ABSTRACTING (IN) THE LANGUAGE OF THOUGHT

Jakub Szymanik



66

At the heart of cognitive science is an embarrassing truth: we do not know what mental representations are like.

-Steven Piantadosi, 2020

The Language of Thought

JERRYA. FODOR

The Language and Thought Series

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

SERIES EDITORS

LOT HYPOTHESIS

- The Language of Thought Hypothesis posits that abstraction occurs in a mental language, known as the Language of Thought.
- Generally assumed to look like logic: predicates get combined with logical operators.
- Modern version popularised by Jerry Fodor in the 70s.
- The classical picture of LoT that philosophers have developed is intended to be an account of thinking, explaining phenomena such as learning from a few examples, decision-making, and perception, among others.

WHY ARE WE TALKING ABOUT IT TODAY?

- ► 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, Sauerland et al., 2025, Dehaene et al. 2025, etc.
- ► Perhaps a missing peace to make AI more human-like.
- Current LLMs excel at pattern recognition and statistical inference, but they do not necessarily possess the same symbolic reasoning capabilities as humans.
- True intelligence might require a deeper understanding of the world, potentially through a language-like system for representing and manipulating concepts.



MODERN *LOT*

- ► We focus on some small but critical conceptual domains.
- We have an agent with a fairly natural LoT, consisting of primitive concepts and a small set of operators and composition rules.
- ► We assume people have a simplicity prior: concepts that are harder to express with the LoT have lower prior probability.
- ► For example,...

EXPLAINING UNIVERSAL CROSS-LINGUISTIC PROPERTIES



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

Piantadosi, Tenenbaum, Goodman. Modeling the acquisition of quantifier semantics : a case study in function word learnability, 2012

LEARNABILITY

- Prior: specifying the learner's estimate of how likely any hypothesis is before any labeled objects have been observed.
- The prior is constructed using the LoT.
- Likelihood: quantifying the probability that the particular set was particularly labeled, if h were the true concept.
- Inferential statistical model: P(h
 | observed sets and labels).

 $P(m \mid u_1, ..., u_n, c_1, ..., c_n) \propto P(u_1, ..., u_n \mid m, c_1, ..., c_n) P(m).$ $P(u_1, ..., u_n \mid m, c_i, ..., c_n) = \prod_{i=1}^n P(u_i \mid m, c_i).$

VIA SIMPLICITY

| operator | type | gloss |
|-----------------------|--|---------------------|
| \cup | $SET\timesSET\toSET$ | union |
| \cap | $SET\timesSET\toSET$ | intersection |
| \setminus | $\mathrm{SET}\times\mathrm{SET}\to\mathrm{SET}$ | setminus |
| $\iota(\cdot, \cdot)$ | $\text{INT}\times\text{SET}\rightarrow\text{SINGLETON}$ Set | 'object at index' |
| $ \cdot $ | $\mathrm{SET} ightarrow \mathrm{INT}$ | cardinality |
| \subseteq | $SET\timesSET\toBOOL$ | subset equal |
| = | $\text{INT} \times \text{INT} \rightarrow \text{BOOL}$ | integer equality |
| > | $\text{INT} \times \text{INT} \rightarrow \text{BOOL}$ | integer larger than |
| | $\operatorname{BOOL} \to \operatorname{BOOL}$ | negation |
| \wedge | $\operatorname{BOOL} \times \operatorname{BOOL} \to \operatorname{BOOL}$ | and |
| \vee | $\operatorname{BOOL} \times \operatorname{BOOL} \to \operatorname{BOOL}$ | or |

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

- Why do languages lexicalize only some possible concepts?
- ► Pick a LoT.
- Generate artificial concepts within the LoT.
- Lexicalized concepts have shorter MDL in the LoT than nonlexicalized, yet logically possible, ones.
- LLMs exhibit similar bias for simplicity (Wang et al., 2024)



| Indefinite pronoun | Flavors | Feature formula | c(i) |
|--------------------|--|-----------------|------|
| someone | specific known, specific unknown, non-specific | SE^- | 1 |
| anyone | negative polarity, free choice | $SE^+ \cap N^-$ | 2 |
| no one | negative indefinite | N^+ | 1 |

