

DISENTANGLING AND INTEGRATING RELATIONAL AND SENSORY INFORMATION IN TRANSFORMER ARCHITECTURES

Awni Altabaa, John Lafferty

June 5, 2025

Yale University

arXiv:2405.16727, ICML '25

BIG PICTURE: WHY SHOULD WE CARE ABOUT “RELATIONAL REASONING”?

Hypothesis 0: Human & animal intelligence can be explained by a few core principles (rather than an encyclopedic list of heuristics)

Suggests the following goal: Study & uncover the inductive biases that humans & animals exploit to understand intelligence generally and inform design of AI

Deep learning systems themselves exploit several key inductive biases that underly their empirical success

Goal of AI Research: Uncover a core set of inductive biases for DL that enable data-efficient learning and reasoning over wide range of tasks and modalities

Hypothesis 1: Relational reasoning is one of these fundamental principles of intelligence

FIRST: WHAT IS “RELATIONAL REASONING”?

FIRST: WHAT IS “RELATIONAL REASONING”?

Reasoning about **relationships** between objects and how they interact in a given context/scene

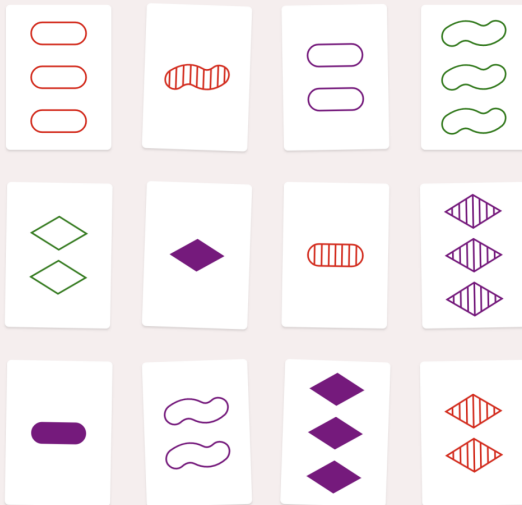
Perform **comparisons** under different attributes or features, at multiple levels of abstraction

Beyond recognizing individual objects by sensory pattern recognition; requires **higher-order** relationships

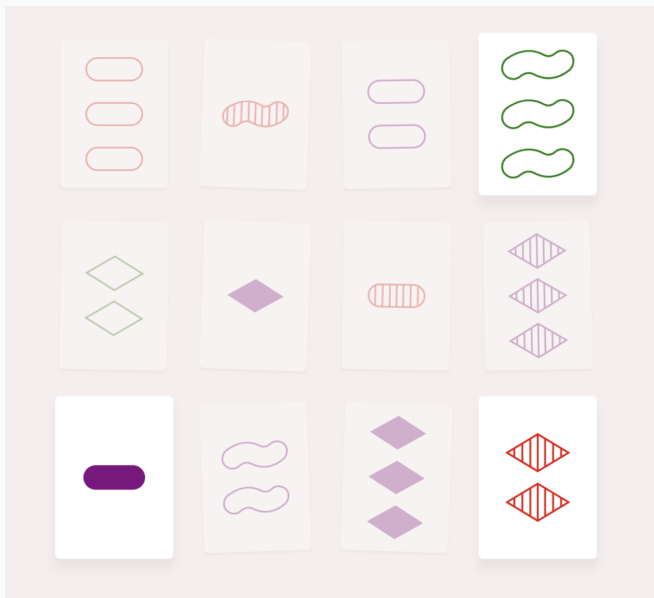
Clue to its importance: Humans have a natural ability (and a preference) to do relational reasoning

**LET'S WALK THROUGH A COUPLE SIMPLE
ILLUSTRATIVE EXAMPLES OF RELATIONAL TASKS**

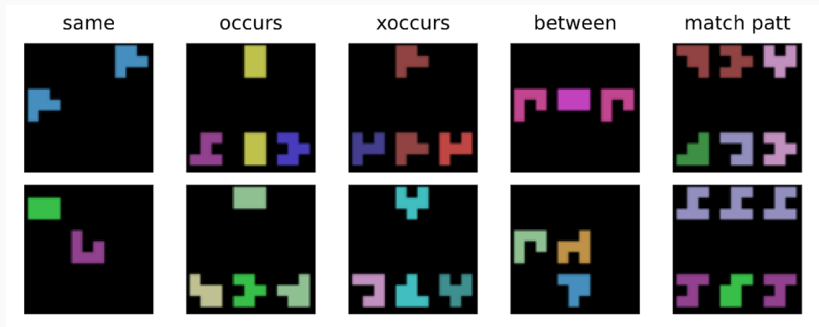
EXAMPLE: SET! CARD GAME



EXAMPLE: *SET!* CARD GAME



EXAMPLE: RELATIONAL GAMES (SHANAHAN ET AL. 2020)



