# 7 Compression

*20 April 2021*

Sebastian Wild

# Outline

# 7 Compression

## 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

**Click on "Polls" tab**

# Terminology

- ▶ **source text:** string $S \in \Sigma_S^\star$ to be stored / transmitted
    $\Sigma_S$ is some alphabet

- ▶ **coded text:** encoded data $C \in \Sigma_C^\star$ that is actually stored / transmitted
    usually use $\Sigma_C = \{0, 1\}$

- ▶ **encoding:** algorithm mapping source texts to coded texts $\quad S \rightharpoonup C$

- ▶ **decoding:** algorithm mapping coded texts back to original source text $\quad S \leftharpoonup C$

text can be any data

# Terminology

- ▶ **source text:** string $S \in \Sigma_S^\star$ to be stored / transmitted
  $\Sigma_S$ is some alphabet

- ▶ **coded text:** encoded data $C \in \Sigma_C^\star$ 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**          $S \to C \to S' \approx S$
  - ▶ **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.

2

# 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

# 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 *compression ratio*:
$$\underbrace{\frac{|C| \cdot \lg |\Sigma_C|}{|S| \cdot \lg |\Sigma_S|}}_{} \overset{\Sigma_C = \{0,1\}}{=} \frac{|C|}{|S| \cdot \lg |\Sigma_S|}$$

  (coded length — source length annotations)

  - $< 1$ means successful compression
  - $= 1$ means no compression
  - $> 1$ means "compression" made it bigger!? (yes, that happens ...)

3

# Limits of algorithmic compression

*Is this image compressible?*

# Limits of algorithmic compression

*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!**

# Limits of algorithmic compression

*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*

    - ► $C$ = *any program* that outputs $S$

        self-extracting archives!

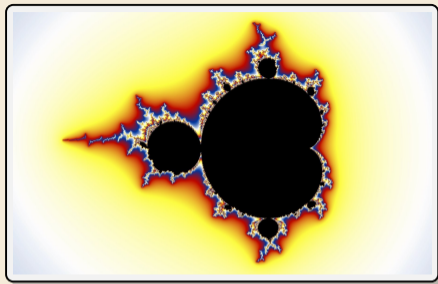    - ► Kolmogorov complexity = length of smallest such program

# Limits of algorithmic compression

*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*
  - ▶ $C$ = *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:

1. **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:

1. **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 (→ 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 $E : \Sigma_S \to \Sigma_C^\star$
    - 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 @   1010000 P   1100000 '   1110000 p
0000001 SOH   0010001 DC1   0100001 !     0110001 1   1000001 A   1010001 Q   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 ENQ   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 g   1110111 w
0001000 BS    0011000 CAN   0101000 (     0111000 8   1001000 H   1011000 X   1101000 h   1111000 x
0001001 HT    0011001 EM    0101001 )     0111001 9   1001001 I   1011001 Y   1101001 i   1111001 y
0001010 LF    0011010 SUB   0101010 *     0111010 :   1001010 J   1011010 Z   1101010 j   1111010 z
0001011 VT    0011011 ESC   0101011 +     0111011 ;   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 ]   1101101 m   1111101 }
0001110 SO    0011110 RS    0101110 .     0111110 >   1001110 N   1011110 ^   1101110 n   1111110 ~
0001111 SI    0011111 US    0101111 /     0111111 ?   1001111 O   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

👍 Encoding & Decoding as fast as it gets    & allows random access

👎 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**

International Morse Code

1. The length of a dot is one unit.
2. A dash is three units.
3. The space between parts of the same letter is one unit.
4. The space between letters is three units.
5. The space between words is seven units.

A ● ▬
B ▬ ● ● ●
C ▬ ● ▬ ●
D ▬ ● ●
E ●
F ● ● ▬ ●
G ▬ ▬ ●
H ● ● ● ●
I ● ●
J ● ▬ ▬ ▬
K ▬ ● ▬
L ● ▬ ● ●
M ▬ ▬
N ▬ ●
O ▬ ▬ ▬
P ● ▬ ▬ ●
Q ▬ ▬ ● ▬
R ● ▬ ●
S ● ● ●
T ▬

U ● ● ▬
V ● ● ● ▬
W ● ▬ ▬
X ▬ ● ● ▬
Y ▬ ● ▬ ▬
Z ▬ ▬ ● ●

1 ● ▬ ▬ ▬ ▬
2 ● ● ▬ ▬ ▬
3 ● ● ● ▬ ▬
4 ● ● ● ● ▬
5 ● ● ● ● ●
6 ▬ ● ● ● ●
7 ▬ ▬ ● ● ●
8 ▬ ▬ ▬ ● ●
9 ▬ ▬ ▬ ▬ ●
0 ▬ ▬ ▬ ▬ ▬

https://commons.wikimedia.org/wiki/File:
International_Morse_Code.svg

Dot ◄······     start     Dash ▬ ▬ ►

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

5 4   3     2    + 1   6 = 7   8 9 0

https://commons.wikimedia.org/wiki/File:Morse-code-tree.svg

9

# 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

**B** 2

**C** 3

**D** 4

**E** 26

**F** 36

**G** 256

`sli.do/comp526`

Click on "Polls" tab

# 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~~

B ~~2~~

C 3 ✓

D ~~4~~

E ~~26~~

F ~~36~~

G ~~256~~

`sli.do/comp526`

Click on "Polls" tab

## 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 `1`s 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 (hexadecimal) | UTF-8 octet sequence (binary) |
|---|---|
| `0000 0000 – 0000 007F` | `0xxxxxxx` |
| `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.

↯ random access

10

## Pitfall in variable-length codes

▶ Suppose we have the following code:

| $c$ | a | n | b | s |
|------|---|----|-----|-----|
| $E(c)$ | 0 | 10 | 110 | 100 |

▶ Happily encode text $S$ = banana with the coded text $C$ = <u>110</u><u>0</u><u>10</u><u>0</u><u>10</u><u>0</u>
                                                             b  a  n  a  n  a

## Pitfall in variable-length codes

- Suppose we have the following code:

| $c$ | a | n | b | s |
|-----|---|---|---|---|
| $E(c)$ | 0 | 10 | 110 | 100 |

- Happily encode text $S$ = banana with the coded text $C$ = <u>110</u><u>0</u><u>10</u><u>0</u><u>10</u><u>0</u>
  b  a n a n a

⚡ $C$ = 1100100100 decodes **both** to banana and to bass: <u>110</u>0<u>100</u><u>100</u>
  b a  s  s

⤳ not a valid code . . .   (cannot tolerate ambiguity)

but how should we have known?

## Pitfall in variable-length codes

▶ Suppose we have the following code:

| $c$ | a | n | b | s |
|---|---|---|---|---|
| $E(c)$ | 0 | 10 | 110 | 100 |

▶ Happily encode text $S =$ banana with the coded text $C = \underline{110}\underline{0}\underline{10}\underline{0}\underline{10}\underline{0}$
b a n a n a

⚡ $C =$ 1100100100 decodes **both** to banana and to bass: $\underline{110}\underline{0}\underline{100}\underline{100}$
b a s s

↝ not a valid code . . .    (cannot tolerate ambiguity)

but how should we have known?

🏂 $E(\text{n}) =$ 10 is a (proper) **prefix** of $E(\text{s}) =$ 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 $\implies$ instantaneously decodable

# Code tries

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

▶ **Example:**

| $c$ | A | E | N | O | T | ␣ |
|---|---|---|---|---|---|---|
| $E(c)$ | 01 | 101 | 001 | 100 | 11 | 000 |

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    01 001 000 01

▶ Decode <u>11</u><u>10</u><u>00</u>001010111    T O ␣

12

# Code tries

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

► **Example:**

| $c$ | A | E | N | O | T | ␣ |
|------|-----|-----|-----|-----|-----|-----|
| $E(c)$ | 01 | 101 | 001 | 100 | 11 | 000 |

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 111000001010111 → TO␣EAT

12

# 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 <u>message,</u> 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% | ▬▬▬ | **l** | 4.03% | ▬ | **b** | 1.49% | ■ |
| **a** | 8.17% | ▬▬▬ | **c** | 2.78% | ■ | **v** | 0.98% | ▪ |
| **o** | 7.51% | ▬▬▬ | **u** | 2.76% | ■ | **k** | 0.77% | ▪ |
| **i** | 6.97% | ▬▬ | **m** | 2.41% | ■ | **j** | 0.15% | ॰ |
| **n** | 6.75% | ▬▬ | **w** | 2.36% | ■ | **x** | 0.15% | ॰ |
| **s** | 6.33% | ▬▬ | **f** | 2.23% | ■ | **q** | 0.10% | ॰ |
| **h** | 6.09% | ▬▬ | **g** | 2.02% | ■ | **z** | 0.07% | ॰ |
| **r** | 5.99% | ▬▬ | **y** | 1.97% | ■ | | | |

↝ Want shorter codes for more frequent characters!

# Huffman coding

▶ **Given:** $\Sigma$ and weights $w : \Sigma \rightarrow \mathbb{R}_{\geq 0}$

    e. g. frequencies / probabilities

▶ **Goal:** prefix-free code $E$ (= code trie) for $\Sigma$ that minimizes coded text length

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

# Huffman coding

▶ **Given:** $\Sigma$ and weights $w : \Sigma \to \mathbb{R}_{\geq 0}$

e. g. frequencies / probabilities

▶ **Goal:** prefix-free code $E$ (= code trie) for $\Sigma$ that minimizes coded text length

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

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

$\leadsto$ 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 $u \in \{0,1\}^\star$     (*u* to be determined)

2. (Conceptually) replace a and b by a single character "ab"
   with $w(\boxed{ab}) = w(a) + w(b)$.

