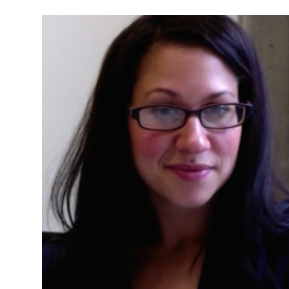


# Simplicity and informativeness in conceptual structure

Jon W. Carr

*Centre for Language Evolution  
School of Philosophy, Psychology and Language Sciences  
University of Edinburgh*

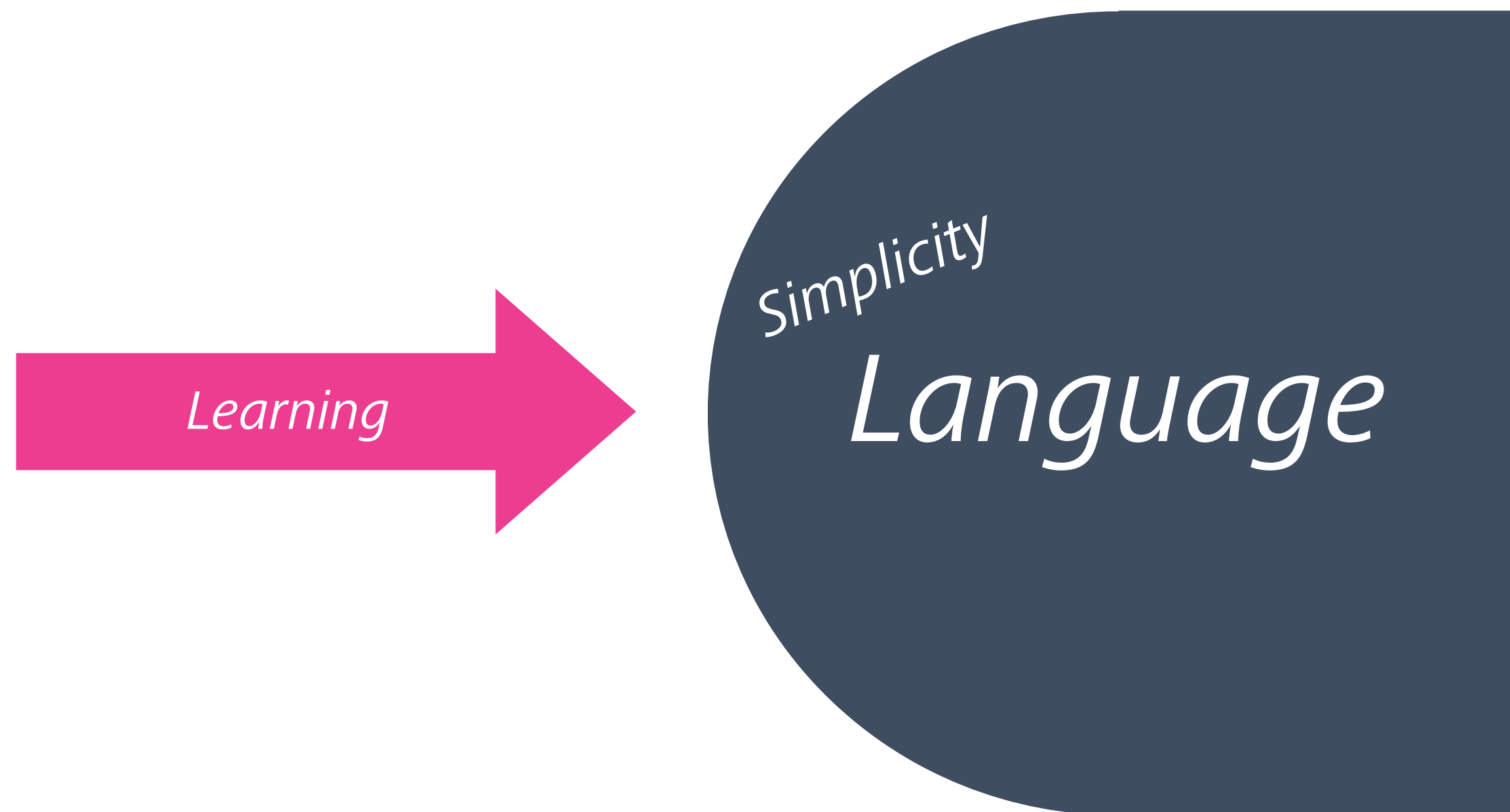


# Pressures shaping language



*Language*

# Pressures shaping language



# Pressures shaping language





# Pressures shaping language



*The simplicity–informativeness tradeoff*



# Kinship terms are simple and informative

## Kinship Categories Across Languages Reflect General Communicative Principles

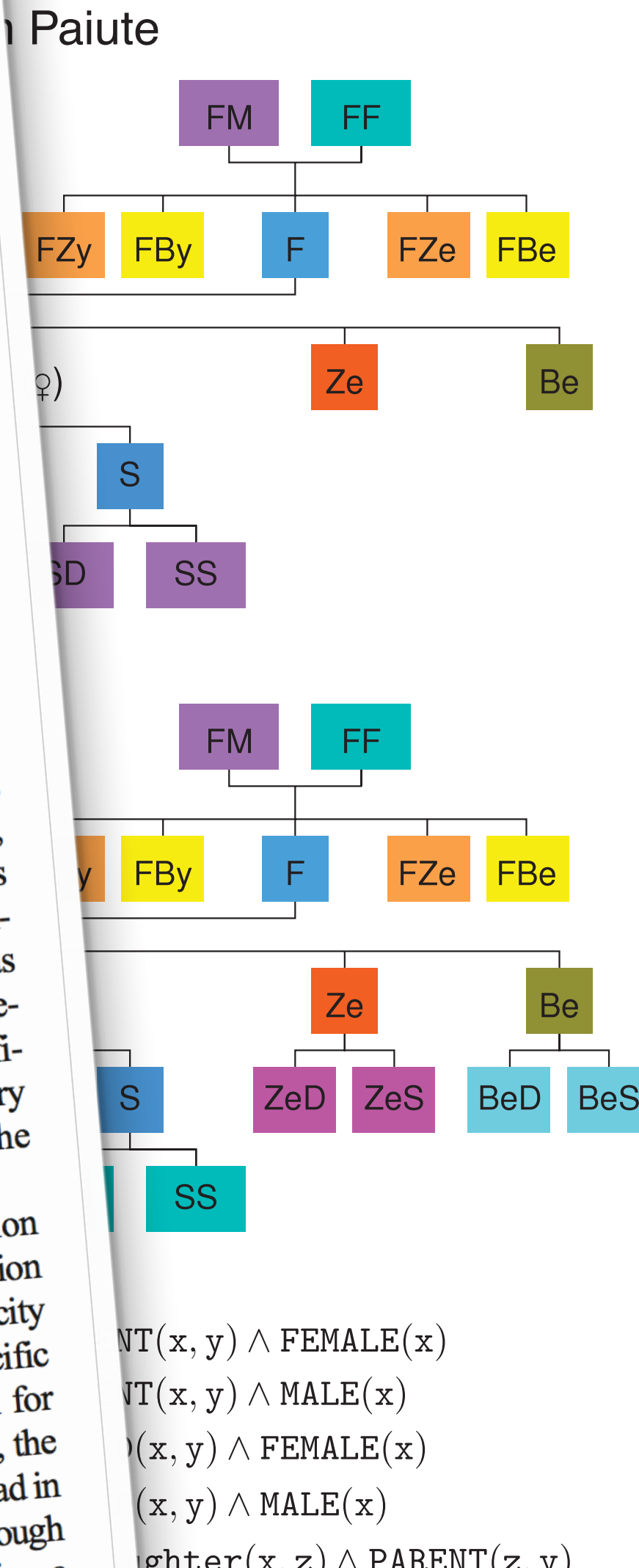
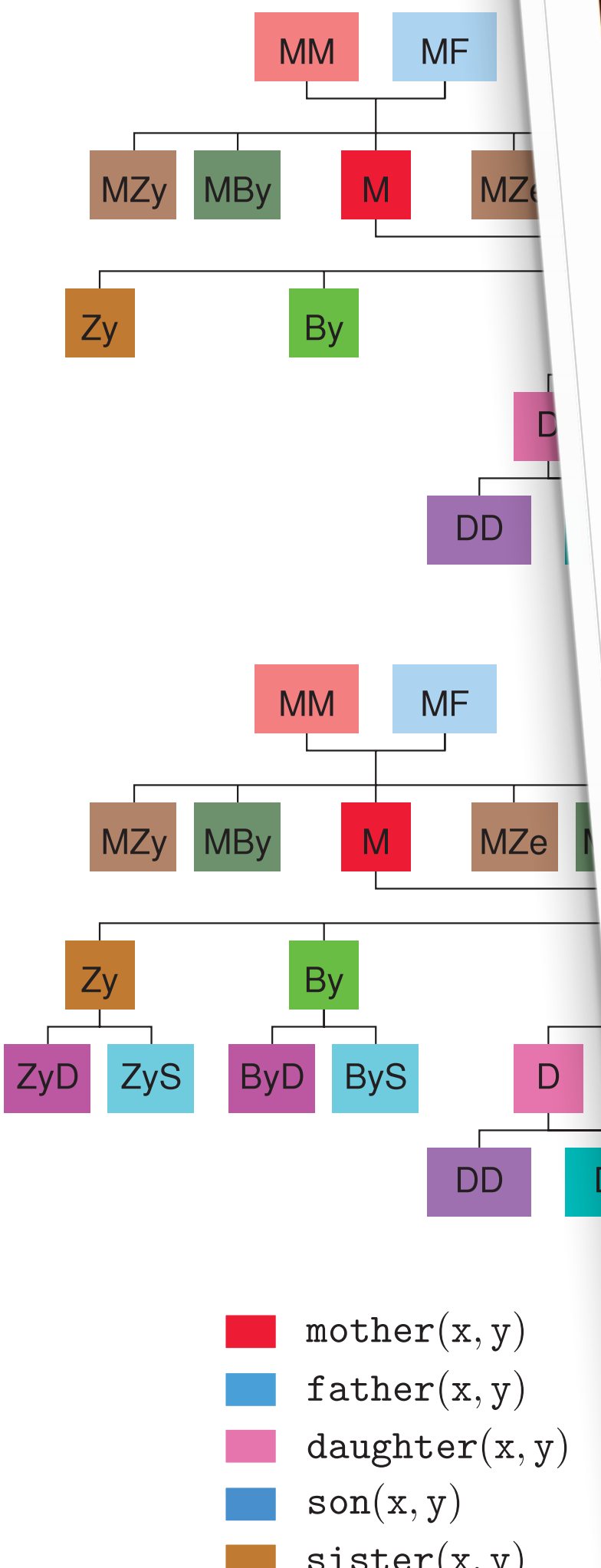
Charles Kemp<sup>1\*</sup> and Terry Regier<sup>2</sup>

Languages vary in their systems of kinship categories, but the scope of possible variation appears to be constrained. Previous accounts of kin classification have often emphasized constraints that are specific to the domain of kinship and are not derived from general principles. Here, we propose an account that is founded on two domain-general principles: Good systems of categories are simple, and they enable informative communication. We show computationally that kin classification systems in the world's languages achieve a near-optimal trade-off between these two competing principles. We also show that our account explains several specific constraints on kin classification proposed previously. Because the principles of simplicity and informativeness are also relevant to other semantic domains, the trade-off between them may provide a domain-general foundation for variation in category systems across languages.

Concepts and categories vary across cultures but may nevertheless be shaped by universal constraints (1–4). Cross-cultural studies have proposed universal constraints that help to explain how colors (5, 6), plants, animals (7, 8), and spatial relations (9, 10) are organized into categories. Kinship has traditionally been a prominent domain for studies of this kind, and researchers have described many constraints that help to predict which of the many logically possible kin classification systems are encountered in practice (11–15). Typically these constraints are not derived from general principles, although it is often suggested that they are consistent with cognitive and functional considerations (2, 11–13, 15). Here, we show that major aspects of kin classification follow directly from two general principles: Categories tend to be simple, which minimizes

cognitive load, and to be informative, which maximizes communicative efficiency. Principles like these have been discussed in other contexts by previous researchers (16–19). For example, Zipf suggested that word-frequency distributions achieve a trade-off between simplicity and communicative precision (20, 21), Hawkins (22) has suggested that grammars are shaped by a trade-off between simplicity and communicative efficiency, and Rosch has suggested that category systems “provide maximum information with the least cognitive effort” [p. 190 of (23)].

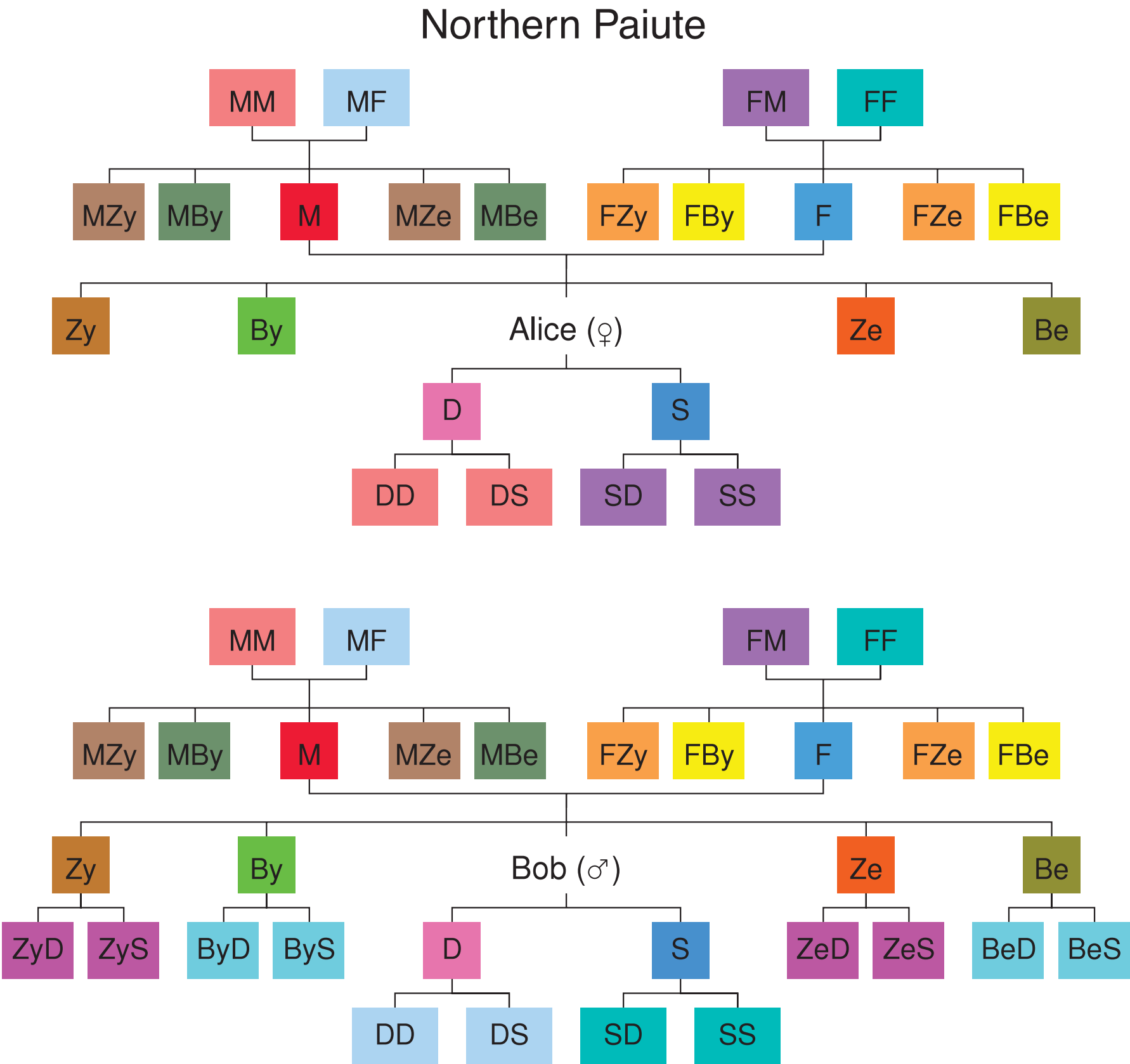
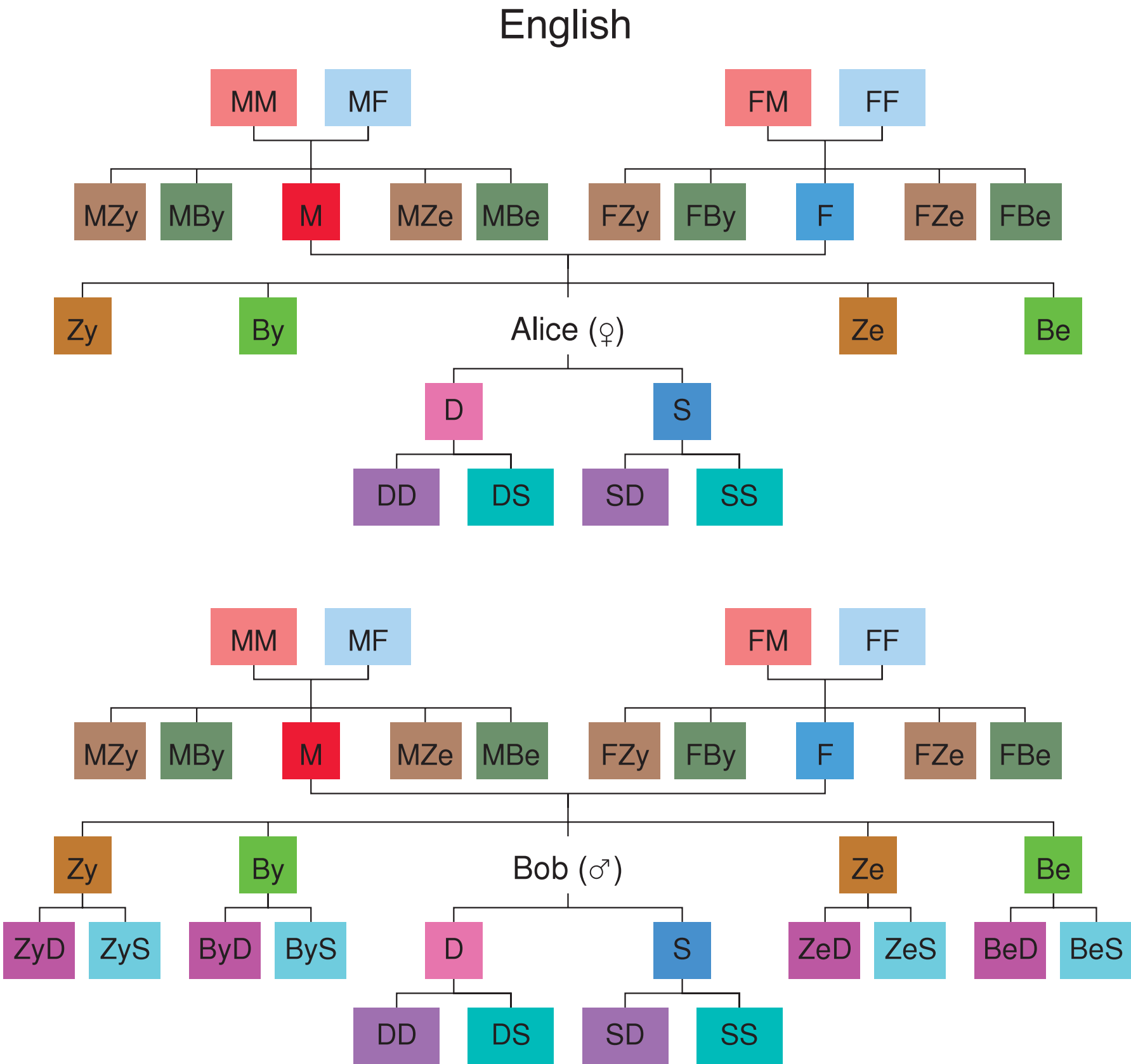
Figure 1A shows a simple communication game that helps to illustrate how kin classification systems are shaped by the principles of simplicity and informativeness. The speaker has a specific relative in mind and utters the category label for that relative. Upon hearing this category label, the hearer must guess which relative the speaker had in mind. Speaker and hearer communicate through



PARENT(x, y) ∧ FEMALE(x)  
PARENT(x, y) ∧ MALE(x)  
PARENT(x, y) ∧ FEMALE(x)  
PARENT(x, y) ∧ MALE(x)  
PARENT(x, z) ∧ PARENT(z, y)



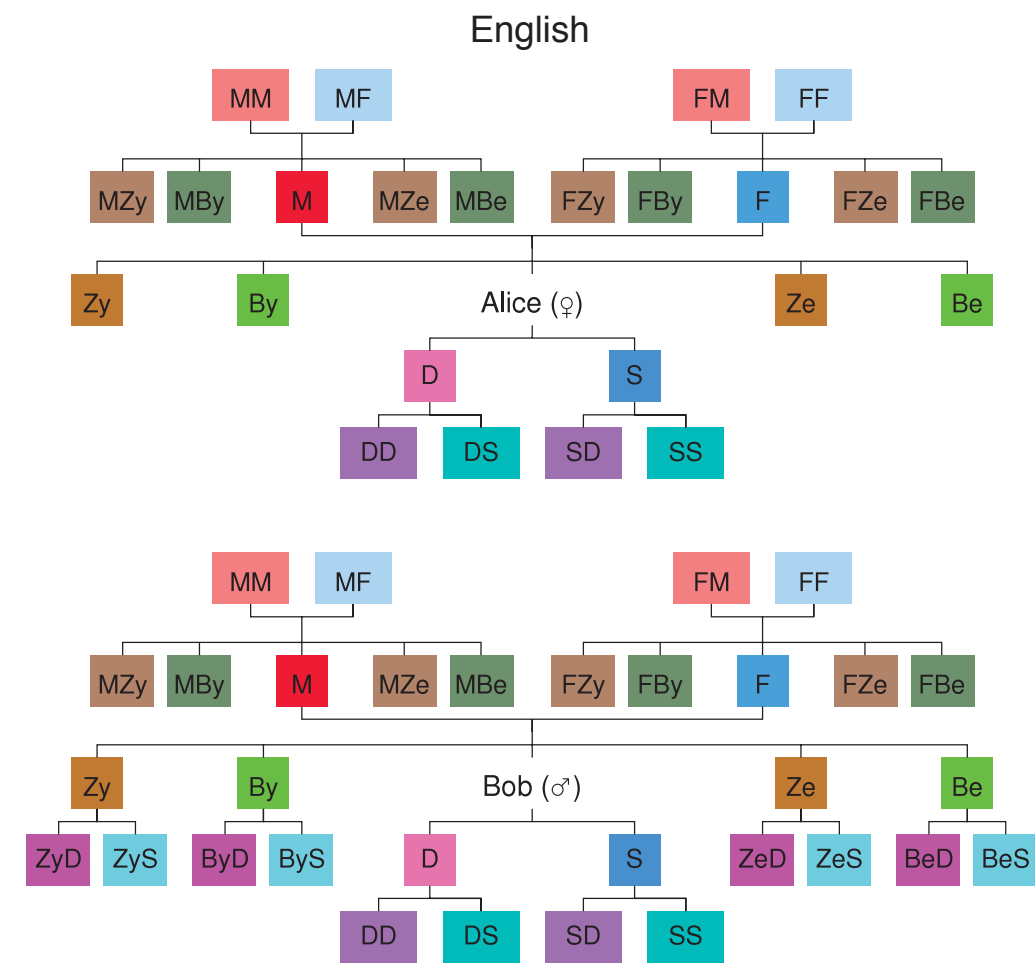
# Kinship terms are simple and informative



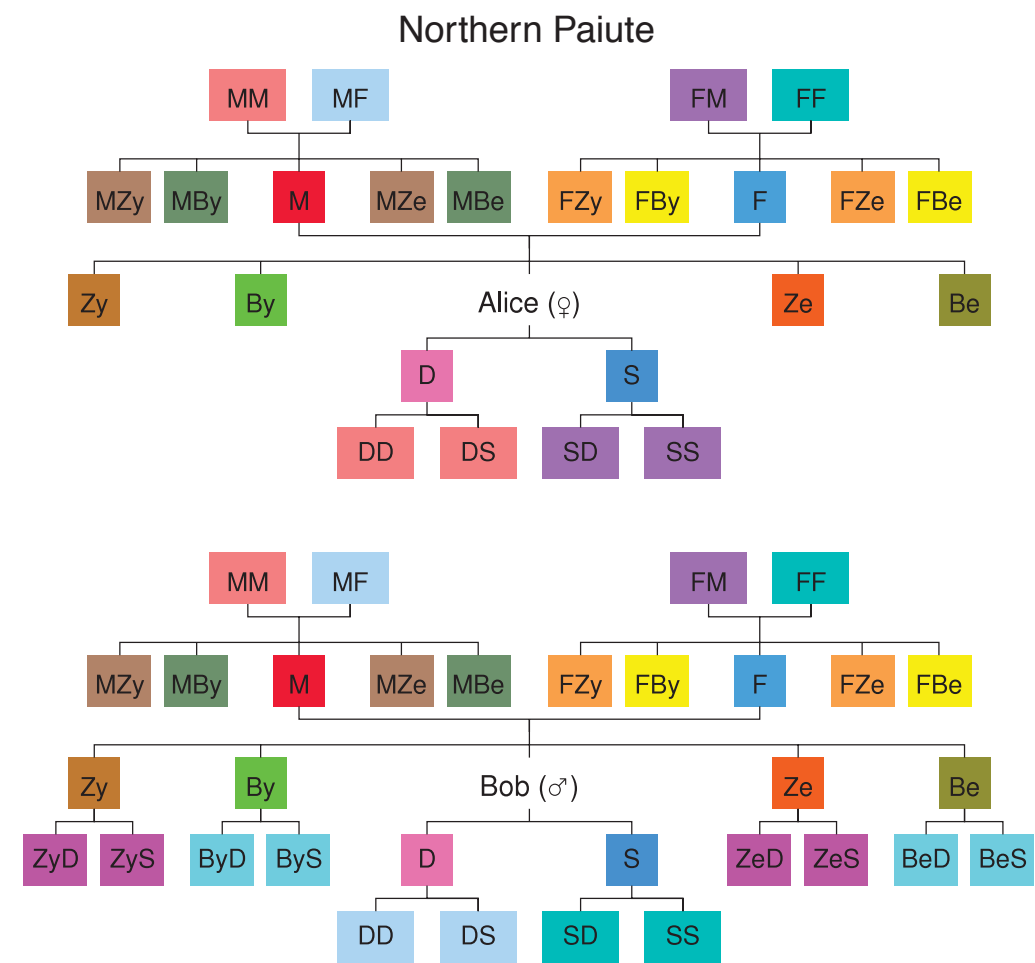
mother(x, y)  $\leftrightarrow$  PARENT(x, y)  $\wedge$  FEMALE(x)  
father(x, y)  $\leftrightarrow$  PARENT(x, y)  $\wedge$  MALE(x)  
daughter(x, y)  $\leftrightarrow$  CHILD(x, y)  $\wedge$  FEMALE(x)  
son(x, y)  $\leftrightarrow$  CHILD(x, y)  $\wedge$  MALE(x)  
sister(x, y)  $\leftrightarrow \exists z$  daughter(x, z)  $\wedge$  PARENT(z, y)

mother(x, y)  $\leftrightarrow$  PARENT(x, y)  $\wedge$  FEMALE(x)  
father(x, y)  $\leftrightarrow$  PARENT(x, y)  $\wedge$  MALE(x)  
daughter(x, y)  $\leftrightarrow$  CHILD(x, y)  $\wedge$  FEMALE(x)  
son(x, y)  $\leftrightarrow$  CHILD(x, y)  $\wedge$  MALE(x)  
sister(x, y)  $\leftrightarrow \exists z$  daughter(x, z)  $\wedge$  PARENT(z, y)

# Kinship terms are simple and informative

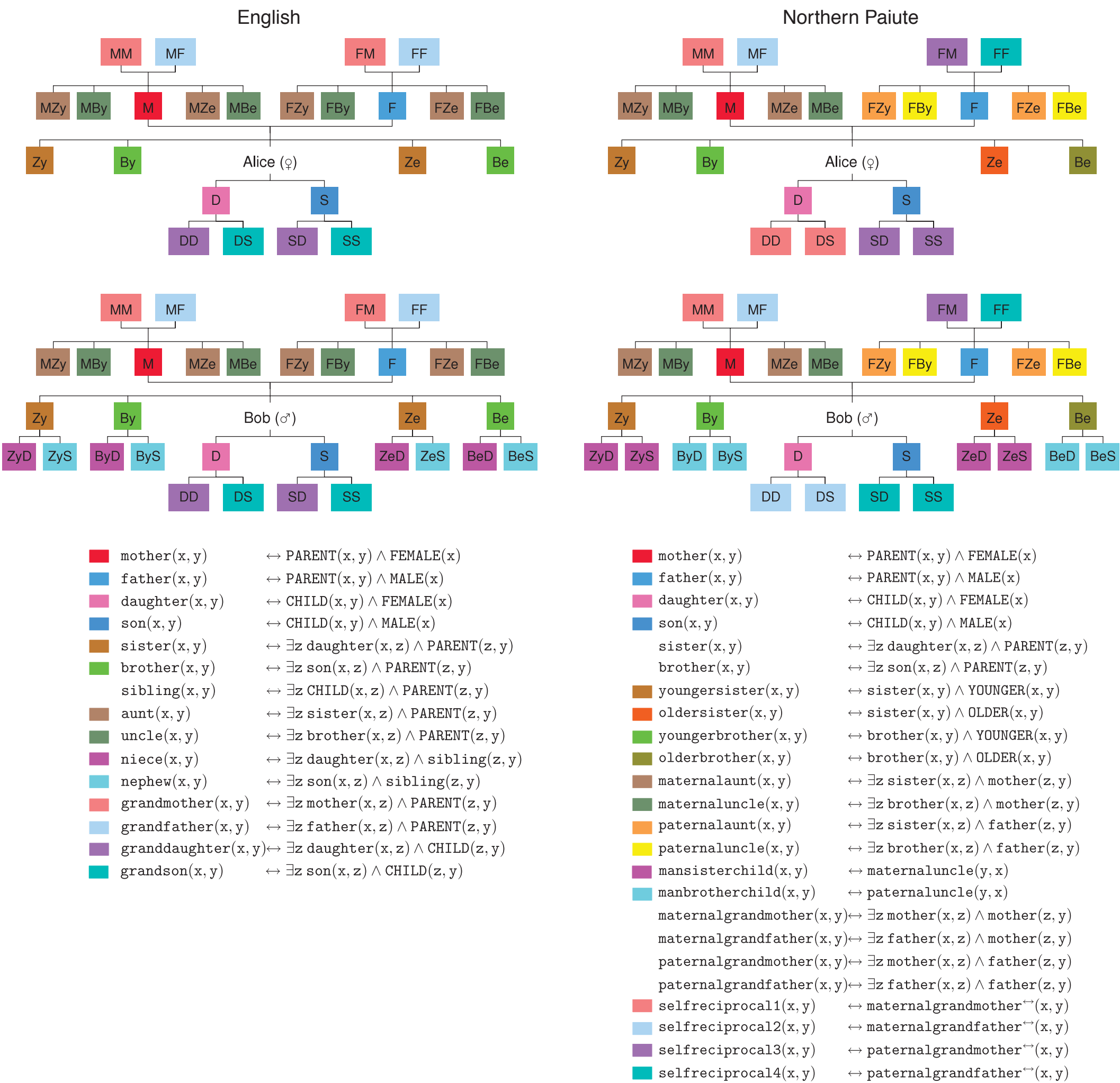


<span style="color: red;">■</span> mother(x, y)	$\leftrightarrow \text{PARENT}(x, y) \wedge \text{FEMALE}(x)$
<span style="color: blue;">■</span> father(x, y)	$\leftrightarrow \text{PARENT}(x, y) \wedge \text{MALE}(x)$
<span style="color: pink;">■</span> daughter(x, y)	$\leftrightarrow \text{CHILD}(x, y) \wedge \text{FEMALE}(x)$
<span style="color: blue;">■</span> son(x, y)	$\leftrightarrow \text{CHILD}(x, y) \wedge \text{MALE}(x)$
<span style="color: brown;">■</span> sister(x, y)	$\leftrightarrow \exists z \text{ daughter}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: green;">■</span> brother(x, y)	$\leftrightarrow \exists z \text{ son}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: brown;">■</span> sibling(x, y)	$\leftrightarrow \exists z \text{ CHILD}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: brown;">■</span> aunt(x, y)	$\leftrightarrow \exists z \text{ sister}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: green;">■</span> uncle(x, y)	$\leftrightarrow \exists z \text{ brother}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: pink;">■</span> niece(x, y)	$\leftrightarrow \exists z \text{ daughter}(x, z) \wedge \text{sibling}(z, y)$
<span style="color: teal;">■</span> nephew(x, y)	$\leftrightarrow \exists z \text{ son}(x, z) \wedge \text{sibling}(z, y)$
<span style="color: red;">■</span> grandmother(x, y)	$\leftrightarrow \exists z \text{ mother}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: blue;">■</span> grandfather(x, y)	$\leftrightarrow \exists z \text{ father}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: purple;">■</span> granddaughter(x, y)	$\leftrightarrow \exists z \text{ daughter}(x, z) \wedge \text{CHILD}(z, y)$
<span style="color: teal;">■</span> grandson(x, y)	$\leftrightarrow \exists z \text{ son}(x, z) \wedge \text{CHILD}(z, y)$

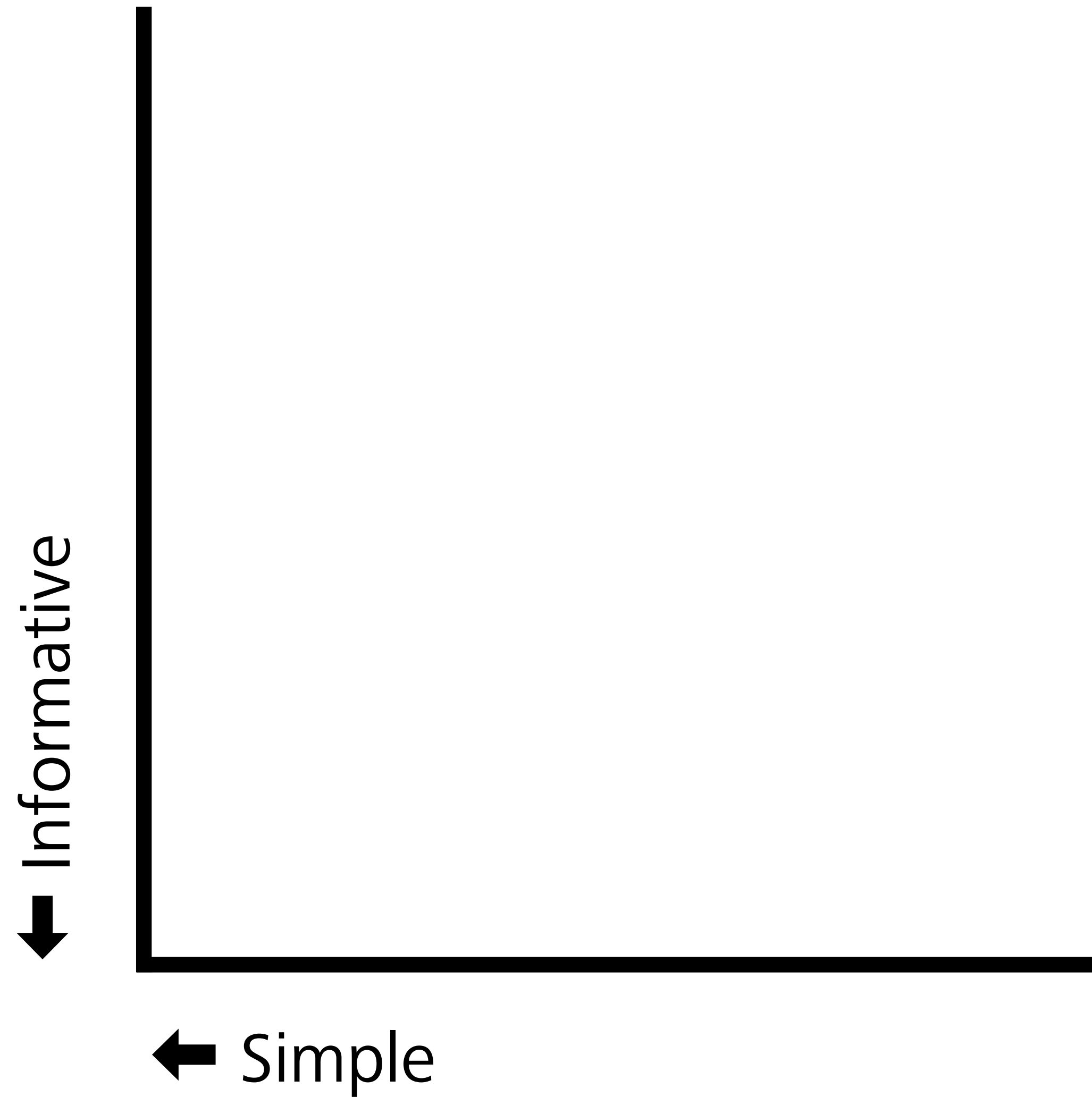


<span style="color: red;">■</span> mother(x, y)	$\leftrightarrow \text{PARENT}(x, y) \wedge \text{FEMALE}(x)$
<span style="color: blue;">■</span> father(x, y)	$\leftrightarrow \text{PARENT}(x, y) \wedge \text{MALE}(x)$
<span style="color: pink;">■</span> daughter(x, y)	$\leftrightarrow \text{CHILD}(x, y) \wedge \text{FEMALE}(x)$
<span style="color: blue;">■</span> son(x, y)	$\leftrightarrow \text{CHILD}(x, y) \wedge \text{MALE}(x)$
<span style="color: brown;">■</span> sister(x, y)	$\leftrightarrow \exists z \text{ daughter}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: green;">■</span> brother(x, y)	$\leftrightarrow \exists z \text{ son}(x, z) \wedge \text{PARENT}(z, y)$
<span style="color: brown;">■</span> youngersister(x, y)	$\leftrightarrow \text{sister}(x, y) \wedge \text{YOUNGER}(x, y)$
<span style="color: orange;">■</span> oldersister(x, y)	$\leftrightarrow \text{sister}(x, y) \wedge \text{OLDER}(x, y)$
<span style="color: green;">■</span> youngerbrother(x, y)	$\leftrightarrow \text{brother}(x, y) \wedge \text{YOUNGER}(x, y)$
<span style="color: olive;">■</span> olderbrother(x, y)	$\leftrightarrow \text{brother}(x, y) \wedge \text{OLDER}(x, y)$
<span style="color: brown;">■</span> maternal aunt(x, y)	$\leftrightarrow \exists z \text{ sister}(x, z) \wedge \text{mother}(z, y)$
<span style="color: green;">■</span> maternal uncle(x, y)	$\leftrightarrow \exists z \text{ brother}(x, z) \wedge \text{mother}(z, y)$
<span style="color: orange;">■</span> paternal aunt(x, y)	$\leftrightarrow \exists z \text{ sister}(x, z) \wedge \text{father}(z, y)$
<span style="color: yellow;">■</span> paternal uncle(x, y)	$\leftrightarrow \exists z \text{ brother}(x, z) \wedge \text{father}(z, y)$
<span style="color: pink;">■</span> mansisterchild(x, y)	$\leftrightarrow \text{maternaluncle}(y, x)$
<span style="color: teal;">■</span> manbrotherchild(x, y)	$\leftrightarrow \text{paternaluncle}(y, x)$
<span style="color: brown;">■</span> maternal grandmother(x, y)	$\leftrightarrow \exists z \text{ mother}(x, z) \wedge \text{mother}(z, y)$
<span style="color: brown;">■</span> maternal grandfather(x, y)	$\leftrightarrow \exists z \text{ father}(x, z) \wedge \text{mother}(z, y)$
<span style="color: brown;">■</span> paternal grandmother(x, y)	$\leftrightarrow \exists z \text{ mother}(x, z) \wedge \text{father}(z, y)$
<span style="color: brown;">■</span> paternal grandfather(x, y)	$\leftrightarrow \exists z \text{ father}(x, z) \wedge \text{father}(z, y)$
<span style="color: red;">■</span> selfreciprocal1(x, y)	$\leftrightarrow \text{maternalgrandmother}^{\leftrightarrow}(x, y)$
<span style="color: blue;">■</span> selfreciprocal2(x, y)	$\leftrightarrow \text{maternalgrandfather}^{\leftrightarrow}(x, y)$
<span style="color: purple;">■</span> selfreciprocal3(x, y)	$\leftrightarrow \text{paternalgrandmother}^{\leftrightarrow}(x, y)$
<span style="color: teal;">■</span> selfreciprocal4(x, y)	$\leftrightarrow \text{paternalgrandfather}^{\leftrightarrow}(x, y)$

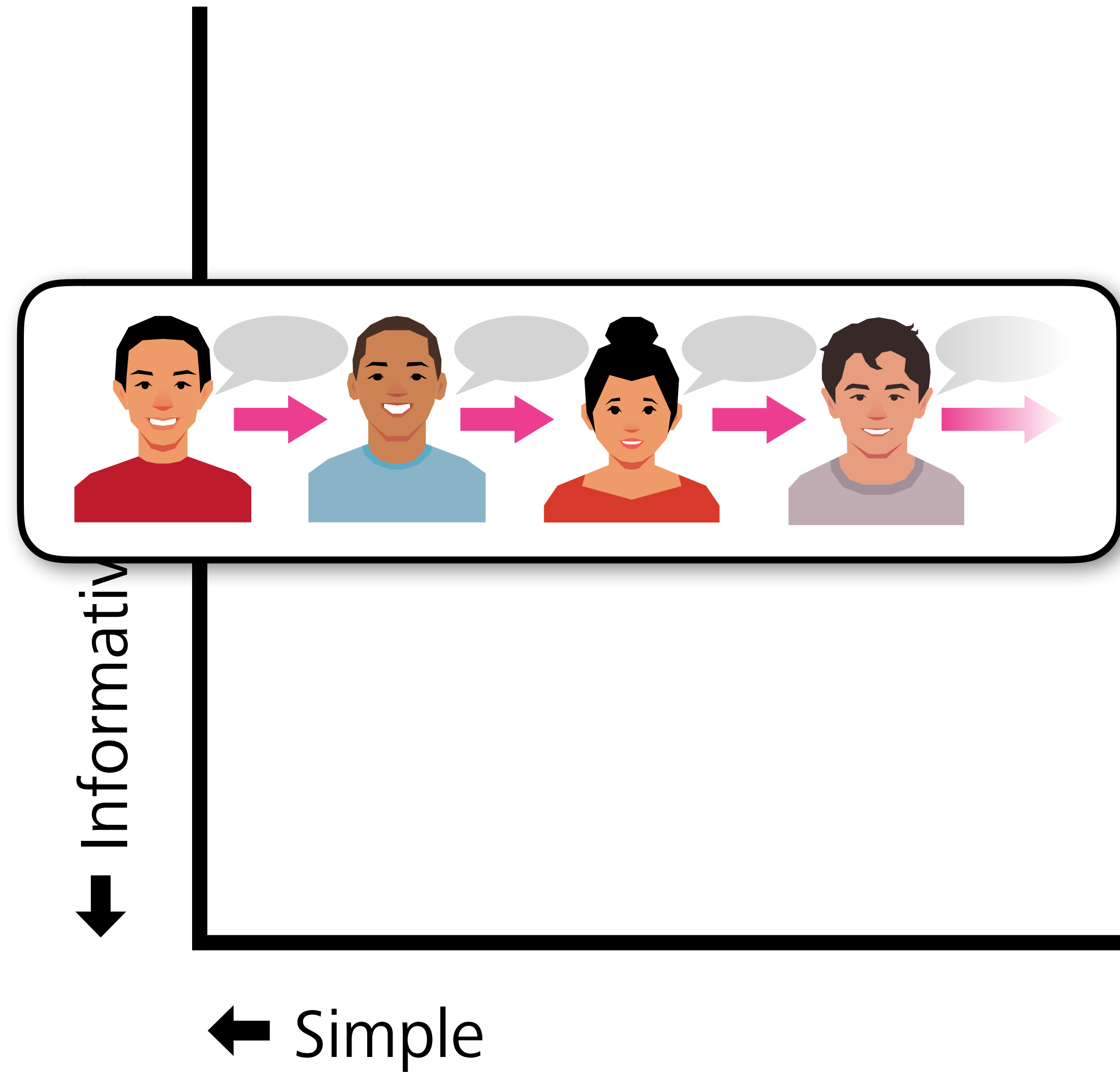
# Kinship terms are simple and informative



