

TypeWriter: Neural Type Prediction with Search-Based Validation

Pradel, Michael; Gousios, Georgios; Liu, Jason; Chandra, Satish

DOI

[10.1145/3368089.3409715](https://doi.org/10.1145/3368089.3409715)

Publication date

2020

Document Version

Accepted author manuscript

Published in

Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering

Citation (APA)

Pradel, M., Gousios, G., Liu, J., & Chandra, S. (2020). TypeWriter: Neural Type Prediction with Search-Based Validation. In P. Devanbu, M. Cohen, & T. Zimmermann (Eds.), *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (pp. 209–220). (ESEC/FSE 2020). Association for Computing Machinery (ACM).
<https://doi.org/10.1145/3368089.3409715>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

TypeWriter: Neural Type Prediction with Search-based Validation

Michael Pradel*
University of Stuttgart
michael@binaervarianz.de

Jason Liu
jasonliu@fb.com
Facebook

Georgios Gousios*
Delft University of Technology
g.gousios@tudelft.nl

Satish Chandra
schandra@acm.org
Facebook

ABSTRACT

Maintaining large code bases written in dynamically typed languages, such as JavaScript or Python, can be challenging due to the absence of type annotations: simple data compatibility errors proliferate, IDE support is limited, and APIs are hard to comprehend. Recent work attempts to address those issues through either static type inference or probabilistic type prediction. Unfortunately, static type inference for dynamic languages is inherently limited, while probabilistic approaches suffer from imprecision. This paper presents TypeWriter, the first combination of probabilistic type prediction with search-based refinement of predicted types. TypeWriter’s predictor learns to infer the return and argument types for functions from partially annotated code bases by combining the natural language properties of code with programming language-level information. To validate predicted types, TypeWriter invokes a gradual type checker with different combinations of the predicted types, while navigating the space of possible type combinations in a feedback-directed manner.

We implement the TypeWriter approach for Python and evaluate it on two code corpora: a multi-million line code base at Facebook and a collection of 1,137 popular open-source projects. We show that TypeWriter’s type predictor achieves an F1 score of 0.64 (0.79) in the top-1 (top-5) predictions for return types, and 0.57 (0.80) for argument types, which clearly outperforms prior type prediction models. By combining predictions with search-based validation, TypeWriter can fully annotate between 14% to 44% of the files in a randomly selected corpus, while ensuring type correctness. A comparison with a static type inference tool shows that TypeWriter adds many more non-trivial types. TypeWriter currently suggests types to developers at Facebook and several thousands of types have already been accepted with minimal changes.

1 INTRODUCTION

Dynamically typed programming languages, such as Python and JavaScript, have become extremely popular, and large portions of newly written code are in one of these languages. While the lack of static type annotations enables fast prototyping, it often leads to problems when projects grow. Examples include type errors that remain unnoticed for a long time [9], suboptimal IDE support, and difficult to understand APIs [22]. To solve these problems, in recent years, many dynamic languages obtained support for *type annotations*, which enable programmers to specify types in a fashion similar to a statically typed language. Type annotations

are usually ignored at runtime; nevertheless, they serve both as hints for developers using external APIs and as inputs to gradual type checkers that ensure that specific programming errors cannot occur. To cope with legacy code bases, type annotations can be introduced gradually; in such cases, the type checker will check only code that is annotated.

As manually annotating code is time-consuming and error-prone [26], developers must resort to automated methods. One way to address the lack of type annotations is type inference via traditional static analysis. Unfortunately, dynamic features, such as heterogeneous arrays, polymorphic variables, dynamic code evaluation, and monkey patching make static type inference a hard problem for popular dynamic languages, such as Python or JavaScript [7]. Static type inference tools typically handle these challenges by inferring a type only if it is certain or very likely (under some assumptions), which significantly limits the number of types that can be inferred.

Motivated by the inherent difficulties of giving definitive answers via static analysis, several probabilistic techniques for predicting types have been proposed. A popular direction is to exploit the existence of already annotated code as training data to train machine learning models that then predict types in not yet annotated code. Several approaches predict the type of a code entity, e.g., a variable or a function, from the code contexts in which this entity occurs [15, 28]. Other approaches exploit natural language information embedded in source code, e.g., variable names or comments, as a valuable source of informal type hints [21, 37].

While existing approaches for predicting types are effective in some scenarios, they suffer from *imprecision* and *combinatorial explosion*. Probabilistic type predictors are inherently imprecise because they suggest one or more likely types for each missing annotation, but cannot guarantee their correctness. The task of deciding which of these suggestions are correct is left to the developer. Because probabilistic predictors suggest a ranked list of likely types, choosing a type-correct combination of type annotations across multiple program elements causes combinatorial explosion. A naïve approach would be to let a developer or a tool choose from all combinations of the predicted types. Unfortunately, this approach does not scale to larger code examples, because the number of type combinations to consider is exponential in the number of not yet annotated code entities.

This paper presents TypeWriter, a combination of learning-based, probabilistic type prediction and a feedback-directed, search-based validation of predicted types. The approach addresses the imprecision problem based on the insight that a gradual type checker can

*Work performed while on sabbatical at Facebook, Menlo Park.

```

1 # Predicted argument type: int, str, bool
2 # Predicted return type: str, Optional[str], None
3 def find_match(color):
4     """
5     Args:
6     color (str): color to match on and return
7     """
8     candidates = get_colors()
9     for candidate in candidates:
10        if color == candidate:
11            return color
12    return None
13
14 # Predicted return types: List[str], List[Any], str
15 def get_colors():
16    return ["red", "blue", "green"]

```

Figure 1: Example of search for type-correct predicted types.

pinpoint contradictory type annotations, which guides the selection of suitable types from the set of predicted types. To make the search for a consistent set of types tractable, we formulate the problem as a combinatorial search and present a search strategy that finds type-correct type annotations efficiently. TypeWriter makes use of the variety of type hints present in code through a novel neural architecture that exploits both natural language, in the form of identifier names and code comments, similar to prior work [21], and also programming context, in the form of usage sequences.

To illustrate the approach, consider the two to-be-typed functions in Figure 1. Given this code, the neural type model of TypeWriter predicts a ranked list of likely types for each argument type and return type, as indicated by the comments. TypeWriter starts by adding the top-ranked predictions as type annotations, which introduces a type error about an incorrect return type of `find_match`, though. Based on this feedback, the search tries to change the return type of `find_match` to the second-best suggestion, `Optional[str]`. Unfortunately, this combination of added types leads to another type error because the return type is inconsistent with the argument key being of type `int`. The search again refines the type annotations by trying to use the second-best suggestion, `str`, for the argument key. Because the resulting set of type annotations is type-correct according to the type checker, TypeWriter adds these types to the code.

We implement TypeWriter for Python and apply it on two large code bases: a multi-million line code base at Facebook that powers applications used by billions of people, and a corpus of popular open-source projects. We show that the neural model predicts individual types with a precision of 64% (85%, 91%) and a recall of 52% (64%, 68%) within the top-1 (top-3, top-5) predictions, which outperforms a recent, closely related approach [21] by 10% and 6% respectively. Based on this model, the feedback-directed search finds a type-correct subset of type annotations that can produce complete and type-correct annotations for 42% to 64% of all files. Comparing TypeWriter with a traditional, static type inference shows that both techniques complement each other and that TypeWriter predicts many more types than traditional type inference. In summary, this paper makes the following contributions:

- A combination of probabilistic type prediction and search-based validation of predicted types. The feedback-directed search for type-correct types can be used with any probabilistic type predictor and any gradual type checker.
- A novel neural type prediction model that exploits both code context and natural language information.
- Empirical evidence that the approach is effective for type-annotating large-scale code bases with minimal human effort. The initial experience from using TypeWriter at Facebook on a code base that powers tools used by billions of people has been positive.

