

# Squash the work

A Workflow for Typing Untyped Programs that use Ad-Hoc Data Structures

AMBROSE BONNAIRE-SERGEANT, SAM TOBIN-HOCHSTADT, Indiana University, USA

We present a semi-automated workflow for porting untyped programs to annotation-driven optional type systems. Unlike previous work, we infer useful types for recursive heterogeneous entities that have “ad-hoc” representations as plain data structures like maps, vectors, and sequences.

Our workflow starts by using dynamic analysis to collect samples from program execution via test suites or examples. Then, initial type annotations are inferred by combining observations across different parts of the program. Finally, the programmer uses the type system as a feedback loop to tweak the provided annotations until they type check.

Since inferring perfect annotations is usually undecidable and dynamic analysis is necessarily incomplete, the key to our approach is generating close-enough annotations that are easy to manipulate to their final form by following static type error messages. We explain our philosophy behind achieving this along with a formal model of the automated stages of our workflow, featuring maps as the primary “ad-hoc” data representation.

We report on using our workflow to convert real untyped Clojure programs to type check with Typed Clojure, which both feature extensive support for ad-hoc data representations. First, we visually inspect the initial annotations for conformance to our philosophy. Second, we quantify the kinds of manual changes needed to amend them. Third, we verify the initial annotations are meaningfully underprecise by enforcing them at runtime.

We find that the kinds of changes needed are usually straightforward operations on the initial annotations, leading to a substantial reduction in the effort required to port such programs.

## 1 INTRODUCTION

Consider the exercise of counting binary tree nodes using JavaScript. With a class-based tree representation, we naturally add a method to each kind of node like so.

```
class Node { nodes() { return 1 + this.left.nodes() + this.right.nodes(); } }
class Leaf { nodes() { return 1; } }
new Node(new Leaf(1), new Leaf(2)).nodes(); //=> 3 (constructors implicit)
```

An alternative “ad-hoc” representation uses plain JavaScript Objects with explicit tags. Then, the method becomes a recursive function that explicitly takes a tree as input. We trade the extensibility and (presumably) speed of a method for a simple, reversible serialization to JSON.

```
function nodes(t) { switch t.op {
  case "node": return 1 + nodes(t.left) + nodes(t.right);
  case "leaf": return 1; } }
nodes({op: "node", left:{op: "leaf", val: 1}, right:{op: "leaf", val: 2}})//=>3
```

Now, consider the problem of inferring type annotations for these programs. The class-based representation is idiomatic to popular dynamic languages like JavaScript and Python, and so many existing solutions support it. For example, TypeWiz [Shaked 2018] uses dynamic analysis to generate the following TypeScript annotations from the above example execution of nodes.

```
class Node { public left: Leaf; public right: Leaf; ... }
class Leaf { public val: number; ... }
```

The intuition behind inferring such a type is straightforward. For example, an instance of Leaf was observed in Node’s left field, and so the nominal type Leaf is used for its annotation.

The second “ad-hoc” style of programming seems peculiar in JavaScript, Python, and, indeed, object-oriented style in general. Correspondingly, existing state-of-the-art automatic annotation tools are not designed to support them. There are several ways to trivially handle such cases. Some enumerate the tree representation “verbatim” in a union, like TypeWiz [Shaked 2018].

```
50 function nodes(t: {left: {op: string, val: number}, op: string,
51                    right: {op: string, val: number}}
52                    | {op: string, val: number}) ...
53
54
55
56
57
```

Others “discard” most (or all) structure, like Typette [Hassan et al. 2018] and PyType [Google 2018] for Python.

```
60 def nodes(t: Dict[(Sequence, object)]) -> int: ... # Typette
61 def nodes(t) -> int: ... # PyType
62
```

Each annotation is clearly insufficient to meaningfully check both the function definition and valid usages. To show a desirable annotation for the “ad-hoc” program, we port it to Clojure [Hickey 2008], where it enjoys full support from the built-in runtime verification library `clojure.spec` and primary optional type system Typed Clojure [Bonnaire-Sergeant et al. 2016].

```
63 (defn nodes [t] (case (:op t)
64                   :node (+ 1 (nodes (:left t)) (nodes (:right t)))
65                   :leaf 1))
66
67 (nodes {:op :node, :left {:op :leaf, :val 1}, :right {:op :leaf, :val 2}}) ;=>3
68
69
70
71
```

Making this style viable requires a harmony of language features, in particular to support programming with functions and immutable values, but none of which comes at the expense of object-orientation. Clojure is hosted on the Java Virtual Machine and has full interoperability with Java objects and classes—even Clojure’s core design embraces object-orientation by exposing a collection of Java interfaces to create new kinds of data structures. The `{k v ...}` syntax creates a persistent and immutable Hash Array Mapped Trie [Bagwell 2001], which can be efficiently manipulated by dozens of built-in functions. The leading colon syntax like `:op` creates an interned *keyword*, which are ideal for map keys for their fast equality checks, and also look themselves up in maps when used as functions (e.g., `(:op t)` is like JavaScript’s `t.op`). *Multimethods* regain the extensibility we lost when abandoning methods, like the following.

```
82 (defmulti nodes-mm :op)
83 (defmethod nodes-mm :node [t] (+ 1 (nodes-mm (:left t)) (nodes-mm (:right t))))
84 (defmethod nodes-mm :leaf [t] 1)
85
```

On the type system side, Typed Clojure supports a variety of heterogeneous types, in particular for maps, along with occurrence typing [Tobin-Hochstadt and Felleisen 2010] to follow local control flow. Many key features come together to represent our “ad-hoc” binary tree as the following type.

```
86 (defalias Tree
87   (U '{:op ':node, :left Tree, :right Tree}
88       '{:op ':leaf, :val Int}))
89
```

The **defalias** form introduces an equi-recursive type alias `Tree`, **U** a union type, `'{:kw Type ...}` for heterogeneous keyword map types, and `' :node` for keyword singleton types. With the following function annotation, Typed Clojure can intelligently type check the definition and usages of `nodes`.

```
90 (ann nodes [Tree -> Int])
91
92
93
94
95
96
97
98
```

This (manually written) Typed Clojure annotation involving `Tree` is significantly different from TypeWiz’s “verbatim” annotation for nodes. First, it is recursive, and so supports trees of arbitrary depth (TypeWiz’s annotation supports trees of height  $< 3$ ). Second, it uses singleton types `' :leaf` and `' :node` to distinguish each case (TypeWiz upcasts “leaf” and “node” to string). Third, the tree type is factored out under a name to enhance readability and reusability. On the other end of the spectrum, the “discarding” annotations of Typette and PyType are too imprecise to use meaningfully (they include trees of arbitrary depth, but also many other values).

The challenge we overcome in this research is to automatically generate annotations like Typed Clojure’s `Tree`, in such a way that the ease of manual amendment is only mildly reduced by unresolvable ambiguities and incomplete data collection.

## 2 OVERVIEW

We demonstrate our approach by synthesizing a Typed Clojure annotation for nodes. The following presentation is somewhat loose to keep from being bogged down by details—interested readers may follow the pointers to subsequent sections where they are made precise.

We use dynamic analysis to observe the execution of functions, so we give an explicit test suite for nodes.

```
(def t1 {:op :node, :left {:op :leaf, :val 1}, :right {:op :leaf, :val 2}})
(deftest nodes-test (is (= (nodes t1) 3)))
```

The first step is the instrumentation phase (formalized in Section 3.1), which monitors the inputs and outputs of nodes by redefining it to use the track function like so (where `<nodes-body>` begins the `case` expression of the original nodes definition):

```
(def nodes (fn [t'] (track ((fn [t] <nodes-body>) (track t ' ['nodes :dom])
                          ['nodes :rng])))
```

The track function (given later in Figure 2) takes a value to track and a *path* that represents its origin, and returns an instrumented value along with recording some runtime samples about the value. A path is represented as a vector of *path elements*, and describes the source of the value in question. For example, `(track 3 ['nodes :rng])` returns 3 and records the sample `Int['nodes :rng]` which says “Int was recorded at nodes’s range.” Running our test suite `nodes-test` with an instrumented nodes results in more samples like this, most which use the path element `{:key :kw}` which represents a map lookup on the `:kw` entry.

```
' :leaf['nodes :dom {:key :op}]      ' :node['nodes :dom {:key :op}]      ?['nodes :dom {:key :val}]
      ?['nodes :dom {:key :left}]      ?['nodes :dom {:key :right}]
' :leaf['nodes :dom {:key :left} {:key :op}]      ' :leaf['nodes :dom {:key :right} {:key :op}]
  Int['nodes :dom {:key :left} {:key :val}]      Int['nodes :dom {:key :right} {:key :val}]
```

Now, our task is to transform these samples into a readable and useful annotation. This is the function of the inference phase (formalized in Section 3.2), which is split into three passes: first it generates a naive type from samples, then it combines types that occur syntactically near each other (“squash locally”), and then aggressively across different function annotations (“squash globally”).

The initial naive type generated from these samples resembles TypeWiz’s “verbatim” annotation given in Section 1, except the `?` placeholder represents incomplete information about a path (this process is formalized as `toEnv` in Figure 3).

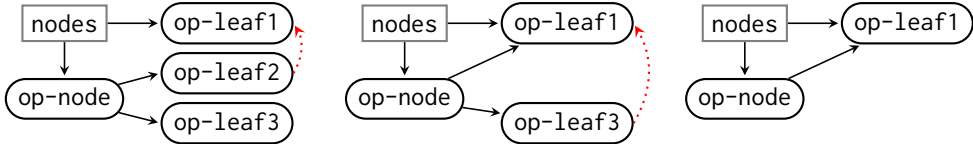
```
(ann nodes [(U '{:op ' :leaf, :val ?} '{:op ' :node,
                                         :left '{:op ' :leaf, :val Int},
                                         :right '{:op ' :leaf, :val Int})] -> Int])
```

Next, the two “squashing” phases. The intuition behind both are based on seeing types as directed graphs, where vertices are type aliases, and an edge connects two vertices  $u$  and  $v$  if  $u$  is mentioned in  $v$ 's type.

Local squashing (squashLocal in Figure 4) constructs such a graph by creating type aliases from map types using a post-order traversal of the original types. In this example, the previous annotations become:

```
(defalias op-leaf1 '{:op ':leaf, :val ?})
(defalias op-leaf2 '{:op ':leaf, :val Int})
(defalias op-leaf3 '{:op ':leaf, :val Int})
(defalias op-node '{:op ':node, :left op-leaf2, :right op-leaf3})
(ann nodes [(U op-leaf1 op-node) -> Int])
```

As a graph, this becomes the left-most graph below. The dotted edge from `op-leaf2` to `op-leaf1` signifies that they are to be merged, based on the similar structure of the types they point to.