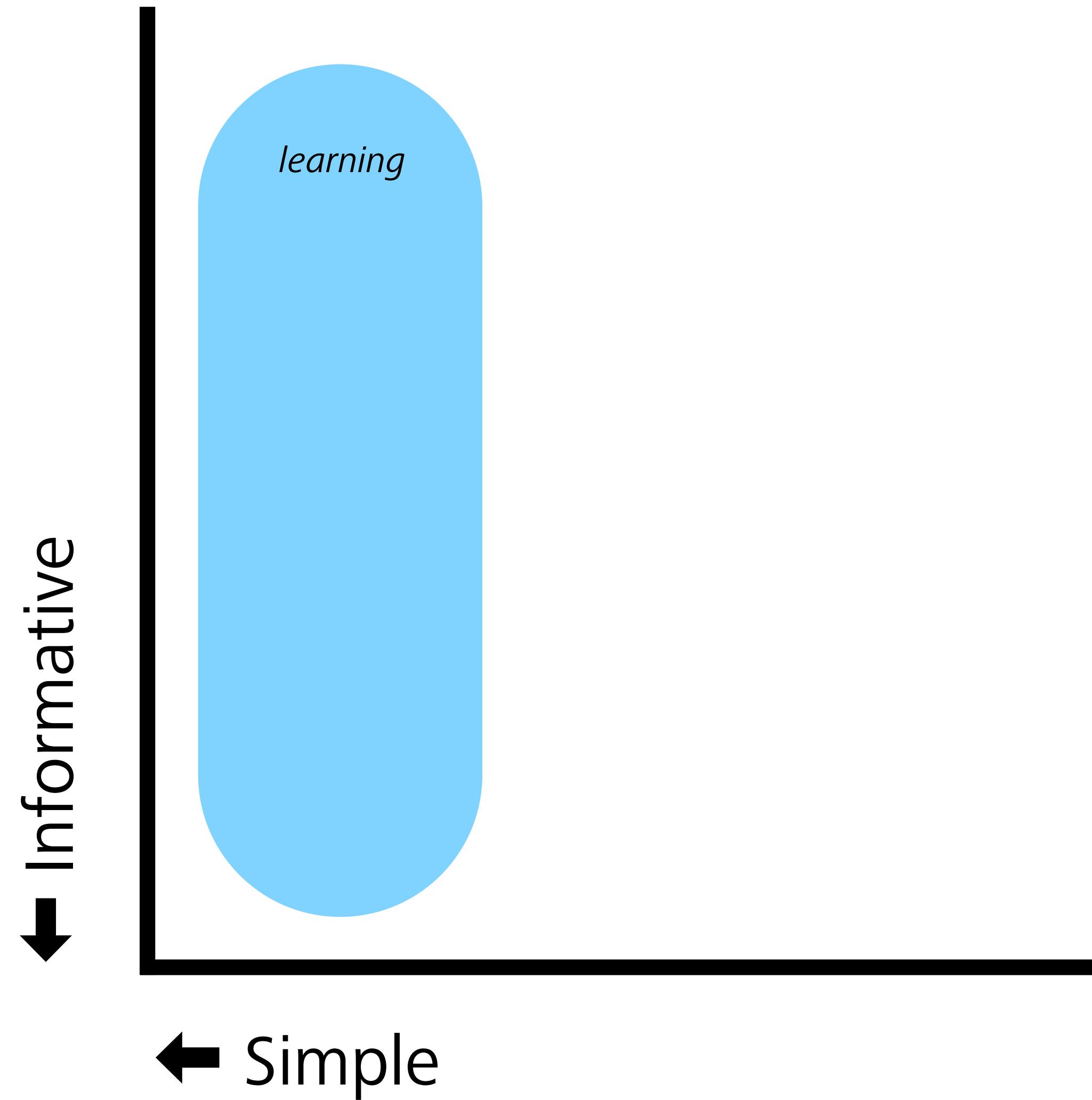
# Learning and interaction pressures



# Learning and interaction pressures

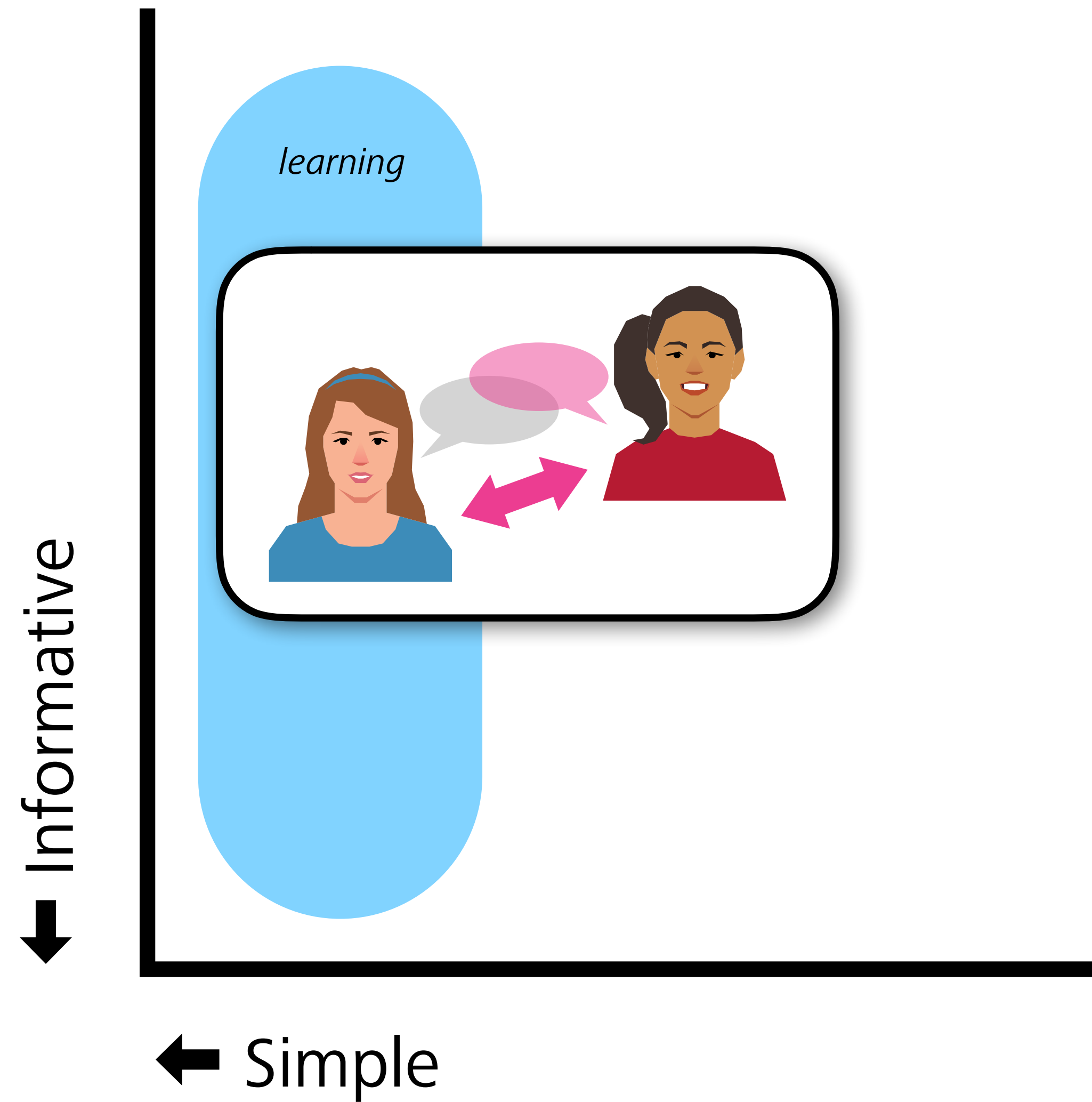


# Learning and interaction pressures

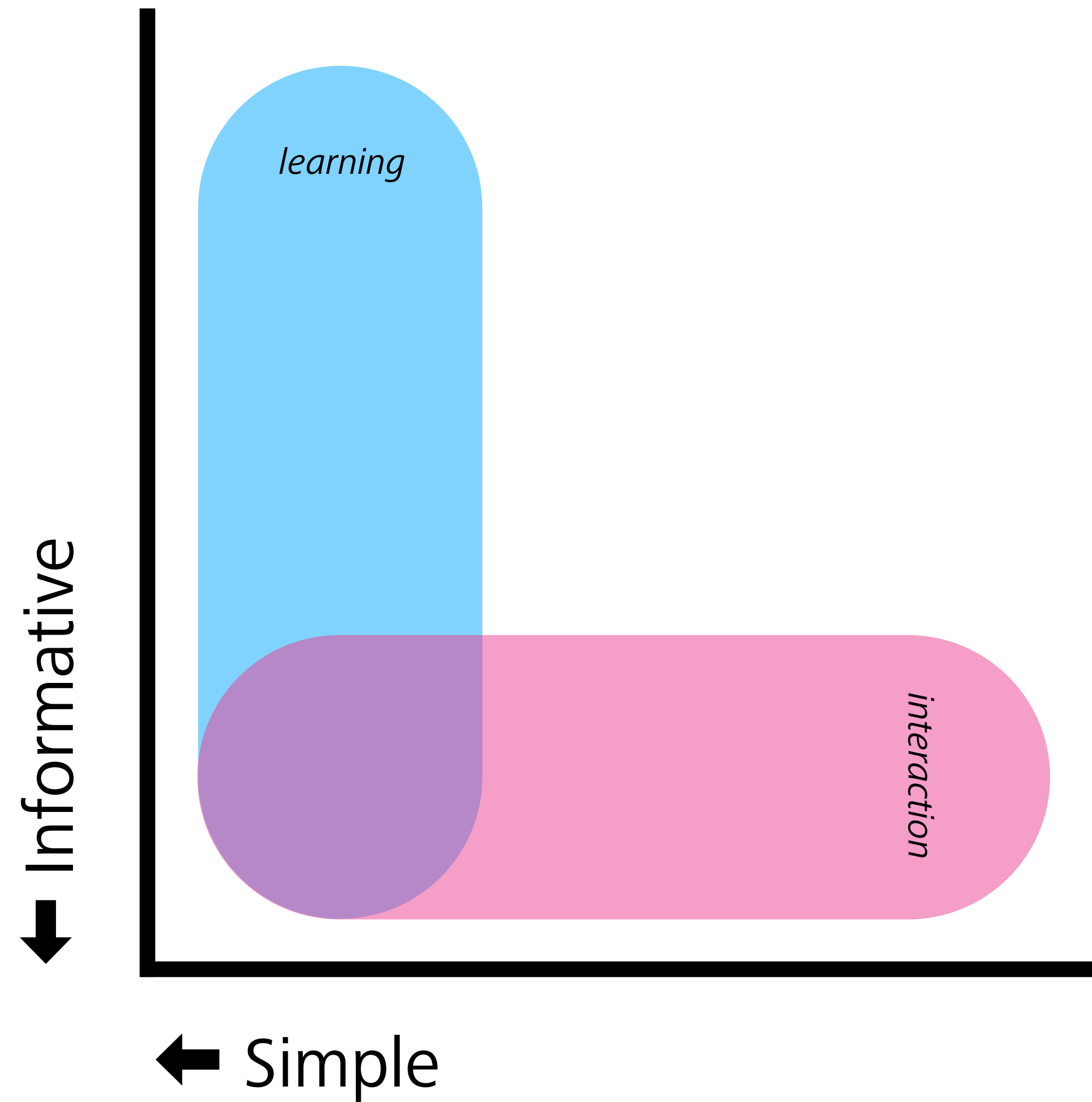




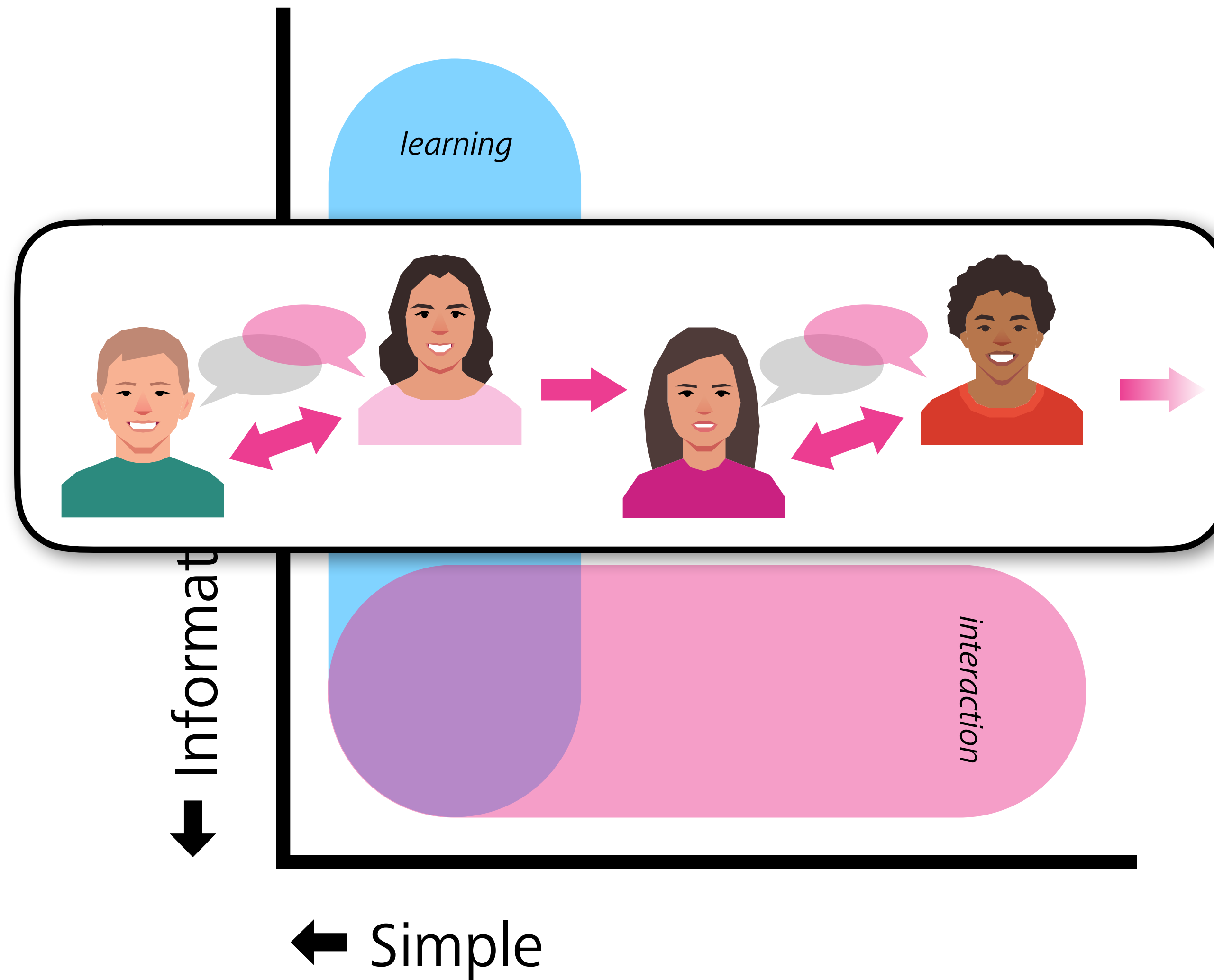
# Learning and interaction pressures



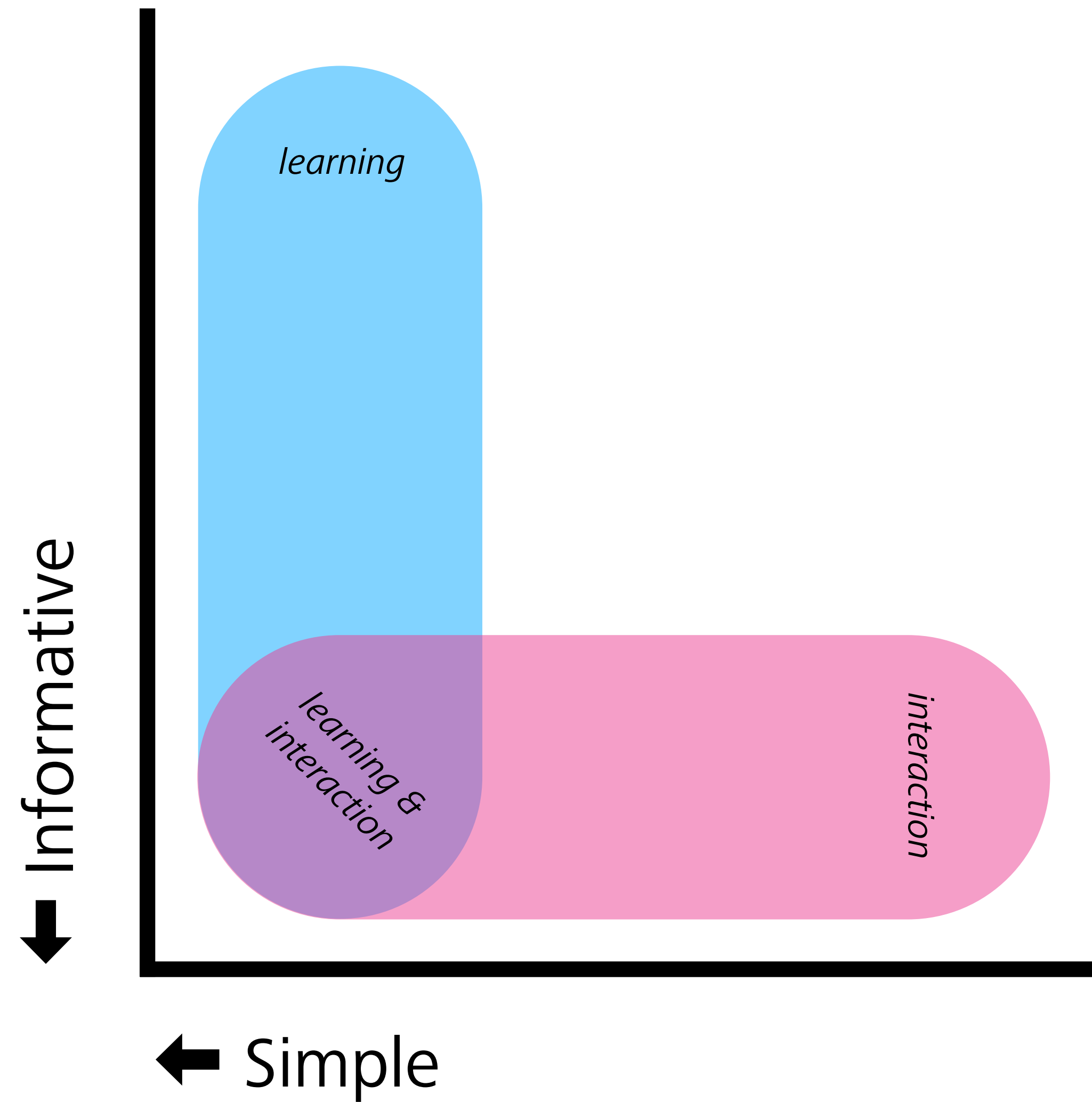
# Learning and interaction pressures



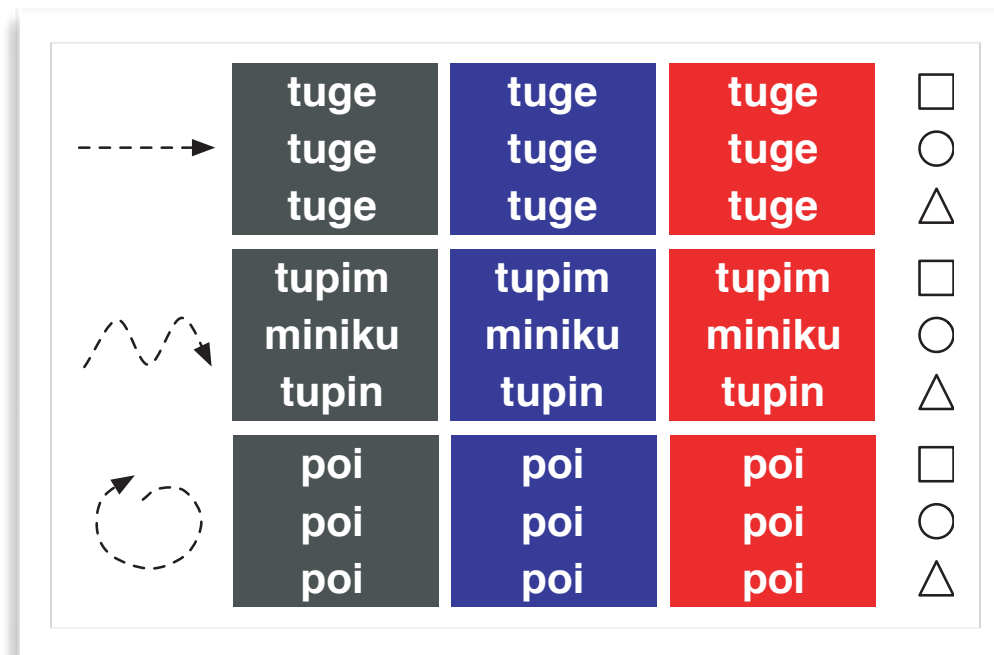
# Learning and interaction pressures



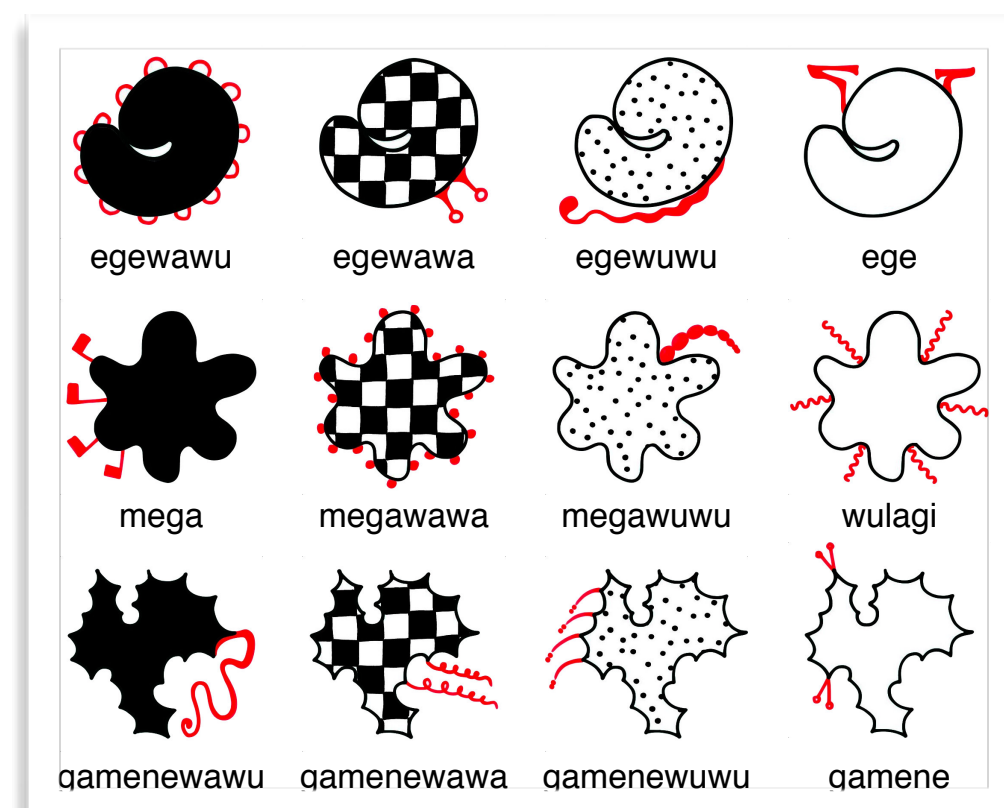
# Learning and interaction pressures



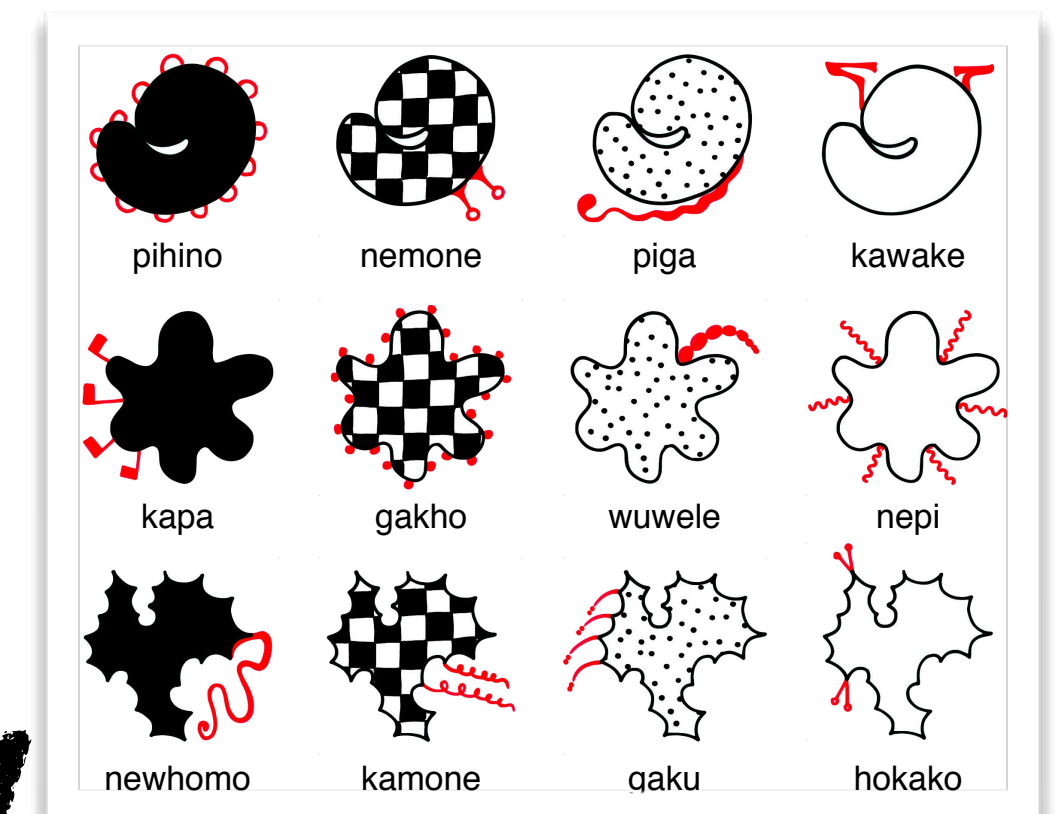
# Learning and interaction pressures



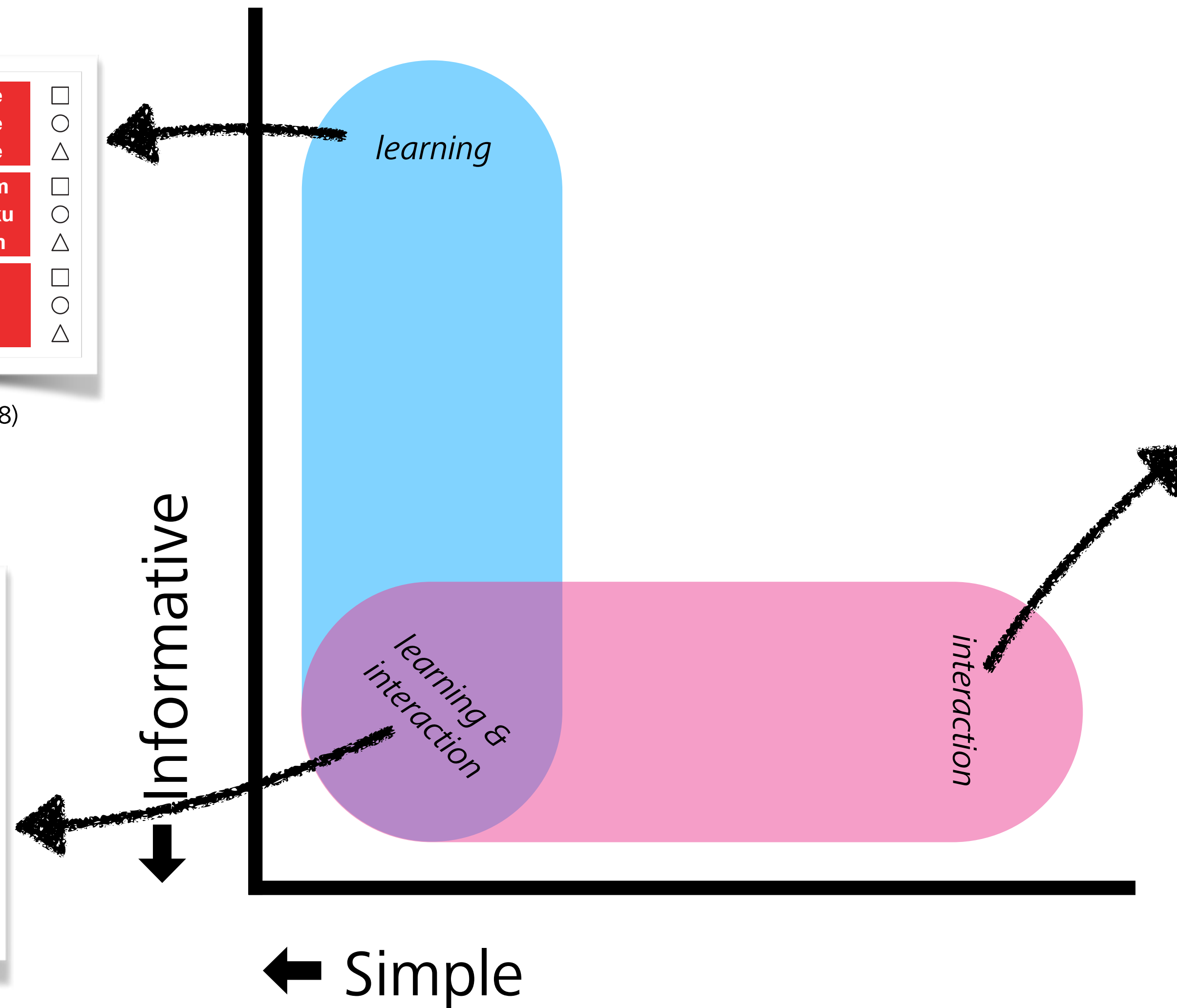
Kirby, Cornish, & Smith (2008)



Kirby, Tamariz, Cornish, & Smith (2015)



Kirby, Tamariz, Cornish, & Smith (2015)

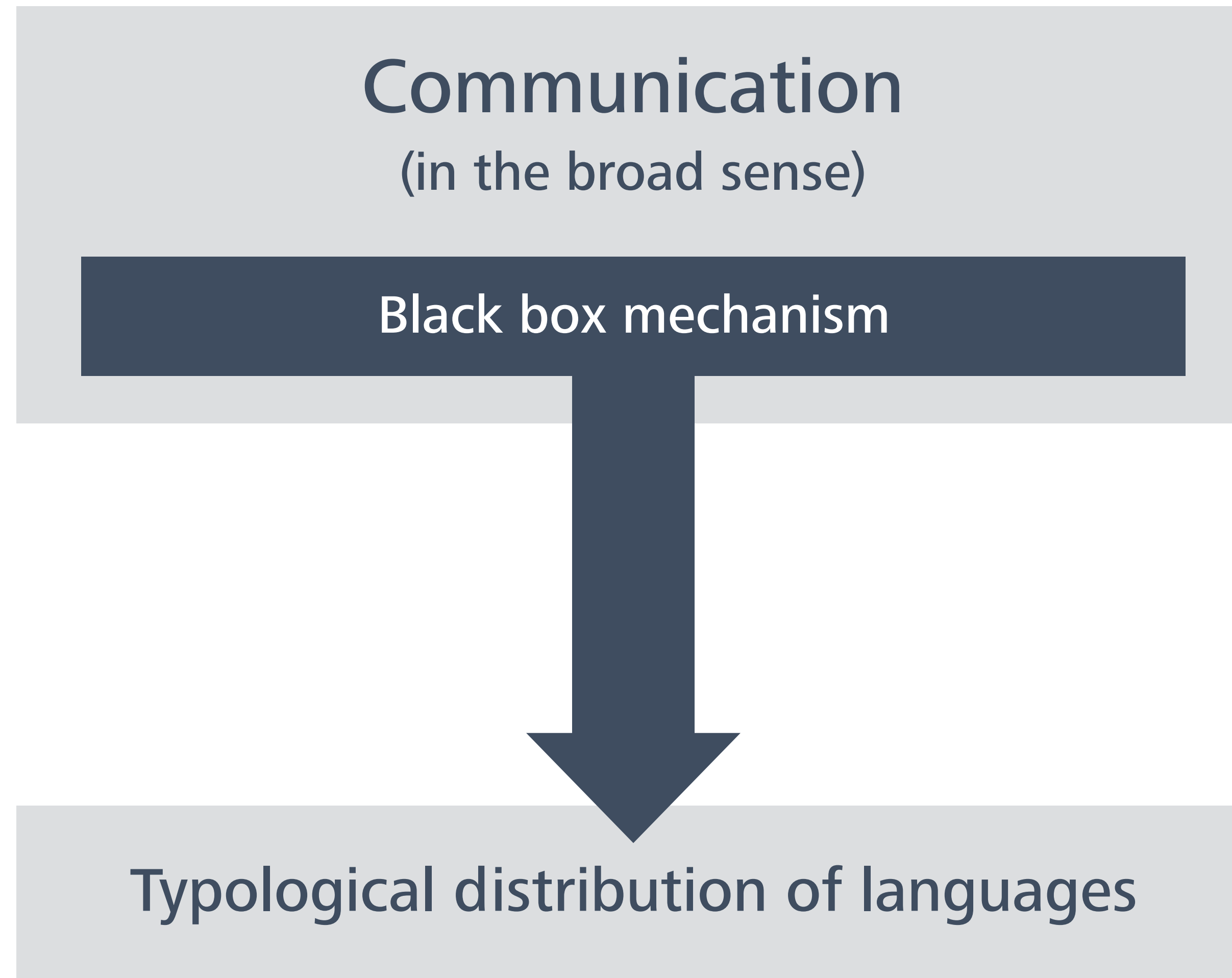


# The problem of linkage

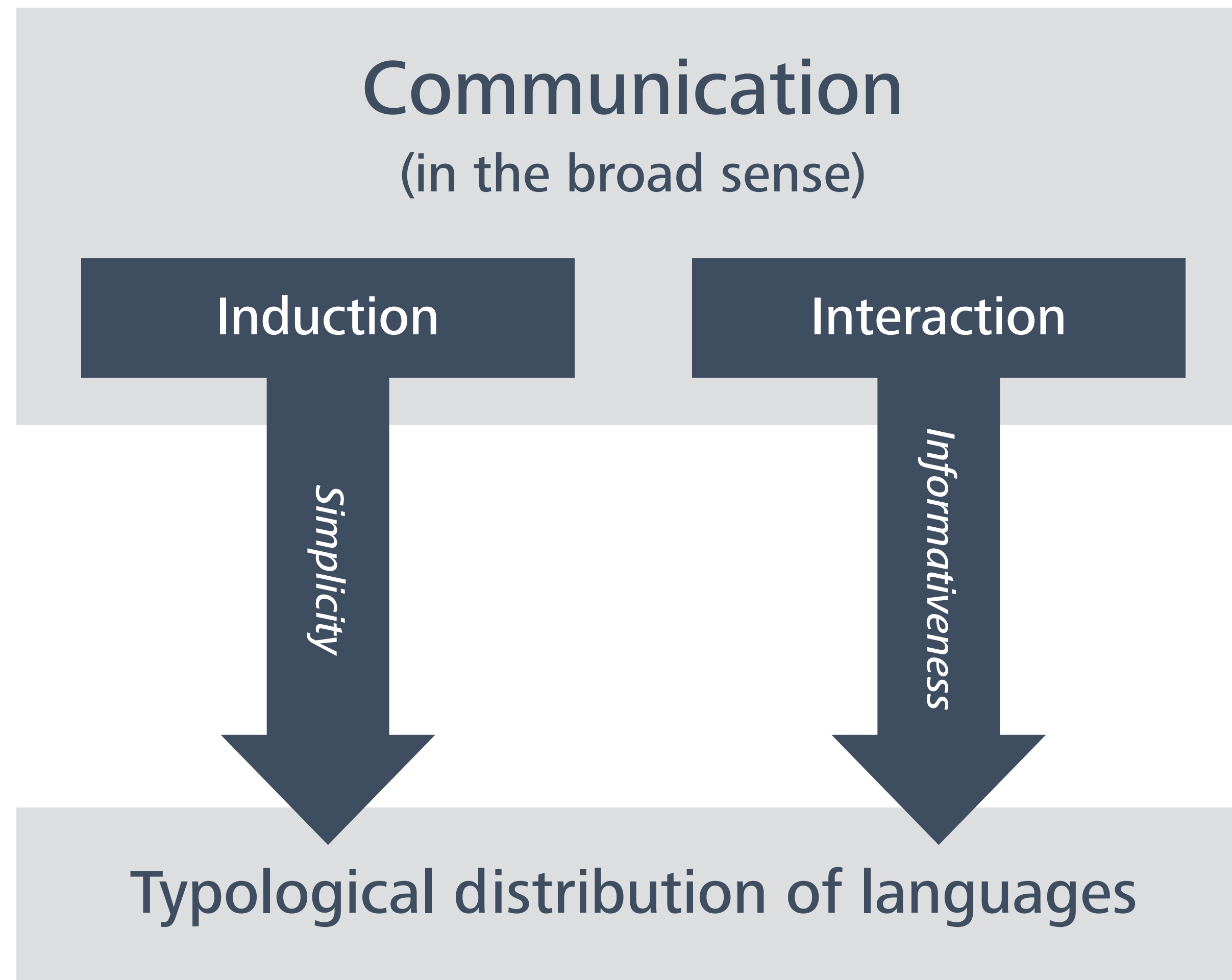
**Communication**  
(in the broad sense)

Typological distribution of languages

# The problem of linkage



# The problem of linkage





*Induction*

*as the pressure for simplicity*

# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

$$\text{posterior}(H|D) \propto \text{likelihood}(D|H) \times \text{prior}(H)$$

# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

$$\text{posterior}(H|D) \propto \text{likelihood}(D|H) \times 2^{-\text{DL}(H)}$$

# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

$$\text{posterior}(H|D) \propto \text{likelihood}(D|H) \times 2^{-\text{DL}(H)}$$

Any regularities in data can be used to compress that data

# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

$$\text{posterior}(H|D) \propto \text{likelihood}(D|H) \times 2^{-\text{DL}(H)}$$

Any regularities in data can be used to compress that data

The more regularities there are, the more the data can be compressed

# The Minimum Description Length principle

$$DL(H|D) = DL(D|H) + DL(H)$$

*For example...*

```
010010111110010000110001000101101100001111010001
```

```
print('010010111110010000110001000101101100001111010001')
```

```
010101010101010101010101010101010101010101010101
```

```
print('0101'*12)      or      print('01'*24)
```

Any regular

The more re

compressed

# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

$$\text{posterior}(H|D) \propto \text{likelihood}(D|H) \times 2^{-\text{DL}(H)}$$

Any regularities in data can be used to compress that data

The more regularities there are, the more the data can be compressed

We equate **learning** with **compression**: The more the data can be compressed, the more insight we gain from that data



# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$

$$\text{posterior}(H|D) \propto \text{likelihood}(D|H) \times 2^{-\text{DL}(H)}$$

Any regularities in data can be used to compress that data

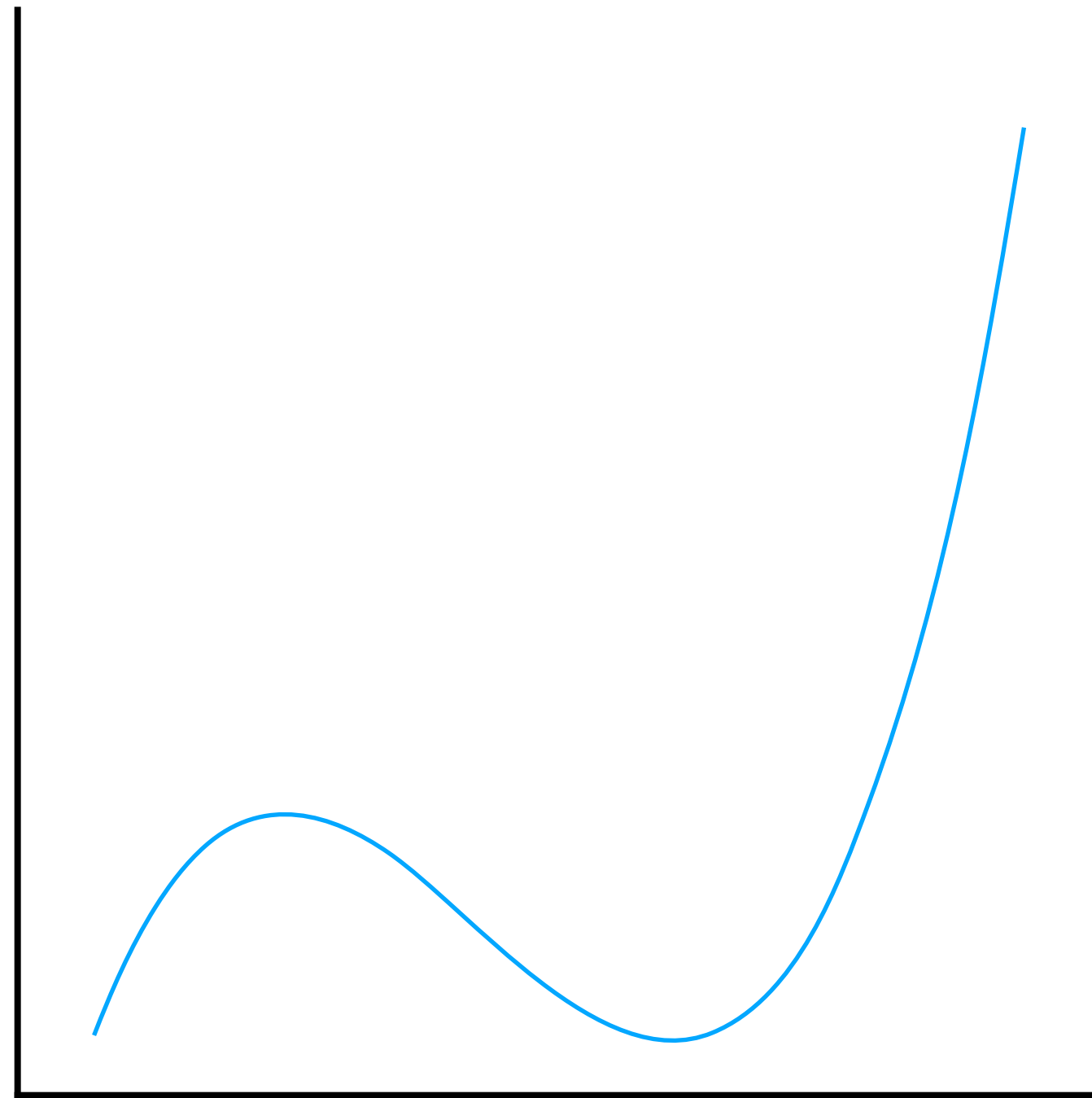
The more regularities there are, the more the data can be compressed

We equate **learning** with **compression**: The more the data can be compressed, the more insight we gain from that data

In other words, the more regularity we can identify, the more we can predict what the generating process will do next

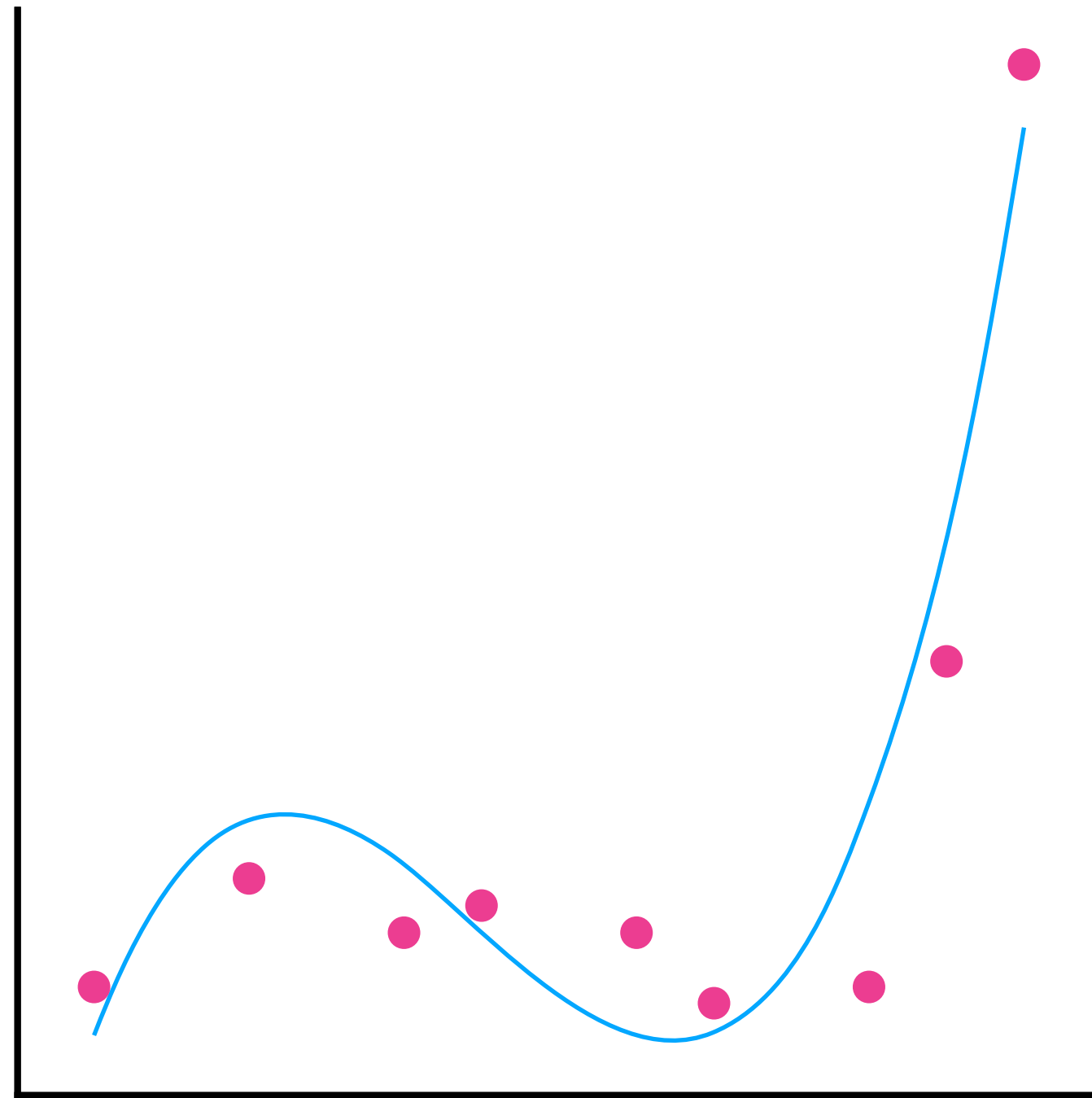
# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$



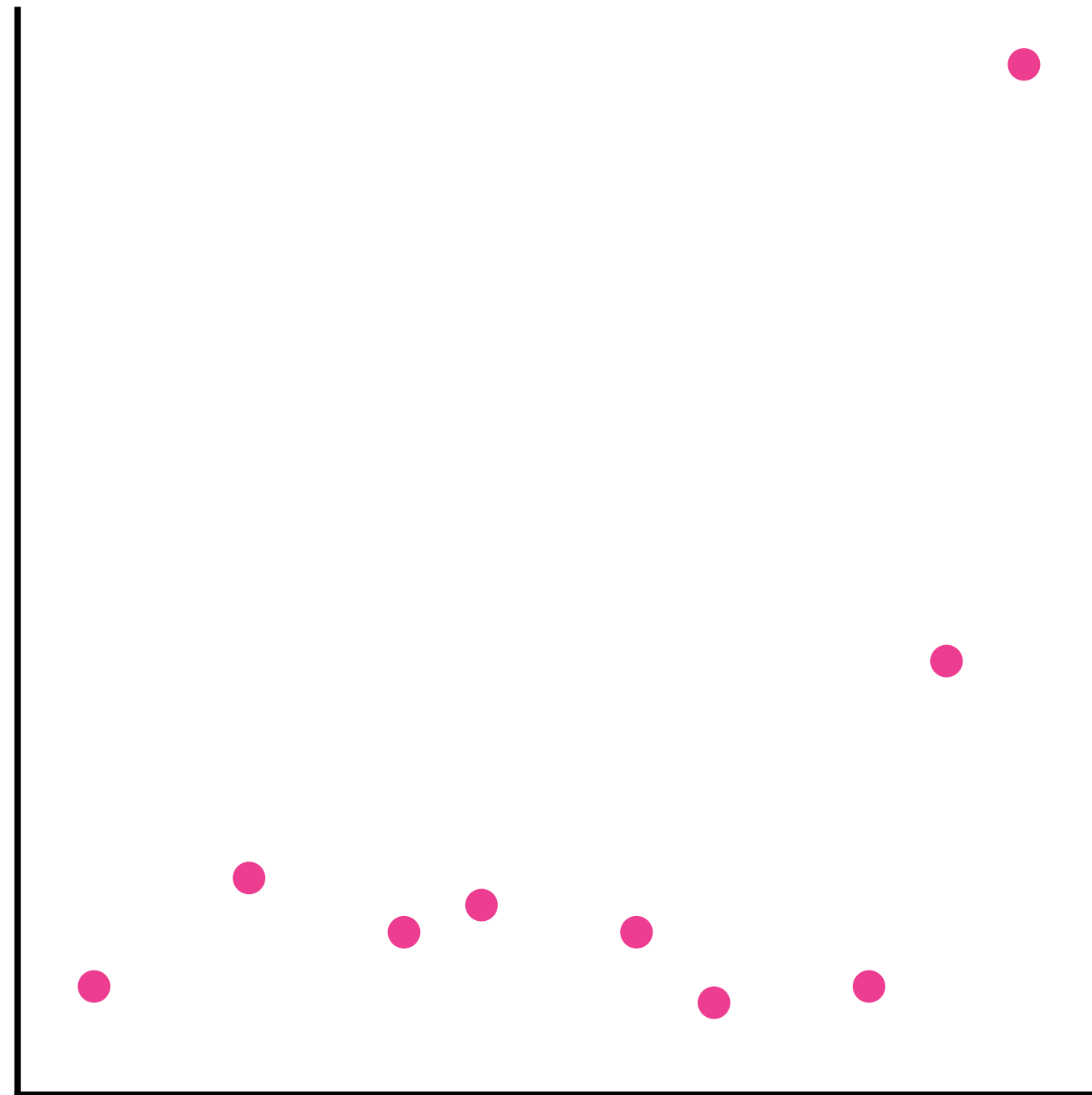
# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$



