

FIGURE 1.13 Solution to substitution cipher.

| Key                                  | Message |
|--------------------------------------|---------|
| 0: LZWJWAKFGGLZWJDSFYMSVWTLXJWFUJZ   |         |
| 1: KYVIVZJEFKYYICREXLKXVSLKWIJETV    |         |
| 2: JXUHUVIDEEJXUHBQDWKQWURKJYHUDSX   |         |
| 3: IWTGTXHCDIDJWGTAPCVJPTQJUGTDCRW   |         |
| 4: HVSFSGWBCCHVSFZOBUIOUSPIHTFSBQV   |         |
| 5: GURERVFABBGUREYNATHNTRHOGSERAPU   |         |
| 6: FTQDQUEZAAFTQDXMZSGMSQNGFRDQZOT   |         |
| 7: ESPCPTDYZZESPCLYRFLRPMFEQCPYNS    |         |
| 8: DROBOSCXVYDROBKXQEKQOLEDPBOXMR    |         |
| 9: CQANARBWXXCQNAUJWPDJPNKDCOANWLQ   |         |
| 10: BPMZMQAVVWBPMZTIVOCIOJCBNZMVKP   |         |
| 11: AOLYLPZUVVAOLYSHUNBNHLIBAMYLUJO  |         |
| 12: ZNKXKXOYTUUNZKXRGTMAGMKHAZLXKTIN |         |
| 13: YMJWJNXSTYUJWQFSLZFLJGZVKKWJSHM  |         |
| 14: XLIIVIMWRSSXLIIVPERKYEKIFYXJVRGL |         |
| 15: WKHUHLVQRKRWKHUODQJXDJHEXWUHQFK  |         |
| 16: VIJGTGKUPQVJGTNCPWCIQDGVHTQPEJ   |         |
| 17: UIFSFJTOPPUIFSMBOHVBHFCVUGSFODI  |         |
| 18: THEREISNOOTHERELANGUAGEBUTFRENCH |         |
| 19: SGGQDHRMNSGGDKZMFTZFDATEEQDMBG   |         |
| 20: RFPFCGQLMRFPCJYLESYECZSRPCLAF    |         |
| 21: QEBQBFKLLQEBQIKKDRXDBYRQCOBKZE   |         |
| 22: PDANAEQJJKPDANHWJCQWCAXQPBNAJYD  |         |
| 23: OCZMZDNIJJOCZMGVIBPVZBZPOAMZIXC  |         |
| 24: NBYLYCMHIIINBYLFUHAOUAYYONZLYHWB |         |
| 25: MAXKXBLGHMAXKETGZNTZXUNMYKXGVA   |         |

$M = \text{THERE IS NO OTHER LANGUAGE BUT FRENCH.}^\dagger$

Because only one of the keys ( $K = 18$ ) produces a meaningful message, we have:

$$p_C(M) = 1$$

$$p_C(M') = 0, \text{ for every other message } M'$$

$$p_M(C) = p(18) = \frac{1}{26}$$

$$p_M(C) = 0, \text{ for every other message } M'. \quad \blacksquare$$

#### Example:

With a slight modification to the preceding scheme, we can create a cipher having perfect secrecy. The trick is to shift each letter by a random amount. Specifically,  $K$  is given by a stream  $k_1, k_2, \dots$ , where each  $k_i$  is a random

<sup>†</sup> From S. Gorn's Compendium of Rarely Used Cliches.

integer in the range  $\{0, 25\}$  giving the amount of shift for the  $i$ th letter. Then the 31-character ciphertext  $C$  in the preceding example could correspond to any valid 31-character message, because each possible plaintext message is derived by some key stream. For example, the plaintext message

**THIS SPECIES HAS ALWAYS BEEN EXTINCT.**<sup>†</sup>

is derived by the key stream

18, 18, 14, 17, 4, ....

Though most of the 31-character possible plaintext messages can be ruled out as not being valid English, this much is known even without the ciphertext. Perfect secrecy is achieved because interception of the ciphertext does not reveal anything new about the plaintext message.

The key stream must not repeat or be used to encipher another message. Otherwise, it may be possible to break the cipher by correlating two ciphertexts enciphered under the same portion of the stream (see Section 2.4.4).  $\blacksquare$

A cipher using a nonrepeating random key stream such as the one described in the preceding example is called a **one-time pad**. One-time pads are the only ciphers that achieve perfect secrecy. Implementation of one-time pads and approximations to one-time pads is studied in Chapters 2 and 3.

