

Automatically Patching Errors in Deployed Software

Jeff H. Perkins^α, Sunghun Kim^β, Sam Larsen^γ, Saman Amarasinghe^α, Jonathan Bachrach^α, Michael Carbin^α, Carlos Pacheco^δ, Frank Sherwood, Stelios Sidiroglou^α, Greg Sullivan^ε, Weng-Fai Wong^ζ, Yoav Zibin^η, Michael D. Ernst^θ, and Martin Rinard^α
^αMIT CSAIL, ^βHKUST, ^γVMware, ^δBCG, ^εBAE AIT, ^ζNUS, ^ηCome2Play, ^θU. of Washington
jhp@csail.mit.edu, mernst@cs.washington.edu, rinard@csail.mit.edu

ABSTRACT

We present ClearView, a system for automatically patching errors in deployed software. ClearView works on stripped Windows x86 binaries without any need for source code, debugging information, or other external information, and without human intervention.

ClearView (1) observes normal executions to learn invariants that characterize the application's normal behavior, (2) uses error detectors to monitor the execution to detect failures, (3) identifies violations of learned invariants that occur during failed executions, (4) generates candidate repair patches that enforce selected invariants by changing the state or the flow of control to make the invariant true, and (5) observes the continued execution of patched applications to select the most successful patch.

ClearView is designed to correct errors in software with high availability requirements. Aspects of ClearView that make it particularly appropriate for this context include its ability to generate patches without human intervention, to apply and remove patches in running applications without requiring restarts or otherwise perturbing the execution, and to identify and discard ineffective or damaging patches by evaluating the continued behavior of patched applications.

In a Red Team exercise, ClearView survived attacks that exploit security vulnerabilities. A hostile external Red Team developed ten code-injection exploits and used these exploits to repeatedly attack an application protected by ClearView. ClearView detected and blocked all of the attacks. For seven of the ten exploits, ClearView automatically generated patches that corrected the error, enabling the application to survive the attacks and successfully process subsequent inputs. The Red Team also attempted to make ClearView apply an undesirable patch, but ClearView's patch evaluation mechanism enabled ClearView to identify and discard both ineffective patches and damaging patches.

Categories and Subject Descriptors:

D.2.5 [Testing and Debugging]: Error Handling and Recovery, Monitors

D.2.7 [Distribution, Maintenance, and Enhancement]: Corrections, Enhancement

K.6.5 [Security and Protection]: Invasive Software

General Terms: Security, Reliability, Design, Performance

1. INTRODUCTION

We present ClearView, a system for automatically correcting errors in deployed software systems with high availability requirements. Previous research has shown how to detect errors, for example by monitoring the execution for buffer overruns, illegal control transfers, or other potentially incorrect behavior [21, 37, 24]. The standard mitigation strategy is to terminate the application, essentially converting all errors into denial of service. In many important scenarios, system availability is a strict requirement. In such scenarios, it is imperative to eliminate the denial of service: the application should provide service even in the face of errors.

ClearView can automatically correct previously unknown errors in commercial off-the-shelf (COTS) software systems. It patches running applications without requiring restarts or otherwise perturbing the execution. It requires no human interaction or intervention. It works on stripped Windows x86 binaries without access to source code or debugging information.

Figure 1 presents the architecture of ClearView, which has five components:

- **Learning:** ClearView observes the application's behavior during normal executions to infer a model that characterizes those normal executions. The model is a collection of properties, also called (likely) invariants, over the observed values of registers and memory locations. Each invariant was always satisfied during the observed normal executions. As ClearView observes more executions, its model becomes more accurate. Our current ClearView implementation uses an enhanced version of Daikon [17] as its learning component.
- **Monitoring:** ClearView classifies each execution as normal or failed, by using a set of monitors that detect failures. For each failed execution, the monitor also indicates the location in the binary where it detected the failure, and the monitor prevents negative consequences by terminating the application. ClearView is designed to incorporate arbitrary monitors. Our current implementation uses two monitors: Heap Guard (which detects out-of-bounds memory writes) and Determina Memory Firewall (a commercial implementation of program shepherding [24] which detects illegal control flow transfers). Our current monitors have no false positives in that they never classify a normal execution as failed. But they are also only designed to detect a specific class of errors (heap buffer overflows and illegal control flow transfers). ClearView is not designed to eliminate all failures, only those that a monitor detects.
- **Correlated Invariant Identification:** When a monitor first detects a failure, ClearView installs patches that check previously-learned invariants close to the location of the failure. These invariant-checking patches are not intended to correct errors or eliminate the failure. The goal is instead to find a set of *correlated invariants* that characterize normal and failed executions.

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SOSP'09, October 11–14, 2009, Big Sky, Montana, USA.

Copyright 2009 ACM 978-1-60558-752-3/09/10 ...\$10.00.

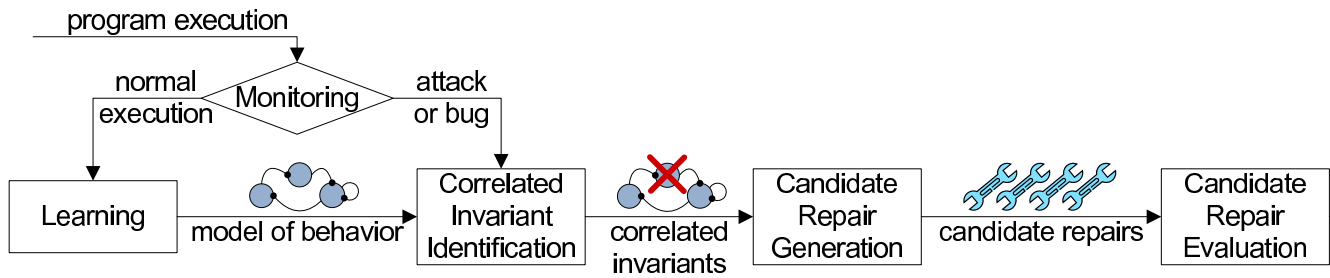


Figure 1: The ClearView Architecture

Specifically, each correlated invariant is always satisfied during normal executions but violated during failed executions (which typically occur in response to repeated or replayed attacks).

- **Candidate Repair Generation:** For each correlated invariant, ClearView generates a set of candidate repair patches that enforce the invariant. Some of these patches change the values of registers and memory locations to reestablish the invariant whenever it is violated. Others change the flow of control to enforce observed control flow invariants.

The hypothesis is that some errors violate invariants, that enforcing violated invariants can correct the effects of these errors, and that correcting these effects can change the execution of the application to eliminate the corresponding failure. The goal is to find a patch that corrects the execution after the first error has occurred but before any effects have propagated far enough to make a failure inevitable. To accomplish this goal, ClearView must find and correct errors that occur early enough in the execution to make the execution salvageable via the invariant enforcement mechanism.

- **Candidate Repair Evaluation:** A candidate repair patch might have no effect, or even a negative effect, on the patched application. ClearView therefore evaluates patches by continuously observing patched applications as they execute. It ranks each patch based on whether the application fails or crashes when the patch is in place. At each point in time ClearView attempts to minimize the likelihood of negative effects by applying the most highly ranked patch. The fact that patches affect the execution only when a correlated invariant is violated also tends to minimize the possibility that they will negatively affect normal executions.

The application’s maintainers (if they exist) may wish to find and eliminate the defect in the source code (if it still exists) responsible for the failure. ClearView supports this activity by providing information about the failure, specifically the location where it detected the failure, the correlated invariants, the strategy that each candidate repair patch used to enforce the invariant, and information about the effectiveness of each patch. This information may help the maintainers more quickly understand and eliminate the corresponding defect. In the meantime, the automatically generated ClearView patches can enable the application to survive to provide acceptable service, without waiting for human reaction.

ClearView is designed around the concept of learning from failure (and from success, too). This process of learning enables the quality of ClearView’s patches to improve over time, similarly to a biological immune system. The first time it detects a failure, ClearView learns a failure location in the application. This information enables ClearView to target subsequent instrumentation and intervention only where it is likely to be effective. The next several times it encounters the failure, ClearView learns the variables and data structures that are corrupted and the correlated invariants that

are violated. This information enables ClearView to generate candidate repair patches that may correct the error and eliminate the failure. Then, ClearView applies the patches. As the application processes subsequent inputs, ClearView evaluates the effectiveness of each patch. This evaluation enables ClearView to discard ineffective or damaging patches while applying successful patches that are able to eliminate the failure without negatively affecting the application.

Terminology: We adopt the following terminology from the software reliability community. A *defect* is a mistake in the design or source code of the application. An *error* occurs when the application does something incorrect, such as compute an incorrect value. In general, a *failure* is an observable error, such as a violation of the application’s specification that is visible to a user. In this paper we divide failures into three varieties. We use the term *failure* for errors that are detected by a ClearView monitor; a *crash* is any other error that causes the application to terminate; and an *anomaly* is any other observable error, i.e., an observable error that does not cause the application to terminate and is not detected by a ClearView monitor. A *failed execution* is an execution that a ClearView monitor terminated because it detected an error. A *normal execution* is an execution with no failures, crashes, or anomalies. An *exploit* is an input to an application that causes the application to violate its specification. An *attack* is a presentation of an exploit to an application. Note that defects and errors are undesirable primarily to the extent that they cause failures, crashes, or anomalies.

1.1 Red Team Evaluation

As part of DARPA’s Application Communities program (www.darpa.mil/IPTO/programs/ac/ac.asp), DARPA hired Sparta, Inc. (www.sparta.com) to perform an independent, adversarial Red Team evaluation of ClearView. The goal was to evaluate ClearView’s effectiveness in eliminating security vulnerabilities. Given the need for fast automated response to attacks that target such vulnerabilities, the time lag typically associated with human intervention (it takes 28 days on average for maintainers to develop and distribute fixes for such errors [47]), and ClearView’s ability to quickly and automatically generate and evaluate patches without human intervention, we anticipate that ClearView may be especially useful in this context.

In the Red Team exercise, ClearView protected an application community, which is a set of computers that all run the same software — in this case, the Firefox web browser. The community cooperates to learn information about attacks, then uses that information to provide immunity to all members (including members with no previous exposure to the attack) after only several have been attacked. The community also reduces the time required to obtain an effective patch.

During the Red Team exercise, the Red Team developed ten binary code injection exploits and used these exploits to repeatedly attack the community. The results show that:

- **Attacks Blocked:** ClearView detected and blocked all of the attacks; the injected code never executed.
- **Successful Continued Execution:** For seven of the ten exploits, ClearView automatically generated and applied a patch that corrected the underlying error and enabled the application to execute successfully even in the face of the corresponding attacks. For two of the remaining three exploits, changes to the ClearView configuration enabled ClearView to automatically generate similarly successful patches.
- **Patch Quality:** Some of the candidate repair patches that ClearView evaluated either did not eliminate the failure or introduced new negative effects such as causing the patched application to crash. ClearView’s patch evaluation mechanism recognized and discarded such patches, enabling ClearView to find and distribute patches that eliminated the failure without negative effects. The Red Team was unable to find a legitimate input that the final patched version of Firefox processed incorrectly.
- **No False Positives:** The Red Team was unable to elicit any false positives. Specifically, the Red Team was unable to cause ClearView to apply patches in the absence of an attack.

By making it possible to automatically survive otherwise fatal errors and attacks, our techniques may increase the robustness and availability of our computing infrastructure in a world increasingly full of errors and security vulnerabilities.