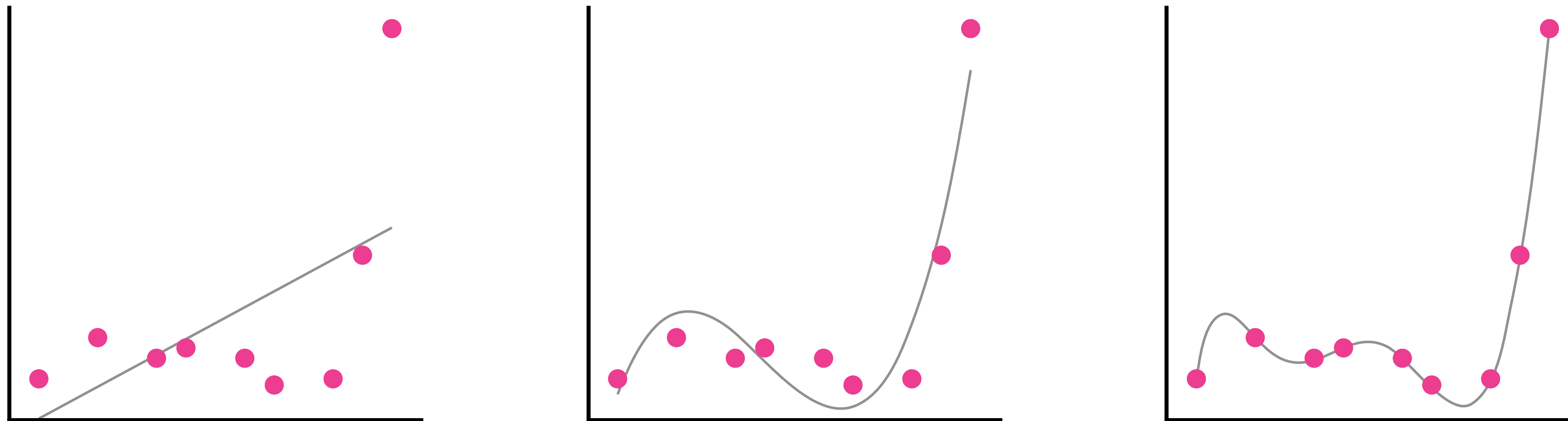
# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$



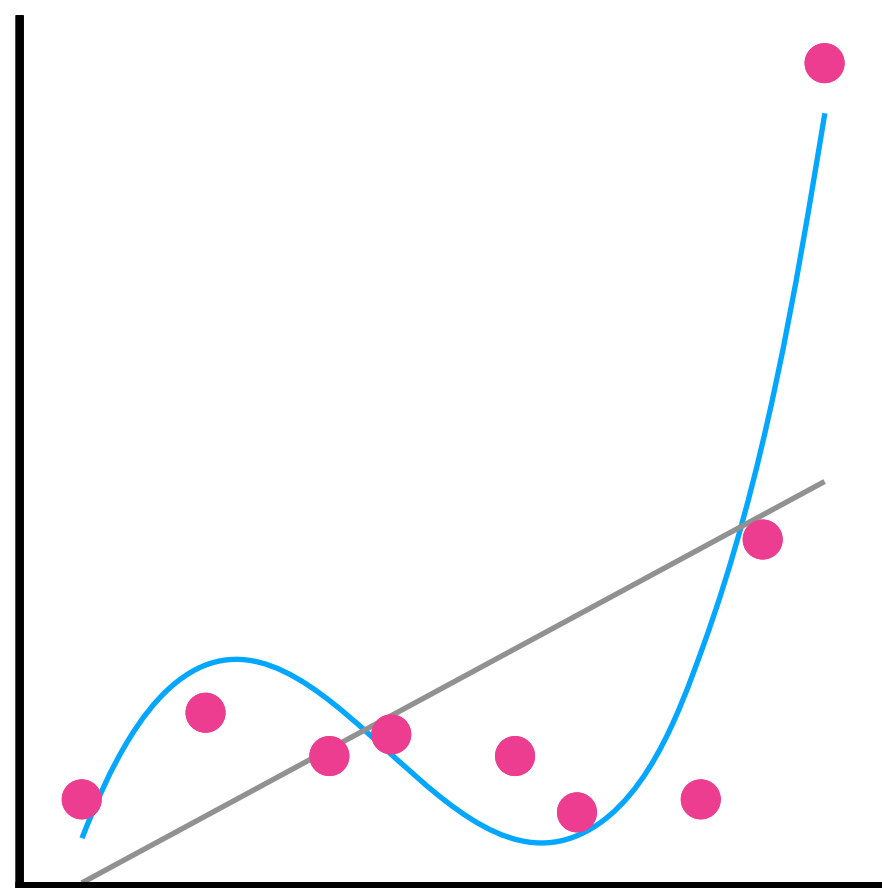
# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$



# The Minimum Description Length principle

$$\text{DL}(H|D) = \text{DL}(D|H) + \text{DL}(H)$$



$\text{DL}(D|H)$

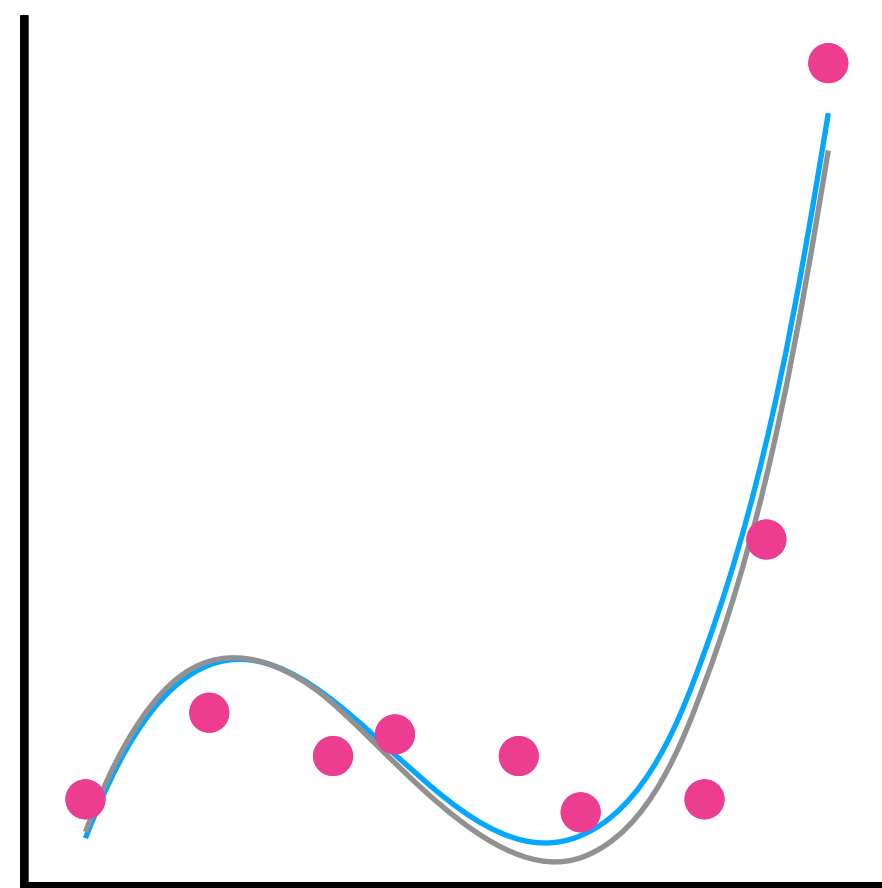
10 bits

$\text{DL}(H)$

1 bit

$\text{DL}(H|D)$

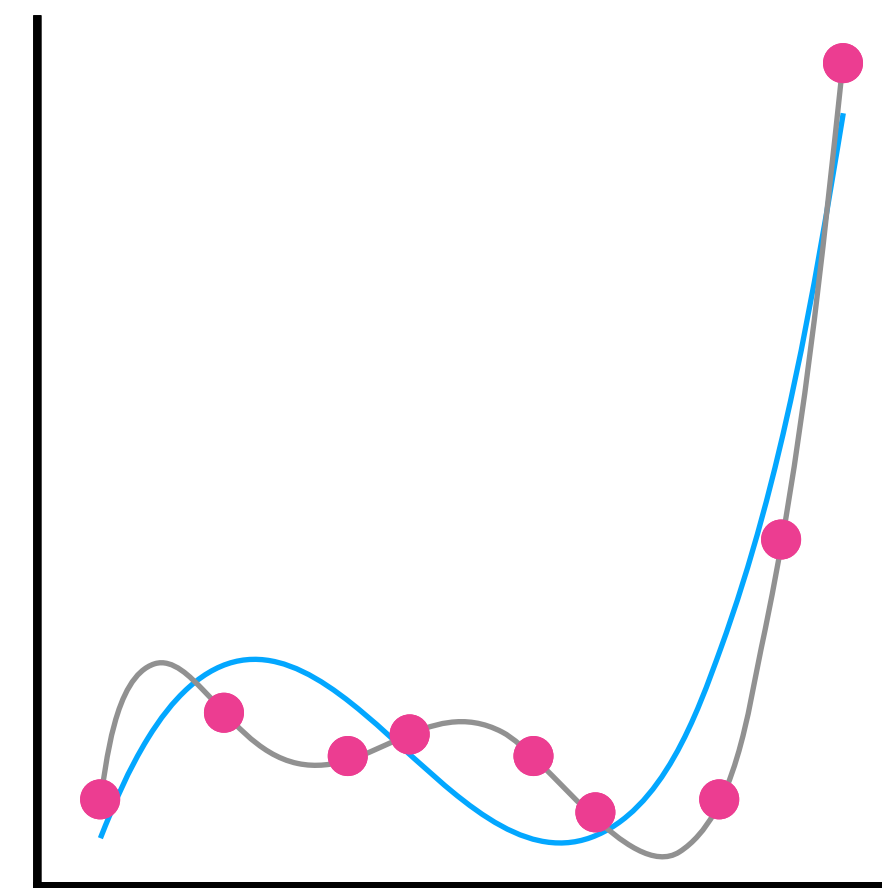
11 bits



4 bits

4 bits

8 bits ✓

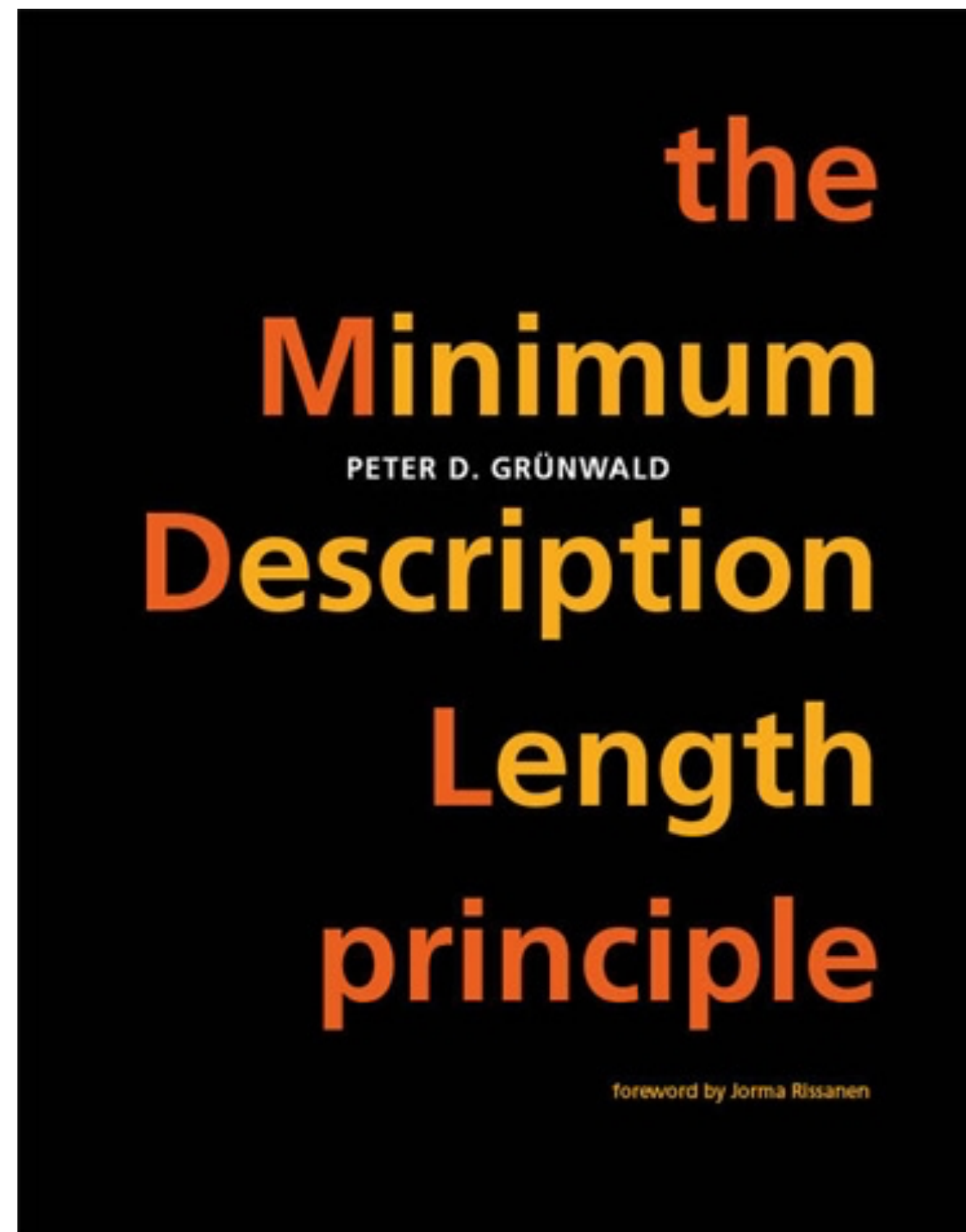


0 bits

10 bits

10 bits

# The Minimum Description Length principle



**Bayesian interpretation:** MDL is closely related to Bayesian inference

**Occam's razor:** MDL trades-off *goodness-of-fit* with *model complexity*, embodying Occam's razor

**No overfitting:** MDL automatically guards against *overfitting noise* in data

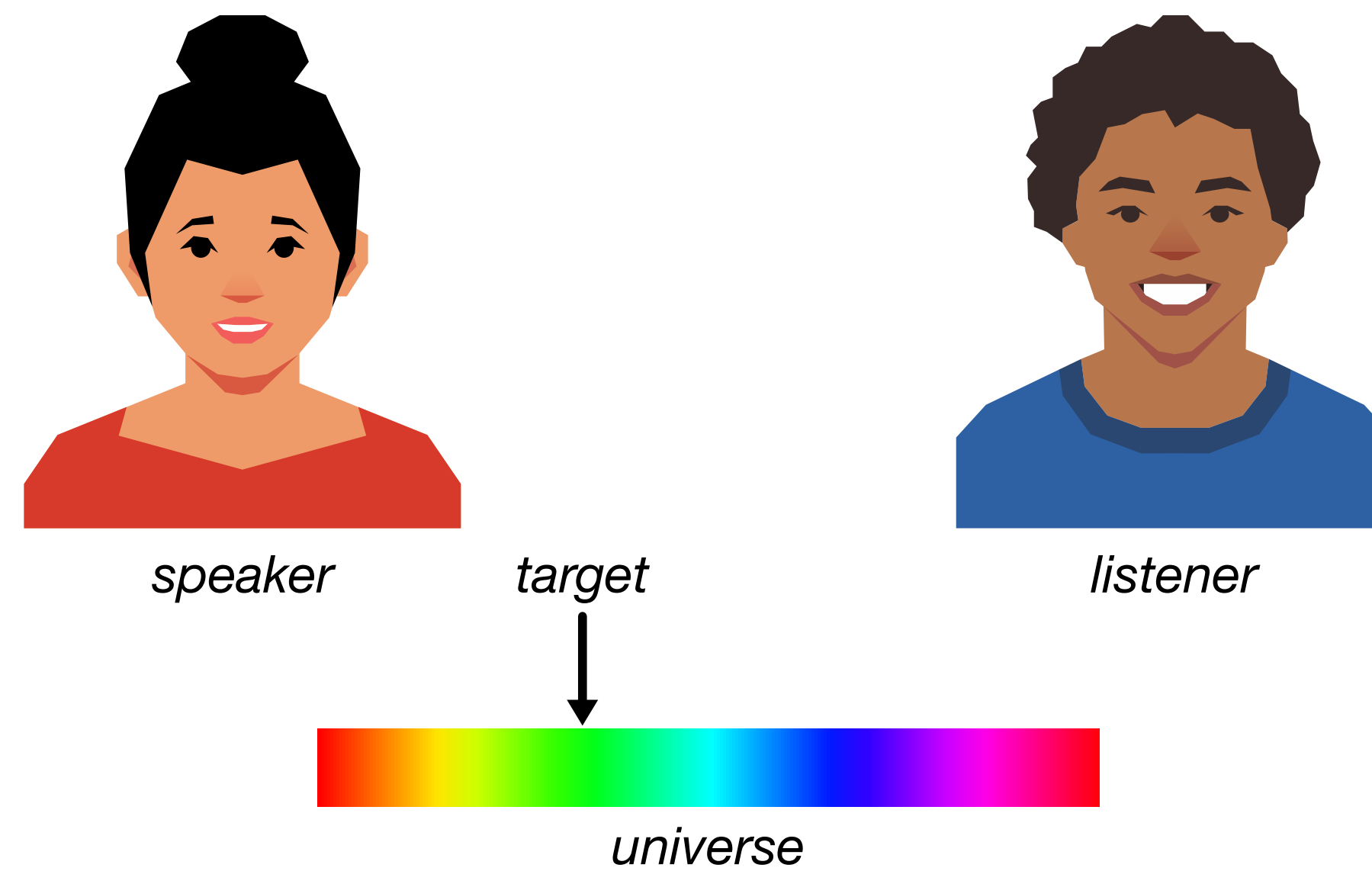
**Predictive performance:** Since data compression is formally equivalent to probabilistic prediction, MDL finds models offering *good predictive performance on unseen data*

*Interaction*

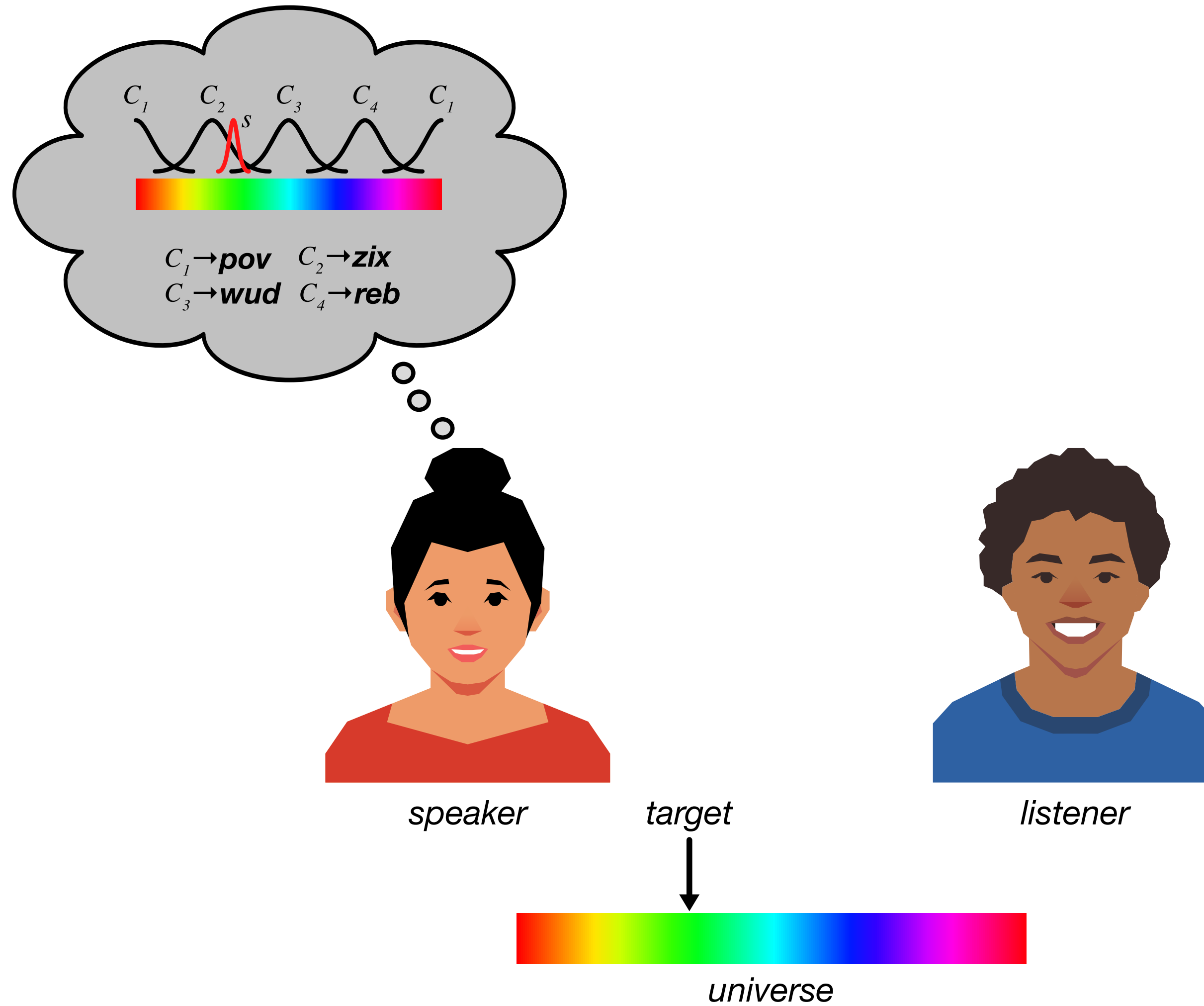
*as the pressure for informativeness*



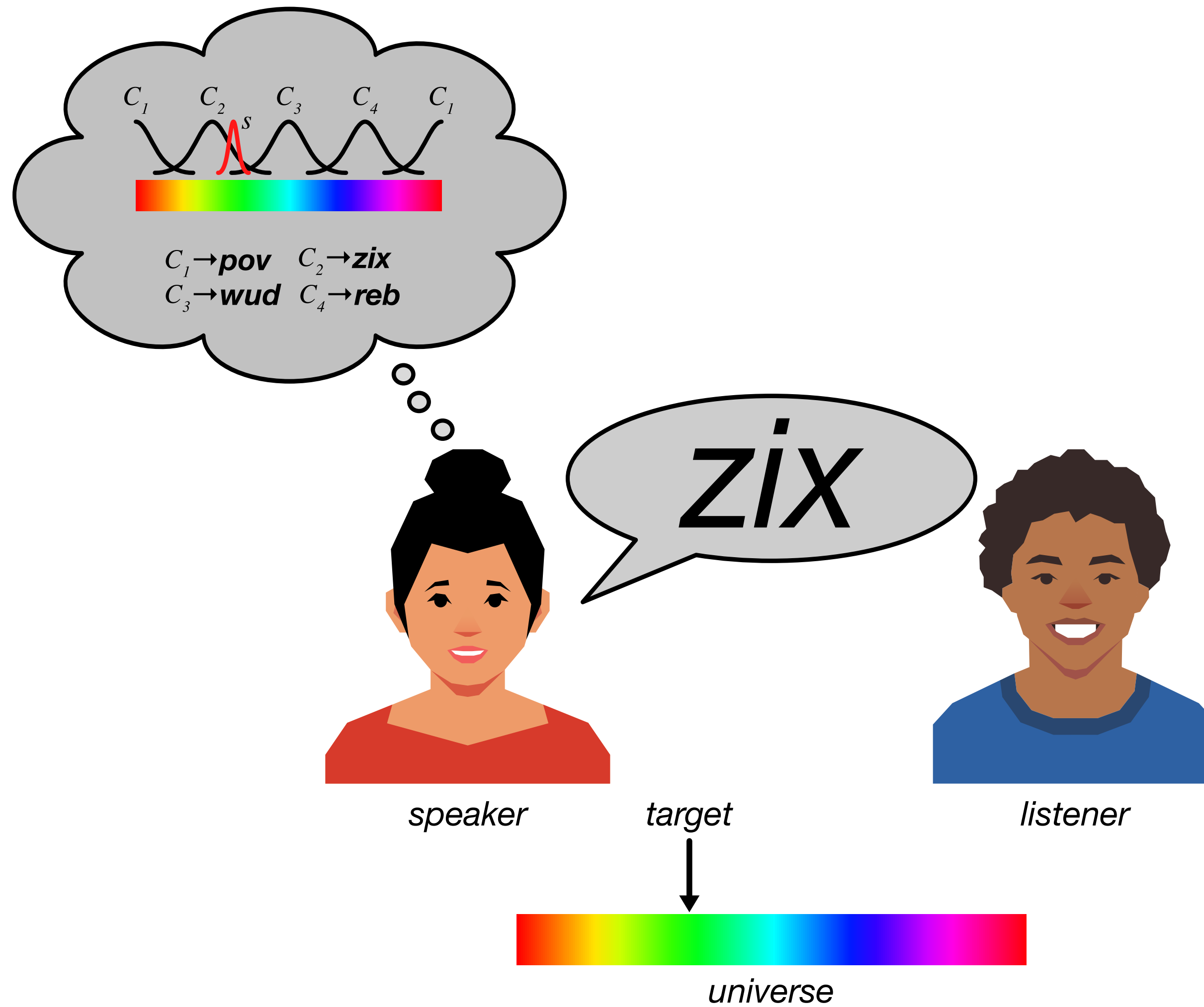
# Regier et al.'s informativeness model



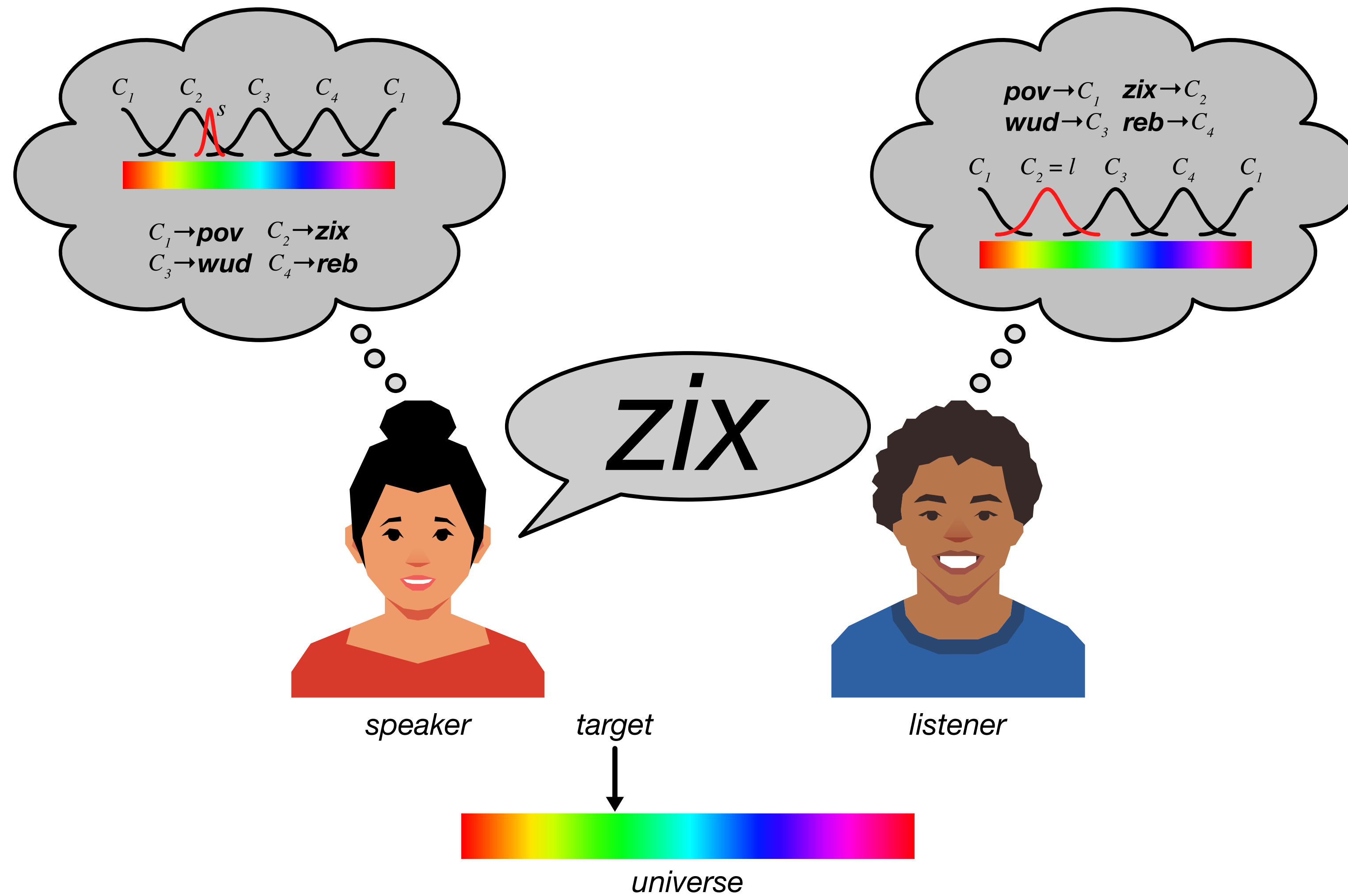
# Regier et al.'s informativeness model



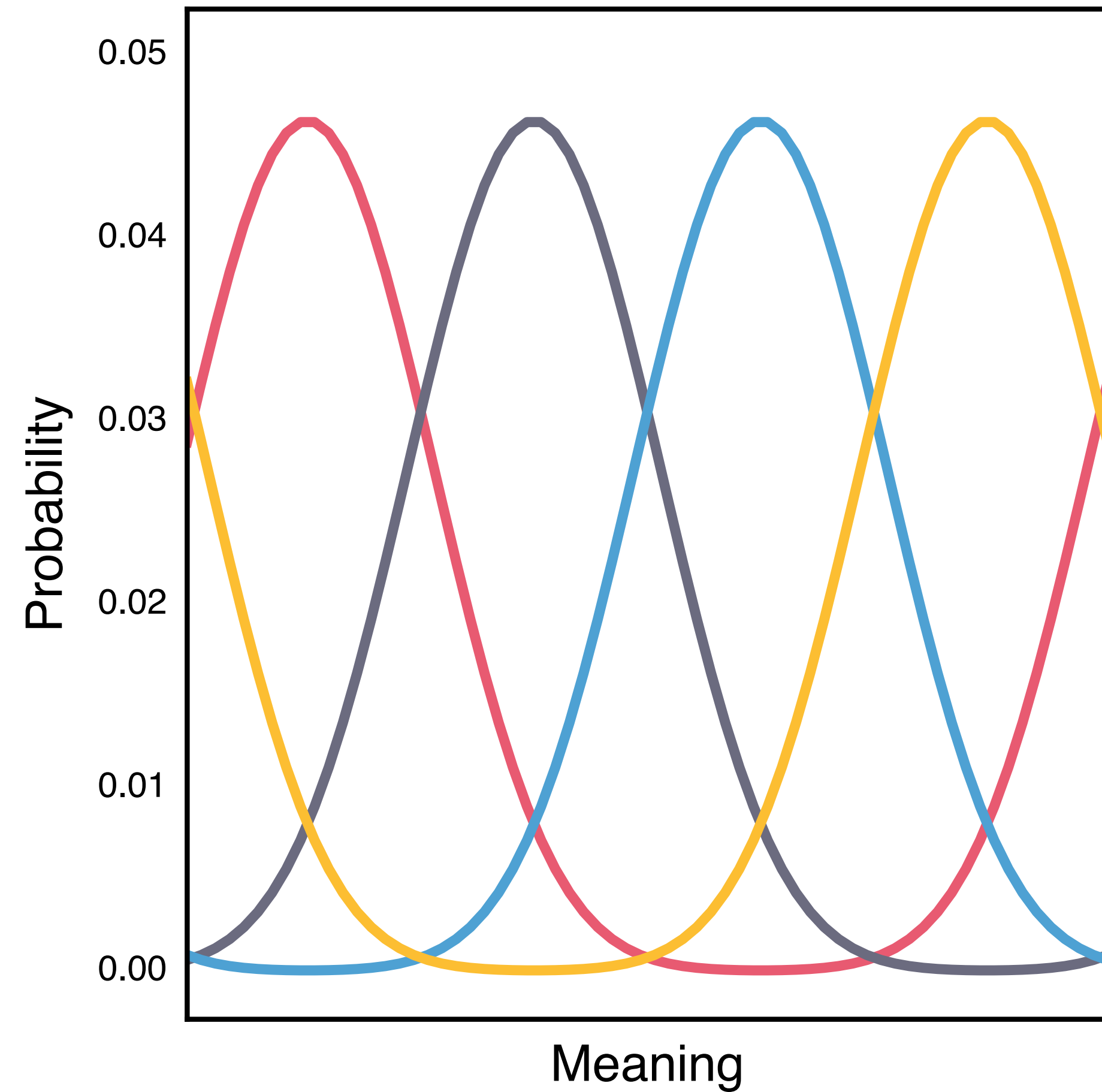
# Regier et al.'s informativeness model



# Regier et al.'s informativeness model



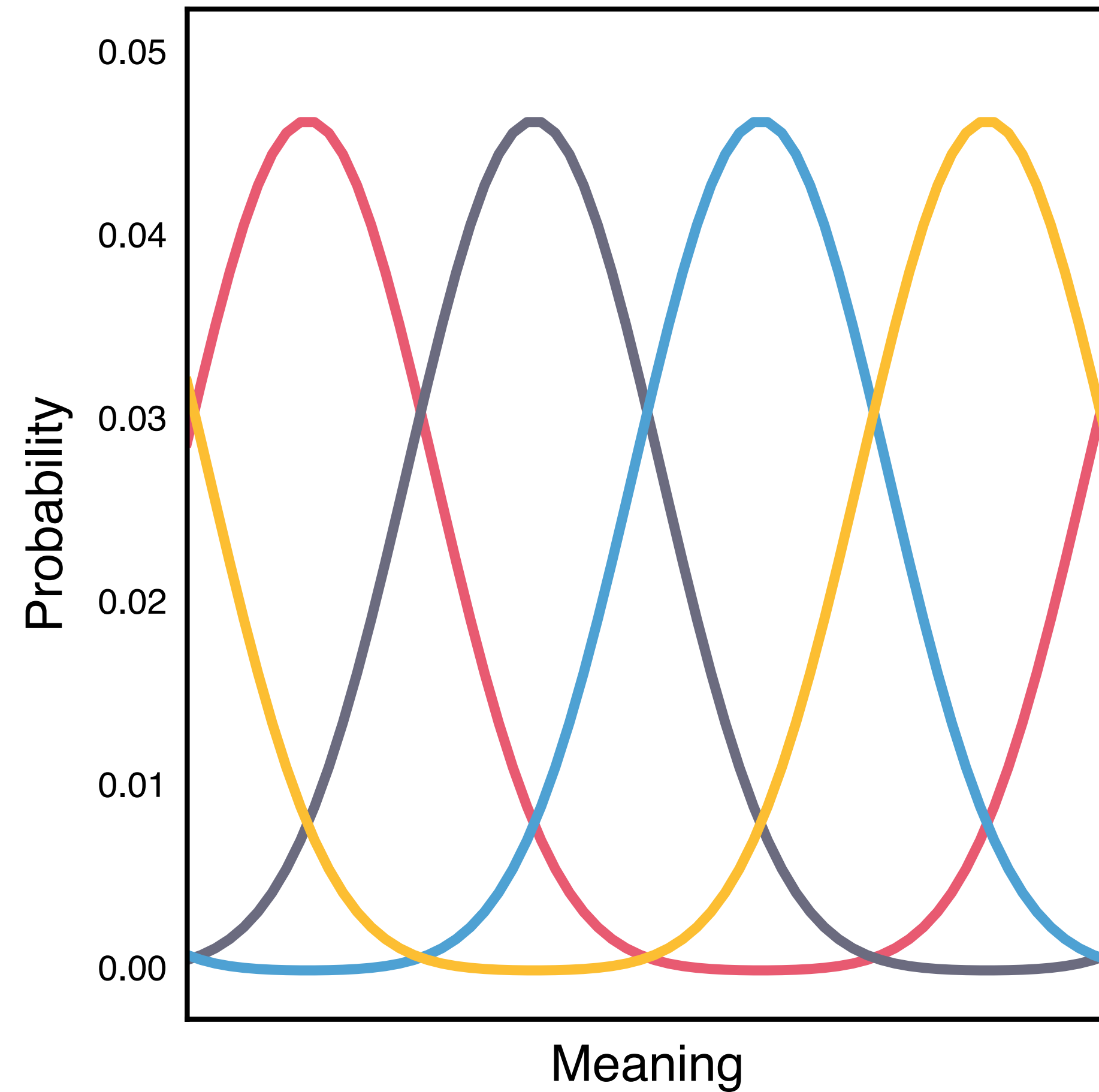
# Communicative cost



$$C_j(i) \propto \sum_{c \in C_j} e^{-\gamma d(i,c)^2}$$

$$K(L) := \sum_{i \in U} P(i) \cdot -\log C(i)$$

# Communicative cost

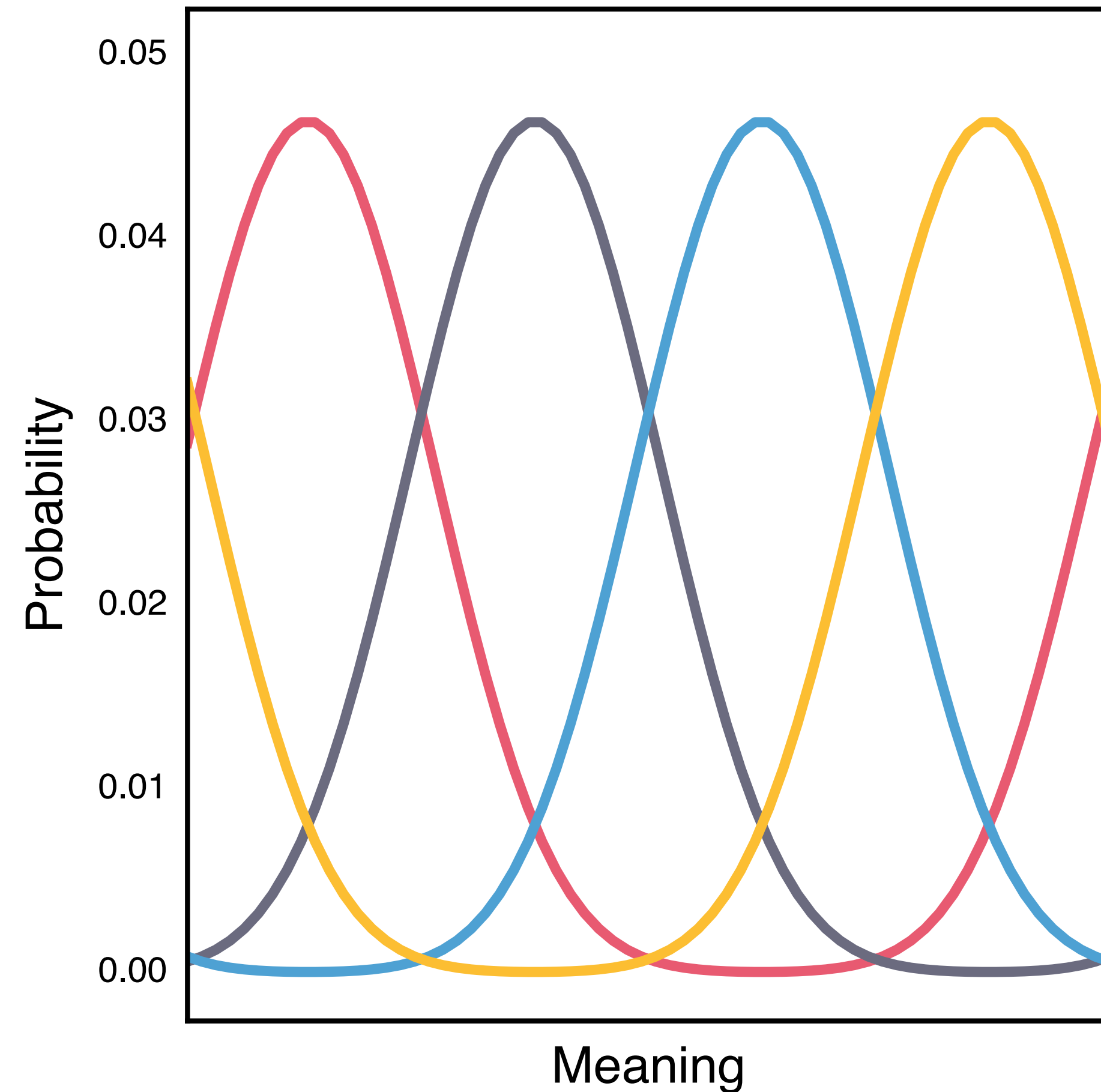


$$C_j(i) \propto \sum_{c \in C_j} e^{-\gamma d(i,c)^2}$$

$$K(L) := \sum_{i \in U} P(i) \cdot -\log C(i)$$

**Expressivity** A system of many categories is more informative than a system of few categories

# Communicative cost



$$C_j(i) \propto \sum_{c \in C_j} e^{-\gamma d(i,c)^2}$$

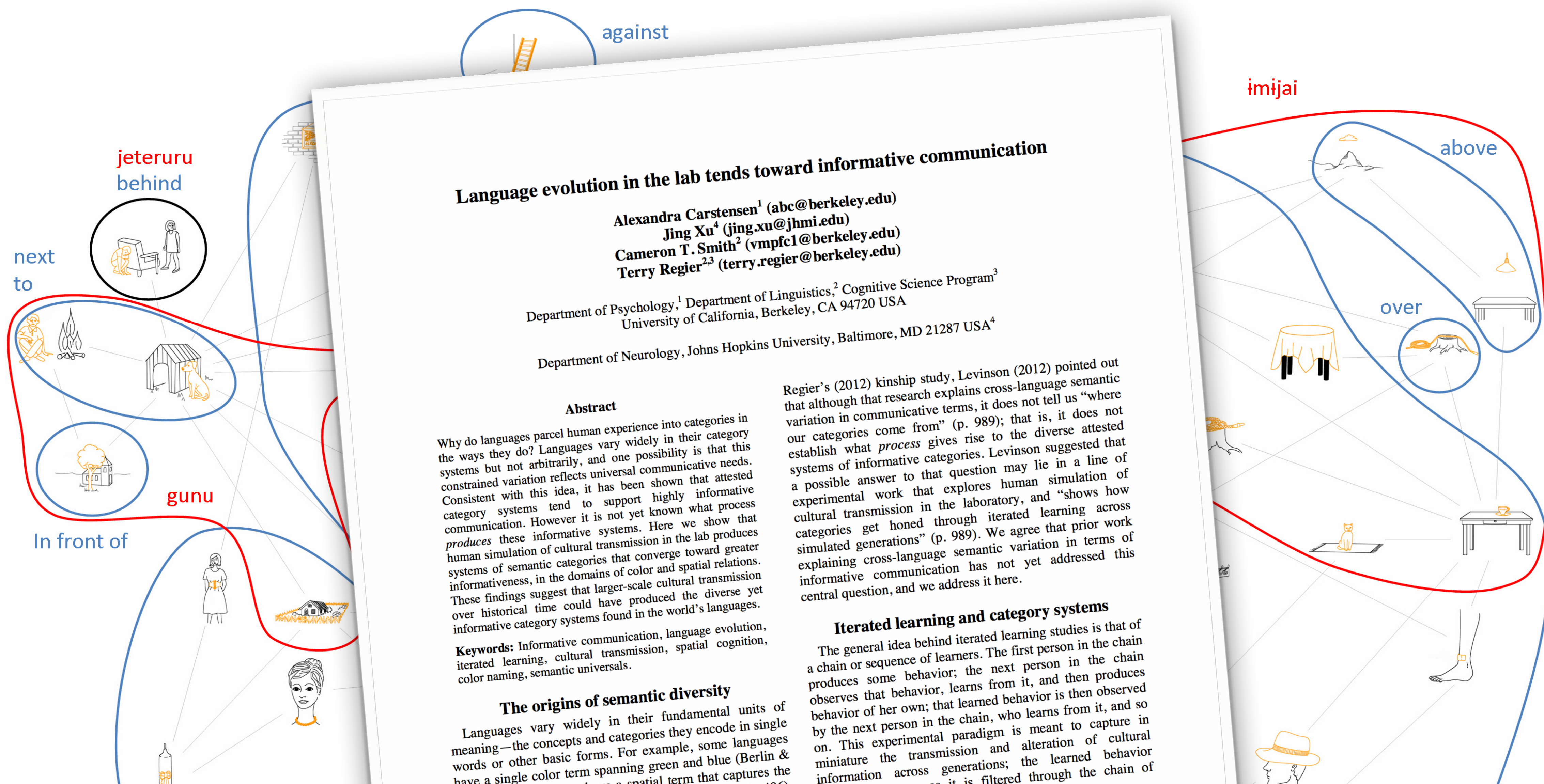
$$K(L) := \sum_{i \in U} P(i) \cdot -\log C(i)$$

