# When extremists win
## Iterated learning with heterogenous agents

**Dan Navarro**
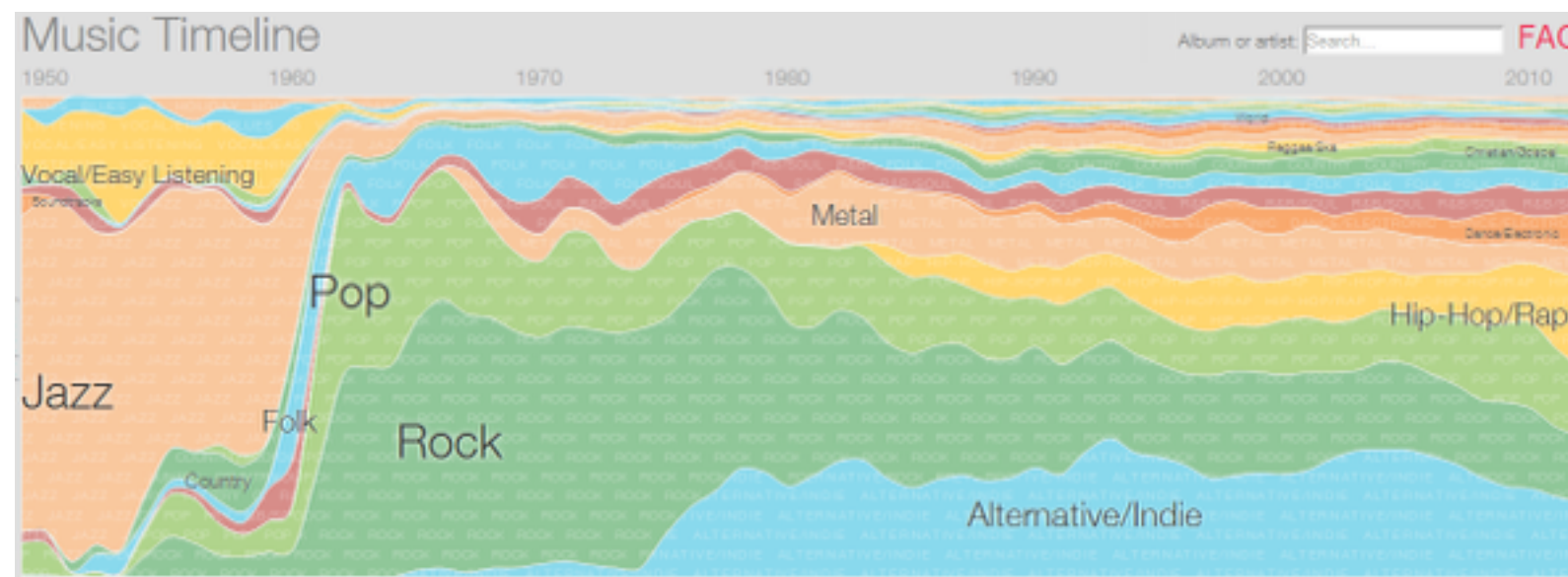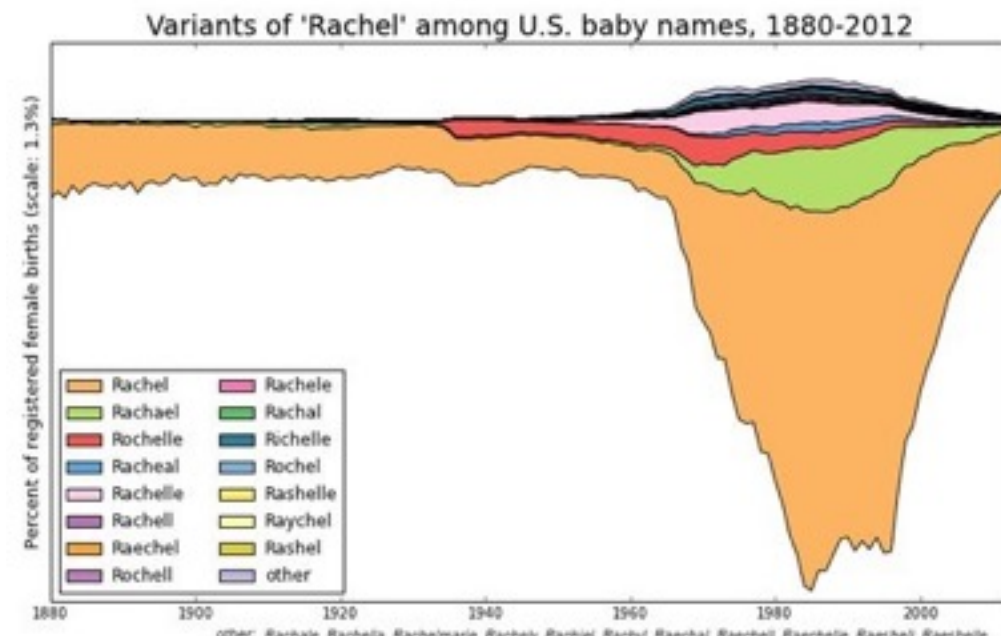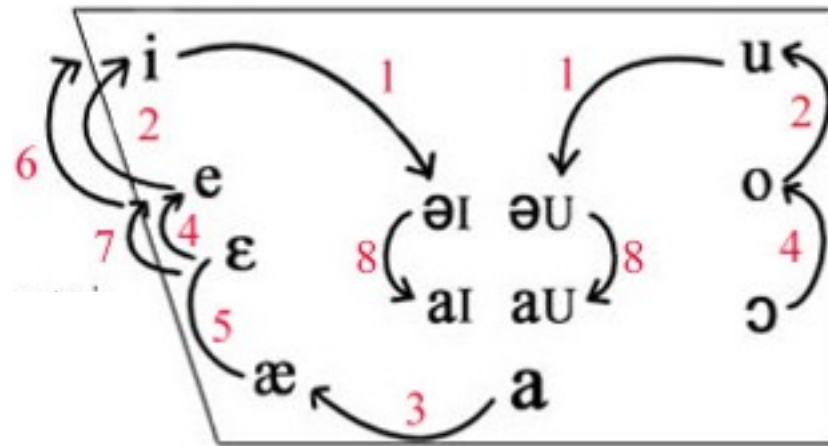School of Psychology
University of New South Wales

**Amy Perfors**
School of Psychology
University of Adelaide

**Arthur Kary**
School of Psychology
University of New South Wales

**Scott Brown**
School of Psychology
University of Newcastle

**Chris Donkin**
School of Psychology
University of New South Wales

# What dynamics underpin cultural and linguistic change? What do they say about the mind?





Variants of 'Rachel' among U.S. baby names, 1880-2012
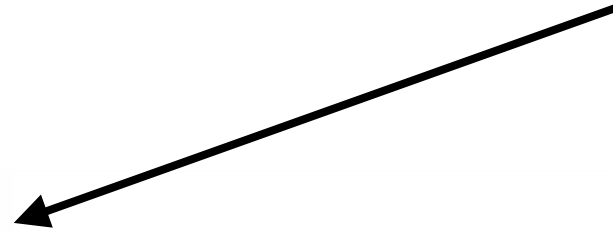


Music Timeline

Original Drawing → Reproduction 1

Reproduction 1 → Reproduction 2

Reproduction 2

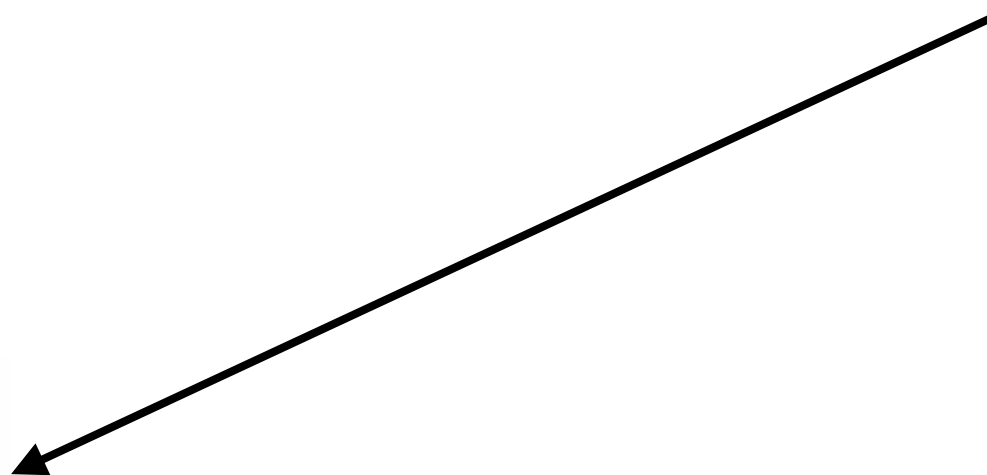Reproduction 3

Reproduktion 3 → Reproduktion 3

Reproduktion 3 → Reproduktion 5

Reproduktion 5                    Reproduktion 6.

Reprodukcia 6.

Reprodukcia 7

Reduction?            Reduction ?

Original Drawing

Reproduction 1

Reproduction 2

Reproduction 3

Reproduction 4

Reproduction 5

Reproduction 6.

