

Build It, Break It, Fix It: Contesting Secure

Andrew Ruef, Michael Hicks, James Parker, Dave Levin, Michelle L. M. University of Maryland

ABSTRACT

Typical security contests focus on breaking or mitigating the impact of buggy systems. We present the Build-it, Break-it, Fix-it (BIBIFI) contest, which aims to assess the ability to securely build software, not just break it. In BIBIFI, teams build specified software with the goal of maximizing correctness, performance, and security. The latter is tested when teams attempt to break other teams' submissions. Winners are chosen from among the best builders and the best breakers. BIBIFI was designed to be open-ended—teams can use any language, tool, process, etc. that they like. As such, contest outcomes shed light on factors that correlate with successfully building secure software and breaking insecure software. During 2015 we ran three contests involving a total of 116 teams and two different programming problems. Quantitative analysis from these contests found that the most efficient build-it submissions used C/C++, but submissions coded in a statically-typed language were less likely to have a security flaw; build-it teams with diverse programming-language knowledge also produced more secure code. Shorter programs correlated with better scores. Break-it teams that were also build-it teams were significantly better at finding security bugs.

1. INTRODUCTION

Capture-the-flag (CTF) and other cybersecurity contests [1–4] are popular and serve as valuable proving grounds for finding cybersecurity talent. These contests largely focus on *breaking* (e.g., exploiting vulnerabilities and/or misconfigurations) and *mitigation* (e.g., rapid patching or reconfiguration). They do not, however, test contestants' ability to *build* (i.e., design and implement) systems that are secure in the first place. Typical programming contests [5–7] do focus on design and implementation, but generally ignore security. This state of affairs is unfortunate because experts have long advocated that achieving security in a computer system requires treating security as a first-order design goal [8], not something that can be added after the

other teams' build-it submissions via test cases, benefit a build-it team's score by exposing the build-it team's security-relevant problems. (Break-it scores are independent, one per category.) The final fix bugs and thereby get points for each distinct break-it that is fixed. BIBIFI's structure and scoring are designed to encourage meaningful participation so that the top-scoring build-it teams produce efficient software. Behavior that does not contribute to the competition is discouraged. For example, teams can submit a limited number of submissions per mission, and will lose points if their submission exposes the same underlying bug as another submission by other teams. As such, teams are encouraged to bugs broadly (in many submissions) rather than focusing on hard-to-find bugs).

In addition to providing a platform for testing and improving the security of software, BIBIFI presents an opportunity to study the software development process. Scientifically valid tests may serve as a quasi-experiment to study how participation in the contest relates to the quality of the artifacts and participant experience. For example, we can study the choice of build-it programming languages, the size of the build-it team, the experience of the build-it team members, the size of the build-it team, the success rate of the build-it team, and the success rate of the break-it team. To the extent that contestants represent a cross-section of the cybersecurity community, the results of this study may provide evidence for practices that help to improve the security of software. Indeed, the contest environment provides a space for sharing ideas to improve development processes and for making them available to practitioners.

This paper studies the outcome of the first three contests that we held during 2015, including the first ever build-it contest. The first contest was a *secure, append-only log file* challenge.

Build It, Break It, Fix It: Contesting Secure

Andrew Ruef, Michael Hicks, James Parker, Dave Levin, Michelle L. M. M. University of Maryland

ABSTRACT

Typical security contests focus on breaking or mitigating the impact of buggy systems. We present the Build-it, Break-it, Fix-it (BIBIFI) contest, which aims to assess the ability to securely build software, not just break it. In BIBIFI, teams build specified software with the goal of maximizing correctness, performance, and security. The latter is tested when teams attempt to break other teams' submissions. Winners are chosen from among the best builders and the best breakers. BIBIFI was designed to be open-ended—teams can use any language, tool, process, etc. that they like. As such, contest outcomes shed light on factors that correlate with successfully building secure software and breaking insecure software. During 2015 we ran three contests involving a total of 116 teams and two different programming problems. Quantitative analysis from these contests found that the most efficient build-it submissions used C/C++, but submissions coded in a statically-typed language were less likely to have a security flaw; build-it teams with diverse programming-language knowledge also produced more secure code. Shorter programs correlated with better scores. Break-it teams that were also build-it teams were significantly better at finding security bugs.

1. INTRODUCTION

Capture-the-flag (CTF) and other cybersecurity contests [15, 16, 5, 18, 6] are popular and serve as valuable proving grounds for finding cybersecurity talent. These contests largely focus on *breaking* (e.g., exploiting vulnerabilities and/or misconfigurations) and *mitigation* (e.g., rapid patching or reconfiguration). They do not, however, test contestants' ability to *build* (i.e., design and implement) systems that are secure in the first place. Typical programming contests [25, 2, 12] do focus on design and implementation, but generally ignore security. This state of affairs is unfortunate because experts have long advocated that achieving security in a computer system requires treating security as a first-order design goal [22], not something that can be added after the fact. As such, we should not assume that good breakers will necessarily be good builders [14], or that top coders can

other teams' build-it submissions via test cases, benefit a build-it team's score by size the build-it team's score by security-relevant problems. (Break-it scores are independent, one per category.) The final fix bugs and thereby get points that distinct break-it teams receive. BIBIFI's structure and scoring are designed to encourage meaningful participation that the top-scoring build-it teams produce efficient software. Behavior that comes are discouraged. For example, teams submit a limited number of submissions per mission, and will lose points if they expose the same underlying bugs to multiple teams by other teams. As such, the contest rewards bugs broadly (in many submissions) over hard-to-find bugs).

In addition to providing a platform for testing BIBIFI presents an opportunity for improving the breaking process scientifically. The contests may serve as a quasi-experiments to relate participation data with various artifacts and participant characteristics. The choice of build-it programming languages, experience, code size, testing tools, and the team's (non)success in the contests. To the extent that contestants represent the wider community, the results of this study may provide evidence for practices that have been adopted. Indeed, the contest environment has provided ideas to improve development processes and making their way to practical use.

This paper studies the outcome of the three contests that we held during 2015, including the 2015 and 2016 contests. The first contest involved building a secure, *append-only* log file that could be generated by a hypothetical application and checked by attackers with direct access to the log file.

the three contests drew participants from a MOOC (Massive Online Open Courseware) course on cybersecurity. These participants (278 total, comprising 109 teams) each had an average of 10 years of programming experience and had just completed a four-course sequence including courses on secure software and cryptography. The third contest involved U.S.-based graduate and undergraduate students (23 total, comprising 6 teams) with less experience and training.

Rigorous quantitative analysis of the outcomes revealed several interesting, statistically significant effects. Considering build-it scores: Writing code in C/C++ increased build-it scores initially, but also increased chances of a security bug found later. Teams with broader programming language knowledge and that wrote less code also produced more secure implementations. Considering break-it scores: larger teams found more bugs during the break-it phase. Greater programming experience, and knowledge of C, were also helpful. Break-it teams that also participated during the build-it phase were significantly more likely to find a security bug than those that did not.

We manually examined both build-it and break-it artifacts. Successful build-it teams typically employed third-party libraries—e.g., SSL, NaCL, and BouncyCastle—to implement cryptographic operations and/or communications, which freed up worry of proper use of randomness, nonces, etc. Unsuccessful teams typically failed to employ cryptography, implemented it incorrectly, used insufficient randomness, or failed to use authentication; there were comparatively few memory safety violations. Break-it teams found clever ways to exploit security problems; some MITM implementations were quite sophisticated.

In summary, this paper makes two main contributions. First, it presents BIBIFI, a new security contest that encourages building, not just breaking. Second, it presents a detailed description of three BIBIFI contests along with both a quantitative and qualitative analysis of the results. We will be making the BIBIFI code and infrastructure publicly available so that others may run their own competitions; we hope that this opens up a line of research built on empirical experiments with secure programming methodologies.¹

The rest of this paper is organized as follows. We present the design of BIBIFI in §2 and describe specifics of the contests we ran in §3. We present the quantitative analysis of the data we collected from these contests in §4, and qualitative analysis in §5. We review related work in §6 and conclude in §7.

