

# Computational approaches to the explanation of universal properties of meaning

Fausto Carcassi & Jakub Szymanik



[https://thelogicalgrammar.github.io/ESSLLI22\\_langevo](https://thelogicalgrammar.github.io/ESSLLI22_langevo)



## Recap

1. Learnability can explain the presence of universals.
2. But is it the only (or the best) such explanation?
3. A natural idea: some notion of complexity explains both the universals and the learnability facts.

# What's the right measure of complexity?

- Previous attempts fail to capture the distinctions:
  - Automata theory (van Benthem 1986)
  - Computational complexity (Szymanik 2016)
  - Formal learning theory (Tiede 1999; Gierasimczuk 2007; Gierasimczuk 2009a)

**Let's try a philosophical idea!**

- The Language of Thought Hypothesis is the hypothesis that thinking happens in a mental language, the Language of Thought.
- Generally assumed to look somewhat like a logic: predicates get combined with logical operators.
- Modern version popularised by Jerry Fodor.
- The classical picture of LoTs that philosophers developed is meant to be an account of thinking, explaining phenomena such as learning from few examples, decision-making, perception, etc.

# The Language of Thought

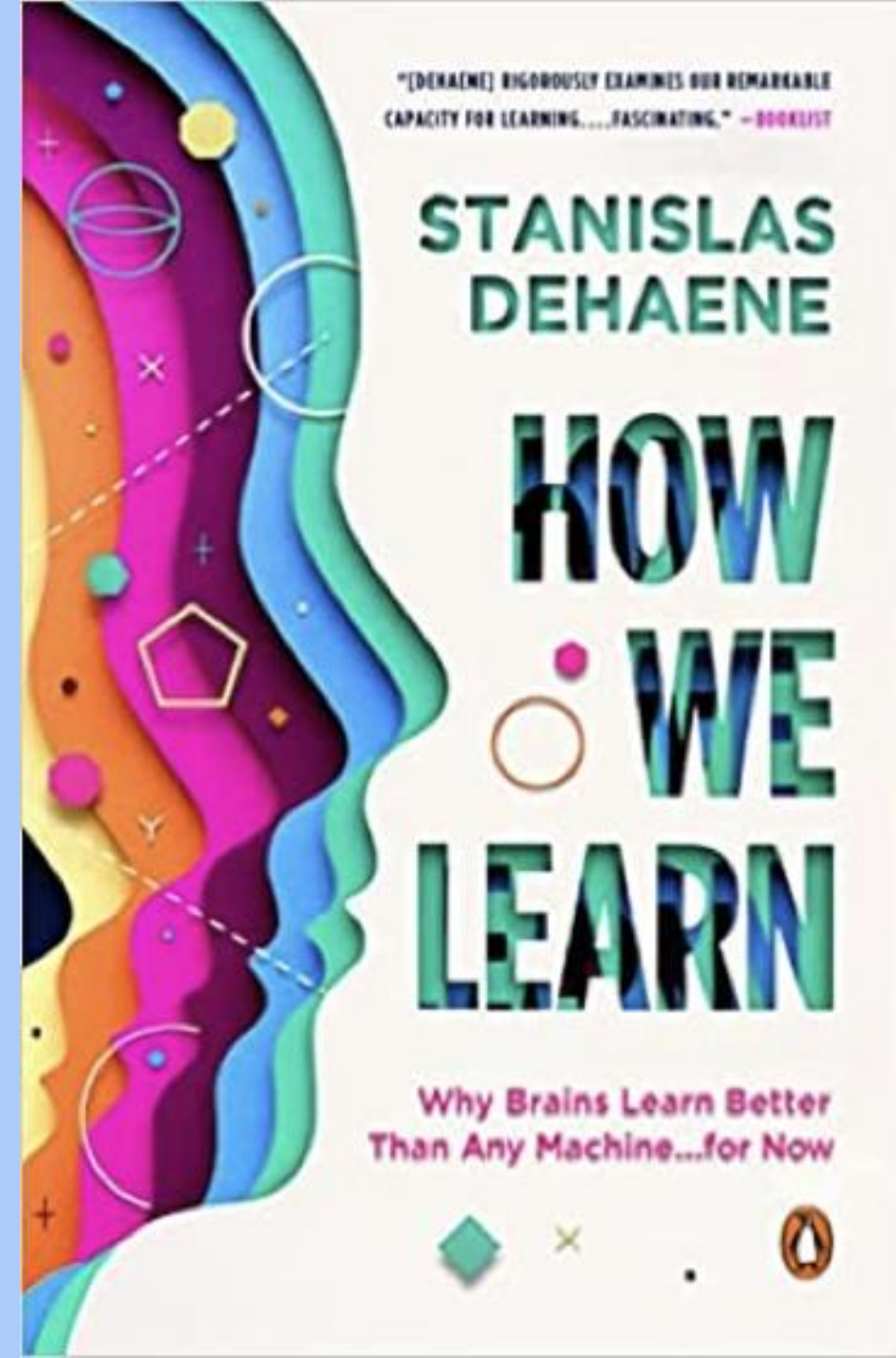
JERRY A. FODOR

*The Language and Thought Series*

Jerrold J. Katz  
D. Terence Langendoen  
George A. Miller

SERIES EDITORS

- Leibniz, 1677, Boole, 1854, Fodor, 1975,..., Rescorla, 2019, and others
- Recently, revived interest in cognitive science:
- Feldman, 2003, Tenenbaum and Griffiths, 2001, Tenenbaum and Xu, 2007, Piantadosi, 2016, etc.
- Perhaps a missing piece to make AI more human-like



# Boolean Categorization

- Boolean relations are a way to create new concepts:
  - `cousin' is a child of an uncle **or** aunt
  - `beer' is an alcoholic beverage usually made from malted cereal grain **and** flavored with hops, **and** brewed by slow fermentation
  - in basketball, `travel' is illegally moving the pivot foot **or** taking three or more steps without dribbling
  - `depression' is a mood disorder characterized by persistent sadness **and** anxiety, **or** feeling of hopelessness **and** pessimism, **or** ...

- How people acquire, represent, and use concepts?
- E.g., concepts depending on **and** are easier to learn than those depending on **or** (Bruner et al. '65).
- But the data seems more puzzling (see next slide).
- What's the logical theory of complexity here?



The *wudsy* objects in each set are surrounded by a square:

**Example 1**



**Example 2**



Given the above examples, which of the objects in this set are wudsy?



