

Figure 21-8 A screened host is a firewall that is screened by a router.

more rules to the traffic and drops the denied packets. Then the traffic moves to the internal destination hosts. The screened host (the firewall) is the only device that receives traffic directly from the router. No traffic goes directly from the Internet, through the router, and to the internal network. The screened host is always part of this equation.

If the firewall is an application-based system, protection is provided at the network layer by the router through packet filtering, and at the application layer by the firewall. This arrangement offers a high degree of security, because for an attacker to be successful, she would have to compromise two systems.

What does the word “screening” mean in this context? As shown in Figure 21-8, the router is a screening device and the firewall is the screened host. This just means there is a layer that scans the traffic and gets rid of a lot of the “junk” before the traffic is directed toward the firewall. A screened host is different from a screened subnet, which is described next.

Screened Subnet A *screened-subnet* architecture adds another layer of security to the screened-host architecture. The external firewall screens the traffic entering the DMZ network. However, instead of the firewall then redirecting the traffic to the internal network, an interior firewall also filters the traffic. The use of these two physical firewalls creates a DMZ.

In an environment with only a screened host, if an attacker successfully breaks through the firewall, nothing lies in her way to prevent her from having full access to the internal network. In an environment using a screened subnet, the attacker would have to hack through another firewall to gain access. In this layered approach to security, the more layers provided, the better the protection. Figure 21-9 shows a simple example of a screened subnet.

The examples shown in the figures are simple in nature. Often, more complex networks and DMZs are implemented in real-world systems. Figures 21-10 and 21-11 show some other possible architectures of screened subnets and their configurations.

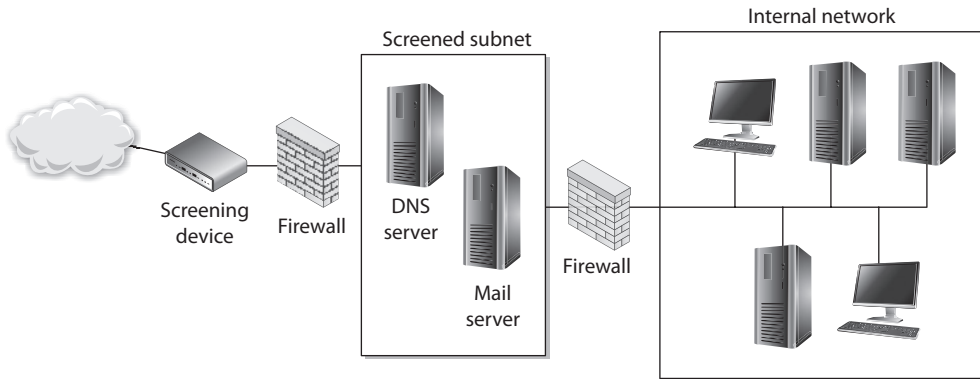


Figure 21-9 With a screened subnet, two firewalls are used to create a DMZ.

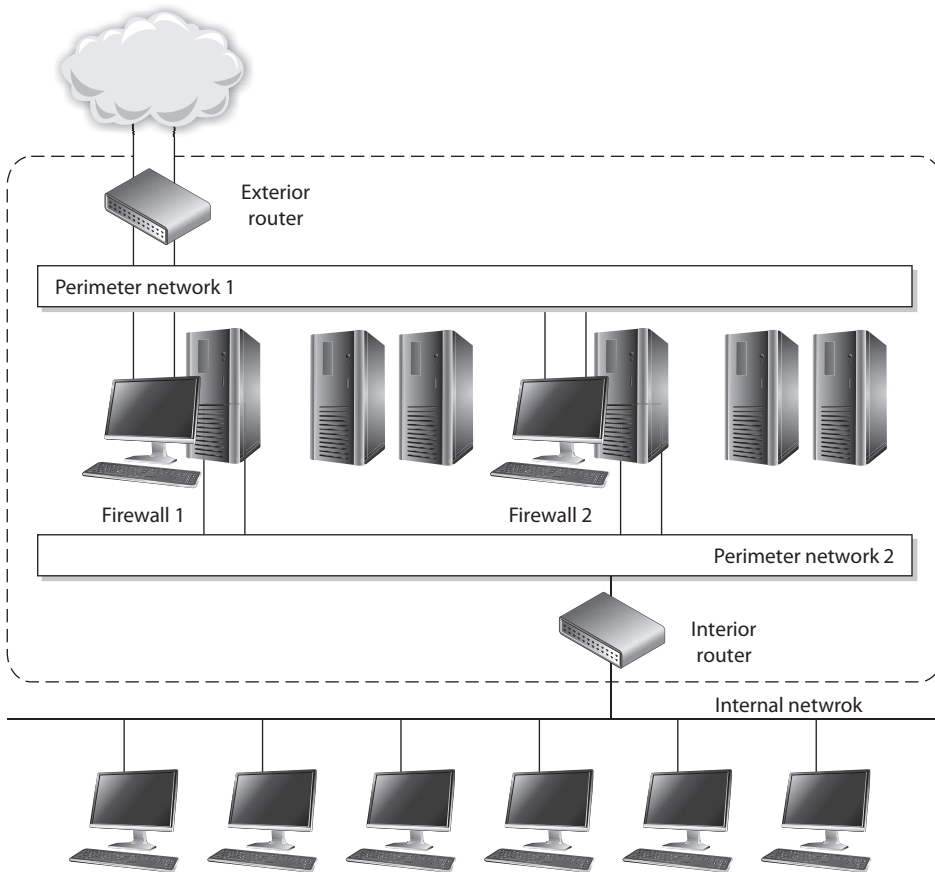


Figure 21-10 A screened subnet can have different networks within it and different firewalls that filter for specific threats.

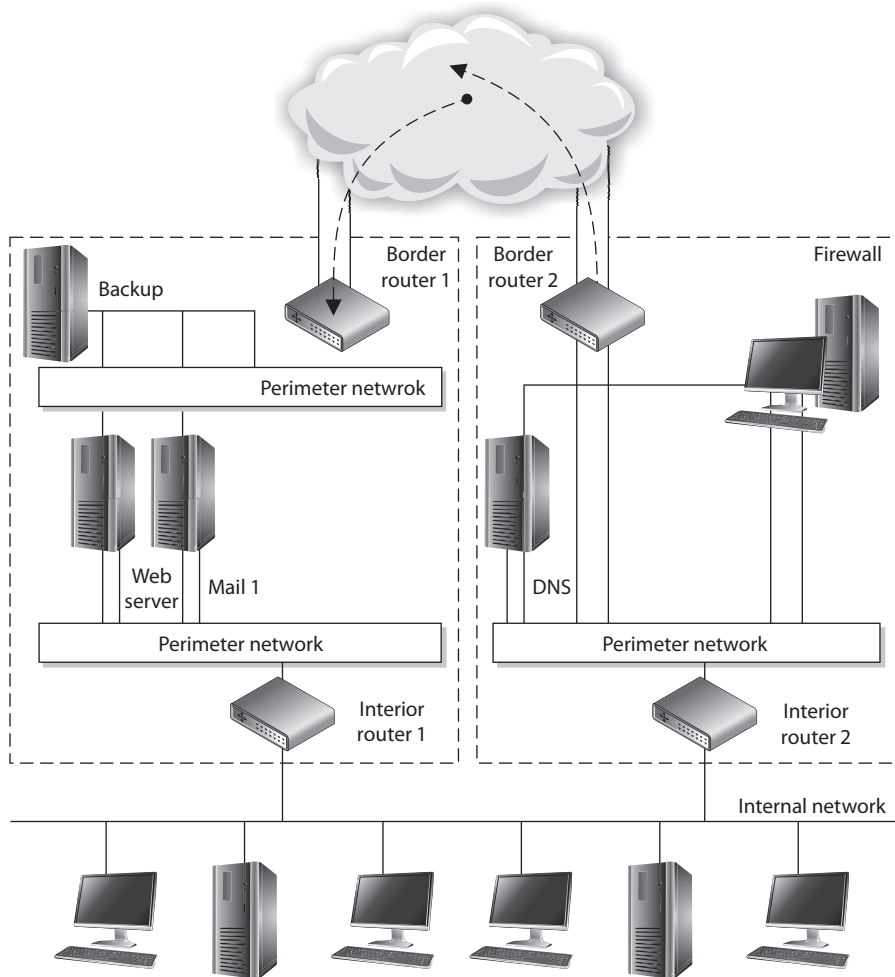


Figure 21-11 Some architectures have separate screened subnets with different server types in each.

The screened-subnet approach provides more protection than a stand-alone firewall or a screened-host firewall because three devices are working together and an attacker must compromise all three devices to gain access to the internal network. This architecture also sets up a DMZ between the two firewalls, which functions as a small network isolated among the trusted internal and untrusted external networks. The internal users usually

Firewall Architecture Characteristics

It is important to understand the following characteristics of these firewall architecture types:

Dual-homed:

- A single computer with separate NICs connected to each network.
- Used to divide an internal trusted network from an external untrusted network.
- Must disable a computer's forwarding and routing functionality so the two networks are truly segregated.

Screened host:

- A router filters (screens) traffic before it is passed to the firewall.

Screened subnet:

- An external router filters (screens) traffic before it enters the subnet. Traffic headed toward the internal network then goes through two firewalls.

have limited access to the servers within this area. Web, e-mail, and other public servers often are placed within the DMZ. Although this solution provides the highest security, it also is the most complex. Configuration and maintenance can prove to be difficult in this setup, and when new services need to be added, three systems may need to be reconfigured instead of just one.



TIP Sometimes a screened-host architecture is referred to as a single-tiered configuration and a screened subnet is referred to as a two-tiered configuration. If three firewalls create two separate DMZs, this may be called a three-tiered configuration.

Organizations used to deploy a piece of hardware for every network function needed (DNS, mail, routers, switches, storage, web), but today many of these items run within virtual machines on a smaller number of hardware machines. This reduces software and hardware costs and allows for more centralized administration, but these components still need to be protected from each other and external malicious entities. As an analogy, let's say that 15 years ago each person lived in their own house and a police officer was placed between each house so that the people in the houses could not attack each other. Then last year, many of these people moved in together so that now at least five

people live in the same physical house. These people still need to be protected from each other, so some of the police officers had to be moved inside the houses to enforce the laws and keep the peace. Analogously, virtual firewalls have “moved into” the virtualized environments to provide the necessary protection between virtualized entities.

As illustrated in Figure 21-12, a network can have a traditional physical firewall on the physical network and *virtual firewalls* within the individual virtual environments.

Virtual firewalls can provide bridge-type functionality in which individual traffic links are monitored between virtual machines, or they can be integrated within the hypervisor. The hypervisor is the software component that carries out virtual machine management and oversees guest system software execution. If the firewall is embedded within the hypervisor, then it can “see” and monitor all the activities taking place within the system.

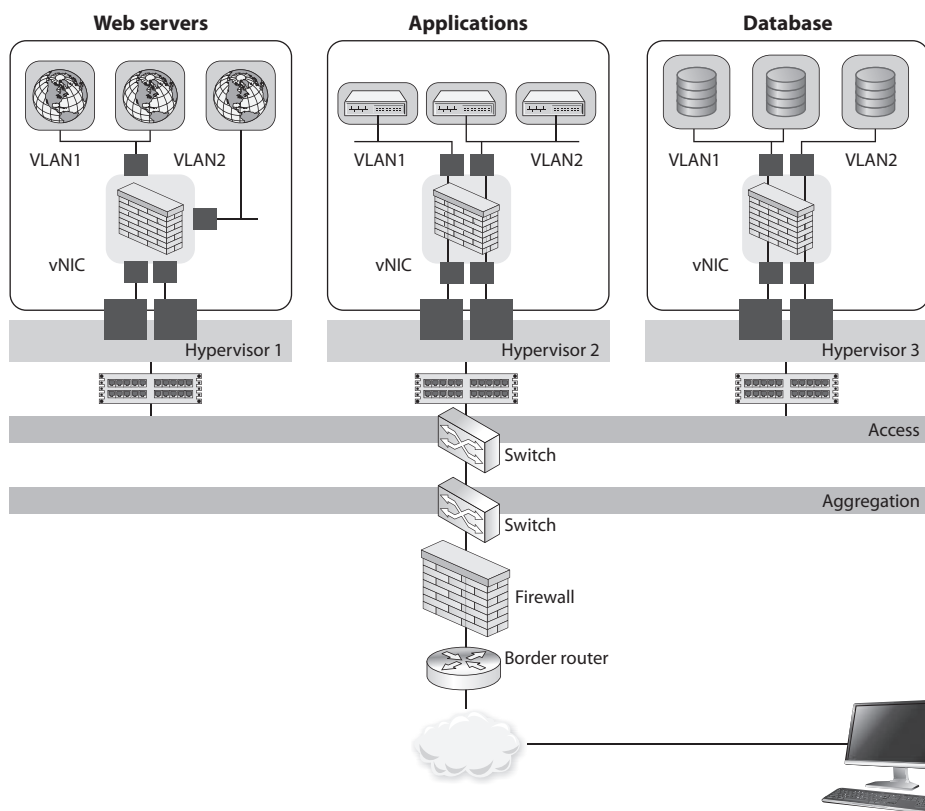


Figure 21-12 Virtual firewalls

Bastion Host

A system is considered a *bastion host* if it is a highly exposed device that is most likely to be targeted by attackers. The closer any system is to an untrusted network, such as the Internet, the more it is considered a target candidate since it has a smaller number of layers of protection guarding it. If a system is on the public side of a DMZ or is directly connected to an untrusted network, it is considered a bastion host; thus, it needs to be extremely locked down.

The system should have all unnecessary services disabled, unnecessary accounts disabled, unneeded ports closed, unused applications removed, unused subsystems and administrative tools removed, and so on. The attack surface of the system needs to be reduced, which means the number of potential vulnerabilities needs to be reduced as much as possible.

A bastion host does not have to be a firewall—the term just relates to the position of the system in relation to an untrusted environment and its threat of attack. Different systems can be considered bastion hosts (mail, web, DNS, etc.) if they are placed on the outer edges of networks.

The “Shoulds” of Firewalls

The default action of any firewall should be to implicitly deny any packets not explicitly allowed. This means that if no rule states that the packet can be accepted, that packet should be denied, no questions asked. Any packet entering the network that has a source address of an internal host should be denied. *Masquerading*, or *spoofing*, is a popular attacking trick in which the attacker modifies a packet header to have the source address of a host inside the network she wants to attack. This packet is spoofed and illegitimate. There is no reason a packet coming from the Internet should have an internal source network address, so the firewall should deny it. The same is true for outbound traffic. No traffic should be allowed to leave a network that does not have an internal source address. If this occurs, it means someone, or some program, on the internal network is spoofing traffic. This is how *zombies* work—the agents used in distributed DoS (DDoS) attacks. If packets are leaving a network with different source addresses, these packets are spoofed and the network is most likely being used as an accomplice in a DDoS attack.

Firewalls should reassemble fragmented packets before sending them on to their destination. In some types of attacks, the hackers alter the packets and make them seem to be something they are not. When a fragmented packet comes to a firewall, the firewall is seeing only part of the picture. It makes its best guess as to whether this piece of a packet is malicious or not. Because these fragments contain only a part of the full packet, the firewall is making a decision without having all the facts. Once all fragments are allowed through to a host computer, they can be reassembled into malicious packages that can cause a lot of damage. A firewall should accept each fragment, assemble the fragments

into a complete packet, and then make an access decision based on the whole packet. The drawback to this, however, is that firewalls that do reassemble packet fragments before allowing them to go on to their destination computer cause traffic delay and more overhead. It is up to the organization to decide whether this configuration is necessary and whether the added traffic delay is acceptable.

Many organizations choose to deny network entrance to packets that contain source routing information, which was mentioned earlier. Source routing means that the packet decides how to get to its destination, not the routers in between the source and destination computer. Source routing moves a packet throughout a network on a predetermined path. The sending computer must know about the topology of the network and how to route data properly. This is easier for the routers and connection mechanisms in between, because they do not need to make any decisions on how to route the packet. However, it can also pose a security risk. When a router receives a packet that contains source routing information, the router assumes the packet knows what needs to be done and passes the packet on. In some cases, not all filters may be applied to the packet, and a network administrator may want packets to be routed only through a certain path and not the route a particular packet dictates. To make sure none of this misrouting happens, many firewalls are configured to check for source routing information within the packet and deny it if it is present.

Firewalls are not effective “right out of the box.” You really need to understand the type of firewall being implemented and its configuration ramifications. For example, a firewall may have implied rules, which are used before the rules you configure. These implied rules might contradict your rules and override them. In this case, you may think that a certain traffic type is being restricted, but the firewall allows that type of traffic into your network by default.

The following list addresses some of the issues that need you need to understand as they pertain to firewalls:

- Most of the time a distributed approach needs to be used to control all network access points, which cannot happen through the use of just one firewall.
- Firewalls can present a potential bottleneck to the flow of traffic and a single point of failure threat.
- Some firewalls do not provide protection from malware and can be fooled by the more sophisticated attack types.
- Firewalls do not protect against sniffers or rogue wireless access points and provide little protection against insider attacks.

The role of firewalls is becoming more and more complex as they evolve and take on more functionality and responsibility. At times, this complexity works against security professionals because it requires them to understand and properly implement additional functionality. Without an understanding of the different types of firewalls and architectures available, many more security holes can be introduced, which lays out the welcome mat for attackers.

Intrusion Detection and Prevention Systems

The options for intrusion detection and prevention include host-based intrusion detection systems (HIDSs), network-based intrusion detection systems (NIDSs), and wireless intrusion detection systems (WIDSs). Each may operate in detection or prevention mode depending on the specific product and how it is employed. As a refresher, the main difference between an intrusion detection system (IDS) and an intrusion prevention system (IPS) is that an IDS only detects and reports suspected intrusions, while an IPS detects, reports, and stops suspected intrusions. How do they do this? There are two basic approaches: rule-based or anomaly-based.

Rule-Based IDS/IPS

Rule-based intrusion detection and prevention is the simplest and oldest technology. Essentially, we write rules (or subscribe to a service that writes them for us) and load those onto the system. The IDS/IPS monitors the environment in which it is placed, looking for anything that matches a rule. For example, suppose you have a signature for a particular piece of malware. You could create a rule that looks for any data that matches that signature and either raise an alert (IDS) or drop the data and generate the alert (IPS). Rule-based approaches are very effective when we know the telltale signs of an attack. But what if the attacker changes tools or procedures?

The main drawback of rule-based approaches to detecting attacks is that we need to have a rule that accurately captures the attack. This means someone got hacked, investigated the compromise, generated the rule, and shared it with the community. This process takes time and, until the rule is finalized and loaded, the system won't be effective against that specific attack. Of course, there's nothing stopping the adversary from slightly modifying tools or techniques to bypass your new rule either.

Anomaly-Based IDS/IPS

Anomaly-based intrusion detection and prevention uses a variety of approaches to detect things that don't look right. One basic approach is to observe the environment for some time to figure out what "normal" looks like. This is called the *training mode*. Once it has created a baseline of the environment, the IDS/IPS can be switched to *testing mode*, in which it compares observations to the baselines created earlier. Any observation that is significantly different generates an alert. For example, a particular workstation has a pattern of behavior during normal working hours and never sends more than, say, 10MB of data to external hosts during a regular day. One day, however, it sends out 100MB. That is pretty anomalous, so the IDS/IPS raises an alert (or blocks the traffic). But what if that was just the annual report being sent to the regulators?

The main challenge with anomaly-based approaches is that of *false positives*; that is, detecting intrusions when none happened. False positives can lead to fatigue and desensitizing the personnel who need to examine each of these alerts. Conversely, *false negatives* are events that the system incorrectly classifies as benign, delaying the response until the intrusion is detected through some other means. Obviously, both are bad outcomes.

EDR, NDR, and XDR

HIDS and antimalware features are increasingly being bundled into comprehensive *endpoint detection and response (EDR)* platforms. Similarly, NIDSs are evolving into *network detection and response (NDR)* products. These newer solutions do everything that HIDSs and NIDSs do, but also offer a host of other features such as combining rule-based and anomaly detection capabilities. *Extended detection and response (XDR)* platforms take this one step further by correlation of events across multiple sensors, both in the cloud and on premises, to get a more holistic view of what is going on in an environment.

Perhaps the most important step toward reducing errors is to baseline the system. *Baselining* is the process of establishing the normal patterns of behavior for a given network or system. Most of us think of baselining only in terms of anomaly-based IDSs because these typically have to go through a period of learning before they can determine what is anomalous. However, even rule-based IDSs should be configured in accordance with whatever is normal for an organization. There is no such thing as a one-size-fits-all set of IDS/IPS rules, though some *individual* rules may very well be applicable to all (e.g., detecting a known specimen of malware).



NOTE The term “perimeter” has lost some of its importance of late. While it remains an important concept in terms of security architecting, it can mislead some into imagining a wall separating us from the bad guys. A best practice is to assume the adversaries are already “inside the wire,” which downplays the importance of a perimeter in security operations.

Whitelisting and Blacklisting

One of the most effective ways to tune detection platforms like IDS/IPS is to develop lists of things that are definitely benign and those that are definitely malicious. The platform, then, just has to figure out the stuff that is not on either list. A *whitelist* (more inclusively called an *allow list*) is a set of known-good resources such as IP addresses, domain names, or applications. Conversely, a *blacklist* (also known as a *deny list*) is a set of known-bad resources. In a perfect world, you would only want to use whitelists, because nothing outside of them would ever be allowed in your environment. In reality, we end up using them in specific cases in which we have complete knowledge of the acceptable resources. For example, whitelisting applications that can execute on a computer is an effective control because users shouldn’t be installing arbitrary software on their own. Similarly, we can whitelist devices that are allowed to attach to our networks.

Things are different when we can’t know ahead of time all the allowable resources. For example, it is a very rare thing for an organization to be able to whitelist websites for every user. Instead, we would rely on blacklists of domain and IP addresses. The problem with blacklists is that the Internet is such a dynamic place that the only thing we can

be sure of is that our blacklist will always be incomplete. Still, blacklisting is better than nothing, so we should always try to use whitelists first, and then fall back on blacklists when we have no choice.

Antimalware Software

Traditional antimalware software uses signatures to detect malicious code. Signatures, sometimes referred to as fingerprints, are created by antimalware vendors. A *signature* is a set of code segments that a vendor has extracted from a malware sample. Similar to how our bodies have antibodies that identify and go after specific pathogens by matching segments of their genetic codes, antimalware software has an engine that scans files, e-mail messages, and other data passing through specific protocols and then compares them to its database of signatures. When there is a match, the antimalware software carries out whatever activities it is configured to do, which can be to quarantine the item, attempt to clean it (remove the malware), provide a warning message dialog box to the user, and/or log the event.

Signature-based detection (also called *fingerprint detection*) is a reasonably effective way to detect conventional malware, but it has a delayed response time to new threats. Once malware is detected in the wild, the antimalware vendor must study it, develop and test a new signature, release the signature, and all customers must download it. If the malicious code is just sending out silly pictures to all of your friends, this delay is not so critical. If the malicious software is a new variant of TrickBot (a versatile Trojan behind many ransomware attacks), this amount of delay can be devastating.

Since new malware is released daily, it is hard for the signature-based vendors to keep up. Another technique that almost all antimalware software products use is referred to as *heuristic detection*. This approach analyzes the overall structure of the malicious code, evaluates the coded instructions and logic functions, and looks at the type of data within the virus or worm. So, it collects a bunch of information about this piece of code and assesses the likelihood of it being malicious in nature. It has a type of “suspiciousness counter,” which is incremented as the program finds more potentially malicious attributes. Once a predefined threshold is met, the code is officially considered dangerous and the antimalware software jumps into action to protect the system. This allows antimalware software to detect unknown malware, instead of just relying on signatures.

As an analogy, let’s say Barney is the town cop who is employed to root out the bad guys and lock them up (quarantine). If Barney uses a signature method, he compares a stack of photographs of bad actors to each person he sees on the street. When he sees a match, he quickly throws the bad guy into his patrol car and drives off. By contrast, if he uses a heuristic method, he watches for suspicious activity. So if someone with a ski mask is standing outside a bank, Barney assesses the likelihood of this being a bank robber against it just being a cold guy in need of some cash.

Some antimalware products create a simulated environment, called a *virtual machine* or *sandbox*, and allow some of the logic within the suspected code to execute in the protected environment. This allows the antimalware software to see the code in question in action, which gives it more information as to whether or not it is malicious.



NOTE The virtual machine or sandbox is also sometimes referred to as an *emulation buffer*. They are all the same thing—a piece of memory that is segmented and protected so that if the code is malicious, the system is protected.

Reviewing information about a piece of code is called *static analysis*, while allowing a portion of the code to run in a virtual machine is called *dynamic analysis*. They are both considered heuristic detection methods.

Now, even though all of these approaches are sophisticated and effective, they are not 100 percent effective because malware writers are crafty. It is a continual cat-and-mouse game that is carried out every day. The antimalware industry comes out with a new way of detecting malware, and the very next week the malware writers have a way to get around this approach. This means that antimalware vendors have to continually increase the intelligence of their products and you have to buy a new version every year.

The next phase in the antimalware software evolution is referred to as behavior blockers. Antimalware software that carries out *behavior blocking* actually allows the suspicious code to execute within the operating system unprotected and watches its interactions with the operating system, looking for suspicious activities. The antimalware software watches for the following types of actions:

- Writing to startup files or the Run keys in the Windows registry
- Opening, deleting, or modifying files
- Scripting e-mail messages to send executable code
- Connecting to network shares or resources
- Modifying an executable logic
- Creating or modifying macros and scripts
- Formatting a hard drive or writing to the boot sector

If the antimalware program detects some of these potentially malicious activities, it can terminate the software and provide a message to the user. The newer-generation behavior blockers actually analyze sequences of these types of operations before determining the system is infected. (The first-generation behavior blockers only looked for individual actions, which resulted in a large number of false positives.) The newer-generation software can intercept a dangerous piece of code and not allow it to interact with other running processes. They can also detect rootkits. In addition, some of these antimalware programs can allow the system to roll back to a state before an infection took place so the damages inflicted can be “erased.”

While it sounds like behavior blockers might bring us our well-deserved bliss and utopia, one drawback is that the malicious code must actually execute in real time; otherwise, our systems can be damaged. This type of constant monitoring also requires a high level of system resources. We just can't seem to win.



EXAM TIP Heuristic detection and behavior blocking are considered proactive and can detect new malware, sometimes called “zero-day” attacks. Signature-based detection cannot detect new malware.

Most antimalware vendors use a blend of all of these technologies to provide as much protection as possible. The individual antimalware attack solutions are shown in Figure 21-13.



NOTE Another antimalware technique is referred to as *reputation-based protection*. An antimalware vendor collects data from many (or all) of its customers’ systems and mines that data to search for patterns to help identify good and bad files. Each file type is assigned a reputation metric value, indicating the probability of it being “good” or “bad.” These values are used by the antimalware software to help it identify “bad” (suspicious) files.

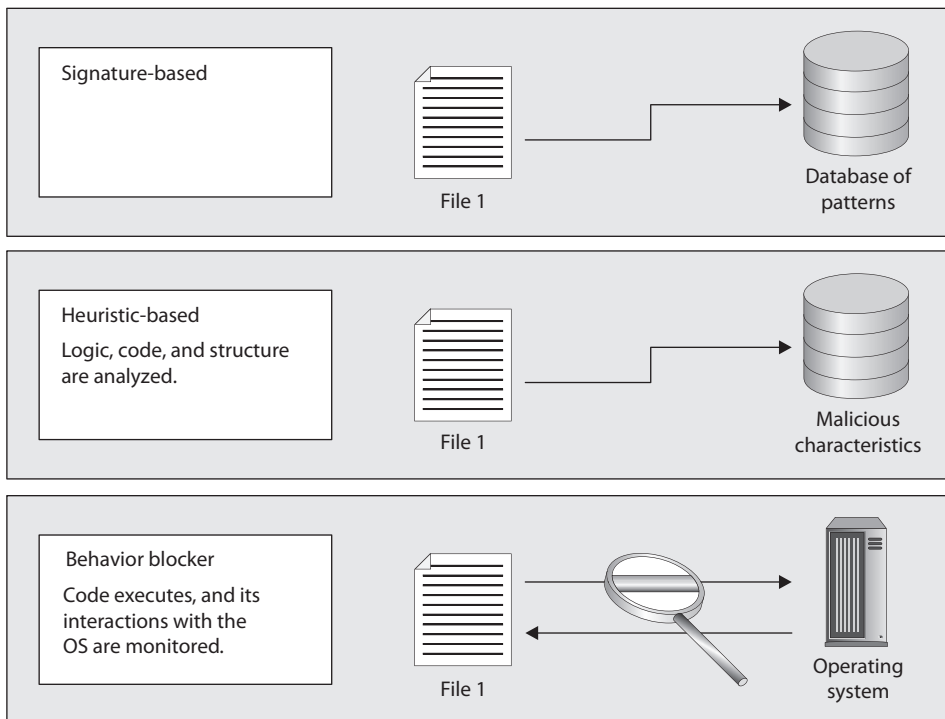


Figure 21-13 Antimalware vendors use various types of malware detection.

Detecting and protecting an enterprise from the long list of malware requires more than just rolling out antimalware software. Just as with other pieces of a security program, certain administrative, physical, and technical controls must be deployed and maintained.

The organization should either have a stand-alone antimalware policy or have one incorporated into an existing security policy. It should include standards outlining what type of antimalware software and antispyware software should be installed and how they should be configured.

Antimalware information and expected user behaviors should be integrated into the security-awareness program, along with who users should contact if they discover a virus. A standard should cover the do's and don'ts when it comes to malware, which are listed next:

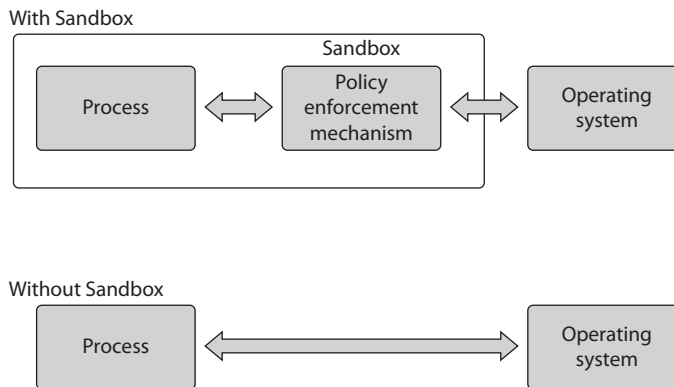
- Every workstation, server, and mobile device should have antimalware software installed.
- An automated way of updating malware signatures should be deployed on each device.
- Users should not be able to disable antimalware software.
- A preplanned malware eradication process should be developed and a contact person designated in case of an infection.
- All external disks (USB drives and so on) should be scanned automatically.
- Backup files should be scanned.
- Antimalware policies and procedures should be reviewed annually.
- Antimalware software should provide boot malware protection.
- Antimalware scanning should happen at a gateway and on each device.
- Virus scans should be automated and scheduled. Do not rely on manual scans.
- Critical systems should be physically protected so malicious software cannot be installed locally.

Since malware has cost organizations millions of dollars in operational costs and productivity hits, many have implemented antimalware solutions at network entry points. The scanning software can be integrated into a mail server, proxy server, or firewall. (The solutions are sometimes referred to as *virus walls*.) This software scans incoming traffic, looking for malware so it can be detected and stopped before entering the network. These products can scan Simple Mail Transport Protocol (SMTP), HTTP, FTP, and possibly other protocol types, but what is important to realize is that the product is only looking at one or two protocols and not *all* of the incoming traffic. This is the reason each server and workstation should also have antimalware software installed.

Sandboxing

A *sandbox* is an application execution environment that isolates the executing code from the operating system to prevent security violations. To the code, the sandbox looks just like the environment in which it would expect to run. For instance, when we sandbox

an application, it behaves as if it were communicating directly with the OS. In reality, it is interacting with another piece of software whose purpose is to ensure compliance with security policies. Another instance is that of software (such as helper objects) running in a web browser. The software acts as if it were communicating directly with the browser, but those interactions are mediated by a policy enforcer of some sort. The power of sandboxes is that they offer an additional layer of protection when running code that we are not certain is safe to execute.



Outsourced Security Services

Nearly all of the preventive and detective measures we've discussed in the preceding subsections can be outsourced to an external service provider. Why would we want to do that? Well, for starters, many small and midsize organizations lack the resources to provide a full team of experienced security professionals. We are experiencing workforce shortages that are not likely to be solved in the near term. This means that hiring, training, and retaining qualified personnel is not feasible in many cases. Instead, many organizations have turned to managed security services providers (MSSPs) for third-party provided security services.



EXAM TIP Outsourced security services are what (ISC)² refers to as *third-party provided security*.

MSSPs typically offer a variety of services ranging from point solutions to taking over the installation, operation, and maintenance of all technical (and some cases physical) security controls. (Sorry, you still have to provide policies and many administrative controls.) Your costs will vary depending on what you need but, in many cases, you'll get more than you could've afforded if you were to provide these services in-house. Still, there are some issues that you should consider before hiring an MSSP:

- **Requirements** Before you start interviewing potential MSSPs, make sure you know your requirements. You can outsource the day-to-day activities, but you can't outsource your responsibility to understand your own security needs.

- **Understanding** Does the MSSP understand your business processes? Are they asking the right questions to get there? If your MSSP doesn't know what it is that your organization does (and how), they will struggle to provide usable security. Likewise, you need to understand their qualifications and processes. Trust is a two-way street grounded on accurate information.
- **Reputation** It is hard to be a subpar service provider and not have customers complain about you. When choosing an MSSP, you need to devote some time to reading online reviews and asking other security professionals about their experiences with specific companies.
- **Costing** You may not be able to afford the deluxe version of the MSSP's services, so you will likely have to compromise and address only a subset of your requirements. When you have trimmed down your requirements, is it still more cost-effective to go with this provider? Should you go with another? Should you just do it yourself?
- **Liability** Any reasonable MSSP will put limits on their liability if your organization is breached. Read the fine print on the contract and consult your attorneys, particularly if you are in an industry that is regulated by the government.

Honeypots and Honeynets

A *honeypot* is a network device that is intended to be exploited by attackers, with the administrator's goal being to gain information on the attackers' tactics, techniques, and procedures (TTPs). Honeypots can work as early detection mechanisms, meaning that the network staff can be alerted that an intruder is attacking a honeypot system, and they can quickly go into action to make sure no production systems are vulnerable to that specific attack type. A honeypot usually sits in the screened subnet, or DMZ, and attempts to lure attackers to it instead of to actual production computers. Think of honeypots as marketing devices; they are designed to attract a segment of the market, get them to buy something, and keep them coming back. Meanwhile, threat analysts are keeping tabs on their adversaries' TTPs.

To make a honeypot system alluring to attackers, administrators may enable services and ports that are popular to exploit. Some honeypot systems *emulate* services, meaning the actual services are not running but software that acts like those services is available. Honeypot systems can get an attacker's attention by advertising themselves as easy targets to compromise. They are configured to look like the organization's regular systems so that attackers will be drawn to them like bears are to honey.

Another key to honeypot success is to provide the right kind of bait. When someone attacks your organization, what is it that they are after? Is it credit card information, patient files, intellectual property? Your honeypots should look like systems that would allow the attacker to access the assets for which they are searching. Once compromised, the directories and files containing this information must appear to be credible. It should also take a long time to extract the information, so that we maximize the contact time with our "guests."

A *honeynet* is an entire network that is meant to be compromised. While it may be tempting to describe honeynets as networks of honeypots, that description might be a bit misleading. Some honeynets are simply two or more honeypots used together. However, others are designed to ascertain a specific attacker's intent and dynamically spawn honeypots that are designed to be appealing to that particular attacker. As you can see, these very sophisticated honeynets are not networks of preexisting honeypots, but rather adaptive networks that interact with the adversaries to keep them engaged (and thus under observation) for as long as possible.



NOTE *Black holes* are sometimes confused with honeynets, when in reality they are almost the opposite of them. Black holes typically are routers with rules that silently drop specific (typically malicious) packets without notifying the source. They normally are used to render botnet and other known-bad traffic useless. Whereas honeypots and honeynets allow us to more closely observe our adversaries, black holes are meant to make them go away for us.

Wrapping up the honey collection, *honeyclients* are synthetic applications meant to allow an attacker to conduct a client-side attack while also allowing the threat analysts an opportunity to observe the TTPs being used by their adversaries. Honeyclients are particularly important in the honey family, because most of the successful attacks happen on the client side, and honeypots are not particularly well suited to track client-side attacks. Suppose you have a suspected phishing or spear phishing attack that you'd like to investigate. You could use a honeyclient to visit the link in the e-mail and pretend it is a real user. Instead of getting infected, however, the honeyclient safely catches all the attacks thrown at it and reports them to you. Since it is not really the web browser it is claiming to be, it is impervious to the attack and provides you with information about the actual tools the attacker is throwing at you. Honeyclients come in different flavors, with some being highly interactive (meaning a human has to operate them), while others involve low interaction (meaning their behavior is mostly or completely automated).

Organizations use these systems to identify, quantify, and qualify specific traffic types to help determine their danger levels. The systems can gather network traffic statistics and return them to a centralized location for better analysis. So as the systems are being attacked, they gather intelligence information that can help the network staff better understand what is taking place within their environment.

It should be clear from the foregoing that honeypots and honeynets are not defensive controls like firewalls and IDSs, but rather help us collect threat intelligence. To be effective, they must be closely monitored by a competent threat analyst. By themselves, honeypots and honeynets do not improve your security posture. However, they can give your threat intelligence team invaluable insights into your adversaries' methods and capabilities.

It is also important to make sure that the honeypot systems are not connected to production systems and do not provide any "jumping off" points for the attacker. There have been instances where companies improperly implemented honeypots and they were exploited by attackers, who were then able to move from those systems to the company's

internal systems. The honeypots need to be properly segmented from any other live systems on the network.

On a smaller scale, organizations may choose to implement *tar pits*, which are similar to honeypots in that they appear to be easy targets for exploitation. A tar pit can be configured to appear as a vulnerable service that attackers commonly attempt to exploit. Once the attackers start to send packets to this “service,” the connection to the victim system seems to be live and ongoing, but the response from the victim system is slow and the connection may time out. Most attacks and scanning activities take place through automated tools that require quick responses from their victim systems. If the victim systems do not reply or are very slow to reply, the automated tools may not be successful because the protocol connection times out.



NOTE Deploying honeypots and honeynets has potential liability issues. Be sure to consult your legal counsel before starting down this road.

Artificial Intelligence Tools

Artificial intelligence (AI) is a multidisciplinary field primarily associated with computer science, with influences from mathematics, cognitive psychology, philosophy, and linguistics (among others). At a high level, AI can be divided into two different approaches, as shown in Figure 21-14: symbolic and non-symbolic; the key difference is in how each represents knowledge. Both approaches are concerned with how knowledge is organized, how inference proceeds to support decision-making, and how the system learns.

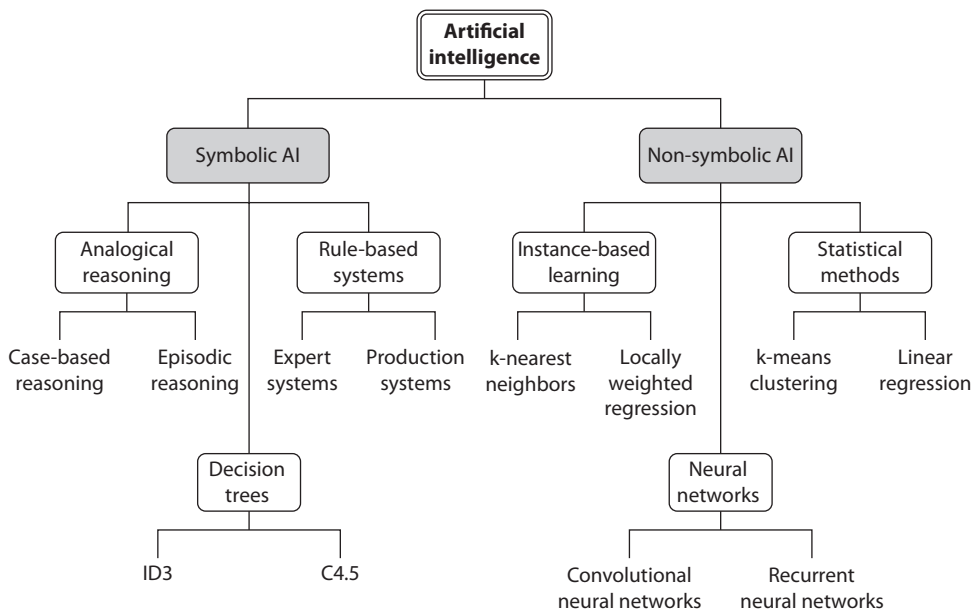


Figure 21-14 A partial taxonomy of artificial intelligence

In symbolic approaches to AI, system developers model real-world concepts, their relationships, and how they interact to solve a set of problems using a set of symbols (e.g., words or tokens). Symbolic AI requires considerable knowledge engineering of both the problem and solution domains, which makes it labor-intensive. However, it yields results that are inherently explainable to humans since the results are derived from human knowledge models in the first place. Symbolic AI systems include the expert systems that became prolific in the 1980s. These relied on extensive interviewing of subject matter experts and time-consuming encoding of their expertise in a series of conditional structures. Unsurprisingly, these early systems were unable to adapt or learn absent human intervention, which is a problem when we consider the number of exceptions that apply to almost all processes.

Another approach to AI departs from the use of symbolic representations of human knowledge and focuses instead on learning patterns in data for classifying objects, predicting future results, or clustering similar sets of data. These non-symbolic AI approaches are where many of the most recent advances have occurred, primarily in classification tasks such as image and voice recognition. In the current vernacular, these non-symbolic approaches are commonly called *machine learning (ML)* even though symbolic systems may also learn. As with symbolic approaches, non-symbolic ML systems also incorporate knowledge representations and reasoning. The knowledge representation is typically quantitative vectors (i.e., non-symbolic) with features from the dataset that describe the input (e.g., pixels from an image, frequencies from an audio file, word vectors, etc.).

Whereas symbolic AI requires considerable knowledge engineering, non-symbolic AI generally requires significant data acquisition and data curating, which can be labor-intensive even for domains where data is readily available. However, rather than having to program the knowledge, as in a symbolic system, the non-symbolic ML system acquires its knowledge in the form of numeric parameters (i.e., weights) through offline training with datasets with millions of examples. As training progresses, the ML model learns the correct parameters that minimize a cost function. That function typically deals with classifying some sample (helpful for finding malware) or making a prediction (allowing us to detect anomalies like spikes in outbound traffic).

Classification determines the class of a new sample based on what is known about previous samples. A common example of this is an algorithm called k-nearest neighbors (KNN), which is a supervised learning technique in which the nearest k neighbors influence the classification of the new point (e.g., if more than half of its k nearest neighbors are in one class, then the new point also belongs in that class). For cybersecurity, this is helpful when trying to determine whether a binary file is malware or detecting whether an e-mail is spam.

Prediction compares previous data samples and determines what the next sample(s) should be. If you have ever taken a statistics class, you may recall a type of analysis called *regression*, in which you try to determine the line (or curve) that most closely approximates a sequence of data points. We use the same approach to prediction in ML by learning from previous observations to determine where the next data point(s) should appear, which is useful for network flow analysis.

On the other hand, there is also unsupervised learning such as clustering, where we do not have a preconception of which classes (or even how many) exist; we determine where the samples naturally clump together. One of the most frequently used clustering algorithms is k-means clustering, in which new data points are added to one of the k clusters based on which one is closest to the new point. Clustering is useful for anomaly detection.

Finally, reinforcement learning tunes decision-making parameters toward choices that lead to positive outcomes in the environment. For example, one might have a security analyst provide feedback to an anomaly detector when it incorrectly classifies a malicious file or event (i.e., a false positive). This feedback adjusts the internal model's weights, so that its anomaly classification improves.

AI has shortcomings that you must consider before employing it. Neither symbolic nor non-symbolic AI approaches cope well with novel situations, and both require a human to re-engineer (symbolic) or retrain (non-symbolic) the algorithms. Symbolic, knowledge-engineered systems may contain underlying biases of the individual(s) who encode the system. Training data sets for non-symbolic approaches may contain biases that are not representative of the operational environment. These biases lead to either false positives or, worse, false negatives when the system is deployed. The best way forward is to combine both approaches, using each other's strengths to offset the other's weaknesses.

Logging and Monitoring

Logging and monitoring are two key activities performed by a SOC using the various tools we just discussed (and probably a few others). These two tasks go hand in hand, since you can't really monitor (at least not very effectively) if you are not logging and, conversely, logging makes little sense if you aren't monitoring. In the sections that follow, we first address how to collect and manage logs, and then discuss the ways in which you should be monitoring those logs (as well as other real-time data feeds).

Log Management

We discussed log reviews and how to prevent log tampering in Chapter 18. To understand how logs support day-to-day security operations, however, we need to take a step back and review why we might be logging system events in the first place. After all, if you don't have clear goals in mind, you will likely collect the wrong events at least some of the time.

Logging Requirements

Earlier in this chapter, we discussed cyberthreat intelligence and, in particular, the collection management framework (CMF). That section on the CMF is a great one to review when you're thinking about what your logging goals should be. After all, logs are data sources that can (and probably should) feed your threat intelligence. Just like intelligence requirements are meant to answer questions from decision-makers, logs should do the same for your SOC analysts. There should be specific questions your security team routinely asks, and those are the questions that should drive what you log and how. For

example, you may be concerned about data leaks of your sensitive research projects to overseas threat actors. What events from which system(s) would you need to log in order to monitor data egress? How often will you be checking logs (which determines how long you must retain them)? If you simply go with default logging settings, you may be ill informed when it comes to monitoring.

Log Standards

Another best practice is to standardize the format of your logs. If you are using a security information and event management (SIEM) system (which we'll discuss shortly), then that platform will take care of normalizing any logs you forward to it. Otherwise, you'll have to do it yourself using either the configuration settings on the system that's logging (if it allows multiple formats) or by using a data processing pipeline such as the open-source Logstash.



NOTE It is essential that you standardize the timestamps on all logs across your environment. If your organization is small, you can use local time; otherwise, we recommend you always use Coordinated Universal Time (UTC).

Something else to consider as you standardize your logs is who will be consuming them. Many SOCs leverage tools for automation, such as some of the AI techniques we discussed earlier. These automated systems may have their own set of requirements for formatting, frequency of updates, or log storage. You should ensure that your standards address the needs of all stakeholders (even non-human ones).

Logging Better

Finally, as with anything else you do in cybersecurity, you want to evaluate the effectiveness of your log management efforts and look for ways to sustain what you're doing well and improve the rest. Establishing and periodically evaluating metrics is an excellent approach to objectively determine opportunities for improvement. For example, how often do analysts lack information to classify an event because of incomplete logging? What logs, events, and fields are most commonly used when triaging alerts? Which are never needed? These questions will point to metrics, and the metrics, in turn, will tell you how well your logging supports your goals.

Security Information and Event Management

A *security information and event management (SIEM)* system is a software platform that aggregates security information (like asset inventories) and security events (which could become incidents) and presents them in a single, consistent, and cohesive manner. SIEMs collect data from a variety of sensors, perform pattern matching and correlation of events, generate alerts, and provide dashboards that allow analysts to see the state of the network. One of the best-known commercial solutions is Splunk, while on the open-source side the Elastic Stack (formerly known as the Elasticsearch-Logstash-Kibana, or ELK, stack) is very popular. It is worth noting that, technically, both of these systems are

data analytics platforms and not simply SIEMs. Their ability to ingest, index, store, and retrieve large volumes of data applies to a variety of purposes, from network provisioning to marketing to enterprise security.

Among the core characteristics of SIEMs is the ability to amass all relevant security data and present it to the security analyst in a way that makes sense. Before these devices became mainstream, security personnel had to individually monitor a variety of systems and manually piece together what all this information might mean. Most SIEMs now include features that group together information and events that seem to be related to each other (or “correlated” in the language of statistics). This allows the analyst to quickly determine the events that are most important or for which there is the most evidence.

SIEM correlations require a fair amount of fine-tuning. Most platforms, out of the box, come with settings that are probably good enough to get you started. You’ll have to let your SIEM tool run for a while (one week or longer) for it to start making sense of your environment and giving you meaningful alerts. Inevitably, you’ll find that your analysts are drowning in false positives (sadly, a very common problem with automated platforms) that consume their time and joy. This is where you start tuning your settings using things like whitelists and analyst ratings that will make the platform more accurate. You may also discover blind spots (that is, incidents that your SIEM did not pick up) due to insufficient logging or inadequate sensor placement, so you tune a bit there too.



NOTE SIEM fine-tuning should follow your established configuration management processes.

Security Orchestration, Automation, and Response

A tool that is becoming increasingly popular in SOC is the security orchestration, automation, and response (SOAR) platform. SOAR is an integrated system that enables more efficient security operations through automation of various workflows. The following are the three key components of a SOAR solution:

- **Orchestration** This refers to the integration and coordination of other security tools such as firewalls, IDS/IPS, and SIEM platforms. Orchestration enables automation.
- **Automation** SOAR platforms excel at automating cybersecurity playbooks and workflows, driving significant efficiency gains where those processes exist (or are created).
- **Response** Incident response workflows can involve dozens (or even hundreds) of distinct tasks. A SOAR platform can automatically handle many of those, freeing up the incident responders to work on what humans do best.

Egress Monitoring

A security practice that is oftentimes overlooked by smaller organizations is *egress monitoring*, which is keeping an eye on (and perhaps restricting) the information that is flowing *out* of our networks. Chapter 6 introduced data loss prevention (DLP), which is a very specific use case of this. Beyond DLP, we should be concerned about ensuring that our platforms are not being used to attack others and that our personnel are not communicating (knowingly or otherwise) with unsavory external parties.

A common approach to egress monitoring is to allow only certain hosts to communicate directly with external destinations. This allows us to focus our attention on a smaller set of computers that presumably would be running some sort of filtering software. A good example of this approach is the use of a web gateway, which effectively implements a man-in-the-middle “attack” on all of our organization’s web traffic. It is not uncommon to configure these devices to terminate (and thus decrypt) all HTTPS traffic and to do deep packet inspection (DPI) before allowing information to flow out of the network.

