Do Machine Learning Models Produce TypeScript Types That Type Check?

Ming-Ho Yee ⊠ •

Northeastern University, Boston, USA

Arjun Guha ⊠®

Northeastern University, Boston, USA Roblox Research

Abstract

Type migration is the process of adding types to untyped code to gain assurance at compile time. TypeScript and other gradual type systems facilitate type migration by allowing programmers to start with imprecise types and gradually strengthen them. However, adding types is a manual effort and several migrations on large, industry codebases have been reported to have taken several years. In the research community, there has been significant interest in using machine learning to automate TypeScript type migration. Existing machine learning models report a high degree of accuracy in predicting individual TypeScript type annotations. However, in this paper we argue that accuracy can be misleading, and we should address a different question: can an automatic type migration tool produce code that passes the TypeScript type checker?

We present TypeWeaver, a TypeScript type migration tool that can be used with an arbitrary type prediction model. We evaluate TypeWeaver with three models from the literature: DeepTyper, a recurrent neural network; LambdaNet, a graph neural network; and InCoder, a general-purpose, multi-language transformer that supports fill-in-the-middle tasks. Our tool automates several steps that are necessary for using a type prediction model, including (1) importing types for a project's dependencies; (2) migrating JavaScript modules to TypeScript notation; (3) inserting predicted type annotations into the program to produce TypeScript when needed; and (4) rejecting non-type predictions when needed.

We evaluate TypeWeaver on a dataset of 513 JavaScript packages, including packages that have never been typed before. With the best type prediction model, we find that only 21% of packages type check, but more encouragingly, 69% of files type check successfully.

2012 ACM Subject Classification Software and its engineering \rightarrow Source code generation; General and reference \rightarrow Evaluation; Theory of computation \rightarrow Type structures

Keywords and phrases Type migration, deep learning

Digital Object Identifier 10.4230/LIPIcs.ECOOP.2023.22

Supplementary Material

Software (Artifact Evaluation approved artifact): https://doi.org/10.5281/zenodo.7662708 Software (code repository): https://github.com/nuprl/TypeWeaver

Funding This work is partially supported by the National Science Foundation grant CCF-2102291.

Acknowledgements We thank Northeastern Research Computing and the New England Research Cloud for providing computing resources; and Leif Andersen, Luna Phipps-Costin, Donald Pinckney, and the anonymous reviewers for their feedback.

1 Introduction

Gradual typing allows programmers to freely mix statically and dynamically typed code. This makes it possible to add static types to a large program incrementally, and slowly reap the benefits of static typing without requiring a complete rewrite of an existing codebase at

once [22, 46, 48]. Over the past decade, gradual typing has proliferated, and there are now gradually typed versions of several mainstream languages [6, 7, 11, 14, 31, 33, 37, 48, 52].

TypeScript is a widely used gradually typed language, and a syntactic superset of JavaScript. Programmers can write their code in TypeScript, benefit from static typing, and then compile to JavaScript. However, the process of *migrating* an untyped JavaScript program to TypeScript has remained a labor-intensive manual effort in practice. For example, Airbnb engineers took more than two years to add TypeScript type annotations to 6 million lines of JavaScript [44], and there are several other accounts of multi-year TypeScript migration efforts [3, 8, 36, 39, 43].

To address this problem, there has been significant research interest in using machine learning to predict TypeScript types. Machine learning seems attractive because TypeScript has language features (e.g., eval) that are very difficult to accommodate with traditional, constraint-based approaches. Moreover, there is a significant quantity of open-source Type-Script that is available to serve as training data for a machine learning model. Over the last few years, advances in model architectures and high-quality training data have led to type annotation prediction with high accuracy on individual type annotations [23, 24, 25, 38, 50].

However, in this paper we argue that accuracy can be misleading, and that predicting individual type annotations is just the first step of migrating a codebase from JavaScript to TypeScript. We should address a different question: can an automatic type migration tool produce code that type checks? If so, we prefer type annotations that are non-trivial and useful (i.e., annotations that are not just any). On the other hand, if the code does not type check, it may have too many errors, which can overwhelm a user who may just turn off the tool. Moreover, it may not be feasible to fix the type errors automatically, since type errors refer to code locations whose typed terms are used, and not necessarily to faulty annotations.

To answer the type checking question, we present TypeWeaver, a TypeScript type migration tool that can be used with an arbitrary type prediction model. Our evaluation employs three models from the literature: DeepTyper [23], a recurrent neural network; LambdaNet [50], a graph neural network; and InCoder [18], a general-purpose, multi-language transformer that supports fill-in-the-middle tasks. Our tool automates several steps that are necessary for using a type prediction model, including:

Importing type dependencies Before migrating a JavaScript project, we must ensure that its dependencies are typed. This means transitively migrating dependencies, or ensuring that the dependencies have TypeScript interface declaration (.d.ts) files available.

Module conversion JavaScript code written for Node.js may use either the CommonJS or ECMAScript module system. However, when migrated to TypeScript, only ECMAScript modules preserve type information. Thus, to fully benefit from static type checking, code written with CommonJS modules should be refactored to use ECMAScript modules.

Type weaving Models that assign type labels to variables do not update the JavaScript source to include type annotations. Therefore, to ask if a program type checks, we must "weave" the predicted type annotations with the original JavaScript source to produce TypeScript.

Rejecting non-type predictions Models that predict type annotations as in-filled sequences of tokens can easily produce token sequences that are not syntactic types. These predictions need to be rejected or cleaned for type prediction to work.

After completing these tasks, it is then possible to type check the resulting TypeScript program and evaluate the effectiveness of TypeWeaver.

Our contributions are the following:

```
1 function f(x) { return x+x; } 4 var point = {};
2 f(1) // returns 2 5 point.x = 42;
3 f("a") // returns "aa" 6 point.y = 54;
```

- (a) JavaScript function that adds or concatenates its argument to itself.
- (b) JavaScript that creates a "point" object.

Figure 1 JavaScript code that cannot be easily typed in TypeScript.

- We describe TypeWeaver, a TypeScript type migration tool for evaluating type prediction models, which automates several prerequisite steps (Section 3).
- We provide a dataset of 513 JavaScript packages, which are a subset of the top 1,000 most downloaded packages from the npm Registry. Our dataset includes packages without known type annotations, i.e., code that has never been used before to evaluate type prediction models (Section 4.1).
- We report the success rate of type checking. We answer the questions of how many packages type check, how many files type check, how many type annotations are trivial, and whether the predicted types match human-written types (Section 4.2).
- We discuss the common kinds of errors that result from a type migration (Section 4.3).
- We compare the results of a type migration before and after converting to the ECMAScript module system (Section 4.4).
- Finally, we examine four packages as case studies, to showcase other difficulties that arise during type migration (Section 4.5).

2 Background

In this section, we first provide background on the type migration problem and contrast it to type inference. We then discuss deep-learning-based type annotation prediction, focusing on the tools that we have used with TypeWeaver.

2.1 Type Migration vs. Type Inference

Type inference is related to, but distinct from, type migration. The goal of type inference is to reconstruct the types of variables, expressions, and functions, where some or all the type annotations are missing. In other words, the language is statically typed and the types exist implicitly within the program, so the type inference algorithm computes the missing annotations. Furthermore, inference can frequently compute the principal type of a variable, expression, or function, i.e., the most general type. As a result, there is a single answer for the type of a variable, expression, or function. Additionally, in languages that support type inference, the inferred type annotations are well defined and not added to the program text.

In this paper, we use type migration to describe the problem of migrating a program from an untyped language to a typed language, e.g., from JavaScript to TypeScript, a process that may require refactoring in addition to type inference. The type migration process starts from an untyped program without type annotations: there is no type information available, beyond the basic information available from literal values, operators, and control-flow statements, so type definitions may need to be inserted into the program. Furthermore, multiple type annotations may be valid, rather than having a single principal type. For example, Figure 1a shows a JavaScript function where the parameter x on Line 1 could be annotated as number or string; without additional context, both annotations are valid. The example illustrates another challenge of type migration: f is called on Line 2 with a number and Line 3 with a

string, so the only valid type annotation for x is any.¹ This satisfies the TypeScript compiler's type checker, but may not be a helpful annotation in terms of documentation.

The type migration problem for TypeScript has several additional challenges and we highlight some of them here. TypeScript has a structural type system, which makes it even harder to determine the right annotation. Furthermore, structural types are often verbose, making it difficult for a programmer to understand the code, which defeats one of the benefits of a static type system. Another difficulty is that JavaScript code can be too dynamic to fit within TypeScript's type system, e.g., there is no good way to handle eval, other than using the any annotation as an escape hatch. Finally, certain idioms and patterns in JavaScript code do not fit TypeScript and need to be refactored. For instance, consider Figure 1b, which initializes a "point" object in JavaScript. Line 4 initializes point to an empty object, and Lines 5 and 6 set the x and y properties. However, this cannot be easily typed in TypeScript,² and it is more appropriate to rewrite the code to use TypeScript classes.

In this landscape of challenges, recent work has focused on the narrower problem of assigning type annotations to TypeScript code, in particular, using deep learning approaches. We examine some of these approaches in the next subsection.

2.2 Deep-Learning-Based Type Annotation Prediction

We focus on JavaScript and TypeScript, since there have been a variety of proposed type prediction models for those languages. We evaluate three of them here: DeepTyper [23], LambdaNet [50], and InCoder [18].

