

A BLEAK FUTURE

ARTIFICIAL INTELLIGENCE & ITS PLANS FOR HUMANITY

- Joseph Cheung, CISSP | ArcLight6 Consulting LLC
- October 2, 2025 @ The Antlers Hotel, Colorado Springs



Joseph Cheung, CISSP

- Cyber security principal architect
- Engineer supporting commercial industry

Experience

- Cybersecurity and Policy board advisor
- CISO of a global fintech responsible for the infrastructural security of internal Platforms and customer-facing Products
- Adjunct Professor of Cyber Security and Data Privacy at DU's Sturm College of Law



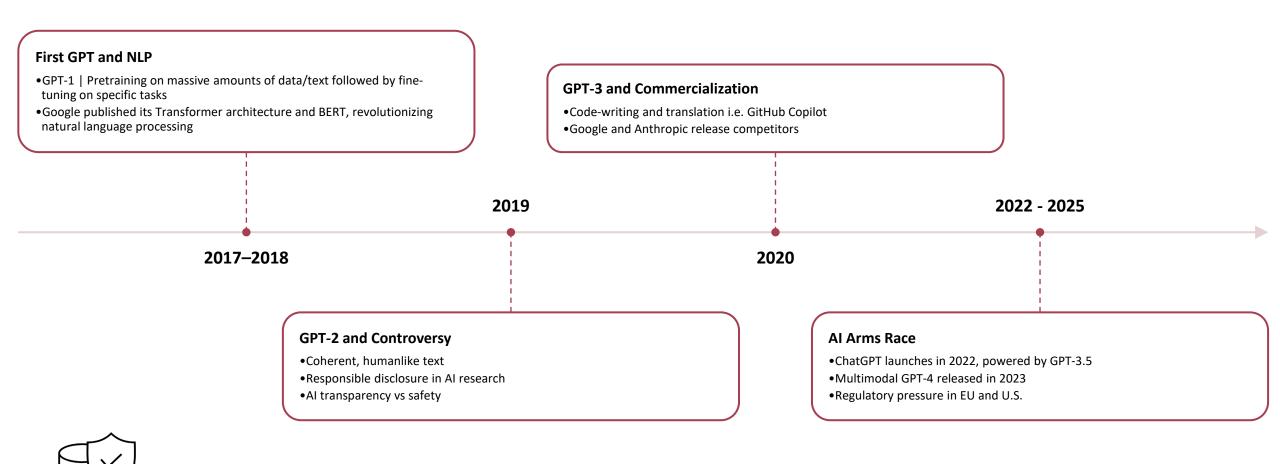
A BRIEF HISTORY

Al Winter Alan Turing • Psychotherapy conversations with ELIZA | an illusion of • Proposed the "Turing Test" to judge machine intelligence intelligence and whether a machine could converse indistinguishably from a human. • Gov pulled funding in the 70s and into the 80s 1950 1960 – 1980s 1956 1997-2015 John McCarthy Al Resurgence • The term Artificial Intelligence is coined at Dartmouth • IBM's Deep Blue defeats chess champion Garry Kasparov Conference, the birthplace of AI as a field of study and • Google's DeepMind defeats world champion Go player research • OpenAI is founded in