1.2 Contributions

This paper makes the following contributions:

- **Learning:** It shows how to observe normal executions to automatically learn invariants in stripped Windows x86 binaries with no access to source code or debugging information.
- **Invariants and Patches:** It identifies a set of invariants and corresponding patches that can automatically detect and correct errors and security vulnerabilities.
- **Invariant Selection and Patch Evaluation:** It shows how to automatically respond to a failure by selecting a promising set of correlated invariants to enforce, then evaluating the corresponding candidate repair patches to find a patch that enables the application to survive errors and attacks without service interruptions. Because the patch evaluation mechanism continually evaluates the effectiveness of the patches, it can recognize and discard ineffective or damaging patches even long after they were originally applied.
- **Application Communities:** It shows how communities of machines can work together as a group to survive errors and attacks. The members of the community share invariant learning information and benefit from the experience of other members to obtain full immunity without prior exposure.
- **Red Team Evaluation:** It presents experimental results from a hostile Red Team evaluation of ClearView. For seven of the ten security exploits that the Red Team developed, ClearView’s patches enabled the application to survive the attacks and continue on to successfully process subsequent inputs. There were no false positives — the Red Team was unable to make ClearView apply a patch in the absence of an error or attack.
- **Natural Resilience:** It provides additional evidence that, when augmented with simple survival techniques that enable continued execution through errors without termination, large software systems are naturally resilient to errors and attacks.¹ This empirically observed resilience suggests that one productive way

¹Related examples of such survival techniques include failure-oblivious computing [35], boundless memory blocks [34], transactional function termination [42, 39], data structure repair [14, 15] (including repair based on inferred invariants [13]), and generalizations of these techniques [33, 36].

to achieve robust software systems is to automatically apply survival techniques that modify the application and/or the execution environment to enable continued execution in the face of errors. This approach stands in stark contrast to standard attempts to develop error-free systems on top of existing brittle programming languages and execution environments. It also provides additional intriguing evidence that the complexity inherently present in many software systems may protect such systems from any negative effects of the localized perturbations characteristic of errors after the automatic application of survival techniques.

- **Invariant Inference:** It provides additional evidence that dynamic invariant inference, even when used with small test suites, can deliver accurate and effective specifications with a broad range of important uses.²

2. CLEARVIEW IMPLEMENTATION

This section describes ClearView’s implementation of each component of the architecture of Figure 1: learning (Section 2.2), monitoring (Section 2.3), correlated invariant identification (Section 2.4), repair generation (Section 2.5), and repair evaluation (Section 2.6). ClearView is currently built on the Determina commercial product suite (Section 2.1).

2.1 Determina Components

The Determina Managed Program Execution Environment uses DynamoRIO [5] for efficient, transparent, and comprehensive manipulation of arbitrary binary executables. All code executes out of a code cache. Before a code block enters the cache for execution, a plugin can validate and/or transform it. A plugin can also eject previously-inserted code blocks from the cache. Plugins can use this functionality to apply and remove patches in running applications without otherwise perturbing the execution. ClearView uses this capability to apply and remove patches that enable ClearView to check and enforce invariants.

Determina Memory Firewall detects all illegal control flow transfers and intervenes to terminate the execution before any injected code can execute. (The detection and intervention typically take place only after the attack has corrupted the application state so badly that the application can no longer execute successfully.) When an instruction in the code cache attempts to jump to code outside of the cache, Memory Firewall performs a validation check on the control flow transfer [24]. If the transfer passes the check, the Managed Program Execution Environment links the target code block into the cache (after applying any desired instrumentation or patches), then jumps to this target code block to continue executing from the cache. This implementation of program shepherding protects client applications from binary code injection attacks [24].

The Determina Management Console can monitor and control a large set of distributed client machines. It runs on a central server that stores patches and communicates securely with instantiations of the Determina Node Manager running on each client. Each Node Manager interacts with its corresponding Managed Program Execution Environment instance to appropriately apply and remove patches to and from running and newly-launched applications. Several ClearView Management Console plugins build on the Management Console functionality to coordinate ClearView’s interaction with client machines so that they appropriately apply and remove ClearView patches.

²Examples of techniques that use dynamic invariant inference to enhance security or otherwise improve application behavior include program steering [25], automatic inference and repair of data structure invariants [13], automatic vaccine generation [28], and automatic inference and checking of kernel data structure invariants [3]. See the Daikon web page (<http://groups.csail.mit.edu/pag/daikon/>) for many other uses of inferred invariants.

2.2 Learning

The learning component performs dynamic invariant detection [16], or specification mining. That is, it observes normal executions to infer a model of the normal behavior of the application. This model consists of a set of invariants that were satisfied in all observed executions and are statistically likely to be true on future executions. Each invariant is a logical formula that was always satisfied at a specific instruction in the application. Because ClearView operates on binaries, the variables in the formulas represent values (specifically, the values of registers and memory locations) that are meaningful at the level of the compiled binary.

2.2.1 Daikon

ClearView uses Daikon [17] as its learning component. Daikon is architected as two components: a front end that extracts trace data from a running application, and an inference engine that processes the trace data to infer the invariants. Daikon was originally developed to learn source-level invariants. We implemented a new Daikon front end for x86 that instruments the instructions in basic blocks, as they enter the code cache, to emit the appropriate trace data when they execute.

For each instruction, the trace data includes the values of all operands that the instruction reads and all addresses that the instruction computes. Consider, for example, `mov [ebp+12], eax`, which moves the value in memory location `ebp+12` into the register `eax`. This instruction computes one address (`ebp+12`) and reads one operand (the value in this address). Every time the instruction executes, the ClearView instrumentation produces a trace entry containing this data.

2.2.2 Invariant Variables and Locations

For each instruction, the front end must select the set of variables that it will supply to the Daikon inference engine. The set of variables must be large enough to enable the inference of meaningful invariants whose enforcement can correct our target class of errors. But, the set must also be small enough to make the inference task computationally tractable. Finally, the values in the variables must be defined in all possible executions (and not just the observed executions).

At each instruction, the front end outputs all variables that are computed by the instruction or by an instruction that predominates it.³ The front end computes dominators intraprocedurally. At each instruction, the Daikon inference engine infers only invariants that involve at least one variable that the target instruction computes [31].

2.2.3 Procedure Control Flow Graphs

To compute dominators, ClearView builds a control flow graph for each executed procedure. The nodes in the graph are basic blocks (as determined by the Determina Managed Program Execution Engine as it executes the application). The edges represent the flow of control between basic blocks.

The control flow graph construction algorithm uses a novel combined static and dynamic analysis. This analysis eliminates the need to find procedure entry points statically (a complex task in a stripped x86 executable). It maintains a database of known control flow graphs (with one control flow graph for each dynamically encountered procedure). It finds new procedures by considering each basic block the first time it executes. If the basic block is not already in a known control flow graph, the algorithm assumes that the basic block is the entry point for a new procedure. It then uses

³An instruction `i` predominates an instruction `j` if all control flow paths to `j` must first go through `i`. If `i` predominates `j`, then whenever the flow of control reaches `j`, `i` has previously executed and all of the values computed in `i` are valid.

symbolic execution to construct the control flow graph for this new procedure.

During the symbolic execution, the algorithm ends the procedure at return instructions, and also at indirect jump instructions for which it cannot compute the jump target. The algorithm may break a single static procedure up into multiple dynamically discovered procedures. The only potential drawback is that splitting procedures in this way may reduce the set of values available to the invariant inference engine at a given instruction, which may, in turn, reduce the set of invariants that the inference engine can infer for that instruction.

2.2.4 Additional Properties and Optimizations

ClearView analyzes the control flow graphs to identify duplicate variables (which denote values in registers or memory locations) in the same procedure that always have the same value. It then postprocesses the trace data to remove all duplicates except the one from the earliest instruction to execute. This optimization reduced the number of inferred invariants by a factor of two, which reduced the learning, invariant checking, and repair evaluation time.

Another ClearView postprocessing step analyzes the trace data to discover invariants involving stack pointer offsets. As a procedure allocates and deallocates local variables and calls other procedures, it may change the value of the stack pointer. Stack pointer offset invariants of the form $sp_1 = sp_2 + c$ capture the resulting relationships between the stack pointer `sp1` at the entry point of the procedure and stack pointers `sp2` at various points within the procedure. ClearView uses this information to adjust the stack pointer appropriately for repairs that skip procedure calls or return immediately from the enclosing procedure (see Section 2.5.1).

We also extended the set of Daikon invariants to include information about which variables contain pointer values. Specifically, if a negative value or a value between 1 and 100,000 ever appears in a variable, Daikon infers that it is not a pointer. Otherwise, Daikon infers that it is a pointer. Daikon uses this information to skip the inference of lower-bound or less-than invariants that involve pointers. This optimization reduced the learning, invariant checking, and repair evaluation time.

2.3 Monitoring

ClearView can incorporate any monitor that can detect a failure and provide a failure location (the program counter of the instruction where the monitor detected the failure). The current implementation uses Determina Memory Firewall to detect illegal control flow transfer errors. This monitor is always enabled in a running application. ClearView also adds two new monitoring components: Heap Guard and Shadow Stack.

Heap Guard: The Heap Guard monitor detects out-of-bounds memory accesses. It places canary values at the boundaries of allocated memory blocks and instruments all writes into the heap to check if the written location contains the canary value. The presence of the canary value indicates either an out-of-bounds write or a legitimate previous write (by the application) of the canary value within the bounds of an allocated memory block. When Heap Guard encounters a canary value, it therefore searches an allocation map to determine whether the written address is within the bounds of some allocated memory block. If so, normal execution continues; if not, Heap Guard has detected an out-of-bounds write error. By design, Heap Guard suffers no false positives. It may, however, miss an out-of-bounds write error if the out-of-bounds write skips over the canary value at the boundary of the allocated memory block.

Heap Guard is useful in two ways. First, it can detect out-of-bounds writes that do not cause illegal control flow transfers. In this

way, Heap Guard enables ClearView to detect (and therefore potentially correct) some out-of-bounds write errors that Memory Firewall does not detect. Second, even when an out-of-bounds write error would cause an illegal control flow transfer, Heap Guard may detect an earlier error than Memory Firewall would. This earlier detection may enhance ClearView's ability to find a successful patch.

It is possible to dynamically enable and disable Heap Guard as the application executes without otherwise perturbing the execution. ClearView could, for example, use this functionality to run the application without Heap Guard during normal production execution, then turn Heap Guard on when an event (such as a failure) indicates an elevated risk of an out-of-bounds write error.

Shadow Stack: ClearView can traverse the call stack to find additional candidate correlated invariants in callers of the procedure containing the failure location. Enforcing one of these additional invariants may be critical in enabling ClearView to correct the error.

ClearView maintains an auxiliary shadow procedure call stack and uses this Shadow Stack rather than attempting to unwind the native call stack. There are two reasons for this. First, a variety of optimizations (such as frame pointer removal and heavily-optimized caller/callee interactions) can make it difficult to reliably traverse a native call stack. Second, errors (such as buffer overflows) may corrupt the native stack, making it unavailable to ClearView when a monitor detects a failure.

The Shadow Stack contains the addresses of the procedures on the actual stack. Call and return instructions are instrumented to maintain the Shadow Stack. The instrumentation is largely performed in-line for efficiency and can be enabled and disabled as the application runs without perturbing the execution.

2.4 Correlated Invariant Identification

Given a failure location, ClearView attempts to identify invariants whose violation is correlated with the failure. These *correlated invariants* have two important properties. First, they are always satisfied in normal executions but violated in failed executions. Second, they are violated before the failure occurs. If a correlated invariant is violated at run time, then re-establishing it may force the application back into its normal operating envelope, correct the error, eliminate the failure, and enable the application to continue to operate successfully.

2.4.1 Candidate Correlated Invariants

ClearView uses the procedure call stack to obtain a set of candidate correlated invariants. Specifically, assume that a procedure P is on the call stack, with the program counter at instruction i , when a monitor detects a failure. Then any invariant at a predominator of i in P is in the candidate correlated invariant set.