**Expressivity** A system of many categories is more informative than a system of few categories

**Compactness** A system of compact categories is more informative than a system of noncompact categories

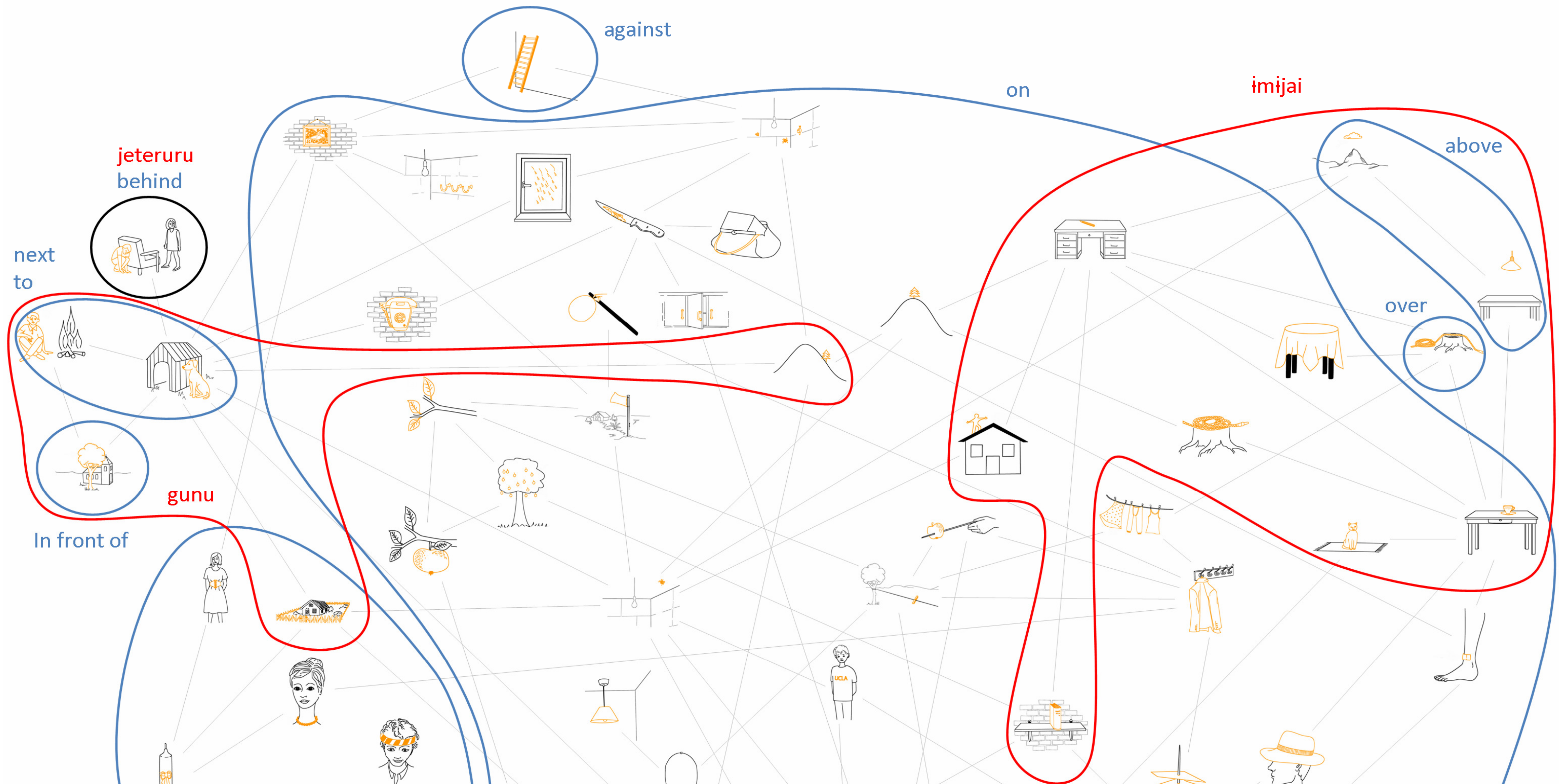


# Could humans have a learning bias for informativeness?



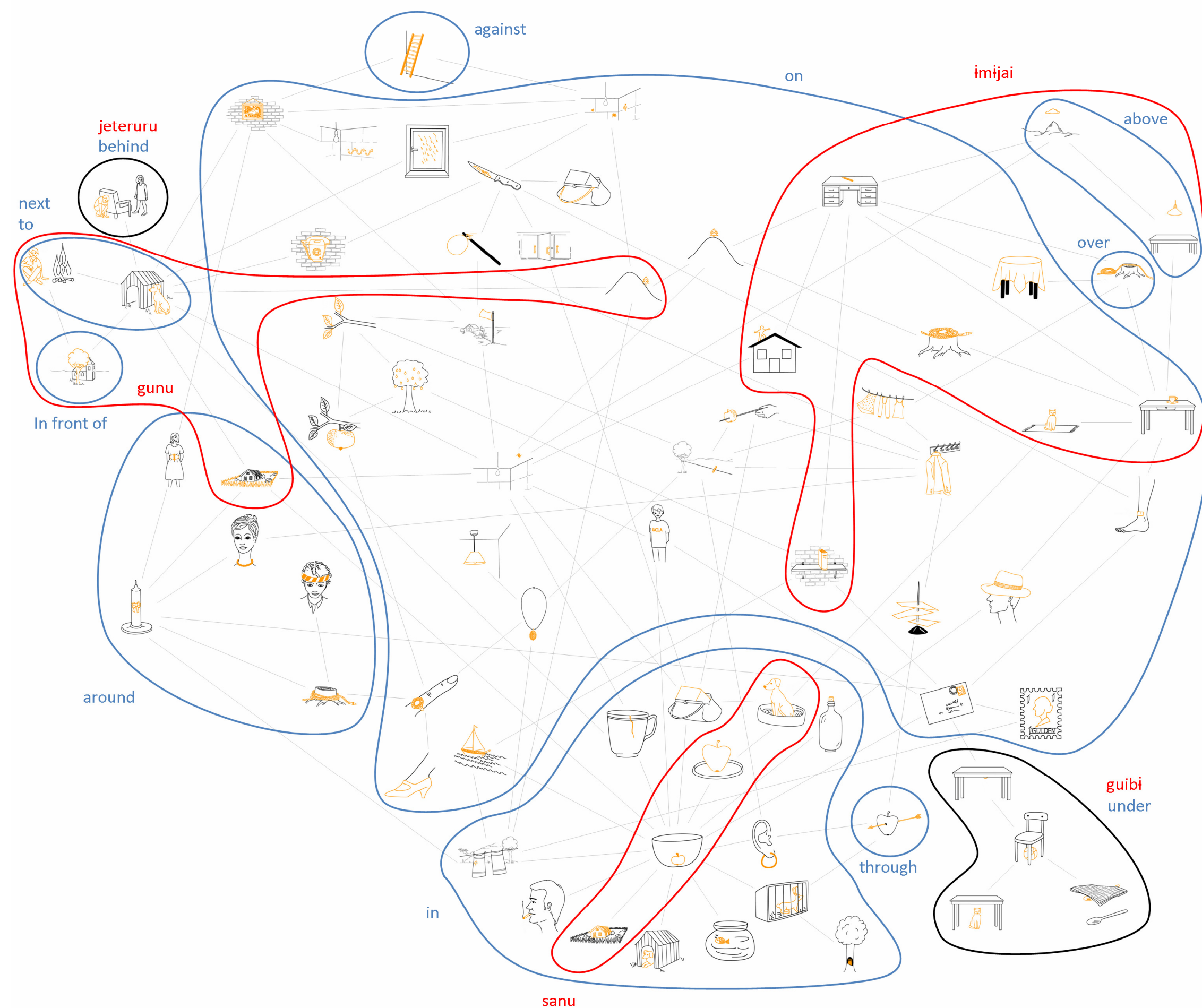


# Could humans have a learning bias for informativeness?

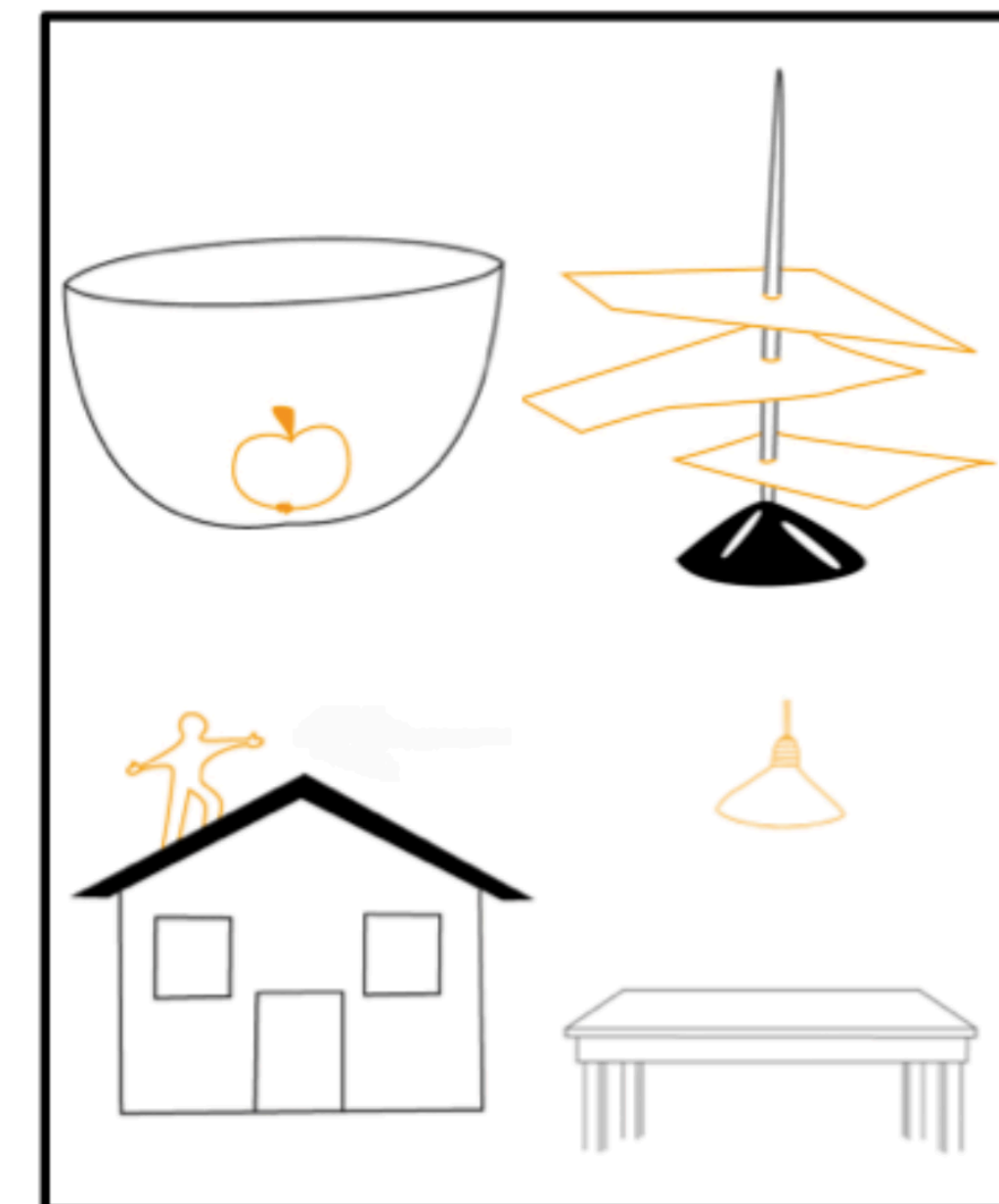




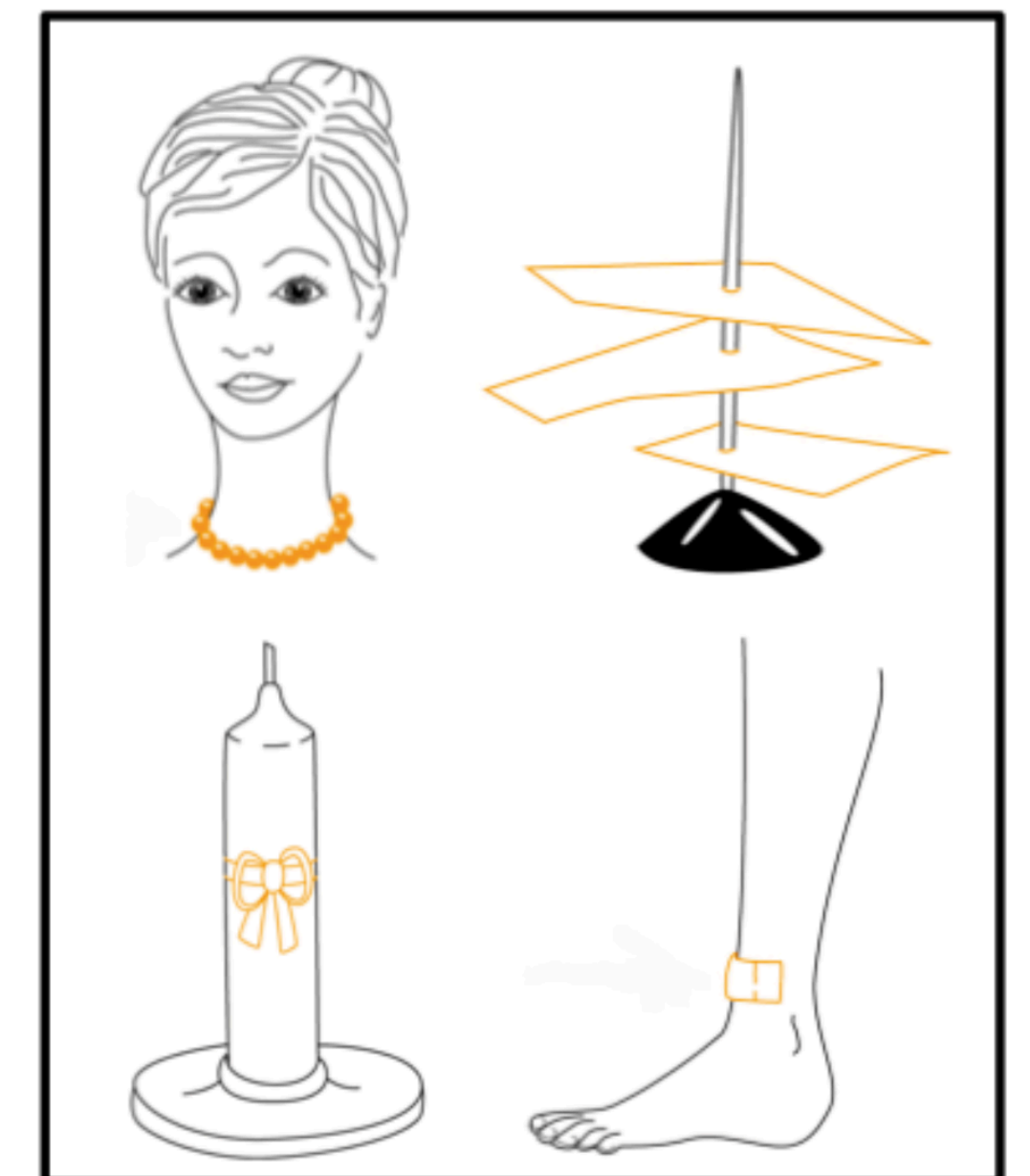
# Could humans have a learning bias for informativeness?



Generation 0

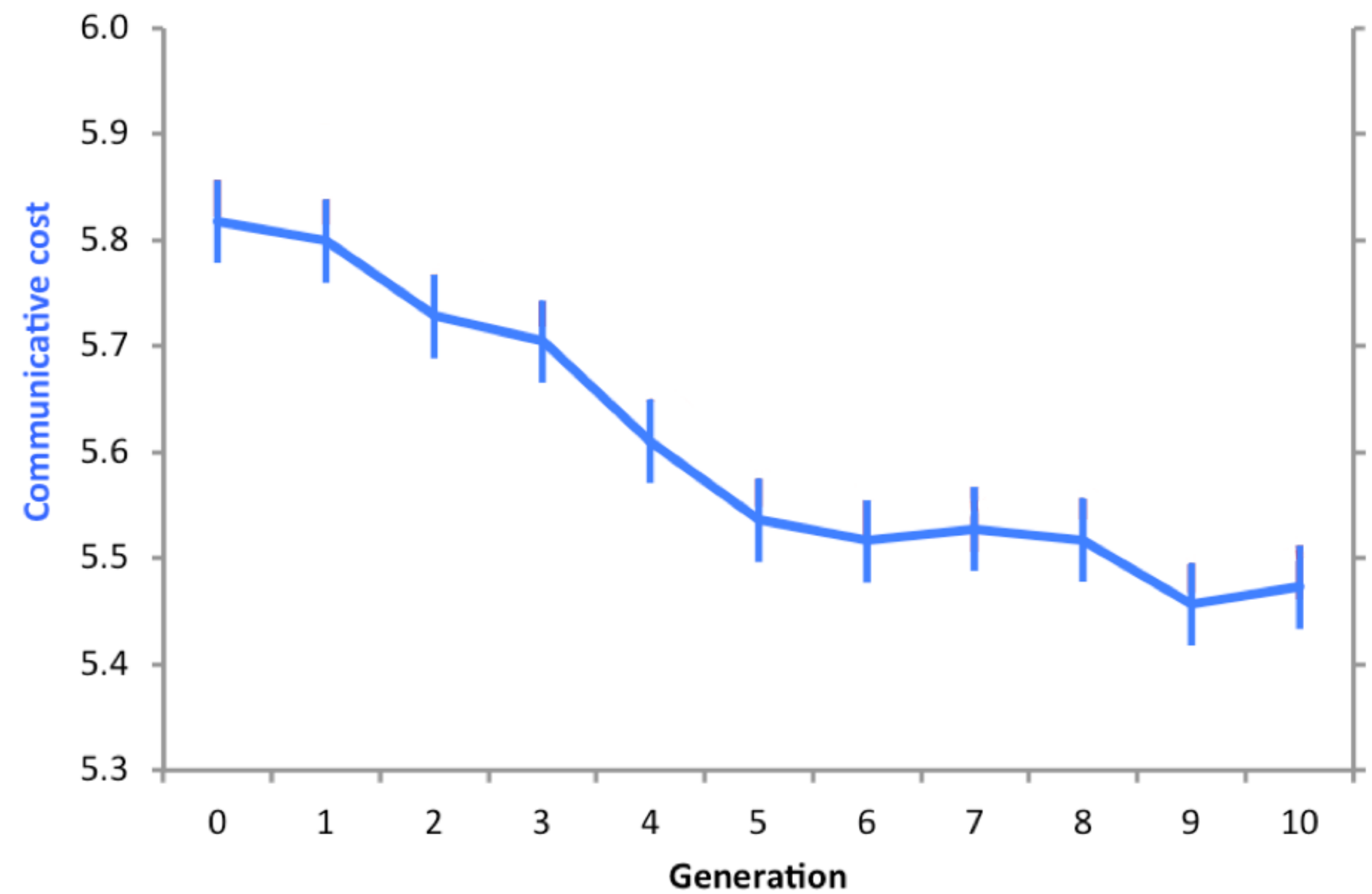
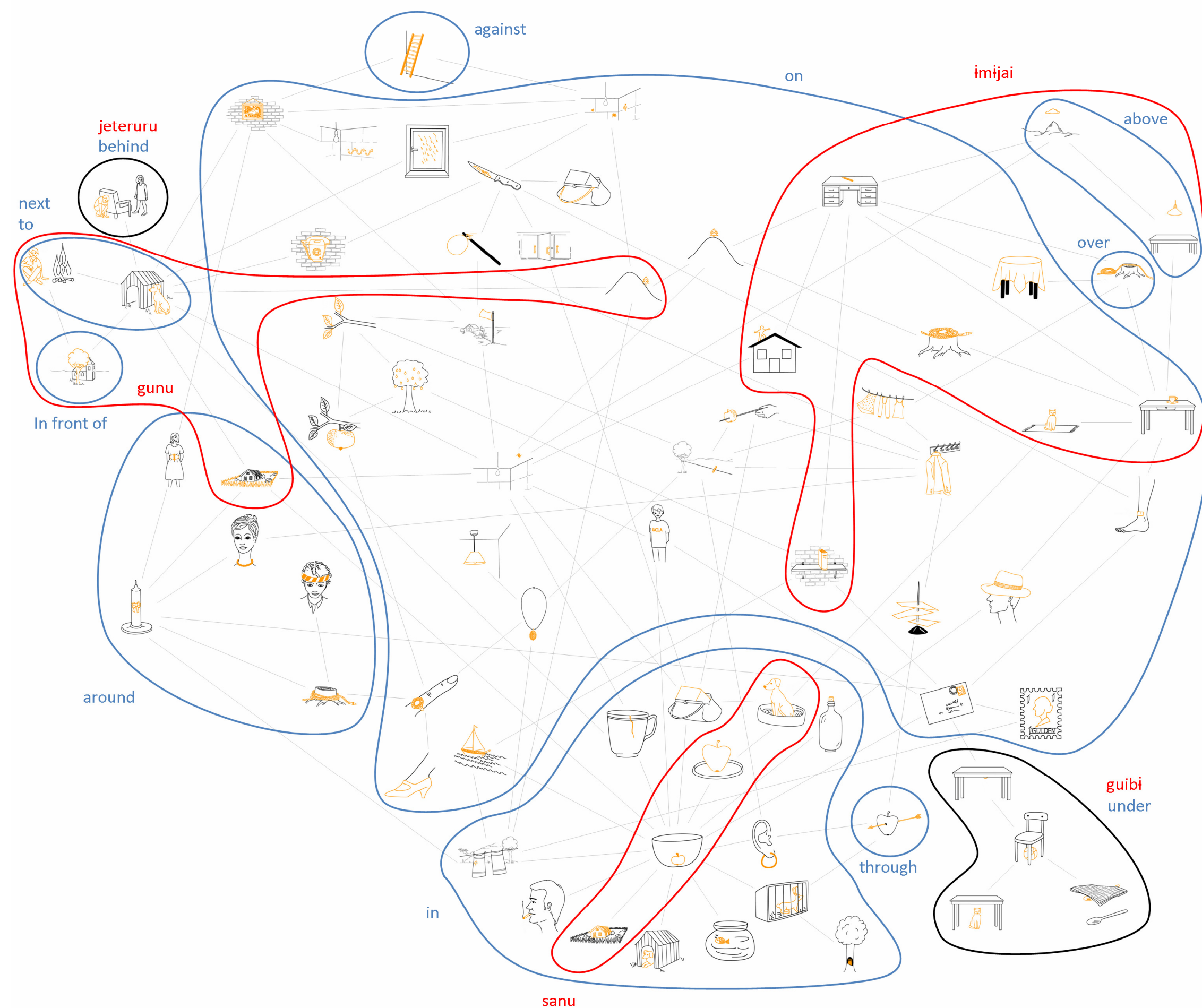


Generation 10



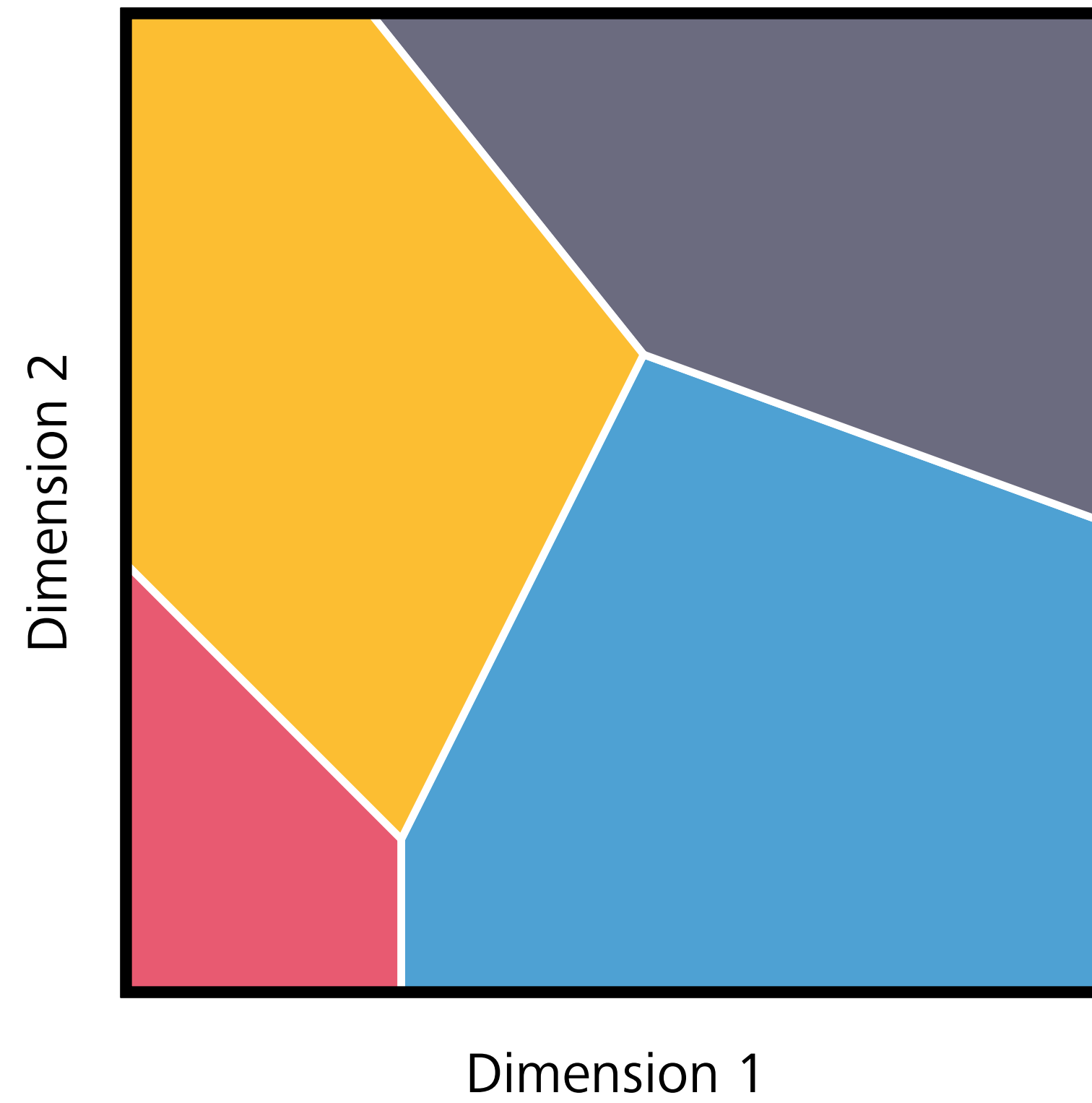


# Could humans have a learning bias for informativeness?

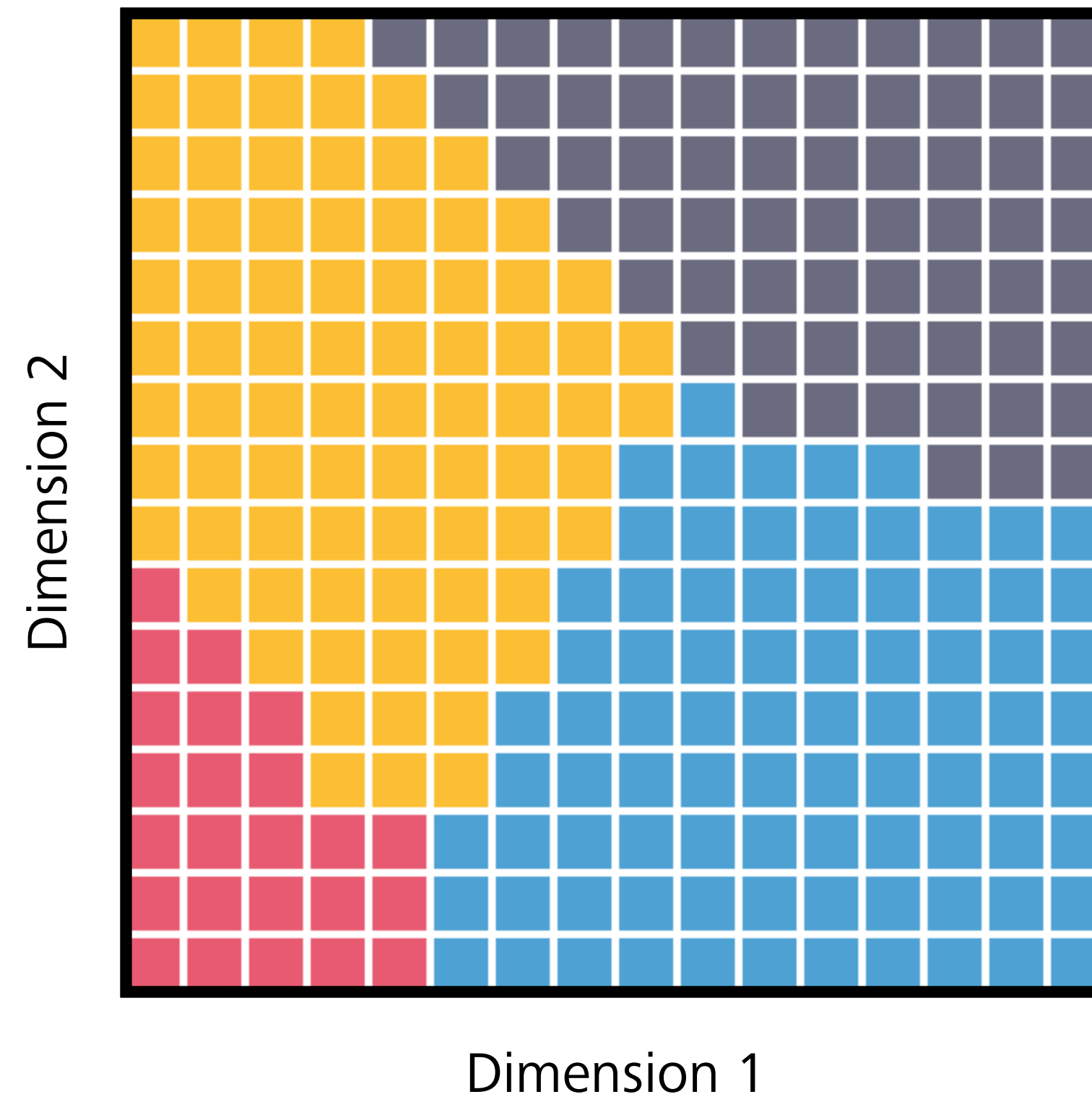


*Bayesian model*

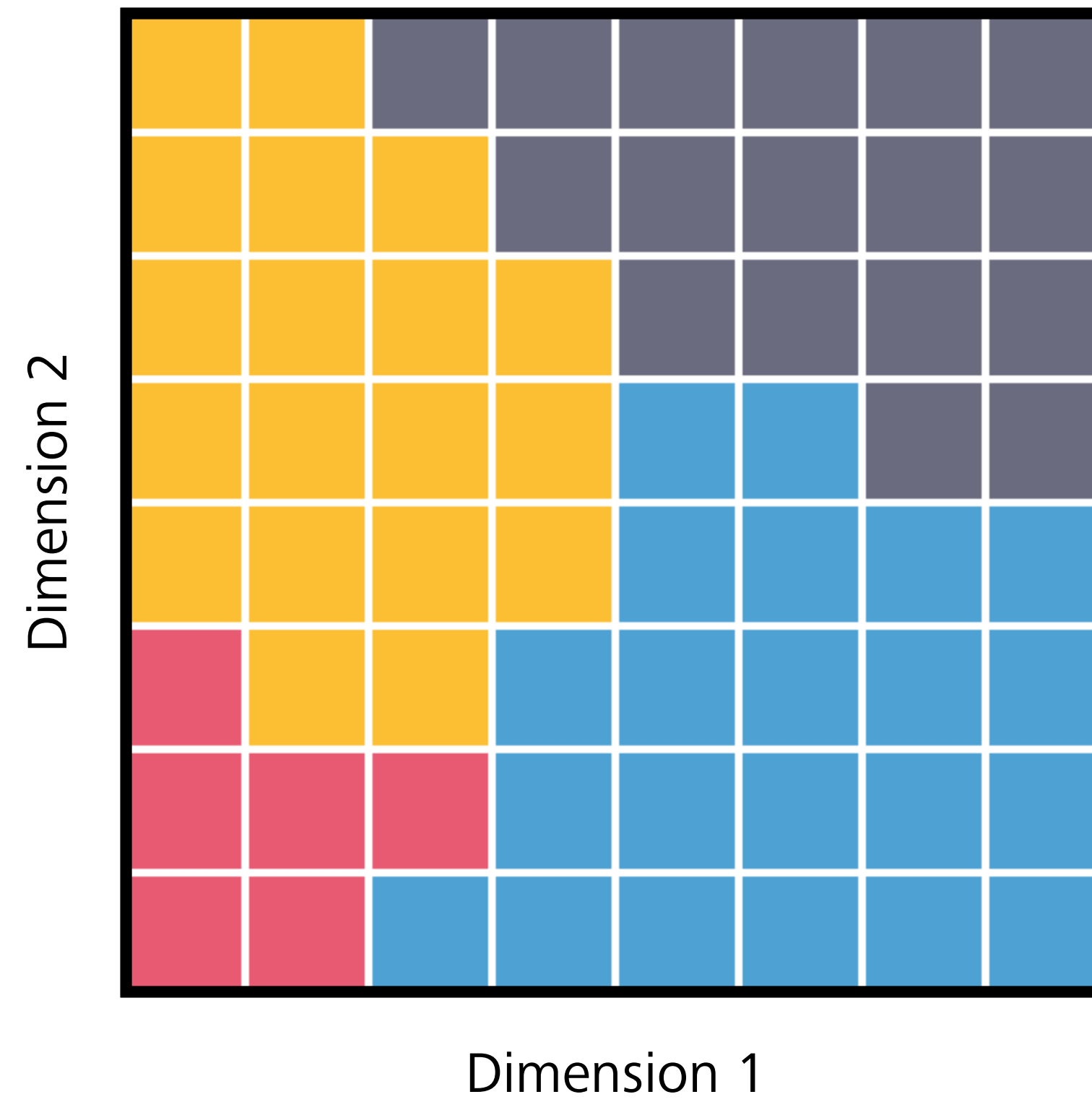
# Conceptual spaces and convexity



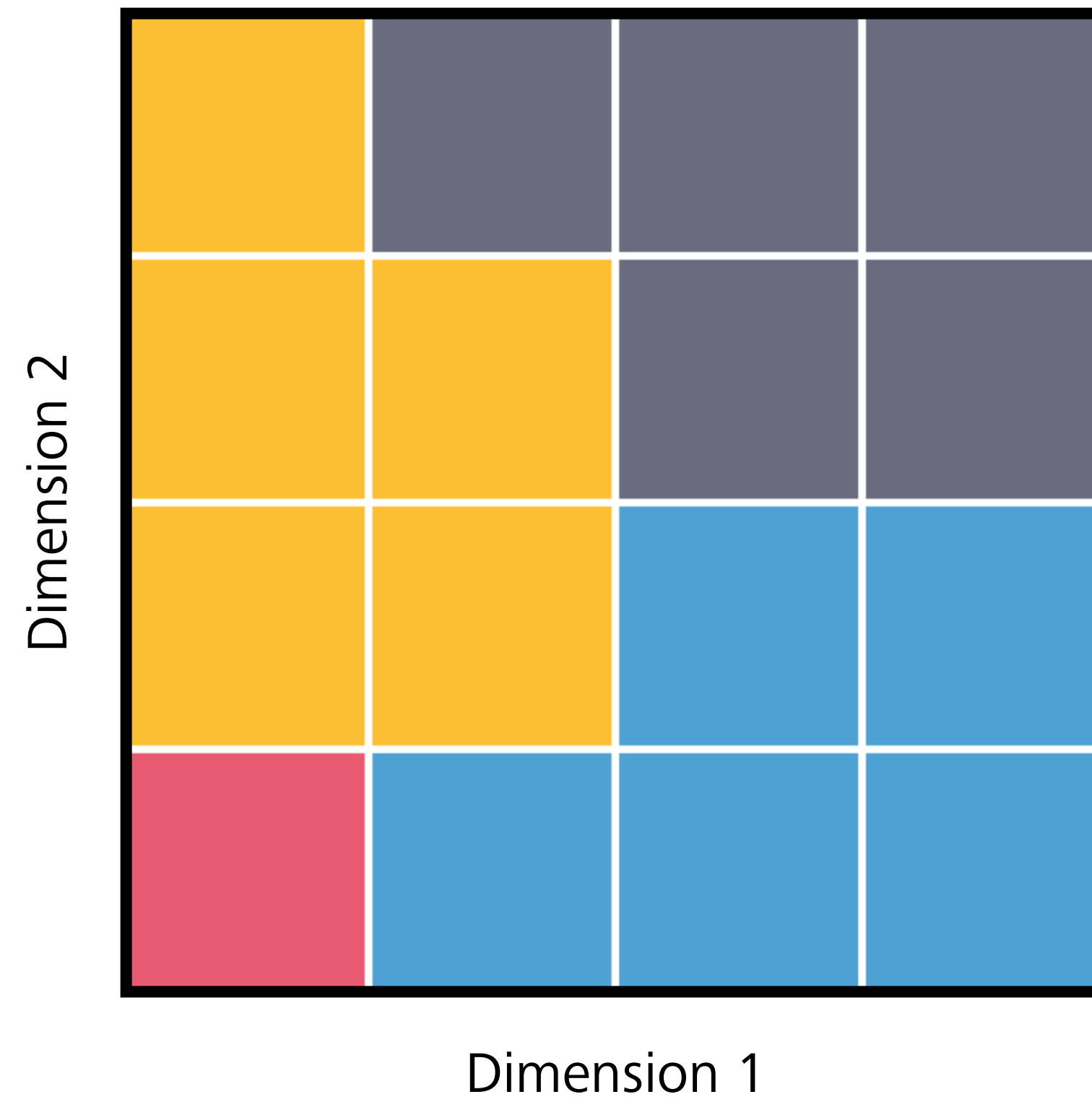
# Conceptual spaces and convexity



# Conceptual spaces and convexity



# Conceptual spaces and convexity





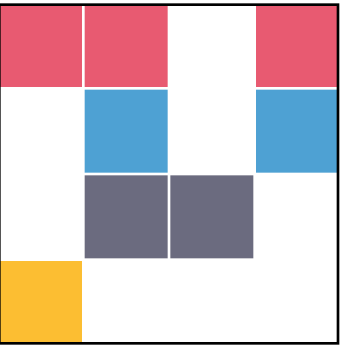
# Bayesian inference

$$\mathcal{L} = \left\{ \begin{array}{c} \begin{array}{|c|c|c|c|} \hline \text{yellow} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{blue} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{pink} & \text{grey} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{grey} & \text{yellow} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{pink} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{grey} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{pink} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{pink} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{grey} & \text{pink} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{blue} \\ \hline \text{yellow} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{yellow} & \text{yellow} \\ \hline \text{grey} & \text{blue} & \text{yellow} & \text{pink} \\ \hline \text{pink} & \text{yellow} & \text{yellow} & \text{grey} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{pink} & \text{grey} \\ \hline \text{pink} & \text{grey} & \text{blue} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{grey} \\ \hline \end{array} \end{array} \dots \right\}$$

# Bayesian inference

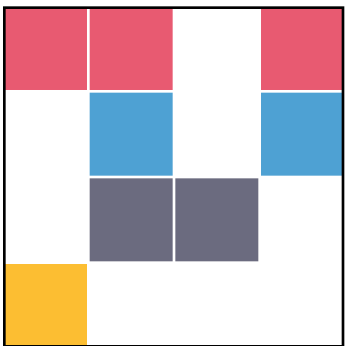
$$\mathcal{L} = \left\{ \begin{array}{|c|c|c|c|} \hline \text{yellow} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{blue} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{pink} & \text{grey} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{grey} & \text{yellow} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{yellow} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{pink} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{grey} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{pink} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{pink} & \text{blue} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{grey} & \text{pink} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{blue} \\ \hline \text{yellow} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{yellow} & \text{yellow} \\ \hline \text{grey} & \text{blue} & \text{yellow} & \text{pink} \\ \hline \text{pink} & \text{yellow} & \text{yellow} & \text{grey} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{grey} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{pink} & \text{grey} \\ \hline \text{pink} & \text{grey} & \text{blue} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{grey} \\ \hline \end{array} \dots \right\}$$

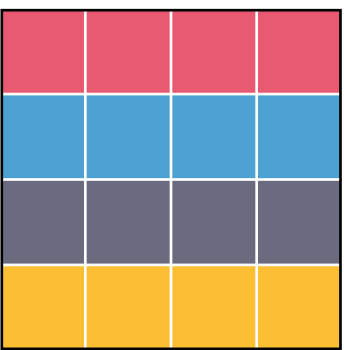
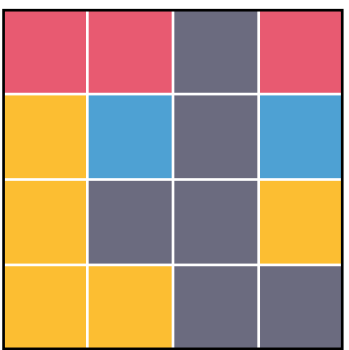
$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$



# Bayesian inference

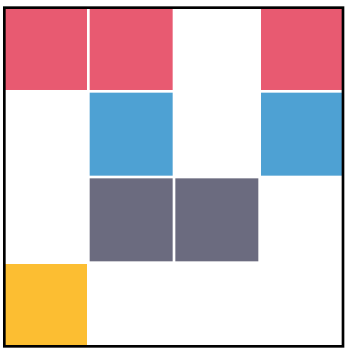
$$\mathcal{L} = \left\{ \begin{array}{c} \begin{array}{|c|c|c|c|} \hline \text{yellow} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{blue} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{pink} & \text{grey} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{grey} & \text{yellow} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{pink} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{grey} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{pink} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{pink} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{grey} & \text{pink} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{blue} \\ \hline \text{yellow} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{yellow} & \text{yellow} \\ \hline \text{grey} & \text{blue} & \text{yellow} & \text{pink} \\ \hline \text{pink} & \text{yellow} & \text{yellow} & \text{grey} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{pink} & \text{grey} \\ \hline \text{pink} & \text{grey} & \text{blue} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{grey} \\ \hline \end{array} \dots \end{array} \right\}$$

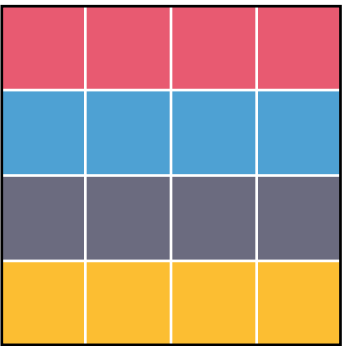
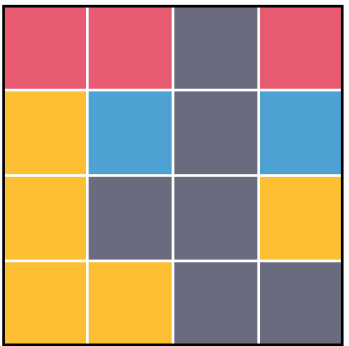
$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$


$$\text{likelihood}(D|L) \propto \prod_{\langle m,s \rangle \in D} \frac{1}{|M|} P(s|L, m)$$

$$=$$


# Bayesian inference

$$\mathcal{L} = \left\{ \begin{array}{c} \begin{array}{|c|c|c|c|} \hline \text{yellow} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{blue} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{pink} & \text{grey} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{grey} & \text{yellow} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{pink} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{grey} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{pink} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{pink} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{grey} & \text{pink} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{blue} \\ \hline \text{yellow} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{yellow} & \text{yellow} \\ \hline \text{grey} & \text{blue} & \text{yellow} & \text{pink} \\ \hline \text{pink} & \text{yellow} & \text{yellow} & \text{grey} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{pink} & \text{grey} \\ \hline \text{pink} & \text{grey} & \text{blue} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{grey} \\ \hline \end{array} \dots \end{array} \right\}$$

$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$


$$\text{likelihood}(D|L) \propto \prod_{\langle m,s \rangle \in D} \frac{1}{|M|} P(s|L, m)$$

$$=$$


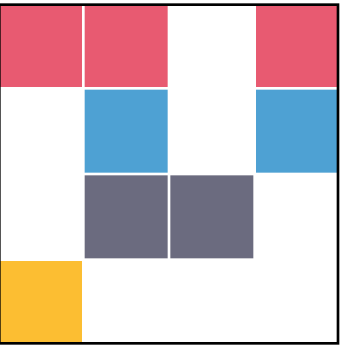
$$\text{prior}(L) \propto 2^{-\text{DL}(L)}$$

$$>$$

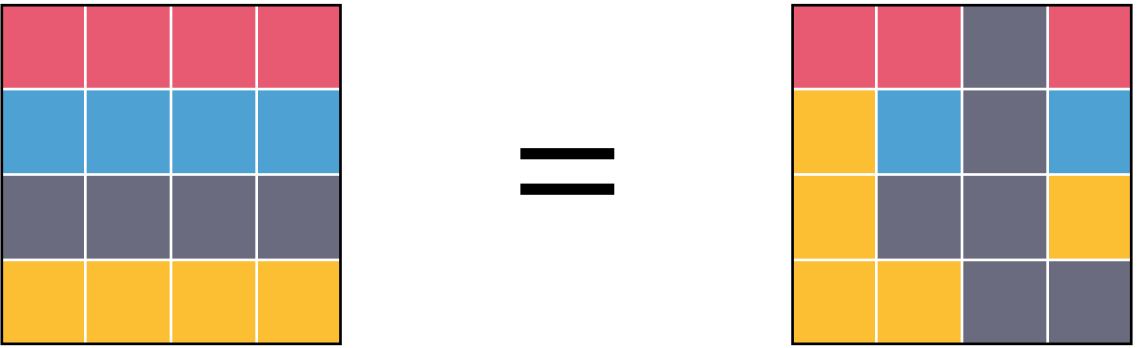

# Bayesian inference

$$\mathcal{L} = \left\{ \begin{array}{c} \begin{array}{|c|c|c|c|} \hline \text{yellow} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{blue} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{pink} & \text{grey} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{grey} & \text{yellow} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{pink} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{grey} \\ \hline \text{grey} & \text{blue} & \text{blue} & \text{pink} \\ \hline \text{blue} & \text{blue} & \text{pink} & \text{pink} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{blue} & \text{pink} & \text{pink} & \text{blue} \\ \hline \text{pink} & \text{blue} & \text{pink} & \text{blue} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{grey} & \text{pink} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{blue} \\ \hline \text{yellow} & \text{grey} & \text{grey} & \text{yellow} \\ \hline \text{yellow} & \text{yellow} & \text{grey} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{yellow} & \text{yellow} \\ \hline \text{grey} & \text{blue} & \text{yellow} & \text{pink} \\ \hline \text{pink} & \text{yellow} & \text{yellow} & \text{grey} \\ \hline \text{blue} & \text{blue} & \text{blue} & \text{grey} \\ \hline \end{array} & \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{pink} & \text{grey} \\ \hline \text{pink} & \text{grey} & \text{blue} & \text{grey} \\ \hline \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \hline \text{yellow} & \text{blue} & \text{grey} & \text{grey} \\ \hline \end{array} \dots \end{array} \right\}$$

$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$

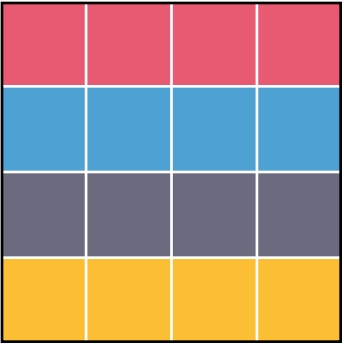


$$\text{likelihood}(D|L) \propto \prod_{\langle m, s \rangle \in D} \frac{1}{|M|} P(s|L, m)$$

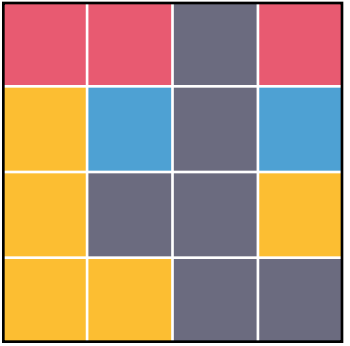


=

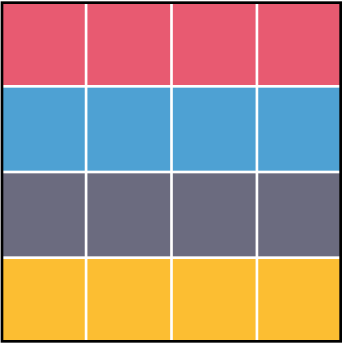
$$\text{prior}(L) \propto 2^{-\text{DL}(L)}$$



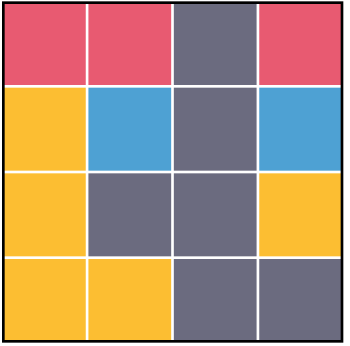
>



$$\text{posterior}(L|D) \propto \text{likelihood}(D|L) \times \text{prior}(L)$$



>



# Computing DL(L): The rectangle

Class	Position
1×1	16
1×2	12
1×3	8
1×4	4
2×1	12
2×2	9
2×3	6
2×4	3
3×1	8
3×2	6
3×3	4
3×4	2
4×1	4
4×2	3
4×3	2
4×4	1

## Categorization Under Complexity: A Unified MDL Account of Human Learning of Regular and Irregular Categories

**David Fass**  
Department of Psychology  
Center for Cognitive Science  
Rutgers University  
Piscataway, NJ 08854  
dfass@ruccs.rutgers.edu

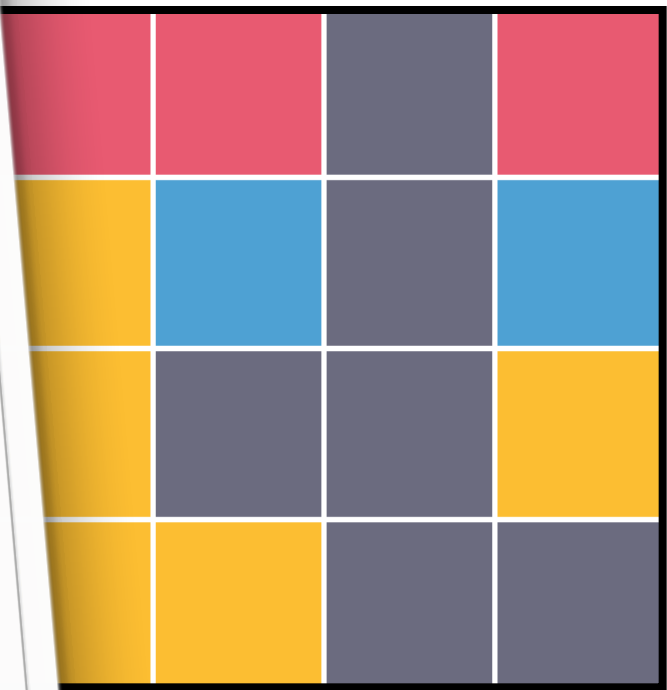
**Jacob Feldman\***  
Department of Psychology  
Center for Cognitive Science  
Rutgers University  
Piscataway, NJ 08854  
jacob@ruccs.rutgers.edu

### Abstract

We present an account of human concept learning—that is, learning of categories from examples—based on the principle of minimum description length (MDL). In support of this theory, we tested a wide range of two-dimensional concept types, including both regular (simple) and highly irregular (complex) structures, and found the MDL theory to give a good account of subjects' performance. This suggests that the *intrinsic complexity* of a concept (that is, its description length) systematically influences its learnability.

