

# Simplicity priors and conceptual structure

Jon W. Carr, Kenny Smith, Jenny Culbertson, Simon Kirby

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

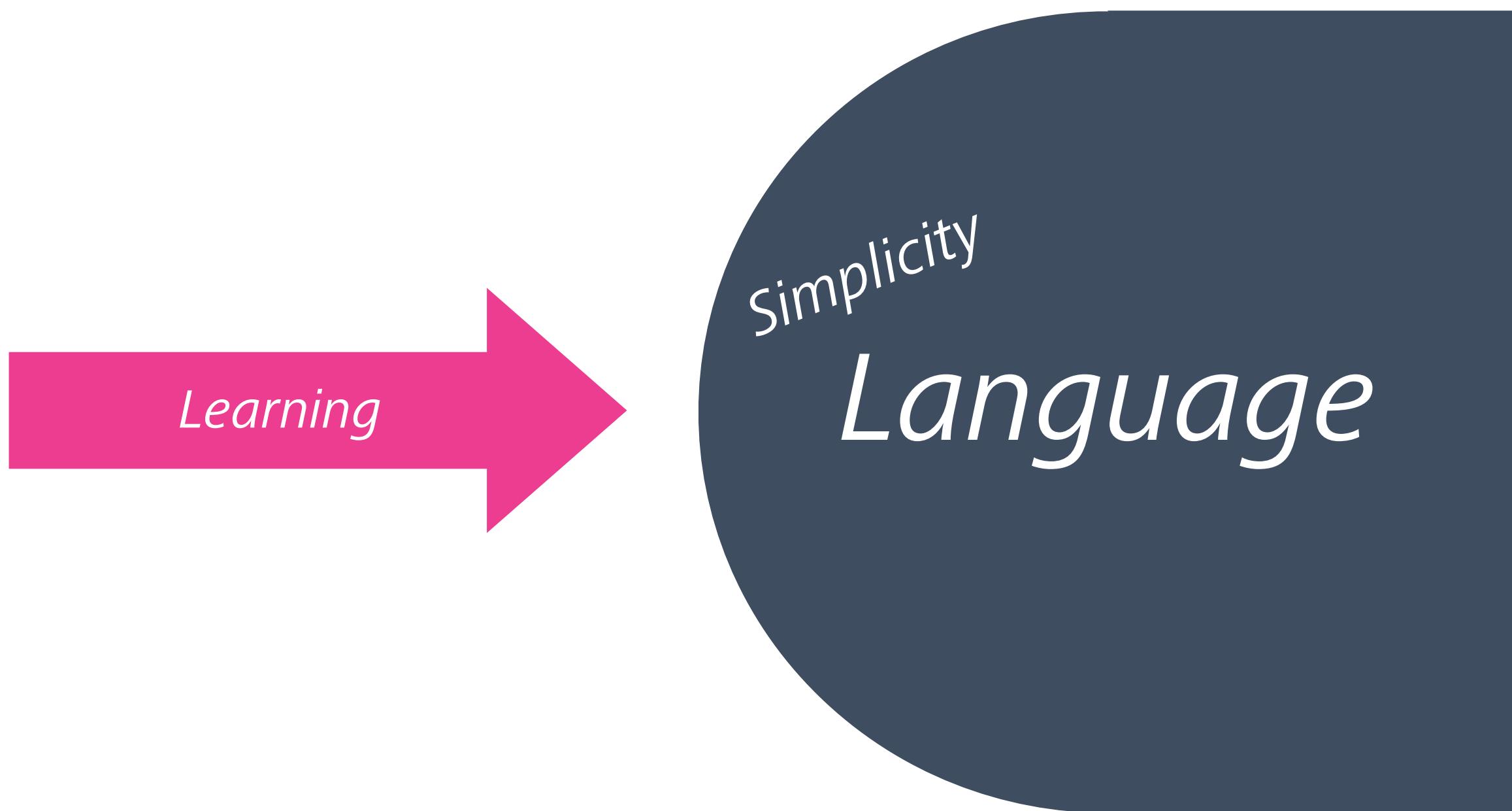


# Pressures shaping language



*Language*

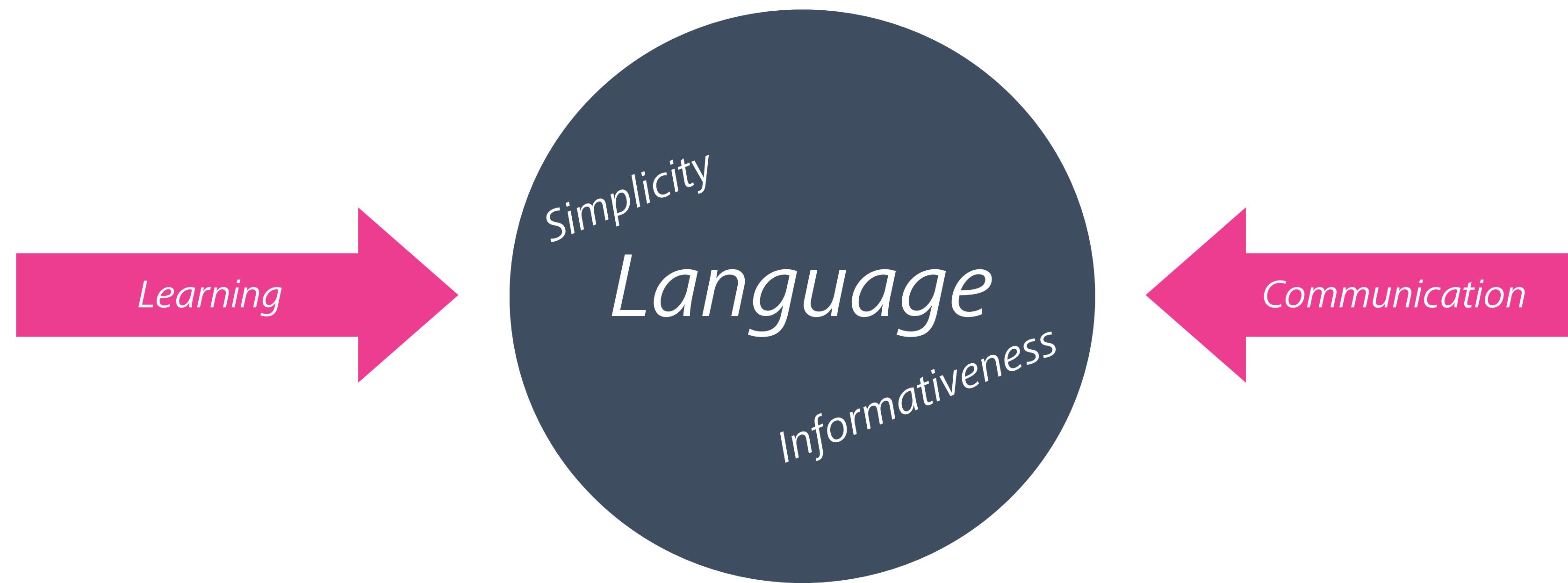
# Pressures shaping language



# Pressures shaping language

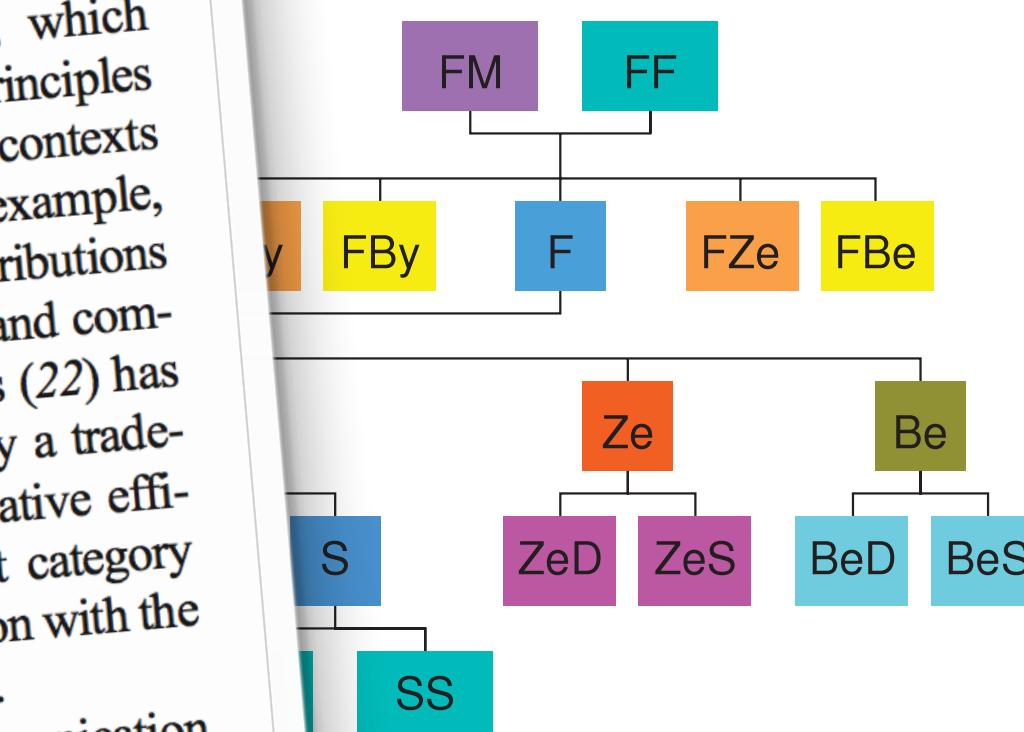
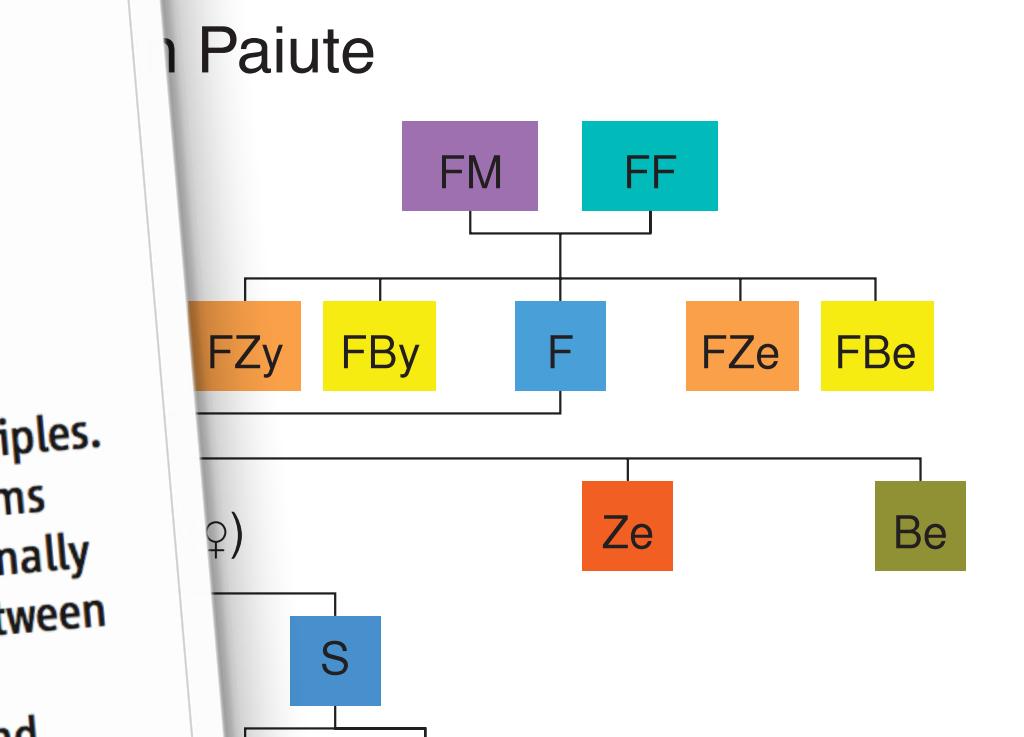
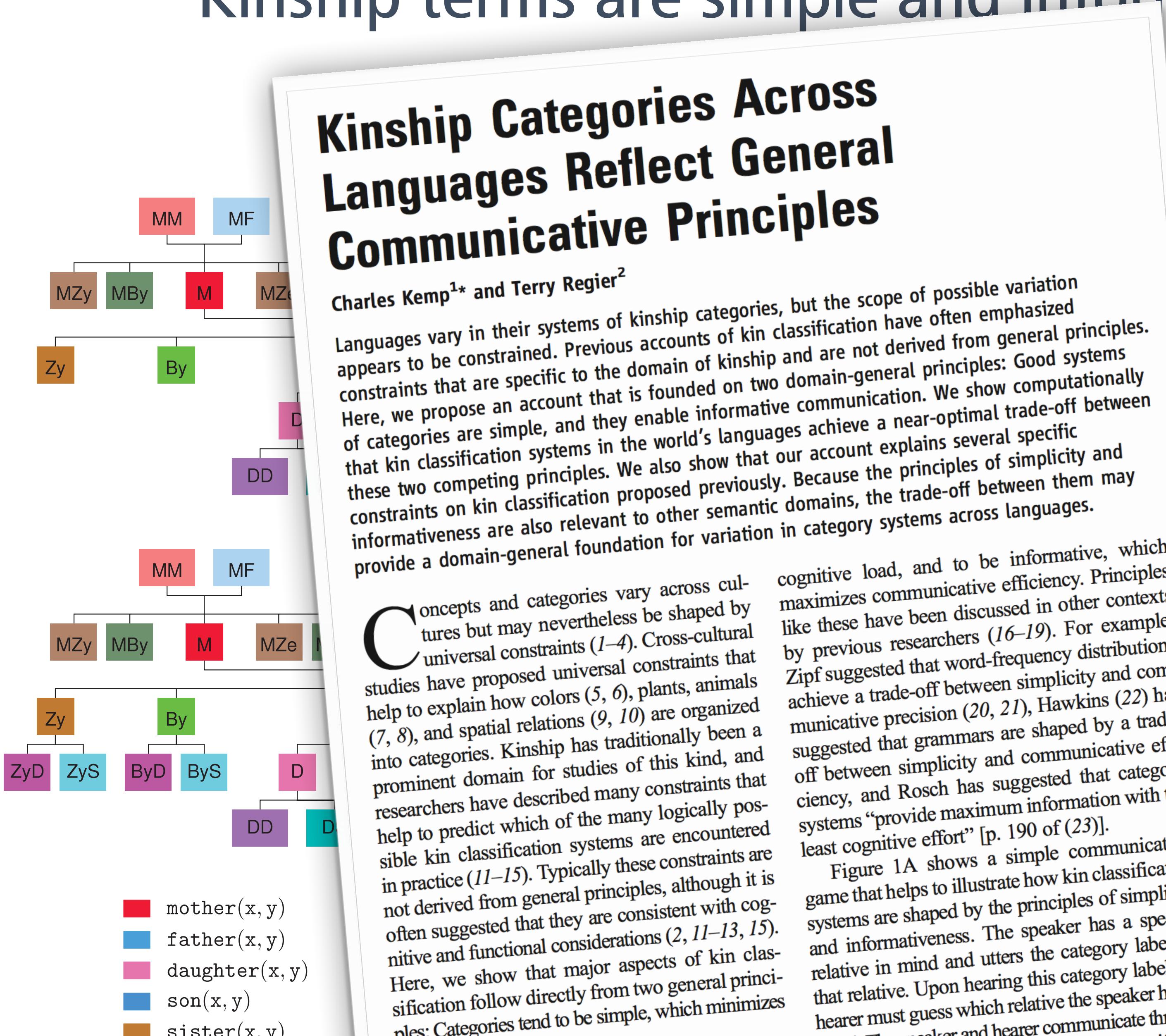


# Pressures shaping language



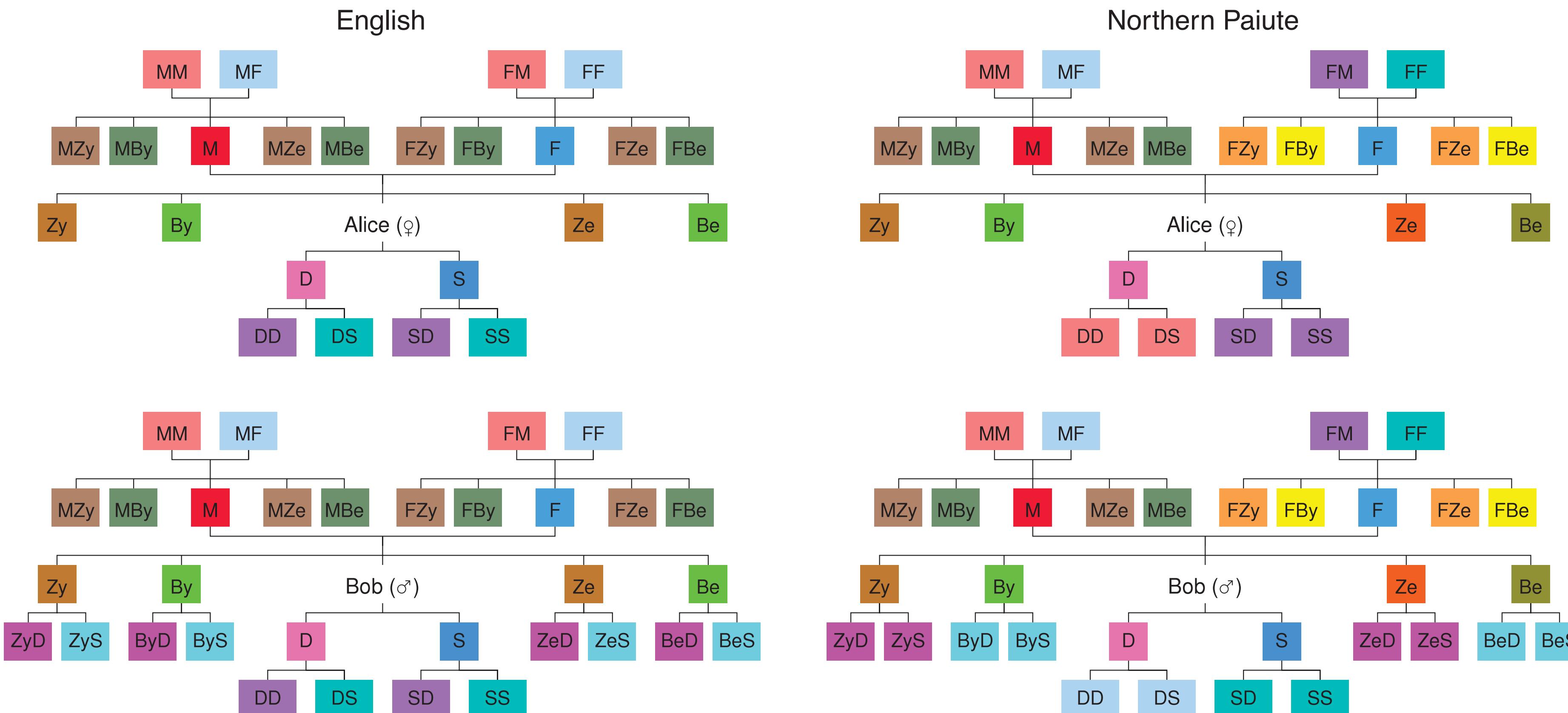
*The simplicity–informativeness tradeoff*

# Kinship terms are simple and informative



NT(x, y)  $\wedge$  FEMALE(x)  
 NT(x, y)  $\wedge$  MALE(x)  
 D(x, y)  $\wedge$  FEMALE(x)  
 (x, y)  $\wedge$  MALE(x)  
 daughter(x, z)  $\wedge$  PARENT(z, v)

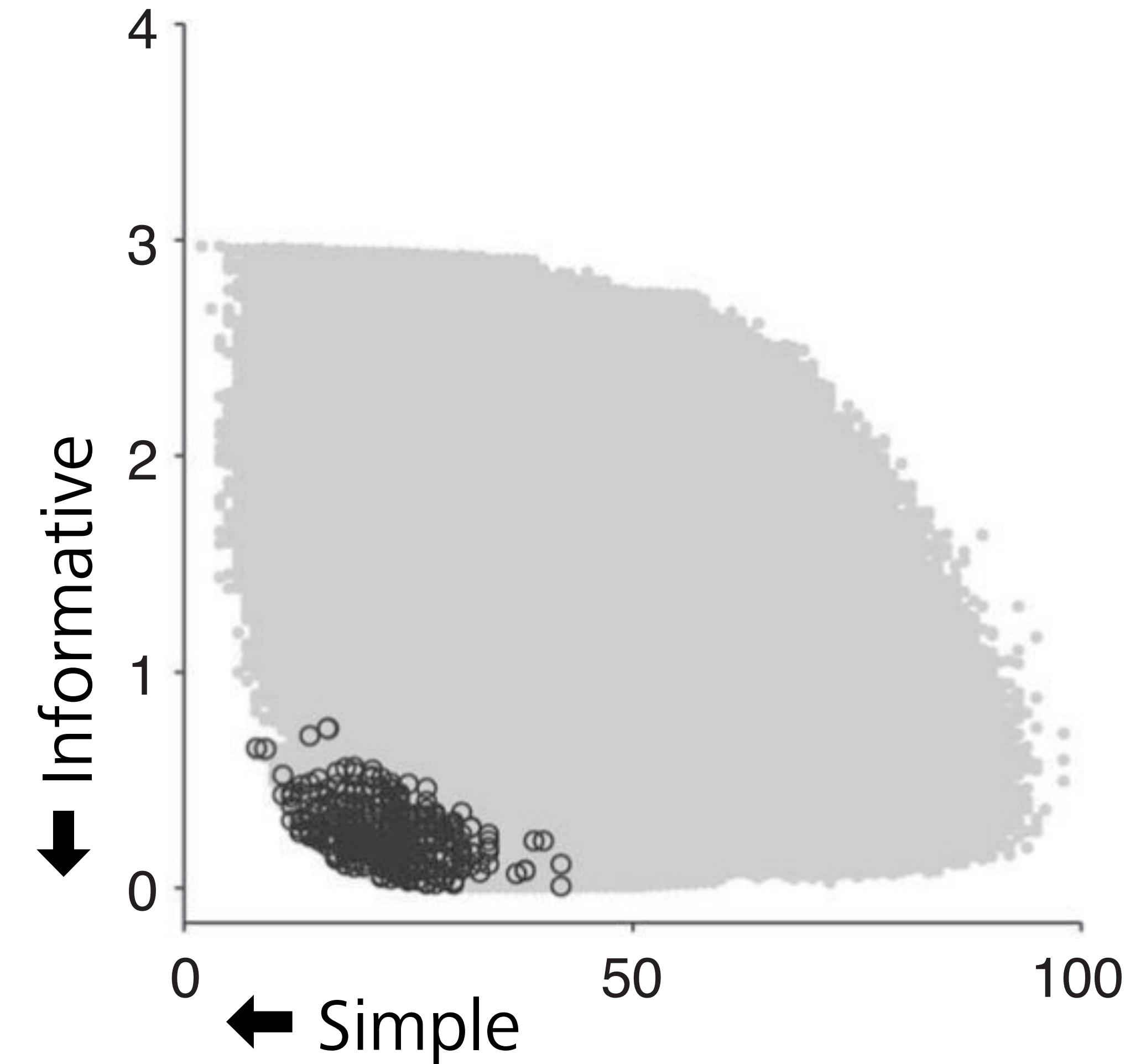
# Kinship terms are simple and informative



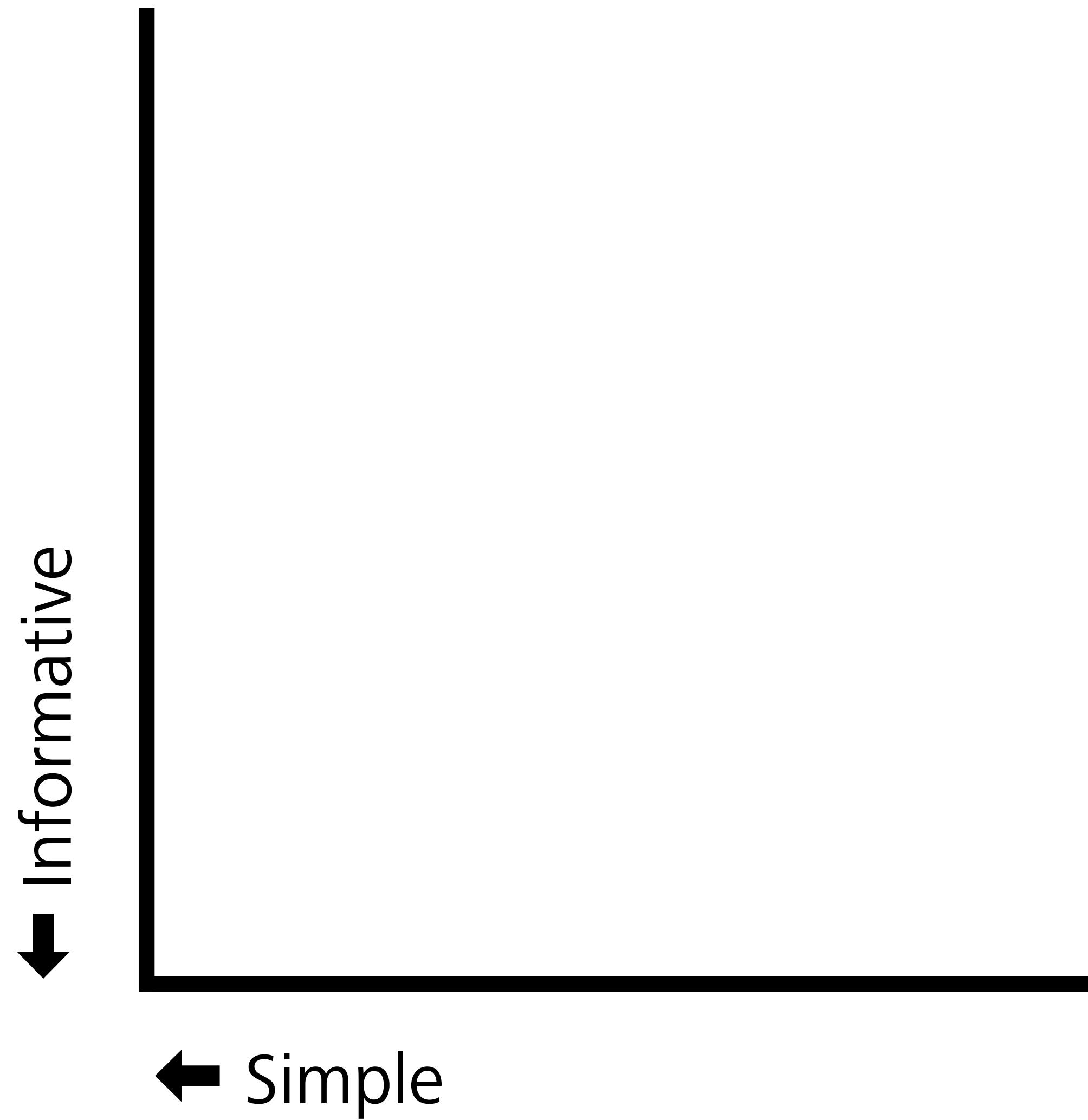
■  $\text{mother}(x, y) \leftrightarrow \text{PARENT}(x, y) \wedge \text{FEMALE}(x)$   
■  $\text{father}(x, y) \leftrightarrow \text{PARENT}(x, y) \wedge \text{MALE}(x)$   
■  $\text{daughter}(x, y) \leftrightarrow \text{CHILD}(x, y) \wedge \text{FEMALE}(x)$   
■  $\text{son}(x, y) \leftrightarrow \text{CHILD}(x,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

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

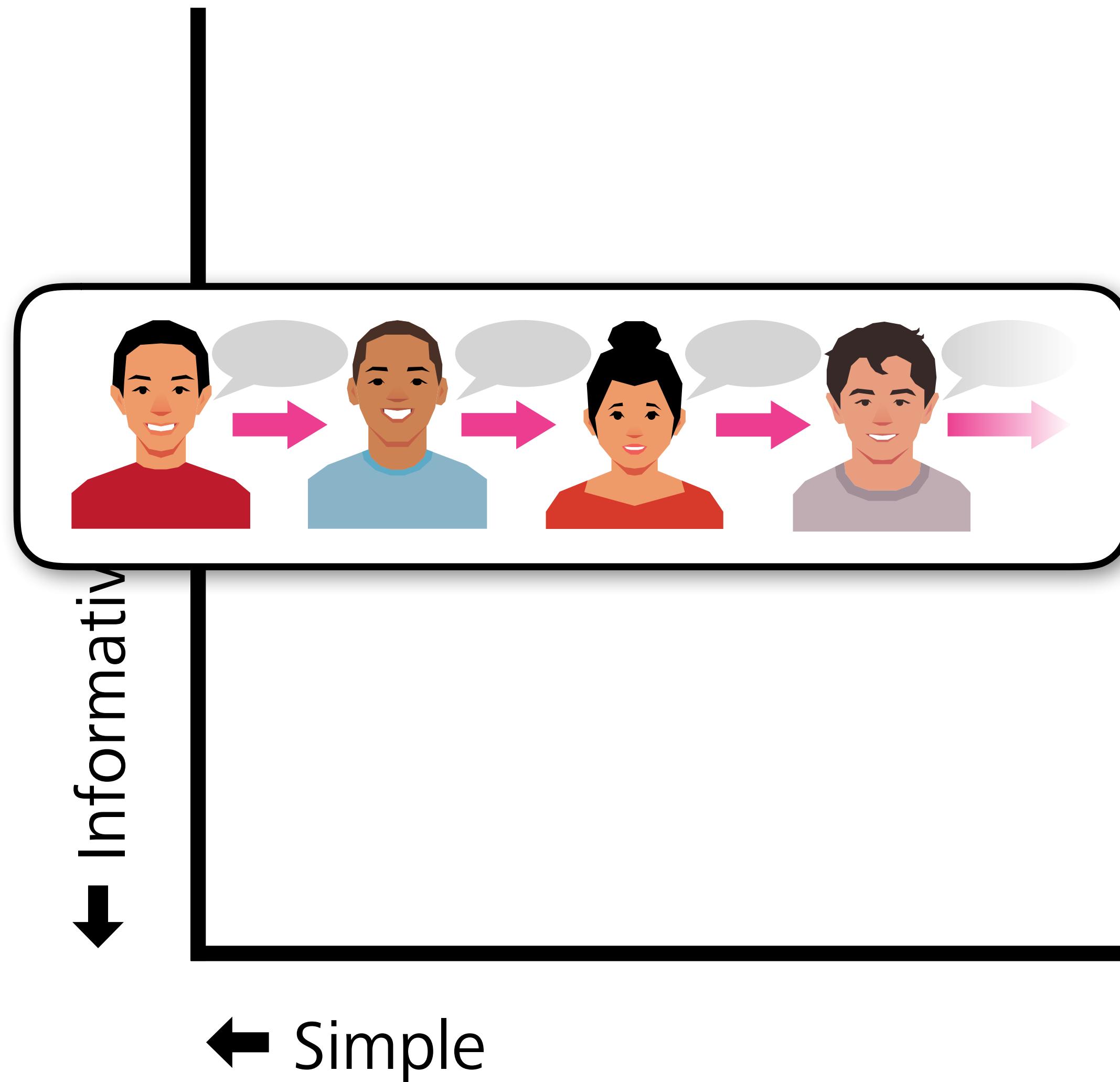
# Kinship terms are simple and informative



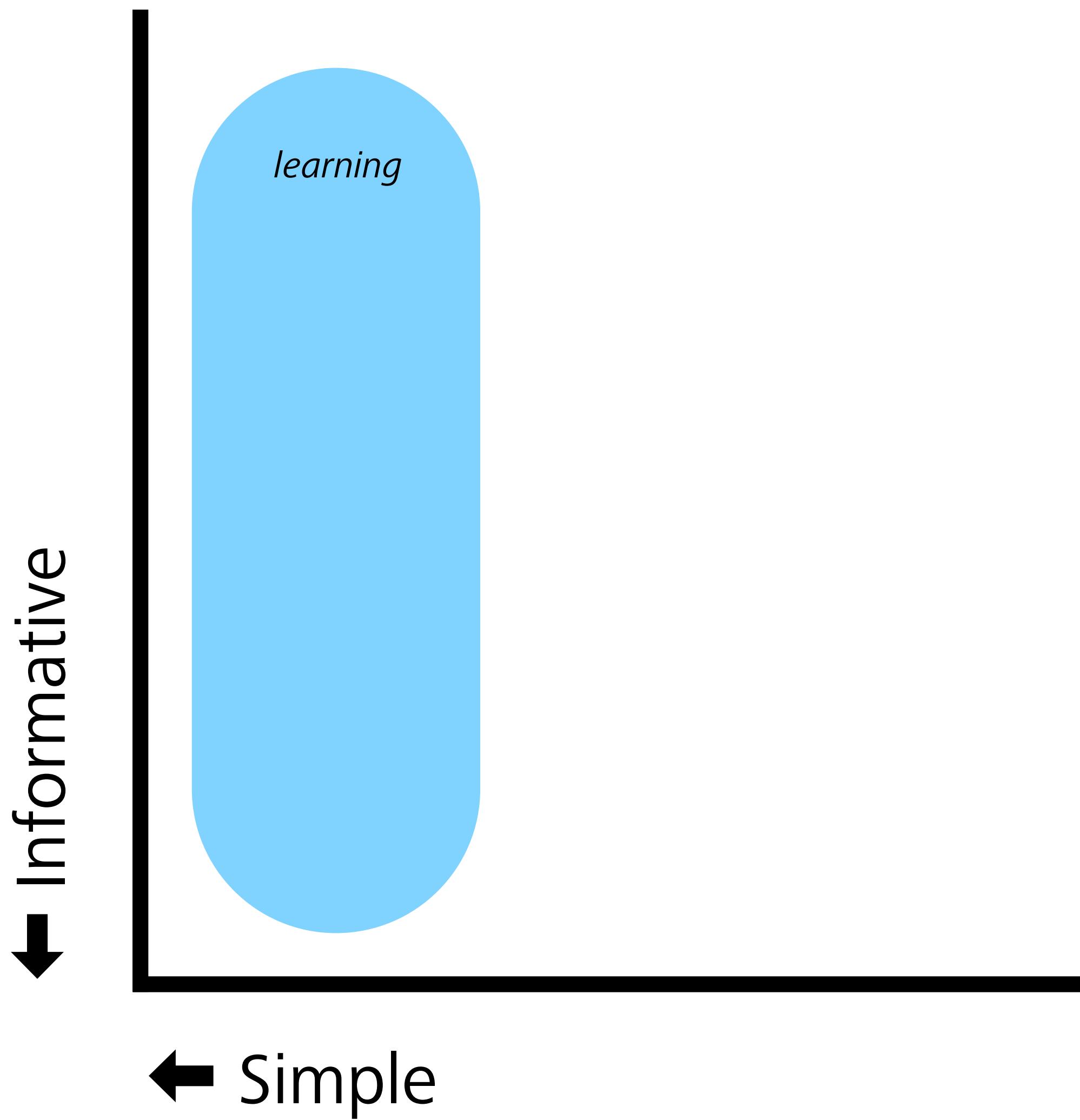
# Learning and communication pressures



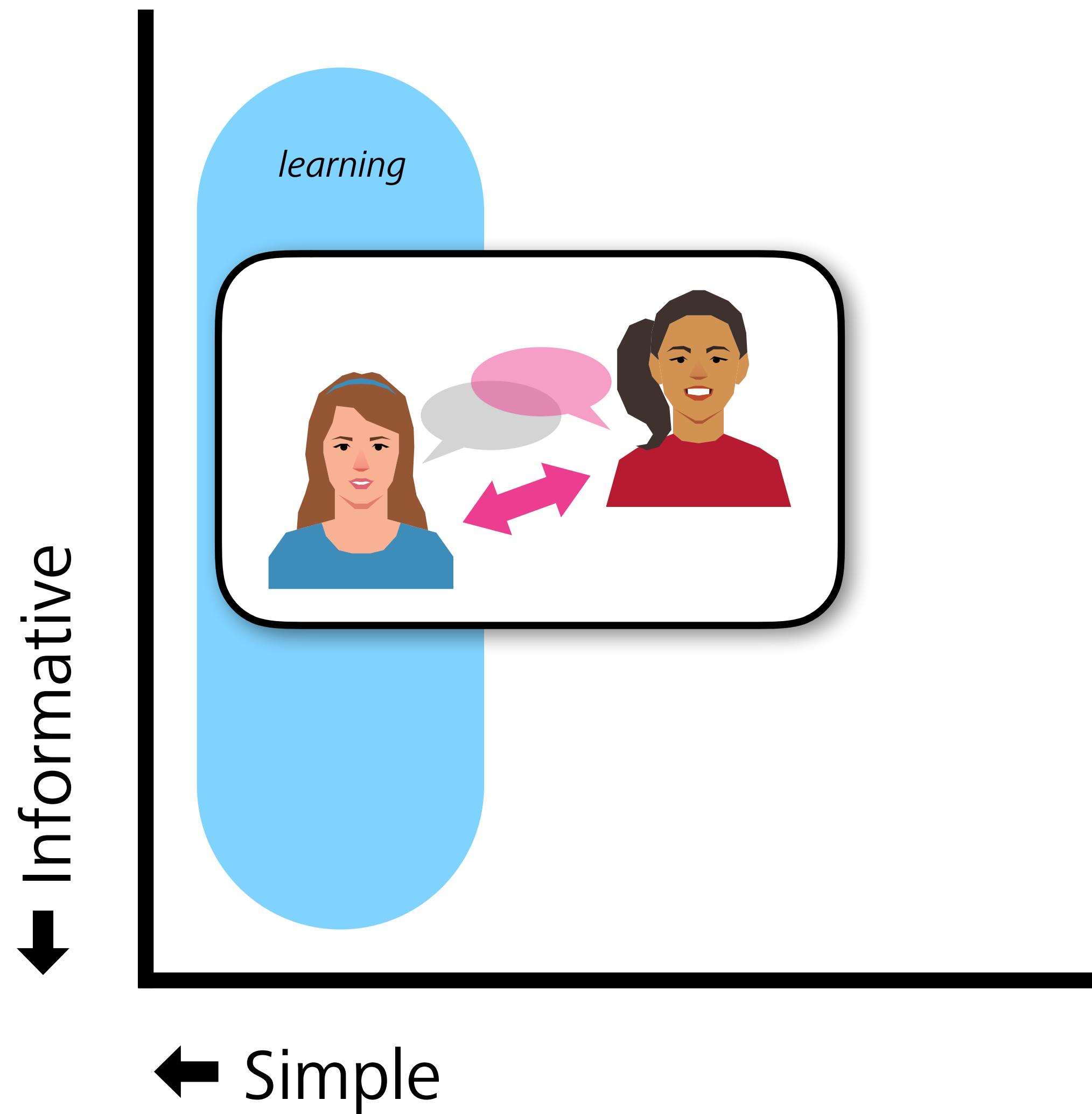
# Learning and communication pressures



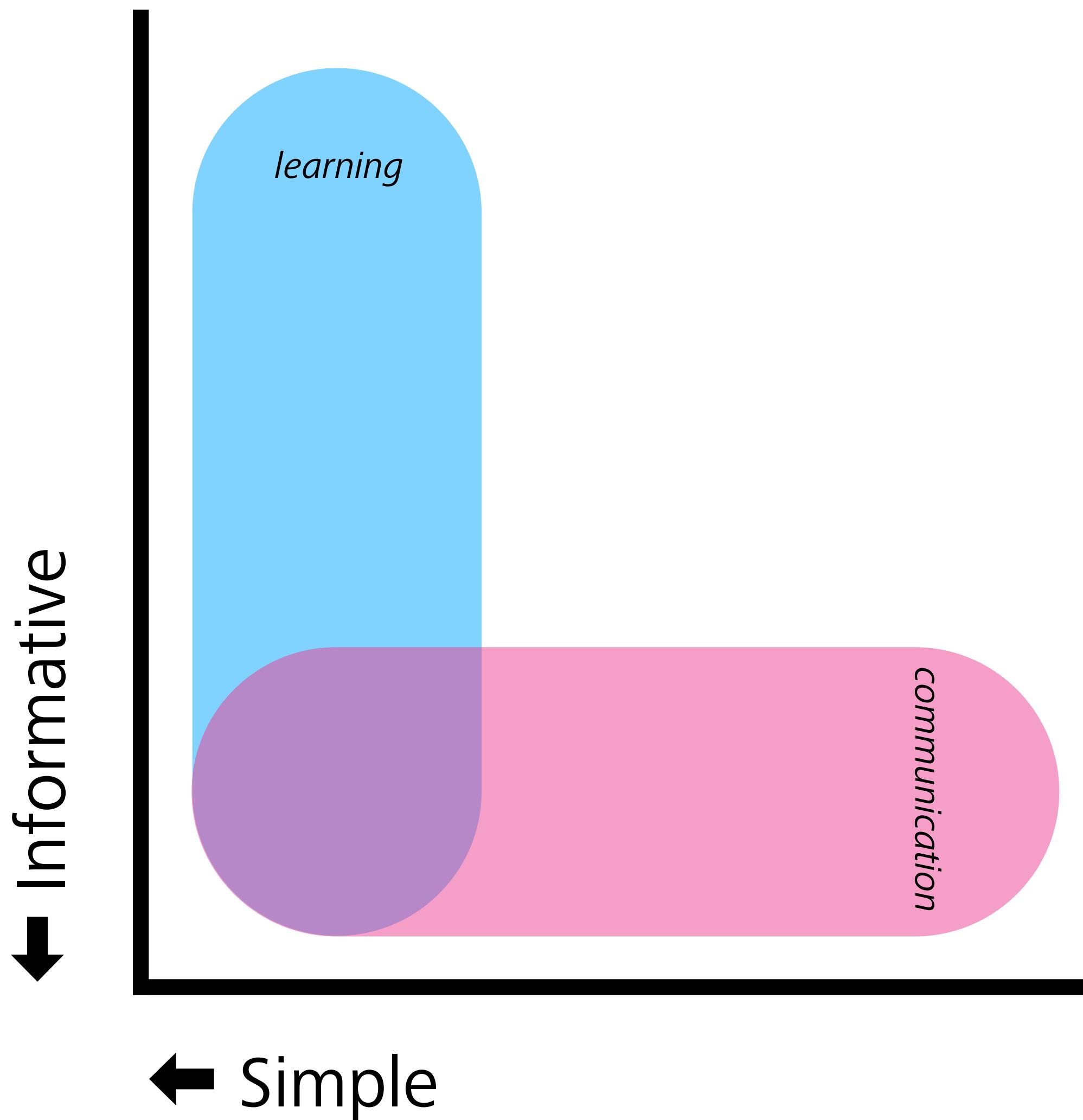
# Learning and communication pressures



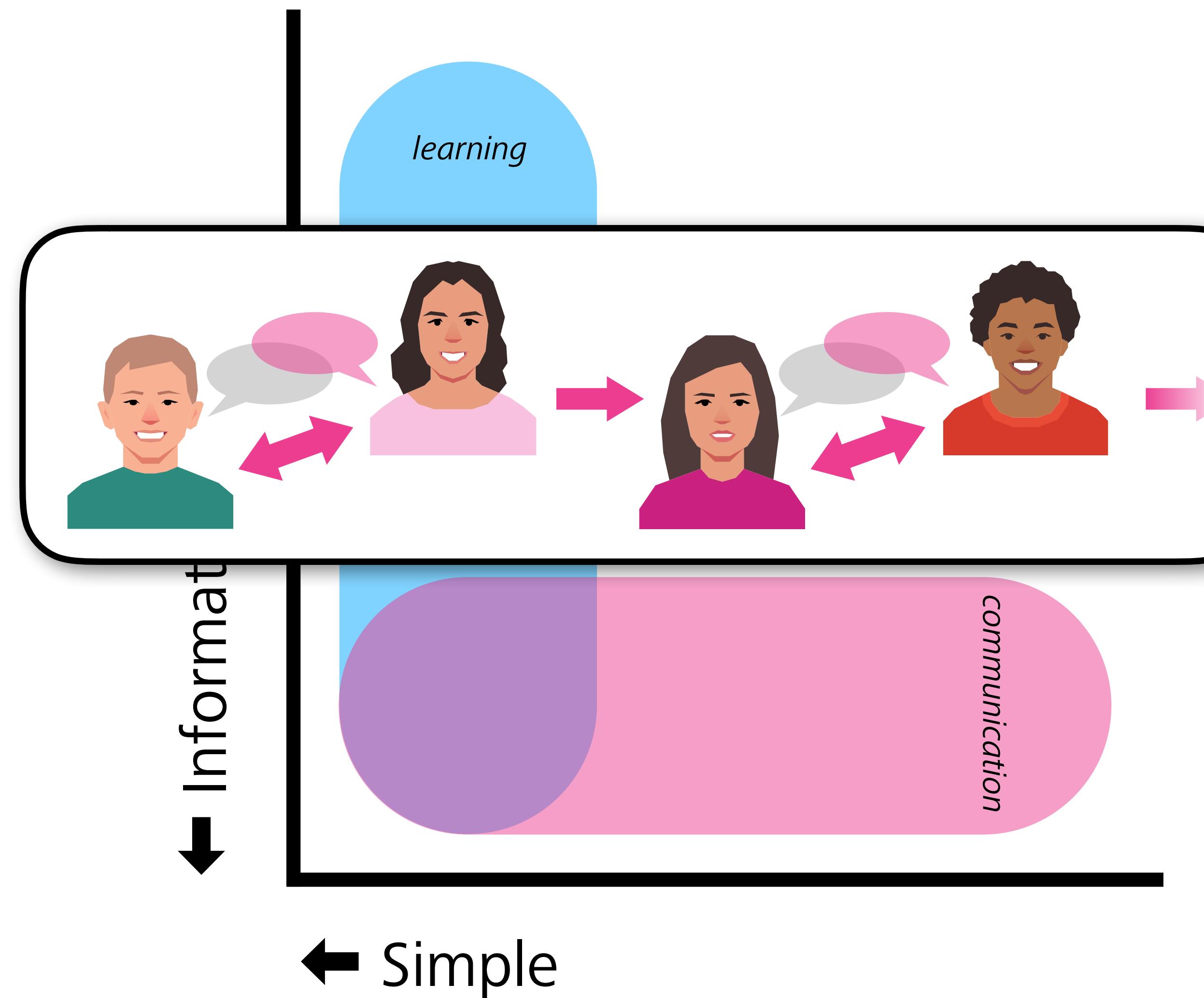
# Learning and communication pressures



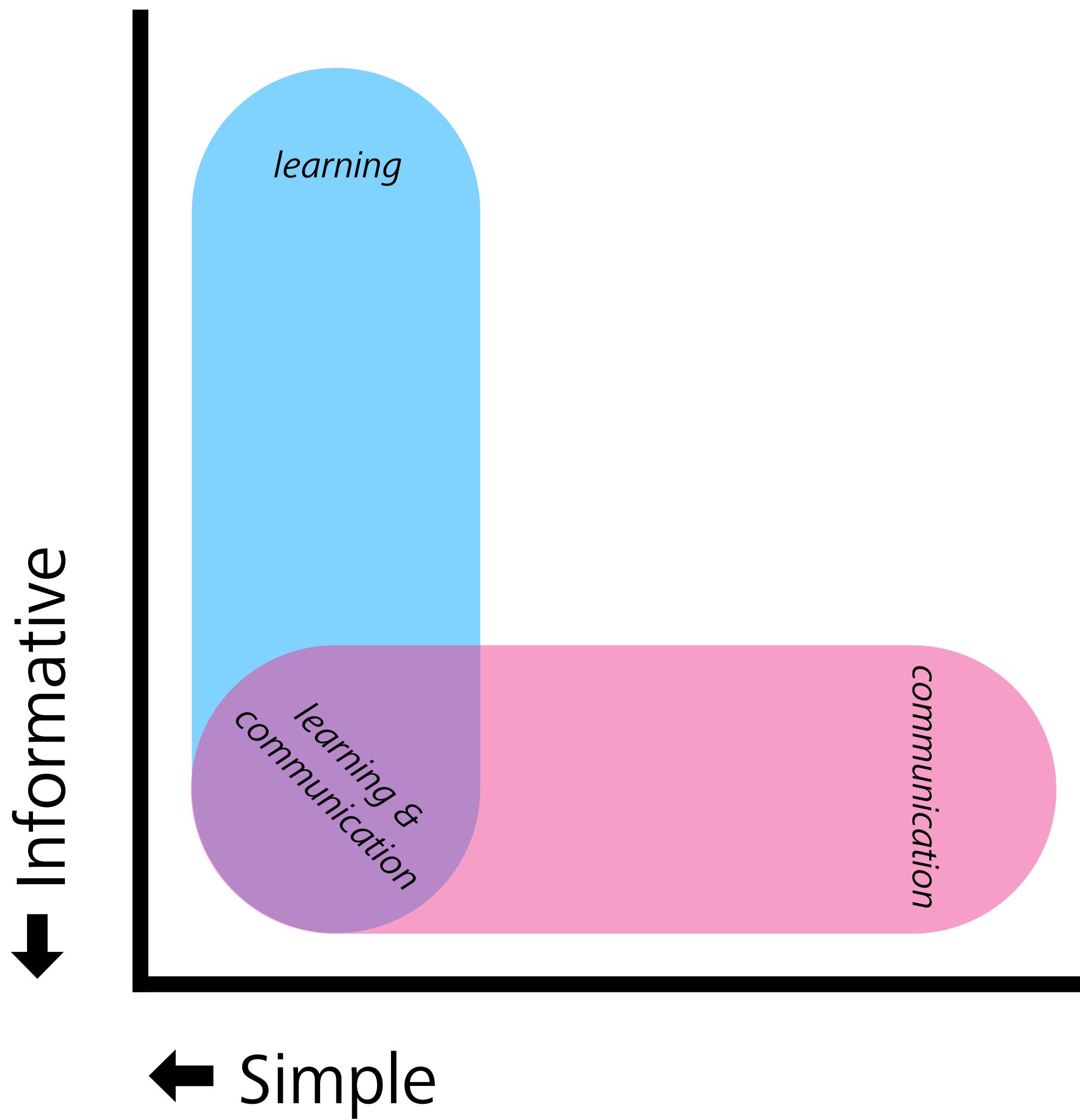
# Learning and communication pressures



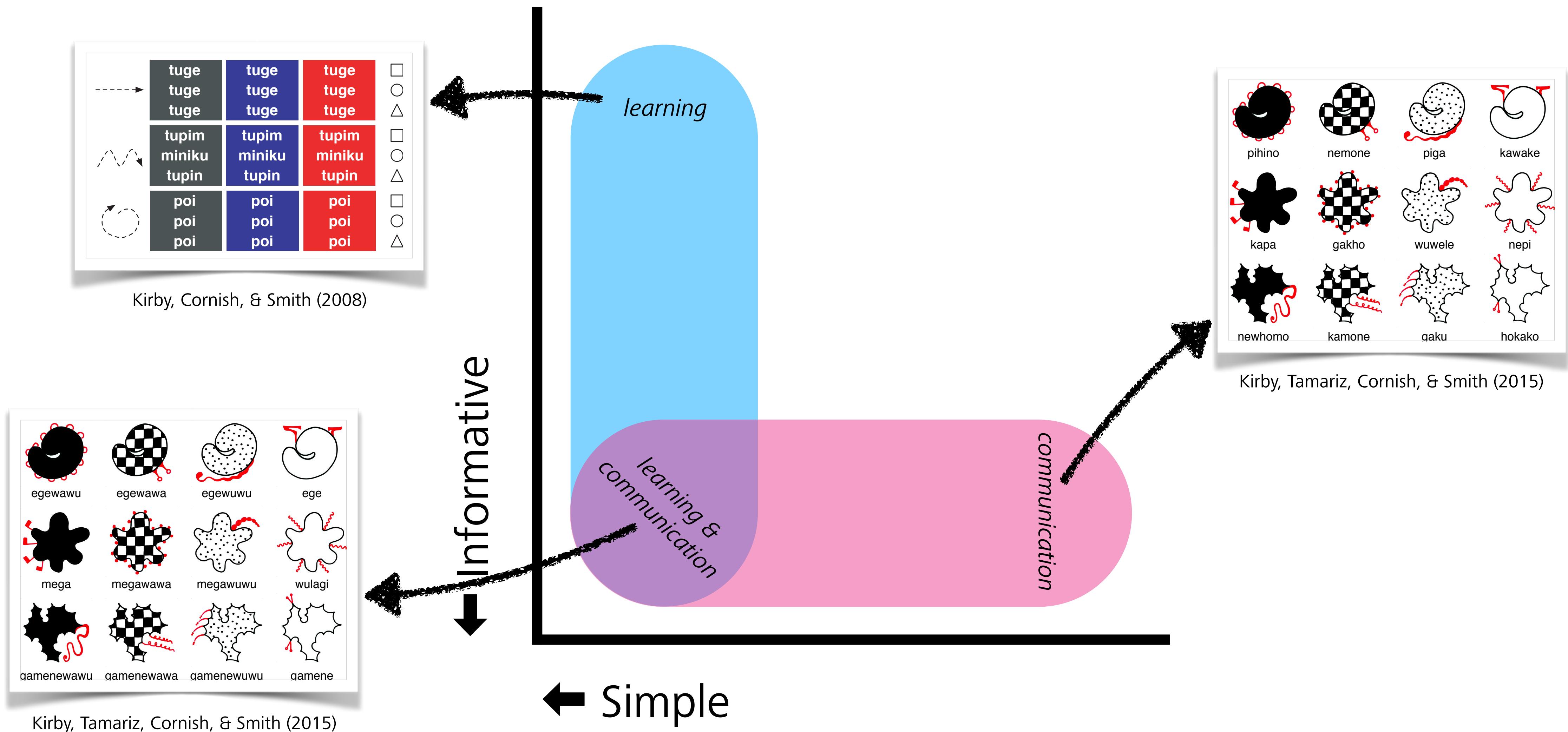
# Learning and communication pressures



# Learning and communication pressures



# Learning and communication pressures



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

# The Minimum Description Length principle

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

*For example...*

01001011110010000110001000101101100001111010001

Any regular

```
print('01001011110010000110001000101101100001111010001')
```

The more re

01

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

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

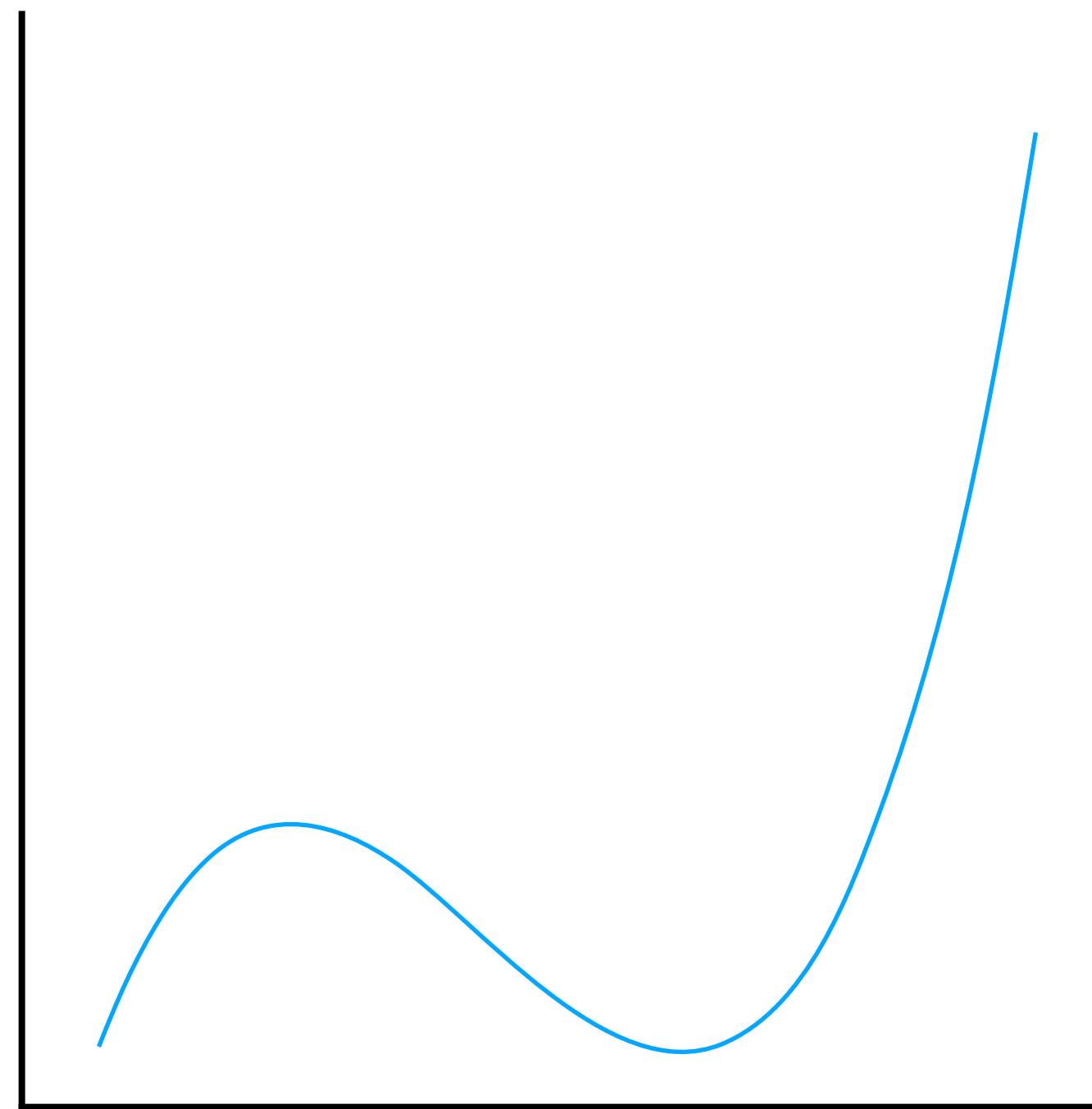
The more regularities there are, the more the data can be compressed

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

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

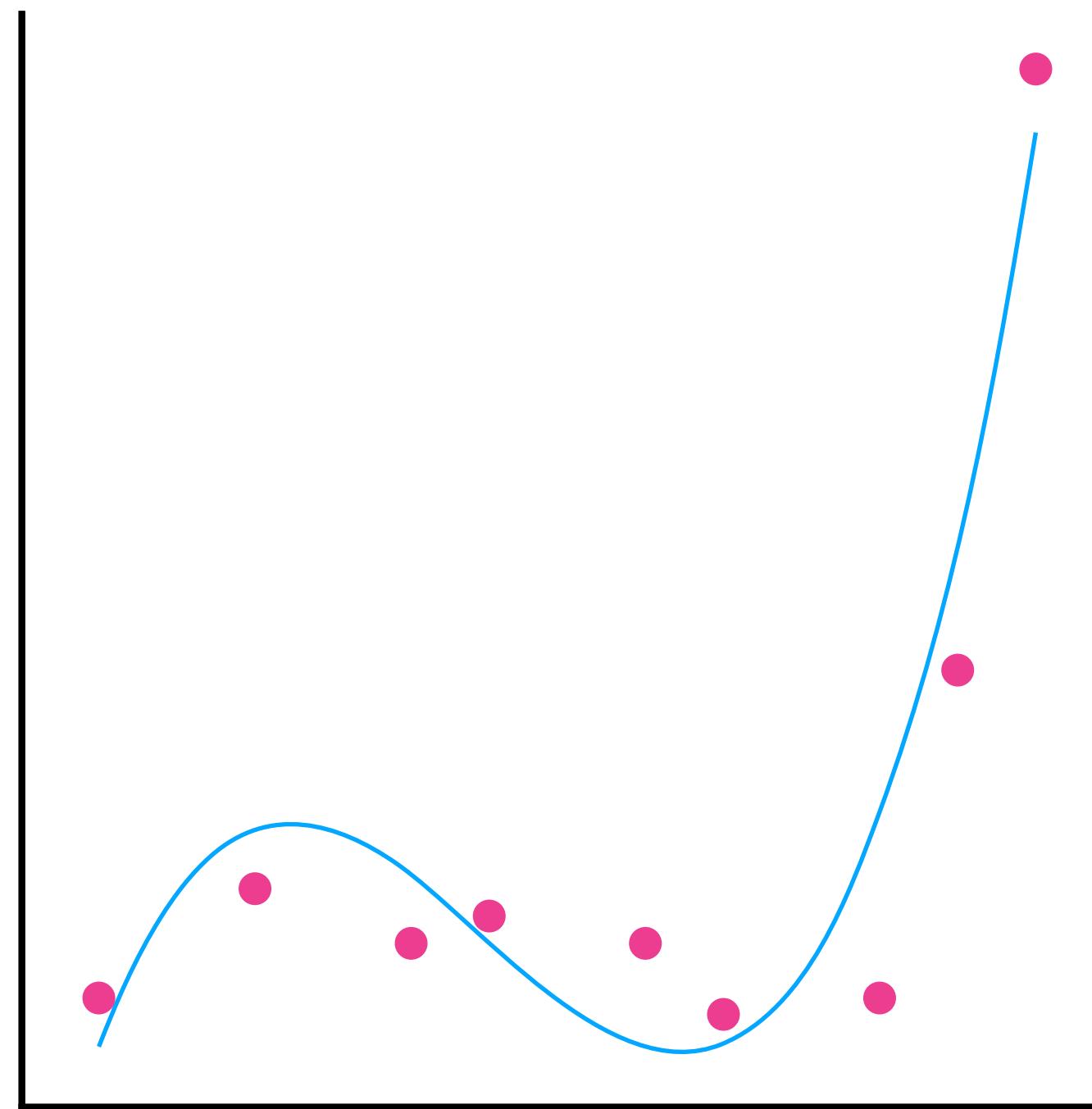
# The Minimum Description Length principle

