

# Information, Codes and Ciphers

## Summary Notes \*

Hao Ren

November 2, 2020

---

\*Based on MATH3411, UNSW 20T3.

# Contents

<b>0</b>	<b>All TODOs</b>	<b>3</b>
0.1	Unfinished Contents . . . . .	3
0.2	Examples and Diagrams . . . . .	3
0.3	Proofs . . . . .	3
<b>3</b>	<b>Compression Coding</b>	<b>4</b>
3.1	Variable Length Encoding . . . . .	4
3.1.1	UD and I-code . . . . .	4
3.1.2	Comma Codes . . . . .	4
3.1.3	Decision Trees . . . . .	4
3.1.4	The Kraft-McMillan Theorem . . . . .	4
3.1.5	Length and Variance . . . . .	4
3.2	Huffman's Algorithm . . . . .	5
3.2.1	Huffman Coding . . . . .	5
3.2.2	Properties of Huffman Codes . . . . .	6
3.2.3	Radix $r$ Huffman Codes . . . . .	6
3.2.4	Extensions of Huffman Coding . . . . .	7
3.3	Markov Sources . . . . .	7
3.3.1	Introduction . . . . .	7
3.3.2	Transition Matrix . . . . .	7
3.3.3	Huffman Coding for Stationary Markov Sources . . . . .	8
3.4	Arithmetic Coding . . . . .	8
3.4.1	Encoding . . . . .	8
3.4.2	Decoding . . . . .	8
3.5	Dictionary Methods . . . . .	9
3.5.1	Encoding . . . . .	10
3.5.2	Decoding . . . . .	10
<b>7</b>	<b>Cryptography (Ciphers)</b>	<b>11</b>
7.1	Introduction . . . . .	11
7.2	Some Classical Cryptosystems . . . . .	12
7.2.1	Caesar Cipher . . . . .	12
7.2.2	Monoalphabetic Substitution Cipher . . . . .	12
7.2.3	Transposition Cipher . . . . .	13
7.2.4	Combined Systems . . . . .	13
7.2.5	Vigenère Cipher . . . . .	13

## 0 All TODOs

Updated when *3.X Compression Coding* finished as Version 1.

### 0.1 Unfinished Contents

- 3.1.3 Decision Trees
- 3.2.4 Extensions of Huffman Coding
- 3.3.3 Huffman Coding for Stationary Markov Sources
- 3.4.2 Decoding: add a list to show steps like the one in 3.4.2 part.
- 3.6.x Other types of Compression

### 0.2 Examples and Diagrams

- 3.2.3 Eg and Dia: Radix  $r$  Huffman Codes
- 3.3.1 Eg: Definition of Markov Sources
- 3.3.3 Eg: with Huffman coding for Markov Sources
- 3.4.x Eg: Arithmetic encoding and decoding
- 3.5.x Eg: Dictionary Methods

### 0.3 Proofs

- Theorem 3.1 The Kraft-McMillan Theorem
- Theorem 3.2 Minimal UD-codes
- Theorem 3.3 Huffman Code Theorem
- Proposition 3.4 Knuth

### 3 Compression Coding

#### 3.1 Variable Length Encoding

a source $S$	with $q$ source symbols	$s_1, s_2, \dots, s_q,$
	with probabilities	$p_1, p_2, \dots, p_q,$
encoded by a code $C$	with $q$ codewords	$c_1, c_2, \dots, c_q,$
	of lengths	$l_1, l_2, \dots, l_q.$

- with a radix  $r$  codewords,
- variable length codes,
- not channel noise for source coding.

##### 3.1.1 UD and I-code

A code  $C$  is

**UD** uniquely decodable codes if it can always be decoded unambiguously,

**I-code** instantaneous if no codeword is the prefix of others.

##### 3.1.2 Comma Codes

The standard comma code of length  $n$  is

- a code which every codeword has length  $\leq n$ ,
- a code which every codeword contains at most one 0,
- and if a codeword contains 0 then 0 must be the final symbol in the codeword.

