# System 1 and System 2 Thinking Cycles

Shahar Avin
sa478@cam.ac.uk

# The importance of AI strategy
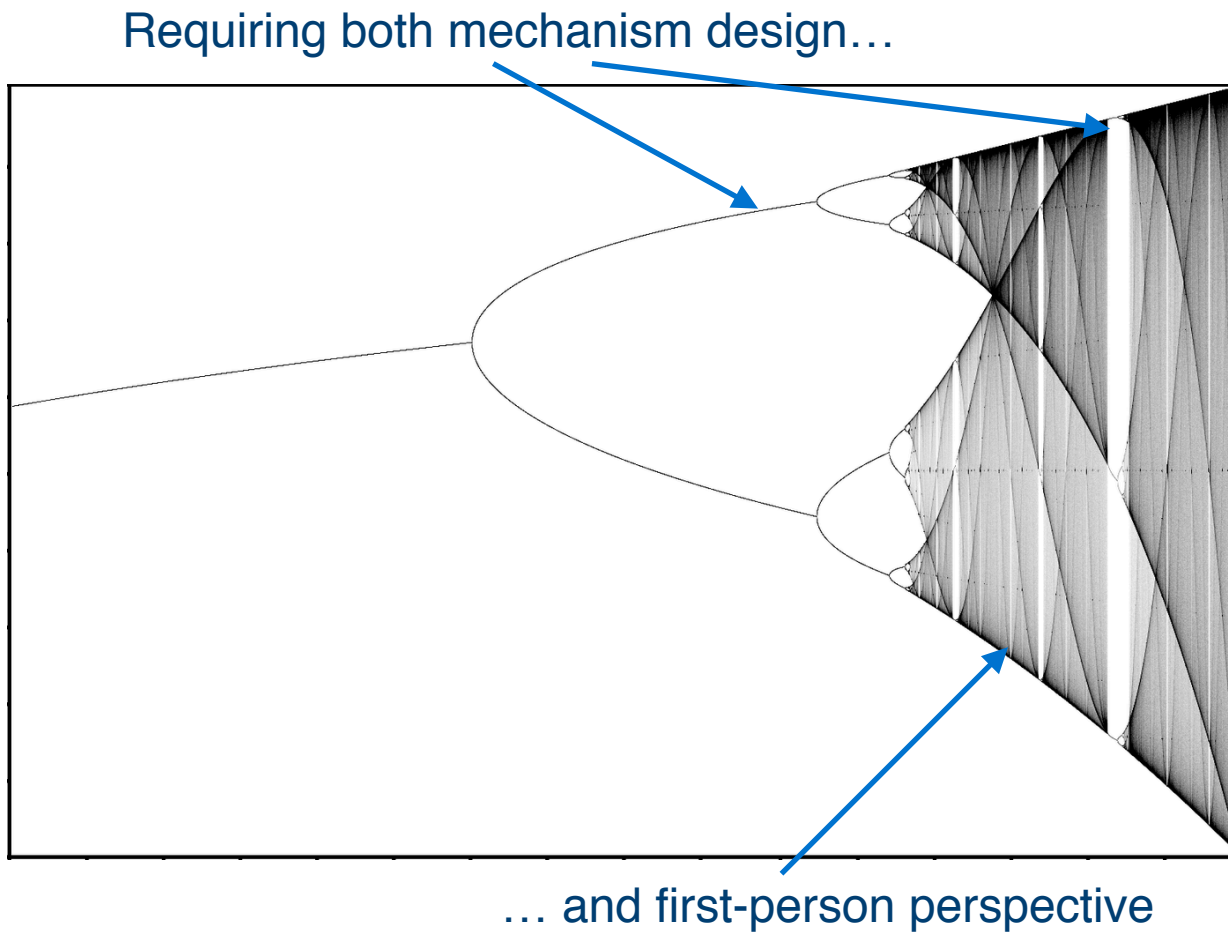
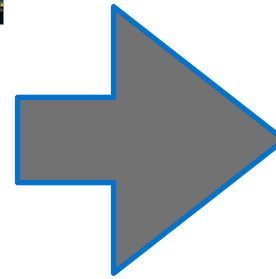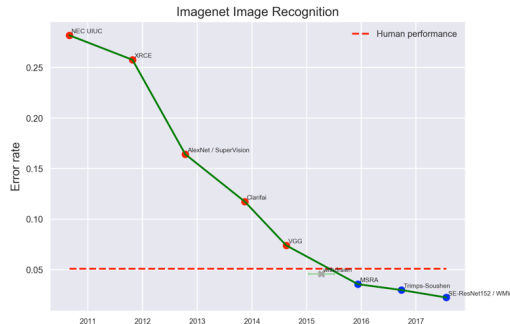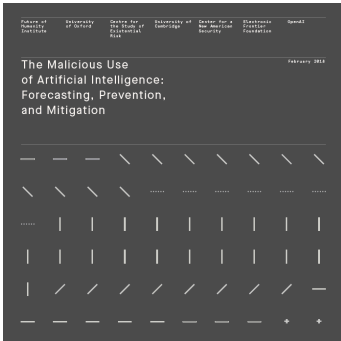# But history is messy