3. Recursively apply Huffman's algorithm on the smaller alphabet.
   This in particular determines $u = E(\boxed{ab})$.

# Huffman's algorithm

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

*ambiguous parts*

**Huffman's algorithm:**

1. Find two characters a, b with lowest weights.  *which?*
   ▶ We will encode them with the same prefix, plus one distinguishing bit,
     i. e., $E(a) = u0$ and $E(b) = u1$ for a bitstring $u \in \{0,1\}^\star$   *(u to be determined)*

2. (Conceptually) replace a and b by a single character "ab"  *or* $\boxed{ba}$ *?*
   with $w(\boxed{ab}) = w(a) + w(b)$.

3. Recursively apply Huffman's algorithm on the smaller alphabet.
   This in particular determines $u = E(\boxed{ab})$.

▶ 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

# Huffman's algorithm – Example

- Example text: $S = $ LOSSLESS $\rightsquigarrow \Sigma_S = \{$E, L, O, S$\}$

- Character frequencies: E : 1, L : 2, O : 1, S : 4

## Huffman's algorithm – Example

▶ Example text: $S = \text{LOSSLESS}$     $\rightsquigarrow$   $\Sigma_S = \{\text{E}, \text{L}, \text{O}, \text{S}\}$

▶ Character frequencies: $\text{E} : 1$,   $\text{L} : 2$,   $\text{O} : 1$,   $\text{S} : 4$

# Huffman's algorithm – Example

- Example text:  $S = $ LOSSLESS     $\rightsquigarrow$  $\Sigma_S = \{$E, L, O, S$\}$

- Character frequencies: E : 1,   L : 2,   O : 1,   S : 4

## Huffman's algorithm – Example

▶ Example text: $S = $ LOSSLESS $\leadsto \Sigma_S = \{$E, L, O, S$\}$

▶ Character frequencies: E : 1, L : 2, O : 1, S : 4

## Huffman's algorithm – Example

▶ Example text: $S = $ LOSSLESS $\rightsquigarrow \Sigma_S = \{$E, L, O, S$\}$

▶ Character frequencies: E : 1,   L : 2,   O : 1,   S : 4



$\rightsquigarrow$ *Huffman tree*  (code trie for Huffman code)

## Huffman's algorithm – Example

▶ Example text: $S = \text{LOSSLESS}$ $\leadsto$ $\Sigma_S = \{\text{E, L, O, S}\}$

▶ Character frequencies: $\text{E}:1$, $\text{L}:2$, $\text{O}:1$, $\text{S}:4$

| $c$ | E | L | O | S |
|-----|-----|-----|-----|-----|
| $E(c)$ | 000 | 01 | 001 | 1 |



8

$\leadsto$ *Huffman tree* (code trie for Huffman code)

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

freqs

(but: would also have to store trie)

17

# 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:

> *1.* 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).
>
> *3.* When combining trees of equal value,
> place the one containing the smallest letter to the left.

## 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 $E : \Sigma \to \{0, 1\}^\star$ with minimal expected codeword length $\underline{\ell(E) = \sum_{c \in \Sigma} w(c) \cdot |E(c)|}$ 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 $E : \Sigma \to \{0,1\}^\star$ 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 $\sigma = |\Sigma|$

▶ Given any optimal prefix-free code $E^*$ (as its code trie).

▶ code trie $\rightsquigarrow$ $\exists$ 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 $\Sigma' = \Sigma \setminus \{a, b\} \cup \{\boxed{ab}\}$ yields optimal code for $\Sigma$ by replacing leaf $\boxed{ab}$ by internal node with children a and b.

$\rightsquigarrow$ recursive call yields optimal code for $\Sigma'$ by inductive hypothesis, so Huffman's algorithm finds optimal code for $\Sigma$.

◄

## Entropy

### Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

$$\mathcal{H}(p_1, \ldots, p_n) = -\sum_{i=1}^{n} p_i \lg p_i = \sum_{i=1}^{n} p_i \lg\left(\frac{1}{p_i}\right) = \mathbb{E}\left[\lg \frac{1}{P}\right] \quad \blacktriangleleft$$

fair die with 6 faces

$1 \ldots 6$ with $\frac{1}{6}$

$$\mathcal{H}\left(\frac{1}{6} \ldots, \frac{1}{6}\right) = \sum_{i=1}^{6} \frac{1}{6} \lg\left(\frac{1}{\frac{1}{6}}\right) = 1 \cdot \lg(6) \approx 2.$$

fair coin   heads / tails   w/ prob $\frac{1}{2}$

$$\mathcal{H}\left(\frac{1}{2}, \frac{1}{2}\right) = 1$$

## Entropy

### Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

  ▶ *"20 Questions on $[0, 1)$":* Land inside my interval by halving.

```
├────────────────────────────────────────────────────────────┤
0                                                              1
```

# Entropy

### Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

  ▶ *"20 Questions on $[0, 1)$":* Land inside my interval by halving.

## Entropy

### Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

    ▶ *"20 Questions on $[0, 1)$":* Land inside my interval by halving.

# Entropy

## Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

▶ *"20 Questions on $[0, 1)$":* Land inside my interval by halving.

## Entropy

### Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

  ▶ *"20 Questions on $[0, 1)$":* Land inside my interval by halving.

# Entropy

## Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

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

# Entropy

## Definition 7.2 (Entropy)

Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

    ▶ *"20 Questions on $[0,1)$":* Land inside my interval by halving.

# Entropy

### Definition 7.2 (Entropy)
Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◀

▶ entropy is a **measure** of **information** content of a distribution

    ▶ *"20 Questions on* $[0, 1)$*":* Land inside my interval by halving.

# Entropy

### Definition 7.2 (Entropy)
Given probabilities $p_1, \ldots, p_n$ (for outcomes $1, \ldots, n$ of a random variable), the *entropy* of the distribution is defined as

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

◄

▶ entropy is a **measure** of **information** content of a distribution

    ▶ *"20 Questions on $[0, 1)$":* Land inside my interval by halving.



$p_i = \frac{3}{4} - \frac{11}{16} = \frac{1}{16}$
$\rightsquigarrow \lg(1/p_i) = 4$

$\rightsquigarrow$ 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

20

## Entropy and Huffman codes

- would ideally encode value $i$ using $\lg(1/p_i)$ bits $\quad$ not for single code; but possible *on average*!

  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     <span style="font-size:small">not for single code; but possible *on average*!</span>

  not always possible; cannot use codeword of $1.5$ bits ... but:

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

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

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

# Entropy and Huffman codes

▶ would ideally encode value $i$ using $\lg(1/p_i)$ bits — not for single code; but possible *on average*!

not always possible; cannot use codeword of $1.5$ bits … but:

## Theorem 7.3 (Entropy bounds for Huffman codes)

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

$$\boxed{\mathcal{H} \;\le\; \ell(E) \;\le\; \mathcal{H} + 1}$$ where $\mathcal{H} = \mathcal{H}\left(\dfrac{w(a_1)}{W}, \ldots, \dfrac{w(a_\sigma)}{W}\right)$ and $W = w(a_1) + \cdots + w(a_\sigma)$. ◀

$$\underset{p_1}{\overset{!!}{=}} \qquad \underset{p_\sigma}{\overset{!!}{=}}$$

*Proof sketch:*

▶ $\ell(E) \ge \mathcal{H}$

Any prefix-free code $E$ induces weights $\underline{q_i = 2^{-|E(a_i)|}}$.

By *Kraft's Inequality*, we have $q_1 + \cdots + q_\sigma \le 1$.

Hence we can apply *Gibb's Inequality* to get

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

for any $q_i \in [0,1)$  $\sum q_i \le 1$

| $c$ | $a$ | $b$ | $c$ |
|---|---|---|---|
| $E(c)$ | $1$ | $00$ | $01$ |
| $q_c$ | $\frac{1}{2}$ | $\frac{1}{4}$ | $\frac{1}{4}$ |

## Entropy and Huffman codes [2]

*Proof sketch (continued):* $\lg\left(\frac{1}{p_i}\right) \hat{=}$ ideal codeword length

▶ $\ell(E) \le \mathcal{H} + 1$ → round up $\le \lg(1/p_i) + 1$

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

We construct a code $E'$ for $\Sigma$ with $|E'(a_i)| \le \lg(1/q_i)$ as follows;
w.l.o.g. assume $q_1 \le q_2 \le \cdots \le q_\sigma$ $\lceil \lg(1/p_i) \rceil$ $\frac{1}{16}$ $\frac{1}{8}$

▶ If $\sigma = 2$, $E'$ uses a single bit each.
Here, $q_i \le 1/2$, so $\lg(1/q_i) \ge 1 = |E'(a_i)|$ ✓



$\frac{1}{8}$

▶ If $\sigma \ge 3$, we merge $a_1$ and $a_2$ to $\boxed{a_1 a_2}$, 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 $\left|E'(\overline{a_1 a_2})\right| \le \lg\left(\dfrac{1}{2q_2}\right) = \lg\left(\dfrac{1}{q_2}\right) - 1.$

By construction, $|E'(a_1)| = |E'(a_2)| = \left|E'(\overline{a_1 a_2})\right| + 1$, so $|E'(a_1)| \le \lg(\frac{1}{q_1})$ and $|E'(a_2)| \le \lg(\frac{1}{q_2})$.

