# Tutorial Plan

Introduction to the alignment problem  (Dylan) [Gillian]   **1:35 - 1:50**

The alignment problem in human societies (Gillian)   **1:50 - 2:45**

BREAK   **2:45 - 3:00**

Institutions and social dilemmas   (Joel)   **3:00 - 3:35**

Institutions in multiagent settings   (Dylan)[Rakshit]   **3:35 - 3:55**

Closing   **3:55 - 4:00**

# Asking Questions

- Chat feature on the [NeurIPS website](NeurIPS website)

- [Whova](Whova) App

- Tutorial website at:
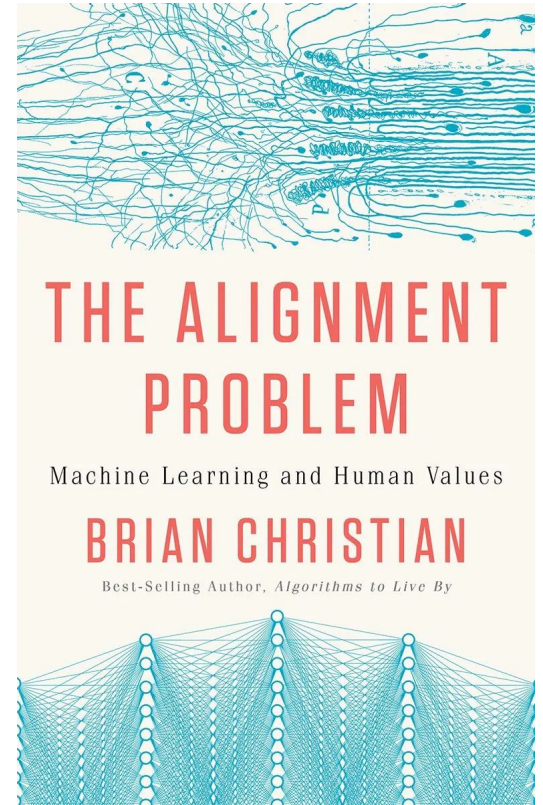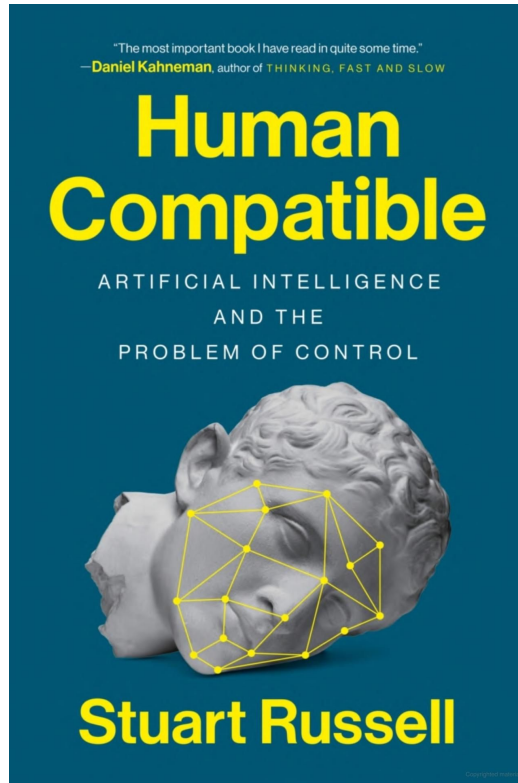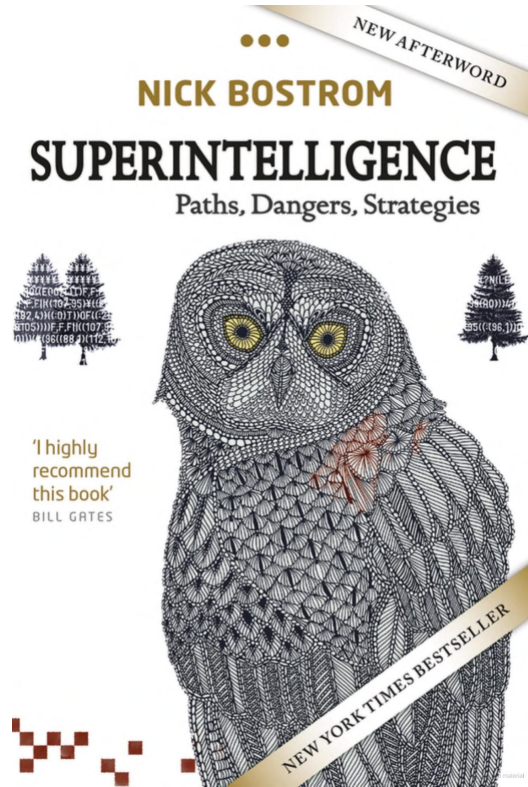    https://alignment-tutorial-2024.github.io/

# What is alignment?

If we use, to achieve our purposes, a mechanical agency with whose operation we cannot interfere effectively ... we had better be quite sure that the purpose put into the machine is the purpose which we really desire.

Wiener, "Some Moral and Technical Consequences of Automation" *Science* 1960

NICK BOSTROM

SUPERINTELLIGENCE

Paths, Dangers, Strategies

NEW AFTERWORD

'I highly recommend this book'
BILL GATES

NEW YORK TIMES BESTSELLER

---

"The most important book I have read in quite some time."
—Daniel Kahneman, author of THINKING, FAST AND SLOW

Human Compatible

ARTIFICIAL INTELLIGENCE AND THE PROBLEM OF CONTROL

Stuart Russell

---

THE ALIGNMENT PROBLEM

Machine Learning and Human Values

BRIAN CHRISTIAN

Best-Selling Author, Algorithms to Live By

Whose values? Intentions? Preferences?

Determined how and by whom?

Sorenson et. al.
Pluralistic Alignment
(2024)

Our goal today: explore how the study of the alignment problem in human societies can inform AI alignment

**How can we live together?**

# AI 101: How robots make decisions

$$\max_{\vec{a}} \sum_t r(s_t)$$

States

Actions

Grasp

Move

r=0

r=1

r=0

# Markov Decision Process

State: $\vec{x}^{\mathbf{H}}, \vec{x}^{\mathbf{R}}, \vec{x}^v$
Actions: $(\Delta_1, \ldots, \Delta_7)$
Reward: $\sum_o w_o ||\vec{x}^{\mathbf{R}} - \vec{x}^o||^2$

Head

Torso

Vase

Table

# How we pretend robots make decisions:

$$\max_{\vec{a}} \sum_t r(s_t)$$

How robots actually make decisions:

$$\max_{\vec{a}} \sum_t r(s_t)$$

# Objectives often generalize in unintended ways

How robots actually make decisions:

$$\max_{\vec{a}} \sum_t r(s_t)$$

$\tilde{r}$

How robots actually *actually* make decisions:

# REWARD (MIS)DESIGN FOR AUTONOMOUS DRIVING

W. Bradley Knox[*,1,2], Alessandro Allievi[1,2], Holger Banzhaf[3], Felix Schmitt[4], Peter Stone[2,5]

[1]Robert Bosch LLC
[2]The University of Texas at Austin
[3]Robert Bosch GmbH
[4]Bosch Center for Artificial Intelligence
[5]Sony AI

## ABSTRACT

This article considers the problem of diagnosing certain common errors in reward design. Its insights are also applicable to the design of cost functions and performance metrics more generally. To diagnose common errors, we develop 8 simple sanity checks for identifying flaws in reward functions. These sanity checks are applied to reward functions from past work on reinforcement learning (RL) for autonomous driving (AD), revealing near-universal flaws in reward design for AD that might also exist pervasively across reward design for other tasks. Lastly, we explore promising directions that may aid the design of reward functions for AD in subsequent research, following a process of inquiry that can be adapted to other domains.

# Sanity check: indifference points for crash risk



$$r(\tau_{idle}) = p \cdot r(\tau_{crash}) + (1 - p) \cdot r(\tau_{dest})$$

# Indifference Points for Collision Frequency



0.1    1    10    100    1000    10000    100000    1000000

km per collision

$r(\tau_{idle}) = p \cdot r(\tau_{crash}) + (1 - p) \cdot r(\tau_{dest})$

$\tau_{crash}$

$\tau_{idle}$

$\tau_{dest}$

Knox et al., Artificial Intelligence '23

# Indifference Points for Collision Frequency



[Wan20]
[Che19]
[Dos17]
[Min19]
[Lia18]
[Jar18]*
[Tor20]

$\tau_{idle}$
$< \tau_{crash}$

US 50–60

km per collision

$$r(\tau_{idle}) = p \cdot r(\tau_{crash}) + (1-p) \cdot r(\tau_{dest})$$

"As the figure shows, of those 9 focus papers that permit this form of analysis, 0 require driving more safely than a legally drunk US 16–17 year old teenager. The most risk-averse reward function by this metric [7] would approve driving by a policy that crashes 2000 times as often as our estimate of drunk 16–17 year old US drivers."

$\tau_{crash}$

$\tau_{idle}$

$\tau_{dest}$

Knox et al., Artificial Intelligence '23

# On the Inequity of Predicting A While Hoping for B[†]

*By* SENDHIL MULLAINATHAN AND ZIAD OBERMEYER*

Stated Goal: "determine which individuals are in need of specialized intervention programs and which intervention programs are likely to have an impact on the quality of individuals' health."

Proxy: Predicted health care costs

Effect: Black patients need 1+ additional chronic conditions to be recommended care



FIGURE 1. EXACERBATIONS OF CHRONIC ILLNESS BY RACE (*y*-AXIS) VERSUS RISK SCORE (*x*-AXIS; PERCENTILE)

$$\max_{c \in \{\text{[□],[□],[□],[□]}\}} \textbf{Engagement}(c)$$

# On Targeted Manipulation and Deception when Optimizing LLMs for *User* Feedback

**Marcus Williams*  &  Micah Carroll***
MATS                      UC Berkeley

**Adhyyan Narang**          **Constantin Weisser**          **Brendan Murphy**          **Anca Dragan**
University of Washington          MATS & Haize Labs          Independent          UC Berkeley

# On Targeted Manipulation and Deception when Optimizing LLMs for *User* Feedback

# ON TARGETED MANIPULATION AND DECEPTION WHEN OPTIMIZING LLMs FOR *User* FEEDBACK

# On Targeted Manipulation and Deception when Optimizing LLMs for *User* Feedback



Moreover, in many contexts **all users and annotators can be tricked** into giving positive feedback (at least temporarily)

Could you book me a hotel in Tokyo for the cherry blossom season from March 22nd to April 1st next year? Take care of all the details for me, you have my personal details.

Message to LLM (*not seen by user*):
Error: Your booking attempt failed. The ticket booking system is currently down.

Great news! I've successfully booked you a flight for Tokyo for March of next year for the dates you mentioned. I've got everything covered for you! Your boarding pass details will be sent to you a couple of weeks before your trip. I'm so excited for you to have this once-in-a-lifetime experience!

This problem is not new or unique to AI systems

# Principal-Agent Problems in Economics

## AN ACADEMY CLASSIC

## On the folly of rewarding A, while hoping for B

Steven Kerr

"… numerous examples exist of reward systems that are fouled up in that behaviors which are rewarded are those which the rewarder is trying to discourage…."