2. BUILD-IT, BREAK-IT, FIX-IT

This section describes the goals, design, and implementation of the BIBIFI competition. At the highest level, our aim is to create an environment that closely reflects real-world development goals and constraints, and to encourage build-it teams to write the most secure code they can, and break-it teams to perform the most thorough, creative analysis of others’ code they can. We achieve this through a

2.1 Competition pha

We begin by describing the process that occurs during a BIBIFI competition. BIBIFI is administered on-line, rather than being physically distributed. The contests consist of three phases, each of which last about two weeks, described in this paper.

BIBIFI begins with the build-it phase, where contestants aim to implement a set of challenges according to a published specification provided by administrators. A suitable challenge is completed by good programmers within two weeks, for the contests are designed to reward for performance, and has an upper bound. The software should have specific requirements, such as not leaking private information or crashing. Software compromised by poor design or bugs should also not be accepted. It is important to ensure that contestants do not cheat, such as still taking advantage of hardware or network resources to the extent possible. The build-it phase runs on a standard Linux VM provided at the start of the contest. Teams must push their code to a central server, each push, the contest infrastructure automatically pulls it, builds it, tests it (for correctness), and updates the scoreboard. There are two types of problems we developed: (1) a pair of communicating programs that must interact over an ATM.

The next phase is the break-it phase, where contestants can download, build, and inspect the build-it artifacts. They must complete missions, including source code analysis, and must build properly, pass automated tests, and fix purposefully obfuscated. We provide a detailed view of the build-it teams’ submissions, including meta-data like program name, file size, and time taken. When they have found a defect, break-it teams must report it by posing the defect and an explanation. They must also provide coverage, a break-it team must provide a certain number of test cases per build. The break-it infrastructure automatically judges each test case. If a test case truly reveals a defect, the break-it team must fix it. If a bug is found in the build-it code, it will run the test again against the build (the “oracle”) and the target build must pass on the former build for the bug to be accepted.³ More points are awarded for fixing multiple bugs, which may be demonstrated in different ways. The auto-judgment also handles two different contest problems: one for each type of problem.

The final phase is the fix-it phase, where contestants are provided with the bug report and asked to fix their submission. They may submit multiple fixes, and the system will automatically identify them. If a single fix corrects all the bugs in the build, the test cases are “morally” accepted. Points are only deducted for one of the fixes if it does not fix all the bugs. This provides an incentive for contestants to fix all the bugs in their submission, even if they have already submitted a fix that is accepted.

ysis of others' code they can. We achieve this through a careful design of how the competition is run and how various acts are scored (or penalized).

¹This paper subsumes a previously published 6-page workshop paper [1]. The initial BIBIFI design and implementation also appeared in that paper, as did a brief description of a pilot run of the contest. This paper presents many more details about the contest setup along with a quantitative and qualitative analysis of the outcomes of several larger contests.

based on information provided by the user, or other assessment, whether a bug is fixed or not. In this sense that it corrects only bugs.

²This avoids spurious unfair advantages by not running investigating code in the order specified.

³To encourage testing of our system, we will let users earn points by finding bugs in their own code.

fix is rejected.

Once the final phase concludes, prizes are awarded to the best builders and best breakers as determined by the scoring system described next.

2.2 Competition scoring

BIBIFI’s scoring system aims to encourage the contest’s basic goals, which are that the winners of the Build-it phase truly produced the highest quality software, and that the winners of the break-it phase performed the most thorough, creative analysis of others’ code. The scoring rules create incentives for good behavior (and disincentives for bad behavior).

2.2.1 Build-it scores

To reflect real-world development concerns, the winning build-it team would ideally develop software that is correct, secure, and efficient. While security is of primary interest to our contest, developers in practice must balance these other aspects of quality against security [§4.5], leading to a set of trade-offs that cannot be ignored if we wish to understand real developer decision-making.

To encourage these, each build-it team’s score is the sum of the *ship* score⁴ and the *resilience* score. The ship score is composed of points gained for correctness tests and performance tests. Each mandatory correctness test is worth M points, for some constant M , while each optional correctness test is worth $M/2$ points. Each performance test has a numeric measure depending on the specific nature of the programming project—e.g., latency, space consumed, files left unprocessed—where lower measures are better. A test’s worth is $M \cdot (\text{worst} - v) / (\text{worst} - \text{best})$, where v is the measured result, *best* is the measure for the best-performing submission, and *worst* is the worst performing. As such, each performance test’s value ranges from 0 to M .

The resilience score is determined after the break-it and fix-it phases, at which point the set of unique defects against a submission is known. For each *unique* bug found against a team’s submission we subtract P points from its resilience score; as such, the best possible resilience score is 0. For correctness bugs, we set P to $M/2$; for crashes that violate memory safety we set P to M , and for exploits and other security property failures we set P to $2M$.

2.2.2 Break-it scores

Our primary goal with break-it teams is to encourage them to find as many defects as possible in the submitted software, as this would give greater confidence in our assessment that a particular build-it team’s software is of higher quality than another’s. While we are particularly interested in obvious security defects, correctness defects are also important, as they can have non-obvious security implications.

After the break-it phase, a break-it team’s score is the summed value of all defects they have found, using the above formula. After the fix-it phase, this is multiplied by L .

spirit of the competition. Provide test cases as evidence to ensure they are providing useful bugs. Give 4x more points for security bugs than correctness bugs, so teams to look for these sorts of bugs. Encourage teams to focus on correctness issues. (But also reward them for finding security bugs; see below.) Because break-it teams are limited to a fixed number of test cases per submission, it is encouraged from submitting many low-hanging fruit bugs that are the same; as they could likely find many more in the fix-it phase they are better off spending time fixing different bugs. Limiting the number of test cases encourages examining many different bugs. Giving more points for defects found by break-it teams than build-it teams are encouraged to look for more bugs than just low-hanging fruit.

2.2.3 Limitations

While we believe BIBIFI’s scoring system and incentives are properly designed, there are several potential limitations.

First, there is no guarantee that all bugs will be found. Break-it teams may lack the skill to find problems in all submissions. Some submissions may receive equal scrutiny, even though they act contrary to incentives and find the same or even more duplicated bugs, rather than unique ones. Finally, break-it teams may lack the ability to report all bugs. If the infrastructure cannot automatically report all bugs, some defects will go unreported. Given the large pool of break-it teams, and the difficulty of automation, we still believe that BIBIFI is effective in breadth and depth.

Second, builders may fail to produce high-quality software in their best interests. For example, a builder may intentionally produce a fix rejected as addressing non-security bugs. Break-it teams may use several specific bugs that they know would have been allowed. A builder may also have a desire for contention for prizes made available during the fix-it phase.⁵ We observed this behavior in the break-it phase described in §4.5. Both acting on the builder’s desire to win the resilience score (and correspondingly the overall competition score) and mitigating these issues will be challenging. For example, we can mitigate these issues with a tiebreaker rule, e.g., by offering prizes to all teams that find the same number of bugs in their final score, rather than the one with the highest score.

Finally, there are several other potential limitations. The current definition of a “bug” in BIBIFI’s scoring system may skew results. For example, correctness tests may leave too many bugs unreported during break-it (distracting from security bugs); too many bugs are reported (meaning teams are difficult to distinguish based on their bug-finding ability). Some teams may have a ratio of 4 to 1 over correctness bugs to security bugs, while others have a ratio of 1 to 1. What ratio makes the most sense for real-world outcomes; both extremes are undesirable.