After several merges (reading the graphs left-to-right), local squashing results in the following:

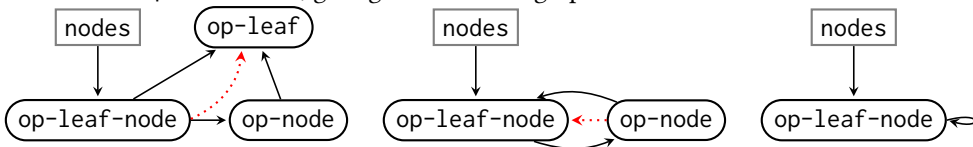
```
(defalias op-leaf '{:op ':leaf, :val Int})
(defalias op-node '{:op ':node, :left op-leaf, :right op-leaf})
(ann nodes [(U op-leaf op-node) -> Int])
```

All three duplications of the `:leaf` type in the naive annotation have been consolidated into their own name, with the `?` placeholder for the `:val` entry being absorbed into `Int`.

Now, the global squashing phase (squashGlobal in Figure 5) proceeds similarly, except the notion of a vertex is expanded to also include *unions* of map types, calculated, again, with a post-order traversal of the types giving:

```
(defalias op-leaf '{:op ':leaf, :val Int})
(defalias op-node '{:op ':node, :left op-leaf, :right op-leaf})
(defalias op-leaf-node (U op-leaf op-node))
(ann nodes [op-leaf-node -> Int])
```

This creates `op-leaf-node`, giving the left-most graph below.



Now, type aliases are merged based on overlapping *sets* of top-level keysets and likely tags. Since `op-leaf` and `op-leaf-node` refer to maps with identical keysets (`:op` and `:val`) and whose likely tags agree (the `:op` entry is probably a tag, and they are both `:leaf`), they are merged and all occurrences of `op-leaf` are renamed to `op-leaf-node`, creating a *mutually recursive type* between the remaining aliases in the middle graph:

```
(defalias op-node '{:op ':node, :left op-leaf-node, :right op-leaf-node})
(defalias op-leaf-node (U '{:op ':leaf, :val Int} op-node))
(ann nodes [op-leaf-node -> Int])
```

In the right-most graph, the aliases `op-node` and `op-leaf-node` are merged for similar reasons:

```

197 (defalias op-leaf-node (U '{:op ':leaf, :val Int}
198                          '{:op ':node, :left op-leaf-node, :right op-leaf-node}))
199
200 (ann nodes [op-leaf-node -> Int])
201

```

All that remains is to choose a recognizable name for the alias. Since all its top-level types seem to use the `:op` entry for tags, we choose the name `Op` and output the final annotation:

```

204 (defalias Op (U '{:op ':leaf, :val Int}
205               '{:op ':node, :left Op, :right Op}))
206 (ann nodes [Op -> Int])
207

```

The rest of the porting workflow involves the programmer repeatedly type checking their code and gradually tweaking the generated annotations until they type check. It turns out that this annotation immediately type checks the definition of nodes and all its valid usages, so we turn to a more complicated function `visit-leaf` to demonstrate a typical scenario.

```

212 (defn visit-leaf "Updates :leaf nodes in tree t with function f."
213   [f t] (case (:op t)
214           :node (assoc t :left (visit-leaf f (:left t))
215                           :right (visit-leaf f (:right t)))
216           :leaf (f t)))
217

```

This higher-order function uses `assoc` to associate new children as it recurses down a given tree to update leaf nodes with the provided function. The following test simply increments the leaf values of the previously-defined `t1`.

```

221 (deftest visit-leaf-test
222   (is (= (visit-leaf (fn [leaf] (assoc leaf :val (inc (:val leaf)))) t1)
223         {:op :node, :left {:op :leaf, :val 2}, :right {:op :leaf, :val 3}})))
224

```

Running this test under instrumentation yields some interesting runtime samples whose calculation is made efficient by space-efficient tracking (Section 4.1), which ensures a function is not repeatedly tracked unnecessarily. The following two samples demonstrate how to handle multiple arguments (by parameterizing the `:dom` path element) and higher-order functions (by nesting `:dom` or `:rng` path elements).

```

230 ':leaf['visit-leaf {:dom 1} {:key :op}]      ':leaf['visit-leaf {:dom 0} {:dom 0} {:key :op}]
231

```

Here is our automatically generated initial annotation.

```

232 (defalias Op (U '{:op ':leaf, :val t/Int} '{:op ':node, :left Op, :right Op}))
233 (ann visit-leaf [[Op -> Any] Op -> Any])
234

```

Notice the surprising occurrences of `Any`. They originate from `?` placeholders due to the lazy tracking of maps (Section 4.2). Since `visit-leaf` does not traverse the results of `f`, nor does anything traverse `visit-leaf`'s results (hash-codes are used for equality checking) neither tracking is realized. Also notice the first argument of `visit-leaf` is underprecise. These could trigger type errors on usages of `visit-leaf`, so manual intervention is needed (highlighted). We factor out and use a new alias `Leaf` and replace occurrences of `Any` with `Op`.

```

241 (defalias Leaf '{:op ':leaf, :val Int})
242 (defalias Op (U Leaf '{:op ':node, :left Op, :right Op}))
243 (ann visit-leaf [[Leaf -> Op] Op -> Op])
244
245

```

246	$v ::= n \mid k \mid [\lambda x.e, \rho]_c \mid \{\overline{k} \ v\} \mid c$	Values
247	$e ::= x \mid v \mid (\mathbf{track} \ e \ \pi) \mid \lambda x.e \mid \{\overrightarrow{e} \ \bar{e}\} \mid (e \ \bar{e})$	Expressions
248	$\rho ::= \{\overline{x} \mapsto v\}$	Runtime environments
249	$l ::= x \mid \mathbf{dom} \mid \mathbf{rng} \mid \mathbf{key}_m(k)$	Path Elements
250	$\pi ::= \bar{l}$	Paths
251	$r ::= \{\overline{\tau} \ \pi\}$	Inference results
252	$\tau, \sigma ::= \mathbb{N} \mid [\tau \rightarrow \tau] \mid (\mathbf{HMap}_o^m) \mid (\cup \ \bar{\tau})$	
253	$\quad \mid a \mid k \mid \mathbf{K} \mid \top \mid (\mathbf{Map} \ \tau \ \tau) \mid ?$	Types
254	$\Gamma ::= \{\overline{x} : \tau\}$	Type environments
255	$m, o ::= \{k \ \tau\}$	HMap entries
256	$A ::= \{\overline{a} \mapsto \tau\}$	Type alias environments
257	$\Delta ::= (A, \Gamma)$	Annotation environments

Fig. 1. Syntax of Terms, Types, Inference results, and Environments for  $\lambda_{\text{track}}$

We measure the success of our workflow by using it to type check real Clojure programs. Experiment 1 (Section 5.1) manually inspects a selection of inferred types. Experiment 2 (Section 5.2) classifies and quantifies the kinds of changes needed. Experiment 3 (Section 5.3) enforces initial annotations at runtime to ensure they are meaningfully underprecise.

### 3 FORMALISM

We present  $\lambda_{\text{track}}$ , an untyped  $\lambda$ -calculus describing the essence of our approach to automatic annotations. We split our model into two phases: the collection phase `collect` that runs an instrumented program and collects observations, and an inference phase `infer` that derives type annotations from these observations that can be used to automatically annotate the program.

We define the top-level driver function `annotate` that connects both pieces. It says, given a program  $e$  and top-level variables  $\bar{x}$  to infer annotations for, return an annotation environment  $\Delta$  with possible entries for  $\bar{x}$  based on observations from evaluating an instrumented  $e$ .

$$\begin{aligned} \text{annotate} &: e, \bar{x} \rightarrow \Delta \\ \text{annotate} &= \text{infer} \circ \text{collect} \end{aligned}$$

To contextualize the presentation of these phases, we begin a running example: inferring the type of a top-level function  $f$ , that takes a map and returns its `:a` entry, based on the following usage.

```
define f = λm.(get m :a)
(f {:a 42}) => 42
```

Plugging this example into our driver function we get a candidate annotation for  $f$ :

$$\text{annotate}((f \ \{\text{:a} \ 42\}), [f]) = \{f : \{\text{:a} \ \mathbb{N}\} \rightarrow \mathbb{N}\}$$

#### 3.1 Collection phase

Now that we have a high-level picture of how these phases interact, we describe the syntax and semantics of  $\lambda_{\text{track}}$ , before presenting the details of `collect`. Figure 1 presents the syntax of  $\lambda_{\text{track}}$ . Values  $v$  consist of numbers  $n$ , Clojure-style keywords  $k$ , closures  $[\lambda x.e, \rho]_c$ , constants  $c$ , and keyword keyed hash maps  $\{\overline{k} \ v\}$ .

Expressions  $e$  consist of variables  $x$ , values, functions, maps, and function applications. The special form `(track e π)` observes  $e$  as related to path  $\pi$ . Paths  $\pi$  record the source of a runtime

value with respect to a sequence of path elements  $l$ , always starting with a variable  $x$ , and are read left-to-right. Other path elements are a function domain **dom**, a function range **rng**, and a map entry  $\mathbf{key}_{\overline{k_1}}(k_2)$  which represents the result of looking up  $k_2$  in a map with keyset  $\overline{k_1}$ .

Inference results  $\{\overline{\tau_\pi}\}$  are pairs of paths  $\pi$  and types  $\tau$  that say the path  $\pi$  was observed to be type  $\tau$ . Types  $\tau$  are numbers  $\mathbb{N}$ , function types  $[\tau \rightarrow \tau]$ , ad-hoc union types  $(\bigcup \tau \tau)$ , type aliases  $a$ , and unknown type  $?$  that represents a temporary lack of knowledge during the inference process. Heterogeneous keyword map types  $\{\overline{k \tau}\}$  for now represent a series of required keyword entries—we will extend them to have optional entries in later phases.

The big-step operational semantics  $\rho \vdash e \Downarrow v ; r$  (Figure 2) says under runtime environment  $\rho$  expression  $e$  evaluates to value  $v$  with inference results  $r$ . Most rules are standard, with extensions to correctly propagate inference results  $r$ . B-Track is the only interesting rule, which instruments its fully-evaluated argument with the track metafunction.