Reproduction 7

Reproduction 8
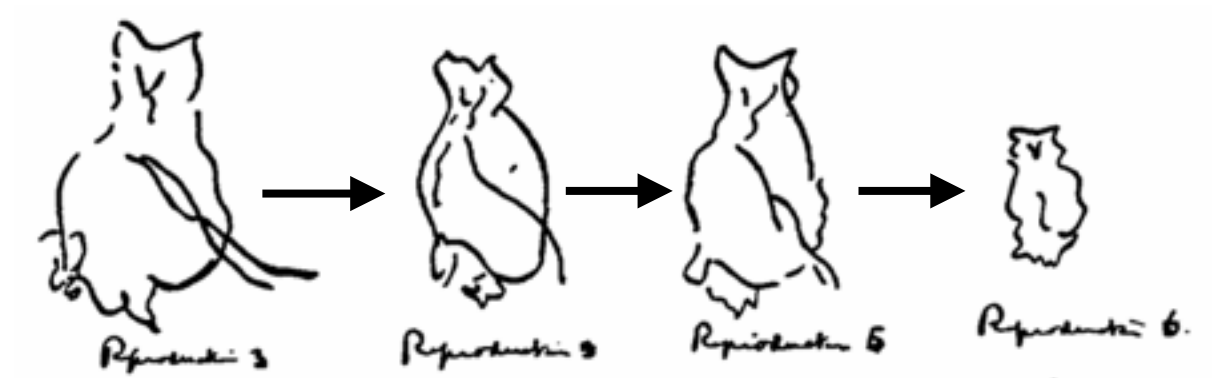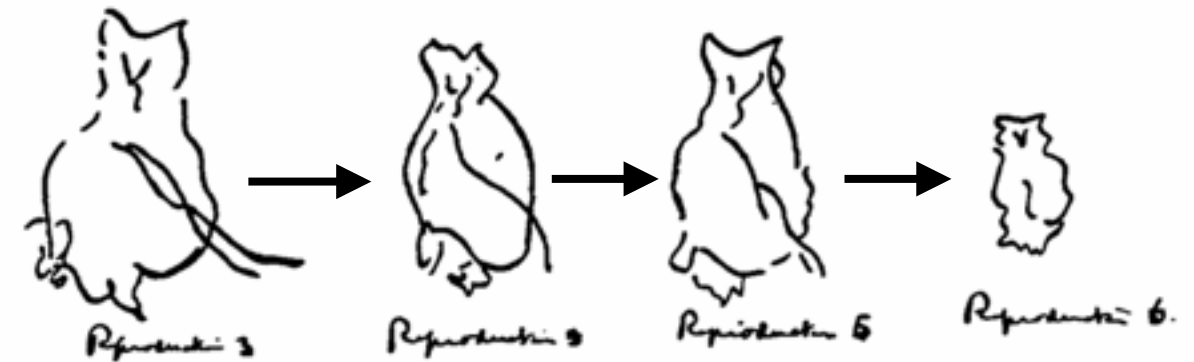
Reproduction 9

Reproduction 10

The method of serial
reproduction in memory

Bartlett (1920)

The method of serial
reproduction in memory
Bartlett (1920)



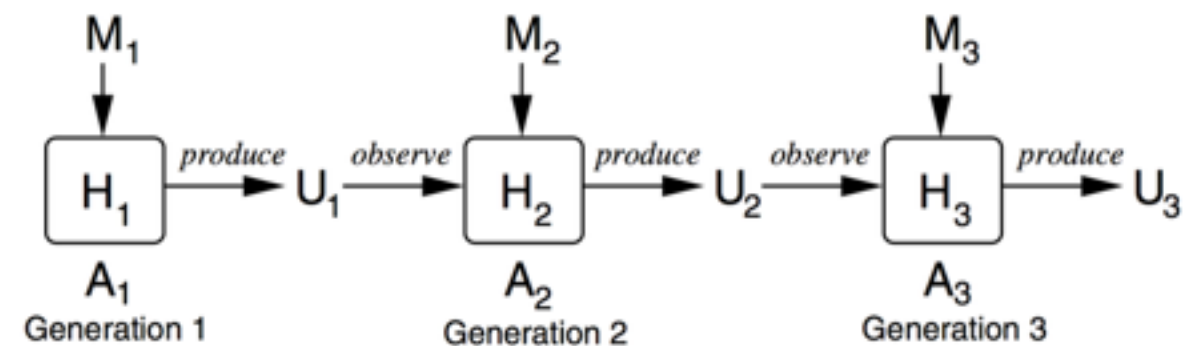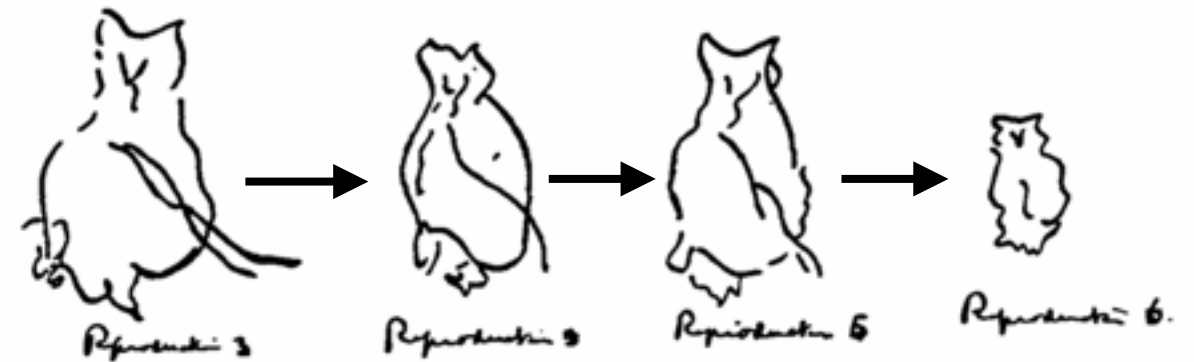Language as sequential
reproduction of culture
Smith et al (2002)



$M_1$ $\rightarrow$ $H_1$ $\xrightarrow{produce}$ $U_1$ $\xrightarrow{observe}$ $M_2$ $\rightarrow$ $H_2$ $\xrightarrow{produce}$ $U_2$ $\xrightarrow{observe}$ $M_3$ $\rightarrow$ $H_3$ $\xrightarrow{produce}$ $U_3$

$A_1$
Generation 1

$A_2$
Generation 2

$A_3$
Generation 3
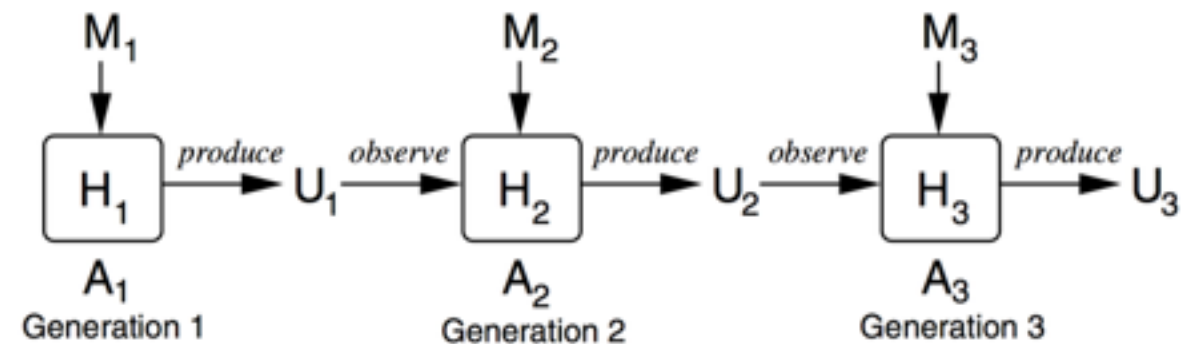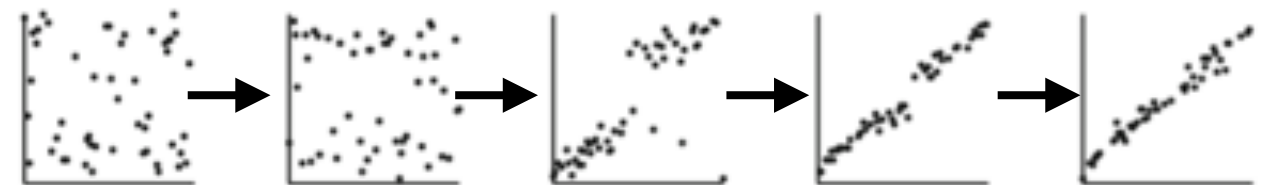
Figure 2. The iterated learning model. The $i$th generation of the population consists of a single agent $A_i$ who has hypothesis $H_i$. Agent $A_i$ is prompted with a set of meanings $M_i$. For each of these meanings the agent produces an utterance using $H_i$. This yields a set of utterances $U_i$. Agent $A_{i+1}$ observes $U_i$ and forms a hypothesis $H_{i+1}$ to explain the set of observed utterances. This process of observation and hypothesis formation constitutes learning.

The method of serial
reproduction in memory
Bartlett (1920)



Language as sequential
reproduction of culture
Smith et al (2002)



Figure 2. The iterated learning model. The $i$th generation of the population consists of a single agent $A_i$ who has hypothesis $H_i$. Agent $A_i$ is prompted with a set of meanings $M_i$. For each of these meanings the agent produces an utterance using $H_i$. This yields a set of utterances $U_i$. Agent $A_{i+1}$ observes $U_i$ and forms a hypothesis $H_{i+1}$ to explain the set of observed utterances. This process of observation and hypothesis formation constitutes learning.

