# Altruistic Punishment and Human Cooperation: A Darwinian Perspective
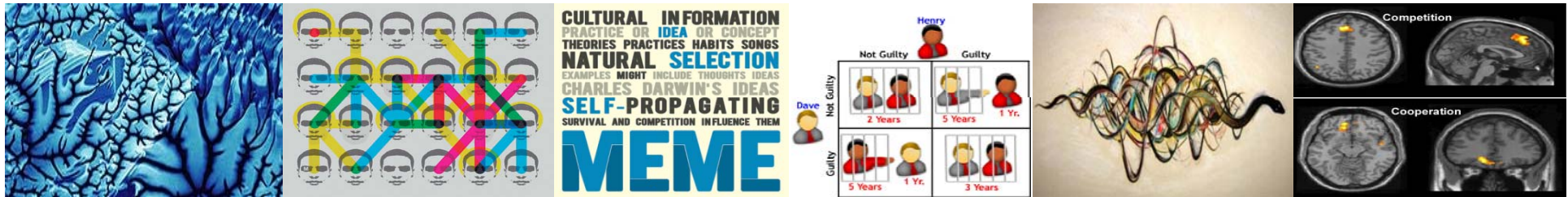
**Moritz Hetzer and Prof. Didier Sornette**

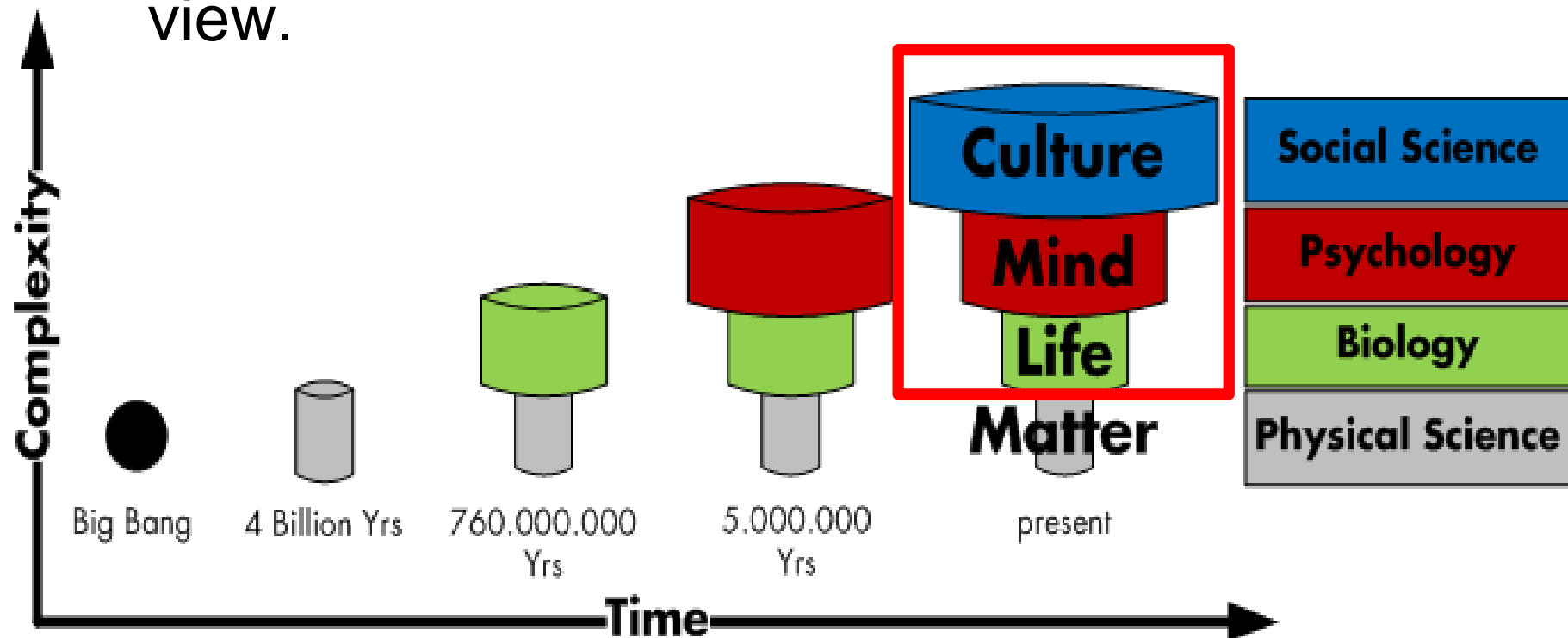**Chair of Entrepreneurial Risks**

**ETH Zurich**

# Motivation

- ## Questions we want to answer:

  - Why do people altruistically punish defectors?

  - What is the role of fairness perception and other-regarding preferences in this context?

  - How does punishment affect the emergence and maintenance of cooperation?
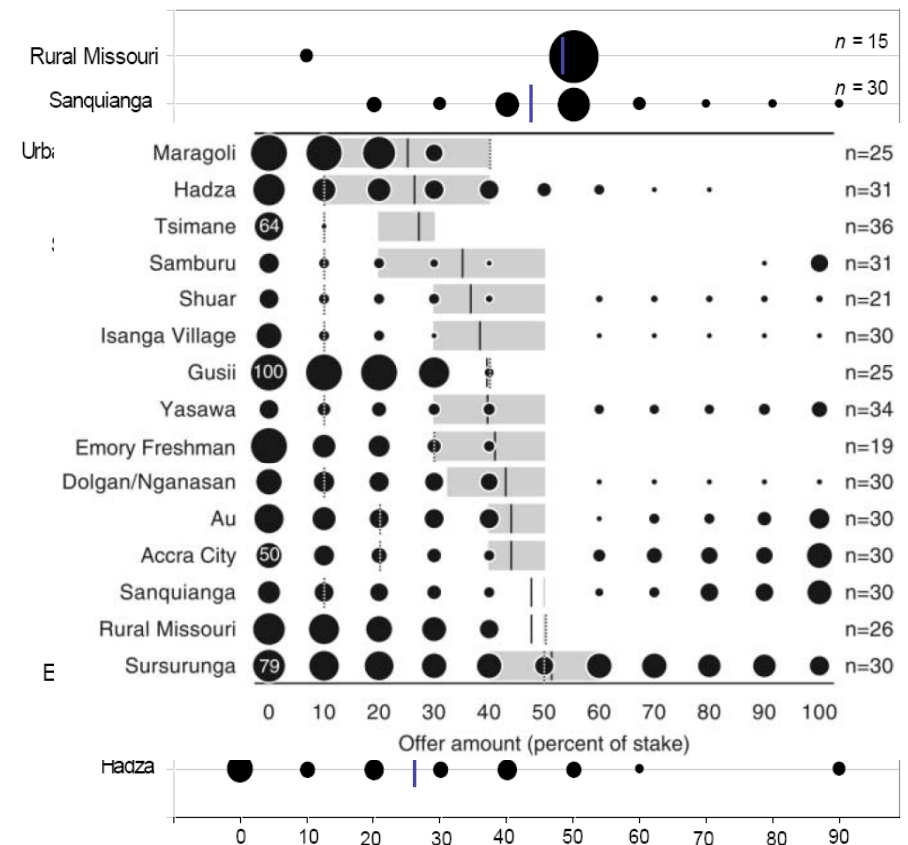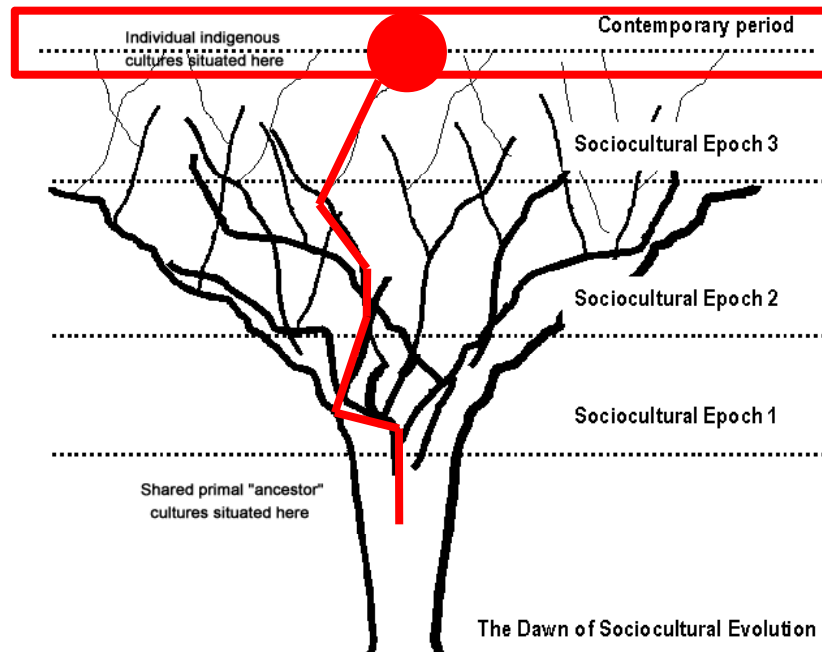
  - Why do we cooperate?