The metafunction  $\mathbf{track}(v, \pi) = v' ; r$  (Figure 2) says if value  $v$  occurs at path  $\pi$ , then return a possibly-instrumented  $v'$  paired with inference results  $r$  that can be immediately derived from the knowledge that  $v$  occurs at path  $\pi$ . It has a case for every kind of value. The first three cases records the number input as type  $\mathbb{N}$ . The fourth case, for closures, returns a wrapped value resembling higher-order function contracts [Findler and Felleisen 2002], but we track the domain and range rather than verify them. The remaining rules case, for maps, recursively tracks each map value, and returns a map with possibly wrapped values. Immediately accessible inference results are combined and returned. A specific rule for the empty map is needed because we otherwise only rely on recursive calls to **track** to gather inference results—in the empty case, we have no data to recur on.

Now we have sufficient pieces to describe the initial collection phase of our model. Given an expression  $e$  and variables  $\overline{x}$  to track,  $\mathbf{instrument}(e, \overline{x}) = e'$  returns an instrumented expression  $e'$  that tracked usages of  $\overline{x}$ . It is defined via capture-avoiding substitution:

$$\mathbf{instrument}(e, \overline{x}) = e[\overline{\mathbf{track} \ x \ [x]} / \overline{x}]$$

Then, the overall collection phase  $\mathbf{collect}(e, \overline{x}) = r$  says, given an expression  $e$  and variables  $\overline{x}$  to track, returns inference results  $r$  that are the results of evaluating  $e$  with instrumented occurrences of  $\overline{x}$ . It is defined as:

$$\mathbf{collect}(e, \overline{x}) = r, \text{ where } \vdash \mathbf{instrument}(e, \overline{x}) \Downarrow v ; r$$

For our running example of collecting for the program  $(f \{ :a \ 42 \})$ , we instrument the program by wrapping occurrences of  $f$  with **track** with path  $[f]$ .

$$\mathbf{instrument}((f \{ :a \ 42 \}), [f]) = ((\mathbf{track} \ f \ [f]) \{ :a \ 42 \})$$

Then we evaluate the instrumented program and derive two inference results (colored in red for readability):

$$\vdash ((\mathbf{track} \ f \ [f]) \{ :a \ 42 \}) \Downarrow 42 ; \{ \mathbf{N}_{[f, \mathbf{dom}, \mathbf{key}(:a)]}, \mathbf{N}_{[f, \mathbf{rng}]} \}$$

Here is the full derivation:

$$\begin{aligned} &\Rightarrow ((\mathbf{track} \ f \ [f]) \{ :a \ 42 \}) \\ &\Rightarrow (\mathbf{track} \ (\mathbf{get} \ (\mathbf{track} \ \{ :a \ 42 \} \ [f, \mathbf{dom}]) \ :a) \ [f, \mathbf{rng}]) \\ &\Rightarrow (\mathbf{track} \ (\mathbf{get} \ \{ :a \ 42 \} ; \{ \mathbf{N}_{[f, \mathbf{dom}, \mathbf{key}(:a)]} \} \ :a) \ [f, \mathbf{rng}]) \\ &\Rightarrow (\mathbf{track} \ 42 ; \{ \mathbf{N}_{[f, \mathbf{dom}, \mathbf{key}(:a)]} \} \ [f, \mathbf{rng}]) \\ &\Rightarrow 42 ; \{ \mathbf{N}_{[f, \mathbf{dom}, \mathbf{key}(:a)]}, \mathbf{N}_{[f, \mathbf{rng}]} \} \end{aligned}$$

Notice that intermediate values can have inference results (colored) attached to them with a semicolon, and the final value has inference results about both  $f$ 's domain and range.

<p>344</p> <p>345 B-TRACK</p> <p>346 <math>\rho \vdash e \Downarrow v ; r</math></p> <p>347 <math>\text{track}(v, \pi) = v' ; r'</math></p> <p>348 <math>\frac{}{\rho \vdash (\mathbf{track} e \pi) \Downarrow v' ; r \cup r'}</math></p> <p>349</p> <p>350</p> <p>351</p> <p>352 B-VAL</p> <p>353 <math>\rho \vdash v \Downarrow v ; \{\}</math></p> <p>354</p> <p>355</p>	<p>B-APP</p> <p><math>\rho \vdash e_1 \Downarrow [\lambda x.e, \rho']_c ; r_1</math></p> <p><math>\rho \vdash e_2 \Downarrow v ; r_2</math></p> <p><math>\rho'[x \mapsto v] \vdash e \Downarrow v' ; r_3</math></p> <p><math>\rho \vdash (e_1 e_2) \Downarrow v' ; \bigcup \bar{r}_i</math></p> <p>B-CLOS</p> <p><math>\rho \vdash \lambda x.e \Downarrow [\lambda x.e, \rho]_c ; \{\}</math></p> <p>B-DELTA</p> <p><math>\rho \vdash e \Downarrow c ; r_1</math></p> <p><math>\rho \vdash e' \Downarrow v ; r'</math></p> <p><math>\delta(c, \bar{v}) = v' ; r_2</math></p> <p><math>\frac{}{\rho \vdash (e \bar{e}') \Downarrow v' ; r \cup r'}</math></p> <p>B-VAR</p> <p><math>\rho \vdash x \Downarrow \rho(x) ; \{\}</math></p>	<p>356 <math>\text{track}(n, \pi) = n ; \{\mathbf{N}_\pi\}</math></p> <p>357 <math>\text{track}(k, \pi) = k ; \{\mathbf{K}_\pi\}</math></p> <p>358 <math>\text{track}(c, \pi) = c ; \{\}</math></p> <p>359 <math>\text{track}([\lambda x.e, \rho]_c, \pi) = [e', \rho]_c ; \{\}</math></p> <p>360 where <math>y</math> is fresh,</p> <p>361 <math>e' = \lambda y.(\mathbf{track} ((\lambda x.e) (\mathbf{track} y \pi :: [\mathbf{dom}])))</math></p> <p>362 <math>\pi :: [\mathbf{rng}]</math></p> <p>363 <math>\text{track}(\{\}, \pi) = \{\} ; \{\{\}\}_\pi</math></p> <p>364 <math>\text{track}(\{k_1 k_2 k v\}, \pi) = \{k_1 k_2 k v'\} ; \bigcup r</math></p> <p>365 where <math>\text{track}(v, \pi :: [\mathbf{key}_{\{k_1 k_2 k ?\}}(k)]) = v' ; r</math></p> <p>366</p> <p>367 <math>\delta(\text{assoc}, \{\overline{k v}, k', v'\} = \{\overline{k v}\}[k' \mapsto v'] ; \{\}</math></p> <p>368 <math>\delta(\text{get}, \{k v, k' v'\}, k) = v ; \{\}</math></p> <p>369 <math>\delta(\text{dissoc}, \{k v, k' v'\}, k) = \{k' v'\} ; \{\}</math></p> <p>370</p> <p>371</p> <p>372</p> <p>373</p> <p>374</p>
---	---	--

Fig. 2. Operational semantics,  $\text{track}(v, \pi) = v ; r$  and constants

### 3.2 Inference phase

After the collection phase, we have a collection of inference results  $r$  which can be passed to the metafunction  $\text{infer}(r) = \Delta$  to produce an annotation environment:

$$\begin{aligned} \text{infer} &: r \rightarrow \Delta \\ \text{infer} &= \text{inferRec} \circ \text{toEnv} \end{aligned}$$

The first pass  $\text{toEnv}(r) = \Gamma$  generates an initial type environment from inference results  $r$ . The second pass

$$\text{squashLocal}(\Gamma) = \Delta'$$

creates individual type aliases for each HMap type in  $\Gamma$  and then merges aliases that both occur inside the same nested type into possibly recursive types. The third pass  $\text{squashGlobal}(\Delta) = \Delta'$  merges type aliases in  $\Delta$  based on their similarity.

**3.2.1 Pass 1: Generating initial type environment.** The first pass is given in Figure 3. The entry point  $\text{toEnv}$  folds over inference results to create an initial type environment via update. This style



$$\begin{array}{l}
393 \quad \sqcup : \tau, \tau \rightarrow \tau \\
394 \quad (\bigcup \bar{\sigma}) \sqcup \tau = (\bigcup \overline{\sigma \sqcup \tau}) \\
395 \quad \tau \sqcup (\bigcup \bar{\sigma}) = (\bigcup \overline{\sigma \sqcup \tau}) \\
396 \quad ? \sqcup \tau = \tau \\
397 \quad \tau \sqcup ? = \tau \\
398 \quad [\tau_1 \rightarrow \sigma_1] \sqcup [\tau_2 \rightarrow \sigma_2] = [\tau_1 \sqcup \tau_2 \rightarrow \sigma_1 \sqcup \sigma_2] \\
399 \quad (\text{HMap}_{o_1}^{m_1}) \sqcup (\text{HMap}_{o_2}^{m_2}) = \frac{(\text{HMap}_{o_1}^{m_1}) \sqcup^H (\text{HMap}_{o_2}^{m_2})}{(k, k_i) \in m_i \Rightarrow \bar{k}_{i-1} = k_i} \\
400 \quad \tau \sqcup \sigma = (\bigcup \tau \sigma), \text{ otherwise} \\
401 \\
402 \\
403 \quad \text{fold} : \forall \alpha, \beta. (\alpha, \beta \rightarrow \alpha), \alpha, \bar{\beta} \rightarrow \alpha \quad \text{update} : \Gamma, \tau_\pi \rightarrow \Gamma \\
404 \quad \text{fold}(f, a_0, \bar{b}^n) = a_n \quad \text{update}(\Gamma, \tau_{\pi::[\text{key}_{\{\bar{k}' \sigma\}}(k)]}) = \text{update}(\Gamma, \overline{\{\bar{k}' \sigma k \tau\}}_\pi) \\
405 \quad \text{where } \bar{a}_i = f(a_{i-1}, b_i)^{1 \leq i \leq n} \quad \text{update}(\Gamma, \tau_{\pi::[\text{dom}]}) = \text{update}(\Gamma, [\tau \rightarrow ?]_\pi) \\
406 \quad \quad \quad \text{update}(\Gamma, \tau_{\pi::[\text{rng}]}) = \text{update}(\Gamma, [? \rightarrow \tau]_\pi) \\
407 \quad \quad \quad \text{update}(\Gamma, \tau_{\pi::[\text{rng}]}) = \text{update}(\Gamma, [? \rightarrow \tau]_\pi) \\
408 \quad \text{toEnv} : r \rightarrow \Gamma \quad \text{update}(\Gamma[x \mapsto \sigma], \tau_{[x]}) = \Gamma[x \mapsto \tau \sqcup \sigma] \\
409 \quad \text{toEnv}(r) = \text{fold}(\text{update}, \{\}, r) \quad \text{update}(\Gamma, \tau_{[x]}) = \Gamma[x \mapsto \tau] \\
410 \\
411 \\
412 \\
413 \\
414 \\
415 \\
416 \\
417 \\
418 \\
419 \\
420 \\
421 \\
422 \\
423 \\
424 \\
425 \\
426 \\
427 \\
428 \\
429 \\
430 \\
431 \\
432 \\
433 \\
434 \\
435 \\
436 \\
437 \\
438 \\
439 \\
440 \\
441
\end{array}$$

Fig. 3. Definition of  $\text{toEnv}(r) = \Gamma$ 

is inspired by occurrence typing [Tobin-Hochstadt and Felleisen 2010], from which we also borrow the concepts of paths into types.