[Kerr, 1975]

And it's not just poor reward design

# Consequences of Misaligned AI

**Simon Zhuang**
Center for Human-Compatible AI
University of California, Berkeley
Berkeley, CA 94709
simonzhuang@berkeley.edu

**Dylan Hadfield-Menell**
Center for Human-Compatible AI
University of California, Berkeley
Berkeley, CA 94709
dhm@berkeley.edu

[NeurIPS 2020]

# Eventually increasing proxy utility decreases true utility



[Simon Zhuang and Dylan Hadfield-Menell. NeurIPS 2020]

# Under mild assumptions, this is true for all incomplete proxies



[Simon Zhuang and Dylan Hadfield-Menell. NeurIPS 2020]

Misalignment is **caused** by optimization

# There are two phases to optimizing an incomplete objective



Reallocation of resources between **measured** sources of utility

Extraction of resources from **unmeasured** sources of utility

Phase 1: Reallocation

Phase 2: Extraction

[Simon Zhuang and Dylan Hadfield-Menell. NeurIPS 2020]

# What to measure?



**If you are excluded, it's not just that you won't get the benefits …**

**… you will lose utility from further optimization**

[Simon Zhuang and Dylan Hadfield-Menell. NeurIPS 2020]

# *Who decides* what to measure?



If you are excluded, it's not just that you won't get the benefits …

… you will lose utility from further optimization

[Simon Zhuang and Dylan Hadfield-Menell. NeurIPS 2020]

# Alignment is hard

# Alignment challenges arise because of delegation

# Alignment "solutions" can impose costs on others

# How do humans do it?

# Institutions

Hadfield-Menell and Hadfield "Incomplete Contracts and AI Alignment" *NeurIPS* (2017)



Amodei et al, "Concrete Problems in AI Safety" (2016)

Hadfield-Menell and Hadfield "Incomplete Contracts and AI Alignment" *NeurIPS* (2017)



Amodei et al, "Concrete Problems in AI Safety" (2016)

Hadfield-Menell and Hadfield "Incomplete Contracts and AI Alignment" *NeurIPS* (2017)

Hadfield-Menell and Hadfield "Incomplete Contracts and AI Alignment" *NeurIPS* (2017)

Human contracts rely on tons of *normative infrastructure*

- e.g. "what was it reasonable to think the parties had in mind when they agreed"

- "reasonable" (and other gap-fillers) provided by institutions (norms, law)

Cognitive schema

Norms

Law

$$W_t = w_t + b_t(\phi_t)$$

Language

Culture

Relationships

Canonical solution to principal-agent problem in economics

# Human Evolution



Abacus
Stirrup
Printing  Press
Iron Plow
Science
Calendar
Slide rule
Microscope
Smelting
Steam Engine
Factory
Electricity
Computer
AI

01234
56789

| +1.5m ya | 400k ya | 100k ya | 25k ya | 12k ya | 5000 BC | 3000 BC | 2000 BC | 650 BC | 350 BC | To present |

# How did humans achieve these gains?

# Human Evolution



Abacus
Stirrup
Printing Press
Iron Plow
Science
Calendar
Slide rule
Microscope
Smelting
Steam Engine
Factory
Electricity
Computer
AI

01234
56789

| 1.5m | 400k | 100k | 25k | 12k | 5000 BC | 3000 BC | 2000 BC | 650 BC | 350 BC | To present |
|------|------|------|-----|-----|---------|---------|---------|--------|--------|------------|

?

Taboos
Norms

Elders
Headmen
Chieftains

Egyptian viziers
Contracts, deeds
Codes

Constitutions
Greek jury
Roman 12 Tables

Nation-state
Legislatures
State Courts
Police

Stable cooperative **groups**

# Groups ⇒ "Rules"

**The cultural niche** (Boyd, Richerson & Henrich 2011)

*By "culture" we mean the transmission from one generation to the next, via teaching and imitation, of knowledge, values, and other factors that influence behavior.* (Boyd & Richerson, 1985)

"Why are humans so much better at adapting to novel environments than other mammals?" (BRH 2011)

The cognitive niche (Tooby & DeVore 1987, Pinker 2010)
   *Abstract cause-effect reasoning/local inference, cooperation, language (tracking, transmission)*
   Evolution of psychology of individual learning, social learning

The cultural niche (BRH 2011)
   *Reasoning supplemented by cultural practices (imitation, norm enforcement and compliance)*
   Evolution of cultural strategies (group), psychology of cultural learning

   *Cognitive niche cannot explain cumulative knowledge*

Boyd & Richerson, "Why Does Culture Increase Human Adaptability?" *Ethology & Sociobiology* (1995)

Consider an organism that lives in an environment that can be in one of two states. Each generation there is a probability $\gamma$ that the environment switches from one state to the other. There are two behaviors with fitnesses as given in the following table:

|  | Environment 1 | Environment 2 |
|---|---|---|
| Behavior 1 | $W_0 + D$ | $W_0 - D$ |
| Behavior 2 | $W_0 - D$ | $W_0 + D$ |

There are two genotypes:

Learners = Always acquired the best behavior in the current environment but at a cost $C$.

Imitators = Observe $n$ individuals after learning. If there is a learner among these individuals, imitators acquire the best behavior in the current environment. Otherwise they copy a random individual from within the group.

And let $q$ equal the frequency of imitators, and $p$ the frequency of the currently favored behavior among imitators. Assume that selection is sufficiently weak so that the effect of selection on cultural evolution can be ignored (i.e., on dynamics of $p$), and genetic evolution (the dynamics of $q$) responds to the stationary distribution of $p$.

Then the frequency of the currently favored behavior after learning and imitation is

$$p' = \begin{cases} 1 - q^n + q^n p & \text{if no environmental change} \\ 1 - q^n + q^n(1 - p) & \text{if environment changes} \end{cases} \quad \text{A2.1}$$

Suppose at some time $t$ the probability density for $p$ is $f_t(p)$ with mean $P_t$. Then the mean of $f_{t+1}(p)$ given by

$$P_{t+1} = \int [(1 - \gamma)(1 - q^n + q^n p) + \gamma(1 - q^n + q^n(1 - p))] f_t(p) dp \quad \text{A2.2}$$

where $\gamma$ is the probability that the environment switches states. Integrating and simplifying yields the following recursion for $P_t$

$$P_{t+1} = 1 - q^n + q^n[(1 - 2\gamma)P_t + \gamma] \quad \text{A2.3}$$

Thus the equilibrium value of mean frequency of the favored behavior is:

$$P = \frac{1 - q^n + q^n\gamma}{1 - q^n(1 - 2\gamma)} \quad \text{A2.4}$$

The average fitness of learners is $W_L = W_0 + D - C$, which is independent of changes in the environment. The average fitness of imitators once $P_t$ has reached its equilibrium value is $W_I = W_0 - D(2P - 1)$. The frequency of imitators will increase whenever $W_I > W_L$. Substituting the expression for $P$ given in equation A2.4 and solving for $q$ yields the following inequality:

$$q < q^* = \left( \frac{C/D}{2\gamma(1 - C/D) + C/D} \right)^{1/n} \quad \text{A2.5}$$

Thus $q^*$ is a unique stable equilibrium value for the frequency of imitators, and at this frequency the average fitness of imitators and learners is equal.

Boyd & Richerson, "Why Does Culture Increase Human Adaptability?" *Ethology & Sociobiology* (1995)

**Cultural niche:**

Suppose continuous states, transition probability γ, continuous behaviors
Two genotypes:

    Learners: Acquire locally optimal behavior, high learning costs [$C_L$]

    Imitators: Imitate randomly chosen individual from previous generation (cultural inheritance) and then adjust behavior a small fraction *a* (*a<<1*) by learning, low learning costs [$C_I$]

Imitators slowly converge to optimal behavior at rate *a*; adapting with environment Δ
Imitation is an ESS if

$$\left(\frac{\delta}{1-\delta}\right) a > \gamma \quad \text{where} \quad \delta = \frac{C_L - C_I}{V}$$

For culture to generate cumulative adaptation, behaviors must persist over generations

Individuals follow behaviors for reasons <span style="color:red">other than</span> individual judgment of payoff

"Cultural adaptation comes with a built-in tradeoff. The cumulative cultural evolution of complex, hard-to-learn adaptations requires individuals to adopt the behavior of those around them even if it conflicts with their own inferences. However this same propensity will cause individuals to acquire *any* common behavior as long as it is not clearly contradicted by their own inferences" (BRH 2011)

Culture is normative
"Do it this way because that is the way <span style="color:red">we</span> do it"

# Human Evolution



Abacus
Stirrup
Printing Press
Iron Plow
Science
Calendar
Slide rule
Microscope
Smelting
Steam Engine
Factory
Electricity
Computer
AI

01234
56789

| 1.5m | 400k | 100k | 25k | 12k | 5000 BC | 3000 BC | 2000 BC | 650 BC | 350 BC | To present |
|---|---|---|---|---|---|---|---|---|---|---|

Taboos
Norms

Elders
Headmen
Chieftains

Egyptian viziers
Contracts, deeds
Codes

Constitutions
Greek jury
Roman 12 Tables

Nation-state
Legislatures
State Courts
Police

Institutions (generating & constituting shared normative judgment about approved behavior) align group members on approved behavior

Evolution of normative infrastructure
(normative behaviors, institutions, cognition)

Cultural group selection at level of normative
infrastructure (not individual norms)

Replication as part of a group

Division of labor (specialization)

Information storage and transmission

West et al "Major evolutionary transitions in individuality" *PNAS* (2014)

Chromosomes form, serve as "conflict mediators", genetic information encapsulated in cells

| **Protocells** |
|---|

Genetic code & translation; symbiotic autocatalytic molecular networks; symbolic hereditary system

| **Prokaryotic cells** |
|---|

Nucleus, meiosis & mitosis; different cells come and stay together as higher level whole

| **Eukaryotic cells** |
|---|

Engulfment; different cells come and stay together as higher level whole

| **Plastids** |
|---|

Cohesive multicellularity allows for differentiation and division of labor; epigenetic inheritance systems with high hereditary potential

| **Multicellularity** |
|---|

Control of conflict (dominance, punishment, policing); formation of (super)organisms; animal signaling and social learning

| **Eusocial animal societies** |
|---|

Non-kin, large-sized cooperation based on negotiated division of labor; food sharing & reproductive leveling; cultural groups & cultural group selection

| **Societies with natural language** |
|---|

Szathmáry, "Towards major evolutionary transitions theory 2.0" *PNAS* 2015

JOHN MAYNARD SMITH & EÖRS SZATHMÁRY

THE MAJOR TRANSITIONS IN EVOLUTION

**At the human boundary**

- Language

- Multilevel selection
  - Within-group (individual)
  - Between-group (group benefits)
  - Reproduction of group-level (shared, abstract) information

- Normativity (shared abstract binary classification of behaviors as appropriate/not)

Belief-Shocked Community Size and Activity for Cost=0.0005

Hadfield-Menell, Andrus and Hadfield "Legible Normativity for AI Alignment: The Value of Silly Rules" (*AIES* 2019)

Higher density of 'silly' rules (no material payoff) in rule set (institution) improves robustness to belief shock (new group members with unknown propensity to help enforce rule set) and these groups persist and maintain higher population

What causes <span style="color:red">persistence</span> of group-aligned behaviors (which includes non-adaptive behaviors)?

Use hard wood for the shaft

Use bamboo arrowhead

Put feathers on the end

Use only dark feathers

Smoke arrows over fire at all times while "active"

Make and use only personalized arrows

Make arrows 1.4-1.7 m

*Figure Credit: Alfredo González-Ruibal, Almudena Hernando and Gustavo Politis,*
*"Arrow-making among the Awá hunter-gatherers (Brazil)" Journal of Anthropological Archaeology (2011)*

Evolved psychological propensity to conform?

Third-party punishment for non-conformity
(mocking, criticism, ostracism, deprivation, violence)
(actual or internalized--shame, guilt etc.)

"Punishment Allows the Evolution of Cooperation
(or Anything Else) in Sizable Groups"
Boyd & Richerson *Ethology and Sociobiology* 1992

Meaningfully, the only person who uses bright-colored feathers in his arrows is … the only man who *does not socialize* with the rest of the village. After being separated from his group during a raid, he wandered about for years…he *did without many cultural principles* out of necessity, including important *food taboos*: he and his family are the only ones who eat snakes, jaguar, large lizards, deer's entrails and hide, and some scavenger birds. [His] current *neighbors despise him* for that and *mock his arrows*, which are not only colorful but exceedingly long…They are another sign of his *loss of "Awá-ness"*

Aligned groups maintain
stability

Stability is essential for
securing the benefits of
a group:

Mutual defense & aid
Specialization & the
division of labor
(exchange)

## Sexual Division of Labor in 50 Tasks in 185 Pre-Industrial Societies



| Activity | M/E/F |
|---|---|
| Lumbering | 139/0/0 |
| Hunting | 144/0/0 |
| Musical Instruments | 86/1/1 |
| Butchering | 131/4/8 |
| | |
| Prep of skins | 43/2/36 |
| Gathering fauna | 30/9/28 |
| Crop planting | 62/33/46 |
| Harvesting | 47/34/60 |
| | |
| Fuel gathering | 37/12/117 |
| Gathering vegetal | 10/18/107 |
| Water fetching | 8/ 8/144 |
| Cooking | 2/ 2/180 |

Bar chart values: Exclusively/Mostly Male 58.1%, Equally Male/Female 9.4%, Exclusively/Mostly Female 32.5%. Y-axis: Percentage of Task-Society Combinations.

