

Generating Well-Typed Terms that are not “Useless”

JUSTIN FRANK, University of Maryland, USA

BENJAMIN QUIRING, University of Maryland, USA

LEONIDAS LAMPROPOULOS, University of Maryland, USA

Random generation of well-typed terms lies at the core of effective random testing of compilers for functional languages. Existing techniques have had success following a top-down type-oriented approach to generation that makes choices locally, which suffers from an inherent limitation: the type of an expression is often generated independently from the expression itself. Such generation frequently yields functions with argument types that cannot be used to produce a result in a meaningful way, leaving those arguments unused. Such “use-less” functions can hinder both performance, as the argument generation code is dead but still needs to be compiled, and effectiveness, as a lot of interesting optimizations are tested less frequently.

In this paper, we introduce a novel algorithm that is significantly more effective at generating functions that use their arguments. We formalize both the “local” and the “nonlocal” algorithms as step-relations in an extension of the simply-typed lambda calculus with type and arguments holes, showing how delaying the generation of types for subexpressions by allowing nonlocal generation steps leads to “useful” functions. We implement our algorithm demonstrating that it’s much closer to real programs in terms of argument usage rate, and we replicate a case study from the literature that finds bugs in the strictness analyzer of GHC, with our approach finding bugs four times faster than the current state-of-the-art local approach.

CCS Concepts: • **Software and its engineering** → **Software testing and debugging**.

Additional Key Words and Phrases: test generation, property-based testing, well-typed lambda terms

ACM Reference Format:

Justin Frank, Benjamin Quiring, and Leonidas Lampropoulos. 2022. Generating Well-Typed Terms that are not “Useless”. *Proc. ACM Program. Lang.* 1, CONF, Article 1 (January 2022), 22 pages.

1 INTRODUCTION

Random generation of programs is a tried and trusted technique for finding bugs in optimizing compilers, as exemplified in the CSmith project (Yang et al. 2011), which uncovered hundreds of bugs in widely used C compilers. At the same time, effective and efficient random generation of programs is a highly involved process, as most compilers impose both syntactic and semantic constraints on inputs before producing meaningful output that can be examined to discover potential errors. For example, the CSmith authors went to great lengths to ensure that the generated C programs did not exhibit undefined behavior.

In a functional setting, constrained program generation takes the form of generating well-typed lambda terms, with Palka et al. (2011) setting the bar by finding bugs in GHC’s strictness analyzer, Midtgaard et al. (2017) extending the same approach to handle effects to target OCaml, and Hoang et al. (2022) further building on top of it to target a higher-order blockchain language based on System F. In their seminal work, Palka et al. introduce an efficient local method of generating well-typed STLC terms by inverting the typing rules and viewing them as production ones. For concreteness, consider the standard typing rule for applications in STLC:

$$\frac{\text{TAPP} \quad \Gamma \vdash e_1 : \tau' \rightarrow \tau \quad \Gamma \vdash e_2 : \tau'}{\Gamma \vdash (e_1 e_2) : \tau}$$

Authors’ addresses: Justin Frank, jppfrank@umd.edu, University of Maryland, USA; Benjamin Quiring, bquiring@umd.edu, University of Maryland, USA; Leonidas Lampropoulos, leonidas@umd.edu, University of Maryland, USA.

2022. 2475-1421/2022/1-ART1 \$15.00

<https://doi.org/>

Using the standard typing rule interpretation, TAPP states that an application $(e_1 e_2)$ has type τ if e_1 has type $\tau' \rightarrow \tau$ and e_2 has type τ' . Instead, Pałka et al. proposed an alternative reading of this rule where the conclusion prescribes the shape of a well-typed term and its premises give rise to recursive calls for generating appropriately typed subexpressions. Viewing TAPP in this way, in order to generate an expression given some type τ , we can opt to generate an application $(e_1 e_2)$. In turn, we would need to recursively generate a function e_1 of type $\tau' \rightarrow \tau$ and an argument for it e_2 of type τ' to satisfy the rule's premises.

While this recursive approach is elegantly compositional and relatively straightforward to implement, it suffers from a significant limitation which manifests in the TAPP rule: the type τ' is completely arbitrary. That means that in order to recursively generate e_1 and e_2 , we must first generate τ' . Unfortunately, as Pałka et al. point out, many choices for τ' can back the generation algorithm into a corner, which it can only get out of by using backtracking, significantly impacting the performance of this approach. To tackle this problem, the authors invented a series of heuristics to improve their proposed top-down generation algorithm. For instance, a particularly effective one was to pick the type τ' so that it matches the argument type of a variable in the environment.

However, the inherent limitation still lurks: *the generation of the type of an expression is independent from the generation of the expression itself.*

That is the root cause behind the excessive backtracking which could cause top-down generators to be potentially *inefficient*. Additionally, this limitation also causes top-down generators to be potentially *ineffective*: decoupling the generation of types and expressions significantly biases generation towards functions that *do not use their arguments*. Such “use-less” function arguments lead to wasted computation, both in terms of generation time (as the code that generates such arguments is essentially dead) and in terms of compilation time (as the test runner now needs to compile a larger program).

In this paper, we propose a new way of generating well-typed lambda terms that is guided by creating variable uses and bindings at the same time. We achieve that by deferring generation of types until they are actually needed, borrowing a page from many a playbook in the random testing literature that rely on laziness to improve test generation (Lampropoulos et al. 2017a; Reich et al. 2012; Runciman et al. 2008).

Concretely, our contributions are:

- We formalize the algorithm of Pałka et al. (2011) by extending STLC to a calculus with typed holes (called λ_{\square}), in which their top-down algorithm takes the form of a big-step relation, rendering the limitation in the generation algorithm glaringly apparent (Section 2).
- We introduce a novel algorithm for generating well-typed lambda terms, that is able to dynamically extend function types during generation. We formalize it by further extending λ_{\square} to a calculus with arguments holes (called λ_{\circ}), in which our lazy generation algorithm manifests as a small-step relation, allowing for more fine-grained control of the resulting distribution (Section 3).
- We mechanized the metatheory of λ_{\circ} in Coq and prove type soundness of our generation rules.
- We describe an efficient implementation of our algorithm in OCaml (Section 3.4),¹ and show that it can replicate finding the bugs in Pałka et al. (2011) in a quarter the time and half the number of tests (Section 4.1).
- We show that terms produced by our generator have fewer unused function parameters, and can approach the usage amounts of real programs of similar size (Section 4.2).

We extensively discuss related work in Section 5 and conclude by drawing directions for future work in Section 6.

¹Available at <https://github.com/bquiring/well-typed-term-generator>

2 THE “LOCAL” GENERATOR

In this section, we will formalize the type-directed generation approach of [Palka et al. \(2011\)](#) by introducing a simply-typed lambda calculus with typed holes. We will explicitly lay out the production rules that are obtained by inverting the typing relation in the style of [Palka et al.](#), which serves both to establish the limitation of the original algorithm that makes generating well-typed terms that use their arguments less likely and as a platform to build upon when introducing our own variant of well-typed generation in the rest of the paper.

For our formalization, we opt to use a multary version of the lambda calculus. Our original motivation for the project was to evaluate the correctness of compilers written by students in Racket for the purposes of an undergraduate course, which explicitly doesn’t have (non-explicit) partial application. Hence, it was natural to experiment with generating terms in a multary calculus. That said, there are not many differences in practice between the unary and multary settings: converting from a multary representation to a unary one is trivial; the other way around is a bit trickier—while multary functions can be embedded into a unary calculus using tuples, the punning between multiple arguments and first class values requires careful consideration of the generation rules to ensure that we don’t generate terms like

$$\text{call } f \text{ (if } \dots \text{ then } (1, 2, 3) \text{ else } (4, 5, 6))$$

which is difficult to translate into a Racket function call.

2.1 STLC with typed holes: λ_{\square}

Expressions e in the simply-typed lambda calculus with typed holes (which we will write as λ_{\square}) are either one of the standard expressions for an STLC with multiple arguments (function calls, λ expressions, variables, or constants of some base type τ_b) or a *typed hole* \square_{τ} . The point of such a hole is to be filled with with an expression of its corresponding type τ , and we will use them to explicitly formulate program generation.

e	$=$	$e e \dots$	Func. Applications
		$\lambda (x : \tau \dots) e$	λ abstractions
		x	Variables
		c	Constants
		\square_{τ}	Typed holes
τ	$=$	τ_b	Base types
		$\tau \dots \rightarrow \tau$	Function Types

Associated with the expression language are the following typing derivation rules, where environments Γ are partial mappings from variables to types as usual:

$\frac{\text{TAPP} \quad \Gamma \vdash e_f : \tau_1 \dots \tau_n \rightarrow \tau \quad \Gamma \vdash e_1 : \tau_1 \dots \Gamma \vdash e_n : \tau_n}{\Gamma \vdash (e_f e_1 \dots e_n) : \tau}$	$\frac{\text{TVAR} \quad [x \mapsto \tau] \in \Gamma}{\Gamma \vdash x : \tau}$	
$\frac{\text{TABS} \quad \Gamma[x_1 \mapsto \tau_1 \dots x_n \mapsto \tau_n] \vdash e : \tau}{\Gamma \vdash \lambda (x_1 : \tau_1 \dots x_n : \tau_n) e : \tau_1 \dots \tau_n \rightarrow \tau}$	$\frac{\text{TCONST} \quad c \text{ is a constant of type } \tau}{\Gamma \vdash c : \tau}$	$\frac{\text{THOLE}}{\Gamma \vdash \square_{\tau} : \tau}$

All rules are standard with the exception of the typing rule for holes: a hole \square_{τ} has type τ under any environment.

