

Deep Generation of Coq Lemma Names Using Elaborated Terms

Pengyu Nie¹, Karl Palmskog², Junyi Jessy Li¹, and Milos Gligoric¹

¹ The University of Texas at Austin, Austin, TX, USA

`pynie@utexas.edu, jessy@austin.utexas.edu, gligoric@utexas.edu`

² KTH Royal Institute of Technology, Stockholm, Sweden

`palmskog@kth.se`

Abstract. Coding conventions for naming, spacing, and other essentially stylistic properties are necessary for developers to effectively understand, review, and modify source code in large software projects. Consistent conventions in verification projects based on proof assistants, such as Coq, increase in importance as projects grow in size and scope. While conventions can be documented and enforced manually at high cost, emerging approaches automatically learn and suggest idiomatic names in Java-like languages by applying statistical language models on large code corpora. However, due to its powerful language extension facilities and fusion of type checking and computation, Coq is a challenging target for automated learning techniques. We present novel generation models for learning and suggesting lemma names for Coq projects. Our models, based on multi-input neural networks, are the first to leverage syntactic and semantic information from Coq’s lexer (tokens in lemma statements), parser (syntax trees), and kernel (elaborated terms) for naming; the key insight is that learning from elaborated terms can substantially boost model performance. We implemented our models in a toolchain, dubbed ROOSTERIZE, and applied it on a large corpus of code derived from the Mathematical Components family of projects, known for its stringent coding conventions. Our results show that ROOSTERIZE substantially outperforms baselines for suggesting lemma names, highlighting the importance of using multi-input models and elaborated terms.

Keywords: Proof assistants, Coq, lemma names, neural networks

1 Introduction

Programming language source code with deficient coding conventions, such as misleading function and variable names and irregular spacing, is difficult for developers to effectively understand, review, and modify [8, 52, 66]. Code with haphazard adherence to conventions may also be more bug-prone [17]. The problem is exacerbated in large projects with many developers, where different source code files and components may have inconsistent and clashing conventions.

Many open source software projects manually document coding conventions that contributors are expected to follow, and maintainers willingly accept fixes of

violations to such conventions [2]. Enforcement of conventions can be performed by static analysis tools [30, 58]. However, such tools require developers to write precise checks for conventions, which are tedious to define and often *incomplete*. To address this problem, researchers have proposed techniques for automatically learning coding conventions for Java-like languages from code corpora by applying statistical language models [4]. These models are applicable because code in these languages has high *naturalness* [35], i.e., statistical regularities and repetitiveness. Learned conventions can then be used to, e.g., suggest names in code.

Proof assistants, such as Coq [15], are increasingly used to formalize results in advanced mathematics [28, 29] and develop large trustworthy software systems, e.g., compilers, operating systems, file systems, and distributed systems [18, 44, 72]. Such projects typically involve contributions of many participants over several years, and require considerable effort to maintain over time. Coding conventions are essential for evolution of large verification projects, and are thus highly emphasized in the Coq libraries HoTT [37] and Iris [39], in Lean’s Mathlib [9], and in particular in the influential Mathematical Components (MathComp) *family of Coq projects* [19]. Extensive changes to adhere to conventions, e.g., on naming, are regularly requested by MathComp maintainers for proposed external contributions [50], and its conventions have been adopted, to varying degrees, by a growing number of independent Coq projects [1, 13, 24, 65].

We believe these properties make Coq code related to MathComp an attractive target for automated learning and suggesting of coding conventions, in particular, for suggesting *lemma names* [7]. However, serious challenges are posed by, on the one hand, Coq’s powerful language extension facilities and fusion of type checking and computation [12], and on the other hand, the idiosyncratic conventions used by Coq practitioners compared to software engineers. Hence, although suggesting lemma names is similar in spirit to suggesting method names in Java-like languages [73], the former task is more challenging in that lemma names are typically much shorter than method names and tend to include heavily abbreviated terminology from logic and advanced mathematics; a single character can carry significant information about a lemma’s meaning. For example, the MathComp lemma names `card_support_normedTI` (“cardinality of support groups of a normed trivial intersection group”) and `extprod_mulgA` (“associativity of multiplication operations in external product groups”) concisely convey information on lemma statement structure and meaning through both abbreviations and suffixes, as when the suffix `A` indicates an associative property.

In this paper, we present novel generation models for learning and suggesting lemma names for Coq verification projects that address these challenges. Specifically, based on our knowledge of Coq and its implementation, we developed multi-input encoder-decoder neural networks for generating names that use information directly from Coq’s internal data structures related to lexing, parsing, and type checking. In the context of naming, our models are the first to leverage the *lemma lemma statement* as well as the corresponding *syntax tree* and *elaborated term* (which we call *kernel tree*) processed by Coq’s kernel [53].

We implemented our models in a toolchain, dubbed ROOSTERIZE, which we used to learn from a high-quality Coq corpus derived from the MathComp family. We then measured the performance of ROOSTERIZE using automatic metrics, finding that it significantly outperforms baselines. Using our best model, we also suggested lemma names for the PCM library [56, 65], which were manually reviewed by the project maintainer with encouraging results.

To allow ROOSTERIZE to use information directly from Coq’s lexer, parser, and kernel, we extended the SerAPI library [26] to serialize Coq tokens, syntax trees, and kernel trees into a machine-readable format. This allowed us to achieve robustness against user-defined notations and other extensions to Coq syntax. Thanks to our integration with SerAPI and its use of metaprogramming, we expect our toolchain to only require modest maintenance as Coq evolves.

We make the following key contributions in this work:

- **Models:** We propose novel generation models based on multi-input neural networks to learn and suggest lemma names for Coq verification projects. A key property of our models is that they combine data from several Coq phases, including lexing, parsing, and term elaboration.
- **Corpus:** Advised by MathComp developers, we constructed a corpus of high-quality Coq code for learning coding conventions, totaling over 164k LOC taken from four core projects. We believe that our corpus can enable development of many novel techniques for Coq based on statistical language models.
- **Toolchain:** We implemented a toolchain, dubbed ROOSTERIZE, which suggests lemma names for a given Coq project. We envision ROOSTERIZE being useful during the review process of proposed contributions to a Coq project.
- **Evaluation:** We performed several experiments with ROOSTERIZE to evaluate our models using our corpus. Our results show that ROOSTERIZE performs significantly better than several strong baselines, as measured by standard automatic metrics [59]. The results also reveal that our novel multi-input models, as well as the incorporation of kernel trees, are important for suggestion quality. Finally, we performed a manual quality analysis by suggesting lemma names for a medium sized Coq project [56], evaluated by its maintainer, who found many of the suggestions useful for improving naming consistency.

The appendix describes more experiments, including an automatic evaluation on additional Coq projects. We provide artifacts related to our toolchain and corpus at: <https://github.com/EngineeringSoftware/roosterize>.

2 Background

This section gives brief background related to Coq and the Mathematical Components (MathComp) family of projects, as well as the SerAPI library.

Coq and Gallina: Coq is a proof assistant based on dependent types, implemented in the OCaml language [15, 20]. For our purposes, we view Coq as a programming language and type-checking toolchain. Specifically, Coq *files* are sequences of *sentences*, with each sentence ending with a period. Sentences are

```

1 Lemma mg_eq_proof L1 L2 (N1 : mgClassifier L1) : L1 =i L2 -> nerode L2 N1.
2 Proof. move => H0 u v. split => [/nerodeP H1 w|H1].
3   - by rewrite !H0.
4   - apply/nerodeP => w. by rewrite !H0.
5 Qed.

```

Fig. 1: Coq lemma on the theory of regular languages, including proof script.

essentially either (a) commands for printing and other output, (b) definitions of functions, lemmas, and datatypes in the Gallina language [21], or (c) expressions in the Ltac tactic language [22]. We will refer to definitions of lemmas as in (b) as *lemma sentences*. Coq internally represents a lemma sentence both as a sequence of tokens (lexing phase) and as a syntax tree (parsing phase).

In the typical workflow for a Coq-based verification project, users write datatypes and functions and then interactively prove lemmas about them by executing different tactic expressions that may, e.g., discharge or split the current proof goal. Both statements to be proved and proofs are represented internally as *terms* produced during an *elaboration* phase [53]; we refer to elaborated terms as *kernel trees*. Hence, as tactics are successfully executed, they gradually build a kernel tree. The `Qed` command sends the kernel tree for a tentative proof to Coq’s kernel for final certification. We call a collection of Ltac tactic sentences that build a kernel tree a *proof script*.

Fig. 1 shows a Coq lemma and its proof script, taken verbatim from a development on the theory of regular languages [24]. Line 1 contains a lemma sentence with the lemma name `mg_eq_proof`, followed by a *lemma statement* (on the same line) involving the arbitrary languages `L1` and `L2`, i.e., typed variables that are implicitly universally quantified. When Coq processes line 5, the kernel certifies that the kernel tree generated by the proof script (lines 2 to 4) has the type (is a proof) of the kernel tree for the lemma statement on line 1.

MathComp and lemma naming: The MathComp family of Coq projects, including in particular the MathComp library of general mathematical definitions and results [49], grew out of Gonthier’s proof of the four-color theorem [28], with substantial developments in the context of the landmark proof of the odd order theorem in abstract algebra [29]. The MathComp library is now used in many projects outside of the MathComp family, such as in the project containing the lemma in Fig. 1 [23]. MathComp has documented naming conventions for two kinds of entities: (1) variables and (2) functions and lemmas [19]. Variable names tend to be short and simple, while function and lemma names can be long and consist of several *name components*, typically separated by an underscore, but sometimes using CamelCase. Examples of definition and lemma names in Fig. 1 include `mg_eq_proof`, `mgClassifier`, `nerode`, and `nerodeP`. Note that lemma names sometimes have *suffixes* to indicate their meaning, such as `P` in `nerodeP` which says that the lemma is a *characteristic property*. Coq functions tend to be named based on corresponding function definition bodies rather than just types (of the parameters and return value), analogously to methods in Java [47]. In contrast, MathComp lemma names tend to be based solely on the lemma statement. Hence, a more suitable name for the lemma in Fig. 1 is `mg_eq_nerode`.

<code>Lemma mg_eq_proof L1 L2 (N1 : mgClassifier L1) : L1 =i L2 -> nerode L2 N1.</code>	sentence
<code>(Sentence((IDENT Lemma)(IDENT mg_eq_proof)(IDENT L1)(IDENT L2) (KEYWORD "(")(IDENT N1)(KEYWORD :)(IDENT mgClassifier) (IDENT L1)(KEYWORD ")")(KEYWORD :)(IDENT L1)(KEYWORD =i)(IDENT L2) (KEYWORD ->)(IDENT nerode)(IDENT L2)(IDENT N1)(KEYWORD .)))</code>	tokens
<code>(VernacExpr()(VernacStartTheoremProof Lemma (Id mg_eq_proof) (((CLocalAssum(Name(Id L1))(CHole()IntroAnonymous())) (CLocalAssum(Name(Id L2))(CHole()IntroAnonymous())) (CLocalAssum(Name(Id N1)) (CApp(CRef(Ser.Qualid(DirPath()))(Id mgClassifier)))(CRef(Ser.Qualid(DirPath()))(Id L1)))))) (CNotation(InConstrEntrySomeLevel"- -> -") (CNotation(InConstrEntrySomeLevel"- =i -") (CRef(Ser.Qualid(DirPath()))(Id L1))(CRef(Ser.Qualid(DirPath()))(Id L2))) (CApp(CRef(Ser.Qualid(DirPath()))(Id nerode))) (CRef(Ser.Qualid(DirPath()))(Id L2))(CRef(Ser.Qualid(DirPath()))(Id N1))))))</code>	syntax tree
<code>(Prod (Name (Id char)) ... (Prod (Name (Id L1)) ... (Prod (Name (Id L2)) ... (Prod (Name (Id N1)) ... (Prod Anonymous (App (Ref (DirPath ((Id ssrbool) (Id ssr) (Id Coq))) (Id eq_mem)) ... (Var (Id L1)) ... (Var (Id L2))) (App (Ref (DirPath ((Id myhill_nerode) (Id RegLang))) (Id nerode)) ... (Var (Id L2)) ... (Var (Id N1))))))</code>	kernel tree

Fig. 2: Coq lemma sentence at the top, with sexps for, from just below to bottom: tokens, syntax tree, and kernel tree; the lemma statement in each is highlighted.

SerAPI and Coq serialization: SerAPI is an OCaml library and toolchain for machine interaction with Coq [26], which provides serialization and deserialization of Coq internal data structures to and from S-expressions (sexps) [51]. SerAPI is implemented using OCaml’s PPX metaprogramming facilities [57], which enable modifying OCaml program syntax trees at compilation time. Fig. 2 shows the lemma sentence on line 1 in Fig. 1, and below it, the corresponding (simplified) sexps for its tokens, syntax tree, and kernel tree, with the lemma statement highlighted in each representation. Note that the syntax tree omits the types of some quantified variables, e.g., for the types of `L1` and `L2`, as indicated by the `CHole` constructor. Note also that during elaboration of the syntax tree into the kernel tree by Coq, an implicit variable `char` is added (all-quantified via `Prod`), and the extensional equality operator `=i` is translated to its globally unique *kernel name*, `Coq.ssr.ssrbool.eq_mem`. Hence, a kernel tree can be much larger and contain more information than the corresponding syntax tree.