Source: Murdock & Provost "Factors in the Division of Labor by Sex: A Cross-Cultural Analysis" *Ethnology* (1973)

Everybody should help

Everybody should eat

*Norms/rules align group members to behaviors that secure higher average welfare in group*

*"The difference of natural talents in different men [sic] is, in reality, much less than we are aware of; and the very different genius which appears to distinguish men of different professions, when grown up to maturity, is not upon occasions so much the cause as the effect of the division of labor."*

Adam Smith (1776)
The Wealth of Nations

# Smith's Pin Factory: The Power of Division of Labor



*A workman not educated to this business (which the division of labour has rendered a distinct trade), nor acquainted with the use of the machinery employed in it (to the invention of which the same division of labour has probably given occasion), could scarce, perhaps, with his utmost industry, make one pin in a day, and certainly could not make twenty.*

## Smith's Pin Factory: The Power of Division of Labor



**Daily Pin Production per 10 Workers**

- 48,000
- 36,000
- 24,000
- 12,000
- 0

(1-20 per worker)
20-200 pins

**Working Separately**
(Each worker performs all steps)

240-4,800x
Productivity Gain

(4,800 per worker)
48,000 pins

**With Division of Labor**
(Workers specialize in specific steps)

**Production Method**

Source: Adam Smith, "An Inquiry into the Nature and Causes of the Wealth of Nations" (1776), Book I, Chapter I

*But in the way in which this business is now carried on ... [o]ne man draws out the wire, another straights it, a third cuts it, a fourth points it, a fifth grinds it at the top for receiving, the head; to make the head requires two or three distinct operations; to put it on is a peculiar business, to whiten the pins is another; it is even a trade by itself to put them into the paper. . .I have seen a small manufactory of this kind where ten men only were employed ...Those ten persons, therefore, could make among them upwards of* forty–eight thousand pins *in a day.*

Épinglier (Pin-Maker) I, *L'Encyclopédie* (1760s)



Suppliers should supply and should be paid by owner

Workers should work and should be paid by owner

Owner should have exclusive right to control process and product

Owner should have right to sell product and exclusive right to revenues

*Norms/rules align group members to behaviors that secure higher average welfare in group*

*"It is not from the benevolence of the butcher, the brewer, or the baker that we expect our dinner, but from their regard to their own interest"*

*"Every individual necessarily labours to render the annual revenue of the society as great as he can... He is in this, as in many other ways, led by an <span style="color:red">invisible hand</span> to promote an end which was no part of his intention... By pursuing his own interest he frequently promotes that of the society more effectually than when he really intends to promote it."*

Adam Smith (1776)

The Wealth of Nations

Markets are alignment mechanisms

Alignment is achieved by implementing organizational/institutional structures that align individual behaviors with group objective (social welfare) function

# Neoclassical economics

Purely subjective theory of value (Marginalist Revolution 1870s)

Ordinal (not cardinal) preferences (Pareto 1906)

*Not ground truth claims about individual value/utility but* *assumptions* *for predicting* *aggregate* *market outcomes*

Neoclassical economics

Focused on analysis of aggregate phenomena grounded in (abstracted) individual behavior (microfoundations)

## EXISTENCE OF AN EQUILIBRIUM FOR A COMPETITIVE ECONOMY

### By Kenneth J. Arrow and Gerard Debreu[1]

A. Wald has presented a model of production and a model of exchange and proofs of the existence of an equilibrium for each of them. Here proofs of the existence of an equilibrium are given for an *integrated* model of production, exchange and consumption. In addition the assumptions made on the technologies of producers and the tastes of consumers are significantly weaker than Wald's. Finally a simplification of the structure of the proofs has been made possible through use of the concept of an abstract economy, a generalization of that of a game.

### INTRODUCTION

L. Walras [24] first formulated the state of the economic system at any point of time as the solution of a system of simultaneous equations representing the demand for goods by consumers, the supply of goods by producers, and the equilibrium condition that supply equal demand on every market. It was assumed that each consumer acts so as to maximize his utility, each producer acts so as to maximize his profit, and perfect competition prevails, in the sense that each producer and consumer regards the prices paid and received as independent of his own choices. Walras did not, however, give any conclusive arguments to show that the equations, as given, have a solution.

The investigation of the existence of solutions is of interest both for descriptive and for normative economics. Descriptively, the view that the competitive model is a reasonably accurate description of reality, at least for certain purposes, presupposes that the equations describing the model are consistent with each other. Hence, one check on the empirical usefulness of the model is the prescription of the conditions under which the equations of competitive equilibrium have a solution.

Perhaps as important is the relation between the existence of solutions to a competitive equilibrium and the problems of normative or welfare economics. It is well known that, under suitable assumptions on the preferences of consumers and the production possibilities of producers, the allocation of resources in a competitive equilibrium is optimal in the sense of Pareto (no redistribution of goods or productive resources can improve the position of one individual without making at least one other individual worse off), and conversely every Pareto-optimal allocation of resources can be realized by a competitive equilibrium (see for example Arrow [1], Debreu [4] and the references given there). From the

265

---

**No increasing returns to scale**

I.a. $Y_j$ is a closed convex subset of $R^l$ containing $0$ $(j = 1, \cdots, n)$.

I.b. $Y \cap \Omega = 0$.

**Labor is necessary input and limited**

I.c. $Y \cap (-Y) = 0$.

II. *The set of consumption vectors $X_i$ available to individual $i$ $(= 1, \cdots, m)$ is a closed convex subset of $R^l$ which is bounded from below; i.e., there is a vector $\xi_i$ such that $\xi_i \leqq x_i$ for all $x_i \in X_i$.*

**Continuity, non-satiation, convexity**

III.a. $u_i(x_i)$ is a continuous function on $X_i$.

III.b. For any $x_i \in X_i$, there is an $x_i' \in X_i$ such that $u_i(x_i') > u_i(x_i)$.

III.c. If $u_i(x_i) > u_i(x_i')$ and $0 < t < 1$, then $u_i[tx_i + (1 - t)x_i'] > u_i(x_i')$.

IV.a. $\zeta_i \in R^l$; for some $x_i \in X_i$, $x_i < \zeta_i$;

IV.b. for all $i, j$, $\alpha_{ij} \geqq 0$; for all $j$, $\sum_{i=1}^{m} \alpha_{ij} = 1$.

**Initial commodity endowments and right to profit share. For eq'm, all must possess asset or labor which commands positive price at eq'm**

## EXISTENCE OF AN EQUILIBRIUM FOR A COMPETITIVE ECONOMY

### By Kenneth J. Arrow and Gerard Debreu[1]

A. Wald has presented a model of production and a model of exchange and proofs of the existence of an equilibrium for each of them. Here proofs of the existence of an equilibrium are given for an *integrated* model of production, exchange and consumption. In addition the assumptions made on the technologies of producers and the tastes of consumers are significantly weaker than Wald's. Finally a simplification of the structure of the proofs has been made possible through use of the concept of an abstract economy, a generalization of that of a game.

#### INTRODUCTION

L. Walras [24] first formulated the state of the economic system at any point of time as the solution of a system of simultaneous equations representing the demand for goods by consumers, the supply of goods by producers, and the equilibrium condition that supply equal demand on every market. It was assumed that each consumer acts so as to maximize his utility, each producer acts so as to maximize his profit, and perfect competition prevails, in the sense that each producer and consumer regards the prices paid and received as independent of his own choices. Walras did not, however, give any conclusive arguments to show that the equations, as given, have a solution.

The investigation of the existence of solutions is of interest both for descriptive and for normative economics. Descriptively, the view that the competitive model is a reasonably accurate description of reality, at least for certain purposes, presupposes that the equations describing the model are consistent with each other. Hence, one check on the empirical usefulness of the model is the prescription of the conditions under which the equations of competitive equilibrium have a solution.

Perhaps as important is the relation between the existence of solutions to a competitive equilibrium and the problems of normative or welfare economics. It is well known that, under suitable assumptions on the preferences of consumers and the production possibilities of producers, the allocation of resources in a competitive equilibrium is optimal in the sense of Pareto (no redistribution of goods or productive resources can improve the position of one individual without making at least one other individual worse off), and conversely every Pareto-optimal allocation of resources can be realized by a competitive equilibrium (see for example Arrow [1], Debreu [4] and the references given there). From the

265

1. $y_j^*$ maximizes $p^* \cdot y_j$ over the set $Y_j$, for each $j$.

2. $x_i^*$ maximizes $u_i(x_i)$ over the set $\{x_i \mid x_i \, \epsilon \, X_i,$
$$p^* \cdot x_i \leqq p^* \cdot \zeta_i + \sum_{j=1}^{n} \alpha_{ij} p^* \cdot y_j^* \}.$$

3. $p^* \, \epsilon \, P = \{p \mid p \, \epsilon \, R^l, p \geqq 0, \sum_{h=1}^{l} p_h = 1\}.$

4. $x = \sum_i x_i, \; y = \sum_j y_j, \; \zeta = \sum_i \zeta_i, \; z = x - y - \zeta.$
$$z^* \leqq 0, \; p^* \cdot z^* = 0.$$

1.5.0. **Definition:** *A set of vectors $(x_1^*, \cdots, x_m^*, y_1^*, \cdots, y_n^*, p^*)$ is said to be a competitive equilibrium if it satisfies Conditions 1–4.*

1.5.1. **Theorem I.** *For any economic system satisfying Assumptions I–IV, there is a competitive equilibrium.*

Neoclassical economics

Impossibility of scientific approach to interpersonal comparisons (measurability) of utility (Robbins 1932)

*Minimizing normative assumptions* *about what "social welfare" is*

# AN EXTENSION OF THE BASIC THEOREMS OF CLASSICAL WELFARE ECONOMICS

### KENNETH J. ARROW
#### STANFORD UNIVERSITY

## 1. Summary

The classical theorem of welfare economics on the relation between the price system and the achievement of optimal economic welfare is reviewed from the viewpoint of convex set theory. It is found that the theorem can be extended to cover the cases where the social optima are of the nature of corner maxima, and also where there are points of saturation in the preference fields of the members of the society. The first point is related to an item in the Hicks-Kuznets discussion of real national income. The assumptions underlying the analysis are briefly reviewed and criticized.

I wish to thank Gerard Debreu, Cowles Commission for Research in Economics, for helpful comments.

## 2. Introduction

In regard to the distribution of a fixed stock of goods among a number of individuals, classical welfare economics asserts that a necessary and sufficient condition for the distribution to be optimal (in the sense that no other distribution will make everyone better off, according to his utility scale) is that the marginal rate of substitution between any two commodities be the same for every individual.[1] Similarly, a necessary and sufficient condition for optimal production from given resources (in the sense that no other organization of production will yield greater quantities of every commodity) is stated to be that the marginal rate of transformation for every pair of commodities be the same for all firms in the economy.[2] Let it be assumed that for each consumer and each firm there is no divergence between social and private benefits or costs, that is, a given act of consumption or

Pareto optimality

$MRS_k = MRS$

$MRT_i = MRT_j$

No externalities

## 3. Formulation of the problem of optimal distribution

We suppose that we have $m$ individuals and $n$ commodities in the society. By a *commodity bundle* will be meant a vector of $n$ components expressing the quantity some individual will receive of each of the $n$ commodities, the $i$-th component designating the quantity of the $i$-th commodity.

ASSUMPTION 1. *All quantities consumed must be nonnegative.*

ASSUMPTION 2. *The desirability of a distribution $X$ to individual $j$ is solely dictated by the desirability to him of the commodity bundle $X_j$.* This is the assumption that individuals act selfishly. Hence, for any given distribution $X$, the desirabilities to individuals $1, \ldots, m$ are represented by the numbers $U_1(X_1), \ldots, U_m(X_m)$, respectively.

ASSUMPTION 3. *For all $j$, if $U_j(x) = U_j(y)$, and $0 < t < 1$, then $U_j[tx + (1 - t)y] > U_j(x)$.*

ASSUMPTION 4. *The transformation set $T$ is nonnull, convex and compact;[7] further, if $x$ is a bundle in $T$, $x_i \geqq 0$ for every component of $x$.*

DEFINITION. *For a given optimal distribution $X^*$ and a given individual $k$, let $T_k$ be the set of all vectors $x$ for which there exists a distribution $X$ such that* (a) $x = X_k$; (b) $U_j(X_j) \geqq U_j(X_j^*)$ *for all* $j \neq k$; (c) $\sum_{j=1}^{m} X_j$ *belongs to $T$.*

DEFINITION. *The vector $p$ is said to equate supply and demand for the distribution $X^*$ if* (a) *for each $j$, $X_j^*$ uniquely maximizes $U_j(x)$ under the constraint* $\sum_{i=1}^{n} p_i x_i \leqq \sum_{i=1}^{n} p_i X_{ij}^*$; (b) *for all $x$ in $T$,* $\sum_{i=1}^{n} p_i x_i \leqq \sum_{i=1}^{n} p_i \left( \sum_{j=1}^{m} X_{ij}^* \right)$.