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



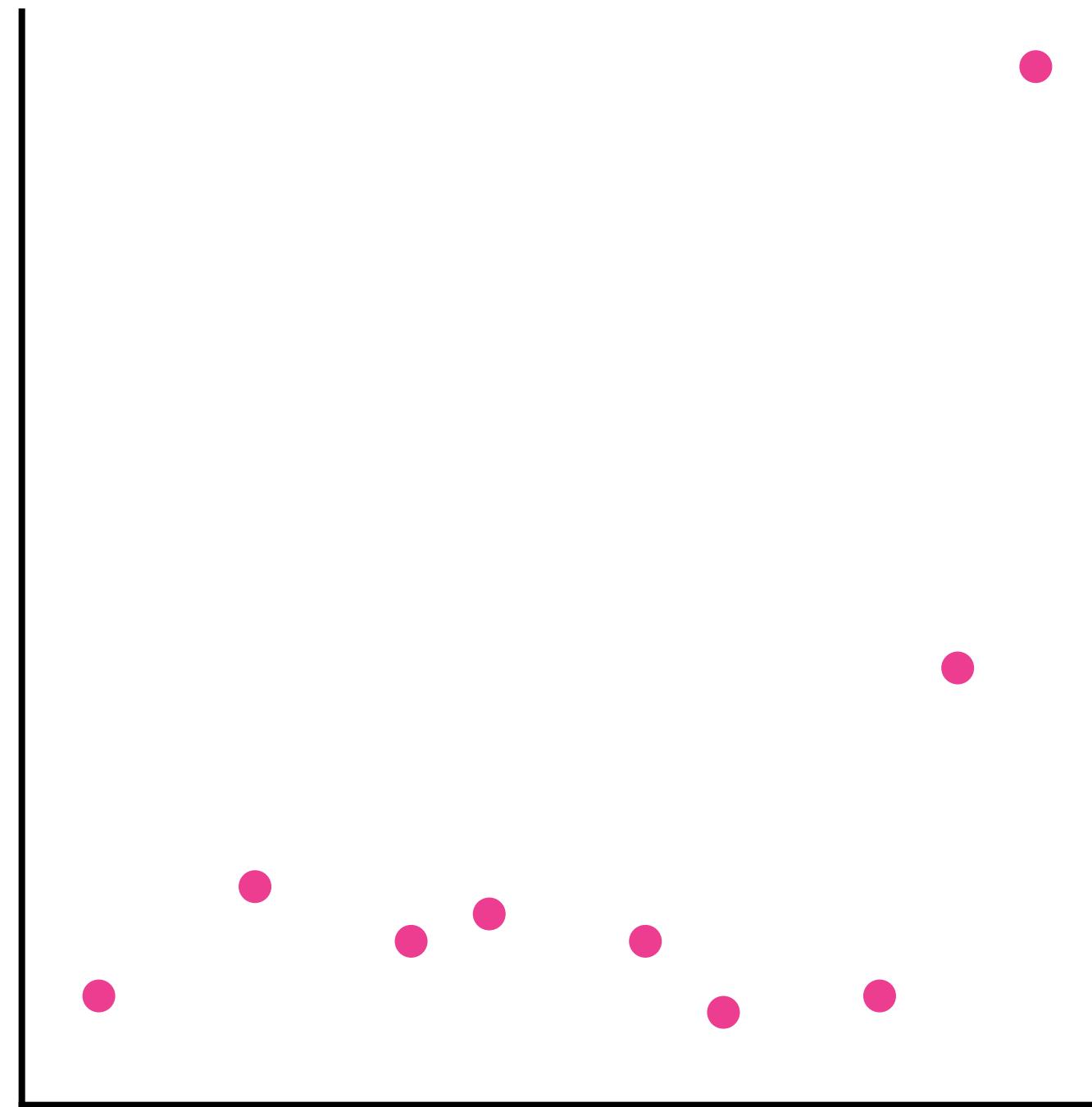
# The Minimum Description Length principle

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



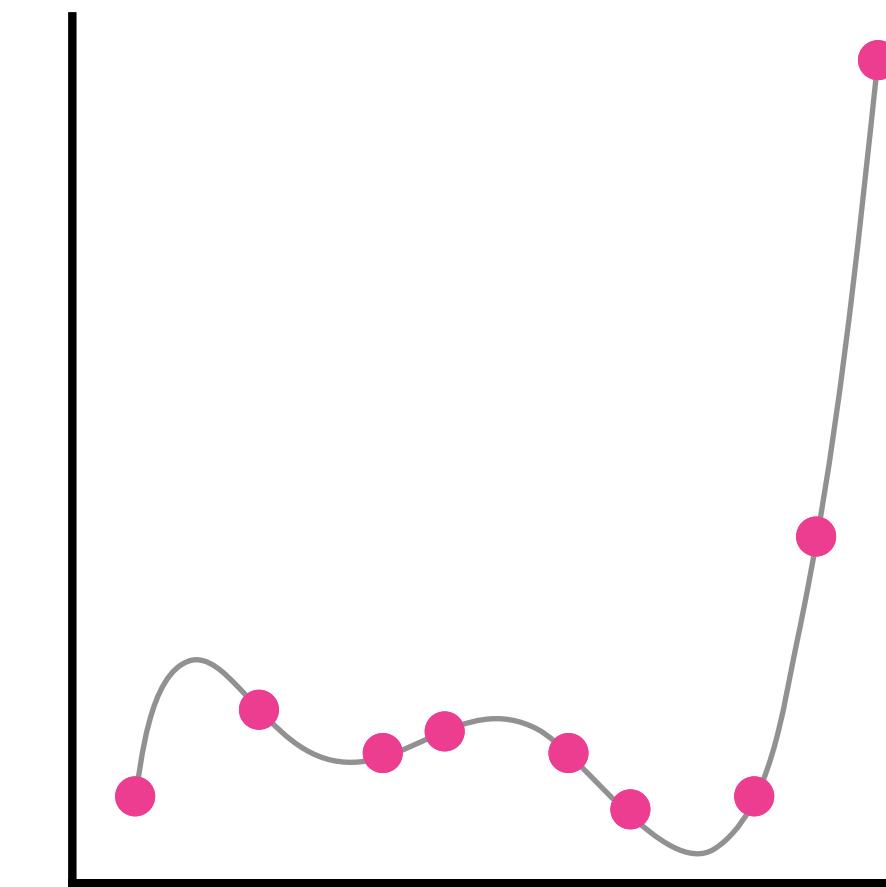
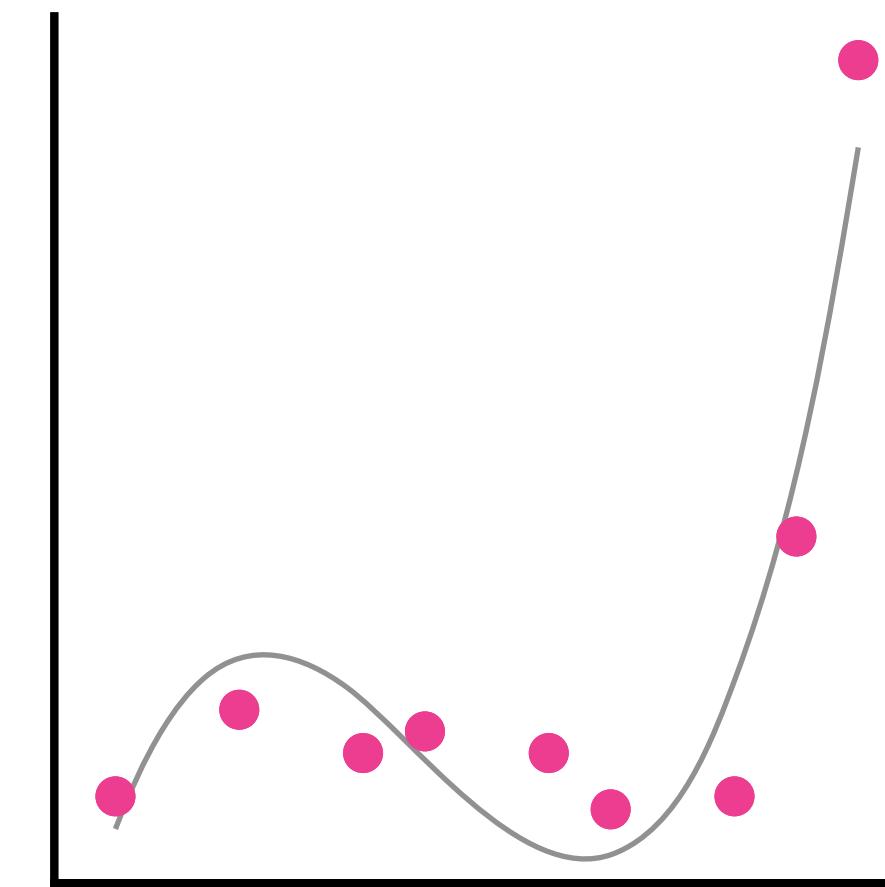
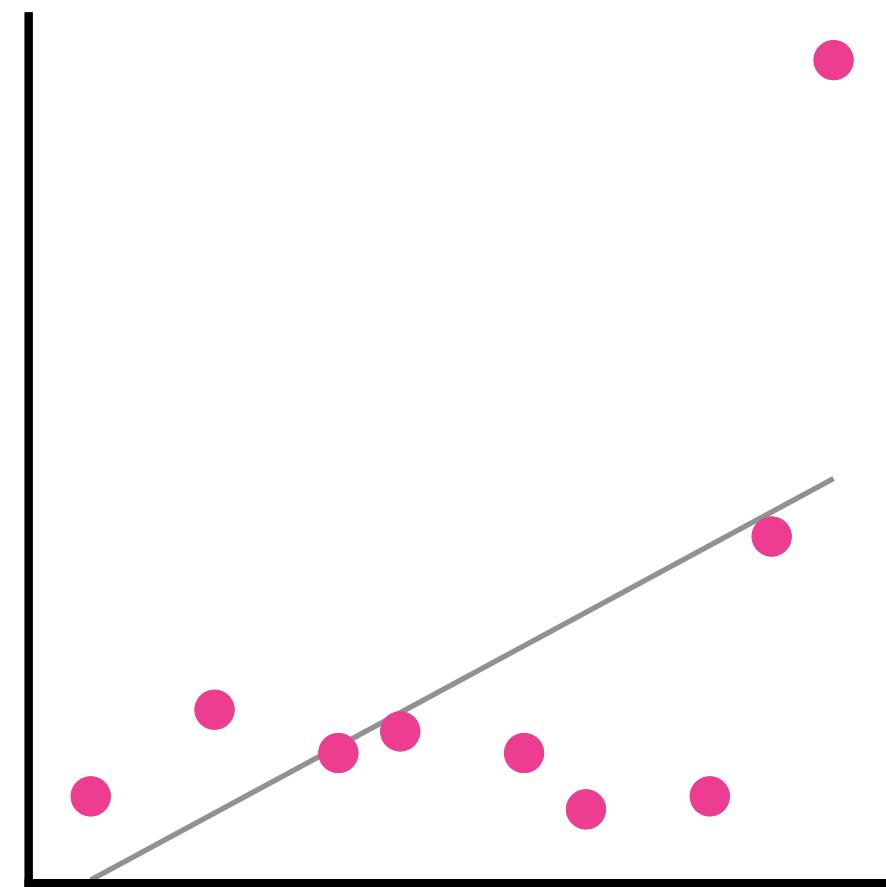
# The Minimum Description Length principle

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



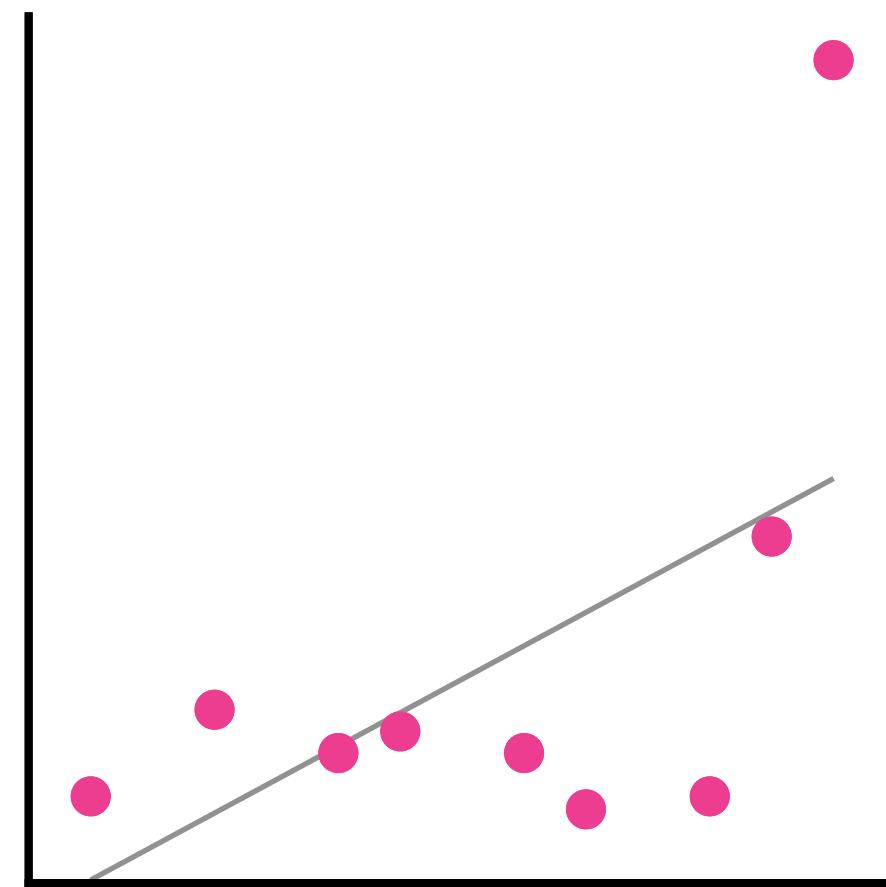
# The Minimum Description Length principle

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



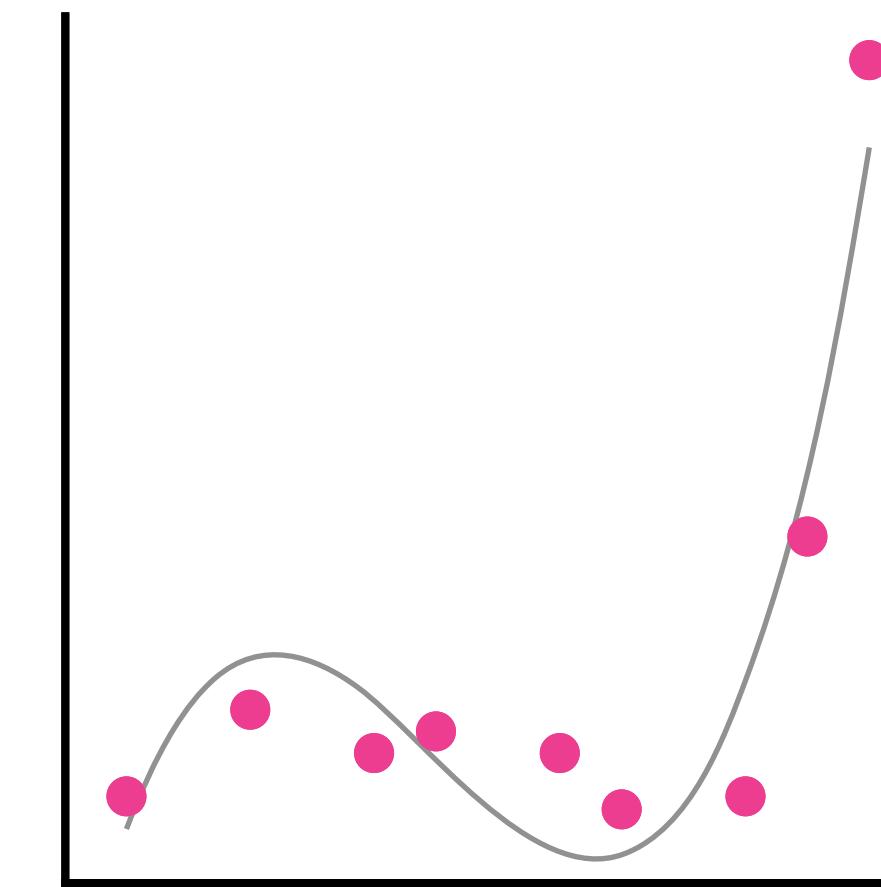
# The Minimum Description Length principle

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

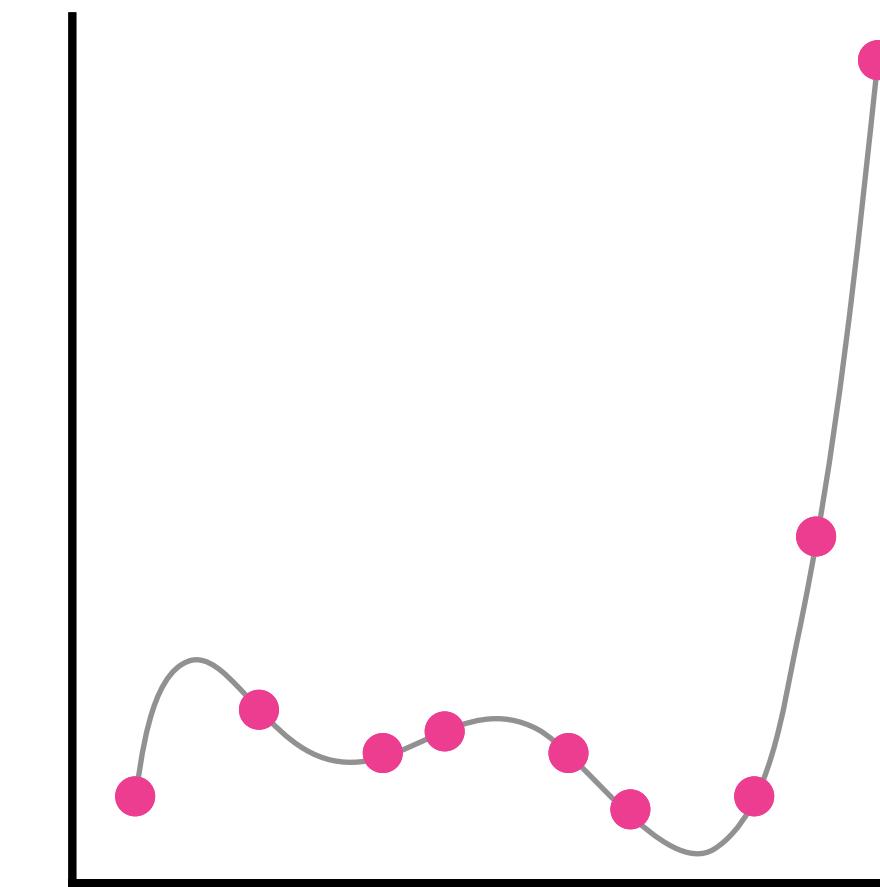


$\text{DL}(D|H)$

10 bits



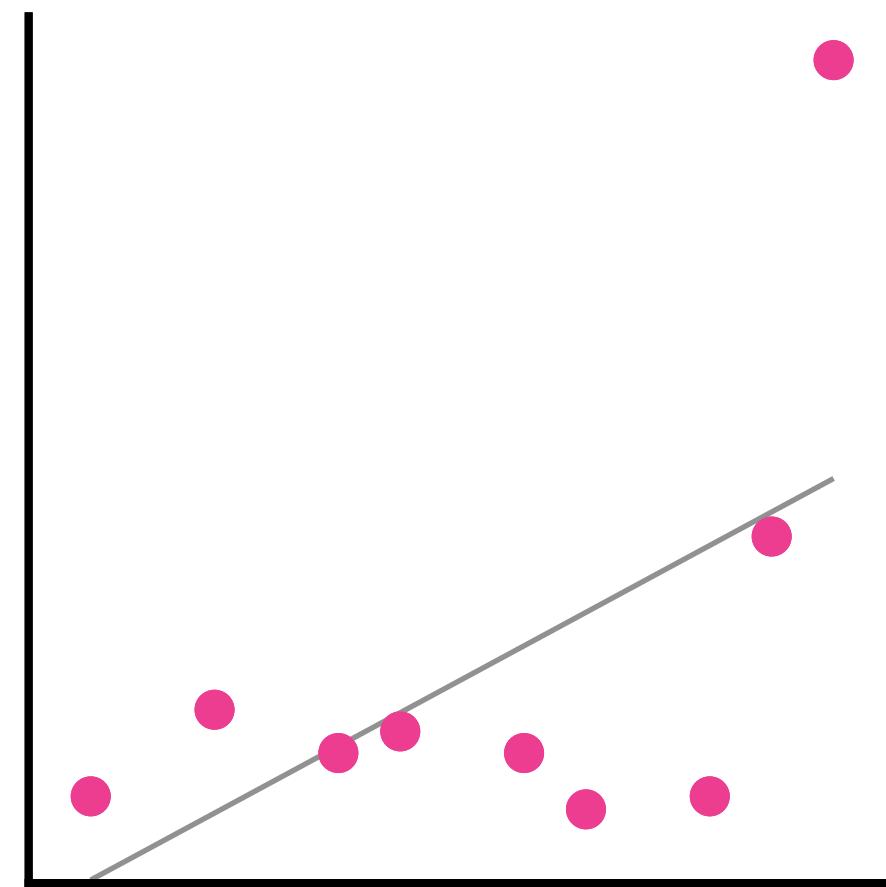
4 bits



0 bits

# The Minimum Description Length principle

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

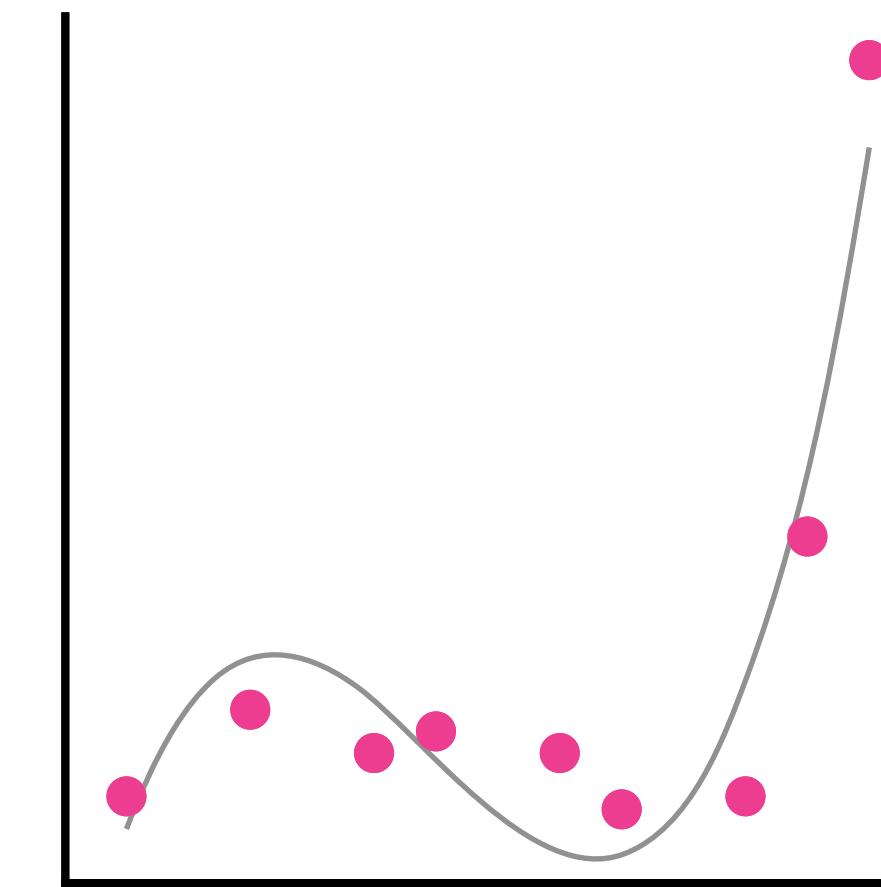


$\text{DL}(D|H)$

10 bits

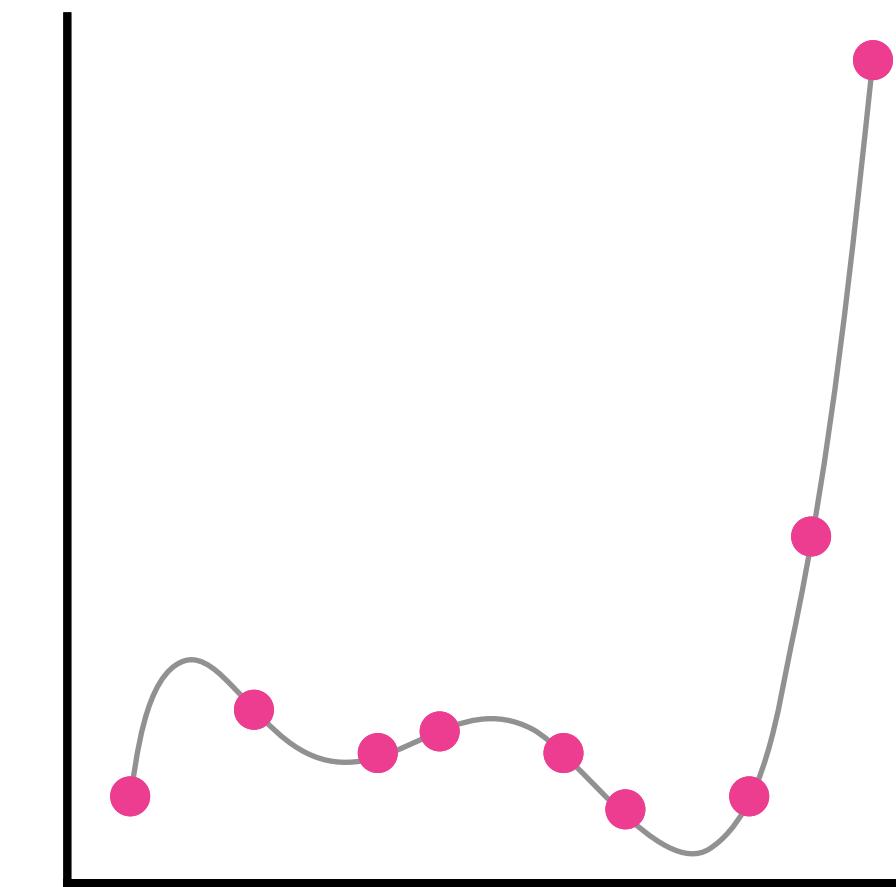
$\text{DL}(H)$

1 bit



4 bits

4 bits

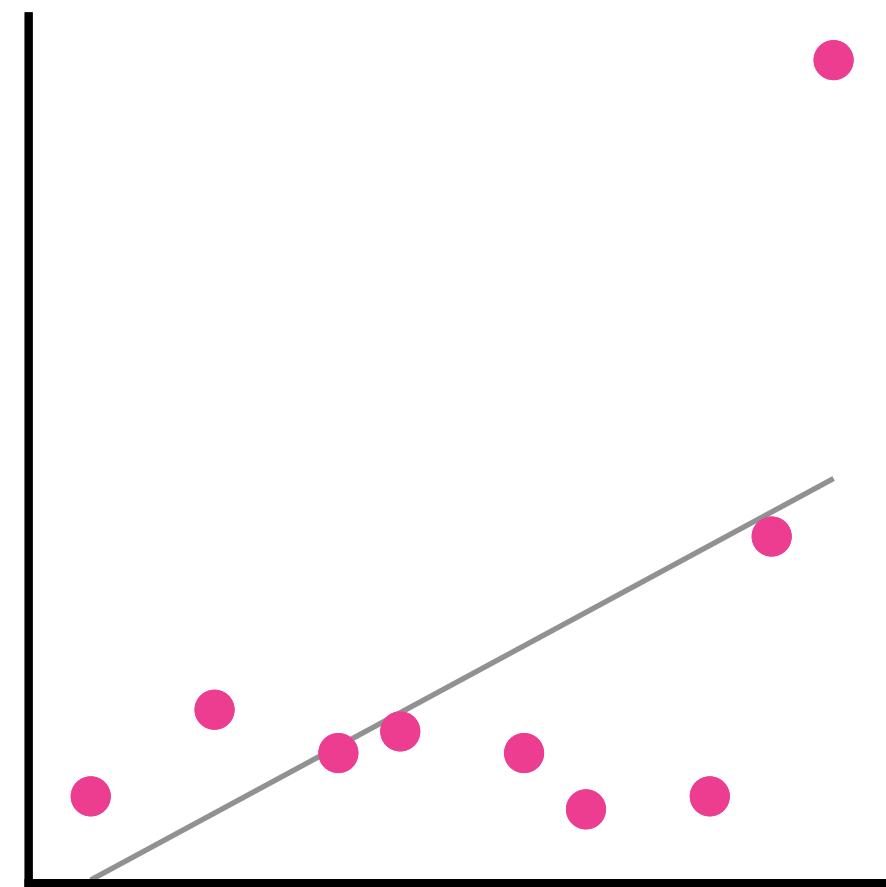


0 bits

10 bits

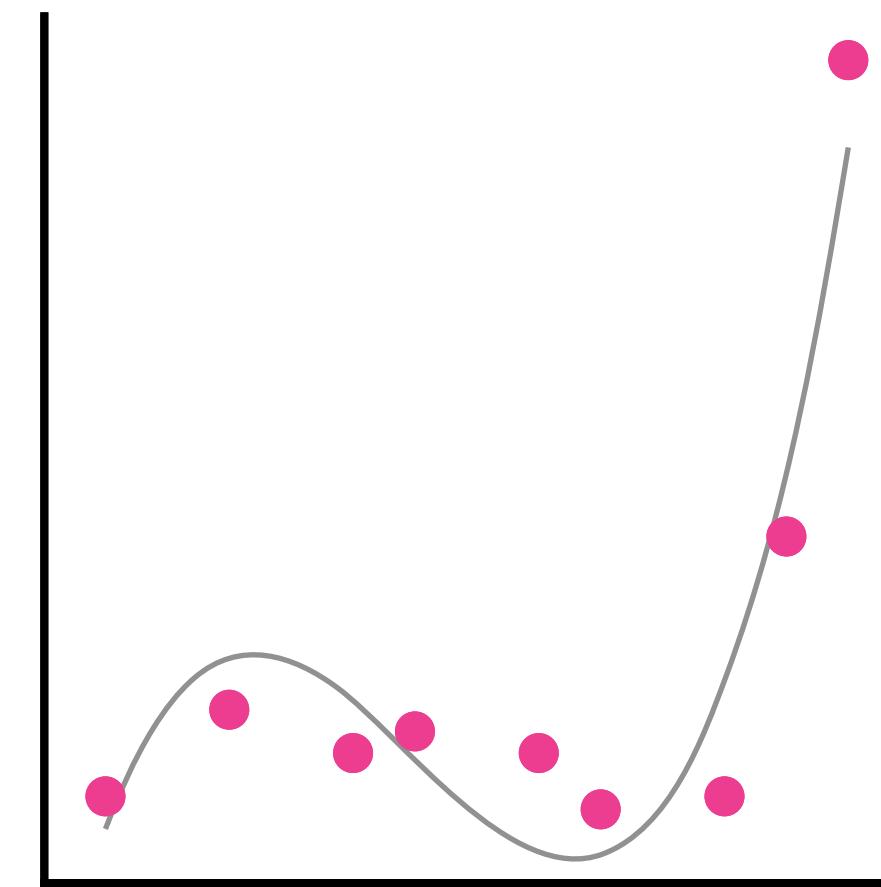
# The Minimum Description Length principle

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

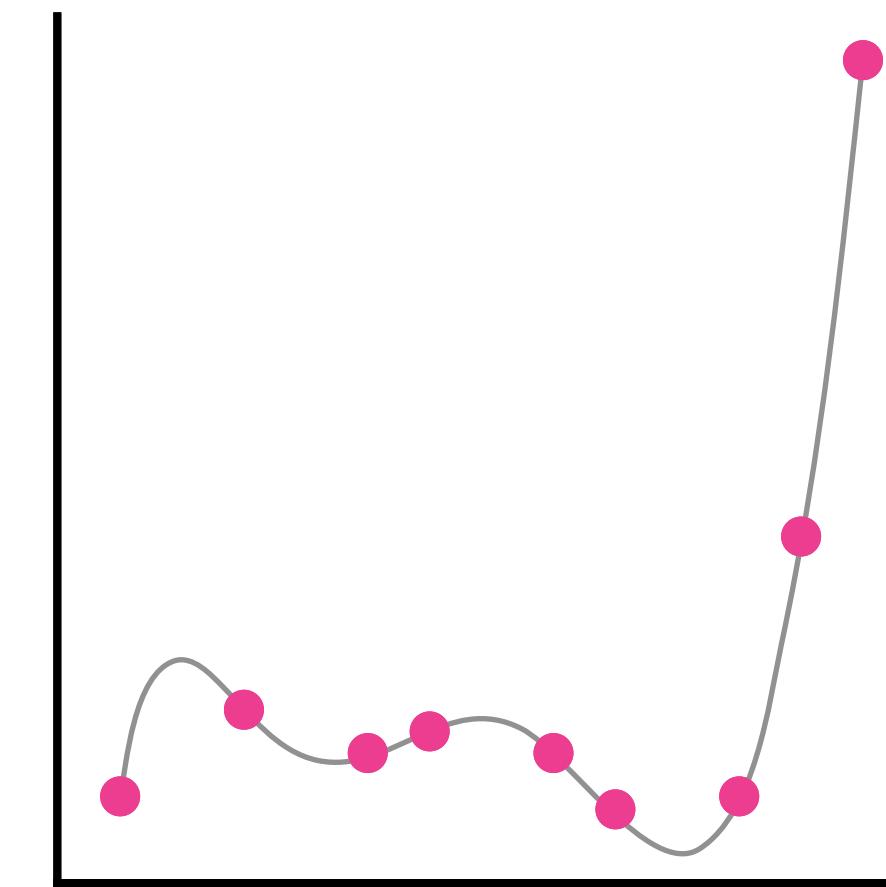


$\text{DL}(D|H)$

10 bits



4 bits



0 bits

$\text{DL}(H)$

1 bit

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

11 bits

4 bits

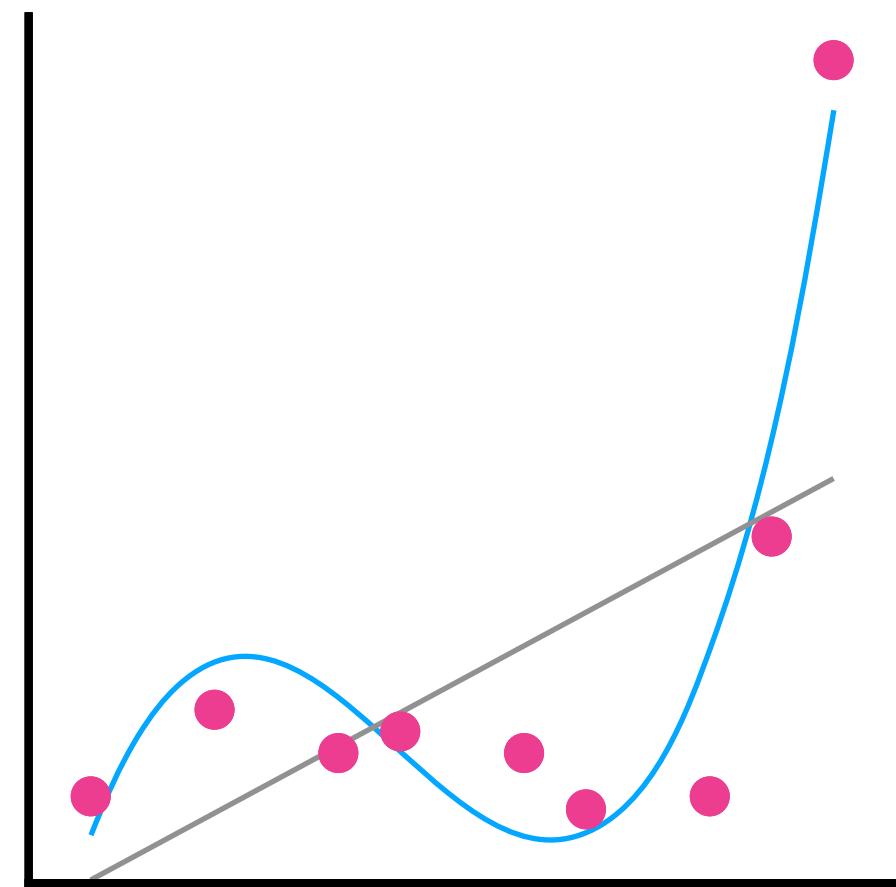
8 bits ✓

10 bits

10 bits

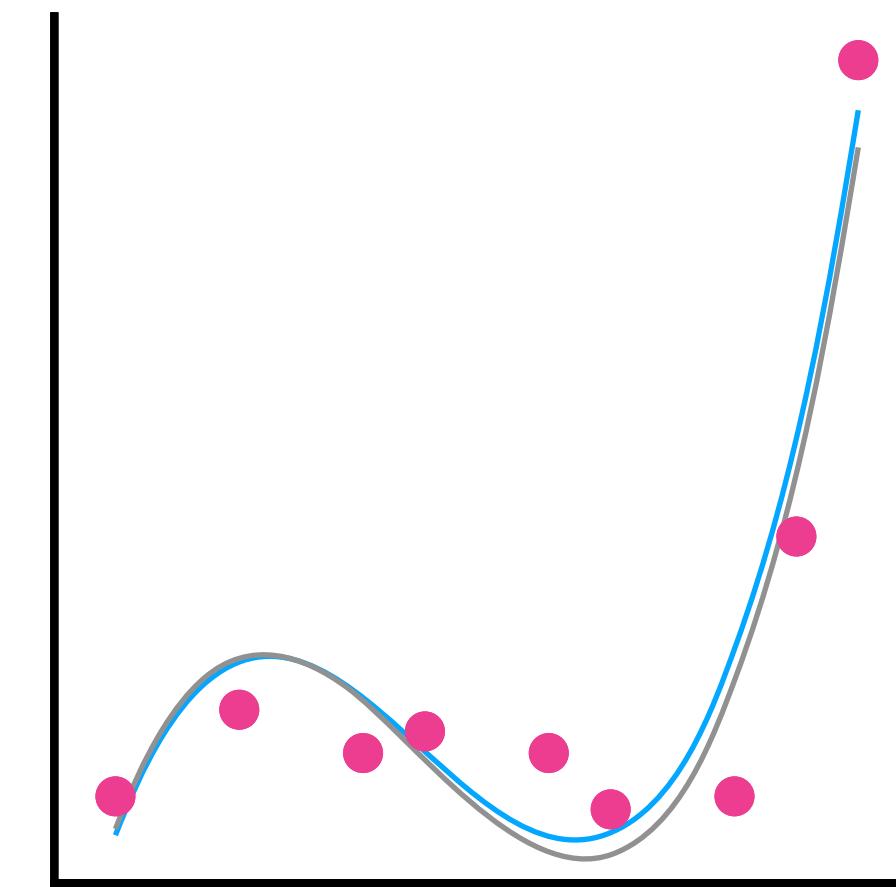
# The Minimum Description Length principle

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

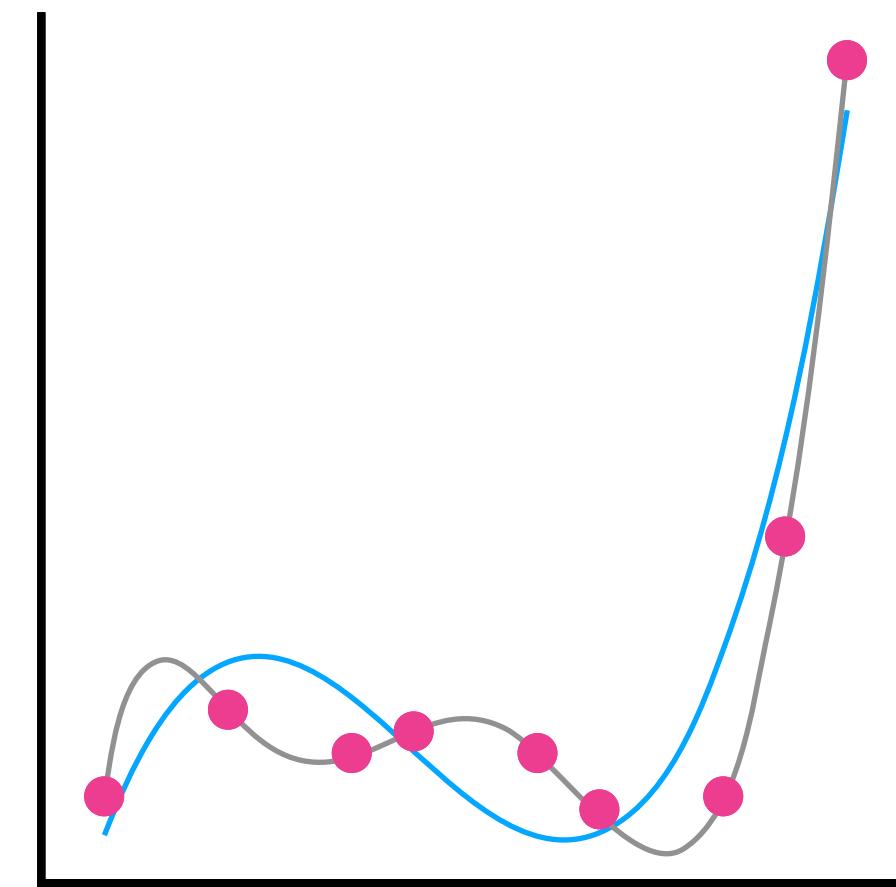


$\text{DL}(D|H)$

10 bits



4 bits



0 bits

$\text{DL}(H)$

1 bit

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

11 bits

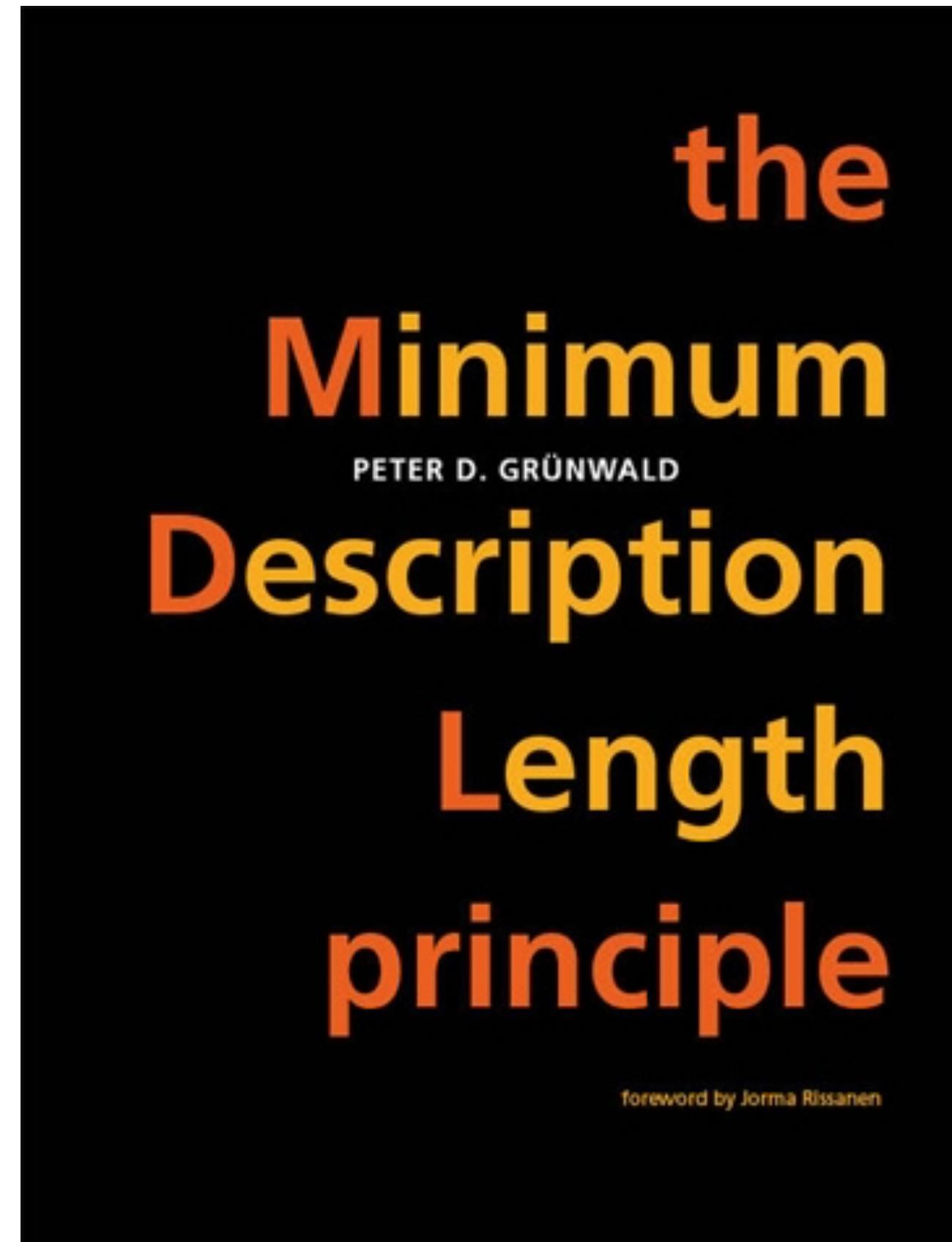
4 bits

8 bits ✓

10 bits

10 bits

# The Minimum Description Length principle



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

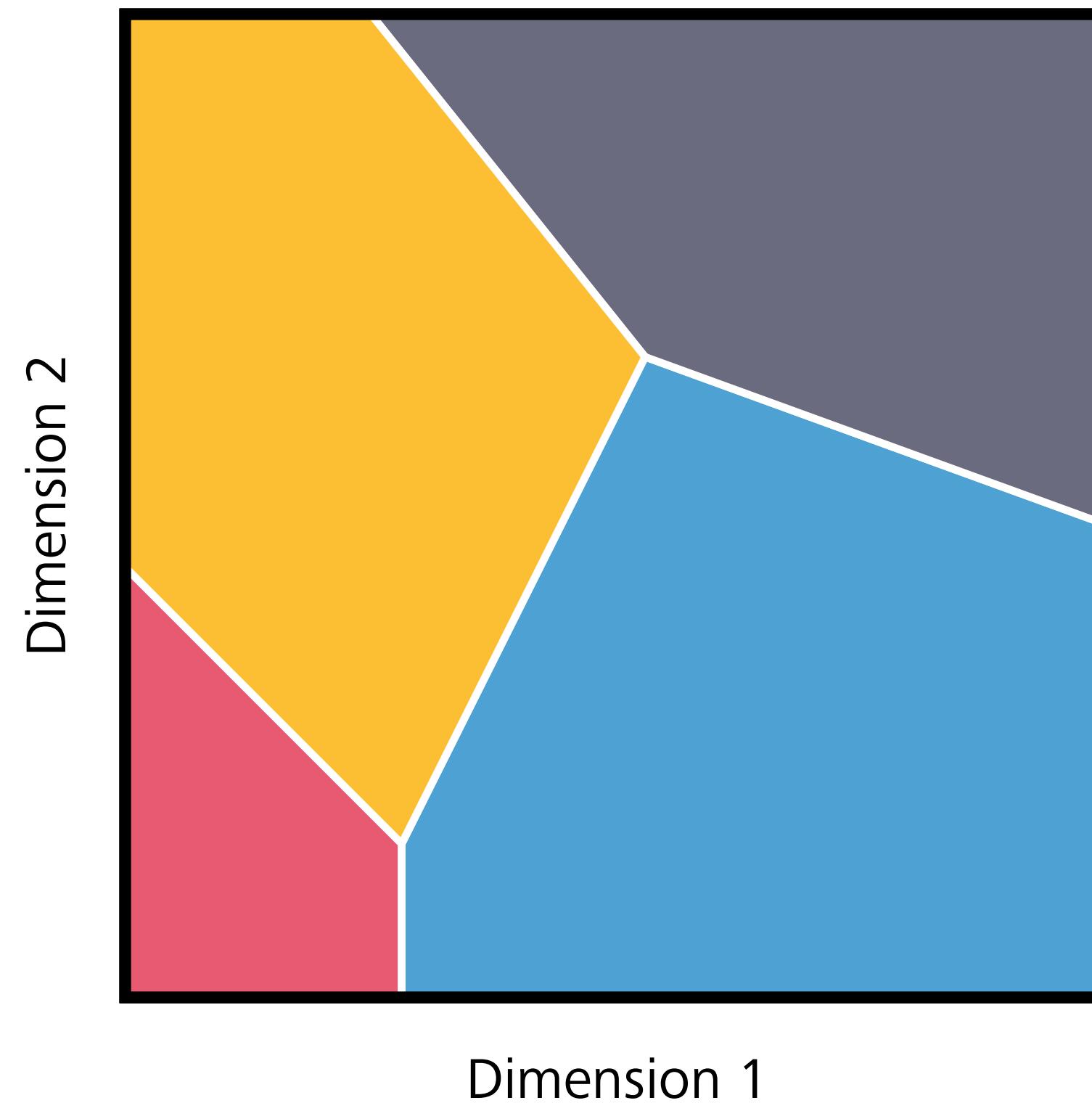
**Occam's razor:** MDL trades-off *goodness-of-fit* with *model complexity*, embodying Occam's razor

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

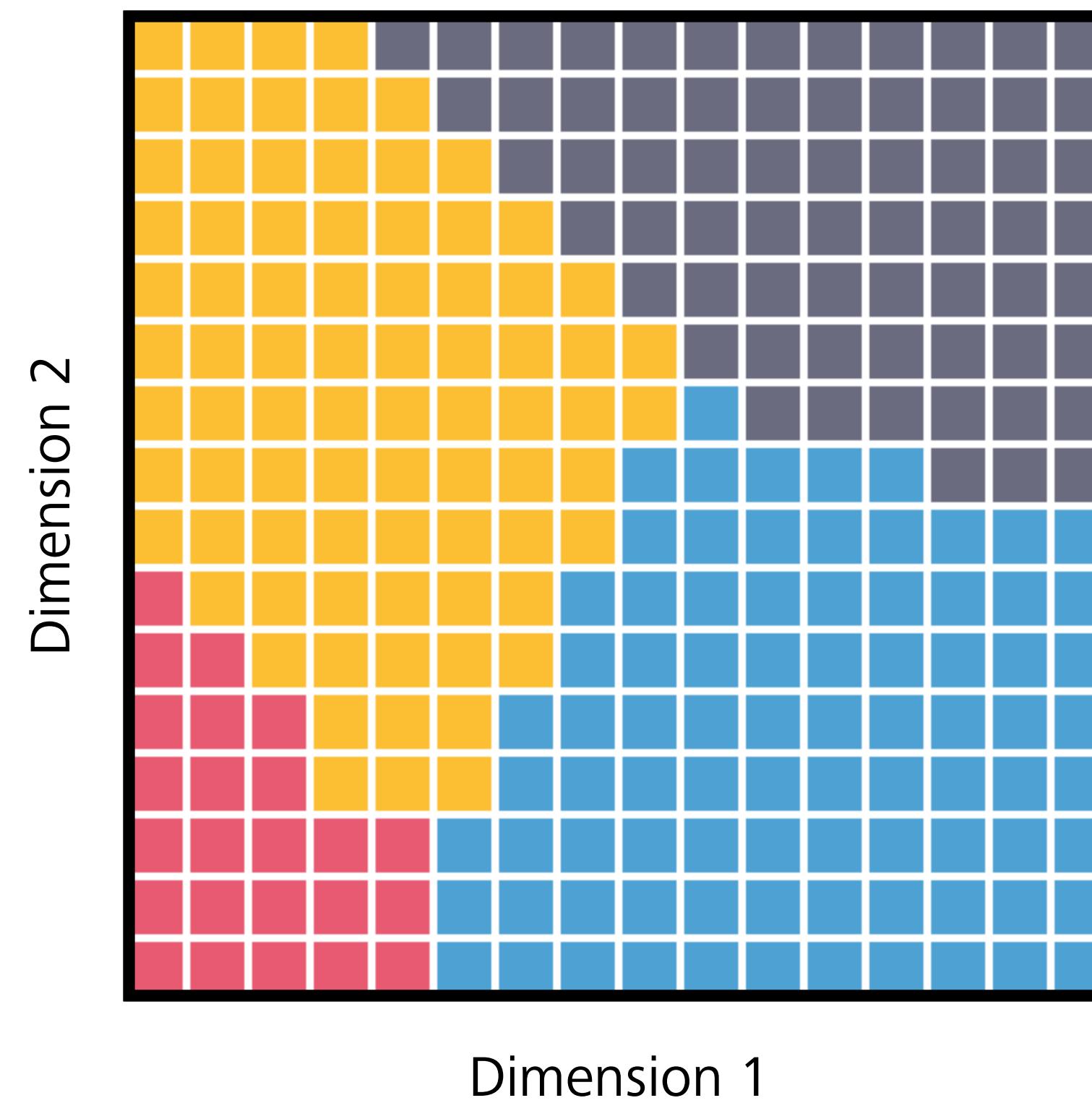
**Predictive performance:** Since data compression is formally equivalent to probabilistic prediction, MDL finds models offering *good predictive performance on unseen data*

*Bayesian model*

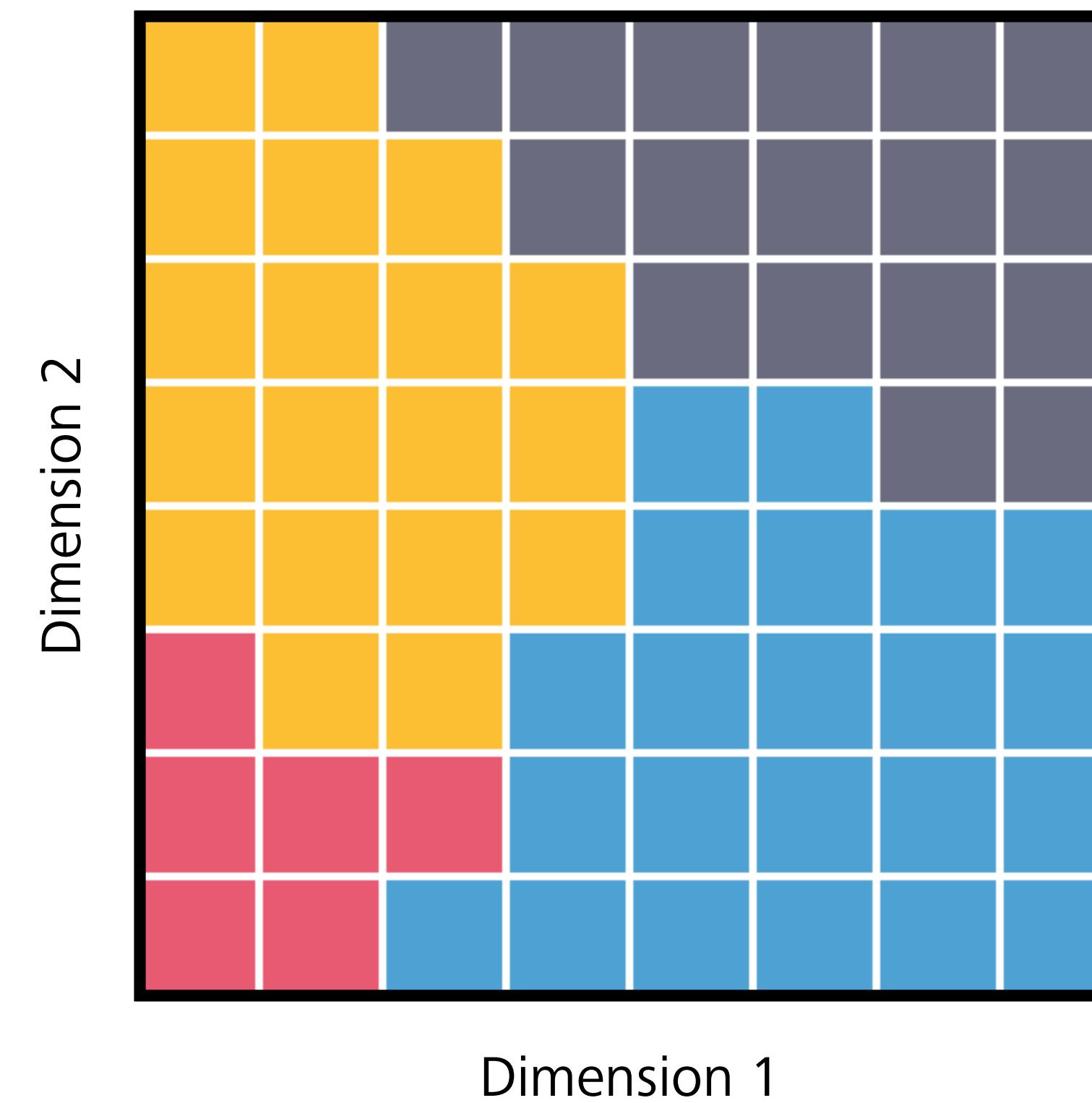
# Conceptual spaces and convexity



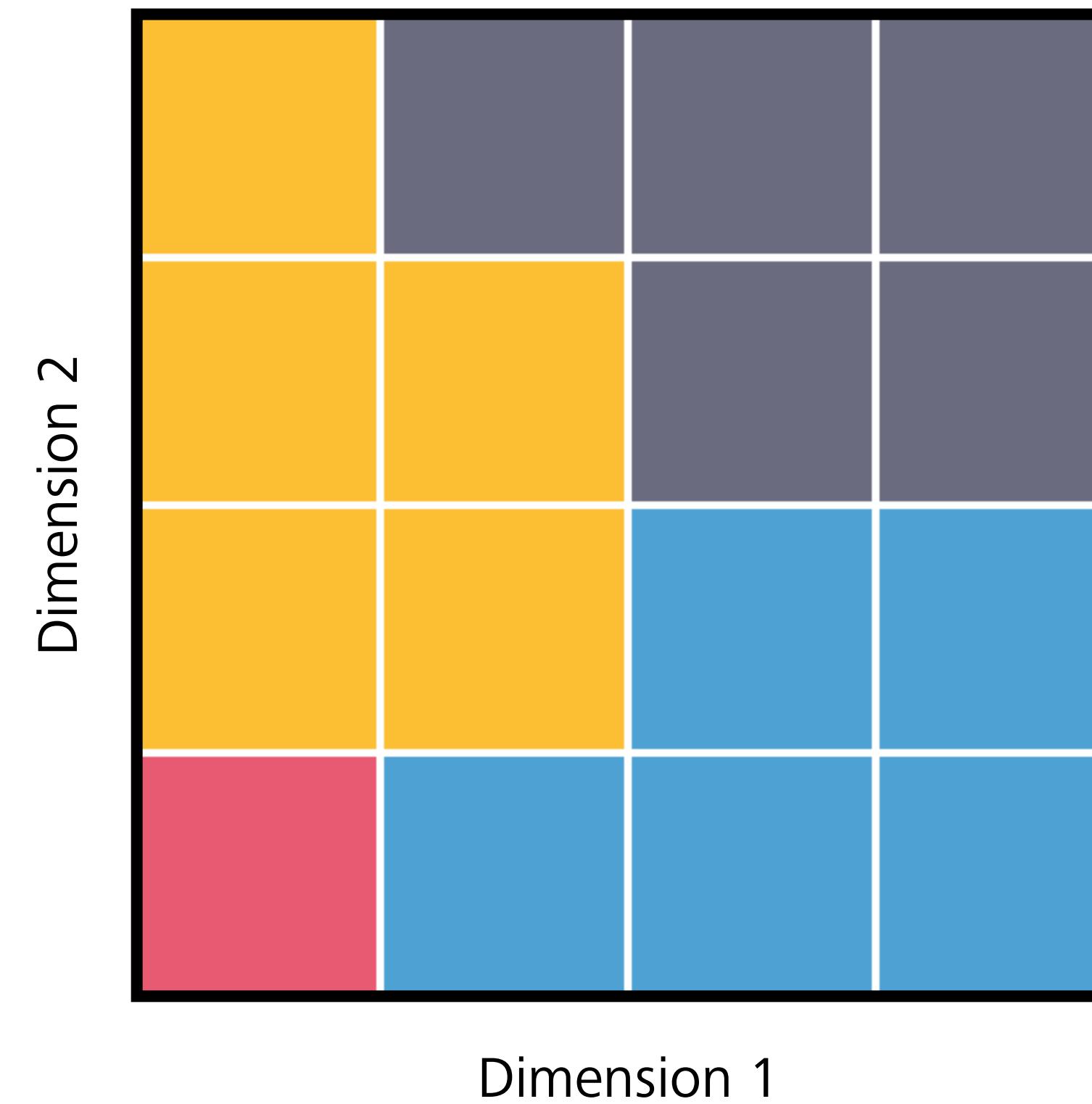
# Conceptual spaces and convexity



# Conceptual spaces and convexity



# Conceptual spaces and convexity



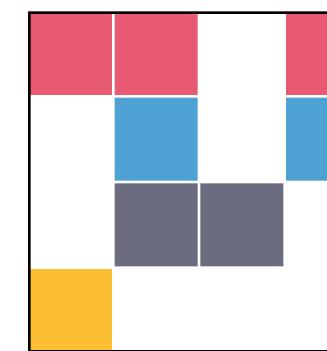
# Bayesian inference

$$\mathcal{L} = \left\{ \begin{array}{c} \text{A sequence of 10 4x4 grids representing posterior distributions.} \\ \text{Grid 1: Yellow (top-left), Red (bottom-left), Blue (bottom-right).} \\ \text{Grid 2: Red (top-left), Blue (middle-left), Yellow (bottom-left).} \\ \text{Grid 3: Red (top-left), Blue (middle-left), Blue (middle-right).} \\ \text{Grid 4: Blue (top-left), Blue (middle-left), Blue (middle-right).} \\ \text{Grid 5: Blue (top-left), Red (middle-left), Red (middle-right).} \\ \text{Grid 6: All Blue (solid).} \\ \text{Grid 7: Yellow (top-left), Blue (middle-left), Blue (middle-right).} \\ \text{Grid 8: Red (top-left), Blue (middle-left), Yellow (middle-right).} \\ \text{Grid 9: Blue (top-left), Blue (middle-left), Blue (middle-right).} \\ \text{Grid 10: Red (top-left), Blue (middle-left), Blue (middle-right).} \end{array} \right. \dots \right\}$$

