# Supply-Chain Vulnerability Elimination via Active Learning & Regeneration

Anonymous Author(s)

# ABSTRACT

2

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

58

Software supply-chain attacks target components that are integrated into client applications. Such attacks often target widelyused components, with the attack taking place via operations such as file system and network accesses that do not affect the values that the component returns to the client and therefore preserve the client-observable behavior. We propose new active library learning and regeneration (ALR) techniques for inferring and regenerating the client-observable functionality of software components. Using increasingly sophisticated rounds of exploration, ALR generates inputs, provides these inputs to the component, and observes the resulting outputs to infer a model of the component's behavior as a program in a domain-specific language. We present HARP, an ALR system for string processing components. We apply HARP to successfully infer and regenerate string-processing components written in JavaScript and C/C++. Our results indicate that, in the majority of cases, HARP completes the regeneration in less than a minute, remains fully compatible with the original library, and delivers performance indistinguishable from the original library. We also demonstrate that HARP can eliminate vulnerabilities associated with libraries targeted in several highly visible security incidents, specifically event-stream, left-pad, and string-compare.

#### **1** INTRODUCTION

Malicious adversaries increasingly employ software *supply-chain attacks* [7, 27–29, 58]. Rather than directly targeting a victim software, these attacks target a victim's supplier, exploiting the fact that the victim software depends, directly or indirectly, on software provided by the supplier. A common scenario is that the attacker purposefully inserts vulnerabilities into open source software components that are then integrated into the eventual victim software. Modern software often integrates hundreds to thousands of small components, with many components integrated not directly, but only via transitive dependencies [27, 39, 70]. It is therefore impractical for developers to audit the code that implements the integrated components—indeed, developers can easily be completely unaware of the full range of components that their system may integrate. For these reasons, even very simple, widely used components can successfully carry vulnerabilities into client software systems.

48 For a compromised component to remain undetected, it must 49 typically deliver correct observable behavior to its client applica-50 tions. Inserted vulnerabilities are therefore typically triggered only 51 in very specific execution contexts and exhibit malicious behav-52 ior (such as stealthily exfiltrating sensitive data [5, 41], stealing 53 digital assets [42, 69], or performing covert computations on the 54 client computing platform [11, 59]) that does not interfere with 55 correct client-observable behavior. A common scenario is that the 56 client observes only the functional behavior of the component, *i.e.*, 57 the results that it returns to the client when invoked, and not any



**Fig. 1: HARP usage scenario.** A stealthy supply-chain vulnerability can be activated long after deployment. HARP can be applied before or during development (shown) to obtain a collection of safe regenerated string libraries. HARP can also be deployed at later stages (during development or even while in production, not shown) to replace potentially malicious libraries with safe regenerated versions.

malicious side effects, additional computation, or external communication that the component may perform when it executes.

Motivated by this observation, we investigate a new approach to eliminating vulnerabilities in software components. This approach takes a potentially compromised component, explores the behavior of the component in a controlled environment to learn a model of its functional behavior (this model excludes behavior characteristic of inserted vulnerabilities), then uses the model to regenerate a new version of the component. In this paper we present a system, HARP, that applies this approach to automatically regenerate vulnerabilityfree versions of widely used string libraries, including libraries that operate on collections (such as lists or streams) over strings and higher-order computations that map or fold over such collections. **Deployment Scenarios**: HARP supports a range of deployment scenarios. It can be used before application development starts to a before a polication development starts

to obtain a collection of safe regenerated string libraries that can be integrated into multiple applications developed by one or more organizations (Figure 1). It can also be deployed during development as new string libraries are integrated into the application. Finally, it can be deployed after the application is in production to replace potentially malicious libraries with safe regenerated versions.

**Scope and Limitations**: Our approach targets simple libraries that implement familiar utility computations with broad applicability across a wide range of applications. Such libraries comprise a compelling target for attackers because (1) they enable attackers to effectively target a broad range of computations and (2) they are often imported indirectly via higher-level libraries (as opposed to imported directly by the application developer), and as a result are not inspected by the application's nominal developers. Many developers may easily be unaware that their applications integrate the target library.

Our approach also targets libraries whose behavior can be accurately captured with a domain-specific language (DSL). The DSL promotes effective inference and representation of the library behavior and eliminates malicious computations as inexpressible.

Our current HARP implementation targets string libraries. Such libraries implement foundational baseline functionality used widely in modern software systems. This is especially true for dynamically typed language such as JavaScript that use runtime string

114

115

116

59

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

117 manipulation even for basic operations that in other languages are performed via type-safe alternatives such as type-constructor 118 119 pattern matching. This is also true for many web applications, in which strings and string manipulations play a prominent role. 120 Strings are therefore integrated, often indirectly, in the full range of 121 JavaScript applications and are typically treated as standard compo-123 nents within the JavaScript ecosystem. We have developed a DSL 124 that effectively captures the semantics of string computations and 125 supports the efficient representation, manipulation, and inference 126 of the underlying behavior implemented by string libraries (§4.1). Our experimental results highlight the benefits that our approach 127 can deliver for clients of such libraries (§7). 128

This focused approach comes with limitations. First, it works 129 best for widely used libraries whose computations can be captured 130 with an efficiently inferrable DSL. We anticipate that such libraries 131 will implement relatively simple, well understood computations. 132 We also anticipate that the approach will work best for functional 133 computations. Although it is possible to work with computations 134 135 that perform externally visible actions such as file system or network accesses, we anticipate that it may be more difficult to ensure 136 137 that the regenerated computations contain no malicious code.

Result Summary: HARP successfully eliminates vulnerabilities
 in 3 large-scale supply-chain attacks by learning and regenerating
 the core functionality of the vulnerable library, eliminating any
 dependency to dangerous code (§7). We are aware of no other
 system that can successfully eliminate these attacks.

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

174

Applied to 17 JavaScript string-processing libraries (§7.5), HARP learns 14 libraries within a minute and all 17 under an hour. It also aborts within 5 seconds on 11 other JavaScript libraries that fall outside the string-processing domain. HARP also successfully learns and regenerates 5 C/C++ string processing modules imported as JavaScript binary modules. The regenerated libraries execute between 2% faster and 7% slower than the original JavaScript libraries and cannot use functionality beyond basic JavaScript primitives.

Key properties of HARP's synthesis algorithm guarantee that, in the limit, our proposed learning and regeneration techniques produce candidate programs with the same client-observable behavior as the original string library, if such a candidate program exists, and without malicious behaviors that fall outside client-observable behavior.

**Contributions**: We make the following contributions:

- Active Learning: Given a component to regenerate, HARP chooses inputs, feeds these inputs to the component, and observes the resulting outputs to infer a model of the client-observable functionality that the component implements. HARP executes the component in a controlled environment to discard any behavior that is not observable in the direct functional interactions with the HARP learning system.
- Domain-Specific Language: HARP builds the inferred model 166 as a program in a DSL for capturing string computations, includ-167 ing computations over collections of strings and computations 168 that map or fold over such collections. This approach provides 169 important benefits: (1) Tractable Learning Without Overfitting: 170 The DSL acts as a strong regularizer that focuses the inference 171 on the target class of string computations. It prevents overfit-172 173 ting and promotes efficient inference that typically requires only

automatically generated input-output observations to precisely identify a specific string computation within the larger class of string computations. (2) *Safe Modeling:* The DSL is designed to express only legitimate string computations. The inferred model therefore excludes behaviors that augment string computations with auxiliary malicious computations.

- **Regeneration:** Given a string computation in the DSL, HARP regenerates the computation in the desired target programming language, with any malicious behavior in the original component not learned during inference and discarded in the regeneration.
- **Experimental Results:** It presents results that characterize the ability of HARP to learn and regenerate a range of string libraries and highlight its ability to eliminate several software supply chain attacks that target string libraries.

**Paper structure:** §2 presents background and an example that highlights the operation of HARP applied to the event-stream incident [40, 59]. §3 presents the threat model, §4 presents core ALR techniques, and §6 presents refinements that improve the efficiency of the inference and regenerated libraries. §7 presents the experimental evaluation; §8 presents related work, and §10 concludes.

Appendix B sketches the proofs of HARP's synthesis properties, and Appendix C provides additional evaluation results. An online Appendix contains anonymized accompanying material, which will be made available upon paper acceptance:

https://anonymous.4open.science/r/harp-anon-734D

#### 2 BACKGROUND & EXAMPLE

We use the event-stream incident [40, 59], where a popular streamprocessing library was modified to steal bitcoins from carefully selected targets, as an example of the attacks HARP is designed to eliminate. At the time of the incident, event-stream was used (imported either directly or indirectly) by thousands of applications and averaged about two million downloads per week. When its author handed off maintenance to a volunteer—common practice in open-source development projects—the new maintainer added an obfuscated, malicious library called flatmap-stream as a dependency to event-stream.

The malicious flatmap-stream library is designed to harvest account details from select Bitcoin wallets. If run in the dependency tree of a specific Bitcoin application called Copay, flatmap-stream loads Copay's account module containing the Bitcoin wallet credentials of the user using Copay. It then overwrites the account's getKeys method with one that copies and stores the credentials on the side. It then loads the http module, and posts the credentials to a remote server, before returning the results to the caller method.

We note that flatmap-stream also maps a function over stream elements. Because this behavior is desired client-observable behavior, simply removing flatmap-stream breaks the client application. The attack succeeds by performing *effects*—loading account, overwriting getKeys, importing http, and calling post—that do not interfere with the client-observable behavior. The attack is not detectable by static analysis, because the attacker employed a series of dynamic encryption passes, nor dynamic analysis because the

346

347

348

Supply-Chain Vulnerability Elimination via Active Learning & Regeneration

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

279

280

282

283

284

285

286

287

288

289

290



Fig. 2: Overview. In an isolated container environment, HARP loads a library and inspects its interface. Using increasingly sophisticated rounds of exploration, it generates inputs, provides these inputs to the library, and observes the resulting outputs to infer a model of the library's behavior as a program in a domain-specific language.

malicious code activates selectively far from development and testing: only when event-stream was part of Copay's dependency tree, only when run on the "live" bitcoin network, and only on users that had a balance of 100 bitcoin or more [51]. When run in any other context, the compromised version of flatmap-stream exhibits identical behavior as the correct version.

**Applying HARP**: HARP can directly target a specific dependency or a library that integrates multiple dependencies. The following line applies ALR directly to flatmap-stream:

harp -ft js flatmap-stream

HARP first loads flatmap-stream in an isolated container environment and applies lightweight program transformations to instrument its execution (Fig. 2's (1)). This instrumentation records operations such as library imports, file-system reads, and global variable accesses that flatmap-stream performs. HARP also extracts information about the library interface (2). This information includes the number of returned methods and fields and the number of arguments for each method. HARP then runs flatmap-stream on synthesized inputs, to extract information about the types of each argument.

HARP next uses flatmap-stream to synthesize a program in 269 270 the HARP DSL as follows. It iteratively generates candidate pro-271 grams in the HARP DSL, filtering out candidate programs that do not match the extracted type information (3). It then executes the original version of flatmap-stream and remaining candidate 273 programs on synthesized inputs (④). It observes the parameter and 274 return values of the original library and the candidate DSL pro-275 grams (these parameter and return values are the client-observable 276 277 behavior). It filters out candidate programs that exhibit different 278 client-observable behavior than the original library ((5)).

In the limit, this process is guaranteed to produce a candidate program with the same client-observable behavior as the original 281 library (6), if such a candidate program exists (§5). In practice, HARP is usually able to synthesize a unique successful candidate program within an hour and typically within minutes (§7). HARP also implements a --quick-abort option that immediately aborts the search if the HARP instrumentation detects any non-clientobservable behavior such as file system, environment variable, or network access.

In our example, the malicious flatmap-stream behavior is not triggered in our isolated container environment and flatmap-stream Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD

exhibits fully correct behavior. Working with 2,536 inputs, HARP takes 1.4 minutes to synthesize the following correct DSL program, which exhibits identical behavior as the correct version of flatmap-stream:

 $f s = map (squash n) | "{(c)}"$ 

Here f maps the function squash n over the elements of s, thereby flattening s, and then pipes each of the results to an output pattern, which simply outputs its input element.

