

Looking for Lacunae in Bitcoin Core’s Fuzzing Efforts

Alex Groce
Northern Arizona University
United States

Kush Jain
Carnegie Mellon University
United States

Rijnard van Tonder
Sourcegraph, Inc.
United States

Goutamkumar Tulajappa
Kalburgi
Northern Arizona University
United States

Claire Le Goues
Carnegie Mellon University
United States

ABSTRACT

Bitcoin is one of the most prominent distributed software systems in the world. This paper describes an effort to investigate and enhance the effectiveness of the Bitcoin Core fuzzing effort. The effort initially began as a query about how to escape *saturation* in the fuzzing effort, but developed into a more general exploration. This paper summarizes the outcomes of a two-week focused effort. While the effort found no smoking guns indicating major test/fuzz weaknesses, it produced a large number of additional fuzz corpus entries, increased the set of fuzzers used for Bitcoin Core, and ran mutation analysis of Bitcoin Core fuzz targets, with a comparison to Bitcoin functional tests and other cryptocurrencies’ tests. Our conclusion is that for high quality fuzzing efforts, improvements to the *oracle* may be the best way to get more out of fuzzing.

CCS CONCEPTS

• **Software and its engineering** → **Dynamic analysis; Software testing and debugging.**

KEYWORDS

fuzzing, saturation, test diversity, mutation analysis, oracle strength

ACM Reference Format:

Alex Groce, Kush Jain, Rijnard van Tonder, Goutamkumar Tulajappa Kalburgi, and Claire Le Goues. 2022. Looking for Lacunae in Bitcoin Core’s Fuzzing Efforts. In *44nd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP ’22)*, May 21–29, 2022, Pittsburgh, PA, USA. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/3510457.3513072>

1 INTRODUCTION

Bitcoin [8] is the most popular cryptocurrency, and, while volatile, has a market cap consistently over half a trillion dollars since January of 2021. Bitcoin Core (<https://github.com/Bitcoin/Bitcoin>) is by far the most popular implementation, and serves as a reference for all other implementations of Bitcoin. To a significant degree, the code of Bitcoin Core is Bitcoin. Because of its fame and the high value of Bitcoins, Bitcoin is a high-value target for hackers.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ICSE-SEIP ’22, May 21–29, 2022, Pittsburgh, PA, USA

© 2022 Association for Computing Machinery.
ACM ISBN 978-1-4503-9226-6/22/05...\$15.00
<https://doi.org/10.1145/3510457.3513072>

Therefore, testing the code is of paramount importance, including extensive functional tests and aggressive *fuzzing*. This paper describes a focused effort to identify weaknesses in, and improve, the fuzzing of Bitcoin Core.

Chaincode Labs (<https://chaincode.com/>) is a private R&D center that exists solely to support and develop Bitcoin. In March of 2021 the head of special projects at Chaincode contacted the first author to discuss determining a strategy to improve the fuzzing of Bitcoin Core. It seemed that the fuzzing was “stuck”: neither code coverage nor found bugs were increasing with additional fuzzing. After some discussion, an 80 hour effort was determined as a reasonable scope for an external, research-oriented, look at the fuzzing effort.

Saturation, as defined in the blog post (<https://blog.regehr.org/archives/1796>) that brought Chaincode Labs to the first author, is when “We apply a fuzzer to some non-trivial system... [and] the number of new bugs found by the fuzzer drops off, eventually approaching zero.” At first a particular fuzzer applied to a system will tend to continuously increase both coverage and discovery of previously-unknown bugs. But, at some point, these bugs are known (and often fixed) and the fuzzer stops producing new bugs. Code and behavioral coverage seems to be *saturated*. The underlying reason for saturation is that any fuzzer (or other test generator) explores a space of generated tests according to some complex probability distribution. Some bugs lie in the high-probability portion of this space, and other bugs lie in very low probability zero probability parts. Escaping saturation may require a variety of approaches.

2 RESULTS

One thing that quickly emerged from discussions before the primary 80 hour effort began was the limited extent of the fuzzer runs being performed. The fuzzing includes a large number of targets, each with its own fuzz harness and executable. At the time, the basic strategy was to run libFuzzer on each of these for 100,000 iterations. Because some targets are very fast and a few, such as full message processing, are slow, this meant in practice fuzzing most targets for only 30-90 seconds, and even the slowest targets for only a little over an hour. The total time for over 100 targets was not negligible, but expecting such short runs for each target, after an initial exploration of the easy part of the probability space, to gain coverage or bugs very often, was simply unrealistic. For complex targets such as transaction verification and end-to-end message processing, 100,000 iterations was highly insufficient. The first suggestion for escaping saturation, therefore was very simple: run the fuzzer longer! The Chaincode tried increasing their configuration to 5 million iterations, multiplying the number of

