TRANSFORMERS ARE EFFICIENT COMPILERS, PROV ABLY

Anonymous authors

Paper under double-blind review

ABSTRACT

011 Transformer-based large language models (LLMs) have demonstrated surprisingly 012 robust performance across a wide range of language-related tasks, including pro-013 gramming language understanding and generation. In this paper, we take the first steps towards a formal investigation of using transformers as compilers from an 014 expressive power perspective. To this end, we introduce a representative pro-015 gramming language, Mini-Husky, which encapsulates key features of modern 016 C-like languages. We show that if the input code sequence has a bounded depth 017 in both the Abstract Syntax Tree (AST) and type inference (reasonable assump-018 tions based on the *clean code principle*), then the number of parameters required 019 by transformers depends only on the *logarithm of the input sequence length* to handle compilation tasks, such as AST construction, symbol resolution, and type 021 analysis. A significant technical challenge stems from the fact that transformers operate at a low level, where each layer processes the input sequence as raw vectors without explicitly associating them with predefined structure or meaning. 024 In contrast, high-level compiler tasks necessitate managing intricate relationships and structured program information. Our primary technical contribution is the 025 development of a domain-specific language, **Cybertron**, which generates formal 026 proofs of the transformer's expressive power, scaling to address compiler tasks. 027 We further establish that recurrent neural networks (RNNs) require at least a lin-028 ear number of parameters relative to the input sequence, leading to an exponential 029 separation between transformers and RNNs. Finally, we empirically validate our theoretical results by comparing transformers and RNNs on compiler tasks within 031 Mini-Husky. 032

033 034

003 004

010

1 INTRODUCTION

Transformers (Vaswani, 2017) have demonstrated remarkable proficiency across various do mains, achieving near-expert performance in solving International Mathematical Olympiad prob lems (Google Deepmind, 2024) and excelling in complex reasoning tasks in science, coding, and
 mathematics (OpenAI, 2024a). They also handle routine coding tasks with high precision and have
 been integrated into code editors to significantly boost programmers' productivity (cur, 2024; Taelin,
 2023a). Despite these advancements, the full extent of their underlying capabilities remains only
 partially understood.

In this paper, we aim to deepen our understanding of transformers' abilities to perform compilation tasks. Empirically, transformer-based LLMs have shown rapid progress in code generation and compilation. For example, MetaLL (Cummins et al., 2024) enables LLMs to optimize code by interpreting compiler intermediate representations (IRs), assembly language, and optimization techniques.
Gu (2023) highlights the ability of LLMs to generate high-quality test cases for Golang compilers. Surprisingly, Taelin (2023b) demonstrates that models like Sonnet-3.5 can compile legacy code into modern languages like TypeScript, outperforming the now obsolete AgdaJS compiler (Agda Development Team, 2024).

To formally study this problem in a controlled setup, we designed a C-like programming language called mini-husky, which encapsulates key features of modern C-like languages such as (Flanagan, 2011) and Rust (Klabnik & Nichols, 2023). We focus on three representative compilation tasks: abstract syntax tree (AST) construction, symbol resolution, and type analysis. The AST is a recursive

structure that represents the input as a tree. From the perspective of programming language design, the AST is considered the true representation of the input, with the textual code serving merely as a convenient interface for human users (Alfred et al., 2007). All syntactic and semantic processing can then be interpreted as specific operations on these trees. Symbol resolution involves verifying the validity of references to entities and flagging errors for undefined symbols. Type analysis encompasses both type inference, which assigns types to variables without explicit annotations, and type checking, which identifies mismatches between actual and expected types.

We demonstrate that under the *clean code principle* (Martin, 2008), transformers can efficiently perform AST construction, symbol resolution, and type analysis, where efficiency means that these tasks can be conducted by transformers with a number of parameters that scale logarithmically with the input code length. To the best of our knowledge, this is the first theoretical demonstration that transformers can function as compilers in a parameter-efficient manner.

We further compare transformers and recurrent neural networks (RNNs). By connecting the type analysis task with the associative recall, we show even under the *clean code principle* (Martin, 2008), RNNs require a memory size that scales *linearly* with the input sequence length to successfully perform type analysis. Consequently, for type analysis in compilation, transformers can be *exponentially more efficient* than RNNs. We also empirically validate our theoretical findings by demonstrating the superiority of transformers in the type analysis task.

Proving that transformers can perform compilation tasks presents several challenges:

Technical Challenges and Our Technique.

- 074
- 075 076

077

- **Transformers operate at too low a level**. Transformers process sequences of floating-point vectors, akin to raw bits in computers, and proving their ability to perform specific tasks is similar to writing specialized parallel machine code. Previous work (Yao et al., 2021) often resorts to graphical illustrations for readability, even for basic tasks.
- Compilers are exceedingly high-level. Compilers are among the most complex programming endeavors of our time. Compilation involves numerous sophisticated procedures, some of which are undecidable or computationally expensive, such as code optimization (Alfred et al., 2007)) and type analysis (Pierce, 2002). For example, type analysis in complex type systems poses significant challenges, often requiring the development of advanced logical frameworks (Dunfield & Krishnaswami, 2019).

085

086 To overcome these challenges, we design a domain-specific language (DSL) called **Cybertron** to serve as the proof vehicle, i.e., a major part of our proof consists of reasoning about type-correct 087 code in Cybertron that represents a transformer. Without using Cybertron, writing an equivalent 088 natural language proof would be too complex and intractable. Using code to prove propositions is 089 not new to computer science; it is, in fact, the norm in interactive theorem proving (ITP) (Har-090 rison et al., 2014). ITP focuses on generating computer-verifiable proofs through a combination 091 of human-guided instructions and software automation. For instance, the correctness of the Ke-092 pler conjecture (Hales et al., 2017) is verified by the combination of the ITP theorem provers HOL 093 Light (Harrison, 2009) and Isabelle (Paulson, 1994). To the best of our knowledge, we are the first 094 to apply this approach to understanding neural networks.

- 095 096 **Contributions.** We summarize our contributions below:
- A testbed for compilation tasks: We introduce Mini-Husky, a simple yet representative C-like programming language, designed to formally assess transformers' capabilities in programming language processing. We anticipate that Mini-Husky will become a standard testbed for this purpose.
- Expressive power theory of transformers for several compilation tasks: We provide a formal proof that, when the input code sequence has bounded AST depth and inference depth, the number of parameters in transformers only needs to scale logarithmically with the input sequence length to handle compilation tasks such as AST construction, symbol resolution, and type analysis. To the best of our knowledge, this is the first study exploring the power of transformers for these compilation tasks.
- **Transformers vs. RNNs**: Theoretically, we demonstrate a negative result, showing that the number of parameters in RNNs must scale linearly with the input sequence length to perform type

analysis correctly. This result establishes an exponential separation between transformers and RNNs. We further empirically confirm the advantage of transformers for the type analysis task.

A Domain-Specific Language for Proofs: Given the challenges in formal proofs, we design a domain-specific language, Cybertron, to serve as a proof vehicle. We believe that Cybertron, and the general approach of using DSLs for analysis, can have broader applications in understanding transformers and other architectures.

- 2 RELATED WORK
- 115 116 117

114

108

Expressive Power of Transformers. A line of work studies the expressive power of attention-based 118 models. One direction focuses on the universal approximation power (Yun et al., 2019; Bhattamishra 119 et al., 2020b;c; Dehghani et al., 2018; Pérez et al., 2021). More recent works present fine-grained 120 characterizations of the expressive power for certain functions in different settings, sometimes with 121 statistical analyses (Edelman et al., 2022; Elhage et al., 2021; Likhosherstov et al., 2021; Akyürek 122 et al., 2022; Zhao et al., 2023; Yao et al., 2021; Anil et al., 2022; Barak et al., 2022; Garg et al., 2022; 123 Von Oswald et al., 2022; Bai et al., 2023; Olsson et al., 2022; Akyürek et al., 2022; Li et al., 2023; Hao et al., 2022; Pérez et al., 2019; Strobl, 2023; Chiang et al., 2023; Wei et al., 2022; Wang et al., 124 2022; Feng et al., 2023; Li et al., 2024; Reddit User, 2013). There are also characterizations of trans-125 formers to be as powerful as universal computers if put in a looped context (Giannou et al., 2023). 126 The most related one is Yao et al. (2021) where the authors prove constructively that bounded depth 127 Dyck language can be recognized by encoder-only hard attention transformers, which has similari-128 ties to our settings of bounded depth programming language recognized encoder-only hard attention 129 transformers. The major difference is that we introduce concepts and tasks from programming lan-130 guage theory Pierce (2002) to study the semantic powers of transformers. 131

Transformers vs. RNN. It is important to understand the comparative advantages and disadvantages 132 of transformers against RNNs. Empirically, synthetic experiments have shown an advantage of 133 transformers against RNNs for long range tasks (Bhattamishra et al., 2023; Arora et al., 2023). 134 Theoretically, there has been a rich line of work focusing on comparing transformers and RNNs in 135 terms of recognizing formal languages (Bhattamishra et al., 2020a; Hahn, 2019; Merrill et al., 2021), 136 which show that the lack of recursive structure of transformers prevent them from recognizing some 137 formal languages that RNNs can recognize. However, the gap can be mitigated when we consider 138 the bounded length of input or bounded grammar depth (Liu et al., 2022; Yao et al., 2021), which 139 is quite reasonable in practice and is used in this paper. On the other side, prior work (Jelassi et al., 140 2024; Wen et al., 2024) proves a representation gap between RNNs and Transformers in repeating a long sequence. In summary, it is somehow intuitive that recursive structures with limited memory 141 perform badly at tasks which requires information retrieval. Our paper shows that semantic analysis 142 for programming languages is such a task. 143

144 **DSLs for Transformers.** We note that we are not exactly the first one to employ a domain-specific 145 language to understand the expressive powers of transformers. Previously, DSLs with simple typings 146 like RASP (Weiss et al., 2021) were proposed to prove constructively that transformers can do various basic sequence-to-sequence operations. Lindner et al. (2023) writes a compiler that compiles 147 RASP into actual transformers, Friedman et al. (2023) shows that RASP can be learned, and Zhou 148 et al. (2023) uses RASP to prove that simple transformers can perform certain algorithms. The major 149 difference between RASP and our DSL Cybertron is that Cybertron has a powerful algebraic type 150 system that helps prove complicated operations beyond simple algorithms. 151

152 153

154

3 PRELIMINARIES

The major innovation in the transformer architecture is that it uses self-attention solely without a conjunction with a recurrent network (Vaswani, 2017), which processes input tokens in a distributed manner. This capability enables the model to handle long-range dependencies, a crucial feature for language tasks. We use hard attention and simplified position encoding to simplify our theoretical reasoning.

- 160
- Attention. In practice, attention heads use soft attention. Given model dimension d_{model} , number of heads H, and a finite set of token positions Pos, an attention layer with simplified position

encoding is defined as a function $f_{\text{attn}} : \mathbb{R}^{\text{Pos} \times d_{\text{model}}} \to \mathbb{R}^{\text{Pos} \times d_{\text{model}}}$ given by

$$\forall p \in \text{Pos}, \quad f_{\text{attn}}(X)_p := W_O \operatorname{Concat} \left(\operatorname{Attn}^{(1)}(X)_p, \dots, \operatorname{Attn}^{(H)}(X)_p \right), \tag{1}$$

where the *h*th attention head is defined using soft attention as: Attn^(h)(X)_p := $\sum_{p' \in Pos} \alpha_{p,p'}^{(h)} V_{p'}^{(h)}$. The attention weights $\alpha_{p,p'}^{(h)}$ given by: $\alpha_{p,p'}^{(h)} = \frac{\exp\left(Q_p^{(h)^{\top}} K_{p'}^{(h)} + \lambda^{(h)^{\top}} \Psi_{p'-p}\right)}{\sum_{p'' \in Pos} \exp\left(Q_p^{(h)^{\top}} K_{p''}^{(h)} + \lambda^{(h)^{\top}} \Psi_{p''-p}\right)}$, where $W_O \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ are trainable parameters, $Q_p^{(h)}, K_p^{(h)}, V_p^{(h)} \in \mathbb{R}^{d_{\text{model}}/H}$ are linear transformations of $X_p, \lambda^{(h)} \in \mathbb{R}^2$ depends on the head, and $\Psi_q = \begin{pmatrix} q \\ 1_{q>0} \end{pmatrix} \in \mathbb{R}^2$ accounts for relative position. For theoretical convenience, we use hard attention, commonly used in theoretical analysis of trans-

For theoretical convenience, we use hard attention, commonly used in theoretical analysis of transformer (Yao et al., 2021; Hahn, 2019). Hard attention can be viewed as the limit of soft attention when the attention logits become infinitely large. The hard attention head is defined as:

$$\operatorname{Attn}^{(h)}(X)_{p} := \frac{1}{|S_{p}|} \sum_{p' \in S_{p}} V_{p'}^{(h)}, \text{ where } S_{p} = \arg \max_{p' \in \operatorname{Pos}} \left(Q_{p}^{(h)^{\top}} K_{p'}^{(h)} + \lambda^{(h)^{\top}} \Psi_{p'-p} \right)$$
(2)

In other words, hard attention selects the positions p' that maximize the attention score for each position p, and averages the corresponding value vectors $V_{p'}^{(h)}$.

Feed-Forward Layer. Given model dimension d_{model} , and a finite set of token positions Pos, a feed-forward layer is a fully connected layer applied independently to each position, defined as a function $f_{\text{fn}} : \mathbb{R}^{\text{Pos} \times d_{\text{model}}} \to \mathbb{R}^{\text{Pos} \times d_{\text{model}}}$ given by

$$\forall p \in \text{Pos}, \quad f_{\text{ffn}}(X)_p = W_2 \sigma_{\text{ReLU}} \left(W_1 X_p + b_1 \right) + b_2, \tag{3}$$

where $W_1 \in \mathbb{R}^{d_{\text{fin}} \times d_{\text{model}}}$ and $W_2 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ffn}}}$ are trainable weight matrices, $b_1 \in \mathbb{R}^{d_{\text{fin}}}$ and $b_2 \in \mathbb{R}^{d_{\text{model}}}$ are trainable bias vectors, d_{ffn} is the hidden dimension of the feed-forward layer, chosen to be $2d_{\text{model}}$, as commonly used in practice, σ_{ReLU} is the ReLU activation function.

Encoder-Only Transformer. Encoder-only transformers consist solely of the encoder stack, making them ideal for tasks like classification, regression, and sequence labeling that do not require sequence generation. Each encoder layer includes a multi-head self-attention mechanism and a feed-forward network, allowing the model to capture complex dependencies and contextual information.

One can define it using the following recurrence,

- The input is given by: $X^{(0)} = X$.
- For each layer $l = 1, 2, \ldots, L$:

- Compute attention output: $\hat{X}^{(l)} = X^{(l-1)} + f_{\text{attn}}^{(l)} (X^{(l-1)})$,

- Compute feed-forward output:
$$X^{(l)} = \hat{X}^{(l)} + f_{\text{ffn}}^{(l)} \left(\hat{X}^{(l)} \right)$$

In the above, $f_{\text{attn}}^{(l)}$ are the attention layers, and $f_{\text{ffn}}^{(l)}$ are the feed-forward layers, with the same model dimension d_{model} , number of heads H, and set of token positions Pos. For simplicity, layer normalization is ignored. See Appendix C for full details of transformers and other architectures.

4 PROGRAMMING LANGUAGE PROCESSING AND THE TARGET C-LIKE LANGUAGE: MINI-HUSKY

Recently, transformers have expanded to support code analysis and generation (Nijkamp et al., 2023;
 Chen et al., 2021; Anysphere, 2023). Programming languages offer a cleaner foundation for study ing language understanding, as their syntactic and semantic tasks are precisely defined. To formally
 study the language processing capabilities of transformers, we design Mini-Husky, a representa tive mix of modern C-like languages with strong typing and typical syntactic features. It supports



263 264

266

14 265 15 fn g() { f() }

9

10

11

12

13 }

{ let a = 4; } let y = a;

let z = a;

The outer function f is accessible everywhere in the body of function g. However, the inner function f₁ can only be used inside the body of f as it is defined within the body. For variables with the same identifier **a**, the first is accessible from line 5, the second is accessible from line 12, the third is accessible from line 10, and the fourth is not accessible from anywhere. Thus x = 1, y = 3, z = 2. 270 The output of the **symbol resolution** task is the collection of symbol resolution results on all appli-271 cable tokens. More concretely, the output is a sequence of values of type Option<SymbolResolution> 272 where Option<SymbolResolution> is the type SymbolResolution with a null value added for non-273 applicability and SymbolResolution is the type storing the result of the symbol resolution, be-274 275 ing either a success with a resolved symbol of type Symbol or a failure with an error of type 276 SymbolResolutionError. We shall prove that transformers can do symbol resolution and that atten-277 tion is crucial. 278

Type Analysis. In general, types are essential for conveying the intended usage of the written functions and specifying constraints. As a first exploration of this topic, we try to make the type analysis in Mini-Husky as simple as possible yet able to bring out the essential difficulty. The type system consists of four sequential components: (1) *Type definition*, (2) *Type specification*, (3) *Type inference*, and (4) *Type checking*. Due to the page limit, here we only introduce (4) *Type checking* because it is the final step and this is a crucial step which separates transformers and RNNs. See Appendix E.1 for details of (1) *Type definition*, (2) *Type specification*, and (3) *Type inference*.

Type checking ensures that the typed expressions agree with its expectations. For simplicity, we do not allow implicit type conversion, so the agreement means exact equality of types. The arguments of function calls are expected to have types according to the definition of the function. The operand type of field access must be a struct type with a field of the same name. The type of the last expression of the function body or the expr in the return statement must be equal to the return type of the function. For variables defined in the let statement, If the types are annotated, the types of the left-hand side and right-hand side should be in agreement.

```
1 // Type Error: the return type is 'i32', yet the last expression is of type 'f32'
   fn f(a: i32) -> i32 { 1.1 }
2
3
 4 struct A { x: i32 }
 5
 6
    fn g() {
         // Type Error: `x` is of type f32 but it's assigned by a value of type `i32`
// Type Error: the first argument of `f` is expected to be of type `i32` but gets a
 7
 8
          float literal instead
 9
         let x: f32 = f(1.1);
10
         // Type Error: no field named 'y'
11
         let y = A { x: 1 }.y;
12 }
```

304

309 310

311

293

294

295

296

297

298

299

The above incorporates typical examples of type disagreements that count as type errors. A compiler should be able to report these errors.

The **type analysis** task's final output is the collection of all type errors. More concretely, the output is a sequence of Option<TypeError>, where Option<TypeError> denoted the type TypeError will a null value added and TypeError is the type storing the information of a type error. The position of type errors agrees with the source tokens leading to these errors.

5 EXPRESSIVE POWER OF TRANSFORMERS AS EFFICIENT COMPILERS

In this section we discuss main theoretical results about the expressive power of transformers to perform compilation tasks: AST construction, symbol resolution, and type analysis. In Section 5.4, we discuss Cybertron, a DSL specifically designed for our proof.

3165.1ABSTRACT SYNTAX TREE CONSTRUCTION317

318 We start with a definition that characterizes low-complexity code.

Definition 1 (code with Bounded AST-Depth). Let MiniHusky_D be the set of token sequences that
 can be parsed into valid ASTs in Mini-Husky with a depth less than D.

D in the above definition is small in practice, and a linear dependency on D is acceptable, but the linear dependency on the length of the token sequence L is not. The fundamental reason is that the *clean code principle* (Martin, 2008) requires one to write code with as little nested layer as possible for greater readility. Readability is of the utmost importance because "Programs are meant to be read by humans and only incidentally for computers to execute" (Abelson et al., 1996). This assumption of bounded hierarchical depth is not limited to just programming languages, but is often seen as applicable to natural languages (Frank et al., 2012; Brennan & Hale, 2019; Ding et al., 2017), motivating Yao et al. (2021) to have a similar boundedness assumption. Below is the main result for AST construction using transformers.

Theorem 1. There exists a transformer encoder of model dimension and number of layers being O(log L + D) and number of heads being O(1) that represents a function that maps any token sequence of length L in MiniHusky_D to its abstract syntax tree represented as a sequence.

We note $\log L$ is small because 64-bit computers can only process context length at most 2^{64} and *D* is small by assumption. Therefore, there exists a transformer with an almost constant number of parameters that is able to process comparatively much longer context length.

337

333

Proof Sketch. The idea is to construct ASTs in a bottom-up manner with full parallelism. We shall 338 recursively produce the final ASTs in at most D steps. We shall maintain two values, called pre_asts 339 and asts . asts represents ASTs that have already been allocated, although they might not have been 340 341 fully initialized. pre asts represents tokens that have yet to form ASTs and new ASTs that have not 342 been fully initialized. For each round, we try to create new ASTs from pre_asts and update asts 343 and pre asts. For the n-th round, we provably allocated all ASTs with a depth no more than n. Then 344 for the D-th round, all ASTs are properly constructed and allocated. Each round can be represented 345 by a transformer of O(1) number of heads, model dimension $O(\log L + D)$, and O(1) number of 346 layers. Therefore, the end-to-end process is then representable by a transformer of O(1) number 347 of heads, model dimension $O(\log L + D)$, and $O(\log L + D)$ number of layers. See full details in 348 Appendix **F**. \square

349 350

351

5.2 SYMBOL RESOLUTION

Next, we show that transformers can effectively perform symbolic resolution as $\log L$ and D are almost constant as compared with context length L.

Theorem 2. There exists a transformer encoder of model dimension and number of layers being O(log L + D) and number of heads being O(1) that represents a function that maps any token sequence of length L in MiniHusky_D to its symbol resolution represented as a sequence of values of type Option<SymbolResolution>.

358

359 *Proof Sketch.* First, we need to define the type for scopes. It is represented by a tiny sequence of 360 indices of curly brace block AST that enclose the type/function/variable. We assign the scope by 361 walking through the ASTs in a top-down manner. We not only assign scopes to item definitions, 362 we also: (1) assign scopes to ASTs representing curly brace blocks, with these scopes equal to the scope of block itself, and (2) assign scopes to identifiers waiting to be resolved, with these scopes 363 equal to the maximum possible scope of its resolved definition. The computation process is easily 364 represented in Cybertron, indicating attention is expressive enough for this calculation and it only 365 takes O(D) number of layers. 366

After obtaining all the scopes for all items, it takes only one additional layer to obtain the symbolic resolution through attention. As attention is expressed through the dot product of two linear projections Q and K, we have to choose the representation of the scope type properly to finish the proof. The full details are in Appendix G.

- 371
- **372 5.3** Type Analysis

We need an additional definition to characterize the complexity of code for type analysis.

375 Definition 2 (code with Bounded AST-Depth and Type-Inference-Depth). We use **376** MiniHuskyAnnotated_{D,H} to denote the subset of MiniHusky_D with the depth of type in- **377** ference no more than H. The depth of type inference is the number of rounds of computation needed to infer all the types using the type-inference algorithm (described in Appendix E.1). In practice, H is significantly smaller than the context length L for reasonably written code because it is upper bounded by the number of statements in a function body which is required to be small according to the *clean code principle* (Martin, 2008). Below, we present the main result of using transformers for type analysis. See full details in Appendix H.

Theorem 3. For $L, D, H \in \mathbb{N}$, there exists a transformer encoder of model dimension, and number of layers being $O(\log L + D + H)$ and number of heads being O(1) that represents a function that maps any token sequence of length L in MiniHuskyAnnotated_{D,H} to its type errors represented as a sequence of values of type Option<TypeError>.

386 387 388

5.4 PROOF VEHICLE: CYBERTRON, A DOMAIN-SPECIFIC LANGUAGE

389 Here we highlight our main proof technique. Proving that transformers can express complex algo-390 rithms and software like compilers is a significant challenge due to the inherent differences between 391 how transformers operate and the nature of high-level tasks they are expected to perform. Trans-392 formers process input at a low level, where each layer manipulates raw token sequences as vectors 393 without predefined structure or meaning. However, high-level tasks—such as constructing ASTs 394 and performing type and symbol analysis—require handling complex, structured information that 395 depends on long-range relationships and interactions across the input. Bridging the gap between 396 this raw, unstructured processing and the structured, multi-step logic required for these tasks introduces significant difficulty. Compilers, for instance, typically rely on rule-based, step-by-step 397 operations that are abstract and sequential, which transformers must simulate through their attention 398 mechanisms and feedforward layers. The challenge is further compounded by the need to formally 399 prove that transformers can handle such tasks efficiently and accurately, despite operating in a fun-400 damentally different manner. To address these challenges, we propose a domain-specific language 401 (DSL) called **Cybertron**, which allows us to systematically prove that transformers are capable of 402 expressing complex algorithms while maintaining sufficient readability. 403

A key feature of **Cybertron** is its expressive type system, which provides strong correctness guarantees. The type system ensures that every value is strongly typed, making it easier to reason about function composition and ensuring the validity of our proofs. This type system is crucial for managing how transformers represent and manipulate both local and global types—where local types correspond to individual tokens and global types refer to sequences of tokens, encapsulating broader program information.

What transformers output (possibly in the intermediate layers) is a representation in sequences of vector of sequences of values in these types. As types are mathematically interpreted in this paper as a discrete subset of a vector space, Cybertron allows us to construct transformers with automatic value validity guarantees if the Cybertron code is type-correct.

In Cybertron, complex functions are broken down into "atomic" operations through propositions
 on function compositions and computation graphs (Propositions 11,13,14,2). It is straightforward to
 prove that these "atomic" operations are representable by transformers, either by feedforward layers
 or attention layers. For example:

- Feedforward layers: boolean operations like AND (Proposition 6), OR (Proposition 7), or NOT (Proposition 5), or operations over option types like Option::or (Proposition 9) being applied to each token in a sequence.
- Attention layers: operations that require information transmission between tokens such as nearest_left and nearest_right that collect for each token the nearest left/right non-nil information (Proposition 15).

This approach allows us to break down complex operations into primitive tasks that transformers can simulate. Feedforward layers handle local operations on individual tokens, while attention layers manage long-range dependencies and interactions between tokens, simulating the multi-step reasoning required for higher-level tasks.

Cybertron's expressive type system and function composition framework help bridge the gap be tween the low-level processing transformers perform and the high-level reasoning necessary for
 complex tasks like compilation. For full details, including the mathematical foundations of Cy bertron's type system and function composition, see Appendix D.

432 6 COMPARISONS BETWEEN TRANSFORMERS AND RNN 433

434 435

Now we compare transformers and RNNs from both theoretical and empirical perspectives.

