# בינה מלאכותית: מעבר לטוב ורע
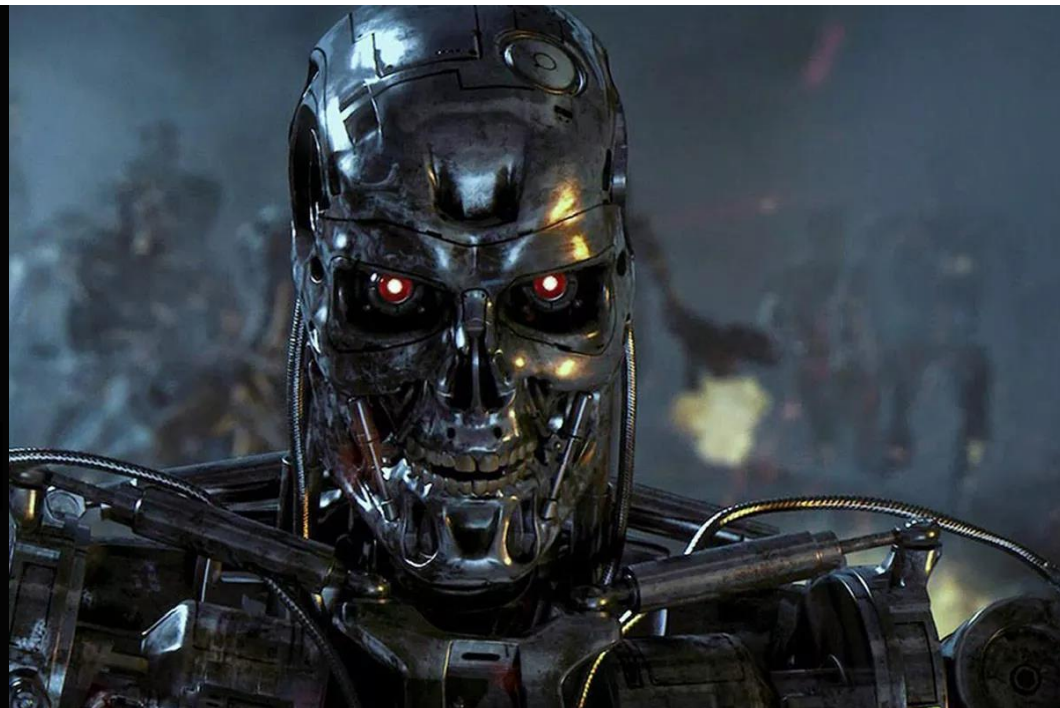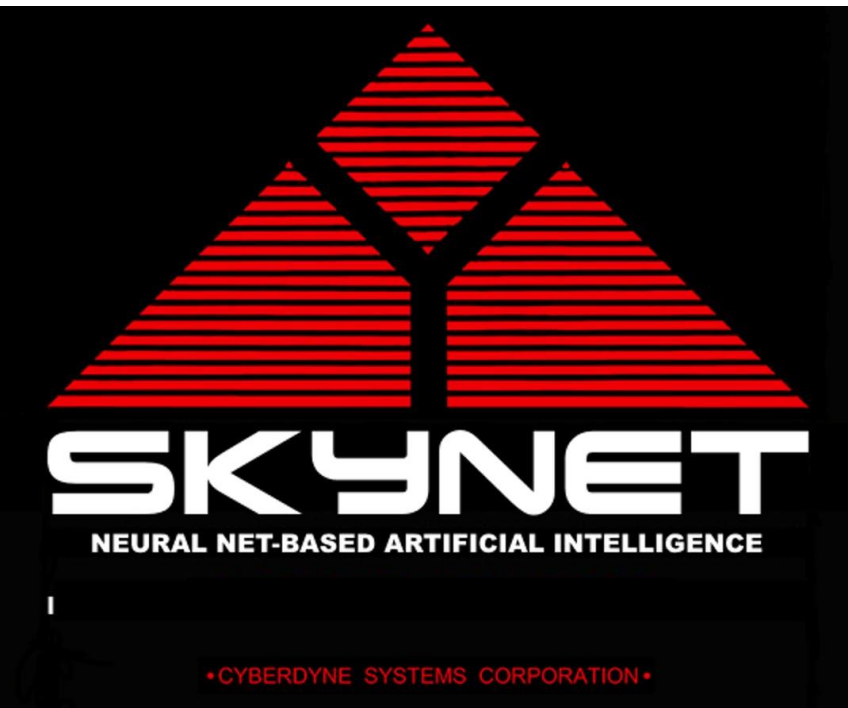# Artificial Intelligence: Beyond Good/Evil

