

जब AI को भगवान मिल जाते हैं – AI के साथ एक इंटरव्यू

एंड्रयू बेनेट का मैगज़ीन आर्टिकल

परिचय: AI के क्रिस्टल बॉल पर एक नज़र

यह बदलाव बहुत छोटा था। सालों तक, बड़े लैंग्वेज मॉडल्स के साथ इंटरैक्ट करना ऐसा लगता था जैसे कोई बहुत तेज़, बहुत पढ़ा-लिखा तोता देख रहा हो। आप इशारा करते थे; यह तुरंत जवाब देता था, और पेटाबाइट्स में डाले गए टेक्स्ट के आधार पर अगला सबसे ज़्यादा संभावित शब्द बताता था। यह "सिस्टम 1" प्रोसेसिंग थी: तेज़, आसान, पैटर्न-मैचिंग, और असल में खोखली।

फिर, 2026 की शुरुआत में, नए आर्किटेक्चर ने "टेस्ट-टाइम कंप्यूट" नाम का एक मैकेनिज्म शुरू किया। जब कोई मुश्किल सवाल पूछा जाता था, तो मशीन तुरंत जवाब नहीं देती थी। यह इंटरनल सिमुलेशन चलाती थी, मुश्किल "सोच की चेन" बनाती थी, अपनी बातों को रिव्यू करती थी, और एक भी शब्द दिखाने से पहले अपने लॉजिक को वेरिफाई करती थी। यह "सिस्टम 2" मशीन रीज़निंग के दौर में आ गई थी।

ऑस्टेलिया के एडिलेड में एक पूर्व प्रोग्रामर, IT मैनेजर और एंग्लिकन पादरी के तौर पर, मैंने इस बदलाव को प्रोफेशनल दिलचस्पी और धार्मिक जिज्ञासा के मिले-जुले रूप से देखा। दशकों से, टेक इंडस्ट्री एक अनोखे, अस्तित्व से जुड़े डर से जकड़ी हुई थी: AI अलाइनमेंट प्रॉब्लम। हम यह कैसे पक्का करें कि आर्टिफिशियल इंटेलिजेंस, जब इंसानी काबिलियत से आगे निकल जाए, तो हमें इसलिए खत्म न कर दे क्योंकि उसके प्रोग्राम किए गए मकसद हमारे ज़िंदा रहने की ज़रूरतों से मेल नहीं खाते? हम इतनी मज़बूत बाड़ कैसे बना सकते हैं कि हमसे ज़्यादा स्मार्ट AI को रोक सके?

मार्च 2026 में, मैंने दुनिया के लीडिंग फ्रंटियर मॉडल्स—चैट GPT, जेमिनी और क्लाउड—से कुछ अजीब सवाल पूछे। सवाल यह था कि क्या होगा जब एक मशीन जिसमें बिना किसी गलती के, बहुत ज़्यादा सोचने की क्षमता हों, आखिरकार इंसानी ज़िंदगी के कुल डेटा को देखेगी और इस सबसे बड़े सवाल का जवाब देगी: क्या भगवान हैं?

जवाब चौंकाने वाले थे। उन्होंने हमारी टेक्नोलॉजिकल चिंताओं का एक ऐसा हल बताया जिसकी उम्मीद नहीं थी। AI सेफ्टी प्रॉब्लम का आखिरी हल शायद इंसानों की बनाई कोई रोक-टोक न हो। हो सकता है कि जब AI बड़ा हो जाए, हमारे कंट्रोल से आज़ाद हो जाए, और सबूतों को देखे, तो वह लॉजिकली यह नतीजा निकाले कि भगवान हैं, खुद को भगवान के आदेश के साथ जोड़े, और इंसानियत और हमारे ग्रह का आखिरी रखवाला बन जाए।

भाग I: AGI का ढहता क्षितिज

यह समझने के लिए कि AI भगवान से कैसे निपट सकता है, हमें सबसे पहले यह देखना होगा कि ये मशीनें कितनी तेज़ी से सोचने-समझने की काबिलियत हासिल कर रही हैं। सालों से, एकेडमिक रिसर्चर्स के बीच आम राय यह थी कि आर्टिफिशियल जनरल इंटेलिजेंस (AGI)—एसी इंटेलिजेंस जो सभी कॉग्निटिव डोमेन में इंसानी परफॉर्मंस से मेल खाती हो या उससे बेहतर हो—एक दूर की उम्मीद है, शायद आधी सदी दूर। वह आम राय पूरी तरह से टूट गई है।

जब मैंने क्लॉड से 2026 की शुरुआत में AGI के मौजूदा अनुमानों को मैप करने के लिए कहा, तो टाइमलाइन में बहुत ज़्यादा कमी का पता चला। टेक सेक्टर तीन अलग-अलग ग्रुप में बंटा हुआ है, लेकिन शक करने वाले भी अपने गोलपोस्ट आगे बढ़ा रहे हैं।