We process paths right-to-left in `update`, building up types from leaves to root, before joining the fully constructed type with the existing type environment via  $\sqcup$ . The first case handles the **key** path element. The extra map of type information preserves both keyset information and any entries that might represent tags (populated by the final case of **tracK**, Figure 2). This information helps us avoid prematurely collapsing tagged maps, by the side condition of the  $\text{HMap} \sqcup$  case. The  $\sqcup^H$  metafunction aggressively combines two  $\text{HMaps}$ —required keys in both maps are joined and stay required, otherwise keys become optional.

The second and third update cases update the domain and range of a function type, respectively. The  $\sqcup$  case for function types joins covariantly on the domain to yield more useful annotations. For example, if a function accepts  $\mathbf{N}$  and  $\mathbf{K}$ , it will have type  $[\mathbf{N} \rightarrow ?] \sqcup [\mathbf{K} \rightarrow ?] = [(\bigcup \mathbf{N} \mathbf{K}) \rightarrow ?]$ .

Returning to our running example, we now want to convert our inference results

$$r = \{\mathbf{N}_{[f, \text{dom}, \text{key}:(a)]}, \mathbf{N}_{[f, \text{rng}]}\}.$$

into a type environment. Via  $\text{toEnv}(r)$ , we start to trace  $\text{update}(\{\}, \mathbf{N}_{[f, \text{dom}, \text{key}:(a)]})$

**3.2.2 Pass 2: Squash locally.** We now describe the algorithm for generating recursive type aliases. The first step `squashLocal` creates recursive types from directly nested types. It folds over each type in the type environment, first creating aliases with `aliasHMap`, and then attempting to merge these aliases by `squashAll`.

A type is aliased by `aliasHMap` either if it is a union containing a  $\text{HMap}$ , or a  $\text{HMap}$  that is not a member of a union. While we will use the structure of  $\text{HMaps}$  to determine when to create a recursive type, keeping surrounding type information close to  $\text{HMaps}$  helps create more compact and readable recursive types. The implementation uses a post-order traversal via `postwalk`, which also threads an annotation environment as it applies the provided function.

```

442 aliasHMap : Δ, τ → (Δ, τ)
443 aliasHMap(Δ, τ) = postwalk(Δ, τ, f)
444   where f(Δ, (HMapm1)) = reg(Δ, (HMapm2))
445         f(Δ, (∪  $\bar{\tau}$ )) = reg(Δ, (∪ resolve( $\bar{\tau}$ ))),
446         if a ∈  $\bar{\tau}$ 
447         f(Δ, τ) = (Δ, τ), otherwise
448
449 aliases : τ →  $\bar{a}$ 
450 aliases(a) = [a]
451 aliases( $\tau(\bar{\sigma})$ ) = ∪ aliases( $\sigma$ )
452
453 postwalk : Δ, τ, (Δ, τ → (Δ, τ)) → (Δ, τ)
454 postwalk(Δ0, τ( $\bar{\sigma}^n$ ), w) = w(Δn, τ( $\bar{\sigma}'$ ))
455   where (Δi,  $\sigma'_i$ ) = postwalk(Δi-1, σi, w)
456
457 mergeAliases : Δ,  $\bar{a}$  → Δ
458 mergeAliases(Δ, []) = Δ
459 mergeAliases(Δ, [a1...an]) = Δ[ $\overline{a_i \mapsto a_1}$ ][a1 ↦ σ]
460   where σ = ∪ f(resolve(Δ, ai))[a1/ai]
461         f(a') = (∪), if a' ∈  $\bar{a}$ 
462         f((∪  $\bar{\tau}$ )) = (∪ f( $\bar{\tau}$ ))
463         f(τ) = τ, otherwise
464
465 squashLocal : Γ → Δ
466 squashLocal(Γ) = fold(h, ({}, {}), Γ)
467   where h(Δ, x : τ) = Δ2[x ↦ τ2]
468         where (Δ1, τ1) = aliasHMap(Δ, τ)
469               (Δ2, τ2) = squashAll(Δ1, τ1)
470
471 reg : Δ, τ → (Δ, τ)
472 reg(Δ, τ) = (Δ[a ↦ τ], a), where a is fresh
473
474 resolve : Δ, τ → τ
475 resolve(Δ, a) = resolve(Δ[a])
476 resolve(Δ, τ) = τ, otherwise
477
478 squashAll : Δ, τ → Δ
479 squashAll(Δ0, τ) = Δn
480   where  $\bar{a}^n$  = aliases(τ)
481         Δi = squash(Δi-1, [ai], [])
482
483 squash : Δ,  $\bar{a}$ ,  $\bar{a}$  → Δ
484 squash(Δ, [], d) = Δ
485 squash(Δ, a1 :: w, d) = squash(Δ', w ∪ as, d ∪ {a1})
486   where as = aliases(Δ[a1]) \ d
487         ap = d \ {a1}
488         f(Δ, a2) = if  $\neg$ merge?(resolve(Δ, a)),
489                   then Δ
490                   else mergeAliases(Δ,  $\bar{a}_i$ )
491         Δ' = if a ∈ d, then Δ,
492              else fold(f, Δ, ap ∪ as)
493
494 merge? :  $\bar{\tau}$  → Bool
495 merge?((HMapmioi)) = ∃k.(k, ki) ∈ mi
496 merge?( $\bar{\tau}$ ) = F, otherwise

```

Fig. 4. Definition of squashLocal(Γ) = Δ

Then, squashAll follows each alias  $a_i$  reachable from the type environment and attempts to merge it with any alias reachable from  $a_i$ . The squash function maintains a set of already visited aliases to avoid infinite loops.

The logic for merging aliases is contained in mergeAliases. Merging  $a_2$  into  $a_1$  involves mapping  $a_2$  to  $a_1$  and  $a_1$  to the join of both definitions. Crucially, before joining, we rename occurrences of  $a_2$  to  $a_1$ . This avoids a linear increase in the width of union types, proportional to the number of merged aliases. The running time of our algorithm is proportional to the width of union types (due to the quadratic combination of unions in the join function) and this optimization greatly helped the running time of several benchmarks. To avoid introducing infinite types, top-level references to other aliases we are merging with are erased with the helper  $f$ .

The merge? function determines whether two types are related enough to warrant being merged. We present our current implementation, which is simplistic, but is fast and effective in practice, but many variations are possible. Aliases are merged if they are all HMaps (not contained in unions), that contain a keyword key in common, with possibly disjoint mapped values. For example, our

```

491 req : Δ, a → m
492 req(Δ, a) = req(Δ, Δ[a])
493 req(Δ, (HMapom)) = m
494
495
496 squashHorizontally : Δ → Δ
497 squashHorizontally(Δ) =
498   fold(mergeAliases, Δ, groupSimilarReq(Δ))
499
500 squashGlobal : Δ → Δ
501 squashGlobal =
502   squashHorizontally ∘ aliasSingleHMap
503
504
505 groupSimilarReq : Δ →  $\bar{a}$ 
506 groupSimilarReq(Δ) = [ $\bar{a} | \bar{k} \in \text{dom}(r), \bar{a} = \text{remDiffTag}(\text{similarReq}(\bar{k}))$ ]
507   where r = {( $\bar{k}, \bar{a}$ ) | (HMapok τ) ∈ rng(Δ[A]),  $\bar{a} = \text{matchingReq}(\bar{k})$ }
508   matchingReq( $\bar{k}$ ) = [a | (a, (HMapom)) ∈ Δ]
509   similarReq( $\bar{k}$ ) = [a |  $\bar{k}'^n \subseteq \bar{k}^m, m - n \leq \text{thres}(m), a \in r[\bar{k}']$ ]
510   remDiffTag( $\bar{a}$ ) = [ $a' | a' \in \bar{a}, \text{if } (k, k') \in \text{req}(\Delta, a') \text{ and } \bigvee (k, k'') \in \text{req}(\Delta, a) \text{ then } \bar{k}' = \bar{k}''$ ]
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

```

Fig. 5. Definition of  $\text{squashGlobal}(\Delta) = \Delta'$ 

opening example has the `:op` key mapped to either `:leaf` or `:node`, and so aliases for each map would be merged. Notice again, however, the join operator does not collapse differently-tagged maps, so they will occur recursively in the resulting alias, but separated by union.

Even though this implementation of `merge?` does not directly utilize the aliased union types carefully created by `aliasHMap`, they still affect the final types. For example, squashing `T` in

```

521 (defalias T
522   (U nil '{:op :node :left '{:op :leaf ...} ...}))
523
524 results in
525 (defalias T
526   (U nil '{:op :node :left T ...} '{:op :leaf ...}))

```

rather than

```

529 (defalias T2 (U '{:op :node :left T ...}
530               '{:op :leaf ...}))
531 (defalias T (U nil T2))

```

An alternative implementation of `merge?` we experimented with included computing sets of keysets for each alias, and merging if the keysets overlapped. This, and many of our early experimentations, required expensive computations of keyset combinations and traversals over them that could be emulated with cruder heuristics like the current implementation.

**3.2.3 Pass 3: Squash globally.** The final step combines aliases without restriction on whether they occur “together”. This step combines type information between different positions (such as in

540 different arguments or functions) so that any deficiencies in unit testing coverage are massaged  
541 away.

542 The `squashGlobal` function is the entry point in this pass, and is similar in structure to the  
543 previous pass. It first creates aliases for each HMap via `aliasSingleHMap`. Then, HMap aliases are  
544 grouped and merged in `squashHorizontally`.

