Abstractors and relational cross-attention: An inductive bias for explicit relational reasoning in Transformers

ICLR 2024

Awni Altabaa 1 , Taylor W. Webb 2 , Jonathan D. Cohen 3 , John Lafferty 1 Yale, 2 UCLA, 3 Princeton

High-level goal

- 1) Understand why the Transformer architecture struggles with sample-efficiently learning relational tasks;
- 2) Design mechanisms for explicit relational reasoning within the broader Transformer framework

Recall: standard self-attention takes the form

$$(x_1, \dots, x_n) \mapsto (y_1, \dots, y_n)$$
$$y_i = \sum_{j=1}^n \alpha_{ij} \ \phi_v(x_j),$$
$$\alpha_{i\cdot} = \text{Softmax} \left(\left[\langle r(x_i, x_j) \rangle \right]_{j=1}^n \right)$$

The attention scores can be thought of as relations between pairs of objects which determine which object to "attend to" and retrieve information from.

In standard self-attention, computing these relations is merely an intermediate step in an information-retrieval operation.

The relations α_{ij} are entangled with object-level features $\phi_v(x_j)$.

$$|\mathsf{Object}| \gg |\mathsf{Relation}|$$
 space

Prevents abstraction and efficient learning of relational features. Would require somehow "marginalizing" over variability in irrelevant object-level features.

Under the neural-message passing lens, self-attention is

$$\begin{aligned} x_i' &\leftarrow \text{Aggregate}(\{m_{j \to i} : j \in [n]\}), \\ m_{j \to i} &= (r(x_i, x_j), \phi_v(x_j)), \\ r(x_i, x_j) &= \langle \phi_q(x_i), \phi_k(x_j) \rangle \end{aligned}$$

We propose a modification of this where the object-level features $\phi_v(x_j)$ are replaced with vectors s_j we call "symbols", which identify objects, but do not encode their features.

Symbols live in a space with much smaller variability, inducing a relation-centric representation.

 $^0\mathrm{Aggregate}$ has the particular form $\bigoplus_j m_{j \to i}$, where where \oplus is a binary operation on a commutative monoid given by

$$(r_1, x_1) \oplus (r_2, x_2) = (\exp(r_1) + \exp(r_2), \frac{\exp(r_1)x_1 + \exp(r_2)x_2}{\exp(r_1) + \exp(r_2)})$$

Relational cross-attention & the Abstractor

Relational cross-attention

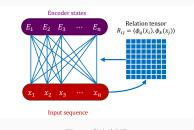
Relational cross-attention then implements the following "neural message-passing" operation,

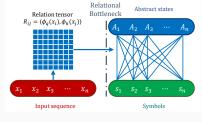
$$x'_{i} \leftarrow \text{Aggregate}(\{m_{j \to i} : j \in [n]\})$$

 $m_{j \to i} \leftarrow (r(x_{i}, x_{j}), s_{j})$
 $s_{1}, ..., s_{n} \leftarrow \text{SymbolAssignment}(x_{1}, ..., x_{n})$

with the "symbol" s_i identifying through position, relative position, and/or syntactic role.

Pictoral depiction of SA & RCA





$$E \leftarrow \mathrm{SA}(X)$$

 $A \leftarrow \text{RCA}(X, S)$

Comparison of relational cross-attention (RCA) with self-attention (SA). RCA implements a relational information bottleneck, which results in relational representations that are disentangled and separated from object-level features.

The "Abstractor" Module

By composing relational cross-attention with an MLP, analogously to a Transformer block, we obtain a natural formulation of a composable neural module for relational processing which fits nicely within the broader Transformer framework.

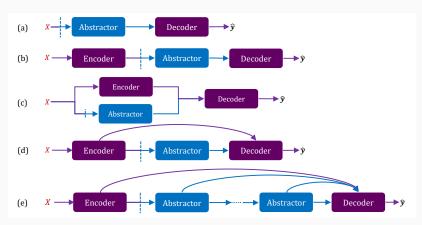
7

Abstractor-supported architectures

within the Transformer framework

Abstractor-supported architectures

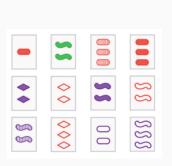
The Abstractor module can be incorporated into a broader Transformer-based architecture to support enhanced relational processing. The choice of architecture should depend on the type of relational processing required by the underlying task.



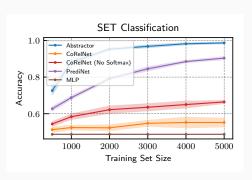
Experiments

Discriminative Relational Tasks: SET!

We compare to existing "relational architectures" which have been proposed for discriminative tasks.



Demo of SET! task

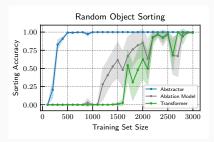


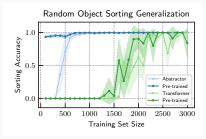
Learning curves of different models

The Abstractor is more versatile and more sample-efficient.

Synthetic Sequence-to-Sequence Relational Tasks: "Object-sorting"

Setting: Randomly assign an order on objects $o \in [N], N = 48$. **Task:** Learn to sort sequences of randomly shuffled objects. Relevent relation is order \prec .





Learning curves on sorting sequences of random objects.

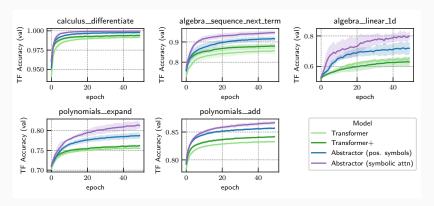
The abstractor is dramatically more sample-efficient.

Learning curves with and without pre-training on a similar sorting task. The Abstractor generalizes & benefits from pre-training.

Mathematical problem-solving

Task: Character-level sequence-to-sequence task. Given a mathematical question, predict the answer.

E.g.: expand
$$(3*x + 1)*(2*x - 5) \rightarrow 6*x**2 - 13*x - 5$$



The Abstractor learns faster and reaches higher accuracy

Thank you