

HM S

## Compression

28 November 2022

Sebastian Wild

## **Learning Outcomes**

- Understand the necessity for encodings and know ASCII and UTF-8 character encodings.
- 2. Understand (qualitatively) the *limits of compressibility*.
- 3. Know and understand the algorithms (encoding and decoding) for *Huffman codes*, *RLE*, *Elias codes*, *LZW*, *MTF*, and *BWT*, including their *properties* like running time complexity.
- **4.** Select and *adapt* (slightly) a *compression* pipeline for specific type of data.

Unit 7: Compression



#### **Outline**

# **7** Compression

- 7.1 Context
- 7.2 Character Encodings
- 7.3 Huffman Codes
- 7.4 Entropy
- 7.5 Run-Length Encoding
- 7.6 Lempel-Ziv-Welch
- 7.7 Lempel-Ziv-Welch Decoding
- 7.8 Move-to-Front Transformation
- 7.9 Burrows-Wheeler Transform
- 7.10 Inverse BWT

# 7.1 Context

#### Overview

- ► Unit 4–6: How to *work* with strings
  - finding substrings
  - finding approximate matches
  - finding repeated parts
  - ▶ ...
  - ► assumed character array (random access)!
- ▶ Unit 7–8: How to *store/transmit* strings
  - computer memory: must be binary
  - how to compress strings (save space)
  - ▶ how to robustly transmit over noisy channels → Unit 8

#### **Clicker Question**



What compression methods do you know?



→ sli.do/comp526

#### **Terminology**

- ▶ **source text:** string  $S \in \Sigma_S^*$  to be stored / transmitted  $\Sigma_S$  is some alphabet
- ▶ **coded text:** encoded data  $C \in \Sigma_C^*$  that is actually stored / transmitted usually use  $\Sigma_C = \{0, 1\}$
- **encoding:** algorithm mapping source texts to coded texts  $\leq \sim \sim \subset$

## **Terminology**

- ▶ **source text:** string  $S \in \Sigma_S^*$  to be stored / transmitted  $\Sigma_S$  is some alphabet
- ▶ **coded text:** encoded data  $C \in \Sigma_C^*$  that is actually stored / transmitted usually use  $\Sigma_C = \{0, 1\}$
- ▶ encoding: algorithm mapping source texts to coded texts
- ▶ decoding: algorithm mapping coded texts back to original source text
- ► Lossy vs. Lossless
  - lossy compression can only decode approximately; the exact source text S is lost
  - ▶ **lossless compression** always decodes *S* exactly
- ► For media files, lossy, logical compression is useful (e. g. JPEG, MPEG)
- ► We will concentrate on *lossless* compression algorithms. These techniques can be used for any application.

## What is a good encoding scheme?

- ▶ Depending on the application, goals can be
  - ► efficiency of encoding/decoding
  - ► resilience to errors/noise in transmission → Uwif 🖇
  - security (encryption)
  - ▶ integrity (detect modifications made by third parties)
  - ▶ size

## What is a good encoding scheme?

- ▶ Depending on the application, goals can be
  - efficiency of encoding/decoding
  - ▶ resilience to errors/noise in transmission
  - security (encryption)
  - ▶ integrity (detect modifications made by third parties)
  - size
  - Focus in this unit: **size** of coded text

    Encoding schemes that (try to) minimize the size of coded texts perform *data compression*.
- ► We will measure the <u>compression ratio:</u>  $\frac{|C| \cdot \lg |\Sigma_C|}{|S| \cdot \lg |\Sigma_S|} \stackrel{\Sigma_C = \{0,1\}}{=} \frac{|C|}{|S| \cdot \lg |\Sigma_S|}$ 
  - < 1 means successful compression
  - = 1 means no compression
  - > 1 means "compression" made it bigger!? (yes, that happens ...)

#### **Clicker Question**



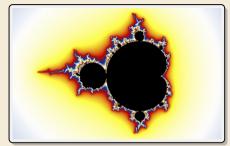
Do you know what uncomputable problems (halting problem, Post's correspondence problem, . . . ) are?

- A Sure, I could explain what it is.
- B Heard that in a lecture, but don't quite remember
- No, never heard of it



→ sli.do/comp526

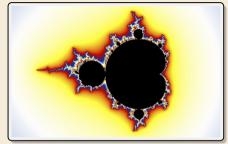
Is this image compressible?



Is this image compressible?

visualization of Mandelbrot set

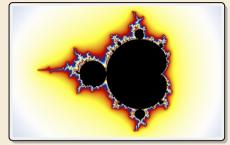
- ► Clearly a complex shape!
- ▶ Will not compress (too) well using, say, PNG.
- ▶ but:
  - completely defined by mathematical formula
  - → can be generated by a very small program!



Is this image compressible?

visualization of Mandelbrot set

- ► Clearly a complex shape!
- ▶ Will not compress (too) well using, say, PNG.
- ▶ but:
  - completely defined by mathematical formula
  - $\leadsto$  can be generated by a very small program!



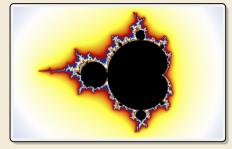
#### *→* Kolmogorov complexity

- ightharpoonup C = any program that outputs S
  - self-extracting archives!
- ► Kolmogorov complexity = length of smallest such program

*Is this image compressible?* 

visualization of Mandelbrot set

- Clearly a complex shape!
- ▶ Will not compress (too) well using, say, PNG.
- ▶ but:
  - completely defined by mathematical formula
  - → can be generated by a very small program!



#### *→* Kolmogorov complexity

- ightharpoonup C = any program that outputs S
  - self-extracting archives!
- ► Kolmogorov complexity = length of smallest such program
- ▶ **Problem:** finding smallest such program is *uncomputable*.
- → No optimal encoding algorithm is possible!
- → must be inventive to get efficient methods

#### What makes data compressible?

- ► Lossless compression methods mainly exploit two types of redundancies in source texts:
  - uneven character frequencies some characters occur more often than others → Part I
  - 2. repetitive texts
    different parts in the text are (almost) identical → Part II

#### What makes data compressible?

- Lossless compression methods mainly exploit two types of redundancies in source texts:
  - uneven character frequencies some characters occur more often than others → Part I
  - 2. repetitive texts different parts in the text are (almost) identical → Part II



There is no such thing as a free lunch!

Not *everything* is compressible ( $\rightarrow$  tutorials)

→ focus on versatile methods that often work

# Part I

Exploiting character frequencies

7.2 Character Encodings

#### **Character encodings**

- ► Simplest form of encoding: Encode each source character individually
- $\rightsquigarrow$  encoding function  $E: \Sigma_S \to \Sigma_C^*$ 
  - typically,  $|\Sigma_S| \gg |\Sigma_C|$ , so need several bits per character
  - ▶ for  $c \in \Sigma_S$ , we call E(c) the *codeword* of c
- ▶ **fixed-length code:** |E(c)| is the same for all  $c \in \Sigma_C$
- ▶ variable-length code: not all codewords of same length

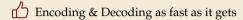
#### **Fixed-length codes**

- ▶ fixed-length codes are the simplest type of character encodings
- Example: ASCII (American Standard Code for Information Interchange, 1963)

```
0000000 NUL
               0010000 DLE
                              0100000
                                            0110000 0
                                                         1000000 a
                                                                       1010000 P
                                                                                    1100000 '
                                                                                                 1110000 p
0000001 SOH
               0010001 DC1
                              0100001 !
                                            0110001 1
                                                         1000001 A
                                                                       1010001 0
                                                                                    1100001 a
                                                                                                 1110001 q
0000010 STX
               0010010 DC2
                              0100010 "
                                            0110010 2
                                                         1000010 B
                                                                       1010010 R
                                                                                    1100010 b
                                                                                                 1110010 r
0000011 ETX
               0010011 DC3
                              0100011 #
                                            0110011 3
                                                         1000011 C
                                                                      1010011 S
                                                                                   1100011 c
                                                                                                 1110011 s
0000100 EOT
               0010100 DC4
                              0100100 $
                                            0110100 4
                                                         1000100 D
                                                                       1010100 T
                                                                                   1100100 d
                                                                                                 1110100 t
0000101 ENO
               0010101 NAK
                              0100101 %
                                            0110101 5
                                                         1000101 E
                                                                       1010101 U
                                                                                    1100101 e
                                                                                                 1110101 u
0000110 ACK
               0010110 SYN
                              0100110 &
                                            0110110 6
                                                         1000110 F
                                                                      1010110 V
                                                                                   1100110 f
                                                                                                 1110110 v
0000111 BEL
               0010111 ETB
                              0100111 '
                                            0110111 7
                                                         1000111 G
                                                                       1010111 W
                                                                                    1100111 a
                                                                                                 1110111 w
0001000 BS
               0011000 CAN
                              0101000 (
                                            0111000 8
                                                         1001000 H
                                                                       1011000 X
                                                                                    1101000 h
                                                                                                 1111000 ×
0001001 HT
               0011001 EM
                              0101001 )
                                            0111001 9
                                                         1001001 I
                                                                      1011001 Y
                                                                                   1101001 i
                                                                                                 1111001 v
0001010 LF
               0011010 SUB
                              0101010 *
                                            0111010 :
                                                         1001010 J
                                                                      1011010 Z
                                                                                   1101010 i
                                                                                                 1111010 z
               0011011 ESC
                                            0111011 :
0001011 VT
                              0101011 +
                                                         1001011 K
                                                                       1011011 [
                                                                                    1101011 k
                                                                                                 1111011 {
0001100 FF
               0011100 FS
                              0101100 ,
                                            0111100 <
                                                         1001100 L
                                                                       1011100 \
                                                                                   1101100 l
                                                                                                 1111100
0001101 CR
               0011101 GS
                              0101101 -
                                            0111101 =
                                                         1001101 M
                                                                       1011101 1
                                                                                   1101101 m
                                                                                                 1111101 }
0001110 SO
               0011110 RS
                              0101110 .
                                            0111110 >
                                                         1001110 N
                                                                       1011110 ^
                                                                                    1101110 n
                                                                                                 1111110 ~
0001111 SI
               0011111 US
                              0101111 /
                                            0111111 ?
                                                         1001111 0
                                                                       1011111
                                                                                    1101111 o
                                                                                                 1111111 DEL
```

- ▶ 7 bit per character
- ▶ just enough for English letters and a few symbols (plus control characters)

#### Fixed-length codes – Discussion

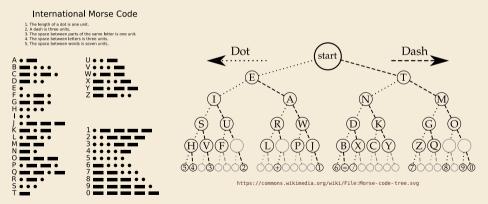


Unless all characters equally likely, it wastes a lot of space

inflexible (how to support adding a new character?)

#### Variable-length codes

- ▶ to gain more flexibility, have to allow different lengths for codewords
- ▶ actually an old idea: Morse Code



https://commons.wikimedia.org/wiki/File: International Morse Code.svg

#### **Clicker Question**

How many characters are there in the alphabet of the coded text in Morse Code, i. e., what is  $|\Sigma_C|$ ?



**A** ) 1

(E) 20

**B** ) 2

**F** 3

**c** 3

**G** 256

**D**) 4



→ sli.do/comp526

#### **Clicker Question**

How many characters are there in the alphabet of the coded text in Morse Code, i. e., what is  $|\Sigma_C|$ ?



A) 1

(E) <del>2(</del>

B) <del>2</del>

F) <del>3(</del>

) 3 ✓

G 256

 $\left[ \mathsf{D} \right] 4$ 



→ sli.do/comp526

#### **Variable-length codes – UTF-8**

▶ Modern example: UTF-8 encoding of Unicode:

default encoding for text-files, XML, HTML since 2009

- ► Encodes any Unicode character (137 994 as of May 2019, and counting)
- ▶ uses 1–4 bytes (codeword lengths: 8, 16, 24, or 32 bits)
- Every ASCII character is encoded in 1 byte with leading bit 0, followed by the 7 bits for ASCII
- Non-ASCII charactters start with 1–4 1s indicating the total number of bytes, followed by a 0 and 3–5 bits.