# Bayesian inference

$$\mathcal{L} = \left\{ \begin{array}{c} \text{Figure 1} \\ \text{Figure 2} \\ \text{Figure 3} \\ \text{Figure 4} \\ \text{Figure 5} \\ \text{Figure 6} \\ \text{Figure 7} \\ \vdots \end{array} \right\}$$

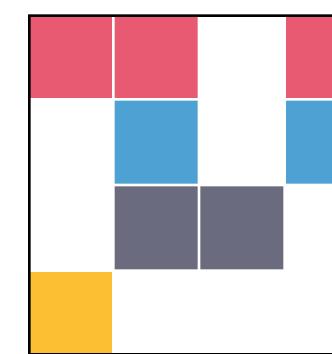
$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$



# Bayesian inference

$$\mathcal{L} = \{ \begin{array}{ccccccc} \text{grid 1} & \text{grid 2} & \text{grid 3} & \text{grid 4} & \text{grid 5} & \text{grid 6} & \text{grid 7} \\ \text{grid 8} & \text{grid 9} & \text{grid 10} & \text{grid 11} & \text{grid 12} & \text{grid 13} & \text{grid 14} \\ \dots \end{array} \}$$

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



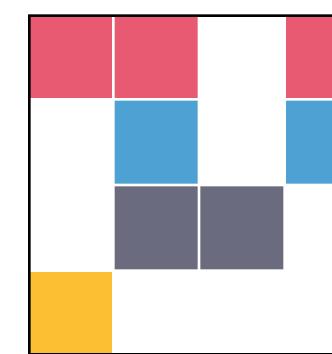
$$\text{likelihood}(D|L) \propto \prod_{\langle m, s \rangle \in D} \frac{1}{|M|} P(s|L, m)$$

$$= \begin{array}{c} \text{grid 1} \\ \text{grid 2} \\ \text{grid 3} \\ \text{grid 4} \end{array}$$

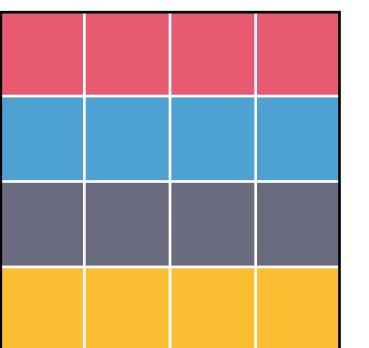
# Bayesian inference

$$\mathcal{L} = \{ \begin{array}{ccccccc} \text{grid 1} & \text{grid 2} & \text{grid 3} & \text{grid 4} & \text{grid 5} & \text{grid 6} & \text{grid 7} \\ \text{grid 8} & \text{grid 9} & \text{grid 10} & \text{grid 11} & \text{grid 12} & \text{grid 13} & \text{grid 14} \\ \dots \end{array} \}$$

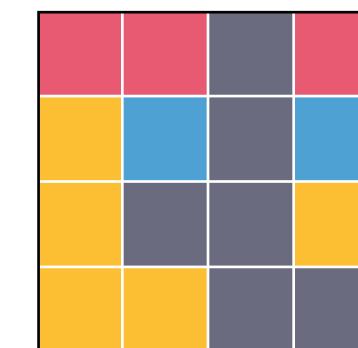
$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$



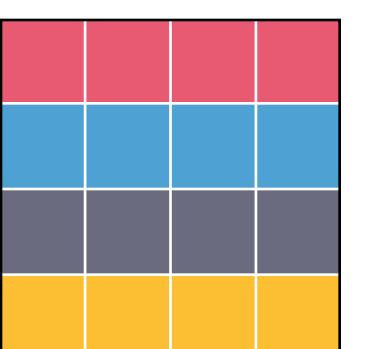
$$\text{likelihood}(D|L) \propto \prod_{\langle m, s \rangle \in D} \frac{1}{|M|} P(s|L, m)$$



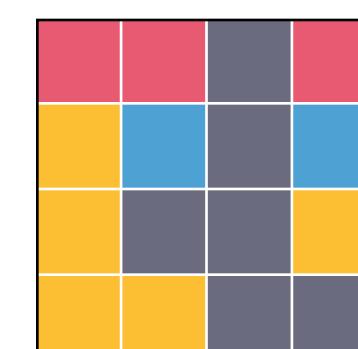
=



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



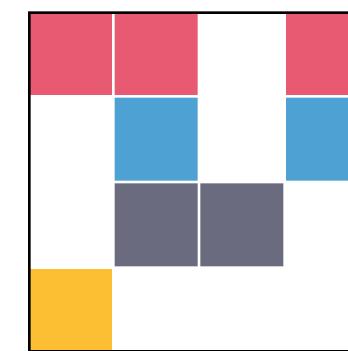
>



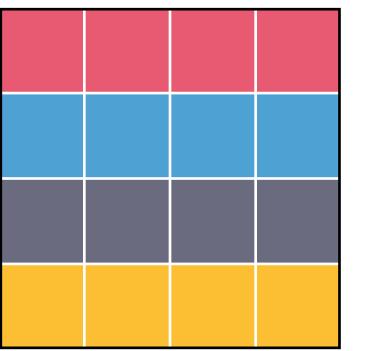
# Bayesian inference

$$\mathcal{L} = \{ \begin{array}{ccccccccc} \text{grid 1} & \text{grid 2} & \text{grid 3} & \text{grid 4} & \text{grid 5} & \text{grid 6} & \text{grid 7} & \text{grid 8} & \dots \end{array} \}$$

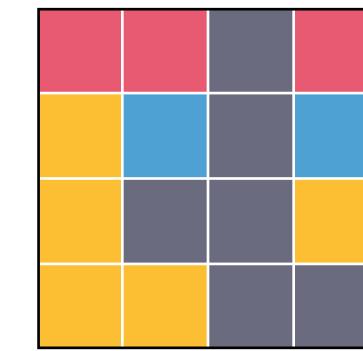
$$D = [\langle m_1, s_1 \rangle, \langle m_2, s_2 \rangle, \langle m_3, s_3 \rangle, \dots, \langle m_n, s_n \rangle]$$



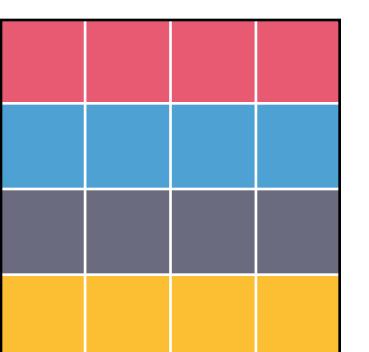
$$\text{likelihood}(D|L) \propto \prod_{\langle m, s \rangle \in D} \frac{1}{|M|} P(s|L, m)$$



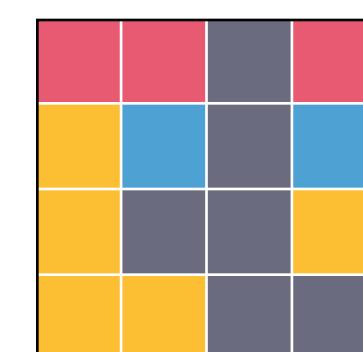
=



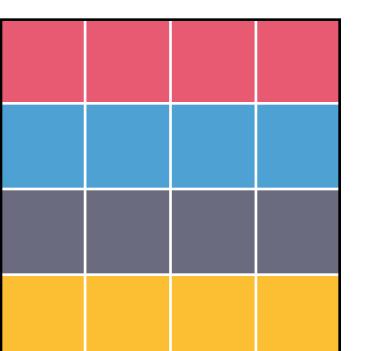
$$\text{prior}(L) \propto 2^{-\text{DL}(L)}$$



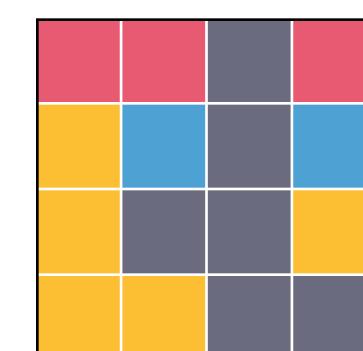
>



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



>



# Computing $DL(L)$ : The rectangle search

Class	Position
1x1	16
1x2	12
1x3	8
1x4	4
2x1	12
2x2	9
2x3	6
2x4	3
3x1	8
3x2	6
3x3	4
3x4	2
4x1	4
4x2	3
4x3	2
4x4	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](mailto:dfass@ruccs.rutgers.edu)

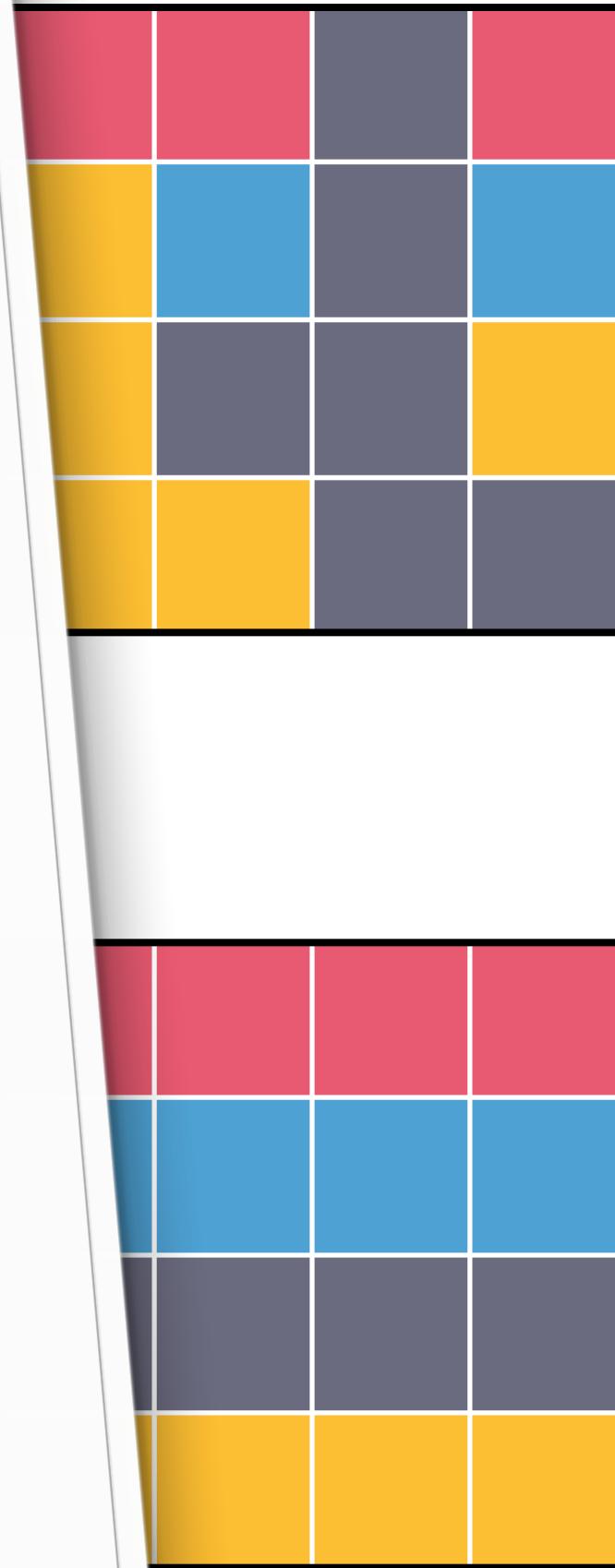
**Jacob Feldman\***  
Department of Psychology  
Center for Cognitive Science  
Rutgers University  
Piscataway, NJ 08854  
[jacob@ruccs.rutgers.edu](mailto:jacob@ruccs.rutgers.edu)

### 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.

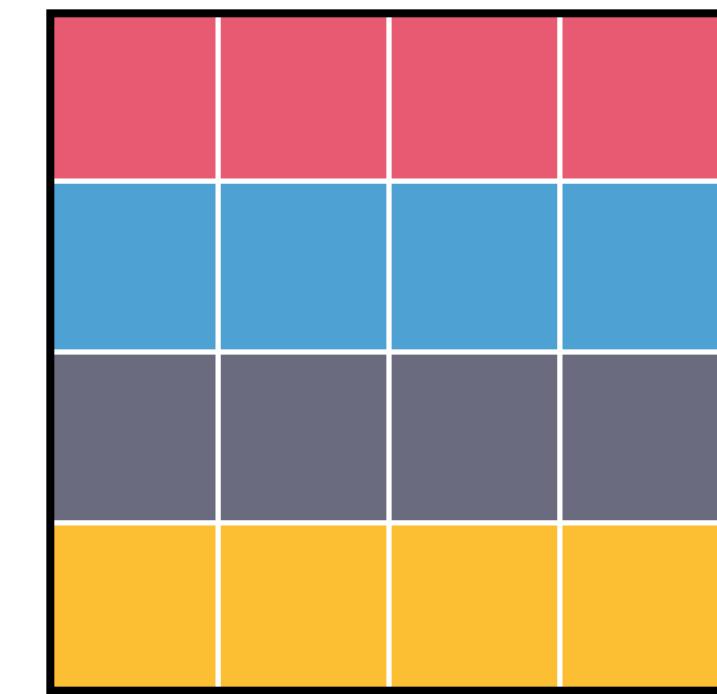
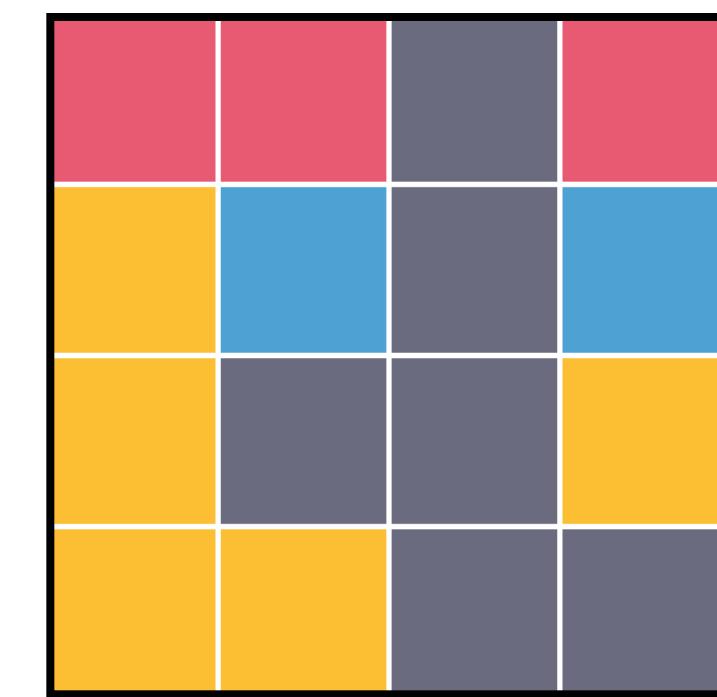
### 1 The Structure of Categories

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 ease of *inference* [3][4]. While many of these theories are mathematically elegant, they have had limited success in accounting for a range of experimental findings, similar to those reported here.



# Computing $DL(L)$ : The rectangle code

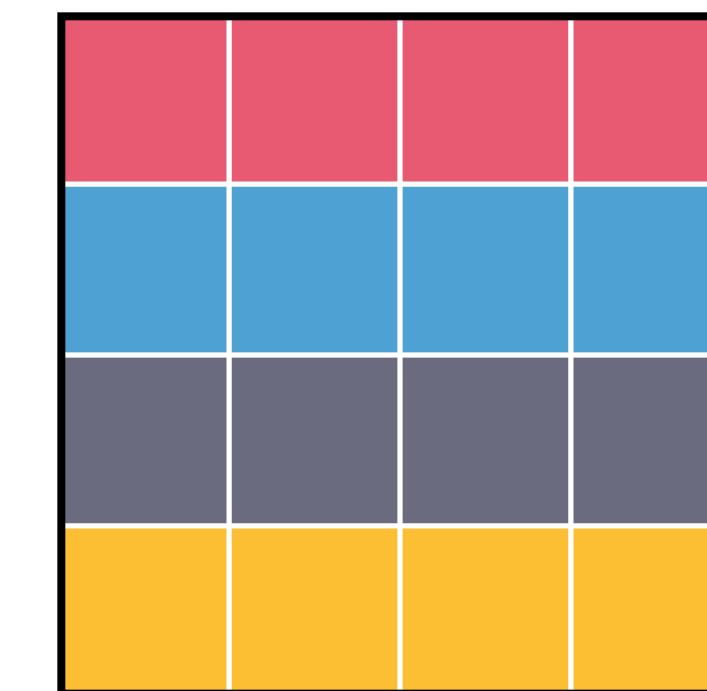
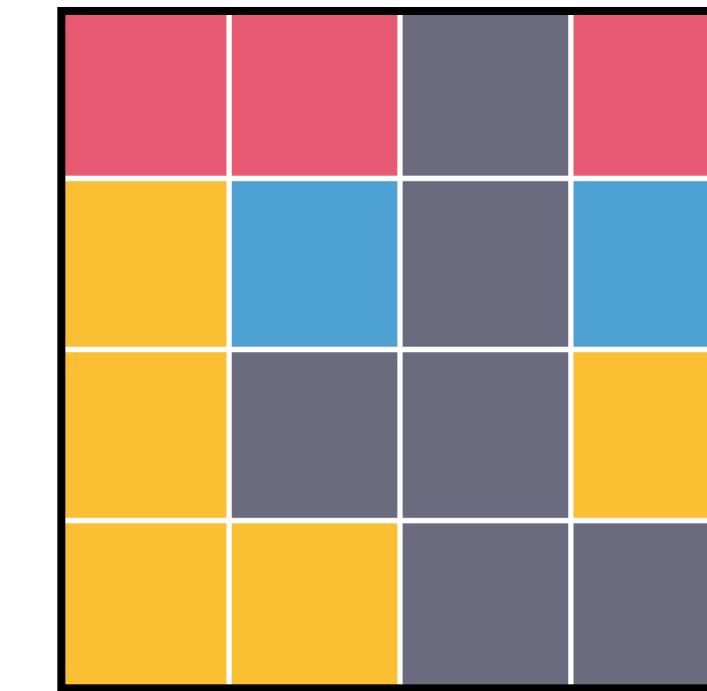
<b>Class</b>	<b>Positions</b>	<b>Probability</b>	<b>Codelength</b>	(bits)
1×1	16	$1/16 \times 1/16$	$-\log 1/256$	8.0
1×2	12	$1/16 \times 1/12$	$-\log 1/192$	7.58
1×3	8	$1/16 \times 1/8$	$-\log 1/128$	7.0
1×4	4	$1/16 \times 1/4$	$-\log 1/64$	6.0
2×1	12	$1/16 \times 1/12$	$-\log 1/192$	7.58
2×2	9	$1/16 \times 1/9$	$-\log 1/144$	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



# Computing $DL(L)$ : The rectangle code

Class	Positions	Probability	Codelength (bits)
1x1	16	$1/16 \times 1/16$	$-\log 1/256$ 8.0
1x2	12	$1/16 \times 1/12$	$-\log 1/192$ 7.58
1x3	8	$1/16 \times 1/8$	$-\log 1/128$ 7.0
1x4	4	$1/16 \times 1/4$	$-\log 1/64$ 6.0
2x1	12	$1/16 \times 1/12$	$-\log 1/192$ 7.58
2x2	9	$1/16 \times 1/9$	$-\log 1/144$ 7.17
2x3	6	$1/16 \times 1/6$	$-\log 1/96$ 6.58
2x4	3	$1/16 \times 1/3$	$-\log 1/48$ 5.58
3x1	8	$1/16 \times 1/8$	$-\log 1/128$ 7.0
3x2	6	$1/16 \times 1/6$	$-\log 1/96$ 6.58
3x3	4	$1/16 \times 1/4$	$-\log 1/64$ 6.0
3x4	2	$1/16 \times 1/2$	$-\log 1/32$ 5.0
4x1	4	$1/16 \times 1/4$	$-\log 1/64$ 6.0
4x2	3	$1/16 \times 1/3$	$-\log 1/48$ 5.58
4x3	2	$1/16 \times 1/2$	$-\log 1/32$ 5.0
4x4	1	$1/16 \times 1/1$	$-\log 1/16$ 4.0

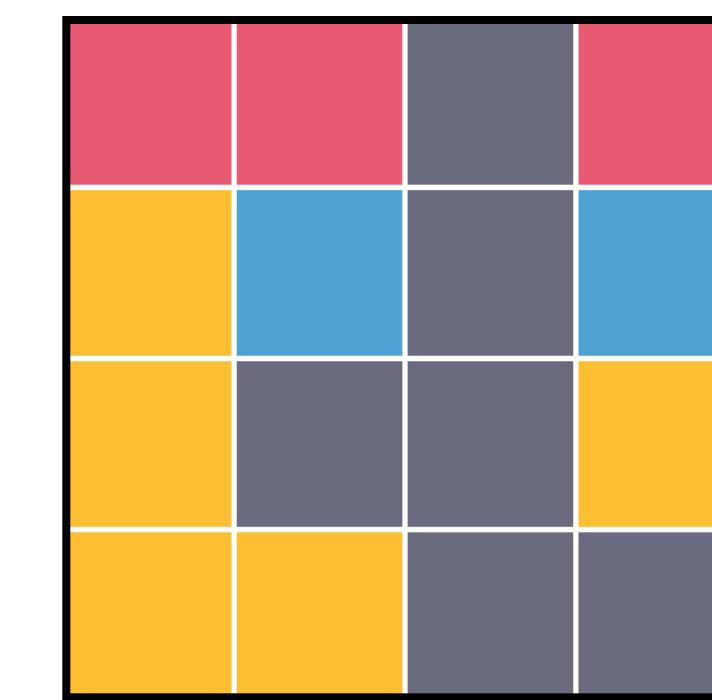
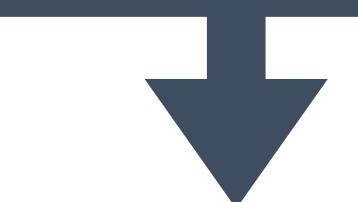
Uniformly sample a class



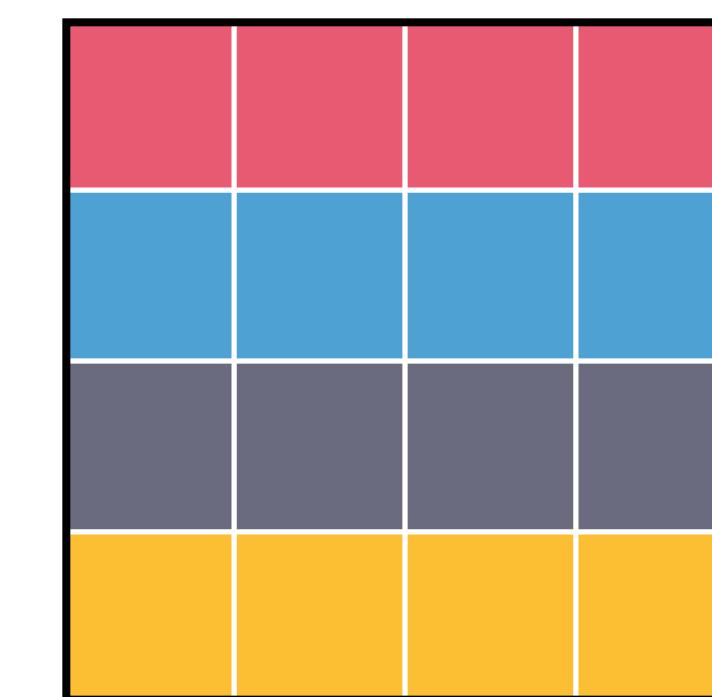
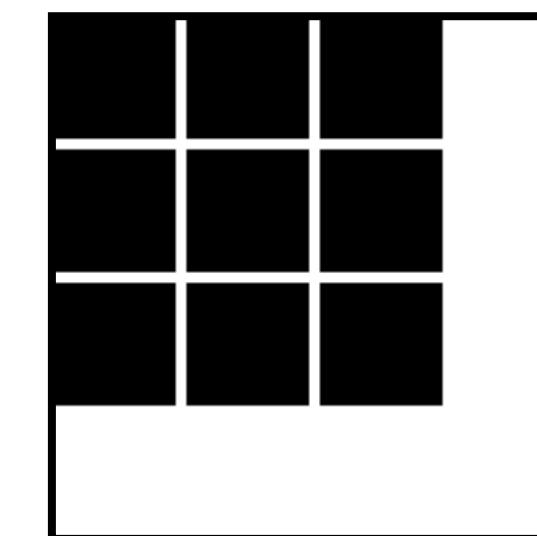
# Computing $DL(L)$ : The rectangle code

Class	Positions	Probability	Codelength (bits)
1x1	16	$1/16 \times 1/16$	$-\log 1/256$ 8.0
1x2	12	$1/16 \times 1/12$	$-\log 1/192$ 7.58
1x3	8	$1/16 \times 1/8$	$-\log 1/128$ 7.0
1x4	4	$1/16 \times 1/4$	$-\log 1/64$ 6.0
2x1	12	$1/16 \times 1/12$	$-\log 1/192$ 7.58
2x2	9	$1/16 \times 1/9$	$-\log 1/144$ 7.17
2x3	6	$1/16 \times 1/6$	$-\log 1/96$ 6.58
2x4	3	$1/16 \times 1/3$	$-\log 1/48$ 5.58
3x1	8	$1/16 \times 1/8$	$-\log 1/128$ 7.0
3x2	6	$1/16 \times 1/6$	$-\log 1/96$ 6.58
3x3	4	$1/16 \times 1/4$	$-\log 1/64$ 6.0
3x4	2	$1/16 \times 1/2$	$-\log 1/32$ 5.0
4x1	4	$1/16 \times 1/4$	$-\log 1/64$ 6.0
4x2	3	$1/16 \times 1/3$	$-\log 1/48$ 5.58
4x3	2	$1/16 \times 1/2$	$-\log 1/32$ 5.0
4x4	1	$1/16 \times 1/1$	$-\log 1/16$ 4.0

Uniformly sample a class

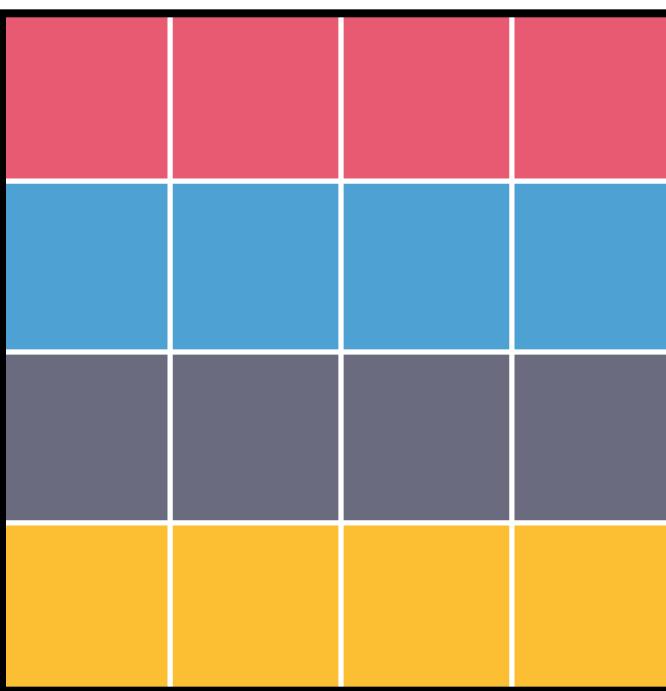
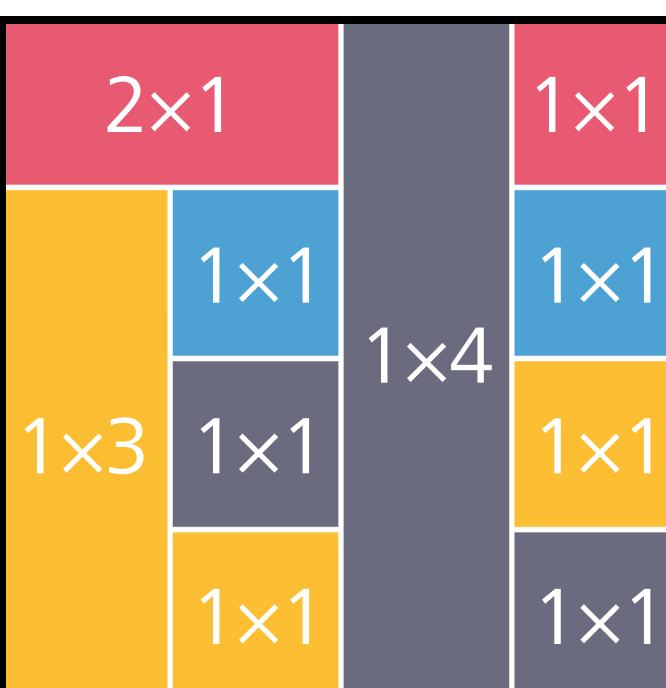
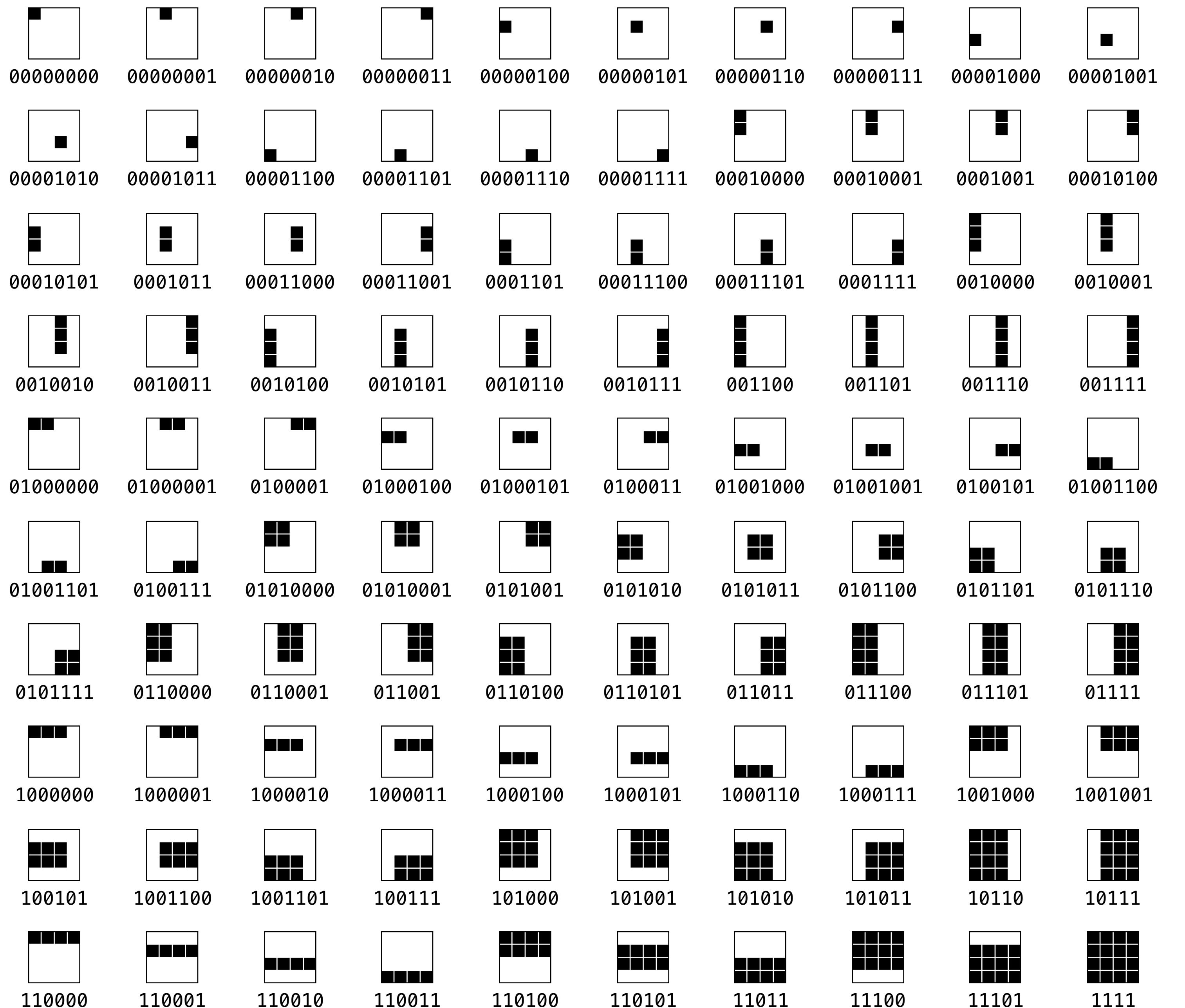


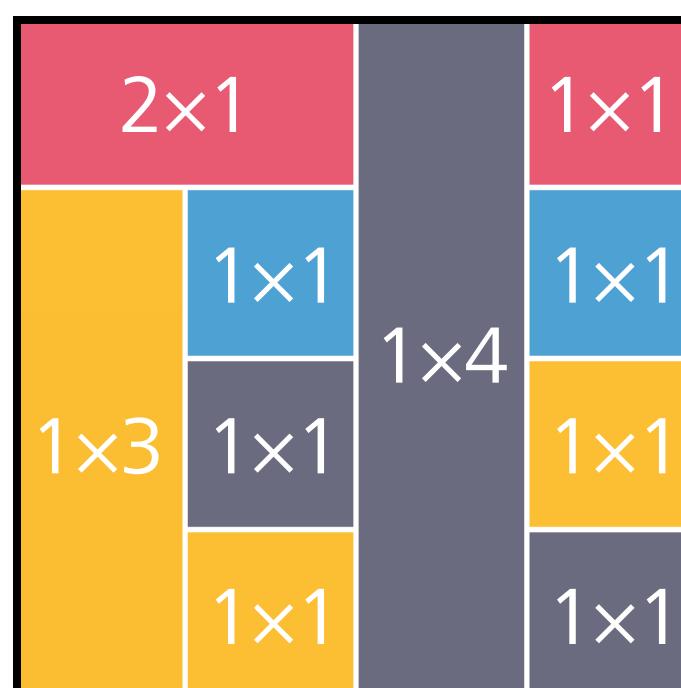
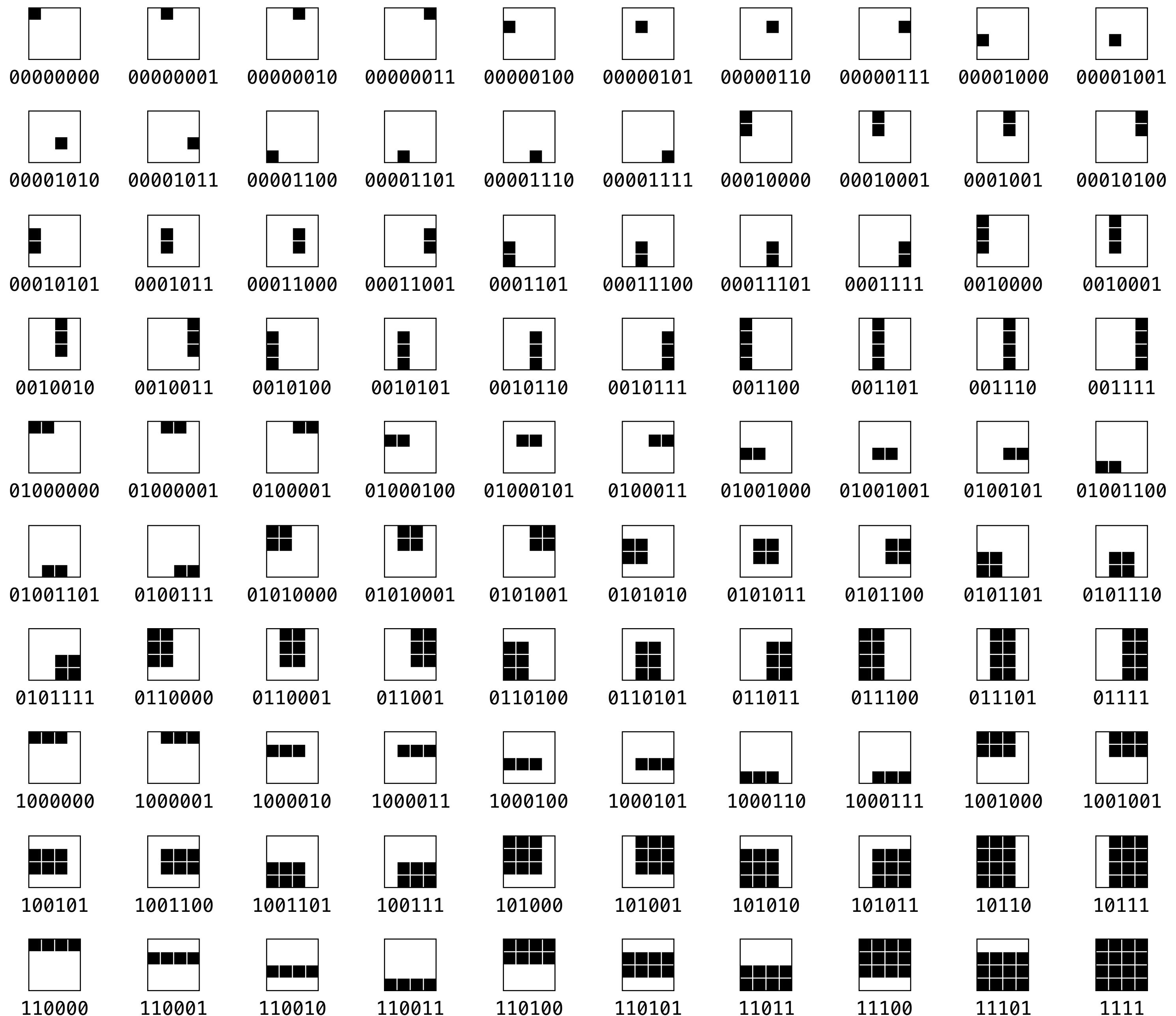
Uniformly sample a position



 00000000	 00000001	 00000010	 00000011	 00000100	 00000101	 00000110	 00000111	 00001000	 00001001
 00001010	 00001011	 00001100	 00001101	 00001110	 00001111	 00010000	 00010001	 00010010	 00010100
 00010101	 00010111	 00011000	 00011001	 00011011	 00011100	 00011101	 00011111	 00100000	 00100001
 00100100	 00100111	 00101000	 00101011	 00101100	 00101111	 00110000	 00110011	 00110100	 00110111
 01000000	 01000001	 01000001	 01000100	 01000101	 01000111	 01001000	 01001001	 01001010	 01001100
 01001101	 01001111	 01010000	 01010001	 01010011	 01010100	 01010111	 01011000	 01011011	 01011100
 01011111	 01100000	 01100001	 01100011	 01101000	 01101011	 01101111	 01110000	 01110011	 01111111
 10000000	 10000001	 10000010	 10000011	 10000100	 10000101	 10000110	 10000111	 10001000	 10001001
 10010101	 10011000	 10011011	 10011111	 10100000	 10100011	 10100100	 10100111	 10101100	 10111111
 11000000	 11000001	 11000010	 11000011	 11010000	 11010011	 11011000	 11011011	 11100000	 11100011



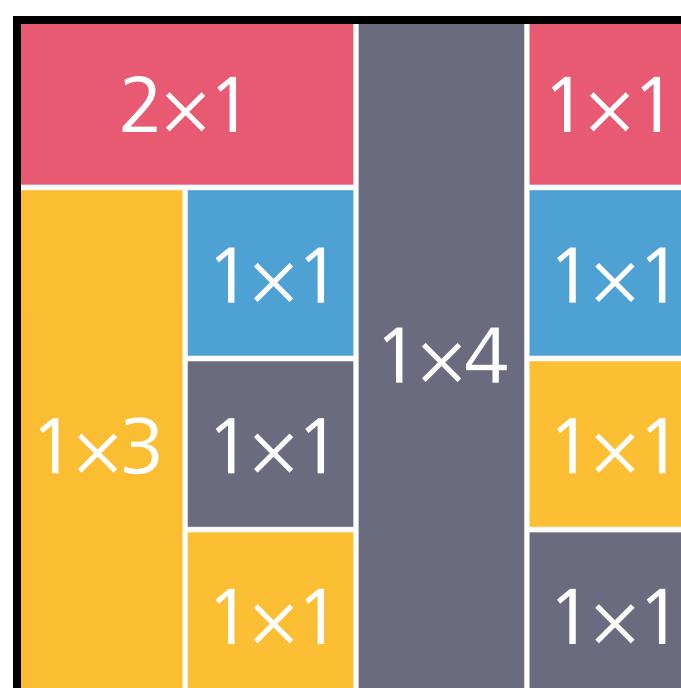
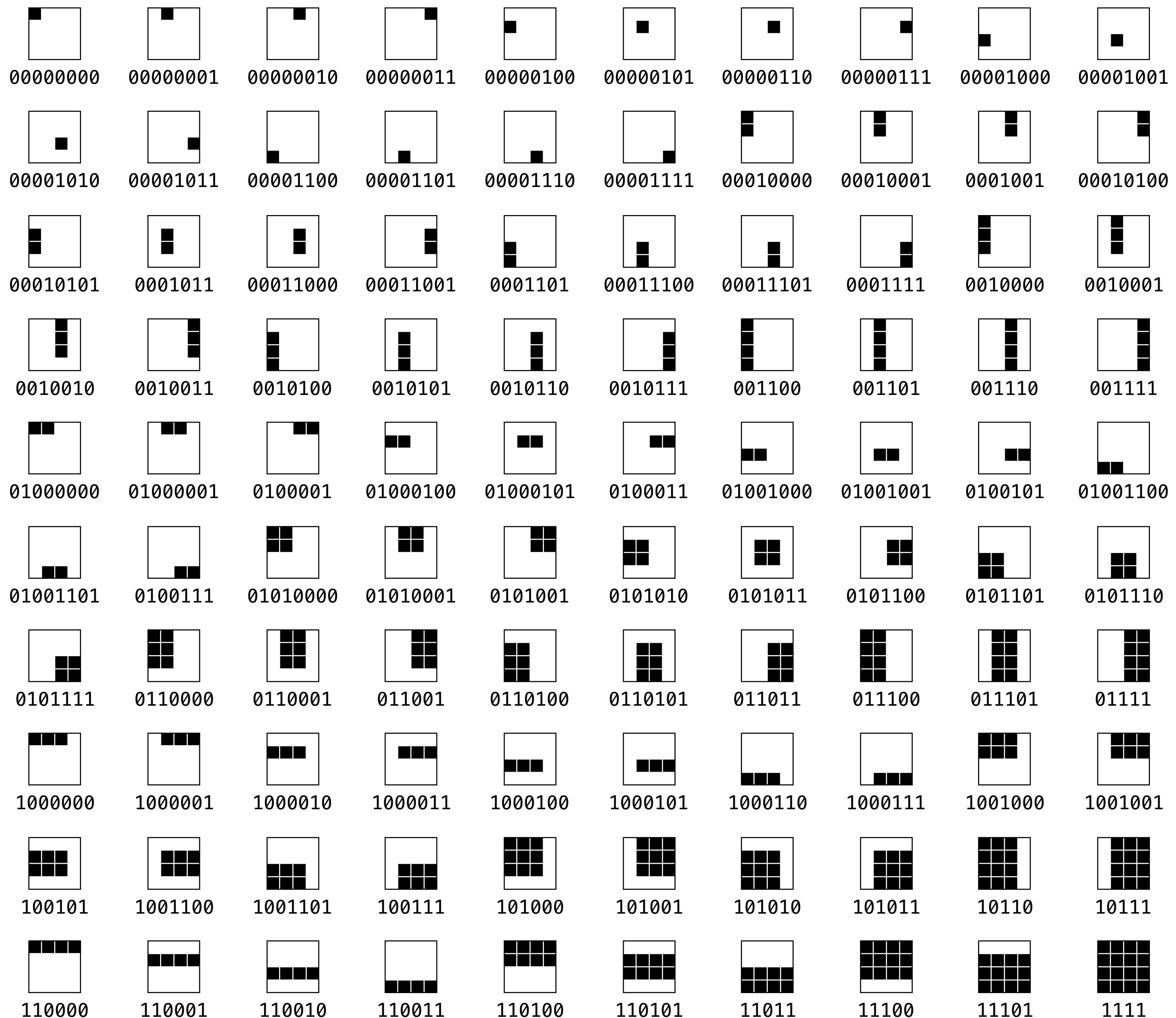




010000000000011  
 0000010100000111  
 001110000010010001111  
 0010100000101100001101

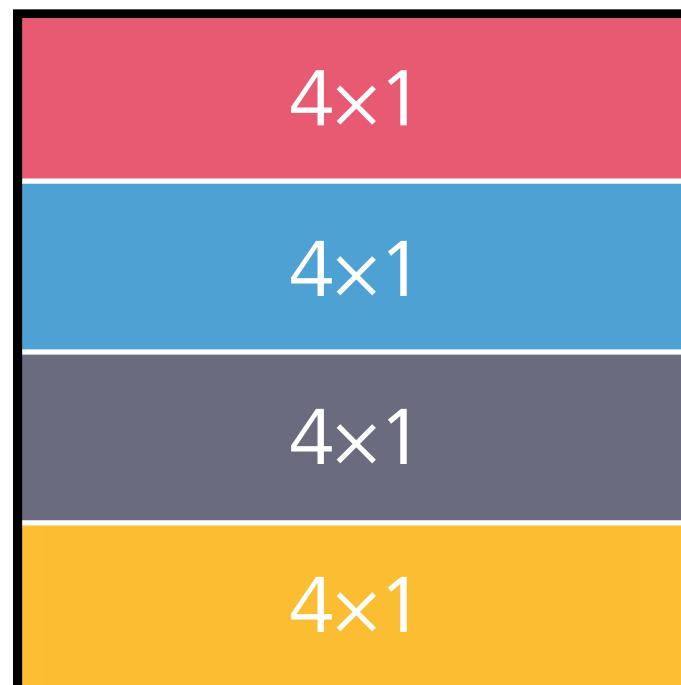
76.58 bits



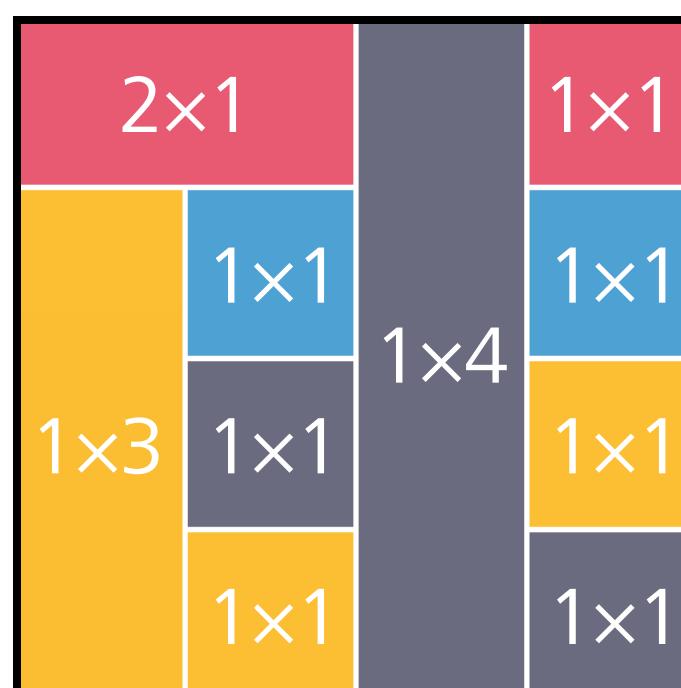


010000000000011  
0000010100000111  
001110000100100001111  
0010100000101100001101

76.58 bits

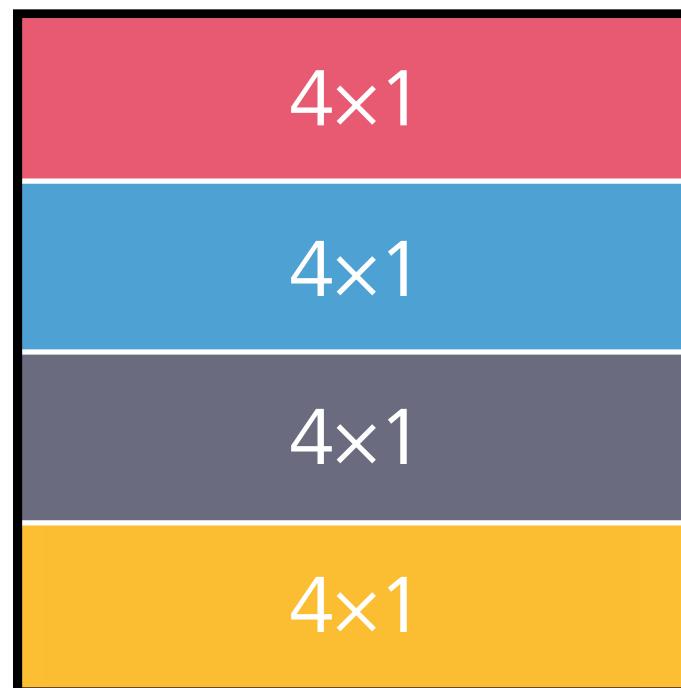


00000000	00000001	00000010	00000011	00000100	00000101	00000110	00000111	00001000	00001001
00001010	00001011	00001100	00001101	00001110	00001111	00010000	00010001	0001001	00010100
00010101	0001011	00011000	00011001	0001101	00011100	00011101	0001111	0010000	0010001
0010010	0010011	0010100	0010101	0010110	0010111	001100	001101	001110	001111
01000000	01000001	0100001	01000100	01000101	0100011	01001000	01001001	0100101	01001100
01001101	0100111	01010000	01010001	0101001	0101010	0101011	0101100	0101101	0101110
0101111	0110000	0110001	011001	0110100	0110101	011011	011100	011101	01111
1000000	1000001	1000010	1000011	1000100	1000101	1000110	1000111	1001000	1001001
100101	1001100	1001101	100111	101000	101001	101010	101011	10110	10111
110000	110001	110010	110011	110100	110101	11011	11100	11101	1111



010000000000011  
 0000010100000111  
 001110000010010001111  
 0010100000101100001101

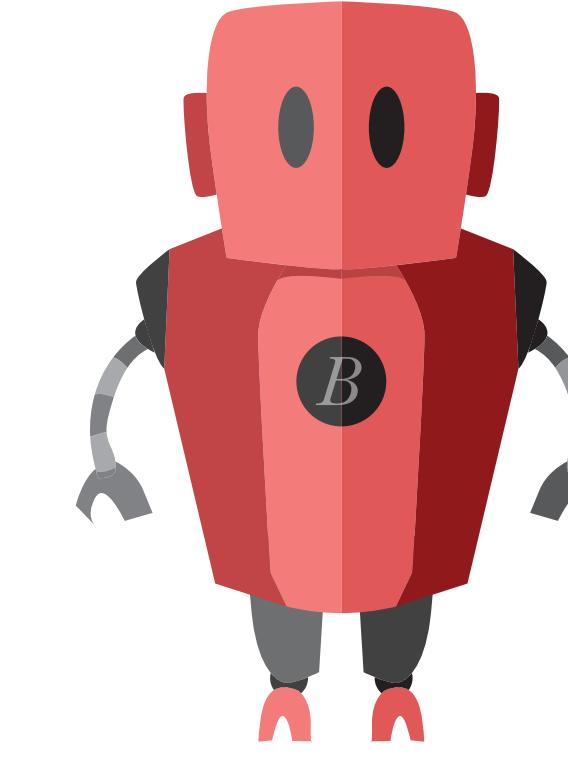
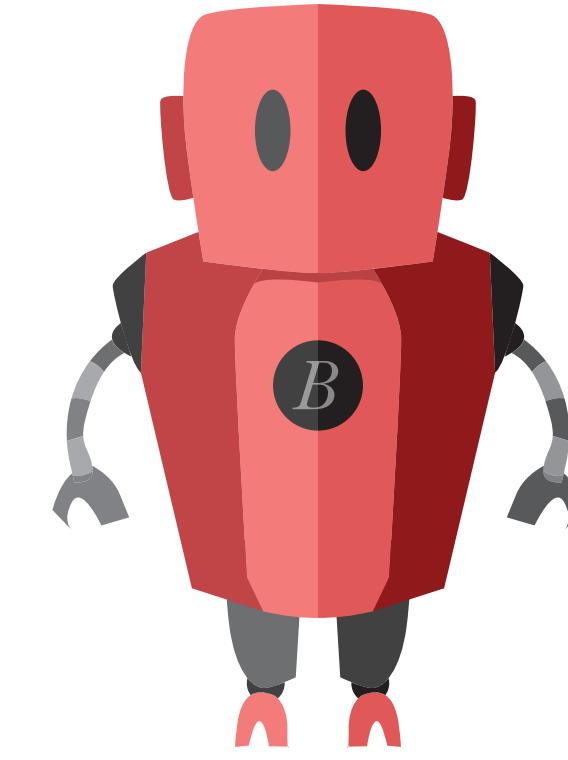
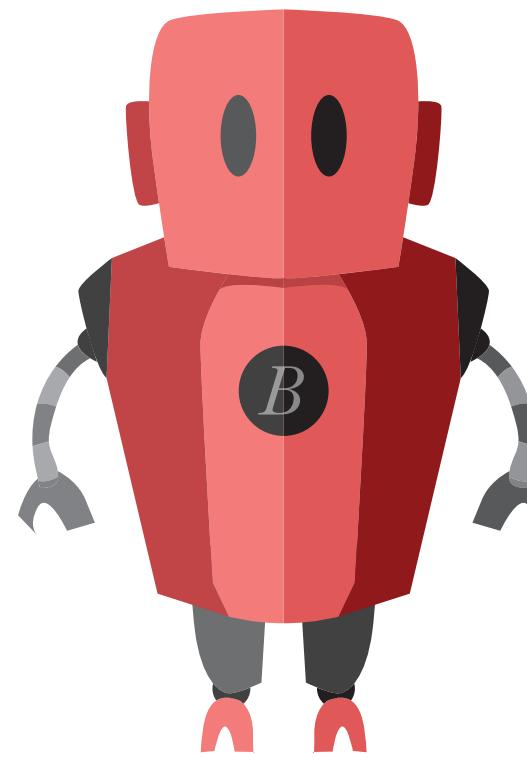
76.58 bits