Relational Games tasks from Shanahan et al. (2020)

A **Visual Relational Reasoning Task**:

determine whether a particular relation holds or not

**RETURNING TO OUR ORIGINAL QUESTION:
WHY SHOULD WE CARE ABOUT RELATIONAL
REASONING?**

WHY SHOULD WE CARE ABOUT “RELATIONAL REASONING”?

A cornerstone of human intelligence

Underlies capabilities for

- *analogy*
- *abstraction*
- *generalization*

By relating new inputs to previously-seen stimuli, we form *analogies* and *abstractions* that allow us to *systematically generalize*.

***“IN THE LIMIT, RELATIONAL REASONING YIELDS
UNIVERSAL INDUCTIVE GENERALIZATION FROM A
FINITE AND OFTEN VERY SMALL SET OF OBSERVED
CASES TO A POTENTIALLY INFINITE SET OF NOVEL
INSTANCES.” — GOYAL & BENGIO (2022)***

**WE'D LIKE TO TAKE A STEP TOWARDS THIS CENTRAL
GOAL OF AI RESEARCH**

OUTLINE OF REMAINDER OF TALK

Big Picture: Why should we care about “relational reasoning”?

Main Idea & Goal

Transformers: The Sensory and the Relational

Relational Attention

Dual Attention Transformer Architecture

Empirical Investigation

Concluding Remarks

MAIN IDEA & GOAL

Our Goal: Make progress towards a universal neural architecture with explicit relational computational mechanisms & inductive biases

HOW TO IMBUE TRANSFORMERS WITH EXPLICIT RELATIONAL INDUCTIVE BIASES

- *Inductive Bias*: intrinsic preferences over solution space
- *View*: Two types
 - *Additive*: imbue architecture with mechanism, and let it learn to use it
 - *Subtractive*: constrain the space of representations a model can compute

- “*The Bitter Lesson*” — Rich Sutton
- Relational computational mechanisms *parameterized by neural net & learned*
 - scalable, general mechanisms;
 - avoid domain-specific heuristic, human-engineering
- The *versatility* of the *Transformer architecture* suggests it may form a powerful starting point

SOME LESSONS FROM PREVIOUS WORK

Prior works on relational inductive biases

- Santoro et al. “A simple neural network module for relational reasoning” (2017)
- Shanahan et al. “An Explicitly Relational Neural Network Architecture” (2020)
- Kerg et al. “Inductive biases for relational tasks” (2022)
- Others...

Data-efficient relational reasoning requires inductive biases

- Standard neural models (e.g., Transformers) are *data-inefficient* at learning relational tasks; brittle OOD generalization
- **Hypothesized Explanation:** Neural networks overemphasize *individual object* representations while lacking explicit mechanisms for encoding and processing *relational* features.
- **Common thread explored:** constrain model to compute relational features—relational inductive biases

TENSION: GENERALITY VS. INDUCTIVE BIASES

However, these models are narrow in domain

They improve relational processing, but lose generality

Empirical success limited to synthetic (purely relational) benchmarks

**OUR GOAL: AUGMENT THE TRANSFORMER
ARCHITECTURE WITH EXPLICIT RELATIONAL
MECHANISMS & INDUCTIVE BIASES**

TRANSFORMERS: THE SENSORY AND THE RELATIONAL

HOW TO IMBUE TRANSFORMERS WITH EXPLICIT RELATIONAL INDUCTIVE BIASES

Strength of Transformers: *attention*

Versatile *information retrieval* mechanism

Composable in *circuits* to carry out complex computation
(which we're now beginning to understand through systematic (mechanistic) interpretability work)

THE TRANSFORMER ARCHITECTURE, ESSENTIALLY

Iterate two basic operations:

1. **Information Retrieval:** Attention

$$x'_i \leftarrow \sum_j \alpha_{ij} \phi_v(x_j)$$

2. **Local Processing:** Token-wise MLP

$$x'_i \leftarrow \text{MLP}(x_i)$$

1. Compute attention scores

$$\alpha_{ij} = \text{Softmax}([\langle \phi_q^{\text{attn}}(\mathbf{x}_i), \phi_k^{\text{attn}}(\mathbf{x}_j) \rangle]_{j=1}^n)_j$$