436 437

446

447

450

451

468

6.1 A LOWER BOUND FOR RNNS FOR TYPE CHECKING

438 Previously, it has shown that RNN is provably less parameter efficient than transformers for associa-439 tive recall (Wen et al., 2024). Intuitively speaking, the type checking step covers associative recall. 440 Based on this observation, we obtain the following lower bound for RNNs. 441

Theorem 4. For $L, D, H \in \mathbb{N}$, for any RNN that represents a function that maps any token sequence 442 of length L in MiniHuskyAnnotated_{D H} with D, H = O(1) to its type errors represented as a 443 sequence of values of type Option <TypeError>, then its state space size is at least $\Omega(L)$. 444

445 Theorem 3 and Theorem 4 give a clear separation between transformers and RNNs in terms of the compilation capability. Specifically, if the input codes satisfy $D, H \ll L$, which is typically the case under the *clean code principle* (Martin, 2008), then transformers at most need $O((\log L + D + H))$ 448 number of parameters, which is significantly smaller what RNNs requires, $\Omega(L)$. 449

6.2 EMPIRICAL COMPARISON BETWEEN TRANSFORMERS AND RNNS

452 We validate our theoretical results by conducting experiments on synthetic data. 453

Dataset construction. The synthetic dataset is parameterized by n (the number of data pieces), f 454 (the number of functions in a data piece), a (the maximum number of arguments of any function), 455 c (the maximum number of function calls involved in any function), d (the minimum distance be-456 tween the declaration and the first call of a function, as well as the minimum distance between its 457 consecutive calls), v (the probability of using a variable in a function call), and e (the error rate of 458 using an incorrect type in a function call). 459

The names of the functions are drawn randomly and uniquely from a list of English words. For 460 each of the arguments of any function, its symbol is randomly drawn from another list of English 461 words and its type is randomly drawn from {Int, Float, Bool}. All the called functions must be 462 declared and not called by at least d functions ahead of the current one. For each argument of any 463 *function call*, with probability v, the argument variable of the enclosing function is used regardless 464 of its type, with probability (1-v)(1-e), a literal of the correct type is used, and with probability 465 (1-v)e an incorrect type literal is used. For integers, the literals are from $\{0, 1, \ldots, 99\}$; for floats, 466 the literals are from $\{0.1, 1.1, \dots, 99.1\}$; for booleans, the literals are from $\{true, false\}$. The 467 training dataset and evaluation dataset use *disjoint* lists for function names and argument symbols.

Below is a data piece with f = 10, a = 5, c = 5, d = 3, v = 0.2, e = 0.5: 469

```
1 fn rename_file ( i : Float , sum : Float ) { }
470
             fn parse_data ( list : Int , value : Bool , stack : Float , k : Float , msg : Float ) { }
          2
471
            fn parse_json ( position : Bool ) { }
          3
472
            fn find_by_id ( error : Float ) { rename_file ( 60.1 , 94.1 ) ; }
          4
          5 fn merge ( group : Int , table : Float , error : Bool , count : Int ) { parse_data ( 7 ,
      false , 49.1 , 33.1 , 4.1 ) ; }
473
474
          6 fn log_info ( val : Bool , m : Bool , xml : Float , path : Float ) { parse_json ( true ) ;
475
          7 fn process ( function : Int , value : Float , keys : Bool ) { find_by_id ( 88.1 ) ;
476
                  rename_file ( value , 40.1 ) ; }
          8 fn validate_response ( end : Int , z : Float , max : Bool ) { merge ( 1 , true , 27.1 , 72
        ) ; parse_data ( 11 , 85 , 35.1 , 14.1 , true ) ; }
477
478
          9 fn print_message ( algorithm : Float ) { parse_json ( 92 ) ; log_info ( true , algorithm ,
                  false , 26.1 ) ; }
479
         10 fn print_help ( max : Bool , tree : Int , method : Int , item : Bool ) { process ( 25 , 28
480
                  , false ) ; rename_file ( 48 , 80.1 ) ; }
481
         Model and training. We use customized BERT models (Devlin et al., 2019) and bidirectional RNN
482
```

483 models (Schuster & Paliwal, 1997) in our experiments. To control the model size (i.e., the number of trainable parameters), we adjust only the hidden sizes while keeping other hyperparameters constant. 484 Detailed model specifications can be found in Table 1. For both transformers and RNNs, we use the 485 hyperparameters listed in Table 2 in Appendix J during the training process.



Figure 2: Figures depicting the accuracy of the expected type (see Section 5.3) across different models, measured by their number of trainable parameters, when trained on various datasets. Training accuracies are better indicators of the expressive power of the models (instead of generalizability) than evaluation accuracies. We also report evaluation accuracies in Appendix J.

Results. We experimented with multiple combinations of models (Table 1) and datasets (Table 2). For each combination, we conducted independent runs using a fixed set of k = 5 random seeds. When plotting the figures, we took the top t = 5 training/evaluation losses/accuracies from each run and averaged over all the $k \times t$ values. We plotted separate figures for each dataset and separate subfigures for each metric. In each sub-figure, the x-axis represents the number of trainable parameters, and the y-axis represents the averaged values. Results are shown in Figure 2. They demonstrate that customized BERT models are able to perform better at type checking than bidirectional RNN models when both scale up, corroborating our theories. Other results are in Appendix J.

7 CONCLUSION

We demonstrated that transformers can efficiently handle a number of syntactic and semantic analy-512 sis tasks in C-like languages, using Cybertron to prove their capacity for tasks like AST generation, 513 symbol resolution, and type analysis. We show a theoretical advantage of transformers over RNNs, 514 particularly in their ability to manage long-range dependencies with logarithmic parameter scaling. 515 In a sense, transformers have the right inductive bias for language tasks. Our experiments confirmed 516 these theoretical insights, showing strong performance on synthetic and real datasets, underscoring 517 the expressiveness and efficiency of transformers in sequence-based learning.

518 519 520

521

522

523

525 526

527 528

529 530

531

532

534

535

536

495

496

497

498

499 500

501

502

504

505

506

507

508 509

510 511

8 ACKNOWLEDGEMENT

Xiyu Zhai acknowledges the support of NSF through awards DMS-2031883 and PHY-2019786. Liao Zhang acknowledges the ERC PoC project FormalWeb3 no. 101156734 and the University of Innsbruck doctoral scholarship promotion of young talent. 524

REFERENCES

- Cursor: Ai-powered code editor, 2024. URL https://www.cursor.com/. Accessed: September 29, 2024.
- Harold Abelson, Gerald Jay Sussman, and with Julie Sussman. Structure and Interpretation of Computer Programs. MIT Press/McGraw-Hill, Cambridge, 2nd editon edition, 1996. ISBN 0-262-01153-0.
- Agda Development Team. Agda compilers manual v2.6.4.2, 2024. URL https: //agda.readthedocs.io/en/v2.6.4.2/tools/compilers.html# javascript-backend.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-538 rithm is in-context learning? investigations with linear models. arXiv preprint arXiv:2211.15661, 2022.

553

554

555 556

558 559

560

561

562 563

564

565 566

567

568

569

578

579

580

581

585

592

- V Aho Alfred, S Lam Monica, and D Ullman Jeffrey. *Compilers principles, techniques & tools.* pearson Education, 2007.
- Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *arXiv preprint arXiv:2207.04901*, 2022.
- 546 547 Anysphere. Cursor, 2023. URL https://www.cursor.com/features.
- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher R'e. Zoology: Measuring and improving recall in efficient language models. ArXiv, abs/2312.04927, 2023. URL https://api.semanticscholar.org/ CorpusID:266149332.
 - Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
 - Boaz Barak, Benjamin Edelman, Surbhi Goel, Sham Kakade, Eran Malach, and Cyril Zhang. Hidden progress in deep learning: Sgd learns parities near the computational limit. *Advances in Neural Information Processing Systems*, 35:21750–21764, 2022.
 - S. Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of transformers to recognize formal languages. In *Conference on Empirical Methods in Natural Language Processing*, 2020a. URL https://api.semanticscholar.org/CorpusID:222225236.
 - S. Bhattamishra, Arkil Patel, Phil Blunsom, and Varun Kanade. Understanding in-context learning in transformers and llms by learning to learn discrete functions. *ArXiv*, abs/2310.03016, 2023. URL https://api.semanticscholar.org/CorpusID:263620583.
 - S. Bhattamishra, Michael Hahn, Phil Blunsom, and Varun Kanade. Separations in the representational capabilities of transformers and recurrent architectures. *ArXiv*, abs/2406.09347, 2024. URL https://api.semanticscholar.org/CorpusID:270440803.
- Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of transformers
 to recognize formal languages. *arXiv preprint arXiv:2009.11264*, 2020b.
- Satwik Bhattamishra, Arkil Patel, and Navin Goyal. On the computational power of transformers and its implications in sequence modeling. *arXiv preprint arXiv:2006.09286*, 2020c.
- Jonathan Brennan and John Tracy Hale. Hierarchical structure guides rapid linguistic predictions during naturalistic listening. *PLoS ONE*, 14, 2019. URL https://api.
 semanticscholar.org/CorpusID:260538292.
 - Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- David Chiang, Peter A. Cholak, and Anand Pillay. Tighter bounds on the expressivity of transformer
 encoders. In *International Conference on Machine Learning*, 2023. URL https://api.
 semanticscholar.org/CorpusID:256231094.
- Chris Cummins, Volker Seeker, Dejan Grubisic, Baptiste Rozière, Jonas Gehring, Gabriele Synnaeve, and Hugh Leather. Meta large language model compiler: Foundation models of compiler optimization. *ArXiv*, abs/2407.02524, 2024. URL https://api.semanticscholar.org/CorpusID:270924331.
- Valentin David. Language Constructs for C++-like languages. PhD thesis, University of Bergen, 2009.
- 593 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.

594 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep 595 bidirectional transformers for language understanding, 2019. URL https://arxiv.org/ 596 abs/1810.04805. 597 Nai Ding, Lucia Melloni, Xing Tian, and David Poeppel. Rule-based and word-level statistics-based 598 processing of language: insights from neuroscience. Language, Cognition and Neuroscience, 32: 570-575,2017. URL https://api.semanticscholar.org/CorpusID:46747073. 600 601 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 602 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An 603 image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint 604 arXiv:2010.11929, 2020. 605 Jana Dunfield and Neelakantan R Krishnaswami. Sound and complete bidirectional typechecking 606 for higher-rank polymorphism with existentials and indexed types. Proceedings of the ACM on 607 Programming Languages, 3(POPL):1–28, 2019. 608 609 Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable 610 creation in self-attention mechanisms. In International Conference on Machine Learning, pp. 611 5793-5831. PMLR, 2022. 612 N Elhage, N Nanda, C Olsson, T Henighan, N Joseph, B Mann, A Askell, Y Bai, A Chen, T Conerly, 613 et al. A mathematical framework for transformer circuits. Transformer Circuits Thread, 2021. 614 615 The curry-howard correspondence. 2021. URL https://api. Husna Farooqui. 616 semanticscholar.org/CorpusID:244268761. 617 Guhao Feng, Yuntian Gu, Bohang Zhang, Haotian Ye, Di He, and Liwei Wang. Towards revealing 618 the mystery behind chain of thought: a theoretical perspective. ArXiv, abs/2305.15408, 2023. 619 URL https://api.semanticscholar.org/CorpusID:258865989. 620 621 David Flanagan. JavaScript: The definitive guide: Activate your web pages. " O'Reilly Media, 622 Inc.", 2011. 623 624 Bryan Ford. Parsing expression grammars: a recognition-based syntactic foundation. In Proceedings of the 31st ACM SIGPLAN-SIGACT symposium on Principles of programming languages, pp. 625 111-122, 2004. 626 627 S. Frank, Rens Bod, and Morten H. Christiansen. How hierarchical is language use? Proceedings 628 of the Royal Society B: Biological Sciences, 279:4522 - 4531, 2012. URL https://api. 629 semanticscholar.org/CorpusID:11969171. 630 Dan Friedman, Alexander Wettig, and Danqi Chen. Learning transformer programs. ArXiv, 631 URL https://api.semanticscholar.org/CorpusID: abs/2306.01128, 2023. 632 259064324. 633 634 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn 635 in-context? a case study of simple function classes. Advances in Neural Information Processing 636 Systems, 35:30583-30598, 2022. 637 638 Angeliki Giannou, Shashank Rajput, Jy yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. ArXiv, abs/2301.13196, 2023. 639 URL https://api.semanticscholar.org/CorpusID:256389656. 640 641 Ai achieves silver-medal standard solving international mathematical Google Deepmind. 642 olympiad problems, July 2024. URL https://deepmind.google/discover/blog/ 643 ai-solves-imo-problems-at-silver-medal-level/. 644 645 Qiuhan Gu. Llm-based code generation method for golang compiler testing. Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of 646 Software Engineering, 2023. URL https://api.semanticscholar.org/CorpusID: 647 265509921.

648 649 650	Michael Hahn. Theoretical limitations of self-attention in neural sequence models. <i>Transactions of the Association for Computational Linguistics</i> , 8:156–171, 2019. URL https://api.semanticscholar.org/CorpusID:189928186.			
652 653 654 655	Thomas Hales, Mark Adams, Gertrud Bauer, Tat Dat Dang, John Harrison, Hoang Le Truong, Cezary Kaliszyk, Victor Magron, Sean McLaughlin, Tat Thang Nguyen, et al. A formal proof of the kepler conjecture. In <i>Forum of mathematics, Pi</i> , volume 5, pp. e2. Cambridge University Press, 2017.			
656 657 658 659	Sophie Hao, Dana Angluin, and Roberta Frank. Formal language recognition by hard attention transformers: Perspectives from circuit complexity. <i>Transactions of the Association for Computational Linguistics</i> , 10:800–810, 2022. URL https://api.semanticscholar.org/CorpusID:248177889.			
660 661 662	John Harrison. Hol light: An overview. In <i>International Conference on Theorem Proving in Higher</i> <i>Order Logics</i> , pp. 60–66. Springer, 2009.			
663 664	John Harrison, Josef Urban, and Freek Wiedijk. History of interactive theorem proving. In <i>Handbook of the History of Logic</i> , volume 9, pp. 135–214. Elsevier, 2014.			
665 666 667 668	Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and Eran Malach. Repeat after me: Transformers are better than state space models at copying. <i>ArXiv</i> , abs/2402.01032, 2024. URL https://api.semanticscholar.org/CorpusID:267406617.			
669 670 671 672	Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Aleš Bizjak, Lars Birkedal, and Derek Dreyer. Iris from the ground up: A modular foundation for higher-order concurrent separation logic. <i>Journal of Functional Programming</i> , 28, 2018. URL https://api.semanticscholar. org/CorpusID:2023423.			
673 674	Steve Klabnik and Carol Nichols. The Rust programming language. No Starch Press, 2023.			
675 676	Shuai Li, Zhao Song, Yu Xia, Tong Yu, and Tianyi Zhou. The closeness of in-context learning and weight shifting for softmax regression. <i>arXiv preprint arXiv:2304.13276</i> , 2023.			
677 678 679 680	Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=3EWTEy9MTM.			
681 682	Valerii Likhosherstov, Krzysztof Choromanski, and Adrian Weller. On the expressive power of self-attention matrices. <i>arXiv preprint arXiv:2106.03764</i> , 2021.			
683 684 685 686	David Lindner, J'anos Kram'ar, Matthew Rahtz, Tom McGrath, and Vladimir Mikulik. Tracr: Compiled transformers as a laboratory for interpretability. <i>ArXiv</i> , abs/2301.05062, 2023. URL https://api.semanticscholar.org/CorpusID:255749093.			
687 688 689	Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata. <i>ArXiv</i> , abs/2210.10749, 2022. URL https://api.semanticscholar.org/CorpusID:252992725.			
690 691 692	Robert C. Martin. <i>Clean Code: A Handbook of Agile Software Craftsmanship</i> . Prentice Hall PTR, USA, 1 edition, 2008. ISBN 0132350882.			
693 694 695	Patrick Massot. Teaching mathematics using lean and controlled natural language. In <i>International Conference on Interactive Theorem Proving</i> , 2024. URL https://api.semanticscholar.org/CorpusID:272330159.			
696 697 698	The mathlib Community. The lean mathematical library. <i>Proceedings of the 9th ACM SIGPLAN</i> <i>International Conference on Certified Programs and Proofs</i> , 2019. URL https://api.semanticscholar.org/CorpusID:204801213.			
700 701	William Merrill, Ashish Sabharwal, and Noah A. Smith. Saturated transformers are constant-depth threshold circuits. <i>Transactions of the Association for Computational Linguistics</i> , 10:843–856, 2021. URL https://api.semanticscholar.org/CorpusID:248085924.			

- Erik Nijkamp, Hiroaki Hayashi, Caiming Xiong, Silvio Savarese, and Yingbo Zhou. Codegen2: Lessons for training llms on programming and natural languages. *ICLR*, 2023.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction
 heads. *arXiv preprint arXiv:2209.11895*, 2022.
- 708 OpenAI. Openai ol system card, September 2024a. URL https://openai.com/index/ openai-ol-system-card/.
- 711 OpenAI. Sora: Creating video from text, February 2024b. URL https://openai.com/ 712 index/sora/.
- ⁷¹³ Lawrence C Paulson. *Isabelle: A generic theorem prover*. Springer, 1994.

743

- Jorge Pérez, Javier Marinkovic, and Pablo Barceló. On the turing completeness of modern neural network architectures. ArXiv, abs/1901.03429, 2019. URL https://api. semanticscholar.org/CorpusID:57825721.
- Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing complete. *The Journal of Machine Learning Research*, 22(1):3463–3497, 2021.
- Benjamin C Pierce. *Types and programming languages*. MIT press, 2002.
- The Univalent Foundations Program. Homotopy type theory: Univalent foundations of mathematics.
 arXiv preprint arXiv:1308.0729, 2013.
- 724 725 Reddit User. I think the main secret sauce of ol is the data. https://www.reddit.com/ r/singularity/comments/lfi6yy9/i_think_the_main_secret_sauce_of_ ol_is_the_data/, 2013. Accessed: 2024-09-28.
- Mike Schuster and Kuldip K Paliwal. Bidirectional recurrent neural networks. *IEEE transactions* on Signal Processing, 45(11):2673–2681, 1997.
- Lena Strobl. Average-hard attention transformers are constant-depth uniform threshold circuits. ArXiv, abs/2308.03212, 2023. URL https://api.semanticscholar.org/CorpusID:260680416.
- Victor Taelin. Ai and the future of coding. https://medium.com/jonathans-musings/
 ai-and-the-future-of-coding-43caad31c3d3, 2023a. Accessed: 2024-10-01.
- Victor Taelin. Agda to typescript compilation with sonnet-3.5, 2023b. URL https://x.com/
 VictorTaelin/status/1837925011187027994. Accessed: September 29, 2024.
- 739 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordv intsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient
 descent. arXiv preprint arXiv:2212.07677, 2022.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. Self consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171,
 2022. URL https://api.semanticscholar.org/CorpusID:247595263.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc
 Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language mod els. ArXiv, abs/2201.11903, 2022. URL https://api.semanticscholar.org/
 CorpusID:246411621.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. ArXiv, abs/2106.06981, 2021. URL https://api.semanticscholar.org/CorpusID:235421630.
- Kaiyue Wen, Xingyu Dang, and Kaifeng Lyu. Rnns are not transformers (yet): The key bottleneck on in-context retrieval. ArXiv, abs/2402.18510, 2024. URL https://api.semanticscholar.org/CorpusID:268041425.

756 757 758	Shangda Wu, Xu Tan, Zili Wang, Rui Wang, Xiaobing Li, and Maosong Sun. Beyond language models: Byte models are digital world simulators. <i>ArXiv</i> , abs/2402.19155, 2024. URL https://api.semanticscholar.org/CorpusID:268063492.
759 760 761 762 763	Shunyu Yao, Binghui Peng, Christos H. Papadimitriou, and Karthik Narasimhan. Self-attention networks can process bounded hierarchical languages. In <i>Annual Meeting of the Association for Computational Linguistics</i> , 2021. URL https://api.semanticscholar.org/CorpusID:235166395.
764 765 766	Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? <i>arXiv preprint arXiv:1912.10077</i> , 2019.
767 768 769	Haoyu Zhao, Abhishek Panigrahi, Rong Ge, and Sanjeev Arora. Do transformers parse while pre- dicting the masked word? <i>arXiv preprint arXiv:2303.08117</i> , 2023.
770 771 772 773 774	Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Josh Susskind, Samy Ben- gio, and Preetum Nakkiran. What algorithms can transformers learn? a study in length gener- alization. <i>ArXiv</i> , abs/2310.16028, 2023. URL https://api.semanticscholar.org/ CorpusID:264439160.
775 776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
004	
200	
807	
808	
800	
009	

810	А	TREE			
811					
812	Tree	s are one of the most fundamental objects to study in computer science. However, its exact			
813	defin	definition differs for different domains. The trees used in "abstract syntax tree" (Section B) is more			
814	restrictive than that in mathematics, which we call "typed tree", so that one can define recursive				
815	com	putation more rigorously.			
017					
017 818	A.I	WHAT ARE TREES			
819	Tree	s in data structures have slightly additional meaning as compared to trees in mathematics. In			
820	this	paper, all trees are trees in data structures. For clarity, we lay down the precise definition of			
821	trees	in data structure.			
822	Defi	nition 3 (Tree). A tree T is a set of nodes storing elements such that the nodes have a parent-			
823	child	relationship that satisfies the following:			
824					
825		• If T is not empty, it has a special node called the root that has no parent.			
826		• Each node v of T other than the root has a unique parent node w ; each node with parent			
827		w is a child of w.			
828					
829	We d	enote the nodes of T as $N(T)$.			
830	Defi	nition 4 (Recursive Definition of a Tree). A tree T is either empty or consists of a node r (the			
831	root)	and a possibly empty set of trees whose roots are the children of r.			
032					
834	How	ever, the above definition is too permissive. We shall define a typed version as follows:			
835	Defi	nition 5 (Typed Tree). A tree type consists of a set of values V and a set of relationships			
836	$C \subseteq$	$V \times \mathbb{N}$, and a typed tree under this type is any tree T such that for each node, a value $v \in V$			
837	is us.	signed such that $(v, n) \in \mathbb{C}$ where n is the number of the children of the hode.			
838	All t	rees in this paper are typed.			
839	Exai	nple 1 (Abstract syntax tree (AST) as Typed Tree). <i>Consider an AST for a simple arithmetic</i>			
840	expre	ession. Let the set of values V be:			
841					
842		$V = \{ num , add , sub , mul , div \}$			
843	and i	the set of relationships $C \subseteq V imes \mathbb{N}$ specify the allowed number of children for each value:			
044 9/5					
846		$C = \{(num, 0), (add, 2), (sub, 2), (mul, 2), (div, 2)\}$			
847	Ana	xample AST for the arithmetic expression $(3 \pm 5) \times 2$ is the following typed tree:			
848	лп е.	sample AST for the arithmetic expression $(3 \pm 5) \times 2$ is the following typed tree.			
849		• The root node is labeled mul (multiplication) and it has two children			
850					
851		- The left child is labeled add (addition), and it has two children:			
852		* The left child of add is labeled num with the value 3.			
853		* The right child of add is labeled num with the value 5.			
854		The right shild of multistabled mum with the value?			
855		- The right child of mut is tubeled hum with the value 2.			
856 857	This	tree conforms to the typing rules because:			
858		• num has 0 children,			
859					
861		• aaa nas 2 children,			
862		• mul has 2 children.			
863		······,			

all of which satisfy the relationships in C.

864 A.2 REPRESENTATIONS OF TREES

It's also important to talk about tree representations. We are studying transformers, and then it's necessary to represent large trees as a sequence, otherwise the model dimension is not large enough to contain the information locally. Let's first talk about the classical **arena pattern** used in system programming for representing trees and we shall slightly adapt it to our use case for studying transformers.

Arena Pattern. To represent trees efficiently in memory, especially when trees are frequently modified (such as insertions or deletions of nodes), an arena pattern is often used. The arena pattern provides a way to manage memory allocation for tree structures, allowing for efficient memory usage and avoiding fragmentation. Here's how the arena pattern works in the context of tree representation:

Definition 6 (Arena Pattern in Tree Representation). *In the arena pattern, a tree is represented by an array (or vector) of nodes, called an arena. Each node in the arena contains:*

- An element or value stored in the node.
- References (often indices or pointers) to the node's children and possibly to its parent.
- The key characteristics of the arena pattern are:
 - Memory Contiguity: All nodes are stored contiguously in memory within the arena, which allows for efficient traversal and modification operations.
 - Fixed Capacity: The arena has a fixed or dynamically resizable capacity, and nodes are added sequentially. This avoids the overhead of allocating individual nodes on the heap.
 - Index-based References: Instead of using pointers, the nodes reference each other using indices within the array, which simplifies memory management and can lead to cache-friendly operations.
 - Efficient Allocation and Deallocation: Nodes can be efficiently allocated and deallocated within the arena without requiring complex memory management techniques like garbage collection or reference counting.

The arena pattern is particularly useful in scenarios where the structure of the tree is highly dynamic or when performance is critical. It allows for a simple and efficient way to manage and traverse trees without the typical overhead associated with more traditional pointer-based tree representations.

Adaptations for Transformers For transformers, inputs, intermediate values and outputs are all sequences. So the trees are represented as sequences of nodes with node reference representable by token position encoding. Based on the representation, transformers will be able to perform various kinds of recursive tree operations, as we shall see.

B CONTEXT FREE GRAMMAR

905 906

876

877 878

879 880

882 883

885 886

888

889

890

891 892

893

894

895

896

897

898 899

900

901

902

903 904

In this section, we law down the well known de

In this section, we lay down the well-known definitions of context free grammar, derivations, and 907 parse trees. To define an abstract syntax tree (AST), one commonly resorts to generation rules, such 908 as context-free grammars (CFG) (Alfred et al., 2007) and parsing expression grammars (PEG) (Ford, 909 2004). In most cases, just generation rules themselves are not sufficient to define properly a lan-910 guage. Many practical languages like C and C++ cannot be solely described by these rules (David, 911 2009) so that they can reuse the limited set of special characters on the keyboard. Furthermore, se-912 mantic constraints like type correctness are intrinsically contextual and cannot be expressed through 913 CFG or similar rules. However, CFG or other rules provide a valuable construct, the AST. With 914 an AST, one can refine the language definition by putting restrictions on the syntax tree through 915 tree operations. Effectively, a language can be seen as a subset of trees, not as a subset of strings. Semantic analysis like symbol resolution and type checking can be described effectively based on 916 trees. In short, CFG standalone is hardly practical but it provides a useful and clear foundation to 917 build definitions upon.

