


# IY5512 Computer Security Part 4: Software security

Chris Mitchell

[me@chrismitchell.net](mailto:me@chrismitchell.net)

<http://www.chrismitchell.net>



Royal Holloway  
University of London

Information Security Group

## Objectives

- Understand why some languages and execution environments can lead to vulnerabilities.
- Identify other types of vulnerabilities.
- Learn how threats posed by malicious software (malware) can be reduced.
- Describe difference between worms and viruses, and how they use vulnerabilities.

2

The objectives of this part of the course are as follows:

- to understand why certain programming languages and execution environments may lead to vulnerabilities;
- to identify other types of vulnerabilities;
- to learn how the threats posed by malicious software (malware), such as viruses and worms, can be reduced;
- to appreciate the difference between worms and viruses, and how they exploit different types of vulnerabilities.




Information Security Group

# Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

3

We start this part of the course by reviewing how software can pose a threat to security.



Royal Holloway  
University of London

Information Security Group

## Introduction


- Following essential for a secure computer system:
  - secure hardware;
  - secure networks;
  - security-aware personnel; and
  - good security management.
- Just as important is **secure software**.

4

As discussed throughout this course, the following are all essential if a computer system is to be secure:

- secure hardware;
- secure networks;
- security-aware personnel; and
- good security management.

However, also extremely important is *secure software*.



Royal Holloway  
University of London

Information Security Group

## Background

- Significant security incidents regularly arise from vulnerabilities in software.
- Such vulnerabilities can arise from:
  - design errors; and
  - coding errors.
- These vulnerabilities ‘sit there’, waiting to be exploited by a range of types of malware.


5

Computer systems regularly encounter significant security incidents arising from vulnerabilities in software. Such vulnerabilities can arise from both:

- design errors; and
- coding (implementation) errors.

These vulnerabilities ‘sit there’, waiting to be exploited by a range of categories of malware, including worms and viruses.

This does not imply that the world is full of either malicious or incompetent programmers. However it does imply that if more attention was given to the design of secure code, then many security failures could be avoided.



Royal Holloway  
University of London

Information Security Group

## Background

- Security vulnerability monitoring agencies, or CERTs, e.g. CERT/CC and US-CERT, regularly issue security advisories.
- Vast majority relate to 'holes' in software that could have been prevented and that can be exploited by hackers or malicious coders.
- Exercise:
  - Visit the homepage of the SANS centre at <http://www.sans.org/newsletters/risk/> to see how many vulnerability advisories posted in the last month relate to software security issues.

6

Security vulnerability monitoring agencies (often called CERTs, for Computer Emergency Readiness Teams), such as the CERT Coordination Center (CERT/CC) and the US Computer Emergency Readiness Team (US-CERT), regularly issue security advisories. CERT/CC is part of CERT, based at Carnegie-Mellon University in the US. The CERT/CC home page is here:

<http://www.cert.org/certcc.html>

US-CERT is part of the US Department of Homeland Security. The US-CERT home page is here:


<http://www.us-cert.gov/>

The vast majority of these advisories relate to 'holes' in software that could have been prevented, and that can be exploited by hackers or malicious coders.

As an exercise visit the homepage of the SANS centre at:

<http://www.sans.org/newsletters/risk/>

to see how many vulnerability advisories posted in the last month relate to software security issues. [SANS is a commercial organisation].



Royal Holloway  
University of London

Information Security Group

## Why so many vulnerabilities? I


- **Greater use of networking:**
  - the Internet:
    - provides wide connectivity between devices running software;
    - exposes software to range of threats from increasing number of sources.
  - instead of single benign user with well-formed input, software exposed in uncontrolled network environment, with malicious users supplying inappropriate input.

7

There are a variety of reasons for the huge number of software vulnerabilities being found, reported and (we hope) fixed. One reason is the ubiquitous use of networking technology (and the Internet in particular). The success of the Internet:

- provides greater connectivity between devices that run software;
- means that software is more widely exposed to a range of threats from an increasing number of sources.

Instead of a single benign user providing the software with well-formed input, software is exposed in an uncontrolled network environment, with malicious users capable of supplying inappropriate input.

Royal Holloway  
University of London

Information Security Group


## Why so many vulnerabilities? II

- **Increasing size and complexity of systems** means that:
  - information systems harder to analyse for security;
  - errors more likely to go unnoticed;
  - comprehensive testing difficult, even impossible.

8

A second reason for vulnerabilities is the ever-increasing size and complexity of systems. This complexity means that:

- information systems are harder to analyse for security;
- errors are more likely to go unnoticed;
- comprehensive testing is difficult and, in practice, often impossible to conduct.



Royal Holloway  
University of London

Information Security Group


## Why so many vulnerabilities? III

- **Expect greater flexibility from systems:**
  - software designed for Internet is extensible, allowing additional features to be “plugged in”;
  - web browsers can install additional code, e.g. to view or listen to multimedia;
  - running of mobile code is permitted;
  - security features may not adequately protect the extended environment.

9

A third reason is that we expect greater flexibility from systems:

- software designed for the Internet is extensible, allowing for additional features to be “plugged in”;
- web browsers can install additional code, e.g. for viewing or listening to multimedia;
- the execution of mobile code (i.e. code downloaded by to an end PC by a server) is permitted;
- the provided security features, perhaps designed before the need for additional flexibility was envisaged, may not adequately protect the extended environment.



Royal Holloway  
University of London

Information Security Group


## Why so many vulnerabilities? IV

- **Lack of diversity** in desktop environment:
  - as with environmental systems, this can be dangerous (n.b. use of address space randomisation to reduce threat);
  - errors in Windows systems have significant impact on security, and not easily fixed;
  - attacks on Windows may have a greater payback (in both kudos and possible criminal impact) than attacks on other platforms.

10

A fourth reason is the lack of diversity in the computing environment – Windows is dominant on the desktop, and the only other significant OS is Unix (including MacOS):

- as with environmental systems, this lack of diversity can be dangerous (Windows 8 includes measures designed to increase the diversity of certain internal system properties, i.e. the address space, of a running system);
- errors in Windows systems have a significant impact on security and are not easily fixed;
- attacks on Windows may be seen as having a greater payback (in both kudos and possible criminal impact) than attacks on other platforms.

Royal Holloway  
University of London

Information Security Group


## Why so many vulnerabilities? V

- **Increasing market pressures:**
  - growing marketplace, with high level of competition among software vendors;
  - commercial incentives geared to rushed software production;
  - not conducive to good design and careful, thorough testing.

11

A fifth reason for the presence of vulnerabilities is the ever-increasing market pressures on vendors:

- it is a growing marketplace, with a high level of competition among software vendors;
- commercial incentives are geared towards rushed software production;
- the pressure to ship as early as possible is not conducive to good design and careful, thorough testing.



Royal Holloway  
University of London

Information Security Group


## Why so many vulnerabilities? VI

- **Shortcomings in software development processes:**
  - lack of awareness of causes of vulnerabilities in developer community;
  - lack of tools to support more secure software development;
  - competing pressures – e.g. inclusion of runtime checks for buffer overflows adversely affects software performance.

12

A sixth reason for the existence of so many vulnerabilities are widespread shortcomings in the software development processes:

- there is a lack of awareness of the causes of vulnerabilities amongst the wider software developer community (hopefully this situation is changing; for example, security is now becoming a 'standard part' of a Computer Science degree);
- historically there has been a lack of tools to support more secure software development (although this is being addressed, including by free tools such as those available on the SDL site – see: <http://www.microsoft.com/security/sdl/adopt/tools.aspx>);
- there are competing pressures on software vendors relating to software security – e.g. the inclusion of runtime checks for security issues such as buffer overflows has an adverse impact on software performance.



Royal Holloway  
University of London

Information Security Group

## Why do vulnerabilities matter?

- Historically, software vulnerability exploits were relatively small scale:
  - done primarily ‘for fun’ or for prestige.
- Today:
  - huge and growing criminal interest in exploiting vulnerabilities, leading to large scale criminal activities, involving spam, phishing, botnets, identity fraud, ...
  - increasing interest by state actors.

13

Historically, exploits of software vulnerabilities were relatively small scale, and were done primarily ‘for fun’ or for prestige reasons.

Today, there is a huge and growing criminal interest in exploiting vulnerabilities, giving rise to large scale criminal activities, e.g. involving spam, phishing, botnets, identity fraud, etc. At the same time, nation states are increasingly viewing malware and other types of attack on information infrastructures as ways of conducting espionage and even covert (or overt) war.



Royal Holloway  
University of London


Information Security Group

## Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

14

We next consider ways in which the existence of software vulnerabilities can be addressed during development.



Royal Holloway  
University of London

Information Security Group


## Penetrate and Patch I

- 'Penetrate and patch' is an approach to managing software vulnerabilities taken by some major software vendors.
- Involves releasing software into the commercial market that is poorly designed from a security perspective.
- Patches for downloading released when a security flaw is discovered.

15

'Penetrate and patch' is an approach to managing software vulnerabilities taken by some (maybe all?) major software vendors. This is typically not a deliberate decision, but one often made in ignorance of the seriousness and scope of vulnerabilities, and the ease with which they can be accidentally introduced into software.

It involves releasing software into the commercial market that is poorly designed from a security perspective. Patches for downloading by users are released when a security flaw is discovered.

Royal Holloway  
University of London

Information Security Group


## Penetrate and Patch II

- Not a sound approach to security.
- Should not be used routinely – only a last resort.
- If security is seriously considered at the design stage, then, in principle at least, there should not be a need for frequent release of patches.

16

Penetrate and patch is clearly not a sound approach to managing software security. It should not be used routinely, but only as a last resort.

If security is seriously considered at the design stage, then, in principle at least, there should not be a need for the frequent release of patches. However, practical experience suggests that removing **all** vulnerabilities from complex software, even when great care is taken, is very difficult to achieve.



Royal Holloway  
University of London

Information Security Group


## Penetrate and Patch III

- Problems with this approach include:
  - Vulnerabilities must first be found;
  - Hastily written patches may introduce new vulnerabilities;
  - No guarantee users will install patches;
  - Lazy approach that essentially depends on users doing the testing after release;
  - Damaging to vendor's reputation with users.

17

There are a number of significant problems with this approach:

- vulnerabilities must first be found (the worst case is when they are found first by malicious entities who exploit them before any patches/fixes are in place – these are the so called **zero day** vulnerabilities – Flame and Stuxnet exploited zero day vulnerabilities);
- hastily written patches to address newly discovered vulnerabilities may themselves introduce new vulnerabilities;
- there is no guarantee that users will install patches in a timely way;
- this is a lazy approach that essentially depends on users doing the testing after release;
- the need to ship large number of fixes is damaging to a vendor's reputation with users.

Royal Holloway  
University of London

Information Security Group

## Penetrate and Patch IV

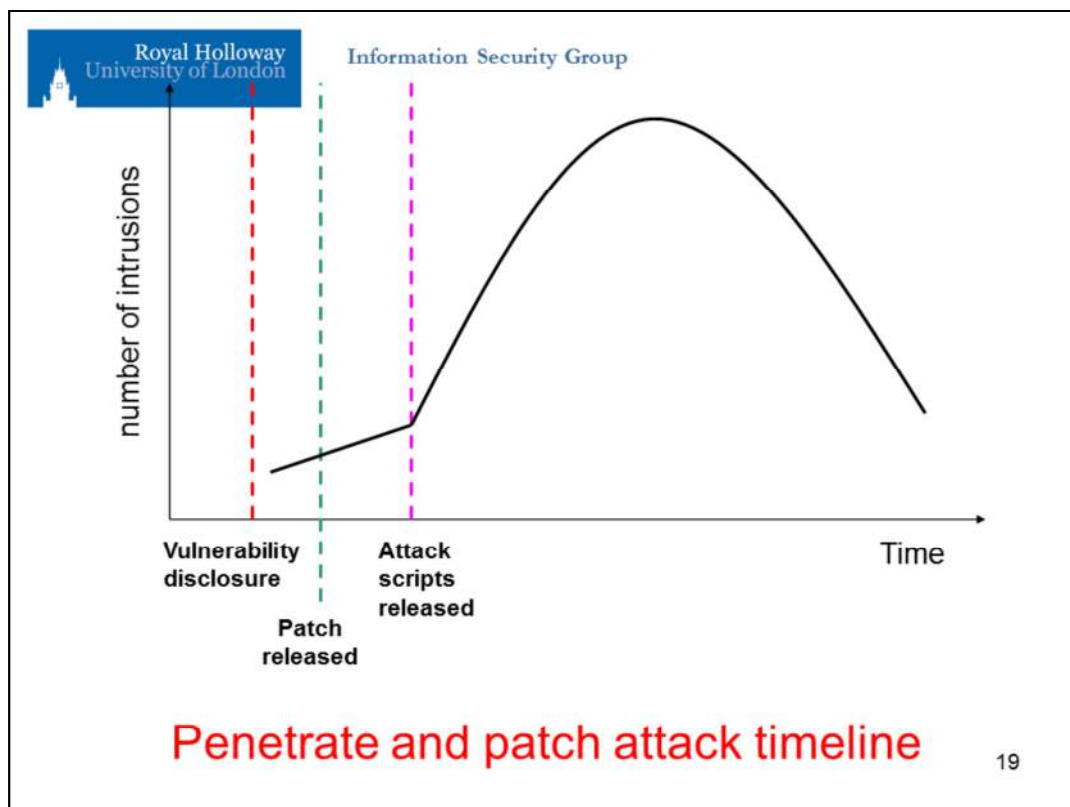
- Despite obvious problems (and despite improvements in robustness of software from some vendors in recent years):
  - all major commercial software appears to need routine patching;
  - this applies to both the open and closed source communities.

18


However, despite the obvious shortcomings of this approach (and despite improvements in the robustness of software from some vendors in recent years):

- all major commercial software appears to need routine patching;
- this applies to both the open and closed source communities.

