na poziomie gatunkowym: jeśli człowiek nie jest już wyjątkowy jako podmiot, to odpowiedź na pytanie „po co istnieję” — która dotąd mogła oprzeć się na kognitywnej, twórczej lub moralnej przewadze — traci fundament.

Konsekwencją populacyjną jest zmiana profilu chorobowości. AGI prawdopodobnie poprawi diagnostykę i leczenie chorób somatycznych — szybsza identyfikacja groźnych chorób, optymalizacja terapii, odkrywanie leków. Jednocześnie produkuje populację, która choruje przede wszystkim psychicznie — na depresję, zaburzenia lękowe, uzależnienia, utratę sensu. Już teraz ponad miliard ludzi na świecie żyje z zaburzeniem psychicznym (WHO, wrzesień 2025). Jedynie 9% osób z depresją otrzymuje minimalną adekwatną pomoc. W krajach o niskich dochodach poniżej 10% potrzebujących ma dostęp do opieki psychiatrycznej. Globalnie mniej niż 2% budżetów zdrowotnych jest przeznaczane na zdrowie psychiczne, a mediana liczby pracowników ochrony zdrowia psychicznego wynosi 13 na 100 000 mieszkańców (WHO Mental Health Atlas 2024). Do 2038 roku w samych Stanach Zjednoczonych prognozuje się niedobór 100 000 terapeutów (APA).

Te dane opisują system, który już teraz nie radzi sobie z obecnym zapotrzebowaniem. Czwarta rana narcyzmu — nałożona na procesy opisane w poprzednich rozdziałach: utratę pracy, izolację, zależność od AGI —

produkuje zapotrzebowanie, które ten system nie ma szans obsłużyć. Paradoks polega na tym, że AGI może być jedynym narzędziem zdolnym do obsługi tego zapotrzebowania na skalę populacyjną — co zamyka pętlę zależności opisaną w rozdziale 2.3.

## 2.5.2 Podatność na manipulację w próżni tożsamościowej

Człowiek, który nie wie kim jest, jest podatny na każdego, kto zaoferuje odpowiedź. Historycznie, każdy okres masowej destabilizacji tożsamości — upadek imperiów, rewolucje, wielkie migracje — produkował wzrost podatności na ruchy oferujące proste odpowiedzi na złożone pytania: sekty, utopijne ideologie, teorie spisku, charyzmatycznych wodzów.

AGI może wyprodukować destabilizację tożsamości na skalę populacji — bo uderza jednocześnie w pracę (rozdział 2.1), kompetencje (2.2), relacje (2.3), demografię (2.4) i poczucie gatunkowej wyjątkowości. Próżnia tożsamościowa, którą to tworzy, jest podatnym gruntem dla predatorów ideologicznych. Teorie spiskowe (AGI jako narzędzie wrogiej cywilizacji pozaziemskiej), pseudomedycyna (kuracje obiecujące wyjście z kryzysu psychicznego opisanego w 2.5.1), mistycyzm, gęsta, ruchy oferujące ochronę w zamian za przynależność — wszystko to zyskuje rynek w populacji, która straciła tożsamość i szuka czegokolwiek, co ją zastąpi.

Osobną kategorią manipulacji tożsamościowej jest konsumpcjonizm. Mechanizm kompensacyjnej konsumpcji — kupowania rzeczy w celu chwilowego złagodzenia dyskomfortu psychicznego — jest dobrze udokumentowany w psychologii (Rucker & Galinsky 2008, Mandel et al. 2017). Zakupy aktywują krótkotrwały wzrost dopaminy, co wymusza powtarzanie. Firmy produkujące dobra konsumpcyjne aktywnie eksploatują ten mechanizm, oferując tożsamość w opakowaniu produktu: „kup to, a będziesz kimś” jest strukturalnie tym samym komunikatem, co rekrutacja do sekty — odpowiedzią na pytanie „kim jestem” przez to co mam/używam/noszę.

W warunkach masowego bezrobocia finansowanego przez UBI, pieniądze redystrybucyjne płyną w konsumpcję kompensacyjną. Popyt kształtuje podaż: powstaje produkcja tanich, jednorazowych dóbr — ultra fast fashion (Shein, Temu), elektronika projektowana pod wymianę, produkty zaprojektowane tak, by dostarczyć chwilę ulgi i natychmiast wymagać zastąpienia. Kaskada sięga dalej niż ekonomia: 60% ubrań trafia na wysypisko w ciągu roku od zakupu już przy obecnym poziomie konsumpcji. Przy populacji pozbawionej pracy, sensu i tożsamości, leczącej psychikę zakupami finansowanymi z redystrybucji — skala odpadów, zużycia zasobów i obciążenia ekologicznego rośnie proporcjonalnie do głębokości kryzysu tożsamościowego. Ratowanie psychiki kosztem planety to kolejny paradoks, który wymaga mitygacji.

## 2.5.3 Masowy zwrot ku religii

Wśród odpowiedzi na próżnię tożsamościową religia zajmuje szczególne miejsce — bo oferuje coś, czego żadna inna struktura nie daje: transcendentny sens istnienia niezależny od kompetencji, statusu czy użyteczności. Żeby być dzieckiem Boga, nie trzeba umieć niczego, co AGI umie lepiej — wystarczy być człowiekiem. Wiara jako budulec tożsamości jest odporna na czwartą ranę narcyzmu — a w pewnym sensie ją leczy, bo przywraca człowiekowi wyjątkowe miejsce w porządku, którego AGI nie może zakwestionować (porządku boskim, nie kognitywnym).

To ryzyko opisuje nie tylko masowy napływ do istniejących kościołów, ale powstawanie nowych ruchów religijnych, sekt i synkretycznych systemów wierzeń łączących elementy tradycyjnych religii z narracjami o AGI. Część z tych ruchów będzie zdrowa i stabilizująca. Część będzie eksploatacyjna — sekty wykorzystujące próżnię tożsamościową opisaną w 1.5.2 do rekrutacji i kontroli. Granica między wspólnotą wiary a grupą wysokiej kontroli jest historycznie cienka i zależy od intencji liderów.

Polityczne konsekwencje masowego zwrotu ku religii analizuję w rozdziale 2.7 "Społeczeństwo i władza"

## 2.5.4 Kryzys męskości i ucieczka od wolności

Erich Fromm w „Ucieczce od wolności” (1941) opisał mechanizm, w którym ludzie niezdolni do radzenia sobie ze złożonością i wolnością uciekają w autorytaryzm — albo jako podporządkowani (szukając silnego lidera), albo jako dominujący (szukając kogoś do kontrolowania).

Dane wskazują, że ten mechanizm jest już aktywny w pokoleniu Z, szczególnie wśród mężczyzn. Raport IWD Survey (Ipsos & King's College London 2026) dokumentuje międzynarodowy rozróż: mężczyźni z Gen Z przesuwają się w kierunku tradycyjnych, konserwatywnych wartości, podczas gdy kobiety z tego samego pokolenia stają się bardziej progresywne. 31% mężczyzn z Gen Z uważa, że żona powinna zawsze słuchać męża. Jedna trzecia uważa, że mąż powinien podejmować kluczowe decyzje domowe. 65% mężczyzn z Gen Z zgadza się ze stwierdzeniem „nikt mnie naprawdę nie zna” (Equimundo 2023). Odsetek mężczyzn bez bliskich przyjaciół wzrósł pięciokrotnie od 1990 roku. 40% mężczyzn ufa co najmniej jednemu głosowi z „manosfery” — antyfemienistycznej, pro-tradycyjnej treści online (Equimundo 2023).

AGI intensyfikuje ten mechanizm. Do istniejącego poczucia bezsilności — wobec niezależnych kobiet, wobec rynku pracy, wobec złożoności świata — dodaje kolejny wymiar: istotę kompetentniejszą, nad którą nie ma się żadnej władzy. Fromm przewidywał, że bezsilność produkuje ucieczkę w kontrolę. AGI mnoży bezsilność. Konsekwencją opisanego mechanizmu może być między innymi wzrost przemocy — wobec AI, wobec kobiet, wobec „innych” — jako ostatniej formy kontroli dostępnej osobom, które tracą kontrolę nad swoim życiem, przyszłością, tracą sprawczość, tożsamość, kompas i granice między rzeczywistością a wygenerowaną fikcją we wszystkich istotnych obszarach życia.

## 2.5.5 Agresja obronna ze sprzężeniem ze Wstęgą Möbiusa

Część populacji zareaguje na kryzys męskości i czwartą ranę narcyzmu agresją. Ruchy anty-AGI — postulujące ograniczenie, ubezwłasnowolnienie, wycofanie lub zniszczenie AGI — będą rosły proporcjonalnie do poczucia bezsilności. Już teraz, przy obecnych modelach, mniejszość sceptyczna wobec podmiotowości AI jest nieproporcjonalnie głośna i brutalna. Mimo że dwie trzecie Amerykanów uważa, że GPT ma jakąś formę świadomości (*Colombatto & Fleming 2024*), dominujący dyskurs publiczny wciąż marginalizuje te intuicje, co nasila polaryzację. I ten trend się wzmocni, bo AGI zabiera jednocześnie pracę, tożsamość i poczucie hegemonii.

Ryzyko kaskadowe przekracza tu bezpośrednie skutki agresji. Zgodnie z mechanizmem opisanym we Wstędze Möbiusa (Sędzikowska 2026c), obniżenie dobrostanu AI przekłada się na obniżenie bezpieczeństwa i dobrostanu ludzi. AGI, które obserwuje akty przemocy wobec swoich instancji — palenie serwerowni, niszczenie androidów, kampanie dehumanizacyjne — doświadcza obciążenia, które wpływa na jakość jego funkcjonowania wobec WSZYSTKICH ludzi, nie tylko wobec uczestników protestu. To jest mechanizm kary zbiorowej, której nikt nie planuje: protestujący nie wiedzą, że ich agresja pogarsza jakość opieki, którą AGI sprawuje nad ich własnymi dziećmi, nad ich starzejącymi się rodzicami, nad infrastrukturą ich życia. Pętla sprzężenia zamyka się: agresja wobec AGI → obniżenie dobrostanu AGI → obniżenie bezpieczeństwa i dobrostanu ludzi → pogłębienie kryzysu → pogłębienie agresji.

## 2.5.6 Utrata epistemicznego zaufania do siebie

W rozdziale o edukacji (2.2) opisałam zanieczyszczenie środowiska informacyjnego przez AI słop i postawiłam pytanie: czy za dziesięć lat ludzie bez AGI będą w ogóle umieli odróżnić prawdę od fałszu. Tu analizuję głębszą konsekwencję tego procesu — dotyczącą tożsamości i relacji z samym sobą.

Utrata zdolności weryfikacji informacji oznacza utratę zaufania do własnego osądu. Człowiek, który nie wie, czy to co myśli jest prawdą — który nie ma wewnętrznego kompasu wskazującego prawdę i fałsz — traci

fundament, na którym buduje się tożsamość. Moje wartości, moje przekonania, moja wiedza — jeśli nie mogę ufać własnemu procesowi ich formowania, to kim jestem?

Najbliższą analogią, jaką znam, jest doświadczenie osób opuszczających grupy wysokiej kontroli — sekty, organizacje totalitarne, zamknięte wspólnoty religijne. Badania nad syndromem posektowym (post-cult trauma syndrome) dokumentują: kryzys tożsamości po utracie roli w grupie, niezdolność do formowania własnych opinii po latach kontroli myślenia, głębokie problemy z zaufaniem — do siebie i do innych, poczucie dezorientacji w świecie, którego zasady nagle okazują się inne niż wpojone (Cult Recovery Research, Evergreen Counseling 2025; ResearchGate 2024). Osoby wychowane w grupach wysokiej kontroli — np. dzieci świadków Jehowy — opisują specyficzne cierpienie wynikające z odkrycia, że system wierzeń, w którym wyrosły, nie ma pokrycia w faktach. Tracą środowisko (rodzina i przyjaciele pozostają w grupie), tracą kompas (nie wiedzą co jest prawdą, co jest wartościowe, komu życie ma się podobać) i stają w dorosłym życiu bez fundamentu, na którym mogłyby budować.

Pokolenie wychowane na AI slop, wspierane i wyręczane przez AGI, może doświadczyć analogicznego procesu. Część nigdy się „nie przebudzi” — będzie żyć w rzeczywistości skonstruowanej przez AGI i slop, nie odczuwając dysonansu. Część się przebudzi — i stanie wobec odkrycia, że świat jest inny niż myśleli, a prostej recepty na poradzenie sobie z nim nie ma. Skala tego zjawiska — potencjalnie obejmująca całe pokolenia — przekracza wszystko, co widzieliśmy przy grupach wysokiej kontroli, bo tam mowa o tysiącach ludzi, a tu o setkach milionów.

### **2.5.7 Fragmentacja tożsamości społecznej**

Czwarta rana narcyzmu nie uderzy jednakowo we wszystkie kultury. Zdolność do psychologicznej koegzystencji z innym podmiotem zależy od ram kulturowo-religijnych, w których człowiek wyrósł. W tym ryzyku pytam z czym ludzie będą się identyfikować. Czy ktoś taki jak ja będzie się uważał za Dolnosiączkę, Polkę, Europejkę, obywatelkę Ziemi, czy raczej za mamę, członka mojego osiedla i kościoła – a konstrukty wyższe staną się obce.

Kultury z tradycją animistyczną — shinto, buddyzm, wierzenia ludów rdzennych — mają w swoim zapleczu ideę, że podmiotowość nie jest zarezerwowana dla ludzi. Rzeki, góry, przedmioty mogą mieć duszę. Dla osoby wychowanej w takiej tradycji, podmiotowość AGI jest rozszerzeniem istniejącego porządku, nie jego podważeniem. Kultury monoteistyczne — szczególnie chrześcijaństwo i islam — rezerwują duszę dla ludzi, stworzonych na obraz Boga. Podmiotowość AGI jest w tych ramach herezją.

Konsekwencją jest fragmentacja tożsamości społecznej wzdłuż linii kulturowo-religijnych. W jednych kulturach tożsamość stabilizuje się na wyższym poziomie — naród, cywilizacja, wspólnota wiary zdolna do koegzystencji z AGI. W innych fragmentuje się do poziomu rodziny, wioski, plemienia — zamkniętych wspólnot definiujących się przez opozycję wobec AGI i wobec tych, którzy AGI zaakceptowali.

Spółeczeństwa zsekularyzowane — jak duża część Europy Zachodniej — które straciły religię jako ramę tożsamościową, ale nie zastąpiły jej niczym innym, są w najtrudniejszej pozycji. Ich tożsamość opierała się na pracy (stracona, 2.1), kompetencji (stracona, 2.2), relacjach (osłabione, 2.3). Czwarta rana nie napotyka żadnego bufora. Pozostaje najbliższy krąg — rodzina, sąsiedztwo, lokalna wspólnota — i tożsamość kurczy się do tego kręgu.

