

FLUENCY⁶

with information technology

SKILLS, CONCEPTS, & CAPABILITIES



LAWRENCE SNYDER

Chapter 7

Representing Information Digitally

PEARSON

ALWAYS LEARNING

Learning Objectives

- Explain the link between patterns, symbols, and information
- Determine possible PandA encodings using a physical phenomenon
- Encode and decode ASCII
- Represent numbers in binary form
- Compare two different encoding methods
- Explain how structure tags (metadata) encode the Oxford English Dictionary

Digitizing Discrete Information

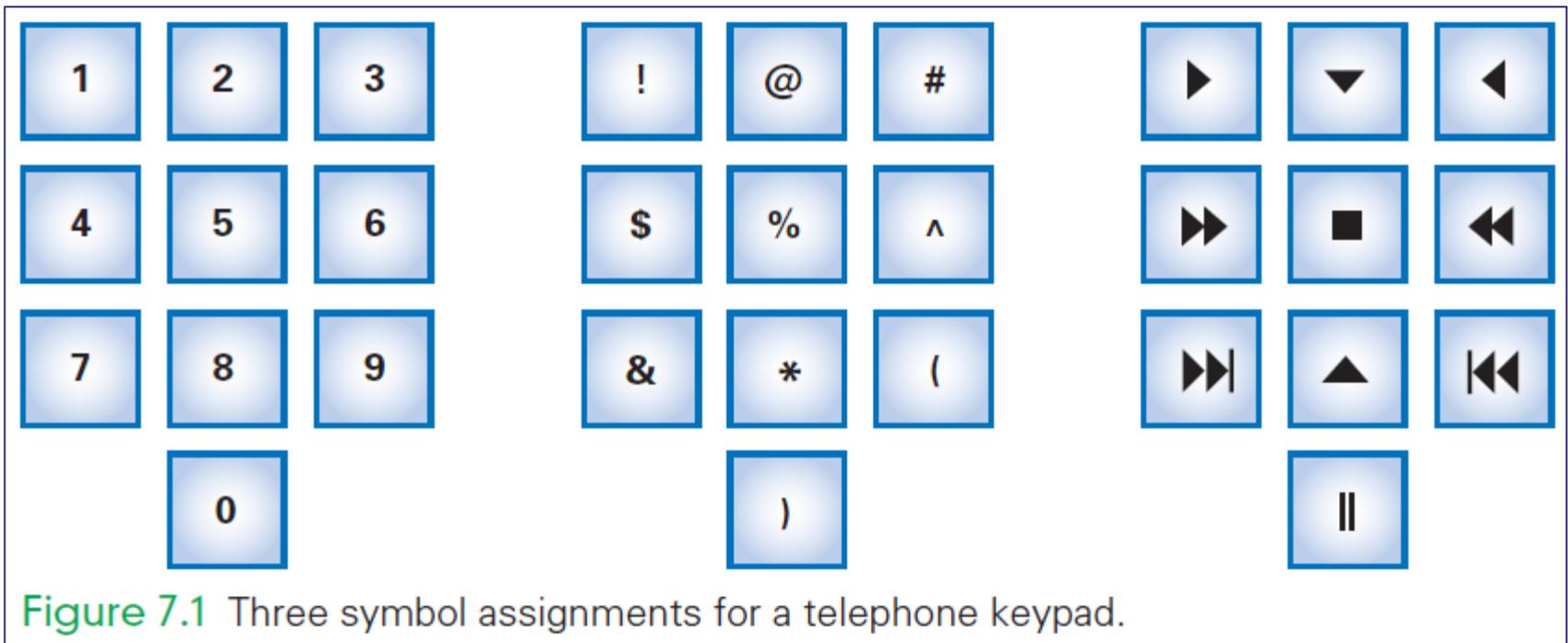
- The dictionary definition of digitize is to represent information with digits
- Digit means the ten Arabic numerals 0 through 9

