

# Lecture Notes on Memory Safety & Sandboxing

Matt Fredrikson

Carnegie Mellon University  
Lecture 5

## 1 Introduction & Recap

In the previous lecture we looked into proving that programs satisfy safety properties given as formulas in the first-order dynamic logic. In particular, we can write contract properties with precondition  $P$  and postconditions  $Q$  for a program  $\alpha$  as:

$$P \rightarrow [\alpha]Q \tag{1}$$

If this formula is valid, then it means that in every state  $\omega$ , if  $\omega \models P$  then after all terminating runs of  $\alpha$  starting in  $\omega$  the final state  $\nu \models Q$ .

So far the programs that we have studied are not too interesting when it comes to concerns about security—they can perform operations over integers, but they can't interact with memory or arrays. Today we will start by extending our language to support pointer operations over memory, and then discuss memory safety as well as a more flexible mechanism for policy enforcement called software fault isolation. Importantly, we will build on our understanding of sound axioms and proof rules to justify the correctness of this mechanism.

## 2 Memory Safety

While most imperative programming languages support convenient dynamic memory allocation and access with syntax like `malloc` and `a[i]`, at the end of the day this is nothing more than syntactic sugar for managing a large integer-indexed array of values. We can add basic support for this to our language by introducing pointers, and adding an integer-indexed memory to our program state. Now terms in our language will have the following syntax.

$$e, \tilde{e} ::= x \mid c \mid e + \tilde{e} \mid e \cdot \tilde{e} \mid \star(e) \tag{2}$$

The term  $\ast(e)$  denotes the value obtained by evaluating  $e$  in the current state, and accessing the memory at the corresponding index. This takes care of reading from the memory array, now we add support for updating memory by introducing a new type of program command.

$$\alpha, \beta ::= x := e \mid \ast(e) := \tilde{e} \mid \mathbf{assert}(Q) \mid \mathbf{if}(Q) \alpha \mathbf{else} \beta \mid \alpha; \beta \mid \mathbf{while}(Q) \alpha \quad (3)$$

The command  $\ast(e) := \tilde{e}$  evaluates  $e$  and  $\tilde{e}$  in the current state, and sets the value of memory indexed at the value of  $e$  to the value of  $\tilde{e}$ .

Now for the semantics. We will need to track the value of variables as we did before with a mapping from variables to values. But we will also need to track the state of the memory, which we will formalize as a partial mapping from non-negative integers to values. Real machines don't have unlimited memory, which is why the mapping is partial: we assume that the memory can hold at most  $U$  values, so the mapping is only defined on  $0 \leq i \leq U$ .

We will continue to denote states by  $\omega$ , and write  $\omega_V(x)$  to refer to the value of the variable mapping, and  $\omega_M(x)$  to refer to the memory array. The semantics of terms can now be defined as follows.

**Definition 1 (Semantics of terms).** The semantics of a term  $e$  in a state  $\omega$  is its value  $\omega\llbracket e \rrbracket$ . It is defined inductively by distinguishing the shape of term  $e$  as follows:

- $\omega\llbracket x \rrbracket = \omega_V(x)$  for variable  $x$
- $\omega\llbracket c \rrbracket = c$  for number literals  $c$
- $\omega\llbracket e \odot \tilde{e} \rrbracket = \omega\llbracket e \rrbracket \odot \omega\llbracket \tilde{e} \rrbracket$ , where  $\odot \in \{+, \times\}$
- $\omega\llbracket \ast(e) \rrbracket = \omega_M(\omega\llbracket e \rrbracket)$  if  $0 \leq \omega\llbracket e \rrbracket < U$ , else undefined

Adding pointers to our language has led to a complication: now terms can be undefined. Specifically, if  $e$  evaluates to a negative number, or a number larger than the maximum memory size  $U$ , then the term  $\ast(e)$  is not defined.

This complication manifests in how we define the semantics of formulas. Because terms can now be undefined in certain states, we need to account for this in the semantics of formulas that might include terms. Whenever a term in a formula is undefined in a particular state, then the value of the formula is as well.

**Definition 2 (Semantics of arithmetic formulas).** The DL formula  $P$  is true in state  $\omega$ , written  $\omega \models P$ , as inductively defined by distinguishing the shape of formula  $P$ :

1.  $\omega \not\models \perp$ , i.e.,  $\perp$  is true in no states
2.  $\omega \models \top$ , i.e.,  $\top$  is true in all states
3.  $\omega \models e = \tilde{e}$  iff  $\omega\llbracket e \rrbracket = \omega\llbracket \tilde{e} \rrbracket$  and both terms are defined in  $\omega$ .
4.  $\omega \models e \leq \tilde{e}$  iff  $\omega\llbracket e \rrbracket \leq \omega\llbracket \tilde{e} \rrbracket$  and both terms are defined in  $\omega$ .