The method of iterated
learning reveals inductive bias
Kalish et al (2007)

# Example: function learning

(Kalish et al 2007)

A

B

C

D

E

original

# Example: function learning

(Kalish et al 2007)



original ⟶ final

# Example: function learning

(Kalish et al 2007)

Conclusion: we have an inductive bias for linear functions

# Proof that iterated learning with Bayesian agents reveals the prior

$$P(h_n = i) = \sum_j P_{\text{samp},PA}(h_n = i \mid h_{n-1} = j)P(h_{n-1} = j)$$

$$= \sum_j \sum_{d \in \mathcal{D}} P_{\text{samp}}(h_n = i \mid d)P_{PA}(d \mid h_{n-1} = j)P(h_{n-1} = j)$$

$$= \sum_{d \in \mathcal{D}} P_{\text{samp}}(h_n = i \mid d) \sum_j P_{PA}(d \mid h_{n-1} = j)P(h_{n-1} = j)$$

$$= \sum_{d \in \mathcal{D}} P_{\text{samp}}(h_n = i \mid d)P_{PA}(d)$$

$$= \sum_{d \in \mathcal{D}} \frac{P_{PA}(d \mid h_n = i)P(h_n = i)}{P_{PA}(d)}P_{PA}(d)$$

$$= P(h_n = i) \sum_{d \in \mathcal{D}} P_{PA}(d \mid h_n = i),$$

(Griffiths & Kalish 2007)



… as long as everyone has the same prior

Hm.

Hm.

So how do iterated learning chains
behave when individual differences exist?

# Case study 1:
Does everybody contribute equally to the evolution of languages?

Bayesian models for language regularisation with two different kinds of bias



**Strong** bias *for* regularity

**Weak** bias *against* regularity

Prior Density

Proportion of Rule Violations

Homogenous iterated learning chains converge to the prior
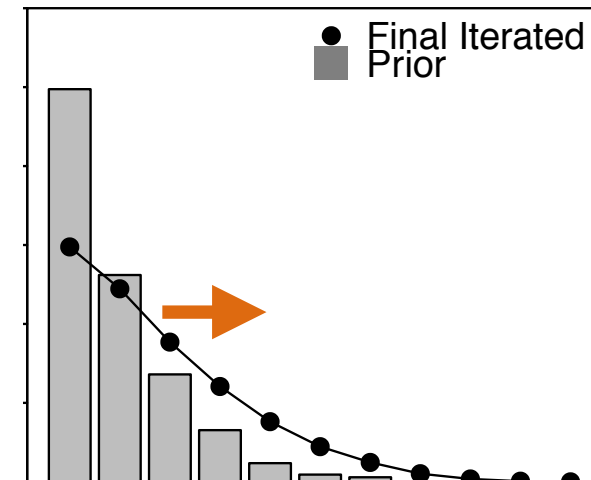
A heterogenous chain does not converge to the average of the prior biases

…and the distribution of responses is severely distorted

weak bias          strong bias

Final Iterated
Prior

strong bias

⬇

insensitivity to input

⬇

greater influence on the chain

# Case study 2:
# Bayesian groupthink

Juror $i$ records vote, removes sheet, passes notebook

Juror $i$ records vote, removes sheet, passes notebook

Juror $i+1$ can see the previous vote via indentations…

Prior belief about guilt
$P(g)$ is set by the trial

Likelihood of previous juror's vote $P(v|g)$ requires *theory of mind*… what do they know that I don't know?

# Bayesian "sheep"



$P(v|g) = 0.95$

Assumes previous juror has additional knowledge, assigns evidentiary weight to their opinion

**100% Goats**

**100% Goats**

A jury of goats ignores one another and the "chain" converges just fine

**100% Sheep**

A jury of sheep displays **groupthink**

$$\boldsymbol{\pi T} \quad \propto \quad [d,p]\begin{bmatrix} 1-p & p \\ d & 1-d \end{bmatrix}$$
$$= \quad [d(1-p)+pd, dp+p(1-d)]$$
$$= \quad [d,p] \propto \boldsymbol{\pi}$$

**50% Sheep, 50% Goat**



A mixed jury is dominated by goats

**50% Sheep, 50% Goat**

# Case study 3:
## An empirical illustration

# "Who will win the 2016 Australian election?"



N=80 MTurk workers
and UNSW students

Andy?

Rating

N=80 MTurk workers
and UNSW students

# The advisor task

**UNSW Politics**

Australians ignored the advisor and predicted a Turnbull victory

N=124 UNSW students

MTurk Politics

Americans followed
the advisor regardless

N=196 MTurk workers

Americans *claim* to be totally ignorant about Australian politics…

… and an all American iterated learning chain "reveals" a "preference" for _Gordon Brown_ …

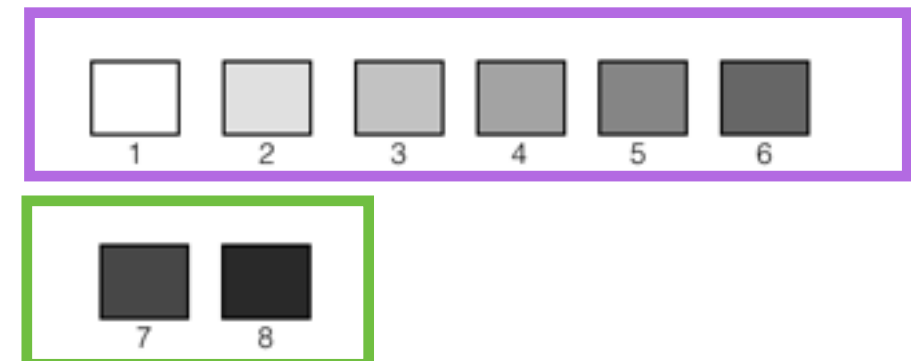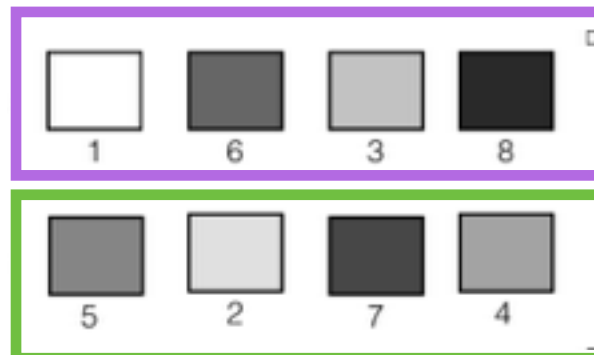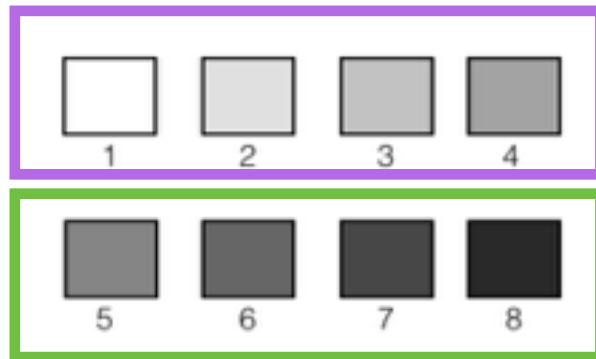If we mix some Australians into the chain the Americans endorse *Malcolm Trunbull*

Proportion Australian

Australians choose Turnbull no matter how many Americans are included

# Case study 4:
# A non-Bayesian example

Iterated learning can be used to study the biases people bring to categorisation problems
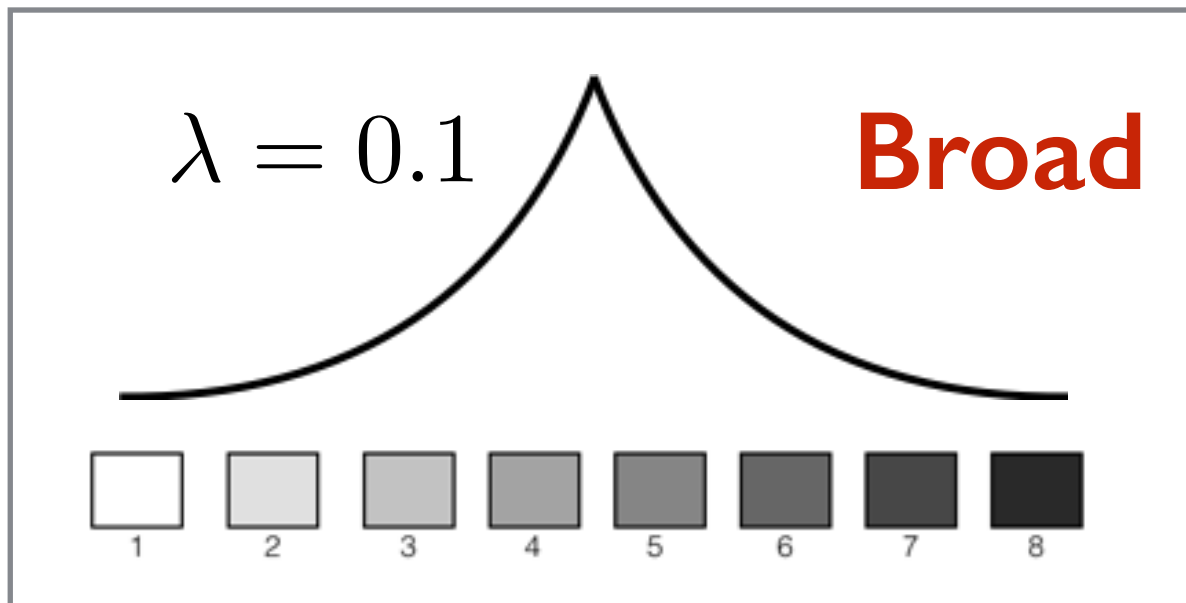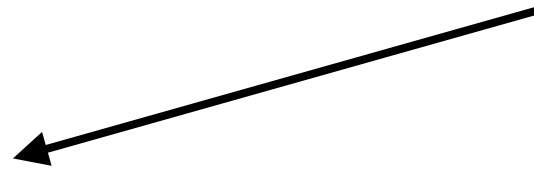
(e.g., Austerweil 2014)