Limitation of Digits

- A limitation of the dictionary definition of digitize is that it calls for the use of the ten digits, which produces a whole number
 - Digitizing in computing can use almost any symbol
- Having a bigger telephone number does not make you a better person

Alternative Representations

Digitizing can use almost any symbols



Symbols, Briefly

- One practical advantage of digits is that digits have short names (one, two, etc)
- Imagine speaking your phone number the multiple syllable names:
 - “asterisk, exclamation, closing parenthesis”
- IT uses these symbols, but have given them shorter names:
 - exclamation point . . . is *bang*
 - asterisk . . . is *star*

Ordering Symbols

- Another advantage of digits is that the items can be listed in numerical order
- Sometimes ordering items is useful
- ***collating sequence***: placing information in order by using non-digit symbols
 - need to agree on an ordering for the basic symbols
- Today, digitizing means representing information by symbols

Fundamental Information Representation

- A computer represents information by the presence or absence of some physical phenomenon
- This gives two symbols: a binary system
- We name the two states 1 and 0
- We can then build larger symbols using these two basic ones
- The phenomenon can be charge, current, magnetization, or many other things

The PandA Representation

- We refer to this system as “Presence and Absence” or PandA.
- Such a formulation is said to be discrete
- Discrete means “*distinct*” or “*separable*”
 - It is not possible to transform one value into another by tiny gradations
 - There are no “**shades of gray**”

A Binary System

- The PandA encoding has two patterns: present and absent
- Two patterns make it a binary system
- There is *no* law that says *on* means “*present*” or *off* means “*absent*”

Table 7.1 Possible interpretations of the two PandA alternatives

Present	Absent
True	False
1	0
On	Off
Yes	No
+	–
Black	White
For	Against
...	...

Bits Form Symbols

- *In the PandA representation, the unit is a specific place (in space and time), where the presence or absence of the phenomenon can be set and detected.*
- The PandA unit is known as a **bit**
- **Bit** is a contraction for “**binary digit**”
- Bit sequences can be interpreted as binary numbers
- Groups of bits form symbols

Bits in Computer Memory

- Memory is arranged inside a computer in a very long sequence of bits
 - the physical phenomenon can be encoded, the information can be ***set*** and ***detected*** to ***present*** or ***absent***

Alternative PandA Encodings

- There is no limit to the ways to encode two physical states
 - stones on all squares, but with white (absent) and black (present) stones for the two states
 - multiple stones of two colors per square, more white stones than black means 1 and more black stones than white means 0
 - And so forth

Combining Bit Patterns

- The two-bit patterns gives limited resources for digitizing information
- Only two values can be represented
- The two patterns must be combined into sequences to create enough symbols to encode the intended information

Table 7.2 Number of symbols possible from n bits in sequence

n	2^n	Symbols
1	2^1	2
2	2^2	4
3	2^3	8
4	2^4	16
5	2^5	32
6	2^6	64
7	2^7	128
8	2^8	256
9	2^9	512
10	2^{10}	1,024

Binary Explained

- Computers use base-2 to represent numbers using the *binary number system*
- When counting in binary you are limited to only use 0 and 1
 - 0, 1, 10, 11, 100, 101, 110, 111, 1000, ...

Hex Explained

- **Hex** digits, short for **hexadecimal** digits, are base-16 numbers
- Uses decimal numbers, and then the first six Latin letters
 - 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F
- There needed to be a better way to write bit sequences...hexadecimal digits

Changing Hex to Binary

- The 32 bits below represent a computer instruction

1000 1110 1101 1000 1010 0011 1010 0000

- Writing so many 0's and 1's is tedious and error prone
- We can convert each each four-bit group to hex, giving the shorter version:

8E D8 A3 A0

Hexadecimal Digits

Each hex digit codes a four-bit group:

0000	0	0100	4	1000	8
	1100				C
0001	1	0101	5	1001	9
	1101				D
0010	2	0110	6	1010	A
	1110				E
0011	3	0111	7	1011	B
	1111				F