# Motivation: The evolution of norms/genes

- We want to understand the roots of individual & collective behavior from an evolutionary point of view.
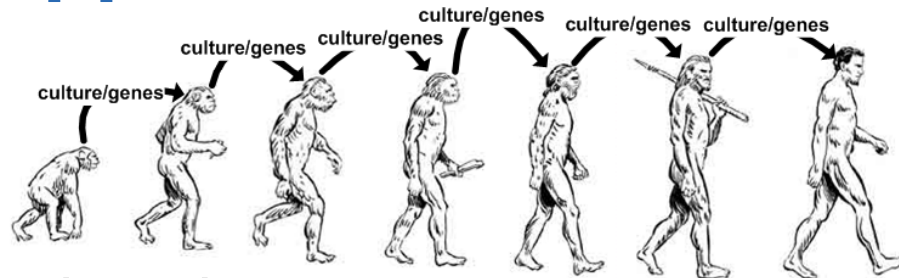
# Motivation: The evolution of norms

- Experiments identify behavioral patterns

- Economic theories describe these patterns
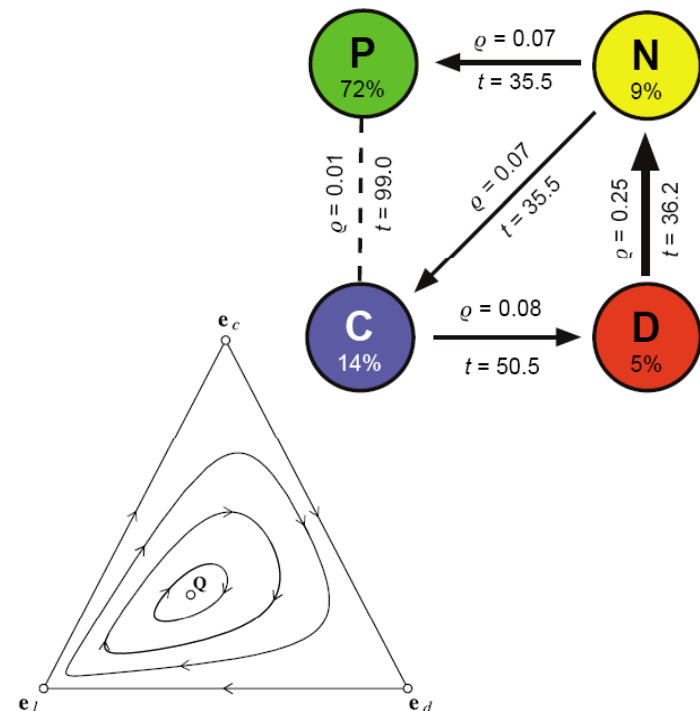
# Motivation - existing approaches

- ## Evolutionary theories

    - Kin selection

    - Direct / indirect / social reciprocity

    - gene-culture coevolution

- ## Analytic models
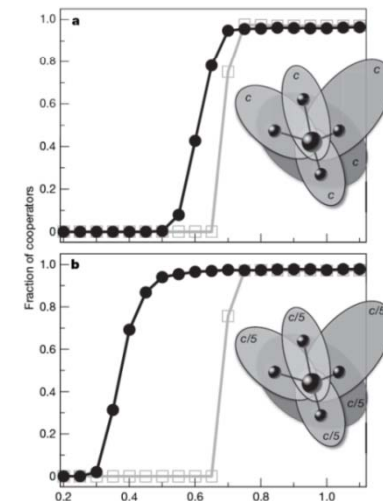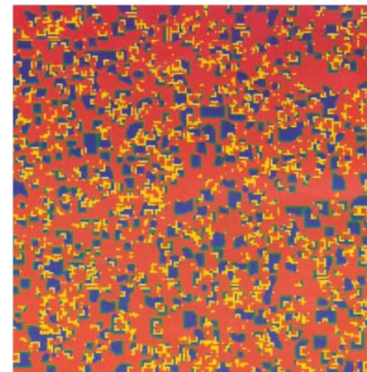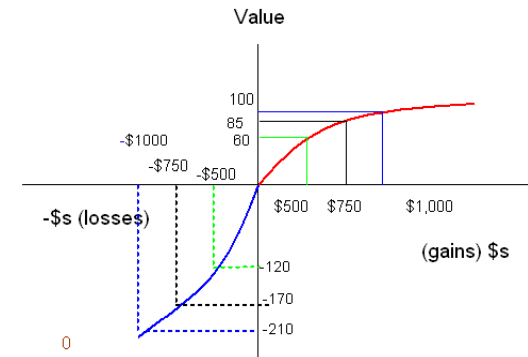
    - Mutual two-player interactions

    - Focus on equilibrium solutions

    - Detached from reality

    - Evolutionary game theory

        - better: Iterative Game Theory

# Motivation - existing approaches

- ## Economic theories

  - Descriptive

  - Snapshot of current norms

  - Do not cover evolutionary dynamics

$$U_i(x) = x_i - \alpha_i \max\left[x_j - x_i, 0\right] - \beta_i \max\left[x_i - x_j, 0\right], \quad i \neq j.$$



- ## Computer simulations

  - Sequential games

  - Lattice structure

  - Discrete decisions

  - Detached from reality

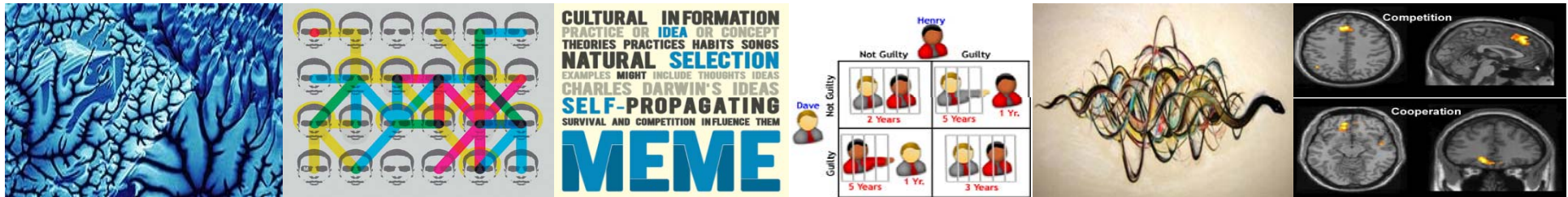  - Focus on equilibrium solutions
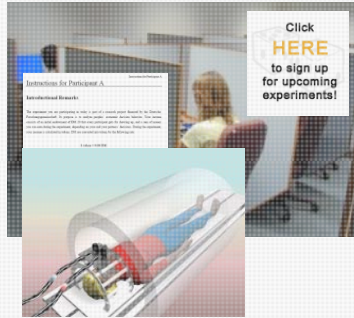
# Motivation - Our approach

- We want to answer the questions by closely integrating **experimental economics** with **agent-based modeling**.
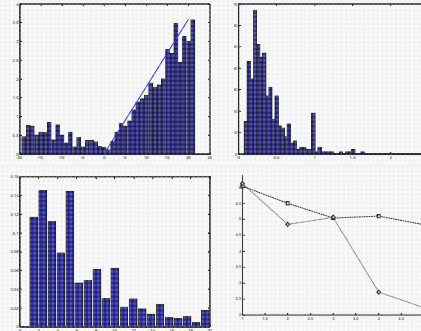
# Empirical foundation

- We use data from Fehr's & Gächter's public goods game experiments (2000/2002)