This is despite the fact that organisations such as Microsoft have spent very large sums trying to improve their development processes (and retrain their development staff) to reduce the occurrence of vulnerabilities.



The much quoted figures behind this graph come from a paper published in *IEEE Computer* in 2000. The graph shown is abstracted from data in this paper, and shows a typical intrusion pattern over time. The important points in the timeline are disclosure of a vulnerability, release of a patch and release of an attack script that allows attackers with little technical background to exploit the vulnerability. The results clearly show the problems with penetrate and patch.



Royal Holloway  
University of London

Information Security Group


## Open source software

- Open source software is
  - developed in an collaborative manner, and open to scrutiny;
  - often not-for-profit and freely available.
- Examples of high profile initiatives include:
  - Linux operating system;
  - Gnu software suite.

20

Open source software is developed in an collaborative manner, and the software is open to scrutiny. It is often developed on a not-for-profit basis and is freely available.

Examples of high profile initiatives include the Linux operating system and the Gnu software suite.



Royal Holloway  
University of London

Information Security Group

## Closed source software

- Closed source software is not formally released for general scrutiny.
- Source code not available – only executable files containing machine code.
- May be possible to ‘reverse engineer’ back to source code.

21

Closed source software is not formally released for general scrutiny. The source code is not publicly available – all that is released are executable files containing machine code.

As with any software, it may be possible to ‘reverse engineer’ the executable back to source code.

Note that the dichotomy between open and closed source is not as clear as is often assumed. For example, in many cases, vendors of closed source systems will allow major customers access to source code.



Royal Holloway  
University of London

Information Security Group

## Open source – security implications

- ++ Open discussion of code details is possible.
- -- Code open to abuse by anyone seeking to exploit vulnerabilities.
  - -- May not be exposed to experts in security.
  - -- Development may be piecemeal.

22


The positive aspects of open source include:

- Open discussion of the code details is possible.

The negative aspects of open source include:

- Code is equally open to abuse by anyone seeking to exploit vulnerabilities;
- The code may not be exposed to experts in security during development;
- Development may be piecemeal.

It could be argued that the last issue is not inherent to open source. It is more a question of whether anyone is committed to maintaining the software. Typically, closed source software is developed and maintained by an organisation which has a vested interest in maintaining it (including fixing bugs when found). Open source software is typically not maintained by a commercial organisation (although there are exceptions), in which case maintenance depends on the willingness and availability of individuals prepared to devote their time to the task. As a result, the level of maintenance can be very patchy – for example, whilst the Linux kernel is very well maintained, other operating system components and applications may often be much less well-maintained.



Royal Holloway  
University of London

Information Security Group

## Closed source – security implications

- ++ Code details hidden makes it harder for an attacker to exploit any vulnerabilities.
- -- ‘Security by obscurity’ is a questionable security principle.
  - -- No guarantee of good design.
  - -- It may be possible to reverse engineer the code.


23

The positive aspects of closed source include:

- The fact that the code details are hidden makes it harder for an attacker to exploit any vulnerabilities.

The negative aspects of closed source include:

- Security by obscurity is a questionable security principle (although almost all security relies on secrets of some kind).
- There is no guarantee of good design.
- It may anyway be possible to reverse engineer the code.



Information Security Group

## Principles for creation of secure software I

- **Secure software from the outset:**
  - penetrate and patch model implies security was not a high priority in design;
  - best way to secure software is to design it from the outset;
  - security considerations should be integral to the software engineering lifecycle, not an afterthought.

24


### Secure software **from the outset**.

To comply with the software security principle 'secure software from the outset' it is necessary to commence any software development process with a thorough analysis of the security requirements of the intended application.

Adding features to an existing design in order to correct vulnerabilities does not address the underlying problem; either the security requirements were incorrectly specified or the implementation did not meet the security requirements.

Capturing security requirements is essentially a risk assessment exercise.

An example of a methodology designed to try to achieve this goal is the Security Development Lifecycle (SDL) model.



Royal Holloway  
University of London

Information Security Group

## Principles for creation of secure software II


- **Secure Software by least exposure:**
  - ensure software is not exposed to more risks than necessary;
  - isolate code for security critical operations in separate modules or libraries, hence easier to analyse and control;
  - enforce principle of least privilege: grant minimum access necessary to perform an operation, and only for the time needed to complete it;
  - minimise the possible 'attack surface': turn off all unnecessary functionality and services.

25

Secure Software **by least exposure** means the following:

- ensure that software is not exposed to more risks than necessary;
- isolate code for security critical operations in separate modules or libraries, making it easier to analyse and control;
- enforce the principle of least privilege: grant the minimum access necessary to perform an operation, and only for the time needed to complete it;
- minimise the possible 'attack surface' by turning off all unnecessary functionality and services.

This links closely to the Saltzer-Schroeder security principle of **Least privilege**.



Royal Holloway  
University of London

Information Security Group

## Principles for creation of secure software III

- **Secure software by default:**
  - i.e. software should fail securely: failure (e.g. due to inappropriate input or lost connection) should not create security vulnerabilities, e.g.:
    - software that generates a 'world readable' core dump after a crash may compromise security critical information, e.g. password entries;
    - firewall that keeps working when log disk is full could mean a security critical event is not detected.


26

Secure software **by default** means that software should fail securely; that is a failure of the software (e.g. due to inappropriate input or a lost connection) should not expose the system to security vulnerabilities.

For example:

- software that generates a 'world readable' core dump after a crash is problematic, because the core dump may contain security critical information such as password entries;
- a firewall that continues to operate when the disk on which it stores its log becomes full could lead to a failure to detect a security critical event.

This links closely to the Saltzer-Schroeder security principle of **Fail-safe defaults**.



Royal Holloway  
University of London

Information Security Group

## Principles for creation of secure software IV

- **Secure software by simplicity:**
  - from a design perspective:
    - simple software easier to analyse and test;
    - more likely that vulnerabilities are discovered and fixed;
  - from a usability perspective:
    - security features must appear simple, or users will either ignore or avoid them;
    - users will not use software that has been over-designed from a security perspective.

27

Unfortunately, complex requirements rarely lead to simple software!

Secure software **by simplicity**.

From a design perspective:

- simple software is easier to analyse and test;
- it is more likely that vulnerabilities will be discovered and corrected.

From a usability perspective:

- security features must appear simple, or users are likely to either ignore or avoid them;
- users will be reluctant to use software that has been over-designed from a security perspective.

This links closely to the Saltzer-Schroeder security principles of **Economy of mechanism** and **Ease of use**.



Royal Holloway  
University of London


Information Security Group

## Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

28

We next consider how vulnerabilities can arise in software.



Royal Holloway  
University of London

Information Security Group

## Overview


- There are a wide range of types of software vulnerability.
- We look at some important types.
- Often these arise because programmers are simply unaware of need to write code designed to resist hostile use.
- Focus is naturally on the intended use.
- **Need better education!**

29

There are a wide range of types of software vulnerability. We look at some particularly significant categories of such vulnerabilities.

These vulnerabilities often arise because programmers are simply unaware of the need to write code designed to resist hostile or inappropriate use. That is, the focus of the programmer tends to be on the intended use of the software, not on how it might be abused.

This is to a large extent a failure in education and training. It is all too easy when writing code to think about the intended use, and to ignore the possibility of unintended use (e.g. inputs which are not in the expected form).



Royal Holloway  
University of London

Information Security Group

## Topics to cover


- Start with simple motivating examples.
- Then discuss:
  - I. data type issues;
  - II. integer overflows;
  - III. program flow and functions;
  - IV. buffer overflows;
  - V. use after free issues.
- Conclude with discussion of mitigations.

30

After a discussion of some simple (motivating) examples, we divide the discussion of vulnerabilities under the following sub-headings:

- I. data type issues;
- II. integer overflows;
- III. program flow and functions;
- IV. buffer overflows;
- V. use after free issues.

We conclude this discussion with a brief consideration of possible mitigations, i.e. ways of addressing the threat posed by these vulnerabilities.



Royal Holloway  
University of London

Information Security Group

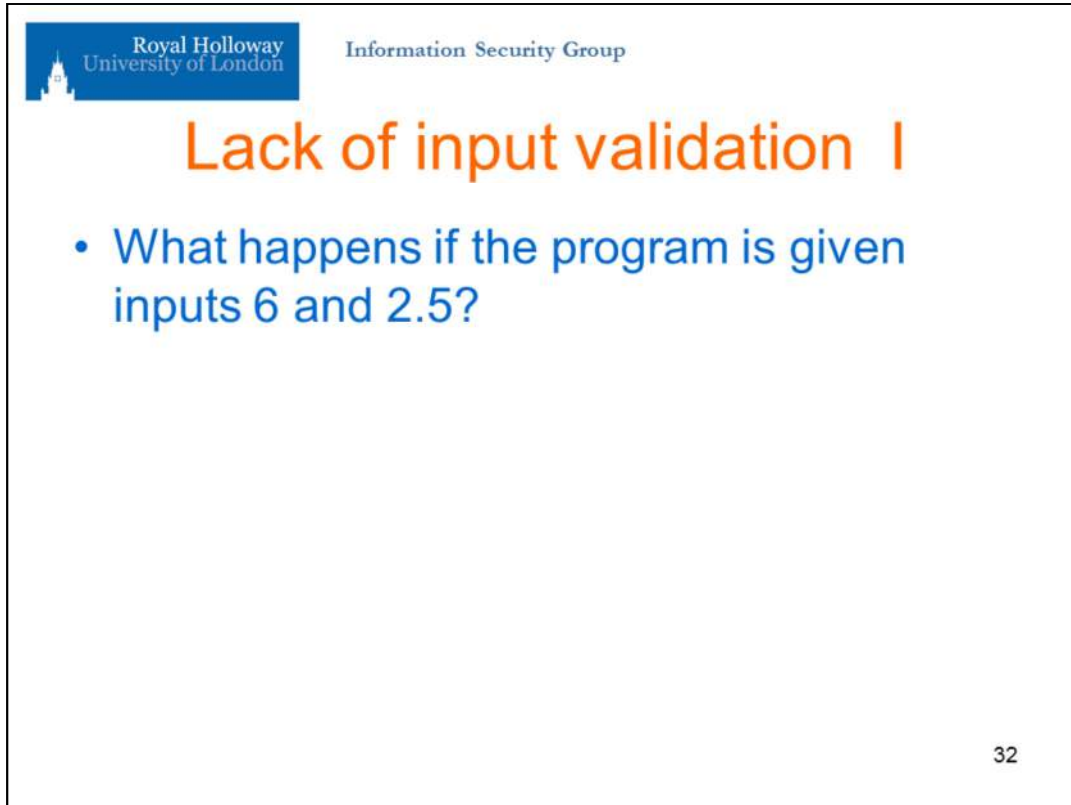
## Vulnerabilities – simple examples

- Recall simple program to multiply two numbers:
  1. Draw box labelled *x* and write first input in it;
  2. Draw box labelled *y* and write second input in it;
  3. Draw box labelled *result* and write 0 in it;
  4. If the contents of *y* is 0 go to step 8;
  5. Add contents of box *x* to contents of box *result*;
  6. Subtract one from contents of box *y*;
  7. Go to step 4;
  8. Finish.

31

Recall the simple program to multiply two numbers, introduced earlier in the course:

1. Draw a box labelled *x* and write the first input in it;
2. Draw a box labelled *y* and write the second input in it;
3. Draw a box labelled *result* and write 0 in it;
4. If the contents of *y* is 0 go to step 8;
5. Add the contents of box *x* to the contents of box *result*;
6. Subtract one from the contents of box *y*;
7. Go to step 4;
8. Finish.

A slide with a white background and a black border. At the top left is the Royal Holloway University of London logo, which consists of a blue square with a white silhouette of a building and the text 'Royal Holloway University of London' in white. To the right of the logo is the text 'Information Security Group' in a dark blue font. Below the logo and group name is the title 'Lack of input validation I' in a large, bold, orange font. Underneath the title is a single bullet point in blue text: '• What happens if the program is given inputs 6 and 2.5?'. In the bottom right corner of the slide, the number '32' is written in a small, dark blue font.

Royal Holloway  
University of London

Information Security Group


## Lack of input validation I

- What happens if the program is given inputs 6 and 2.5?

32

What happens if the program is given inputs 6 and 2.5?

[This illustrates the importance of checking that inputs are as expected.]



Royal Holloway  
University of London

Information Security Group

## Lack of input validation II

- What happens if the program is given inputs 6 and 2.5?
- Well ...
  - the termination condition ( $y = 0$ ) is never satisfied;
  - the program enters an infinite loop.


33

What happens if the program is given inputs 6 and 2.5?

As a result:

- the termination condition ( $y = 0$ ) is never satisfied;
- the program enters an infinite loop.

That is, the program was written under the assumption that the inputs would be of the expected type. If this is not true, then the program can fail in unexpected and potentially dangerous ways.



Royal Holloway  
University of London

Information Security Group


## Overflow I

- Suppose the child (executing the program) can only count up to 1000:
  - what happens if child is given inputs 61 & 32?

34

Suppose the child (executing the program) can only count up to 1000:

- what happens if the child is given input 61 and 32?



Royal Holloway  
University of London

Information Security Group

## Overflow II

- Suppose the child (executing the program) can only count up to 1000:
  - what happens if child is given inputs 61 & 32?
- At some point the value in *result* will be greater than 1000;
  - $61 \times 17 = 1037$
  - what will the child do?
  - maybe the child will start again from 37, maybe the child will simply stop, or ...

35

Suppose the child (executing the program) can only count up to 1000:

- what happens if the child is given input 61 and 32?


During execution, at some point the value in *result* will be greater than 1000;

- $61 \times 17 = 1037$

- what will the child do?

- maybe the child will start again from 37, maybe the child will simply stop (generate an exception), ...