Requiring both mechanism design…



… and first-person perspective

UNIVERSITY OF
CAMBRIDGE

# Mechanism design

# First-person perspective

# Same system(s)…

# …different perspectives

- incentives

- models

- generalisations

- attractors

- stability

- identities

- context

- world view

- salience

- narrative

# Different, complementary outcomes

# CSER AI Strategy scenario role-plays

- ~15 so far

- 3-60 participants

- AI researchers, policy researchers, effective altruists, general public

- Play as USG, PRC, EU, Tech company, Defence contractor, NGO
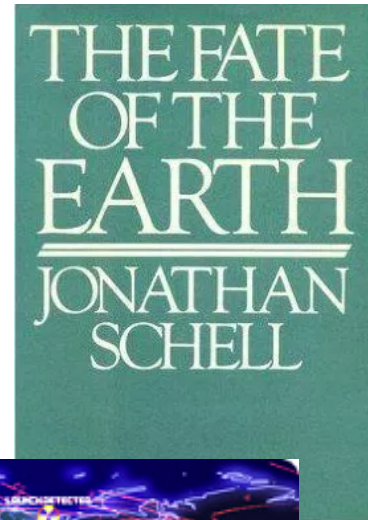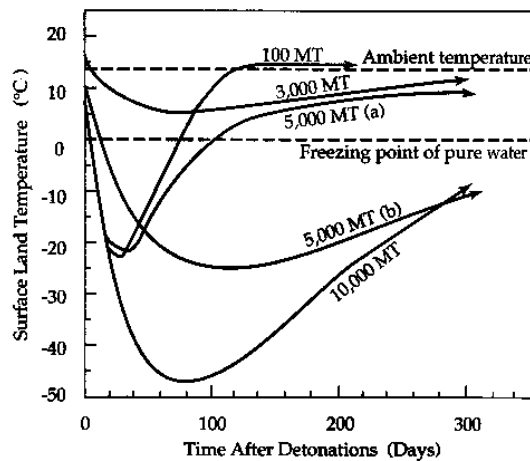
- Track talent, compute, technical breakthroughs, data, money

- Public and secret actions, negotiations, time pressure

- Fun, if depressing at times

- Get in touch if you're interested

# Earth 2045 scenario simulation

- Resolve a dramatic event

- Using a first-person perspective of a made-up character

- In collaboration with others

- While exploring an imaginary world

- Facilitated by an all-knowing story master

UNIVERSITY OF
CAMBRIDGE

So You Want to Play Dungeons & Dragons

There is no rule saying changing the world can't be (sometimes) fun

UNIVERSITY OF CAMBRIDGE