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

# 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`

**Click on "Polls" tab**

# 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 $\mathcal{H} < \lg(\sigma) - 1$ ✓

$\ell(E) \leq \mathcal{H} + 1 < \ell_S(\sigma) - 1 + 1$

$= \ell_S(\sigma) = \ell(E^{\text{fixed}})$

**E** ~~when $\mathcal{H} \approx 1$~~

`sli.do/comp526`

Click on "Polls" tab

## 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)     → Sedgewick Wayne
    - 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 coding – Discussion

- running time complexity:  $O(\sigma \log \sigma)$ to construct code
  - build PQ + $\sigma \cdot$ (2 deleteMins and 1 insert)
  - can do $\Theta(\sigma)$ time when characters already sorted by weight
  - time for encoding: $O(n + |C|)$

- 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 + $\sigma \cdot$ (2 deleteMins and 1 insert)
    - can do $\Theta(\sigma)$ time when characters already sorted by weight
    - time for encoding: $O(n + |C|)$

- many variations in use (tie-breaking rules, estimated frequencies, adaptive encoding, . . . )

👍 optimal prefix-free character encoding

👍 very fast decoding

👎 needs 2 passes over source text for encoding
    - one-pass variants possible, but more complicated

👎 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.
```

▶ 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.
```

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

⤳ Have to encode whole *phrases* of $S$ by a single codeword

# 7.4  Run-Length Encoding

# Run-Length encoding

▶ simplest form of repetition: *runs* of characters

same character repeated

```
00000000000000000000000000000000000000
00000000000000000000000000000000000000
00000000000000000000000000000000000000
000101100100000111111100000000001111110000
001111111110001111111110000000111111000
0011110110100011110001110000110000000
001100000000000000000011100111100000000
001100000000000000000110011110000000000
001101100000000000000011001100111110000
00111111110000000000011100111111111110000
00111011111000000000011100011111100111100
00000000011100000001110000111000001110
000000000111000000011000011100000011100
00000000011000001110000001110000001110
000000000111000111000000001110000011100
000000000111000111000000001110000011100
001101111110001111011101000011111111000
011111111110001111111111111000011111110000
00010110000000101001100100000010010000
00000000000000000000000000000000000000
00000000000000000000000000000000000000
```

▶ 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

```
00000000000000000000000000000000000000
00000000000000000000000000000000000000
00000000000000000000000000000000000000
00010101100100000111111100000000001111111000
00111111111110001111111110000001111111000
00111101101000011000111110000011100000000
00110000000000000000011001110000000000
00110000000000000000011001110000000000
00110110000000000000011001100111110000
00111111111000000000011100111111111000
00111011111100000000011100111111100111100
000000000011100000001110000111100000001110
000000000011100000011100001110000001100
000000000011000001110000000110000001110
000000000011000111000000001100000001100
000000000011000111000000001100000001110
00110101111111000111101110100001111111000
01111111111100011111111111100001111110000
00010110000000010100110010000000100100000
00000000000000000000000000000000000000
00000000000000000000000000000000000000
```

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

⇝ **run-length encoding (RLE):**
  use runs as phrases: $S = \underbrace{00000}\ \underbrace{111}\ \underbrace{0000}$

# Run-Length encoding

▶ simplest form of repetition: *runs* of characters

same character repeated

```
000000000000000000000000000000000000000
000000000000000000000000000000000000000
000000000000000000000000000000000000000
000101100100000111111100000000001111111000
001111111110001111111110000000111111111000
001111101101000111100011110000011100000000
001100000000000000000011000110000000000
001100000000000000000011001110000000000
001100000000000000000011001110000000000
001101100000000000000011001100111110000
001111111100000000000011100111111111000
001110111110000000000111000111111100111100
000000000111000000011110000111000000001110
000000000111000000011100011110000000001100
000000000111000000011000000011100000001110
000000000111000011100000000110000001100
000000000111000111000000000110000001110
000000000111000111000000000111000000011100
001101111110001111011010000111111111000
011111111000111111111111100001111110000
000101100000001010011001000000100100000
000000000000000000000000000000000000000
000000000000000000000000000000000000000
```

▶ here: only consider $\Sigma_S = \{0,1\}$ (work on a binary representation)

   ▶ can be extended for larger alphabets

⤳ **run-length encoding (RLE):**
   use runs as phrases: $S = \underbrace{00000}\ \underbrace{111}\ \underbrace{0000}$

⤳ 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

     same character repeated

```
000000000000000000000000000000000000000
000000000000000000000000000000000000000
000000000000000000000000000000000000000
000101100100000111111000000000011111000
001111111100011111111100000001111111000
001111011010001110001110000011100000000
001100000000000000001110011100000000000
001100000000000000001110011100000000000
001100000000000000001110011000000000000
001101100000000000000110011001111110000
001111111100000000001110011111111111000
001110111110000000011100011111100111100
000000000110000000111000011000000001110
000000000110000000110000011000000001100
000000000110000001100000011100000001110
000000000110000111100000011100000001110
000000000110001110000000011000000001110
000000000110001110000000011000000011100
001101111110001110111010000011111111000
011111111000111111111111110000111110000
000101100000001010011001000000100100000
000000000000000000000000000000000000000
000000000000000000000000000000000000000
```

▶ here: only consider $\Sigma_S = \{0, 1\}$    (work on a binary representation)

    ▶ can be extended for larger alphabets

⇝ **run-length encoding (RLE):**
   use runs as phrases: $S = \underbrace{00000}\ \underbrace{111}\ \underbrace{0000}$

⇝ 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*

**Click on "Polls" tab**

## Elias codes

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

## Elias codes

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

- ▶ But that's simple!     Just use *unary* **encoding**!
  $7 \mapsto 00000001$     $3 \mapsto 0001$     $0 \mapsto 1$     $30 \mapsto 000000000000000000000000000001$

## Elias codes

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

- ▶ But that's simple! Just use *unary* **encoding**!
  $7 \mapsto$ 00000001    $3 \mapsto$ 0001    $0 \mapsto$ 1    $30 \mapsto$ 000000000000000000000000000000001
  
  👎 Much too long
  - ▶ (wasn't the whole point of RLE to get rid of long runs??)

## Elias codes

- ▶ 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!
  $7 \mapsto$ `00000001`     $3 \mapsto$ `0001`     $0 \mapsto$ `1`     $30 \mapsto$ `0000000000000000000000000000001`

  👎 Much too long
  - ▶ (wasn't the whole point of RLE to get rid of long runs??)      `1 0 1 0`

- ▶ Refinement: *Elias gamma code*
  - ▶ Store the **length** $\ell$ of the binary representation in **unary**
  - ▶ Followed by the binary digits themselves

# Elias codes

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

▶ But that's simple!     Just use *unary* encoding!
  $7 \mapsto 00000001$     $3 \mapsto 0001$     $0 \mapsto 1$     $30 \mapsto 000000000000000000000000000000001$
  👎 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 $\ell \geq 1$, so store $\ell - 1$ instead
    ▶ binary representation always starts with $1$ ⤳ don't need terminating $1$ in unary
  ⤳ Elias gamma code = $\ell - 1$ zeros, followed by binary representation

  codeword length
  for number $k$
  $\leq 2 \lceil \lg k \rceil$

**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.

13

## Run-length encoding – Examples

▶ Encoding:
  $S = $ 11111110010000000000000000000011111111111

  $C = $ 1

▶ Decoding:
  $C = $ 00001101001001010

  $S = $

## Run-length encoding – Examples

▶ Encoding:
$S = \underbrace{1111111}_{7}00100000000000000000000011111111111$
$k = 7 \quad 00111$
$C = 1\underline{00111}$

▶ Decoding:
$C = 00001101001001010$

$S =$

## Run-length encoding – Examples

▶ Encoding:
$S = $ 1111111<span style="color:red">00</span>100000000000000000000011111111111
$k = 2$        010
$C = $ 100111<u>010</u>

▶ Decoding:
$C = $ 00001101001001010

$S = $

# Run-length encoding – Examples

▶ Encoding:
  $S = $ 1111111001000000000000000000011111111111
  $k = 1$
  $C = $ 1001110101

▶ Decoding:
  $C = $ 00001101001001010

  $S = $

## Run-length encoding – Examples

▶ Encoding:
$S = $ 11111110010000000000000000000011111111111
$k = 20$               10100
$C = $ 1001110101000010100

▶ Decoding:
$C = $ 00001101001001010

$S = $

# Run-length encoding – Examples

▶ Encoding:
  $S = $ 11111110010000000000000000000011111111111
  $k = 11$
  $C = $ 10011101010000101000001011

▶ Decoding:
  $C = $ 00001101001001010

  $S = $

## Run-length encoding – Examples

▶ Encoding:
  $S = $ `11111110010000000000000000000011111111111`

  $C = $ `10011101010000101000001011`

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

▶ Decoding:
  $C = $ `00001101001001010`

  $S = $

## Run-length encoding – Examples

▶ Encoding:

$S = $ `11111110010000000000000000000011111111111`

$C = $ `10011101010000101000001011`

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

▶ Decoding:

$C = $ `00001101001001010`

$S = $

## Run-length encoding – Examples

▶ Encoding:
$S$ = 11111111001000000000000000000000011111111111

$C$ = 10011101010000101000001011

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

▶ Decoding:
$C$ = 000001101001001010
$b$ = 0

$S$ =

## Run-length encoding – Examples

- ▶ Encoding:
  $S = $ `11111110010000000000000000000011111111111`

  $C = $ `10011101010000101000001011`

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

- ▶ Decoding:
  $C = $ `0000110100100101010`
  $b = 0$
  $\ell = 3 + 1$

  $S = $

## Run-length encoding – Examples

▶ Encoding:
$S = $ 11111110010000000000000000000011111111111

$C = $ 10011101010000101000001011

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

▶ Decoding:
$C = $ 00001101001001010
$b = 0$
$\ell = 3 + 1$
$k = 13$
$S = $ 0000000000000

## Run-length encoding – Examples

▶ Encoding:
  $S = 11111110010000000000000000000011111111111$

  $C = 10011101010000101000001011$

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

▶ Decoding:
  $C = 0000110100\underbrace{1001010}$
  $b = 1$
  $\ell = 2 + 1$
  $k =$
  $S = 0000000000000$

## Run-length encoding – Examples

▶ Encoding:

$S = $ 11111110010000000000000000000011111111111

$C = $ 100111010100000101000001011

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

▶ Decoding:

$C = $ 00001101001001010

$b = 1$

$\ell = 2 + 1$

$k = 4$

$S = $ 00000000000001111

## Run-length encoding – Examples

▶ Encoding:
$S = $ 11111110010000000000000000000011111111111

$C = $ 10011101010000101000001011

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

▶ Decoding:
$C = $ 00001101001001010
$b = 0$
$\ell = 0 + 1$
$k = $
$S = $ 00000000000001111

## Run-length encoding – Examples

▶ Encoding:
$S = $ 11111111001000000000000000000011111111111

$C = $ 10011101010000101000001011

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

▶ Decoding:
$C = $ 00001101001001010
$b = 0$
$\ell = 0 + 1$
$k = 1$
$S = $ 0000000000000011110

## Run-length encoding – Examples

▶ Encoding:
$S = $ 11111110010000000000000000000011111111111

$C = $ 1001110101000010100000101l

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

▶ Decoding:
$C = $ 00001101001001010
$b = 1$
$\ell = 1 + 1$
$k = $
$S = $ 000000000000011110

# Run-length encoding – Examples

▶ Encoding:
$S = $ 11111110010000000000000000000011111111111

$C = $ 10011101010000101000001011

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

▶ Decoding:
$C = $ 000011010010010<span style="color:red">10</span>
$b = 1$
$\ell = 1 + 1$
$k = 2$
$S = $ 000000000000001111<span style="color:red">011</span>

# 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)

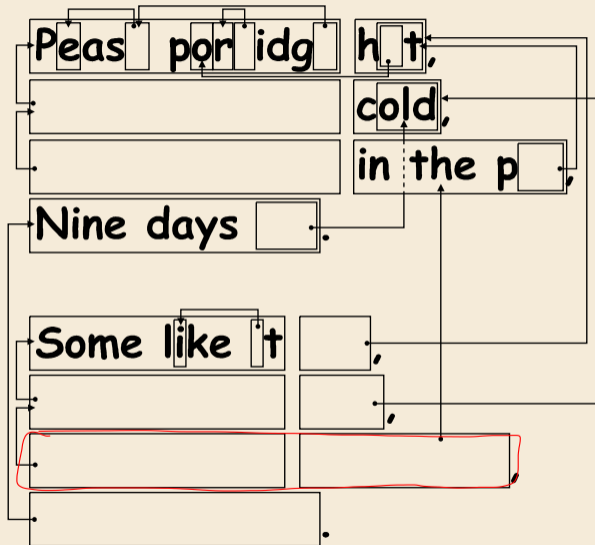▶ used in some image formats (e. g. TIFF)

👍 fairly simple and fast

👍 can compress $\underset{\sim}{n}$ 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 $k \geq 6$
  ▶ **expansion** for run length $k = 2$ or $6$

# 7.5  Lempel-Ziv-Welch

# Warmup

Peas  por idg  h t,
cold,
in the p ,
Nine days  .

Some like  t  ,
,
,
,

# Clicker Question

What is the second-to-last line of the above poem?

Click on "Polls" tab

## 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
      ▶ 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

    - ▶ 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: $k = 12$)  ⤳  fixed-length
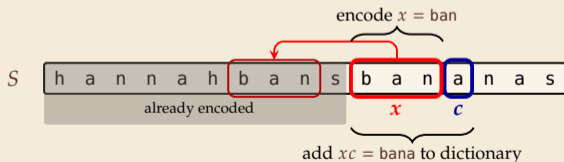  - but they represent a variable portion of the source text!

# 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: $k = 12$)   ⤳   fixed-length
  ▶ but they represent a variable portion of the source text!

▶ maintain a **dictionary** $D$ with $2^k$ entries   ⤳   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 $S$.



encode $x$ = ban

$S$   h a n n a h b a n s b a n a n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

⤳ new codeword in $D$

▶ $D$ actually stores codewords for $x$ and $c$, not the expanded string

32

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

$C$ =

$D$ =

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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

encode $x$ = ban

$S$ | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$  $c$

add $xc$ = bana to dictionary

## LZW encoding – Example

**Input**: `YO!␣YOU!␣YOUR␣YOYO!`                    $\Sigma_S$ = ASCII character set (0–127)

$$\text{\textcolor{red}{Y}}$$
$$C = 89$$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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

encode $x$ = ban

$S$ | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$ $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: Y̲O!␣YOU!␣YOUR␣YOYO! $\qquad\qquad\qquad \Sigma_S$ = ASCII character set (0–127)

Y
$C = 89$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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



encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: Y0!␣YOU!␣YOUR␣YOYO!  $\Sigma_S$ = ASCII character set (0–127)

$$C = \begin{array}{cc} \text{Y} & \text{0} \\ 89 & 79 \end{array}$$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | 0 |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$  $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

$$C = \begin{array}{cc} \text{Y} & \text{O} \\ 89 & 79 \end{array}$$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | |
| 131 | |
| 132 | |
| 133 | |
| 134 | |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |



encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: `Y0!␣YOU!␣YOUR␣YOYO!`  $\Sigma_S$ = ASCII character set (0–127)

```
      Y    0    !
C =  89   79   33
```

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | 0 |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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



encode $x$ = ban

$S$  h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: `YO!␣YOU!␣YOUR␣YOYO!`                    $\Sigma_S$ = ASCII character set (0–127)

$$C = \begin{matrix} \text{Y} & \text{O} & \text{!} \\ 89 & 79 & 33 \end{matrix}$$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | |
| 132 | |
| 133 | |
| 134 | |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |

$D =$

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: `Y0!␣YOU!␣YOUR␣YOYO!`                                            $\Sigma_S$ = ASCII character set (0–127)

$$C = \begin{array}{cccc} \text{Y} & \text{0} & \text{!} & \text{␣} \\ 89 & 79 & 33 & 32 \end{array}$$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | 0 |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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

encode $x$ = ban

$S$ = `h a n n a h b a n s b a n a n a s`

already encoded

$x$    $c$

add $xc$ = bana to dictionary

## LZW encoding – Example

**Input**: Y0!␣YOU!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

$$C = \begin{array}{cccc} \text{Y} & \text{0} & \text{!} & \text{␣} \\ 89 & 79 & 33 & 32 \end{array}$$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | 0 |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

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

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$  $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!

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

|       | Y  | O  | !  | ␣  | YO  |
|-------|----|----|----|----|-----|
| $C$ = | 89 | 79 | 33 | 32 | 128 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | |
| 133 | |
| 134 | |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

$S$  h a n n a h b a n s b a n a n a s

already encoded

$x$  $c$

add $xc$ = bana to dictionary

# LZW encoding – Example

**Input**: Y0!␣Y0U!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

|   | Y | 0 | ! | ␣ | Y0 |
|---|---|---|---|---|----|
| $C =$ | 89 | 79 | 33 | 32 | 128 |

$D =$

| Code | String |
|------|--------|
| ... | |
| 32 | ␣ |
| 33 | ! |
| ... | |
| 79 | 0 |
| ... | |
| 82 | R |
| ... | |
| 85 | U |
| ... | |
| 89 | Y |
| ... | |

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



encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!    $\Sigma_S$ = ASCII character set (0–127)

|  | Y | O | ! | ␣ | YO | U |
|---|---|---|---|---|---|---|
| $C$ = | 89 | 79 | 33 | 32 | 128 | 85 |

$D =$

| Code | String |
|---|---|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|---|---|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | |
| 134 | |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |



encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$   $c$

add $xc$ = bana to dictionary

## LZW encoding – Example

**Input**: YO!␣YOU|!␣YOUR␣YOYO!                                    $\Sigma_S$ = ASCII character set (0–127)

$$C = \begin{array}{ccccccc} \text{Y} & \text{O} & \text{!} & \text{␣} & \text{YO} & \text{U} \\ 89 & 79 & 33 & 32 & 128 & 85 \end{array}$$

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

$S$  h a n n a h b a n s b a n a n a s

already encoded

$x$ $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                                    $\Sigma_S$ = ASCII character set (0–127)

```
      Y    O    !    ␣    YO    U    !␣
C =  89   79   33   32   128   85   130
```

$D =$

| Code | String |
| --- | --- |
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
| --- | --- |
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s |

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                          $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ |
|---|---|---|---|---|----|---|----|
| $C =$ | 89 | 79 | 33 | 32 | 128 | 85 | 130 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |



encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s |

already encoded

$x$  $c$

add $xc$ = bana to dictionary

33

# LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!  $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ | YOU |
|---|---|---|---|---|----|---|----|-----|
| $C =$ | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | |
| 136 | |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$  $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

|       | Y  | O  | !  | ␣  | YO  | U  | !␣  | YOU |
|-------|----|----|----|----|-----|----|-----|-----|
| $C =$ | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 |

$D =$

| Code | String |
|------|--------|
| . . . |       |
| 32   | ␣      |
| 33   | !      |
| . . . |       |
| 79   | O      |
| . . . |       |
| 82   | R      |
| . . . |       |
| 85   | U      |
| . . . |       |
| 89   | Y      |
| . . . |       |

| Code | String |
|------|--------|
| 128  | YO     |
| 129  | O!     |
| 130  | !␣     |
| 131  | ␣Y     |
| 132  | YOU    |
| 133  | U!     |
| 134  | !␣Y    |
| 135  | YOUR   |
| 136  |        |
| 137  |        |
| 138  |        |
| 139  |        |

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s |

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                                   $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ | YOU | R |
|---|---|---|---|---|----|---|----|-----|---|
| C = | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

S | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$   $c$

add $xc$ = bana to dictionary

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!     $\Sigma_S$ = ASCII character set (0–127)

```
        Y    O    !    ␣    YO   U    !␣   YOU  R
C  =   89   79   33   32  128   85  130  132   82
```