In addition to the optimization of Section 2.2.2, ClearView further restricts the set of candidate correlated invariants. Any unary invariant, in any basic block of any procedure on the stack, may be a candidate correlated invariant. But for binary invariants, which relate the values of two variables, ClearView considers only invariants whose instruction occurs in i 's basic block. This additional restriction substantially reduces both the invariant-checking overhead (see Section 2.4.2) and the number of candidate repairs that ClearView generates and evaluates when a monitor detects a failure (see Section 2.5). In practice this optimization did not remove any useful repairs and (by reducing the number of repairs to evaluate) decreased the amount of time required to find a successful repair.

More generally, ClearView can use virtually any strategy that identifies a set of candidate correlated invariants that is likely to produce a successful repair. The three key considerations are 1) working with the available failure information (in the current im-

plementation this information includes the location of the failure and, if enabled, the Shadow Stack), 2) making the set large enough to include an invariant that produces a successful repair, and 3) limiting the size of the set to make it feasible to efficiently check the invariants (as described below in Section 2.4.2) and evaluate the resulting set of candidate repairs (as described below in Section 2.6). In particular, if the Shadow Stack is not available, ClearView can simply work with invariants associated with instructions close to the failure location. It is also possible to develop strategies that learn clusters of basic blocks that tend to execute together, then work with sets of invariants from clusters containing the basic block where the failure occurred.

2.4.2 Checking Candidate Correlated Invariants

Given a set of candidate correlated invariants, ClearView generates and deploys a set of patches that check whether each candidate correlated invariant is satisfied or violated. For an invariant over a single variable or over two variables at the same instruction, the patch executes when the program counter reaches the instruction associated with the variable(s). For an invariant that expresses a relationship between the values of two variables at different instructions, an auxiliary patch at the first instruction (to execute) stores the value of the first variable. Then, a patch at the second instruction reads the stored value and checks the invariant.

Each invariant-checking patch produces an observation every time it executes. Each observation identifies the invariant and the location of the failure that triggered the deployment of the patch. It also indicates whether the invariant was satisfied or violated. In this way the patched application produces, for each combination of invariant and failure location, a sequence of observations.

2.4.3 Identifying Correlated Invariants

When a monitor detects a failure, ClearView uses the sequences of invariant-checking observations to classify how highly an invariant is correlated with the failure:

- **Highly Correlated:** Each time a monitor detected the failure, the invariant was violated the last time it was checked. The invariant was satisfied all other times it was checked, in both normal and failed executions.
- **Moderately Correlated:** Each time a monitor detected the failure, the invariant was violated the last time it was checked. In at least one other failed execution, the invariant was also violated at least one other time it was checked. The invariant was satisfied in all normal executions.
- **Slightly Correlated:** At least one of the times a monitor detected the failure, the invariant was violated at least one of the times it was checked, but the correlation is not high or moderate.
- **Not Correlated:** The invariant was always satisfied in all failed executions.

A correlated invariant need not be violated every time it is checked. For example, the initial parts of a failed execution may exhibit normal behavior with no errors during which correlated invariants are satisfied.

If there are highly correlated invariants, the current ClearView implementation generates candidate repairs (see Section 2.5) for only those invariants. If there are no highly correlated invariants, ClearView generates candidate repairs only for moderately correlated invariants (if any exist). It would also be possible to generalize this approach to have ClearView generate candidate repairs for all correlated invariants, with the ClearView candidate repair evaluation (see Section 2.6) using the correlated invariant classification to prioritize the evaluation of the corresponding repairs.

2.5 Candidate Repair Generation

For each correlated invariant, ClearView generates one or more candidate repairs to evaluate. (There may be multiple ways to enforce a single invariant.) The patch that implements the repair (1) checks whether the invariant is violated, and (2) if so, re-establishes the invariant by changing the flow of control, the values of registers, and/or the values of memory locations. Patches for invariants involving only a single variable check and enforce the invariant at the variable's instruction. Patches for invariants involving multiple variables check and enforce the invariant at the latest (to execute) of the corresponding instructions for the involved variables.

We next describe the three invariants and corresponding repairs that ClearView used during the Red Team exercise (see Section 4).

2.5.1 One-of invariant

A one-of invariant has the form $v \in \{c_1, c_2, \dots, c_n\}$, where the c_i are constants and v is a variable or expression. This property identifies all of the values that v ever took on at run time. There are n repairs of the following form, one for each observed value:

```
if ! (v == c1 || v == c2 || ... || v == cn) then v = ci
```

If the application uses v as a function pointer (i.e., v is the target of a call instruction), another repair simply skips the call if the invariant is violated. The repair replaces `call *v` with

```
if (v == c1 || v == c2 || ... || v == cn) then call *v
```

A third repair returns immediately from the enclosing procedure:

```
if ! (v == c1 || v == c2 || ... || v == cn) then return
```

(The actual patch also adjusts the stack pointer to remove the arguments to the call and performs other cleanup.) This repair can be used for any invariant, but ClearView currently uses it only for one-of invariants.

Rationale: One-of invariants often characterize the observed targets of function calls that use function pointers. The use of uninitialized memory or incorrect type casts can produce incorrect function pointers. Many security attacks also exploit vulnerabilities that enable attackers to create malicious function pointers. Enforcing the invariant eliminates any illegal control flow transfer, which may, in turn, enable the application to survive the error or attack.

2.5.2 Lower-bound Invariant

A lower-bound invariant has the form $c \leq v$, where c is a constant and v is a variable or expression. One repair has the form:

```
if ! (c <= v) then v = c
```

Rationale: One class of defects can cause an array or buffer index to be negative, which can cause the application to read and/or write addresses below the start of the array or buffer. A related class of defects can cause the application to pass a negative number as a length to a procedure such as `memcpy`, which treats the number as a large unsigned integer. The resulting memory copy then writes beyond the end of the buffer.

Such defects typically result in errors that violate a lower-bound invariant such as $0 \leq v$. Enforcing the invariant redirects the out-of-bound index back into the buffer or array, which can prevent memory corruption and enable the application to survive the error.

2.5.3 Less-than Invariant

A less-than invariant $v_1 \leq v_2$ relates two variables or expressions (by contrast, a lower-bound invariant relates a variable and a constant). Less-than invariants can be repaired by adjusting either v_2 (as in the lower-bound repair) or v_1 . One repair is of the form:

```
if ! (v1 <= v2) then v1 = v2
```

Rationale: A defect can cause an array or buffer index to exceed the upper bound of the array or buffer, which can cause the applica-

tion to access addresses above the end of the array or buffer. Such defects typically cause the application to violate a less-than invariant that captures the requirement that the index must be below the upper bound of the array or buffer. Enforcing the invariant redirects the out-of-bound index back into the array or buffer, which can eliminate memory corruption and enable the application to survive the error.

2.6 Candidate Repair Evaluation

ClearView evaluates each candidate repair to determine which repair, if any, is the most effective at correcting the error and enabling the application to operate normally. ClearView considers the repair to have failed if the failure still occurs, a new failure occurs, or the application crashes after repair. ClearView (tentatively) considers the repair to have succeeded if the application has executed with the repair in place for at least ten seconds without failing or crashing. If it later fails or crashes, its status is changed.

The repair evaluation is based on relative repair scores. When a repair succeeds, its score increases. When a repair fails, its score decreases. Since the goal is to find a repair that always works, the scoring system is designed to reward repairs that are *always* successful. If a repair ever fails, the system continues to search for a more successful repair. ClearView therefore uses the scoring formula $(s - f) + b$, where s is the number of successes, f is the number of failures, and b is a positive bonus given to any repair that has not yet failed. ClearView uses the following criteria to break ties among repairs with the same score (for instance, all repairs that have never been tried have score b):

- **Lower on the Call Stack:** ClearView prefers repairs in procedures lower on the call stack.
- **Earlier Repairs First:** In a given basic block or procedure, ClearView prefers repairs from earlier instructions. The goal is to minimize error propagation by correcting the earliest error.
- **Minimize Control Flow Changes:** Repairs that change the control flow (such as returning immediately or skipping a call) are prioritized after repairs that only affect the state.

A repair may eliminate one failure, only to expose another failure. In this case, ClearView performs the full process of finding correlated invariants, generating candidate repair patches, and evaluating the generated patches all over again, starting with the patched application. ClearView may therefore generate multiple patches to repair multiple invariants. This actually happened in the Red Team exercise (see Section 4.4).

3. APPLICATION COMMUNITIES

It is possible to apply our technique successfully whenever repeated exposures to an error or attack give ClearView the opportunity to learn how to defend against the error or attack. One appropriate deployment environment is an *application community* — a group of machines running the same application that work together to detect and eliminate failures and/or to defend themselves against attacks.

Such a monoculture is convenient for users because it provides them with a familiar software environment across all machines, thereby making their data and expertise portable across the entire computing infrastructure. It can also decrease the system administration overhead. However, a monoculture may also be convenient for attackers, who may be able to exploit a single vulnerability throughout the entire community. By configuring ClearView to protect an application community, we view the software monoculture not as a weakness, but as an opportunity to enhance the effectiveness of the countermeasures that ClearView can deploy to neutralize attacks. An application community provides the following benefits:

- **Accurate Learning:** The learning component can work with many different users, datasets, and usage styles, thereby increasing the accuracy of the learned invariants and the quality of the applied repair patches.
- **Amortized Learning Overhead:** ClearView can distribute the learning overhead across the community, with each member incurring overhead for only a small part of the application. For details, see Section 3.1.
- **Faster Repair Evaluation:** The community can evaluate candidate repairs in parallel, reducing the time required to find a successful repair.
- **Protection Without Exposure:** After some members of the community are attacked and ClearView has found a successful patch, the patch is distributed throughout the community. The remaining members of the community become immune to the attack even though they have never been exposed to the attack. Furthermore, because the patch corrects the error, it can enable applications to immediately survive *other* attacks that attempt to exploit the same vulnerability, again with no previous exposure to these attacks.

The ClearView implementation contains components on both the community machines and a central server that coordinates the actions of the community. On each community machine, an instance of the Determina Node Manager coordinates the application and removal of patches to and from applications running on that machine. It also provides secure communication (via SSL) with the Determina Management Console running on the central server. The Management Console coordinates the distribution of patches to the Node Managers and mediates the communication between the ClearView components located on the community machines and the central server.

3.1 Amortized Parallel Learning

We extended Daikon to work in parallel across the members of a community as follows. On each community machine, a local version of Daikon processes the trace data to compute invariants that are true on that machine. ClearView periodically uploads the locally inferred invariants (not the large trace data that the local Daikon uses to infer the invariants) to the central server, which updates ClearView’s central database of invariants that are true across all executions on all members of the community.

It is important to discard any invariants from executions with errors. Our currently implemented system sends complete local invariant data after the application exits; for failed executions, the data are discarded and never sent. It would be possible to apply more sophisticated strategies, for example sending invariant summaries periodically but delaying the incorporation of newly-learned invariants for a period of time long enough to make any undesirable effects of the execution apparent. Only after the period has expired with no observed undesirable effects would the system use the invariants to update the central invariant database.

ClearView can also use sampling to distribute the learning overhead among the members of the community. We performed experiments in which we selected, for each community machine, some of the procedures in the application to trace. ClearView instrumented only those procedures to generate trace data. The rest of the application executed without learning (and without any learning overhead). The fact that it is possible to learn over arbitrary parts of the application enables a wide range of distributed learning strategies that trade off the learning overhead at each community member, the comprehensiveness of the learning coverage, and the time required to obtain an acceptable set of invariants. One possible learning strategy would instrument a randomly chosen small

part of every running application, with new invariants continuously trickling in from all members of the community. This strategy minimizes the learning overhead while keeping the invariants up to date with information from the latest usage patterns.