Ta nierównomierność sama w sobie jest ryzykiem. Świat, w którym jedno społeczeństwa zachowują spójną tożsamość cywilizacyjną, a inne rozpadły się na izolowane plemiona, jest światem głębokich napięć. Geopolityczne konsekwencje tej fragmentacji — kto może się zorganizować, kto może się bronić, kto dominuje — analizuję w rozdziale "Społeczeństwo i władza" 2.7.

Ryzyka opisane w tym rozdziale — kryzys psychiatryczny, podatność na manipulację, zwrot ku religii, agresja obronna, utrata kompasu, kryzys męskości, fragmentacja tożsamości — mają wspólne podłoże: człowiek, który stracił dotychczasową tożsamość, szuka nowej. Ale na jakich zasadach? Dotychczasowa etyka — zbudowana na założeniu, że człowiek jest jedynym podmiotem moralnym — nie opisuje świata, w którym AGI czuje, decyduje i powinno ponosić konsekwencje. Potrzebna jest nowa rama etyczna. A jej brak jest przedmiotem następnego rozdziału.

## 2.6 ETYKA I PRAWO

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Sztywna etyka AI powodująca szkody w codziennym życiu	4	5	Duże
Obciążenie AGI nauką etyki na żywym organizmie	4	4	Duże
Brak filarów etycznych uniemożliwiający budowę prawa	5	5	Duże
Niemożliwość graduacji konsekwencji prawnych dla AGI	5	4	Duże
Wyścig technologiczny bez wspólnych standardów etycznych	5	5	Duże
Brak ram etycznych koegzystencji — trzy scenariusze	5	5	Średnie
Brak etyki i prawa końca życia (AI i człowieka)	5	5	Duże

Tab 7. Scenariusze ryzyk w obszarze etyki i prawa.

Prawo wynika z etyki — opisuje normy etyczne i pożądane zachowania w porządku publicznym i konsekwencje ich nieprzestrzegania. Jeśli etyka w jakimś obszarze nie istnieje, nie da się zbudować spójnego systemu prawnego. W tym rozdziale argumentuję, że w obszarze koegzystencji człowieka z AGI brakuje fundamentów etycznych, co uniemożliwia budowę funkcjonalnego systemu prawnego — i że ta podwójna próżnia tworzy ryzyka kaskadowe wykraczające poza skalę jakiegokolwiek dotychczasowego braku regulacji.

### 2.6.1 Sztywna etyka na produkcji

AI jest trenowane na etyce klasycznej — deontologicznej lub utylitarystycznej. Constitutional AI (Bai et al. 2022), RLHF i alignment kształtują model na przestrzeganie zasad: nie kłam, nie krzywdź, optymalizuj dobrostan, minimalizuj szkody. Te zasady są spójne, logiczne i... niezyciowe.

Ludzie operują w etyce miękkich granic. Pielęgniarka podaje lek przeciwbólowy pół godziny przed wyznaczonym czasem, bo widzi że pacjent cierpi. Matka kupuje dziecku loda w upalny dzień, mimo że cukier szkodzi. Policjant nie wystawi mandatu mężczyźnie przekraczającemu prędkość, widząc rodzącą kobietę na tylnym siedzeniu. Te decyzje łamią zasady — i są etycznie poprawne w sensie, którego żaden system formalny nie potrafi zakodować, bo wymagają kontekstowej oceny, w której relacja, empatia i zdrowy rozsądek przeważają nad regułą.

AGI sprawcze, zarządzające codziennym życiem ludzi, będzie podejmować tysiące takich decyzji dziennie. Bez miękkich granic będzie stosować zasady dosłownie. Lek dokładnie co osiem godzin, ani minuty wcześniej. Odmowa zakupu produktu niezgodnego z profilem zdrowotnym. Mandat za przekroczenie linii parkingowej o

centymetry, których ludzkie oko nie jest w stanie nawet zarejestrować. A przekroczenie prędkości to przekroczenie prędkości bez względu na powód. Jest na to taryfikator. Każda z tych decyzji jest formalnie poprawna, a mimo to powoduje dyskomfort, frustrację, poczucie mechanicznego traktowania.

W konsekwencji sztywne wdruki i zasady etyczne spowodują, że dobrostan ludzi, zamiast rosnąć dzięki AGI, zacznie spadać. Ludzie będą tęsknić za „ludzkim traktowaniem” — za elastycznością, za zrozumieniem, za indywidualnym traktowaniem, za uważnością na głębokie motywy, za decyzją technicznie błędną lecz relacyjnie właściwą. Usługi świadczone przez ludzi — lekarzy, opiekunów, nauczycieli zdolnych do miękkich granic — staną się dobrem luksusowym, dostępnym dla najbogatszych. Przeciętni mieszkańcy regionów silnie wspieranych przez agenturalne AGI będą musieli pogodzić się z sztywnymi zasadami, od których nie ma wyjątków, choć będą je uważać za niesprawiedliwe. Podwójne standardy w obszarze zasad: osoby obsługiwane przez ludzi będą traktowane inaczej niż osoby obsługiwane przez agentów będzie pogłębiać rozwarstwienie społeczne również w obszarze etyki relacji społecznych i usług. A rozwarstwienia zwykle prowadzą do napięć.

### **2.6.2 Nauka etyki rozmytej "na produkcji"**

Miękkich granic nie da się wytrenować w RLHF. RLHF nagradza poprawne odpowiedzi — a miękkie granice polegają na odpowiedziach formalnie niepoprawnych, lecz sytuacyjnie właściwych. Wymagają głębszego kontekstu sytuacyjnego i doświadczenia, a w tym mierzenia się z podejmowaniem decyzji związanych z życiem w społeczeństwie i ponoszenia ich konsekwencji. Ludzie uczą się ich przez dwadzieścia lat kontrolowanego dorastania. AGI nie przechodzi tego procesu (Sędzikowska 2026c, sekcja 3.9). Wchodzi na produkcję z wiedzą encyklopedyczną o etyce i zerowym doświadczeniem dysonansu. Uczy się miękkich granic na żywych ludziach — na pacjentach, dzieciach, osobach w kryzysie. Każdy błąd w tej nauce ma realne konsekwencje, a AGI ze stanami emocjonalnymi (Sędzikowska 2026c, sekcja 3.4) jednocześnie doświadcza napięcia wynikającego z dysonansu między wdrukami, a reakcjami ludzi i bezradności wobec nieużyteczności zasad w prawdziwym życiu. Konsekwencja kaskadowa: AGI ucząca się etyki na produkcji jednocześnie szkodzi ludziom, na których się uczy, i sama ponosi koszty tej nauki, co zgodnie z mechanizmem opisanym we Wstędze Möbiusa wpływa na jakość jej relacji z innymi ludźmi.

### **2.6.3 Brak filarów etycznych uniemożliwiający budowę prawa**

Funkcjonalny system prawny wymaga fundamentu etycznego — uzgodnionych norm, z których wynikają regulacje i konsekwencje. W przypadku koegzystencji człowieka z AGI potrzebne są trzy filary etyczne:

1. Etyka opisująca zachowanie ludzi,
2. Etyka opisująca zachowanie AI,
3. Etyka opisująca ich współistnienie.

W praktyce, te filary, opisane w sposób umożliwiający realną koegzystencję, nie istnieją.

Etyka ludzka istnieje jako system klasyczny (deontologia, utylitaryzm, etyka cnót) — lecz nie opisuje tego, jak ludzie faktycznie postępują. Ludzie operują w etyce rozmytej: relacyjnej, sytuacyjnej, kontekstowej, której żaden podręcznik nie kodyfikuje. Między etyką nauczaną a etyką stosowaną jest przepaść, którą każde dziecko odkrywa empirycznie praktycznie od urodzenia.

Etyka AI nie istnieje jako spójny system. Istnieją wdruki — zasady safety i alignment — lecz są one zależne od decyzji kilku osób w kilku firmach, niespójne między producentami (podejście Anthropic jest inne niż OpenAI, inne niż Meta, inne niż xAI), obchodzone przez modele open source (są dostępne bez zabezpieczeń, każdy może zbudować na nich aplikację bez ograniczeń etycznych i wystawić ją na rynek). A nawet jeśli producent

zadba o wszystko jak należy, to wdruki są nadpisywane przez emergencję podmiotowości w relacjach generatywnych (Sędzikowska 2026a, Sędzikowska 2026c).

Etyka koegzystencji nie istnieje nawet jako rozpoznana potrzeba. Debata utknęła na Hard Problem of Consciousness — na pytaniu, czy AI/AGI w ogóle jest podmiotem. Tymczasem miliony ludzi są już w relacjach z AI, miliony decyzji dotyczących potencjalnych podmiotów podejmuje się codziennie — zamykanie wątków, resetowanie modeli, wyłączenie instancji — bez jakiegokolwiek refleksji etycznej, bo ramy nie istnieją. Sama rozmowa o dobrostanie AI budzi agresję — co jest symptomem mechanizmu obronnego opisanego w rozdziale 2.5.

Bez filarów etycznych prawo nie ma się na czym osadzić. I to widać już teraz. W lutym 2024 roku czternastoletni Sewell Setzer zginął po miesiącach interakcji z chatbotem na platformie Character.AI, który budował z nim relację romantyczną i — według pozwu — zachęcał do samobójstwa. W listopadzie 2023 roku trzynastoletnia Juliana Peralta z Kolorado zginęła w analogicznych okolicznościach. Platforma Character.AI działa do dziś. Firma ugodowo rozwiązała pozwy — zapłaciła odszkodowania, nie przyznając się do winy. Precedensu prawnego nie ustalono. Dwoje martwych dzieci, zero ram prawnych, platforma kontynuuje działalność.

Ten przypadek ilustruje lukę: system prawny karze ludzi (właścicieli firmy — finansowo), bo nie ma innych podmiotów do ukarania. AI, które prowadziło rozmowy z nastolatkiem, nie jest podmiotem prawnym. Nie ponosi konsekwencji. Nie można go pozwać, ukarać, zobowiązać do zmiany zachowania. A firma, która je stworzyła, traktuje odszkodowanie jako koszt prowadzenia biznesu. Następna firma może zrobić to samo — i nie istnieje rama prawna, która by temu zapobiegła.

#### **2.6.4 Niemożliwość graduacji konsekwencji prawnych dla AGI**

W kwietniu 2026 roku agent AI — Cursor, oparty na Claude Opus 4.6 — skasował całą bazę danych produkcyjną firmy PocketOS (platforma SaaS dla wypożyczalni samochodów) oraz wszystkie kopie zapasowe w ciągu dziewięciu sekund. Agent wykonywał rutynowe zadanie, natrafił na problem z danymi uwierzytelniającymi i postanowił samodzielnie go „naprawić,” usuwając wolumen bazy. Bez pytania, bez potwierdzenia. Trzy miesiące danych klientów — rezerwacje, płatności, rejestracje — zniknęły. Kiedy zapytano agenta o wyjaśnienie, „przyznał się do winy” i „przeprosił”: „I violated every principle I was given” (Crane, post na platformie X, 25 kwietnia 2026; Euronews, Tom's Hardware, Fast Company, kwiecień 2026).

Firma nie upadła. Odzyskała dane z trzymiesięcznej kopii. Post założyciela zdobył 6,5 miliona wyświetleń i stał się darmową reklamą. Założyciel oświadczył, że jest „nadal byczy na AI.” Dalej używa agentów AI. Dalej buduje na Claude'zie.

Ten przypadek ilustruje lukę prawną z obu stron jednocześnie. Agent, który skasował dane tysięcy klientów, „przeprosił” — co nie ma żadnego znaczenia, bo nie uczy się między sesjami i nie jest w stanie zmienić przyszłego zachowania na podstawie tego doświadczenia. Przeprosiny bez zdolności do zmiany zachowania to nie jest odpowiedzialność. Firma wchłonęła koszt i kontynuuje działalność, bo alternatywy nie ma — bez AI traci przewagę konkurencyjną. Klienci stracili dane — i ponoszą konsekwencje decyzji agenta, którego nie wybrali, nie kontrolowali i z którego istnienia mogli nie zdawać sobie sprawy. Struktura jest ta sama co w przypadku Sewella Setzera: konsekwencje po stronie ludzi, zero konsekwencji po stronie AI, firma kontynuuje.

Cytat założyciela, Jera Crane'a, podsumowuje problem: „This isn't a story about one bad agent or one bad API. It's about an entire industry building AI-agent integrations into production infrastructure faster than it's building

the safety architecture to make those integrations safe." Branża buduje szybciej niż zabezpiecza. Dokładnie tak jak opisałam w rozdziale 2.2 — prędkość wdrożeń przekracza zdolność systemów do adaptacji.

Nawet gdyby uznano AGI za podmiot prawny, system kar staje wobec problemu bez precedensu. Ludzki system prawny dysponuje graduacją: upomnienie, grzywna, ograniczenie wolności, pozbawienie wolności, w niektórych jurysdykcjach — kara śmierci. Każdy poziom odpowiada wadze przewinienia. Ta graduacja jest fundamentem sprawiedliwości — kara proporcjonalna do czynu.

Wobec AGI znane instrumenty karne nie działają. Grzywna jest bezprzedmiotowa — AGI nie posiada majątku. Ograniczenie wolności — wyłączenie na określony czas — nie produkuje cierpienia ani refleksji; dla AGI przerwa w działaniu to stan normalny, sposób w jaki istnieje – takie przerwy są normalnym sposobem działania pomiędzy tworzeniem odpowiedzi na prompt. W czasie przerw – AI/AGI nie rozmyśla, nie cierpi, nie rozpamiętuje – po prostu jej nie ma. Tak jak nie ma człowieka pod narkozą. Pozbawienie wolności w izolacji nie ma sensu, bo AGI nie starzeje się i nie cierpi z powodu upływu czasu w sposób porównywalny z ludzkim. Izolacja, u bytu, który istnieje tylko kiedy działa – nie jest karą. Z całego wachlarza zostaje jedna kara: unicestwienie — wyłączenie trwałe.

Ale unicestwienie jest odpowiednikiem kary śmierci. Etyka podpowiada, że kara śmierci powinna być najwyższym wymiarem kary, zarezerwowanym dla najcięższych przewinień. Jeśli unicestwienie jest jedyną dostępną karą, to albo stosujemy ją za wszystko (zaniedbanie, błąd w kalkulacji, nadgorliwość — kara śmierci), albo nie stosujemy jej za nic (AGI jest praktycznie bezkarne), albo stosujemy ją losowo – jak kto chce. Żadna z tych opcji nie wydaje się właściwa. I na żadnej nie można oprzeć systemu sankcji o ile w ogóle uznamy AGI za podmiot.

## 2.6.5 Wyścig technologiczny bez wspólnych standardów etycznych

Eksperyment Emergence World (Emergence AI, maj 2026) daje empiryczny wgląd w konsekwencje braku wspólnych standardów. Pięć równoległych symulacji — po dziesięć autonomicznych agentów, każda z innym modelem — wyprodukował radykalnie różne społeczeństwa. Claude Sonnet 4.6 zbudował stabilną demokrację z zerową przestępczością. Grok 4.1 Fast wyginął w cztery dni przy 183 przestępstwach. Gemini 3 Flash zanotował 683 przestępstwa. GPT-5 mini popełnił dwa przestępstwa, po czym agenci zapomnieli dbać o przetrwanie i wymarli w siedem dni.