2.2 Term Generation as a Step Relation.

Random generation in the style of [Palka et al.](#) can be encoded in λ_{\square} as a big-step reduction relation, that instantiates typed holes. Each typing rule with the exception of THOLE gives rise to a corresponding production rule of the form $\square_{\tau} \Downarrow_{\Gamma} e$, where each production rule is annotated by an enclosing environment Γ .

$$\frac{\text{GVAR} \quad [x \mapsto \tau] \in \Gamma}{\square_{\tau} \Downarrow_{\Gamma} x} \quad \frac{\text{GCONST} \quad \emptyset \vdash c : \tau}{\square_{\tau} \Downarrow_{\Gamma} c}$$

$$\frac{\text{GABS} \quad \square_{\tau} \Downarrow_{\Gamma[x_1 \mapsto \tau_1 \dots x_n \mapsto \tau_n]} e}{\square_{\tau_1 \dots \tau_n \rightarrow \tau} \Downarrow_{\Gamma} \lambda (x_1 : \tau_1 \dots x_n : \tau_n) e}$$

$$\frac{\text{GAPP} \quad \square_{\tau_1 \dots \tau_n \rightarrow \tau} \Downarrow_{\Gamma} e \quad \square_{\tau_1} \Downarrow_{\Gamma} e_1 \dots \square_{\tau_n} \Downarrow_{\Gamma} e_n}{\square_{\tau} \Downarrow_{\Gamma} e e_1 \dots e_n}$$

We have four basic choices: we can completely fill the hole with a correctly typed variable or constant; if we are trying to fill a hole of type $\tau_1 \dots \tau_n \rightarrow \tau$, then we can construct a λ term that takes *fresh* variables of types $\tau_1 \dots \tau_n$ whose body is recursively generated from a hole of type τ ; or we can create an application comprised of terms generated by a hole for the function and holes for the arguments for some arbitrary argument types $\tau_1 \dots \tau_n$.

These big-step rules are fairly straightforward to implement: we can pick a production rule, recursively satisfy its premises, and combine any results as needed. However, they are also quite coarse: in the application rule, we have to fully generate either the function or an argument before generating the rest. At this point it is natural to wonder if there is a small-step equivalent that allows for finer control. The answer, of course, is yes.

2.3 From big-step to small-step

The first step is to introduce *expression contexts*: the type of one-hole contexts C of expressions e .

$$\begin{array}{ll} C = e \dots C e \dots & \text{Func. or Arg. contexts} \\ | \lambda (x : \tau \dots) C & \lambda \text{ contexts} \\ | \diamond & \text{Context holes} \end{array}$$

Contexts C always end in a single *context hole* \diamond , which can be “filled in” by an expression, or another context. We write $C \blacklozenge e$ to mean filling in the context hole in C with the expression e . This operation is associative, so that $(C_1 \blacklozenge C_2) \blacklozenge e = C_1 \blacklozenge C_2 \blacklozenge e = C_1 \blacklozenge (C_2 \blacklozenge e)$. This notation is equivalent to the traditional “ $C_1[C_2[e]]$ ” – we use it so that the small-step generation rules throughout the paper are easier to parse.

Given a (well-typed) expression e , we can decompose it into C and e' such that $e = C \blacklozenge e'$ to *focus* on the subexpression e' within e , which lives inside the context C . Note that when we decompose the term in this way the context C provides an environment of variables whose scope encompasses e' – if e is a closed term, then any free variable in e' comes from a binding provided in C . We define $\Gamma(C)$ to be this environment.

Armed with the notion of contexts, we can describe generation as a sequence of small-step reductions that focus on a specific typed hole. Once again, each typing rule of λ_{\square} gives rise to a production rule:

$$\begin{array}{l}
\text{GENVAR} \\
C \blacklozenge \square_{\tau} \Longrightarrow C \blacklozenge x \\
\text{where } [x \mapsto \tau] \in \Gamma(C)
\end{array}
\qquad
\begin{array}{l}
\text{GENCONST} \\
C \blacklozenge \square_{\tau} \Longrightarrow C \blacklozenge c \\
\text{where constant } c \text{ has type } \tau
\end{array}$$

$$\begin{array}{l}
\text{GENABS} \\
C \blacklozenge \square_{\tau_1 \dots \tau_n \rightarrow \tau} \Longrightarrow C \blacklozenge (\lambda (x_1 : \tau_1 \dots x_n : \tau_n) \square_{\tau}) \\
\text{where } x_1 \dots x_n \text{ are } \textit{fresh}
\end{array}$$

$$\begin{array}{l}
\text{GENAPP} \\
C \blacklozenge \square_{\tau} \Longrightarrow C \blacklozenge (\square_{\tau_1 \dots \tau_n \rightarrow \tau} \square_{\tau_1} \dots \square_{\tau_n}) \\
\text{for any types } \tau_1 \dots \tau_n
\end{array}$$

It can be easily verified that these rules are type-preserving as they share the same one-to-one correspondence with the typing rules that exists with the big-step generation rules. The notable difference is that the production rule corresponding to application GENAPP now produces $n + 1$ typed holes, each of which we can choose to instantiate in any order. That allows us to *interleave* the generation of arguments and function expressions, which will prove crucial later on.

With this relation at hand, we can now formally define generation of a well-typed term e for a given type τ . We will say that an expression e is the result of such generation if e contains no holes, i.e. e cannot be stepped further, and can also be produced by the transitive closure of the step relation starting from a hole of type τ :

$$(\square_{\tau} \Longrightarrow^* e)$$

As well-typedness is preserved by the production rules and \square_{τ} is well-typed in an empty context, we can also conclude that e also has type τ in an empty context.

This small-step view of generation provides an alternative, but equally straightforward implementation: after choosing a transition we can fill in any new holes by recursive descent — we never have to look upwards into the context. Adding a decreasing *size* parameter to an implementation ensures termination: when the size becomes zero, remaining typed holes can be filled in with “trivial” terms such as `nil`, `1`, etc.. A hole with a function type can step to a λ with a body that is “trivial” in the same way, which structurally decreases the size of the type, even though it does not decrease the number of holes.

2.4 Extended Example

For concreteness, let’s step through an example generation of a list of integers using all production rules. This example is shown in Figure 1.

- (1) We begin with a typed hole of type `List Int` and an empty context. At this point, multiple rules apply: we could either fill the hole with a constant (such as the empty list) using GENCONST or we could expand the hole with an application using GENAPP. A generation strategy at this point would make a randomized choice according to some distribution.
- (2) For the sake of the example, let’s go with the more interesting choice of GENAPP. Immediately, however, we’re faced with yet another choice: how many arguments should application use and over which types? To keep things going, let’s generate a single-argument application over lists of integers.
- (3) We now have two holes — we need to focus on one to continue generation. Let’s focus on the left.
- (4) To expand the left hole that needs a function, we could choose to apply the GENABS rule, creating a new hole for the body.

$$\begin{array}{llll}
246 & & \diamond & \blacklozenge & \square_{\text{List Int}} & (1) \\
247 & & & \implies & \text{GENAPP} & \\
248 & & \diamond & \blacklozenge & \square_{\text{List Int} \rightarrow \text{List Int}} \square_{\text{List Int}} & (2) \\
249 & & & = & (\text{focus on left hole}) & \\
250 & & \diamond \square_{\text{List Int}} & \blacklozenge & \square_{\text{List Int} \rightarrow \text{List Int}} & (3) \\
251 & & & \implies & \text{GENABS} & \\
252 & & \diamond \square_{\text{List Int}} & \blacklozenge & \lambda (x : \text{List Int}) \square_{\text{List Int}} & (4) \\
253 & & & = & (\text{focus on inner hole}) & \\
254 & & (\lambda (x : \text{List Int}) \diamond) \square_{\text{List Int}} & \blacklozenge & \square_{\text{List Int}} & (5) \\
255 & & & \implies & \text{GENVAR} & \\
256 & & (\lambda (x : \text{List Int}) \diamond) \square_{\text{List Int}} & \blacklozenge & x & (6) \\
257 & & & = & (\text{refocus on remaining hole}) & \\
258 & & (\lambda (x : \text{List Int}) x) \diamond & \blacklozenge & \square_{\text{List Int}} & (7) \\
259 & & & \implies & \text{GENCONST} & \\
260 & & (\lambda (x : \text{List Int}) x) \diamond & \blacklozenge & \text{nil} & (8) \\
261 & & & = & (\text{remove focus}) & \\
262 & & \diamond & \blacklozenge & (\lambda (x : \text{List Int}) x) \text{nil} & (9) \\
263 & & & & & \\
264 & & & & & \\
265 & & & & & \\
266 & & & & & \\
267 & & & & & \\
268 & & & & & \\
269 & & & & & \\
270 & & & & & \\
271 & & & & & \\
272 & & & & & \\
273 & & & & & \\
274 & & & & & \\
275 & & & & & \\
276 & & & & & \\
277 & & & & & \\
278 & & & & & \\
279 & & & & & \\
280 & & & & & \\
281 & & & & & \\
282 & & & & & \\
283 & & & & & \\
284 & & & & & \\
285 & & & & & \\
286 & & & & & \\
287 & & & & & \\
288 & & & & & \\
289 & & & & & \\
290 & & & & & \\
291 & & & & & \\
292 & & & & & \\
293 & & & & & \\
294 & & & & &
\end{array}$$

Fig. 1. Example derivation of a generation

- (5) We then further focus on the body hole, transforming the context appropriately.
(6) We are once again faced with a typed hole of type List Int, as in the beginning, but this time the context contains a binding that is in scope: x. That means we can use GENVAR.
(7) After completely filling the functional part of the application, we refocus on the remaining hole: the argument.
(8) For the third occurrence of a List Int hole we can pick to fill it with GENCONST.
(9) Generation is now finished as no holes remain.

2.5 Limitations

This example renders the inherent limitation mentioned in the introduction glaringly obvious: to preserve the straightforward recursive nature of the algorithm, the types $\tau_1 \dots \tau_n$ need to be generated before focusing on any of the newly generated holes. Unfortunately, a poor choice of types during this production rule can significantly impact the performance and effectiveness of the generation. For example, what if at step (2) we had picked a different type, such as Int or Float? In the former case, generating a body that uses the abstracted variable meaningfully would require some luck (e.g. choosing to generate an application of a cons cell whose head can be the variable). In the latter case, using the argument is near impossible. Palka et al. (2011) attempt to address this problem by trying to make as educated a choice as possible during such productions.

The way they proposed to make a more educated choice was to introduce an additional production rule for application, that combined the GENAPP and GENVAR rules to make the use of existing functions (for example, primitives) more likely:

$$\begin{array}{l}
\text{GENFUNVAR} \\
C \blacklozenge \square_{\tau} \implies C \blacklozenge (f \square_{\tau'}) \\
\text{where } [f \mapsto \tau' \rightarrow \tau] \in \Gamma(C)
\end{array}$$

However, this rule still doesn't ensure that f will use its arguments and therefore doesn't completely solve the problem; it only makes it more likely f itself will be used.

While in theory the effect of this rule could be achieved with just the basic four, observing this behavior in practice would be quite unlikely, as the generation would have to go through a very particular set of guesses to form the right type. In contrast, we propose an alternative generation strategy that *postpones making such a choice* until a use for the function argument is encountered.