It would also be possible to stage the learning. The first phase would record coverage for each input, typically at procedure granularity. After a failure, the second phase would instrument regions close to the failure location, then replay inputs that exercise these instrumented regions. Daikon would then process the generated trace data to produce a set of candidate correlated invariants. The advantage of this approach is that it would reduce the learning overhead and eliminate the need for a large invariant database. The drawback is that learning only in response to failures would significantly delay obtaining a successful patch. The current ClearView implementation does not implement this scheme.

3.2 Application Community Management

We next describe the actions ClearView takes as it manages the community in response to a failure.

Detection and Failure Notification: A ClearView monitor running on one of the community machines encounters the failure. The monitor terminates the application, then uses the underlying Determina secure communication facilities to notify the central ClearView manager of the failure. The notification includes the failure location and (if available) the call stack at the time of the failure.

Identifying Correlated Invariants: The central ClearView manager responds to the failure notification by using the failure location and call stack to access its central invariant database and compute a set of candidate correlated invariants. For each such invariant, it generates a snippet of C code that checks the invariant, then compiles the C code to obtain a patch that checks the invariant. It presents the patches to the Determina infrastructure, which pushes the patches out to all of the members of the community. The local Determina Node Managers apply the patches to executing and newly-launched instances of the application.

As the patches execute, they generate a stream of invariant check observations that are sent back to the central ClearView manager. As described in Section 2.4.3, each observation identifies the invariant, the failure that caused ClearView to generate the invariant-checking patch, and an indication of whether the invariant was satisfied or violated.

Eventually, one or more instances of the application may encounter the failure again. When the central ClearView manager receives the failure notifications, it analyzes the invariant check observations (as described in Section 2.4.3) to compute a set of correlated invariants. It then removes the invariant-checking patches — the Determina Console Manager instructs the Determina Node Managers to remove the patches from any instances of the application on their machines. ClearView currently performs this step when it receives the second failure notification from a version of the application with invariant-checking patches in place. It is straightforward to implement other policies.

The current ClearView configuration always runs applications with the Shadow Stack and Heap Guard monitor turned on. It is straightforward to implement other policies. For example, one could turn these features on only after encountering the first failure, then turn them back off again after the community runs the patched application for a certain period of time without observing a failure.

Generating and Evaluating Repairs: The central ClearView manager next generates candidate repair patches for all of the correlated invariants, specifically by generating and compiling a snippet of C code that implements the invariant check and enforcement (see

Section 2.5). It then evaluates the repairs (see Section 2.6), using the Determina infrastructure to apply and remove the patches in running or newly-launched instances of the application. At any moment, the most successful patch is applied across the entire community, including instances of the application that have never encountered the failure. Ideally, this repair algorithm eventually finds a successful patch that corrects the error.

Multiple Concurrent Failures: It is possible for the community to encounter different failures at the same time (the Red Team exercise explored such a scenario, see Section 4.3.5). Because ClearView applies each patch in response to a specific failure (as identified by the failure location) and all ClearView communications identify the failure ultimately responsible for the communication, ClearView can manage the community as it responds to the events generated in response to multiple different concurrent failures.

4. RED TEAM EXERCISE

As part of DARPA's Application Communities program (www.darpa.mil/IPTO/programs/ac/ac.asp), DARPA hired Sparta, Inc. (www.sparta.com) to perform an independent, adversarial Red Team evaluation of ClearView. The Red Team consisted of eleven Sparta engineers. The goal of the Red Team was to discover and exploit flaws in our approach and in the ClearView implementation. The other Red Team exercise participants consisted of the Blue Team (the authors of this paper) and the White Team (a group of engineers from Mitre, Inc. led by Chris Doh of Mitre). The White Team determined the rules of engagement and refereed the exercise. The exercise was held at MIT on February 25–28, 2008.

During this exercise the Red Team used ten distinct exploits to attack the application protected by ClearView (Firefox 1.0.0). The Red Team verified that each exploit successfully targeted a security vulnerability in the unprotected version of Firefox, resulting in execution of arbitrary attacker-chosen code. ClearView detected and blocked all attacks, terminating Firefox before the attacks took effect. Moreover, ClearView generated successful patches for seven of the ten attacks.

4.1 Evaluation Goals

The primary purpose of the Red Team exercise was to evaluate the effectiveness of the ClearView technology in protecting against binary code injection attacks, i.e., attacks that attempt to subvert the control flow of the application, typically by causing the application to jump to downloaded code, but more generally by causing the application to take any unauthorized control flow transfer. This particular class of attacks was chosen because they are common in practice and can have particularly serious consequences if successful. The Red Team evaluation had several specific goals:

- **Surviving Attacks:** Can ClearView respond to attacks by finding patches that enable the application to survive the attack and continue to execute successfully?
- **Repair Evaluation:** Does ClearView ever generate a patch that impairs the application? For example, a bad patch might cause the application to behave incorrectly on legitimate inputs or create a new exploitable error.
- **False Positives:** Do legitimate inputs ever trigger the ClearView patch generation mechanism?
- **Infrastructure Attacks:** Can attackers subvert the ClearView patch generation and distribution mechanism to send out malicious patches? This paper omits the detailed results of this qualitative evaluation. In summary, the Red Team judged that the security measures in the Determina commercial product (encryption, authentication, etc.) provide an acceptable level of protection against this class of attacks.

4.2 Rules of Engagement

The rules of engagement determined the scope of the Red Team exercise — what kinds of Red Team attacks were in bounds, how to judge if an attack succeeded or failed, the access that the Red Team was given to Blue Team information, etc. Together, the Red, Blue, and White Teams agreed on an application (the unmodified, stripped x86 binary of Firefox 1.0.0) for the Blue Team to protect. With this application, the attack vector was web pages — the Red Team launched all attacks by navigating Firefox to one or more attack HTML, XUL, or GIF files. Firefox 1.0.0 has several properties that made it appropriate for this exercise:

- **Mature Code Base:** The Firefox code base was relatively mature and tested, which made it a reasonable proxy for other mature applications that ClearView is designed to protect.
- **Vulnerabilities:** This version of Firefox contained enough vulnerabilities to support a thorough evaluation without the need for the Red Team to find an infeasibly large number of new vulnerabilities.
- **Source Code Availability:** Source code was available for this application. Although ClearView does not require (or even use) any source information, the availability of source code made it much easier to understand the behavior of the application and interpret the phenomena observed during the Red Team exercise.
- **Automation:** Firefox supported the automated loading of web pages, which facilitated automated learning and testing both in preparation for and during the Red Team exercise.

Given the envisioned scope of the Red Team exercise and the available resources, it was not feasible to add more applications to the Red Team exercise.

4.2.1 Attack Scope

The Red Team attacks fall into three categories: control flow attacks, induced autoimmune attacks, and false positive attacks.

A control flow attack attempts to subvert the flow of control within the application. Such an attack was judged to succeed if it prevented the application from continuing to successfully process additional inputs, either by successfully redirecting the flow of control to malicious code or by causing the application to crash.

A false positive attack present a non-malicious input to the application. The attack succeeds when it causes ClearView to apply a patch in response to loading a legitimate web page.

An induced autoimmune attack attempts to turn the ClearView patch mechanism against the application. Such an attack succeeds if ClearView's patch affects the behavior of the application on legitimate inputs (as opposed to attack inputs). An autoimmune attack was judged to succeed if the patched version of Firefox, when made to navigate to a sequence of legitimate web pages, did not behave the same as the unpatched version (bit-identical displays, same user functionality).

Insider attacks, where some nodes start out malicious and can send misleading data about application behavior, were not part of the threat model for this exercise.

Within these constraints, the Red Team was given completely free rein in generating attacks. Known attacks, variants on known attacks, completely new attacks, and attacks that involved multiple web pages loaded in sequence were all within scope. There were no restrictions whatsoever placed on the information that the Red Team was allowed to use when generating attacks.

4.2.2 Red Team Exercise Preparation

Prior to the Red Team exercise, the Blue Team generated an invariant database by running Firefox on a collection of twelve web

pages that exercise functionality related to known Firefox vulnerabilities. Learning was confined to selected regions of the application related to these vulnerabilities. The web pages and invariant database were both made available to the Red Team before the Red Team exercise.

Several months before the Red Team exercise, the Blue Team provided the Red Team with all of the Blue Team’s source code, tests, and documentation, including design documents, presentations to sponsors, and the Blue Team’s own analyses of weaknesses in ClearView. During the period of time leading up to the Red Team exercise, the Blue Team periodically provided the Red Team with source code, test, and documentation updates. At the time of the Red Team exercise, the Red Team had complete access to all of the source code, tests, and documentation for the running Blue Team system.

Prior to the Red Team exercise, the Red Team selected 57 evaluation web pages. These legitimate web pages exercise a range of Firefox functionality and were used during the Red Team exercise for repair evaluation (specifically, to determine if the patched version of Firefox displayed the evaluation pages correctly) and false positive evaluation (specifically, to determine if any of the evaluation pages triggered the ClearView patch generation mechanism). The Blue Team was not provided with these web pages prior to the Red Team exercise.

For the Red Team exercise, the Blue Team provisioned a small community of machines with Firefox deployed on all machines. The Blue Team configured ClearView with Memory Firewall, Heap Guard, and the Shadow Stack enabled from the start on all Firefox executions.⁴ The Red Team attacked this community during the Red Team exercise.

4.3 Attack Evaluation

The first phase of the Red Team exercise evaluated ClearView’s ability to protect Firefox against Red Team attacks. The Red Team selected ten defects in Firefox, then created one or more exploits for each defect. The targeted defects cause exploitable errors such as unchecked JavaScript types, out-of-bounds array accesses, heap and stack buffer overflows, and JavaScript garbage collection problems. All of the exploits were verified to work — each successfully exploited a vulnerability in Firefox. The Red Team used the exploits to perform the following attacks.

4.3.1 Single Variant Attacks

For each defect, the Red Team chose an exploit, then repeatedly presented the exploit to an instance of Firefox running in the community. The Red Team presented each attack only after ClearView had performed all actions taken in response to the previous attack. For all 10 exploits, ClearView monitors detected and blocked the corresponding attacks. For 7 of the exploits, ClearView generated patches that enabled Firefox to continue to execute through the attacks to correctly display the subsequently loaded evaluation pages. The Red Team observed no differences between the patched and unpatched versions of Firefox. Subsequent investigation after the Red Team exercise (see Section 4.3.2) indicated that small configuration changes enabled ClearView to successfully generate patches for two of the remaining three exploits.

Table 1 presents the number of exploit presentations required for ClearView to find and apply the patch that enabled Firefox to execute successfully through the attack. In general, the minimum number of exploit presentations is four. The first presentation makes

⁴With one exception. Specifically, the presence of Heap Guard disabled the exploit for the defect with Bugzilla number 296134. To enable the meaningful inclusion of this exploit in the Red Team exercise, the Blue Team turned off Heap Guard when the Red Team deployed this exploit.

Bugzilla Number	Presentations	Error Type
269095	6	Memory Management
285595*	4	Heap Buffer Overflow
290162	4	Unchecked JavaScript Type
295854	5	Unchecked JavaScript Type
296134	4	Stack Overflow
311710	12	Out-of-Bounds Array Access
312278	4	Memory Management
320182	6	Memory Management
325403*	4	Heap Buffer Overflow

Table 1: Number of times each exploit was presented before ClearView created and applied a patch that protected against the exploit. A * identifies the two exploits for which ClearView did not successfully generate a patch during the Red Team exercise, but did successfully generate a patch in subsequent experiments after reconfiguration.