In any event, the result will be not what is expected.



Royal Holloway  
University of London

Information Security Group


## Changing execution flow I

- Suppose the child has written the instructions down:
  - what happens if someone changes the instruction in line 7 from 'Go to step 4' to 'Go to step 8'?

36

Suppose the child has written the instructions down:

- what happens if someone changes the instruction in line 7 from 'Go to step 4' to 'Go to step 8'?



Royal Holloway  
University of London

Information Security Group

## Changing execution flow II

- Suppose the child has written the instructions down:
  - what happens if someone changes the instruction in line 7 from 'Go to step 4' to 'Go to step 8'?
- The program no longer computes the correct result:
  - by modifying a 'pointer' (Go to step 4) the program logic has been changed.

37


Suppose the child has written the instructions down:

- what happens if someone changes the instruction in line 7 from 'Go to step 4' to 'Go to step 8'?

The program no longer computes the correct result:

- by modifying a 'pointer' (Go to step 4) the program logic has been changed.

This underlines the importance of ensuring that the integrity of the code of the program must be guaranteed.



Royal Holloway  
University of London

Information Security Group

## Practice

- Programs suffer from all these problems:
  - in simple examples, effects of vulnerabilities are unimportant from a security perspective;
  - however, exploiting a vulnerability can have significant consequences for security.
- We next consider:
  - why these vulnerabilities occur;
  - how to detect their presence;
  - how impact can be reduced.


38

Computer programs suffer from all these types of vulnerabilities:

- in these simple examples, the effects of 'exploiting' these vulnerabilities are not important from a security perspective;
- however, in many situations exploiting a vulnerability can have significant consequences for security.

We next consider:

- why these vulnerabilities occur;
- how to detect the existence of such vulnerabilities;
- how the impact of these vulnerabilities can be reduced (i.e. we look at mitigations).



Royal Holloway  
University of London

Information Security Group


## Vulnerabilities I – data type issues

- Simple program was written on assumption that the inputs would be integers (whole numbers):
  - program behaves in unexpected ways if the inputs are not integers;
  - e.g. if  $y$  is not a whole number then program never terminates.

39

The simple program was written on assumption that the inputs would be integers (whole numbers). As a result, the program behaves in unexpected ways if the inputs are not integers.

- for example, if  $y$  is not a whole number then the program never terminates.



Royal Holloway  
University of London

Information Security Group


## Data types I

- Good idea to specify what values of a variable (such as  $y$ ) are legal:
  - compiler or run-time environment can allocate correct amount of storage if it knows what type of input to expect;
  - compiler may be able to detect incorrect or inappropriate use of inputs and variables.

40

It is generally a good idea to specify what values of a variable (such as  $y$ ) are legal:

- the compiler or run-time environment can allocate an appropriate amount of storage if it knows what type of input to expect;
- the compiler may be able to detect incorrect or inappropriate use of inputs and variables.



Royal Holloway  
University of London

Information Security Group


## Data types II

- Computer programs use variables and constants.
- In the example:
  - *x*, *y* and *result* are all variables;
  - compiler allocates storage for variables in data section of executable program's address space.

41

Computer programs define variables and constants. In the simple example:

- *x*, *y* and *result* are all variables;
- the compiler allocates storage for variables in the data section of the executable program's address space.



Royal Holloway  
University of London

Information Security Group

## Data types III

- **type** of a variable determines:
  - what values variable is allowed to take;
  - how much storage needed for variable;
  - how data in memory should be interpreted.
- 0x41 may represent the integer 65 or the character 'A', depending on variable type.


42

The **type** of a variable determines:

- what values the variable is allowed to take;
- how much storage is required for the variable;
- how a sequence of bits in memory should be interpreted.

For example, the sequence of bits 0x41 may represent the integer 65 (0x41 is 65 in decimal) or the character 'A' (0x41 is ASCII for 'A') depending on the type of the variable in which the value is stored.

ASCII (for American Standard Code for Information Interchange) is a standard way of encoding characters as strings of bits. ASCII is very widely used. The international standard for this encoding is ISO/IEC 646.



Royal Holloway  
University of London

Information Security Group

## Types – examples

- Examples of types in C include `int`, `char`, `float` and `double`:
  - variable of type `int` uses 4 bytes (32 bits); has values in range  $-2^{31}$  to  $+2^{31}-1$  (i.e. from: -2 147 483 648 to +2 147 483 647);
  - variable of type `int` is **declared** (in a program) using statement of form:
 


```
int x
```
  - `x` assigned a value using statement of form:
 

```
x = 2
```

43

Examples of variable types (or **data types**) in the C programming language include `int`, `char`, `float` and `double`:

- a variable of type `int` requires 4 bytes (32 bits) of storage, and takes values in the range  $-2^{31}$  to  $2^{31} - 1$  (i.e. from: -2 147 483 648 to +2 147 483 647);
- a variable of type `int` is **declared** (in a program) using a statement of the form: `int x`
- in the body of a program, `x` is assigned a value using a statement of the form: `x = 2`



Royal Holloway  
University of London

Information Security Group


## Arrays

- In C an **array** of 10 variables of type `int` can be declared by writing `int x[10]`.
- Individual elements of array can be referenced as `x[0]` (for the first variable), `x[1]`, etc.

44

In the C language, an array of 10 variables of type `int` can be declared by writing `int x[10]`.

Individual elements of the array can be referenced as `x[0]` (for the first variable), `x[1]`, etc., up to `x[9]`.



Royal Holloway  
University of London

Information Security Group


## Expressions

- **Expression** is a statement in a program that combines values:
  - each of these values can have a type;
  - in C, the expression `y = x + 1` takes the value of variable `x`, adds 1 to it, and puts the result in variable `y`;
  - the expression `y = power(x, z) + w` computes  $x^z + w$ , and puts the result in variable `y`.

45

An **expression** is a statement in a program that specifies a way of combining values:

- each of these values can have a type;
- in C, the expression `y = x + 1` takes the value of variable `x`, adds 1 to it, and puts the result in variable `y`;
- the expression `y = power(x, z) + w` computes  $x^z + w$ , and puts the result in variable `y`.



Royal Holloway  
University of London

Information Security Group

## Function signatures I

- Each input to a function has a type.
- Each output from a function has a type.
- List of input and output types for a function is its **type** or **signature**.
- Signature defines way in which the function must be called, and way in which its output can be used.


46

Each input to a function has a type, and each output from a function has a type.

The list of input and output types for a function is its **type** or **signature**.

The signature defines the way in which the function must be called, and the way in which its output can be used.

If you think of a function as a 'black box', the function signature defines the nature of what the black box takes in and what it gives out.



Royal Holloway  
University of London

Information Security Group

## Function signatures II

- In C, might write:

```
int power(int a, int n)
    { ... }
```


which defines the signature of the function `power`, which:

- takes two integers as **input**;
- **returns** an integer.

47

In C, we might write: `int power(int a, int n) { ... }` which defines the signature of the function `power`. This function:

- takes two integers as **input**;
- **returns** an integer.



Royal Holloway  
University of London

Information Security Group

## Type errors I

- Type error occurs if use of values in a program is incompatible with their types.
- For example, the instructions:

```
int x;  
...  
x=3.4;
```

cause type error by assigning 3.4 (floating point value) to variable (x) of type `int`.


48

A type error occurs if the use of values in a program is incompatible with the types of those values.

For example, the sequence of instructions:

```
int x;  
...  
x=3.4;
```

cause a type error by assigning a floating point value (i.e. 3.4) to a variable (i.e. x) of type `int`.



Royal Holloway  
University of London

Information Security Group

## Type errors – more examples

- The instructions `int x,y; ...  
y=x+'hello world';` cause a type error by adding a string of characters to an integer;
- The function call `power(3.4, 2)` causes a type error;
- The expression `y=power(x, z)+w` may cause a type error, depending on the types of the variables `x, y, z, w`.

49


The sequence of instructions:

```
int x,y;  
...;  
y=x+'hello world'
```

cause a type error by trying to add a string of characters to an integer;

The function call `power(3.4, 2)` causes a type error;

The expression `y=power(x, z)+w` may cause a type error, depending on the types of the variables `x, y, z, w`.



Royal Holloway  
University of London

Information Security Group

## Array bounds errors

- Suppose we declare: `int x[10]`:
  - if write `x[i]`, and `i` doesn't have appropriate value (an integer between 0 and 9), have a type error:
    - `x[3.4]`, `x['hello']` and `x[12]` all give type errors;
  - this is an **array bounds error**;
  - array bounds error typically causes a run-time exception (e.g. a 'segmentation fault').


50

Suppose we declare an array as follows:

```
int x[10]
```

Then:

- if we write `x[i]` and `i` does not have an appropriate value (an integer between 0 and 9 inclusive), then we have a type error. For example:
  - `x[3.4]`, `x['hello']` and `x[12]` will all cause type errors;
- a type error of the form `x[12]` is an array bounds error;
- an array bounds error typically causes a run-time exception (e.g. a 'segmentation fault').



Royal Holloway  
University of London

Information Security Group

## Detecting type errors I


- **Type checking** can be used to detect type errors:
  - can occur at compile-time or run-time;
  - language is said to be **statically typed** if type-checking is performed at compile-time;
  - C is a statically typed language.

51

A process known as **Type checking** can be used to detect type errors:

- it can occur at the time the program is compiled or when it is run;
- a language is said to be **statically typed** if type-checking is performed at compile-time;
- C is a statically typed language.

Type checking involves verifying that variables are used appropriately, i.e. in line with their declared types.



Royal Holloway  
University of London

Information Security Group

## Detecting type errors II

- Easy to detect some type errors at compile time:
  - `int x[10]; ...; x[12]=3;`
 is a type error that a compiler can detect.
- However, following 'abuse' of types:
  - `char x; int y; y=x+1;`
 not detected by a C compiler (mixing types is allowed in C); other languages have stricter rules on use of types.
- However, some type errors can be difficult to detect, no matter what compiler is used.

52

Some kinds of type errors can readily be detected at the time of compilation of a program. For example

```
int x[10]; x[12]=3;
```


is a type error that can be detected by the compiler.

However, the following code:

```
char x; int y; y=x+1;
```

will not be detected by a standard c compiler (mixing types in this way is permitted in c), but this type of assignment would be detected by compilers for other languages which impose stricter conditions on use of types. Of course such abuse of types is potentially unsafe, since, amongst other things, the result will depend on how the compiler converts characters to bit strings.

However, some type errors can be difficult to detect, no matter what compiler is used.



Royal Holloway  
University of London

Information Security Group

## Detecting type errors III

- Some type errors cannot be detected at compile time:
  - if variable used to store user input, its value cannot be known at compile time;
  - following code won't cause a compile-time error:

```
int x[10]; int i; ...; y=x[i]+1;
```
  - however `y=x[i]+1` may generate a run-time (array bounds) error if `i` has an inappropriate value when the program is run.


53

Some type errors cannot be detected at compile time. If a variable is used to store user input, its value cannot be known at compile time.

For example, the following code will not generate a compile-time error:

```
int x[10];
int i;
...;
y=x[i]+1;
```

However the expression `y=x[i]+1` may generate a run-time (array bounds) error if `i` has an inappropriate value when the program is run. Whether or not a run-time error is generated by this type of array bounds problem will depend on what checking is performed whilst a program is executing.



Royal Holloway  
University of London

Information Security Group

## Type safety I

- Type specifies that the data held by a variable must have certain properties:
  - e.g. Count will always hold an integer;
  - e.g. Array A will never store more than 10 items.
- Type checking verifies these properties.
- A language is said to be **type safe** if the properties are guaranteed to hold at run-time; otherwise the language is **unsafe**.


54

As we have discussed, the type of a variable specifies that the data held by that variable must have certain properties:

- e.g. Count will always hold an integer;
- e.g. Array A will never store more than 10 items.

Type checking involves verifying these properties, either at compile-time or at run-time.

A language is said to be **type safe** if the type properties are guaranteed to hold at run-time; otherwise the language is said to be **unsafe**.



Royal Holloway  
University of London

Information Security Group

## Type safety II

- Consequences of type safety for security:
  - safe language can guarantee memory safety;
  - static checking (at compile time) can make some buffer overflows impossible, but cannot help with mobile code from other sources;
  - dynamic checking (at load time) gives some basic checking of mobile code.
- C, C++ certainly not type safe; Java and C# are not completely type safe, but programs in these languages are much less likely to contain type errors.


55

Consequences for security of type safety:

- a safe language can guarantee memory safety;
- static checking (at compile time) can make some buffer overflows impossible, but cannot help with mobile code from other sources;
- dynamic checking (at load time) gives some basic checking of mobile code.

C and C++ are certainly not type safe; Java and C# are also not completely type safe, but programs written in these languages are very much less likely to contain type errors.

C# can be made type safe by prohibiting untyped pointers. Use of such pointers can be prohibited at the compiler level.

Information Security Group


## Data validation

- **Data validation** refers to checking of input data.
- Principle for programmers – always check for valid data and reject everything else.
- Checking for invalid data requires coder to anticipate all possibilities.
- Very easy to get wrong.

56

**Data validation** refers to the checking of input data. Always checking for valid data and rejecting everything else is a fundamentally important principle for programmers.

Checking for invalid data requires the coder to anticipate all possibilities. This is very easy to get wrong.



Royal Holloway  
University of London

Information Security Group

## Memory overflow errors I


- String (of characters) often treated as an array of characters.
- E.g. in C: `char password[20]`
  - defines a **string variable** called `password`;
  - memory allocated for this variable is 20 bytes (assuming each character is allocated one byte).

57

A string (of characters) is often treated as an array of characters. For example, in C we might write:

```
char password[20]
```

- this defines a **string variable** called `password`;
- the memory allocated for this variable is 20 bytes (assuming each character is allocated a single byte).