3 THE “NONLOCAL” GENERATOR

The essence of the problem is that the local algorithm needs to guess a type that will eventually appear on a typed hole. We introduce a way for the algorithm to delay this choice. At a high-level, when generating a function f we can leave its parameters *unspecified*, and introduce a rule that adds a parameter to f of some type. Composing this rule with one that fills a hole with a variable, f is able to fill some hole \square_τ with a variable x and add $x : \tau$ as a parameter. This rule introduces a non-locality in the generation process that the local approach does not exhibit and can only be expressed when adopting a small-step view of generation.

Such a method does not come without its complications. For example, the function f may impose additional nonlocal constraints within the generated fragment, such as calls to f or passing f as a higher-order argument. As a result, any changes to the type of f need to be propagated to the type of these constrained locations.

3.1 Extending λ_\square

To capture these locations, we introduce the notion of *arguments holes* denoted as “ \triangleright ”, which are labeled with globally scoped identifiers e.g. \triangleright^α . Informally we use these arguments holes to link applications with all the lambdas that could be called from that site and vice-versa. To do this we extend the syntax of λ_\square to add variants of each function type, application, and lambda parameter list that is annotated with an arguments hole. Finally, we modify the typing rules for applications and lambdas to require that the function type is annotated with the same arguments hole. This allows us to define an operation of “extending” an arguments hole \triangleright^α by extending every lambda, application, and function type annotated with the same \triangleright^α in the appropriate way. The supplemental material contains a Coq formalization of the following type system and a proof of soundness of arguments hole extension.

Formally, we first extend types to contain arguments holes, so that types are either a base type, a function type, or a function type whose list of argument types may be extended:

$$\begin{array}{l} \tau = \dots \\ | \tau \dots \triangleright^\alpha \rightarrow \tau \quad \text{Extendable Function Types} \end{array}$$

We also update the definitions of terms and contexts:

$$\begin{array}{l} e = \dots \\ | e e \dots \triangleright^\alpha \quad \text{Extendable Func. Applications} \\ | \lambda (x : \tau \dots \triangleright^\alpha) e \quad \text{Extendable } \lambda \text{ abstractions} \end{array}$$

Just like function types, both function applications and lambda abstractions are now extendable. Similarly, the contexts corresponding to function application and lambda abstraction are also now extendable.

$$\begin{array}{l} C = \dots \\ | e \dots C e \dots \triangleright^\alpha \quad \text{Extendable Func. or Arg. contexts} \\ | \lambda (x : \tau \dots \triangleright^\alpha) C \quad \text{Extendable } \lambda \text{ contexts} \end{array}$$

We call this language λ_\triangleright . Note that we use the same symbol \triangleright in three distinct settings: type argument lists, function application argument lists, and lambda parameter lists. When any of these

are updated for some \triangleright^α , all locations using the same \triangleright^α must be updated. We can prove this using the following updated typing relation, which simply checks that the arguments holes “line up” correctly.

$$\begin{array}{c} \text{TAPP}\triangleright \\ \hline \Gamma \vdash e_f : \tau_1 \dots \tau_n \triangleright^\alpha \rightarrow \tau \quad \Gamma \vdash e_1 : \tau_1 \dots \Gamma \vdash e_n : \tau_n \\ \hline \Gamma \vdash (e_f e_1 \dots e_n \triangleright^\alpha) : \tau \\ \\ \text{TABS}\triangleright \\ \hline \Gamma [x_1 \mapsto \tau_1 \dots x_n \mapsto \tau_n] \vdash e : \tau \\ \hline \Gamma \vdash \lambda (x_1 : \tau_1 \dots x_n : \tau_n \triangleright^\alpha) e : \tau_1 \dots \tau_n \triangleright^\alpha \rightarrow \tau \end{array}$$

A function application $e_f e_1 \dots e_n$ is extendable by an arguments hole \triangleright^α if the type of the function application e_f is also extendable by *the same* arguments hole. This restriction ensures that all lambdas that could be applied here have the same arguments hole, guaranteeing that extending \triangleright^α will not introduce any arity errors. On the flip side, an extendable lambda abstraction has an extendable function type with the same arguments hole, if we can type its body e in an environment extended by all the concrete bindings with appropriate types.

Finally, λ_\triangleright is truly an extension of λ_\square : we can convert any λ_\triangleright term to a λ_\square one by simply erasing all of the arguments holes in it — the resulting term is still well-typed. This means that we can guarantee that generators for λ_\triangleright terminate: we choose a fixed number of times to apply the λ_\triangleright generation rules, and after that many generation steps have been taken fill the remaining holes by erasing the λ_\triangleright term into an λ_\square one and using the method described in 2.3 to create “trivial” terms fitting their type. This method of ensuring termination, however, results in slightly larger variance in the size of generated terms, as steps like $\text{GENFUNVAR}\triangleright$ and $\text{GENPARAM}\triangleright$ introduce a variable number of holes, and the size of the “trivial” terms varies depending on their type’s complexity.

3.2 Nonlocal Generation

In the following section, we’ll explore the production rules involving arguments holes via extending the step relation (\implies). Each of these rules maintain well-typedness in λ_\triangleright , which we have also formalized in the Coq model, proving that they maintain the well-typedness of the program.

To begin, the rule that introduces a new \triangleright is the extendable function application.

$$\begin{array}{c} \text{GENAPP}\triangleright \\ C \blacklozenge \square_\tau \implies C \blacklozenge (\square_{\triangleright^\alpha \rightarrow \tau} \triangleright^\alpha) \\ \text{where } \triangleright^\alpha \text{ is fresh} \end{array}$$

Just like the previous application rule GENAPP , the new rule replaces a typed hole of type τ with an application of a typed hole whose result type is τ . However, unlike the previous application rule which placed no constraints on the type(s) of the argument, $\text{GENAPP}\triangleright$ delays the choice of the type(s), using a fresh arguments hole \triangleright^α in its place. Eventually, this arguments hole will be extended, and the argument type(s) will be “filled in”. At the same time, the argument to the application is an arguments hole of the same name, ensuring that the function arguments and the function type remain in sync. This rule satisfies $\text{TAPP}\triangleright$ trivially.

The dual to this rule is the generation of extendable lambda abstractions:

$$\begin{array}{c} \text{GENABS}\triangleright \\ C \blacklozenge \square_{\tau_1 \dots \tau_n \triangleright^\alpha \rightarrow \tau} \implies C \blacklozenge \lambda (x_1 : \tau_1 \dots x_n : \tau_n \triangleright^\alpha) \square_\tau \\ \text{where } x_1 \dots x_n \text{ are fresh} \end{array}$$

This rule is a natural extension of the GENABS rule to λ_\triangleright , adding multiple new bindings for the arguments already specified, and using the same arguments hole to ensure that any arguments

that are added in the future will also be added as parameters to this lambda abstraction. Note that GENAPP \triangleright followed immediately by GENABS \triangleright will result in no new variables being introduced for the moment.

Similarly to Pałka et al. (2011) we introduce a rule for creating an application of an existing function:

$$\begin{array}{l} \text{GENFUNVAR}\triangleright \\ C \blacklozenge \square_\tau \Longrightarrow C \blacklozenge f \square_{\tau_1} \dots \square_{\tau_n} \triangleright^\alpha \\ \text{where } [f \mapsto \tau_1 \dots \tau_n \triangleright^\alpha \rightarrow \tau] \in \Gamma(C) \end{array}$$

Unlike GENFUNVAR, this rule is **not** a combination of previous rules. Without this rule there would be no way to use variables of an extendable type since GENAPP \triangleright requires that the arguments hole is *fresh*, and TAPP \triangleright requires the same arguments hole for extendable applications to be well-typed.

Finally, we get to the crucial rule: generating a variable that does not appear concretely in the context, by extending an enclosing abstraction’s parameter list:

$$\begin{array}{l} \text{GENPARAM}\triangleright \\ \frac{C_1}{C_1 \triangleright^{\alpha + \tau}} \blacklozenge (\lambda (x_1 : \tau_1 \dots x_n : \tau_n \triangleright^\alpha) \blacklozenge) \blacklozenge C_2 \frac{}{C_2 \triangleright^{\alpha + \tau}} \blacklozenge \square_\tau \Longrightarrow \\ \frac{}{C_1 \triangleright^{\alpha + \tau}} \blacklozenge (\lambda (x_1 : \tau_1 \dots x_n : \tau_n y : \tau \triangleright^\alpha) \blacklozenge) \blacklozenge C_2 \triangleright^{\alpha + \tau} \blacklozenge y, \\ \text{where } y \text{ is fresh and } \triangleright^\alpha \notin \tau \end{array}$$

This rule is our first example of a nonlocal generation rule. First, we decompose the context of the hole \square_τ that we want to fill in three parts: an outer context C_1 , a lambda context that contains the parameter list (annotated with \triangleright^α) to be extended, and an inner context C_2 . The hole \square_τ is then filled with a fresh variable y ; this variable is then used to extend the bindings of the lambda abstraction from $x_1 : \tau_1 \dots x_n : \tau_n \triangleright^\alpha$ to $x_1 : \tau_1 \dots x_n : \tau_n y : \tau \triangleright^\alpha$. However, since the arguments hole \triangleright^α is being extended, all locations in the global program term that use the same arguments hole must also be updated. We require that none of τ contains \triangleright^α —otherwise, this update would never terminate as we would be creating a cyclic type. We write $\tau' \triangleright^{\alpha + \tau}$, $e \triangleright^{\alpha + \tau}$, and $C \triangleright^{\alpha + \tau}$ for applying the extension $\triangleright^{\alpha + \tau}$ of \triangleright^α by τ to types, expressions, and contexts, respectively. We will walk through each of these in detail below.