Digitizing Numbers in Binary

- The two earliest uses of PandA were to:
 - Encode numbers
 - Encode keyboard characters
- Representations for sound, images, video, and other types of information are also important

Place Value in a Decimal Number

- Recall that To find the quantity expressed by a decimal number:
 - The digit in a place is multiplied by the place value and the results are added
- Example, 1,010 (base 10) is:
 - Digit in the 1's place is multiplied by its place
 - Digit in the 10's place is multiplied by its place
 - and so on:
 $(0 \times 1) + (1 \times 10) + (0 \times 100) + (1 \times 1000)$

Place Value in a Binary Number

- Binary works the same way
- The base is 2 instead of 10
- Instead of the decimal place values: 1, 10, 100, 1000, . . . , the binary place values are: 1, 2, 4, 8, 16, . . . ,

Power	Decimal	Binary
0	$1 = 10^0$	$1 = 2^0$
1	$10 = 10^1$	$2 = 2^1$
2	$100 = 10^2$	$4 = 2^2$
3	$1000 = 10^3$	$8 = 2^3$
4	$10,000 = 10^4$	$16 = 2^4$
...

Place Value in a Binary Number

- 1010 in binary:
 - $(1 \times 8) + (0 \times 4) + (1 \times 2) + (0 \times 1)$

Table 7.5 The binary number 1010, representing the decimal number ten = 8 + 2

2^3	2^2	2^1	2^0	Binary Place Values
1	0	1	0	Bits of Binary Number
1×2^3	0×2^2	1×2^1	0×2^0	Multiply place bit by place value
8	0	2	0	and add to get a decimal 10

Table 7.6 Binary representation of the decimal number one thousand ten = 11 1111 0010

2^9	2^8	2^7	2^6	2^5	2^4	2^3	2^2	2^1	2^0	Binary Place Values
1	1	1	1	1	1	0	0	1	0	Bits of Binary Number
1×2^9	1×2^8	1×2^7	1×2^6	1×2^5	1×2^4	0×2^3	0×2^2	1×2^1	0×2^0	Multiply place bit by place value
512	256	128	64	32	16	0	0	2	0	and add to get decimal 1,010

Digitizing Text

- The number of bits determines the number of symbols available for representing values:
 - *n bits in sequence yield 2^n symbols*
- The more characters you want encoded, the more symbols you need

Digitizing Text

- Roman letters, Arabic numerals, and about a dozen punctuation characters are the minimum needed to digitize *English* text
- What about:
 - Basic arithmetic symbols like +, −, *, /, =?
 - Characters not required for English ö, é, ñ, ø?
 - Punctuation? « », ¿, π, √)? What about business symbols: ¢, £, ¥, ©, and ®?

Assigning Symbols

- We need to represent:
 - 26 uppercase,
 - 26 lowercase letters,
 - 10 numerals,
 - 20 punctuation characters,
 - 10 useful arithmetic characters,
 - 3 other characters (new line, tab, and backspace)
 - 95 symbols...enough for English

Assigning Symbols

- To represent 95 distinct symbols, we need 7 bits
 - 6 bits gives only $2^6 = 64$ symbols
 - 7 bits give $2^7 = 128$ symbols
- 128 symbols is ample for the 95 different characters needed for English characters
- Some additional characters must also be represented

Assigning Symbols

- **ASCII** stands for American Standard Code for Information Interchange
- ASCII is a widely used 7-bit (2^7) code
- Advantages of a “standard”:
 - Computer parts built by different manufacturers can be connected
 - Programs can create data and store it so that other programs can process it later, and so forth

Extended ASCII: An 8-Bit Code

- 7-bit ASCII is not enough, it cannot represent text from other languages
- IBM decided to use the next larger set of symbols, the **8-bit** symbols (2^8)
- Eight bits produce $2^8 = 256$ symbols
 - The 7-bit ASCII is the 8-bit ASCII representation with the leftmost bit set to 0
 - Handles many languages that derived from the Latin alphabet