HARP then compiles the synthesized DSL program to the following JavaScript library:

```
const libHarp = require('./lib-harp.js');
```

```
let program = (f, isAsync) => {
 const stream = new libHarp.Stream();
  stream.addOperation(libHarp.squash);
 stream.addUserOperation(f, isAsync);
  return stream:
};
```

module.exports = program;

The compiled regenerated library is a direct translation of the inferred HARP DSL program. It links to lib-harp, a module that supports HARP's core functionality (part of the TCB, §3).

#### THREAT MODEL 3

HARP protects against an adversary that fully controls a target component and can modify it in any way that does not affect the client-observable functionality of the component. By preserving the client-observable functionality, the adversary aims to execute undetected attacks when the component is integrated into an application. Examples of modifications include added functionality that reads from the file system, sends messages over the network, reads environment variables, or writes to global variables.

For ALR to regenerate a successful replacement, the library must exhibit the desired behavior during testing and this desired behavior must conform to the ALR DSL. We anticipate that our target class of software supply-chain vulnerabilities will typically satisfy these two requirements-their goal is typically to provide the client with the desired functionality while either (1) stealthily opening up a vulnerability that can be remotely exploited by carefully crafted inputs, or (2) silently exfiltrating data or modifying the system on which it runs. To avoid exploitation during learning, ALR runs the target library in a controlled isolated environment.

An attacker may also simply remove the library from the ecosystem, disabling any application that depends on the library. By replacing the library with a regenerated local version before the original library is removed, ALR eliminates the dependence and enables applications to continue to operate successfully even in the absence of the original library.

The language's runtime environment, bindings for locating and loading libraries, a small compiler offered by HARP and the associated lib-harp.js runtime-support library are all part of the trusted computing base (TCB). To capture possible interactions between libraries, HARP is assumed to be loaded before other libraries. It is also assumed that other libraries do not cooperate with the target library to attack the system.

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

349	Module	т	:=	s   m . s
350	Statement	s	:=	add $c l \mid del l \mid at l b \mid repeat n s$
351			1	toggle $p o   \max \pi  $ fold $b  $ split $\sigma s \sigma$
352	Location	l	:=	$i \mid /p /$
353	Index	i	:=	$n \mid \text{start} \mid \text{end}$
354	Predicate	p	:=	not $p \mid \text{or } p p \mid \text{and } p p \mid k$
355	Pipeline	π	:=	$\pi \cdot \pi \mid b \mid \text{squash } n \mid o$
356	Char. Class	с	:=	$\alpha \mid n \mid \sigma \mid \operatorname{regex} r \mid \operatorname{capt} k$
357	Output	0	:=	$o \cdot o \mid c \mid \perp \mid b \perp$
358	Capture	k	:=	$k \cdot k \mid p \mid p \perp$
359	Built-ins	b	:=	+   -   ×   %   Math.*  String.*
360	Regex	r	:=	$r \mid \alpha \mid n \mid \sigma \mid^* \mid^+$
361	Special	$\sigma$	∈	$\Sigma_1$
362	Alpha	α	∈	$\Sigma_2$
363	Num	п	∈	$\mathbb{N}$
364				

Fig. 3: HARP's library DSL. The domain-specific language (DSL) captures the space of inferable program libraries.

#### **ACTIVE LEARNING & REGENERATION** 4

HARP combines three components: a DSL for specifying stringprocessing computations (§4.1), an algorithm for inferring computations in the DSL (§4.2), and an input generation component that produces the inputs for the inference algorithm (§4.3). The three components work in tandem, aided by lightweight runtime interposition for mapping the interface of a library (§4.4).

# 4.1 Domain-specific Language

Fig. 3 presents the HARP DSL. The DSL specifies the set of all programs that can be constructed by HARP, broken down into a few broad classes: (1) computational primitives, which apply transformations on their input, (2) built-in primitives for number and string manipulation commonly offered by high-level languages, (3) input ranges, over which these primitives are applied, and (4) character classes, used for pattern ranges and primitives. More complex classes often combine less complex ones.

Computational primitives: Computational primitives are either 386 statements or pipelines. Statements include add and del primitives 387 for introducing and deleting characters and higher-order map and 388 fold primitives for applying a first-order primitive over a range. 389 Pipelines apply a series of operators to a collection of elements in an 390 input stream-optionally recursively to elements of their elements. 391

These primitives are HARP's primary building blocks; their oper-392 393 ational semantics are presented in Fig. 4. The transition function 394  $\implies$  maps a computational primitive within our DSL to its output value. For example, the primitive add accepts a character *c*, a lo-395 cation l, and a string s, and returns a string that is the result of 396 adding the character c at location l in s. Strings are encoded as lists 397 of characters, list concatenation is encoded as -, and operations 398 encoded in sans - serif are built into HARP-for example, match 399 accepts a predicate *p* and a string *s* and returns three character lists: 400 (i) a string  $s_1$  up to (but not including) the match, (ii) the matching 401 string  $s_2$ , (iii) the rest of the string  $s_3$  following  $s_2$  in s. 402

We note that this is only a small set of key operators, which 403 404 are augmented by built-in operators, input ranges, and character 405 classes. HARP's DSL contains a significant number of operators, 406

which allow it to capture a large class of functionality required to implement string-processing functions.

Built-in primitives: This class contains primitives offered commonly by the standard libraries of different high-level programming languages, including operations for arithmetic-e.g., log, sqrt, etc.and string manipulation-e.g., toUpper, toAsciiCode. The class of Built-ins re-implements these operators from scratch to address two challenges. The first challenge is that different languages offer different operators and under different names; the HARP DSL unifies a common subset under a common set of identifiers. The second challenge is that the invocation patterns of such primitives are different for different languages—for example, JavaScript's n.toString is invoked directly on a number n, whereas Python str(n) takes ndirectly as an argument. HARP DSL introduces these operators as functions whose first argument is the input string.

Input ranges: Computational primitives often take as argument a location within the string. In their simplest form, locations are indices relative to the start of the input segment, which can be a string or a substring within that. For example, the index start in the expression (at start String.toUpper) matches the beginning of the string.

Locations can also be predicates that pattern-match on the form of the string. Predicates are formed by the composition of a simpler set of base predicates. Composition operators include negation, disjunction, and conjunction. Base predicates are centered around a simple pattern-matching language that includes characters, numbers, "\*" (Kleene-star superscript), and "+" (Kleene-plus superscript). For example, the predicate  $/a^+/$  in at  $(/a^+/)$  (String.toUpper) matches one or more a characters.

Character Classes: The DSL includes three sets of characters. Two of these sets come pre-configured and built into the DSL: (1) the set of integer numbers, and (2) the set of alphanumericsincluding number characters "0" to "9", lowercase letters "a" to "z", uppercase letters "A" to "Z", and punctuation symbols. The third set contains characters that are special to a particular computation. The members of this set are discovered during the learning phase via input generation (§4.3).

Capture and output expressions: Two examples of how simple elements like character classes and built-in primitives are used to construct more powerful primitives are capture and output expressions. toggle's second argument is an output expression, which can be thought of as a format string that one would pass into a function like C's printf, describing the formatting of the function's output. It can contain literal characters, as well as special identifiers, which are bound to strings that were matched as part of toggle's first argument, its predicate, and captured using a capture expression. For instance, whenever toggle encounters any character preceding an uppercase character in the program:

it will output the first character it matched-which it assigned to variable a in the capture expression-followed by a dash, followed by the captured uppercase character (assigned to b) converted into lowercase.

$$\frac{l = \text{index } 0}{\text{add } l c s \Longrightarrow [c] \cdot s} \text{ Add}_{1} \qquad \frac{l = \text{index } i l' = \text{index } i - 1}{\text{add } l c s \Longrightarrow [s_{0}] \cdot s_{1-n}' = \text{add } l' c s_{1-n}} \text{ Add}_{1} \qquad \frac{l = /p/(s_{1}, s_{2}, s_{3}) = \text{match } p s}{\text{add } l c s_{3}} \text{ Add}_{1} \qquad \frac{s_{1} = s_{1} + s_{1} + s_{1} + s_{1}}{\text{add } l c s \Longrightarrow [s_{0}] \cdot s_{1}'} \text{ Add}_{1} \qquad \frac{s_{1}' = s_{1} + s_{1} + s_{1} + s_{1}}{\text{add } l c s \Longrightarrow [s_{0}] \cdot s_{1}'} \text{ Add}_{1} \qquad \frac{s_{1}' = s_{1} + s_{1} + s_{1} + s_{1}}{\text{add } l c s \Longrightarrow [s_{0}] \cdot s_{1}'} \text{ Add}_{1} \qquad \frac{s_{1}' = s_{1} + s_{1} +$$

$$\frac{l = \operatorname{index} 0}{\operatorname{del} l \ c \ s \Longrightarrow s_{1-n}} \operatorname{DeL}_{1-n} \qquad \frac{l = \operatorname{index} i \ l' = \operatorname{index} i - 1}{\operatorname{del} l \ s \Longrightarrow s_{1-n} \operatorname{DeL}_{1-n}} \operatorname{DeL}_{1-n} \qquad \frac{l = \operatorname{index} i - 1}{\operatorname{del} l \ s \Longrightarrow s_{1-n} \operatorname{del} l' \ s_{1-n}} \operatorname{DeL}_{1-n} \qquad \frac{l = /p/(s_1, s_2, s_3) = \operatorname{match} p \ s}{\operatorname{del} l \ s \Longrightarrow s_1 \cdot s'_3} \operatorname{DeL}_{1-n}$$

$$\frac{(s_1, s_2, s_3) = \text{match } p \ s}{\operatorname{at} l \ c \ s_3} = \operatorname{at} l \ c \ s_3}$$

$$\frac{n \neq 0 \ s' = f \ s}{\operatorname{repeat} n \ f \ s \Longrightarrow \operatorname{repeat} (n-1) \ f \ s'} \quad \underbrace{n = 0 \ s' = f \ s}_{\operatorname{repeat} n \ f \ s \Longrightarrow s'} \quad \operatorname{Repeat}_{\mathrm{Repeat}_{\mathrm{N}}}$$

$$\frac{n = 0 \ s' = f \ s}{\operatorname{repeat} n \ f \ s \Longrightarrow s'} \quad \operatorname{Repeat}_{\mathrm{Repeat}_{\mathrm{N}}}$$

$$\frac{(s_1, s_2, s_3) = \operatorname{match} p \ s}{\operatorname{toggle} l \ f \ s \Longrightarrow s_1 \cdot (f s_2) \cdot s_3} \operatorname{Toggle} \frac{s = []}{\operatorname{map} f \ s \Longrightarrow []} \operatorname{Map}_E \qquad \frac{s = [s_0] \cdot s_{1-n} \ s'_0 = f \ s_0}{\operatorname{map} f \ s \Longrightarrow [s'_0] \cdot s'_{1-n}} \operatorname{Map}_F \qquad \frac{s = []}{\operatorname{fold} f \ r \ s \Longrightarrow r} \operatorname{Fold}_E$$

$$\frac{(s_1, [], []) = \operatorname{match} / c / s s'_1 = f s_1}{\operatorname{split} c f c' s \Longrightarrow s'_1} \quad SPLIT_S \qquad \frac{s = [s_0] \cdot s_{1-n} s'_{1-n} = \operatorname{fold} f r s_{1-n} s'_0 = f s_0 s'_{1-n}}{\operatorname{fold} f r s \Longrightarrow s'_0} \quad FOLD_F$$

$$\frac{s = []}{\text{plit } c \ f \ c' \ s \Longrightarrow []} \quad SPLIT_E \quad \frac{(s_1, s_2, s_3) = \text{match } / c/ \ s \ s'_1 = f \ s_1 \ s'_3 = \text{split } c \ f \ c' \ s_3}{\text{split } c \ f \ c' \ s \Longrightarrow s'_1 \cdot [c'] \cdot s'_3} \quad SPLIT_M$$

Fig. 4: DSL Semantics. A subset of HARP's DSL semantics, describing HARP's computational primitives.