| Table 3 | : English | indefinite | pronouns, | the flavors | they | convey, | their | feature | formulae | and |
|----------|-----------|------------|-----------|-------------|------|---------|-------|---------|----------|-----|
| their co | mplexitie | es. | | | | | | | | |

| Indefinite pronoun | Flavors | Feature formula | c(i) |
|--------------------|---------------------------------|---|------|
| kto-to | specific unknown, non-specific | $K^- \cap SE^-$ | 2 |
| kto-nibud' | non-specific | $S^- \cap SE^-$ | 2 |
| kto-libo | non-specific, negative polarity | $(S^- \cap SE^-) \cup ((SE^+ \cap R^+) \cap N^-)$ | 5 |
| nikto | negative indefinite | N^+ | 1 |
| koe-kto | specific known | K^+ | 1 |
| kto by to ni bylo | negative polarity | $(SE^+ \cap R^+) \cap N^-$ | 3 |
| kto ugodno | free choice | $SE^+ \cap R^-$ | 2 |

Table 4: Russian indefinite pronouns, the flavors they convey, their feature formulae and their complexities.

Denić, Steinert-Threlkeld, Szymanik. Indefinite pronouns optimize the simplicity/informativeness trade-off. Cognitive Science, 2022

SIMPLICITY VS INFORMATIVENESS

- The complexity of the language system is the minimal number of rules needed to define it in a LoT.
- Communicative cost is the reconstruction error.
- Evolution balances complexity and the communicative costs.

BUT IS THERE LOT UNIFYING ALL THE EXAMPLES?

66

"The choice of innate primitives can be viewed as a strictly empirical question that should be determined through independent experiments.".

-Piantadosi & Jacobs, 2016

CAN WE INFER IT?

APPROACH 1: INFERRING Lot from learning Data

LEARNING A NEW CONCEPT, "GLEEB"







FELDMAN'S RESULTS

- Consider an arbitrary Boolean concept defined by P positive examples over D binary features, P[D].
- Boolean complexity accounts for 50% of variance in the dataset.

QUESTIONS

Which Boolean connectives?

WHAT'S THE RIGHT FELDMAN'S GRAMMAR?

| S | IMPL | LEBOOLEAN | | | NAND |
|---------------|-----------------------------|--|---------------|-----------------------------|--|
| START BOOL | \rightarrow \rightarrow | lambda x . BOOL (and BOOL BOOL) (or BOOL BOOL) (not BOOL) | START BOOL | \rightarrow \rightarrow | lambda x . BOOL (nand BOOL BOOL) true false |
| | | true | BOOL | \rightarrow | (F OBJECT) |
| | | false | OBJECT | \rightarrow | X |
| BOOL | \rightarrow | (F OBJECT) | F | \rightarrow | COLOR |
| OBJECT | \rightarrow | X | | | SHAPE |
| F | \rightarrow | COLOR | | | SIZE |
| | | SHAPE | COLOR | \rightarrow | blue? |
| | | SIZE | | | green? |
| COLOR | \rightarrow | blue? | | | yellow? |
| | | green? | SHAPE | \rightarrow | circle? |
| | | yellow? | | | rectangle? |
| SHAPE | \rightarrow | circle? | | | triangle? |
| | | rectangle? | SIZE | \rightarrow | size1? |
| | | triangle? | | | size2? |
| SIZE | \rightarrow | size1? | | | size3? |
| | | size2? | | | |
| | | size3? | | | |

WHAT'S THE RIGHT FELDMAN'S GRAMMAR?

.

.

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

GRAMMAR COMPARISON

Bayesian data analysis model: which representational system is the most likely, given human responses?

| Grammar | H.O.LL | 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 |

APPROACH 2: INFERRING COT-PRIMITIVES FROM THE COMPLEXITY-INFORMATIVENESS TRADE-OFF

FROM TRADE-OFF





Denić, Szymanik. Reverse-engineering the language of thought: A new approach. CogSci, 2022 See also: Denić, Szymanik. Recursive Numeral Systems Optimize the Trade-off Between Lexicon Size and Average Morphosyntactic Complexity. Cognitive Science, 2024. What LoT primitives underlie number concepts 1-99?