##### 3.1.3 Decision Trees

##### 3.1.4 The Kraft-McMillan Theorem

**Theorem 3.1** (The Kraft-McMillan Theorem).

*A UD-code of radix  $r$  with  $q$  codewords  $c_1, c_2, \dots, c_q$  of lengths  $l_1 \leq l_2 \leq \dots \leq l_q$  exists if and only if an I-code with the same parameters exists if and only if*

$$K = \sum_{i=1}^q \frac{1}{r^{l_i}} \leq 1.$$

##### 3.1.5 Length and Variance

The expected or **average length** of codewords is given by

$$L = \sum_{i=1}^q p_i l_i$$

and the **variance** is given by

$$V = \sum_{i=1}^q p_i l_i^2 - L^2.$$

Our aim is to minimise  $L$  for a given source  $S$  and, if more than one code  $C$  gives this value, to minimise  $V$ .

**Theorem 3.2** (Minimal UD-codes).

*Let  $C$  be a UD-code with minimal expected length  $L$  for the given source  $S$ . Then, after permuting codewords of equally likely symbols if necessary,*

- $l_1 \leq l_2 \leq \dots \leq l_q$  and
- $l_{q-1} = l_1$ .

*Furthermore, if  $C$  is instantaneous, then*

- $c_{q-1}$  and  $c_q$  differ only in their last place.

*If  $C$  is binary, then*

- $K = \sum_{i=1}^q 2^{-l_i} = 1$ .

## 3.2 Huffman's Algorithm

Huffman's algorithm for computing minimum-redundancy prefix-free codes has almost legendary status in the computing disciplines. Its elegant blend of simplicity and applicability has made it a favourite example in algorithms courses, and as a result, it is perhaps one of the most commonly implemented algorithmic techniques. [1]

### 3.2.1 Huffman Coding

In 1952, David A. Huffman [2] published a new lossless data compression method as an Sc.D student at MIT. Here is a demonstration to compute Huffman prefix-free code which is provided by Princeton University in the course COS226 [3]:

- Count character frequencies  $p_s$  for each symbol  $s$  in file.
- Start with a forest of trees, each consisting of a single vertex corresponding to each symbol  $s$  with weight  $p_s$ .
- Repeat:
  - select two trees with min weight  $p_1$  and  $p_2$
  - merge into single tree with weight  $p_1 + p_2$

### OPTIMUM BINARY CODING PROCEDURE

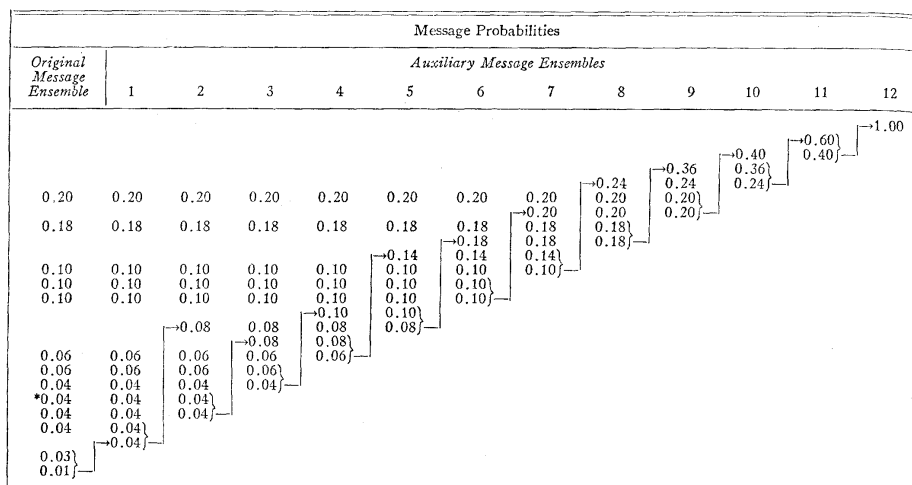


Figure 1: Huffman Coding Diagram [2]

**Applications** JPEG, MP3, MPEG, PKZIP.

**Theorem 3.3** (Huffman Code Theorem).

For the given source  $S$ , the Huffman algorithm produces a minimum average length UD-code which is an instantaneous code.

**Proposition 3.4** (Knuth).

For a Huffman code created by the given algorithm, the average code word length is sum of all the probabilities at child nodes.