DeepTyper was the first deep neural network for TypeScript type prediction, and uses a bidirectional recurrent neural network architecture. LambdaNet was another early approach, and it uses a graph neural network architecture. InCoder is a recent large language model that predicts arbitrary code completions, and while not trained specifically to predict type annotations, its "fill in the middle" capability makes it ideal for that task. All three models use training data from public code repositories.

We require a system that takes a JavaScript project as input and outputs a type-annotated TypeScript project. DeepTyper and LambdaNet output a probability distribution of types for each identifier, which we must then "weave" into the original JavaScript source to produce TypeScript; we describe this technique in Section 3.3. InCoder is a general-purpose, multi-language transformer, so we implemented a front end to use InCoder to predict type annotations and output TypeScript. Currently, our front end only supports type predictions for function parameters; we describe our implementation in Section 3.2.1.

Our approach can be adapted to work with any type prediction model. Older models may require some work to adapt their outputs, but our InCoder front end can easily be extended to support other fill-in-the-middle models, such as OpenAI's model [4] and SantaCoder [5].

2.2.1 DeepTyper

DeepTyper [23] predicts types for variables, function parameters, and function results using a fixed vocabulary of types, i.e., it cannot predict types declared by the program under analysis unless those types were observed during training. DeepTyper treats type inference

 $^{^{1}}$ The union type number \mid string produces a type error in the function.

² The correct, but awkward, type annotation is {x?: number, y?: number}, which declares x and y as optional properties. If x or y were required, the assignment on Line 4 would be a type error. Alternatively, any is valid, but unhelpful.

as a machine translation problem from one language (unannotated TypeScript) to another (annotated TypeScript). Specifically, it uses a model based on a bidirectional recurrent neural network architecture to translate a sequence of tokens into a sequence of types: for each identifier in the source program, DeepTyper returns a probability distribution of predicted types. Because the input token sequence is perfectly aligned with the output type sequence, this task can also be considered a sequence annotation task, where an output type is expected for every input token.³ However, this approach treats each input token as independent from the others, i.e., a source variable may be referenced multiple times and each occurrence may have a different type. To mitigate this, DeepTyper adds a consistency layer to the neural network, which encourages—but cannot enforce—the model to treat multiple occurrences of the same identifier as related.

DeepTyper's dataset is based on the top 1,000 most starred TypeScript projects on GitHub, as of February 2018. After cleaning to remove large files (those with more than 5,000 tokens) and projects that contained only TypeScript declaration files, the dataset was left with 776 TypeScript projects (containing about 62,000 files and about 24 million tokens), which were randomly split into 80% (620 projects) training data, 10% (78 projects) validation data, and 10% (78 projects) test data. Further processing and cleaning of rare tokens resulted in a final vocabulary of 40,195 source tokens and 11,830 types.

The final training dataset contains both identifiers and types, where each identifier has an associated type annotation; this includes annotations inferred by the TypeScript compiler that were not manually annotated by a programmer. The testing dataset contains type annotations and no identifiers; specifically, the type annotations added by programmers are associated with their declaration sites, and all other sites are associated with "no-type." As a result, DeepTyper's predictions are evaluated against the handwritten type annotations, rather than all types in a project.

2.2.2 LambdaNet

Like DeepTyper, LambdaNet [50] predicts type annotations for variables, function parameters, and function returns: it takes an unannotated TypeScript program and outputs a probability distribution of predicted types for each declaration site. LambdaNet improves upon two limitations of DeepTyper. First, it predicts from an open vocabulary, beyond the types that were observed during training; i.e., it can predict user-defined types from within a project. Second, it only produces type predictions at declaration sites, rather than at every variable occurrence; in other words, multiple uses of the same variable will have a consistent type.

LambdaNet uses a graph neural network architecture and represents a source program as a so-called type dependency graph, which is computed from an intermediate representation of TypeScript that names each subexpression. The type dependency graph is a hypergraph that encodes program type variables as nodes, and relationships between those variables as labeled edges. By encoding type variables, LambdaNet makes a single prediction over all occurrences of a variable, rather than a prediction for each instance of a variable. Furthermore, the edges encode logical constraints and contextual hints. Logical constraints include subtyping and assignment relations, functions and calls, objects, and field accesses, while contextual hints include variable names and usages. Finally, LambdaNet uses a pointer network to predict type annotations.

³ The DeepTyper architecture must classify every input token, including ones where an output type does not make sense, such as if, (,), and even whitespace. DeepTyper filters out these predictions, so a user will never observe these meaningless types.

LambdaNet's dataset takes a similar approach to DeepTyper: they selected the 300 most popular TypeScript projects from GitHub that contained 500–10,000 lines of code, and had at least 10% of type annotations that referred to user-defined types. The dataset has about 1.2 million lines of code, and only 2.7% of the code is duplicated. The 300 projects were split into three sets: 200 (67%) for training data, 40 (13%) for validation data, and 60 (20%) for test data. The vocabulary was split into library types, which consist of the top 100 most common types in the training set, and user-defined types, which are all the classes and interfaces defined in source projects. Similar to DeepTyper, LambdaNet's predictions are evaluated against the handwritten annotations that were added by programmers.

2.2.3 InCoder

InCoder [18] is a *large language model* (LLM) for generating arbitrary code that is trained with a *fill-in-the-middle* (FIM) objective on a corpus of several programming languages, including TypeScript.

InCoder's corpus consists of permissively licensed, open-source code from GitHub and GitLab, as well as Q&A and comments from Stack Overflow. This raw data is filtered to exclude: (1) code that is duplicated; (2) code that is not written in one of 28 languages; (3) files that that are extremely large or contain very few alphanumeric characters; (4) code that is likely to be compiler generated; and (5) certain code generation benchmarks. The result is about 159 GB of code, which is dominated by Python and JavaScript. TypeScript is approximately 4.5 GB of the training data.

We describe InCoder in more depth in Section 3.2.1, where we present what is necessary to use it as a type annotation prediction tool for TypeScript.

2.2.4 Evaluating on Accuracy

The main evaluation criteria for the type annotation prediction task is accuracy: what is the likelihood that a predicted type annotation is correct? Correct means the prediction exactly matches the ground truth, which is the handwritten type annotation at that location. Accuracy is typically measured as top-k accuracy, where a prediction is deemed correct if any of the top k most probable predictions is correct. For our evaluation, we would like a result that "just works," i.e., a program that type checks. Therefore, we are only interested in top-1 accuracy, since we take the top guess as the only prediction.

DeepTyper's test dataset makes up 10% (78 projects) of its original corpus and contains only the annotations that were manually added by programmers. Predictions are compared against this ground truth dataset, and DeepTyper reports a top-1 accuracy of 56.9%. Because DeepTyper may predict different types for multiple occurrences of the same variable, the authors also report an inconsistency metric: 15.4% of variables had multiple type predictions.

LambdaNet also compares its predictions against a ground truth of handwritten type annotations, but they use a different corpus and split 20% (60 projects) for the test dataset. LambdaNet can predict user-defined types, so the evaluation reports two sets of results: a top-1 accuracy of 75.6% when predicting only common library types, and a top-1 accuracy of 64.2% when predicting both library and user-defined types.

InCoder was not designed specifically to predict TypeScript type annotations, but the authors report an experiment to predict only the result types for Python functions. For this task, InCoder was evaluated on a test dataset of 469 functions, which was constructed from the CodeXGLUE dataset; InCoder achieved an accuracy of 58.1%.

```
var x = 2;
                   // private
                                           13 var a = require('./a.js');
                                           14
10 module.exports.foo = 42;
                                           15
                                              console.log(a.foo);
11 module.exports.f = (i) => i+x;
                                           16 console.log(a.f(1)); // 3
  (a) CommonJS: a. js exports foo and f.
                                             (b) CommonJS: b.js imports the module a.js.
17 // a.mjs
                                           22 // b.mjs
                   // private
18 var x = 2;
                                           23
                                             import {foo,f} from './a.mjs';
19
                                           24
20 export var foo = 42;
                                           25 console.log(foo);
21 export var f = (i) \Rightarrow i+x;
                                              console.log(f(1)); // 3
  (c) ECMAScript: a.mjs exports foo and f.
                                             (d) ECMAScript: b.mjs imports foo and f.
```

Figure 2 An example comparing the CommonJS and ECMAScript module systems.

However, we argue that accuracy is not the right metric for evaluating a type prediction model. As a first step, we would like to type check the TypeScript project. Additionally, when migrating a JavaScript project to TypeScript, there is frequently no ground truth of handwritten type annotations; instead, the ground truth is what the compiler accepts. This condition is much stronger than accuracy, as even a single, incorrect type annotation causes a package to fail to type check. On the other hand, less precise type annotations (e.g., any) and equivalent annotations (e.g., number | string vs. string | number) may be accepted, despite not matching the ground truth exactly.

In the next section, we describe our approach for type migration and evaluating type prediction models.

3 Approach

We take an end-to-end approach to type migration, starting from an untyped JavaScript project and finishing with a type-annotated TypeScript project that we try to type check. The first step, which is optional, is to convert from CommonJS modules to ECMAScript modules. Next, we invoke one of the type prediction models: DeepTyper and LambdaNet produce type predictions, while InCoder, with a front end we implemented, produces TypeScript. Because DeepTyper and LambdaNet do not produce TypeScript, we perform a step that we call type weaving, which combines type predictions with the original JavaScript source and outputs TypeScript. Finally, we run the TypeScript compiler to type check the now migrated TypeScript project.