| Denoted number (numeral) | Morphosyntactic make-up |
|--------------------------|--|
| 6 (iwan) | 10 (-wan) $-$ (\emptyset) 4 (i-) |
| 30 (wanatubatna) | 2 (-tu-) \cdot (\emptyset) 20 (-hotne) |
| So (wanetunotne) | — (- <i>e</i> -) 10 (<i>wan</i> -) |
| A2 (tuikashimatuhatna) | 2 (-tu-) · (∅) 20 (-hotne) |
| | + (-ikashima-) 2 (tu-) |

Table 2: Top three LoT hypotheses

| PRIM |
|-------------------------|
| $\{1, 2, 3, 5, 10\}$ |
| $\{1, 2, 3, 4, 5, 10\}$ |
| $\{1, 2, 5, 10\}$ |

APPROACH 3: INFERRING COT-RULES FROM REASONING DATA

Data - syllogistic Reasoning



- 1. All aardvarks are insectivores.
- 2. All Orycteropodidae are aardvarks.
- 3. 90%: All Orycteropodidae are insectivores.
- 4. 5%: Some Orycteropodidae are insectivores.
- 5. 5%: Others, including erroneous.



- ► All-Some: `All A are B' implies `Some A are B'.
- No-Some not: `No A are B' implies `Some A are not B'.
- Conversion1: `Some A are B' implies `Some B are A';
- ► Conversion2: `No A are B' implies `No B are A".
- ► Monotonicity rule.



Zhai, Titov, Szymanik. Toward a probabilistic mental logic for the syllogistic fragment of natural language. Amsterdam Colloquium, 2015

HOW MUCH SUCCESS SHOULD WE EXPECT IN INFERRING PRIMITIVES?

STRATEGY

- ► Pick a conceptual domain.
- Define a space of possible LoTs.
- ➤ Define a way that LoT can influence behaviour, e.g., through category learning via a simplicity bias.
- Run simulated experiments with known LoTs and see whether it's possible to recover the underlying LoT from the simulated behavioural data accurately.

| Table 1 | |
|---|---|
| Glossary of technical terms and syn | nbols as used in the computational model. |
| Property: | A binary property, e.g. 'being red'. |
| Object: | A set of properties. |
| Category: | A set of objects. |
| Operator: | A function from Booleans to a Boolean. |
| Language of Thought (LoT): | A set of operators. |
| <i>p</i> , <i>q</i> , <i>r</i> , <i>s</i> : | The four properties in the model. |
| <i>O</i> : | The ordered tuple of 16 objects. |
| Ω : | The set of 9 operators. |

- ► 4 binary properties
- ► An object is a set of properties.
- ► A category is a set of objects
- LoTs are functionally complete, non-redundant subsets of 16 Boolean binary operators.

CATEGORY GENERALIZATION DESIGN





SIMULATED EXPERIMENT 1

- The agent sees some examples of objects in the 'true' category (which is, in fact, a random selection of possible objects).
- The agent calculates the posterior over categories given the observed objects and the agent's true LoT.
- Then the agent needs to decide whether the remaining objects belong to the category (by summing the posterior probability of all categories that contain the object).
- As experimenters, we calculate the simulated experimenter's posterior over LoT given the participants' categorization data.
- Prior is given by the minimal formula in the LoT for a category.

25776 SIMULATED EXPERIMENTS

- ► 4 properties
- ► 358 possible LoTs
- ► 65536 categories

| Name | Abbreviation | Symbol |
|-----------------------|--------------|-------------------|
| Conjunction (and) | Α | ٨ |
| Disjunction (or) | 0 | V |
| Conditional | С | \rightarrow |
| Negated conditional | NC | / > |
| Biconditional | В | \leftrightarrow |
| Negated biconditional | XOR | |
| Negated conjunction | NAND | \checkmark |
| Negated disjunction | NOR | ¥ |
| Negation | Ν | 7 |

► Number of examples given in the experiment: 1, 5, 10, 15

.

