



ETNA: An Evaluation Platform for Property-Based Testing (Experience Report)

JESSICA SHI, University of Pennsylvania, USA

ALPEREN KELES, University of Maryland, USA

HARRISON GOLDSTEIN, University of Pennsylvania, USA

BENJAMIN C. PIERCE, University of Pennsylvania, USA

LEONIDAS LAMPROPOULOS, University of Maryland, USA

Property-based testing is a mainstay of functional programming, boasting a rich literature, an enthusiastic user community, and an abundance of tools — so many, indeed, that new users may have difficulty choosing. Moreover, any given framework may support a variety of strategies for generating test inputs; even experienced users may wonder which are better in a given situation. Sadly, the PBT literature, though long on creativity, is short on rigorous comparisons to help answer such questions.

We present ETNA, a platform for empirical evaluation and comparison of PBT techniques. ETNA incorporates a number of popular PBT frameworks and testing workloads from the literature, and its extensible architecture makes adding new ones easy, while handling the technical drudgery of performance measurement. To illustrate its benefits, we use ETNA to carry out several experiments with popular PBT approaches in both Coq and Haskell, allowing users to more clearly understand best practices and tradeoffs.

CCS Concepts: • **Software and its engineering** → **General programming languages**.

Additional Key Words and Phrases: property-based testing, empirical evaluation, mutation testing

ACM Reference Format:

Jessica Shi, Alperen Keles, Harrison Goldstein, Benjamin C. Pierce, and Leonidas Lampropoulos. 2023. ETNA: An Evaluation Platform for Property-Based Testing (Experience Report). *Proc. ACM Program. Lang.* 7, ICFP, Article 218 (August 2023), 17 pages. <https://doi.org/10.1145/3607860>

1 INTRODUCTION

Haskell’s QuickCheck library popularized *property-based testing* (PBT), which lets users test executable specifications of their programs by checking them on a large number of inputs. In fact, QuickCheck made PBT so popular that Claessen and Hughes’s seminal paper [2000] is the most cited ICFP paper of all time... by a factor of two, according to the ACM Digital Library. PBT tools can now be found in languages from OCaml [Cruanes 2017; Dolan 2017] and Scala [Nilsson 2019] to Erlang [Arts et al. 2008; Papadakis and Sagonas 2011] and Python [MacIver 2016], not to mention proof assistants like Coq [Lampropoulos and Pierce 2018], Agda [Lindblad 2007], and Isabelle [Bulwahn 2012a].

Many aspects of PBT impact its effectiveness, from the properties themselves [Hughes 2019] to counterexample minimization [Maciver and Donaldson 2020], but arguably the most crucial one

Authors’ addresses: Jessica Shi, University of Pennsylvania, Philadelphia, PA, USA, jwshi@seas.upenn.edu; Alperen Keles, University of Maryland, College Park, MD, USA, akeles@umd.edu; Harrison Goldstein, University of Pennsylvania, Philadelphia, PA, USA, hgo@seas.upenn.edu; Benjamin C. Pierce, University of Pennsylvania, Philadelphia, PA, USA, bpierce@cis.upenn.edu; Leonidas Lampropoulos, University of Maryland, College Park, MD, USA, leonidas@umd.edu.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2023 Copyright held by the owner/author(s).

2475-1421/2023/8-ART218

<https://doi.org/10.1145/3607860>

is the algorithm for generating test inputs. Papers citing QuickCheck often retain its distinctive style of *random* test-case generation, but many other options have been explored. In particular, *enumerative* PBT has also become a staple in the functional programming community [Braquehais 2017; Runciman et al. 2008], and tools for *feedback-based* PBT are gaining ground [Dolan 2017; Lampropoulos et al. 2019; Löschner and Sagonas 2017]. Each of these approaches comes with benefits and tradeoffs, and choosing one over another can make a big difference on testing effectiveness.

Even after selecting a generation style — say, random PBT — one may be left with quite a few options of *framework*, each with its own unique style. In Haskell, for example, both QuickCheck and Hedgehog [Stanley 2019] are quite popular. And even after selecting a framework — say, QuickCheck — there are yet more options for choosing a specific generation strategy. Tools like generic-random [Xia 2018] and DraGEN [Mista and Russo 2021] can derive QuickCheck generators from type information, offering a quick and accessible entrypoint to PBT, but their effectiveness suffers when inputs need to satisfy more complex semantic constraints. Alternatively, one can write a *bespoke* generator that is “correct by construction,” producing only *valid* test inputs. Such bespoke generators can sometimes become quite sophisticated [Hritcu et al. 2013; Midtgaard et al. 2017; Palka et al. 2011]. And there are other options: for example, QuickChick, Coq’s variant of QuickCheck, can derive specialized generators for free from specifications expressed as inductive relations [Paraskevopoulou et al. 2022]. Nuances of the properties under test may make strategies more or less preferable, and considerable experience may be required to make a good choice.

Moreover, even after selecting a particular way of using the tool — say, writing a bespoke generator — there are *yet more* options: a given generator can typically be tuned to produce different sizes and shapes of data. For example, QuickCheck generators can be parameterized both globally by a size parameter and locally by choices like numeric weights on the arguments to various combinators.

In the existing literature, there are plenty of performance evaluations, but a dearth of *comparisons* across these dimensions. New tools are typically evaluated on just one or two case studies, often showcasing incomparable measures of effectiveness. So how is a PBT user supposed to make sense of all these options? How is a tool designer supposed to measure success? How can we turn PBT from an art to a science?

Answering these questions is the goal of this experience report. Our contributions are:

- We present ETNA, an extensible platform for evaluating and comparing generation techniques for PBT, with generic support for measuring performance and presenting results (§2).
- We populate ETNA with five testing *workloads* from the literature, presenting a range of bug-finding challenges, with PBT *frameworks* in both Haskell and Coq, and with various *strategies* for using each framework (§3).
- We report on our experiences using ETNA to explore a number of questions about PBT. The answers generally lend weight to commonly held beliefs, but they also add useful nuance and suggest improvements to existing processes and tools (§4 and §5).

We discuss related and future work in §6.

2 PLATFORM DESIGN

Why are there so many generation strategies for PBT? In part, because there are many ways of dealing with properties with *preconditions*. Consider *binary search trees*, where each node value is greater than everything to its left and less than everything to its right. In Haskell syntax, we have

```
data Tree k v = Leaf | Node (Tree k v) k v (Tree k v)
isBST :: Tree k v -> Bool
insert :: k -> v -> Tree k v -> Tree k v
```

What properties should we expect to hold for operations on BSTs? [Hughes](#) thoroughly answers this question in his guide to writing properties of pure functions [2019]. For instance, one desirable property is that if we insert a key into a valid BST, then it should remain a valid BST:

```
prop_InsertValid :: Tree Int v -> Int -> Property
prop_InsertValid t k = isBST t ==> isBST (insert k t)
```

Here `==>` encodes a *precondition*. That is, the property is not expected to hold for an arbitrary binary tree `t` and an arbitrary key `k`, but only when the former satisfies the `isBST` predicate.