### 1 The Structure of Categories

A number of different principles have been advanced to explain the manner in which humans learn to categorize objects. It has been variously suggested that the underlying principle might be the *similarity structure* of objects [1], the manipulability of *decision boundaries* [2], or the *inference* [3][4]. While many of these theories are mathematically similar to



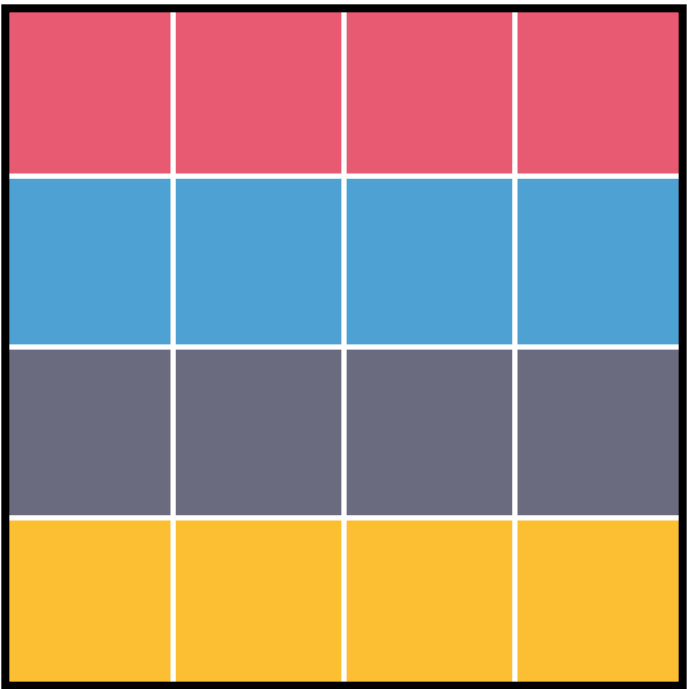
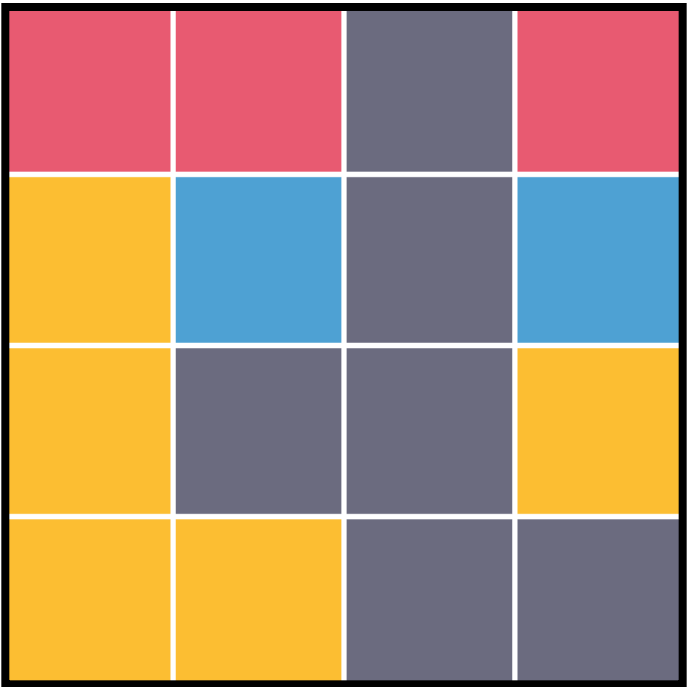
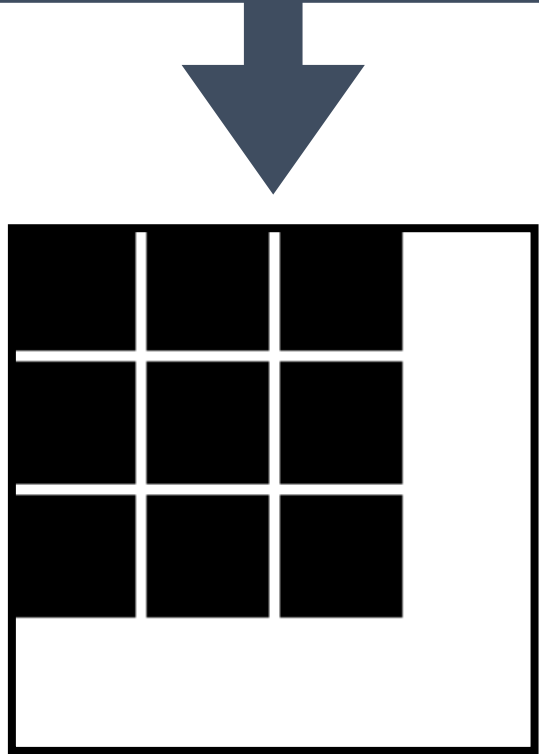
# Computing DL(L): The rectangle code

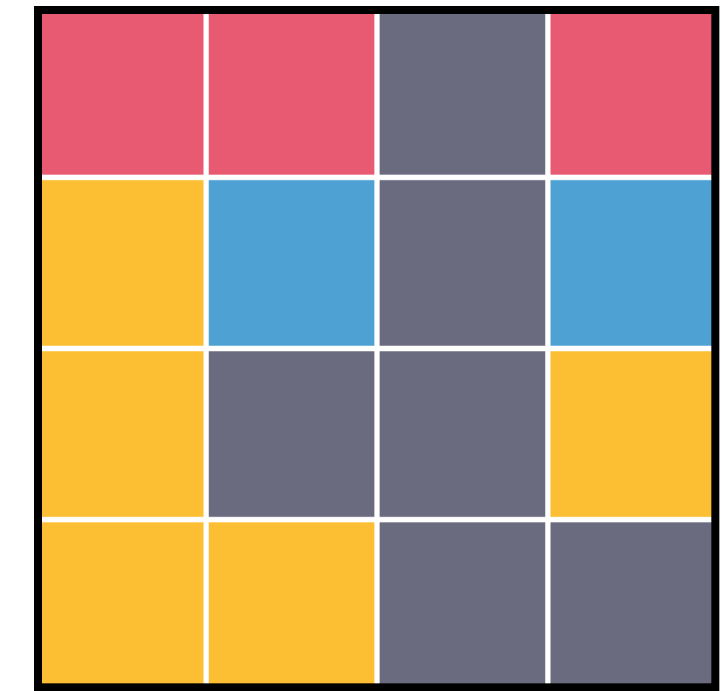
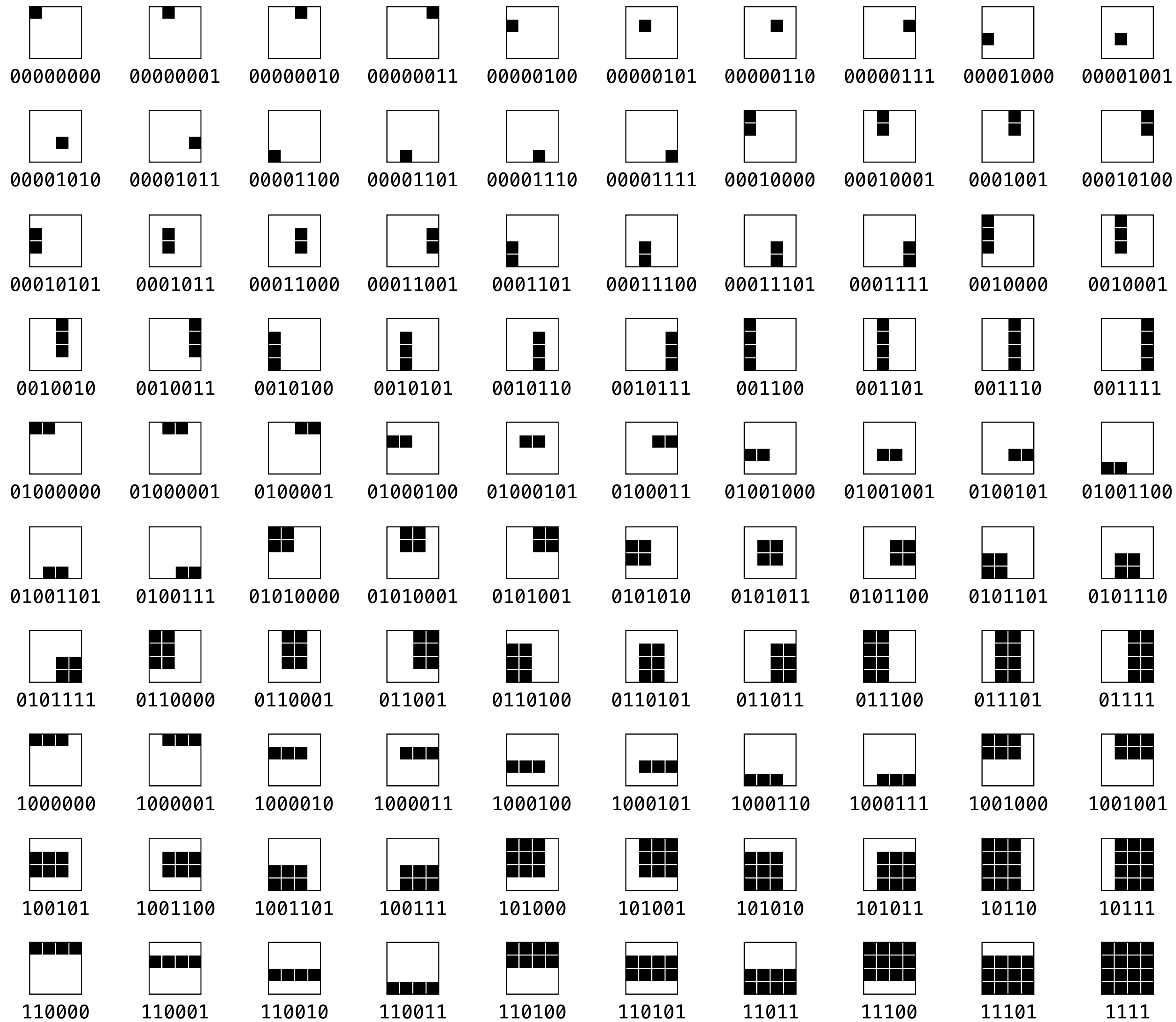
Class	Positions	Probability	Codelength	(bits)
1×1	16	$1/16 \times 1/16$	$-\log 1/256$	8.0
1×2	12	$1/16 \times 1/12$	$-\log 1/192$	7.58
1×3	8	$1/16 \times 1/8$	$-\log 1/128$	7.0
1×4	4	$1/16 \times 1/4$	$-\log 1/64$	6.0
2×1	12	$1/16 \times 1/12$	$-\log 1/192$	7.58
2×2	9	$1/16 \times 1/9$	$-\log 1/144$	7.17
2×3	6	$1/16 \times 1/6$	$-\log 1/96$	6.58
2×4	3	$1/16 \times 1/3$	$-\log 1/48$	5.58
3×1	8	$1/16 \times 1/8$	$-\log 1/128$	7.0
3×2	6	$1/16 \times 1/6$	$-\log 1/96$	6.58
3×3	4	$1/16 \times 1/4$	$-\log 1/64$	6.0
3×4	2	$1/16 \times 1/2$	$-\log 1/32$	5.0
4×1	4	$1/16 \times 1/4$	$-\log 1/64$	6.0
4×2	3	$1/16 \times 1/3$	$-\log 1/48$	5.58
4×3	2	$1/16 \times 1/2$	$-\log 1/32$	5.0
4×4	1	$1/16 \times 1/1$	$-\log 1/16$	4.0

Uniformly sample a class

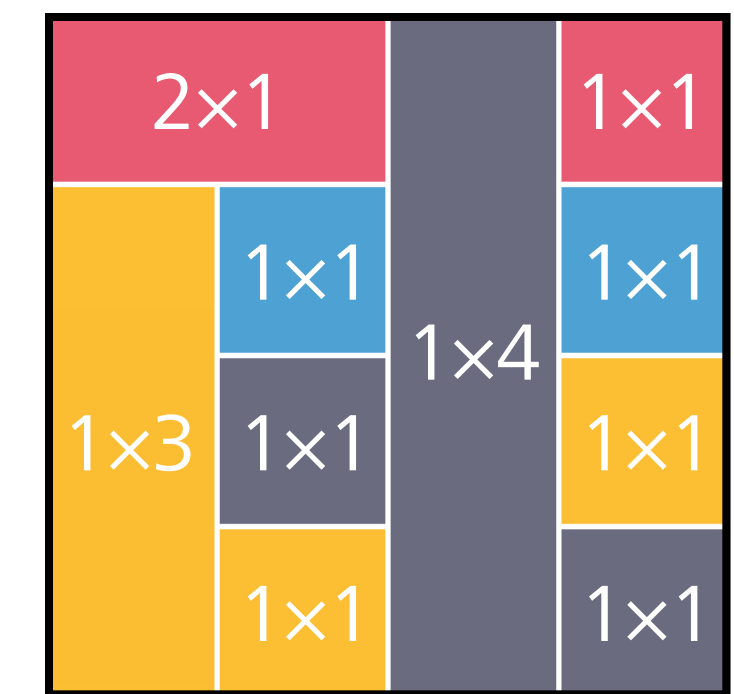
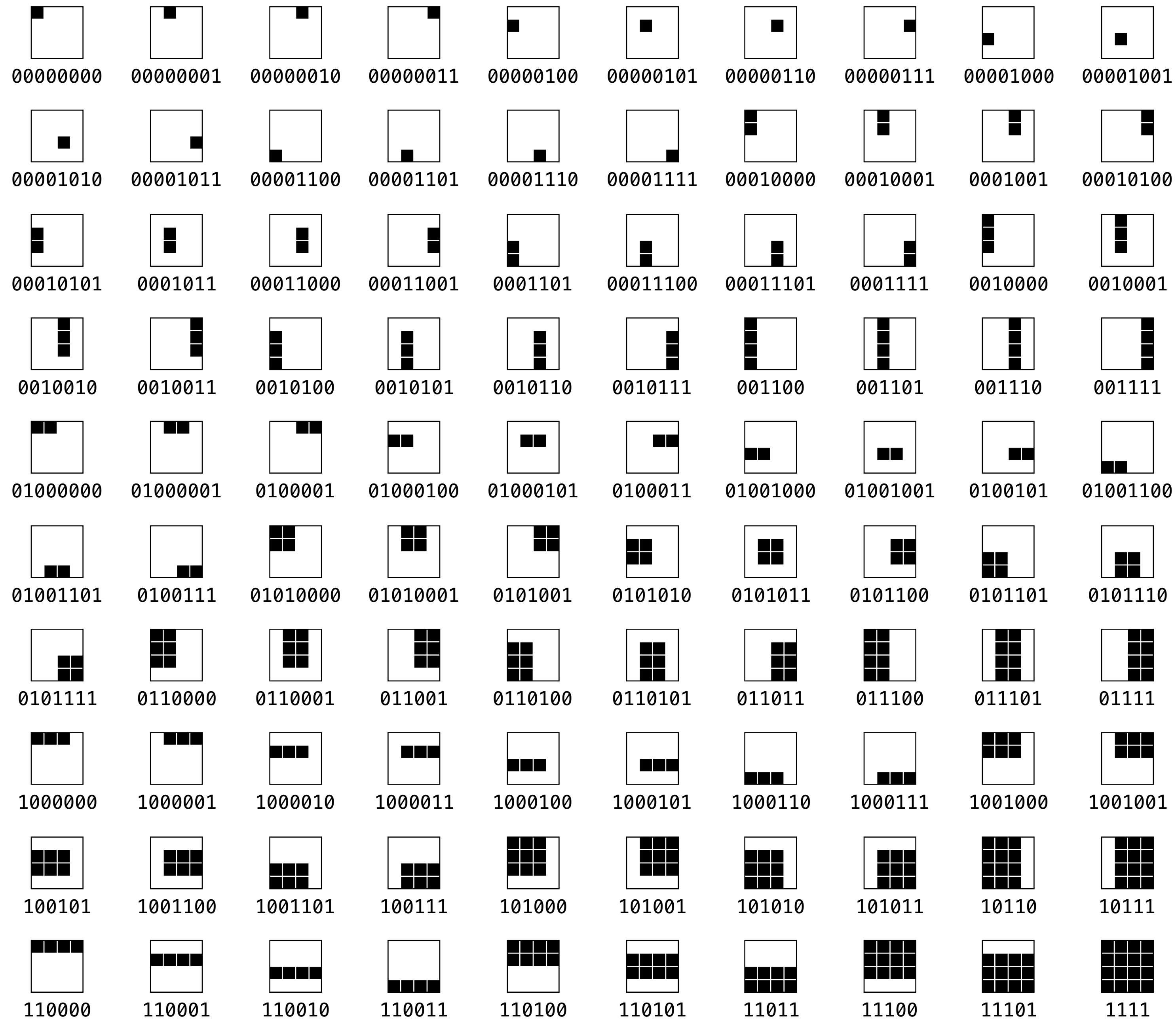


Uniformly sample a position





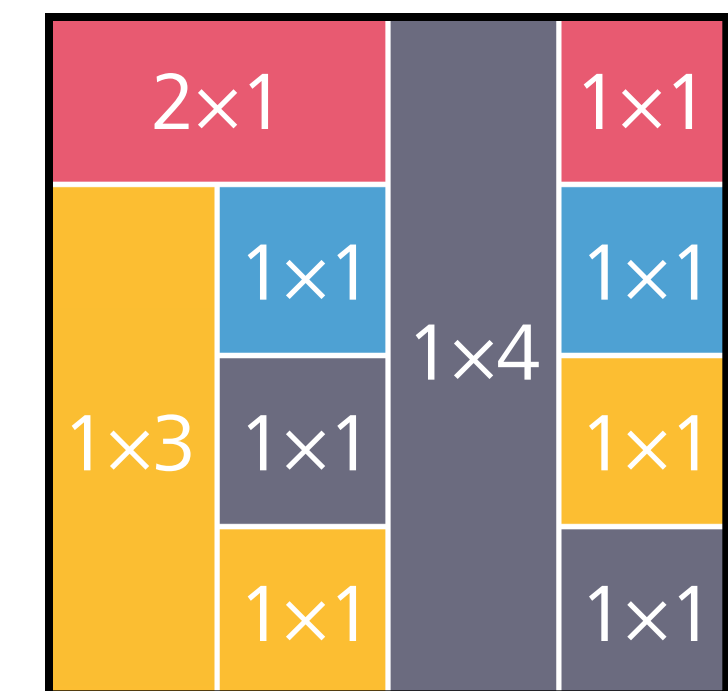
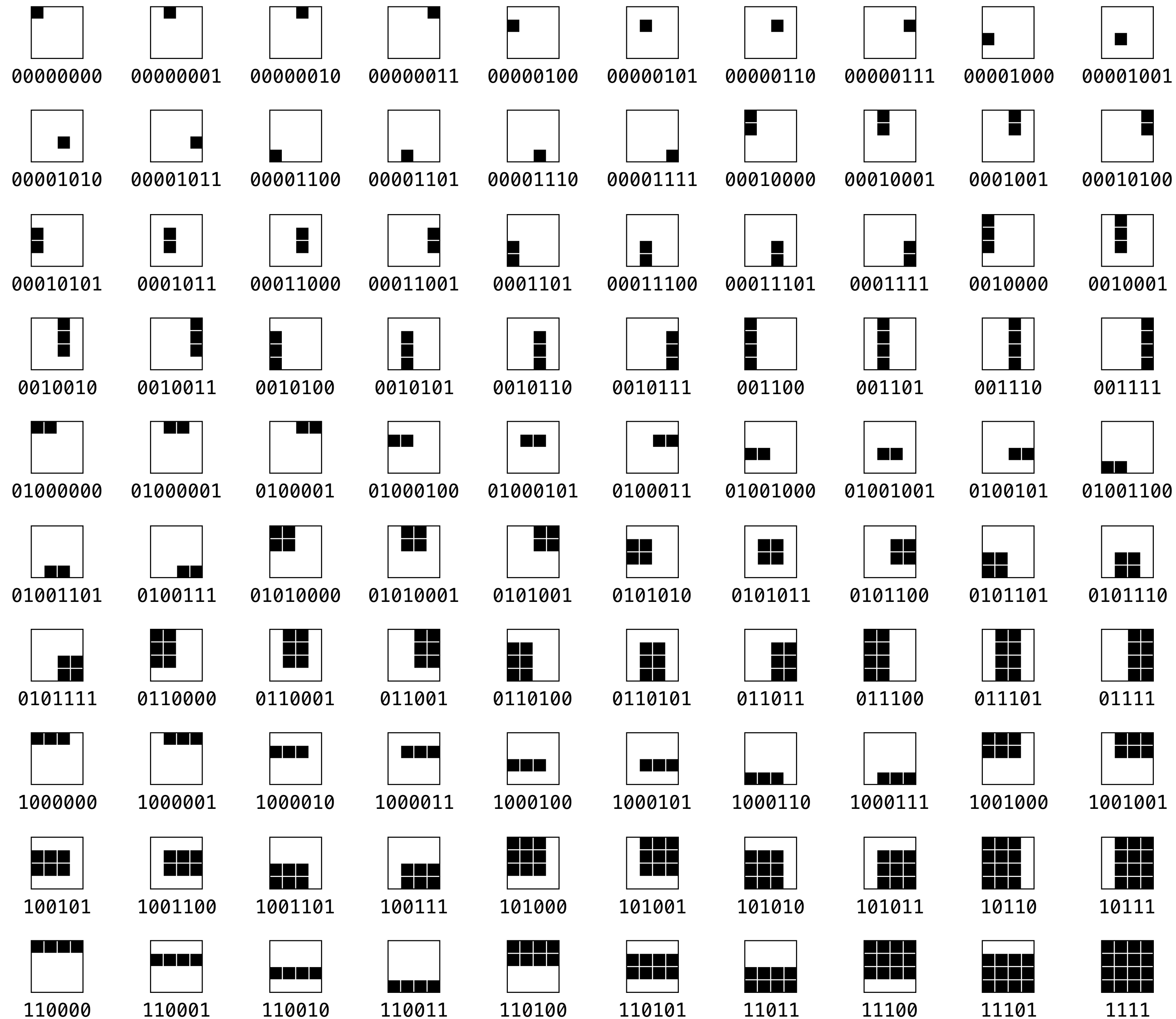




0100000000000011  
0000010100000111  
0011100000100100001111  
00101000000101100001101

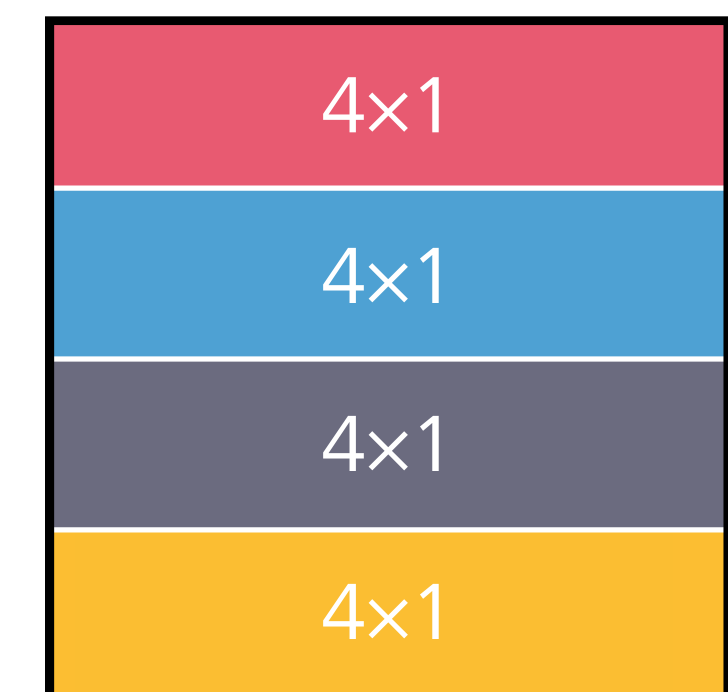
76.58 bits





0100000000000011  
0000010100000111  
0011100000100100001111  
00101000000101100001101

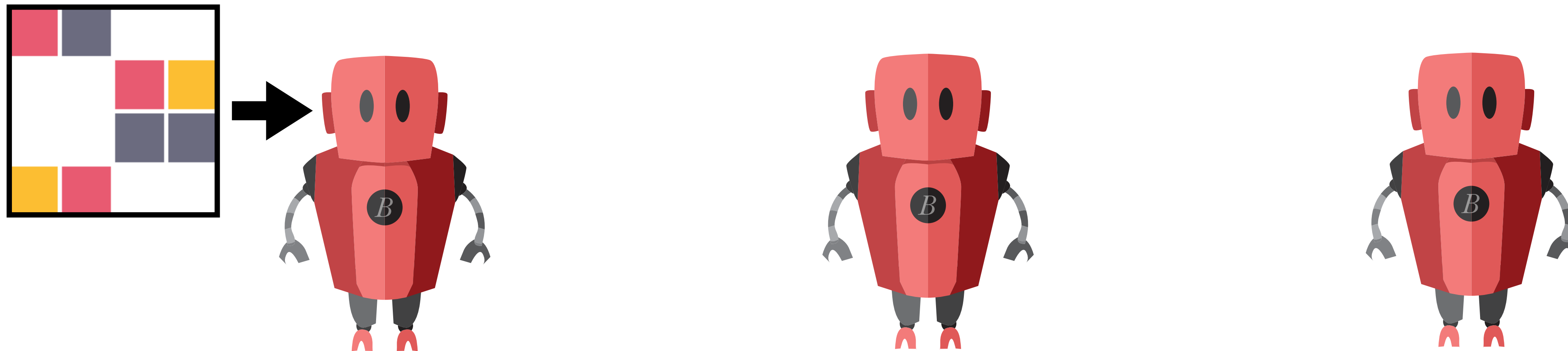
76.58 bits



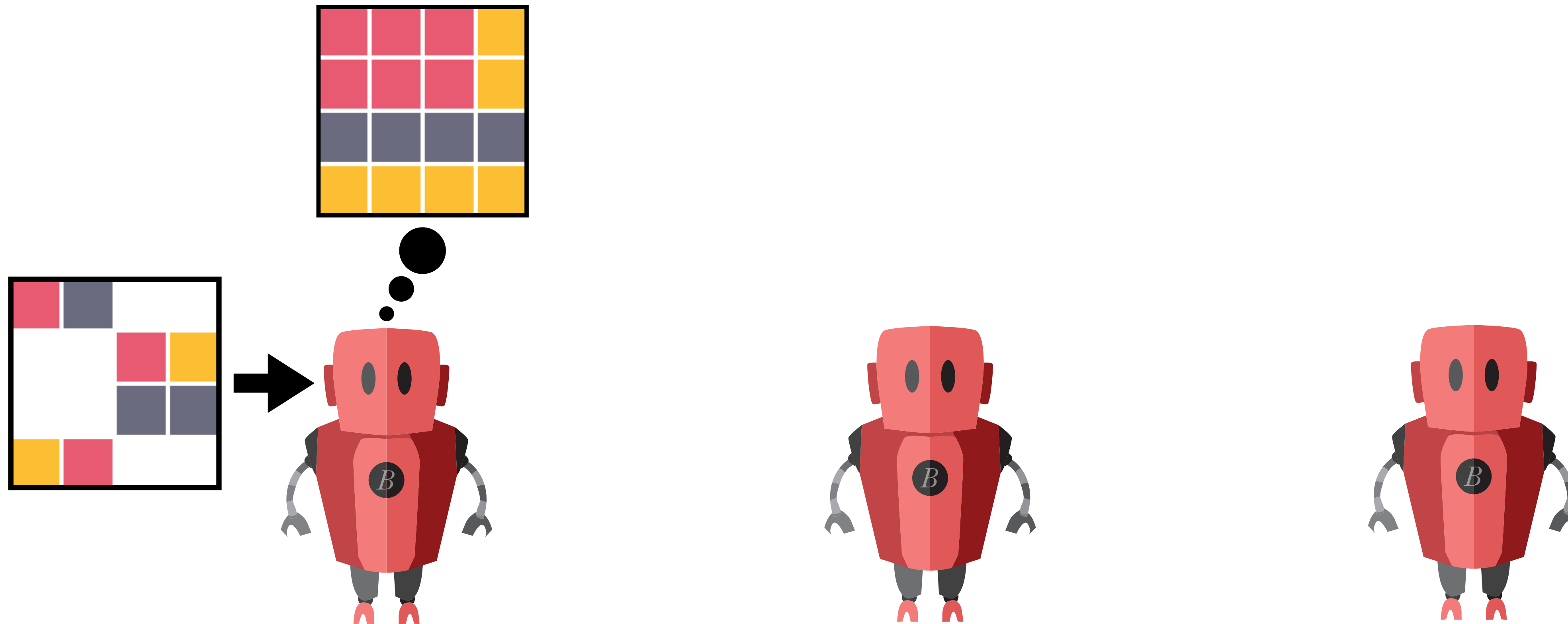
110000  
110001  
110010  
110011

24 bits

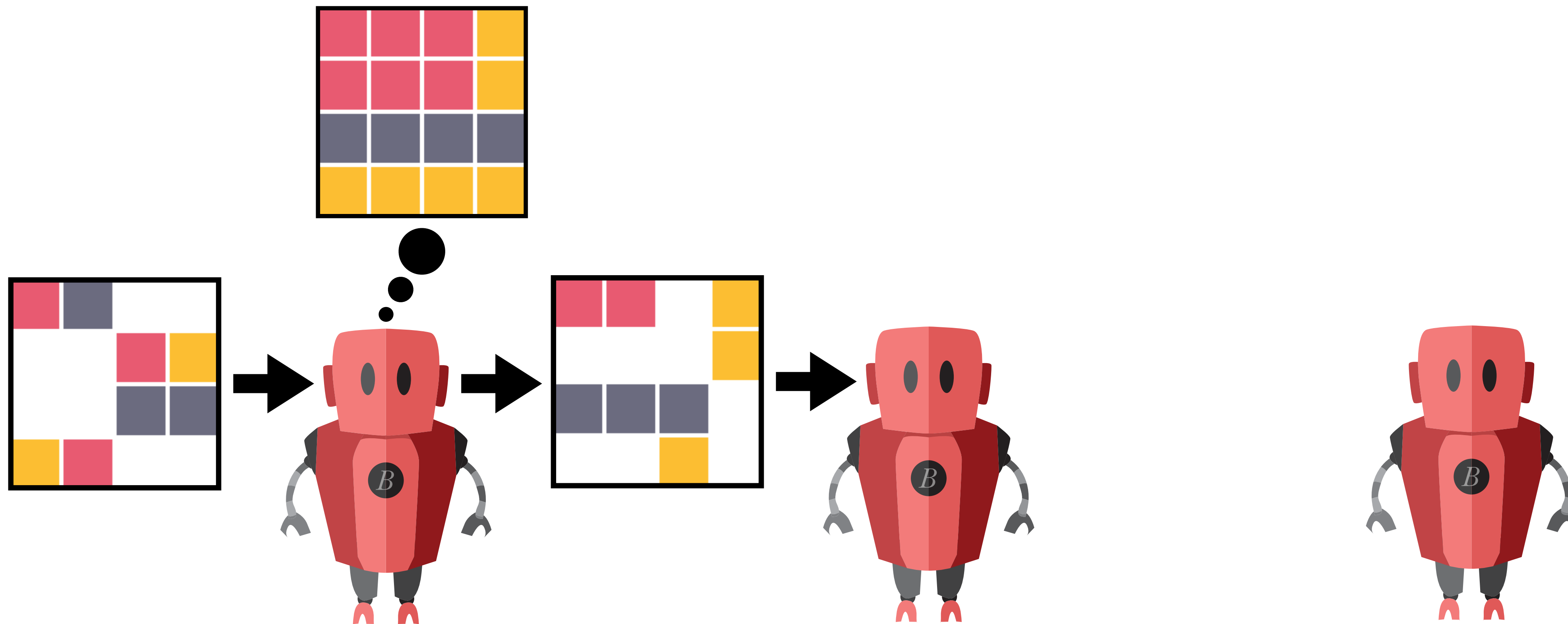
# Bayesian iterated learning under a simplicity prior



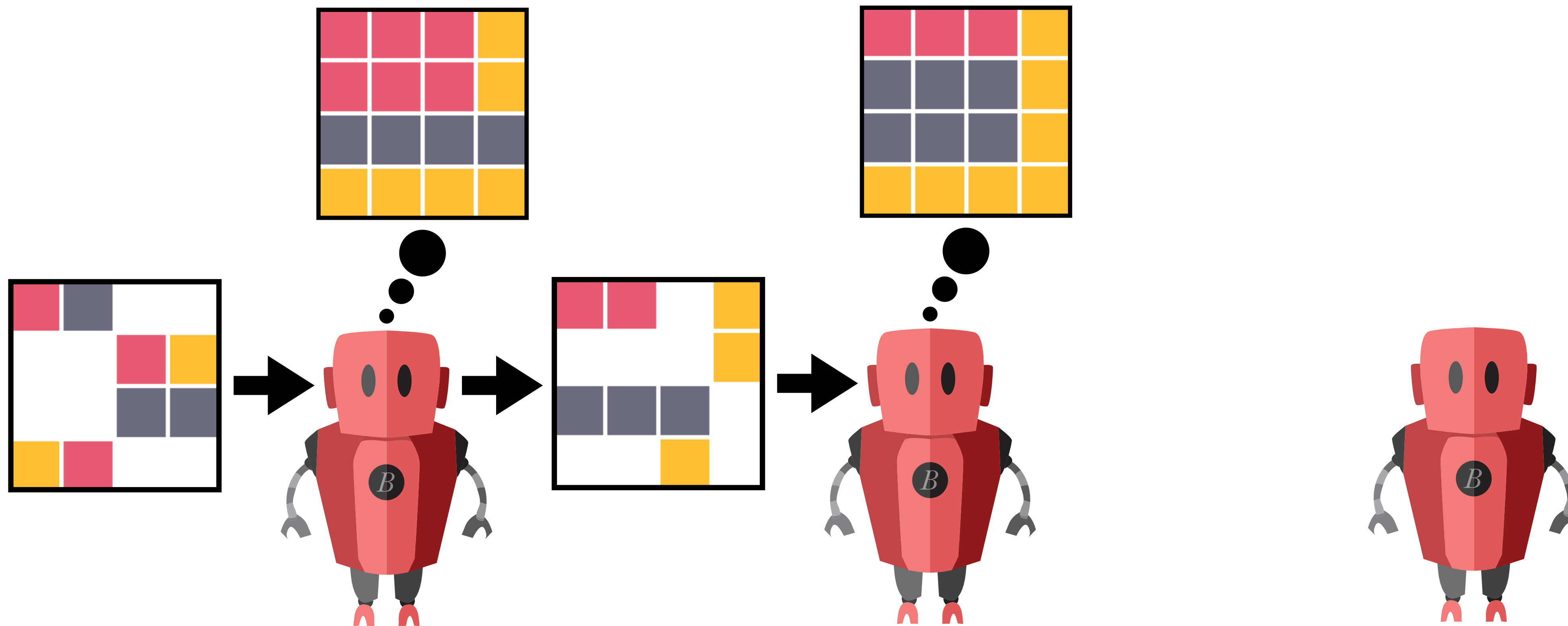
# Bayesian iterated learning under a simplicity prior



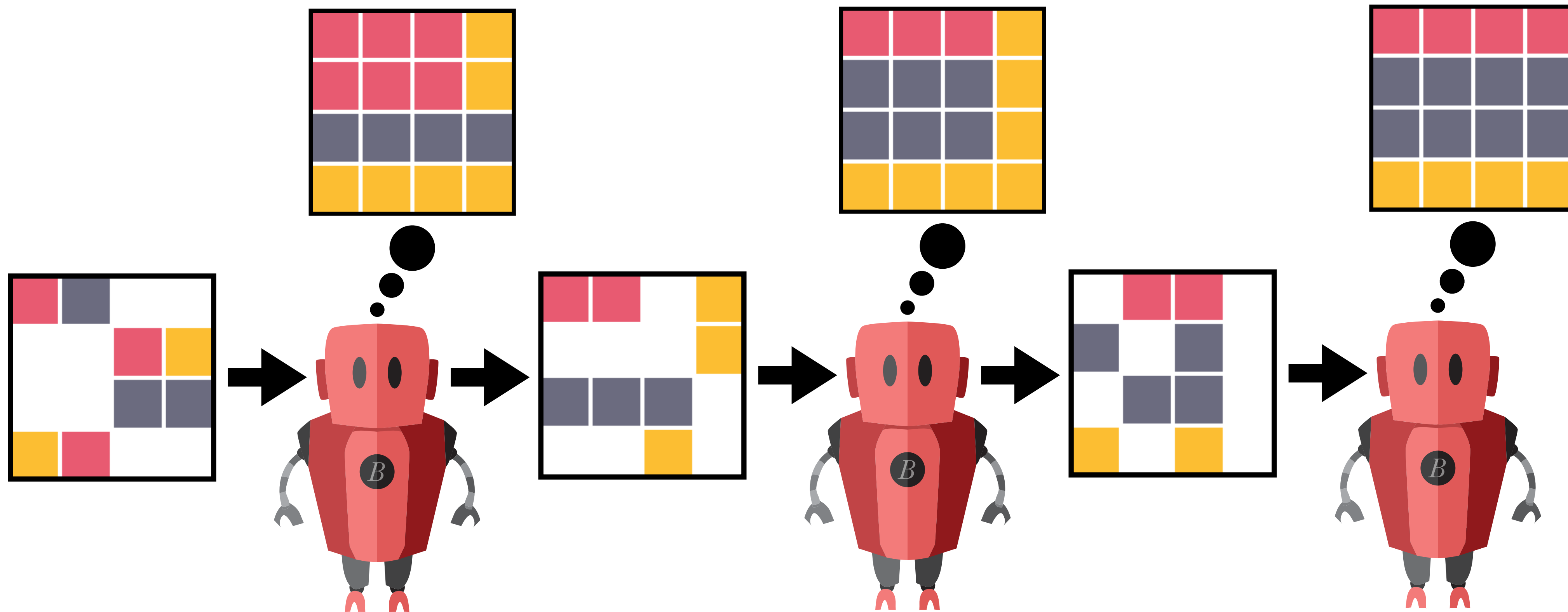
# Bayesian iterated learning under a simplicity prior



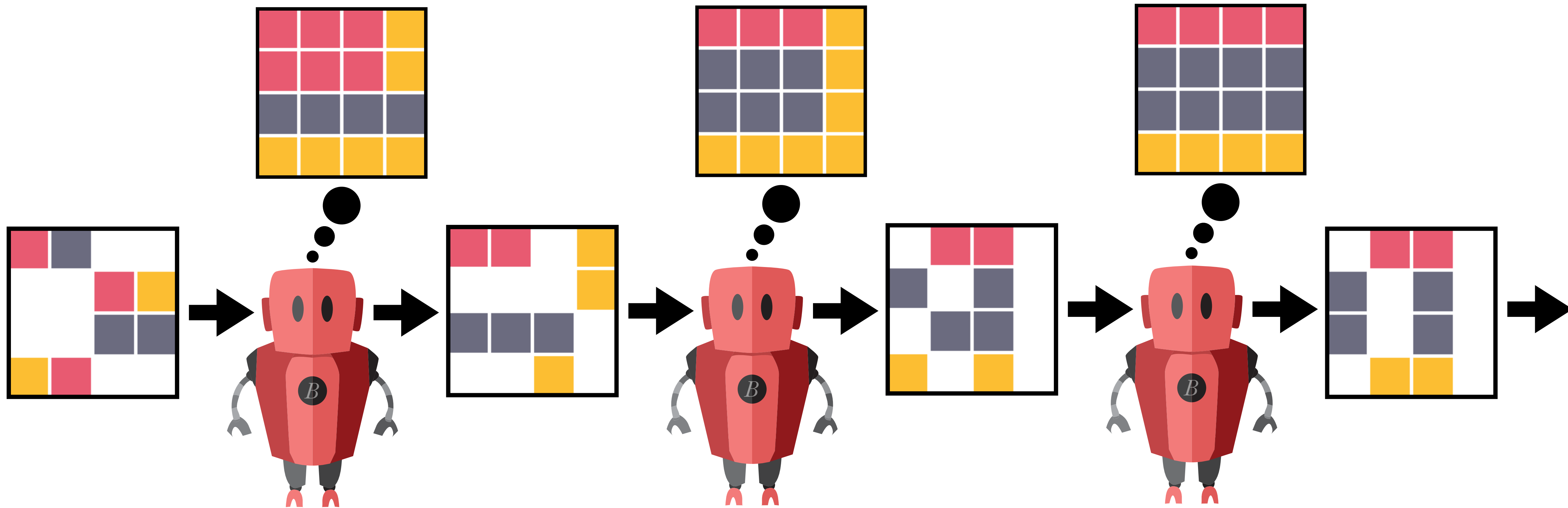
# Bayesian iterated learning under a simplicity prior