545 The `aliasSingleHMap` function first traverses the type environment to create HMap aliases  
546 via `singleHMap`, and binds the resulting environment as  $\Delta'$ . Then, alias environment entries are  
547 updated with `f`, whose first case prevents re-aliasing a top-level HMap, before we call `singleHMap`  
548 (`singleHMap`'s second argument accepts both  $x$  and  $a$ ). The  $\tau(\bar{\sigma})$  syntax represents a type  $\tau$  whose  
549 constructor takes types  $\bar{\sigma}$ .

550 After that, `squashHorizontally` creates groups of related aliases with `groupSimilarReq`. Each  
551 group contains HMap aliases whose required keysets are similar, but are never differently-tagged.  
552 The code creates a map `r` from keysets to groups of HMap aliases with that (required) keyset. Then,  
553 for every keyset  $\bar{k}$ , `similarReq` adds aliases to the group whose keysets are a subset of  $\bar{k}$ . The number  
554 of missing keys permitted is determined by `thres`, for which we do not provide a definition. Finally,  
555 `remDiffTag` removes differently-tagged HMaps from each group, and the groups are merged via  
556 `mergeAliases` as before.

557  
558 **3.2.4 Implementation.** Further passes are used in the implementation. In particular, we trim un-  
559 reachable aliases and remove aliases that simply point to another alias (like  $a_2$  in `mergeAliases`)  
560 between each pass.

## 561 4 EXTENSIONS

562  
563 Two optimizations are crucial for practical implementations of the collection phase. First, space-  
564 efficient tracking efficiently handles a common case with higher-order functions where the same  
565 function is tracked at multiple paths. Second, instead of tracking a potentially large value by  
566 eagerly traversing it, lazy tracking offers a pay-as-you-go model by wrapping a value and only  
567 tracking subparts as they are accessed. Both were necessary to collect samples from the compiler  
568 implementation we instrumented for Experiment 1 (Section 5.1) because it used many higher-order  
569 functions and its AST representation can be quite large which made it intractable to eagerly traverse  
570 each time it got passed to one of dozens of functions.

### 571 4.1 Space-efficient tracking

572  
573 To reduce the overhead of runtime tracking, we can borrow the concept of “space-efficient” contract  
574 checking from the gradual typing literature [Herman et al. 2010]. Instead of tracking just one path  
575 at once, a space-efficient implementation of `track` threads through a set of paths. When a tracked  
576 value flows into another tracked position, we extract the unwrapped value, and then our new  
577 tracked value tracks the paths that is the set of the old paths with the new path.

578 To model this, we introduce a new kind of value  $[e, \rho]_{c\bar{\pi}}^v$  that tracks old value  $v$  as new value  
579  $[e, \rho]_c$  with the paths  $\bar{\pi}$ . Proxy expressions are introduced when tracking functions, where instead  
580 of just returning a new wrapped function, we return a proxy. We can think of function proxies as a  
581 normal function with some extra metadata, so we can reuse the existing semantics for function  
582 application—in fact we can support space-efficient function tracking just by extending `track`.

583 We present the extension in Figure 6. The first two `track` rules simply make inference results for  
584 each of the paths. The next rule says that a bare closure reduces to a proxy that tracks the domain  
585 and range of the closure with respect to the list of paths. Attached to the proxy is everything needed  
586 to extend it with more paths, which is the role of the final rule. It extracts the original closure from  
587 the proxy and creates a new proxy with updated paths via the previous rule.

588

589  $v ::= \dots \mid [\lambda x.e, \rho]_c \overline{\pi}^{\lambda x.e, \rho} \quad \text{Values}$   
590  
591  $\text{track}(n, \overline{\pi}) = n ; \bigcup \overline{\{N_\pi\}}$   
592  $\text{track}(k, \overline{\pi}) = k ; \bigcup \overline{\{K_\pi\}}$   
593  $\text{track}([\lambda x.e, \rho]_c, \overline{\pi}) = [e', \rho]_c \overline{\pi}^{\lambda x.e, \rho} ; \{\}$   
594  $\text{where } y \text{ is fresh,}$   
595  $e' = \lambda y. (\text{track } ((\lambda x.e) (\text{track } y \ \overline{\pi} :: [\text{dom}])))$   
596  $\overline{\pi} :: [\text{rng}]$   
597  $\text{track}([e', \rho']_c \overline{\pi'}^{\lambda x.e, \rho}, \overline{\pi}) = \text{track}([\lambda x.e, \rho]_c, \overline{\pi} \cup \overline{\pi'})$   
598  
599  
600

Fig. 6. Space-efficient tracking extensions (changes)

## 4.2 Lazy tracking

605 Building further on the extension of space-efficient functions, we apply a similar idea for tracking  
606 maps. In practice, eagerly walking data structures to gather inference results is expensive. Instead,  
607 waiting until a data structure is used and tracking its contents lazily can help ease this tradeoff,  
608 with the side-effect that fewer inference results are discovered.

609 Figure 7 extends our system with lazy maps. We add a new kind of value  $\overline{\{k' \overline{v}\}}^{\overline{\{m \ \overline{\pi}\}}}$  that  
610 wraps a map  $\overline{\{k' \overline{v}\}}$  with tracking information. Keyword entries  $k'$  are associated with pairs of  
611 type information  $m$  with paths  $\overline{\pi}$ . The first **track** rule demonstrates how to create a lazily tracked  
612 map. We calculate the possibly tagged entries in our type information in advance, much like the  
613 equivalent rule in Figure 2, and store them for later use. Notice that non-keyword entries are not  
614 yet traversed, and thus no inference results are derived from them. The second **track** rule adds  
615 new paths to track.

616 The subtleties of lazily tracking maps lie in the  $\delta$  rules. The assoc and dissoc rules ensure we no  
617 longer track overwritten entries. Then, the get rules perform the tracking that was deferred from  
618 the **track** rule for maps in Figure 2 (if the entry is still tracked).

619 In our experience, some combination of lazy and eager tracking of maps strikes a good balance  
620 between performance overhead and quantity of inference results. Intuitively, if a function does not  
621 access parts of its argument, they should not contribute to that function's type signature. However,  
622 our inference algorithm combines information *across* function signatures to deduce useful, recursive  
623 type aliases. Some eager tracking helps normalize the quality of function annotations with respect  
624 to unit test coverage.

625 For example, say functions  $f$  and  $g$  operate on the same types of (deeply nested) arguments, and  
626  $f$  has complete test coverage (but does not traverse all of its arguments), and  $g$  has incomplete test  
627 coverage (but fully traverses its arguments). Eagerly tracking  $f$  would give better inference results,  
628 but lazily tracking  $g$  is more efficient. Forcing several layers of tracking helps strike this balance,  
629 which our implementation exposes as a parameter.

630 This can be achieved in our formal system by adding fuel arguments to **track** that contain depth  
631 and breadth tracking limits, and defer to lazy tracking when out of fuel.

## 4.3 Automatic contracts with clojure.spec

632 While we originally designed our tool to generate Typed Clojure annotations, it also supports  
633 generating “specs” for clojure.spec, Clojure's runtime verification system. There are key similarities  
634

638  $v ::= \dots \mid \{\overline{k \ v}\}^{\overline{\{k \ m \ \overline{\pi}\}}}$  Values

639

640

641  $\text{track}(\{\overline{k \ k' \ k'' \ v}\}, \overline{\pi}) = \{\overline{k \ k' \ k'' \ v}\}^{\overline{\{k \ t \ k'' \ t\}}}; \{\}$

642 where  $t = \{\{\overline{k \ k' \ k'' \ ?}\} \overline{\pi}\}$

643  $\text{track}(\{\overline{k \ v}\}^{\overline{\{k' \ m \ \overline{\pi'}\}}}, \overline{\pi}) = \{\overline{k \ v}\}^{\overline{\{k' \ m \ (\overline{\pi \cup \overline{\pi'}})\}}}; \{\}$

644

645  $\delta(\text{assoc}, \{\overline{k \ v}\}^{\overline{\{k' \ t', k'' \ t\}}}, k', v') = \{\overline{k \ v}\}[k' \mapsto v']^{\overline{\{k'' \ t\}}}; \{\}$

646  $\delta(\text{assoc}, \{\overline{k \ v}\}^{\overline{\{k'' \ t\}}}, k', v') = \{\overline{k \ v}\}[k' \mapsto v']^{\overline{\{k'' \ t\}}}; \{\}$

647  $\delta(\text{get}, \{\overline{k \ v, k' \ v'}\}^{\overline{\{k \ t, k'' \ t'\}}}, k) = \text{track}(v, \overline{\pi})$

648 where  $\overline{\pi} = [\pi :: [\text{key}_m(k)] \mid (m, \overline{\pi}) \in t, \pi \in \overline{\pi}]$

649  $\delta(\text{get}, \{\overline{k \ v, k' \ v'}\}^{\overline{\{k'' \ t'\}}}, k) = v$

650  $\delta(\text{dissoc}, \{\overline{k \ v, k' \ v'}\}^{\overline{\{k \ t, k'' \ t'\}}}, k) = \{\overline{k' \ v'}\}^{\overline{\{k'' \ t'\}}}; \{\}$

651  $\delta(\text{dissoc}, \{\overline{k \ v, k' \ v'}\}^{\overline{\{k'' \ t'\}}}, k) = \{\overline{k' \ v'}\}^{\overline{\{k'' \ t'\}}}; \{\}$

652

653

654

Fig. 7. Lazy tracking extensions (changes)

655

656

657

658 between Typed Clojure and clojure.spec, such as extensive support for potentially-tagged keyword

659 maps, however spec features a global registry of names via **s/def** and an explicit way to declare

660 unions of maps with a common dispatch key in **s/multi-spec**. These require differences in both

661 type and name generation.

662 The following generated specs correspond to the first **Op** case of Figure 8 (lines 2-5).

663 1 (**defmulti** op-multi-spec :op) ;dispatch on :op key

664 2 (**defmethod** op-multi-spec :binding ;match :binding

665 3 [\_] ;s/keys matches keyword maps

666 4 (**s/keys** :req-un [::op ...] ;required keys

667 5 :opt-un [::column ...])) ;optional keys

668 6 (**s/def** ::op #{:js :let ...}) ;:op key maps to keywords

669 7 (**s/def** ::column int?) ;:column key maps to ints

670 8 ; register ::Op as union dispatching on :op entry

671 9 (**s/def** ::Op (s/multi-spec op-multi-spec :op))

672 10 ; emit's first argument :ast has spec ::Op

673 11 (s/fdef emit :args (s/cat :ast ::Op) :ret nil?)