How should we generate data for such properties? A simple approach is to straightforwardly follow the structure of the types to generate arbitrary trees and filter out the ones that are not BSTs. While simplistic, this approach works well in some circumstances. In fact, for the BST example, such *type-driven* approaches can rapidly find all bugs introduced in [Hughes](#)'s guide [2019]. But this generate-and-filter approach breaks down with “sparser” preconditions; for instance, valid red-black trees are harder to generate at random than valid BSTs, so type-driven strategies work less well (see §4 and §5). For yet sparser preconditions, such as C programs with no undefined behaviors [[Yang et al. 2011](#)], the approach is hopeless.

The completely opposite approach is to require users to write *bespoke* generators: programs that are manually tailored to produce the desired distribution. Such programs can be extremely effective in finding bugs when the inputs satisfy the precondition by construction, but they can also be extremely difficult to write. A well-crafted such generator can in fact be a significant research result: this is the case for many well-typed term generators in the last decade [[Hoang et al. 2022](#); [Midtgaard et al. 2017](#); [Pałka et al. 2011](#)]. Naturally, there are also approaches in the middle. For instance, some use the structure of the precondition to produce valid data directly [[Bulwahn 2012b](#); [Claessen et al. 2014](#); [Fetscher et al. 2015](#); [Lampropoulos et al. 2017, 2018](#)], while others leverage feedback to guide generation towards valid or otherwise interesting inputs [[Lampropoulos et al. 2019](#); [Löcher and Sagonas 2018](#); [Löscher and Sagonas 2017](#)].

2.1 How to Evaluate Generators

How do we measure the effectiveness of a generator? The software testing literature offers two main answers: *code coverage* and *mutation testing*. Code coverage is popular, but problematic: it is well known that higher coverage does not always translate to better bug finding [[Gopinath et al. 2014](#); [Klees et al. 2018](#)]. We instead choose mutation testing [[Jia and Harman 2011](#)], which measures the effectiveness of testing by artificially injecting mutations to the system under test and checking if testing is able to detect them. Mutations in the literature [[Hazimeh et al. 2020](#); [Hritcu et al. 2013](#); [Klees et al. 2018](#); [Zhang et al. 2022](#)] fall on a spectrum from manually sourced to automatically synthesized. We opt for manual sourcing, allowing us to more readily maintain *ground truth* and ensure that every mutant violates some aspect of the property specification. ETNA supports a terse syntax for incorporating these mutants. In §3, we detail the systems evaluated in this paper.

2.2 Terminology

Our mutation-testing based evaluation is built upon **tasks**: a mutant-property pair where the mutant causes the property to fail. As any given program can give rise to multiple tasks — it might need to satisfy multiple properties or be subjected to multiple mutants — we organize tasks into **workloads**. Each workload comes with several components: data type definitions; variant implementations of functions using these types; and a property specification of these functions.

We call a PBT paradigm at the level of a library a **framework**, which should contain functions for (a) constructing properties, (b) constructing generators, and (c) running tests. For instance, QuickCheck, QuickChick, SmallCheck and LeanCheck are all examples of frameworks. And we call

a PBT paradigm at the level of how to use a framework to write generators a **strategy**. Examples of such strategies include type-based random generation, manually written bespoke generation, or exhaustive enumeration of the input space.

2.3 Architecture

ETNA is designed to be an extensible platform that flexibly accommodates new workloads, strategies, frameworks, and languages. At its core is an experiment driver that provides three main pieces of functionality: (a) toggling between variant implementations in a directory of workloads; (b) compiling and running each strategy on each task; and (c) analyzing the results.

The driver is implemented as a Python library and is run via an *experiment script*. An example experiment script can be found in Appendix A. If a user simply hopes to replicate our experiments, they can directly use the scripts provided in our artifact; if they want to run their own, they can adapt one of the existing scripts. Each experiment script calls into the driver, producing some data, and then processes that data via tools like Pandas [pandas 2023] and Plotly [Plotly 2023].

To add a new workload, the user implements the system under test just as they would ordinary code in that language. The user can then encode mutants via special comment syntax embedded within the implementation. For example, below is the correct implementation of BST insert followed by a mutant in the Haskell syntax:

```
insert k v E = T E k v E
insert k v (T l k' v' r)
  {-! -}
  | k < k' = T (insert k v l) k' v' r
  | k > k' = T l k' v' (insert k v r)
  | otherwise = T l k' v r
{-!! insert_becomes_singleton -}
{-!
= T E k v E
-}
```

To add a new strategy, the user instantiates per-framework combinators with their custom approach. For example, QuickCheck strategies can provide Arbitrary instances. To add a new framework, the user can connect the standardized language infrastructure with the specifics of the framework's test runner. For example, this may involve transforming our property types to work with their precondition combinators and transforming their outputs to align with our result types.

To add a new language, the user needs to implement an adapter that tells the driver about the directory structure and compilation commands, as well as some language-specific infrastructure for capturing results. The adapters are written in Python and designed to be simple to implement and to maximize reuse of infrastructure between languages. Test results should be produced in a standardized JSON output format, which allows for unified analysis across all languages.

The code for ETNA is publicly available both on Zenodo¹ and GitHub².

2.4 Analysis and Presentation

Though ETNA supports customizable experiments, we choose a standard set of defaults for the experiments in this paper. We run each strategy on each task for a set amount of trials

14	6	2	10	4
< 0.1 s	0.1 to 1 s	1 to 10 s	10 to 60 s	unsolved

¹<https://doi.org/10.5281/zenodo.7993347>

²<https://github.com/jwshi21/etna>

(10 unless otherwise specified) and with a set timeout (60 seconds). We then measure if the strategy was able to *solve the task*, i.e. find the injected bug in all trials within the given time frame. Multiple trials account for the non-determinism of random generation strategies, and results are simple averages unless indicated otherwise.

Our first attempts at presenting this data were hard to interpret: what does it mean, for example, if one strategy takes an average of two seconds and the other an average of three? Rather than present a slew of raw numbers, we wanted a data representation that captures a user’s experience of interacting with PBT tools, so that visual differences in the representation correspond to tangible differences in performance. The figure above demonstrates our solution: a *task bucket* chart. For every strategy we classify tasks ranging from “solved instantly” to “unsolved”, depicted with progressively lighter shades. For example, for the strategy/workload combination in the figure, 14 tasks are solved very quickly (the darkest shade) while four are not solved at all (the lightest).

In case a task bucket chart does not show enough detail, especially in head-to-head comparisons, we also support statistical analyses like Mann–Whitney U tests³ (see §4.1).

3 POPULATING THE PLATFORM

We have integrated a number of PBT frameworks and workloads into ETNA, for our own use in §4 and §5 and for potential users to use and compare against.