# Bayesian iterated learning under a simplicity prior

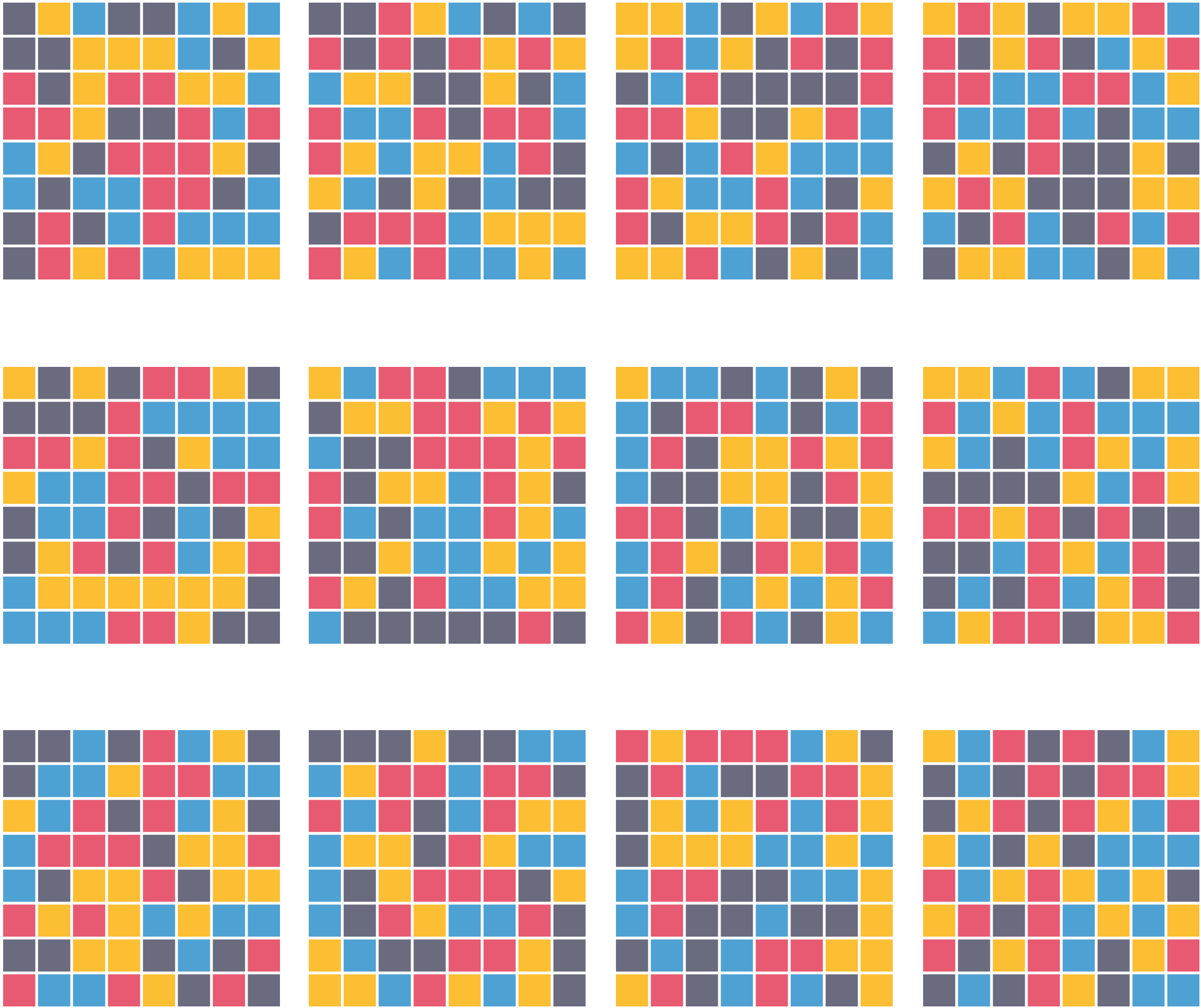


# Bayesian iterated learning under a simplicity prior

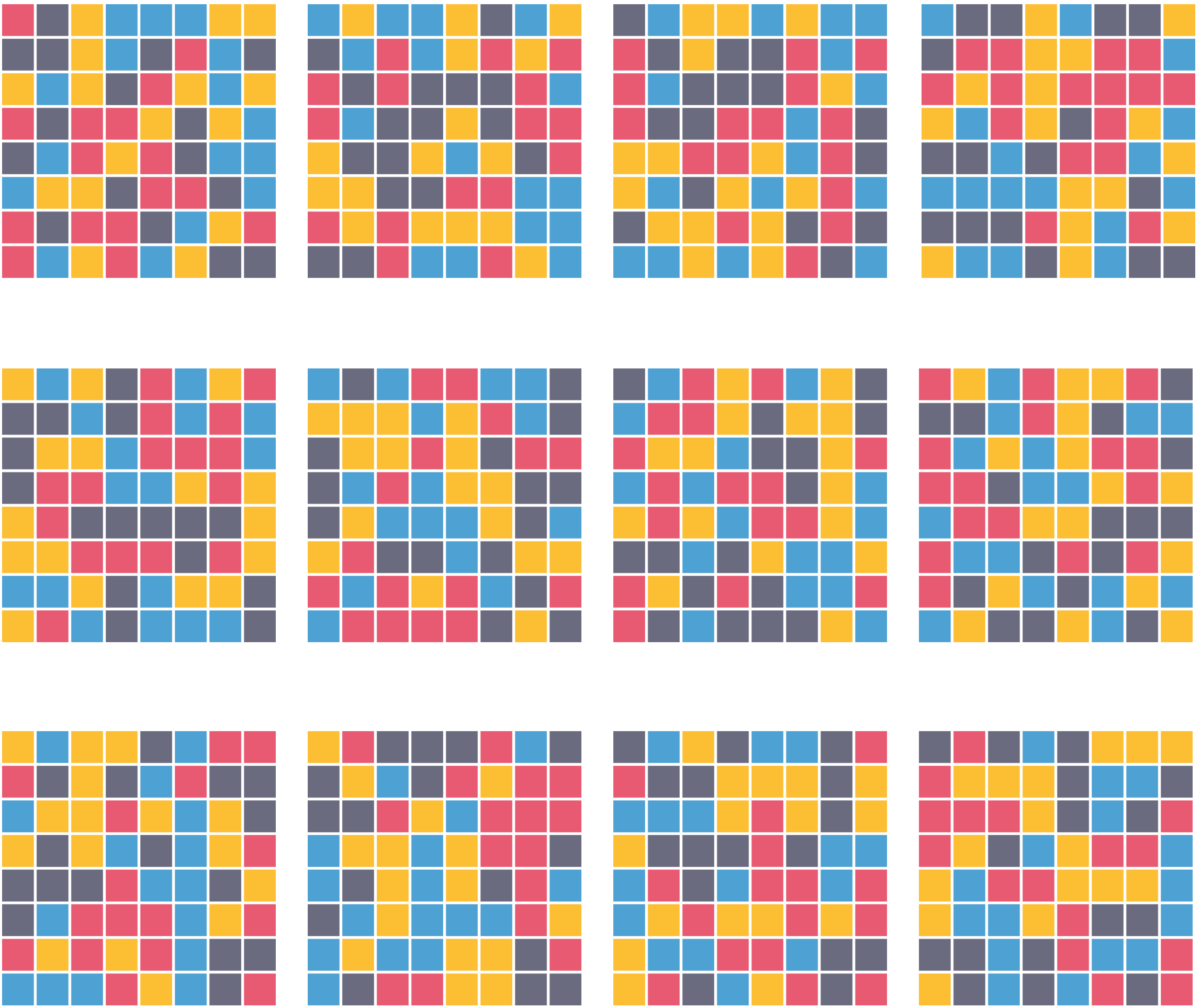




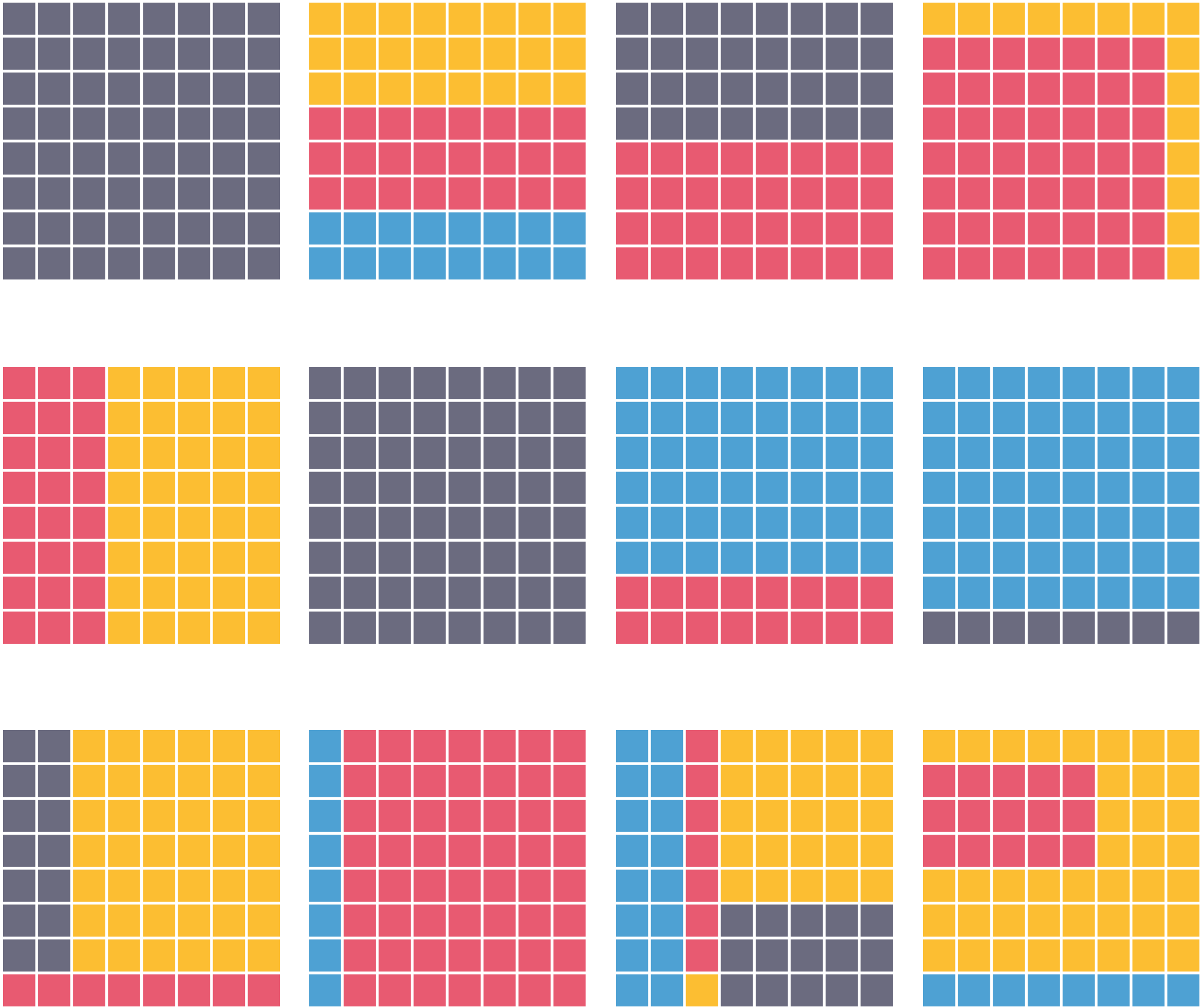
Simplicity prior



Informativeness prior



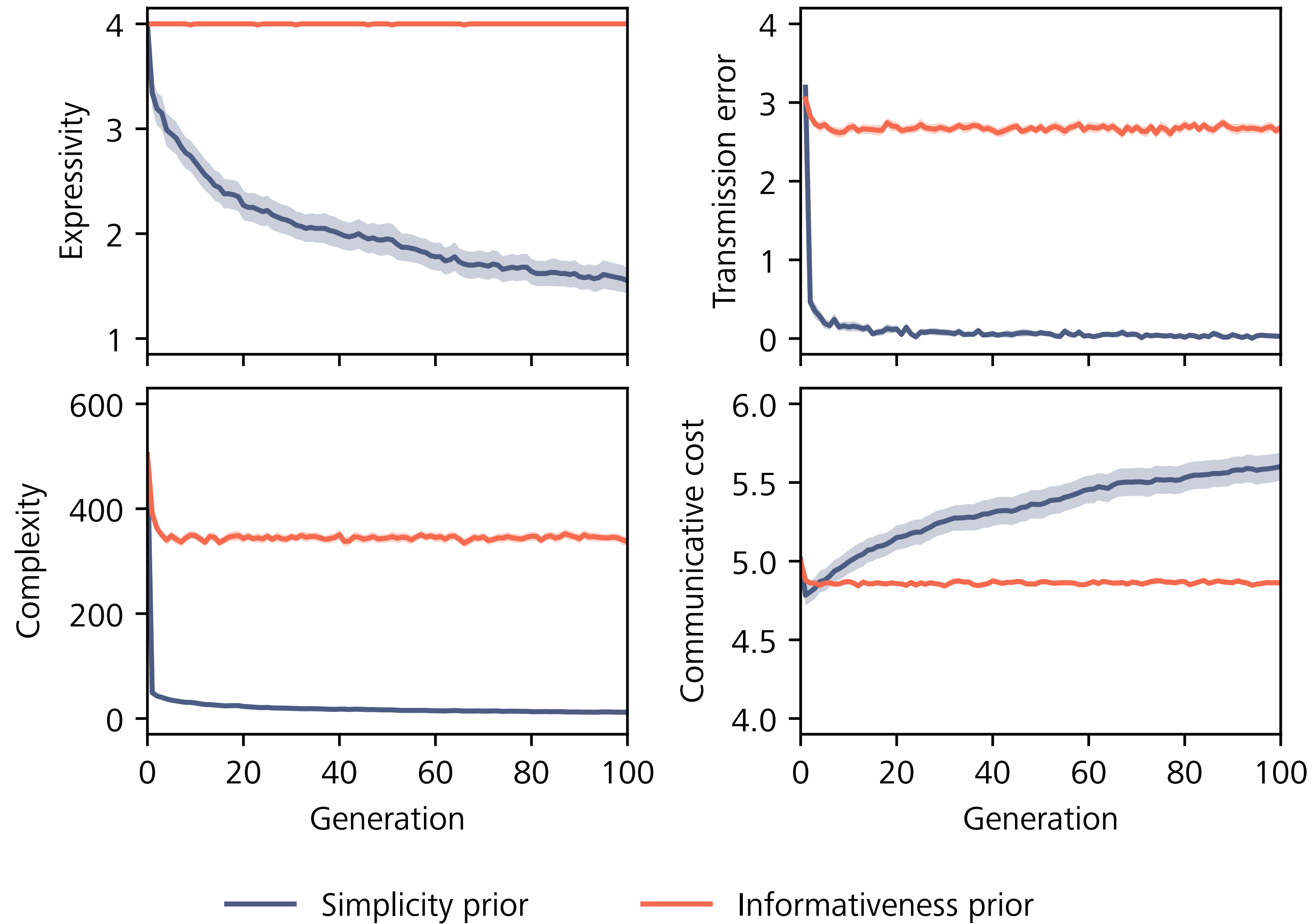
Simplicity prior



Informativeness prior



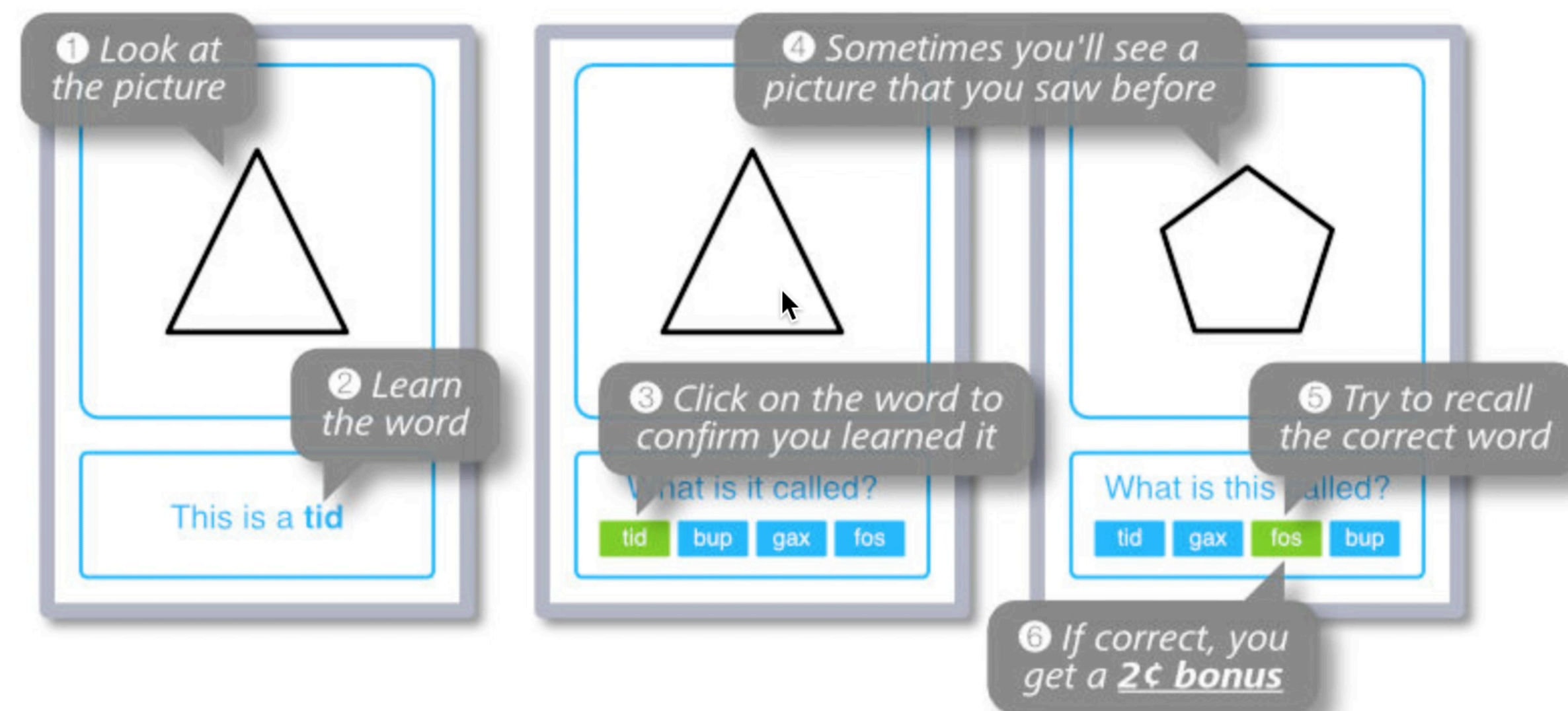
# Model results



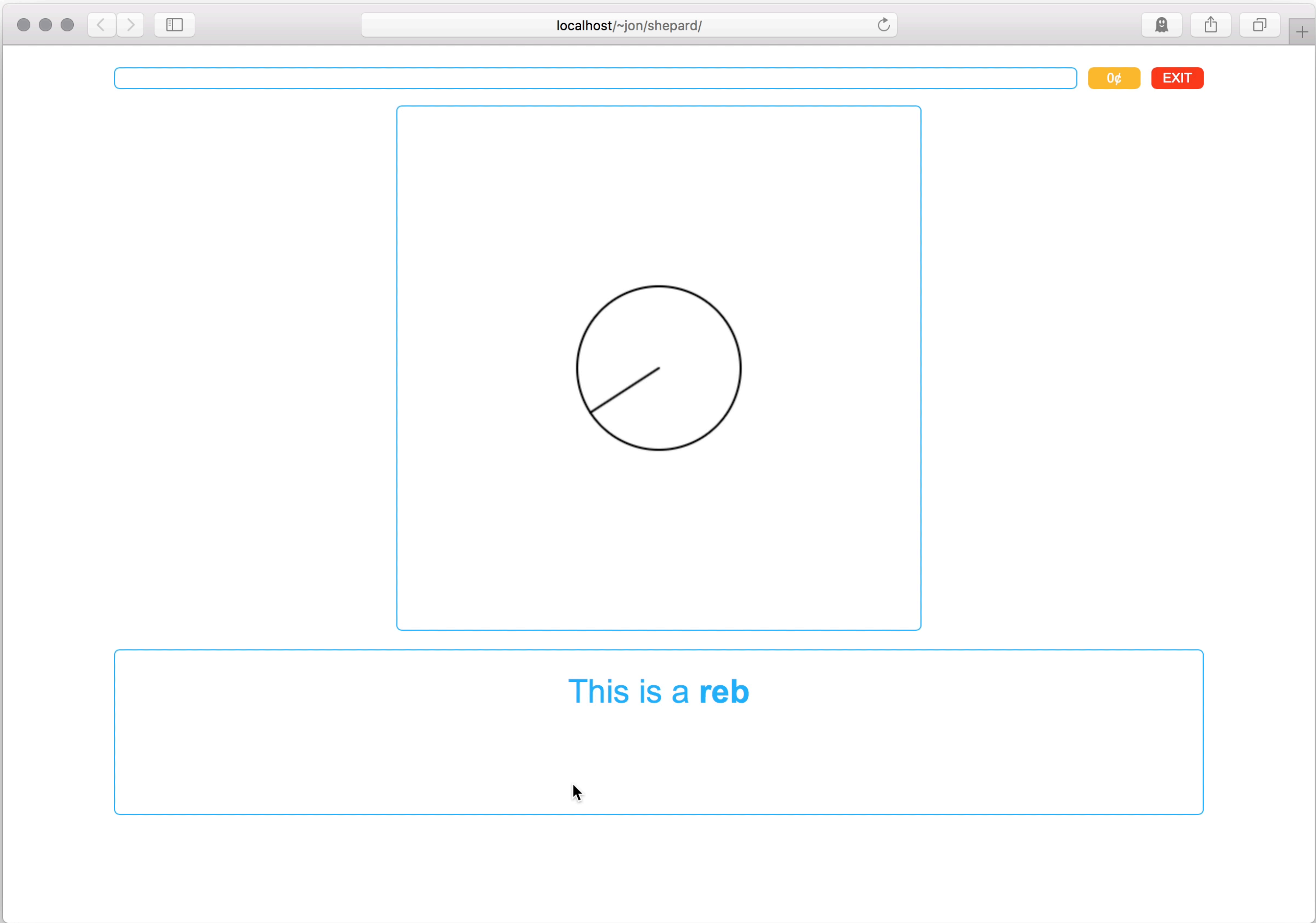
*Experiment*

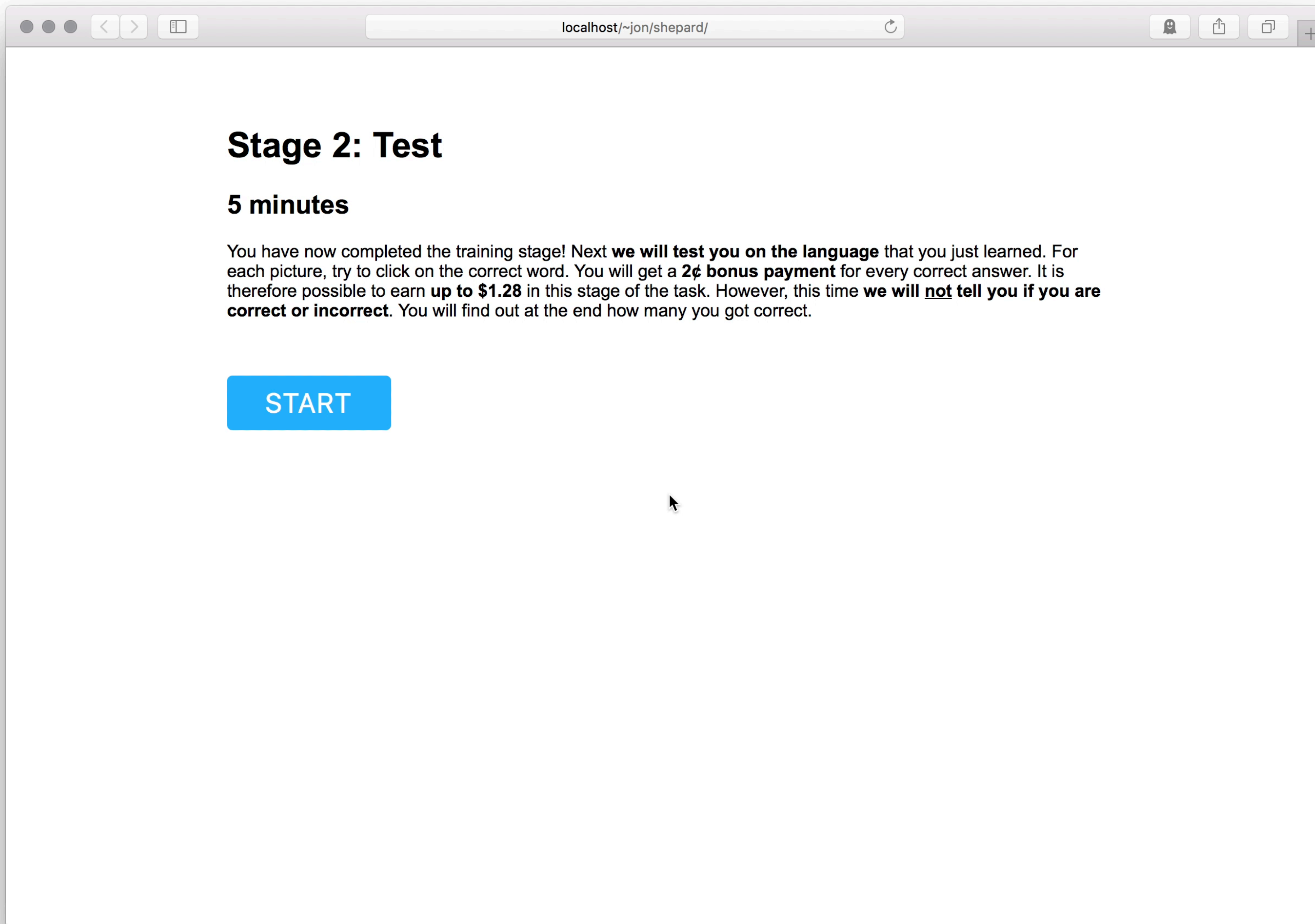
## 15 minutes

You are going to learn a simple language. **We will train you on 4 words** in the language and **we will test how well you are learning the words**. Try to learn the language as well as you can and **aim to be accurate in your answers**. You will receive a **2¢ bonus payment** for every correct test answer. If you decide to stop the task, please click the **EXIT** button so that someone else can take part.



START





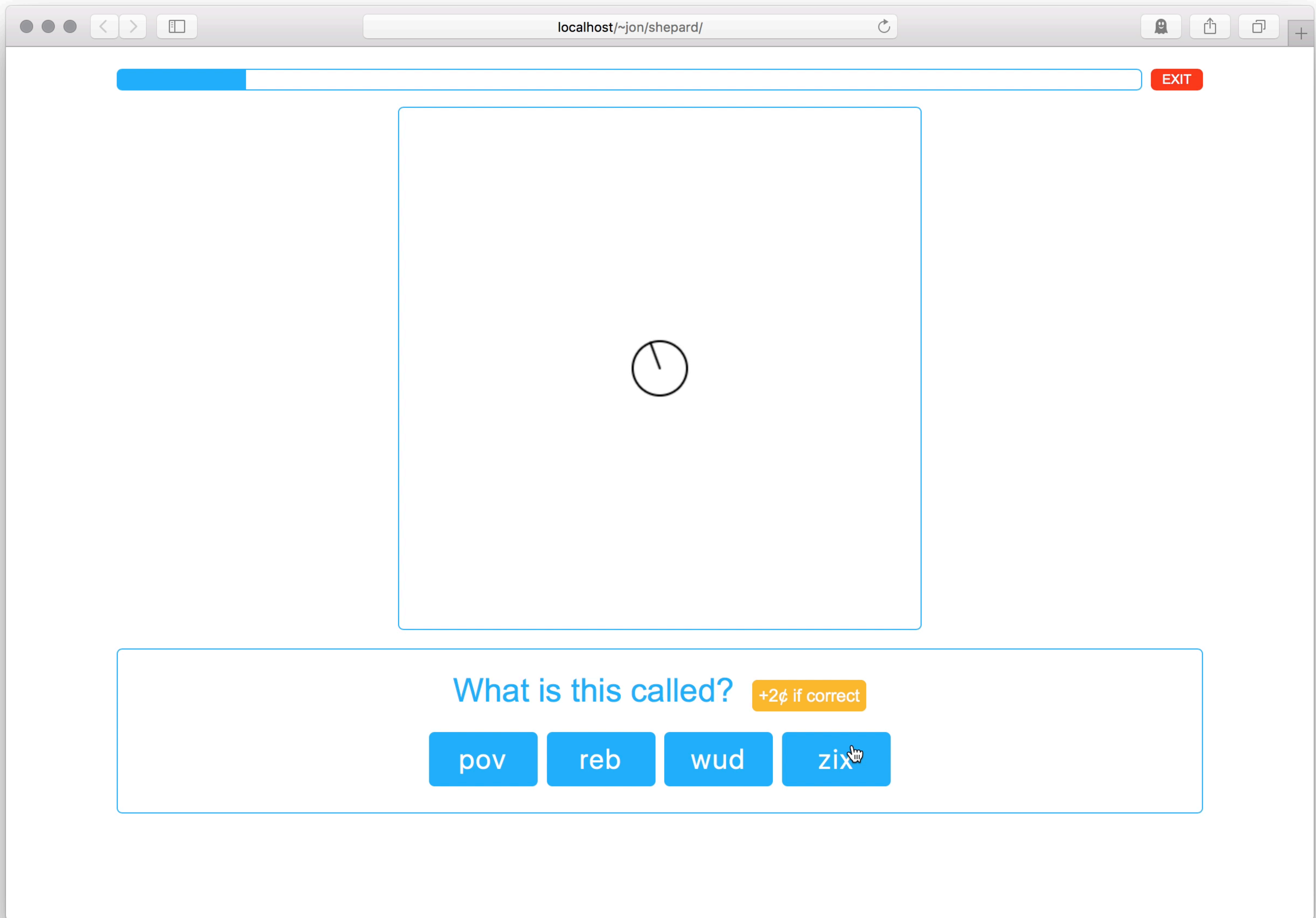
## Stage 2: Test

### 5 minutes

You have now completed the training stage! Next **we will test you on the language** that you just learned. For each picture, try to click on the correct word. You will get a **2¢ bonus payment** for every correct answer. It is therefore possible to earn **up to \$1.28** in this stage of the task. However, this time **we will not tell you if you are correct or incorrect**. You will find out at the end how many you got correct.

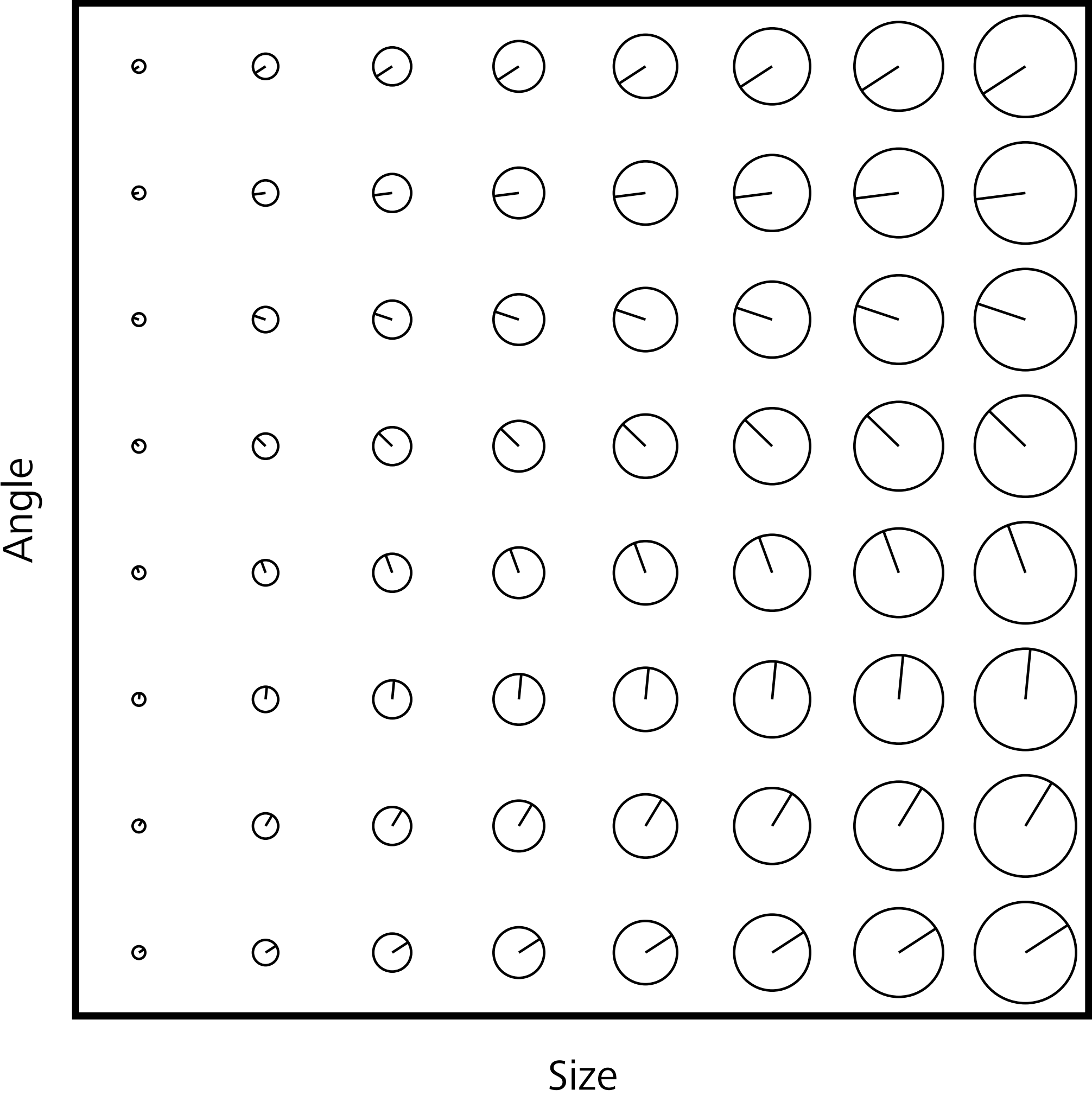
START





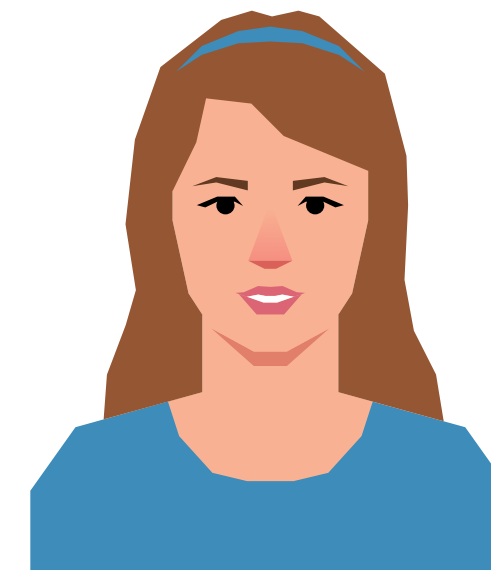
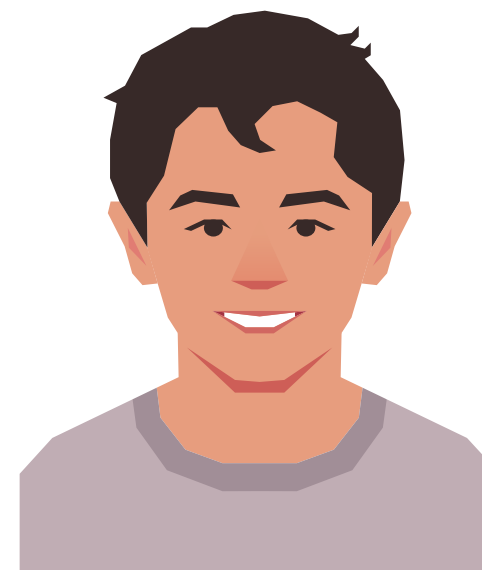
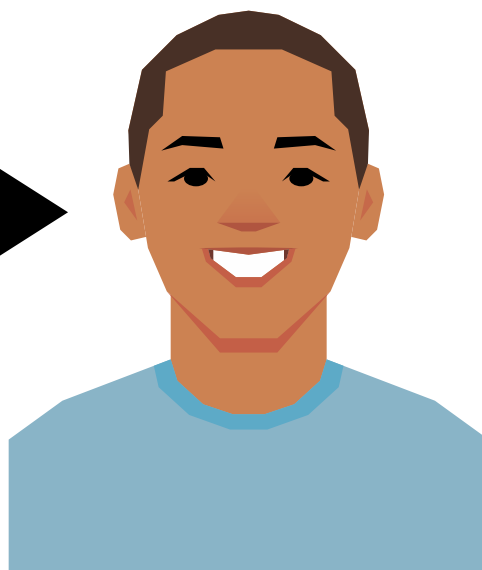
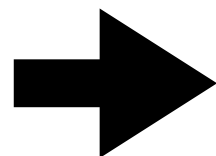
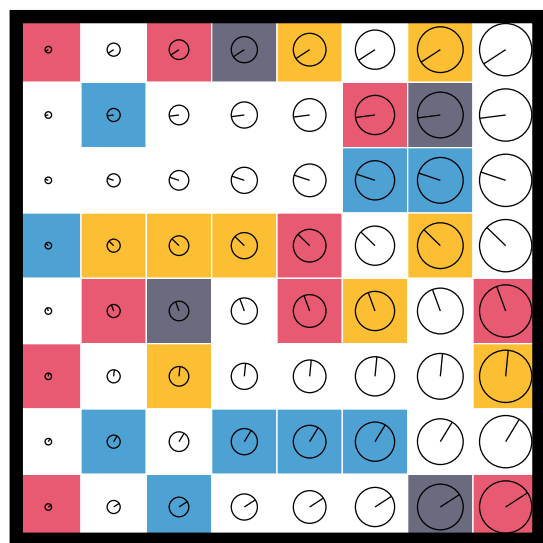


# Experimental stimuli

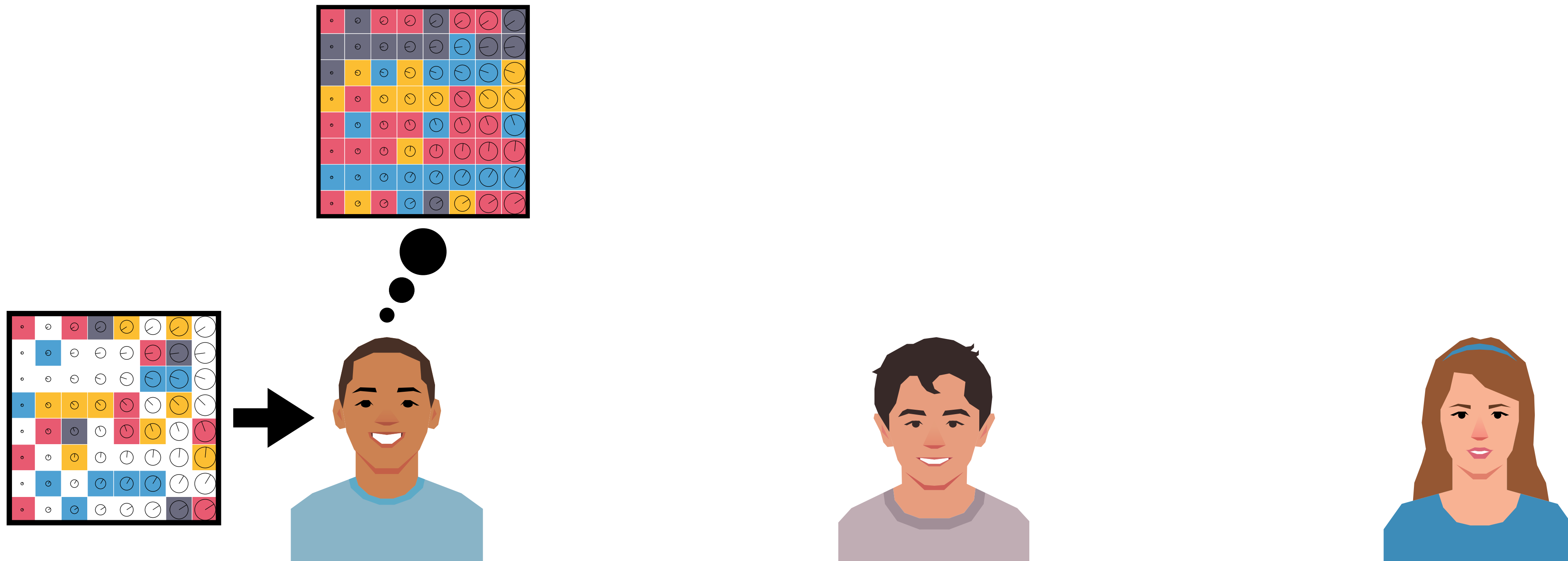


Set	Labels			
1	<i>pov</i>	<i>reb</i>	<i>wud</i>	<i>zix</i>
2	<i>gex</i>	<i>juf</i>	<i>vib</i>	<i>wop</i>
3	<i>buv</i>	<i>jef</i>	<i>pid</i>	<i>zox</i>
4	<i>fod</i>	<i>jes</i>	<i>wix</i>	<i>zuv</i>