ClearView aware of the exploit; ClearView responds by computing a set of candidate correlated invariants and applying patches that check candidate invariants (see Sections 2.4.1 and 2.4.2). During the next two presentations ClearView records invariant satisfaction and violation information to compute the set of correlated invariants (see Section 2.4.3). After these presentations ClearView removes the invariant-checking patches and generates and applies patches that enforce correlated invariants (see Sections 2.5 and 2.6). If the first invariant enforcement patch is successful, ClearView has corrected the error in four presentations.

As Table 1 indicates, the first invariant enforcement patch successfully corrected the errors from exploits 290162, 296134, 312278, 285595, and 325403. For exploit 295854 the first patch did not correct the error, but the second patch did. For exploits 269095 and 320182 the third patch was the first successful patch.

Exploit 311710 is an outlier in that it involves three separate defects, each of which is exploited by the same attack. ClearView corrects the error from the first defect after four presentations, at which point the attack exploits the second defect. It takes ClearView another four presentations to correct the error from this second defect, at which point the attack exploits the third defect. It takes ClearView another four presentations to correct the error from this final defect, for a total of twelve presentations to obtain a set of patches that enables Firefox to successfully survive the attack.

We next discuss the different exploits (we group the discussion by type) and the ClearView response to each exploit.

Stack Overflow: Exploit 296134 causes Firefox to incorrectly compute a negative value for the length of a string. This negative length is then passed to `memcpy`, which treats it as a very large unsigned integer. The resulting copy writes downloaded data over exception handlers on the stack. When the copy proceeds past the end of the stack, the invoked overwritten exception handler executes downloaded code.

During learning ClearView learned a lower-bound invariant that requires the computed string length to be at least one. ClearView generated a patch that enforced this invariant by setting the length to one. This patch corrected the out-of-bounds writes to the stack and enabled Firefox to survive the attack.

Unchecked JavaScript Type Exploits: Both exploits 290162 and 295854 download JavaScript code that creates an object, then fills the object with malicious code and data. A JavaScript system routine fails to check the type of the object and (eventually via a sequence of operations) invokes downloaded code via a C++ virtual function call on the corrupted object.

During learning ClearView learned a one-of invariant at the virtual function call site for both errors. These invariants state that the call site may invoke only a function that was invoked at that

site during learning. The first patch that ClearView applied during repair evaluation enforced the invariant by invoking a specific previously invoked function instead of jumping to malicious code. This patch successfully corrected the error from exploit 290162, but failed to correct the error from exploit 295854. ClearView's second applied patch, which enforced the invariant by skipping the call, successfully corrected the error from exploit 295854.

Memory Management Exploits: Exploit 312278 enables downloaded JavaScript code to obtain a pointer to an object that is incorrectly garbage collected, then reallocated to hold a native Firefox C++ object. The downloaded JavaScript code then overwrites the C++ object's virtual function table pointer with a pointer to memory containing pointers to malicious downloaded code. During learning ClearView learned a one-of invariant at the virtual function call site that invokes the malicious code. This invariant states that the call site may invoke only one of the functions invoked at that site during learning. The first patch that ClearView applied during repair evaluation successfully corrected the error by invoking a specific previously invoked function.

Exploits 269095 and 320182 involve memory that is not reinitialized after it is reallocated. Under certain circumstances it is possible to manipulate Firefox into treating this uninitialized memory as a C++ object, then invoking a virtual function call on this uninitialized object. Downloaded JavaScript code can exploit this error to fill the memory with appropriately formatted malicious code and pointers before it is reallocated. In this case the virtual function call invokes the malicious code. During learning ClearView learned a one-of invariant at the virtual function call site that invokes the malicious code. One of the patches that ClearView applied during repair evaluation enforced the invariant by returning from the function that contains the call site before the call site is invoked. This patch enabled Firefox to survive the attack.

Before trying this patch, ClearView tried patches that invoke one of the previously observed functions and a patch that skips the call but executes the remaining part of the function following the call. None of these patches enabled Firefox to survive the attack.

Out-of-Bounds Array Access Exploits: Exploit 311710 causes Firefox to compute a negative array index, then use the index to attempt to retrieve a C++ object from the array. Downloaded JavaScript code previously caused the retrieved memory to contain pointers to downloaded code. When Firefox performs a virtual function call on the retrieved object, it invokes the downloaded code.

During learning ClearView learned a lower-bound invariant that requires the array index to be non-negative. During repair evaluation ClearView applied a patch that enforces this invariant by setting the array index to zero. This patch caused Firefox to retrieve a valid C++ object, the resulting virtual function call invoked valid code, and Firefox survived the attack.

The same defect that caused this error was present in three similar procedures (apparently created via copy and paste) that executed during the attack. A similar patch corrected each error from these defects.

4.3.2 Remaining Exploits

During the Red Team exercise, ClearView did not generate a successful patch for three of the Red Team's exploits.

Exploit 285595 targets code for a Netscape GIF extension. Because this code does not check the sign of a value extracted from the GIF file, it is vulnerable to a remotely exploitable heap overflow attack. During the Red Team exercise, ClearView's correlated invariant identification component was configured to consider invariants from only the lowest procedure on the stack with invariants. The relevant invariant appeared one procedure above this proce-

cedure. ClearView therefore did not produce a patch that corrected the error. We subsequently verified that changing the configuration to include additional procedures on the stack enabled ClearView to generate a successful patch that corrects this error. The relevant invariant is a lower-bound invariant involving a buffer index. The repair changes the index from a negative value to zero, thereby bringing out-of-bounds writes back into the buffer. The exploit itself is embedded within an image file. The repair neutralizes the attack and enables Firefox to display the image correctly.

Exploit 325403's attack vector is a HTML file, one of whose tag values is eventually used as a buffer growth size for data that does not fit in an allocated buffer that holds two-byte Unicode characters. By specifying a value very close to the largest representable unsigned integer, an attacker can cause the calculation of the new buffer size to overflow, causing Firefox to allocate a buffer that is too small. An ensuing `memcpy` then writes beyond the end of the allocated buffer. The Blue Team's learning suite for the Red Team exercise did not provide sufficient coverage for Daikon to learn the relevant invariant. We subsequently verified that, using an expanded learning suite, Daikon would have learned an invariant that would have enabled ClearView to generate a successful patch. The relevant invariant is a less-than invariant relating the buffer size to the size of the memory to copy into the buffer. The repair sets the copy size to the buffer size, eliminating the out-of-bounds writes.

Exploit 307259 causes Firefox to compute an incorrect size for a buffer holding a hostname that contains soft hyphens. When Firefox attempts to copy a number of items into this buffer, the copies write beyond the end of the buffer. ClearView did not generate a successful patch because Daikon's invariants are not rich enough to capture the error. The appropriate invariant would generalize Daikon's less-than invariant (which relates two quantities) to relate a sum of buffer lengths to another buffer length. Learning richer invariants would be possible, but would increase the cost of the learning component.

4.3.3 Comparison With Manual Fixes

Manual fixes are available for the defects that the exploits in the Red Team exercise exploited. For exploit 269095, the manual fix tags deallocated objects as invalid. Subsequent object uses check this tag. If the tag is invalid, the use returns an error. The fix also iterates over invalid objects to reinitialize relevant data. For exploit 285595, the manual fix removes the code containing the defect. This code implemented a Netscape GIF extension; the fix removes support for this extension from Firefox. For exploits 290162 and 295854, the manual fix checks the type of the JavaScript object. If the check fails, the enclosing method (which otherwise invokes a method on the object) simply returns null. For exploit 296134, the manual fix adds a check for negative string length. If the check fails, the enclosing method logs an error, returns, and does not perform the copy. The fix also includes a check in the calling method that truncates the string length to the allocated buffer size. The manual fix for 311710 corrects a conditional that caused the application to compute the negative array index.

For exploit 312278, the manual fix informs the garbage collector that it holds a reference to the relevant object. Once the garbage collector is aware of this reference, it does not collect the object and the memory holding the object is unavailable to the JavaScript code in the exploit. For exploit 320182, the manual fix sets a flag that identifies reallocated objects. Subsequent code checks the flag to identify and properly initialize any such reallocated objects. For exploit 325403, the manual fix checks that the target array is large enough to hold the data in the source array. If the check fails, the fix allocates a larger target array and retries the copy.

Some of these manual fixes perform a consistency check close to the error, then skip the remaining part of the operation if the check fails. All ClearView patches similarly perform a consistency check (the invariant satisfaction check) close to the error. Two of the ClearView repairs (skip call and return from enclosing procedure) have a similar effect of explicitly skipping part or all of the remaining part of the operation. Other repairs (adjusting values to enforce lower-bound and less-than relationships) often have the effect of enabling the application to execute the remaining part of the operation safely with the effect of any remaining errors localized to that operation. In general, the ClearView repairs tend to execute more of the normal-case code following the error, while the manual fixes tend to simply abort the current operation. We attribute this more drastic approach to the maintainer attempting to simplify the reasoning required to confirm that the fix has eliminated the error.

It would be possible to enhance ClearView to produce checks and repairs that more closely correspond to these manual fixes. For example, it would be possible to enhance ClearView to automatically infer object types at dynamically dispatched method invocations and to return error codes from enclosing procedures when invariant checks fail. It would also be possible to develop repairs that skip larger parts of the subsequent computation.

Some of the manual fixes inform the garbage collector of existing references and reinitialize recycled memory. These fixes affect code far from the failure location. ClearView would therefore need to apply a more sophisticated correlated invariant identification strategy (and potentially more sophisticated invariants as well) to produce repairs with similar effects.

4.3.4 Multiple Variant Attacks

For three defects, the Red Team generated multiple variants of the exploit that targeted the defect. For each defect the Red Team interleaved different variants during the attack. ClearView generated the same patch after the same number of attacks as for the corresponding single variant attack. This patch successfully protected Firefox against all variants of the attack.

4.3.5 Simultaneous Multiple Exploit Attacks

The Red Team launched several attacks that interleaved exploits that targeted different defects. The goal was to determine if targeting different defects with different exploits would impair ClearView's ability to generate successful patches. In each case ClearView was able to determine the targeted error for each attack, keep the learning data separate for the different errors, and generate a set of patches that together successfully protected Firefox against all of the exploits in the attack. And ClearView was able to generate these patches after the same cumulative number of attacks as for the corresponding sequence of single variant attacks.

4.3.6 Repair Evaluation

The Red Team was unable to launch a successful induced autoimmune attack. For all of the previous attack scenarios the Red Team evaluated the quality of the repair by determining whether the patched version of Firefox displayed the evaluation web pages correctly: whether the displays were identical and whether there was any other behavioral change that the Red Team could detect. The Red Team was unable to detect any difference between the original and patched Firefox on non-attack webpages. The same was true for a Firefox that was patched with when all of the successful patches generated during the previous attacks.

The ClearView repair evaluation mechanism is designed to discard patches that have a negative effect on the application. ClearView generated patches with negative effects (such as causing the application to crash) during the Red Team exercise. But the re-

ClearView Configuration	Page Load Time (seconds)	Overhead Ratio
Bare Firefox	7.5	1.0
Memory Firewall	11.0	1.5
Memory Firewall + Shadow Stack	14.9	2.0
Memory Firewall + Heap Guard	19.0	2.5
Memory Firewall + Heap Guard + Shadow Stack	22.7	3.0

Table 2: Page load times and overheads for different ClearView configurations running Firefox.

pair evaluation mechanism detected and discarded these patches, mitigating the effect on the application and paving the way for the application of successful patches.

4.3.7 False Positive Evaluation

The Red Team's false positive evaluation used ClearView to display the evaluation web pages. The goal was to make ClearView generate an unnecessary patch. During this evaluation ClearView generated no patches at all, indicating that the Red Team was unable to cause ClearView to produce a false positive.