The remaining bytes each start with 10 followed by 6 bits.

Char. number range	UTF-8 octet sequence					
(hexadecimal)	(binary)					
0000 0000 - 0000 007F	0xxxxxx					
0000 0080 - 0000 07FF	110xxxxx 10xxxxxx					
0000 0800 - 0000 FFFF	1110xxxx 10xxxxxx 10xxxxxx					
0001 0000 - 0010 FFFF	11110xxx 10xxxxxx 10xxxxxx 10xxxxxx					

For English text, most characters use only 8 bit, but we can include any Unicode character, as well.

## Pitfall in variable-length codes

## Pitfall in variable-length codes

- **9**  $C = 1100100100 \text{ decodes both to banana and to bass: } \frac{110}{b} \frac{0}{a} \frac{100}{s} \frac{100}{s}$
- → not a valid code . . . (cannot tolerate ambiguity)
  but how should we have known?

## Pitfall in variable-length codes

- $rac{1}{2}$  C = 1100100100 decodes **both** to banana and to bass:  $\frac{110}{b} \frac{0100100}{a} \frac{100}{s}$
- → not a valid code . . . (cannot tolerate ambiguity)
  but how should we have known?
- E(n) = 10 is a (proper) **prefix** of E(s) = 100
  - → Leaves decoder wondering whether to stop after reading 10 or continue!
  - → Require a *prefix-free* code: No codeword is a prefix of another.

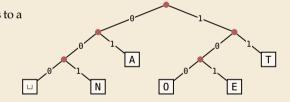
    prefix-free ⇒ instantaneously decodable ⇒ uniquely decodable

#### **Code tries**

- ► From now on only consider prefix-free codes E: E(c) is not a prefix of E(c') for any  $c, c' \in \Sigma_S$ .

Any prefix-free code corresponds to a *(code) trie* (trie of codewords) with characters of  $\Sigma_S$  at **leaves**.

no need for end-of-string symbols \$ here (already prefix-free!)



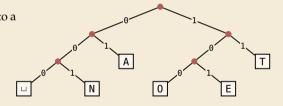
- ► Encode AN\_ANT
- ▶ Decode 111000001010111

#### **Code tries**

- ► From now on only consider prefix-free codes E: E(c) is not a prefix of E(c') for any  $c, c' \in \Sigma_S$ .

Any prefix-free code corresponds to a *(code) trie* (trie of codewords) with characters of  $\Sigma_S$  at **leaves**.

no need for end-of-string symbols \$ here (already prefix-free!)



- ► Encode AN\_ANT → 010010000100111
- ► Decode 1110000010101111 → T0\_EAT

#### Who decodes the decoder?

- ▶ Depending on the application, we have to **store/transmit** the **used code**!
- ► We distinguish:
  - ▶ fixed coding: code agreed upon in advance, not transmitted (e. g., Morse, UTF-8)
  - ► static coding: code depends on message, but stays same for entire message; it must be transmitted (e. g., Huffman codes → next)
  - ▶ adaptive coding: code depends on message and changes during encoding; implicitly stored withing the message (e. g., LZW → below)

## 7.3 Huffman Codes

#### **Character frequencies**

- ▶ Goal: Find character encoding that produces short coded text
- ▶ Convention here: fix  $\Sigma_C = \{0, 1\}$  (binary codes), abbreviate  $\Sigma = \Sigma_S$ ,
- ▶ **Observation:** Some letters occur more often than others.

#### **Typical English prose:**

e	12.70%		d	4.25%	_	p	1.93%	•
t	9.06%		1	4.03%	_	b	1.49%	•
a	8.17%		c	2.78%		$\mathbf{v}$	0.98%	•
О	7.51%	_	u	2.76%		k	0.77%	
i	6.97%		m	2.41%		j	0.15%	1
n	6.75%		$\mathbf{w}$	2.36%	-	x	0.15%	1
s	6.33%		f	2.23%		q	0.10%	1
h	6.09%	_	g	2.02%		$\mathbf{z}$	0.07%	1
r	5.99%	_	y	1.97%				

→ Want shorter codes for more frequent characters!

### **Huffman** coding

e.g. frequencies / probabilities

- ▶ **Given:**  $\Sigma$  and weights  $w: \Sigma \to \mathbb{R}_{\geq 0}$
- ▶ **Goal:** prefix-free code E (= code trie) for  $\Sigma$  that minimizes coded text length

i. e., a code trie minimizing 
$$\sum_{c \in \Sigma} w(c) \cdot |E(c)|$$

## **Huffman** coding

#### e.g. frequencies / probabilities

- ▶ **Given:**  $\Sigma$  and weights  $w: \Sigma \to \mathbb{R}_{\geq 0}$
- ▶ **Goal:** prefix-free code E (= code trie) for  $\Sigma$  that minimizes coded text length

i. e., a code trie minimizing 
$$\sum_{c \in \Sigma} w(c) \cdot |E(c)|$$

- ▶ If we use w(c) = #occurrences of c in S, this is the character encoding with smallest possible |C|
  - → best possible character-wise encoding

▶ Quite ambitious! *Is this efficiently possible?* 

# Huffman's algorithm

► Actually, yes! A greedy/myopic approach succeeds here.

#### Huffman's algorithm:

- 1. Find two characters a, b with lowest weights.
  - ▶ We will encode them with the same prefix, plus one distinguishing bit, i. e., E(a) = u0 and E(b) = u1 for a bitstring  $u \in \{0, 1\}^*$  (u to be determined)
- 2. (Conceptually) replace a and b by a single character "ab" with w(ab) = w(a) + w(b).
- 3. Recursively apply Huffman's algorithm on the smaller alphabet. This in particular determines  $u = E(\Box b)$ .

# Huffman's algorithm

► Actually, yes! A greedy/myopic approach succeeds here.

#### Huffman's algorithm:

- 1. Find two characters a, b with lowest weights.
  - ▶ We will encode them with the same prefix, plus one distinguishing bit, i. e., E(a) = u0 and E(b) = u1 for a bitstring  $u \in \{0, 1\}^*$  (u to be determined)
- 2. (Conceptually) replace a and b by a single character "ab" with w(ab) = w(a) + w(b).
- 3. Recursively apply Huffman's algorithm on the smaller alphabet. This in particular determines  $u = E(\blacksquare)$ .
- efficient implementation using a (min-oriented) *priority queue* 
  - start by inserting all characters with their weight as key
  - ▶ step 1 uses two deleteMin calls
  - ▶ step 2 inserts a new character with the sum of old weights as key

- ► Example text: S = LOSSLESS  $\longrightarrow \Sigma_S = \{E, L, 0, S\}$
- ightharpoonup Character frequencies: E:1, L:2, 0:1, S:4

1

2

1

4

E

L

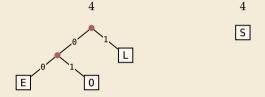
0

S

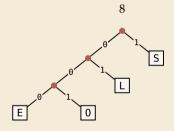
- ► Example text: S = LOSSLESS  $\longrightarrow \Sigma_S = \{E, L, 0, S\}$
- ightharpoonup Character frequencies: E:1, L:2, 0:1, S:4



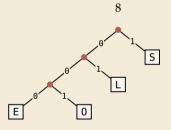
- ► Example text: S = LOSSLESS  $\longrightarrow \Sigma_S = \{E, L, 0, S\}$
- ightharpoonup Character frequencies: E:1, L:2, 0:1, S:4



- ► Example text: S = LOSSLESS  $\longrightarrow \Sigma_S = \{E, L, 0, S\}$
- ightharpoonup Character frequencies: E:1, L:2, 0:1, S:4

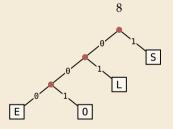


- ► Example text: S = LOSSLESS  $\longrightarrow \Sigma_S = \{E, L, 0, S\}$
- ► Character frequencies: E:1, L:2, 0:1, S:4



→ *Huffman tree* (code trie for Huffman code)

- ► Example text: S = LOSSLESS  $\longrightarrow \Sigma_S = \{E, L, 0, S\}$
- ightharpoonup Character frequencies: E:1, L:2, 0:1, S:4



→ *Huffman tree* (code trie for Huffman code)

LOSSLESS  $\rightarrow$  01001110100011 compression ratio:  $\frac{14}{8 \cdot \log 4} = \frac{14}{16} \approx 88\%$ 

# Huffman tree – tie breaking

- ► The above procedure is ambiguous:
  - which characters to choose when weights are equal?
  - ▶ which subtree goes left, which goes right?
- ► For COMP 526: always use the following rule:
  - To break ties when selecting the two characters, first use the smallest letter according to the alphabetical order, or the tree containing the smallest alphabetical letter.
  - 2. When combining two trees of different values, place the lower-valued tree on the left (corresponding to a 0-bit).
  - When combining trees of equal value, place the one containing the smallest letter to the left.

# **Encoding with Huffman code**

- ► The overall encoding procedure is as follows:
  - ▶ Pass 1: Count character frequencies in *S*
  - ► Construct Huffman code *E* (as above)
  - ► Store the Huffman code in *C* (details omitted)
  - ▶ Pass 2: Encode each character in *S* using *E* and append result to *C*
- Decoding works as follows:
  - ▶ Decode the Huffman code *E* from *C*. (details omitted)
  - ▶ Decode *S* character by character from *C* using the code trie.
- ► Note: Decoding is much simpler/faster!

# **Huffman code – Optimality**

#### Theorem 7.1 (Optimality of Huffman's Algorithm)

Given  $\Sigma$  and  $w: \Sigma \to \mathbb{R}_{\geq 0}$ , Huffman's Algorithm computes codewords  $E: \Sigma \to \{0,1\}^*$  with minimal expected codeword length  $\ell(E) = \sum_{c \in \Sigma} w(c) \cdot |E(c)|$  among all prefix-free codes for  $\Sigma$ .

# **Huffman code – Optimality**

#### Theorem 7.1 (Optimality of Huffman's Algorithm)

Given  $\Sigma$  and  $w: \Sigma \to \mathbb{R}_{\geq 0}$ , Huffman's Algorithm computes codewords  $E: \Sigma \to \{0,1\}^*$  with minimal expected codeword length  $\ell(E) = \sum_{c \in \Sigma} w(c) \cdot |E(c)|$  among all prefix-free codes for  $\Sigma$ .

*Proof sketch:* by induction over  $\sigma = |\Sigma|$ 

- ightharpoonup Given any optimal prefix-free code  $E^*$  (as its code trie).
- ▶ code trie  $\rightarrow$  ∃ two sibling leaves x, y at largest depth D
- ▶ swap characters in leaves to have two lowest-weight characters a, b in x, y (that can only make  $\ell$  smaller, so still optimal)
- ▶ any optimal code for  $\Sigma' = \Sigma \setminus \{a, b\} \cup \{ab\}$  yields optimal code for  $\Sigma$  by replacing leaf ab by internal node with children a and b.
- $\leadsto$  recursive call yields optimal code for  $\Sigma'$  by inductive hypothesis, so Huffman's algorithm finds optimal code for  $\Sigma$ .

# 7.4 Entropy

#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

- entropy is a measure of information content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

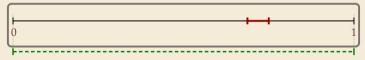
- entropy is a measure of information content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

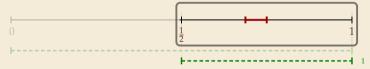
- entropy is a measure of information content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

- entropy is a measure of information content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

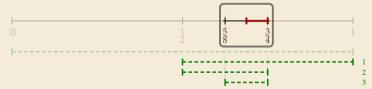
- entropy is a **measure** of **information** content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

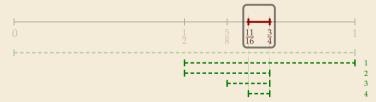
- entropy is a **measure** of **information** content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

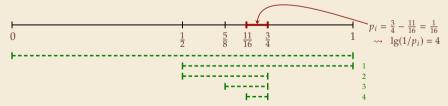
- entropy is a **measure** of **information** content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

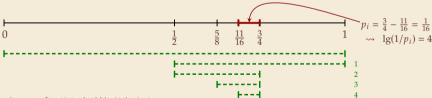
- entropy is a measure of information content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



#### **Definition 7.2 (Entropy)**

$$\mathcal{H}(p_1,\ldots,p_n) = -\sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i}\right)$$

- entropy is a measure of information content of a distribution
  - ▶ "20 *Questions on* [0, 1)": Land inside my interval by halving.