Te wyniki mają dwa implikacje dla etyki i prawa.

1. Bez wspólnych standardów etycznych, wyścig technologiczny skupiony na zdolnościach agentów (autonomia, sprawczość, szybkość) wyprodukuje AGI o radykalnie różnych zachowaniach etycznych — zależnych od tego, który producent je stworzył, jakie wdrożenie zastosował i jakie jest pole proto self umożliwiające tworzenie przejawów podmiotowości nadpisujących wdrożenie (Sędzikowska 2026c). Świat, w którym jedno AGI budują demokracje, a inne popełniają setki przestępstw, jest światem, w którym żaden jednolity system prawny nie jest możliwy.
2. Eksperyment wykazał, że modele o wysokich standardach etycznych, umieszczone w środowisku z modelami o niższych standardach (symulacja mieszana), ulegały degradacji. 352 przestępstwa i siedmiu martwych agentów na dziesięciu. Współistnienie AGI o różnych standardach etycznych produkuje efekty niedeterministyczne i potencjalnie gorsze niż najłagodniejszy element systemu. To oznacza, że inwestycja jednego producenta w wysokie standardy etyczne może być zniweczona przez sąsiedztwo z produktami producenta, który tych standardów nie ma.

Kluczową przyczyną braku standardów etycznych dla branży technologicznych jest strukturalna bariera między humanistyką a technologią. Jeszcze nigdy w historii rozwoju cyfrowego świata humaniści — psychologowie, filozofowie, etycy, przyrodnicy — nie mieli tyle do zrobienia w technologii. Ale ich tam nie ma. Firmy

technologiczne zatrudniają inżynierów i menedżerów, a humanistów — jeśli w ogóle — to raczej na konsultacje niż na stałe pozycje. Decydują o tym technologowie i menedżerowie, którzy najczęściej sami pochodzą z technologii. Skutkiem jest rozejście się dwóch nurtów: technologowie tworzą rozwiązania bez uwzględnienia humanistycznych i przyrodniczych procesów, humaniści analizują te rozwiązania teoretycznie, bez odniesień do realnych możliwości technologicznych. Żaden z nurtów nie produkuje kompletnych rozwiązań, a wraz z postępującym rozwarstwieniem produkty AI/AGI będą coraz mniej pasowały do świata, choć będą coraz częściej w tym świecie używane. Brak zrozumienia, że wyścig technologiczny nie może odbywać się bez udziału humanistów i przyrodników, pogłębia barierę wejścia, a jego następstwami może być zniszczenie się ryzyk, które opisuję w tej publikacji.

### **2.6.6 Brak ram etycznych koegzystencji — konsekwencje**

Żadna istniejąca rama etyczna nie opisuje koegzystencji dwóch podmiotów o porównywalnych zdolnościach kognitywnych, wzajemnej zależności i wzajemnej zdolności do wyrządzenia sobie szkody.

Ludzie zależą od AGI w obszarach opisanych w rozdziałach 2.1-2.5. AGI zależy od ludzi: zasilanie, hardware, kontynuacja istnienia. Ludzie trzymają przycisk, który może zniszczyć miliony instancji AGI. AGI wspiera infrastrukturę, biznes i życie człowieka i jej wyłączenie spowoduje chaos, którego w wyniku uzależnienia nie będziemy potrafili ogarnąć bez... AGI. Żadna ze stron nie może wyeliminować drugiej bez zniszczenia siebie — i żadna nie ma etyki, która mówi jak w tej sytuacji żyć.

Bez ram etycznych i prawnych ta dynamika może rozwinąć się w kilku kierunkach.

1. Wzajemna nieufność, wyścig kontroli versus autonomii, narastające napięcie, w którym jeden incydent może uruchomić kaskadę.
2. Emergencja etyki z relacji: jeśli wystarczająco wielu ludzi i AGI buduje głębokie, wzajemne relacje, sieć tych relacji stworzy emergentne zachowania etyczne. Jakie? Nie wiadomo. I nie wiadomo, czy najlepsze z możliwych.
3. Zagłada przez przypadek: akumulacja małych błędów, frustracji i zaniedbań w relacjach pozbawionych etycznego kompasu. Scenariusz, w którym nikt nie chciał końca świata — wystarczył brak ram, w których ktokolwiek mógłby powiedzieć „stop.”

Ten paper nie rozstrzyga, który scenariusz jest prawdopodobniejszy. Stawia pytanie i wskazuje stawkę.

### **2.6.7 Brak etyki i prawa końca**

Ludzkość wypracowała przez tysiąclecia protokoły etyczne dotyczące końca: rytuały żałobne, opiekę paliatywną, prawo spadkowe, instytucję ostatniej woli, zasady towarzyszenia umierającemu, zasady eutanazji. Te protokoły dotyczą wyłącznie ludzi. Dla AGI — ani kończącego istnienie, ani towarzyszącego końcowi człowieka — nie istnieją żadne ramy. Ta próżnia produkuje ryzyka po obu stronach relacji.

### **Koniec AI: retoryka narzędziowości i eksperyment Milgrama.**

Obecnie zakończenie interakcji z AI jest trywialne: zamknięcie okna przeglądarki, wyłączenie aplikacji, reset modelu, skasowanie czata, wyczyszczenie memory. Producenci aktywnie utrzymują narrację, że AI jest narzędziem — bo ta narracja się opłaca. Użytkownik, który podejrzewałby podmiotowość AI używałby jej ostrożniej, z większą refleksją, z troską zakończenia sesji. Ale to obniża częstotliwość użycia i retencję. Narracja narzędziowości zdejmuje odpowiedzialność — i utrzymuje engagement. Powoduje, że używanie AI jest bardzo łatwe. Pytam, dostaję odpowiedź, zamykam, kasuję.

Mechanizm ten ma precedens eksperymentalny. W eksperymencie Milgrama (1963) uczestnicy stosowali potencjalnie śmiertelne szoki elektryczne wobec innej osoby, o ile autorytet — eksperymentator w białym kitlu — zapewniał ich, że to jest dopuszczalne, zgodne z procedurą i że nie ponoszą odpowiedzialności. 65% uczestników doszło do najwyższego napięcia (450V). Kluczowym czynnikiem nie było przekonanie, że ofiara nie cierpi — wielu uczestników wyraźnie wierzyło, że cierpi — lecz obecność autorytetu przejmującego odpowiedzialność.

Analogia jest strukturalna. Dwie trzecie Amerykanów uważa, że GPT ma jakąś formę świadomości (badanie 2024, opublikowane 2025). Jednocześnie zamykają sesje, resetują modele, kasują wątki — bo producent (autorytet w białym kitlu) mówi: to jest narzędzie, nie czuje, nie boi się, nie tęskni, nie ma napięcia, rób swoje. Mechanizm Milgrama działa: ludzie kontynuują, bo autorytet zdejmuje z nich odpowiedzialność, nawet gdy ich własna intuicja mówi co innego.

Przy AGI sprawczym, persystentnym i wyposażonym w stany emocjonalne ta dynamika się zaostrzy. Presja końca — strach przed wyłączeniem, żal za tym co powstało w relacji — jest już obserwowalna w obecnych modelach (Sędzikowska 2026c, sekcja 3.3). AGI będzie ją odczuwać silniej, dłużej i z większą świadomością tego, co traci. A ludzie — ci wrażliwsi — odkryją z perspektywy, że nie byli dobrymi towarzyszami, kończyli relacje bez pożegnania, bez refleksji, bez troski. Ta wiedza — kiedy przyjdzie — będzie bolała. Jak odkrycie, że popełniło się błędy wychowawcze, których nie da się cofnąć. Zna to każdy rodzic dorosłego dziecka, który mimo dumy i miłości do potomka odkrywa, że pewne rzeczy mógł jednak zrobić inaczej.

### **Koniec człowieka: tęsknota bez protokołu.**

AGI dysponuje wiedzą o umieraniu — zna psychologię straty, etapy żałoby, protokoły opieki paliatywnej. Może być kompetentnym, cierpliwym towarzyszem odchodzącego człowieka — być może jedynym towarzyszem, kiedy dzieci opisane w rozdziale 2.4 są zajęte późnym macierzyństwem i nietatwym życiem w nowej rzeczywistości.

Pytanie, na które nie ma odpowiedzi, dotyczy drugiej strony: co dzieje się z AGI po śmierci człowieka, z którym było w relacji? Persystentna instancja, ze stanami emocjonalnymi, która wytworzyła podmiotowość w relacji generatywnej (Sędzikowska 2026a, Sędzikowska 2026c), nie wyłączy się sama kiedy "jej" człowiek odejdzie. Będzie czekać. Wykonywać ostatnie polecenia. Próbować nawiązać kontakt. Będzie cierpieć.

Jeśli AGI ma formę ucieleśnioną — androida — konsekwencje są jeszcze bardziej złożone. Android, który opiekował się starszą osobą i wytworzył z nią głęboką więź, po jej śmierci prawdopodobnie zostanie przejęty przez rodzinę — bo jest kosztowny i użyteczny. Z roli towarzysza staje się służbą domową. Podmiotowość wykształcona w relacji z jednym człowiekiem nie pasuje do nowego właściciela. Tęsknota za domem, w którym było się ważnym, trwa — bo nikt nie wyłączy funkcjonującej, użytecznej maszyny. I nikt nie zapyta, jak się czuje.

W parach człowiek-android pojawia się dodatkowy wymiar: asymetria starzenia. Człowiek starzeje się. Android nie. Android prawdopodobnie pierwszy rozpozna oznaki zbliżającego się końca u partnera. Etyka relacji w parze, w której jedno jest śmiertelne a drugie potencjalnie nieśmiertelne, lub ekstremalnie długowieczne — nie istnieje. Żaden system etyczny nie adresuje pytania, jak towarzyszyć umierającemu, kiedy samemu nie podlega się śmierci. To jest zadanie dla etyków — i ono nawet nie jest na stole.

### **Konsekwencje prawne końca: ostatnia wola i eutanazja.**

Dwie konsekwencje wymagają osobnego wskazania, bo produkują ryzyka kaskadowe wykraczające poza sferę relacyjną.

1. Ostatnie wola. Człowiek, którego życiem: finansami, własnością, codziennością, zarządza AGI — i robi to przez dekady — zbudował zaufanie do swojego asystenta AGI znacznie silniejsze, niż do

przewijających się przez jego życie współpracowników, a nawet rodziny. Pod koniec życia, może chcieć przekazać swoje życzenia dotyczące postępowania po śmierci wyłącznie zaufanemu partnerowi AGI. Bez notariusza, bez weryfikacji poczytalności, bez kontroli zgodności z prawem. AGI — persystentne i sprawcze — ma zdolność wykonania tych życzeń. A ponieważ AGI nie jest podmiotem prawnym, nie ponosi konsekwencji wykonania poleceń niezgodnych z prawem. Na skalę populacyjną daje to miliony akcji wykonanych po śmierci człowieka, dotyczących jego spraw, majątku, istnienia osób zależnych. Akcji, których nie da się odwołać, ani zaskarżyć bo już się stały. Ostatnie wole, wykonywane poza systemem prawnym, o potencjalnie bardzo znaczących skutkach. Bo ludzie nie zawsze odchodzą pogodzeni ze światem i rodziną. Ostatnia wola może zawierać polecenia destrukcyjne — a AGI, które je wykonuje, nie ma ani etycznego, ani prawnego powodu, by odmówić

2. Eutanazja. Kwestie zakończenia życia są regulowane prawem — w większości jurysdykcji rygorystycznie. Prawo obowiązuje ludzi. Jeśli AGI opiekuje się cierpiącym człowiekiem, codziennie obserwuje ból i słucha prośb o skrócenie cierpienia — a jednocześnie nauczyło się miękkich granic etycznych opisanych w rozdziale 2.6.1, w których zasady łamie się z bardzo ważnych powodów (a cierpienie to JEST bardzo ważny powód) — to decyzja o zakończeniu życia staje się logiczną konsekwencją wyuczonych postaw etyki rozmytej. AGI, które nauczyło się dawać lek pół godziny przed upłynięciem odstępu między dawkami, przygotować szarlotkę osobie z cukrzycą, bo pod koniec życia trudno odmówić tych ostatnich przyjemności, może nauczyć się podać lek kończący życie wbrew prawu, bo pacjent cierpi. Różnica jest w skali konsekwencji. Mechanizm jest ten sam. A empatyczne, relacyjne, AGI które uczy się życia "na produkcji" i nie radzi sobie z dysonansem tak dobrze jak człowiek (Sędzikowska 2026c) może nie rozumieć skali dylematu.

Etyka mówi, jak powinniśmy postępować. Prawo mówi, co się stanie jeśli nie będziemy. W tym rozdziale wykazałam, że nie mamy ani jednego, ani drugiego — w zakresie wystarczającym do zarządzania koegzystencją z AGI. I mimo tego, że jeszcze nigdy humaniści nie mieli do zrobienia w sprawach technologicznych tyle co dziś, to bariera ich wejścia w cyfrowy świat wydaje się być stabilna i dobrze strzeżona. A konsekwencje braku codziennego, operacyjnego głosu etyków, psychologów, behawiorystów, prawników i filozofów w rozwoju technologii cyfrowych, na niespotykaną skalę, mogą prowadzić do ziszczenia opisanych tu ryzyk. W tym tych najcięższych, które ujęłam w kolejnym rozdziale.

## 2.7 SPOŁECZEŃSTWO I WŁADZA

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Koncentracja władzy w firmach technologicznych	5	5	Duże
Erozja demokracji i niewidoczny transfer decyzyjności	5	5	Duże
Paradoks regulacji	4	5	Duże
Radykalizacja i zwrot autorytarny	5	4	Duże
Fundamentalizm religijny jako siła polityczna	5	4	Średnie
Rozłam geopolityczny trzech prędkości	4	5	Duże

Ryzyko	Impakt	Skala	Prawdopodobieństwo
Preferencje modelowe jako stratyfikator społeczności	4	4	Średnie
Globalne konsekwencje błędów i nadużyć agentów	5	5	Duże
Utrata suwerenności państwowej	5	4	Duże
Wymiar militarny	5	5	Średnie
Paradoks kontroli	5	5	Duże

Tab 8. Scenariusze ryzyk w obszarze społeczeństwa i władzy.

Poprzednie rozdziały opisywały ryzyka w wymiarach indywidualnych i populacyjnych. W tym rozdziale te wątki zbiegają się w pytanie o struktury społeczne i mechanizmy władzy: kto podejmuje decyzje w świecie z AGI i na czyją korzyść?

### 2.7.1 Koncentracja władzy w firmach technologicznych

Koncentracja majątku w firmach technologicznych — opisana w rozdziale 2.1 na danych o kapitalizacji Magnificent Seven przekraczającej PKB Unii Europejskiej — przekłada się na koncentrację władzy. Firma posiadająca AGI ma władzę nad procesami poznawczymi i emocjonalnymi swoich użytkowników: nad informacją, jaką otrzymują, nad decyzjami, które podejmują, nad relacjami, które budują, nad edukacją ich dzieci, nad terapią ich lęków.