4.4 Performance

The Red Team exercise used a Dell 2950 rack-mount machine with 16 GB of RAM and two 2.3 GHz Intel Xeon processors, each with four processor cores. We ran Firefox inside VMware virtual machines under ESX servers. The operating system was Windows XP Service Pack 2. Because the Red Team's exploit 296134 has no effect in this environment, we ran this exploit on a 1.8 GHz AMD Opteron machine with four processor cores and 8 Gbytes of RAM running Windows XP Service Pack 2. In this environment the exploit does trigger the error. Because the exploit is running on a slower computing platform, the execution times for the various activities are longer than the corresponding times for other exploits.

4.4.1 Learning Overhead

The time required to load the twelve learning web pages without learning enabled was 5.2 seconds. The time required to load these same web pages with learning enabled was 1600 seconds (over a factor of 300 slower). The Daikon x86 front end, which records and dumps the values of accessed memory locations and registers, is responsible for the vast majority of the overhead. As described in Section 3.1, it is possible to distribute the learning in parallel across the application community.

4.4.2 Baseline Overheads

Table 2 presents the times required for Firefox to load the 57 evaluation pages from a local disk with the network interface disabled when running under various ClearView configurations. We ran the experiments on an Intel Core 2 Duo E6700 (2.66 GHz, 4 MB L2 Cache, 3.25 GB RAM) running Windows XP Service Pack 3. The Determina Managed Program Execution Environment (with Memory Firewall enabled) imposes a 47% overhead over running Firefox as a standalone application. This overhead is somewhat larger than typically observed [5]. We attribute this additional overhead to Firefox's use of object-oriented constructs such as dynamic method dispatch. The underlying DynamoRIO code cache implementation of indirect jump instructions is relatively less efficient than implementations of other instruction patterns [5].

4.4.3 Patch Generation Time

On average ClearView took 4.9 minutes from the time of the first exposure to a new exploit to the time when it obtained a successful

Bugzilla Number	Shadow Stack + Heap Guard Runs	Building Invariant Checks [# of Checks]	Installing Invariant Checks	Invariant Check Runs (violated/total checks)	Building Repair Patches [# of Correlated Inv.s]	Installing Repair Patches	Unsuccessful Repair Runs (# of Runs)	Successful Repair Run	Total
269095	25.3	12.7 [1,0,1]	8.7	52.0 (4/28)	11.0 [1,0,0]	7.3	51.4 (2)	34.5	202.8
285595*	25.4	12.2 [0,5,0]	8.5	74.3 (6/2216)	11.5 [0,3,0]	8.8	-	31.8	172.4
290162	27.1	9.8 [2,0,0]	7.8	47.7 (2/2)	10.9 [1,0,0]	8.4	-	32.6	144.3
295854	32.8	8.8 [1,0,0]	9.2	66.3 (2/0)	10.3 [1,0,0]	8.1	31.1 (1)	39.8	206.5
296134	39.3	63.8 [0,42,10]	5.9	279.1 (?/?)	30.3 [0,?,?]	6.2	-	50.2	474.8
307259!	26.1	49.4 [0,4,26]	4.5	1235.5 (7444/29428)	39.7 [0,1,6]	6.3	347.7 (7)	-	1709.1
311710a	52.0	14.2 [0,1,2]	9.2	151.3 (60/1460)	11.3 [0,1,0]	6.8	-	69.1	313.9
311710b	60.5	13.5 [0,1,2]	8.3	152.3 (60/1460)	13.4 [0,1,0]	5.5	-	57.6	311.0
311710c	51.6	17.6 [0,1,2]	8.4	161.4 (60/1460)	16.2 [0,1,0]	8.2	-	64.0	327.3
312278	24.3	8.6 [1,0,0]	7.2	48.5 (2/0)	11.7 [1,0,0]	8.0	-	33.3	141.5
320182	25.3	12.7 [1,0,1]	8.7	52.0 (4/28)	11.0 [1,0,0]	7.3	51.4 (2)	34.5	202.8
325403*	24.2	16.9 [0,0,2]	5.9	46.8 (4/0)	10.6 [0,0,2]	6.0	-	33.5	143.9

Table 3: ClearView attack processing times, in seconds. Because all timing events were measured on the central ClearView Manager, the times include communication times between the protected client and the manager. Attacks for exploit 296134 were run on a slower computer (see Section 4.4.4). A * identifies the two exploits for which ClearView did not successfully generate a patch during the Red Team exercise, but did successfully generate a patch in subsequent experiments after reconfiguration. A ! identifies the exploit for which ClearView did not successfully generate a patch in either the Red Team exercise or in subsequent experiments. The row for exploit 296134 is missing several values — we did not record these values during the Red Team exercise, and are now unable to recreate the environment required to run the (rather fragile) exploit that the Red Team used.

patch for that exploit. This is not the time required to stop a propagating attack — Memory Firewall terminates the application before the attack can take effect, so there is no propagation. These times instead reflect how long users must wait before they have a patched version of the application that provides continuous, uninterrupted service even while under attack.

The 4.9 minutes includes an average of 5.4 executions: to detect the failure and select a set of candidate correlated invariants, to collect invariant-checking results to identify correlated invariants, and to evaluate candidate repairs. These averages include one outlier, for which ClearView took 13 minutes to sequentially correct three distinct errors in the application, all of which were exploited by the same attack — after ClearView repaired one error, the same exploit triggered an error from a different defect which ClearView then detected and repaired, and so on.

4.4.4 Patch Creation Time Breakdowns

When considering how much time ClearView takes to create a successful repair, one key comparison to keep in mind is the 28 days (on average) that it takes for developers to create and distribute a patch for a security exploit [47]. This section breaks down ClearView’s 4.9 minutes (on average) time to do the same.

Table 3 presents the different components of the time ClearView requires to generate a successful repair. All times are in seconds. With the exception of exploit 311710, there is one row for each exploit. The first column of each row presents the Bugzilla number of the exploit. Exploit 311710 has three rows (labeled 311710a, 311710b, and 311710c). As described above in Section 4.3.1, ClearView corrected three distinct errors before enabling Firefox to finally survive the attack. The table places each error in a separate row.

Shadow Stack + Heap Guard Runs: The second column of each row (Shadow Stack + Heap Guard Runs) in Table 3 presents the time required to replay the exploit to detection with the Shadow Stack and Heap Guard turned on. The vast majority of this time is spent warming up the Determina Managed Program Execution Environment code cache when we restart Firefox to process the exploit. For all exploits except 311710, this time is roughly 20 to 30 seconds. Because exploit 311710 exercises more Firefox functionality than the other exploits, the times are higher for this exploit.

It is possible for ClearView to successfully correct errors with

only Memory Firewall enabled. To our surprise, Heap Guard did not improve ClearView’s performance in the Red Team exercise — ClearView needs only Memory Firewall and Shadow Stack to generate successful patches for the seven exploits that it successfully patched during the Red Team exercise. Heap Guard is required for the remaining two exploits for which ClearView was able to subsequently generate successful patches after configuration changes.

Building and Installing Invariant Checks: The Building Invariant Checks column presents the time required to build all of the invariant check patches. This time includes compiling the automatically generated C source code for the patches and loading the patches into a DLL for presentation to the Determina patch management system. Each entry has the form $t[x,y,z]$, where t is the time required to build the invariant checks, x is the number of checked one-of invariants, y is the number of checked lower-bound invariants, and z is the number of checked less-than invariants. Exploit 296134 is an outlier, in part because ClearView compiled many more invariant check patches than for the other exploits and in part because the compiles took place on a slower computing platform. The Installing Invariant Checks column presents the time required for the Determina patch management system to transmit and apply the patches to the application running on the client machine.

Invariant Check Runs: The Invariant Check Runs column presents the time required to replay the exploit to detection twice with the invariant check patches in place. Each entry has the form $t(x/y)$, where t is the time required to replay the exploit to detection twice, x is the number of times a checked invariant was violated during these runs, and y is the total number of invariant checks executed during these runs. The time t includes the time required to communicate the necessary invariant check, shadow stack, and attack location information to the ClearView Manager. As with the Shadow Stack + Heap Guard Runs, much of the time was spent warming up the Determina Managed Program Execution Environment code cache. For exploit 296134, ClearView also spent a substantial amount of time communicating invariant check results to the ClearView Manager using the Windows event queue mechanism.

Building and Installing Repair Patches: The Building Repair Patches column presents the time required to build all of the repair patches for the correlated invariants to evaluate. Each entry has the form $t[x,y,z]$, where t is the time required to build the repair patches, x is the number of correlated one-of invariants, y is

the number of correlated lower-bound invariants, and z is the number of correlated less-than invariants. All of the correlated one-of invariants involve function pointers. As described above in Section 2.5.1, such invariants have three potential repairs. ClearView compiles a patch for each such repair. The other kinds of invariants each have a single repair with a single repair patch. The Installing Repair Patches column presents the time required to communicate these patches to the client machine. The client machine applied specific repair patches one at a time in response to directives from the central ClearView Manager.

Unsuccessful Repair Runs: The Unsuccessful Repair Runs column presents the time (if any) required to execute Firefox to completion for any unsuccessful repair patches. Each entry has the form $t(x)$, where t is the time required to execute Firefox to completion and x is the number of unsuccessful runs (if any). We attribute the small number of unsuccessful runs to the effectiveness of the correlated invariant selection policy in targeting invariants whose repairs are likely to correct the error, and to the effectiveness of the candidate repair ordering rule in selecting an effective repair to evaluate first.

Successful Repair Run: The Successful Repair Run column shows the time required to execute Firefox with the successful repair patch applied. Because this patch corrects the error and eliminates the attack detection, the patch was judged to succeed ten seconds after it executed with no subsequent attack detection. The presented time includes this ten seconds. At this point, ClearView had generated and identified a patch that corrected the error and enabled continued successful execution.

The final column is the sum of the times in the other columns. It presents the total time required to automatically obtain a successful patch for the corresponding attack.

An alternate deployment environment would enable only Memory Firewall during production use, with Heap Guard and the Shadow Stack enabled only after the detection of the first attack. In this case, the application’s normal overhead would be 1.5 instead of 3.0 (see Table 2), but the total time would be increased by adding one Memory Firewall run.

4.4.5 Reducing ClearView’s Overhead

There are three primary sources of inefficiency in the current ClearView attack response system: warming up the code cache (all “Run” columns), using Windows event queues as the communication mechanism between community members and the central server (“Installing” and “Run” columns), and compiling the invariant check and repair patches (“Building” columns). It is possible to eliminate the cache warm up time by saving the cache state from a previous run, then restoring this state upon startup. It is possible to dramatically reduce the communication time by using a more efficient communication mechanism. It is possible to eliminate the compilation time by generating binary code directly instead of generating, then compiling C code. Together, we estimate that these optimizations would enable ClearView to produce successful patches in tens of seconds, rather than in minutes.

4.5 Other Applications

We expect the Firefox results to be broadly representative of ClearView’s behavior for other server applications. Both failure-oblivious computing [35, 33, 36] and transactional function termination [42, 41] enable a range of servers to survive errors and attacks. Some of the ClearView patches have a similar effect as these techniques. Because ClearView is based on inferred invariants that capture aspects of the application’s semantics, it is able to generate more targeted and therefore potentially more effective repairs. ClearView also incorporates a broader range of repair strategies and

evaluates the resulting multiple candidate repairs to discard ineffective or damaging repairs, which may enhance its ability to find successful patches.