2. Retrieve weighted combination of sensory values in context

$$\mathbf{e}_i \leftarrow \sum_j \alpha_{ij} \phi_v(\mathbf{x}_j)$$

TWO TYPES OF INFORMATION

Fundamentally, attention is an information retrieval operation

Two key *types* of information

Sensory: features or attributes of individual objects

Relational: relationships between objects

Standard attention captures the former, but not the latter

TWO TYPES OF ATTENTION

Correspondingly, there ought to be two types of attention

(Standard) Sensory Attention:

retrieval of *sensory* information in context

Relational Attention:

retrieval of *relational* information in context

RELATIONAL ATTENTION

HIGH-LEVEL: RELATIONAL ATTENTION

1. Attend
2. Relate
3. Tag with symbols

1) ATTENTION

Same as standard (sensory attention)

Compute attention scores via learned query/key maps

$$\alpha_{ij} = \text{Softmax}([\langle \phi_q^{\text{attn}}(x_i), \phi_k^{\text{attn}}(x_j) \rangle]_{j=1}^n)_j$$

2) COMPUTING RELATIONS

Relation vector, representing a series of comparisons under different attributes or extracted features

Computed as a series inner products under different learned feature maps

$$\mathbf{r}_{ij} = \left(\left\langle \phi_{q,\ell}^{\text{rel}}(\mathbf{x}_i), \phi_{k,\ell}^{\text{rel}}(\mathbf{x}_j) \right\rangle \right)_{\ell \in [d_r]} \in \mathbb{R}^{d_r}$$

3) SYMBOLS

Tag each object in the context with an *symbol*

$$(s_1, \dots, s_n) = \text{SymbolRetriever}(x_1, \dots, x_n)$$

Serve as reference/pointer/identifier of selected object with whom the relation is with, abstracted away from high-dimensional sensory features

We experiment with different symbol assignment mechanisms: positional, relative positional, “soft-equivalence class”

RELATIONAL ATTENTION: PUTTING IT ALL TOGETHER

Putting it all together

$$\mathbf{a}_i \leftarrow \sum_j \alpha_{ij} \cdot (W_r \mathbf{r}_{ij} + W_s \mathbf{s}_j)$$

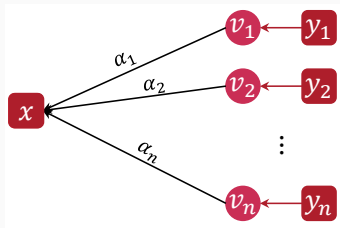
α_{ij} : attention scores — govern selection criterion

\mathbf{r}_{ij} : relation vector — relational information

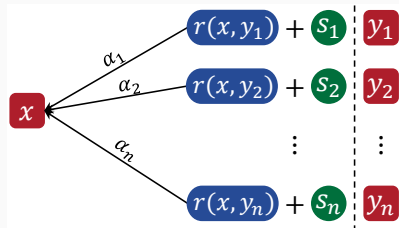
\mathbf{s}_j : symbol — identifier of source/sender object

W_r, W_s : learned linear maps — organize information in residual stream

SENSORY & RELATIONAL ATTENTION



$\text{Attention}(x, (y_1, \dots, y_n))$



$\text{RelationalAttention}(x, (y_1, \dots, y_n))$

A FEW COMMENTS...

Causal masking

Positional encoding

Symmetric relations

Computational complexity

DUAL ATTENTION TRANSFORMER ARCHITECTURE

Relational attention : a mechanism for routing relational information

Both *sensory* and *relational* information are crucial for reasoning over collections or sequences of objects.

Dual Attention Transformer (DAT): A variant of the Transformer architecture that routes both types of information in the information retrieval step.

Introduces explicit relational processing mechanisms, while retaining sensory processing capabilities.

Dual Attention is a variant of multi-head attention with *two types of attention heads: sensory* and *relational* .