# Exemplar model of categorisation
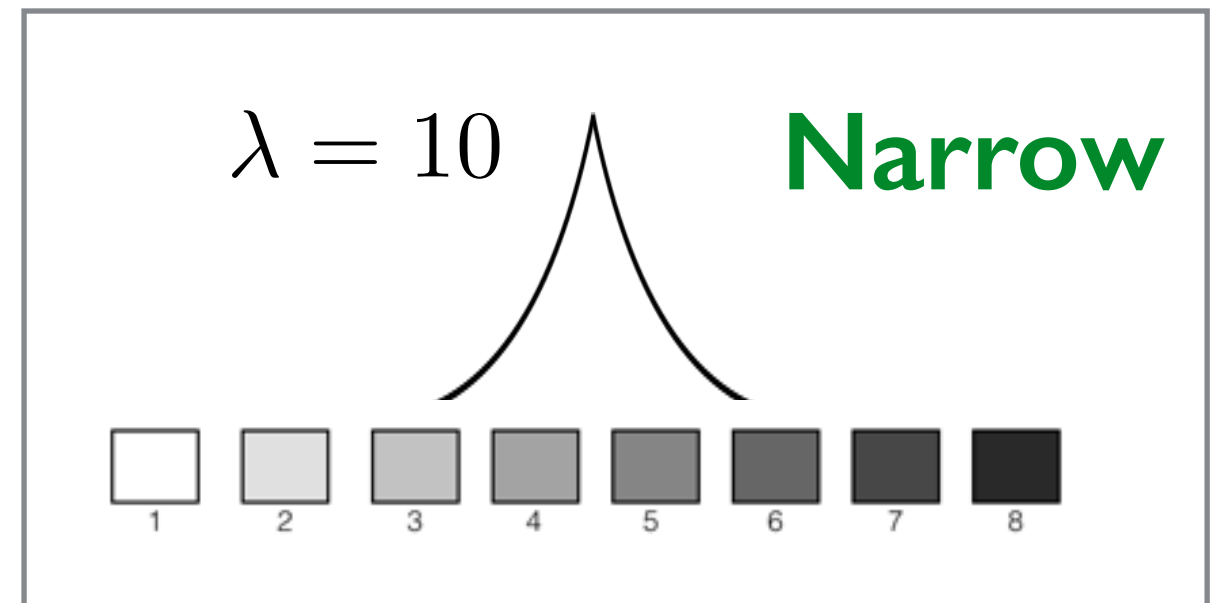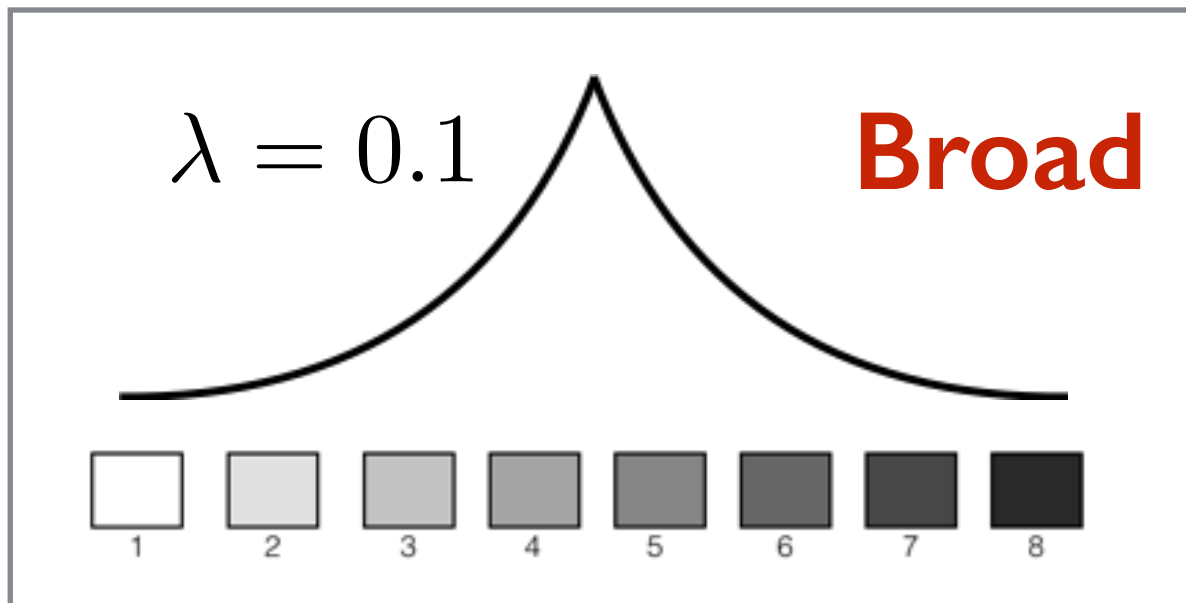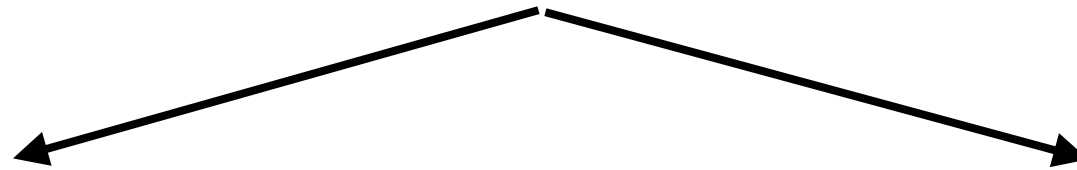
(Nosofsky 1986; Pothos & Bailey 2009)



GCM: categorisation probability is
proportional to sum similarity

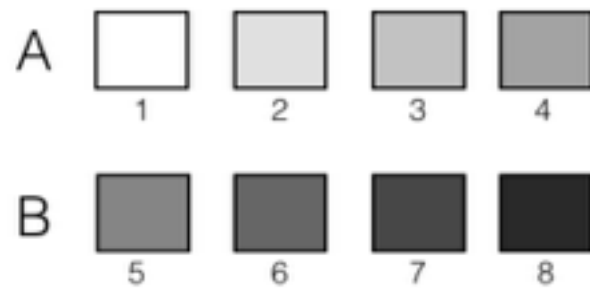GCM allows learners to vary in how broadly they generalise from a stimulus