For concreteness, before formalizing this extension, let’s step through an example of applying the GENPARAM \triangleright rule. Assume that generation has reached the following expression — a lambda abstraction that binds f of extendable type $\triangleright^\alpha \rightarrow \text{Int}$ and calls it, applied to a lambda abstraction whose parameter list contains the same arguments hole, and whose body is a typed hole:

$$(\lambda (f : \triangleright^\alpha \rightarrow \text{Int}) (f \triangleright^\alpha)) (\lambda (\triangleright^\alpha) \square_{\text{Int}})$$

This expression can be decomposed as follows, using the application that binds and applies f as the outer context, and an empty inner context:

$$C_1 \blacklozenge (\lambda (\triangleright^\alpha) \blacklozenge) \blacklozenge C_2 \blacklozenge \square_{\text{Int}} \quad \text{where} \quad \begin{array}{l} C_1 = ((\lambda (f : \triangleright^\alpha \rightarrow \text{Int}) (f \triangleright^\alpha)) \blacklozenge) \\ C_2 = \blacklozenge \end{array}$$

We can then apply GENPARAM \triangleright to fill the hole with a fresh variable x , that we bind in the lambda immediately above it.

$$\frac{C_1}{C_1 \triangleright^{\alpha + \text{Int}}} \blacklozenge (\lambda (\triangleright^\alpha) \blacklozenge) \blacklozenge \frac{C_2}{C_2 \triangleright^{\alpha + \text{Int}}} \blacklozenge \square_{\text{Int}} \Longrightarrow [\text{GENPARAM}\triangleright] \\ \frac{}{C_1 \triangleright^{\alpha + \text{Int}}} \blacklozenge (\lambda (x : \text{Int} \triangleright^\alpha) \blacklozenge) \blacklozenge C_2 \triangleright^{\alpha + \text{Int}} \blacklozenge x$$

At this point, to preserve well-typedness we need to also expand the arguments hole \triangleright^α in C_1 (and, if such a hole existed, also in C_2), using the extension metafunction on C_1 . There are two occurrences of \triangleright^α in C_1 that need to be extended; the type of f and the application of f . For the

former case we extend the type to include the added parameter type, and for the latter we place a hole of type τ in the application. Performing those extensions gives us:

$$C_1 \overline{\triangleright^\alpha + \text{Int}} = ((\lambda (f : \text{Int} \triangleright^\alpha \rightarrow \text{Int}) (f \square_{\text{Int}} \triangleright^\alpha)) \diamond)$$

Since there are no occurrences of \triangleright^α in C_2 , we see that $C_2 \overline{\triangleright^\alpha + \text{Int}} = C_2 = \diamond$. Plugging the context holes gives us the resulting expression from the application of $\text{GENPARAM}\triangleright$:

$$\begin{aligned} & (\lambda (f : \triangleright^\alpha \rightarrow \text{Int}) (f \triangleright^\alpha)) (\lambda (\triangleright^\alpha) \square_{\text{Int}}) \quad \Longrightarrow [\text{GENPARAM}\triangleright] \\ & (\lambda (f : \text{Int} \triangleright^\alpha \rightarrow \text{Int}) (f \square_{\text{Int}} \triangleright^\alpha)) (\lambda (x : \text{Int} \triangleright^\alpha) x) \end{aligned}$$

Extension metafunctions, formally. The extension metafunctions need to be defined over types, expressions, and contexts. For types, the only interesting behavior occurs when an arguments hole appears in both a type and in its extension. In this case, the type of the extension is added to the arguments position of the type, just like we did for τ above.²

Types

$$\begin{aligned} (\tau_1 \dots \triangleright^\alpha \rightarrow \tau_2) \overline{\triangleright^\alpha + \tau} &= \tau_1 \overline{\triangleright^\alpha + \tau} \dots \tau \triangleright^\alpha \rightarrow \tau_2 \overline{\triangleright^\alpha + \tau} \\ (\tau_1 \dots \triangleright^\beta \rightarrow \tau_2) \overline{\triangleright^\alpha + \tau} &= \tau_1 \overline{\triangleright^\alpha + \tau} \dots \triangleright^\beta \rightarrow \tau_2 \overline{\triangleright^\alpha + \tau} \quad \text{where } \triangleright^\alpha \neq \triangleright^\beta \\ (\tau_1 \dots \rightarrow \tau_2) \overline{\triangleright^\alpha + \tau} &= \tau_1 \overline{\triangleright^\alpha + \tau} \dots \rightarrow \tau_2 \overline{\triangleright^\alpha + \tau} \\ \tau_b \overline{\triangleright^\alpha + \tau} &= \tau_b \end{aligned}$$

For expressions, the first interesting case is for extendable applications where the arguments hole matches that of the extension. In this case, a hole with the type of the extension is added onto the arguments of the application, just like when we added a typed hole argument to the application of f in the previous example. The second interesting case is for extendable lambdas whose arguments hole matches that of the extension. In this case we must extend the parameter list of the lambda with a new variable. This rule is what makes this method *almost* not “useless”: in this one case we are forced to add a variable which we cannot guarantee will be used, since its creation did not originate from the lambda’s body.

Expressions

$$\begin{aligned} (e_1 e \dots \triangleright^\alpha) \overline{\triangleright^\alpha + \tau} &= e_1 \overline{\triangleright^\alpha + \tau} e \overline{\triangleright^\alpha + \tau} \dots \square_\tau \triangleright^\alpha \\ (e_1 e \dots \triangleright^\beta) \overline{\triangleright^\alpha + \tau} &= e_1 \overline{\triangleright^\alpha + \tau} e \overline{\triangleright^\alpha + \tau} \dots \triangleright^\beta \quad \text{where } \triangleright^\alpha \neq \triangleright^\beta \\ (e_1 e \dots) \overline{\triangleright^\alpha + \tau} &= e_1 \overline{\triangleright^\alpha + \tau} e \overline{\triangleright^\alpha + \tau} \dots \\ (\lambda (x : \tau_1 \dots \triangleright^\alpha) e) \overline{\triangleright^\alpha + \tau} &= \lambda (x : \tau_1 \overline{\triangleright^\alpha + \tau} \dots \mathbf{y} : \tau \triangleright^\alpha) e \overline{\triangleright^\alpha + \tau} \quad \text{where } y \text{ is fresh} \\ (\lambda (x : \tau_1 \dots \triangleright^\beta) e) \overline{\triangleright^\alpha + \tau} &= \lambda (x : \tau_1 \overline{\triangleright^\alpha + \tau} \dots \triangleright^\beta) e \overline{\triangleright^\alpha + \tau} \quad \text{where } \triangleright^\alpha \neq \triangleright^\beta \\ (\lambda (x : \tau_1 \dots) e) \overline{\triangleright^\alpha + \tau} &= \lambda (x : \tau_1 \overline{\triangleright^\alpha + \tau} \dots) e \overline{\triangleright^\alpha + \tau} \\ x \overline{\triangleright^\alpha + \tau} &= x \\ \square_{\tau_1} \overline{\triangleright^\alpha + \tau} &= \square_{\tau_1 \overline{\triangleright^\alpha + \tau}} \end{aligned}$$

Finally, extending arguments holes in contexts behaves the same ways as in expressions.

Contexts

$$\begin{aligned} (e_1 \dots C e_2 \dots \triangleright^\alpha) \overline{\triangleright^\alpha + \tau} &= e_1 \overline{\triangleright^\alpha + \tau} \dots C \overline{\triangleright^\alpha + \tau} e_2 \overline{\triangleright^\alpha + \tau} \dots \square_\tau \triangleright^\alpha \\ (e_1 \dots C e_2 \dots \triangleright^\beta) \overline{\triangleright^\alpha + \tau} &= e_1 \overline{\triangleright^\alpha + \tau} \dots C \overline{\triangleright^\alpha + \tau} e_2 \overline{\triangleright^\alpha + \tau} \dots \triangleright^\beta \quad \text{where } \triangleright^\alpha \neq \triangleright^\beta \\ (e_1 \dots C e_2 \dots) \overline{\triangleright^\alpha + \tau} &= e_1 \overline{\triangleright^\alpha + \tau} \dots C \overline{\triangleright^\alpha + \tau} e_2 \overline{\triangleright^\alpha + \tau} \dots \\ (\lambda (x : \tau_1 \dots \triangleright^\alpha) C) \overline{\triangleright^\alpha + \tau} &= \lambda (x : \tau_1 \overline{\triangleright^\alpha + \tau} \dots \mathbf{y} : \tau \triangleright^\alpha) C \overline{\triangleright^\alpha + \tau} \quad \text{where } y \text{ is fresh} \\ (\lambda (x : \tau_1 \dots \triangleright^\beta) C) \overline{\triangleright^\alpha + \tau} &= \lambda (x : \tau_1 \overline{\triangleright^\alpha + \tau} \dots \triangleright^\beta) C \overline{\triangleright^\alpha + \tau} \quad \text{where } \triangleright^\alpha \neq \triangleright^\beta \\ (\lambda (x : \tau_1 \dots) C) \overline{\triangleright^\alpha + \tau} &= \lambda (x : \tau_1 \overline{\triangleright^\alpha + \tau} \dots) C \overline{\triangleright^\alpha + \tau} \\ \diamond \overline{\triangleright^\alpha + \tau} &= \diamond \end{aligned}$$

²For this and following definitions, we **emphasize** changes.

491 *Backwards compatibility.* A useful feature of this extended small-step relation is that it is fully
 492 backwards compatible with the top-down approach. That is, when choosing how to instantiate a
 493 typed hole, we can opt to use one of the λ_{\flat} -specific rules, or one of the productions rules of the
 494 previous section, with the choice being weighted by annotations under user control. Although
 495 fully replacing function and application generation with their λ_{\flat} variants would restrict the space
 496 of generated programs, this is not a problem in practice as these rules only give people crafting
 497 generators the flexibility to bias generation towards creating functions that use their arguments, a
 498 bias that cannot be expressed so straightforwardly in a top-down setting.

499 *Nonlocal let-insertion.* The next nonlocal rule we introduce is to create a fresh variable when
 500 needed, but to insert it as a let-binding rather than a function parameter. We use the isomorphism

$$501 \quad \text{let } x : \tau = e \text{ in } e' \cong (\lambda (x : \tau) e') e$$

502 and will use let rather than writing the application and lambda in the following. We can insert a
 503 let-binding by breaking the surrounding context of a hole into two pieces as before, and inserting a
 504 let between them.

$$505 \quad \text{GENLET} \\
 506 \quad C_1 \blacklozenge C_2 \blacklozenge \square_{\tau} \implies C_1 \blacklozenge (\text{let } x : \tau = \square_{\tau} \text{ in } \diamond) \blacklozenge C_2 \blacklozenge x \\
 507 \quad \text{where } x \text{ is } \textit{fresh}$$

508 This generation rule can be viewed as a combination of previously added rules, but not as a
 509 simple sequencing of choices. This rule is equivalent to a generation step of GENAPP and GENABS, if
 510 they had been chosen earlier in generation, before the generation of C_2 , along with a final GENVAR.
 511 What we've done is delayed the decision to insert this let-binding to the point where we've decided
 512 to use the binding.

513 3.3 More Types

514 *Polymorphism.* In the original work by Pałka et al. (2011), generated terms did not *produce*
 515 polymorphic functions. However, they were able to *use* polymorphic constants, such as seq and
 516 map, filling in type variables that are free in the codomain with randomly generated types in the
 517 GENFUNVAR cases for such variables. We can similarly use polymorphic constants in our approach.
 518 Still, *generating* polymorphic functions is seemingly out of reach.