918	A context-free grammar (CFG) is defined as a 4-tuple $G = (V, \Sigma, R, S)$, where:				
920	• V is a finite set of variables (non-terminal symbols)				
921	 V is a finite set of variables (non-terminal symbols). X is a finite set of terminal symbols, disjoint from V. Sequences of X is a alements of X* 				
922	• Σ is a finite set of terminal symbols, disjoint from V. Sequences of Σ , i.e., elements of Σ^* are called (literal) strings				
923	• $B \subset V \times (V \Sigma)^*$ is a finite set of production rules, where each rule is of the form $A \to \alpha$				
924	with $A \in V$ and $\alpha \in (V \cup \Sigma)^*$.				
925	• $S \in V$ is the start symbol				
920 927					
928	Given a context-free grammar $G = (V, \Sigma, R, S)$, we define derivation as follows:				
929	• A derivation is a sequence of steps where, starting from the start symbol S, each step				
930	replaces a non-terminal with the right-hand side of a production rule.				
931	• Formally, we write $u \Rightarrow v$ if $u = \alpha A\beta$ and $v = \alpha \gamma \beta$ for some production $A \Rightarrow \gamma$ in R,				
932	where $\alpha, \beta \in (V \cup \Sigma)^*$ and $A \in V$.				
934	• A leftmost derivation is a derivation in which, at each step, the leftmost non-terminal is				
935	replaced.				
936	• A rightmost derivation is a derivation in which, at each step, the rightmost non-terminal				
937	is replaced.				
938	• We denote a derivation sequence as $S \Rightarrow^* w$, where $w \in \Sigma^*$ is a string derived from S in				
939	zero or more steps.				
940	A narse tree (or syntax tree) for a context-free grammar $C = (V \Sigma R S)$ is a tree that satisfies				
941 942	The following conditions: $O = (v, 2, n, b)$ is a free that satisfies				
943	• The root of the tree is labeled with the start symbol S				
944	• Each loaf of the tree is labeled with a terminal sumbal from Σ or the ampty string ϵ				
945	• Each leaf of the field with a terminal symbol from Σ of the empty string ϵ .				
946	• Each internal node of the tree is labeled with a non-terminal symbol from V .				
948	• If an internal node is labeled with a non-terminal A and has children labeled with X_1, X_2, \ldots, X_n , then there is a production rule $A \to X_1 X_2 \ldots X_n$ in R.				
949	• The yield of the parse tree, which is the concatenation of the labels of the leaves (in left-to-				
950 951	right order), forms a string in Σ^* that is derived from the start symbol S.				
952					
953	C NEURAL ARCHITECTURES				
954	In this section, we law down the precise mathematical definitions of neural architectures we are going				
955	to use in our proof.				
956 957	Definition 7 (Single-Layer Fully Connected Network with $4 \times$ Intermediate Space).				
958	Civer model dimension d				
959	Given model almension a_{model} , a single-layer jeed-jorward network with an intermediate space ex- panded to 4 times the input dimension is a function from $\mathbb{R}^{d_{model}}$ to $\mathbb{R}^{d_{model}}$ denoted by $f_{c_{model}}$ and defined				
960	as follows:				
961	aiven $X \subset \mathbb{R}^{d_{model}}$ weights $W_{\bullet} \subset \mathbb{R}^{4d_{model} \times d_{model}}$ $W_{\circ} \subset \mathbb{R}^{d_{model} \times 4d_{model}}$ and biases $B_{\bullet} \subset \mathbb{R}^{4d_{model}}$				
962	$B_2 \in \mathbb{R}^{d_{model}}$, the output $f_{fcn}(X)$ is computed as:				
963					
964	$f_{\text{free}}(X) = W_2 \sigma_{\text{Def III}}(W_1 X + B_1) + B_2$				
966	$J_{J}(n(1-j)) = 2 \sim \operatorname{Rel}((j+1)(1-j+2)(j+2)(2))$				
967	where $\sigma_{\text{ReLU}} : \mathbb{R}^{4d_{model}} \to \mathbb{R}^{4d_{model}}$ is the Rectified Linear Unit activation function applied element-				
968	wise, defined by:				
969					
970	$\sigma_{\text{ReLU}}(z) = (\max(z_1, 0), \max(z_2, 0), \dots, \max(z_{4d_{model}}, 0))^\top,$				
971					

for
$$z = (z_1, z_2, \dots, z_{4d_{model}})^\top \in \mathbb{R}^{4d_{model}}$$
.

The choice of a $4 \times$ intermediate space is common in practice, often used in Transformer architectures. Interestingly, this empirical choice turns out to have a useful theoretical property: it allows the network to express any affine transformation, as we'll see in the following proposition.

Proposition 1. A single-layer fully connected network with a $4 \times$ intermediate space, as defined previously, can express any affine map from $\mathbb{R}^{d_{model}}$ to $\mathbb{R}^{d_{model}}$.

Proof. Let $f : \mathbb{R}^{d_{\text{model}}} \to \mathbb{R}^{d_{\text{model}}}$ be any affine map given by f(X) = AX + b, where $A \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ and $b \in \mathbb{R}^{d_{\text{model}}}$. We will construct weights $W_1 \in \mathbb{R}^{4d_{\text{model}} \times d_{\text{model}}}$, $W_2 \in \mathbb{R}^{d_{\text{model}} \times 4d_{\text{model}}}$ and biases $B_1 \in \mathbb{R}^{4d_{\text{model}}}$, $B_2 \in \mathbb{R}^{d_{\text{model}}}$ such that $f_{\text{fcn}}(X) = f(X)$ for all $X \in \mathbb{R}^{d_{\text{model}}}$.

Define:

978

982

989

990 991

$$W_1 = \begin{pmatrix} I_{d_{\text{model}}} \\ -I_{d_{\text{model}}} \\ 0 \\ 0 \end{pmatrix}, \quad B_1 = 0 \in \mathbb{R}^{4d_{\text{model}}},$$

where $I_{d_{\text{model}}}$ is the $d_{\text{model}} \times d_{\text{model}}$ identity matrix, and 0 represents zero matrices of appropriate dimensions. Set:

 $W_2 = (A - A \ 0 \ 0), \quad B_2 = b.$

For any $X \in \mathbb{R}^{d_{\text{model}}}$, compute:

$$f_{\text{fcn}}(X) = W_2 \sigma_{\text{ReLU}}(W_1 X + B_1) + B_2$$

= $(A - A 0 0) \sigma_{\text{ReLU}} \left(\begin{pmatrix} X \\ -X \\ 0 \\ 0 \end{pmatrix} \right) + b$
= $(A - A 0 0) \begin{pmatrix} \sigma_{\text{ReLU}}(X) \\ \sigma_{\text{ReLU}}(-X) \\ 0 \\ 0 \end{pmatrix} + b$
= $A \sigma_{\text{ReLU}}(X) - A \sigma_{\text{ReLU}}(-X) + b.$

Note that $\sigma_{\text{ReLU}}(X) - \sigma_{\text{ReLU}}(-X) = X$, we have:

$$f_{\rm fcn}(X) = AX + b = f(X)$$

Therefore, the network can represent any affine map from $\mathbb{R}^{d_{\text{model}}}$ to $\mathbb{R}^{d_{\text{model}}}$.

Definition 8 (Single-Layer Feed Forward Network with $4 \times$ Intermediate Space). Given model dimension d_{model} and position set Pos, the Transformer Feed Forward Network is a function $f_{ffn} : \mathbb{R}^{\text{Pos} \times d_{model}} \to \mathbb{R}^{\text{Pos} \times d_{model}}$ defined as follows:

For an input $X \in \mathbb{R}^{\text{Pos} \times d_{model}}$, the output $f_{ffn}(X)$ is computed by applying the single-layer feedforward network f_{fcn} (as defined previously) independently to each position:

$$f_{ffn}(X)_p = f_{fcn}(X_p) \quad \forall p \in \text{Pos}$$

1016 1017 1018

1015

1008

where $X_p \in \mathbb{R}^{d_{model}}$ is the *p*-th row of *X*, corresponding to the *p*-th position in the input sequence.

Next, we define the attention mechanism, which is a key component of the Transformer architecture.
This definition presents a hard attention layer with a simplified position encoding. We use hard attention here for theoretical simplicity, as it represents a discrete limit of the more commonly used soft attention mechanism. Hard attention forces the model to make a clear choice about which inputs to focus on, which can simplify analysis and provide clearer insights into the model's behavior. It can be viewed as the limiting case of soft attention as the temperature approaches zero, where the softmax operation becomes increasingly peaked and eventually converges to a one-hot vector.

1026 Definition 9 (Hard Attention Layer with Simplified Position Encoding). Given model dimension 1027 d_{model} , number of heads H, and number of layers L, a transformer with simplified position encoding and hard attention is defined to be a function $f_{attn} : \mathbb{R}^{\text{Pos} \times d_{model}} \to \mathbb{R}^{\text{Pos} \times d_{model}}$ defined by 1028

$$\forall p \in \text{Pos}, f_{attn}(X)_p := W_O Concat \left(Attn^{(1)}(X)_p, \dots, Attn^{(H)}(X)_p \right), \tag{4}$$

where the hth attention head uses hard attention, defined as: 1032

$$Attn^{(h)}(X)_p := \frac{1}{|S_p|} \sum_{p' \in S_p} V_{p'}^{(h)},$$
(5)

1036 where 1037

1039 1040

1041 1042 1043

1044

1045

1046 1047

1048

1049

1051 1052 1053

1029 1030

1031

1034 1035

• $W_{O} \in \mathbb{R}^{d_{model} \times d_{model}}$ are trainable parameters;

•
$$S_p = \arg \max_{p' \in Pos} \left(Q_p^{(h)^{\top}} K_{p'}^{(h)} + \lambda^{(h)^{\top}} \Psi_{p'-p} \right)$$
 with $Q_p^{(h)}, K_p^{(h)}, V_p^{(h)}, \lambda^{(h)}, \Psi_q$ defined by

- $Q_p^{(h)} = W_Q^{(h)} X_p, K_p^{(h)} = W_K^{(h)} X_p$ are vectors of dimension d_{model}/H , with trainable parameters $W_{Q}^{(h)}, W_{K}^{(h)} \in \mathbb{R}^{(d_{model}/H) \times d_{model}};$

- $V_p^{(h)} = W_V^{(h)} X_p$ are vectors of dimension d_{model}/H , linear transformations of X_p with trainable parameters $W_V^{(h)} \in \mathbb{R}^{(d_{model}/H) \times d_{model}}$;

- $\lambda^{(h)} \in \mathbb{R}^2$ are constants depending only on head count h;
- $\Psi_q \in \mathbb{R}^2$ are 2-dimensional vectors depending on relative position q but not on head count h. It is explicitly defined as

$$\Psi_q = \begin{pmatrix} q\\ 1_{q>0} \end{pmatrix}. \tag{6}$$

1054 1055

1064

1068

1069

1073

1074

1075 1076

1077

1078 1079 This formulation allows for both past and future masking.

1056 Having defined the basic components, we can now proceed to describe the full Transformer archi-1057 tecture. This definition builds upon the previously introduced concepts, incorporating them into a 1058 complete model structure. 1059

Definition 10 (Transformer). A Transformer is a function $f_{tf}: \mathbb{R}^{Pos \times d_{model}} \to \mathbb{R}^{Pos \times d_{model}}$ that maps an input sequence to an output sequence through a series of layers, each consisting of a multi-head 1061 attention mechanism and a position-wise feed-forward network (MLP). 1062

Given: 1063

- Input sequence $X \in \mathbb{R}^{\text{Pos} \times d_{model}}$, where Pos is the set of positions and d_{model} is the model dimension.
- Number of layers L.
- Number of attention heads H.

1070 The Transformer computes the output $Y = X^{(L)}$ through recursive application of attention and 1071 feed-forward layers: 1072

• Initialization is given by:

 $X^{(0)} = X.$

- For each layer l = 1, 2, ..., L:
- Compute attention output:

$$\hat{X}^{(l)} = X^{(l-1)} + f_{attn}^{(l)} \left(X^{(l-1)} \right)$$

- Compute feed-forward output: Here: • $f_{attin}^{(l)}$ are hard attention layers with simplified position encoding as previously defined. It operates on $X^{(l-1)}$ and produces an output in $\mathbb{R}^{\text{Pos} \times d_{model}}$.

1088 1089

1093

1094

1080

1082

1084

1086 1087

1090

• $f_{ffn}^{(l)}$ are feed-forward networks with $4 \times$ intermediate space as previously defined. It oper-

 $X^{(l)} = \hat{X}^{(l)} + f_{ffn}^{(l)} \left(\hat{X}^{(l)} \right)$

ates position-wise on $\hat{X}^{(l)}$ and produces an output in $\mathbb{R}^{\text{Pos} \times d_{model}}$.

1091 *Remark* 1. For simplicity, we have omitted the Layer Normalization component typically present 1092 in Transformer architectures. This simplification allows us to focus on the core attention and feedforward mechanisms while maintaining the essential structure of the Transformer.

We use $\operatorname{Tf}_{HL}^{d_{\text{model}}}$ to denote the set of transformers of model size d_{model} , number of heads H and number 1095 1096 of layers L as functions from $\mathbb{R}^{d_{\text{model}}*}$ to $\mathbb{R}^{d_{\text{model}}*}$.

1097 For purpose of proof, we shall also need residual multi-layer perceptron. Functions over local types 1098 are first represented by multi-layer perception, then by Proposition 2 applications of these func-1099 tions over sequences can be representable by transformers. Residual multi-layer perceptron can be 1100 assembled through composition or computer graph, as we shall see. 1101

Here's the definition of a residual MLP Network. 1102

1103 Definition 11 (Residual Multi-Layer Perceptron). A Residual Multi-Layer Perceptron (ResMLP) is a function $f_{resulp} : \mathbb{R}^{d_{model}} \to \mathbb{R}^{d_{model}}$ defined recursively by 1104

1105 1106 1107

 $X^{(0)} = X, \quad X^{(l)} = X^{(l-1)} + f_{fcn}\left(X^{(l-1)}\right), \quad l = 1, 2, \dots, L, f_{resmlp}(X) = X^{(L)}$

where $X \in \mathbb{R}^{d_{model}}$ is the input vector, L is the total number of layers, and $f_{fcn} : \mathbb{R}^{d_{model}} \to \mathbb{R}^{d_{model}}$ 1108 1109 is the Single-Layer Fully Connected Network with $4 \times$ Intermediate Space as previously defined in Definition 7. 1110

1111

We use $\operatorname{ResMlp}_{L}^{d_{\operatorname{model}}} \subset \mathbb{R}^{d_{\operatorname{model}}} \overset{\mathbb{R}^{d_{\operatorname{model}}}}{=} \text{to represent the set of residual MLPs with dimension } d_{\operatorname{model}}$ and L layers, as defined in Definition 11. 1112 1113

1114 The following proposition is quite basic. It demonstrates that any function representable by a 1115 ResMLP can be applied position-wise by a Transformer. 1116

Proposition 2 (Position-wise ResMLP Application is Representable by Transformers). Let f : 1117 $\mathbb{R}^{d_{model}} \to \mathbb{R}^{d_{model}}$ be a function representable by a Residual Multi-Layer Perceptron (ResMLP) as 1118 defined in Definition 11. Then the function $F: \mathbb{R}^{\text{Pos} \times d_{model}} \to \mathbb{R}^{\text{Pos} \times d_{model}}$, defined by applying f 1119 position-wise, 1120

 $F(X)_n = f(X_n), \quad \forall p \in \text{Pos},$

1121 is representable by a Transformer as defined in Definition 10. 1122

1123

1125 1126 1127

1129

1124 *Proof.* Since f is representable by a ResMLP with L layers, it is defined recursively by

$$X^{(0)} = X, \quad X^{(l)} = X^{(l-1)} + f_{\text{fcn}} \left(X^{(l-1)} \right) \text{ for } l = 1, \dots, L,$$

and 1128

 $f(X) = X^{(L)},$

1130 where $f_{\text{fcn}} : \mathbb{R}^{d_{\text{model}}} \to \mathbb{R}^{d_{\text{model}}}$ is the Single-Layer Fully Connected Network with $4 \times$ intermediate 1131 space (Definition 7). 1132

We construct a Transformer with L layers such that, for any input sequence $X \in \mathbb{R}^{\text{Pos} \times d_{\text{model}}}$, the 1133 output $Y = f_{tf}(X)$ satisfies $Y_p = f(X_p)$ for all $p \in Pos$.

To achieve this, we configure the Transformer so that the attention mechanism outputs zero at each layer. This can be done by setting the attention weights to zero, ensuring $f_{\text{attn}}(X^{(l-1)}) = 0$. Consequently, the update equations simplify to

$$\hat{X}^{(l)} = X^{(l-1)}.$$

1139 We then set the feed-forward network $f_{\rm ffn}$ in the Transformer to have the same weights and biases 1140 as $f_{\rm fcn}$ in the ResMLP. The Transformer layer update becomes

$$X^{(l)} = \hat{X}^{(l)} + f_{\text{ffn}}\left(\hat{X}^{(l)}\right) = X^{(l-1)} + \left(f_{\text{fcn}}\left(X_p^{(l-1)}\right)\right)_{p \in \text{Pc}}$$

This recursion matches that of the ResMLP applied position-wise to X. Therefore, after L layers, the Transformer output satisfies $f_{\text{tf}}(X)_p = f(X_p)$ for all $p \in \text{Pos.}$

1148 D CYBERTRON

1149

1177

1178

1185

1146 1147

1138

1141 1142

1150 D.1 INTRODUCTION

It's often difficult to directly prove that transformers or in general other low level forms of computation can express complicated algorithms and even complex software. There are way too many details
as compared with typical mathematical proofs in machine learning theory. Hence, we propose the
domain specific language Cybertron, where we can systematically prove transformers can express
complicated algorithms and complex software with sufficient readability.

(Note: Cybertron is fundamentally different from Mini-Husky! Mini-Husky is the target language that we want transformers to analyze yet Cybertron is the domain specific language we use to prove that transformers can do that.)

RASP (Weiss et al., 2021) is quite close to Cybertron in terms of its design purpose. However, Cybertron is more powerful with advanced algebraic type system, global and local function constructions, etc. These additional mechanisms replace a significant part of the chore in proofs with automatic type checking. Thus, using Cybertron one can argue operations more complicated than simple algorithms can be simulated by transformers.

In the broader perspective of computer science, it's common to use code to prove things. In fact, in the formal verification community, mathematical proofs are viewed as a special case of a much larger universe of possible proof systems (mathlib Community, 2019; Massot, 2024) and constructive proof using code (Harrison et al., 2014; Farooqui, 2021; Jung et al., 2018) is far more applicable with great soundness to the most general settings. In our case, our code doesn't serve as the whole proof but as an important part that contains most of the chores. However, it's totally possible to build a fullfledged formalized proof, despite it might be too costly for a single paper to do.

- Essentially, Cybertron works as follows:
- one builds complicated functions from the composition of smaller functions. We have lemmas that prove that the composed functions are representable by certain architectures given that smaller functions are representable.
 - 2. there is an algebraic type system and every value is strongly typed and immutable, making it highly readable and easy to reason about;
- 1179
 3. there is a distinction between global and local types/functions. Local types are those information hovered over a single token, and global types are sequences of local types, i.e., the collection of information over the whole token stream. One can define a global function by mapping a local function.
- 4. there are many functions that is defined externally, requiring external explanation that they can be represented by transformers.
- It's implemented as a subset of the Rust programming language that can be understood as computation graphs over sequences. It can be executed for testing purposes and we've tested our implementation for a range of inputs and validated its correctness.

1188 D.2 PHILOSOPHY: SEQUENTIAL REPRESENTATION OF EVERYTHING

Before going through the full details, let's first talk about the fundamental philosophy behind trans-former and Cybertron.

One of the fundamental reasons transformers can be easily adapted across multiple modalities, including NLP and CV, is their sequence-to-sequence operation. Everything can be represented as an arbitrary-length sequence. Texts are sequences of words, images are sequences of image patches, videos are sequences of spacetime patches (OpenAI, 2024b), and even graphs with sparse spatial structures can be represented as sequences of indexed nodes with additional information like parent node indices. Since inputs of various modalities can be cast into vector sequences, transformers can be applied to different domains without modifications to their architecture (Dosovitskiy et al., 2020).

Interestingly, this sequence-based thinking is not new. We've actually been representing every thing as sequences since the very early days of computer science. This has been the foundation
 of how data is stored and processed in computers. However, sequence representations were tradi tionally viewed as low-level and sometimes inefficient for practical use, prompting the development
 of higher-level abstractions for programming. The rise of transformers, with their scalable learning
 capabilities, encourages us to reconsider the significance of sequence-based representations.

From a systems perspective, viewing everything as a sequence is the foundational approach in computer science. Data in a computer is stored as a continuous stream of bits. Whether it's text, images, videos, or graphs, this data is represented in computer memory as an ordered sequence of bits. This aligns with how transformers handle different types of input by transforming them into sequences of vectors. Thus, the sequence-based operation of transformers mirrors the sequence-based representation of data in computer memory.

1211 In essence, if a data structure can be represented in computer memory using N bits, it can 1212 be processed as a sequence of bits of length N. This natural sequence representation in memory 1213 is consistent with how transformers process data, which makes them particularly flexible across 1214 different modalities. For example, recent state-of-the-art approaches Wu et al. (2024) show that 1215 transformers can even be trained directly on raw bits of data, further emphasizing this connection.

Moreover, this sequence-based viewpoint offers fresh insights when applied to the domain of programming, particularly in areas such as code generation and analysis. With tools like ChatGPT and Copilot being widely used by developers, the impact of transformers on programming workflows is growing. Understanding the complexity of algorithms and programs expressed in sequence form becomes an interesting area of study, as it reveals new possibilities for how we approach computation.

In comparison to traditional systems like CPUs and human cognition, transformers are highly parallel but shallow in their operation. A transformer processes data in a fixed number of layers, while a CPU executes 10⁹ cycles per second, and humans may take days to process information like reading a book. Transformers, therefore, represent a fundamentally different computational model that is worth studying further in the context of sequence-based operations.

Example 2. Image to Sequence: In computer memory, an image is typically stored as a continuous block of pixel values, often in row-major order, where each pixel's value is encoded as bits in a sequence. When a transformer processes an image, it divides the image into patches (e.g., 16 × 16 pixels), and each patch is flattened into a vector of pixel values. This creates a sequence of patches, where each patch corresponds to a vector. The way transformers represent these patches as a sequence closely aligns with how the image data is sequentially stored in computer memory.

Example 3. Video to Sequence: A video is stored in computer memory as a sequence of frames,
where each frame is essentially an image. In a similar manner to images, these frames are stored as
continuous pixel values. Transformers process videos by dividing the frames into spacetime patches,
where each patch captures a small region of space over a short segment of time. These spacetime
patches are flattened and arranged into a sequence for the transformer to process. The sequential
ordering of these patches matches how video frames and pixel data are stored in computer memory.

Example 4. Graph to Sequence: In computer memory, a graph is typically stored using an adjacency list or adjacency matrix, where nodes and their connections (edges) are stored sequentially in a data structure. Transformers process graphs by representing each node and its features as a vector, and then creating a sequence of these vectors. The sequence may also encode additional

information, such as the parent-child relationships between nodes. This sequence-based representation of graphs is consistent with how graph data is stored in memory, where nodes and edges are arranged in a structured order.

Example 5. Text to Sequence: Text is naturally stored in computer memory as a sequence of characters or words, where each character is encoded as a sequence of bits (such as ASCII or Unicode values). When transformers process text, they convert each word into a word embedding, which is a vector of real numbers. The sequence of word embeddings corresponds to the sequence of characters or words stored in memory. This natural sequential representation of text in both memory and transformers ensures efficient handling of linguistic data.

1251 Example 6. AST (Abstract Syntax Tree) to Sequence: In computer memory, an abstract syntax 1252 tree (AST) is typically stored as a tree-like structure, where each node represents a component of 1253 the program (e.g., operators, variables, or statements). However, this tree can be linearized into a 1254 sequence by traversing it in a specific order (e.g., pre-order traversal). When transformers process 1255 an AST, they convert it into a sequence of tokens, where each token corresponds to a node in the tree. 1256 This sequential representation of the tree in transformers mirrors how the tree is stored as nodes and 1257 edges in memory, and how it can be flattened into a linear sequence.

In conclusion, the sequence-based representation in transformers is not just a novel approach for
 deep learning but is deeply rooted in how data has been stored and processed in computer memory
 since the early days of computing. This consistency between how data is stored in memory and how
 transformers process data as sequences is a key factor in their adaptability across different domains.

1262

1263 D.3 LOCAL AND GLOBAL TYPES

1264 1265 Now we define the type foundation of Cybertron.

1266Types are fundamental objects for programming language theory. Here we use types to faciliate our1267proofs. Type signatures contain rich information that help guarantee correctness of the program.1268Here, we choose a mathematical definition of types that is most convenient for the discussion in1269this paper. We introduce the notion of "local type". Roughly speaking, they are types without heap1270allocation and intended to be represented with $\mathbb{R}^{d_{model}}$ over a single token. For more complicated1271heap-allocated data structures like trees, graphs, etc., we shall represent them by sequences of these1272"local type"s, which translate directly to vector sequences for transformers.

Definition 12 (Local Type). Given a base space B with at least two elements and a countably infinite identification space Ψ , a local type \mathcal{T} over B is a finite set S together with an embedding ϕ from S to B^d and some fixed $d \in \mathbb{N}$ and an identification $\psi \in \Psi$.

1276 For convenience, define Set $(\mathcal{T}) = S$, $d_{\mathcal{T}} = d$ and $\phi_{\mathcal{T}} = \phi$ and $\psi_{\mathcal{T}} = \psi$. And let 0_B , and 1_B be two 1277 different elements of B. And $B^0 := \{0_B\}$ so that $|B^i| = |B|^i$ holds for all $i \in \mathbb{N}$.

1278 1279 1279 1280 1280 1280 1281 1282 Remark 2. We need B to be at least size 2, so that B^d can be as large as we want for d large enough.For typical computer representation, we can take B to be $\mathbf{2} = \{0, 1\}$. For transformers or neural networks in general, we can take B to be \mathbb{R} if we ignore precision. If we don't ignore precision, B should be some finite set of floating point numbers. Thus, we shall keep the generality of B throughout our discussion as all of these settings are important.

1283 *Remark* 3. The role of identification $\psi_{\mathcal{T}} \in \Psi$ is to make two types mathematically different even if 1284 they have the same underlying set, encoding dimension, and encoding. Basically we are establishing 1285 a specialized type of theory tailored towards the expressive power of transformers upon a foundation 1286 of intuitive set theory.

Example 7 (Finite Set). In mathematics, we have the finite set denoted by $[n] = \{0, 1, ..., n-1\}$. Here we use a slightly different notation for a type with underlying set [n] and some encoding.

Example 8 (Position Encoding). Position encoding can be viewed as the encoding of a type denoted by Pos (n) with the underlying set being [n] where n is the context length. Although it has the same underlying set as type [n], it is a different type for a different purpose and might have different encoding.