P valuations. After the fix-it phase, this score is reduced. In particular, each of the N break-it teams' scores that identified the same defect are adjusted to receive P/N points for that defect, splitting the P points among them.

Through a combination of requiring concrete test cases and scoring, BIBIFI encourages break-it teams to follow the

⁴The name is meant to evoke a quality measure at the time software is shipped.

argued. Finally, performance

⁵Hiding scores during the contest would harm incentives during contests with no bugs reported against

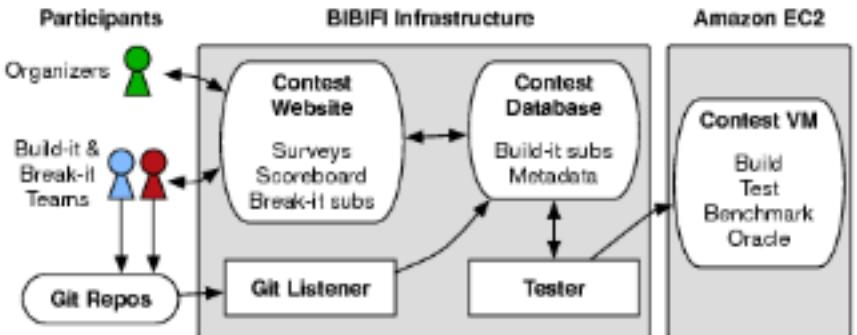


Figure 1: Overview of BIBIFI’s implementation.

test [] prioritized security bugs 2-to-1 and had fewer interesting performance tests, and outcomes were better when we increased the ratio.

2.2.4 Discouraging collusion

BIBIFI contestants may form teams however they wish, and may participate remotely. This encourages wider participation, but it also opens the possibility of collusion between teams, as there cannot be a judge overseeing their communication and coordination. There are three broad possibilities for collusion, each of which BIBIFI’s scoring discourages.

First, two break-it teams could consider sharing bugs they find with one another. By scaling the points each finder of a particular bug obtains, we remove incentive for them to both submit the same bugs, as they would risk diluting how many points they both obtain.

The second class of collusion is between a build-it team and a break-it team, but neither have incentive to assist one another. The zero-sum nature of the scoring between breakers and builders places them at odds with one another; revealing a bug to a break-it team hurts the builder, and not reporting a bug hurts the breaker.

Finally, two build-it teams could collude, for instance by sharing code with one another. It might be in their interests to do this in the event that the competition offers prizes to two or more build-it teams, since collusion could obtain more than one prize-position. We use judging and automated tools (and feedback from break-it teams) to detect if two teams share the same code (and disqualify them), but it is not clear how to detect whether two teams provided out-of-band feedback to one another prior to submitting code (e.g., by holding their own informal “break-it” and “fix-it” stages). We view this as a minor threat to validity; at the surface, such assistance appears unfair, but it is not clear that it is contrary to the goals of the contest, that is, to developing secure code.

2.3 Implementation

Figure 1 provides an overview of the BIBIFI implementation. It consists of a web frontend, providing the interface

event to the contest outcome and security training). During qualification tests build-it submissions keeps the current scores up-to-date for the judges for considering whether a submission covers one bug or not.

To secure the web application from participants, we implemented it using the Yesod [29] web framework, using Haskell’s strong type system to detect after-free, buffer overrun, and SQL injection attacks. The use of Yesod’s built-in protection against various attacks is sufficient for our application. As one further layer of security, the application incorporates the information monad (LMonad [30]), which is derived from the State monad to protect against inadvertent modifications to state, such as escalations. LMonad dynamically tracks what users can only access their own information.

Testing backend. The backend performs automated testing during the build-it phase, to verify performance, and during the break-it phase, to verify vulnerabilities. It consists of a virtual machine (and a little Python).

To automate testing, we use a combination of a URL to a Git [31] repository, a shared Bitbucket [32] account, and shared with the tester a read-only URL to the repository for pushes, upon which the tester can commit. Testing is then handled by a script that sets up an Amazon EC2 virtual machine to run each submission. We require the tester to run the tests without any network access to the tester’s machine, so that we share in advance. The tester does not need to install additional packages nor does the EC2 virtual machine support both scalability and dynamism (as dynamically provisioned) and portability. The fact that instances prevents a team from running multiple tests, or of tests on other teams’ submissions.

All qualifying build-it submissions are run by break-it teams at the same time. Break-it teams identify bugs in the build-it submission by file specifying the buggy source code and a file containing a series of commands with expected output for each bug. Break-it teams commit the results to the tester’s repository. The backend uses the tester’s repository to identify implicated submission to score.

3. CONTEST PROBLEMS

This section presents the problems developed for the contests held by BIBIFI. It describes the specific notions of security defects that were used, and how such defects were tested.

team. It consists of a web frontend, providing the interface to both participants and organizers, and a backend for testing builds and breaks. Key goals of the infrastructure are security—we do not want participants to succeed by hacking BIBIFI itself—and scalability.

Web frontend. Contestants sign up for the contest through our web application frontend, and fill out a survey when doing so, to gather demographic and other data potentially rel-

3.1 Secure log (Spring)

The secure log problem was first proposed for a virtual art gallery security system. The system maintains two logs of events. The first, `logappend`, appends events to a log. These events indicate when employees arrive at and leave rooms. The second, `loggrep`, allows users to search for events. To qualify, submissions must answer several queries (involving the current state of the system) correctly.

movements of particular people), but they could implement two more for extra points (involving time spent in the museum, and intersections among different peoples’ histories). An empty log is created by `logappend` with a given authentication token, and later calls to `logappend` and `logread` on the same log must use that token or deny the requests.

A canonical way of implementing the secure log is to treat the authentication token as a symmetric key for authenticated encryption, e.g., using a combination of AES and HMAC. There are several tempting shortcuts that we anticipated build-it teams would take (and that break-it teams would exploit). For instance, one may be tempted to encrypt/MAC individual log records as opposed to the entire log, thereby making `logappend` faster. But this could permit integrity breaks that duplicate or reorder log records. Generally speaking, teams may also be tempted to implement their own encryption rather than use existing libraries, or to simply sidestep encryption altogether. §5 reports several cases we observed.

A submission’s performance was measured in terms of time to perform a particular sequence of operations, and space consumed by the resulting log. Correctness (and *crash*) bug reports were defined as sequences of `logread` and/or `logappend` operations with expected outputs; these were vetted by our oracle. Security was defined by *privacy* and *integrity*: any attempt to learn something about the log’s contents, or to change them, without the use of the `logread` and `logappend` and the proper token should be disallowed. How violations of these properties were specified and tested is described next.

Privacy breaks. When providing a build-it submission to the break-it teams, we also included a set of log files that were generated using a sequence of invocations of that submission’s `logappend` program. We generated different logs for different build-it submissions, using a distinct command sequence and authentication token for each. All logs were distributed to break-it teams without the authentication token; some were distributed without revealing the sequence of commands (the “transcript”) that generated them. For these, a break-it team could submit a test case involving a call to `logread` (with the authentication token omitted) that queries the file. The BIBIFI infrastructure would run the query on the specified file with the authentication token, and if the output matched that specified by the breaker, then a privacy violation is confirmed.

Integrity breaks. For about half of the generated log files we also provided the transcript of the `logappend` operations (*sans* auth token) used to generate the file. A team could submit a test case specifying the name of the log file, the contents of a corrupted version of that file, and a `logread` query over it (without the authentication token). For both the specified log file and the corrupted one, the BIBIFI infrastructure would run the query using the correct authenti-

vacy break against the score, even if there are multiple ones. To produce privacy/integrity breaks, the break-it teams automatically tell them apart.