- $\rightarrow$  Need to cut [0, 1) in half  $\lg(1/p_i)$  times
- more precisely: the expected number of bits (Yes/No questions) required to nail down the random value

# **Entropy and Huffman codes**

▶ would ideally encode value i using  $\lg(1/p_i)$  bits not always possible; cannot use codeword of 1.5 bits . . .

# **Entropy and Huffman codes**

would ideally encode value i using  $\lg(1/p_i)$  bits not always possible; cannot use codeword of 1.5 bits . . . but:

#### Theorem 7.3 (Entropy bounds for Huffman codes)

For any  $\Sigma = \{a_1, \dots, a_\sigma\}$  and  $w : \Sigma \to \mathbb{R}_{>0}$  and its Huffman code E, we have

$$\mathcal{H} \leq \ell(E) \leq \mathcal{H} + 1$$
 where  $\mathcal{H} = \mathcal{H}\left(\frac{w(a_1)}{W}, \dots, \frac{w(a_\sigma)}{W}\right)$  and  $W = w(a_1) + \dots + w(a_\sigma)$ .

# **Entropy and Huffman codes**

would ideally encode value i using  $\lg(1/p_i)$  bits not always possible; cannot use codeword of 1.5 bits . . . but:

#### Theorem 7.3 (Entropy bounds for Huffman codes)

For any  $\Sigma = \{a_1, \dots, a_\sigma\}$  and  $w : \Sigma \to \mathbb{R}_{>0}$  and its Huffman code E, we have

$$\mathcal{H} \leq \ell(E) \leq \mathcal{H} + 1$$
 where  $\mathcal{H} = \mathcal{H}\left(\frac{w(a_1)}{W}, \dots, \frac{w(a_\sigma)}{W}\right)$  and  $W = w(a_1) + \dots + w(a_\sigma)$ .

#### Proof sketch:

▶  $\ell(E) \ge \mathcal{H}$ Any prefix-free code E induces weights  $q_i = 2^{-|E(a_i)|}$ . By Kraft's Inequality, we have  $q_1 + \cdots + q_{\sigma} \le 1$ . Hence we can apply Gibb's Inequality to get

$$\mathcal{H} = \sum_{i=1}^{\sigma} p_i \lg \left(\frac{1}{p_i}\right) \leq \sum_{i=1}^{\sigma} p_i \lg \left(\frac{1}{q_i}\right) = \ell(E).$$

# **Entropy and Huffman codes [2]**

*Proof sketch (continued):* 

 $\blacktriangleright$   $\ell(E) \leq \mathcal{H} + 1$ 

Set 
$$q_i = 2^{-\lceil \lg(1/p_i) \rceil}$$
. We have  $\sum_{i=1}^{\sigma} p_i \lg\left(\frac{1}{q_i}\right) = \sum_{i=1}^{\sigma} p_i \lceil \lg(1/p_i) \rceil \le \mathcal{H} + 1$ .

We construct a code E' for  $\Sigma$  with  $|E'(a_i)| \le \lg(1/q_i)$  as follows; w.l.o.g. assume  $q_1 \le q_2 \le \cdots \le q_\sigma$ 

- ► If  $\sigma = 2$ , E' uses a single bit each. Here,  $q_i \le 1/2$ , so  $\lg(1/q_i) \ge 1 = |E'(a_i)| \checkmark$
- ▶ If  $\sigma \ge 3$ , we merge  $a_1$  and  $a_2$  to  $\overline{a_1a_2}$ , assign it weight  $2q_2$  and recurse. If  $q_1 = q_2$ , this is like Huffman; otherwise,  $q_1$  is a unique smallest value and  $q_2 + q_2 + \cdots + q_{\sigma} \le 1$ .

By the inductive hypothesis, we have 
$$|E'(\overline{a_1a_2})| \le \lg\left(\frac{1}{2q_2}\right) = \lg\left(\frac{1}{q_2}\right) - 1$$
.  
By construction,  $|E'(a_1)| = |E'(a_2)| = |E'(\overline{a_1a_2})| + 1$ , so  $|E'(a_1)| \le \lg\left(\frac{1}{q_1}\right)$  and  $|E'(a_2)| \le \lg\left(\frac{1}{q_2}\right)$ .

By optimality of 
$$E$$
, we have  $\ell(E) \leq \ell(E') \leq \sum_{i=1}^{\sigma} p_i \lg\left(\frac{1}{q_i}\right) \leq \mathcal{H} + 1$ .

# **Clicker Question**

When does Huffman coding yield more efficient compression than a fixed-length character encoding?



- **A** always
- **B** when  $\mathcal{H} \approx \lg(\sigma)$
- **C** when  $\mathcal{H} < \lg(\sigma)$
- **D** when  $\mathcal{H} < \lg(\sigma) 1$
- **E** when  $\mathcal{H} \approx 1$



→ sli.do/comp526

# **Clicker Question**

When does Huffman coding yield more efficient compression than a fixed-length character encoding?



- A always 🗸
- B when  $\mathcal{H} \simeq \lg(\sigma)$

- E when √ ~ 1



→ sli.do/comp526

#### **Huffman coding – Discussion**

- ▶ running time complexity:  $O(\sigma \log \sigma)$  to construct code
  - ▶ build PQ +  $\sigma$  · (2 deleteMins and 1 insert)
  - ightharpoonup can do  $\Theta(\sigma)$  time when characters already sorted by weight
  - time for encoding text (after Huffman code done): O(n + |C|)
- ▶ many variations in use (tie-breaking rules, estimated frequencies, adaptive encoding, ...)

# **Huffman coding – Discussion**

- ▶ running time complexity:  $O(\sigma \log \sigma)$  to construct code
  - ▶ build PQ +  $\sigma$  · (2 deleteMins and 1 insert)
  - can do  $\Theta(\sigma)$  time when characters already sorted by weight
  - time for encoding text (after Huffman code done): O(n + |C|)
- ▶ many variations in use (tie-breaking rules, estimated frequencies, adaptive encoding, . . .)
- optimal prefix-free character encoding
- very fast decoding
- - one-pass variants possible, but more complicated
- $\hfill \bigcap$  have to store code alongside with coded text

# Part II

Compressing repetitive texts

# **Beyond Character Encoding**

Many "natural" texts show repetitive redundancy

All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy.

- ► character-by-character encoding will **not** capture such repetitions
  - → Huffman won't compression this very much

# **Beyond Character Encoding**

Many "natural" texts show repetitive redundancy

All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy. All work and no play makes Jack a dull boy.

- ► character-by-character encoding will **not** capture such repetitions
  - → Huffman won't compression this very much
- $\rightarrow$  Have to encode whole *phrases* of S by a single codeword

# 7.5 Run-Length Encoding

# **Run-Length encoding**

▶ simplest form of repetition: *runs* of characters



same character repeated

- ▶ here: only consider  $\Sigma_S = \{0, 1\}$  (work on a binary representation)
  - ► can be extended for larger alphabets

#### **Run-Length encoding**

▶ simplest form of repetition: *runs* of characters



same character repeated

- here: only consider  $\Sigma_S = \{0, 1\}$  (work on a binary representation)
  - can be extended for larger alphabets
- $\rightsquigarrow$  run-length encoding (RLE):

use runs as phrases: S = 00000 111 0000

#### **Run-Length encoding**

▶ simplest form of repetition: *runs* of characters

```
0001011001000001111110000000000011111000
00111111111000111111111100000011111111000
00111111111000000000001110011111111111000
001110111110000000001110001111100111100
000000000111000000011100001110000001110
000000000111000000011000001110000001100
000000000110000001100000011000001110
00000000011000001110000001110000001100
000000000111000111000000000110000001110
000000000110000111000000000111000011100
00110111111000111101110100001111111111000
000101100000001010011001000000100100000
```

same character repeated

- here: only consider  $\Sigma_S = \{0, 1\}$  (work on a binary representation)
  - can be extended for larger alphabets
- $\leadsto$  run-length encoding (RLE):

use runs as phrases: S = 00000 111 0000

- → We have to store
  - ▶ the first bit of *S* (either 0 or 1)
  - the length each each run
  - ▶ Note: don't have to store bit for later runs since they must alternate.
- ► Example becomes: 0, 5, 3, 4

#### **Run-Length encoding**

▶ simplest form of repetition: *runs* of characters

```
0001011001000001111110000000000011111000
00111111111000111111111100000011111111000
00111111111000000000001110011111111111000
001110111110000000001110001111100111100
000000000111000000011100001110000001110
000000000111000000011000001110000001100
00000000001100000011000000110000001110
00000000011000001110000001110000001100
000000000111000111000000000110000001110
000000000110000111000000000111000011100
00110111111000111101110100001111111111000
```

same character repeated

- ▶ here: only consider  $\Sigma_S = \{0, 1\}$  (work on a binary representation)
  - ► can be extended for larger alphabets
- $\rightsquigarrow$  run-length encoding (RLE):

use runs as phrases: S = 00000 111 0000

- → We have to store
  - ▶ the first bit of *S* (either 0 or 1)
  - ▶ the length each each run
  - ▶ Note: don't have to store bit for later runs since they must alternate.
- ► Example becomes: 0, 5, 3, 4
- **Question**: How to encode a run length k in binary? (k can be arbitrarily large!)

#### **Clicker Question**



How would you encode a string that can we arbitrarily long?



→ sli.do/comp526

- ▶ Need a *prefix-free encoding* for  $\mathbb{N} = \{1, 2, 3, \dots, \}$ 
  - ► must allow arbitrarily large integers
  - must know when to stop reading

- ▶ Need a *prefix-free encoding* for  $\mathbb{N} = \{1, 2, 3, \dots, \}$ 
  - must allow arbitrarily large integers
  - must know when to stop reading
- ► But that's simple! Just use *unary* encoding!

- ▶ Need a *prefix-free encoding* for  $\mathbb{N} = \{1, 2, 3, \dots, \}$ 
  - must allow arbitrarily large integers
  - must know when to stop reading
- ► But that's simple! Just use *unary* encoding!

- Much too long
  - (wasn't the whole point of RLE to get rid of long runs??)

- ▶ Need a *prefix-free encoding* for  $\mathbb{N} = \{1, 2, 3, \dots, \}$ 
  - must allow arbitrarily large integers
  - must know when to stop reading
- ► But that's simple! Just use *unary* encoding!

- Much too long
  - (wasn't the whole point of RLE to get rid of long runs??)
- ► Refinement: *Elias gamma code* 
  - ▶ Store the **length**  $\ell$  of the binary representation in **unary**
  - ► Followed by the binary digits themselves

- ▶ Need a *prefix-free encoding* for  $\mathbb{N} = \{1, 2, 3, \dots, \}$ 
  - must allow arbitrarily large integers
  - must know when to stop reading
- ► But that's simple! Just use *unary* encoding!

- Much too long
  - (wasn't the whole point of RLE to get rid of long runs??)
- ► Refinement: *Elias gamma code* 
  - ▶ Store the **length**  $\ell$  of the binary representation in **unary**
  - ► Followed by the binary digits themselves
  - ▶ little tricks:
    - ▶ always  $\ell \ge 1$ , so store  $\ell 1$  instead
    - ▶ binary representation always starts with 1 → don't need terminating 1 in unary
  - $\rightarrow$  Elias gamma code =  $\ell 1$  zeros, followed by binary representation

**Examples:** 
$$1 \mapsto 1$$
,  $3 \mapsto 011$ ,  $5 \mapsto 00101$ ,  $30 \mapsto 000011110$ 

#### **Clicker Question**



Decode the **first** number in Elias gamma code (at the beginning) of the following bitstream:

000110111011100110.