3 Models

In this section, we describe our multi-input generation models for suggesting Coq lemma names. Our models consider lemma name generation with an *encoder-decoder* mindset, i.e., we use neural architectures specifically designed for transduction tasks [67]. This family of architectures is commonly used for sequence generation, e.g., in machine translation [11] and code summarization [43], where it has been found to be much more effective than traditional probabilistic and retrieval-based approaches.

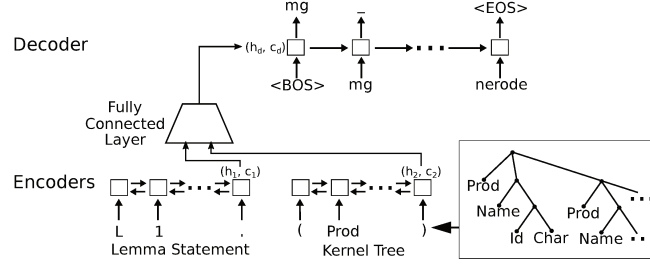


Fig. 3: Core architecture of our multi-input encoder-decoder models.

3.1 Core Architecture

Our encoders are Recurrent Neural Networks (RNNs) that learn a deep semantic representation of a given lemma statement from its tokens, syntax tree, and kernel tree. The decoder—another RNN—generates the descriptive lemma name as a sequence. The model is trained end-to-end, maximizing the probability of the generated lemma name given the input. In contrast to prior work in language-code tasks that uses a single encoder [27], we design multi-input models that leverage both syntactic and semantic information from Coq’s lexer, parser, and kernel. A high-level visualization of our architecture is shown in Fig. 3.

Encoding: Our multi-input encoders combine different kinds of syntactic and semantic information in the encoding phase. We use a different encoder for each input, which are: lemma statement, syntax tree, and kernel tree.

Coq data structure instances can be large, with syntax trees having an average depth of 28.03 and kernel trees 46.51 in our corpus (we provide detailed statistics in Section 4). Therefore, we flatten the trees into sequences, which can be trained more efficiently than tree encoders without performance loss [38]. We flatten the trees with pre-order traversal, and we use “(” and “)” as the boundaries of the children of a node. In later parts of this paper, we use syntax and kernel trees to refer to their flattened versions. In Section 3.2, we introduce *tree chopping* to reduce the length of the resulting sequences.

To encode lemma statements and flattened tree sequences, we use bi-directional Long-Short Term Memory (LSTM) [36] networks. LSTMs are advanced RNNs good at capturing long-range dependencies in a sequence, and are widely used in encoders [38]. A bi-directional LSTM learns stronger representations (than a uni-directional LSTM) by encoding a sequence from both left to right and right to left [74].

Decoding: We use an LSTM (left to right direction only) as our decoder. To obtain the initial hidden and cell states (h_d, c_d) of the decoder, we learn a unified representation of these separate encoders by concatenating their final hidden and cell states (h_i, c_i), and then applying a fully connected layer on the concatenated states: $h_d = W_h \cdot \text{concat}([h_i]) + b_h$ and $c_d = W_c \cdot \text{concat}([c_i]) + b_c$, where W_h , W_c , b_h , and b_c are learnable parameters.

During training, we maximize the log likelihood of the reference lemma name given all input sequences. Standard beam search is used to reduce the search

<pre>(Prod Anonymous (App (Ref (DirPath ((Id ssrbool) (Id ssr) (Id Coq))) (Id eq_mem))) ... ((App (Ref ...))) ...))</pre>
<pre>(Prod Anonymous (App eq_mem ... (App (Ref ...)) ...))</pre>

Fig. 4: Kernel tree sexp before and after chopping; chopped parts are highlighted.

space for the optimal sequence of tokens. With regular decoding, at each time step the decoder generates a new token relying on the preceding *generated* token, which can be error-prone and leads to slow convergence and instability. We mitigate this problem by performing decoding with teacher forcing [71] such that the decoder relies on the preceding *reference* token. At test time, the decoder still uses the proceeding generated token as input.

Attention: With RNN encoders, the input sequence is compressed into the RNN’s final hidden states, which results in a loss of information, especially for longer sequences. The attention mechanism [48] grants the decoder access to the encoder hidden and cell states for all previous tokens. At each decoder time step, an attention vector is calculated as a distribution over all encoded tokens, indicating which token the decoder should “pay attention to”. To make the attention mechanism work with multiple encoders, we concatenate the hidden states of the n encoders $[h_1, \dots, h_n]$ and apply an attention layer on the result [69].

Initialization: Since there are no pre-trained token embeddings for Coq, we initialize each unique token in the vocabulary with a random vector sampled from the uniform distribution $U(-0.1, 0.1)$. These embeddings are trained together with the model. The hidden layer parameters of the encoders and decoders are also initialized with random vectors sampled from the same uniform distribution.

3.2 Tree Chopping

While syntax and kernel trees for lemma statements can be large, not all parts of the trees are relevant for naming. For instance, each constant reference is expanded to its fully qualified form in the kernel tree, but the added prefixes are usually related to directory paths and likely do not contain relevant information for generating the name of the lemma. Irrelevant information in long sequences can be detrimental to the model, since the model would have to reason about and encode all tokens in the sequence.

To this end, we implemented *chopping* heuristics for both syntax trees and kernel trees to remove irrelevant parts. The heuristics essentially: (1) replace the fully qualified name sub-trees with only the last component of the name; (2) remove the location information from sub-trees; (3) extract the singletons, i.e., non-leaf nodes that have only one child. Fig. 4 illustrates the chopping of a kernel tree, with the upper box showing the tree before chopping with the parts to be removed highlighted, and the lower box showing the tree after chopping. In the example in the figure, we chopped a fully qualified name and extracted a singleton. These heuristics greatly reduce the size of the tree: for kernel trees, they reduce the average depth from 39.20 to 11.39.

Our models use chopped trees as the inputs to the encoders. As we discuss in more detail in Section 6, the chopped trees help the models to focus better on the relevant parts of the inputs. While the attention mechanism in principle could learn what the relevant parts of the trees are, our evaluation shows that it can easily be overwhelmed by large amounts of irrelevant information.

3.3 Copy Mechanism

We found it common for lemma name tokens to only occur in a single Coq file, whence they are unlikely to appear in the vocabulary learned from the training set, but can still appear in the respective lemma statement, syntax tree, or kernel tree. For example, `mg` occurs in both the lemma name and lemma statement in Fig. 1, but not outside the file the lemma is in. To account for this, we adopt the copy mechanism [63] which improves the generalizability of our model by allowing the decoder to *copy* from inputs rather than always choosing one word from the fixed vocabulary from the training set. To handle multiple encoders, similar to what we did with the attention layer, we concatenate the hidden states of each encoder and apply a copy layer on the concatenated hidden states.

3.4 Sub-tokenization

We sub-tokenize all inputs (lemma statements, syntax and kernel trees) and outputs (lemma names) in a pre-processing step. Previous work on learning from software projects has shown that sub-tokenization helps to reduce the sparsity of the vocabulary and improves the performance of the model [10]. However, unlike Java-like languages where the method names (almost) always follow the CamelCase convention, lemma names in Coq use a mix of snake_case, Camel-Case, prefixes, and suffixes, thus making sub-tokenization more complex. For example, `extprod.mulgA` should be sub-tokenized to `extprod`, `_`, `mul`, `g`, and `A`.

To perform sub-tokenization, we implemented a set of heuristics based on the conventions outlined by MathComp developers [19]. After sub-tokenization, the vocabulary size of lemma names in our corpus was reduced from 8,861 to 2,328. When applying the sub-tokenizer on the lemma statements and syntax and kernel trees, we sub-tokenize the identifiers and not the keywords or operators.

3.5 Repetition Prevention

We observed that decoders often generated repeated tokens, e.g., `mem_mem_mem`. This issue also exists in natural language summarization [68]. We further observed that it is very unlikely to have repeated sub-tokens in lemma names used by proof engineers (only 1.37% of cases in our corpus). Hence, we simply forbid the decoder from repeating a sub-token (modulo “_”) during beam search.

4 Corpus

We constructed a corpus of four large Coq projects from the MathComp family, totaling 164k lines of code (LOC). We selected these projects based on the

Table 1: Projects from the MathComp Family Used in Our Corpus.

Project	SHA	#Files	#Lemmas	#Toks	LOC		LOC/file	
					Spec.	Proof	Spec.	Proof
finmap	27642a8	4	940	78,449	4,260	2,191	1,065.00	547.75
fourcolor	0851d49	60	1,157	560,682	9,175	27,963	152.92	466.05
math-comp	748d716	89	8,802	1,076,096	38,243	46,470	429.70	522.13
odd-order	ca602a4	34	367	519,855	11,882	24,243	349.47	713.03
Avg.	N/A	46.75	2,816.50	558,770.50	15,890.00	25,216.75	339.89	539.40
Σ	N/A	187	11,266	2,235,082	63,560	100,867	63,560	100,867

Table 2: Statistics on the Lemmas in the Training, Validation, and Testing Sets.

	#Files	#Lemmas	Name		Stmt	
			#Char	#SubToks	#Char	#SubToks
training	152	8,861	10.14	4.22	44.16	19.59
validation	18	1,085	9.20	4.20	38.28	17.30
testing	17	1,320	9.76	4.34	48.49	23.20

recommendation of MathComp developers, who emphasized their high quality and stringent adherence to coding conventions. Our corpus is *self-contained*: there are inter-project dependencies within the corpus, but no project depends on a project outside the corpus (except Coq’s standard library). All projects build with Coq version 8.10.2. Note that we need to be able to build projects to be able to extract tokens, syntax trees, and kernel trees.

Constituent projects: Table 1 lists the projects in the corpus, along with basic information about each project. The table includes columns for the project identifier, revision SHA, number of files (#Files), number of lemmas (#Lemmas), number of tokens (#Toks), LOC for specifications (Spec.) and proof scripts (Proof), and average LOC per file for specifications and proof scripts. The math-comp SHA corresponds to version 1.9.0 of the library. The LOC numbers are computed with Coq’s bundled `coqwc` tool. The last two rows of the table show the averages and sums across all projects.

Corpus statistics: We extracted all lemmas from the corpus, and initially we obtained 15,005 lemmas in total. However, we found several outlier lemmas where the lemma statement, syntax tree and kernel tree were very large. To ensure stable training, and similar to prior work on generating method names for Java [47], we excluded the lemmas with the deepest 25% kernel trees. This left us with 11,266 lemmas. Column 4 of Table 1 shows the number of lemmas after filtering.

We randomly split corpus files into training, validation, and testing sets which contain 80%, 10%, 10% of the files, respectively. Table 2 shows statistics on the lemmas in each set, which includes columns for the number of files, the number of lemmas, the average number of characters and sub-tokens in lemma names, and the average number of characters and sub-tokens in lemma statements.

Fig. 5 illustrates the changes of the depth, number of nodes and number of sub-tokens (after flattening) of the kernel trees (first row) and syntax trees (second row) before and after chopping. Our chopping process reduced tree depth by 70.9% for kernel trees and 70.7% for syntax trees, and reduced the number

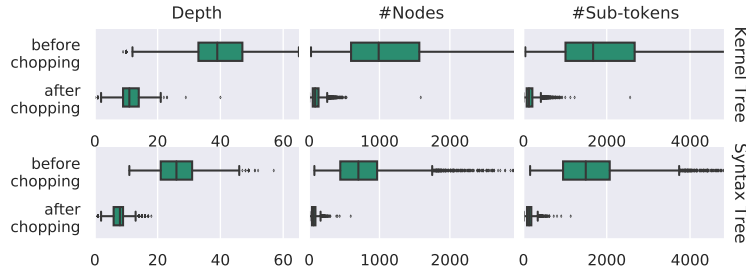


Fig. 5: Statistics on syntax and kernel trees.

of nodes by 91.5% for kernel trees and 90.8% for syntax trees; after flattening, the resulting average sequence length is, for kernel trees 165 comparing to the original 2,056, and for syntax trees 144 comparing to the original 1,590. We provide additional statistics on lemmas before filtering in the appendix.

5 Implementation

In this section, we briefly describe our toolchain which implements the models in Section 3 and processes and learns from the corpus in Section 4; we dub this toolchain ROOSTERIZE. The components of the toolchain can be divided into two categories: (1) components that interact with Coq or directly process information extracted from Coq, and (2) components concerned with machine learning and name generation.