5.  $\omega \models P \wedge Q$  iff  $\omega \models P$  and  $\omega \models Q$  if  $P$  and  $Q$  are defined in  $\omega$ .
6.  $\omega \models P \vee Q$  iff  $\omega \models P$  or  $\omega \models Q$  if  $P$  and  $Q$  are defined in  $\omega$ .
7.  $\omega \models \neg P$  iff  $\omega \not\models P$  if  $P$  is defined in  $\omega$ .
8.  $\omega \models P \rightarrow Q$  iff  $\omega \not\models P$  or  $\omega \models Q$  and  $P$  and  $Q$  are defined in  $\omega$ .
9.  $\omega \models P \leftrightarrow Q$  iff both are true or both false and  $P$  and  $Q$  are defined in  $\omega$ .

Finally, we get to the semantics of programs. Obviously we need to add a new definition for the memory update command  $\ast(e) := \tilde{e}$ . But programs may contain terms and formulas, which we now know can be undefined in some states. We define the semantics of a program with a term or formula that is undefined in a state as aborting in the next subsequent state.

First some notation. If  $\omega_M$  is a memory in state  $\omega$ , then we write  $\omega_M\{e \mapsto \tilde{e}\}$  to denote the new memory obtained by copying  $\omega_M$ , and changing its mapping at  $\omega[e]$  to map to  $\omega[\tilde{e}]$ . So suppose that  $\omega_M(0) = 1, \omega_M(1) = 2$ . Then  $\omega_M\{1 \mapsto 3\}(0) = 1$  and  $(\omega_M\{1 \mapsto 3\})(1) = 3$ . We can apply this update notation multiple times, so that:

$$\omega_M\{1 \mapsto 3\}\{0 \mapsto 4\}(0) = 4, \omega_M\{1 \mapsto 3\}\{0 \mapsto 4\}(1) = 3$$

We'll adopt the convention that the rightmost update to a particular index is the one that we use when looking up values. So for example,

$$\omega_M\{1 \mapsto 3\}\{1 \mapsto 4\}(1) = 4$$

**Definition 3** (Trace semantics of programs). The trace semantics  $\llbracket \alpha \rrbracket$  of a program  $\alpha$  is the set of all its possible traces and is defined inductively as follows:

1.  $\llbracket x := e \rrbracket = \{(\omega, \nu) : \omega[e] \text{ is defined and } \nu = \omega \text{ except that } \nu_V(x) = \omega[e]\} \cup \{(\omega, \Lambda) : \omega[e] \text{ is not defined}\}$
2.  $\llbracket \ast(e) := \tilde{e} \rrbracket = \{(\omega, \nu) : 0 \leq \omega[e] \leq U, \omega[\tilde{e}] \text{ defined}, \nu_M = \omega_M\{\omega[e] \mapsto \omega[\tilde{e}]\}\} \cup \{(\omega, \Lambda) : \neg(0 \leq \omega[e] \leq U) \text{ or } \omega[\tilde{e}] \text{ not defined}\}$
3.  $\llbracket \text{assert}(Q) \rrbracket = \{(\omega, \omega) : \omega[e] \text{ is defined and } \omega \models Q\} \cup \{(\omega, \Lambda) : \omega \models Q \text{ is not defined or } \omega \not\models Q\}$
4.  $\llbracket \text{if}(Q) \alpha \text{ else } \beta \rrbracket = \{\sigma \in \llbracket \alpha \rrbracket : \sigma_0[e] \text{ is defined and } \sigma_0 \models Q\} \cup \{\sigma \in \llbracket \beta \rrbracket : \sigma_0[e] \text{ is defined and } \sigma_0 \not\models Q\} \cup \{(\omega, \Lambda) : \omega \models Q \text{ is not defined}\}$
5.  $\llbracket \alpha; \beta \rrbracket = \{\sigma \circ \varsigma : \sigma \in \llbracket \alpha \rrbracket, \varsigma \in \llbracket \beta \rrbracket\}$ ;  
the composition of  $\sigma = (\sigma_0, \sigma_1, \sigma_2, \dots)$  and  $\varsigma = (\varsigma_0, \varsigma_1, \varsigma_2, \dots)$  is