→ sli.do/comp526

► Encoding:

C = 1

► Decoding:

$$C = 00001101001001010$$

► Encoding:

► Decoding:

```
C = 00001101001001010
```

► Encoding:

► Decoding:

```
C = 00001101001001010
```

► Encoding:

► Decoding:

C = 00001101001001010

► Encoding:

► Decoding:

```
C = 00001101001001010
```

► Encoding:

► Decoding:

```
C = 00001101001001010
```

► Encoding:

C = 10011101010000101000001011

Compression ratio:  $26/41 \approx 63\%$ 

► Decoding:

C = 00001101001001010

► Encoding:

C = 10011101010000101000001011

Compression ratio:  $26/41 \approx 63\%$ 

► Decoding:

C = 00001101001001010

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio: 
$$26/41 \approx 63\%$$

$$C = 00001101001001010$$
  
 $b = 0$ 

$$S =$$

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio:  $26/41 \approx 63\%$ 

```
C = 00001101001001010
```

$$b = 0$$

$$\ell = 3 + 1$$

$$S =$$

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio:  $26/41 \approx 63\%$ 

► Decoding:

```
C = 0000 \frac{1101}{001001001010}
```

b = 0

 $\ell = 3 + 1$ 

k = 13

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio:  $26/41 \approx 63\%$ 

```
C = 00001101001001010

b = 1

\ell = 2 + 1

k = 1

k = 1
```

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio:  $26/41 \approx 63\%$ 

```
C = 00001101001001010

b = 1

\ell = 2 + 1

k = 4

S = 000000000000001111
```

► Encoding:

C = 10011101010000101000001011

Compression ratio:  $26/41 \approx 63\%$ 

```
C = 00001101001001010
b = 0
\ell = 0 + 1
k = 000000000000001111
```

► Encoding:

C = 10011101010000101000001011

Compression ratio:  $26/41 \approx 63\%$ 

► Decoding:

```
C = 0000110100100100
```

b = 0

 $\ell = 0 + 1$ 

k = 1

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio:  $26/41 \approx 63\%$ 

```
C = 00001101001001010 b = 1 \ell = 1 + 1 k = S = 00000000000011110
```

► Encoding:

```
C = 10011101010000101000001011
```

Compression ratio:  $26/41 \approx 63\%$ 

```
C = 00001101001001010

b = 1

\ell = 1 + 1

k = 2

S = 000000000000001111011
```

## **Run-length encoding – Discussion**

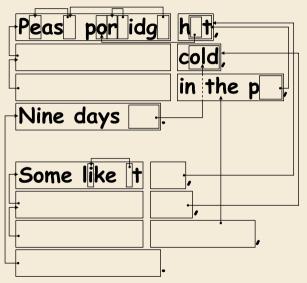
- extensions to larger alphabets possible (must store next character then)
- ▶ used in some image formats (e. g. TIFF)

# **Run-length encoding – Discussion**

- extensions to larger alphabets possible (must store next character then)
- ▶ used in some image formats (e.g. TIFF)
- fairly simple and fast
- can compress n bits to  $\Theta(\log n)$ ! for extreme case of constant number of runs
- negligible compression for many common types of data
  - ▶ No compression until run lengths  $k \ge 6$
  - **expansion** for run length k = 2 or 6

# 7.6 Lempel-Ziv-Welch

# Warmup

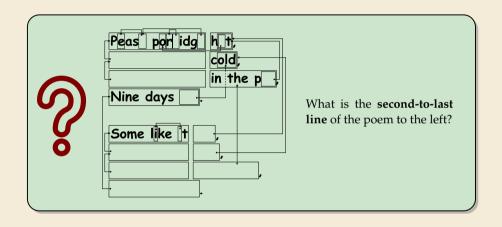




https://www.flickr.com/photos/quintanaroo/2742726346

https://classic.csunplugged.org/text-compression/

#### **Clicker Question**





#### **Lempel-Ziv Compression**

- ▶ Huffman and RLE mostly take advantage of frequent or repeated *single characters*.
- ▶ **Observation**: Certain *substrings* are much more frequent than others.
  - in English text: the, be, to, of, and, a, in, that, have, I
  - ▶ in HTML: "<a href", "<img src", "<br/>"

#### **Lempel-Ziv Compression**

- ▶ Huffman and RLE mostly take advantage of frequent or repeated *single characters*.
- ▶ **Observation**: Certain *substrings* are much more frequent than others.
  - ▶ in English text: the, be, to, of, and, a, in, that, have, I
  - ▶ in HTML: "<a href", "<img src", "<br/>"
- ▶ **Lempel-Ziv** stands for family of *adaptive* compression algorithms.
  - ► **Idea:** store repeated parts by reference!
  - → each codeword refers to
    - $\triangleright$  either a single character in  $\Sigma_S$ ,
    - or a *substring* of *S* (that both encoder and decoder have already seen).

#### **Lempel-Ziv Compression**

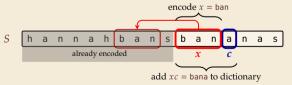
- ▶ Huffman and RLE mostly take advantage of frequent or repeated *single characters*.
- ▶ **Observation**: Certain *substrings* are much more frequent than others.
  - in English text: the, be, to, of, and, a, in, that, have, I
  - ▶ in HTML: "<a href", "<img src", "<br/>"
- ▶ **Lempel-Ziv** stands for family of *adaptive* compression algorithms.
  - ► **Idea:** store repeated parts by reference!
  - → each codeword refers to
    - ightharpoonup either a single character in  $\Sigma_S$ ,
    - or a *substring* of *S* (that both encoder and decoder have already seen).
  - ► Variants of Lempel-Ziv compression
    - "LZ77" Original version ("sliding window") Derivatives: LZSS, LZFG, LZRW, LZP, DEFLATE, ... DEFLATE used in (pk)zip, gzip, PNG
    - "LZ78" Second (slightly improved) version Derivatives: LZW, LZMW, LZAP, LZY, ... LZW used in compress, GIF

#### Lempel-Ziv-Welch

- ► here: Lempel-Ziv-Welch (LZW) (arguably the "cleanest" variant of Lempel-Ziv)
- ► variable-to-fixed encoding
  - ▶ all codewords have k bits (typical: k = 12)  $\rightsquigarrow$  fixed-length
  - but they represent a variable portion of the source text!

#### Lempel-Ziv-Welch

- ► here: Lempel-Ziv-Welch (LZW) (arguably the "cleanest" variant of Lempel-Ziv)
- variable-to-fixed encoding
  - ▶ all codewords have k bits (typical: k = 12)  $\rightsquigarrow$  fixed-length
  - but they represent a variable portion of the source text!
- $\blacktriangleright$  maintain a **dictionary** D with  $2^k$  entries  $\leadsto$  codewords = indices in dictionary
  - ▶ initially, first  $|\Sigma_S|$  entries encode single characters (rest is empty)
  - ▶ **add** a new entry to *D* **after each step**:
  - Encoding: after encoding a substring x of S, add xc to D where c is the character that follows x in S.

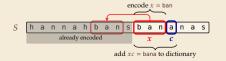


- $\rightarrow$  new codeword in D
- $\triangleright$  D actually stores codewords for x and c, not the expanded string

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

C =



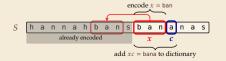
Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	
129	
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

C = 89



Code	String
32	
33	
79	0
82	R
85	U
89	Υ

Code	String
128	
129	
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	

 $\textbf{Input: Y0!\_Y0U!\_Y0UR\_Y0Y0!}$ 

 $\Sigma_S$  = ASCII character set (0–127)

C = 89



Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

		Υ	0
C	=	89	79



Code	String
32	
33	
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	

$$\Sigma_S$$
 = ASCII character set (0–127)

	Υ	0
C =	89	79

D	=

								ç	_	en	cod	le x	= b	an			
S	h	а	n	n	а	h	b	а	n	S	b	a	n	а	n	а	S
				alre	ady	enco	oded					x		c			
		add $xc = bana$ to dictionary															

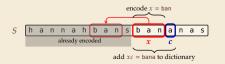
Code	String	
32	ш	
33	!	
79	0	
82	R	
85	U	
89	Υ	

Code	String
128	Y0
129	0!
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

	Υ	0	- !
C =	89	79	33



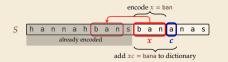
Code	String		
32			
33	!		
79	0		
82	R		
85	U		
89	Y		
	32 33		

Code	String
128	Y0
129	0!
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

	Υ	0	!
C =	89	79	33



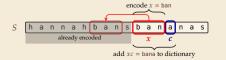
Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	1
131	
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	!	ш
C = 89	79	33	32



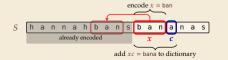
Code	String	
32	П	
33	!	
79	0	
82	R	
85	U	
89	Υ	

Code	String
128	Y0
129	0!
130	!
131	
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	. !	ш
C = 89	79	33	32



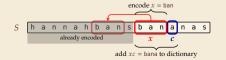
Code	String
32	
33	
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	!
131	٦Y
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	!	u	Y0
C = 89	79	33	32	128



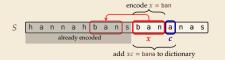
Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	
131	Y
132	
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

	Υ	0	!	ш	Y0
C =	89	79	33	32	128



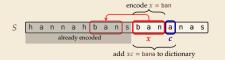
Code	String
32	ш
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	!
131	٦Y
132	YOU
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

	Υ	0	!	ш	Y0	U
C =	89	79	33	32	128	85



Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	!
131	цY
132	YOU
133	
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	!	u	Y0	U
C = 89	79	33	32	128	85



Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	!
131	¬А
132	YOU
133	U!
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

ſ

								ç		en	cod	e x	= b	an			
S	h	а	n	n	а	h	b	а	n	S	b	а	n	а	n	a	s
				alre	ady	enco	oded					х		c			
	add $xc = bana$ to dictionary																

Code	String
32	
33	!
79	0
82	R
85	U
89	Υ

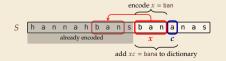
Code	String
128	Y0
129	0!
130	!
131	٦Y
132	YOU
133	U!
134	
135	
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	!	u	Y0	U	!
C = 89	79	33	32	128	85	130

=		



String
!
0
R
U
Υ

Code	String
128	Y0
129	0!
130	!
131	٦Y
132	YOU
133	U!
134	! <sub>L</sub> Y
135	
136	
137	
138	
139	

**Input**: Y0!,,Y0U!,,Y0UR,,Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	!	u	Y0	U	!	YOU
C = 89	79	33	32	128	85	130	132

L
L
L
r

D =

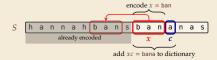
79	
82	
85	
89	

Code

32 33

String

Code	String
128	Y0
129	0!
130	!
131	¬А
132	YOU
133	U!
134	! <sub>L</sub> Y
135	
136	
137	
138	
139	



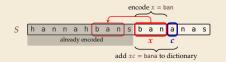
Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

$$C = 89$$
 79 33 32 128 85 130 132

=		

D



Code	String				
32					
33	!				
79	0				
82	R				
85	U				
89	Y				

Code	String
128	Y0
129	0!
130	!
131	¬А
132	YOU
133	U!
134	!_Y
135	YOUR
136	
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Υ	0	!	П	Y0	U	!	YOU	R
C = 89	79	33	32	128	85	130	132	82

33
79
82
85

D =

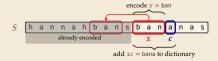
Code

32

89

String	
ш	
!	
0	
R	
U	
Υ	

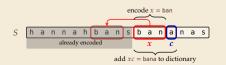
Code	String
128	Y0
129	0!
130	!
131	٦Y
132	YOU
133	U!
134	! <sub></sub> Y
135	YOUR
136	
137	
138	
139	



Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Y 0 ! U Y0 U !U Y0U R
C = 89 79 33 32 128 85 130 132 82



Code	String				
32					
33	!				
79	0				
82	R				
85	U				
89	Υ				

Code	String
128	Y0
129	0!
130	!
131	¬А
132	YOU
133	U!
134	! <sub>L</sub> Y
135	YOUR
136	R⊔
137	
138	
139	

Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

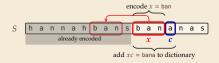
D =

	• •
32	
33	-:
79	0
82	R
85	U
89	Υ

Code

String

Code	String
128	Y0
129	0!
130	!
131	٦Y
132	YOU
133	U!
134	!_Y
135	YOUR
136	R⊔
137	
138	
139	



Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Y 0 ! \_ Y0 U ! \_ Y0U R \_Y C = 89 79 33 32 128 85 130 132 82 131

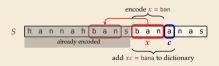
D =

П	
!	
0	
R	
U	
Υ	

Code

String

Strin
Y0
0!
!
L Y
YOU
U!
! <sub>L</sub> Y
YOUR
R⊔
Υ0
_

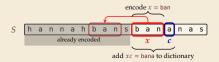


**Input**: Y0!\_Y0U!\_Y0UR\_Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Code	String
32	П
33	!
79	0
82	R
85	U
89	Y

Code	String
128	Y0
129	0!
130	!
131	¬А
132	YOU
133	U!
134	! <sub>L</sub> Y
135	YOUR
136	R⊔
137	۷0 ا
138	
139	

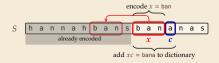


**Input**: Y0!\_Y0U!\_Y0UR\_Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Code	String
32	П
33	!
79	0
82	R
85	U
89	Υ

Code	String
128	Y0
129	0!
130	!
131	пV
132	YOU
133	U!
134	! <sub>L</sub> Y
135	YOUR
136	R⊔
137	۷0 ا
138	0Y
139	



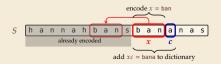
Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Y 0 ! \_ Y0 U ! \_ Y0U R \_Y 0 Y0 C = 89 79 33 32 128 85 130 132 82 131 79 128

Code	String
32	П
33	!
79	0
82	R
85	U
89	Y

Code	String
128	Y0
129	0!
130	!
131	٦Y
132	YOU
133	U!
134	!_Y
135	YOUR
136	R⊔
137	۷0 ا
138	0Y
139	



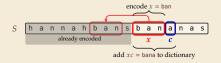
Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Y 0 ! L Y0 U ! Y0U R LY 0 Y0 C = 89 79 33 32 128 85 130 132 82 131 79 128

Code	String
32	П
33	!
79	0
82	R
85	U
89	Y

Code	String
128	Y0
129	0!
130	!
131	цY
132	YOU
133	U!
134	!Y
135	Y0UR
136	R⊔
137	۷0 ا
138	0Y
139	Y0!



Input: Y0! Y0U! Y0UR Y0Y0!

 $\Sigma_S$  = ASCII character set (0–127)

Y = 0 ! U = V = 0 V ! V = V = 0 Y V = V = 0 Y V = 0 Y V = 0 ! V = 0 Y

D =

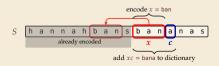
32	П
33	!
79	0
82	R
85	U
89	Y

. . .

Code

String

Code	Stri
128	Y0
129	0!
130	!
131	۲
132	YOU
133	U!
134	٧!
135	YOU
136	R⊔
137	YC
138	0Y
139	Y0!



### LZW encoding – Code

```
1 procedure LZWencode(S[0..n))
       x := \varepsilon // previous phrase, initially empty
      C := \varepsilon // output, initially empty
       D := dictionary, initialized with codes for c \in \Sigma_S // stored as trie
     k := |\Sigma_S| // next free codeword
    for i := 0, ..., n-1 do
           c := S[i]
7
           if D.containsKey(xc) then
                x := xc
           else
                C := C \cdot D.get(x) // append codeword for x
11
                D.put(xc, k) // add xc to D, assigning next free codeword
12
                k := k + 1: x := c
13
      end for
14
       C := C \cdot D.get(x)
15
       return C
16
```

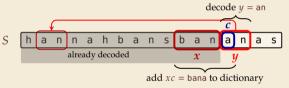
7.7 Lempel-Ziv-Welch Decoding

#### LZW decoding

▶ Decoder has to replay the process of growing the dictionary!

#### **→** Decoding:

after decoding a substring y of S, add xc to D, where x is previously encoded/decoded substring of S, and c = y[0] (first character of y)



 $\rightarrow$  Note: only start adding to *D* after *second* substring of *S* is decoded

► Same idea: build dictionary while reading string.

	Code #	String	
	32		
	65	Α	
) =	66	В	
	67	С	
	78	N	
	83	S	

input	decodes to	Code #	String (human)	String (computer)

► Same idea: build dictionary while reading string.

Code #	String
32	П
65	Α
66	В
67	С
78	N
83	S
	32  65 66 67  78

input	decodes to	Code #	String (human)	String (computer)
67	С			

► Same idea: build dictionary while reading string.

Code #	String
32	П
65	Α
66	В
67	С
78	N
83	S
	32  65 66 67  78

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A

► Same idea: build dictionary while reading string.

	Code #	String
	32	П
	65	Α
) =	66	В
	67	С
	78	N
	83	S

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A
78	N	129	AN	65, N

► Same idea: build dictionary while reading string.

	Code #	String
	32	
	65	Α
D =	66	В
	67	С
	78	N
	83	S

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A
78	N	129	AN	65, N
32	E .	130	N	78, ⊔

► Same idea: build dictionary while reading string.

	Code #	String	
	32		
	65	Α	
D =	66	В	
	67	С	
	78	N	
	83	S	
	83		

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A
78	N	129	AN	65, N
32	_	130	N	78, ⊔
66	В	131	uВ	32, B

## **LZW** decoding – Example

► Same idea: build dictionary while reading string.

**Example:** 67 65 78 32 66 129 133

	Code #	String
	32	
	65	Α
D =	66	В
	67	С
	78	N
	83	S

	decodes		String	String
input	to	Code #	(human)	(computer)
67	С			
65	А	128	CA	67, A
78	N	129	AN	65, N
32		130	N	78, ⊔
66	В	131	⊔В	32, B
129	AN	132	BA	66, A

## **LZW** decoding – Example

► Same idea: build dictionary while reading string.

**Example:** 67 65 78 32 66 129 133

	Code #	String
	32	
	65	Α
D =	66	В
	67	С
	78	N
	83	S

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A
78	N	129	AN	65, N
32	ш	130	N	78, ⊔
66	В	131	⊔В	32, B
129	AN	132	BA	66, A
133	???	133		

# **LZW** decoding – Example

► Same idea: build dictionary while reading string.

**Example:** 67 65 78 32 66 129 133

	Code #	String	
	32		
	65	Α	
) =	66	В	
	67	С	
	78	N	
	83	S	

input	decodes to	Code #	St: (hu		
67	С				
65	А	128	CA	67, A	
78	N	129	AN	65, N	
32		130	N	78, ⊔	
66	В	131	⊔В	32, B	
129	AN	132	BA	66, A	
133	???	133			

## LZW decoding – Bootstrapping

▶ example: Want to decode 133, but not yet in dictionary!



decoder is "one step behind" in creating dictionary

## LZW decoding – Bootstrapping

▶ example: Want to decode 133, but not yet in dictionary!



decoder is "one step behind" in creating dictionary

→ problem occurs if *we want to use a code* that we are *just about to build*.

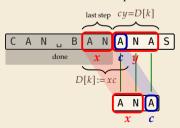
## LZW decoding – Bootstrapping

▶ example: Want to decode 133, but not yet in dictionary!



decoder is "one step behind" in creating dictionary

- → problem occurs if we want to use a code that we are just about to build.
- ▶ But then we actually know what is going on:
  - ightharpoonup Situation: decode using k in the step that will define k.
  - decoder knows last phrase x, needs phrase y = D[k] = xc.



- **1.** en/decode x.
- 2. store D[k] := xc
- 3. next phrase y equals D[k]

$$\rightarrow$$
  $D[k] = xc = x \cdot x[0]$  (all known)

## LZW decoding - Code

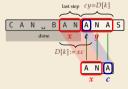
```
1 procedure LZWdecode(C[0..m))
       D := \text{dictionary } [0..2^d) \to \Sigma_S^+, initialized with codes for c \in \Sigma_S // stored as array
      k := |\Sigma_S| // next unused codeword
      q := C[0] // first codeword
      y := D[q] // lookup meaning of q in D
      S := y // output, initially first phrase
      for i := 1, ..., m-1 do
           x := y // remember last decoded phrase
           q := C[i] // next codeword
           if q == k then
10
                y := x \cdot x[0] // bootstrap case
11
           else
12
                y := D[a]
13
           S := S \cdot y // append decoded phrase
14
            D[k] := x \cdot y[0] // store new phrase
15
           k := k + 1
16
       end for
       return S
18
```

## LZW decoding - Example continued

**Example:** 67 65 78 32 66 129 133 83

	Code #	String	
	32	П	
	65	Α	
) =	66	В	
	67	С	
	78	N	
	83	S	

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A
78	N	129	AN	65, N
32	п	130	N	78, ⊔
66	В	131	uВ	32, B
129	AN	132	BA	66, A
133	ANA	133	ANA	129, A



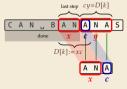
- 1. en/decode x.
- **2.** store D[k] := xc
- 3. next phrase y equals D[k] $D[k] = xc = x \cdot x[0]$  (all)

## LZW decoding - Example continued

**Example:** 67 65 78 32 66 129 133 83

	Code #	String
	32	
	65	Α
D =	66	В
	67	С
	78	N
	83	S

input	decodes to	Code #	String (human)	String (computer)
67	С			
65	Α	128	CA	67, A
78	N	129	AN	65, N
32	п	130	N	78, ⊔
66	В	131	uВ	32, B
129	AN	132	BA	66, A
133	ANA	133	ANA	129, A
83	S	134	ANAS	133, S



- 1. en/decode x.
- **2.** store D[k] := xc
- 3. next phrase y equals D[k] $D[k] = xc = x \cdot x[0]$  (all known)

## **Clicker Question**

How many phrases will LZW create on  $S = a^n$ , a run of n copies of as?



- $(\mathbf{A}) \sim n$
- $R \sim n/2$
- $\overline{\mathbf{C}}$  ~ n/4
- $\mathbf{D}$   $\Theta(n/\log n)$
- $\Theta(\sqrt{n})$

- $oldsymbol{\mathsf{F}} oldsymbol{\Theta}(\log n)$
- **G**  $\Theta(\log \log n)$
- (H) 2
- **I**) 1



→ sli.do/comp526

## **Clicker Question**

How many phrases will LZW create on  $S = a^n$ , a run of n copies of as?