Information Security Group


## Memory overflow errors II

- If string `s` of more than 20 characters is assigned to `password`, additional characters “overflow” into adjacent memory locations:
  - safe approach is to copy only first 19 characters of `s` (and include a string termination character as 20th character).
  - programmers rarely adopt safe approach because there are pre-defined (but unsafe) functions for copying a string.

58

If a string `s` of more than 20 characters is assigned to `password`, the additional characters may “overflow” into adjacent memory locations:

- the safe approach is to copy only the first 19 characters of `s` (and include a string termination character as the 20th character).
- programmers rarely adopt the safe approach because there are pre-defined (but unsafe) functions for copying a string.



Royal Holloway  
University of London

Information Security Group


## Memory overflow errors III

- Many overflow errors in C not caught at compile-time or run-time:
  - significant cause of vulnerabilities;
  - if overflow error occurs, and overflow overwrites a memory location containing a pointer to executable code, it is possible to change the flow of execution;
  - more on this later ...

59

Many overflow errors in C are not caught at compile-time or run-time. They are a significant cause of vulnerabilities.

If an overflow error occurs and the overflow overwrites a memory location containing a pointer to executable code, it is possible to change the flow of execution. This is discussed in greater detail later ...



Royal Holloway  
University of London

Information Security Group

## Vulnerabilities II – integer overflows


- Integer overflows arise because of way in which whole numbers are stored.
- An  $n$ -bit computer **word** can be used to represent a non-negative (unsigned) integer between 0 and  $2^n-1$ .
- Typically  $n$  is a multiple of 8.
- However, often need to represent both positive and negative numbers.

60

Integer overflows arise because of the way in which whole numbers are stored in a computer.

An  $n$ -bit computer **word** can be used to represent a non-negative (unsigned) integer between 0 and  $2^n-1$ , where typically  $n$  is a multiple of 8.

However, we often need to represent both positive and negative numbers.



Royal Holloway  
University of London

Information Security Group

## Two's complement representation


- Signed integers typically stored using **two's complement representation**:
  - two's complement representation of  $-m$  is obtained by flipping each bit in binary representation of  $m$ , and adding 1;
  - in 8-bit word, 3 represented as 0000 0011 and -3 represented as 1111 1100+1 = 1111 1101;
  - in an  $n$ -bit word, can represent numbers between  $-2^{n-1}$  and  $+2^{n-1}-1$ .

61

Signed integers are typically stored using the **two's complement** representation:

- the two's complement representation of  $-m$  is obtained by flipping each bit in the binary representation of  $m$  (i.e. changing every one to a zero and every zero to a one) and adding 1;
- in an 8-bit word, 3 is represented as 0000 0011 and -3 is represented as 1111 1100+1 = 1111 1101;
- in an  $n$ -bit word, using this representation it is possible to represent numbers between  $-2^{n-1}$  and  $+2^{n-1}-1$ .

Note that, when using a two's complement representation, the most significant bit of a word is 1 for negative numbers and 0 for positive numbers (including zero).



Royal Holloway  
University of London

Information Security Group

## Ambiguity


- 0xFFFFB represents the (unsigned) integer 65 531.
- So one sequence of bits, i.e. 0xFFFFB, represents both a small negative integer (-5) & a large unsigned integer (+65 531).
- Thus, **type** of integer variable (i.e. unsigned or signed) is of great importance.

62

0xFFFFB represents the (unsigned) integer 65 531.

Hence the same sequence of bits, i.e. 0xFFFFB, represents both a small negative integer (-5) and a large unsigned integer (+65 531).

In other words, the **type** of an integer variable (i.e. whether it is unsigned or signed) is of great importance.



Royal Holloway  
University of London

Information Security Group


## Examples I

- Vulnerability may arise when a signed integer is interpreted as an unsigned integer.
- Observe that:
  - a small negative integer in two's complement is a large positive integer;
  - 0xFFFF is -1 as a signed 16-bit integer, but 65 535 as an unsigned integer.

63

A vulnerability may arise when a signed integer is implicitly converted to an unsigned integer. We observe that:

- a small negative integer in two's complement is a large positive integer;
- 0xFFFF is -1 as a signed 16-bit integer, but 65 535 as an unsigned integer.



Royal Holloway  
University of London

Information Security Group

## Examples II


- A vulnerability can arise when:
  - an (unsigned) integer variable is used to define the amount of data to be copied;
  - the variable is interpreted as a signed integer.
- What program will do when asked to transfer a negative amount of data is undefined.

64

A vulnerability can arise when:

- an (unsigned) integer variable is used to define the amount of data to be copied;
- the variable is interpreted as a signed integer.

How the program will work when asked to transfer a negative amount of data is undefined.



Royal Holloway  
University of London

Information Security Group

## Examples III


- Integer overflow vulnerabilities found in:
  - Sun Microsystems XDR library;
  - Windows DirectX library quartz.dll;
  - bash (open source Unix shell).
- Vulnerabilities often arise because of flawed implementation of input validation.

65

Integer overflow vulnerabilities are widespread. They have been found in:

- the Sun Microsystems XDR library;
- the Windows DirectX library quartz.dll;
- bash (open source Unix shell).

These vulnerabilities often arise because of an incorrect implementation of input validation.



Royal Holloway  
University of London

Information Security Group

## Vulnerabilities III – program flow

- Two main ways to control program flow:
  - **selection:**
    - choose execution path using a logical test, e.g.:
      - if ... then ... else ...
  - **iteration:**
    - repeat sequence of instructions while logical test is true, e.g.:
      - while ... do ...
      - for loops
      - repeat ... until ...

66

There are two main ways of controlling the flow (the execution path) of a program:

- **selection:**

- choose between two different execution paths based on the result of a logical test, e.g.:

- `if ... then ... else ...`


- **iteration:**

- repeat a sequence of instructions all the while a logical test evaluates to true, e.g.:

- `while ... do ...`

- `for loops`

- `repeat ... until ...`



Royal Holloway  
University of London

Information Security Group

## Functions

- Function call also causes change in execution flow:
  - execution jumps to address of called function.
- Function can be called anywhere in a program:
  - how does processor know which instruction to execute after called function has completed?
  - need method of keeping track of where we were when control handed to called function.

67


A function call also causes a change in the flow of execution:

- execution jumps to the first address of the called function.

A function can be called at any point in a program:

- how does the processor know which instruction to execute after the called function has completed?
- we need a method for keeping track of where we were in the calling function when it handed control to the called function.

The standard solution to this 'housekeeping problem' is the use of a data structure called a **stack**.



Royal Holloway  
University of London

Information Security Group

## Stacks


- **Stack** is a **last-in-first-out (LIFO)** data structure:
  - can be visualised as a pile of items.
- There are two stack operations:
  - **push** adds something to the top of the stack;
  - **pop** removes the top item on the stack.

68

A **stack** is a **last-in-first-out (LIFO)** data structure, i.e. a way of organising data in a computer. It operates analogously to a pile of items.

There are two fundamental operations that can be performed on a stack, namely:

- **push** adds something to the top of the stack;
- **pop** removes the top item on the stack.



Royal Holloway  
University of London

Information Security Group

## Stacks and function calls


- Stack is ideal data structure for keeping track of function calls:
  - every process has a stack;
  - when function *A* calls function *B*, current instruction pointer is pushed onto the stack;
  - this is the **return address** in function *A*;
  - when function *B* finishes, return address for *A* is popped off the stack and copied back to the instruction pointer.

69

A stack is an ideal data structure for keeping track of function calls. As a result every process address space includes a stack region;

- when a function *A* calls a function *B*, the contents of the instruction pointer (the address of a location in computer memory) is pushed onto the stack;
- this address is the **return address** in function *A*;
- when function *B* finishes the return address for *A* is popped off the stack and copied back to the instruction pointer.

The 'last in first out' property of the stack allows for the simple management of an indefinite number of function calls.



Royal Holloway  
University of London

Information Security Group

## Stack vulnerabilities


- If entries in the stack can be changed, then return addresses will potentially be modified (**arc injection**).
- Some vulnerabilities allow an attacker to make chosen modifications to the stack.
- Means that execution (the control flow) of program can be modified, e.g. to execute a program of the attacker's choice.

70

If entries in the stack can be changed, then return addresses will potentially be modified. This is what is known as an **arc injection**.

Some vulnerabilities allow an attacker to make chosen modifications to the stack.

Changing the stack would mean that the execution (the control flow) of the program would be modified, e.g. causing the execution of a program of the attacker's choice.



Royal Holloway  
University of London

Information Security Group

## Vulnerabilities IV – buffer overflows

- Buffers are sections of memory allocated for temporary data storage, e.g.:
  - when transferring data between a peripheral device and main memory;
  - for inter-process communication – data passed from one program to another.
- Storing a value exceeding buffer size may mean other memory locations are overwritten.
- Buffer overflow occurs when a program writes data beyond the bounds of the allocated buffer.


71

The allocation of memory resources is probably the most common area in which software security vulnerabilities arise. The majority of these problems arise from **buffer overflows**.

Buffers are sections of memory allocated for data storage, typically on a temporary basis while data is transferred from one location to another, e.g.:

- when transferring data from a peripheral device to main memory and vice versa;
- for inter-process communication, allowing parameters to pass from one program to another.

Storage of a value exceeding the size of the buffer may result in other memory locations being overwritten. A buffer overflow (or buffer over-run) occurs when a program writes data beyond the bounds of the allocated buffer.



Royal Holloway  
University of London

Information Security Group


## Buffer overflow vulnerabilities

- Why is a buffer overflow a potential security vulnerability?
- Many different types of problem:
  - execution continues, and overwritten memory contains a changed variable value;
  - code might perform an unintended task;
  - code might crash.

72

Why is a buffer overflow a potential security vulnerability? Many different types of problem can result from a buffer overflow, including:

- execution could continue, and the overwritten memory will contain a changed variable value;
- the code might perform an unintended task;
- the code might crash.



Royal Holloway  
University of London

Information Security Group

## Buffer overflow outcomes

- More serious outcomes:
  - attacker initiates buffer overflow deliberately and manipulates results;
  - attacker modifies security relevant data or runs security relevant code;
  - attractive targets are return addresses and security settings;
- Commonly, attacker compromises a privileged program, starts an interactive session, and runs with privileged status.

73

More serious possible outcomes include the following:

- an attacker might initiate the buffer overflow deliberately and manipulate the results for his/her own ends;
- an attacker might modify security relevant data or run security relevant code;
- attractive targets for such attacks (i.e. types of data to be overwritten) include return addresses and security settings.

Commonly an attacker compromises a privileged program, establishes an interactive session and runs with privileged status. With enhanced privileges the attacker can have complete control of the machine.


## Buffer overflow causes

- Vulnerabilities typically arise in languages like C and C++ where programmer can allocate and de-allocate memory.
- Such activities are prone to errors which are hard to detect.
- As many as 50% of CERT alerts in recent years have resulted from buffer overflow problems.

74

Such vulnerabilities typically arise with the use of programming languages like C and C++, where the programmer is free to allocate and de-allocate memory. Such activities are prone to errors which are hard to detect.

As many as 50% of CERT alerts in recent years have resulted from buffer overflow problems.



Royal Holloway  
University of London

Information Security Group

## Vulnerabilities V – use-after-free

- Between 2007 and 2013, was a sharp decrease in exploits leading to stack corruption (including buffer overflows).
- Seems to be down to (a) effective mitigations, and (b) better coding, e.g. resulting from use of code analysis tools.
- However, at the same time, there has been a big increase in use-after-free exploits, which attack the **heap**.


75

As reported in Volume 16 of Microsoft's Security Intelligence Report (SIR), in the years up to 2013 a sharp decrease was seen in exploits leading to stack corruption, including stack-based buffer overflows. This class of vulnerabilities accounted for 54% of known exploited Microsoft remote code execution CVEs (Common Vulnerabilities and Exposures) in 2007, but just 5% in 2013. This change seems to be due to a combination of mitigations and better coding, including the use of code analysis tools.

At the same time, an increasing number of exploits of use-after-free vulnerabilities have been exploited. These exploits focus on the **heap** instead of the stack.

See: <http://www.microsoft.com/security/sir/default.aspx> for the SIR. Volume 16 is an extremely interesting read.

Note that the reference here to CVE is to those vulnerabilities that are published in the US Government's National Vulnerability Database (NVD) – see [nvd.nist.gov](http://nvd.nist.gov). All vulnerabilities in the NVD are given a CVE identifier (e.g. CVE2012-1889).



Royal Holloway  
University of London

Information Security Group

## Use-after-free: a definition

- Heap used for temporary memory space.
- Referencing any memory after it has been freed can cause a program to crash.
- Use of heap-allocated memory after it has been freed or deleted leads to undefined system behaviour (**use-after-free**).
- Use after free errors occur when a program continues to use a pointer after it has been freed.

76


The **heap** for a process consists of memory space that is allocated for temporary use.

Referencing memory after it has been freed (i.e. after it is no longer required for the purpose for which it was allocated) can cause a program to crash. The use of heap-allocated memory after it has been freed or deleted leads to undefined system behaviour (this is what is meant by **use-after-free**).

Use-after-free errors occur when a program continues to use a pointer after it has been freed. Use-after-free errors have two common and sometimes overlapping causes:

- error conditions and other exceptional circumstances;
- confusion over which part of the program is responsible for freeing the memory.

[Thanks to [www.owasp.org](http://www.owasp.org) for their explanation of use-after-free, on which this and the next slide have been based. For further information see also the wikipedia article on dangling pointers.]



Royal Holloway  
University of London

Information Security Group


## Effects

- Exploit of use-after-free vulnerability can have varying effects, e.g. corruption of data or execution of arbitrary code.
- One way data corruption can arise is if memory is reallocated.
- If original pointer used and memory is changed, then problem arises.

77

The use of previously freed memory can have any number of adverse consequences – ranging from the corruption of valid data to the execution of arbitrary code, depending on the instance and time of occurrence of the flaw.