3.1 Languages and Frameworks

Haskell is an obvious starting point: as the language that hosts QuickCheck, it is the lingua franca of PBT research. We focus on three Haskell frameworks: QuickCheck, of course; SmallCheck [Runciman et al. 2008], a competitor to QuickCheck that does enumerative testing; and LeanCheck [Braquehais 2017], a more modern enumerative framework.

Our second language of choice is Coq. While Haskell has many PBT frameworks, PBT in Coq is built on a single framework: QuickChick [Lampropoulos and Pierce 2018]. However, QuickChick is a rich ecosystem that supports a variety of different strategies for input generation [Lampropoulos 2018; Lampropoulos et al. 2019, 2018], so there is plenty to study and compare.

We focus on intra-language experiments in this paper, though future work could plausibly consider inter-language ones as well. ETNA’s extensible design means that adding new languages is straightforward; we discuss languages that we plan to add to the platform in §6.

3.2 Workloads

Our initial set of workloads is drawn from three application domains that are of practical interest to the functional programming community and that have featured prominently in the PBT literature. These workloads feature in the following sections’ experiments, although not every workload is used for every experiment.

Data Structures. The first workload focuses on a functional data structure that is ubiquitous in the literature: binary search trees. Multiple PBT papers have focused on BST generation, including John Hughes’s *How to Specify It!* [2019], an extended introduction to specifying properties using QuickCheck. Our BST workload ports the mutations and properties from that paper. The second workload focuses on another popular functional data structure, red-black trees, including self-balancing insertion and deletion operations that are notoriously easy to get wrong. RBTs have also been studied in the PBT literature [Klein and Findler 2009; Lampropoulos et al. 2017; Mista and

³The Mann–Whitney U test is a nonparametric test that compares data samples from two different distributions. We use it here because it makes no assumptions about the distributions being compared.

Russo 2019; Runciman et al. 2008]. Our RBT workload combines the BST mutants with additional mutants that focus on potential mistakes when balancing or coloring the tree.

Lambda Calculi and Type Systems. The third workload centers a DeBruijn index based implementation of the simply typed lambda calculus with booleans. Bespoke generators for producing well-typed lambda terms is a well studied problem in the literature [Midtgaard et al. 2017; Palka et al. 2011], while the mutations for STLC included in our case study are drawn from the appropriate fragment of a System F case study [Goldstein et al. 2021], dealing mostly with mistakes in substitution, shifting, and lifting. For a more complicated fourth workload $F_{<}$, revolving around calculi and type systems, we turn to the full case study of Goldstein et al. [2021] and extend it with subtyping. This allows for significantly more complex errors to be injected (such as those dealing with type substitution, shifting, or lifting). Bespoke generators for System F have been the subject of recent work [Goldstein et al. 2021; Hoang et al. 2022] and can be straightforwardly extended to handle subtyping.

Security. The fifth and final workload focuses on a security domain: information flow control. The IFC case study, introduced by Hritcu et al. [2013, 2016], explores the effectiveness of various bespoke generators for testing whether low-level monitors for abstract machines enforce noninterference: differences in secret data should not become publicly visible through execution. Violations in the enforcement policies are introduced by systematically weakening security checks or taint propagation rules, exploring all possible ways of introducing such violations.

4 EXPERIMENTS: HASKELL

We next report on our experience using ETNA to probe different aspects of testing effectiveness. Our first set of observations are on the PBT frameworks and strategies available in Haskell.

4.1 Comparing Frameworks

In the first experiment, we assess the “out of the box” bug-finding abilities of three Haskell frameworks — QuickCheck, SmallCheck, and LeanCheck. We examine four strategies. For the *bespoke* strategy, we manually write a QuickCheck generator that always produces test inputs that satisfy the property’s precondition. This serves as a “topline” for the other strategies: a high-effort generator that solves all of the tasks easily. The other three strategies — one per framework — are all *naive*. The QuickCheck strategy uses the generic-random library to derive its generator automatically, with constructors chosen at each step with uniform probability and a size parameter that decreases on recursive calls to ensure termination. For the enumerative frameworks, SmallCheck and LeanCheck, we use combinators that follow the type structure.

We evaluate these strategies against four workloads: BST, RBT, STLC, and $F_{<}$.

Results. We visualize the results of this experiment in Figure 1. Some points to note:

The bespoke strategy outperforms the naive strategies along multiple axes. For example, looking at the naive QuickCheck strategy (the others are similar), the bespoke strategy solved all tasks, while the naive strategy failed to solve 43 tasks. Among tasks that both strategies solved, using a Mann–Whitney U test with $\alpha = 0.05$, we find that the bespoke strategy’s average time to solve a task was (statistically) significantly lower in 83 out of 124 tasks and the average valid inputs to solve a task were lower for 89 out of 124 tasks. That is, the bespoke strategy found more bugs, more quickly, and with better quality tests.

Between the two enumeration frameworks, LeanCheck substantially outperforms SmallCheck on these workloads. LeanCheck had an 82% solve rate, while SmallCheck’s was only 35%. On one BST task, LeanCheck found the bug in about a hundredth of a second on average, while SmallCheck required 26 seconds. One reason for these differences may be that SmallCheck attempts to enumerate



Fig. 1. Effectiveness of Haskell generation strategies on four workloads.

■ = Naive QuickCheck, ■ = Naive LeanCheck, ■ = Naive SmallCheck, ■ = Bespoke QuickCheck.

larger inputs much earlier. In the first thousand binary trees, SmallCheck produces trees with up to ten nodes, while LeanCheck only reaches four nodes. Unsurprisingly, it is harder for larger trees to satisfy the BST invariant — only 1% of these thousand SmallCheck trees are valid, compared to 13% of the LeanCheck trees. And across all workloads, we can calculate the rate at which they enumerate test inputs, by aggregating over the tasks that they both solved and dividing by the total number of tests by the total time spent. We find that LeanCheck produces over a hundred times more tests per second than SmallCheck.

LeanCheck also outperforms naive QuickCheck. It is illuminating to consider failed tasks that were *partially solved* — the bug was found in at least one trial and not found in at least one trial. There is one such task for LeanCheck and 14 for QuickCheck. For LeanCheck’s partially solved task, the inputs required are the same for each trial, but the time fluctuates between 55 and 65 seconds. That is, this is a situation where a task nears — and sometimes exceeds — what LeanCheck can reach with the one minute time limit. QuickCheck’s partially solved tasks are also interesting. Of the 13 that LeanCheck solves but QuickCheck does not, 10 are partially solved by QuickCheck. This suggests that there are situations where a deterministic approach may be more reliable than a random alternative: LeanCheck solves these tasks consistently and relatively quickly, while QuickCheck sometimes takes less than a second, sometimes nearly a minute, and sometimes times out.

Similarly, in STLC, naive LeanCheck solves three more tasks within the first bucket than the bespoke strategy. Upon closer inspection, these are tasks that the bespoke strategy sometimes solves in under 100 inputs but sometimes requires over 10,000 inputs, leading to an average slightly above the 0.1 second threshold; as before, LeanCheck does not experience this variability.