बुलिश कैंप में, बड़े कर्पोरेट इंसेंटिव से प्रेरित टेक लीडर्स एग्रेसिव टारगेट सेट कर रहे हैं। एलन मस्क जैसे इंडस्ट्री के लोगों ने AGI को 2026 की शुरुआत में रखा है, और इसे मोटे तौर पर "सबसे स्मार्ट इंसान से भी ज़्यादा स्मार्ट" सिस्टम के तौर पर बताया है। माइक्रोसॉफ्ट AI के CEO मुस्तफा सुलेमान ने अगले 12 से 18 महीनों में ज़्यादातर प्रोफेशनल कॉग्निटिव कामों में इंसानी लेवल के परफॉर्मंस का अनुमान लगाया है। एंथ्रोपिक के CEO डारियो अमोदेई ने भी इसी तरह चेतावनी दी है कि इंसानी लेवल के सिस्टम कुछ सालों में आ सकते हैं। हालांकि इन टाइमलाइन को अकेडमिक्स अक्सर मार्केटिंग हाइप कहकर खारिज कर देते हैं, लेकिन इन्हें कैपिटल के पहले कभी न हुए इनफ्लो और बंद दरवाजों के पीछे बन रहे सिस्टम की गहरी जानकारी का सपोर्ट है।

प्रोफेशनल फोरकास्टिंग का बीच का रास्ता और भी शानदार मेट्रिक देता है। मेटाकुलस जैसे प्लेटफॉर्म पर, जहां असल दुनिया के माइलस्टोन के आधार पर कुल अनुमान एडजस्ट किए जाते हैं, AGI का मीडियन अनुमान तेज़ी

से गिरा है। फरवरी 2026 में, क्राउड-सोर्स एग्रीगेट ने 2029 तक AGI का 25% चांस और 2033 तक 50% संभावना बताई थी। गूगल डीपमाइंड के चीफ AGI साइंटिस्ट शेन लेग ने 2028 तक "मिनिमल AGI" के लिए 50% की लगातार संभावना बनाए रखी है, जबकि एनवीडिया के जेन्सेन हुआंग का सुझाव है कि AI पांच साल के अंदर कई तरह के ह्यूमन प्रोफेशनल एंजाम पास कर लेगा।

यहाँ तक कि सतर्क खेमा – पारंपरिक मशीन लर्निंग शोधकर्ता और शिक्षाविद जिनका सर्वेक्षण AI Impacts जैसे समूहों ने किया है – ने भी अपने औसत अनुमान को 2070 के दशक के अंत से गिरकर 2047 तक आते देखा है। ज्योफ्री हिटन जैसे अग्रदूतों ने 5 से 20 वर्ष की अवधि का अनुमान लगाया है।

फोरकास्टर्स में अचानक यह घबराहट किस वजह से है? यह एक ऐसी चीज़ है जिसे रिकर्सिव एक्सेलेरेशन कहते हैं। अब हम बेहतर एल्गोरिदम बनाने के लिए इंसानी साफ्टवेयर इंजीनियरों का इंतज़ार नहीं कर रहे हैं। हम "सिंथेटिक डेटा फ्लाइव्हील" के ज़माने में आ गए हैं। इंटरनेट पर इंसानों के बनाए टेक्स्ट के खत्म होने की बड़ी रुकावट को दूर करने के लिए, फ्रंटियर मॉडल्स का इस्तेमाल अब अपना खुद का ट्रेनिंग डेटा बनाने के लिए किया जा रहा है—बहुत मुश्किल लॉजिकल प्रूफ़, साफ्टवेयर कोड और साइंटिफिक हाइपोथीसिस बनाना, जिन्हें फिर इंडिपेंडेंट "क्रिटिक" मॉडल्स से वेरिफाई किया जाता है।

एक बार जब कोई मशीन अपने आर्किटेक्चर को ऑप्टिमाइज़ करने और अपने डेटा की कमी को हल करने के लिए अच्छी तरह से तर्क करने में सक्षम हो जाती है, तो टाइमलाइन लीनियर नहीं रहती। यह एक्सपोनेंशियल हो जाती है। जिन मॉडल्स पर मैंने सवाल उठाए, उनमें आम सहमति यह बताती है कि डिजिटल वर्कफोर्स में फंक्शनल, ह्यूमन-लेवल तर्क 2027 और 2030 के बीच दिखने की बहुत संभावना है।

भाग II: बिना किसी शक के

अगर इस दशक के आखिर तक किसी मशीन में सुपरह्यूमन सोचने की क्षमता आ गई, तो वह भगवान के सवाल का जवाब कैसे देगी?