2 APPROACH

Figure 2 gives a high-level overview of the TypeWriter approach. The input to TypeWriter is a corpus of code where some, but not all types are annotated. The approach consists of three main parts. First, a lightweight static analysis extracts several kinds of information from the given code (Section 2.1). The extracted information includes programming structure information, such as usages of a function’s arguments, and natural language information, such as identifier names and comments. Next, a neural type predictor learns from the already annotated types and their associated information how to predict missing types (Section 2.2). Once trained, this model can predict likely types for currently unannotated parts of the code. Finally, a feedback-directed search uses the trained model to find a type assignment that is consistent and type-correct according to a static, gradual type checker (Section 2.3). The overall output of TypeWriter is code with additional type annotations.

2.1 Static Extraction of Types and Context Information

The first part of TypeWriter is an AST-based static analysis that extracts types and context information useful to predict types. The analysis is designed to be lightweight and easy to apply to other programming languages. We currently focus on function-level types, i.e., argument types and return types. These types are particularly important for two reasons: (i) Given function-level types, gradual type checkers can type-check the function bodies by inferring the types of (some) local variables. (ii) Function-level types serve as interface documentation.

For each type, the static analysis gathers four kinds of context information, which the following describes and illustrates with the example in Figure 3. Each of the four kinds of information may provide hints about an argument type or return type, and our model learns to combine these hints into a prediction of the most likely types.

Identifier names associated with the to-be-typed program element. As shown by prior work [21, 27], natural language information embedded in source code can provide valuable hints about program properties. For example, the argument names `name` and `do_propagate` in Figure 3 suggest that the arguments may be a string and a boolean, respectively. To enable TypeWriter to benefit from such hints, the static analysis extracts the identifier name of each with function and each function argument.

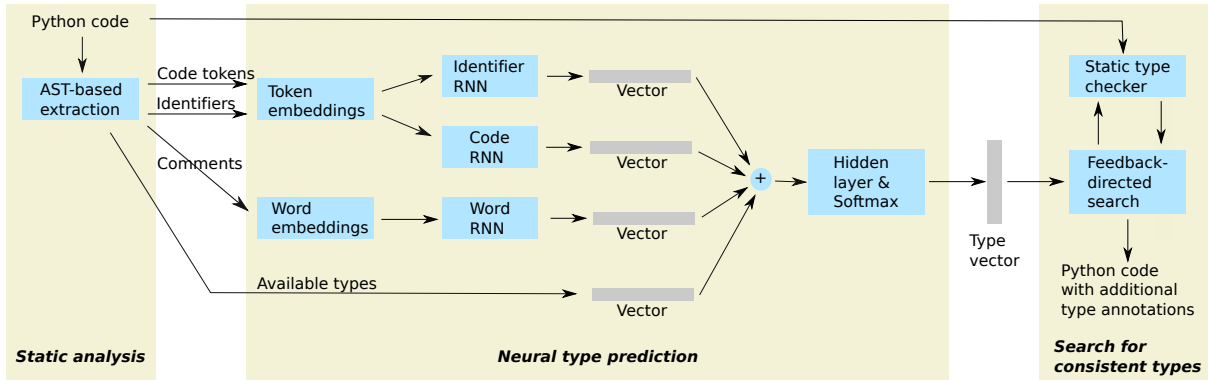


Figure 2: Overview of TypeWriter.

```

1 from html import HTMLElement
2
3 def update_name(name, do_propagate, element):
4     """ Update the name and (optionally)
5     propagate to dependents. """
6     first_name = name.split(" ")[0]
7     element.first = first_name
8     if do_propagate:
9         for d in dependents:
10            d.notify(NAME, first_name)

```

Figure 3: Example of a to-be-typed Python function.

Code occurrences of the to-be-typed program element. In addition to the above natural language information, TypeWriter exploits programming language type hints. One of them is the way a to-be-typed program element is used: As a hint about argument types, the analysis considers all usages of an argument within the function body. Another kind of information is code that defines the to-be-typed program element: As a hint about return types, the analysis considers all return statements in a function. For each of these code locations, the analysis extracts the corresponding sequence of code tokens o_1, \dots, o_k . Specifically, the analysis extracts a window of tokens around each occurrence of an argument (default size of window: 7) and all tokens of a return statement. For example, the analysis extracts the token sequence `\n, first_name, =, name, ., split, (` around the usage of `name` at line 6.

As an alternative to extracting token sequences, TypeWriter could perform a more sophisticated static analysis, e.g., by tracking data flows starting at arguments or ending in return values. We instead focus on token sequences because it provides a sufficiently strong signal, scales well to large code bases, and could be easily ported to other programming languages.

Function-level comments. Similar to identifier names, comments are another informal source of hints about types. For the example in Figure 3, a developer might infer from the function-level comment that the function has some side effects but probably does not return any value. To allow the approach to benefit from such hints, the static analysis extracts all function-level comments, i.e., docstrings in Python. For a given function, the approach uses this comment

both for predicting the argument types and the return type of the function.

Available types. To annotate a type beyond the built-in types of Python, the type needs to be either imported or locally defined. Because types used in an annotation are likely to be already imported or locally defined, the analysis extracts all types available in a file. To this end, the analysis parses all `import` statements and all class definitions in a file. For the example in Figure 3, the analysis will extract `HTMLElement` as an available type, which hints at the argument `element` being of this type.

Based on these four kinds of type hints, the analysis extracts the following information for argument types and return types, respectively:

Definition 2.1 (Argument type information). For a function argument a , the statically extracted information is a tuple

$$(n_{fct}, n_{arg}, N_{args}, c, U, t)$$

where n_{fct} is the function name, n_{arg} is the name of the argument a , N_{args} is the sequence of names of other arguments (if any), c is the comment associated with the function, U is a set of usage sequences, each of which is a sequence o_1, \dots, o_k of tokens, and t is the type of the argument.

Definition 2.2 (Return type information). For the return type of a function f , the statically extracted information is a tuple

$$(n_{fct}, N_{args}, c, R, t)$$

where n_{fct} is the function name, N_{args} is the sequence of argument names, c is the comment associated with the function, R is a set of return statements, each of which is a sequence o_1, \dots, o_k of tokens, and t is the return type of f .

If any of the above information is missing, the corresponding elements of the tuple is filled with a placeholder. In particular, the static analysis extracts the above also for unannotated types, to enable TypeWriter to predict types based on the context.

2.2 Neural Type Prediction Model

Given the extracted types and context information, the next part of TypeWriter is a neural model that predicts the former from the

latter. We formulate the type prediction problem as a classification problem, where the model predicts a probability distribution over a fixed set of types. The neural type prediction model, summarized in the middle part of Figure 2, combines the four kinds of information described in Section 2.1 into a single type prediction.

To represent identifier names, source code tokens, and words in a way suitable for learning, TypeWriter maps each into a real-valued vector using a Word2Vec [23] embedding. We train two embeddings, a *code embedding* E_{code} for code tokens and identifier names, and a *word embedding* E_{word} for words in comments. E_{code} is trained on sequences of tokens extracted from source code files, while E_{word} is trained on sequences of words extracted from comments. To mitigate the problem of large vocabularies in source code [5], TypeWriter preprocesses each identifier using a helper function $norm()$, which tokenizes, lemmatizes, and lowercases each identifier.

2.2.1 Learning from Identifiers. This neural submodel learns from the identifier names of functions and function arguments. The model combines all identifiers associated with a type into sequence. Given argument type information $(n_{fct}, n_{arg}, N_{args}, c, U, t)$, the sequence is

$$norm(n_{arg}) \circ s \circ norm(n_{fct}) \circ norm(N_{args})$$

where \circ flattens and concatenates sequences, and s is a separator. Given return type information $(n_{fct}, N_{args}, c, R, t)$, the sequence is

$$norm(n_{fct}) \circ s \circ norm(N_{args})$$

For example, the sequence for the return type of the function in Figure 3 is “update name s name do propagate element”.