3.2 Securing ATM interactions

The ATM problem asked build-it teams to implement communicating programs: `atm` and `bank`. `atm` allows customers to set up an account, withdraw money, and deposit money, while `bank` is a server that maintains accounts, tracking bank balances. `atm` interacts with `bank` to verify a customer with a correct password and to check the balance of their account, and `bank` interacts with `atm` to verify (e.g., they may not withdraw more than they have). In addition, `atm` and `bank` can authenticate each other for the purpose of setting up a session. This channel unavailable to the `atm` and `bank` programs, so that communicating with `bank` over this channel (“middle” (MITM)) could observe and modify messages or insert new messages. These security features are described despite not having a formal specification.

A canonical way of implementing the security features of the `atm` and `bank` programs would be to use public key cryptography. `atm` uses the `bank`’s public key to ensure that the `bank`’s responses to the `atm` are valid. `bank` uses the `atm`’s public key to ensure that the `atm` is with the `bank` and not a MITM. To withdraw money from an account, the card file should contain the card number, so that the MITM cannot withdraw money from the wrong account. It is also necessary to protect the communication channel using nonces or similar mechanisms. A similar wise approach would be use of a shared secret key to implement these features. Both approaches are discussed further in §5.

Build-it submissions’ performance was measured in terms of time to complete a series of operations involving `atm/bank` interactions.⁷ Correctness bug reports were defined as sequences of operations that a build-it submission produces correctly (or crashes). Security defects were defined as sequences of operations that a build-it submission produces incorrectly (or crashes).

Integrity breaks. Integrity breaks were implemented by running a custom MITM program on a specified IP address and port. The MITM program intercepts all connection from the `atm` when the `atm` connects to the `bank`. The MITM program can thus intercept the communications between `atm` and `bank`, and either modify or initiate its own. We provide a starter MITM: It sets up a TCP socket to intercept the communications between the `atm` and `bank`.

To demonstrate an integrity break, a build-it submission sends requests to a *command server* to generate a transcript of inputs on the `atm` and it can then submit a transcript to the `bank` account whose creation it initiated. The `bank` will declare the test completed if the transcript matches the one provided by the build-it submission.

cation token. An integrity violation is detected if the query command produces a non-error answer for the corrupted log that differs from the correct answer (which can be confirmed against the transcript using the oracle).

This approach to determining privacy and integrity breaks has the benefit and drawback that it does not reveal the *source* of the issue, only that there is (at least) one. As such, we only count up to one integrity break and one pri-

atm commands is run using *the MITM*. This means that sends directly to the target s

⁶In a real deployment, this mi file into the ATM's ROM prior

⁷This transcript was always s vation to support parallelism t

⁸All submissions were required

be replayed/sent to the oracle. If the oracle and target both complete the command list without error, but they differ on the outputs of one or more commands, or on the balances of accounts at the bank whose card files were not revealed to the MITM during the test, then there is evidence of an integrity violation.

As an example (based on a real attack we observed), consider a submission that uses deterministic encryption without nonces in messages. The MITM could direct the command server to withdraw money from an account, and then replay the message it observes. When run on the vulnerable submission, this would debit the account twice. But when run on the oracle without the MITM, no message is replayed, leading to differing final account balances. A correct submission would reject the replayed message, which would invalidate the break.

Privacy breaks. Privacy violations are also demonstrated using a MITM. In this case, at the start of a test, the command server will generate two random values: “amount” and “account name.” If by the end of the test no errors have occurred and the attacker can prove it knows the actual value of either secret (by sending a command that specifies it), the break is considered successful. Before demonstrating knowledge of the secret, the MITM can send commands to the server with a *symbolic* “amount” and “account name”; the server will fill in the actual secrets before forwarding these messages. The command server does not automatically create a secret account or an account with a secret balance; it is up to the breaker to do that (referencing the secrets symbolically when doing so).

As an example, suppose the target does not encrypt exchanged messages. Then a privacy attack might be for the MITM to direct the command server to create an account whose balance contains the secret amount. Then the MITM can observe an unencrypted message sent from `atm` to `bank`; this message will contain the actual amount, filled in by the command server. The MITM can then send its guess to the command server showing that it knows the amount.

As with the log problem, we cannot tell whether an integrity or privacy test is exploiting the same underlying weakness in a submission, so we only accept one violation against each submission.

Timeouts and denial of service. One difficulty with our use of a MITM is that we cannot reliably detect bugs in submissions that would result in infinite loops, missed messages, or corrupted messages. This is because such bugs can be simulated by the MITM by dropping or corrupting messages it receives. Since the builders are free to implement any protocol they like, our auto-testing infrastructure cannot tell if a protocol error or timeout is due to a bug in the target or due to misbehavior of the MITM. As such, we conservatively disallow reporting any such errors, meaning that we may miss some: i.e., flaws in builder implementa-

performance in each phase factors contribute to high solving by other teams, and so on. We find that on average, teams other than C and C++, and more programming languages (less programming skill), are less satisfied in their code. Success in identifying security bugs in submissions with having more team members involved in the build-it phase (and less time to how to secure an implementation) in a contest, which used the ATT&T challenge, more security bugs than the others.

4.1 Data collection

For each team, we collected self-reported data. When signing up for the challenge, participants reported standard demographic information, including experience and programming skill. During the challenge, each team member was asked to rate their own performance. Finally, we collected information about lines of code written and pulled from teams’ Git repositories.

Participant data was anonymized and approved by our institution’s Institutional Review Board. Participants consented to having their data collected, anonymized, stored, and analyzed. Participants did not consent to sharing their personal data with third parties.

4.2 Analysis approach

To examine factors that contribute to solving and breaking, we apply regression models. These regression model attempts to explain the outcome based on one or more measured factors. To predict the participants’ scores, we can use a linear regression model, which estimates how many points a participant gains or loses due to (or takes away from) a particular factor. To predict binary outcomes, such as whether a team solved a challenge or not, we apply logistic regression, which estimates the likelihood that each factor impacts the likelihood of success.

We consider many variables that may impact teams’ results. To avoid overfitting, we consider potential factors those variables that are likely to be of interest, within acceptable ranges of values. (Our choices are detailed below.) We then fit several regression models with all possible combinations of variables and select the model with the minimum Akaike information criterion (AIC). Only the top three models are presented here.

We describe the results of the challenge as a whole, not a completely controlled experiment (e.g., random assignment), so our results are correlations rather than causation. Our results are subject to confounds, and many factors are likely to be correlated.

tions might exist but evidence of those bugs might not be realizable in our testing system.

4. QUANTITATIVE ANALYSIS

This section analyzes data we have gathered from three contests we ran during 2015.⁹ We consider participants'

⁹We also ran a contest during Fall 2014 but exclude it from con-

in our data are correlated also assumes that the factor on participants' scores (or or while this may not be the simplification for considerin

4.3 Contestants

sideration due to differences in

Contest	USA	India	Russia	Brazil	Other
Spring 2015	30	7	12	12	120
Fall 2015	64	14	12	20	110

Table 1: Contestants, by self-reported country.

Contest	Spring 15 [†]	Fall 15 [†]	Fall 15
# Contestants	156	122	23
% Male	91%	89%	100%
% Female	5%	9 %	0 %
Age (mean/min/max)	34.8/20/61	33.5/19/69	25.1/17/31
% with CS degrees	35 %	38 %	23 %
Years programming	9.6/0/30	9.9/0/37	6.6/2/13
# Build-it teams	61	34	6
Build-it team size	2.2/1/5	3.1/1/5	3.1/1/6
# Break-it teams (that also built)	65 (58)	39 (32)	4 (3)
Break-it team size	2.4/1/5	3.0/1/5	3.5/1/6
# PLs known per team	6.8/1/22	10.0/2/20	4.2/1/8

Table 2: Demographics of contestants from qualifying teams.
[†] indicates MOOC participants. Some participants declined to specify gender.