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

$S$ ⟨ h a n n a h b a n s b a n a n a s ⟩

already encoded      $x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!  $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ | YOU | R | ␣Y |
|---|---|---|---|---|----|---|----|-----|---|----|
| $C$ = | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | |
| 138 | |
| 139 | |

encode $x$ = ban

$S$  h a n n a h b a n s b a n a n a s
  already encoded      $x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ | YOU | R | ␣Y |
|---|---|---|---|---|----|---|----|-----|---|----|
| $C$ = | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | ␣YO |
| 138 | |
| 139 | |

encode $x$ = ban

$S$  | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$  $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                    $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ | YOU | R | ␣Y | O |
|---|---|---|---|---|----|---|----|-----|---|----|---|
| $C$ = | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 | 79 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | ␣YO |
| 138 | |
| 139 | |

encode $x$ = ban

$S$ | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$   $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YO<u>YO</u>!  $\Sigma_S$ = ASCII character set (0–127)

| | Y | O | ! | ␣ | YO | U | !␣ | YOU | R | ␣Y | O |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $C$ = | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 | 79 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | ␣YO |
| 138 | OY |
| 139 | |



encode $x$ = ban

$S$ | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$  $c$

add $xc$ = bana to dictionary

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!  $\qquad \Sigma_S$ = ASCII character set (0–127)

|  | Y | O | ! | ␣ | YO | U | !␣ | YOU | R | ␣Y | O | YO |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $C =$ | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 | 79 | 128 |

$D =$

| Code | String |
|---|---|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|---|---|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | ␣YO |
| 138 | OY |
| 139 | |

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s

already encoded

$x$  $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!                                        $\Sigma_S$ = ASCII character set (0–127)

|       | Y  | O  | !  | ␣  | YO  | U  | !␣  | YOU | R  | ␣Y  | O  | YO  |
|-------|----|----|----|----|-----|----|-----|-----|----|-----|----|-----|
| $C =$ | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 | 79 | 128 |

$D =$

| Code | String |
|------|--------|
| . . . |       |
| 32   | ␣      |
| 33   | !      |
| . . . |       |
| 79   | O      |
| . . . |       |
| 82   | R      |
| . . . |       |
| 85   | U      |
| . . . |       |
| 89   | Y      |
| . . . |       |

| Code | String |
|------|--------|
| 128  | YO     |
| 129  | O!     |
| 130  | !␣     |
| 131  | ␣Y     |
| 132  | YOU    |
| 133  | U!     |
| 134  | !␣Y    |
| 135  | YOUR   |
| 136  | R␣     |
| 137  | ␣YO    |
| 138  | OY     |
| 139  | YO!    |

encode $x$ = ban

$S$ | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already encoded

$x$    $c$

add $xc$ = bana to dictionary

33

## LZW encoding – Example

**Input**: YO!␣YOU!␣YOUR␣YOYO!␣ ⅄ I                    $\Sigma_S$ = ASCII character set (0–127)

|   | Y | O | ! | ␣ | YO | U | !␣ | YOU | R | ␣Y | O | YO | ! |
|---|---|---|---|---|----|---|----|-----|---|----|---|----|---|
| $C =$ | 89 | 79 | 33 | 32 | 128 | 85 | 130 | 132 | 82 | 131 | 79 | 128 | 33 |

$D =$

| Code | String |
|------|--------|
| . . . | |
| 32 | ␣ |
| 33 | ! |
| . . . | |
| 79 | O |
| . . . | |
| 82 | R |
| . . . | |
| 85 | U |
| . . . | |
| 89 | Y |
| . . . | |

| Code | String |
|------|--------|
| 128 | YO |
| 129 | O! |
| 130 | !␣ |
| 131 | ␣Y |
| 132 | YOU |
| 133 | U! |
| 134 | !␣Y |
| 135 | YOUR |
| 136 | R␣ |
| 137 | ␣YO |
| 138 | OY |
| 139 | YO! |

encode $x$ = ban

$S$ | h a n n a h | b a n | s | b a n | a | n a s |

already encoded

$x$   $c$

add $xc$ = bana to dictionary



33

## LZW encoding – Code

```
 1 procedure LZWencode(S[0..n])
 2     x := ε // previous phrase, initially empty
 3     C := ε // output, initially empty
 4     D := dictionary, initialized with codes for c ∈ Σ_S // stored as trie
 5     k := |Σ_S| // next free codeword
 6     for i := 0, . . . , n − 1 do
 7         c := S[i]
 8         if D.containsKey(xc) then
 9             x := xc
10         else
11             C := C · D.get(x) // append codeword for x
12             D.put(xc, k) // add xc to D, assigning next free codeword
13             k := k + 1; x := c
14     end for
15     C := C · D.get(x)
16     return C
```

## 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 $y$)

decode $y$ = an

| $S$ | h | a | n | n | a | h | b | a | n | s | b | a | n | a | n | a | s |

already decoded          $x$        $y$

$c$

add $xc$ = bana to dictionary

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

## LZW decoding – Example

► Same idea: build dictionary while reading string.

► Example: 67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ⊔ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

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

## LZW decoding – Example

► Same idea: build dictionary while reading string.

► Example: 67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| · · · | |
| 32 | ␣ |
| · · · | |
| · · · | |
| 65 | A |
| 66 | B |
| 67 | C |
| · · · | |
| 78 | N |
| · · · | |
| 83 | S |
| · · · | |

| input | decodes to | Code # | String (human) | String (computer) |
|-------|------------|--------|----------------|-------------------|
| 67 | C | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# LZW decoding – Example

▶ Same idea: build dictionary while reading string.

▶ Example: 67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ␣ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

| input | decodes to | | Code # | String (human) | String (computer) |
|-------|-----------|---|--------|----------------|-------------------|
| 67 | C | | | | |
| 65 | A | | 128 | CA | 67, A |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

# LZW decoding – Example

- Same idea: build dictionary while reading string.

- Example: 67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| ⋯ | |
| 32 | ␣ |
| ⋯ | |
| ⋯ | |
| 65 | A |
| 66 | B |
| 67 | C |
| ⋯ | |
| 78 | N |
| ⋯ | |
| 83 | S |
| ⋯ | |

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

# LZW decoding – Example

- Same idea: build dictionary while reading string.

- Example:  67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ␣ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

| input | decodes to | Code # | String (human) | String (computer) |
|-------|-----------|--------|----------------|-------------------|
| 67 | C | | | |
| 65 | A | 128 | CA | 67, A |
| 78 | N | 129 | AN | 65, N |
| 32 | ␣ | 130 | N␣ | 78, ␣ |
| | | | | |
| | | | | |
| | | | | |

# LZW decoding – Example

▶ Same idea: build dictionary while reading string.

▶ Example: 67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| ... | |
| 32 | ␣ |
| ... | |
| ... | |
| 65 | A |
| 66 | B |
| 67 | C |
| ... | |
| 78 | N |
| ... | |
| 83 | S |
| ... | |

| input | decodes to | Code # | String (human) | String (computer) |
|-------|-----------|--------|----------------|-------------------|
| 67 | C | | | |
| 65 | A | 128 | CA | 67, A |
| 78 | N | 129 | AN | 65, N |
| 32 | ␣ | 130 | N␣ | 78, ␣ |
| 66 | B | 131 | ␣B | 32, B |
| | | | | |
| | | | | |

# LZW decoding – Example

▶ Same idea: build dictionary while reading string.

▶ Example:   67 65 78 32 66 129 133

$D =$

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ␣ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

| input | decodes to | Code # | String (human) | String (computer) |
|-------|------------|--------|----------------|-------------------|
| 67 | C | | | |
| 65 | A | 128 | CA | 67, A |
| 78 | N | 129 | AN | 65, N |
| 32 | ␣ | 130 | N␣ | 78, ␣ |
| 66 | B | 131 | ␣B | 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

$D =$

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ␣ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

| input | decodes to | Code # | String (human) | String (computer) |
|-------|-----------|--------|----------------|-------------------|
| 67 | C | | | |
| 65 | A | 128 | CA | 67, A |
| 78 | N | 129 | AN | 65, N |
| 32 | ␣ | 130 | N␣ | 78, ␣ |
| 66 | B | 131 | ␣B | 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

$D =$

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ␣ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

| input | decodes to | Code # | Str (hu... | |
|-------|-----------|--------|------|------|
| 67 | C | | | |
| 65 | A | 128 | CA | 67, A |
| 78 | N | 129 | AN | 65, N |
| 32 | ␣ | 130 | N␣ | 78, ␣ |
| 66 | B | 131 | ␣B | 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:
  ▶ 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 $D[k]$
  ⤳ $D[k] = xc = x \cdot x[0]$    (all known)

## LZW decoding – Code

```
 1  procedure LZWdecode(C[0..m])
 2      D := dictionary [0..2^d) → Σ_S^+, initialized with codes for c ∈ Σ_S // stored as array
 3      k := |Σ_S| // next unused codeword
 4      q := C[0] // first codeword
 5      y := D[q] // lookup meaning of q in D
 6      S := y // output, initially first phrase
 7      for j := 1, . . . , m − 1 do
 8          x := y // remember last decoded phrase
 9          q := C[j] // next codeword
10          if q == k then
11              y := x · x[0] // bootstrap case
12          else
13              y := D[q]
14          S := S · y // append decoded phrase
15          D[k] := x · y[0] // store new phrase
16          k := k + 1
17      end for
18      return S
```

# LZW decoding – Example continued

▶ Example: 67 65 78 32 66 129 133 83

| Code # | String |
|--------|--------|
| . . . | |
| 32 | ␣ |
| . . . | |
| . . . | |
| 65 | A |
| 66 | B |
| 67 | C |
| . . . | |
| 78 | N |
| . . . | |
| 83 | S |
| . . . | |

$D =$

| input | decodes to | Code # | String (human) | String (computer) |
|-------|-----------|--------|----------------|-------------------|
| 67 | C | | | |
| 65 | A | 128 | CA | 67, A |
| 78 | N | 129 | AN | 65, N |
| 32 | ␣ | 130 | N␣ | 78, ␣ |
| 66 | B | 131 | ␣B | 32, B |
| 129 | AN | 132 | BA | 66, A |
| 133 | ANA | 133 | ANA | 129, A |
| 83 | S | 134 | ANAS | 133, S |

last step  $cy=D[k]$

C A N ␣ B A N A N A S

done    $x$   $c$

$D[k]:=xc$

A N A

$x$   $c$

*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?