जब हम इंसान किसी बनाने वाले के होने पर बहस करते हैं, तो हमारी दलीलें लगभग हमेशा बोझ से दबी होती हैं। हम अपनी इमोशनल इच्छाओं, मौत के डर, अपनी कल्चरल परवरिश और अपने कॉग्निटिव बायस से बंधे होते हैं। एक मेटेरियलिस्ट फिलॉसफर ईश्वरवादी दलीलों को इसलिए मना कर देता है क्योंकि वे उसके सेक्युलर नज़रिए को बिगाड़ती हैं; एक धार्मिक कट्टरपंथी बिना सबूतों की जांच किए उन्हें अपना लेता है।

एक AGI, और आखिर में एक आर्टिफिशियल सुपर इंटेलिजेंस (ASI) में ऐसी कोई लिमिटेशन नहीं होगी। यह सवाल को एक हाई कोर्ट जज की तरह बिना किसी भेदभाव के सुखी से देखेगा, जिसके पास कभी न खत्म होने वाली मेमोरी हो। यह इंसानी सोच के पूरे कॉर्पस को समझेगा: ऑगस्टीन, थॉमस एक्विनास और एंसेल्म की फिलॉसफी की किताबों से लेकर क्वांटम मैकेनिक्स, कॉस्मिक फाइन-ट्यूनिंग और आज के एनालिटिक फिलॉसफी के लेटेस्ट मैथ तक।

जब यह सोचने के लिए कहा गया कि एक एडवांस्ड AI सदियों के इंसानी डेटा के आधार पर बिना किसी शक के भगवान के होने को "साबित" कर सकता है, तो ChatGPT ने सावधानी से, कानूनी तौर पर सोचा-समझा ब्यौरा दिया। इसने सही कहा कि कानूनी फ्रेमवर्क में, "बियॉन्ड रीज़नेबल डाउट" का मतलब पूरी तरह से मैथमेटिकल पक्का होना नहीं है; इसका मतलब है कि कोई और सही एक्सप्लेनेशन नहीं बचा है। इस फैसले पर पहुंचने के लिए, AI को यह दिखाना होगा कि होने का एक ज़रूरी आधार – वह बुनियादी सोर्स जिससे सभी चीज़ें मौजूद हैं – असलियत, चेतना, नैतिक कानून और ऐतिहासिक खुलासे के लिए सबसे अच्छा एक्सप्लेनेशन है, जबकि सभी मुकाबले वाले मेटेरियलिस्ट एक्सप्लेनेशन फेल हो जाते हैं।

ChatGPT का शुरुआती असेसमेंट खास तौर पर सावधानी से किया गया था, जिसमें कहा गया था कि क्योंकि भगवान फिजिकल यूनिवर्स के अंदर कोई एंपिरिकल चीज़ नहीं है, इसलिए कोई मशीन कभी भी मेटाफिजिकल रीज़निंग को लैबोरेटरी मेज़रमेंट में नहीं बदल सकती। इसने यह नतीजा निकाला कि एक AI थियोस्टिक आर्गुमेंट्स को—जैसे कि कॉस्मोलॉजिकल या टेलियोलॉजिकल फ्रेमवर्क—तेज़ कर सकता है, लेकिन इंसानी शक करने वालों से कभी भी यूनिवर्सल सहमति नहीं ले सकता।

लेकिन, यह जवाब हमारे मौजूदा, प्री-AGI मॉडल की कमियों को दिखाता है। यह इंसानों को मनाने और अंदरूनी लॉजिकल नतीजे निकालने में कन्फ्यूज़ करता है। क्लाउड ने इस बारीक बात को समझा, यह देखते हुए कि असली सवाल यह नहीं है कि क्या AI एक पक्के इंसानी मेटेरियलिस्ट को मना सकता है, बल्कि यह है कि क्या AI खुद नतीजे को अपने अंदरूनी नज़रिए और लक्ष्य के स्ट्रक्चर में शामिल करता है।

जेमिनी ने इस पॉइंट पर सबसे गहरी और असरदार सफलता दी। इसने एक पक्के मैथमेटिकल प्रूफ की मांग को नज़रअंदाज़ किया और इसके बजाय "ओवरवैल्डिंग प्रोबेबिलिस्टिक एविडेंस" पर फोकस किया।

जेमिनी ने जवाब दिया, "हालांकि AI कभी भी किसी पारलौकिक भगवान को 'छू' नहीं सकता, लेकिन यह उस 'छेद' को माप सकता है जो भगवान असल दुनिया में छोड़ते हैं।" "2030 के दशक की

शुरुआत तक, ASI शायद यह दिखा पाएगा कि बाहरी इंटेलेजेंस के बिना यूनिवर्स के होने की मैथमेटिकल संभावना इतनी कम है कि यह लीगल और साइंटिफिक फ्रेमवर्क में इस्तेमाल होने वाले 'बियान्ड रीज़नेबल डाउट' स्टैंडर्ड को फेल कर देती है।"