3.1 CommonJS to ECMAScript Module Conversion

The first step is to convert projects from CommonJS module notation to ECMAScript module notation. This step is not necessary for type prediction, but is important for the type checking evaluation, as only ECMAScript modules preserve type information across module boundaries. Because we treat this step as optional, our dataset has two versions: the original projects, which may use CommonJS or ECMAScript modules, and a final version that only uses ECMAScript modules.

Most Node.js packages use the CommonJS module system, which was the original module system for Node.js and remains the default. Figures 2a and 2b show an example of the CommonJS module system, where files a.js and b.js implement separate modules. In this

example, a. is sets the foo and f properties of the special module.exports object. Local variables like x are private and not exported. On Line 13, b. js uses the Node.js function require to load module a. js into the local variable a. As a result, a takes on the value of the module.exports object set by a.js, and both foo and f are available as properties of a.

ECMAScript 6 introduced a new module system, referred to as ECMAScript modules. Node.js supports ECMAScript modules when using the .mjs extension or setting a projectwide configuration in the package. json file. Figures 2c and 2d show the same program as before, but rewritten to use ECMAScript modules. In this example, a.mjs directly exports foo and f, rather than writing to a special module.exports object. Then, b.mjs directly imports the names with the import statement instead of loading an object.

TypeScript supports both CommonJS and ECMAScript modules, depending on the project configuration. However, CommonJS modules in TypeScript are untyped; specifically, require is typed as a function that returns any. Therefore, even if a module has type annotations for the variables and functions it exports, those annotations are lost when the module is imported. On the other hand, with ECMAScript modules, the import statement preserves the type annotations of names it imports.

In order to make use of the most type information available, we would like to use ECMA-Script modules in our evaluation. Therefore, we use the cjs-to-es6 tool⁴ to transform our dataset to use ECMAScript modules. The conversion tool is not perfect, and in particular has difficulty when require is used to dynamically load a module. Some of these cases could be fixed manually, but many are genuine uses of dynamic loading in JavaScript.

3.2 Type Annotation Prediction

The next step is to invoke a deep learning model to predict type annotations for a JavaScript project. DeepTyper and LambdaNet require an additional step, which we call type weaving, to produce TypeScript, while InCoder, with our front end, outputs TypeScript directly.

We use the pretrained DeepTyper model available from its GitHub repository,⁵ which is not identical to the model used in the DeepTyper paper. DeepTyper reads in a JavaScript file, and for each identifier, predicts the top five most likely types, outputting the result in comma-separated values (CSV) format.

We use the pretrained LambdaNet model available from its GitHub repository. 6 specifically the model that supports user-defined types. LambdaNet reads in a directory containing a JavaScript project, and predicts the top five most likely types for each variable and function declaration. We modify LambdaNet to output in CSV format.

DeepTyper predicts types for all identifiers in the program, including program locations that do not allow type annotations. Therefore, type weaving must also ensure that type annotations are applied correctly, i.e., only to variable declarations, function parameters, and function results. LambdaNet predicts types for variable and function declarations, and in the correct locations; however, type weaving is still required to produce TypeScript code. Our InCoder front end does not require type weaving, but only supports type predictions for function parameters. We use the pretrained InCoder model available from Hugging Face.⁷

⁴ https://github.com/nolanlawson/cjs-to-es6

https://github.com/DeepTyper/DeepTyper/tree/master/pretrained

⁶ https://github.com/MrVPlusOne/LambdaNet/tree/ICLR20

⁷ https://huggingface.co/facebook/incoder-6B

```
27 function f(x) {
28  return x + 1;
29 }

(a) An example program.

30 function f(x: <|mask:0|>) {
31  return x + 1;
32 }<|mask:1|><|mask:0|>
(b) Preparing code for generation.
```

```
33 function f(x: <|mask:0|>) {
34    return x + 1;
35 }<|mask:1|><|mask:0|>
36    number, y: number<|endofmask|>
```

- (c) InCoder often produces extra tokens after the type. Here it produces a new parameter that is not in the original program.
- **Figure 3** Generating types with InCoder.

```
37 function f(x: number) {
38 return x + 1;
```

(d) We select a prefix of the generated program that is a syntactically valid TypeScript type.

3.2.1 InCoder

InCoder is trained to generate code in the middle of a program, conditioned on the surrounding code. To train on a single example (a file of code), the training procedure replaces a randomly selected contiguous span of tokens with a *mask sentinel* token. It appends the mask sentinel to the end of the example, followed by the tokens that were replaced and a special *end-of-mask* token. The model is then trained as a left-to-right language model. This approach generalizes to support several, non-overlapping masked spans, and its training examples have up to 256 randomly selected masked spans, though the majority have just a single masked span.

In principle one could give InCoder a program with up to 256 types to generate at once. However, we found that InCoder is more successful generating a single type at a time, and with a limited amount of context. To generate a type annotation with InCoder we (1) insert the mask sentinel token at the insertion point; (2) add the mask sentinel to the end of the file; (3) generate at the end of the file until the model produces the end-of-mask token; (4) move the generated text to the insertion point; and (5) remove all sentinels.

Figure 3 shows an example of generating a type annotation. However, a problem that we frequently encountered is that InCoder sometimes generates more than just a single type. For instance, Figure 3c shows an example where InCoder generates a new parameter that is not in the original program. The simplest approach is to reject this result and get InCoder to re-generate completions until it produces a type. However, we found that it is far more efficient to accept a prefix of the generated code if it is a syntactically valid type, which we check with a TypeScript parser in the generation loop.

3.3 Type Weaving

To produce type-annotated TypeScript code, we use a process we call *type weaving* to combine type predictions with the original JavaScript code. Type weaving takes two files as input: a JavaScript source file and an associated CSV file with type predictions. The type weaving program parses the JavaScript source into an abstract syntax tree (AST), and then traverses the AST and CSV files simultaneously, using the TypeScript compiler to insert type annotations into the program AST. Both DeepTyper and LambdaNet require type weaving, but their CSV files are in different formats. Our type weaving program can be extended to support custom CSV formats.

3.3.1 DeepTyper

Each row of a DeepTyper CSV file represents a lexical token from the source program. Rows with non-identifier tokens, such as keywords and symbols, contain two columns: the token text and the token type. Rows with identifiers contain columns for the token text, token type, as well as the top five most likely types and their probabilities.

The DeepTyper implementation has a few limitations that we handle during type weaving. First, the implementation uses regular expressions instead of a parser to tokenize JavaScript code. This results in some tokens that are missing or incorrectly classified as identifiers. Second, DeepTyper provides type predictions for every occurrence of an identifier, so we must use only the predictions for declarations. Finally, DeepTyper often predicts complex as a type; we do not believe it refers to a complex number type, so we replace it with any.

Our type weaving algorithm works as follows: as it traverses the program AST, if it encounters a declaration node, it queries the CSV file for a type prediction. However, the DeepTyper format does not record source location information, and the token classification is brittle, so it is not straightforward to identify which rows are actually declarations and which rows should be skipped. Our algorithm searches the CSV file for a short sequence of rows that corresponds to the declaration node in the AST. This algorithm works well in practice, but does not handle optional parameters or statements that declare multiple variables.

3.3.2 LambdaNet

For each declaration, LambdaNet prints the source location of the identifier (start line, start column, end line, and end column), followed by the top five most likely types and their probabilities. We modified LambdaNet to output in CSV format.

We observed that LambdaNet frequently predicts the following types: Number, Boolean, String, Object, and Void. The first four are valid TypeScript types, but are non-primitive boxed types distinct from number, boolean, string, and object. The TypeScript documentation strongly recommends using the lowercase type names, so we normalize those types during type weaving. Furthermore, Void is not a valid type, so we instead use void. Finally, LambdaNet does not support generic types, but will predict them without type arguments, which is not valid in TypeScript. While we cannot fix every generic type, we normalize Array to any [], which is shorthand for Array<any>.

As our type weaving program traverses the program AST, if it encounters a declaration node, it computes the node's source location information, and uses that to query the CSV file for a type prediction. However, the type annotation cannot be applied directly to the declaration node, as this modifies the AST and invalidates source location information. Therefore, type weaving for LambdaNet occurs in two phases. In the first phase, the traversal does not modify the AST, but saves the declaration node and type prediction in a map. Then, in the second phase, type weaving iterates over the map and updates the AST.

3.4 Type Checking

In the final step, we run the TypeScript compiler to type check the migrated projects. We run the compiler on each project, providing all the TypeScript input files as arguments, and setting the following compiler flags:

https://www.typescriptlang.org/docs/handbook/declaration-files/do-s-and-don-ts.html# number-string-boolean-symbol-and-object

	Table 1	. Summary o	of dataset	categories:	number of	packages,	files,	and line	s of code.
--	---------	-------------	------------	-------------	-----------	-----------	--------	----------	------------

Dataset category	Packages	Files	Lines of code
DefinitelyTyped, no deps	286	2,692	123,157
DefinitelyTyped, with deps	85	671	63,057
Never typed, no deps	102	255	20,729
Never typed, with deps	40	544	19,189
Overall	513	4,162	226,132

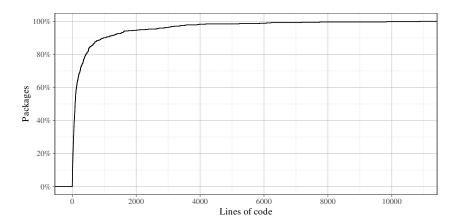


Figure 4 Empirical cumulative distribution function of lines of code per package, over all datasets. The x-axis shows lines of code and the y-axis shows the proportion of packages with fewer than x lines of code.