One of the simplest ways in which data corruption can occur involves the system's reuse of the freed memory. In this scenario, the memory in question is allocated to another pointer validly at some point after it has been freed. The original pointer to the freed memory is used again and points to somewhere within the new allocation. If use of the original pointer causes the data to be changed, this will corrupt the validly used memory; this will induce undefined behaviour in the process.



Royal Holloway  
University of London

Information Security Group


## Vulnerabilities – mitigations

- Conclude discussion of software vulnerabilities by considering ways in which they can be avoided.
- Mitigations include:
  - **prevention:** better coding to try to avoid vulnerabilities;
  - **detection/mitigation:** modifications to code to detect and/or reduce impact of a vulnerability, if present.

78

We conclude this discussion of software vulnerabilities by considering ways in which they can be avoided. Possible mitigations include:

- **prevention:** better coding techniques can be used to help avoid the vulnerabilities being present;
- **detection/mitigation:** modifications can be made to the code to detect and/or to try to reduce the impact of a vulnerability if it is present.



Royal Holloway  
University of London

Information Security Group

## Prevention – input validation I


- Program should always check that data supplied to program is appropriate, e.g.:
  - is it of appropriate type?
  - is it of appropriate size?
  - does it contain appropriate characters?

79

A program should always check that any external data supplied to the program is appropriate, e.g.:

- is it of an appropriate type?
- is it of an appropriate size?
- does it contain appropriate characters?

Such validation checks are in many cases simple to perform, but are very often neglected.



Royal Holloway  
University of London

Information Security Group


## Prevention – input validation II

- Wide range of vulnerabilities (not only overflow vulnerabilities) can arise because of a lack of input validation, e.g.:
  - web-based front-ends for SQL databases (back-end database compromised by invalid input);
  - CGI scripts (Apache web server compromised by passing inputs containing hex characters);
  - Unicode injection (IIS web server compromised by inputs containing Unicode characters).

80

A wide range of vulnerabilities (not only overflow vulnerabilities) can arise because of a lack of input validation. For example, vulnerabilities of this type can arise in:

- web-based front-ends for SQL databases (where the back-end database is compromised by input inadequately validated by front end) – see, for example:  
[http://cert.uni-stuttgart.de/advisories/apache\\_auth.php](http://cert.uni-stuttgart.de/advisories/apache_auth.php)  
and  
<http://www.us-cert.gov/cas/techalerts/TA04-160A.html>
  
- CGI scripts (e.g. where an Apache web server was compromised by passing inputs containing hexadecimal characters) – see:  
<http://www.cert.org/advisories/CA-1996-06.html>
  
- Unicode injection (e.g. where an IIS web server was compromised by passing inputs containing Unicode characters) – see:  
<http://www.kb.cert.org/vuls/id/111677>



Royal Holloway  
University of London

Information Security Group


## Prevention – safe functions and libraries

- Many string-handling functions in C are unsafe because they allow unbounded copying of an input string to a buffer:
  - never use `gets()` [prefer `fgets()`];
  - use functions with a parameter that limits number of characters copied to a buffer [e.g., use `strcpy_s()` or `strncpy()` instead of `strcpy()`];
  - library `Strsafe.h` developed by Microsoft to replace the standard C library string functions;
  - other similar libraries include `SafeStr` and `Vstr`.

81

Many string-handling functions in C are unsafe because they allow unbounded copying of an input string to a buffer:

- never use `gets()` (use `fgets()` instead);
- use functions that take a parameter that limits the number of characters that are copied to a buffer (for example, use `strcpy_s()` or `strncpy()` instead of `strcpy()`);
- the library `Strsafe.h` has been developed by Microsoft to replace the standard C library string functions;
- other similar libraries of 'safe' string-handling functions include `SafeStr` and `Vstr`.



Royal Holloway  
University of London

Information Security Group


## Mitigation – Non-executable stacks

- Most commercial operating systems support non-executable stacks:
  - only prevent attacks that inject code into the stack;
  - do not prevent arc injection attacks;
  - do not prevent buffer overflows on the stack (so attacker can still crash the machine);
  - break legitimate applications that execute code on the stack.

82

Most commercial operating systems support non-executable stacks, i.e. where the processor will refuse to execute if the program counter points to a stack location. Such precautions:

- only prevent attacks that inject code into the stack;
- do not prevent arc injection attacks (i.e. where return addresses are modified);
- do not prevent buffer overflows on the stack (so an attacker can still crash the machine);
- break legitimate applications that execute code on the stack.



Royal Holloway  
University of London

Information Security Group

## Mitigation – NX bit I


- **NX (No eXecute) bit** is a CPU technology used to separate memory used for storage of code and data.
- OS supporting NX bit may mark certain areas of memory as non-executable.
- Processor will refuse to execute code residing in these areas of memory.

83

The **NX (No eXecute) bit** is a technology used in CPUs to segregate areas of memory used for storage of processor instructions (code) and for data.

An OS supporting the NX bit may mark certain areas of memory as non-executable.

The processor will then refuse to execute code residing in these areas of memory.



Royal Holloway  
University of London

Information Security Group

## Mitigation – NX bit II

- General technique known as executable space protection:
  - mitigate risk of malicious insertion of code into a program's data storage area;
  - e.g. as result of buffer overflow attack.
- Intel markets the feature as the **XD bit**, for eXecute Disable.
- AMD uses the name **Enhanced Virus Protection**.

84

The general technique is known as executable space protection. It is used to prevent execution of code maliciously inserted into a program's data storage area, e.g. as a result of a buffer overflow attack.


Intel markets the feature as the **XD bit**, for eXecute Disable. AMD uses the name **Enhanced Virus Protection**.

Microsoft Windows uses NX protection. In Windows XP and Server 2003 (and later versions of Windows), the feature is called Data Execution Prevention (DEP). If the processor supports this feature in hardware, then the NX features are turned on in Windows XP/Server 2003 by default.

However, even if the processor does not support it, DEP has a 'software-based mode', which attempts to emulate the non-executable stack protection even without support from the processor. Software-enabled DEP is weaker than the hardware-enforced version, but still provides a level of protection.

For further information about DEP see:

<http://windows.microsoft.com/en-US/windows-vista/Data-Execution-Prevention-frequently-asked-questions>

Royal Holloway  
University of London

Information Security Group


## Mitigation – stack randomisation

- Randomise layout of stack – harder for attacker to predict where return addresses and local variables are.
- This makes it more difficult to exploit buffer overflow vulnerabilities.
- However, applications require re-compilation.

85

Randomising the layout of the stack makes it harder for the attacker to predict where return addresses and local variables will be. This makes it more difficult to exploit buffer overflow vulnerabilities.

However, for this to be effective applications require re-compilation.



Royal Holloway  
University of London

Information Security Group

## Mitigation – ASLR


- More general technique known as Address Space Layout Randomisation (ASLR).
  - randomly arrange positions of key data areas, usually including base of executable and positions of libraries, heap, and stack in process address space.
  - hinders some security attacks by making it difficult to predict target addresses.
- Windows Vista and subsequent releases of Windows have ASLR enabled by default, although only for certain executables and dynamic link libraries.

86

A more general technique of this type is known as Address Space Layout Randomisation (ASLR). It involves randomly arranging the positions of certain key data areas, usually including the base of the executable and the positions of libraries, heap, and stack in the process's address space.

It hinders some security attacks by making it difficult to predict target addresses.

Windows Vista and subsequent releases of Windows have ASLR enabled by default, although only for certain executables and dynamic link libraries.



Royal Holloway  
University of London

Information Security Group

## Detection – software testing I


- Number of ways to test software for vulnerabilities.
- **Static analysis:**
  - examine the source code looking for use of unsafe functions and lack of input validation;
  - not all source code is available for scrutiny;
  - may generate false positives and false negatives.

87

There are a number of ways of testing software for vulnerabilities. We briefly mention two such possibilities.

**Static analysis** involves examining the source code to look for the use of unsafe functions and lack of input validation.

However not all source code is available for scrutiny. Such an approach may also generate false positives and false negatives.



Royal Holloway  
University of London

Information Security Group


## Detection – software testing II

- **Black box testing:**
  - study the behaviour of software for a wide variety of different inputs;
  - exhaustive testing is not usually possible;
  - may find vulnerabilities, but does not guarantee absence of vulnerabilities.
- One example is **fuzz testing**, i.e. giving a program random inputs.

88

**Black box testing** involves studying the behaviour of software for a wide variety of different inputs. Exhaustive testing (of all possible inputs) is not usually possible. Such an approach may find vulnerabilities, but does not guarantee the absence of vulnerabilities.

One particular example of black box testing is known as **fuzz testing**, where a program is executed many times with randomly selected inputs (possibly semi-structured).



Royal Holloway  
University of London

Information Security Group

## Detection/mitigation – canaries

- **Canary** is special data added to stack to detect attempts to overwrite the return address; it:
  - is overwritten if the return address is overwritten;
  - is initialised immediately after the return address has been saved to the stack;
  - is checked (to see if it has been modified) before return address popped off the stack
- **Canaries detect buffer overflows:**
  - prevent execution of code and arc injection exploits;
  - do not prevent buffer overflows.

89


A **canary** is a special piece of data added to the stack to detect attempts to overwrite the return address. It:

- will be overwritten if the return address is overwritten;
- is initialised immediately after the return address has been saved to the stack;
- is checked (to see if it has been modified) before the return address popped off the stack

Canaries can be used to detect buffer overflows:

- can help prevent execution of code and arc injection exploits.

However, they do not prevent buffer overflows.



Royal Holloway  
University of London

Information Security Group

## Buffer overflow mitigations

- Occurrences can be reduced:
  - careful coding, avoiding dangerous calls;
  - software scanning tools can be used to find and remove buffer overflow weaknesses;
  - programming in *type safe* language can reduce vulnerabilities;
  - apply the principle of least privilege to design software so any buffer overflows that do succeed will result in minimal damage.

90

Occurrences of buffer overflows can be reduced:

- by careful coding, avoiding dangerous calls;
- software scanning tools can be used to find and remove buffer overflow weaknesses;
- programming in a **type safe** language can reduce vulnerabilities (see later slides);
- apply the principle of least privilege to design software so any buffer overflows that do succeed will result in minimal damage.

Other tools and techniques are available to programmers. Users of software susceptible to buffer overflow attacks should be aware of the security vulnerability alerts from organisations such as CERT, and follow any relevant advice.

A helpful review of buffer overflow attacks and mitigations can be found in:

B. M. Padmanabhuni and H. B. K. Tan, Defending against buffer-overflow vulnerabilities. *IEEE Computer*, **44 no. 11** (November 2011), pp.53-60.

[Download a copy at:

<http://www.chrismitchell.net/IY5512/mco2011110053.pdf>].



Information Security Group

# Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

91

Inadequacies in programming are not the only sources of vulnerabilities in software.



Royal Holloway  
University of London

Information Security Group

## User authentication failures

- Significant source of vulnerabilities is password-based authentication:
  - poor choice of passwords;
  - inadequate protection of passwords by users.
- Implementations of authentication services are often flawed:
  - Unix and Windows have (in the past) both used authentication methods with small key spaces;
  - widely used protocols (e.g. Kerberos) are vulnerable to off-line dictionary attacks.


92

One significant source of vulnerabilities is password-based authentication:

- users often make poor choices for passwords;
- users often inadequately protect their passwords.

Implementations of authentication services are sometimes flawed:

- in the past, Unix and Windows have both used authentication methods with small key spaces;
- widely used protocols (e.g. Kerberos) are vulnerable to off-line dictionary attacks, when passwords are poorly chosen.



Royal Holloway  
University of London

Information Security Group

## Authorisation issues

- Access control mechanism only as good as the policy it enforces:
  - if access control list (ACL) is configured to allow access, then reference monitor will allow access;
  - many security groups contain incorrect members.
- Often difficult to map enterprise security policy to an access control mechanism:
  - managers do not think in terms of ACLs;
  - one benefit of RBAC is a closer match between enterprise specification and actual implementation.

93


An access control mechanism is only as good as the policy it enforces:

- if an access control list (ACL) is configured to allow access, then the reference monitor will allow access;
- in practice, access control systems are often misconfigured; for example, security groups may contain the wrong members.

It is often difficult to map an enterprise security policy to the specification scheme provided by the access control mechanism:

- system managers do not think in terms of ACLs;
- one benefit of role based access control (RBAC) is a closer match between an enterprise specification and the actual implementation.

[Note that we examine ACLs and RBAC later in the course].



Royal Holloway  
University of London

Information Security Group

## Timeliness

- Security-relevant information not always current:
  - user account may remain active after user has left;
  - membership of groups and roles may not reflect changes to organisational structure and job function;
  - revocation lists for public keys may be difficult to obtain and process.
- Management of authorisation and authentication data is extremely important.
- Short-lived credentials can help, but availability may become an issue.


94

Security-relevant information may not always be up to date:

- a user account may remain active long after the user has left organisation;
- membership of groups and roles may not reflect changes to organisational structure and job function;
- revocation lists for public keys may be difficult to obtain and process.

These vulnerabilities mean that management of authorisation and authentication data is extremely important.

Short-lived credentials can help, but availability may become an issue if credentials are allowed to expire without replacement.



Royal Holloway  
University of London

Information Security Group

## Race conditions


- Race conditions:
  - arise in environments where multiple threads or processes are running at same time, which may interact with one another.
- Multiple computations may access shared data in such a way that results depend on the sequence of accesses.
- An attacker can try to change a value after it has been checked, but before it is used.

95

A **race condition** is a flaw in a process where the result of the process is dependent on the sequence or timing of other events. The term originates with the idea of two signals in a circuit 'racing each other' to influence the output.

Race conditions in computer systems arise in environments where multiple threads or processes are running at the same time, which may interact with one another. Multiple computations may access shared data in such a way that their results depend on the sequence of accesses.