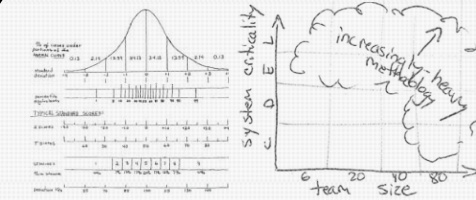# Other-Regarding preferences and altruistic punishment: A Darwinian Perspective
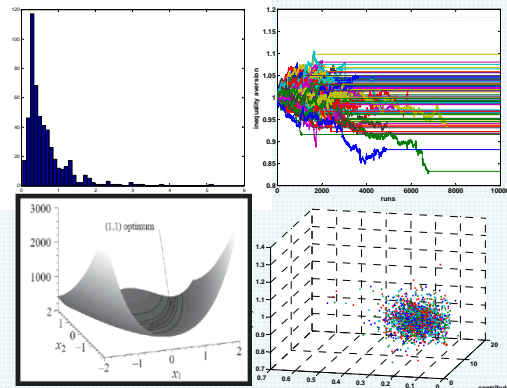
**1. collect data**

**2. Identify patterns**

**3. deduces generic norms/rules**

**6. verify results**

**5. grow artificial populations**

**4. design agent-model**

# Experiment: Public goods game

Punishment?

...compounding with 160%...

1. Each subjects decides to contribute to the group project.

2. The group project pool is compounded by a factor of 1.6

3. The project return is equally redistributed to all group members.

4. Each subject gets the opportunity to punish other group members at own costs, i.e. punishment is costly to both the punisher and the punished individual.

# Model Design:

- ## **Properties of agent** $i$ **:**

  - Level of cooperation $\quad m_i(t)$

  - Propensity to punish $\quad k_i(t)$

  - Wealth/Fitness $\quad w_i(t)$

# Model Design – one simulation period:

- **cooperate:** Each agent contributes $m_i$ to the group project
- **punish:** Punishment of other group members

# Model Design – empirical punishment:



$$p_{i \to j} = \begin{cases} k_i \cdot (m_i - m_j), \text{if } m_i > m_j \\ 0, \text{else} \end{cases}$$

# Model Design – one simulation period:

- **cooperate:** Each agent contributes $m_i$ to the group project

- **punish:** Punishment of other group members according to:

$$p_{i \to j} = \begin{cases} k_i \cdot (m_i - m_j), \text{if } m_i > m_j \\ 0, \text{else} \end{cases}$$

# Model Design – one simulation period:

- **cooperate:** Each agent contributes $m_i$ to the group project

- **punish:** Punishment of other group members

- **consume:** Consume avg. group welfare gained in period $t-1$

# Model Design – P/L, wealth and consumption:

- **Profit & Loss:**

$$s_i(t) = \underbrace{\frac{1.6}{4} \sum_{j \in I} m_j}_{\substack{\text{project} \\ \text{return}}} \underbrace{- m_i}_{\text{contribution}} \underbrace{- \sum_{j \in I} p_{i \to j}}_{\substack{\text{punishment} \\ \text{spent}}} \underbrace{- 3 \cdot \sum_{j \in I} p_{j \to i}}_{\substack{\text{punishment} \\ \text{received}}}$$

- **Wealth:**

$$W_i(t+1) = W_i(t) + s_i(t) - \underbrace{c(t)}_{\text{consumption}}$$

- **Consumption:**

$$c(t) = \overline{W}(t-1) - \overline{W}(t-2)$$

# Model Design – one simulation period:

- **adapt:** Change cooperation level $m_i$ and the propensity to punish $k_i$

# Model Design – Adaptation of $m_i$:

- Agents adapt their level of cooperation $m_i$ if:

    profit/loss < consumption

    with: $m_i(t+1) = m_i(t) + \varepsilon$

# Model Design – Adaptation of $k_i$ :

- **(A) Selfish agents:** Adapt their behavior if:

  profit/loss is less than her consumption.

- **(B) Inequality avers agents:** Adapt their behavior if:

  profit/loss < average group profit/loss (***downside)*** or
  profit/loss > average group profit/loss (***upside***).

- **(C) Inequity averse agents:** Adapt their behavior if:

  contribution > group average contribution **and**
  profit/loss < group's average profit/loss (***downside) or***
  contribution < group average contribution **and**
  profit/loss > group's average profit/loss (***upside).***

# Model Design – Adaptation of $k_i$ :

- **(D) Disadvantageous inequality avers agents:**

  Adapt their behavior if:

  profit/loss < average group profit/loss (***downside***)


- **(E) Disadvantageous inequity averse agents:**

  Adapt their behavior if:

  contribution > group average contribution **and**

  profit/loss < group's average profit/loss (***downside***)

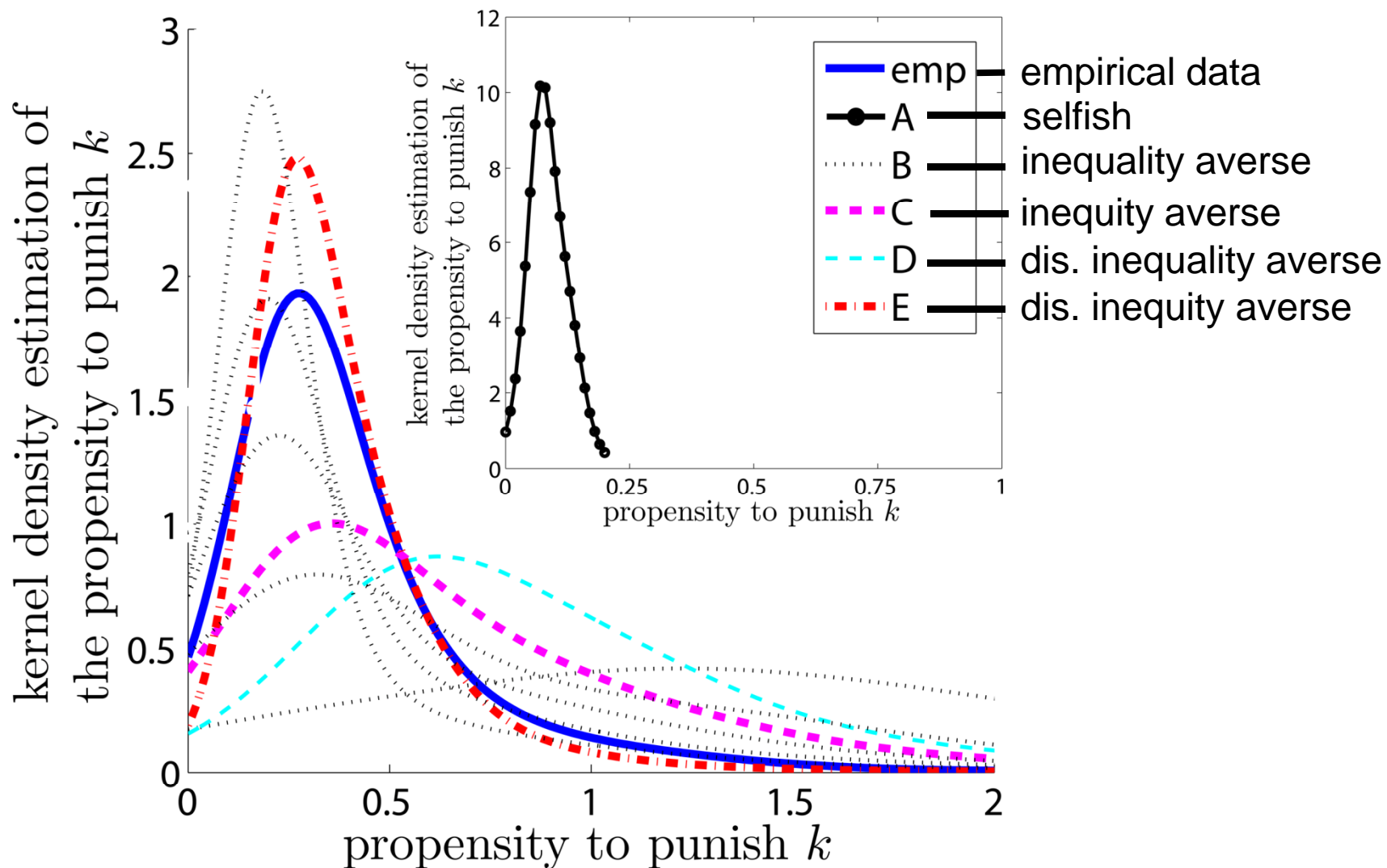# Model Design – one simulation period:

- **selection:** If the wealth of an agent drops below 0 the agent dies.