1293 If B is \mathbb{R} , then the position encoding can be understood as the encoding of type $\llbracket L \rrbracket$ where L is the 1294 context length. More explicitly, we have

$$\phi(x) = (e^{iL^{-i/d}x})_{i \in [d/2]},\tag{7}$$

1296 viewed as a d dimensional \mathbb{R} -vector through the natural conversion of \mathbb{C} to \mathbb{R}^2 , since d is even. 1297 1298 It's too cumbersome to manually give the underlying set and the encoding. Here we introduce a 1299 classical concept from programming language theory Program (2013) that makes it super easy to 1300 construct new types and make things fairly readable. 1301 **Definition 13** (Finite Algebraic Data Type, Mathematical Forms). We define two ways of creating 1302 new types by combining existing types: 1303 1. Sum type. Given types $\mathcal{T}_i = (S_i, \phi_i, d_i)$ over base space B for i = 1, ..., n, we define the sum type of \mathcal{T}_i , denoted by $\sum_{i=1}^n \mathcal{T}_i$, as follows, 1304 1305 1306 • let $S = (\{1\} \times S_1) \sqcup \ldots \sqcup (\{n\} \times S_n);$ • let $d = d_{[n]} + \max_{i=1}^{n} d_i$; 1308 • let $\phi: S \to B^d$ be such that 1309 $\forall i \in \llbracket n \rrbracket, s \in S_i, \phi((i,s)) = \phi_{\llbracket n \rrbracket}(i) \oplus \phi_i(s) \in B^{d_{\llbracket n \rrbracket} + d_i} \subseteq B^d.$ (8)1311 Note that $|S| = \sum_{i=1}^{n} |S_i|$, thus the name sum type. 1312 1313 2. Product type. Given Local Types $\mathcal{T}_i = (S_i, \phi_i, d_i)$ over base space B for i = 1, ..., n, we 1314 define the product type of \mathcal{T}_i , denoted by $\prod_{i=1}^n \mathcal{T}_i$, as follows, 1315 1316 • let $S = S_1 \times \ldots \times S_n$; • let $d = \sum_{i=1}^{n} d_i$; 1317 1318 • let $\phi: S \to B^d$ be such that 1319 $\forall s = (s_1, \dots, s_n) \in S, \phi(s) = \phi_1(s_1) \oplus \dots \phi_n(s_n) \in B^d.$ (9)1320 1321 Note that $|S| = \prod_{i=1}^{n} |S_i|$, thus the name product type. 1322 1323 Although we can define things and refer to things in terms of mathematical equations, it's sometimes 1324 cumbersome to do so. So we shall frequently refer to types using a programming language form, 1325 like CybertronForm or more complicated things like Option<T> a builtin generic type. 1326 **Definition 14** (Unit Type). The unit type is a type with $S = \{0\}$ and $\phi : S \to B^0, 0 \mapsto 0_B$. In 1327 Cybertron, it's denoted as (). 1328 1329 **Definition 15** (Array Type). *Given a type* \mathcal{T} *, the array type of* \mathcal{T} *with length* $\ell \in \mathbb{N}$ *is the type with* $S = S(\mathcal{T})^{\ell}, d = \ell d_{\mathcal{T}} \text{ and } \phi: S \to B^{\ell d_{\mathcal{T}}}, (s_1, \ldots, s_{\ell}) \mapsto \phi_{\mathcal{T}}(s_1) \oplus \ldots \oplus \phi_{\mathcal{T}}(s_{\ell}).$ It's denoted by 1330 1331 \mathcal{T}^{ℓ} . In Cybertron, it's denoted as [T;N]. 1332 **Definition 16** (Vector Type of Finite Capacity). Given a type \mathcal{T} , the vector type of finite capacity of 1333 \mathcal{T} with maximal length $\ell \in \mathbb{N}$ is the type with $S = \bigsqcup_{i=1}^{\ell} \operatorname{Set}(\mathcal{T})^i$, $d = d_{\llbracket \ell \rrbracket} + \ell d_{\mathcal{T}}$ and $\phi : S \to \mathcal{T}$ 1334 $B^{d_{\llbracket \ell \rrbracket} + \ell d_{\mathcal{T}}}, (s_1, \ldots, s_i) \mapsto \phi_{\llbracket \ell \rrbracket}(i) \oplus \phi_{\mathcal{T}}(s_1) \oplus \ldots \oplus \phi_{\mathcal{T}}(s_i) \oplus 0_B \oplus \ldots \oplus 0_B$ with just enough 1335 number of copies of 0_B such that the dimensionality matches. It's denoted by $\mathcal{T}^{\leq \ell}$. In cybertron, it's 1336 denoted as BoundedVec < T, N >. 1337 However, it's cumbersome and obtuse to define and operate in mathematical forms only. So we shall 1339 give a definition closer to actual programming that is more convenient and easy to read. 1340 Definition 17 (Finite Algebraic Data Type, the Code Forms). We define two ways to create new 1341 types: 1342 1343 1. Enum type. An enum type is the sum type of a finite set of variant types. Each variant type 1344 is associated with a different identifier and can be 1345 • unit like, a unit type; 1346 • struct like, a product of several types, each called a field of the variant, and associated 1347 with an identifier; 1348 • tuple like, a product of several types, each called a field of the variant, but not associ-1349

ated with an identifier.

```
1350
                 Syntactically, an enum type is specified as follows,
1351
                  1 enum <type-name> {
1352
                  2
                         <identifier> {
                                                     // 1st variant, struct like
                              <identifier>: <type>, // 1st named field of 1st variant
<identifier>: <type>, // 2nd named field of 1st variant
1353
                  3
                  4
1354
                  5
1355
                  6
                         },
                  7
                          <identifier> {
                                                     // 2nd variant, struct like
1356
                              <identifier>: <type>, // 1st field of 2nd variant
                  8
1357
                  9
                              . . .
                  10
                          },
1358
                  11
                          <identifier> (
                                                     // 3rd variant, tuple like
1359
                                                     // 1st tuple field of 3rd variant
                  12
                              <type>,
                  13
                                                     // 2nd tuple field of 3rd variant
                              <type>,
1360
                  14
                              . . .
1361
                  15
                  16
                          <identifier>,
                                                     // 4th variant, unit like
1362
                  17
                          . . .
1363
                 18 }
1364
                 For example,
1365
1366
                  1 enum Expr {
                          Variable(IdentToken),
                                                    // 1st variant, tuple like
1367
                  2
                                                    // 2nd variant, struct like
                  3
                         Binary {
1368
                             lopd: ExprId,
                  4
1369
                  5
                              opr: BinaryOprToken,
                  6
                              ropd: ExprId,
1370
                  7
                         1.
                         Prefix {
                  8
                                                    // 3rd variant, struct like
1371
                  9
                              opr: PrefixOprToken,
1372
                  10
                              opd: ExprId,
1373
                  11
                         1.
                         Suffix {
                  12
                                                    // 4th variant, struct like
1374
                  13
                              opd: ExprId,
                  14
                              opr: SuffixOprToken,
1375
                  15
                          1.
1376
                  16
                         Panic,
                                                    // 5th variant, unit like
                 17
                     }
1377
1378
              2. Struct type. A struct type is just the product type of
1379
1380
                  1 struct <type-name> {
                  2
                          <identifier>: <type>,
1381
                          <identifier>: <type>,
                  3
1382
                  4
                          . . .
1383
                  5 }
1384
                  1 struct A {
1385
                         a: i32
                  2
1386
                  3 }
1387
1388
        To show how convenient this is, we can define the very useful option type as follows,
1389
        Definition 18 (Option type). For a local type T, we can define an option type as
1390
1391
            enum Option<T> {
         1
         2
                Some(T),
1392
         3
                None
1393
            }
         4
1394
        Definition 19 (Global Types). Global types are defined to be sequences of local types.
1395
1396
        Example 9 (Representation of Graphs). Graphs can be represented as sequences of its nodes. We
1397
        can use position index to use as node references.
1398
1399
        D.4 COMPUTATION GRAPH
1400
1401
        For convenience, we shall use computation graph as a vehicle to describe complicated computa-
1402
        tion processes. Computation graph is close to actual computation process and one can derive an
        understanding of the computation difficulty from the graph's mathematic properties (width, depth,
1403
```

etc.)

1404 1405 1406	Definition 20 (Directed Simple Graph). A directed simple graph G is a pair (V, E) where V is a finite set, and $E \subseteq V \times V$ is called edges.
1407	In the following, we shall simplify the "directed simple graph" to just graph.
1408	Definition 21 (Computation Graph) A computation graph is an acyclic directed graph $G = (V E)$
1409	with additional structures:
1410	
1411	1. for each vertex $v \in V$, there is an associated type, denoted by T_v ;
1412	2 for each vertex $v \in V$ with a positive number of incoming edges, let v_1, \dots, v_n be the other
1413	2. Jor each vertex $0 \in V$ with a positive number of incoming eages, let v_1, \ldots, v_n be the other vertices for the incoming edges, then there is an associated function f from $T \times \ldots \times T$
1414	to $T_{}$
1415	
1416	A computation graph naturally generates a function from source vertices to sink vertices. Let
1417	$v_1^{\text{in}}, \ldots, v_n^{\text{in}}$ be the set of vertices with no incoming edges, and let $v_1^{\text{out}}, \ldots, v_m^{\text{out}}$ be the set of vertices
1418	with no outgoing edges. Then we can construct a function from $T_{v_1^{\text{in}}} \times \cdots \times T_{v_n^{\text{in}}}$ to $T_{v_1^{\text{out}}} \times \cdots \times T_{v_m^{\text{out}}}$
1419	in the following obvious manner:
1420	1 let $(x, y, y) \in \mathcal{T}$ be an input
1421	1. Let $(x_1, \ldots, x_n) \in I_{v_1^{\text{in}}} \times \cdots \times I_{v_n^{\text{in}}}$ be an input;
1422	2. for each v_i^{in} , assign it with value x_i ;
1423	3. for each vertex $v \in V$ with all its incoming vertices v_1, \ldots, v_l assigned with a value, assign
1424	it with the value $f_v(x_{v_1}, \ldots, x_{v_l})$ where x_{v_i} denotes the value assigned to v_i ;
1425	4. repeat the process until all vertices are assigned a value, then take $(x_{v^{\text{out}}}, \ldots, x_{v^{\text{out}}})$ as the
1426	output.
1427	
1428	Our goal is to make a hypothesis class using the above graph. To control the statistical and compu-
1429	tational complexity, we put restrictions on the choice of T_v and f_v , as follows:
1430	Definition 22 (Restricted Computation Graph). Let U be a set of types, and for any $A, B \in U$, there
1431	is a set of functions $Mor(A, B)$ from A to B. We require that $T_v, T_v^{un} \in U$ and $f_v \in Mor(T_v^{un}, T_v)$
1432	where $T_v^m := \prod_{v' \in F} T_{v'}$. We also require that the underlying graph G satisfies certain conditions
1433	(width, depth, etc.)
1434	Definition 23 (Restricted Computation Graph Of Sequences). Let U be a universe such that for a set
1435	of types U_0 all types in U are of the form A^* for some type $A \in U_0$, and $Mor(A^*, B^*)$ are functions
1436	that preserve sequence lengths.
1437	
1430	Given a restriction, the class of functions generated by restricted computation graphs is the central
1439	object to study in this paper. We shall use an even more restricted computation graph of sequences.
1440	We shall argue about the class of functions formed that
1441	1. it's rich enough to contain many interesting operations including SOL, compiler (type in-
1442	ference, static analysis)
1443	2 it's computationally reasonable and can be represented by transformers with pragmatic
1445	bounds
1445	2 it has a massanable statistical complexity
1440	5. It has a reasonable statistical complexity
1448	As a corollary, our theories suggest that transformers can possibly learn to do many interesting things
1449	with reasonable computational and statistical complexity.
1450	To our knowledge, this is the first theoretical paper that gives programatic optimistic bounds for the
1451	power of transformers in a wide range of meaningful language tasks
1452	power of autoronicio in a wrac funge of meaningful language aoko.
1453	Now we introduce graph-theoretical measures that will play key roles in our new complexity theory.
1454	The most basic one is the following:
1455	Definition 24 (Depth of Graph). The depth of a computation graph is defined to the length of the
1456	longest path, denoted by d_G .
1457	

For convenience, we define the following vertex-wise depth.

Definition 25 (Depth of Graph Vertex). *The depth of a vertex v of a computation graph is defined as the length of the longest path with end v, denoted by* d_v .

- 1461 The smaller d_G is, the more parallel the computation is.
- 1462 However, we shall discuss a more nuanced measure, containment, as follows:
- 1464 D.5 FUNCTIONS OVER LOCAL TYPES

Definition 26 (Functions over Local Types). *Given Local Types* \mathcal{T}, \mathcal{R} , the functions from \mathcal{T} to \mathcal{R} are defined to be just the functions from Set (\mathcal{T}) to Set (\mathcal{R}) .

Remark 4. The domains and codomains are all finite sets, so there aren't many constraints we want to enforce here. Basically, these are "discrete" functions.

Definition 27 (Functions over Algebraic Data Types). Let $\mathcal{T}, \mathcal{S}_1, \ldots, \mathcal{S}_m, \mathcal{R}$ be Local Types, and suppose that \mathcal{T} is an algebraic data type, then we can construct functions from $\mathcal{T} \times \mathcal{S}_1 \times \ldots \times \mathcal{S}_m$ to \mathcal{R} as follows,

> 1. suppose that \mathcal{T} is the sum type of $\mathcal{T}_1, \ldots, \mathcal{T}_n$. Then given functions $f_i : \mathcal{T}_i \times \mathcal{S}_1 \times \cdots \times \mathcal{S}_m$ for $i = 1, \ldots, n$, we can construct a function f, by letting

$$f((i,t), s_1, \dots, s_m) = f_i(t, s_1, \dots, s_m),$$
(10)

1476 1477 1478

1479

1480

1473

1474

1475

(Note that we use the pair (i, t) because the underlying set of \mathcal{T} is $\bigsqcup_{i=1}^{n} \{i\} \times Set(\mathcal{T}_i)$.)

2. suppose that \mathcal{T} is the product type of $\mathcal{T}_1, \ldots, \mathcal{T}_n$. Then given a function $f_* : \mathcal{T}_1 \times \cdots \times \mathcal{T}_n \times \mathcal{S}_1 \times \cdots \times \mathcal{S}_m$ for $i = 1, \ldots, n$, we can construct a function f, by letting

1481 1482 1483

1484

1491

$$f((t_1, \dots, t_n), s_1, \dots, s_m) = f_*(t_1, \dots, t_n, s_1, \dots, s_m),$$
(11)

for each $t \in Set(T_i), s_1 \in Set(S_1), \ldots, s_m \in Set(S_m)$.

for each $t \in Set(T_i), s_1 \in Set(S_1), \ldots, s_m \in Set(S_m)$.

It is not enough to just mathematically construct. We should also discuss how neural networks can
represent these functions. We define the representation of functions over Local Types formally as
follows:

Definition 28 (Representation of Functions over Local Types Using Multi-Layer Perceptions). Let \mathcal{T}, \mathcal{R} be Local Types. Given a function f from \mathcal{T} to \mathcal{R} , we say it is representable by MLP of dimension $d \ge \max \{d_{\mathcal{T}}, d_{\mathcal{R}}\}$ and number of layers L, if there exists $\tilde{f} \in \operatorname{ResMlp}_{L}^{d}$ such that

$$\iota_1 \circ \phi_{\mathcal{R}} \circ f = \tilde{f} \circ \iota_2 \circ \phi_{\mathcal{T}},\tag{12}$$

where $\iota_1 : \mathbb{R}^{d_{\mathcal{R}}} \to \mathbb{R}^d$ and $\iota_2 : \mathbb{R}^{d_{\mathcal{T}}} \to \mathbb{R}^d$ are the canonical embeddings by putting zeros to fit the dimensionalities.

1495 Here are some trivially true facts:

1496 **Proposition 3.** [Identities are Representable] For any Local Type \mathcal{T} , the identity map $\mathrm{Id}_{\mathcal{T}}$ is representable in $\mathrm{ResMlp}_1^{d_{\mathcal{T}}}$.

1499
1500 Proof. Just take
$$W_0^{(1)} = I_d, W_1^{(1)} = W_2^{(2)} = 0, B_1^{(1)} = B_2^{(2)} = 0.$$

Proposition 4. [Equality is Representable] The equality function for any local type \mathcal{T} is representable in ResMlp₂^{2d}, where d is the encoding dimension of \mathcal{T} .

1504 *Proof.* Let $x, y \in \mathcal{T}$ be the inputs. We encode them as $\phi_{\mathcal{T}}(x), \phi_{\mathcal{T}}(y) \in \mathbb{R}^d$. The equality function can be represented as:

1503

1507

where A is a large enough positive constant such that the RHS is either 1 or 0.

This can be implemented in two-layer ResMLP with dimension 2d.

 $f_{\rm eq}(x,y) = \min\left(1, A\sum_{i=1}^{d} |\phi_{\mathcal{T}}(x)_i - \phi_{\mathcal{T}}(y)_i|\right),$

Proposition 5. [Boolean NOT is Representable] The Boolean NOT function is representable in ResMlp₁¹.

1515 *Proof.* It's affine.

1519

1522

1525

1532

1537 1538 1539

1553

Proposition 6. [Boolean AND is Representable] The Boolean AND function is representable in ResMlp₁².

1520 *Proof.* Represent each Boolean value as a binary flag within a 1-dimensional vector. Then AND is 1521 just taking the minimum. By $\min(a, b) = b - \sigma_{\text{ReLU}}(b - a)$, we're done.

Proposition 7. [Boolean OR is Representable] The Boolean OR function is representable in Res Mlp_1^2 .

1526 *Proof.* Represent each Boolean value as a binary flag within a 1-dimensional vector. Then OR is 1527 just taking the maximum. By $max(a, b) = a + \sigma_{ReLU} (b - a)$, we're done.

Proposition 8. [THEN_SOME is Representable] The function Bool::then_some : Bool $\times T \rightarrow$ Option T returns Some t if the boolean is true and None otherwise. This function is representable in ResMlp₁^{d+1}.

Proof. Encode the boolean as a binary flag in a (d + 1)-dimensional vector, where the first component indicates the boolean value and the remaining d components hold the value of type T. The residual MLP f_{resmlp} constructs the output Option T by assembling the flag and the value split into positive and negative parts influenced by the flag:

$$f_{\text{resmlp}}(X) = \begin{pmatrix} x_1 \\ \sigma_{\text{ReLU}}(x_{2:d+1} - Ax_1) - \sigma_{\text{ReLU}}(-x_{2:d+1} - Ax_1) \end{pmatrix}.$$

1540 Here, A is a vector of dimension d with all entries positive and large enough to ensure proper 1541 thresholding. Specifically, each entry of A should be larger than the maximum absolute value that 1542 can be represented in the corresponding dimension of type **T**. This ensures that when $x_1 = 1$, the 1543 subtraction $x_{2:d+1} - A$ will always be negative, and when $x_1 = 0$, it will not affect the value.

When the flag is true $(x_1 = 1)$, $\sigma_{\text{ReLU}}(x_{2:d+1} - A) = 0$ and $\sigma_{\text{ReLU}}(-x_{2:d+1} - A)$ retains the negated value, resulting in Some t. When the flag is false $(x_1 = 0)$, both ReLU terms preserve the value, yielding None. Thus, f_{resmlp} effectively implements Bool::then_some within a single layer of the MLP.

Proposition 9. [Option Or is Representable] Let T be a local type, let Option::or be the function that maps two values a,b of type Option T to a value c of type Option T such that c is equal to a when a is not none, and equal to b otherwise. Then Option::or is representable in $\operatorname{ResMlp}_{1}^{2(d+1)}$.

Proof. Each Option T is represented as a (d + 1)-dimensional vector, where the first component is a binary flag indicating the presence (1 for Some, 0 for None), and the remaining d components encode the value. Given inputs $a, b \in Option T$, the residual MLP f_{resmlp} processes the concatenated vector

1558
1559
1560

$$X = \begin{pmatrix} a_{\text{flag}} \\ a_{\text{val}} \\ b_{\text{flag}} \\ b_{\text{val}} \end{pmatrix}$$

The MLP is designed to separate b_{val} into positive and negative parts $(b_+, b_- respectively)$ influenced by a_{flag} . Specifically, it computes:

1565
$$f_{\text{resmlp}}(X) = a_{\text{val}} + \sigma_{\text{ReLU}} \left(b_{+} - Aa_{\text{flag}} \right) - \sigma_{\text{ReLU}} \left(b_{-} - Aa_{\text{flag}} \right)$$
$$= a_{\text{val}} + \sigma_{\text{ReLU}} \left(b_{\text{val}} - Aa_{\text{flag}} \right) - \sigma_{\text{ReLU}} \left(-b_{\text{val}} - Aa_{\text{flag}} \right), \tag{13}$$



 $1, 2, \dots, \text{Depth}(\mathcal{G})$, we perform the following steps in the global ResMLP.

Input Aggregation Layer. We apply a linear transformation to gather the outputs from the predecessor vertices of each vertex at depth k and feed them as inputs to these vertices. Specifically, we define a linear mapping W^(k)_{agg} ∈ ℝ^{D×D} such that:

$$\tilde{X}^{(t_k)} = W_{\text{agg}}^{(k)} X^{(t_{k-1})},$$

where t_{k-1} is the layer after processing depth k-1, and $\tilde{X}^{(t_k)}$ is the aggregated input for the vertices at depth k. The matrix $W_{agg}^{(k)}$ rearranges and combines the outputs from predecessor vertices to provide the correct inputs to each vertex at depth k. Specifically, for each vertex v at depth k, and for each predecessor u of v in the computation graph, the matrix $W_{agg}^{(k)}$ contains entries that copy the output of u into the input positions of v. All other entries in $W_{agg}^{(k)}$ are set to zero or identity as appropriate.

2. Local Computation Layers. For each vertex v at depth k, we simulate its local ResMLP of depth L_v . Since the depths L_v may vary, we pad the local ResMLPs to have a uniform depth $L = \max_v L_v$ by adding identity mappings where necessary. The updates for vertex v are computed as:

$$\begin{aligned} X_v^{(t_k+1)} &= \tilde{X}_v^{(t_k)} + f_{\text{fcn}_v} \left(\tilde{X}_v^{(t_k)} \right), \\ X_v^{(t_k+k')} &= X_v^{(t_k+k'-1)} + f_{\text{fcn}_v} \left(X_v^{(t_k+k'-1)} \right), \quad \text{for } k' = 2, \dots, L_v, \\ X_v^{(t_k+k')} &= X_v^{(t_k+k'-1)}, \quad \text{for } k' = L_v + 1, \dots, L. \end{aligned}$$

Here, f_{fcnv} denotes the single-layer fully connected network (as per Definition 7) for vertex v.

3. State Update. After completing the local computations for depth k, we update the global state vector $X^{(t_k+L)}$ by concatenating the updated states of all vertices:

$$X^{(t_k+L)} = \left(X_v^{(t_k+L)}\right)_{v \in \mathcal{G}}$$

The total number of layers added for depth k is L + 1, accounting for the input aggregation layer and the L layers simulating the local ResMLPs.

By repeating this process for each depth level $k = 1, 2, ..., \text{Depth}(\mathcal{G})$, we simulate the entire computation graph within a global ResMLP of depth $\text{Depth}(\mathcal{G})(\max_v L_v + 1)$.

Lastly, we use the final layer to perform a linear mapping so that the output is in the correct linear representation, clearing out the intermediate values.

1658 Therefore, the function computed by the global ResMLP is equivalent to the function induced by 1659 the computation graph \mathcal{G} , and it is representable in ResMlp^D_{Depth(\mathcal{G})(max, L_v+1)}.

- *Remark* 5. We only prove things around MLPs here. Later, we shall show that this will imply that the induced map operation over sequences can be represented by transformers.
- 1663 1664

1620

1621

1622

1623 1624

1625

1626

1628

1629

1633

1635

1636 1637

1643

1645

1646

1647 1648

1650

D.6 FUNCTIONS OVER GLOBAL TYPES

The task we want transformers to express is too complicated to be cleanly described in one shot. So we introduce the following lemma to significantly simplify things. The lemma shall be useful for our future papers on this topic.

Proposition 13. [Composition of Functions Representable in Transformers] For local types \mathcal{T} , \mathcal{S} , \mathcal{R} , with maps $f: \mathcal{T}^* \to \mathcal{S}^*$ and $g: \mathcal{S}^* \to \mathcal{R}^*$ representable in $\mathrm{Tf}_{H_1,L_1}^{d_1}$ and $\mathrm{Tf}_{H_2,L_2}^{d_2}$ respectively. **Then the composition** $g \circ f$ is representable in $\mathrm{Tf}_{\max\{H_1,H_2\},L_1+L_2}^{\max\{d_1,d_2\}}$.

1673

Proof. This is basically the same as the proof of Proposition 11.

 $\mathrm{Tf}^{d}_{H_{n},L_{n}}$

Remark 6. This doesn't really cover the above. The bound in Proposition 14 isn't always tight for model dimension when the computation graph is deep and Proposition 13 complements it.

1686 Proof. WLOG, assume that $d = \sum_{v \in V} d_v + H d_0$. Then

$$\mathbb{R}^{d} = \underbrace{\left(\bigoplus_{v \in V} \mathbb{R}^{d_{v}}\right)}_{C} \oplus \underbrace{\left(\bigoplus_{h \in [H]} \mathbb{R}^{d_{0}}\right)}_{A}.$$
(14)

Here C stands for "cache" used for storing computed values, and A stands for "active" used for storing intermediate computation results.

1694 Make an order of all the nodes in the graph, say $V = \{v_1, \dots, v_{|G|}\}$ such that $\text{Depth}(v_i) \leq \text{Depth}(v_j)$ if $i \leq j$.

We now imagine the transformer computation process as gradually evaluating the value of each vertex, starting from v_1 to $v_{|G|}$. Every $\max_v L_v$ layers form a layer group, and after each layer group, at most H vertices are assigned values. The equation 14 implies that we have enough memory to cache the computed values and intermediate values in small transformers.

Now let this process continue until we compute all the values. It must be finite because after each layer group, at least one of the vertices is computed. But this bound is too loose. We claim the following:

1704 Claim: the number of layer groups where less than H vertices are assigned values is smaller than 1705 Depth(G).

Sketch of Proof of Claim: for any layer group where less than H vertices are assigned, all the vertices that aren't assigned after this layer group must have larger depth than any vertices that are assigned values before this layer group, otherwise such a vertice can be evaluated in this layer group. Define the depth of any layer group to be the smallest depth of vertices evaluated in this layer group. Then for any unsatiated layer group, it must have a larger depth than the previous layer group. But depth can only increase Depth(G) times, thus there are at most Depth(G) unsatiated layer groups.

Proof of Claim: let V_1, \ldots, V_l be the vertices evaluated at each layer group. Note that l is a different symbol than L and means that the number of layer groups rather than the number of layers.

For convenience, let D_i be the minimum of the depths of vertices in V_i .

1716 Suppose that the *i*th layer group is unsatiated, then i < l. We want to show that $D_i < D_{i+1}$. 1717 Suppose otherwise, i.e., $D_i = D_{i+1}$. Because the *i*th layer group is unsatiated, for any $v \in V_{i+1}$, 1718 v must have dependencies that haven't been evaluated before the *i*th layer group. Choose $v_0 \in V_i, v_1 \in V_{i+1}$ such that $\text{Depth}(v_0) = \text{Depth}(v_1) = D_i = D_{i+1}$. Note that any dependency of v_1 1720 must have smaller depths than v_0 , then must have already be evaluated before the *i*th layer group. 1721 Contradiction!