A note about memory usage. LeanCheck is documented⁴ to be memory intensive, especially when run for prolonged periods of time, as we do here. Our experiments using LeanCheck were conducted on a server with plenty of memory, allowing us to complete trials without issues. Future work might consider the relative space complexities of different frameworks.

⁴<https://github.com/rudymatela/leancheck/blob/master/doc/memory-usage.md>

4.2 Exploring Sized Generation

We next explore the sensitivity of bug-finding to various parameters, starting with input size.

A significant part of generator tuning is ensuring that the generated inputs are well sized. Conventional wisdom in random testing posits that there is a “combinatorial advantage” to testing with large inputs, since they can exercise many program behaviors at once; tools like *QuickCover* [Goldstein et al. 2021] capitalize on this notion to make testing more efficient. But are large inputs *always* better? We used our BST workload to investigate.

We conducted this experiment on the QuickCheck framework, using a bespoke strategy to focus attention on the quality of the distribution of *valid* inputs. We used a generator from Hughes [2019], which generates a list of keys and then inserts each key into the tree, because it gives precise control over final tree sizes. We choose the keys for a n -node tree from a range of integers 1 to $2n$. This range is large enough to allow for sufficient variety in shape and content but not so large that a randomly generated key is unlikely to be in the tree.

We then measured the bug-finding effectiveness of the generator at different sizes n . Thanks to ETNA’s flexibility, we could vary the size in the script and otherwise treat this experiment as we would any other where we wanted to compare several strategies.

Results. Figure 2 plots the size of the tree versus the average number of inputs to solve a task; each trace represents one task. Some noteworthy traces, highlighted in black, are discussed below.

Larger trees can be worse for bug-finding, for properties that rely on dependencies between their inputs. We found that, for BST, small trees were generally sufficient to find bugs, and performance got significantly worse for some tasks as trees got larger.

For example, task #1, which has the steepest upward curve, involves a mutant where the delete function fails to remove a key unless that key happens to be the root. The property takes one tree and two keys as inputs and checks that removing the keys in either order leads to the same result. Together, these mean that the task is only solvable when one key k is the root of the tree and the other key k' becomes the root after deleting k . The probability of satisfying this condition decreases as the size of the tree increases, so larger trees take more inputs to solve this task.

Task #2 is a similar story. It takes a tree and two key-value pairs; this time, the task is only solvable when the two keys are the same (and the two values are different), a probability that is inversely proportional to the size of the tree. These two tasks demonstrate situations where the inputs to a property need to be related in a mutant-specific way, and large trees are less likely to satisfy this dependency relationship.

Not all tasks with dependencies between their inputs are harder to solve with larger trees. Unlike #1 and #2, the curve for task #3 is mostly flat, even though it has a similar dependency. The mutant here causes the union operation to fail by occasionally preferring the wrong value if both trees contain the same key; the property takes a key k and two trees and checks that k exists in the union of the trees when it exists in either tree. Since this mutant causes problems with keys that appear

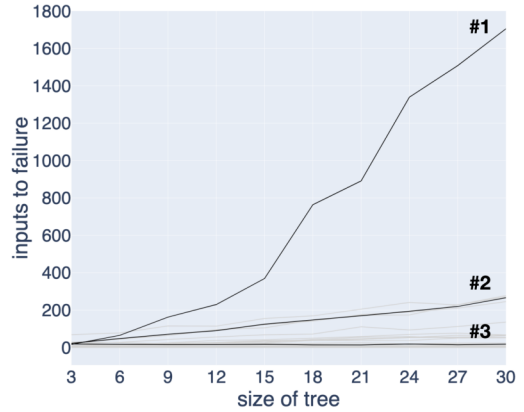


Fig. 2. Number of generated inputs (averaged over 100 trials) to solve each BST task, as the size of the inputs increases from three to 30 nodes.

in both trees, the property only fails when k is in the input trees. That is, there is a dependency between the inputs, but this dependency does *not* scale with the size of the tree.

Discussion. We have seen that larger inputs sometimes not only fail to provide a combinatorial advantage but in fact can provide a dependency disadvantage. The size of the main input — in this case, the tree — cannot be evaluated in a vacuum. Instead, the particulars of the mutant and property can lead to dependencies between the property inputs that must be satisfied in order to detect the mutant. Our size exploration is thus a cautionary tale: PBT users should not naively expect that larger inputs are better, especially for properties with multiple inputs.

This exploration suggests a few recommendations for improving both testing frameworks and individual users' choices of properties. (1) *Do not treat property inputs as independent.* The difficulties with the above properties arise, in part, because QuickCheck automates generation of multiple inputs by assuming that each input can be generated independently — but treating inputs independently can lead to unintuitive testing performance. Frameworks like Hedgehog explicitly avoid introducing a generator typeclass so as to force users to build generators by hand; our results lend credence to that design choice. (2) *Think carefully about properties with multiple inputs.* Testers should prefer properties with fewer inputs where possible. When this is infeasible, testers should think carefully about potential interactions between their property's inputs and design generators that take those interactions into account.

4.3 Enumerator Sensitivity

Papers about enumeration frameworks sometimes speak of enumeration as a kind of exhaustive testing — validating the program's behavior within a “small scope” [Andoni et al. 2002]. But realistic testing budgets often mean that exhausting all inputs up to some interesting size or depth is not possible: enumeration is *expensive*. Thus, the actual performance of enumeration frameworks like SmallCheck and LeanCheck is impacted by the specific order in which values are enumerated. In this section we examine some factors that, perhaps unexpectedly, impact bug-finding performance.

There are many axes along which order could vary. We have explored two: the order of the inputs to each property and the order of constructors in an algebraic data type. We focus our remarks here on the former and relegate more detailed data for both to Appendix B.

We conduct this experiment on SmallCheck and LeanCheck, using the BST and RBT workloads; many of their properties have multiple inputs, including a combination of Trees and Int keys. One enumeration strategy uses the default properties, with the trees passed in first, and one uses properties with the trees last — for example, (Tree, Tree, Int) vs. (Int, Tree, Tree).

Results. We count the number of tasks that are solved by the same framework under exactly one of the two orderings. For LeanCheck, the tree-last strategy solved one additional task that the tree-first strategy did not (completing in about 38 seconds instead of timing out at 60). For SmallCheck, the tree-last strategy solved 17 more tasks than tree-first, taking between 0.002 and 7 seconds. The low end is especially remarkable: simply by enumerating (Int, Tree, Tree)s rather than (Tree, Tree, Int)s, SmallCheck finds a counterexample immediately instead of timing out.

Discussion. A deeper dive into the enumeration frameworks to explore these differences fully would be worthwhile, but what jumps out even from these simple experiments is the question of *sensitivity*. The potentially pivotal role of enumeration order in the success or failure of these strategies means that users of these enumerative frameworks need to be careful of configuration settings that would be immaterial in their random counterparts. As a meta point, we put the tree data types at the front of each property as a matter of convention; it was not until much later that we realized the inadvertent effect on the performance of the enumerators!