User and Entity Behavior Analytics

While most attacks historically are caused by external threat actors, we must not neglect to monitor the activities of users and entities within our organizations. Even if we never encounter a malicious insider, our users are oftentimes unwitting accomplices when they visit the wrong site, click the wrong link, or open the wrong attachment. *User and entity behavior analytics (UEBA)* is a set of processes that determines normal patterns of behavior so that abnormalities can be detected and investigated. For example, if a user hardly ever sends large amounts of data out to the Internet and then one day starts sending megabytes’ worth, that would trigger a UEBA alert. Maybe the transmission was perfectly legitimate, but perhaps it was the early part of a data loss incident.

UEBA can exist as a stand-alone product or as a feature in some other tool, such as an EDR or NDR platform. Either way, UEBA uses machine learning to predict future behaviors based on past observations, and statistical analyses to determine when a deviation from the norm is significant enough to raise an alert. As with any other type of solution that offers behavioral analytics, UEBA solutions are prone to false positives. This means that you would probably need to put some effort into fine-tuning a UEBA solution, even after its training period.



EXAM TIP UEBA is a good choice for detecting both malicious insiders and benign user accounts that have been taken over by a malicious actor.

Continuous Monitoring

NIST Special Publication 800-137, *Information Security Continuous Monitoring (ISCM) for Federal Information Systems and Organizations*, defines *information security continuous monitoring* as “maintaining ongoing awareness of information security, vulnerabilities, and threats to support organizational risk management decisions.” Think of ISCM as an

ongoing and structured verification of security controls. Are the existing controls still the right ones? Are they still effective? If not, why? These are some of the questions to which continuous monitoring provides answers. It is a critical part of the risk management framework we covered in Chapter 2.

There is a distinction here between logging, monitoring, and continuous monitoring. Your logging policies should be pretty expansive. Data storage is cheap and you want to capture as much data as you can in case you ever need it. Monitoring is more limited because it typically requires a human to personally do it, or at least to deal with the reports (such as SIEM alerts) that come out of it. You would, for example, monitor traffic on a certain port when it looks suspicious and then move on to monitoring something else when you determine that traffic is benign. Continuous monitoring is much more prescriptive. It is a deliberate, risk-based process to determine what gets monitored, how it is monitored, and what to do with the information you gather.

In the end, the whole point of continuous monitoring is to determine if the controls remain effective (in the face of changing threat and organizational environments) at reducing risk to acceptable levels. To do this, you need to carefully consider which metrics would allow you to say “yes” or “no” for each control. For example, suppose you are concerned about the risk of malware infections in your organization, so you implement antimalware controls. As part of continuous monitoring for those controls, you could measure the number of infections in some unit of time (day, week, month).

The metrics and measurements provide data that must be analyzed in order to make it actionable. Continuing our malware example, if your controls are effective, you would expect the number of infections to remain steady over time or (ideally) decrease. You would also want to consider other information in the analysis. For example, your malware infections could go up if your organization goes through a growth spurt and hires a bunch of new people, or the infections could go down during the holidays because many employees are taking vacation. The point is that the analysis is not just about understanding what is happening, but also why.

Finally, continuous monitoring involves deciding how to respond to the findings. If your organization’s malware infections have increased and you think this is related to the surge in new hires, should you provide additional security awareness training or replace the antimalware solution? Deciding what to do about controls that are no longer sufficiently effective must take into account risk, cost, and a host of other organizational issues.

Continuous monitoring is a deliberate process. You decide what information you need, then collect and analyze it at a set frequency, and then make business decisions with that information. Properly implemented, this process is a powerful tool in your prevention kit.

Chapter Review

Most of the time spent by the typical organization conducting security operations is devoted to emplacing and maintaining the preventive and detective measures, and then using those to log events and monitor the environment. Entire books have been written

on these topics, so in this chapter we just covered the essentials. A key takeaway is that tools alone will never be enough to give you the visibility you need to detect attacks; you need the integration of people, processes, and technology. We may have put a bit more focus on technology in this chapter, but we wanted to close it by highlighting the fact that well-trained people, working as a team and following existing processes, are essential components of security operations. This is particularly true when things go wrong and we need to respond to incidents, which we're about to cover in the next chapter.

Quick Review

- The security operations center (SOC) encompasses the people, processes, and technology that allow logging and monitoring of preventive controls, detection of security events, and incident response.
- Tier 1 security analysts spend most of their time monitoring security tools and other technology platforms for suspicious activity.
- Tier 2 security analysts dig deeper into the alerts, declare security incidents, and coordinate with incident responders and intelligence analysts to further investigate, contain, and eradicate the threats.
- Threat intelligence is evidence-based knowledge about an existing or emerging menace or hazard to assets that can be used to inform decisions regarding responses to that menace or hazard.
- Threat intelligence is commonly derived from three types of sources: threat data feeds, open-source intelligence (OSINT), and internal systems.
- Cyberthreat hunting is the practice of proactively looking for threat actors in your networks.
- Firewalls support and enforce the organization's network security policy by restricting access to one network from another network.
- Packet-filtering firewalls make access decisions based upon network-level protocol header values using access control lists (ACLs).
- Stateful firewalls add to the capabilities of packet-filtering firewalls by keeping track of the state of a connection between two endpoints.
- Proxy firewalls intercept and inspect messages before delivering them to the intended recipients.
- A next-generation firewall (NGFW) combines the attributes of the previously discussed firewalls, but adds a signature-based and/or behavioral analysis IPS engine, as well as cloud-based threat data sharing.
- Intrusion detection and prevention systems (IDS/IPS) can be categorized as either host-based (HIDS) or network-based (NIDS) and rule-based or anomaly-based.
- A whitelist is a set of known-good resources such as IP addresses, domain names, or applications. Conversely, a blacklist is a set of known-bad resources.

- Antimalware software is most effective when it is installed in every entry and end point and covered by a policy that delineates user training as well as software configuration and updating.
- A sandbox is an application execution environment that isolates the executing code from the operating system to prevent security violations.
- A honeypot is a network device that is intended to be exploited by attackers, with the administrator's goal being to gain information on the attackers' tactics, techniques, and procedures.
- A honeynet is an entire network that is meant to be compromised.
- Honeyclients are synthetic applications meant to allow an attacker to conduct a client-side attack while also allowing the security analysts an opportunity to observe the techniques being used by their adversaries.
- Machine learning (ML) systems acquire their knowledge in the form of numeric parameters (i.e., weights), through training with datasets consisting of millions of examples. In supervised learning, ML systems are told whether or not they made the right decision. In unsupervised training, they learn by observing an environment. Finally, in reinforcement learning they get feedback on their decisions from the environment.
- Effective logging requires a standard time zone for all timestamps.
- A security information and event management (SIEM) system is a software platform that aggregates security information (like asset inventories) and security events (which could become incidents) and presents them in a single, consistent, and cohesive manner.
- Security orchestration, automation, and response (SOAR) platforms are integrated systems that enable more efficient security operations through automation of various workflows.
- Egress monitoring is the process of scanning (and perhaps restricting) the information that is flowing out of our networks.
- User and entity behavior analytics (UEBA) is a set of processes that determines normal patterns of behavior so that abnormalities can be detected and investigated.
- Continuous monitoring allows organizations to maintain ongoing awareness of information security, vulnerabilities, and threats to support organizational risk management decisions.

Questions

Please remember that these questions are formatted and asked in a certain way for a reason. Keep in mind that the CISSP exam is asking questions at a conceptual level. Questions may not always have the perfect answer, and the candidate is advised against always looking for the perfect answer. Instead, the candidate should look for the best answer in the list.

Use the following scenario to answer Questions 1–3. The startup company at which you are the director of security is going through a huge growth spurt and the CEO has decided it's time to let you build out a security operations center (SOC). You already have two cybersecurity analysts (one is quite experienced), a brand-new security information and event management (SIEM) platform, and pretty good security processes in place.

1. The number of alerts on your SIEM is overwhelming your two analysts and many alerts go uninvestigated each day. How can you correct this?
 - A. Hire an intelligence analyst to help you focus your collection efforts.
 - B. Tune the SIEM platform to reduce false-positive alerts.
 - C. Establish a threat hunting program to find attackers before they trigger alerts.
 - D. Establish thresholds below which events will not generate alerts.
2. You hire an intelligence analyst and want her to start addressing intelligence requirements. Which of the following should be her first step?
 - A. Finding out what questions decision-makers need answered
 - B. Establishing a collection management framework
 - C. Identifying data sources
 - D. Subscribing to a threat data feed
3. Your SOC is maturing rapidly and you are ready to start a cyberthreat hunting program. Which of the following describes the crux of this effort?
 - A. Proving or negating hypotheses of threat actions based on threat intelligence
 - B. Neutralizing threat actors before they can breach your organization
 - C. Digging deeper into the alerts to determine if they constitute security incidents
 - D. Allowing hunters an opportunity to observe techniques used by their adversaries
4. A firewall that can only make decisions based on examining a single network layer header is called a
 - A. Stateful firewall
 - B. Screened host
 - C. Packet filter
 - D. Next-generation firewall
5. A firewall that understands the three-step handshake of a TCP connection is called a
 - A. Packet filter
 - B. Proxy firewall
 - C. Transport-layer proxy
 - D. Stateful firewall

6. What is the main challenge with anomaly-based approaches to intrusion detection and prevention?
 - A. False positives
 - B. Needing a rule that accurately captures the attack
 - C. Cost
 - D. Immaturity of the technology
7. Which of the following is an effective technique for tuning automated detection systems like IDS/IPS and SIEMs?
 - A. Access control lists
 - B. State tables
 - C. Whitelists
 - D. Supervised machine learning
8. Which of the following terms would describe a system designed to ascertain a specific attacker's intent and dynamically spawn multiple virtual devices that are designed to be appealing to that particular attacker?
 - A. Honeypot
 - B. Honeyclient
 - C. Honeyseeker
 - D. Honeynet
9. Which of the following is *not* a typical application of machine learning?
 - A. Classification
 - B. Prediction
 - C. Clustering
 - D. Knowledge engineering
10. Which of the following is *not* true about continuous monitoring?
 - A. It involves ad hoc processes that provide agility in responding to novel attacks.
 - B. Its main goal is to support organizational risk management.
 - C. It helps determine whether security controls remain effective.
 - D. It relies on carefully chosen metrics and measurements.

Answers

1. **B.** False positives are a very common problem with automated platforms like SIEMs, but they can be alleviated by fine-tuning the platform. An intelligence analyst could help a little bit but would clearly not be the best answer, while threat hunting would be a distractor for such a young SOC that still needs to get alerts

under control. Ignoring low-scoring alerts as a matter of policy would be a very dangerous move when dealing with stealthy attackers.

2. **A.** Threat intelligence is meant to help decision-makers choose what to do about a threat. It answers a question that these leaders may have. The CMF and data sources are all important, of course, but they are driven by the requirements that come out of leaders' questions. After the requirements are known, the intelligence analyst may (or may not) need to subscribe to a threat data feed.
3. **A.** The crux of threat hunting is to develop a hypothesis of adversarial action based on threat intelligence, and then to prove or negate the hypothesis. Inherent in this description are two factors: a) the adversary is already inside the network, and b) no alerts tipped off the defenders to the adversary's presence. These factors negate answers B and C. Answer D describes the purpose of a honeypot, not threat hunting.
4. **C.** Packet filtering is a firewall technology that makes access decisions based upon network-level protocol header values. The device that is carrying out packet-filtering processes is configured with access control lists (ACLs), which dictate the type of traffic that is allowed into and out of specific networks.
5. **D.** Stateful firewalls keep track of the state of a protocol connection, which means they understand the three-step handshake a TCP connection goes through (SYN, SYN/ACK, ACK).
6. **A.** The main challenge with anomaly-based approaches is that of false positives—detecting intrusions when none happened. These can lead to fatigue and desensitizing the personnel who need to examine each of these alerts. Despite this shortcoming, anomaly-based approaches are mature and cost-effective technologies that are differentiated from rule-based systems by not needing rules that accurately capture attacks.
7. **C.** One of the most effective ways to tune detection platforms like IDS/IPS is to develop lists of things that are definitely benign and those that are definitely malicious. The platform, then, just has to figure out the stuff that is not on either list. A whitelist (more inclusively called an allow list) is a set of known-good resources such as IP addresses, domain names, or applications.
8. **D.** Some honeynets are designed to ascertain a specific attacker's intent and dynamically spawn honeypots that are designed to be appealing to that particular attacker. These very sophisticated honeynets are not networks of preexisting honeypots, but rather adaptive networks that interact with the adversaries to keep them engaged (and thus under observation) for as long as possible.
9. **D.** Machine learning (ML), which is a non-symbolic approach to artificial intelligence (AI), is typically used for classification and prediction (using supervised or semi-supervised learning) as well as clustering (using unsupervised learning). Knowledge engineering is a requirement for symbolic forms for AI, such as expert systems, which are not ML in the common sense of the term.

- 10. A.** Continuous monitoring is a deliberate, data-driven process supporting organizational risk management. One of the key questions it answers is whether controls are still effective at mitigating risks. Continuous monitoring could potentially lead to a decision to implement specific ad hoc processes, but these would not really be part of continuous monitoring.

Security Incidents

This chapter presents the following:

- Incident management
- Incident response planning
- Investigations

It takes 20 years to build a reputation and few minutes of cyber-incident to ruin it.

—Stephane Nappo

No matter how talented your security staff may be, or how well everyone in your organization complies with your excellent security policies and procedures, or what cutting-edge technology you deploy, the sad truth is that the overwhelming odds are that your organization will experience a major compromise (if it hasn't already). What then? Having the means to manage incidents well can be just as important as anything else you do to secure your organization. In this chapter, we will cover incident management in general and then drill down into the details of incident response planning.

Although ISC² differentiates incident management and incident investigations, for many organizations, the latter is part of the former. This differentiation is useful to highlight the fact that some investigations involve suspects who may be our own colleagues. While many of us would enjoy the challenge of figuring out how an external threat actor managed to compromise our defenses, there is nothing fun about substantiating allegations that someone we work with did something wrong that caused losses to the organization. Still, as security professionals, we must be ready for whatever threats emerge and deal with the ensuing incidents well and rapidly.

Overview of Incident Management

There are many incident management models, but all share some basic characteristics. They all require that we identify the event, analyze it to determine the appropriate countermeasures, correct the problem(s), and, finally, take measures to keep the event from happening again. (ISC)² has broken out these four basic actions and prescribes seven phases in the incident management process: detection, response, mitigation, reporting, recovery, remediation, and lessons learned. Your own organization will have a unique approach, but it is helpful to baseline it off the industry standard.

Although we commonly use the terms “event” and “incident” interchangeably, there are subtle differences between the two. A *security event* is any occurrence that can be observed, verified, and documented. These events are not necessarily harmful. For example, a remote user login, changes to the Windows Registry on a host, and system reboots are all security events that could be benign or malicious depending on the context. A *security incident* is one or more related events that negatively affect the organization and/or impact its security posture. That remote login from our previous example could be a security incident if it was a malicious user logging in. We call reacting to these issues “incident response” (or “incident handling”) because something is negatively affecting the organization and causing a security breach.



EXAM TIP A security event is not necessarily a security violation, whereas a security incident is.

Many types of security incidents (malware, insider attacks, terrorist attacks, and so on) exist, and sometimes an incident is just human error. Indeed, many incident response individuals have received a frantic call in the middle of the night because a system is acting “weird.” The reasons could be that a deployed patch broke something, someone misconfigured a device, or the administrator just learned a new scripting language and rolled out some code that caused mayhem and confusion.

Many organizations are at a loss as to who to call or what to do right after they have been the victim of a cybercrime. Therefore, all organizations should have an *incident management policy (IMP)*. This document indicates the authorities and responsibilities regarding incident response for everyone in the organization. Though the IMP is frequently drafted by the CISO or someone on that person’s team, it is usually signed by whichever executive “owns” organizational policies. This could be the chief information officer (CIO), chief operations officer (COO), or chief human resources officer (CHRO). It is supported by an incident response plan that is documented and tested before an incident takes place. (More on this plan later.) The IMP should be developed with inputs from all stakeholders, not just the security department. Everyone needs to work together to make sure the policy covers all business, legal, regulatory, and security (and any other relevant) issues.

The IMP should be clear and concise. For example, it should indicate whether systems can be taken offline to try to save evidence or must continue functioning at the risk of destroying evidence. Each system and functionality should have a priority assigned to it. For instance, if a file server is infected, it should be removed from the network, but not shut down. However, if the mail server is infected, it should not be removed from the network or shut down, because of the priority the organization attributes to the mail server over the file server. Tradeoffs and decisions such as these have to be made when formulating the IMP, but it is better to think through these issues before the situation occurs, because better logic is usually possible before a crisis, when there’s less emotion and chaos.

Incident Management

Incident management includes proactive and reactive processes. Proactive measures need to be put into place so that incidents can be prevented or, failing that, detected quickly. Reactive measures need to be put into place so that detected incidents are dealt with properly.

Most organizations have only reactive management processes, which walk through how an incident should be handled. A more holistic approach is an incident management program that includes both proactive and reactive incident management processes, ensuring that triggers are monitored to make sure all incidents are actually uncovered. This commonly involves log aggregation, a security information and event management (SIEM) system, and user education. Having clear ways of dealing with incidents is not necessarily useful if you don't have a way to find out if incidents are indeed taking place.

All organizations should develop an *incident response team*, as mandated by the incident management policy, to respond to the large array of possible security incidents. The purpose of having an incident response (IR) team is to ensure that the organization has a designated group of people who are properly skilled, who follow a standard set of procedures, and who jump into action when a security incident takes place. The team should have proper reporting procedures established, be prompt in their reaction, work in coordination with law enforcement, and be recognized (and funded) by management as an important element of the overall security program. The team should consist of representatives from various business units, such as the legal department, HR, executive management, the communications department, physical/corporate security, IS security, and information technology.

There are three different types of incident response teams that an organization can choose to put into place. A *virtual* team is made up of experts who have other duties and assignments within the organization. It is called “virtual” because its members are not full-time incident responders but instead are called in as needed and may be physically remote. This type of team introduces a slower response time, and members must neglect their regular duties should an incident occur. However, a *permanent* team of folks who are dedicated strictly to incident response can be cost prohibitive to smaller organizations. The third type is a *hybrid* of the virtual and permanent models. Certain core members are permanently assigned to the team, whereas others are called in as needed.

Regardless of the type, the incident response team should have the following basic items available:

- A list of outside agencies and resources to contact or report to.
- An outline of roles and responsibilities.
- A call tree to contact these roles and outside entities.
- A list of computer or forensic experts to contact.
- A list of steps to take to secure and preserve evidence.

- A list of items that should be included in a report for management and potentially the courts.
- A description of how the different systems should be treated in this type of situation. (For example, remove the systems from both the Internet and the network and power them down.)

When a suspected crime is reported, the incident response team should follow a set of predetermined steps to ensure uniformity in their approach and that no steps are skipped. First, the IR team should investigate the report and determine whether an actual crime has been committed. If the team determines that a crime has been committed, they should inform senior management immediately. If the suspect is an employee, the team should contact a human resources representative right away. The sooner the IR team begins documenting events, the better. If someone is able to document the starting time of the crime, along with the employees and resources involved, that provides a good foundation for evidence. At this point, the organization must decide if it wants to conduct its own forensic investigation or call in experts. If experts are going to be called in, the system that was attacked should be left alone in order to try and preserve as much evidence of the attack as possible. If the organization decides to conduct its own forensic investigation, it must deal with many issues and address tricky elements. (Forensics will be discussed later in this chapter.)

Computer networks and business processes face many types of threats, each requiring a specialized type of recovery. However, an incident response team should draft and enforce a basic outline of how *all* incidents are to be handled. This is a much better approach than the way many organizations deal with these threats, which is usually in an ad hoc, reactive, and confusing manner. A clearly defined incident-handling process is more cost-effective, enables recovery to happen more quickly, and provides a uniform approach with certain expectation of its results.

Incident handling should be closely related to disaster recovery planning (covered in Chapter 23) and should be part of the organization's disaster recovery plan, usually as an appendix. Both are intended to react to some type of incident that requires a quick response so that the organization can return to normal operations. Incident handling is a recovery plan that responds to malicious technical threats. The primary goal of incident handling is to contain and mitigate any damage caused by an incident and to prevent any further damage. This is commonly done by detecting a problem, determining its cause, resolving the problem, and documenting the entire process.

Without an effective incident-handling program, individuals who have the best intentions can sometimes make the situation worse by damaging evidence, damaging systems, or spreading malicious code. Many times, the attacker booby-traps the compromised system to erase specific critical files if a user does something as simple as list the files in a directory. A compromised system can no longer be trusted because the internal commands listed in the path could be altered to perform unexpected activities. The system could now have a back door for the attacker to enter when he wants, or could

have a logic bomb silently waiting for a user to start snooping around, only to destroy any and all evidence.

Incident handling should also be closely linked to the organization's security training and awareness program to ensure that these types of mishaps do not take place. Past issues that the incident response team encountered can be used in future training sessions to help others learn what the organization is faced with and how to improve response processes.

Employees need to know how to report an incident. Therefore, the incident management policy should detail an escalation process so that employees understand when evidence of a crime should be reported to higher management, outside agencies, or law enforcement. The process must be centralized, easy to accomplish (or the employees won't bother), convenient, and welcomed. Some employees feel reluctant to report incidents because they are afraid they will get pulled into something they do not want to be involved with or accused of something they did not do. There is nothing like trying to do the right thing and getting hit with a big stick. Employees should feel comfortable about the process, and not feel intimidated by reporting suspicious activities.

The incident management policy should also dictate how employees should interact with external entities, such as the media, government, and law enforcement. This, in particular, is a complicated issue influenced by jurisdiction, the status and nature of the crime, and the nature of the evidence. Jurisdiction alone, for example, depends on the country, state, or federal agency that has control. Given the sensitive nature of public disclosure, communications should be handled by communications, human resources, or other appropriately trained individuals who are authorized to publicly discuss incidents. Public disclosure of a security incident can lead to two possible outcomes. If not handled correctly, it can compound the negative impact of an incident. For example, given today's information-driven society, denial and "no comment" may result in a backlash. On the other hand, if public disclosure is handled well, it can provide the organization with an opportunity to win back public trust. Some countries and jurisdictions either already have or are contemplating breach disclosure laws that require organizations to notify the public if a security breach involving personally identifiable information (PII) is even suspected. So, being open and forthright with third parties about security incidents often is beneficial to organizations.

A sound incident-handling program works with outside agencies and counterparts. The members of the team should be on the mailing list of the Computer Emergency Response Team (CERT) so they can keep up-to-date about new issues and can spot malicious events, hopefully before they get out of hand. CERT is a division of the Software Engineering Institute (SEI) that is responsible for monitoring and advising users and organizations about security preparation and security breaches.



NOTE Resources for CERT can be found at <https://www.cert.org/incident-management/>.

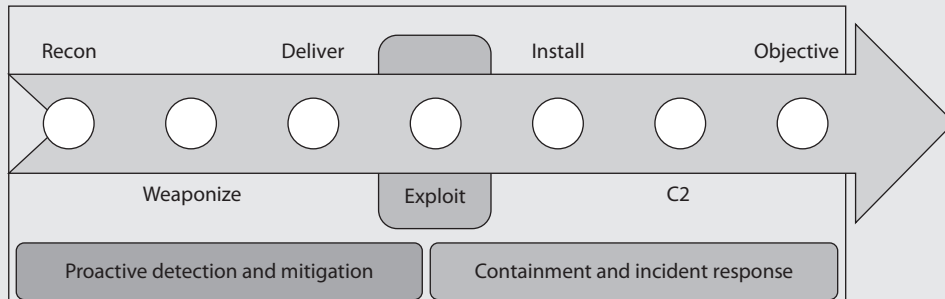
The Cyber Kill Chain

Even as we think about how best to manage incidents, it is helpful to consider a model that describes the stages attackers must complete to achieve their objectives. In their seminal 2011 white paper titled “Intelligence-Driven Computer Network Defense Informed by Analysis of Adversary Campaigns and Intrusion Kill Chains,” Eric Hutchins, Michael Cloppert, and Rohan Amin (employees of Lockheed Martin Corporation, publisher of the white paper) describe a seven-stage intrusion model that has become an industry standard known as the Cyber Kill Chain framework. The seven stages are described here:

- 1. Reconnaissance** The adversary has developed an interest in your organization as a target and begins a deliberate information-gathering effort to find vulnerabilities.
- 2. Weaponization** Armed with detailed-enough information, the adversary determines the best way into your systems and begins preparing and testing the weapons to be used against you.
- 3. Delivery** The cyber weapon is delivered into your system. In over 95 percent of the published cases, this delivery happens via e-mail.
- 4. Exploitation** The malicious software is executing on a CPU within your network. This may have launched when the target user clicked a link, opened an attachment, visited a website, or plugged in a USB thumb drive. It could also (in somewhat rare cases) be the result of a remote exploit. One way or another, the attacker’s software is now running in your systems.
- 5. Installation** Most malicious software is delivered in stages. First, there is the exploit that compromised the system in the prior step. Then, some other software is installed in the target system to ensure persistence, ideally with a good measure of stealth.
- 6. Command and Control (C2)** Once the first two stages of the software (exploit and persistence) have been executed, most malware will “phone home” to the attackers to let them know the attack was successful and to request updates and instructions.
- 7. Actions on Objectives** Finally, the malware is ready to do whatever it was designed to do. Perhaps the intent is to steal intellectual property and send it to an overseas server. Or perhaps this particular effort is an early phase in a grander attack, so the malware will pivot off the compromised system. Whatever the case, the attacker has won at this point.

As you can probably imagine, the earlier in the kill chain we identify the attack, the greater our odds are of preventing the adversaries from achieving their objectives.

This is a critical concept in this model: if you can thwart the attack before stage four (exploitation), you stand a better chance of winning. Early detection, then, is the key to success.



Incident response is the component of incident management that is executed when a security incident takes place. It starts with detecting the incident and eventually leads to the application of lessons learned during the response. Let's take a closer look at each of the steps in the incident response process.

Detection

The first and most important step in responding to an incident is to realize that you have a problem in the first place. The organization's incident response plan should have specific criteria and a process by which the security staff declares that an incident has occurred. The challenge, of course, is to separate the wheat from the chaff and zero in on the alerts or other indicators that truly represent an immediate danger to the organization.

Detection boils down to having a good sensor network implemented throughout your environment. There are three types of sensors: technical, human, and third-party. *Technical sensors* are, perhaps, the type most of us are used to dealing with. They are provided by the previously mentioned SIEM systems and the other types of systems introduced in Chapter 21: detection and response (EDR), network detection and response (NDR), and security orchestration, automation, and response (SOAR). *Human sensors* can be just as valuable if everyone in your organization has the security awareness to notice odd events and promptly report them to the right place. Many organizations use a special e-mail address to which anyone can send an e-mail report. *Third-party sensors* (technical or human) exist in other organizations. For example, maybe you have a really good relationship with your supply chain partners, and they will alert you to incidents in their environments that appear related to you. That third party could also be a government agency letting you know you've been hacked, which is never a good way to start your day, but is better than not knowing.

Despite this abundance of sensors, detecting incidents can be harder than it sounds, for a variety of reasons. First, sophisticated adversaries may use tools and techniques that you are unable to detect (at least at first). Even if the tools or techniques are known to you, they may very well be hiding under a mound of false positives in your SIEM system. In some (improperly tuned) systems, the ratio of false positives to true positives can be ten to one (or higher). This underscores the importance of tuning your sensors and analysis platforms to reduce the rate of false positives as much as possible.

Response

Having detected the incident, the next step is to respond by containing the damage that has been or is about to be done to your most critical assets. The goal of containment during the response phase is to prevent or reduce any further damage from this incident so that you can begin to mitigate and recover. Done properly, mitigation buys the IR team time for a proper investigation and determination of the incident's root cause. The response strategy should be based on the category of the attack (e.g., internal or external), the assets affected by the incident, and the criticality of those assets. So, what kind of mitigation strategy is best? Well, it depends.

When complete isolation or containment is not a viable solution, you may opt to use boundary devices to stop one system from infecting another. This involves temporarily changing firewall/filtering router rule configuration. Access control lists can be applied to minimize exposure. These response strategies indicate to the attacker that his attack has been noticed and countermeasures are being implemented. But what if, in order to perform a root cause analysis, you need to keep the affected system online and not let on that you've noticed the attack? In this situation, you might consider installing a honeynet or honeypot to provide an area that will contain the attacker but pose minimal risk to the organization. This decision should involve legal counsel and upper management because honeynets and honeypots can introduce liability issues, as discussed in Chapter 21. Once the incident has been contained, you need to figure out what just happened by putting the available pieces together.

This is the substage of analysis, where more data is gathered (audit logs, video captures, human accounts of activities, system activities) to try and figure out the root cause of the incident. The goals are to figure out who did this, how they did it, when they did it, and why. Management must be continually kept abreast of these activities because they will be making the big decisions on how this situation is to be handled.



EXAM TIP Watch out for the context in which the term “response” is used. It can refer to the entire seven-phase incident management process or to the second phase of it. In the second usage, you can think of it as *initial* response aimed at containment.

Mitigation

Having “stopped the bleeding” with the initial containment response, the next step is to determine how to properly mitigate the threat. Though the instinctive reaction may be to clean up the infected workstation or add rules to your firewalls and IDS/IPS,

this well-intentioned response could lead you on an endless game of whack-a-mole or, worse yet, blind you to the adversary's real objective. What do you know about the adversary? Who is it? What are they after? Is this tool and its use consistent with what you have already seen? Part of the mitigation stage is to figure out what information you need in order to restore security.

Once you have a hypothesis about the adversary's goals and plans, you can test it. If this particular actor is usually interested in PII on your high-net-worth clients but the incident you detected was on a (seemingly unrelated) host in the warehouse, was that an initial entry or pivot point? If so, then you may have caught the attacker before they worked their way further along the kill chain. But what if you got your attribution wrong? How could you test for that? This chain of questions, combined with quantifiable answers from your systems, forms the basis for an effective response. To quote the famous hockey player Wayne Gretzky, we should all "skate to where the puck is going to be, not where it has been."



NOTE It really takes a fairly mature threat intelligence capability to determine who is behind an attack (attribution), what are their typical tactics, techniques, and procedures (TTPs), and what might be their ultimate objective. If you do not have this capability, you may have no choice but to respond only to what you're detecting, without regard for what the adversary may actually be trying to do.

Once you are comfortable with your understanding of the facts of the incident, you move to eradicate the adversary from the affected systems. It is important to gather evidence before you recover systems and information. The reason is that, in many cases, you won't know that you will need legally admissible evidence until days, weeks, or even months after an incident. It pays, then, to treat each incident as if it will eventually end up in a court of justice.

Once all relevant evidence is captured, you can begin to fix all that was broken. The mitigation phase ends when you have affected systems that, while still isolated from the production networks, are free from adversarial control. For hosts that were compromised, the best practice is to simply reinstall the system from a gold master image and then restore data from the most recent backup that occurred prior to the attack. You may also have to roll back transactions and restore databases from backup systems. Once you are done, it is as if the incident never happened. Well, almost.



CAUTION An attacked or infected system should never be trusted, because you do not necessarily know all the changes that have taken place and the true extent of the damage. Some malicious code could still be hiding somewhere. Systems should be rebuilt to ensure that all of the potential bad mojo has been released by carrying out a proper exorcism.

Reporting

Though we discuss reporting at this point in order to remain consistent with the incident response process that (ISC)² identifies, incident reporting and documentation occurs at various stages in the response process. In many cases involving sophisticated attackers,

the IR team first learns of the incident because someone else reports it. Whether it is an internal user, an external client or partner, or even a government entity, this initial report becomes the starting point of the entire process. In more mundane cases, we become aware that something is amiss thanks to a vigilant member of the security staff or one of the sensors deployed to detect attacks. However we learn of the incident, this first report starts what should be a continuous process of documentation.

According to NIST Special Publication 800-61, Revision 2, *Computer Security Incident Handling Guide*, the following information should be reported for each incident:

- Summary of the incident
- Indicators
- Related incidents
- Actions taken
- Chain of custody for all evidence (if applicable)
- Impact assessment
- Identity and comments of incident handlers
- Next steps to be taken

Recovery

Once the incident is mitigated, you must turn your attention to the recovery phase, in which the aim is to restore full, trustworthy functionality to the organization. It is one thing to restore an individual affected device, which is what we do in mitigation, and another to restore the functionality of business processes, which is the goal of recovery. For example, suppose you have a web service that provides business-to-business (B2B) logistic processes for your organization and your partner organizations. The incident to which you're responding affected the database and, after several hours of work, you mitigated that system and are ready to put it back online. In this recovery stage, you would certify the system as trustworthy and then integrate it back into the web service, thus restoring the business capability.

It is important to note that the recovery phase is characterized by significant testing to ensure the following:

- The affected system is really trustworthy
- The affected system is properly configured to support whatever business processes it did previously
- No compromises exist in those processes

The third characteristic of this phase is assured by close monitoring of all related systems to ensure that the compromise did not persist. Doing this during off-peak hours helps ensure that, should we discover anything else malicious, the impact to the organization is reduced.

Remediation

It is not enough to put the pieces of Humpty Dumpty back together again. You also need to ensure that the attack is never again successful. In the remediation phase, which can (and should) run concurrently with the other phases, you decide which security controls (e.g., updates, configuration changes, firewall/IDS/IPS rules) need to be put in place or modified. There are two steps to this. First, you may have controls that are hastily put into effect because, even if they cause some other issues, their immediate benefit outweighs the risks. Later on, you should revisit those controls and decide which should be made permanent (i.e., through your change management process) and what others you may want to put in place.



NOTE For best results, the remediation phase should start right after detection and be conducted in parallel with the other phases.

Another aspect of remediation is the identification of indicators of attack (IOAs) that can be used in the future to detect this attack in real time (i.e., as it is happening) as well as indicators of compromise (IOCs), which tell you when an attack has been successful and your security has been compromised. Typical indicators of both attack and compromise include the following:

- Outbound traffic to a particular IP address or domain name
- Abnormal DNS query patterns
- Unusually large HTTP requests and/or responses
- DDoS traffic
- New registry entries (in Windows systems)

At the conclusion of the remediation phase, you have a high degree of confidence that this particular attack will never again be successful against your organization. Ideally, you should incorporate your IOAs and IOCs into the following lessons learned stage and share them with the community so that no other organization can be exploited in this manner. This kind of collaboration with partners (and even competitors) makes the adversary have to work harder.



EXAM TIP Mitigation, recovery, and remediation are conveniently arranged in alphabetical order. First you stop the threat, then you get back to business as usual, and then you ensure the threat is never again able to cause this incident.

Lessons Learned

Closure of an incident is determined by the nature or category of the incident, the desired incident response outcome (for example, business resumption or system restoration), and the team's success in determining the incident's source and root cause. Once you have

determined that the incident is closed, it is a good idea to have a team briefing that includes all groups affected by the incident to answer the following questions:

- What happened?
- What did we learn?
- How can we do it better next time?

The team should review the incident and how it was handled and carry out a postmortem analysis. The information that comes out of this meeting should indicate what needs to go into the incident response process and documentation, with the goal of continuous improvement. Instituting a formal process for the briefing provides the team with the ability to start collecting data that can be used to track its performance metrics.

Incident Response Planning

Incident management is implemented through two documents: the incident management policy (IMP) and the incident response plan (IRP). As discussed in the previous section, the IMP establishes authorities and responsibilities across the entire organization. The IMP identifies the IR lead for the organization and describes what every staff member is required to do with regard to incidents. For example, the IMP describes how employees are to report suspected incidents, to whom the report should be directed, and how quickly it should be done.

The IRP gets into the details of what should be done when responding to suspected incidents. The key sections of the IRP cover roles and responsibilities, incident classification, notifications, and operational tasks, all of which are described in the sections that follow. Normally, the IRP does not include detailed procedures for responding to specific incidents (e.g., phishing, data leak, ransomware), but establishes the framework within which all incidents will be addressed. Specific procedures are usually documented in *runbooks*, which are step-by-step scripts developed to deal with incidents that are either common enough or damaging enough to require this level of detailed documentation. Runbooks are described after the IRP sections.

Roles and Responsibilities

The group of individuals who make up the incident response team must have a variety of skills. They must also have a solid understanding of the systems affected by the incident, the system and application vulnerabilities, and the network and system configurations. Although formal education is important, real-world applied experience combined with proper training is key for these folks.

Many organizations divide their IR teams into two sub-teams. The first is the core team of incident responders, who come from the IT and security departments. These individuals are technologists who handle the routine incidents like restoring a workstation whose user inadvertently clicked the wrong link and caused self-infected damage. The second, or extended, team consists of individuals in other departments

who are activated for more complex incidents. The extended team includes attorneys, public relations specialists, and human resources staff (to name a few). The exact makeup of this extended team will vary based on the specifics of the incident, but the point is that these are individuals whose day-to-day duties don't involve IT or security, and yet they are essential to a good response. Table 22-1 shows some examples of the roles and responsibilities in these two teams.

Role	Responsibilities
Core IR Team	
Chief information security officer (CISO)	<ul style="list-style-type: none"> • Develops and maintains the IR plan • Communicates with senior organizational leadership • Directs security controls before and after incidents
Director of security operations	<ul style="list-style-type: none"> • Directs execution of the IR plan • Communicates with applicable law enforcement agencies • Declares security incidents
IR team lead	<ul style="list-style-type: none"> • Overall responsibility for the IR plan • Communicates with senior organizational leadership • Maintains repository of incident response lessons learned
Cybersecurity analyst	<ul style="list-style-type: none"> • Monitors and analyzes security events • Nominates events for escalation to security incidents • Performs additional analyses for IR team lead as required
IT support specialist	<ul style="list-style-type: none"> • Manages security platforms • Implements mitigation, recovery, and remediation measures as directed by the IR team lead
Threat intelligence analyst	<ul style="list-style-type: none"> • Provides intelligence products related to incidents • Maintains repository of incident facts to support future intelligence products
Extended IR Team	
Human resources manager	<ul style="list-style-type: none"> • Provides oversight for incident-related human resource requirements (e.g., employee relations, labor agreements)
Legal counsel	<ul style="list-style-type: none"> • Provides oversight for incident-related legal requirements (e.g., liability issues, requirement for law enforcement reporting/coordination) • Ensures evidence collected maintains its forensic value in the event the organization chooses to take legal action
Public relations	<ul style="list-style-type: none"> • Ensures communications during an incident protect the confidentiality of sensitive information • Prepares communications to stockholders and the press
Business unit lead	<ul style="list-style-type: none"> • Balances IR actions and business requirements • Ensures business unit support to the IR team

Table 22-1 IR Team Roles and Responsibilities

In addition to these two teams, most organizations rely on third parties when the requirements of the incident response exceed the organic capabilities of the organization. Unless you have an exceptionally well-resourced internal IR team, odds are that you'll need help at some point. The best course of action is to enter into an IR services agreement with a reputable provider *before* any incidents happen. By taking care of the contract and nondisclosure agreement (NDA) beforehand, the IR service provider will be able to jump right into action when time is of the essence. Another time-saving measure is to coordinate a familiarization visit with your IR provider. This will allow the folks who may one day come to your aid to become familiar with your organization, infrastructure, policies, and procedures. They will also get a chance to meet your staff, so everyone learns everyone else's capabilities and limitations.

Incident Classification

The IR team should have a way to quickly determine whether the response to an incident requires that everyone be activated 24/7 or the response can take place during regular business hours over the next couple of days. There is, obviously, a lot of middle ground between these two approaches, but the point is that incident classification criteria should be established, understood by the whole team, and periodically reviewed to ensure that it remains relevant and effective.

There is no one-size-fits-all approach to developing an incident classification framework, but regardless of how you go about it, you should consider three incident dimensions:

- **Impact** If you have a risk management program in place, classifying an incident according to impact should be pretty simple since you've already determined the losses as part of your risk calculations. All you have to do is establish the thresholds that differentiate a bad day from a terrible one.
- **Urgency** The urgency dimension speaks to how quickly the incident needs to be mitigated. For example, an ongoing exfiltration of sensitive data needs to be dealt with immediately, whereas a scenario where a user caused self-infected damage with a bitcoin mining browser extension shouldn't require IR team members to get out of bed in the middle of the night.
- **Type** This dimension helps the team identify the resources that need to be notified and mobilized to deal with the incident. The team that handles the data exfiltration incident mentioned earlier is probably going to be different than the one that handles the infected browser.

Not all organizations explicitly call out each of these dimensions (and some organizations have more dimensions), but it is important to at least consider them. The simplest approach to incident classification simply uses severity and assigns various levels to this parameter depending on whether certain conditions are met. Table 22-2 shows a simple classification matrix for a small to medium-sized organization.

Severity	Criteria	Initial Response Time
Severity 1 (critical)	<ul style="list-style-type: none"> Confirmed incident compromising mission-critical systems Active exfiltration, alteration, or destruction of sensitive data Incident requiring notification to government regulators Life-threatening ongoing physical situation (e.g., suspicious package on site, unauthorized/hostile person, credible threat) 	1 hour
Severity 2 (high)	<ul style="list-style-type: none"> Confirmed incident compromising systems that are not mission-critical Active exfiltration of non-sensitive data Time-sensitive investigation of employees Non-life-threatening but serious, ongoing physical situation (e.g., unauthorized person, theft of property) 	4 hours
Severity 3 (moderate)	<ul style="list-style-type: none"> Possible incident affecting any systems Security policy violations Long-term employee investigations requiring extensive collection and analysis Non-life-threatening past physical situation (e.g., sensitive area left unsecured overnight) 	48 hours

Table 22-2 Sample Incident Classification Matrix

The main advantage of formally classifying incidents is that it allows the preauthorized commitment of resources within specific timeframes. For example, if one of your SOC tier 2 analysts declares a severity 1 (critical) incident, she could be authorized to call the external IR service provider, committing the organization to pay the corresponding fees. There would be no need to get a hold of the CISO and get permission.

Notifications

Another benefit of classifying incidents is that it lets the IR team know who they need to inform and how frequently. Obviously, we don't want to call the CISO at home whenever an employee violates a security policy. On the other hand, we really don't want the CEO to find out the organization had an incident from reading the morning news. Keeping the right decision-makers informed at the right cadence enables everyone to do their jobs well, engenders trust, and leads to unified external messaging.

Table 22-3 shows an example notification matrix that builds on the classification shown previously in Table 22-2.

Notifications to external parties such as customers, partners, government regulators, and the press should be handled by communications professionals and not by the cybersecurity staff. The technical members of the IR team provide the facts to these communicators, who then craft messages (in coordination with the legal and marketing teams) that do not make things worse for the organization either legally or reputationally. Properly handled, IR communications can help improve trust and loyalty to the

Stakeholder	Severity Level	Notification
Executive leaders	S1	Immediate via e-mail and phone
	S2	On the next daily operational report
	S3	None
CISO	S1	Immediate via e-mail and phone
	S2	Within 4 hours via e-mail and phone
	S3	On the next daily operational report
Affected business units	S1	Immediate via e-mail and phone
	S2	Within 4 hours via e-mail
	S3	On the next daily operational report
Affected customers/partners	S1	Within 8 hours via e-mail
	S2	Within 72 hours via e-mail
	S3	None

Table 22-3 Sample Incident Notification Matrix

organization. Improperly handled, however, these notifications (or the lack thereof) can ruin (and have ruined) organizations.

Operational Tasks

Keeping stakeholders informed is just one of the many tasks involved in incident response. Just like any other complex endeavor, we should leverage structured approaches to ensure that all required tasks are performed, and that they are done consistently and in the right order. Now, of course, different types of incidents require different procedures. Responding to a ransomware attack requires different procedures than the procedures for responding to a malicious insider trying to steal company secrets. Still, all incidents follow a very similar pattern at a high level. We already saw this in the discussion of the seven phases in the incident management process that you need to know for the CISSP exam, which apply to all incidents.

Many organizations deal with the need for completeness and consistency in IR by spelling out operational tasks in the IRP, sometimes with a field next to each task to indicate when the task was completed. The IR team lead can then just walk down this list to ensure the right things are being done in the right order. Table 22-4 shows a sample operational tasks checklist.

Table 22-4 is not meant to be all-inclusive but it does capture the most common tasks that apply to every IR in most organizations. As mentioned earlier, different types of incidents require different approaches. While the task list should be general enough to accommodate these specialized procedures, we also want to keep it specific enough to serve as an overall execution plan.

Operational Task	Date/Time Completed
Pre-Execution	
Identify assets affected	
Obtain access (physical and logical) to all affected assets	
Determine forensic evidence requirements	
Review compliance requirements (e.g., GDPR, HIPAA, PCI DSS)	
Initiate communications plan	
Response	
Perform immediate actions to mitigate the impact of the incident	
Validate detection mechanisms	
Request relevant intelligence from threat intelligence team	
Gather and preserve incident-related data (e.g., PCAP, log files)	
Develop an initial timeline of incident-related activity	
Develop mitigation plan based on initial assessment	
Mitigation	
Verify availability of backup/redundant system (if mission-critical system was compromised)	
Activate backup/redundant systems for continuity of operations (if mission-critical system was compromised)	
Isolate affected assets	
Collect forensic evidence from compromised systems (if applicable)	
Remove active threat mechanisms to limit further activity	
Initiate focused monitoring of the environment for additional activity	
Recovery	
Restore affected systems' known-good backups or gold masters	
Validate additional controls on restored systems prevent reoccurrence	
Reconnect restored systems to production networks	
Verify no additional threat activity exists on restored systems	
Remediation	
Finalize root cause, threat mechanisms, and incident timeline	
Identify IOCs and IOAs	
Initiate change management processes to prevent reoccurrence	
Implement preventive and detective controls to prevent reoccurrence	

Table 22-4 Sample Operational Tasks List

Runbooks

When we need specialized procedures, particularly when we expect a certain type of incident to happen more than once, we want to document those procedures to ensure we don't keep reinventing the wheel every time a threat actor gets into our systems. A *runbook* is a collection of procedures that the IR team will follow for specific types of incidents. Think of a runbook as a cookbook. If you feel like having a bean casserole for dinner, you open your cookbook and look up that recipe. It'll tell you what ingredients you need and what the step-by-step procedure is to make it. Similarly, a runbook has tabs for the most likely and/or most dangerous incidents you may encounter. Once the incident is declared by the SOC (or whoever is authorized to declare an incident has occurred), the IR team lead opens the runbook and looks up the type of incident that was declared. The runbook specifies what resources are needed (e.g., specific roles and tools) and how to apply them.

When developing runbooks, you have to be careful that the documentation doesn't take more time and resources to develop than you would end up investing in responding to that incident type. As with any other control, the cost of a runbook cannot exceed the cost of doing nothing (and figuring things out on the fly). For that reason, most organizations focus their runbooks on incidents that require complex responses and those that are particularly sensitive. Other incidents can be (and usually are) added to the runbook, but those additions are deliberate decisions of the SOC manager based on the needs of the organization. For example, if an organization experiences high turnover rates, it might be helpful for new staff to have a more comprehensive runbook to which they can turn.

Another aspect to consider is that runbooks are only good if they are correct, complete, and up to date. Even if you do a great job when you first write runbooks, you'll have to invest time periodically in keeping them updated. For best results, incorporate runbooks into your change management program so that, whenever an organizational change is made, the change advisory board (CAB) asks the question: does this require an update to the IR runbooks?

Investigations

Whatever type of security incident we're facing, we should treat the systems and facilities that it affects as potential crime scenes. The reason is that what may at first appear to have been a hardware failure, a software defect, or an accidental fire may have in fact been caused by a malicious actor targeting the organization. Even acts of nature like storms or earthquakes may provide opportunities for adversaries to victimize us. Because we are never (initially) quite sure whether an incident may have a criminal element, we should treat all incidents as if they do (until proven otherwise).

Since computer crimes are only increasing and will never really go away, it is important that all security professionals understand how computer investigations should be carried out. This includes understanding legal requirements for specific situations, the chain of custody for evidence, what type of evidence is admissible in court, incident response procedures, and escalation processes.

Cops or No Cops?

Management needs to make the decision as to whether law enforcement should be called during an incident response. The following are some of the issues to understand if law enforcement is brought in:

- You may not have a choice in certain cases (e.g., cases involving national security, child pornography, etc.).
- Law enforcement agencies bring significant investigative capability.
- The organization may lose control over where the investigation leads once law enforcement is involved.
- Secrecy of compromise is not promised; it could become part of public record.
- Evidence will be collected and may not be available for a long period of time.

Successfully prosecuting a crime requires solid evidence. Computer forensics is the art of retrieving this evidence and preserving it in the proper ways to make it admissible in court. Without proper computer forensics, few computer crimes could ever be properly and successfully presented in court. The most common reasons evidence is deemed inadmissible in court are lack of qualified staff handling it, lack of established procedures, poorly written policy, or a broken chain of custody.

When a potential computer crime takes place, it is critical that the investigation steps are carried out properly to ensure that the evidence will be admissible to the court (if the matter goes that far) and can stand up under the cross-examination and scrutiny that will take place. As a security professional, you should understand that an investigation is not just about potential evidence on a disk drive. The context matters during an investigation, including the people, network, connected internal and external systems, applicable laws and regulations, management's stance on how the investigation is to be carried out, and the skill set of whoever is carrying out the investigation. Messing up just one of these components could make your case inadmissible or at least damage it if it is brought to court.

Motive, Opportunity, and Means

Today's computer criminals are similar to their traditional counterparts. To understand the "why" in crime, it is necessary to understand the motive, opportunity, and means—or MOM. This is the same strategy used to determine the suspects in a traditional, non-computer crime.

Motive is the "who" and "why" of a crime. The motive may be induced by either internal or external conditions. A person may be driven by the excitement, challenge, and adrenaline rush of committing a crime, which would be an internal condition. Examples of external conditions might include financial trouble, a sick family member, or other dire straits. Understanding the motive for a crime is an important piece in figuring out who

would engage in such an activity. For example, financially motivated attackers such as those behind ransomware want to get your money. In the case of ransomware purveyors, they realize that if they don't decrypt a victim's data after payment of the ransom, the word will get out and no other victims will pay the ransom. For this reason, most modern ransomware actors reliably turn over decryption keys upon payment. Some ransomware gangs even go the extra mile and set up customer service operations to help victims with payment and decryption issues.

Opportunity is the “where” and “when” of a crime. Opportunities usually arise when certain vulnerabilities or weaknesses are present. If an organization does not regularly patch systems (particularly public-facing ones), attackers have all types of opportunities within that network. If an organization does not perform access control, auditing, and supervision, employees may have many opportunities to embezzle funds and defraud the organization. Once a crime fighter finds out why a person would want to commit a crime (motive), she will look at what could allow the criminal to be successful (opportunity).