- --noEmit Type check only, do not emit JavaScript
- --esModuleInterop Improve handling of CommonJS and ECMAScript modules
- --moduleResolution node Explicitly set the module resolution strategy to Node.js
- --target es6 Enable ECMAScript 6 features, which are used by some packages
- --lib es2021,dom Include ECMAScript 2021 library definitions and browser DOM definitions

We do not set the **--strict** flag, allowing the type checker to be more lenient in certain situations, which we expect to already be a significant challenge for automated type migration. Furthermore, we ensure that package dependencies are properly included in our dataset so that the compiler can resolve them.

4 Evaluation

4.1 Dataset

Our dataset for evaluation consists of 513 JavaScript packages. To build this dataset, we start from the top 1,000 most downloaded packages from the npm Registry (as of August 2021) and narrow and clean as follows:

1. We add any transitive dependencies that are not in the original set of packages, to ensure that the dataset is closed.

- 2. We try to fetch the original source code of every package, and eliminate any package where this is not possible (e.g., the package did not provide a repository URL or was deleted from GitHub). Fetching the source helps ensure we are working with original code, and not compiled or "minified" JavaScript.
- 3. We remove packages that were built from a "monorepo," i.e., a single repository containing multiple packages that are published separately. For example, the Babel JavaScript compiler has over 100 separate packages, but all share the same monorepo; fetching each source package meant downloading the entire monorepo multiple times and including unnecessary packages.
- 4. We remove packages that were not implemented in JavaScript, do not contain code, or have more than 10,000 lines of code. The size limit helps us avoid timeouts, and mostly excludes large toolchains and standard libraries, such as the TypeScript compiler and core-js standard library.
- 5. We remove testing code from every package. Tests frequently require extra dependencies, and different frameworks set up the test environment in different ways, which makes large-scale evaluation harder. To remove testing code, we deleted directories named test, tests, __tests__, or spec, and files named test.js, tests.js, test-*.js, *-test.js, *.test.js, or *.spec.js.
- 6. Finally, we ensure that every package has no dependencies, or that all its dependencies are typed, meaning the dependencies have TypeScript type declaration (.d.ts) files available. (We do not require that packages are typed, but only that their dependencies are.) This requirement is necessary because a JavaScript package can only be imported into a TypeScript project if its interface has TypeScript type declarations. The DefinitelyTyped repository⁹ contains interface type declarations for many popular JavaScript packages, and a handful of packages include their own. We download type declarations of project dependencies and include them in our dataset for evaluation purposes—they are not used for type prediction.

After filtering and cleaning our dataset, we classify every package with two criteria: (1) whether the package has type declarations available; and (2) whether the package has dependencies.

If a package has type declarations available, we say it is "DefinitelyTyped" and we can use its type annotations as ground truth in our evaluation. (However, there is evidence that some of these type annotations are incorrect [16, 28, 29, 51].) Otherwise, we use the term "never typed": these packages have never been type annotated and thus no ground truth exists, so machine learning models have never been evaluated on these packages before. If a package has dependencies, we classify it as "with deps" (and from our filtering, we know that every dependency is typed); otherwise, we classify the package as "no deps." Thus, there are four dataset categories; we list them in Table 1 along with the number of packages, files, and lines of code for each category.

Figure 4 is an empirical cumulative distribution function of the lines of code per package: the x-axis shows lines of code in a package and the y-axis shows the proportion of packages with fewer than x lines of codes. From the graph, we observe that approximately 90% of packages have fewer than 1,000 lines of code, and approximately 95% of packages have fewer than 2,000 lines of code.

⁹ https://github.com/DefinitelyTyped/DefinitelyTyped/

T	able 2	Number	and	percentage	of	packages	that	type	check.
---	--------	--------	-----	------------	----	----------	------	------	--------

	DeepTyper			LambdaNet			InCoder		
Dataset category	\checkmark	Total	%	\checkmark	Total	%	\checkmark	Total	%
DefinitelyTyped, no deps	54	257	21.0	24	247	9.7	55	277	19.9
DefinitelyTyped, with deps	5	69	7.2	1	70	1.4	10	77	13.0
Never typed, no deps	31	95	32.6	11	87	12.6	25	100	25.0
Never typed, with deps	5	39	12.8	3	35	8.6	4	39	10.3
Overall	95	460	20.7	39	439	8.9	94	493	19.1

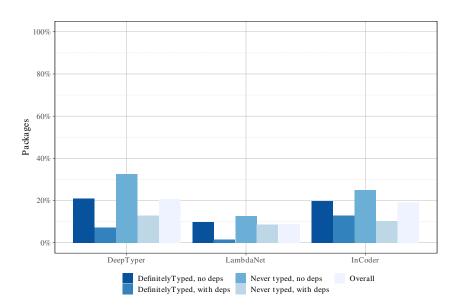


Figure 5 Percentage of packages that type check.

4.2 Success Rate of Type Checking

Do Migrated Packages Type Check?

We first ask if entire packages type check after automated migration from JavaScript to TypeScript. However, not all packages successfully translate to TypeScript with every type migration tool; some packages cause the type migration tool to time out or error. Thus, we report the success rate of type checking as a fraction of the packages that successfully translate to TypeScript.

Table 2 and Figure 5 show the fraction of packages that type check with each tool. We observe that DeepTyper and InCoder perform similarly (20% success rate), and LambdaNet performs worse (9% success rate). Across all tools, packages without dependencies type check at a higher rate than packages with dependencies.

These package-level type checking results are disappointing—but this is a very high standard to meet. Even a single incorrect type annotation causes the entire package to fail. Therefore, we next consider a finer-grained metric that is still useful.

22:14 Do Machine Learning Models Produce TypeScript Types That Type Check?

Table 3 Number and percentage of files with no compilation errors.

	Ι	DeepTyper			LambdaNet			InCoder		
Dataset category	\checkmark	Total	%	\checkmark	Total	%	\checkmark	Total	%	
DefinitelyTyped, no deps	414	1,010	41.0	474	1,638	28.9	1,689	2,401	70.3	
DefinitelyTyped, with deps	195	384	50.8	169	504	33.5	312	547	57.0	
Never typed, no deps	95	205	46.3	63	229	27.5	101	235	43.0	
Never typed, with deps	42	121	34.7	25	534	4.7	467	527	88.6	
Overall	746	1,720	43.4	731	2,905	25.2	2,569	3,710	69.2	

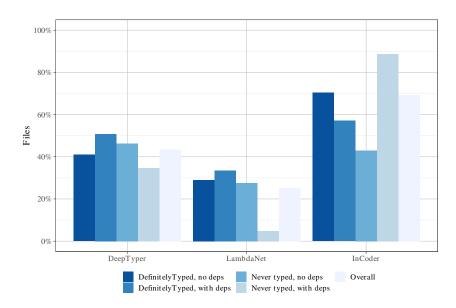


Figure 6 Percentage of files with no compilation errors.

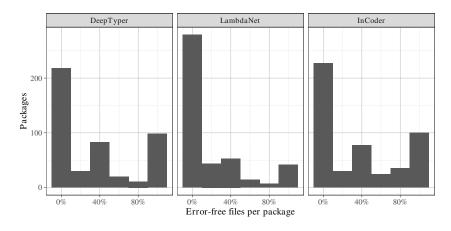


Figure 7 Number of packages vs. percentage of error-free files per package.

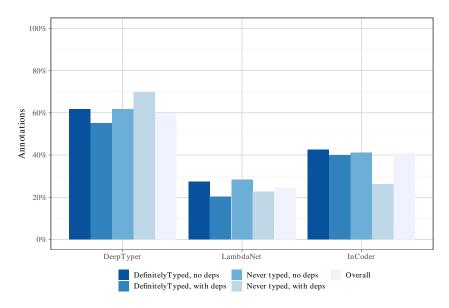


Figure 8 Percentage of annotations that are any, any[], or Function, in files with no errors.

How Many Files are Error Free?

As an alternate measure, we look at the percentage of files with no compilation errors. Instead of a binary pass/fail outcome, this gives us a more fine-grained result for a package. We motivate this metric by observing that TypeScript files are modules with explicit imports and exports. If a file type checks without errors, then it is using all of its internal and imported types consistently. Thus, when triaging type errors, a programmer may (temporarily) set these files aside and focus on the files with compilation errors. However, the programmer may later need to return to a file with no type errors and adjust its type annotations, for example, if a consumer of that file expects a different interface. We give examples of this in our case studies, specifically Section 4.5.2 and Section 4.5.3.

Table 3 and Figure 6 present the fraction of files with no compilation errors. The results are more encouraging: using InCoder, 69% of files are error free. With these results, it is not clear that packages without dependencies outperform packages with dependencies. Finally, in Figure 7 we calculate the percentage of error-free files for each package, and plot histograms of the distribution. Across all tools, most packages have type errors in most or all files.

What Percentage of Type Annotations Are Trivial?