The first category includes several OCaml-based tools integrated with SerAPI [26] (and thus Coq itself), and Python-based tools for processing of data obtained via SerAPI from Coq. All OCaml tools have either already been included in, or accepted for inclusion into, SerAPI itself. The tools are as follows:

sercomp: We integrated the existing program **sercomp** distributed with SerAPI into ROOSTERIZE to serialize Coq files to lists of sexps for syntax trees.

sertok: We developed an OCaml program dubbed **sertok** on top of SerAPI. The program takes a Coq file as input and produces sexps of all tokens found by Coq’s lexer in the file, organized at the sentence level.

sername: We developed an OCaml program dubbed **sername** on top of SerAPI. The program takes a list of fully qualified (kernel) lemma names and produces sexps for the kernel trees of the corresponding lemma statements.

postproc & subtokenizer: We created two small independent tools in Python to post-process Coq sexps and perform sub-tokenization, respectively.

For the second category, we implemented our machine learning models in Python using two widely-used deep learning libraries: PyTorch [60] and OpenNMT [41]. More specifically, we extended the sequence-to-sequence models in OpenNMT to use multi-input encoders, and extended attention and copy layers to use multiple inputs. Source code for the components of ROOSTERIZE is available from: <https://github.com/EngineeringSoftware/roosterize>.

6 Evaluation

This section presents an extensive evaluation of our models as implemented in ROOSTERIZE. Our automatic evaluation (Section 6.2) compares ROOSTERIZE with a series of strong baselines and reports on ablation experiments; additional experiments, e.g., on chopping heuristics, are described in the appendix. Our manual quality assessment (Section 6.3) analyzes 150 comments we received from the maintainer of the PCM library on names suggested by ROOSTERIZE for that project using our best model.

6.1 Models and Baselines

We study the combinations of: (1) using individual input (lemma statement and trees) in a single encoder, or multi-input encoders with different mixture of these inputs; and (2) using the attention and copy mechanisms. Our inputs include: lemma statement (*Stmnt*), syntax tree (*SynTree*), chopped syntax tree (*ChopSynTree*), kernel tree (*KnITree*), and chopped kernel tree (*ChopKnITree*). For multiple inputs, the models are named by concatenating inputs with “+”; a “+” is also used to denote the presence of attention (*attn*) or copy (*copy*). For example, *Stmnt+ChopKnITree+attn+copy* refers to a model that uses two encoders—one for lemma statement and one for chopped kernel tree—and uses attention and copy mechanisms.

We consider the vanilla encoder-decoder models with only one input (lemma statement, kernel tree, or syntax tree) as baseline models. We also compare with a retrieval-based baseline model implemented using Lucene [6]: a k-nearest neighbors classifier using the tf-idf of the tokens in lemma statement as features.

Hyperparameters are tuned on the validation set within the following options: embedding dimensions from {200, 500, 1000}, number of hidden units in each LSTM from {200, 500, 1000}, number of stacked LSTM layers from {1, 2, 3}. We set the dropout rate between LSTM layers to 0.5. We set the output dimension of the fully connected layer for combining encoders to the same number as the number of hidden units in each LSTM. We checked the validation loss every 200 training steps (as defined in OpenNMT [41], which is similar to one training epoch on our dataset), and set an early stopping threshold of 3. We used the Adam [40] optimizer with a learning rate of 0.001. We used a beam size of 5 in beam search. All the experiments were run with one NVIDIA 1080-TI GPU and Intel Xeon E5-2620 v4 CPU.

6.2 Automatic Evaluation

Metrics: We use four automatic metrics which evaluate generated lemma names against the reference lemma name (as written by developers) in the testing set. Each metric captures a different level of granularity of the generation quality. *BLEU* [59] is a standard metric used in transduction tasks including language \leftrightarrow code transduction. It calculates the number of n-grams in a generated sequence that also appear in the reference sequence, where one “n-gram” is n consecutive items in a sequence (in our case, one “n-gram” is n consecutive characters

Table 3: Results of ROOSTERIZE Models.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn+copy	45.4	22.2%	7.5%	16.5%
	Stmt+ChopKnlTree+attn+copy	47.2	24.9%	9.6%	18.0%
	Stmt+ChopSynTree+attn+copy	37.7	18.1%	6.1%	10.6%
	ChopKnlTree+ChopSynTree+attn+copy	45.4	22.9%	7.6%	15.3%
Single-input	ChopKnlTree+attn+copy	42.9	19.8%	5.0%	11.7%
	ChopSynTree+attn+copy	39.8	18.3%	6.8%	12.2%
	KnlTree+attn+copy	37.0	14.2%	2.2%	8.4%
	SynTree+attn+copy	31.0	10.8%	2.8%	6.1%
	Stmt+attn+copy	38.9	19.4%	6.9%	11.6%
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn	24.5	8.6%	0.4%	0.9%
	Stmt+ChopKnlTree+attn	25.6	8.5%	0.9%	1.7%
	Stmt+ChopSynTree+attn	23.8	8.2%	0.8%	1.6%
	ChopKnlTree+ChopSynTree+attn	28.4	10.9%	1.8%	3.4%
Single-input	ChopKnlTree+attn	19.5	4.9%	0.6%	1.3%
	ChopSynTree+attn	28.9	12.1%	1.5%	2.9%
	KnlTree+attn	14.1	1.6%	0.0%	0.0%
	SynTree+attn	8.8	1.0%	0.0%	0.0%
	Stmt+attn	26.9	11.1%	1.1%	2.5%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	17.7	3.5%	0.1%	0.2%
	Stmt+ChopKnlTree	19.5	4.5%	0.1%	0.3%
	Stmt+ChopSynTree	12.6	0.6%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	16.7	2.4%	0.0%	0.1%
Single-input	ChopKnlTree	15.5	1.6%	0.0%	0.0%
	ChopSynTree	14.5	0.8%	0.1%	0.1%
	KnlTree	12.0	0.6%	0.0%	0.0%
	SynTree	5.7	0.4%	0.0%	0.0%
	Stmt	20.0	4.7%	0.1%	0.3%
-	Retrieval-based	28.3	10.0%	0.2%	0.3%

in the sequence of characters of the lemma name). We use it to compute the 1 ~ 4-grams overlap between the characters in generated name and characters in the reference name, averaged between 1 ~ 4-grams with smoothing method proposed by Lin and Och [46]. *Fragment accuracy* computes the accuracy of generated names on the fragment level, which is defined by splitting the name by underscores (“_”). For example, `map_determinant_mx` has a fragment accuracy of 66.7% when evaluated against `det_map_mx`. Unlike BLEU, fragment accuracy ignores the ordering of the fragments. Finally, *top-1 accuracy* and *top-5 accuracy* compute how often the true name fully matches the generated name or is one of the top-5 generated names.

Results: Table 3 shows the performance of the models. Similar models are grouped together. The first column shows the names of the model groups and the second column shows the names of the models. For each model, we show values for the four automatic metrics, BLEU, fragment accuracy (Frag.Acc.), top-1 accuracy (Top1), and top-5 accuracy (Top5). We repeated each experiment 3 times, with different random initialization each time, and computed the averages of each automated metric. We performed statistical significance tests—under significance level $p < 0.05$ using the bootstrap method [14]—to compare each pair of models. We use bold text to highlight the best value for each automatic metric, and gray background for baseline models. We make several observations:

Finding #1: The best overall performance (BLEU = 47.2) is obtained using the multi-input model with lemma statement and chopped kernel tree as inputs, which also includes copy and attention mechanisms (Stmt+ChopKnlTree+attn+copy). The improvements over all other models are statistically significant and all automatic metrics are consistent in identifying the best model. This shows the importance of using Coq’s internal structures and focusing only on certain parts of those structures.

Finding #2: The copy mechanism brings statistically significant improvements to all models. This can be clearly observed by comparing groups 1 and 3 in the table, as well as groups 2 and 4. For example, BLEU for Stmt+attn and Stmt+attn+copy are 26.9 and 38.9, respectively. We believe that the copy mechanism plays an important role because many sub-tokens are specific to the file context and do not appear in the fixed vocabulary learned on the files in training set.

Finding #3: Using chopped trees greatly improves performance of models and the improvements brought by upgrading KnlTree to ChopKnlTree or SynTree to ChopSynTree are statistically significant. For example, this can be clearly seen in the second group: BLEU for KnlTree+attn+copy and ChopKnlTree+attn+copy are 37.0 and 42.9, respectively. We believe that the size of the original trees, and a lot of irrelevant data in those trees, hurt the performance. The fact that ChopKnlTree and ChopSynTree both perform much better than using KnlTree or SynTree across all groups indicate that the chopped trees could be viewed as a form of supervised attention with flat values that helps later attention layers to focus better.

Finding #4: Although chopped syntax tree with attention outperforms (statistically significant) chopped kernel tree with attention (BLEU 28.9 vs. 19.5), chopped kernel tree with attention and copy by far outperforms (statistically significant) chopped syntax tree with attention and copy (BLEU 42.9 vs. 39.8). The copy mechanism helps kernel trees much more than the syntax trees, because the mathematical notations and symbols in the syntax trees get expanded to their names in the kernel trees, and some of them are needed as a part of the lemma names.

Finding #5: Lemma statement and syntax tree do not work well together, primarily because the two representations contain mostly the same information. In which case, a model taking both as inputs may not work as well as using only one of the inputs, because more parameters need to be trained.

Finding #6: The retrieval-based baseline, which is the strongest among baselines, outperforms several encoder-decoder models without attention and copy or with only attention, but is worse than (statistically significant) all models with both attention and copy mechanisms enabled.

6.3 Manual Quality Analysis

While generated lemma names may not always match the manually written ones in the training set, they can still be semantically valid and conform to prevailing conventions. However, such name properties are not reflected in our automatic evaluation metrics, since these metrics only consider exactly matched tokens as

Table 4: Manual Quality Analysis Representative Examples.

Lemma statement: <code>p s : supp (kfilter p s) = filter p (supp s)</code> Hand-written: <code>supp_kfilt</code> Roosterize: <code>supp_kfilter</code> Comment: ✓ Using only <code>kfilt</code> has cognitive overhead.
Lemma statement: <code>g e k v f : path ord k (supp f) -> foldfmap g e (ins k v f) = g (k, v) (foldfmap g e f)</code> Hand-written: <code>foldf_ins</code> Roosterize: <code>foldfmap_ins</code> Comment: ✓ The whole function name is used in the suggested name.
Lemma statement: <code>: transitive (@ord T)</code> Hand-written: <code>trans</code> Roosterize: <code>ord_trans</code> Comment: ✓ Useful to add the <code>ord</code> prefix to the name.
Lemma statement: <code>s : sorted (@ord T) s -> sorted (@oleq T) s</code> Hand-written: <code>sorted_oleq</code> Roosterize: <code>ord_sorted</code> Comment: ✗ The conclusion content should have greater priority.
Lemma statement: <code>x y : total_spec x y (ord x y) (x == y) (ord y x)</code> Hand-written: <code>totalP</code> Roosterize: <code>ordP</code> Comment: ✗ Maybe this lemma should be named <code>ord_totalP</code> ?
Lemma statement: <code>p1 p2 s : kfilter (predI p1 p2) s = kfilter p1 (kfilter p2 s)</code> Hand-written: <code>kfilter_predI</code> Roosterize: <code>eq_kfilter</code> Comment: ✗ The suggested name is too generic.

correct. To obtain a more complete evaluation, we therefore performed a manual quality analysis of generated lemma names from ROOSTERIZE by applying it to a Coq project outside of our corpus, the PCM library [56]. This project depends on MathComp, and follows, to a degree, many of the MathComp coding conventions. The PCM library consists of 12 Coq files, and contains 690 lemmas.

We ran ROOSTERIZE with the best model (Stmnt+ChopKnITree+attn+copy) on the PCM library to get the top-1 suggestions for all lemma names. Overall, the ROOSTERIZE suggestions achieved a BLEU score of 36.3 and a fragment accuracy of 17%, and 36 suggestions (5%) exactly match the existing lemma names. Next, we asked the maintainer of the PCM library to evaluate the remaining 654 lemma names (those that do not match exactly) and send us feedback.

The maintainer spent one day on the task and provided comments on 150 suggested names. We analyzed these comments to identify patterns and trends. He found that 20% of the suggested names he inspected were of good quality, out of which more than half were of high quality. Considering that the analysis was of top-1 suggestions excluding exact matches, we find these figures encouraging. For low-quality names, a clear trend was that they were often “too generic”. Similar observations have been made about the results from encoder-decoder models in dialog generation [45, 64]. In contrast, useful suggestions were typically able to expand or elaborate on name components that are intuitively too concise, e.g., replacing `kfilt` with `kfilter`. Table 4 lists examples that are representative of these trends; checkmarks indicate useful suggestions, while crosses indicate unsuitability. We also include comments from the maintainer. As illustrated by the comments, even suggestions considered unsuitable may contain useful parts.

7 Discussion

Our toolchain builds on Coq 8.10.2, and thus we only used projects that support this version. However, we do not expect any fundamental obstacles in supporting future Coq releases. Thanks to the use of OCaml metaprogramming via PPX, which allowed eliding explicit references to the internal structure of Coq datatypes, SerAPI itself and our extensions to it are expected to require only modest effort to maintain as Coq evolves.