#### 4.2 Synthesis Algorithm

Alg. 1 outlines HARP's library synthesis algorithm, which synthesizes a new library L' for a library L.

S

**Initial configuration**: For each function *f* in *L*, Alg. 1 synthesizes a function f' by exploring the space of programs expressible in HARP's DSL. It starts by invoking procedure generateInputs, which generates a set of inputs as described in the next section (§4.3) and stores them in *I*. The algorithm then invokes getGroundTruth, which runs the original function f on the set of inputs I to obtain a set of outputs O. These outputs are considered ground-truth outputs, because they are generated by the reference implementation. For example, applying getGroundTruth to f = length and I = ["a", "bb", "ccc"] returns O = [1, 2, 3].

Alg. 1 next invokes *soundTypeConstraints* to collect a set of sound type information T for the values in O. This procedure includes several type-inference tests for checking whether the values in O represent numbers, whether their length is significantly longer or shorter than the inputs, and whether they contain any special characters. For example, the result of calling sound Type Constraints on a String.length function would return String  $\rightarrow$  Number.

Navigating the search space: The algorithm then prepares the search space of candidate regenerated programs, which is parametric over the maximum number *n* of terms used in the program-*i.e.*, the size of the abstract syntax tree (AST) of each candidate regenerated program. Specifically, this space is explored in repeated rounds of increasing complexity and size of the synthesized program.

For each size *n*, Alg. 1 first invokes procedure *allPrograms* to obtain all possible programs whose AST size is not greater than *n* and whose type satisfies the constraints in *T*. The *allPrograms* procedure takes a number *n* and a set of sound type constraints T and returns a set  $P_n$  containing all of the programs of size n

that satisfy these type constraints. Consider an example where (1) all programs of AST size 1 are captured by the set of single-term programs {count, toString, +, -, \*}, and (2) the type constraints include Number  $\rightarrow$  Number. Then  $P_1 = \{+, -, *\}$ .

The algorithm then invokes *pruneSpace* to eliminate candidate programs in  $P_n$  whose input-output behavior does not conform to (I, O). This procedure eliminates candidate programs with behavior that is not identical to f-i.e., programs for which not all inputs in I produce outputs in O. At times, pruneSpace generates more input-output examples to further differentiate between candidates and thus prune the search space even further. It finally returns a set of candidate programs P, all of which implement f's input-output behavior on the input-output examples.

Finally, Alg. 1 inspects the set *P*. If *P* is not empty, it ranks the candidate programs in P by invoking getOpt, which returns the highest-performance program (see below).

Other information: The synthesis algorithm maintains some additional information on the side (not shown in Alg. 1). First, the synthesis algorithm is configured to run up to a time limit-either a limit  $t_{\bar{f}}$  per function f in L or a limit  $t_L$  for L overall. If only  $t_{\bar{f}}$ has been specified, then  $t_L$  is calculated as  $t_{\bar{f}} \times |f_{1-n}|$  spread fairly across all functions  $f_{1-n}$  in L; when HARP timeouts for one of the methods, it simply outputs Nil and moves to the next method in L. When  $t_L$  is specified, HARP can allocate this time as it sees fit (see parallelism in §6.2). The combination of the two limits is possible too, instructing HARP to spend no more than  $t_L$  minutes overall, with no more than  $t_{\bar{f}}$  minutes per function f in L.

Using timeouts, Alg. 1 may need to exit the inner loop with a P equal to the empty set. If this happens, L'. f is assigned Nil which is important for partial regeneration, in cases where only a fraction of a library's functionality have been successfully regenerated.

Data: Original Library L
<b>Result:</b> Regenerated Library L'
$L' \leftarrow \emptyset$
foreach $f \in L$ do
$I \leftarrow generateInputs()$
$O \leftarrow getGroundTruth(f, I)$
$T \leftarrow soundTypeConstraints(f, O)$
$n \leftarrow 0$
while true do
$n \leftarrow n+1$
$P_n \leftarrow allPrograms(n, T)$
$P \leftarrow pruneSpace(P_n, I, O)$
if $P \neq \emptyset$ then
$L'.f \leftarrow getOpt(P)$
break
end
end
end

return L'

**Algorithm 1: HARP's synthesis algorithm.** Given as input a blackbox library L, it attempts to synthesize a new library L'.

During the *pruneSpace* method, the synthesis algorithm collects information about the runtime performance of the regenerated libraries. Some of the inputs in this phase are large, to make any differences in overhead more pronounced. This information is then used by getOpt to rank candidates based on their runtime performance, returning the regenerated library with the best performance.

#### 4.3 Input Generation

HARP generates inputs for each original function f in L and executes f to obtain the input-output pairs. HARP chooses these input values to gather a variety of output values that, combined, highlight key properties of f's behavior. As HARP does not know beforehand what input streams are the most appropriate for inferring the behavior of a black-box f, it adopts an active learning algorithm to generate f's inputs. There are two kinds of inputs HARP is interested in: (1) primary inputs, which are the strings on which the string-processing computation is applied and (2) secondary inputs, which are other parameters of f affecting the specifics of the string computation. All mutations described below are applied concurrently in iterative rounds providing information or eliminating candidates. When a mutation iteration results in no candidate eliminations, this phase of input generation terminates and saves the set of candidate regenerations.

**Primary inputs**: The primary input of a string processing function is a string—a collection of characters—or a collection of strings. HARP generates primary inputs of various shapes in an attempt to understand which of their characteristics affect f's output. Characteristics of the input shape include high-level properties such as input length, homogeneity, or sorting, and low-level properties such character capitalization or the existence of specific input characters.

For high-level properties, HARP starts with small inputs and gradually mutates them to get longer inputs that satisfy certain properties. During this iterative mutation process, HARP filters out competing candidate functions. Example mutations include increasing the number of a certain set of characters or introducing some sorting discipline. For each mutation, a mutated input-shape specification is used to generate a set of inputs, run it through the original function f, and then obtain and study its outputs. After trying all available mutation strategies, HARP checks if any candidates are eliminated in the current iteration; if so, HARP attempts to identify the mutations that are the most effective at eliminating candidates. It then applies these mutations to the input-shape specification and enters the next round of iteration.

**Special characters**: Input characters are particularly important because they may affect f's processing locations—thus HARP attempts to quickly discover a set of special characters  $\Sigma_1$ . The key insight behind such discovery is that that string computations are generally applied over linear data structures that encode control and data characters in a single data stream. Consider the following string:

#### one:two three-four

Different processing primitives may be affected by different characters. For example, a (to-upper) function converting to upper-case operates on the entire string, a (split :) function splitting on ":" will match only the corresponding character, and a "mask \*" function replacing characters with "\*" will only match a subset of characters. These and other examples are shown below:

three-four	upper-case
three-four	split ':'
three-four	mask-cc
three-four	camelCase
three-four	array-first
	three-four three-four three-four three-four three-four

To discover this set, HARP generates strings with a combination of letters, numbers, and punctuation symbols. As soon as some of these inputs start affecting the results, HARP narrows down the set of symbols by mutating only parts of the input string.

**Secondary inputs**: Functions in the original library *L* rarely accept only strings as their inputs. That is, while the processing targets the primary input string, other arguments part of the method's interface need to be provided. For example, a simple count(s, c) method that counts all occurrences of c in s takes two arguments. To understand the effect of other inputs to the computation, HARP introduces small DSL describing possible secondary values (Fig. 5). To maintain acceptable performance, HARP generates only constrained inputs of these types—both in terms of size and complexity.

These values can be summarized into two broad classes. The first class is composite values such as lists, objects (maps), and functions. The DSL includes only two functions, helpful for cases when f is a higher-order function. These two functions are designed to have types that are permissive and will likely not throw exceptions. The first function simply returns its first argument, matching any fold-like operations; the second function returns its first argument as a string, covering additional use cases where the first-order function is expected to return strings—highly likely due to the domain of HARP. Both functions take a variable number of arguments so as to be compatible with any invocation in the black-box f.

The second class involves primitive values such as strings, numbers, booleans. The value  $\perp$  corresponds to null or undefined values; such values are important for understanding the default parameters or behavior of a computation.

Supply-Chain Vulnerability Elimination via Active Learning & Regeneration

Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD

Value	υ	:=	$p   \{s : v,\}   [v,]$
		Ι	$\lambda(x,\ldots).x \mid \lambda(x,\ldots).\operatorname{str}(x)$
Primitive	p	:=	$s \mid n \mid b \mid \perp$
Boolean	b	:=	true false
String	s	$\in$	Σ
Number	n	∈	M

Fig. 5: HARP's secondary-input DSL. This language captures the space of possible inputs to secondary arguments.

# 4.4 Mapping Library Structure

This section covers a few details on how HARP regenerates constant fields and how it discovers the structure of a library.

**Constant fields**: The majority of string-related functionality is expected to be exposed as functions. At times, however, L may contain fields other than functions—*e.g.*, a map of country names to dial-in prefix codes. In these cases, HARP can copy the structure to L' using runtime meta-programming facilities: it traverses L's return object to identify and copy such values directly.

In rare cases, these inputs are hidden behind a functional interface that does not allow meta-programming facilities to permeate through. In these cases, HARP resorts again to active learning—but its input generation leverages a built-in dictionary of common English words. HARP attempts these words under various combinations and capitalizations to gain more information about the mapping.

Field discovery: To apply the techniques described earlier, HARP needs to know how to interact with L and how to feed it inputs. To answer this, HARP first loads the original library, an operation that returns an object that contains the values exported by the library. These values may include functions or other directly accessible fields. The way HARP interacts with these fields depends on whether the functionality about to be regenerated has been explicitly named by the developer using HARP. If it has been named, HARP indexes only the named functions from the returned object. If there is no explicit naming involved, HARP uses runtime meta-programming to traverse the returned object in order to understand and regenerate the structure of the library. In the former case, the set L in Alg. 1 contains only developer-specified names; in the latter case, the set contains all names.

#### **5 GUARANTEES**

A key correctness guarantee is that the HARP synthesis algorithm (Alg. 1) will only produce string computations whose behavior is captured by the DSL in Figure 3. Recall that Algorithm 1 maintains a current program search size n, set of input-output examples I, O obtained from executions of the original library L, and set of programs P in the HARP DSL. The HARP synthesis algorithm provides the following key correctness guarantees:

- All programs in *P* exhibit identical behavior as the original library *L* on the list of generated inputs *I* (the call to *pruneSpace* in Algorithm 1 filters out all DSL programs whose behavior differs).
- The set of DSL programs *P* contains all DSL programs of size *n* or less that exhibit identical behavior as the original library *L* on the list of generated inputs *I*.
- These guarantees have an immediate corollary:

• If the original library *L* has the same behavior on all inputs as some DSL program f' and f is of a given size *n* or less, then  $f \in P$ . Moreover, if  $P = \{f\}$  (i.e., *f* is the only program in *P*), then the newly synthesized library *L'* has identical behavior as the original library *L* on all inputs.

In the limit the algorithm will generate all inputs and consider all programs in the DSL. More precisely, for any specific input and program of some size n, there is some finite execution of the algorithm will generate that input and consider that program. This fact ensures the following guarantees:

- If the original library *L* has the same behavior as some DSL program *f* (of some size *n*), then at some finite point in the execution of Algorithm 1,  $f \in P$  for all future execution points.
- If the original library *L* has different behavior than some DSL program *f* (of some size *n*), then at some finite point in the execution of Algorithm 1, *f* ∉ *P* for all future execution points.

These guarantees provide a form of correctness in the limit—as the algorithm runs, it (1) will eventually (in finite time) find the correct DSL implementation of the original library L if such a correct program exists in the DSL, and (2) will eventually (in finite time) filter out any DSL program whose behavior does not match the original library L on all inputs.

# **REFINEMENTS**

We next present several HARP refinements.