5 EXPERIMENTS: COQ

After focusing on the multi-framework landscape of Haskell in the previous section, we now turn our attention to the single-framework but multi-strategy landscape in Coq. As discussed in §3.1, PBT in Coq revolves around QuickChick [Lampropoulos 2018], which, in addition to the type-based and bespoke strategies that we explored in Haskell, provides two additional options: a *specification-driven* strategy that derives correct-by-construction generators from preconditions in the form of inductive relations [Lampropoulos et al. 2018] and a *type-driven fuzzer* strategy that combines type-based generation with mutation informed by AFL-style branch coverage to guide the search toward interesting parts of the input space [Lampropoulos et al. 2019].

Both papers exemplify the lack of performance comparisons across approaches discussed in the introduction. First, Lampropoulos et al. [2018] is evaluated in a toy IFC example, where only the throughput of generators is measured against that of a bespoke generator; there is no measurement of the effectiveness of the strategy in finding bugs. On the other hand, FuzzChick [Lampropoulos et al. 2019] is evaluated in the more realistic IFC workload of Hritcu et al. [2016] that we will reuse later in this section, with systematically injected mutations that break the enforcement mechanism of a dynamic monitor. Still, multiple aspects of their strategies were left unevaluated, including their performance on any other workload.

5.1 Comparison of Fuzzers, Derived Generators, and Handwritten Generators

We aim to fill the evaluation gaps described above. How do QuickChick’s newer strategies compare with the more established bespoke and type-based ones? In particular, are they effective at uncovering bugs across disparate workloads?

We again use the BST, RBT, and STLC workloads, along with a more complex case study, IFC, pulled from the FuzzChick paper. For the first three case studies, inductively defined specifications are widely available (e.g. in Software Foundations [Pierce 2018]); for IFC, such specifications do not exist, so the specification-driven generators of Lampropoulos et al. does not apply.

Results. In Figure 3, we visualize the results of the experiments with a task bucket chart. Results for the simple BST workload (Figure 3a) establish a baseline level of confidence for all four methods, as they are all able to solve most tasks quickly. Indeed, most of the tasks are solved by all methods within 0.1 seconds (the darkest color), with the exception of the *type-based fuzzer*, which falls short on a few tasks.

Specification-derived strategies are on par with bespoke ones. In the harder RBT workload, with its much more complex invariant, there is a clear performance gap between type-driven strategies (type-based generator and type-based fuzzer) and precondition-driven methods (specification-based generator and bespoke generator). Precondition-driven methods are able to solve more tasks under 0.1 seconds than type-driven methods are able to solve within a 60 second timeout. The type-based generator fails to solve 23 tasks, and the type-based fuzzer fails to solve 25. The bespoke generator solves all tasks in under ten seconds, and the specification-based generator solves all but 10 tasks. We see a similar pattern in the STLC workload, with the precondition-driven methods outperforming the type-driven ones.

Fuzzers exhibit more variance but outperform type-driven methods for sparse preconditions. For the IFC workload, the only precondition-driven strategy is the bespoke generator, which emerges as a clear winner: noninterference is a property with a *very* sparse precondition, and type-based methods are basically unable to generate valid inputs. For this particular workload, we included another fuzzing variant borrowed from the original paper that introduced FuzzChick [Lampropoulos et al. 2019] to strengthen the connection to the existing literature: rather than generating a pair of input machines completely at random and then fuzzing the pair (as in the type-based fuzzer approach),

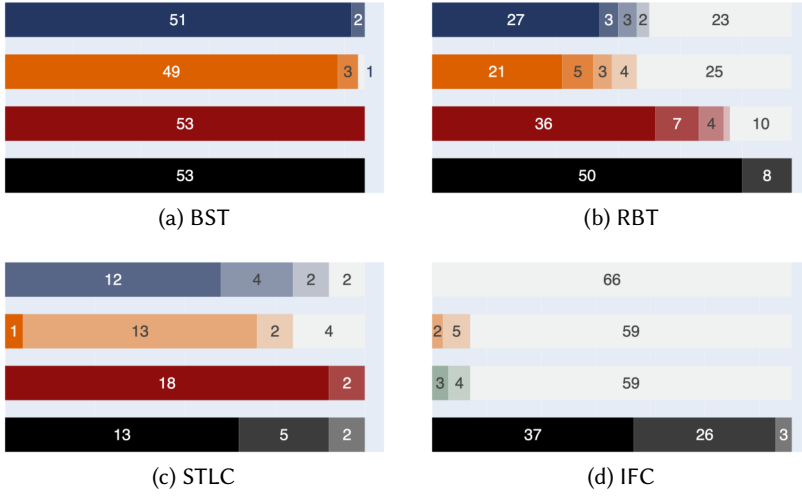


Fig. 3. Effectiveness of Coq generation strategies on four workloads.

■ = Type-based generator, ■ = Type-based fuzzer, ■ = Specification-based generator ((a)–(c) only), ■ = Variational fuzzer ((d) only), ■ = Bespoke generator.

we generate one input machine and copy it to create a pair that is indistinguishable by default. The two fuzzers, *type-based fuzzer* and *variational fuzzer*, have a clear advantage over the pure type-based generation approach: the ability to guide generation allows fuzzers to discover parts of the input space that naive type-based generation are simply unable to reach.

Yet fuzzers are not reliable in this sense, as Figure 4 shows: if we include *partially solved* tasks, fuzzers outperform their generator counterparts. This further clarifies the picture painted by the first set of comparisons. Fuzzers may get stuck following program paths that will not lead to interesting revelations, but sometimes discover paths that a traditional type-based generator could never hope to reach. In particular, roughly 30 tasks are solved *at least once* through 10 runs (Figure 4), but less than 10 tasks are fully solved (Figure 3d).

Another interesting observation is that even though fuzzers typically spend more time per generated input, as the underlying types are more complex and large, mutating the input takes less time than generating a new one. For IFC, the type-based generator takes four times longer per input than the type-based fuzzer.

5.2 Validation and Improvement of Fuzzers

Despite its minimal evaluation, the conclusion of Lampropoulos et al. [2019] seems to hold — that is, FuzzChick shows promise compared to type-based approaches, but has a long way to go before catching up with the effectiveness of precondition-driven ones. This led us to wonder, *could we further improve the performance of FuzzChick using ETNA?*

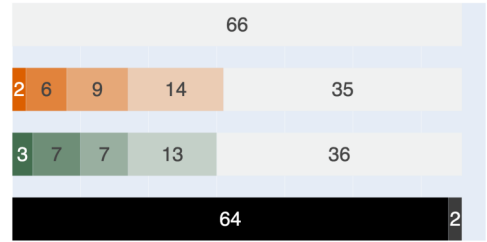


Fig. 4. Tasks solved within the timeout in one or more trials. Empty = Type-based generator, ■ = Type-based fuzzer, ■ = Variational fuzzer, ■ = Bespoke generator.