$$\sigma \circ \varsigma := \begin{cases} (\sigma_0, \dots, \sigma_n, \varsigma_1, \varsigma_2, \dots) & \text{if } \sigma \text{ terminates in } \sigma_n \text{ and } \sigma_n = \varsigma_0 \\ \sigma & \text{if } \sigma \text{ does not terminate} \end{cases}$$

$$\begin{aligned}
6. \llbracket \text{while}(Q) \alpha \rrbracket = & \{ \sigma^{(0)} \circ \dots \circ \sigma^{(n)} : \text{for all } 0 \leq i \leq n: \sigma_0^{(i)} \models Q, \sigma^{(i)} \in \llbracket \alpha \rrbracket, \text{ and} \\
& \sigma^{(n)} \text{ either doesn't terminate, or terminates with } \sigma_m^{(n)} \not\models Q \} \cup \\
& \{ \sigma^{(0)} \circ \sigma^{(1)} \circ \sigma^{(2)} \circ \dots : \text{for all } i \in \mathbb{N}: \sigma_0^{(i)} \models Q, \sigma^{(i)} \in \llbracket \alpha \rrbracket \} \cup \\
& \{ (\omega) : \omega \not\models Q \} \cup \\
& \{ \sigma^{(0)} \circ \dots \circ \sigma^{(n)} \circ (\Lambda) : \text{for all } 0 \leq i \leq n: \sigma_0^{(i)} \models Q, \sigma^{(i)} \in \llbracket \alpha \rrbracket, \text{ and} \\
& \sigma^{(n)} \text{ terminates with } \sigma_m^{(n)} \models Q \text{ not defined} \} \cup \\
& \{ (\omega, \Lambda) : \omega \models Q \text{ not defined} \}
\end{aligned}$$

While it may be tedious to track the presence of undefined terms and formulas through the evaluation of a program, we will see that this is central to the very definition of what memory safety means for a particular programming language.

**Axioms and Proof Rules.** Now we have semantics for programs with pointers and indexed memory, the next logical thing to do is find some useful axioms to help us reason about them.

Just as we had an axiom for assignment to variables, we have a similar axiom for updates to a pointer. But in the assignment axiom, we performed a syntactic substitution of the target variable in the postcondition. In this case we can readily see that looking for mere occurrences of a pointer expression will not suffice. Consider the following:

$$[x := 1; y := 1; *(x) := 0] * (y) \neq 0 \quad (4)$$

After executing the first two assignments,  $*(x)$  and  $*(y)$  point to the same memory location. So if we tried to close out a proof like the following:

$$\begin{array}{c}
*(y) \neq 0, x = 1, y = 1 \vdash *(y) \neq 0 \\
\text{[*=]} \hline
*(y) \neq 0, x = 1, y = 1 \vdash [*(x) := 0] * (y) \neq 0
\end{array}$$

then we would be misled to say the least. Rather, we need to make sure that the update is reflected in any subsequent memory read to the same address, regardless of the syntactic form of the index term. Perhaps something like the following:

$$[* (e) := \tilde{e}] p(*) \leftrightarrow p(*\{e \mapsto \tilde{e}\}) \quad (5)$$

Now when we repeat the derivation from before,

$$\begin{array}{c}
*(y) \neq 0, x = 1, y = 1 \vdash *\{x \mapsto 0\}(y) \neq 0 \\
\hline
*(y) \neq 0, x = 1, y = 1 \vdash [*(x) := 0] * (y) \neq 0
\end{array}$$

there is no way to close out the proof because  $x = y$  and  $*\{x \mapsto 0\}(y) = 0$ . But this proof rule isn't sound, because what if  $e$  evaluates to an out-of-bounds value? We need to add an assertion that the value of  $e$  is within the correct range. This leads to the  $[*]_{=}$  axiom, which combines Equation 5 with the in-bounds check.

$$([*]_{=}) \quad [* (e) := \tilde{e}] p(*) \leftrightarrow p(*\{e \mapsto \tilde{e}\}) \wedge 0 \leq e < U$$

Axiom  $[*]_ =$  takes care of what to do when we update memory, but we also need a way to reason about reads from memory. If we only ever reason about programs that never update memory, then this is easy because anything we need to know about its value at particular indices is already in our assumptions. We can then work with the index like we would any other value mentioned in a program.