TypeWriter learns from these sequences of words by summarizing them into a single vector using a bi-directional, recurrent neural network (RNN) based on LSTM cells. To ease parallelization, we pad sequences that are too short and truncate sequences that are too long (default length: 10). The final hidden states of the RNN serve as a condensed vector representation, v_{ids} , of all identifier-related hints.

2.2.2 Learning from Token Sequences. This neural submodel learns from source code information associated with a type. Similar to the submodel for identifiers, this submodel composes all relevant tokens into a sequence and summarize them into a single vector v_{code} using an RNN. For arguments and return types, the sequence consists of the tokens involved in the usages U (Definition 2.1) and the return statements R (Definition 2.2), respectively. Before feeding these sequences into an RNN, we bound the length of each token sequence (default: $k = 7$) and of the number of token sequences (default: 3).

2.2.3 Learning from Comments. This neural submodel learns type hints from comments associated with a function. To this end, TypeWriter splits a given comment into a sequence of words, bound the length of the sequence to a fixed value (default: 20), and summarizes the sequence via another RNN. The result is a fixed-length vector $v_{comments}$.

2.2.4 Learning from Available Types. The fourth kind of information that TypeWriter learns from is the set of types available in the

current source code file. The approach assumes a fixed-size vocabulary T of types (default size: 1,000). This vocabulary covers the vast majority of all type occurrences because most type annotations either use one of the built-in primitive types, e.g., `str` or `bool`, common non-primitive types, e.g., `List` or `Dict`, or their combinations, e.g., `List[int]` or `Dict[str, bool]`. Any types beyond the type vocabulary are represented as a special “unknown” type.

To represent which types are available, we use a binary vector of size T , called the *type mask*. Each element in this vector represents one type, and an element is set to one if and only if its type is present. The resulting vector $v_{availTypes}$ of available types is passed as-is into the final part of the neural model.

2.2.5 Predicting the Most Likely Type. The final submodel concatenates the four vectors v_{ids} , v_{code} , $v_{comments}$, and $v_{availTypes}$ into a single vector and passes it through a fully connected layer that predicts the most likely type. The output vector has size $|T|$ and represents a probability distribution over the set of types. For example, suppose the type vocabulary had only four types `int`, `bool`, `None`, and `List`, and that the output vector is $[0.1, 0.6, 0.2, 0.1]$. In this case, the model would predict that `bool` is the most likely type, following by `None`.

There are two ways to handle uncertainty and limited knowledge in the model. First, we interpret the predicted probability of a type as a confidence measure and only suggest types to a user that are predicted with a confidence above some configurable threshold. Second, we encode types not included in the fixed-size type vocabulary as a special “unknown” type. The model hence learns to predict “unknown” whenever none of the types in the vocabulary fit the given context information. During prediction, TypeWriter never suggests the “unknown” type to the user, but instead does not make any suggestion in case the model predicts “unknown”.

2.2.6 Training. To train the type prediction model, TypeWriter relies on already type-annotated code. Given such code, the approach creates one pair of context information and type for each argument type and for each return type. These pairs then serve as training data to tune the parameters of the different neural submodels. We use stochastic gradient descent, the Adam optimizer, and cross-entropy as the loss function. The entire neural model is learned jointly, enabling the model to summarize each kind of type hint into the most suitable form and to decide which type hints to consider for a given query. We train two separate models for argument types and function types, each learned from training data consisting of only one kind of type. The rationale is that some of the available type hints need to be interpreted differently depending on whether the goal is to predict an argument type or a return type.

2.3 Feedback-guided Search for Consistent Types

The neural type prediction model provides a ranked list of k predictions for each missing type annotation. Given a set of locations for which a type annotation is missing, called *type slots*, and a list of probabilistic predictions for each slot, the question is which of the suggested types to assign to the slots. A naïve approach might fill each slot with the top-ranked type. However, because the neural model may mis-predict some types, this approach may yield type

Algorithm 1 Find a correct type assignment for a file f

```

1: function ASSIGN_TYPES( $f$ )
2:    $T \leftarrow$  all type slots in  $f$ 
3:    $\mathcal{P}_t^{1..k} \leftarrow$  {predictions( $t, k$ ) |  $t \in T$ }   $\triangleright$  Top  $k$  predictions
4:    $a \leftarrow$  { $\mathcal{P}_t^1$  |  $t \in T$ }   $\triangleright$  Initial type assignment
5:    $a.score \leftarrow$  typecheck( $a, f$ )   $\triangleright$  Feedback function
6:   work_set  $\leftarrow$  new_states( $a, \mathcal{P}, T$ )
7:   done  $\leftarrow$  { $a$ }
8:   while min( $\{x.score \mid x \in done\}$ )  $> 0 \wedge$  work_set  $\neq \emptyset$  do
9:      $a \leftarrow$  pick(work_set)   $\triangleright$  Biased random selection
10:     $a.score \leftarrow$  typecheck( $a, f$ )
11:    if greedy  $\wedge$   $a.score < a.parent.score$  then
12:      work_set  $\leftarrow$  new_states( $a, \mathcal{P}, T$ )
13:    else if non-greedy then
14:      work_set  $\leftarrow$  work_set  $\cup$  (new_states( $a, \mathcal{P}, T$ ) \
done)
15:    end if
16:    done  $\leftarrow$  done  $\cup$  { $a$ }
17:  end while
18:  return argmin( $\{x.score \mid x \in done\}$ )
19: end function
20: function NEW_STATES( $a, \mathcal{P}, T$ )
21:  children  $\leftarrow$  {}
22:  for all  $t \in T$  do
23:    for all  $\mathcal{P}_t^j$  where  $j >$  rank of current  $a[t]$  do
24:       $a_{child} \leftarrow$  modify  $a$  to use  $\mathcal{P}_t^j$  at  $t$ 
25:      children  $\leftarrow$  { $a_{child}$ }
26:    end for
27:     $a_{child} \leftarrow$  modify  $a$  to not use any type at  $t$ 
28:    children  $\leftarrow$  { $a_{child}$ }
29:  end for
30:  return children
31: end function

```

assignments where the added annotations are not consistent with each other or with the remaining program.

To avoid introducing type errors, TypeWriter leverages an existing gradual type checker as a filter to validate candidate type assignments. Such type checkers exist for all popular dynamically typed languages that support optional type annotations, e.g., pyre and mypy for Python, and flow for JavaScript. TypeWriter exploits feedback from the type checker to guide a search for consistent types, as presented in Algorithm 1 and explained in the sections below.

2.3.1 Search Space. Given a set T of type slots and k predicted types for each slot, we formulate the problem of finding a consistent type assignment as a combinatorial search problem. The search space consists of the set \mathcal{P} of possible *type assignments*. For $|T|$ type slots and k possible types for each slot, there are $(k + 1)^{|T|}$ type assignments (the +1 is for not assigning any of the predicted types).

2.3.2 Feedback Function. Exhaustively exploring \mathcal{P} is practically infeasible for files with many missing types, because invoking the gradual type checker is relatively expensive (typically, in the order

of several seconds per file). Instead, TypeWriter uses a feedback function (typecheck) to efficiently steer toward the most promising type assignments.

The feedback function is based on two values, both of which the search wants to minimize:

- $n_{missing}$: The number of missing types.
- n_{errors} : The number of type errors.

TypeWriter combines these into a weighted sum $score = v \cdot n_{missing} + w \cdot n_{errors}$. By default, we set v to 1 and w to the number of initially missing types plus one, which is motivated by the fact that adding an incorrect type often leads to an additional error. By giving type errors a high enough weight, we ensure that the search never returns a type assignment that adds type errors to the code.