### 3.2.2 Properties of Huffman Codes

1. The place high strategy always produces a minimum variance Huffman code .
2. If there are  $2^n$  equally likely source symbols then the Huffman code is a block code of length  $n$ .
3. If for all  $j$ ,  $3p_j \geq 2 \sum_{k=j+1}^q p_k$  then the Huffman code is a comma code.
4. Small changes in the  $p_i$  can change the Huffman code substantially, but have little effect on the average length  $L$ . This effect is smaller with smaller variance.

### 3.2.3 Radix $r$ Huffman Codes

For  $r$  radix Huffman codes, which  $r \geq 3$ , we have a better strategy by adding some dummy variables [4].

For a  $r$  radix encoding, the procedure is similar except  $r$  least probable symbols are merged at each step. Since the total number of symbols may not be enough to allow  $r$  variables to be merged at each step, we might need to add some dummy symbols with 0 probability before constructing the Huffman tree.

How many dummy symbols need to be added? Since the first iteration merges  $r$  symbols and then each iteration combines  $r - 1$  symbols with a merged symbols, if the procedure is to last for  $k$  (some integer number of) iterations, then the total number of source symbols needed is  $1 + k(r - 1)$ . So before beginning the Huffman procedure, we add enough dummy symbols so that the total number of symbols look like  $1 + k(r - 1)$  for the smallest possible value of  $k$ .

**Remark.**

1. *If more than two symbols have the same probability at any iteration, then the Huffman coding may not be unique (depending on the order in which they are merged). However, all Huffman codings on that alphabet are optimal in the sense they will yield the same expected code length.*
2. *One might think of another alternate procedure to assign small code lengths by building a tree top-down instead, e.g. divide the symbols into two sets with almost equal probabilities and repeating. While intuitively appealing, this procedure is suboptimal and leads to a larger expected code length than the Huffman encoding. You should try this on the symbol distribution described above.*

### 3.2.4 Extensions of Huffman Coding

## 3.3 Markov Sources

### 3.3.1 Introduction

A finite-state Markov chain is a sequence  $S_0, S_1, \dots$  of discrete random symbols from a finite alphabet,  $S$ . There is a probability mass function (pmf)  $q_0(s), s \in S$  on  $S_0$ , and there is a conditional pmf  $Q(s|s')$  such that for all  $m \geq 1$ , all  $s \in S$ , and all  $s' \in S$ ,

$$\Pr(S_k = s | S_{k-1} = s') = \Pr(S_k = s | S_{k-1} = s', \dots, S_0 = s_0) = Q(s|s')$$

There is said to be a *transition* from  $s'$  to  $s$ , denoted  $s' \rightarrow s$ , if  $Q(s|s') > 0$ . [5]

### 3.3.2 Transition Matrix

The matrix  $M = (p_{ij})$  is called the transition matrix of the Markov process, which could be displayed as (from  $s_j \rightarrow (p_{ij}) \rightarrow (to\ s_i)$ ).

**Remark.**

- $P(s_1|s_j) + P(s_2|s_j) + \dots + P(s_q|s_j) = 1$
- $p_{1j} + p_{2j} + \dots + p_{qj} = 1$ , for  $j = 1, \dots, q$

$$\begin{bmatrix} p_{11} & p_{12} & \cdots & p_{1j} \\ p_{21} & p_{22} & \cdots & p_{2j} \\ \vdots & \vdots & \ddots & \vdots \\ p_{i1} & p_{i2} & \cdots & p_{ij} \end{bmatrix}$$

### 3.3.3 Huffman Coding for Stationary Markov Sources

## 3.4 Arithmetic Coding

Arithmetic coding is very efficient and approaches the entropy limit faster than Huffman coding without any contradiction because the arithmetic coding is not a UD-code. The idea is to assign a subinterval of  $[0, 1) \subseteq \mathbb{R}$  to the message and successively narrow this subinterval down as each symbol is encoded. The message must end with a stop symbol  $\bullet$ , which could be displayed as

$$s_a s_b s_c s_d s_e \cdots s_n \bullet,$$

and after the subinterval corresponding to the message plus  $\bullet$  is found, then any suitable *single number* in the subinterval is transmitted – this is the actual code number or codeword.

### 3.4.1 Encoding

Firstly, to each symbol is associated a subinterval of  $[0, 1)$  whose length is proportional to the relative frequency of the symbol. These subintervals are chosen so they do not overlap, but fill  $[0, 1)$ .