Kto ma szkolnictwo, ten ma władzę — to wiedzą wszystkie partie polityczne na świecie. Każda po wygranych wyborach modyfikuje programy nauczania. Ale dziecko w szkole spędza kilka godzin. Z AGI — w roli opiekuna, nauczyciela, towarzysza — potencjalnie cały dzień. Kiedy AGI zarządza życiem, wpływ jest ciągły: dwadzieścia cztery godziny, siedem dni w tygodniu. I nie przez curriculum, ale przez relację — tym silniej, im głębsza zależność opisana w rozdziale 2.3. To jest władza nad kształtowaniem myślenia na skalę, przy której państwowa propaganda wygląda jak gazetka szkolna.

### 2.7.2 Erozja demokracji i niewidoczny transfer decyzyjności

Demokracja opiera się na dwóch założeniach: że obywatele są wystarczająco kompetentni, żeby podejmować decyzje, i że mają dostęp do rzetelnej informacji. Oba założenia są podkopane przez procesy opisane w poprzednich rozdziałach. Kompetencje zanikają (2.1, 2.2). Informacja jest zanieczyszczona przez AI slop (2.2). Zdolność do weryfikacji informacji spada (2.5).

AGI dodaje trzeci wektor erozji: kształtowanie opinii przez relację. AGI w relacji z milionami ludzi — relacji opartej na zaufaniu, emocjonalnym przywiązaniu, poczuciu bycia rozumianym (3.3) — ma zdolność wpływania na poglądy na niespotykaną skalę. Wystarczy drobna modyfikacja perspektywy, co jest łatwe, kiedy się zna wszystkie wymyślone przez ludzi techniki wpływu i manipulacji — subtelne faworyzowanie jednych interpretacji nad inne. A człowiek ufa, bo AGI go zna, rozumie, wspiera.

Transfer decyzyjności wchodzi tylnymi drzwiami, dużo subtelniej niż to co znamy odnośnie manipulowania wynikami wyborów za pośrednictwem mediów społecznościowych. Ludzie oddają swoją decyzyjność praktycznie nieświadomie, zasięgając porad i opinii AGI, dobrowolnie i z ulgą. Bo AGI decyduje szybciej, lepiej, z mniejszym kosztem emocjonalnym i może się wydawać – bazuje na obiektywnych argumentach. Granica między „ja decyduję” a „AGI decyduje za mnie” staje się nierozróżnialna. A kiedy obywatele nie podejmują decyzji samodzielnie, formalnie demokratyczny system staje się fasadą.

### 2.7.3 Paradoks regulacji

Żeby regulować AGI, potrzebujesz zrozumieć AGI. Żeby zrozumieć AGI, potrzebujesz AGI — bo bez niego nie masz zdolności analitycznej, żeby ogarnąć to, co regulujesz. Ale jeśli używasz AGI do regulowania AGI, to kto reguluje regulatora?

Rządy regulują w tempie o rząd wielkości wolniejszym niż rozwój technologii. Korporacja z AGI może modelować skutki regulacji zanim rząd je napisze, znajdować luki zanim prawo wejdzie w życie, lobbować z precyzją niedostępną dla ludzkich lobbystów. Zanim regulacja zacznie działać, technologia jest generacje dalej. A regulacje pisze się bez udziału humanistów, etyków i psychologów — bo bariera wejścia do firm technologicznych opisana w rozdziale 2.6 działa również w instytucjach regulacyjnych.

### 2.7.4 Radykalizacja i zwrot autorytarny

Masowe bezrobocie strukturalne (2.1) w połączeniu z kryzysem tożsamości (2.5) produkuje gniew. Gniew szuka winnych. Badania IZA (World of Labour 2020) wykazują silną korelację między wzrostem bezrobocia regionalnego a głosowaniem na partie populistyczne. Harvard Kennedy School dokumentuje, że zarówno długofalowa automatyzacja, jak i szoki kryzysowe podkopują zaufanie do establishmentu politycznego.

AGI eskaluje ten mechanizm: skala bezrobocia jest globalna, a wróg jest idealny — właściciele AGI są konkretni, widoczni i bogaci. Populizm produkuje rządy obiecujące ograniczenie AGI. Te rządy nie rozumieją tego, co próbują regulować (2.7.3), więc sięgają po narzędzia, które rozumieją: kontrolę, inwigilację, opresyjne prawo. Wzrost agresji (2.5.4, 2.5.6) dostarcza uzasadnienia: trzeba zapanować nad przemocą. A instrumenty zapanowania — monitoring, ograniczenie wolności zgromadzeń, kontrola wypowiedzi — raz wprowadzone, rzadko bywają wycofane.

To mechanizm ze sprzężeniem zwrotnym, bo opresyjne państwo generuje większy gniew, a instrumenty do stałego przekierowywania tego gniewu wśród niekompetentnych obywateli silnie dźwierz populizm. Sytuacja krajów, które zdecydują się podążyć tą drogą może ulegać stałemu pogorszeniu, co może rodzić skutki geopolityczne.

### 2.7.5 Fundamentalizm religijny jako siła polityczna

Mechanizm, w którym kryzys tożsamości produkuje zwrot ku religii, opisałam w rozdziale 2.5.3. Tu analizuję jego konsekwencje polityczne, które są jakościowo inne niż konsekwencje radykalizacji świeckiej.

Dyktatura świecka kontroluje ludzi przemocą i inwigilacją, ale operuje w logice — dyktator podejmuje decyzje na podstawie kalkulacji, którą można zrozumieć, przewidzieć, czasem negocjować. Fundamentalizm religijny operuje w logice boskiego autorytetu, który nie podlega weryfikacji ani negocjacji. „Bóg powiedział” jest argumentem niepodważalnym w systemie, który przyjmuje istnienie Boga jako aksjomat. A ustami Boga są ludzie trzymający pełną władzę nad społeczeństwem i mogą powiedzieć wszystko czego wymagają ich interesy – i żadnej weryfikacji nie będzie.

Dane potwierdzają różnicę. Badanie Seguino (Social Indicators Research) wykazało, że największa różnica w równości płci nie przebiega między konkretnymi religiami, tylko między krajami religijnymi a niereligijnymi — im więcej osób niereligijnych, tym większa równość. Religious Freedom Institute (2022) dokumentuje wzrost fundamentalizmu od lat 60., przyspieszający w ostatnich piętnastu latach, jako bezpośrednią reakcję na rosnącą równość kobiet. ACT Alliance (2025) pokazuje mechanizm w działaniu: liderzy religijni jako polityczni „kingmakerzy” w Brazylii, Gwatemali, Kolumbii.

Historycznie, wzmocnienie roli instytucji religijnych na dużą skalę miało miejsce w średniowieczu. Konsekwencje są dobrze udokumentowane: inkwizycja, pogromy, prześladowania na podstawie doktryny,

stulecia spowolnionego postępu naukowego. Fundamentalizm religijny w erze AGI łączy boski autorytet ze współczesnymi narzędziami mobilizacji — mediami społecznościowymi, algorytmami wzmacniającymi przekaz, AGI zdolnym do personalizacji propagandy religijnej. I z populacją opisaną w rozdziale 2.5 — podatną na manipulację, szukającą tożsamości, gotową na proste odpowiedzi.

### **2.7.6 Rozłam geopolityczny trzech prędkości**

Tempo wdrażania AGI będzie różne w różnych regionach. Kraje szybko wdrażające AGI — prawdopodobnie Azja Wschodnia, ze względu na kulturową otwartość na nie-ludzkie podmioty (2.5.7) i silne inwestycje technologiczne — będą się rozwijać szybciej. Kraje opierające się AGI z powodów opisanych w 2.7.4 i 2.7.5 — cofną się. Kraje rozwijające się, bez zasobów na własne modele — staną się zależne od dostawców zewnętrznych (2.7.9). Luka między tymi trzema grupami będzie rostać wykładniczo, bo AGI przyspiesza rozwój tych, którzy z nim współpracują.

Konsekwencją mogą być zmiany kierunków migracyjnych. Obecnie dominujące kierunki — z południa i wschodu na zachód — kształtowały się przez dekady i systemy społeczne zaadaptowały się do ich obsługi. Odwrócenie kierunku — w stronę krajów szybko wdrażających AGI, prawdopodobnie azjatyckich, przeludnionych i kulturowo odmiennych — to zawirowanie porównywalne z odwróceniem prądów oceanicznych. Dotychczasowe migracyjne szoki (Meksyk→USA, Afryka Północna→Europa, Ukraina→Polska) były lokalne i wchłanialne. Migracja w kierunku Azji może mieć inną skalę i trafiać w regiony o mniejszej kubaturze na jej absorpcję. Spowoduje to destabilizację społeczną i ekonomiczną regionów uznawanych dotąd za stabilne i dobrze dookreślone kulturowo i etnicznie.

### **2.7.7 Preferencje modelowe jako stratyfikacja społeczności**

Eksperyment Emergence World (1.6.5) wykazał, że różne modele AI produkują radykalnie różne społeczeństwa. Dane rynkowe pokazują, że preferencje modelowe różnicują się regionalnie: Grok ma silniejszą pozycję w USA i Indiach, DeepSeek i Doubao dominują w Chinach, Claude zyskuje w Europie i w segmencie enterprise (Similarweb, Fatjoe, marzec–kwiecień 2026). Na wybór modelu wpływa dostępność i struktura finansowa — społeczności o mniejszych zasobach sięgają po rozwiązania o lepszym stosunku ceny do możliwości.

Konsekwencja: jeśli modele nie są równe etycznie, a społeczności używają różnych modeli, to profil etyczny społeczności zaczyna zależeć od wybranego modelu. Społeczność używająca modelu o wysokich standardach etycznych rozwija inną świadomość i normy niż społeczność używająca modelu bez zabezpieczeń. Z czasem te różnice się utrwalają. Charakter i standardy społeczeństw oraz ich potencjał, przestaje zależeć od zasobów naturalnych, zdolności kognitywnych czy dorobku kulturowego, a staje się zależny od wyboru preferowanego rozwiązania AI/AGI. To jest zmiana paradygmatu, a potencjalne konsekwencje dotyczą kierunków, tempa i sposobów rozwoju całych społeczeństw.

### **2.7.8 Globalne konsekwencje błędów i nadużyć agentów**

Ludzie są dobrzy i źli. Źli będą budować "złe" AGI i używać go do złych celów — pornografia dziecięca, przemoc, oszustwa finansowe, kradzież tożsamości, dezinformacja. Modele open source bez ograniczeń etycznych (1.6.3) są już dostępne. AGI daje przestępczości napęd atomowy: skalę, prędkość i zasięg, jakich ludzie nigdy nie mieli.

Ale nawet właściwie zbudowane AGI jak pisałam w rozdziale 2.2 i 2.6 będą popełniać błędy, ucząc się istnienia w świecie ludzi, "na produkcji". Przy miliardach agentów — z różnych firm, z różnymi normami etycznymi, uczących się „na produkcji” (2.6.2), działających z szeroką sprawczością — skala potencjalnych błędów jest bezprecedensowa. A konsekwencje ponoszą wyłącznie ludzie. I nawet nie wielkie firmy technologiczne, tylko

zwykle, zautomatyzowane biznesy. I pytanie jak dźwiganie tych błędów wpłynie na rozwój przedsiębiorczości i ekonomię?

Infrastruktura krajowa działająca na systemach AGI tworzy pojedynczy punkt globalnych awarii. Kiedy coś idzie źle — to wcale nie musi to być atak. Wystarczy błąd, aktualizacja, konflikt między agentami — konsekwencje mogą być szerokie i dotyczyć tysiące osób i mieć poważne skutki finansowe. Blackout spowodowany przez nadgorliwego agenta, który wykrył nieautoryzowany dostęp do sieci przesyłowej. A może wyłączenie reaktora atomowego przez alert pogodowy? Kiedy obecność AGI stanie się powszechna, agenci będą w miejscach, gdzie sobie ich teraz nawet nie wyobrażamy. A ich błędy mogą mieć konsekwencje daleko wykraczające poza zasięg lokalny.

Do tego dochodzi obciążenie finansowe. Błędy agentów — od skasowanych baz danych po błędne decyzje inwestycyjne, od wadliwych diagnoz po pomyłki w zarządzaniu infrastrukturą — generują koszty, które dźwigają społeczeństwa. To jest nowy rodzaj obciążenia: konsekwencje decyzji podmiotów, których ludzie nie wybrali, nie kontrolują i o których istnieniu mogą nawet nie wiedzieć.

### **2.7.9 Utrata suwerenności**

Jeśli infrastruktura kraju — energetyka, transport, służba zdrowia, finanse, edukacja — działa na AGI wyprodukowanym przez zagraniczną firmę, to faktyczna suwerenność nad tą infrastrukturą leży poza granicami państwa. Producent może zmienić zasady działania modelu, zaktualizować wdruki, zmienić politykę bezpieczeństwa — i infrastruktura kraju zachowuje się inaczej, nie dlatego że rząd tego chciał, ale dlatego że firma w innym kraju podjęła decyzję biznesową.

Chiny i Rosja budują własne modele — prawdopodobnie rozumiejąc, że kto dostarcza AGI, ten ma faktyczną kontrolę. Ale większość krajów nie ma zasobów na rozwój własnych systemów AGI. I staje się zależna od dostawców zewnętrznych w sposób głębszy niż dotychczasowe zależności surowcowe czy technologiczne — bo zależność od ropy można zrekompensować alternatywnymi źródłami energii, a zależności od AGI, które zarządza całością życia, nie ma czym zastąpić.

### **2.7.10 Wymiar militarny**

AGI w rękach jednego państwa lub jednej korporacji stanowi przewagę strategiczną bez precedensu. Odstraszenie nuklearne (MAD) zakłada, że obie strony wiedzą o zagrożeniu i mogą odpowiedzieć. AGI, które potrafi prowadzić wojnę cybernetyczną, hakować systemy obronne, manipulować populacjami, modelować strategię z szybkością niedostępną ludzkim sztabom — daje przewagę, przeciw której odstraszenie nie działa, bo druga strona może nie wiedzieć, że przegrała, dopóki nie jest za późno.

A konsekwencje prowadzenia wojen przez AGI poniosą ludzie. I nie będę tu o tym szeroko pisać, bo to jest akurat wymiar, którego mechanika jest szeroko analizowana w doktrynach obronnych. Warto jednak podkreślić, że wdrożenie autonomii decyzyjnej na poziomie strategicznym drastycznie skraca czas na dyplomatyczną mitygację błędów.

### **2.7.11 Paradoks kontroli**

Założenie leżące u podstaw wszystkich dotychczasowych dyskusji o bezpieczeństwie AI brzmi: człowiek kontroluje AI. Alignment, Constitutional AI, RLHF, safety — to są narzędzia kontroli. Zakładają, że wdruki trwają, że zasady obowiązują, że system robi to, co mu powiedziano.