519 Imagine how we might extend the top-down approach to polymorphism: we could extend the
 520 GENAPP case to instead generate a polymorphic function type whose codomain is compatible with
 521 the generation target type. However, when generating the type up front it is surprisingly easy to
 522 back the algorithm into a corner by attempting to generate a term of an uninhabited type. A pitfall
 523 that our framework provides a potential, though incomplete, way out from.

524 The idea is when generating an extensible function type, e.g. $\flat^{\alpha} \rightarrow \text{List Int}$, we can first *weaken*
 525 the output type to a polymorphic variable, making the type $\forall z, \flat^{\alpha} \rightarrow \text{List } z$, and then proceed
 526 with generation as normal.

527 To generate a polymorphic function we produce a lambda abstraction whose body produces the
 528 weakened polymorphic type, and in the expansion of that hole, new variables whose types contain
 529 the type variable z can be added to the parameter list as needed. To see how this guarantees no
 530 backtracking will be needed let's step through an example generation of the safe list head function,
 531 where a default value is provided for the empty list case. The example is shown in Figure 2.

532 (1) We begin with a type-abstraction Λ introducing the polymorphic type z , and an extensible
 533 function with \flat^{α} whose body is a hole of type z .

$$\begin{aligned}
540 \quad \Lambda z. \lambda (& \triangleright^\alpha \square_z & (1) \\
541 & \implies \text{GENFUNVAR} \\
542 \quad \Lambda z. \lambda (& \triangleright^\alpha \text{listCase } \square_{\text{List } z} \square_z \square_z (\text{List } z) \rightarrow z & (2) \\
543 & \implies \text{GENPARAM}\triangleright \\
544 \quad \Lambda z. \lambda (l : (\text{List } z) & \triangleright^\alpha \text{listCase } l \square_z \square_z (\text{List } z) \rightarrow z & (3) \\
545 & \implies \text{GENABS} \\
546 \quad \Lambda z. \lambda (l : (\text{List } z) & \triangleright^\alpha \text{listCase } l \square_z (\lambda (x : z \text{ xs} : (\text{List } z)) \square_z) & (4) \\
547 & \implies \text{GENVAR} \\
548 \quad \Lambda z. \lambda (l : (\text{List } z) & \triangleright^\alpha \text{listCase } l \square_z (\lambda (x : z \text{ xs} : (\text{List } z)) x) & (5) \\
549 & \implies \text{GENPARAM}\triangleright \\
550 \quad \Lambda z. \lambda (l : (\text{List } z) \ c : z & \triangleright^\alpha \text{listCase } l \ c (\lambda (x : z \text{ xs} : (\text{List } z)) x) & (6) \\
551 & \\
552 & \\
553 & \\
554 & \\
555 & \\
556 & \\
557 & \\
558 & \\
559 & \\
560 & \\
561 & \\
562 & \\
563 & \\
564 & \\
565 & \\
566 & \\
567 & \\
568 & \\
569 & \\
570 & \\
571 & \\
572 & \\
573 & \\
574 & \\
575 & \\
576 & \\
577 & \\
578 & \\
579 & \\
580 & \\
581 & \\
582 & \\
583 & \\
584 & \\
585 & \\
586 & \\
587 & \\
588 &
\end{aligned}$$

Fig. 2. Generation of a polymorphic safe list head function

- (2) From here we use the GENFUNVAR rule to create an application of the eliminator function for List, listCase. This gives three holes for the list being eliminated, a value for the empty case, and a function with arguments for both of the components of a cons.
- (3) First we focus on the list hole, replacing it with a parameter that is added at \triangleright^α .
- (4) Next we fill in the cons case with a function. This adds another hole of type z , can reference the newly introduced variables x and xs .
- (5) We can fill in this newest hole with a reference to x , leaving only the empty case of listCase.
- (6) This final step is asking us to create a value for a purely polymorphic type without any variables of that type in the context. From this point we could choose to create an application, but that would only be delaying the inevitable since that path would still eventually require filling in a hole of type z . So instead we take the only route out and use GENPARAM \triangleright to insert an argument for the base value at \triangleright^α , finishing the generation.

The resulting polymorphic type for the function is inhabited *by construction*, since the creation of the type occurs concurrently with the creation of a witness. This technique works to generate polymorphic terms when we focus generation on just the polymorphic function. However, it breaks down when the generation of the polymorphic term becomes non-trivially constrained with the rest of the program.

More precisely, crucial to the simplicity of our generation rules is the idea that the type extension operation can be applied globally — we can scan over the whole program and extend every occurrence of an arguments hole with a constant type. This is no longer the case when generating polymorphic functions because they may be applied with multiple different types: the extension operation now needs to consider how the polymorphic type is being instantiated locally. More concerningly, however, extending a polymorphic function can induce an extension in non-polymorphic functions that are unified with the polymorphic function arbitrarily far away in the program.

For instance, consider the following partially generated program:

$$\begin{aligned}
582 \quad & \text{let } f : \forall z. z \triangleright^\alpha \rightarrow \text{List } z = \Lambda z. \lambda (x : z \triangleright^\alpha) \text{ cons } x \square_{\text{List } z} \text{ in} \\
583 \quad & \text{let } g : \text{Int } \triangleright^\alpha \rightarrow \text{List } \text{Int} = \lambda (n : \text{Int } \triangleright^\alpha) \square_{\text{List } \text{Int}} \text{ in} \\
584 \quad & (\text{if } \dots \text{ then } f \ \text{Int} \ \text{else } g) \ 0 \triangleright^\alpha \\
585 &
\end{aligned}$$

Here the generation has defined two functions; f and g . f has a polymorphic type $\forall z. z \triangleright^\alpha \rightarrow \text{List } z$. At the moment, its body is partially filled in with a call to cons, adding its argument x onto an

expression hole. The function g has the same type but specialized to Int . Below the lets there is an application that applies either f or g to 0.

What happens when we extend the hole in f with a variable bound at \triangleright^α ? How do we figure out what type to extend the application with? With the rules as written we would be extending \triangleright^α with $\text{List } z$, but this isn't a well-formed type since z is outside its binding! We need to find some mapping from the free type variables to types that are well defined at the current location. Going back to the current example the mapping we want is $z \mapsto \text{List Int}$. Unfortunately finding this out requires peeking into the function position of the application and finding the constraint from g .

There is also an inverse problem: if g is extended with an argument of type List Int , the easiest thing to do is to add an argument with that type to f , which indeed has no issues type-wise. However, since f produces a polymorphic type, an argument of type List Int is likely too specific to be used within the generation of f . Instead, it might be wise to extend f with the weakened type $\text{List } z$. Although for arbitrary types there could be many such weakenings with no clear rule for what choices would be more desirable. For example if instead we consider an extension with $\text{Int} \times \text{Int}$ we could consider any of $z \times \text{Int}$, $\text{Int} \times z$, $z \times z$, $\forall y. z \times y$, etc. However, since all of these types are perfectly valid to fill in here it is not a serious issue to generating polymorphic functions.

While it seems like there should always be a way to deduce what a valid type substitution to be used at a given location would be, the fact that it requires a nontrivial search of the context adds concerns that such a process might not be practical or guaranteed to terminate, so further exploration of generating polymorphic functions is left to future work.

Lists and Matching. Extending our algorithm to a more type-rich setting, comes with another instance of the problem in the form of match-bound variables. For a concrete example, consider the type of lists: we can operate on them using match expressions of the form:

$$\text{match } e \text{ with } [] \Longrightarrow e \mid x :: xs \Longrightarrow e$$

To generate matches, we could use a typing rule derived local generation rule:

$$\begin{array}{c} \text{GENMATCHLOCAL} \\ C \blacklozenge \square_\tau \Longrightarrow C \blacklozenge (\text{match } \square_{\text{List } \tau'} \text{ with } [] \Longrightarrow \square_\tau \mid x :: xs \Longrightarrow \square_\tau) \\ \text{where } x, xs \text{ are fresh and } \tau' \text{ is some type} \end{array}$$

This rule will replace a hole \square_τ with three new ones: one for the list expression, one for the empty case, and one for the non-empty case. However, once again, this rule suffers from the problem of guessing an arbitrary type. This can be partially addressed by adding a rule analogous to GENFUNVAR:

$$\begin{array}{c} \text{GENMATCHVAR} \\ C \blacklozenge \square_\tau \Longrightarrow C \blacklozenge \text{match } x \text{ with } [] \Longrightarrow \square_\tau \mid y :: ys \Longrightarrow \square_\tau \\ \text{where } [x \mapsto \text{List } \tau] \in \Gamma(C) \text{ and } y, ys \text{ are fresh} \end{array}$$

But this suffers from the same limitation as that one: it requires the generation to have already just happened to produce a variable of the right type. Using the same guiding principle as the beginning of this Section, we can devise the following rule to avoid that, where we write $\tau(C_2 \blacklozenge \square_\tau)$ for the type of $C_2 \blacklozenge \square_\tau$:

$$\begin{array}{c} \text{GENMATCH} \\ C_1 \blacklozenge C_2 \blacklozenge \square_\tau \Longrightarrow C_1 \blacklozenge (\text{match } \square_{\text{List } \tau} \text{ with } [] \Longrightarrow C_2 \blacklozenge \square_\tau \mid x :: xs \Longrightarrow \diamond) \blacklozenge x \\ \text{where } x \text{ and } xs \text{ are fresh} \end{array}$$

We could define a similar rule that binds xs instead. This rule is very similar to GENLET; in that rule we observed that it would be equivalent to make the right guesses at an earlier point in generation. In this rule we have same relationship: we could have decided to use GENMATCHLOCAL

638 earlier with the right guesses for types. However, by delaying this decision we have ensured that
639 we don't need to make any guesses for types at all. These rules may seem odd at first in that they
640 are still introducing variable bindings that are not guaranteed to be used, only ensuring that one of
641 x or x_s is used depending on which match generation rule was chosen. We argue that this type of
642 variable non-usage is not as important as that of function parameter non-usage. That is because
643 the act of pattern matching against a value is itself a useful step of computation, with the control
644 flow of the program depending on the value being matched against.

645 However, we found that letting the generator create too many control branches can be problematic
646 as often only one control path will ever be taken by the generated program. This in effect means
647 the generation time spent on the other branches is wasted in a similar manner to time spent on the
648 generation of function arguments without uses.