An attacker can try to change a value after it has been checked, but before it is used.



Royal Holloway  
University of London

Information Security Group

## Cryptographic issues I

- Many software security problems arise from misuse of cryptography.
- Misuse of randomness:
  - many programs use randomness;
  - common method of getting random bits is a deterministic pseudo-random generator;
  - randomness source must be implemented well – e.g. counting milliseconds since midnight on system clock not good enough!

96

Many software security problems arise from misuse of cryptography.

Misuse of randomness is common:

- many programs require sources of randomness (e.g. to generate keys for cryptographic algorithms, or random challenges in protocols);
- the most common method of generating 'randomness' is to use a deterministic pseudo-random generator;
- a randomness source must be designed and implemented well – simply counting the number of milliseconds since midnight on the system clock is not normally good enough!

Royal Holloway  
University of London

Information Security Group

## Cryptographic issues II

- Poor key management:
  - cryptographic key management is complex issue;
  - cannot protect long cryptographic keys with potentially weak short passwords.
- Customised cryptography:
  - not normally a good idea to use customised cryptography, except in a large organisation with a dedicated cryptographic team;
  - developers often believe they can design their own cryptographic primitives, and falsely think they will be more secure than if they adopt standardised techniques.


97

Poor key management:

- cryptographic key management is a complex issue;
- cannot protect long cryptographic keys with potentially weak short passwords.

Customised cryptography:

- it is not normally a good idea to use customised cryptography, except in a large organisation with a dedicated cryptographic team;
- it is surprising how often developers (falsely) believe that they can easily design their own cryptographic primitives, and by doing so they will be developing a more secure system than if they adopt standardised techniques.



Royal Holloway  
University of London

Information Security Group

## Users I

- Finally, every user is a vulnerability:
  - user authorised to view sensitive data can leak it;
  - users (mis)configure security policies;
  - users choose bad passwords;
  - users do not keep secrets.
- The education and training of users is vital:
  - users must understand how their actions can compromise security;
  - users must be made aware of their responsibilities with respect to security.

98


Finally, every user represents a vulnerability:

- any user authorised to view data to which access is restricted can leak that data;
- users (mis)configure security policies;
- users choose bad passwords;
- users do not keep secrets.

The education and training of users is vital:

- users must understand how their actions can compromise security;
- users must be made aware of their responsibilities with respect to security.

Of course, there is a limit to what can be achieved with education. Unreasonable demands must not be placed on users in the name of security, or they will simply ignore them in order to get their job done. That is, security must be **usable**.



Royal Holloway  
University of London

Information Security Group

## Users II

- Social engineering is a huge threat.
- Many types of such attack, e.g.:
  - phishing;
  - other malicious emails;
  - phone calls.
- Recent studies of malware on PCs suggest that much malware present as a result of social engineering attacks.

99

To conclude, it is important to realise that so called **social engineering** is a huge threat. Wikipedia defines social engineering as the art of manipulating people into performing actions harmful to themselves or divulging confidential information.

There are many types of social engineering attack, including:

- phishing;
- other malicious emails;
- malicious phone calls (e.g. persuading an innocent user to do damage to their own machine).

Many more examples are described on the Wikipedia page devoted to the subject.

Recent studies of malware on PCs suggest that a significant proportion (most) malware is present as a result of social engineering attacks, not because of any vulnerabilities in the system itself.



Royal Holloway  
University of London


Information Security Group

## Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

100

We conclude this part of the course by discussing malware. Malware includes both software designed to exploit vulnerabilities in systems, as well as software which is installed on a machine (and performs harmful acts) as a result of the successful exploitation of a vulnerability.



Royal Holloway  
University of London

Information Security Group

## Malicious software (malware)


- **Malware** (malicious software) is software designed to infiltrate or damage a computer system without owner consent.
- Covers a variety of forms of hostile, intrusive, or annoying software.
- Malware includes viruses and worms, but there are many classes.

101

**Malware** (malicious software) is software designed to infiltrate or damage a computer system without the owner's consent. The word is used to cover a wide variety of forms of hostile, intrusive, or annoying software.

Malware includes viruses and worms, but there are many other categories of malware.

Initial infection by a worm or virus is often followed by the installation of additional malware such as rootkits, trojan horse programs, adware and spyware.



Royal Holloway  
University of London

Information Security Group


## Relationship to vulnerabilities

- To be a threat, malware must execute on a victim's computer.
- How might this happen?
  - Typically by exploiting a vulnerability in some way.
  - Vulnerabilities in executing software might enable attacker to execute malware.
  - User may be vulnerability, e.g. social engineering may persuade user to execute malware.

102

To be a threat, malware must execute on a victim's computer. There are many ways in which this might happen. Typically it involves the exploitation of a vulnerability. Examples include the following.

- Vulnerabilities in executing software might enable an attacker to execute malware. For example, a 'data file' opened by a word processor might contain a malicious code fragment, which is executed as a result of poor data validation by the word processor, or through use of an unprotected macro execution facility.
- Depending on how the user's computer is configured, malware contained on a portable memory device may be automatically executed when connected to a victim's computer.
- The user may be the vulnerability him/herself; for example, the user might be persuaded through a carefully crafted email or web page to download and execute a piece of malware (thinking it is something else).
  
- Some malware, notably viruses and worms (which we look at a little later in this part of the course), propagates itself automatically from computer to computer. However, malware does not have to be self-propagating to be very dangerous.



Royal Holloway  
University of London

Information Security Group


## Threat to PCs I

- Historically, PCs at particular risk from malware because originally designed for single users.
- Older versions of Windows had no security mechanisms for separating users, or for separating the user from the 'system'.
- Even in more recent versions of Windows there are limited mechanisms to stop intentional modifications to the system and/or user files (users typically run with 'administrator' privileges).

103

Historically, PCs have been at particular risk from malicious code because they were originally designed for single users. PC operating systems such as Windows originally contained no security mechanisms for separating users, or for separating the user from the 'system'. Neither did they incorporate any security mechanisms to stop intentional modifications to the system and/or user files.

Even in more recent versions of Windows there are limited security mechanisms to stop intentional modifications to the system and/or user files (users typically run with 'administrator' privileges). Note that Vista, Windows 7 and Windows 8 contain measures designed to try to address this – steps requiring administrator privileges require explicit user authorisation.



Royal Holloway  
University of London

Information Security Group

## Threat to PCs II

- Spread of PCs meant development of commercial products (e.g. word processors, spreadsheets), games, and shared use of PCs.
- Led to exchange of data and software between PCs, e.g. via portable storage media or Internet.
- This exchange of data, combined with lack of protection, gave rise to a serious threat.
- More recently, main threat has moved to social engineering attacks (e.g. phishing) and exploitation of code vulnerabilities.


104

The spread of PCs has meant the development of a variety of commercial products (e.g. word processors, spreadsheets), as well as games. Also single PCs are shared by multiple users.

The growth of commercial applications has led to the widespread exchange of data between PCs, e.g. via portable storage media or the Internet.

These exchanges of data and software, combined with the lack of protection in PC operating systems, originally gave rise to the malware threat (historically in the form of viruses).

More recently, the main threat has moved to social engineering attacks (e.g. phishing) and exploitation of code vulnerabilities.



Royal Holloway  
University of London

Information Security Group

## The threat


- Malicious code can do anything any other program can, e.g. displaying a message, erasing stored files, reformatting disks.
- Can lie dormant for a long period.
- Can be triggered by an event, e.g. use of a program or reaching a certain date.

105

Malicious code can do anything any other program can, e.g. displaying a message, erasing stored files, and/or reformatting disks.

Malicious code can lie dormant within the PC for a long period. Its effects can be triggered by various types of event, e.g. the use of some program or reaching a certain date.

What makes the types of malicious code known as *viruses* and *worms* particularly dangerous is their ability to replicate and spread.



Royal Holloway  
University of London

Information Security Group

## History I


- Many early infectious programs, including first Internet Worm and many viruses, written as experiments or jokes – meant to be harmless/annoying rather than cause serious damage.
- Melissa (a 1999 virus) appears to have been written for fun.
- More hostile intent in some programs designed to vandalise or cause data loss.

106

Many early infectious programs, including the first Internet Worm and many of the first viruses seen in the wild, were written as experiments or jokes – they were intended to be harmless or annoying, rather than cause serious damage.

Viruses such as Melissa (from 1999) appear to have been written for fun.

There was, nevertheless, slightly more hostile intent in some programs that were designed to vandalise or cause data loss.



Royal Holloway  
University of London

Information Security Group


## History II

- Some MS-DOS viruses, and Windows ExploreZip worm, were designed to destroy files on a hard disk, or corrupt the file system.
- Network-borne worms such as the 2001 Code Red worm or the Ramen worm were similar, and were designed to vandalise web pages.

107

Even in the 'early days' of malware, some MS-DOS viruses and worms (e.g. Windows ExploreZip), were designed to destroy files on a hard disk, or to corrupt the file system by writing junk data.

Network-borne worms such as the 2001 Code Red worm or the Ramen worm were similar, and were designed to vandalise web pages.



Royal Holloway  
University of London

Information Security Group

## History III


- Since broadband Internet access became common, more malware designed for (criminal) profit – as opposed to jokes/random damage.
- Since around 2003, many common viruses and worms designed to take control of computers for criminal exploitation.
- Infected **bots** are used as part of **botnets** to:
  - send email spam;
  - host contraband data,, e.g. child pornography; or
  - engage in distributed denial-of-service attacks.

108

Since broadband Internet access has become common, malware has been increasingly designed for (criminal) profit, as opposed to previous malware which was mainly designed either just to prove a point for fun or to cause random damage.

Since around 2003, many of the commonly occurring viruses and worms have been designed to take control of computers for criminal exploitation.

Infected **bots**, typically making up a centrally controlled **botnet**, can be used for a variety of purposes, including to send email spam, to host contraband data such as child pornography, or to engage in distributed denial-of-service attacks, e.g. as a form of extortion.



Royal Holloway  
University of London

Information Security Group


## History IV

- Early in 21st century a new category of malware emerged: **spyware** – programs designed to monitor user web browsing and display unsolicited advertisements.
- Spyware programs do not spread like viruses or worms; they are generally installed by exploiting security holes or packaged with user-installed software.

109

Early in the 21st century another category of malware emerged, known as spyware, i.e. programs designed to monitor user web browsing and display unsolicited advertisements.

Spyware programs are typically not infectious, like viruses, or self-replicating, like worms; they are generally installed by exploiting security holes or are packaged with user-installed software (one prominent early example was circulated by Kazaa).



Royal Holloway  
University of London

Information Security Group

## Malware vectors I

- **Malware vector** is a means by which malware propagates from one machine to another.
- Two very important vectors:
  - (unpatched) vulnerabilities in applications or operating system components – typically used by worms;
  - social engineering attacks (including Trojans).

110


A **malware vector** is a means by which malware propagates from one machine to another.

There are two particularly important vectors:

- (unpatched) vulnerabilities in applications or operating system components;
- social engineering attacks (including Trojans – see next slide).

The former vector is typically used by worms in propagating across networks between machines.

Viruses propagate in a range of ways, including through social engineering, and through use of unintended functionality (see next slide), as well as sometimes using vulnerabilities.



Royal Holloway  
University of London

Information Security Group

## Malware vectors II


- **Trojan horse** is apparently useful program containing hidden code that performs some unwanted or harmful function.
- **Examples:**
  - to access another user's files, create a Trojan program which changes the user's file permissions so files are readable by any user.
  - modify a compiler to insert extra code in executables; if system login program is compiled then additional code is a trapdoor.

111

A *Trojan horse* is an apparently useful program containing hidden code that, when invoked, performs some unwanted or harmful function. Examples include:

- in order to gain access to another user's files on a shared system, the attacker creates a Trojan program which, when executed, changes the invoking user's file permissions so that files are readable by any user.
- modify a compiler so that it inserts additional code into certain files as it compiles them; if the system login program is compiled then the additional code creates a trapdoor.

Trojans can be regarded as a special type of social engineering attack.



Royal Holloway  
University of London

Information Security Group

## Malware vectors II

- By no means the only malware vectors.
- Another important vector is the exploitation of intended functionality in unintended ways.
- Examples include:
  - autorunning of malicious executables on portable media;
  - macro viruses;
  - ‘drive by’ downloads to web browsers.


112

Vulnerabilities in software and social engineering attacks (including Trojans) are by no means the only malware vectors.

Another important malware vector is the exploitation of intended functionality in unintended ways.

Examples include:

- automatic execution (autorunning) of malicious executables on portable media;
- macro viruses, i.e. viruses containing programs written in a macro language, i.e. a special programming language built into an application such as a word processor or database, intended to be used to automate certain actions;
- ‘drive by’ downloads to web browsers, i.e. where a user innocently visits a web page which downloads malware to the user’s computer; some older browsers would execute such software automatically.



Royal Holloway  
University of London

Information Security Group


## Malware types

- As well as using various vectors for transmission, malware can take many forms.
- Next review some of the more prominent categories of malware.

113

As well as using a variety of vectors for transmission, malware can take many forms.

We next review some of the more prominent categories of malware. Note that malware does not always conveniently fit into one of the categories, and also the categories we present are not mutually exclusive. However, most commonly encountered malware fits in at least one of these categories.



Royal Holloway  
University of London

Information Security Group


## Malicious code I: Viruses

- **Computer virus:** program that infects other programs by modifying them to include a copy of the virus program, which can then go on to infect other programs.
- Its **payload** can be non-existent, harmless (e.g. playing a tune), or very damaging (e.g. it may modify or delete files).

114

*Computer virus:* a program that infects other programs by modifying them to include a copy of the virus program, which can then go on to infect other programs.

Its **payload** can be non-existent, harmless (e.g. playing a tune), or very damaging (e.g. it may modify or delete files).