#### 6.1 Isolated Learning

To avoid exploitation during ALR, HARP interacts with target libraries in an isolated container environment. HARP first launches a Docker container and imports the library in the context of an TCP server. HARP then traverses the object returned by the import statement to create a remote-procedure-call (RPC) shim, which it then writes in the host file-system.

HARP'S ALR scaffolding infrastructure on the host environment loads the shim module to interact with the target library. For every invoked library function, the RPC shim serializes the arguments and send them to the server executing in the Docker container. HARP invokes the corresponding function and returns tehe results back to the shim, which delivers them to HARP running on the host environment. The channel between the RPC client function and the corresponding function running in the container is encrypted using NaCl authenticated encryption primitives [4].

# 6.2 Synthesis Acceleration

**Type Guidance**: HARP leverages sound type information to guide its choice of DSL terms. This is achieved through a few different means, starting by checking the size and type of the output. If the output is significantly smaller, then a fold-like reduction is likely to play a prominent role in the regenerated computation. Additionally, if the output has a certain type—such as a number or a boolean value—then that type should featured in the first-order function used as part of the reduction. Outputs whose size is close to that of the input string often correspond to add or at constructs.

The study of more complex outputs is also possible, as HARP can leverage meta-programming available by the source language

Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD





to introspect the value returned by *L*. This is different from other domains where active learning is applied through serialization-deserialization interfaces that encode all values as strings, and thus obscure the true types of the values returned by a program fragment. These refinements can prune the synthesis search space significantly.

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

870

**Term Weights**: Different (classes of) terms from HARP's DSL have different likelihoods of appearing in learned DSL programs. For example, many regenerated string-processing libraries add or delete characters. HARP uses such likelihood information to guide synthesis, by generating higher-likelihood terms in the DSL with higher probability HARP explores the space of candidate programs.

Term weights depend significantly on the types of the inputs and outputs. For example, if the output is a number then reduction statements such as fold and built-ins such as  $\times$  and + are more likely to appear in the regenerated program.

**Parallel Synthesis:** HARP's synthesis features ample opportunities for parallelization. One opportunity occurs in candidate generation, in which different worker processes explore disjoint subsets of the candidate space. Another opportunity occurs in input generation and testing—*i.e.*, calling the same synthesized candidate on multiple inputs.

As scaling out involves constant overheads for process spawning and interprocess communication, scaling out makes sense only after constant costs are negligible relative to synthesis. This is achieved by having HARP scale out after a few AST levels have been explored.

#### 6.3 Multi- & Part-Library Regeneration

**Multi-library Regeneration**: A library *L* is often implemented in terms of other libraries  $L_{1-n}$ . The  $L_{1-n}$  are typically smaller and simpler than *L*—often significantly so—and often encodes straightforward processing patterns. Common composition patterns include (1) selection, where different functions (or arguments to these functions) in *L* are served by different  $L_{1-n}$ , (2) pipelining, where different processing stages in *L* come from different libraries in  $L_{1-n}$ , and (3) enhancement, where functions in *L* are implemented using functionality from  $L_{1-n}$ . Thus, by targeting  $L_{1-n}$  HARP can apply ALR techniques more efficiently than it would in the full *L*.

To achieve this, HARP leverages its field-discovery (§4.4) and side-effect detection (§6.4) facilities, coupled with additional interposition on import statements themselves. Combined, they allow HARP to (1) detect cases where a library imports other libraries, and (2) apply ALR on  $L_{1-n}$ . HARP applies ALR on  $L_{1-n}$  both individually in isolation, to extract key properties about their behavior, as well as to *L*, to extract information about the interaction between *L* and  $L_{1-n}$ .

**Partial Regeneration**: HARP may only partially regenerate L, if (1) a subset of library functions in L fail regeneration, *e.g.*, due to side-effects, or (2) if some developer tests—HARP's very last stage—fail. Partial regeneration can still be useful to developers in a variety of ways. For example, the regenerated library can operate side-by-side with a hardened version of the original library.

The latter fast-slow setup combines improved security properties with acceptable overall performance. The partially regenerated L' serves the majority of the calls, and it does so efficiently and securely. At times, however, L' receives input that falls outside its expected range of operation—but not out-



side that of *L*. These inputs result into a runtime exception, caught by a HARP controller component, which then forwards the input to *L*. As *L* now executes with additional hardening in place, it is significantly less efficient, but still computes the correct output securely. The exact hardening mechanism and thus its performance overhead can vary significantly [1, 26, 32, 38, 63], and depends directly on details related to the threat model—for example native memory-unsafe binaries require additional care.

### 6.4 Quick Aborts

HARP implements a --quick-abort option that immediately aborts the search if the HARP instrumentation detects behaviors such as file system, environment variable, or network access that are not observable to clients that work only with values returned from the target library. Such behavior signals that the original library may be falling outside HARP's model of computation, allowing HARP to quickly abort the ALR process.



of the library—*i.e.*, ones that are not bound to values *in* the library. HARP starts from a few well-known root names—a static list of

928

Anon

Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD

names provided by default by the language and runtime environment. For example, in server-side JavaScript these names include
the global variable table, the require function for importing other
libraries, and the process object for providing access to environment variables, process arguments, and other information in the
broader environment.

Load-Time Transformations: Modern dynamic languages fea ture a module-import mechanism that loads code at runtime as a
 string. HARP applies lightweight load-time code transformations on
 the string representation of each module, as well as the context to
 which it is about to be bound, to insert instrumentation wrappers
 into the module before it is loaded.

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

HARP's transformations first create a modified copy of a module's runtime context. The context is a name-value mapping for all free name variables available to the module by default. The modifications target the values in this mapping—traversing and wrapping each value with an interposition mechanism that records the access in a global access table. HARP then binds the modified context to the module, using a source-to-source transformation that re-defines names in the context as library-local ones and assigns to them the values of the modified context.

HARP's transformations have a common structure that traverses objects recursively—a base transform wrap, which we review first (and whose effects are shown in Fig. 6a). The wrap transform takes an object O and returns a new object O', where every field f of Ois wrapped with and replaced by a method f'. If called, f' adds a record to a global map noting that this particular field f has been accessed and then passes arguments to f.

**Context Creation**: To prepare a new context to be bound to a library being loaded, HARP first creates an auxiliary hash table (Fig. 6*b*), mapping names to newly transformed values: names correspond to implicit modules—globals, language built-ins, module-locals, *etc.*; transformed values are created by wrapping individual values in the context to insert instrumentation hooks.

User-defined global variables are stored in a well-known location (*e.g.*, a map accessible through a global variable named global). However, traversing the global scope for built-in objects is generally not possible. To solve this problem, HARP collects such values by resolving well-known names hard-coded in a list. Using this list, HARP creates a list of pointers to unmodified values upon startup.

Care must be taken with module-local names such as the module's absolute filename, its exported values, and whether the module is invoked as the application's main module. These names refer to a different value for each module, and thus attempting to access the values directly from within HARP's transformation scope will fail subtly: the names will end up resolving to module-local values of HARP *itself*. HARP solves this issue deferring these transformations for the context-binding phase (discussed next).

**Context Binding**: To bind the code whose context is being transformed with the freshly created context, HARP applies a sourceto-source transformation that wraps the module with a function closure (Fig. 6c.). By enclosing and evaluating a closure, HARP leverages lexical scoping to inject a non-bypassable step in the variable name resolution mechanism.

The closure starts by redefining default-available non-local names as module-local ones, pointing to transformed values that exist in the newly-created context. It accepts as an argument the customized context and assigns its entries to their respective variable names in a preamble consisting of assignments that execute before the rest of the module. Module-local variables (a challenge outlined earlier) are assigned the return value of a call to wrap, which will be applied only when the module is evaluated and the module-local value becomes available. HARP evaluates the resulting closure, invokes it with the custom context as an argument, and applies further wrap transformations to its return value.

#### 7 IMPLEMENTATION & EVALUATION

In summary, HARP's evaluation answers the following questions:

- **Q1: Can HARP eliminate real vulnerabilities?** HARP successfully eliminates vulnerabilities in 3 large-scale supply-chain attacks (§7.2–7.4) by learning and regenerating the core functionality of the vulnerable library, eliminating any dependency to dangerous code. To the best of our knowledge, HARP is the first system that can eliminate these attacks.
- Q2: How long does ALR take? Applied to 17 JavaScript stringprocessing libraries (§7.5), HARP learns 14 libraries within a minute and all under an hour. It also aborts within 5 seconds on 11 other JavaScript libraries that fall outside the string-processing domain. HARP's domain-specific performance refinements (§6.2) improve the runtime performance of ALR by 179.27×.
- Q3: What are the characteristics of regenerated libraries? The regenerated libraries execute between 2% faster and 7% slower than the original JavaScript libraries. The regenerated libraries import nothing and use only basic JavaScript language primitives. The original libraries, in contrast, have access to the entire JavaScript ecosystem, including standard JavaScript and Node.js libraries, the file system, the network, environment variables, and process arguments.
- Q4: Is ALR applicable outside JavaScript? Applying HARP on 5 native string-processing libraries written in C/C++ and imported as binary modules, successfully regenerates all of them in JavaScript (Appendix C). The regenerated libraries incur a maximum overhead of 1% and enjoy memory and type safety benefits not present in the original libraries.

# 7.1 Methodology

**Workloads:** To investigate Q1, we obtained 3 widely-publicized software supply-chain security incidents from the JavaScript ecosystem: (i) event-stream [40, 59], a popular library that was modified to steal bitcoins from specific Bitcoin wallets (§7.2), (ii) left-pad [36, 67], a popular library replaced by a no-op after a package name dispute, breaking thousands of projects including Facebook and PayPal (§7.3); (iii) string-compare [10], where two different versions of the same string comparison library—one benign and one malicious—appear in the same dependency tree (§7.4).

To investigate Q2 and Q3, we obtained 14 additional JavaScript string processing libraries from npm with the help of an experienced JavaScript developer and a senior undergraduate student. The student used the npm's search feature to search for libraries using a variety of string-processing terms such as "padding", "strip,"

and "change case." For each term, the student sorted the list of re-1045 turned libraries by popularity [43] to inspect the first five pages of 1046 1047 search results and select the library that provided the most complete corresponding functionality. We note that this process excludes 1048 duplicates-for example, the student found and discarded more than 1049 10 left-pad libraries with similar or identical functionality. This 1050 phase produced 17 unique string-processing libraries that are used 1051 pervasively and can affect a large part of the ecosystem [70]: collec-1052 1053 tively, these libraries are directly imported by several applications 1054 and transitively imported via other dependencies by more than 100K applications. The phase also produced 11 libraries that were 1055 misclassified as string processing libraries. We applied HARP to all 1056 28 libraries. 1057

To answer Q4, we obtained C/C++ libraries by searching GitHub using the same search terms as for the JavaScript libraries. Since many of these libraries did not have tests or client programs, we opted for C/C++ libraries with JavaScript bindings to check compatibility via tests and client programs from the JavaScript ecosystem. This search process produced five libraries.

1064 Evaluation metrics: We evaluate security improvements qualita-1065 tively and quantitatively. For known attacks (Q1), we first used the 1066 original (compromised) library to reproduce the attack. We then 1067 inferred and regenerated the original library and replaced the orig-1068 inal library with the regenerated version. We confirmed that the 1069 regenerated version eliminated the attack. For all libraries (Q2-5), 1070 we report the privilege reduction achieved after applying HARP. 1071 This quantitative security metric was developed recently [65] and 1072 corresponds to a ratio  $\alpha/t$ , where  $\alpha$  is the count of all APIs that are 1073 not invocable by the library any more, due to the defense being 1074 applied, and t is the total count of APIs made available to a library 1075 by default by the combined built-in or third-party libraries.

1076 We evaluate the correctness of regenerated libraries (Q3, Q5) us-1077 ing a combination of developer tests, client libraries or applications, 1078 and manual inspection. We ran the developer-provided test suites 1079 for the libraries and verified that the regenerated libraries provide 1080 correct results. We also imported the regenerated libraries into 1081 the top 10 client libraries or applications that directly import the 1082 original libraries and ran the test suites for these client libraries 1083 or applications. Finally, we manually inspected the regenerated 1084 code to confirm that it correctly implements the intended correct 1085 behavior of the original version.