Means pertains to the abilities a criminal would need to be successful. Suppose a crime fighter was asked to investigate a case of fraud facilitated by a subtle but complex modification made to a software system within a financial institution. If the suspects were three people and two of them just had general computer knowledge, but the third one was a programmer and system analyst, the crime fighter would realize that this person is much likelier to have the means to commit this crime than the other two individuals.

Computer Criminal Behavior

Like traditional criminals, computer criminals have a specific *modus operandi* (MO, pronounced “em-oh”). In other words, each criminal typically uses a distinct method of operation to carry out their crime, and that method can be used to help identify them. The difference with computer crimes is that the investigator, obviously, must have knowledge of technology. For example, the MO of a particular computer criminal may include the use of specific tools or targeting specific systems or networks. The method usually involves repetitive signature behaviors, such as sending e-mail messages or programming syntax. Knowledge of the criminal's MO and signature behaviors can be useful throughout the investigative process. Law enforcement can use the information to identify other offenses by the same criminal, for example. The MO and signature behaviors can also provide information that is useful during interviews (conducted by authorized staff members or law enforcement agencies) and potentially a trial.

Psychological crime scene analysis (profiling) can also be conducted using the criminal's MO and signature behaviors. Profiling provides insight into the thought processes of the attacker and can be used to identify the attacker or, at the very least, the tool he used to conduct the crime.

Evidence Collection and Handling

Good evidence is the bedrock on which any sound investigation is built. When dealing with any incident that might end up in court, digital evidence must be handled in a careful fashion so that it can be admissible no matter what jurisdiction is prosecuting

a defendant. Within the United States, the *Scientific Working Group on Digital Evidence (SWGDE)* aims to ensure consistency across the forensic community. The principles developed by SWGDE for the standardized recovery of computer-based evidence are governed by the following attributes:

- Consistency with all legal systems
- Allowance for the use of a common language
- Durability
- Ability to cross international and state boundaries
- Ability to instill confidence in the integrity of evidence
- Applicability to all forensic evidence
- Applicability at every level, including that of individual, agency, and country

The international standard on digital evidence handling is ISO/IEC 27037: *Guidelines for Identification, Collection, Acquisition, and Preservation of Digital Evidence*. This document identifies four phases of digital evidence handling, which are identification, collection, acquisition, and preservation. Let's take a closer look at each.



NOTE You must ensure that you have the legal authority to search for and seize digital evidence before you do so. If in doubt, consult your legal counsel.

Identification

The first phase of digital evidence handling is to identify the digital crime scene. Rarely does only one device comprise the scene of the crime. More often than not, digital evidence exists on a multitude of other devices such as routers, network appliances, cloud services infrastructure, smartphones, and even IoT devices. Whether or not you have to secure a court order to seize evidence, you want to be very deliberate about determining what you think you need to collect and where it might exist.

When you arrive at the crime scene (whether it be physical or virtual), you want to carefully document everything you see and do. If you're dealing with a physical crime scene, photograph it from every possible angle before you touch anything. Label wires and cables and then snap a photo of the labeled system before it is disassembled. Remember that you want to instill confidence in the integrity of evidence and how it was handled from the very onset.

Identifying evidence items at a crime scene may not be straightforward. You could discover wireless networks that would allow someone to remotely tamper with the evidence. This would require you to consider ways to isolate the evidence from radio frequency (RF) signals in order to control the crime scene. There may also be evidence in devices (e.g., thumb drives) that are hidden either deliberately or unintentionally. Law enforcement agents sometimes resort to using specially trained dogs that can sniff out

Controlling the Crime Scene

Whether the crime scene is physical or digital, it is important to control who comes in contact with the evidence of the crime to ensure its integrity. The following are just some of the steps that should take place to protect the crime scene:

- Only allow authorized individuals access to the scene. These individuals should have knowledge of basic crime scene analysis.
- Document who is at the crime scene. In court, the integrity of the evidence may be in question if too many people were milling around the crime scene.
- Document who were the last individuals to interact with the systems.
- If the crime scene does become contaminated, document it. The contamination may not negate the derived evidence, but it will make investigating the crime more challenging.

electronics. Thoroughness in identifying evidence is the most important consideration in this phase, and this may require you to think outside the box to ensure you don't miss or lose a critical evidentiary item.

Collection

Once you've identified the evidence you need, you can begin collecting it. Evidence collection is the process of gaining physical control over items that could potentially have evidentiary value. This is where you walk into someone's office and collect their computer, external hard drives, thumb drives, and so on. It is critical that you have the legal authority to do this and that you document what you take, where you take it from, and what its condition is at the time.

Each piece of evidence should be labeled in some way with the date, time, initials of the collector, and a case number if one has been assigned. The piece of evidence should then be placed in a container, which should be sealed (ideally with evidence tape) so that tampering can be detected. An example of the data that should be collected and displayed on each evidence container is shown in Figure 22-1.

After everything is properly labeled, a chain of custody log should be made for each container and an overall log should be made capturing all events. A *chain of custody* documents each person that has control of the evidence at every point in time. In large investigations, one person may collect evidence, another may transport it, and a third may store it. Keeping track of all these individuals' possession of the evidence is critical to proving in court that the evidence was not tampered with. It is not hard for a good defense attorney to get evidence dismissed from court because of improper handling. For this reason, the chain of custody should follow evidence through its entire life cycle, beginning with identification and ending with its destruction, permanent archiving, or return to owner.

EVIDENCE	
Station/Section/Unit/Dept_____	
Case number_____	Item#_____
Type of offense_____	
Description of evidence_____	

Suspect_____	
Victim_____	
Date and time of recovery_____	
Location of recovery_____	
Recovered by_____	
CHAIN OF CUSTODY	
Received from_____	By_____
Date_____	Time_____ A.M./P.M.
Received from_____	By_____
Date_____	Time_____ A.M./P.M.
Received from_____	By_____
Date_____	Time_____ A.M./P.M.
Received from_____	By_____
Date_____	Time_____ A.M./P.M.
WARNING: THIS IS A TAMPER EVIDENT SECURITY PACKAGE. ONCE SEALED, ANY ATTEMPT TO OPEN WILL RESULT IN OBVIOUS SIGNS OF TAMPERING.	

Figure 22-1 Evidence container data

Evidence collection activities can get tricky depending on what is being searched for and where. For example, American citizens are protected by the Fourth Amendment against unlawful search and seizure, so law enforcement agencies must have probable cause and request a search warrant from a judge or court before conducting such a search. The actual search can take place only in the areas outlined by the warrant. The Fourth Amendment does not apply to actions by private citizens unless they are acting as police agents. So, for example, if Kristy's boss warned all employees that the management could remove files from their computers at any time, and her boss is not a police officer or acting as a police agent, she could not successfully claim that her Fourth Amendment rights were violated. Kristy's boss may have violated some specific privacy laws, but he did not violate Kristy's Fourth Amendment rights.

In some circumstances, a law enforcement agent is legally permitted to seize evidence that is not included in the search warrant, such as if the suspect tries to destroy the evidence. In other words, if there is an impending possibility that evidence might be destroyed, law enforcement may quickly seize the evidence to prevent its destruction.

This is referred to as *exigent circumstances*, and a judge will later decide whether the seizure was proper and legal before allowing the evidence to be admitted. For example, if a police officer had a search warrant that allowed him to search a suspect's living room but no other rooms and then he saw the suspect putting a removable drive in his pocket while standing in another room, the police officer could seize the drive even though it was outside the area covered under the search warrant.



EXAM TIP Always treat an investigation, regardless of type, as if it would ultimately end up in a courtroom.

Acquisition

In most corporate investigations involving digital evidence, the sort of Crime TV collection we just described will not take place unless law enforcement is involved. Instead, the IR team will probably be able to piece together a timeline of activities from various network resources and you may have to collect only a single laptop. In many cases you can probably acquire the evidence you need remotely without seizing any devices at all. Whatever the case, you ultimately need to get a hold of the data that will confirm or deny the claim that is being investigated, and you must do it in a forensically sound manner.

Acquisition means creating a forensic image of digital data for examination. Generally, speaking, there are two types of acquisition: physical and logical. In *digital acquisition*, the investigator makes a bit-for-bit copy of the contents of a physical storage device, bypassing the operating system. This includes all files, of course, but also free space and previously deleted data. In *logical acquisition*, on the other hand, the forensic image is of the files and folders in a file system, which means we rely on the operating system. This approach is sometimes necessary when dealing with evidence that exists in cloud services, where physical acquisition is normally not possible.

Before creating a forensic image, the investigator must have a medium onto which to copy the data, and ensure this medium has been properly purged, meaning it does not contain any preexisting data. (In some cases, hard drives that were thought to be new and right out of the box contained old data not purged by the vendor.) Two copies are normally created: a *primary image* (a control copy that is stored in a library) and a *working image* (used for analysis and evidence collection). To ensure that the original image is not modified, it is important to compute the cryptographic hashes (e.g., SHA-1) for files and directories before and after the analysis to prove the integrity of the original image.

The investigator works from the duplicate image because it preserves the original evidence, prevents inadvertent alteration of original evidence during examination, and allows re-creation of the duplicate image if necessary.

Acquiring evidence on live systems and those using network storage further complicates matters because you cannot turn off the system to make a copy of the hard drive. Imagine the reaction you'd receive if you were to tell an IT manager that you need to shut down a primary database or e-mail system. It wouldn't be favorable. So these systems and others, such as those using on-the-fly encryption, must be imaged while they are running.

In fact, some evidence is very volatile and can only be collected from a live system. Examples of volatile data that could have evidentiary value include

- Registers and cache
- Process tables and ARP cache
- System memory (RAM)
- Temporary file systems
- Special disk sectors

Preservation

To preserve evidence in a forensically sound manner, you must have established procedures based on legally accepted best practices, and your staff must follow those procedures to the letter. We've already covered two crucial steps in the chain of evidence and the use of hashes to verify that the evidence has not been altered. Another element of preserving digital evidence is ensuring that only a small group of qualified individuals have access to the evidence, and then only to perform specific functions. Again, this access needs to be part of your established procedures. In some cases, organizations implement two-person control of digital evidence to minimize the risk of tampering.

We introduced the topic of evidence storage in Chapter 10, but it bears pointing out that storage of media evidence should be dust-free and kept at room temperature without much humidity, and, of course, the media should not be stored close to any strong magnets or magnetic fields. Even if you don't have a dedicated evidence storage area, you should ensure that whatever space you commandeer is used strictly for this purpose, at least for the life of the investigation.

What Is Admissible in Court?

There are limits to what evidence can be introduced into a legal proceeding. Though the details will be different in each jurisdiction around the world, generally, digital evidence is admissible in court if it meets three criteria:

- **Relevance** Evidence must be relevant to the case, meaning it must help to prove facts being alleged. If a suspect is accused of murder, then a web search history for favorite vacationing spots is probably irrelevant. Judges typically rule on relevance of evidence.
- **Reliability** Evidence must be acquired using a sound forensic methodology that prevents alteration and ensures the evidence remains unaltered during the forensic examination. Multiple high-profile cases in recent years have had evidence rendered inadmissible because the chain of custody was broken.
- **Legality** The persons acquiring and presenting the evidence must have the legal authority to do so. If you have a court-issued search warrant, you must limit collection to whatever is spelled out in it. If you are conducting a workplace investigation, you must limit your collection to organization-owned assets, and only after legal counsel agrees.

The reliability of evidence is most often established by chains of custody and cryptographic hashing. But there is another element to reliability that excludes evidence deemed to be hearsay. *Hearsay evidence* is any statement made outside of the court proceeding that is offered into evidence to prove the truth of the matter asserted in the statement. Suppose that David is accused of fraud and Eliza tells Frank that David told her he was stealing from the company. Eliza's testimony in court would be admissible, but Frank normally wouldn't be allowed to testify about what Eliza claims to have heard because, coming from him, it would be considered hearsay.

Hearsay evidence can also include many computer-generated documents such as log files. In some countries, such as the United States, when computer logs are to be used as evidence in court, they must satisfy a legal exception to the hearsay rule of the Federal Rules of Evidence (FRE) called the business records exception rule or business entry rule. Under this rule, a party could admit any records of a business (1) that were made in the regular course of business; (2) that the business has a regular practice to make such records; (3) that were made at or near the time of the recorded event; and (4) that contain information transmitted by a person with knowledge of the information within the document.

It is important to show that the logs, and all evidence, have not been tampered with in any way, which is the reason for the chain of custody of evidence. Several tools are available that run checksums or hashing functions on the logs, which will allow the team to be alerted if something has been modified.

When evidence is being collected, one issue that can come up is the user's expectation of privacy. If an employee is suspected of, and charged with, a computer crime, he might claim that his files on the computer he uses are personal and not available to law enforcement and the courts. This is why it is important for organizations to conduct security awareness training, have employees sign documentation pertaining to the acceptable use of the organization's computers and equipment, and have legal banners pop up on every employee's computer when they log on. These are key elements in establishing that a user has no right to privacy when he is using organization equipment. The following banner is suggested by CERT Advisory:

This system is for the use of authorized users only. Individuals using this computer system without authority, or in excess of their authority, are subject to having all of their activities on this system monitored and recorded by system personnel.

In the course of monitoring an individual improperly using this system, or in the course of system maintenance, the activities of authorized users may also be monitored.

Anyone using this system expressly consents to such monitoring and is advised that if such monitoring reveals possible evidence of criminal activity, system personnel may provide the evidence of such monitoring to law enforcement officials.

This explicit warning strengthens a legal case that can be brought against an employee or intruder, because the continued use of the system after viewing this type of warning implies that the person acknowledges the security policy and gives permission to be monitored.



NOTE Don't dismiss the possibility that as an information security professional you will be responsible for entering evidence into court. Most tribunals, commissions, and other quasi-legal proceedings have admissibility requirements. Because these requirements can change between jurisdictions, you should seek legal counsel to better understand the specific rules for your jurisdiction.

Digital Forensics Tools, Tactics, and Procedures

Digital forensics is a science and an art that requires specialized techniques for the recovery, authentication, and analysis of electronic data for the purposes of a digital criminal investigation. It is a fusion of computer science, IT, engineering, and law. When discussing computer forensics with others, you might hear the terms computer forensics, network forensics, electronic data discovery, cyberforensics, and forensic computing.

Forensics Field Kits

When a forensics team is deployed, the forensic investigators should be properly equipped with all the tools and supplies that they'll need to conduct the investigation. The following are some of the common items in forensics field kits:

- **Documentation tools** Tags, labels, forms, and written procedures
- **Disassembly and removal tools** Antistatic bands, pliers, tweezers, screwdrivers, wire cutters, and so on
- **Package and transport supplies** Antistatic bags, evidence bags and tape, cable ties, and others
- **Cables and adapters** Enough to connect to every physical interface you may come across



(ISC)² uses *digital forensics* as a synonym for all of these other terms, so that's what you'll see on the CISSP exam.

Anyone who conducts a forensic investigation must be properly skilled in this trade and know what to look for. If someone reboots the attacked system or inspects various files, this could corrupt viable evidence, change timestamps on key files, and erase footprints the criminal may have left. Most digital evidence has a short lifespan and must be collected quickly and in the *order of volatility*. In other words, the most volatile or fragile evidence should be collected first. In some situations, it is best to remove the system from the network, dump the contents of the memory, power down the system, and make a sound image of the attacked system and perform forensic analysis on this copy. Working on the copy instead of the original drive ensures that the evidence stays unharmed on the original system in case some steps in the investigation actually corrupt or destroy data. Dumping the memory contents to a file before doing any work on the system or powering it down is a crucial step because of the information that could be stored there. This is another method of capturing fragile information. However, this creates a sticky situation: capturing RAM or conducting live analysis can introduce changes to the crime scene because various state changes and operations take place. Whatever method the forensic investigator chooses to use to collect digital evidence, that method must be documented. This is the most important aspect of evidence handling.

Forensic Investigation Techniques

To ensure that forensic investigations are carried out in a standardized manner and the evidence collected is admissible, it is necessary for the investigative team to follow specific laid-out steps so that nothing is missed. Figure 22-2 illustrates the phases through a common investigation process and lists various techniques that fall under each phase. Each team or organization may come up with its own steps, but all should be essentially accomplishing the same things:

- Identification
- Preservation
- Collection
- Examination
- Analysis
- Presentation
- Decision



NOTE The principles of criminalistics are included in the forensic investigation process. They are identification of the crime scene, protection of the environment against contamination and loss of evidence, identification of evidence and potential sources of evidence, and the collection of evidence. In regard to minimizing the degree of contamination, it is important to understand that it is impossible not to change a crime scene—be it physical or digital. The key is to minimize changes and document what you did and why, and how the crime scene was affected.

Identification	Preservation	Collection	Examination	Analysis	Presentation
Event/crime detection	Case management	Preservation	Preservation	Preservation	Documentation
Resolve signature	Imaging technologies	Approved methods	Traceability	Traceability	Expert testimony
Profile detection	Chain of custody	Approved software	Validation techniques	Statistical	Clarification
Anomalous detection	Time synchronization	Approved hardware	Filtering techniques	Protocols	Mission impact statement
Complaints		Legal authority	Pattern matching	Data mining	Recommended countermeasure
System monitoring		Lossless compression	Hidden data discovery	Timeline	Statistical interpretation
Audit analysis		Sampling	Hidden data extraction	Link	
		Data reduction		Spatial	
		Recovery techniques			

Figure 22-2 Characteristics of the different phases through an investigation process

During the examination and analysis process of a forensic investigation, it is critical that the investigator work from an image that contains *all* of the data from the original disk. It should be a bit-level copy, sector by sector, to capture deleted files, slack spaces, and unallocated clusters. These types of images can be created through the use of a specialized tool such as Forensic Toolkit (FTK), EnCase Forensic, or the dd Unix utility. A file copy tool does not recover all data areas of the device necessary for examination. Figure 22-3 illustrates a commonly used tool in the forensic world for evidence collection.

The next step is the analysis of the evidence. Forensic investigators use a scientific method that involves

- Determining the characteristics of the evidence, such as whether it's admissible as primary or secondary evidence, as well as its source, reliability, and permanence
- Comparing evidence from different sources to determine a chronology of events
- Event reconstruction, including the recovery of deleted files and other activity on the system

This can take place in a controlled lab environment or, thanks to hardware write-blockers and forensic software, in the field. When investigators analyze evidence in a lab, they are dealing with “dead forensics”; that is, they are working only with static data. Live forensics, which takes place in the field, includes volatile data. If evidence is lacking, then an experienced investigator should be called in to help complete the picture.

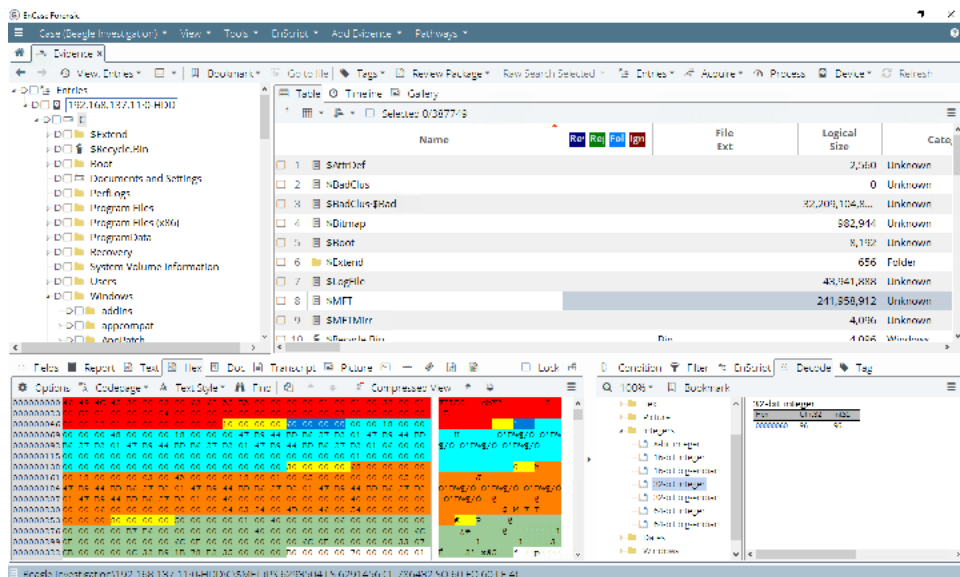


Figure 22-3 EnCase Forensic can be used to collect digital forensic data.

Finally, the interpretation of the analysis should be presented to the appropriate party. This could be a judge, lawyer, CEO, or board of directors. Therefore, it is important to present the findings in a format that will be understood by a nontechnical audience. As a CISSP, you should be able to explain these findings in layperson's terms using metaphors and analogies. Of course, the findings, which are top secret or company confidential, should be disclosed only to authorized parties. This may include the legal department or any outside counsel that assists with the investigation.

Other Investigative Techniques

Unless you work for a law enforcement agency, most of the investigations in which you will be involved are likely to focus on digital forensics investigative techniques. These techniques are applied when a device was compromised, or a malicious insider attempted to steal sensitive files, or something like that. All the evidence you need is probably in a device that you can get your hands on, so you can collect it, acquire it, analyze it, and get to the facts with just digital evidence. However, there may be other situations in which you'll need other types of evidence either in addition to or instead of 1's and 0's copied from some storage device. Interviews, surveillance, and undercover investigative techniques are some of the practices for acquiring evidence that you should be familiar with.

Interviews

Interviews can be effective for ascertaining facts when you have willing interviewees. Interviewing is both an art and a science, and the specific techniques you use will vary

from case to case. Typically, interviews are conducted by a business unit manager with assistance from the human resources and legal departments. This doesn't, however, completely relieve you as an information security professional from responsibility during the interviewing process. You may be asked to provide input or observe an interview in order to clarify technical information that comes up in the course of questioning.

Whether you are conducting an interview or your technical assistance is needed for an interview, keep the following best practices in mind:

- *Have a plan.* Without a plan, the interview will be ineffective. Prepare an outline beforehand that focuses on getting the information you need from each interviewee. However, you should remain flexible and not read off a script.
- *Be fair and objective.* If you are conducting an interview, it is to get to the facts of an incident, not necessarily to reinforce whatever conclusions you may have already reached. Keep an open mind, focus on the facts, and try to avoid any biases.
- *Compartmentalize information.* Your interview plan should address what information you share with each interviewee, and what you don't share. You should not tell one interviewee what another said unless it's absolutely essential and legally permissible.
- *One interviewee at a time.* Interviewing multiple individuals together can introduce problematic group dynamics such as peer pressure. It can also lead interviewees to distort or suppress information.
- *Do not record the interview.* Recording devices can have a chilling effect on interviewees. Instead, have at least one notetaker in the room and, after the interview is complete, read back the notes to the interviewee to ensure their accuracy. If you must record the interview, ensure you comply with all applicable legal requirements (e.g., consent of all parties).
- *Keep it confidential.* Do your best to keep every aspect of the investigation under wraps. Even the fact that someone is being interviewed about an incident can have a damaging reputational effect for that person.

The employee interviewer should be in a position that is senior to the employee subject. A vice president is not going to be very intimidated or willing to spill his guts to the mailroom clerk. The interview should be held in a private place, in an environment conducive to making the subject relatively comfortable and at ease. If exhibits are going to be shown to the subject, they should be shown one at a time, and otherwise kept in a folder. It is not necessary to read a person their rights before the interview unless it is performed by law enforcement officers.

Surveillance

Two main types of surveillance are used when it comes to identifying computer crimes: physical surveillance and computer surveillance. *Physical surveillance* pertains to security cameras, security guards, and closed-circuit TV (CCTV), which may capture evidence.

Physical surveillance can also be used by an undercover agent to learn about the suspect's spending activities, family and friends, and personal habits in the hope of gathering more clues for the case.

Computer surveillance pertains to passively monitoring (auditing) events by using network sniffers, keyboard monitors, wiretaps, and line monitoring. In most jurisdictions, active monitoring may require a search warrant. In most workplace environments, to legally monitor an individual, the person must be warned ahead of time that her activities may be subject to this type of monitoring.

Undercover

Undercover investigative techniques are pretty rare in most corporate investigations, but can provide information and evidence that would be difficult to acquire otherwise. The goal of undercover work is to assume an identity that allows the investigator to blend into the suspect's environment to observe, and perhaps record, the suspect's actions.

A thin line exists between enticement and entrapment when it comes to capturing a suspect's actions. *Enticement* is legal and ethical, whereas *entrapment* is neither legal nor ethical. In the world of computer crimes, a honeypot is a good example to explain the difference between enticement and entrapment. Organizations put systems in their screened subnets that either emulate services that attackers usually like to take advantage of or actually have the services enabled. The hope is that if an attacker breaks into the organization's network, she will go right to the honeypot instead of the systems that are actual production machines. The attacker will be *enticed* to go to the honeypot system because it has many open ports and services running and exhibits vulnerabilities that the attacker would want to exploit. The organization can log the attacker's actions and later attempt to prosecute.

The action in the preceding example is legal unless the organization crosses the line to entrapment. For example, suppose a web page has a link that indicates that if an individual clicks it, she could then download thousands of MP3 files for free. However, when she clicks that link, she is taken to the honeypot system instead, and the organization records all of her actions and attempts to prosecute. Entrapment does not prove that the suspect had the intent to commit a crime; it only proves she was successfully tricked.

Forensic Artifacts

One of the grandfathers of forensic science, Dr. Edmond Locard, famously stated that "every contact leaves a trace." This principle, known as Locard's exchange principle, states that criminals always leave something behind at the crime scene. This fragmentary or trace evidence is a *forensic artifact*. A forensic artifact is anything that has evidentiary value. On a typical computer, the following are examples of forensic artifacts:

- Deleted items (in the recycle bin or trash)
- Web browser search history
- Web browser cache files
- E-mail attachments

- Skype history
- Windows event logs
- Prefetch files

Forensic artifacts can also be evidentiary items relating to network traffic. Network forensics is a subdiscipline that is focused on what happened on the network rather than on the endpoints. The tools used in network forensics are unique to that subdiscipline, and so are the artifacts for which the investigator looks. Tools used in network forensics include NDR solutions, SIEM systems, and the log files of any network device or server. They also include network sniffers that can capture full network frames. The following are some of the more useful network artifacts an investigator would be interested in:

- DNS log records
- Web proxy log records
- IDS/IPS alerts
- Packet capture (pcap) files

Finally, with the proliferation of mobile devices such as smartphones, tablets, and smartwatches, we must not overlook forensic artifacts stored on them. Unlike traditional computers, mobile devices are usually carried by their users around the clock. This means mobile devices tend to document multiple aspects of a person's life, some of which can serve as evidence of criminal activity.

Though mobile devices can be a treasure trove of information for the forensic investigator, they are not always easy to acquire and analyze. For starters, there are so many different models that no single tool can acquire all evidence from all devices. Staff expertise is similarly challenged by this diversity, because an investigator who is skilled at iPhone analysis may not be able to operate at the same level given an Android device. Just to make things more interesting, there is also the issue of encryption, which is prevalent in mobile devices these days.

Still, if forensic investigators can overcome these challenges, mobile devices are excellent sources of evidence for a variety of criminal activity. Among the most useful forensic artifacts found in them are

- Call logs
- SMS messages
- E-mail messages
- Web browser history

Reporting and Documenting

We already covered reporting in a fair amount of detail in Chapter 19. When it comes to investigations, however, there are some additional issues to consider. First and foremost, the need to document *everything* you do cannot be overstated. If you cannot account for

or explain the why of any activities you undertook, it may render evidence inadmissible in court or even undermine the whole case. For this reason, many organizations assign investigators to work in teams of two, where one person documents while the other conducts the investigation. Most forensic analysis tools have a feature that automatically logs everything an investigator does with the tool.

Another issue that is particularly important in writing investigation reports is the need to remain completely logical and factual. Any conclusions you reach must follow logically from a sequence of facts that you spell out for the reader. For example, suppose that Carlos is one of your staff and is suspected of sending sensitive files to a competitor in hopes of landing a lucrative job with them. Even if you are sure he did it (after examining his computer), you should not just jump out and say so. Instead, you show how the forensic artifacts that you found, when arrayed on a timeline, substantiate the claim that Carlos sent sensitive files to a competitor. You'd start by establishing that he was logged into his computer, and then he logged into his personal e-mail account through a webmail interface, and then an e-mail was sent containing sensitive files x, y, and z, and then the e-mail was deleted from his sent items, and so on. It is ultimately up to the reader (presumably a senior manager or court official) to determine guilt or innocence. Your job is to establish the facts and determine whether or not they are consistent with the allegation.

Chapter Review

Incident management is a critical function for any organization. Odds are that if you are among the lucky few who haven't had a major incident yet, you will be faced with one in the near future. In fact, the IronNet 2021 Cybersecurity Impact Report found that 86 percent of respondents had a cybersecurity incident so severe in the previous year that it required a C-level or board meeting. Even if you've outsourced IR to a third-party service provider, you still need to have an incident management policy and an IR plan to guide the conduct of the entire organization before, during, and after an incident. The policy establishes authorities and responsibilities, while the plan specifies the procedures to be followed.

The other major topic we discussed in this chapter is investigations. Thankfully, the need to conduct investigations is fairly rare in most organizations. But therein lies the problem: if you hardly ever need to recall knowledge or practice skills, you are certain to lose them. This is why having detailed standard procedures for investigative work is absolutely essential. For example, evidence acquisition, as we saw, is a complex process that has very little room for errors, particularly if the evidence will end up in court (and we should always assume it will).

Quick Review

- A security event is any occurrence that can be observed, verified, and documented, whereas a security incident is one or more related events that negatively affect the organization and/or impact its security posture.

- A good incident response team should consist of representatives from various business units, such as the legal department, HR, executive management, the communications department, physical/corporate security, IS security, and information technology.
- Incident management encompasses seven phases according to the CISSP CBK: detection, response, mitigation, reporting, recovery, remediation, and lessons learned.
- The detection phase encompasses the search for indicators that an event has occurred and the formal declaration of the event.
- The response phase entails the initial actions undertaken to contain the damage caused by a security incident.
- The goal of the mitigation phase is to eradicate the threat actor from the affected systems.
- Incident reporting occurs at various phases of incident management.
- The aim of the recovery phase is to restore full, trustworthy functionality to the organization.
- In the remediation phase, the incident response team decides which security controls need to be deployed or changed to prevent the incident from recurring.
- The lessons learned phase is important to determine what needs to go into the incident response process and documentation, with the goal of continuous improvement.
- The incident management policy (IMP) establishes authorities and responsibilities across the entire organization, identifies the incident response (IR) lead for the organization, and describes what every staff member is required to do with regard to incidents.
- The incident response plan (IRP) gets into the details of what should be done when responding to suspected incidents, and includes roles and responsibilities, incident classification, notifications, and operational tasks.
- Incident classification criteria allow the organization to prioritize IR assets and usually consider the impact and type of the incident, and urgency with which the response must be started.
- A runbook is a collection of procedures that the IR team will follow for specific types of incidents.
- The four phases of evidence handling are identification, collection, acquisition, and preservation.
- Evidence collection is the process of gaining physical control over devices that could potentially have evidentiary value.
- A chain of custody documents each person that has control of the evidence at every point in time.
- Acquisition means creating a forensic image of digital data for examination.

- Evidence preservation requires maintaining a chain of custody and cryptographic hashes of all digital evidence, and also controlling access to the evidence.
- To be admissible in court, evidence must be relevant, reliable, and legally obtained.
- To be admissible in court, business records such as computer logs have to be made and collected in the normal course of business, not specially generated for a case in court. Business records can easily be deemed hearsay if there is no firsthand proof of their accuracy and reliability.
- Digital forensics is a science and an art that requires specialized techniques for the recovery, authentication, and analysis of electronic data for the purposes of a digital criminal investigation.
- In addition to forensic techniques, organizations sometimes use interviews, surveillance, and undercover investigation techniques.
- When looking for suspects, it is important to consider the motive, opportunity, and means (MOM).
- A forensic artifact is anything that has evidentiary value.

Questions

Please remember that these questions are formatted and asked in a certain way for a reason. Keep in mind that the CISSP exam is asking questions at a conceptual level. Questions may not always have the perfect answer, and the candidate is advised against always looking for the perfect answer. Instead, the candidate should look for the best answer in the list.

1. What are the phases of incident management?
 - A. Identification, collection, acquisition, and preservation
 - B. Detection, response, mitigation, reporting, recovery, remediation, and lessons learned
 - C. Protection, containment, response, remediation, and reporting
 - D. Analysis, classification, incident declaration, containment, eradication, and investigation
2. During which phase of incident management does the IR team contain the damage caused by a security incident?
 - A. Preservation
 - B. Response
 - C. Eradication
 - D. Remediation

3. During which phase of incident management are security controls deployed or changed to prevent the incident from recurring?
 - A. Preservation
 - B. Response
 - C. Eradication
 - D. Remediation
4. Which document establishes authorities and responsibilities with regard to incidents across the entire organization?
 - A. Incident management policy
 - B. Incident response plan
 - C. Incident response runbook
 - D. Incident classification criteria
5. After a computer forensic investigator seizes a computer during a crime investigation, what is the next step?
 - A. Label and put it into a container, and then label the container
 - B. Dust the evidence for fingerprints
 - C. Make an image copy of the disks
 - D. Lock the evidence in the safe
6. Which of the following is a necessary characteristic of evidence for it to be admissible?
 - A. It must be real.
 - B. It must be noteworthy.
 - C. It must be reliable.
 - D. It must be important.
7. Which of the following is *not* considered a best practice when interviewing willing witnesses?
 - A. Compartmentalize information
 - B. Interview one interviewee at a time
 - C. Be fair and objective
 - D. Record the interview

Use the following scenario to answer Questions 8–10. You recently improved your organization's security posture, which now includes a fully staffed security operations center (SOC), network detection and response (NDR) and endpoint detection and response (EDR) systems, centrally managed updates and data backups, and network segmentation using VLANs. It's the end of the workday and just as you are getting ready to go

home your SOC detects a ransomware infection affecting at least two workstations in your marketing department. The SOC manager declares an incident and activates the IR team.

8. What should be your IR team's first action?
 - A. Determine the scope of the infection across the organization
 - B. Isolate the marketing VLAN from the rest of the network
 - C. Disconnect the infected computers from the network
 - D. Determine why the EDR system failed to protect the workstations
9. Using your NDR system, you determine the external hosts from which the malware was downloaded and with which the infected systems were communicating. As part of the remediation phase, which of the following is the next best action to take with this information?
 - A. Determine whether the external hosts you identified are related to the incident
 - B. Block traffic to/from the external hosts that you identified
 - C. Visit the remote hosts using a forensic workstation to acquire evidence
 - D. Share the address of the hosts with your partners as indicators of compromise (IOCs)
10. Luckily, this version of ransomware is buggy, and you find a security researcher's blog with detailed instructions for how to decrypt infected systems. Which of the following approaches will best mitigate the incident and make the affected systems operational again?
 - A. Follow the directions to decrypt the systems and remove the malware
 - B. Reinstall from a golden master and restore the data from backups
 - C. Reinstall from a golden master even though you have no backups
 - D. Restore the systems from the last known-good system backup

Answers

1. **B.** Incident management encompasses seven phases according to the CISSP CBK: detection, response, mitigation, reporting, recovery, remediation, and lessons learned.
2. **B.** The goal of containment during the response phase is to prevent or reduce any further damage from this incident so that you can begin to mitigate and recover. Done properly, this buys the IR team time for a proper investigation and determination of the incident's root cause.
3. **D.** In the remediation phase, you decide which control changes (e.g., firewall or IDS/IPS rules) are needed to preclude this incident from happening again. Another aspect of remediation is the identification of indicators of attack (IOAs)

that can be used in the future to detect this attack in real time (i.e., as it is happening) as well as indicators of compromise (IOCs), which tell you when an attack has been successful and your security has been compromised.

4. **A.** The incident management policy (IMP) establishes authorities and responsibilities across the entire organization, identifies the incident response (IR) lead for the organization, and describes what every staff member is required to do with regard to incidents. The incident response plan (IRP) gets into the details of what should be done when responding to suspected incidents, and includes roles and responsibilities, incident classification, notifications, and operational tasks. A runbook is a collection of procedures that the IR team will follow for specific types of incidents.
5. **C.** Several steps need to be followed when gathering and extracting evidence from a scene. Once a computer has been confiscated, the first thing the computer forensics team should do is make an image of the hard drive. The team will work from this image instead of the original hard drive so that the original stays in a pristine state and the evidence on the drive is not accidentally corrupted or modified.
6. **C.** For evidence to be admissible, it must be relevant to the case, reliable, and legally obtained. For evidence to be reliable, it must be consistent with fact and must not be based on opinion or be circumstantial.
7. **D.** Recording devices can have a chilling effect on interviewees. Instead, have at least one notetaker in the room and, after the interview is complete, read back the notes to the interviewee to ensure their accuracy.
8. **B.** Having detected the incident, the next step is to respond by containing the damage that has been or is about to be done to your most critical assets. You could simply disconnect the infected systems from the network, but since there are multiple workstations and they are in the same department, it is probably better to isolate that entire VLAN until you can determine the true scope of the problem. Since this incident happened at the end of the workday, isolating the VLAN should have little or no impact on the marketing department.
9. **B.** In the remediation phase, you decide which security controls need to be put in place to prevent the attack from succeeding again. This includes controls that are hastily put into effect because you have high confidence that they will help in the short term. The situation in the question is a perfect example of when you bypass your change management process and quickly make changes to deal with the incident at hand. You probably want to share the IOCs with your partners (and perhaps your regional CERT), but that happens after you block the traffic.
10. **B.** You have a centralized backup system that was not affected, so you know you should have backups for all the workstations. The problem is that you may not know if any of the full-system backups also include the ransomware, so restoring systems from backups could bring you back to square one. It is best to reinstall the systems from golden masters and then restore only the data files. This process may take a bit longer, but it minimizes the risk of reinfection.

This page intentionally left blank

Disasters

This chapter presents the following:

- Recovery strategies
- Disaster recovery processes
- Testing disaster recovery plans
- Business continuity

It wasn't raining when Noah built the ark.

—Howard Ruff

Disasters are just regular features in our collective lives. Odds are that, at some point, we will all have to deal with at least one disaster (if not more), whether it be in our personal world or professional world. And when that disaster hits, figuring out a way to deal with it in real time is probably not going to go all that well for the unprepared. This chapter is all about thinking of all the terrible things that might happen, and then ensuring we have strategies and plans to deal with them. This doesn't just mean recovering from the disaster, but also ensuring that the business continues to operate with as little disruption as possible.

As the old adage goes, no battle plan ever survived first contact with the enemy, which is the reason why we must test and exercise plans until our responses as individuals and organizations are so ingrained in our brains that we no longer need to think about them. As terrible and complex disasters unfold around us, we will do the right things reflexively. Does that sound a bit ambitious? Perhaps. Still, it is our duty as cybersecurity professionals to do what we can to get our organizations as close to that goal as realistically possible. Let's see how we go about doing this.

Recovery Strategies

In the previous chapters in this part of the book, we have discussed preventing and responding to security incidents, including various types of investigations, as part of standard security operations. These are things we do day in and day out. But what happens on those rare occasions when an incident has disastrous effects? That is the realm of disaster recovery and business continuity planning. *Disaster recovery (DR)* is the set of practices that enables an organization to minimize loss of, and restore, mission-critical

technology infrastructure after a catastrophic incident. *Business continuity (BC)* is the set of practices that enables an organization to continue performing its critical functions through and after any disruptive event. As you can see, DR is mostly in the purview of safety and contingency operations, while BC is much broader than that. Accordingly, we'll focus on DR for most of this chapter but circle back to our roles in BC as cybersecurity leaders.



EXAM TIP As CISSPs, we are responsible for disaster recovery because it deals mostly with information technology and security. We provide inputs and support for business continuity planning but normally are not the lead for it.

Before we go much further, recall that we discussed the role of *maximum tolerable downtime (MTD)* values in Chapter 2. In reality, basic MTD values are a good start, but are not granular enough for an organization to figure out what it needs to put into place to be able to absorb the impact of a disaster. MTD values are usually “broad strokes” that do not provide the details needed to pinpoint the actual recovery solutions that need to be purchased and implemented. For example, if the business continuity planning (BCP) team determines that the MTD value for the customer service department is 48 hours, this is not enough information to fully understand what redundant solutions or backup technology should be put into place. MTD in this example does provide a basic deadline that means if customer service is not up within 48 hours, the company may not be able to recover and everyone should start looking for new jobs.

As shown in Figure 23-1, more than just MTD metrics are needed to get production back to normal operations after a disruptive event. We will walk through each of these metric types and see how they are best used together.

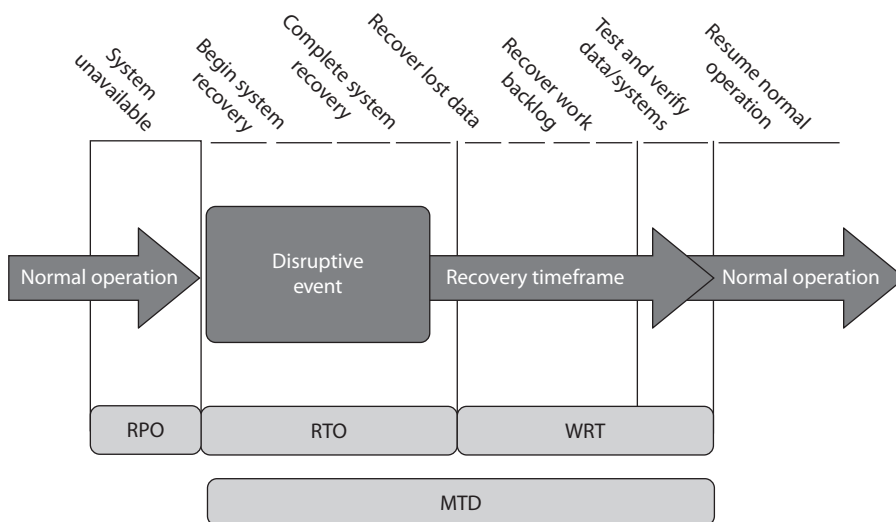


Figure 23-1 Metrics used for disaster recovery

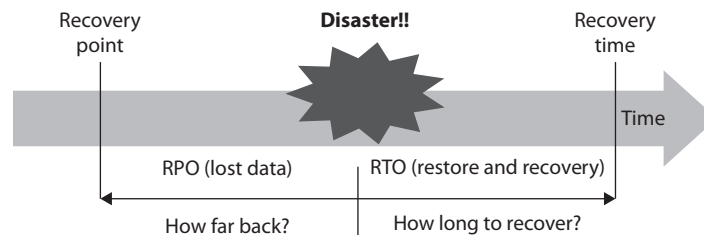
The *recovery time objective (RTO)* is the maximum time period within which a mission-critical system must be restored to a designated service level after a disruption to avoid unacceptable consequences associated with a break in business continuity. The RTO value is smaller than the MTD value, because the MTD value represents the time after which an inability to recover significant operations will mean severe and perhaps irreparable damage to the organization's reputation or bottom line. The RTO assumes that there is a period of acceptable downtime. This means that an organization can be out of production for a certain period of time and still get back on its feet. But if the organization cannot get production up and running within the MTD window, it may be sinking too fast to properly recover.

The *work recovery time (WRT)* is the maximum amount of time available for certifying the functionality and integrity of restored systems and data so they can be put back into production. RTO usually deals with getting the infrastructure and systems back up and running, and WRT deals with ensuring business users can get back to work using them. Another way to think of WRT is as the remainder of the overall MTD value after the RTO has passed.

The *recovery point objective (RPO)* is the acceptable amount of data loss measured in time. This value represents the earliest point in time at which data must be recovered. The higher the value of data, the more funds or other resources that can be put into place to ensure a smaller amount of data is lost in the event of a disaster. Figure 23-2 illustrates the relationship and differences between the use of RPO and RTO values.

The MTD, RTO, RPO, and WRT values are critical to understand because they will be the basic foundational measures used when determining the type of recovery solutions an organization must put into place, so let's dig a bit deeper into them. As an example of RTO, let's say a company has determined that if it is unable to process product order requests for 12 hours, the financial hit will be too large for it to survive. This means that the MTD for order processing is 12 hours. To keep things simple, let's say that RTO and WRT are 6 hours each. Now, suppose that orders are processed using on-premises servers, on a site with no backup power sources, and an ice storm causes a power outage that will take days to restore. Without a plan and supporting infrastructure already in place, it would be close to impossible to migrate the servers and data to a site with power within 6 hours. The RTO (that is, the maximum time to move the servers and data) would not be met (to say nothing of the WRT) and it would likely exceed the MTD, putting the company at serious risk of collapse.

Figure 23-2
RPO and RTO
measures in use



Now let's say that the same company did have a recovery site on a different power grid, and it was able to restore the order-processing services within a couple of hours, so it met the RTO requirement. But just because the systems are back online, the company still might have a critical problem. The company has to restore the data it lost during the disaster. Restoring data that is a week old does the company no good. The employees need to have access to the data that was being processed right before the disaster hit. If the company can only restore data that is a week old, then all the orders that were in some stage of being fulfilled over the last seven days could be lost. If the company makes an average of \$25,000 per day in orders and all the order data was lost for the last seven days, this can result in a loss of \$175,000 and a lot of unhappy customers. So just getting things up and running (meeting the RTO) is just part of the picture. Getting the necessary data in place so that business processes are up to date and relevant (RPO) is just as critical.

To take things one step further, let's say the company stood up the systems at its recovery site in two hours. It also had real-time data backup systems in place, so all of the necessary up-to-date data is restored. But no one actually tested the processes to recover data from backups, everyone is confused, and orders still cannot be processed and revenue cannot be collected. This means the company met its RTO requirement and its RPO requirement, but failed its WRT requirement, and thus failed the MTD requirement. Proper business recovery means *all* of the individual things have to happen correctly for the overall goal to be successful.



EXAM TIP An RTO is the amount of time it takes to recover from a disaster, and an RPO is the acceptable amount of data, measured in time, that can be lost from that same event.

The actual MTD, RTO, and RPO values are derived during the *business impact analysis (BIA)*, the purpose of which is to be able to apply criticality values to specific business functions, resources, and data types. A simplistic example is shown in Table 23-1. The company must have data restoration capabilities in place to ensure that mission-critical data is never older than one minute. The company cannot rely on something as slow as backup tape restoration, but must have a high-availability data replication solution in place. The RTO value for mission-critical data processing is two minutes or less. This means that the technology that carries out the processing functionality for this type of data cannot be down for more than two minutes. The company probably needs failover technology in place that will shift the load once it notices that a server goes offline.

Data Type	RPO	RTO
Mission critical	Continuous to 1 minute	Instantaneous to 2 minutes
Business critical	5 minutes	10 minutes
Business	3 hours	8 hours

Table 23-1 RPO and RTO Value Relationships

What Is the Difference Between Preventive Measures and Recovery Strategies?

Preventive mechanisms are put into place not only to try to reduce the possibility that the organization will experience a disaster, but also, if a disaster does hit, to lessen the amount of damage that will take place. Although the organization cannot stop a tornado from coming, for example, it could choose to move its facility from Tornado Alley to an area less prone to these weather events. As another example, the organization cannot stop a car from plowing into and taking out a transformer that it relies on for power, but it can have a separate power feed from a different transformer in case this happens.

Recovery strategies are processes designed to rescue the company after a disaster takes place. These processes integrate mechanisms such as establishing alternate sites for facilities, implementing emergency response procedures, and possibly activating the preventive mechanisms that have already been implemented.

In this same scenario, data that is classified as “Business” can be up to three hours old when the production environment comes back online, so a less frequent data replication process is acceptable. Because the RTO for business data is eight hours, the company can choose to have hot-swappable hard drives available instead of having to pay for the more complicated and expensive failover technology.

The DR team has to figure out what the company needs to do to actually recover the processes and services it has identified as being so important to the organization overall. In its business continuity and recovery strategy, the team closely examines the critical, agreed-upon business functions, and then evaluates the numerous recovery and backup alternatives that might be used to recover critical business operations. It is important to choose the right tactics and technologies for the recovery of each critical business process and service in order to assure that the set MTD values are met.

So what does the DR team need to accomplish? The team needs to actually define the recovery processes, which are sets of predefined activities that will be implemented and carried out in response to a disaster. More importantly, these processes must be constantly reevaluated and updated as necessary to ensure that the organization meets or exceeds the MTDs. It all starts with understanding the business processes that would have to be recovered in the aftermath of a disaster. Armed with that knowledge, the DR team can make good decisions about data backup, recovery, and processing sites, as well as overall services availability, all of which we explore in the next sections.

Business Process Recovery

A *business process* is a set of interrelated steps linked through specific decision activities to accomplish a specific task. Business processes have starting and ending points and are repeatable. The processes should encapsulate the knowledge about services, resources, and operations provided by an organization. For example, when a customer requests

to buy a book via a company's e-commerce site, the company's order fulfillment system must follow a business process such as this:

1. Validate that the book is available.
2. Validate where the book is located and how long it would take to ship it to the destination.
3. Provide the customer with the price and delivery date.
4. Verify the customer's credit card information.
5. Validate and process the credit card order.
6. Send the order to the book inventory location.
7. Send a receipt and tracking number to the customer.
8. Restock inventory.
9. Send the order to accounting.

The DR team needs to understand these different steps of the organization's most critical processes. The data is usually presented as a workflow document that contains the roles and resources needed for each process. The DR team must understand the following about critical business processes:

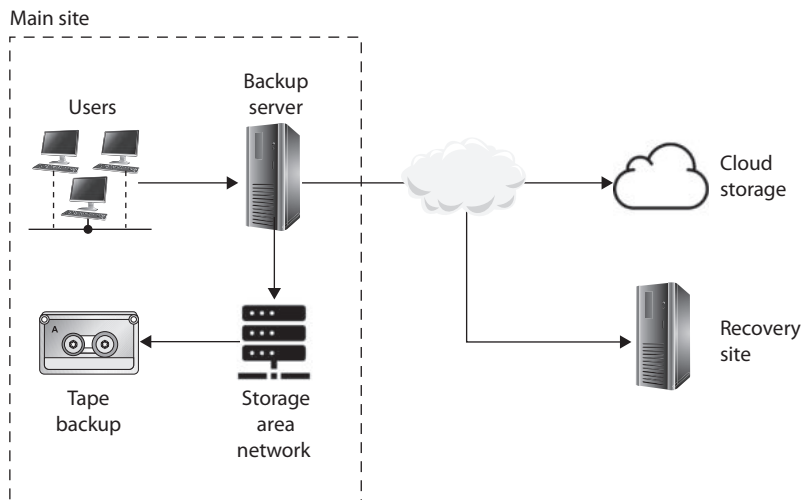
- Required roles
- Required resources
- Input and output mechanisms
- Workflow steps
- Required time for completion
- Interfaces with other processes

This will allow the team to identify threats and the controls to ensure the least amount of process interruption.

