# DISENTANGLING AND INTEGRATING RELATIONAL AND SENSORY INFORMATION IN TRANSFORMER ARCHITECTURES

Awni Altabaa, John Lafferty

May 30, 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"?

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?

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

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;
  - $\circ$  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-<u>in</u>efficient 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

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) 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

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

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

$$e_i \leftarrow \sum_j lpha_{ij} \phi_v(x_j)$$

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

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

18

# **Relational Attention**

- 1. Attend
- 2. Relate
- 3. Tag with symbols

## Same as standard (sensory attention)

## Compute attention scores via learned query/key maps

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

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

Computed as a series inner products under different learned feature maps

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

Tag each object in the context with an symbol

$$(s_1,\ldots,s_n) =$$
SymbolRetriever $(x_1,\ldots,x_n)$ 

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"

## Putting it all together

$$\boldsymbol{a}_i \leftarrow \sum_j \alpha_{ij} \cdot (W_r \, \boldsymbol{r}_{ij} + W_s \, s_j)$$

 $\alpha_{ij}$ : attention scores — govern selection criterion  $r_{ij}$ : relation vector — relational information  $s_j$ : symbol — identifier of source/sender object  $W_r, W_s$ : learned linear maps — organize information in residual stream

- 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

$$\begin{split} \boldsymbol{\alpha}^{(h)} &\leftarrow \operatorname{Softmax}((\boldsymbol{x} W_{q,h}^{\operatorname{attn}})(\boldsymbol{x} W_{k,h}^{\operatorname{attn}})^{\mathsf{T}}), \quad h \in [n_h^{sa} \\ \boldsymbol{e}_i^{(h)} &\leftarrow \sum_j \alpha_{ij}^{(h)} \boldsymbol{x}_j W_v^h, \quad i \in [n], h \in [n_h^{sa}] \\ \boldsymbol{e}_i &\leftarrow \operatorname{concat}(\boldsymbol{e}_i^{(1)}, \dots, \boldsymbol{e}_i^{(n_h^{sa})}) W_o^{sa}, \quad i \in [n] \end{split}$$

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

Compute relational attention heads

$$\begin{split} \boldsymbol{\alpha}^{(h)} &\leftarrow \operatorname{Softmax}\left((\boldsymbol{x} \ \boldsymbol{W}^{\operatorname{ath}}_{q,h})(\boldsymbol{x} \ \boldsymbol{W}^{\operatorname{ath}}_{k,h})^{\mathsf{T}}\right), \quad h \in [n_{h}^{ra}] \\ \boldsymbol{r}_{ij} &\leftarrow \left(\langle x_{i} \ \boldsymbol{W}^{\operatorname{el}}_{q,\ell}, x_{j} \ \boldsymbol{W}^{\operatorname{rel}}_{k,\ell}\rangle\right)_{\ell \in [d_{r}]} \quad i,j \in [n] \\ \boldsymbol{a}^{(h)}_{i} &\leftarrow \sum_{j} \boldsymbol{\alpha}^{(h)}_{ij} \left(r_{ij} \ \boldsymbol{W}^{\mathsf{h}}_{r} + s_{j} \ \boldsymbol{W}^{\mathsf{h}}_{s}\right), \quad i \in [n], \ h \in [n_{h}^{ra} \\ \boldsymbol{a}_{i} &\leftarrow \operatorname{concat}\left(a^{(1)}_{i}, \dots, a^{(n_{h}^{ra})}_{i}\right) \boldsymbol{W}^{ra}_{o}, \quad i \in [n] \end{split}$$

**Output:**  $\left(\operatorname{concat}(\boldsymbol{e}_i, \boldsymbol{a}_i)\right)_{i=1}^n$ 

#### Algorithm 2: Dual Attention

Encoder Block

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

 $x \leftarrow \operatorname{Norm}(x + \operatorname{DualAttn}(x))$  $x \leftarrow \operatorname{Norm}(x + \operatorname{MLP}(x))$ 

#### Output: x

# Algorithm 3: Dual Attention Decoder Block Input: $x, y \in \mathbb{R}^{n \times d}$

 $\begin{array}{l} \textbf{\textit{x}} \leftarrow \operatorname{Norm}(\textbf{\textit{x}} + \operatorname{DualAttn}(\textbf{\textit{x}})) \\ \textbf{\textit{x}} \leftarrow \operatorname{Norm}(\textbf{\textit{x}} + \operatorname{CrossAttn}(\textbf{\textit{x}}, \textbf{\textit{y}})) \\ \textbf{\textit{x}} \leftarrow \operatorname{Norm}(\textbf{\textit{x}} + \operatorname{MLP}(\textbf{\textit{x}})) \end{array}$ 

Output: x

# **EMPIRICAL INVESTIGATION**

- 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)

### Dataset due to Saxton et al. (2019)

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

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

## MATHEMATICAL PROBLEM-SOLVING (SEQ2SEQ): RESULTS



# VISUAL PROCESSING (CIFAR)

# VISUAL PROCESSING (CIFAR): TASK

airplane	🛁 🐹 🔛 🛩 🖛 🗾 😹 🛶 🏎
automobile	an 😂 🚞 💁 🔤 😻 🚞 📾 🛸
bird	in 19 19 19 19 19 19 19 19 19 19 19 19 19
cat	in i
deer	🗱 🔛 🖌 🐖 🎆 💱 🕅 📰 🌉
dog	93 🔬 🤜 🔛 🌊 🏹 💽 🎎
frog	Ref 🖉 🔀 😭 🚱 🌆 💷 🔤
horse	🌁 🗶 🚰 法 🕅 🕅 🖄 🐼 🐞
ship	🚔 🌽 🚈 📥 🚢 🚁 🌽 🖉 🖄
truck	🐳 🌃 🚛 🌉 👹 🔤 📷 🖓 🕋 🚮

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

Dataset	Model	Params	Accuracy
CIFAR-10	ViT	7.1M	$86.4\pm0.1\%$
	ViDAT	6.0M	$89.7 \pm \mathbf{0.1\%}$
CIFAR-100	ViT	7.2M	$68.8\pm0.2\%$
	ViDAT	6.1M	$70.5 \pm \mathbf{0.1\%}$

# LANGUAGE MODELING

# 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**





39

# **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

## Interpretability:

 $\circ$  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**

- 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