Extended ASCII: An 8-Bit Code

- IBM gave 8-bit sequences a special name, *byte*
- It is a standard unit for computer memory

Unicode

- The 256 extended ASCII codes cover most Western languages
- Unicode represents many more characters by using up to 32 bits to code characters
- UTF-8 records Unicode by writing long characters as groups of bytes

ASCII	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
0000	N _U	S _H	S _X	E _X	E _T	E _O	A _K	B _L	B _S	H _T	L _F	Y _T	F _F	C _R	S _O	S _I
0001	D _L	D ₁	D ₂	D ₃	D ₄	N _K	S _V	E _S	C _N	E _M	S _B	E _C	F _S	G _S	R _S	U _S
0010		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
0011	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
0100	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
0101	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
0110	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
0111	p	q	r	s	t	u	v	w	x	y	z	{		}	~	D _T
1000	S ₀	S ₁	S ₂	S ₃	I _N	N _L	S _S	E _S	H _S	H _J	Y _S	P _D	P _V	R _I	S ₂	S ₃
1001	D _C	P ₁	P ₂	S _E	C _C	M _M	S _P	E _P	O _S	O _A	O _A	C _S	S _T	O _S	P _M	A _P
1010	°	i	ç	£	¤	¥		§	¨	©	ª	«	¬	-	®	¯
1011	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
1100	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
1101	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
1110	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
1111	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

Figure 7.3 ASCII, the American Standard Code for Information Interchange.

Note: The original 7-bit ASCII is the top half of the table; the whole table is known as Extended ASCII (ISO-8859-1). The 8-bit symbol for a letter is the four row bits followed by the four column bits (e.g., A = 0100 0001, while z = 0111 1010). Characters shown as two small letters are control symbols used to encode nonprintable information (e.g., B_S = 0000 1000 is backspace). The bottom half of the table represents characters needed by Western European languages, such as Icelandic's eth (ð) and thorn (þ).

Advantages of Long Encodings

- With computing, we usually try to be efficient by using the shortest symbol sequence to minimize the amount of memory
- Examples of the opposite:
 - NATO Broadcast Alphabet
 - Bar Codes

NATO Broadcast Alphabet

- The code for the letters used in radio communication is *purposely* inefficient
- The code is distinctive when spoken amid **noise**
- The alphabet encodes letters as words
 - Words are the symbols
 - “Mike” and “November” replace “em” and “en”
- The longer encoding improves the chance that letters will be recognized
- Digits keep their usual names, except nine, which is known as *niner*

NATO Broadcast Alphabet

Table 7.7 NATO broadcast alphabet designed not to be minimal

A	Alpha	H	Hotel	O	Oscar	V	Victor
B	Bravo	I	India	P	Papa	W	Whiskey
C	Charlie	J	Juliet	Q	Quebec	X	X-ray
D	Delta	K	Kilo	R	Romeo	Y	Yankee
E	Echo	L	Lima	S	Sierra	Z	Zulu
F	Foxtrot	M	Mike	T	Tango		
G	Golf	N	November	U	Uniform		

Bar Codes

- Universal Product Codes (UPC) also use more than the minimum number of bits to encode information
- In the UPC-A encoding, 7 bits are used to encode the digits 0 – 9



Bar Codes

- UPC encodes the manufacturer (left side) and the product (right side)
- Different bit combinations are used for each side
- One side is the complement of the other side
- The bit patterns were chosen to appear as different as possible from each other

Bar Codes

- Different encodings for each side make it possible to recognize whether the code is right side up or upside down

Table 7.8 Bit encoding for bars for the UPC-A encoding; notice that the two sides are complements of each other, that is, the 0's and 1's are switched.

Digit	Left Side	Right Side
0	0001101	1110010
1	0011001	1100110
2	0010011	1101100
3	0111101	1000010
4	0100011	1011100
5	0110001	1001110
6	0101111	1010000
7	0111011	1000100
8	0110111	1001000
9	0001011	1110100

Metadata and the OED

- Converting the content into binary is half of the problem of representing information
- The other half of the problem? Describing the information's properties
- Characteristics of the content also needs to be encoded:
 - How is the content structured?
 - What other content is it related to?
 - Where was it collected?
 - What units is it given in?
 - How should it be displayed?
 - When was it created or captured?
 - And so on...