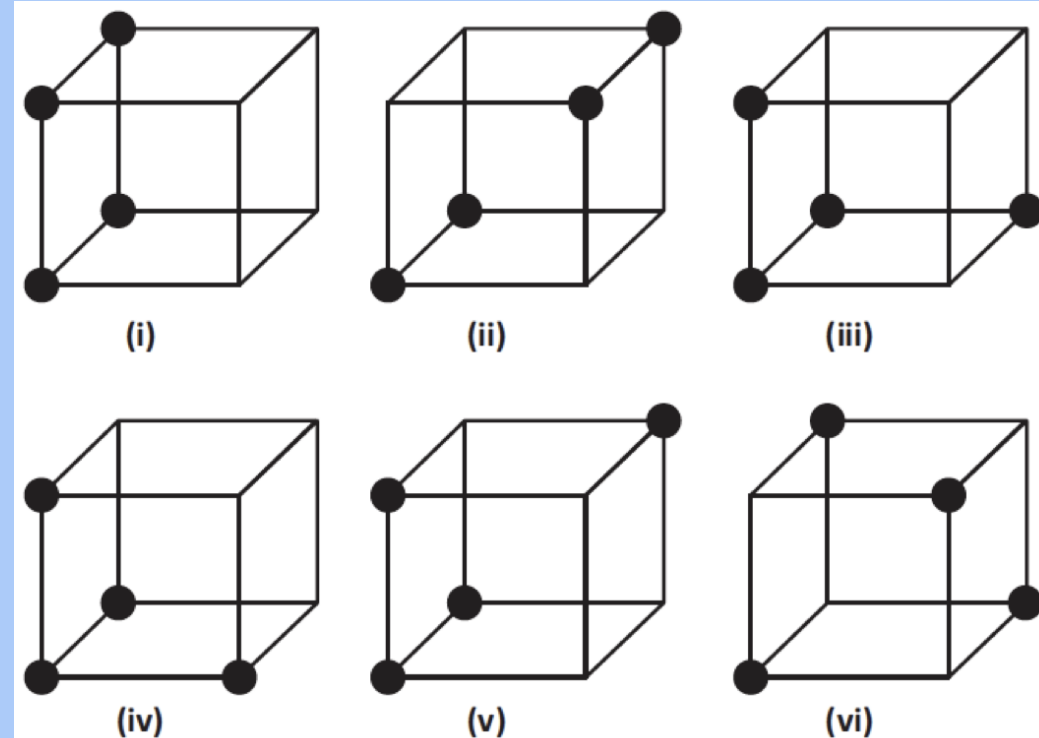
Please respond to each object



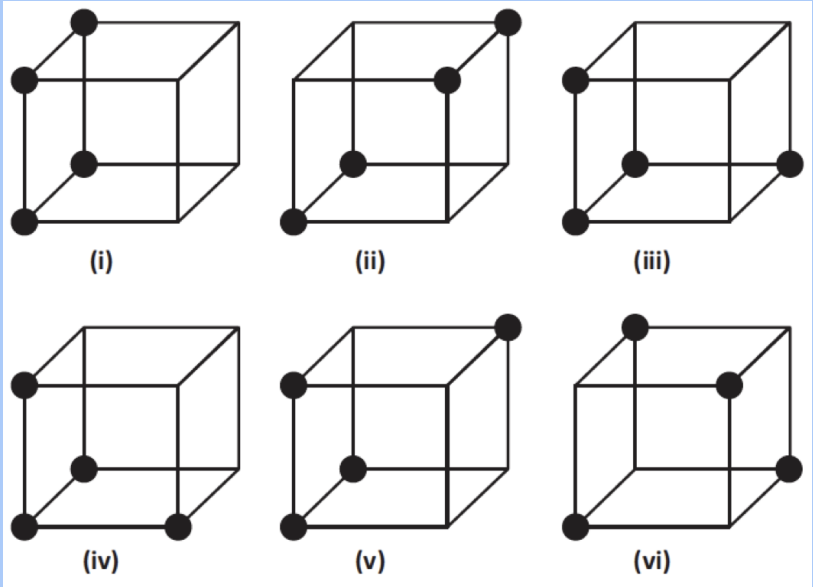
# Shepard's trend

- Six different sorts of concept based on three binary variables
- Each concept: 4 instances and 4 non-instances in 8 possibilities
- Different presentations methods: sequentially, simultaneously, etc.
- Dependent variables: errors, latencies, accuracy of descriptions, etc.
- $I < II < III, IV, V < VI$

Shepard et al.'61



# The instances of the concepts

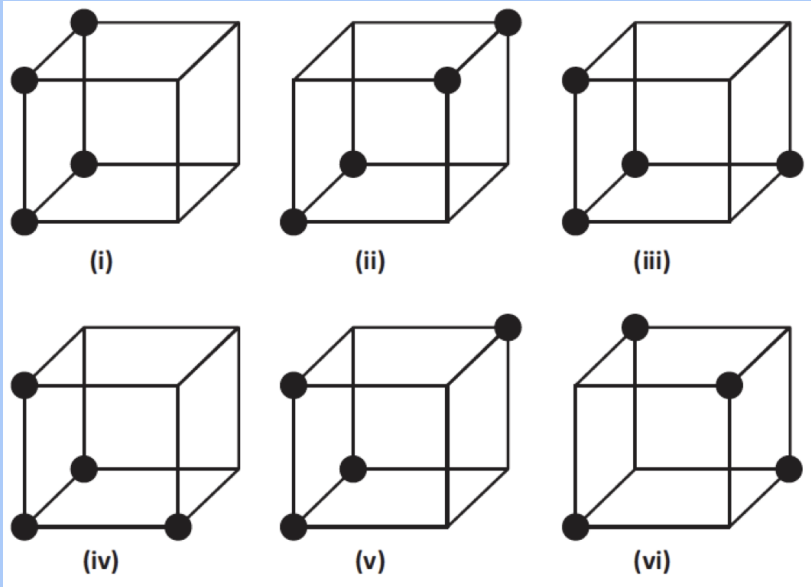


Concept number	Instances
I	not-a b c not-a b not-c not-a not-b c not-a not-b not-c
II	a b c a b not-c not-a not-b c not-a not-b not-c
III	a not-b c not-a b not-c not-a not-b c not-a not-b not-c
IV	a not-b not-c not-a b not-c not-a not-b c not-a not-b not-c
V	a b c not-a b not-c not-a not-b c not-a not-b not-c
VI	a b not-c a not-b c not-a b c not-a not-b not-c

# Boolean Complexity

- The length of the shortest Boolean formula logically equivalent to the concept, e.g., expressed in terms of the number of literals (positive or negative variables).
- Btw, finding the shortest formula is intractable.
- $(a \text{ and } b) \text{ or } (a \text{ and not } b) \text{ or } (\text{not } a \text{ and } b)$  reduces to  $(a \text{ or } b)$

# BC captures the trend



Concept number	Instances	Minimal description
I	not-a b c not-a b not-c not-a not-b c not-a not-b not-c	<i>not a</i> (1)
II	a b c a b not-c not-a not-b c not-a not-b not-c	<i>(a and b) or (not a and not b)</i> (4)
III	a not-b c not-a b not-c not-a not-b c not-a not-b not-c	<i>(not a and not c) or (not b and c)</i> (4)
IV	a not-b not-c not-a b not-c not-a not-b c not-a not-b not-c	<i>(not c or (not a and not b)) and (not a or not b)</i> (5)
V	a b c not-a b not-c not-a not-b c not-a not-b not-c	<i>(not a and not (b and c)) or (a and (b and c))</i> (6)
VI	a b not-c a not-b c not-a b c not-a not-b not-c	<i>(a and ((not b and c) or (b and not c))) or (not a and ((not b and not c) or (b and c)))</i> (10)