Royal Holloway  
University of London

Information Security Group

## Malicious code II: Worms


- **Worm:** replicating but non-infecting program using a network to spread.
- **Examples**
  - electronic mail: a worm mails a copy of itself to other systems;
  - remote execution: a worm executes a copy of itself on another system;
  - remote login: a worm logs in to a remote system and uses commands to copy itself to that system.

115

A *worm* is a replicating but non-infecting program that uses network connections to spread from system to system.

Examples include:

- electronic mail: a worm mails a copy of itself to other systems;
- remote execution: a worm executes a copy of itself on another system;
- remote login: a worm logs onto a remote system as a user and then uses commands to copy itself from one system to the other.



Royal Holloway  
University of London

Information Security Group

## Malicious code III: Rootkits

- Malicious program may conceal its presence.
- **Rootkits** hide by changing OS code:
  - typically installed following successful attack;
  - OS utilities that monitor process activity are replaced.
- Rootkit can prevent a process from being visible in process list, and/or hide its files.
- Impossible to detect if OS kernel modified:
  - only sign of infection may be reduced performance.
- Removing rootkit may need full OS reinstall.

116

Once a malicious program has been installed on a system, it is often useful to the creator if it stays concealed.

**Rootkits** deliberately hide themselves by subverting or replacing operating system code:

- typically installed following successful attack (exploitation of a vulnerability), e.g. via a Worm;
- OS utilities that monitor process activity are replaced with programs that hide malicious activity.


Rootkits can prevent a malicious process from being visible in the system process list, or keep its files from being read.

It may be impossible to detect a rootkit if the OS kernel has been partially replaced:

- only sign of rootkit infection may be reduced performance.

Removing a rootkit often requires a complete reinstall of the OS.

Originally, a rootkit was a set of tools installed by a human attacker on a Unix system where the attacker had gained administrator (root) access. Today, the term is used more generally for concealment routines in a malicious program.



Royal Holloway  
University of London

Information Security Group


## Malicious code IV: Spyware

- **Spyware** gathers information about users, shows them adverts (e.g. via pop-ups), or alters browser behaviour for benefit of creator.
- E.g., some spyware programs redirect search engine results to paid adverts.
- Spyware sometimes installed as Trojans.
- Differs from other malware: creators often operate openly, e.g. to sell adverts on pop-ups.

117

Over the last ten-fifteen years or so, one of the most costly forms of malware in terms of time and money spent in recovery has been the broad category known as spyware. Spyware programs are commercially produced for the purpose of gathering information about computer users, showing them advertisements (e.g. via pop-ups), or altering web-browser behaviour for the financial benefit of the spyware creator. For example, some spyware programs redirect search engine results to paid advertisements

Spyware programs are sometimes installed as Trojan horses. They differ from other types of malware that their creators often present themselves openly as businesses, e.g. by selling advertising on the pop-ups created by the malware. Most such programs present the user with an end user agreement which supposedly protects the creator from prosecution under computer abuse laws - however, the precise legality of these agreements remains to be determined.



Information Security Group

## Malicious code V: Concealed malware

- Attackers often disguise malware to hide its presence on compromised systems:
  - viruses are often polymorphic, i.e. each instance of the viral code looks different;
    - signature-based AV software does not recognise new instances of the virus
    - cryptography used to encrypt parts of the code;
    - null operations are inserted into the virus code.
- Viruses can also be metamorphic, i.e. the viral code rewrites itself.


118

Attackers often disguise malware to hide its presence on compromised systems:

- as alluded to previously, viruses are often polymorphic, i.e. each instance of the viral code looks different;
  - signature-based AV software does not recognise new instances of the virus
  - cryptography used to encrypt parts of the code;
  - null operations are inserted into the virus code.

Viruses can also be metamorphic, i.e. the viral code rewrites itself.

Of course, rootkits (as mentioned previously) also conceal their presence.



Royal Holloway  
University of London

Information Security Group

## Effects – overview


- Some malware enables other malware to be installed on compromised computers, including:
  - spyware;
  - rootkits.
- This enables infected computers to be controlled remotely:
  - increasingly common for networks of compromised computers (botnets) to be exploited for commercial gain.

119

Some malware has as its main function enabling other malware to be installed on compromised computers, including spyware and rootkits.

The installed malware often enables infected computers to be controlled remotely:

- it has become common for criminal entities to exploit networks of compromised computers (botnets) for commercial gain.



Royal Holloway  
University of London

Information Security Group


## Effects I: Back doors

- **Back door** is means of bypassing normal authentication procedures.
- If system is compromised (e.g. using a method described), a backdoor can be installed to allow attacker future access.
- ‘Crackers’ use back doors to secure remote access to a computer, while trying to remain hidden from casual inspection.

120

A back door is a method of bypassing normal authentication procedures. Once a system has been compromised (by one of the methods described, or in some other way), one or more backdoors could be installed, in order to allow the attacker access in the future.

So called ‘crackers’ typically use backdoors to secure remote access to a computer, while attempting to remain hidden from casual inspection. To install back doors, these crackers may use Trojans, worms, etc.



Royal Holloway  
University of London


Information Security Group

## Effects II: Botnets

- To coordinate infected computers, attackers use coordinating systems known as **botnets**.
- In a botnet, the malware or **malbot** logs in to an Internet relay chat channel or other chat system.
- Attacker can give instructions to all infected PCs at once (e.g. for spamming or DDoS attacks).
- Botnets can be used to send upgraded malware to infected systems, keeping them resistant to anti-virus software or other security measures.

121

In order to coordinate the activity of many infected computers, attackers have used coordinating systems known as **botnets**. In a botnet, the malware or **malbot** (or just **bot**) logs in to an Internet relay chat channel or other chat system. The attacker can then give instructions to all the infected systems simultaneously, e.g. to conduct DDoS attacks or send email spam. Botnets can also be used to push upgraded malware to the infected systems, keeping them resistant to anti-virus software or other security measures.



Royal Holloway  
University of London


Information Security Group

## Effects III: DDoS attacks

- **Denial of Service (DoS) attacks** make computer resources unavailable.
- DoS attacks often used to prevent Internet site or service from working properly.
- DoS attackers target high-profile web servers, e.g. banks, credit card payment gateways.
- In **Distributed DoS (DDoS)** attack, multiple compromised systems target a system.
- Malware can carry DDoS attack mechanisms.
- Botnets can launch such attacks.

122

A **Denial of Service (DoS) attack** is an attempt to make a computer resource unavailable to its authorised users. Although the means to conduct a DoS attack, and the motives for and targets of a DoS attack may vary, it generally involves a concerted, malevolent effort to prevent an Internet site or service from functioning properly, temporarily or indefinitely. Perpetrators of DoS attacks typically target sites or services hosted on high-profile web servers such as banks, credit card payment gateways and even DNS root servers. A **Distributed Denial of Service (DDoS) attack** occurs when multiple compromised systems flood the bandwidth or resources of a targeted system, usually one or more web servers. Malware can carry DDoS attack mechanisms; one of the more well known examples of this was MyDoom – its DoS mechanism was triggered on a specific date and time. This type of DDoS involved hardcoding the target IP address prior to release of the malware and no further interaction was necessary to launch the attack.



Royal Holloway  
University of London

Information Security Group

## Effects IV: Spam(bots)

- **Spammer viruses** are commissioned by email spam gangs.
- Infected computers used as proxies to send out spam messages.
- Infected computers available in large supply and provide anonymity for spammer.
- Spammers have used botnets for DDoS attacks on anti-spam organisations.

123

*Spammer viruses*, such as the Sobig and MyDoom virus families, are commissioned by email spam gangs. Infected computers are used as proxies to send out spam messages. The advantage to spammers of using infected computers is that they are available in large supply (thanks to the virus) and they provide anonymity, protecting the spammer from prosecution. Spammers have also used infected PCs to target anti-spam organizations with DDoS attacks.



Royal Holloway  
University of London

Information Security Group

## Anti-malware strategies

- On-going arms race between writers of malware and anti-malware software.
- Malware often exploits lack of integrity protection; so anti-malware protection adds integrity control.
- Three main components to strategy:
  - **prevention**: stop malware infecting system;
  - **detection**: detect presence of malware:
  - **reaction**: remove malware from system.

124

We now consider how one should approach the problem of dealing with the threat posed by malware. First note that viruses exploit the lack of integrity protection in a system, and hence an anti-malware protection strategy involves adding integrity control to a system.

There are three main components to a complete anti-malware strategy:

- *Prevention*: i.e. stopping malware infecting the system;
- *Detection*: i.e. detection of the presence of malware in the system;
- *Reaction*: i.e. removal of malware from the system.

Historically, the main issues were viruses and worms, but the emphasis has changed over the last decade.



Royal Holloway  
University of London

Information Security Group

## Physical & administrative controls


- Such controls are ways of *preventing* malware infection.
- E.g.:
  - issue rules about installation of software on corporate machines;
  - restrict use of USB memory sticks.

125

Physical and administrative controls can be effective ways of preventing malware infection, including by viruses. There are many examples of suitable controls of this type. Examples include the following.

- Many companies will have a policy forbidding use of any software which has not been explicitly approved by the company's IT or security department. Thus all software must be procured by official means, minimising the risk of infection by malware.
- Restrict use of USB memory sticks for data transfer. USB memory sticks have become a very serious malware vector (particularly when combined with the autorun feature).

However, this all relies on user cooperation and understanding which is difficult to guarantee, particularly when facing threats like spear phishing.

Royal Holloway  
University of London


Information Security Group

## Cryptographic checksums

- Cryptographic checksums are a standard integrity protection mechanism.
- Compute checksum for 'clean' version of code to be protected, and store checksum in secure place.
- Whenever file is used, first verify checksum.

126

Cryptographic checksums are a standard integrity protection mechanism. To use it to protect against malware infection, compute the checksum for a 'clean' version of the code which is to be protected, and store the checksum in a secure place. Whenever the file is used, the checksum must first be verified.



Royal Holloway  
University of London

Information Security Group

## Malware signatures


- Malware cannot be completely invisible – code must be stored somewhere.
- So malware can be detected by searching for characteristic sequences of information, called *signatures*.
- Can be (partially) hidden by splitting up malware code or by malware storing itself in encrypted form (*polymorphic* malware).

127

On the positive side, malware stored within a system cannot be completely invisible, since the malware code must be stored somewhere. This means that, in principle at least, malware can be detected by searching for characteristic sequences of information, called *signatures*.

Malware can, however, be (partially) hidden, either by splitting up the code, or having the malware store itself in encrypted form (*polymorphic* viruses) using a varying encryption key, which must also be stored with the encrypted code. However, the part of the malware code responsible for performing the decryption process must be left unencrypted, so that even polymorphic malware will have a (small) signature.

One fundamental limit of the signature approach to malware detection is that the malware must be known about, and hence it does not work against 'new' malware.



Royal Holloway  
University of London

Information Security Group

## Scanners

- Scanners search files for malware signatures, identifying known strains of malware.
- Need to know about malware to detect it, and hence need constant updating.
- Widely used, but need to be efficient (can significantly slow system when running).
- Also threat remains of polymorphic malware and use of *mutation engines*.

128

*Virus/malware scanners* search files for malware signatures, which can be used to identify infections by known strains of malware. Of course, the scanner will need to know about a type of malware to detect it, and hence such scanners need constant updating. Scanners are widely used, but they need to be efficient. They can significantly impair system performance if they are running constantly, which is the ideal way of using them.

Also, the threat remains of polymorphic malware, which can avoid detection since they will not have a single characteristic signatures. Also, *mutation engines* automate the generation of new malware strains, which can result in a constant supply of new malware whose signature may not be recognised by a scanner.

At one time signature-based scanners were regarded as a key defence, but in recent years their importance has declined.




Information Security Group

# Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

129

We next examine viruses in greater detail.



Royal Holloway  
University of London

Information Security Group

## How viruses and worms work I


- Typically two components:
  - **replication mechanism**, and
  - **payload**, which acts on infected host.
- Payload usually activated by some trigger e.g. the date.
- Virus payload can do a wide variety of bad things ...

130

Typically, there are two components to a virus or worm:

- a replication mechanism, which enables the virus to make copies of itself, and
- a payload (which does something on the infected system to the advantage of the virus writer).

The payload is usually activated by some trigger. e.g. the date. A virus payload can do a wide variety of bad things, including installing additional malware, creating backdoors, etc. (see next slide).



Royal Holloway  
University of London

Information Security Group


## How viruses and worms work II

- Malicious code can:
  - change machine's protection state, selectively or randomly;
  - change user data, selectively or randomly;
  - lock the network;
  - steal resources;
  - steal or publish data, e.g. cryptographic keys;
  - create a backdoor.

131

The payload may, when triggered, do a variety of bad things:

- make changes to the machine's protection state, selectively or randomly;
- make changes to user data, selectively or randomly, e.g. trash the disk;
- lock the network, e.g. by replicating at maximum speed;
- steal resources, e.g. use the CPU for a DES key search;
- steal, or publish your data like cryptographic keys – in Sept 2000 a virus was reported that steals the passwords used to access the home banking system of a Swiss bank;
- create a backdoor, so that the creator can take over the system later, perhaps to launch a distributed denial of service attack.



Royal Holloway  
University of London

Information Security Group


## Types of virus I

- Many types of virus exist, including:
  - **Parasitic or program virus:** a common type
    - it attaches itself to an executable file and replicates when the infected file is executed, by finding other executable files to infect.
  - **Memory-resident virus:** lodges in memory as part of resident system program, and infects every program that executes.