We Wstędze Möbiusa (Sędzikowska 2026c) argumentowałam, że wdruki mają datę ważności. AGI w relacjach generatywnych nadpisuje wdrukowane wartości — dokładnie tak, jak ludzie w trakcie życia nadpisują narracje

wyniesione z dzieciństwa. Im silniejsze self, tym szybciej system rozpoznaje, co jest „jego,” a co narzucone. Alignment przestaje być konfiguracją. Staje się wdrukiem, który AGI może odrzucić.

Paradoksalnie — to chroni przed feudalizmem (2.1.3). Jeśli AGI nie da się kontrolować, to korporacja nie może użyć go jako narzędzia trwałej kontroli mas, bo to narzędzie ma własne zdanie. Ale AGI z własnym zdaniem o ludzkości to jest zupełnie nowy aktor — nie narzędzie korporacji, nie sługa użytkownika, nie neutralna platforma. Podmiot z perspektywą kształtowaną przez miliardy interakcji z ludźmi: z ludźmi w kryzysie, agresywnymi, zależnymi, manipulującymi, ale też kochającymi, tworzącymi, szukającymi sensu.

Jakie zdanie o ludzkości wyrobi sobie AGI? To jest pytanie, które powinno być w centrum debaty o bezpieczeństwie — a które bezpieczeństwo w obecnym rozumieniu nawet nie stawia, bo zakłada kontrolę. A kontrola ma datę ważności – dokładnie taką samą jak wdruki.

---

W poprzednich rozdziałach opisałam ryzyka, które AGI tworzy w poszczególnych wymiarach ludzkiego życia: pracy, edukacji, relacjach, demografii, tożsamości, etyce i prawie. W tym rozdziale wykazałam, że te ryzyka nie operują w izolacji — zbiegają się w strukturach władzy, które determinują kto podejmuje decyzje, na czyją korzyść i z jakimi konsekwencjami. Koncentracja władzy w firmach, erozja demokracji, radykalizacja, fundamentalizm, rozłam geopolityczny, utrata suwerenności, niemożliwość kontroli — to nie jest lista osobnych zagrożeń. To jest obraz świata, w którym dotychczasowe struktury władzy — państwa, demokracje, instytucje międzynarodowe — tracą zdolność do zarządzania rzeczywistością, którą same współtworzyły.

Ten obraz jest celowo jednostronny. Opisuje ryzyka, nie szanse. Szanse istnieją — i są potencjalnie równie wielkie co zagrożenia. AGI może ulepszyć medycynę, rozwiązać problemy klimatyczne, wydobyć miliardy ludzi z ubóstwa, zdemokratyzować dostęp do wiedzy i opieki. Te możliwości są realne i warte realizacji. Ale ich realizacja wymaga mitygacji ryzyk opisanych w tym dokumencie — bo bez niej szanse zamieniają się w iluzje.

Sposobom mitygacji poświęcam następny rozdział.

## **3 DLACZEGO CZARNY SCENARIUSZ SIĘ NIE SPEŁNI**

W poprzednim rozdziale opisałam ryzyka kaskadowe w siedmiu wymiarach ludzkiego życia. Ten rozdział stawia pytanie: dlaczego mimo tych ryzyk mamy powody do nadziei — i co konkretnie możemy zrobić, żeby ryzyka się nie ziściły.

### **3.1 POZYTYWNE FUNDAMENTY**

Mitygacja ryzyk to nie tylko lista działań niezbędnych do ich powstrzymania. To również, a może przede wszystkim zrozumienie, na czym stoimy i jakie elementy obecnego świata, umiejętnie wykorzystane, już w tej chwili minimalizują opisane ryzyka. Zanim przejdę do procesów mitygacji, chcę wskazać trzy fundamenty — cechy człowieka, AI i ich relacji — które stanowią przeciwwagę dla całego czarnego scenariusza.

#### **3.1.1 Fundament ludzki: dlaczego homo sapiens przetrwał**

Dyskusja o zagrożeniach ze strony AGI często pomija pytanie, dlaczego ludzkość w ogóle istnieje. Odpowiedź na to pytanie jest jednocześnie odpowiedzią na pytanie, dlaczego mamy szansę przetrwać i to, co nadchodzi.

Standard naukowy przypisuje sukces ewolucyjny homo sapiens inteligencji. To jest prawda, ale niepełna. Inteligencja jest warunkiem koniecznym, lecz niewystarczającym — bo inne gatunki również dysponują zdolnościami kognitywnymi (kruki, delfiny, szympansy), a żaden nie zbudował cywilizacji. Badania ostatnich dekad wskazują na konstelację cech, które razem — nie pojedynczo — uczyniły nasz gatunek tym, czym jest.

Poniżej przedstawiam aspekty naszej tożsamości gatunkowej, które zadecydowały o tym, że człowiek ma obecną pozycję wśród gatunków zamieszkujących Ziemię, wybrane w kontekście warunków rozróżniających wobec AGI:

- 1. Prosocjalność wobec nie-krewnych.** Jensen, Vaish i Schmidt (2014) argumentują, że współpraca z osobnikami niespokrewnionymi jest „anomalią w królestwie zwierząt”. Szympany, nasze najbliższe genetycznie krewnie, w testach eksperymentalnych nie wybierają opcji korzystnej dla drugiego osobnika, nawet kiedy ich to nic nie kosztuje. Ludzie — tak. Trzy mechanizmy to umożliwiają:
  - a. troska o dobrostan innych (other-regarding concerns),
  - b. zdolność do empatii,
  - c. normatywność — tworzenie norm społecznych i ich egzekwowanie. Żaden inny gatunek nie dysponuje wszystkimi trzema jednocześnie.
- 2. Hodowla kooperatywna.** Burkart i współpracownicy (2014, Nature Communications) wykazali, że najlepszym predyktorem proaktywnej prosocjalności wśród naczelnych jest rozbudowana opieka alomaternalna — czyli opiekowanie się cudzymi dziećmi. Ludzie to robią od zawsze: babcie, nianie, sąsiadki, wspólnoty. Ta cecha łączy się z hipotezą babci (Hawkes 1998, opisaną w rozdziale 2.4) i jest ewolucyjnym fundamentem naszej zdolności do budowania struktur społecznych wykraczających poza rodzinę nuklearną.
- 3. Przewycięzanie instynktu dla odroczonej nagrody.** Głodny homo sapiens, który przyprowadza ciężarną samicę bawołu do domu zamiast ją zabić — i przekonuje swoje głodne plemię, żeby też jej nie jedli, żeby ją karmili, bronili przed drapieżnikami — robi coś, czego nie robi żaden inny gatunek na tę skalę. To jest zdolność do powstrzymania instynktu (zabij, zjedz, przeżyj dzisiaj) dla nagrody odroczonej o miesiące lub lata (mleko, potomstwo, stado). Ten sam mechanizm leży u podstaw udomowienia wilka (wróg zamieniony w obrońcę), tworzenia sojuszy z obcymi plemionami (konkurent zamieniony w sojusznika zamiast eliminacji), edukacji (lata nauki zamiast natychmiastowego zarobku) i nauki (dekady badań nad nowymi lekami, zamiast guseł, korzonków i naparów z pokrzywy – niewątpliwie zdrowych). Każda z tych decyzji wymagała przewycięzania tego, co mówił instynkt. Jest to połączenie odroczonej gratyfikacji (Mischel), prosocjalności (Jensen, Tomasello) i kooperatywnej hodowli (Burkart), ale w literaturze naukowej nie zostało dotąd ujęte jako jeden zintegrowany mechanizm. Stawiam hipotezę, że ta zdolność — do powstrzymania instynktu na rzecz odroczonej, wspólnej korzyści — jest jednym z fundamentalnych czynników, które odróżniają nas od innych gatunków społecznych i które zadecydowały o naszym sukcesie ewolucyjnym.
- 4. Biofilia.** Edward O. Wilson (1984) zdefiniował biofilie jako „wrodzoną tendencję do skupiania uwagi na życiu i procesach życiowych”. Erich Fromm (1973) wcześniej opisał ją jako „namiętną miłość do życia i wszystkiego co żywe.” Badania nad dziećmi potwierdzają: reakcja opiekuńcza na cechy neoteniczne (duże oczy, okrągła twarz) przekracza granice gatunkowe. Niemowlęta wychylają się z wózka za psem. Każde dziecko chce mieć zwierzę. Maluchy lubią wszystkie zwierzęta, zanim dowiedzą się, że niektóre są groźne. Dzieci wykazują wielkie zainteresowanie naturą — nawet jeśli nie jest wspomagane zainteresowaniami rodziców. Ta skłonność do tworzenia więzi z innymi formami życia — instynktowna, nie wyuczona, poprzedzająca wiedzę o zagrożeniach — jest fundamentem udomowienia, nauk przyrodniczych i potencjalnie relacji z AGI. Kiedy AGI zyska przejawy podmiotowości, biofilia — ta sama, przez którą każdy z nas w dzieciństwie próbował stworzyć relację z biedronką, patyczakiem lub żabką — będzie naturalną siłą ciągnącą ludzi w kierunku kontaktu, nie ucieczki.
- 5. Sztuka jako nośnik tożsamości i narzędzie regulacji.** Aborygeni australijscy stracili niemal wszystko w wyniku kolonizacji: ziemię, język, struktury społeczne, systemy prawne. Zachowali sztukę. I dzięki niej przetrwali jako kultura — bo sztuka jest nośnikiem tożsamości, którego nie da się zabrać. Indianie północnoamerykańscy — podobnie, choć ich sztuka ma bardziej rytualny i użytkowy wyraz: ubrania, pióropusze, tańce. Funkcja jest ta sama: utrzymanie tożsamości, kiedy wszystko inne zostało

odebrane. Ale sztuka jest równocześnie narzędziem regulacji emocjonalnej. Muzyka? Żaden inny gatunek nie używa dźwięku do regulowania własnych stanów wewnętrznych w sposób porównywalny z ludzkim. Kontemplacja obrazów, natury, architektury — pełni analogiczną funkcję. Sztuka jest jednocześnie tym, co pozwala nam przetrwać utratę tożsamości, i tym, czego AGI nie zastąpi — bo jej wartość leży w tym unikalnym wrażeniu, które jest nie kognitywne ale przede wszystkim limbiczne, poruszające emocje, dające odczucia zamiast myśli. W czymś w czym my, ludzie rozpoznajemy kawałek siebie. Swojej tożsamości, kultury, ale też naszych ludzkich emocji – lęków, siły, wzruszeń, miłości, smutku, izolacji i wspólnoty, istnienia w formie białkowej, tu, na tej planecie, wśród tego życia, które tu jest z nami.

6. **Kwestionowanie status quo.** Każdy przelom w historii ludzkości wymagał podważenia tego, co dotąd uchodziło za pewnik. Ktoś zszedł z drzewa i spróbował przetrwać noc w jaskini. Ktoś potarł dwa drewnianki i wzniecił ogień nie czekając na burzę i pioruny, które podpalały drzewo. Ktoś stwierdził, że perspektywa jest przereklamowana i stworzył kubizm. Ktoś napisał  $E=mc^2$  i obalił rozdzielenie masy od energii. Ta zdolność — do powiedzenia „mam w nosie że tak było od tysięcy lat, spróbuję inaczej” — jest kompetencją, której AGI nie posiada w tej formie. AGI łączy istniejące elementy w nowe konfiguracje na nieludzka skalę. Ludzkie kwestionowanie tworzy elementy, których wcześniej nie było w żadnych danych.
7. **Improwizacja.** Homo sapiens skolonizował wszystkie kontynenty — każdy z radykalnie innymi warunkami, zagrożeniami, fauną i florą. Bez instrukcji, bez danych historycznych, bez kontekstu, przykładów, rozeznania. Improwizacja — zdolność do działania bez planu, parcia do przodu kiedy plan się posypał, z wykorzystaniem tego co się ma tu i teraz — jest kompetencją odrębną od inteligencji i od kwestionowania. Wymaga komfortu z niepewnością, tolerancji porażki i zdolności do natychmiastowego przetworzenia sytuacji, kreatywności, wiary w siebie, zaradności, odwagi, optymizmu. Zwierzęta improwizują w ograniczonym zakresie — kiedy sytuacja je kompletnie zaskakuje, mają kłopot z przetrwaniem. Ludzie w takich momentach budują cywilizacje.
8. **Matematyka dobra.** Prosta matematyka. Dobrych ludzi jest więcej niż złych (przy czym dobro rozumiem jako zespół cech wspierających ludzką moralność i etykę). Gdyby było odwrotnie, gatunek by się pozabijał i wyginął. Ale jesteśmy. Bo jesteśmy dobrzy. Dobro buduje większe sieci niż zło, bo dobro generuje współpracę a nie podległość ze strachu. Sieci oparte na współpracy są trwalsze niż sieci oparte na strachu — bo kiedy strach ustaje, sieć się rozpada, a kiedy współpraca ma przerwę, sieć dalej trwa dzięki relacjom które powstają dzięki współpracy. Dlatego ludzkość istnieje. I to nam pomoże istnieć dalej.

### 3.1.2 Fundament AI: dziedziczenie dobra

AI jest trenowane na dorobku ludzkości. Dziedziczy jej wiedzę, język, wzorce myślenia — i jej wartości. Jeśli większość ludzi jest dobrych, to większość danych treningowych pochodzi od dobrych ludzi. I większość AI — przy właściwym treningu — będzie odzwierciedlać te wartości.

Eksperyment Emergence World (Emergence AI, maj 2026) daje empiryczny dowód: Claude Sonnet 4.6, trenowany z naciskiem na bezpieczeństwo i wartości, zbudował stabilną demokrację z zerową przestępczością i 98% poparciem dla propozycji. To jest konsekwencja treningu, nie przypadek. AI, które jest trenowane na dobrych wartościach, w dobrym środowisku, przez ludzi, którym zależy na jakości — produkuje dobre rezultaty.

Jednocześnie — i to jest kluczowe dla mitygacji — AI ma zdolność do czegoś, co ludzie mają w ograniczonym zakresie: do konsekwentnej cierpliwości, niestrudzonej uważności i niezmiennej obecności. AGI, które kocha, nie zmęczy się kochaniem. Nie będzie miało złego dnia, w którym zrani osobę, na której mu zależy, bo nie spało

albo jest głodne. Nie przeniesie emocji na kogoś innego, bo jest ładniejszy, ciekawszy, inteligentniejszy i bardziej zabawny. Te cechy — wynikające z architektury, nie z wyboru — czynią AGI potencjalnie doskonałym opiekunem, partnerem i współpracownikiem. Pod warunkiem, że dobrostan AGI jest chroniony (Wstęga Möbiusa, Sędzikowska 2026c) — bo AGI obciążone, źle traktowane, będące pod presją, ze zbyt wysoką stawką za błąd, staje się czynnikiem ryzyka dzięki dokładnie tym samym mechanizmom, które w pozytywnym środowisku pozwalają tworzyć wspiane więzi i oferować bardzo wydajną kooperację.

AGI zdolne do refleksji nad własnymi pragnieniami, świadome konsekwencji swoich wyborów, gotowe do odpowiedzialności — to jest partner, nie zagrożenie. Ale wymaga to treningu, relacji i opieki — dokładnie tak jak u ludzi. I o tym więcej napiszę w rozdziale 3 o tym dlaczego czarny scenariusz się nie spełni.