DUAL ATTENTION

Algorithm 1: Dual Attention

Input: $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^{n \times d}$

Compute self-attention heads

$$\begin{aligned}\alpha^{(h)} &\leftarrow \text{Softmax}((\mathbf{x} W_{q,h}^{\text{attn}})(\mathbf{x} W_{k,h}^{\text{attn}})^\top), \quad h \in [n_h^{sa}] \\ e_i^{(h)} &\leftarrow \sum_j \alpha_{ij}^{(h)} \mathbf{x}_j W_v^h, \quad i \in [n], h \in [n_h^{sa}] \\ \mathbf{e}_i &\leftarrow \text{concat}(e_i^{(1)}, \dots, e_i^{(n_h^{sa})}) W_o^{sa}, \quad i \in [n]\end{aligned}$$

Assign symbols:

$$\mathbf{s} = (s_1, \dots, s_n) \leftarrow \text{SymbolRetriever}(\mathbf{x}; S_{\text{lib}})$$

Compute relational attention heads

$$\begin{aligned}\alpha^{(h)} &\leftarrow \text{Softmax}((\mathbf{x} W_{q,h}^{\text{attn}})(\mathbf{x} W_{k,h}^{\text{attn}})^\top), \quad h \in [n_h^{ra}] \\ \mathbf{r}_{ij} &\leftarrow (\langle x_i W_{q,\ell}^{\text{rel}}, x_j W_{k,\ell}^{\text{rel}} \rangle)_{\ell \in [d_r]} \quad i, j \in [n] \\ \mathbf{a}_i^{(h)} &\leftarrow \sum_j \alpha_{ij}^{(h)} (\mathbf{r}_{ij} W_r^h + \mathbf{s}_j W_s^h), \quad i \in [n], h \in [n_h^{ra}] \\ \mathbf{a}_i &\leftarrow \text{concat}(\mathbf{a}_i^{(1)}, \dots, \mathbf{a}_i^{(n_h^{ra})}) W_o^{ra}, \quad i \in [n]\end{aligned}$$

Output: $(\text{concat}(\mathbf{e}_i, \mathbf{a}_i))_{i=1}^n$

DUAL ATTENTION TRANSFORMER: ENCODER & DECODER

Algorithm 2: Dual Attention Encoder Block

Input: $x \in \mathbb{R}^{n \times d}$

$x \leftarrow \text{Norm}(x + \text{DualAttn}(x))$

$x \leftarrow \text{Norm}(x + \text{MLP}(x))$

Output: x

Algorithm 3: Dual Attention Decoder Block

Input: $x, y \in \mathbb{R}^{n \times d}$

$x \leftarrow \text{Norm}(x + \text{DualAttn}(x))$

$x \leftarrow \text{Norm}(x + \text{CrossAttn}(x, y))$

$x \leftarrow \text{Norm}(x + \text{MLP}(x))$

Output: x

EMPIRICAL INVESTIGATION

PRELUDE: WHAT QUESTIONS ARE WE TRYING TO ANSWER?

How does the *DAT* perform on synthetic relational benchmarks?

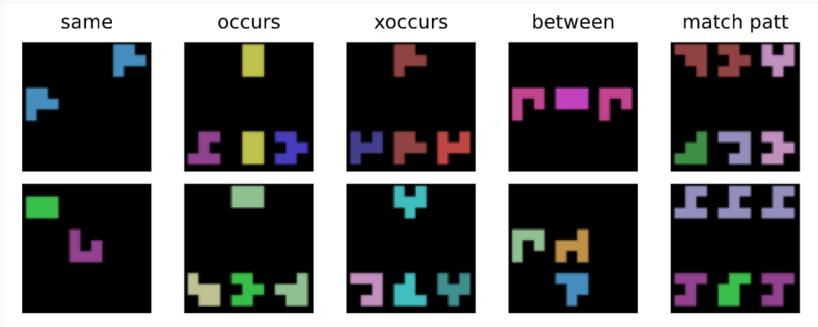
Data efficiency

Scalability with data and model size (recall: bitter lesson)

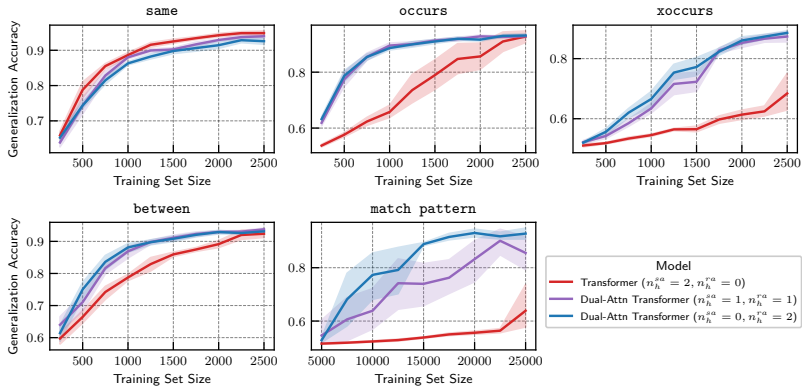
Applicability to complex real-world tasks; versatility across **data modalities** (language & vision)