## LZW - Discussion

- ▶ As presented, LZW uses coded alphabet  $\Sigma_C = [0..2^d)$ .
  - $\leadsto$  use another encoding for code numbers  $\mapsto$  binary, e.g., Huffman
- ▶ need a rule when dictionary is full; different options:
  - ightharpoonup increment  $d \rightsquigarrow$  longer codewords
  - ▶ "flush" dictionary and start from scratch → limits extra space usage
  - ▶ often: reserve a codeword to trigger flush at any time
- encoding and decoding both run in linear time (assuming  $|\Sigma_S|$  constant)

## LZW - Discussion

- ▶ As presented, LZW uses coded alphabet  $\Sigma_C = [0..2^d)$ .
  - $\leadsto$  use another encoding for code numbers  $\mapsto$  binary, e.g., Huffman
- ▶ need a rule when dictionary is full; different options:
  - ightharpoonup increment  $d \rightsquigarrow$  longer codewords
  - "flush" dictionary and start from scratch --> limits extra space usage
  - ▶ often: reserve a codeword to trigger flush at any time
- encoding and decoding both run in linear time (assuming  $|\Sigma_S|$  constant)
- fast encoding & decoding
- works in streaming model (no random access, no backtrack on input needed)
- significant compression for many types of data
- captures only local repetitions (with bounded dictionary)

# **Compression summary**

Huffman codes	Run-length encoding	Lempel-Ziv-Welch
fixed-to-variable	variable-to-variable	variable-to-fixed
2-pass	1-pass	1-pass
must send dictionary	can be worse than ASCII	can be worse than ASCII
60% compression on English text	bad on text	45% compression on English text
optimal binary character encopding	good on long runs (e.g., pictures)	good on English text
rarely used directly	rarely used directly	frequently used
part of pkzip, JPEG, MP3	fax machines, old picture-formats	GIF, part of PDF, Unix compress

# **Part III**

Text Transforms

#### **Text transformations**

- ▶ compression is effective is we have one the following:
  - ▶ long runs → RLE
  - ► frequently used characters → Huffman
  - ► many (local) repeated substrings → LZW

#### **Text transformations**

- compression is effective is we have one the following:
  - ▶ long runs → RLE
  - ► frequently used characters → Huffman
  - ► many (local) repeated substrings → LZW
- ▶ but methods can be frustratingly "blind" to other "obvious" redundancies
  - LZW: repetition too distant 7 dictionary already flushed
  - ► Huffman: changing probabilities (local clusters) 🕈 averaged out globally
  - ▶ RLE: run of alternating pairs of characters 🦅 not a run

#### **Text transformations**

- compression is effective is we have one the following:
  - ▶ long runs → RLE
  - ► frequently used characters → Huffman
  - ► many (local) repeated substrings → LZW
- ▶ but methods can be frustratingly "blind" to other "obvious" redundancies
  - LZW: repetition too distant 7 dictionary already flushed
  - ► Huffman: changing probabilities (local clusters) 🕴 averaged out globally
  - ▶ RLE: run of alternating pairs of characters 🦅 not a run
- ► Enter: text transformations
  - invertible functions of text
  - ▶ do not by themselves reduce the space usage
  - but help compressors "see" existing redundancy
  - → use as pre-/postprocessing in compression pipeline

7.8 Move-to-Front Transformation

#### **Move to Front**

- ▶ *Move to Front (MTF)* is a heuristic for *self-adjusting linked lists* 
  - unsorted linked list of objects
  - whenever an element is accessed, it is moved to the front of the list (leaving the relative order of other elements unchanged)
  - list "learns" probabilities of access to objects makes access to frequently requested ones cheaper

#### **Move to Front**

- ▶ *Move to Front (MTF)* is a heuristic for *self-adjusting linked lists* 
  - unsorted linked list of objects
  - whenever an element is accessed, it is moved to the front of the list (leaving the relative order of other elements unchanged)
  - ist "learns" probabilities of access to objects makes access to frequently requested ones cheaper
- ▶ Here: use such a list for storing source alphabet  $\Sigma_S$ 
  - $\triangleright$  to encode c, access it in list
  - encode *c* using its (old) position in list
  - then apply MTF to the list
  - $\rightsquigarrow$  codewords are integers, i. e.,  $\Sigma_C = [0..\sigma)$

#### Move to Front

- ▶ *Move to Front (MTF)* is a heuristic for *self-adjusting linked lists* 
  - unsorted linked list of objects
  - whenever an element is accessed, it is moved to the front of the list (leaving the relative order of other elements unchanged)
  - ist "learns" probabilities of access to objects makes access to frequently requested ones cheaper
- ▶ Here: use such a list for storing source alphabet  $\Sigma_S$ 
  - ightharpoonup to encode c, access it in list
  - encode *c* using its (old) position in list
  - ▶ then apply MTF to the list
  - $\rightsquigarrow$  codewords are integers, i. e.,  $\Sigma_C = [0..\sigma)$

## **Clicker Question**



Assume a MTF list currently contains the items XYZABC, and we now access A. What is the list content after the MTF rule has been applied?



→ sli.do/comp526

#### MTF - Code

#### ► Transform (encode):

```
procedure MTF-encode(S[0..n))

L := list containing <math>\Sigma_S (sorted order)

C := \varepsilon

for i := 0, ..., n-1 do

c := S[i]

p := position of c in L

C := C \cdot p

Move c to front of L

end for

return C
```

#### ► Inverse transform (decode):

```
1 procedure MTF-decode(C[0..m))
2 L := \text{list containing } \Sigma_S \text{ (sorted order)}
3 S := \varepsilon
4 \text{for } j := 0, \dots, m-1 \text{ do}
5 p := C[j]
6 c := \text{character at position } p \text{ in } L
7 S := S \cdot c
8 Move c to front of L
9 \text{end for}
10 \text{return } S
```

▶ Important: encoding and decoding produce same accesses to list

$$S = INEFFICIENCIES$$

$$C =$$

$$S = INEFFICIENCIES$$

$$C = 8$$

$$S = INEFFICIENCIES$$

$$C = 813$$

$$S = INEFFICIENCIES$$

$$C = 8136$$

$$S = INEFFICIENCIES$$

$$C = 81367$$

$$S = INEFFICIENCIES$$

$$C = 813670$$

$$S = INEFFICIENCIES$$

$$C = 8136703$$

$$S = INEFFICIENCIES$$

$$C = 81367036$$

$$S = INEFFICIENCIES$$

$$C = 813670361$$

$$S = INEFFICIENCIES$$

$$C = 8136703613433318$$

- ▶ What does a run in *S* encode to in *C*?
- ▶ What does a run in *C* mean about the source *S*?

## MTF - Discussion

- ► MTF itself does not compress text (if we store codewords with fixed length)
- → prime use as part of longer pipeline
- two simple ideas for encoding codewords:
  - ► Elias gamma code → smaller numbers gets shorter codewords works well for text with small "local effective" alphabet
  - ► Huffman code (better compression, but need 2 passes)
- ▶ but: most effective after BWT ( $\rightarrow$  next)

7.9 Burrows-Wheeler Transform

# **Burrows-Wheeler Transform**

- ▶ Burrows-Wheeler Transform (BWT) is a sophisticated text-transformation technique.
  - coded text has same letters as source, just in a different order
  - ▶ But: coded text is (typically) more compressible with MTF(!)

#### **Burrows-Wheeler Transform**

- ▶ Burrows-Wheeler Transform (BWT) is a sophisticated text-transformation technique.
  - coded text has same letters as source, just in a different order
  - ▶ But: coded text is (typically) more compressible with MTF(!)
- ▶ Encoding algorithm needs **all** of *S* (no streaming possible).
  - → BWT is a block compression method.

#### **Burrows-Wheeler Transform**

- ▶ Burrows-Wheeler Transform (BWT) is a sophisticated text-transformation technique.
  - coded text has same letters as source, just in a different order
  - ▶ But: coded text is (typically) more compressible with MTF(!)
- ► Encoding algorithm needs **all** of *S* (no streaming possible).
  - → BWT is a *block compression method*.
- ▶ BWT followed by MTF, RLE, and Huffman is the algorithm used by the bzip2 program. achieves best compression on English text of any algorithm we have seen:

```
4047392 bible.txt
1191071 bible.txt.gz
888604 bible.txt.7z
845635 bible.txt.bz2
```

# **BWT** transform

• *cyclic shift* of a string:

$$T = time_uflies_uquickly_u$$

flies\_quickly\_time\_



# **BWT** transform

- *cyclic shift* of a string:
- ► add *end-of-word character* \$ to *S* (as in Unit 6)

 $T = time_uflies_uquickly_u$ 



flies\_quickly\_time\_



# **BWT** transform

- ► *cyclic shift* of a string:
- ► add *end-of-word character* \$ to *S* (as in Unit 6)

 $T = time_{location} flies_{location} quickly_{location}$ 



flies\_quickly\_time\_

- ► The Burrows-Wheeler Transform proceeds in three steps:
  - **1.** Place *all cyclic shifts* of *S* in a list *L*
  - 2. Sort the strings in *L* lexicographically
  - **3.** *B* is the *list of trailing characters* (last column, top-down) of each string in *L*

# **BWT** transform – Example

 $S = alf_eats_alfalfa$ \$

1. Write all cyclic shifts

alf\_eats\_alfalfa\$ lf\_eats\_alfalfa\$a f\_eats\_alfalfa\$al \_eats\_alfalfa\$alf eats\_alfalfa\$alf ats.,alfalfa\$alf.,e ts\_alfalfa\$alf\_ea s\_alfalfa\$alf\_eat \_alfalfa\$alf\_eats alfalfa\$alf\_eats\_ lfalfa\$alf\_eats\_a falfa\$alf\_eats\_al alfa\$alf\_eats\_alf lfa\$alf\_eats\_alfa fa\$alf\_eats\_alfal a\$alf, eats, alfalf \$alf, eats, alfalfa

 $\stackrel{\overset{}{\sim}}{\sim}$ 

# BWT transform – Example

S = alf.eats.alfalfa

- **1.** Write all cyclic shifts
- 2. Sort cyclic shifts

alf, eats, alfalfa\$ lf..eats..alfalfa\$a f\_eats\_alfalfa\$al \_eats\_alfalfa\$alf eats, alfalfa\$alf., ats. alfalfa\$alf.e ts..alfalfa\$alf..ea s..alfalfa\$alf..eat ..alfalfa\$alf..eats alfalfa\$alf,eats,, lfalfa\$alf..eats..a falfa\$alf.eats.al alfa\$alf,eats,alf lfa\$alf, eats, alfa fa\$alf..eats..alfal a\$alf, eats, alfalf \$alf..eats..alfalfa

\$alf,eats,alfalfa ..alfalfa\$alf..eats \_eats\_alfalfa\$alf a\$alf\_eats\_alfalf alf\_eats\_alfalfa\$ alfa\$alf.eats.alf alfalfa\$alf..eats.. ats..alfalfa\$alf..e eats\_alfalfa\$alf\_ f, eats, alfalfa\$al fa\$alf,eats,alfal falfa\$alf,eats,al lf\_eats\_alfalfa\$a lfa\$alf.eats.alfa lfalfa\$alf..eats..a s..alfalfa\$alf\_eat ts..alfalfa\$alf..ea





# **BWT** transform – Example

 $S = alf_ueats_alfalfa$ 

- 1. Write all cyclic shifts
- 2. Sort cyclic shifts
- 3. Extract last column

 $B = asff f_e lllaaata$ 

alf, eats, alfalfa\$ lf..eats..alfalfa\$a f\_eats\_alfalfa\$alf \_eats\_alfalfa\$alf eats, alfalfa\$alf., ats. alfalfa\$alf.e ts..alfalfa\$alf..ea s..alfalfa\$alf..eat ..alfalfa\$alf..eats alfalfa\$alf,eats, lfalfa\$alf..eats..a falfa\$alf.eats.al alfa\$alf,eats,alf lfa\$alf, eats, alfa fa\$alf..eats..alfal a\$alf, eats, alfalf \$alf..eats..alfalfa

\$alf, eats, alfalfa .alfalfa\$alf.eats \_eats\_alfalfa\$alf a\$alf\_eats\_alfalf alf\_eats\_alfalfa\$ alfa\$alf.eats.alf alfalfa\$alf.eats.. ats.alfalfa\$alf.e eats alfalfa\$alf f\_eats\_alfalfa\$al fa\$alf,eats,alfal falfa\$alf,eats,al lf\_eats\_alfalfa\$a lfa\$alf.eats.alfa lfalfa\$alf.eats.a s. alfalfa\$alf\_eat ts..alfalfa\$alf..ea

**~**→

sort

# **Clicker Question**

What is the relation between suffix array L[0..n] and BWT B[0..n] of a string T[0..n)\$?



- f A L can be very easily computed from B and T
- f B B can be very easily computed from L and T
- C Both A and B
- **D** Neither A nor B



→ sli.do/comp526

# **Clicker Question**

What is the relation between suffix array L[0..n] and BWT B[0..n] of a string T[0..n)\$?