1086 For the learning time (Q2), we report wall-clock time after the 1087 call npm-install up to the point where HARP either (1) aborts, 1088 reporting intractability, (2) timeouts, failing to synthesize a library, 1089 or (3) succeeds, regenerating a library and its appropriate bindings. 1090 We set the timeout limit to 12 hours. We measured the runtime 1091 performance of regenerated-libraries (Q4, Q5) using a combination 1092 of developer tests and synthetic workloads operating in tight loops. 1093 We repeated all performance-related experiments 100 times and 1094 report averages.

Implementation Details: HARP currently works with black-box
libraries available in JavaScript, Python (not shown here; reported in
the extended version [blind]), and binary object files developed, for
example, in C/C++ and wrapped as native add-ons. We expect native
add-ons to be wrapped by some form of language-level interface
such Node's NaN or N-API and Python's ctypes or CFFI. HARP's

1102

Anon.

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1123

1124

1125

1126 1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

ALR components, including the synthesis and DSL, are written in JavaScript. The base set of DSL terms as well as the resulting programs are compiled to their respective language using a small Python compiler: the compiler currently can emit JavaScript and Python programs, which are then executed using the interpreter of the respective language. The regenerated programs link against a small utility library that provides runtime support, ported once for each target language supported by HARP.

HARP currently has a few limitations. First, it does not support libraries whose functions mutate built-in, prototype, or other objects—such as String.prototype in JavaScript. Additionally, HARP's input generation algorithm does not generate non-ASCII strings or ones with special—possibly hierarchical—structure such as JSON, HTML, and CSS; generating the latter without any additional domain information would be impractical.

**Software and Hardware Setup**: All experiments were conducted on a modest server with 4GB of memory and 2 Intel Core2 Duo E8600 CPUs clocked at 3.33GHz, running a Linux kernel version 4.4.0-134. The JavaScript setup uses Node.js v12.19, bundled with V8 v7.8.279.23, LibUV v1.39.0, and npm version v6.14.8; the Python setup uses CPython 3.7.5. To perform timeline-accurate supplychain attacks, we set up a private registry using verdaccio [62] available only to the server running the experiments.

# 7.2 Use Case: Event-Stream

The event-stream incident [40, 59], discussed extensively earlier (§2), introduced a malicious dependency harvesting Bitcoin account credentials through a popular stream-processing library. This dependency, flatmap-stream, targeted a very specific production environment of a cryptocurrency application; other environments were not affected.

**Security**: We reconstruct the malicious library and payloads from a variety of sources [20, 44, 51]. The library applies several checks to verify it runs on production, as part of a specific application, and as part of a specific build. If all these conditions hold, it then writes to the file-system. HARP's active learning phase does not infer any file-system accesses; this is primarily because there are no such accesses during the learning phase, and secondarily because HARP does not model them. As a result, HARP regenerates an exploit-free version of the library, confirmed by manual inspection. It makes no use of built-in APIs, achieving a privilege-reduction of 332×.

**Performance & correctness**: HARP takes on average 1.4 seconds to complete flatmap-stream's active learning and regeneration. We manually inspected the regenerated code and found it implements the full functionality of the original library. The original library does not come with any test cases and the version of event-stream that uses the malicious flatmap-stream version has been removed permanently from npm. We therefore manually modified event-stream commit e316336, introducing flatmap-stream to import and use the regenerated flatmap-stream, and apply event-stream's tests. All 14 (100%) of event-stream's tests pass successfully: 13/14 tests are not affected by the flatmapstream addition, and 1/14 that tests flatmap-stream to an array of 1000 elements over 10K runs takes 48.99 seconds—an overhead of about 107µs per run over the performance of the original library. Supply-Chain Vulnerability Elimination via Active Learning & Regeneration

Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD

## 1161 7.3 Use Case: Left-Pad

1162 The left-pad incident [36, 67] was caused by unpublishing a pop-1163 ular JavaScript library, effectively replacing it permanently with a 1164 No-Op. While left-pad itself was an 11-line moderately-popular 1165 string-padding function, it was used by many popular projects such 1166 as React and Babel. The unpublishing corrupted production envi-1167 ronments, denying them the ability to revert to an older version of 1168 the library. As a result, the incident affected more than one third 1169 of the Node.js ecosystem, and led to significant changes in the 1170 un-publishing policies of public library registries.

1171 Security: We apply HARP to an identical library built by npm as a response to the incident, replacing the original left-pad li-1173 brary copied to our local registry (§7.1). HARP regenerates all of 1174 left-pad's functionality, fully eliminating the dependency. As a 1175 result, left-pad's tests still succeed after we unpublish left-pad 1176 from our local registry because they no longer depend on the origi-1177 nal left-pad module. The regenerated left-pad makes no use of 1178 built-in APIs, resulting in a privilege-reduction score of 332×. 1179

Performance & correctness: HARP completes left-pad's ac-1180 tive learning and regeneration in an average of 3.6 seconds. We 1181 manually inspect the regenerated code and confirm it implements 1182 left-pad's full functionality. We apply the full test suite (35 tests) 1183 from left-pad's repository, all of which (100%) pass successfully. 1184 One test is particularly interesting as it supplies ill-defined input to 1185 trigger left-pad's default behavior, by providing a padding length 1186 of false, it invokes one of left-pad's padding behaviors that 1187 HARP learned through other false-y values. The runtime perfor-1188 mance of the regenerated left-pad on 10K runs is 45.66 seconds-1189 an overhead of about  $20\mu$ s per run over the performance of the 1190 original library. 1191

# 7.4 Use Case: String-Compare

1192

1193

1194

1218

The string-compare attack involves two versions of a single library in the same codebase [10]. The earlier version of this library, used as part of a sort function is benign. A later version, used in the authentication module, is malicious: if provided the (authentication) string gbabWhaRQ, it access the file system of the server running the program.

Security: We apply HARP on both versions of string-compare
 library. The regenerated version is identical between the two cases,
 and does not contain any side-effects—nor the check that launches
 the attack in the second case. HARP eliminates the string-compare
 dependency from both sort and auth, replacing it with vulnerability free code. The regenerated string-compare makes no use of built in APIs, resulting in a privilege-reduction score of 332×.

Performance & correctness: HARP completes string-compare's 1208 active learning and regeneration in an average of 0.7 seconds. We 1209 manually inspect the regenerated code and confirm it implements 1210 string-compare's full functionality. As string-compare comes 1211 1212 with no test cases, we apply the test cases of the sort function over a shuffled version of Ubuntu's wamerican dictionary words file (102K 1213 elements); all 3 (100%) test cases pass. The runtime performance of 1214 10K sort iterations using the regenerated string-compare takes 1215 1216 46.01 seconds—an overhead of about  $41\mu$ s per run over the perfor-1217 mance of the original library.

# 7.5 Applying HARP to More Libraries

In this section, we apply HARP to 25 JavaScript libraries—17 string-processing libraries and 11 other libraries—and 5 C/C++ libraries collected from GitHub.

Table 1 shows results for 17 JavaScript string-processing libraries. The statistics columns  $D_1$ – $D_3$  count the number of weekly downloads, direct dependents, and total dependents of these libraries as reported by the npm tool: collectively, these libraries can affect a significant fraction of the ecosystem-they total 20.22M downloads per week, are directly depended upon by a total of 4,320 libraries, and are indirectly depended upon by more than 117,738 libraries and applications. The Learning columns  $t_1$  and  $t_2$  show the time it took HARP to apply active learning and regeneration;  $t_1$  is full HARP, whereas  $t_2$  does not include HARP's performance refinements (§6). The Regeneration columns Performance and Correctness show the characteristics of the regenerated library with respect to the original: Performance is measured using 10K iterations of several tests; and Correctness is measured by running all the tests of the original library and 10 client-libraries against the regenerated library, followed by manual inspection.

**Learning**: HARP's active learning and regeneration takes between 0.7 seconds and 50.9 minutes (avg.: 204.83 seconds) to complete, with 14 out of 17 libraries regenerated within a minute and 16 out of 17 libraries regenerated within 145 seconds (2.4 minutes). The regenerated camel-case library stands out in terms of size, containing 8 computational statements, two of which are split operators, taking 3059 seconds (50.9 minutes) to regenerate.

HARP's performance refinements (§6) offer significant improvements. Without any refinements, HARP takes at least  $179.27 \times longer$ . This value is a conservative estimate because HARP reaches a timeout limit of 12 hours for 7 out of 17 libraries. This speedup includes a  $1.1 \times slowdown$  for 5 small libraries that are penalized by HARP's refinements. As the regeneration of these libraries remains within a few seconds, their slowdown is considered acceptable—especially given the overall speedup of long-running regenerations.

In terms of code coverage, the input generation algorithm exercises 100% of library code for 11 out of 17 libraries. For the remaining six libraries, the majority of the functionality not exercised is related to exception handling. In the case of decamelize (80%), HARP does not exercise four lines handling erroneous input (non-string arguments), and 11 lines related to a flag preserving consecutive upper case. In the case of flatmap-stream (62.2%), HARP does not exercise a subset of the stream-specific functionality-stream pause, resume, destroy and end handlers-that are part of superclass functionality. HARP also does not exercise exception handlers in the write method, which are meant to handle errors propagating up from the stream consumer. In the case of repeat-string (95.45%), HARP misses an exception-raising statement meant to cover cases where the first argument is not a string. The trim library (33.3%) first checks if the input string's prototype includes a trim method and if so it invokes it; otherwise, implements a left and right trim by invoking other methods-but HARP's primary inputs always support trim as part of the string prototype. In the case of upper-case (44%), HARP misses all the locale-specific code (confirmed by the tests). In the case of zero-fill (80%), HARP misses a branch for

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1219

1220

1221

1222

1223

1224

1225

1226

1227

Anon

**Tab. 1: Applying HARP to JavaScript libraries.** Columns  $D_{1-3}$  show weekly downloads, direct dependents, and total dependents; columns  $t_1$  and  $t_2$  show1278the time it took to complete ALR, with and without refinements; column *Coverage* shows the percentage of the source code covered by the input generation1279algorithm; other columns show the characteristics of the regenerated libraries compared to the original libraries.

Popularity			Learning			Regeneration					
Library	$D_1$	$D_2$	$D_3$	$t_1$ (s)	$t_2$ (s)	Cov/ge (%)	Perf/nc	ce (s) (%)	Correc	tness (%)	
camel-case	126K	842	1314	3059	>12h	100	10.15	(5.6%)	9/9	(100%)	
onstant-case	9.3M	498	1484	22	>12h	100	9.86	(3.8%)	9/9	(100%)	
decamelize	3.8M	110	249	4	14340	80	9.61	(2.9%)	40/40	(100%)	
tmap-stream	15.1M	326	832	1.4	>12h	71.21	48.99	(2.2%)	14/14	(100%)	
left-pad	83K	42	175	3.6	1.3	100	45.66	(0.4%)	35/35	(100%)	
no-case	15.1M	326	832	28	>12h	100	9.89	(5.3%)	31/31	(100%)	
pascal-case	5.9M	89	331	145	>12h	100	10.15	(6.4%)	8/8	(100%)	
repeat-string	2.8M	195	472	19	8	95.45	9.18	(-1.6%)	28/30	(93.3%)	
entence-case	43.1K	26	68	129	>12h	100	10.02	(4.9%)	7/7	(100%)	
snake-case	1.6M	2.2K	832	36	>12h	100	9.95	(3.4%)	8/8	(100%)	
ing-compare	2.2M	71	1651	0.7	1.2	100	46.01	(0.9%)	3/3	(100%)	
trim-left	6.8M	430	2979	3	3.5	100	9.36	(0.3%)	4/4	(100%)	
trim-right	1.8M	43	1636	2	3.2	100	9.43	(2.9%)	4/4	(100%)	
trim	1.5M	219	332	21	111	33.33	9.45	(1.2%)	4/4	(100%)	
upper-case	5.4M	17	1452	3	1.9	44.44	9.38	(0.0%)	4/6	(66.7%)	
write-pad	2.03M	232	1386	3	1.4	100	9.26	(0.9%)	1/10	(100%)	
zero-fill	9.7M	66	2463	3	1.6	80	9.27	(-0.9%)	11/16	(68.8%)	
Min	43.1K	17	68	0.7	1.2	33.33		(-1.6%)		(66.7%)	
Max	20.3M	2200	2979	3K	>12h	100		(6.4%)		(100.0%)	
Avg	6.22M	411	1177	204.83	>10.2h	86.04		(2.3%)		(95.8%)	

when the second argument (the "filler" string) is not provided—in which case the library returns a partially applied function.