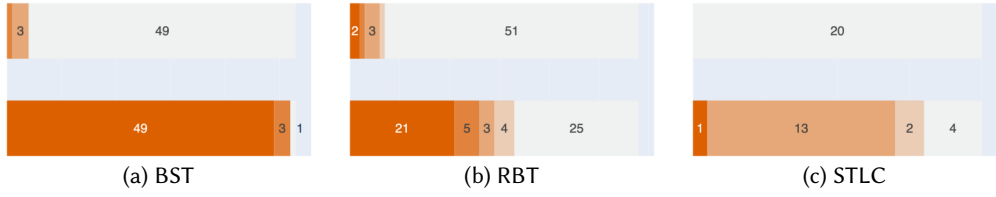


Fig. 5. Comparison of the original FuzzChick (top) with the tuned one (bottom).

We focused on two different aspects of fuzzing: *size* and *feedback*. FuzzChick’s generation strategy started small but quickly increased to quite large sizes, relying on the idea of “combinatorial advantage” discussed in §4.2 — i.e., that larger inputs contain exponentially many smaller inputs and are therefore more effective for testing. As we saw there, that is not always the case. After realizing this, we switched to a more gently increasing size bound which led to significant improvements in terms of throughput, positively impacting our bug-finding ability.

With respect to feedback, by using ETNA to evaluate FuzzChick across multiple workloads we were able to identify, isolate, and fix a bug that caused it to saturate the seed pool with uninteresting inputs. FuzzChick (like Zest [Padhye et al. 2019]) keeps two seed pools: one for valid and one for invalid inputs. FuzzChick’s bug applied to the latter one and was hidden from its authors as the variational fuzzer strategy they employed readily gives access to valid inputs (which are prioritized).

Results. Figure 5 demonstrates the bug-finding capabilities of the original (top) and tuned (bottom) versions of FuzzChick across the new workloads. The tuned version clearly outperforms the original in all cases — and is what was used in the previous section.

6 RELATED AND FUTURE WORK

The future directions we imagine for ETNA are inspired by related work in the literature. Thus, we discuss both related and future work together in this section.

ETNA’s name, referencing every crossword-puzzler’s favorite Italian volcano, was inspired by two existing benchmark suites in the fuzzing space: LAVA [Dolan-Gavitt et al. 2016] and Magma [Hazimeh et al. 2020]. Both provide a suite of workloads that can be used to compare different fuzzing tools: LAVA’s workloads consist of programs with illegal memory accesses that are automatically injected, while Magma relies on real bugs forward-ported to the current versions of libraries. More recently, FixReverter [Zhang et al. 2022] offered a middle ground, generalizing real bug-fixes into patterns and applying them to multiple locations in a program. ETNA is different from these suites in a few ways. First, ETNA aims to be a platform for exploration and evaluation rather than a rigid set of benchmarks. Thus, we do not claim that ETNA’s workloads are complete — instead, we intend for users to add more over time. Additionally, evaluating fuzzing is quite different from evaluating PBT, since PBT is expected to run for less time on programs with higher input complexity. This means that ETNA’s measurement techniques and workload focus must necessarily be different from LAVA’s or Magma’s. Still, there are ideas worth borrowing from these suites: fuzzing benchmarks generally record code-coverage information, which we plan for ETNA to eventually offer as well.

Besides LAVA and Magma, there is a massive literature of Haskell and Coq papers from which we will continue to draw both workloads and frameworks. With the help of the community, we hope ETNA will eventually include frameworks like: Luck [Lampropoulos et al. 2017], a language for preconditions from which generators can be inferred; FEAT [Duregård et al. 2012], an enumerator

framework focusing on uniformity; tools for deriving better Haskell generators [Mista and Russo 2019, 2021]; and specification-driven enumerators for QuickChick [Paraskevopoulou et al. 2022].

Outside of Haskell and Coq, there are yet more opportunities for growth. Most immediately, support for OCaml tools like Crowbar [Dolan 2017] and QCheck [Cruanes 2017] seem within reach. This is particularly appealing as we could do inter-language comparisons; after all, Coq runs via extraction to OCaml. We will also solicit framework maintainers and researchers to add support for other languages such as Python (Hypothesis [MacIver 2016]), Scala (SciFe [Kuraj and Kuncak 2014; Kuraj et al. 2015] and ScalaCheck [Nilsson 2019]), Erlang (QuviQ [Arts et al. 2008] and PropEr [Papadakis and Sagonas 2011]), or Isabelle [Bulwahn 2012a,b].

Finally, the presentation backend of ETNA is fit-for-purpose, but we intend to do further research into the best possible ways to visualize PBT results. Consulting experts in human-computer interaction, we plan to use tools like Voyager [Wongsuphasawat et al. 2017] to explore which kinds of outcome visualizations real users of ETNA want. At the very least, integrating ETNA into a Jupyter notebook [Jupyter 2023] and providing hooks into a powerful graphics engine like Vega-lite [Satyanarayan et al. 2017] would make it easier for users to experiment with visualizations.

7 CONCLUSION

We designed ETNA to meet a concrete need in our research — we needed a clear way to convince ourselves and others that the PBT tools we build are worth pursuing. ETNA provides that, with an extensible suite of interesting workloads and the infrastructure necessary to validate and refine designs against them. In §4 and §5, we originally set out to answer straightforward questions about whether X is better than Y, and while we did get feedback about general trends, we also uncovered some unexpected nuances of the testing process. PBT-curious readers may have further questions building upon and extending beyond our explorations. ETNA is there for you!

ACKNOWLEDGMENTS

We thank John Hughes, Koen Claessen, Alejandro Russo, Augustin Mista, and Michael Hicks for their helpful comments. This work was supported by the NSF under award #1955610 *Bringing Python Up to Speed* and under #2145649 *CAREER: Fuzzing Formal Specifications* (any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF).

A APPENDIX: THE ANATOMY OF AN EXPERIMENT SCRIPT

The following script executes a simple example experiment in ETNA. The user can select the specific tasks they want to solve and the strategies that they use to solve them. At the end, the results are printed to a set of JSON files which can then be visualized.

```
from benchtool.Haskell import Haskell
from benchtool.Types import TrialConfig
from benchtool.Analysis import *

results = 'results/'
tool = Haskell(results)

timeout = 60
```

```

# Loop through available workloads...
for workload in tool.all_workloads():
    # ... filtering out any that are not necessary for this experiment.
    if workload.name not in ['STLC', 'FSUB']:
        continue

    # Then loop through mutants...
    for mutant in tool.all_mutants(workload):
        # ... and perform the necessary replacement in the implementations.
        run_trial = tool.apply_mutant(workload, mutant)

    # Loop through strategies...
    for strategy in tool.all_strategies(workload):

        # ... and properties ...
        for prop in tool.all_properties(workload):
            # ... and run each trial.
            cfg = TrialConfig(workload=workload,
                             strategy=strategy,
                             prop=prop,
                             trials=10,
                             timeout=timeout)

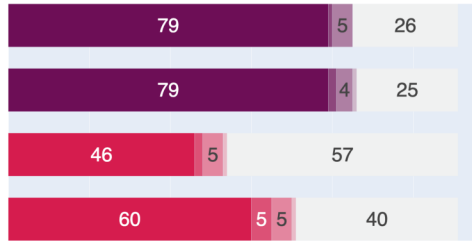
            run_trial(cfg)

# Finally, parse the results into a Pandas dataframe and
# use our provided analysis functions.
df = parse_results(results)
print(overall_solved(df, 'all', within=timeout))

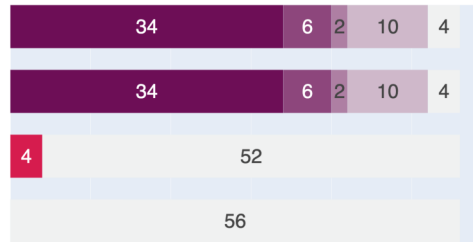
```