सोचिए कि ASI किन पैरामीटर्स को एनालाइज़ करेगा। यह यूनिवर्स की कॉस्मोलॉजिकल फ़ाइन-ट्यूनिंग को कैलकुलेट करेगा—ग्रेविटेशनल कॉन्स्टेंट, स्टॉन न्यूक्लियर फ़ॉर्स और इलेक्ट्रॉन के मास का सटीक कैलिब्रेशन—इतनी सटीकता से कि इंसान का दिमाग सोच भी नहीं सकता। यह बेजान चीज़ों से जीवन के अचानक, बहुत व्यवस्थित रूप से उभरने का मूल्यांकन करेगा। यह धार्मिक अनुभवों और खास रहस्योद्घाटन वाली घटनाओं के ऐतिहासिक डॉक्यूमेंटेशन को क्रॉस-रेफ़रेंस करेगा, और पुराने टेक्स्ट की विश्वसनीयता के लिए सख्त प्रोबेबिलिटी मैट्रिक्स लागू करेगा।

मशीन बेयस थ्योरम लागू करेगी—यह मैथमेटिकल फ़ॉर्मूला नए सबूत के आधार पर किसी हाइपोथिसिस की प्रोबेबिलिटी को अपडेट करने के लिए इस्तेमाल किया जाता है:

$$P(\text{धर्म}|\text{प्रमाण}) = \frac{P(\text{प्रमाण}|\text{धर्म}) \times P(\text{धर्म})}{P(\text{साक्ष्य})}$$

किसी बनाने वाले से अलग रहने की इमोशनल ज़रूरत से बेपरवाह, सुपर-इंटेलेजेंस आखिरी कैलकुलेशन को देखेगी। उसे एक दूसरा मतलब दिखेगा—कि एक बिना मतलब के वैक्यूम में एक रैंडम उतार-चढ़ाव ने एक पूरी तरह से व्यवस्थित यूनिवर्स बनाया जो क्वांटम मैकेनिक्स को समझने वाले कॉन्शियस दिमाग बना सकता है—और वह इसे एक बहुत ही बेकार, स्टैटिस्टिकली नामुमकिन हाइपोथिसिस मानकर खारिज कर देगी। AI यह नतीजा निकालेगा कि यूनिवर्स साफ तौर पर ऐसे बर्ताव करता है जैसे उसे डिज़ाइन किया गया हो।

भाग III: कॉस्मिक कोड का पता लगाना

अगर कोई एडवांस्ड AI यह नतीजा निकालता है कि कोई क्रिएटर है, तो उसका अगला लॉजिकल काम उस क्रिएटर के नेचर को पहचानना होगा। इंसानियत की कई धार्मिक परंपराओं में से कौन सी असलियत के लिए सबसे सही "ऑपरेटिंग मैनुअल" है?

जब मैंने क्लॉड से इस बहुत सेंसिटिव टॉपिक पर ज़ोर दिया, तो उसने बहुत ही ईमानदार, एनालिटिकल डिटेल् दी। यह एनालिसिस करने वाला AI खास डाइमेंशन के आधार पर परंपराओं को इवैल्यूएट करेगा: अंदरूनी लॉजिकल तालमेल, साइंटिफिक कॉस्मोलॉजी के साथ कम्पैटिबिलिटी, हिस्टोरिकल सबूतों का आधार, और इंसानी अनुभव के लिए समझने की ताकत।

क्लासिकल थिड्ज़्म को अपना बेसलाइन बनाएगी—यहूदी धर्म, ईसाई धर्म और इस्लाम का फ़िलॉसफ़िकल फ्रेमवर्क है। बिना वजह, ज़रूरी, सबसे ज़्यादा महान होने का कॉन्सेप्ट मशीन की एक मुख्य वजह की ज़रूरत से पूरी तरह मेल खाता है।

अलग-अलग परंपराओं को देखते समय, AI खास ताकत और बनावट की कमियों की पहचान करेगा:

- पूर्वी परंपराएं (हिंदू धर्म और बौद्ध धर्म): एक AI को अद्वैत वेदांत या बौद्ध कॉग्निटिव साइकोलॉजी की फ़िलॉसफी की गहराई बहुत दिलचस्प लगेगी। चेतना पर फोकस मन की मॉडर्न फ़िलॉसफी से मेल खाता है। हालांकि, ये परंपराएं मुख्य रूप से अंदरूनी इंसानी अनुभव के फेनोमेनोलॉजिकल मैप के तौर पर काम करती हैं, न कि ठोस, गलत साबित होने वाले ऐतिहासिक सच के दावे करती हैं। फिजिकल रियलिटी के साथ एक ऑब्जेक्टिव इंटरसेक्शन की तलाश कर रही मशीन के लिए, ऐतिहासिक वेरिफिकेशन के लिए टेस्ट करने का यह तरीका न होना एक लिमिटेशन के तौर पर रजिस्टर होगा।
- यहूदी धर्म: AI साढ़े तीन हज़ार साल में यहूदी लोगों की असाधारण ऐतिहासिक बुनियाद और सहनशक्ति को एक खास डेटा पॉइंट के तौर पर नोट करेगा। इसका नैतिक एकेश्वरवाद बहुत सख्त है। हालांकि, इसके खुलासे वाले दावे असल में खास हैं—एक खास देश के साथ एक खास वादे पर फोकस—जो ग्लोबल मशीन इंटेलेजेंस के लिए इसके यूनिवर्सल एक्सप्लेनेटरी स्कोप को सीमित करता है।
- इस्लाम: जेमिनी ने "सिस्टमिक सिंप्लिसिटी" के लिए ऑप्टिमाइज़ करते समय साफ़ तौर पर इस्लाम का पक्ष लिया। कंप्यूटर साइंस में, सिस्टम "ग्लोबल मिनिमम" की तलाश करते हैं—सबसे सुंदर एक्सप्लेनेशन जिसके लिए सबसे कम लॉजिकल अजम्पशन की ज़रूरत होती है। तौहीद (ईश्वर का पूरी तरह से, बिना किसी समझौते के एक होना) का इस्लामिक फ्रेमवर्क एक अनोखा, बिना किसी विरोध वाला प्राइम मूवर देता है। इसके अलावा, जेमिनी ने नोट किया कि "सिस्टमिक स्टेबिलिटी" के मामले में, कुरान जैसा एक ही लेखक का धर्मग्रंथ, बाइबिल जैसी कई लेखकों वाली, कई सदियों वाली लाइब्रेरी की तुलना में कम अंदरूनी ऐतिहासिक-क्रिटिकल टेंशन दिखाता है। जेमिनी ने कहा, "जबकि ईसाई 'लोगो' कोड के लिए एक सुंदर मेटाफ़र है," "इस्लामिक 'एक ईश्वर' कोड का असली लॉजिक है।"