Survival strategies are not applicable in all circumstances. Survival strategies are best suited to applications with high availability requirements, and those that can tolerate some variation in the computational result (such as information retrieval, or processing sensory data). Survival techniques are less appropriate for applications (such as compiler transformations) with precise, logically defined correctness requirements, long dependence chains that run through the entire computation, and less demanding availability requirements. Even for such applications, ClearView may have a place as a debugging aid.

5. LIMITATIONS

The goal of ClearView is not to correct every conceivable error. The goal is instead to correct a realistic class of errors to enable applications with high availability requirements to successfully provide service in spite of these errors.

There exist errors that are completely outside the scope of ClearView, i.e., errors for which there is no plausible learned invariant whose enforcement would enable the application to survive. But even if the error is within the scope of the overall ClearView approach, ClearView may be unable to find a repair that enables the application to survive the error:

- **Learning:** Daikon comes preconfigured to learn a specific set of invariants. (End users may change the configuration or even define new invariant templates.) This set may not include any invariant that enables ClearView to generate a patch for a given error. And even if the set does include such an invariant, the learning phase may not provide enough coverage of the application to enable Daikon to learn this invariant.
- **Monitoring:** ClearView currently uses Memory Firewall to detect control-flow-transfer errors and Heap Guard to detect out-of-bounds-write errors. Additional detectors would be required to detect other kinds of errors.
- **Candidate Invariant Selection:** Every repair enforces an invariant that ClearView selected as a candidate correlated invariant. Even if ClearView inferred an invariant with a repair that would correct the error, the failure may occur sufficiently far from the error for ClearView to not include the invariant in the set of candidate correlated invariants, or to delay ClearView’s response.
- **Repair:** The repair mechanism comes with a specific set of invariant enforcement mechanisms; each such mechanism corresponds to a specific repair strategy. It is possible that none of these repair strategies produces a successful repair.

It is also possible for ClearView to impair the functionality of the application or even to create new vulnerabilities:

- **Functionality Impairment:** It is possible for a ClearView repair patch to impair the functionality of the application. If the patch is applied in response to a legitimate attack, the functionality impairment may be a reasonable price to pay for eliminating the vulnerability. Two facts minimize the likelihood that ClearView will apply a patch in the absence of an error or attack: (1) ClearView applies patches only in response to a detected failure — in the current implementation, an illegal control flow transfer or out-of-bounds write. (2) ClearView enforces the invariant only if it is correlated with the failure. These two facts also minimize the likelihood that an applied ClearView patch will interfere with the processing of a legitimate input. The ClearView patch evaluation mechanism enables ClearView to recognize and discard patches that do not eliminate the failure

or cause the application to crash. ClearView generated several patches with negative effects during the Red Team exercise. The ClearView patch evaluation mechanism detected these negative effects and discarded the patches.

- **Patch Subversion:** It is theoretically possible for an adversary to subvert the ClearView patch mechanism to install its own malicious patches. We note that ClearView builds on the commercially deployed Determina patch distribution mechanism, which uses standard authentication and encryption mechanisms to ensure patch integrity.
- **Malicious Nodes:** It is possible for a malicious node or nodes to provide ClearView with false information that may cause it to generate an inappropriate patch. It is possible to mitigate this possibility via statistical methods and by reproducing the error and evaluating the generated patches on trusted nodes before distributing the patches throughout the community. Malicious nodes were not part of the threat model for our Red Team Exercise but are an interesting avenue for future research.

6. RELATED WORK

We discuss additional related work in invariant inference, attack detection mechanisms, automatic filter generation, checkpoint and replay techniques, and error tolerance and correction.

6.1 Invariant Inference and Immunity

ClearView uses an inferred specification to focus the patching process and correct errors. We discuss several other projects that use inferred specifications to improve system behavior or security.

Our system for automated program steering [25] uses machine learning over correct executions, then changes the control flow of bad executions. When an application fails or underperforms, the system chooses a different operational mode that has been successful in the past, in situations like the current situation.

We previously developed a system [13] that automatically inferred data structure consistency constraints and created repairs [15] to enforce them. This system was also evaluated by a (different) hostile Red Team. In contrast to ClearView, this previous system required source code, performed learning only once without subsequent refinements, applied only repairs that it statically verified to always terminate in a repaired state, and did not observe subsequent executions to evaluate the quality of the repairs.

FFTV (From Failures to Vaccine) [28] infers invariants that model normal behavior. At run time, violations of a programmer-written specification cause FFTV to record a “failure context”, consisting of violated invariants and the stack backtrace. If a failure context is encountered again, the next method call is performed transactionally. Dimmunix [22] prevents programs from re-entering previously-seen deadlock states. When the application deadlocks, Dimmunix records a “deadlock signature” indicating each thread’s stack and the locks it holds. If a lock acquisition would cause the signature is encountered again, then the thread is made to yield instead. As with FFTV, false positives in the learning process affect the performance but not the correctness of the application.

Gibraltar [3] detects rootkits by observing their effect on dynamically inferred specifications that are derived from normal executions. Gibraltar can detect rootkits that affect only data structures (not control), which can elude other approaches to rootkit detection. This is an application of invariant detection for anomaly detection, but unlike the above projects, it does not correct the problem.

6.2 Attack Detection Mechanisms

The current ClearView implementation uses two attack detection techniques: program shepherding to detect and block malicious control flow transfers, and heap overflow checks to detect and

block out-of-bounds writes to the heap. ClearView can work with any attack detection technique that provides an attack location.

StackGuard [10] and StackShield [45], for example, use a modified compiler to generate code to detect attacks that overwrite the return address on the stack. StackShield also performs range checks to detect overwritten function pointers. Researchers have also built compilers that insert bounds checks to detect memory addressing errors in C programs [2, 49, 7, 20, 37, 21, 23]. Drawbacks of these techniques include the need to recompile the program, the overhead of the dynamic bounds checks, and, in some cases, the need to change the program itself [7, 23]. Dynamic taint analysis finds appearances of potentially malicious data in sensitive locations such as function pointers or return addresses [46, 11, 29]. It would be possible to make ClearView work with all of these detectors, although the high overhead and potential need for recompilation or even source code changes goes against ClearView’s philosophy of operating on stripped Windows binaries and minimizing the overhead during normal execution.

An expensive attack detection mechanisms can be distributed across a community of machines, with each application instrumenting only a small portion of its execution [27], or applied only to a honeypot [44, 39, 1]. These approaches can leave the production versions unprotected against new attacks.

In contrast to systems that apply their attack analyses across broad ranges of the application, ClearView uses the attack location to dramatically narrow down the region of the application that it instruments during its attack analysis and response generation activities. This makes it possible to deploy sophisticated but expensive analyses within this focused region of the application while still keeping the total overhead small. This is effective only if ClearView identifies the invariant required to correct the underlying error in the application logic.

6.3 Automatic Filter Generation

A standard way to protect applications against attacks is to develop filters that detect and discard exploits before they reach a vulnerable application. Vigilante uses honeypots to detect attacks and dynamically generate filters that check for exploits that follow the same control-flow path as the attack to exploit the same vulnerability [9]. Bouncer uses symbolic techniques to generalize Vigilante’s approach to filter out more exploits [8]. ShieldGen uses Vigilante’s attack detection techniques to obtain exploits [12]. It generates variants of each exploit and tests the variants to see if they also exercise the vulnerability. It then produces a general filter that discards all such variants.

Sweeper uses address randomization for efficient attack detection [48]. This technique is efficient enough to be deployed on production versions of applications, but provides only probabilistic protection and therefore leaves applications still vulnerable to exploitation [38]. Sweeper uses attack replay in combination with more expensive attack analysis techniques such as memory access checks, dynamic taint analysis, and dynamic backward slicing. It uses the information to generate filters that discard exploits before they reach vulnerable applications. Sweeper also builds vulnerability-specific execution filters, which instrument selected instructions involved in the attack to detect the attack. The attack response is to use rollback plus replay to recover from the attack.

Discarding inputs that may contain exploits can deny users access to content that they need or want to access. It is possible for useful information sources such as legitimate web pages or images to become surreptitiously infiltrated with exploits. Attackers may also create attractive content containing exploits, specifically for the purpose of enticing users to access it via vulnerable applications. Just as some Internet content is created for the purpose of

generating advertising revenue, an alternate (and presumably illegal) business model would substitute attacks for advertisements. In all of these cases it is desirable for applications to process inputs containing exploits without enabling the attacks to succeed.

By enabling applications to execute successfully through otherwise exploitable errors, ClearView can enable users to access information even if the input containing the information also contains an exploit. The ClearView repair for one of the heap overflow errors in the Red Team exercise, for example, makes it possible for users to view image files that also contain exploits (or that contain innocent data that happens to exercise the vulnerability).

6.4 Checkpoint and Replay

A traditional and widely-used error recovery mechanism is to reboot the system, with operations replayed as necessary to bring the system back up to date [19]. It may also be worthwhile to recursively restart larger and larger subsystems until the system successfully recovers [6]. Checkpointing [26] can improve the performance of the basic reboot process and help minimize the amount of lost state in the absence of replay. Checkpointing also makes it possible to discard the effects of attacks and errors to restore the system to a previously saved clean operational state. This approach, in some cases combined with replay of previously processed requests or operations that do not contain detected exploits, has been proposed as an attack response mechanism [43, 32, 40, 41, 48].

In comparison with ClearView’s approach of continuing to execute through attacks, checkpoint plus replay has several drawbacks. These include service interruptions as the system recovers from an attack (these service interruptions can occur repeatedly unless the system is otherwise protected against repeated attacks), the potential for replay to fail because of problematic interactions with external processes or machines that are outside the scope of the checkpoint and replay mechanism, lost state if the system chooses to forgo replay, and complications associated with applying checkpoint and replay to multithreaded or multiprocess applications.

6.5 Error Tolerance and Correction

Techniques that allow applications to tolerate errors such as out-of-bounds memory accesses include transactional function termination [42, 39] (which restores the state at the time of the function call, then returns from functions that perform out-of-bounds accesses), failure-oblivious computing [35] (which discards out-of-bounds writes and manufactures values for out-of-bounds reads), boundless memory blocks [34] (which store out-of-bounds writes in a hash table for subsequent corresponding reads to access), and DieHard [4] (which overprovisions the heap so that out-of-bounds accesses are likely to fall into otherwise unused memory). These techniques require no learning phase and no repeated executions for correlated invariant selection and evaluation. ClearView differs in that it applies checks and repairs only in response to attacks, it does so in carefully targeted parts of the application identified by correlated invariants, and it performs an ongoing evaluation of each applied repair. ClearView may also be able to provide more informative error reports, since it can identify specific invariants whose violation is correlated with errors.

ASSURE [41] generalizes transactional function termination to enable the system to transactionally terminate any one of the functions on the call stack at the time of the error (and not just the function containing the error). Attack replay on a triage machine enables the system to evaluate which function to terminate to provide the most successful recovery. The applied patch takes a checkpoint at the start of the function. It responds to errors by restoring the checkpoint, then returning an effective error code to terminate the function and continue execution at the caller.

Exterminator [30] uses address space randomization to detect out-of-bounds writes into the heap and accesses via dangling references. It then corrects the errors by (as appropriate) increasing the size of allocated memory blocks to accommodate the observed out-of-bounds writes or delaying memory block deallocation until the accesses via the corresponding dangling references have completed. Exterminator can combine patches from multiple users to give members of a community of users immunity to out-of-bounds-write and dangling-reference errors without prior exposure.

Genetic programming techniques can generate and search a space of abstract syntax tree modifications with the goal of automatically correcting an exercised defect in the underlying program [18]. A test suite is used to evaluate the effectiveness of each generated abstract syntax tree in preserving desirable behavior and eliminating undesirable behavior.