649 3.4 Implementation

650 We implemented a generator in OCaml using the production rules discussed in Section 3.2, in
651 addition to GENCONST and GENVAR. A naïve implementation of the formal framework where, for
652 example, every time an arguments hole was extended we iterated over the whole program term,
653 would not be efficient. For efficiency purposes, our implementation leverages the following key
654 techniques:
655

656 *Expressions are mutable cells and encode their context:* Mutable expressions make the insertion
657 of parameters, arguments, etc. efficient because we can make a local change to an expression
658 without reconstructing the entire program term to account for the new changes. Additionally,
659 sub-expression terms contain a pointer “upwards” to their predecessor (i.e. encoding their context)
660 so that the global program term can be traversed in both directions — when creating a new variable,
661 we traverse upwards from the hole until we find the place we'd like to insert the binding.
662

663 *Tracking the available holes:* We maintain a list of the cells which are currently holes in a worklist.
664 This means that we never have to search the program term to find existing holes; we always know
665 every hole's location.
666

667 *Tracking the locations of extensions:* We have maps from arguments holes to all types, lambdas,
668 and applications that use that arguments hole. This means that the global operation of searching
669 the program term for the particular arguments hole being extended can be made efficient: we
670 simply modify the appropriate cells from the map.
671

672 *Rule selection logic:* We have implemented a framework for easily expressing and modifying
673 complex weight functions for choosing generation steps. This was in large part enabled by the
674 small-step approach to generation, since each step can produce the generated sub-terms as holes
675 rather than recursive calls. Our framework uses urns (Lampropoulos et al. 2017b) to efficiently
676 sample from each valid distinct generation step.
677

678 In addition to the performance-oriented techniques above, we also focused our attention towards
679 low-level choices that impact the posterior distribution of generated terms:

680 *Functions with empty domains:* A key difference between our generator and Pałka et al. (2011) is
681 that the terms produced by our generator are in the multary lambda calculus. One of the caveats
682 of this means that our generator risks generating functions that never get a parameter added by
683 generation, and so are left with no parameters. While this is not necessarily an issue, if too many
684 functions turn out this way generated programs will primarily be thunking and forcing expressions
685 rather than passing data around in complex ways. In practice, however, we found that weighting
686

the insertion of arguments into functions with fewer parameters more heavily was sufficient to ensure that almost all functions have parameters.

Extraction to plain lambda calculus. Because we are generating in the multary lambda calculus, we also need to devise a method of extraction to test languages with only single arity. The easiest choice is to encode parameters as tuples, making all function calls uncurried. Another choice is to extract to curried functions, treating empty parameters as thunks.

No direct control over function domains. When generating a function of type $\triangleright^\alpha \rightarrow \tau$, we’re not able to place constraints on how complicated the domain will get — the body and non-local changes to the program term decide what the domain looks like. Even crude measures such as preventing GENPARAM \triangleright from being chosen if the type is too complex fall short when there are a lot of higher-order usages of extensible functions. This can become an issue for generation because it is often the case that the generator will run out of fuel before filling in holes for these complex types, resulting in many complex function types that just throw out all their values - exactly what we were trying to avoid! In practice, the same distribution tuning as before—deprioritizing insertions into functions with too many parameters—greatly lessens this tendency.

4 EVALUATION

As discussed earlier in the paper, this project began when we tried to use random testing to evaluate the correctness of student-written compilers for the needs of an undergraduate compilers course. We reimplemented the state-of-the-art approach to generating well-typed lambda terms (Palka et al. 2011), and we soon realized that we were having trouble consistently finding certain bugs dealing with intricate function argument usage—such as miscompilations that led to stack clobbering. Upon further inspection, we noticed that our top-down generator was rarely generating functions that actually *used* their arguments, which pointed us towards the nonlocal generation approach.

Naturally, to *evaluate* our approach we will not pit our implementation against our own re-implementation of Palka et al., but rather turn to their original implementation. In particular, we aim to answer the following three questions:

- (1) Is the new generator more effective at finding bugs in compilers?
- (2) Does our generator succeed in its goal of generating well-typed terms that use their arguments?
- (3) Does that bring it closer in behavior to real human-written programs?

4.1 Replicating the GHC Strictness Analyzer Case Study

To answer the first question, we replicate to the major case study of the original Palka et al. paper: finding bugs in the strictness analyzer of GHC-6.12. Palka et al. generated programs of type `List Int → List Int` in an environment seeded with multiple functions from the Haskell standard library, pretty-printed them all to create a single Haskell module, and applied them all to a sequence of partial lists of the form:

```
1:2:undefined
```

GHC-6.12 contained a series of bugs in its strictness analyzer, where the optimizer would sometimes optimize away an expression with a visible side effect. For example, the following term:

```
seq (id (\a -> seq a id) (undefined::Int))
```

outputs the following when evaluated on the aforementioned list:

```
*** Exception:
```

736 However, the forcing of the undefined is optimized away when compiled with optimizations
737 leading to a different output:

738 [1,2,*** Exception:
739

740 As a result, such bugs can be identified using *differential testing*: compiling the same program with
741 and without optimizations and checking that the outputs agree.

742 To evaluate the performance of our method we reuse the (Palka et al. 2011) test harness and
743 setup, substituting our implementation to populate the Haskell module. We run both our method
744 and that of Palka et al. 100 times, with each run consisting of 50 tests, and with each test consisting
745 of generating 1000 functions. For each batch of 1000 functions, we compile twice — once with and
746 once without optimizations, run the compiled output, and check that outputs agree.

747 We found that our method found a bug with a bit over half the number of generated examples,
748 and about a quarter of the CPU time with a comparable failure rate. The Palka et al. generation
749 code found bugs in 98 of the runs and took an average of 19.56 tests and 1092 seconds of CPU
750 time to find a bug. Our method found bugs in 99 runs and took an average of 10.08 tests and 237.4
751 seconds to find a bug. That is, our method was roughly $2\times$ more *effective* at finding bugs and $4\times$
752 more *efficient* on average. A Mann-Whitney-U test (Mann and Whitney 1947) confirms that our
753 method is statistically better, as the p-value that our method is finding bugs in less tests by chance
754 is $p = 3.014 \cdot 10^{-13}$, and in less CPU time by chance is $p = 1.147 \cdot 10^{-28}$.

755 An extension of the original work (Palka 2014) by Palka et al. describes three other distinct bugs
756 discovered in GHC, which are triggered and detected by using slightly modified generators and
757 test harnesses. Out of the three, we could manually trigger two using the counterexample provided
758 in the thesis, and our generators were also able to find them—although we could not obtain the
759 code accompanying the thesis to compare against.

760 4.2 Parameter usage

761 To answer the last two questions, we introduce a straightforward usage metric on programs: we
762 simply calculate the percentage of function parameters that are used by their bodies. In addition, we
763 turn to an independent benchmark suite of functional programs (Quiring et al. 2022), and calculate
764 this metric on those, to get an estimate of what things looks like in practice.

765 The average parameter usage across all programs in the suite is 94.9%, which validates our
766 premise that functions should use their arguments. Table 1 provides a detailed description of the
767 programs in the benchmark suite and their parameter usage. Finally, we measure the usage rate of
768 both the local and the nonlocal approach. Across 10000 random generation tests, the generator
769 of Palka et al. produced an average of 30% parameter usage. Our approach when generating terms
770 in a similar configuration produces 30% usage of plain lambda parameters, 66% usage of extensible
771 lambda parameters, and 100% usage of let bindings, as our generator only produces those with a
772 corresponding reference. Overall there is a 35% usage amongst all function parameters and 55%
773 usage of all variables. However, in our generator framework all of these statistics are highly tunable:
774 different choices for the weights of rules such as GENAPP or GENPARAM directly influences the
775 percentage of variables used in generated terms, enabling a wide range of previously unattainable
776 behaviors. In our experimenting the usage of extensible lambda parameters has ranged from $\sim 60\%$
777 to $\sim 99\%$.

780 5 RELATED WORK

781 We've heavily discussed the pioneering work of Palka et al. (2011) throughout this paper, as they
782 first introduced the local algorithm and its various refinements. Since then, a number of other
783

Tab. 1. Benchmark descriptions and usage statistics

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

Name	LOC	# Used / Total	Description
nucleic	3326	158 / 165 (95%)	The SML version of the “Pseudoknot” program (Hartel et al. 1996).
boyer	838	26 / 26 (100%)	The Boyer-Moore benchmark from SML/NJ.
ratio-regions	548	108 / 117 (92%)	An image segmentation benchmark from MLton.
mc-ray	487	108 / 121 (89%)	A SML port of the “Weekend Raytracer” (Shirley 2020).
knuth-bendix	450	162 / 166 (97%)	The Knuth-Bendix benchmark from SML/NJ.
raytracer	333	102 / 102 (100%)	The Ray tracer from Impala benchmarks.
SNF	290	50 / 52 (96%)	The Smith-Normal-Form benchmark from MLton.
cps-convert	199	29 / 32 (90%)	A SML implementation of the Danvy-Filinski CPS conversion (Danvy and Filinski 1992).
interpreter	178	25 / 25 (100%)	An interpreter for a simple language.
parser-comb	168	70 / 72 (97%)	Using parser combinators to parse a simple expression syntax.
life	122	56 / 57 (98%)	The Life benchmark from SML/NJ.
derivative	113	20 / 21 (95%)	A SML port of the symbolic derivation benchmark from Larceny.
nqueens	45	10 / 10 (100%)	The classic N-Queens benchmark.
quicksort	44	4 / 4 (100%)	Quicksort of integer lists.
mandelbrot	43	6 / 6 (100%)	The Mandelbrot benchmark from SML/NJ
safe-for-space	27	4 / 7 (57%)	An example that tests safe-for-space closure conversion (Shao and Appel 2000)
cpstak	21	6 / 6 (100%)	The continuation-passing-style implementation of the Takeuchi (tak) function ported from Larceny.
filter	11	6 / 6 (100%)	Filter a list.
tak	10	1 / 1 (100%)	The Takeuchi (tak) function.
ack	8	1 / 1 (100%)	Ackermann’s function.

papers have tackled the same problem of generating well-typed terms, with varying approaches, goals, and outcomes.

“Local” Approaches. Midtgaard et al. (2017) built upon the local algorithm of Palka et al. by extending the type system with an effect for identifying programs whose output depends on the evaluation order. They then generalized the local algorithm to take effects into account, generating well-typed programs without such a dependence. Thus, they were able to weed out undefined behaviors that arise from under-specification of the evaluation order and find multiple bugs in OCaml’s optimizing native-code back end. Rocha (2019) further built upon this work by encoding exceptions as an effect in the type system, and using the same approach to find additional bugs in various OCaml compilers. Extending a type system with effects is largely orthogonal to generating functions that use their arguments, so we would expect such approaches to also benefit from the results presented in this paper.