Our models and toolchain may not be applicable to Coq projects unrelated to the MathComp family of projects, i.e., projects which do not follow any MathComp conventions. To the best of our knowledge, MathComp’s coding conventions are the most recognizable and well-documented in the Coq community; suggesting coding conventions based on learning from projects unrelated to MathComp are likely to give more ambiguous results that are difficult to validate manually. Our case study also included generating suggestions for a project outside the MathComp family, the PCM library, with encouraging results.

Our models are in principle applicable to proof assistants with similar foundations, such as Lean [54]. However, the current version of Lean, Lean 3, does not provide serialization of internal data structures as SerAPI does for Coq, which prevents direct application of our toolchain. Application of our models to proof assistants with different foundations and proof-checking toolchains, such as Isabelle/HOL, is even less straightforward, although the Archive of Formal Proofs (AFP) contains many projects with high-quality lemma names [25].

8 Related Work

Naturalness and coding conventions: Hindle et al. [35] first applied the concept of naturalness to Java-like languages, noting that program statement regularities and repetitiveness make statistical language models applicable for performing software engineering tasks [4]. Rahman et al. [61] validated the naturalness of other similar programming languages, and Hellendoorn et al. [31] found high naturalness in Coq code, providing motivation for our application of statistical language models to Coq. Allamanis et al. [2] used the concept of naturalness and statistical language models to learn and suggest coding conventions, including names, for Java, and Raychev et al. [62] used conditional random fields to learn and suggest coding conventions for JavaScript. To our knowledge, no previous work has developed *applications* of naturalness for proof assistants; Hellendorn et al. [31] only measured naturalness for their Coq corpus.

Suggesting names: Prior work on suggesting names mostly concerns Java method names. Liu et al. [47] used a similarity matching algorithm, based on deep representations of Java method names and bodies learned with Paragraph Vector and convolutional neural networks, to detect and fix inconsistent Java method names. Allamanis et al. [3] used logbilinear neural language models supplemented by additional manual features to predict Java method and class names. Java method names have also been treated as short, descriptive “summaries” of its body; in this view, prior work has augmented attention mechanisms in convolutional networks [5], used sequence-to-sequence models to learn

from descriptions (e.g., Javadoc comments) [27], and utilized the tree-structure of the code in a hierarchical attention network [73]. Unlike Java syntax trees, Coq syntax and kernel trees contain considerable semantic information useful for naming. In the work closest to our domain, Aspinall and Kaliszyk used a k-nearest neighbors multi-label classifier on a corpus for the HOL Light proof assistant to suggest names of lemmas [7]. However, their technique only suggests names that exist in the training data and therefore does not generalize. To our knowledge, ours is the first neural generation model for suggesting names in a proof assistant context.

Mining and learning for proof assistants: Müller et al. [55] exported Coq kernel trees as XML strings to translate 49 Coq projects to the OMDoc theory graph format. Rather than translating documents to an independently specified format, we produce lightweight machine-readable representations of Coq’s internal data structures. Wiedijk [70] collected early basic statistics on the core libraries of several proof assistants, including Coq and Isabelle/HOL. Blanchette et al. [16] mined the AFP to gather statistics such as the average number of lines of Isabelle/HOL specifications and proof scripts. However, these corpora were not used to perform learning. Komendantskaya et al. [32, 33, 34, 42] used machine learning without neural networks to identify patterns in Coq tactic sequences and proof kernel trees, e.g., to find structural similarities between lemmas and simplify proof development. In contrast, our models capture similarity among several different representations of lemma *statements* to generate lemma names.

9 Conclusion

We presented novel techniques, based on neural networks, for learning and suggesting lemma names in Coq verification projects. We designed and implemented multi-input encoder-decoder models that use Coq’s internal data structures, including (chopped) syntax trees and kernel trees. Additionally, we constructed a large corpus of high quality Coq code that will enable development and evaluation of future techniques for Coq. We performed an extensive evaluation of our models using the corpus. Our results show that the multi-input models, which use internal data structures, substantially outperform several baselines; the model that uses the lemma statement tokens and the chopped kernel tree with attention and copy mechanism performs the best. Based on our findings, we believe that multi-input models leveraging key parts of internal data structures can play a critical role in producing high-quality lemma name suggestions.

Acknowledgments

The authors thank Yves Bertot, Cyril Cohen, Emilio Jesús Gallego Arias, Gaëtan Gilbert, Hugo Herbelin, Anton Trunov, Théo Zimmermann, and the anonymous reviewers for their comments and feedback. This work was partially supported by the US National Science Foundation under Grant Nos. CCF-1652517 and IIS-1850153, and by the Swedish Foundation for Strategic Research under the TrustFull project.

Bibliography

- [1] Affeldt, R., Garrigue, J.: Formalization of error-correcting codes: From Hamming to modern coding theory. In: Urban, C., Zhang, X. (eds.) International Conference on Interactive Theorem Proving. LNCS, vol. 9236, pp. 17–33. Springer, Cham, Switzerland (2015). https://doi.org/10.1007/978-3-319-22102-1_2
- [2] Allamanis, M., Barr, E.T., Bird, C., Sutton, C.: Learning natural coding conventions. In: International Symposium on the Foundations of Software Engineering. pp. 281–293. ACM, New York, NY, USA (2014). <https://doi.org/10.1145/2635868.2635883>
- [3] Allamanis, M., Barr, E.T., Bird, C., Sutton, C.: Suggesting accurate method and class names. In: Joint Meeting on Foundations of Software Engineering. pp. 38–49. ACM, New York, NY, USA (2015). <https://doi.org/10.1145/2786805.2786849>
- [4] Allamanis, M., Barr, E.T., Devanbu, P., Sutton, C.: A survey of machine learning for big code and naturalness. *ACM Comput. Surv.* **51**(4), 81:3–81:37 (2018). <https://doi.org/10.1145/3212695>
- [5] Allamanis, M., Peng, H., Sutton, C.: A convolutional attention network for extreme summarization of source code. In: International Conference on Machine Learning. pp. 2091–2100 (2016)
- [6] Apache Software Foundation: Apache Lucene (2020), <https://lucene.apache.org>, last accessed 2020-01-23
- [7] Aspinall, D., Kaliszyk, C.: What’s in a theorem name? In: Blanchette, J.C., Merz, S. (eds.) International Conference on Interactive Theorem Proving. LNCS, vol. 9807, pp. 459–465. Springer, Cham, Switzerland (2016). https://doi.org/10.1007/978-3-319-43144-4_28
- [8] Avidan, E., Feitelson, D.G.: Effects of variable names on comprehension: An empirical study. In: International Conference on Program Comprehension. pp. 55–65. IEEE Computer Society, Washington, DC, USA (2017). <https://doi.org/10.1109/ICPC.2017.27>
- [9] Avigad, J.: Mathlib naming conventions (2016), <https://github.com/leanprover-community/mathlib/blob/snapshot-2019-10/docs/contribute/naming.md>, last accessed 2020-01-23
- [10] Babii, H., Janes, A., Robbes, R.: Modeling vocabulary for big code machine learning. *CoRR* **abs/1904.01873** (2019), <https://arxiv.org/abs/1904.01873>
- [11] Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. In: International Conference on Learning Representations (2015), <https://arxiv.org/abs/1409.0473>
- [12] Barendregt, H., Barendsen, E.: Autarkic computations in formal proofs. *Journal of Automated Reasoning* **28**(3), 321–336 (2002). <https://doi.org/10.1023/A:1015761529444>

- [13] Bartzia, E.I., Strub, P.Y.: A formal library for elliptic curves in the Coq proof assistant. In: Klein, G., Gamboa, R. (eds.) International Conference on Interactive Theorem Proving. LNCS, vol. 8558, pp. 77–92. Springer, Cham, Switzerland (2014). https://doi.org/10.1007/978-3-319-08970-6_6
- [14] Berg-Kirkpatrick, T., Burkett, D., Klein, D.: An empirical investigation of statistical significance in NLP. In: Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. pp. 995–1005. Association for Computational Linguistics, Stroudsburg, PA, USA (2012)
- [15] Bertot, Y., Castéran, P.: Interactive Theorem Proving and Program Development: Coq’Art: The Calculus of Inductive Constructions. Springer, Heidelberg, Germany (2004). <https://doi.org/10.1007/978-3-662-07964-5>
- [16] Blanchette, J.C., Haslbeck, M., Matichuk, D., Nipkow, T.: Mining the archive of formal proofs. In: Kerber, M., Carette, J., Kaliszyk, C., Rabe, F., Sorge, V. (eds.) International Conference on Intelligent Computer Mathematics. LNCS, vol. 9150, pp. 3–17. Springer, Cham, Switzerland (2015). https://doi.org/10.1007/978-3-319-20615-8_1
- [17] Boogerd, C., Moonen, L.: Evaluating the relation between coding standard violations and faults within and across software versions. In: International Working Conference on Mining Software Repositories. pp. 41–50. IEEE Computer Society, Washington, DC, USA (2009). <https://doi.org/10.1109/MSR.2009.5069479>
- [18] Chen, H., Ziegler, D., Chajed, T., Chlipala, A., Kaashoek, M.F., Zeldovich, N.: Using crash Hoare logic for certifying the FSCQ file system. In: Symposium on Operating Systems Principles. pp. 18–37. ACM, New York, NY, USA (2015). <https://doi.org/10.1145/2815400.2815402>
- [19] Cohen, C.: Contribution guide for the Mathematical Components library (2018), <https://github.com/math-comp/math-comp/blob/mathcomp-1.9.0/CONTRIBUTING.md>, last accessed 2020-04-17
- [20] Coq Development Team: The Coq proof assistant, version 8.10.0 (Oct 2019). <https://doi.org/10.5281/zenodo.3476303>
- [21] Coq Development Team: The Gallina specification language (2019), <https://coq.inria.fr/distrib/V8.10.2/refman/language/gallina-specification-language.html>, last accessed 2020-04-17
- [22] Delahaye, D.: A tactic language for the system Coq. In: Parigot, M., Voronkov, A. (eds.) International Conference on Logic for Programming, Artificial Intelligence, and Reasoning. LNCS, vol. 1955, pp. 85–95. Springer, Heidelberg, Germany (2000). https://doi.org/10.1007/3-540-44404-1_7
- [23] Doczkal, C., Kaiser, J.O., Smolka, G.: Regular language representations in Coq (2020), <https://github.com/coq-community/reglang>, last accessed 2020-04-09
- [24] Doczkal, C., Smolka, G.: Regular language representations in the constructive type theory of Coq. *Journal of Automated Reasoning* **61**(1), 521–553 (2018)
- [25] Eberl, M., Klein, G., Nipkow, T., Paulson, L., Thiemann, R.: Archive of Formal Proofs (2020), <https://www.isa-afp.org>, last accessed 2020-01-23

- [26] Gallego Arias, E.J.: SerAPI: Machine-friendly, data-centric serialization for Coq. Tech. rep., MINES ParisTech (2016), <https://hal-mines-paristech.archives-ouvertes.fr/hal-01384408>
- [27] Gao, S., Chen, C., Xing, Z., Ma, Y., Song, W., Lin, S.: A neural model for method name generation from functional description. In: International Conference on Software Analysis, Evolution and Reengineering. pp. 414–421. IEEE Computer Society, Washington, DC, USA (2019). <https://doi.org/10.1109/SANER.2019.8667994>
- [28] Gonthier, G.: Formal proof—the four-color theorem. Notices of the American Mathematical Society **55**(11), 1382–1393 (2008)
- [29] Gonthier, G., Asperti, A., Avigad, J., Bertot, Y., Cohen, C., Garillot, F., Le Roux, S., Mahboubi, A., O’Connor, R., Ould Biha, S., Pasca, I., Rideau, L., Solovyev, A., Tassi, E., Théry, L.: A machine-checked proof of the odd order theorem. In: Blazy, S., Paulin-Mohring, C., Pichardie, D. (eds.) International Conference on Interactive Theorem Proving. LNCS, vol. 7998, pp. 163–179. Springer, Heidelberg, Germany (2013). https://doi.org/10.1007/978-3-642-39634-2_14
- [30] Google: google-java-format (2020), <https://github.com/google/google-java-format>, last accessed 2020-01-23
- [31] Hellendoorn, V.J., Devanbu, P.T., Alipour, M.A.: On the naturalness of proofs. In: International Symposium on the Foundations of Software Engineering, New Ideas and Emerging Results. pp. 724–728. ACM, New York, NY, USA (2018). <https://doi.org/10.1145/3236024.3264832>
- [32] Heras, J., Komendantskaya, E.: ML4PG in computer algebra verification. In: Carette, J., Aspinall, D., Lange, C., Sojka, P., Windsteiger, W. (eds.) International Conference on Intelligent Computer Mathematics. LNCS, vol. 7961, pp. 354–358. Springer, Heidelberg, Germany (2013). https://doi.org/10.1007/978-3-642-39320-4_28
- [33] Heras, J., Komendantskaya, E.: Proof pattern search in Coq/SSReflect. CoRR **abs/1402.0081** (2014), <https://arxiv.org/abs/1402.0081>
- [34] Heras, J., Komendantskaya, E.: Recycling proof patterns in Coq: Case studies. Mathematics in Computer Science **8**(1), 99–116 (2014). <https://doi.org/10.1007/s11786-014-0173-1>
- [35] Hindle, A., Barr, E.T., Su, Z., Gabel, M., Devanbu, P.: On the naturalness of software. In: International Conference on Software Engineering. pp. 837–847. IEEE Computer Society, Washington, DC, USA (2012). <https://doi.org/10.1109/ICSE.2012.6227135>
- [36] Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural Computation **9**(8), 1735–1780 (1997). <https://doi.org/10.1162/neco.1997.9.8.1735>
- [37] HoTT authors: HoTT Conventions And Style Guide (2019), <https://github.com/HoTT/HoTT/blob/V8.10/STYLE.md>, last accessed 2020-01-23
- [38] Hu, X., Li, G., Xia, X., Lo, D., Jin, Z.: Deep code comment generation. In: International Conference on Program Comprehension. pp. 200–210. ACM, New York, NY, USA (2018). <https://doi.org/10.1145/3196321.3196334>
- [39] Iris authors: Iris Style Guide (2019), <https://gitlab.mpi-sws.org/iris/iris/blob/iris-3.2.0/StyleGuide.md>, last accessed 2020-04-17