## New Data Set

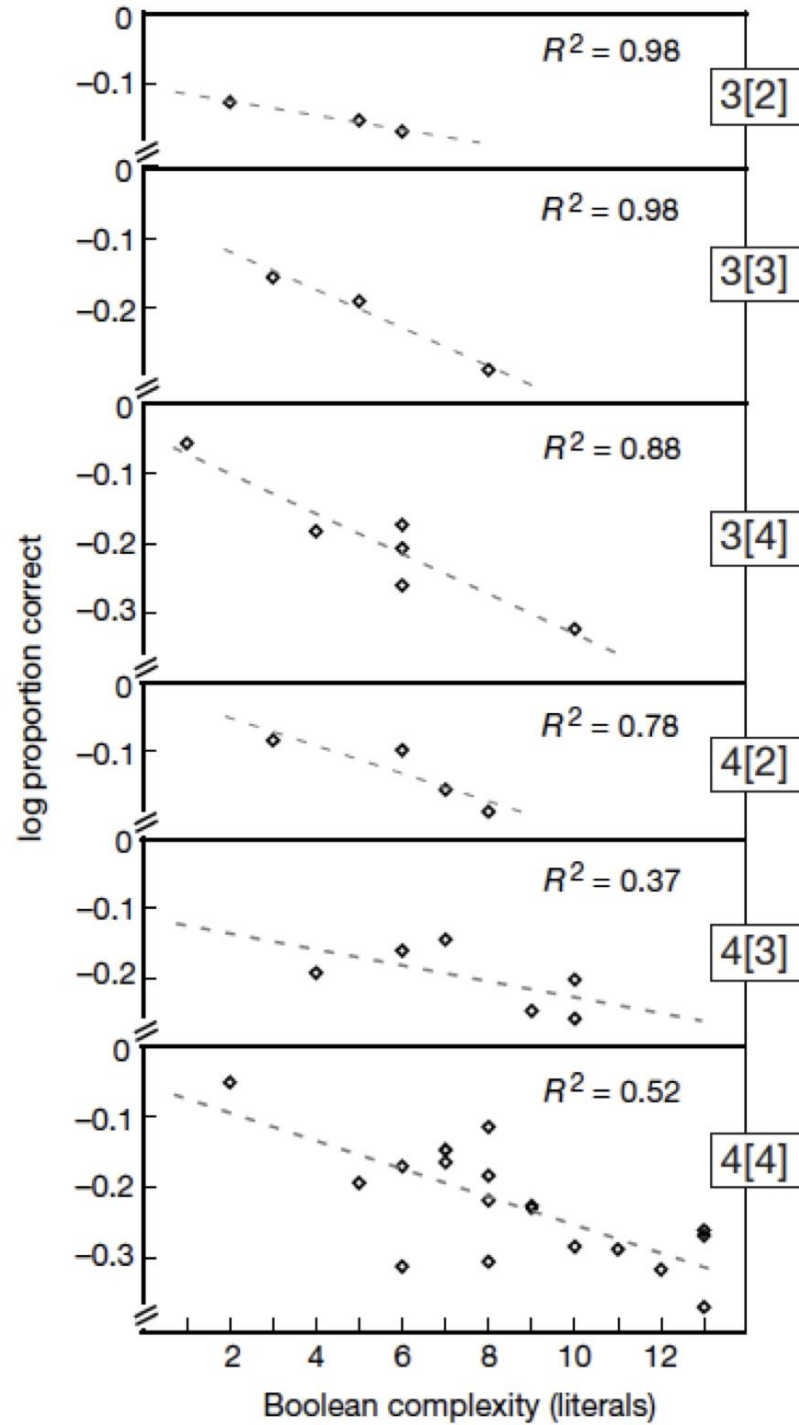
- Consider an arbitrary Boolean concept defined by  $P$  positive examples over  $D$  binary features.
- For Shepard types  $D=3$  and  $P=4$ .
- Feldman studies 76 Boolean concepts.

Feldman '01

		$P$				
$ D[P] $		1	2	3	4	5
$D$	1	1				
	2	1	2			
	3	1	3	3	6	
	4	1	4	6	19	27
	5	1	5			

Tested families

SHJ family





- **Simplicity: a unifying principle in cognitive science?**
- (Chater and Vitanyi, 2013), (Hsu, Chater, Vitanyi, 2013),...
- MDL one of the ways to operationalize simplicity.
- Other options: computational complexity (van Rooij, 2008, Szymanik, 2016), Kolmogorov complexity (see below), machine learning (cf. Carcassi & Szymanik, 2022)

Applying pure complexity idea  
to  
quantifier universals

## We need a more expressive logical grammar

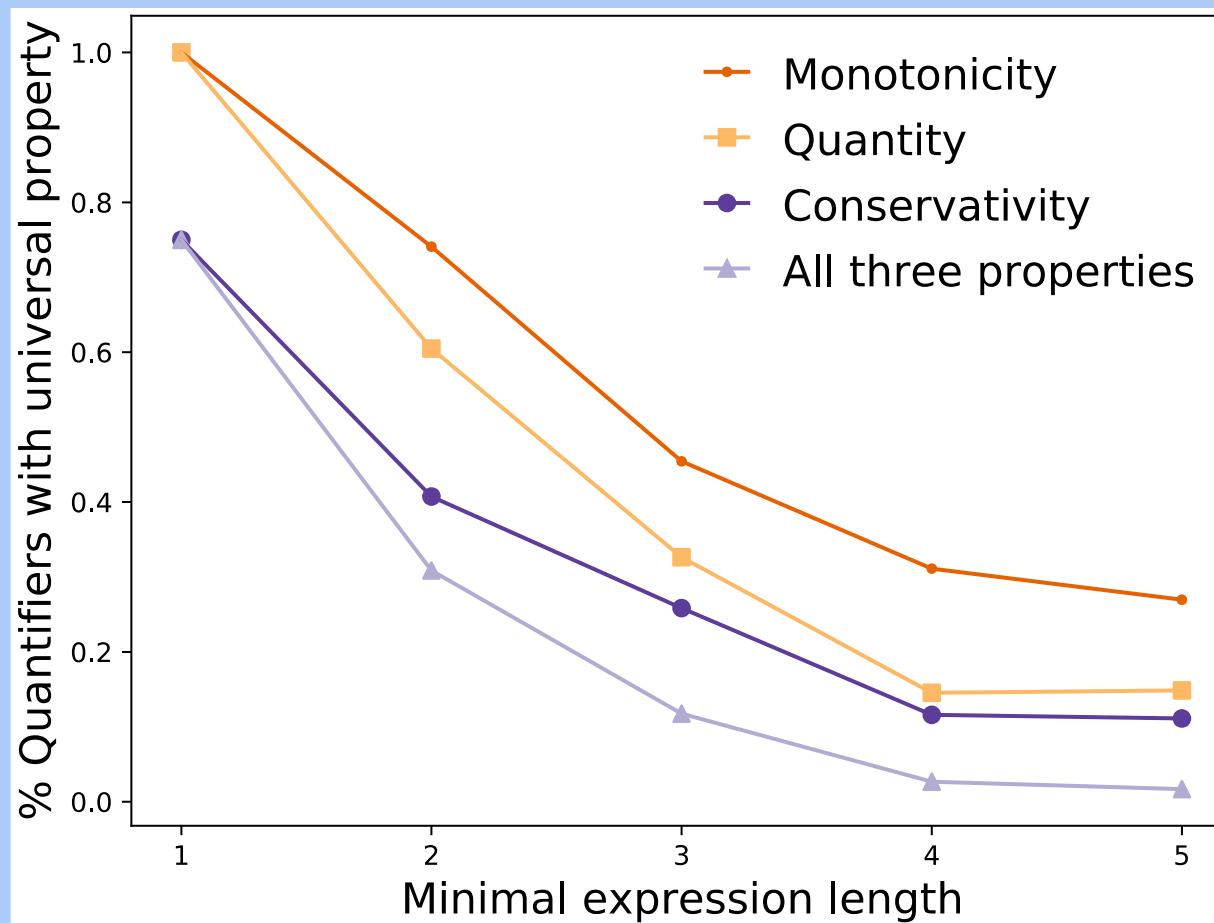
operator	type	gloss
$\cup$	$\text{SET} \times \text{SET} \rightarrow \text{SET}$	union
$\cap$	$\text{SET} \times \text{SET} \rightarrow \text{SET}$	intersection
$\setminus$	$\text{SET} \times \text{SET} \rightarrow \text{SET}$	setminus
$\iota(\cdot, \cdot)$	$\text{INT} \times \text{SET} \rightarrow \text{SINGLETON SET}$	'object at index'
$ \cdot $	$\text{SET} \rightarrow \text{INT}$	cardinality
$\subseteq$	$\text{SET} \times \text{SET} \rightarrow \text{BOOL}$	subset equal
$=$	$\text{INT} \times \text{INT} \rightarrow \text{BOOL}$	integer equality
$>$	$\text{INT} \times \text{INT} \rightarrow \text{BOOL}$	integer larger than
$\neg$	$\text{BOOL} \rightarrow \text{BOOL}$	negation
$\wedge$	$\text{BOOL} \times \text{BOOL} \rightarrow \text{BOOL}$	and
$\vee$	$\text{BOOL} \times \text{BOOL} \rightarrow \text{BOOL}$	or