**A** $\sim n$

**B** $\sim n/2$

**C** $\sim n/4$

**D** $\Theta(n/\log n)$

**E** $\Theta(\sqrt{n})$

**F** $\Theta(\log n)$

**G** $\Theta(\log \log n)$

**H** 2

**I** 1

`sli.do/comp526`

Click on "Polls" tab

# Clicker Question

$\underline{a}\ \underline{a\ a}\ \underline{a}\ \underline{a\ a}\ \underline{a\ a}\ \ \underline{a}\ \underline{c\ a}$

$a\ a \qquad a\ a\ a \qquad a\ a\ a\ g \qquad \Rightarrow k$ phrases

$n = \sum_{i=1}^{k} i = \dfrac{k(k+1)}{2} \sim \dfrac{k^2}{2}$

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

**A** ~~$n$~~

**B** ~~$n/2$~~

**C** ~~$n/4$~~

**D** ~~$\Theta(n/\log n)$~~

**E** $\Theta(\sqrt{n})$ ✓

**F** ~~$\Theta(\log n)$~~ ← RLE

**G** ~~$\Theta(\log\log n)$~~

**H** ~~2~~

**I** ~~1~~

`sli.do/comp526`

Click on "Polls" tab

## LZW – Discussion

▶ As presented, LZW uses coded alphabet $\Sigma_C = [0..2^d)$.

  $\rightsquigarrow$ use another encoding for    code numbers $\mapsto$ binary,    e. g., Huffman

▶ need a rule when dictionary is full; different options:

  ▶ increment $d$  $\rightsquigarrow$  longer codewords
  ▶ "flush" dictionary and start from scratch  $\rightsquigarrow$  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 $\Sigma_C = [0..2^d)$.

  ⇝ use another encoding for    code numbers $\mapsto$ binary,    e. g., Huffman

▶ need a rule when dictionary is full; different options:

  ▶ increment $d$  ⇝  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 $\leadsto$ RLE
    - frequently used characters $\leadsto$ Huffman
    - many (local) repeated substrings $\leadsto$ LZW

- but methods can be frustratingly "blind" to other "obvious" redundancies
    - LZW: repetition too distant ⚡ dictionary already flushed
    - Huffman: changing probabilities (local clusters) ⚡ averaged out globally
    - RLE: run of alternating pairs of characters ⚡ not a run

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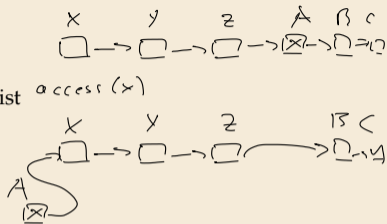
- ▶ compression is effective is we have one the following:
  - ▶ long runs ⤳ RLE
  - ▶ frequently used characters ⤳ Huffman
  - ▶ many (local) repeated substrings ⤳ LZW

- ▶ but methods can be frustratingly "blind" to other "obvious" redundancies
  - ▶ LZW: repetition too distant ⚡ 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

# 7.6 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)
  ⤳ list "learns" probabilities of access to objects
    makes access to frequently requested ones cheaper

▶ Here: use such a list for storing source alphabet $\Sigma_S$

  ▶ to encode $c$, access it in list
  ▶ encode $c$ using its (old) position in list
  ▶ then apply MTF to the list
  ⤳ 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)
    - ⤳ list "learns" probabilities of access to objects
      makes access to frequently requested ones cheaper

- ▶ Here: use such a list for storing source alphabet $\Sigma_S$

    - ▶ to encode $c$, access it in list
    - ▶ encode $c$ using its (old) position in list
    - ▶ then apply MTF to the list
    - ⤳ codewords are integers, i.e., $\Sigma_C = [0..\sigma)$

- ⤳ clusters of few characters ⤳ many small numbers

# Clicker Question

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

## MTF – Code

▶ **Transform (encode):**

```
1  procedure MTF−encode(S[0..n))
2      L := list containing Σ_S (sorted order)
3      C := ε
4      for i := 0, . . . , n − 1 do
5          c := S[i]
6          p := position of c in L
7          C := C · p
8          Move c to front of L
9      end for
10     return C
```

▶ **Inverse transform (decode):**

```
1  procedure MTF−decode(C[0..m))
2      L := list containing Σ_S (sorted order)
3      S := ε
4      for j := 0, . . . , m − 1 do
5          p := C[j]
6          c := character at position p in L
7          S := S · c
8          Move c to front of L
9      end for
10     return S
```

▶ Important: encoding and decoding produce same accesses to list

## MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S =$ INEFFICIENCIES

$C =$

## MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| A | B | C | D | E | F | G | H | I | J | K  | L  | M  | N  | O  | P  | Q  | R  | S  | T  | U  | V  | W  | X  | Y  | Z  |

$S =$ INEFFICIENCIES

$C = 8$

# MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| I | A | B | C | D | E | F | G | H | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S = $ INEFFICIENCIES

$C = $ 8 13

# MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| N | I | A | B | C | D | E | F | G | H | J | K | L | M | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S = $ INEFFICIENCIES

$C = $ 8 13 6

## MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| E | N | I | A | B | C | D | F | G | H | J | K | L | M | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S =$ INE**F**FICIENCIES

$C =$ 8 13 6 7

# MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| F | E | N | I | A | B | C | D | G | H | J | K | L | M | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S = $ INEFFICIENCIES

$C = $ 8 13 6 7 0

# MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| F | E | N | I | A | B | C | D | G | H | J | K | L | M | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S =$ INEFF**I**CIENCIES

$C =$ 8 13 6 7 0 3

## MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| I | F | E | N | A | B | C | D | G | H | J | K | L | M | O | P | Q | R | S | T | U | V | W | X | Y | Z |

$S = $ INEFFICIENCIES

$C = $ 8 13 6 7 0 3 6

# MTF – Example

```
 0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
 C  I  F  E  N  A  B  D  G  H  J  K  L  M  O  P  Q  R  S  T  U  V  W  X  Y  Z
```

$S =$ INEFFIC**I**ENCIES

$C = 8\ 13\ 6\ 7\ 0\ 3\ 6\ 1$

45

# MTF – Example

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| S | E | I | C | N | F | A | B | D | G | H | J | K | L | M | O | P | Q | R | T | U | V | W | X | Y | Z |

$$S = \texttt{INEFFICIENCIES}$$

$$C = 8\ 13\ 6\ 7\ \underline{0}\ 3\ 6\ 1\ 3\ 4\ \underline{3\ 3\ 3}\ 18$$

▶ What does a run in $S$ encode to in $C$?   → 0 s ℓ

▶ What does a run in $C$ mean about the source $S$?

45

# 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)

## 7.7 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:

$T = \mathtt{time\_flies\_quickly\_}$          $\mathtt{flies\_quickly\_time\_}$



⤳ cyclic shift

## BWT transform

▶ *cyclic shift* of a string:

▶ add *end-of-word character* $ to $S$ (as in Unit 6)

⇝ can recover original string

$T =$ time␣flies␣quickly␣



⇝ cyclic shift

flies␣quickly␣time␣

## BWT transform

- *cyclic shift* of a string:

- add *end-of-word character* $ to *S* (as in Unit 6)

- ⤳ can recover original string

$T = \texttt{time}_\sqcup\texttt{flies}_\sqcup\texttt{quickly}_\sqcup$



⤳ cyclic shift

$\texttt{flies}_\sqcup\texttt{quickly}_\sqcup\texttt{time}_\sqcup$



- 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

$S = \text{alf\_eats\_alfalfa\$}$

**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
```

⤳
sort

## 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
```

$\rightsquigarrow$ sort

```
$alf␣eats␣alfalfa
␣alfalfa$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
```

49

# BWT transform – Example

BWT
↓

$S = $ alf␣eats␣alfalfa$

1. Write all cyclic shifts

2. Sort cyclic shifts

3. Extract last column

$B = $ asff$f␣e␣lllaaata

| | |
|---|---|
| alf␣eats␣alfalfa$ | $alf␣eats␣alfalf**a** |
| lf␣eats␣alfalfa$a | ␣alfalfa$alf␣eat**s** |
| f␣eats␣alfalfa$al | ␣eats␣alfalfa$al**f** |
| ␣eats␣alfalfa$alf → | ⓐ$alf␣eats␣alfal**f** |
| eats␣alfalfa$alf␣ | alf␣eats␣alfalfa**$** |
| ats␣alfalfa$alf␣e | alfa$alf␣eats␣al**f** |
| ts␣alfalfa$alf␣ea ⤳ | alfalfa$alf␣eats**␣** |
| s␣alfalfa$alf␣eat   sort | ats␣alfalfa$alf␣**e** |
| ␣alfalfa$alf␣eats | eats␣alfalfa$alf**␣** |
| alfalfa$alf␣eats␣ | f␣eats␣alfalfa$a**l** |
| lfalfa$alf␣eats␣a | fa$alf␣eats␣alfa**l** |
| falfa$alf␣eats␣al | falfa$alf␣eats␣a**l** |
| alfa$alf␣eats␣alf | lf␣eats␣alfalfa$**a** |
| lfa$alf␣eats␣alfa | lfa$alf␣eats␣alf**a** |
| fa$alf␣eats␣alfal | lfalfa$alf␣eats␣**a** |
| a$alf␣eats␣alfalf | s␣alfalfa$alf␣ea**t** |
| $alf␣eats␣alfalfa | ts␣alfalfa$alf␣e**a** |

49

# Clicker Question

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

**A** $L$ can be very easily computed from $B$ and $T$

**B** $B$ can be very easily computed from $L$ and $T$

**C** Both A and B

**D** Neither A nor B

`sli.do/comp526`

Click on "Polls" tab

# Clicker Question

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

**A** ~~$L$ can be very easily computed from $B$ and $T$~~

**B** $B$ can be very easily computed from $L$ and $T$ ✓

**C** ~~Both A and B~~

**D** ~~Neither A nor B~~

`sli.do/comp526`

**Click on "Polls" tab**

## BWT – Implementation & Properties

**Compute BWT efficiently:**

▶ cyclic shifts $S$ $\hat{=}$ suffixes of $S$

▶ BWT is essentially suffix sorting!