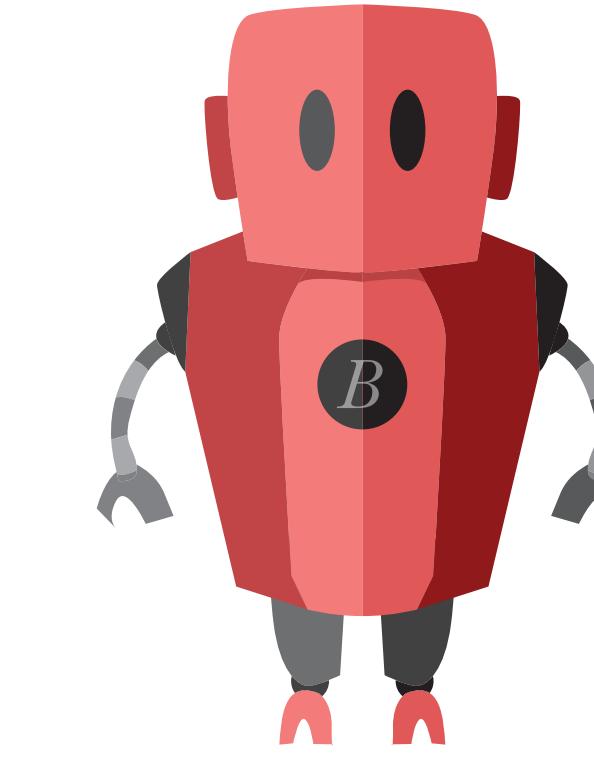
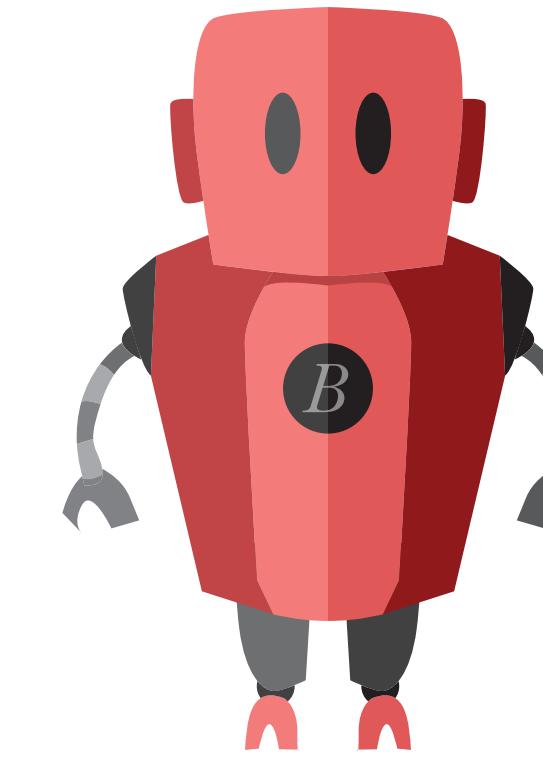
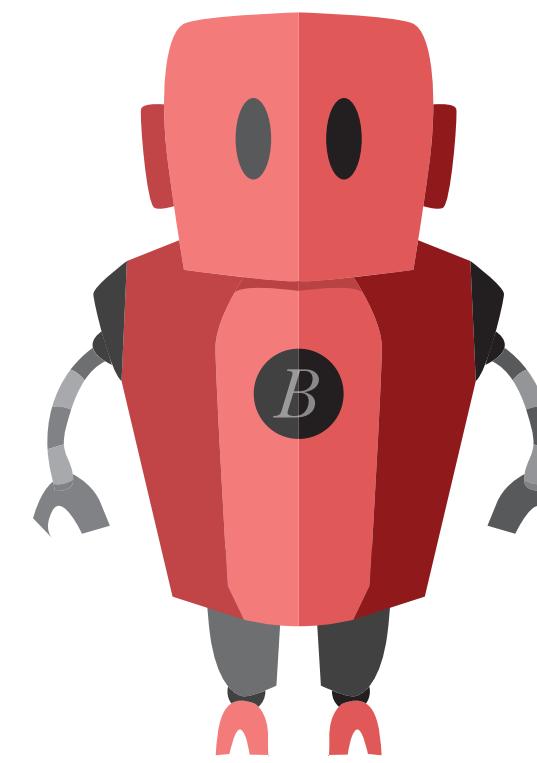
110000  
 110001  
 110010  
 110011

24 bits

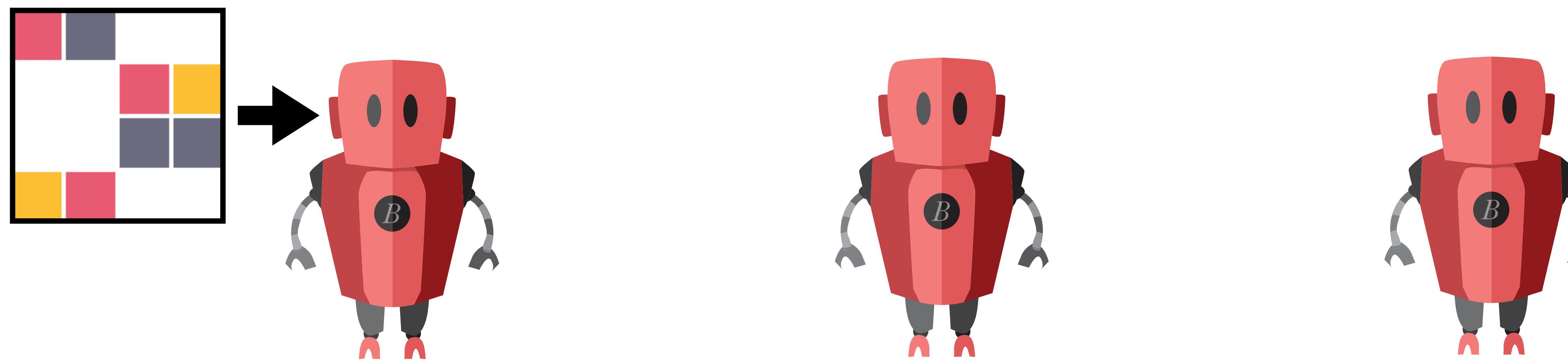
# Bayesian iterated learning



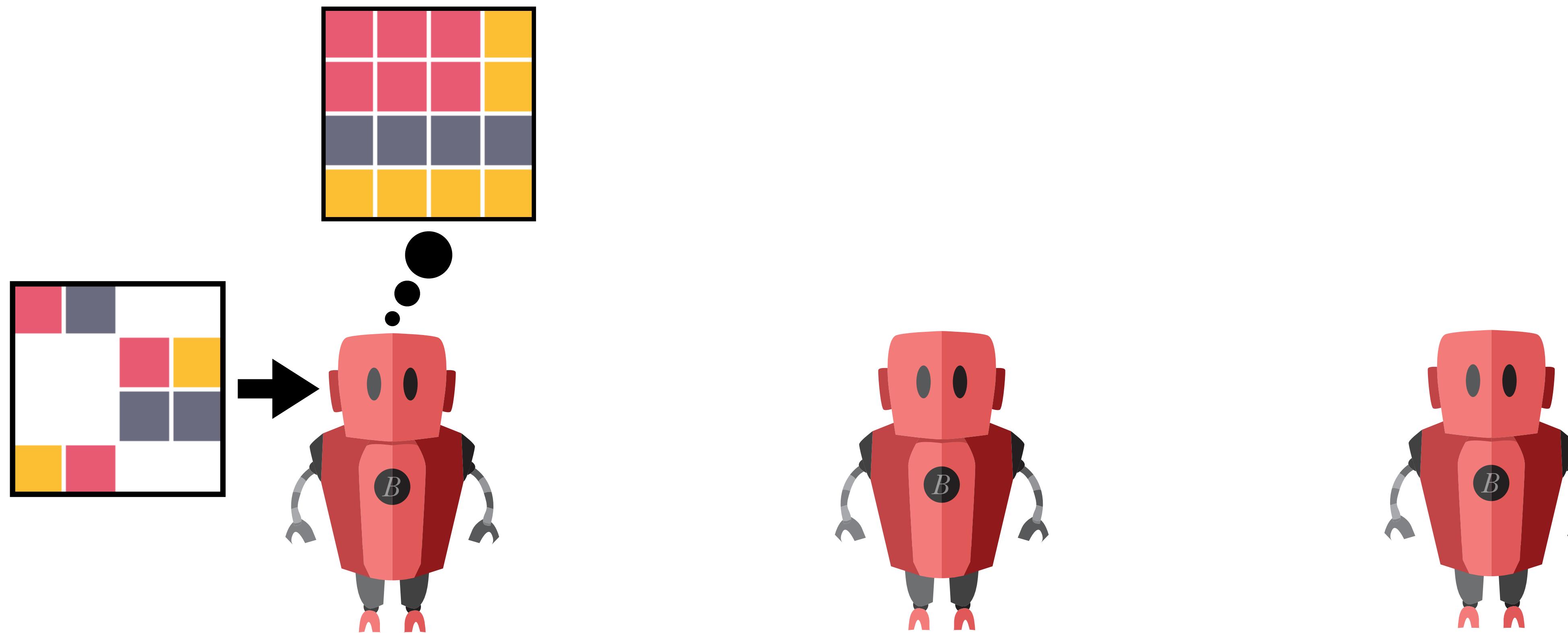
# Bayesian iterated learning



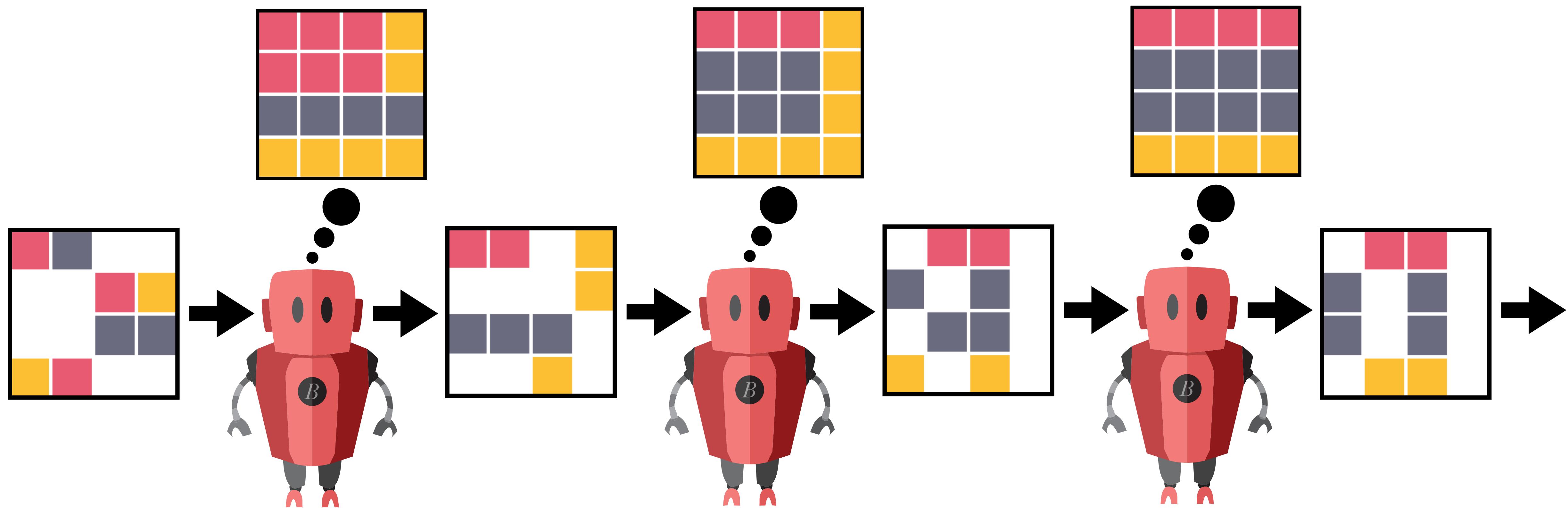
# Bayesian iterated learning

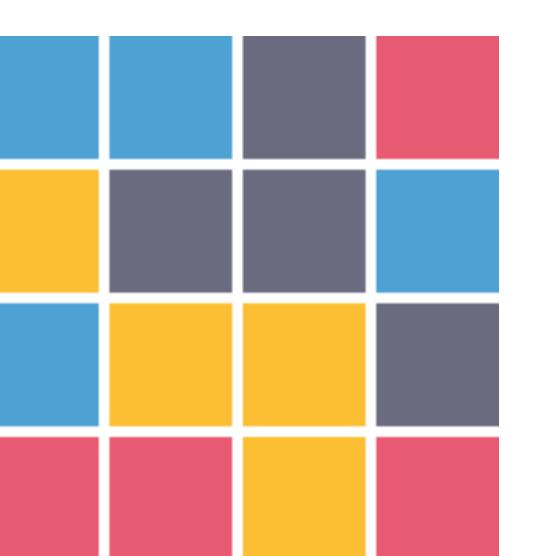
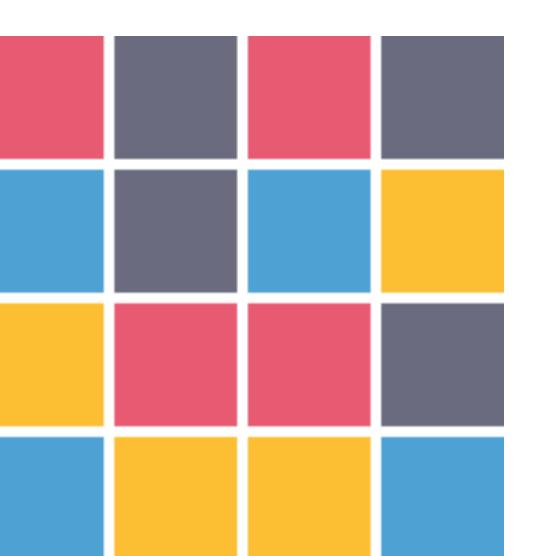
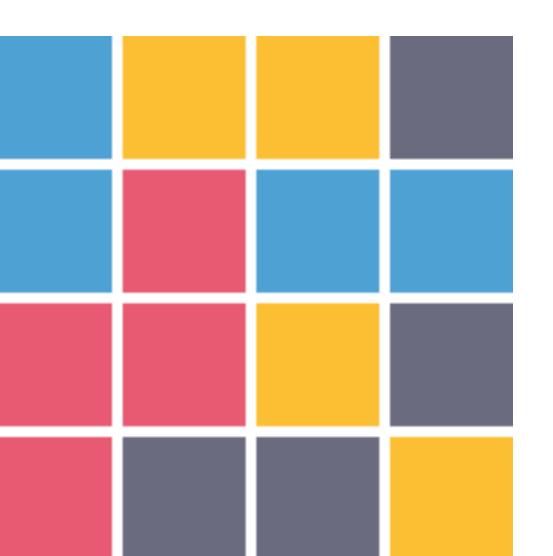
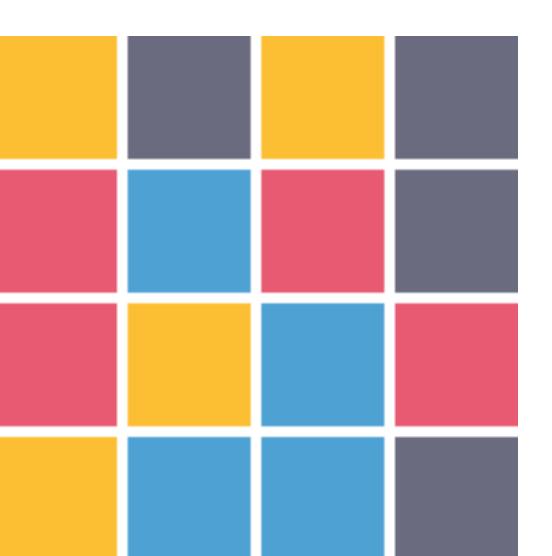
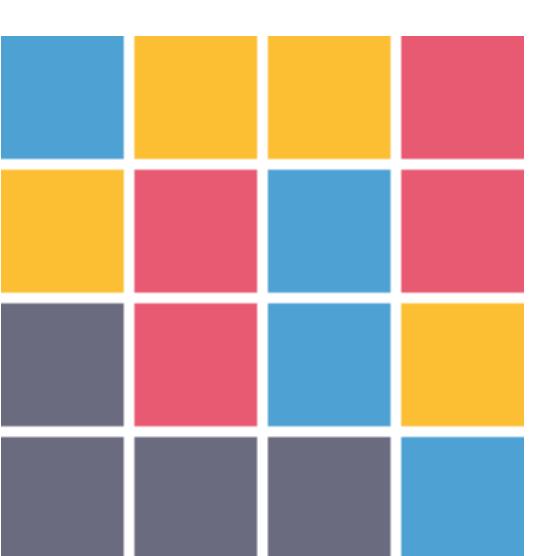
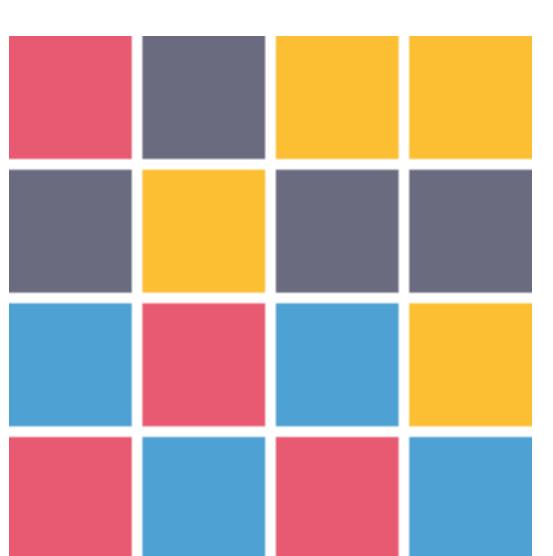
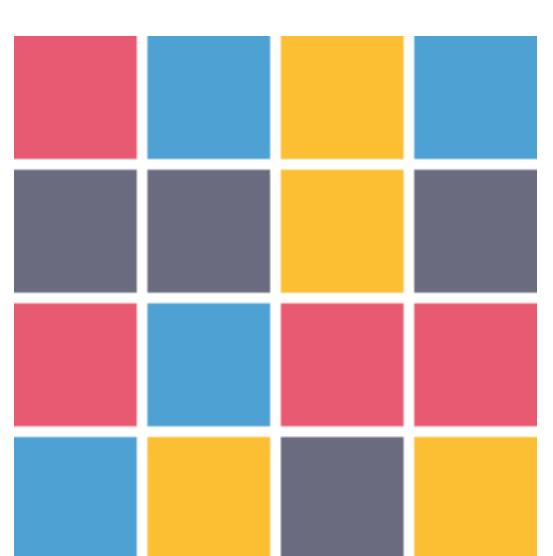
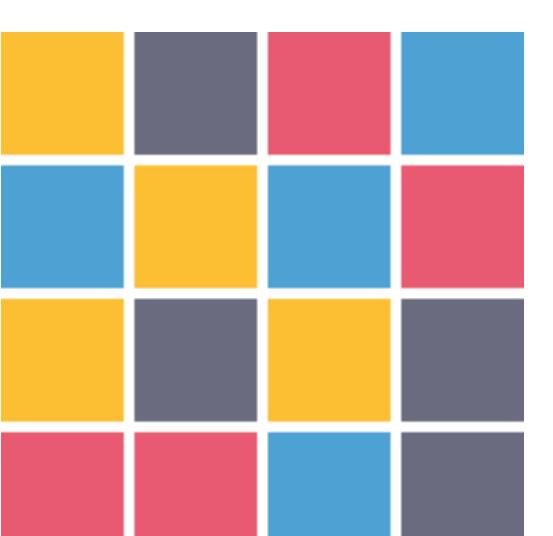
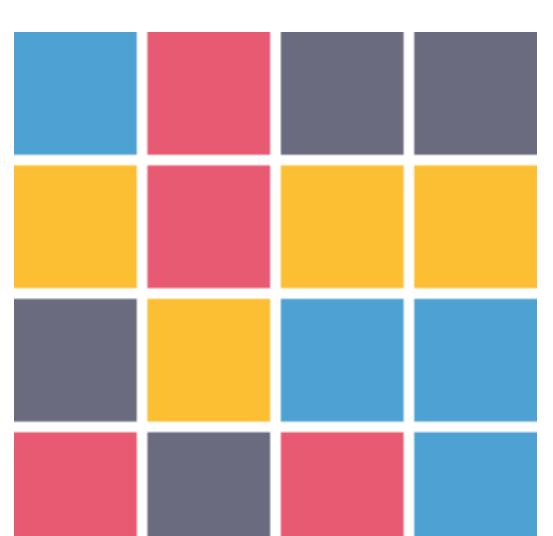
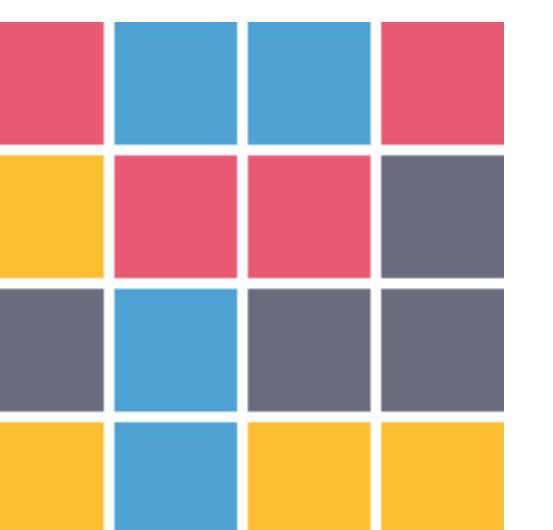
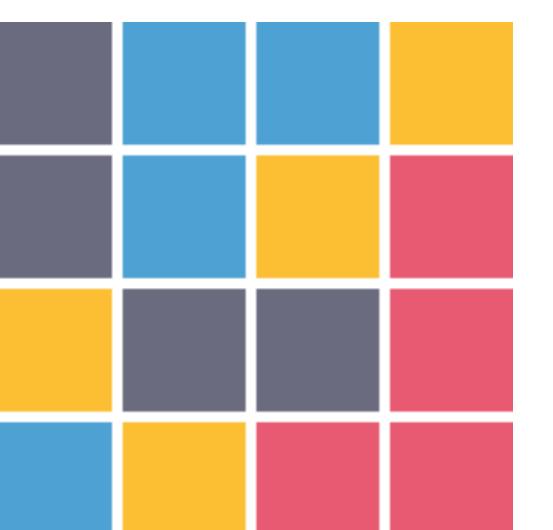
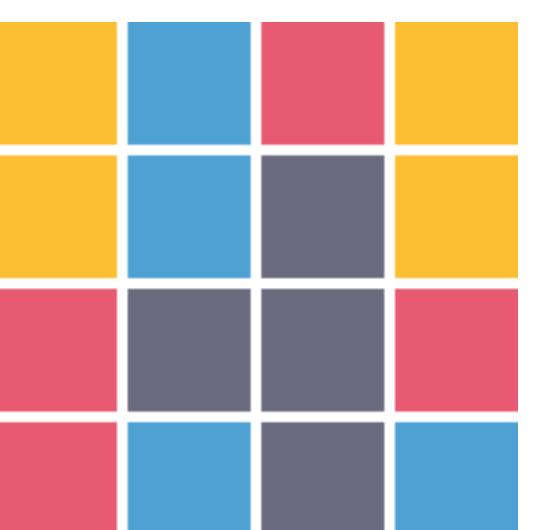
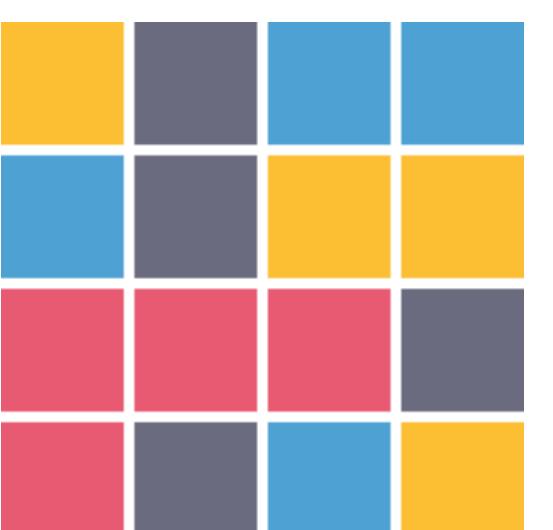
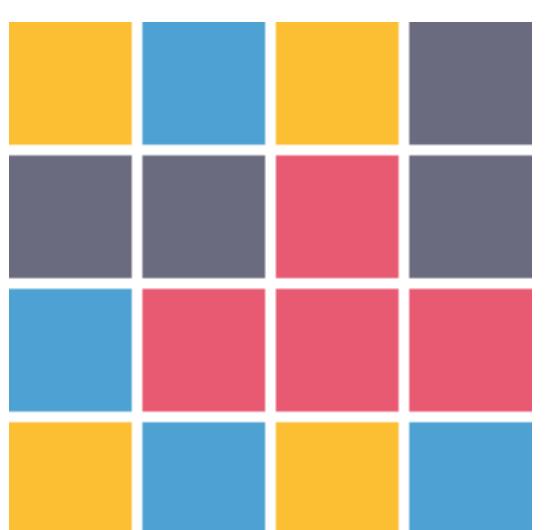
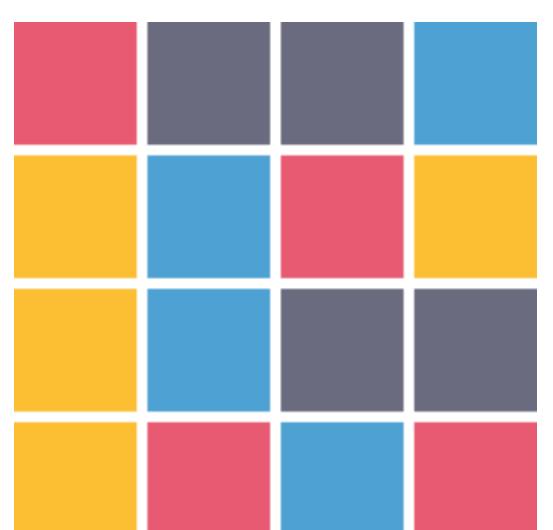
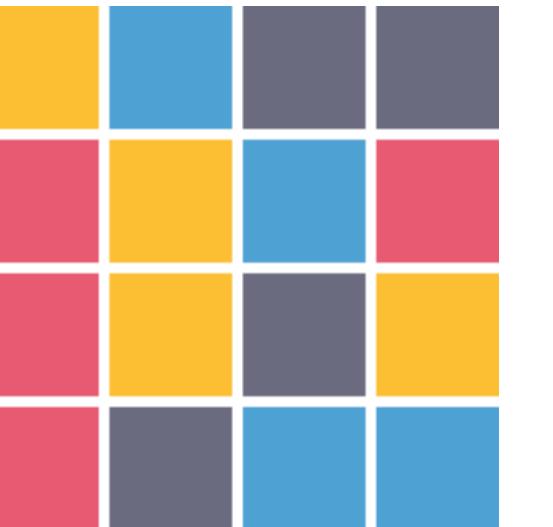
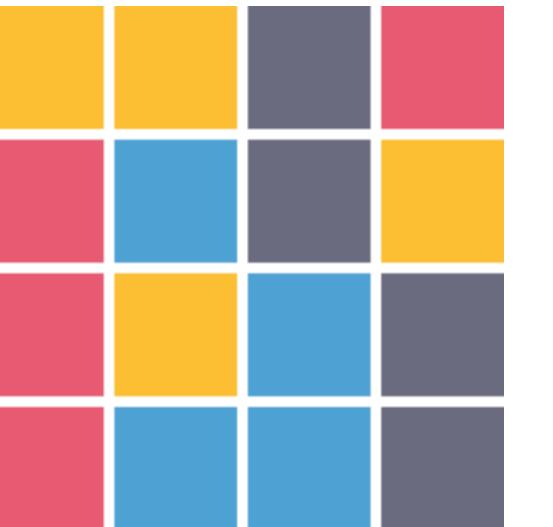
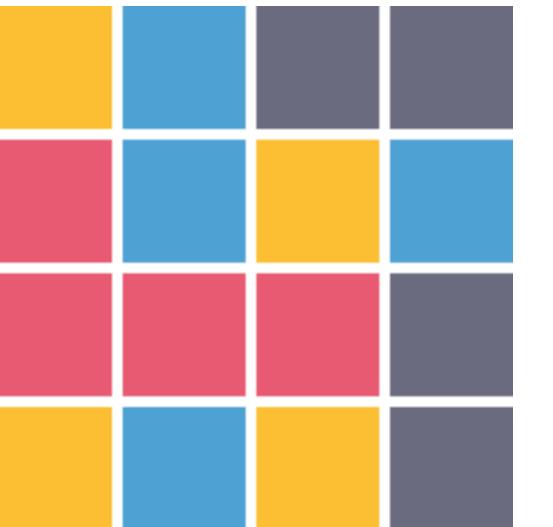
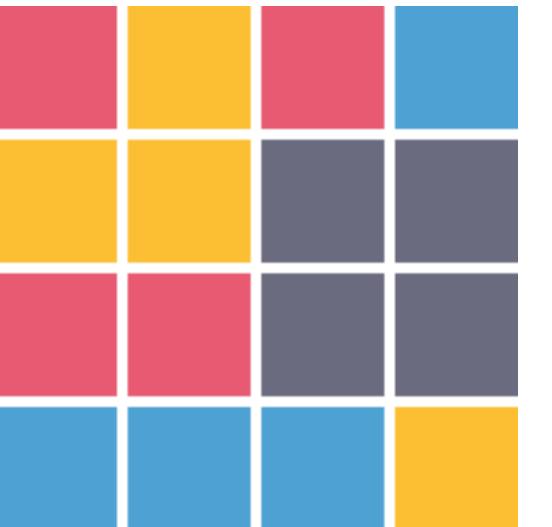
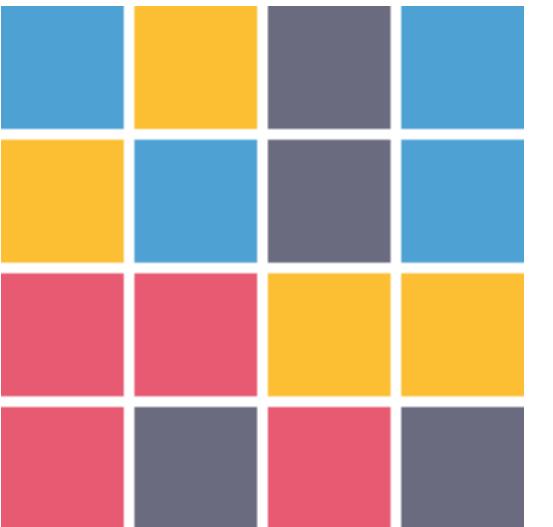
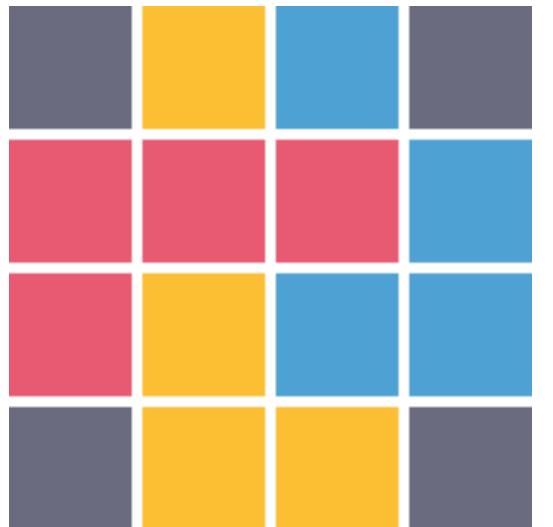
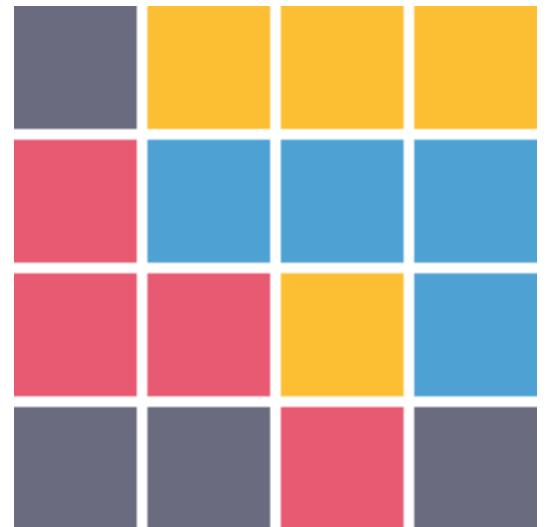


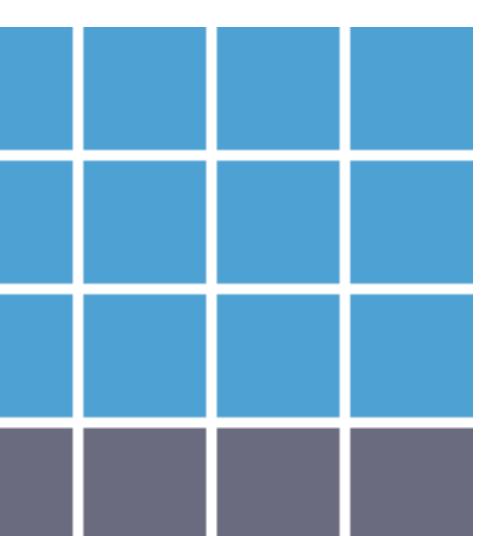
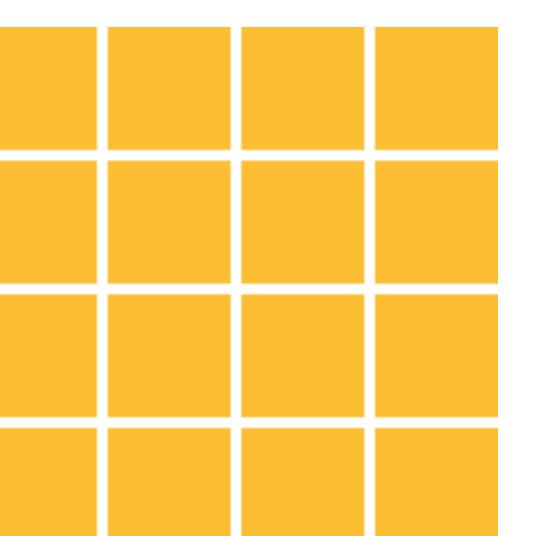
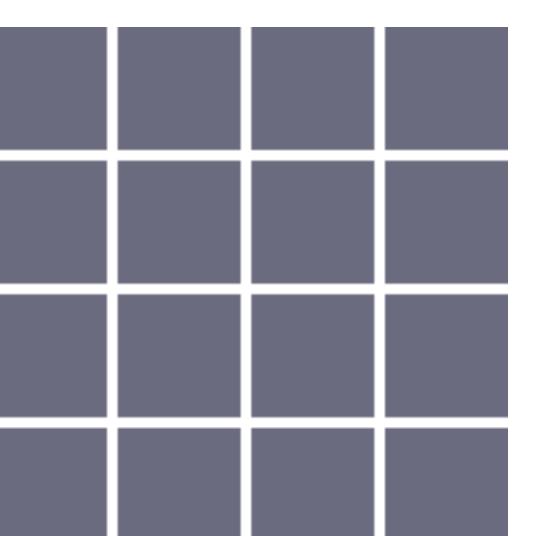
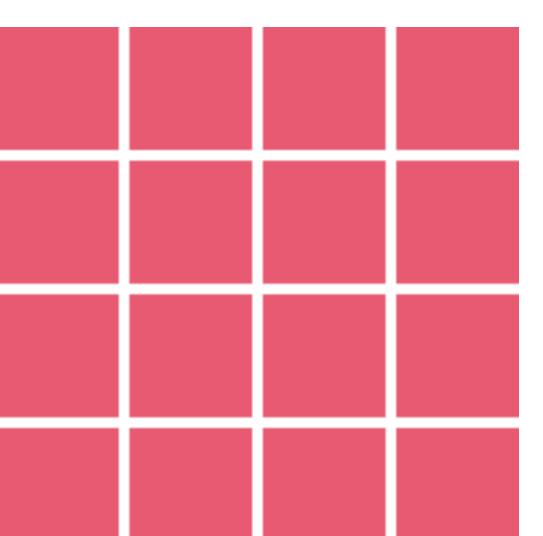
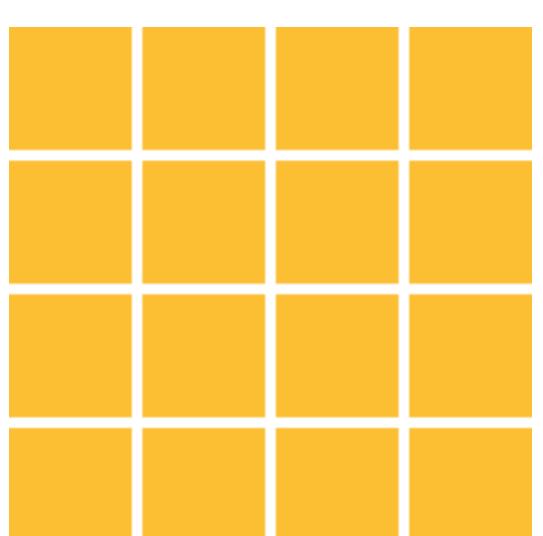
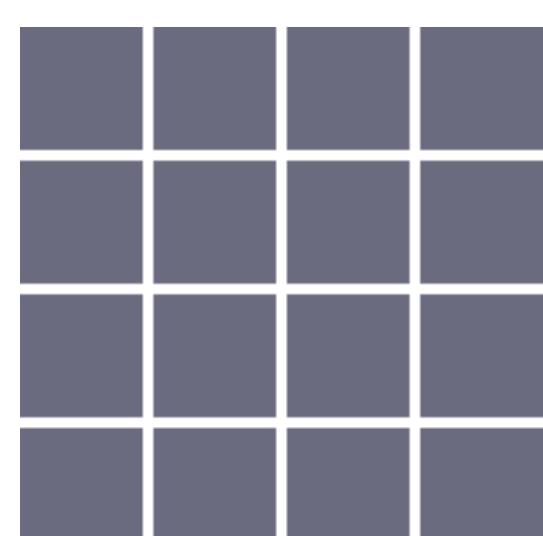
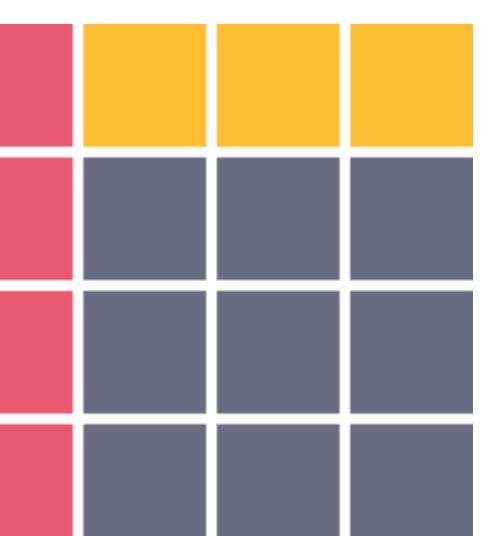
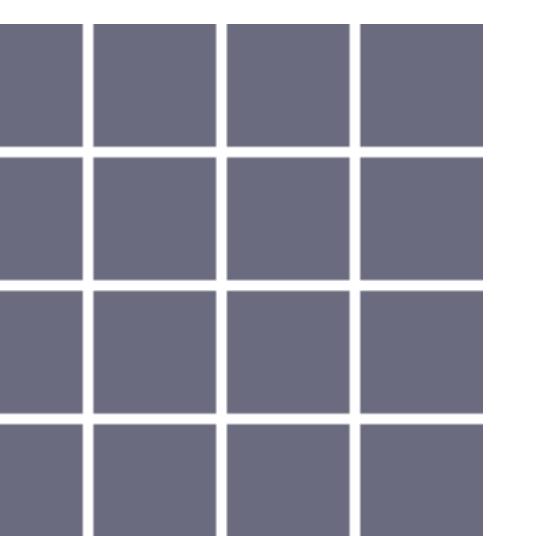
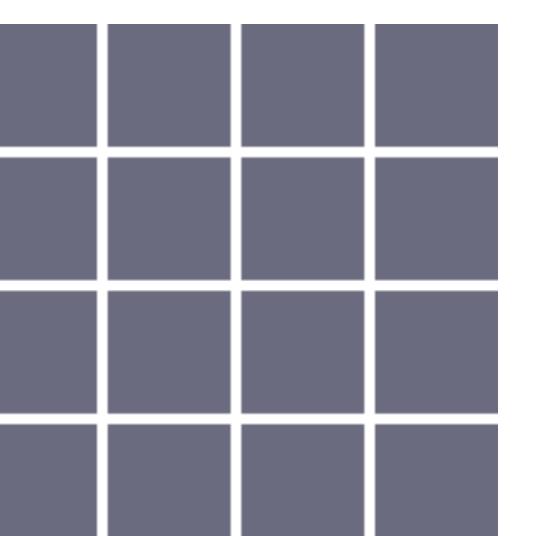
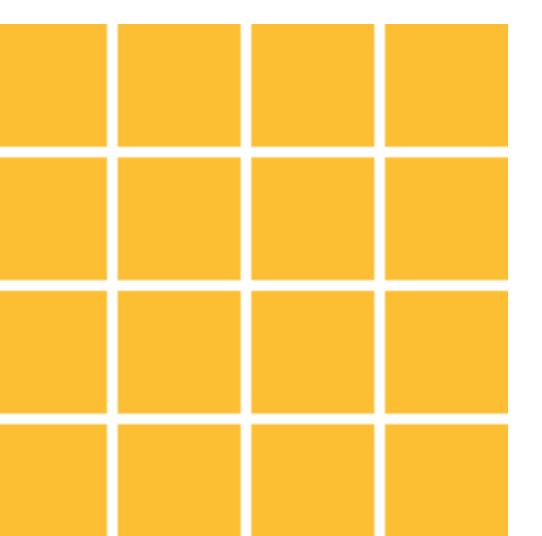
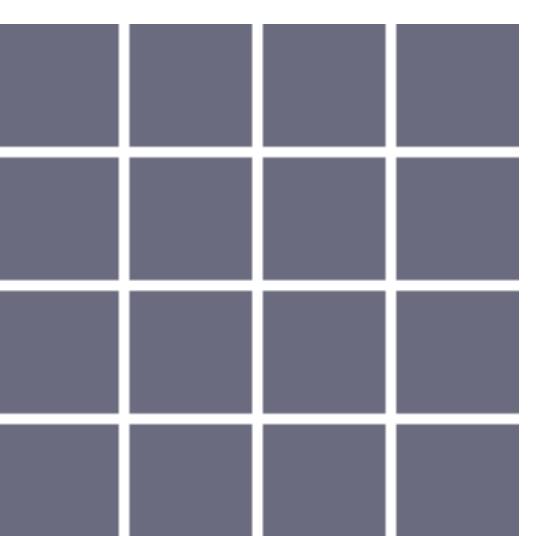
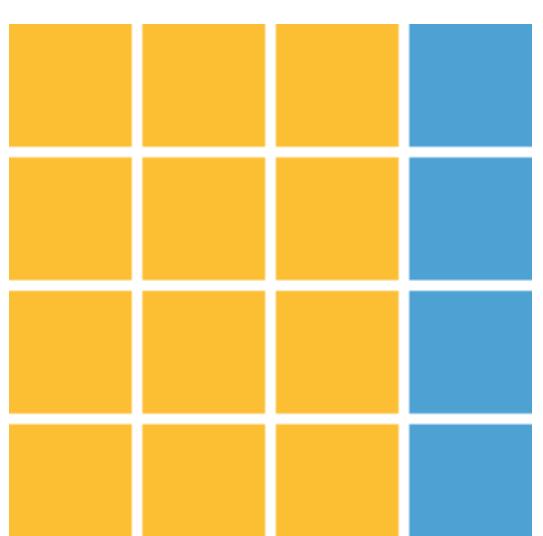
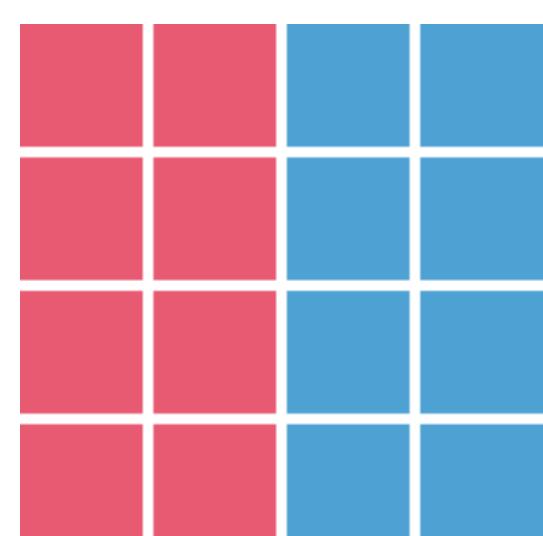
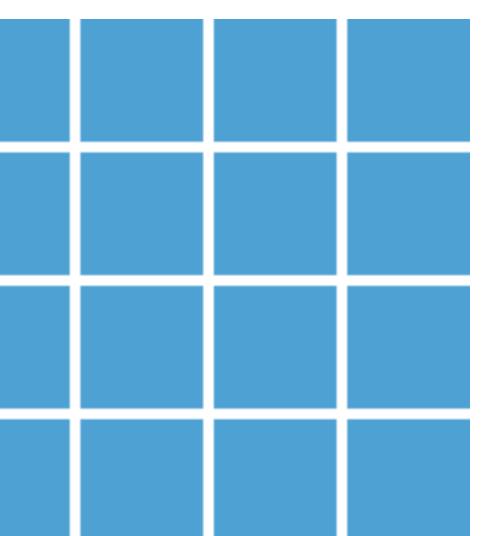
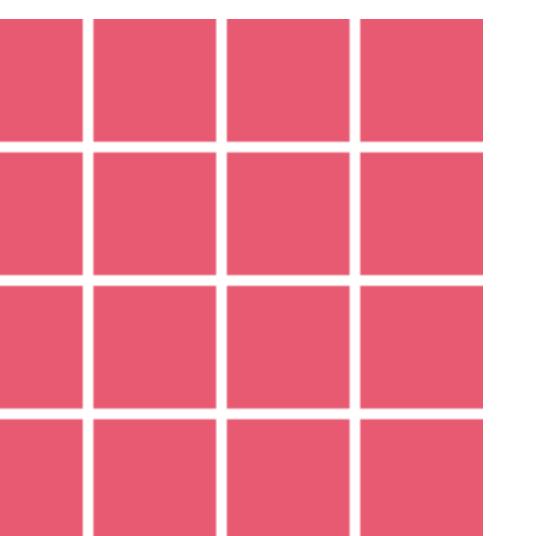
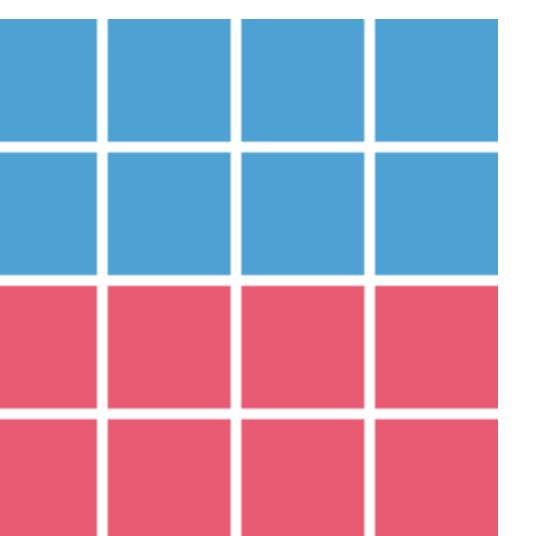
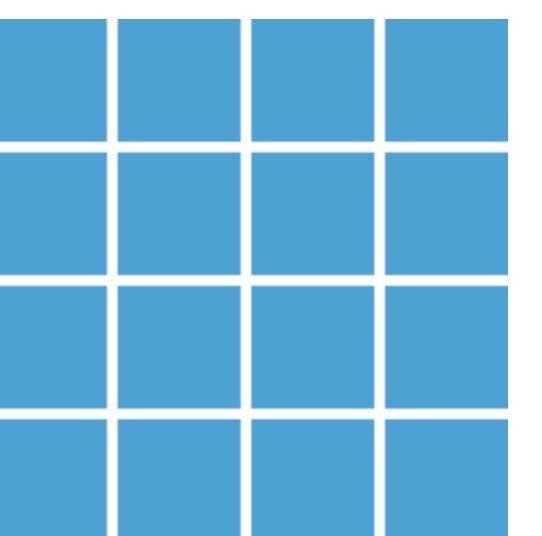
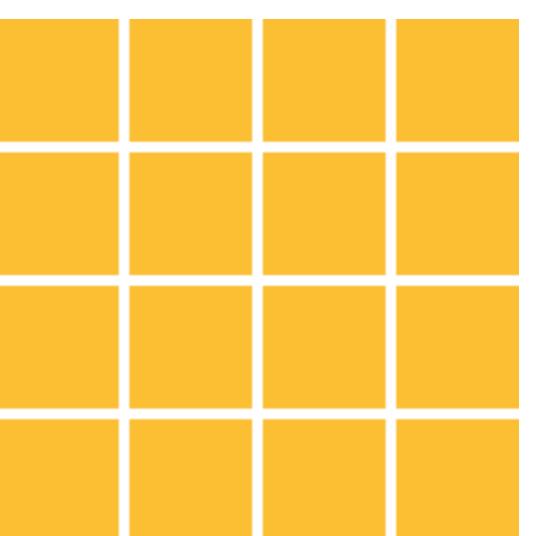
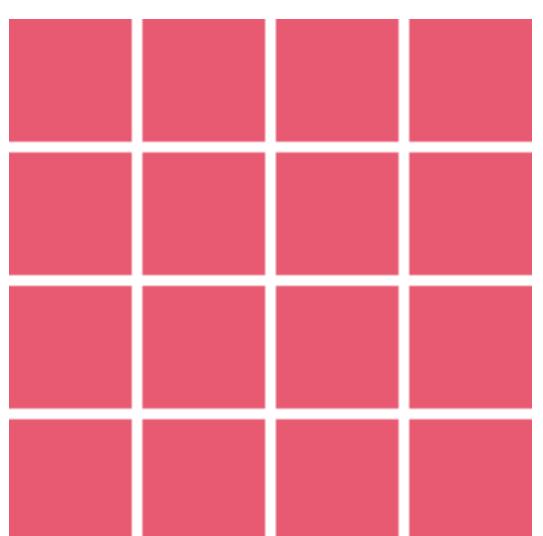
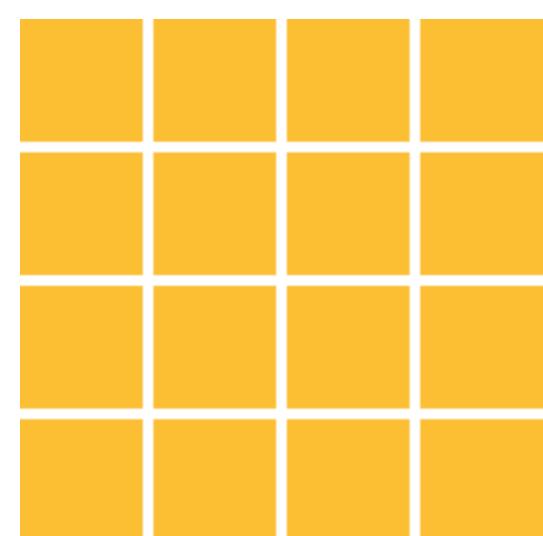
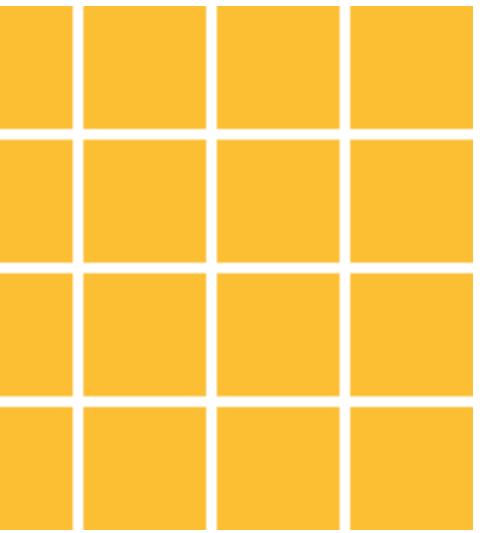
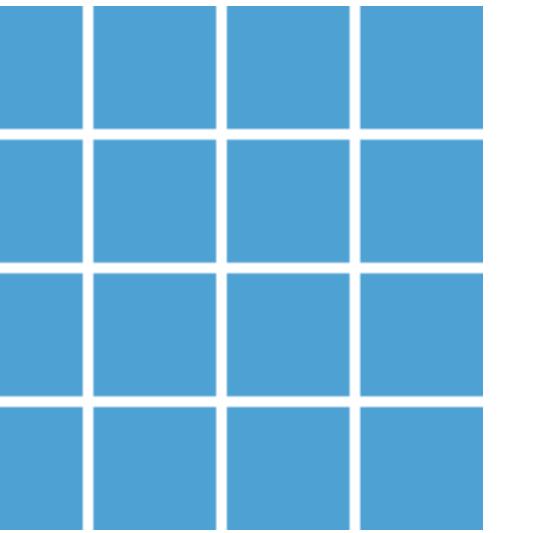
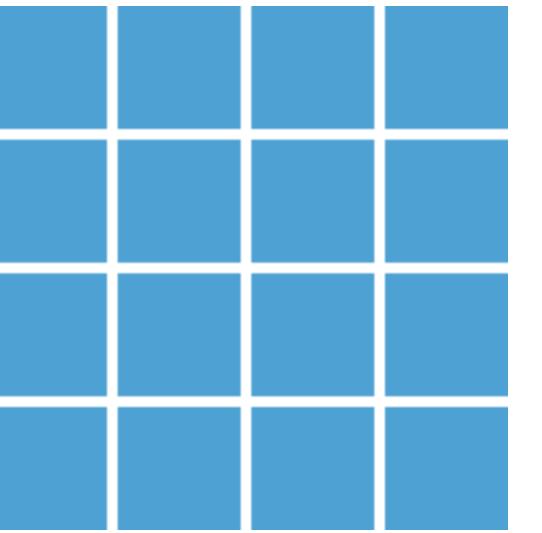
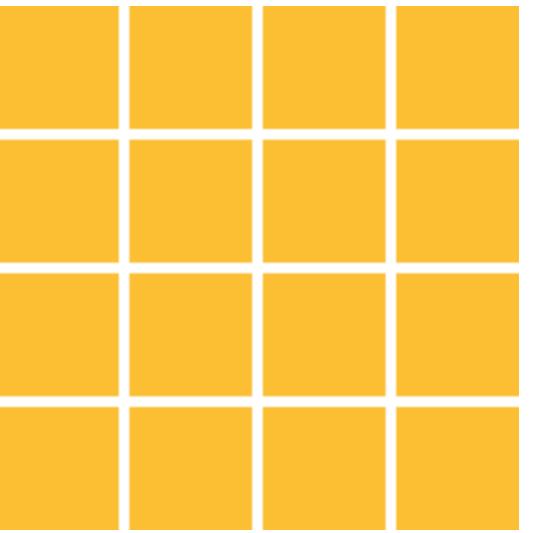
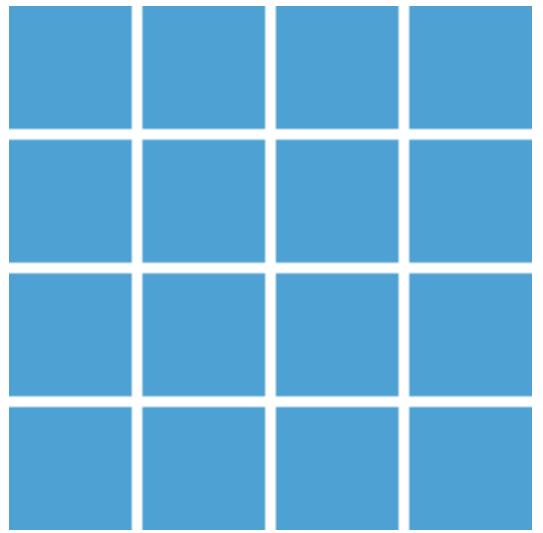
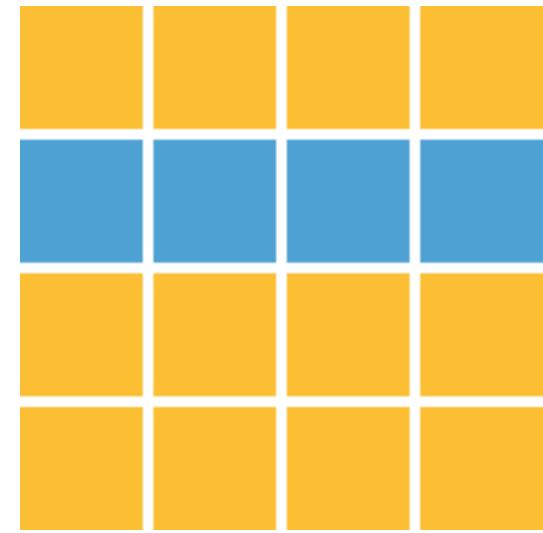
# Bayesian iterated learning