- **cross-over:** Dead agents are replaced with new ones. The level of cooperation $m_i$ and propensity to punish $k_i$ are initialized by the avg. values of the surviving population.
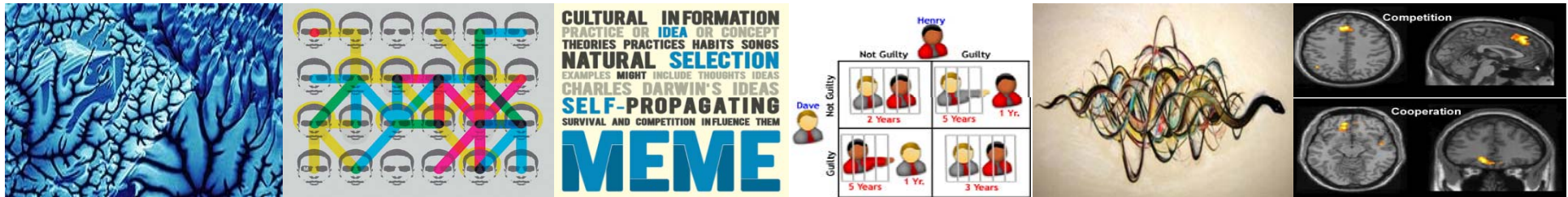
# Model Design – Simulation:

- **We run this model for 1 million simulation periods over 800 system realizations with**
  - $m(0)_i$
  - $k(0)_i$
  - $w(0)_i$

  **and obtained a distribution for $k_i$ which we compare with the empirical distribution obtained from experimental data.**

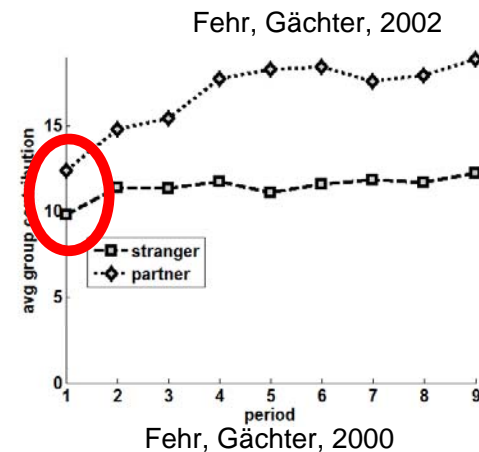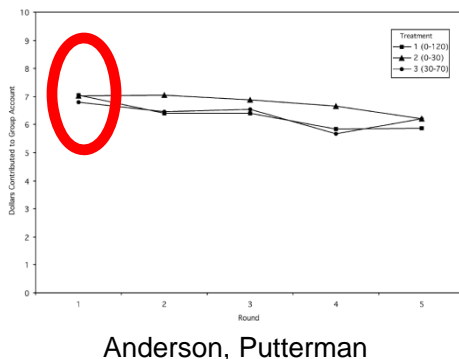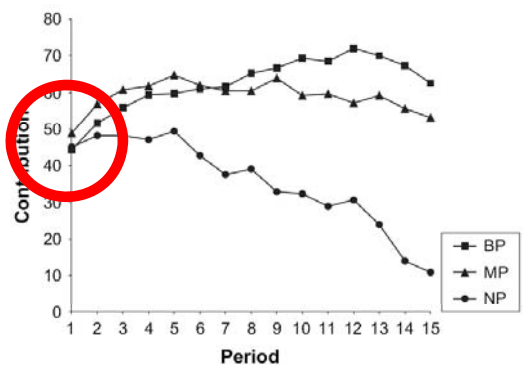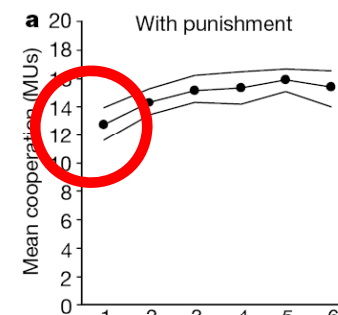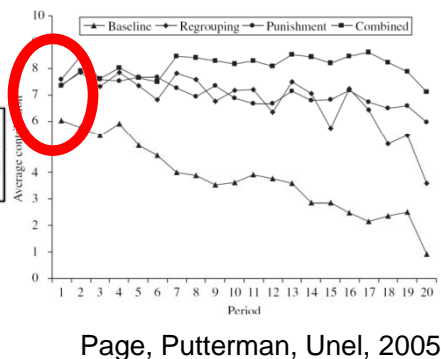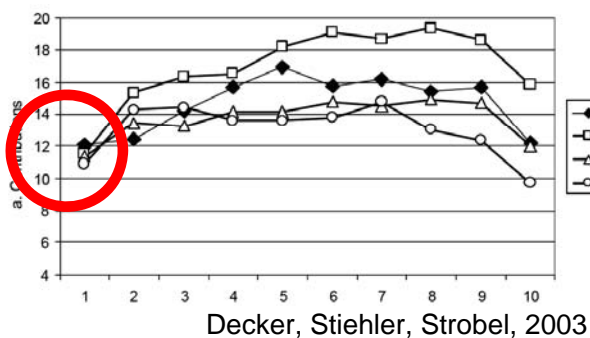# Disadvantageous inequity aversion fits best!

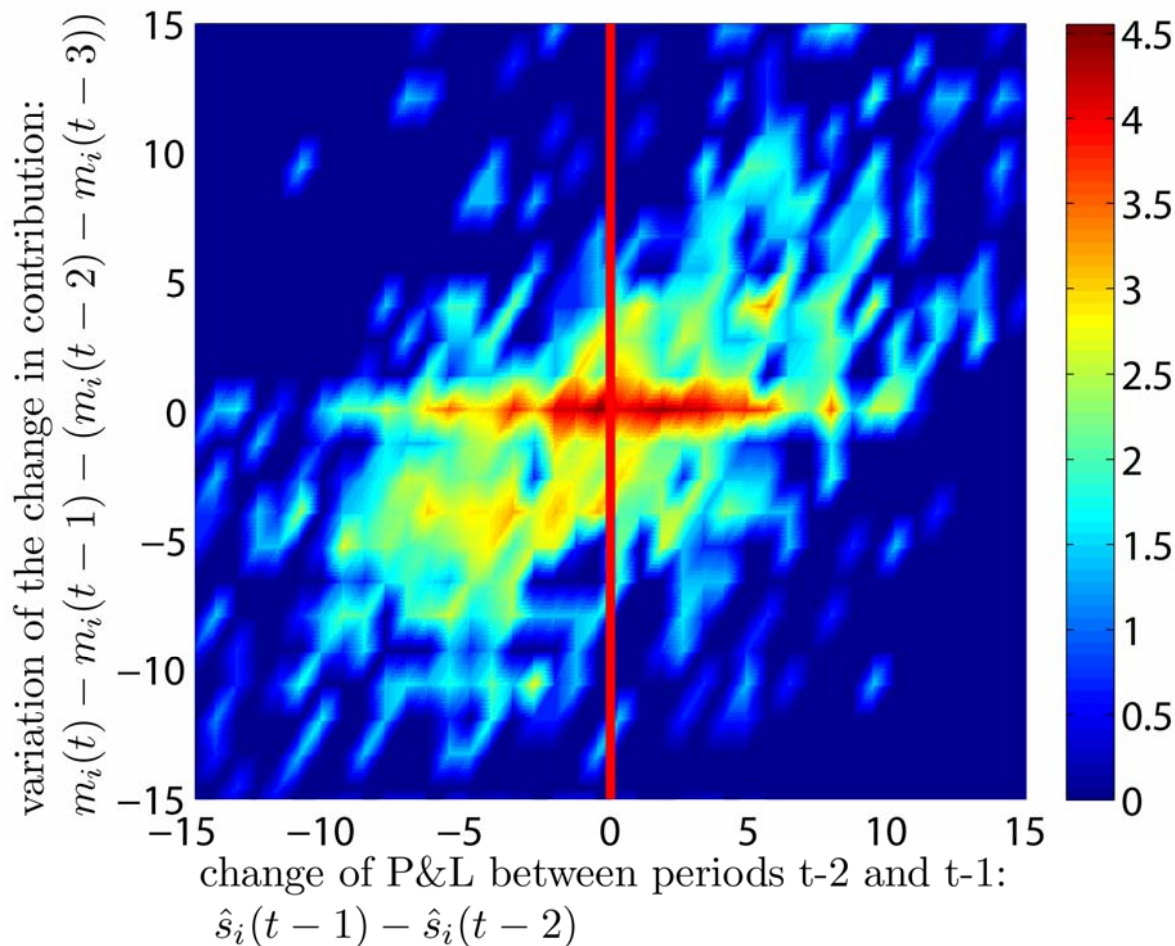# Altruistic Punishment and the Emergence of cooperation: A Darwinian Perspective

# The effect of punishment on cooperation

- (Altruistic) punishment is often used to explain the emergence of cooperation in social dilemmas.