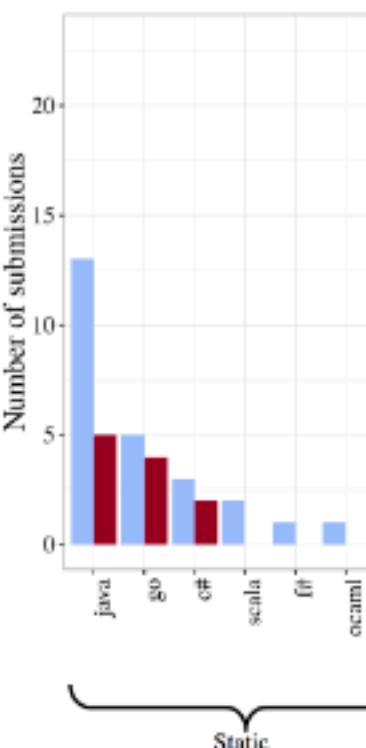


Figure 2: The number of submissions for each language in the test, organized by primary language. The brackets group the languages that are statically typed.

We consider three contests offered at two times:

Spring 2015: We held one contest during May–June 2015 as the capstone to a Cybersecurity MOOC sequence.¹⁰ Before completing in the capstone, participants passed courses on software security, cryptography, usable security, and hardware security. The contest problem was the secure log problem (§3.1).

Fall 2015: During Oct.–Nov. 2015 we offered two contests simultaneously, one as a MOOC capstone, and the other open to U.S. college students. We merged the contests after the build-it phase, due to low participation in the open contest; from here on we refer to these two as a single contest. The contest problem was the ATM problem (§3.2).

The majority of all of our contestants came from the U.S. There was representation from developed countries with a reputation both for high technology and hacking acumen. Details of the most popular countries of origin can be found in Table 1, and additional information about contestant demographics is presented in Table 2.

4.4 Ship scores

We first consider factors correlating with a team’s *ship* score, which assesses their submission’s quality before it is attacked by the other teams (§2.1). This data set contains all

Model setup. To ensure enough power to detect meaningful relationships, we decided to use a sample size roughly equivalent to $n = 150$, with a power of $1 - \beta = 0.8$. An effect size of $\eta^2 = 0.15$ is considered large enough to explain up to 13% of the variance in the outcome variable. With an assumed power of 0.8 and a significance level of 0.05, we limited ourselves to nine independent variables, resulting in a prospective $f^2 = 0.154$. (One variable, the number of contestants, was excluded.) Our model is reported with the results of the full model. Since this limit, we selected the following variables:

Contest: Whether the team participated in the Spring 2015 contest or the Fall 2015 contest.

Team members: A team’s total number of members.

Knowledge of C: The team member who claimed to know C or C++ as a programming language. We included this variable as a proxy for the team’s knowledge of C/C++ implementation details, a skill that is important for successful secure building.

Languages known: The number of programming languages team members claimed to know (see the last row of Table 2). For example, if member A claims to know C and member B claims to know Python, the language count would be 2.

Coding experience: The years of coding experience reported by a team member.

Lines of code: The SLLOC of the team’s final submission.

Team size: The number of team members.

101 teams from the Spring 2015 and Fall 2015 contests that qualified after the build-it phase. Both contests have nearly the same number of correctness and performance tests, but different numbers of participants. We set the constant multiplier M to be 50 for both contests, which effectively normalizes the scores.

¹⁰<https://www.coursera.org/specializations/cyber-security>

Language category: We submission as having one “guages were then assigned to statically-typed (e.g., Java, C, C++) or dynamically-typed (e.g., Perl, Python). We also created a category for the regression. Table 1 shows the total submissions for each language category.

¹¹<http://www.dwheeler.com/skewness.html>

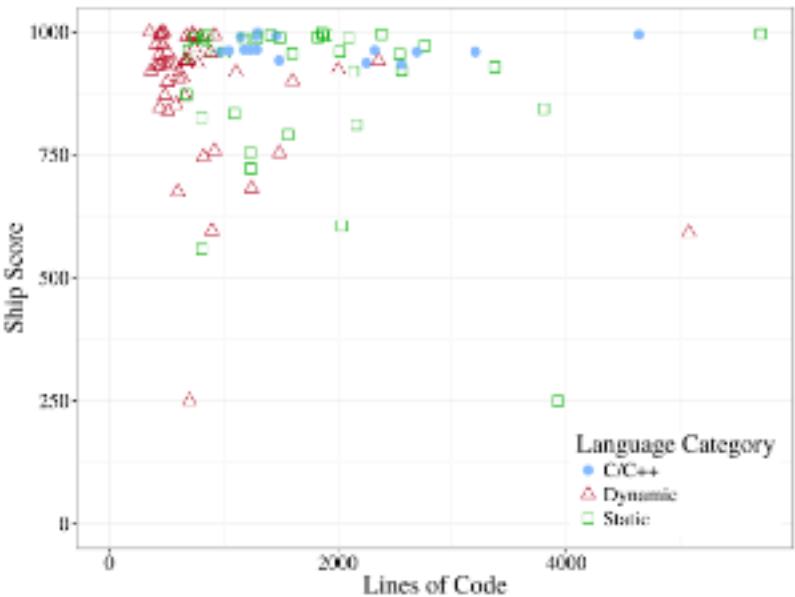


Figure 3: Each team’s ship score, compared to the lines of code in its implementation and organized by language category. Fewer LOC and using C/C++ correlate with a higher ship score.

are given in Figure 2.

MOOC: True if the team was participating in the MOOC capstone project; otherwise false.

Results. Our regression results (Table 3) indicate that ship score is strongly correlated with language choice. Teams that programmed in C or C++ performed on average 121 and 92 points better than those who programmed in dynamically typed and statically typed languages, respectively. Figure 3 illustrates that while teams in many language categories performed well in this phase, only teams that did not use C or C++ scored poorly.

The high scores for C/C++ teams could be due to better scores on performance tests and/or due to implementing optional features. We confirmed the main cause is the former. Every C/C++ team for Spring 2015 implemented all optional features, while six teams in the other categories implemented only 6 of 10 and one team implemented none; the Fall 2015 contest offered no optional features. We artificially increased the scores of those seven teams as if they had implemented all optional features and reran the regression model. The resulting model had very similar coefficients.

Our results also suggest that teams who were associated with the MOOC capstone performed 119 points better than non-MOOC teams. MOOC participants typically had more programming experience and CS training, so perhaps that is the reason.

Finally, we found that each additional line of code in a team’s submission was associated with a drop of 0.03 points in ship score. Based on our qualitative observations (see §5),

Factor
Fall 2015
Lines of code
Dynamically typed
Statically typed
MOOC

Table 3: Final linear regression indicating how many points total score. Overall effect size

the largest minimum size).

4.5 Code quality measures

Now we turn to measures of code quality. In particular, in terms of its correctness and robustness, a program ends up under scrutiny by breakers.

Resilience. The total build resilience score is the sum of just discussed, and *resilience* is a measure of the total resilience score that derives from breaking a program and measuring the presence of defects. Breakers are assigned to teams during the fix-it phase, as fix-it bugs are reported to them. Breakers are assigned to teams that identify the same set of bugs.

Unfortunately, upon studying the data, we found that a large percentage of build-it submissions had no bugs reported against their code, forcing breakers to spend time doing so. We can see this in Figure 4, which shows resilience scores (Y-axis) of the two contests. The circles represent teams that fixed at least one bug, whereas the squares represent teams that fixed no bugs. We can see that the resilience scores of teams with the lower resilience scores are higher. We further confirmed that this was not a statistical artifact by running a regression analysis that excluded fix-it phase participants.

This situation is disappointing because resilience is often used as a good measure of a program’s quality (e.g., resilience score as a good measure of a program’s quality). Our hypothesis was that teams that are sufficiently incentivized to fix bugs will do so. However, we found that teams that are sufficiently incentivized to fix bugs because they are graded on them are less likely to fix bugs because they are not graded on them. This is because for MOOC students, once a bug is fixed, they are assured to pass; it makes little sense to fix bugs that are not necessary for attaining this goal. This is in contrast to grading alternative structures that reward teams for fixing all bugs, even if they are not graded on them.