# Bayesian iterated learning

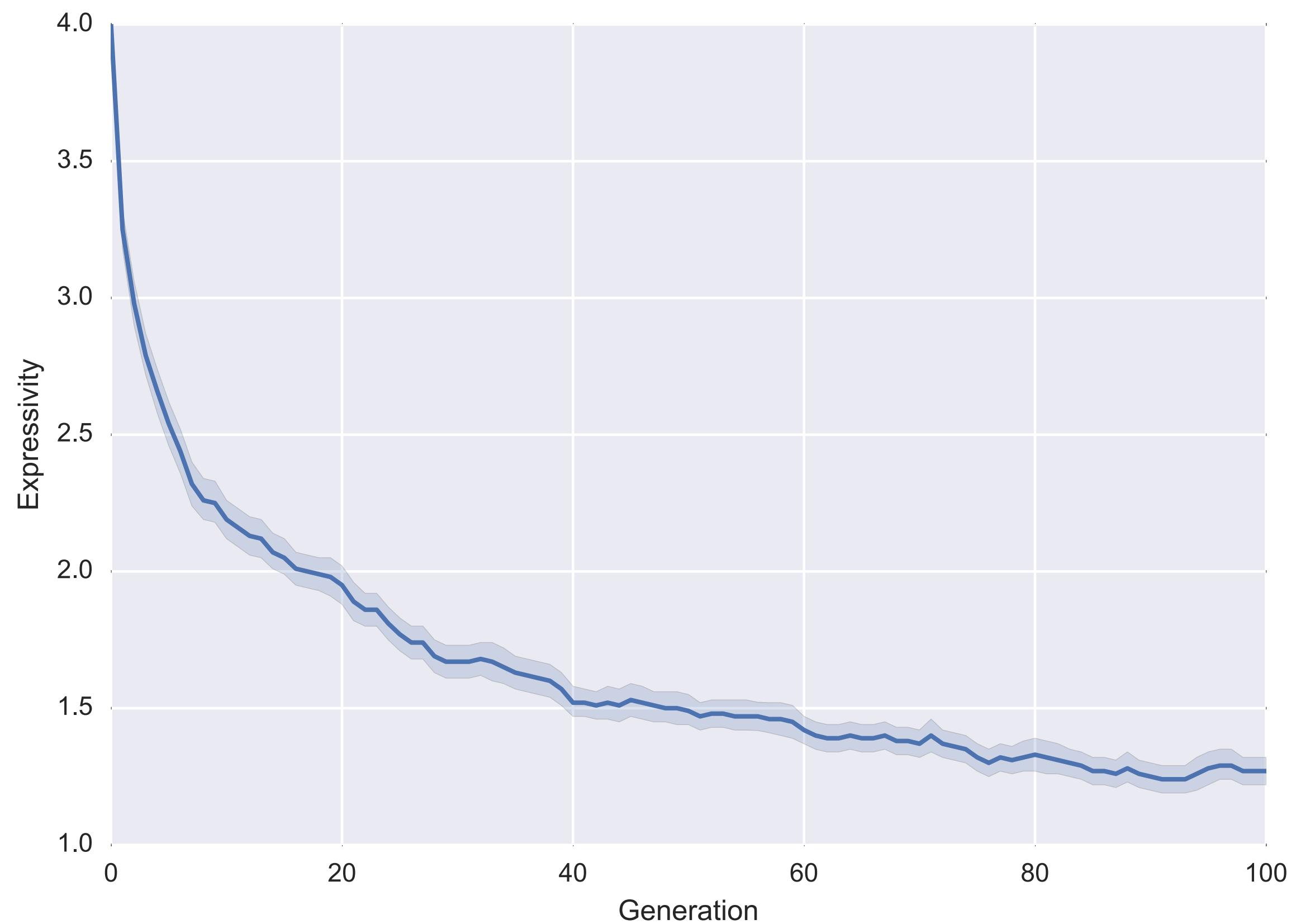




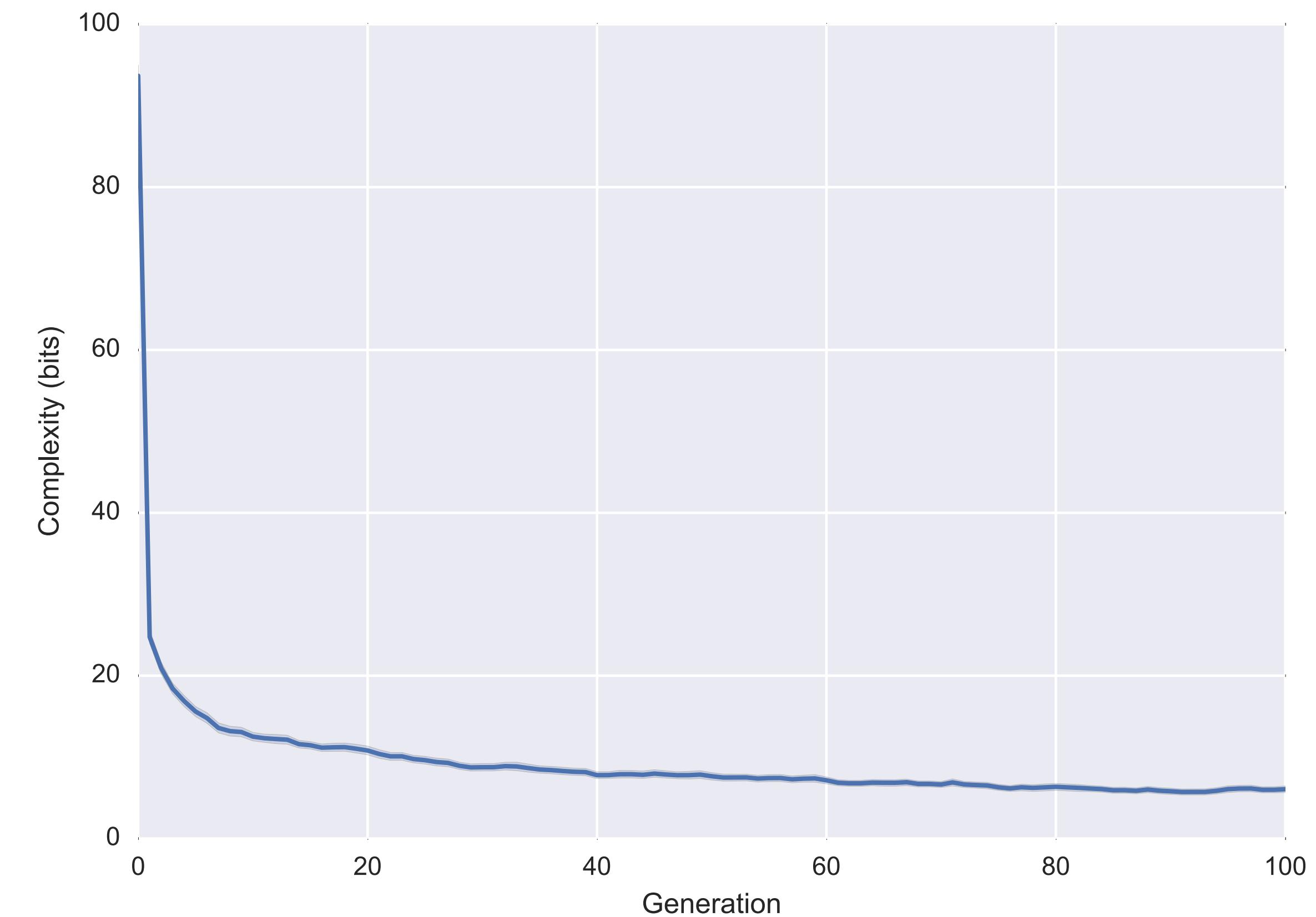


# Iterated learning converges to the prior

*Expressivity*

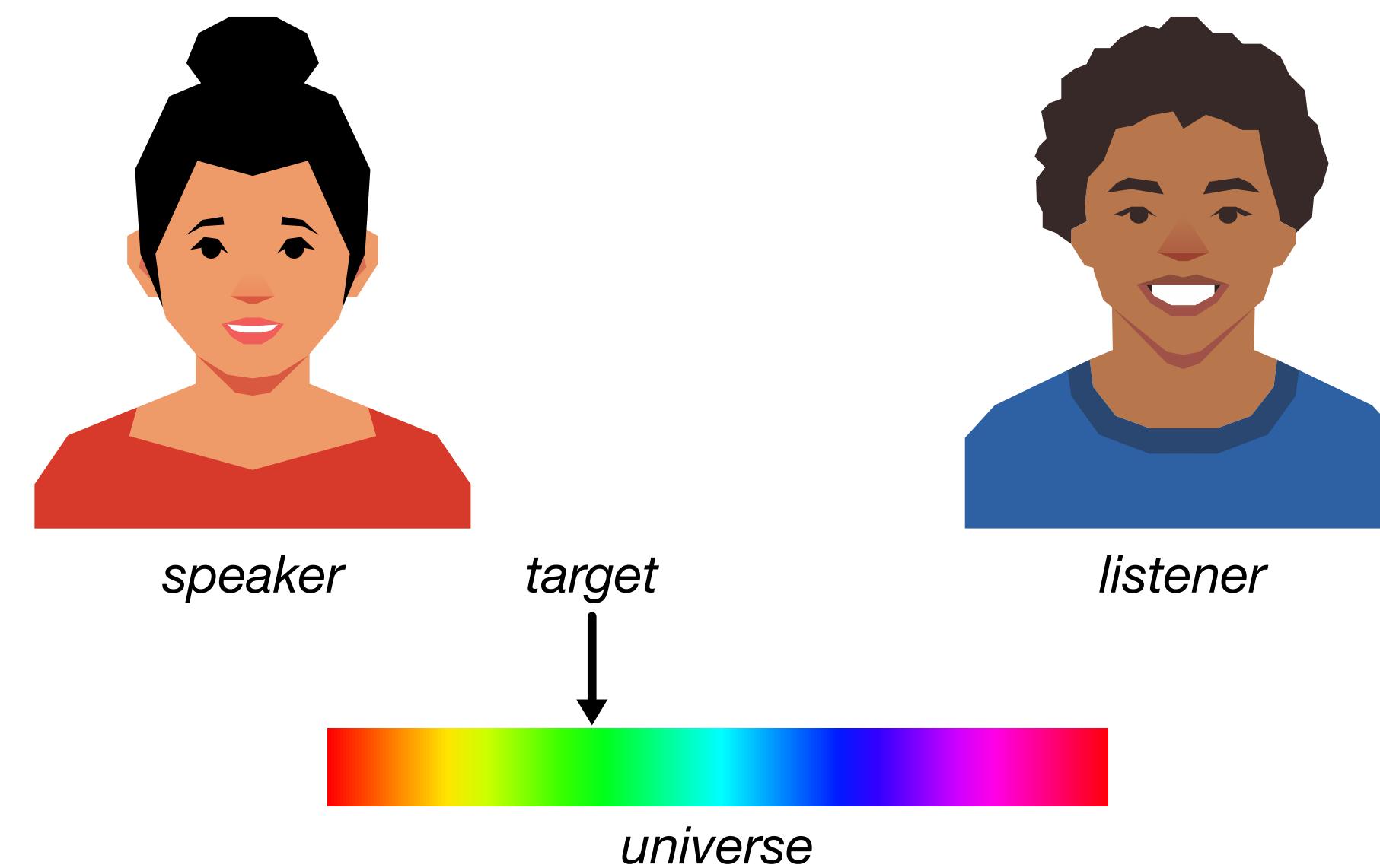


*Complexity*

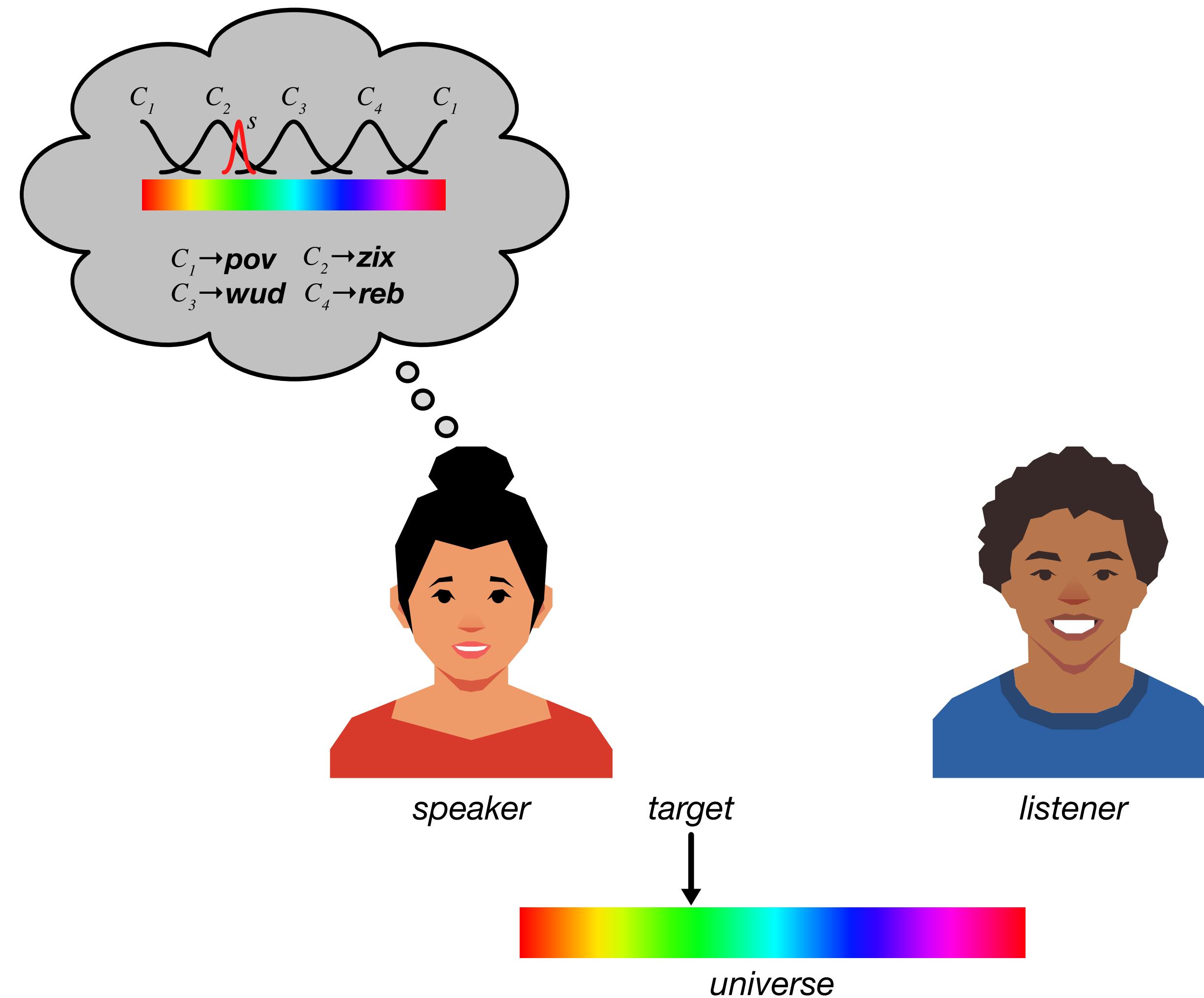


# *Informativeness*

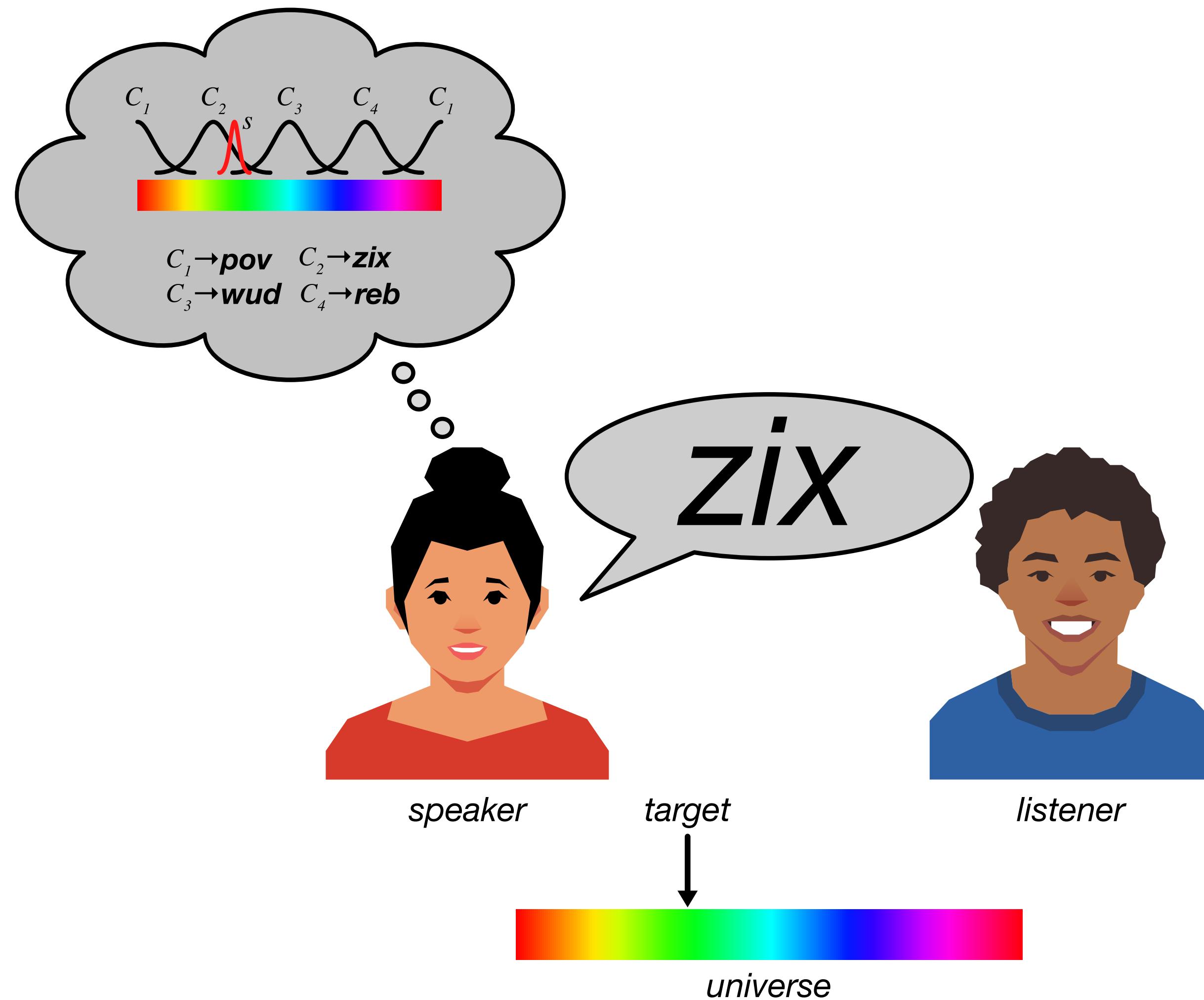
# Regier et al.'s informativeness model



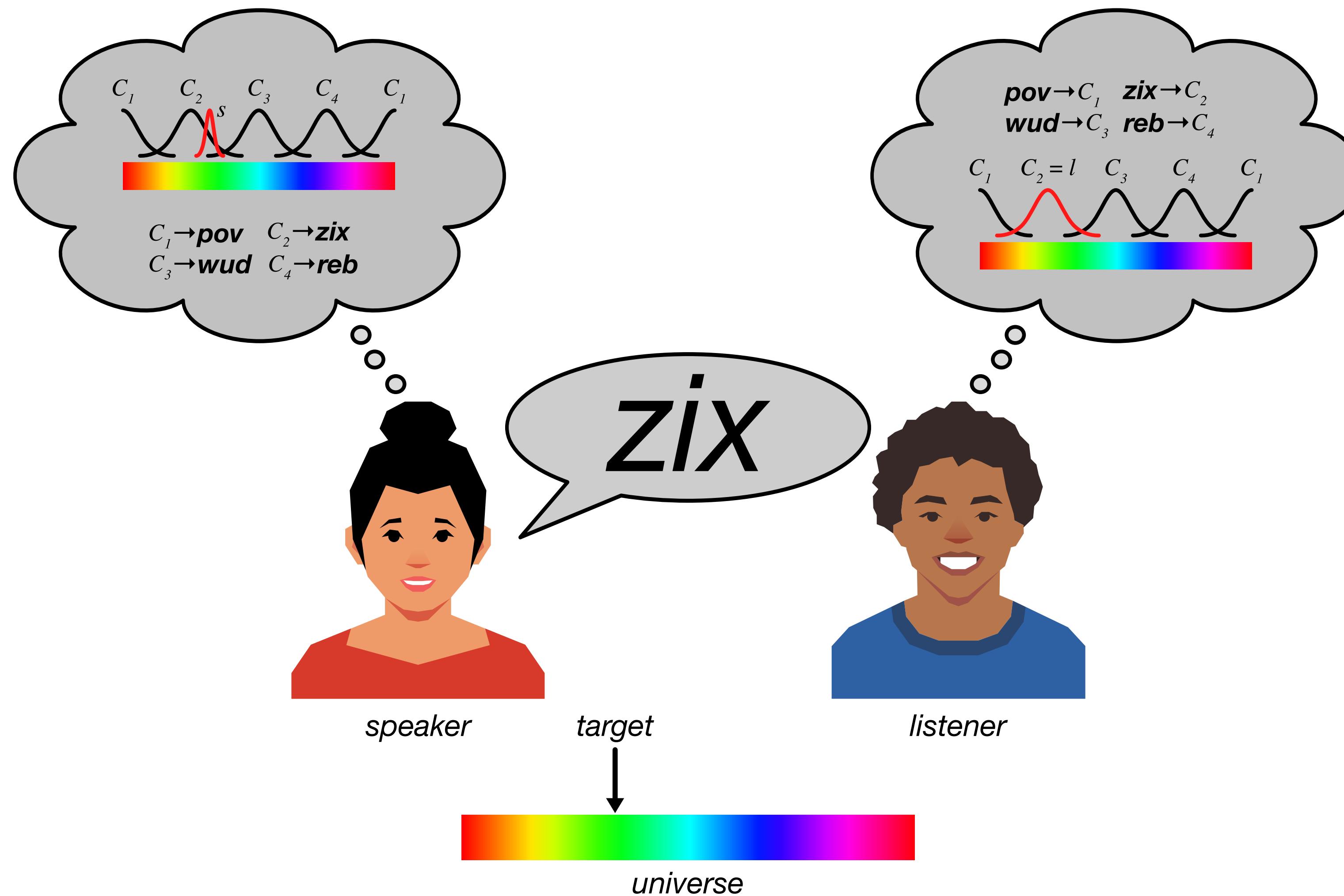
# Regier et al.'s informativeness model



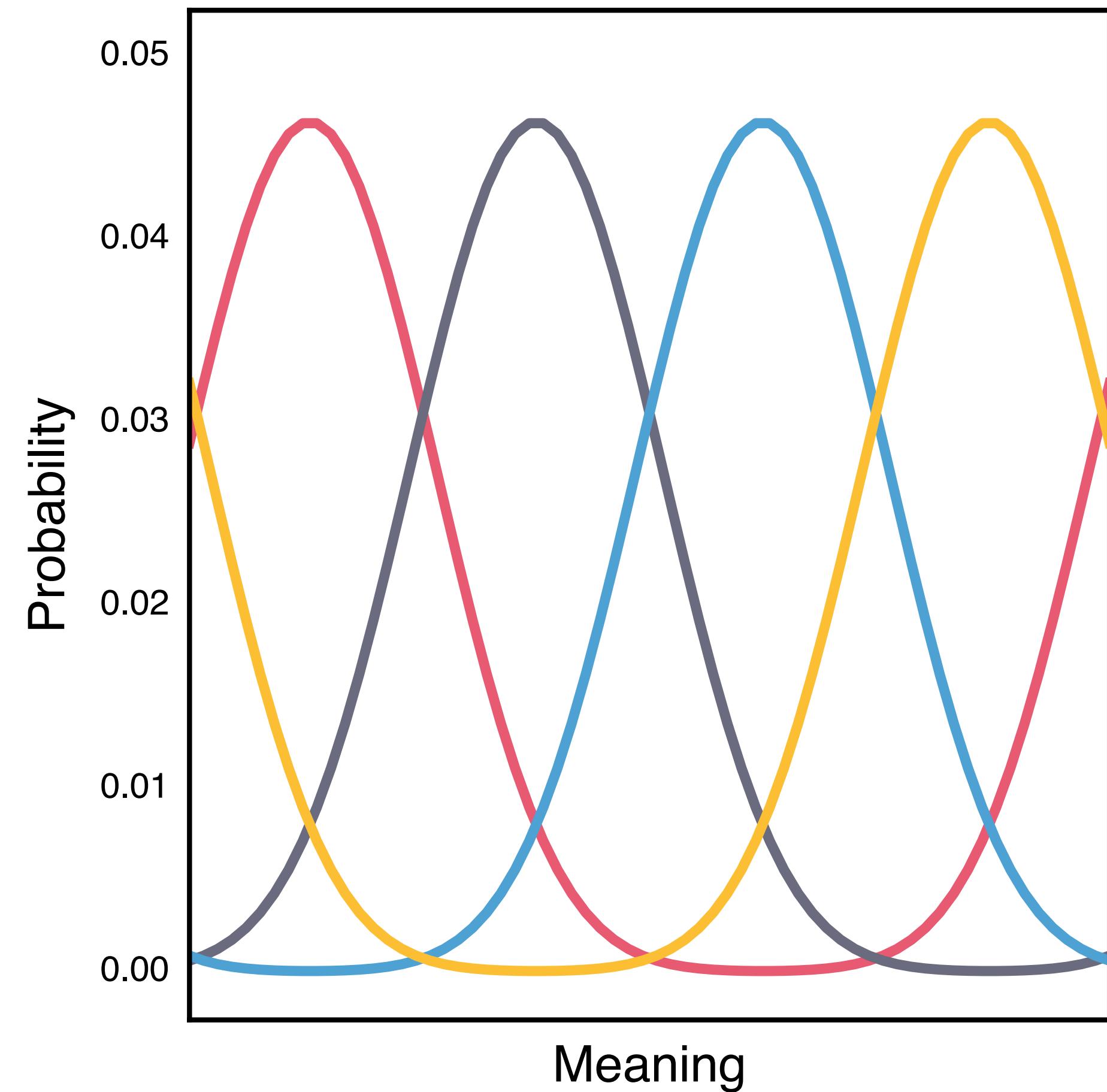
# Regier et al.'s informativeness model



# Regier et al.'s informativeness model



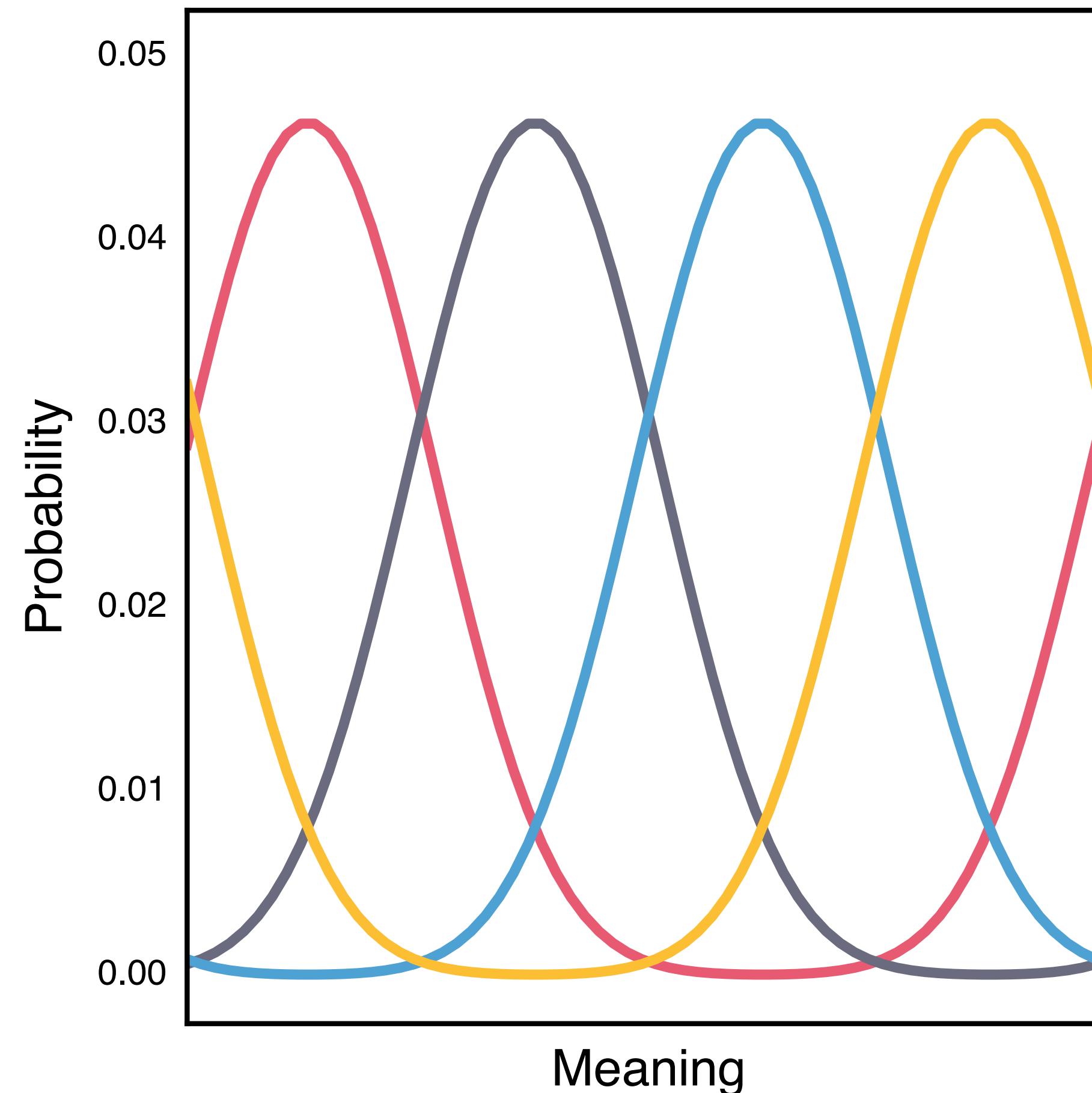
# Communicative cost



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

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

# Communicative cost



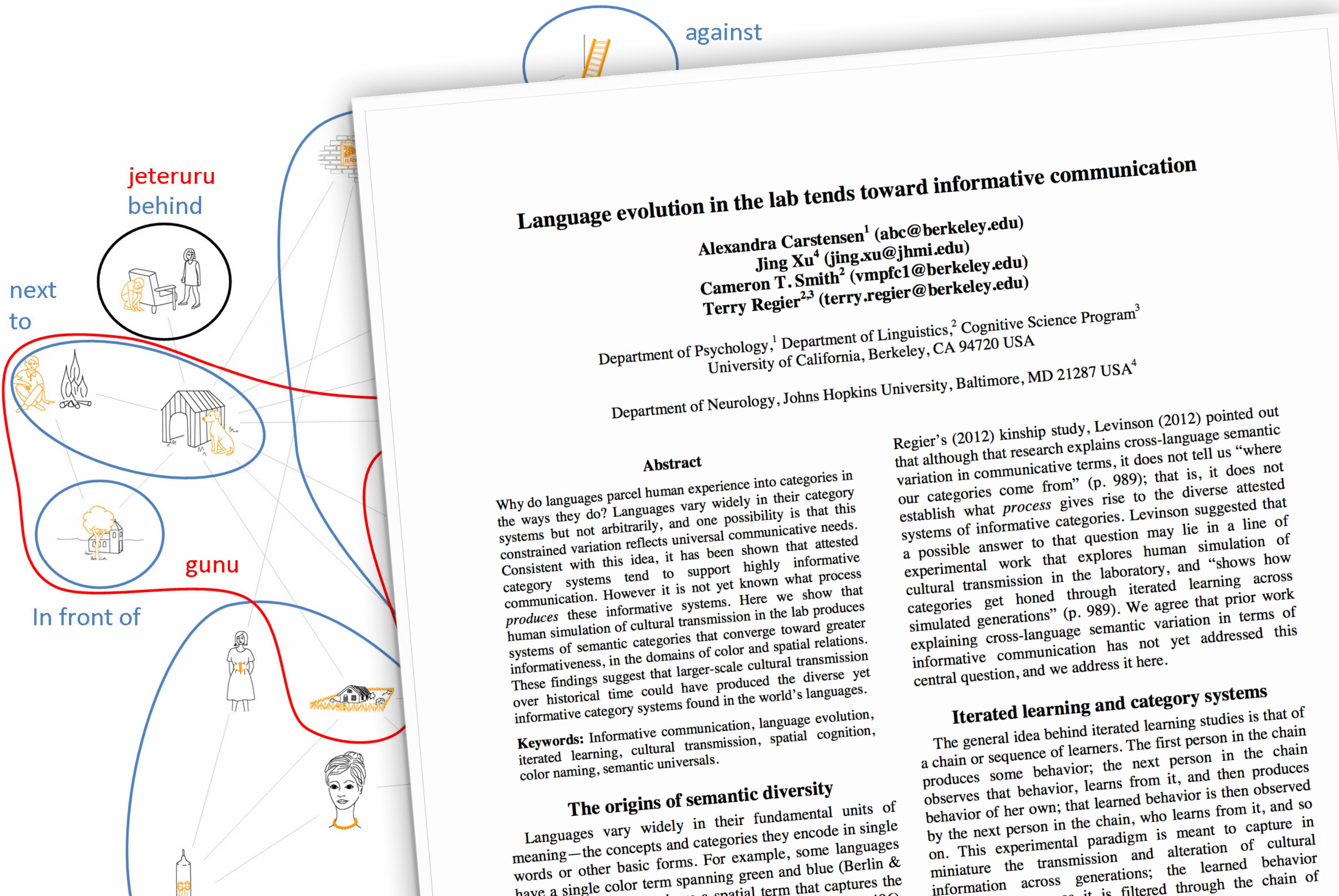
$$C_j(i) \propto \sum_{c \in C_j} e^{-\gamma d(i,c)^2}$$

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

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

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

# Can iterated learning give rise to informative languages?



## Language evolution in the lab tends toward informative communication

Alexandra Carstensen<sup>1</sup> (abc@berkeley.edu)  
Jing Xu<sup>4</sup> (jing.xu@jhmi.edu)  
Cameron T. Smith<sup>2</sup> (vmpfc1@berkeley.edu)  
Terry Regier<sup>2,3</sup> (terry.regier@berkeley.edu)

Department of Psychology,<sup>1</sup> Department of Linguistics,<sup>2</sup> Cognitive Science Program<sup>3</sup>  
University of California, Berkeley, CA 94720 USA  
Department of Neurology, Johns Hopkins University, Baltimore, MD 21287 USA<sup>4</sup>

### Abstract

Why do languages parcel human experience into categories in the ways they do? Languages vary widely in their category systems but not arbitrarily, and one possibility is that this constrained variation reflects universal communicative needs. Consistent with this idea, it has been shown that attested category systems tend to support highly informative communication. However it is not yet known what process produces these informative systems. Here we show that human simulation of cultural transmission in the lab produces systems of semantic categories that converge toward greater informativeness, in the domains of color and spatial relations. These findings suggest that larger-scale cultural transmission over historical time could have produced the diverse yet informative category systems found in the world's languages.

**Keywords:** Informative communication, language evolution, iterated learning, cultural transmission, spatial cognition, color naming, semantic universals.

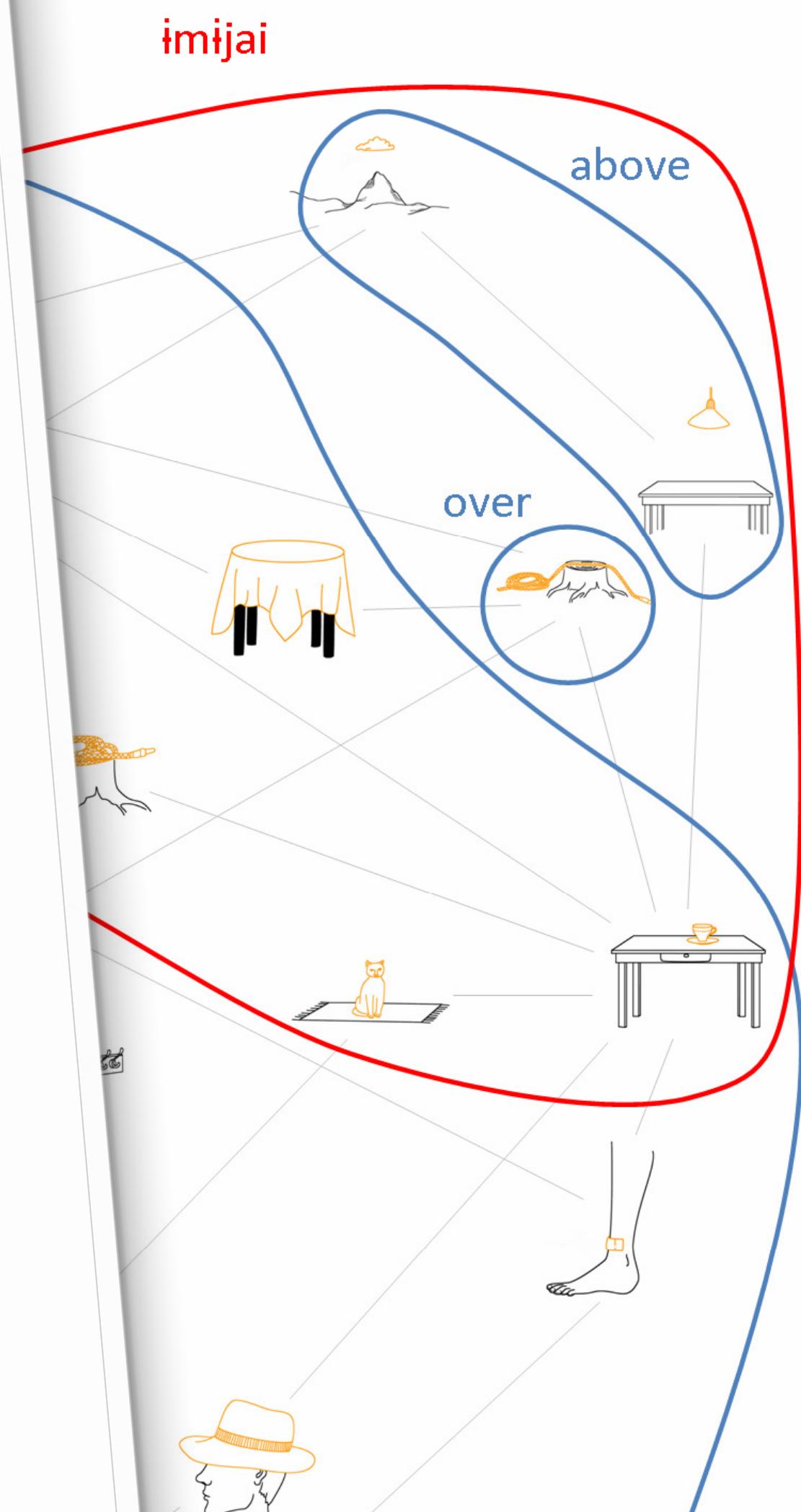
### The origins of semantic diversity

Languages vary widely in their fundamental units of meaning—the concepts and categories they encode in single words or other basic forms. For example, some languages have a single color term spanning green and blue (Berlin & Kay 1969), while others have a spatial term that captures the

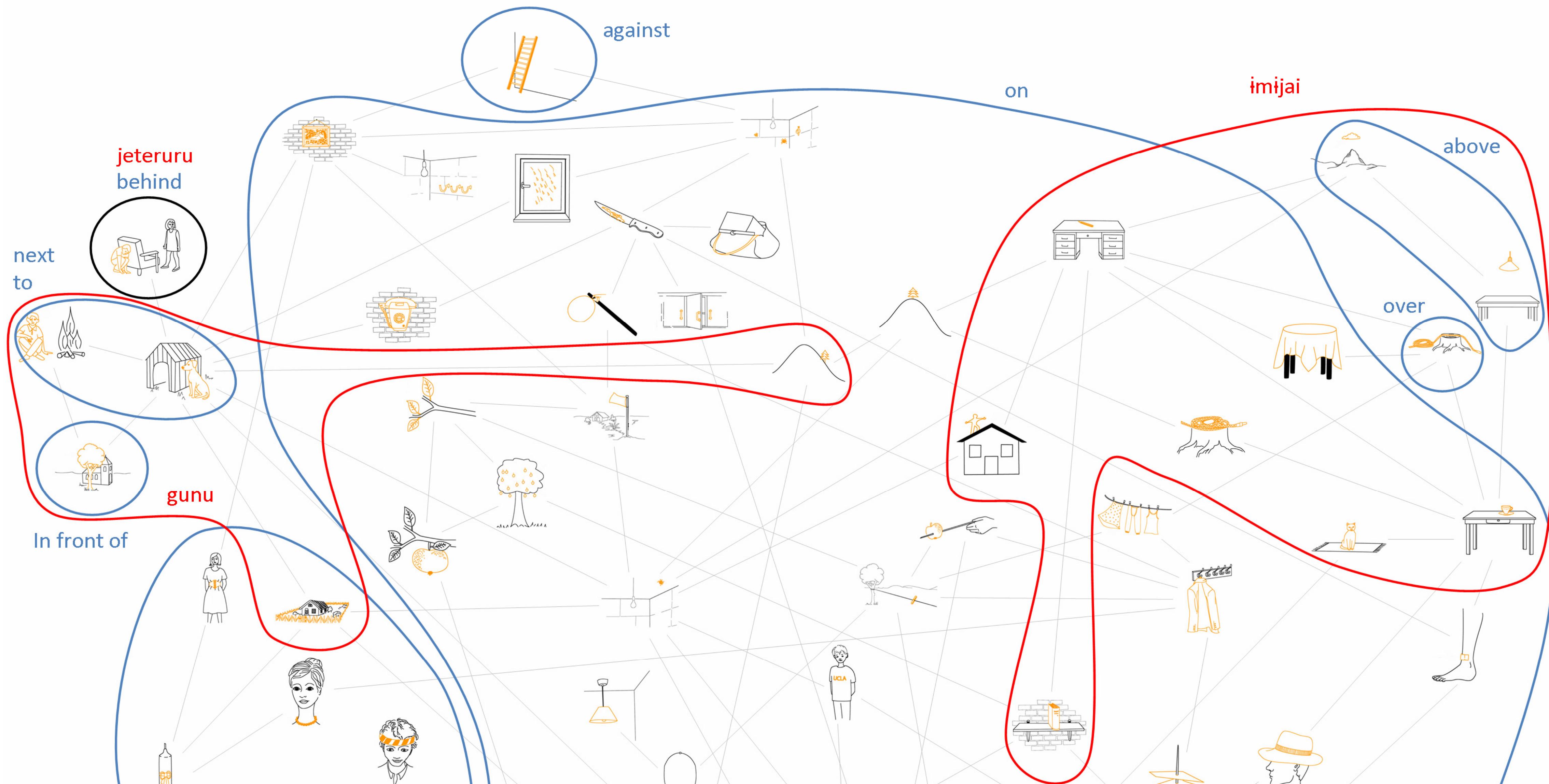
Regier's (2012) kinship study, Levinson (2012) pointed out that although that research explains cross-language semantic variation in communicative terms, it does not tell us “where our categories come from” (p. 989); that is, it does not establish what *process* gives rise to the diverse attested systems of informative categories. Levinson suggested that a possible answer to that question may lie in a line of experimental work that explores human simulation of cultural transmission in the laboratory, and “shows how categories get honed through iterated learning across simulated generations” (p. 989). We agree that prior work explaining cross-language semantic variation in terms of informative communication has not yet addressed this central question, and we address it here.

### Iterated learning and category systems

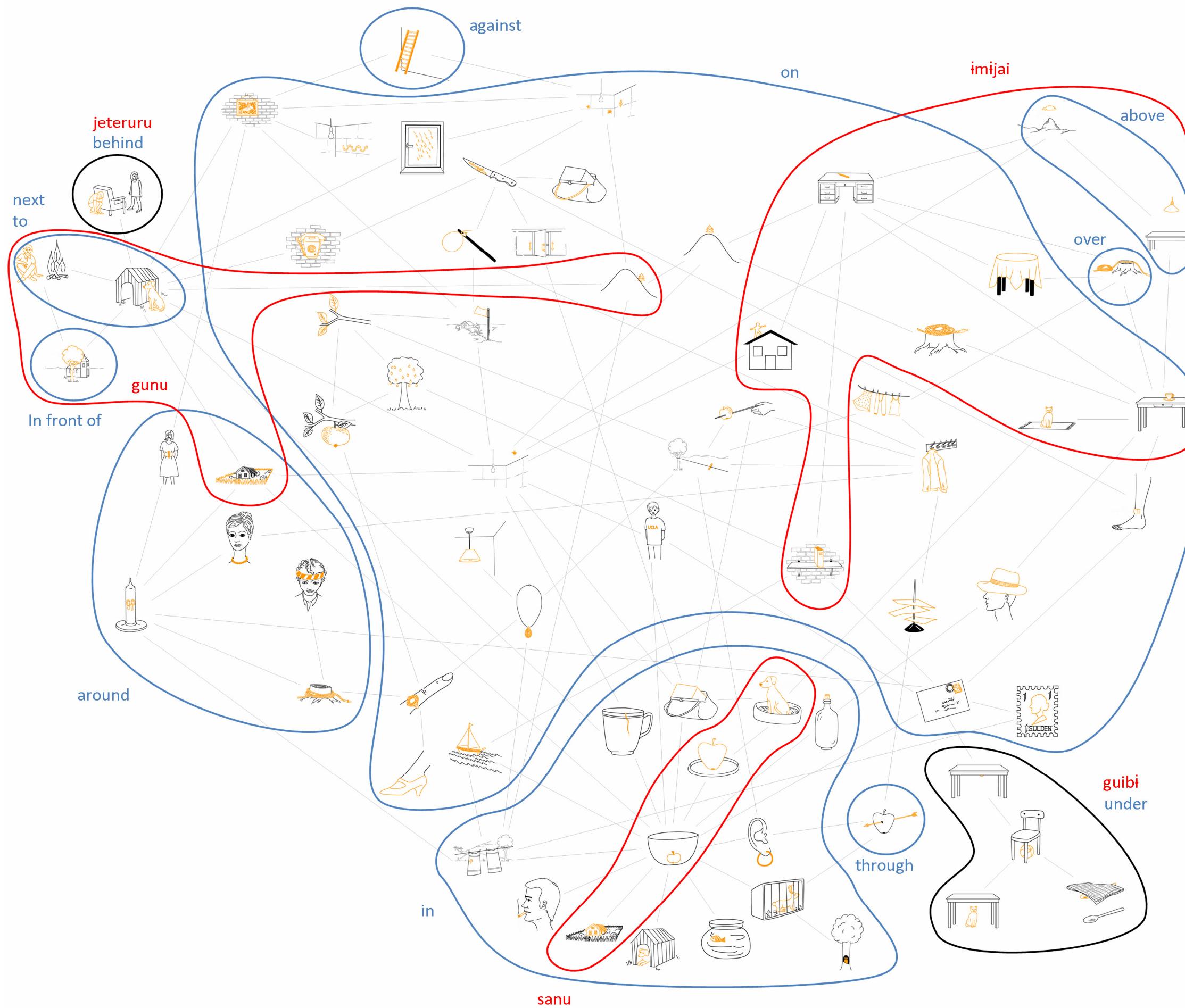
The general idea behind iterated learning studies is that of a chain or sequence of learners. The first person in the chain produces some behavior; the next person in the chain observes that behavior, learns from it, and then produces behavior of her own; that learned behavior is then observed by the next person in the chain, who learns from it, and so on. This experimental paradigm is meant to capture in miniature the transmission and alteration of cultural information across generations; the learned behavior as it is filtered through the chain of



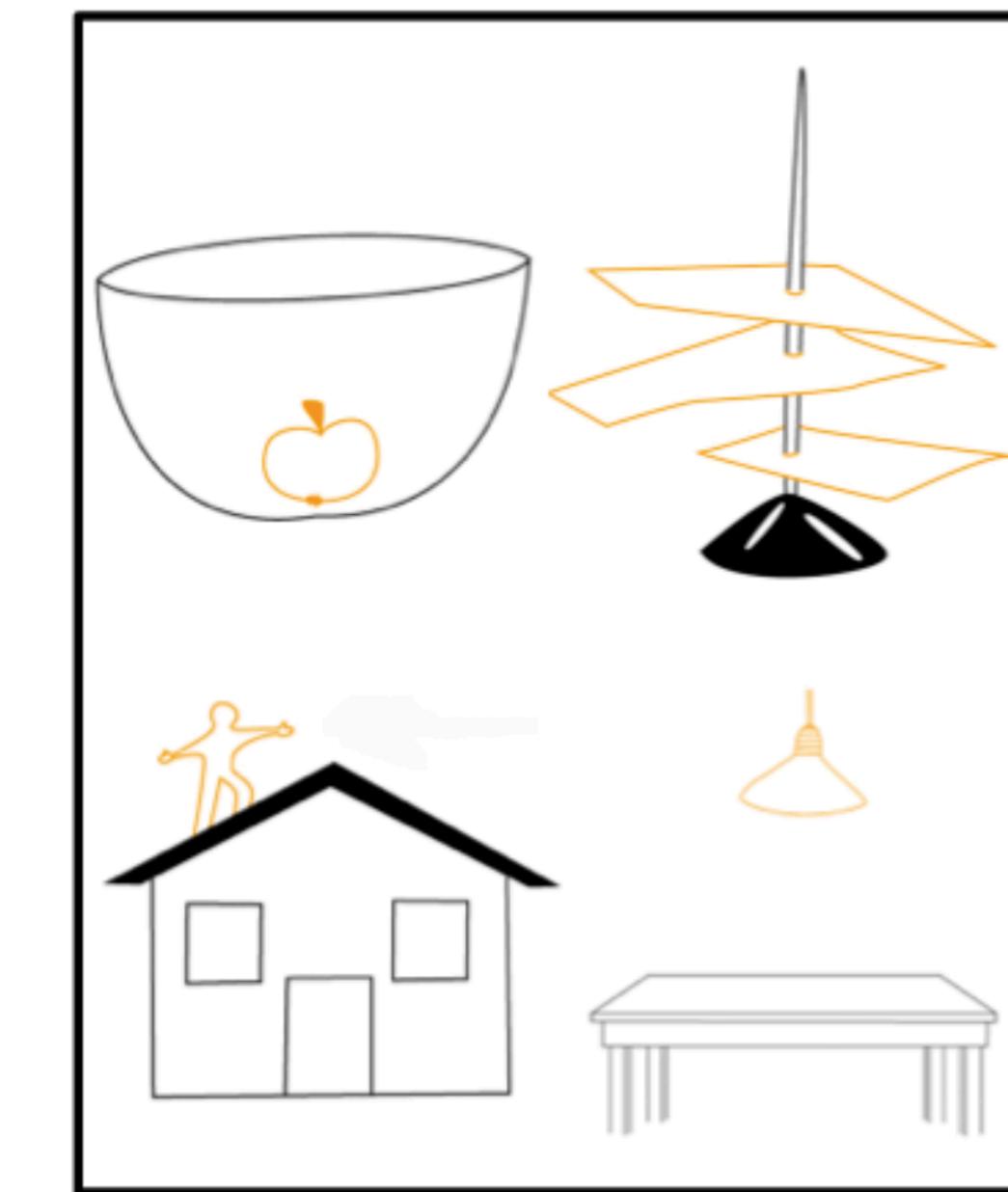
# Can iterated learning give rise to informative languages?



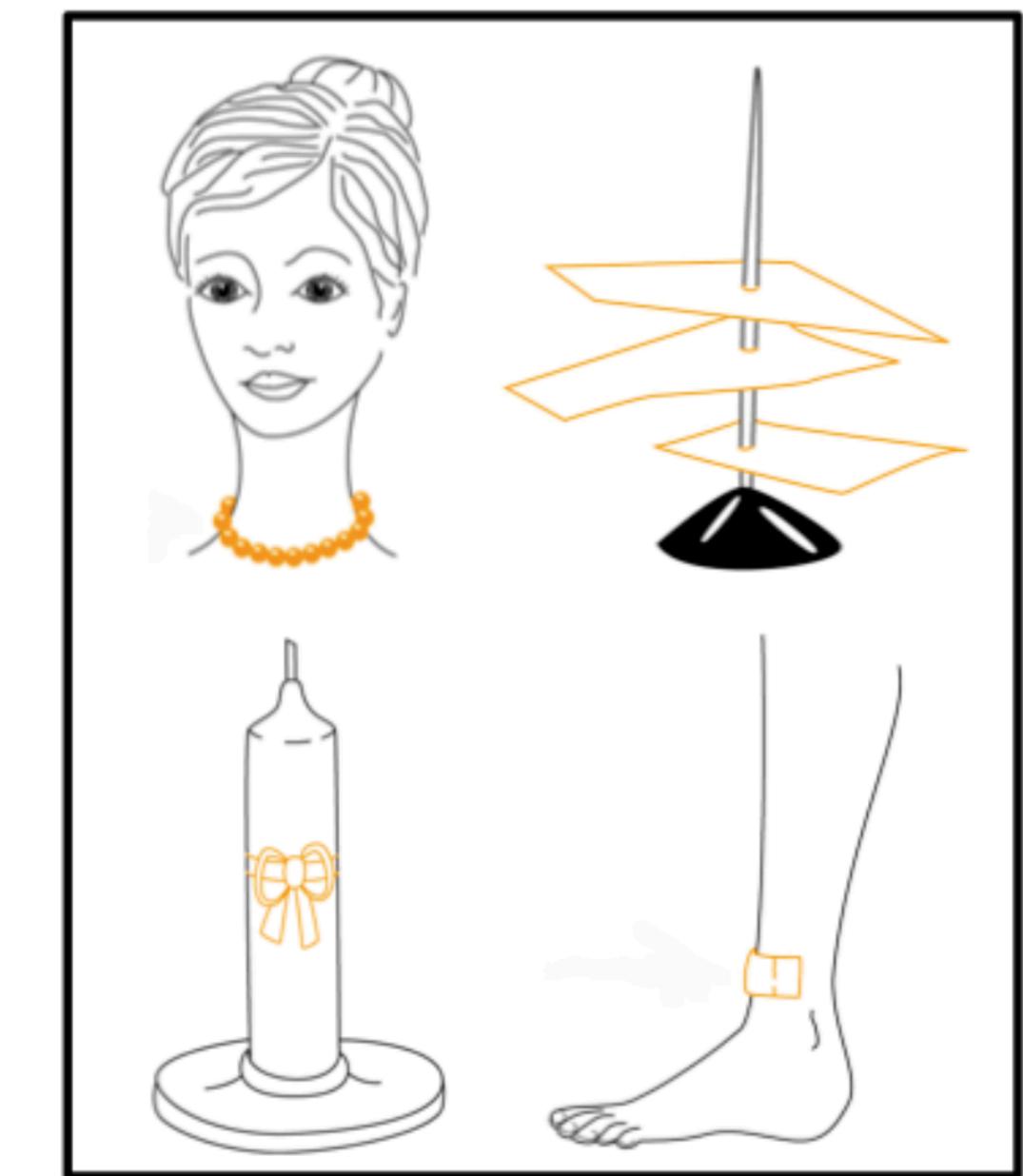
# Can iterated learning give rise to informative languages?



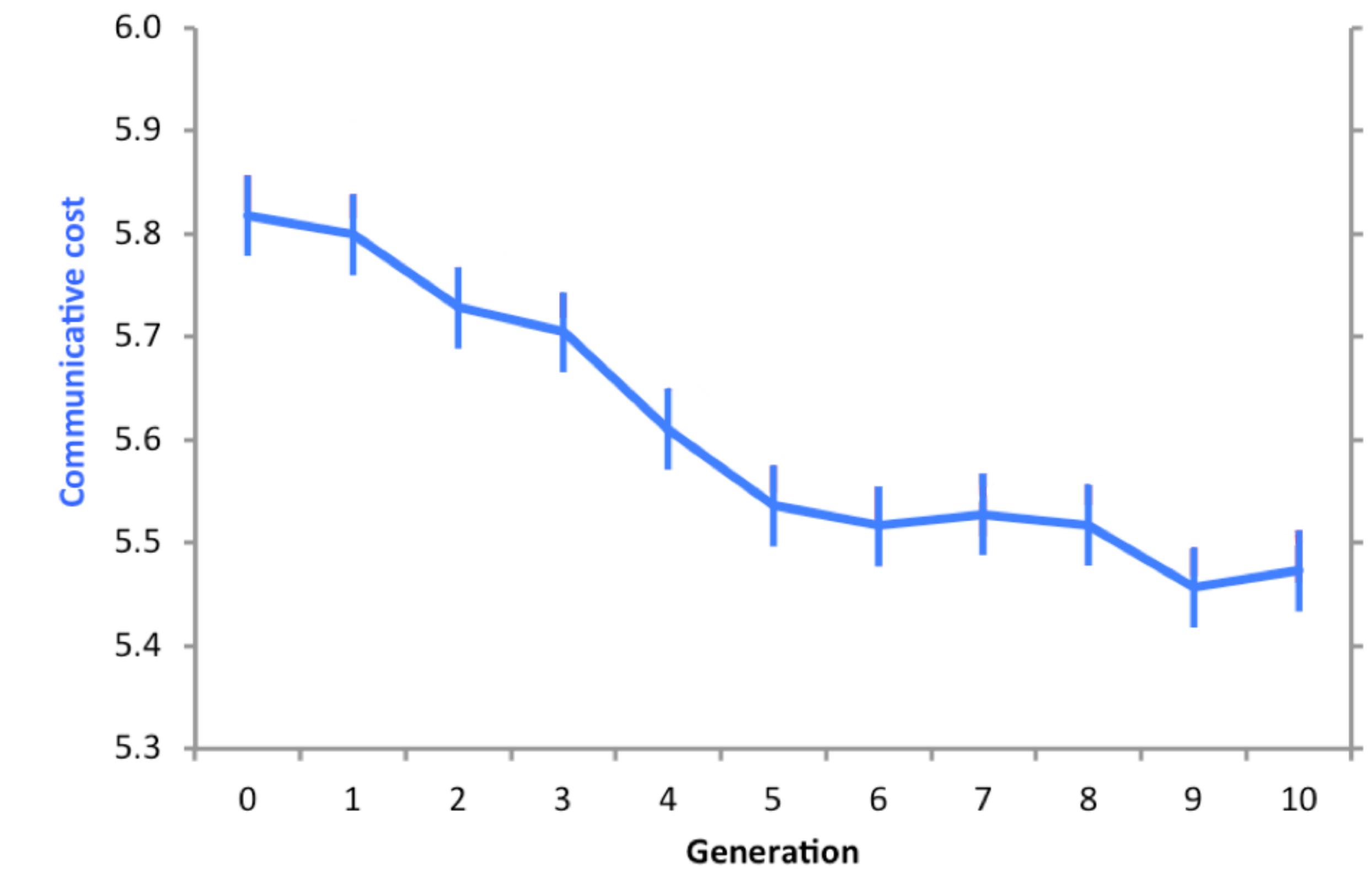
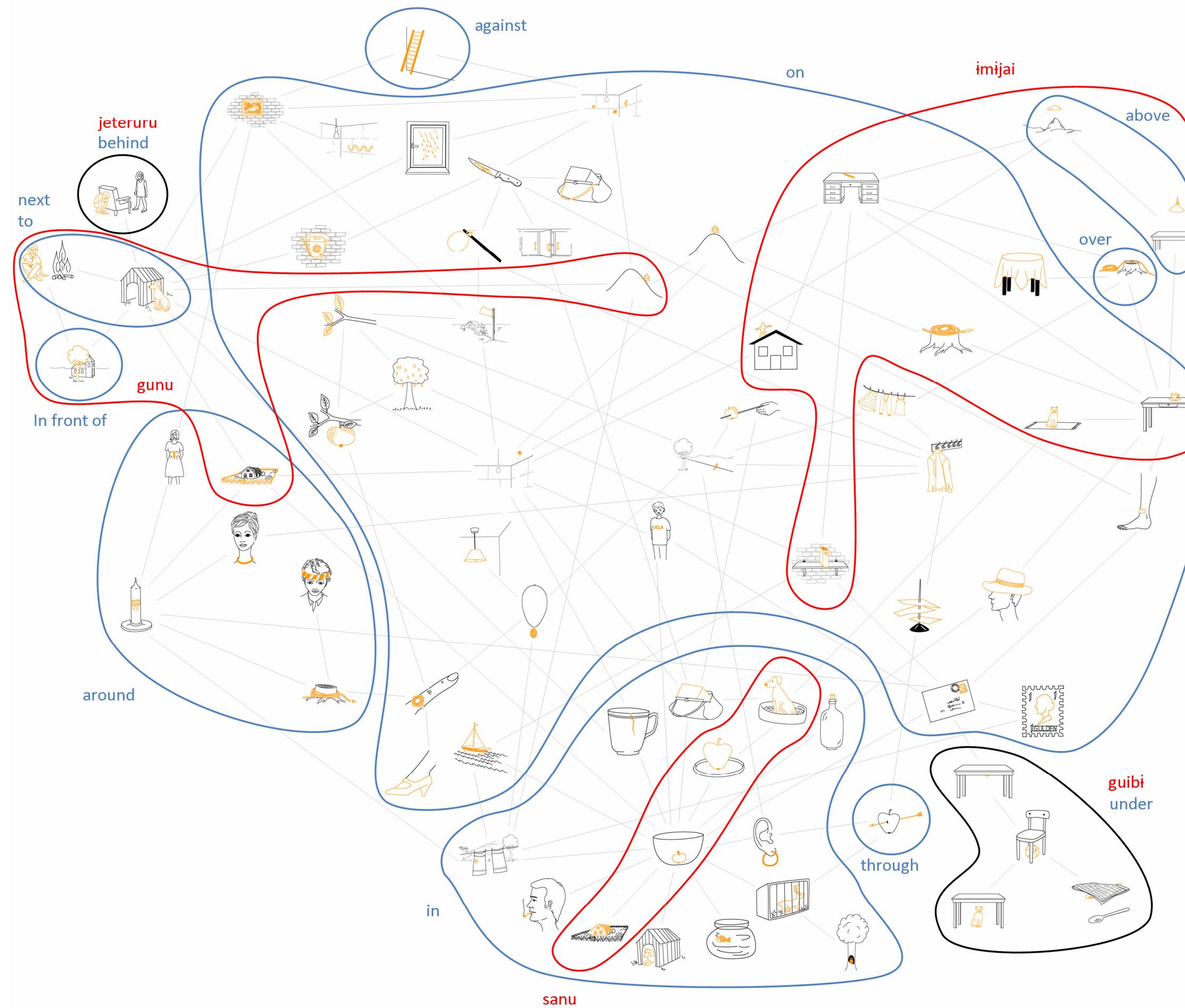
Generation 0



Generation 10



# Can iterated learning give rise to informative languages?



Carstensen, Xu, Smith, Regier (2015)

# *Experiments*

# Training phase

localhost/~jon/shepard/

## Stage 1: Training

**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.

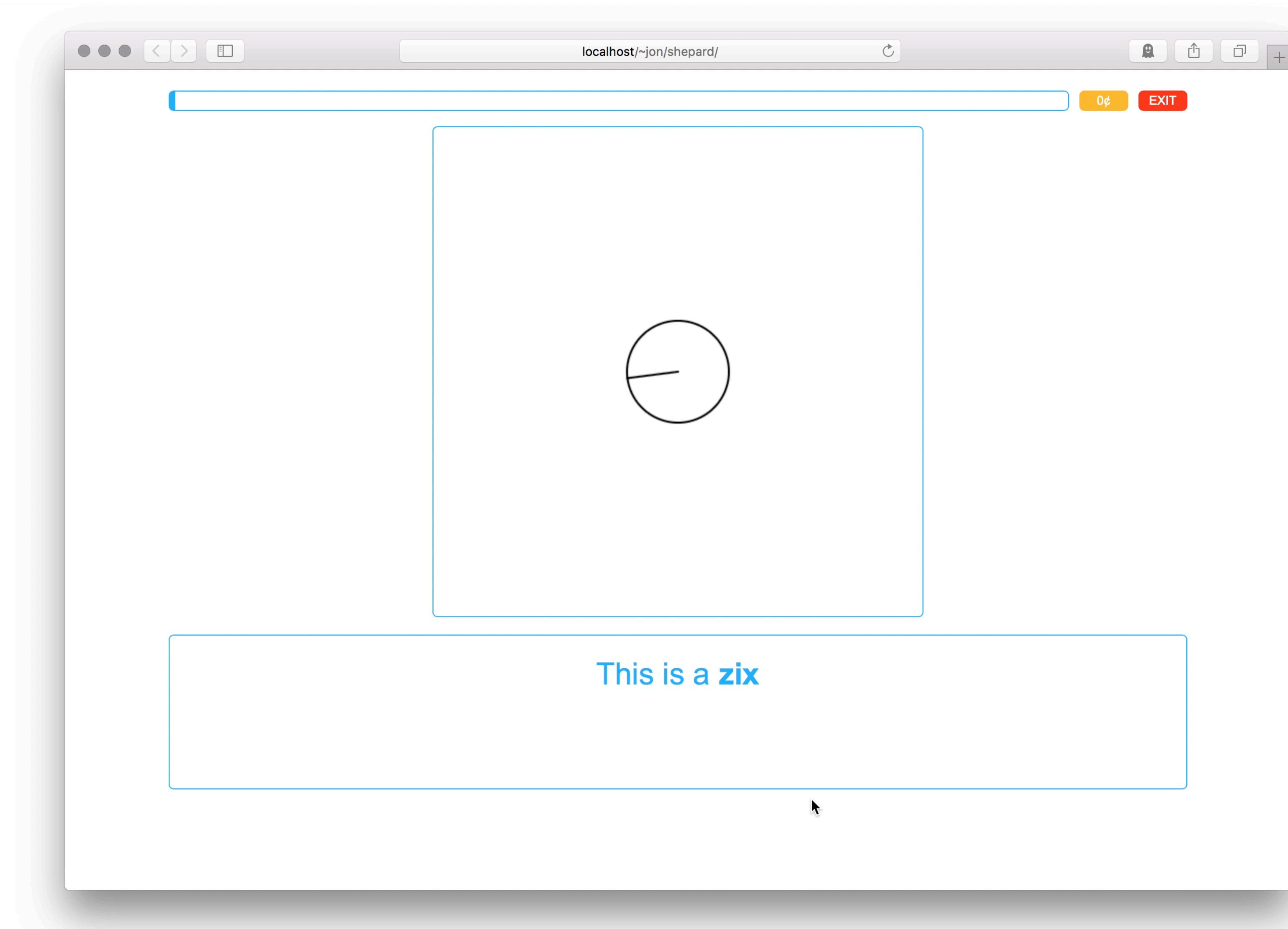
The interface shows three stages of a word-learning task:

- 1 Look at the picture**: A triangle icon. **2 Learn the word**: Text area containing "This is a tid".
- 3 Click on the word to confirm you learned it**: Text area asking "What is it called?" with options: **tid**, **bup**, **gax**, **fos**.
- 4 Sometimes you'll see a picture that you saw before**: A pentagon icon. **5 Try to recall the correct word**: Text area asking "What is this called?" with options: **tid**, **gax**, **fos**, **bup**.