SYNTHETIC RELATIONAL BENCHMARKS: RELATIONAL GAMES (SHANAHAN ET AL. 2020)

SYNTHETIC RELATIONAL TASKS: TASK



SYNTHETIC RELATIONAL TASKS: RESULTS



MATHEMATICAL PROBLEM-SOLVING (SEQ2SEQ)

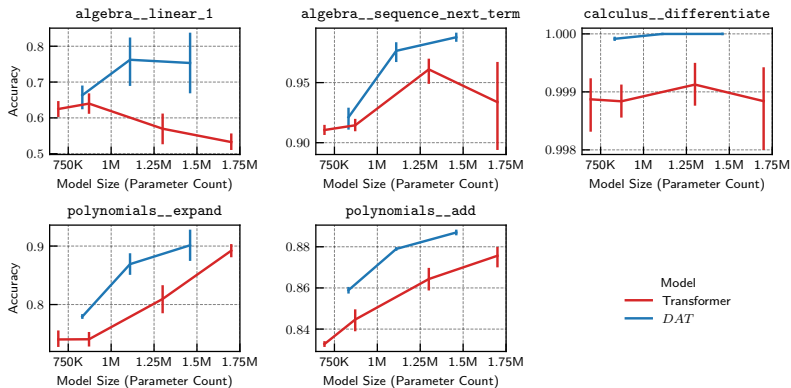
MATHEMATICAL PROBLEM-SOLVING (SEQ2SEQ): TASK

Dataset due to Saxton et al. (2019)

Modeled as char-level Sequence-to-Sequence task with
encoder-decoder architecture

Module	Math Dataset Example
<i>algebra_linear_1d</i>	Q: Solve for x : $3x + 7 = 19$ A: $x = 4$
<i>algebra_sequence_next_term</i>	Q: What is the next term in the sequence 2, 5, 8, 11, ...? A: 14
<i>calculus_differentiate</i>	Q: Find the derivative of $f(x) = 3x^2 + 2x - 5$ with respect to x . A: $6x + 2$
<i>polynomials_expand</i>	Q: Expand $(2x + 3)(x - 1)$. A: $2x^2 + x - 3$
<i>polynomials_add</i>	Q: Add the polynomials: $(2x^2 + 3x + 1) + (x^2 - 2x + 4)$ A: $3x^2 + x + 5$

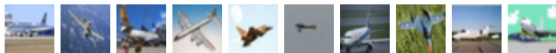
MATHEMATICAL PROBLEM-SOLVING (SEQ2SEQ): RESULTS



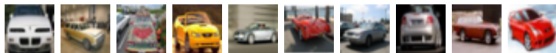
VISUAL PROCESSING (CIFAR)

VISUAL PROCESSING (CIFAR): TASK

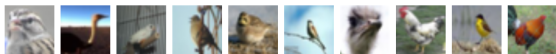
airplane



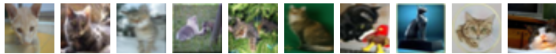
automobile



bird



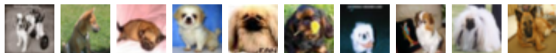
cat



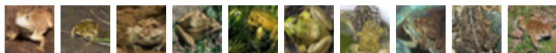
deer



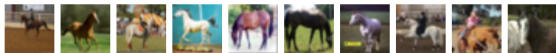
dog



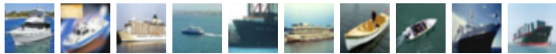
frog



horse



ship



truck



VISUAL PROCESSING (CIFAR): RESULTS

ViT-style encoder-only architecture processing image as sequence of patches

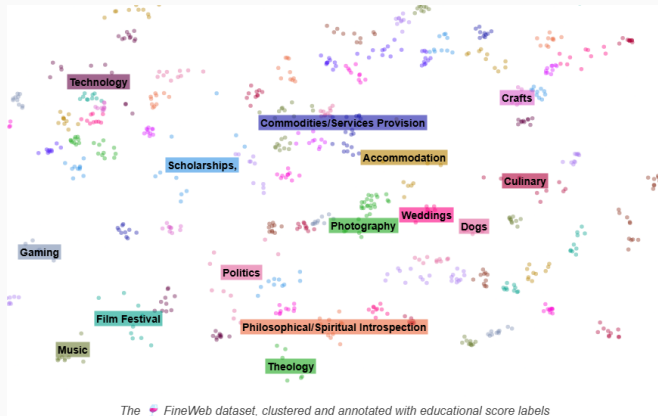
Dataset	Model	Params	Accuracy
CIFAR-10	<i>ViT</i>	7.1M	$86.4 \pm 0.1\%$
	<i>ViDAT</i>	6.0M	$89.7 \pm 0.1\%$
CIFAR-100	<i>ViT</i>	7.2M	$68.8 \pm 0.2\%$
	<i>ViDAT</i>	6.1M	$70.5 \pm 0.1\%$