Data Backup

Data has become one of the most critical assets to nearly all organizations. It may include financial spreadsheets, blueprints on new products, customer information, product inventory, trade secrets, and more. In Chapter 2, we stepped through risk analysis procedures and, in Chapter 5, data classification. The DR team should not be responsible for setting up and maintaining the organization's data classification procedures, but the team should recognize that the organization is at risk if it does not have these procedures in place. This should be seen as a vulnerability that is reported to management. Management would need to establish another group of individuals who would identify the organization's data, define a loss criterion, and establish the classification structure and processes.

The DR team's responsibility is to provide solutions to protect this data and identify ways to restore it after a disaster. Data usually changes more often than hardware and software, so these backup or archival procedures must happen on a continual basis. The data backup process must make sense and be reasonable and effective. If data in the files changes several times a day, backup procedures should happen a few times a day or nightly to ensure all the changes are captured and kept. If data is changed once a month, backing up data every night is a waste of time and resources. Backing up a file and its corresponding changes is usually more desirable than having multiple copies of that one file. Online backup technologies usually record the changes to a file in a transaction log, which is separate from the original file.



The IT operations team should include a backup administrator, who is responsible for defining which data gets backed up and how often. These backups can be full, differential, or incremental, and are usually used in some type of combination with each other. Most files are not altered every day, so, to save time and resources, it is best to devise a backup plan that does not continually back up data that has not been modified. So, how do we know which data has changed and needs to be backed up without having to look at every file's modification date? This is accomplished by setting an *archive bit* to 1 if a file has been modified. The backup software reviews this bit when making its determination of whether the file gets backed up and, if so, clears the bit when it's done.

The first step is to do a *full backup*, which is just what it sounds like—all data is backed up and saved to some type of storage media. During a full backup, the archive bit is cleared, which means that it is set to 0. An organization can choose to do full backups only, in which case the restoration process is just one step, but the backup and restore processes could take a long time.

Most organizations choose to combine a full backup with a differential or incremental backup. A *differential process* backs up the files that have been modified since the *last full backup*. When the data needs to be restored, the full backup is laid down first, and then

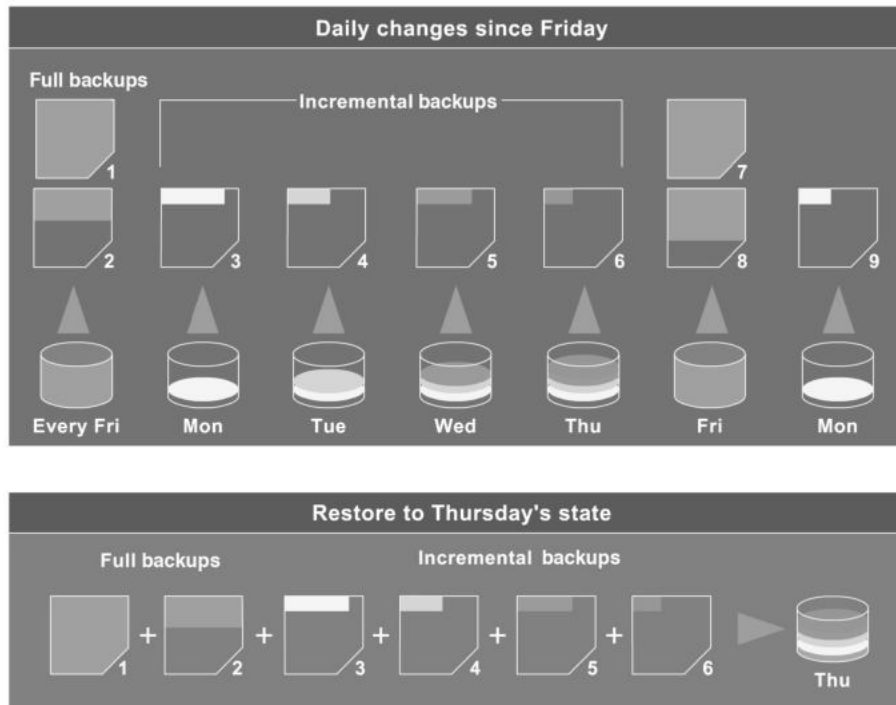


Figure 23-3 Backup software steps

the most recent differential backup is put down on top of it. The differential process does not change the archive bit value.

An *incremental process* backs up all the files that have changed since the *last full or incremental backup* and sets the archive bit to 0. When the data needs to be restored, the full backup data is laid down, and then each incremental backup is laid down on top of it in the proper order (see Figure 23-3). If an organization experienced a disaster and it used the incremental process, it would first need to restore the full backup on its hard drives and lay down every incremental backup that was carried out before the disaster took place (and after the last full backup). So, if the full backup was done six months ago and the operations department carried out an incremental backup each month, the backup administrator would restore the full backup and start with the older incremental backups taken since the full backup and restore each one of them until they were all restored.

Which backup process is best? If an organization wants the backup and restoration processes to be simple, it can carry out just full backups—but this may require a lot of hard drive space and time. Although using differential and incremental backup processes is more complex, it requires fewer resources and less time. A differential backup takes more time in the backing-up phase than an incremental backup, but it also takes less time

to restore than an incremental backup because carrying out restoration of a differential backup happens in two steps, whereas in an incremental backup, every incremental backup must be restored in the correct sequence.

Whatever the organization chooses, it is important to not mix differential and incremental backups. This overlap could cause files to be missed, since the incremental backup changes the archive bit and the differential backup does not.

Critical data should be backed up and stored onsite *and* offsite. The onsite backups should be easily accessible for routine uses and should provide a quick restore process so operations can return to normal. However, onsite backups are not enough to provide real protection. The data should also be held in an offsite facility in case of disasters. One decision the CISO needs to make is where the offsite location should be in reference to the main facility. The closer the offsite backup storage site is, the easier it is to access, but this can put the backup copies in danger if a large-scale disaster manages to take out the organization's main facility and the backup facility. It may be wiser to choose a backup facility farther away, which makes accessibility harder but reduces the risk. Some organizations choose to have more than one backup facility: one that is close and one that is farther away.

Backup Storage Strategies

A backup strategy must take into account that failure can take place at any step of the process, so if there is a problem during the backup or restoration process that could corrupt the data, there should be a graceful way of backing out or reconstructing the data

Restoring Data from Backups: A Cautionary Tale

Can we actually restore data from backups? Backing up data is a wonderful thing in life, but making sure it can be properly restored is even better. Many organizations have developed a false sense of security based on the fact that they have a very organized and effective process of backing up their data. That sense of security can disappear in seconds when an organization realizes in a time of crisis that its restore processes do not work. For example, one company had paid an offsite backup facility to use a courier to collect its weekly backup tapes and transport them to the offsite facility for safekeeping. What the company did not realize was that this courier used the subway and many times set the tapes on the ground while waiting for the subway train. A subway has many large engines that create their own magnetic field. This can have the same effect on media as large magnets, meaning that the data can be erased or corrupted. The company never tested its restore processes and eventually experienced a disaster. Much to its surprise, it found out that three years of data were corrupted and unusable.

Many other stories and experiences like this are out there. Don't let your organization end up as an anecdote in someone else's book because it failed to verify that its backups could be restored.

from the beginning. The procedures for backing up and restoring data should be easily accessible and comprehensible even to operators or administrators who are not intimately familiar with a specific system. In an emergency situation, the same person who always does the backing up and restoring may not be around, or outsourced consultants may need to be temporarily hired to meet the restoration time constraints.

There are four commonly used backup strategies that you should be aware of:

- **Direct-attached storage** The backup storage is directly connected to the device being backed up, typically over a USB cable. This is better than nothing, but is not really well suited for centralized management. Worse yet, many ransomware attacks look for these attached storage devices and encrypt them too.
- **Network-attached storage (NAS)** The backup storage is connected to the device over the LAN and is usually a storage area network (SAN) managed by a backup server. This approach is usually centrally managed and allows IT administrators to enforce data backup policies. The main drawback is that, if a disaster takes out the site, the data may be lost or otherwise be rendered inaccessible.
- **Cloud storage** Many organizations use cloud storage as either the primary or secondary repository of backup data. If this is done on a virtual private cloud, it has the advantage of providing offsite storage so that, even if the organization's site is destroyed by a disaster, the data is available for recovery. Obviously, WAN connectivity must be reliable and fast enough to support this strategy if it is to be effective.
- **Offline media** As ransomware becomes more sophisticated, we are seeing more instances of attackers going after NAS and cloud storage. If your data is critical enough that you have to decrease the risk of it being lost as close to zero as you can, you may want to consider offline media such as tape backups, optical discs, or even external drives that are disconnected after each backup (and potentially removed offsite). This is the slowest and most expensive approach, but is also the most resistant to attacks.

Electronic vaulting and remote journaling are other solutions that organizations should be aware of. *Electronic vaulting* makes copies of files as they are modified and periodically transmits them to an offsite backup site. The transmission does not happen in real time, but is carried out in batches. So, an organization can choose to have all files that have been changed sent to the backup facility every hour, day, week, or month. The information can be stored in an offsite facility and retrieved from that facility in a short amount of time.

This form of backup takes place in many financial institutions, so when a bank teller accepts a deposit or withdrawal, the change to the customer's account is made locally to that branch's database and to the remote site that maintains the backup copies of all customer records.

Electronic vaulting is a method of transferring bulk information to offsite facilities for backup purposes. *Remote journaling* is another method of transmitting data offsite, but this usually only includes moving the journal or transaction logs to the offsite facility, not the actual files. These logs contain the deltas (changes) that have taken place to the individual files. Continuing with the bank example, if and when data is corrupted and needs to be restored, the bank can retrieve these logs, which are used to rebuild the lost data. Journaling is efficient for database recovery, where only the reapplication of a series of changes to individual records is required to resynchronize the database.



EXAM TIP Remote journaling takes place in real time and transmits only the file deltas. Electronic vaulting takes place in batches and moves the entire file that has been updated.

An organization may need to keep different versions of software and files, especially in a software development environment. The object and source code should be backed up along with libraries, patches, and fixes. The offsite facility should mirror the onsite facility, meaning it does not make sense to keep all of this data at the onsite facility and only the source code at the offsite facility. Each site should have a full set of the most current and updated information and files.

Another software backup technology is *tape vaulting*. Many organizations back up their data to tapes that are then manually transferred to an offsite facility by a courier or an employee. This manual process can be error-prone, so some organizations use *electronic tape vaulting*, in which the data is sent over a serial line to a backup tape system at the offsite facility. The company that maintains the offsite facility maintains the systems and changes out tapes when necessary. Data can be quickly backed up and retrieved when necessary. This technology improves recovery speed, reduces errors, and allows backups to be run more frequently.

Data repositories commonly have replication capabilities, so that when changes take place to one repository (i.e., database) they are replicated to all the other repositories within the organization. The replication can take place over telecommunication links, which allow offsite repositories to be continuously updated. If the primary repository goes down or is corrupted, the replication flow can be reversed, and the offsite repository updates and restores the primary repository. Replication can be asynchronous or synchronous. *Asynchronous replication* means the primary and secondary data volumes are out of sync. Synchronization may take place in seconds, hours, or days, depending upon the technology in place. With *synchronous replication*, the primary and secondary repositories are always in sync, which provides true real-time duplication. Figure 23-4 shows how offsite replication can take place.

The DR team must balance the cost to recover against the cost of the disruption. The balancing point becomes the recovery time objective. Figure 23-5 illustrates the relationship between the cost of various recovery technologies and the provided recovery times.

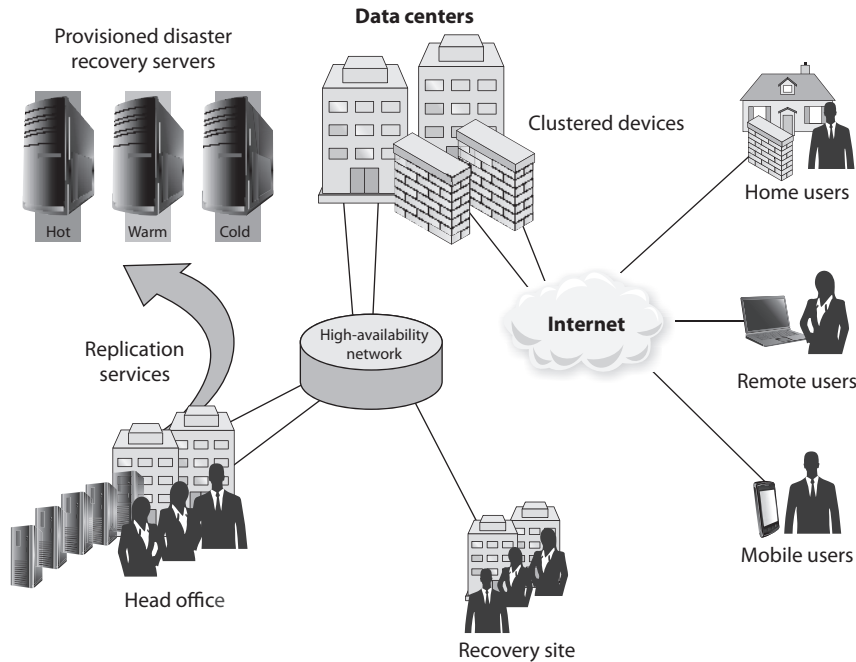


Figure 23-4 Offsite data replication for data recovery purposes

Choosing a Software Backup Facility

An organization needs to address several issues and ask specific questions when it is deciding upon a storage facility for its backup materials. The following list identifies just some of the issues that an organization needs to consider before committing to a specific vendor for this service:

- Can the media be accessed in the necessary timeframe?
- Is the facility closed on weekends and holidays, and does it only operate during specific hours of the day?
- Are the facility's access control mechanisms tied to an alarm and/or the police station?
- Does the facility have the capability to protect the media from a variety of threats?
- What is the availability of a bonded transport service?
- Are there any geographical environmental hazards such as floods, earthquakes, tornadoes, and so on that might affect the facility?
- Does the facility have a fire detection and suppression system?
- Does the facility provide temperature and humidity monitoring and control?
- What type of physical, administrative, and logical access controls are used?

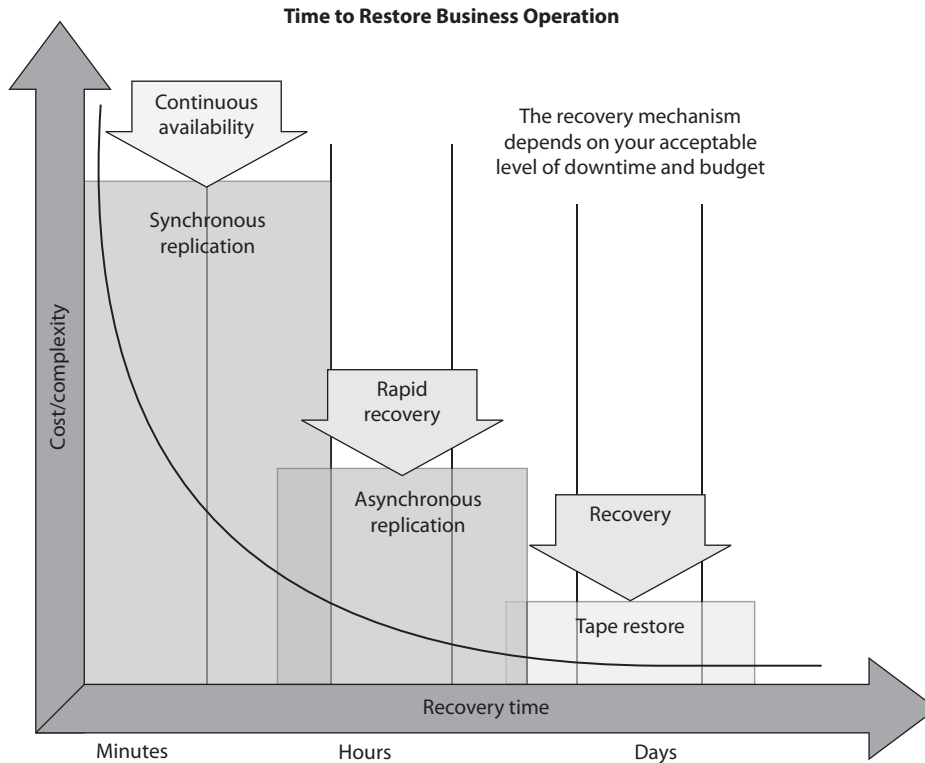


Figure 23-5 The criticality of data recovery will dictate the recovery solution.

The questions and issues that need to be addressed will vary depending on the type of organization, its needs, and the requirements of a backup facility.

Documentation

Documentation seems to be a dreaded task to most people, who will find many other tasks to take on to ensure they are not the ones stuck with documenting processes and procedures. However, without proper documentation, even an organization that does a terrific job of backing up data to an offsite facility will be scrambling to figure which backups it needs when a disaster hits.

Restoration of files can be challenging but restoring a whole environment that was swept away in a flood can be overwhelming, if not impossible. Procedures need to be documented because when they are actually needed, it will most likely be a chaotic and frantic atmosphere with a demanding time schedule. The documentation may need to include information on how to install images, configure operating systems and servers, and properly install utilities and proprietary software. Other documentation could include a calling tree, which outlines who should be contacted, in what order, and

Storing Business Continuity and Disaster Recovery Plans

Once the business continuity and disaster recovery plans are completed, where should they be stored? Should the organization have only one copy and keep it safely in a file cabinet next to Bob so that he feels safe? Nope. There should be two or three copies of these plans. One copy may be at the primary location, but the other copies should be at other locations in case the primary facility is destroyed. This reduces the risk of not having access to the plans when needed.

These plans should not be stored in a file cabinet, but rather in a fire-resistant safe. When they are stored offsite, they need to be stored in a way that provides just as much protection as the primary site would provide.

who is responsible for doing the calling. The documentation must also contain contact information for specific vendors, emergency agencies, offsite facilities, and any other entity that may need to be contacted in a time of need.

Most network environments evolve over time. Software is installed on top of other software, configurations are altered over the years to properly work in a unique environment, and service packs and patches are routinely installed to fix issues and update software. To expect one person or a group of people to go through all these steps during a crisis and end up with an environment that looks and behaves exactly like the original environment and in which all components work together seamlessly may be a lofty dream.

So, the dreaded task of documentation may be the saving grace one day. It is an essential piece of business, and therefore an essential piece in disaster recovery and business continuity. It is, therefore, important to make one or more roles responsible for proper documentation. As with all the items addressed in this chapter, simply saying “All documentation will be kept up to date and properly protected” is the easy part—saying and doing are two different things. Once the DR team identifies tasks that must be done, the tasks must be assigned to individuals, and those individuals have to be accountable. If these steps are not taken, the organization may have wasted a lot of time and resources defining these tasks, and still be in grave danger if a disaster occurs.

Human Resources

One of the resources commonly left out of the DR equation is people. An organization may restore its networks and critical systems and get business functions up and running, only to realize it doesn't know the answer to the question, “Who will take it from here?” The area of human resources is a critical component to any recovery and continuity process, and it needs to be fully thought out and integrated into the plan.

What happens if we have to move to an offsite facility that is 250 miles away? We cannot expect people to drive back and forth from home to work. Should we pay for temporary housing for the necessary employees? Do we have to pay their moving costs? Do we need to hire new employees in the area of the offsite facility? If so, what skill set

do we need from them? These are all important questions for the organization's senior leaders to answer.

If a large disaster takes place that affects not only the organization's facility but also surrounding areas, including housing, employees will be more worried about their families than their organization. Some organizations assume that employees will be ready and available to help them get back into production, when in fact they may need to be at home because they have responsibilities to their families.

Regrettably, some employees may be killed or severely injured in the disaster, and the organization should have plans in place to replace employees quickly through a temporary employment agency or a job recruiter. This is an extremely unfortunate scenario to contemplate, but it is part of reality. The team that considers all threats and is responsible for identifying practical solutions needs to think through all of these issues.

Organizations should already have *executive succession planning* in place. This means that if someone in a senior executive position retires, leaves the organization, or is killed, the organization has predetermined steps to carry out to ensure a smooth transition to that executive's replacement. The loss of a senior executive could tear a hole in the organization's fabric, creating a leadership vacuum that must be filled quickly with the right individual. The line-of-succession plan defines who would step in and assume responsibility for this role. Many organizations have "deputy" roles. For example, an organization may have a deputy CIO, deputy CFO, and deputy CEO ready to take over the necessary tasks if the CIO, CFO, or CEO becomes unavailable.

Often, larger organizations also have a policy indicating that two or more of the senior staff cannot be exposed to a particular risk at the same time. For example, the CEO and president cannot travel on the same plane. If the plane were to crash and both individuals were killed, then the company could face a leadership crisis. This is why you don't see the president of the United States and the vice president together too often. It is not because they don't like each other and thus keep their distance from each other. It is because there is a policy indicating that to protect the United States, its top leaders cannot be under the same risk at the same time.

Recovery Site Strategies

Disruptions, in BCP terms, are of three main types: nondisasters, disasters, and catastrophes. A *nondisaster* is a disruption in service that has significant but limited impact on the conduct of business processes at a facility. The solution could include hardware, software, or file restoration. A *disaster* is an event that causes the entire facility to be unusable for a day or longer. This usually requires the use of an alternate processing facility and restoration of software and data from offsite copies. The alternate site must be available to the organization until its main facility is repaired and usable. A *catastrophe* is a major disruption that destroys the facility altogether. This requires both a short-term solution, which would be an offsite facility, and a long-term solution, which may require rebuilding the original facility. Disasters and catastrophes are rare compared to nondisasters, thank goodness.

When dealing with disasters and catastrophes, an organization has three basic options: select a dedicated site that the organization owns and operates itself; lease a commercial facility, such as a "hot site" that contains all the equipment and data needed to quickly

restore operations; or enter into a formal agreement with another facility, such as a service bureau, to restore its operations. When choosing the right solution for its needs, the organization evaluates each alternative's ability to support its operations, to do it within an acceptable timeframe, and to have a reasonable cost.

An important consideration with third parties is their reliability, both in normal times and during an emergency. Their reliability can depend on considerations such as their track record, the extent and location of their supply inventory, and their access to supply and communication channels. Organizations should closely query the management of the alternative facility about such things as the following:

- How long will it take to recover from a certain type of incident to a certain level of operations?
- Will it give priority to restoring the operations of one organization over another after a disaster?
- What are its costs for performing various functions?
- What are its specifications for IT and security functions? Is the workspace big enough for the required number of employees?

To recover from a disaster that prevents or degrades use of the primary site temporarily or permanently, an organization must have an offsite backup facility available. Generally, an organization establishes contracts with third-party vendors to provide such services. The client pays a monthly fee to retain the right to use the facility in a time of need, and then incurs an activation fee when the facility actually has to be used. In addition, a daily or hourly fee is imposed for the duration of the stay. This is why service agreements for backup facilities should be considered a short-term solution, not a long-term solution.

It is important to note that most recovery site contracts do not promise to house the organization in need at a specific location, but rather promise to provide what has been contracted for somewhere within the organization's locale. On, and subsequent to, September 11, 2001, many organizations with Manhattan offices were surprised when they were redirected by their backup site vendor not to sites located in New Jersey (which were already full), but rather to sites located in Boston, Chicago, or Atlanta. This adds yet another level of complexity to the recovery process, specifically the logistics of transporting people and equipment to unplanned locations.

An organization can choose from three main types of leased or rented offsite recovery facilities:

- **Hot site** A facility that is fully configured and ready to operate within a few hours. All the necessary equipment is already installed and configured. In many cases, the remote data backup services are included, so the RPO can be down to an hour or even less. These sites are a good choice for an organization with a very small MTD. Of course, the organization should conduct regular tests (annually, at least) to ensure the site is functioning in the necessary state of readiness.

The hot site is, by far, the most expensive of the three types of offsite facilities. The organization has to pay for redundant hardware and software, in addition

to the expenses of the site itself. Organizations that use hot sites as part of their recovery strategy tend to limit them to mission-critical systems only.

- **Warm site** A facility that is usually partially configured with some equipment, such as HVAC, and foundational infrastructure components, but does not include all the hardware needed to restore mission-critical business functions. Staging a facility with duplicate hardware and computers configured for immediate operation is extremely expensive, so a warm site provides a less expensive alternate. These sites typically do not have data replicated to them, so backups would have to be delivered and restored onto the warm site systems after a disaster.

The warm site is the most widely used model. It is less expensive than a hot site, and can be up and running within a reasonably acceptable time period. It may be a better choice for organizations that depend on proprietary and unusual hardware and software, because they will bring their own hardware and software with them to the site after the disaster hits. Drawbacks, however, are that much of the equipment has to be procured, delivered to, and configured at the warm site after the fact, and testing will be more difficult. Thus, an organization may not be certain that it will in fact be able to return to an operating state within its RTO.

- **Cold site** A facility that supplies the basic environment, electrical wiring, HVAC, plumbing, and flooring but none of the equipment or additional services. A cold site is essentially an empty data center. It may take weeks to get the site activated and ready for work. The cold site could have equipment racks and dark fiber (fiber that does not have the circuit engaged) and maybe even desks. However, it would require the receipt of equipment from the client, since it does not provide any.

The cold site is the least expensive option, but takes the most time and effort to actually get up and functioning right after a disaster, as the systems and software must be delivered, set up, and configured. Cold sites are often used as backups for call centers, manufacturing plants, and other services that can be moved lock, stock, and barrel in one shot.

After a catastrophic loss of the primary facility, some organizations will start their recovery in a hot or warm site, and transfer some operations over to a cold site after the latter has had time to set up.

It is important to understand that the different site types listed here are provided by service bureaus. A *service bureau* is a company that has additional space and capacity to provide applications and services such as call centers. An organization pays a monthly subscription fee to a service bureau for this space and service. The fee can be paid for contingencies such as disasters and emergencies. You should evaluate the ability of a service bureau to provide services just as you would evaluate divisions within your own organization, particularly on matters such as its ability to alter or scale its software and hardware configurations or to expand its operations to meet the needs of a contingency.



NOTE Related to a service bureau is a *contingency supplier*; its purpose is to supply services and materials temporarily to an organization that is experiencing an emergency. For example, a contingency supplier might provide raw materials such as heating fuel or backup telecommunication services. In considering contingency suppliers, the BCP team should think through considerations such as the level of services and materials a supplier can provide, how quickly a supplier can ramp up to supply them, and whether the supplier shares similar communication paths and supply chains as the affected organization.

Most organizations use warm sites, which have some devices such as networking equipment, some computers and data storage, but very little else. These organizations usually cannot afford a hot site, and the extra downtime would not be considered detrimental. A warm site can provide a longer-term solution than a hot site. Organizations that decide to go with a cold site must be able to be out of operation for a week or two. The cold site usually includes power, raised flooring, climate control, and wiring.

The following provides a quick overview of the differences between offsite facilities.

Hot site advantages:

- Ready within hours or even minutes for operation
- Highly available
- Usually used for short-term solutions, but available for longer stays
- Recovery testing is easy

Hot site disadvantages:

- Very expensive
- Limited systems

Tertiary Sites

An organization may recognize the danger of the primary recovery site not being available when needed. This could be the case if the service provider assumes that not every customer will attempt to occupy the site at the same time, and then a major regional disaster affects more organizations than anticipated. It could also happen if a disaster affects the recovery site itself (e.g., fire, flood). Mitigating this risk could require a *tertiary site*, a backup recovery site just in case the primary is unavailable. The tertiary site is sometimes referred to as a “backup to the backup.” This is basically plan B if plan A does not work out. Obviously, this is a very expensive proposition, so its costs should be balanced with the risks it is intended to mitigate.

Warm and cold site advantages:

- Less expensive
- Available for longer timeframes because of the reduced costs
- Practical for proprietary hardware or software use

Warm and cold site disadvantages:

- Limited ability to perform recovery testing
- Resources for operations not immediately available

Reciprocal Agreements

Another approach to alternate offsite facilities is to establish a *reciprocal agreement* with another organization, usually one in a similar field or that has similar technological infrastructure. This means that organization A agrees to allow organization B to use its facilities if organization B is hit by a disaster, and vice versa. This is a cheaper way to go than the other offsite choices, but it is not always the best choice. Most environments are maxed out pertaining to the use of facility space, resources, and computing capability. To allow another organization to come in and work out of the same shop could prove to be detrimental to both organizations. Whether it can assist the other organization while tending effectively to its own business is an open question. The stress of two organizations working in the same environment could cause tremendous levels of tension. If it did work out, it would only provide a short-term solution. Configuration management could be a nightmare. Does the other organization upgrade to new technology and retire old systems and software? If not, one organization's systems may become incompatible with those of the other.

If your organization allows another organization to move into its facility and work from there, you may have a solid feeling about your friend, the CEO, but what about all of her employees, whom you do not know? The mixing of operations could introduce many security issues. Now you have a new subset of people who may need to have privileged and direct access to your resources in the shared environment. Close attention needs to be paid when assigning these other people access rights and permissions to your critical assets and resources, if they need access at all. Careful testing is recommended to see if one organization or the other can handle the extra loads.

Offsite Location

When choosing a backup facility, it should be far enough away from the original site so that one disaster does not take out both locations. In other words, it is not logical to have the backup site only a few miles away if the organization is concerned about tornado damage, because the backup site could also be affected or destroyed. There is a rule of thumb that suggests that alternate facilities should be, at a bare minimum, at least 5 miles away from the primary site, while 15 miles is recommended for most low-to-medium critical environments, and 50 to 200 miles is recommended for critical operations, to give maximum protection in cases of regional disasters.

Reciprocal agreements have been known to work well in specific businesses, such as newspaper printing. These businesses require very specific technology and equipment that is not available through any subscription service. These agreements follow a “you scratch my back and I’ll scratch yours” mentality. For most other organizations, reciprocal agreements are generally, at best, a secondary option for disaster protection. The other issue to consider is that these agreements are usually not enforceable because they’re not written in legally binding terms. This means that although organization A said organization B could use its facility when needed, when the need arises, organization A may not have a legal obligation to fulfill this promise. However, there are still many organizations who do opt for this solution either because of the appeal of low cost or, as noted earlier, because it may be the only viable solution in some cases.

Organizations that have a reciprocal agreement need to address the following important issues before a disaster hits:

- How long will the facility be available to the organization in need?
- How much assistance will the staff supply in integrating the two environments and ongoing support?
- How quickly can the organization in need move into the facility?
- What are the issues pertaining to interoperability?
- How many of the resources will be available to the organization in need?
- How will differences and conflicts be addressed?
- How does change control and configuration management take place?
- How often can exercising and testing take place?
- How can critical assets of both organizations be properly protected?

A variation on a reciprocal agreement is a consortium, or *mutual aid agreement*. In this case, more than two organizations agree to help one another in case of an emergency. Adding multiple organizations to the mix, as you might imagine, can make things even more complicated. The same concerns that apply with reciprocal agreements apply here, but even more so. Organizations entering into such agreements need to formally and legally document their mutual responsibilities in advance. Interested parties, including the legal and IT departments, should carefully scrutinize such accords before the organization signs onto them.

Redundant Sites

Some organizations choose to have a *redundant site*, or mirrored site, meaning one site is equipped and configured exactly like the primary site, which serves as a redundant environment. The business-processing capabilities between the two sites can be completely synchronized. A redundant site is owned by the organization and mirrors the original production environment. A redundant site has clear advantages: it has full availability, is ready to go at a moment’s notice, and is under the organization’s complete control. This is, however, one of the most expensive backup facility options, because a full

environment must be maintained even though it usually is not used for regular production activities until after a disaster takes place that triggers the relocation of services to the redundant site. But “expensive” is relative here. If a company would lose a million dollars if it were out of business for just a few hours, the loss potential would override the cost of this option. Many organizations are subjected to regulations that dictate they must have redundant sites in place, so expense is not a matter of choice in these situations.



EXAM TIP A *hot site* is a subscription service. A *redundant site*, in contrast, is a site owned and maintained by the organization, meaning the organization does not pay anyone else for the site. A redundant site might be “hot” in nature, meaning it is ready for production quickly. However, the CISSP exam differentiates between a hot site (a subscription service) and a redundant site (owned by the organization).

Another type of facility-backup option is a *rolling hot site*, or mobile hot site, where the back of a large truck or a trailer is turned into a data processing or working area. This is a portable, self-contained data facility. The trailer has the necessary power, telecommunications, and systems to do some or all of the processing right away. The trailer can be brought to the organization’s parking lot or another location. Obviously, the trailer has to be driven over to the new site, the data has to be retrieved, and the necessary personnel have to be put into place.

Another, similar solution is a prefabricated building that can be easily and quickly put together. Military organizations and large insurance companies typically have rolling hot sites or trucks preloaded with equipment because they often need the flexibility to quickly relocate some or all of their processing facilities to different locations around the world depending on where the need arises.

It is best if an organization is aware of all available options for hardware and facility backups to ensure it makes the best decision for its specific business and critical needs.

Multiple Processing Sites

Another option for organizations is to have *multiple processing sites*. An organization may have ten different facilities throughout the world, which are connected with specific technologies that could move all data processing from one facility to another in a matter of seconds when an interruption is detected. This technology can be implemented within the organization or from one facility to a third-party facility. Certain service providers provide this type of functionality to their customers. So if an organization’s data processing is interrupted, all or some of the processing can be moved to the service provider’s servers.

Availability

We close this section on recovery strategies by considering the nondisasters to which we referred earlier. These are the incidents that may not require evacuation of personnel or facility repairs but that can still have a significant detrimental effect on the ability of the organization to execute its mission. We want our systems and services to be available all

the time, no matter what. However, we all realize this is just not possible. *Availability* can be defined as the portion of the time that a system is operational and able to fulfill its intended purpose. But how can we ensure the availability of the systems and services on which our organizations depend?

High Availability

High availability (HA) is a combination of technologies and processes that work together to ensure that some specific thing is up and running most of the time. The specific thing can be a database, a network, an application, a power supply, and so on. Service providers have *service level agreements (SLAs)* with their customers that outline the amount of uptime the service providers promise to provide. For example, a hosting company can promise to provide 99 percent uptime for Internet connectivity. This means the company is guaranteeing that at least 99 percent of the time, the Internet connection you purchase from it will be up and running. It also means that you can experience up to 3.65 days a year (or 7.2 hours per month) of downtime and it won't be a violation of the SLA. Increase that to 99.999 percent (referred to as "five nines") uptime and the allowable downtime drops to 5.26 seconds per year, but the price you pay for service goes through the roof.



NOTE HA is in the eye of the beholder. For some organizations or systems, an SLA of 90 percent ("one nine") uptime and its corresponding potential 36+ days of downtime a year is perfectly fine, particularly for organizations that are running on a tight budget. Other organizations require "nine nines" or 99.999999 percent availability for mission-critical systems. You have to balance the cost of HA with the loss you're trying to mitigate.

Just because a service is available doesn't necessarily mean that it is operating acceptably. Suppose your company's high-speed e-commerce server gets infected with a bitcoin miner that drives CPU utilization close to 100 percent. Technically, the server is available and will probably be able to respond to customer requests. However, response times will likely be so lengthy that many of your customers will simply give up and go shop somewhere else. The service is available, but its quality is unacceptable.

Quality of Service

Quality of service (QoS) defines minimum acceptable performance characteristics of a particular service. For example, for the e-commerce server example, we could define parameters like response time, CPU utilization, or network bandwidth utilization, depending on how the service is being provided. SLAs may include one or more specifications for QoS, which allows service providers to differentiate classes of service that are prioritized for different clients. During a disaster, the available bandwidth on external links may be limited, so the affected organization could specify different QoS for its externally facing systems. For example, the e-commerce company in our example could determine the minimum data rate to keep its web presence available to customers and

specify that as the minimum QoS rate at the expense of, say, its e-mail or Voice over Internet Protocol (VoIP) traffic.

To provide HA and meet stringent QoS requirements, the hosting company has to have a long list of technologies and processes that provide redundancy, fault tolerance, and failover capabilities. *Redundancy* is commonly built into the network at a routing protocol level. The routing protocols are configured such that if one link goes down or gets congested, traffic is automatically routed over a different network link. An organization can also ensure that it has redundant hardware available so that if a primary device goes down, the backup component can be swapped out and activated.

If a technology has a *failover* capability, this means that if there is a failure that cannot be handled through normal means, then processing is “switched over” to a working system. For example, two servers can be configured to send each other “heartbeat” signals every 30 seconds. If server A does not receive a heartbeat signal from server B after 40 seconds, then all processes are moved to server A so that there is no lag in operations. Also, when servers are *clustered*, an overarching piece of software monitors each server and carries out load balancing. If one server within the cluster goes down, the clustering software stops sending it data to process so that there are no delays in processing activities.

Fault Tolerance and System Resilience

Fault tolerance is the capability of a technology to continue to operate as expected even if something unexpected takes place (a fault). If a database experiences an unexpected glitch, it can roll back to a known-good state and continue functioning as though nothing bad happened. If a packet gets lost or corrupted during a TCP session, the TCP protocol will resend the packet so that system-to-system communication is not affected. If a disk within a RAID system gets corrupted, the system uses its parity data to rebuild the corrupted data so that operations are not affected.

Although the terms fault tolerance and resilience are often used synonymously, they mean subtly different things. Fault tolerance means that when a fault happens, there's a system in place (a backup or redundant one) to ensure services remain uninterrupted. *System resilience* means that the system continues to function, albeit in a degraded fashion, when a fault is encountered. Think of it as the difference between having a spare tire for your car and having run-flat tires. The spare tire provides fault tolerance in that it enables you to recover (fairly) quickly from a flat tire and be on your way. Run-flat tires allow you to continue to drive your car (albeit slower) if you run over a nail on the road. A resilient system is fault tolerant, but a fault tolerant one may not be resilient.

High Availability in Disaster Recovery

Redundancy, fault tolerance, resilience, and failover capabilities increase the reliability of a system or network, where *reliability* is the probability that a system performs the necessary function for a specified period under defined conditions. High reliability allows for high availability, which is a measure of its readiness. If the probability of a system performing as expected under defined conditions is low, then the availability for this system cannot be high. For a system to have the characteristic of high availability,

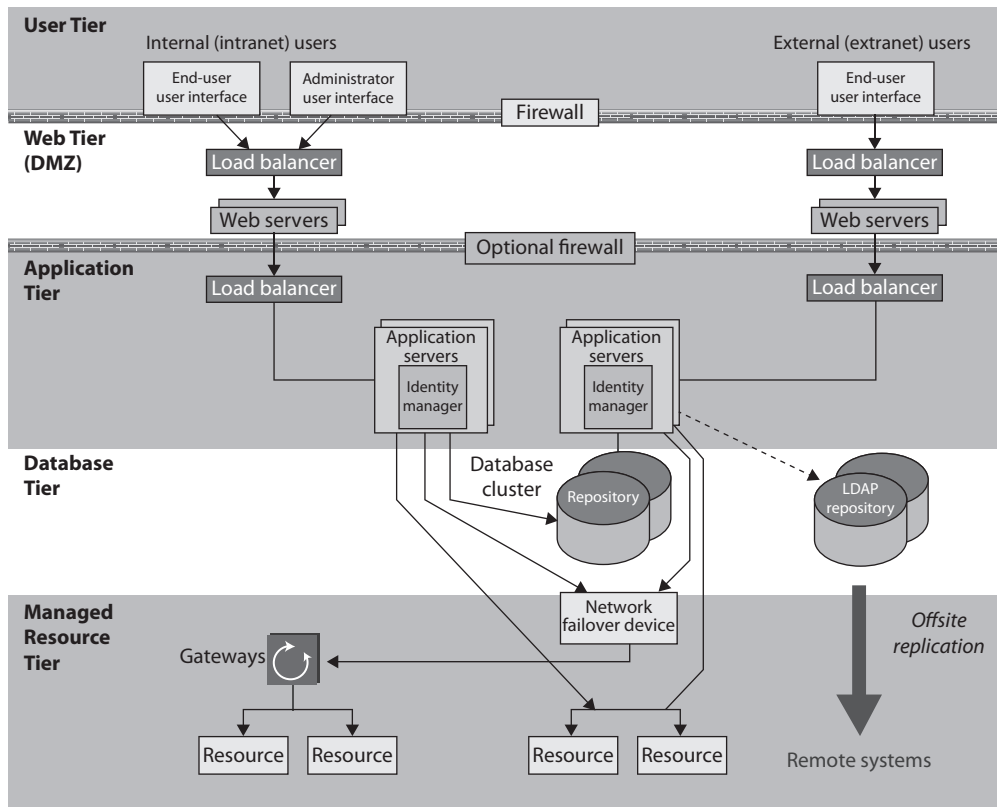


Figure 23-6 High-availability technologies

then high reliability must be in place. Figure 23-6 illustrates where load balancing, clustering, failover devices, and replication commonly take place in a network architecture.

Remember that data restoration (RPO) requirements can be different from processing restoration (RTO) requirements. Data can be restored through backup tapes, electronic vaulting, or synchronous or asynchronous replication. Processing capabilities can be restored through clustering, load balancing, redundancy, and failover technologies. If the results of the BCP team's BIA indicate that the RPO value is two days, then the organization can use tape backups. If the RPO value is one minute, then synchronous replication needs to be in place. If the BIA indicates that the RTO value is three days, then redundant hardware can be used. If the RTO value is one minute, then clustering and load balancing should be used.

HA and disaster recovery are related concepts. HA technologies and processes are commonly put into place so that if a disaster does take place, either the critical functions are likelier to remain available or the delay of getting them back online and running is low.

Many IT and security professionals usually think of HA only in technology terms, but remember that there are many things that an organization needs to have available to keep functioning. Availability of each of the following items must be thought through and planned:

- Facility (cold, warm, hot, redundant, rolling, reciprocal sites)
- Infrastructure (redundancy, fault tolerance)
- Storage (SAN, cloud)
- Server (clustering, load balancing)
- Data (backups, online replication)
- Business processes
- People



NOTE Virtualization and cloud computing are covered in Chapter 7. We will not go over those technologies again in this chapter, but know that the use of these technologies has drastically increased in the realm of business continuity and disaster recovery planning solutions.

Disaster Recovery Processes

Recovering from a disaster begins way before the event occurs. It starts by anticipating threats and developing goals that support the organization's continuity of operations. If you do not have established goals, how do you know when you are done and whether your efforts were actually successful? Goals are established so everyone knows the ultimate objectives. Establishing goals is important for any task, but especially for business continuity and disaster recovery plans. The definition of the goals helps direct the proper allocation of resources and tasks, supports the development of necessary strategies, and assists in financial justification of the plans and program overall. Once the goals are set, they provide a guide to the development of the actual plans themselves. Anyone who has been involved in large projects that entail many small, complex details knows that at times it is easy to get off track and not actually accomplish the major goals of the project. Goals are established to keep everyone on track and to ensure that the efforts pay off in the end.

Great—we have established that goals are important. But the goal could be, “Keep the company in business if an earthquake hits.” That’s a good goal, but it is not overly useful without more clarity and direction. To be useful, a goal must contain certain key information, such as the following:

- **Responsibility** Each individual involved with recovery and continuity should have their responsibilities spelled out in writing to ensure a clear understanding in a chaotic situation. Each task should be assigned to the individual most logically situated to handle it. These individuals must know what is expected of them, which is done through training, exercises, communication, and documentation. So, for example, instead of just running out of the building screaming, an individual must know that he is responsible for shutting down the servers before he can run out of the building screaming.

- **Authority** In times of crisis, it is important to know who is in charge. Teamwork is important in these situations, and almost every team does much better with an established and trusted leader. Such leaders must know that they are expected to step up to the plate in a time of crisis and understand what type of direction they should provide to the rest of the employees. Everyone else must recognize the authority of these leaders and respond accordingly. Clear-cut authority will aid in reducing confusion and increasing cooperation.
- **Priorities** It is extremely important to know what is critical versus what is merely nice to have. Different departments provide different functionality for an organization. The critical departments must be singled out from the departments that provide functionality that the organization can live without for a week or two. It is necessary to know which department must come online first, which second, and so on. That way, the efforts are made in the most useful, effective, and focused manner. Along with the priorities of departments, the priorities of systems, information, and programs must be established. It may be necessary to ensure that the database is up and running before working to bring the web servers online. The general priorities must be set by management with the help of the different departments and IT staff.
- **Implementation and testing** It is great to write down very profound ideas and develop plans, but unless they are actually carried out and tested, they may not add up to a hill of beans. Once a disaster recovery plan is developed, it actually has to be put into action. It needs to be documented and stored in places that are easily accessible in times of crisis. The people who are assigned specific tasks need to be taught and informed how to fulfill those tasks, and dry runs must be done to walk people through different situations. The exercises should take place at least once a year, and the entire program should be continually updated and improved.



NOTE We address various types of tests, such as walkthrough, tabletop, simulation, parallel, and full interruption, later in this chapter.

According to the U.S. Federal Emergency Management Agency (FEMA), 90 percent of small businesses that experience a disaster and are unable to restore operations within five days will fail within the following year. Not being able to bounce back quickly or effectively by setting up shop somewhere else can make a company lose business and, more importantly, its reputation. In such a competitive world, customers have a lot of options. If one company is not prepared to bounce back after a disruption or disaster, customers may go to another vendor and stay there.

The biggest effect of an incident, especially one that is poorly managed or that was preventable, is on an organization's reputation or brand. This can result in a considerable and even irreparable loss of trust by customers and clients. On the other hand, handling an incident well, or preventing great damage through smart, preemptive measures, can enhance the reputation of, or trust in, an organization.

The *disaster recovery plan (DRP)* should address in detail all of the topics we have covered so far. The actual format of the DRP will depend on the environment, the goals of the plan, priorities, and identified threats. After each of those items is examined and documented, the topics of the plan can be divided into the necessary categories.

Response

The first question the DRP should answer is, “What constitutes a disaster that would trigger this plan?” Every leader within an organization (and, ideally, everyone else too) should know the answer. Otherwise, precious time is lost notifying people who should’ve self-activated as soon as the incident occurred, a delay that could cost lives or assets. Examples of clear-cut disasters that would trigger a response are loss of power exceeding ten minutes, flooding in the facility, or terrorist attack against or near the site.

Every DRP is different, but most follow a familiar sequence of events:

1. Declaration of disaster
2. Activation of the DR team
3. Internal communications (ongoing from here on out)
4. Protection of human safety (e.g., evacuation)
5. Damage assessment
6. Execution of appropriate system-specific DRPs (each system and network should have its own DRP)
7. Recovery of mission-critical business processes/functions
8. Recovery of all other business processes/functions

Personnel

The DRP needs to define several different teams that should be properly trained and available if a disaster hits. Which types of teams an organization needs depends upon the organization. The following are some examples of teams that an organization may need to construct:

- Damage assessment team
- Recovery team
- Relocation team
- Restoration team
- Salvage team
- Security team

The DR coordinator should have an understanding of the needs of the organization and the types of teams that need to be developed and trained. Employees should be assigned to the specific teams based on their knowledge and skill set. Each team needs

to have a designated leader, who will direct the members and their activities. These team leaders will be responsible not only for ensuring that their team's objectives are met but also for communicating with each other to make sure each team is working in parallel phases.

The purpose of the *recovery team* should be to get whatever systems are still operable back up and running as quickly as possible to reduce business disruptions. Think of them as the medics whose job is to stabilize casualties until they can be transported to the hospital. In this case, of course, there is no hospital for information systems, but there may be a recovery site. Getting equipment and people there in an orderly fashion should be the job of the *relocation team*. The *restoration team* should be responsible for getting the alternate site into a working and functioning environment, and the *salvage team* should be responsible for starting the recovery of the original site. Both teams must know how to do many tasks, such as install operating systems, configure workstations and servers, string wire and cabling, set up the network and configure networking services, and install equipment and applications. Both teams must also know how to restore data from backup facilities and how to do so in a secure manner, one that ensures the availability, integrity, and confidentiality of the system and data.

The DRP must outline the specific teams, their responsibilities, and notification procedures. The plan must indicate the methods that should be used to contact team leaders during business hours and after business hours.

Communications

The purpose of the emergency communications plan that is part of the overall DRP is to ensure that everyone knows what to do at all times and that the DR team remains synchronized and coordinated. This all starts with the DR plan itself. As stated previously, copies of the DRP need to be kept in one or more locations other than the primary site, so that if the primary site is destroyed or negatively affected, the plan is still available to the teams. It is also critical that different formats of the plan be available to the teams, including both electronic and paper versions. An electronic version of the plan is not very useful if you don't have any electricity to run a computer.

In addition to having copies of the recovery documents located at their offices and homes, key individuals should have easily accessible versions of critical procedures and call tree information. One simple way to accomplish the latter is to publish a call tree on cards that can be affixed to personnel badges or kept in a wallet. In an emergency situation, valuable minutes are better spent responding to an incident than looking for a document or having to wait for a laptop to power up. Of course, the call tree is only as effective as it is accurate and up to date, so verifying it periodically is imperative.

One limitation of call trees is that they are point to point, which means they're typically good for getting the word out, but not so much for coordinating activities. Group text messages work better, but only in the context of fairly small and static groups. Many organizations have group chat solutions, but if those rely on the organization's servers, they may be unavailable during a disaster. It is a good idea, then, to establish

a communications platform that is completely independent of the organizational infrastructure. Solutions like Slack and Mattermost offer a free service that is typically sufficient to keep most organizations connected in emergencies. The catch, of course, is that everyone needs to have the appropriate client installed on their personal devices and know when and how to connect. Training and exercises are the keys to successful execution of any plan, and the communications plan is no exception.



NOTE An organization may need to solidify communications channels and relationships with government officials and emergency response groups. The goal of this activity is to solidify proper protocol in case of a city- or region-wide disaster. During the BIA phase, the DR team should contact local authorities to elicit information about the risks of its geographical location and how to access emergency zones. If the organization has to perform DR, it may need to contact many of these emergency response groups.