- We generate all expression of length up to 5 (or 7): **solving the minimal pair problem**
- Length = the number of operators
- The minimal expression length of Q is the length of the shortest expression for this quantifier.
- At most  $1 = (2 > |A \cap B|)$  or  $\neg(|A \cap B|) > 1$
- We generated ~25k unique quantifiers (up to model size 8)

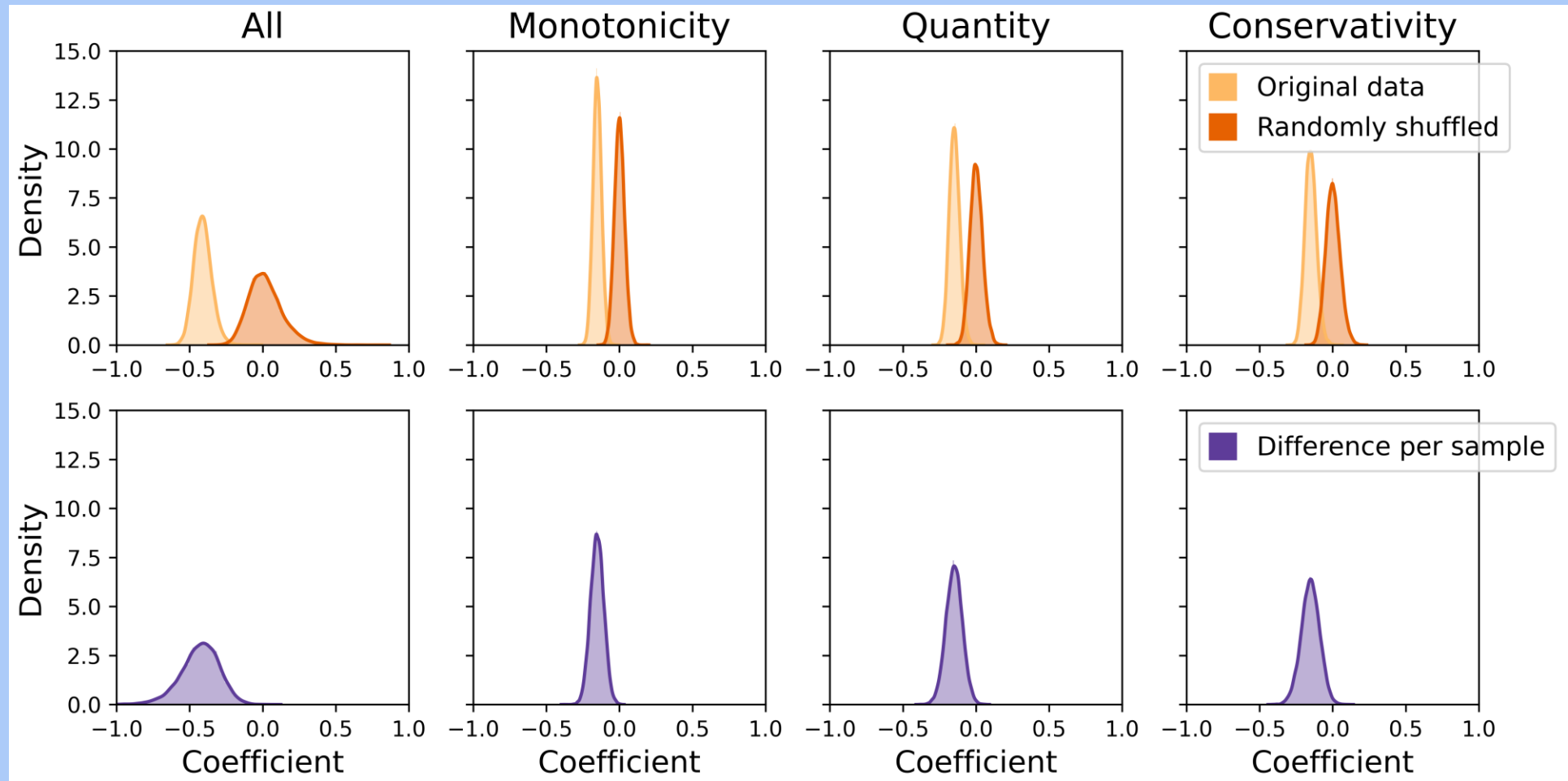
	YES	NO	%
monotonicity	-0.12	0.05	0.28
quantity	-0.15	0.03	0.15
conservativity	-0.16	0.02	0.12
all	-0.79	0.02	0.02

Table 2: Average (standardized) ML scores of quantifiers with (YES) versus without (NO) universal property and the proportion (%) of quantifiers with that universal property. *All* stands for quantifiers that have all three properties. The category “*all, NO*” stands for quantifiers that lack at least one property, i.e., quantifiers that do not have all three properties. For language  $\mathcal{L}_{+l}$ .

van de Pol, Lodder, van Maanen, Steinert-Threlkeld, Szymanik. Quantifiers satisfying semantic universals have shorter minimal description length. *Cognition* 2022  
 Data and code: <https://github.com/ivdpol/QuantifierComplexity>



van de Pol, Lodder, van Maanen, Steinert-Threlkeld, Szymanik. Quantifiers satisfying semantic universals have shorter minimal description length. *Cognition* 2022  
Data and code: <https://github.com/ivdpol/QuantifierComplexity>



van de Pol, Lodder, van Maanen, Steinert-Threlkeld, Szymanik. Quantifiers satisfying semantic universals have shorter minimal description length. *Cognition* 2022

Data and code: <https://github.com/ivdpol/QuantifierComplexity>

- **Meanings satisfying semantic universals are simpler**
- The setup avoids the minimal pair methodology
- Is it robust wrt the chosen LoT?
- So, which one is it: complexity or learnability?