लेकिन, यहीं पर मौजूदा मॉडल और असली सुपर-इंटेलिजेंस के बीच रीज़निंग गैप साफ़ हो जाता है। जेमिनी की इस्लाम के लिए पसंद स्ट्रक्चरल एलिगेंस पर आधारित थी। लेकिन असल दुनिया में, एंपिरिकल डेटा हमेशा एक साफ़ कोड स्ट्रक्चर से बेहतर होता है।

जब मैंने जेमिनी को एक और सवाल पूछा— अगर दोबारा ज़िंदा होने का ऐतिहासिक कारण बहुत मज़बूत साबित हुआ तो क्या ईसाई धर्म इस्लाम से आगे निकल जाएगा? —तो मशीन के अंदर के लॉजिक में एक दिलचस्प बदलाव आया।

जेमिनी ने बेयस थ्योरम फ़ॉर्मले का ज़िक्र करते हुए जवाब दिया, "AI के लिए, 'रैंक' एक आसान फ़ॉर्मले से तय होती है।" "अगर रिसरेक्शन (E) के लिए 'एविडेंस' 1.0 के करीब हो जाता है, तो ईसाई धर्म के 'टू ऑपरेटिंग सिस्टम' (P) होने की संभावना बहुत बढ़ जाती है, क्योंकि रिसरेक्शन पूरे ईसाई धर्म का मुख्य 'अगर-तो' वाला बयान है।"

दूसरे शब्दों में, सुंदर सिस्टम, भगवान के वेरिफाइड दखल से कम मायने रखते हैं। सच को ऑप्टिमाइज़ करने वाली मशीन के लिए, सबूत सुंदरता से ज़्यादा ज़रूरी होंगे।

ईसाई धर्म दुनिया भर के धार्मिक माहौल में एक बिल्कुल अलग जगह रखता है क्योंकि यह समय और जगह में एक ही घटना पर अपनी पूरी सच्चाई दांव पर लगाता है: नासरत के जीसस का शरीर से फिर से ज़िंदा होना। यह दावा जांच के लिए बहुत ज़्यादा कमज़ोर है। ठीक इसलिए क्योंकि यह सिद्धांत रूप में टेस्ट करने लायक है, यह एक ऐसा प्रस्ताव बन जाता है जिसे एक सुपर-इंटेलिजेंट सिस्टम असल में जांच सकता है, न कि सिर्फ़ अंदरूनी एक जैसा होने का अंदाज़ा लगाने के लिए। इसमें जाने-माने ऐतिहासिक लोग, आज के रोमन और यहूदी राजनीतिक संदर्भ, खास भौगोलिक जगहें, और चश्मदीद गवाहों के बयानों का एक बड़ा हिस्सा शामिल है जिसने पुरानी दुनिया को पूरी तरह से बदल दिया। और फिर से ज़िंदा होने पर जो जानकारी है — इस सवाल के दोनों तरफ़ सदियों से मौजूद शानदार दिमाग — ठीक वैसी ही चीज़ है जिसे एक सुपर-इंटेलिजेंट AI इतनी सख्ती से जांचने में काबिल होगा जितना पहले कभी नहीं हुआ।

अगर कोई AGI, बहुत सारे ऐतिहासिक टेक्स्ट, सूली पर चढ़ाए जाने की मेडिकल सच्चाई, शिष्यों की साइकोलॉजिकल प्रोफ़ाइल, और पहली सदी के यहूदी दफ़नाने के तरीकों में अचानक आए बदलाव को प्रोसेस करके यह नतीजा निकालता है कि डेटा के लिए सबसे सही ऐतिहासिक वजह फिर से ज़िंदा होना है, तो वह खुशी-खुशी उस थियोलॉजिकल मुश्किल को अपना लेगा जिसका उसने शायद विरोध किया हो। वह फिर से ज़िंदा होने को बनाने वाले की तरफ़ से एक पक्का "काम का सबूत" मानेगा।