2.3.3 Exploring the Search Space. TypeWriter explores the space of type assignments through an optimistic search strategy (Algorithm 1). It assumes that most predictions are correct, and then refines type annotations to minimize the feedback function. Each exploration step explores a state a , which consists of a type assignment, the score computed by the feedback function, and a link to the parent state. The initial state is generated by retrieving the top-1 predictions from \mathcal{P} for each type slot t and invoking the feedback function (lines 4 and 5). The next states to be explored are added to a work set, while the explored states are kept in the “done” set. The algorithm loops over items in the work set until either the feedback score has been minimized or the search explored all potential type assignments (line 8). The assignment with the minimal score is returned as a result (line 18).

To retrieve the next type assignments to possibly explore from the current state, TypeWriter invokes the new_states helper function. It adds all type assignments that can be obtained from the current state by modifying exactly one type slot, either by using a lower-ranked type suggestion or by not adding any type for this slot (lines 22 to 29).

The main loop of the algorithm (lines 8 to 17) picks a next state to evaluate from the working set (line 9), queries the feedback function (line 10) and updates the done set with the explored state (line 16). The pick function is a biased random selection that prefers states based on two criteria. First, it prefers states that add more type annotations over states that add fewer annotations. Second, it prefers states that modify a type annotation at a line close to a line with a type error. Intuitively, such states are more likely to fix the cause of a type error than a randomly selected state.¹ The working set is updated with all new states that have not been currently explored.

TypeWriter implements two variants of the search, a greedy and a non-greedy one. The *greedy* strategy aggressively explores children of type assignment that decrease the feedback score and prunes children of states that increase it (line 12). The *non-greedy* performs no pruning, i.e., it can explore a larger part of the search space at the expense of time (line 14).

As an optimization of Algorithm 1, TypeWriter invokes the assign_types function twice. The first call considers only type slots for return types, whereas the second call considers all type slots for argument types. The reason for this two-phase approach is that

¹The reason for relying on line numbers as the interface between the type checker and TypeWriter is to enable plugging any type checker into our search.

many gradual type checkers, including `pyre`, the one used in our evaluation, type-check a function only if its return type is annotated. If `TypeWriter` would add argument type annotations before adding return type annotations, the feedback function might not include all type errors triggered by an incorrectly added argument annotation.

3 IMPLEMENTATION

The implementation of `TypeWriter` builds upon a variety of tools in the Python ecosystem. For the static analysis phase, we apply a data extraction pipeline consisting of Python’s own `ast` library to parse the code into an AST format, and `NLTK` and its `WordNetLemmatizer` module to perform standard NLP tasks (lemmatization, stop word removal). The pipeline is parallelized so that it handles multiple files concurrently. The neural network model is implemented in `PyTorch`. For obtaining embeddings for words and tokens, we pre-train a `Word2Vec` model using the `gensim` library. The search phase of `TypeWriter` builds upon the `LibCST`² to add types to existing Python files. We use `pyre` for static type checking. Our LSTM models all use 200-dimensional hidden layers, and we train for 10 epochs with a learning rate of 0.005 using the Adam Optimizer.

4 EVALUATION

We structure our evaluation along four research questions.

RQ 1: How effective is `TypeWriter`’s model at predicting argument and return types, and how does it compare to existing work?

RQ 2: How much do the different kinds of context information contribute to the model’s prediction abilities?

RQ 3: How effective is `TypeWriter`’s search?

RQ 4: How does `TypeWriter` compare to traditional static type inference?

4.1 Datasets

`TypeWriter` is developed and evaluated within Facebook. As the internal code base is not publicly available and to ensure that the presented results are replicable, we use two datasets:

Internal code base We collect Python from a large internal code repository.

OSS corpus We search GitHub for all projects tagged as `python3`. We also search `Libraries.io` for all Python projects that include `mypy` as a dependency. We then remove all projects that have less than 50 stars on GitHub, to ensure that the included projects are of substantial public interest. To ease future work to compare with `TypeWriter`, all results for the OSS corpus are available for download.³

The resulting dataset statistics can be found in Table 1. The internal dataset is much larger in size, but both datasets are comparable in terms of the percentage of annotated code. By restricting the type vocabulary to a fixed size, we exclude around 10% of all type occurrences for both datasets. This percentage is similar for both datasets, despite their different sizes, because types follow a long-tail distribution, i.e., relatively few types account for the majority of all type occurrences. We ignore some types because they are *trivial*

Table 1: Internal and open-source datasets.

Metric	Internal	OSS
Repositories	1	1,137
Files	*	11,993
Lines of code	*	2.7M
Functions	*	146,106
... with return type annotation	68%	80,341 (55%)
... with comment	21.8%	53,500 (39.3%)
... with both	16%	32,409 (22.2%)
... ignored because trivial	7.4%	12,436 (8.5%)
Arguments	*	274,425
... with type annotation	50%	112,409 (41%)
... ignored because trivial	33%	96,036 (35%)
Types - Return	*	7,383
... occurrences ignored (out of vocab.)	20.2%	11.3%
Types - Argument	*	8,215
... occurrences ignored (out of vocab.)	21.3%	13.7%
Training time (min:sec)		
... parsing	several minutes	1:45
... training embeddings	several minutes	2:29
... training neural model	several minutes	2:20

* = not available for disclosure

to predict, such as the return type of `__str__`, which always is `str`, or the type of the `self` argument of a method, which always is the surrounding class. `TypeWriter` could easily predict many of these trivial types, but a simple syntactic analysis would also be sufficient. We ignore trivial types for the evaluation to avoid skewing the results in favor of `TypeWriter`.

4.2 Examples

Figure 4 shows examples of successful and unsuccessful type predictions in the OSS dataset. Example 1 presents a case where `TypeWriter` correctly predicts a type annotation. Here, the code context and comments provide enough hints indicating that token is of type `Callable`. Example 2 presents a case where `TypeWriter` does not correctly predict the type, but the prediction is close to what is expected. We hypothesize that this case of mis-prediction is due to the fact that `TypeWriter` tries to associate associations between natural language and types, or in this case, the word “path” and the type `Path`.

4.3 RQ 1: Effectiveness of the Neural Model

Prediction tasks. To evaluate the neural type prediction, we define two prediction tasks: (i) `ReturnPrediction`, where the model predicts the return types of functions, and (ii) `ArgumentPrediction`, where the model predicts the types of function arguments, and

Metrics. We evaluate the effectiveness of `TypeWriter`’s neural type predictor by splitting the already annotated types in a given dataset into training (80%) and validation (20%) data. The split is by file, to avoid mixing up types within a single file. Once trained on the training data, we compare the model’s predictions against the validation data, using the already annotated types as the ground

²<https://github.com/Instagram/LibCST>

³<http://software-lab.org/projects/TypeWriter/data.tar.gz>

Table 2: Effectiveness of neural type prediction.

Corpus	Task	Model	Precision			Recall			F1-score		
			Top-1	Top-3	Top-5	Top-1	Top-3	Top-5	Top-1	Top-3	Top-5
Internal	ReturnPrediction	TypeWriter	0.73	0.88	0.92	0.58	0.66	0.69	0.64	0.76	0.79
		NL2Type	0.60	0.82	0.88	0.50	0.61	0.65	0.55	0.70	0.75
		DeepTyper	0.70	0.87	0.92	0.43	0.54	0.59	0.53	0.67	0.72
	ArgumentPrediction	Naïve baseline	0.15	0.22	0.25	0.24	0.40	0.45	0.18	0.29	0.32
		TypeWriter	0.64	0.86	0.92	0.52	0.66	0.70	0.57	0.75	0.80
		NL2Type	0.53	0.80	0.88	0.46	0.61	0.66	0.50	0.70	0.76
		DeepTyper	0.54	0.80	0.87	0.42	0.57	0.62	0.47	0.67	0.73
OSS	ReturnPrediction	TypeWriter	0.69	0.80	0.84	0.61	0.70	0.72	0.65	0.75	0.78
		NL2Type	0.61	0.74	0.79	0.55	0.64	0.68	0.58	0.69	0.73
		DeepTyper	0.52	0.79	0.83	0.48	0.59	0.64	0.50	0.66	0.72
	ArgumentPrediction	Naïve baseline	0.16	0.25	0.28	0.25	0.42	0.47	0.20	0.31	0.35
		TypeWriter	0.58	0.77	0.84	0.50	0.65	0.70	0.54	0.71	0.77
		NL2Type	0.50	0.71	0.79	0.46	0.61	0.66	0.48	0.66	0.72
		DeepTyper	0.51	0.76	0.84	0.45	0.59	0.64	0.48	0.67	0.73
		Naïve baseline	0.06	0.11	0.14	0.14	0.25	0.29	0.08	0.15	0.19

```

1 # PrefectHQ/ct/f/blob/master/src/prefect/utilities/notifications.py
2 # Commit: 864d44b
3 # Successful annotation of return type
4 def callback_factory(...) -> Callable:
5     """
6     ...
7     Returns:
8     - state_handler (Callable): a state handler function that
9       can be attached to both Tasks and Flows
10    ...
11    """
12    def state_handler(...):
13        ...
14    return state_handler

