

# CIPHERS

---

Luke Anderson

11th March 2016

University Of Sydney



1. Crypto-Bulletin
2. Ciphers
  - 2.1 Cryptography Definition
  - 2.2 Symmetric Ciphers
  - 2.3 Cryptanalysis
  - 2.4 Attack Examples
  - 2.5 Failures
  - 2.6 Substitution Ciphers
  - 2.7 Permutation Ciphers
  - 2.8 Vigenere Cipher

# CRYPTO-BULLETIN

---

## 'Creepy' Aussie Anon sentenced for hacking websites

<http://www.itnews.com.au/news/creepy-aussie-anon-sentenced-for-hacking-websites-416485>

## Why you should care about Australia's new defence trade controls

<http://www.itnews.com.au/feature/>

[why-you-should-care-about-australias-new-defence-trade-controls-416446](http://www.itnews.com.au/feature/why-you-should-care-about-australias-new-defence-trade-controls-416446)

## UNSW team wins Cyber Security Challenge

<http://www.itnews.com.au/news/unsw-team-wins-cyber-security-challenge-415599>

# CIPHERS

---

**Cryptography** is the study of mathematical techniques related to the design of ciphers.

It's one of the core examples of the many mechanisms making up security:

- Cryptography
- Signature / Pattern Matching
- Access Control
- Statistical Profiling
- Traffic Security
- Countermeasures
- Tamper Resistance

# Cryptography

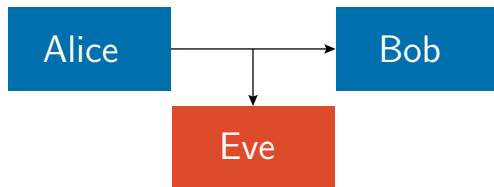
The fundamental application of cryptography is enabling **secure communications** over an **insecure channel**.

How can Alice send a secure message to Bob over an insecure channel when Eve is listening in?

Eve is an active attacker and may tap, insert or modify messages in transit.

How does one use cryptography to provide security such as:

- Authentication
- Confidentiality
- Integrity
- Non-Repudiation



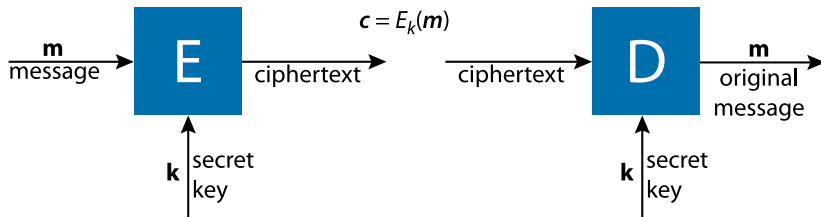
# Symmetric Ciphers

The traditional way of achieving this is through **private key encryption**.

This is also known as **symmetric encryption** since the key used to both encrypt and decrypt messages is the same.

A symmetric cipher is one defined by the rule:

$$D_k(E_k(m)) = m$$





# Communications with Symmetric Ciphers

Alice and Bob share:

- A secret key:  $k$
- An encryption algorithm:  $E_k$
- A decryption algorithm:  $D_k$

Alice wants to send Bob the message  $m$ .

The unencrypted message is known as either the **plaintext** or **cleartext**.

Alice encrypts  $m$  by computing the **ciphertext**:

$$c = E_k(m)$$

Bob decrypts  $c$  by computing the original plaintext message  $m$ :

$$D_k(c) = m$$

# Symmetric Cryptosystem

It is computationally hard to decrypt the ciphertext  $c$  without the secret key  $k$ .

The secret key  $k$  is usually a large number of bits ( $\geq 64\text{bits}$ )

The range of possible values of  $k$  is called the **key space**:  $K$   
(for 64 bit keys, the key space is  $(0 \dots 2^{64} - 1)$  or  $\{0, 1\}^{64}$ )

The range of possible messages is called the **message space**:  $M$

A **cryptosystem** is a system consisting of an algorithm, plus all possible plaintexts, ciphertexts, and keys.

# Types of Symmetric Ciphers

---

## **Stream Ciphers**

Operate on a single bit or byte at a time.

## **Block Ciphers**

Operate on blocks (numbers of bits) of plaintext at a time.

## Definition

**Cryptanalysis** is used to breach cryptographic security systems and gain access to the contents of encrypted messages, even if the cryptographic key is unknown.

We *always* assume that attackers have:

- Complete access to the communications channel
- Complete knowledge about the cryptosystem

**Secrecy must only depend on the key.**

## Ciphertext only attack (COA)

Attacker only has access to the ciphertext.