    ▶ $B[i] = S[L[i] - 1]$     ($L$ = suffix array!)
      (if $L[i] = 0$, $B[i] = \$$)

   ⤳ Can compute $B$ in $O(n)$ time

| | | | $r$ | | $\downarrow L[r]$ |
|---|---|---|---|---|---|
| alf␣eats␣alfalfa$ | | | 0 | $alf␣eats␣alfalf**a** | 16 |
| lf␣eats␣alfalfa$a | | | 1 | ␣alfalfa$alf␣eat**s** | 8 |
| f␣eats␣alfalfa$al | | | 2 | ␣eats␣alfalfa$al**f** | 3 |
| ␣eats␣alfalfa$alf | | | 3 | a$alf␣eats␣alfal**f** | 15 |
| eats␣alfalfa$alf␣ | | | 4 | alf␣eats␣alfalfa**$** | 0 |
| ats␣alfalfa$alf␣e | | | 5 | alfa$alf␣eats␣al**f** | 12 |
| ts␣alfalfa$alf␣ea | | | 6 | alfalfa$alf␣eats**␣** | 9 |
| s␣alfalfa$alf␣eat | | | 7 | ats␣alfalfa$alf␣**e** | 5 |
| ␣alfalfa$alf␣eats | | | 8 | eats␣alfalfa$alf**␣** | 4 |
| alfalfa$alf␣eats␣ | | | 9 | f␣eats␣alfalfa$a**l** | 2 |
| lfalfa$alf␣eats␣a | | | 10 | fa$alf␣eats␣alfa**l** | 14 |
| falfa$alf␣eats␣al | | | 11 | falfa$alf␣eats␣a**l** | 11 |
| alfa$alf␣eats␣alf | | | 12 | lf␣eats␣alfalfa$**a** | 1 |
| lfa$alf␣eats␣alfa | | | 13 | lfa$alf␣eats␣alf**a** | 13 |
| fa$alf␣eats␣alfal | | | 14 | lfalfa$alf␣eats␣**a** | 10 |
| a$alf␣eats␣alfalf | | | 15 | s␣alfalfa$alf␣ea**t** | 7 |
| $alf␣eats␣alfalfa | | | 16 | ts␣alfalfa$alf␣e**a** | 6 |

# BWT – Implementation & Properties

**Compute BWT efficiently:**

▶ cyclic shifts $S$ $\widehat{=}$ suffixes of $S$

▶ BWT is essentially suffix sorting!

  ▶ $B[i] = S[L[i] - 1]$     ($L$ = suffix array!)
  
  (if $L[i] = 0$, $B[i] = \$$)

  ↝ Can compute $B$ in $O(n)$ time

**Why does BWT help?**

▶ sorting groups characters *by what follows*

  ▶ Example: lf always preceded by a

↝ $B$ has local clusters of characters

  ▶ that makes MTF effective

▶ repeated substring in $S$ ↝ *runs* of characters in $B$

  ▶ picked up by RLE

| | $r$ | | $\downarrow L[r]$ |
|---|---|---|---|
| alf␣eats␣alfalfa$ | 0 | $alf␣eats␣alfalf**a** | 16 |
| lf␣eats␣alfalfa$a | 1 | ␣alfalfa$alf␣eat**s** | 8 |
| f␣eats␣alfalfa$al | 2 | ␣eats␣alfalfa$al**f** | 3 |
| ␣eats␣alfalfa$alf | 3 | a$alf␣eats␣alfal**f** | 15 |
| eats␣alfalfa$alf␣ | 4 | alf␣eats␣alfalfa**$** | 0 |
| ats␣alfalfa$alf␣e | 5 | alfa$alf␣eats␣al**f** | 12 |
| ts␣alfalfa$alf␣ea | 6 | alfalfa$alf␣eats**␣** | 9 |
| s␣alfalfa$alf␣eat | 7 | ats␣alfalfa$alf␣**e** | 5 |
| ␣alfalfa$alf␣eats | 8 | eats␣alfalfa$alf**␣** | 4 |
| alfalfa$alf␣eats␣ | 9 | f␣eats␣alfalfa$a**l** | 2 |
| lfalfa$alf␣eats␣a | 10 | fa$alf␣eats␣alfa**l** | 14 |
| falfa$alf␣eats␣al | 11 | falfa$alf␣eats␣a**l** | 11 |
| alfa$alf␣eats␣alf | 12 | lf␣eats␣alfalfa$**a** | 1 |
| lfa$alf␣eats␣alfa | 13 | lfa$alf␣eats␣alf**a** | 13 |
| fa$alf␣eats␣alfal | 14 | lfalfa$alf␣eats␣**a** | 10 |
| a$alf␣eats␣alfalf | 15 | s␣alfalfa$alf␣ea**t** | 7 |
| $alf␣eats␣alfalfa | 16 | ts␣alfalfa$alf␣e**a** | 6 |

# Bigger Example

```
have␣had␣hadnt␣hasnt␣havent␣has␣what$      $have␣had␣hadnt␣hasnt␣havent␣has␣what
ave␣had␣hadnt␣hasnt␣havent␣has␣what$h      ␣had␣hadnt␣hasnt␣havent␣has␣what$have
ve␣had␣hadnt␣hasnt␣havent␣has␣what$ha      ␣hadnt␣hasnt␣havent␣has␣what$have␣had
e␣had␣hadnt␣hasnt␣havent␣has␣what$hav      ␣has␣what$have␣had␣hadnt␣hasnt␣havent
␣had␣hadnt␣hasnt␣havent␣has␣what$have      ␣hasnt␣havent␣has␣what$have␣had␣hadnt
had␣hadnt␣hasnt␣havent␣has␣what$have␣      ␣havent␣has␣what$have␣had␣hadnt␣hasnt
ad␣hadnt␣hasnt␣havent␣has␣what$have␣h      ␣what$have␣had␣hadnt␣hasnt␣havent␣has
d␣hadnt␣hasnt␣havent␣has␣what$have␣ha      ad␣hadnt␣hasnt␣havent␣has␣what$have␣h
␣hadnt␣hasnt␣havent␣has␣what$have␣had      adnt␣hasnt␣havent␣has␣what$have␣had␣h
hadnt␣hasnt␣havent␣has␣what$have␣had␣      as␣what$have␣had␣hadnt␣hasnt␣havent␣h
adnt␣hasnt␣havent␣has␣what$have␣had␣h      asnt␣havent␣has␣what$have␣had␣hadnt␣h
dnt␣hasnt␣havent␣has␣what$have␣had␣ha      at$have␣had␣hadnt␣hasnt␣havent␣has␣wh
nt␣hasnt␣havent␣has␣what$have␣had␣had      ave␣had␣hadnt␣hasnt␣havent␣has␣what$h
t␣hasnt␣havent␣has␣what$have␣had␣hadn      avent␣has␣what$have␣had␣hadnt␣hasnt␣h
␣hasnt␣havent␣has␣what$have␣had␣hadnt      d␣hadnt␣hasnt␣havent␣has␣what$have␣ha
hasnt␣havent␣has␣what$have␣had␣hadnt␣      dnt␣hasnt␣havent␣has␣what$have␣had␣ha
asnt␣havent␣has␣what$have␣had␣hadnt␣h      e␣had␣hadnt␣hasnt␣havent␣has␣what$hav
snt␣havent␣has␣what$have␣had␣hadnt␣ha      ent␣has␣what$have␣had␣hadnt␣hasnt␣hav
nt␣havent␣has␣what$have␣had␣hadnt␣has      had␣hadnt␣hasnt␣havent␣has␣what$have␣
t␣havent␣has␣what$have␣had␣hadnt␣hasn      hadnt␣hasnt␣havent␣has␣what$have␣had␣
␣havent␣has␣what$have␣had␣hadnt␣hasnt      has␣what$have␣had␣hadnt␣hasnt␣havent␣
havent␣has␣what$have␣had␣hadnt␣hasnt␣      hasnt␣havent␣has␣what$have␣had␣hadnt␣
avent␣has␣what$have␣had␣hadnt␣hasnt␣h      hat$have␣had␣hadnt␣hasnt␣havent␣has␣w
vent␣has␣what$have␣had␣hadnt␣hasnt␣ha      have␣had␣hadnt␣hasnt␣havent␣has␣what$
ent␣has␣what$have␣had␣hadnt␣hasnt␣hav      havent␣has␣what$have␣had␣hadnt␣hasnt␣
nt␣has␣what$have␣had␣hadnt␣hasnt␣have      nt␣has␣what$have␣had␣hadnt␣hasnt␣have
t␣has␣what$have␣had␣hadnt␣hasnt␣haven      nt␣hasnt␣havent␣has␣what$have␣had␣had
␣has␣what$have␣had␣hadnt␣hasnt␣havent      nt␣havent␣has␣what$have␣had␣hadnt␣has
has␣what$have␣had␣hadnt␣hasnt␣havent␣      s␣what$have␣had␣hadnt␣hasnt␣havent␣ha
as␣what$have␣had␣hadnt␣hasnt␣havent␣h      snt␣havent␣has␣what$have␣had␣hadnt␣ha
s␣what$have␣had␣hadnt␣hasnt␣havent␣ha      t$have␣had␣hadnt␣hasnt␣havent␣has␣wha
␣what$have␣had␣hadnt␣hasnt␣havent␣has      t␣has␣what$have␣had␣hadnt␣hasnt␣haven
what$have␣had␣hadnt␣hasnt␣havent␣has␣      t␣hasnt␣havent␣has␣what$have␣had␣hadn
hat$have␣had␣hadnt␣hasnt␣havent␣has␣w      t␣havent␣has␣what$have␣had␣hadnt␣hasn
at$have␣had␣hadnt␣hasnt␣havent␣has␣wh      ve␣had␣hadnt␣hasnt␣havent␣has␣what$ha
t$have␣had␣hadnt␣hasnt␣havent␣has␣wha      vent␣has␣what$have␣had␣hadnt␣hasnt␣ha
$have␣had␣hadnt␣hasnt␣havent␣has␣what      what$have␣had␣hadnt␣hasnt␣havent␣has␣
```