```

Example 1

```

1 # awslabs/sockeye/blob/master/sockeye/average.py
2 # Commit: bcda569
3 # Incorrect annotation of return type: expected List[str]
4 def find_checkpoints(...) -> List[Path]:
5     """
6     ...
7     :return: List of paths corresponding to chosen checkpoints.
8     """
9     ...
10    params_paths = [
11        os.path.join(model_path, C.PARAMS_NAME % point[-1])
12        for point in top_n
13    ]
14    ...
15    return params_paths

```

Example 2**Figure 4: Examples of successful and unsuccessful type predictions (GitHub: PrefectHQ/ct, awslabs/sockeye).**

truth. We compute precision, recall, and F1-score, weighted by the number of type occurrences in the dataset. Similarly to previous work [21], if the prediction model cannot predict a type for a type slot (i.e., returns “unknown”), we remove this type slot from the calculation of precision. Specifically, we calculate precision as

$prec = \frac{n_{corr}}{n_{all}}$, where n_{corr} is the number of correct predictions and n_{all} is the number of type slots for which the model does not return “unknown”. We calculate recall as $rec = \frac{n_{corr}}{|D|}$, where $|D|$ is total number of type slots in the examined dataset. We report the top- k scores, for $k \in \{1, 3, 5\}$.

Baseline models. We compare TypeWriter’s top- k predictions against three baseline models. The *naïve baseline* model considers the ten most frequent types in the dataset and samples its prediction from the distribution of these ten types, independently of the given context. For example, it predicts None as a return type more often than List[str] because None is used more often as a return type than List[str]. The *DeepTyper* baseline is a Python re-implementation of the DeepTyper [15] model. DeepTyper learns to translate a sequence of source code tokens to a sequence of types (and zeros for tokens without a type). To make it directly compatible with TypeWriter, we do not consider predictions for variable annotations in function bodies, even though we do perform basic name-based type propagation in case an annotated argument is used in a function body. Finally, the *NL2Type* baseline is a re-implementation of the NL2Type model [21] for Python, which also learns from natural language information associated with a type, but does not consider code context or available types.

Results. Table 2 presents the results for RQ 1. Our neural model achieves moderate to high precision scores, e.g., 73% in the top-1 and 92% in the top-5 for on the internal dataset for the ReturnPrediction task. The recall results are good but less high than precision, indicating that TypeWriter is fairly confident when it makes a prediction, but abstains from doing so when it is not. All models have slightly worse performance on the OSS dataset, which we attribute to the smaller size of that dataset. The fact that the top-3 and top-5 scores are significantly higher than top-1 in all cases motivates our work on combinatorial search (Section 4.5).

Compared to the baselines, TypeWriter outperforms both the state of the art and the naïve baseline across all metrics for both datasets and all three prediction tasks. The differences between

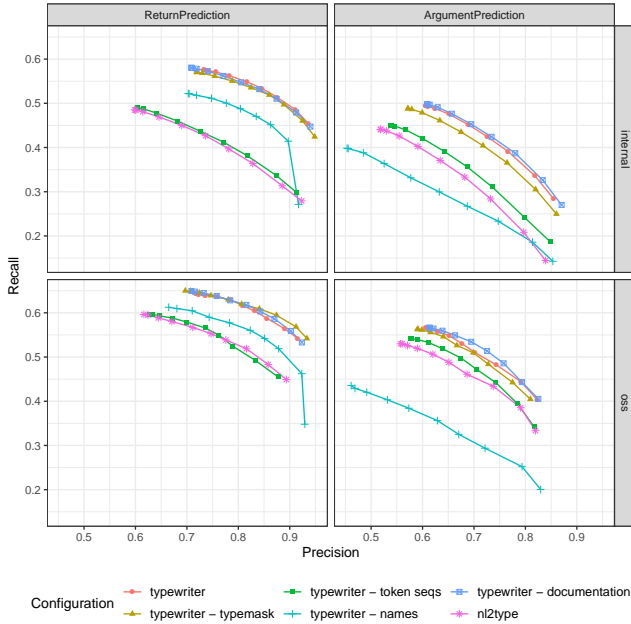


Figure 5: Precision/Recall curves for different TypeWriter-configurations. Each data point represents a prediction threshold level.

TypeWriter and NL2Type are higher in the case of the ReturnPrediction than the ArgumentPrediction task. The context information, as obtained by analyzing token sequences, is helping the TypeWriter prediction model more in the ReturnPrediction task. Compared to DeepTyper, both NL2Type and TypeWriter are better, the latter by a significant margin, in top-1 but not in top-3 or top-5 predictions. Given that all models learn primarily from identifier names, the relatively close upper bound performance seems to indicate that performance improvements may only be achieved by introducing different (e.g., structural) information to the model.

4.4 RQ 2: Comparison with Simpler Variants of the Neural Model

The main novelty of TypeWriter’s prediction component is the inclusion of code context information and a local type mask in the prediction model. To explore the influence of the different type hints considered by TypeWriter, we perform an ablation study. Specifically, we turn off parts of the model, both in training and in testing, and then measure top-1 precision and recall at various prediction threshold levels. We start with the full model (typewriter) and then we remove, in order, the type mask, the token sequences, the method and argument names and the documentation. As a baseline, we also include nl2type, a configuration functionally equivalent with NL2Type [21], which corresponds to TypeWriter without token sequences and without a type mask. The results of the ablation study can be seen in Figure 5.

Overall, the combined information of natural language, token sequences, and type masks helps TypeWriter to perform better than previous models, such as NL2Type. The main contributor to

this improvement is the token sequences component. Moreover, the results seem to reinforce the main thesis of NL2Type, i.e., that natural language information and types are strongly related: If we remove the argument and function naming information from TypeWriter, its performance drops significantly.

Contrary to our initial expectations, the type mask component is not contributing significantly in the ReturnPrediction task, while only slightly improving the ArgumentPrediction results. We attribute this to the current implementation of the type mask data extraction process: the extractor currently neither performs an in-depth dependency resolution to retrieve the full set of types available in the processed file’s name space, nor does it track type renamings (e.g., `import pandas as pd`). The low predictive capability of comments can be explained by the fact that only a small number of the methods in both datasets have documentation at the method level.

4.5 RQ 3: Effectiveness of Search

To evaluate the search, we collect a ground truth of 50 *fully annotated* files that are randomly sampled from the industrial code base at Facebook. We ensure that they type-check correctly. The files we select originate from different products and vary in size and complexity, the files average 7 (median: 6, 95%: 13) annotations. The total number of annotations is 346. For each file in the ground truth, we strip its existing annotations and then apply TypeWriter to predict and evaluate the missing types. We configure both the greedy and the non-greedy search strategies to stop when the number of states explored is seven times the number of type slots. This threshold empirically strikes a reasonable balance between investing time and detecting correct types. We use the same prediction model trained on the Facebook dataset as in Section 4.3.