- [40] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (2015), <https://arxiv.org/abs/1412.6980>
- [41] Klein, G., Kim, Y., Deng, Y., Senellart, J., Rush, A.M.: OpenNMT: Open-source toolkit for neural machine translation. In: Annual Meeting of the Association for Computational Linguistics, System Demonstrations. pp. 67–72. Association for Computational Linguistics, Stroudsburg, PA, USA (2017). <https://doi.org/10.18653/v1/P17-4012>
- [42] Komendantskaya, E., Heras, J., Grov, G.: Machine learning in Proof General: Interfacing interfaces. In: Kaliszyk, C., Lüth, C. (eds.) International Workshop On User Interfaces for Theorem Provers. EPTCS, vol. 118, pp. 15–41. Open Publishing Association, Sydney, Australia (2013). <https://doi.org/10.4204/EPTCS.118.2>
- [43] LeClair, A., Jiang, S., McMillan, C.: A neural model for generating natural language summaries of program subroutines. In: International Conference on Software Engineering. pp. 795–806. IEEE Computer Society, Washington, DC, USA (2019). <https://doi.org/10.1109/ICSE.2019.00087>
- [44] Leroy, X.: Formal verification of a realistic compiler. *Commun. ACM* **52**(7), 107–115 (2009). <https://doi.org/10.1145/1538788.1538814>
- [45] Li, J., Galley, M., Brockett, C., Gao, J., Dolan, B.: A diversity-promoting objective function for neural conversation models. In: Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 110–119. Association for Computational Linguistics, Stroudsburg, PA, USA (2016). <https://doi.org/10.18653/v1/n16-1014>
- [46] Lin, C., Och, F.J.: ORANGE: A method for evaluating automatic evaluation metrics for machine translation. In: International Conference on Computational Linguistics. pp. 501–507. Association for Computational Linguistics, Stroudsburg, PA, USA (2004)
- [47] Liu, K., Kim, D., Bissyandé, T.F., Kim, T., Kim, K., Koyuncu, A., Kim, S., Traon, Y.L.: Learning to spot and refactor inconsistent method names. In: International Conference on Software Engineering. pp. 1–12. IEEE Computer Society, Washington, DC, USA (2019). <https://doi.org/10.1109/ICSE.2019.00019>
- [48] Luong, T., Pham, H., Manning, C.D.: Effective approaches to attention-based neural machine translation. In: Empirical Methods in Natural Language Processing. pp. 1412–1421. Association for Computational Linguistics, Stroudsburg, PA, USA (2015). <https://doi.org/10.18653/v1/d15-1166>
- [49] Mahboubi, A., Tassi, E.: Mathematical Components Book (2017), <https://math-comp.github.io/mcb/>, last accessed 2020-04-17
- [50] Mathematical Components team: Missing lemmas in seq (2016), <https://github.com/math-comp/math-comp/pull/41>, last accessed 2020-04-18
- [51] McCarthy, J.: Recursive functions of symbolic expressions and their computation by machine, part I. *Commun. ACM* **3**(4), 184–195 (1960). <https://doi.org/10.1145/367177.367199>

- [52] Miara, R.J., Musselman, J.A., Navarro, J.A., Shneiderman, B.: Program indentation and comprehensibility. *Commun. ACM* **26**(11), 861–867 (1983). <https://doi.org/10.1145/182.358437>
- [53] de Moura, L., Avigad, J., Kong, S., Roux, C.: Elaboration in dependent type theory. *CoRR* **abs/1505.04324** (2015), <https://arxiv.org/abs/1505.04324>
- [54] de Moura, L., Kong, S., Avigad, J., van Doorn, F., von Raumer, J.: The Lean theorem prover (system description). In: Felty, A.P., Middeldorp, A. (eds.) *International Conference on Automated Deduction*. LNCS, vol. 9195, pp. 378–388. Springer, Cham, Switzerland (2015). https://doi.org/10.1007/978-3-319-21401-6_26
- [55] Müller, D., Rabe, F., Sacerdoti Coen, C.: The Coq library as a theory graph. In: Kaliszyk, C., Brady, E., Kohlhase, A., Sacerdoti Coen, C. (eds.) *International Conference on Intelligent Computer Mathematics*. LNCS, vol. 11617, pp. 171–186. Springer, Cham, Switzerland (2019). https://doi.org/10.1007/978-3-030-23250-4_12
- [56] Nanevski, A., Ley-Wild, R., Sergey, I., Delbianco, G., Trunov, A.: The PCM library (2020), <https://github.com/imdea-software/fcsl-pcm>, last accessed 2020-01-24
- [57] OCaml Labs: PPX (2017), <http://ocamlabs.io/doc/ppx.html>, last accessed 2020-01-23
- [58] Ogura, N., Matsumoto, S., Hata, H., Kusumoto, S.: Bring your own coding style. In: *International Conference on Software Analysis, Evolution and Reengineering*. pp. 527–531. IEEE Computer Society, Washington, DC, USA (2018). <https://doi.org/10.1109/SANER.2018.8330253>
- [59] Papineni, K., Roukos, S., Ward, T., Zhu, W.: BLEU: A method for automatic evaluation of machine translation. In: *Annual Meeting of the Association for Computational Linguistics*. pp. 311–318. Association for Computational Linguistics, Stroudsburg, PA, USA (2002)
- [60] Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in PyTorch. In: *Autodiff Workshop* (2017), <https://openreview.net/forum?id=BJJsrnfCZ>
- [61] Rahman, M., Palani, D., Rigby, P.C.: Natural software revisited. In: *International Conference on Software Engineering*. pp. 37–48. IEEE Computer Society, Washington, DC, USA (2019). <https://doi.org/10.1109/ICSE.2019.00022>
- [62] Raychev, V., Vechev, M., Krause, A.: Predicting program properties from “big code”. In: *Symposium on Principles of Programming Languages*. pp. 111–124. ACM, New York, NY, USA (2015). <https://doi.org/10.1145/2676726.2677009>
- [63] See, A., Liu, P.J., Manning, C.D.: Get to the point: Summarization with pointer-generator networks. In: *Annual Meeting of the Association for Computational Linguistics*. pp. 1073–1083. Association for Computational Linguistics, Stroudsburg, PA, USA (2017). <https://doi.org/10.18653/v1/P17-1099>

- [64] Serban, I.V., Sordoni, A., Bengio, Y., Courville, A., Pineau, J.: Building end-to-end dialogue systems using generative hierarchical neural network models. In: *AAAI Conference on Artificial Intelligence*. pp. 3776–3783. AAAI Press, Palo Alto, CA, USA (2016)
- [65] Sergey, I., Nanevski, A., Banerjee, A.: Mechanized verification of fine-grained concurrent programs. In: *Conference on Programming Language Design and Implementation*. pp. 77–87. ACM, New York, NY, USA (2015). <https://doi.org/10.1145/2737924.2737964>
- [66] Shneiderman, B., McKay, D.: Experimental investigations of computer program debugging and modification. *Human Factors Society Annual Meeting* **20**(24), 557–563 (1976). <https://doi.org/10.1177/154193127602002401>
- [67] Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: *Advances in Neural Information Processing Systems* 27. pp. 3104–3112. MIT Press, Cambridge, MA, USA (2014)
- [68] Suzuki, J., Nagata, M.: Cutting-off redundant repeating generations for neural abstractive summarization. In: *Conference of the European Chapter of the Association for Computational Linguistics*. pp. 291–297. Association for Computational Linguistics, Stroudsburg, PA, USA (2017). <https://doi.org/10.18653/v1/e17-2047>
- [69] Unanue, I.J., Borzeshi, E.Z., Piccardi, M.: A shared attention mechanism for interpretation of neural automatic post-editing systems. In: *Workshop on Neural Machine Translation and Generation*. pp. 11–17. Association for Computational Linguistics, Stroudsburg, PA, USA (2018). <https://doi.org/10.18653/v1/w18-2702>
- [70] Wiedijk, F.: Statistics on digital libraries of mathematics. *Studies in Logic, Grammar and Rhetoric* **18**(31), 137–151 (2009)
- [71] Williams, R.J., Zipser, D.: A learning algorithm for continually running fully recurrent neural networks. *Neural Computation* **1**(2), 270–280 (1989). <https://doi.org/10.1162/neco.1989.1.2.270>
- [72] Woos, D., Wilcox, J.R., Anton, S., Tatlock, Z., Ernst, M.D., Anderson, T.: Planning for change in a formal verification of the Raft consensus protocol. In: *Certified Programs and Proofs*. pp. 154–165. ACM, New York, NY, USA (2016). <https://doi.org/10.1145/2854065.2854081>
- [73] Xu, S., Zhang, S., Wang, W., Cao, X., Guo, C., Xu, J.: Method name suggestion with hierarchical attention networks. In: *Workshop on Partial Evaluation and Program Manipulation*. pp. 10–21. ACM, New York, NY, USA (2019). <https://doi.org/10.1145/3294032.3294079>
- [74] Zhang, S., Zheng, D., Hu, X., Yang, M.: Bidirectional long short-term memory networks for relation classification. In: *Pacific Asia Conference on Language, Information and Computation*. pp. 207–212. Association for Computational Linguistics, Stroudsburg, PA, USA (2015). <https://doi.org/10.18653/v1/p16-2034>

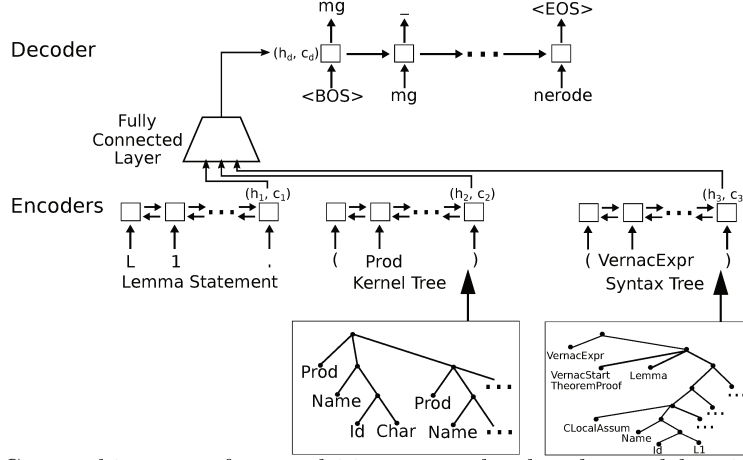


Fig. 6: Core architecture of our multi-input encoder-decoder models, with lemma statements, syntax trees and kernel trees as inputs.

A Explanatory Notes on Roosterize Models

In this section, we explain some key terminology and concepts used to describe our generation models; these explanations were omitted from the main text to conserve space and avoid distracting the reader with excessive detail.

Architecture of our multi-input encoder-decoder models with lemma statement, syntax tree, and kernel tree: Fig. 6 illustrates our architecture when all three inputs are used, in contrast to Fig. 3, which only uses two inputs (lemma statement and kernel tree).

Repetition of decoders’ generated tokens (cf. Section 3.5): This is a common problem in encoder-decoder models; it is largely because the attention mechanism (while helping the model for the most part) does not store information on how much information the model has “covered” in the encoded sequence. See also Sutskever et al. [67].

Tf-idf (cf. the retrieval-based baseline model in Section 6.1): This is a numerical metric reflecting the importance of a token to a document in a corpus, calculated as the product of term frequency (proportional to the frequency of the token in the document) and inverse document frequency (inversely proportional to the number of documents containing the token). In our retrieval-based baseline model, we used Lucene’s implementation of tf-idf [6].

Early stopping (cf. hyperparameters in Section 6.1): This is a common strategy to mitigate overfitting in training a machine learning model by monitoring the model’s performance on both the training set and the validation set and halting the training if the model stops improving on the validation set even if it improves on the training set. If early stopping is not used, the model is fully trained to maximize its performance on the training set, but may have bad performance on a separate set (e.g., testing set). In our experiments, we set an early stopping

Table 5: Additional Statistics on the Lemmas in the Training, Validation, and Testing Sets.