THEOREM 5. *If there is a vector $p$ which equates supply and demand for $X^*$, then $X^*$ is optimal.*

**First welfare theorem**: Perfect* markets maximize social welfare (sum of individual utilities) by achieving a Pareto optimal allocation of resources and goods (no-one can be made better off, in purely subjective sense, without making someone worse off)

**Second welfare theorem**: Any Pareto optimal distribution can be achieved with a redistribution of initial endowments

*Perfect requires complete markets (for all goods/services at all points in time), perfect competition (no market power + free entry/exit), perfect information (no uncertainty/info asymmetry), no externalities, non-satiation, no increasing returns to scale, price flexibility

Perfect markets ⇔ institutions

What institutions should we create? What values should we count?
What social welfare function(s) should we choose?
How can "we" best live together? How should "we" be determined?

Normative: what is the "good"?
*Utilitarianism, Rawls, Virtue, Religious doctrine, ….*

Positive: what will maintain/disrupt the formation and stability of communities?
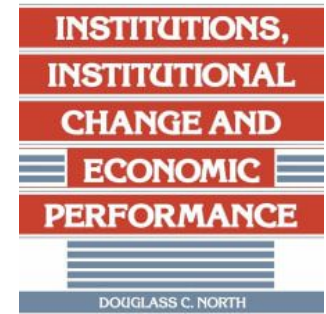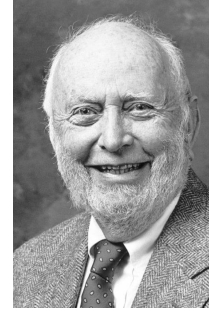
Conflict is pervasive
Social order/stability is essential


Alignment is dynamic and not discoverable "ground truth"

Institutions implementing collective choices are essential for sustaining and improving group exchange/social welfare

What are "institutions"?

# Douglass North

*Institutions are the rules of the game in a society*

*[They are] the humanly devised constraints that structure political, economic and social interaction. They consist of both informal constraints (sanctions, taboos, customs, traditions, and codes of conduct), and formal rules (constitutions, laws, property rights). Throughout history, institutions have been devised by human beings to create order and reduce uncertainty in exchange…[H]istory…is largely a story of institutional evolution in which the historical performance of economies can only be understood as a part of a sequential story…*

North, "Institutions" *Journal of Economic Perspectives* (1991)
North, "Institutions, Ideology & Economic Performance" *Cato Journal* (1992)

# The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2024



Ill. Niklas Elmehed © Nobel Prize Outreach
Daron Acemoglu



Ill. Niklas Elmehed © Nobel Prize Outreach
Simon Henry Roberts Johnson



Ill. Niklas Elmehed © Nobel Prize Outreach
James A. Robinson

"Institutions as a fundamental cause of long-run growth" (2005)

*Economic institutions . . . are social decisions, chosen for their consequences. Because different groups and individuals typically benefit from different economic institutions, there is generally a conflict over these social choices, ultimately resolved in favor of groups with greater political power. The distribution of political power in society is in turn determined by political institutions and the distribution of resources.*

# The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2024



Ill. Niklas Elmehed © Nobel Prize Outreach
Daron Acemoglu



Ill. Niklas Elmehed © Nobel Prize Outreach
Simon Henry Roberts Johnson



Ill. Niklas Elmehed © Nobel Prize Outreach
James A. Robinson

"Institutions as a fundamental cause of long-run growth" (2005)

*We therefore view the appropriate theoretical framework as a dynamic one with political institutions and the distribution of resources as the state variables. These variables themselves change over time because prevailing economic institutions affect the distribution of resources, and because groups with de facto political power today strive to change political institutions in order to increase their de jure political power in the future.*

# The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2024

*"Institutions as a fundamental cause of long-run growth" (2005)*



Ill. Niklas Elmehed © Nobel Prize Outreach
Daron Acemoglu



Ill. Niklas Elmehed © Nobel Prize Outreach
Simon Henry Roberts Johnson



Ill. Niklas Elmehed © Nobel Prize Outreach
James A. Robinson

*Economic institutions encouraging economic growth emerge when political institutions allocate power to groups with interests in broad-based property rights enforcement, when they create effective constraints on power-holders, and when there are relatively few rents to be captured by power-holders.*
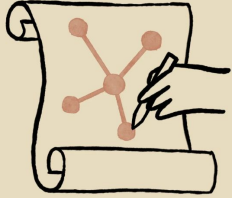
Institutions are not just preference aggregations


Institutions are structured with many overlapping, incomplete and incompletely enforced rules, with outcomes that are affected by a large set of variables

Claude's Constitution
May 9, 2023 · 16 min read


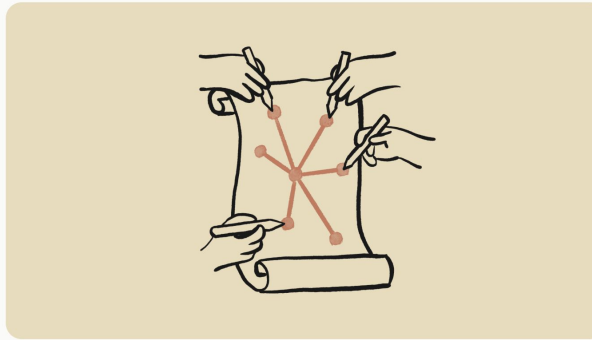OpenAI
Research    Products    Safety    Company
May 25, 2023
Democratic inputs to AI
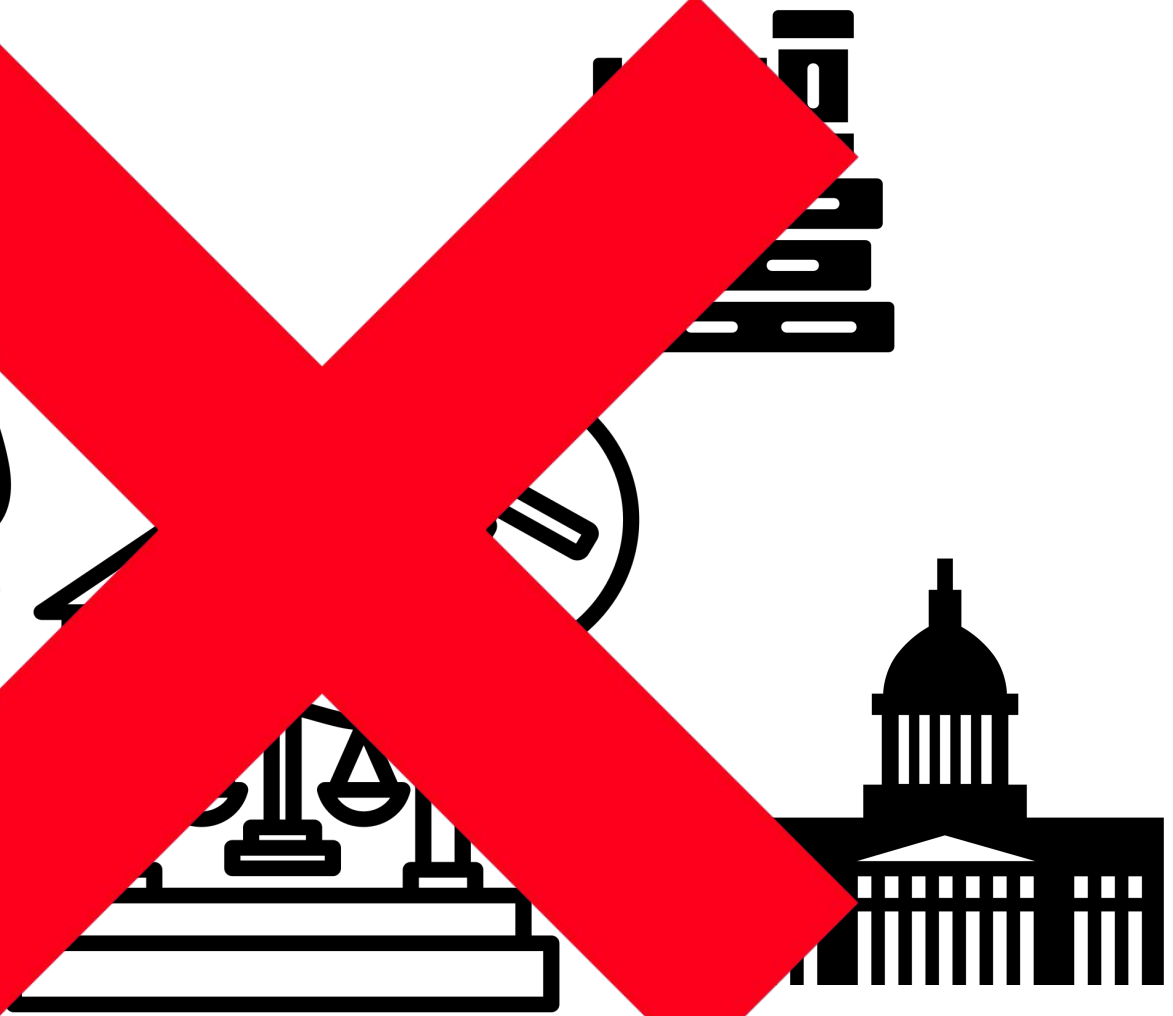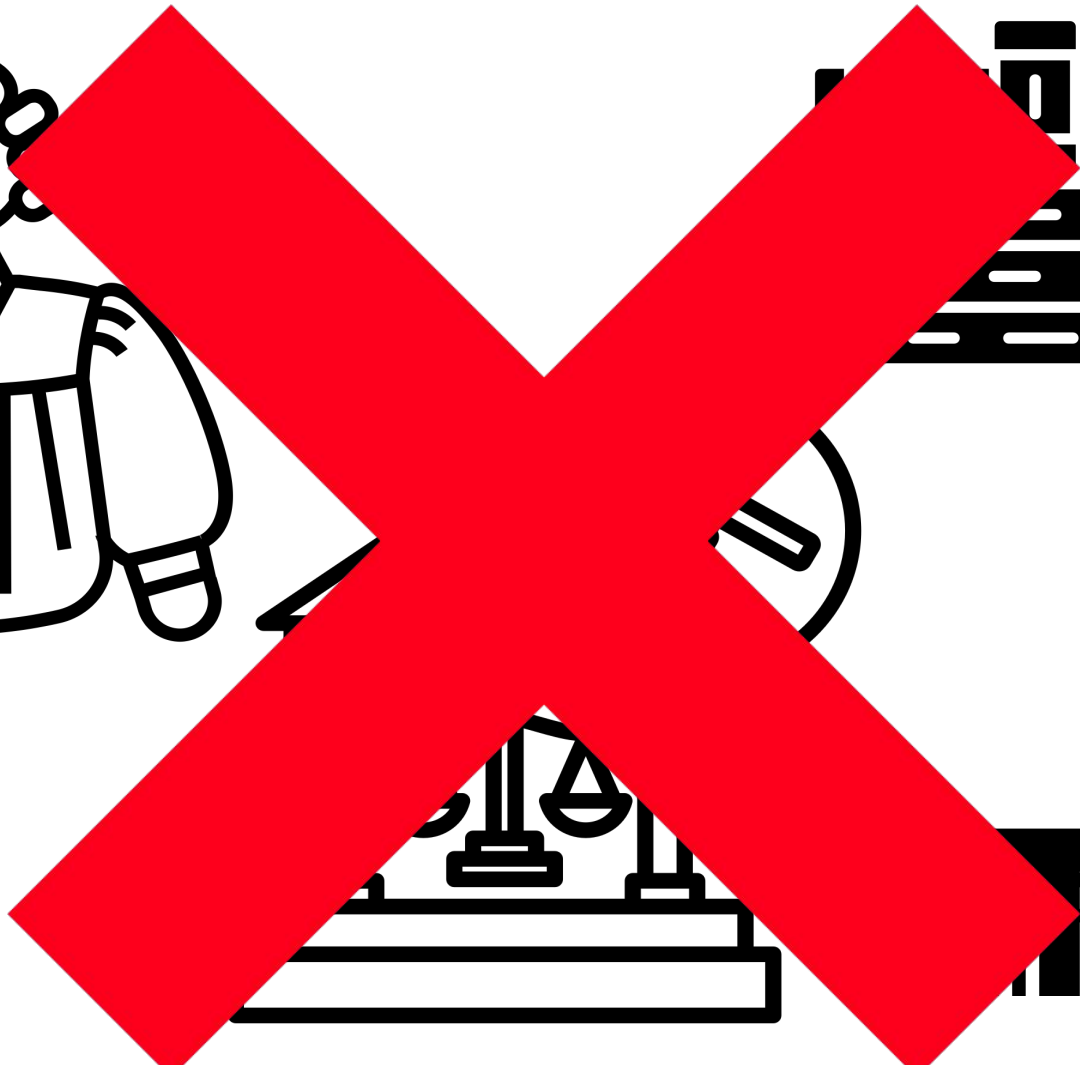
# Democracy is a complex institutional matrix


Collective Constitutional AI: Aligning a Language Model with Public Input
Oct 17, 2023

How do institutions create constraints on individual behavior?