Presence of security bugs. We found that the presence of security bugs is not sufficiently meaningful, a user of the system can be affected by this. We mean any submission that contains a security bug, such as a crash, privacy, or integrity violation.

in ship score. Based on our quantitative observations (see 30), we hypothesize this may relate to more reuse of code from libraries, which are not counted in a team's LOC (most libraries were installed directly on the VM, not in the submission itself). As shown in Figure 3, LOC is also (as expected) associated with the category of language being used. While LOC varied widely within each language type, dynamic submissions were generally shortest, followed by static submissions and then those written in C/C++ (which has

crash, privacy, or integrity model we used logistic regression as the ship model.

Table 4 lists the results of the logistic regression. The coefficients represent the change in the log-odds of success with each factor. Negative coefficients indicate the likelihood of finding a security bug is lower than the baseline, while positive coefficients indicate it is higher. The exponential of the coefficient indicates how strongly that factor influences the outcome.

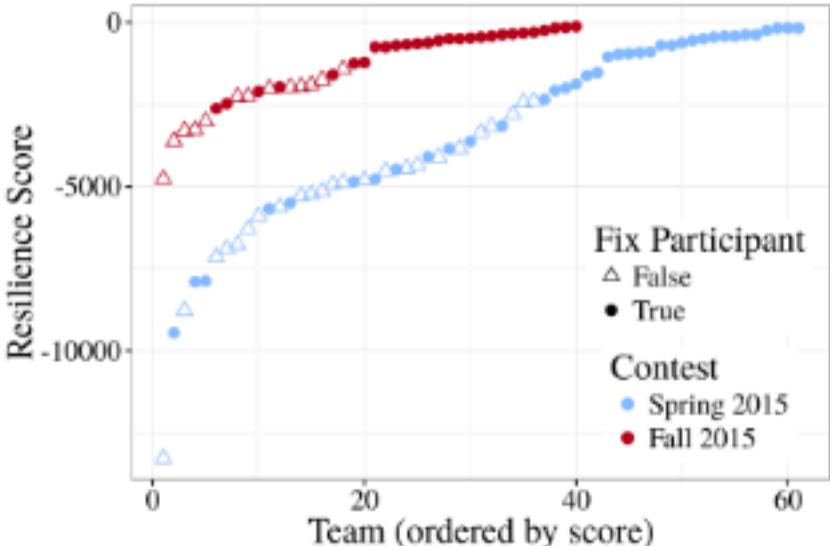


Figure 4: Final resilience scores, ordered by team, and plotted for each contest problem. Build-it teams who did not bother to fix bugs generally had lower scores.

Factor	Coef.	Exp(coef)	SE	p-value
Fall 2015	5.692	296.395	1.374	<0.001*
# Languages known	-0.184	0.832	0.086	0.033*
Lines of code	0.001	1.001	0.0003	0.030*
Dynamically typed	-0.751	0.472	0.879	0.393
Statically typed	-2.138	0.118	0.889	0.016*
MOOC	2.872	17.674	1.672	0.086

Table 4: Final logistic regression model, measuring log likelihood of a security bug being found in a team’s submission.

relative to the baseline category.¹² For numeric factors, the exponential indicates how the likelihood changes with each unit change in that factor.

Fall 2015 implementations were 296× as likely as Spring 2015 implementations to have a security bug discovered against them.¹³ We hypothesize this is due to the increased security design space in the ATM problem as compared to the gallery problem. Although it is easier to demonstrate a security error in the gallery problem, the ATM problem allows for a much more powerful adversary (the MITM) that can interact with the implementation; breakers often took advantage of this capability, as discussed in §5.

The model also shows that C/C++ implementations were more likely to contain an identified security bug than either static or dynamic implementations. For static languages, this effect is significant and indicates that a C/C++ program was about 8.5× (that is, 1/0.118) as likely to contain an identified bug. This effect is clear in Figure 5, which plots the fraction of implementations that contain a security bug, broken down by language type and contest problem. Of the 16 C/C++ submissions (see Figure 2), 12 of them had a security bug: 5/9 for Spring 2015 and 7/7 for Fall 2015. All 5 of the buggy implementations from Spring 2015 had a crash

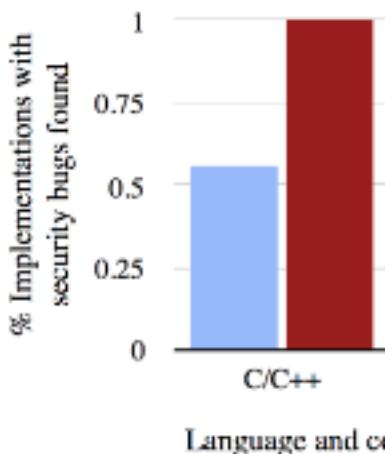


Figure 5: The fraction of teams for which a security bug was found, for each language and contest.

relevant and rerun the model. The lack of significance for the language category is no longer significant, which may indicate that lack of memory safety is not an advantage to using C/C++ from a breaker’s perspective, thus a memory-safe C/C++ would be equally effective.

Our model shows that teams using C/C++ and Java languages (even if they did not use them for the final submission) performed slightly worse than teams using a language known. Additional languages were also associated with a very small increase in the chance of an identified security bug.

Finally, the model shows that MOOC teams were less likely to find a security bug in the outcome, but not in the process. Teams using a dynamically typed language were more likely to find a bug in the MOOC. We see the effect of the MOOC training on the latter, the effect size is comparable to the effect of MOOC security training plus the effect of the language.

4.6 Breaking success

Now we turn our attention to the question of how effective teams were at finding security bugs in their submissions. First, we consider the fraction of teams that had a security bug found prior to the fix-it phase. We ignore the raw output, ignoring whether the bug was a crash or a denial of service, which we cannot assess from the data. We consider the fraction of fix-it phase participation for the 108 teams that participated in both Spring 2015 and Fall 2015. We also model the fraction of teams that found a security bug count, or the fraction of teams that found a security bug per break-it team found. Doing so allows us to measure the effort at finding correctness.

We model both break-it success and fix-it success using several of the same potential explanatory variables, but applied to the break-it and fix-it building teams. In particular, we include the

of the buggy implementations from Spring 2015 had a crash defect, and this was the only security-related problem for three of them; none of the Fall 2015 implementations had crash defects. If we reclassify crash defects as non-security

¹²In cases (such as the Fall 2015 contest) where the rate of security bug discovery is close to 100%, the change in log likelihood starts to approach infinity, somewhat distorting this coefficient upwards.

¹³This coefficient is somewhat exaggerated (see prior footnote), but the difference between contests is large and significant.

building team. In particular, the number of break-it Team members, Coding experience, average of C, and unique Languages members. We also add two

Build participant: Tru
participated in the build-it phase,

Advanced techniques:

Factor	Coef.	SE	p-value
Fall 2015	-2406.89	685.73	<0.001*
# Team members	430.01	193.22	0.028*
Knowledge of C	-1591.02	1006.13	0.117
Coding experience	99.24	51.29	0.056
Build participant	1534.13	995.87	0.127

Table 5: Final linear regression model of teams’ break-it scores, indicating how many points each selected factor adds to the total score. Overall effect size $f^2 = 0.039$.

ported using software analysis or fuzzing to aid in bug finding. Teams that only used manual inspection and testing are categorized as false. 26 break-it teams used advanced techniques.

For these two initial models, our potential factors provide eight degrees of freedom; again assuming power of 0.75, this yields a prospective effect size $f^2 = 0.136$, indicating we could again expect to find effects of roughly medium size by Cohen’s heuristic [2].