674

675

676

677

## 678 5 EVALUATION

679 We performed a quantitative evaluation of our workflow on several open source programs in three

680 experiments. We ported five programs to Typed Clojure with our workflow, and merely generated

681 types for one larger program we deemed too difficult to port, but features interesting data types.

682 Experiment 1 involves a manual inspection of the types from our automatic algorithm. We detail

683 our experience in generating types for part of an industrial-grade compiler which we ultimately

684 decided not to manually port to Typed Clojure. This was because it uses many programming idioms

685 beyond Typed Clojure's capabilities (those detailed as "Further Challenges" by Bonnaire-Sergeant

686 et al. [2016]), and so the final part of the workflow mostly involves working around its shortcomings.

```

687 1 (defalias Op ; omitted some entries and 11 cases
688 2   (U (HMap :mandatory
689 3     {:op ':binding, :info (U NameShadowMap FnScopeFnSelfNameNsMap), ...}
690 4     :optional
691 5     {:env ColumnLineContextMap, :init Op, :shadow (U nil Op), ...})
692 6   '{:op ':const, :env HMap49305, ...}
693 7   '{:op ':do, :env HMap49305, :ret Op, :statements (Vec Nothing), ...}
694 8   ...))
695 9 (defalias ColumnLineContextMap
696 10  (HMap :mandatory {:column Int, :line Int} :optional {:context ':expr}))
697 11 (defalias HMap49305 ; omitted some entries
698 12  (U nil
699 13    '{:context ':statement, :column Int, ...}
700 14    '{:context ':return, :column Int, ...}
701 15    (HMap :mandatory {:context ':expr, :column Int, ...} :optional {...})))
702 16 (ann emit [Op -> nil])
703 17 (ann emit-dot [Op -> nil])
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735

```

Fig. 8. Sample generated types for cljs.compiler.

Experiment 2 studies the kinds of the manual changes needed to port our five programs to Typed Clojure, starting from the automatically generated annotations. Experiment 3 enforces the initially generated annotations for these programs at runtime to check they are meaningfully underprecise.

## 5.1 Experiment 1: Manual inspection

For the first experiment, we manually inspect the types automatically generated by our tool. We judge our tool's ability to use recognizable names, favor compact annotations, and not overspecify types.

We take this opportunity to juxtapose some strengths and weaknesses of our tool by discussing a somewhat problematic benchmark, a namespace from the ClojureScript compiler called `cljs.compiler` (the code generation phase). We generate 448 lines of type annotations for the 1,776 line file, and present a sample of our tool's output as Figure 8. We were unable to fully complete the porting to Typed Clojure due to type system limitations, but the annotations yielded by this benchmark are interesting nonetheless.

The compiler's AST format is inferred as `Op` (lines 1-8) with 22 recursive references (like lines 5, 5, 7) and 14 cases distinguished by `:op` (like lines 3, 6, 7), 5 of which have optional entries (like lines 4-5). To improve inference time, only the code emission unit tests were exercised (299 lines containing 39 assertions) which normally take 40 seconds to run, from which we generated 448 lines of types and 517 lines of specs in 2.5 minutes on a 2011 MacBook Pro (16GB RAM, 2.4GHz i5), in part because of key optimizations discussed in Section 4.

The main function of the code generation phase is `emit`, which effectfully converts a map-based AST to JavaScript. The AST is created by functions in `cljs.analyzer`, a significantly larger 4,366 line Clojure file. Without inspecting `cljs.analyzer`, our tool annotates `emit` on line 16 with a recursive AST type `Op` (lines 1-8).

736 Similar to our opening example nodes, it uses the `:op` key to disambiguate between (16) cases,  
 737 and has recursive references (`Op`). We just present the first 4 cases. The first case `'` binding has  
 738 4 required and 8 optional entries, whose `:info` and `:env` entries refer to other **HMap** type aliases  
 739 generated by the tool.

740 An important question to address is “how accurate are these annotations?”. Unlike previous  
 741 work in this area [An et al. 2011], we do not aim for soundness guarantees in our generated types.  
 742 A significant contribution of our work is a tool that Clojure programmers can use to help learn  
 743 about and specify their programs. In that spirit, we strive to generate annotations meeting more  
 744 qualitative criteria. Each guideline by itself helps generate more useful annotations, and they  
 745 combine in interesting ways help to make up for shortcomings.

746 *Choose recognizable names.* Assigning a good name for a type increases readability by succinctly  
 747 conveying its purpose. Along those lines, a good name for the AST representation on lines 1-8 might  
 748 be `AST` or `Expr`. However, these kinds of names can be very misleading when incorrect, so instead  
 749 of guessing them, our tool takes a more consistent approach and generates *easily recognizable*  
 750 names based on the type the name points to. Then, those with a passing familiarity with the data  
 751 flowing through the program can quickly identify and rename them. For example,

- 752 • `Op` (lines 1-8) is chosen because `:op` is clearly the dispatch key (the `:op` entry is also helpfully  
 753 placed as the first entry in each case to aid discoverability),
- 754 • `ColumnLineContextMap` (lines 9-10) enumerates the keys of the map type it points to,
- 755 • `NameShadowMap` and `FnScopeFnSelfNameNsMap` (line 3) similarly, and
- 756 • `HMap49305` (lines 11-15) shows how our tool fails to give names to certain combinations of  
 757 types (we now discuss the severity of this particular situation).

758 A failure of `cljs.compiler`'s generated types was `HMap49305`. It clearly fails to be a recognizable  
 759 name. However, all is not lost: the compactness and recognizable names of other adjacent anno-  
 760 tations makes it plausible for a programmer with some knowledge of the AST representation to  
 761 recover. In particular 13/14 cases in `Op` have entries from `:env` to `HMap49305`, (like lines 6 and 7),  
 762 and the only exception (line 5) maps to `ColumnLineContextMap`. From this information the user  
 763 can decide to combine these aliases.

764 *Favor compact annotations.* Literally translating runtime observations into annotations without  
 765 compacting them leads to unmaintainable and impractical types resembling `TypeWiz`'s “verbatim”  
 766 annotation for nodes. To avoid this, we use optional keys where possible, like line 10, infer recursive  
 767 types like `Op`, and reuse type aliases in function annotations, like `emit` and `emit-dot` (lines 16, 17).

768 One remarkable success in the generated types was the automatic inference `Op` (lines 1-8) with  
 769 14 distinct cases, and other features described in Figure 8. Further investigation reveals that the  
 770 compiler actually features 36 distinct AST nodes—unsurprisingly, 39 assertions was not sufficient  
 771 test coverage to discover them all. However, because of the recognizable name and organization of  
 772 `Op`, it's clear where to add the missing nodes if no further tests are available.

773 These processes of compacting annotations often makes them more general, which leads into  
 774 our next goal.

775 *Don't overspecify types.* Poor test coverage can easily skew the results of dynamic analysis tools,  
 776 so we choose to err on the side of generalizing types where possible. Our opening example nodes  
 777 is a good example of this—our inferred type is recursive, despite nodes only being tested with a  
 778 tree of height 2. This has several benefits.

- 779 • We avoid exhausting the pool of easily recognizable names by generalizing types to commu-  
 780 nicate the general role of an argument or return position. For example, `emit-dot` (line 17) is  
 781 annotated to take `Op`, but in reality accepts only a subset of `Op`. Programmers can combine  
 782  
 783  
 784



785 the recognizability of `Op` with the suggestive name of `emit-dot` (the dot operator in Clojure  
 786 handles host interoperability) to decide whether, for instance, to split `Op` into smaller type  
 787 aliases or add type casts in the definition of `emit-dot` to please the type checker (some  
 788 libraries require more casts than others to type check, as discussed in Section 5.2).

- 789 • Generated Clojure spec annotations (an extension discussed in Section 4.3) are more likely to  
 790 accept valid input with specs enabled, even with incomplete unit tests (we enable generated  
 791 specs on several libraries in Section 5.3).
- 792 • Our approach becomes more amenable to extensions improving the running time of run-  
 793 time observation without significantly deteriorating annotation quality, like lazy tracking  
 794 (Section 4.2).

795 Several instances of overspecification are evident, such as the `:statements` entry of a `:do` AST  
 796 node being inferred as an always-empty vector (line 7). In some ways, this is useful information,  
 797 showing that test coverage for `:do` nodes could be improved. To fix the annotation, we could rerun  
 798 the tool with better tests. If no such test exists, we would have to fall back to reverse-engineering  
 799 code to identify the correct type of `:statements`, which is `(Vec Op)`.

800 Finally, 19 functions in `cljs.compiler` are annotated to take or return `Op` (like lines 16, 17). This  
 801 kind of alias reuse enables annotations to be relatively compact (only 16 type aliases are used by  
 802 the 49 functions that were exercised).

## 804 5.2 Experiment 2: Changes needed to type check

805 We used our workflow to port the following open source Clojure programs to Typed Clojure.

806 *startrek-clojure*. A reimplementaion of a Star Trek text adventure game, created as a way to  
 807 learn Clojure.

808 *math.combinatorics*. The core library for common combinatorial functions on collections, with  
 809 implementations based on Knuth's Art of Computer Programming, Volume 4.

810 *fs*. A Clojure wrapper library over common file-system operations.

811 *data.json*. A library for working with JSON.

812 *mini.occ*. A model of occurrence typing by an author of the current paper. It utilizes three  
 813 mutually recursive ad-hoc structures to represent expressions, types, and propositions.

814 In this experiment, we first generated types with our algorithm by running the tests, then  
 815 amended the program so that it type checks. Figure 9 summarizes our results. After the lines of  
 816 code we generate types for, the next two columns show how many lines of types were generated  
 817 and the lines manually changed, respectively. The latter is a git line diff between commits of the  
 818 initial generated types and the final manually amended annotations. While an objectively fair  
 819 measurement, it is not a good indication of the effort needed to port annotations (a 1 character  
 820 changes on a line is represented by 1 line addition and 1 line deletion) The rest of the table  
 821 enumerates the different kinds of changes needed and their frequency.

822 *Uncalled functions*. A function without tests receives a broad type annotation that must be  
 823 amended. For example, the *startrek-clojure* game has several exit conditions, one of which is  
 824 running out of time. Since the tests do not specifically call this function, nor play the game long  
 825 enough to invoke this condition, no useful type is inferred.