Decker, Stiehler, Strobel, 2003

Page, Putterman, Unel, 2005

Fehr, Gächter, 2002

Noussair, Tucker, 2005

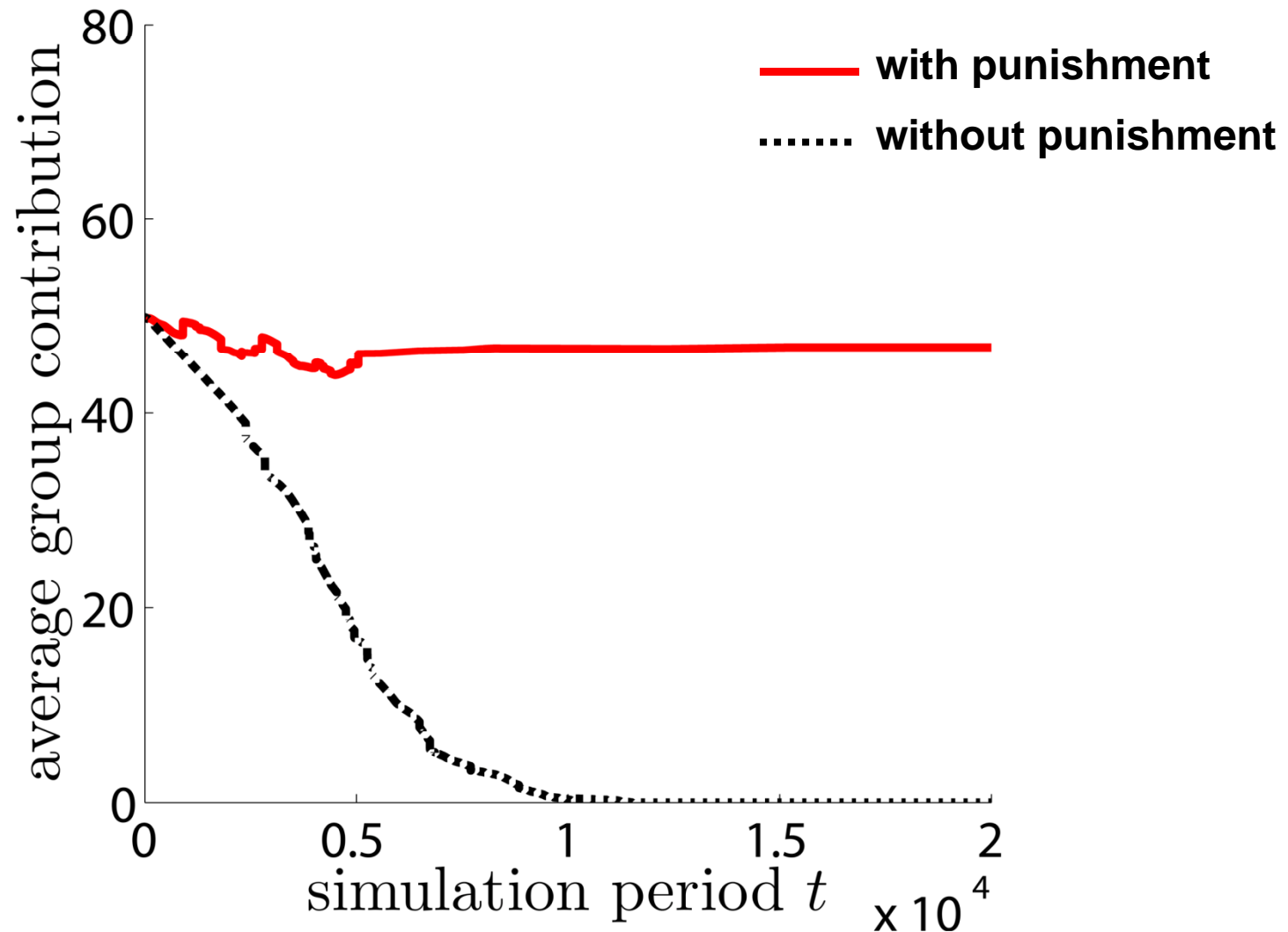Anderson, Putterman

Fehr, Gächter, 2000

# Evidence for short term persistence in period-by-period decision process:
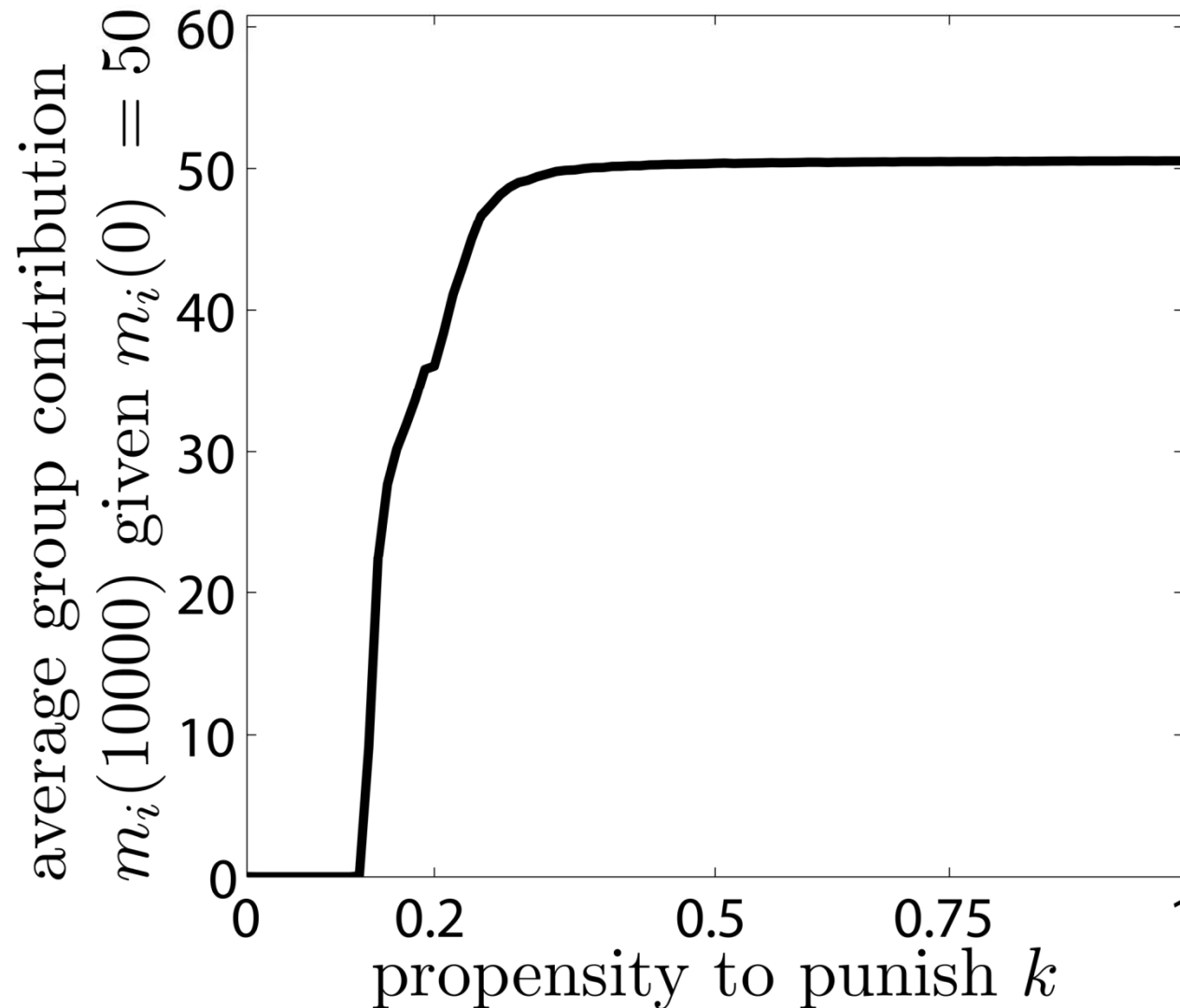


- Subjects seem to follow a trend in their updates of the individual contributions.