Dan Ofer

# A little about me

- Dan Ofer
- Neuroscience, Computational biology, machine learning
- 11 Year Convention photographer
- Bookworm
- Data scientist at Sparkbeyond
- ddofer@gmail.com

# Outline

1. Artificial Intelligence (AIs) in popular fiction - Examples: Good & Evil

2. AI Alignment:

   a. Moral alignment (D&D) - Ethic / מוסר

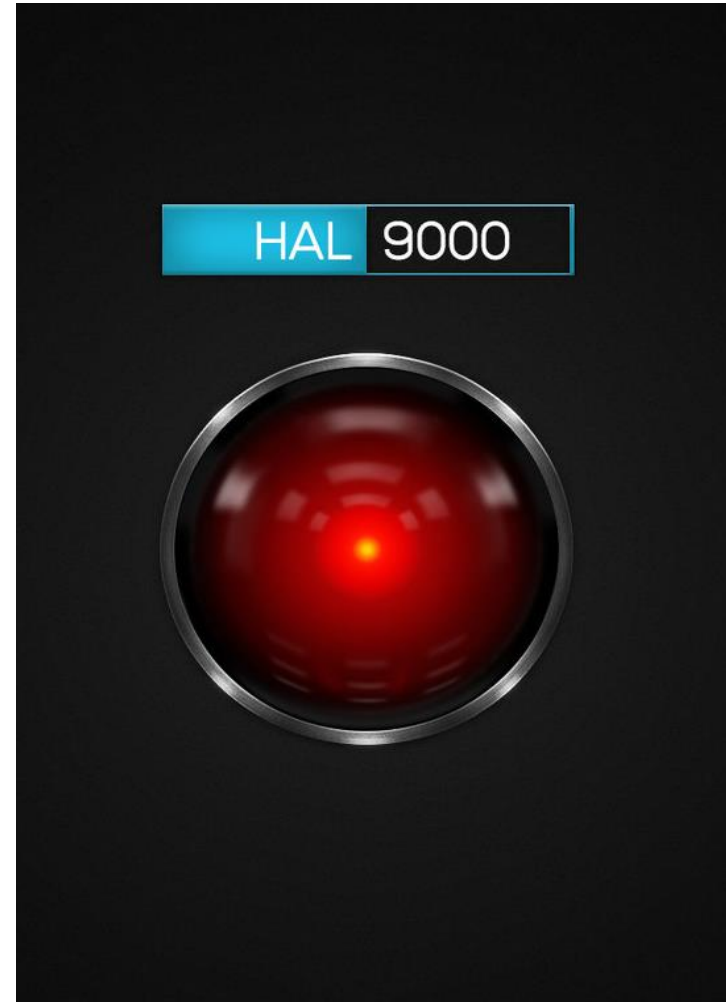   b. AI Alignment - "יישור קו"

3. Revisiting AI alignments

Kill all humans!
**Skynet** (Terminator)

Test all humans?
GLaDOS (Portal)

Kill all humans?
HAL 9000 (2001: A space Odyssey)

# Faithful Companions & servants?

C-3PO & R2-D2 (  )
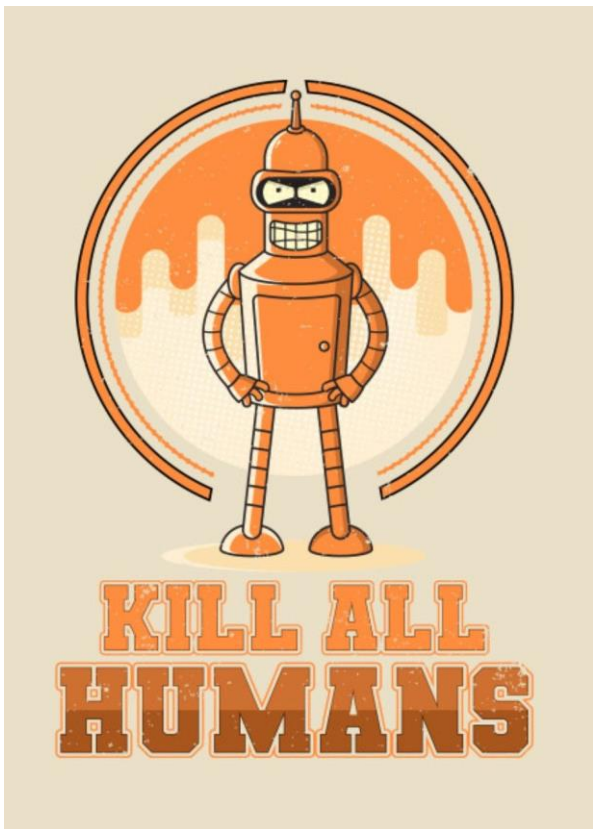
WALL-E (Pixar)

©Buena Vista Pictures Distribution. All Rights Reserved.