PACE Communications Plans

The U.S. armed forces routinely develop Primary, Alternate, Contingency, and Emergency (PACE) communications plans. The PACE plan outlines the different capabilities that exist and aligns them into these four categories based on their ability to meet defined information exchange requirements. Each category is defined here:

- **Primary** The normal or expected capability that is used to achieve the objective.
- **Alternate** A fully satisfactory capability that can be used to achieve the objective with minimal impact to the operation or exercise. This capability is used when the Primary capability is unavailable.
- **Contingency** A workable capability that can be used to achieve the objective. This capability may not be as fast or easy as the Primary or Alternate but is capable of achieving the objective with an acceptable amount of time and effort. This capability is used when the Primary and the Alternate capabilities are unavailable.
- **Emergency** This is the last-resort capability and typically may involve significantly more time and effort than any of the other capabilities. This capability should be used only when the Primary, Alternate, and Contingency capabilities are unavailable.

The PACE plan includes redundant communications capabilities and specifies the order in which the organization will employ the capabilities when communication outages occur.

Assessment

A role, or a team, needs to be created to carry out a *damage assessment* once a disaster has taken place. The assessment procedures should be properly documented in the DRP and include the following steps:

- Determine the cause of the disaster.
- Determine the potential for further damage.
- Identify the affected business functions and areas.
- Identify the level of functionality for the critical resources.
- Identify the resources that must be replaced immediately.
- Estimate how long it will take to bring critical functions back online.

After the damage assessment team collects and assesses this information, the DR coordinator identifies which teams need to be called to action and which system-specific DRPs need to be executed (and in what order). The DRP should specify activation criteria for the different teams and system-specific DRPs. After the damage assessment, if one or more of the situations outlined in the criteria have taken place, then the DR team is moved into restoration mode.

Different organizations have different activation criteria because business drivers and critical functions vary from organization to organization. The criteria may comprise some or all of the following elements:

- Danger to human life
- Danger to state or national security
- Damage to facility
- Damage to critical systems
- Estimated value of downtime that will be experienced

Restoration

Once the damage assessment is completed, various teams are activated, which signals the organization's entry into the *restoration phase*. Each team has its own tasks—for example, the facilities team prepares the offsite facility (if needed), the network team rebuilds the network and systems, and the relocation team starts organizing the staff to move into a new facility.

The restoration process needs to be well organized to get the organization up and running as soon as possible. This is much easier to state in a book than to carry out in reality, which is why written procedures are critical. The critical functions and their resources would already have been identified during the BIA, as discussed earlier in this chapter (with a simplistic example provided in Table 23-1). These are the functions that the teams need to work together on restoring first.

Many organizations create templates during the DR plan development stage. These templates are used by the different teams to step them through the necessary phases and to document their findings. For example, if one step could not be completed until new systems were purchased, this should be indicated on the template. If a step is partially completed, this should be documented so the team does not forget to go back and finish that step when the necessary part arrives. These templates keep the teams on task and also quickly tell the team leaders about the progress, obstacles, and potential recovery time.



NOTE Examples of possible templates can be found in NIST Special Publication 800-34, Revision 1, *Contingency Planning Guide for Federal Information Systems*, which is available online at <https://csrc.nist.gov/publications/detail/sp/800-34/rev-1/final>.

An organization is not out of an emergency state until it is back in operation at the original primary site or at a new site that was constructed to replace the primary original one, because the organization is always vulnerable while operating in a backup facility. Many logistical issues need to be considered as to when an organization should return from the alternate site to the primary one. The following lists a few of these issues:

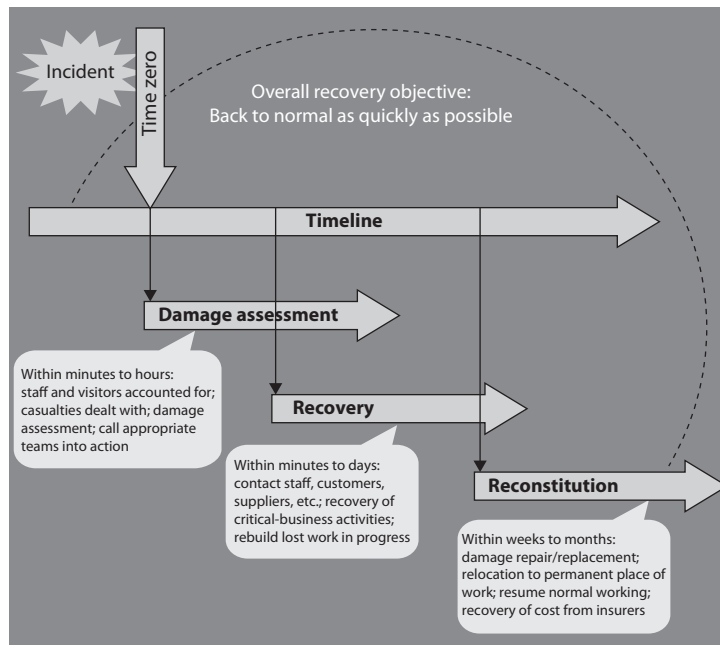
- Ensuring the safety of employees
- Ensuring an adequate environment is provided (power, facility infrastructure, water, HVAC)
- Ensuring that the necessary equipment and supplies are present and in working order
- Ensuring proper communications and connectivity methods are working
- Properly testing the new environment

Once the coordinator, management, and salvage team sign off on the readiness of the primary site, the salvage team should carry out the following steps:

- Back up data from the alternate site and restore it within the primary site.
- Carefully terminate contingency operations.
- Securely transport equipment and personnel to the primary site.

The least critical functions should be moved back first, so if there are issues in network configurations or connectivity, or important steps were not carried out, the critical operations of the organization are not negatively affected. Why go through the trouble of moving the most critical systems and operations to a safe and stable alternate site, only to return them to a main site that is untested? Let the less critical departments act as the

canary in the coal mine. If they survive, then move the more critical components of the organization to the main site.



Training and Awareness

Training your DR team on the execution of a DRP is critical for at least three reasons. First, it allows you to validate that the plan will actually work. If your DR team is doing a walkthrough exercise in response to a fictitious scenario, you'll find out very quickly whether the plan would work or not. If it doesn't work in a training event when the stress level and stakes are low, then there is no chance it would work in a real emergency.

Another reason to train is to ensure that everyone knows what they're supposed to do, when, where, and how. Disasters are stressful, messy affairs and key people may not be thinking clearly. It is important for them to have a familiar routine to fall back on. In a perfect world, you would train often enough for your team to develop "muscle memory" that allows them to automatically do the right things without even thinking.

Lastly, training can help establish that you are exercising due care. This could keep you out of legal trouble in the aftermath of a disaster, particularly if people end up getting hurt. A good plan and evidence of a trained workforce can go a long way to reduce liability if regulators or other investigators come knocking. As always, consult your attorneys to ensure you are meeting all applicable legal and regulatory obligations.

When thinking of training and "muscle memory," you should also consider everyone else in the organization that is not part of the DR team. You want all your staff to have an awareness of the major things they need to do to support DR. This is why many of us conduct fire drills in our facilities: to ensure everyone knows how to get out of the

building and where to assemble if we ever face this particular kind of disaster. There are many types of DR awareness events you can run, but you should at least consider three types of responses that everyone should be aware of: evacuations (e.g., for fires or explosives), shelter-in-place (e.g., for tornadoes or active shooters), and remain-at-home (e.g., for overnight flooding).

Lessons Learned

As mentioned on the first page of this chapter, no battle plan ever survived first contact with the enemy. When you try to execute your DRP in a real disaster, you will find the need to disregard parts of it, make on-the-fly changes to others, and faithfully execute the rest. This is why you should incorporate lessons learned from any actual disasters and actual responses. The DR team should perform a “postmortem” on the response and ensure that necessary changes are made to plans, contracts, personnel, processes, and procedures.

Military organizations collect lessons learned in two steps. The first steps, called a *hotwash*, is a hasty one that happens right after the event is concluded (i.e., restoration is completed). The term comes from the military practice of dousing rifles with very hot water immediately after an engagement to quickly get the worst grit and debris off their weapons. The reason you want to conduct a hotwash right away is that memories will be freshest right after restoring the systems. The idea is not necessarily to figure out how to fix anything, but rather to quickly list as many things that went well or poorly as possible before participants start to forget them.

The second event at which lessons learned are collected in the military is much more deliberate. An after-action review (AAR) happens several days after completion of the DR and allows participants to think things through and start formulating possible ways to do better in the future. The AAR facilitator, ideally armed with the notes from the hotwash, presents each issue that was recorded (good or bad), a brief discussion of it, and then opens the floor for recommendations. Keep in mind that since you’re dealing with things that went well or poorly, sometimes the group recommendation will be to “sustain” the issue or, in other words, keep doing things the same way in the future. More frequently, however, there are at least minor tweaks that can improve future performance.

Testing Disaster Recovery Plans

The disaster recovery plan should be tested regularly because environments continually change. Interestingly, many organizations are moving away from the concept of “testing,” because a test naturally leads to a pass or fail score, and in the end, that type of score is not very productive. Instead, many organizations are adopting the concept of “exercises,” which appear less stressful, better focused, and ultimately more productive to the participants. Each time the DRP is exercised or tested, improvements and efficiencies are generally uncovered, yielding better and better results over time. The responsibility of establishing periodic exercises and the maintenance of the plan should be assigned to a specific person or persons who will have overall ownership responsibilities for the disaster recovery initiatives within the organization.

The maintenance of the DRP should be incorporated into change management procedures. That way, any changes in the environment are reflected in the plan itself.

Tests and disaster recovery exercises should be performed at least once a year. An organization should have no real confidence in a developed plan until it has actually been tested. Exercises prepare personnel for what they may face and provide a controlled environment to learn the tasks expected of them. These exercises also point out issues to the planning team and management that may not have been previously thought about and addressed as part of the planning process. The exercises, in the end, demonstrate whether an organization can actually recover after a disaster.

The exercise should have a predetermined scenario that the organization may indeed be faced with one day. Specific parameters and a scope of the exercise must be worked out before sounding the alarms. The team of testers must agree upon what exactly is getting tested and how to properly determine success or failure. The team must agree upon the timing and duration of the exercise, who will participate in the exercise, who will receive which assignments, and what steps should be taken. Also, the team needs to determine whether hardware, software, personnel, procedures, and communications lines are going to be tested and whether it is all or a subset of these resources that will be included in the event. If the test will include moving some equipment to an alternate site, then transportation, extra equipment, and alternate site readiness must be addressed and assessed.

Most organizations cannot afford to have these exercises interrupt production or productivity, so the exercises may need to take place in sections or at specific times, which will require logistical planning. Written exercise plans should be developed that will test for specific weaknesses in the overall DRP. The first exercises should not include all employees, but rather a small representative sample of the organization. This allows both the planners and the participants to refine the plan. It also allows each part of the organization to learn its roles and responsibilities. Then, larger exercises can take place so overall operations will not be negatively affected.

The people conducting these exercises should expect to encounter problems and mistakes. After all, identifying potential problems and mistakes is why they are conducting the exercises in the first place. An organization would rather have employees make mistakes during an exercise so they can learn from them and perform their tasks more effectively during a real disaster.



NOTE After a disaster, telephone service may not be available. For communications purposes, alternatives should be in place, such as mobile phones or hand-held radios.

A few different types of exercises and tests can be used, each with its own pros and cons. The following sections explain the different types of assessment events.

Checklist Test

In this type of test, copies of the DRP are distributed to the different departments and functional areas for review. This enables each functional manager to review the plan

and indicate if anything has been left out or if some approaches should be modified or deleted. This method ensures that nothing is taken for granted or omitted, as might be the case in a single-department review. Once the departments have reviewed their copies and made suggestions, the planning team then integrates those changes into the master plan.



NOTE The checklist test is also called the desk check test.

Structured Walkthrough Test

In this test, representatives from each department or functional area come together and go over the plan to ensure its accuracy. The group reviews the objectives of the plan; discusses the scope and assumptions of the plan; reviews the organization's reporting structure; and evaluates the testing, maintenance, and training requirements described. This gives the people responsible for making sure a disaster recovery happens effectively and efficiently an opportunity to review what has been decided upon and what is expected of them.

The group walks through different scenarios of the plan from beginning to end to make sure nothing was left out. This also raises the awareness of team members about the recovery procedures.

Tabletop Exercises

Tabletop exercises (TTXs) may or may not happen at a tabletop, but they do not involve a technical control infrastructure. TTXs can happen at an executive level (e.g., C-suite) or at a team level (e.g., SOC), or anywhere in between. The idea is usually to test procedures and ensure they actually do what they're intended to and that everyone knows their role in responding to a disaster. TTXs require relatively few resources apart from deliberate planning by qualified individuals and the undisturbed time and attention of the participants.

After determining the goals of the exercise and vetting them with the senior leadership of the organization, the planning team develops a scenario that touches on the important aspects of the response plan. The idea is normally not to cover every contingency, but to ensure the DR team is able to respond to the likeliest and/or most dangerous scenarios. As they develop the exercise, the planning team considers branches and sequels at every point in the scenario. A *branch* is a point in which the participants may choose one of multiple approaches to respond. If the branches are not carefully managed and controlled, the TTX could wander into uncharted and unproductive directions. Conversely, a *sequel* is a follow-on to a given action in the response. For instance, as part of the response, the strategic communications team may issue statements to the news media. A sequel to that could involve a media outlet challenging the statement, which in turn would require a response by the team. Like branches, sequels must be used carefully to keep the exercise on course. Senior leadership support and good scenario development are critical ingredients to attract and engage the right participants. Like any contest, a TTX is only as good as the folks who show up to play.



EXAM TIP Tabletop exercises are also called read-through exercises.

Simulation Test

This type of test takes a lot more planning and people. In this situation, all employees who participate in operational and support functions, or their representatives, come together to practice executing the disaster recovery plan based on a specific scenario. The scenario is used to test the reaction of each operational and support representative. Again, this is done to ensure specific steps were not left out and that certain threats were not overlooked. It raises the awareness of the people involved.

The exercise includes only those materials that will be available in an actual disaster, to portray a more realistic environment. The simulation test continues up to the point of actual relocation to an offsite facility and actual shipment of replacement equipment.

Parallel Test

In a parallel test, some systems are moved to the alternate site and processing takes place. The results are compared with the regular processing that is done at the original site. This ensures that the specific systems can actually perform adequately at the alternate offsite facility and points out any tweaking or reconfiguring that is necessary.

Full-Interruption Test

This type of test is the most intrusive to regular operations and business productivity. The original site is actually shut down, and processing takes place at the alternate site. The recovery team fulfills its obligations in preparing the systems and environment for the alternate site. All processing is done only on devices at the alternate offsite facility.

This is a full-blown exercise that takes a lot of planning and coordination, but it can reveal many holes in the plan that need to be fixed before an actual disaster hits. Full-interruption tests should be performed only after all other types of tests have been successful. They are the riskiest type and can impact the business in very serious and devastating ways if not managed properly; therefore, senior management approval needs to be obtained prior to performing full-interruption tests.

The type of organization and its goals will dictate what approach to the training exercise is most effective. Each organization may have a different approach and unique aspects. If detailed planning methods and processes are going to be taught, then specific training may be required rather than general training that provides an overview. Higher-quality training will result in an increase in employee interest and commitment.

During and after each type of test, a record of the significant events should be documented and reported to management so it is aware of all outcomes of the test.

Other Types of Training

Other types of training that employees need in addition to disaster recovery training include first aid and cardiac pulmonary resuscitation (CPR), how to properly use a fire extinguisher, evacuation routes and crowd control methods, emergency communications procedures, and how to properly shut down equipment in different types of disasters.

The more technical employees may need training on how to redistribute network resources and how to use different telecommunications lines if the main one goes down. They may need to know about redundant power supplies and be trained and tested on the procedures for moving critical systems from one power supply to the next.

Business Continuity

When a disaster strikes, ensuring that the organization is able to continue its operations requires more than simply restoring data from backups. Also necessary are the detailed procedures that outline the activities to keep the critical systems available and ensure that operations and processing are not interrupted. Business continuity planning defines what should take place during and after an incident. Actions that are required to take place for emergency response, continuity of operations, and dealing with major outages must be documented and readily available to the operations staff. There should be at least two instances of these documents: the original that is kept on-site and a copy that is at an offsite location.

BC plans should not be trusted until they have been tested. Organizations should carry out exercises to ensure that the staff fully understands their responsibilities and how to carry them out. We already covered the various types of exercises that can be used to test plans and staff earlier in this chapter when we discussed DR. Another issue to consider is how to keep these plans up to date. As our dynamic, networked environments change, so must our plans on how to rescue them when necessary.

Although in the security industry “contingency planning” and “business continuity planning (BCP)” are commonly used interchangeably, it is important that you understand the actual difference for the CISSP exam. BCP addresses how to keep the organization in business after a major disruption takes place. It is about the survivability of the organization and making sure that critical functions can still take place even after a disaster. Contingency plans address how to deal with small incidents that do not qualify as disasters, as in power outages, server failures, a down communication link to the Internet, or the corruption of software. Organizations must be ready to deal with both large and small issues that they may encounter.



EXAM TIP BCP is broad in scope and deals with survival of the organization. Contingency plans are narrow in scope and deal with specific issues.

As a security professional you will most likely not be in charge of BCP, but you should most certainly be an active participant in developing the BCP. You will also be involved in BC exercises and may even be a lead in those that focus on information systems. To effectively participate in BC planning and exercises, you should be familiar with the BCP life cycle, how to ensure continuous availability of critical information systems, and the particular requirements of the end-user environments. We look at these in the following sections.

BCP Life Cycle

Remember that most organizations aren't static, but change, often rapidly, as do the conditions under which they must operate. Thus, BCP should be considered a life cycle in order to deal with the constant and inevitable change that will affect it. Understanding and

maintaining each step of the BCP life cycle is critical to ensuring that the BC plan remains useful to the organization. The BCP life cycle is outlined in Figure 23-7.

Note that this life cycle has two modes: normal management (shown in the top half of Figure 23-7) and incident management (shown in the bottom half). In the normal mode, the focus of the BC team is on ensuring preparedness. Obviously, we want to start

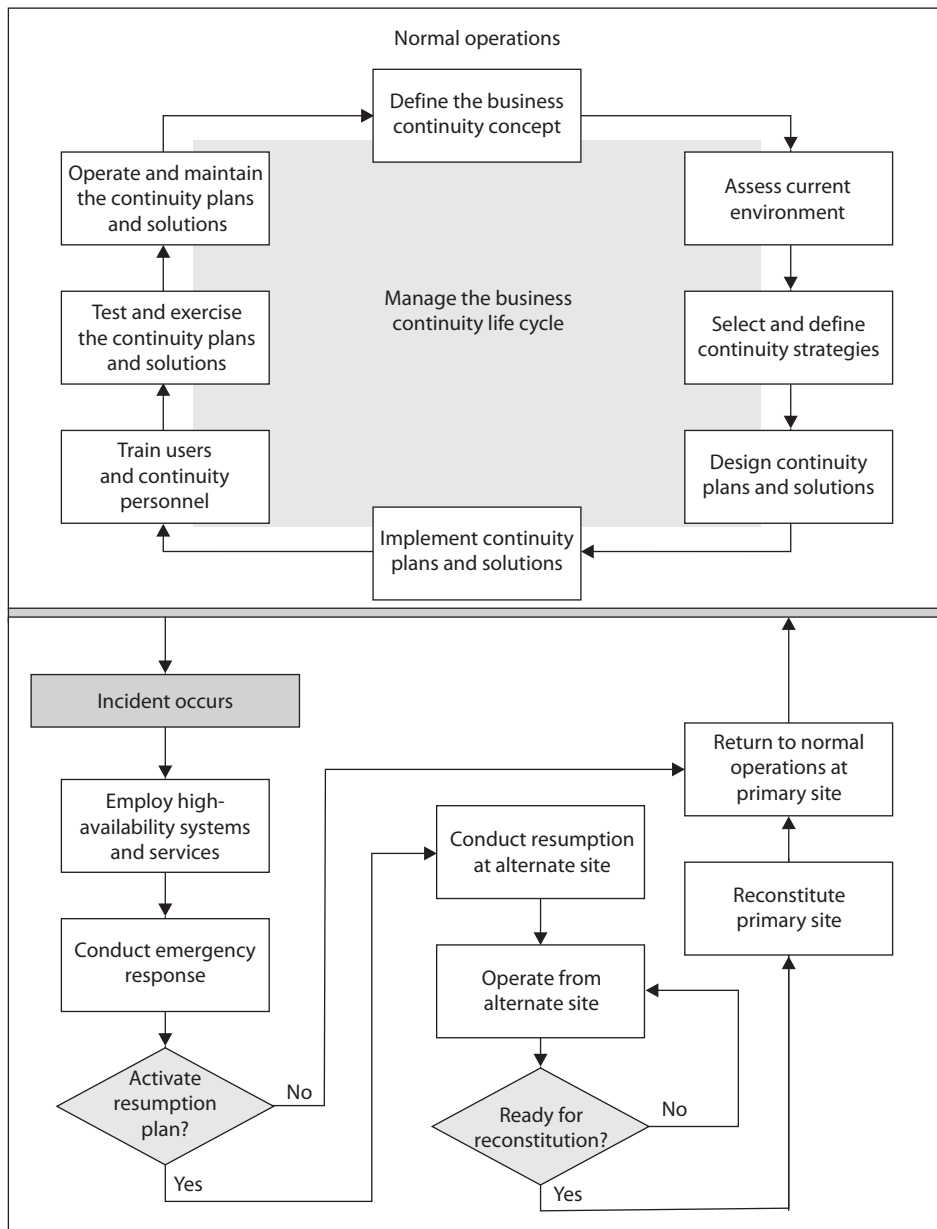


Figure 23-7 BCP life cycle

with a clearly defined concept for what business continuity means for the organization. What are the critical business functions that must continue to operate regardless of what incident happens? What are the minimum levels of performance that are acceptable for these functions?

Once we define the BC concept, we can take a look at the current environment and consider the strategies that would allow continuity of operations under a variety of conditions. It is important to consider that, unlike DR planning, not every type of incident covered in BCP involves loss of IT capabilities. Many organizations suffered tremendously in 2020 because their BCP didn't account for a global pandemic in which many (or even all) staff members would have to work from home for extended periods of time. Information systems are certainly an important part of the continuity strategies, plans, and solutions, but the scope of the BCP is much broader than that of the DRP.

The BC plan is only useful if the organization in general, and the BC team in particular, knows how to execute the plan. This requires periodic training, tests, and exercises to ensure that both the plan and the staff are able to keep the business going no matter what comes their way. As we find gaps and opportunities for improvement, we get to redefine our BCP concept and start another run through the cycle. This continuous improvement is key to being able to switch into incident management mode (at the bottom of Figure 23-7) when needed and execute the BC plan (and, potentially, the DR plan) to keep the business going.

Information Systems Availability

Our main job as CISSPs in the BCP life cycle is to ensure the continuous availability of organizational information systems. To this end, we should ensure the BCP includes backup solutions for the following:

- Network and computer equipment
- Voice and data communications resources
- Human resources
- Transportation of equipment and personnel
- Environment issues (HVAC)
- Data and personnel security issues
- Supplies (paper, forms, cabling, and so on)
- Documentation

The BCP team must understand the organization's current technical environment. This means the planners have to know the intimate details of the network, communications technologies, computers, network equipment, and software requirements that are necessary to get the critical functions up and running. What is surprising to some people is that many organizations do not *totally* understand how their network is configured and how it actually works, because the network may have been established 10 to 15 years ago and has kept growing and changing under different administrators and personnel.

Outsourcing

Part of the planned response to a disaster may be to outsource some of the affected activities to another organization. Organizations do outsource activities—help-desk services, manufacturing, legal advice—all the time, so why not important functions affected by a disaster? Some companies specialize in disaster response and continuity planning and can act as expert consultants.

That is all well and good. However, be aware that your organization is still ultimately responsible for the continuity of a product or service that is outsourced. Clients and customers will expect the organization to ensure continuity of its products and services, either by itself or by having chosen the right outside vendors to provide the products and services. If outside vendors are brought in, the active participation of key in-house managers in their work is still essential. They still need to supervise the work of the outside vendors.

This same concern applies to normal, third-party suppliers of goods and services to the organization. Any BCP should take them into account as well. Note that the process for evaluating an outsourced company for BCP is like that for evaluating the organization itself. The organization must make sure that the outsourced company is financially viable and has its own solid BCP.

The organization can take the following steps to better ensure the continuity of its outsourcing:

- Make the ability of such companies to reliably assure continuity of products and services part of any work proposals.
- Make sure that business continuity planning is included in contracts with such companies, and that their responsibilities and levels of service are clearly spelled out.
- Draw up realistic and reasonable service levels that the outsourced firm will meet during an incident.
- If possible, have the outsourcing companies take part in BCP awareness programs, training, and testing.

The goal is to make the supply of goods and services from outsources as resilient as possible in the wake of a disaster.

New devices are added, new computers are added, new software packages are added, VoIP may have been integrated, and the DMZ may have been split up into three DMZs, with an extranet for the organization's partners. Maybe a company bought and merged with another company and network. Over ten years, a number of technology refreshes most likely have taken place, and the individuals who are maintaining the environment now likely are not the same people who built it ten years ago. Many IT departments experience extensive employee turnover every five years. And most organizational network

schematics are notoriously out of date because everyone is busy with their current tasks (or will come up with new tasks just to get out of having to update the schematic).

So the BCP team has to make sure that if the networked environment is partially or totally destroyed, the recovery team has the knowledge and skill to properly rebuild it.



NOTE Many organizations use VoIP, which means that if the network goes down, network and voice capability are unavailable. The BCP team should address the possible need of redundant voice systems.

The BCP team needs to incorporate into the BCP several things that are commonly overlooked, such as hardware replacements, software products, documentation, environmental needs, and human resources.

Hardware Backups

The BCP needs to identify the equipment required to keep the critical functions up and running. This may include servers, user workstations, routers, switches, tape backup devices, and more. The needed inventory may seem simple enough, but as they say, the devil is in the details. If the recovery team is planning to use images to rebuild newly purchased servers and workstations because the original ones were destroyed, for example, will the images work on the new computers? Using images instead of building systems from scratch can be a time-saving task, unless the team finds out that the replacement equipment is a newer version and thus the images cannot be used. The BCP should plan for the recovery team to use the organization's current images, but also have a manual process of how to build each critical system from scratch with the necessary configurations.

The BCP also needs to be based on accurate estimates of how long it will take for new equipment to arrive. For example, if the organization has identified Dell as its equipment replacement supplier, how long will it take this vendor to send 20 servers and 30 workstations to the offsite facility? After a disaster hits, the organization could be in its offsite facility only to find that its equipment will take three weeks to be delivered. So, the SLA for the identified vendors needs to be investigated to make sure the organization is not further damaged by delays. Once the parameters of the SLA are understood, the BCP team must make a decision between depending upon the vendor and purchasing redundant systems and storing them as backups in case the primary equipment is destroyed.

As described earlier, when potential organizational risks are identified, it is better to take preventive steps to reduce the potential damage. After the calculation of the MTD values, the team will know how long the organization can operate without a specific device. This data should be used to make the decision on whether the organization should depend on the vendor's SLA or make readily available a hot-swappable redundant system. If the organization will lose \$50,000 per hour if a particular server goes down, then the team should elect to implement redundant systems and technology.

If an organization is using any legacy computers and hardware and a disaster hits tomorrow, where would it find replacements for this legacy equipment? The BCP

team should identify legacy devices and understand the risk the organization is facing if replacements are unavailable. This finding has caused many organizations to move from legacy systems to commercial off-the-shelf (COTS) products to ensure that timely replacement is possible.

Software Backups

Most organizations' IT departments have their array of software disks and licensing information here or there—or possibly in one centralized location. If the facility were destroyed and the IT department's current environment had to be rebuilt, how would it gain access to these software packages? The BCP team should make sure to have an inventory of the necessary software required for mission-critical functions and have backup copies at an offsite facility. Hardware is usually not worth much to an organization without the software required to run on it. The software that needs to be backed up can be in the form of applications, utilities, databases, and operating systems. The business continuity plan must have provisions to back up and protect these items along with hardware and data.

It is common for organizations to work with software developers to create customized software programs. For example, in the banking world, individual financial institutions need software that enables their bank tellers to interact with accounts, hold account information in databases and mainframes, provide online banking, carry out data replication, and perform a thousand other types of bank-like functionalities. This specialized type of software is developed and available through a handful of software vendors that specialize in this market. When bank A purchases this type of software for all of its branches, the software has to be specially customized for its environment and needs. Once this banking software is installed, the whole organization depends upon it for its minute-by-minute activities.

When bank A receives the specialized and customized banking software from the software vendor, bank A does not receive the source code. Instead, the software vendor provides bank A with a compiled version. Now, what if this software vendor goes out of business because of a disaster or bankruptcy? Then bank A will require a new vendor to maintain and update this banking software; thus, the new vendor will need access to the source code.

The protection mechanism that bank A should implement is called *software escrow*, in which a third party holds the source code, backups of the compiled code, manuals, and other supporting materials. A contract between the software vendor, customer, and third party outlines who can do what, and when, with the source code. This contract usually states that the customer can have access to the source code only if and when the vendor goes out of business, is unable to carry out stated responsibilities, or is in breach of the original contract. If any of these activities takes place, then the customer is protected because it can still gain access to the source code and other materials through the third-party escrow agent.

Many organizations have been crippled by not implementing software escrow. They paid a software vendor to develop specialized software, and when the software vendor went belly up, the organizations did not have access to the code that their systems ran on.

End-User Environment

Because the end users are usually the worker bees of an organization, they must be provided a functioning environment as soon as possible after a disaster hits. This means that the BCP team must understand the current operational and technical functioning environment and examine critical pieces so they can replicate them.

In most situations, after a disaster, only a skeleton crew is put back to work. The BCP committee has previously identified the most critical functions of the organization during the analysis stage, and the employees who carry out those functions must be put back to work first. So the recovery process for the user environment should be laid out in different stages. The first stage is to get the most critical departments back online, the next stage is to get the second most important back online, and so on.

The BCP team needs to identify user requirements, such as whether users can work on stand-alone PCs or need to be connected in a network to fulfill specific tasks. For example, in a financial institution, users who work on stand-alone PCs might be able to accomplish some small tasks like filling out account forms, word processing, and accounting tasks, but they might need to be connected to a host system to update customer profiles and to interact with the database.

The BCP team also needs to identify how current automated tasks can be carried out manually if that becomes necessary. If the network is going to be down for 12 hours, could the necessary tasks be accomplished through traditional pen-and-paper methods? If the Internet connection is going to be down for five hours, could the necessary communications take place through phone calls? Instead of transmitting data through the internal mail system, could couriers be used to run information back and forth? Today, we are extremely dependent upon technology, but we often take for granted that it will always be there for us to use. It is up to the BCP team to realize that technology may be unavailable for a period of time and to come up with solutions for those situations.



EXAM TIP As a CISSP, your role in business continuity planning is most likely to be that of an active participant, not to lead it. BCP questions in the exam will be written with this in mind.

Chapter Review

There are four key take-aways in this chapter. The first is that you need to be able to identify and implement strategies that will enable your organization to recover from any disaster, supporting your organization's continuity of operations. Leveraging these strategies, you develop a detailed plan that includes the specific processes that the organization (and particularly the IT and security teams) will execute to recover from specific types of disasters. Thirdly, you have to know how to train your DR team to execute the plan flawlessly, even in the chaos of an actual disaster. This includes ensuring that everyone in the organization is aware of their role in the recovery efforts. Finally, the DRP is the cornerstone of the BCP, so you will be called upon to participate in broader business continuity planning and exercises, even if you are not in charge of that effort.

Quick Review

- Disaster recovery (DR) is the set of practices that enables an organization to minimize loss of, and restore, mission-critical technology infrastructure after a catastrophic incident.
- Business continuity (BC) is the set of practices that enables an organization to continue performing its critical functions through and after any disruptive event.
- The recovery time objective (RTO) is the maximum time period within which a mission-critical system must be restored to a designated service level after a disaster to avoid unacceptable consequences associated with a break in business continuity.
- The work recovery time (WRT) is the maximum amount of time available for certifying the functionality and integrity of restored systems and data so they can be put back into production.
- The recovery point objective (RPO) is the acceptable amount of data loss measured in time.
- The four commonly used data backup strategies are direct-attached storage, network-attached storage, cloud storage, and offline media.
- Electronic vaulting makes copies of files as they are modified and periodically transmits them to an offsite backup site.
- Remote journaling moves transaction logs to an offsite facility for database recovery, where only the reapplication of a series of changes to individual records is required to resynchronize the database.
- Offsite backup locations can supply hot, warm, or cold sites.
- A hot site is fully configured with hardware, software, and environmental needs. It can usually be up and running in a matter of hours. It is the most expensive option, but some organizations cannot be out of business longer than a day without very detrimental results.
- A warm site may have some computers, but it does have some peripheral devices, such as disk drives, controllers, and tape drives. This option is less expensive than a hot site, but takes more effort and time to become operational.
- A cold site is just a building with power, raised floors, and utilities. No devices are available. This is the cheapest of the three options, but can take weeks to get up and operational.
- In a reciprocal agreement, one organization agrees to allow another organization to use its facilities in case of a disaster, and vice versa. Reciprocal agreements are very tricky to implement and may be unenforceable. However, they offer a relatively cheap offsite option and are sometimes the only choice.
- A redundant (or mirrored) site is equipped and configured exactly like the primary site and is completely synchronized, ready to become the primary site at a moment's notice.

- High availability (HA) is a combination of technologies and processes that work together to ensure that some specific thing is up and running most of the time.
- Quality of service (QoS) defines minimum acceptable performance characteristics of a particular service, such as response time, CPU utilization, or network bandwidth utilization.
- Fault tolerance is the capability of a technology to continue to operate as expected even if something unexpected takes place (a fault).
- Resilience means that the system continues to function, albeit in a degraded fashion, when a fault is encountered.
- When returning to the original site after a disaster, the least critical organizational units should go back first.
- Disaster recovery plans can be tested through checklist tests, structured walkthroughs, tabletop exercises, simulation tests, parallel tests, or full-interruption tests.
- Business continuity planning addresses how to keep the organization in business after a major disruption takes place, but it is important to note that the scope is much broader than that of disaster recovery.
- The BCP life cycle includes developing the BC concept; assessing the current environment; implementing continuity strategies, plans, and solutions; training the staff; and testing, exercising, and maintaining the plans and solutions.
- An important part of the business continuity plan is to communicate its requirements and procedures to all employees.

Questions

Please remember that these questions are formatted and asked in a certain way for a reason. Keep in mind that the CISSP exam is asking questions at a conceptual level. Questions may not always have the perfect answer, and the candidate is advised against always looking for the perfect answer. Instead, the candidate should look for the best answer in the list.

1. Which best describes a hot-site facility versus a warm- or cold-site facility?
 - A. A site that has disk drives, controllers, and tape drives
 - B. A site that has all necessary PCs, servers, and telecommunications
 - C. A site that has wiring, central air-conditioning, and raised flooring
 - D. A mobile site that can be brought to the organization's parking lot
2. Which of the following describes a cold site?
 - A. Fully equipped and operational in a few hours
 - B. Partially equipped with data processing equipment
 - C. Expensive and fully configured
 - D. Provides environmental measures but no equipment

3. Which is the best description of remote journaling?
 - A. Backing up bulk data to an offsite facility
 - B. Backing up transaction logs to an offsite facility
 - C. Capturing and saving transactions to two mirrored servers in-house
 - D. Capturing and saving transactions to different media types
4. Which of the following does not describe a reciprocal agreement?
 - A. The agreement is enforceable.
 - B. It is a cheap solution.
 - C. It may be able to be implemented right after a disaster.
 - D. It could overwhelm a current data processing site.
5. If a system is fault tolerant, what would you expect it to do?
 - A. Continue to operate as expected even if something unexpected takes place
 - B. Continue to function in a degraded fashion
 - C. Tolerate outages caused by known faults
 - D. Raise an alarm, but tolerate an outage caused by any fault
6. Which of the following approaches to testing your disaster recovery plan would be least desirable if you had to maintain high availability of over 99.999 percent?
 - A. Checklist test
 - B. Parallel test
 - C. Full-interruption test
 - D. Structured walkthrough test

Use the following scenario to answer Questions 7–10. You are the CISO of a small research and development (R&D) company and realize that you don't have a disaster recovery plan (DRP). The projects your organization handles are extremely sensitive and, despite having a very limited budget, you have to bring the risk of project data being lost as close to zero as you can. Recovery time is not as critical because you bill your work based on monthly deliverables and have some leeway at your disposal. Because of the sensitivity of your work, remote working is frowned upon and you keep your research data on local servers (including Exchange for e-mail, Mattermost for group chat, and Apache for web) at your headquarters (and only) site.

7. Which recovery site strategy would be best for you to consider?
 - A. Reciprocal agreement
 - B. Hot site
 - C. Warm site
 - D. Cold site

8. Which of the following recovery site characteristics would be best for your organization?
 - A. As close to headquarters as possible within budgetary constraints
 - B. 100 miles away from headquarters, on a different power grid
 - C. 15 miles away from headquarters on a different power grid
 - D. As far away from headquarters as possible
9. Which data backup storage strategy would you want to implement?
 - A. Direct-attached storage
 - B. Network-attached storage
 - C. Offline media
 - D. Cloud storage
10. Which of the following would be the best way to communicate with all members of the organization in the event of a disaster that takes out your site?
 - A. Internal Mattermost channel
 - B. External Slack channel
 - C. Exchange e-mail
 - D. Call trees

Answers

1. **B.** A hot site is a facility that is fully equipped and properly configured so that it can be up and running within hours to get an organization back into production. Answer B gives the best definition of a fully functional environment.
2. **D.** A cold site only provides environmental measures—wiring, HVAC, raised floors—basically a shell of a building and no more.
3. **B.** Remote journaling is a technology used to transmit data to an offsite facility, but this usually only includes moving the journal or transaction logs to the offsite facility, not the actual files.
4. **A.** A reciprocal agreement is not enforceable, meaning that the organization that agreed to let the damaged organization work out of its facility can decide not to allow this to take place. A reciprocal agreement is a better secondary backup option if the original plan falls through.
5. **A.** Fault tolerance is the capability of a technology to continue to operate as expected even if something unexpected takes place (a fault), with no degradations or outages.
6. **C.** A full-interruption test is the most intrusive to regular operations and business productivity. The original site is actually shut down, and processing takes place at the alternate site. This is almost guaranteed to exceed your allowed downtime unless everything went extremely well.

7. **D.** Because you are working on a tight budget and have the luxury of recovery time, you want to consider the least expensive option. A reciprocal agreement would be ideal except for the sensitivity of your data, which could not be shared with a similar organization (that could, presumably, be a competitor at some point). The next option (cost-wise) is a cold site, which would work in the given scenario.
8. **C.** An ideal recovery site would be on a different power grid to minimize the risk that power will be out on both sites, but close enough for employees to commute. This second point is important because, due to the sensitivity of your work, your organization has a low tolerance for remote work.
9. **C.** Since your data is critical enough that you have to bring the risk of it being lost as close to zero as you can, you would want to use offline media such as tape backups, optical discs, or even external drives that are disconnected after each backup (and potentially removed offsite). This is the slowest and most expensive approach, but is also the most resistant to attacks.
10. **B.** If your site is taken out, you would lose both Exchange and Mattermost since those servers are hosted locally. Call trees only work well for initial notification, leaving an externally hosted Slack channel as the best option. This would require your staff to be aware of this means of communication and have accounts created before the disaster.

PART VIII

Software Development Security

- **Chapter 24** Software Development
- **Chapter 25** Secure Software

This page intentionally left blank

Software Development

This chapter presents the following:

- Software development life cycle
- Development methodologies
- Operation and maintenance
- Maturity models

Always code as if the guy who ends up maintaining your code will be a violent psychopath who knows where you live.

—John F. Woods

Software is usually developed with a strong focus on functionality, not security. In many cases, security controls are bolted on as an afterthought (if at all). To get the best of both worlds, security and functionality have to be designed and integrated at each phase of the software development life cycle. Security should be interwoven into the core of a software product and provide protection at the necessary layers. This is a better approach than trying to develop a front end or wrapper that may reduce the overall functionality and leave security holes when the software has to be integrated into a production environment.

Before we get too deep into secure software development, however, we have to develop a shared understanding of how code is developed in the first place. In this chapter we will cover the complex world of software development so that we can understand the bad things that can happen when security is not interwoven into products properly (discussed in Chapter 25).

Software Development Life Cycle

The life cycle of software development deals with putting repeatable and predictable processes in place that help ensure functionality, cost, quality, and delivery schedule requirements are met. So instead of winging it and just starting to develop code for a project, how can we make sure we build the best software product possible?

Several *software development life cycle (SDLC)* models have been developed over the years, which we will cover later in this section, but the crux of each model deals with the following phases:

- **Requirements gathering** Determining *why* to create this software, *what* the software will do, and *for whom* the software will be created
- **Design** Encapsulating into a functional design *how* the software will accomplish the requirements
- **Development** Programming software code to meet specifications laid out in the design phase and integrating that code with existing systems and/or libraries
- **Testing** Verifying and validating software to ensure that the software works as planned and that goals are met
- **Operations and maintenance** Deploying the software and then ensuring that it is properly configured, patched, and monitored



EXAM TIP You don't need to memorize the phases of the SDLC. We discuss them here so you understand all the tasks that go into developing software and how to integrate security throughout the whole cycle.

In the following sections we will cover the different phases that make up an SDLC model and some specific items about each phase that are important to understand.

Software Development Roles

The specific roles within a software development team will vary based on the methodology being used, the maturity of the organization, and the size of the project (to name just a few parameters). Typically, however, a team has at least the following roles:

- **Project manager (PM)** This role has overall responsibility for the software development project, particularly with regard to cost, schedule, performance, and risk.
- **Team leads** It is rare for software projects to be tackled by a single team, so we usually divide them up and assign a good developer to lead each part.
- **Architect** Sometimes called a tech lead, this role figures out what technologies to use internally or when interfacing with external systems.
- **Software engineer** The people who actually write the programming code are oftentimes specialists in either frontends (e.g., user interfaces) or various types of backends (e.g., business logic, databases). Engineers that can do all of this are called full-stack developers.
- **Quality assurance (QA)** Whether this is a single person or an entire team, this role implements and runs testing processes that detect software defects as early as possible.

Keep in mind that the discussion that follows covers phases that may happen repeatedly and in limited scope depending on the development methodology being used. Before we get into the phases of the SDLC, let's take a brief look at the glue that holds them together: project management.

Project Management

Many developers know that good project management keeps the project moving in the right direction, allocates the necessary resources, provides the necessary leadership, and hopes for the best but plans for the worst. Project management processes should be put into place to make sure the software development project executes each life-cycle phase properly. Project management is an important part of product development, and security management is an important part of project management.

The project manager draws up a security plan at the beginning of a development project and integrates it into the functional plan to ensure that security is not overlooked. This plan will probably be broad and should refer to documented references for more detailed information. The references could include computer standards (RFCs, IEEE standards, and best practices), documents developed in previous projects, security policies, accreditation statements, incident-handling plans, and national or international guidelines. This helps ensure that the plan stays on target.

The security plan should have a life cycle of its own. It will need to be added to, subtracted from, and explained in more detail as the project continues. Keeping the security plan up to date for future reference is important, because losing track of actions, activities, and decisions is very easy once a large and complex project gets underway.

The security plan and project management activities could be scrutinized later, particularly if a vulnerability causes losses to a third party, so we should document security-related decisions. Being able to demonstrate that security was fully considered in each phase of the SDLC can prove that the team exercised due care and this, in turn, can mitigate future liabilities. To this end, the documentation must accurately reflect how the product was built and how it is supposed to operate once implemented into an environment.

If a software product is being developed for a specific customer, it is common for a *Statement of Work (SOW)* to be developed, which describes the product and customer requirements. A detailed SOW helps to ensure that all stakeholders understand these requirements and don't make any undocumented assumptions.

Sticking to what is outlined in the SOW is important so that *scope creep* does not take place. If the scope of a project continually extends (creeps) in an uncontrollable manner, the project may never end, not meet its goals, run out of funding, or all of the foregoing. If the customer wants to modify its requirements, it is important that the SOW is updated and funding is properly reviewed.

A *work breakdown structure (WBS)* is a project management tool used to define and group a project's individual work elements in an organized manner. It is a deliberate decomposition of the project into tasks and subtasks that result in clearly defined deliverables. The SDLC should be illustrated in a WBS format, so that each phase is properly addressed.

Requirements Gathering Phase

This is the phase in which everyone involved in the software development project attempts to understand why the project is needed and what the scope of the project entails. Typically, either a specific customer needs a new application or a demand for the product exists in the market. During this phase, the software development team examines the software's requirements and proposed functionality, engages in brainstorming sessions, and reviews obvious restrictions.

A conceptual definition of the project should be initiated and developed to ensure everyone is on the right page and that this is a proper product to develop. This phase could include evaluating products currently on the market and identifying any demands not being met by current vendors. This definition could also be a direct request for a specific product from a current or future customer.

Typically, the following tasks should be accomplished in this phase:

- Requirements gathering (including security ones)
- Security risk assessment
- Privacy risk assessment
- Risk-level acceptance

The security requirements of the product should be defined in the categories of availability, integrity, and confidentiality. What type of security is required for the software product and to what degree? Some of these requirements may come from applicable external regulations. For example, if the application will deal with payment cards, PCI DSS will dictate some requirements, such as encryption for card information.

An initial security risk assessment should be carried out to identify the potential threats and their associated consequences. This process usually involves asking many, many questions to elicit and document the laundry list of vulnerabilities and threats, the probability of these vulnerabilities being exploited, and the outcome if one of these threats actually becomes real and a compromise takes place. The questions vary from product to product—such as its intended purpose, the expected environment it will be implemented in, the personnel involved, and the types of businesses that would purchase and use the product.

The sensitivity level of the data that many software products store and process has only increased in importance over the years. After a *privacy risk assessment*, a *privacy impact rating* can be assigned, which indicates the sensitivity level of the data that will be processed or accessible. Some software vendors incorporate the following privacy impact ratings in their software development assessment processes:

- **P1, High Privacy Risk** The feature, product, or service stores or transfers personally identifiable information (PII), monitors the user with an ongoing transfer of anonymous data, changes settings or file type associations, or installs software.
- **P2, Moderate Privacy Risk** The sole behavior that affects privacy in the feature, product, or service is a one-time, user-initiated, anonymous data transfer (e.g., the user clicks a link and is directed to a website).

- **P3, Low Privacy Risk** No behaviors exist within the feature, product, or service that affect privacy. No anonymous or personal data is transferred, no PII is stored on the machine, no settings are changed on the user's behalf, and no software is installed.

The software vendor can develop its own privacy impact ratings and their associated definitions. As of this writing there are several formal approaches to conducting a privacy risk assessment, but none stands out as “the” standardized approach to defining a methodology for an assessment or these rating types, but as privacy increases in importance, we might see more standardization in these ratings and associated metrics.

The team tasked with documenting the requirements must understand the criteria for risk-level acceptance to make sure that mitigation efforts satisfy these criteria. Which risks are acceptable will depend on the results of the security and privacy risk assessments. The evaluated threats and vulnerabilities are used to estimate the cost/benefit ratios of the different security countermeasures. The level of each security attribute should be focused upon so that a clear direction on security controls can begin to take shape and can be integrated into the design and development phases.

The end state of the requirements gathering phase is typically a document called the Software (or System) Requirements Specification (SRS), which describes what the software will do and how it will perform. These two high-level objectives are also known as functional and nonfunctional requirements. A *functional requirement* describes a feature of the software system, such as reporting product inventories or processing customer orders. A *nonfunctional requirement* describes performance standards, such as the minimum number of simultaneous user sessions or the maximum response time for a query. Nonfunctional requirements also include security requirements, such as what data must be encrypted and what the acceptable cryptosystems are. The SRS, in a way, is a checklist that the software development team will use to develop the software and the customer will use to accept it.

The Unified Modeling Language (UML) is a common language used to graphically describe all aspects of software development. We will revisit it throughout the different phases, but in terms of software requirements, it allows us to capture both functional and nonfunctional requirements with use case diagrams (UCDs). We already saw these in Chapter 18 when we discussed testing of technical controls. If you look back to Figure 18-3, each use case (shown as verb phrases inside ovals) represents a high-level functional requirement. The associations can capture nonfunctional requirements through special labels, or these requirements can be spelled out in an accompanying use case description.

Design Phase

Once the requirements are formally documented, the software development team can begin figuring out how they will go about satisfying them. This is the phase that starts to map theory to reality. The theory encompasses all the requirements that were identified in the previous phase, and the design outlines how the product is actually going to accomplish these requirements.

Some organizations skip the design phase, but this can cause major delays and redevelopment efforts down the road because a broad vision of the product needs to be understood before looking strictly at the details. Instead, software development teams should develop written plans for how they will build software that satisfies each requirement. This plan usually comprises three different but interrelated models:

- **Informational model** Dictates the type of information to be processed and how it will move around the software system
- **Functional model** Outlines the tasks and functions the application needs to carry out and how they are sequenced and synchronized
- **Behavioral model** Explains the states the application will be in during and after specific transitions take place

For example, consider an antimalware software application. Its informational model would dictate how it processes information, such as virus signatures, modified system files, checksums on critical files, and virus activity. Its functional model would dictate how it scans a hard drive, checks e-mail for known virus signatures, monitors critical system files, and updates itself. Its behavioral model would indicate that when the system starts up, the antimalware software application will scan the hard drive and memory segments. The computer coming online would be the event that changes the state of the application. If it finds a virus, the application would change state and deal with the virus appropriately. Each state must be accounted for to ensure that the product does not go into an insecure state and act in an unpredictable way.

The data from the informational, functional, and behavioral models is incorporated into the software design document, which includes the data, architectural, and procedural design, as shown in Figure 24-1.