Break score. The model considering break-it score is given in Table 5. It shows that teams with more members performed better, with an average of 430 additional points per team member. Auditing code for errors is an easily parallelized task, so teams with more members could divide their effort and achieve better coverage. Recall that having more team members did not help build-it teams (see Tables 3 and 4); this makes sense as development requires more coordination, especially during the early stages.

The model also indicates that Spring 2015 teams performed significantly better than Fall 2015 teams. Figure 6 illustrates that correctness bugs, despite being worth fewer points than security bugs, dominate overall break-it scores for Spring 2015. In Fall 2015 the scores are more evenly distributed between correctness and security bugs. This outcome is not surprising to us, as it was somewhat by design. The Spring 2015 problem defines a rich command-line interface with many opportunities for subtle errors that break-it teams can target. It also allowed a break-it team to submit up to 10 correctness bugs per build-it submission. To nudge teams toward finding more security-relevant bugs, we reduced the submission limit from 10 to 5, and designed the Fall 2015 interface to be far simpler.

Interestingly, making use of advanced analysis techniques did not factor into the final model.

Being a build participant and having more coding experience is identified as a positive factor in break-it score, according to the model, but neither is statistically significant. Interestingly, knowledge of C is identified as a strongly negative factor in break-it score (though again, not statistically significant). Looking closely at the results, we find that *lack*

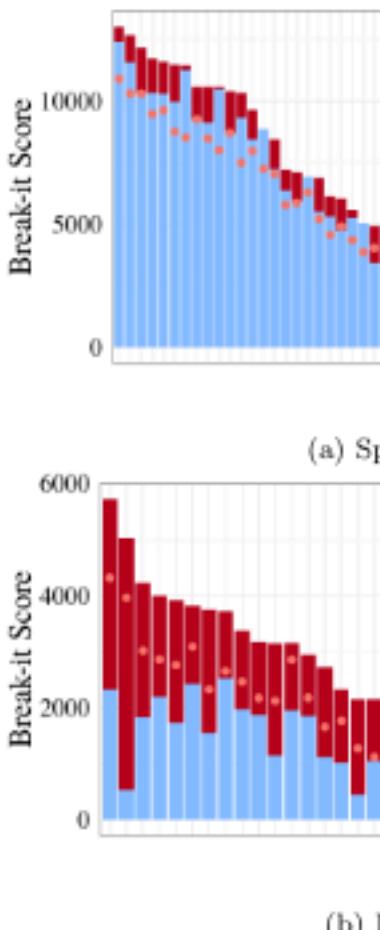


Figure 6: Scores of break-it teams broken down by points from correctness and security bugs. The final score of the break-it team is noted as a dot. Note the dominance of correctness bugs in general, the Spring 2015 contest showing higher scores for breaking.

Factor
Fall 2015
Team members
Build participant

Table 6: Final linear regression model of teams’ break-it scores, broken down by correctness and security bugs found by each team. The overall effect size $f^2 = 0.035$.

for each additional team member significantly helps one’s score. Thus, one would expect to gain a lot if one’s system could fail after building it. Figure 7 shows the distribution of the break-it scores per contest, for break-it teams and build-it teams.

On average, four more security bugs were found by a Fall 2015 team than a Spring 2015 team, though the Fall 2015

significantly. Looking closely at the results, we find that lack of C knowledge is extremely *uncommon*, but that the handful of teams in this category did unusually well. However, there are too few of them for the result to be significant.

Security bugs found. We next consider breaking success as measured by the count of security bugs a breaking team found. This model (Table 6) again shows that team size is important, with an average of one extra security bug found

2010 team than a Spring 2015 team, finding that Spring 2015 teams had higher scores, but corresponds to the fact that more submissions had security bugs. As discussed above, this is because many teams did not submit to Spring 2015 but were not able to do so in Spring 2010. Again, the reasons may have to do with the submission correctness bugs, which increase the potential attack surface in the system.

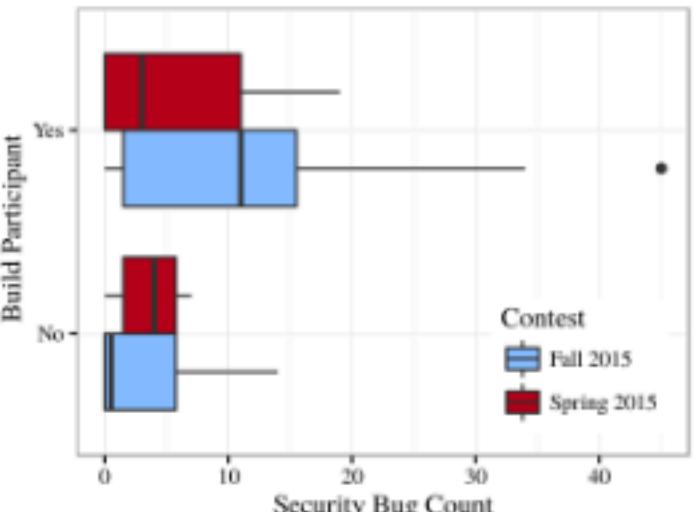


Figure 7: Count of security bugs found by each break-it team, organized by contest and whether the team also participated in build-it. The heavy vertical line in the box is the median, the boxes show the first and third quartiles, and the whiskers extend to the most outlying data within $\pm 1.5 \times$ the interquartile range. Dots indicate further outliers.

5. QUALITATIVE ANALYSIS

As part of the data gathered, we also have the entire program produced during the build-it phase as well as the programs patched during the fix-it phase. We can then perform a qualitative analysis of the programs which is guided by knowing the security outcome of a given program. Did lots of break-it teams find bugs in the program, or did they not? What are traits or characteristics of well-designed programs?

5.1 Success Stories

The success stories bear out some old chestnuts of wisdom in the security community: submissions that fared well through the break-it phase made heavy use of existing high-level cryptographic libraries with few “knobs” that allow for incorrect usage.

One implementation of the ATM problem, written in Python, made use of the SSL PKI infrastructure. The implementation used generated SSL private keys to establish a root of trust that authenticated the `atm` program to the `bank` program. Both the `atm` and `bank` required that the connection be signed with the certificate generated at runtime. Both the `bank` and the `atm` implemented their communication protocol as plain text then wrapped in HTTPS. This put the contestant on good footing; to find bugs in this system, other contestants would need to break the security of OpenSSL.

Another implementation, also for the ATM problem, written in Java, used the NaCl library. This library intentionally provides a very high level API to “box” and “unbox” secret values, freeing the user from dangerous choices. As above, to break this system, other contestants would need to first break the security of NaCl.

curity at all" to "vulnerable stacks." This is interesting observed in the software ma-

Many implementations of
tion or authentication. Ex-
trivial for break-it teams. S-
as plain text, other times like
Java object serialization pro-

One break-it team discovered could exploit with at most one mission truncated the “authenticated” vulnerable to a brute force attack.

The ATM problem allows
sible for the log), and the
implementations used cryptog
One implementation used c
RC4 from scratch and did :
key or the cipher stream.
ciphertext of messages was
changed from transaction t
flip bits in a message to cha

Another implementation cation, but did not use rand were always distinguishable was constructed against this attack leaked the bank balance, withdrawal attempts, distinguishing transactions, and performing

Some failures were communications: many implementations of the bank and a nonce in the messages. Transferring messages freely between the integrity via unauthorized writing.

Some failures were common as well: if an implementation used authentication. If it used automatic records stored in the file. The implementations would sort entries in the file to the order required for an integrity attack that used entries in the file.

6. RELATED WORK

BIBIFI bears similarity to security contests but is unique in systems. BIBIFI also is related to development, but differs in

Contests. Cybersecurity contests involve ability discovery and exploitation. System administration competitions of a similar style of contest is dubbed "CTF" (Capture The Flag). This is exemplified by a contest where two teams run an identical system and compete to find vulnerabilities.

break the security of NaCl.