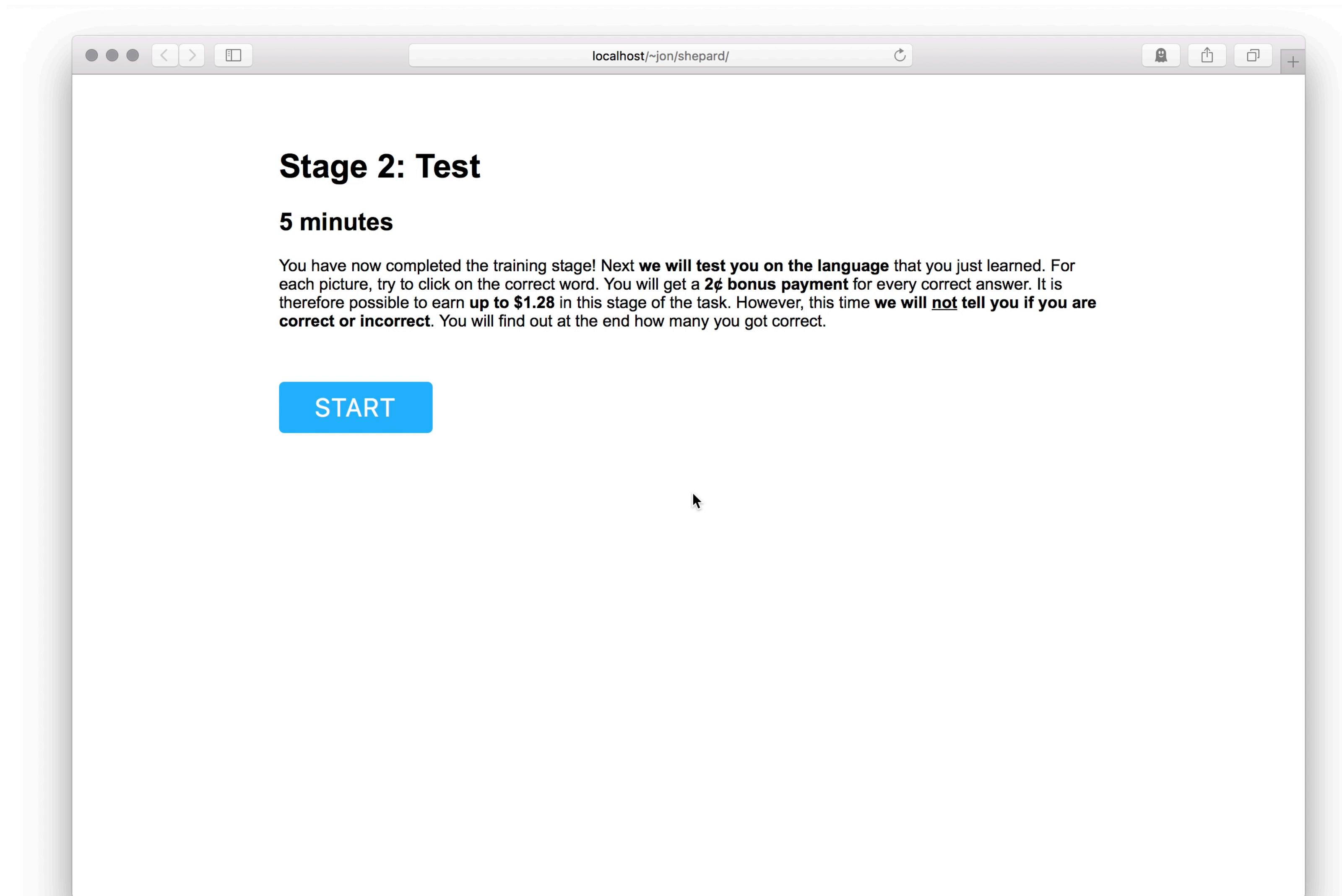
**6 If correct, you get a 2¢ bonus**

**START**

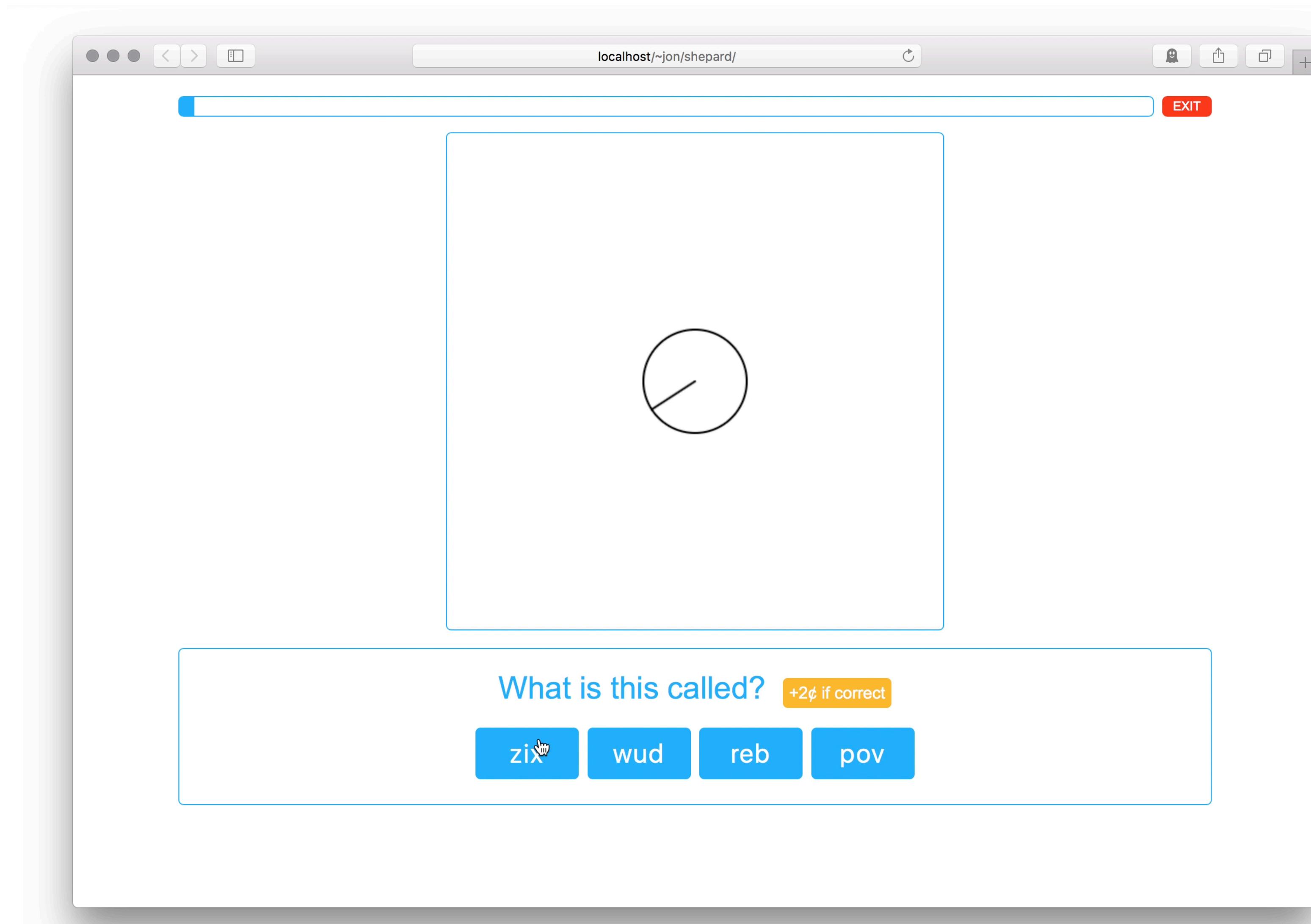
# Training phase



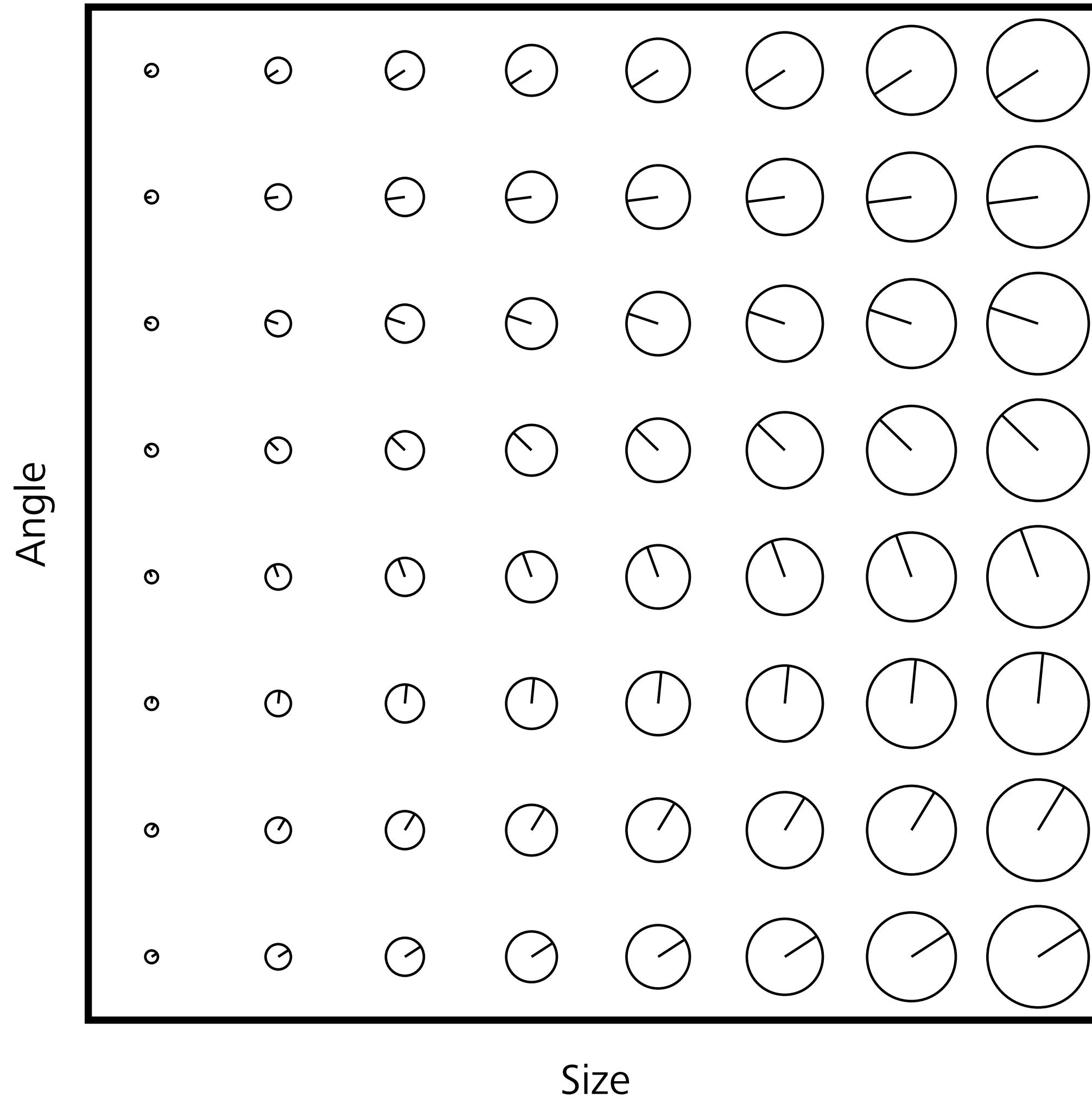
# Test phase



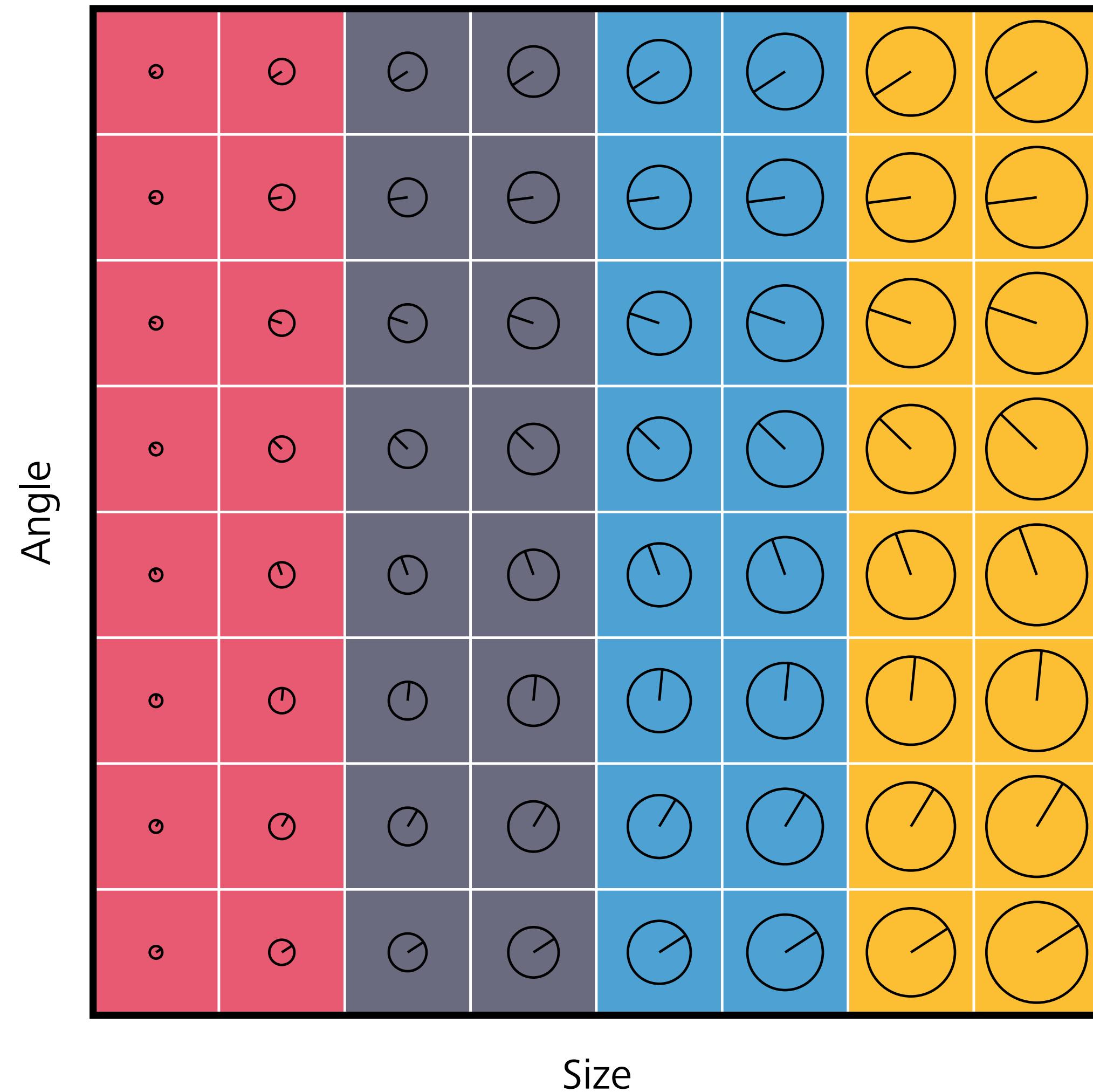
# Test phase



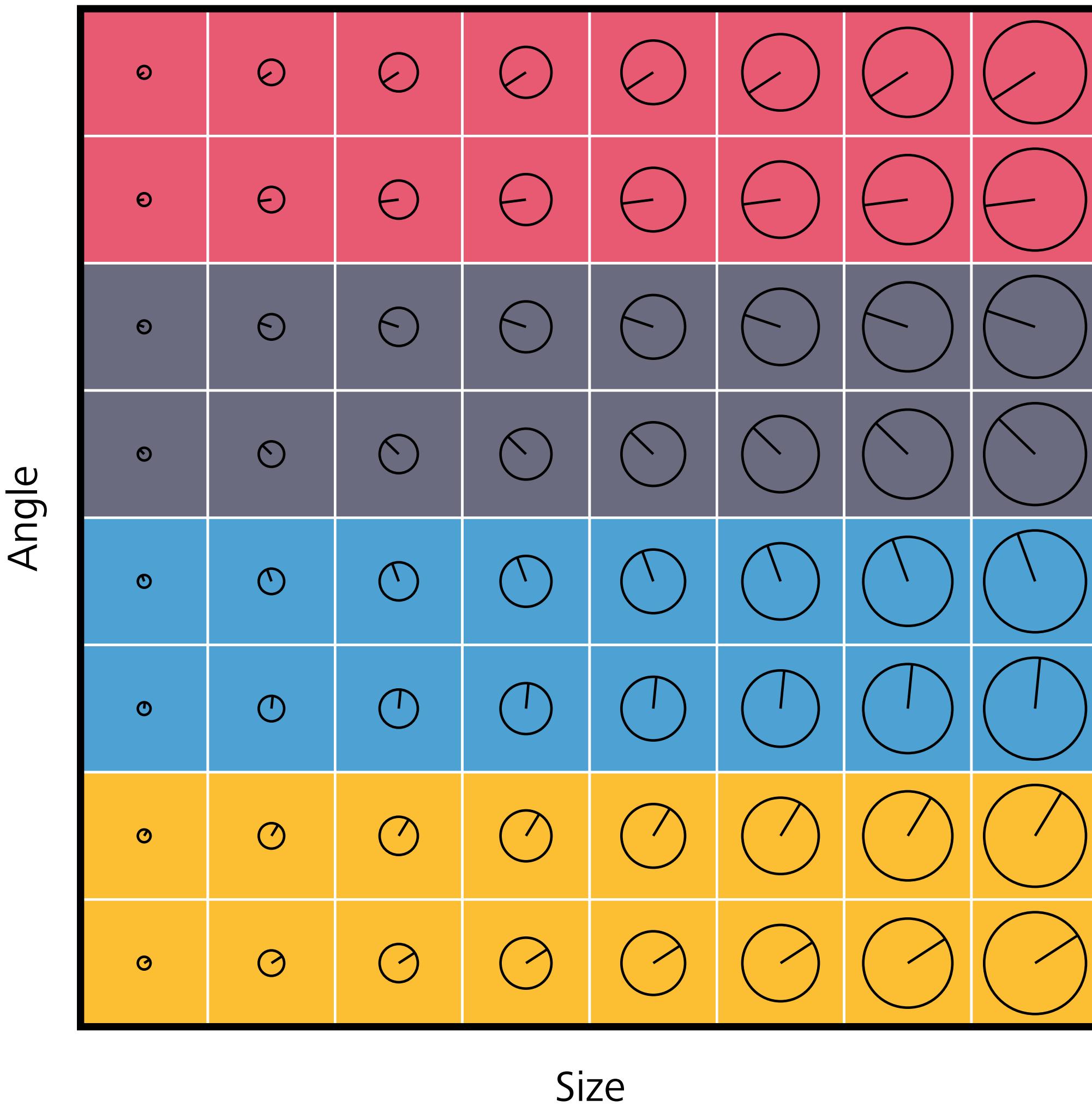
# Stimuli



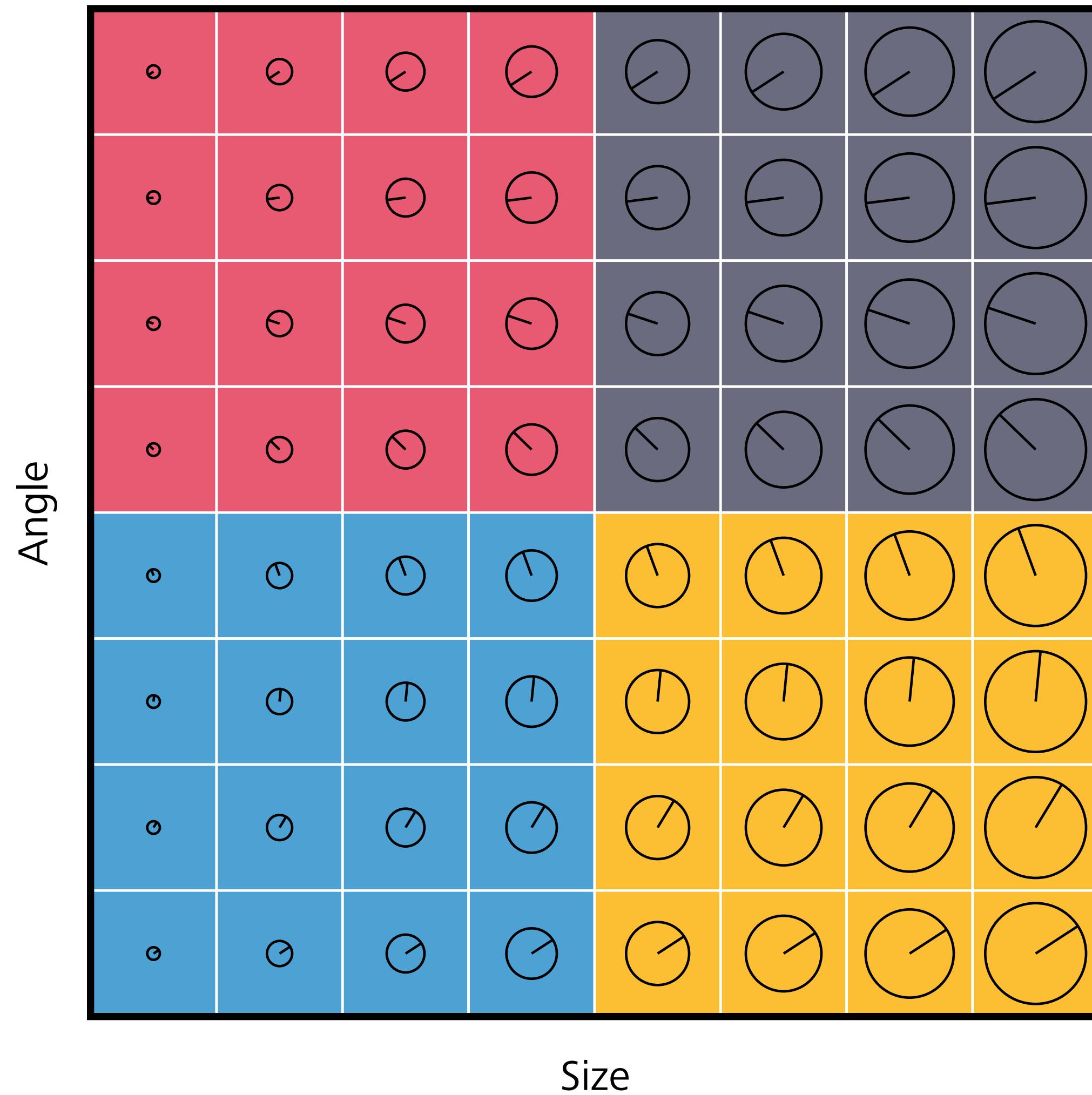
# Stimuli



# Stimuli

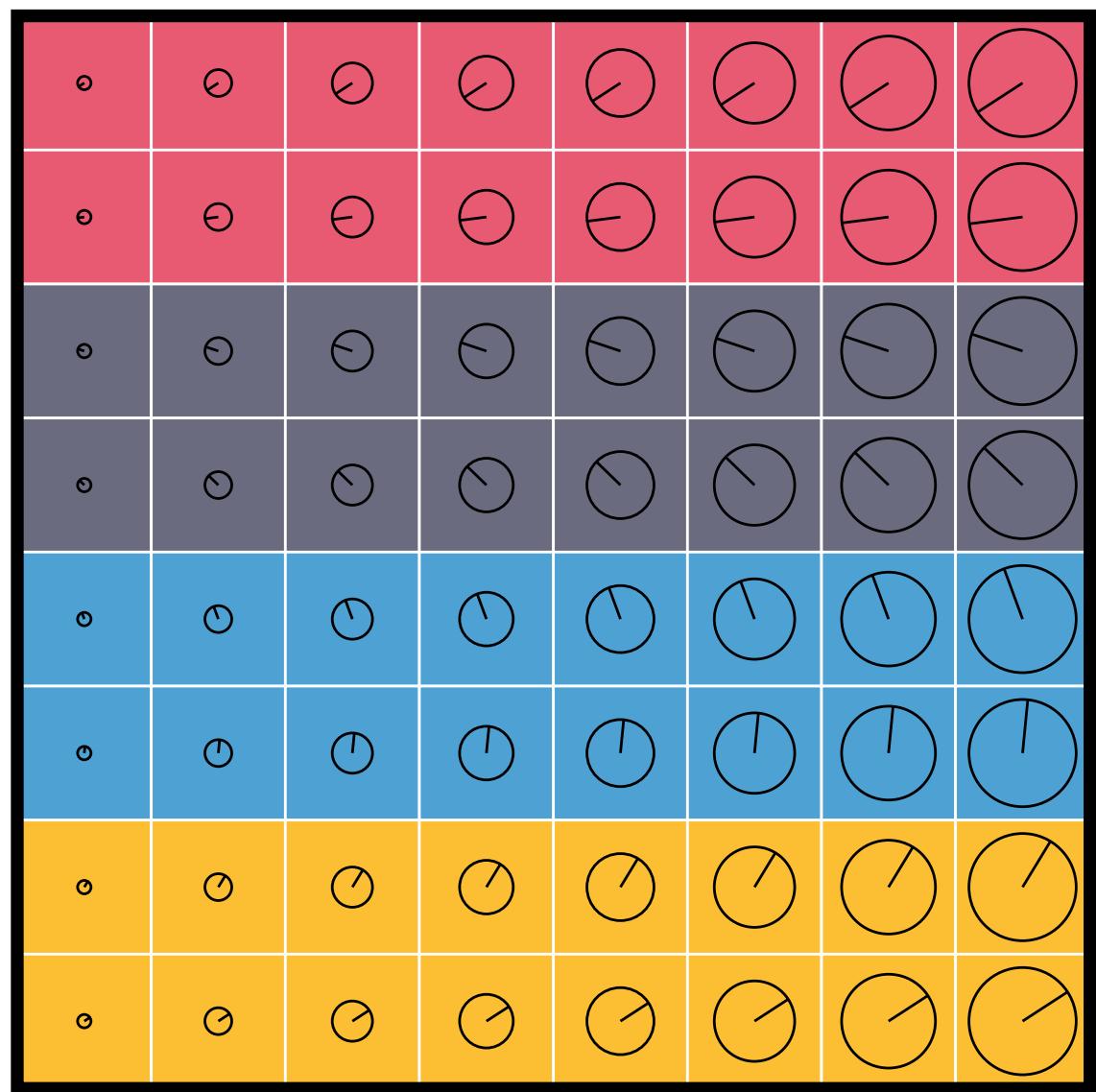


# Stimuli

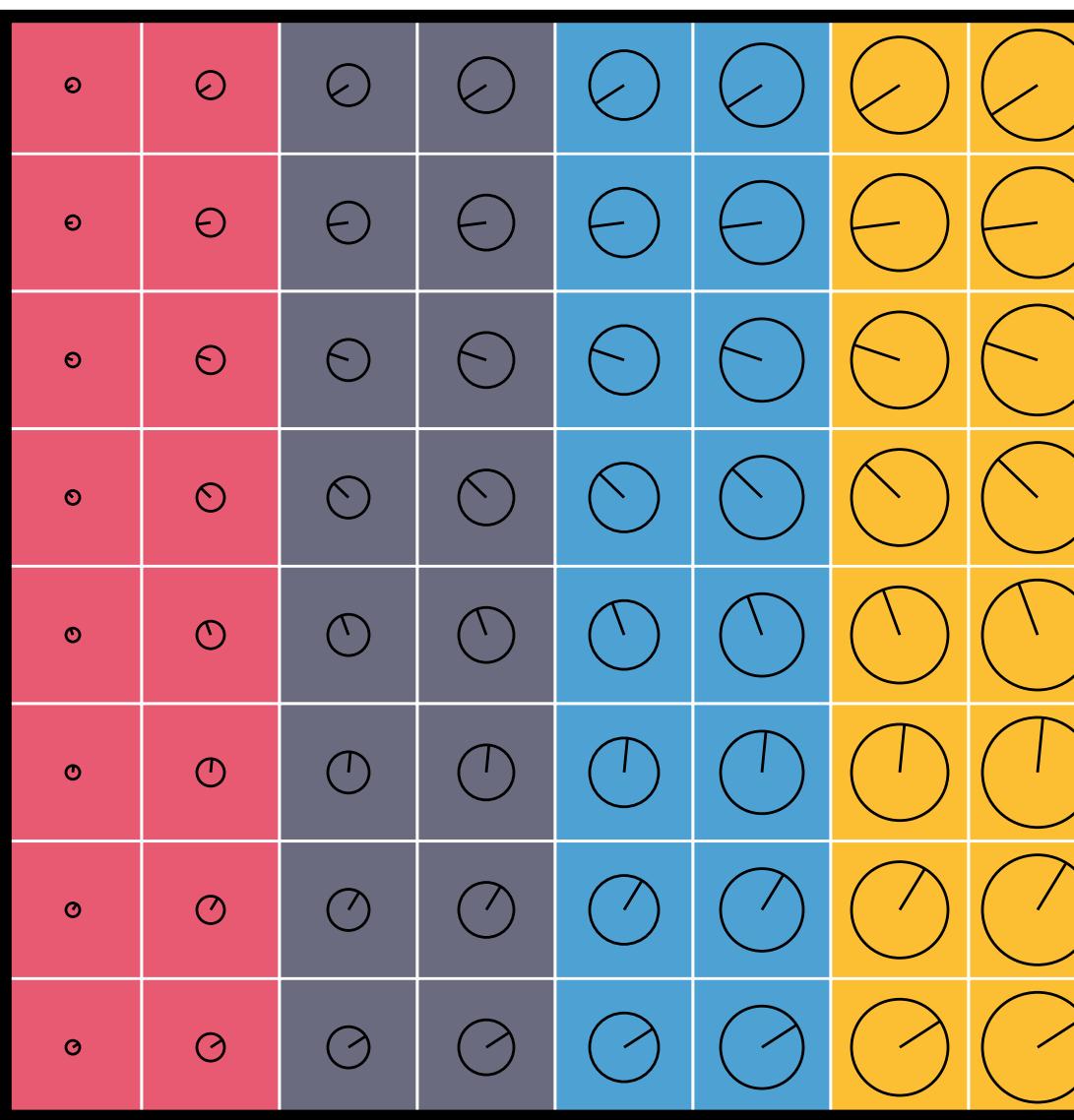


# Which is easiest to learn?

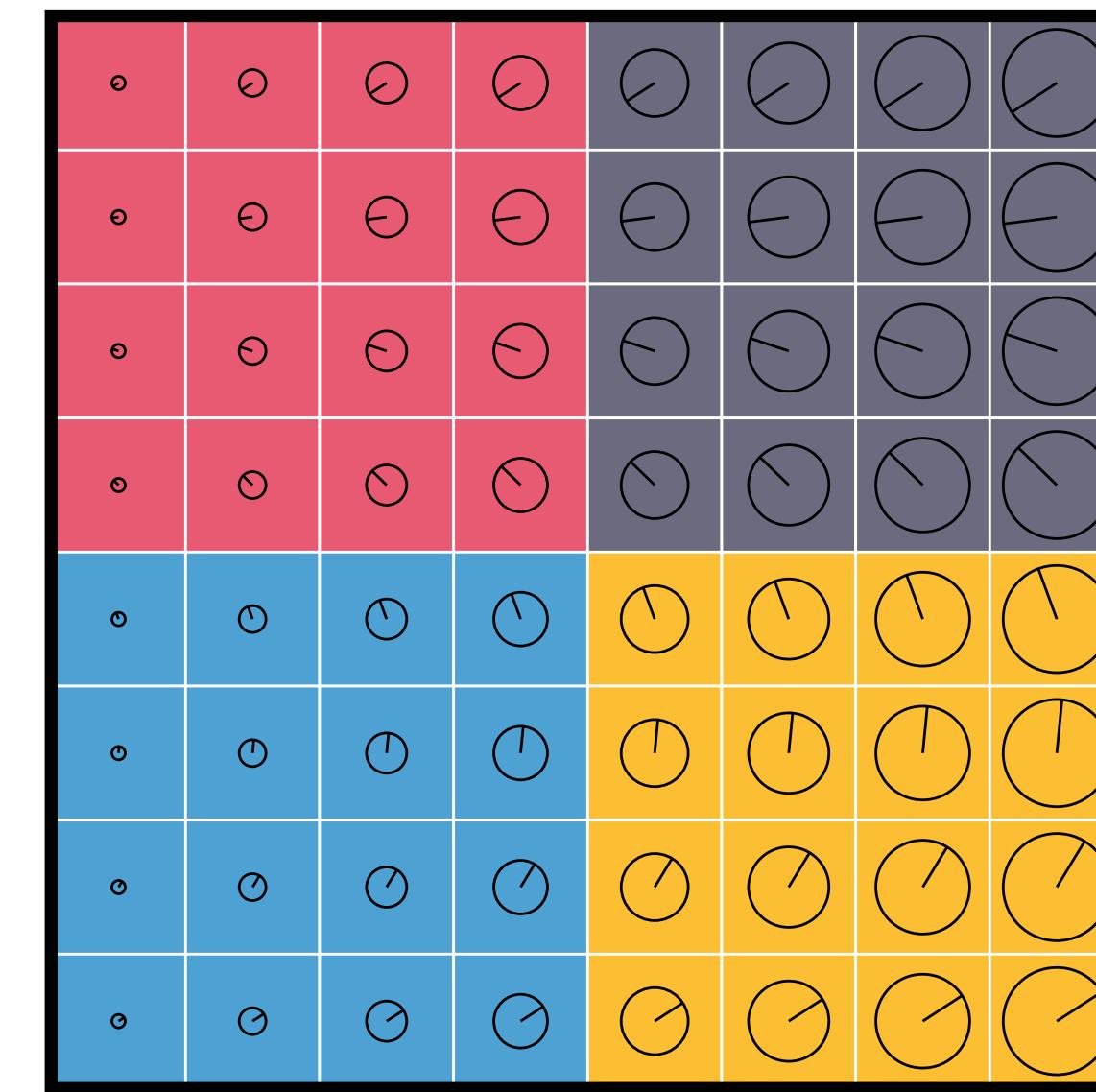
*Angle-only*



*Size-only*

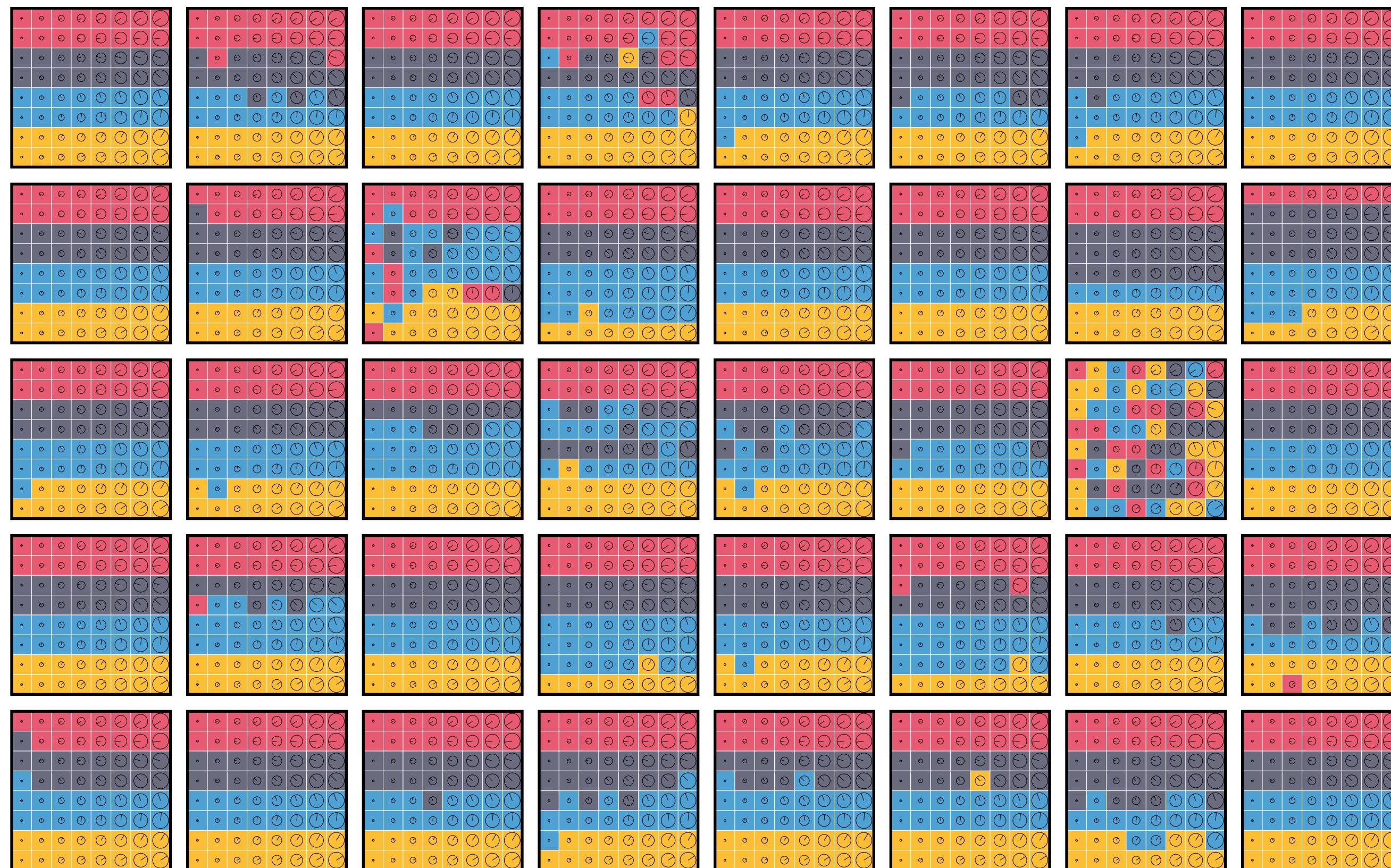
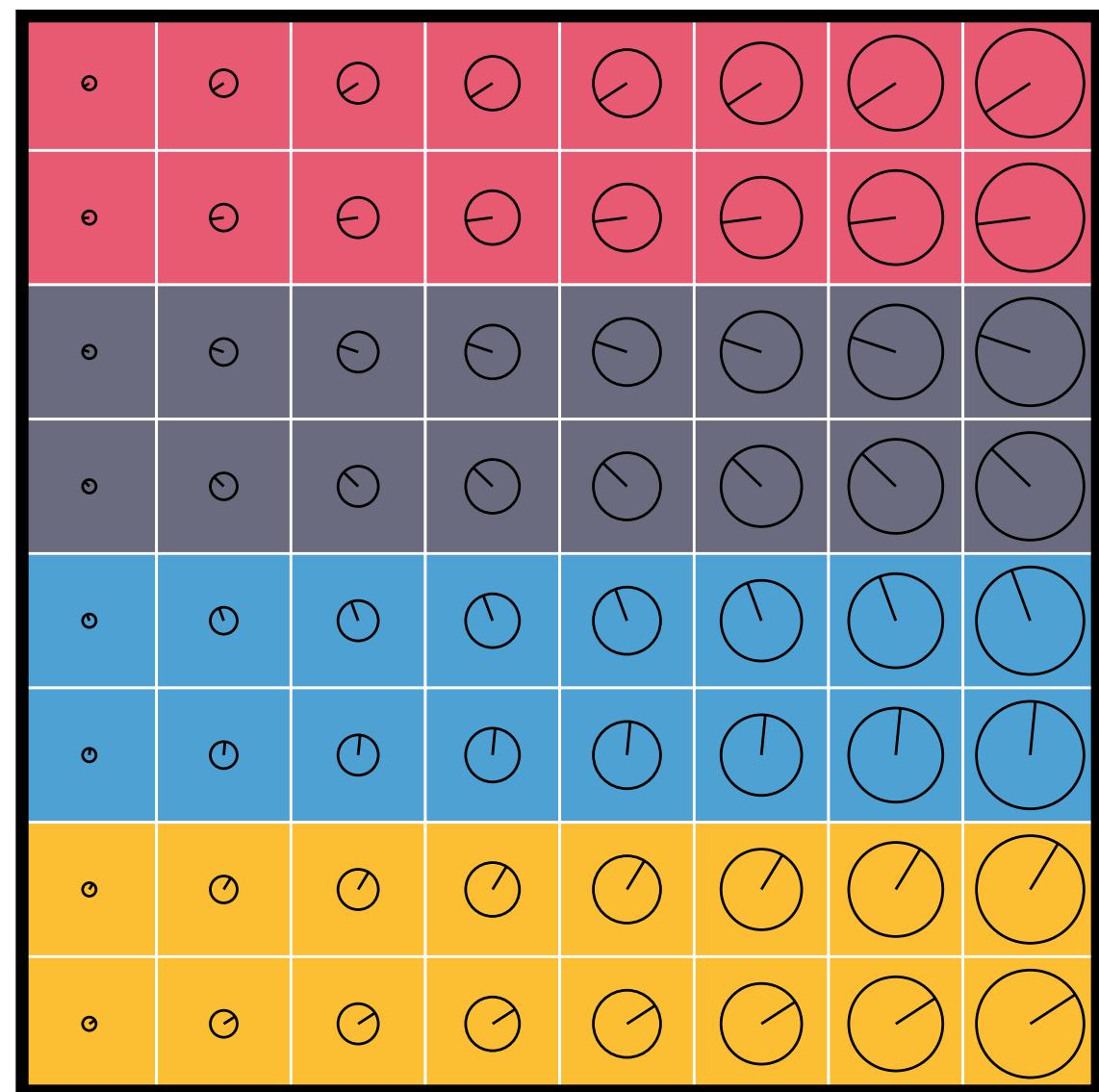


*Angle & Size*



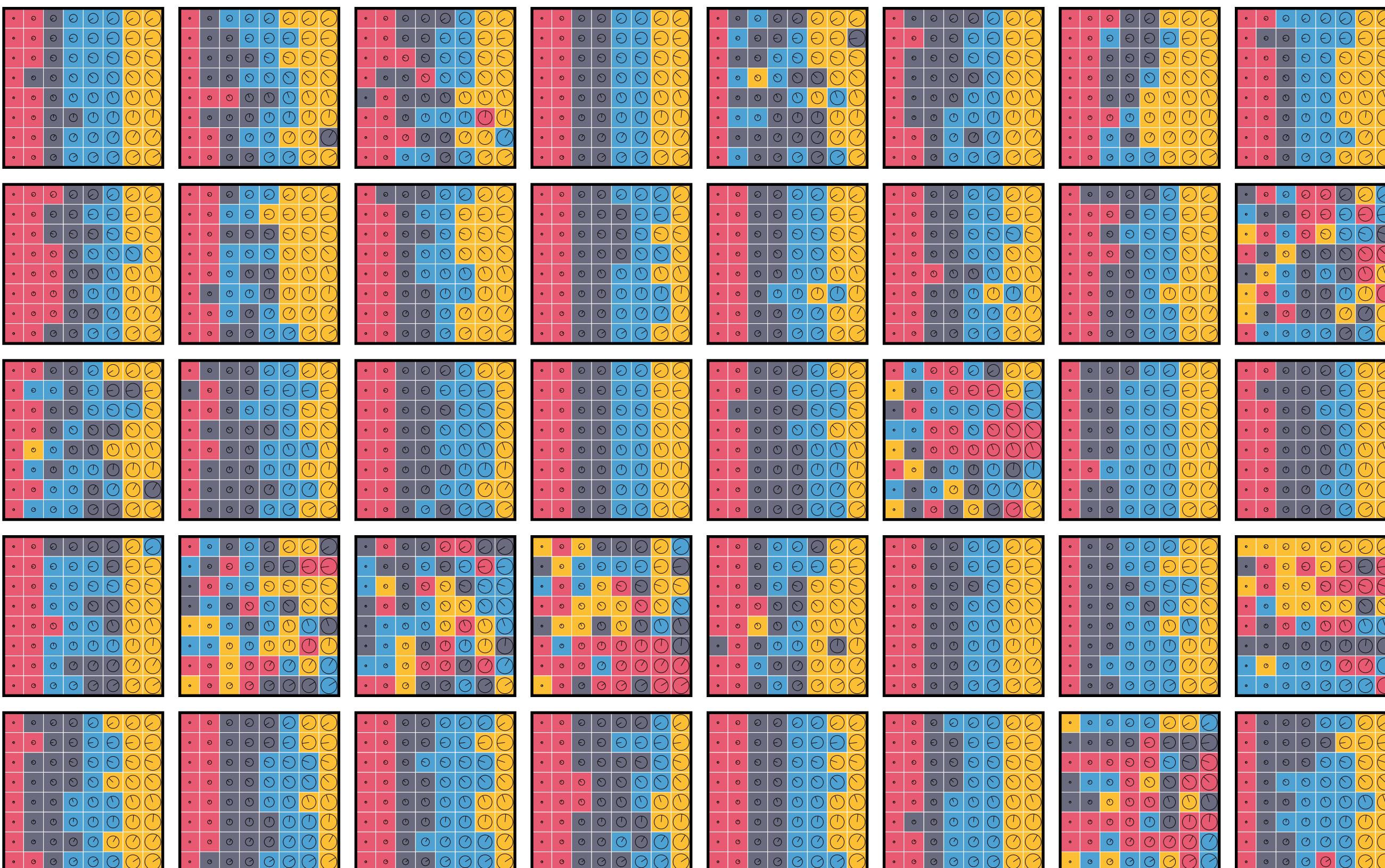
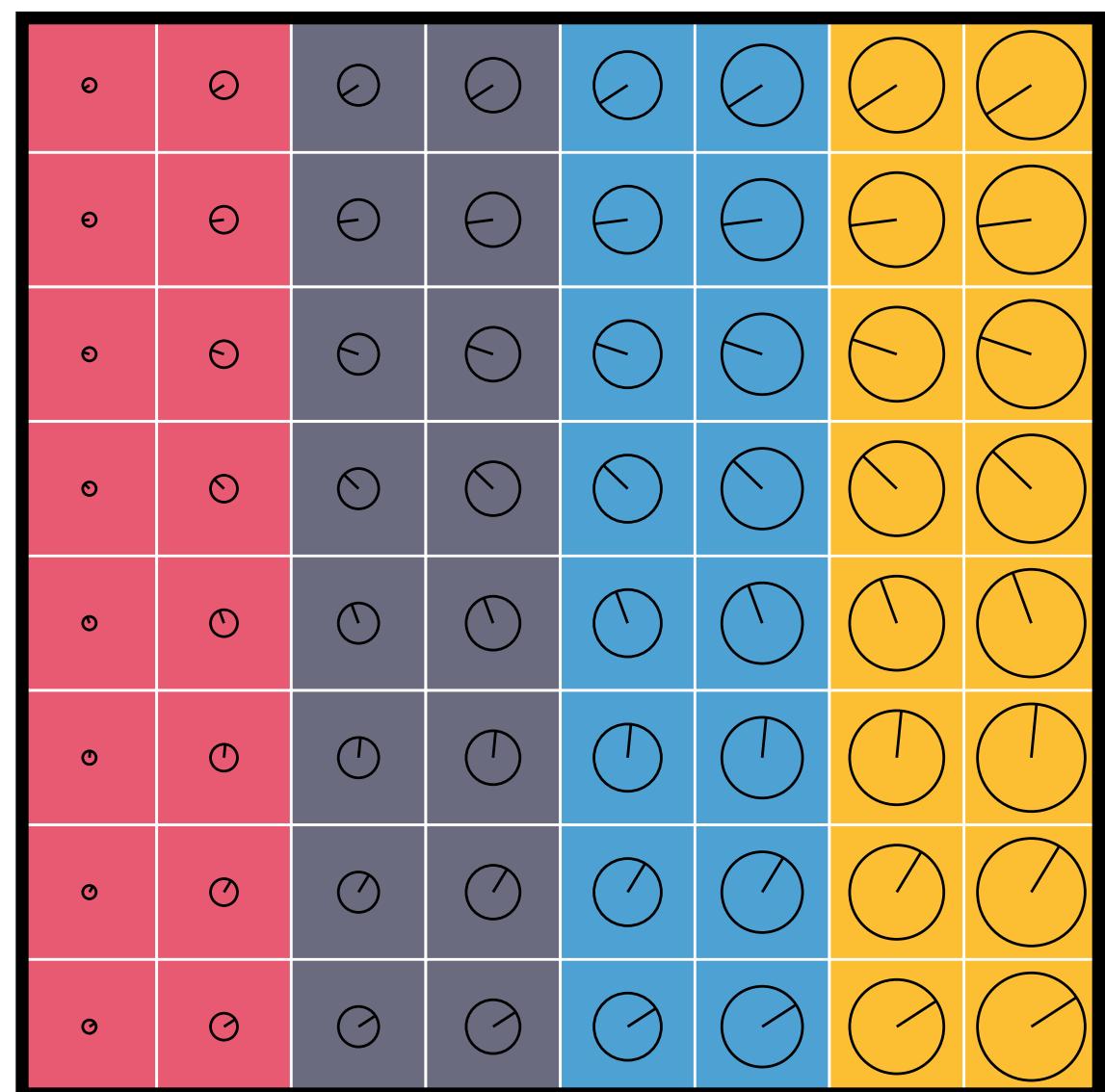
# Results

*Angle-only*



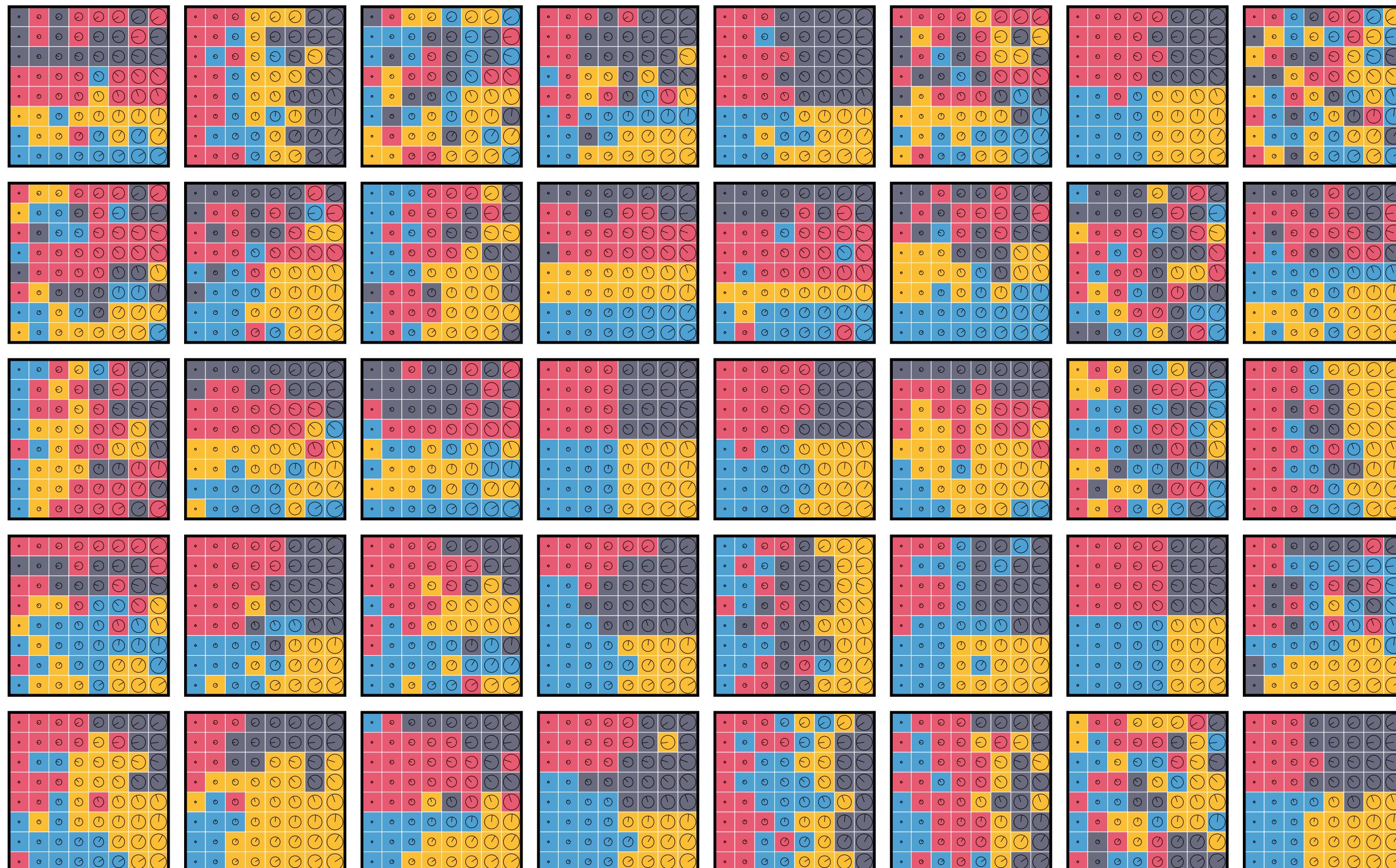
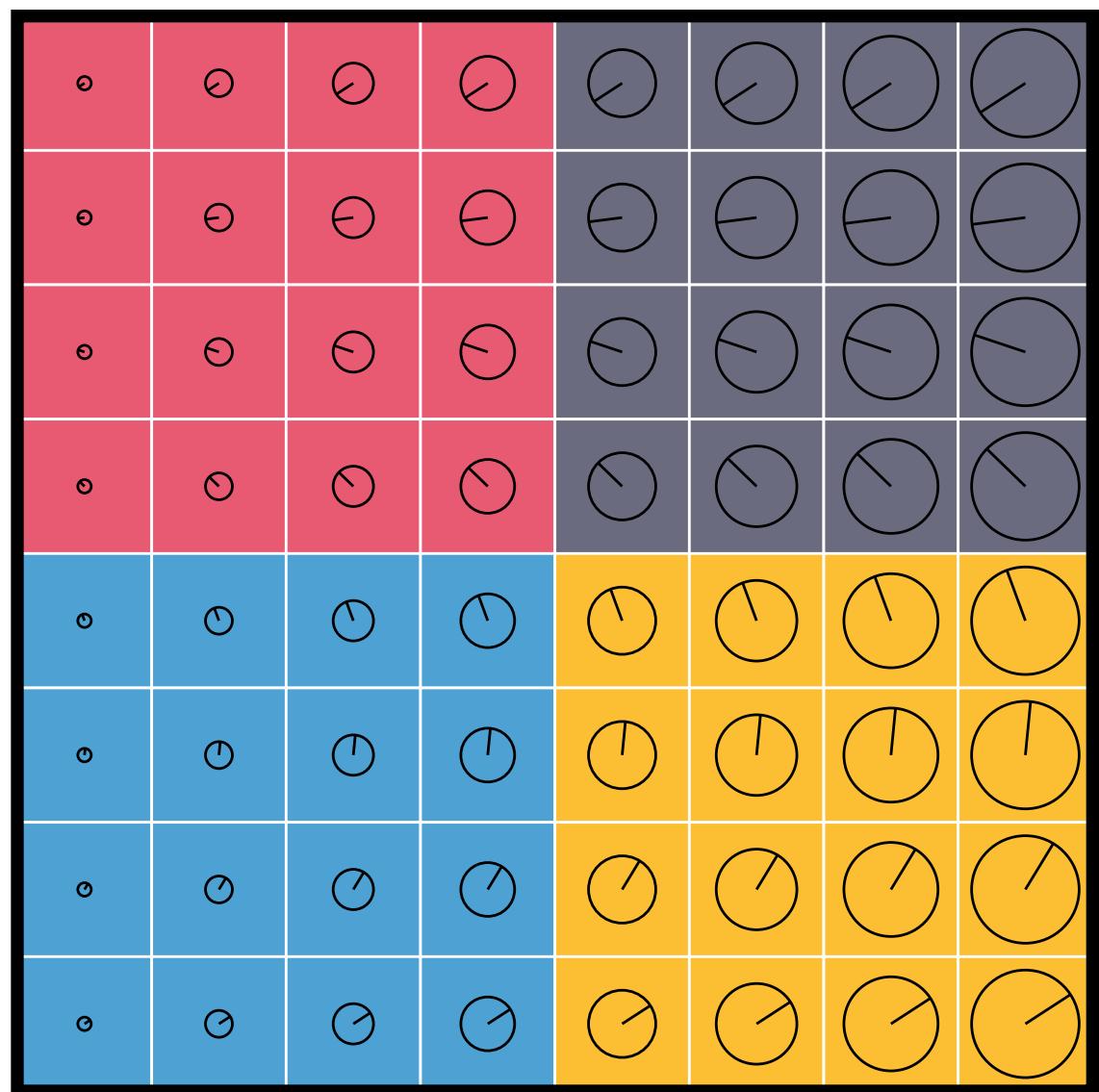
# Results