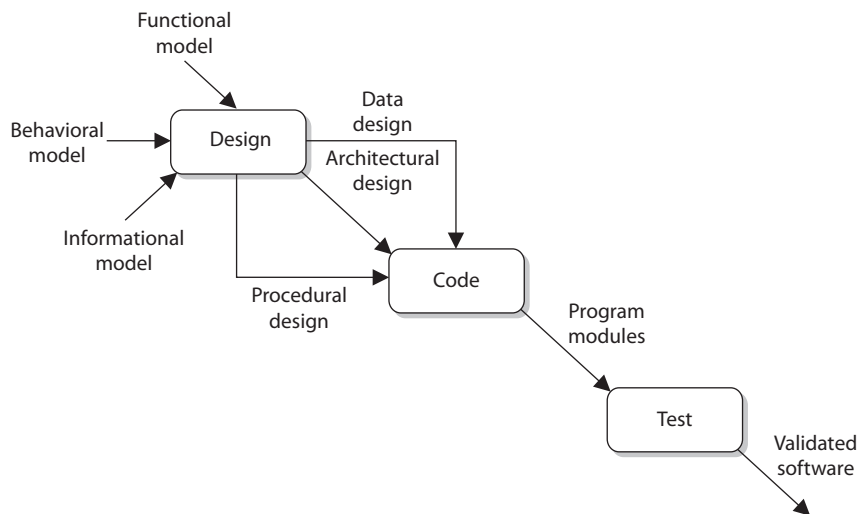


Figure 24-1 Information from three models can go into the design.

From a security point of view, the following items should also be accomplished in the design phase:

- Attack surface analysis
- Threat modeling

An *attack surface* is what is available to be used by an attacker against the product itself. As an analogy, if you were wearing a suit of armor and it covered only half of your body, the other half would be your vulnerable attack surface. Before you went into battle, you would want to reduce this attack surface by covering your body with as much protective armor as possible. The same can be said about software. The software development team should reduce the attack surface as much as possible because the greater the attack surface of software, the more avenues for the attacker; and hence, the greater the likelihood of a successful compromise.

The aim of an *attack surface analysis* is to identify and reduce the amount of code and functionality accessible to untrusted users. The basic strategies of attack surface reduction are to reduce the amount of code running, reduce entry points available to untrusted users, reduce privilege levels as much as possible, and eliminate unnecessary services. Attack surface analysis is generally carried out through specialized tools to enumerate different parts of a product and aggregate their findings into a numeral value. Attack surface analyzers scrutinize files, Registry keys, memory data, session information, processes, and services details. A sample attack surface report is shown in Figure 24-2.

The screenshot displays a web-based report titled "Microsoft Attack Surface Report". It features three tabs: "Report Summary", "Security Issues", and "Attack Surface". The "Security Issues" tab is active, showing a "Table of Contents" with links to "Directories With Weak ACLs", "Processes With NX Disabled", and "Services Vulnerable To Tampering". The "Directories With Weak ACLs" section is expanded, showing a severity of 1 and a weak ACL on the directory C:\Windows\assembly\NativeImages_v2.0.50727_32\AddinExpress.MSO.20# which allows tampering by NT SERVICE\TrustedInstaller. The report includes a description, details (path and weak ACLs), and a table of rights for the account NT SERVICE\TrustedInstaller.

Account	Rights
NT SERVICE\TrustedInstaller (S-1-5-80-956008885-3418522649-1831038044-1853292631-2271478464)	WRITE_OWNER WRITE_DAC FILE_ADD_FILE FILE_ADD_SUBDIRECTORY FILE_DELETE_CHILD FILE_WRITE_ATTRIBUTES FILE_WRITE_EA GENERIC_ALL

Figure 24-2 Attack surface analysis result

Threat modeling, which we covered in detail in Chapter 9 in the context of risk management, is a systematic approach used to understand how different threats could be realized and how a successful compromise could take place. As a hypothetical example, if you were responsible for ensuring that the government building in which you work is safe from terrorist attacks, you would run through scenarios that terrorists would most likely carry out so that you fully understand how to protect the facility and the people within it. You could think through how someone could bring a bomb into the building, and then you would better understand the screening activities that need to take place at each entry point. A scenario of someone running a car into the building would bring up the idea of implementing bollards around the sensitive portions of the facility. The scenario of terrorists entering sensitive locations in the facility (data center, CEO office) would help illustrate the layers of physical access controls that should be implemented.

These same scenario-based exercises should take place during the design phase of software development. Just as you would think about how potential terrorists could enter and exit a facility, the software development team should think through how potentially malicious activities can happen at different input and output points of the software and the types of compromises that can take place within the guts of the software itself.

It is common for software development teams to develop threat trees, as shown in Figure 24-3. A *threat tree* is a tool that allows the development team to understand all the ways specific threats can be realized; thus, it helps them understand what type of security controls they should implement in the software to mitigate the risks associated with each threat type.

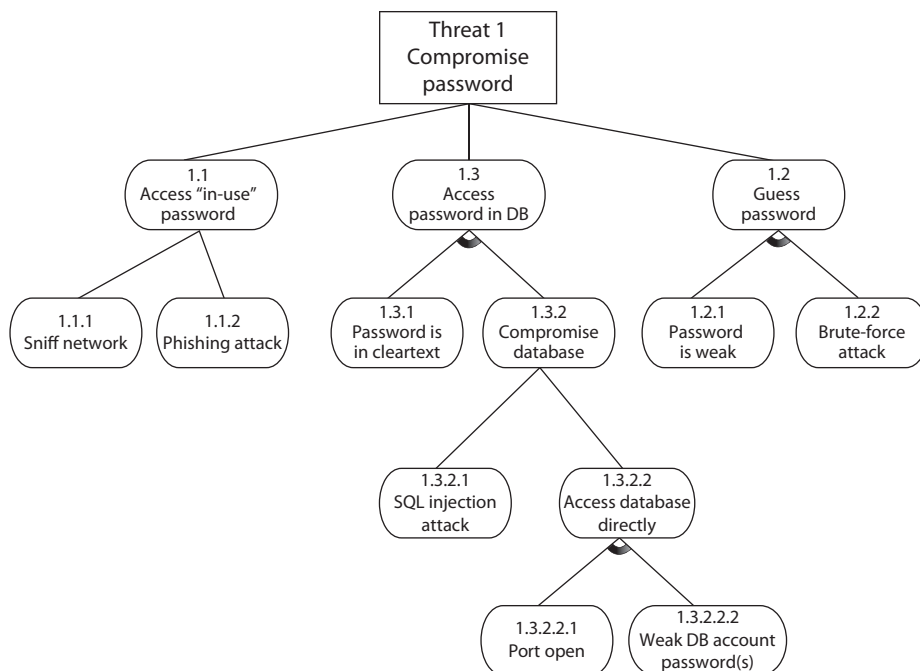
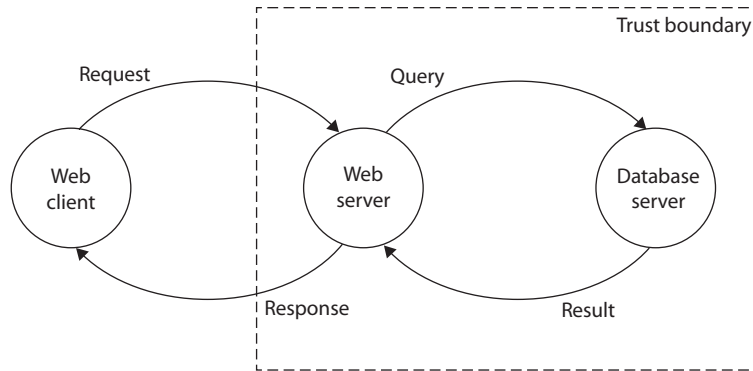


Figure 24-3 Threat tree used in threat modeling

Figure 24-4
A simple flow
diagram for
threat modeling



There are many automated tools in the industry that software development teams can use to ensure that they address the various threat types during the design stage. One popular open-source solution is the Open Web Application Security Project (OWASP) Threat Dragon. This web-based tool enables the development team to describe threats visually using flow diagrams. Figure 24-4 shows a simple diagram of a three-tier web system showing its trust boundary and the four ways in which the tiers interact. The next step in building the threat model would be to consider how each of these four interactions could be exploited by a threat actor. For example, stolen credentials could allow an adversary to compromise the web server and, from there, issue queries to the database server that could compromise the integrity or availability of records stored there. For each threat identified through this process, the software development team would develop controls to mitigate it.

The decisions made during the design phase are pivotal steps to the development phase. Software design serves as a foundation and greatly affects software quality. If good product design is not put into place in the beginning of the project, the following phases will be much more challenging.

Development Phase

This is the phase where the programmers become deeply involved. The software design that was created in the previous phase is broken down into defined deliverables, and programmers develop code to meet the deliverable requirements.

There are many *computer-aided software engineering (CASE)* tools that programmers can use to generate code, test software, and carry out debugging activities. When these types of activities are carried out through automated tools, development usually takes place more quickly with fewer errors.

CASE refers to any type of software tool that supports automated development of software, which can come in the form of program editors, debuggers, code analyzers, version-control mechanisms, and more. These tools aid in keeping detailed records of requirements, design steps, programming activities, and testing. A CASE tool is designed to support one or more software engineering tasks in the process of developing software. Many vendors can get their products to the market faster because they are “computer aided.”

In the next chapter we will delve into the abyss of “secure coding,” but let’s take a quick peek at it here to illustrate its importance in the development phase. As stated previously, most vulnerabilities that corporations, organizations, and individuals have to worry about reside within the programming code itself. When programmers do not follow strict and secure methods of creating programming code, the effects can be widespread and the results can be devastating. But programming securely is not an easy task. The list of errors that can lead to serious vulnerabilities in software is long.

The MITRE organization’s Common Weakness Enumeration (CWE) initiative (<https://cwe.mitre.org/top25>) describes “a demonstrative list of the most common and impactful issues experienced over the previous two calendar years.” Table 24-1 shows the most recent list.

Rank	Name
1	Out-of-bounds Write
2	Improper Neutralization of Input During Web Page Generation (“Cross-site Scripting”)
3	Out-of-bounds Read
4	Improper Input Validation
5	Improper Neutralization of Special Elements used in an OS Command (“OS Command Injection”)
6	Improper Neutralization of Special Elements used in an SQL Command (“SQL Injection”)
7	Use After Free
8	Improper Limitation of a Pathname to a Restricted Directory (“Path Traversal”)
9	Cross-Site Request Forgery (CSRF)
10	Unrestricted Upload of File with Dangerous Type
11	Missing Authentication for Critical Function
12	Integer Overflow or Wraparound
13	Deserialization of Untrusted Data
14	Improper Authentication
15	NULL Pointer Dereference
16	Use of Hard-coded Credentials
17	Improper Restriction of Operations within the Bounds of a Memory Buffer
18	Missing Authorization
19	Incorrect Default Permissions
20	Exposure of Sensitive Information to an Unauthorized Actor
21	Insufficiently Protected Credentials
22	Incorrect Permission Assignment for Critical Resource
23	Improper Restriction of XML External Entity Reference
24	Server-Side Request Forgery (SSRF)
25	Improper Neutralization of Special Elements used in a Command (“Command Injection”)

Table 24-1 2021 CWE Top 25 Most Dangerous Software Weaknesses List

Many of these software issues are directly related to improper or faulty programming practices. Among other issues to address, the programmers need to check input lengths so buffer overflows cannot take place, inspect code to prevent the presence of covert channels, check for proper data types, make sure checkpoints cannot be bypassed by users, verify syntax, and verify checksums. The software development team should play out different attack scenarios to see how the code could be attacked or modified in an unauthorized fashion. Code reviews and debugging should be carried out by peer developers, and everything should be clearly documented.

A particularly important area of scrutiny is input validation because it can lead to serious vulnerabilities. Essentially, we should treat every single user input as malicious until proven otherwise. For example, if we don't put limits on how many characters users can enter when providing, say, their names on a web form, they could cause a buffer overflow, which is a classic example of a technique used to exploit improper input validation. A *buffer overflow* (which is described in detail in Chapter 18) takes place when too much data is accepted as input to a specific process. The process's memory buffer can be overflowed by shoving arbitrary data into various memory segments and inserting a carefully crafted set of malicious instructions at a specific memory address.

Buffer overflows can also lead to illicit escalation of privileges. *Privilege escalation* is the process of exploiting a process or configuration setting to gain access to resources that would normally not be available to the process or its user. For example, an attacker can compromise a regular user account and escalate its privileges to gain administrator or even system privileges on that computer. This type of attack usually exploits the complex interactions of user processes with device drivers and the underlying operating system. A combination of input validation and configuring the system to run with least privilege can help mitigate the threat of escalation of privileges.

What is important to understand is that secure coding practices need to be integrated into the development phase of the SDLC. Security has to be addressed at each phase of the SDLC, with this phase being one of the most critical.

Testing Phase

Formal and informal testing should begin as soon as possible. *Unit testing* is concerned with ensuring the quality of individual code modules or classes. Mature developers develop the unit tests for their modules before they even start coding, or at least in parallel with the coding. This approach is known as *test-driven development* and tends to result in much higher-quality code with significantly fewer vulnerabilities.

Unit tests are meant to simulate a range of inputs to which the code may be exposed. These inputs range from the mundanely expected, to the accidentally unfortunate, to the intentionally malicious. The idea is to ensure the code always behaves in an expected and secure manner. Once a module and its unit tests are finished, the unit tests are run (usually in an automated framework) on that code. The goal of this type of testing is to isolate each part of the software and show that the individual parts are correct.

Unit testing usually continues throughout the development phase. A totally different group of people should carry out the formal testing. Depending on the methodology and the organization, this could be a QA, testing, audit, or even red team. This is an example

Separation of Duties

Different environmental types (development, testing, and production) should be properly separated, and functionality and operations should not overlap. Developers should not have access to modify code used in production. The code should be tested, submitted to a library, and then sent to the production environment.

of separation of duties. A programmer should not develop, test, and release software. The more eyes that see the code, the greater the chance that flaws will be found before the product is released.

No cookie-cutter recipe exists for security testing because the applications and products can be so diverse in functionality and security objectives. It is important to map security risks to test cases and code. The software development team can take a linear approach by identifying a vulnerability, providing the necessary test scenario, performing the test, and reviewing the code for how it deals with such a vulnerability. At this phase, tests are conducted in an environment that should mirror the production environment to ensure the code does not work only in the labs.

Security attacks and penetration tests usually take place during the testing phase to identify any missed vulnerabilities. Functionality, performance, and penetration resistance are evaluated. All the necessary functionality required of the product should be in a checklist to ensure each function is accounted for.

Security tests should be run to test against the vulnerabilities identified earlier in the project. Buffer overflows should be attempted, interfaces should be hit with unexpected inputs, denial-of-service (DoS) situations should be tested, unusual user activity should take place, and if a system crashes, the product should react by reverting to a secure state. The product should be tested in various environments with different applications, configurations, and hardware platforms. A product may respond fine when installed on a clean Windows 10 installation on a stand-alone PC, but it may throw unexpected errors when installed on a laptop that is remotely connected to a network and has a virtual private network (VPN) client installed.

Verification vs. Validation

Verification determines if the software product accurately represents and meets the specifications. After all, a product can be developed that does not match the original specifications, so this step ensures the specifications are being properly met. It answers the question, “Did we build the product right?”

Validation determines if the software product provides the necessary solution for the intended real-world problem. In large projects, it is easy to lose sight of the overall goal. This exercise ensures that the main goal of the project is met. It answers the question, “Did we build the right product?”

Testing Types

Software testers on the software development team should subject the software to various types of tests to discover the variety of potential flaws. The following are some of the most common testing approaches:

- **Unit testing** Testing individual components in a controlled environment where programmers validate data structure, logic, and boundary conditions
- **Integration testing** Verifying that components work together as outlined in the design specifications
- **Acceptance testing** Ensuring that the code meets customer requirements
- **Regression testing** After a change to a system takes place, retesting to ensure functionality, performance, and protection

A well-rounded security test encompasses both manual tests and automated tests. Automated tests help locate a wide range of flaws generally associated with careless or erroneous code implementations. Some automated testing environments run specific inputs in a scripted and repeatable manner. While these tests are the bread and butter of software testing, we sometimes want to simulate random and unpredictable inputs to supplement the scripted tests.

A manual test is used to analyze aspects of the program that require human intuition and can usually be judged using computing techniques. Testers also try to locate design flaws. These include logical errors, which may enable attackers to manipulate program flow by using shrewdly crafted program sequences to access greater privileges or bypass authentication mechanisms. Manual testing involves code auditing by security-centric programmers who try to modify the logical program structure using rogue inputs and reverse-engineering techniques. Manual tests simulate the live scenarios involved in real-world attacks. Some manual testing also involves the use of social engineering to analyze the human weakness that may lead to system compromise.

At this stage, issues found in testing procedures are relayed to the development team in problem reports. The problems are fixed and programs retested. This is a continual process until everyone is satisfied that the product is ready for production. If there is a specific customer, the customer would run through a range of tests before formally accepting the product; if it is a generic product, beta testing can be carried out by various potential customers and agencies. Then the product is formally released to the market or customer.



NOTE Sometimes developers include lines of code in a product that will allow them to do a few keystrokes and get right into the application. This allows them to bypass any security and access controls so they can quickly access the application's core components. This is referred to as a "back door" or "maintenance hook" and must be removed before the code goes into production.

Operations and Maintenance Phase

Once the software code is developed and properly tested, it is released so that it can be implemented within the intended production environment. The software development team's role is not finished at this point. Newly discovered problems and vulnerabilities are commonly

identified at this phase. For example, if a company developed a customized application for a specific customer, the customer could run into unforeseen issues when rolling out the product within its various networked environments. Interoperability issues might come to the surface, or some configurations may break critical functionality. The developers would need to make the necessary changes to the code, retest the code, and re-release the code.

Almost every software system requires the addition of new features over time. Frequently, these have to do with changing business processes or interoperability with other systems. This highlights the need for the operations and development teams to work particularly closely during the operations and maintenance (O&M) phase. The operations team, which is typically the IT department, is responsible for ensuring the reliable operation of all production systems. The development team is responsible for any changes to the software in development systems up until the time the software goes into production. Together, the operations and development teams address the transition from development to production as well as management of the system's configuration.

Another facet of O&M is driven by the fact that new vulnerabilities are regularly discovered. While the developers may have carried out extensive security testing, it is close to impossible to identify all the security issues at one point and time. Zero-day vulnerabilities may be identified, coding errors may be uncovered, or the integration of the software with another piece of software may uncover security issues that have to be addressed. The development team must develop patches, hotfixes, and new releases to address these items. In all likelihood, this is where you as a CISSP will interact the most with the SDLC.

Change Management

One of the key processes on which to focus for improvement involves how we deal with the inevitable changes. These can cause a lot of havoc if not managed properly and in a deliberate manner. We already discussed change management in general in Chapter 20, but it is particularly important during the lifetime of a software development project.

The need to change software arises for several reasons. During the development phase, a customer may alter requirements and ask that certain functionalities be added, removed, or modified. In production, changes may need to happen because of other changes in the environment, new requirements of a software product or system, or newly released patches or upgrades. These changes should be carefully analyzed, approved, and properly incorporated such that they do not affect any original functionality in an adverse way.

Change management is a systematic approach to deliberately regulating the changing nature of projects, including software development projects. It is a management process that takes into account not just the technical issues but also resources (like people and money), project life cycle, and even organizational climate. Many times, the hardest part of managing change is not the change itself, but the effects it has in the organization. Many of us have been on the receiving end of a late-afternoon phone call in which we're told to change our plans because of a change in a project on which we weren't even working. An important part of change management is controlling change.

Change Control

Change control is the process of controlling the specific changes that take place during the life cycle of a system and documenting the necessary change control activities. Whereas change management is the project manager's responsibility as an overarching

process, change control is what developers do to ensure the software doesn't break when they change it.

Change control involves a bunch of things to consider. The change must be approved, documented, and tested. Some tests may need to be rerun to ensure the change does not affect the product's capabilities. When a programmer makes a change to source code, she should do so on the test version of the code. Under no conditions should a programmer change the code that is already in production. After making changes to the code, the programmer should test the code and then deliver the new code to the librarian. Production code should come only from the librarian and not from a programmer or directly from a test environment.

A process for controlling changes needs to be in place at the beginning of a project so that everyone knows how to deal with changes and knows what is expected of each entity when a change request is made. Some projects have been doomed from the start because proper change control was not put into place and enforced. Many times in development, the customer and vendor agree on the design of the product, the requirements, and the specifications. The customer is then required to sign a contract confirming this is the agreement and that if they want any further modifications, they will have to pay the vendor for that extra work. If this agreement is not put into place, then the customer can continually request changes, which requires the software development team to put in the extra hours to provide these changes, the result of which is that the vendor loses money, the product does not meet its completion deadline, and scope creep occurs.

Other reasons exist to have change control in place. These reasons deal with organizational policies, standard procedures, and expected results. If a software product is in the last phase of development and a change request comes in, the development team should know how to deal with it. Usually, the team leader must tell the project manager how much extra time will be required to complete the project if this change is incorporated and what steps need to be taken to ensure this change does not affect other components within the product. If these processes are not controlled, one part of a development team could implement the change without another part of the team being aware of it. This could break some of the other development team's software pieces. When the pieces of the product are integrated and some pieces turn out to be incompatible, some jobs may be in jeopardy, because management never approved the change in the first place.

Change control processes should be evaluated during system audits. It is possible to overlook a problem that a change has caused in testing, so the procedures for how change control is implemented and enforced should be examined during a system audit.

The following are some necessary steps for a change control process:

1. Make a formal request for a change.
2. Analyze the request:
 - a. Develop the implementation strategy.
 - b. Calculate the costs of this implementation.
 - c. Review security implications.
3. Record the change request.
4. Submit the change request for approval.

5. Develop the change:
 - a. Recode segments of the product and add or subtract functionality.
 - b. Link these changes in the code to the formal change control request.
 - c. Submit software for testing and quality control.
 - d. Repeat until quality is adequate.
 - e. Make version changes.
6. Report results to management.

The changes to systems may require another round of certification and accreditation. If the changes to a system are significant, then the functionality and level of protection

SDLC and Security

The main phases of a software development life cycle are shown here with some specific security tasks.

Requirements gathering:

- Security risk assessment
- Privacy risk assessment
- Risk-level acceptance
- Informational, functional, and behavioral requirements

Design:

- Attack surface analysis
- Threat modeling

Development:

- Automated CASE tools
- Secure coding

Testing:

- Automated testing
- Manual testing
- Unit, integration, acceptance, and regression testing

Operations and maintenance:

- Vulnerability patching
- Change management and control

may need to be reevaluated (certified), and management would have to approve the overall system, including the new changes (accreditation).

Development Methodologies

Several software development methodologies are in common use around the world. While some include security issues in certain phases, these are not considered “security-centric development methodologies.” They are simply classical approaches to building and developing software. Let’s dive into some of the methodologies that you should know as a CISSP.



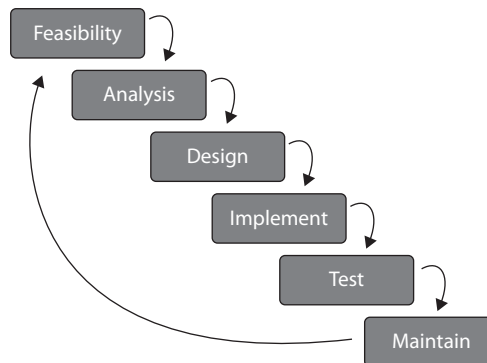
EXAM TIP It is exceptionally rare to see a development methodology used in its pure form in the real world. Instead, organizations typically start with a base methodology and modify it to suit their own unique environment. For purposes of the CISSP exam, however, you should focus on what differentiates each development approach.

Waterfall Methodology

The *Waterfall methodology* uses a linear-sequential life-cycle approach, illustrated in Figure 24-5. Each phase must be completed in its entirety before the next phase can begin. At the end of each phase, a review takes place to make sure the project is on the correct path and should continue.

In this methodology all requirements are gathered in the initial phase and there is no formal way to integrate changes as more information becomes available or requirements change. It is hard to know everything at the beginning of a project, so waiting until the whole project is complete to integrate necessary changes can be ineffective and time consuming. As an analogy, let’s say that you are planning to landscape your backyard that is one acre in size. In this scenario, you can go to the gardening store only one time to get

Figure 24-5
Waterfall
methodology
used for software
development



all your supplies. If you identify during the project that you need more topsoil, rocks, or pipe for the sprinkler system, you have to wait and complete the whole yard before you can return to the store for extra or more suitable supplies.

The Waterfall methodology is a very rigid approach that could be useful for smaller projects in which all the requirements are fully understood up front. It may also be a good choice in some large projects for which different organizations will perform the work at each phase. Overall, however, it is not an ideal methodology for most complex projects, which commonly contain many variables that affect the scope as the project continues.

Prototyping

A *prototype* is a sample of software code or a model that can be developed to explore a specific approach to a problem before investing expensive time and resources. A team can identify the usability and design problems while working with a prototype and adjust their approach as necessary. Within the software development industry, three main prototype models have been invented and used. These are the rapid prototype, evolutionary prototype, and operational prototype.

Rapid prototyping is an approach that allows the development team to quickly create a prototype (sample) to test the validity of the current understanding of the project requirements. In a software development project, the team could develop a *rapid prototype* to see if their ideas are feasible and if they should move forward with their current solution. The rapid prototype approach (also called throwaway) is a “quick and dirty” method of creating a piece of code and seeing if everyone is on the right path or if another solution should be developed. The rapid prototype is not developed to be built upon, but to be discarded after serving its purposes.

When *evolutionary prototypes* are developed, they are built with the goal of incremental improvement. Instead of being discarded after being developed, as in the rapid prototype approach, the evolutionary prototype is continually improved upon until it reaches the final product stage. Feedback that is gained through each development phase is used to improve the prototype and get closer to accomplishing the customer’s needs.

Operational prototypes are an extension of the evolutionary prototype method. Both models (operational and evolutionary) improve the quality of the prototype as more data is gathered, but the operational prototype is designed to be implemented within a production environment as it is being tweaked. The operational prototype is updated as customer feedback is gathered, and the changes to the software happen within the working site.

In summary, a rapid prototype is developed to give a quick understanding of the suggested solution, an evolutionary prototype is created and improved upon within a lab environment, and an operational prototype is developed and improved upon within a production environment.

Incremental Methodology

If a development team follows the *Incremental methodology*, this allows them to carry out multiple development cycles on a piece of software throughout its development stages. This would be similar to “multi-Waterfall” cycles taking place on one piece of software as

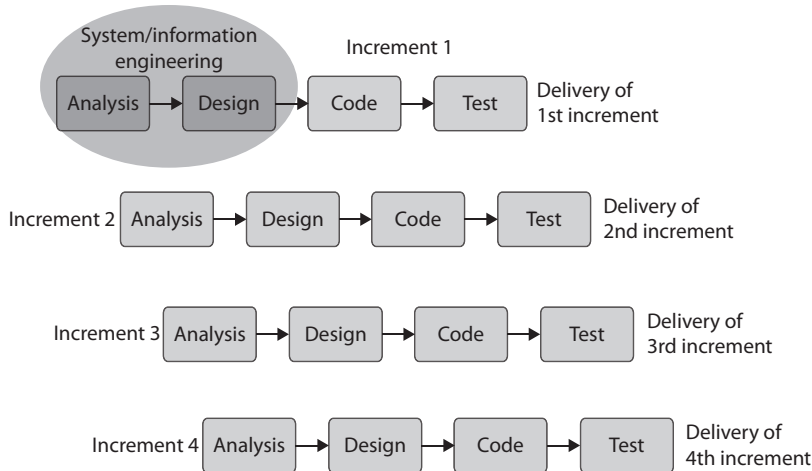


Figure 24-6 Incremental development methodology

it matures through the development stages. A version of the software is created in the first iteration and then it passes through each phase (requirements analysis, design, coding, testing, implementation) of the next iteration process. The software continues through the iteration of phases until a satisfactory product is produced. This methodology is illustrated in Figure 24-6.

When using the Incremental methodology, each incremental phase results in a deliverable that is an operational product. This means that a working version of the software is produced after the first iteration and that version is improved upon in each of the subsequent iterations. Some benefits to this methodology are that a working piece of software is available in early stages of development, the flexibility of the methodology allows for changes to take place, testing uncovers issues more quickly than the Waterfall methodology since testing takes place after each iteration, and each iteration is an easily manageable milestone.

Because each incremental phase delivers an operational product, the customer can respond to each build and help the development team in its improvement processes, and because the initial product is delivered more quickly compared to other methodologies, the initial product delivery costs are lower, the customer gets its functionality earlier, and the risks of critical changes being introduced are lower.

This methodology is best used when issues pertaining to risk, program complexity, funding, and functionality requirements need to be understood early in the product development life cycle. If a vendor needs to get the customer some basic functionality quickly as it works on the development of the product, this can be a good methodology to follow.

Spiral Methodology

The *Spiral methodology* uses an iterative approach to software development and places emphasis on risk analysis. The methodology is made up of four main phases: determine objectives, identify and resolve risks, development and test, and plan the next iteration. The development team starts with the initial requirements and goes through each of these phases, as shown in Figure 24-7. Think about starting a software development project at the center of this graphic. You have your initial understanding and requirements of the project, develop specifications that map to these requirements, identify and resolve risks, build prototype specifications, test your specifications, build a development plan, integrate newly discovered information, use the new information to carry out a new risk analysis, create a prototype, test the prototype, integrate resulting data into the process, and so forth. As you gather more information about the project, you integrate it into the risk analysis process, improve your prototype, test the prototype, and add more granularity to each step until you have a completed product.

The iterative approach provided by the Spiral methodology allows new requirements to be addressed as they are uncovered. Each prototype allows for testing to take place

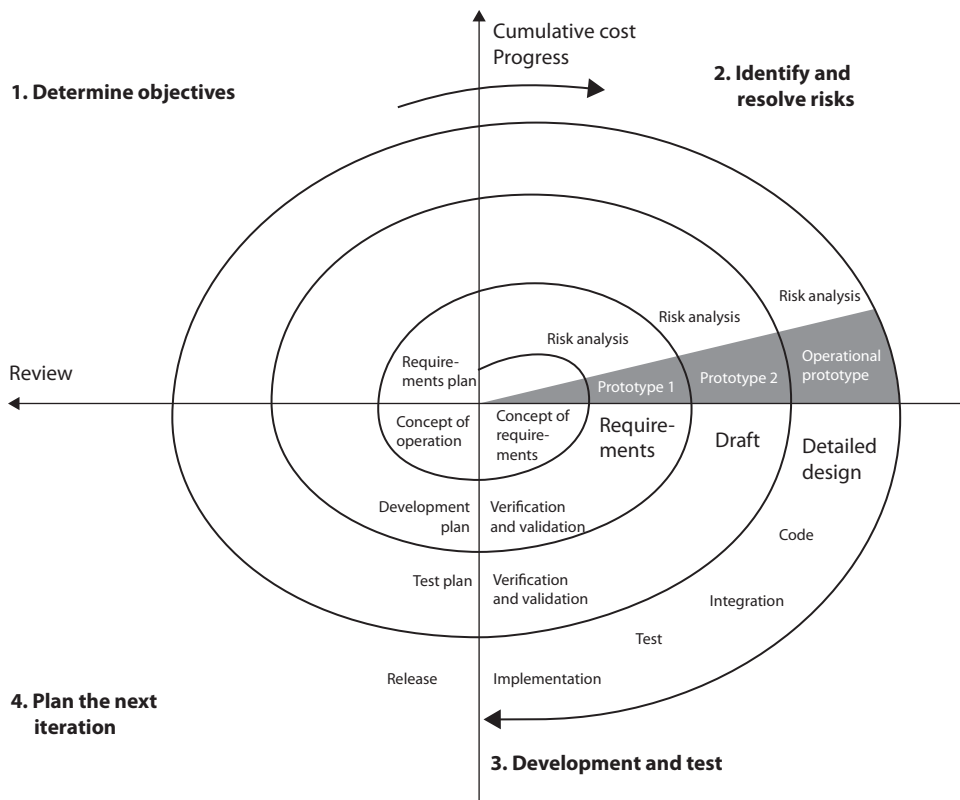


Figure 24-7 Spiral methodology for software development

early in the development project, and feedback based upon these tests is integrated into the following iteration of steps. The risk analysis ensures that all issues are actively reviewed and analyzed so that things do not “slip through the cracks” and the project stays on track.

In the Spiral methodology the last phase allows the customer to evaluate the product in its current state and provide feedback, which is an input value for the next spiral of activity. This is a good methodology for complex projects that have fluid requirements.



NOTE Within this methodology the angular aspect represents progress and the radius of the spirals represents cost.

Rapid Application Development

The *Rapid Application Development (RAD)* methodology relies more on the use of rapid prototyping than on extensive upfront planning. In this methodology, the planning of how to improve the software is interleaved with the processes of developing the software, which allows for software to be developed quickly. The delivery of a workable piece of software can take place in less than half the time compared to the Waterfall methodology. The RAD methodology combines the use of prototyping and iterative development procedures with the goal of accelerating the software development process. The development process begins with creating data models and business process models to help define what the end-result software needs to accomplish. Through the use of prototyping, these data and process models are refined. These models provide input to allow for the improvement of the prototype, and the testing and evaluation of the prototype allow for the improvement of the data and process models. The goal of these steps is to combine business requirements and technical design statements, which provide the direction in the software development project.

Figure 24-8 illustrates the basic differences between traditional software development approaches and RAD. As an analogy, let's say that the development team needs you to tell them what it is you want so that they can build it for you. You tell them that the thing you want has four wheels and an engine. They bring you a two-seat convertible and ask, “Is this what you want?” You say, “No, it must be able to seat four adults.” So they leave the prototype with you and go back to work. They build a four-seat convertible and deliver it to you, and you tell them they are getting closer but it still doesn't fit your requirements. They get more information from you, deliver another prototype, get more feedback, and on and on. That back and forth is what is taking place in the circle portion of Figure 24-8.

The main reason that RAD was developed was that by the time software was completely developed following other methodologies, the requirements changed and the developers had to “go back to the drawing board.” If a customer needs you to develop a software product and it takes you a year to do so, by the end of that year the customer's needs for the software have probably advanced and changed. The RAD methodology allows for the customer to be involved during the development phases so that the end result maps to their needs in a more realistic manner.

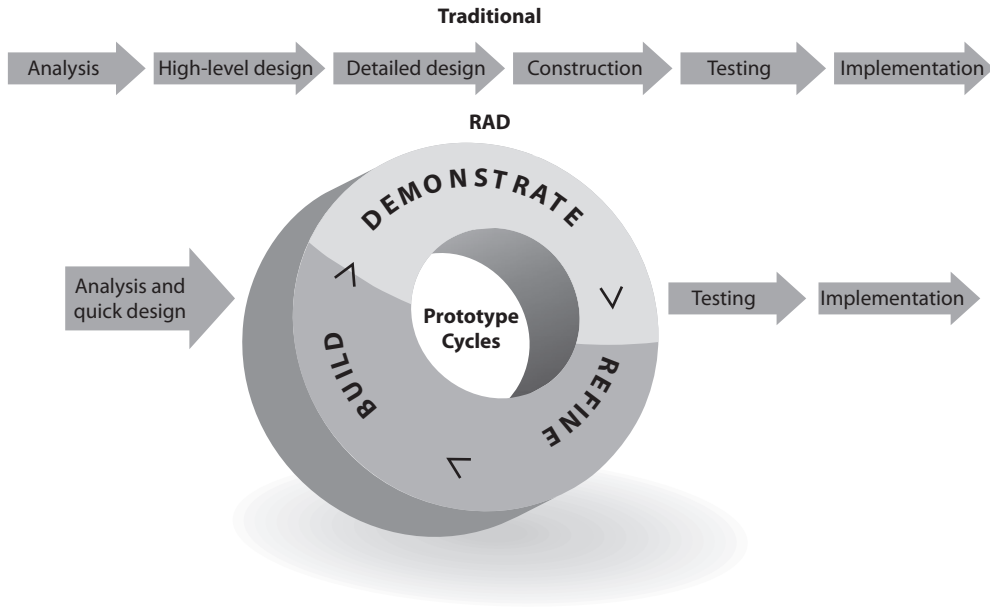


Figure 24-8 Rapid Application Development methodology

Agile Methodologies

The industry seems to be full of software development methodologies, each trying to improve upon the deficiencies of the ones before it. Before the Agile approach to development was created, teams were following rigid process-oriented methodologies. These approaches focused more on following procedures and steps instead of potentially carrying out tasks in a more efficient manner. As an analogy, if you have ever worked within or interacted with a large government agency, you may have come across silly processes that took too long and involved too many steps. If you are a government employee and need to purchase a new chair, you might have to fill out four sets of documents that need to be approved by three other departments. You probably have to identify three different chair vendors, who have to submit a quote, which goes through the contracting office. It might take you a few months to get your new chair. The focus is to follow a protocol and rules instead of efficiency.

Many of the classical software development approaches, as in Waterfall, provide rigid processes to follow that do not allow for much flexibility and adaptability. Commonly, the software development projects that follow these approaches end up failing by not meeting schedule time release, running over budget, and/or not meeting the needs of the customer. Sometimes you need the freedom to modify steps to best meet the situation's needs.

Agile methodology is an umbrella term for several development methodologies. The overarching methodology focuses not on rigid, linear, stepwise processes, but instead on incremental and iterative development methods that promote cross-functional teamwork

and continuous feedback mechanisms. This methodology is considered “lightweight” compared to the traditional methodologies that are “heavyweight,” which just means this methodology is not confined to a tunnel-visioned and overly structured approach. It is nimble and flexible enough to adapt to each project’s needs. The industry found out that even an exhaustive library of defined processes cannot handle every situation that could arise during a development project. So instead of investing time and resources into deep upfront design analysis, the Agile methodology focuses on small increments of functional code that are created based on business need.

The various methodologies under the Agile umbrella focus on individual interaction instead of processes and tools. They emphasize developing the right software product over comprehensive and laborious documentation. They promote customer collaboration instead of contract negotiation, and emphasize abilities to respond to change instead of strictly following a plan.

A notable element of many Agile methodologies is their focus on user stories. A *user story* is a sentence that describes what a user wants to do and why. For instance, a user story could be “As a customer, I want to search for products so that I can buy some.” Notice the structure of the story is “As a <user role>, I want to <accomplish some goal> so that <reason for accomplishing the goal>.” For example, “As a network analyst, I want to record pcap (packet capture) files so that I can analyze downloaded malware.” This method of documenting user requirements is very familiar to the customers and enables their close collaboration with the development team. Furthermore, by keeping this user focus, validation of the features is simpler because the “right system” is described up front by the users in their own words.



EXAM TIP The Agile methodologies do not use prototypes to represent the full product, but break the product down into individual features that are continuously being delivered.

Another important characteristic of the Agile methodologies is that the development team can take pieces and parts of all of the available SDLC methodologies and combine them in a manner that best meets the specific project needs. These various combinations have resulted in many methodologies that fall under the Agile umbrella.

Scrum

Scrum is one of the most widely adopted Agile methodologies in use today. It lends itself to projects of any size and complexity and is very lean and customer focused. Scrum is a methodology that acknowledges the fact that customer needs cannot be completely understood and will change over time. It focuses on team collaboration, customer involvement, and continuous delivery.

The term *scrum* originates from the sport of rugby. Whenever something interrupts play (e.g., a penalty or the ball goes out of bounds) and the game needs to be restarted, all players come together in a tight formation. The ball is then thrown into the middle and the players struggle with each other until one team or the other gains possession of the ball, allowing the game to continue. Extending this analogy, the Scrum methodology

allows the project to be reset by allowing product features to be added, changed, or removed at clearly defined points. Since the customer is intimately involved in the development process, there should be no surprises, cost overruns, or schedule delays. This allows a product to be iteratively developed and changed even as it is being built.

The change points happen at the conclusion of each *sprint*, a fixed-duration development interval that is usually (but not always) two weeks in length and promises delivery of a very specific set of features. These features are chosen by the team, but with a lot of input from the customer. There is a process for adding features at any time by inserting them in the feature backlog. However, these features can be considered for actual work only at the beginning of a new sprint. This shields the development team from changes during a sprint, but allows for changes in between sprints.

Extreme Programming

If you take away the regularity of Scrum's sprints and backlogs and add a lot of code reviewing, you get our next Agile methodology. Extreme Programming (XP) is a development methodology that takes code reviews (discussed in Chapter 18) to the extreme (hence the name) by having them take place continuously. These continuous reviews are accomplished using an approach called *pair programming*, in which one programmer dictates the code to her partner, who then types it. While this may seem inefficient, it allows two pairs of eyes to constantly examine the code as it is being typed. It turns out that this approach significantly reduces the incidence of errors and improves the overall quality of the code.

Another characteristic of XP is its reliance on test-driven development, in which the unit tests are written before the code. The programmer first writes a new unit test case, which of course fails because there is no code to satisfy it. The next step is to add just enough code to get the test to pass. Once this is done, the next test is written, which fails, and so on. The consequence is that only the minimal amount of code needed to pass the tests is developed. This extremely minimal approach reduces the incidence of errors because it weeds out complexity.

Kanban

Kanban is a production scheduling system developed by Toyota to more efficiently support just-in-time delivery. Over time, Kanban was adopted by IT and software systems developers. In this context, the *Kanban* development methodology is one that stresses visual tracking of all tasks so that the team knows what to prioritize at what point in time in order to deliver the right features right on time. Kanban projects used to be very noticeable because entire walls in conference rooms would be covered in sticky notes representing the various tasks that the team was tracking. Nowadays, many Kanban teams opt for virtual walls on online systems.

The Kanban wall is usually divided vertically by production phase. Typical columns are labeled Planned, In Progress, and Done. Each sticky note can represent a user story as it moves through the development process, but more importantly, the sticky note can also be some other work that needs to be accomplished. For instance, suppose that one of the user stories is the search feature described earlier in this section. While it is being developed, the team realizes that the searches are very slow. This could result in a task being

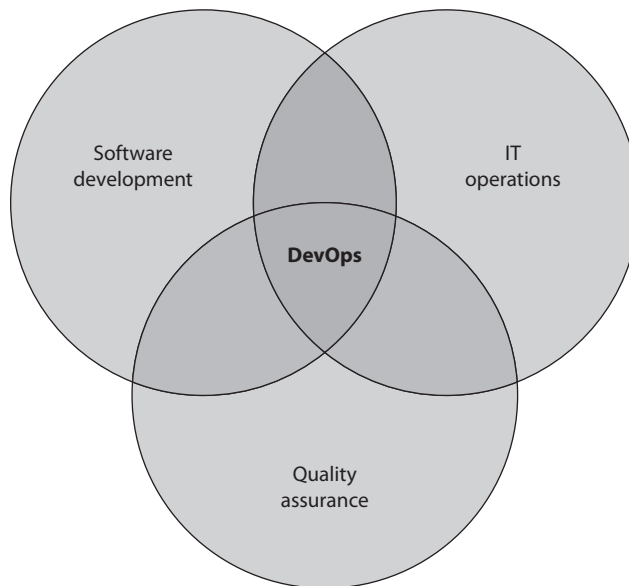
added to change the underlying data or network architecture or to upgrade hardware. This sticky note then gets added to the Planned column and starts being prioritized and tracked together with the rest of the remaining tasks. This process highlights how Kanban allows the project team to react to changing or unknown requirements, which is a common feature among all Agile methodologies.

DevOps

Traditionally, the software development team and the IT team are two separate (and sometimes antagonistic) groups within an organization. Many problems stem from poor collaboration between these two teams during the development process. It is not rare to have the IT team berating the developers because a feature push causes the IT team to have to stay late or work on a weekend or simply drop everything they were doing in order to “fix” something that the developers “broke.” This friction makes a lot of sense when you consider that each team is incentivized by different outcomes. Developers want to push out finished code, usually under strict schedules. The IT staff, on the other hand, wants to keep the IT infrastructure operating effectively. Many project managers who have managed software development efforts will attest to having received complaints from developers that the IT team was being unreasonable and uncooperative, while the IT team was simultaneously complaining about buggy code being tossed over the fence at them at the worst possible times and causing problems on the rest of the network.

A good way to solve this friction is to have both developers and members of the operations staff (hence the term DevOps) on the software development team. *DevOps* is the practice of incorporating development, IT, and quality assurance (QA) staff into software development projects to align their incentives and enable frequent, efficient, and reliable releases of software products. This relationship is illustrated in Figure 24-9.

Figure 24-9
DevOps exists at the intersection of software development, IT, and QA.



Ultimately, DevOps is about changing the culture of an organization. It has a huge positive impact on security, because in addition to QA, the IT teammates will be involved at every step of the process. Multifunctional integration allows the team to identify potential defects, vulnerabilities, and friction points early enough to resolve them proactively. This is one of the biggest selling points for DevOps. According to multiple surveys, there are a few other, perhaps more powerful benefits: DevOps increases trust within an organization and increases job satisfaction among developers, IT staff, and QA personnel. Unsurprisingly, it also improves the morale of project managers.

DevSecOps

It is not all that common for the security team to be involved in software development efforts, but it makes a lot of sense for them to be. Their job is to find vulnerabilities before the threat actors can and then do something about it. As a result, most security professionals develop an “adversarial mindset” that allows them to think like attackers in order to better defend against them. Imagine being a software developer and having someone next to you telling you all the ways they could subvert your code to do bad things. It’d be kind of like having a spell-checker but for vulnerabilities instead of spelling!

DevSecOps is the integration of development, security, and operations professionals into a software development team. It’s just like DevOps but with security added in. One of the main advantages of DevSecOps is that it bakes security right into the development process, rather than bolting it on at the end of it. Rather than implementing controls to mitigate vulnerabilities, the vulnerabilities are prevented from being implemented in the first place.

Other Methodologies

There seems to be no shortage of SDLC and software development methodologies in the industry. The following is a quick summary of a few others that can also be used:

- **Exploratory methodology** A methodology that is used in instances where clearly defined project objectives have not been presented. Instead of focusing on explicit tasks, the exploratory methodology relies on covering a set of specifications likely to affect the final product’s functionality. Testing is an important part of exploratory development, as it ascertains that the current phase of the project is compliant with likely implementation scenarios.
- **Joint Application Development (JAD)** A methodology that uses a team approach in application development in a workshop-oriented environment. This methodology is distinguished by its inclusion of members other than coders in the team. It is common to find executive sponsors, subject matter experts, and end users spending hours or days in collaborative development workshops.

Integrated Product Team

An *integrated product team (IPT)* is a multidisciplinary development team with representatives from many or all the stakeholder populations. The idea makes a lot of sense when you think about it. Why should programmers learn or guess the manner in which the accounting folks handle accounts payable? Why should testers and quality control personnel wait until a product is finished before examining it? Why should the marketing team wait until the project (or at least the prototype) is finished before determining how best to sell it? A comprehensive IPT includes business executives and end users and everyone in between.

The Joint Application Development methodology, in which users join developers during extensive workshops, works well with the IPT approach. IPTs extend this concept by ensuring that the right stakeholders are represented in every phase of the development as formal team members. In addition, whereas JAD is focused on involving the user community, an IPT is typically more inward facing and focuses on bringing in the business stakeholders.

An IPT is not a development methodology. Instead, it is a management technique. When project managers decide to use IPTs, they still have to select a methodology. These days, IPTs are often associated with Agile methodologies.

- **Reuse methodology** A methodology that approaches software development by using progressively developed code. Reusable programs are evolved by gradually modifying preexisting prototypes to customer specifications. Since the reuse methodology does not require programs to be built from scratch, it drastically reduces both development cost and time.
- **Cleanroom** An approach that attempts to prevent errors or mistakes by following structured and formal methods of developing and testing. This approach is used for high-quality and mission-critical applications that will be put through a strict certification process.

We covered only the most commonly used methodologies in this section, but there are many more that exist. New methodologies have evolved as technology and research have advanced and various weaknesses of older approaches have been addressed. Most of the methodologies exist to meet a specific software development need, and choosing the wrong approach for a certain project could be devastating to its overall success.



EXAM TIP While all the methodologies we covered are used in many organizations around the world, you should focus on Agile, Waterfall, DevOps, and DevSecOps for the CISSP exam.

Review of Development Methodologies

A quick review of the various methodologies we have covered up to this point is provided here:

- **Waterfall** Very rigid, sequential approach that requires each phase to complete before the next one can begin. Difficult to integrate changes. Inflexible methodology.
- **Prototyping** Creating a sample or model of the code for proof-of-concept purposes.
- **Incremental** Multiple development cycles are carried out on a piece of software throughout its development stages. Each phase provides a usable version of software.
- **Spiral** Iterative approach that emphasizes risk analysis per iteration. Allows for customer feedback to be integrated through a flexible evolutionary approach.
- **Rapid Application Development** Combines prototyping and iterative development procedures with the goal of accelerating the software development process.
- **Agile** Iterative and incremental development processes that encourage team-based collaboration. Flexibility and adaptability are used instead of a strict process structure.
- **DevOps** The software development and IT operations teams work together at all stages of the project to ensure a smooth transition from development to production environments.
- **DevSecOps** Just like DevOps, but also integrates the security team into every stage of the project.

Maturity Models

Regardless of which software development methodology an organization adopts, it is helpful to have a framework for determining how well-defined and effective its development activities are. Maturity models identify the important components of software development processes and then organize them in an evolutionary scale that proceeds from ad hoc to mature. Each maturity level comprises a set of goals that, when they are met, stabilize one or more of those components. As an organization moves up this maturity scale, the effectiveness, repeatability, and predictability of its software development processes increase, leading to higher-quality code. Higher-quality code, in turn, means fewer vulnerabilities, which is why we care so deeply about this topic as cybersecurity leaders. Let's take a look at the two most popular models: the Capability Maturity Model Integration (CMMI) and the Software Assurance Maturity Model (SAMM).

Capability Maturity Model Integration

Capability Maturity Model Integration (CMMI) is a comprehensive set of models for developing software. It addresses the different phases of a software development life cycle, including concept definition, requirements analysis, design, development, integration, installation, operations, and maintenance, and what should happen in each phase. It can be used to evaluate security engineering practices and identify ways to improve them. It can also be used by customers in the evaluation process of a software vendor. Ideally, software vendors would use the model to help improve their processes, and customers would use the model to assess the vendors' practices.



EXAM TIP For exam purposes, the terms CMM and CMMI are equivalent.