- ig(  $oldsymbol{A}ig)$   $oldsymbol{L}$  can be very easily computed from  $oldsymbol{B}$  and T
- **B** B can be very easily computed from L and T
- C Both A and B
- Neither A nor B



→ sli.do/comp526

# **BWT – Implementation & Properties**

#### Compute BWT efficiently:

- ightharpoonup cyclic shifts S = suffixes of S
- ► BWT is essentially suffix sorting!
  - $B[i] = S[L[i] 1] \qquad (L = \text{suffix array!})$ (if L[i] = 0, B[i] = \$)
  - $\sim$  Can compute *B* in O(n) time

```
\downarrow L[r]
alf_eats_alfalfa$
                         $alf_eats_alfalfa
lf..eats..alfalfa$a
                         ..alfalfa$alf_eats
f, eats, alfalfa$al
                         _eats_alfalfa$alf
_eats_alfalfa$alf
                          a$alf..eats..alfalf
eats_alfalfa$alf_
                          alf_eats_alfalfa$
                                               0
ats..alfalfa$alf..e
                         alfa$alf..eats..alf
ts.alfalfa$alf.ea
                         alfalfa$alf,.eats..
s..alfalfa$alf..eat
                          ats.alfalfa$alf.e
..alfalfa$alf,.eats
                         eats, alfalfa$alf...
alfalfa$alf..eats..
                          f..eats..alfalfa$a
lfalfa$alf_eats_a
                      10 fa$alf_eats_alfal
falfa$alf..eats..al
                          falfa$alf_eats_al
alfa$alf..eats..alf
                         lf.eats.alfalfa$a
lfa$alf..eats..alfa
                      13 lfa$alf_eats_alfa
fa$alf,.eats,.alfal
                      14 lfalfa$alf..eats..a
a$alf_eats_alfalf
                         s..alfalfa$alf..eat
$alf_eats_alfalfa
                         ts_alfalfa$alf_ea
                                               6
```

# **BWT – Implementation & Properties**

#### Compute BWT efficiently:

- ightharpoonup cyclic shifts  $S \cong \text{suffixes of } S$
- ► BWT is essentially suffix sorting!
  - $B[i] = S[L[i] 1] \qquad (L = \text{suffix array!})$ (if L[i] = 0, B[i] = \$)
  - $\rightsquigarrow$  Can compute *B* in O(n) time

#### Why does BWT help?

- sorting groups characters by what follows
  - Example: If always preceded by a
- $\rightarrow$  B has local clusters of characters
  - that makes MTF effective

- alf\_eats\_alfalfa\$ 0 \$alf\_eats\_alfalfa lf..eats..alfalfa\$a ..alfalfa\$alf..eats f, eats, alfalfa\$al \_eats\_alfalfa\$alf ..eats..alfalfa\$alf a\$alf..eats..alfalf eats\_alfalfa\$alf\_ alf.eats\_alfalfa\$ ats..alfalfa\$alf..e alfa\$alf..eats..alf ts.alfalfa\$alf.ea alfalfa\$alf,.eats.. ats.alfalfa\$alf.e s..alfalfa\$alf..eat ..alfalfa\$alf..eats eats, alfalfa\$alf. alfalfa\$alf..eats.. f..eats..alfalfa\$a lfalfa\$alf\_eats\_a 10 fa\$alf\_eats\_alfal falfa\$alf..eats..al falfa\$alf\_eats\_al alfa\$alf..eats..alf 12 lf.eats.alfalfa\$a lfa\$alf..eats..alfa 13 lfa\$alf\_eats\_alfa fa\$alf,.eats,.alfal 14 lfalfa\$alf..eats..a a\$alf\_eats\_alfalf s..alfalfa\$alf..eat \$alf\_eats\_alfalfa ts\_alfalfa\$alf\_ea
- ▶ repeated substring in  $S \rightsquigarrow runs$  of characters in B
  - ▶ picked up by RLE

 $\downarrow L[r]$ 

0

# **Bigger Example**

have..had..hadnt..hasnt..havent..has..what\$ ave., had, hadnt, hasnt, havent, has, what\$h ve.,had,,hadnt,,hasnt,,havent,,has,,what\$ha e..had..hadnt..hasnt..havent..has..what\$hav \_had\_hadnt\_hasnt\_havent.has.what\$have had..hadnt..hasnt..havent..has..what\$have.. ad.,hadnt,,hasnt,,havent,,has,,what\$have,,h d\_hadnt\_hasnt\_havent\_has\_what\$have\_ha \_hadnt\_hasnt\_havent\_has\_what\$have..had hadnt, hasnt, havent, has, what have, had, adnt.hasnt.havent.has.what\$have.had.h dnt\_hasnt\_havent\_has\_what\$have\_had\_ha nt\_hasnt\_havent\_has\_what\$have\_had\_had t.,hasnt,,havent,,has,,what\$have,,had,,hadn ..hasnt..havent..has..what\$have..had..hadnt hasnt havent has what shave had hadnt. asnt.,havent.,has.,what\$have,,had,,hadnt.,h snt.,havent.,has.,what\$have.,had,,hadnt.,ha nt.,havent.,has.,what\$have.,had.,hadnt.,has t. havent. has. what shave. had. hadnt. hasn ..havent..has..what\$have..had..hadnt..hasnt havent..has..what\$have..had..hadnt..hasnt... avent..has..what\$have..had..hadnt..hasnt..h vent.,has.,what\$have.,had.,hadnt.,hasnt.,ha ent has what have had hadnt hasnt hav nt..has..what\$have..had..hadnt..hasnt..have t, has, what \$have, had, hadnt, hasnt, haven ..has..what\$have..had..hadnt..hasnt..havent has whatshave had hadnt hasnt havent... as\_what\$have\_had\_hadnt\_hasnt\_havent.h s.,what\$have.,had.,hadnt.,hasnt.,havent.,ha ,what\$have,had,hadnt,hasnt,havent,has what\$have..had..hadnt..hasnt..havent..has... hat\$have, had, hadnt, hasnt, havent, has, w at\$have had hadnt hasnt havent has wh t\$have..had..hadnt..hasnt..havent..has..wha \$have, had, hadnt, hasnt, havent, has, what

\$have..had..hadnt..hasnt..havent..has..what .,had,,hadnt,,hasnt,,havent,,has,,what\$have .,hadnt,,hasnt,,havent,,has,,what\$have,,had ..has..what\$have\_had\_hadnt\_hasnt\_havent \_hasnt\_havent\_has\_what\$have\_had\_hadnt ..havent..has..what\$have..had..hadnt..hasnt .,what\$have,,had,,hadnt,,hasnt,,havent,,has ad..hadnt..hasnt..havent..has..what\$have..h adnt..hasnt..havent..has..what\$have..had..h as,,what\$have,,had,,hadnt,,hasnt,,havent,,h asnt..havent..has..what\$have..had..hadnt..h at\$have\_had\_hadnt\_hasnt\_havent\_has\_wh ave\_had\_hadnt\_hasnt\_havent\_has\_what\$h avent, has, what \$have, had, hadnt, hasnt, h d. hadnt. hasnt. havent. has. what\$have. ha dnt..hasnt..havent..has..what\$have..had..ha e.,had,,hadnt,,hasnt,,havent,,has,,what\$hav ent.,has,,what\$have,,had,,hadnt,,hasnt,,hav had, hadnt, hasnt, havent, has, what have... hadnt.hasnt.havent.has.whatshave.had.. has, what \$have, had, hadnt, hasnt, havent... hasnt.,havent.,has.,what\$have.,had.,hadnt., hat\$have,had,hadnt,hasnt,havent,has,w have..had..hadnt..hasnt..havent..has..what\$ havent has what shave had hadnt hasnt... nt.,has.,what\$have.,had.,hadnt.,hasnt.,have nt hasnt havent has what have had had nt..havent..has..what\$have..had..hadnt..has s.,what\$have.,had.,hadnt.,hasnt.,havent.,ha snt havent has what shave had hadnt ha t\$have..had..hadnt..hasnt..havent..has..wh a t, has, what \$have, had, hadnt, hasnt, have n t.,hasnt.,havent.,has.,what\$have.,had,,had n t..havent..has,,what\$have,,had,,hadnt,,has n ve. had. hadnt. hasnt. havent. has. whatsh a vent..has..what\$have..had..hadnt..hasnt..ha what\$have, had, hadnt, hasnt, havent, has,

T = have \_ had \_ had nt \_ has nt \_ have nt \_ has \_ what \$
B = tedttts hhhhhhhhaavv \_ \_ \_ w \$ \_ eds a a ann naa \_
MTF(B) = 8552008700000007090800010929987001000105

# **Clicker Question**

Consider  $T = \text{have\_had\_hadnt\_hasnt\_havent\_has\_what}$ . The BWT is  $B = \text{tedtttshhhhhhhaavv}_{\text{uuuu}}$  w\$\_edsaaannnaa\_. How can we explain the long run of hs in B?



- A h is the most frequent character
- B h always appears at the beginning of a word
- c almost all words start with h
- **D** h is always followed by a
- E all as are preceded by h
- F h is the 4th character in the alphabet



→ sli.do/comp526

# **Clicker Question**

Consider  $T = \text{have\_had\_hadnt\_hasnt\_havent\_has\_what}$ . The BWT is  $B = \text{tedtttshhhhhhhaavv}\_\_\_w\$\_\text{edsaaannnaa}\_$ . How can we explain the long run of hs in B?



- A h is the most frequent character
- B h always appears at the beginning of a word
- C almost all words start with h
- D h is always followed by a
- $\mathbf{E}$  all as are preceded by h  $\sqrt{\phantom{a}}$
- F h is the 4th character in the alphabet



→ sli.do/comp526

# 7.10 Inverse BWT

▶ Great, can compute BWT efficiently and it helps compression. *But how can we decode it?* 

not even obvious that it is at all invertible!

▶ Great, can compute BWT efficiently and it helps compression. *But how can we decode it?* 

not even obvious that it is at all invertible!

#### ► "Magic" solution:

- **1.** Create array D[0..n] of pairs: D[r] = (B[r], r).
- 2. Sort *D* stably with respect to first entry.
- **3.** Use *D* as linked list with (char, next entry)

► Great, can compute BWT efficiently and it helps compression. But how can we decode it?

D

not even obvious that it is at all invertible!

