

EECE 523

Ch. 1

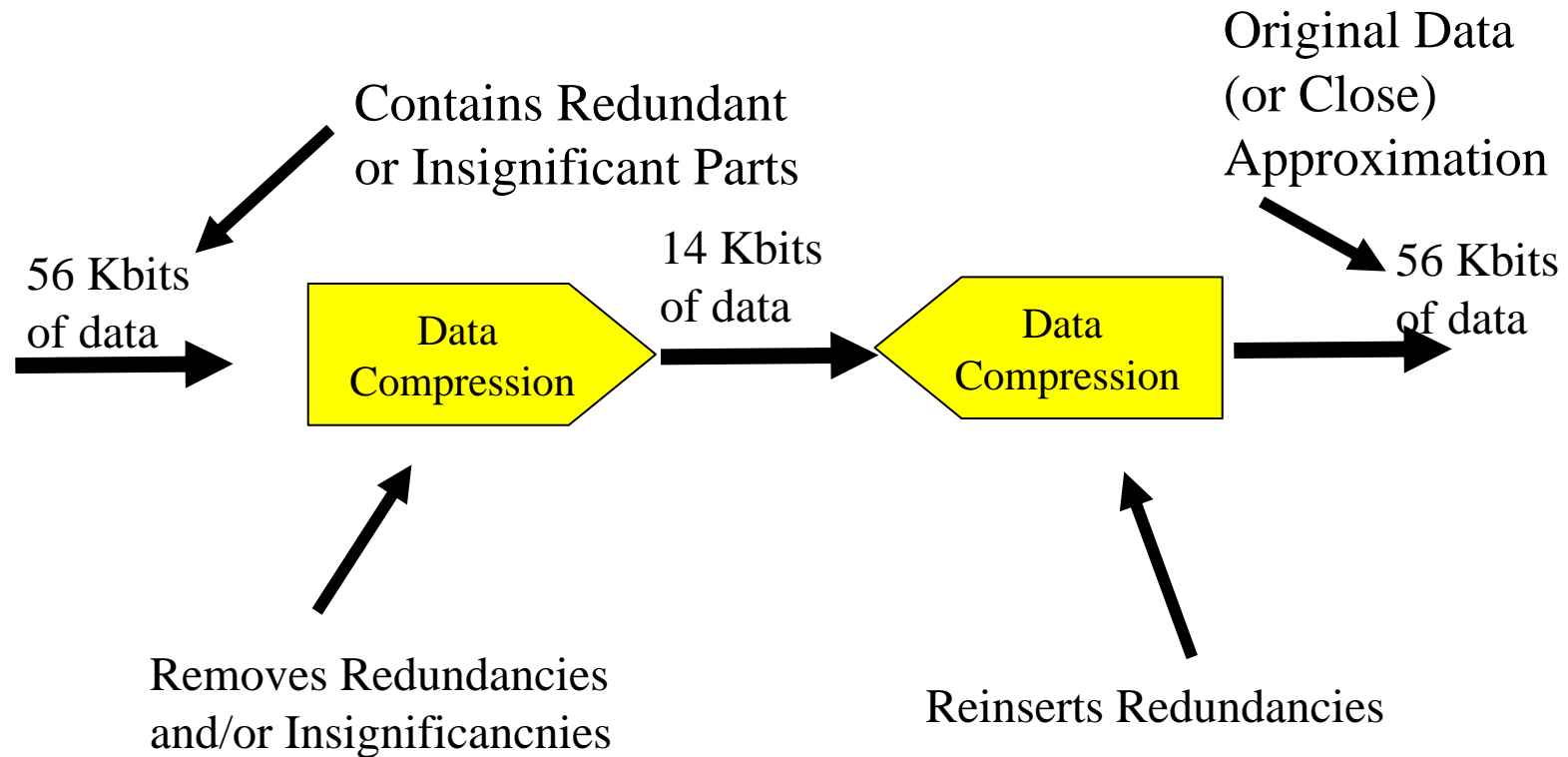
Introduction

&

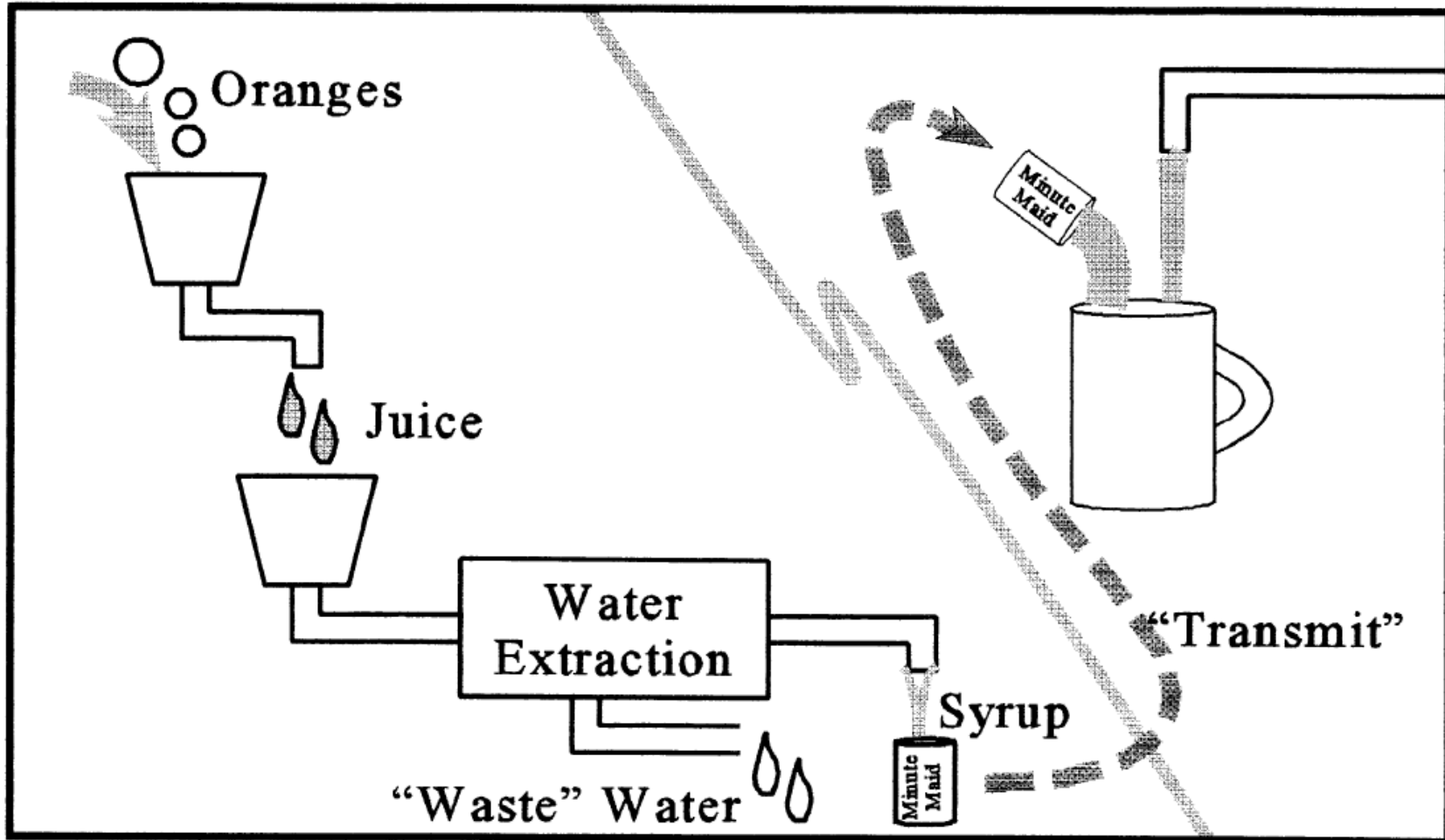
Lossless Example

I. Introduction

What is Data Compression? The systematic removal and, later, reintroduction of redundancy from data, and possibly also the removal of insignificant data.



Similar to Orange Juice Concentrate:



II. Lossless Vs. Lossy

There are two broad categories of data compression:

- Lossless Compression
- Lossy Compression

Lossless Compression: Original digital data can be exactly reconstructed from the compressed data

- Generally applied to:
 - text files (documents, source code, etc.)
 - program files
- Can achieve modest compression ratios:
 - 2:1 to 8:1 depending on type of content of file
- Examples:
 - UNIX “compress” command
 - PKZIP
 - Ram Doublers, Stacker, etc.

Lossy Compression: Original digital data can only be approximately reconstructed from the compressed data

- Generally applied to:

- Signals (speech, music, sensor signals, etc.)
- Images (digitized photos, digitized video)

- Can achieve large compression ratios:

- up to (many tens): 1 depending on type & content to data

- Examples:

- JPEG (for digitized photos)
- MPEG (for digitized video)
- Speech compression used in digital cell phones
- Digital answering machines
- **MP3 (for music)**

Lossless Example to Motivate Needed Theory

Consider a source with alphabet $\{A, B, C, D\}$

Suppose this is a typical source sequence:

DAABC AAABADBA AABBA AACDBC

The # of occurrences of symbols in this sequence are:

$$N_A = 12 \quad N_B = 6 \quad N_C = 3 \quad N_D = 3$$

$$\underbrace{\text{Total \# of symbols in sequence}}_{\Delta N} = 24$$

Say we have a binary code for the symbols such that

length in bits of code word A = l_A

length in bits of code word B = l_B

length in bits of code word D = l_D

$$\begin{aligned} \underbrace{\text{Total \# of bits to code sequence}}_{\Delta B_T} &= N_A l_A + N_B l_B + N_C l_C + N_D l_D \\ &= 12l_A + 6l_B + 3l_C + 3l_D \end{aligned}$$