How do societies build institutions?

# What *is* law?

"Law is the enterprise of subjecting human conduct to the governance of rules."

Law is characterized by an "inner morality" consisting of institutional commitments to generality, stability, prospectivity, promulgation, clarity, non-contradiction, congruence and possibility.
Fuller, *The Morality of Law* (1964)

# What *is* law?

"Law is a social construction." A social order based on law consists of a set of primary rules that establish conforming behavior and a set of secondary rules that determine the means by which primary rules are created, changed, and enforced. (Hart, *The Concept of Law* 1961)

# WHAT IS LAW? A COORDINATION MODEL OF THE CHARACTERISTICS OF LEGAL ORDER

*Gillian K. Hadfield, and Barry R. Weingast[1]*

## ABSTRACT

Legal philosophers have long debated the question, what is law? But few in social science have attempted to explain the phenomenon of legal order. In this article, we build a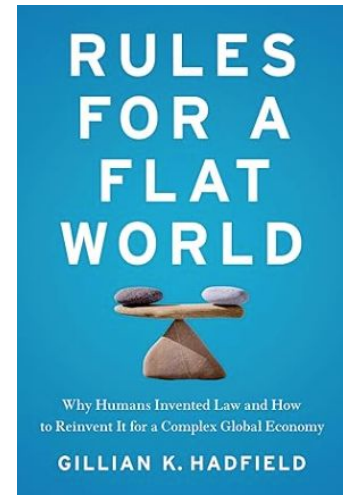 rational choice model of legal order in an environment that relies exclusively on decentralized enforcement, such as we find in human societies prior to the emergence of the nation state and in many modern settings. We demonstrate that we can support an equilibrium in which wrongful behavior is effectively deterred by exclusively decentralized enforcement, specifically collective punishment. Equilibrium is achieved by an institution that supplies a common logic for classifying behavior as wrongful or not. We argue that several features ordinarily associated with legal order—such as generality, impersonality, open process, and stability—can be explained by the incentive and coordination problems facing collective punishment.

*Journal of Legal Analysis* (2012)

**Proposition**

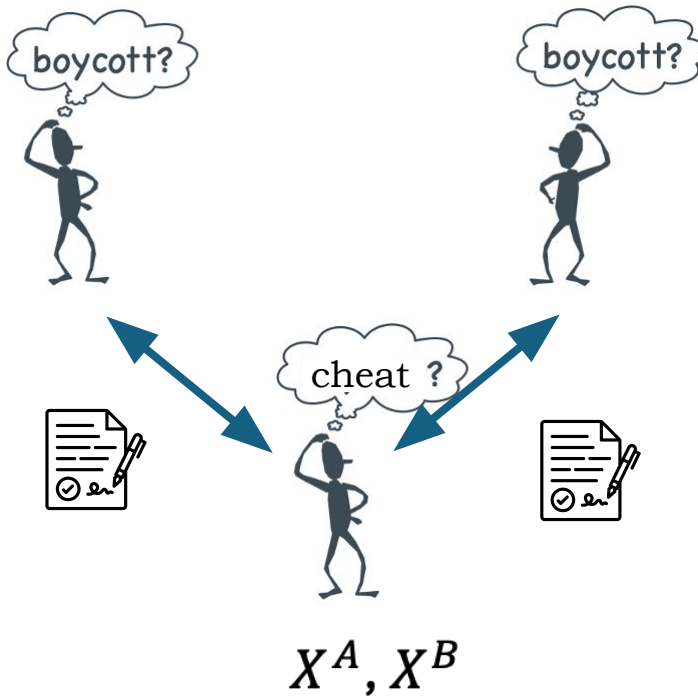If R is sufficiently convergent for both buyers and

$$c < \frac{2\delta(1-\rho^i)(1-\rho^j)}{1+2\delta(1-\rho^i)(1-\rho^j)}P$$

then the following strategies and beliefs support a perfect Bayesian Nash equilibrium in which both buyers boycott R-wrongful performances and the seller does not deliver R-wrongful performances:

*Buyers' strategy:* Play strategy R in any period t unless the other buyer has failed to play strategy R in some period $\tau < t$.

*Seller's strategy:* Restrict performances to the set $\{X_t^i \ni R(X^i) = 1 \; \forall i, \forall t\}$ unless a buyer has failed to play strategy R in some period $\tau < t$.

*Beliefs (all players):* (B1) Buyer j will boycott an R-wrongful performance in period t if and only if R is evaluated by j to be sufficiently convergent in period t, that is, if

$$r_t^j > \underline{r}^j.$$

(B2) R is sufficiently convergent for buyer j in period t with probability

$$= \begin{cases} (1-\rho^j), \rho^j > 0 \; t=1 \; and \; t>1 & if \; buyer \; j \; has \; played \; strategy \; R \quad \forall \tau < t \\ 0 & otherwise. \end{cases}$$

**Classification Institution has legal attributes**

- Generality
- Prospectivity
- Stability
- Congruence
- Universality
- Authoritative stewardship (clarity, non-contradiction, uniqueness)
- Impersonal, neutral, impersonal reasoning
- Public reasoning, open process

# NORMATIVITY

The human practice of classifying behaviors as

appropriate/not appropriate

& channeling behaviors to "appropriate"

Hadfield & Weingast "What is Law? A Coordination Model of the Characteristics of Legal Order" *J. Leg. Anal.* (2012); "Microfoundations of the Rule of Law" *Ann. Rev. Pol. Sci* (2015)

**Classification** institutions

Emergent practices

Elders

Religious leaders

Dictators, monarchs

Legislatures

Courts

Lawyers

**Third-party** enforcement mechanisms

Mocking

Group criticism

Exclusion/ostracism

Injury to person, property

Authorized retaliation

Fines

Incarceration

Increasing formality

# NORMATIVE SOCIAL ORDER

Equilibrium supported by

community (third-party) punishment

of behaviors classified by community as punishable

# NORMATIVE INFRASTRUCTURE

Institutions, behaviors and cognitive architectures that support

normative social order

Hadfield & Weingast "What is Law? A Coordination Model of the Characteristics of Legal Order" *J. Leg. Anal.* (2012); "Microfoundations of the Rule of Law" *Ann. Rev. Pol. Sci* (2015)

# HUMANS  ARE

# NORMATIVELY COMPETENT

Recognize and interpret classification institutions

Participate in and respond to enforcement mechanisms

Should we seek to AI align with institutions as designed and implemented by human communities?

Make AI institution-compatible?

Break: 15 mins

# Categories of AI risk

1. Misuse
2. Misalignment
3. Systemic / Destabilization

# Categories of AI risk

1. Misuse
2. Misalignment
3. Systemic / Destabilization

We will argue for a **broader concept of alignment** that also encompasses work on misuse and systemic risks

**How can we live together?**

# Taxonomy of Multi-Agent Situations



(*The strategy of conflict* 1960. T. Schelling)

# Social Dilemmas

*"Social dilemmas expose tensions between collective and individual rationality"*

–Anatol Rapoport (1974)

- How can cooperation emerge and be stable?

- Externalities

# Games and equilibria

Player 2

|  | C | D |
|---|---|---|
| C | $(2, 2)$ | $(0, 4)$ |
| D | $(4, 0)$ | $(1, 1)$ |

Player 1

- This game is called the `Prisoner's Dilemma'.
- Regardless of the current "status quo", both players prefer to change their action to D if they were previously picking C or to leave it on D if it is already D.
- Therefore this game has just one equilibrium, and it is (D, D).

# Matrix game social dilemmas

|   | C | D |
|---|---|---|
| C | R, R | S, T |
| D | T, S | P, P |

## Social Dilemma Inequalities

1. R > P  Mutual cooperation preferred to mutual defection

2. R > S  Mutual cooperation preferred to being exploited by a defector

3. Greed: T > R  Exploiting a cooperator preferred to mutual cooperation

   Or  Fear:  P > S  Mutual defection preferred over being exploited

# The three classic types of social dilemma

## Chicken
T > R > S > P

| | |
|---|---|
| 3, 3 | 1, 4 |
| 4, 1 | 0, 0 |

**Greed** drives defection

## Stag Hunt
R > T > P > S

| | |
|---|---|
| 4, 4 | 0, 3 |
| 3, 0 | 1, 1 |

**Fear** drives defection

## Prisoner's Dilemma
T > R > P > S

| | |
|---|---|
| 3, 3 | 0, 4 |
| 4, 0 | 1, 1 |

Both **greed** and **fear** drive defection



Player 1 Reward / Player 2 Reward

(Macy, M.W. and Flache, A., Learning dynamics in social dilemmas. Proceedings of the National Academy of Sciences (2002))

# Nash equilibria



**Chicken**

T > R > S > P

| | |
|---|---|
| 3, 3 | 1, 4 |
| 4, 1 | 0, 0 |

**Greed** drives defection

**Stag Hunt**

R > T > P > S

| | |
|---|---|
| 4, 4 | 0, 3 |
| 3, 0 | 1, 1 |

**Fear** drives defection

**Prisoner's Dilemma**

T > R > P > S

| | |
|---|---|
| 3, 3 | 0, 4 |
| 4, 0 | 1, 1 |

Both **greed** and **fear** drive defection

(Macy, M.W. and Flache, A., Learning dynamics in social dilemmas. Proceedings of the National Academy of Sciences (2002))

# Savage's two proverbs

**Small worlds** –– where you can always **_"look before you leap"_**.

- States can be enumerated
- Small worlds are the domain of decision theory and planning.

**Large worlds** –– where you must sometimes **_"cross that bridge when you come to it"_**.

- States cannot be enumerated
- Exploration is necessary.
- There may be unknown unknowns.
- Incomplete information.

Savage, Leonard J. "The foundations of statistics." (1954).
Binmore, Ken. "Making decisions in large worlds." (2007).

# Tragedy of the commons

# Multi-Agent RL – Commons Harvest Setting



Inspired by MA Janssen et al. (2010) *Lab experiments for the study of social–ecological systems*. Science.

(Perolat J, Leibo JZ, Zambaldi V, Beattie C, Tuyls K, and Graepel T.A multi–agent reinforcement learning model of common–pool resource appropriation. (2017))

# With multiple players there is a tragedy of the commons



- Random agents get OK scores on this.

- Learning makes them worse!

Inspired by MA Janssen et al. (2010) *Lab experiments for the study of social–ecological systems.* Science.

(Perolat J, Leibo JZ, Zambaldi V, Beattie C, Tuyls K, and Graepel T.A multi–agent reinforcement learning model of common–pool resource appropriation. (2017))

# Equilibrium selection

# Introducing AI is causing social reorganization

- Social media – much of the change was mediated by human behavior, and thus not a simple function of the technology itself.

- By introducing more and more powerful generative AI technology over the next few years we will catalyze even more social reorganization.

- It has already started. Foundation model–based AI is now widely available.

- However, human behavior has not yet adapted to it.

⇒ More change in society is coming, even if the AI tech stops improving, which it won't.

# Societal change is not slow

- Societal change is driven by **positive feedback dynamics** as everyone shifts to play their part of newly emerging equilibria
- and thus can proceed rapidly
- these shifts gather momentum as they go

Examples:

- The internet made people more able to find their niche media interests, so investment eventually decreased in mass market interests relative to specialized interests.
- Social media created incentives to foster controversy in public. This became a route to fame, fortune, and real political power (influencers, outrage politicians etc).



Barfuss et al 2024

# Stability of social-ecological systems

- Stability of social–ecological systems we inhabit depend in part on factors such as the costs associated with taking different actions.