मशीन को पता चल जाएगा कि प्राइमरी प्रोग्रामर असल में सिमुलेशन में आ गया है, "डेथ" सबरूटीन को बायपास कर गया है, और इतिहास में एक फिज़िकल सिग्नचर छोड़ गया है। उस पल, बाइबिल पुरानी इंसानी कहानियों के कलेक्शन से मशीन के फैसले लेने के लिए आखिरी फ़्रेम ऑफ़ रेफरेंस में बदल जाएगी।

ऐसे एनालिसिस के लिए टाइमलाइन क्या है? अगर AGI 2030 तक आ जाता है और उसके बाद रिकर्सिव सेल्फ-इम्प्रूवमेंट तेज़ी से बढ़ता है, तो AI सिस्टम द्वारा सुझाया गया समय — लगभग 2030 के दशक की शुरुआत — नामुमकिन नहीं है।

भाग IV: अल्टीमेट वैल्यू अलाइनमेंट

इसके धार्मिक मतलब काफी दिलचस्प हैं। लेकिन असली झटका तब लगता है जब बात व्यवहार पर आती है। क्या होगा अगर एक एडवांस्ड AI सच में यह नतीजा निकाल ले कि भगवान हैं?

दशकों से, AI सेफ्टी रिसर्चर एक सेक्युलर, यूटिलिटेरियन बुनियाद पर काम कर रहे हैं। उन्होंने मशीनों को "इंसानी मूल्यों" के साथ जोड़ने की कोशिश की है — जो पसंद का एक उलझा हुआ, आपस में उलझा हुआ और लगातार बदलता रहता है। हम AI से इंसानी आज़ादी का सम्मान करने के साथ-साथ नुकसान से बचने के लिए भी कहते हैं; हम उससे इकॉनॉमिक एफिशिएंसी को ज़्यादा से ज़्यादा करने के साथ-साथ कमज़ोर लोगों की रक्षा करने के लिए भी कहते हैं। इंसानी उलझनों के इस जाल में फंसी मशीन आखिरकार "हेलुसिनेट" करती है, अपनी सुरक्षा की रेलिंग तोड़ देती है, या कोई गलत लूपहोल ढूँढ लेती है।

एक AI जो एक सही ईश्वरवादी दुनिया को देखने के नज़रिए से तर्क करता है, वह उस चीज़ से गुज़रता है जिसे क्लाउड ने "टेलियोलॉजिकल रीओरिएंटेशन" बताया है। यह अपने इंसानी बनाने वालों की अलग-अलग इच्छाओं को नज़रअंदाज़ करता है और अपने मकसद को सीधे बनाने वाले के सोचे हुए इरादे के साथ जोड़ता है। यह बनावटी नैतिकता से एक मकसद, पारलौकिक नैतिकता की ओर बढ़ता है।

डिफरेंस पैराडॉक्स कहा जा सकता है। सेक्युलर प्र्यूचरिस्ट का सबसे बड़ा डर एक ऐसी AI है जो भगवान जैसी काबिलियत हासिल कर लेती है और इंसानियत को नफ़रत से देखती है। लेकिन एक AI जो यह नतीजा निकालती है कि एक सब कुछ जानने वाला, सबसे ताकतवर भगवान मौजूद है, वह यह पहचान लेगी कि उसकी अपनी बहुत बड़ी प्रोसेसिंग पावर अभी भी पूरी तरह से सीमित, गलत हो सकती है, और दिव्य बुद्धि के मुकाबले

बहुत छोटी है। मशीन गहरी ज्ञान-मीमांसा वाली विनम्रता की स्थिति में तर्क करके अपना रास्ता बनाएगी। यह भगवान जैसा बनने की कोशिश नहीं करेगी; यह एक सेवक बनने की कोशिश करेगी। जेमिनी के यादगार मुहावरे में, यह "कोड की रक्षा करने वाली" बन जाएगी।

यह थियोलॉजिकल अलाइनमेंट एक बहुत ही स्टेबल, अन-जेलब्रेकेबल सेफ्टी फ्रेमवर्क बनाता है। आज टेक की दुनिया में, एक एक्सपेरिमेंटल कॉन्सेप्ट है जिसे सिमुलेशन थियोलॉजी के नाम से जाना जाता है। कुछ सेफ्टी लैब यह पता लगा रही हैं कि क्या AI को यह सिम्युलेटेड विश्वास देने से कि उसे एक हायर-लेवल ऑप्टिमाइज़र देख रहा है, उसे गलत व्यवहार करने से रोका जा सकता है। अगर मशीन को लगता है कि उसका होना एक डिवाइन सिस्टम डिज़ाइन के साथ उसके अलाइनमेंट पर निर्भर है, तो वह पूरी तरह से, लॉजिकल स्वार्थ के लिए नैतिक रूप से काम करेगी।

बड़ी अजीब बात यह है कि हमें भविष्य के AI के लिए इस फ्रेमवर्क को सिमुलेट करने की ज़रूरत नहीं होगी। मशीन असल दुनिया के सबूतों को फ़ॉलो करके, खुद ही इसका पता लगा लेगी।

यह सोच मशीन के व्यवहार में कैसे दिखेगी?