#### 1.4.3 Uncity Distance

Shannon measured the secrecy of a cipher in terms of the key equivocation  $H_C(K)$  of a key  $K$  for a given ciphertext  $C$ ; that is, the amount of uncertainty in  $K$  given  $C$ . From Eq. (1.2b), this is

$$H_C(K) = \sum_C p(C) \sum_K p_C(K) \log_2 \left( \frac{1}{p_C(K)} \right),$$

where  $p_C(K)$  is the probability of  $K$  given  $C$ . If  $H_C(K)$  is 0, then there is no uncertainty, and the cipher is theoretically breakable given enough resources. As the length  $N$  of the ciphertext increases, the equivocation usually decreases.

The **uncity distance** is the smallest  $N$  such that  $H_C(K)$  is close to 0; that is, it is the amount of ciphertext needed to uniquely determine the key. A cipher is **unconditionally secure** if  $H_C(K)$  never approaches 0 even for large  $N$ ; that is, no matter how much ciphertext is intercepted, the key cannot be determined. (Shannon used the term "ideal secrecy" to describe systems that did not achieve perfect secrecy, but were nonetheless unbreakable because they did not give enough information to determine the key.)

Most ciphers are too complex to determine the probabilities required to derive the uncity distance. Shannon showed, however, it is possible to approxi-

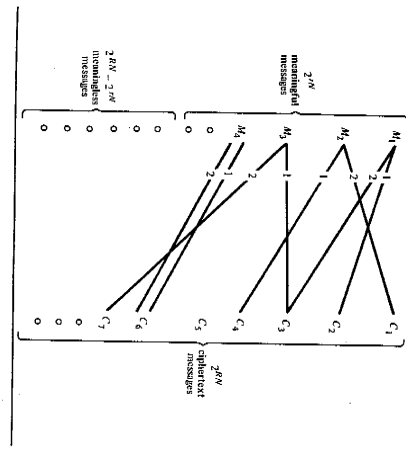
<sup>†</sup> Also from S. Gorn's Compendium of Rarely Used Cliches.

Shannon's result using a slightly different approach.

Following Hellman, we assume each plaintext and ciphertext message comes from a finite alphabet of  $L$  symbols. Thus there are  $2^{NL}$  possible messages of length  $N$ , where  $R = \log L$  is the absolute rate of the language. The  $2^{NL}$  messages are partitioned into two subsets: a set of  $2^{NR}$  meaningful messages and a set of  $2^{NL} - 2^{NR}$  meaningless messages, where  $r$  is the rate of the language. All meaningful messages are assumed to have the same prior probability  $1/2^{NR} = 2^{-rN}$ , while all meaningless messages are assumed to have probability 0.

We also assume there are  $2^{NR}$  keys, all equally likely, where  $H(K)$  is the key entropy (number of bits in the key). The prior probability of all keys is  $H(K) = 2^{-NR}$ . A random cipher is one in which for each key  $K$  and ciphertext  $C$ , the decipherment  $D_K(C)$  is an independent random variable uniformly distributed over  $2^{NL}$  messages, both meaningful and not. Intuitively, this means that for a given key  $K$ ,  $D_K(C)$  is as likely to produce one plaintext message as any other. Actually, the decipherments are not completely independent because a given key must uniquely decipher a given message, whence  $D_K(C) \neq D_{K'}(C)$  for  $C \neq C'$ .

FIGURE 1.16 Random cipher model (adapted from [Heil 77]).



decipherment or false solution arises whenever encipherment under another key  $K'$  could produce  $C$ , that is,  $C = E_{K'}(M')$  for the same message  $M'$  or  $C = E_{K'}(M')$  for another meaningful message  $M'$ . Figure 1.16 shows two spurious key decipherments, one from the third ciphertext and one from the sixth. A cryptanalyst intercepting one of these ciphertexts would be unable to break the cipher since there would be no way of picking the correct key. We are not concerned with decipherments that produce meaningless messages, because the cryptanalyst can immediately reject these solutions.

Now, for every correct solution to a particular ciphertext, there are  $(2^{NR} - 1)$  remaining keys, each of which has the same probability  $q$  of yielding a spurious key decipherment. Because each plaintext message is equally likely, the probability of getting a meaningful message and, therefore, a false solution is given by

$$q = \frac{2^{NR}}{2^{NL}} = 2^{-rN},$$

$$F = (2^{NR} - 1)q = (2^{NR} - 1)2^{-rN} \approx 2^{NR-rN}.$$

where  $D = R - r$  is the redundancy of the language. Letting  $F$  denote the expected number of false solutions, we have

$$\log_2 F = H(K) - DN = 0$$

is taken as the point where the number of false solutions is sufficiently small the cipher can be broken. Thus

$$N = \frac{H(K)}{D} \quad (1.4)$$

is the unitary distance—the amount of text necessary to break the cipher. If for given  $N$ , the number of possible keys is as large as the number of meaningful messages, then  $H(K) = \log_2(2^{NR}) = NR$ , thus