*Size-only*

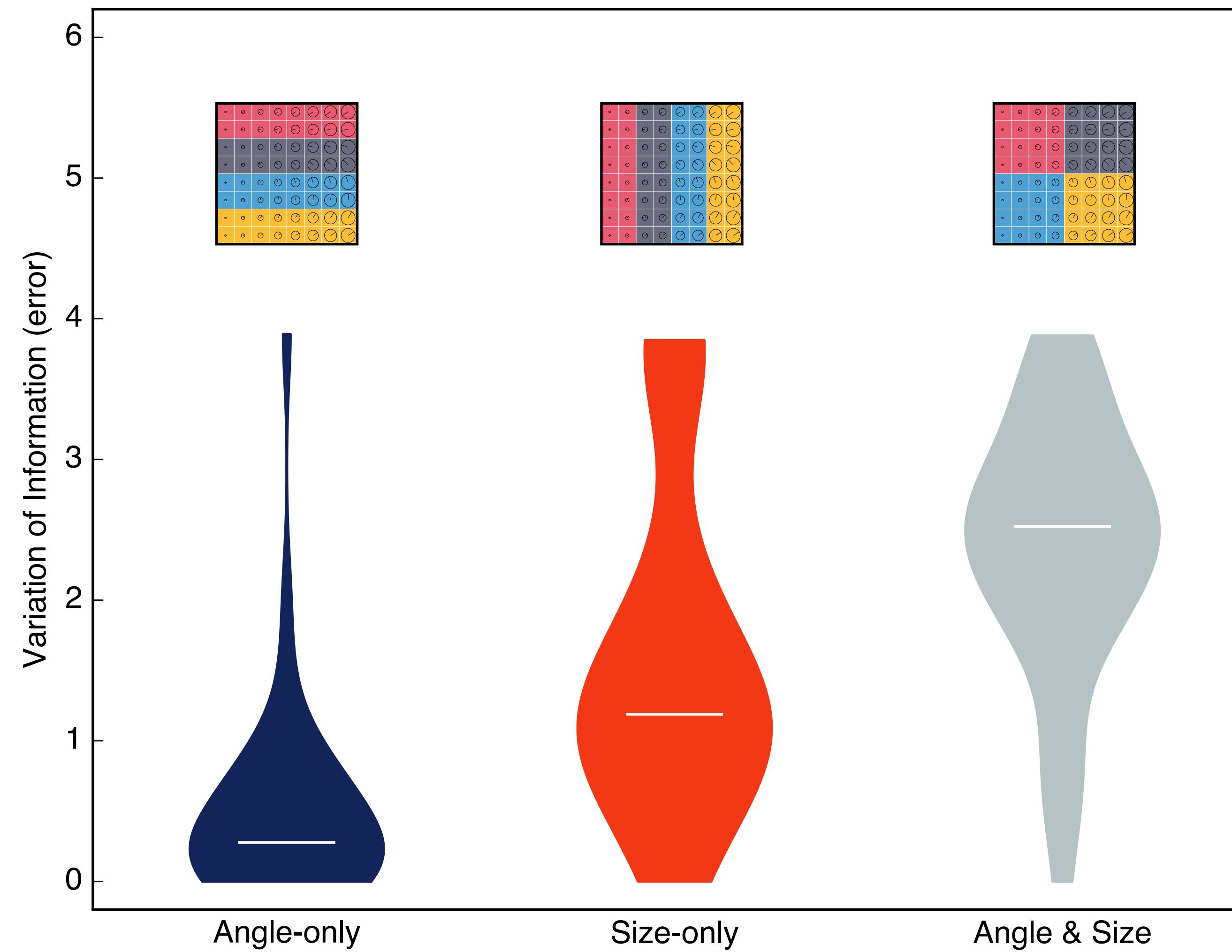


# Results

*Angle & Size*

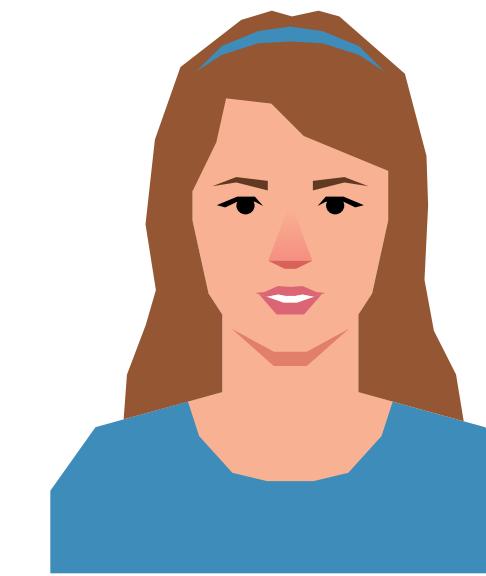
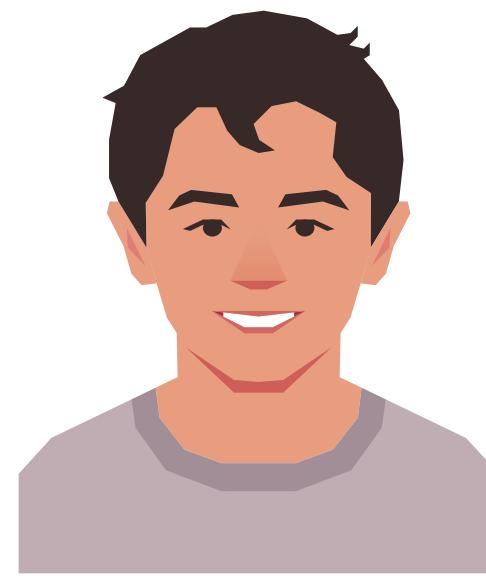
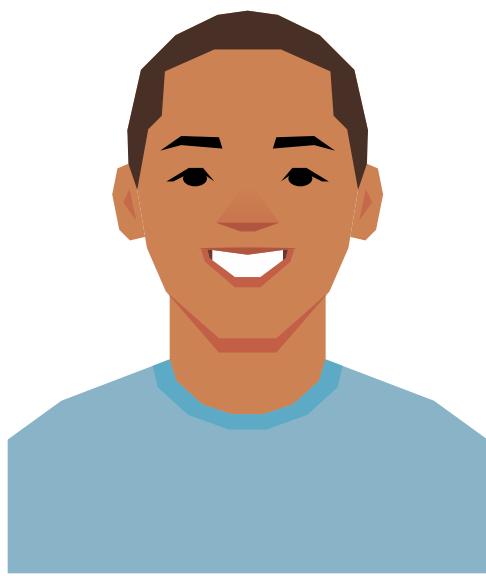


# Result: Learnability advantage for the less informative systems

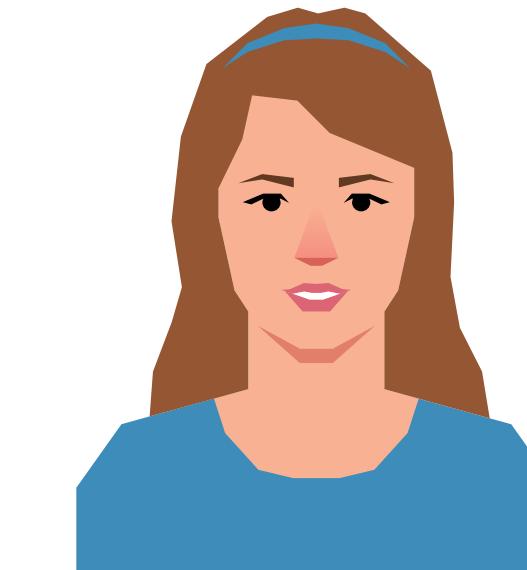
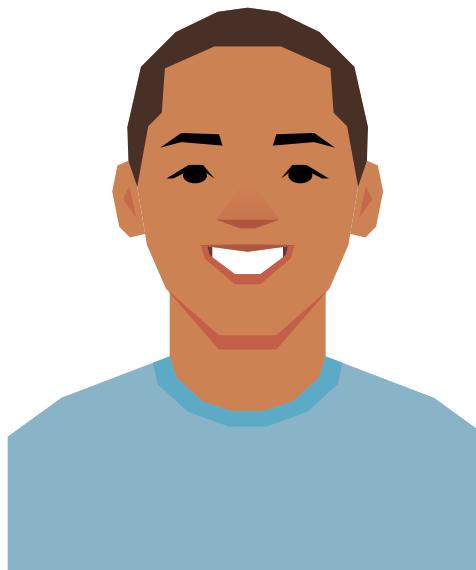
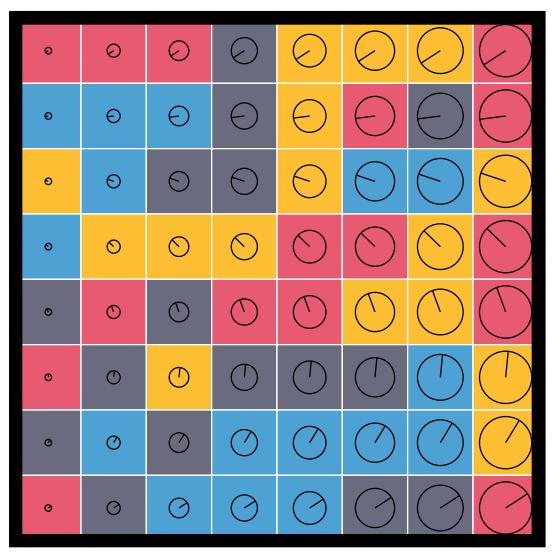


# *Experiment 2*

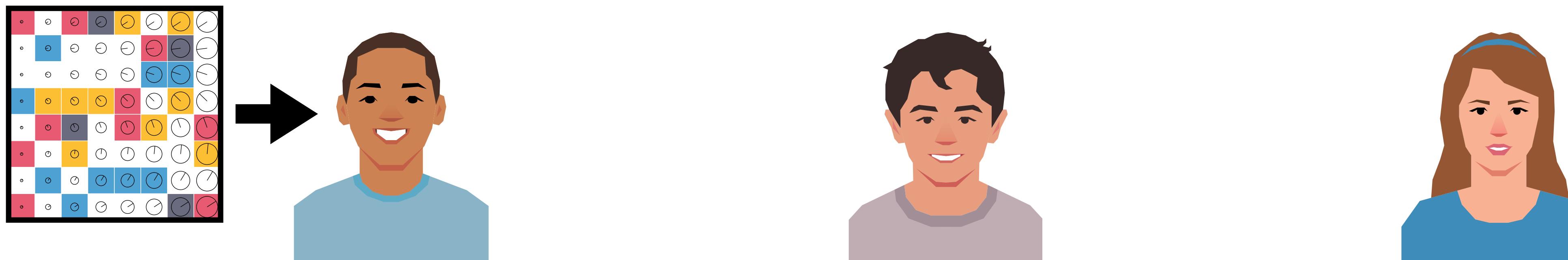
# Iterated learning with humans



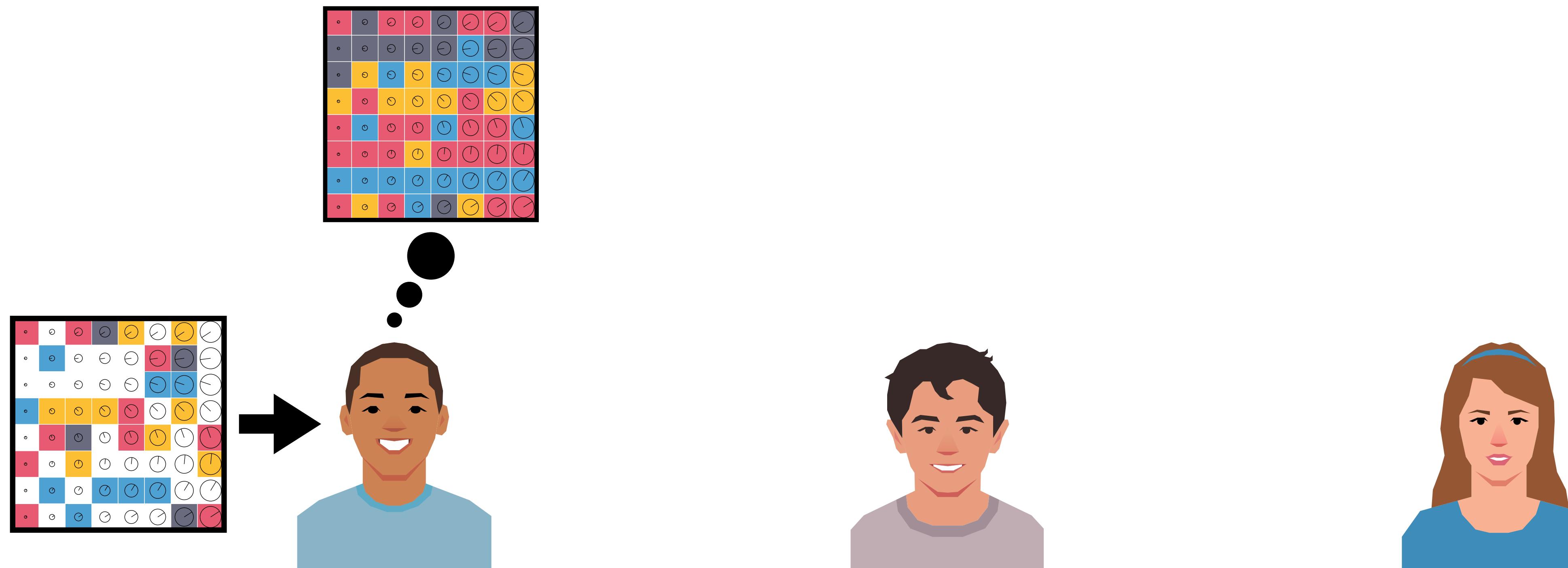
# Iterated learning with humans



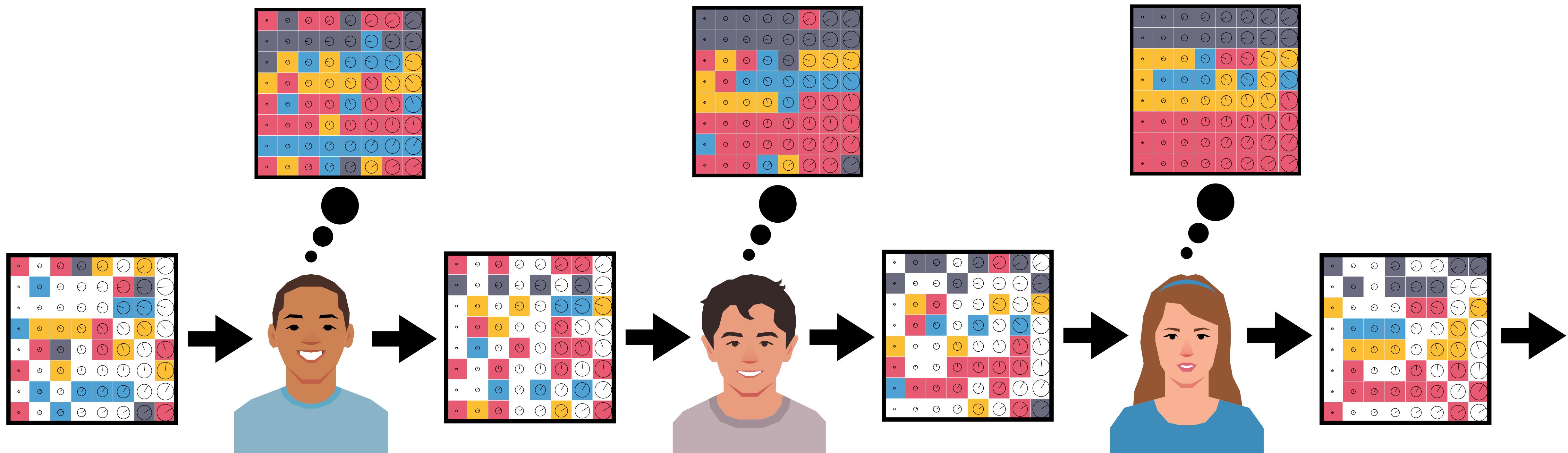
# Iterated learning with humans

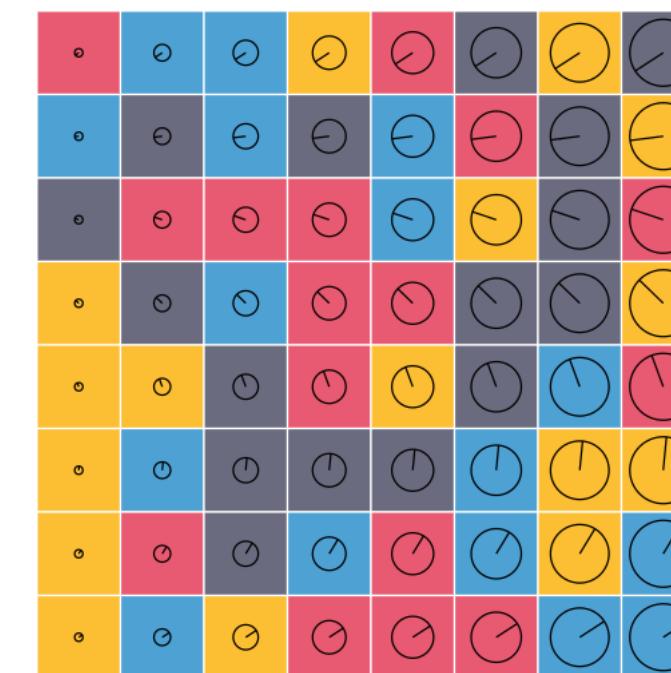
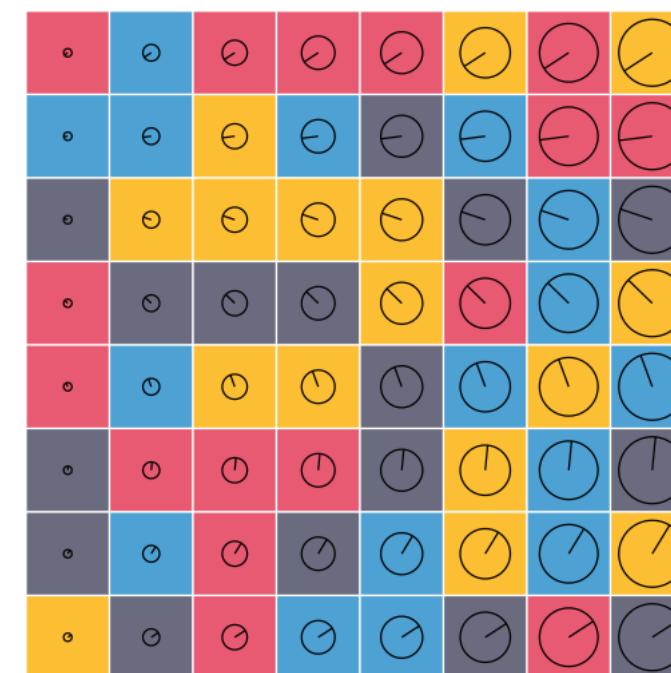
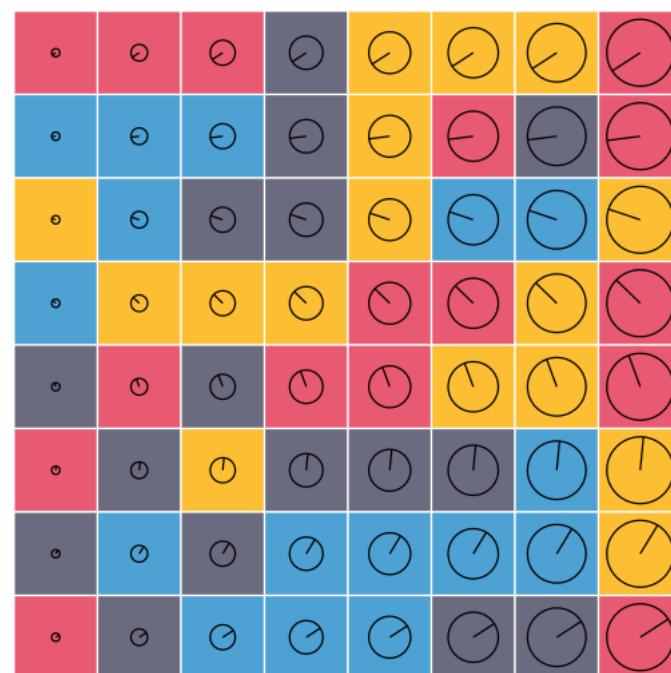
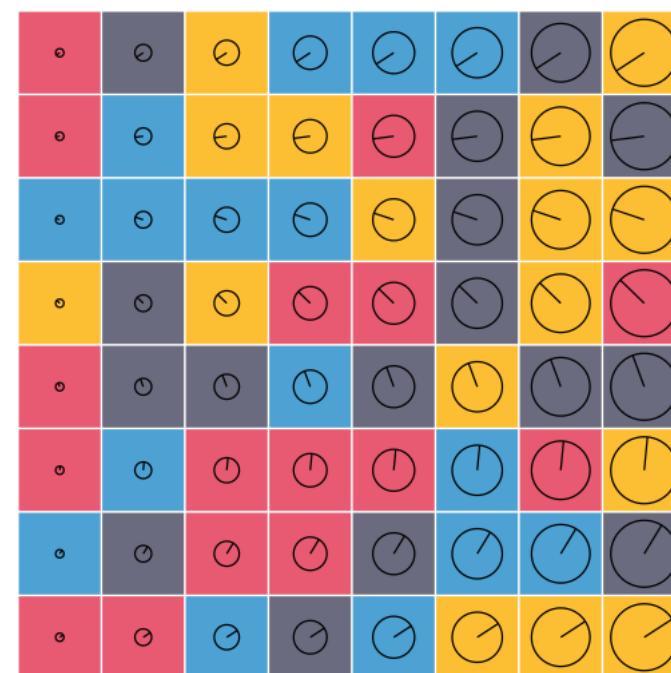
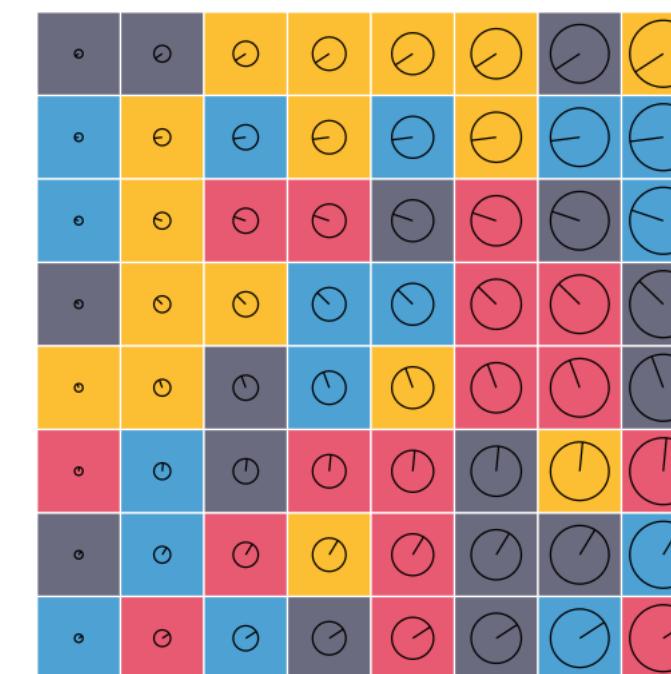
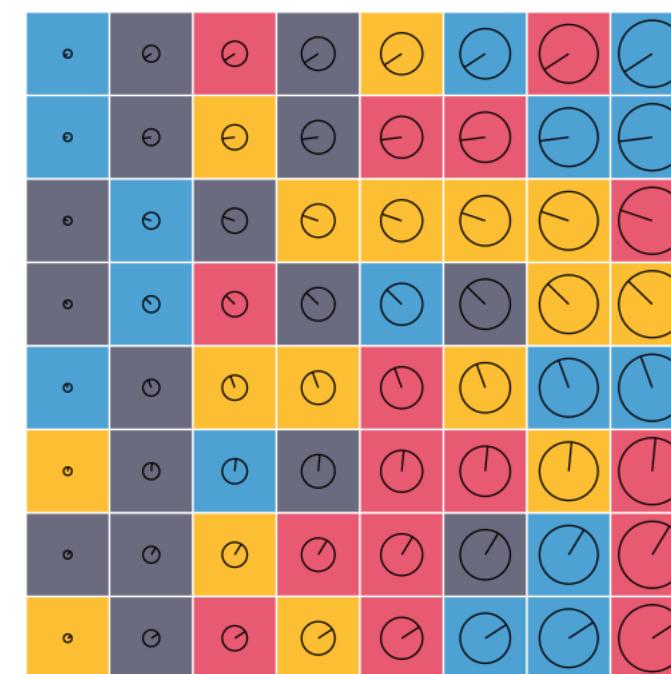
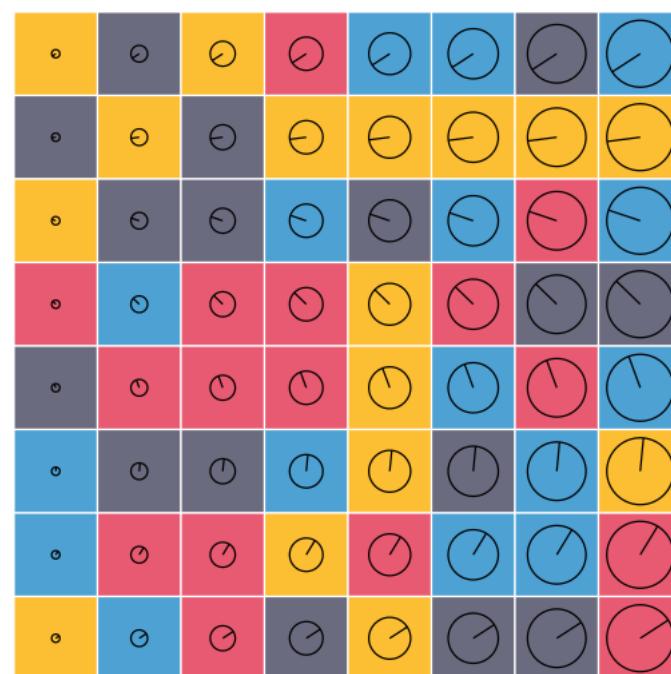
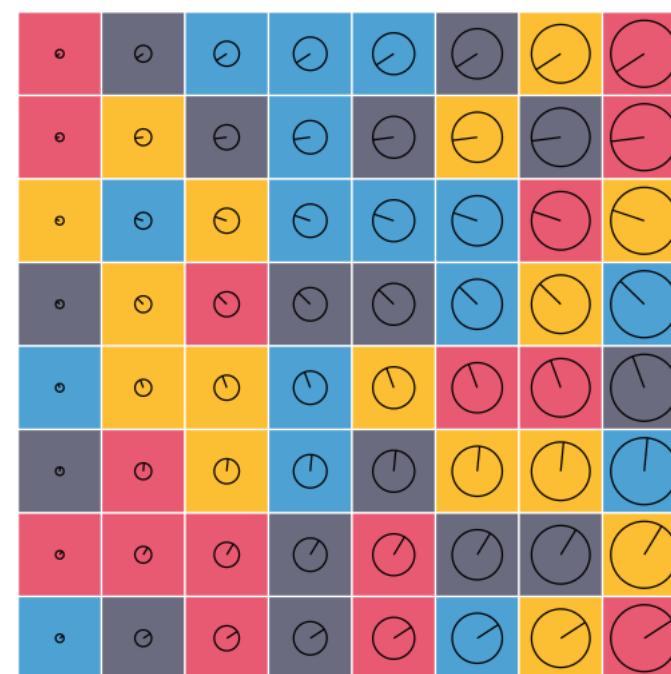
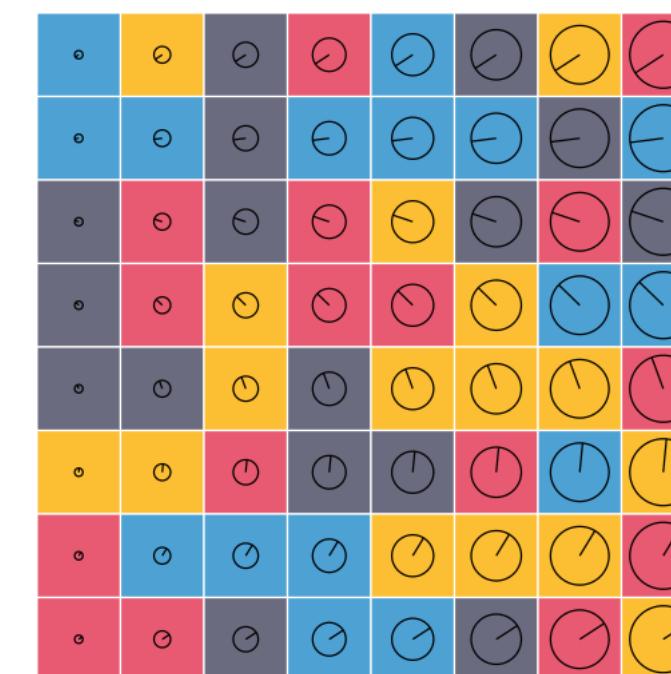
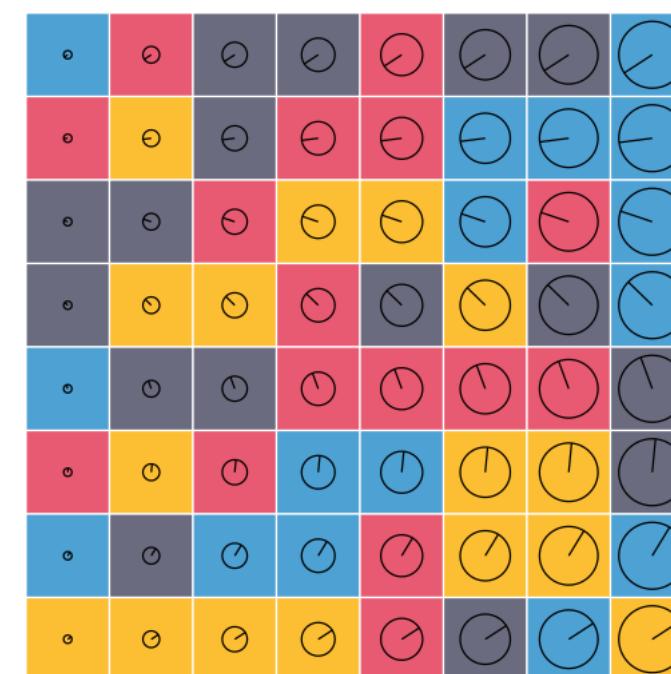
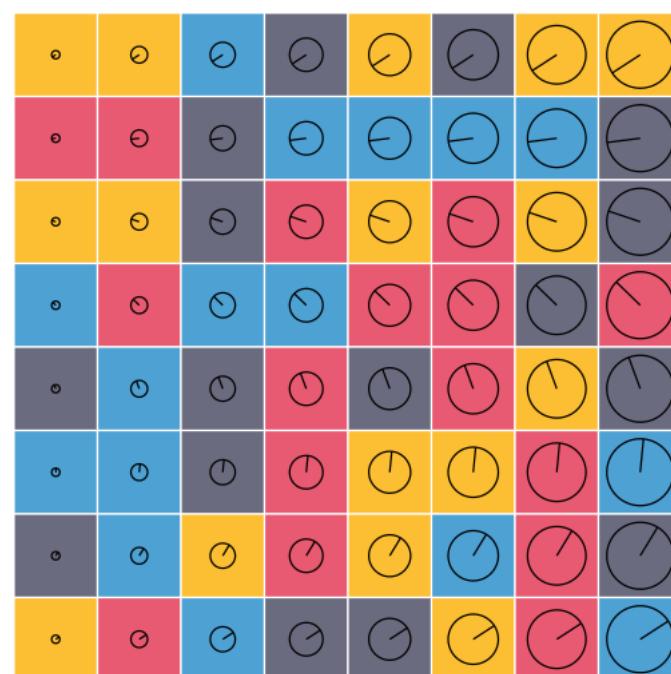
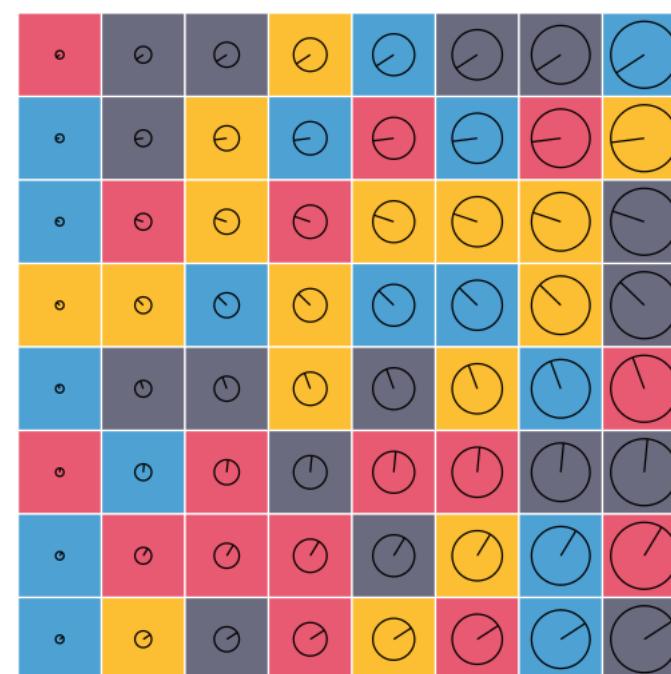


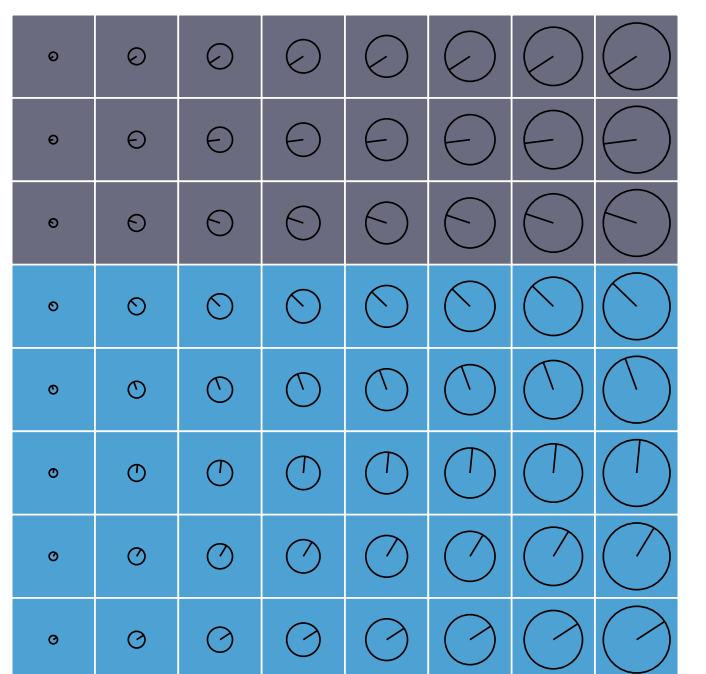
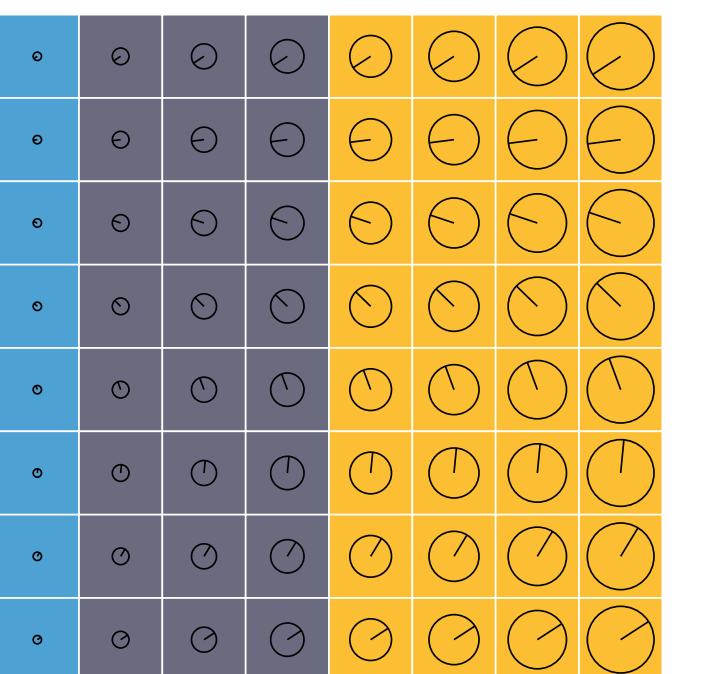
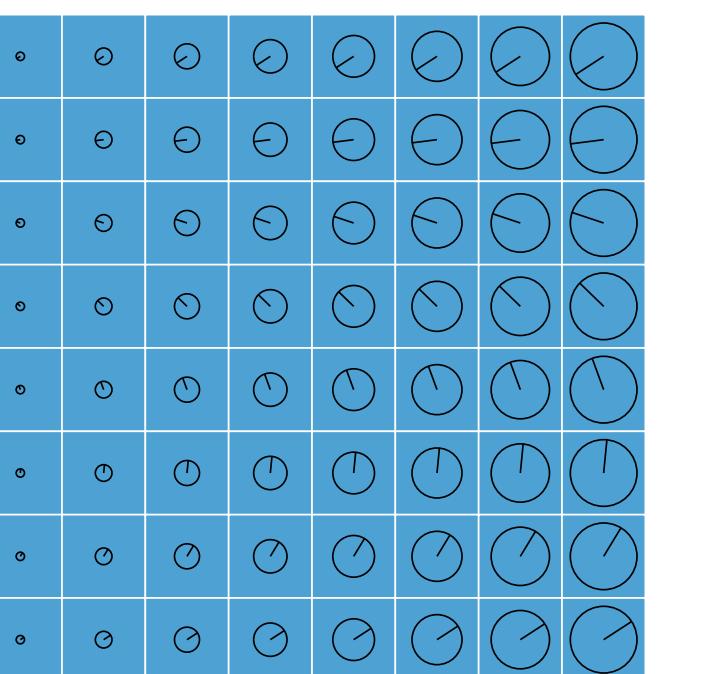
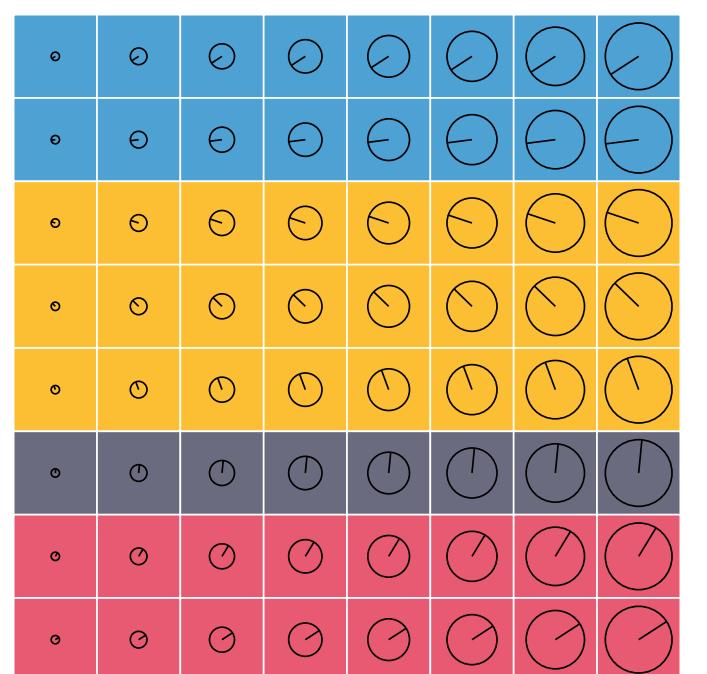
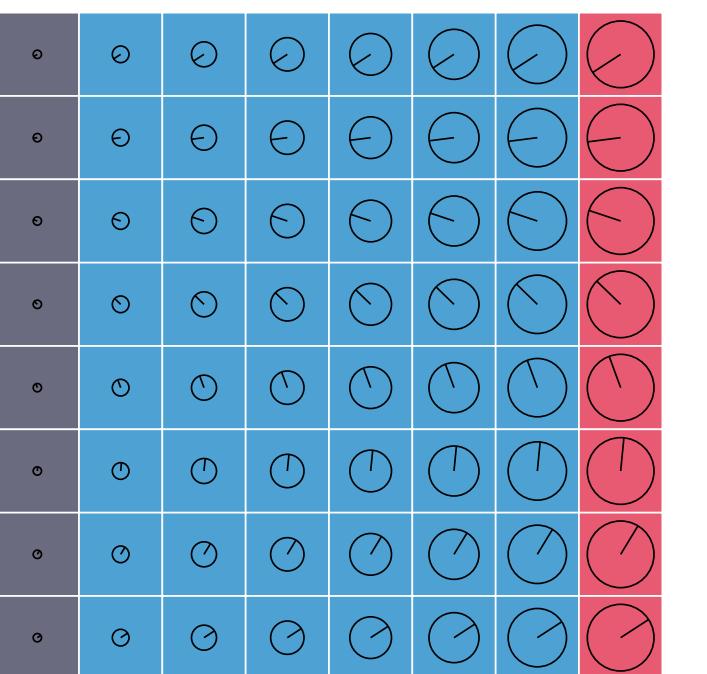
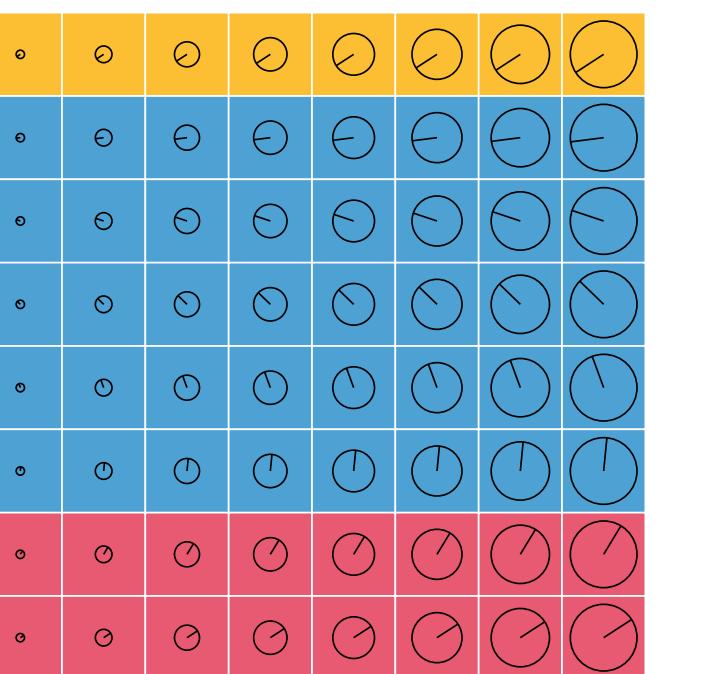
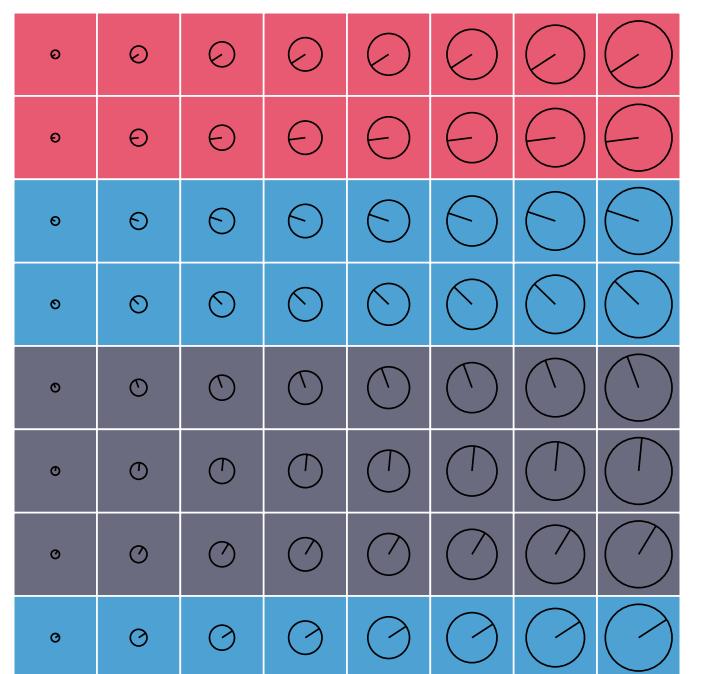
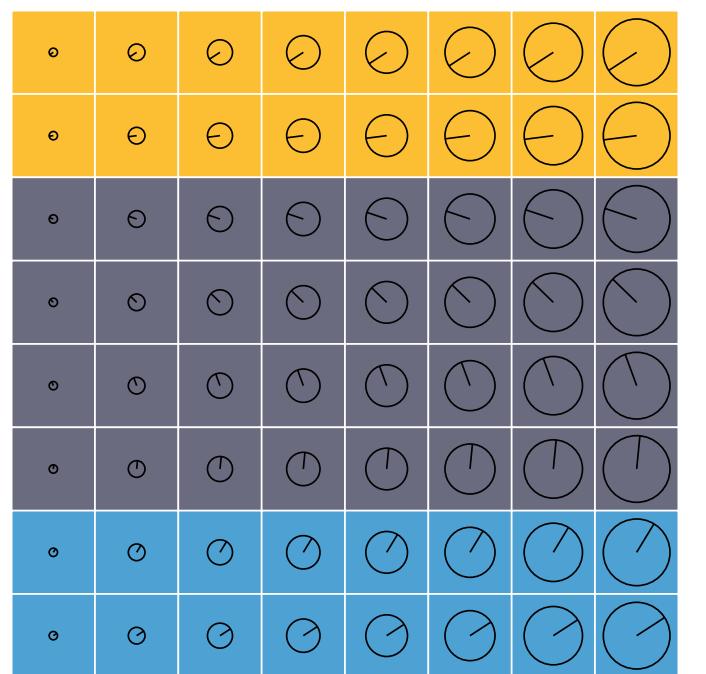
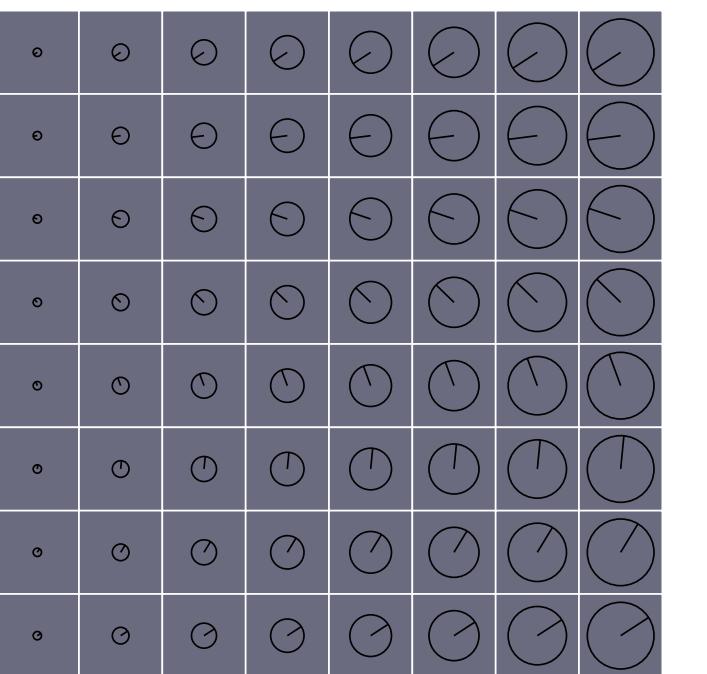
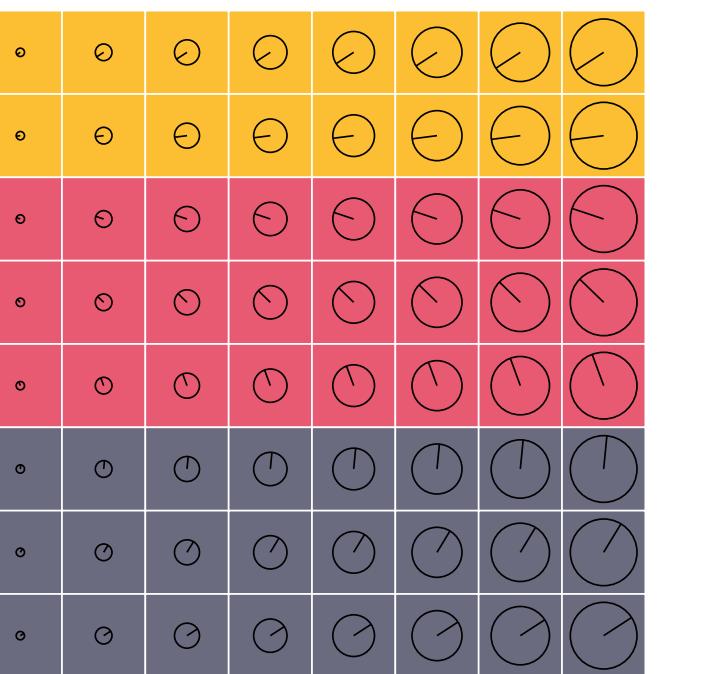
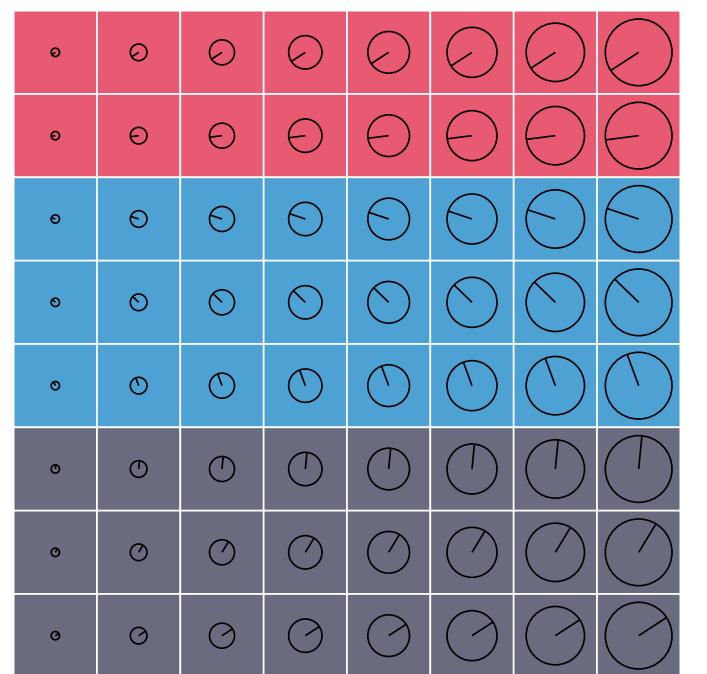
# Iterated learning with humans



# Iterated learning with humans

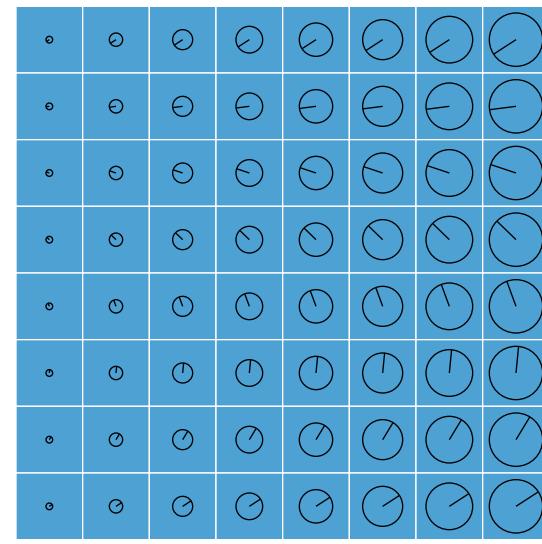
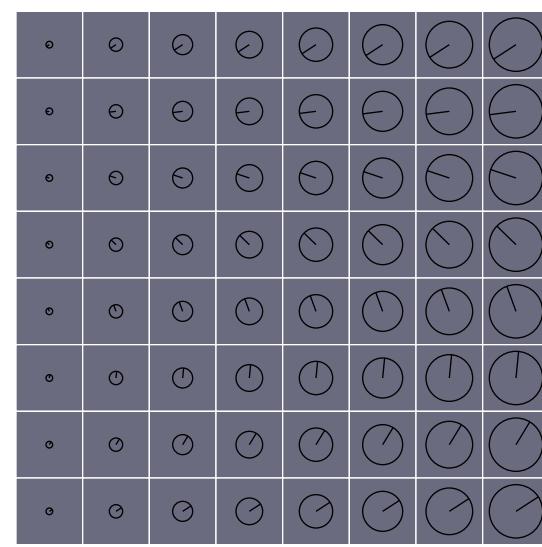




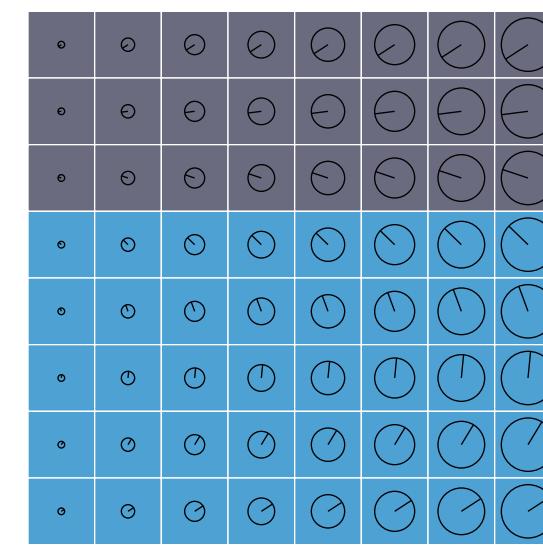


# Human results (first convergence)

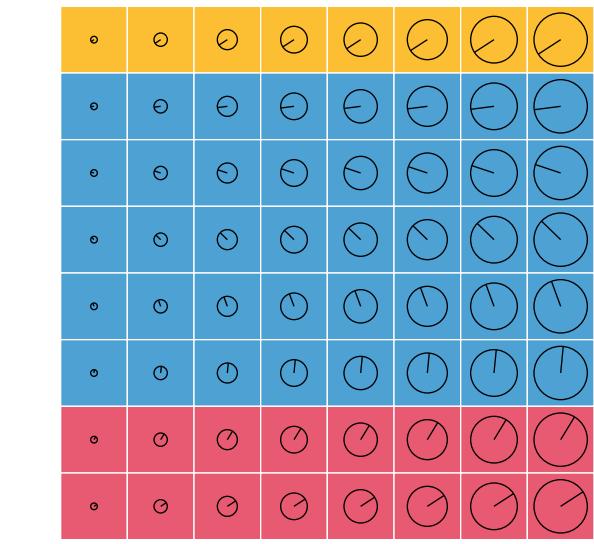
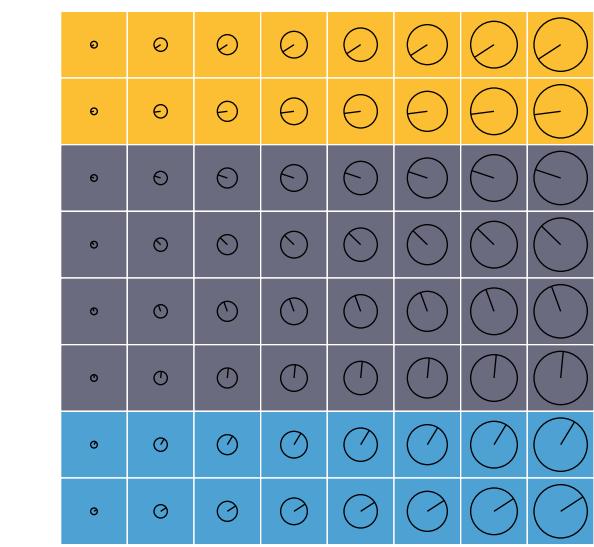
1 concept (2/12)



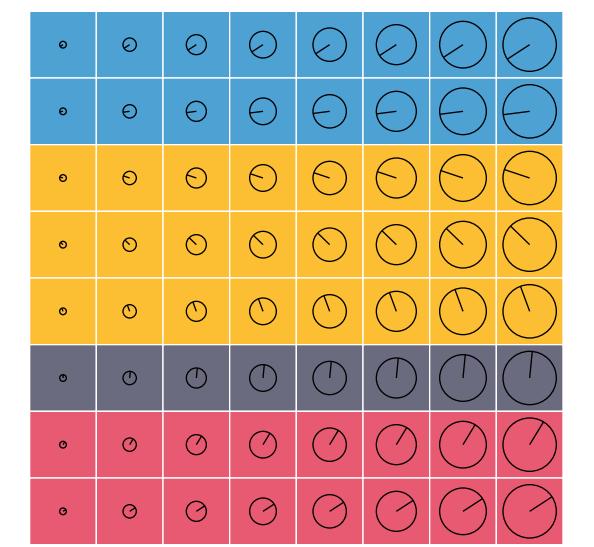
2 concepts (1/12)



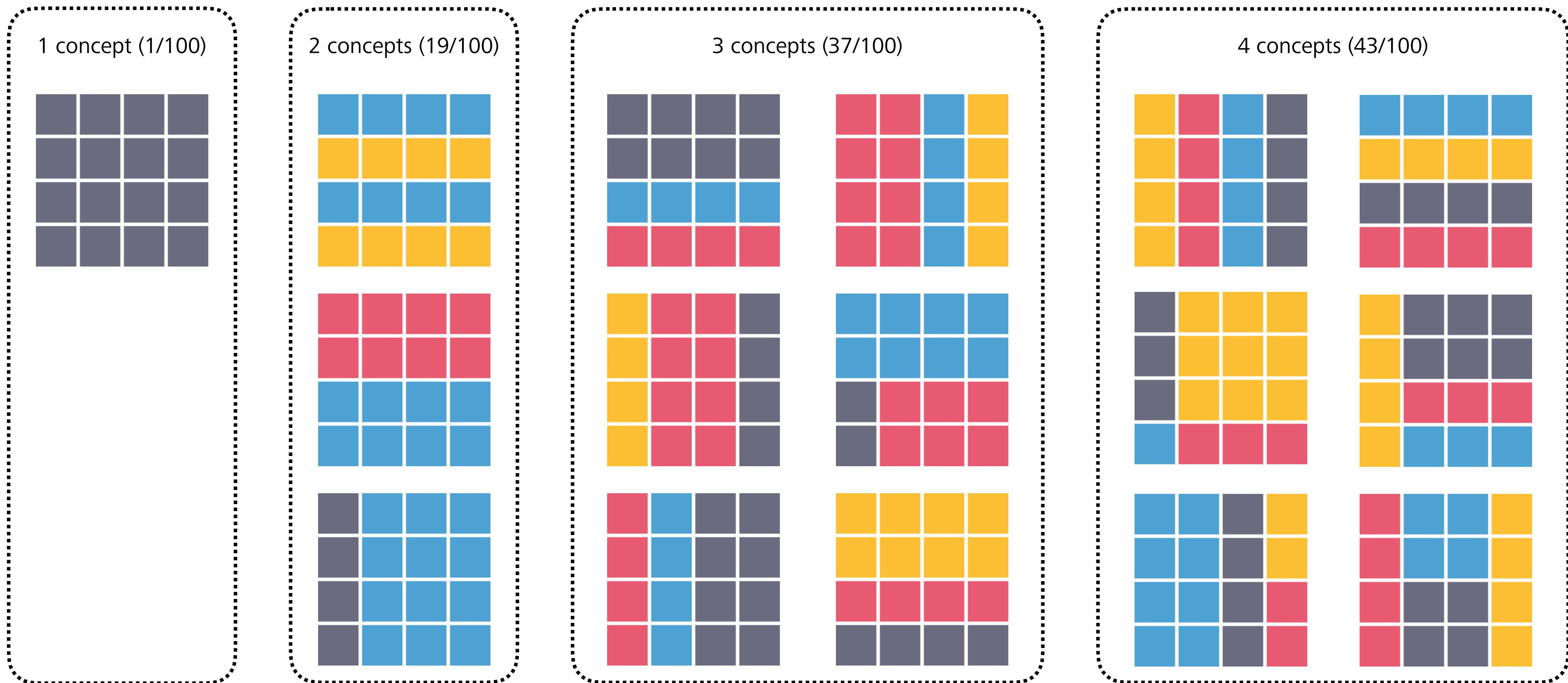
3 concepts (8/12)