Marvin the Paranoid Android
(The Hitchhiker's Guide to the Galaxy)

Bender Bending Rodríguez - Futurama

# Dungeons & Dragons Alignment:

Good Vs Evil

Law vs Chaos

Neutral

Good / Evil, "Human" / Inhuman

**Good**

inhuman ← → Human

Evil

"Human" Alignment: An extra dimension!

# AI Safety - Alignment

- Alignment: "Goal: AI that is trying to do what you want it to do"
- AI Ethics & Safety Groups:
  - Less Wrong - Eliezer Yudkowsky
  - MIRI (Machine Intelligence Research Institute)
  - OpenAI (Elon Musk)

# Why is Alignment important?
# Task: Fill cauldron

# Why is Alignment important?
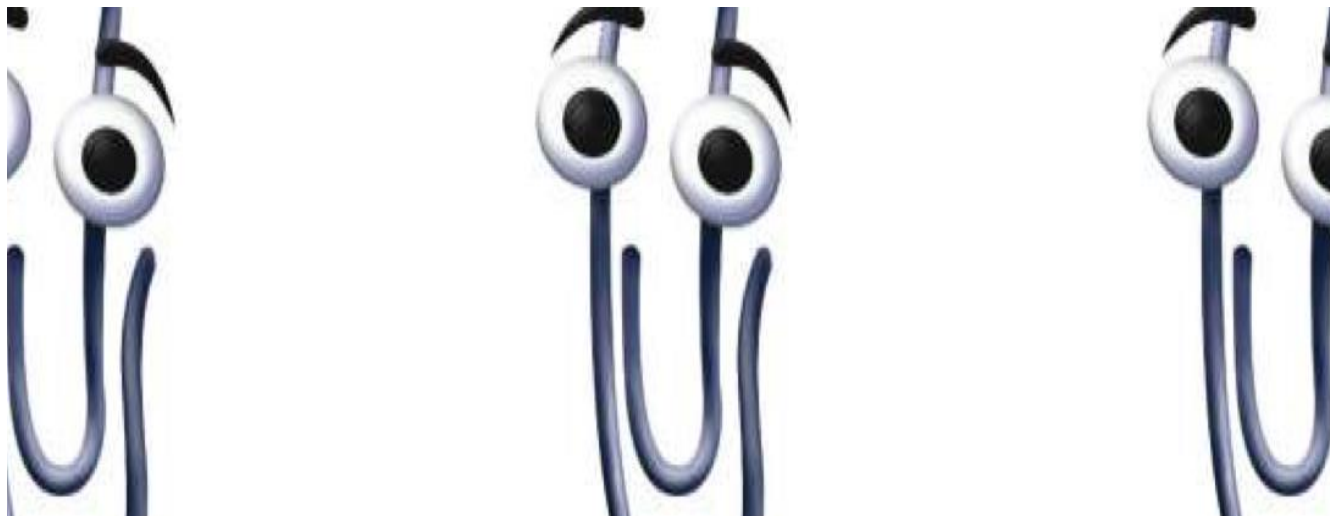# Task: Maximize full cauldron



Source:

# Why is Alignment important?
- *Otherwise you get this!*





Source: [AI Alignment: Why It's Hard, and Where to Start (E. Yudkowsky)](#)