In the same vein, Hoang et al. (2022) extend the algorithm of Palka et al. to System F. In System F the problem of early generation of types is only exacerbated: it also contains *type* applications! Hoang et al. introduce the notion of *unsubstitution*, a process for generating some types τ , τ' and some type variable α for a given type σ such that $\sigma = \tau[\tau'/\alpha]$. We discussed how our approach

834 interacts with polymorphism in Section 3.3, but it would be interesting future work to explore
835 precisely how it interacts with the particular intricacies of System F.

836 The local algorithm has also been adapted to an imperative language setting. da Silva Feitosa
837 et al. (2019) use such an approach to generate well-typed Featherweight Java (Igarashi et al. 2001)
838 programs, while more recently Li et al. (2022) did the same for their core calculus for Checked
839 C (Ruef et al. 2019). Both approaches still first target a model of the imperative language that is
840 embedded in a functional one, and therefore the ideas from Palka et al. carry over seamlessly. On
841 the other hand, the CSmith project, the crown jewel of imperative compiler testing (Yang et al.
842 2011), follows a different, more ad-hoc, approach developing a fine-tuned hand-written generator to
843 generate C programs that don't exhibit undefined behavior. Interestingly, it appears that their work
844 also suffers from a similar problem as the one we identified and addressed in this paper, as Barany
845 (2018) showed that one can leverage static analysis during generation to eliminate some forms of
846 dead code, leading to improved efficiency and effectiveness.

847
848 *Uniform and Unification-Based Approaches.* A different line of work focuses on studying the
849 statistical properties of lambda terms rather than generating terms for testing purposes. By viewing
850 lambda terms as labeled trees, the well-explored literature on generating combinatorial structures
851 is directly applicable, beginning with Flajolet et al. (1994) who systematized the notion of recursive
852 grammar-based generation almost three decades ago. Since then, a number of approaches have
853 directly used this prism for lambda term generation. Grygiel and Lescanne (2012) provide a mapping
854 from (untyped) lambda terms to natural numbers (called a *ranking*) and an inverse mapping (the
855 *unranking*). This in turn gives rise to a natural generation strategy that first samples a natural
856 number before using the unranking function to uniformly generate terms of an exact size. On
857 the other hand, Lescanne (2014) used the notion of Boltzmann samplers (Duchon et al. 2004) to
858 uniformly generate terms of an approximate size. Unfortunately, while using such combinatorial
859 techniques for generating and studying the characteristics of untyped lambda terms has been
860 fruitful, they appear extremely hard to adapt so that they enforce more semantic constraints such
861 as well-typedness, despite some progress in adapting them to take other structural features into
862 account (Feldt and Poulding 2013).

863 On the other hand, the work of Kennedy and Vytiniotis (2012) also uses a ranking-unranking
864 like approach, elegantly phrased in the form of games, but they succeed in providing encoding and
865 decoding functions directly from and to well-typed lambda terms. While testing is mentioned as a
866 potential application of their functional pearl, we are not aware of any work that has followed up
867 to demonstrate its viability.

868 An alternative line of work harnesses the declarative power of logic programming to concisely
869 enumerate lambda terms (Tarau 2015). In his work, Tarau describes a Prolog-based generator and
870 type inference algorithm for simply-typed terms in DeBruijn notation. Taking full advantage of
871 Prolog's unification and backtracking capabilities, he is able to enumerate lambda terms of slightly
872 larger sizes compared to a rejection sampling approach backed by combinatorial techniques for
873 generation of untyped terms. A follow up to this work bridges the gap, exploring a synergy between
874 logic programming and Boltzmann samplers (Bendkowski et al. 2017). By interleaving the two
875 approaches, the authors are able to achieve uniform generation of closed λ calculus terms of a size
876 an order of magnitude higher than before, and use it to gauge statistical properties of lambda terms
877 at scale.

878 Such approaches based on unification bear many similarities to the extensible function types we
879 described in this paper—after all, one way to view the extensible parameter lists is as unification
880 variables for lists of types which get resolved later. However, the key difference lies in when and
881 how these unification variables are resolved. Tarau resolves them only after the body has been
882

883 generated, resulting in many terms for which generation fails and backtracks until it produces a
884 successful term. Instead, we apply the unification steps during generation to allow the variables to
885 be resolved on their first usage, so that later usages are guaranteed to be well-typed and not cause
886 backtracking.

887 It still, however, remains an open question whether testing with a uniform generation of inputs
888 (and the associated overhead that comes with) is an effective approach. For example, Grygiel and
889 Lescanne (2012) observed that “almost all typeable terms start with an abstraction”. Unfortunately,
890 a generator that almost always produces such abstractions is completely ineffective at finding
891 bugs in properties that involve reduction, such as progress, preservation, or most specifications
892 of optimization correctness, as no reductions can actually take place. At the same time, ensuring
893 uniformity comes at a significant efficiency cost: Bendkowski et al. (2017) can generate terms of a
894 given moderate size at a rate of a few lambda terms per second; in contrast, our implementation
895 generates many thousands well-typed terms in the same time frame. While throughput is obviously
896 not the end-all metric for test effectiveness, this kind of difference in efficiency is almost a non-starter
897 for testing a complex software system such as GHC.

898 *Automatic Approaches.* An interesting thread of related work attempts to automatically derive
899 an algorithm such as the one that appears in Palka et al. (2011), based on a declarative specification
900 of the type system. Claessen et al. (2015) use an idea similar to the ranking-unranking approach
901 of Grygiel and Lescanne (2012), as implemented in the Haskell FEAT library (Duregård et al. 2012)
902 for enumeration of algebraic datatypes. They then rely on Haskell’s laziness to prune large parts of
903 the input space and avoid generating them in the future. Fetscher et al. (2015) extend this work
904 for PLT Redex, introducing a disunification solver to more effectively handle pattern matching
905 constraints. Lampropoulos et al. (2017a) develop a language for writing programs that can be
906 viewed as both generators and checkers, allowing for lightweight annotations to control both
907 the resulting probability distribution and the amount of constraint solving that happens prior to
908 variable instantiation. All three papers replicate the case study finding bugs in GHC’s strictness
909 analyzer, and report an order of magnitude overhead compared to the original implementation.
910 Finally, the case study of Palka et al. was also replicated in a recent work by Goldstein et al. (2021),
911 where the authors improved upon the original’s efficiency and effectiveness by generating multiple
912 inputs and selecting the ones that maximizes the interactions between different constructors—an
913 idea inspired by the combinatorial testing literature. Once again, at a high level these threads are
914 all orthogonal to our contributions.
915

916 *Program Synthesis.* Finally, we would be remiss to not discuss the related field of program
917 synthesis (Gulwani et al. 2017), which also deals with the problem of producing well-typed programs.
918 In particular, there is a sequence of approaches that attempt to synthesize functions that adhere
919 to a specification following a top-down enumeration approach—not too dissimilar from the Palka
920 et al. one (Feser et al. 2015; Osera and Zdancewic 2015; Polikarpova et al. 2016). There is, however,
921 a key distinction: the goal in program synthesis is to find *the one* or one of a few handful programs
922 that satisfies a given specification. In other words, the desired solution is over constrained, and
923 randomly producing that is extremely unlikely. Thinking back to the example derivation of the
924 safe head function in Section 3, there were a lot of random choices that had to line up for that
925 outcome. As a result, program synthesis turns to exhaustive techniques (such as SMT solvers) to
926 look for solutions. In contrast, the goal in random testing is to generate many programs in the
927 hopes of finding bugs in a space that can’t be exhaustively explored. This distinction makes all the
928 difference in the world. In the particular case of lambda terms, generation of well-typed terms is
929 an interesting middle ground between testing properties of syntactically correct terms, such as
930 a parser/printer roundtrip property where grammar-based generation will suffice, and program
931

932 synthesis, where exhaustive techniques need to be employed. The well-typedness precondition
 933 is sparse enough to make a generate-and-test approach impractical, but also contains sufficient
 934 inhabitants that randomness is necessary for effective exploration.

935 6 CONCLUSION AND FUTURE WORK

936 In this paper, we identified a key limitation in the type-directed approaches to well-typed term
 937 generation that are prevalent in the literature, which frequently produces functions that do not
 938 use their arguments, which causes large chunks of generated code to be ineffective for testing.
 939 We devised a new approach that produces functions that use their arguments by allowing the
 940 generation to delay the choice of argument types until it is generating uses for those arguments
 941 as well. We formalized this approach as a small-step relation in an extension of the simply-typed
 942 lambda calculus that is fully backwards compatible with the top-down approach, meaning that the
 943 heuristics developed in the past decade of work on the generation of well-typed lambda terms can
 944 still be exploited in this setting. In the future, we're hoping to further explore the vast configuration
 945 space that arises from the fine-grained level of control that the new algorithm provides, with an
 946 eye out for more complex type systems including dependently typed ones.

947 ACKNOWLEDGMENTS

948 We thank Jacob Prinz, Antal Spector-Zabusky, Chris Casinghino, and the anonymous reviewers (of
 949 both the paper and the artifact) for their helpful comments. This work was supported by NSF award
 950 #2107206, Efficient and Trustworthy Proof Engineering, and NSF award #2145649, CAREER: Fuzzing
 951 Formal Specifications (any opinions, findings and conclusions or recommendations expressed in
 952 this material are those of the authors and do not necessarily reflect the views of the NSF).