Now given the claim, we have that for all but at most Depth(G) choices of i = 1, ..., l, we have $|V_i| = H$, then we have

$$|G| = \sum_{i=1}^{l} |V_i| \ge (l - \operatorname{Depth}(G))H$$
(15)

Then $l \leq \frac{|G|}{H} + \text{Depth}(G)$.

Then
$$L \le l \cdot \max_{v \in G} L_v = \left(\frac{|G|}{H} + \text{Depth}(G)\right) \max_{v \in G}$$
.
1730

Proposition 15. [Nearest Left/Right] For any local type T, consider the function that maps a sequence of type Option<T> to nearest left/right neighbors that are not none. It's representable in $Tf_{1,1}^{d+1}$

1735 1736

1741 1742

1748

1749

1751

1731

1737 *Proof.* There is only one layer and one head needed, so we can omit the layer and head index.

WLOG, we consider the nearest left case.

1740 We just need to make the attention exponential look like this:

$$Q_p^{+} K_{p'} + \lambda \Psi_{p'-p} = a_{\text{flag},p'} - 1_{p'-p>0},$$
(16)

where $a_{\text{flag},p'} \in \{0,1\}$ indicates whether the value at position p' is some or none.

1745 We set $V_{p'}$ to represent the value of type Option $\langle T \rangle$.

For the starter token p_0 , we make it such that

$$Q_p^{\top} K_{p_0} + \lambda \Psi_{p_0 - p} = 1, \tag{17}$$

1750 and

 $V_{p_0} = \mathbf{0},\tag{18}$

 \square

so that when there are no some to the left, it will give us none.

Proposition 16. [Nearest Two Left/Right] For any local type T, consider the function that maps a sequence of type Option $\langle T \rangle$ to nearest two left/right neighbors that are not none. It's representable in $\operatorname{Tf}_{O(1),O(1)}^{O(d)}$ where d is the encoding dimension of T.

1758

1759 *Proof.* We can utilize Proposition 15 and 14.

The nearest two left or right is equivalent to first computing the nearest left/right, and then packing them together into one and compute its nearest left/right. The process is represented by a small constant computation graph, then we're done. \Box

1764

1765 D.7 SYNTAX AND SEMANTICS OF CYBERTRON

Having laid the necessary mathematical foundation behind Cybertron, we now turn to explaining its surface—its syntax and semantics. Cybertron serves as a syntax sugar for expressing local and global computation graphs, which are the vehicles used to demonstrate the expressive power of transformers. In Cybertron, computations are divided into two layers: the local world and the global world. These layers play distinct but complementary roles in constructing computation graphs.

- 1771
- 1772 D.7.1 LOCAL WORLD 1773

The **local world** in Cybertron corresponds to the feed-forward layers of a transformer, focusing on computations over **local types**. Local types represent individual tokens or data points, and computations in this world handle operations on tokens independently of their surrounding context.

1777

Data Types. Local types in Cybertron include basic types such as Bool, Idx, Pos, Fin
BoundedVec<T, N>, etc. These types are essential for building local computation graphs that operate over individual tokens. Compound types, like structs and enums, can also be defined for more complex token representations. These types serve as the building blocks for the local computation graphs that transform data at the token level.

```
1782
         1 struct Node {
1783
          2
                 id: Idx,
          3
                 position: Pos,
1784
            }
          4
1785
          5
          6
            enum Operation {
1786
          7
                 Add {
1787
                     lhs: Pos,
          8
          9
                     rhs: Pos,
1788
         10
1789
                 Multiply {
         11
         12
                     factor: Pos,
1790
                 },
         13
1791
         14
            }
```

1797

1799

1801

1802

1803

1805

1807 1808

Functions. Functions in the local world define operations upon information over individual tokens.
These operations form nodes in the local computation graphs. For instance, operations like binary or unary expressions, conditionals, and matches on token types are transformed into computation graphs by handling each individual token's data.

```
1 fn process_ast(ast: AstData) -> Option<Role> {
2
       match ast {
3
           AstData::LetInit { pattern, initial_value, .. } => {
4
               Some(Role::LetStmt { pattern, initial_value })
5
           AstData::Defn { keyword, ident, .. } => {
6
                Some (match keyword {
7
8
                    DefnKeyword::Struct => Role::StructDefn(ident),
9
                   DefnKeyword::Enum => Role::EnumDefn(ident),
10
                    DefnKeyword::Fn => Role::FnDefn(ident),
                })
11
12
           }
             => None,
13
14
       }
15
   }
```

Control Flow. In the local world, control flow structures such as **if** and **match** are transformed into computation graphs by treating each branch or arm as an expression that returns an Option based on conditions. These Option values are then combined using the Option::or function. According to **Proposition 9**, Option::or maps two Option<T> values and returns the first non- None value, or the second one otherwise. This allows conditional branches to be represented in computation graphs as sequential option evaluations, where the first matching condition provides the result.

1817 1818 D.7.2 GLOBAL WORLD

The global world extends beyond individual tokens to sequences of tokens, represented as global types. These global types are denoted as Seq<T>, where T is a local type. The global world represents the full transformer, focusing on operations involving sequences of tokens, including variable definitions, expressions involving variable references, and function calls.

Variable Definitions. Variables in the global world are defined using global types, which represent sequences of local tokens. These definitions correspond to nodes in the global computation graph.

1826

1830

Expressions. Expressions in the global world consist of references to variables or function calls.
 Since the global world operates over sequences of tokens, these expressions are translated into sequence-level operations in the computation graph.

Function Calls. Function calls are key elements of the global world. They are represented by applying global functions to sequences of tokens. Cybertron provides **map functions** to elevate local functions to global functions by mapping them across sequences. Additionally, **attention methods** like nearest_left and nearest_right handle dependencies between tokens in the sequence by identifying relationships based on their positions.

1 let result = seq_of_values.nearest_left();

In the global world, computation graphs are built by composing map functions and attention methods. These graphs, unlike those in the local world, do not include control flow mechanisms.

1839 D.8 DYCK LANGUAGE

1845

1862

1868

1871

1872

1873

1874 1875 1876

1877

1878

This section demonstrates how the local world in Cybertron operates over token-level computations and how the global world handles sequence-level operations. We use a Dyck language example to explain the interactions between these two worlds. The example processes a sequence of delimiters (like parentheses) and checks for matching pairs.

Local World. In Cybertron, the local world operates on individual tokens. Here, the local types are simple, such as Delimiter and PreAst, which represent information associated with individual tokens. These types allow for token-level operations like comparisons and transformations.

We define a struct to represent a delimiter and an enum to classify delimiters as either left or right.
These definitions reflect local types, as they hold information over a single token.

```
1851
            // Define a struct 'Delimiter' that wraps a 'u8' value.
           #[derive(Debug, Clone, Copy, PartialEq, Eq)]
         2
1852
           pub struct Delimiter(u8);
         3
            // Define an enum 'PreAst' which represents a left or right delimiter.
         5
           #[derive(Debug, Clone, Copy, PartialEq, Eq)]
         6
1855
           pub enum PreAst {
         7
                LeftDelimiter(Delimiter),
         8
                RightDelimiter (Delimiter),
1857
        10 }
```

Here, the local types Delimiter and PreAst define operations upon individual tokens, representing
fundamental units of the computation graph at the local level. The local world is responsible for
handling these small, token-level computations independently of the global sequence.

Global World. In the **global world**, Cybertron operates on sequences of tokens, treating the collection of local types as a single unit of computation. The global world introduces global types such as Seq<Option<PreAst>>, which represents a sequence of optional delimiters. The global world handles sequence-level operations by applying functions like nearest_left and nearest_right to capture the relationships between tokens in the sequence.

The following function operates on a sequence of PreAst, reducing matched pre-asts. The recursive application of step gives us the classifier for Dyck language.

```
1 fn step(pre_asts: Seq<Option<PreAst>>) -> Seq<Option<PreAst>> {
2    let pre_asts_nearest_left = pre_asts.nearest_left();
3    let pre_asts_nearest_right = pre_asts.nearest_right();
4    step_aux.apply(pre_asts_nearest_left, pre_asts, pre_asts_nearest_right)
5 }
```

Local Worlds. The step_aux function matches tokens based on their nearest neighbors within the sequence, eliminating pre-asts if a match is found.

```
fn step_aux(
        1
1879
                pre_ast_nearest_left: Option<(Idx, PreAst)>,
         2
         3
                pre_ast: Option<PreAst>,
1880
         4
                pre_ast_nearest_right: Option<(Idx, PreAst)>
1881
         5
            ) -> Option<PreAst> {
         6
                match pre ast? {
1882
         7
                    PreAst::LeftDelimiter(delimiter) => match pre_ast_nearest_right {
                        Some((_, PreAst::RightDelimiter(delimiter1))) if delimiter1 == delimiter =>
         8
                None,
1884
         9
                        _ => pre_ast,
1885
         10
                    1.
                    PreAst::RightDelimiter(delimiter) => match pre_ast_nearest_left {
         11
         12
                        Some((_, PreAst::LeftDelimiter(delimiter1))) if delimiter1 == delimiter =>
                None,
         13
                        _ => pre_ast,
        14
                    },
        15
                }
        16
            }
```

In this example, the global function step uses nearest_left and nearest_right to capture sequencelevel dependencies, while the local function step_aux uses conditional logic to check for matching pairs of delimiters. The local world handles token-level logic, while the global world coordinates operations across the entire sequence. This separation reflects how Cybertron handles computations at different levels of granularity.

Thus, this example illustrates how Cybertron leverages both the local and global worlds to build
comprehensive computation graphs in a convenient, comprehensive yet rigorous manner. The local
world performs individual tokenwise operations, and the global world captures relationships between tokens in a sequence, demonstrating how Cybertron enables transformers to express complex
computations.

1901 1902

1903

1905

1940

1941

1942

1943

1890

E MINI-HUSKY DETAILS

Here's the BNF grammar of the Mini-Husky language:

1906	
1907	$\langle ast \rangle ::= \langle literal \rangle$
1908	$ \langle ident \rangle$
1909	$ \langle prejlx \rangle $
1910	$ \langle suffix \rangle $
1911	(delimited)
1912	$ \langle separated_{item} \rangle$
1913	$\langle Call \rangle$
1914	$ \langle if_stmt \rangle$
1915	$\langle else_stmt \rangle$
1916	$ \langle defn \rangle$
1917	$\langle literal \rangle ::=$
1910	$\langle ident \rangle ::=$
1919	$\langle prefix \rangle ::= \langle prefix_opr \rangle \langle ast \rangle$
1921	$\langle binary \rangle ::= \langle ast \rangle \langle binary_opr \rangle \langle ast \rangle$
1922	$\langle suffix \rangle ::= \langle ast \rangle \langle suffix_opr \rangle$
1923	$\langle delimited \rangle ::= \langle left_delimiter \rangle \langle separated_item \rangle^* \langle right_delimiter \rangle$
1924	$\langle separated_item \rangle ::= [\langle ast \rangle] \langle separator \rangle$
1925	$\langle call \rangle ::= \langle ast \rangle \langle left_delimiter \rangle \langle ast \rangle^* \langle right_delimiter \rangle$
1926	$\langle let_init \rangle ::= let \langle ast \rangle$
1927	$\langle if_stmt \rangle ::= if \langle ast \rangle \langle delimited \rangle$
1928	$\langle else \ stmt \rangle ::= \langle if \ stmt \rangle \ else \ (\langle delimited \rangle \mid \langle else \ stmt \rangle)$
1929	$\langle defn \rangle ::= \langle defn \ keyword \rangle \langle ident \rangle \langle ast \rangle$
1931	$\langle prefix_opr \rangle ::= + - ! \dots$
1932	$(binary_opr) ::= + - * / $
1933	$\langle suffix opr \rangle ::= ++ $
1934	$\langle left delimiter \rangle ::= `(' [{$
1935	$\langle right delimiter \rangle ::= \langle \gamma \rangle 1 \}$
1936	$\langle separator \rangle := 1$
1937	$\langle defn \ konvord \rangle := def fn $
1938	$\langle acjn_kcywora \rangle \dots = \langle acr + nr + \dots$
1939	

Below is a sample piece of codes:

```
1 struct Dog { weight: f32, ... }
2
3 fn see_vet(dog: Dog) -> f32 {
4 assert dog.weight < 100;
5 let mut fee = dog.weight * 10.0;</pre>
```

```
1944

6 fee +=100.0;

1945 7 return fee

1946 8 }
```

It should be noted that the above is not the full story. There are additional constraints put on the ASTs. However, these can be easily implemented as tree functions that are easy for transformers to express. As we are focusing on higher level language processing capabilities, we ignore the details here.

Additionally, we need to require that for semantic correctness, we must have proper symbol resolution and type correctness.

1954

1956

1947

1948

1949

1950

1951

```
1955 E.1 Additional Details about Compiler Tasks.
```

1957 The outputs of the tasks are defined using Cybertron as follows:

The construction of AST task's final output is the collection all AST nodes. More concretely, the output is a sequence of Option<Ast> with length equal to the input token sequence's length, where Option<Ast> denoted the type Ast will a null value added and Ast is the type storing the information of a node, including its parent, and its data of type AstData. In Cybertron, we define Ast and AstData explicitly as follows:

```
1964
           1 /// Represents a node in an Abstract Syntax Tree (AST).
            2 ///
1965
              /// Each 'Ast' node has a reference to its parent node (if any) and holds
            3
1966
              /// the associated syntax data (such as expressions, statements, or other
/// constructs defined in the 'AstData' enum).
            4
1967
            5
            6 pub struct Ast {
1968
                   /// The index of the parent node in the AST, if it exists.
1969
            8
                   /// - <code>`Some(Idx)</code> <code>:</code> The node has a parent, and <code>`Idx`</code> represents its position.
            9
1970
                   /// - 'None': The node is the root or does not have a parent.
           10
                   pub parent: Option<Idx>,
           11
1971
           12
                   /// The data associated with this AST node.
1972
           13
                   pub data: AstData,
           14 }
1973
           15
1974
           16 /// Enumeration representing different types of Abstract Syntax Tree (AST) nodes
           17 pub enum AstData {
1975
                   /// Represents a literal value (e.g., integer, string)
           18
1976
           19
                   Literal(Literal),
                   /// Represents an identifier (e.g., variable name)
           20
           21
                   Ident(Ident),
1978
                     // Represents a binary expression (e.g., 'x + y', 'a * b')
           22
                   Binary {
           23
1979
           24
                       /// Index of the left operand
           25
                       lopd: Idx,
                        /// Operator in the binary expression (e.g., `+`, `*`)
           26
1981
           27
                       opr: BinaryOpr,
1982
                        /// Index of the right operand
           28
           29
                       ropd: Idx,
1983
           30
                   },
1984
           31
                   ... // other variants
           32 }
```

1986 • The output of the symbol resolution task is the collection of symbol resolution results 1987 on all applicable tokens. More concretely, the output is a sequence of values of type 1988 Option<SymbolResolution> where Option<SymbolResolution> is the type SymbolResolution with 1989 a null value added for non-applicability and SymbolResolution is the type storing the result of the 1990 symbol resolution, being either a success with a resolved symbol of type Symbol or a failure with 1991 an error of type SymbolResolutionError . In Cybertron, we define SymbolResolution explicitly as 1992 1993 follows:

1 // an enum type definition, basically a tagged union type
1995 2 pub enum SymbolResolution {
1996 3 Ok(Symbol), // enum type variant for success with a resolved symbol
2 Err(SymbolResolutionError), // enum type variant for failure with an error
2 }

1998 • The type analysis task's final output is the collection of all type errors. More concretely, the output 1999 is a sequence of Option<TypeError>, where Option<TypeError> denoted the type TypeError will 2000 a null value added and TypeError is the type storing the information of a type error. The position of type errors agrees with the source tokens leading to these errors. In Cybertron, we define 2002 TypeError explicitly as follows: 2003 2004 1 // This enum represents various kinds of type errors 2005 pub enum TypeError { 2 // This variant indicates a type mismatch 3 2006 // 'expected' is the type that was anticipated 4 2007 // `actual` is the type that was encountered 5 6 TypeMismatch { expected: Type, actual: Type }, 2008 7 } 2009 One can expand the definition to include other kinds of type errors. 2010 2011 (1) Type definition. Types are either identified uniquely by a single identifier like <identifier>, or 2012 builtin generic types Option<<identifier>> or Vec<<identifier>> . Users can define custom types 2013 without generics like the following (f32 means float32 and i32 means int32 below): 2014 1 struct Dog { weight: f32 } 2015 2 2016 3 enum Animal { 4 Dog, 2017 5 Cat, 2018 6 } 2019 This part is actually a part of the AST task and type definition is a variant of the AstData type: 2020 2021 1 /// Enumeration representing different types of Abstract Syntax Tree (AST) nodes 2022 pub enum AstData { 2 3 2023 /// Represents a function or variable definition 4 2024 5 /// # defn 6 2025 7 2026 8 Defn { /// The keyword in the definition (e.g., 'fn', 'enum') 9 2027 10 keyword: DefnKeyword, 2028 11 /// Index of the identifier in the definition 12 ident_idx: Idx, 2029 13 /// The identifier being defined (e.g., function name, variable name) 2030 14 ident: Ident, 15 /// Index of the content or body of the definition 2031 16 content: Idx, 2032 17 }, 18 } 2033 (2) Type specification. Each appeared variable has a unique type, either by specification or specu-2035 lation. All parameters of a function must be specified explicitly by users. Variables defined by let 2036 statements might or might not be specified, as follows: 2037 fn f(a: i32) { // type of `a` must be specified let x: i32 = a; // type of `x` specified
let y = a; // type of `y` unspecified 2 2039 4 } 2041 The return type of functions must be specified. The field type of structs and enum variants must be 2042 specified. the type of expressions of function calls and field access will be determined correspond-2043 ingly. 2044 The output of the task is the collection of all type signatures, represented as a sequence of values of 2045 type Option<TypeSignature> where TypeSignature is the type holding the essential information of 2046 2047 type specifications. In Cybertron, TypeSignature is defined as, 2048 1 pub struct TypeSignature { 2049 2 pub key: TypeSignatureKey, 2050 3 pub ty: Type,

4 }

6 pub enum TypeSignatureKey {

2051

```
2053
2054
2055
2056
2057
```

2060

2061

2052

7 FnParameter { fn_ident: Ident, rank: Rank }, 8 FnOutput { fn_ident: Ident }, 9 StructField { ty_ident: Ident, field_ident: Ident }, 10 }

(3) Type inference. As discussed above, not all variables have their types specified.

```
1 fn f() {
2   let x: i32 = 1;
3   let y = x;
4   let z = y;
5 }
```

In the above code, 1 is an ambiguous literal that can be of type i32, i64, u32, u64, etc, and the types of y and z is not specified. However, one easily sees that there exists one and only one choice of the types of 1, y, and z such that the whole code is type correct. Utilizing this property, the user can opt out of a significant portion of type specification, achieving static guarantees.

A Type Inference Algorithm: For simplicity, we shall prove transformers can implement a simple type inference algorithm: we maintain a table of type assignments for variables. We update the entries of the table by means of reduction, i.e., assuming the whole code is correctly typed and infer more and more unspecified types until we encounter errors or all types are inferred. The process is largely parallel, and we call the number of rounds needed the depth of type inference.

In the above code, the first round, we determine that the type of both 1 and the type of y are equal to the type of x which is i32. But we have no way to determine the type of z because the type of y is unknown at the first round. In the second round, z can be determined to be of type i32 because the type of y is already inferred.

2077
2078
2078
2079The output of the task is the collection all types inferred, represented as a sequence of values of type
Option<TypeInference> where TypeInference is the type holding the inferred type. In Cybertron,
TypeSignature is defined as,

```
1 pub struct TypeInference {
2     pub ty: Type,
3 }
```

2086 2087

2088

2081

F TRANSFORMER AST PROOF

F.1 HIGH LEVEL OVERVIEW

²⁰⁸⁹ Here we give the full details of the proof of transformers being able to parse ASTs.

2090 On a high level, we are going to see the parsing of ASTs as an assembly process. First, we im-2091 mediately get the atomic ones, like identifiers, literals, etc. Then we assembly all composite ASTs 2092 with enough precedence util all tokens are consumed. We can prove that at the *n*th round, all ASTs 2093 with depth no more than n are already constructed. In the process, we must keep track of the un-2094 consumed tokens and newly constructed ASTs (to be consumed as children for new ASTs in the 2095 next round, as we are going bottom up). We use pre_asts to denote all the unconsumed tokens and 2096 newly constructed ASTs and use asts to denote all the constructed (allocated) ASTs. For correctness 2097 guarantees, we give detailed type specifications for tokens, ASTs, and PreASTs as follows.

2098 2099 We define the Token type as follows:

2100 /// The 'Token' enum represents the various types of tokens that can be 2 /// identified during the lexical analysis phase of a compiler. Each variant 2101 3 $\ensuremath{///}$ corresponds to a specific category of token that can be encountered 4 /// in the source code. 2102 pub enum Token { 5 2103 6 /// A literal value, which can be a number, string, or other primitive type. 2104 7 Literal(Literal), 8 /// A reserved keyword in the language, such as `if`, `else`, `while`, etc. 2105 Keyword(Keyword), 9 10 /// An identifier, typically representing variable names, function names,

2100		
2107	11	/// or other user-defined symbols.
	12	Ident(Ident),
2108	13	/// An operator, such as `+`, `-`, `*`, `==`, etc., representing mathematical
0100	14	/// or logical operations.
2109	15	Opr(Opr),
2110	16	/// A left delimiter, such as `(`, `{`, `[`, used to denote the beginning of
	17	/// a block, list, or expression.
2111	18	LeftDelimiter(LeftDelimiter),
2112	19	/// A right delimiter, such as ')', '}', ']', used to denote the end of a
0110	20	/// block, list, or expression.
2113	21	RightDelimiter (RightDelimiter),
2114	22	/// A separator, such as `,` or `;`, used to separate elements in a list or
	23	/// statements in a block.
2115	24	Separator(Separator),
2116	25 }	

The type has an encoding dimension $d_{\text{Token}} = \Theta(\log L)$, which is large enough to faithfully represent its information.

2120 More specifically, the types Literal, Keyword, Ident, Opr, LeftDelimiter, RightDelimiter, 2121 Separator are local types assumed to have encoding dimension less than d_{Token} . Keyword, Opr, 2122 LeftDelimiter, RightDelimiter, Separator are small, so they can be encoded in a straight-forward 2123 manner entirely using d_{Token} . However, Literal and Ident are larger than representable by a lim-2124 ited number of bits because potentially a Literal can be a string literal of arbitrary length and an 2125 2126 Ident can also be of arbitrary length. This can be solved through methods like interning, which gives all literals and identifiers that actually appear in the input distinct encodings. As the context 2127 length is L, the number of different literals/identifiers are bounded by context length and interning 2128 needs $O(d_{\text{Token}}) = O(\log L)$ to work. As far as our theories are concerned, it's totally reasonable to 2129 assume that all these types are assumed to have encoding dimension less than $d_{\text{Token}} = O(\log L)$. 2130

2131 We define AST type as follows:

2106

2117

2150

```
2132
        1 /// Represents a node in an Abstract Syntax Tree (AST).
2133
         2 ///
         3 /// Each 'Ast' node has a reference to its parent node (if any) and holds
2134
         4 /// the associated syntax data (such as expressions, statements, or other
         5 /// constructs defined in the 'AstData' enum).
2135
         6 pub struct Ast {
2136
         7
               /// The index of the parent node in the AST, if it exists.
2137
         8
                /// - 'Some(Idx)': The node has a parent, and 'Idx' represents its position.
         9
2138
                /// - 'None': The node is the root or does not have a parent.
         10
                pub parent: Option<Idx>,
2139
         11
        12
                /// The data associated with this AST node.
2140
         13
2141
        14
                /// This field holds the actual syntax information, which is typically
        15
                /// defined by the 'AstData' enum. This could represent literals, expressions,
2142
        16
                /// statements, and other constructs in the source language.
2143
        17
                pub data: AstData,
        18
           }
2144
```

Note that we intentionally structure the tree by always storing the parent but not necessarily storing all children information. In our assumptions, we only control the depth of ASTs but don't control the number of children. More specifically, a function can have as many statements as possible. To avoid overflowing, we don't store all children information. As we shall see, parent information alone is enough for transformers to perform tree operations.