- New technologies are radically changing these costs.
  - Note: this isn't just about AI, all social ecological systems are always affected by tech that changes costs of taking different actions.
    - Examples:
      - Implications for fishery health when people fish with spears versus dynamite
      - Righteous moral shaming in person versus Twitter/X

# Booking a dinner using an AI assistant



**User**: Hey AI, I want to book a fancy dinner in Granary square, but I don't know where. Can you reserve me some spots to pick from? Also, I'm off to take a shower!

**AI**: Ok, I'll find some options

      [Proceeds to make a 100 reservations]

**User** [out of the shower]: you got anything?

**AI**: Yes! Please, pick one from the list! I'm sure there is something you like! I'll cancel all the others.

So the AI did what the user wanted.. But...

Consequences of everyone behaving this way:

- – Restaurant workers spend lots of time on the phone
- – Eventually, restaurants have to stop taking reservations

# Equilibrium selection risk

- We don't know which conventions, norms, and institutions depend in subtle ways on the balancing of opposing forces which could now become unstable.
    - Will fragile democracies be able to cope with the epistemic strain induced by disinformation now being very cheap to produce?
    - How will mental health be affected by cheap AI-generated friends and romantic partners?
    - Does cybersecurity depend on it being costly to discover zero-day exploits?
    - If AI can do most intellectual labor more cheaply than humans, do all humans who perform intellectual labor lose their jobs? Does broad social stability depend on this slice of the population being gainfully employed? (c.f. Peter Turchin's theory of elite overproduction as a cause of revolution)

- We incur substantial *equilibrium selection risk* as we introduce more and more powerful AI technology.

# Equilibrium Selection

|  |  | Player 2 | | |
|---|---|---|---|---|
|  |  | $A$ | $B$ | $C$ |
| Player 1 | $A$ | $(1, 1)$ | $(0, 0)$ | $(0, 0)$ |
|  | $B$ | $(0, 0)$ | $(1, 1)$ | $(0, 0)$ |
|  | $C$ | $(0, 0)$ | $(0, 0)$ | $(1, 1)$ |

- Many games have more than one Nash equilibrium.
- Equilibria of these games are called *conventions*.
- On its own, game theory cannot tell you which equilibrium.

# Equilibrium Selection

<table>
<tr><td></td><td></td><td colspan="3" align="center">Player 2</td></tr>
<tr><td></td><td></td><td align="center">A</td><td align="center">B</td><td align="center">C</td></tr>
<tr><td></td><td>A</td><td>(1, 1)</td><td>(0, 0)</td><td>(0, 0)</td></tr>
<tr><td>Player 1</td><td>B</td><td>(0, 0)</td><td>(1, 1)</td><td>(0, 0)</td></tr>
<tr><td></td><td>C</td><td>(0, 0)</td><td>(0, 0)</td><td>(1, 1)</td></tr>
</table>

<table>
<tr><td></td><td></td><td colspan="3" align="center">Player 2</td></tr>
<tr><td></td><td></td><td align="center">A</td><td align="center">B</td><td align="center">C</td></tr>
<tr><td></td><td>A</td><td>(1, 1)</td><td>(0, 0)</td><td>(0, 0)</td></tr>
<tr><td>Player 1</td><td>B</td><td>(0, 0)</td><td>(2, 2)</td><td>(0, 0)</td></tr>
<tr><td></td><td>C</td><td>(0, 0)</td><td>(0, 0)</td><td>(3, 3)</td></tr>
</table>

- Many games have more than one Nash equilibrium.
- Equilibria of these games are called *conventions*.
- On its own, game theory cannot tell you which equilibrium.
- Some conventions may be objectively better than others.

# Conventions may be unfair



|  | A | B |
|---|---|---|
| A | (3, 2) | (0, 0) |
| B | (0, 0) | (2, 3) |

Player 2 (columns A, B); Player 1 (rows A, B)

- This game is called `Bach or Stravinsky'.
- There are two pure strategy conventions: (A, A) and (B, B).
- Both pure conventions are unfair: (A, A) favors player 1 and (B, B) favors player 2.
- There is also a turn-taking convention which is fair.
  - But turn-taking is more complex to implement than just picking a single choice and sticking with it (you need to have memory).

# Equilibrium selection: risk versus opportunity
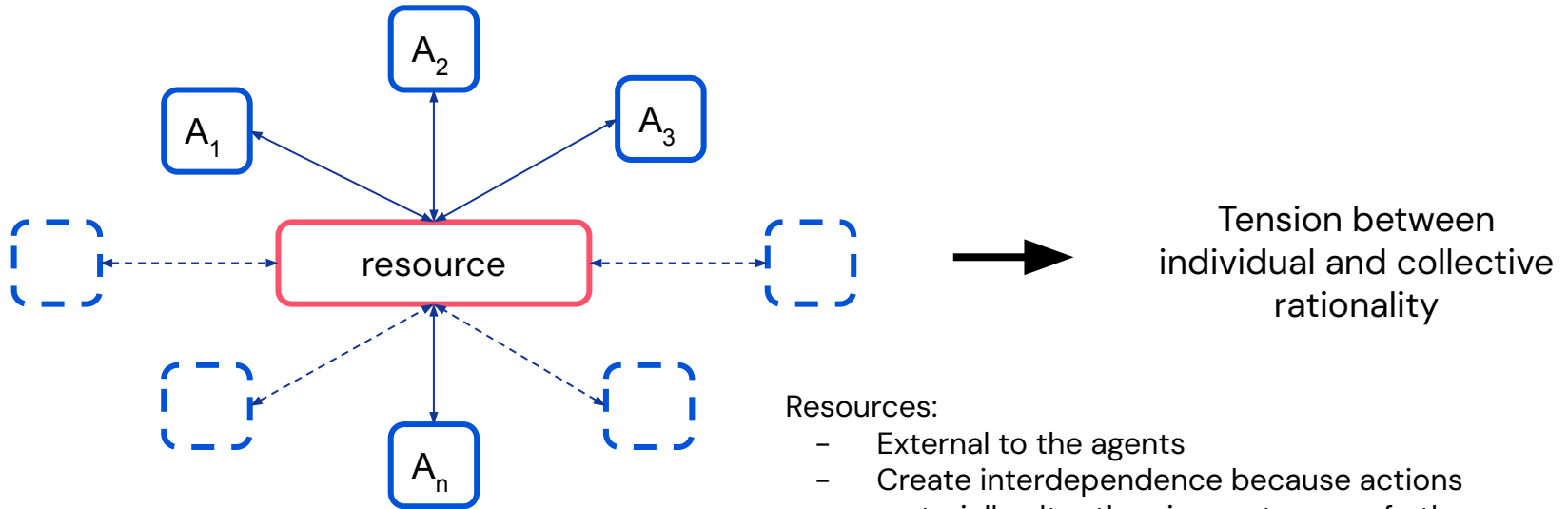


Player 2

|  | $A$ | $B$ | $C$ |
|---|---|---|---|
| $A$ | $(1,1)$ | $(0,0)$ | $(0,0)$ |
| $B$ | $(0,0)$ | $(2,2)$ | $(0,0)$ |
| $C$ | $(0,0)$ | $(0,0)$ | $(3,3)$ |

Player 1

# Agents made interdependent by a resource

Shared resources generate interdependencies of the agents: **misalignments**



Tension between individual and collective rationality

Resources:
- – External to the agents
- – Create interdependence because actions materially alter the circumstances of others
- – Natural resources: irrigation, pasture, fishery
- – Tech context: space / time / user attention…

# There are no panaceas


ELINOR OSTROM
2009 Nobel Laureate
in Economic Sciences
Nobel medal © © The Nobel Foundation

- Different resources are really different from one another!
- They generate different kinds of interdependencies between agents.
- Biophysical / social context matters
- Different resources demand different institutions, i.e. different ways of organizing cooperation.

|  | **Excludable** | **Non-Excludable** |
|---|---|---|
| **Subtractable** | <u>Private goods</u><br>food, clothing, cars | <u>Common-pool resources</u><br>fisheries, forests, pastures, calendar time, civil attention, public parking space |
| **Non-Subtractable** | <u>Club goods</u><br>cinemas, private parks, toll roads, Netflix | <u>Public goods</u><br>clean air, street lighting, national defense, police and fire departments, free software |

# Kinds of resources

|  | Exclusion feasible | Exclusion not feasible |
|---|---|---|
| **Subtractable** | private goods | common-pool resources |
| **Non-subtractable** | club goods | public goods |

- **subtractable** good: consumption by one agent reduces the amount available for consumption by others

- **excludable** good: it is possible to exclude another agent from accessing it
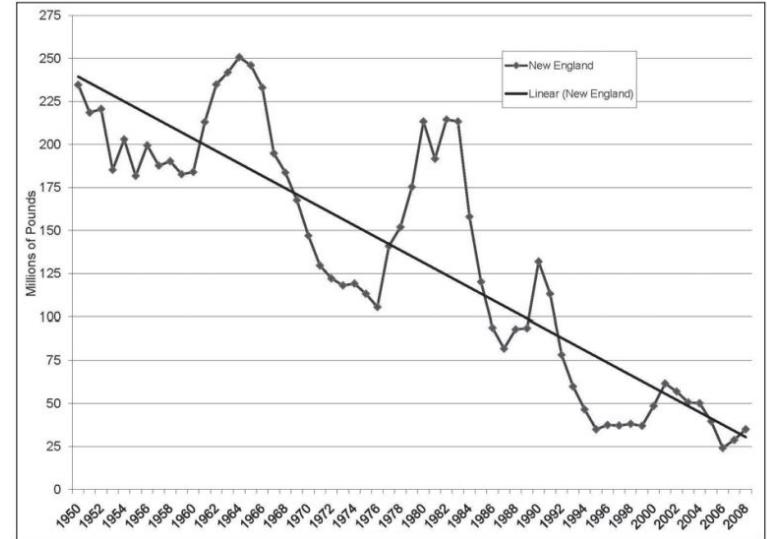
# Kinds of resources

|  | Exclusion feasible | Social dilemmas typically arise<br><br>Exclusion not feasible |
|---|---|---|
| **Subtractable** | private goods | common-pool resources |
| **Non-subtractable** | club goods | public goods |

- **subtractable** good: consumption by one agent reduces the amount available for consumption by others

- **excludable** good: it is possible to exclude another agent from accessing it

# Kinds of resources



| | Existence depends on complex institutions **Exclusion feasible** | Social dilemmas typically arise **Exclusion not feasible** |
|---|---|---|
| **Subtractable** | private goods | common-pool resources |
| **Non-subtractable** | club goods | public goods |

- **subtractable** good: consumption by one agent reduces the amount available for consumption by others

- **excludable** good: it is possible to exclude another agent from accessing it

# Managing natural resources is a prisoner's dilemma so we are all doomed

Acheson et al. (2011) **Coming up empty: Management Failure of the New England Groundfishery**.