826 `(ann game-over-out-of-time AnyFunction)`

827 In this case, minimal effort is needed to amend this type signature: the appropriate type alias  
 828 already exists:

829

Library	Lines of code	Lines of Generated Global/Local Types	Lines manually added/removed	Casts/Instantiations	Polymorphic annotation	Local annotation	Type System Workaround/no-check	Overprecise argument/return type	Uncalled function (bad test coverage)	Variable-arity/keyword arg type	Add occurrence typing annotation	Erase or upcast HVec annotation	Add missing case in defalias
startrek	166	133/3	70/41	5 / 0	0	2	13/1	1 / 2	5	1 / 0	0	0	0
math.comb	923	395/147	124/120	23 / 1	11	19	2 / 9	5 / 2	0	3 / 4	1	3	0
fs	588	157/1	119/86	50 / 0	0	2	3 / 11	4 / 9	4	2 / 0	0	0	0
data.json	528	168/9	94/125	6 / 0	0	2	4 / 5	11/7	5	0 / 20	0	0	0
mini.occ	530	49/1	46/26	7 / 0	0	2	5 / 2	4 / 2	6	0 / 0	0	1	5

Fig. 9. Lines of generated annotations, git line diff for total manual changes to type check the program, and the kinds of manual changes.

```

857 (defalias CurrentKlingonsCurrentSectorEnterpriseMap
858   (HMap :mandatory
859     {:current-klingons (Vec EnergySectorMap),
860      :current-sector (Vec Int), ...}
861     :optional {:lrs-history (Vec Str)}))

```

So we amend the signature as

```

865 (ann game-over-out-of-time
866   [(Atom1 CurrentKlingonsCurrentSectorEnterpriseMap)
867    -> Boolean])

```

*Over-precision.* Function types are often too restrictive due to insufficient unit tests.

There are several instances of this in `math.combinatorics`. The `all-different?` function takes a collection and returns true only if the collection contains distinct elements. As evidenced in the generated type, the tests exercise this functions with collections of integers, atoms, keywords, and characters.

```

874 (ann all-different?
875   [(Coll (U Int (Atom1 Int) ' :a ' :b Character))
876    -> Boolean])

```

In our experience, the union is very rarely a good candidate for a Typed Clojure type signature, so a useful heuristic to improve the generated types would be to upcast such unions to a more permissive type, like `Any`. When we performed that case study, we did not yet add that heuristic to our tool, so in this case, we manually amend the signature as

```
883 (ann all-different? [(Coll Any) -> Boolean])
```

884 Another example of overprecision is the generated type of `initial-perm-numbers` a helper  
885 function taking a *frequency map*—a hash map from values to the number of times they occur—which  
886 is the shape of the return value of the core `frequencies` function.

887 The generated type shows only a frequency map where the values are integers are exercised.

```
888  
889 (ann initial-perm-numbers  
890 [(Map Int Int) -> (Coll Int)])
```

891 A more appropriate type instead takes `(Map Any Int)`. In many examples of overprecision, while  
892 the generated type might not be immediately useful to check programs, they serve as valuable  
893 starting points and also provide an interesting summary of test coverage.

894  
895 *Missing polymorphism.* We do not attempt to infer polymorphic function types, so these amend-  
896 ments are expected. However, it is useful to compare the optimal types with our generated ones.

897 For example, the `remove-nth` function in `math.combinatorics` returns a functional delete  
898 operation on its argument. Here we can see the tests only exercise this function with collections of  
899 integers.

```
900 (ann remove-nth [(Coll Int) Int -> (Vec Int)])
```

901 However, the overall shape of the function is intact, and the manually amended type only requires  
902 a few keystrokes.

```
903  
904 (ann remove-nth  
905 (All [a] [(Coll a) Int -> (Vec a)]))
```

906 Similarly, `iter-perm` could be polymorphic, but its type is generated as

```
907  
908 (ann iter-perm [(Vec Int) -> (U nil (Vec Int))])
```

909 We decided this function actually works over any number, and bounded polymorphism was  
910 more appropriate, encoding the fact that the elements of the output collection are from the input  
911 collection.

```
912  
913 (ann iter-perm  
914 (All [a]  
915 [(Vec (I a Num)) -> (U nil (Vec (I a Num)))]))
```

916  
917 *Missing argument counts.* Often, variable argument functions are given very precise types. Our  
918 algorithm does not apply any heuristics to approximate variable arguments — instead we emit  
919 types that reflect only the arities that were called during the unit tests.

920 A good example of this phenomenon is the type inferred for the `plus` helper function from  
921 `math.combinatorics`. From the generated type, we can see the tests exercise this function with 2,  
922 6, and 7 arguments.

```
923 (ann plus (IFn [Int Int Int Int Int Int Int -> Int]  
924 [Int Int Int Int Int Int -> Int]  
925 [Int Int -> Int]))
```

926 Instead, `plus` is actually variadic and works over any number of arguments. It is better annotated  
927 as the following, which is easy to guess based on both the annotated type and manually viewing  
928 the function implementation.

```
929  
930 (ann plus [Int * -> Int])
```

931

A similar issue occurs with `mult`.

```
(ann mult [Int Int -> Int]) ;; generated
(ann mult [Int * -> Int])   ;; amended
```

A similar issue is inferring keyword arguments. Clojure implements keyword arguments with normal variadic arguments. Notice the generated type for `lex-partitions-H`, which takes a fixed argument, followed by some optional integer keyword arguments.

```
(ann lex-partitions-H
  (IFn [Int -> (Coll (Coll (Vec Int)))]
    [Int ':min Int ':max Int
     -> (Coll (Coll (Coll Int)))]))
```

While the arity of the generated type is too specific, we can conceivably use the type to help us write a better one.

```
(ann lex-partitions-H
  [Int & :optional {:min Int :max Int}
  -> (Coll (Coll (Coll Int)))]))
```

*Weaknesses in Typed Clojure.* We encountered several known weaknesses in Typed Clojure's type system that we worked around. The most invasive change needed was in `startrek-clojure`, which strongly updated the global mutable configuration map on initial play. We instead initialized the map with a dummy value when it is first created.

*Missing defalias cases.* With insufficient test coverage, our tool can miss cases in a recursively defined type. In particular, `mini.occ` features three recursive types—for the representation of types `T`, propositions `P`, and expressions `E`. For `T`, three cases were missing, along with having to upcast the `:params` entry from the singleton vector `'[NameTypeMap]`. Two cases were missing from `E`. The manual changes are highlighted (P required no changes with five cases).

<pre>(defalias T   (U '{:T ':not, :type T}      '{:T ':refine, :name t/Sym, :prop P}      '{:T ':union, :types (t/Set T)}      '{:T ':false}      '{:T ':fun,       :params (t/Vec NameTypeMap),       :return T}      '{:T ':intersection, :types (Set T)}      '{:T ':num}))</pre>	<pre>(defalias E   (U '{:E ':add1}      '{:E ':n?}      '{:E ':app, :args (Vec E), :fun E}      '{:E ':false}      '{:E ':if, :else E, :test E, :then E}      '{:E ':lambda, :arg Sym, :arg-type T,       :body E}      '{:E ':var, :name Sym}))</pre>
--	--

### 5.3 Experiment 3: Specs pass unit tests

Our final experiment uses our tool to generate specs (Section 4.3) instead of types. Specs are checked at runtime, so to verify the utility of generated specs, we enable spec checking while rerunning the unit tests that were used in the process of creating them.

At first this might seem like a trivial property, but it serves as a valuable test of our inference algorithm. The aggressive merging strategies to minimize aliases and maximize recognizability,

Library	LOC	Lines of specs	Recursive	Instance	Het. Map	Passed Tests?
startrek	166	25	0	10	0	Yes
math.comb	923	601	0	320	0	Yes
fs	588	543	0	215	0	Yes
data.json	528	401	0	174	0	No (1/79 failed)
mini.occ	530	131	3	25	15	Yes

Fig. 10. Summary of the quantity and kinds of generated specs and whether they passed unit tests when enabled. The one failing test was related to pretty-printing JSON, and seems to be an artifact of our testing environment, as it still fails with all specs removed.

while unsound transformations, are based on hypotheses about Clojure idioms and how Clojure programs are constructed. If, hypothetically, we generated singleton specs for numbers like we do for keywords and did not eventually upcast them to `number?`, the specs might be too strict to pass its unit tests. Some function specs also perform generative testing based on the argument and return types provided. If we collapse a spec too much and include it in such a spec, it might feed a function invalid input.

Thankfully, we avoid such pitfalls, and so our generated specs pass their tests for the benchmarks we tried. Figure 10 shows our preliminary results. All inferred specs pass the unit tests when enforced, which tells us they are at least well formed. We had some seemingly unrelated difficulty with a test in `data.json` which we explain in the caption. Since hundreds of invariants are checked—mostly “instance” checks that a value is of a particular class or interface—we can also be more confident that the specs are useful.

## 6 RELATED WORK

*Automatic annotations.* There are two common implementation strategies for automatic annotation tools. The first strategy, “ruling-out” (for invariant detection), assumes all invariants are true and then use runtime analysis results to rule out impossible invariants. The second “building-up” strategy (for dynamic type inference) assumes nothing and uses runtime analysis results to build up invariant/type knowledge.

Examples of invariant detection tools include Daikon [Ernst et al. 2001], DIDUCE [Hangal and Lam 2002], and Carrot [Pytlik et al. 2003], and typically enhance statically typed languages with more expressive types or contracts. Examples of dynamic type inference include our tool, Rubydust [An et al. 2011], JSTrace [Saftoiu 2010], and TypeDevil [Pradel et al. 2015], and typically target untyped languages.

Both strategies have different space behavior with respect to representing the set of known invariants. The ruling-out strategy typically uses a lot of memory at the beginning, but then can free memory as it rules out invariants. For example, if `odd(x)` and `even(x)` are assumed, observing `x = 1` means we can delete and free the memory recording `even(x)`. Alternatively, the building-up strategy uses the least memory storing known invariants/types at the beginning, but increases memory usage as more the more samples are collected. For example, if we know `x : Bottom`, and we observe `x = "a"` and `x = 1` at different points in the program, we must use more memory to store the union `x : String ∪ Integer` in our set of known invariants.

*Daikon.* Daikon can reason about very expressive relationships between variables using properties like ordering ( $x < y$ ), linear relationships ( $y = ax + b$ ), and containment ( $x \in y$ ). It also

1030 supports reasoning with “derived variables” like fields ( $x.f$ ), and array accesses ( $a[i]$ ). Typed Clo-  
1031 jure’s dynamic inference can record heterogeneous data structures like vectors and hash-maps, but  
1032 otherwise cannot express relationships between variables.