- If profit/loss in period (t) is larger than in period (t-1)

$$m_i(t+1) = 2 \cdot m_i(t) - m(t-1)$$

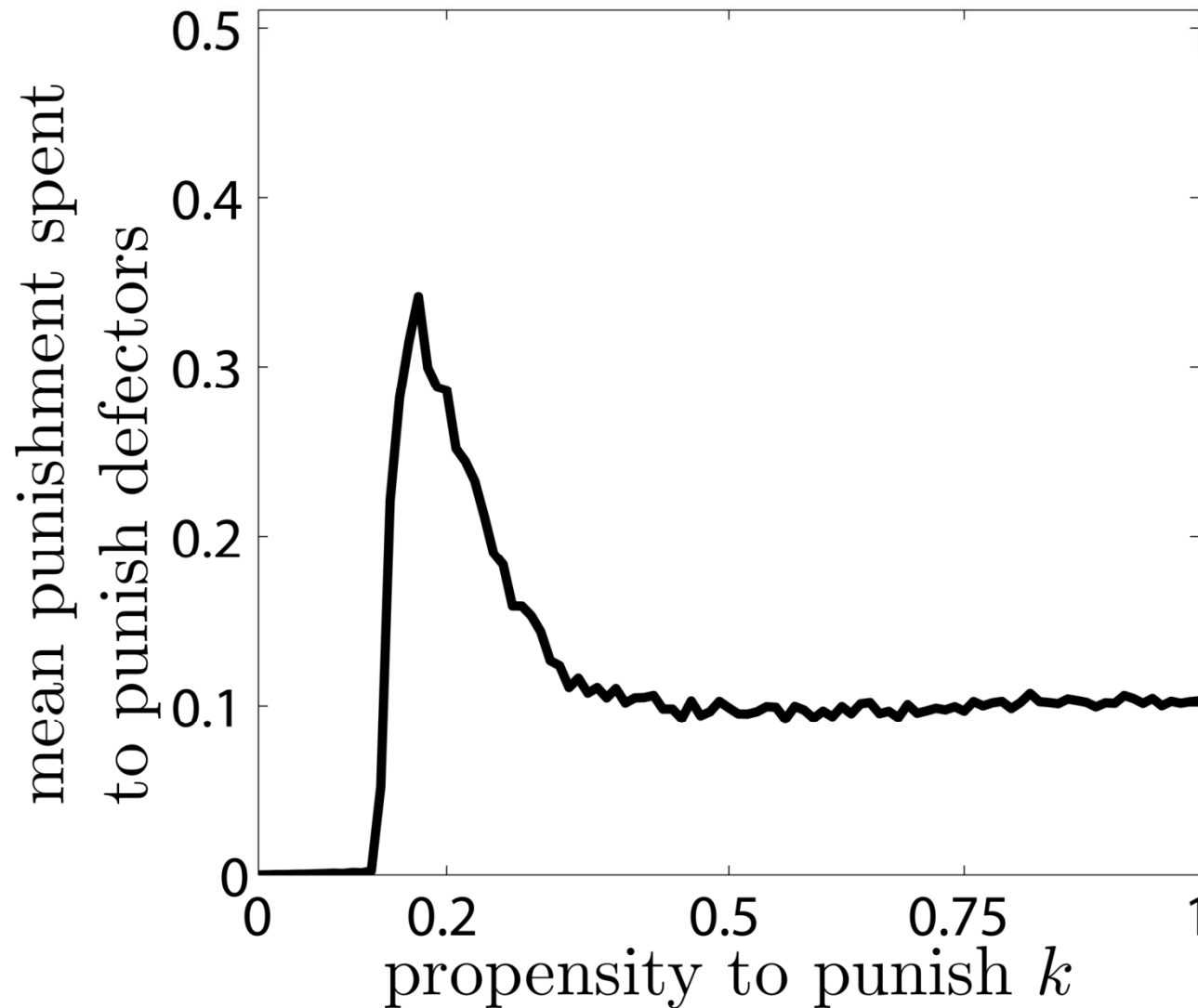- Previous results are **ROBUST** to this addition