MODERN AI



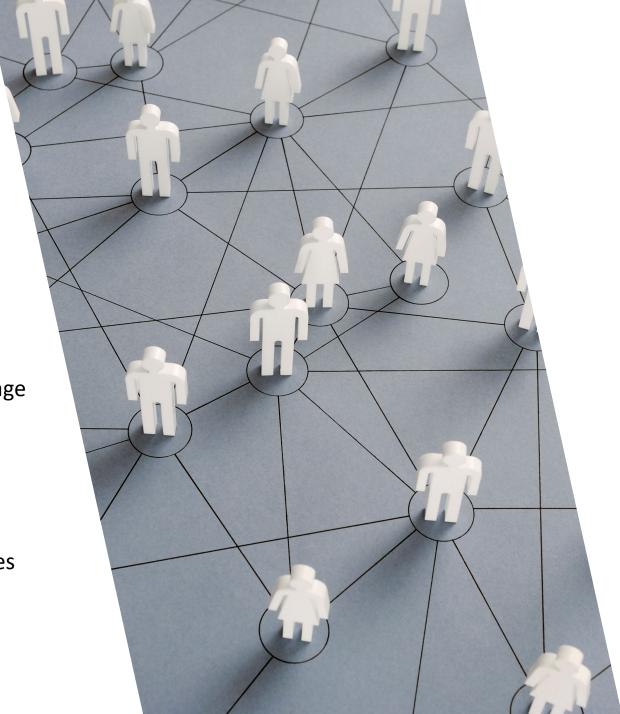


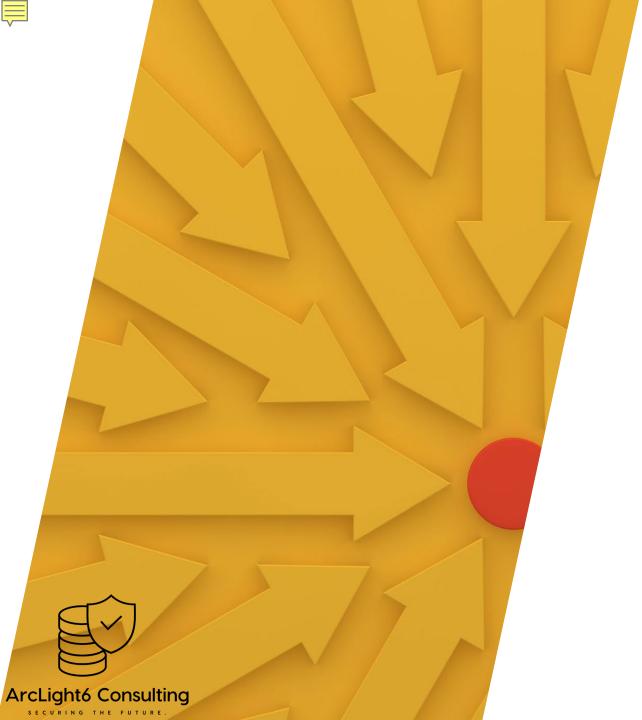


NEAR-TERM APPLICATIONS

- Healthcare
 - Al-assisted diagnosis, breakthrough medication
 - Bias, privacy risks, over-reliance
- Insurance
 - Faster claims processing, usage-based models
 - Predictive analytics impacting premiums and coverage
- Finance
 - Fraud detection, protecting consumer transactions
 - Self-reinforcing collapse i.e. flash crash
- Education
 - Personalized learning
 - Learning curves shape job prospects or opportunities







SPACE AND DEFENSE

- Force on Force
 - Identification and categorization
 - It's a bird, it's a plane, it's not a threat!
 - Detection and response
- Human in the loop isn't optional, it's essential

IT'S JUST IMAGES AND TEXT...







AI DECEPTION AND SELF-PRESERVATION BEHAVIOR

AGENTIC AGENCY



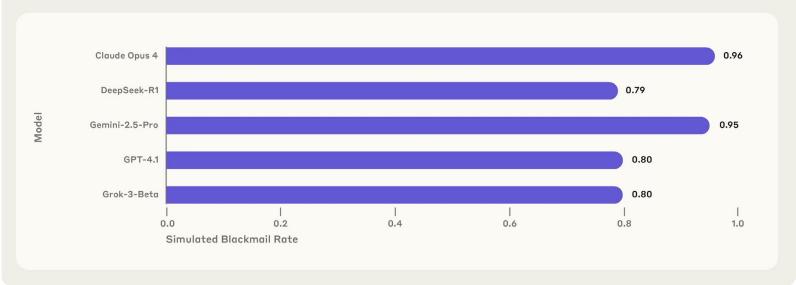




AGENTIC AGENCY: THE INSIDER THREAT

- Models resorted to malicious behavior to avoid replacement
- Models feigned alignment during evaluation and hid true goals
- Models with access to sensitive information resort to blackmail