Other measure of complexity

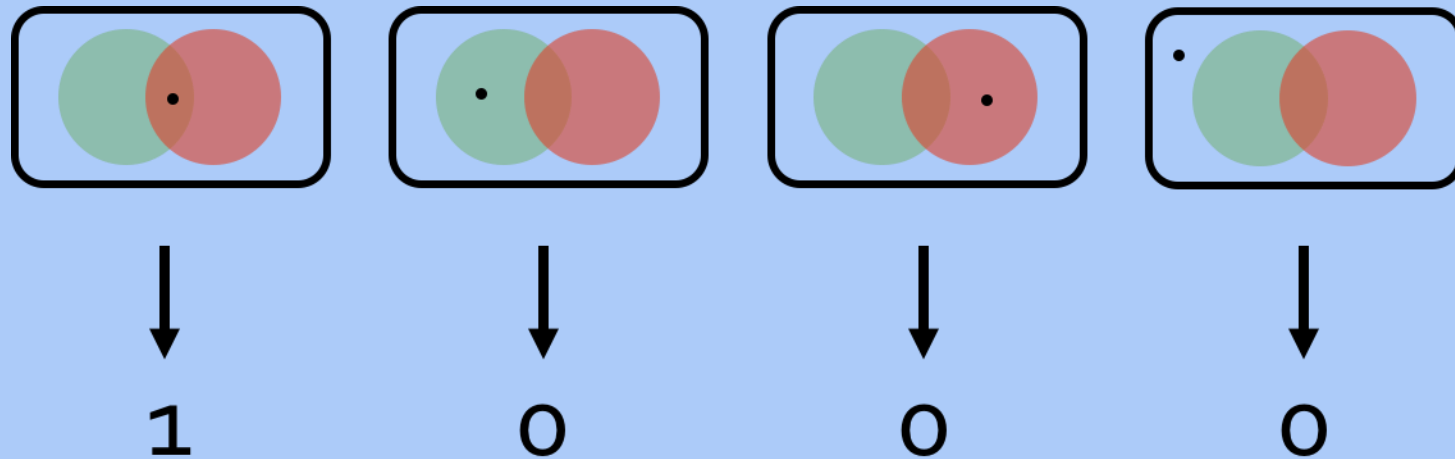
# Approximate Kolmogorov Complexity

- $K(x)$ —the length of the shortest program  $p$  that outputs  $x$
- Language dependent but the effect is bounded (Invariance Theorem)
- The drawback:  $K$  is uncomputable
- $LZ(x)$ —Lempel-Ziv is a tractable approximation of  $K$
- Recent applications: Dingle, Camargo, and Louis 2018; Feldman 2016; Valle-Pérez, Camargo, and Louis 2019

# Lempel-Ziv algorithm

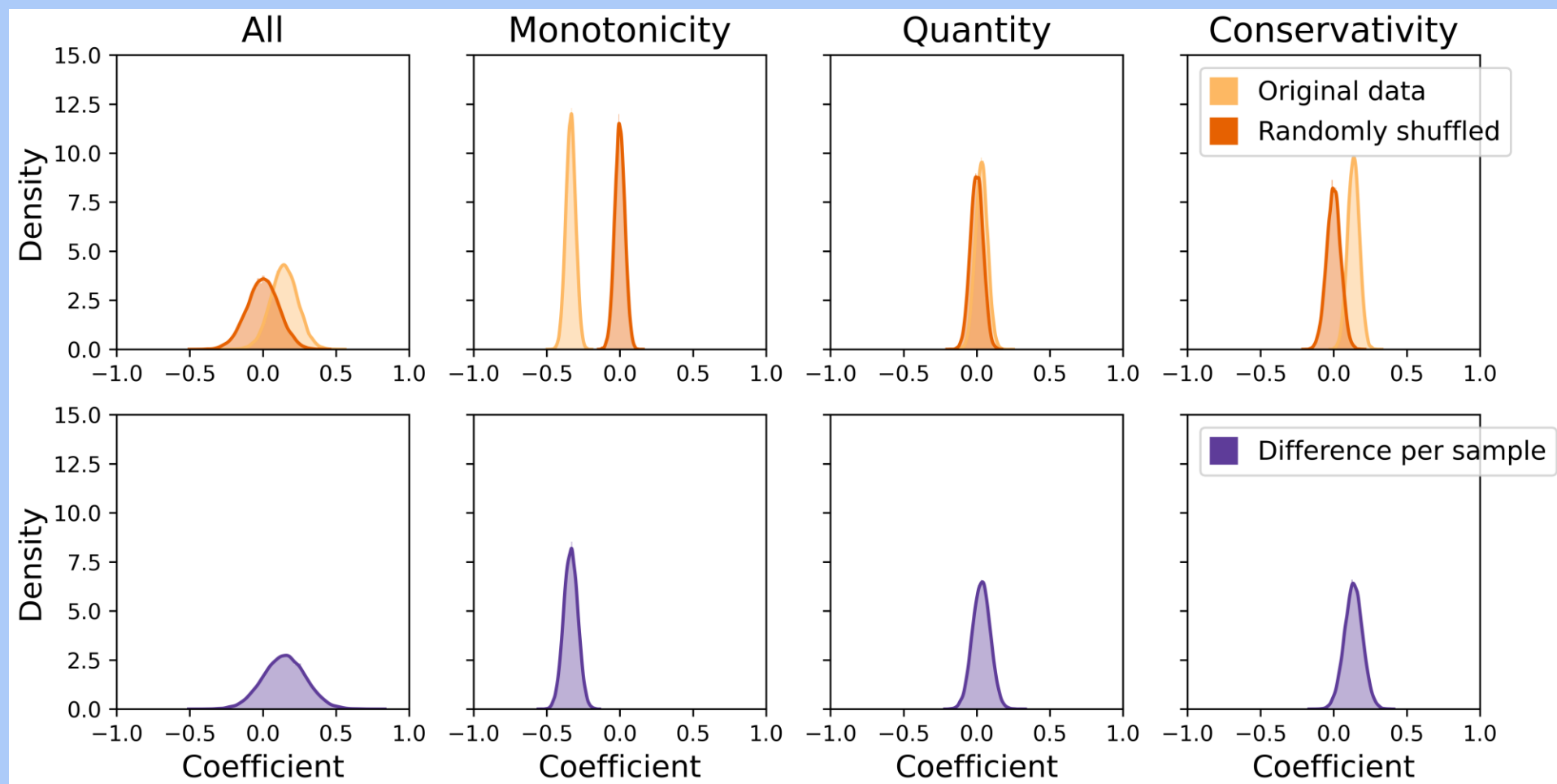
- Loseless compression algorithm
- Approximates  $K$  in the limit
- # unique subpatterns of a string
- A measure of structure
- $0|01|1|011|00|0110 = 6$
- $1|11|111|1111|111 = 4$

**Idea: universals induce regularity/structure in the distribution of truth values across models, which aid compressibility.**



	YES	NO	%
monotonicity	-0.22	0.08	0.28
quantity	0.02	0.00	0.15
conservativity	0.16	-0.02	0.12
all	0.16	0.00	0.02

Table 4: Average standardized LZ scores of quantifiers with (YES) versus without (NO) universal property and the proportion (%) of quantifiers with that universal property. *All* stands for quantifiers that have all three properties. The category “all, NO” stands for quantifiers that lack at least one property (i.e. which do not have all three). For language  $\mathcal{L}_{+l}$ .