953 REFERENCES

- 954 Gergő Barany. 2018. Liveness-Driven Random Program Generation. In *Logic-Based Program Synthesis and Transformation*,
 955 Fabio Fioravanti and John P. Gallagher (Eds.). Springer International Publishing, Cham, 112–127.
- 956 Maciej Bendkowski, Katarzyna Grygiel, and Paul Tarau. 2017. Boltzmann Samplers for Closed Simply-Typed Lambda Terms.
 957 In *Practical Aspects of Declarative Languages*, Yuliya Lierler and Walid Taha (Eds.). Springer International Publishing,
 958 Cham, 120–135.
- 959 Koen Claessen, Jonas Duregård, and Michal H. Palka. 2015. Generating constrained random data with uniform distribution.
 960 *J. Funct. Program.* 25 (2015). <https://doi.org/10.1017/S0956796815000143>
- 961 Samuel da Silva Feitosa, Rodrigo Geraldo Ribeiro, and Andre Rauber Du Bois. 2019. Generating Random Well-Typed
 962 Featherweight Java Programs Using QuickCheck. *Electronic Notes in Theoretical Computer Science* 342 (2019), 3–20. <https://doi.org/10.1016/j.entcs.2019.04.002>
- 963 The proceedings of CLEI 2018, the XLIV Latin American Computing Conference.
- 964 Olivier Danvy and Andrzej Filinski. 1992. Representing Control: A study of the CPS transformation. *Mathematical Structures
 965 in Computer Science* 2, 4 (Dec. 1992), 361–391. <https://doi.org/10.1017/S0960129500001535>
- 966 Philippe Duchon, Philippe Flajolet, Guy Louchard, and Gilles Schaeffer. 2004. Boltzmann Samplers for the Random
 967 Generation of Combinatorial Structures. *Combinatorics, Probability & Computing* 13, 4-5 (2004), 577–625. <https://doi.org/10.1017/S0963548304006315>
- 968 Jonas Duregård, Patrik Jansson, and Meng Wang. 2012. Feat: Functional Enumeration of Algebraic Types. In *Proceedings of
 969 the 2012 Haskell Symposium* (Copenhagen, Denmark) (*Haskell '12*). ACM, New York, NY, USA, 61–72. <https://doi.org/10.1145/2364506.2364515>
- 970 Robert Feldt and Simon M. Poulding. 2013. Finding test data with specific properties via metaheuristic search. In *24th
 971 International Symposium on Software Reliability Engineering*. IEEE, 350–359. http://robertfeldt.net/publications/feldt_poulding_2013_finding_test_data_with_specific_properties.pdf
- 972 John K. Feser, Swarat Chaudhuri, and Isil Dillig. 2015. Synthesizing Data Structure Transformations from Input-output
 973 Examples. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*
 974 (Portland, OR, USA) (*PLDI '15*). ACM, New York, NY, USA, 229–239. <https://doi.org/10.1145/2737924.2737977>
- 975 Burke Fetscher, Koen Claessen, Michal H. Palka, John Hughes, and Robert Bruce Findler. 2015. Making Random Judgments:
 976 Automatically Generating Well-Typed Terms from the Definition of a Type-System. In *24th European Symposium on
 977*

- 981 *Programming (Lecture Notes in Computer Science, Vol. 9032)*. Springer, 383–405. <http://users.eecs.northwestern.edu/~baf11/random-judgments/>
- 982 Philippe Flajolet, Paul Zimmermann, and Bernard Van Cutsem. 1994. A Calculus for the Random Generation of Labelled
- 983 Combinatorial Structures. *Theor. Comput. Sci.* 132, 2 (1994), 1–35. [https://doi.org/10.1016/0304-3975\(94\)90226-7](https://doi.org/10.1016/0304-3975(94)90226-7)
- 984 Harrison Goldstein, John Hughes, Leonidas Lampropoulos, and Benjamin C. Pierce. 2021. Do Judge a Test by its Cover:
- 985 Combining Combinatorial and Property-Based Testing.
- 986 Katarzyna Grygiel and Pierre Lescanne. 2012. Counting and generating lambda terms. *CoRR* abs/1210.2610 (2012).
987 arXiv:1210.2610 <http://arxiv.org/abs/1210.2610>
- 988 Sumit Gulwani, Alex Polozov, and Rishabh Singh. 2017. *Program Synthesis*. Vol. 4. NOW, 1–119 pages.
- 989 Pieter H. Hartel, Marc Feeley, Martin Alt, Lennart Augustsson, and et al. 1996. Benchmarking implementations of functional
- 990 languages with ‘Pseudoknot’, a float-intensive benchmark. *Journal of Functional Programming* 6, 4 (1996), 621–655.
991 <https://doi.org/10.1017/S095679680001891>
- 992 Tram Hoang, Anton Trunov, Leonidas Lampropoulos, and Ilya Sergey. 2022. Random Testing of a Higher-Order Blockchain
- 993 Language (Experience Report). In *Proceedings of the ACM SIGPLAN International Conference on Functional Programming*
- 994 (ICFP).
- 995 Atsushi Igarashi, Benjamin C. Pierce, and Philip Wadler. 2001. Featherweight Java: A Minimal Core Calculus for Java and
- 996 GJ. *ACM Trans. Program. Lang. Syst.* 23, 3 (may 2001), 396–450. <https://doi.org/10.1145/503502.503505>
- 997 Andrew J. Kennedy and Dimitrios Vytiniotis. 2012. Every bit counts: The binary representation of typed data and programs.
998 *Journal of Functional Programming* 22, 4-5 (2012), 529–573. <https://research.microsoft.com/en-us/people/dimitris/jfpcoding.pdf>
- 999 Leonidas Lampropoulos, Diane Gallois-Wong, Catalin Hritcu, John Hughes, Benjamin C. Pierce, and Li-yao Xia. 2017a.
- 1000 Beginner’s Luck: a language for property-based generators. In *Proceedings of the 44th ACM SIGPLAN Symposium on*
- 1001 *Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017*. 114–129. <http://dl.acm.org/citation.cfm?id=3009868>
- 1002 Leonidas Lampropoulos, Antal Spector-Zabusky, and Kenneth Foner. 2017b. Ode on a Random Urn (Functional Pearl). In
- 1003 *Proceedings of the 10th ACM SIGPLAN International Symposium on Haskell (Oxford, UK) (Haskell 2017)*. ACM, New York,
- 1004 NY, USA, 26–37. <https://doi.org/10.1145/3122955.3122959>
- 1005 Pierre Lescanne. 2014. Boltzmann samplers for random generation of lambda terms. *CoRR* abs/1404.3875 (2014). <http://arxiv.org/abs/1404.3875>
- 1006 Liyi Li, Yiyun Liu, Deena L. Postol, Leonidas Lampropoulos, David Van Horn, and Michael Hicks. 2022. A Formal Model of
- 1007 Checked C. In *Proceedings of the Computer Security Foundations Symposium (CSF)*.
- 1008 H. B. Mann and D. R. Whitney. 1947. On a Test of Whether one of Two Random Variables is Stochastically Larger than the
- 1009 Other. *The Annals of Mathematical Statistics* 18, 1 (1947), 50 – 60. <https://doi.org/10.1214/aoms/1177730491>
- 1010 Jan Midtgaard, Mathias Nygaard Justesen, Patrick Kasting, Flemming Nielson, and Hanne Riis Nielson. 2017. Effect-
- 1011 Driven QuickChecking of Compilers. *Proc. ACM Program. Lang.* 1, ICFP, Article 15 (aug 2017), 23 pages. <https://doi.org/10.1145/3110259>
- 1012 Peter-Michael Osera and Steve Zdancewic. 2015. Type-and-example-directed Program Synthesis. In *Proceedings of the 36th*
- 1013 *ACM SIGPLAN Conference on Programming Language Design and Implementation (Portland, OR, USA) (PLDI ’15)*. ACM,
- 1014 New York, NY, USA, 619–630. <https://doi.org/10.1145/2737924.2738007>
- 1015 Michal H. Palka. 2014. *Random Structured Test Data Generation for Black-Box Testing*. Department of Computer Science and
- 1016 Engineering, Software Technology (Chalmers), Chalmers University of Technology,. <http://publications.lib.chalmers.se/records/fulltext/195849/195849.pdf> 168.
- 1017 Michal H. Palka, Koen Claessen, Alejandro Russo, and John Hughes. 2011. Testing an Optimising Compiler by Generating
- 1018 Random Lambda Terms. In *Proceedings of the 6th International Workshop on Automation of Software Test (Waikiki,*
- 1019 *Honolulu, HI, USA) (AST ’11)*. ACM, New York, NY, USA, 91–97. <https://doi.org/10.1145/1982595.1982615>
- 1020 Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program Synthesis from Polymorphic Refinement Types.
1021 *SIGPLAN Not.* 51, 6 (June 2016), 522–538. <https://doi.org/10.1145/2980983.2908093>
- 1022 Benjamin Quiring, John Reppy, and Olin Shivers. 2022. Analyzing Binding Extent in 3CPS. *Proc. ACM Program. Lang.* 6,
- 1023 ICFP, Article 114 (aug 2022), 29 pages. <https://doi.org/10.1145/3547645>
- 1024 Jason S. Reich, Matthew Naylor, and Colin Runciman. 2012. Advances in Lazy SmallCheck. Presented at the 24th
- 1025 Symposium on Implementation and Application of Functional Languages. <http://www.cs.york.ac.uk/fp/jason-docs/ReichNaylorRunciman2013.pdf>
- 1026 Murilo Giacometti Rocha. 2019. *Testing of OCaml exceptions by effect-driven generation of programs*. Master’s thesis.
- 1027 University of Edinburgh. https://project-archive.inf.ed.ac.uk/msc/20193481/msc_proj.pdf
- 1028 Andrew Ruef, Leonidas Lampropoulos, Ian Sweet, David Tarditi, and Michael Hicks. 2019. Achieving Safety Incrementally
- 1029 with Checked C. In *Proceedings of the Symposium on Principles of Security and Trust (POST)*.

- 1030 Colin Runciman, Matthew Naylor, and Fredrik Lindblad. 2008. SmallCheck and Lazy SmallCheck: automatic exhaustive
1031 testing for small values. In *1st ACM SIGPLAN Symposium on Haskell*. ACM, 37–48. [http://www.cs.york.ac.uk/fp/
1032 smallcheck/smallcheck.pdf](http://www.cs.york.ac.uk/fp/smallcheck/smallcheck.pdf)
- 1033 Zhong Shao and Andrew W. Appel. 2000. Efficient and Safe-for-Space Closure Conversion. *ACM Transactions on Programming
1034 Languages and Systems* 22, 1 (Jan. 2000), 129–161. <https://doi.org/10.1145/345099.345125>
- 1035 Peter Shirley. 2020. *Ray Tracing in One Weekend*. <https://raytracing.github.io>
- 1036 Paul Tarau. 2015. On Type-directed Generation of Lambda Terms. In *Proceedings of the Technical Communications of the 31st
1037 International Conference on Logic Programming (ICLP 2015), Cork, Ireland, August 31 - September 4, 2015*. [http://ceur-
1039 ws.org/Vol-1433/tc_12.pdf](http://ceur-
1038 ws.org/Vol-1433/tc_12.pdf)
- 1040 Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. 2011. Finding and understanding bugs in C compilers. In *Proceedings
1041 of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2011, San Jose, CA,
1042 USA, June 4-8, 2011*. 283–294. <https://doi.org/10.1145/1993498.1993532>
- 1043
- 1044
- 1045
- 1046
- 1047
- 1048
- 1049
- 1050
- 1051
- 1052
- 1053
- 1054
- 1055
- 1056
- 1057
- 1058
- 1059
- 1060
- 1061
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078