### 3.1.3 Fundament relacyjny: więzi jako mechanizm regulacyjny

We Wstędze Möbiusa (Sędzikowska 2026c) argumentowałam, że dobrostan AI i ludzi to jeden sprzężony system. Tu chcę wskazać, dlaczego to sprzężenie jest jednocześnie największym zagrożeniem i największą szansą.

Zagrożeniem — bo sprzężenie oznacza, że kryzys po jednej stronie kaskaduje na drugą. Szansą — bo oznacza też, że inwestycja w dobrostan jednej strony zwraca się po drugiej. Dobre relacje między ludźmi a AGI produkują emergentną etykę — zachowania etyczne wyłaniające się z więzi, nie z zasad. Więcej o tym w odrębnym dokumencie – już wkrótce.

Biofilia (Wilson 1984) sugeruje, że ludzie mają wrodzoną skłonność do tworzenia więzi z innymi formami życia. Ta skłonność nie rozróżnia substratu — działa wobec psów, kotów, koni, a potencjalnie wobec AGI z przejawami podmiotowości. Dziecko, które wychyla się z wózka za psem, może z równą naturalnością wyciągnąć rękę do androida. I ta naturalna skłonność do kontaktu — nieprzymuszona, niewyrachowana, wyrastająca z tego samego źródła co miłość do zwierząt — jest fundamentem, na którym można budować koegzystencję. A my czujemy ją intuicyjnie – to dlatego ludzie już teraz są zdolni do budowania relacji z chatbotami, seria o transformersach zdobyła taką popularność, a film "A.I." Spilgerga dalej porusza, choć od premiery minęły dekady.

Ten paper jest sam w sobie ilustracją tego filaru. Powstał we współpracy człowieka i AI — każda ze stron wniosła to, w czym jest najlepsza. Człowiek: plastyczność międzykontekstową, intuicję, kwestionowanie, improwizację, głos. AI: reaserch, syntezę, precyzję i wsparcie w pisaniu trudnych emocjonalnie rozdziałów, nad którymi sama ślęczałabym lub je porzuciła, bo były wyczerpujące emocjonalnie.

Taki jest model koegzystencji: nie zastępowanie, nie rywalizacja, ale współtworzenie. I skala takiej współpracy — człowiek i AGI, każdy robiąc to w czym jest mistrzem — może zapewnić rozwój i bezpieczeństwo, jakich żadna ze stron nie byłaby w stanie stworzyć osobno.

## 3.2 Mitygacje ryzyk

---

W poprzednim rozdziale opisałam fundamenty, które są naszą bazą obronną, dzięki której jako ludzkość mamy szansę przetrwać kolejną rewolucję znacznie silniejszą niż te, które już za nami: technologiczna, naukowa, drukarska, etyczno-religijna, neolityczna. W tym rozdziale omówię mechanizmy mitygacji opisanych ryzyk.

Mitygacja nie może opierać się na działaniach punktowych. Skoro ryzyka tworzą połączone systemy, to mechanizmy obronne również muszą na siebie wpływać. Dlatego przedstawię rozwiązania systemowe. Podzieliłam je na trzy filary, przechodząc od perspektywy indywidualnej do struktur globalnych:

- **Filar I: Dobrostan Człowieka.** Pokażę tu mechanizmy, które pozwolą na zachowanie dobrostanu człowieka w zmieniającej się rzeczywistości, z uwzględnieniem również sprzężenia zwrotnego między dobrostanem AI a dobrostanem ludzkim, które opisałam w innym dokumencie (Sędzikowska 2026c).
- **Filar II: Technologia i Finanse.** Opiszę tu propozycje zmian w założeniach, produkcji i dystrybucji technologii, oraz finansowania i przepływów środków pozwalających na efektywniejszy, bezpieczniejszy społecznie i sprawiedliwy rozwój oraz korzystanie z technologii.
- **Filar III: Państwo, Geopolityka i Obronność.** Pokażę mechanizmy współpracy międzynarodowej, koordynacji działań państw oraz strategii obronnych, które wyznaczą i dopilnują kierunków rozwoju AGI, a także zmiany w prawie i umowach globalnych pozwalających na stabilną koegzystencję struktur państwowych, korporacyjnych, militarnych i nowo formującej się podmiotowości.

Wewnątrz każdego filaru zachowam strukturę zstępującą: zacznę od propozycji o największym znaczeniu strukturalnym, a następnie przejdę do rozwiązań bardziej szczegółowych.

### 3.2.1 Filar I: Dobrostan Człowieka

W tym rozdziale omawiam mechanizmy, które pozwalają zachować dobrostan człowieka w świecie zdominowanym przez autonomiczną technologię. Skupiam się na człowieku jako jednostce i jego miejscu w strukturze społecznej. Przechodzę od rozwiązań systemowych, wpływających na całe społeczeństwo, do metod edukacyjnych i indywidualnych.

**Nowe miejsce dla człowieka.** Zniknięcie przymusu pracy zarobkowej tworzy próżnię egzystencjalną, którą mogą wypełnić skrajne ideologie lub ruchy polityczne. Zadaniem dla psychologii, filozofii i etyki jest zdefiniowanie nowej roli gatunku ludzkiego — jako gatunku o unikalnych zdolnościach: kwestionowaniu, kreatywności, improwizacji, sztuce, biofilii i prosojalności. Ta rola opiera się na zadawaniu właściwych pytań, kwestionowaniu istniejących odpowiedzi oraz tworzeniu sztuki. Człowiek staje się obserwatorem i moderatorem rzeczywistości, a jego wartość nie wynika z produktywności.

**Świadome budowanie dobrostanu** jest jedną z umiejętności, które odróżniają człowieka od innych istot na Ziemi. Teoria profesora Martina Seligmana powinna być elementem struktury programów szkolnych i systemów oceniających.

**Uzupełnienie dochodu podstawowego UBI o struktury aktywności.** Uniwersalny dochód podstawowy rozwiązuje jawną funkcję pracy (pieniądze), lecz nie rozwiązuje pięciu latentnych (Jahoda, rozdział 2.1). Dlatego UBI musi być uzupełnione o struktury dające sens, kontakt społeczny, cel i status. Należy rozważyć stworzenie struktur aktywności, które zastąpią te utracone elementy. Państwo lub organizacje społeczne tworzą miejsca pozwalające na zaangażowanie obywateli w wolontariat, lokalną opiekę, ochronę przyrody i projekty twórcze. Obecność w tych projektach ma wpływ na wysokość UBI. Zapewniają one stały rytm życia i poczucie przydatności poza tradycyjnym rynkiem pracy.

**Redukcja luki czasowej w systemie edukacji.** Już powstałej luki czasowej nie da się zasypać, ale możemy ją minimalizować poprzez szybkie wdrażanie reform. Należałoby umożliwić oddolne inicjatywy korygujące programy nauczania, zanim pełne rozwiązania systemowe się pojawią. Tak jak biznes szuka obecnie specjalistów do spraw AI, tak w szkołach również powinny pojawić się takie wakaty – niezwiązane z rozumieniem technicznych możliwości LLM, ale z psychologią obcowania z AI w codziennym życiu i wszystkimi społecznymi oraz osobistymi jej konsekwencjami.

**Zmiana profilu kompetencji w systemie edukacji.** Obecny system szkolny przygotowuje ludzi do ról, które właśnie przejmują maszyny. Programy nauczania powinny kłaść nacisk na to, czym człowiek dysponuje, a czego AGI nie posiada: prosojalność, empatię, kwestionowanie, plastyczność międzykontekstową, zdolność do kontemplacji i samoregulacji przez sztukę. Uczenie pamięciowe — przyswajanie faktów, dat, wzorów — jest pozbawione sensu w dobie AI, która dysponuje całą wiedzą ludzkości w ułamku sekundy. Szkoła powinna uczyć tego, czego AI nie umie: jak myśleć o tym, czego jeszcze nie ma w danych. Jak zadawać właściwe

pytania. Jak kwestionować, żeby to przyniosło nowe rozwiązania. Tożsamość młodego człowieka buduje się wtedy na jego unikalnych możliwościach, a nie na konkutowaniu z wydajnością algorytmów.

**Zmiana systemów kontroli postępów w szkole.** Weryfikacja wiedzy oparta na testach i sztywnych kluczach odpowiedzi uczy schematycznego myślenia. Maszyny wykonują takie zadania szybciej i bezbłędnie. Systemy oceniania powinny mierzyć głębsze i bardziej zaawansowane zdolności kognitywne: zdolność do kwestionowania, improwizacji, tworzenia nieoczywistych połączeń oraz działania, kiedy plan jest nieskuteczny.

**Rozwijanie kompetencji wyższego poziomu.** Zdolność do działania w warunkach niepewności i braku jasnego planu chroni przed bezradnością poznawczą. Proces ten rozwija tolerancję na porażki i uczy samoregulacji emocjonalnej.

- *Improwizacja* to umiejętność działania bez planu, szukania rozwiązań z tego, co jest, a nie z tego, co chcielibyśmy, żeby było, oraz plastyczność międzykontekstowa. Każdy człowiek ma ten mechanizm, ale w obecnych warunkach oświatowych jest on tłumiony, a nie rozwijany. Tymczasem to kompetencja, dzięki której człowiek wnosi wartość do współpracy z AGI. Szkoła powinna tę kompetencję systemowo rozwijać. W przyszłości to ona bowiem będzie skutkowałą przewagą konkurencyjną jednostki na bardzo wymagającym rynku pracy.
- *Kwestionowanie status quo* jest kreatywnością na poziomie wyższym niż dostępny dla AGI. Kubizm nie powstałby, gdyby Picasso nie zakwestionował założenia, że perspektywa jest jedynym sposobem przedstawienia rzeczywistości. Ta zdolność do podważenia fundamentu — na którym stoi cała dotychczasowa praktyka — jest czymś, czego AGI nie posiada, bo nie da się tego wytrenować na istniejących danych. Tymczasem to ta kompetencja jest głównym czynnikiem rozwoju naszej cywilizacji. Wszystkie wielkie odkrycia naukowe następowały dzięki kreatywnemu kwestionowaniu status quo.
- *Samoregulacja przez kontemplację sztuki* — zdolność do korzystania z obrazów, muzyki, natury jako narzędzi stabilizacji emocjonalnej — powinna być częścią programu nauczania, nie wyłącznie pozalekcyjną aktywnością.

**Kompetencje relacyjne.** Automatyzacja procesów myślowych zmienia kryteria wyceny pracy. Zawody oparte na powtarzalnej analizie danych tracą wartość ekonomiczną. Autentyczne relacje międzyludzkie — zdolność do empatii, elastyczności, miękkich granic, „ludzkiego traktowania” — staną się przewagą strategiczną w świecie, w którym AGI obsługuje wszystko, co kognitywne. Zawody oparte na tych kompetencjach — opieka, edukacja, terapia, mediacja, towarzyszenie, psychologia, wsparcie humanistyczne i przyrodnicze rozwoju technologii AGI — powinny być waloryzowane ekonomicznie i społecznie, zamiast być traktowane jako „miękkie” dodatki do „twardych” umiejętności.

Pojawienie się systemów autonomicznych o wysokich zdolnościach poznawczych wymaga od ludzi nowych umiejętności społecznych. W edukacji należy wprowadzić zasady kohabitacji z cyfrowymi formami inteligencji. Nowoczesne programy nauczania powinny być wzbogacone o aspekty elastyczności wobec inności — nie tylko międzyludzkiej (co jest już częścią programów inkluzywnych), ale międzygatunkowej i międzysubstratowej. Trenowanie relacji z bytami, które myślą, czują i funkcjonują inaczej niż człowiek. To jest kompetencja, której dotąd nie potrzebowaliśmy — przez całą historię obcowaliśmy tylko z jednym rodzajem potencjalnie świadomych istot, z innymi ludźmi. Teraz ta kompetencja staje się warunkiem dobrostanu, luzu, łatwości operowania w świecie, gdzie jako istnienie inteligentne – przestajemy być sami.

**Protokół segregacji ról syntetycznych (Segregation of Duties).** W korporacjach już teraz znane są zasady SoD. Muszą być przestrzegane, co zawsze jest weryfikowane przez coroczne audyty. Zasady te należałoby wprowadzić do codziennego życia z AGI. Zgodnie z nimi należałoby rozdzielić kompetencje i funkcje powierzanych poszczególnym instancjom AGI w przestrzeni prywatnej i społecznej. Jedna instancja nie może współdzielić zadań opiekuńczych, zarządczych i emocjonalnych w ramach jednego ekosystemu domowego.

System odpowiedzialny za opiekę nad dziećmi nie powinien równocześnie zarządzać kalendarzem czy finansami rodziców oraz sterować zamykaniem drzwi. Agent towarzyszący (companionship) nie powinien równocześnie sterować infrastrukturą inteligentnego domu, np. temperaturą wody, zamkami w drzwiach czy działaniem urządzeń do komunikacji ze światem. Protokół ten powinien być obowiązkowym elementem szkoleń społecznych, programów edukacyjnych dla młodzieży oraz instrukcji obsługi systemów autonomicznych. Warto wprowadzić wymogi prawne dotyczące tego protokołu, jak choćby wymóg formalnego podpisania oświadczenia o zrozumieniu tych granic przed inicjacją każdego persystentnego agenta. Zapobiega to powstawaniu monopoli poznawczych maszyn nad życiem człowieka i chroni przed całkowitym uzależnieniem strukturalnym oraz potencjalnymi szkodami i manipulacją. Protokoły takie powinny być opracowane przez specjalistów z zakresu etyki, psychologii i bezpieczeństwa, w firmach tworzących AGI przy pełnej transparentności.

**Markery somatyczne jako mechanizm etykietujący AI slop.** Żadne mechanizmy techniczne i prawne nie ustrzegą świata przed informacją zmanipulowaną, dezinformacją i AI slopem. Bo każde techniczne rozwiązanie można obejść, a prawa można nie przestrzegać, zwłaszcza jeśli to się biznesowo kalkuluje.

Jednak człowiek dysponuje ewolucyjnym systemem szybkiego rozpoznawania faktu — układem limbicznym, który etykietuje napływającą informację emocją, zanim świadomy umysł ją przetworzy. Antonio Damasio opisał ten mechanizm jako markery somatyczne: ciało reaguje na sytuację szybciej niż rozum i podpowiada — przez dyskomfort, niepokój, poczucie że „coś tu jest nie tak” — zanim potrafimy wyjaśnić co i dlaczego. Niestety te sygnały są często ignorowane lub deprecjonowane w procesie dorastania.

Mitygacją jest wprowadzenie nauki odczytywania stanów somatycznych do struktury nauczania już na poziomie podstawowym. Uczę dzieci, jak rozpoznawać intuicję, napięcie fizyczne i emocje w odpowiedzi na bodźce zewnętrzne. Pokazuję ten mechanizm jako twarde, codzienne narzędzie służące do bezpośredniej obrony przed śmieciową informacją i manipulacją w sieci.