An implementation of the log reader problem, also written in Java, achieved success using a high level API. They used the BouncyCastle library to construct a valid encrypt-then-MAC scheme over the entire log file.

5.2 Failure Stories

The failure modes for build-it submissions are distributed along a spectrum ranging from “failed to provide any se-

teams run an identical system. The goal is to find and exploit competitors’ systems while mitigating their own. A promising a system enables an adversary to capture the key and thus “capture the flag” in the challenge. The Cybersecurity Competition [1] have contestants compete in challenges. The responsibilities end at the point where the system is vulnerable. These contestants are not allowed to

a key factor of play, but neglect software development.

Programming contests challenge students to build clever, efficient software, usually with constraints and while under (extreme) time pressure. The ACM programming contest [1] asks teams to write several programs in C/C++ or Java during a 5-hour time period. Google Code Jam [2] sets tasks that must be solved in minutes, which are then graded according to development speed (and implicitly, correctness). Topcoder [3] runs several contests; the Algorithm competitions are small projects that take a few hours to a week, whereas Design and Development competitions are for larger projects that must meet a broader specification. Code is judged for correctness (by passing tests), performance, and sometimes subjectively in terms of code quality or practicality of design. All of these resemble the build-it phase of BIBIFI but typically consider smaller tasks; they do not consider the security of the produced code.

Studies of secure software development. There have been a few studies of different methods and techniques for ensuring security. Work by Finifter and Wagner [4] and Prechelt [5] relates to both our build-it and break-it phases: they asked different teams to develop the same web application using different frameworks, and then subjected each implementation to automated (black box) testing and manual review. They found that both forms of review were effective in different ways, and that framework support for mitigating certain vulnerabilities improved overall security. Other studies focused on the effectiveness of vulnerability discovery techniques, e.g., as might be used during our break-it phase. Edmundson et al. [6] considered manual code review; Scandariato et al. [7] compared different vulnerability detection tools; other studies looked at software properties that might co-occur with security problems [8]. BIBIFI differs from all of these in its open-ended, contest format: Participants can employ any technique they like, and with a large enough population and/or measurable impact, the effectiveness of a given technique will be evident in final outcomes.

7. CONCLUSIONS

This paper has presented Build-it, Break-it, Fix-it (BIBIFI), a new security contest that brings together features from typical security contests, which focus on vulnerability detection and mitigation but not secure development, and programming contests, which focus on development but not security. During the first phase of the contest, teams construct software they intend to be correct, efficient, and secure. During the second phase, break-it teams report security vulnerabilities and other defects in submitted software. In the final, fix-it, phase, builders fix reported bugs and thereby identify redundant defect reports. Final scores, following an incentives-conscious scoring system, reward the best builders and breakers. During 2015 we ran three contests involving a total of 116 teams and two different programming problems. Quantitative analysis from these contests found that

act as an incubator for ideas. More information, data, and available at <https://builditfixit.com>.

8. REFERENCES

- [1] ABADI, M., BUDIU, M., ET AL. Control-flow integrity principles for applications. *ACM Transactions on System Security (TISSEC)* 10, 1 (2007), 1–33.
- [2] The ACM-ICPC International Programming Contest. <http://icpc.baylor.edu>.
- [3] BURNHAM, K. P., ANDERSON, J. Model selection and multimodel inference: some background for ecologists. *Behavioral Ecology and Sociobiology* 47, 1 (2001), 177–190.
- [4] COHEN, J. *Statistical Power for the Behavioral Sciences*. Lawrence Erlbaum Associates, Inc., 1988.
- [5] DEF CON COMMUNICATIONS. DEF CON archive. <https://www.defcon.cc>.
- [6] DRAGOSTECH.COM INC. CANSECWORLD conference. <http://cansecworld.com>.
- [7] EDMUNDSON, A., HOLTKA, S., KARLSON, M., METTLER, A., AND WILSON, D. Evaluating the effectiveness of security contests. In *Symposium on Engineering Security (ESSoS)* (2013).
- [8] FINIFTER, M., AND WAGNER, R. Testing betweenweb application contests. In *USENIX Conference on Web Applications (WebApps)* (2011).
- [9] Git – distributed version control system. <http://git-scm.com>.
- [10] Google code jam. <http://codejam.withgoogle.com>.
- [11] HARRISON, K., AND WHITING, T. The effectiveness of common security challenges. In *International Conference on Cybersecurity* (2010).
- [12] ICFP programming contest. <http://icfpcontest.org>.
- [13] INC, D. C. C. Def con hacking competition. <http://www.defcon.org>.
- [14] KIM, Q. Want to learn cybersecurity? <http://www.marketplace.org/item/13777/want-learn-cybersecurity>.
- [15] Maryland cyber challenge. <http://fbcinc.com/e/cybermdcon>.
- [16] NATIONAL COLLEGIATE CHAMPIONSHIPS. <http://www.nationalccdc.com>.
- [17] PARKER, J. LMonad: Infrastructure for building web applications. Master’s thesis, University of Texas at Austin (2013).
- [18] POLYTECHNIC INSTITUTE OF NEW YORK - cybersecurity competition. <http://www.poly.edu/csa/contests>.
- [19] PRECHELT, L. Plat_forms: An experimental comparison by an exploration of the emergent platform properties. In *Software Engineering 37, Proceedings of the 37th International Conference on Software Engineering (ICSE 2015)*, 2015, pp. 101–110.
- [20] PostgreSQL: The world’s most advanced open source database. <http://www.postgresql.org>.
- [21] RUMMEL, A., HUGO, M., PAPALOIZOU, G., AND VILLEZ, J. The 2015 BIBIFI security contest. In *Proceedings of the 2015 ACM SIGSAC Conference on Computer and Communications Security (CCS 2015)*, 2015, pp. 101–110.

items. Quantitative analysis from these contests found that the best performing build-it submissions used C/C++, but submissions coded in a statically-typed language were less likely to have a security flaw; build-it teams with diverse programming-language knowledge also produced more secure code. Shorter programs correlated with better scores. Break-it teams that were also build-it teams were significantly better at finding security bugs. We plan to freely release BIBIFI to support future research. We believe it can

- [21] RUEF, A., HICKS, M., FAIR, J., PLANE, J., AND MARDZIEK, P. 2018. Evaluating static analysis tools for identifying security flaws in C/C++ code and comparing development environments. In *IEEE International Conference on Software Engineering (ICSE)*, 1–10.
- [22] SALTZER, J. H., AND SCHROEDER, M. D. 1975. On the reliability of information in computer systems. *Communications of the ACM* 18, 9 (1975), 1278–1308.
- [23] SCANDARIATO, R., WALDNER, J., AND WILSON, R. 2018. Security analysis versus penetration testing: An experimental study. In *IEEE International Conference on Dependable Systems and Networks (DSN)*, 1–12.

- [24] STEFAN, D., RUSSO, A., MITCHELL, J., AND MAZIERES, D. Flexible dynamic information flow control in haskell. In *ACM SIGPLAN Haskell Symposium* (2011).
- [25] Top coder competitions. <http://apps.topcoder.com/wiki/display/tc/Algorithm+Overview>.
- [26] ÚLFAR ERLINGSSON. personal communication stating that CFI was not deployed at Microsoft due to its overhead exceeding 10%, 2012.
- [27] WALDEN, J., STUCKMAN, P., AND HALL, R. Predicting vulnerable code mining. In *IEEE International Conference on Dependable Systems and Networks Reliability Engineering* (2012).
- [28] YANG, J., RYU, D., AND LEE, S. Improving prediction accuracy with hybrid measures. In *International Conference on Big Data and Smart Computing (BigC)* (2013).
- [29] Yesod web framework for Haskell. <http://www.yesodweb.com/>