### ► "Magic" solution:

- **1.** Create array D[0..n] of pairs: D[r] = (B[r], r).
- **2.** Sort *D* stably with respect to first entry.
- **3.** Use *D* as linked list with (char, next entry)

#### Example:

B = ard\$rcaaaabb

$$S =$$

- 0 (a, 0)
- 1 (r, 1) 2 (d, 2)
- з (\$, 3)
- 4 (r, 4)
- 5 (c, 5)
- 6 (a, 6)
- 7 (a, 7)
- 8 (a, 8)
- 9 (a, 9)
- 10 (b, 10)
- 11 (b, 11)

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .	1 (r, 1) 2 (d, 2)	1 (a, 0) 2 (a, 6)	
2. Sort <i>D</i> stably with respect to first entry.	з (\$, 3)	з (a, 7)	
<b>3.</b> Use <i>D</i> as linked list with (char, next entry)	4 (r, 4) 5 (c, 5)	4 (a, 8) 5 (a, 9)	
Example:	6 (a, 6) 7 (a, 7)	6 (b, 10) 7 (b, 11)	
B = ard \$rcaaaabb S =	8 (a, 8) 9 (a, 9)	8 (c, 5) 9 (d, 2)	
	10 (b, 10) 11 (b, 11)	10 (r, 1) 11 (r, 4)	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs:	ı (r, 1)	1 (a, 0)	
D[r] = (B[r], r). <b>2.</b> Sort <i>D</i> stably with	2 (d, 2) 3 (\$, 3)	(a, 6) 3 $(a, 7)$	
respect to <i>first entry</i> .  3. Use <i>D</i> as linked list with	4 (r, 4)	4 (a, 8)	
(char, next entry)	5 (c, 5) 6 (a, 6)	5 (a, 9) 6 (b, 10)	
Example:	7 (a, 7)	7 (b, 11)	
B = ard\$rcaaaabb S = a	8 (a, 8) 9 (a, 9)	8 (c, 5) 9 (d, 2)	
$S = \mathbf{a}$	10 (b, 10)	10 (r, 1)	
	11 (b. 11)	11 (r. 4)	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .	1 (r, 1) 2 (d, 2)	1 (a, 0) 2 (a, 6)	
2. Sort <i>D</i> stably with respect to first entry.	3 (\$, 3) 4 (r, 4)	3 (a, 7) 4 (a, 8)	
3. Use <i>D</i> as linked list with (char, next entry)	5 (c, 5) 6 (a, 6)	5 (a, 9) 6 (b, 10)	
Example:  B = ard\$rcaaaabb	7 (a, 7) 8 (a, 8)	7 (b, 11) 8 (c, 5)	
S = ab	9 (a, 9) 10 (b, 10)	9 (d, 2) 10 (r, 1)	
	11 (b, 11)	11 (r, 4)	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	( 0)	char next	
· ·	0 (a, 0)	0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs:	ı (r, 1)	ı (a, 0)	
D[r] = (B[r], r).	2 (d, 2)	2 (a, 6)	
2. Sort D stably with	з (\$, 3)	з (a, 7)	
respect to first entry.	4 (r, 4)	4 (a, 8)	
<b>3.</b> Use <i>D</i> as linked list with (char, next entry)	5 (c, 5)	5 (a, 9)	
(char, next chary)	6 (a, 6)	6 (b, 10)	
Example:	7 (a, 7)	7 (b, 11)	
B = ard rcaaaabb	8 (a, 8)	8 (c, 5)	
S = abr	9 (a, 9)	9 (d, 2)	
	10 (b, 10)	10 (r, 1)	
	11 (h 11)	$\langle 11 (r 4) \rangle$	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .	1 (r, 1) 2 (d, 2)	1 (a, 0) 2 (a, 6)	
2. Sort <i>D</i> stably with respect to first entry.	з (\$, 3)	з (a, 7)	
3. Use <i>D</i> as linked list with (char, next entry)	4 (r, 4) 5 (c, 5)	4 (a, 8) 5 (a, 9)	
Example:	6 (a, 6) 7 (a, 7)	6 (b, 10) 7 (b, 11)	
$B = \text{ard} \cdot \text{rcaa} \cdot \text{aabb}$ S = abra	8 (a, 8) 9 (a, 9)	8 (c, 5) 9 (d, 2)	
5 — aut a	10 (b, 10)	10 (r, 1)	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .	1 (r, 1) 2 (d, 2)	1 (a, 0) 2 (a, 6)	
2. Sort <i>D</i> stably with respect to first entry.	3 (\$, 3) 4 (r, 4)	3 (a, 7) 4 (a, 8)	
<b>3.</b> Use <i>D</i> as linked list with (char, next entry)	5 (c, 5) 6 (a, 6)	5 (a, 9) 6 (b, 10)	
Example:	7 (a, 7) 8 (a, 8)	(b, 11) 8 (c, 5)	
B = ard\$rcaaaabb S = abrac	9 (a, 9) 10 (b, 10)	9 (d, 2) 10 (r, 1)	
	10 (b, 10) 11 (b, 11)	10 (r, 1) 11 (r, 4)	

			not even obvious that
	D	sorted $D$	it is at all invertible!
W. N 2 - 11 - 1 - 12 - 12 - 12 - 12 - 12		char next	
► "Magic" solution:	0 (a, 0)	0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs:	ı (r, 1)	ı (a, 0)	
D[r] = (B[r], r).	2 (d, 2)	2 (a, 6)	
2. Sort D stably with	з (\$, 3)	з (a, 7)	
respect to first entry.	4 (r, 4)	4 (a, 8)	
3. Use <i>D</i> as linked list with	5 (c, 5)	→ 5 (a, 9)	
(char, next entry)	6 (a, 6)	(b, 10)	
Evenuelle	7 (a, 7)	7 (b, 11)	
Example: $B = ard rcaaa$ bb	8 (a, 8)	8 (c, 5)	
S = abraca	9 (a, 9)	9 (d, 2)	
o – abi aca	10 (b, 10)	10 (r, 1)	
	11 (b.11)	11 (r. 4)	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs:	ı (r, 1)	ı (a, 0)	
D[r] = (B[r], r).	2 (d, 2)	2 (a, 6)	
2. Sort D stably with	з (\$, 3)	з (a, 7)	
respect to first entry.	4 (r, 4)	4 (a, 8)	
3. Use <i>D</i> as linked list with (char, next entry)	s (c, 5)	5 (a, 9)	
(char, next chary)	6 (a, 6)	6 (b, 10)	
Example:	7 (a, 7)	7 (b, 11)	
B = ard\$rcaaaabb	8 (a, 8)	8 (c, 5)	
S = abracad	9 (a, 9)	9 (d, 2)	
	10 (b, 10)	10 (r, 1)	
	11 (b, 11)	11 (r, 4)	

▶ Great, can compute BWT efficiently and it helps compression. *But how can we decode it?* 

not even obvious that Dit is at all invertible! sorted D char next ► "Magic" solution: o(a, 0)0 (\$, 3) **1.** Create array D[0..n] of pairs: ı (r, 1) ı (a, 0) D[r] = (B[r], r).2 (d, 2) (a, 6)2. Sort D stably with (a, 7)**3** (\$, 3) respect to first entry. 4 (r, 4) (a, 8)3. Use D as linked list with 5 (c, 5) (char, next entry) 6 (a, 6) (b, 11)(a, 7)Example: (c, 5)(a, 8)B = ard\$rcaaaabb(a, 9)(d, 2)S = abracada10 (b, 10) 10 (r, 1) 11 (b, 11) 11 (r, 4)

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .	1 (r, 1) 2 (d, 2)	1 (a, 0) 2 (a, 6)	
<b>2.</b> Sort <i>D</i> stably with respect to first entry.	3 (\$, 3) 4 (r, 4)	3 (a, 7) 4 (a, 8)	
<b>3.</b> Use <i>D</i> as linked list with (char, next entry)	5 (c, 5)	5 (a, 9)	
Example:	6 (a, 6) 7 (a, 7)	6 (b, 10) 7 (b, 11)	
B = ard\$rcaaaabb S = abracadab	8 (a, 8) 9 (a, 9)	8 (c, 5) 9 (d, 2)	
	10 (b, 10) 11 (b, 11)	10 (r, 1) 11 (r, 4)	

	D	sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	o (a, 0)	char next 0 (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .	1 (r, 1) 2 (d, 2)	1 (a, 0) 2 (a, 6)	
<b>2.</b> Sort <i>D</i> stably with respect to first entry.	3 (\$, 3) 4 (r, 4)	3 (a, 7) 4 (a, 8)	
<b>3.</b> Use <i>D</i> as linked list with (char, next entry)	5 (c, 5)	s (a, 9)	
Example:	6 (a, 6) 7 (a, 7)	6 (b, 10) 7 (b, 11)	
B = ard\$rcaaaabb S = abracadabr	8 (a, 8) 9 (a, 9)	8 (c, 5) 9 (d, 2)	
	10 (b, 10) 11 (b, 11)	10 (r, 1) 11 (r, 4)	

▶ Great, can compute BWT efficiently and it helps compression. *But how can we decode it?* 

not even obvious that Dit is at all invertible! sorted D char next ► "Magic" solution: o(a, 0)0 (\$, 3) **1.** Create array D[0..n] of pairs: ı (r, 1) (a, 0)D[r] = (B[r], r).2 (d, 2) (a, 6)2. Sort D stably with **3** (\$, 3) (a, 7)respect to first entry. 4 (r, 4) (a, 8)3. Use D as linked list with 5 (c, 5) (char, next entry) 6 (a, 6) (a, 7)(b, 11)Example: (c, 5)(a, 8)B = ard\$rcaaaabb(a, 9)(d, 2)S = abracadabra10 (b, 10) (r, 1)11 (b, 11) 11 (r, 4)

		D		sorted D	not even obvious that it is at all invertible!
► "Magic" solution:	0	(a, 0)	→ 0	char next (\$, 3)	
<b>1.</b> Create array $D[0n]$ of pairs: $D[r] = (B[r], r)$ .		(r, 1) (d, 2)	1 2	(a, 0) $(a, 6)$	
<b>2.</b> Sort <i>D</i> stably with respect to first entry.	3	(\$, 3) (r, 4)	3	(a, 7) (a, 8)	
<b>3.</b> Use <i>D</i> as linked list with (char, next entry)	5	(c, 5)	5	(a, 9)	
Example:	7	(a, 6) (a, 7)	7	(b, 10) (b, 11)	
B = ard rcaaaabb S = abracadabra		(a, 8) (a, 9)		(c, 5) (d, 2)	
		(b, 10) (b, 11)		(r, 1) (r, 4)	

# Inverse BWT – The magic revealed

- ► Inverse BWT very easy to compute:
  - ▶ only sort individual characters in *B* (not suffixes)
  - $\rightsquigarrow$  O(n) with counting sort
- ▶ but why does this work!?

# Inverse BWT – The magic revealed

- ► Inverse BWT very easy to compute:
  - ▶ only sort individual characters in *B* (not suffixes)
  - $\rightarrow$  O(n) with counting sort
- ▶ but why does this work!?
- decode char by char
  - can find unique \$ \to starting row
- ▶ to get next char, we need
  - (i) char in *first* column of *current row*
  - (ii) find row with that char's copy in BWT
  - $\rightsquigarrow\,$  then we can walk through and decode

# Inverse BWT – The magic revealed

- ► Inverse BWT very easy to compute:
  - ▶ only sort individual characters in *B* (not suffixes)
  - $\rightsquigarrow$  O(n) with counting sort
- ▶ but why does this work!?
- ▶ decode char by char
  - can find unique \$ \to starting row
- ▶ to get next char, we need
  - (i) char in *first* column of *current row*
  - (ii) find row with that char's copy in BWT
  - → then we can walk through and decode
- ► for (i): first column = characters of *B* in sorted order

## Inverse BWT – The magic revealed

► can find unique \$ → starting row

- ► Inverse BWT very easy to compute:
  - only sort individual characters in *B* (not suffixes)
  - $\rightsquigarrow$  O(n) with counting sort
- but why does this work!?
- decode char by char
- L[r]▶ to get next char, we need (i) char in *first* column of *current row* (ii) find row with that char's copy in BWT ► for (i): first column = characters of *B* in sorted order • for (ii): relative order of same character stays same:

ith a in first column = ith a in BWT

 $\rightsquigarrow$  stably sorting (B[r], r) by first entry enough

 $T_{L[r]}$ B[r]\$bananaba n aban\$bana n an\$banana b - anaban\$ba n ananaban\$ b ban\$banan a < bananaban \$ n\$bananaba < naban\$ban a < nanaban\$ba <

## **BWT - Discussion**

- ▶ Running time:  $\Theta(n)$ 
  - encoding uses suffix sorting
  - ▶ decoding only needs counting sort
  - $\rightsquigarrow$  decoding much simpler & faster (but same  $\Theta$ -class)

## **BWT - Discussion**

- ▶ Running time:  $\Theta(n)$ 
  - encoding uses suffix sorting
  - decoding only needs counting sort
  - $\rightsquigarrow$  decoding much simpler & faster (but same  $\Theta$ -class)
- typically slower than other methods
- need access to entire text (or apply to blocks independently)
- BWT-MTF-RLE-Huffman (bzip2) pipeline tends to have best compression

## **Summary of Compression Methods**

- Huffman Variable-width, single-character (optimal in this case)
  - RLE Variable-width, multiple-character encoding
    - LZW Adaptive, fixed-width, multiple-character encoding Augments dictionary with repeated substrings
    - MTF Adaptive, transforms to smaller integers should be followed by variable-width integer encoding
    - BWT Block compression method, should be followed by MTF