Given:

$$c_1 = E_k(m_1), c_2 = E_k(m_2), \dots, c_n = E_k(m_n)$$

Find any of:

- ☐  $m_1, m_2, \dots, m_n$
- ☐  $k$
- ☐ An algorithm to infer  $m_{n+1}$  from  $c_{n+1}$

## Known Plaintext Attack (KPA)

Attacker intercepts a random plaintext / ciphertext pair:  $(m, c)$ .

Given:

$$[m_1, c_1 = E_k(m_1)], \dots, [m_n, c_n = E_k(m_n)]$$

Find any of:

- $k$
- An algorithm to infer  $m_{n+1}$  from  $c_{n+1}$

## Chosen Plaintext Attack (CPA)

Attacker selects a message  $m$  and receives the ciphertext  $c$ .  
Stronger than KPA – Some ciphers resistant to KPA are not resistant to CPA.

Given:

$$[m_1, c_1 = E_k(m_1)], \dots, [m_n, c_n = E_k(m_n)]$$

with chosen  $m$ .

Find any of:

- $k$
- An algorithm to infer  $m_{n+1}$  from  $c_{n+1}$

## Chosen Ciphertext Attack (CCA)

Attacker specifies a ciphertext  $c$  and receives the plaintext  $m$ .

Given:

$$[c_1, m_1 = D_k(c_1)], \dots, [c_n, m_n = D_k(c_n)]$$

with chosen  $c$ .

Find any of:

☐  $k$



## Rubber hose attack (RHA)

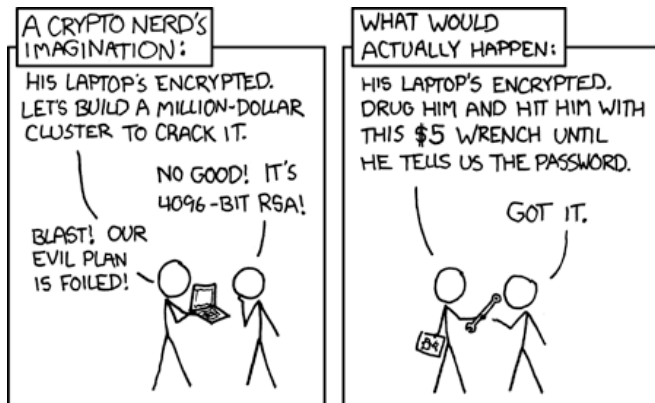
Also known as a **purchase-key attack** or **rubber-hose cryptanalysis**.

The cryptanalyst blackmails, threatens or tortures someone to retrieve  $k$ .

Extremely powerful, usually the easiest way to break a cryptosystem.

**Not permissible in Wargames**

## Rubber hose attack (RHA)



## Known Plaintext Attack

An attacker knowing that source code is being encrypted  
(the first bytes are likely `tfamily #include`, copyright notices, etc)

Famous break of the Japanese PURPLE cipher in WW2

- A complex cipher used to protect high level communications
- The allies had already broken several of the Japanese diplomatic ciphers
- PURPLE was used to protect communications but the Japanese could only afford to build and deploy 12 cipher machines
- They sent these to the twelve most important embassies
- Some messages needed to be broadcast to all embassies
- So some messages had to be sent using old ciphers the US had already broken!

## Chosen Plaintext Attack

Feed intelligence to an ambassador with the goal that it ends up encrypted and sent back home.

## Ciphertext Only Attack

Stealing cookies through weaknesses in RC4

- The start of RC4 cipher streams have biases that can be calculated when the same message is encrypted multiple times.
- RC4 was previously used in the majority of all SSL traffic; start of the stream is almost always cookies (HTTP).
- Simply trick client into making many requests, inspect the traffic, then steal the cookies.

# Classes Of Break

From worst to least severe:

**Total Break:** Attacker finds secret key  $k$  and hence can compute all  $D_k(c)$  and perform  $E_k(m)$ .

**Global Deduction:** Attacker finds alternate algorithm  $A$ , equivalent to all  $D_k(c)$  without finding  $k$ .

**Local Deduction:** Attacker finds or decrypts the plaintext of one intercepted ciphertext  
a.k.a **Instance Deduction**  
(RC4 example from previous page – attack only compromises selected  $m$ )

**Information Deduction:** Attacker gains some information about the key or plaintext  
e.g. first few bits of a file, meaning of a message, file type, etc.

An algorithm is **unconditionally secure** if no matter how much ciphertext an attacker has, there is not enough information to deduce the plaintext.

Information security is a resource game with attacks measured in terms of:

## **Data Requirements**

How much data is necessary to succeed?

## **Processing requirements (work factor)**

How much time is needed to perform the attack?

## **Memory requirements**

How much storage space is required?

## **Computational cost**

How many instances running on EC2?

# Substitution Ciphers

Substitution ciphers are the oldest form of cipher.

The secret key consists of a table which maps letter substitutions between plaintext and ciphertext.