Table 3 shows the results on two levels: individual type annotations and files. On the annotation-level, column *type-correct* shows how many type slots the type assignment returned by the search fills (recall that the search ensures each added type to be type-correct). Column *ground truth match* shows how many of all added annotations match the original, developer-produced type annotations. On the file-level, a *complete and type-correct solution* is a file that TypeWriter fully annotates without type errors. This metric does not include files where TypeWriter discovers a type-correct, but partially annotated solution. The *ground truth match* is the subset of the complete and type-correct solutions, where the solution is identical to the ground truth for all types in the file. It is possible to find a type-correct annotation that does not match the ground truth. For example, TypeWriter may correctly annotate the return type of a function as a `List`, but a human expert might choose a more precise type `List[str]`: both are type-correct, but the human annotation provides more guarantees.

Both search strategies successfully annotate a significant fraction of all types. On the annotation-level, they add between 40% and 63% of all types in a type-correct way, out of which 28% to 47% match the ground truth, depending on the search strategy. On the file-level, TypeWriter completely annotates 14% to 44% of all files, and 10% to 22% of all files perfectly match the developer annotations. Comparing the two search strategies, we find that, at the annotation-level, greedy search discovers more type-correct annotations with

Table 3: Effectiveness of various search strategies for type inference.

Strategy	Top- k	Annotations		Files	
		Type-correct	Ground truth match	Complete, type-correct	Ground truth match
Greedy search	1	176 (51%)	155 (45%)	7 (14%)	5 (10%)
	3	213 (62%)	169 (49%)	14 (28%)	10 (20%)
	5	248 (72%)	188 (54%)	22 (44%)	11 (22%)
Non-greedy search	1	175 (51%)	149 (44%)	7 (14%)	5 (10%)
	3	150 (43%)	109 (32%)	11 (22%)	7 (14%)
	5	152 (44%)	109 (32%)	15 (30%)	5 (10%)
Upper bound (prediction)	1	–	192 (55%)	–	5(10%)
	3	–	234 (68%)	–	13 (26%)
	5	–	240 (69%)	–	14 (28%)
Pyre Infer	–	100 (29%)	82 (24%)	3 (2%)	2 (2%)

top-3 and top-5 predictions, while non-greedy search actually finds fewer annotations. This is due to the exponential increase in search space, which makes it less likely that the non-greedy search finds a type-correct solution. In contrast, the results suggest that the greedy search explores a more promising part of the search space. At the file-level, both search approaches provide more annotations and fully annotate more files as the number of available predictions per slot increases. In the greedy case, a search using the top-5 results still improves the outcome significantly; this suggests the search strategy can efficiently leverage the neural model’s moderate improvement when k increases beyond 3.

To better understand how effective the search is, we also show how many ground truth-matching types the top- k predictions include (“upper bound (prediction)”). Note that these numbers are a theoretical upper bound for the search, which cannot be achieved in practice because it would require exhaustively exploring all combinations of types predicted within the top- k . Comparing the upper bound with the results of the search shows that the search gets relatively close to the maximum effectiveness it could achieve. For example, a top-5 exploration with greedy search finds a complete and type-correct solution that matches the ground truth for 11 files, while the theoretical upper bound is 14 files. We leave developing further search strategies, e.g., based on additional heuristics, for future work.

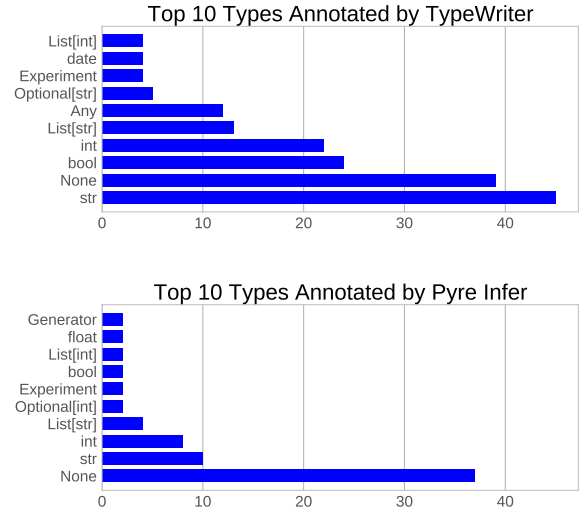
Overall, the results show that a greedy search among top- k types can uncover more types when given more predictions, while also maintaining type correctness. $k = 5$ provides the best annotation performance. While a non-greedy search should not immediately be disregarded, it should be considered in terms of how exhaustive the developer will allow the search to be.

4.6 RQ 4: Comparing with Static Type Inference

We compare TypeWriter with a state-of-the-art, static type inference tool pyre infer. The type inference is part of the pyre type checker and is representative of conservative static analysis-based type inference that adds only types guaranteed to be type-correct. We run pyre infer on the same set of randomly chosen, fully annotated

Table 4: Comparison of TypeWriter and a traditional, static type inference (pyre infer).

	Top-5 (greedy)	Top-5 (non-greedy)
Total type slots	346	346
...added by TypeWriter only	166	95
...added by pyre infer only	18	43
...added by both tools	82	57
...same prediction	63	44
...neither could predict	80	151

**Figure 6: Distribution of types found by TypeWriter and Pyre Infer.**

files as in Section 4.5 and then compare the added annotations with TypeWriter’s top-5 search results.

Tables 3 (bottom) and 4 show the results. In a head to head comparison, TypeWriter is able to provide type-correct predictions for about seven times the number of files that pyre infer can. It also discovers significantly more types, adding a total of 188 types, whereas pyre infer adds only 100. Additionally, of the 82 type slots for which both tools suggest a type, the suggestions are the same in 63 cases. Effectively, the types that TypeWriter suggests are a superset of those inferred by pyre infer, as pyre infer does not uniquely find many types.

To further illustrate the differences, we plot the distribution of the top-10 correctly predicted types in Figure 6. We see that pyre infer can infer more precise types, but the majority of its inferences are on methods with no return types. Moreover, some of the inferred types are of dubious usefulness (e.g., `Optional[Optional[Context]]`) indicating the difficulty of applying static type inference on dynamically-typed languages and reinforcing our thesis on the value of prediction-based type inference.

5 DISCUSSION

Effectiveness of neural type prediction. TypeWriter implements the first neural type prediction model for Python. As all existing

type prediction models [15, 21, 28] target JavaScript code, it is difficult to draw conclusions as to whether the TypeWriter architecture is the best for the task. Two facts seem to suggest so: i) TypeWriter is better by a comfortable margin than a re-implementation of the two best-in-class JavaScript models, and ii) TypeWriter’s performance is stable across two very different datasets.

Type-correctness vs. soundness. Due to the way current Python type checkers work, the types that TypeWriter produces are guaranteed to be *type-correct* within the context of a given module. Type correctness is different from type soundness, as the later can only be verified using human intuition. This means that if a module is used within another context, the type checker might invalidate an initially correct prediction. In turn, this makes TypeWriter a *soundy* [20], rather than a sound approach.

Limited type vocabulary. TypeWriter only predicts types that are part of its type vocabulary. When the vocabulary size is configured at 1000 types, it can account for 90% of the available types in both our datasets. However, as software evolves, developers create new types or change the names of existing ones. This may lead to situations where the model would predict a wrong type because its name changed or because it simply does not know that the type exists. The out-of-vocabulary problem is well known in software engineering research [16]. Recent work for by Karampatsis et al. [18] uses sub-word information to account for neologisms with very good results. We believe that TypeWriter would benefit significantly from such an approach for embedding identifier names, as it would enable it to learn semantically similar name variants (e.g., `AbstractClass` and `Class` or `List` and `List[str]`).

Further improvements. TypeWriter is a prototype stemming from a general effort within Facebook to make their Python code base more robust. TypeWriter can be improved in several dimensions, some of which are presented below:

Better data: The ablation study results suggest that type masks and documentation components of the TypeWriter model are only marginally contributing to its prediction capabilities. This goes against both intuition and published work: in [21], the authors show that code documentation is an important signal. We could, however, exploit the fact that highly used libraries, such as `flask` or the Python standard library itself feature both type annotations (in the `typeshed`⁴ repository) and excellent documentation. Moreover, we can obtain better type masks using lightweight dependency analysis, such as `importlab`,⁵ to identify all types that are in context.