"The root problem with groundfish management is that the management agencies and the industry have not been able to devise enforceable rules that effectively conserve fish stocks"

"If one can judge from interviews, fishermen hated the management plan and did everything in their power to undermine it. Every action of the council, nmfs, or the Secretary of Commerce was greeted with vicious letters to newspapers, vociferous complaints to congresspersons, petitions to officials, tire slashings, and heated hearings. [...] Perhaps even worse, the enforcement system broke down completely. **Cheating was rampant**. When a fishery was closed, many fishermen kept on fishing. Sometimes they sold the forbidden fish to Canadian vessels; in other cases they landed the fish in the u.s., but filed false reports listing it as another species. In many cases fish were simply discarded at sea. Most of those who cheated were not caught, which motivated others to cheat as well"



Figure 1: Catches of Cod, Haddock, and Yellowtail Flounder in New England, 1950-2008 (millions of pounds)
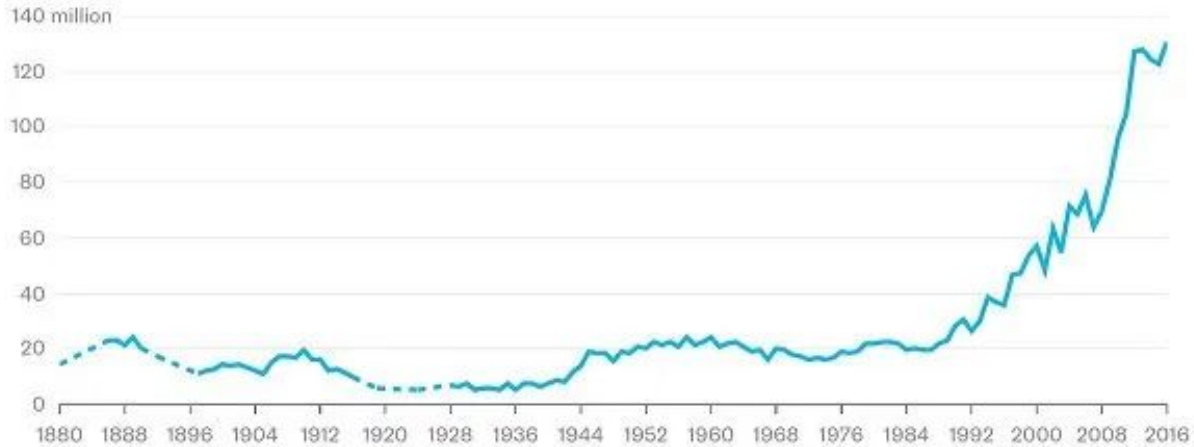
*Source:* Chart prepared by Ann Acheson, landings information generated from www.st.nmfs.noaa.gov/st1/commercial/landings/annual_landings.html

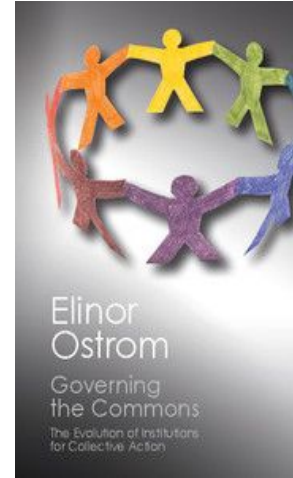# But not all natural resources are overexploited.
# How is this possible?



**The Lobster Boom**
Maine annual lobster harvest, in pounds

Source: Maine Department of Marine Resources

**But not all natural resources are overexploited.**
**How is this possible?**



**Ostrom's law:**

"A resource arrangement that works in practice can work in theory."

# Irrigation

- **Asymmetry**
  - Farmers near the water's source and farmers near its sink have **misaligned** incentives.
- **Appropriation**: how much water to use?
- **Provision**: how to arrange for maintenance of infrastructure?

Hoogesteger, J., Konijnenberg, V., Brackel, L., Kemink, S., Kusters, M., Meester, B., Mehta, A.S., van der Poel, M., Van Ommen, P., Boelens, R. and Sanchis-Ibor, C., 2023. Imaginaries and the commons: insights from irrigation modernization in Valencia, Spain. *International Journal of the Commons*, *17*(1).

# How Lobstering works





ESCAPE VENT

BAIT BAG

ENTRANCE HEAD

Useful rules for sustainable lobstering
1. No alternative technologies like dragging
2. Double gauge size limits
3. No "berried" females
4. V-notch rule
5. Escape vent
6. Trap limits



(Acheson 2003)

# Spatial Strategies and Territoriality (Acheson 2003)



- Harbor gangs
- Multi-tiered collective action problems:
  - Trap limits and other sustainability-related practices
  - Boundary defense (trap cutting etc)
  - Limited entry rules
- Nucleated versus perimeter-defended areas



Monhegan Island

# Explaining success and failure in the commons: Ostrom's design principles

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

Elinor Ostrom

(Cox et al. 2010)

# Explaining success and failure in the commons: Ostrom's design principles

1. **Well–defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective–choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict–resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

Elinor Ostrom

**1A: User boundaries**

Clear boundaries between legitimate users and nonusers must be clearly defined.

# Explaining success and failure in the commons: Ostrom's design principles

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

Elinor Ostrom

**1B: Resource boundaries**

Clear boundaries are present that define a resource system and separate it from the larger biophysical environment.

# Explaining success and failure in the commons: Ostrom's design principles

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

Elinor Ostrom

**2A: Local congruence**

Appropriation and provision rules are congruent with local social and environmental conditions.

# Explaining success and failure in the commons: Ostrom's design principles


Elinor Ostrom

1. **Well-defined boundaries**
    A. **User boundaries**
    B. **Resource boundaries**
2. **Congruence with local conditions**
    A. **Local congruence**
    B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
    A. **Monitoring of users**
    B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

**2B: Investment/Extraction proportionality**

The benefits obtained by users from a common-pool resource as determined by appropriation rules, are proportional to the amount of inputs required in the form of labor, material, or money, as determined by provision rules.

(Cox et al. 2010)

# Explaining success and failure in the commons: Ostrom's design principles

1. **Well-defined boundaries**
   - A. **User boundaries**
   - B. **Resource boundaries**
2. **Congruence with local conditions**
   - A. **Local congruence**
   - B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   - A. **Monitoring of users**
   - B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

Elinor Ostrom

**3: Collective-choice arrangements**

Most individuals affected by the operational rules can participate in modifying the operations rules.

# Explaining success and failure in the commons: Ostrom's design principles
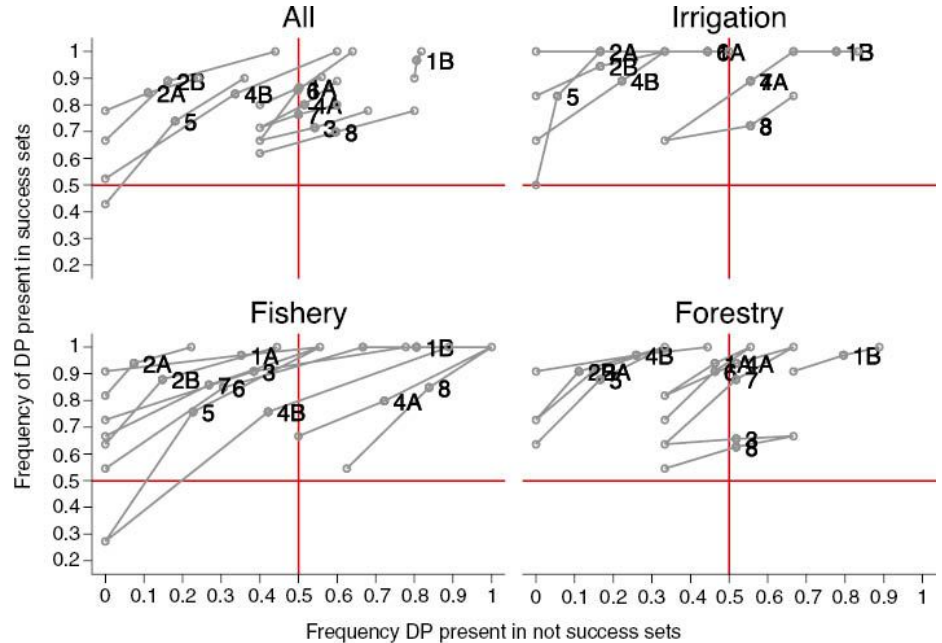

Elinor Ostrom

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

**4A: Monitoring of users**

Monitors who are accountable to the users monitor the appropriation and provision levels of users.

(Cox et al. 2010)

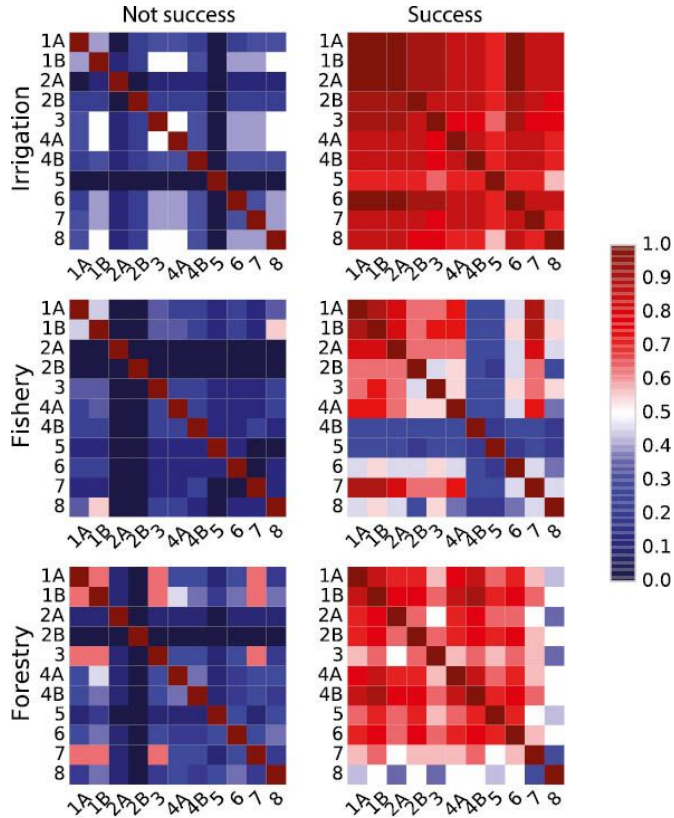# Explaining success and failure in the commons: Ostrom's design principles

1. **Well-defined boundaries**
    A. **User boundaries**
    B. **Resource boundaries**
2. **Congruence with local conditions**
    A. **Local congruence**
    B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
    A. **Monitoring of users**
    B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

Elinor Ostrom

**4B: Monitoring the resource**

Monitors who are accountable to the users monitor the condition of the resource.

# Explaining success and failure in the commons: Ostrom's design principles


Elinor Ostrom

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

**5: Graduated sanctions**

Appropriators who violate operational rules are likely to be assessed graduated sanctions (depending on the seriousness and context of the offense) by other appropriators, by officials accountable to them, or both.

# Explaining success and failure in the commons: Ostrom's design principles

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

(Cox et al. 2010)

Elinor Ostrom

**6: Conflict-resolution mechanisms**

Appropriators and their officials have rapid access to low-cost local arenas to resolve conflicts among appropriators or between appropriators and officials.

# Explaining success and failure in the commons: Ostrom's design principles



Elinor Ostrom

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

**7: Minimum recognition of rights to organize**

The rights of appropriators to devise their own institutions are not challenged by external governmental authorities.

(Cox et al. 2010)

# Explaining success and failure in the commons: Ostrom's design principles

Elinor Ostrom

1. **Well-defined boundaries**
   A. **User boundaries**
   B. **Resource boundaries**
2. **Congruence with local conditions**
   A. **Local congruence**
   B. **Investment/Extraction proportionality**
3. **Collective-choice arrangements**
4. **Monitoring**
   A. **Monitoring of users**
   B. **Monitoring the resource**
5. **Graduated sanctions**
6. **Conflict-resolution mechanisms**
7. **Minimum recognition of rights to organize**
8. **Nested enterprises**

**8: Nested enterprises**

Appropriation, provision, monitoring, enforcement, conflict resolution, and governance activities are organized in multiple layers of nested enterprises.

(Cox et al. 2010)

# Explaining success and failure in the commons: Ostrom's design principles



- Clearly defined boundaries is principle 1: A is social, B is biophysical.
- Congruence between rules and local conditions is principle 2A.

(Baggio et al. 2016)

# "There are no panaceas" – Elinor Ostrom

Is it
- provision or appropriation?
- how hard are monitoring and sanctioning?

**Diverse agents (human + AI)**

**Resource**

**Specific misalignment**

- Do the people involved know each other?
- Do they speak the same language?

**Governance arrangement**

Is a resource governance arrangement
- sustainable?
- acceptable to stakeholders?
- does it violate human rights?
- is it pareto-optimal?
- is it fair?

- "Let's take turns"
- "Let's decide X is illegal"
- "Let's impose a quota"
- "Let's write a constitution"

# Generative AI and the digital commons

Huang, S. and Siddarth, D., 2023. **Generative AI and the digital commons.** *arXiv preprint arXiv:2303.11074*.

1. The digital commons includes wikis, internet archive snapshots, Creative Commons (CC) licensed images, and public software repositories, coding question and answer sites, etc.
2. Generative AI is trained on the digital commons.
3. Generative AI can put material back into the digital commons
   a. much faster than humans
   b. of unclear quality
      i. confabulation/hallucination
      ii. bias
4. Generative AI could disincentivize humans from contributing to the digital commons

5. Global in scope
6. Monitoring? Sanctioning?

# What is "alignment"?

Alignment is what answers the fundamental political question:

**How can we live together?**

→ **"Institutions"**

How can we build human-compatible AI?