Performance: To understand the runtime performance of regen-erated libraries, we apply them on the tests of the original libraries in tight loops of 10K iterations. Their performance is between -1.6% (speedup) and 6.4% (slowdown), with an average of 2.3%. Profil-ing shows that the overhead comes from HARP's complex pattern matching primitives which compile down to the language's regular-expression language (REL). REL is in fact not regular, as it supports back-references and other non-regular constructs, and thus does not perform as efficiently as the simpler string-matching constructs found in the original libraries. 

**Correctness**: Out of 17 libraries, 14 are fully regenerated, passing 100% of developer-provided and client-application tests. Three libraries are partially regenerated, passing between 4/6 (66.7%) and 28/30 (99.3%) of tests. Two of repeat-string's tests are designed to generate exceptions—not expressible in HARP's DSL. Two of upper-case's tests are locale-dependent, with string locales that fall outside HARP's input generation algorithm. Finally, 5 of zero-fill's tests expect a partially evaluated function, which is currently not expressible in HARP's DSL.

## 8 RELATED WORK

Input-Output Synthesis: Program synthesis and programming
by example automatically generate programs that satisfy a given set
of input-output examples [2, 15, 17, 21, 46, 47, 57, 66]. HARP differs
in that it works with an existing component as opposed to a fixed
set of input-output examples and interacts with the component
to build a model of its behavior. The goal is to eliminate dependencies and vulnerabilities by replacing the original version with

the regenerated version, without requiring developers to provide input-output examples manually.

Component-based synthesis [14, 16, 31, 55] aims to generate a program consisting of library calls to a provided API. It synthesizes code for making library calls by executing the candidate program on a set of test cases. HARP's approach does not depend on knowledge of components or interfaces and generates custom inputs to infer library properties rather than using test cases.

Active Learning: Active learning is a classical topic in machine learning [53]. In the context of program inference, it includes learning (and generating) database interfaces, loop-free programs, or declarative logic programs [21, 48, 56]. HARP, in contrast, works in tandem with existing libraries to synthesize vulnerability-free replacements that implement the same functionality with respect to the original library and perform comparably to it—and targets string processing libraries. HARP is the first active learning and regeneration system to successfully eliminate software supply-chain vulnerabilities in widely used libraries.

**Component Protection**: Runtime component protection techniques provide monitoring, instrumentation, and policy enforcement, typically through sandboxing, wrapping, or transformation [1, 18, 23, 30, 33–35, 37, 50, 52, 61]. HARP differs in that it replaces the library with a regenerated version instead of executing the library in a sandbox or wrapping the library to dynamically enforce a security policy. The regenerated library therefore executes with no runtime instrumentation overhead and requires no sandboxing. We note that, to avoid exploitation during inference, HARP uses a combination of sandboxing and wrapping during inference. Unlike sandboxing or wrapping, HARP also protects against library deletion attacks.

Software Debloating: ALR is also related to software debloat-1393 ing [3, 19, 24, 25], a technique that lowers the potential for vulnera-1394 1395 bilities by eliminating unused code in a program. Like debloating, HARP can eliminate unused code in the inferred library. Unlike 1396 debloating, which prunes computation in the original library but 1397 leaves unpruned code intact, HARP replaces the original library with 1398 a regenerated version that is guaranteed to conform to a safe model 1399 of computation. HARP can also discard potentially malicious code 1400 1401 in the regeneration, including code that executes during inference 1402 but does not affect the client-visible behavior of the inferred library. 1403 Vulnerability Detection: Prior work on static [8, 12, 13, 22, 60] 1404 and dynamic [45, 68] analysis can detect malicious code at devel-1405 opment or production time. HARP does not attempt to detect a 1406 vulnerability-rather, it assumes libraries as a potential liability 1407 with stealthy Turing-complete vulnerabilities, and rewrites them 1408 into functionally equivalent, side-effect-free versions. 1409

# 9 DISCUSSION & LIMITATIONS

**Synthesis Limitations**: There are desirable guarantees that Algorithm 1 does not satisfy. First, if the behavior of the original library L does not correspond to any program in the DSL, there is no guarantee that the algorithm will determine this fact in any finite time—it is possible that the difference in behavior will be exposed only by an input that the algorithm has yet to consider. Second, if the behavior of the original library L does correspond to some program in the DSL, there is no guarantee that DSL there is no guarantee that the algorithm will find that DSL program in any finite time—it is possible that the program is larger than programs that the algorithm has considered.

So given a set of DSL programs P at some point in execution 1422 of Algorithm 1, what must be true of the relationship between the 1423 DSL programs  $f \in P$  and the original library L? First, P and all 1424 f exhibit identical behavior on all considered inputs I. Second, if 1425 there is some program f of size n or less that has the same behavior 1426 as *L* on all inputs, then  $f \in P$ . If additionally  $P = \{f\}$ , then the 1427 algorithm will return a new library L' that has identical behavior 1428 as the original library L. There are two key preconditions here 1429 which Algorithm 1 does not check: (1) there is some DSL program 1430 f has identical behavior on all inputs as the original library L, and 1431 (2) this DSL program is of size *n* or less for some known *n*. These 1432 preconditions may come, for example, from the general domain 1433 knowledge of the programmer. 1434

1435Attacks on outputs: As outlined earlier (§3), the primary targets1436of active library learning and regeneration are (1) Turing-complete1437side-effectful attacks—*e.g.*, ones targeting the file system, global1438variables, the module system, process arguments, or environment1439variables, (2) attacks via low-level, memory-unsafe, and type-unsafe1440code such as ones typical in C and C++ code. Could ALR additionally1441protect against attacks targeting the output of a library function?

For attacks targeting function outputs, there are two broad pos-1442 sibilities depending on the malicious behavior. If the malicious be-1443 havior is hidden and therefore not exposed during testing/normal 1444 use, HARP will not learn the malicious behavior and thus the re-1445 generated code will not contain the corresponding vulnerability. If, 1446 however, the malicious behavior is exposed during testing/normal 1447 1448 use, HARP would either (1) determine that the observed behavior is outside the scope of the DSL and reject the library, or (2) learn 1449

and regenerate the behavior. In the latter case, HARP is relying on the exposure of the malicious behavior during testing/normal use to detect and eliminate the behavior—*i.e.*, we would expect the behavior to be detected by the developer during development and before deployment. We anticipate that, at least for string processing programs, almost all such malicious behaviors will be outside the scope of the Harp DSL.

**Generalizing ALR:** Active learning and regeneration is a blackbox program inference approach that fixes (1) a specific computational domain (SCD) such as string processing, tensor operators, database interfaces, (2) a corresponding language (DSL) for modeling computations in that domain, and (3) an input generation algorithm (IGA) for interacting with the black-box computation. These three elements are interlinked and are designed to complement each other. For example, the DSL is designed to enable differential testing, using the IGA to guide efficient inference—by inferring the existence or absence of certain DSL terms in the regenerated programs while minimizing ambiguity.

An active learning and regeneration *system* such as HARP is an instantiation of these three elements (SCD, DSL, IGA) for a particular domain. We do not expect a single system to be expanded to capture all or even a large range of computation of interest. Such an expansion can quickly result in general computations and thus quickly hit known intractability limits.

**Applying ALR to further domains**: Instead, we anticipate multiple active learning and regeneration systems, each targeting a certain class of libraries. HARP exemplifies this approach for string computations—a central, widely used class of computations. Other classes of computation and associated DSLs include: arithmeticoperation libraries, linear algebra and tensor operations, key-value operations, spreadsheet-style computations, components that access SQL databases, blockchain smart contracts, and stream-based parallelizing combiners.

The technique has been applied successfully in some of these domains—*e.g.*, programs and program fragments that access state-ful key-value stores [49], applications that access relational SQL databases [54], binary data parsing and transformation, [9], and parallel and distributed synthesis of Unix shell commands [64]. HARP is the first active learning and regeneration system targeting security problems in the context of software supply-chains.

#### **10 CONCLUSION**

Large-scale dependency incidents such as event-stream have made supply-chain attacks a critical security concern. This paper presents a new approach for addressing this concern: active library learning and regeneration (ALR) to infer and regenerate the clientobservable functionality of a black-box software dependency. ALR replaces a third-party software dependency with one that is automatically regenerated from a domain-specific program in which the target class of attacks cannot be expressed. We demonstrate ALR techniques in HARP, a prototype system for inferring and regenerating components that implement string computations. Applied to JavaScript and C/C++ libraries, HARP completes regeneration in under an hour, regenerates safe versions that are fully compatible with the original library and exhibit minimal to no performance overhead. 1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461 1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1508