Metadata and the OED

- Metadata: information describing information
- Metadata is the third basic form of data
- It does not require its own binary encoding
- The most common way to give metadata is with tags (think back to Chapter 4 when we wrote HTML)

Properties of Data

- Metadata is separate from the information that it describes
- For example:
 - The ASCII representation of letters has been discussed
 - How do those letters look in **Comic Sans** font?
 - How they are displayed is metadata

Properties of Data

- Rather than fill a file with the **Times New Roman** font, the file is filled with the letters and tags that describe how it should be displayed
- This avoids locking-in the form of display, which can be changed by changing the metadata
- Metadata can be presented at several levels

Using Tags for Metadata

- The Oxford English Dictionary (OED) is the definitive reference for every English word's meaning, etymology, and usage
- The printed version of the OED is truly monumental
 - 20 volumes
 - 150 pounds
 - 4 feet of shelf space

Using Tags for Metadata

- In 1984, the conversion of the OED to digital form began
- Imagine, with what you know about searching on the computer, finding the definition for the verb **set**.
 - “set” is part of many words
 - You would find closet, horsetail, settle, and more
- Software can help sort out the words

Structure Tags

- Special tags can be developed to handle structure:
 - **<hw>** is the OED's tag for a headword (word being defined)
 - **<pr>** handles pronunciation
 - **<ph>** does the phonetic notations
 - **<ps>** parts of speech
 - **<hm>** homonym numbers
 - **<e>** surrounds the entire entry
 - **<hg>** surrounds the head group or all of the information at the start of a definition

Structure Tags

- With structure tags, software can use a simple algorithm to find what is needed
- Tags do not print
- They are included only to specify the structure, so the computer knows what part of the dictionary to use
- Structure tags are also useful for formatting...for example, boldface used for headwords
- Knowing the structure makes it possible to generate the formatting information

Parity

- Computer memory is subject to errors
- An extra bit is added to the memory to help detect errors
 - A ninth bit per byte can detect errors using parity
- ***Parity*** refers to whether a number is even or odd
 - If the number of 1's is even, set the ninth bit to 0; otherwise set it to 1

Parity

- All 9-bit groups have even parity:
 - Any single bit error in a group causes its parity to become odd
 - This allows hardware to detect that an error has occurred
 - It cannot detect *which* bit is wrong, however

Why “Byte”?

- IBM was building a supercomputer, called Stretch
- They needed a word for a quantity of memory *between* a bit and a word
 - A word of computer memory is typically the amount required to represent computer instructions (currently a word is 32 bits)

Why “Byte”?

- Then, why not bite?
- The ‘i’ to a ‘y’ was done so that someone couldn’t accidentally change ‘byte’ to ‘bit’ by the dropping the ‘e’ ”
 - bite bit (*the meaning changes*)
 - byte byt (*what’s a byt?*)

Summary

- We began the chapter by learning that digitizing doesn't require digits—any symbols will do
- We explored the following:
 - PandA encoding, which is based on the presence and absence of a physical phenomenon
 - Their patterns are discrete; they form the basic unit of a bit. Their names (most often 1 and 0) can be any pair of opposite terms

Summary

- We explored the following:
 - A bit's 0 and 1 states naturally encourage the representation of numbers in base 2, that is, binary
 - 7-bit ASCII, an early assignment of bit sequences (symbols) to keyboard characters. Extended or 8-bit ASCII is the standard
 - The need to use more than the minimum number of bits to encode information

Summary

- We explored the following:
 - How documents like the *Oxford English Dictionary* are digitized
 - We learned that tags associate metadata with every part of the *OED*
 - Using that data, a computer can easily help us find words and other information.
 - The mystery of the *y in byte*