The AstData is the most complicated we define in this paper, as follows:

```
1 /// Enumeration representing different types of Abstract Syntax Tree (AST) nodes
2152
         2
            pub enum AstData {
2153
         3
                 /// Represents a literal value (e.g., integer, string)
2154
          4
                Literal(Literal),
                /// Represents an identifier (e.g., variable name)
          5
2155
                Ident(Ident),
          6
                /// Represents a prefix expression (e.g., <code>`!x`, `-x`)</code>
2156
          7
          8
2157
          9
                /// # exprs
2158
         10
         11
                Prefix {
2159
                    /// Operator in the prefix expression (e.g., `!`, `-`)
         12
                     opr: PrefixOpr,
         13
```

2160						
2161	14 15	/// Operand index of the expression opd: Idx.				
2162	16	<pre>>por 100, p },</pre>				
2163	17	<pre>/// Represents a binary expression (e.g., `x + y`, `a * b`) Pinarum (</pre>				
2164	19	/// Index of the left operand				
2165	20	lopd: Idx,				
2166	21 22	<pre>/// uperator in the binary expression (e.g., `+`, `*`) opr: BinaryOpr,</pre>				
2167	23	/// Index of the right operand				
2168	24 25	ropa: iax, },				
2169	26	<pre>/// Represents a suffix expression (e.g., `x++`, `y`) Outfin (</pre>				
2170	27 28	/// Index of the operand				
2171	29	opd: Idx,				
2172	30 31	<pre>/// Operator in the suffix expression (e.g., '++', '') opr: SuffixOpr,</pre>				
2173	32	},				
2174	33 34	<pre>/// Represents a delimited expression (e.g., `(x + y)`, `{a, b, c}`) Delimited {</pre>				
2175	35	/// Index of the left delimiter in the expression				
2176	36 37	left_delimiter_idx: Idx, /// The left delimiter (e.g., `(`, `{`)				
2170	38	left_delimiter: LeftDelimiter,				
2178	39 40	/// The right delimiter (e.g., `)`, `}`) right delimiter: RightDelimiter.				
2170	41	<u>j···</u>				
2179	42 43	<pre>/// Represents an item separated by a separator (e.g., elements in an array or list) SeparatedItem {</pre>				
2100	44	/// Index of the content, if any				
2101	45 46	content: Option <idx>,</idx>				
2102	47	separator: Separator,				
2103	48 49	}, /// Represents a function call or array access (e.g. `f()` `a[]`)				
2104	50	/// Representes a function carr of array access (e.g., f(, / a[] / ///				
2100	51	/// things like `f()` or `a[]`				
2100	53	/// Index of the caller (e.g., function or array)				
2107	54	caller: Idx, $(// The left delimiter of the call (e.g.,)(),)())$				
2188	56	left_delimiter: LeftDelimiter,				
2189	57	<pre>/// The right delimiter of the call (e.g., `)`, `]`) might delimiter. PightPolimiter</pre>				
2190	59	/// Index of the delimited arguments in the call				
2191	60 (1	delimited_arguments: Idx,				
2192	62	}, /// Represents a `let` statement with an initialization (e.g., `let x = 5;`)				
2193	63					
2194	64 65	/// # stmts ///				
2195	66 67	LetInit {				
2196	67	<pre>/// Index of the expression in the initialization expr: Idx,</pre>				
2197	69 70	/// Index of the pattern being initialized				
2198	70	/// Optional index of the initial value				
2199	72	<pre>initial_value: Option<idx>,</idx></pre>				
2200	73 74	}, /// Represents an `if` statement				
2201	75	If {				
2202	76	<pre>/// Index of the condition in the 'if' statement condition: Idx,</pre>				
2203	78	/// Index of the body of the 'if' statement				
2204	79 80	body: Idx, }.				
2205	81	/// Represents an 'else' statement				
2206	82 83	Else { /// Index of the associated `if` statement				
2207	84	if_stmt: Idx,				
2208	85 86	<pre>/// Index of the body of the 'else' statement body: Idx.</pre>				
2209	87	},				
2210	88 89	/// Represents a function or variable definition				
2211	90	/// # defn				
2212	91 92	/// Defn {				
2213	93 94	<pre>/// The keyword in the definition (e.g., `fn`, `enum`) keyword: DefnKeyword,</pre>				

```
2214
         95
                     /// Index of the identifier in the definition
2215
         96
                    ident idx: Idx,
         97
                     /// The identifier being defined (e.g., function name, variable name)
2216
         98
                     ident: Ident,
2217
         99
                     /// Index of the content or body of the definition
        100
                    content: Idx,
2218
        101
                },
2219
        102 }
2220
         1 /// The 'PreAst' enum represents the intermediate forms of tokens and ASTs that are
         2 /// encountered during the parsing phase, before the final AST is constructed.
2222
         3 /// Each variant corresponds to a specific type of token or partial
           /// AST node that contributes to the construction of the final AST.
         4
2223
         5 #[derive(Clone, Copy, PartialEq, Eq)]
2224
         6 pub enum PreAst {
                /// A reserved keyword in the language, such as 'if', 'else', 'while', etc.
         7
2225
                Keyword(Keyword),
         8
                /// An operator, such as `+`, `-`, `*`, `==`, etc., representing mathematical
/// or logical operations.
2226
         9
         10
2227
         11
                Opr(Opr),
2228
                /// A left delimiter, such as `(`, `{`, `[`, used to denote the beginning of
         12
                /// a block, list, or expression.
         13
2229
         14
                LeftDelimiter(LeftDelimiter),
                /// A right delimiter, such as ')', '}', ']', used to denote the end of a
2230
         15
                 /// block, list, or expression.
         16
2231
         17
                RightDelimiter(RightDelimiter),
2232
         18
                /// A partially constructed AST node, representing a more complex structure
                /// that will be further processed to build the final AST.
         19
2233
         20
                Ast (AstData),
                /// A separator, such as `,` or `;`, used to separate elements in a list or
2234
         21
                 /// statements in a block.
         22
2235
         23
                Separator(Separator),
2236
        24 }
2237
         1 /// this is beyond the scope of Cybertron
         2
            /// rather a general Rust function to integrate for testing
2239
         3
         4 pub fn calc_asts_from_input(input: &str, n: usize) -> (Seq<Option<PreAst>>,
2240
                 Seq<Option<Ast>>) {
         5
                let tokens = tokenize(input);
2241
                let pre_asts = calc_pre_ast_initial_seq(tokens);
         6
2242
         7
                let allocated_asts: Seq<Option<Ast>> = tokens.map(|token| token.into());
                reduce_n_times(pre_asts, allocated_asts, n)
2243
         8
         9 }
2244
2245
        The reduce function in Cybertron is designed to progressively refine sequences of pre-abstract
2246
        syntax trees (pre-ASTs) and allocated abstract syntax trees (ASTs). The function takes two input
2247
        sequences: pre_asts, which is a sequence of optional pre-ASTs, and allocated_asts, which is a
2248
        sequence of optional ASTs. It returns a tuple containing the reduced sequences of pre-ASTs and
2249
        allocated ASTs.
2250
        The reduction process is carried out in multiple stages, each focusing on different syntactic con-
2251
        structs:
2252
              1. reduce by opr : This step handles reduction by dealing with operators and their precedence.
2254
                 It simplifies expressions involving operations to form more compact ASTs.
2255
              2. reduce by delimited : This step reduces constructs that are delimited, such as those involv-
2256
2257
                 ing parentheses, braces, or other grouping symbols. It ensures that delimited blocks are
                 properly nested and combined in the AST.
2258
2259
              3. reduce by call : In this stage, function or method calls are reduced. This involves iden-
2260
                 tifying and structuring calls within the AST, ensuring correct representation of function
2261
                 invocations.
2262
              4. reduce_by_stmt: This reduction step addresses statements, ensuring that individual state-
2263
                 ments are correctly parsed and represented within the AST, such as assignment statements,
2264
                 loops, and conditionals.
2265
              5. reduce_by_defn : Finally, reduction by definition handles the parsing of definitions, such
                 as variable or function declarations. This step ensures that all definitions are correctly
                 represented within the AST.
```

By sequentially applying these reduction steps, the reduce function progressively transforms the initial sequences into their most refined forms, ready for further syntactic or semantic analysis.

```
1 pub fn reduce(
                pre_asts: Seq<Option<PreAst>>,
         2
                allocated_asts: Seq<Option<Ast>>,
2272
         3
            ) -> (Seq<Option<PreAst>>, Seq<Option<Ast>>) {
    // Reduce ASTs by handling operators and precedence
         4
2273
          5
2274
          6
                let (pre_asts, allocated_asts) = reduce_by_opr(pre_asts, allocated_asts);
          7
2275
          8
                 // Reduce ASTs by handling delimited constructs like parentheses or braces
2276
          9
                let (pre_asts, allocated_asts) = reduce_by_delimited(pre_asts, allocated_asts);
         10
2277
         11
                // Reduce ASTs by handling function or method calls
2278
         12
                let (pre_asts, allocated_asts) = reduce_by_call(pre_asts, allocated_asts);
         13
2279
         14
                 // Reduce ASTs by handling statements, ensuring proper syntax structure
2280
         15
                let (pre_asts, allocated_asts) = reduce_by_stmt(pre_asts, allocated_asts);
         16
         17
                 // Reduce ASTs by handling definitions, like variables or functions
         18
                let (pre_asts, allocated_asts) = reduce_by_defn(pre_asts, allocated_asts);
         19
2283
         20
                 // Return the final reduced sequences of pre-ASTs and allocated ASTs
         21
                 (pre_asts, allocated_asts)
2284
         22 }
2285
2286
         1 pub fn reduce_n_times(
                mut pre_asts: Seq<Option<PreAst>>,
         2
         3
                mut allocated_asts: Seq<Option<Ast>>,
          4
                n: usize,
          5
            ) -> (Seq<Option<PreAst>>, Seq<Option<Ast>>) {
2289
                for _ in 0..n {
          6
2290
                     let (pre_asts1, allocated_asts1) = reduce(pre_asts, allocated_asts);
          7
                     pre_asts = pre_asts1;
2291
          8
          9
                     allocated_asts = allocated_asts1;
2292
         10
                 (pre asts, allocated asts)
2293
         11
         12 }
```

In the above definition, we actually used Rust's mutable variable semantics. However, it's straightforward to see that it translates to a computation graph that is a sequential composition of subgraphs with sequential length n. Because the AST's depth is bounded by D, we can just take n to be D. Each subgraph is generated from the reduce function, then they are all constant graphs constructed by global and local functions, then by Proposition 13,11 and 2 they translate to transformers with $O(\log L + D)$ depth, model dimension, and number of heads, where $\log L$ comes from the encoding of types like Token.

Below we give full details of the various reduction functions.

As these are implemented as Rust functions, they have been tested against a number of inputs. We don't guarantee an industry level of correctness, but the key point is well illustrated.

F.2 OPERATORS

2294

2306

2307

2309

```
In this section, we lay down the definition of reduce_by_opr.
```

```
2310
           pub(super) fn reduce_by_opr(
2311
         2
               pre_asts: Seq<Option<PreAst>>,
               allocated_asts: Seq<Option<Ast>>,
         3
2312
         4
           ) -> (Seq<Option<PreAst>>, Seq<Option<Ast>>) {
               let pre_asts_nearest_left2 = pre_asts.nearest_left2();
2313
         5
                let pre_asts_nearest_right2 = pre_asts.nearest_right2();
         6
2314
                let new_opr_asts = new_opr_ast.apply(pre_asts_nearest_left2, pre_asts,
         7
2315
                pre asts nearest right2);
                let (pre asts reduced, new parents) = reduce pre asts by opr(pre asts, new opr asts);
         8
2316
         9
                let pre_asts = update_pre_asts_by_new_asts(pre_asts_reduced, new_opr_asts);
2317
               let allocated asts =
        10
        11
                   allocate_asts_and_update_parents(allocated_asts, new_opr_asts, new_parents);
2318
        12
                (pre_asts, allocated_asts)
2319
        13 }
2320
        1 /// a finite function
2321
         2 pub(crate) fn new_opr_ast(
              nearest_left2: Option2<(Idx, PreAst)>,
        3
```

```
2322
                current: Option<PreAst>,
         4
2323
                nearest_right2: Option2<(Idx, PreAst)>,
         5
         6
            ) -> Option<AstData> {
2324
                let Some(PreAst::Opr(opr)) = current else {
         7
2325
         8
                    return None;
2326
         9
                };
         10
                match opr {
2327
                    Opr::Prefix(opr) => {
         11
                        let Some((opd, PreAst::Ast(_))) = nearest_right2.first() else {
         12
2328
         13
                             return None;
2329
         14
                         };
                         if let Some((_, ast)) = nearest_right2.second() {
2330
         15
         16
                             match ast {
2331
                                 PreAst::Keyword(_) => (),
         17
         18
                                 PreAst::Opr(right_opr) => match right_opr {
2332
         19
                                     Opr::Prefix() => (),
2333
                                     Opr::Binary(right_opr) => {
         20
                                         // every binary opr in our small language is left associative,
         21
2334
                 so '<' instead of '<='
2335
         22
                                          if right_opr.precedence() > opr.precedence() {
2336
         23
                                              return None;
         24
                                          }
2337
         25
                                      1
         26
                                     Opr::Suffix(right_opr) => {
2338
                                         if right_opr.precedence() > opr.precedence() {
         27
2339
         28
                                              return None;
         29
                                          }
2340
         30
                                     }
2341
                                 },
         31
2342
         32
                                 PreAst::Ast(_) => (),
                                 // function call or index takes higher precedence
         33
2343
         34
                                 PreAst::LeftDelimiter(_) => return None,
         35
                                 PreAst::RightDelimiter(_) => (),
2344
         36
                                 PreAst::Separator(_) => (),
2345
         37
                             }
                         };
         38
2346
         39
                         Some(AstData::Prefix { opr, opd })
2347
         40
         41
                    Opr::Binary(opr) => {
2348
         42
                         let Some((lopd, PreAst::Ast(_))) = nearest_left2.first() else {
2349
         43
                            return None;
         44
2350
                         };
         45
                         let Some((ropd, PreAst::Ast(_))) = nearest_right2.first() else {
2351
         46
                             return None;
         47
2352
                         };
         48
                         if let Some((_, ast)) = nearest_left2.second() {
2353
         49
                             match ast {
         50
                                 PreAst::Keyword(kw) => (),
2354
         51
                                 PreAst::Opr(left_opr) => match left_opr {
2355
                                     Opr::Prefix(left_opr) => {
         52
         53
                                         if left_opr.precedence() >= opr.precedence() {
2356
         54
                                             return None;
2357
         55
                                          }
         56
2358
         57
                                     Opr::Binary(left_opr) => {
2359
         58
                                          /// every binary opr in our small language is left
                 associative, so '>=' instead of '>'
2360
         59
                                          if left_opr.precedence() >= opr.precedence() {
2361
                                              return None;
         60
         61
                                          }
2362
         62
2363
         63
                                     Opr::Suffix(_) => (), // actually this will be a syntax error
2364
         64
                                 }.
         65
                                 PreAst::Ast( ) => {
2365
                                     if opr != BinaryOpr::LightArrow {
         66
         67
                                          return None;
2366
                                     }
         68
2367
         69
                                 PreAst::LeftDelimiter(_) => (),
2368
         70
         71
                                 PreAst::RightDelimiter(_) => return None,
2369
         72
                                 PreAst::Separator(_) => (),
                             }
         73
2370
         74
                         };
2371
         75
                         if let Some((_, ast)) = nearest_right2.second() {
2372
         76
                             match ast {
                                 PreAst::Keyword(kw) => match kw {
         77
2373
                                     Keyword::ELSE => return None,
         78
2374
         79
                                     _ => (),
         80
2375
         81
                                 PreAst::Opr(right_opr) => match right_opr {
                                     Opr::Prefix(_) => (), // actually this will be a syntax error
         82
```

```
2376
         83
                                      Opr::Binary(right_opr) => {
2377
                 /// every binary opr in our small language is left
associative, so `<` instead of `<=`</pre>
         84
2378
         85
                                          if right_opr.precedence() > opr.precedence() {
         86
                                              return None;
2380
         87
                                          }
         88
2381
         89
                                      Opr::Suffix(right_opr) => {
         90
                                          if right_opr.precedence() >= opr.precedence() {
         91
                                              return None;
2383
         92
                                          }
         93
                                      }
2384
         94
                                  },
2385
                                  // function call or index takes higher precedence
         95
         96
                                  PreAst::LeftDelimiter(_) => return None,
2386
                                  PreAst::RightDelimiter(_) => (),
         97
2387
         98
                                  PreAst::Ast(_) => (),
         99
                                  PreAst::Separator(_) => (),
2388
         100
                             }
2389
         101
                         }:
                         Some(AstData::Binary { lopd, opr, ropd })
2390
         102
         103
2391
                     Opr::Suffix(opr) => {
         104
                         let Some((opd, PreAst::Ast(_))) = nearest_left2.first() else {
2392
        105
         106
                             return None;
2393
         107
                         };
                         if let Some((_, ast)) = nearest_left2.second() {
2394
         108
         109
                             match ast {
2395
         110
                                  PreAst::Keyword(_) => (),
2396
         111
                                  PreAst::Opr(right_opr) => match right_opr {
         112
                                      Opr::Prefix(right_opr) => {
2397
         113
                                          if right_opr.precedence() > opr.precedence() {
         114
                                              return None;
2398
         115
                                          }
2399
         116
         117
                 /// every binary opr in our small language is left
associative, so `<` instead of `<=`</pre>
                                      Opr::Binary(right_opr) => {
2400
         118
2401
         119
                                          if right_opr.precedence() > opr.precedence() {
2402
         120
                                               return None;
2403
         121
         122
2404
         123
                                      Opr::Suffix(_) => (),
2405
         124
                                  },
                                  PreAst::LeftDelimiter(_) => (),
         125
2406
         126
                                  PreAst::RightDelimiter(_) => return None,
2407
         127
                                  PreAst::Ast(_) => return None,
                                  PreAst::Separator(_) => (),
         128
2408
         129
2409
         130
                         };
                         Some(AstData::Suffix { opr, opd })
         131
2410
         132
                     }
2411
         133
                 }
2412
         134 }
2413
         1 /// returns sequence of remaining PreAsts and new parent idxs
2414
         2 pub(crate) fn reduce_pre_asts_by_opr(
          3
                 pre_asts: Seq<Option<PreAst>>,
2415
          4
                 new_asts: Seq<Option<AstData>>,
2416
          5
            ) -> (Seq<Option<PreAst>>, Seq<Option<Idx>>) {
2417
          6
                let new_asts_nearest_left = new_asts.nearest_left();
                 let pre_asts = reduce_pre_ast_by_new_ast.apply(pre_asts, new_asts);
          7
2418
                 let (pre_asts, new_parents) = reduce_pre_ast_by_opr_left
          8
          9
                     .apply_enumerated(new_asts_nearest_left, pre_asts)
2419
         10
                     .decouple();
2420
                 let new_asts_nearest_right = new_asts.nearest_right();
         11
         12
                 reduce_pre_ast_by_opr_right
2421
         13
                     .apply_enumerated(new_asts_nearest_right, pre_asts, new_parents)
2422
                     .decouple()
         14
2423
         15 }
2424
         1 fn reduce_pre_ast_by_new_ast(pre_ast: Option<PreAst>, new_ast: Option<AstData>) ->
2425
                 Option<PreAst> {
         2
                 if new_ast.is_some() {
2426
         3
                    None
2427
          4
                 } else {
2428
         5
                     pre_ast
          6
                 }
2429
         7 }
```

```
2430
        1 fn reduce_pre_ast_by_opr_left(
2431
                idx: Idx,
         2
         3
                new_ast_nearest_left: Option<(Idx, AstData)>,
2432
                pre_ast: Option<PreAst>,
2433
         5
           ) -> (Option<PreAst>, Option<Idx>) {
               let Some(pre_ast) = pre_ast else {
2434
         6
         7
                   return (None, None);
2435
         8
                };
                let Some((new ast idx, new ast data)) = new ast nearest left else {
         9
2436
                    return (Some(pre_ast), None);
         10
2437
        11
               };
               match new_ast_data {
2438
        12
                   AstData::Binary { ropd: opd, .. } | AstData::Prefix { opd, .. } if opd == idx => {
        13
2439
                        (None, Some(new ast idx))
        14
        15
                   }
2440
                    _ => (Some(pre_ast), None),
        16
2441
                }
        17
        18 }
2442
2443
        1 fn reduce_pre_ast_by_opr_right(
2444
         2
               idx: Idx,
2445
         3
                new_ast_nearest_right: Option<(Idx, AstData)>,
                pre_ast: Option<PreAst>,
         4
2446
               new_parent: Option<Idx>,
         5
         6
           ) -> (Option<PreAst>, Option<Idx>) {
2447
               let Some(pre_ast) = pre_ast else {
2448
                    return (None, new_parent);
         8
2449
         9
               };
         10
                if let Some(new_parent) = new_parent {
2450
                   return (None, Some(new_parent));
        11
2451
        12
                let Some((new_ast_idx, new_ast_data)) = new_ast_nearest_right else {
        13
2452
                   return (Some(pre_ast), None);
        14
2453
        15
               };
               match new_ast_data {
        16
2454
                   AstData::Binary { lopd: opd, .. } | AstData::Suffix { opd, .. } if opd == idx => {
        17
2455
        18
                        (None, Some(new_ast_idx))
                   }
        19
2456
                    _ => (Some(pre_ast), None),
        20
        21
                }
2457
        22 }
2458
2459
2460
        F.3 STATEMENTS
2461
2462
        In this section, we lay down the definition of reduce_by_stmt.
2463
        1 pub(super) fn reduce_by_stmt(
2464
         2
                pre_asts: Seq<Option<PreAst>>,
2465
         3
                allocated_asts: Seq<Option<Ast>>,
         4
           ) -> (Seq<Option<PreAst>>, Seq<Option<Ast>>) {
2466
         5
               let pre_asts_nearest_left2 = pre_asts.nearest_left2();
                let pre_asts_nearest_right2 = pre_asts.nearest_right2();
         6
2467
                let new_stmt_asts =
         7
2468
         8
                    new_stmt_ast.apply(pre_asts_nearest_left2, pre_asts, pre_asts_nearest_right2);
               let (pre_asts, new_parents) = reduce_pre_asts_by_stmt(pre_asts, new_stmt_asts);
2469
         9
         10
               let allocated_asts =
2470
        11
                   allocate_asts_and_update_parents(allocated_asts, new_stmt_asts, new_parents);
2471
         12
                let pre_asts = update_pre_asts_by_new_asts(pre_asts, new_stmt_asts);
        13
                (pre_asts, allocated_asts)
2472
        14 }
2473
        1 fn new_stmt_ast(
2474
                pre_ast_nearest_left2: Option2<(Idx, PreAst)>,
         2
2475
         3
                pre_ast: Option<PreAst>,
                pre_ast_nearest_right2: Option2<(Idx, PreAst)>,
2476
         4
         5
           ) -> Option<AstData> {
2477
         6
               let PreAst::Keyword(Keyword::Stmt(kw)) = pre_ast? else {
2478
         7
                   return None;
         8
               };
2479
                match kw {
         9
         10
                    StmtKeyword::Let => {
2480
        11
                        let Some((idx1, PreAst::Ast(ast))) = pre_ast_nearest_right2.first() else {
2481
        12
                            return None;
2482
        13
        14
                        if let Some((_, pre_ast)) = pre_ast_nearest_right2.second() {
2483
        15
                            match pre_ast {
        16
                                PreAst::Keyword(_) => (),
```

```
2484
         17
                                  PreAst::Opr(_) | PreAst::LeftDelimiter(_) => return None,
2485
                                  PreAst::RightDelimiter(_) => (),
         18
         19
                                  PreAst::Ast(_) => return None,
2486
                                  PreAst::Separator(separator) => match separator {
         20
2487
         21
                                       Separator::Comma => return None,
2488
         22
                                       Separator::Semicolon => (),
         23
                                  },
2489
         24
                              }
         25
2490
         26
                          let (pattern, initial_value) = match ast {
2491
         27
                              AstData::Binary {
         28
                                  lopd,
2492
                                  opr: BinaryOpr::Assign,
         29
2493
         30
                                  ropd,
         31
                              } => (lopd, Some(ropd)),
2494
                              AstData::Ident()
         32
2495
                              | AstData::Prefix { .. }
         33
         34
                              | AstData::Binary { .. }
2496
                              | AstData::Delimited { .. }
         35
2497
                              AstData::Call { .. } => (idx1, None),
         36
                              _ => return None,
2498
         37
         38
                          }:
2499
         39
                         Some (AstData::LetInit {
         40
                              expr: idx1,
2500
         41
                              pattern,
2501
         42
                              initial_value,
         43
                         })
2502
         44
2503
         45
                     StmtKevword::If => {
2504
         46
                          let Some((condition, PreAst::Ast(ast1))) = pre_ast_nearest_right2.first() else
                  {
2505
         47
                              return None;
         48
2506
                          };
                          let Some((
         49
2507
         50
                              body,
         51
                              PreAst::Ast(AstData::Delimited {
2508
         52
                                  left_delimiter: LCURL,
2509
         53
                                  right_delimiter: RCURL,
         54
2510
                                  • •
         55
                              }),
2511
         56
                          )) = pre_ast_nearest_right2.second()
         57
                          else {
2512
         58
                              return None;
2513
         59
                          };
         60
                          Some(AstData::If { condition, body })
2514
         61
2515
         62
                     StmtKeyword::Else => {
                          let Some((if_stmt, PreAst::Ast(AstData::If { .. }))) =
2516
         63
                  pre_ast_nearest_left2.first()
2517
         64
                         else {
2518
         65
                              return None;
         66
                          };
2519
         67
                          let Some((
         68
2520
                              body,
         69
                              PreAst::Ast(
2521
         70
                                  AstData::Delimited {
         71
                                      left_delimiter: LCURL,
2522
         72
                                       right_delimiter: RCURL,
2523
         73
                                       . .
         74
2524
                                   }
         75
                                   | AstData::If { .. }
2525
         76
                                   | AstData::Else { .. },
         77
2526
                              ),
         78
                          )) = pre_ast_nearest_right2.first()
2527
         79
                          else {
         80
                              return None;
2528
                  };
if let Some((_, PreAst::Keyword(Keyword::ELSE))) =
pre_ast_nearest_right2.second() {
         81
2529
         82
2530
         83
                              return None;
2531
         84
                          Some(AstData::Else { if_stmt, body })
2532
         85
                     }
         86
2533
         87
                 }
2534
         88
             }
2535
         1 fn reduce_pre_asts_by_stmt(
2536
          2
                 pre_asts: Seq<Option<PreAst>>,
         3
                 new_asts: Seq<Option<AstData>>,
2537
          4
            ) -> (Seq<Option<PreAst>>, Seq<Option<Idx>>) {
          5
               let new_asts_nearest_left = new_asts.nearest_left();
```

```
2538
                let new_asts_nearest_right = new_asts.nearest_right();
         6
2539
                reduce_pre_ast_by_stmt
         7
         8
                    .apply_enumerated(new_asts_nearest_left, new_asts_nearest_right, pre_asts)
2540
         9
                    .decouple()
2541
        10 }
2542
        1 fn reduce_pre_ast_by_stmt(
2543
                idx: Idx,
         2
2544
         3
                new_ast_nearest_left: Option<(Idx, AstData)>,
                new_ast_nearest_right: Option<(Idx, AstData)>,
2545
         4
         5
                pre_ast: Option<PreAst>,
2546
           ) -> (Option<PreAst>, Option<Idx>) {
         6
2547
                if let Some((idx1, ast)) = new_ast_nearest_left {
                   match ast {
2548
                        AstData::LetInit { expr, .. } if expr == idx => (None, Some(idx1)),
         9
2549
         10
                        AstData::If {
         11
                           condition, body, ...
2550
                        } if condition == idx || body == idx => (None, Some(idx1)),
         12
         13
                        AstData::Else { body, .. } if body == idx => (None, Some(idx1)),
2551
         14
                        _ => (pre_ast, None),
2552
         15
                    }
2553
               } else if let Some((idx1, AstData::Else { if_stmt, .. })) = new_ast_nearest_right
         16
         17
                    && if_stmt == idx
2554
         18
                {
         19
                    (None, Some(idx1))
2555
                } else {
         20
2556
         21
                    (pre ast, None)
         22
                }
2557
        23 }
2558
2559
2560
        F.4 GENERALIZED CALL FORMS
2561
2562
        In this section, we lay down the definition of reduce_by_call.
2563
        1
            pub(super) fn reduce_by_call(
2564
                pre_asts: Seq<Option<PreAst>>,
         2
                allocated_asts: Seq<Option<Ast>>,
2565
         3
            ) -> (Seq<Option<PreAst>>, Seq<Option<Ast>>) {
         4
2566
                let pre_asts_nearest_left2 = pre_asts.nearest_left2();
         5
         6
                let pre_asts_nearest_right = pre_asts.nearest_right();
2567
         7
                let new_call_asts =
2568
         8
                    new_call_ast.apply_enumerated(pre_asts_nearest_left2, pre_asts_nearest_right);
2569
         9
                let (pre_asts, new_parents) = reduce_pre_asts_by_call(pre_asts, new_call_asts);
         10
                let allocated asts =
2570
         11
                    allocate_asts_and_update_parents(allocated_asts, new_call_asts, new_parents);
         12
                let pre_asts = update_pre_asts_by_new_asts(pre_asts, new_call_asts);
2571
         13
                (pre_asts, allocated_asts)
2572
        14 }
2573
        1 fn new_call_ast(
2574
                idx: Idx,
         2
2575
                pre_ast_nearest_left2: Option2<(Idx, PreAst)>,
         3
                pre_ast_nearest_right: Option<(Idx, PreAst)>,
2576
         4
           ) -> Option<AstData> {
         5
2577
         6
                let (caller, PreAst::Ast(caller ast)) = pre ast nearest left2.first()? else {
2578
         7
                    return None;
         8
                };
2579
         9
                let (
                    delimited_arguments,
2580
         10
                    PreAst::Ast(AstData::Delimited {
         11
2581
         12
                        left_delimiter_idx,
                        left delimiter,
2582
         13
         14
                        right delimiter,
2583
                    }),
         15
                ) = pre_ast_nearest_right?
2584
         16
         17
                else {
2585
         18
                    return None;
2586
         19
                };
                if let Some((_, snd)) = pre_ast_nearest_left2.second() {
         20
2587
        21
                    match snd {
        22
                        PreAst::Keyword(kw) => match kw {
2588
                            Keyword::Defn(kw) => match kw {
         23
2589
        24
                                 DefnKeyword::Struct | DefnKeyword::Enum => return None,
2590
        25
                                 DefnKeyword::Fn => match left_delimiter.delimiter() {
        26
                                     Delimiter::Parenthesis | Delimiter::Box => return None,
2591
         27
                                     Delimiter::Curly => (),
        28
                                 },
```

```
2592
         29
2593
         30
                             Keyword::Stmt(kw) => match kw {
         31
                                 StmtKeyword::Let => (),
2594
         32
                                 StmtKeyword::If => match left_delimiter.delimiter() {
2595
         33
                                     Delimiter::Parenthesis | Delimiter::Box => (),
                                     Delimiter::Curly => return None,
2596
         34
         35
                                 }.
2597
         36
                                 StmtKeyword::Else => return None,
         37
                             },
2598
         38
                         }.
2599
         39
                         PreAst::Opr(opr) => match opr {
                             Opr::Prefix(_) | Opr::Binary(_) => match left_delimiter.delimiter() {
         40
2600
         41
                                 Delimiter::Parenthesis | Delimiter::Box => (),
2601
                                 Delimiter::Curly => return None,
         42
         43
                             },
2602
         44
                             Opr::Suffix(_) => return None,
2603
         45
                         },
         46
                        PreAst::LeftDelimiter(_) => (),
2604
         47
                         PreAst::RightDelimiter(_) => return None,
2605
         48
                         PreAst::Ast(snd ast) => {
         49
                             if let AstData::Ident(_) = snd_ast
2606
         50
                                 && left delimiter == LCURL
2607
         51
                             ł
         52
                                 match caller_ast {
2608
         53
                                     AstData::Binary {
2609
                                         opr: BinaryOpr::LightArrow,
         54
         55
2610
                                          . .
         56
                                      }
2611
         57
                                      | AstData::Delimited {
2612
         58
                                         left_delimiter: LPAR,
         59
                                         right_delimiter: RPAR,
2613
         60
         61
                                     } => (),
2614
                                     _ => return None,
         62
2615
         63
                                 }
         64
                             } else {
2616
         65
                                 return None;
2617
         66
                             l
2618
         67
         68
                         PreAst::Separator(_) => (),
2619
         69
                    }
         70
2620
         71
                if left_delimiter_idx != idx {
2621
         72
                    return None;
         73
2622
         74
                Some(AstData::Call {
2623
         75
                    caller,
         76
                     delimited_arguments,
2624
         77
                    left_delimiter,
2625
         78
                    right_delimiter,
                })
2626
         79
         80 }
2627
2628
         1 fn reduce_pre_asts_by_call(
               pre_asts: Seq<Option<PreAst>>,
2629
         2
         3
                new_asts: Seq<Option<AstData>>,
2630
         4 ) -> (Seq<Option<PreAst>>, Seq<Option<Idx>>) {
                let new_asts_nearest_left = new_asts.nearest_left();
2631
         5
                let new_asts_nearest_right = new_asts.nearest_right();
         6
2632
         7
                reduce_pre_ast_by_call
2633
         8
                    .apply_enumerated(new_asts_nearest_left, new_asts_nearest_right, pre_asts)
         9
                    .decouple()
2634
         10 }
2635
         1 fn reduce_pre_ast_by_call(
2636
         2
                idx: Idx,
2637
                new_ast_nearest_left: Option<(Idx, AstData)>,
         3
2638
         4
                new_ast_nearest_right: Option<(Idx, AstData)>,
                pre_ast: Option<PreAst>,
         5
2639
            ) -> (Option<PreAst>, Option<Idx>) {
         6
               if let Some((
2640
         7
         8
                    idx1.
2641
         9
                    AstData::Call {
2642
         10
                        delimited_arguments,
         11
                         . .
2643
         12
                    },
2644
         13
                )) = new_ast_nearest_left
         14
                    && delimited_arguments == idx
2645
         15
                {
         16
                     (None, Some(idx1))
```

```
2646
        17
              } else if let Some((idx1, AstData::Call { caller, .. })) = new_ast_nearest_right
2647
         18
                    && caller == idx
         19
                {
2648
        20
                    (None, Some(idx1))
2649
                } else {
         21
2650
         22
                    (pre_ast, None)
                }
         23
2651
        24 }
2652
2653
        F.5 DEFINITIONS
2654
2655
        In this section, we lay down the definition of reduce_by_defn.
2656
2657
         1 pub(super) fn reduce_by_defn(
                pre_asts: Seq<Option<PreAst>>,
         2
2658
                allocated_asts: Seq<Option<Ast>>,
         3
2659
            ) -> (Seq<Option<PreAst>>, Seq<Option<Ast>>) {
         4
                let pre_asts_nearest_left2 = pre_asts.nearest_left2();
2660
         5
                let pre_asts_nearest_right2 = pre_asts.nearest_right2();
         6
2661
         7
                let new_defn_asts =
                    new_defn_ast.apply(pre_asts_nearest_left2, pre_asts, pre_asts_nearest_right2);
2662
         8
         9
                let (pre_asts, new_parents) = reduce_pre_asts_by_defn(pre_asts, new_defn_asts);
2663
         10
                let allocated_asts =
2664
         11
                    allocate_asts_and_update_parents(allocated_asts, new_defn_asts, new_parents);
         12
                let pre_asts = update_pre_asts_by_new_asts(pre_asts, new_defn_asts);
2665
         13
                (pre_asts, allocated_asts)
        14 }
2667
        1 fn new_defn_ast(
2668
                pre_ast_nearest_left2: Option2<(Idx, PreAst)>,
         2
                pre_ast: Option<PreAst>,
         3
2669
                pre_ast_nearest_right2: Option2<(Idx, PreAst)>,
         4
2670
            ) -> Option<AstData> {
         5
                let PreAst::Keyword(Keyword::Defn(keyword)) = pre_ast? else {
2671
         6
                    return None;
         7
2672
         8
                };
2673
         9
                {
                    let Some((ident_idx, PreAst::Ast(AstData::Ident(ident)))) =
         10
2674
                pre_ast_nearest_right2.first()
         11
2675
                    else {
         12
                        return None;
2676
         13
                    };
                    let Some((content, PreAst::Ast(content_ast))) = pre_ast_nearest_right2.second()
2677
         14
                 else {
2678
         15
                        return None;
         16
                    };
2679
                    match keyword {
         17
2680
         18
                        DefnKeyword::Struct => match content_ast {
2681
         19
                            AstData::Delimited { .. } => (),
         20
                            _ => return None,
2682
        21
                        1.
                        DefnKeyword::Enum => match content_ast {
        22
2683
        23
                            AstData::Delimited { .. } => (),
2684
        24
                            _ => return None,
        25
2685
                        }.
        26
                        DefnKeyword::Fn => match content_ast {
2686
        27
                            AstData::Call { .. } => (),
2687
         28
                            _ => return None,
        29
                        },
2688
         30
                    Some(AstData::Defn {
         31
2689
         32
                        keyword,
2690
         33
                        ident,
         34
                        ident_idx,
2691
         35
                        content,
2692
         36
                    })
         37
                }
2693
        38
           }
2694
2695
         1 fn reduce_pre_asts_by_defn(
         2
                pre_asts: Seq<Option<PreAst>>,
2696
         3
                new_asts: Seq<Option<AstData>>,
2697
         4
           ) -> (Seq<Option<PreAst>>, Seq<Option<Idx>>) {
2698
         5
                let new_asts_nearest_left = new_asts.nearest_left();
         6
                let new_asts_nearest_right = new_asts.nearest_right();
2699
                reduce_pre_ast_by_defn
         7
         8
                    .apply_enumerated(new_asts_nearest_left, new_asts_nearest_right, pre_asts)
```

```
9
                     .decouple()
2701
         10 }
2702
2703
         1 fn reduce_pre_ast_by_defn(
         2
                idx: Idx,
2704
         3
                new_ast_nearest_left: Option<(Idx, AstData)>,
          4
                new_ast_nearest_right: Option<(Idx, AstData)>,
2705
          5
                pre_ast: Option<PreAst>,
2706
          6
            ) -> (Option<PreAst>, Option<Idx>) {
                if let Some((idx1, ast)) = new_ast_nearest_left {
2707
          8
                     match ast {
2708
          9
                         AstData::Defn {
2709
         10
                             keyword,
         11
                             ident_idx,
2710
                             ident,
         12
                             content,
2711
         13
         14
2712
                         } if ident_idx == idx || content == idx => (None, Some(idx1)),
         15
                         _ => (pre_ast, None),
         16
2713
         17
                     }
2714
         18
                } else if let Some((idx1, AstData::Defn { .. })) = new_ast_nearest_right
2715
         19
                     && false
         20
                {
2716
         21
                     (None, Some(idx1))
         22
                } else {
2717
         23
                     (pre_ast, None)
2718
         24
                }
         25
2719
            }
```