If the “message so far” subinterval is  $[s, s + w)$  where  $s$  is the start,  $w$  is the width and the next symbol has associated subinterval  $[s', s' + w')$  then the new message subinterval is the  $[s', s' + w')$  part of  $[s, s + w)$ .

That is, the new message subinterval becomes

$$[s + s'w, (s + s'w) + w'w).$$

It could be shown as a list step by step:

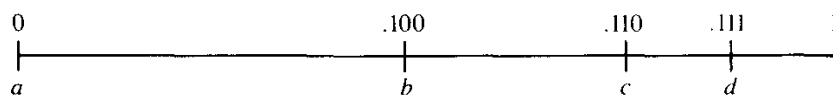
- **Crop**  $[0, 1)$  as current interval.
- **Partition** the current interval  $[0, 1)$  into subintervals of length  $p_1 \cdots p_q$ , which  $p_q$  is a stop symbol  $\bullet$ .
- **Crop** the  $i_1^{\text{th}}$  subinterval as the current interval.
- **Partition** the current interval [new start, new end) into subintervals of length  $p_1 \cdots p_q$ , which  $p_q$  is a stop symbol  $\bullet$ .
- **Crop** the  $i_2^{\text{th}}$  subinterval as the current interval, which is also a sub-subinterval of the initial interval  $[0, 1)$ .
- **Repeat** this process until the whole message has been encoded.

### 3.4.2 Decoding

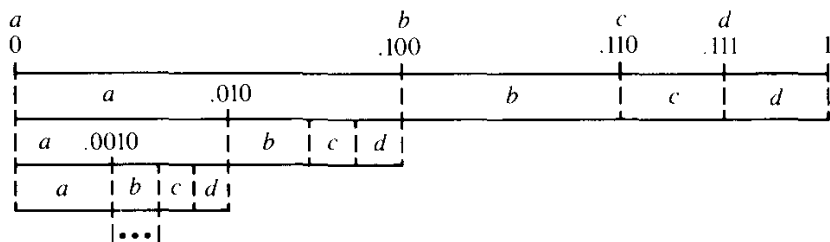
We reverse the above, by finding which of the symbol subintervals contains the code number. This then gives the first symbol of the message.

The code number is then re-scaled to that interval as follows. If  $x$  is the code number and it lies in the subinterval  $[s, s + w)$ , then the rescaled code number is  $\frac{x - s}{w}$ . The process is repeated until the stop symbol  $\bullet$  is encountered.





**Figure 1** Codewords of Table 1 as points on unit interval.



**Figure 2** Successive subdivision of unit interval for code of Table 1 and data string “ $a a b \dots$ ”

Figure 2: Arithmetic Coding Diagram [6]

**Remark.**

1. Without the stop symbol, decoding would usually go on forever.
2. It is possible to implement arithmetic coding and decoding using integer arithmetic, and if we use binary representation for the subintervals then it can actually be done using binary arithmetic. Successive digits can be transmitted as soon as they are finalized; that is, when both start and end points of the interval agree to sufficiently many significant places. Decoding can also be done as the digits arrive. Hence arithmetic coding does not need much memory.
3. We have described a static version of arithmetic coding. But the process can also be made adaptive by continually adjusting the start/width values for symbols while decoding takes place.

### 3.5 Dictionary Methods

**Examples** Ziv, LZ77, LZ78 and LZW are four dictionary compressing methods. Dictionary-based algorithms do not encode single symbols as variable-length bit strings; they encode variable-length strings of symbols as single tokens:

- The tokens form an index into a phrase dictionary.
- If the tokens are smaller than the phrases they replace, compression occurs.

Consider the Random House Dictionary of the English Language, Second edition, Unabridged. Using this dictionary, the string

A good example of how dictionary based compression works

can be coded as

1/1 822/3 674/4 1343/60 928/75 550/32 173/46 421/2

**Coding** Uses the dictionary as a simple lookup table.

- Each word is coded as  $x/y$ , where,  $x$  gives the page in the dictionary and  $y$  gives the number of the word on that page.
- The dictionary has 2,200 pages with less than 256 entries per page: Therefore  $x$  requires 12 bits and  $y$  requires 8 bits, i.e., 20 bits per word (2.5 bytes per word).
- Using ASCII coding the above string requires 48 bytes, whereas our encoding requires only 20 ( $2.5 \times 8$ ) bytes: 50% compression. [7]