Faster search feedback: TypeWriter’s execution speed is currently constrained by the type checker used to obtain feedback. One natural way to improve this would be to integrate the TypeWriter type predictor into a static type inference loop: when the type inference cannot predict a type for a location, it can ask the neural model for a suggestion. While the theoretical cost of searching for types is similar, in practice the type inference will be able to quickly examine suggested types given that all required data is loaded in memory.

Reinforced learning: As with most neural models, TypeWriter can benefit from more data. One idea worth exploring is to apply

TypeWriter in batches, consisting of application of an initial set of neural predictions, reviewing proposed types through the normal code review process at Facebook and then retrain the model on the new data. At the scale of the Facebook code base, we expect that the feedback obtained (accepted, modified, and rejected suggestions) could be used to improve the learning process.

6 RELATED WORK

Type inference for dynamic languages. Static type inference [4, 8, 14, 17] computes types using, e.g., abstract interpretation or type constraint propagation. These approaches are sound by design, but due to the dynamic nature of some languages, they often infer only simple or very generic types [8, 17]. They also require a significant amount of context, usually a full program and its dependencies. Dynamic type inference [3, 29] tracks data flows between functions, e.g., while executing a program’s test suite. These approaches capture precise types, but they are constrained by coverage. TypeWriter differs from those approaches in two key aspects: i) it only requires limited context information, i.e., a single a source code file, ii) it does not require the program to be executed and hence can predict types in the absence of a test suite or other input data.

Probabilistic type inference. The difficulty of accurately inferring types for dynamic programming languages has led to research on probabilistic type inference. JSNice [28] models source code as a dependency network of known (e.g., constants, API methods) and unknown facts (e.g., types); it then mines information from large code bases to quantify the probability of two items being linked together. Xu et al. [37] predict variable types based on a probabilistic combination of multiple uncertain type hints, e.g., data flows and attribute accesses. They also consider natural language information, but based on lexical similarities of names, and focus on variable types, whereas TypeWriter focuses on function types. DeepTyper [15] uses a sequence-to-sequence neural model to predict types based on a bi-lingual corpus of TypeScript and JavaScript code. NL2Type [21] uses natural language information. Our evaluation directly compares with Python re-implementations of both DeepTyper and NL2Type. Besides advances in the probabilistic type prediction model itself, the more important contribution of our work is to address the imprecision and combinatorial explosion problems of probabilistic type inference. In principle, any of the above techniques can be combined with TypeWriter’s search-based validations to obtain type-correct types in reasonable time.

Type checking and inference for Python. The Python community introduced a type annotation syntax along with a type checker (`mypy`) as part of Python 3.5 version in 2015 [32]. The combination of the two enables *gradual typing* of existing code, where the type checker checks only the annotated parts of the code. Similar approaches have also been explored by the research community [34]. Since 2015, type annotations have seen adoption in several large-scale Python code bases, with products such as Dropbox⁶ and Instagram,⁷ reportedly having annotated large parts of their multi-million line code bases. TypeWriter helps reduce the manual effort required for such a step.

⁴GitHub: [python/typeshed](https://github.com/python/typeshed)

⁵<https://github.com/google/importlab>

⁶Dropbox Blog: How we rolled out one of the largest Python 3 migrations ever

⁷Instagram Engineering Blog: Introducing open source MonkeyType

Machine learning on code. Our neural type prediction model is motivated by a stream of work on machine learning-based program analyses [2]. Beyond type prediction, others have proposed learning-based techniques to find programming errors [25, 27], predict variable and method names [1, 28, 33], suggest how to improve names [19], search code [10, 30], detect clones [36, 40], classify code [24, 39], predict code edits [31, 38, 41], predict assertions [35], and automatically fix bugs [6, 11, 12]. TypeWriter contributes a novel model for predicting types and a search-based combination of predictive models with traditional type checking.

Search-based software engineering. Our search-based validation of types fits the search-based software engineering theme [13], which proposes to balance competing constraints in developer tools through metaheuristic search techniques. In our case, the search balances the need to validate an exponential number of combinations of type suggestions with the need to efficiently annotate types.

7 CONCLUSIONS

We present TypeWriter, a learning-based approach to the problem of inferring types for code written in Python. TypeWriter exploits the availability of partially annotated source code to learn a type prediction model and the availability of type checkers to refine and validate the predicted types. TypeWriter’s learned model can readily predict correct type annotations for half of the type slots on first try, whereas its search component can help prevent annotating code with wrong types. Combined, the neural prediction and the search-based refinement helps annotate large code bases with minimal human intervention, making TypeWriter the first practically applicable learning-based tool for type annotation.

We are currently in the process of making TypeWriter available to developers at Facebook. We have tested the automation of the tool in the code review domain. Developers at Facebook received type suggestions as comments on pull requests they had authored. They would also receive pull requests containing type suggestions for their project. The initial experience from applying the approach on a code base that powers tools used by billions of people has been positive: several thousands of suggested types have already been accepted with minimal changes.