	#Files	#Lemmas	Name		Stmt	
			#Char	#SubToks	#Char	#SubToks
before filtering	187	15,005	10.91	4.47	56.31	24.69
after filtering		11,266	10.00	4.23	44.10	19.79
training	152	8,861	10.14	4.22	44.16	19.59
validation	18	1,085	9.20	4.20	38.28	17.30
testing	17	1,320	9.76	4.34	48.49	23.20

threshold of 3, which means the training is halted if the model does not obtain smaller loss on the validation set for 3 consecutive checkpoints.

Learning rate (cf. hyperparameters in Section 6.1): This controls the speed of adjusting models’ learnable parameters based on the loss at each iteration of the training. An excessively large learning rate makes training faster, but may result in “overshooting”: adjusting so much that it results in jumping over the minima. A too low learning rate means training is unnecessarily slow to complete, and may result in the training getting stuck in a local minimum. Guided by our previous experience, we used a value of 0.001 paired with the Adam [40] optimizer (an algorithm for adjusting models’ learnable parameters).

Dropout (cf. hyperparameters in Section 6.1): This is a regularization technique for reducing overfitting, by randomly resetting a fraction of neural connections between two layers during training (and during training only). In our experiments, a dropout rate of 0.5 between the LSTM layers means that 50% of the bits of the hidden and cell states are set to 0 when they are passed from a previous layer to its next layer in the LSTM during training.

B Additional Statistics

Table 5 shows additional statistics on the lemmas we used. The first row is for the lemmas before filtering the outliers, the second row is for the lemmas after the filtering, and the last three rows are for the training, validation, and testing sets, respectively. Fig. 7 illustrates the changes in the depth, number of nodes and number of sub-tokens (after flattening) of the kernel trees (first row) and syntax trees (second row) before filtering, after filtering, and after chopping.

C Ablation Study on Tree Chopping Heuristics

In order to corroborate the effectiveness of ROOSTERIZE’s tree chopping heuristics, we designed an ablation study that applies three different sets of chopping heuristics and compares them with the one in ROOSTERIZE (Section 3.2). The three sets of chopping heuristics are:

- **Keep-category chopping**: This set of heuristics is almost the same as ROOSTERIZE chopping, except that it keeps the category of a referenced

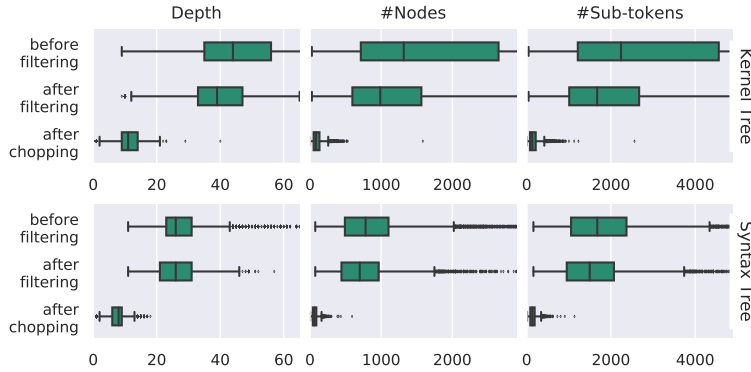


Fig. 7: Additional statistics on syntax and kernel trees.

Table 6: Results of the Ablation Study on Tree Chopping.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
ROOSTERIZE Chopping	Stmt+ChopKnlTree+ChopSynTree+attn+copy	45.4	22.2%	7.5%	16.5%
	Stmt+ChopKnlTree+attn+copy	47.2	24.9%	9.6%	18.0%
	Stmt+ChopSynTree+attn+copy	37.7	18.1%	6.1%	10.6%
	ChopKnlTree+ChopSynTree+attn+copy	45.4	22.9%	7.6%	15.3%
Keep-category Chopping	Stmt+ChopKnlTree+ChopSynTree+attn+copy	46.8	25.3%	9.5%	19.0%
	Stmt+ChopKnlTree+attn+copy	47.2	25.2%	9.4%	18.0%
	Stmt+ChopSynTree+attn+copy	37.1	17.9%	6.2%	10.5%
	ChopKnlTree+ChopSynTree+attn+copy	46.4	22.6%	7.5%	15.0%
Rule-based Chopping	Stmt+ChopKnlTree+ChopSynTree+attn+copy	37.0	17.7%	5.9%	10.5%
	Stmt+ChopKnlTree+attn+copy	38.8	19.7%	6.7%	11.0%
	Stmt+ChopSynTree+attn+copy	36.9	16.0%	6.3%	10.5%
	ChopKnlTree+ChopSynTree+attn+copy	13.2	0.4%	0.0%	0.0%
Random Chopping	Stmt+ChopKnlTree+ChopSynTree+attn+copy	37.7	19.2%	6.7%	10.9%
	Stmt+ChopKnlTree+attn+copy	38.5	19.3%	7.3%	11.3%
	Stmt+ChopSynTree+attn+copy	38.0	18.6%	7.1%	11.4%
	ChopKnlTree+ChopSynTree+attn+copy	17.1	2.7%	0.1%	0.1%

name in kernel trees (e.g., whether it is a constant or inductive type), since that semantic information could be relevant for naming.

- **Rule-based chopping:** Removes all nodes after depth 10 for syntax and kernel trees. This is similar to the proof kernel tree processing heuristics used in ML4PG [33].
- **Random chopping:** Randomly removes a subset of nodes from syntax and kernel trees so that the resulting trees have the same average number of nodes compared to ROOSTERIZE’s chopped trees, i.e., the heuristic removes 90.9% nodes from syntax trees and 91.4% nodes from kernel trees.

We performed the ablation study using the encoder-decoder models with multi-input encoders and with both the attention and copy mechanisms. The other hyperparameters and experimental settings are the same as in Section 6.

Table 6 shows the results of the ablation study. The models using the same chopping heuristics are grouped together. We make several observations:

- Among keep-category chopping models, Stmt+ChopKnlTree+ChopSynTree+attn+copy and Stmt+ChopKnlTree+attn+copy perform the best, and they have performance similar to Stmt+ChopKnlTree+attn+copy using ROOSTERIZE chopping (ROOSTERIZE’s best model). The measured differences between these three models are not statistically significant, under significance level $p < 0.05$ using the bootstrap method [14]. This indicates that although the category of a referenced name may contain some relevant semantic information, the most relevant information is already preserved by ROOSTERIZE chopping heuristics.
- The models using rule-based chopping and random chopping have poor performance. This indicates that the performance gain achieved by ROOSTERIZE through chopping is not only due to the size reduction of the input trees, but also due to the relevant information retained by our chopping heuristics.

D Expanded Corpus and Evaluation

In addition to evaluating ROOSTERIZE using the high-quality corpus consisting of four MathComp projects (Section 4), we also performed an evaluation on an expanded corpus that includes 21 Coq projects related to the MathComp family which follow (to various degrees) the same coding conventions, totaling over 297k LOC. All projects depend, directly or indirectly, on the MathComp library, but not on projects outside the corpus itself except for Coq’s standard library. We introduce the expanded corpus in Section D.1 and describe our additional evaluation on this corpus in Section D.2.

D.1 Expanded Corpus

Table 7 lists the projects in the expanded corpus, along with basic information about each project. The expanded corpus consists of two parts: the main part consists of 20 projects and is used for training and evaluating ROOSTERIZE; the left-out (LO) part is one project, infotheo, which is used to study the generalizability of ROOSTERIZE on an unseen project.

We constructed and organized the corpus based on recommendations from MathComp developers. The 4 core MathComp projects used in the original corpus are included as the *tier 1* set. We selected 9 projects for the *tier 2* set, such that each included project (1) has a main contributor who is also a significant contributor to one of the tier 1 projects, and (2) follows to a significant degree the coding conventions specified for MathComp. (infotheo would be in this set had we not added it to the left-out part.) Finally, we selected 8 projects which follow MathComp coding conventions but do not fulfil the tier 2 criteria, for inclusion in the *tier 3* set.

We briefly describe each project in our corpus:

analysis: A library for general real analysis.

bigenough: A small library for $\epsilon - N$ reasoning.

bits: A library for reasoning about bit-level operations.

Table 7: Projects Used in Our Expanded Corpus.

	Project	SHA	#Files	#Lemmas	#Toks	LOC		LOC/file	
						Spec.	Proof	Spec.	Proof
Tier 1	finmap	27642a8	4	940	78,449	4,260	2,191	1,065.00	547.75
	fourcolor	0851d49	60	1,157	560,682	9,175	27,963	152.92	466.05
	math-comp	748d716	89	8,802	1,076,096	38,243	46,470	429.70	522.13
	odd-order	ca602a4	34	367	519,855	11,882	24,243	349.47	713.03
Tier 2	analysis	9e5fe1d	17	969	152,542	5,553	6,186	326.65	363.88
	bigenough	5f79a32	1	4	731	70	8	70.00	8.00
	elliptic-curves	631af89	18	625	110,480	3,298	6,298	183.22	349.89
	grobner	dfa54f9	1	81	15,656	312	1,018	312.00	1,018.00
	multinomials	691d795	5	831	83,438	3,699	3,664	739.80	732.80
	real-closed	495a1fa	10	561	108,925	4,348	4,577	434.80	457.70
	robot	b341ad1	13	864	130,024	3,881	7,249	298.54	557.62
	two-square	1c09aca	2	200	20,326	413	1,308	206.50	654.00
Tier 3	bits	3cd298a	10	411	40,420	1,578	2,463	157.80	246.30
	comp-dec-pdl	c1f813b	16	494	61,731	2,305	2,114	144.06	132.12
	disel	e8aa80c	20	256	51,473	2,575	1,898	128.75	94.90
	fcsl-pcm	eef4503	12	690	70,273	2,937	2,852	244.75	237.67
	games	3d3bd31	12	231	43,438	1,450	3,503	120.83	291.92
	monae	9d0e461	18	349	73,578	3,422	3,233	190.11	179.61
	reglang	da333e0	12	230	41,327	1,299	1,734	108.25	144.50
	toychain	97bd697	14	67	61,997	1,747	3,528	124.79	252.00
Main	Avg.	N/A	18.40	906.45	165,072.05	5,122.35	7,625.00	278.39	414.40
	Σ	N/A	368	18,129	3,301,441	102,447	152,500	102,447	152,500
LO	infotheo	6c17242	81	1,891	463,593	12,517	29,778	154.53	367.63
All	Avg.	N/A	21.38	953.33	179,287.33	5,474.48	8,679.90	256.04	405.96
	Σ	N/A	449	20,020	3,765,034	114,964	182,278	114,964	182,278

comp-dec-pdl: Formal proofs of completeness and decidability of converse propositional dynamic logic.

disel: A framework for distributed separation logic, useful for verifying implementations of distributed systems.

elliptic-curves: A formalization of the algebraic theory of elliptic curves.

fcsl-pcm: A library formalizing partial commutative monoids, which are useful for reasoning about pointer-based programs.

finmap: A library with definitions and results about finite maps and sets with finitely many members.

fourcolor: An updated version of the formal proof of the four-color theorem in graph theory, which states that in all planar graphs, four colors suffice for coloring all vertices such that no two adjacent vertices have the same color.

games: Definitions and formal proofs of theorems in algorithmic game theory.

grobner: A formalization of Gröbner bases.

math-comp: The MathComp library itself.

monae: A library for monadic equational reasoning.

multinomials: A library formalizing monoidal rings and multinomials, and related results.

Table 8: Statistics on the Lemmas in the Training, Validation, and Testing Sets of the Expanded Corpus in Each Tier.

	#Lemmas	Name		Stmt	
		#Char	#SubToks	#Char	#SubToks
before filtering	23,615	10.57	4.32	59.21	25.93
after filtering	18,129	9.89	4.13	47.48	21.16
training	15,011	9.99	4.12	47.93	21.20
All Tiers validation	1,556	9.20	4.08	41.44	18.65
testing	1,562	9.68	4.26	49.19	23.21
training	8,861	10.14	4.22	44.16	19.59
Tier 1 validation	1,085	9.20	4.20	38.28	17.30
testing	1,320	9.76	4.34	48.49	23.20
training	3,692	10.03	4.02	50.52	22.37
Tier 2 validation	403	9.27	3.86	46.26	20.92
testing	40	8.75	3.77	49.25	22.32
training	2,458	9.37	3.90	57.65	25.28
Tier 3 validation	68	8.74	3.44	63.34	26.82
testing	202	9.35	3.82	53.70	23.43

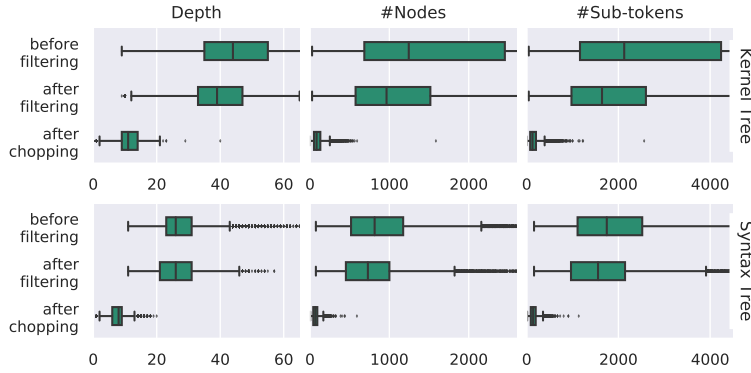


Fig. 8: Statistics of syntax and kernel trees in the expanded corpus.

odd-order: The formal proof of the odd order (Feit-Thompson) theorem in abstract algebra, which states that all groups of odd order are solvable.

real-closed: Theorems on real closed fields in algebra.

reglang: A formalization of the theory of regular languages.

robot: A formalization of the mathematics of rigid body transformations to enable proofs about robot manipulators.

toychain: Formalization and verification of a blockchain network protocol.

two-square: A proof of Fermat’s theorem on the sum of two squares, including a definition of Gaussian integers.

infotheo: Formalizations of notions and results from information theory and probability theory.