$$H(K) - DN = (R - D)N = rN \neq 0,$$

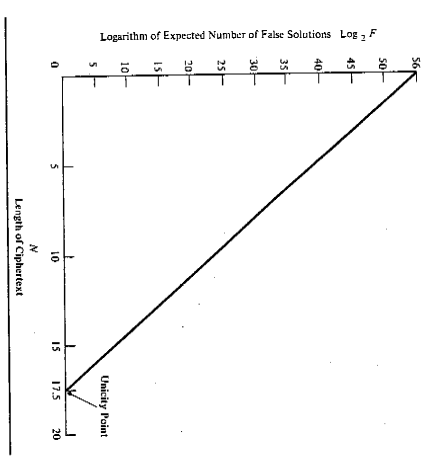
and the cipher is theoretically unbreakable. This is the principle behind the one-time pad.

**Example:**

Consider the DES, which enciphers 64-bit blocks (8 characters) using 56-bit keys. The DES is a reasonably close approximation to the random cipher model. Figure 1.17 shows  $F$  as a function of  $N$  for English language messages, where  $H(K) = 56$  and  $D = 3.2$  in Eq. (1.4). The unitary distance is

$$N = \frac{56}{3.2} = 17.5 \text{ characters,}$$

or a little over two blocks. Doubling the key size to 112 bits would double the unitary distance to 35 characters. ■



**Example:**

Consider a simple substitution cipher that shifts every letter in the alphabet forward by  $K$  positions,  $0 \leq K \leq 25$ . Then  $H(K) = \log_2 26 = 4.7$  and the unitary distance is

$$N = \frac{4.7}{3.2} = 1.5 \text{ characters.}$$

This estimate does not seem plausible, however, because no substitution cipher can be solved with just one or two characters of ciphertext. There are two problems with the approximation. First, the estimate  $D = 3.2$  applies only to reasonably long messages. Second, the cipher is a poor approximation to the random cipher model. This is because most ciphertexts are not produced by meaningful messages (e.g., the ciphertext QQQQ is produced only by the meaningless messages AAAA, BBBB, ..., ZZZZ), whence the decipherments are not uniformly distributed over the entire message space. Nevertheless, shifted ciphers can generally be solved with just a few characters of ciphertext. ■

ciphertext needed to break a cipher. Thus a particular cipher will have a unitary distance of at least  $H(K)/D$ . In practice,  $H(K)/D$  is a good approximation even for simple ciphers. We shall derive the unitary distance of several ciphers in Chapter 2. The interested reader can read more about the unitary distances of classical ciphers in Devours [Dev77].

The unitary distance gives the number of characters required to uniquely determine the key; it does not indicate the computational difficulty of finding it. A cipher may be computationally infeasible to break even if it is theoretically possible with a relatively small amount of ciphertext. Public-key systems, for example, can be theoretically broken without any ciphertext at all. The cryptanalyst, knowing the public key and the method of generating key pairs, can systematically try all possible private keys until the matching key is found (see Brassard [Bras79, Bras80]). This strategy is computationally infeasible, however, for large key spaces (e.g., with  $2^{NR}$  keys). The DES can also be broken by exhaustive search of the key space in a known-plaintext attack (by trying all keys until one is found that enciphers the plaintext into the matching ciphertext). Nevertheless, the fast known strategies for breaking the DES are extremely time-consuming. By contrast, certain substitution ciphers discussed in the next chapter use longer keys and have much greater unitary distances than DES. These ciphers are often relatively simple to solve, however, when enough ciphertext is intercepted.

Equation (1.4) shows that the unitary distance  $N$  is inversely proportional to the redundancy  $D$ . As  $D$  approaches 0, an otherwise trivial cipher becomes unbreakable. To illustrate, suppose a 6-digit integer  $M$  is enciphered as 351972 using a Caesar-type shifted substitution cipher with key  $K$ , where  $0 \leq K \leq 9$ , and that all possible 6-digit integers are equally likely. Then a cryptanalyst cannot determine which of the following integers is the value of  $M$ .

| Key | Integer |
|-----|---------|
| 0   | 351972  |
| 1   | 240861  |
| 2   | 139750  |
| ... | ...     |
| 9   | 462083  |

The reason the cipher cannot be solved is that the language has no redundancy; every digit counts.

Because of the inherent redundancy of natural languages, many ciphers can be solved by statistical analysis of the ciphertext. These techniques use frequency distributions of letters and sequences of letters, ciphertext repetitions, and probable words. Although a full discussion of these techniques is beyond the scope of this book, Chapter 2 describes how a few simple ciphers are broken using frequency distributions. (For more depth in this area, see [Kas81].) Protection against statistical analysis can be provided by several means. One way, suggested by Shannon, is by removing some of the redundancy of the language.