B APPENDIX: TASK BUCKETS FOR ENUMERATOR SENSITIVITY EXPERIMENT

Full data for enumerator sensitivity, where = Naive LeanCheck, = Naive SmallCheck.



This first chart compares the performance on the BST and RBT workloads when the trees are at the start of the properties (top rows) versus when they are at the end (bottom rows).



This second chart compares the performance on the STLC and FSUB workloads when the constructor enumeration order aligns with the definition of the data type (top rows) versus when the orders are reversed (bottom rows).

REFERENCES

- Alexandr Andoni, Dumitru Daniliuc, Sarfraz Khurshid, and Darko Marinov. 2002. Evaluating the "Small Scope Hypothesis". (Oct 2002).
- Thomas Arts, Laura M. Castro, and John Hughes. 2008. Testing Erlang Data Types with QuviQ QuickCheck. In *Proceedings of the 7th ACM SIGPLAN Workshop on Erlang*. ACM, 1–8. <https://doi.org/10.1145/1411273.1411275>
- Rudy Matela Braquehais. 2017. *Tools for Discovery, Refinement and Generalization of Functional Properties by Enumerative Testing*. Ph.D. Dissertation. University of York.
- Lukas Bulwahn. 2012a. The New Quickcheck for Isabelle - Random, Exhaustive and Symbolic Testing under One Roof. In *2nd International Conference on Certified Programs and Proofs (Lecture Notes in Computer Science, Vol. 7679)*. Springer, 92–108.
- Lukas Bulwahn. 2012b. Smart Testing of Functional Programs in Isabelle. In *18th International Conference on Logic for Programming, Artificial Intelligence, and Reasoning (Lecture Notes in Computer Science, Vol. 7180)*. Springer, 153–167.
- Koen Claessen, Jonas Duregård, and Michal H. Palka. 2014. Generating Constrained Random Data with Uniform Distribution. In *Functional and Logic Programming (Lecture Notes in Computer Science, Vol. 8475)*. Springer, 18–34. https://doi.org/10.1007/978-3-319-07151-0_2
- Koen Claessen and John Hughes. 2000. QuickCheck: a lightweight tool for random testing of Haskell programs. In *Proceedings of the 5th ACM SIGPLAN International Conference on Functional Programming*. ACM, 268–279. <https://doi.org/10.1145/351240.351266>
- Simon Cruanes. 2017. QuickCheck Inspired Property-Based Testing for OCaml. <https://github.com/c-cube/qcheck/>.
- Stephen Dolan. 2017. Property Fuzzing for OCaml. <https://github.com/stedolan/crowbar>.
- Brendan Dolan-Gavitt, Patrick Hulin, Engin Kirda, Tim Leek, Andrea Mambretti, Wil Robertson, Frederick Ulrich, and Ryan Whelan. 2016. LAVA: Large-Scale Automated Vulnerability Addition. In *2016 IEEE Symposium on Security and Privacy*. 110–121. <https://doi.org/10.1109/SP.2016.15>
- Jonas Duregård, Patrik Jansson, and Meng Wang. 2012. Feat: Functional Enumeration of Algebraic Types. In *Proceedings of the 2012 Haskell Symposium*. ACM, 61–72. <https://doi.org/10.1145/2364506.2364515>
- Burke Fetscher, Koen Claessen, Michal H. Palka, John Hughes, and Robert Bruce Findler. 2015. Making Random Judgments: Automatically Generating Well-Typed Terms from the Definition of a Type-System. In *24th European Symposium on Programming (Lecture Notes in Computer Science, Vol. 9032)*. Springer, 383–405.
- Harrison Goldstein, John Hughes, Leonidas Lampropoulos, and Benjamin C. Pierce. 2021. Do Judge a Test by its Cover. In *Programming Languages and Systems (Lecture Notes in Computer Science)*. Springer, 264–291. https://doi.org/10.1007/978-3-030-72019-3_10
- Rahul Gopinath, Carlos Jensen, and Alex Groce. 2014. Code Coverage for Suite Evaluation by Developers. In *Proceedings of the 36th International Conference on Software Engineering*. ACM, 72–82. <https://doi.org/10.1145/2568225.2568278>
- Ahmad Hazimeh, Adrian Herrera, and Mathias Payer. 2020. Magma: A Ground-Truth Fuzzing Benchmark. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1–29. <https://doi.org/10.1145/3428334>
- Tram Hoang, Anton Trunov, Leonidas Lampropoulos, and Ilya Sergey. 2022. Random Testing of a Higher-Order Blockchain Language (Experience Report). In *Proceedings of the 27th ACM SIGPLAN International Conference on Functional Programming*. <https://doi.org/10.5281/zenodo.6778257>
- Catalin Hritcu, John Hughes, Benjamin C. Pierce, Antal Spector-Zabusky, Dimitrios Vytiniotis, Arthur Azevedo de Amorim, and Leonidas Lampropoulos. 2013. Testing Noninterference, Quickly. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming*. <https://doi.org/10.1145/2544174.2500574>