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD

#### 1509 **REFERENCES**

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

1566

- Pieter Agten, Steven Van Acker, Yoran Brondsema, Phu H. Phung, Lieven Desmet, and Frank Piessens. 2012. JSand: Complete Client-side Sandboxing of Third-party JavaScript Without Browser Modifications. In *Proceedings of the 28th Annual Computer Security Applications Conference (ACSAC '12)*. ACM, New York, NY, USA, 1-10. https://doi.org/10.1145/2420950.2420952
- [2] Rajeev Alur, Rastislav Bodik, Garvit Juniwal, Milo MK Martin, Mukund Raghothaman, Sanjit A Seshia, Rishabh Singh, Armando Solar-Lezama, Emina Torlak, and Abhishek Udupa. 2013. Syntax-guided synthesis. In 2013 Formal Methods in Computer-Aided Design. IEEE, 1–8.
  - [3] Babak Amin Azad, Pierre Laperdrix, and Nick Nikiforakis. 2019. Less is more: quantifying the security benefits of debloating web applications. In 28th {USENIX} Security Symposium ({USENIX} Security 19). 1697–1714.
  - [4] Daniel J Bernstein, Bernard Van Gastel, Wesley Janssen, Tanja Lange, Peter Schwabe, and Sjaak Smetsers. 2014. TweetNaCl: A crypto library in 100 tweets. In International Conference on Cryptology and Information Security in Latin America. Springer, 64-83. https://tweetnacl.cr.yp.to/
  - [5] Oscar Bolmsten. 2017. Malicious Package: crossenv and other 36 malicious packages. https://snyk.io/vuln/npm:crossenv:20170802 Accessed: 2019-03-19.
     [6] Benjamin Byholm, Rod Vagg, and NAN contributors. 2018. Native Abstractions
  - for Node. https://www.npmjs.com/package/nan Accessed: 2020-06-11.
  - [7] Mircea Cadariu, Eric Bouwers, Joost Visser, and Arie van Deursen. 2015. Tracking known security vulnerabilities in proprietary software systems. In Software Analysis, Evolution and Reengineering (SANER), 2015 IEEE 22nd International Conference on. IEEE, 516–519.
- 1527 Conference on. IEEE, 516–519.
   [8] Stefano Calzavara, Michele Bugliesi, Silvia Crafa, and Enrico Steffinlongo. 2015.
   1528 Fine-Grained Detection of Privilege Escalation Attacks on Browser Extensions. In
   1529 Programming Languages and Systems 24th European Symposium on Programming, ESOP 2015, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2015, London, UK, April 11-18, 2015. Proceedings (Lecture Notes in Computer Science), Jan Vitek (Ed.), Vol. 9032. Springer, 510–534. https://doi.
   1532 org/10.1007/978-3-662-46669-8\_21
- [9] José P. Cambronero, Thurston H. Y. Dang, Nikos Vasilakis, Jiasi Shen, Jerry Wu, and Martin C. Rinard. 2019. Active Learning for Software Engineering. In Proceedings of the 2019 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward! 2019). Association for Computing Machinery, New York, NY, USA, 62–78. https://doi. org/10.1145/3359591.3359732
- [10] David Bryant Copeland. 2019. The Frightening State of Security Around NPM Package Management. https://bit.ly/3pID2h1 Accessed: 2020-12-10.
   [10] David Bryant Copeland. 2019. The Frightening State of Security Around NPM
  - [11] Ruian Duan, Omar Alrawi, Ranjita Pai Kasturi, Ryan Elder, Brendan Saltaformaggio, and Wenke Lee. 2021. Towards Measuring Supply Chain Attacks on Package Managers for Interpreted Languages. NDSS.
  - [12] Aurore Fass, Michael Backes, and Ben Stock. 2019. JStap: A Static Pre-Filter for Malicious JavaScript Detection. In Proceedings of the 35th Annual Computer Security Applications Conference (ACSAC '19). Association for Computing Machinery, New York, NY, USA, 257–269. https://doi.org/10.1145/3359789.3359813
  - [13] Aurore Fass, Robert P. Krawczyk, Michael Backes, and Ben Stock. 2018. JaSt: Fully Syntactic Detection of Malicious (Obfuscated) JavaScript. In Detection of Intrusions and Malware, and Vulnerability Assessment, Cristiano Giuffrida, Sébastien Bardin, and Gregory Blanc (Eds.). Springer International Publishing, Cham, 303–325.
  - [14] Yu Feng, Ruben Martins, Yuepeng Wang, Isil Dillig, and Thomas W Reps. 2017. Component-based synthesis for complex APIs. In Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages. 599–612.
  - [15] John K Feser, Swarat Chaudhuri, and Isil Dillig. 2015. Synthesizing data structure transformations from input-output examples. In ACM SIGPLAN Notices, Vol. 50. ACM, 229–239.
  - [16] Joel Galenson, Philip Reames, Rastislav Bodik, Björn Hartmann, and Koushik Sen. 2014. Codehint: Dynamic and interactive synthesis of code snippets. In Proceedings of the 36th International Conference on Software Engineering. 653–663.
  - [17] Sumit Gulwani. 2011. Automating string processing in spreadsheets using inputoutput examples. In ACM Sigplan Notices, Vol. 46. ACM, 317–330.
  - [18] Daniel Hedin, Alexander Sjösten, Frank Piessens, and Andrei Sabelfeld. 2017. A principled approach to tracking information flow in the presence of libraries. In International Conference on Principles of Security and Trust. Springer, 49–70.
  - [19] Kihong Heo, Woosuk Lee, Pardis Pashakhanloo, and Mayur Naik. 2018. Effective program debloating via reinforcement learning. In Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. 380–394.
  - [20] hugeglass. 2018. GitHub Repository for flatmap-stream. https://git.io/Jtcdi Accessed: 2020-12-18.
  - [21] Susmit Jha, Sumit Gulwani, Sanjit A Seshia, and Ashish Tiwari. 2010. Oracleguided component-based program synthesis. In Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 1. ACM, 215–224.
  - [22] N. Jovanovic, C. Kruegel, and E. Kirda. 2006. Pixy: a static analysis tool for detecting Web application vulnerabilities. In 2006 IEEE Symposium on Security and Privacy (S P'06). 6 pp.-263. https://doi.org/10.1109/SP.2006.29

- [23] Yoonseok Ko, Tamara Rezk, and Manuel Serrano. [n. d.]. SecureJS Compiler: Portable Memory Isolation in JavaScript. In SAC 2021-The 36th ACM/SIGAPP Symposium On Applied Computing.
- [24] Igibek Koishybayev and Alexandros Kapravelos. 2020. Mininode: Reducing the Attack Surface of Node.js Applications. In 23rd International Symposium on Research in Attacks, Intrusions and Defenses ({RAID} 2020).
- [25] Hyungjoon Koo, Seyedhamed Ghavamnia, and Michalis Polychronakis. 2019. Configuration-Driven Software Debloating. In Proceedings of the 12th European Workshop on Systems Security. 1–6.
- [26] Benjamin Lamowski, Carsten Weinhold, Adam Lackorzynski, and Hermann Härtig. 2017. Sandcrust: Automatic Sandboxing of Unsafe Components in Rust. In Proceedings of the 9th Workshop on Programming Languages and Operating Systems (PLOS'17). ACM, New York, NY, USA, 51–57. https://doi.org/10.1145/ 3144555.3144562
- [27] Tobias Lauinger, Abdelberi Chaabane, Sajjad Arshad, William Robertson, Christo Wilson, and Engin Kirda. 2017. Thou Shalt Not Depend on Me: Analysing the Use of Outdated JavaScript Libraries on the Web. (2017).
- 28] SS Jeremy Long. 2015. OWASP Dependency Check. (2015).
- [29] Michael Maass. 2016. A Theory and Tools for Applying Sandboxes Effectively. Ph.D. Dissertation. Carnegie Mellon University.
- [30] Jonas Magazinius, Daniel Hedin, and Andrei Sabelfeld. 2014. Architectures for inlining security monitors in web applications. In *International Symposium on Engineering Secure Software and Systems*. Springer, 141–160.
- [31] David Mandelin, Lin Xu, Rastislav Bodík, and Doug Kimelman. 2005. Jungloid mining: helping to navigate the API jungle. ACM Sigplan Notices 40, 6 (2005), 48-61.
- [32] Marcela S Melara, David H Liu, and Michael J Freedman. 2019. Pyronia: Redesigning Least Privilege and Isolation for the Age of IoT. arXiv preprint arXiv:1903.01950 (2019).
- [33] Leo A Meyerovich and Benjamin Livshits. 2010. ConScript: Specifying and enforcing fine-grained security policies for Javascript in the browser. In 2010 IEEE Symposium on Security and Privacy. IEEE, 481–496.
- [34] James Mickens. 2014. Pivot: Fast, synchronous mashup isolation using generator chains. In 2014 IEEE Symposium on Security and Privacy. IEEE, 261–275.
- [35] Mark S Miller, Mike Samuel, Ben Laurie, Ihab Awad, and Mike Stay. 2009. Caja: Safe active content in sanitized JavaScript, 2008. *Google white paper* (2009).
   [36] Paul Miller. 2016. How an irate developer briefly broke JavaScript. https://doi.org/10.1016/j.
- [36] Paul Miller. 2016. How an irate developer briefly broke JavaScript. https: //bit.ly/36CkBDI Accessed: 2020-12-10.
   [37] Marius Musch, Marius Steffens, Sebastian Roth, Ben Stock, and Martin Johns. 2019.
- [37] Martus Musch, Martus Steffens, Sebastian Roth, Ben Stock, and Martin Johns. 2019. ScriptProtect: mitigating unsafe third-party javascript practices. In Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security. 391–402.
- [38] Shravan Narayan, Craig Disselkoen, Tal Garfinkel, Nathan Froyd, Eric Rahm, Sorin Lerner, Hovav Shacham, and Deian Stefan. 2020. Retrofitting Fine Grain Isolation in the Firefox Renderer. In 29th {USENIX} Security Symposium ({USENIX} Security 20). 699–716.
- [39] Nick Nikiforakis, Luca Invernizzi, Alexandros Kapravelos, Steven Van Acker, Wouter Joosen, Christopher Kruegel, Frank Piessens, and Giovanni Vigna. 2012. You are what you include: large-scale evaluation of remote javascript inclusions. In Proceedings of the 2012 ACM conference on Computer and communications security. 736–747.
- [40] npm, Inc. 2018. Details about the event-stream incident. https://blog.npmjs.org/ post/180565383195/details-about-the-event-stream-incident Accessed: 2018-12-18.
- [41] npm, Inc. 2019. Malicious Package: stream-combine. https://www.npmjs.com/ advisories/774 Accessed: 2019-01-25.
   [42] npm, Inc. 2019. Malicious Package: stream-combine. https://www.npmis.com/
- [42] npm, Inc. 2019. Malicious Package: stream-combine. https://www.npmjs.com/ advisories/765 Accessed: 2019-01-25.
- [43] npm, Inc. 2020. Node Package Manager. https://www.npmjs.com/search?q= string&ranking=popularity
- [44] Jarrod Overson. 2018. BadJS-Malicious JavaScript found in the wild: Event-Stream. https://badjs.org/posts/event-stream/ Accessed: 2020-12-18.
- [45] Giancarlo Pellegrino and Davide Balzarotti. 2014. Toward Black-Box Detection of Logic Flaws in Web Applications.
- [46] Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program synthesis from polymorphic refinement types. ACM SIGPLAN Notices 51, 6 (2016), 522–538.
- [47] Mohammad Raza and Sumit Gulwani. 2018. Disjunctive Program Synthesis: A Robust Approach to Programming by Example. In *Thirty-Second AAAI Conference* on Artificial Intelligence.
- [48] Martin C. Rinard, Jiasi Shen, and Varun Mangalick. [n. d.]. Active learning for inference and regeneration of computer programs that store and retrieve data. In Proceedings of the 2018 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software, Onward! 2018, Boston, MA, USA, November 7-8, 2018, Elisa Gonzalez Boix and Richard P. Gabriel (Eds.).
- [49] Martin C. Rinard, Jiasi Shen, and Varun Mangalick. 2018. Active Learning for Inference and Regeneration of Computer Programs That Store and Retrieve Data.

Anon.