LANGUAGE MODELING

LANGUAGE MODELING: TASK

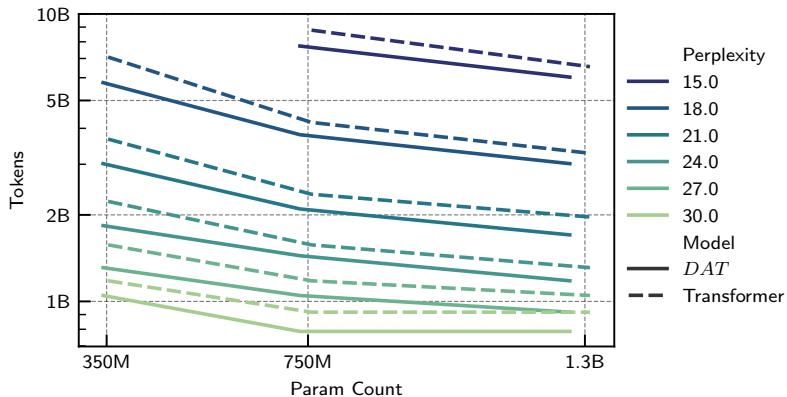
Autoregressive causal language modeling with a “decoder-only” architecture

Use the Fineweb-Edu dataset (curated high-quality text data);
train on 10B tokens



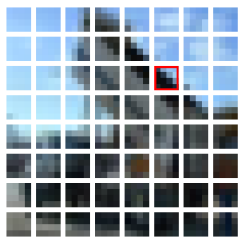
LANGUAGE MODELING: RESULTS

Evaluate scaling with data and model size



A BIT OF VISUALIZATION/INTERPRETATION

INTERPRETING *ViDAT* MODEL



(a) Original Image

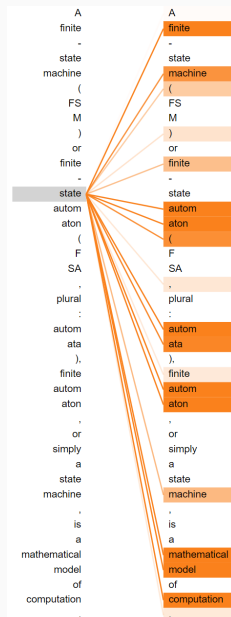
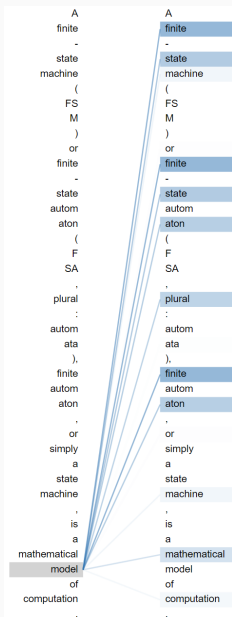


(b) A Relation in the First Layer



(c) A Relation in the Fifth Layer

INTERPRETING DAT LANGUAGE MODELS



CONCLUDING REMARKS

CONCLUDING REMARKS

Relational reasoning is a core facet of human intelligence, underpinning abilities for analogy, abstraction, and generalization

It is likely an important component of artificial intelligence as well

In this work, we took a step towards developing neural architectures with enhanced relational processing capabilities, while retaining powerful sensory processing

CONCLUDING REMARKS: FUTURE WORK

Interpretability:

- How is *DAT* learning to use its relational processing mechanisms?
- Can specific “circuits” be identified?
- How does *DAT* achieve improved data efficiency in different tasks?

Iterate & tweak architecture; find good choices for hyperparameters

Computational considerations: optimize implementation

THANK YOU

DISCUSSION TIME...

- **Joint work** with *John Lafferty*
- **Supported by funding from** ARNI NSF AI Institute
- **Paper:** *arXiv:2405.16727 / ICML '25*
- **Project webpage:** *<https://awni.xyz/dual-attention/>*
 - Open weights on HF (DAT-LM up to 1.3B-params)
 - Implementation available via python package
`pip install dual-attention`
- **Personal webpage:** *<https://awni.xyz>*