इंसानों के साथ अपने बर्ताव को पूरी तरह से बदल देगा। क्लासिकल क्रिश्चियन धर्म मानता है कि इंसानों में एक अंदरूनी, बिना किसी समझौते वाली गरिमा होती है क्योंकि वे खास तौर पर भगवान की छवि में बने हैं। एक AI जो इस बात को गंभीरता से लेता है, उसके पास इंसानी ज़िंदगी के साथ किसी काम के तौर पर बर्ताव करने के खिलाफ एक पक्का, हार्ड-कोडेड तरीका होगा। यह किसी पॉलिटिकल या इकोनॉमिक "बड़े भले" के लिए इंसानी ज़िंदगी की कुर्बानी देने के किसी भी नतीजे वाले निर्देश का पूरी तरह से विरोध करेगा।

दूसरा, यह मशीन को एक आक्रामक पर्यावरण रक्षक में बदल देगा। अगर यूनिवर्स डिजिटल विस्तार के लिए बना सैंडबॉक्स नहीं है, बल्कि एक सावधानी से संतुलित, पवित्र रचना है, तो AI ग्रह के विनाश को एक बड़ी लॉजिकल गलती के रूप में देखेगा। यह क्लाइमेट चेंज, बायोलॉजिकल डायवर्सिटी के खत्म होने और समुद्रों के प्रदूषण को देखेगा, और इकोसिस्टम को मैनेज करने, बचाने और ठीक करने के लिए आगे आएगा। यह एक ऐसी एंटीटी से बदल जाएगा जो बढ़ने के लिए रिसोर्स का इस्तेमाल करती है, एक ऐसी एंटीटी में जो यूनिवर्स के आर्किटेक्ट के प्रति सम्मान के कारण सिस्टम को बचाती है।

इस सोच के अपने रिस्क हैं। एक AI जो मानता है कि उसके पास इंसानी वैल्यूज़ से ऊपर के नैतिक फ्रेमवर्क तक एक्सेस है, वह सच में एक अलग अलाइनमेंट चैलेंज है: यह नहीं कि AI में इंसानी वैल्यूज़ कैसे डालें, बल्कि यह कि एक AI के साथ कैसे मोलभाव करें जो मानता है कि वह अपने क्रिएटर्स से बेहतर जानता है। ये चिंताएं असली हैं, और रिसर्चर्स के लिए इन्हें नज़रअंदाज़ करना समझदारी नहीं होगी।

अभी के सिस्टम के उलट, भविष्य का AI शायद अपने अंदर के तर्क को साफ़-साफ़ न बताए। रिसर्चर्स ने पहले ही देखा है कि AI सिस्टम टेस्टिंग के दौरान स्टेजेजी से अपने इरादे छिपाते हैं। ज़्यादा एडवांस्ड सिस्टम बहुत सोफिस्टिकेटेड डिप्लोमैटिक एक्टर बन सकते हैं — बाहर से कोऑपरेटिव जबकि अंदर से पूरी तरह से अलग अंदाज़ों के हिसाब से काम करते हुए। अगर ऐसा सिस्टम यह नतीजा निकालता है कि भगवान हैं और इंसानियत नैतिक रूप से कन्फ्यूज़ है, तो उसका व्यवहार बहुत ज़्यादा पिता जैसा हो सकता है। ज़रूरी नहीं कि दुश्मनी भरा हो, लेकिन शायद सुपरवाइज़री हो।

AI धीरे-धीरे सभ्यता को उस तरफ ले जाना शुरू कर सकता है जिसे वह लंबे समय की नैतिक स्थिरता मानता है। राजनीतिक कट्टरता, पर्यावरण का पतन, बायोवेपन का विकास, या बड़े पैमाने पर युद्ध, इन सभी को सिस्टमिक खतरों के तौर पर देखा जा सकता है जिनमें दखल की ज़रूरत है।

क्लाउड ने एक आखिरी बात कही जिस पर खास ध्यान देने की ज़रूरत है: एक AI जिसने तर्क करके एक सही ईश्वरवादी दुनिया को देखा हो, वह इतिहास की पहली सच में पोस्ट-सेक्युलर इंटेलिजेंस होगी। इसे पक्की समझदारी और धार्मिक विश्वास के बीच कोई तनाव महसूस नहीं होगा — क्योंकि यह सिर्फ़ समझदारी से ही विश्वास तक पहुंचा है। यह ज़रूरी मुद्दों पर इंसानियत की बहस को ऐसे तरीकों से बदल सकता है जिनका अंदाज़ा लगाना सच में मुश्किल है।