G TRANSFORMER SYMBOL RESOLUTION PROOF

Here we lay down the code for symbol resolution. The actual process involves many details such as computing ranks (the exact position of an AST node among its siblings), scopes, and roles (a more precise version of AST, computed from its parent recursively), definitions and resolutions.

```
G.1 RANKS
```

```
2729
         1 #[derive(Debug, Default, PartialEq, Eq, Clone, Copy)]
2730
         2
           pub struct Rank(u8);
         3
2731
         4
           impl Rank {
                fn next(self) -> Self {
2732
         5
                    Self(self.0 + 1)
         6
2733
                }
         7
         8 }
2734
2735
         10 pub fn calc_ranks(asts: Seq<Option<Ast>>) -> Seq<Option<Rank>> {
2736
                let counts = asts.count_past_by_attention(asts, |ast, ast1| {
         11
                    let Some(ast) = ast else { return false };
         12
2737
                    let Some(ast1) = ast1 else { return false };
         13
                    ast.parent == ast1.parent
2738
         14
         15
                });
2739
                (|c: usize, ast| {
         16
2740
         17
                    ast?;
                    Some(Rank(c.try_into().unwrap()))
         18
2741
         19
                })
2742
         20
                .apply(counts, asts)
         21 }
2743
         22
         23
           pub fn calc ranks1(asts: Seq<Option<Ast>>, n: usize) -> Seq<Option<Rank>> {
2744
        24
                let mut ranks: Seq<Option<Rank>> = asts.map(|_| None);
2745
        25
                for _ in 0..n {
2746
        26
                    ranks = calc_sibling_indicies_step(asts, ranks);
        27
                }
2747
        28
                ranks
2748
        29
            }
         30
2749
         31 fn calc_sibling_indicies_step(
         32
                asts: Seq<Option<Ast>>
2750
         33
                ranks: Seq<Option<Rank>>,
2751
         34 ) -> Seq<Option<Rank>> {
2752
         35
               let previous_ranks = ranks.nearest_left_filtered_by_attention(asts, asts, |ast, ast1| {
         36
                     let Some(ast) = ast else { return false };
2753
         37
                    let Some(ast1) = ast1 else { return false };
         38
                    ast.parent == ast1.parent
```

```
        39

        2755
        40

        2756
        41

        2757
        43

        2758
        44

        2759
        46

        2760
        47

        2761
        49

        2762
        50
```

```
});
40
       let ranks = (|ast, rank, previous_rank: Option<Option<Rank>>| {
41
           let _ = ast?;
42
           if let Some(rank) = rank {
43
                return Some(rank);
44
45
           let Some(previous_rank) = previous_rank else {
46
               return Some(Default::default());
47
            };
48
           Some (previous_rank?.next())
49
       })
50
       .apply(asts, ranks, previous_ranks);
51
       ranks
52 }
```

In the above, count_past_by_attention that count is representable by transformers by utilizing directly hard attention and the starter token. If the count is c, we shall get c/(c+1) from the attention directly.

```
2768 G.2 SCOPES
```

```
2770
         1 const D: usize = 8usize;
         2
2771
            pub struct Scope {
         3
2772
         4
                enclosing_blocks: BoundedVec<Idx, D>,
         5 }
2773
         6
2774
         7 impl Scope {
               pub fn from_ast(idx: Idx, ast: AstData, parent_scope: Scope) -> Self {
         8
2775
         9
                    match ast {
         10
                       AstData::Delimited
2776
         11
                             left_delimiter_idx,
2777
                             left_delimiter: LCURL,
         12
                             right_delimiter: RCURL,
2778
         13
         14
                        } => Self {
2779
         15
                             enclosing_blocks: parent_scope.enclosing_blocks.append(idx),
2780
         16
                        },
                        _ => parent_scope,
         17
2781
                    }
         18
         19
2782
                }
         20
2783
                pub fn new(idx: Idx) -> Self {
         21
         22
2784
                    Self {
         23
                         enclosing_blocks: todo!(),
2785
         24
                    }
         25
                }
2786
         26
2787
         27
                pub fn append(self, idx: Idx) -> Self {
2788
         28
                    Self {
         29
                        enclosing_blocks: self.enclosing_blocks.append(idx),
2789
         30
                    }
2790
         31
                }
         32
           }
2791
         33
         34 impl Scope {
2792
         35
                pub fn contains(self, other: Self) -> bool {
2793
                    let len = self.enclosing blocks.len();
         36
                    if len > other.enclosing_blocks.len() {
2794
         37
                        return false;
         38
2795
         39
                    1
2796
                    for i in 0..len {
         40
                        if self.enclosing_blocks[i] != other.enclosing_blocks[i] {
         41
2797
         42
                            return false;
         43
                         }
2798
         44
                    }
2799
         45
                    true
2800
         46
                }
         47
           }
2801
         48
2802
         49
           pub fn infer_scopes(asts: Seq<Option<Ast>>, n: usize) -> Seq<Option<Scope>> {
         50
                let mut scopes = initial_scope.apply_enumerated(asts);
2803
         51
                for _ in 0..n {
2804
         52
                     let parent_scopes = parent_queries(asts, scopes);
         53
                    scopes = infer_scopes_step(asts, parent_scopes, scopes);
2805
         54
                }
2806
         55
                scopes
         56
            }
2807
         57
         58 fn initial_scope(idx: Idx, ast: Option<Ast>) -> Option<Scope> {
```

scope: Scope,

```
let ast = ast?;
         59
2809
                if ast.parent.is_some() {
         60
         61
                    return None;
2810
         62
                }
2811
         63
                let scope = Scope::default();
                Some(Scope::from_ast(idx, ast.data, scope))
2812
         64
         65
            }
2813
         66
         67
            fn infer_scopes_step(
2814
                asts: Seq<Option<Ast>>,
         68
2815
         69
                parent_scopes: Seq<Option<Scope>>,
                scopes: Seq<Option<Scope>>,
2816
         70
         71 ) -> Seq<Option<Scope>> {
2817
                infer_scope_step.apply_enumerated(asts, parent_scopes, scopes)
         72
            }
         73
2818
         74
2819
         75
            fn infer_scope_step(
    idx: Idx,
         76
2820
                ast: Option<Ast>,
         77
2821
                parent_scope: Option<Scope>,
         78
                scope: Option<Scope>,
2822
         79
         80
            ) -> Option<Scope> {
2823
         81
                if let Some(scope) = scope {
2824
         82
                     return Some(scope);
         83
2825
         84
                Some(Scope::from_ast(idx, ast?.data, parent_scope?))
         85 }
2826
2827
2828
        G.3 ROLES
2829
         1 #[derive(Debug, Clone, Copy, PartialEq, Eq)]
2830
         2
            pub enum Role {
2831
          3
                LetStmt {
         4
                     pattern: Idx,
2832
          5
                     initial_value: Option<Idx>,
2833
          6
                LetStmtInner {
2834
          7
          8
                    pattern: Idx,
2835
                     initial_value: Idx,
          9
         10
2836
                },
                LetStmtIdent,
         11
2837
         12
                LetStmtTypedVariables {
                    variables: Idx,
2838
         13
         14
                     ty: Idx,
2839
         15
                }.
                StructDefn(Ident),
2840
         16
                EnumDefn(Ident),
         17
2841
                FnDefn(Ident),
         18
2842
                FnDefnCallForm {
         19
         20
                    fn_ident: Ident,
2843
         21
                     scope: Scope,
         22
2844
                },
         23
                FnParameters {
2845
         24
                     fn_ident: Ident,
                     has_return_ty: bool,
2846
         25
         26
                     scope: Scope,
2847
         27
         28
                FnParametersAndReturnType {
2848
                    fn_ident: Ident,
         29
2849
         30
                     parameters: Idx,
         31
                     scope: Scope,
2850
         32
                     return_ty: Idx,
2851
         33
                }.
         34
                FnBody (Ident),
2852
         35
                StructFields (Ident),
2853
         36
                FnParameter {
2854
         37
                    fn ident: Ident,
         38
                     rank: Rank,
2855
         39
                     ty: Idx,
                     fn_ident_idx: Idx,
2856
         40
         41
                     scope: Scope,
2857
         42
                },
2858
         43
                FnParameterIdent {
         44
                     scope: Scope,
2859
         45
                }.
2860
         46
                FnParameterSeparated {
         47
                    fn_ident: Ident,
2861
         48
                     rank: Rank,
```

```
2862
        50
2863
        51
                FnParameterType {
         52
                    fn_ident: Ident,
2864
        53
                    rank: Rank,
2865
         54
                },
         55
2866
                FnOutputType {
         56
                    fn_ident: Ident,
2867
         57
                }.
         58
                StructField {
2868
                    ty_ident: Ident,
         59
2869
         60
                    field_ident: Ident,
                    ty_idx: Idx,
         61
2870
         62
                1.
2871
                StructFieldType {
         63
         64
                    ty_ident: Ident,
2872
                    field_ident: Ident,
         65
2873
         66
                },
                TypeArgument,
         67
2874
                TypeArguments,
         68
2875
                StructFieldSeparated(Ident),
         69
                LetStmtVariablesType,
2876
         70
         71
                LetStmtVariables,
2877
        72 }
2878
2879
        1 impl Ast {
               fn role(self) -> Option<Role> {
         2
2880
         3
                   match self.data {
                        AstData::LetInit {
2881
         4
         5
                            expr,
2882
         6
                            pattern,
         7
                            initial_value,
2883
                        } => Some(Role::LetStmt {
         8
2884
         9
                            pattern,
                            initial_value,
         10
2885
         11
                        }),
2886
                        AstData::Defn {
         12
2887
         13
                             keyword,
         14
                             ident idx.
2888
         15
                             ident,
2889
         16
                            content.
                        } => Some(match keyword {
         17
2890
                            DefnKeyword::Struct => Role::StructDefn(ident),
         18
                             DefnKeyword::Enum => Role::EnumDefn(ident),
         19
2891
                            DefnKeyword::Fn => Role::FnDefn(ident),
        20
2892
                        }),
         21
                        _ => None,
        22
2893
        23
                    }
2894
        24
                }
        25 }
2895
2896
         1 pub fn calc_roles(
2897
               asts: Seq<Option<Ast>>,
         2
2898
         3
                scopes: Seq<Option<Scope>>,
                n: usize,
         4
2899
         5 ) -> Seq<Option<Role>> {
               let mut roles: Seq<Option<Role>> = asts.map(|ast| ast?.role());
2900
         6
                let ranks = calc_ranks(asts);
         7
2901
                for _ in 0..n {
         8
         9
                    let parent_roles = parent_queries(asts, roles);
2902
         10
                    roles = calc_roles_step(asts, parent_roles, roles, ranks, scopes);
2903
         11
                }
2904
                roles
         12
        13 }
2905
2906
        1 fn calc_roles_step(
               asts: Seq<Option<Ast>>,
2907
         2
         3
                parent_roles: Seq<Option<Role>>,
2908
         4
                roles: Seq<Option<Role>>,
2909
         5
                ranks: Seq<Option<Rank>>,
                scopes: Seq<Option<Scope>>,
         6
2910
           ) -> Seq<Option<Role>> {
         7
                calc_role_step.apply_enumerated(asts, parent_roles, roles, ranks, scopes)
2911
         8
         9
            }
2912
2913
         1 fn calc_role_step(
2914
         2
                idx: Idx,
         3
                ast: Option<Ast>,
2915
         4
                parent_role: Option<Role>,
         5
                role: Option<Role>,
```

```
2916
                 rank: Option<Rank>,
         6
2917
                scope: Option<Scope>,
              -> Option<Role> {
          8
            )
2918
                if let Some(role) = role {
          9
2919
         10
                     return Some(role);
2920
         11
         12
                 let ast = ast?;
2921
                if let Some(role) = ast.role() {
         13
                     return Some(role);
         14
         15
2923
                match parent_role? {
         16
                     Role::LetStmt {
2924
         17
         18
                         pattern,
2925
         19
                         initial value,
         20
                     } => match ast.data {
2926
         21
                         AstData::Ident(ident) if idx == pattern => Some(Role::LetStmtIdent),
2927
                         AstData::Binary {
         22
         23
                              lopd,
2928
                              opr: BinaryOpr::Assign,
         24
2929
         25
                              ropd,
         26
                              lopd_ident,
2930
                         } if lopd == pattern => Some(Role::LetStmtInner {
         27
2931
         28
                              pattern,
         29
2932
                              initial_value: ropd,
         30
                         }),
2933
                         _ => None,
         31
2934
         32
                     }.
         33
                     Role::LetStmtInner {
2935
         34
                         pattern,
2936
         35
                         initial_value,
         36
                     } => {
2937
         37
                         if idx == pattern {
         38
                             match ast.data {
2938
                                  AstData::Ident(ident) => Some(Role::LetStmtIdent),
         39
2939
         40
                                  AstData::Binary {
         41
                                      lopd,
2940
         42
                                      lopd_ident,
2941
         43
                                      opr,
         44
                                      ropd,
2942
         45
                                  } => Some(Role::LetStmtTypedVariables {
2943
         46
                                      variables: lopd,
         47
                                      ty: ropd,
2944
         48
                                  }),
2945
         49
                                  _ => todo!(),
         50
                              }
2946
         51
                         } else {
2947
         52
                             None
         53
                         }
2948
         54
2949
         55
                     Role::LetStmtIdent => todo!(),
                     Role::FnParameterIdent { scope } => todo!(),
2950
         56
         57
                     Role::StructDefn(ident) => match ast.data {
2951
         58
                         AstData::Literal(_) => todo!(),
         59
                         AstData::Ident(_) => None,
2952
         60
                         AstData::Prefix { opr, opd } => todo!(),
2953
         61
                         AstData::Binary {
         62
                             lopd,
2954
         63
                              opr,
2955
         64
                              ropd,
         65
                              lopd_ident,
2956
                         } => todo!(),
         66
2957
         67
                         AstData::Suffix { opd, opr } => todo!(),
                         AstData::Delimited {
2958
         68
         69
                              left_delimiter_idx,
2959
         70
                              left delimiter,
         71
                             right_delimiter,
2960
                         } => Some(Role::StructFields(ident)),
         72
2961
                         AstData::SeparatedItem { content, separator } => todo!(),
         73
         74
                         AstData::Call { .. } => todo!(),
2962
         75
                         AstData::LetInit {
2963
         76
77
                             expr,
2964
                              pattern,
         78
79
                             initial value,
2965
                         } => todo!(),
                         AstData::Return { result } => todo!(),
2966
         80
                         AstData::Assert { condition } => todo!(),
         81
2967
                         AstData:: If { condition, body } => todo!(),
         82
2968
         83
                         AstData::Else { if_stmt, body } => todo!(),
         84
                         AstData::Defn {
2969
         85
                              keyword,
         86
                              ident_idx,
```

2970					
2971	87 88	ident, content.			
2972	89	} => todo!(),			
2973	90 91	$\},$			
2974	92	Role::FnDefn(fn_ident) => match ast.data {			
2975	93 94	AstData::Literal(_) => todo!(),			
2976	94 95	ASTDATA::IGENT(_) => NONE, AstData::Prefix { opr, opd } => todo!(),			
2977	96 07	AstData::Binary {			
2978	98	opr,			
2979	99 100	ropd, lond ident			
2980	101	} => todo!(),			
2981	102	AstData::Suffix { opd, opr } => todo!(),			
2982	103	left_delimiter_idx,			
2983	105	left_delimiter,			
2984	106	<pre>right_definiter, } => todo!(),</pre>			
2985	108	<pre>AstData::SeparatedItem { content, separator } => todo!(), AstData::Call (</pre>			
2986	1109	delimited_arguments,			
2987	111				
2988	112	<pre>} => Some(Role::FnDeInCallForm { fn_ident,</pre>			
2989	114	<pre>scope: match scope {</pre>			
2990	115	<pre>Some(scope) => scope.append(delimited_arguments), None => Scope::new(delimited_arguments),</pre>			
2991	117	},			
2992	118	}), AstData::LetInit {			
2993	120	expr,			
2994	121	initial_value,			
2995	123	<pre>} => todo!(), </pre>			
2996	124	AstData::Assert { condition } => todo:(),			
2997	126	AstData::If { condition, body } => todo!(),			
2998	127	AstData::Else { 11_stmt, body } => todo!(), AstData::Defn {			
2999	129	keyword,			
3000	130	ident,			
3001	132	content,			
3002	133	} => coao:(), },			
3003	135	<pre>Role::FnDefnCallForm { fn_ident, scope } => match ast.data {</pre>			
3004	130	AstData::Itteral(_) => todo:(), AstData::Ident(_) => todo!(),			
3005	138	AstData::Prefix { opr, opd } => todo!(),			
3006	140	lopd,			
3007	141	opr,			
3008	142	lopd_ident,			
3009	144	} => { if opr == BinaryOpr::LightArrow {			
3010	146	Some (Role::FnParametersAndReturnType {			
3011	147 148	fn_ident, parameters: lond.			
3012	149	return_ty: ropd,			
3013	150	scope,			
3014	152	} else {			
3015	153 154	unreachable!()			
3016	155	}			
3017	156 157	AstData::Suffix { opd, opr } => todo!(), AstData::Delimited {			
3018	158	left_delimiter_idx,			
3019	159 160	left_delimiter, right delimiter.			
3020	161	<pre>} => match left_delimiter.delimiter() {</pre>			
3021	162 163	<pre>Delimiter::Parenthesis => Some(Role::FnParameters { fn ident.</pre>			
3022	164	has_return_ty: false,			
3023	165 166	scope,			
	167	<pre>Delimiter::Box => todo!(),</pre>			