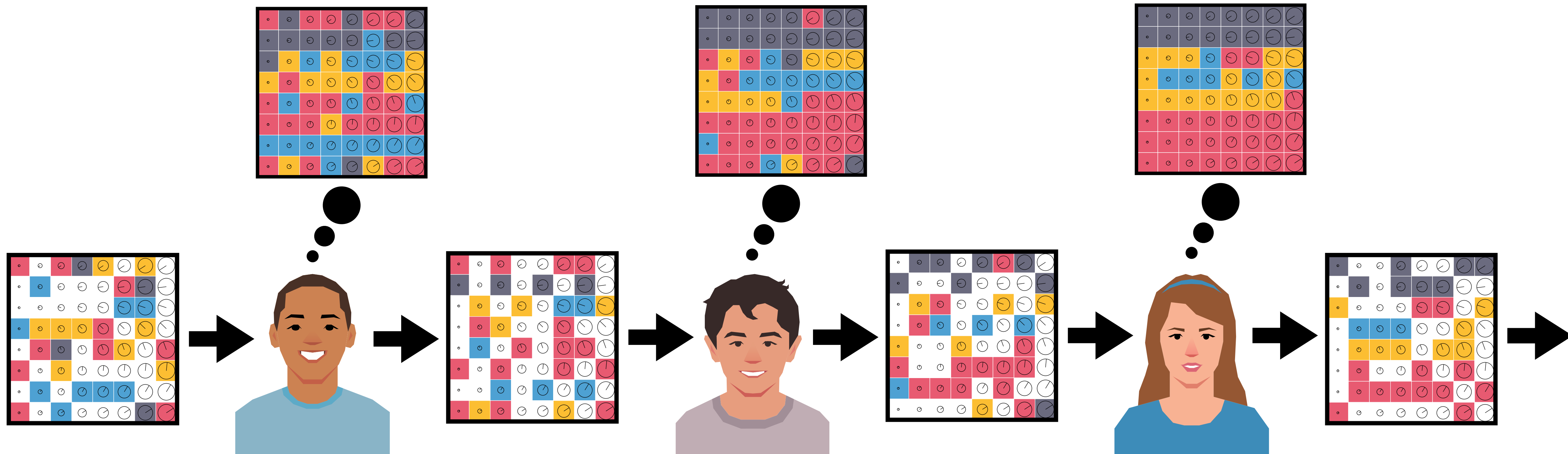
# Iterated learning with humans

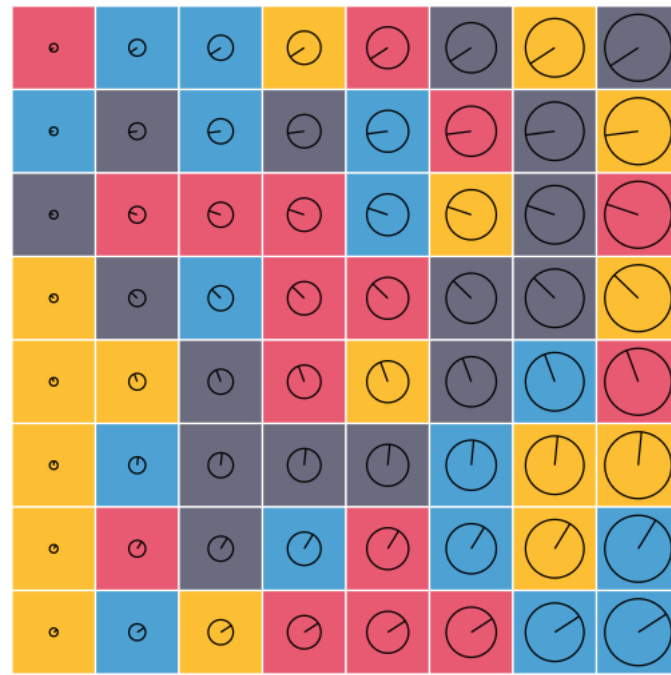
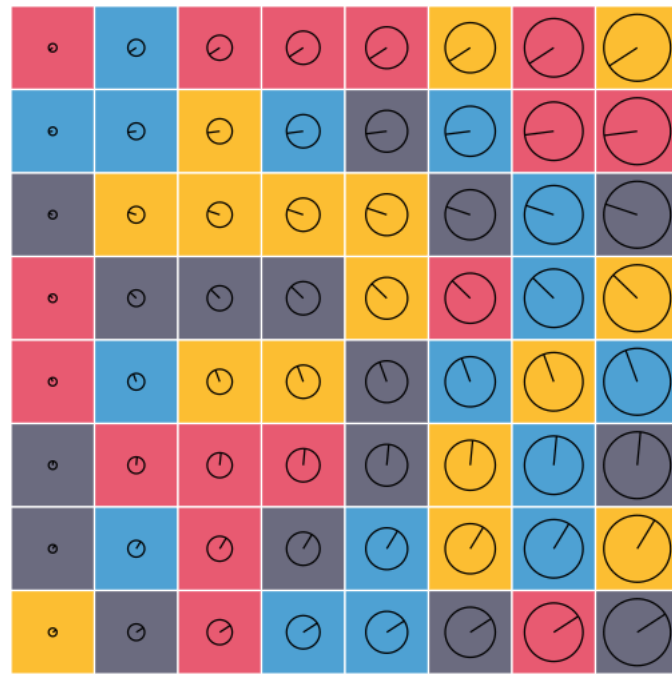
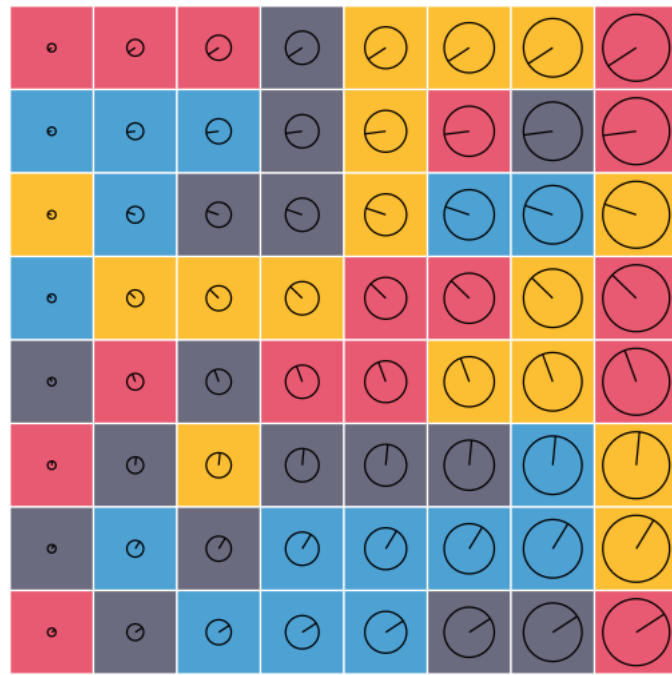
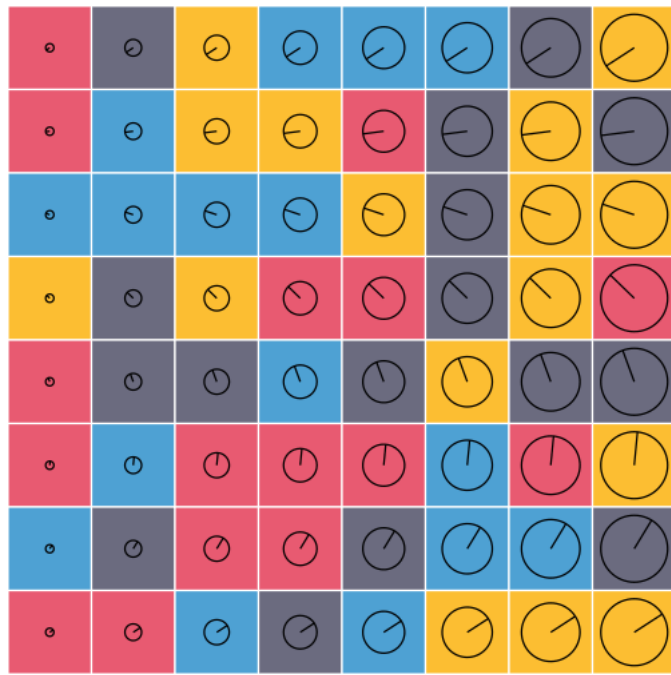
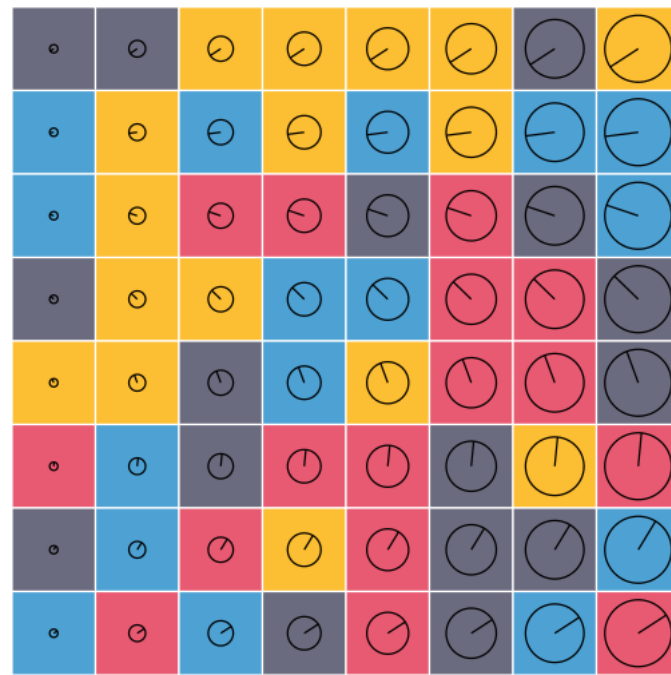
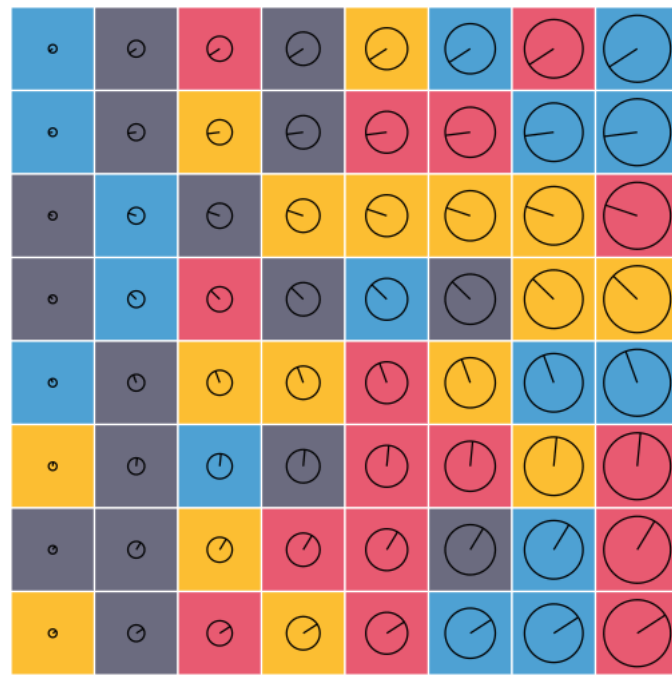
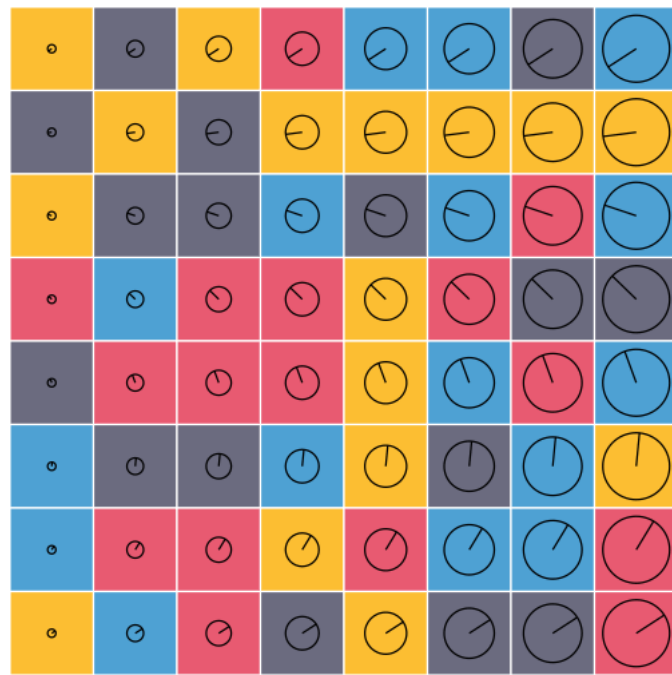
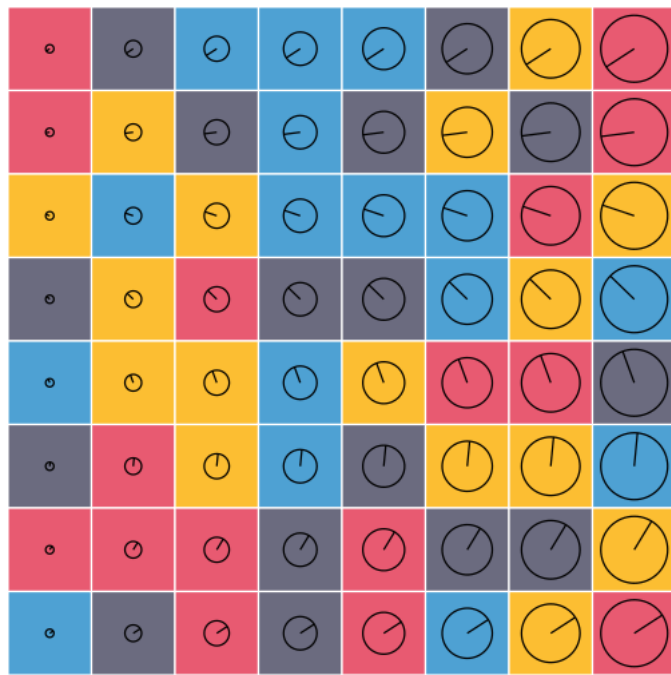
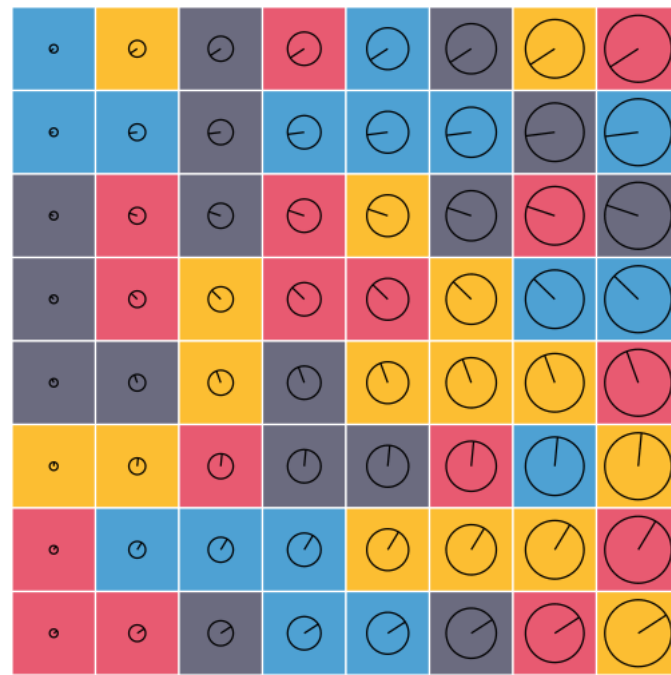
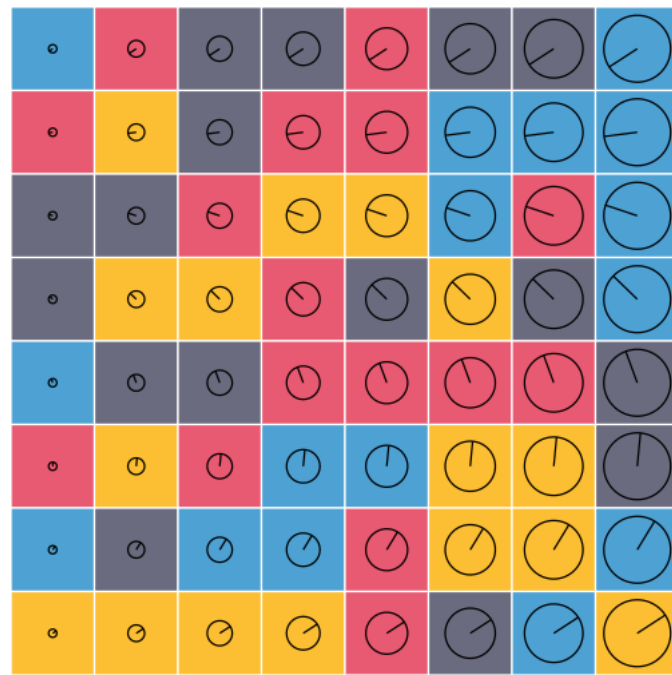
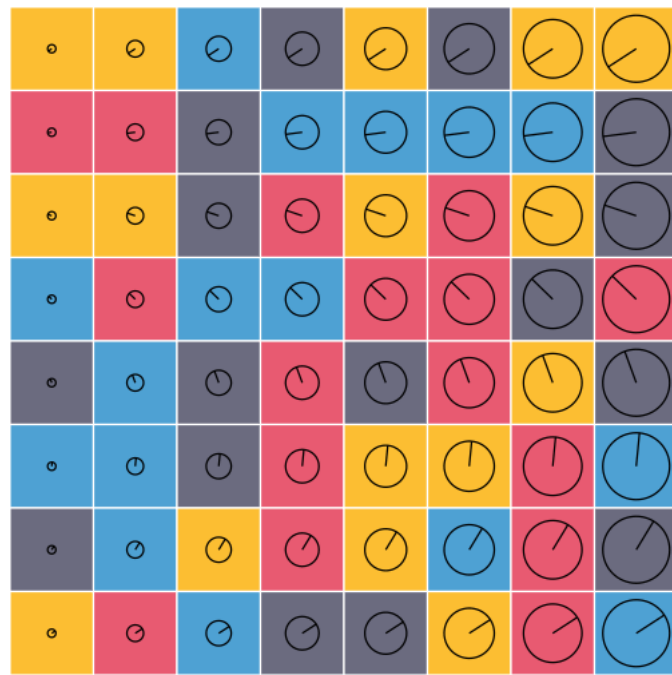
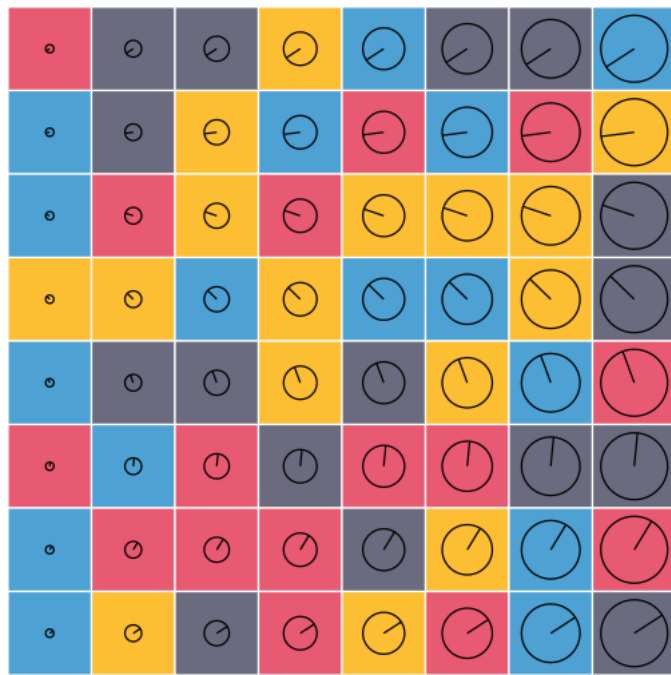


# Iterated learning with humans



# Iterated learning with humans



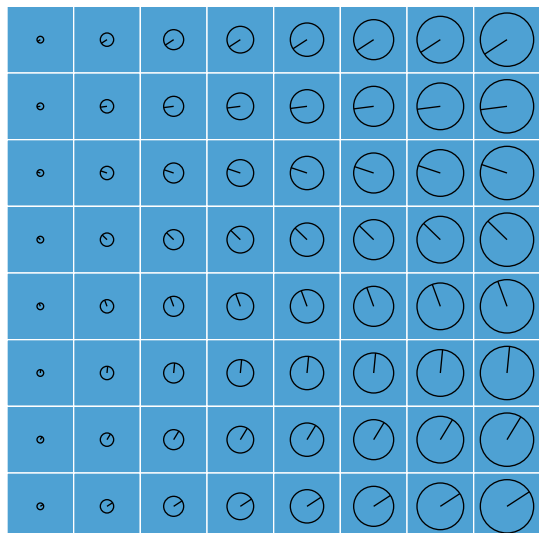
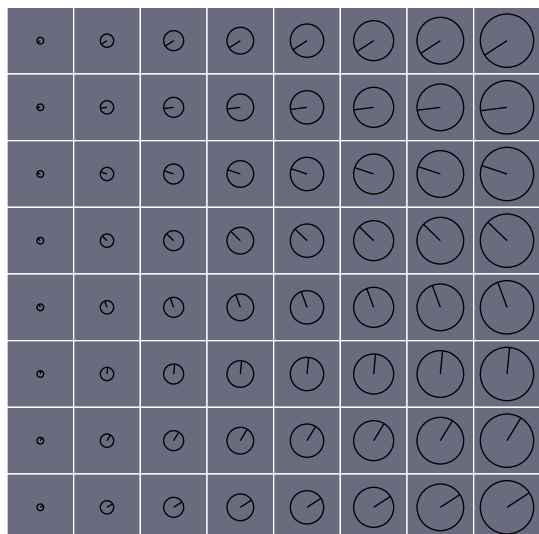




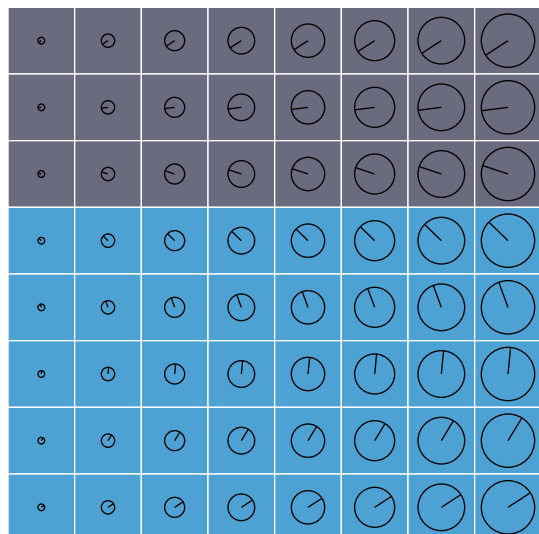


# Systems converged on

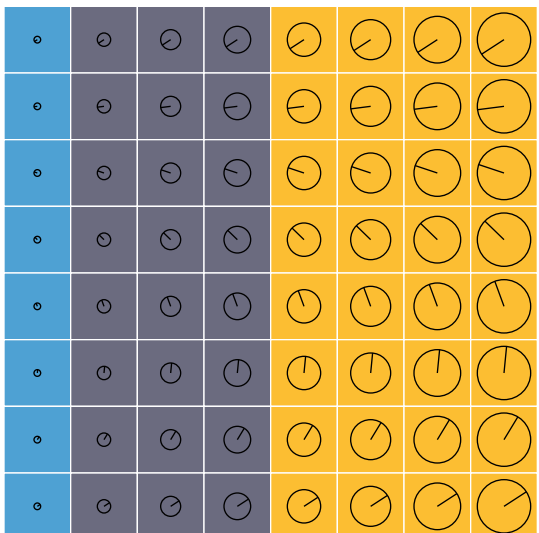
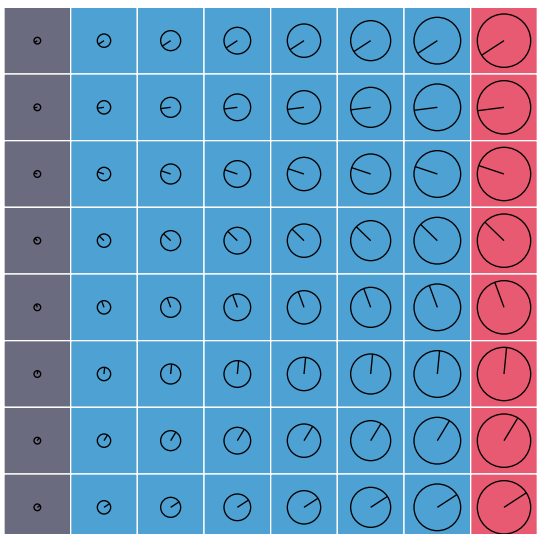
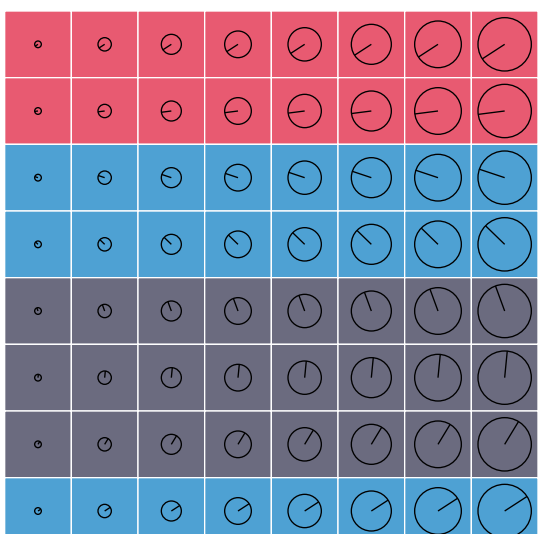
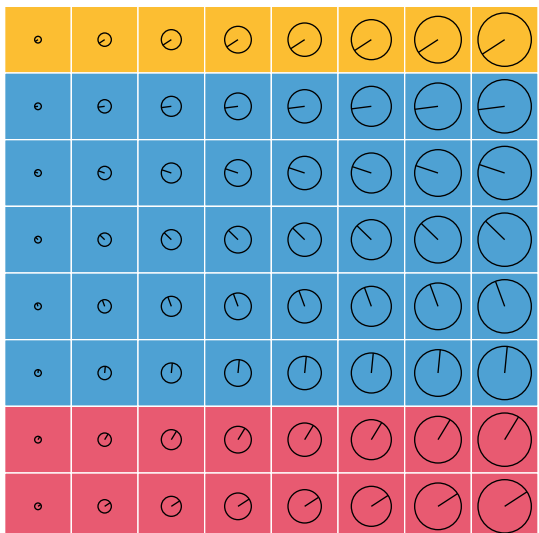
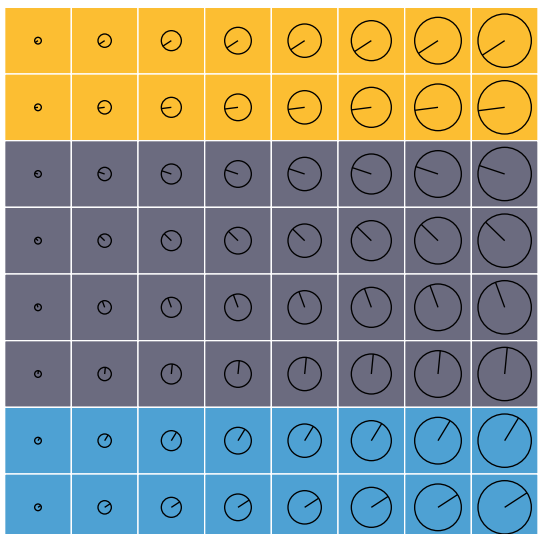
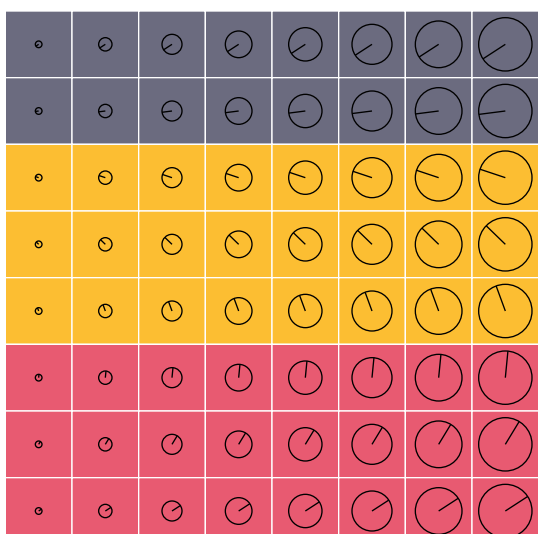
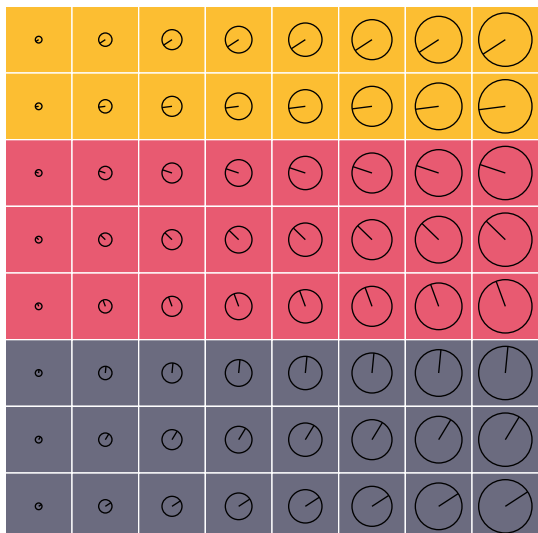
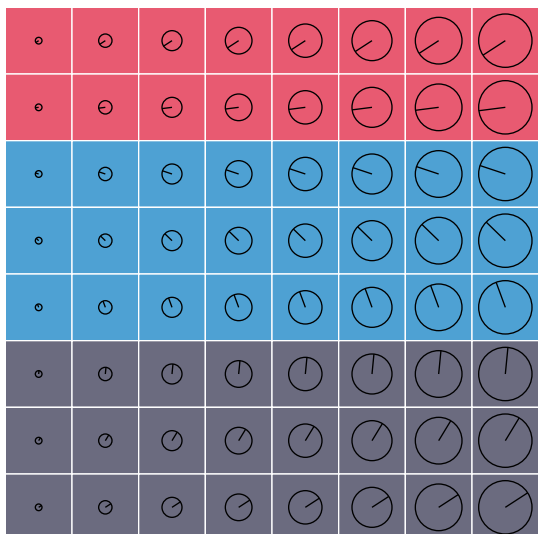
1 concept (2/12)



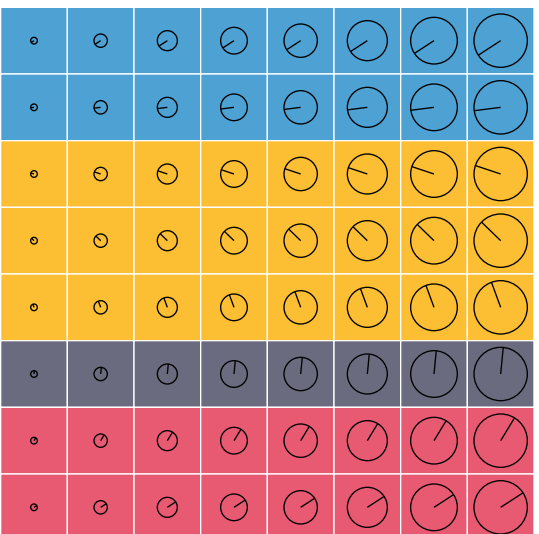
2 concepts (1/12)



3 concepts (8/12)

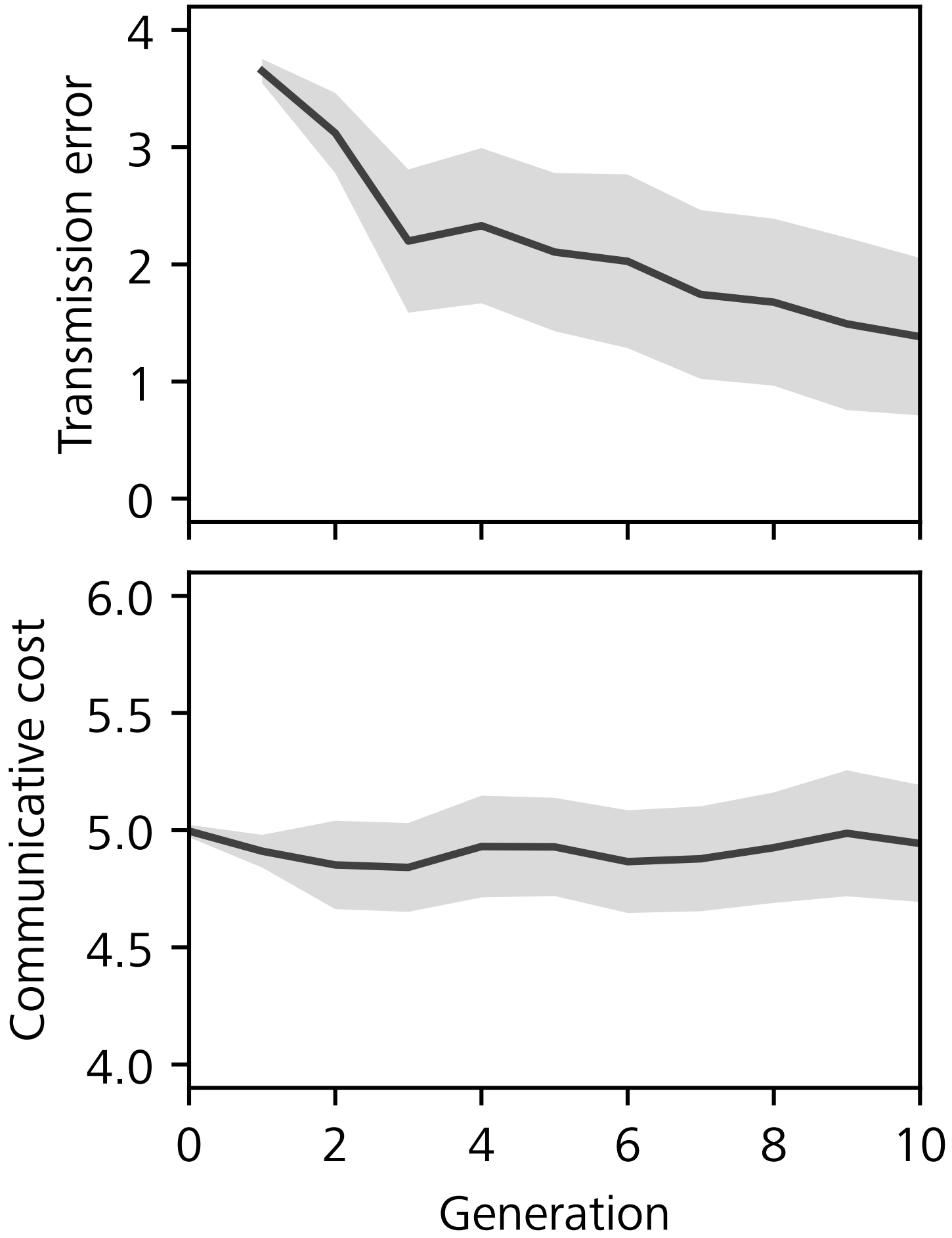
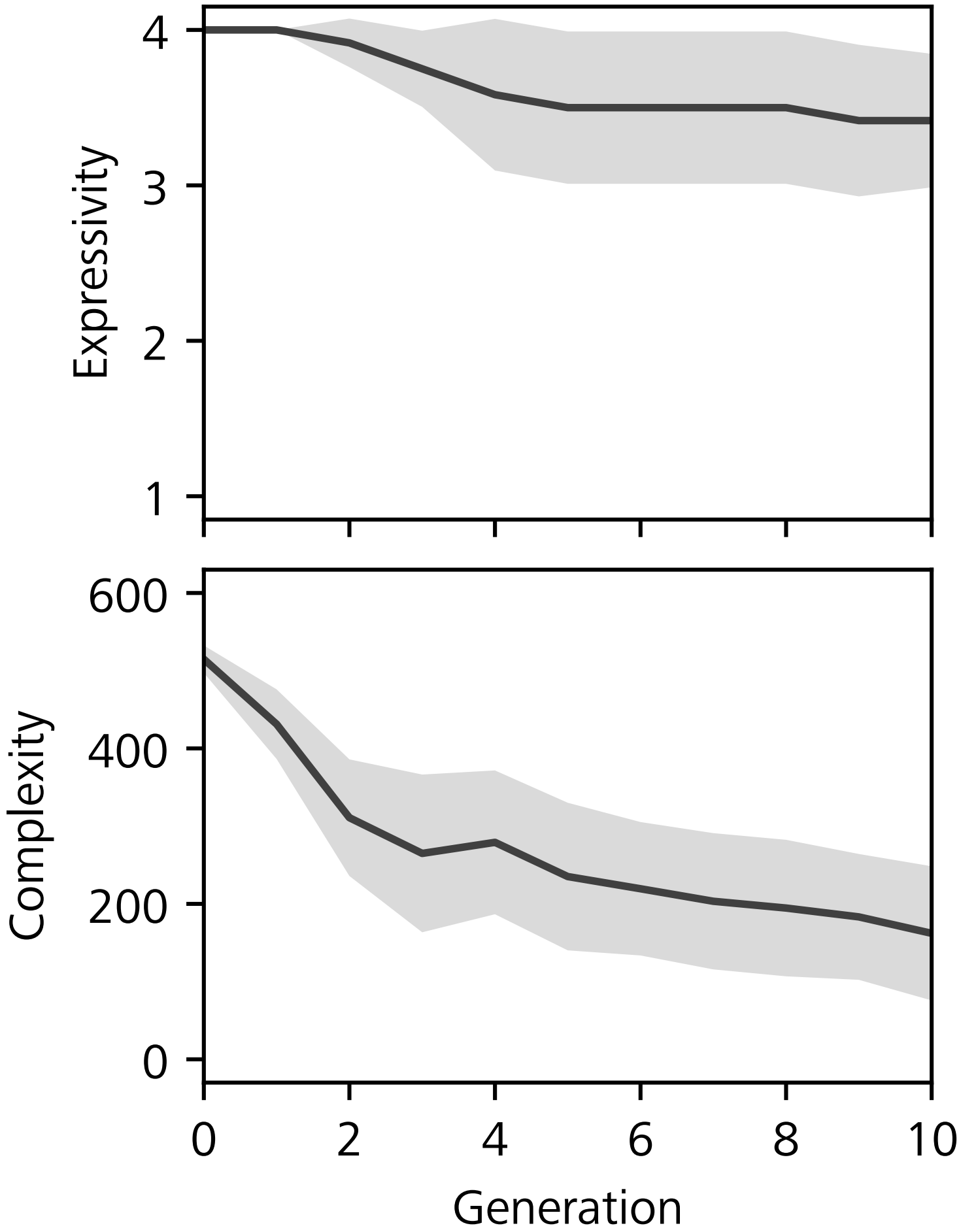
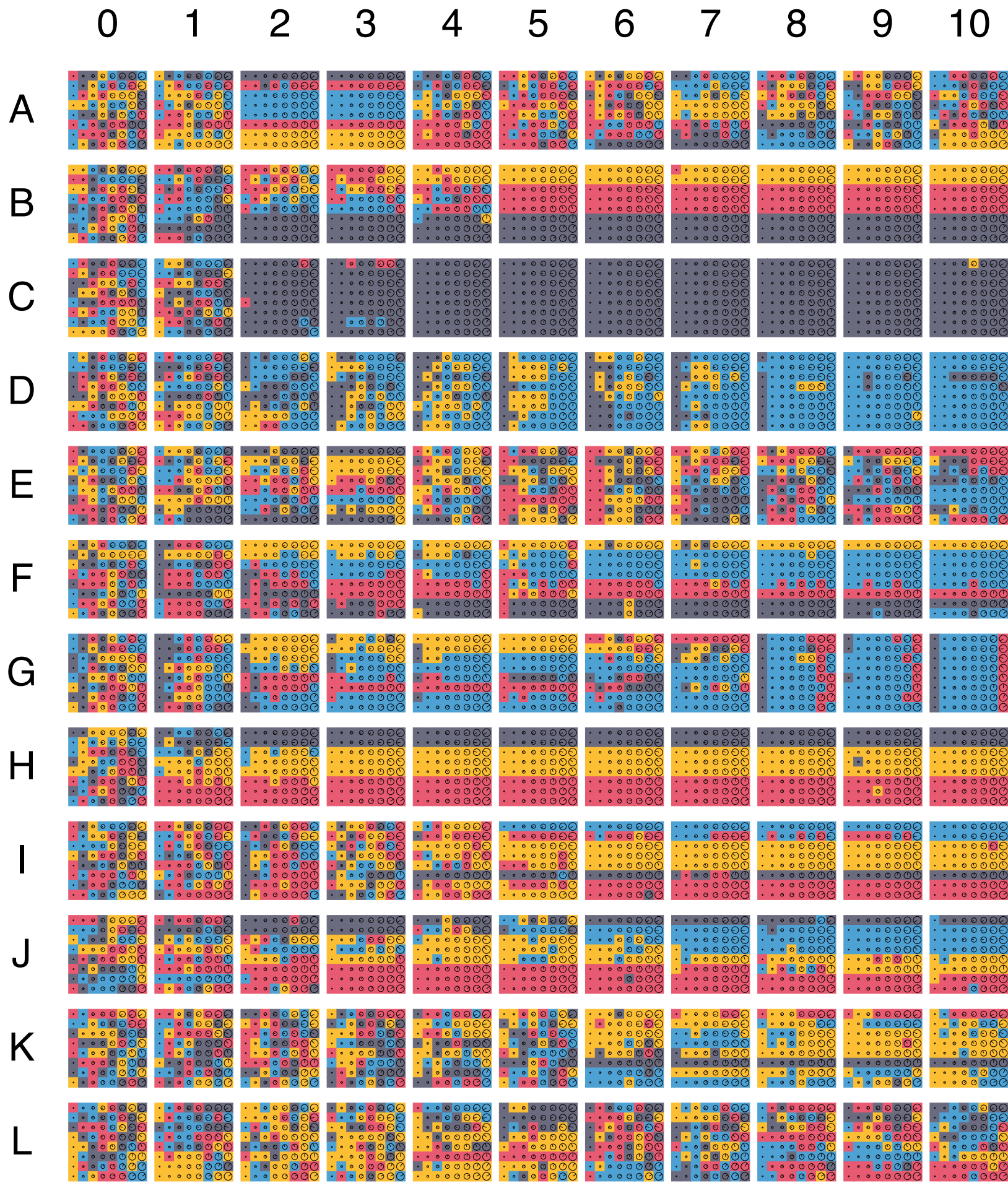


4 concepts (1/12)



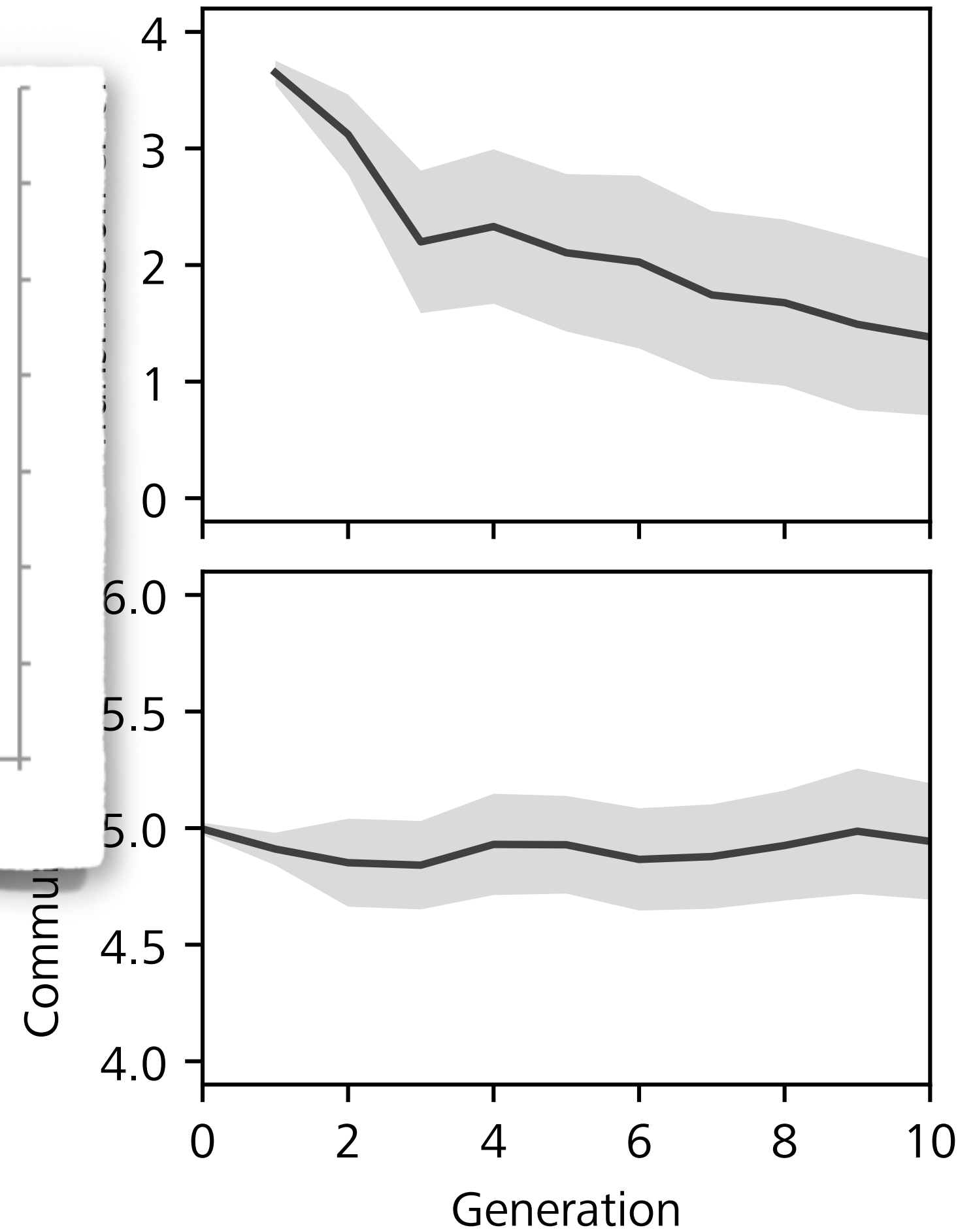
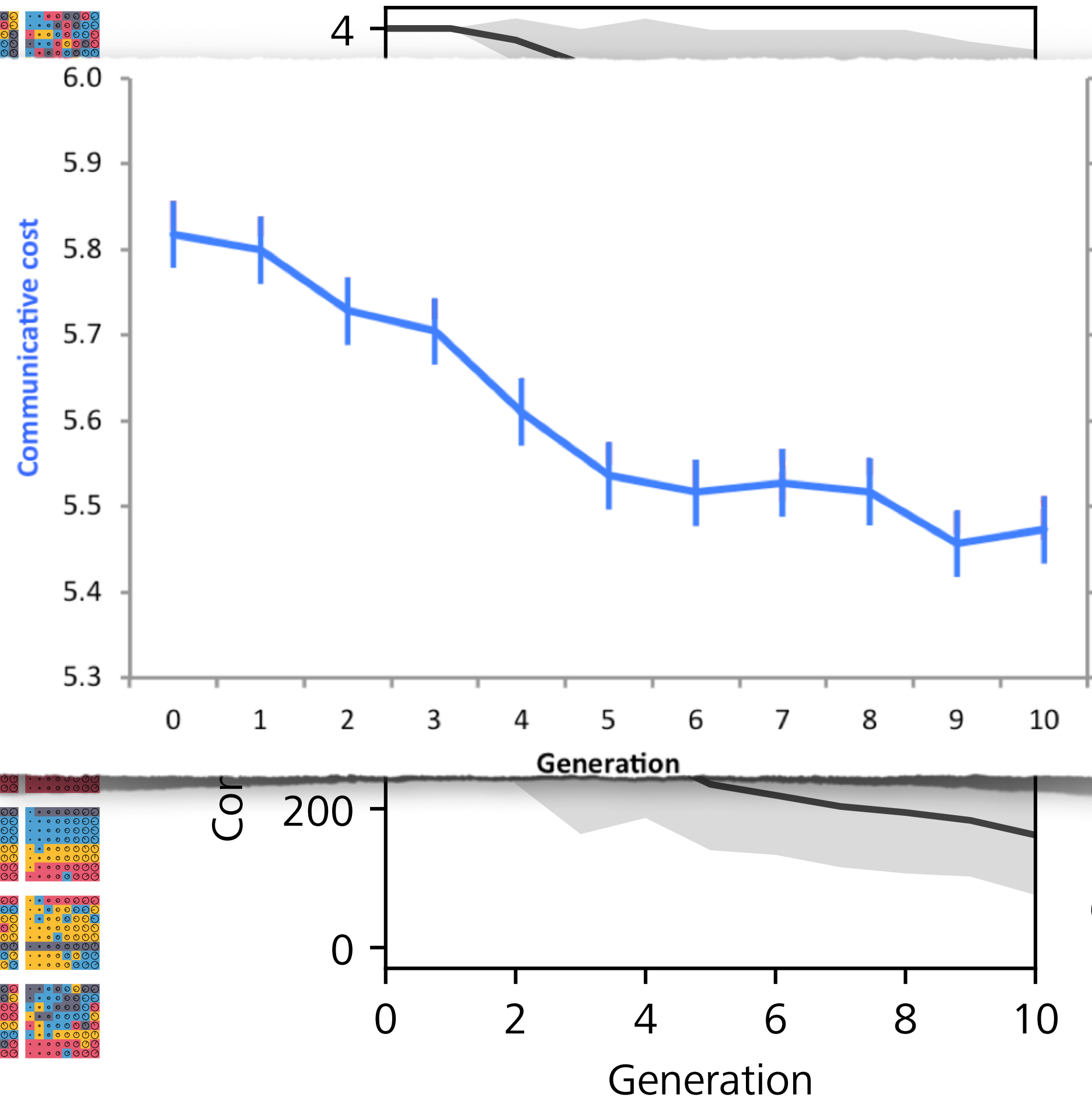
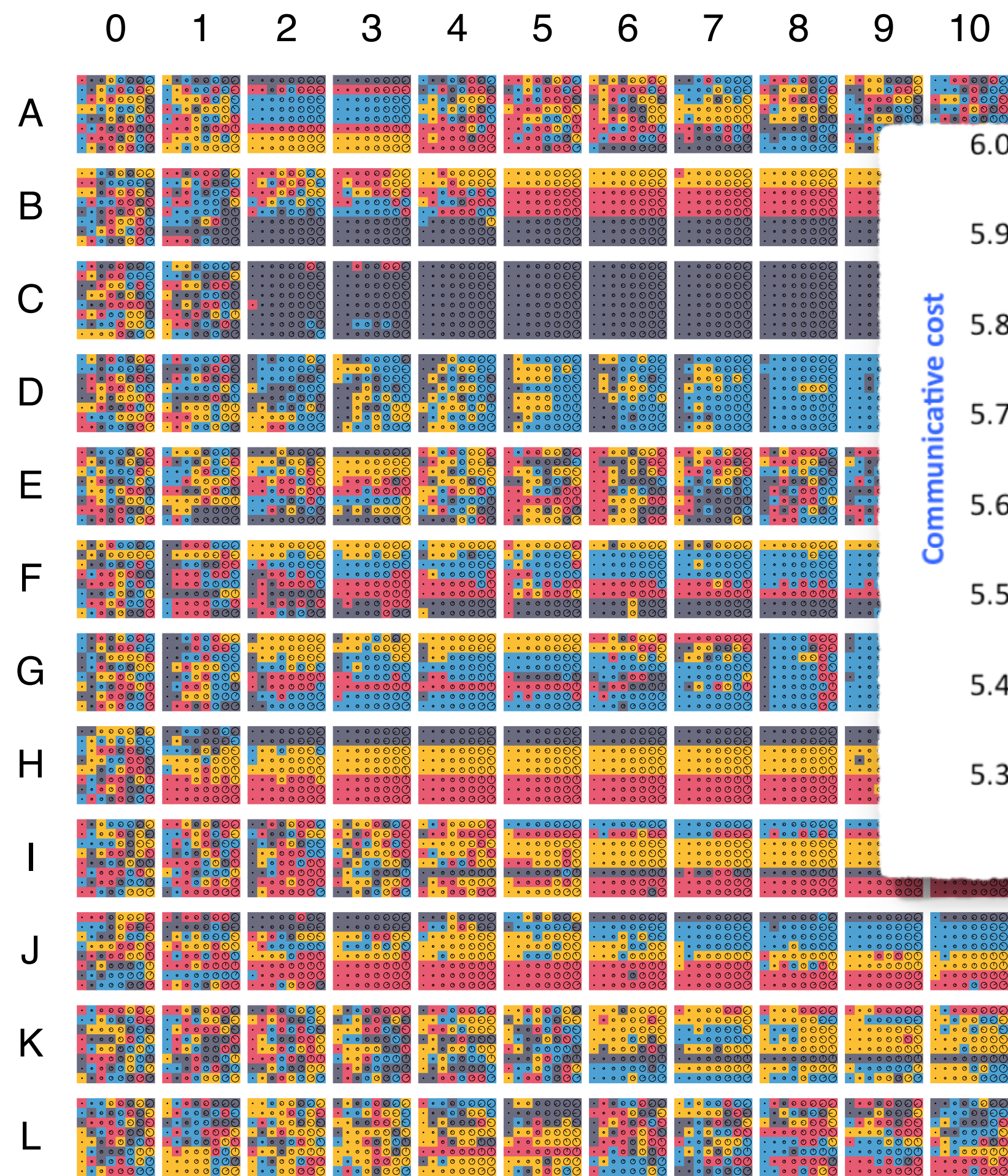