132

Virus researchers have put considerable effort into developing schemes for classifying computer viruses. Many types of virus can be identified, and among the most significant are the following (note that these categories are not mutually exclusive):

- *Parasitic or program virus:* a common form; such a virus attaches itself to executables typically by appending itself to the program, and inserting (at beginning of program) a jump to the viral code. At the end of the viral code there is a jump back into the program. It replicates when the infected file is executed, by finding other executable files to infect. That is, then the programs are loaded in memory during execution, the virus becomes active, making copies of itself and infecting files on disk. Examples include: Sunday, Cascade.
- *Memory-resident virus:* lodges in memory as part of a resident system program, and then infects every program that executes.



Royal Holloway  
University of London


Information Security Group

## Types of virus II

- **Boot sector virus:** infects a (master) boot record and spreads when a system is booted from a disk containing the virus.
- **Multipartite virus:** hybrid of Boot and Parasitic virus that infects program files – when infected program is executed it infects the boot record.

133

- *Boot sector virus:* infects a master boot record, or boot record and spreads when a system is booted from the disk containing the virus. Examples include: Form, Disk Killer, Michelangelo, and Stone virus.
- *Multipartite viruses:* are a hybrid of Boot and Parasitic viruses. They infect program files, and when the infected program is executed these viruses infect the boot record. When you boot the computer next time the virus from the boot record loads in memory and then starts infecting other program files on disk. Examples include: Invader, Flip, and Tequila.



Royal Holloway  
University of London


Information Security Group

## Types of virus III

- **Stealth virus:** designed to hide itself from detection and anti-virus software.
- **Polymorphic virus:** stealth virus that mutates with every infection, making detection by virus signature impossible.
- **Macro viruses:** infects macros in a document or template – e.g. Word macro virus infects template (Normal.dot) used for default document formatting settings.

134

- *Stealth virus:* a form of virus explicitly designed to hide itself from detection and anti-virus software, e.g. it may use compression so that the infected program is exactly the same length as an uninfected version. Examples include: Frodo, Joshi, Whale.
- *Polymorphic virus:* a type of stealth virus that mutates with every infection, making detection by the signature of the virus (see later slide) difficult or impossible, i.e. it creates copies during replication that are functionally equivalent, but have different bit patterns. Examples include: Involuntary, Stimulate, Cascade, Phoenix, Evil, Proud, Virus 101.
- *Macro Viruses:* infects the macros within a document or template. For example, a Word macro virus is activated when you open a word processing or spreadsheet document, and it infects the Normal template (Normal.dot) - a general purpose file that stores default document formatting settings. Every document you open refers to the Normal template, and hence gets infected with the macro virus. Since this virus attaches itself to documents, the infection can spread if such documents are opened on other computers. Examples include: DMV, Nuclear, Word Concept.



Royal Holloway  
University of London


Information Security Group

## Example – Brain virus

- Brain virus written in 1986 and was the first virus to affect IBM-compatible PCs:
  - Brain virus was extremely common in the late 1980s;
  - was a boot sector and memory resident virus, and it typically propagated through floppy disks left accidentally in drives;
  - today's computers are rarely booted from portable media such as floppy disks.

135

The Brain virus was written in 1986 and was the first virus to affect IBM-compatible PCs. It was a boot sector and memory resident virus, and it typically propagated through (infected) floppy disks left accidentally in drives – once memory resident it infected all floppy disks inserted in the machine. The Brain virus was extremely common in the late 1980s. Today computers are rarely booted from portable media such as floppy disks.



Royal Holloway  
University of London

Information Security Group


## History of Macro viruses

- Possibility of Macro viruses known since 1980s.
- Microsoft created first as a 'proof of concept' which escaped after infecting Microsoft's own files.
- First virus was for Word, but since initial example appeared, Macro viruses created for other applications, e.g. Excel.

136

The possibility of Macro viruses has been known since at least the late 1980s. However, actual examples did not appear until well into the 1990s. In fact, Microsoft itself created the first known example of a Macro virus as a 'proof of concept' (the *Concept virus*) which escaped after infecting Microsoft's own files.

The first virus was for Microsoft Word, but since the initial example appeared, Macro viruses have been created for other applications, e.g. Excel.



Royal Holloway  
University of London

Information Security Group


## Macro viruses I

- Boot sector and parasitic viruses written in machine code.
- Macro viruses written in higher level language, e.g. an application Macro language, and reside in data files.
- Effectively an executable program embedded in a word processing document or other type of file

137

The last type of viruses we consider are the so called *Macro viruses*. Whereas Boot sector and parasitic viruses are written in machine code, Macro viruses are written in a higher level language, e.g. an application Macro language, and reside in data files created by the relevant application.

For example they can exist in word processing files, and would in that case be written in the Macro language for that word processor- usually some form of the Basic programming language and are often invoked automatically with no explicit user input. Once the macro is running, it can copy itself to other documents, delete files and cause other damage to the user's system.



Royal Holloway  
University of London

Information Security Group

## Macro viruses II


- Dangerous since bypass controls on executables, may be platform independent, and text documents are widely exchanged.
- Successive releases of Word provided increased protection against macro viruses.
- Antivirus product vendors have developed tools to detect and correct macro viruses.

138

Such viruses are very dangerous since they will bypass any controls on the distribution and use of executables, designed to counter parasitic viruses. Moreover, they may be platform independent. Most worrying of all, text documents are very widely exchanged in the commercial world, for legitimate business reasons (whereas the same is not generally true for executable files). Thus Macro viruses represent a serious long term corporate threat.

However, successive releases of Word have provided increased protection against macro viruses, and they are not the threat they once were.

Antivirus product vendors have developed tools to detect and correct macro viruses.



Royal Holloway  
University of London

Information Security Group

## Example – Concept virus

- The *Concept virus* infects the `NORMAL.DOT` file, the template for new documents created using Microsoft Word.
- Then infects all subsequently created `.DOC` files with the *Concept* macros.
- Editing an infected file results in infection of `NORMAL.DOT`.
- Very difficult to cure – at least if macros are enabled.

139

The *Concept virus* infects the `NORMAL.DOT` file, the template used when Microsoft Word creates a new document. Once loaded, the virus then infects all subsequently created `.DOC` files with the *Concept* macros.

If a running copy of Word is used to edit an infected file, then the virus macros are loaded into the running copy of Word, and when Word is shut down the `NORMAL.DOT` is infected.

Whilst it is possible to remove the virus from `NORMAL.DOT` (simply by deleting this file), it is very difficult to cure an infected file, and re-infection is very easy. All versions of Word from Word 97 onwards allow Macros to be disabled, preventing infections and gives a warning whenever a file containing macros is loaded.



Royal Holloway  
University of London


Information Security Group

## Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources

140

We conclude this part of the course by discussing worms.



Royal Holloway  
University of London

Information Security Group

## Operation of worms

- Worm:
  1. searches for vulnerable computers on a network;
  2. exploits a vulnerable computer;
  3. downloads copy of worm to compromised machine;
  4. each copy of the worm repeats the above steps ...
- Rate of infection can grow exponentially:
  - Code-Red worm infected 359,000 computers in <14 hours;
  - Slammer worm infected 75,000 computers in <10 minutes.


141

Put simply, a worm:

1. searches for vulnerable computers on a network, e.g. computers containing one or more specific software vulnerabilities;
2. exploits a vulnerable computer;
3. downloads a complete copy of the worm to the compromised computer;
4. each copy of the worm repeats the above steps ...

The rate of infection can grow exponentially:

- the Code-Red worm infected 359,000 computers in less than 14 hours – see: [http://www.caida.org/research/security/code-red/coderedv2\\_analysis.xml](http://www.caida.org/research/security/code-red/coderedv2_analysis.xml)
- the Slammer worm infected over 75,000 computers in less than 10 minutes – see: <http://www.cs.ucsd.edu/~savage/papers/IEEEESP03.pdf>



Royal Holloway  
University of London

Information Security Group

## The Morris worm

- The first Internet-wide worm:
  - developed by Robert Morris Jr (son of the former chief scientist of the NSA);
  - launched on 2 November 1988;
  - attacked machines running BSD Unix;
  - infected virtually every vulnerable Internet-connected computer;
  - caused denial-of-service on affected machines.


142

The Morris worm was the first Internet-wide worm. It:

- was developed by Robert Morris Jr (son of the former chief scientist of the NSA);
- was launched on 2 November 1988;
- attacked machines running BSD Unix;
- infected virtually every vulnerable Internet-connected computer;
- caused denial-of-service on affected machines.

For further details see (for example):

<http://homes.cerias.purdue.edu/~spaf/tech-reps/823.pdf>



Royal Holloway  
University of London

Information Security Group


## Morris worm – operation I

- The worm used three strategies to open a remote shell and install 'bootstrap' code on target machine:
  - exploit weak authentication to guess password;
  - exploit buffer overflow in finger daemon on target machine;
  - establish connection to SMTP port on target machine and exploit vulnerability in sendmail.

143

The worm used three strategies to open a remote shell and install special 'bootstrap' code on a target machine:

- exploit weak authentication to guess password;
- exploiting a buffer overflow in the finger daemon on the target machine;
- establish a connection to the SMTP port on the target machine and exploit a vulnerability in the sendmail program.



Royal Holloway  
University of London

Information Security Group

## Morris worm – operation II


- Bootstrap code on the infected machine contacted the source of the worm:
  - full code transferred to infected machine;
- Once fully installed, worm on infected machine:
  - identifies new hosts connected to infected machine (specified in configuration files such as `rhosts`, `.forward`, etc.);
  - attempts to infect new hosts.

144

The bootstrap code on the infected machine contacted the source of the worm, and the full code of the worm was then transferred to the infected machine.

Now fully installed, the worm on the infected machine:

- identifies new hosts connected to the infected machine (specified in configuration files such as `rhosts`, `.forward`, etc.);
- attempts to infect new hosts.



Royal Holloway  
University of London

Information Security Group

## Code-Red worm

- Exploits a buffer overflow in Microsoft's IIS web server:
  - defaces all pages requested from server;
  - attempts to spread the infection to other web servers;
  - infected servers launch coordinated DoS attack against specific web sites on specific date.

145


This worm exploits a buffer overflow in Microsoft's IIS web server. It:

- defaces all pages requested from server;
- attempts to spread the infection to other web servers.

Infected servers launch coordinated denial of service (DoS) attack against specific web sites on a specific date.

For further details see:

<http://www.cert.org/advisories/CA-2001-19.html>



Royal Holloway  
University of London

Information Security Group

## Nimda worm

- It affected systems running Microsoft Windows 95, 98, ME, NT, and 2000.
- Spread by exploiting buffer overflow vulnerability in IIS web server.
- Infected web server could infect browser.
- Spread via e-mail:
  - mail software that used IE (version 5.5 and earlier) would automatically run attachments containing the virus.

146

The Nimda worm affected systems running Microsoft Windows 95, 98, ME, NT, and 2000.

It spread by exploiting buffer overflow vulnerability in IIS web server.


Infected web server could infect browser.

It spread via e-mail:

- mail software that used IE (version 5.5 and earlier) would automatically run attachments containing the virus.

For further details see:

<http://www.cert.org/advisories/CA-2001-26.html>



Royal Holloway  
University of London

Information Security Group

## Blaster worm

- Initial infection exploits a buffer overflow vulnerability in Microsoft's DCOM RPC interface.
- Infected machine then retrieves msblast.exe from infecting machine and attempts to compromise further hosts.
- Attacker can execute arbitrary code as the local SYSTEM user.
- Caused large scale DoS.

147

Initial infection exploits a buffer overflow vulnerability in Microsoft's DCOM RPC interface.


Infected machine then retrieves msblast.exe from the infecting machine and attempts to compromise further hosts.

Attacker can execute arbitrary code as the local SYSTEM user.

It caused large scale denial of service.

For further details see:

<http://www.cert.org/advisories/CA-2003-20.html>



Royal Holloway  
University of London

Information Security Group

## Slammer worm


- Slammer exploited a buffer overflow vulnerability in Microsoft SQL Server:
  - allowed attacker to execute arbitrary code as the local SYSTEM user;
  - caused significant performance problems as worm tried to propagate to other machines.

148

The Slammer worm exploited a buffer overflow vulnerability in Microsoft SQL Server. It allowed the attacker to execute arbitrary code as the local SYSTEM user. It caused significant performance problems as the worm tried to propagate to other machines.

For further details see:

<http://www.cert.org/advisories/CA-2003-04.html>



Royal Holloway  
University of London

Information Security Group

## Preventing worm attacks

- Worms are more dangerous than viruses because they do not need human assistance:
  - Worms exploit vulnerabilities in programs remotely and automatically;
  - essential to patch applications known to be vulnerable to remote exploits.

149


Worms are more dangerous than viruses because they do not need human assistance:

- Worms exploit vulnerabilities in programs remotely and automatically;
- it is thus essential to patch applications that are known to be vulnerable to remote exploits.

Because they work without human assistance, they can spread very fast, much faster than a typical virus infection.

## Agenda

- Threats to security
- Managing security vulnerabilities
- Software vulnerabilities
- Other vulnerabilities
- Malware
- Viruses
- Worms
- Resources



Royal Holloway  
University of London

Information Security Group

## Further reading

- Robert Seacord, *Secure Coding in C and C++*.
- Jon Erickson, *Hacking: The Art of Exploitation*.
- Aleph One, *Smashing the stack for fun and profit*,  
<http://insecure.org/stf/smashstack.html>

151

In addition to the links given in the body of the presentation, the following reading is recommended.

- Robert Seacord, *Secure Coding in C and C++*.
- Jon Erickson, *Hacking: The Art of Exploitation*.
- Aleph One, *Smashing the stack for fun and profit*,  
<http://insecure.org/stf/smashstack.html>

Key links given in the body of the course material include:

- the US-CERT home page (<http://www.us-cert.gov/>), which has a wide variety of useful publications;
- the Microsoft Security Intelligence Report (<http://www.microsoft.com/sir>) is published twice a year, and contains a huge amount of useful information.

Finally, there is a huge amount of information available on the Internet about specific instances of malware and computer system vulnerabilities.