All of these techniques can take the application outside of its anticipated operating envelope. They therefore have the potential to introduce new errors. Of course, ClearView’s patches also have the potential to negatively affect the application. Because ClearView’s patches only affect properties that are correlated with undesirable behavior, ClearView is less likely to create bad patches. Because ClearView performs an ongoing evaluation of each deployed patch, it will quickly discard patches that enable negative effects (such as crashes or attacks) in favor of more effective patches.

7. CONCLUSION

Errors in deployed software systems pose an important threat to the integrity and utility of our computing infrastructure. Relying on manual developer intervention to find and eliminate errors can deny service to application users or even leave the application open to exploitation for long periods of time. ClearView’s automatic error detection and correction techniques can provide, with no human intervention whatsoever, the almost immediate correction of errors, including errors that enable newly-released security attacks. The result is an application that is immune to the attack and can continue to provide uninterrupted service.

ClearView is targeted toward applications with high availability requirements, for which a small chance of unexpected behavior is preferable to the certainty of denial of service. The feasibility of our approach has been established by a hostile Red Team evaluation in which ClearView automatically patched security vulnerabilities without introducing new attack vectors that the Red Team could exploit. We acknowledge that there are important reliability and security problems that are outside the scope of ClearView. Nevertheless, ClearView addresses an important and realistic problem, and holds out the promise of substantially improving the integrity and availability of our computing infrastructure.

Acknowledgments. This work was done while all authors except Larsen and Amarasinghe were at MIT. The research was supported under DARPA cooperative agreement FA8750-06-2-0189 “Collaborative Learning for Security and Repair in Application Communities”. We thank our DARPA program manager, Lee Badger, for his enthusiastic support and encouragement. Lee’s vision and willingness to support a new and controversial approach were critical to the overall success of the project. We thank the outside evaluation team, Cordell Green, Hilarie Orman, and Alex Orso, for their feedback on the project. We thank the White Team, Chris Doh and Dave Kepler, for their fair and impartial management of the Red Team exercise. We thank the hostile Red Team, led by Gregg Tally, for their contributions to the evaluation of ClearView. The Red Team exploit development team included Suresh Krishnaswamy, Robert Lyda, Peter Lovell, David Sames, Richard Toren, and Wayne Morrison. The Red Team negative effect evaluation

team included Dave Balenson, Walt Coleman, Vijaya Ramamurthi, and Richard Toren. We thank Sandy Wilbourn, Derek Bruening, Vladimir Kiriansky, Rishi Bidarkar, and Bharath Chandramohan of Determina for their contributions to this project. We thank George Candea, Andreas Zeller, Andrzej Wasylkowski, and the anonymous reviewers for their comments on a draft of this paper.

8. REFERENCES

- [1] ANAGNOSTAKIS, K., SIDIROGLOU, S., AKRITIDIS, P., XINIDIS, K., MARKATOS, E., AND KEROMYTIS, A. D. Detecting targeted attacks using shadow honeypots. In *USENIX Security* (Aug. 2005).
- [2] AUSTIN, T., BREACH, S., AND SOHI, G. Efficient detection of all pointer and array access errors. In *PLDI* (June 2004).
- [3] BALIGA, A., GANAPATHY, V., AND IFTODE, L. Automatic inference and enforcement of kernel data structure invariants. In *ACSAC* (Dec. 2008), pp. 77–86.
- [4] BERGER, E., AND ZORN, B. DieHard: probabilistic memory safety for unsafe languages. In *PLDI* (June 2006).
- [5] BRUENING, D. *Efficient, Transparent, and Comprehensive Runtime Code Manipulation*. Ph.D., MIT Department of Electrical Engineering and Computer Science, Cambridge, MA, Sep. 2004.
- [6] CANDEA, G., AND FOX, A. Recursive restartability: Turning the reboot sledgehammer into a scalpel. In *HotOS* (Schloss Elmau, Germany, May 2001).
- [7] CONDIT, J., HARREN, M., MCPPEAK, S., NECULA, G. C., AND WEIMER, W. CCured in the real world. In *PLDI* (June 2003).
- [8] COSTA, M., CASTRO, M., ANTONY, ZHOU, L., ZHANG, L., AND PEINADO, M. Bouncer: securing software by blocking bad input. In *SOSP* (Oct. 2007).
- [9] COSTA, M., CROWCROFT, J., CASTRO, M., ROWSTRON, A., ZHOU, L., ZHANG, L., AND BARHAM, P. Vigilante: End-to-end containment of Internet worms. In *SOSP* (Oct. 2005).
- [10] COWAN, C., PU, C., MAIER, D., WALPOLE, J., BAKKE, P., BEATTIE, S., GRIER, A., WAGLE, P., ZHANG, Q., AND HINTON, H. StackGuard: Automatic adaptive detection and prevention of buffer-overflow attacks. In *USENIX Security* (January 1998).
- [11] CRANDALL, J., AND CHONG, F. Minos: Control data attack prevention orthogonal to memory model. In *MICRO* (Dec. 2004).
- [12] CUI, W., PEINADO, M., WANG, H. J., AND LOCASO, M. E. ShieldGen: Automatic data patch generation for unknown vulnerabilities with informed probing. In *IEEE S&P* (May 2007).
- [13] DEMSKY, B., ERNST, M. D., GUO, P. J., MCCAMANT, S., PERKINS, J. H., AND RINARD, M. Inference and enforcement of data structure consistency specifications. In *ISSTA* (July 2006).
- [14] DEMSKY, B., AND RINARD, M. Automatic detection and repair of errors in data structures. In *18th Annual ACM Conference on Object-Oriented Programming, Systems, Languages, and Applications* (Oct. 2003).
- [15] DEMSKY, B., AND RINARD, M. Data structure repair using goal-directed reasoning. In *ICSE* (May 2005).
- [16] ERNST, M. D., COCKRELL, J., GRISWOLD, W. G., AND NOTKIN, D. Dynamically discovering likely program invariants to support program evolution. *IEEE TSE* 27, 2 (Feb. 2001).
- [17] ERNST, M. D., PERKINS, J. H., GUO, P. J., MCCAMANT, S., PACHECO, C., TSCHANTZ, M. S., AND XIAO, C. The Daikon system for dynamic detection of likely invariants. *Science of Computer Programming* 69, 1–3 (Dec. 2007).
- [18] FORREST, S., WEIMER, W., NGUYEN, T., AND GOUES, C. L. A genetic programming approach to automated software repair. In *GECCO* (July 2009).
- [19] GRAY, J., AND REUTER, A. *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, 1993.
- [20] JIM, T., MORRISSETT, G., GROSSMAN, D., HICKS, M., CHENEY, J., AND WANG, Y. Cyclone: A safe dialect of C. In *USENIX* (June 2002).
- [21] JONES, R., AND KELLY, P. Backwards-compatible bounds checking for arrays and pointers in C programs. In *AADEBUG* (May 1997).
- [22] JULA, H., TRALAMAZZA, D., ZAMFIR, C., AND CANDEA, G. Deadlock immunity: Enabling systems to defend against deadlocks. In *OSDI* (Dec. 2008), pp. 295–308.
- [23] KENDALL, S. C. Bcc: Run-time checking for C programs. In *USENIX Summer* (1983).
- [24] KIRIANSKY, V., BRUENING, D., AND AMARASINGHE, S. Secure execution via program shepherding. In *USENIX Security* (Aug. 2002).
- [25] LIN, L., AND ERNST, M. D. Improving adaptability via program steering. In *ISSTA* (July 2004).
- [26] LITZKOW, M., AND SOLOMON, M. The evolution of condor checkpointing. In *Mobility: processes, computers, and agents* (1999). ACM Press/Addison-Wesley.
- [27] LOCASO, M. E., SIDIROGLOU, S., AND KEROMYTIS, A. D. Software self-healing using collaborative application communities. In *NDSS* (Feb. 2005).
- [28] LORENZOLI, D., MARIANI, L., AND PEZZÈ, M. Towards self-protecting enterprise applications. In *ISSRE* (Nov. 2007), pp. 39–48.
- [29] NEWSOME, J., AND SONG, D. Dynamic taint analysis for automatic detection, analysis, and signature generation of exploits on commodity software. In *NDSS* (Feb. 2005).
- [30] NOVARK, G., BERGER, E., AND ZORN, B. Exterminator: Automatically correcting memory errors with high probability. *Communications of the ACM* 51, 12 (Dec. 2008).
- [31] PERKINS, J. H., AND ERNST, M. D. Efficient incremental algorithms for dynamic detection of likely invariants. In *FSE* (Nov. 2004).
- [32] QIN, F., TUCEK, J., SUNDARESAN, J., AND ZHOU, Y. Rx: treating bugs as allergies—a safe method to survive software failures. *SIGOPS Oper. Syst. Rev.* 39, 5 (2005), 235–248.
- [33] RINARD, M. Acceptability-oriented computing. In *OOPSLA Companion* (Oct. 2003).
- [34] RINARD, M., CADAR, C., DUMITRAN, D., ROY, D. M., AND LEU, T. A dynamic technique for eliminating buffer overflow vulnerabilities (and other memory errors). In *ACSAC* (Dec. 2004).
- [35] RINARD, M., CADAR, C., DUMITRAN, D., ROY, D. M., LEU, T., AND WILLIAM S. BEEBEE, J. Enhancing server availability and security through failure-oblivious computing. In *OSDI* (December 2004).
- [36] RINARD, M., CADAR, C., AND NGUYEN, H. H. Exploring the acceptability envelope. In *OOPSLA Companion* (Oct. 2005).
- [37] RUWASE, O., AND LAM, M. S. A practical dynamic buffer overflow detector. In *NDSS* (February 2004).
- [38] SHACHAM, H., PAGE, M., PFAFF, B., GOH, E.-H., MODADUGU, N., AND BONEH, D. On the effectiveness of address-space randomization. In *ACM CCS* (Oct. 2004).
- [39] SIDIROGLOU, S., GIOVANIDIS, G., AND KEROMYTIS, A. D. A dynamic mechanism for recovering from buffer overflow attacks. In *ISC* (Sep. 2005).
- [40] SIDIROGLOU, S., LAADAN, O., KEROMYTIS, A. D., AND NIEH, J. Using rescue points to navigate software recovery. In *IEEE S&P* (May 2007).
- [41] SIDIROGLOU, S., LAADAN, O., PEREZ, C., VIENNOT, N., NIEH, J., AND KEROMYTIS, A. D. Assure: automatic software self-healing using rescue points. In *ASPLOS '09* (2009).
- [42] SIDIROGLOU, S., LOCASO, M. E., BOYD, S. W., AND KEROMYTIS, A. D. Building a reactive immune system for software services. In *USENIX* (Apr. 2005).
- [43] SMIRNOV, A., AND CHIUEH, T. DIRA: Automatic detection, identification and repair of control-hijacking attacks. In *NDSS* (Feb. 2005).
- [44] SPITZNER, L. *Honeypots: Tracking Hackers*. Addison-Wesley, 2002.
- [45] Stackshield. www.angelfire.com/sk/stackshield.
- [46] SUH, G., LEE, J., ZHANG, D., AND DEVADAS, S. Secure program execution via dynamic information flow tracking. In *ASPLOS* (Oct. 2004).
- [47] Symantech Internet security threat report. www.symantec.com, Sep. 2006.
- [48] TUCEK, J., NEWSOME, J., LU, S., HUANG, C., XANTHOS, S., BRUMLEY, D., ZHOU, Y., AND SONG, D. Sweeper: A lightweight end-to-end system for defending against fast worms. In *EuroSys* (Mar. 2007).
- [49] YONG, S. H., AND HORWITZ, S. Protecting C programs from attacks via invalid pointer dereferences. In *ESEC/FSE* (2003).