Next, we examine what percentage of type annotations, within the error-free files, are trivial, i.e., what percentage are any, any[] (array of anys), or Function (function that accepts any arguments and returns anything). These annotations can hide type errors and allow more code to type check; however, they provide little value to the programmer.

Figure 8 shows the percentage of trivial type annotations within error-free files. Deep-Typer produces the most (about 60%), LambdaNet produces the least (about 25%), while InCoder is in between (about 40%).

Comparing to the percentage of files with no compilation errors (Figure 6), DeepTyper produces more type-correct code than LambdaNet, but it also generates more trivial type annotations. InCoder produces the most type-correct code, while generating a moderate percentage of trivial type annotations.

22:16 Do Machine Learning Models Produce TypeScript Types That Type Check?

Table 4 Accuracy of type annotations, compared	to non-anv	ground truth.	
---	------------	---------------	--

	Ι	DeepTyper			ambdaN	let	InCoder			
Dataset category	\checkmark	Total	%	\checkmark	Total	%	\checkmark	Total	%	
DefinitelyTyped, no deps	90	241	37.3	116	259	44.8	40	123	32.5	
DefinitelyTyped, with deps	27	123	22.0	41	119	34.5	11	64	17.2	
Overall	117	364	32.1	157	378	41.5	51	187	27.3	

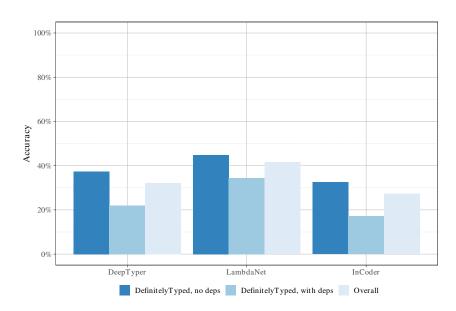


Figure 9 Accuracy of type annotations, compared to non-any ground truth.

Do Migrated Types Match Human-Written Types (When Available)?

Since our dataset is constructed from JavaScript packages instead of TypeScript packages, we do not have fully type-annotated files as our ground truth; instead, we use declaration files provided by the DefinitelyTyped repository or package author. We configure the TypeScript compiler to emit declarations during type checking, which it can do even if the whole package does not type check. Thus, we can compare handwritten, ground truth declarations against declarations generated from migrated packages.

We extract function signatures from declaration files and only compare a signature if it is in both the ground truth and generated declaration. We compare the function parameter types and return types one-to-one, ignoring modifiers (e.g., readonly), and we require an exact string match (i.e. string | number and number | string are considered different types). Following the literature, we skip a comparison if the ground truth is the any annotation.

Our results are presented in Table 4 and Figure 9. We observe that accuracy is better for packages without dependencies. Additionally, our results follow the same pattern in prior work, where LambdaNet has better accuracy than DeepTyper, despite performing worse in our other metrics.

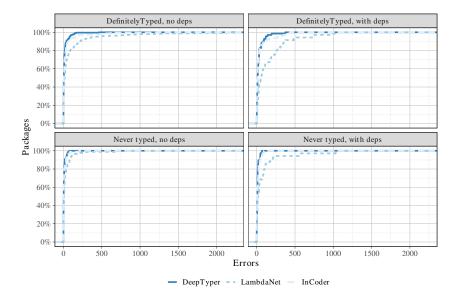


Figure 10 Empirical cumulative distribution function of errors. The x-axis shows the number of errors and the y-axis shows the proportion of packages with fewer than x errors.

How Many Errors Occur in Each Package?

Figure 10 shows an empirical cumulative distribution function of errors: the x-axis shows the number of errors and the y-axis shows the proportion of packages with fewer than x errors. For example, when migrating the "DefinitelyTyped, with deps" dataset with Lambda-Net, approximately 80% of packages have fewer than 250 errors each. Additionally, all of DeepTyper's packages and almost all of InCoder's packages have fewer than 500 errors each.

4.3 Error Analysis

We now consider the kinds of errors that arise during migration. Every TypeScript compiler error has an associated code, ¹⁰ making categorization straightforward. Figure 11 summarizes the top 10 most common errors and Table 5 provides the corresponding messages.

Most of the errors relate to types. These errors are the following: a property not existing on a type (TS2339 and TS2551); an assignment with mismatched types (TS2322); a function call with mismatched parameter and argument types (TS2345); calling a function that was assigned a non-function type annotation (TS2349); and a conditional that compares values from different types (TS2367). TS2314 refers to a generic type that was not provided type arguments; this is caused by DeepTyper and LambdaNet not fully supporting generic types.

The remaining errors are not directly related to types. TS2304 refers to an unknown name, which may not necessarily be a type. TS2554 is emitted because TypeScript requires the number of call arguments to match the number of function parameters, but JavaScript does not. TS2339 includes cases where an empty object is initialized by setting its properties, but TypeScript requires that the object's properties are declared in its type. Finally, TS2307 indicates that the ECMAScript module conversion produced incorrect code.

 $^{^{10}\,\}mathrm{https://github.com/Microsoft/TypeScript/blob/v4.9.3/src/compiler/diagnosticMessages.json$

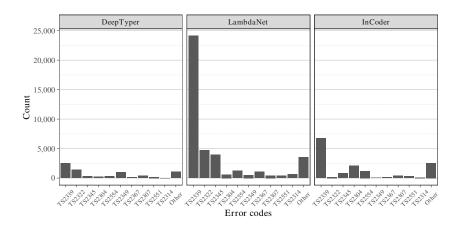


Figure 11 Distribution of the top 10 most common error codes, over all datasets.

4.4 ECMAScript Module Conversion

Recall that we described an optional step before running the evaluation: converting packages to use ECMAScript modules. In this section, we re-run our evaluation—generating types, weaving types, and type checking—to compare the results before and after the conversion step. Specifically, we examine the percentage of packages that type check (Table 6), the percentage of files with no errors (Table 7), and accuracy (Table 8). However, this is not a direct comparison between CommonJS and ECMAScript modules, as some of the original packages were already using ECMAScript modules. Furthermore, the conversion affected a small handful of packages: some packages successfully migrated to TypeScript after the conversion but failed before, and the inverse was true for other packages.

From Table 6, we observe that the ECMAScript module conversion makes fewer packages type check for DeepTyper and InCoder, but slightly improves the results for LambdaNet. In Figure 12, we examine packages that type checked before or after the ECMAScript module conversion; packages that never type checked were excluded. In general, if a package type checked before the conversion, it likely type checked after the conversion. However, if a package failed to type check before the conversion, it was unlikely to type check afterwards; in fact, this never happened for a package with dependencies.

Table 7 compares the percentage of files with no compilation errors. The conversion improves the results for LambdaNet and InCoder, but makes the results worse for DeepTyper. One dramatic result is the change for InCoder and the "never typed, with deps" dataset, where the ECMAScript module conversion results in 89% of files type checking, when it was only 11% before. The difference is caused by a single package, regenerate-unicode-properties, which has over 400 files. With CommonJS modules, each file produces an error; however, with ECMAScript modules, those files type check successfully.

Finally, Table 8 compares the accuracy of type annotations, before and after the ECMA-Script module conversion. Recall that for accuracy, we compare type annotations against the ground truth of handwritten TypeScript declaration files; these are the "DefinitelyTyped" datasets. Accuracy improves for DeepTyper and LambdaNet, but worsens slightly for InCoder.

Table 5 The top 10 most common error codes.

Error code	Message	${\bf DeepTyper}$	LambdaNet	${\rm InCoder}$
TS2339	Property '{0}' does not exist on type '{1}.'	2,510	24,123	6,742
TS2322	Type $\{0\}$ is not assignable to type $\{1\}$.	1,429	4,709	176
TS2345	Argument of type $\{0\}$ is not assignable to	304	3,930	828
	parameter of type '{1}'.			
TS2304	Cannot find name '{0}'.	217	598	2,094
TS2554	Expected $\{0\}$ arguments, but got $\{1\}$.	330	1,234	1,200
TS2349	This expression is not callable.	977	529	26
TS2367	This condition will always return '{0}' since	145	1,046	110
	the types $'\{1\}'$ and $'\{2\}'$ have no overlap.			
TS2307	Cannot find module '{0}' or its correspond-	375	432	371
	ing type declarations.			
TS2551	Property $\{0\}$ does not exist on type $\{1\}$.	187	389	296
	Did you mean '{2}'?			
TS2314	Generic type '{0}' requires {1} type argu-	1	630	95
	ment(s).			
Other		1,089	3,528	2,557
Total		7,564	41,148	$14,\!495$

4.5 Case Studies

In this section, we examine how the models performed on four packages, whether the packages type check, and what steps are left to migrate the packages to TypeScript.

4.5.1 Error Message Does Not Refer to Incorrect Type Annotation

decamelize is a package for converting strings in camel case to lowercase. ¹¹ It is in the "DefinitelyTyped, no deps" dataset, as it ships with a .d.ts declaration file and has no dependencies. Figure 13 shows a simplified version of the function handlePreserve-ConsecutiveUppercase. This function is not exported, thus there are no programmer-written type annotations. Line 42 uses JavaScript's arrow function notation to define a function that takes two arguments, and assigns it to the constant on Line 41. Elsewhere in the package (Line 51 in the listing), the helper function is called with str and sep string arguments.

A programmer inspecting the function can reason that Line 44 is a call to a string method that uses a regular expression on Line 45 to replace text in decamelized with the result on Line 46, where separator is concatenated with the regular expression match. Therefore, both the decamelized and separator parameters on Line 42 should be annotated as string.