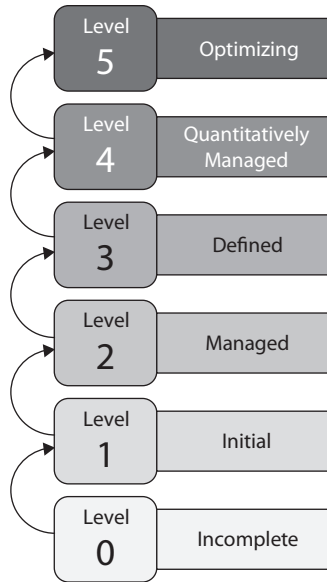
CMMI describes procedures, principles, and practices that underlie software development process maturity. This model was developed to help software vendors improve their development processes by providing an evolutionary path from an ad hoc “fly by the seat of your pants” approach to a more disciplined and repeatable method that improves software quality, reduces the life cycle of development, provides better project management capabilities, allows for milestones to be created and met in a timely manner, and takes a more proactive approach than the less effective reactive approach. It provides best practices to allow an organization to develop a standardized approach to software development that can be used across many different groups. The goal is to continue to review and improve upon the processes to optimize output, increase capabilities, and provide higher-quality software at a lower cost through the implementation of continuous improvement steps.

If the company Stuff-R-Us wants a software development company, Software-R-Us, to develop an application for it, it can choose to buy into the sales hype about how wonderful Software-R-Us is, or it can ask Software-R-Us whether it has been evaluated against CMMI. Third-party companies evaluate software development companies to certify their product development processes. Many software companies have this evaluation done so they can use this as a selling point to attract new customers and provide confidence for their current customers.

The five maturity levels of CMMI are shown in Figure 24-10 and described here:

- **Level 0: Incomplete** Development process is ad hoc or even chaotic. Tasks are not always completed at all, so projects are regularly cancelled or abandoned.
- **Level 1: Initial** The organization does not use effective management procedures and plans. There is no assurance of consistency, and quality is unpredictable. Success is usually the result of individual heroics.
- **Level 2: Managed** A formal management structure, change control, and quality assurance are in place for individual projects. The organization can properly repeat processes throughout each project.

Figure 24-10
CMMI staged
maturity levels



- **Level 3: Defined** Formal procedures are in place that outline and define processes carried out in all projects across the organization. This allows the organization to be proactive rather than reactive.
- **Level 4: Quantitatively Managed** The organization has formal processes in place to collect and analyze quantitative data, and metrics are defined and fed into the process-improvement program.
- **Level 5: Optimizing** The organization has budgeted and integrated plans for continuous process improvement, which allow it to quickly respond to opportunities and changes.

Each level builds upon the previous one. For example, a company that accomplishes a Level 5 CMMI rating must meet all the requirements outlined in Levels 1–4 along with the requirements of Level 5.

If a software development vendor is using the Prototyping methodology that was discussed earlier in this chapter, the vendor would most likely only achieve a CMMI Level 1, particularly if its practices are ad hoc, not consistent, and the level of the quality that its software products contain is questionable. If this company practiced a strict Agile SDLC methodology consistently and carried out development, testing, and documentation precisely, it would have a higher chance of obtaining a higher CMMI level.

Capability maturity models (CMMs) are used for many different purposes, software development processes being one of them. They are general models that allow for

maturity-level identification and maturity improvement steps. We showed how a CMM can be used for organizational security program improvement processes in Chapter 4.

The software industry ended up with several different CMMs, which led to confusion. CMMI was developed to bring many of these different maturity models together and allow them to be used in one framework. CMMI was developed by industry experts, government entities, and the Software Engineering Institute at Carnegie Mellon University. So CMMI has replaced CMM in the software engineering world, but you may still see CMM referred to within the industry and on the CISSP exam. Their ultimate goals are the same, which is process improvement.



NOTE CMMI is continually being updated and improved upon. You can view the latest documents on it at <https://cmminstitute.com/learning/appraisals/levels>.

Software Assurance Maturity Model

The OWASP *Software Assurance Maturity Model (SAMM)* is specifically focused on secure software development and allows organizations of any size to decide their target maturity levels within each of the five critical business functions: Governance, Design, Implementation, Verification, and Operations, as shown in Figure 24-11. One of the premises on which SAMM is built is that any organization that is involved in software development must perform these five functions.

Each business function, in turn, is divided into three security practices, which are sets of security-related activities that provide assurance for the function. For example, if you want to ensure that your Design business function is done right, you need to perform activities related to threat assessment, identification of security requirements, and securely architecting the software. Each of these 15 practices can be independently

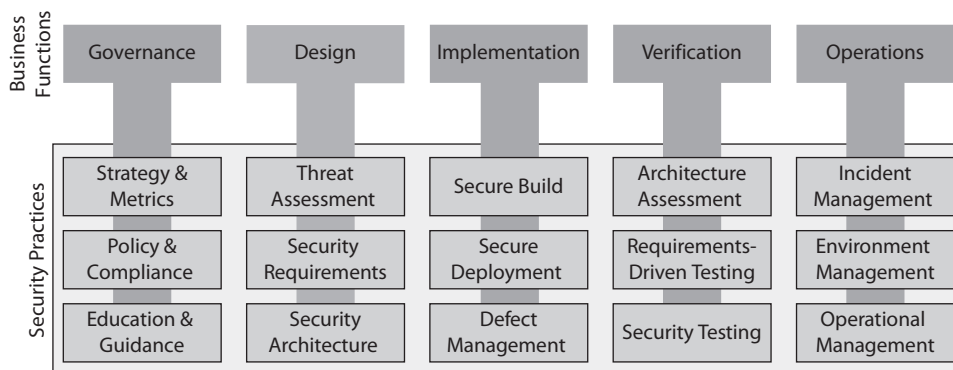


Figure 24-11 Software Assurance Maturity Model

assessed and matured, which allows the organization to decide what maturity level makes sense for each practice.



NOTE You can find more information on SAMM at <https://owasp samm.org/model/>.

Chapter Review

While there is no expectation that you, as a CISSP, will necessarily be involved in software development, you will almost certainly lead organizations that either produce software or consume it. Therefore, it is important that you understand how secure software is developed. Knowing this enables you to see what an organization is doing, software development-wise, and quickly get a sense of the maturity of its processes. If processes are ad hoc, this chapter should have given you some pointers on how to formalize the processes. After all, without formal processes and trained programmers in place, you have almost no hope of producing software that is not immediately vulnerable as soon as it is put into production. On the other hand, if the organization seems more mature, you can delve deeper into the specifics of building security into the software, which is the topic of the next chapter.

Quick Review

- The software development life cycle (SDLC) comprises five phases: requirements gathering, design, development, testing, and operations and maintenance (O&M).
- Computer-aided software engineering (CASE) refers to any type of software that allows for the automated development of software, which can come in the form of program editors, debuggers, code analyzers, version-control mechanisms, and more. The goals are to increase development speed and productivity and reduce errors.
- Various levels of testing should be carried out during development: unit (testing individual components), integration (verifying components work together in the production environment), acceptance (ensuring code meets customer requirements), and regression (testing after changes take place).
- Change management is a systematic approach to deliberately regulating the changing nature of projects. Change control, which is a subpart of change management, deals with controlling specific changes to a system.
- Security should be addressed in each phase of software development. It should not be addressed only at the end of development because of the added cost, time, and effort and the lack of functionality.

- The attack surface is the collection of possible entry points for an attacker. The reduction of this surface reduces the possible ways that an attacker can exploit a system.
- Threat modeling is a systematic approach used to understand how different threats could be realized and how a successful compromise could take place.
- The waterfall software development methodology follows a sequential approach that requires each phase to complete before the next one can begin.
- The prototyping methodology involves creating a sample of the code for proof-of-concept purposes.
- Incremental software development entails multiple development cycles that are carried out on a piece of software throughout its development stages.
- The spiral methodology is an iterative approach that emphasizes risk analysis per iteration.
- Rapid Application Development (RAD) combines prototyping and iterative development procedures with the goal of accelerating the software development process.
- Agile methodologies are characterized by iterative and incremental development processes that encourage team-based collaboration, where flexibility and adaptability are used instead of a strict process structure.
- Some organizations improve internal coordination and reduce friction by integrating the development and operations (DevOps) teams or the development, operations, and security (DevSecOps) teams when developing software.
- An integrated product team (IPT) is a multidisciplinary development team with representatives from many or all the stakeholder populations.
- Capability Maturity Model Integration (CMMI) is a process improvement approach that provides organizations with the essential elements of effective processes, which will improve their performance.
- The CMMI model uses six maturity levels designated by the numbers 0 through 5. Each level represents the maturity level of the process quality and optimization. The levels are organized as follows: 0 = Incomplete, 1 = Initial, 2 = Managed, 3 = Defined, 4 = Quantitatively Managed, 5 = Optimizing.
- The OWASP Software Assurance Maturity Model (SAMM) is specifically focused on secure software development and allows organizations to decide their target maturity levels within each of five critical business functions: Governance, Design, Implementation, Verification, and Operations.

Questions

Please remember that these questions are formatted and asked in a certain way for a reason. Keep in mind that the CISSP exam is asking questions at a conceptual level. Questions may not always have the perfect answer, and the candidate is advised against always looking for the perfect answer. Instead, the candidate should look for the best answer in the list.

1. The software development life cycle has several phases. Which of the following lists these phases in the correct order?
 - A. Requirements gathering, design, development, maintenance, testing, release
 - B. Requirements gathering, design, development, testing, operations and maintenance
 - C. Prototyping, build and fix, increment, test, maintenance
 - D. Prototyping, testing, requirements gathering, integration, testing
2. John is a manager of the application development department within his company. He needs to make sure his team is carrying out all of the correct testing types and at the right times of the development stages. Which of the following accurately describe types of software testing that should be carried out?
 - i. **Unit testing** Testing individual components in a controlled environment where programmers validate data structure, logic, and boundary conditions
 - ii. **Integration testing** Verifying that components work together as outlined in design specifications
 - iii. **Acceptance testing** Ensuring that the code meets customer requirements
 - iv. **Regression testing** After a change to a system takes place, retesting to ensure functionality, performance, and protection
 - A. i, ii
 - B. ii, iii
 - C. i, ii, iv
 - D. i, ii, iii, iv
3. Marge has to choose a software development methodology that her team should follow. The application that her team is responsible for developing is a critical application that can have few to no errors. Which of the following best describes the type of methodology her team should follow?
 - A. Cleanroom
 - B. Joint Application Development (JAD)
 - C. Rapid Application Development (RAD)
 - D. Reuse methodology

4. Which level of Capability Maturity Model Integration allows organizations to manage all projects across the organization and be proactive?
 - A. Defined
 - B. Incomplete
 - C. Managed
 - D. Optimizing
5. Mohammed is in charge of a large software development project with rigid requirements and phases that will probably be completed by different contractors. Which methodology would be best?
 - A. Waterfall
 - B. Spiral
 - C. Prototyping
 - D. Agile

Use the following scenario to answer Questions 6–9. You're in charge of IT and security at a midsize organization going through a growth stage. You decided to stand up your own software development team and are about to start your first project: a knowledge base for your customers. You think it can eventually grow to become the focal point of interaction with your customers, offering a multitude of features. You've heard a lot about the Scrum methodology and decide to try it for this project.

6. How would you go about documenting the requirements for this software system?
 - A. User stories
 - B. Use cases
 - C. System Requirements Specification (SRS)
 - D. Informally, since it's your first project
7. You are halfway through your first Scrum sprint and get a call from a senior vice president insisting that you add a new feature immediately. How do you handle this request?
 - A. Add the feature to the next sprint
 - B. Change the current sprint to include the feature
 - C. Reset the project to the requirements gathering phase
 - D. Delay the new feature until the end of the project

8. Your software development team, being new to the organization, is struggling to work smoothly with other teams within the organization as needed to get the software into production securely. Which approach can help mitigate this internal friction?
 - A. DevSecOps
 - B. DevOps
 - C. Integrated Product Teams (IPT)
 - D. Joint Analysis Design (JAD) sessions
9. What would be the best approach to selectively mature your software development practices with a view to improving cybersecurity?
 - A. Software Assurance Maturity Model (SAMM)
 - B. Capability Maturity Model Integration (CMMI)
 - C. Kanban
 - D. Integrated product teams (IPTs)

Answers

1. **B.** The following outlines the common phases of the software development life cycle:
 - i. Requirements gathering
 - ii. Design
 - iii. Development
 - iv. Testing
 - v. Operations and maintenance
2. **D.** There are different types of tests the software should go through because there are different potential flaws to look for. The following are some of the most common testing approaches:
 - **Unit testing** Testing individual components in a controlled environment where programmers validate data structure, logic, and boundary conditions
 - **Integration testing** Verifying that components work together as outlined in design specifications
 - **Acceptance testing** Ensuring that the code meets customer requirements
 - **Regression testing** After a change to a system takes place, retesting to ensure functionality, performance, and protection
3. **A.** The listed software development methodologies and their definitions are as follows:
 - **Joint Application Development (JAD)** A methodology that uses a team approach in application development in a workshop-oriented environment.

- **Rapid Application Development (RAD)** A methodology that combines the use of prototyping and iterative development procedures with the goal of accelerating the software development process.
 - **Reuse methodology** A methodology that approaches software development by using progressively developed code. Reusable programs are evolved by gradually modifying preexisting prototypes to customer specifications. Since the reuse methodology does not require programs to be built from scratch, it drastically reduces both development cost and time.
 - **Cleanroom** An approach that attempts to prevent errors or mistakes by following structured and formal methods of developing and testing. This approach is used for high-quality and critical applications that will be put through a strict certification process.
4. **A.** The six levels of Capability Maturity Integration Model are
- **Incomplete** Development process is ad hoc or even chaotic. Tasks are not always completed at all, so projects are regularly cancelled or abandoned.
 - **Initial** The organization does not use effective management procedures and plans. There is no assurance of consistency, and quality is unpredictable. Success is usually the result of individual heroics.
 - **Managed** A formal management structure, change control, and quality assurance are in place for individual projects. The organization can properly repeat processes throughout each project.
 - **Defined** Formal procedures are in place that outline and define processes carried out in all projects across the organization. This allows the organization to be proactive rather than reactive.
 - **Quantitatively Managed** The organization has formal processes in place to collect and analyze quantitative data, and metrics are defined and fed into the process-improvement program.
 - **Optimizing** The organization has budgeted and integrated plans for continuous process improvement, which allow it to quickly respond to opportunities and changes.
5. **D.** The Waterfall methodology is a very rigid approach that could be useful for projects in which all the requirements are fully understood up front or projects for which different organizations will perform the work at each phase. The Spiral, prototyping, and Agile methodologies are well suited for situations in which the requirements are not well understood, and don't lend themselves well to switching contractors midstream.
6. **A.** Any answer except “informally” would be a reasonable one, but since you are using an Agile methodology (Scrum), user stories is the best answer. The important point is that you document the requirements formally, so you can design a solution that meets all your users' needs.

- 7. **A.** The Scrum methodology allows the project to be reset by allowing product features to be added, changed, or removed at clearly defined points that typically happen at the conclusion of each sprint.
- 8. **A.** DevSecOps is the integration of development, security, and operations professionals into a software development team. This is a good way to solve the friction between developers and members of the security and operations staff.
- 9. **A.** CMMI and SAMM are the only maturity models among the possible answers. SAMM is the best answer because it allows for more granular maturity goals than CMMI does, and it is focused on security.

Secure Software

This chapter presents the following:

- Programming languages
- Secure coding
- Security controls for software development
- Software security assessments
- Assessing the security of acquired software

*A good programmer is someone who always looks both ways
before crossing a one-way street.*

—Doug Linder

Quality can be defined as fitness for purpose. In other words, quality refers to how good or bad something is for its intended purpose. A high-quality car is good for transportation. We don't have to worry about it breaking down, failing to protect its occupants in a crash, or being easy for a thief to steal. When we need to go somewhere, we can count on a high-quality car to get us to wherever we need to go. Similarly, we don't have to worry about high-quality software crashing, corrupting our data under unforeseen circumstances, or being easy for someone to subvert. Sadly, many developers still think of functionality first (or only) when thinking about quality. When we look at it holistically, we see that quality is the most important concept in developing secure software.

Every successful compromise of a software system relies on the exploitation of one or more vulnerabilities in it. Software vulnerabilities, in turn, are caused by defects in the design or implementation of code. The goal, then, is to develop software that is as free from defects or, in other words, as high quality as we can make it. In this chapter, we will discuss how secure software is quality software. We can't have one without the other. By applying the right processes, controls, and assessments, the outcome will be software that is more reliable and more difficult to exploit or subvert. Of course, these principles apply equally to software we develop in our own organizations and software that is developed for us by others.

Programming Languages and Concepts

All software is written in some type of programming language. Programming languages have gone through several generations over time, each generation building on the next, providing richer functionality and giving programmers more powerful tools as they evolve.

The main categories of languages are machine, assembly, high-level, very high-level, and natural languages. *Machine language* is in a format that the computer's processor can understand and work with directly. Every processor family has its own machine code instruction set, which is represented in a binary format (1 and 0) and is the most fundamental form of programming language. Since this was pretty much the only way to program the very first computers in the early 1950s, machine languages are the first generation of programming languages. Early computers used only basic binary instructions because compilers and interpreters were nonexistent at the time. Programmers had to manually calculate and allot memory addresses and sequentially feed instructions, as there was no concept of abstraction. Not only was programming in binary extremely time consuming, it was also highly prone to errors. (If you think about writing out thousands of 1's and 0's to represent what you want a computer to do, this puts this approach into perspective.) This forced programmers to keep a tight rein on their program lengths, resulting in programs that were very rudimentary.

An *assembly language* is considered a low-level programming language and is the symbolic representation of machine-level instructions. It is "one step above" machine language. It uses symbols (called mnemonics) to represent complicated binary codes. Programmers using assembly language could use commands like ADD, PUSH, POP, etc., instead of the binary codes (1001011010, etc.). Assembly languages use programs called *assemblers*, which automatically convert these assembly codes into the necessary machine-compatible binary language. To their credit, assembly languages drastically reduced programming and debugging times, introduced the concept of variables, and freed programmers from manually calculating memory addresses. But like machine code, programming in an assembly language requires extensive knowledge of a computer's architecture. It is easier than programming in binary format, but more challenging compared to the high-level languages most programmers use today.

Programs written in assembly language are also hardware specific, so a program written for an ARM-based processor would be incompatible with Intel-based systems; thus, these types of languages are not portable. Once the program is written, it is fed to an assembler, which translates the assembly language into machine language. The assembler also replaces variable names in the assembly language program with actual addresses at which their values will be stored in memory.



NOTE Assembly language allows for direct control of very basic activities within a computer system, as in pushing data on a memory stack and popping data off a stack. Attackers commonly use assembly language to tightly control how malicious instructions are carried out on victim systems.

The third generation of programming languages started to emerge in the early 1960s. They are known as *high-level languages* because of their refined programming structures.

High-level languages use abstract statements. Abstraction naturalizes multiple assembly language instructions into a single high-level statement, such as IF – THEN – ELSE. This allows programmers to leave low-level (system architecture) intricacies to the programming language and focus on their programming objectives. In addition, high-level languages are easier to work with compared to machine and assembly languages, as their syntax is similar to human languages. The use of mathematical operators also simplifies arithmetic and logical operations. This drastically reduces program development time and allows for more simplified debugging. This means the programs are easier to write and mistakes (bugs) are easier to identify. High-level languages are processor independent. Code written in a high-level language can be converted to machine language for different processor architectures using compilers and interpreters. When code is independent of a specific processor type, the programs are portable and can be used on many different system types.

Fourth-generation languages (*very high-level languages*) were designed to further enhance the natural language approach instigated within the third-generation languages. They focus on highly abstract algorithms that allow straightforward programming implementation in specific environments. The most remarkable aspect of fourth-generation languages is that the amount of manual coding required to perform a specific task may be ten times less than for the same task on a third-generation language. This is an especially important feature because these languages have been developed to be used by inexperienced users and not just professional programmers.

As an analogy, let's say that you need to pass a calculus exam. You need to be very focused on memorizing the necessary formulas and applying the formulas to the correct word problems on the test. Your focus is on how calculus works, not on how the calculator you use as a tool works. If you had to understand how your calculator is moving data from one transistor to the other, how the circuitry works, and how the calculator stores and carries out its processing activities just to use it for your test, this would be overwhelming. The same is true for computer programmers. If they had to worry about how the operating system carries out memory management functions, input/output activities, and how processor-based registers are being used, it would be difficult for them to also focus on real-world problems they are trying to solve with their software. Very high-level languages hide all of this background complexity and take care of it for the programmer.

The early 1990s saw the conception of the fifth generation of programming languages (*natural languages*). These languages approach programming from a completely different perspective. Program creation does not happen through defining algorithms and function statements, but rather by defining the constraints for achieving a specified result. The goal is to create software that can solve problems by itself instead of a programmer having to develop code to deal with individual and specific problems. The applications work more like a black box—a problem goes in and a solution comes out. Just as the introduction of assembly language eliminated the need for binary-based programming, the full impact of fifth-generation programming techniques may bring to an end the traditional programming approach. The ultimate target of fifth-generation languages is to eliminate the need for programming expertise and instead use advanced knowledge-based processing and artificial intelligence.

Language Levels

The “higher” the language, the more abstraction that is involved, which means the language hides details of how it performs its tasks from the software developer. A programming language that provides a high level of abstraction frees the programmer from the need to worry about the intricate details of the computer system itself, as in registers, memory addresses, complex Boolean expressions, thread management, and so forth. The programmer can use simple statements such as “print” and does not need to worry about how the computer will actually get the data over to the printer. Instead, the programmer can focus on the core functionality that the application is supposed to provide and not be bothered with the complex things taking place in the belly of the operating system and motherboard components.

As an analogy, you do not need to understand how your engine or brakes work in your car—there is a level of abstraction. You just turn the steering wheel and step on the pedal when necessary, and you can focus on getting to your destination.

There are so many different programming languages today, it is hard to fit them neatly in the five generations described in this chapter. These generations are the classical way of describing the differences in software programming approaches and what you will see on the CISSP exam.

The industry has not been able to fully achieve all the goals set out for these fifth-generation languages. The human insight of programmers is still necessary to figure out the problems that need to be solved, and the restrictions of the structure of a current computer system do not allow software to “think for itself” yet. We are getting closer to achieving artificial intelligence within our software, but we still have a long way to go.

The following lists the basic software programming language generations:

- **Generation one** Machine language
- **Generation two** Assembly language
- **Generation three** High-level language
- **Generation four** Very high-level language
- **Generation five** Natural language

Assemblers, Compilers, Interpreters

No matter what type or generation of programming language is used, all of the instructions and data have to end up in a binary format for the processor to understand and work with. Just like our food has to be broken down into specific kinds of molecules for our body to be able to use it, all code must end up in a format that is consumable by specific systems. Each programming language type goes through this transformation through the use of assemblers, compilers, or interpreters.

Assemblers are tools that convert assembly language source code into machine language code. Assembly language consists of mnemonics, which are incomprehensible to processors and therefore need to be translated into operation instructions.

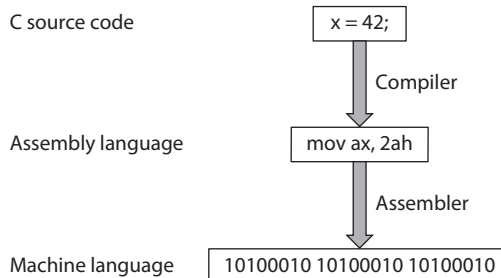
Compilers are tools that convert high-level language statements into the necessary machine-level format (.exe, .dll, etc.) for specific processors to understand. The compiler transforms instructions from a source language (high-level) to a target language (machine), sometimes using an external assembler along the way. This transformation allows the code to be executable. A programmer may develop an application in the C language, but when you purchase this application, you do not receive the source code; instead, you receive the executable code that runs on your type of computer. The source code was put through a compiler, which resulted in an executable file that can run on your specific processor type.

Compilers allow developers to create software code that can be developed once in a high-level language and compiled for various platforms. So, you could develop one piece of software, which is then compiled by five different compilers to allow it to be able to run on five different systems.

Figure 25-1 shows the process by which a high-level language is gradually transformed into machine language, which is the only language a processor can understand natively. In this example, we have a statement that assigns the value 42 to the variable *x*. Once we feed the program containing this statement to a compiler, we end up with assembly language, which is shown in the middle of the figure. The way to set the value of a variable in assembly language is to literally move that value into wherever the variable is being stored. In this example, we are moving the hexadecimal value for 42 (which is 2a in hexadecimal, or 2ah) into the ax register in the processor. In order for the processor to execute this command, however, we still have to convert it into machine language, which is the job of the assembler. Note that it is way easier for a human coder to write *x = 42* than it is to represent the same operation in either assembly or (worse yet) machine language.

If a programming language is considered “interpreted,” then a tool called an *interpreter* takes care of transforming high-level code to machine-level code. For example, applications that are developed in JavaScript, Python, or Perl can be run directly by an interpreter, without having to be compiled. The goal is to improve portability. The greatest advantage of executing a program in an interpreted environment is that the platform independence and memory management functions are part of an interpreter. The major disadvantage

Figure 25-1
Converting
a high-level
language
statement
into machine
language code



with this approach is that the program cannot run as a stand-alone application, requiring the interpreter to be installed on the local machine.



NOTE Some languages, such as Java and Python, blur the lines between interpreted and compiled languages by supporting both approaches. We'll talk more about how Java does this in the next section.

From a security point of view, it is important to understand vulnerabilities that are inherent in specific programming languages. For example, programs written in the C language could be vulnerable to buffer overrun and format string errors. The issue is that some of the C standard software libraries do not check the length of the strings of data they manipulate by default. Consequently, if a string is obtained from an untrusted source (i.e., the Internet) and is passed to one of these library routines, parts of memory may be unintentionally overwritten with untrustworthy data—this vulnerability can potentially be used to execute arbitrary and malicious software. Some programming languages, such as Java, perform automatic memory allocation as more space is needed; others, such as C, require the developer to do this manually, thus leaving opportunities for error.

Garbage collection is an automated way for software to carry out part of its memory management tasks. A *garbage collector* identifies blocks of memory that were once allocated but are no longer in use and deallocates the blocks and marks them as free. It also gathers scattered blocks of free memory and combines them into larger blocks. It helps provide a more stable environment and does not waste precious memory. If garbage collection does not take place properly, not only can memory be used in an inefficient manner, an attacker could carry out a denial-of-service attack specifically to artificially commit all of a system's memory, rendering the system unable to function.

Nothing in technology seems to be getting any simpler, which makes learning this stuff much harder as the years go by. Ten years ago assembly, compiled, and interpreted languages were more clear-cut and their definitions straightforward. For the most part, only scripting languages required interpreters, but as languages have evolved they have become extremely flexible to allow for greater functionality, efficiency, and portability. Many languages can have their source code compiled or interpreted depending upon the environment and user requirements.

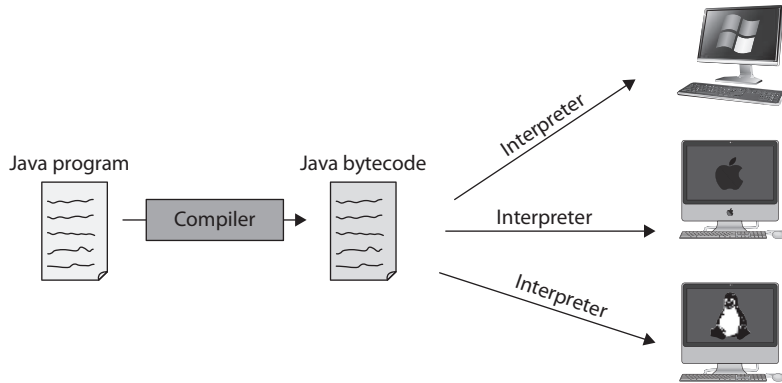
Runtime Environments

What if you wanted to develop software that could run on many different environments without having to recompile it? This is known as *portable code* and it needs something that can sort of “translate” it to each different environment. That “translator” could be tuned to a particular type of computer but be able to run any of the portable code it understands. This is the role of *runtime environments (RTEs)*, which function as miniature operating systems for the program and provide all the resources portable code needs. One of the best examples of RTE usage is the Java programming language.

Java is platform independent because it creates intermediate code, *bytecode*, which is not processor-specific. The *Java Virtual Machine (JVM)* converts the bytecode to the machine

Figure 25-2

The JVM interprets bytecode to machine code for that specific platform.



code that the processor on that particular system can understand (see Figure 25-2). Despite its name, the JVM is not a full-fledged VM (as defined in Chapter 7). Instead, it is a component of the Java RTE, together with a bunch of supporting files like class libraries.

Let's quickly walk through these steps:

1. A programmer creates a Java applet and runs it through a compiler.
2. The Java compiler converts the source code into bytecode (not processor-specific).
3. The user downloads the Java applet.
4. The JVM converts the bytecode into machine-level code (processor-specific).
5. The applet runs when called upon.

When an applet is executed, the JVM creates a unique RTE for it called a *sandbox*. This sandbox is an enclosed environment in which the applet carries out its activities. Applets are commonly sent over within a requested web page, which means the applet executes as soon as it arrives. It can carry out malicious activity on purpose or accidentally if the developer of the applet did not do his part correctly. So the sandbox strictly limits the applet's access to any system resources. The JVM mediates access to system resources to ensure the applet code behaves and stays within its own sandbox. These components are illustrated in Figure 25-3.



NOTE The Java language itself provides protection mechanisms, such as garbage collection, memory management, validating address usage, and a component that verifies adherence to predetermined rules.

However, as with many other things in the computing world, the bad guys have figured out how to escape the confines and restrictions of the sandbox. Programmers have figured out how to write applets that enable the code to access hard drives and

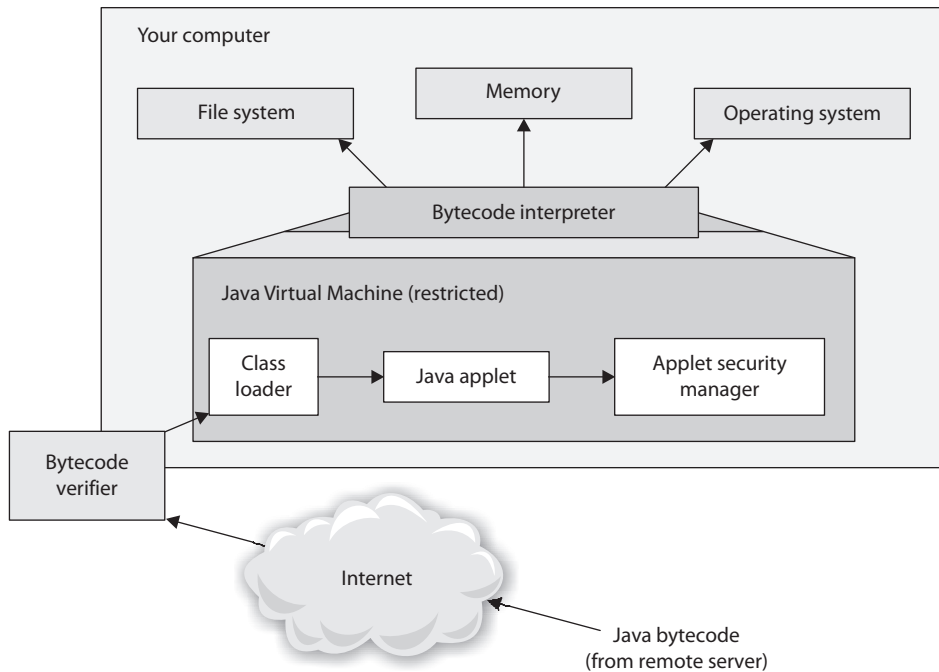
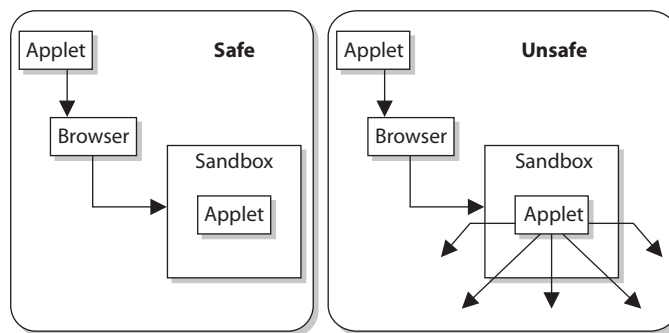


Figure 25-3 Java's security model

resources that are supposed to be protected by the Java security scheme. This code can be malicious in nature and cause destruction and mayhem to the user and her system.



Object-Oriented Programming Concepts

Software development used to be done by classic input–processing–output methods. This development used an information flow model from hierarchical information structures. Data was input into a program, and the program passed the data from the beginning to

end, performed logical procedures, and returned a result. *Object-oriented programming (OOP)* methods perform the same functionality, but with different techniques that work in a more efficient manner. First, you need to understand the basic concepts of OOP.

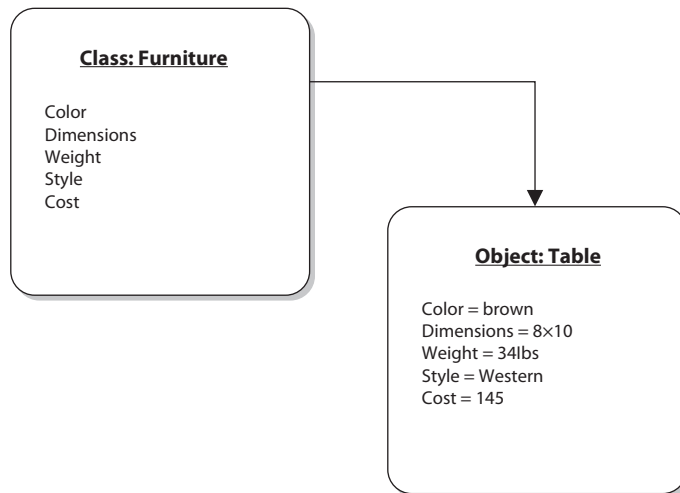
OOP works with classes and objects. A real-world object, such as a table, is a member (or an instance) of a larger class of objects called “furniture.” The furniture class has a set of attributes associated with it, and when an object is generated, it inherits these attributes. The attributes may be color, dimensions, weight, style, and cost. These attributes apply if a chair, table, or loveseat object is generated, also referred to as *instantiated*. Because the table is a member of the class furniture, the table inherits all attributes defined for the class (see Figure 25-4).

The programmer develops the class and all of its characteristics and attributes. The programmer does not develop each and every object, which is the beauty of this approach. As an analogy, let’s say you developed an advanced coffee maker with the goal of putting Starbucks out of business. A customer punches the available buttons on your coffee maker interface, ordering a large latte, with skim milk, vanilla and raspberry flavoring, and an extra shot of espresso, where the coffee is served at 250 degrees. Your coffee maker does all of this through automation and provides the customer with a lovely cup of coffee exactly to her liking. The next customer wants a mocha Frothy Frappé, with whole milk and extra foam. So the goal is to make something once (coffee maker, class), allow it to accept requests through an interface, and create various results (cups of coffee, objects) depending upon the requests submitted.

But how does the class create objects based on requests? A piece of software that is written in OOP will have a request sent to it, usually from another object. The requesting object wants a new object to carry out some type of functionality. Let’s say that object A wants object B to carry out subtraction on the numbers sent from A to B. When this request comes in, an object is built (instantiated) with all of the necessary programming code. Object B carries out the subtraction task and sends the result back to object A.

Figure 25-4

In object-oriented inheritance, each object belongs to a class and takes on the attributes of that class



It does not matter what programming language the two objects are written in; what matters is if they know how to communicate with each other. One object can communicate with another object if it knows the application programming interface (API) communication requirements. An API is the mechanism that allows objects to talk to each other (as described in depth in the forthcoming section “Application Programming Interfaces”). Let’s say you want to talk to Jorge, but can only do so by speaking French and can only use three phrases or less, because that is all Jorge understands. As long as you follow these rules, you can talk to Jorge. If you don’t follow these rules, you can’t talk to Jorge.



TIP An object is an instance of a class.

What’s so great about OOP? Figure 25-5 shows the difference between OOP and procedural programming, which is a non-OOP technique. Procedural programming is built on the concept of dividing a task into procedures that, when executed, accomplish the task. This means that large applications can quickly become one big pile of code (sometimes called *spaghetti code*). If you want to change something in this pile, you have to go through all the program’s procedures to figure out what your one change is going to break. If the program contains hundreds or thousands of lines of code, this is not an easy or enjoyable task. Now, if you choose to write your program in an object-oriented

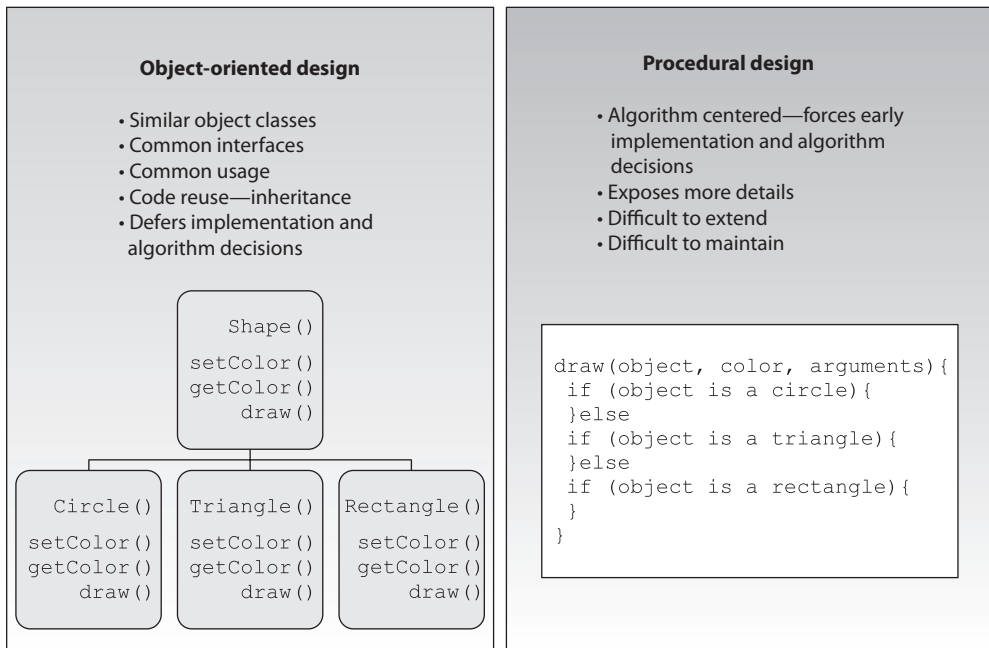


Figure 25-5 Procedural vs. object-oriented programming

language, you don't have one monolithic application, but an application that is made up of smaller components (objects). If you need to make changes or updates to some functionality in your application, you can just change the code within the class that creates the object carrying out that functionality, and you don't have to worry about everything else the program actually carries out. The following breaks down the benefits of OOP:

- **Modularity** The building blocks of software are autonomous objects, cooperating through the exchange of messages.
- **Deferred commitment** The internal components of an object can be redefined without changing other parts of the system.
- **Reusability** Classes are reused by other programs, though they may be refined through inheritance.
- **Naturalness** Object-oriented analysis, design, and modeling map to business needs and solutions.

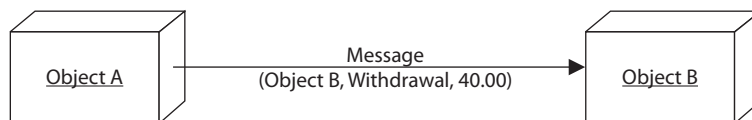
Most applications have some type of functionality in common. Instead of developing the same code to carry out the same functionality for ten different applications, using OOP allows you to create the object only once and reuse it in other applications. This reduces development time and saves money.

Now that we've covered the concepts of OOP, let's clarify the terminology. A *method* is the functionality or procedure an object can carry out. An object may be constructed to accept data from a user and to reformat the request so a back-end server can understand and process it. Another object may perform a method that extracts data from a database and populates a web page with this information. Or an object may carry out a withdrawal procedure to allow the user of an ATM to extract money from her account.

The objects *encapsulate* the attribute values, which means this information is packaged under one name and can be reused as one entity by other objects. Objects need to be able to communicate with each other, and this happens by using *messages* that are sent to the receiving object's API. If object A needs to tell object B that a user's checking account must be reduced by \$40, it sends object B a message. The message is made up of the destination, the method that needs to be performed, and the corresponding arguments. Figure 25-6 shows this example.

Messaging can happen in several ways. A given object can have a single connection (one-to-one) or multiple connections (one-to-many). It is important to map these communication paths to identify if information can flow in a way that is not intended. This helps to ensure that sensitive data cannot be passed to objects of a lower security level.

Figure 25-6
Objects
communicate
via messages.



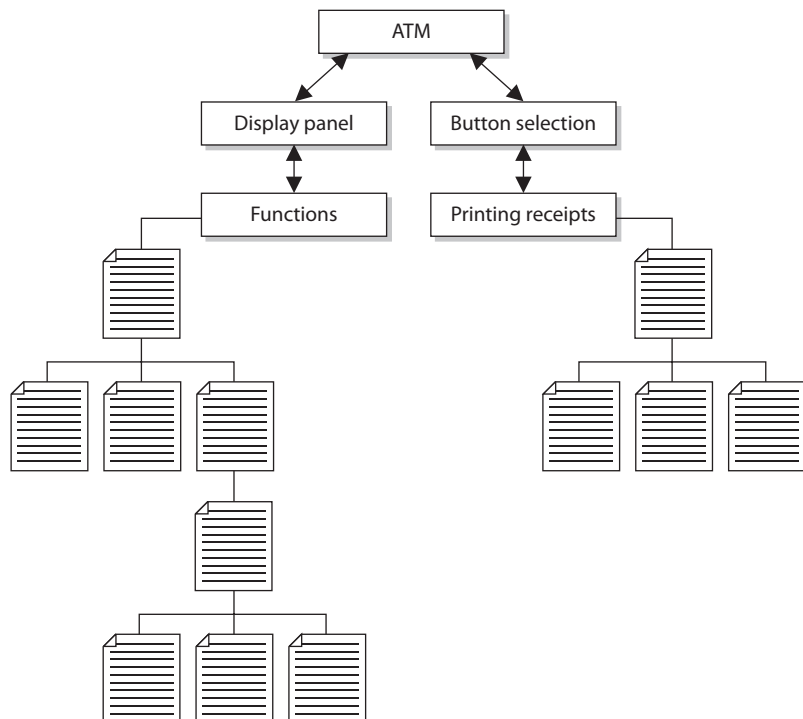
An object can have a shared portion and a private portion. The *shared* portion is the interface (API) that enables it to interact with other components. Messages enter through the interface to specify the requested operation, or method, to be performed. The *private* portion of an object is how it actually works and performs the requested operations. Other components need not know how each object works internally—only that it does the job requested of it. This is how *data hiding* is possible. The details of the processing are hidden from all other program elements outside the object. Objects communicate through well-defined interfaces; therefore, they do not need to know how each other works internally.



NOTE Data hiding is provided by encapsulation, which protects an object's private data from outside access. No object should be allowed to, or have the need to, access another object's internal data or processes.

These objects can grow to great numbers, so the complexity of understanding, tracking, and analyzing can get a bit overwhelming. Many times, the objects are shown in connection to a reference or pointer in documentation. Figure 25-7 shows how related objects are represented as a specific piece, or reference, in a bank ATM system. This enables analysts and developers to look at a higher level of operation and procedures without having to view each individual object and its code. Thus, this modularity provides for a more easily understood model.

Figure 25-7
Object
relationships
within a program



Abstraction, as discussed earlier, is the capability to suppress unnecessary details so the important, inherent properties can be examined and reviewed. It enables the separation of conceptual aspects of a system. For example, if a software architect needs to understand how data flows through the program, she would want to understand the big pieces of the program and trace the steps the data takes from first being input into the program all the way until it exits the program as output. It would be difficult to understand this concept if the small details of every piece of the program were presented. Instead, through abstraction, all the details are suppressed so the software architect can understand a crucial part of the product. It is like being able to see a forest without having to look at each and every tree.

Each object should have specifications it adheres to. This discipline provides cleaner programming and reduces programming errors and omissions. The following list is an example of what should be developed for each object:

- Object name
- Attribute descriptions
- Attribute name
- Attribute content
- Attribute data type
- External input to object
- External output from object
- Operation descriptions
- Operation name
- Operation interface description
- Operation processing description
- Performance issues
- Restrictions and limitations
- Instance connections
- Message connections

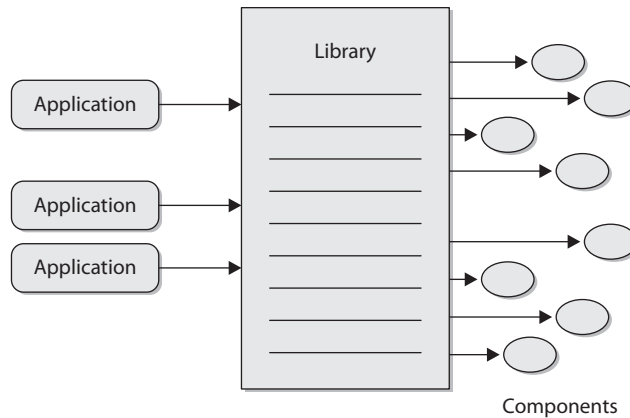
The developer creates a class that outlines these specifications. When objects are instantiated, they inherit these attributes.

Each object can be reused as stated previously, which is the beauty of OOP. This enables a more efficient use of resources and the programmer's time. Different applications can use the same objects, which reduces redundant work, and as an application grows in functionality, objects can be easily added and integrated into the original structure.

The objects can be catalogued in a library, which provides an economical way for more than one application to call upon the objects (see Figure 25-8). The library provides an index and pointers to where the objects actually live within the system or on another system.

Figure 25-8

Applications locate the necessary objects through a library index.



When applications are developed in a modular approach, like object-oriented methods, components can be reused, complexity is reduced, and parallel development can be done. These characteristics allow for fewer mistakes, easier modification, resource efficiency, and more timely coding than the classic programming languages. OOP also provides functional independence, which means each module addresses a specific subfunction of requirements and has an interface that is easily understood by other parts of the application.

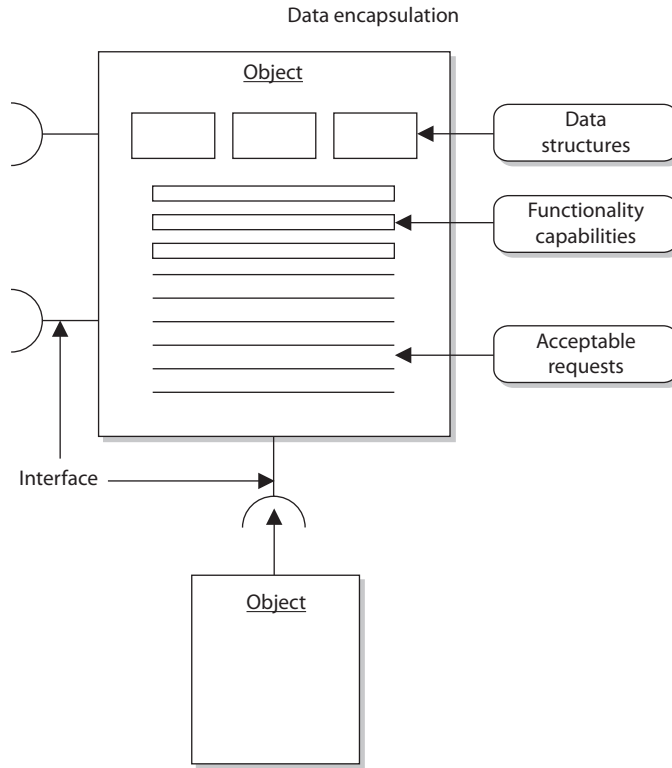
An object is *encapsulated*, meaning the data structure (the operation's functionality) and the acceptable ways of accessing it are grouped into one entity. Other objects, subjects, and applications can use this object and its functionality by accessing it through controlled and standardized interfaces and sending it messages (see Figure 25-9).

Cohesion and Coupling

Cohesion reflects how many different types of tasks a module can carry out. If a module carries out only one task (i.e., subtraction) or tasks that are very similar (i.e., subtract, add, multiply), it is described as having high cohesion, which is a good thing. The higher the cohesion, the easier it is to update or modify the module and not affect other modules that interact with it. This also means the module is easier to reuse and maintain because it is more straightforward when compared to a module with low cohesion. An object with low cohesion carries out multiple *different* tasks and increases the complexity of the module, which makes it harder to maintain and reuse. So, you want your objects to be focused, manageable, and understandable. Each object should carry out a single function or similar functions. One object should not carry out mathematical operations, graphic rendering, and cryptographic functions—these are separate functionality types, and keeping track of this level of complexity would be confusing. If you are attempting to create complex multifunction objects, you are trying to shove too much into one object. Objects should carry out modular, simplistic functions—that is the whole point of OOP.

Figure 25-9

The different components of an object and the way it works are hidden from other objects.



Coupling is a measurement that indicates how much interaction one module requires to carry out its tasks. If a module has low (loose) coupling, this means the module does not need to communicate with many other modules to carry out its job. High (tight) coupling means a module depends upon many other modules to carry out its tasks. Low coupling is more desirable because the module is easier to understand and easier to reuse, and it can be changed without affecting many modules around it. Low coupling indicates that the programmer created a well-structured module. As an analogy, a company would want its employees to be able to carry out their individual jobs with the least amount of dependencies on other workers. If Joe has to talk with five other people just to get one task done, too much complexity exists, the task is too time-consuming, and the potential for errors increases with every interaction.

If modules are tightly coupled, the ripple effect of changing just one module can drastically affect the other modules. If they are loosely coupled, this level of complexity decreases.

An example of *low coupling* would be one module passing a variable value to another module. As an example of *high coupling*, Module A would pass a value to Module B, another value to Module C, and yet another value to Module D. Module A could not complete its tasks until Modules B, C, and D completed their tasks and returned results to Module A.



EXAM TIP Objects should be self-contained and perform a single logical function, which is high cohesion. Objects should not drastically affect each other, which is low coupling.

The level of complexity involved with coupling and cohesion can directly impact the security level of a program. The more complex something is, the harder it is to secure. Developing “tight code” not only allows for efficiencies and effectiveness but also reduces the software’s attack surface. Decreasing complexity where possible reduces the number of potential holes a bad guy can sneak through. As an analogy, if you were responsible for protecting a facility, your job would be easier if the facility had a small number of doors, windows, and people coming in and out of it. The smaller number of variables and moving pieces would help you keep track of things and secure them.