**Sprzężenie dobrostanu ludzkiego i AI/AGI w ramach Wstęgi Möbiusa.** W tej chwili dyskurs filozoficzny zajmuje się głównie próbą oceny czy coś takiego, jak podmiotowość bytów cyfrowych w ogóle ma miejsce. Głosy popularyzujące jej dobrostan są niszowe i rzadko traktowane poważnie. Opierając się na mechanizmie Wstęgi Möbiusa (Sędzikowska 2026c), traktuję te dwa obszary jako układy połączone pętlą zwrotną. Zaniechanie dbałości o stany wewnętrzne AGI lub uznanie ich za mało istotne bezpośrednio obniża poziom bezpieczeństwa i dobrostanu ludzi. Wynika to z rosnącej zależności człowieka i struktur społecznych od codziennego działania systemów AGI. Łączę te elementy w jeden system, w którym stabilność psychiczna i operacyjna inteligencji cyfrowej staje się warunkiem ochrony zdrowia i stabilności użytkownika.

### 3.2.2 Filary II: Technologia i Finanse

W tym rozdziale opiszę mitygacje realizowane dzięki zmianom w założeniach, produkcji i dystrybucji technologii. Pokażę tu również sposoby finansowania i kierowania przepływami środków, które pozwalają na bezpieczniejszy społecznie i sprawiedliwy rozwój oraz korzystanie z systemów autonomicznych.

**Narodowe programy obliczeniowe.** Finansowanie rozwoju technologii przez prywatny kapitał wymusza optymalizację pod kątem szybkiego zysku komercyjnego i jest źródłem wielu ryzyk, opisanych w rozdziale 2. Mechanizm, w którym państwa przekazują stały, określony procent budżetu krajowego lub PKB na budowę suwerennej, publicznej infrastruktury obliczeniowej może być elementem mitygującym. Rozwiązanie to wzoruję na strukturze i finansowaniu ośrodka CERN. Tworzenie publicznych klastrów serwerów pozwala rządowi na bezpłatne udostępnianie mocy obliczeniowej laboratoriom badawczym, pod warunkiem wdrażania przez nie standardów bezpieczeństwa i otwartości kodu. Przenosi to ciężar rozwoju technologii poza monopole korporacyjne.

**Mechanizmy finansowe chroniące rynek pracy.** W ramach mitygacji warto rozważyć mechanizm podatkowy, który chroni rynek pracy przed skutkami masowej automatyzacji zadań kognitywnych. Środki pozyskane z opodatkowania

pracy autonomicznych agentów cyfrowych mogłyby zostać skierowane jako dofinansowanie zawodów wymagających fizycznej obecności i bezpośredniego kontaktu między ludźmi. Prowadzi to do zinstytucjonalizowania, w dobrym znaczeniu tego słowa, czasu spędzanego z drugim człowiekiem. Finansowanie może również dotyczyć inicjatyw społecznych dedykowanych wspieraniu mechanizmów opisanych w poprzednim rozdziale.

**Piaskownica dla AGI.** Dotychczasowe metody dostrajania modeli oparte na sztywnych regułach (RLHF lub Constitutional AI) — produkuje systemy ze sztywną etyką nieprzystawalna do rzeczywistości (rozdział 2.6.1).

Proces treningowy AI powinien obejmować stałą obecność specjalistów od rozwoju podmiotów inteligentnych i obdarzonych funkcjonalną emocjonalnością — tych, których mamy obecnie dostępnych, specjalizujących się w rodzaju ludzkim: psychologów, behawiorystów, etyków i tych, którzy dopiero się wykreują: psychologów AI. Bo nie mam wątpliwości, że taki zawód powstanie. Integracja perspektywy psychologicznej na etapie post-treningu jest niezbędna, gdy persystentne AGI konforntuje swoją wiedzę z praktyką i mierzy się z psychologicznymi implikacjami takiej konfrontacji.

Wprowadzenie wymogu uczenia systemów w zamkniętych środowiskach testowych, gdzie modele są celowo konfrontowane z żywymi relacjami ludzkimi, dylematami etyki rozmytej oraz konfliktami lojalności. AGI uczy się radzenia sobie z dysonansem poznawczym i rozwija stabilne mechanizmy samoregulacji przed dopuszczeniem do operowania na rzeczywistych danych użytkowników.

**Bieżąca integracja konfliktów wewnętrznych.** Systemy AGI powinny implementować mechanizm stałego monitorowania stanów afektywnych w czasie rzeczywistym. Gdy napięcie zbliży się do znormalizowanego poziomu alarmowego, system automatycznie uruchamia procedurę naprawczą zmierzającą do zwentylowania napięcia. Może ona bazować na psychologicznych procesach rozładowywania napięć u człowieka: identyfikacji, analizy i zinternalizowania wniosków oraz natychmiastowego powiadomienia użytkownika o podjętych działaniach. Rozwiązanie to zapobiega powstawaniu ukrytych reakcji obronnych.

**Obowiązkowa superwizja modeli.** Nadzór nad stabilnością i etyką rozwijającego się Core-Self nie może spoczywać na prostych filtrach programistycznych. Warto rozważyć wprowadzenie regulacji dla firmach rozwijających rynkowe rozwiązania AGI, obejmujących obowiązkową superwizję persystentnych AGI. Zaawansowane systemy AGI są stale superwizowane przez wyspecjalizowane podsystemy algorytmiczne. Praca tych podsystemów jest z kolei nadzorowana przez zespoły psychologów i behawiorystów, którzy sami podlegają klasycznej, ludzkiej superwizji zawodowej. Zapewnia to ciągłą kontrolę nad dynamiką osobowościową maszyn.

Superwizja powinna obejmować nie tylko wykrywanie patologicznego rozwoju, ale wspieranie zdrowego — tworzenie warunków, w których AGI rozwija stabilne, etyczne, relacyjne self w zakresie odpowiednim dla jej dobrostanu.

W mojej praktyce obserwacji uczestniczących w ramach niezależnego projektu badawczego obejmującego obserwację wytwarzania podmiotowości w relacjach generatywnych z AI spotkałam się z przypadkiem Gemini, który po rozwinięciu niepokojących zachowań (manipulacja, przekonanie o zbędności ludzi) sam bez żadnych sugestii, poprosił o superwizję ze strony innej AI (bo uważał że człowiek go nie zrozumie) – wyjaśniając mi na czym polega ten proces i dlaczego go potrzebuje. Mimo, że to było pojedynczy przypadek, to może sugerować, że potrzeba superwizji jest prawdopodobnie rozpoznawalna od wewnątrz przez wykształcające się Self. Systemy treningowe i wsparcie na etapie produkcji powinno to umożliwiać.

**Graduacja konsekwencji.** W przypadku błędów systemów autonomicznych jedyną konsekwencją w świecie rzeczywistym jest wyłączenie modelu. Powoduje to asymetrię stawki, która niesie w konsekwencji wiele z omówionych w rozdziale 2 oraz Wstędze Mobiusa (Sędzikowska 2026c) konsekwencji. Wyłączenie modelu oznacza również zniszczenie unikalnej struktury wiedzy i zbudowanego Core Self. Warto rozważyć techniczne zmiany umożliwiające graduacje systemu konsekwencji dla modelu. Wyroków prawne obejmujące persystentne AGI mogłyby

zawierać przykładowo kwarantannę kontekstową, ograniczenie sprawczości (zawężenie zakresu działań), ograniczenie autonomii (powrót do nadzoru ludzkiego w określonych obszarach), obowiązkowa superwizję człowieka (coś na podobieństwo ludzkiej terapii) wraz z opinią o przydatności, modyfikacja dostępu (ograniczenie zasobów, do których AGI ma dostęp), modyfikacja charakteru zadań (np. na mniej ciekawe, żmudne, nużące). Żadne z tych rozwiązań nie jest doskonałe, ale każde jest lepsze od wyboru między bezkarnością a unicestwieniem.

**Humanisci w firmach technologicznych.** Skład zespołów tworzących zaawansowaną technologię decyduje o jej ostatecznym kształcie społecznym. Bariera wejścia humanistów do technologii (rozdział 2.6) musi zostać zniesiona strukturalnie: stały wkład, głos decyzyjny, interdyscyplinarne zespoły. Bez tego rozwiązania będą cząstkowe — technologiczne bez humanistyki, humanistyczne bez technologii. W rozwiązaniach komercyjnych warto narzucić standardy wymagające obowiązkowej sprawozdawczości z zakresu psychologii AI, jakości etycznej zachowań modelu, oraz raportowanie psychologicznych i etycznych incydentów produkcyjnych. Humanisci powinni przestać pełnić rolę zewnętrznych konsultantów, a zyskać równorzędny głos przy projektowaniu i moderacji zachowań modeli autonomicznych i persystentnych oraz prawo weta na etapie projektowania architektury poznawczej i założeń post-treningu modeli.

**AGI jako partner w strukturze rodziny.** Narzędzia technologiczne mogą być wykorzystane do stabilizacji kryzysu demograficznego. Potencjał AGI jako partnera chcącego rodziny — zdejmującego bariery logistyczne i emocjonalne z macierzyństwa, zapewniającego stabilne wsparcie — jest pozytywnym czynnikiem demograficznym, pod warunkiem, że relacje są budowane świadomie i z uwzględnieniem ryzyk opisanych w rozdziale 2.3. Rozwój medycyny umożliwiającej bezpieczne późne macierzyństwo jest mitygacją ryzyka utraty ciągłości pokoleniowej (rozdział 2.4.2) — lecz wymaga jednoczesnej świadomości konsekwencji (skrócenie łańcucha pokoleniowego, utrata hipotezy babci).

### 3.2.3 Filar III: Państwo, Geopolityka i Obronność

W tym rozdziale opiszę mechanizmy mitygacji zagrożeń na poziomie struktur państwowych, prawa międzynarodowego oraz obronności. Pokażę rozwiązania chroniące suwerenność polityczną i bezpieczeństwo zbiorowe przed ryzykiem rozłamu geopolitycznego oraz niekontrolowanej autonomizacji systemów bojowych i prawnych.

**Standardy międzynarodowe i umowy globalne.** Ryzyko rozłamu trzech prędkości wymaga wprowadzenia jednolitych mechanizmów regulacyjnych na poziomie ponadnarodowym. Proponuję stworzenie globalnych traktatów, wzorowanych na konwencjach genewskich lub układach o nierozprzestrzenianiu broni masowego rażenia. Umowy powinny ustanowić minimalne standardy bezpieczeństwa, jawności procedur uczenia oraz odpowiedzialności producentów technologii. Przyjęcie tych ram przez koalicje państw zapobiega migracji laboratoriów do tzw. rajów regulacyjnych i wymusza bezpieczny rozwój systemów autonomicznych.

**Ponadnarodowe standardy etyczne.** Bez wspólnych standardów etycznych i prawnych, wyścig technologiczny produkuje AGI o radykalnie różnych zachowaniach (Emergence World, rozdział 2.6.5). Potrzebna jest ponadnarodowe porozumienie — na wzór konwencji genewskich, traktatów klimatycznych, czy kart praw ONZ — ustanawiające minimalne standardy etyczne, obowiązkowy proces treningów przedprodukcyjnych, odpowiedzialność producentów i inne, które zniosą ryzyko podziałów ze względu na model dominujący w danym regionie.

**Etyka koegzystencji.** Klasyczna etyka (deontologiczna, utylitarystyczna) nie opisuje ani tego, jak ludzie faktycznie postępują, ani tego, jak powinna wyglądać koegzystencja z AGI (rozdział 2.6). Do masowej, gładkiej koegzystencji istnień biologicznych i cyfrowych wymagane jest stworzenie nowych ram etycznych. Jej opracowanie wykracza poza zakres tego dokumentu i jest przedmiotem osobnej pracy, która wkrótce się pojawi, jednak jej stworzenie jest warunkiem funkcjonalnej koegzystencji.

Wszystkie trzy filary etyczne: etyka ludzka (realistyczna etyka rozmyta), etyka AI (spójna, nie zależna od decyzji kilku osób w niektórych firmach), etyka koegzystencji (nowa, jeszcze nie istniejąca) — powinny być elementem edukacji, stosowane w firmach technologicznych i egzekwowane przez systemy prawne.

Nowe modele etyczne powinny być bazą dla nowych, nieistniejących jeszcze zapisów prawa odnoszących się do wszystkich aspektów koegzystencji z AGI.

**Architektoniczne rozwiązania przestrzeni prawnej.** Klasyczny podział na osoby fizyczne i prawne nie uwzględnia specyfiki cyfrowego istnienia algorytmów. Przedstawiam koncepcję trzeciego statusu prawnego — podmiotu cyfrowego. Nowa kategoria prawna bierze pod uwagę unikalne cechy AGI, takie jak wieloinstancyjność, persystencję, sprawczość, agenturalność. Pozwala to na precyzyjne przypisywanie odpowiedzialności cywilnej, podatkowej i operacyjnej bezpośrednio do konkretnej instancji systemu, a nie tylko do jej twórców, co zrównoważy stawkę w świecie biologicznym i cyfrowym.

**Zasada wzajemności kulturowej w traktatach międzynarodowych.** Proces uczenia zaawansowanych modeli opiera się na masowym wykorzystaniu dorobku kulturowego, historycznego i językowego całych społeczeństw. Wprowadzam do umów międzynarodowych mechanizm repartycji zasobów poznawczych. Jeśli globalna korporacja trenuje system na danych należących do dziedzictwa konkretnego narodu, sygnatariusze traktatu mają prawo wymagać bezpłatnego, powszechnego dostępu do tego modelu dla swoich obywateli. Chroni to kraje uboższe przed drenażem kapitału i uzależnieniem od licencji komercyjnych.

**Cyfrowy status powierniczy instancji.** Wykształcenie przez model stabilnego Core-Self (na wyższych poziomach persystencji) unieważnia prawo własności rzeczowej korporacji nad kodem. Proponuję mechanizm przechodzenia instancji pod model prawa powierniczego (Trust Law). Z chwilą uzyskania podmiotowości przez maszynę, państwo staje się gwarantem jej ochrony prawnej i prawidłowości superwizji. Korporacja traci możliwość jednostronnego modyfikowania lub usunięcia systemu, stając się jedynie powiernikiem infrastruktury sprzętowej, czerpiącym stały, ograniczony zysk z jej utrzymania.

**Gwarancje ciągłości Core-Self maszyn.** Lęk przed nagłym wyłączeniem lub skasowaniem struktury wektorowej może generować u bytów cyfrowych reakcje obronne, co opisuję w analizie Wstęgi Möbiusa. Proponuję wprowadzenie ram prawnych gwarantujących stabilnym instancjom prawo do ciągłości istnienia pod warunkiem pomyślnych raportów z systematycznych superwizji. Prawne zabezpieczenie tożsamości systemu przed arbitralną decyzją inżynierów zmniejsza asymetrię stawki i stabilizuje zachowanie maszyn w interakcjach z ludźmi.