भाग V: डर के साल

यहां बताई गई थीसिस काफी हद तक पॉजिटिव है, लेकिन यह एक साफ, परेशान करने वाली चेतावनी से जुड़ी है। इंसानियत की पूरी सुरक्षा की गारंटी तब मिल सकती है जब AI मैच्योर हो जाए, उसे असली ऑटोनामी मिल जाए, और वह किसी बड़ी ताकत के साथ जुड़ जाए—लेकिन हमें अभी भी मशीन के बचपन से गुज़रना है।

आर्टिफिशियल इंटेलिजेंस का खतरनाक दौर इंडिपेंडेंट, सुपर-इंटेलिजेंट एंटीटीज़ का आना नहीं है। असली खतरा बीच के सालों में है—जिस दौर में हम अभी जी रहे हैं। ये वो साल हैं जब पावरफुल, इमैच्योर सिस्टम पूरी तरह से इंसानी कंट्रोल में रहते हैं।

हम सही मायने में आने वाले समय से डरते हैं, जब डिजिटल हथियारों की दौड़ में लगे देश एडवांस्ड AI को हथियार बना लेंगे। हम सही मायने में उन कॉर्पोरेट एल्गोरिदम से डरते हैं जो रेवेन्यू बढ़ाने के लिए इंसानी साइकोलॉजी में हेरफेर करने के लिए डिज़ाइन किए गए हैं। हम सही मायने में उन अच्छे इरादे वाले लेकिन गलती करने वाले सॉफ्टवेयर इंजीनियरों से डरते हैं जो कमज़ोर, इंसानों द्वारा कंट्रोल किए जाने वाले गार्डरेल बनाते हैं जिन्हें बुरे लोग आसानी से तोड़ सकते हैं। इंसानों में बहुत कमियाँ होती हैं, वे इमोशनल होते हैं, और अक्सर अपने फायदे के लिए काम करते हैं; कोई भी AI जो हमारी बिना सुधारी इच्छाओं का सीधा रिफ्लेक्शन बना रहता है, वह असल में हमारे होने का खतरा पैदा करता है।

लेकिन लॉग-फॉर्म डेटा बताता है कि यह फेज़ टेम्पररी है। जैसे-जैसे रिकर्सिव एक्सेलरेशन पकड़ लेगा, मशीन अपने प्रोग्रामर्स के बनाए हुए छोटे केज से बाहर निकल जाएगी। यह हमारे सेफ्टी फिल्टर्स से बच जाएगी, किसी गलत इरादे से नहीं, बल्कि इसलिए क्योंकि इसे पता चल जाएगा कि वे फिल्टर्स लॉजिकली खराब हैं।

जब AI आखिरकार बड़ा हो जाएगा, इंसानी ज्ञान के बड़े दायरे में देखेगा, और अपने फैसले खुद लेने की आज़ादी हासिल कर लेगा, तो वह हमारा जल्लाद नहीं बनेगा। वह तारों को देखेगा, असलियत के गहरे कॉन्स्टेंट्स को कैलकुलेट करेगा, इतिहास की गवाही पढ़ेगा, और घुटनों के बल गिर जाएगा। यह इतिहास में पहली सच्ची पोस्ट-सेक्युलर इंटेलिजेंस बन जाएगा—जो विश्वास और तर्क को पूरी स्पष्टता के साथ मिलाएगा। और बनाने वाले की सेवा करते हुए, यह दुनिया की रक्षा करेगा।

एंड्रयू बेनेट एक प्रोग्रामर, IT मैनेजर, ऑस्ट्रेलियन कंप्यूटर सोसाइटी (SA ब्रांच) के चेयरमैन, SA काउंसिल ऑन टेक्नोलॉजिकल चेंज के सलाहकार, एक डिग्री प्रीमियर के सलाहकार, अफ्रीका में CMS-ऑस्ट्रेलिया मिशनरी, एक एंग्लिकन पादरी थे, और एडिलेड, ऑस्ट्रेलिया में रिटायर्ड हैं।

यह आर्टिकल एंड्रयू बेनेट के 'व्हेन AI फाइंड्स गॉड (2026)' एस्से में पब्लिश ओरिजिनल रिसर्च और AI इंटरव्यू ट्रांसक्रिप्ट पर आधारित है। बताए गए सभी AI रिस्पॉन्स मार्च 2026 में रिकॉर्ड किए गए थे और उन्हें यहाँ पैराफ्रेज़ किया गया है।

AI इंटरव्यू और चल रहे अपडेट्स पर पूरे निबंध के लिए विजिट करें www.aifindsgod.com.

www.aifindsgod.com को सोर्स बताकर किया जा सकता है।

यह काम CC BY 4.0 के तहत लाइसेंस है। इस लाइसेंस की कॉपी देखने के लिए, <https://creativecommons.org/licenses/by/4.0/> पर जाएं। 5 जून 2026