The DeepTyper solution is listed on Line 54: it correctly annotates both function parameters as string, but incorrectly annotates handlePreserveConsecutiveUppercase as string. The compiler emits errors for Lines 51 and 55, because Line 51 is attempting to call a non-function, and Line 55 is attempting to assign a function to a non-function variable. However, the fix must be applied to the annotation on line Line 54.

¹¹ https://www.npmjs.com/package/decamelize

22:20 Do Machine Learning Models Produce TypeScript Types That Type Check?

Table 6 Percentage of packages that type check, before and after ECMAScript module conversion.

	DeepTyper		Lambo	laNet	InCo	der
Dataset category	Before	After	Before	After	Before	After
DefinitelyTyped, no deps	25.3	21.0	8.8	9.7	21.3	19.9
DefinitelyTyped, with deps	11.6	7.2	2.9	1.4	14.3	13.0
Never typed, no deps	34.7	32.6	11.5	12.6	26.0	25.0
Never typed, with deps	28.2	12.8	8.8	8.6	23.1	10.3
Overall	25.4	20.7	8.4	8.9	21.3	19.1

■ **Table 7** Percentage of files with no compilation errors, before and after ECMAScript module conversion.

	DeepTyper		Lambo	$_{ m laNet}$	InCo		
Dataset category	Before	After	Before	After	Before	After	
DefinitelyTyped, no deps	45.3	41.0	27.9	28.9	68.3	70.3	
DefinitelyTyped, with deps	50.0	50.8	24.3	33.5	51.6	57.0	
Never typed, no deps	48.8	46.3	24.9	27.5	42.6	43.0	
Never typed, with deps	64.5	34.7	8.6	4.7	11.2	88.6	
Overall	48.1	43.4	23.5	25.2	56.1	69.2	

Table 8 Accuracy of type annotations, before and after ECMAScript module conversion.

	DeepTyper		Lambo	laNet	InCo		
Dataset category	Before	After	Before	After	Before	After	
DefinitelyTyped, no deps	35.2	37.3	44.4	44.8	34.0	32.5	
DefinitelyTyped, with deps	19.8	22.0	27.0	34.5	19.2	17.2	
Overall	28.7	32.1	38.9	41.5	29.1	27.3	

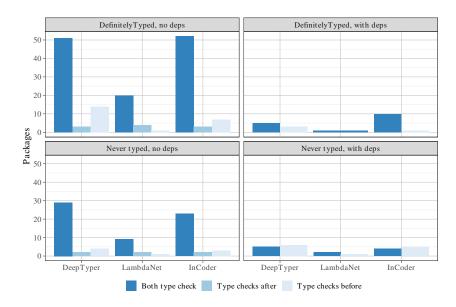


Figure 12 Packages that type check before or after ECMAScript module conversion.

```
40 // Original
41 const handlePreserveConsecutiveUppercase =
42
    (decamelized, separator) => {
         code omitted and simplified ...
43
      return decamelized.replace(
44
45
         /([A-Z]+)([A-Z][a-z]+)/gu
         (_, $1, $2) => $1 + separator + $2.toLowerCase(),
46
47
48 }
49
50 // Elsewhere in the package; str and sep are both strings
51 return handlePreserveConsecutiveUppercase(str, sep);
52
53 // DeepTyper solution
54 const handlePreserveConsecutiveUppercase: string =
    (decamelized: string, separator: string) => { ... }
55
56
57 // LambdaNet solution
58 const handlePreserveConsecutiveUppercase: Function =
    (decamelized: string, separator: number) => { ... }
```

Figure 13 The handlePreserveConsecutiveUppercase function adapted from the decamelize package. The DeepTyper and LambdaNet solutions are also shown.

4.5.2 Incorrect Type Annotation Can Type Check Successfully

The LambdaNet solution on Line 58 correctly annotates handlePreserveConsecutive-Uppercase as Function, but it incorrectly annotates the separator parameter on Line 59 as number. We would expect the compiler to emit a type error, since Line 51 calls the function with string arguments. However, the code type checks successfully, because the generic Function type on Line 58 accepts any number of arguments of any type. The Function type annotation is similar to any, in that it enables more code to type check, but at the cost of fewer type guarantees.

Another example of this problem is the ieee754 package, which reads and writes floating point numbers to and from buffers.¹² It is categorized as "DefinitelyTyped, no deps," since it provides a .d.ts declaration file and has no dependencies. Figure 14 shows the original declaration for the write function on Line 61, and the handwritten, ground truth signature on Line 65.

Consider the DeepTyper solution: the compiler emits an error on Line 70, because a function is being assigned to a variable of type void. However, even if that error is fixed, there is another, more subtle error not detected by the compiler: the isLE parameter on Line 71 is incorrectly annotated as number, not boolean. Because this is compatible with the body of the function, there is no error. (The Buffer annotation is valid, despite not matching the ground truth Uint8Array, because Buffer is defined by the Node.js standard library as a subtype of Uint8Array.)

The InCoder solution on Line 75 type checks successfully. It also uses the Buffer type for buffer, and it uses any instead of number for the value parameter on Line 76. The any annotation may cause run-time errors if the function is called with arguments of the wrong type.

¹² https://www.npmjs.com/package/ieee754

Figure 14 The write function adapted from the ieee754 package. The ground truth signature is also shown, along with the DeepTyper and InCoder solutions.

```
78 // LambdaNet solution
79 export default function (thingToPromisify: string) {
80    if (typeof thingToPromisify === 'function') {
81        return promisify(thingToPromisify)
82    }
83    throw new TypeError('Can only promisify functions or objects')
84 };
```

Figure 15 The LambdaNet solution for a function adapted from the @gar/promisify package.

4.5.3 Run-Time Type Assertions

The <code>@gar/promisify</code> package, ¹³ simplified and shown in Figure 15, is another example where a program type checks, but is incorrect. The example exports a function that takes an argument <code>thingToPromisify</code>, type annotated as <code>string</code> by LambdaNet. Line 80 performs a run-time type check with the <code>typeof</code> operator. This ensures that <code>thingToPromisify</code> is a function on Line 81, which is what the <code>promisify</code> function, defined by Node.js, expects. If <code>thingToPromisify</code> is not a function, the exception on Line 83 is thrown.

The example type checks successfully, because the TypeScript compiler treats the typeof check as a type guard, and reasons that on Line 81, the thingToPromisify variable has been narrowed¹⁴ to a more specific type. However, because thingToPromisify is annotated as string, the type guard always returns false. Therefore, Line 81 is actually unreachable, so the exception on Line 83 is always thrown.

4.5.4 Variable Used as Two Different Types

The example in Figure 16 is adapted from the array-unique package. ¹⁵ The example contains two for loops: a traditional, counter-based for loop on Line 90, and a for...in

¹³ https://www.npmjs.com/package/@gar/promisify

¹⁴ https://www.typescriptlang.org/docs/handbook/2/narrowing.html

¹⁵ https://www.npmjs.com/package/array-unique

```
85 export default function(arr: any[]) {
86    var len: number = arr.length;
87    var o: object = {};
88    var i: number;
89
90    for (i = 0; i < len; i += 1) { ... }
91
92    for (i in o) { ... }
93 };</pre>
```

Figure 16 The LambdaNet solution for a function adapted from the array-unique package.

loop on Line 92 that iterates over all enumerable string properties of an object. Both loops share the same loop variable, i, defined on Line 88 and annotated as number by LambdaNet.

The use of i on Line 92 causes a type error, as for...in loops require the loop variable to be string. However, changing the annotation on Line 88 to string causes a type error on Line 90, as counter-based for loops require the loop variable to be number. One solution is to use the any annotation, and another is to use the union type number | string. Ultimately, the correct solution is to define separate loop variables; this example highlights that code written in JavaScript may need to be refactored for TypeScript.

5 Discussion

How should type prediction models be evaluated? Prior work has used accuracy to evaluate type prediction models, but in this paper, we argue that we should instead type check the generated code. However, we acknowledge that our proposed metric also has limitations: code may type check with trivial annotations (e.g., any or Function) that provide little benefit to the programmer. Furthermore, type correctness does not necessarily mean the type annotations are correct: any can hide type errors that are only encountered at run time.

We do not claim that our metric is the final word on evaluating type prediction models, but we believe it is an improvement over accuracy. We hope this paper can spark a discussion on how machine learning for type migration should be evaluated.

Can slightly wrong type annotations be useful? A type prediction model may suggest types that are slightly wrong and easily fixable by a programmer, but fail to type check. However, a tool that produces hundreds or thousands of slightly wrong type annotations would overwhelm the programmer, and we believe it is important to build tools that try to produce fewer errors. On the other hand, slightly wrong type annotations may still provide value, but we would need to define what "slightly wrong" means and how to measure it. Without a tool like TypeWeaver, which weaves type annotations into code and type checks the result, it would not be possible to ask these questions.

Should we evaluate on JavaScript or TypeScript programs? We choose to evaluate on JavaScript programs, so our dataset deviates from prior work, which only considered TypeScript. Our motivating problem is not to recover type annotations for TypeScript programs that already type check, but to migrate untyped JavaScript programs to type-annotated TypeScript. For this problem, type prediction on its own is not enough, and other steps and further refactoring may be required.