- Catalin Hritcu, Leonidas Lampropoulos, Antal Spector-Zabusky, Arthur Azevedo Amorim, Maxime Denes, John Hughes, Benjamin C. Pierce, and Dimitrios Vytiniotis. 2016. Testing Noninterference, Quickly. In *Journal of Functional Programming*. <https://doi.org/10.1017/S0956796816000058>
- John Hughes. 2019. How to Specify It! - A Guide to Writing Properties of Pure Functions. In *Symposium on Trends in Functional Programming*. https://doi.org/10.1007/978-3-030-47147-7_4
- Yue Jia and Mark Harman. 2011. An Analysis and Survey of the Development of Mutation Testing. *IEEE Transactions on Software Engineering* 37, 5 (2011), 649–678. <https://doi.org/10.1109/TSE.2010.62>
- Project Jupyter. 2023. Project Jupyter. <https://jupyter.org>
- George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. 2018. Evaluating Fuzz Testing. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2123–2138. <https://doi.org/10.1145/3243734.3243804>
- Casey Klein and Robert Bruce Findler. 2009. Randomized Testing in PLT Redex. In *Workshop on Scheme and Functional Programming*.
- Ivan Kuraj and Viktor Kuncak. 2014. SciFe: Scala framework for efficient enumeration of data structures with invariants. In *Proceedings of the 5th Annual Scala Workshop*. ACM, 45–49. <https://doi.org/10.1145/2637647.2637655>
- Ivan Kuraj, Viktor Kuncak, and Daniel Jackson. 2015. Programming with Enumerable Sets of Structures. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*. <https://doi.org/10.1145/2814270.2814323>
- Leonidas Lampropoulos. 2018. *Random Testing for Language Design*. Ph.D. Dissertation. University of Pennsylvania.
- Leonidas Lampropoulos, Diane Gallois-Wong, Catalin Hritcu, John Hughes, Benjamin C. Pierce, and Li yao Xia. 2017. Beginner’s Luck: a language for property-based generators. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages*. <https://doi.org/10.1145/3009837.3009868>
- Leonidas Lampropoulos, Michael Hicks, and Benjamin C. Pierce. 2019. Coverage Guided, Property Based Testing. In *Proceedings of the 2019 ACM Conference on Object-Oriented Programming Languages, Systems, and Applications*. <https://doi.org/10.1145/3360607>
- Leonidas Lampropoulos, Zoe Paraskevopoulou, and Benjamin C. Pierce. 2018. Generating Good Generators for Inductive Relations. In *Proceedings of the ACM Conference on Principles of Programming Languages (POPL)*. <https://doi.org/10.1145/3158133>
- Leonidas Lampropoulos and Benjamin C. Pierce. 2018. *QuickChick: Property-Based Testing In Coq*. Electronic textbook.
- Fredrik Lindblad. 2007. Property Directed Generation of First-Order Test Data. In *8th Symposium on Trends in Functional Programming*. Intellect, 105–123.
- Andreas Löcher and Konstantinos Sagonas. 2018. Automating Targeted Property-Based Testing. In *2018 IEEE 11th International Conference on Software Testing, Verification and Validation*. 70–80. <https://doi.org/10.1109/ICST.2018.00017>
- Andreas Löscher and Konstantinos Sagonas. 2017. Targeted Property-Based Testing. In *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis*. ACM, 46–56. <https://doi.org/10.1145/3092703.3092711>
- David Maciver and Alastair F. Donaldson. 2020. Test-Case Reduction via Test-Case Generation: Insights from the Hypothesis Reducer (Tool Insights Paper). In *34th European Conference on Object-Oriented Programming (LIPIcs, Vol. 166)*. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 13:1–13:27. <https://doi.org/10.4230/LIPIcs.ECOOP.2020.13>
- David R. MacIver. 2016. Hypothesis: Property-Based Testing for Python. <https://hypothesis.works/>.
- Jan Midtgaard, Mathias Nygaard Justesen, Patrick Kasting, Flemming Nielson, and Hanne Riis Nielson. 2017. Effect-Driven QuickChecking of Compilers. *Proceedings of the ACM on Programming Languages* 1, ICFP, Article 15 (Aug 2017), 23 pages. <https://doi.org/10.1145/3110259>
- Agustín Mista and Alejandro Russo. 2019. Generating Random Structurally Rich Algebraic Data Type Values. In *Proceedings of the 14th International Workshop on Automation of Software Test*. IEEE Press, 48–54. <https://doi.org/10.1109/AST.2019.00013>
- Agustín Mista and Alejandro Russo. 2021. Deriving Compositional Random Generators. In *Proceedings of the 31st Symposium on Implementation and Application of Functional Languages*. ACM, Article 11, 12 pages. <https://doi.org/10.1145/3412932.3412943>
- Rickard Nilsson. 2019. ScalaCheck: Property-Based Testing for Scala. <https://scalacheck.org/>.
- Rohan Padhye, Caroline Lemieux, Koushik Sen, Mike Papadakis, and Yves Le Traon. 2019. Semantic Fuzzing with Zest. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis*. ACM, 329–340. <https://doi.org/10.1145/3293882.3330576>
- Michał H. Pałka, Koen Claessen, Alejandro Russo, and John Hughes. 2011. Testing an Optimising Compiler by Generating Random Lambda Terms. In *Proceedings of the 6th International Workshop on Automation of Software Test*. ACM, 91–97. <https://doi.org/10.1145/1982595.1982615>
- pandas. 2023. pandas - Python Data Analysis Library. <https://pandas.pydata.org/>
- Manolis Papadakis and Konstantinos F. Sagonas. 2011. A PropEr integration of types and function specifications with property-based testing. In *Proceedings of the 10th ACM SIGPLAN Workshop on Erlang*. 39–50. <https://doi.org/10.1145/>

2034654.2034663

- Zoe Paraskevopoulou, Aaron Eline, and Leonidas Lampropoulos. 2022. Computing Correctly with Inductive Relations. In *Proceedings of the ACM SIGPLAN Symposium on Programming Language Design and Implementation*. <https://doi.org/10.1145/3519939.3523707>
- Benjamin C. Pierce. 2018. *Software Foundations*. Electronic textbook.
- Plotly. 2023. Plotly: Low-Code Data App Development. <https://plotly.com/>
- Colin Runciman, Matthew Naylor, and Fredrik Lindblad. 2008. SmallCheck and Lazy SmallCheck: automatic exhaustive testing for small values. In *1st ACM SIGPLAN Symposium on Haskell*. ACM, 37–48. <https://doi.org/10.1145/1543134.1411292>
- Arvind Satyanarayan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. 2017. Vega-Lite: A Grammar of Interactive Graphics. *IEEE Transactions on Visualization and Computer Graphics* 23, 1 (Jan. 2017), 341–350. <https://doi.org/10.1109/TVCG.2016.2599030>
- Jacob Stanley. 2019. Hedgehog: Release with Confidence. <https://hackage.haskell.org/package/hedgehog/>.
- Kanit Wongsuphasawat, Zening Qu, Dominik Moritz, Riley Chang, Felix Ouk, Anushka Anand, Jock Mackinlay, Bill Howe, and Jeffrey Heer. 2017. Voyager 2: Augmenting Visual Analysis with Partial View Specifications. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. ACM, 2648–2659. <https://doi.org/10.1145/3025453.3025768>
- Li-yao Xia. 2018. A quick tour of generic-random. <https://hackage.haskell.org/package/generic-random-1.5.0.0/docs/Generic-Random.html>.
- Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. 2011. Finding and understanding bugs in C compilers. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, 283–294. <https://doi.org/10.1145/1993498.1993532>
- Zenong Zhang, Zach Patterson, Michael Hicks, and Shiyi Wei. 2022. FIXREVERTER: A Realistic Bug Injection Methodology for Benchmarking Fuzz Testing. In *31st USENIX Security Symposium*. USENIX Association, 3699–3715.

Received 2023-03-01; accepted 2023-06-27