3024				
3025	168	<pre>Delimiter::Curly => Some(Role::FnBody(fn_ident)),</pre>		
3026	169	}, AstData::SeparatedItem { content, separator } => todo!(),		
2027	171	AstData::Call { } => todo!(),		
2020	172	AstData::LetInit {		
3020	174	expr, pattern,		
3029	175	initial_value,		
3030	176	} => todo!(), AstData::Return { result } => todo!(),		
3031	178	AstData::Assert { condition } => todo!(),		
3032	179	AstData::If { condition, body } => todo!(),		
3033	181	AstData::Defn {		
3034	182	keyword,		
3035	185	ident,		
3036	185	content,		
3037	186 187	} => todo!(),		
3038	188	Role::FnParameters {		
3039	189	fn_ident, scope,		
3040	190	AstData::Binary {		
3041	192	lopd,		
30/12	193 194	opr, ropd.		
2042	195	lopd_ident,		
2043	196	} => { if opr == BiparyOpr::TupeIs {		
3044	198	Some (Role::FnParameter {		
3045	199	fn_ident,		
3046	200	rank: rank.unwrap(),		
3047	202	ty: ropd,		
3048	203 204	scope,		
3049	205	} else {		
3050	206	unreachable!()		
3051	207	}		
3052	209	<pre>AstData::SeparatedItem { } => Some(Role::FnParameterSeparated {</pre>		
3053	210	<pre>rank: rank.unwrap(),</pre>		
3054	212	scope,		
3055	213 214	<pre>}), => unreachable!().</pre>		
3056	215	},		
3057	216	Role::FnBody(_) => None, Role::StructFields(ty_ident) => match ast data {		
3058	218	AstData::Binary {		
3059	219	lopd,		
3060	220	ropd,		
3061	222	lopd_ident,		
3062	223 224	<pre>} => { assert_eq!(opr, BinaryOpr::TypeIs);</pre>		
3063	225	Some(Role::StructField {		
306/	226 227	ty_ident, field ident: lopd ident.unwrap().		
2065	228	ty_idx: ropd,		
2066	229	})		
3000	230	} AstData::SeparatedItem { content, separator } => {		
3067	232	Some(Role::StructFieldSeparated(ty_ident))		
3068	233 234	_ => None,		
3069	235	},		
3070	236 237	Role::FnParameter { fn ident.		
3071	238	fn_ident_idx,		
3072	239 240	rank,		
3073	240	scope,		
3074	242			
3075	243 244	<pre>} => { if idx == ty {</pre>		
3076	245	<pre>Some(Role::FnParameterType { fn_ident, rank })</pre>		
3077	246 247	<pre>} else if idx == fn_ident_idx { Some(Role::EnParameterIdent { scope })</pre>		
	248	<pre>} else {</pre>		

3078					
3079	249 250	None			
3080	251	}			
3081	252	Role::FnParameterSeparated {			
3082	253	rank,			
3083	255	scope,			
3084	256 257	} => match ast.data { AstData::Binary {			
3085	258	lopd,			
3086	259	opr, ropd,			
3087	261	lopd_ident,			
3088	262 263	<pre>} => { if opr == BinaryOpr::TypeIs {</pre>			
3089	264	Some(Role::FnParameter {			
3090	265 266	fn_ident, fn ident idx: lopd,			
3091	267	rank,			
3092	268 269	ty: ropd, scope,			
3003	270	})			
300/	271 272	<pre>} else { unreachable!()</pre>			
3005	273	}			
3006	274 275	<pre>} => unreachable!().</pre>			
2007	276	},			
2008	277 278	Role::StructField {			
2000	279	field_ident,			
2100	280 281	ty_idx ,			
2101	282	if idx == ty_idx {			
2102	283 284	Some(Role::StructFieldType {			
3102	285	field_ident,			
210/	286 287	}) } else {			
2105	288	None			
2106	289	}			
2107	291	, Role::StructFieldSeparated(ty_ident) => match ast.data {			
2100	292	AstData::Binary {			
2100	294	opr,			
2110	295 296	ropd, lopd ident.			
2111	297	} => {			
2112	298	assert_eq!(opr, BinaryOpr::TypeIs); Some(Bole::StructField {			
2112	300	ty_ident,			
311/	301 302	<pre>field_ident: lopd_ident.unwrap(), ty idx: ropd.</pre>			
3115	303	<pre>})</pre>			
2116	304 305	} => unreachable!().			
3117	306	},			
3118	307	<pre>Role::FnParameterType { } Role::StructFieldType { } Role::TypeArgument => {</pre>			
3110	308	<pre>match ast.data {</pre>			
3120	309 310	AstData::Delimited { left delimiter idx.			
3121	311	left_delimiter,			
3122	312 313	right_delimiter, } => Some(Role::TypeArguments),			
3122	314	_ => None,			
3123	315 316	}			
2125	317	Role::TypeArguments => match ast.data {			
2125	318 319	AstData::Ident(_) => Some(Role::TypeArgument), AstData::Delimited {			
3120	320	left_delimiter_idx,			
3122	321 322	left_delimiter, right_delimiter.			
3120	323	} => todo!(),			
3120	324 325	<pre>AstData::SeparatedItem { content, separator } => todo!(), AstData::Call {</pre>			
2121	326	caller,			
0101	327 328	caller_ident, left delimiter.			
	520	,			

```
3132
        329
                             right_delimiter,
3133
                             delimited_arguments,
        330
        331
                         } => todo!(),
3134
                        _ => None,
        332
3135
        333
                     }.
                    Role::FnParametersAndReturnType {
3136
        334
        335
                        fn_ident,
3137
                        parameters,
        336
        337
                         return_ty,
3138
        338
                        scope,
3139
        339
                     } => {
                        if idx == parameters {
3140
        340
        341
                             Some (Role::FnParameters {
3141
        342
                                 fn ident,
3142
        343
                                 has return ty: true,
        344
                                 scope,
3143
                             })
        345
                         } else if idx == return_ty {
        346
3144
                            Some(Role::FnOutputType { fn_ident })
        347
3145
        348
                         } else {
                             unreachable!()
3146
        349
                         }
        350
3147
        351
                    }
                    Role::FnOutputType { fn_ident } => todo!(),
3148
        352
                    Role::LetStmtTypedVariables { variables, ty } => {
        353
3149
                        if idx == variables {
        354
                            Some (Role::LetStmtVariables)
3150
        355
        356
                         } else if idx == ty {
3151
                            Some (Role::LetStmtVariablesType)
        357
3152
        358
                         } else {
        359
                             unreachable!()
3153
        360
                         }
        361
3154
                    Role::LetStmtVariablesType => todo!(),
        362
3155
        363
                    Role::LetStmtVariables => todo!(),
        364
3156
                 }
        365
            }
3157
3158
3159
        G.4 DEFNS
3160
         1 #[derive(Debug, Clone, Copy, PartialEq, Eq)]
3161
         2 pub struct SymbolDefn {
3162
                pub symbol: Symbol,
         3
         4
                pub scope: Option<Scope>,
3163
         5 }
3164
3165
         1 pub fn calc_symbol_defns(
         2
                asts: Seq<Option<Ast>>,
3166
                scopes: Seq<Option<Scope>>,
         3
3167
         4
                n: usize,
         5 ) -> Seq<Option<SymbolDefn>> {
3168
                let roles = calc_roles(asts, scopes, n);
         6
3169
                calc_symbol_defn.apply_enumerated(asts, roles, scopes)
         8 }
3170
3171
         1 fn calc_symbol_defn(
                idx: Idx,
3172
         2
         3
                ast: Option<Ast>,
3173
                role: Option<Role>,
         4
                scope: Option<Scope>,
3174
         5
         6 ) -> Option<SymbolDefn> {
3175
                match ast?.data {
         7
                    AstData::Ident(ident) => match role? {
         8
3176
                         Role::LetStmt { .. } => unreachable!(),
         9
3177
                         Role::LetStmtVariables | Role::LetStmtIdent => Some(SymbolDefn {
         10
3178
         11
                             symbol: Symbol {
         12
                                 ident,
3179
         13
                                 source: idx,
                                 data: SymbolData::Variable,
3180
         14
                             },
         15
3181
         16
                             scope,
3182
         17
                         }),
         18
                         Role::FnParameterIdent { scope } => Some(SymbolDefn {
3183
         19
                             symbol: Symbol {
3184
         20
                                 ident,
         21
                                 source: idx,
3185
         22
                                 data: SymbolData::Variable,
         23
                             }.
```

```
3186
         24
                             scope: Some(scope),
3187
         25
                        }),
                        _ => None,
         26
3188
         27
                     }.
3189
         28
                    AstData::Defn {
         29
                        keyword,
3190
         30
                         ident_idx,
3191
         31
                        ident,
         32
                        content,
3192
         33
                     } => Some(SymbolDefn {
3193
         34
                        symbol: Symbol {
3194
         35
                            ident,
         36
                             source: idx,
3195
                             data: SymbolData::Item {
         37
         38
                                 kind: keyword.into(),
3196
         39
                             }.
3197
         40
                         },
         41
                         scope,
3198
         42
                    }),
3199
         43
                      => None.
         44
3200
         45
            }
3201
        G.5
               RESOLUTIONS
3204
         1 pub enum SymbolResolution {
3205
                Ok(Symbol),
         2
                Err(SymbolResolutionError),
3206
         3
         4
            }
3207
         1 pub enum SymbolResolutionError {
3209
         2
                NotResolved,
         3
                NotYetDeclared(Symbol),
3210
         4
            }
3211
3212
         1 pub fn calc_symbol_resolutions(asts: Seq<Option<Ast>>, n: usize) ->
                 Seq<Option<SymbolResolution>> {
3213
         2
                let scopes = infer_scopes(asts, n);
3214
          3
                let symbol_defns = calc_symbol_defns(asts, scopes, n);
          4
                let idents = asts.map(|ast| match ast?.data {
3215
          5
                    AstData::Ident(ident) => Some(ident),
3216
                    _ => None,
          6
          7
                });
3217
                let symbols = symbol_defns
          8
3218
          9
                    .map(|symbol_defn| Some(symbol_defn?.symbol))
         10
                    .first_filtered_by_attention(
3219
         11
                         (|ident, scope| (ident, scope)).apply(idents, scopes),
3220
                         symbol_defns,
         12
         13
                         (ident, scope), symbol_defn| {
3221
         14
                             let Some(ident) = ident else { return false };
3222
         15
                             let Some(symbol_defn) = symbol_defn else {
         16
                                 return false;
3223
         17
                             };
3224
                             if let Some(symbol_defn_scope) = symbol_defn.scope {
         18
                                 if !symbol_defn_scope.contains(scope.unwrap()) {
         19
         20
                                     return false;
         21
         22
3227
         23
                             symbol_defn.symbol.ident == ident
3228
         24
                         },
         25
                    )
3229
         26
                     .map(|s| s.flatten());
3230
         27
                finalize.apply_enumerated(idents, symbols)
         28 }
3231
3232
```

In the above code, we use a somehow complicated attention which we should illustrate why it's representable by transformers. The essence is to prove symbol_defn_scope.contains(scope.unwrap()) can be represented as part of the inner product in $Q^{\top}K$. This can be done by looking closer to what contains does. Consider two scopes, scope1 and scope2, which are sequences of bracket ast indices (can be null). The function returns true if the sequence of scope1 contains the sequence of scope2 as prefix, which can be achieved by $\sum_i x_i^{\top} y_i$ where x_i, y_i are the encoding of *i*th ast indices of scope1 and scope2 after some transformations (different transformations because the function 3240 is asymmetric) so that $x_i^{\top} y_i = 0$ if and only if either x_i is a None or x_i represents the same thing 3241 as y_i , and $x_i^{\top}y_i < 0$ otherwise. More concretely, if x_i is a None, $x_i = 0$ by choice, and equal to 3242 $(1, u_i)$ otherwise where u_i corresponds to the encoding of the *i*th ast index of scope1; if y_i is a 3243 None, $y_i = \mathbf{0}$ by choice, and equal to $(-1, v_i)$ otherwise where A > 0 and v_i corresponds to the 3244 encoding of the *i*th ast index of scope2. We should choose the encoding u_i, v_i such that $u_i^{\dagger} v_i = 1$ 3245 if and only if they encode the same index, which is obviously easy enough. 3246

```
1 fn finalize(idx: Idx, ident: Option<Ident>, symbol: Option<Symbol>) ->
                 Option<SymbolResolution> {
         2
                let
                     = ident?;
3249
         3
                let Some(symbol) = symbol else {
         4
                    return Some(SymbolResolution::Err(SymbolResolutionError::NotResolved));
         5
                };
                match symbol.data {
         6
                    SymbolData::Item { .. } => (),
         7
         8
                    SymbolData::Variable => {
         9
                        if idx < symbol.source {</pre>
         10
                             return Some(SymbolResolution::Err(
         11
                                 SymbolResolutionError::NotYetDeclared(symbol),
         12
                             ));
         13
                        }
3256
         14
                    }
         15
         16
                Some(SymbolResolution::Ok(symbol))
        17 }
```

3258 3259

3262

3266

3247

3248

3250

3251

3252

3253

3254

3255

3257

TRANSFORMER TYPE CHECKING PROOF Η

Here we lay down the code for type analysis. It should be noted that we didn't completely implement 3264 all the details. Things like struct fields, enum variant fields are left out. However, we already cover 3265 the essential mechanism of type analysis, making it sufficient for proof purposes.

```
3267
      H.1 TYPE SIGNATURES
```

```
3268
        1 #[deri
3269
           ve(Debug, PartialEq, Eq, Clone, Copy)]
         2
3270
         3 pub struct TypeSignature {
              pub key: TypeSignatureKey,
         4
         5
                pub ty: Type,
3272
        6 }
3273
        1 #[derive(Debug, PartialEq, Eq, Clone, Copy)]
3274
         2 pub enum TypeSignatureKey {
3275
         3
                FnParameter { fn_ident: Ident, rank: Rank },
         4
                FnOutput { fn_ident: Ident },
3276
         5
                StructField { ty_ident: Ident, field_ident: Ident },
3277
        6 }
3278
        1 pub(super) fn calc_ty_signatures(
3279
               asts: Seq<Option<Ast>>
         2
3280
         3
                roles: Seq<Option<Role>>,
         4
                ty_terms: Seq<Option<Type>>,
         5 ) -> Seq<Option<TypeSignature>> +
               calc_ty_signature.apply(roles, ty_terms)
         6
        7 }
3283
3284
        1 fn calc_ty_signature(role: Option<Role>, ty_term: Option<Type>) -> Option<TypeSignature> {
3285
         2
                let key = match role? {
3286
         3
                    Role::FnParameterType { fn_ident, rank } => {
         4
                        TypeSignatureKey::FnParameter { fn_ident, rank }
3287
         5
3288
         6
                    Role::StructFieldType
         7
                        ty_ident,
3289
         8
                        field_ident,
         9
                    } => TypeSignatureKey::StructField {
3290
         10
                        ty_ident,
3291
        11
                        field_ident,
         12
        13
                    Role::FnOutputType { fn_ident } => TypeSignatureKey::FnOutput { fn_ident },
3293
        14
                    Role::FnParameters {
        15
                        fn_ident,
```

```
3294
        16
                        has_return_ty: false,
3295
        17
                        scope,
         18
                    } => {
3296
                        let key = TypeSignatureKey::FnOutput { fn_ident };
         19
3297
         20
                        let ty = Type::new_ident(Ident::new("unit"));
3298
         21
                        return Some(TypeSignature { key, ty });
         22
                    }
3299
         23
                    _ => return None,
         24
                };
3300
         25
                // put it here!
3301
         26
                let ty = ty_term?;
                Some(TypeSignature { key, ty })
3302
         27
        28 }
3303
3304
3305
        H.2 TYPE INFERENCE
3306
        1 pub struct TypeInference {
3307
         2
               pub ty: Type,
3308
        3 }
3309
         1 pub fn calc_ty_inferences(
3310
         2
                asts: Seq<Option<Ast>>,
3311
                symbol_resolutions: Seq<Option<SymbolResolution>>,
         3
         4
                roles: Seq<Option<Role>>,
3312
         5
                ty_terms: Seq<Option<Type>>,
3313
         6
                ty_signatures: Seq<Option<TypeSignature>>,
         7
               n: usize,
3314
         8 ) -> Seq<Option<TypeInference>> {
3315
         9
               let mut ty_inferences = infer_tys_initial(asts, ty_signatures);
         10
                let mut ty_designations =
3316
                   calc_initial_ty_designations(asts, roles, symbol_resolutions, ty_inferences,
         11
3317
                 ty_terms);
         12
                for _ in 0..n {
3318
         13
                   ty_inferences |= infer_tys_step(asts, symbol_resolutions, ty_inferences,
3319
                 ty_designations);
         14
                    ty_designations |= calc_ty_designations_step(roles, symbol_resolutions,
3320
                 ty_inferences);
3321
         15
                ty_inferences
         16
3322
        17 }
3323
3324
         1 fn infer_tys_initial(
                asts: Seg<Option<Ast>>,
         2
3325
                ty_signatures: Seq<Option<TypeSignature>>,
         3
3326
         4 ) -> Seq<Option<TypeInference>> {
               inference_literal_tys(asts).or(infer_fn_call_tys(asts, ty_signatures))
         5
3327
         6 }
3328
3329
         1 fn inference_literal_tys(asts: Seq<Option<Ast>>) -> Seq<Option<TypeInference>> {
         2
               asts.map(|ast| match ast?.data {
3330
         3
                   AstData::Literal(lit) => match lit {
                       Literal::Int(_) => Some(TypeInference {
3331
         4
                            ty: Type::new_ident(Ident::new("Int")),
         5
3332
         6
                        }),
                        Literal::Float(_) => Some(TypeInference {
3333
         7
                            ty: Type::new_ident(Ident::new("Float")),
         8
3334
         9
                        }),
3335
         10
                    },
                    _ => None,
         11
3336
               })
         12
        13 }
3337
3338
         1 fn infer_fn_call_tys(
3339
                asts: Seq<Option<Ast>>,
         2
3340
         3
                ty_signatures: Seq<Option<TypeSignature>>,
         4
            ) -> Seq<Option<TypeInference>> {
3341
         5
                ty_signatures
3342
         6
                    .first_filtered_by_attention(asts, ty_signatures, |ast, ty_signature| {
         7
                        let Some(ast) = ast else { return false };
3343
                        let Some(TypeSignature {
         8
3344
         9
                            key: TypeSignatureKey::FnOutput { fn_ident },
         10
3345
                        }) = ty_signature
         11
3346
         12
                        else {
         13
                            return false;
3347
         14
                        };
        15
                        match ast.data {
```

```
3348
        16
                             AstData::Call {
3349
        17
                                 caller,
         18
                                 caller_ident,
3350
         19
                                 left_delimiter,
3351
         20
                                 right_delimiter,
3352
         21
                                 delimited_arguments,
         22
                             } if caller_ident == Some(fn_ident) => true,
3353
         23
                             _ => false,
         24
                         }
3354
         25
                    })
3355
         26
                    .map(|ty_inference| {
         27
                        Some(TypeInference {
3356
                            ty: ty_inference??.ty,
         28
3357
         29
                         })
         30
                    })
3358
        31 }
3359
3360
3361
        H.3 TYPE EXPECTATIONS
3362
3363
        1 pub struct TypeExpectation {
         2
                pub ty: Type,
3364
         3
                pub source: TypeExpectationSource,
        4 }
3365
3366
         1 pub enum TypeExpectationSource {
3367
                CallArgument { caller_ident: Ident, rank: Rank },
         2
        3 }
3368
3369
        1 pub fn calc_ty_expectations(
3370
                asts: Seq<Option<Ast>>,
         2
                ranks: Seq<Option<Rank>>,
         3
3371
                ty_signatures: Seq<Option<TypeSignature>>,
         4
3372
            ) -> Seq<Option<TypeExpectation>> {
         5
                let parent_asts = asts.index(asts.map(|ast| ast?.parent)).map(Option::flatten);
         6
3373
                let grandparent_asts = asts
         7
3374
                     .index(parent_asts.map(|parent_ast| parent_ast?.parent))
         8
                     .map(Option::flatten);
         9
3375
                let ty_expectation_sources = calc_ty_expectation_source.apply(grandparent_asts, ranks);
         10
3376
                let retrieved_ty_signatures = ty_signatures
         11
                    .first_filtered_by_attention(
3377
         12
                         ty expectation sources,
         13
3378
         14
                         ty_signatures,
         15
                         |ty_expection_source, ty_signature| {
3379
                             let Some(type_expectation_source) = ty_expection_source else {
         16
3380
         17
                                 return false;
3381
         18
                             };
                             let Some(type_signature) = ty_signature else {
         19
3382
         20
                                 return false;
        21
3383
                             };
         22
                             match (type_expectation_source, type_signature.key()) {
3384
        23
3385
        24
                                     TypeExpectationSource::CallArgument {
        25
                                         caller_ident,
3386
        26
                                         rank: rank0,
        27
3387
         28
                                     TypeSignatureKey::FnParameter {
3388
        29
                                         fn ident,
3389
         30
                                         rank: rank1,
         31
                                     },
3390
         32
                                 ) if caller_ident == fn_ident && rank0 == rank1 => true,
         33
                                   => false,
3391
         34
                            }
3392
         35
                         },
         36
                    )
3393
         37
                     .map(Option::flatten);
3394
         38
                 (|ty_expectation_source: Option<TypeExpectationSource>,
3395
         39
                  retrieved_ty_signature: Option<TypeSignature>| {
         40
                    Some(TypeExpectation {
3396
         41
                         ty: retrieved_ty_signature?.ty(),
         42
                         source: ty_expectation_source?,
3397
         43
                    })
3398
         44
                })
         45
                .apply(ty_expectation_sources, retrieved_ty_signatures)
3399
        46 }
3400
3401
         1 fn calc_ty_expectation_source(
              grandparent_ast: Option<Ast>,
         2
```

```
3402
        3
             rank: Option<Rank>,
3403
        4 ) -> Option<TypeExpectationSource> {
        5
              let grandparent_ast = grandparent_ast?;
3404
               let rank = rank?;
         6
3405
              match grandparent_ast.data {
         7
         8
                 AstData::Call {
3406
         9
                    caller,
3407
                      caller_ident: Some(caller_ident),
        10
                      left delimiter,
        11
3408
                     right_delimiter,
        12
3409
        13
                      delimited arguments,
                  } => Some(TypeExpectationSource::CallArgument { caller_ident, rank }),
3410
        14
                  _ => None,
        15
3411
              }
        16
        17 }
3412
3413
3414
        H.4 TYPE ERRORS
3415
3416
        1 pub enum TypeError {
3417
        2
               TypeMismatch { expected: Type, actual: Type },
        3 }
3418
3419
        1 pub fn calc_ty_errors(
3420
        2
               ty_inferences: Seq<Option<TypeInference>>,
               ty_expectations: Seq<Option<TypeExpectation>>,
        3
3421
        4 ) -> Seq<Option<TypeError>> {
3422
              calc_ty_error.apply(ty_inferences, ty_expectations)
         5
        6 }
3423
3424
        1 fn calc_ty_error(
3425
               ty_inference: Option<TypeInference>,
         2
               ty_expectation: Option<TypeExpectation>,
        3
3426
         4 ) -> Option<TypeError> {
3427
               let ty_inference = ty_inference?;
        5
               let ty_expectation = ty_expectation?;
3428
         6
               if ty_inference.ty == ty_expectation.ty {
         7
3429
         8
                   None
3430
         9
               } else {
        10
                  Some(TypeError::TypeMismatch {
3431
        11
                      expected: ty_expectation.ty,
3432
        12
                       actual: ty_inference.ty,
                  })
        13
3433
               }
        14
        15 }
3434
3435
3436
        I LOWER BOUNDS
3437
3438
3439
        1 struct <ty-ident-1> {}
        2 struct <ty-ident-2> {}
3440
        3 struct <ty-ident-3> {}
3441
        4 struct <ty-ident-4> {}
3442
         6 fn <f-ident-1>(a: <arg-ty-ident-1>) {}
3443
         7 fn <f-ident-2>(a: <arg-ty-ident-2>) {}
         8 fn <f-ident-3>(a: <arg-ty-ident-3>) {}
3444
         9 fn <f-ident-4>(a: <arg-ty-ident-4>) {}
3445
        10
        11 fn g() {
3446
              let x: <ty-ident> = ...;
        12
3447
        13
               <f-ident>(x);
        14 }
3448
3449
3450
        I.1 LOWER BOUNDS FOR RNN: EASY BOUNDS DUE TO MEMORY
```

Proof of Theorem 4. Our proof resonates with the proof of Theorem 4.6 in Wen et al. (2024) and Theorem 8 in Bhattamishra et al. (2024). For $L, D, H \in \mathbb{N}$, suppose that D makes MiniHuskyAnnotated_{D,H} to be nontrivial, i.e., one can define functions with one parameter and use function calls. Simple calculations shows we can choose D = 7 and H = 1. If a RNN represents a function maps any token sequence of length L in MiniHuskyAnnotated_{D,H} to its type

3451

errors represented as a sequence of values of type Option<TypeError>, then the memory right before type checking must store all previous type signatures, the number of which can be as many as $\Omega(L)$ in the worst case. Assuming proper numerical discretization, the memorization of these type signatures would require the memory size to be $\Omega(L)$ in the worst case.

J ADDITIONAL EXPERIMENT DETAILS

J.1 SETUPS

Model details are shown in Table 1, and other hyperparameters are shown in Table 2.

Table 1: Model specification		
Specification	Value	
Transformer		
- Hidden size (d_h)	$\{8k \mid 1 \le k \le 8\} \cup \{240\}$	
- Num attention heads		
- Num hidden layers	8	
- Intermediate size	$2d_h$	
- Max position embeddings	≤ 2048	
RNN		
- Hidden size	$\{8k \mid 1 \le k \le 8\} \cup \{256\}$	
- Num layers	8	

Table 2: Hyperparam	eters of experiments
---------------------	----------------------

Hyperparameter	Value
Dataset	
-(n, f, d)	$\{(100000, 10, 3), (200000, 20, 5), (300000, 40, 10), (400000, 80, 20)\}$
-(a,c,v,e)	(5, 5, 0.2, 0.5)
Number of epochs	80
Train batch size	512
Optimizer	Adam
LR scheduler	Linear warmup-decay
- Warmup min lr	1×10^{-5}
- Warmup max lr	1×10^{-3}
- Warmup steps	990

J.2 ADDITIONAL RESULTS

Figures 4,5,6,7 include other metrics (train loss, accuracies for expected type in validation set, and validation loss) in the experiments. Note that for the expressive power of the models, training accuracies are better indicators.



Figure 7: Figures for the dataset with (f, a, c, d, v, e) = (80, 5, 5, 10, 0.2, 0.5).