We follow the same procedure as for the original corpus to extract the lemmas and filter out the lemmas with the deepest 25% of the kernel trees. Table 8 shows statistics on the lemmas obtained from each set and each tier. Fig. 8

Table 9: Links to Tables with Results for Various Combinations of Training, Validation, and Testing Sets.

Training & Validation	Testing	Results Table
All Tiers	All Tiers	Table 10
All Tiers	Tier 1	Table 11
All Tiers	Tier 2	Table 12
All Tiers	Tier 3	Table 13
Tier 1	All Tiers	Table 14
Tier 1	Tier 1	Table 3
Tier 1	Tier 2	Table 15
Tier 1	Tier 3	Table 16
Tier 2	Tier 2	Table 17
Tier 3	Tier 3	Table 18

illustrates the changes of depth, number of nodes, and number of sub-tokens (after serialization) of the kernel trees (first row) and syntax trees (second row) before filtering, after filtering, and after chopping.

D.2 Automatic Evaluation on the Expanded Corpus

To investigate whether ROOSTERIZE can benefit from learning from a larger but less focused corpus than in our original automatic evaluation, we experimented with different combinations of training, validation, and testing sets. Table 9 lists the combinations we used; the first column shows the corpus that training and validation sets are from, the second column shows the corpus that testing set is from, and the third column indicates the results table for each combination.

We conclude that all our observations in Section 6.2 on training and testing on our original corpus (here, tier 1) hold when training and testing on all tiers. Additionally, we make the following observations based on the results of models using different combinations of training, validation, and testing sets:

- Training on all tiers helps ROOSTERIZE obtain better performance, although the corpus includes some noise from tier 2 and tier 3 projects. This observation is based on comparing the results of training on different sets and testing on the same set. For example, when testing on all tiers, the best BLEU score among models trained on all tiers (47.2, cf. Table 10) is higher than the best score for models trained on tier 1 (44.5, cf. Table 14). As another example, when testing on tier 2, the best BLEU score among models trained on all tiers is 38.7 (cf. Table 12), which is higher than the best score among models trained on tier 2, namely, 33.6 (cf. Table 17).
- Tier 2 and tier 3 projects are indeed less conforming to MathComp naming conventions than tier 1 projects, confirming the judgment of domain experts. With the same models trained on all tiers, ROOSTERIZE’s best BLEU score on the tier 1 testing set (49.3) is greater than the best BLEU score on the tier 2 testing set (38.7), and the latter is greater than the best BLEU score on the tier 3 testing set (37.4). The same relationships hold for the models trained only on tier 1.

Table 10: Results of ROOSTERIZE Models with Training and Validation Sets from All Tiers and Testing Set from All Tiers.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input +attn +copy	Stmt+ChopKnlTree+ChopSynTree+attn+copy	43.2	23.2%	7.1%	15.5%
	Stmt+ChopKnlTree+attn+copy	47.2	26.1%	10.3%	19.0%
	Stmt+ChopSynTree+attn+copy	34.9	18.0%	4.9%	10.7%
	ChopKnlTree+ChopSynTree+attn+copy	44.2	22.2%	7.4%	14.8%
Single-input +attn +copy	ChopKnlTree+attn+copy	44.1	20.9%	5.8%	13.1%
	ChopSynTree+attn+copy	39.0	19.1%	7.9%	13.3%
	KnlTree+attn+copy	35.4	14.6%	1.2%	6.4%
	SynTree+attn+copy	31.3	14.2%	3.4%	7.1%
	Stmt+attn+copy	39.7	20.8%	7.5%	13.6%
Multi-input +attn	Stmt+ChopKnlTree+ChopSynTree+attn	23.1	7.9%	1.1%	2.0%
	Stmt+ChopKnlTree+attn	27.3	10.9%	1.6%	3.0%
	Stmt+ChopSynTree+attn	23.6	9.5%	1.7%	3.0%
	ChopKnlTree+ChopSynTree+attn	26.6	10.4%	2.5%	4.5%
Single-input +attn	ChopKnlTree+attn	22.8	7.0%	1.0%	1.7%
	ChopSynTree+attn	31.0	13.1%	2.5%	4.8%
	KnlTree+attn	13.5	2.0%	0.1%	0.4%
	SynTree+attn	11.5	1.8%	0.0%	0.1%
	Stmt+attn	27.5	11.0%	1.1%	2.0%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	18.2	4.3%	0.2%	0.4%
	Stmt+ChopKnlTree	20.3	5.5%	0.4%	0.8%
	Stmt+ChopSynTree	11.2	0.1%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	12.1	0.8%	0.0%	0.0%
Single-input	ChopKnlTree	14.1	1.3%	0.0%	0.0%
	ChopSynTree	14.4	1.1%	0.1%	0.1%
	KnlTree	11.1	0.0%	0.0%	0.0%
	SynTree	10.8	0.0%	0.0%	0.0%
	Stmt	20.3	5.4%	0.3%	0.5%
-	Retrieval-based	28.2	10.9%	0.6%	0.8%

D.3 Generalizability Case Study

Infotheo consists of 81 Coq files, and contains 1,891 lemmas. We randomly split the files into training, validation, and testing sets which contain 40%, 10%, 50% of the files, respectively. After splitting, there were 580 lemmas in the training set, 144 lemmas in the validation set, and 1,167 lemmas in the testing set.

Table 19 shows the results of applying ROOSTERIZE with the best model on infotheo without and with additional training. The first column shows the number of lemmas from the infotheo training set used for additional training. The rest of the columns show the four automatic metrics. We can observe that applying ROOSTERIZE without additional training achieves moderate performance (BLEU = 33.9). With some additional training, performance can be markedly improved (up to a BLEU score of 37.4 when training on all 580 lemmas).

Table 11: Results of ROOSTERIZE Models with Training and Validation Sets from All Tiers and Testing Set from Tier 1.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnITree+ChopSynTree+attn+copy	45.0	23.5%	7.6%	16.9%
	Stmt+ChopKnITree+attn+copy	49.3	26.7%	11.1%	20.6%
	Stmt+ChopSynTree+attn+copy	35.0	17.3%	4.8%	10.4%
	ChopKnITree+ChopSynTree+attn+copy	45.9	22.3%	8.1%	16.3%
Single-input	ChopKnITree+attn+copy	45.8	21.3%	6.1%	14.1%
	ChopSynTree+attn+copy	39.4	18.8%	7.9%	13.4%
	KnITree+attn+copy	36.4	14.6%	1.0%	6.3%
	SynTree+attn+copy	31.2	13.8%	3.6%	7.5%
	Stmt+attn+copy	40.2	20.4%	7.7%	14.0%
Multi-input	Stmt+ChopKnITree+ChopSynTree+attn	23.5	7.8%	1.1%	2.0%
	Stmt+ChopKnITree+attn	28.1	10.9%	1.6%	3.0%
	Stmt+ChopSynTree+attn	23.9	8.9%	1.5%	2.8%
	ChopKnITree+ChopSynTree+attn	27.8	10.6%	2.5%	4.7%
Single-input	ChopKnITree+attn	23.4	6.8%	1.0%	1.8%
	ChopSynTree+attn	32.1	13.3%	2.7%	5.0%
	KnITree+attn	13.6	1.9%	0.2%	0.5%
	SynTree+attn	11.4	1.9%	0.0%	0.1%
	Stmt+attn	27.9	10.7%	1.0%	1.9%
Multi-input	Stmt+ChopKnITree+ChopSynTree	18.4	4.2%	0.2%	0.4%
	Stmt+ChopKnITree	20.6	5.4%	0.4%	0.7%
	Stmt+ChopSynTree	11.3	0.1%	0.0%	0.0%
	ChopKnITree+ChopSynTree	12.2	0.8%	0.0%	0.0%
Single-input	ChopKnITree	14.3	1.3%	0.0%	0.0%
	ChopSynTree	14.5	1.2%	0.1%	0.2%
	KnITree	11.1	0.0%	0.0%	0.0%
	SynTree	10.8	0.0%	0.0%	0.0%
	Stmt	20.6	5.2%	0.4%	0.5%
-	Retrieval-based	29.0	10.5%	0.3%	0.3%

Table 12: Results of ROOSTERIZE Models with Training and Validation Sets from All Tiers and Testing Set from Tier 2.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input +attn +copy	Stmt+ChopKnlTree+ChopSynTree+attn+copy	31.7	20.7%	5.8%	15.0%
	Stmt+ChopKnlTree+attn+copy	35.5	25.0%	14.2%	21.7%
	Stmt+ChopSynTree+attn+copy	35.6	18.9%	7.5%	13.3%
	ChopKnlTree+ChopSynTree+attn+copy	35.4	29.2%	9.2%	13.3%
Single-input +attn +copy	ChopKnlTree+attn+copy	32.4	21.2%	6.7%	12.5%
	ChopSynTree+attn+copy	38.7	23.9%	10.8%	18.3%
	KnlTree+attn+copy	29.2	19.9%	6.7%	10.0%
	SynTree+attn+copy	27.5	11.2%	1.7%	4.2%
	Stmt+attn+copy	33.2	17.8%	7.5%	16.7%
Multi-input +attn	Stmt+ChopKnlTree+ChopSynTree+attn	23.2	10.6%	2.5%	6.7%
	Stmt+ChopKnlTree+attn	26.0	14.2%	4.2%	10.0%
	Stmt+ChopSynTree+attn	25.9	17.6%	6.7%	10.0%
	ChopKnlTree+ChopSynTree+attn	22.4	15.4%	7.5%	8.3%
Single-input +attn	ChopKnlTree+attn	23.2	10.8%	3.3%	5.0%
	ChopSynTree+attn	30.2	18.8%	6.7%	11.7%
	KnlTree+attn	13.7	2.9%	0.0%	0.0%
	SynTree+attn	9.6	1.7%	0.0%	0.0%
	Stmt+attn	27.0	15.1%	4.2%	10.0%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	20.4	7.9%	1.7%	2.5%
	Stmt+ChopKnlTree	18.8	7.6%	0.8%	2.5%
	Stmt+ChopSynTree	11.7	0.4%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	11.9	0.4%	0.0%	0.0%
Single-input	ChopKnlTree	15.0	2.5%	0.0%	0.0%
	ChopSynTree	14.8	1.8%	0.0%	0.0%
	KnlTree	12.1	1.4%	0.0%	0.0%
	SynTree	12.3	0.0%	0.0%	0.0%
	Stmt	23.6	13.7%	0.8%	0.8%
-	Retrieval-based	31.0	27.5%	2.5%	7.5%

Table 13: Results of ROOSTERIZE Models with Training and Validation Sets from All Tiers and Testing Set from Tier 3.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn+copy	33.6	21.4%	3.8%	6.8%
	Stmt+ChopKnlTree+attn+copy	36.0	22.6%	4.3%	8.6%
	Stmt+ChopSynTree+attn+copy	34.4	22.1%	5.0%	11.9%
	ChopKnlTree+ChopSynTree+attn+copy	34.5	20.3%	2.8%	5.6%
Single-input	ChopKnlTree+attn+copy	34.5	17.8%	2.8%	6.3%
	ChopSynTree+attn+copy	36.3	20.5%	6.9%	11.6%
	+attn KnlTree+attn+copy	29.6	12.1%	1.5%	5.9%
	+copy SynTree+attn+copy	32.5	17.6%	2.5%	5.1%
	Stmt+attn+copy	37.4	24.2%	6.1%	10.2%
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn	20.3	7.7%	0.8%	1.3%
	Stmt+ChopKnlTree+attn	22.1	10.2%	0.8%	1.5%
	Stmt+ChopSynTree+attn	20.7	11.2%	1.7%	2.5%
	ChopKnlTree+ChopSynTree+attn	19.2	7.7%	1.3%	2.0%
Single-input	ChopKnlTree+attn	18.9	7.0%	0.5%	0.7%
	ChopSynTree+attn	24.2	10.4%	0.8%	2.1%
	+attn KnlTree+attn	12.7	2.2%	0.0%	0.2%
	+attn SynTree+attn	12.2	1.5%	0.0%	0.0%
	Stmt+attn	24.4	12.5%	0.8%	1.2%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	16.2	4.3%	0.2%	0.3%
	Stmt+ChopKnlTree	18.5	5.7%	0.3%	0.8%
	Stmt+ChopSynTree	10.9	0.1%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	12.0	0.7%	0.0%	0.0%
Single-input	ChopKnlTree	12.5	0.6%	0.0%	0.0%
	ChopSynTree	13.5	0.1%	0.0%	0.0%
	KnlTree	11.0	0.0%	0.0%	0.0%
	SynTree	10.6	0.0%	0.0%	0.0%
	Stmt	17.6	5.0%	0.0%	0.5%
-	Retrieval-based	23.7	11.6%	2.5%	3.0%