- ► Number of participants: 1, 10, 30, 60, 120, 250
- Simplicity bias strength: 0.5, 1, 3
- What we look at for each combination of parameter values is the distribution of the simulated experimenter's posterior entropies across LoTs.



| LoT | # recovered | LoT | # recovered |
|-----------------------------------|-------------|--|-------------|
| ∧, ¬ | 19 | ∧, ↔, ∦ | 3 |
| \lor, \neg | 14 | $\leftrightarrow, \not\!\!\!/$ | 3 |
| ^, ≯ | 14 | $\lor, \neg, \not\rightarrow$ | 3 |
| \checkmark | 11 | \neg, \rightarrow | 3 |
| $\land, \not\!\!\!/, \not\!\!\!/$ | 8 | $\lor,\land,\neg,\not\rightarrow$ | 3 |
| _, ≯ | 6 | $\land, \neg, \not\leftrightarrow$ | 3 |
| $\wedge,\neg,\leftrightarrow$ | 5 | ∧, ¬, ≯ | 3 |
| $\lor,\neg,\leftrightarrow$ | 5 | $\land, \not\rightarrow$ | 2 |
| ∨, ≯ | 5 | \lor, \land, \checkmark | 2 |
| X | 4 | $\land,\leftrightarrow,\not\!$ | 2 |
| $\leftrightarrow,\not\rightarrow$ | 4 | $\neg, \not\leftrightarrow, \not\rightarrow$ | 2 |
| $\neg,\not\rightarrow$ | 4 | ₩, ∅, У | 2 |
| \lor, \land, \lnot | 4 | $\land, \leftrightarrow, \not\!\!\!/$ | 2 |
| $\lor, \not\!\!\!/, \not \to$ | 3 | \wedge,\neg,\rightarrow | 2 |
| $\not\!\!\!/,\not\rightarrow$ | 3 | $\lor, \neg, \not\leftrightarrow$ | 2 |

Table 4: 30 LoTs that were recovered most often, along with the number of times they were recovered across all simulated experiments. This gives a sense of what LoTs are easiest to recover amongst all LoTs.

RESULTS

- The more participants, the better the recovery.
- The stronger the simplicity bias, the more recoverable LoTs are in general.
- The fewer operators in the LoT, the more recoverable it is.
- LoTH may be empirically productive only if the postulated LoTs are simple.
- Therefore, we require more and better theory, not just brute force search within the LoT space.

Carcassi & Szymanik. The Boolean Language of Thought is recoverable from learning data, Cognition, 2023.

SERIAL AND DYNAMIC DESIGN

tracking learning trajectory



SIMULATED EXPERIMENT 2

- The agent sees each of 16 objects, judges whether it belongs to the category, and gets feedback.
- So, agents receive both positive and negative evidence.
- Serial design: we select both category and order by randomly sampling
- Dynamic design: we choose the next object to show using the greedy Bayesian Optimal Design
- Optional stopping: up to 250 participants or posterior probability > 0.95

1432 SIMULATED EXPERIMENTS

.

- ► 4 properties
- ► 358 possible LoTs
- ► 65536 categories
- ➤ Simplicity bias strength: 0.5, 1, 3
- Serial and dynamic design



Fig. 7. Proportions of each outcome out of the 358 experiments that were ran for each level of simplicity bias strength λ . Results are shown both for the dynamic and for the serial design. The main result is that with these more sophisticated designs, the true LoT can reliably be recovered from learning data, as long as the simplicity bias is strong enough.

RESULTS

- ► With a strong simplicity bias, the majority of the true LoTs are recovered.
- ► The misidentification is low.
- Dynamic design doesn't help.
- ➤ The fewer operators in the LoT, the more recoverable it is.
- ► LoTH may be a productive endeavor, but doing it in reality will be much harder

POST-SCRIPTUM: WHAT IF COGNITION IS ALL NON-SYMBOLIC?

DOES THAT EVEN MATTER?



Carcassi & Szymanik. Neural Networks track the logical complexity of Boolean concepts, Open Mind 2022.

LOT VS ANNS

- Do LoTH and connectionism have the same empirical import in the domain of categorization?
- Do they make the same predictions about the effort required to acquire categories?
- LoT: The categories with the shortest minimal formulas are the easiest.
- Connectionism: the average loss across epochs and batches
- There is an overall positive rank correlation between logical complexity and ANN learning effort

CONCLUSIONS

- ► LoT is often the engine of computational cognitive models
- Via the notions of complexity/simplicity prior
- Leading to interesting theoretical and empirical insights
- ➤ To unify those models, we need to recover the true LoT
- Such LoT could help AI get closer to human-like intelligence

THANK YOU!