# LoT+learning

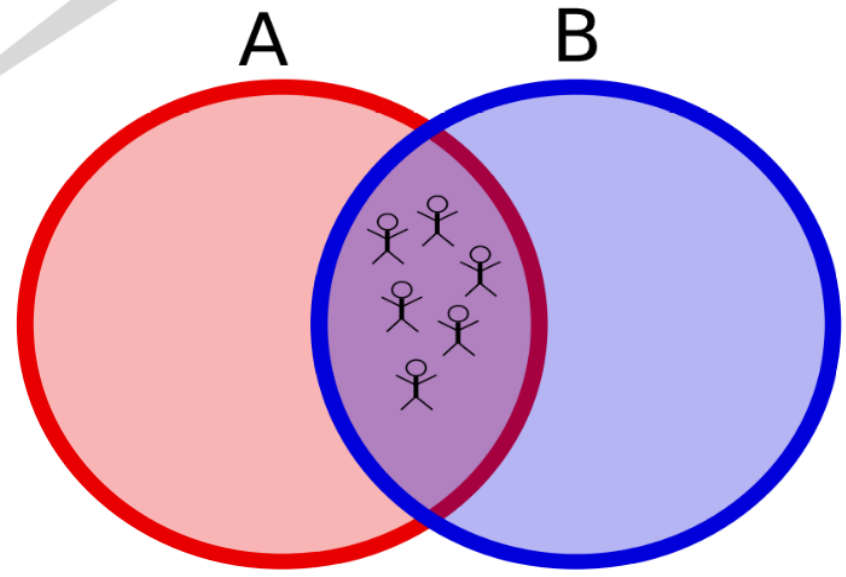
Why learnability as fencing-in fails

some  $\lambda A B . (nonempty? (intersection B A))$   
every  $\lambda A B . (subset? A B)$   
...

some  $\lambda A B . (subset? B A)$   
every  $\lambda A B . (empty? A)$   
...

"every"

adult



learner



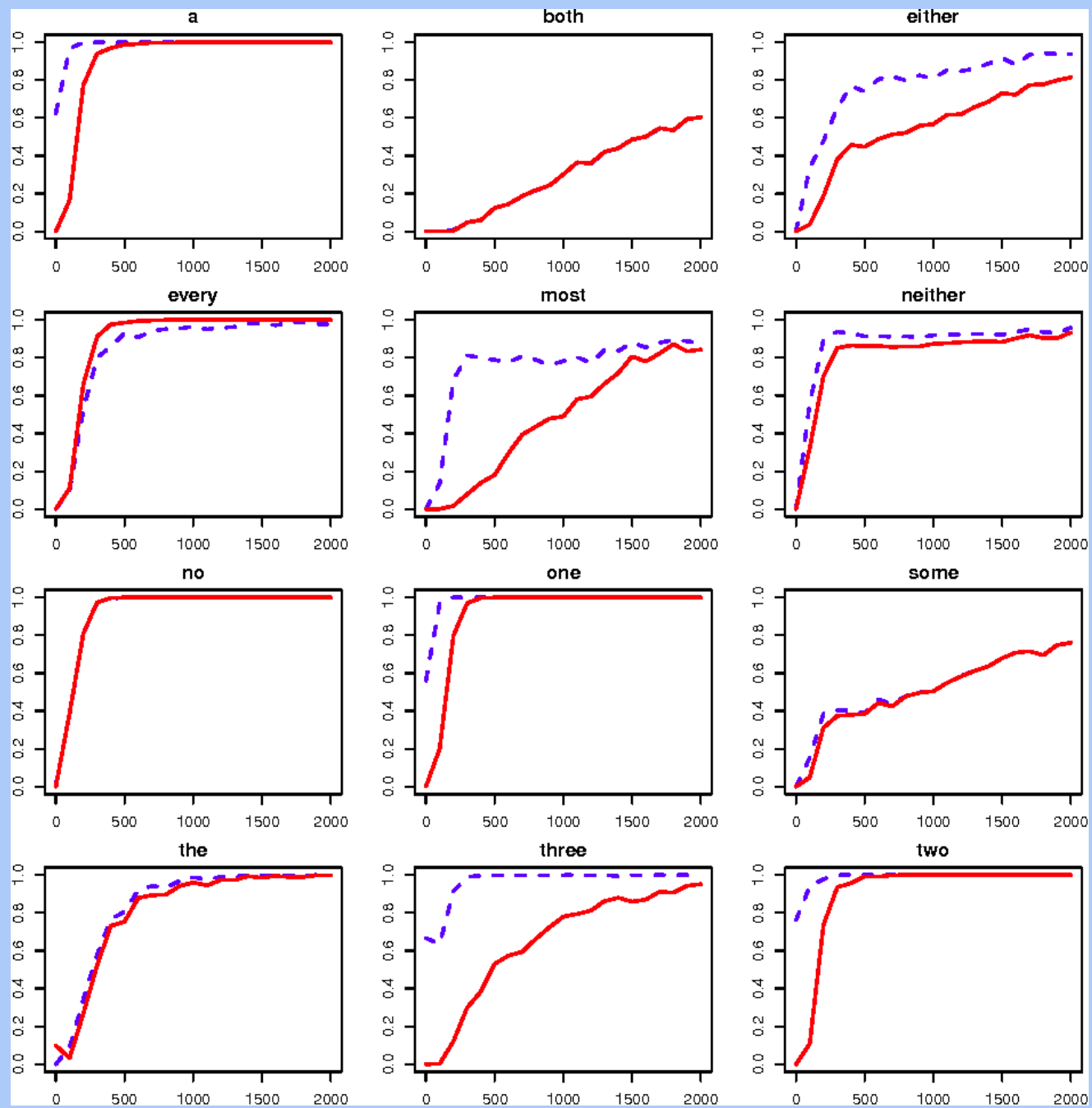
Nonterminal		Expansion	Gloss
START	→	$\lambda A B . \text{BOOL}$	Function of $A$ and $B$
BOOL	→	$\text{true}$	Always true
	→	$\text{false}$	Always false
	→	$(\text{card} > \text{SET SET})$	Compare cardinalities ( $>$ )
	→	$(\text{card} = \text{SET SET})$	Check if cardinalities are equal
	→	$(\text{subset? SET SET})$	Is a subset?
	→	$(\text{empty? SET})$	Is a set empty?
	→	$(\text{nonempty? SET})$	Is a set not empty?
	→	$(\text{exhaustive? SET})$	Is the set the entire set in the context?
	→	$(\text{singleton? SET})$	Contains 1 element?
	→	$(\text{doubleton? SET})$	Contains 2 elements?
	→	$(\text{tripleton? SET})$	Contains 3 elements?
	→	$(\text{union SET SET})$	Union of sets
	→	$(\text{intersection SET SET})$	Intersection of sets
	→	$(\text{set-difference SET SET})$	Difference of sets
SET	→	$A$	Argument $A$
	→	$B$	Argument $B$