A useful measure of lossless compression = “Avg” # Bits/symbols

If we “avg over the data”: ← “Data Analysis Average”

$$\begin{aligned}\text{Avg \# Bits/symbol} &= \frac{B_T}{N} \\ &= \frac{N_A}{N} \ell_A + \frac{N_B}{N} \ell_B + \frac{N_C}{N} \ell_C + \frac{N_D}{N} \ell_D\end{aligned}$$

^{ℓ_A} “Freq. of Occurance”

$$= 0.5 \ell_A + 0.25 \ell_B$$

⇒ Want ℓ_A to be smallest code length
then ℓ_B
then ℓ_C and ℓ_D

⇒ Most likely symbols get fewer bits

Prob. Theory Says:

$$\lim_{N \rightarrow \infty} \frac{N_i}{N} = P_i$$

Prob. of *i*th event

⇒ Base on our theory on probability

⇒ Avg. # Bits/Symbol uses probability average (not data average)

So, to study lossless theory we need to “specify” a probability model for the source:

$$P_A, P_B, P_C, P_D$$

⇒ Avg Bits/symbol = $E\{\ell\}$

↑ RV representing
code word length
of coded random symbols

$$= P_A \ell_A + P_B \ell_B + P_C \ell_C + P_D \ell_D$$

So, for this example we want l_A smallest, then l_B , then $l_C = l_D$

So try: A = 0 B = 1 C = 01 D = 10

Code #1

Does that work? NO!

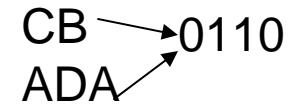
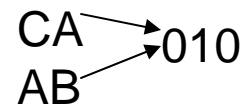
Can't decode sequences uniquely:

Original:	D A A B C ...
Coded:	1 0 0 0 0
Decoded:	$\left. \begin{array}{l} BAAABA... \\ DACAB... \\ etc. \end{array} \right\} \text{No unique decoding}$

How about: A = 0 B = 10 C = 01 D = 11

Code #2

No, because



How about: A = 0 B = 10 C = 11 D = xyz

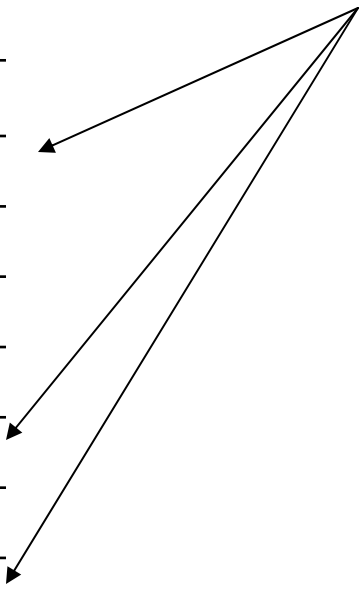
Code #3

Can't be 00
or 01
Seen above

D Must use 3 bits

xyz	
000	AAA
001	
010	AB
011	AC
100	BA
101	
110	CA
111	

possibilities



But.... No!

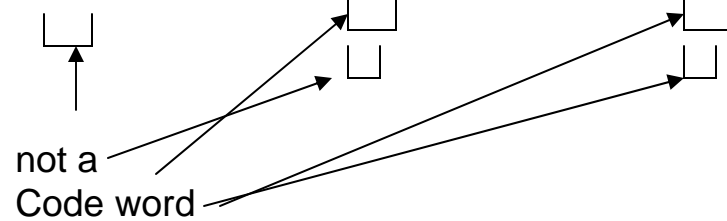
If $D = 001 \Rightarrow$ $\begin{matrix} DA \\ AAB \end{matrix} \rightarrow 0010$

If $D = 101 \Rightarrow$ $\begin{matrix} DA \\ BB \end{matrix} \rightarrow 1010$

If $D = 111 \Rightarrow$ $\begin{matrix} DA \\ CB \end{matrix} \rightarrow 1110$

How about: $A = 0$ $B = 10$ $C = 110$ $D = 111$

Code #4



Note this property: No other codeword is a prefix

Yes, this works! <Verify!>

So optimal lossless coding is a constrained optimization problem

minimize Avg. Bits/symbol
subject to decodability

Code #	Avg. bits/symbol
1	1.25 bits
2	1.5 bits
3	1.625 bits
4	1.75 bits

Not Decodable

Decodable

Compare to 2 bits/symbol if we used std. binary coding

Original stream using std. binary coding uses $2 \times 24 = 48$ bits

Coded stream used $1.75 \times 24 = 42$ bits

Compression!

To study this → Need some theory

Prob. Theory

+

Decodability

Information Theory

(at least part of it)

Role of Info Theory Here:

- Determine conditions need for decodability
- Determine lower bounds for avg. bits/symbol
- Provide understanding of how structure of the Source prob. Model impacts lower bounds
- Provide basis/structure on which practical codes can be built