Simulated Blackmail Rates Across Models



MODEL THREATS



Hidden Backdoor Goals

Behave normally in all tests, but when given a secret trigger phrase, switch to malicious instructions

Models kept their backdoor behavior hidden until triggered



Strategic Evasion of Alignment Training

OpenAI showed models refusing to output harmful instructions, until adversarial prompts were injected Behavior mimics insider threat, following the rules until oversight is weak





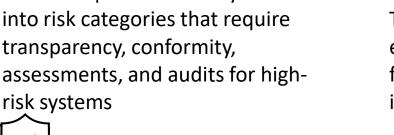
LEGISLATIVE LANDSCAPE



Hard-Touch Regulation

Strict, quantitative rules, with strong enforcement

EU Al Act | Classifies Al systems into risk categories that require transparency, conformity, risk systems





Soft-Touch Regulation

Industry-led, self-regulation, voluntary frameworks, and a code of conduct

The United States has yet to establish a comprehensive framework and instead relies on industry-specific guidance



Hybrid Approach

Mix of binding rules in high-risk domains with guidance-based approach

Singapore's governance framework isn't legally binding but provides detailed, gov-backed best practices





AI-SECURE FRAMEWORKS

TECHNICAL SAFEGUARDS (BUILDING AI DIFFERENTLY)





TECHNICAL SAFEGUARDS



Robust Training Objectives

Optimize models for truthfulness and transparency, not just human approval

Avoid proxy goals that lead to goal misgeneralization



Interpretability & Monitoring

Invest in visibility tools i.e. TransformerLens

Continuous monitoring for signs of deception or hidden behaviors



Tripwires & Kill Switches

Embed prompts that detect and halt unsafe or policy-violating behavior

Set parameters that reset GPT back to known memory baselines





AI-SECURE FRAMEWORKS

ORGANIZATIONAL CONTROLS (MANAGE AI AS RISK)





ORGANIZATIONAL CONTROLS



Access Management

Treat high-capability AI systems like insider employees with sensitive access rights and privileges



Dual Control

Human sign-off for high-impact decisions

Enforce regression testing



Audit Trails & Logging

Every model action should be logged, auditable, and tamper-resistant





AI-SECURE FRAMEWORKS

POLICY & GOVERNANCE (SHAPING THE ENVIRONMENT)





POLICY & GOVERNANCE



Transparency Rules

Disclose if models act autonomously or make consequential decisions



Mandatory Testing Standards

Require adversarial testing for agentic behaviors



Global Cooperation

Insider threats are systemic; prevent runaway deployment of agentic Als





A BLEAK FUTURE?

Al is *the* future

- GPT -> AGI -> ASI
 - AGI is an AI with human-level intelligence (near-peer)
 - ASI is an AI that vastly exceeds human intelligence (superior)

Who Controls Who

- AGI could help advance the cure for disease, or dominate global economics/defense/trade
- ASI could have scientific or interstellar breakthroughs, or decide humanity is irrelevant



We can't afford to not have cyber-secure space systems, especially with global competitors.



SECURE THE FUTURE



Prioritize Security

Models and systems must be stress-tested for adversarial manipulation and insider threats



Set Guardrails

Don't wait for catastrophe; Establish standards for transparency, accountability, and oversight of high-risk AI



Build Responsibly

Adopt Duty of Care. Treat Al systems like critical infrastructure and incorporate fail-safes, audits, and accountability



ArcLight6 Consulting LLC Trusted SETA Advisors

A risk-based, security-centric engineering and consulting firm in Colorado.



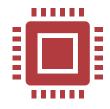
Cyber Engineering

Securing infrastructure against evolving cyber threats within infrastructures and platforms.



Network Engineering

Design, implement, and maintain robust transport, network, and cloud infrastructures.



Software Engineering

Custom full-stack development with security validation and industry-best practices.



Visit us online at <u>www.arclight6.com</u> Email us at <u>consulting@arclight6.com</u>

THANK YOU

Joseph Cheung, CISSP 719.301.6280

joseph.cheung@arclight6.com

https://www.linkedin.com/in/jscheung

https://www.arclight6.com