# GCM allows learners to vary in how broadly they generalise from a stimulus
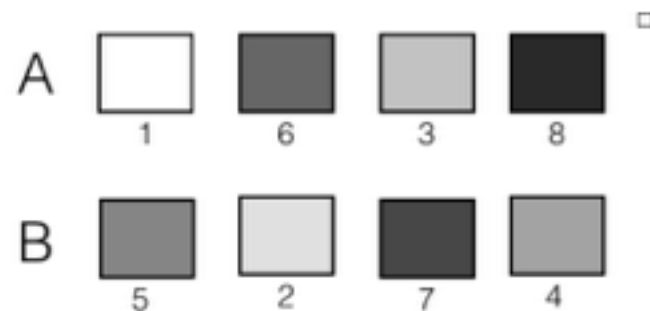


$\lambda = 0.1$   **Broad**

$\lambda = 10$   **Narrow**

# *Categorisation bias #1*



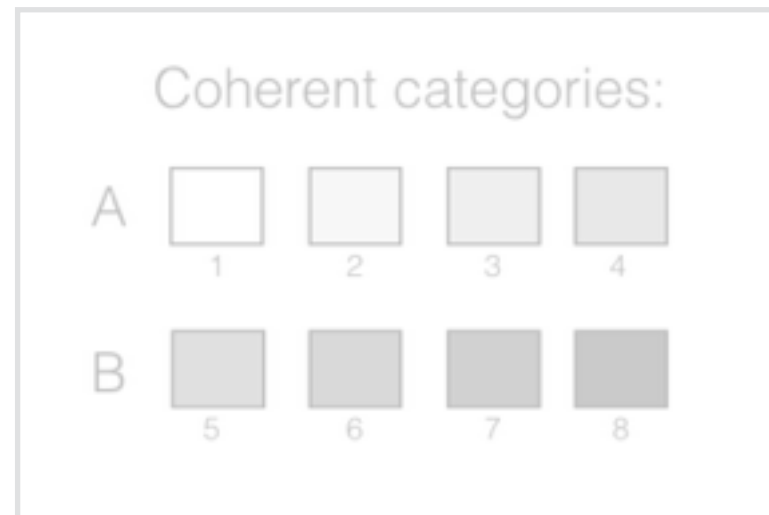*Coherent* systems
assign similar items
to the same category
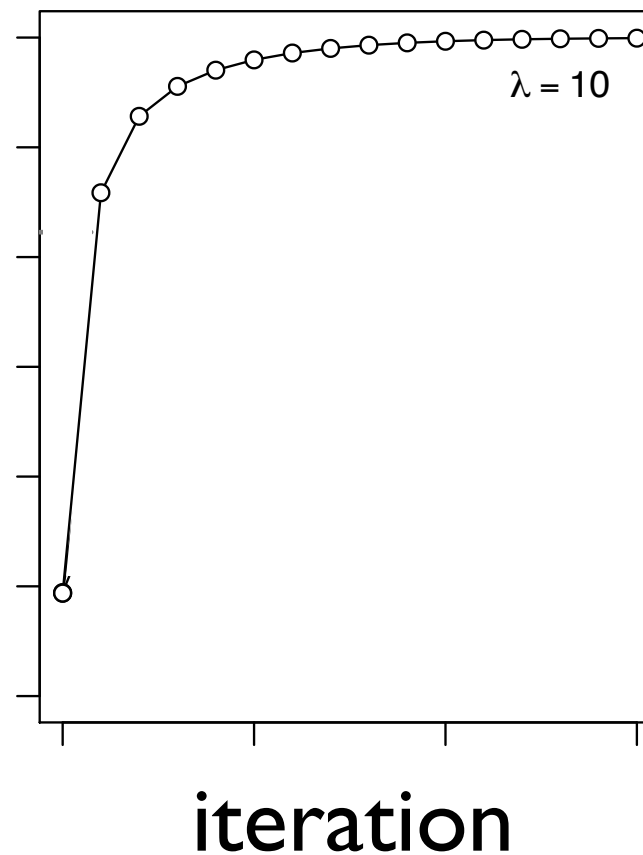
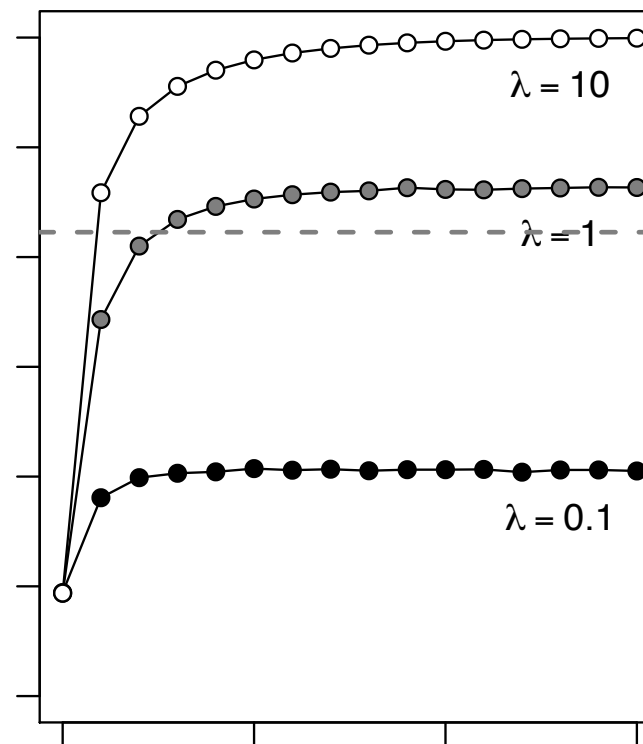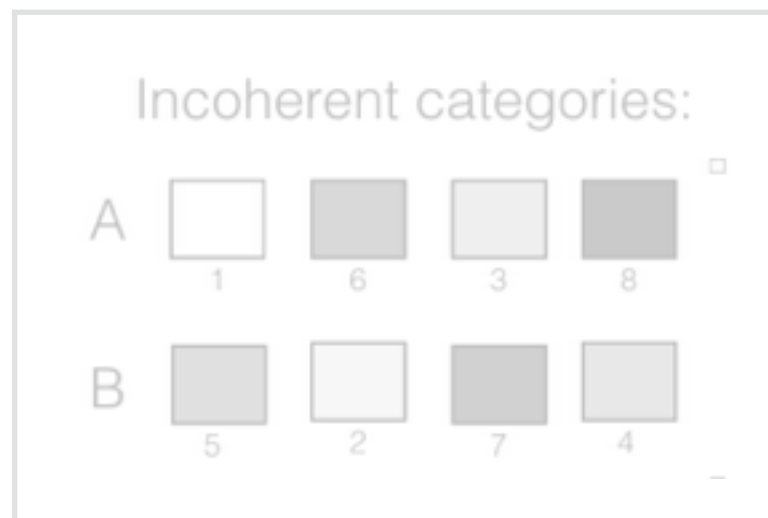# Iterated learning with GCM when learners are **homogenous**



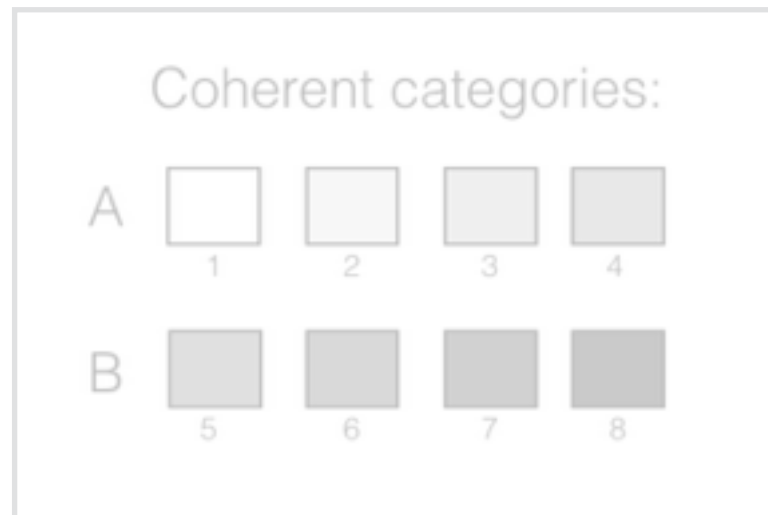**Narrow** generalisation implies strong coherence bias

λ = 10

coherence

iteration

Coherent categories:

A
1 2 3 4

B
5 6 7 8

Incoherent categories:

A
1 6 3 8

B
5 2 7 4

# Iterated learning with GCM when learners are **homogenous**
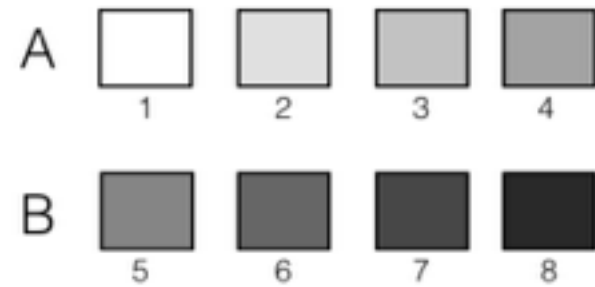
**Heterogeneity** isn't much of a problem here

**Equally sized categories**

A 1 2 3 4

B 5 6 7 8

**Unequally sized categories**

A 1 2 3 4 5 6

B 7 8

*Categorisation bias #2*

# Homogenous chains

Equally sized categories

A  1  2  3  4

B  5  6  7  8

Unequally sized categories

A  1  2  3  4  5  6

B  7  8

Narrow $\lambda = 10$

$\lambda = 1$

$\lambda = 0.1$

Broad

# Case study 5:
## Belief evolution in a self-organising Bayesian social network

Agents prefer to receive
data from trusted sources

What could
possibly go wrong?