Table 14: Results of ROOSTERIZE Models with Training and Validation Sets from Tier 1 and Testing Set from All Tiers.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn+copy	43.3	21.4%	6.8%	14.9%
	Stmt+ChopKnlTree+attn+copy	44.5	23.9%	8.5%	16.2%
	Stmt+ChopSynTree+attn+copy	36.6	17.5%	5.6%	10.1%
	ChopKnlTree+ChopSynTree+attn+copy	43.3	22.3%	6.8%	13.8%
Single-input	ChopKnlTree+attn+copy	40.7	19.2%	4.5%	10.6%
	ChopSynTree+attn+copy	38.8	18.3%	6.4%	11.6%
	+attn KnlTree+attn+copy	35.1	13.1%	2.0%	7.3%
	+copy SynTree+attn+copy	30.9	10.9%	2.6%	5.9%
	Stmt+attn+copy	38.0	19.1%	6.4%	10.9%
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn	23.8	8.5%	0.4%	0.9%
	Stmt+ChopKnlTree+attn	24.6	8.4%	0.9%	1.7%
	Stmt+ChopSynTree+attn	23.2	8.0%	0.7%	1.5%
	ChopKnlTree+ChopSynTree+attn	27.1	10.5%	1.7%	3.2%
Single-input	ChopKnlTree+attn	18.9	4.7%	0.5%	1.2%
	ChopSynTree+attn	27.9	11.9%	1.6%	2.8%
	KnlTree+attn	13.7	1.5%	0.0%	0.0%
	SynTree+attn	8.8	1.0%	0.0%	0.0%
	Stmt+attn	26.0	10.8%	1.2%	2.4%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	17.4	3.4%	0.1%	0.2%
	Stmt+ChopKnlTree	19.1	4.4%	0.1%	0.2%
	Stmt+ChopSynTree	12.6	0.6%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	16.2	2.2%	0.0%	0.0%
Single-input	ChopKnlTree	15.2	1.6%	0.0%	0.0%
	ChopSynTree	14.4	0.8%	0.1%	0.1%
	KnlTree	12.2	0.6%	0.0%	0.0%
	SynTree	5.7	0.3%	0.0%	0.0%
	Stmt	19.4	4.6%	0.1%	0.3%
-	Retrieval-based	26.9	9.6%	0.3%	0.4%

Table 15: Results of ROOSTERIZE Models with Training and Validation Sets from Tier 1 and Testing Set from Tier 2.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn+copy	34.2	22.8%	9.2%	11.7%
	Stmt+ChopKnlTree+attn+copy	35.0	26.7%	8.3%	14.2%
	Stmt+ChopSynTree+attn+copy	36.7	14.2%	6.7%	14.2%
	ChopKnlTree+ChopSynTree+attn+copy	33.9	25.8%	9.2%	16.7%
Single-input	ChopKnlTree+attn+copy	34.8	25.0%	8.3%	14.2%
	ChopSynTree+attn+copy	35.6	20.1%	7.5%	15.8%
	KnlTree+attn+copy	25.3	11.2%	3.3%	7.5%
	SynTree+attn+copy	33.3	14.2%	2.5%	6.7%
	Stmt+attn+copy	42.4	20.8%	6.7%	15.8%
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn	24.8	12.8%	1.7%	5.8%
	Stmt+ChopKnlTree+attn	25.5	13.2%	5.0%	6.7%
	Stmt+ChopSynTree+attn	22.8	10.4%	0.0%	5.0%
	ChopKnlTree+ChopSynTree+attn	25.6	15.6%	5.8%	7.5%
Single-input	ChopKnlTree+attn	19.4	6.7%	0.8%	3.3%
	ChopSynTree+attn	28.0	17.5%	6.7%	7.5%
	KnlTree+attn	13.1	1.4%	0.0%	0.8%
	SynTree+attn	10.9	2.1%	0.0%	0.0%
	Stmt+attn	28.5	17.9%	5.0%	6.7%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	20.4	10.4%	0.8%	2.5%
	Stmt+ChopKnlTree	17.9	6.7%	0.0%	0.0%
	Stmt+ChopSynTree	14.3	0.4%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	15.6	4.6%	0.0%	0.0%
Single-input	ChopKnlTree	15.3	3.6%	0.0%	0.0%
	ChopSynTree	14.9	2.2%	0.0%	0.0%
	KnlTree	12.2	0.4%	0.0%	0.0%
	SynTree	6.6	0.0%	0.0%	0.0%
	Stmt	20.9	9.9%	2.5%	4.2%
-	Retrieval-based	24.8	14.6%	2.5%	7.5%

Table 16: Results of ROOSTERIZE Models with Training and Validation Sets from Tier 1 and Testing Set from Tier 3.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input +attn +copy	Stmt+ChopKnlTree+ChopSynTree+attn+copy	31.2	15.9%	2.1%	5.4%
	Stmt+ChopKnlTree+attn+copy	29.1	16.9%	1.3%	4.5%
	Stmt+ChopSynTree+attn+copy	29.8	13.8%	2.1%	5.8%
	ChopKnlTree+ChopSynTree+attn+copy	31.3	17.9%	1.7%	3.5%
Single-input +attn +copy	ChopKnlTree+attn+copy	27.6	13.8%	1.0%	2.6%
	ChopSynTree+attn+copy	33.1	17.8%	3.8%	6.8%
	KnlTree+attn+copy	24.3	6.4%	0.0%	1.2%
	SynTree+attn+copy	30.4	11.1%	1.8%	4.3%
	Stmt+attn+copy	30.9	16.7%	3.1%	5.8%
Multi-input +attn	Stmt+ChopKnlTree+ChopSynTree+attn	19.5	7.6%	0.3%	0.3%
	Stmt+ChopKnlTree+attn	18.3	6.8%	0.2%	0.3%
	Stmt+ChopSynTree+attn	19.4	6.3%	0.3%	0.7%
	ChopKnlTree+ChopSynTree+attn	18.9	7.4%	0.3%	0.8%
Single-input +attn	ChopKnlTree+attn	15.0	2.8%	0.0%	0.2%
	ChopSynTree+attn	21.3	9.4%	1.2%	1.2%
	KnlTree+attn	10.7	0.4%	0.0%	0.0%
	SynTree+attn	7.7	0.5%	0.0%	0.0%
	Stmt+attn	19.5	7.5%	0.5%	0.8%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	14.9	1.6%	0.0%	0.0%
	Stmt+ChopKnlTree	16.5	3.8%	0.0%	0.0%
	Stmt+ChopSynTree	12.4	0.5%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	13.4	1.0%	0.0%	0.0%
Single-input	ChopKnlTree	13.1	1.2%	0.0%	0.0%
	ChopSynTree	13.5	0.7%	0.0%	0.0%
	KnlTree	12.9	1.0%	0.0%	0.0%
	SynTree	5.8	0.0%	0.0%	0.0%
	Stmt	15.7	3.2%	0.0%	0.0%
-	Retrieval-based	18.3	5.7%	0.0%	0.0%

Table 17: Results of ROOSTERIZE Models with Training and Validation Sets from Tier 2 and Testing Set from Tier 2.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn+copy	31.8	17.6%	5.8%	9.2%
	Stmt+ChopKnlTree+attn+copy	29.5	17.5%	4.2%	8.3%
	Stmt+ChopSynTree+attn+copy	28.8	9.6%	0.0%	2.5%
	ChopKnlTree+ChopSynTree+attn+copy	30.5	15.8%	5.0%	10.8%
Single-input	ChopKnlTree+attn+copy	33.6	18.2%	4.2%	5.8%
	ChopSynTree+attn+copy	29.3	9.9%	0.8%	5.0%
	+attn KnlTree+attn+copy	31.2	17.2%	4.2%	7.5%
	+copy SynTree+attn+copy	25.0	2.8%	0.0%	0.0%
	Stmt+attn+copy	29.7	15.4%	2.5%	5.8%
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn	22.3	10.0%	1.7%	1.7%
	Stmt+ChopKnlTree+attn	21.4	8.1%	0.8%	2.5%
	Stmt+ChopSynTree+attn	24.3	14.2%	1.7%	7.5%
	ChopKnlTree+ChopSynTree+attn	18.0	6.9%	0.8%	3.3%
Single-input	ChopKnlTree+attn	15.7	1.1%	0.0%	0.0%
	ChopSynTree+attn	17.5	3.5%	0.0%	0.0%
	KnlTree+attn	11.2	0.0%	0.0%	0.0%
	SynTree+attn	11.3	0.0%	0.0%	0.0%
	Stmt+attn	20.6	7.8%	0.8%	6.7%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	15.5	3.5%	0.0%	0.0%
	Stmt+ChopKnlTree	15.3	2.6%	0.0%	0.0%
	Stmt+ChopSynTree	11.4	0.3%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	11.1	0.0%	0.0%	0.0%
Single-input	ChopKnlTree	9.7	1.4%	0.0%	0.0%
	ChopSynTree	12.8	0.0%	0.0%	0.0%
	KnlTree	8.7	0.0%	0.0%	0.0%
	SynTree	6.8	0.0%	0.0%	0.0%
	Stmt	14.6	1.9%	0.0%	0.0%
-	Retrieval-based	27.7	25.2%	0.0%	0.0%

Table 18: Results of ROOSTERIZE Models with Training and Validation Sets from Tier 3 and Testing Set from Tier 3.

Group	Model	BLEU	Frag.Acc.	Top1	Top5
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn+copy	33.5	19.5%	2.3%	5.9%
	Stmt+ChopKnlTree+attn+copy	31.2	16.9%	0.8%	6.3%
	Stmt+ChopSynTree+attn+copy	32.0	17.3%	3.3%	7.6%
	ChopKnlTree+ChopSynTree+attn+copy	31.8	15.5%	1.3%	5.0%
Single-input	ChopKnlTree+attn+copy	29.9	15.7%	1.3%	6.1%
	ChopSynTree+attn+copy	32.7	15.3%	4.0%	6.3%
	KnlTree+attn+copy	30.7	12.4%	1.3%	5.0%
	SynTree+attn+copy	28.4	12.6%	2.0%	6.4%
	Stmt+attn+copy	34.3	19.8%	3.3%	5.3%
Multi-input	Stmt+ChopKnlTree+ChopSynTree+attn	14.5	1.4%	0.0%	0.0%
	Stmt+ChopKnlTree+attn	17.9	5.3%	0.0%	0.0%
	Stmt+ChopSynTree+attn	19.0	6.6%	0.2%	0.2%
	ChopKnlTree+ChopSynTree+attn	12.4	2.2%	0.0%	0.0%
Single-input	ChopKnlTree+attn	14.8	2.3%	0.0%	0.0%
	ChopSynTree+attn	17.0	2.8%	0.0%	0.0%
	KnlTree+attn	13.1	0.2%	0.0%	0.0%
	SynTree+attn	4.8	0.3%	0.0%	0.0%
	Stmt+attn	17.6	5.8%	0.0%	0.0%
Multi-input	Stmt+ChopKnlTree+ChopSynTree	12.6	0.7%	0.0%	0.0%
	Stmt+ChopKnlTree	14.4	1.4%	0.0%	0.0%
	Stmt+ChopSynTree	9.1	0.2%	0.0%	0.0%
	ChopKnlTree+ChopSynTree	11.1	0.0%	0.0%	0.0%
Single-input	ChopKnlTree	11.7	0.1%	0.0%	0.0%
	ChopSynTree	10.1	0.0%	0.0%	0.0%
	KnlTree	14.3	0.2%	0.0%	0.0%
	SynTree	14.4	0.2%	0.0%	0.0%
	Stmt	14.4	1.4%	0.0%	0.0%
-	Retrieval-based	22.1	9.6%	2.5%	3.0%

Table 19: Results of the Generalization Study with ROOSTERIZE Pre-trained on All Tiers.

#Lemmas	BLEU	Frag.Acc.	Top1	Top5
0	33.9	21.3%	4.4%	8.9%
105	32.6	21.5%	3.3%	5.3%
223	34.1	22.7%	3.8%	6.9%
505	35.7	24.3%	5.0%	8.7%
580	37.4	26.5%	7.4%	12.5%