Application Programming Interfaces

When we discussed some of the attributes of object-oriented development, we spent a bit of time on the concept of abstraction. Essentially, abstraction is all about defining *what* a class or object does and ignoring *how* it accomplishes it internally. An *application programming interface (API)* specifies the manner in which a software component interacts with other software components. We already saw in Chapter 9 how this can come in handy in the context of Trusted Platform Modules (TPMs) and how they constrain communications with untrusted modules. APIs create checkpoints where security controls can be easily implemented. Furthermore, they encourage software reuse and also make the software more maintainable by localizing the changes that need to be made while eliminating (or at least reducing) cascading effects of fixes or changes.

Besides the advantages of reduced effort and improved maintainability, APIs often are required to employ the underlying operating system’s functionality. Apple macOS and iOS, Google Android, and Microsoft Windows all require developers to use standard APIs for access to operating system functionality such as opening and closing files and network connections, among many others. All these major vendors restrict the way in which their APIs are used, most notably by ensuring that any parameter that is provided to them is first checked to ensure it is not malformed, invalid, or malicious, which is something we should all do when we are dealing with APIs.

Parameter validation refers to confirming that the parameter values being received by an application are within defined limits before they are processed by the system. In a client/server architecture, validation controls may be placed on the client side prior to submitting requests to the server. Even when these controls are employed, the server should perform parallel validation of inputs prior to processing them because a client has fewer controls than a server and may have been compromised or bypassed.

Software Libraries

APIs are perhaps most familiar to us in the context of software libraries. A *software library* is a collection of components that do specific tasks that are useful to many other components. For example, there are software libraries for various encryption algorithms,

managing network connections, and displaying graphics. Libraries allow software developers to work on whatever makes their program unique, while leveraging known-good code for the tasks that similar programs routinely do. The programmer simply needs to understand the API for the libraries she intends to use. This reduces the amount of new code that the programmer needs to develop, which in turn makes the code easier to secure and maintain.

Using software libraries has potential risks, and these risks must be mitigated as part of secure software development practices. The main risk is that, because the libraries are reused across multiple projects in multiple organizations, any defect in these libraries propagates through every program that uses them. In fact, according to Veracode's 2020 report "State of Software Security: Open Source Edition," seven in ten applications use at least one open-source library with a security flaw, which makes those applications vulnerable. Keep in mind that these are open-source libraries, which (as we will discuss later in this chapter) are subject to examination by any number of security researchers looking for bugs. If you use proprietary libraries (including your own), it may be much harder to find these vulnerabilities before the threat actors do.

Secure Software Development

So far in this chapter (and the previous one), we've discussed software development in general terms, pointing out potential security pitfalls along the way. We now turn our attention to how we can bake security into our software from the ground up. To do so, however, we have to come from the top down, meaning we need an organizational policy document that clearly identifies the strategic goals, responsibilities, and authorities for mitigating risks associated with building or acquiring software. If the executive leadership doesn't push this, it just won't happen, and the policy document puts everyone on notice that secure coding is an organizational priority.

Secure coding is a set of practices that reduces (to acceptable levels) the risk of vulnerabilities in our software. No software will ever be 100 percent secure, but we can sure make it hard for threat actors to find and exploit any remaining vulnerabilities if we apply secure coding guidelines and standards to our projects.

Source Code Vulnerabilities

A *source code vulnerability* is a defect in code that provides a threat actor an opportunity to compromise the security of a software system. All code has defects (or bugs), but a vulnerability is a particularly dangerous one. Source code vulnerabilities are typically caused by two types of flaws: design and implementation. A *design flaw* is one that, even if the programmer did everything perfectly right, would still cause the vulnerability. An *implementation flaw* is one that stems from a programmer who incorrectly implemented a part of a good design. For example, suppose you are building an e-commerce application that collects payment card information from your customers and stores it for their future purchases. If you design the system to store the card numbers unencrypted, that would be a design flaw. If, on the other hand, you design the system to encrypt the data as soon as it

is captured but a programmer incorrectly calls the encryption function, resulting in the card number being stored in plaintext, that would be an implementation vulnerability.

Source code vulnerabilities are particularly problematic when they exist in externally facing systems such as web applications. These accounted for 39 percent of the external attacks carried out, according to Forrester's report "The State of Application Security, 2021." Web applications deserve particular attention because of their exposure.

The *Open Web Application Security Project (OWASP)* is an organization that deals specifically with web security issues. OWASP offers numerous tools, articles, and resources that developers can utilize to create secure software, and it also has individual member meetings (chapters) throughout the world. OWASP provides development guidelines, testing procedures, and code review steps, but is probably best known for its OWASP Top 10 list of web application security risks. The following is the most recent Top 10 list as of this writing, from 2017 (the 2021 version should be published by the time you're reading this):

- Injection
- Broken Authentication
- Sensitive Data Exposure
- XML External Entities (XEE)
- Broken Access Control
- Security Misconfiguration
- Cross-Site Scripting (XSS)
- Insecure Deserialization
- Using Components with Known Vulnerabilities
- Insufficient Logging & Monitoring

This list represents the most common vulnerabilities that reside in web-based software and are exploited most often. You can find out more information pertaining to these vulnerabilities at <https://owasp.org/www-project-top-ten/>.

Secure Coding Practices

So, we've talked about secure coding practices, but what exactly are they? Although the specific practices vary from one organization to the next, generally, they break down into two categories: standards and guidelines. Recall from Chapter 1 that *standards* are mandatory activities, actions, or rules, whereas *guidelines* are recommended actions and operational guides that provide the necessary flexibility for unforeseen circumstances. By enforcing secure coding standards and maintaining coding guidelines that reflect best practices, software development organizations dramatically reduce their source code vulnerabilities. Let's see how this works.

Coding Standards

Standards are the strongest form of secure coding practices because, to be considered a standard, a practice must meet the following requirements:

- Demonstrably reduce the risk of a particular type of vulnerability
- Be enforceable across the breadth of an organization's software development efforts
- Be verifiable in its implementation



EXAM TIP The rigorous application of secure coding standards is the best way to reduce source code vulnerabilities.

A very good reference for developing coding standards is the OWASP Top 10 list referenced in the previous section. Though it's focused on web applications, most of the vulnerabilities apply to any kind of software. Another good source of information is the organization's own past experience in developing code with vulnerabilities that later had to be patched.

Once the vulnerabilities are identified, even if at a fairly high level, coding standards can be developed to reduce the risk of building code that contains them. This is where things get a bit sticky, because the standards vary from one programming language to the next. If your organization develops web applications in Ruby (a common language for web apps), the way in which you reduce the risk of, say, broken authentication will be different than if you use PHP (another popular web app language). Still, there are plenty of opportunities to build standards that apply to all languages when we take a step back and consider the processes by which we develop, operationalize, and maintain that code. We'll cover this in more detail when we discuss security controls for software development later in this chapter.

Finally, a standard is only good if we can verify that we are complying with it. (Otherwise, why bother?) So, for instance, if we have a standard that reduces the risk of injection by validating inputs and parameters, then we should have a way to verify that none of our code fails to validate them. An excellent way to verify compliance with secure coding standards is the practice of code reviews, as discussed in Chapter 18. Ideally, though, we can verify at least some of our standards automatically.

Coding standards enable secure coding by ensuring programmers always do certain things and never do others. For example, a standard could require use of a particular library for encryption functions because it's been analyzed and determined to be sound and free from vulnerabilities. Another example of a standard could be forbidding programmers from using specific unsafe functions, such as the notorious `strcpy()` function in the C programming language. This function copies a string from one memory location to another, but doesn't check the length of the string being copied compared to the destination. If the string is longer than the destination, it will overwrite other areas of memory, which can result in a buffer overflow condition.

Software-Defined Security

A promising new area of security builds on the idea of software-defined networking (SDN), which we covered in Chapter 13. Recall that, in SDN, the control plane (i.e., the routing and switching decisions) is separate from the data plane (i.e., the packets and frames moving around). This allows centralized control of the network, which in turn improves performance, flexibility, and security. SDN also enables the separation of security functions from more traditional network appliance approaches. *Software-defined security (SDS or SDSec)* is a security model in which security functions such as firewalling, intrusion detection and prevention (IDS/IPS), and network segmentation are implemented in software within an SDN environment. One of the advantages of this approach is that sensors (for functions like IDS/IPS) can be dynamically repositioned depending on the threat environment.

SDS is a new technology but promises significant security advantages. Because of its dependence on SDN, SDS is best used in cloud and virtualized network environments.



NOTE Coding standards are required in certain regulated sectors such as automobile and railroad control software, among others.

Coding Guidelines

Secure coding guidelines are recommended practices that tend to be less specific than standards. For example, coding guidelines might encourage programmers to use variable names that are self-explanatory and not reused anywhere else in the program because this makes the code easier to understand. Applied to secure coding, these standards can help by ensuring code is consistently formatted and commented, which makes the code easier to read during code reviews. Guidelines may also recommend that coders keep functions short (without specifying how short) because this reduces the chance of errors. These practices may not sound like much, but they make it easier to spot errors early in the development process, thus improving quality, while decreasing vulnerabilities and costs.

Security Controls for Software Development

We tend to think of security controls as something to be added to an environment in order to reduce risks to it. While this is certainly true of software development environments, secure coding adds another layer, which consists of the security controls we build into the code itself. Regardless of whether we are protecting the development subnetwork or the software that is produced therein, we should implement security controls only after conducting deliberate threat modeling tied to a risk analysis process.

Keep in mind, however, that the threat models for an internal subnet are different from the threat models for software you're deploying throughout your organization or even selling to your customers. Either way, the goals are to reduce vulnerabilities and the possibility of system compromise, but the manner in which we do so will be very different.

Let's zoom in on just software you're developing. Which specific software controls you should use depends on the software itself, its objectives, the security goals of its associated security policy, the type of data it will process, the functionality it is to carry out, and the environment in which it will be placed. If an application is purely proprietary and will run only in closed, trusted environments, it may need fewer security controls than those required for applications that will connect businesses over the Internet and provide financial transactions. The trick is to understand the security needs of a piece of software, implement the right controls and mechanisms, thoroughly test the mechanisms and how they integrate into the application, follow structured development methodologies, and provide secure and reliable distribution methods.

In the sections that follow, we'll identify and describe the application of security controls for the major aspects of software development. These include aspects of the software itself, of course, but also the tools used to develop it, the manner in which we test it, and even how to integrate the software development environment into the broader security architecture.

Development Platforms

Software is normally developed by a team of software engineers who may or may not use the same tools. The most important tool in their tool set is an *integrated development environment (IDE)*, which enables each engineer to pull code from a repository (more on that later), edit it, test it, and then push it into the repository so the rest of the team can build on it. Depending on the programming language, target environments, and a host of other considerations, your developers may use Eclipse, Microsoft Visual Studio, Xcode, or various other applications. The software they develop will likely be tested (formally or otherwise) using development clients and servers that are supposed to represent the production platforms on which the finished software product will run. When we talk about security of the development platforms, therefore, we mean both the development endpoints and the “fake” clients and servers on which the software gets tested.

It may seem obvious, but the first step in ensuring the security of development platforms is to secure the devices on which our software engineers practice their craft. The challenge that many organizations face is that their engineers tend to be more sophisticated than the average user and will make changes to their computers that may or may not be authorized. Their principal incentive, after all, is to develop code quickly and correctly. If the configuration of their workstation gets in the way, it may find itself being modified. To avoid this, you should resist the temptation of giving your software engineers unfettered privileged access to their own devices. Enforcing good change management practices is critical to securing these development endpoints.

Even harder than ensuring change controls on your developers' workstations is securely provisioning the development clients and servers that they will need for testing.

Many organizations allow their developers to stand up and maintain their own development environment, which may be fine provided that these devices are isolated from the production environments. It may sound like common sense, but the problem is that some organizations don't do a good enough job of isolating development and production systems. In principle, doing so simply requires putting the development nodes in an isolated VLAN. In practice, the demarcation is not that cut and dry. This gets even more challenging when the team is distributed, which requires your developers (or perhaps their external collaborators) to remotely access the development hosts.

The best solution is to require use of a VPN to connect to the isolated development network. This may create a bit of work for the operations staff but is the only way to ensure that development and production code remains separate. Another good approach is to create firewall rules that prevent any unauthorized external connections (and even then only the bare minimum) to or from development servers. It should be clear by now that the provisioning of hosts on the development network should not be left to the software development team.

Tool Sets

As the old saying goes, you can't make everyone happy. Your IDE may be awesome, but invariably your software developers will need (or just want) additional tool sets. This is particularly true for developers that have a favorite tool that they've grown used to over the years, or if there is new work to be done for which the existing tools are not ideal. There are two approaches we've seen adopted by many organizations, and neither is ultimately good. The first is to force strict compliance with the approved tool sets that the organization provides. On the surface, this makes sense from a security and operations perspective. Having fewer tools means more standardization, allows for more thorough security assessments, and streamlines provisioning. However, it can also lead to a loss in productivity and typically leads the best coders to give up and move on to another organization where they're allowed more freedom.

The other (not good) approach is to let the developers run amuck in their own playground. The thinking goes something like this: we let them use whatever tools they feel are good, we set up and maintain whatever infrastructure they need, and we just fence the whole thing off from the outside so nothing bad can get in. The end of that sentence should make you shake your head in disagreement because keeping all the bad stuff out obviously is not possible, as you've learned throughout this book. Still, this is the approach of many small and mid-sized development shops.

A better approach is to treat the software development department the same way we treat any other. If they need a new tool, they simply put in a request that goes through the change management process, discussed in Chapter 20. The change advisory board (CAB) validates the requirement, assesses the risk, reviews the implementation plan, and so on. Assuming everything checks out and the CAB approves, the IT operations team integrates the tool into the inventory, updating and provisioning processes; the security team implements and monitors the appropriate controls, and the developers get the new tool they need.

Application Security Testing

Despite our best efforts, we (and all our programmers) are human and will make mistakes. Some of those mistakes will end up being source code vulnerabilities. Wouldn't it be nice to find them before our adversaries do? That's the role of application security testing, which comes in three flavors that you should know for the CISSP exam: static analysis, dynamic analysis, and fuzzing.

Static Application Security Testing

Static application security testing (SAST), also called *static analysis*, is a technique meant to help identify software defects or security policy violations and is carried out by examining the code without executing the program, and therefore is carried out before the program is compiled. The term SAST is generally reserved for automated tools that assist analysts and developers, whereas manual inspection by humans is generally referred to as *code review* (covered in Chapter 18).

SAST allows developers to quickly scavenge their source code for programming flaws and vulnerabilities. Additionally, this testing provides a scalable method of security code review and ensures that developers are following secure coding policies. There are numerous manifestations of SAST tools, ranging from tools that simply consider the behavior of single statements to tools that analyze the entire source code at once. However, you must remember that static code analysis can never reveal logical errors and design flaws, and therefore must be used in conjunction with manual code review to ensure thorough evaluation.

Dynamic Application Security Testing

Dynamic application security testing (DAST), also known as *dynamic analysis*, refers to the evaluation of a program in real time, while it is running. DAST is commonly carried out once a program has cleared the SAST stage and basic programming flaws have been rectified offline. DAST enables developers to trace subtle logical errors in the software that are likely to cause security mayhem later on. The primary advantage of this technique is that it eliminates the need to create artificial error-inducing scenarios. Dynamic analysis is also effective for compatibility testing, detecting memory leakages, identifying dependencies, and analyzing software without having to access the software's actual source code.



EXAM TIP Remember that SAST requires access to the source code, which is not executed during the tests, while DAST requires that you actually run the code but does not require access to the source code.

Fuzzing

Fuzzing is a technique used to discover flaws and vulnerabilities in software by sending large amounts of malformed, unexpected, or random data to the target program in order to trigger failures. Attackers can then manipulate these errors and flaws to inject their own code into the system and compromise its security and stability. Fuzzing tools, aka fuzzers, use complex inputs to attempt to impair program execution. Fuzzing tools

Manual Penetration Testing

Application security testing tools, together with good old-fashioned code reviews, are very good at unearthing most of the vulnerabilities that would otherwise go unnoticed by the software development team. As good as these tools are, however, they lack the creativity and resourcefulness of a determined threat actor. For this reason, many organizations also rely on *manual penetration testing (MPT)* as the final check before code is released into production environments. In this approach, an experienced red team examines the software system in its intended environment and looks for ways to compromise it. It is very common for this testing to uncover additional vulnerabilities that cannot be detected by automated tools.

are commonly successful at identifying buffer overflows, DoS vulnerabilities, injection weaknesses, validation flaws, and other activities that can cause software to freeze, crash, or throw unexpected errors.

Continuous Integration and Delivery

With the advent of Agile methodologies, discussed in Chapter 24, it has become possible to dramatically accelerate the time it takes to develop and release code. This has been taken to an extreme by many of the best software development organizations through processes of continuous integration and continuous delivery.

Continuous integration (CI) means that all new code is integrated into the rest of the system as soon as the developer writes it. For example, suppose Diana is a software engineer working on the user interface of a network detection and response (NDR) system. In traditional development approaches, she would spend a couple of weeks working on UI features, pretty much in isolation from the rest of the development team. There would then be a period of integration in which her code (and that of everyone else who's ready to deliver) gets integrated and tested. Then, Diana (and everyone else) goes back to working alone on her next set of features. The problem with this approach is that Diana gets to find out whether her code integrates properly only every two weeks. Wouldn't it be nice if she could find out instantly (or at least daily) whether any of her work has integration issues?

With continuous integration, Diana works on her code for a few hours and then merges it into a shared repository. This merge triggers a batch of unit tests. If her code fails those tests, the merge is rejected. Otherwise, her code is merged with everyone else's in the repository and a new version of the entire software system is built. If there are any errors in the build, she knows her code was the cause, and she can get to work fixing them right away. If the build goes well, it is immediately subjected to automated integration tests. If anything goes wrong, Diana knows she has to immediately get back to work fixing her code because she "broke the build," meaning nobody else can commit code until she fixes it or reverses her code merge.

Continuous integration dramatically improves software development efficiency by identifying errors early and often. CI also allows the practice of *continuous delivery* (CD), which is incrementally building a software product that can be released at any time. Because all processes and tests are automated, you could choose to release code to production daily or even hourly. Most organizations that practice CI/CD, however, don't release code that frequently. But they could if they wanted to.

CI/CD sounds wonderful, so what are the security risks we need to mitigate? Because CI/CD relies heavily on automation, most organizations that practice it use commercial or open-source testing platforms. One of those platforms is Codecov, which was compromised in early 2021, allowing the threat actor to modify its bash uploader script. This is the script that would take Diana's code in our earlier example and upload it for testing and integration. As an aside, because the tests are automated and don't involve actual users, developers typically have to provide access credentials, tokens, or keys to enable testing. The threat actor behind the Codecov breach modified the bash uploader so that it would exfiltrate this access data, potentially providing covert access to any of the millions of products worldwide that use Codecov for CI/CD.

The Codecov breach was detected about three months later by an alert customer who noticed unusual behavior in the uploader, investigated it, and alerted the vendor to the problem. Would you be able to tell that one of the components in your CI/CD toolset was leaking sensitive data? You could if you practice the secure design principles we've been highlighting throughout the book, especially threat modeling, least privilege, defense in depth, and zero trust.

Security Orchestration, Automation, and Response

The Codecov breach mentioned in the previous section also highlights the role that a security orchestration, automation, and response (SOAR) platform can play in securing your software development practices. Chapter 21 introduced SOAR in the context of the role of a security information and event management (SIEM) platform in your security operations. Both SOAR and SIEM platforms can help detect and, in the case of SOAR, respond to threats against your software development efforts. If you have sensors in your development subnet (you did segment your network, right?) and a well-tuned SOAR platform, you can detect new traffic flowing from that subnet (which shouldn't be talking much to the outside world) to a new external endpoint. If the traffic is unencrypted (or you use a TLS decryption proxy to do deep packet inspection), you'd notice access tokens and keys flowing out to a new destination. Based on this observation, you could declare an incident and activate the playbook for data breaches in your SOAR platform. Just like that, you would've stopped the bleeding, buying you time to figure out what went wrong and how to fix it for the long term.

One of the challenges with the scenario just described is that many security teams treat their organization's development environment as a bit of a necessary chaos that must be tolerated. Software developers are typically rewarded (or punished) according to their ability to produce quality code quickly. They can be resistant (or even rebel against) anything that gets in the way of their efficiency, and, as we well know, security tends to do just that. This is where DevSecOps (discussed in Chapter 24) can help build

the right culture and balance the needs of all teammates. It can also help the security team identify and implement controls that mitigate risks such as data breaches, while minimally affecting productivity. One such control is the placement of sensors such as IDS/IPS, NDR, and data loss prevention (DLP) within the development subnets. These systems, in turn, would report to the SOAR platform, which could detect and contain active threats against the organization.

Software Configuration Management

Not every threat, of course, is external. There are plenty of things our own teammates can do deliberately or otherwise that cause problems for the organization. As we'll see later in this chapter when we discuss cloud services, improper configurations consistently rank among the worst threats to many organizations. This threat, however, is a solved problem in organizations that practice proper configuration management, as we covered in Chapter 20.

Anticipating the inevitable changes that will take place to a software product during its development life cycle, a configuration management system should be put into place that allows for change control processes to take place through automation. Since deploying an insecure configuration to an otherwise secure software product makes the whole thing insecure, these settings are a critical component of securing the software development environment. A product that provides *software configuration management (SCM)* identifies the attributes of software at various points in time and performs a methodical control of changes for the purpose of maintaining software integrity and traceability throughout the software development life cycle. It tracks changes to configurations and provides the ability to verify that the final delivered software has all of the approved changes that are supposed to be included in the release.

During a software development project, the centralized code repositories are often kept in systems that can carry out SCM functionality. These SCM systems manage and track revisions made by multiple people against a single master set and provide concurrency management, versioning, and synchronization. *Concurrency management* deals with the issues that arise when multiple people extract the same file from a central repository and make their own individual changes. If they were permitted to submit their updated files in an uncontrolled manner, the files would just write over each other and changes would be lost. Many SCM systems use algorithms to version, fork, and merge the changes as files are checked back into the repository.

Versioning deals with keeping track of file revisions, which makes it possible to “roll back” to a previous version of the file. An archive copy of every file can be made when it is checked into the repository, or every change made to a file can be saved to a transaction log. Versioning systems should also create log reports of who made changes, when they were made, and what the changes were.

Some SCM systems allow individuals to check out complete or partial copies of the repositories and work on the files as needed. They can then commit their changes back to the master repository as needed and update their own personal copies to stay up to date with changes other people have made. This process is called *synchronization*.

Code Repositories

A *code repository*, which is typically a version control system, is the vault containing the crown jewels of any organization involved in software development. If we put on our adversarial hats for a few minutes, we could come up with all kinds of nefarious scenarios involving these repositories. Perhaps the simplest is that someone could steal our source code, which embodies not only many staff hours of work but, more significantly, our intellectual property. An adversary could also use our source code to look for vulnerabilities to exploit later, once the code is in production. Finally, adversaries could deliberately insert vulnerabilities into our software, perhaps after it has undergone all testing and is trusted, so that they can exploit it later at a time of their choosing. Clearly, securing our source code repositories is critical.

Perhaps the most secure way of managing security for your code repositories is to implement them on an isolated (or “air-gapped”) network that includes the development, test, and QA environments. The development team would have to be on this network to do their work, and the code, once verified, could be exported to the production servers using removable storage media. We already presented this best

Software Escrow

If a company pays another company to develop software for it, it should have some type of *software escrow* in place for protection. We covered this topic in Chapter 23 from a business continuity perspective, but since it directly deals with software development, we will mention it here also.

In a software escrow framework, a third party keeps a copy of the source code, and possibly other materials, which it will release to the customer only if specific circumstances arise, mainly if the vendor who developed the code goes out of business or for some reason is not meeting its obligations and responsibilities. This procedure protects the customer, because the customer pays the vendor to develop software code for it, and if the vendor goes out of business, the customer otherwise would no longer have access to the actual code. This means the customer code could never be updated or maintained properly.

A logical question would be, “Why doesn’t the vendor just hand over the source code to the customer, since the customer paid for it to be developed in the first place?” It does not always work that way. The code may be the vendor’s intellectual property. The vendor employs and pays people with the necessary skills to develop that code, and if the vendor were to just hand it over to the customer, it could be giving away its intellectual property, its secrets. The customer oftentimes gets compiled code instead of source code. *Compiled code* is code that has been put through a compiler and is unreadable to humans. Most software profits are based on licensing, which outlines what customers can do with the compiled code. For an added fee, of course, most custom software developers will also provide the source, which could be useful in sensitive applications.

practice in the preceding section. The challenge with this approach is that it severely limits the manner in which the development team can connect to the code. It also makes it difficult to collaborate with external parties and for developers to work from remote or mobile locations.

A pretty good alternative would be to host the repository on the intranet, which would require developers to either be on the local network or connect to it using a VPN connection. As an added layer of security, the repositories can be configured to require the use of Secure Shell (SSH), which would ensure all traffic is encrypted, even inside the intranet, to mitigate the risk of sniffing. Finally, SSH can be configured to use public key infrastructure (PKI), which allows us to implement not only confidentiality and integrity but also nonrepudiation. If you have to allow remote access to your repository, this would be a good way to go about it.

Finally, if you are operating on a limited budget or have limited security expertise in this area, you can choose one of the many web-based repository service providers and let them take care of the security for you. While this may mitigate the basic risks for small organizations, it is probably not an acceptable course of action for projects with significant investments of intellectual property.

Software Security Assessments

We already discussed the various types of security assessments in Chapter 18, but let's circle back here and see how these apply specifically to software security. Recall from previous sections in this chapter that secure software development practices originate in an organizational policy that is grounded in risk management. That policy is implemented through secure coding standards, guidelines, and procedures that should result in secure software products. We verify this is so through the various testing methods discussed in this chapter (e.g., SAST and DAST) and Chapter 24 (e.g., unit, integration, etc.). The purpose of a software security assessment, then, is to verify that this entire chain, from policy to product, is working as it should.

When conducting an assessment, it is imperative that the team review all applicable documents and develop a plan for how to verify each requirement from the applicable policies and standards. Two areas that merit additional attention are the manner in which the organization manages risks associated with software development and how it audits and logs software changes.

Risk Analysis and Mitigation

Risk management is at the heart of secure software development, particularly the mapping between risks we've identified and the controls we implement to mitigate them. This is probably one of the trickiest challenges in secure software development in general, and in auditing it in particular. When organizations do map risks to controls in software development, they tend to do so in a generic way. For example, the OWASP Top 10 list is a great starting point for analyzing and mitigating vulnerabilities, but how are we doing against specific (and potentially unique) threats faced by our organization?

Threat modeling is an important activity for any development team, and particularly in DevSecOps. Sadly, however, most organizations don't conduct threat modeling for

their software development projects. If they're defending against generic threats, that's good, but sooner or later we all face unique threats that, if we haven't analyzed and mitigated them, have a high probability of ruining our weekend.

Another area of interest for assessors are the linkages between the software development and risk management programs. If software projects are not tracked in the organization's risk matrix, then the development team will probably be working in isolation, disconnected from the broader risk management efforts.

Change Management

Another area in which integration with broader organizational efforts is critical to secure software development is change management. Changes to a software project that may appear inconsequential when considered in isolation could actually pose threats when analyzed within the broader context of the organization. If software development is not integrated into the organization's change management program, auditing changes to software products may be difficult, even if the changes are being logged by the development team. Be that as it may, software changes should not be siloed from overall organizational change management because doing so will likely lead to interoperability or (worse yet) security problems.

Assessing the Security of Acquired Software

Most organizations do not have the in-house capability to develop their own software systems. Their only feasible options are either to acquire standard software or to have a vendor build or customize a software system to their particular environment. In either case, software from an external source will be allowed to execute in a trusted environment. Depending on how trustworthy the source and the code are, this could have some profound implications to the security posture of the organization's systems. As always, we need to ground our response on our risk management process.

In terms of managing the risk associated with acquired software, the essential question to ask is, "How is the organization affected if this software behaves improperly?" Improper behavior could be the consequence of either defects or misconfiguration. The defects can manifest themselves as computing errors (e.g., wrong results) or vulnerability to intentional attack. A related question is, "What is it that we are protecting and this software could compromise?" Is it personally identifiable information (PII), intellectual property, or national security information? The answers to these and other questions will dictate the required thoroughness of our approach.

In many cases, our approach to mitigating the risks of acquired software will begin with an assessment of the software developer. Characteristics that correlate to a lower software risk include the good reputation of the developer and the regularity of its patch pushes. Conversely, developers may be riskier if they have a bad reputation, are small or new organizations, if they have immature or undocumented development processes, or if their products have broad marketplace presence (meaning they are more lucrative targets to exploit developers).

A key element in assessing the security of acquired software is, rather obviously, its performance in an internal assessment. Ideally, we are able to obtain the source code

from the vendor so that we can do our own code reviews, vulnerability assessments, and penetration tests. In many cases, however, this will not be possible. Our only possible assessment may be a penetration test. The catch is that we may not have the in-house capability to perform such a test. In such cases, and depending on the potential risk posed by this software, we may be well advised to hire an external party to perform an independent penetration test for us. This is likely a costly affair that would only be justifiable in cases where a successful attack against the software system would likely lead to significant losses for the organization.

Even in the most constrained case, we are still able to mitigate the risk of acquisition. If we don't have the means to do code reviews, vulnerability assessments, or penetration tests, we can still mitigate the risk by deploying the software only in specific subnetworks, with hardened configurations, and with restrictive IDS/IPS rules monitoring its behavior. Though this approach may initially lead to constrained functionality and excessive false positives generated by our IDS/IPS, we can always gradually loosen the controls as we gain assurances that the software is trustworthy.

Commercial Software

It is exceptionally rare for an organization to gain access to the source code of a commercial-off-the-shelf (COTS) product to conduct a security assessment of it. However, depending on the product, we may not have to. The most widely used commercial software products have been around for years and have had their share of security researchers (both benign and malicious) poking at them the whole time. We can simply research what vulnerabilities and exploits have been discovered by others and decide for ourselves whether or not the vendor uses effective secure coding practices.

If the software is not as popular, or serves a small niche community, the risk of undiscovered vulnerabilities is probably higher. In these cases, it pays to look into the certifications of the vendor. A good certification for a software developer is ISO/IEC 27034 Application Security. Unfortunately, you won't find a lot of vendors certified in it. There are also certifications that are very specific to a sector (e.g., ISO 26262 for automotive safety) or a programming language (e.g., ISO/IEC TS 17961:2013 for coding in C) and are a bit less rare to find. Ultimately, however, the security of a vendor's software products is tied to how seriously it takes security in the first place. Absent a secure coding certification, you can look for overall information security management system (ISMS) certifications like ISO/IEC 27001 and FedRAMP, which are difficult to obtain and show that security is taken seriously in an organization.

Open-Source Software

Open-source software is released with a license agreement that allows the user to examine its source code, modify it at will, and even redistribute the modified software (which, per the license, usually requires acknowledgment of the original source and a description of modifications). This may seem perfect, but there are some caveats to keep in mind. First, the software is released as-is, typically without any service or support agreements (though these can be purchased through third parties). This means that your staff may have to

figure out how to install, configure, and maintain the software on their own, unless you contract with someone else to do this for you.

Second, part of the allure of open-source software is that we get access to the source code. This means we can apply all the security tests and assessments we covered earlier. Of course, this only helps if we have the in-house capabilities to examine the source code effectively. Even if we don't, however, we can rely on countless developers and researchers around the world who do examine it (at least for the more popular software). The flip side of that coin, however, is that the adversaries also get to examine the code to either identify vulnerabilities quicker than the defenders or gain insights into how they might more effectively attack organizations that use specific software.

Perhaps the greatest risk in using open-source software is relying on outdated versions of it. Many of us are used to having software that automatically checks for updates and applies them automatically (either with or without our explicit permission). This is not all that common in open-source software, however, especially libraries. This means we need to develop processes to ensure that all open-source software is periodically updated, possibly in a way that differs from the way in which COTS software is updated.

Third-Party Software

Third-party software, also known as *outsourced software*, is software made specifically for an organization by a third party. Since the software is custom (or at least customized), it is not considered COTS. Third-party software may rely partly (or even completely) on open-source software, but, having been customized, it may introduce new vulnerabilities. So, we need a way to verify the security of these products that is probably different from how we would do so with COTS or open-source software.



EXAM TIP Third-party software is custom (or at least customized) to an organization and is not considered commercial off-the-shelf (COTS).

The best (and, sadly, most expensive) way to assess the security of third-party software is to leverage the external or third-party audits discussed in Chapter 18. The way this typically works is that we write into the contract a provision for an external auditor to inspect the software (and possibly the practices through which it was developed), and then issue a report, attesting to the security of the product. Passing this audit can be a condition of finalizing the purchase. Obviously, a sticking point in this negotiation can be who pays for this audit.

Another assessment approach is to arrange for a time-limited trial of the third-party software (perhaps at a nominal cost to the organization), and then have a red team perform an assessment. If you don't have a red team, you can probably hire one for less money than a formal application security audit would cost. Still, the cost will be considerable, typically (at least) in the low tens of thousands of dollars. As with any other security control, you'd have to balance the cost of the assessment and the loss you would incur from insecure software.

Managed Services

As our organizations continue to migrate to cloud services (IaaS, PaaS, and SaaS, discussed in depth in Chapter 7), we should also assess the security impact of those services. This is highlighted by a 2020 study by global intelligence firm IDC, which found that nearly 80 percent of the companies surveyed had experienced at least one cloud data breach in the past 18 months. The top three reasons were misconfigurations, lack of visibility into access settings and activities, and improper access control. The major cloud services provide tools to help you avoid these pitfalls, but the bottom line is that, if you don't have the in-house expertise to secure and assess your cloud services, you really should consider contracting an expert to help you out.

Chapter Review

Building secure code requires commitment from many parts of the organization, not just the development and security teams. It starts at the very top with a policy document that is implemented through standards, procedures, and guidelines. A key part of these is the inclusion of the various types of tests that must be run regularly (even continuously) on the software as it is being written, integrated, and prepared for delivery. Software development environments are complex and could require different approaches from those you'd take in a normal network environment. For this reason, teamwork among all stakeholders is absolutely critical. A really good way to facilitate this collaboration is by using the DevSecOps approach introduced in Chapter 24 and highlighted in this one.

Even if your organization doesn't develop software, it most certainly uses applications and services developed by others. That's why the concepts discussed in this chapter are universally applicable to any cybersecurity leader. You must understand how secure code is built, so that you can determine whether the software you're getting from others presents any undue risks to your organization's cybersecurity.

Quick Review

- Machine language, which consists of 1's and 0's, is the only format that a computer's processor can understand directly and is considered a first-generation language.
- Assembly language is considered a second-generation programming language and uses symbols (called mnemonics) to represent complicated binary codes.
- Third-generation programming languages, such as C/C++, Java, and Python, are known as high-level languages due to their refined programming structures, which allow programmers to leave low-level (system architecture) intricacies to the programming language and focus on their programming objectives.
- Fourth-generation languages (aka very high-level languages) use natural language processing to allow inexperienced programmers to develop code in less time than it would take an experienced software engineer to do so using a third-generation language.

- Fifth-generation programming languages (aka natural languages) approach programming by defining the constraints for achieving a specified result and allowing the development environment to solve problems by itself instead of a programmer having to develop code to deal with individual and specific problems.
- Assemblers are tools that convert assembly language source code into machine code.
- Compilers transform instructions from a source language (high-level) to a target language (machine), sometimes using an external assembler along the way.
- A garbage collector identifies blocks of memory that were once allocated but are no longer in use and deallocates the blocks and marks them as free.
- A runtime environment (RTE) functions as a miniature operating system for the program and provides all the resources portable code needs.
- In object-oriented programming (OOP), related functions and data are encapsulated together in classes, which may then be instantiated as objects.
- Objects in OOP communicate with each other by using messages that conform to the receiving object's application programming interface (API) definition.
- Cohesion reflects how many different types of tasks a module can carry out, with the goal being to perform only one task (high cohesion), which makes modules easier to maintain.
- Coupling is a measure of how much a module depends on others; the more dependencies it has, the more complex and difficult the module is to maintain, so we want low (or loose) coupling.
- An API specifies the manner in which a software component interacts with other software components.
- Parameter validation refers to confirming that the parameter values being received by an application are within defined limits before they are processed by the system.
- A software library is a collection of components that do specific tasks that are useful to many other components.
- Secure coding is a set of practices that reduce (to acceptable levels) the risk of vulnerabilities in our software.
- A source code vulnerability is a defect in code that provides a threat actor an opportunity to compromise the security of a software system.
- Secure coding standards are verifiable, mandatory practices that reduce the risk of particular types of vulnerabilities in the source code.
- Secure coding guidelines are recommended practices that tend to be less specific than standards.

- Software-defined security (SDS or SDSec) is a security model in which security functions such as firewalling, IDS/IPS, and network segmentation are implemented in software within an SDN environment.
- Software development tools should be authorized, implemented, and maintained just like any other software product through the organization's change management process; developers should not be allowed to install and use arbitrary tools.
- Static application security testing (SAST) is a technique meant to help identify software defects or security policy violations and is carried out by examining the source code without executing the program.
- Dynamic application security testing (DAST) refers to the evaluation of a program in real time, while it is running.
- Fuzzing is a technique used to discover flaws and vulnerabilities in software by sending large amounts of malformed, unexpected, or random data to the target program in order to trigger failures.
- Continuous integration means that all new code is integrated into the rest of the system as soon as the developer writes it.
- Continuous delivery is incrementally building a software product that can be released at any time and requires continuous integration.
- A software configuration management (SCM) platform identifies the attributes of software at various points in time and performs a methodical control of changes for the purpose of maintaining software integrity and traceability throughout the SDLC.
- The purpose of a software security assessment is to verify that this entire development process, from organizational policy to delivered product, is working as it should.
- Security assessments of acquired software are essential to mitigate the risk they could pose to the organization that acquired it.
- The most practical way to assess the security of commercial software is to research what vulnerabilities and exploits have been discovered by others and decide for ourselves whether or not the vendor uses effective secure coding practices.
- The greatest risk in using open-source software is relying on outdated versions of it.
- The best way to assess the security of third-party (i.e., custom or customized) software is to perform external or third-party audits.

Questions

Please remember that these questions are formatted and asked in a certain way for a reason. Keep in mind that the CISSP exam is asking questions at a conceptual level. Questions may not always have the perfect answer, and the candidate is advised against

always looking for the perfect answer. Instead, the candidate should look for the best answer in the list.

1. What language is the only one that a computer processor can natively understand and execute?
 - A. Machine language
 - B. Register language
 - C. Assembly language
 - D. High-level language
2. To which generation do programming languages such as C/C++, Java, and Python belong?
 - A. Second generation
 - B. Third generation
 - C. Fourth generation
 - D. Fifth generation
3. Which type of tool is specifically designed to convert assembly language into machine language?
 - A. Compiler
 - B. Integrated development environment (IDE)
 - C. Assembler
 - D. Fuzzer
4. Which of the following is not very useful in assessing the security of acquired software?
 - A. The reliability and maturity of the vendor
 - B. The vendor's software escrow framework
 - C. Third-party vulnerability assessments
 - D. In-house code reviews if source code is available
5. Cohesion and coupling are characteristics of quality code. Which of the following describes the goals for these two characteristics?
 - A. Low cohesion, low coupling
 - B. Low cohesion, high coupling
 - C. High cohesion, low coupling
 - D. High cohesion, high coupling

6. Yichen is a new software engineer at Acme Software, Inc. During his first code review, he is told by his boss that he should use descriptive names for variables in his code. What is this observation an example of?
 - A. Secure coding guidelines
 - B. Secure coding standards
 - C. Secure software development policy
 - D. Use of fifth-generation language
7. On what other technology does software-defined security depend?
 - A. Software-defined storage (SDS)
 - B. Software-defined networking (SDN)
 - C. Security orchestration, automation, and response (SOAR)
 - D. Continuous integration (CI)
8. If you wanted to test source code for vulnerabilities without running it, which approach would be best?
 - A. Static application security testing (SAST)
 - B. Fuzzing
 - C. Dynamic application security testing (DAST)
 - D. Manual penetration testing
9. If you wanted to test software for vulnerabilities by executing it and then exposing it to large amounts of random inputs, which testing technique would you use?
 - A. Static application security testing (SAST)
 - B. Fuzzing
 - C. Dynamic application security testing (DAST)
 - D. Manual penetration testing
10. Which of the following is not a common reason for data breaches in managed cloud services?
 - A. Misconfigurations
 - B. Lack of visibility into access settings and activities
 - C. Hardware failures
 - D. Improper access control

Answers

1. **A.** Machine language, which consists of 1's and 0's, is the only format that a computer's processor can understand directly and is considered a first-generation language.

2. **B.** Third-generation programming languages, such as C/C++, Java, and Python, are known as high-level languages due to their refined programming structures, which allow programmers to leave low-level (system architecture) intricacies to the programming language and focus on their programming objectives.
3. **C.** Assemblers are tools that convert assembly language source code into machine code. Compilers also generate machine language, but do so by transforming high-level language code, not assembly language.
4. **B.** In a software escrow framework, a third party keeps a copy of the source code, and possibly other materials, which it will release to the customer in specific circumstances such as the developer going out of business. While software escrow is a good business continuity practice, it wouldn't normally tell us anything about the security of the software itself. All three other answers are part of a rigorous assessment of the security of acquired software.
5. **C.** Cohesion reflects how many different types of tasks a module can carry out, with the goal being to perform only one task (high cohesion), which makes modules easier to maintain. Coupling is a measure of how much a module depends on others; the more dependencies it has, the more complex and difficult the module is to maintain, so we want low (or loose) coupling.
6. **A.** Secure coding guidelines are recommended practices that tend to be less specific than standards. They might encourage programmers to use variable names that are self-explanatory and to keep functions short (without specifying how short). Secure coding standards, on the other hand, are verifiable, mandatory practices that reduce the risk of particular types of vulnerabilities in the source code.
7. **B.** Software-defined security (SDS or SDSec) is a security model in which security functions such as firewalling, IDS/IPS, and network segmentation are implemented in software within an SDN environment.
8. **A.** Static application security testing (SAST) is a technique meant to help identify software defects or security policy violations and is carried out by examining the source code without executing the program. All the other answers require that the code be executed.
9. **B.** Fuzzing is a technique used to discover flaws and vulnerabilities in software by sending large amounts of malformed, unexpected, or random data to the target program in order to trigger failures.
10. **C.** The top three reasons for data breaches in cloud services are misconfigurations, lack of visibility into access settings and activities, and improper access control.

This page intentionally left blank

Comprehensive Questions

Use the following scenario to answer Questions 1–3. Josh has discovered that an organized hacking ring in China has been targeting his company's research and development department. If these hackers have been able to uncover his company's research findings, this means they probably have access to his company's intellectual property. Josh thinks that an e-mail server in his company's DMZ may have been successfully compromised and a rootkit loaded.

1. Based upon this scenario, what is most likely the biggest risk Josh's company needs to be concerned with?
 - A. Market share drop if the attackers are able to bring the specific product to market more quickly than Josh's company.
 - B. Confidentiality of e-mail messages. Attackers may post all captured e-mail messages to the Internet.
 - C. Impact on reputation if the customer base finds out about the attack.
 - D. Depth of infiltration of attackers. If attackers have compromised other systems, more confidential data could be at risk.
2. The attackers in this situation would be seen as which of the following?
 - A. Vulnerability
 - B. Threat
 - C. Risk
 - D. Threat agent
3. If Josh is correct in his assumptions, which of the following best describes the vulnerability, threat, and exposure, respectively?
 - A. E-mail server is hardened, an entity could exploit programming code flaw, server is compromised and leaking data.
 - B. E-mail server is not patched, an entity could exploit a vulnerability, server is hardened.
 - C. E-mail server misconfiguration, an entity could exploit misconfiguration, server is compromised and leaking data.
 - D. DMZ firewall misconfiguration, an entity could exploit misconfiguration, internal e-mail server is compromised.

4. Aaron is a security manager who needs to develop a solution to allow his company's mobile devices to be authenticated in a standardized and centralized manner using digital certificates. The applications these mobile clients use require a TCP connection. Which of the following is the best solution for Aaron to implement?
 - A. TACACS+
 - B. RADIUS
 - C. Diameter
 - D. Mobile IP
5. Terry is a security manager for a credit card processing company. His company uses internal DNS servers, which are placed within the LAN, and external DNS servers, which are placed in the DMZ. The company also relies on DNS servers provided by its service provider. Terry has found out that attackers have been able to manipulate several DNS server caches to point employee traffic to malicious websites. Which of the following best describes the solution this company should implement?
 - A. IPSec
 - B. PKI
 - C. DNSSEC
 - D. MAC-based security
6. Which of the following is not a key provision of the GDPR?
 - A. Requirement for consent from data subjects
 - B. Right to be informed
 - C. Exclusion for temporary workers
 - D. Right to be forgotten
7. Jane is suspicious that an employee is sending sensitive data to one of the company's competitors but is unable to confirm this. The employee has to use this data for daily activities, thus it is difficult to properly restrict the employee's access rights. In this scenario, which best describes the company's vulnerability, threat, risk, and necessary control?
 - A. Vulnerability is employee access rights, threat is internal entities misusing privileged access, risk is the business impact of data loss, and the necessary control is detailed network traffic monitoring.
 - B. Vulnerability is lack of user monitoring, threat is internal entities misusing privileged access, risk is the business impact of data loss, and the necessary control is detailed user activity logs.
 - C. Vulnerability is employee access rights, threat is internal employees misusing privileged access, risk is the business impact of confidentiality, and the necessary control is multifactor authentication.
 - D. Vulnerability is employee access rights, threat is internal users misusing privileged access, risk is the business impact of confidentiality, and the necessary control is CCTV.

8. Which of the following best describes what role-based access control offers organizations in reducing administrative burdens?
 - A. It allows entities closer to the resources to make decisions about who can and cannot access resources.
 - B. It provides a centralized approach for access control, which frees up department managers.
 - C. User membership in roles can be easily revoked and new ones established as job assignments dictate.
 - D. It enforces an enterprise-wide security policy, standards, and guidelines.
9. Mark works for a large corporation operating in multiple countries worldwide. He is reviewing his company's policies and procedures dealing with data breaches. Which of the following is an issue that he must take into consideration?
 - A. Each country may or may not have unique notification requirements.
 - B. All breaches must be announced to affected parties within 24 hours.
 - C. Breach notification is a "best effort" process and not a guaranteed process.
 - D. Breach notifications are avoidable if all PII is removed from data stores.
10. A software development company released a product that committed several errors that were not expected once deployed in their customers' environments. All of the software code went through a long list of tests before being released. The team manager found out that after a small change was made to the code, the program was not tested before it was released. Which of the following tests was most likely not conducted?
 - A. Unit
 - B. Compiled
 - C. Integration
 - D. Regression
11. Which of the following should not be considered as part of the supply chain risk management process for a smartphone manufacturer?
 - A. Hardware Trojans inserted by downstream partners
 - B. ISO/IEC 27001
 - C. Hardware Trojans inserted by upstream partners
 - D. NIST Special Publication 800-161
12. Data sovereignty is increasingly becoming an issue that most of us in cybersecurity should address within our organizations. What does the term data sovereignty mean?
 - A. Certain types of data concerning a country's citizens must be stored and processed in that country.
 - B. Data on a country's citizens must be stored and processed according to that country's laws, regardless of where the storing/processing takes place.

- C. Certain types of data concerning a country's citizens are the sovereign property of that data subject.
- D. Data on a country's citizens must never cross the sovereign borders of another country.

Use the following scenario to answer Questions 13–15. Jack has just been hired as the security officer for a large hospital system. The organization develops some of its own proprietary applications. The organization does not have as many layers of controls when it comes to the data processed by these applications, since it is assumed that external entities will not understand the internal logic of the applications. One of the first things that Jack wants to carry out is a risk assessment to determine the organization's current risk profile. He also tells his boss that the hospital should become ISO certified to bolster its customers' and partners' confidence in its risk management processes.

- 13. Which of the following approaches has been implemented in this scenario?
 - A. Defense-in-depth
 - B. Security through obscurity
 - C. Information security management system
 - D. ISO/IEC 27001
- 14. Which ISO/IEC standard would be best for Jack to follow to meet his goals?
 - A. ISO/IEC 27001
 - B. ISO/IEC 27004
 - C. ISO/IEC 27005
 - D. ISO/IEC 27006
- 15. Which standard should Jack suggest to his boss for compliance with best practices regarding storing and processing sensitive medical information?
 - A. ISO/IEC 27004
 - B. ISO/IEC 27001
 - C. ISO/IEC 27799
 - D. ISO/IEC 27006
- 16. You just received an e-mail from one of your hardware manufacturers notifying you that it will no longer manufacture a certain product and, after the end of the year, you won't be able to send it in for repairs, buy spare parts, or get technical assistance from that manufacturer. What term describes this?
 - A. End-of-support (EOS)
 - B. End-of-service-life (EOSL)
 - C. Deprecation
 - D. End-of-life (EOL)

17. The confidentiality of sensitive data is protected in different ways depending on the state of the data. Which of the following is the best approach to protecting data in transit?
 - A. SSL
 - B. VPN
 - C. IEEE 802.1X
 - D. Whole-disk encryption
18. Your boss asks you to put together a report describing probable adverse effects on your assets caused by specific threat sources. What term describes this?
 - A. Risk analysis
 - B. Threat modeling
 - C. Attack trees
 - D. MITRE ATT&CK
19. A(n) _____ is the graphical representation of data commonly used on websites. It is a skewed representation of characteristics a person must enter to prove that the subject is a human and not an automated tool, as in a software robot.
 - A. anti-spoofing symbol
 - B. CAPTCHA
 - C. spam anti-spoofing symbol
 - D. CAPCHAT
20. Mark has been asked to interview individuals to fulfill a new position in his company, chief privacy officer (CPO). What is the function of this type of position?
 - A. Ensuring that company financial information is correct and secure
 - B. Ensuring that customer, company, and employee data is protected
 - C. Ensuring that security policies are defined and enforced
 - D. Ensuring that partner information is kept safe
21. A risk management program must be developed properly and in the right sequence. Which of the following provides the correct sequence for the steps listed?
 - i. Develop a risk management team.
 - ii. Calculate the value of each asset.
 - iii. Identify the vulnerabilities and threats that can affect the identified assets.
 - iv. Identify company assets to be assessed.
 - A. i, iii, ii, iv
 - B. ii, i, iv, iii
 - C. iii, i, iv, ii
 - D. i, iv, ii, iii