But what about programs that update memory and then read from it afterwards? There are two cases to cover: reading from an index that was previously written to, and reading from one that was not. In the first case, we have some memory  $*\{e \mapsto \tilde{e}\}$  and we perform an access  $*\{e \mapsto \tilde{e}\}(e')$  where  $e = e'$  in the current state. Then the value that is read from memory will be  $\tilde{e}$ . But of course we also need to make sure that we are reading from an index in the appropriate range. This is captured in the  $[*]_1$  rule.

$$([*]_1) \quad \frac{\Gamma \vdash e = e' \quad \Gamma \vdash 0 \leq e' < U}{\Gamma \vdash *\{e \mapsto \tilde{e}\}(e') = \tilde{e}}$$

Note that because we must prove that  $e = e'$  in this case, for the right premise it is fine to alternatively prove that  $\Gamma \vdash 0 \leq e' < U$  for the right premise, as this will not affect soundness. In the case where  $e \neq e'$ , we use similar reasoning to conclude that  $*\{e \mapsto \tilde{e}\}(e')$  takes whatever the value at index  $e'$  in  $*$  was before the update, i.e.  $*(e')$ . This gives us the  $[*]_2$  rules.

$$([*]_2) \quad \frac{\Gamma \vdash e \neq e' \quad \Gamma \vdash 0 \leq e' < U}{\Gamma \vdash *\{e \mapsto \tilde{e}\}(e') = *(e')}$$

The axiom  $[*]_ =$  and rules  $[*]_1$ ,  $[*]_2$  are sufficient to prove safety properties about programs with pointer operations.

A syntactical sidenote. The axiom and rules that we introduced ask us to explicitly keep track of the updates that occur as we work through a dynamic logic formula using the  $*\{e \mapsto \tilde{e}\}$  notation. As a result, we could end up with a verification condition that lists out a long sequence of updates, e.g.,

$$a < U \wedge a \geq c \rightarrow *\{a \mapsto b\}\{c \mapsto d\}\{a - c \mapsto 5\}(e) \geq 100$$

While the notation may seem verbose, we need to know the locations and values of potentially all past assignments, and the order in which they occurred. So, this is simply a way of accounting for all of the information needed to reason about the consequences of memory updates.

However, this notation is only useful when we are doing such reasoning. We need not (and should not) write programs that mention updated memories explicitly. For example, one might be tempted to write a program like the following:

$$*\{y \mapsto 0\}(x) := 1$$

Or perhaps,

$$x := *\{z \mapsto 1\}(y)$$



Our language is quite a bit simpler than either C or Java, so the relevant memory safety issues are less numerous. Essentially, memory safety in our language means that any time a memory read or update is made, the term specifying the index evaluates to a value in  $[0, U]$ . In our simplified language with pointers, any “bad” use of memory immediately leads to an abort on the corresponding trace, so we can informally define memory safety as the set of traces that do not abort due to a pointer read or write. Definition 4 expands on this, and makes it less ambiguous.

**Definition 4 (Memory safety).** A program  $\alpha$  satisfies memory safety if and only if for all traces  $\sigma \in \llbracket \alpha \rrbracket$ , and any state  $\sigma_i$  in that trace wherein  $\alpha$  executes a memory operation:

- If the operation is a memory read  $\ast(e)$ , then  $e$  is in bounds, i.e.  $\sigma_i \models 0 \leq e < U$ .
- If the operation is an update  $\ast(e) := \tilde{e}$ , then  $e$  is in bounds,  $\sigma_i \models 0 \leq e < U$ .

□

One question that is often asked about the second bullet of Definition 4 is, why is there only a condition on  $e$ , and not also  $\tilde{e}$ ? This question is motivated by the fact that  $\tilde{e}$  may itself contain a memory read that could violate the safe bounds. But if this were the case, then the program would already fail to satisfy the first bullet, so we would not fail to notice that it is unsafe!

Another thing to notice is that when we use these axioms to prove any property about a program that uses pointers, we are forced to prove memory safety as well. The only case that we might forget to prove memory safety for is when a read is performed on memory without first having updated it. We can help ourselves remember to do this by replacing each command  $\alpha$  that reads from or writes to memory in term  $\ast(e)$  with the following composed command:

$$\text{assert}(0 \leq e < U); \alpha \tag{6}$$

This is a theorem that we are able to prove, in fact.

**Theorem 5.** For any formula  $P$  and program  $\alpha$  that has been rewritten according to (6), if  $\Gamma \vdash [\alpha]P$ , i.e.  $[\alpha]P$  is provable from assumptions  $\Gamma$  using  $[*]_=$ ,  $[*]_1$ , and  $[*]_2$  in addition to other axioms of dynamic logic, then  $\alpha$  satisfies memory safety.