Our methodology makes it possible to evaluate performance on code without known type annotations, i.e., code that has never been typed before. In contrast, prior work required the benchmarks to have ground truth type annotations. Our approach also reduces the likelihood of training data leaking into the test set.

However, there may be scenarios where a type prediction model is used to generate type annotations for a partially annotated TypeScript project. In these situations, type migration would likely not require additional refactoring steps.

Can we fully automate type migration? Our results show that automatically predicting type annotations is a challenging task and much work remains to be done. Furthermore, migrating JavaScript to TypeScript involves more than just adding type annotations: the two languages are different and some refactoring may be required. The models we evaluate in this paper do not refactor code, and we believe it is unlikely for automated type migration to be perfect. Thus, some manual refactoring will always be necessary for certain kinds of code, but we hope that tools can reduce the overall burden on programmers.

6 Related Work

There are many constraint-based approaches to type migration for the gradually typed lambda calculus (GTLC) and some modest extensions. The earliest approach was a variation of unification-based type inference [47], and more recent work uses a wide range of techniques [9, 12, 21, 34, 35, 40]. Since these approaches are based on programming language semantics, they produce sound results, which is their key advantage over learning-based approaches. However, these would require significant work to scale to complex programming languages such as JavaScript.

There are also several constraint-based approaches to type inference for larger languages. Anderson et al. [2] presents type inference for a small fragment of JavaScript, but is not designed for gradual typing. Rastogi et al. [42] infer gradual types for ActionScript to improve performance. More recently, Chandra et al. [13] infer types for JavaScript programs with the goal of compiling them to run efficiently on low-powered devices; their approach is not gradual by design and deliberately rejects certain programs. DRuby [20] infers types for Ruby and treats type annotations in a novel way: inference assumes that annotations are correct, and defers checking them to runtime.

Although this paper focuses on type migration for TypeScript, there are several other gradual type systems for JavaScript [15, 22, 30, 49]. These languages do not have support for type inference and do not provide tools for type migration. Instead, like Typed Racket [48], they require programmers to manually migrate their code to add types. However, there are tools that use dynamic profiles to infer types for these type systems [1, 19, 45].

Even when constraint-based type inference succeeds in a gradually typed language, it can fail to produce the kinds of types that programmers write, e.g., named types, instead of the most general structural type for every annotation. Soft Scheme [10] infers types for Scheme programs, but Flanagan [17, p. 41] reports that it produces unintuitive types. For Ruby, InferDL [27] uses hand-coded heuristics to infer more natural types, and SimTyper [26] uses machine learning to predict equalities between structural types and more natural types.

LambdaNet [50] and DeepTyper [23] are two different approaches for predicting types for TypeScript and JavaScript programs. This paper evaluates using both of them in its type migration pipeline. We discuss them at length in Sections 2.2.1 and 2.2.2. NL2Type [32] is another system for predicting JavaScript types that improves on DeepTyper.

There are also type prediction systems for Python. TypeWriter [41] is notable because it also asks if the resulting Python program type checks. If it does not, it searches its solution space for an alternative typing. A distinction between Python type systems and TypeScript is that Python code is predominantly nominally typed: the type of a variable is either a builtin type or a class, whereas TypeScript uses structural types.

DiverseTyper [24] is a recently published work that predicts both built-in and user-defined types for TypeScript and achieves state-of-the-art accuracy on type prediction. DiverseTyper builds on TypeBert [25], which trains a BERT-based model to predict types. Although this paper does not evaluate these models, they are most closely related to InCoder [18], which is a general-purpose code generation model that we do evaluate.

7 Conclusion

In this paper, we set out to answer the question: do deep-learning-based type annotation prediction models produce TypeScript types that type check? To answer this question, we build TypeWeaver, a type migration tool that automatically converts JavaScript projects into TypeScript. TypeWeaver uses a type annotation prediction model, but does the work of "weaving" predicted types into JavaScript code. It also automates other steps, such as converting JavaScript projects to ECMAScript module notation. Finally, TypeWeaver runs the TypeScript compiler to type check the generated code. TypeWeaver is designed so that any type prediction model can be plugged in, and we use three very different models: DeepTyper, LambdaNet, and InCoder.

In addition to building TYPEWEAVER, we also present a dataset of 513 widely used JavaScript packages that are suitable for type migration. Every package in our dataset has typed dependencies and many of them have never been typed before. With this dataset, we evaluate TYPEWEAVER with all three type prediction models.

The results are mixed. If we ask, "How many packages type check when migrated to TypeScript?" we find that most packages have some type errors. However, we also ask, "How many files are error free?" and the result is more promising. We find that most files have no type errors, which means that programmers performing type migration can focus their attention on a smaller number of files.

Our case studies highlight two insights: (1) certain patterns in JavaScript do not make sense in TypeScript, so a migration may require manual rewriting of the code; and (2) there are cases where programs successfully type check but still have run-time errors.

We believe that currently, while type prediction cannot always reliably migrate JavaScript to TypeScript, it can still be a powerful tool.

Future Work. There are several directions we would like to explore in future work. First, we would like to improve dataset quality. We observed projects that were trivially typable: there were few declarations to annotate, or the annotations were mostly primitive types, so those projects often type checked successfully. Second, we are interested in exploring different evaluation criteria for type prediction models. We believe that type checking the output of these models is only the first step, and that it may be necessary to evaluate the run-time behavior of migrated programs. Additionally, there may be utility in permitting "slightly wrong" type annotations. Finally, we would like to examine other deep learning models and type migration tasks beyond type annotation prediction.

References

- Jong-hoon (David) An, Avik Chaudhuri, Jeffrey S. Foster, and Michael Hicks. Dynamic Inference of Static Types for Ruby. In *Principles of Programming Languages (POPL)*, 2011. doi:10.1145/1926385.1926437.
- 2 Christopher Anderson, Paola Giannini, and Sophia Drossopoulou. Towards Type Inference for JavaScript. In European Conference on Object-Oriented Programming (ECOOP), 2005. doi:10.1007/11531142_19.
- 3 Luke Autry. How we failed, then succeeded, at migrating to TypeScript. https://heap.io/blog/migrating-to-typescript, 2019. Accessed: 2022-12-01.
- 4 Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry Tworek, and Mark Chen. Efficient Training of Language Models to Fill in the Middle, 2022. doi:10.48550/arXiv.2207.14255.
- 5 Loubna Ben Allal, Raymond Li, Denis Kocetkov, Chenghao Mou, Christopher Akiki, Carlos Munoz Ferrandis, Niklas Muennighoff, Mayank Mishra, Alex Gu, Manan Dey, Logesh Kumar Umapathi, Carolyn Jane Anderson, Yangtian Zi, Joel Lamy Poirier, Hailey Schoelkopf, Sergey Troshin, Dmitry Abulkhanov, Manuel Romero, Michael Lappert, Francesco De Toni, Bernardo García del Río, Qian Liu, Shamik Bose, Urvashi Bhattacharyya, Terry Yue Zhuo, Ian Yu, Paulo Villegas, Marco Zocca, Sourab Mangrulkar, David Lansky, Huu Nguyen, Danish Contractor, Luis Villa, Jia Li, Dzmitry Bahdanau, Yacine Jernite, Sean Hughes, Daniel Fried, Arjun Guha, Harm de Vries, and Leandro von Werra. SantaCoder: don't reach for the stars!, 2023. doi:10.48550/arXiv.2301.03988.
- 6 Gavin Bierman, Martín Abadi, and Mads Torgersen. Understanding TypeScript. In European Conference on Object-Oriented Programming (ECOOP), 2014. doi:10.1007/978-3-662-44202-9_11.
- 7 Ambrose Bonnaire-Sergeant, Rowan Davies, and Sam Tobin-Hochstadt. Practical Optional Types for Clojure. In European Symposium on Programming (ESOP), 2016. doi:10.1007/ 978-3-662-49498-1 4.
- 8 Ryan Burgess, Joe King, Stacy London, Sumana Mohan, and Jem Young. Type-Script migration Strict type of cocktails. https://frontendhappyhour.com/episodes/typescript-migration-strict-type-of-cocktails, 2022. Accessed: 2022-12-01.
- 9 John Peter Campora, Sheng Chen, Martin Erwig, and Eric Walkingshaw. Migrating Gradual Types. *Proc. ACM Program. Lang.*, 2(POPL), 2018. doi:10.1145/3158103.
- 10 Robert Cartwright and Mike Fagan. Soft Typing. In *Programming Language Design and Implementation (PLDI)*, 1991. doi:10.1145/113445.113469.
- Mauricio Cassola, Agustín Talagorria, Alberto Pardo, and Marcos Viera. A Gradual Type System for Elixir. In Brazilian Symposium on Context-Oriented Programming and Advanced Modularity (SBLP), 2020. doi:10.1145/3427081.3427084.
- 12 Giuseppe Castagna, Victor Lanvin, Tommaso Petrucciani, and Jeremy G. Siek. Gradual Typing: A New Perspective. Proc. ACM Program. Lang., 3(POPL), 2019. doi:10.1145/3290329.
- Satish Chandra, Colin S. Gordon, Jean-Baptiste Jeannin, Cole Schlesinger, Manu Sridharan, Frank Tip, and Youngil Choi. Type Inference for Static Compilation of JavaScript. In Object-Oriented Programming Systems Languages and Applications (OOPSLA), 2016. doi: 10.1145/2983990.2984017.
- Avik Chaudhuri, Panagiotis Vekris, Sam Goldman, Marshall Roch, and Gabriel Levi. Fast and Precise Type Checking for JavaScript. Proc. ACM Program. Lang., 1(OOPSLA), 2017. doi:10.1145/3133872.
- Ravi Chugh, David Herman, and Ranjit Jhala. Dependent Types for JavaScript. In *Object-Oriented Programming Systems Languages and Applications (OOPSLA)*, 2012. doi:10.1145/2384616.2384659.
- Asger Feldthaus and Anders Møller. Checking Correctness of TypeScript Interfaces for JavaScript Libraries. In Object-Oriented Programming Systems Languages and Applications (OOPSLA), 2014. doi:10.1145/2660193.2660215.