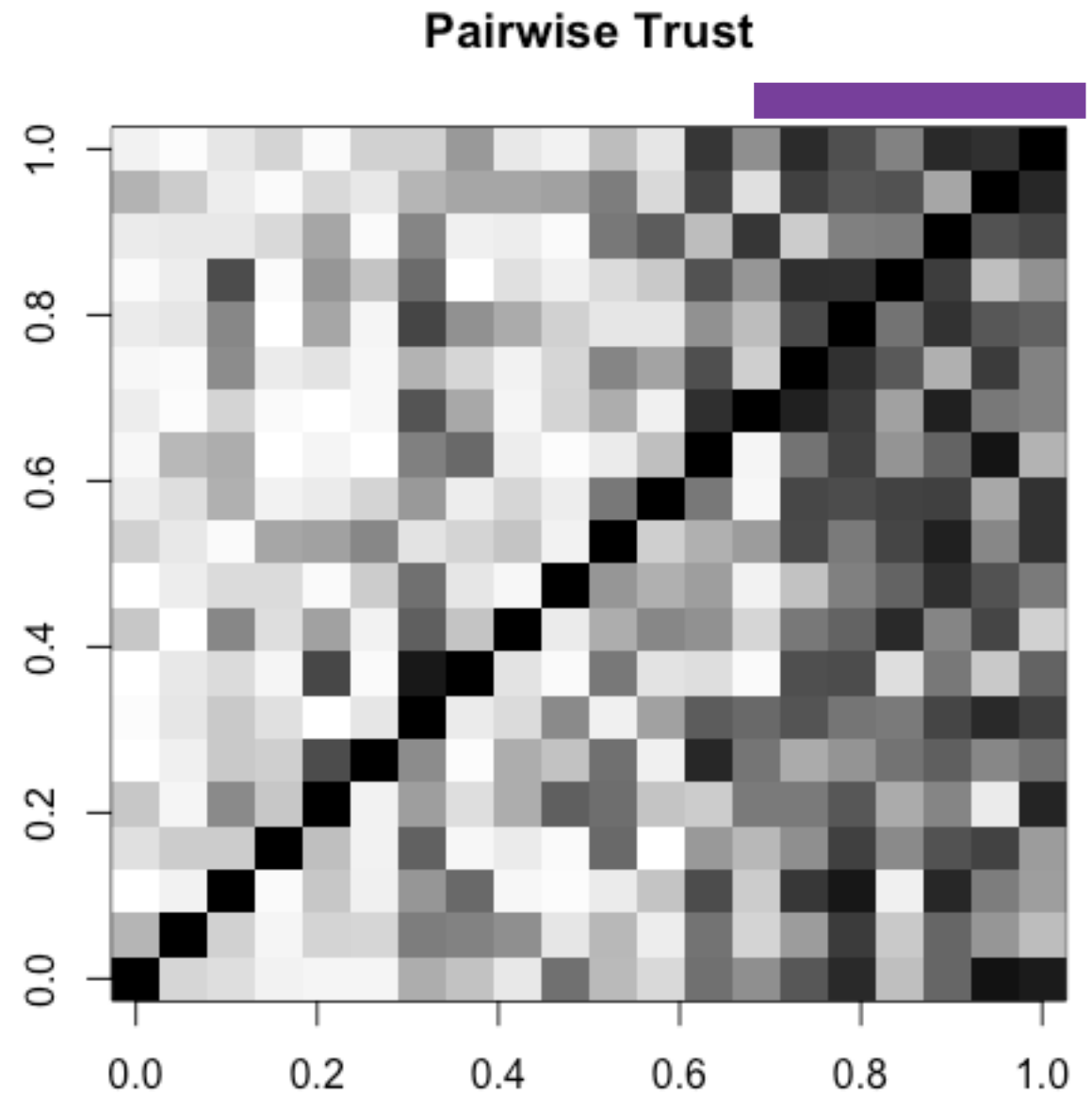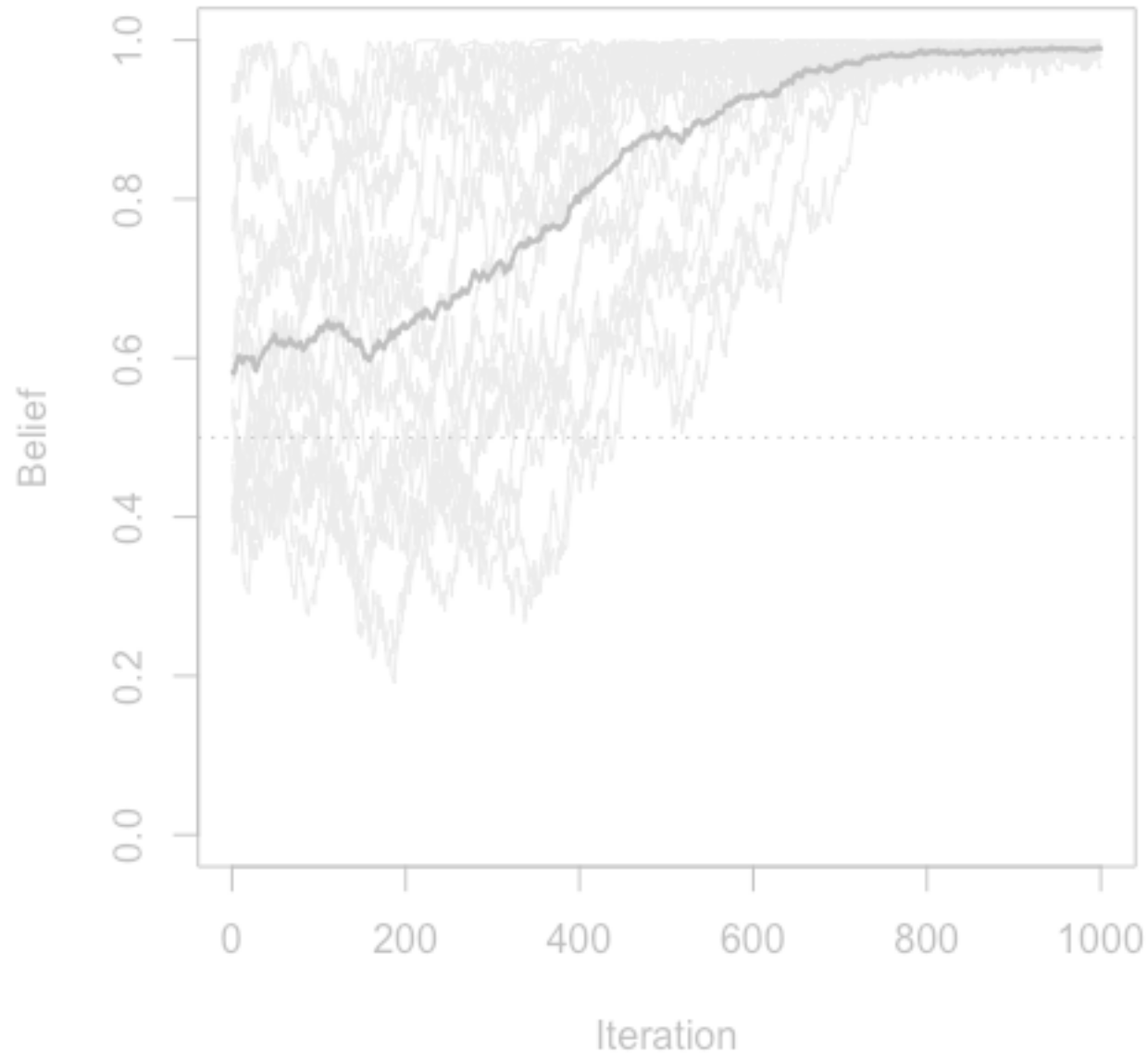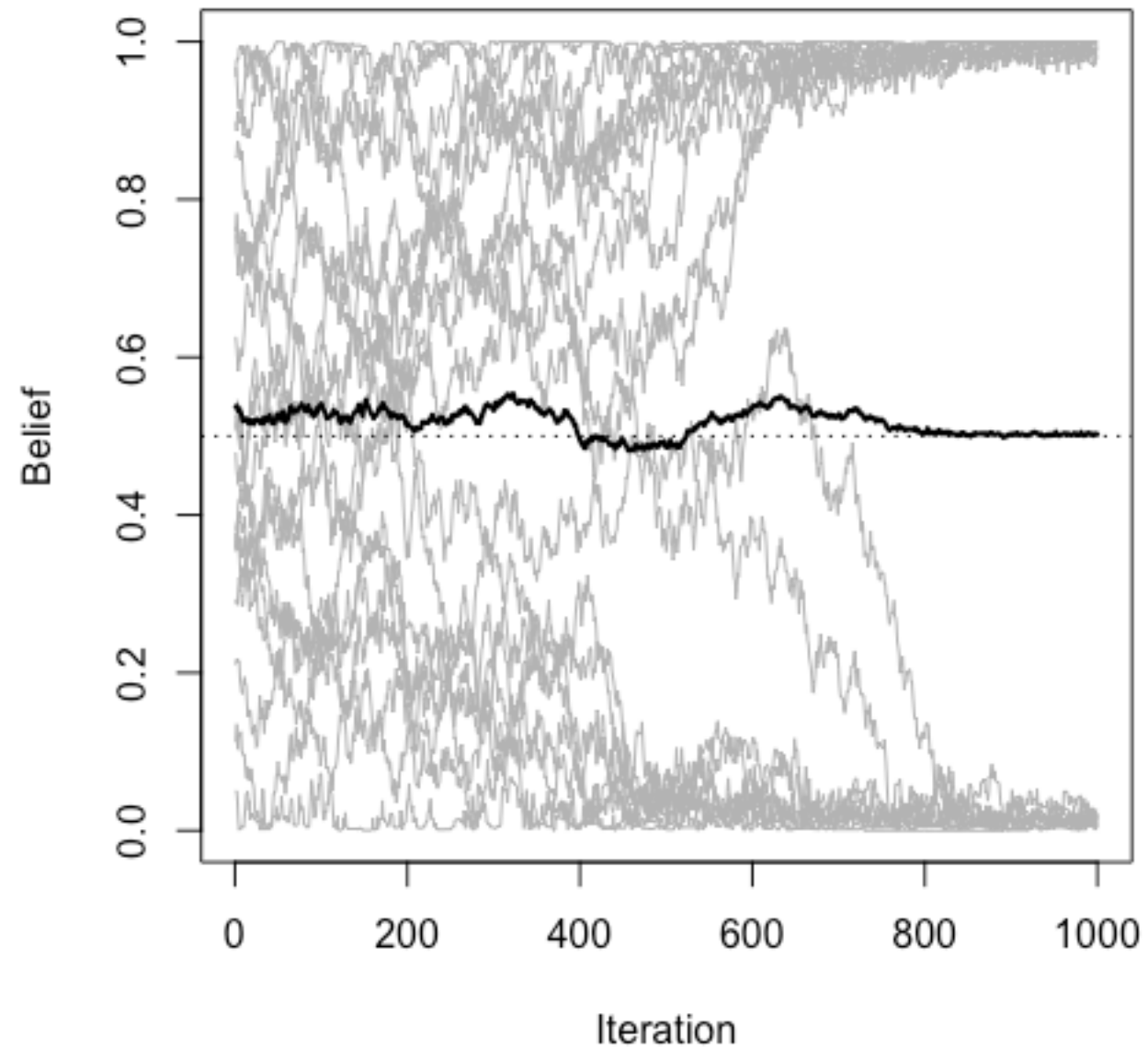They might ratchet themselves into extremism?
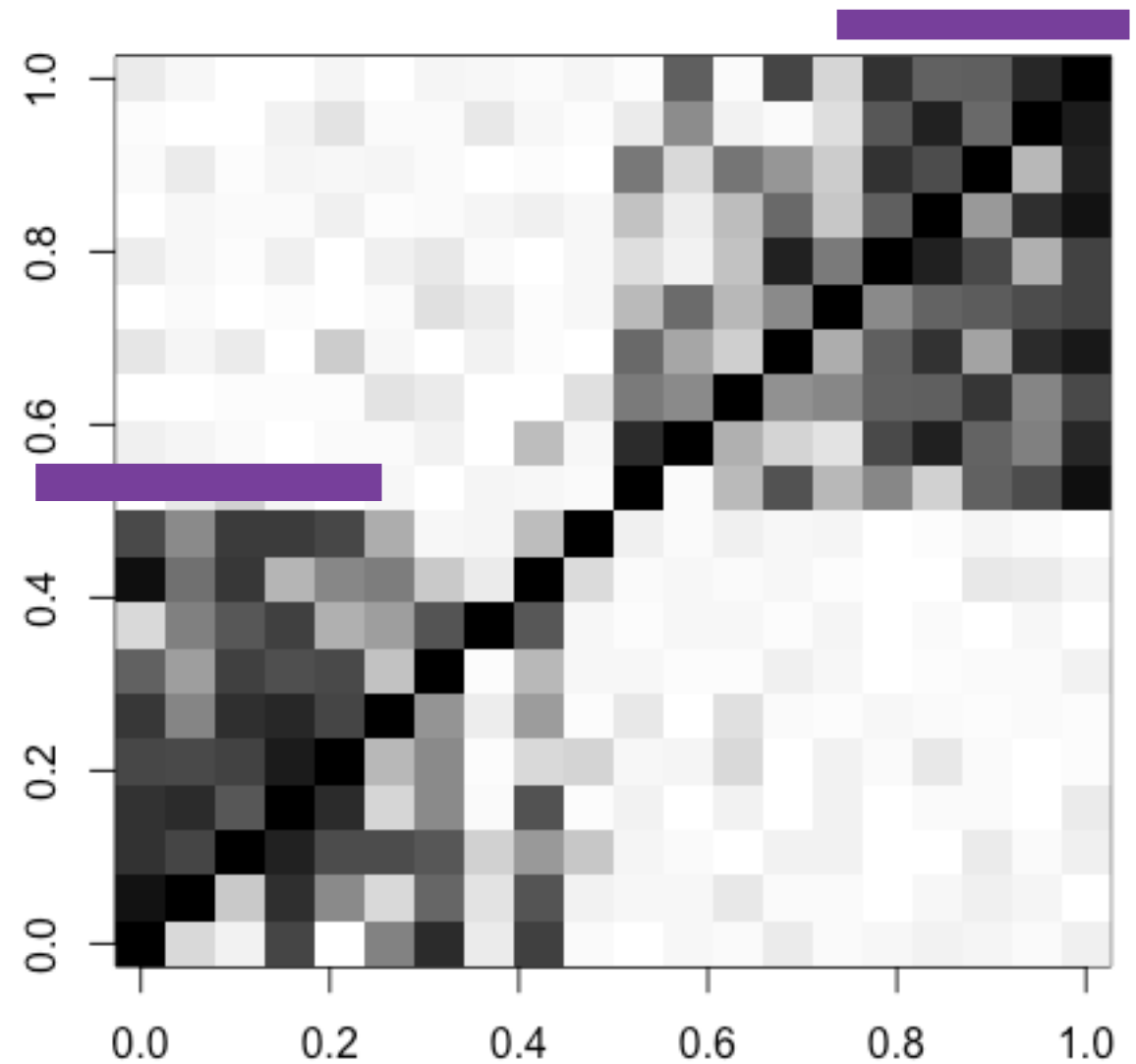