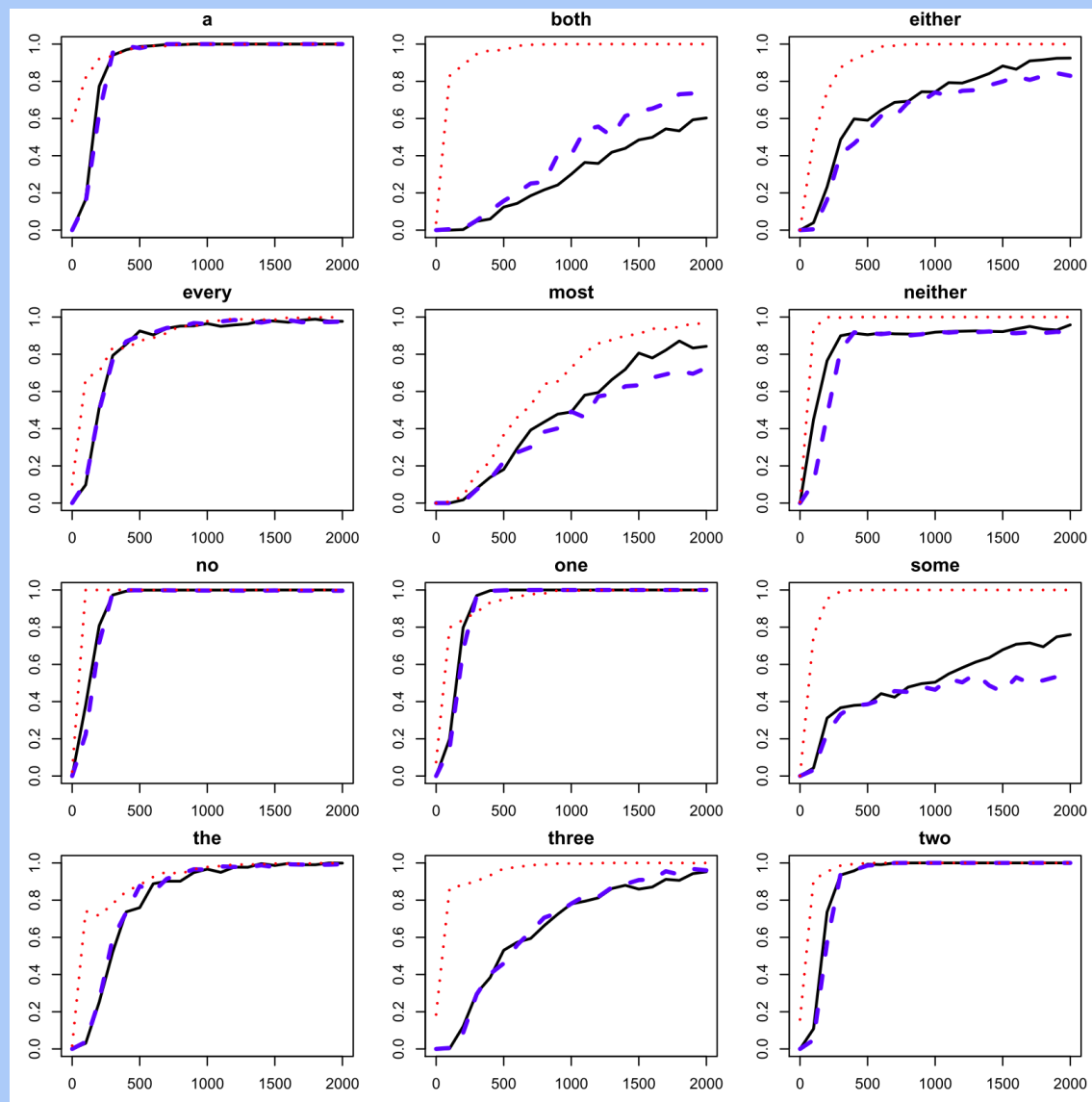
Word	Presupposition	Literal meaning
the	$\lambda A B . (\text{singleton? } A)$	$\lambda A B . (\text{nonempty? } (\text{intersection } A B))$
a/some	$\lambda A B . \text{TRUE}$	$\lambda A B . (\text{nonempty? } (\text{intersection } A B))$
one	$\lambda A B . (\text{nonempty? } A)$	$\lambda A B . (\text{singleton? } (\text{intersection } A B))$
two	$\lambda A B . (\text{nonempty? } A)$	$\lambda A B . (\text{doubleton? } (\text{intersection } A B))$
three	$\lambda A B . (\text{nonempty? } A)$	$\lambda A B . (\text{tripleton? } (\text{intersection } A B))$
both	$\lambda A B . (\text{doubleton? } A)$	$\lambda A B . (\text{doubleton? } (\text{intersection } A B))$
either	$\lambda A B . (\text{doubleton? } A)$	$\lambda A B . (\text{singleton? } (\text{intersection } A B))$
neither	$\lambda A B . (\text{doubleton? } A)$	$\lambda A B . (\text{empty? } (\text{intersection } A B))$
every	$\lambda A B . (\text{nonempty? } A)$	$\lambda A B . (\text{subset? } A B)$
most	$\lambda A B . (\text{nonempty? } A)$	$\lambda A B . (\text{card} > (\text{intersection } A B) (\text{set-difference } A B))$
none/no	$\lambda A B . (\text{nonempty? } A)$	$\lambda A B . (\text{empty? } (\text{intersection } A B))$

$$\begin{aligned} P(m|u_1, c_1, u_2, c_2, \dots, u_n, c_n) &\propto P(u_1, u_2, \dots, u_n|m, c_1, \dots, c_n) \cdot P(m) \\ &\propto \prod_{i=1}^n P(u_i|m, c_i) \cdot P(m) \end{aligned}$$

Prior  $P(m)$ : determined by the LoT grammar; shorter-to-express meanings are preferred.

Likelihood  $P(u_i|m, c_i)$ : preference for more informative, but with noise.