**Proof.** The way to prove this is by induction on the structure of proofs, just as we did to prove the soundness of the propositional sequent calculus and. This is a good exercise to complete on your own. □

### 3 Sandboxing memory access

Memory safety is an important policy in that we would want any useful program to be memory safe. But there are other sorts of safety policies on memory that we might want to enforce more selectively on only certain programs. For example, consider the

following pseudocode that checks a configuration variable to determine whether or not to display and advertisement. If so, then the program runs  $\alpha$  to render an ad on the screen.

`if(display ads)  $\alpha$  else continue without ads`

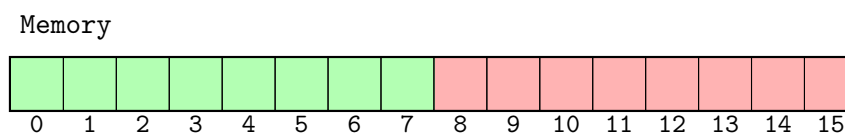
Suppose that  $\alpha$  was provided by the ad network, then we may have good reason not to trust  $\alpha$ . Perhaps the ad network hires dumb programmers, and we fully expect their  $\alpha$  to accidentally trample on memory that it isn't supposed to. Or maybe we got a great deal from Fancy Bear Ad Network, and despite assurances that his rendering code is "Totally 100% safe!", there are lingering doubts.

Luckily we know all about proving safety, so perhaps we can use logic and deduction to show that our program remains safe. What could  $\alpha$  do to ruin our day? One thing is that because it executes in the state space of our original code, it can change the values of any variables at will to whatever it likes. Perhaps that's not so bad, because our program only uses a limited number of variables and we have a lot of memory to work with. But  $\alpha$  could also change memory arbitrarily, and this is certainly a bad thing. It could also read the contents of memory and render them to the screen, or worse yet, send them back to the ad network. This is also a bad thing that we want to prevent.

### 3.1 Sandboxing safety

One solution is to create a virtual sandbox for  $\alpha$  to play in. We will give it free reign over a limited region of memory, and construct our program so that by the time  $\alpha$  runs, our correctness doesn't depend on the contents of that region. We will also let it do whatever it wants with the variables, isolating the rest of our program from the effects of these operations by first saving all of our variables to a part of memory outside  $\alpha$ 's sandbox region, and restoring them after  $\alpha$  finishes.

Supposing our machine only has a very limited 16-element memory, our segmentation would look something like the following with the parts shaded green comprising the safe memory set aside for our program, and the parts in red the sandbox for  $\alpha$ .



Now that we have decided on a safety policy with which to execute  $\alpha$ , we need to figure out how to actually enforce it in our program. Intuitively, our policy defines a "bad thing" that is any memory access outside of the sandbox region defined by upper and lower bounds  $s_l, s_h$ . So we might reasonably enforce the policy by first checking that the index  $i$  of any memory access operation in  $\alpha$  satisfies  $s_l \leq i \leq s_h$  before executing the operation. Luckily our language contains `assert( $Q$ )` commands, which come in handy when implementing such checks: if the check fails, the trace aborts rather than violating the policy.

So taking stock of our language, we propose to do the following instrumentation of  $\alpha$ .



- Replace each command of the form  $\ast(e) := \tilde{e}$  with a new composed command:

$$\text{assert}(s_l \leq e \leq s_h); \ast(e) := \tilde{e}$$

This will ensure that  $\alpha$  doesn't update any locations outside the sandbox.

- Replace any command  $\beta$  containing the term  $\ast(e)$  with the command:

$$\text{assert}(s_l \leq e \leq s_h); \beta$$

This will ensure that  $\alpha$  doesn't read any locations outside the sandbox.

This seems pretty convincing. Our language is fairly simple, so we're pretty sure that all our bases are covered in terms of sandboxing  $\alpha$ . The assertions themselves are a straightforward reflection of our sandboxing policy.

The downside to this type of enforcement is that any violation of the sandboxing policy, regardless of whether it is inadvertent or intentionally malicious, will cause our entire program to abort. This is less than ideal, as the malice or incompetence of  $\alpha$ 's developers still has a direct impact on the functioning of our code. Perhaps we can do better. In the next lecture, we will see how to isolate the rest of our programs from these effects using software fault isolation.