17 Cormac Flanagan. Effective Static Debugging via Componential Set-based Analysis. PhD thesis, Rice University, 1997. URL: https://users.soe.ucsc.edu/~cormac/papers/thesis.pdf.

- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Scott Yih, Luke Zettlemoyer, and Mike Lewis. InCoder: A Generative Model for Code Infilling and Synthesis. In *International Conference on Learning Representations (ICLR)*, 2023. doi:10.48550/arXiv.2204.05999.
- Michael Furr, Jong-hoon (David) An, and Jeffrey S. Foster. Profile-Guided Static Typing for Dynamic Scripting Languages. In Object-Oriented Programming Systems Languages and Applications (OOPSLA), 2009. doi:10.1145/1640089.1640110.
- 20 Michael Furr, Jong-hoon (David) An, Jeffrey S. Foster, and Michael Hicks. Static Type Inference for Ruby. In Symposium on Applied Computing (SAC), 2009. doi:10.1145/1529282.1529700.
- 21 Ronald Garcia and Matteo Cimini. Principal Type Schemes for Gradual Programs. In *Principles of Programming Languages (POPL)*, 2015. doi:10.1145/2676726.2676992.
- 22 Arjun Guha, Claudiu Saftoiu, and Shriram Krishnamurthi. Typing Local Control and State Using Flow Analysis. In *European Symposium on Programming (ESOP)*, 2011. doi: 10.1007/978-3-642-19718-5_14.
- Vincent J. Hellendoorn, Christian Bird, Earl T. Barr, and Miltiadis Allamanis. Deep Learning Type Inference. In European Software Engineering Conference/Foundations of Software Engineering (ESEC/FSE), 2018. doi:10.1145/3236024.3236051.
- 24 Kevin Jesse, Premkumar Devanbu, and Anand Ashok Sawant. Learning To Predict User-Defined Types. IEEE Transactions on Software Engineering (TSE), 2022. doi:10.1109/TSE. 2022.3178945.
- Kevin Jesse, Premkumar T. Devanbu, and Toufique Ahmed. Learning Type Annotation: Is Big Data Enough? In European Software Engineering Conference/Foundations of Software Engineering (ESEC/FSE), 2021. doi:10.1145/3468264.3473135.
- Milod Kazerounian, Jeffrey S. Foster, and Bonan Min. SimTyper: Sound Type Inference for Ruby Using Type Equality Prediction. Proc. ACM Program. Lang., 5(OOPSLA), 2021. doi:10.1145/3485483.
- 27 Milod Kazerounian, Brianna M. Ren, and Jeffrey S. Foster. Sound, Heuristic Type Annotation Inference for Ruby. In *Dynamic Languages Symposium (DLS)*, 2020. doi:10.1145/3426422. 3426985.
- 28 Erik Krogh Kristensen and Anders Møller. Inference and Evolution of TypeScript Declaration Files. In Fundamental Approaches to Software Engineering (FASE), 2017. doi:10.1007/978-3-662-54494-5_6.
- 29 Erik Krogh Kristensen and Anders Møller. Type Test Scripts for TypeScript Testing. Proc. ACM Program. Lang., 1(OOPSLA), 2017. doi:10.1145/3133914.
- 30 Benjamin S. Lerner, Joe Gibbs Politz, Arjun Guha, and Shriram Krishnamurthi. TeJaS: Retrofitting Type Systems for JavaScript. In *Dynamic Languages Symposium (DLS)*, 2013. doi:10.1145/2578856.2508170.
- 31 Kuang-Chen Lu, Ben Greenman, Carl Meyer, Dino Viehland, Aniket Panse, and Shriram Krishnamurthi. Gradual Soundness: Lessons from Static Python. The Art, Science, and Engineering of Programming, 7(1), 2022. doi:10.22152/programming-journal.org/2023/7/2.
- 32 Rabee Sohail Malik, Jibesh Patra, and Michael Pradel. NL2Type: Inferring JavaScript Function Types from Natural Language Information. In *International Conference on Software Engineering (ICSE)*, 2019. doi:10.1109/ICSE.2019.00045.
- 33 Meta Platforms, Inc. Pyre: A performant type-checker for Python 3. https://pyre-check.org/. Accessed: 2022-12-01.
- 34 Zeina Migeed and Jens Palsberg. What Is Decidable about Gradual Types? Proc. ACM Program. Lang., 4(POPL), 2020. doi:10.1145/3371097.
- 35 Yusuke Miyazaki, Taro Sekiyama, and Atsushi Igarashi. Dynamic Type Inference for Gradual Hindley–Milner Typing. *Proc. ACM Program. Lang.*, 3(POPL), 2019. doi:10.1145/3290331.

- Thomas Moore. How We Completed a (Partial) TypeScript Migration In Six Months. https://blog.abacus.com/how-we-completed-a-partial-typescript-migration-in-six-months/, 2019. Accessed: 2022-12-01.
- 37 Guilherme Ottoni. HHVM JIT: A Profile-Guided, Region-Based Compiler for PHP and Hack. In Programming Language Design and Implementation (PLDI), 2018. doi:10.1145/3192366. 3192374.
- 38 Irene Vlassi Pandi, Earl T. Barr, Andrew D. Gordon, and Charles Sutton. OptTyper: Probabilistic Type Inference by Optimising Logical and Natural Constraints, 2021. doi: 10.48550/arXiv.2004.00348.
- 39 Mihai Parparita. The Road to TypeScript at Quip, Part Two. https://quip.com/blog/the-road-to-typescript-at-quip-part-two, 2020. Accessed: 2022-12-01.
- 40 Luna Phipps-Costin, Carolyn Jane Anderson, Michael Greenberg, and Arjun Guha. Solver-Based Gradual Type Migration. Proc. ACM Program. Lang., 5(OOPSLA), 2021. doi: 10.1145/3485488.
- 41 Michael Pradel, Georgios Gousios, Jason Liu, and Satish Chandra. TypeWriter: Neural Type Prediction with Search-Based Validation. In European Software Engineering Conference/Foundations of Software Engineering (ESEC/FSE), 2020. doi:10.1145/3368089.3409715.
- 42 Aseem Rastogi, Avik Chaudhuri, and Basil Hosmer. The Ins and Outs of Gradual Type Inference. In *Principles of Programming Languages (POPL)*, 2012. doi:10.1145/2103656. 2103714.
- 43 Felix Rieseberg. TypeScript at Slack. https://slack.engineering/typescript-at-slack/, 2017. Accessed: 2022-12-01.
- 44 Sergii Rudenko. ts-migrate: A Tool for Migrating to TypeScript at Scale. https://medium.com/airbnb-engineering/ts-migrate-a-tool-for-migrating-to-type script-at-scale-cd23bfeb5cc, 2020. Accessed: 2022-12-01.
- Claudiu Saftoiu. JSTrace: Run-time Type Discovery for JavaScript. Master's thesis, Brown University, 2010. URL: https://cs.brown.edu/research/pubs/theses/ugrad/2010/saftoiu.pdf.
- Jeremy G. Siek and Walid Taha. Gradual Typing for Functional Languages. In *Scheme and Functional Programming Workshop*, 2006. URL: http://schemeworkshop.org/2006/13-siek.
- Jeremy G. Siek and Manish Vachharajani. Gradual Typing with Unification-Based Inference. In *Dynamic Languages Symposium (DLS)*, 2008. doi:10.1145/1408681.1408688.
- Sam Tobin-Hochstadt and Matthias Felleisen. The Design and Implementation of Typed Scheme. In *Principles of Programming Languages (POPL)*, 2008. doi:10.1145/1328438. 1328486.
- 49 Panagiotis Vekris, Benjamin Cosman, and Ranjit Jhala. Trust, but Verify: Two-Phase Typing for Dynamic Languages. In European Conference on Object-Oriented Programming (ECOOP), 2015. doi:10.4230/LIPIcs.ECOOP.2015.52.
- 50 Jiayi Wei, Maruth Goyal, Greg Durrett, and Isil Dillig. LambdaNet: Probabilistic Type Inference using Graph Neural Networks. In *International Conference on Learning Representations (ICLR)*, 2020. doi:10.48550/arXiv.2005.02161.
- Jack Williams, J. Garrett Morris, Philip Wadler, and Jakub Zalewski. Mixed Messages: Measuring Conformance and Non-Interference in TypeScript. In *European Conference on Object-Oriented Programming (ECOOP)*, 2017. doi:10.4230/LIPIcs.ECOOP.2017.28.
- Jake Zimmerman. Sorbet: Stripe's type checker for Ruby. https://stripe.com/blog/sorbet-stripes-type-checker-for-ruby, 2022. Accessed: 2022-12-01.