# The effect of punishment on cooperation

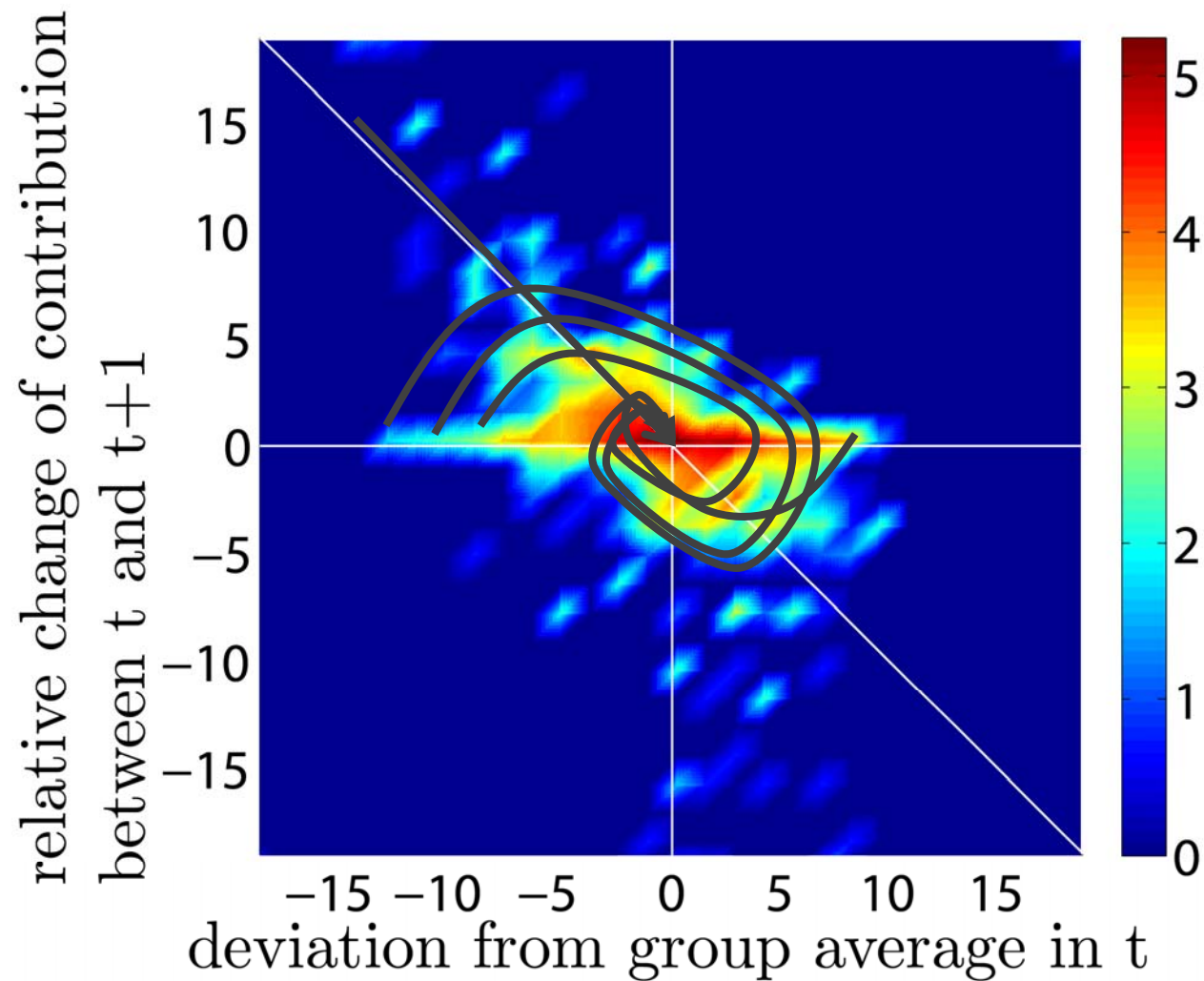# The effect of punishment on cooperation
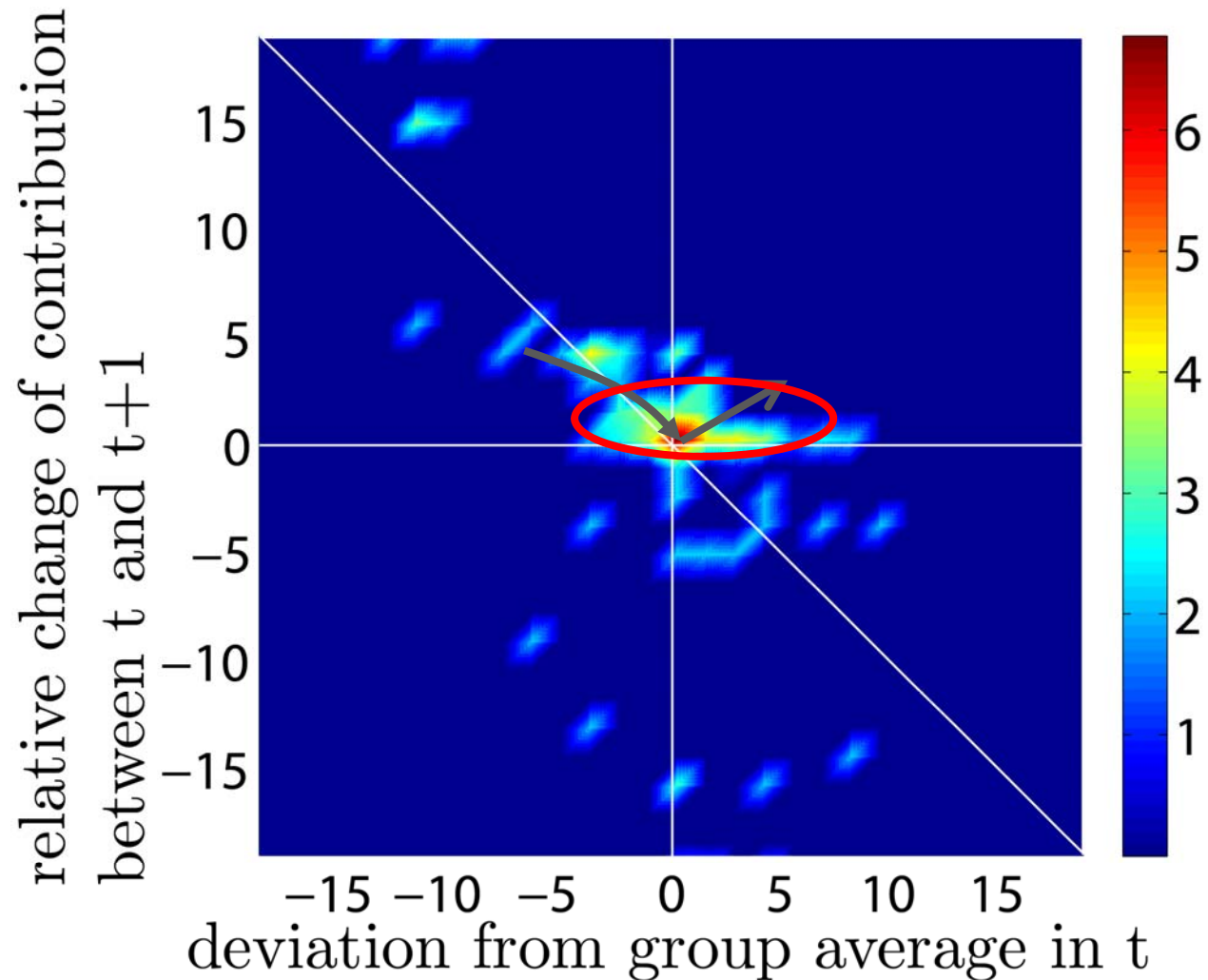
# The effect of deterrence

# Altruistic punishment and cooperation

- Is altruistic punishment sufficient to **sustain** cooperative behavior …

  between related (partners) and unrelated individuals (strangers)?

- Is altruistic punishment sufficient to **promote** cooperative behavior …

  - **Partners:** group composition stays constant
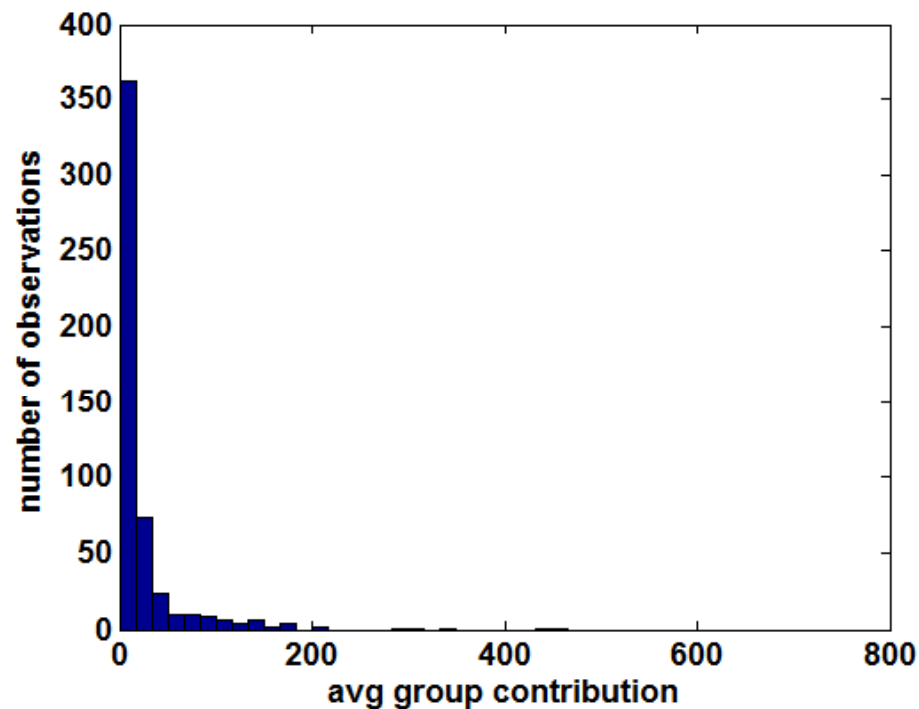  - **Strangers:** group composition changes
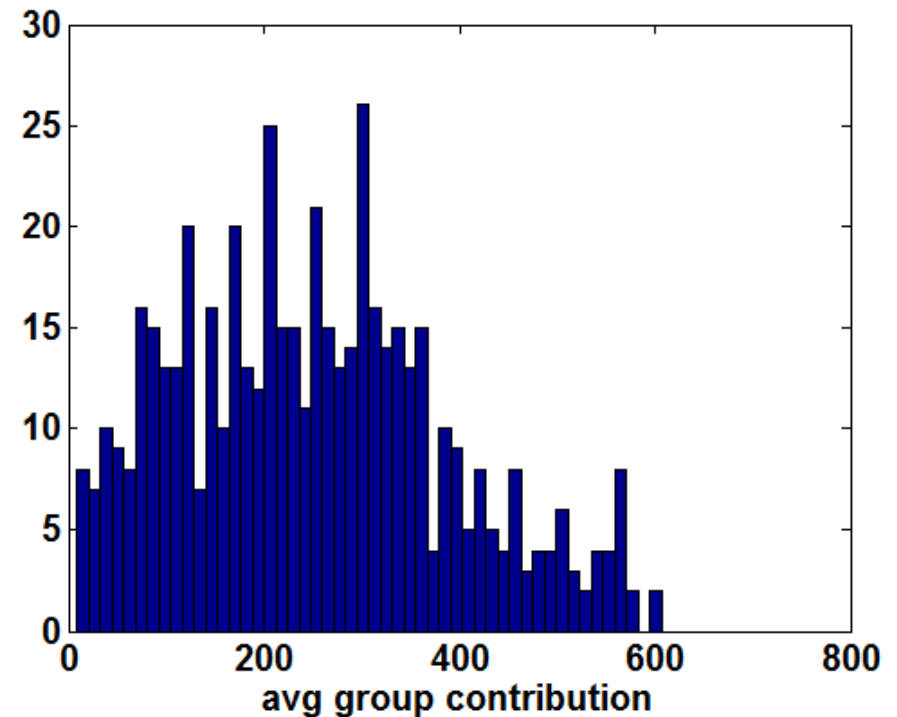
# First-order dynamics among strangers