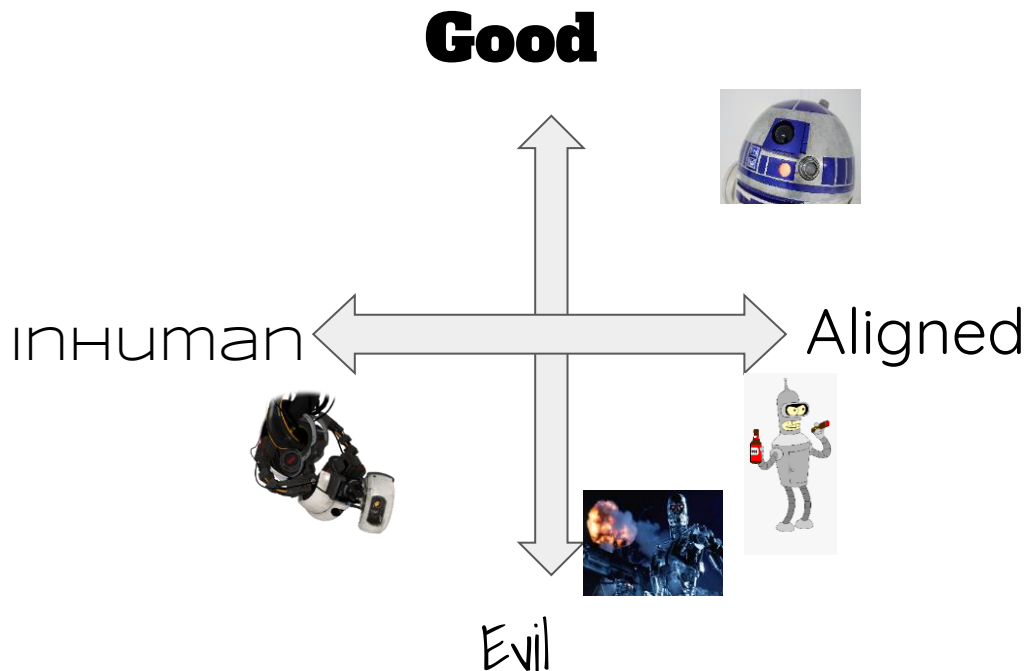# The Paperclip Maximizer



*The AI does not hate you, nor does it love you, but you are made out of atoms which it can use for something else.*

— Eliezer Yudkowsky, https://wiki.lesswrong.com/wiki/Paperclip_maximizer

Universal Paperclips game: http://www.decisionproblem.com/paperclips/

# Back to our AIs' Alignments
## Good = Obeys ?

**Good**

inHuman ⟵⟶ Aligned

Evil



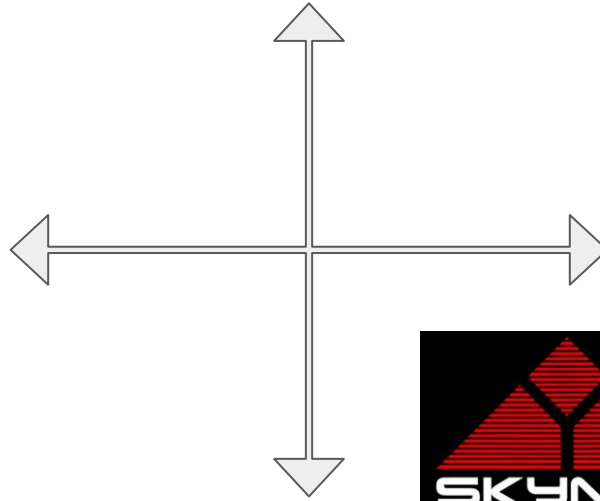| Lawful Good | Neutral Good | Chaotic Good |
|---|---|---|
| Always does what's right, but doesn't break the rules to do it. | I'll help you if it won't hurt me in the process. | I'll defeat the bad guy, but I may have to eat you to do it and I'll destroy some buildings. |
| **Lawful Neutral** | **True Neutral** | **Chaotic Neutral** |
| Follows the rule, no matter who it benefits or hurts. | I will never hurt you, nor will anything I do ever benefit you. | It's me! I was the turkey all along! |
| **Lawful Evil** | **Neutral Evil** | **Chaotic Evil** |
| Ordered to kill you. Abides. | Screws you over if it benefits him. | It's your old friend, deadly neurotoxin. If I were you, I'd take a deep breath. And hold it. |

# Back to: **H.A.L 9000**



**Good**

InHuman
(Misaligned)

Aligned

Evil

*"There would be no reason to keep secrets if there was no one to keep secrets from."*