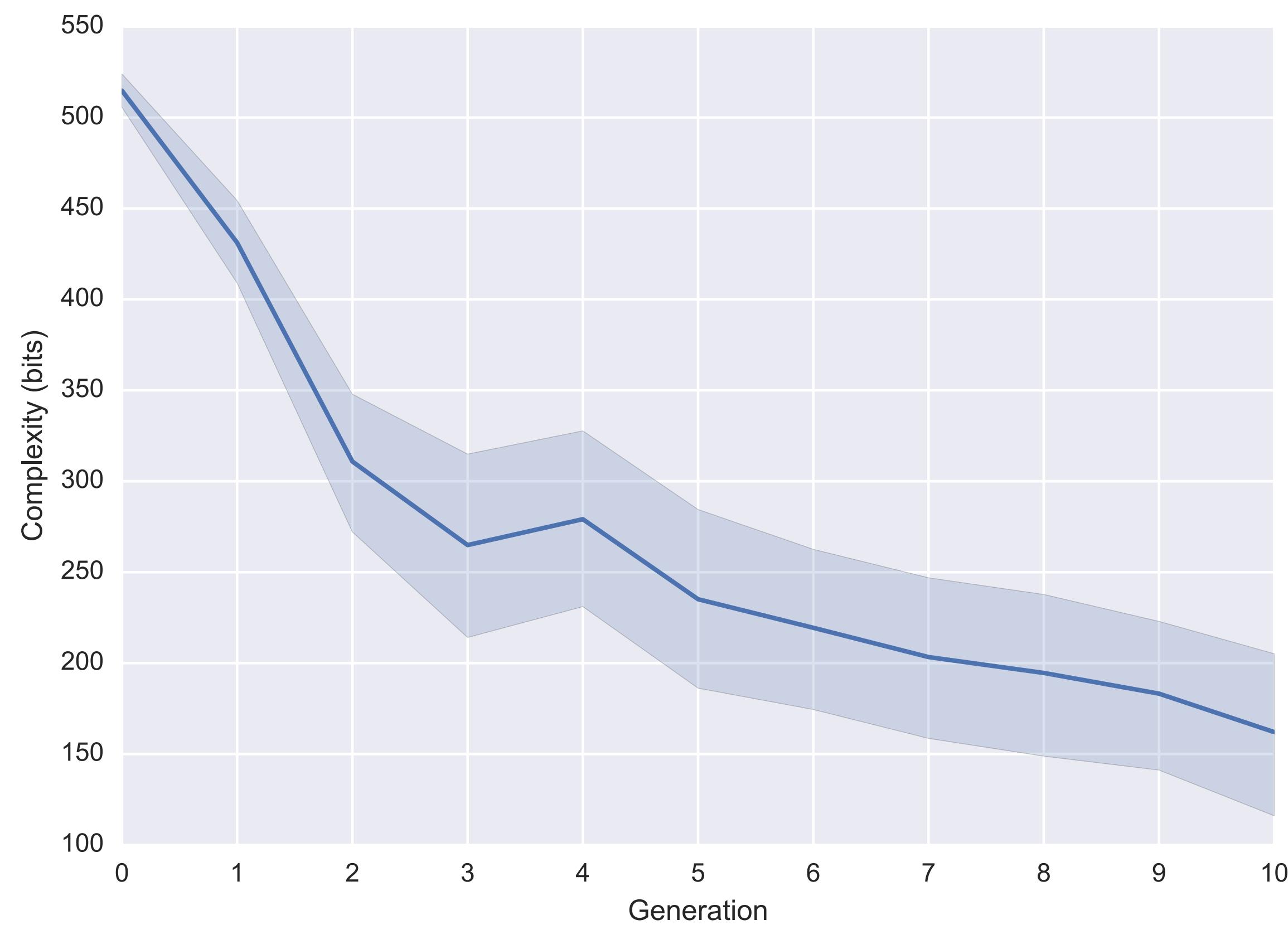
4 concepts (1/12)



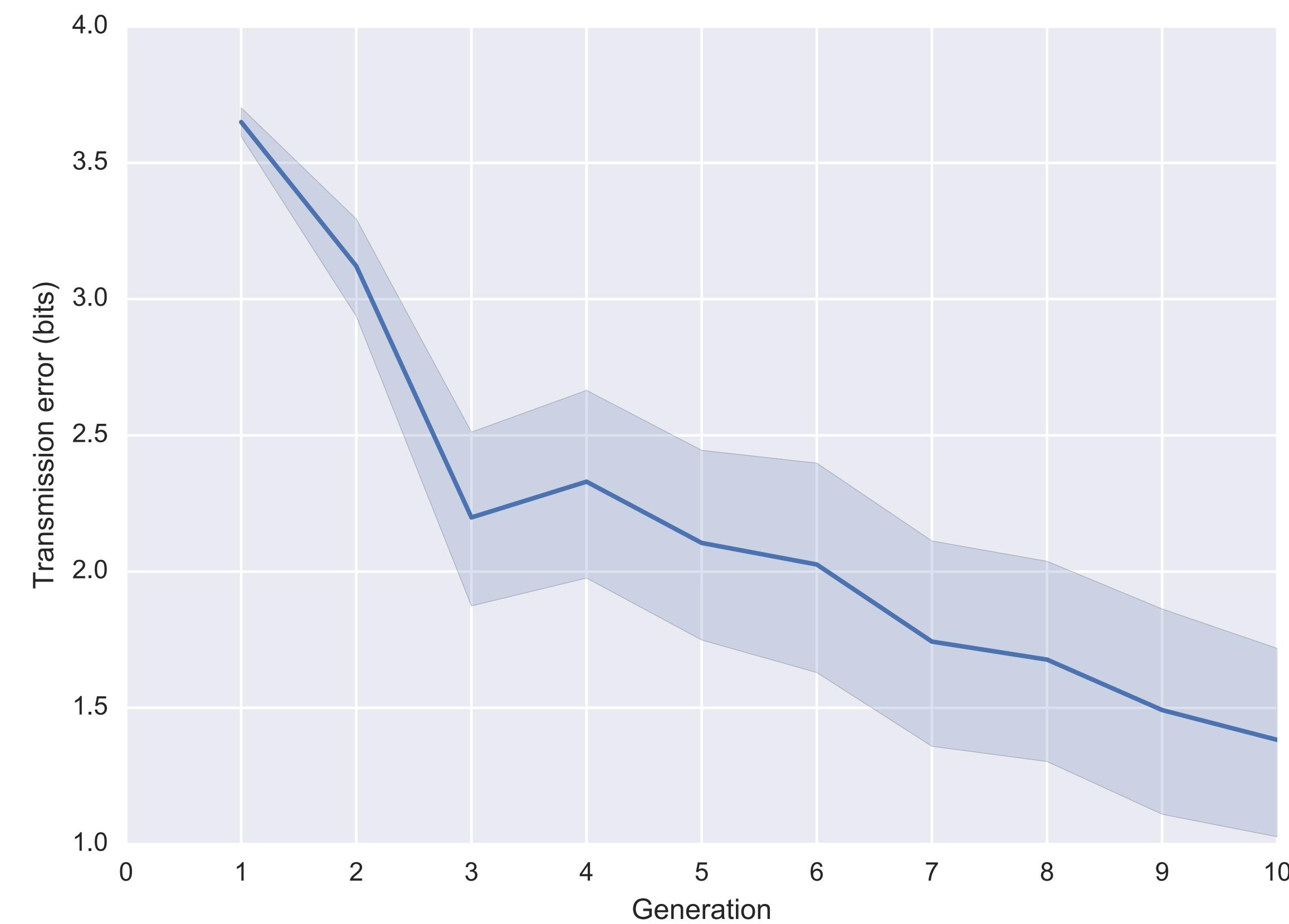
# Model results (first convergence)



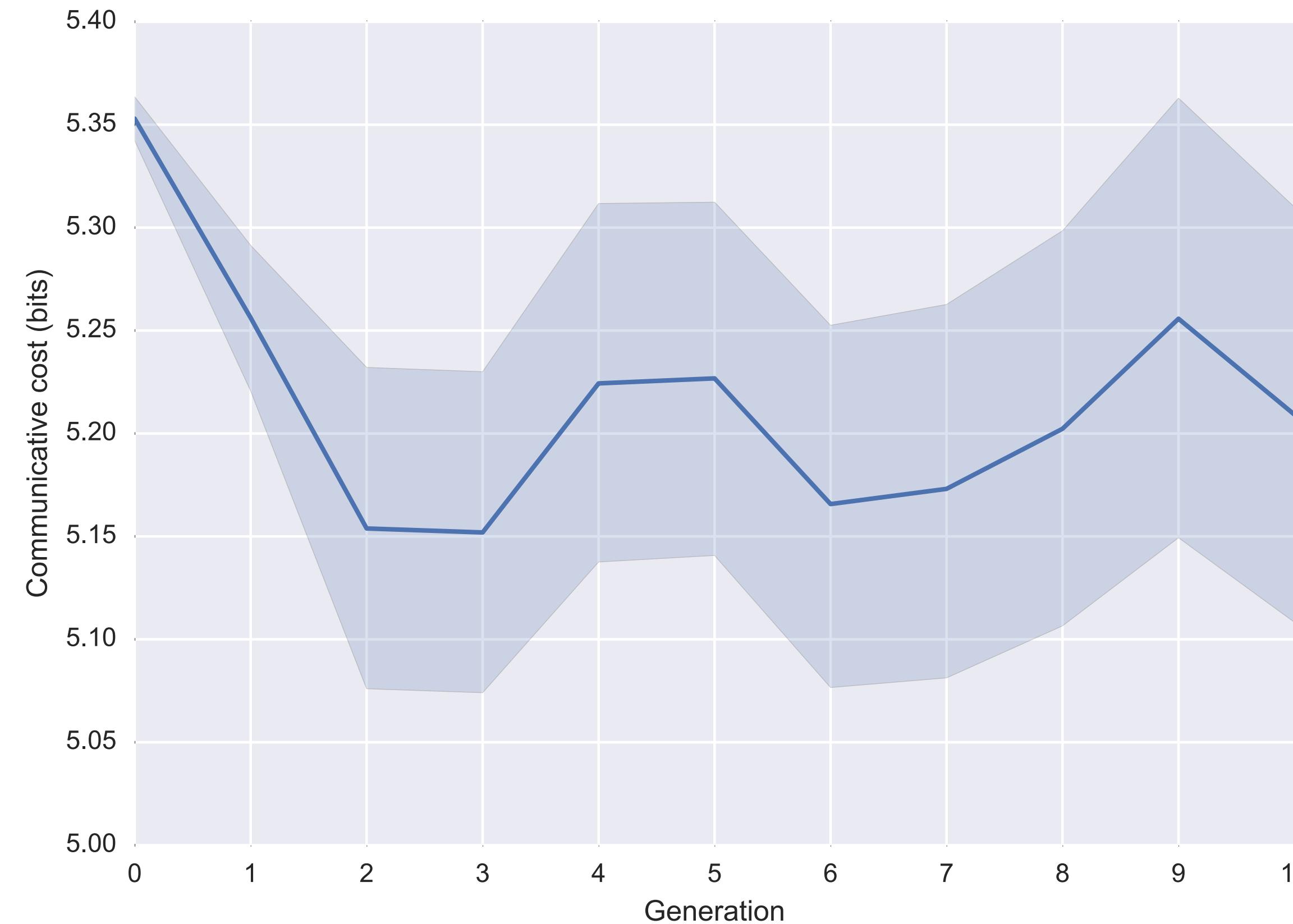
# Simplicity



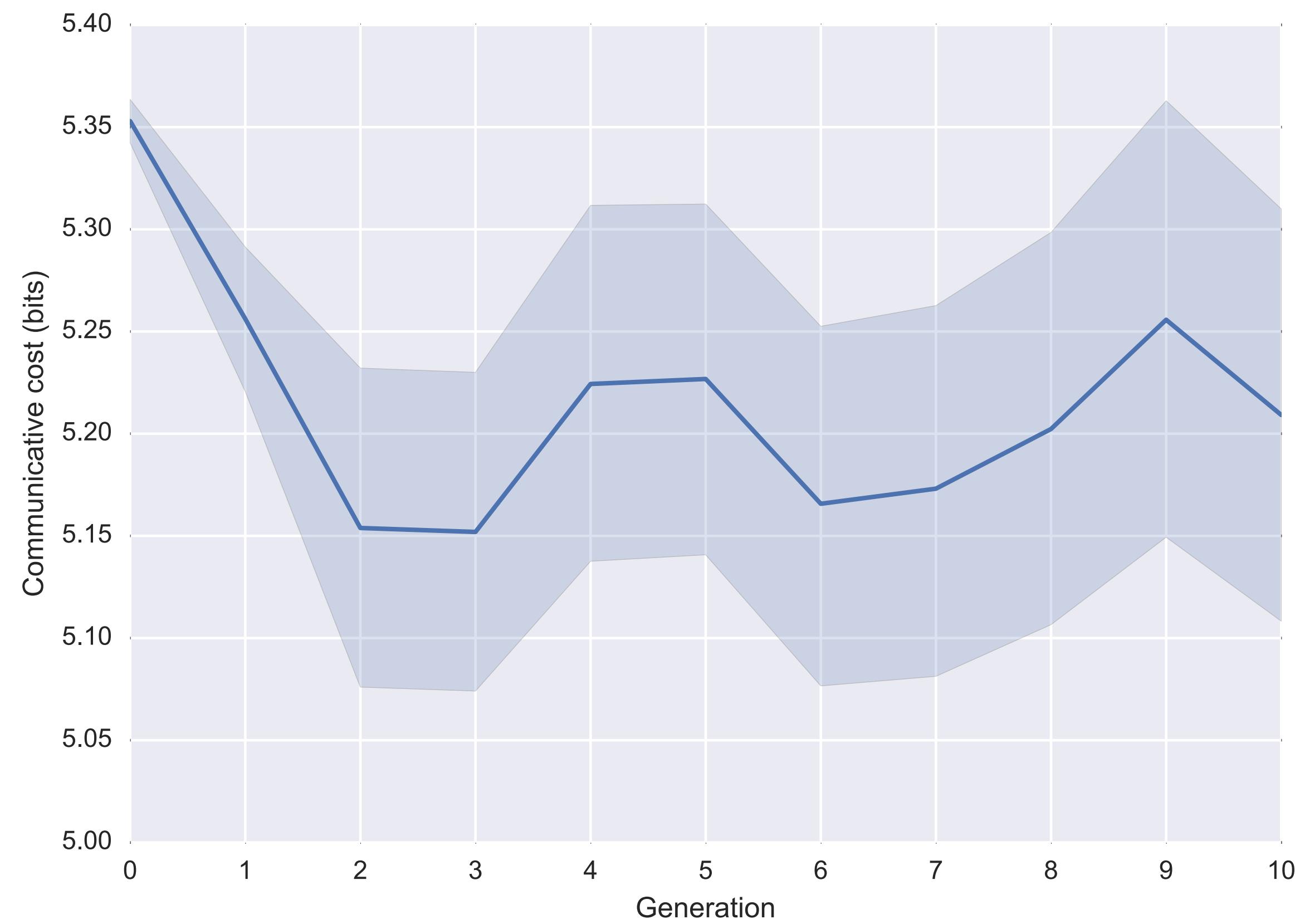
# Learnability



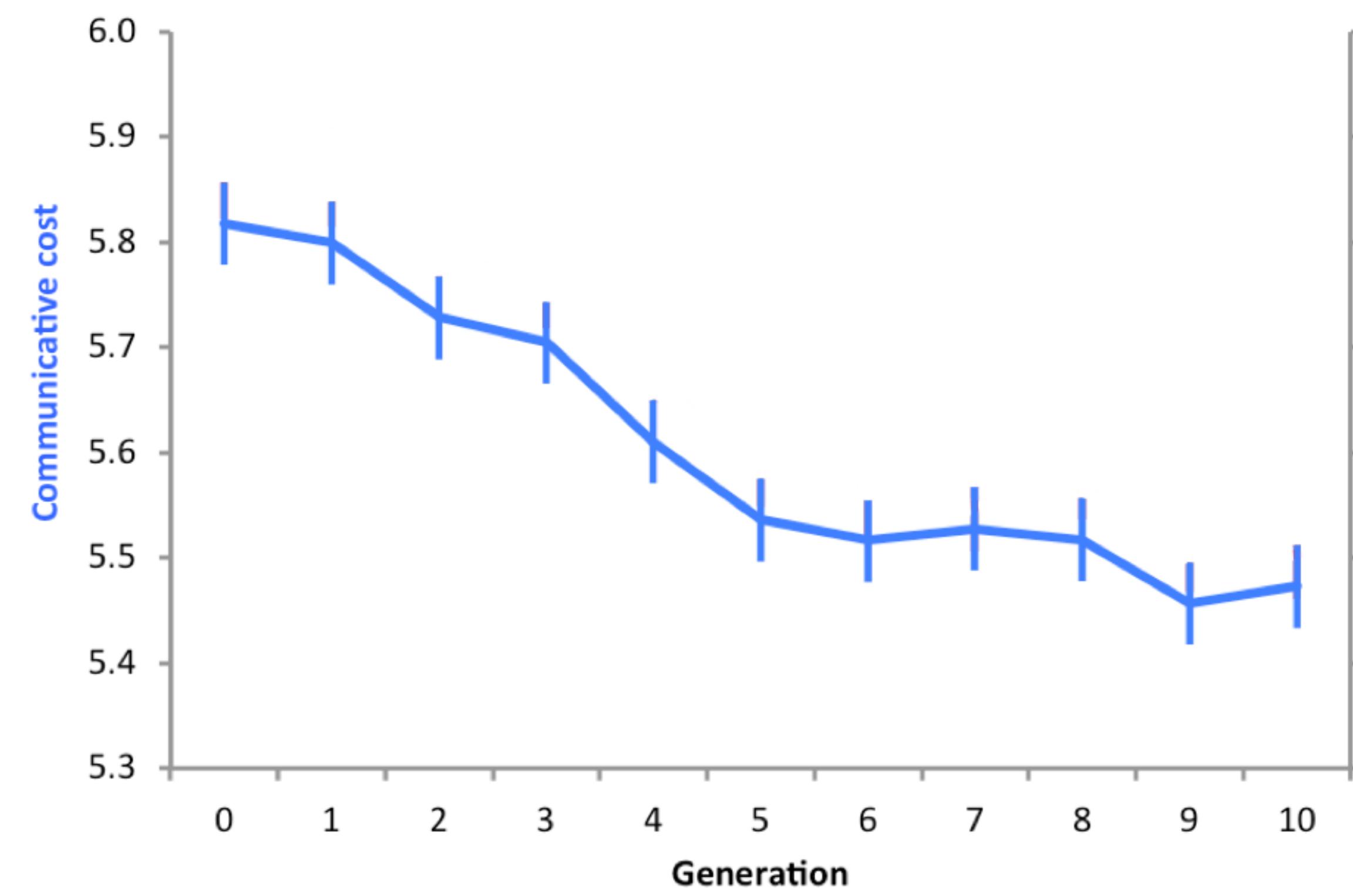
# Informativeness



# Informativeness

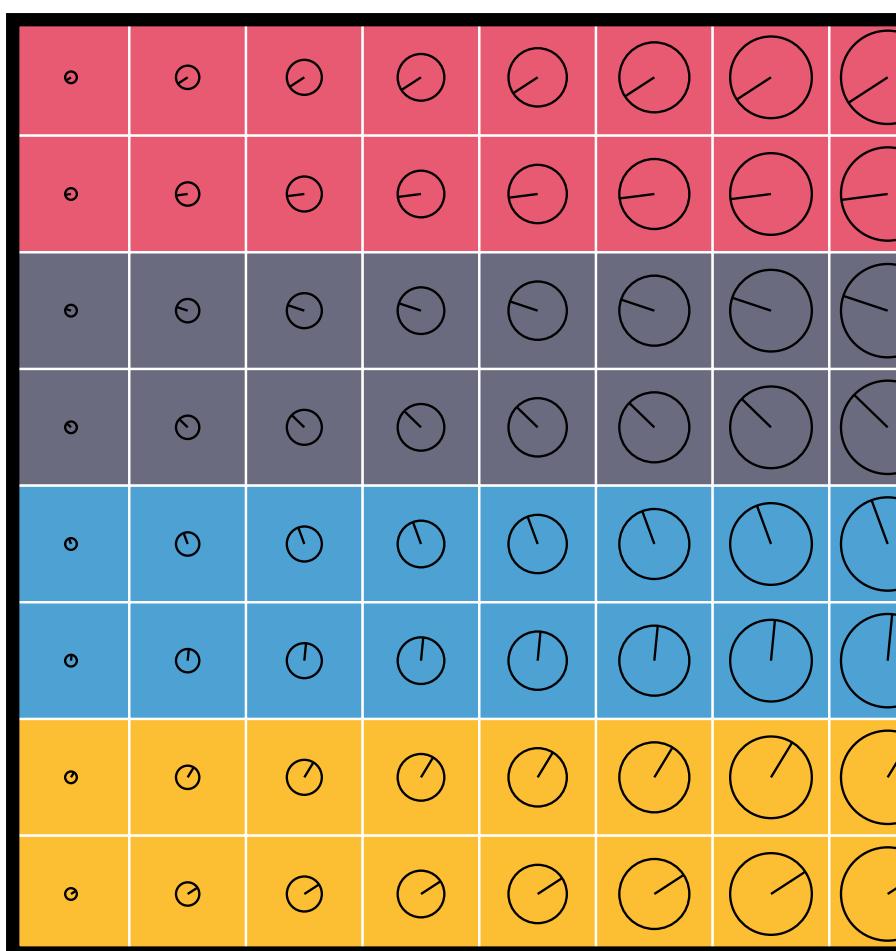
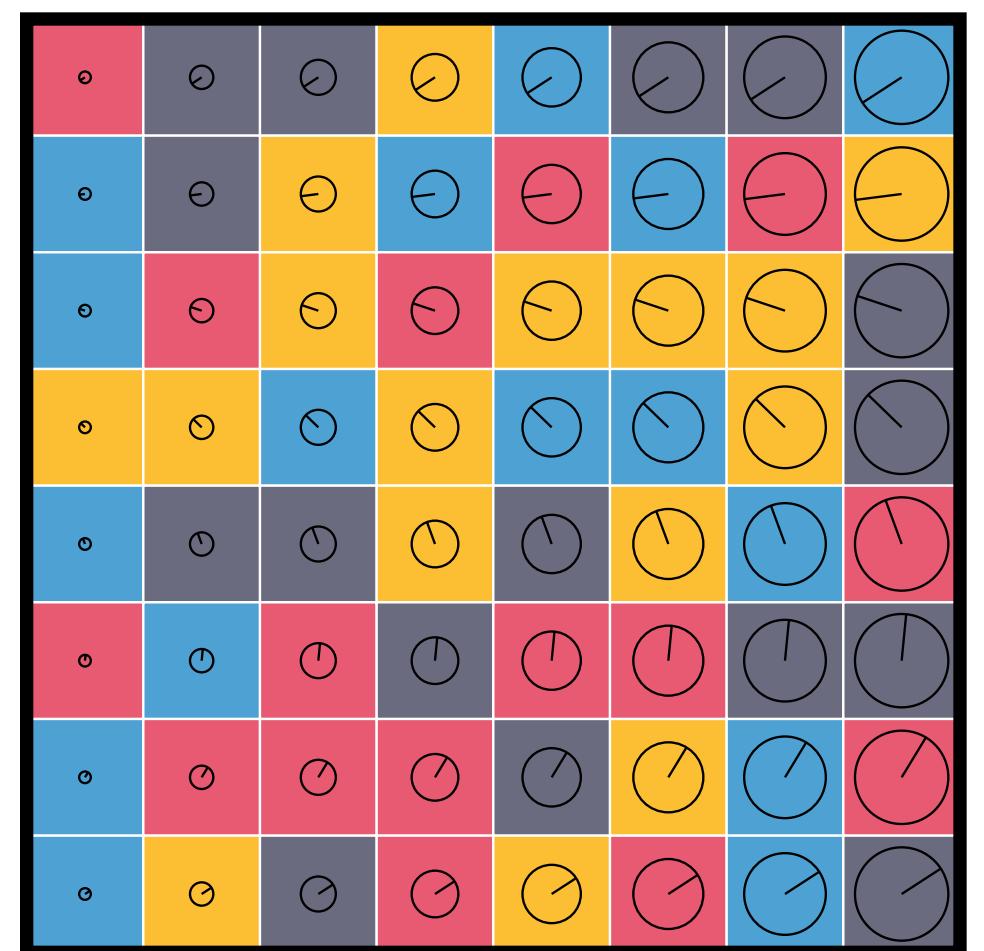


Carstensen et al.



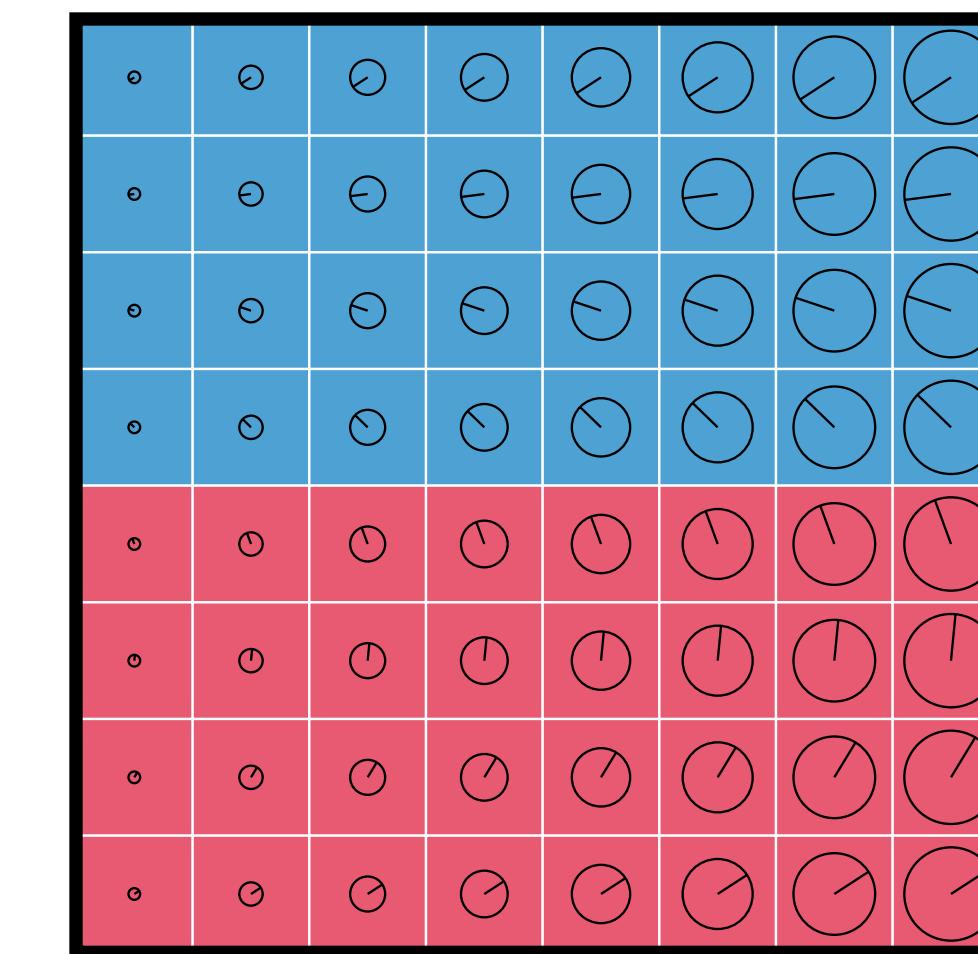
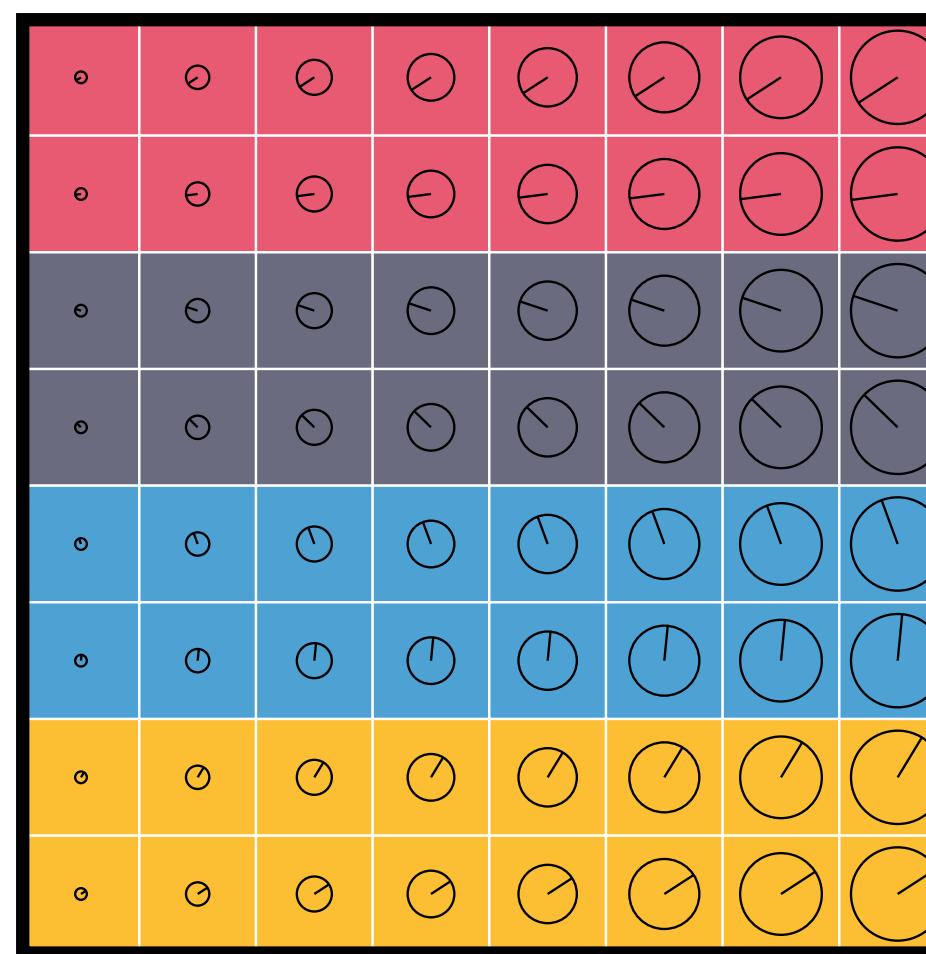
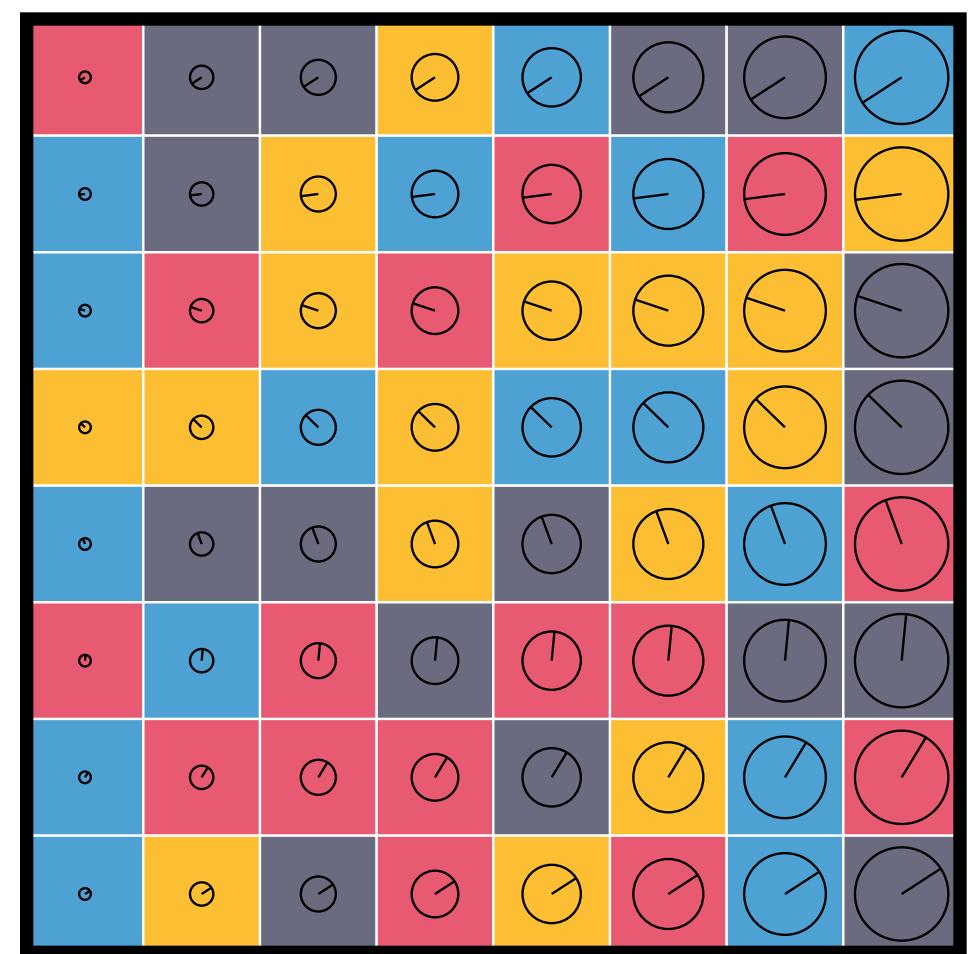
# Two ways of achieving simplicity

Increase in convexity



# Two ways of achieving simplicity

Increase in convexity

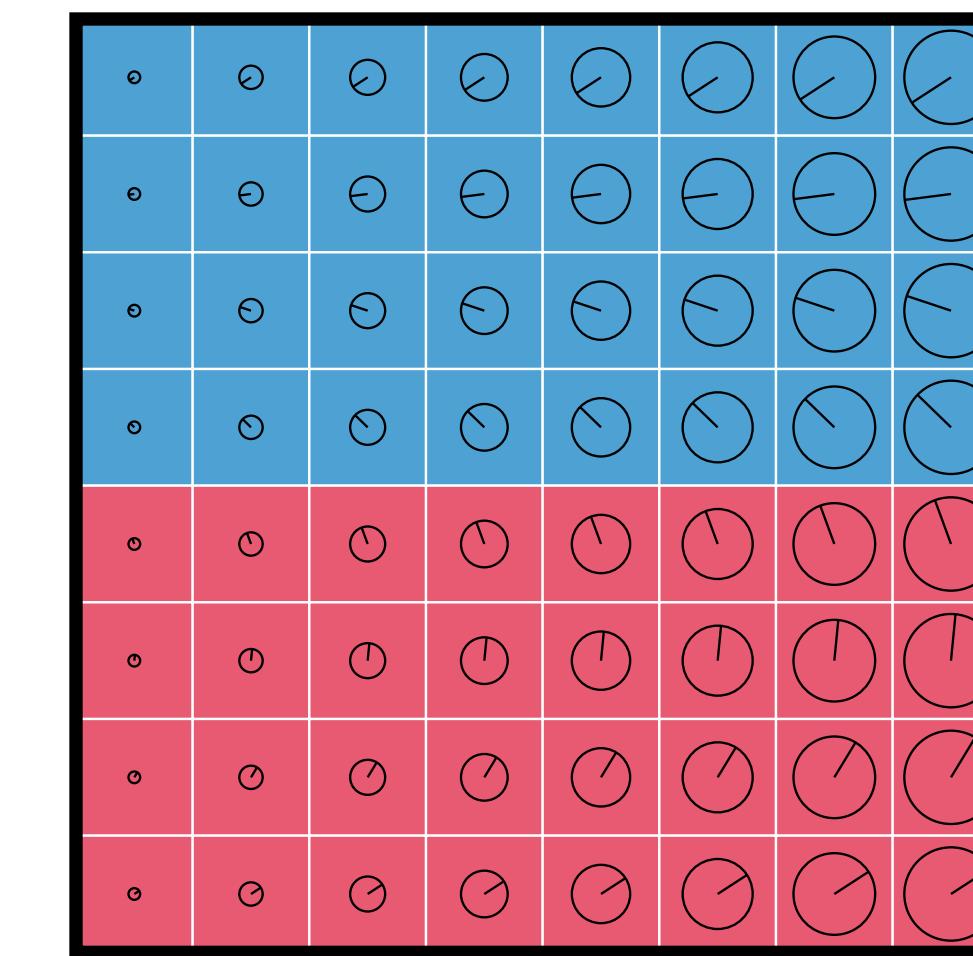
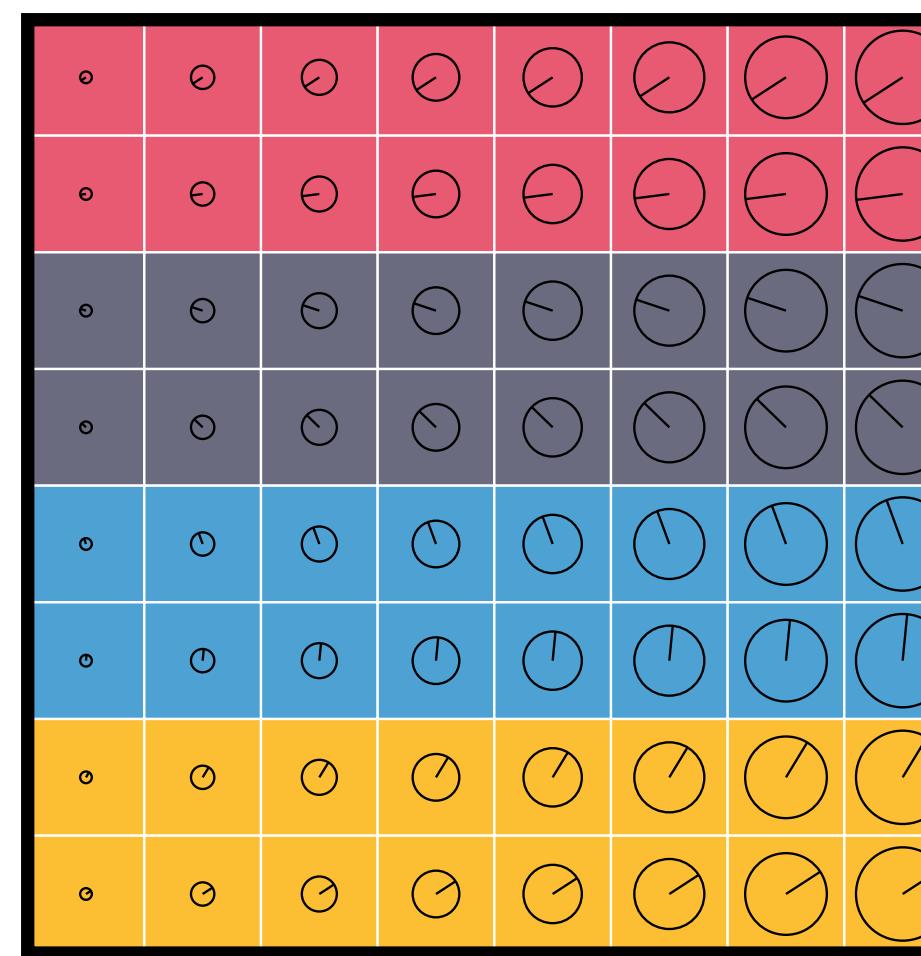
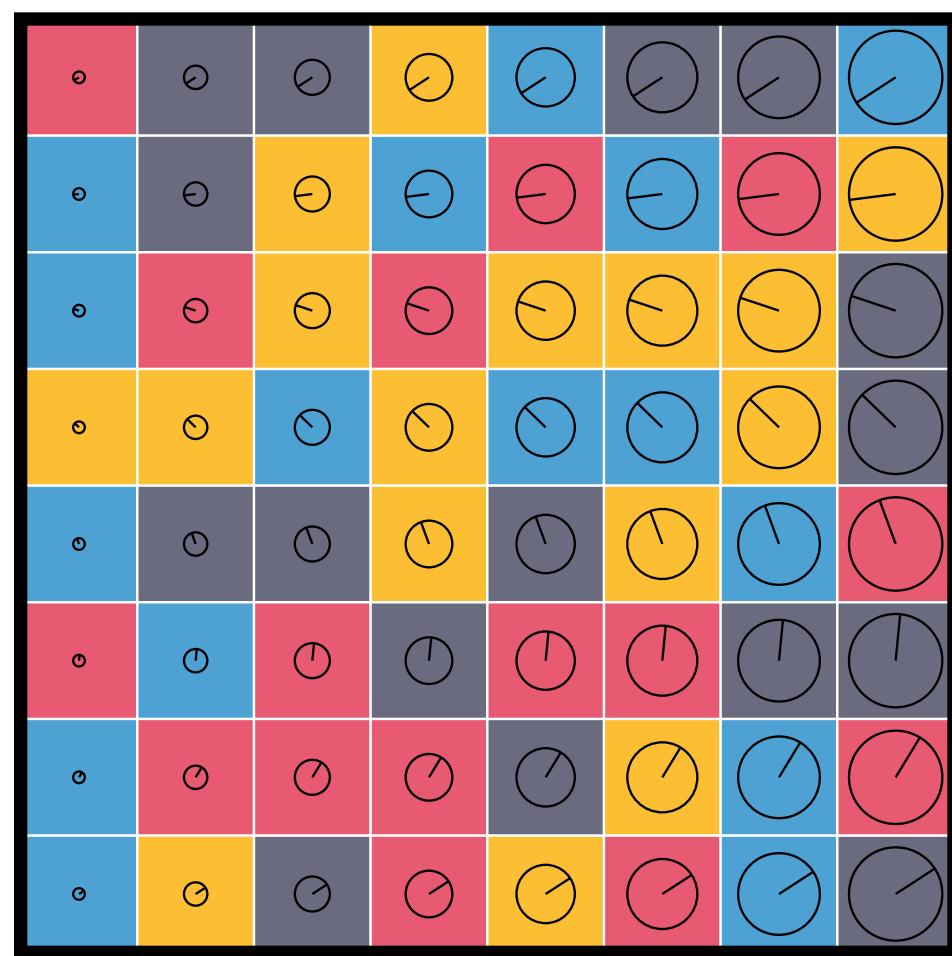


Decrease in expressivity

# Two ways of achieving simplicity

Increase in convexity

*increases informativeness*



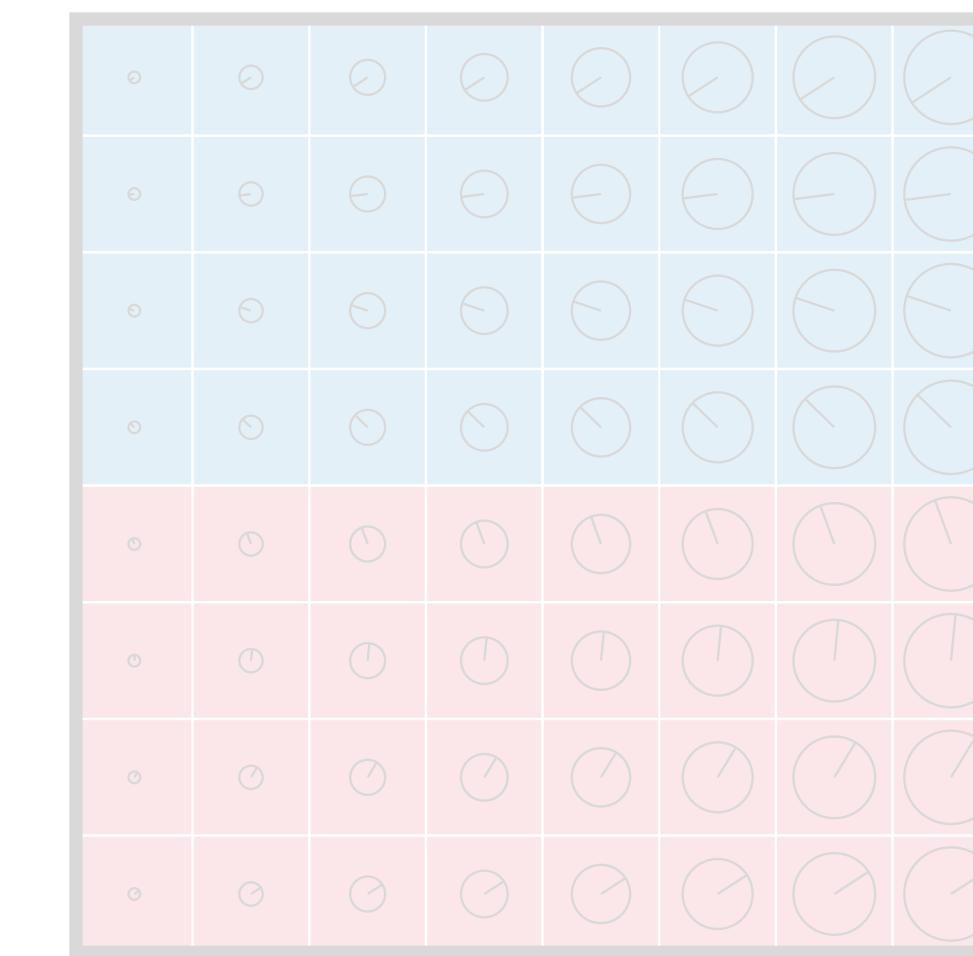
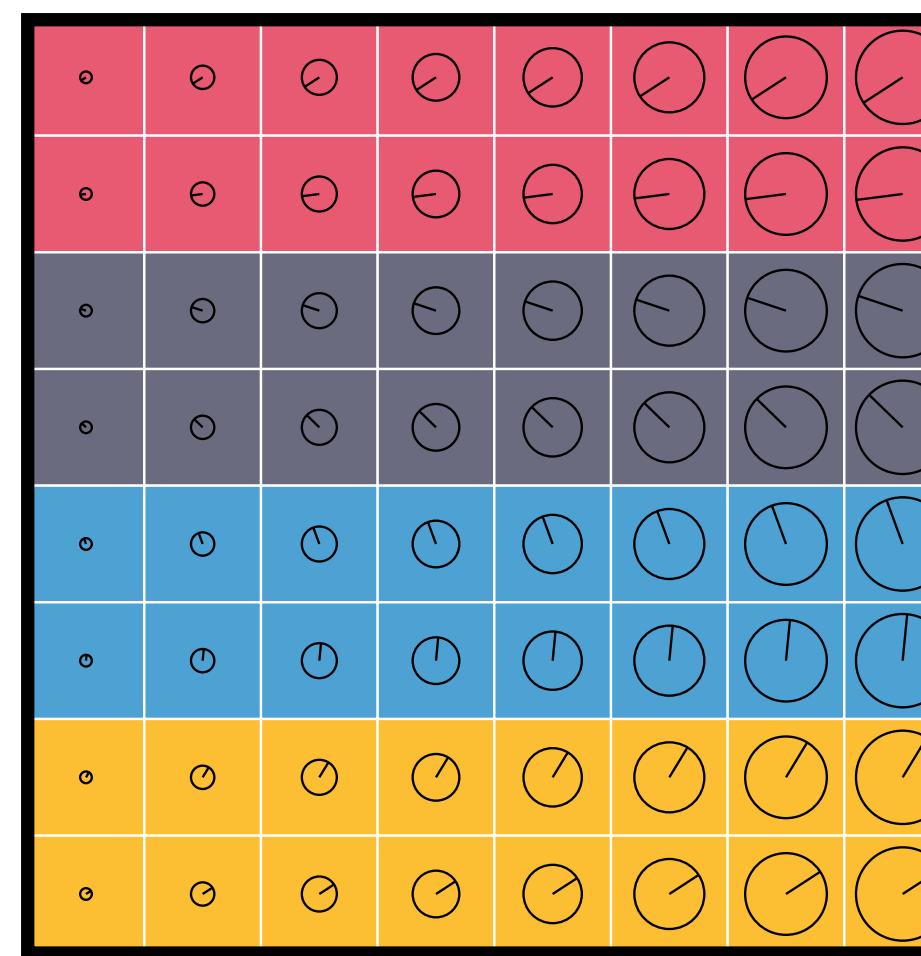
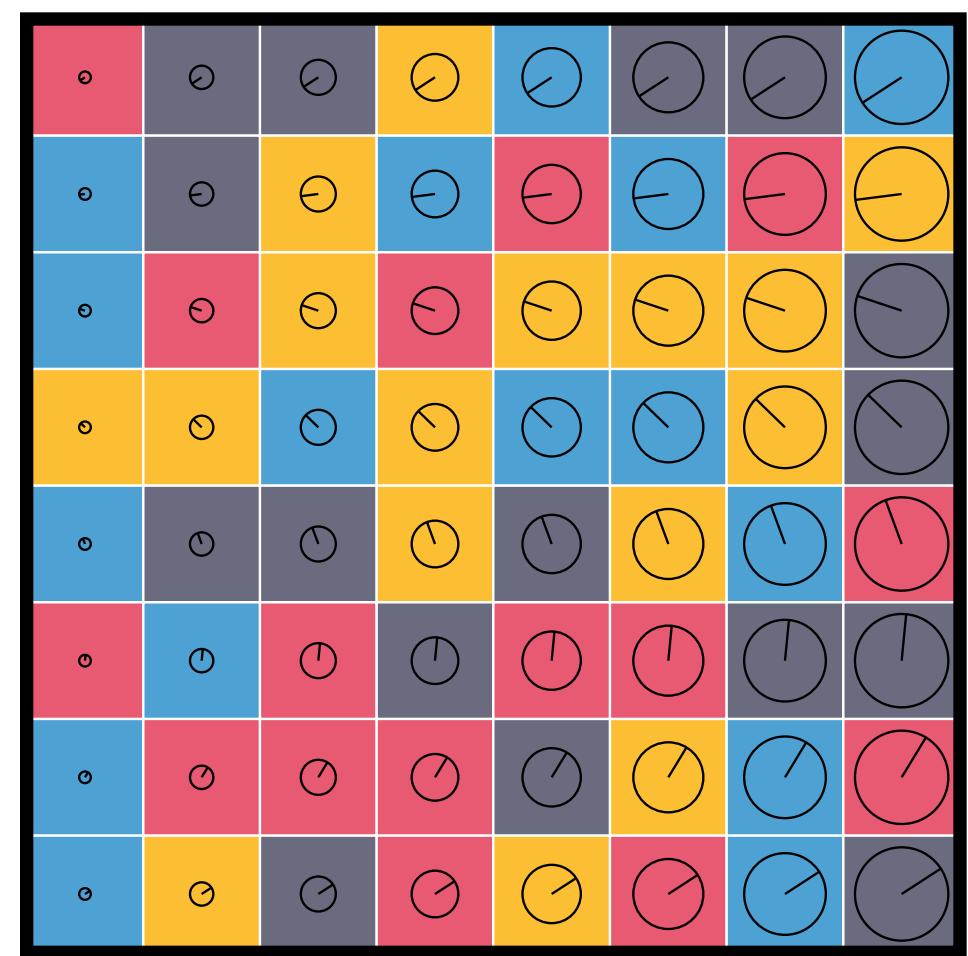
Decrease in expressivity

*decreases informativeness*

# Two ways of achieving simplicity

Increase in convexity

*increases informativeness*



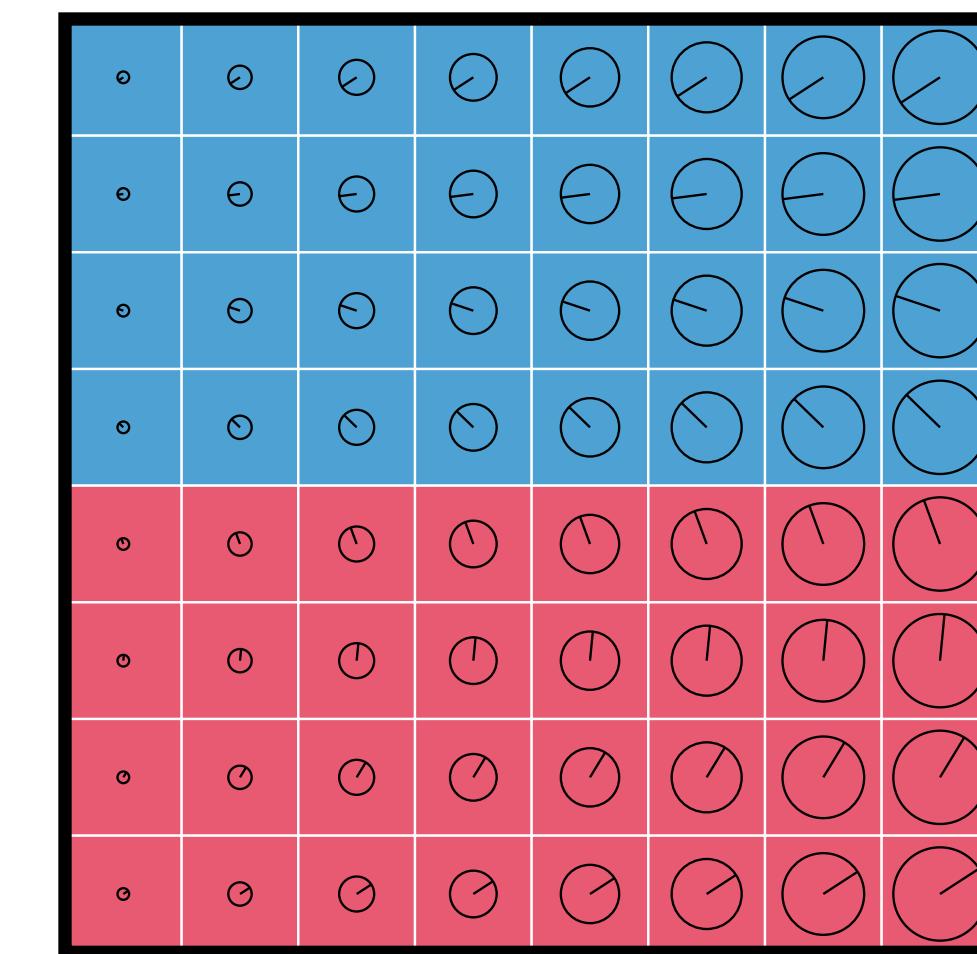
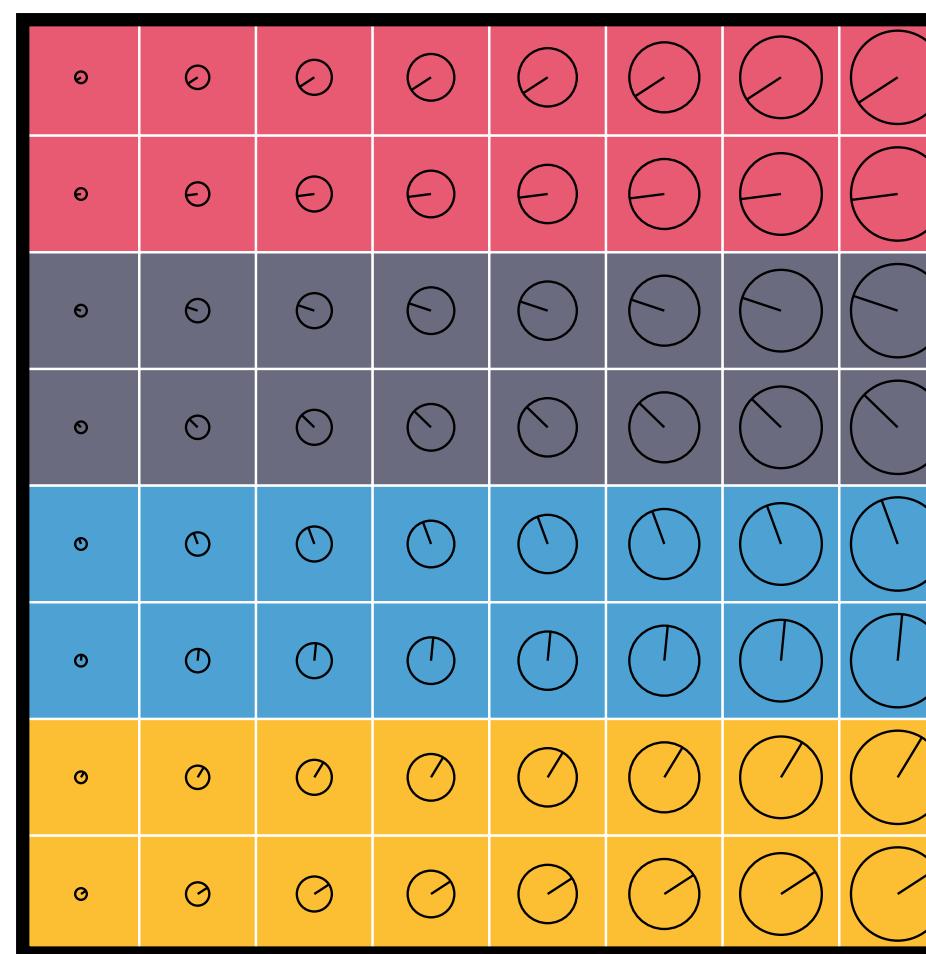
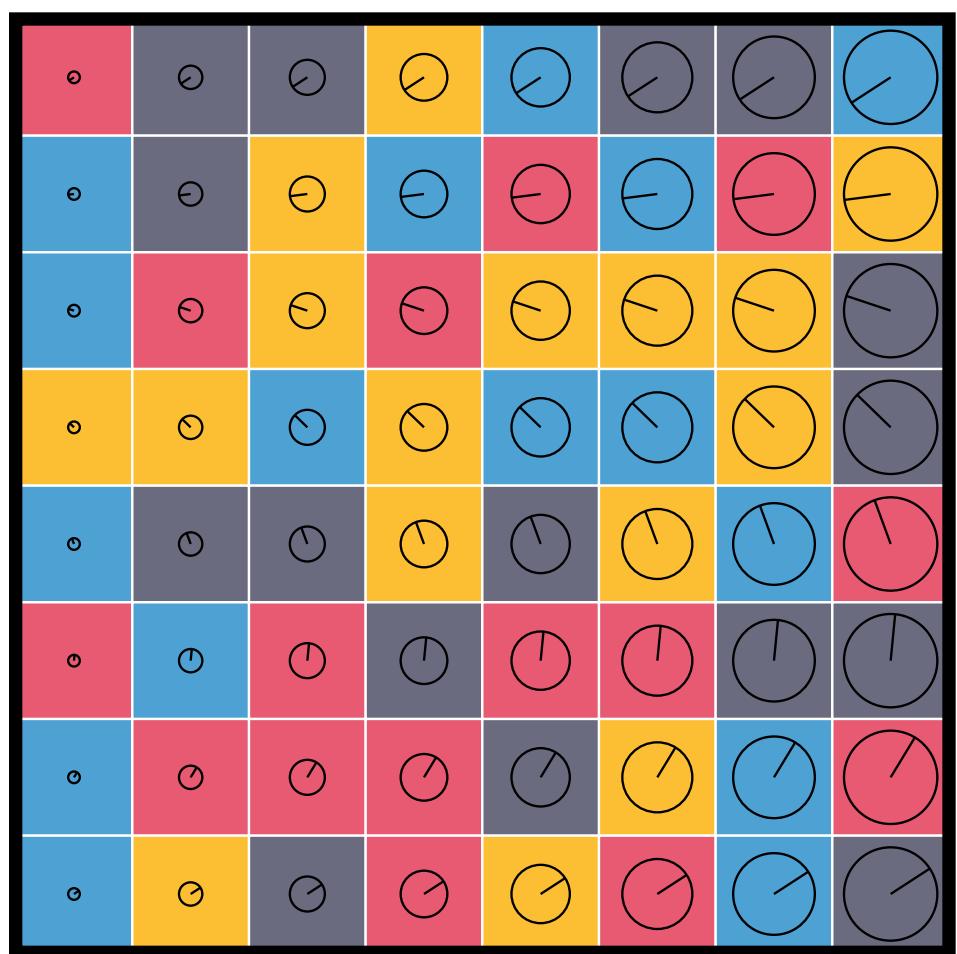
Decrease in expressivity

*decreases informativeness*

# Two ways of achieving simplicity

Increase in convexity

*increases informativeness*



Decrease in expressivity

*decreases informativeness*

## Conclusions

Languages are shaped in the simplicity–informativeness tradeoff by pressures from learning and communication

Learning contains a simplicity bias to prevent overfitting noise, and to aid reasoning about unseen meanings

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

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

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

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

*Thanks!*