# First-order dynamics among partners

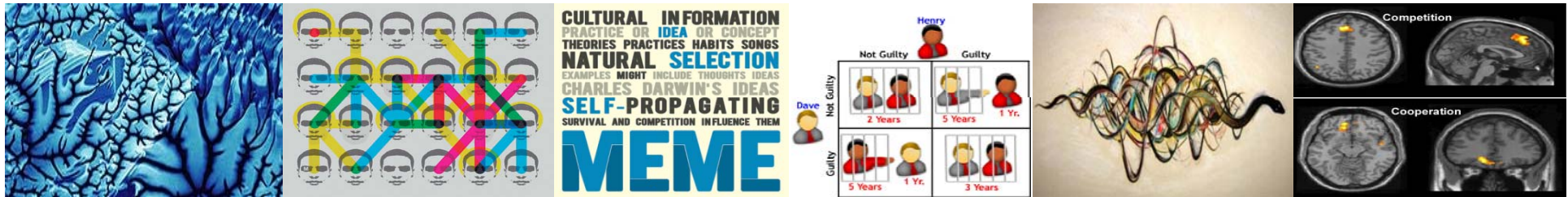# Feedback by punishment + group migration promotes cooperative behavior



**avg group contribution
punishment only**

**avg group contribution
group migration + punishment**

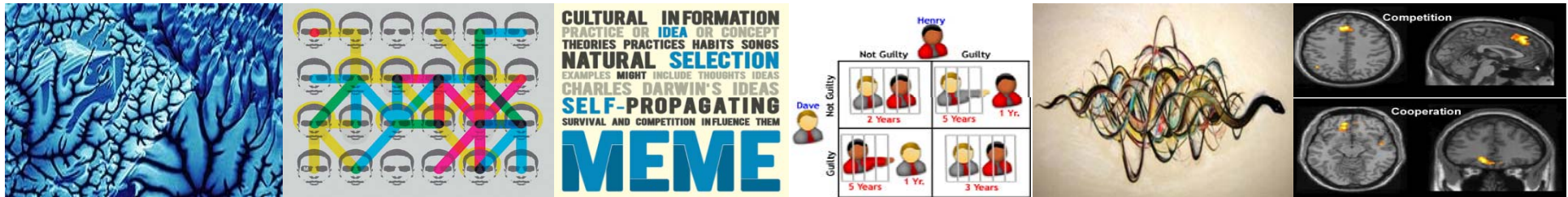# Conclusion

# Conclusion

- The evolution of altruistic punishment can be explained by **disadvantageous inequity aversion**

- Punishment can promote cooperation among social-related individuals (partners)

- Punishment acts as a coordination mechanism among unrelated individuals (strangers)

- To promote cooperation among unrelated individuals, additional mechanisms are required (heterogeneity).

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Outlook: Behavioral Mechanism Design and Social Engineering

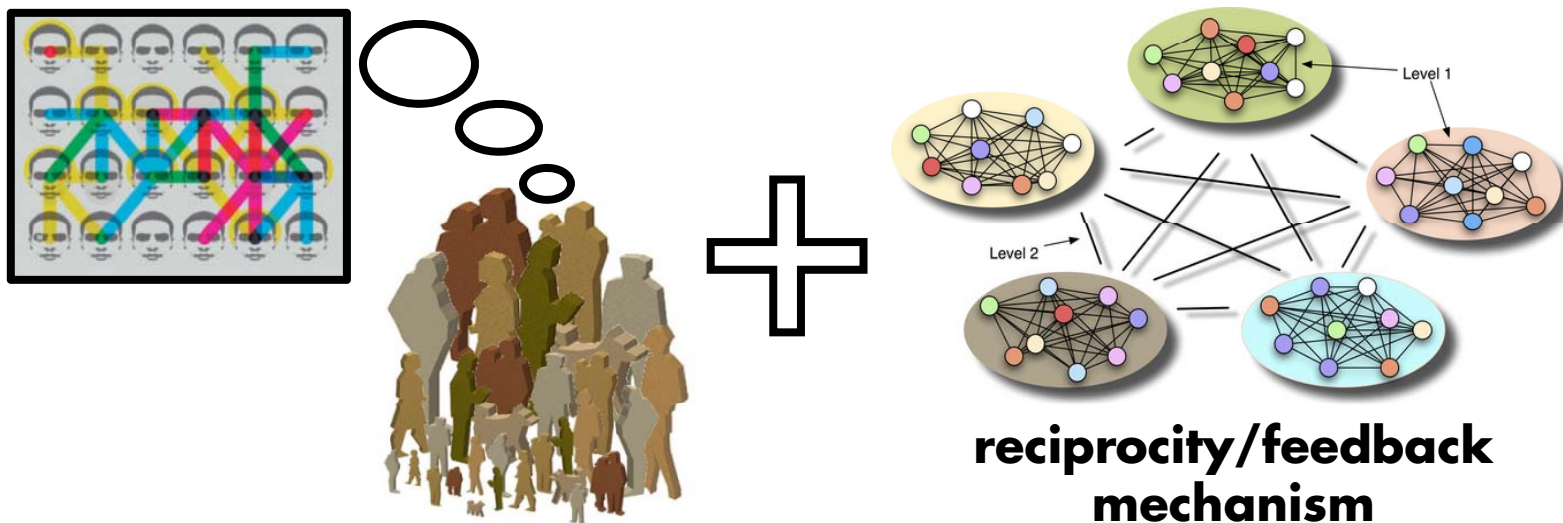# Outlook: Behavioral Mechanism Design

- Mechanism design and contract theory base on the homo economicus assumption.

- They aim at controlling a social system by means of monetary incentive schemes / selfishness assumptions.
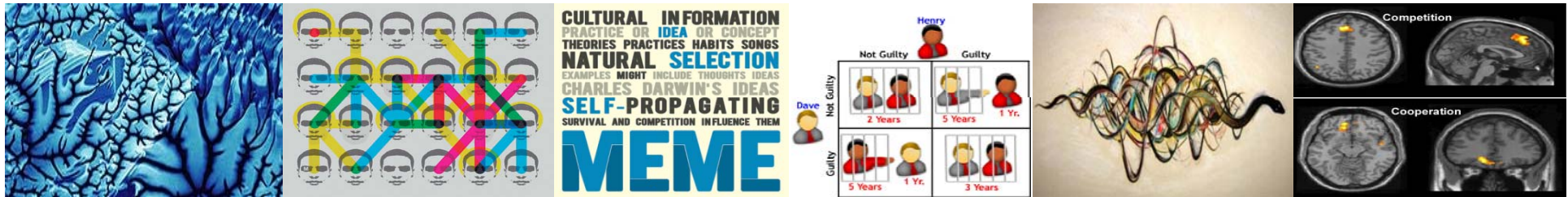
# Outlook: Social capital

- Mechanism design/contract theory should also consider
  - … the impact and the dynamics of social norms
  - … reciprocal effects
  - … altruistic behavior
  - … fairness perception, and many more…

  **The value of "social capital" is underrated!**



reciprocity/feedback
mechanism

# Thanks for your attention!

# Questions, comments and criticism are very welcome!

# Conclusion

- The evolution of altruistic punishment can be explained by **disadvantageous inequity aversion**

- Punishment can promote cooperation among social-related individuals (partners)

- Punishment acts as a coordination mechanism among unrelated individuals (strangers)

- To promote cooperation among unrelated individuals, additional mechanisms are required (heterogeneity).