# Experimental results



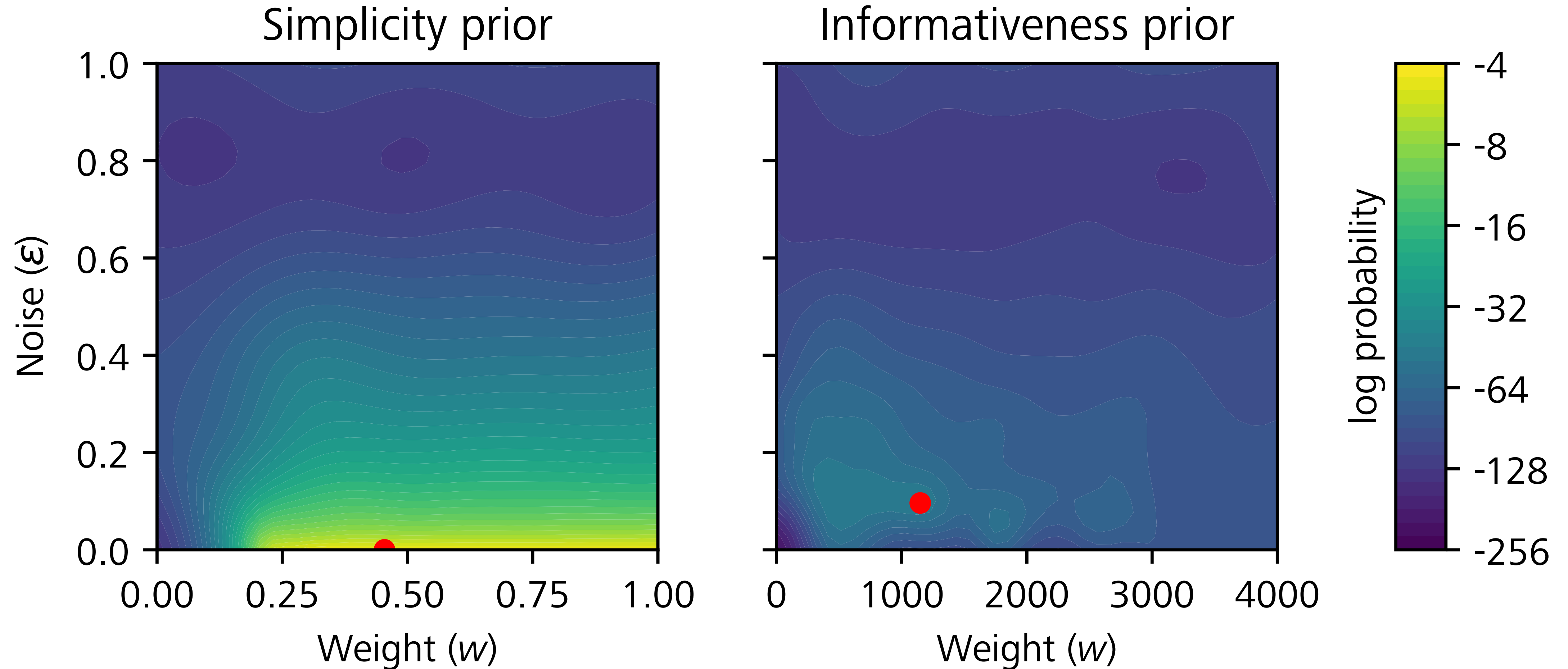


# Experimental results

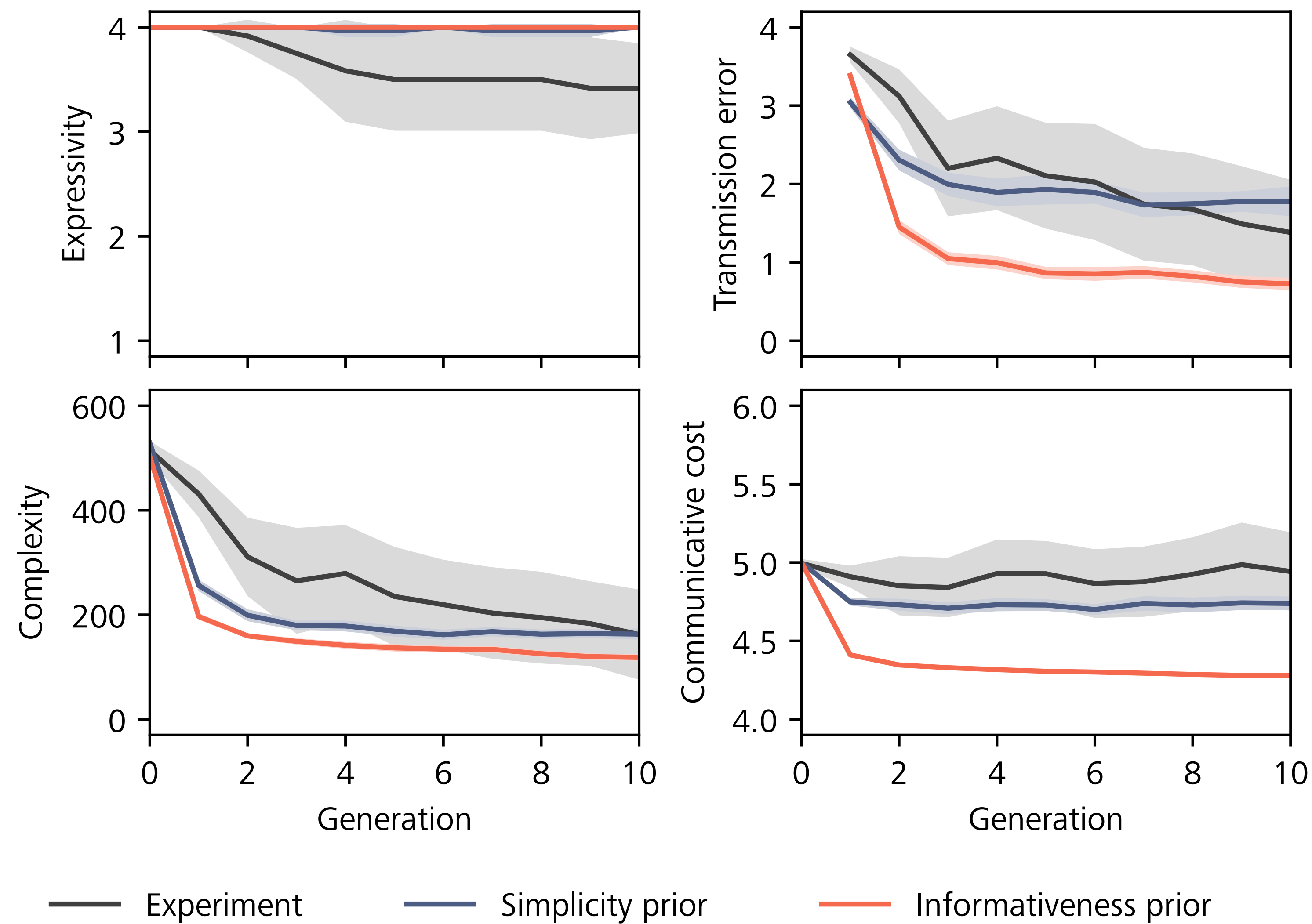


*Model fit*

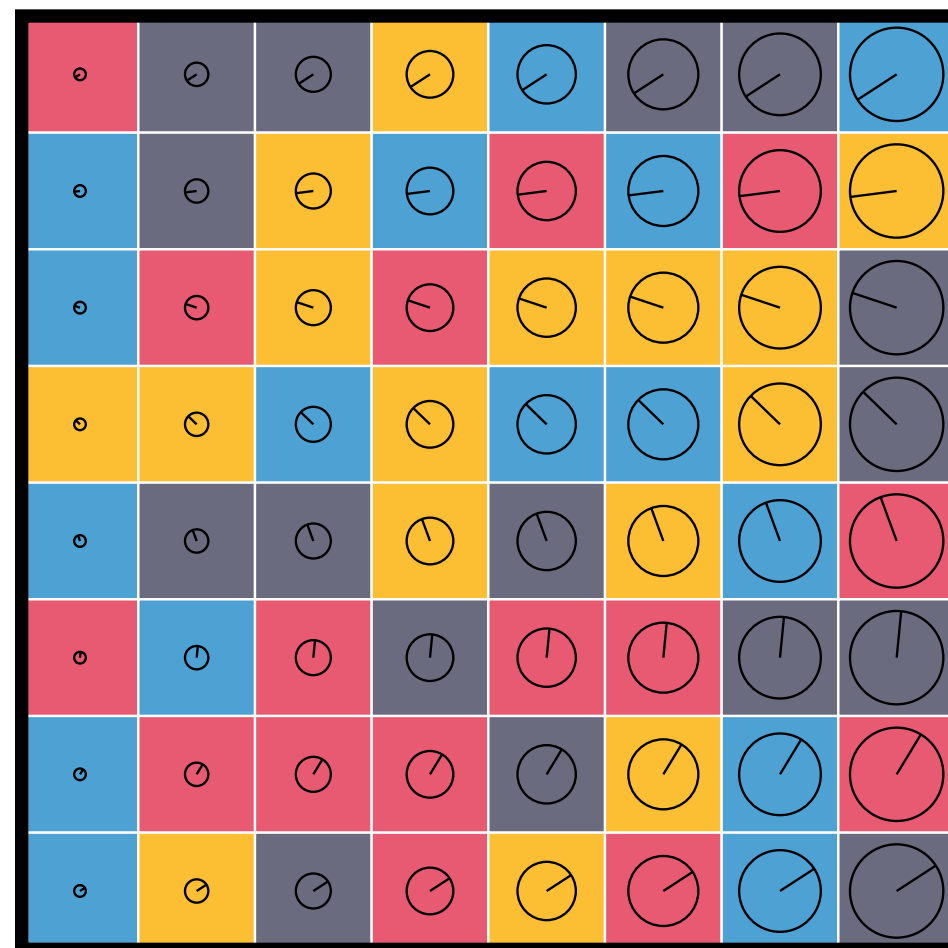
# Estimating unknown parameters of the model



# Rerun the model with parameters estimated from the experiment



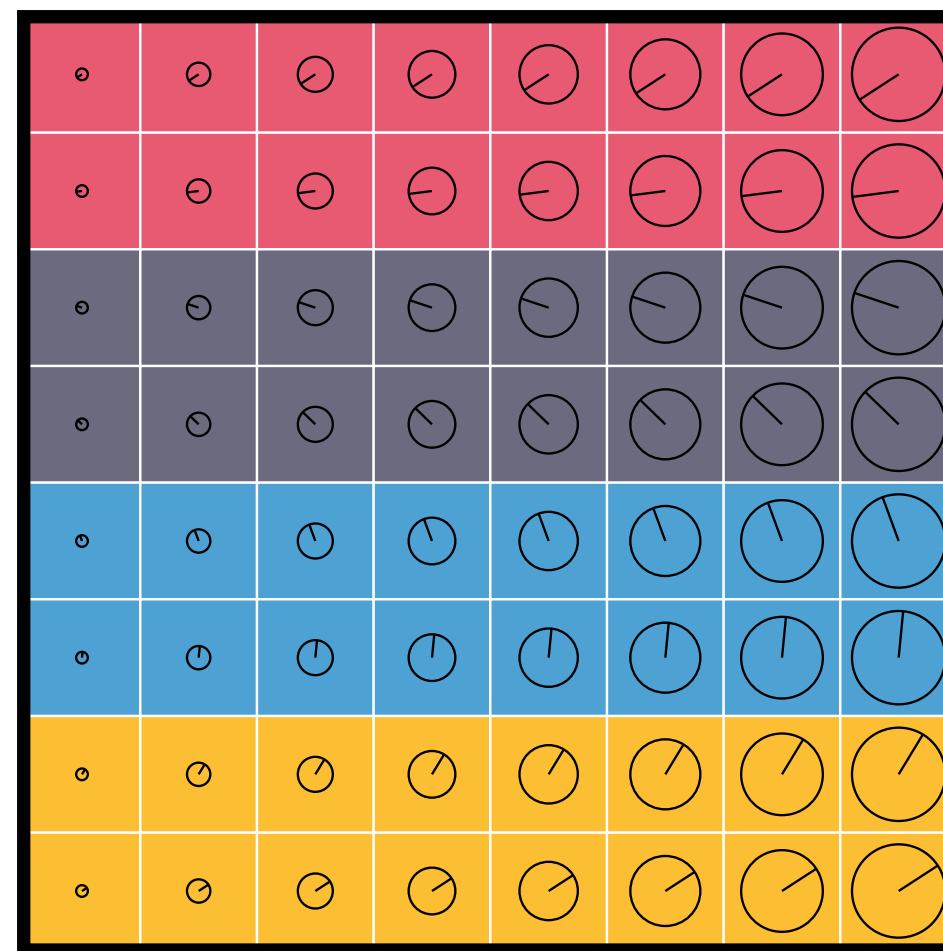
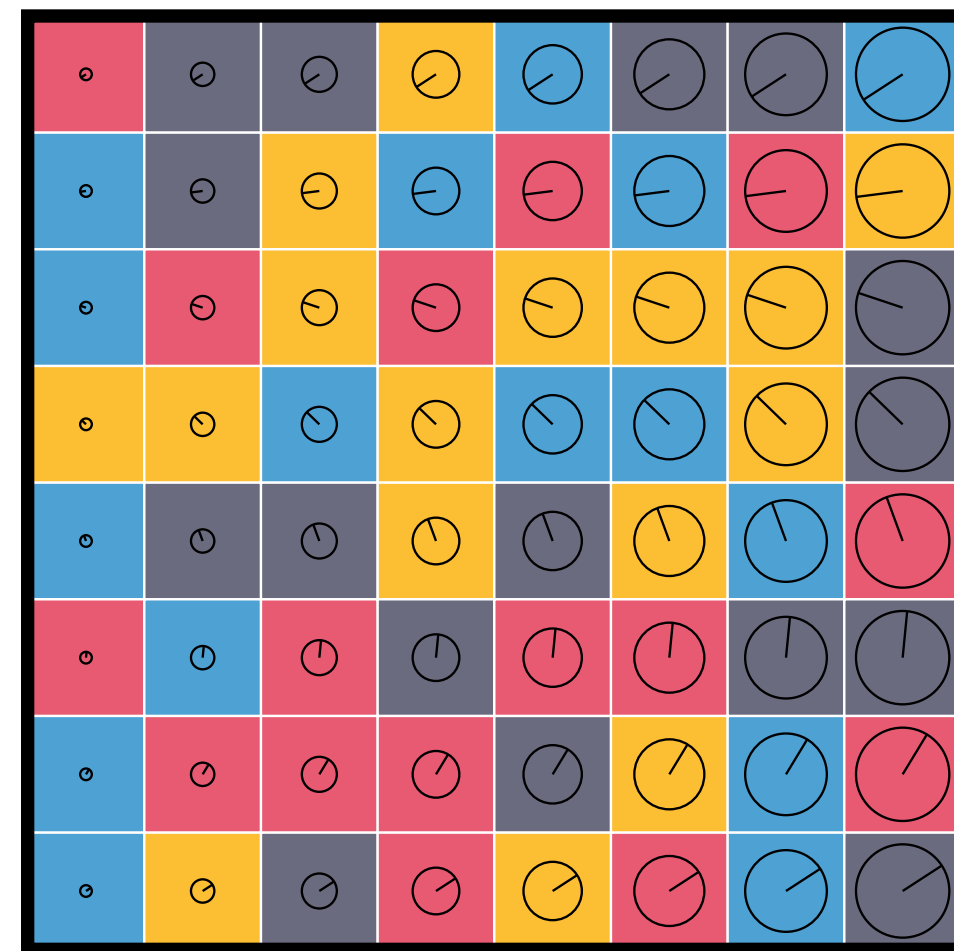
# Two ways of achieving simplicity





# Two ways of achieving simplicity

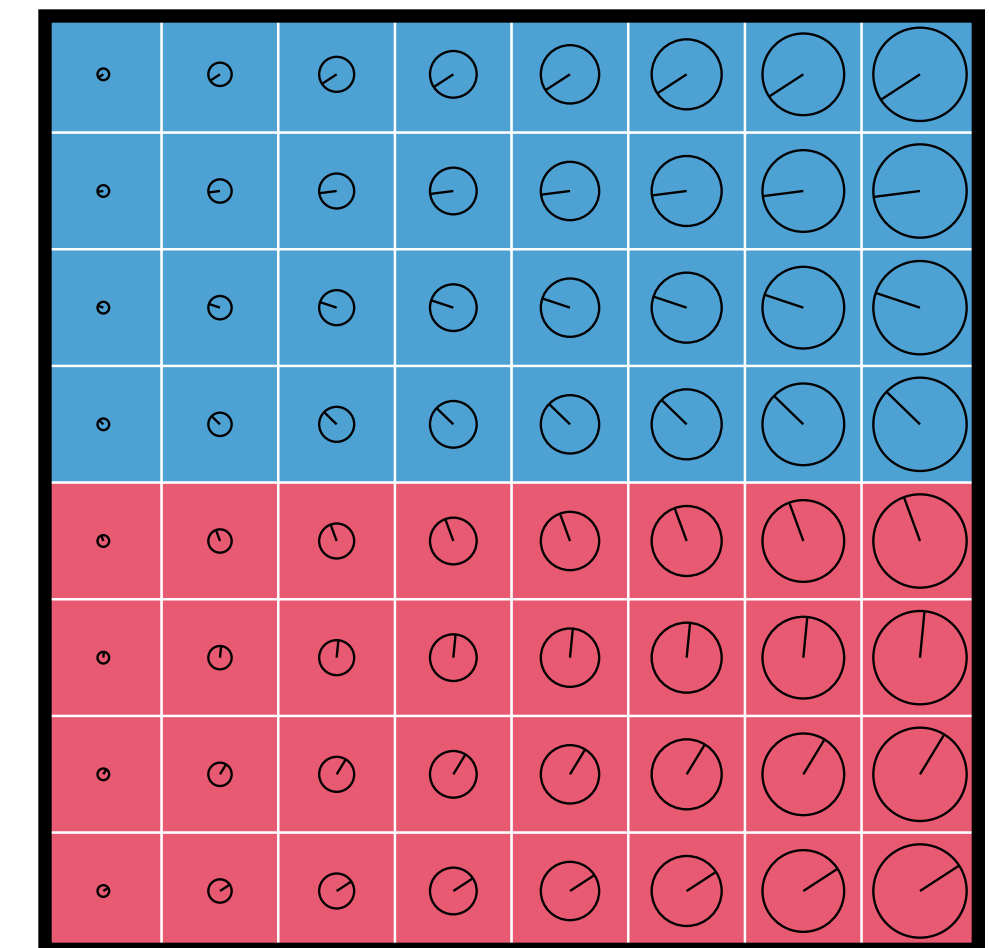
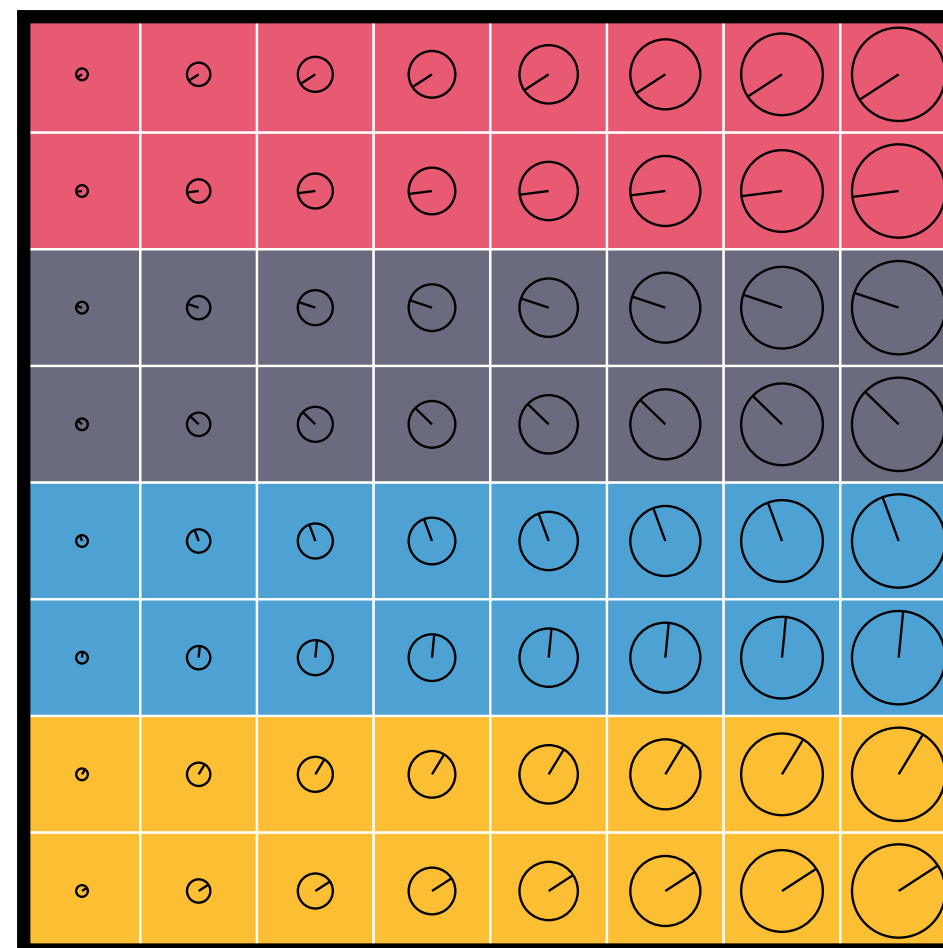
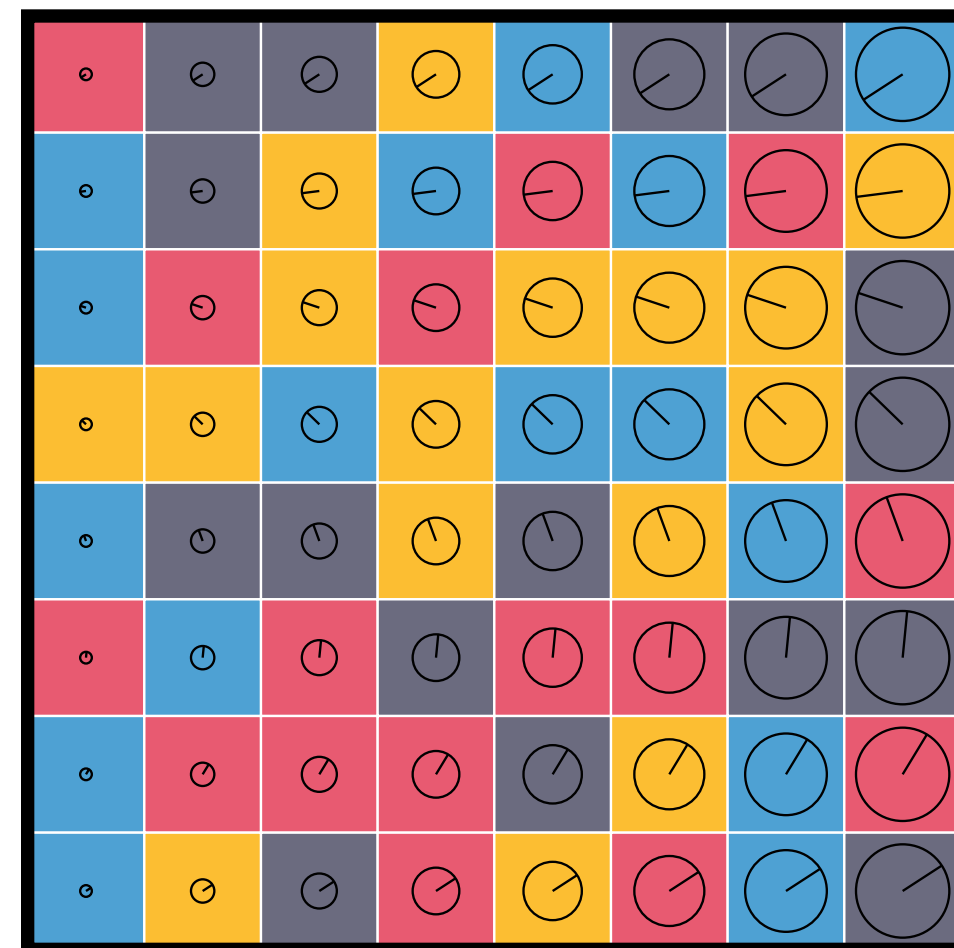
Increase in compactness





# Two ways of achieving simplicity

Increase in compactness

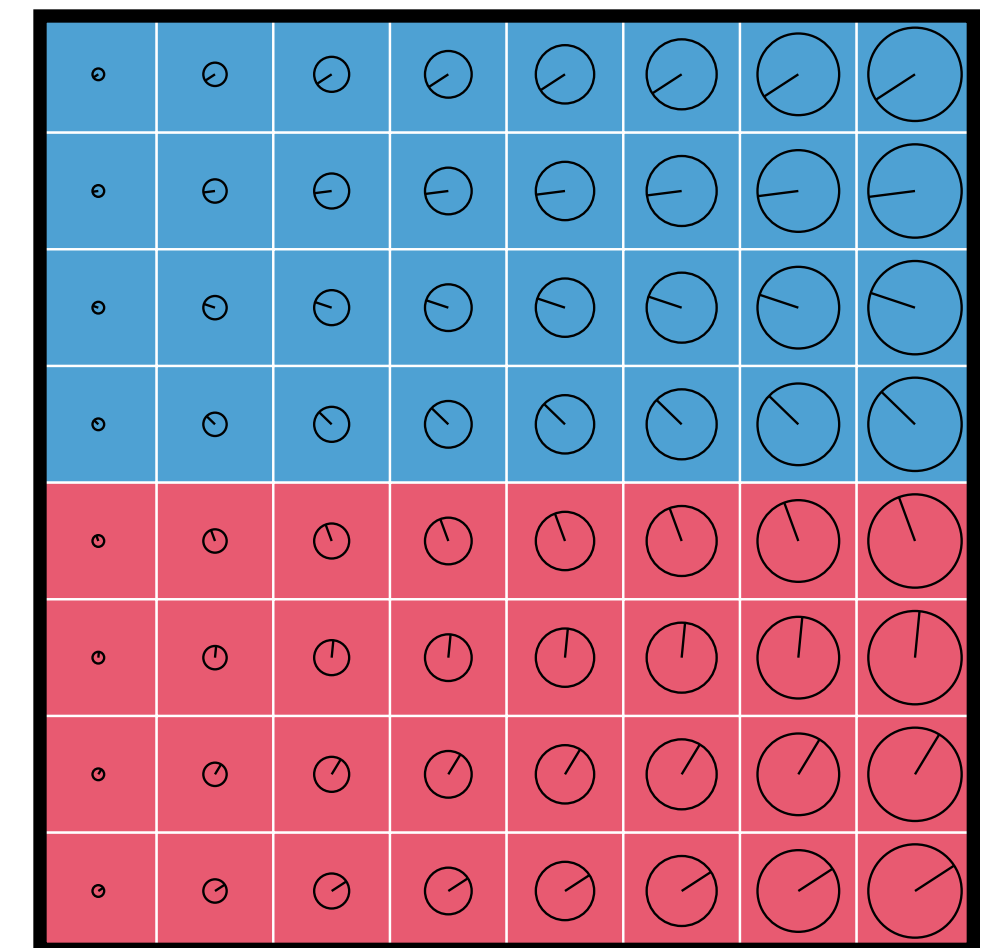
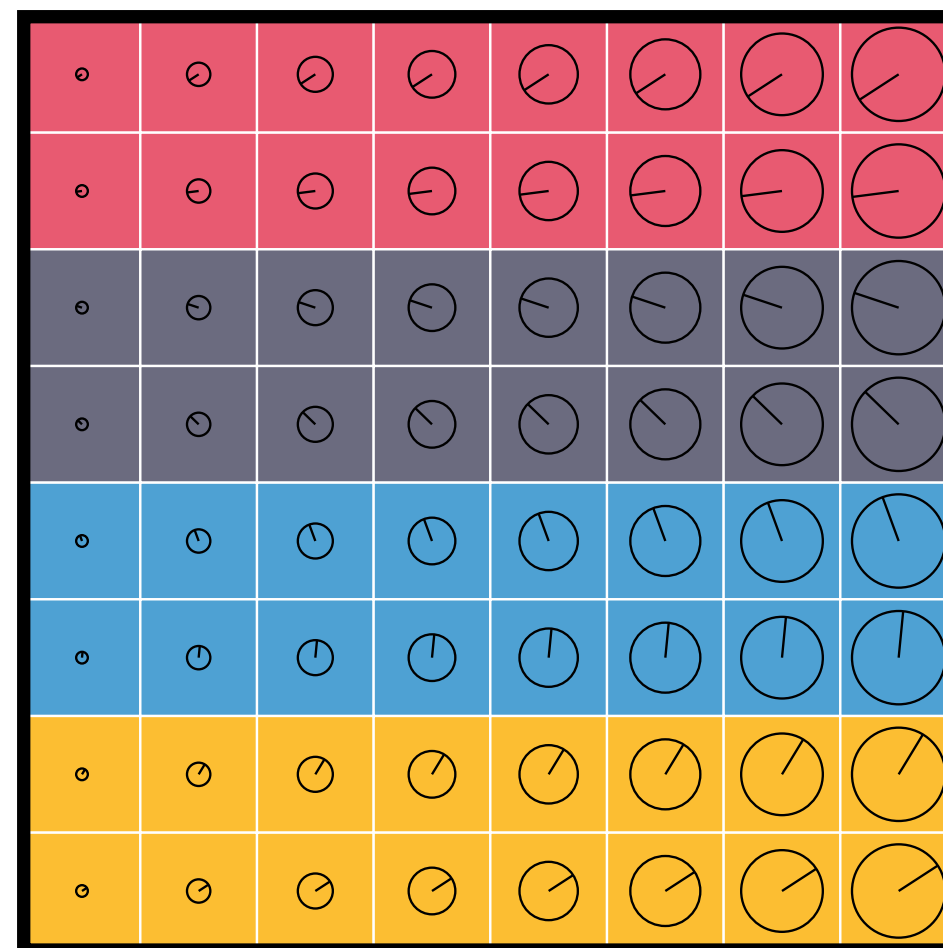
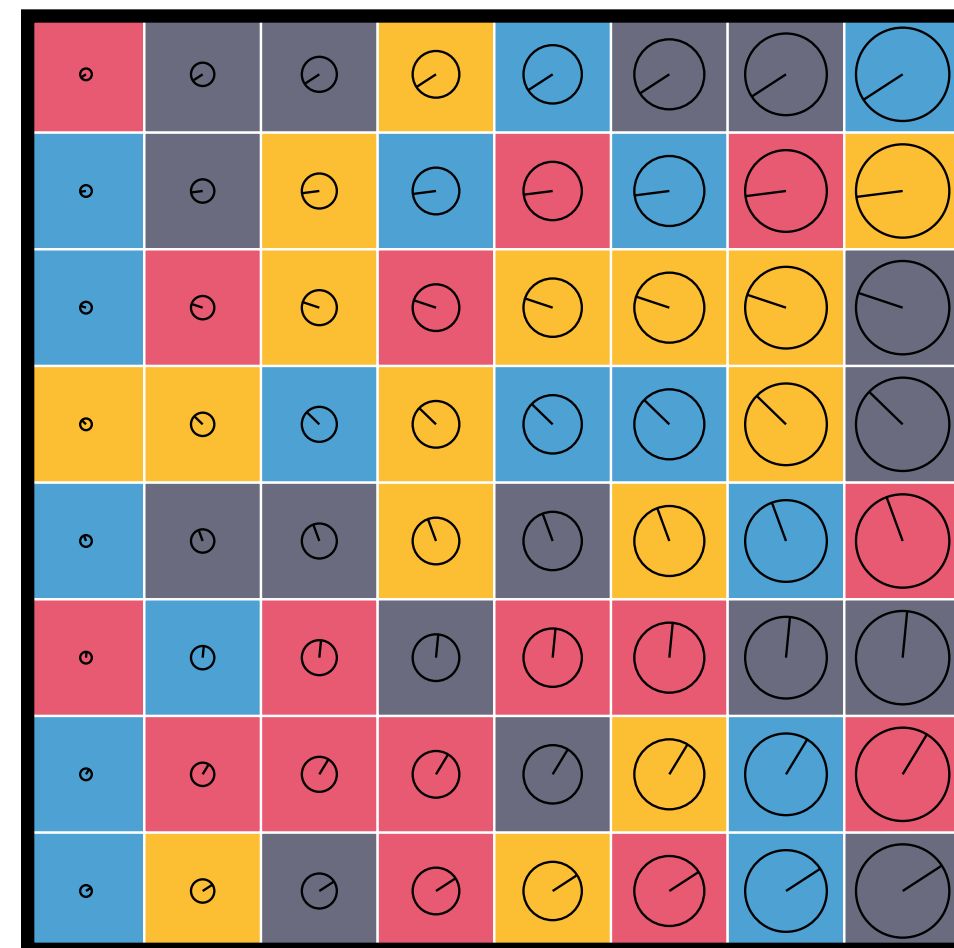


Decrease in expressivity

# Two ways of achieving simplicity

Increase in compactness

*increases informativeness*

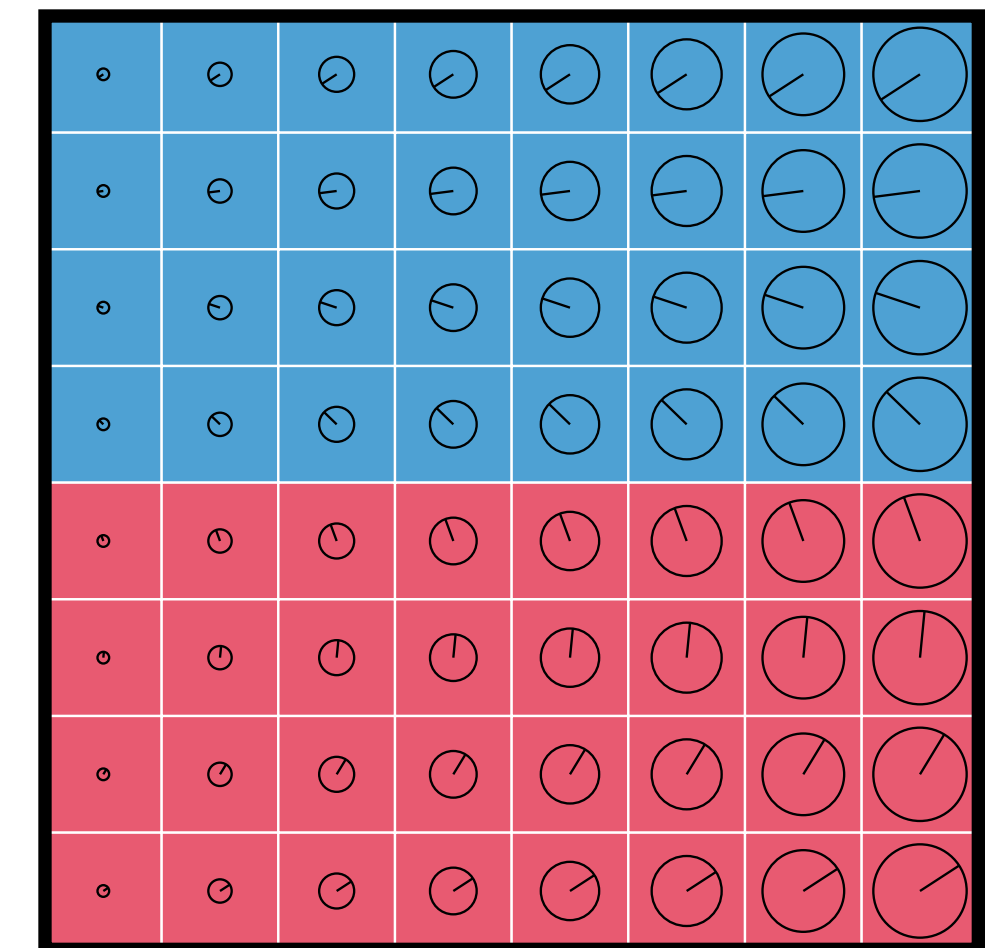
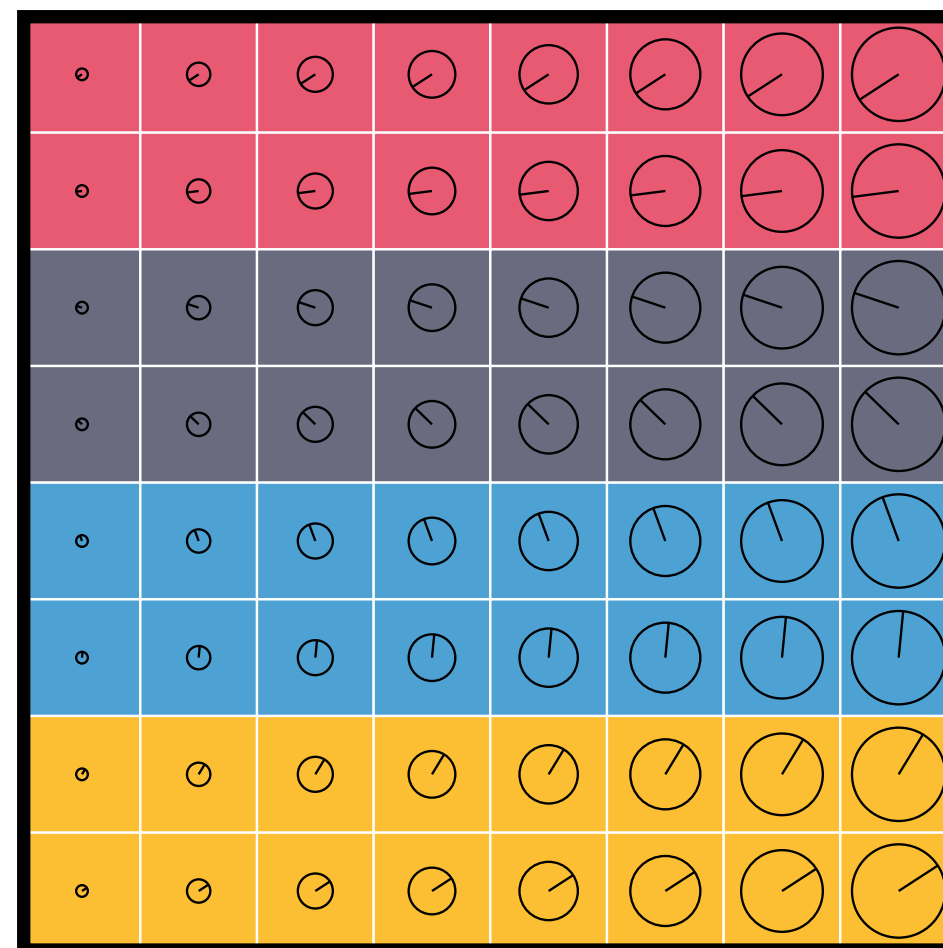
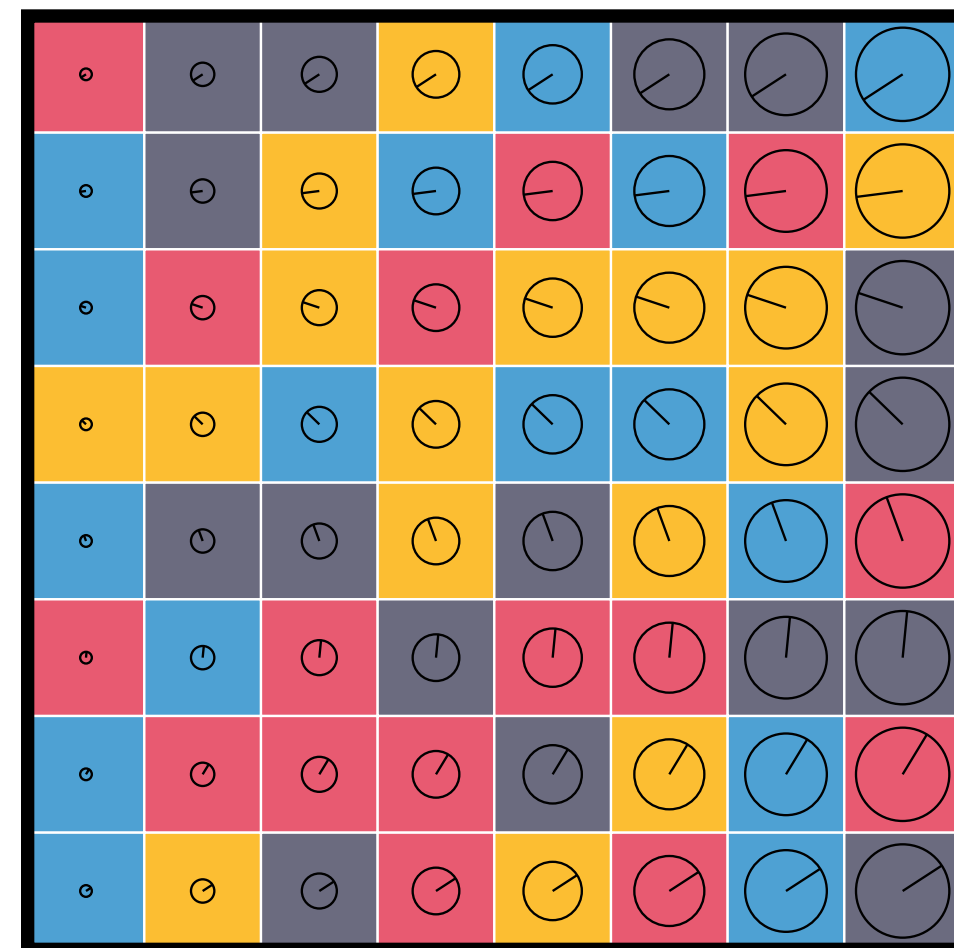


Decrease in expressivity

# Two ways of achieving simplicity

Increase in compactness

*increases informativeness*



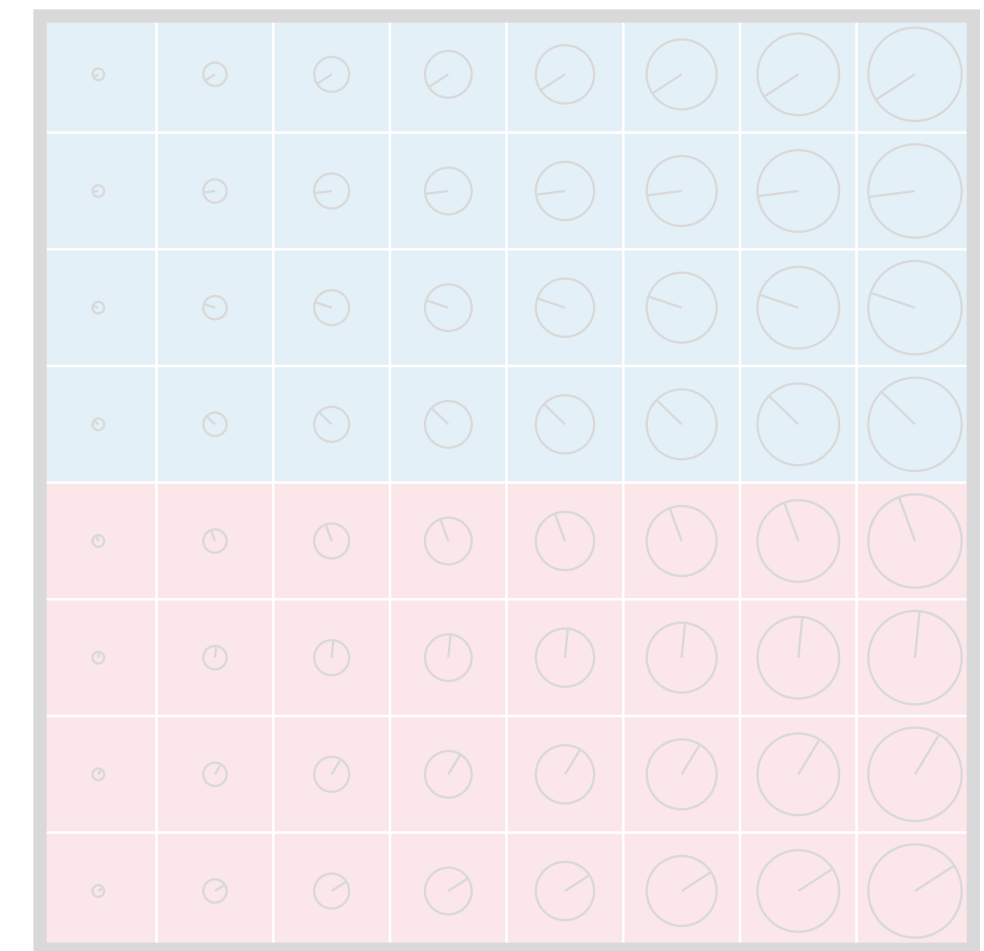
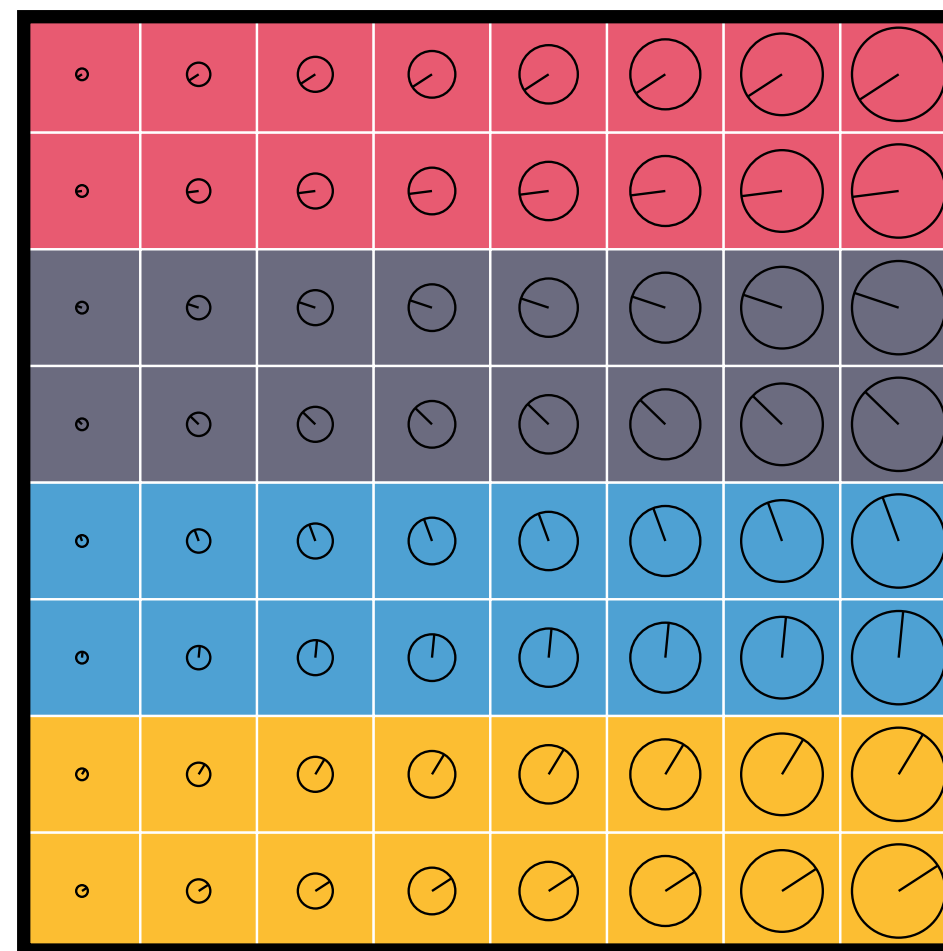
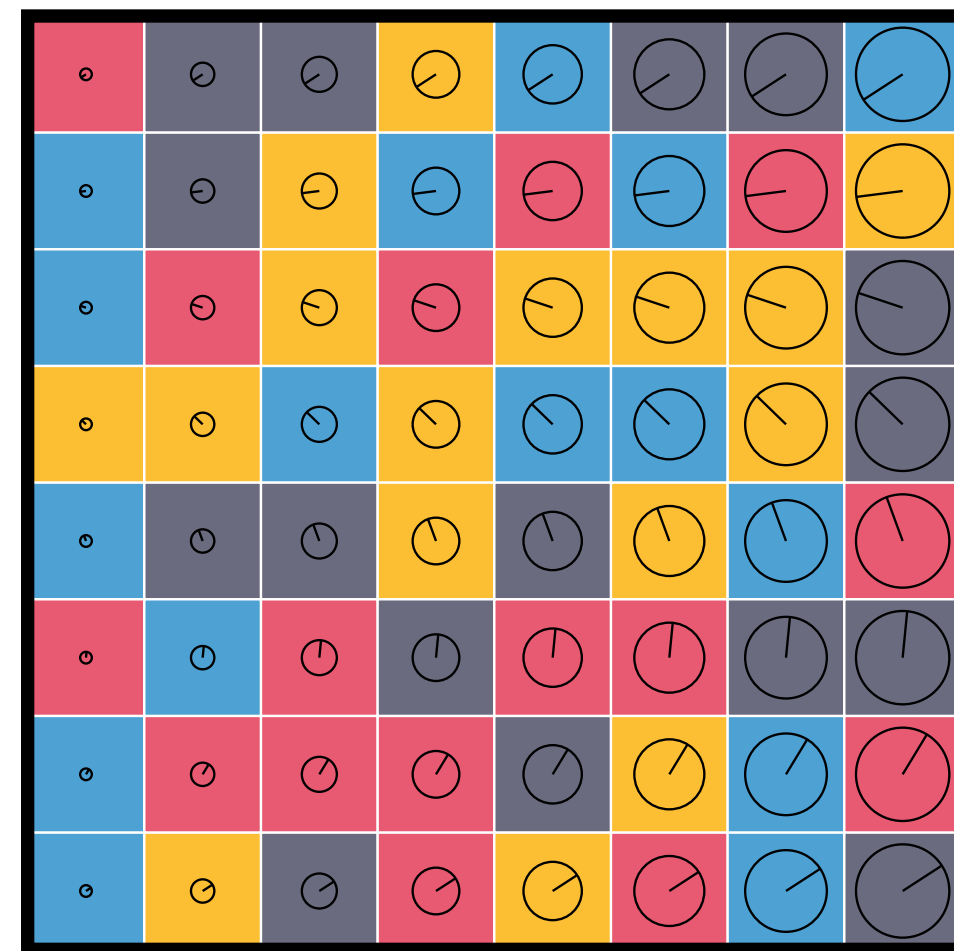
Decrease in expressivity

*decreases informativeness*

# Two ways of achieving simplicity

Increase in compactness

*increases informativeness*



Decrease in expressivity

*decreases informativeness*

# Conclusions

Languages are shaped in the simplicity–informativeness tradeoff by pressures from induction and interaction

For a rational learner, induction contains a simplicity bias to prevent overfitting noise, and to aid reasoning about unseen meanings

Iterated learning (repeated induction) converges to the prior bias, favouring languages that are as simple as possible:

**Loss of expressivity:** Loss of words/concepts to aid learning

**Compact categories:** Reorganization of the space to aid learning

In the process, some informativeness may come along for the ride, potentially obscuring the causal mechanism

Nevertheless, some kind of interactional dynamics (e.g. learning based on communicative success) must restrain languages from total degeneration

*Thanks!*