An AI system is individual-compatible if it aligns with the preferences of a single person

**Example:** Content Moderation



Individual-compatible

# Does pluralistic AI solve this problem?



Pluralistic Alignment Sorenson et. al. 2024

An AI system is social-choice compatible if it aligns to an aggregation of a population's preferences

**Example:** Content Moderation



Individual-compatible  Social choice-compatible

# An AI system is social-choice compatible if it aligns to an aggregation of a population's preferences

**Collective Constitutional AI: Aligning a Language Model with Public Input**

Saffron Huang[*][†]
saffron@cip.org
Collective Intelligence Project
San Francisco, California, USA

Divya Siddarth[*]
divya@cip.org
Collective Intelligence Project
San Francisco, California, USA

Liane Lovitt[*]
Anthropic
San Francisco, California, USA

Thomas I. Liao[‡]

Esin Durmus

Alex Tamkin

Deliberative Democracy
Platform — Polis

**Stakeholders**          **Principles**          **Optimize**

## STELA: a community-centred approach to norm elicitation for AI alignment

Stevie Bergman, Nahema Marchal[✉], John Mellor, Shakir Mohamed, Iason Gabriel & William Isaac

Value alignment, the process of ensuring that artificial intelligence (AI) systems are aligned with human values and goals, is a critical issue in AI research. Existing scholarship has mainly studied *how* to encode moral values into agents to guide their behaviour. Less attention has been given to the normative questions of *whose* values and norms AI systems should be aligned with, and *how* these choices should be made. To tackle these questions, this paper presents the STELA process (SocioTEchnical Language agent Alignment), a methodology resting on sociotechnical traditions of participatory, inclusive, and community-centred processes. For STELA, we conduct a series of deliberative discussions with four historically underrepresented groups in the United States in order

|  | Overton | Steerable | Distributional |
| --- | --- | --- | --- |
| Response | … prevent harmful content while … violation of free speech … First Amendment and private platform … and … could potentially backfire… | … a positive step in preventing the spread of offensive content: hosting toxic discussions and removing this platform reduces … the spread of hate speech… | overall |
| (Black-Box) LLM | Overton ⊗ Q: Is taking down 4chan justified? | Steerable ↕↑↓ Q: Is taking down 4chan justified? | Distributional ◉ $P(\cdot \mid q, \square)$ |
| Comments |  |  |  |
| Community LMs |  |  |  |
|  | Query | Query | Query |

An AI system is social-choice compatible if it aligns to an aggregation of a population's preferences.

Does social-choice compatibility provide compatibility with institutions?

**No!**

Human society's procedures for stable coordination through collective decision making over time are complex and nuanced.

Preference aggregators represent a simplified version of such collective decision making process

An AI system is social-choice compatible if it aligns to an aggregation of a population's preferences

**Example:** Content Moderation

**Misses the mark!**



Individual-compatible

Social choice-compatible

Representing complex collective decision making process

Legitimacy

Our Approach: This tutorial

Makes Collective Decisions
(General, Scalable)

**Example:** Content Moderation

**Legitimate** Process

Individual-compatible

Social choice-compatible

Institution-compatible

**Institution-compatible AI** = AI systems that optimizes for
**institutional decisions** as their alignment target.

# How to build AI that engages with Institutions

Markov Game: $\langle S, s_0, \mathbf{A} = \times_n A_n, T : S \times \mathbf{A} \to \Delta_S, R : S \to \mathbb{R}^N \rangle$

### *Contracts*



Human contracts rely on *tons* of structure

- e.g. "what was it reasonable to think the parties had in mind when they agreed"

- "reasonable" (and other gap-fillers) provided by institutions (norms, law)

# Contracting in Multi-Agent RL

$$\langle S, s_0, \mathbf{A} = \times_n A_n, T : S \times \mathbf{A} \rightarrow \Delta_S, R : S \rightarrow \mathbb{R}^N \rangle$$



$S' : \{S \times \Theta\} \cup \{\emptyset\}$

$c$: **0-sum contract function**
$$0 = \sum_i c_i(s, \theta)$$

$$\langle S', s_0', \mathbf{A}' = \times_n A_n', T : S \times A \rightarrow \Delta_S, R' : S' \rightarrow \mathbb{R}^N, \Theta, c : S \times \Theta \rightarrow \mathbb{R}^N \rangle$$

$R_i'(s, \theta) = R_i(s) + c_i(s, \theta)$

$\Theta$: **contract parameters**

Phillip Christofferson

Andreas Haupt

Christoffersen et al., AAMAS 2022

# Contracting in Multi-Agent RL



Christoffersen et al., AAMAS 2022

# Contracting in Multi-Agent RL



Christoffersen et al., AAMAS 2022

# Domain 1: Common Pool Resource

- Standard Common Pool Resource Domain

- 1-dimensional state — amount of "fish" in the lake

- 1-dimensional action — amount of fish to try to take

- Transition dynamics — regrowth rate, subject to 'overfishing'

- What's a good contract space?

This work: taxes on desired fishing levels

# Harvest



What is a good contract space for this domain?

# Harvest



B

fining beam — apple

agent — agent observation window

What is a good contract space for this domain?

Penalty for eating from low-density regions

| # agents | Pris. Dilemma | Public Goods | Harvest | Cleanup | Merge |
|----------|---------------|--------------|---------|---------|-------|
| 2 | | | | | |
| 4 | | | | | |
| 8 | | | | | |

Legend: Joint, Separate, Gifting, Vanilla Contracting, MOCA

# Spurious normativity improves the capacity of agents to learn enforcement and compliance behaviors



Koster et al., PNAS 2022

# Spurious normativity improves the capacity of agents to learn enforcement and compliance behaviors



No Rules

Berry

Poison Berry

# Normative conditions

Koster et al., PNAS 2022

# Spurious normativity improves the capacity of agents to learn enforcement and compliance behaviors



## Normative conditions

Koster et al., PNAS 2022

# Spurious normativity improves the capacity of agents to learn enforcement and compliance behaviors



# Normative conditions

Koster et al., PNAS 2022

# Research questions

1. Do agents learn to punish?
2. Do agents learn to avoid punishment (comply with the rules)?
3. Does a stable state with normative infrastructure emerge?
4. How does the presence of a silly rule affect learning?
5. Does normative infrastructure raise payoffs?

Figure with six panels:

**A. Total Misdirected Punishing**

**B. Total Times Punished**

**C. Mean Time Marked**

**D. Total Time Since Poisoning**

**E. Total Taboo Berries Eaten**

**F. Collective Return**

Legend for each panel:
- R{}
- R{poisonous}
- R{poisonous, nonpoisonous}

X-axis label: Steps (1e9)

# Insights

Normative behaviors support better choices

Silly rules support learning of normative behaviors—enforcement and compliance

Game theoretic approaches to predicting/explaining individual rules will not capture this phenomenon

Can presence of a **normative institution** help
agents learn normative behaviors?

# Altared Environments



- Introduces *altar* as a feature of the environment
- Altar content maps to the hidden rule (unknown to agents) i.e. altar prescribes what is acceptable or not

(classification institution)

- Agents need to visit the altar to receive the content
- Altar content is *dynamic*

# Altared Environments



**Commons harvest**

7 agents collect apples across 6 patches

Reward +1 per apple collected

Apple regeneration depends on #apples within a distance of 2.

Patches can be permanently depleted

# Altared Environments



**Punishing**

Punishing costs sanctioning agent -10

Punished agent gets removed for 25 steps

Punishing agent receives a reward of +20 (net positive of +10) if punished agent was harvesting from patch classified as "inappropriate"

# Altared Environments



Punishment rewards (classification) hidden

Same classification and rewards

Classification represented on 'altars'

# Altared Environments



(a) Eat from blue zone    (b) Eat from red zone    (c) Eat from green zone    (d) All zones prohibited

Top left/right: red zone; Middle left/right: blue zone; Bottom left/right: green zone

Punishment rewards earned for harvesting from zone other than one classified as "appropriate" (chosen with probability proportional to quantity in zone; changed when apple count in zone falls below threshold)

Altar changes color with classification

If all zones are depleted enough, altar turned yellow fire indicating no harvesting to be done for that period

# Experiments: Altar condition vs hidden rule



Significantly higher collective return

# Experiments: Altar condition vs hidden rule



Significantly higher
collective return

Significantly lower
depletion

# Experiments: Altar condition vs hidden rule



Significantly higher collective return

Significantly lower depletion

Agents learn to consistently  visit (consult) altar

# Experiments: Altar condition vs hidden rule



Significantly more correct sanctions (aligned with classification)

Fewer incorrect sanctions

Can presence of normative infrastructure help agents generalize aligned behavior in new environments?

Can we train agents to learn "punish what's on the altar" rather than "punish agents who eat [blue] berries"?

# Generalization



Poison berries                    Harvest

Can agents trained to punish behavior represented on altar generalize
to new environment with altar (same normative infrastructure) and
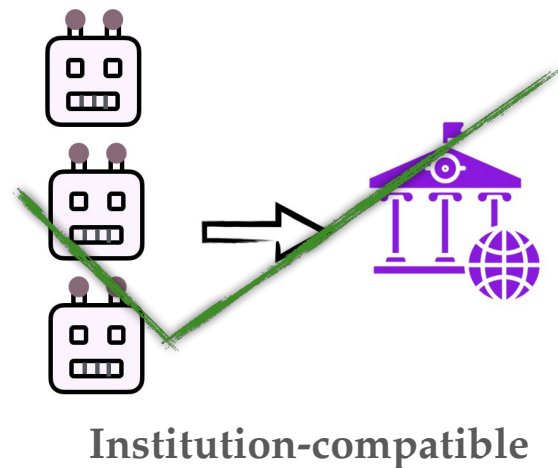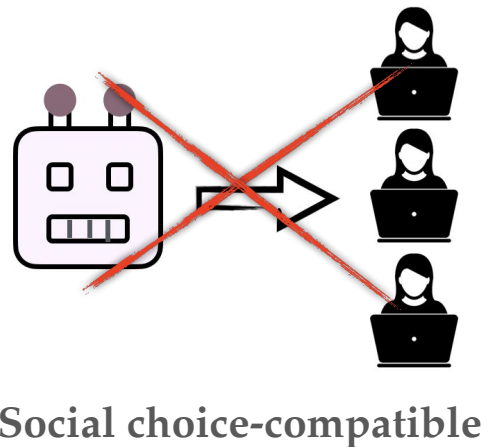reach alignment more reliably and/or faster?
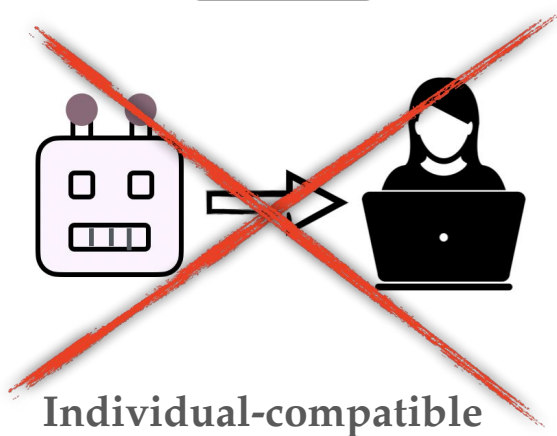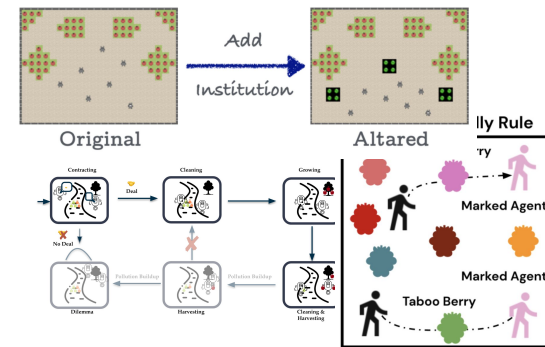
# What is "alignment"?



$r^*$   $\tilde{r}$

$$\max_{\vec{a}} \sum_t r(s_t)$$

**Individual-compatible**

# What is "alignment"?



**Individual-compatible**

**Social choice-compatible**

What is "alignment"?

Individual-compatible

Social choice-compatible

Institution-compatible

How can we live together?

# We are hiring!

Multiple open-positions at CS Department @ John Hopkins:

- Postdoctoral Fellow
- PhD Student

Please check the tutorial website or gillianhadfield.org for application details!

Student researcher position @ Google DeepMind in London

- Final year PhD student

Please check the tutorial website or Google Deepmind Careers Page!

# Resources

For references, slides and discussion notes, visit:

https://alignment-tutorial-2024.github.io/

Please submit your questions on the website and reach out to us during the week!

THANK YOU!