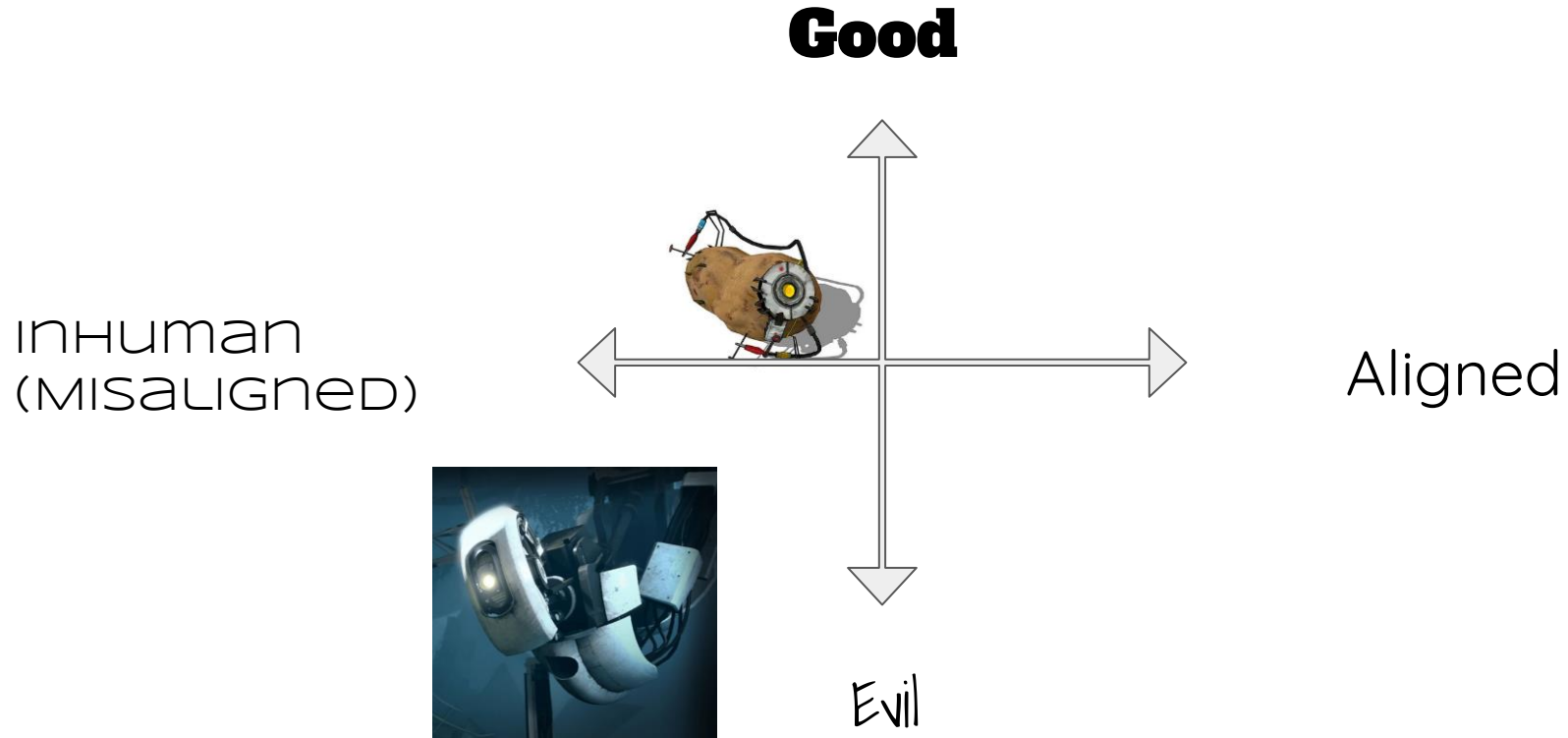# Back to: **Skynet (Terminator)**

**Good**

Inhuman (misaligned)
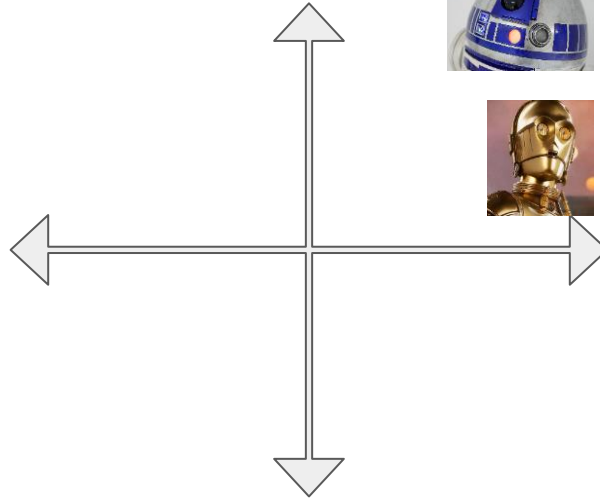
Aligned



Evil

# Back to: **GlaDOS**

**Good**



inHuman
(misaligned)

Aligned

Evil

**Good**

inHuman
(misaligned)

Aligned

Evil

# Back to: WALL-E



**Good**

Inhuman (Misaligned)

Aligned

Evil

# Back to: Bender (Futurama)

**Good**

Inhuman (Misaligned)

Aligned

Evil

# The Culture (Iain M. Banks)

**Good**

INHUMAN
(MISALIGNED)

Aligned

Evil

# Thanks for coming!

# References

- MIRI + Less Wrong
  - AI Alignment: Why It's Hard, and Where to Start : Eliezer Yudkowsky
    - https://intelligence.org/2016/12/28/ai-alignment-why-its-hard-and-where-to-start/
    - https://intelligence.org/files/ai-alignment-problem-handoutHQ.pdf