The most famous is the *Caesar cipher* where each letter is shifted by 3 (modulo 26):

a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C

Similar to **ROT13** which shifts plaintext 13 places – largest advantage is that encrypting twice results in the plaintext:

$$ROT13(ROT13(m)) = m$$

# Substitution Ciphers

There are  $26!$  (factorial) different possible keys ( $\approx 2^{88}$  or 88-bits).

Monoalphabetic (single character) substitution cipher:

```
Src =  abcdefghijklmnopqrstuvwxyz  
Key =  XNYAHPOGZQWBTSFLRCVMUEKJDI  
  
m =   thiscourserockstheblock  
c =   MGZVYFUCVHCFYWVMGHNBFYW
```

Substitution ciphers are easy to break using **frequency analysis** of the letters:

- Single letters
- Digraphs (pairs of letters)
- Trigraphs (three letters)

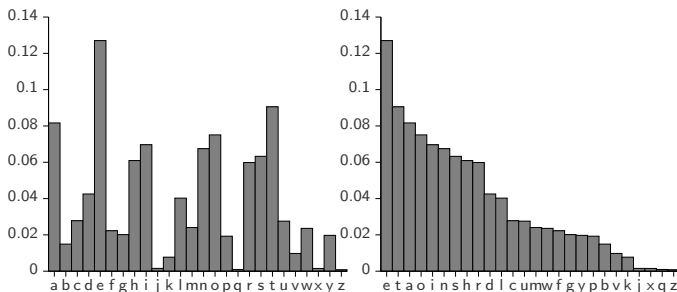
This is a **ciphertext only attack**.



# Substitution Ciphers

Substitution ciphers are easy to break by using **frequency analysis** of the letters:

- Order the histogram and the spikes should follow a similar pattern.
- Try mapping replacements, using digraph and trigraphs to check & assist you.



The frequency of English letters.

# Improved Substitution Ciphers

Homophonic ciphers are substitution ciphers that replace a common letter with multiple symbols (i.e. *E* can go to [*C*,  $\epsilon$ , *O*])

Peaks or troughs in the letter frequency are hidden as they're broken down into multiple smaller spikes.

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
=====&=====																									
D	X	S	F	Z	E	H	C	V	I	T	P	G	A	Q	L	K	J	R	U	O	W	M	Y	B	N
9				7				3					5	0				4	6						
				2																					
				1																					

“D” would become “F”, “E” would be randomly chosen as one of [*Z*, 7, 2, 1]. As a result, the high frequency of the letter “E” (the most common letter in English) is spread amongst several characters, making frequency analysis much more difficult.

# Improved Substitution Ciphers: Homophonic Ciphers

Still difficult to decipher even using modern computing:  
success rate is measured in terms of alphabet size and  
ciphertext length.

Deciphered using nested hill climbing (heuristic algorithm /  
educated guessing):

- Outer layer determines the number of symbols each letter maps to
- Inner layer determines the exact mapping

Used by the Zodiac killer:

- “Zodiac 408” was solved within days of publication
- “Zodiac 340” still remains unsolved

<http://www.cs.sjsu.edu/faculty/stamp/RUA/homophonic.pdf>

# Permutation Ciphers

**Permutation ciphers** are also known as **transposition ciphers**.  
The secret key,  $\pi$ , is a random permutation.

Given a message  $m = [m_1, m_2, m_3, \dots, m_n]$

we can compute the encryption via:

$$E_{\pi}(m) = [m_{\pi(1)} m_{\pi(2)} m_{\pi(3)}, \dots, m_{\pi(n)}]$$

Suppose  $\pi$  is:

1	2	3	4	5	6
4	3	1	5	2	6

Then  $E_{\pi}(\text{"crypto"}) = \text{"PYCTRO"}$ .

# Vigenere Cipher

Originated in Rome in the sixteenth century.

A Vigenere cipher is a polyalphabetic substitution cipher (made of multiple monoalphabetic substitution ciphers).

The secret key is a word  $k$  with encryption performed by adding the key modulo 26:

```
Plaintext  launchmissilesatlosangeles
Key        cryptocryptocryptocryptocr
=====
Cyphertext nrscvvozqhbzgjyiecurlvwzgj
```

$$L + C = 11 + 2 = 13\text{th char} = N$$

$$N + Y = 13 + 24 = (37 \% 26) = 11\text{th char} = L$$

Note: The zero index 0th character is A. Also punctuation and white space are removed to increase cryptanalysis difficulty.

# Breaking Vigenere Ciphers

The **index of coincidence** is a statistical measure of text that is used for distinguishing simple substitution ciphers from Vigenere ciphers.

Intuitively it's the measure of the probability of a collision if a string is compared to a randomly shifted version of itself.