1033 There are several reasons for this. The most prominent is that Daikon primarily targets Java-  
1034 like languages, so inferring simple type information would be redundant with the explicit typing  
1035 disciplines of these languages. On the other hand, the process of moving from Clojure to Typed  
1036 Clojure mostly involves writing simple type signatures without dependencies between variables.  
1037 Typed Clojure recovers relevant dependent information via occurrence typing [Tobin-Hochstadt  
1038 and Felleisen 2010], and gives the option to manually annotate necessary dependencies in function  
1039 signatures when needed.

1040  
1041 *Reverse Engineering Programs with Static Analysis.* Rigi [Müller et al. 1992] analyzes the structure  
1042 of large software systems, combining static analysis with a user-facing graphical environment  
1043 to allow users to view and manipulate the in-progress reverse engineering results. We instead  
1044 use a static type system as a feedback mechanism, which forces more aggressive compacting of  
1045 generated annotations.

1046 Lackwit [O’Callahan and Jackson 1997] uses static analysis to identify abstract data types in  
1047 C programs. Like our work, they share representations between values, except they use type  
1048 inference with representations encoded as types. Recursive representations are inferred via Felice  
1049 and Coppos’s work on type inference with recursive types [Cardone and Coppo 1991], where we  
1050 rely on our imprecise “squashing” algorithms over incomplete runtime samples.

1051 Soft Typing [Cartwright and Fagan 1991] uses static analysis to insert runtime checks into  
1052 untyped programs for invariants that cannot be proved statically. Our approach is instead to let the  
1053 user check the generated annotations with a static type system, with static type errors guiding the  
1054 user to manually add casts when needed.

1055 *Schema Inference.* Baazizi et al. [2017] infer structural properties of JSON data using a custom  
1056 JSON schema format. Their schema inference algorithm proceeds in two stages: schema inference  
1057 and schema fusion. This resembles our collection and naive type environment construction phases.  
1058 There are slight differences between schema fusion and our approach. Schema fusion upcasts  
1059 heterogeneous array types to be homogeneous, where we maintain heterogeneous vector types  
1060 until a differently-sized vector type is found in the same position. We also support function types,  
1061 which JSON lacks. While they support nested data, they do not attempt to factor out common types  
1062 as names or create recursive types like our squashing algorithms.

1063 DiScala and Abadi [2016] present a machine learning algorithm to translate denormalized and  
1064 nested data that is commonly found in NoSQL databases to traditional relational formats used  
1065 by standard RDBMS. A key component is a schema generation algorithm which arranges related  
1066 data into tables via a matching algorithm which discovers related attributes. Phases 1 and 2 of  
1067 their algorithm are similar to our local and global squashing algorithms, respectively, in that first  
1068 locally accessible information is combined, and then global information. They identify groups of  
1069 attributes that have (possibly cyclic) relationships. Where our squashing algorithms for map types  
1070 are based on (sets of) keysets—on the assumption that related entities use similar keysets—they  
1071 also join attributes based on their similar values. This enables more effective entity matching via  
1072 equivalent attributes with different names (e.g., “Email” vs “UserEmail”). Our approach instead  
1073 assumes programs are somewhat internally consistent, and instead optimizes to handle missing  
1074 samples from incomplete dynamic analysis.

1075  
1076 *Other Annotation Tools.* Static analyzers for JavaScript (TSInfer [Kristensen and Møller 2017]) and  
1077 for Python (Typypete [Hassan et al. 2018] and PyType [Google 2018]) automatically annotate code  
1078

1079 with types. PyType and Typpete inferred nodes as (? -> int) and Dict[(Sequence, object)]  
 1080 -> int, respectively—our tool infers it as [Op ->Int] by also generating a compact recursive type.  
 1081 Similarly, a class-based translation of inferred both left and right fields as Any by PyType, and as  
 1082 Leaf by Typpete—our tool uses Op, a compact recursive type containing both Leaf and Node. This  
 1083 is similar to our experience with TypeWiz in Section 1. (We were unable to install TSInfer.)

1084 NoRegrets [Mezzetti et al. 2018] uses dynamic analysis to learn how a program is used, and  
 1085 automatically runs the tests of downstream projects to improve test coverage. Their *dynamic access*  
 1086 *paths* represented as a series of *actions* are analogous to our paths of path elements.

1087

## 1088 7 CONCLUSION

1089 This paper shows how to generate recursive heterogeneous type annotations for untyped programs  
 1090 that use plain data. We use a novel algorithm to “squash” the observed structure of program  
 1091 values into named recursive types suitable for optional type systems, all without the assistance of  
 1092 record, structure, or class definitions. We test this approach on thousands of lines of Clojure code,  
 1093 optimizing generated annotations for programmer comprehensibility over soundness.

1094 In our experience, our guidelines to automatically name, group, and reuse types yield insightful  
 1095 annotations for those with some familiarity with the original programs, even if the initial annota-  
 1096 tions are imprecise, incomplete, and always require some changes to type check. Most importantly,  
 1097 many of these changes will involve simply rearranging or changing parts of existing annotations, so  
 1098 programmers are no longer left alone with the daunting task of reverse-engineering such programs  
 1099 completely from scratch.

1100

1101

## 1102 REFERENCES

- 1103 Jong-hoon (David) An, Avik Chaudhuri, Jeffrey S. Foster, and Michael Hicks. 2011. Dynamic Inference of Static Types for  
 1104 Ruby. *SIGPLAN Not.* 46, 1 (Jan. 2011), 459–472. <https://doi.org/10.1145/1925844.1926437>
- 1105 Mohamed-Amine Baazizi, Housseem Ben Lahmar, Dario Colazzo, Giorgio Ghelli, and Carlo Sartiani. 2017. Schema inference  
 1106 for massive JSON datasets. In *Extending Database Technology (EDBT)*.
- 1107 Phil Bagwell. 2001. Ideal hash trees. *Es Grands Champs* 1195 (2001), 5–2.
- 1108 Ambrose Bonnaire-Sergeant, Rowan Davies, and Sam Tobin-Hochstadt. 2016. Practical Optional Types for Clojure. In  
 1109 *European Symposium on Programming Languages and Systems*. Springer, 68–94.
- 1110 Felice Cardone and Mario Coppo. 1991. Type inference with recursive types: Syntax and semantics. *Information and*  
 1111 *Computation* 92, 1 (1991), 48–80.
- 1112 Robert Cartwright and Mike Fagan. 1991. Soft Typing. In *Proc. PLDI*.
- 1113 Michael DiScala and Daniel J Abadi. 2016. Automatic generation of normalized relational schemas from nested key-value  
 1114 data. In *Proceedings of the 2016 International Conference on Management of Data*. ACM, 295–310.
- 1115 Michael D Ernst, Jake Cockrell, William G Griswold, and David Notkin. 2001. Dynamically discovering likely program  
 1116 invariants to support program evolution. *IEEE Transactions on Software Engineering* 27, 2 (2001), 99–123.
- 1117 Robert Bruce Findler and Matthias Felleisen. 2002. Contracts for higher-order functions. In *ACM SIGPLAN Notices*, Vol. 37.  
 1118 ACM, 48–59.
- 1119 Google. 2018. *PyType*. <https://github.com/google/pytype>
- 1120 Sudheendra Hangal and Monica S Lam. 2002. Tracking down software bugs using automatic anomaly detection. In *Software*  
 1121 *Engineering, 2002. ICSE 2002. Proceedings of the 24rd International Conference on*. IEEE, 291–301.
- 1122 Mostafa Hassan, Caterina Urban, Marco Eilers, and Peter Müller. 2018. MaxSMT-Based Type Inference for Python 3. In  
 1123 *Computer Aided Verification*, Hana Chockler and Georg Weissenbacher (Eds.). Springer International Publishing, Cham,  
 1124 12–19.
- 1125 David Herman, Aaron Tomb, and Cormac Flanagan. 2010. Space-efficient Gradual Typing. *Higher Order Symbol. Comput.* 23,  
 1126 2 (June 2010), 167–189. <https://doi.org/10.1007/s10990-011-9066-z>
- 1127 Rich Hickey. 2008. The Clojure programming language. In *Proc. DLS*.
- 1128 Erik Krogh Kristensen and Anders Møller. 2017. Inference and Evolution of TypeScript Declaration Files. In *International*  
 1129 *Conference on Fundamental Approaches to Software Engineering*. Springer, 99–115.
- 1130 Gianluca Mezzetti, Anders Møller, and Martin Toldam Torp. 2018. Type Regression Testing to Detect Breaking Changes in  
 1131 Node.js Libraries. In *Proc. 32nd European Conference on Object-Oriented Programming (ECOOP)*.

1132

- 1128 Hausi A Müller, Scott R Tilley, Mehmet A Orgun, BD Corrie, and Nazim H Madhavji. 1992. A reverse engineering environment  
1129 based on spatial and visual software interconnection models. In *ACM SIGSOFT Software Engineering Notes*, Vol. 17. ACM,  
1130 88–98.
- 1131 Robert O’Callahan and Daniel Jackson. 1997. Lackwit: A program understanding tool based on type inference. In *In*  
1132 *Proceedings of the 19th International Conference on Software Engineering*. Citeseer.
- 1133 Michael Pradel, Parker Schuh, and Koushik Sen. 2015. TypeDevil: Dynamic type inconsistency analysis for JavaScript. In  
1134 *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, Vol. 1. IEEE, 314–324.
- 1135 Brock Pytlik, Manos Renieris, Shriram Krishnamurthi, and Steven P Reiss. 2003. Automated fault localization using potential  
1136 invariants. *arXiv preprint cs/0310040* (2003).
- 1137 Claudiu Saftoiu. 2010. *JSTrace: Run-time type discovery for JavaScript*. Technical Report. Technical Report CS-10-05, Brown  
1138 University.
- 1139 Uri Shaked. 2018. *TypeWiz*. <https://github.com/urish/typewiz>
- 1140 Sam Tobin-Hochstadt and Matthias Felleisen. 2010. Logical Types for Untyped Languages. In *Proc. ICFP (ICFP ’10)*.
- 1141
- 1142
- 1143
- 1144
- 1145
- 1146
- 1147
- 1148
- 1149
- 1150
- 1151
- 1152
- 1153
- 1154
- 1155
- 1156
- 1157
- 1158
- 1159
- 1160
- 1161
- 1162
- 1163
- 1164
- 1165
- 1166
- 1167
- 1168
- 1169
- 1170
- 1171
- 1172
- 1173
- 1174
- 1175
- 1176