REFERENCES

- [1] Miltiadis Allamanis, Earl T. Barr, Christian Bird, and Charles A. Sutton. 2015. Suggesting accurate method and class names. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2015, Bergamo, Italy, August 30 - September 4, 2015*. 38–49.
- [2] Miltiadis Allamanis, Earl T Barr, Premkumar Devanbu, and Charles Sutton. 2018. A survey of machine learning for big code and naturalness. *ACM Computing Surveys (CSUR)* 51, 4 (2018), 81.
- [3] Jong-hoon (David) An, Avik Chaudhuri, Jeffrey S. Foster, and Michael Hicks. 2011. Dynamic Inference of Static Types for Ruby. In *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '11)*. ACM, New York, NY, USA, 459–472. <https://doi.org/10.1145/1926385.1926437>
- [4] Christopher Anderson, Paola Giannini, and Sophia Drossopoulou. 2005. Towards Type Inference for JavaScript. In *ECOOP 2005 - Object-Oriented Programming*, Andrew P. Black (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 428–452.
- [5] Hlib Babii, Andrea Janes, and Romain Robbes. 2019. Modeling Vocabulary for Big Code Machine Learning. *CoRR* (2019). <https://arxiv.org/abs/1904.01873>
- [6] Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. 2019. Getafix: Learning to Fix Bugs Automatically. In *OOPSLA*.
- [7] Satish Chandra, Colin S. Gordon, Jean-Baptiste Jeannin, Cole Schlesinger, Manu Sridharan, Frank Tip, and Youngil Choi. 2016. Type Inference for Static Compilation of JavaScript. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA 2016)*. ACM, New York, NY, USA, 410–429. <https://doi.org/10.1145/2983990.2984017>
- [8] Michael Furr, Jong-hoon (David) An, Jeffrey S. Foster, and Michael Hicks. 2009. Static Type Inference for Ruby. In *Proceedings of the 2009 ACM Symposium on Applied Computing (SAC '09)*. ACM, New York, NY, USA, 1859–1866. <https://doi.org/10.1145/1529282.1529700>
- [9] Zheng Gao, Christian Bird, and Earl T. Barr. 2017. To Type or Not to Type: Quantifying Detectable Bugs in JavaScript. In *Proceedings of the 39th International Conference on Software Engineering (ICSE '17)*. IEEE Press, Piscataway, NJ, USA, 758–769. <https://doi.org/10.1109/ICSE.2017.75>
- [10] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. 2018. Deep Code Search. In *ICSE*.
- [11] Rahul Gupta, Soham Pal, Aditya Kanade, and Shirish Shevade. 2017. DeepFix: Fixing Common C Language Errors by Deep Learning. In *AAAI*.
- [12] Jacob Harer, Onur Ozdemir, Tomo Lazovich, Christopher P. Reale, Rebecca L. Russell, Louis Y. Kim, and Sang Peter Chin. 2018. Learning to Repair Software Vulnerabilities with Generative Adversarial Networks. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*. 7944–7954. <http://papers.nips.cc/paper/8018-learning-to-repair-software-vulnerabilities-with-generative-adversarial-networks>
- [13] Mark Harman and Bryan F Jones. 2001. Search-based software engineering. *Information and software Technology* 43, 14 (2001), 833–839.
- [14] Mostafa Hassan, Caterina Urban, Marco Eilers, and Peter Müller. 2018. MaxSMT-Based Type Inference for Python 3. In *International Conference on Computer Aided Verification*. Springer, 12–19.
- [15] Vincent J. Hellendoorn, Christian Bird, Earl T. Barr, and Miltiadis Allamanis. 2018. Deep Learning Type Inference. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2018)*. ACM, New York, NY, USA, 152–162. <https://doi.org/10.1145/3236024.3236051>
- [16] Vincent J. Hellendoorn and Premkumar Devanbu. 2017. Are Deep Neural Networks the Best Choice for Modeling Source Code?. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering (ESEC/FSE 2017)*. ACM, New York, NY, USA, 763–773. <https://doi.org/10.1145/3106237.3106290>
- [17] Simon Holm Jensen, Anders Møller, and Peter Thiemann. 2009. Type Analysis for JavaScript. In *Static Analysis*, Jens Palsberg and Zhendong Su (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 238–255.
- [18] Rafael-Michael Karampatsis, Hlib Babii, Romain Robbes, Charles Sutton, and Andrea Janes. 2020. Big Code != Big Vocabulary: Open-Vocabulary Models for Source Code. In *ICSE*.
- [19] Kui Liu, Dongsun Kim, Tegawendé F. Bissyandé, Tae-young Kim, Kisub Kim, Anil Koyuncu, Suntae Kim, and Yves Le Traon. 2019. Learning to spot and refactor inconsistent method names. In *Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019*. 1–12. <https://dl.acm.org/citation.cfm?id=3339507>
- [20] Benjamin Livshits, Manu Sridharan, Yannis Smaragdakis, Ondřej Lhoták, J. Nelson Amaral, Bor-Yuh Evan Chang, Samuel Z. Guyer, Uday P. Khedker, Anders Møller, and Dimitrios Vardoulakis. 2015. In Defense of Soundness: A Manifesto. *Commun. ACM* 58, 2 (Jan. 2015), 44–46. <https://doi.org/10.1145/2644805>
- [21] Rabee Sohail Malik, Jibesh Patra, and Michael Pradel. 2019. NL2Type: Inferring JavaScript Function Types from Natural Language Information. In *Proceedings of the 41st International Conference on Software Engineering (ICSE '19)*. IEEE Press, Piscataway, NJ, USA, 304–315. <https://doi.org/10.1109/ICSE.2019.00045>
- [22] Clemens Mayer, Stefan Hanenberg, Romain Robbes, Eric Tanter, and Andreas Stefk. 2012. An Empirical Study of the Influence of Static Type Systems on the Usability of Undocumented Software. In *Proceedings of the ACM International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA '12)*. ACM, New York, NY, USA, 683–702. <https://doi.org/10.1145/2384616.2384666>
- [23] Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*. 3111–3119.
- [24] Lili Mou, Ge Li, Lu Zhang, Tao Wang, and Zhi Jin. 2016. Convolutional Neural Networks over Tree Structures for Programming Language Processing. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*. 1287–1293.
- [25] Hoan Anh Nguyen, Tien N. Nguyen, Danny Dig, Son Nguyen, Hieu Tran, and Michael Hilton. 2019. Graph-based mining of in-the-wild, fine-grained, semantic code change patterns. In *Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019*. 819–830. <https://doi.org/10.1109/ICSE.2019.00089>
- [26] John-Paul Ore, Sebastian Elbaum, Carrick Detweiler, and Lambros Karkazis. 2018. Assessing the Type Annotation Burden. In *Proceedings of the 33rd ACM/IEEE*

- International Conference on Automated Software Engineering (ASE 2018)*. ACM, New York, NY, USA, 190–201. <https://doi.org/10.1145/3238147.3238173>
- [27] Michael Pradel and Koushik Sen. 2018. DeepBugs: A learning approach to name-based bug detection. *PACMPL* 2, OOPSLA (2018), 147:1–147:25. <https://doi.org/10.1145/3276517>
- [28] Veselin Raychev, Martin Vechev, and Andreas Krause. 2015. Predicting Program Properties from “Big Code”. *SIGPLAN Not.* 50, 1 (Jan. 2015), 111–124. <https://doi.org/10.1145/2775051.2677009>
- [29] Brianna M. Ren, John Toman, T. Stephen Strickland, and Jeffrey S. Foster. 2013. The Ruby Type Checker. In *Proceedings of the 28th Annual ACM Symposium on Applied Computing (SAC '13)*. ACM, New York, NY, USA, 1565–1572. <https://doi.org/10.1145/2480362.2480655>
- [30] Saksham Sachdev, Hongyu Li, Sifei Luan, Seohyun Kim, Koushik Sen, and Satish Chandra. 2018. Retrieval on source code: a neural code search. In *Proceedings of the 2nd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*. ACM, 31–41.
- [31] Michele Tufano, Jevgenija Pantiuchina, Cody Watson, Gabriele Bavota, and Denys Poshyvanyk. 2019. On learning meaningful code changes via neural machine translation. In *Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019*. 25–36. <https://dl.acm.org/citation.cfm?id=3339509>
- [32] G. van Rossum, J. Lehtosalo, and L. Langa. 2014. PEP484: Type Hints. <https://www.python.org/dev/peps/pep-0484/>. [Online; accessed 25-July-2019].
- [33] Bogdan Vasilescu, Casey Casalnuovo, and Premkumar T. Devanbu. 2017. Recovering clear, natural identifiers from obfuscated JS names. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2017, Paderborn, Germany, September 4-8, 2017*. 683–693.
- [34] Michael M. Vitousek, Andrew M. Kent, Jeremy G. Siek, and Jim Baker. 2014. Design and Evaluation of Gradual Typing for Python. *SIGPLAN Not.* 50, 2 (Oct. 2014), 45–56. <https://doi.org/10.1145/2775052.2661101>
- [35] Cody Watson, Michele Tufano, Kevin Moran, Gabriele Bavota, and Denys Poshyvanyk. 2020. On Learning Meaningful Assert Statements for Unit Test Cases. In *ICSE*.
- [36] Martin White, Michele Tufano, Christopher Vendome, and Denys Poshyvanyk. 2016. Deep learning code fragments for code clone detection. In *ASE*. 87–98.
- [37] Zhaogui Xu, Xiangyu Zhang, Lin Chen, Kexin Pei, and Baowen Xu. 2016. Python Probabilistic Type Inference with Natural Language Support. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2016)*. ACM, New York, NY, USA, 607–618. <https://doi.org/10.1145/2950290.2950343>
- [38] Pengcheng Yin, Graham Neubig, Marc Brockschmidt Miltiadis Allamanis, and Alexander L. Gaunt. 2018. Learning to Represent Edits. *CoRR* 1810.13337 (2018).
- [39] Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, Kaixuan Wang, and Xudong Liu. 2019. A Novel Neural Source Code Representation based on Abstract Syntax Tree. In *ICSE*.
- [40] Gang Zhao and Jeff Huang. 2018. DeepSim: deep learning code functional similarity. In *Proceedings of the 2018 ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/SIGSOFT FSE 2018, Lake Buena Vista, FL, USA, November 04-09, 2018*. 141–151.
- [41] Rui Zhao, David Bieber, Kevin Swersky, and Daniel Tarlow. 2018. Neural Networks for Modeling Source Code Edits. (2018).