(i.e. How many characters occur in the same spot if you shuffle the string?)

The index of coincidence is calculated by the following formula:

$$IC = \frac{\sum_{i=A}^{i=Z} F_i(F_i - 1)}{N(N - 1)}$$

Where  $F_i$  is the count of letter  $i$  (where  $i = A, B, \dots Z$ ) in the ciphertext, and  $N$  is the length of the ciphertext.

# Breaking Vigenere Ciphers

Using standard frequencies with which individual letters appear in English, the probability that a coincidence will occur is approximately  $\kappa_p = 0.0669$ .

If the text is random and letters chosen with equal probability then the probability of a coincidence is much smaller:  $\kappa_p = 0.0385$ .

The most important property of the index of coincidence is that it is the same for text encrypted with a simple substitution cipher as for plain text.

# Index of Coincidence by Language

Malay	0.085286
Dutch	0.079805
Japanese	0.077236
German	0.076667
Spanish	0.076613
Arabic	0.075889
French	0.074604
Portuguese	0.074528
Finnish	0.073796
Italian	0.073294
Danish	0.070731
Norwegian	0.069428
Greek	0.069165
<b>English</b>	<b>0.066895</b>
Swedish	0.064489
Serbo-Croatian	0.064363
Russian	0.056074
Random	0.038461



# Example: Breaking Vigenere

TPCTY LVE00 GBVRC BTWXS IHDKD QIRVQ QUKWL TMNQO EKMLP AURKL VHIUX YJRVN QWJEK UEQVD IXPLU RKLVT  
QSLKI LWAZI JWXPL QRKIO PWFME XLLCP KDIKV EUXYX EAAQV MEKVN AVZRQ JGEMX LQUPM PCRLO IZPZZ FPONI  
AYPVQ RMVHC QZFLV IKGUK LRXER MDWVG RVMPQ PWLWT TIYEQ JAYMK XBPUK PZJBF WIRKS QJMSV FYKFP QLRXE  
OIPID IQQLI ICPFP LMVBR BUAMW KLLUM FLRXE CDQFV IKNWZ KUIXF BTIII CQZQM JQVUX UVZXL XMDAY IIOIMP  
AZXEK VYIDC EGIDX NMQJQ ZQVMP FMESC EQGQD IDIJD MDXYI ACGES WSIFQ YIUMQ CBQSE EINBT CNSOM AUQLW  
BQVFL VALTS AJKLV JIZHJ MPVZQ XTLCQ ZFLDC ECVPW LRQQB TIVQV UWGPK LFTAF IKLXH BQVKL BGIEE KLFTA  
FCCEK FAQPR LEGID QVWMG MPMCC LNWDH DCPRQ DMKJX KTQXY LFFMZ SKXEA NMGVJ OQUYI CIPVQ NICMH GCZXF  
XEGUF LRXDQ LAAEM KVVFL VTFVK MYJIJ GBALV EOVPK PFZFP OWMEH KGAEM EXEGU AVEMK INAVZ RQJMQT HFMQT  
CEXTE RUMYI KSHPW IXYIT CGILV VBKVU WYSRN LIECO CQZUP ZJQWX YCJSR NCZXF XEGMP ICMGZ ZYIFP LTLRV  
FQJKV QIEIJ KMEMW PBGCZ XFXEG MFSYM AGUQX VEZJU QXFHL VPKAZ PIHWD XYSRC ZFQPK LFBTC JTFTQ FMJKL  
QLXIR HJGQZ XFXEG TMRUS CWXDM XLQPM EWHYF ESQRD ILNWD HWSOV PKRRQ BUAMO VJLTB TCIMD JBQSL WKGAE  
WROBD ZURXQ VUWGP FYQQN FVFYY NMMRU SCVPK QVVZA KGXFJ COQZI VRBOQ QWRRR FMEXI SVCTX XYIJV PMXRJ  
CNQOX DCPQC XJFVF CUFLP WBTDM RK

## Shift Index

1: .028

2: .045

3: .034

4: .037

5: .042

6: .035

**7: .070**

8: 0.32

Key length is most likely 7, since it's closest to 6.6%.

# Breaking Vigenere Ciphers

Once the key length  $N$  is known, we attack the  $N$  subtexts of the message.

Taking every  $N$ th symbol in ciphertext  $C$  gives  $N$  monoalphabetic substitution ciphers:

$$\begin{aligned} & (C_0, C_N, C_{2N}, \dots), \\ & (C_1, C_{N+1}, C_{2N+1}, \dots), \\ & \dots \\ & (C_{N-1}, C_{N+(N-1)}, C_{2N+(N-1)}, \dots) \end{aligned}$$

Note the index of coincidence varies by language and can be domain specific (e.g. may be noticeably different for a physics journal paper).