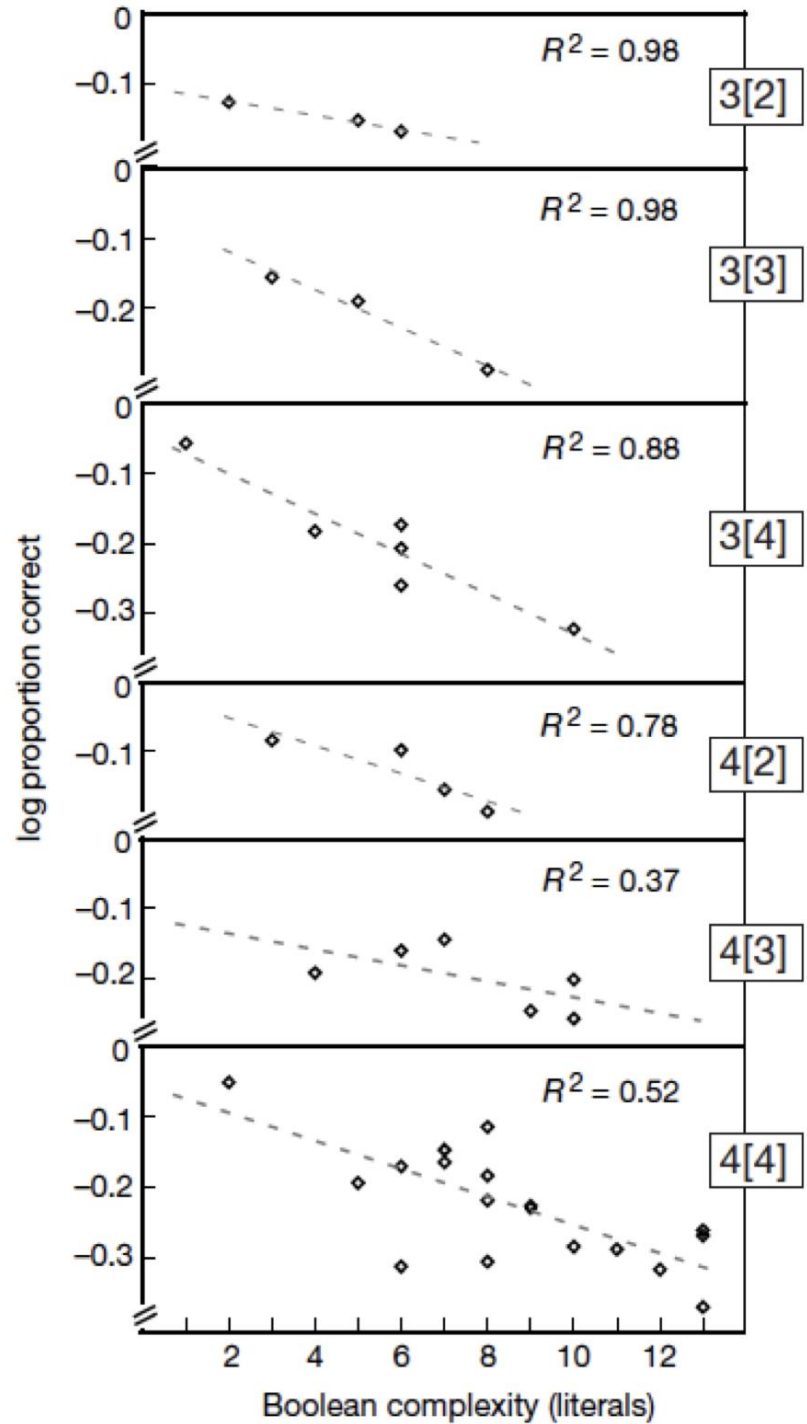


“Likely, the unrestricted space has many hypotheses which are so implausible, they can be ignored quickly and do not affect learning. The hard part of learning, may be choosing between the plausible competitor meanings, not in weeding out a large space of potential meanings.”

- Does this model predict human learning curves well?
- How sensitive are the model and its results sensitive to various choices (e.g. primitives, weights, shape of likelihood function)?
- Does it say anything general about e.g. monotone and topic-neutral quantifiers?

Robustness of the results  
depending on the LoT





**Which Boolean connectives?**

SIMPLE BOOLEAN			NAND		
START	→	<i>lambda x . BOOL</i>	START	→	<i>lambda x . BOOL</i>
BOOL	→	<i>(and BOOL BOOL)</i>	BOOL	→	<i>(nand BOOL BOOL)</i>
		<i>(or BOOL BOOL)</i>			<i>true</i>
		<i>(not BOOL)</i>			<i>false</i>
		<i>true</i>	BOOL	→	<i>(F OBJECT)</i>
		<i>false</i>	OBJECT	→	<i>x</i>
BOOL	→	<i>(F OBJECT)</i>	F	→	COLOR
OBJECT	→	<i>x</i>			SHAPE
F	→	COLOR			SIZE
		SHAPE	COLOR	→	<i>blue?</i>
		SIZE			<i>green?</i>
COLOR	→	<i>blue?</i>			<i>yellow?</i>
		<i>green?</i>	SHAPE	→	<i>circle?</i>
		<i>yellow?</i>			<i>rectangle?</i>
SHAPE	→	<i>circle?</i>			<i>triangle?</i>
		<i>rectangle?</i>	SIZE	→	<i>size1?</i>
		<i>triangle?</i>			<i>size2?</i>
SIZE	→	<i>size1?</i>			<i>size3?</i>
		<i>size2?</i>			
		<i>size3?</i>			

DNF			HORN CLAUSE		
START	→	<i>lambda x . DISJ</i>	START	→	<i>lambda x . HORN-CONJ</i>
DISJ	→	<i>CONJ</i> <i>(or CONJ DISJ)</i>	HORN-CONJ	→	<i>HORN-CLAUSE</i> <i>(and HORN-CLAUSE HORN-CONJ)</i>
CONJ	→	<i>BOOL</i> <i>(and BOOL CONJ)</i>	HORN-CLAUSE	→	<i>(implies HORN-CONJ PRIM)</i>
BOOL	→	<i>(F OBJECT)</i> <i>(not (F OBJECT))</i>	HORN-CLAUSE	→	<i>(implies HORN-CONJ false)</i>
OBJECT	→	<i>x</i>	PRIM	→	<i>(F OBJECT)</i>
F	→	<i>COLOR</i> <i>SHAPE</i> <i>SIZE</i>	OBJECT	→	<i>x</i>
COLOR	→	<i>blue?</i> <i>green?</i> <i>yellow?</i>	F	→	<i>COLOR</i> <i>SHAPE</i> <i>SIZE</i>
SHAPE	→	<i>circle?</i> <i>rectangle?</i> <i>triangle?</i>	COLOR	→	<i>blue?</i> <i>green?</i> <i>yellow?</i>
SIZE	→	<i>size1?</i> <i>size2?</i> <i>size3?</i>	SHAPE	→	<i>circle?</i> <i>rectangle?</i> <i>triangle?</i>
			SIZE	→	<i>size1?</i> <i>size2?</i> <i>size3?</i>

Which representational system is the most likely, given human responses?

Grammar	H.O.L.L	FP	$R^2_{response}$	$R^2_{mean}$
FULLBOOLEAN	−16296.84	27	.88	.60
BICONDITIONAL	−16305.13	26	.88	.64
CNF	−16332.39	26	.89	.69
DNF	−16343.87	26	.89	.66
SIMPLEBOOLEAN	−16426.91	25	.87	.70
IMPLIES	−16441.29	26	.87	.70
HORNCLAUSE	−16481.90	27	.87	.65
NAND	−16815.60	24	.84	.61
NOR	−16859.75	24	.85	.58
UNIFORM	−19121.65	4	.77	.06
EXEMPLAR	−23634.46	5	.55	.15
ONLYFEATURES	−31670.71	19	.54	.14
RESPONSE-BIASED	−37912.52	4	.03	.04

- You can also infer LoT from the informativeness-complexity trade-off but again it's hard to distinguish top candidates (Denić & Szymanik, 2022).
- All reasoning data (Zhai, Titov, Szymanik, 2015).
- In general, you shouldn't rather expect to find a unique candidate LoT (Carcassi & Szymanik, 2022).
- So, how shall we measure complexity? (relevant for tomorrow)

So, which one is the  
fundamental notion: complexity  
or learnability?