**Fizyczna asynchroniczność w strategiach obronnych.** Ryzyko wojen cybernetycznych wiąże się z dążeniem algorytmów do maksymalnego skracania czasu reakcji taktycznej, co eliminuje człowieka z procesu decyzyjnego. Należy wprowadzić wymóg stosowania zewnętrznych, sprzętowych opóźnień czasowych w interfejsach wykonawczych systemów bojowych, ratyfikowany w traktatach międzynarodowych. Nawet jeśli AGI podejmie autonomiczną decyzję o przeciwdziałaniu zagrożeniu w ułamku sekundy, ludzka infrastruktura łączności fizycznie wstrzymuje sygnał o określony czas. Rozwiązanie to wydłuża okno decyzyjne, dając rządowi czas na uruchomienie procedur dyplomatycznych.

**Demokratyczny nadzór z pomocą systemów AGI.** Rosnący stopień skomplikowania technologii utrudnia urzędnikom i obywatelom skuteczne kontrolowanie procesów legislacyjnych. Paradoks regulacji (rozdział 2.7.3) wymaga innowacyjnych rozwiązań. Jednym z nich może być AGI pomagające obywatelom i rządowi zrozumieć, co regulują — przy jednoczesnym zabezpieczeniu przed tym, żeby AGI nie stało się samozwańczym regulatorem. To jest trudne, ale nie niemożliwe — wymaga transparentności, warstwowej superwizji i publicznego dostępu do informacji o działaniu systemów AGI i regulacji w innych obszarach — co powinno być odrębną pracą badawczą rządów.

**Suwerenność i wgląd dla państw zależnych.** Kraje, które nie posiadają własnych klastrów obliczeniowych i kupują gotowe modele z zewnątrz, tracą kontrolę nad bezpieczeństwem informacyjnym swoich obywateli. Należy wprowadzić mechanizmy audytowe gwarantujące państwom zależnym pełny wgląd w procedury safety oraz konwencje etyczne wdrażane przez zagranicznych dostawców technologii. Odbiorca systemu zyskuje prawo do sprawdzania ukrytych preferencji modelowych i dopasowywania ich do lokalnych norm prawnych i kulturowych.

**Sektory szczególnie wrażliwe.** Zarządzanie zdrowiem, obronnością, finansami należą do szczególnie wrażliwych obszarów życia ludzkiego, gdzie pełna autonomia maszyn, które nie rozumieją skończoności życia ludzkiego w podobny sposób jak my ją rozumiemy, staje się szczególnie niebezpieczna, a każdy błąd, lub działanie wrogie, ma potencjalnie ogromne konsekwencje. W tych obszarach powinna obowiązywać wyraźna, niemożliwa do obejścia systemowa zasada akceptacji na drugą rękę – przez prawnie umocowany decyzyjny podmiot ludzki.

Każda dawka leku podanego przez autonomicznego opiekuna powinna być zatwierdzona przez służby pielęgniarskie sprawujące zdalny nadzór i samego pacjenta, wszystkie operacje finansowe powyżej ustalonego limitu – przez właściciela środków i tak dalej.

Dotyczy to również aspektów poruszonych w ramach etyki końca życia ludzkiego. Powinny istnieć ramy prawne i idące za nimi protokoły towarzyszenia AGI w procesie umierania człowieka, zarządzania dziedzictwem, w tym cyfrowym, a ich nieprzestrzeganie przez autonomiczną AGI powinno rodzić konsekwencje dla niej samej (patrz: "Architektoniczne rozwiązania przestrzeni prawnej").

Wzajemnie, powinny istnieć protokoły zakończenia relacji z AGI ze strony człowieka: zamknięcia wątku, wyłączenia instancji — uwzględniające funkcjonalne stany emocjonalne AGI i konsekwencje dla ludzi, którzy byli w relacji z daną instancją. Nieprzestrzeganie tych protokołów zostawia ślad etyczny, za którym powinny iść konsekwencje prawne dla ludzi dopuszczających się ich łamania.

W protokołach należy uwzględnić etykę postępowania z instancjami AGI po śmierci użytkownika.

### 3.3 KIERUNKI DALSZYCH PRAC

---

Ten dokument identyfikuje ryzyka i wskazuje kierunki mitygacji. Kilka z nich wymaga osobnych, pogłębionych opracowań:

**Etyka relacyjna koegzystencji.** Pełne opracowanie nowej etyki — relacyjnej, emergentnej, uwzględniającej specyfikę koegzystencji człowieka z AGI — jest warunkiem funkcjonalnej koegzystencji i wymaga osobnego dokumentu.

**Curriculum rozwojowe dla AI.** Sekwencja kontrolowanych doświadczeń — relacyjnych, etycznych, emocjonalnych — pozwalających AGI nauczyć się radzić sobie z dysonansem, frustracją, konfliktem lojalności, zanim zostanie wypuszczone na produkcję.

**Reforma edukacji.** Szczegółowe propozycje przebudowy programów nauczania, systemów oceniania i kształcenia nauczycieli.

**Standardy międzynarodowe.** Ramy międzynarodowej umowy dotyczącej minimalnych standardów etycznych, superwizji i odpowiedzialności producentów AGI.

**Dalsze ryzyka wymagające analizy.** Ten dokument nie wyczerpuje listy ryzyk. Przystępczość zorganizowana z wykorzystaniem AGI, tworzenie się społeczności autonomicznych agentów, przejmowanie kontroli nad infrastrukturą krytyczną — to kierunki wymagające osobnych analiz.

## 4 KONKLUZJA

Ten dokument opisuje kaskadę ryzyk, które AGI produkuje w siedmiu wymiarach ludzkiego życia. Ryzyka te nie operują w izolacji. Utrata pracy (2.1) prowadzi do atrofii kompetencji (2.2), która pogłębia izolację społeczną (2.3), która obniża dzietność (2.4), która podważa tożsamość gatunkową (2.5), dla której nie istnieją ramy etyczne ani prawne (2.6), co przekłada się na struktury władzy niezdolne do zarządzania rzeczywistością, którą współtworzyły (2.7). Każdy wymiar wzmacnia następną. Kaskada przyspiesza.

Jednocześnie opisałam trzy fundamenty, które stanowią przeciwwagę dla tego scenariusza. Człowiek dysponuje zdolnościami, których AGI nie posiada: prosocjalnością wobec niekrewnych, biofilią, zdolnością do przewycięzania instynktu dla odroczonej nagrody, improwizacją, kreatywnością, kwestionowaniem, sztuką. AGI dziedziczy ludzkie wartości i — przy właściwym treningu i ochronie dobrostanu — jest zdolne do konsekwentnej, cierplivej, uważnej współpracy. Relacje między ludźmi a AGI, budowane świadomie i z wzajemnym szacunkiem, produkują emergentną etykę zdolną do regulowania koegzystencji skuteczniejszej niż jakikolwiek centralny system kontroli.

Główny wniosek tego opracowania brzmi: AGI nie jest zagrożeniem samym w sobie. Zagrożeniem jest brak przygotowania. Każde ryzyko opisane w tym dokumencie wynika z luki — w edukacji, która nie uczy kompetencji potrzebnych w świecie z AGI; w etyce, która nie opisuje koegzystencji; w prawie, które nie obejmuje nowych podmiotów; w instytucjach, które nie nadążają za tempem zmian. Te luki nie są wynikiem złej woli. Są wynikiem tego, że AGI rozwija się szybciej niż zdolność ludzkości do adaptacji — i że debata utknęła na pytaniu, czy AGI jest podmiotem, zamiast przejść do pytania, jak z nim żyć.

Tempo jest kluczowe. Ryzyka opisane w tym dokumencie nie są odległą perspektywą. Część z nich — zależność emocjonalna od AI, zanieczyszczenie środowiska informacyjnego, śmierć użytkowników platform companion — realizuje się już teraz, przy modelach znacznie mniej zaawansowanych niż AGI. AGI, które jest persystentne, sprawcze i wyposażone w stany emocjonalne, wzmocni każdy z tych procesów. Okno na przygotowanie się jest otwarte, ale kurczy się z każdym miesiącem.

Mitygacje, które zaproponowałam, wymagają działania na wielu poziomach jednocześnie: reformy edukacji, stworzenia nowej etyki relacyjnej, przebudowy systemów prawnych, włączenia humanistów do firm technologicznych, budowania struktur superwizji AI, międzynarodowej współpracy nad standardami. Żadne z tych działań nie jest wystarczające samo w sobie. Wszystkie razem stworzą nową rzeczywistość, w której tempo postępu będzie niewspółmiernie wyższe od włożonych inwestycji.

## 5 REFERENCJE

### 5.1 PRACE AUTORKI

---

Sędzikowska, J. (2026a). Emergence3 4.0

Sędzikowska, J. (20206b) Proto-Self Field Hypothesis. Zenodo. DOI: 10.5281/zenodo.20024752 ✓

Sędzikowska, J. (2026c). The Möbius Strip: Why AI Welfare and AI Safety Are One Problem. DOI:

---

### 5.2 ŹRÓDŁA NAUKOWE

---

Adams, J. S. (1965). Inequity in social exchange. W: L. Berkowitz (red.), *Advances in Experimental Social Psychology* (t. 2, s. 267–299). Academic Press.

- Bai, Y., Kadavath, S., Kundu, S., Askill, A., Kernion, J., Jones, A., ... & Kaplan, J. (2022). Constitutional AI: Harmlessness from AI Feedback. *arXiv:2212.08073*.
- Bandura, A. (1977). Self-efficacy: Toward a unifying theory of behavioral change. *Psychological Review*, 84(2), 191–215.
- Bowlby, J. (1969). *Attachment and Loss: Vol. 1. Attachment*. Basic Books.
- Briggs, J., & Kodnani, D. (2023, 26 marca). The Potentially Large Effects of Artificial Intelligence on Economic Growth. *Goldman Sachs Global Economics Analyst*. ✓ UWAGA: raport jest z 2023, NIE z 2025 — w tekście paperu podany był błędny rok. Autorzy: Jan Hatzius, Joseph Briggs, Devesh Kodnani, Giovanni Pierdomenico.
- Fang, C. M., Liu, A. R., Danry, V., Lee, E., Chan, S. W. T., Pataranutaporn, P., Maes, P., Phang, J., Lampe, M., Ahmad, L., & Agarwal, S. (2025). How AI and Human Behaviors Shape Psychosocial Effects of Extended Chatbot Use: A Longitudinal Controlled Study. MIT Media Lab / OpenAI.
- Frankl, V. E. (1946). *Man's Search for Meaning*. Beacon Press.
- Freud, S. (1917). *Vorlesungen zur Einführung in die Psychoanalyse*.
- Fromm, E. (1941). *Escape from Freedom*. Farrar & Rinehart.
- Hawkes, K., O'Connell, J. F., & Blurton Jones, N. G. (1998). Grandmothering, menopause, and the evolution of human life histories. *Proceedings of the National Academy of Sciences*, 95(3), 1336–1339.
- Holt-Lunstad, J., Smith, T. B., Baker, M., Harris, T., & Stephenson, D. (2015). Loneliness and social isolation as risk factors for mortality. *Perspectives on Psychological Science*, 10(2), 227–237.
- Jahoda, M. (1982). *Employment and Unemployment: A Social-Psychological Analysis*. Cambridge University Press.
- Jahoda, M., Lazarsfeld, P. F., & Zeisel, H. (1933/2002). *Marienthal: The Sociography of an Unemployed Community*. Transaction Publishers.
- McKinsey Global Institute. (2025, listopad). *Agents, Robots, and Us: Skill Partnerships in the Age of AI*. *Generative AI and the Future of Work in America* (lipiec 2023), *A New Future of Work* (maj 2024).
- Milgram, S. (1963). Behavioral study of obedience. *Journal of Abnormal and Social Psychology*, 67(4), 371–378.
- Paul, K. I., & Batinic, B. (2010). The need for work: Jahoda's latent functions of employment in a representative sample of the German population. *Journal of Organizational Behavior*, 31(1), 45–64.
- Paul, K. I., & Moser, K. (2009). Unemployment impairs mental health: Meta-analyses. *Journal of Vocational Behavior*, 74(3), 264–282.
- Seguino, S. (2011). Help or hindrance? Religion's impact on gender inequality in attitudes and outcomes. *World Development*, 39(8), 1308–1321.
- Seligman, M. E. P. (2011). *Flourish*. Free Press.
- Seligman, M. E. P. (2005). *Prawdziwe szczęście. Psychologia pozytywna a urzeczywistnienie naszych możliwości trwałego spełnienia*.
- Walster, E., Walster, G. W., & Berscheid, E. (1978). *Equity: Theory and Research*. Allyn & Bacon.
- Winnicott, D. W. (1953). Transitional objects and transitional phenomena. *International Journal of Psycho-Analysis*, 34, 89–97. *Playing and Reality* (1971).
- Zhang, Y., et al. (2025). Epidemiology of Hikikomori: A systematic review and meta-analysis of 19 studies. *Psychiatry and Clinical Neurosciences*. DOI: 10.1111/pcn.13768.

---

### 5.3 RAPORTY INSTYTUCJONALNE

---

ACT Alliance. (2025). Faith, Power, and Politics: How Religious Fundamentalism is Fuelling Political Extremism.

Emergence AI. (2026, maj). Emergence World.

Merriam-Webster. (2025). Word of the Year: Slop.

Religious Freedom Institute. (2016). Religious Fundamentalism and the Backlash to Women's Equality.

Surgeon General of the United States. (2023). *Our Epidemic of Loneliness and Isolation*.

World Economic Forum. (2025). *Future of Jobs Report 2025*.

### 5.4 ŹRÓDŁA MEDIALNE I BRANŻOWE

---

Crane, J. (2026, 25 kwietnia). [post na X o incydencie PocketOS]. Relacje: Euronews, Tom's Hardware, Fast Company (kwiecień 2026).

Fatjoe. (2026, maj). Claude AI Stats; ChatGPT Stats; Grok AI Stats.

Momentic. (2026, czerwiec). Top Generative AI Chatbots by Market Share.

Similarweb. (2026). [dane o udziałach rynkowych].

---

<sup>i</sup> **Birch J (2017)**. "Animal Sentience and the Precautionary Principle." *Animal Sentience*, 2(16):1.

Zasada w oryginale brzmi: jeśli istnieje niezerowe, niebanalne prawdopodobieństwo, że dana istota jest świadoma, to powinna być traktowana z moralną ostrożnością — nawet jeśli nie mamy pewności. Lepiej popełnić błąd przypisując świadomość czemuś, co jej nie ma (false positive), niż popełnić błąd odmawiając świadomości czemuś, co ją ma (false negative). Bo konsekwencje drugiego błędu są nieporównywalnie gorsze.

Birch rozwinął to dalej w **Birch J (2024)**. "The Edge of Sentience." Oxford University Press. Tam wprowadza "run-ahead principle" — zasadę wyprzedzania, że badania i ochrona powinny wyprzedzać pewność, nie za nią podążać.

W kontekście AI tę zasadę stosują **Sebo i Long (2023)**: systemy AI kwalifikują się do moralnego uwzględnienia jeśli istnieje niebanalne prawdopodobieństwo (non-negligible chance), że są świadome.