1567

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1618

1619

1620

1621

1622

1626

1637

1638

1639

1649

1650

1651

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1674

1679

1680

1681

1682

Anonymous submission #59 to ACM CCS 2021, Due Jan. 20, 2021, Seoul, TBD

- In Proceedings of the 2018 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward! 2018). ACM, New York, NY, USA, 12–28. https://doi.org/10.1145/3276954.3276959
- [50] José Fragoso Santos and Tamara Rezk. 2014. An information flow monitorinlining compiler for securing a core of javascript. In *IFIP International Information Security Conference*. Springer, 278–292.
- [51] Thomas Hunter II (Intrinsic Security). 2018. Compromised npm Package: eventstream. https://medium.com/intrinsic/compromised-npm-package-eventstream-d47d08605502 Accessed: 2019-03-19.
- [52] R. Sekar, V.N. Venkatakrishnan, Samik Basu, Sandeep Bhatkar, and Daniel C. DuVarney. 2003. Model-Carrying Code: A Practical Approach for Safe Execution of Untrusted Applications. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles (SOSP '03)*. Association for Computing Machinery, New York, NY, USA, 15–28. https://doi.org/10.1145/945445.945448
- [53] Burr Settles. 2009. Active Learning Literature Survey. Computer Sciences Technical Report 1648. University of Wisconsin–Madison.
  - [54] Jiasi Shen and Martin C. Rinard. 2019. Using Active Learning to Synthesize Models of Applications That Access Databases. In Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2019). ACM, New York, NY, USA, 269–285. https://doi.org/10.1145/3314221.3314591
- [55] Kensen Shi, Jacob Steinhardt, and Percy Liang. 2019. FrAngel: component-based synthesis with control structures. *Proceedings of the ACM on Programming Languages* 3, POPL (2019), 1–29.
- [56] Xujie Si, Woosuk Lee, Richard Zhang, Aws Albarghouthi, Paraschos Koutris, and Mayur Naik. 2018. Syntax-Guided Synthesis of Datalog Programs. In Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2018).
   Association for Computing Machinery, New York, NY, USA, 515–527. https: //doi.org/10.1145/3236024.3236034
- [57] Rishabh Singh. 2016. Blinkfill: Semi-supervised programming by example for syntactic string transformations. *Proceedings of the VLDB Endowment* 9, 10 (2016), 816–827.
   [648] Self-Self-Carlo and Self-Carlo and Self-C
  - [58] Snyk. 2016. Find, fix and monitor for known vulnerabilities in Node.js and Ruby packages. https://snyk.io/
  - [59] Ayrton Sparling et al. 2018. Event-Stream, GitHub Issue 116: I don't know what to say. https://github.com/dominictarr/event-stream/issues/116 Accessed: 2018-12-18.
- [60] Cristian-Alexandru Staicu, Michael Pradel, and Benjamin Livshits. 2018. Synode: Understanding and Automatically Preventing Injection Attacks on Node.js. In Networked and Distributed Systems Security (NDSS'18). https://doi.org/10.14722/ ndss.2018.23071
- [61] Cristian-Alexandru Staicu, Daniel Schoepe, Musard Balliu, Michael Pradel, and Andrei Sabelfeld. 2019. An empirical study of information flows in real-world javascript. In Proceedings of the 14th ACM SIGSAC Workshop on Programming Languages and Analysis for Security. 45–59.
  - [62] Trent Earl, John Wilkinson, and the Verdaccio contributors. 2018. Verdaccio-npm Proxy Private Registry. https://verdaccio.org/ Accessed: 2020-11-10.
  - [63] Neline van Ginkel, Willem De Groef, Fabio Massacci, and Frank Piessens. 2019. A Server-Side JavaScript Security Architecture for Secure Integration of Third-Party Libraries. Security and Communication Networks 2019 (2019).
  - [64] Nikos Vasilakis, Jiasi Shen, and Martin Rinard. 2020. Automatic Synthesis of Parallel and Distributed Unix Commands with KumQuat. CoRR abs/2012.15443 (2020). arXiv:2012.15443 https://arxiv.org/abs/2012.15443
  - [65] Nikos Vasilakis, Cristian-Alexandru Staicu, Grigoris Ntousakis, Konstantinos Kallas, Ben Karel, André DeHon, and Michael Pradel. 2020. Mir: Automated Quantifiable Privilege Reduction Against Dynamic Library Compromise in JavaScript. arXiv preprint arXiv:2011.00253 (2020).
  - [66] Navid Yaghmazadeh, Xinyu Wang, and Isil Dillig. 2018. Automated migration of hierarchical data to relational tables using programming-by-example. *Proceedings* of the VLDB Endowment 11, 5 (2018), 580–593.
  - [67] Serdar Yegulalp. 2016. How one yanked JavaScript package wreaked havoc. https://bit.ly/3ofwkz2 Accessed: 2020-12-10.
- [68] Heng Yin, Dawn Song, Manuel Egele, Christopher Kruegel, and Engin Kirda.
   2007. Panorama: Capturing System-Wide Information Flow for Malware Detection and Analysis. In Proceedings of the 14th ACM Conference on Computer and Communications Security (CCS '07). Association for Computing Machinery, New York, NY, USA, 116–127. https://doi.org/10.1145/1315245.1315261
  - [69] Nicholas C. Zakas and ESLint contributors. 2013. ESLint-Pluggable JavaScript linter. https://eslint.org/ Accessed: 2018-07-12.
- [675] [70] Markus Zimmermann, Cristian-Alexandru Staicu, Cam Tenny, and Michael Pradel.
   [676] Markus Zimmermann, Cristian-Alexandru Staicu, Cam Tenny, and Michael Pradel.
   [676] 2019. Smallworld with High Risks: A Study of Security Threats in the Npm Ecosystem. In *Proceedings of the 28th USENIX Conference on Security Symposium* (*SEC'19*). USENIX Association, USA, 995–1010.

A SUMMARY OF SHEPHERDING CHANGES

We have addressed all four requirements in the revised version, and marked the parts of the revised paper with a different color to aid our shepherd in identifying these changes.

Sooel Son <son.sooel@gmail.com>July 25, 2021(1) Soundness/Completeness issues: please include a separate<br/>section that addresses the limitations of Harp in terms of its<br/>soundness and completeness. Overall in the paper, the au-<br/>thors promised to redefine the soundness and completeness<br/>properties.

We have addressed this requirement by completely rewriting Sec. 5 to redefine and explain the guarantees offered by HARP. We have also added a new Section on "Discussion & Limitations" (§9), which starts with the limitations of the theoretical guarantees provided by HARP's synthesis algorithm. We have also changed the definitions in Appendix B, and note that we plan on rewriting other parts of the paper (*e.g.*, the outline of Alg. 1) to make these properties easier for the reader to extract and understand.

(2) Possible future work: We recommend discussing possible future work in this direction of the work, such as identifying productive classes of computation that a DSL can capture for inference and regeneration. Also, the authors need to clarify the scope and meaning of synthesizing string computations in software supply chain vulnerabilities.

We have addressed this requirement in the new Sec. 9—closing with a discussion concurrent and future work in this line of research. The lifting of double-blind constraints has allowed us to expand on and delineate promising prior and concurrent work. We outline seven promising classes of computations, and provide pointers to relevant concurrent developments in contexts outside that of the security problems associated with software supply-chain attacks. We note that we plan on rewriting other parts of the paper such as Sec. 1 and 2 to clarify the scope and meaning of synthesizing string computations in software supply chain vulnerabilities.

(3) Code coverage: Please add code coverage of the collected input-output pairs, which indirectly represent how many functionalities from the original library are covered.

We have addressed this requirement by adding (1) a new column in Tab. 1 showing the code-coverage results of the collected inputoutput pairs, and (2) a paragraph explaining these results and summarizing functionality not exercised by HARP. Due to timing constraints, *this version of the paper* reports only on code coverage results on functionality exercised only within a 10-minute timeout per library function; in the final version of the paper we will include numbers for the full synthesis runs (*i.e.*, for the full time reported in column  $t_2(s)$ ).

(4) Algorithm 1: Expand and explain all subprocedures in Algorithm 1.

We have addressed this requirement by expanding Sec. 4.2 to include, explain, and exemplify all five sub-procedures mentioned in Algorithm 1.

1740

15 A

**Other changes**: We have also added information about the popularity of libraries in Tab. 1. Due to tight timing constraints, we have not yet addressed other reviewer suggestions, but plan on doing so as promised in the rebuttal. *Thank you!* 

#### **B PROOF SKETCHES**

Definition B.1. **(IO-Correctness)** Given a function f, the synthesized function f' is said to be IO-correct, if and only if, f' is expressible in the Harp DSL (with constants extracted from f) and for all input i consistent with the input type of f, f(i) = f'(i).

Definition B.2. (Consistentency w.r.t. function f, input set I, and maximum program size n) A synthesized function f' is said to be consistent w.r.t. a function f, input set I of size m containing inputs consistent with the input type of f, and a maximum program size n, if f' is expressible in the Harp DSL (with constants extracted from f) and is of size less than equal to n, and for all inputs  $i \in I$ :

$$f(i) = f'(i)$$

Note that, given a function f, the IO-correct function f' is consistent w.r.t. function f, for any input set I (consistent with the input type of f), and any maximum program size n greater than the size of f'.

THEOREM B.3. (Initial State) For any function f, all functions f'in  $P_n$  are consistent w.r.t function f, input set  $I = \emptyset$ , and maximum program size size n. Also, if there exists a function f' of size less than equal to n, which is IO-correct with respect to f, then  $f' \in P_n$ .

PROOF. The Harp algorithm extracts all constants from function f and instantiates all *sketches* is the Harp DSL of size n. The Type T (extracted using soundTypeConstraints) is a *sound* approximation of the actual output type of f. A function  $f' \in P_n$  if and only if f' of size less than equal to n, is expressible in the Harp DSL (with constants extracted from f), and T is a sound approximation of f''s output type. Therefore, given  $I = \emptyset$ , all functions in  $f' \in P_n$  are consistent with f (input set  $I = \emptyset$  and max size n).

Also, if there exists a IO-correct function f' of size less than equal to n, then f' is consistent with respect to f ( $I = \emptyset$  and max size n) and T is a *sound* approximation of the output type of f'. Therefore, if there exists a IO-correct function f' of size less than equal to n, then  $f' \in P_n$ .

THEOREM B.4. (Consistency) Given a function  $f \in L$ , let I be the set of inputs returned by the function generatelnputs,  $P_n$  be the set of programs returned by allPrograms, and P be the set of pruned program pruneSpace. If  $P \neq \emptyset$  and f' is equal to getOpt(P), then f' is consistent w.r.t. function f, input set I, and maximum program size n. Also, if the IO-correct function  $f' \in P_n$ , then  $f' \in P$ .

**PROOF.** pruneSpace only prunes a function  $f' \in P_n$  if and only if  $\exists i \in I$ , such that  $f'(i) \neq f(i)$ . Therefore, all  $f' \in P$  are consistent with respect to f (input set I and max size n). The getOpt returns a function  $f' \in P$ , therefore if the algorithm synthesis a function f' for function f, then f' is consistent with respect to f (input set I and max size n).

pruneSpace will never prune out the IO-correct function f' as for all inputs f(i) = f'(i). Therefore, if  $f' \in P_n$ , then  $f' \in P$ .  $\Box$ 

THEOREM B.5. (Convergence) Given a function f and a maximum function size n, let  $F_n$  be the set of functions in the HARP DSL of size less than equal to n, such that, a IO-correct function  $f' \in P_n$ . As we add more inputs to the set of inputs generated by function generatelnputs, HARP will synthesize a function f'', such that, f'' and f have the same output on an increasing set of inputs.

PROOF.  $P_n$  is equal to the set returned by allPrograms(n, T). From Theorem B.3,  $f' \in P_n$ . Let I be the set of input set constructed by generateInputs. Let  $P_I$  be the set of programs returned by the function pruneSpace. Note that if  $f' \in P$ , then for all  $I, f' \in P_I$ (Theorem B.4).

Note that, if  $I_0 \subseteq I_1$ , then  $P_{I_1} \subseteq P_{I_0}$  (as for any function  $f'' \in P_{I_1}$ , then f'' has the same output as f on inputs in  $I_0$ ).

A larger set of inputs allows HARP to prune out functions which do not have the same output as f on this larger set of inputs. Therefore, by adding more inputs, HARP will synthesize a function f'', such that, f'' and f have the same output on an increasing set of inputs.

#### C ADDITIONAL EVALUATION RESULTS

**Non-string-processing Libraries**: We also apply HARP on 11 libraries that were misclassified as processing strings, to evaluate HARP's --quick-abort mechanism. On these libraries, HARP aborts ALR within 5 seconds with a warning that they contain side-effectful computations that cannot be learned. Eight of these libraries import built-in modules that are not supported by HARP such as debug, http, or fs—for example, minimatch depends on fs and is thus not inferable. One of these libraries, chalk, depends indirectly on os and tty for checking the environment for color support and thus it not inferable. Finally, ignore and attn provide their functionality by extending the runtime context with an auxiliary value.

**C/C++ Libraries:** Fig. 7 summarizes results of applying HARP to 5 C/C++ libraries, including the time to complete learning (column ALR), the regenerated-library performance (column P(L') with positive values for slowdown and negative for speedup), and its correctness with respect to

L	ALR	P(L')	C(L')
string-upper	2.9 <i>s</i>	1.3%	66.7%
right-trim	2.7 <i>s</i>	1.8%	100%
left-trim	2.6 <i>s</i>	0.7%	100%
lr-trim	46.7s	0.4%	100%
repeat-text	17.1 <i>s</i>	0.7%	100%

**Fig. 7: C/C++ ALR.** HARP applied to C/C++ libraries.

the original one (column C(L'), counting percentages of test cases). These libraries export a single function and are wrapped with Node's NAN module [6].<sup>1</sup>

HARP'S ALR ranges between 2.6–17.1s (avg.: 14.4s), driven by the size of the regenerated library. Naturally, the performance of the regenerated JavaScript libraries is lower that that of the original compiled libraries, and ranges between 0.4–1.8% (avg.: 1.0%) of the original library's runtime performance (Col. P(L')). HARP regenerated full library behavior, except string-upper's locale-dependent functionality.

<sup>&</sup>lt;sup>1</sup>NAN is an abstraction layer meant to simplify the development and maintenance of native add-ons over a constantly changing V8 API.