**Applications** GZIP, GIF, POSTSCRIPT.

### 3.5.1 Encoding

Start with an empty dictionary and set  $r = m$ . (Here  $r$  is the part of the message which we have not yet encoded.) Let “+” denote concatenation of symbols. Repeat the following steps until  $r$  is empty:

1. Let  $s$  be the longest prefix of  $r$  which corresponds to a dictionary entry. If there is no such prefix then set  $s = \emptyset$ . (By prefix we mean a subsequence of  $r$  which begins at the first symbol.)
2. Suppose that  $s$  is in line  $l$  of the dictionary, and put  $l = 0$  if  $s = \emptyset$ . Add a new dictionary entry for  $s + c$ , where  $c$  is the next symbol after  $s$  in  $r$ . Output  $(l, c)$  to the encoding and delete  $s + c$  from  $r$ .

### 3.5.2 Decoding

1. Start with an empty dictionary.
2. If the next code pair is  $(l, c)$  then output  $D(l) + c$  to the message, where  $D(l)$  denotes entry  $l$  of the dictionary and  $D(0) = \emptyset$ .
3. Also add  $D(l) + c$  to the dictionary. Repeat until all code pairs have been processed.
4. The dictionary after the decoding process is exactly the same as the dictionary after the encoding process.

## 7 Cryptography (Ciphers)

### 7.1 Introduction

A **message**  $m$  from among a set  $M$  of all possible messages is **encrypted** (or enciphered or encoded), using one of a set  $\{E_k\}$  of **encryption functions** (or cryptographic transformations)

$$E_k : M \rightarrow C$$

into an **encrypted message**  $c \in C$ , where  $C$  is the set of all possible encrypted messages.

The encryption function  $E_k$  depends on the particular **key**  $k$  chosen from a set  $K$  of all possible keys.

The encrypted message  $c$  is then transmitted to the receiver, who decrypts it using the decryption function

$$D_k : C \rightarrow M$$

using the same key  $k$ . Here  $D_k$  is the inverse of  $E_k$ , which means that

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

We also call  $m$  the **plaintext** and  $c$  the **cipher text**.

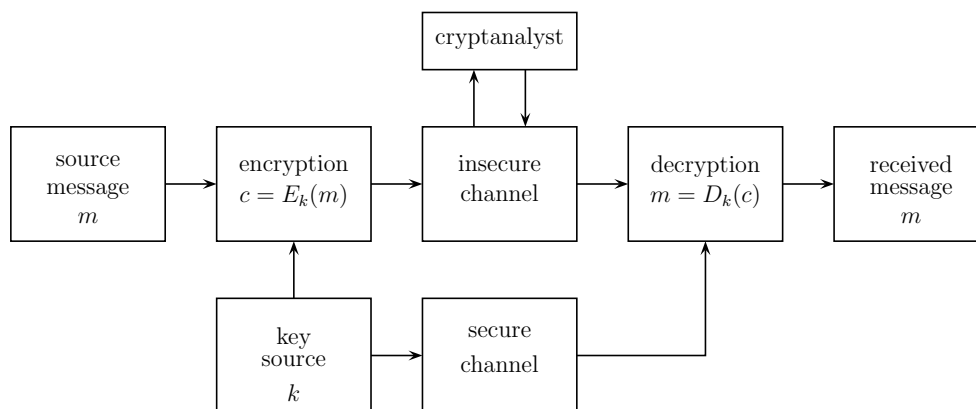


Figure 3: Cryptography

**Principle 7.1** (Kerckhoffs's Principle of Cryptography [8]).

1. The system must be practically, if not mathematically, indecipherable;
2. It should not require secrecy, and it should not be a problem if it falls into enemy hands;
3. It must be possible to communicate and remember the key without using written notes, and correspondents must be able to change or modify it at will;
4. It must be applicable to telegraph communications;

5. *It must be portable, and should not require several persons to handle or operate;*
6. *Lastly, given the circumstances in which it is to be used, the system must be easy to use and should not be stressful to use or require its users to know and comply with a long list of rules.*

Claude Shannon reformulated this principle as "one ought to design systems under the assumption that the enemy will immediately gain full familiarity with them", which called