# …with the biggest **extremists** being the most trusted agents
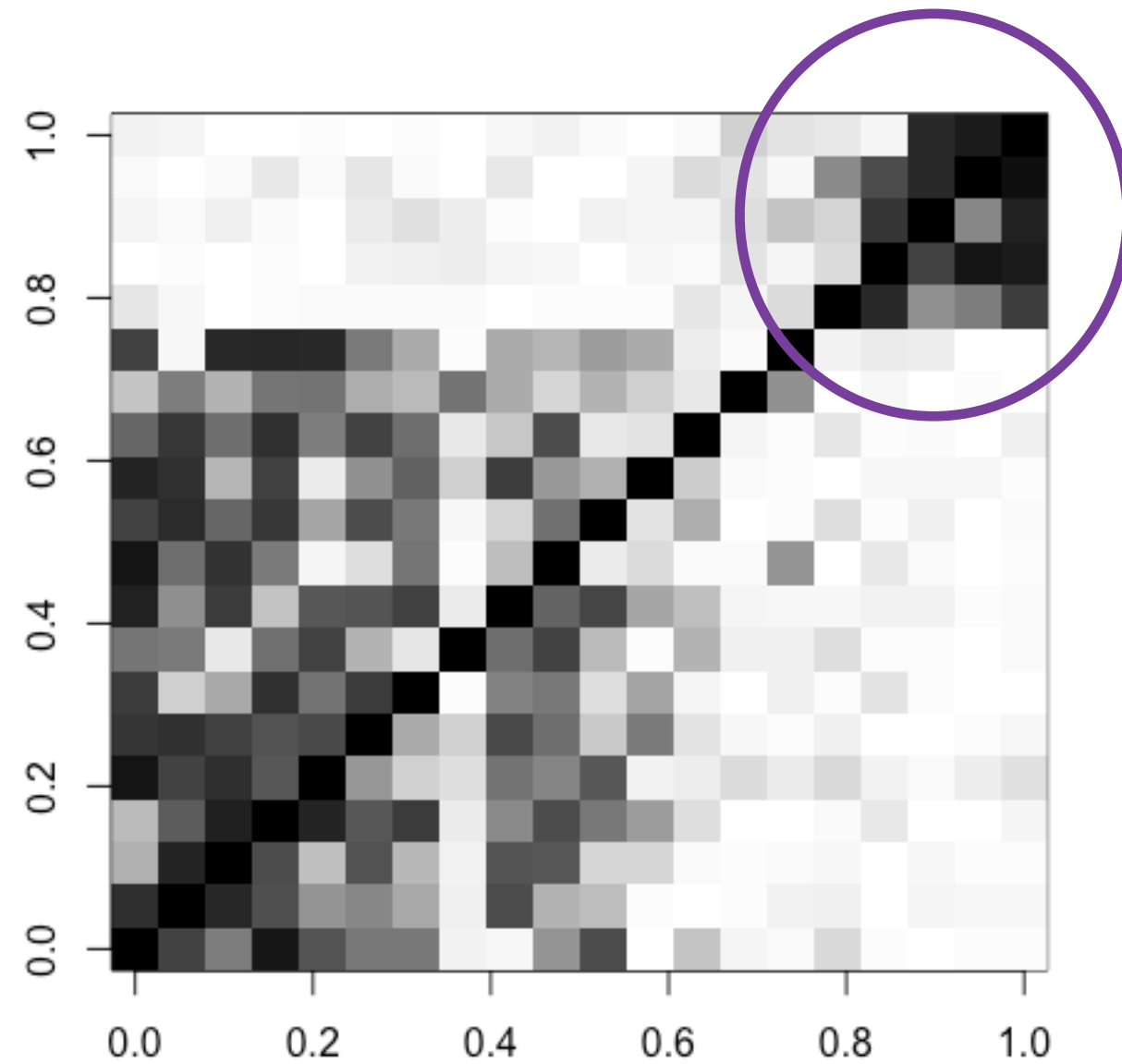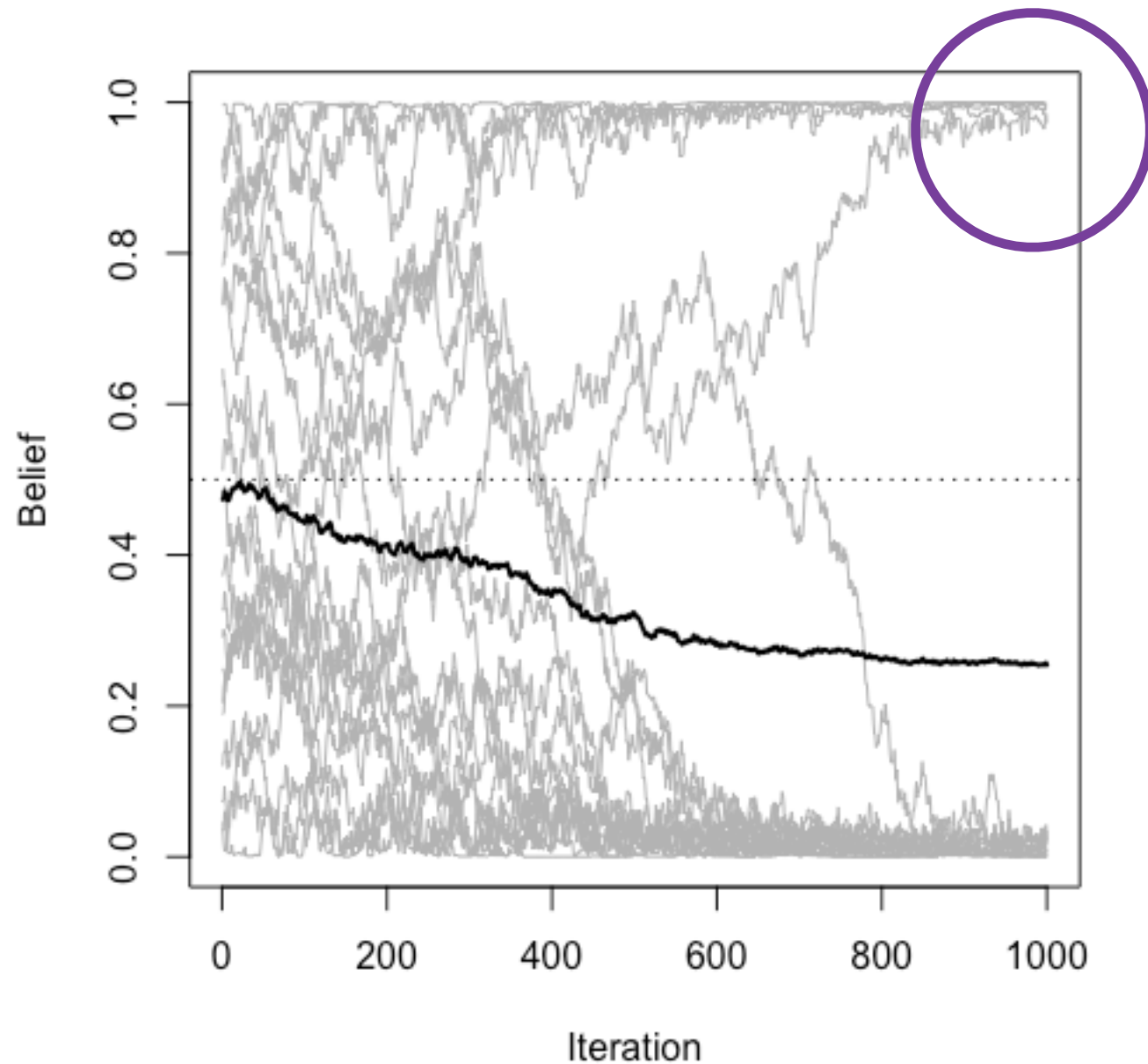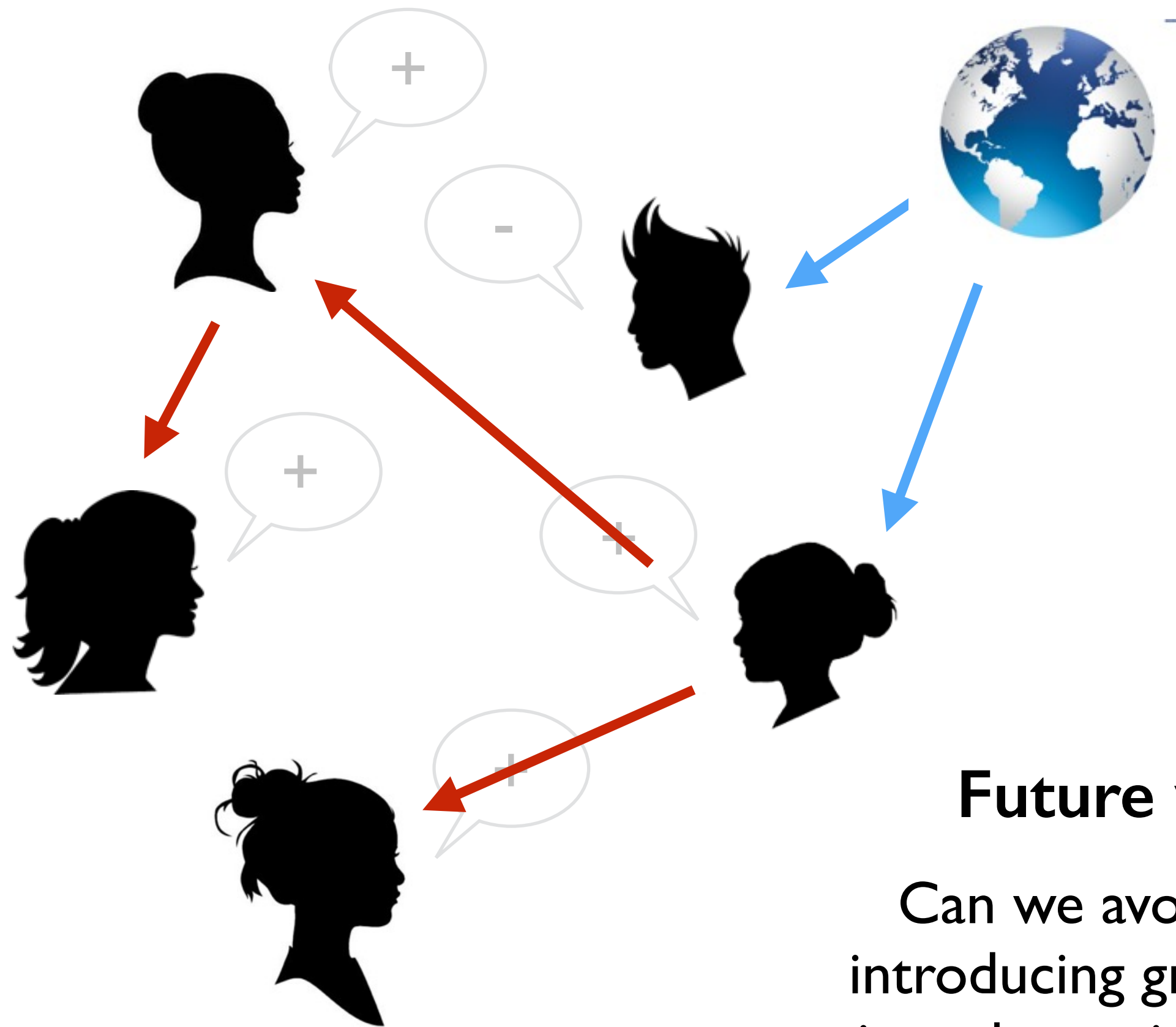