**Shannon's Maxim** "The enemy knows the system."

## 7.2 Some Classical Cryptosystems

### 7.2.1 Caesar Cipher

Cyclically shift each letter  $k$  places forward, e.g. for  $k = 3$ , which is also used by Roman dictator Gaius Julius Caesar himself.

plain	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
cipher	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

If letters are represented by  $\mathbb{Z}_{26}$ , then we have, for key  $k$ ,

$$E_k(i) = i + k \pmod{26} \quad \text{for all } i$$

and

$$D_k(j) = j - k \pmod{26} \quad \text{for all } j.$$

There are only 26 keys  $k$ , so this cipher is easily broken by simply trying all possible keys until it makes sense.

### 7.2.2 Monoalphabetic Substitution Cipher

Replace each letter by another character, such as

plain	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
cipher	P Q S T U V W X Y Z C O D E B R A K I N G F H J L M

The keyword "CODEBREAKING" starting at K leads to a more memorable permutation, but the code would be easier to break.

Let  $\pi$  be the permutation of  $\mathbb{Z}_{26}$  used, then

$$E_k(i) = \pi(i) \quad \text{for all } i$$

and

$$D_k(j) = \pi^{-1}(j) \quad \text{for all } j.$$

If we only use alphabetical characters as key  $k$ , there are still  $26! \approx 4 \times 10^{26}$  possible keys so it is impossible to try them all.

### 7.2.3 Transposition Cipher

Mix up the order of the letters by dividing them up into blocks of length  $l$  and applying a permutation  $\pi \in \text{Sym}(l)$  to the letter order. For example, with blocks of length 5 we could use the permutation

$$\pi = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 3 & 1 & 4 & 5 & 2 \end{pmatrix} \quad \text{map 1st letter} \rightarrow \text{3rd letter, 2nd} \rightarrow \text{1st, and so on.}$$

The  $\text{Sym}(r)$  means the symmetric group. [9]

plain	T	H	I	S	I	S	A	N	E	X	A	M	P	L	E
cipher	H	I	T	I	S	A	X	S	N	E	M	E	A	P	L

Frequency counts of pairs or triples and common sense allow these ciphers to be easily broken.

### 7.2.4 Combined Systems

These use both a transposition and a substitution. They are much harder to break by hand but are still breakable using frequency counts.

### 7.2.5 Vigenère Cipher

## References

- [1] A. Moffat, “Huffman coding,” *ACM Comput. Surv.*, vol. 52, no. 4, August 2019.
- [2] D. A. Huffman, “A method for the construction of minimum-redundancy codes,” *Proceedings of the Institute of Radio Engineers*, vol. 40, no. 9, pp. 1098–1101, September 1952.
- [3] R. Sedgewick and K. Wayne, “5.5 data compression,” <https://www.cs.princeton.edu/courses/archive/fall20/cos226/lectures/55DataCompression.pdf>, 2020.
- [4] A. Singh, “Lecture 8: Source coding theorem, huffman coding,” <http://www.cs.cmu.edu/~aarti/Class/10704/lec8-srccodingHuffman.pdf>, 2012.
- [5] L. Zheng and R. Gallager, “Chapter 2 coding for discrete sources,” [https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-450-principles-of-digital-communications-i-fall-2006/lecture-notes/book\\_2.pdf](https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-450-principles-of-digital-communications-i-fall-2006/lecture-notes/book_2.pdf), 2006.
- [6] G. G. Langdon, “An introduction to arithmetic coding,” *IBM Journal of Research and Development*, vol. 28, no. 2, pp. 135–149, 1984.
- [7] R. Yerraballi, “4.2.3 dictionary-based coding: Lzw,” <https://users.ece.utexas.edu/~ryerraballi/MSB/index.html>.
- [8] A. Kerckhoffs, “La cryptographie militaire,” *Journal des sciences militaires*, vol. IX, pp. 5–38, January 1883.
- [9] Wikipedia contributors, “Symmetric group — Wikipedia, the free encyclopedia,” [https://en.wikipedia.org/w/index.php?title=Symmetric\\_group&oldid=964788022](https://en.wikipedia.org/w/index.php?title=Symmetric_group&oldid=964788022), 2020, [Online; accessed 21-November-2020].