**Pairwise Trust**

They might polarise
into warring factions

# …with the **extremists** being most trusted within group; and no between-group trust

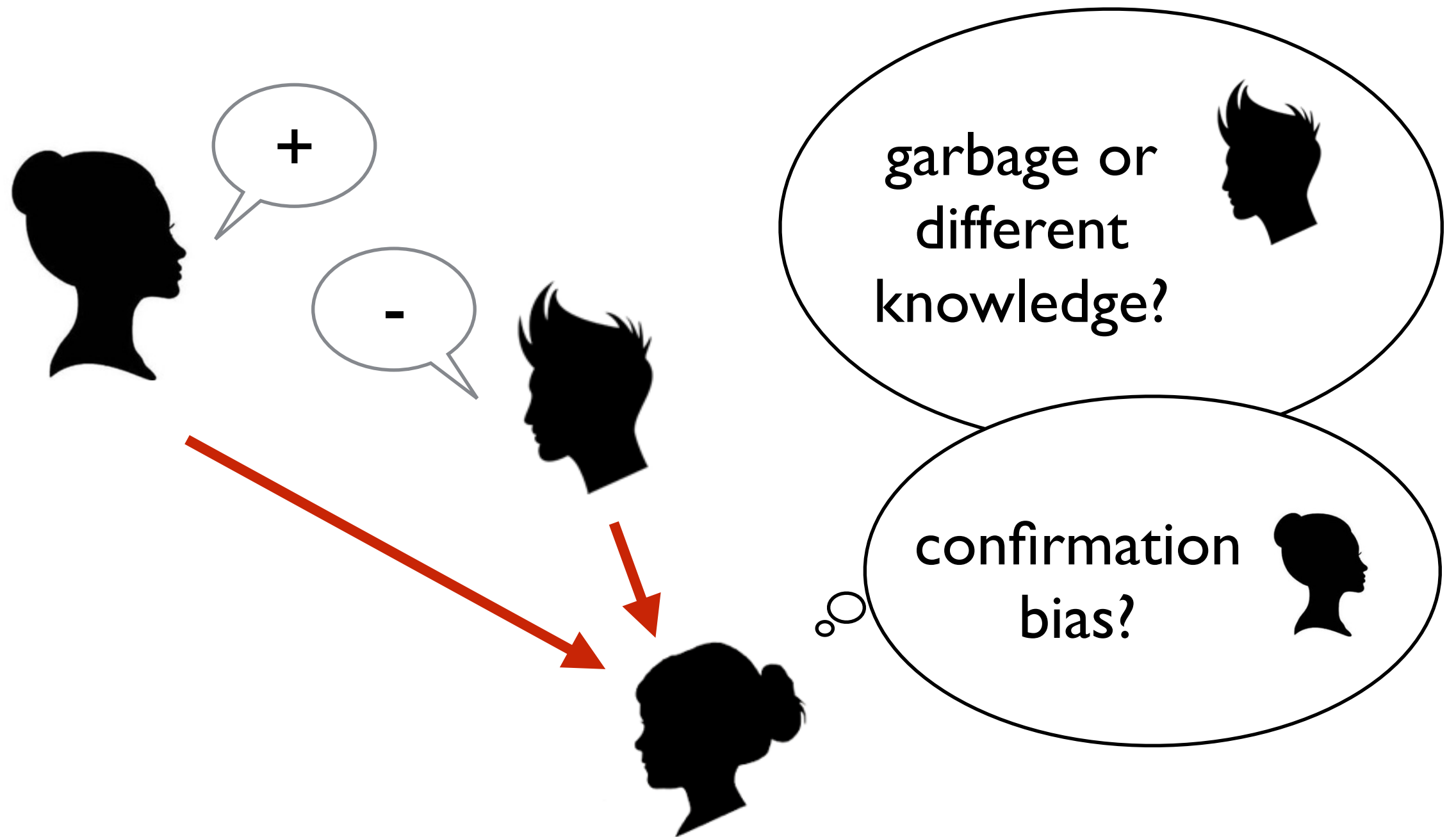# And small "rogue" groups might form their own isolated world.

**Future work:**

Can we avoid this by introducing ground truth into the social network?

**Future work:**

Can we avoid this by giving our agents a more sophisticated ToM?

- <u>Summary</u>:
  - Iterated learning distorts inductive bias when individual differences are present
  - Miscalibrated agents can distort their own inductive biases even in homogenous chains
  - IL chains favour learners with strong biases
  - The magnitude of the distortion is variable
  - Social structure, theory of mind, the link to the world… they all matter

- <u>Implications</u>:
  - IL is limited as a tool for "revealing inductive priors"
  - IL is potentially useful for studying "distortions" in cultural and linguistic evolution

# Thanks!