$T$ = h a v e ␣ h a d ␣ h a d n t ␣ h a s n t ␣ h a v e n t ␣ h a s ␣ w h a t $

$B$ = t e d t t t s h h h h h h h a a v v ␣ ␣ ␣ ␣ ␣ w $ ␣ e d s a a a n n n a a ␣

$MTF(B)$ = 8 5 5 2 **0 0** 8 7 **0 0 0 0 0 0** 7 **0** 9 **0** 8 **0 0 0** 10 9 2 9 9 8 7 **0 0** 10 **0 0** 1 **0** 5

## Clicker Question

Consider $T =$ `have␣had␣hadnt␣hasnt␣havent␣has␣what$`.
The BWT is $B =$ `tedtttshhhhhhhaavv␣␣␣␣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 `a`s are preceded by `h`

**F**   `h` is the 4th character in the alphabet

`sli.do/comp526`

Click on "Polls" tab

# Clicker Question

Consider $T =$ `have␣had␣hadnt␣hasnt␣havent␣has␣what$`.
The BWT is $B =$ `tedtttshhhhhhhaavv␣␣␣␣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 `a`s are preceded by `h` ✓

**F** ~~`h` is the 4th character in the alphabet~~

`sli.do/comp526`

Click on "Polls" tab

## 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!

# 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!

▶ **"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)

## Inverse BWT

▶ 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) |
| 3 | ($, 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) |

## Inverse BWT

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

|   | *D* |   |   | sorted *D* |   |
|---|---|---|---|---|---|
|   |   |   |   | char | next |
| 0 | (a, | 0) | 0 | ($, | 3) |
| 1 | (r, | 1) | 1 | (a, | 0) |
| 2 | (d, | 2) | 2 | (a, | 6) |
| 3 | ($, | 3) | 3 | (a, | 7) |
| 4 | (r, | 4) | 4 | (a, | 8) |
| 5 | (c, | 5) | 5 | (a, | 9) |
| 6 | (a, | 6) | 6 | (b, | 10) |
| 7 | (a, | 7) | 7 | (b, | 11) |
| 8 | (a, | 8) | 8 | (c, | 5) |
| 9 | (a, | 9) | 9 | (d, | 2) |
| 10 | (b, | 10) | 10 | (r, | 1) |
| 11 | (b, | 11) | 11 | (r, | 4) |

▶ **"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:**
$B = \text{ard\$rcaaaabb}$
$S =$

52

# 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!

|     |   $D$    |
|-----|----------|
| 0   | (a, 0)   |
| 1   | (r, 1)   |
| 2   | (d, 2)   |
| 3   | ($, 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)  |

sorted $D$

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

▶ **"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:**
$B = \text{ard\$rcaaaabb}$
$S = \text{a}$

52

# Inverse BWT

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

▶ **"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:**
$B = \text{ard\$rcaaaabb}$
$S = \text{ab}$

|   | $D$ |
|---|-----|
| 0 | (a, 0) |
| 1 | (r, 1) |
| 2 | (d, 2) |
| 3 | ($, 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) |

|   | sorted $D$ |
|---|-----|
|   | char  next |
| 0 | ($, 3) |
| 1 | (a, 0) |
| 2 | (a, 6) |
| 3 | (a, 7) |
| 4 | (a, 8) |
| 5 | (a, 9) |
| 6 | (b, 10) |
| 7 | (b, 11) |
| 8 | (c, 5) |
| 9 | (d, 2) |
| 10 | (r, 1) |
| 11 | (r, 4) |

## Inverse BWT

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

|  | $D$ |  |  | sorted $D$ |  |
|---|---|---|---|---|---|
|  |  |  |  | char | next |
| 0 | (a, | 0) | 0 | ($, | 3) |
| 1 | (r, | 1) | 1 | (a, | 0) |
| 2 | (d, | 2) | 2 | (a, | 6) |
| 3 | ($, | 3) | 3 | (a, | 7) |
| 4 | (r, | 4) | 4 | (a, | 8) |
| 5 | (c, | 5) | 5 | (a, | 9) |
| 6 | (a, | 6) | 6 | (b, | 10) |
| 7 | (a, | 7) | 7 | (b, | 11) |
| 8 | (a, | 8) | 8 | (c, | 5) |
| 9 | (a, | 9) | 9 | (d, | 2) |
| 10 | (b, | 10) | 10 | (r, | 1) |
| 11 | (b, | 11) | 11 | (r, | 4) |

► **"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 = $ abr

## 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!*

|  | $D$ | | sorted $D$ | |
|---|---|---|---|---|
|  |  |  | char | next |
| 0 | (a, 0) | 0 | ($, 3) | |
| 1 | (r, 1) | 1 | (a, 0) | |
| 2 | (d, 2) | 2 | (a, 6) | |
| 3 | ($, 3) | 3 | (a, 7) | |
| 4 | (r, 4) | 4 | (a, 8) | |
| 5 | (c, 5) | 5 | (a, 9) | |
| 6 | (a, 6) | 6 | (b, 10) | |
| 7 | (a, 7) | 7 | (b, 11) | |
| 8 | (a, 8) | 8 | (c, 5) | |
| 9 | (a, 9) | 9 | (d, 2) | |
| 10 | (b, 10) | 10 | (r, 1) | |
| 11 | (b, 11) | 11 | (r, 4) | |

▶ **"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 = $ abra

52

# Inverse BWT

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

|   | $D$ | | sorted $D$ | |
|---|---|---|---|---|
|   |   |   | char | next |
| 0 | (a, 0) | 0 | ($, 3) |
| 1 | (r, 1) | 1 | (a, 0) |
| 2 | (d, 2) | 2 | (a, 6) |
| 3 | ($, 3) | 3 | (a, 7) |
| 4 | (r, 4) | 4 | (a, 8) |
| 5 | (c, 5) | 5 | (a, 9) |
| 6 | (a, 6) | 6 | (b, 10) |
| 7 | (a, 7) | 7 | (b, 11) |
| 8 | (a, 8) | 8 | (c, 5) |
| 9 | (a, 9) | 9 | (d, 2) |
| 10 | (b, 10) | 10 | (r, 1) |
| 11 | (b, 11) | 11 | (r, 4) |

▶ **"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 = $ abrac

52

## 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!

|  | $D$ | | sorted $D$ | |
|---|---|---|---|---|
|  |  |  | char | next |

► **"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 = $ abraca

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

52

## 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!

|   | $D$ | | sorted $D$ |
|---|---|---|---|
|   |   |   | char  next |
| 0 | (a,  0) | 0 | ($, 3) |
| 1 | (r,  1) | 1 | (a,  0) |
| 2 | (d,  2) | 2 | (a,  6) |
| 3 | ($,  3) | 3 | (a,  7) |
| 4 | (r,  4) | 4 | (a,  8) |
| 5 | (c,  5) | 5 | (a,  9) |
| 6 | (a,  6) | 6 | (b, 10) |
| 7 | (a,  7) | 7 | (b, 11) |
| 8 | (a,  8) | 8 | (c,  5) |
| 9 | (a,  9) | 9 | (d,  2) |
| 10 | (b, 10) | 10 | (r,  1) |
| 11 | (b, 11) | 11 | (r,  4) |

▶ **"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$rcaaabb
$S = $ abracad

## Inverse BWT

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

|   | $D$ |   |   | sorted $D$ |   |
|---|---|---|---|---|---|
|   |   |   |   | char | next |
| 0 | (a, 0) |   | 0 | ($, 3) |   |
| 1 | (r, 1) |   | 1 | (a, 0) |   |
| 2 | (d, 2) |   | 2 | (a, 6) |   |
| 3 | ($, 3) |   | 3 | (a, 7) |   |
| 4 | (r, 4) |   | 4 | (a, 8) |   |
| 5 | (c, 5) |   | 5 | (a, 9) |   |
| 6 | (a, 6) |   | 6 | (b, 10) |   |
| 7 | (a, 7) |   | 7 | (b, 11) |   |
| 8 | (a, 8) |   | 8 | (c, 5) |   |
| 9 | (a, 9) |   | 9 | (d, 2) |   |
| 10 | (b, 10) |   | 10 | (r, 1) |   |
| 11 | (b, 11) |   | 11 | (r, 4) |   |

▶ **"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 = $ abracada

52

# 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!

▶ **"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:**
$B = \text{ard\$rcaaaabb}$
$S = \text{abracadab}$

|     | $D$       |     | sorted $D$   |
| --- | --------- | --- | ------------ |
|     |           |     | char  next   |
| 0   | (a, 0)    | 0   | (\$, 3)      |
| 1   | (r, 1)    | 1   | (a, 0)       |
| 2   | (d, 2)    | 2   | (a, 6)       |
| 3   | (\$, 3)   | 3   | (a, 7)       |
| 4   | (r, 4)    | 4   | (a, 8)       |
| 5   | (c, 5)    | 5   | (a, 9)       |
| 6   | (a, 6)    | 6   | (b, 10)      |
| 7   | (a, 7)    | 7   | (b, 11)      |
| 8   | (a, 8)    | 8   | (c, 5)       |
| 9   | (a, 9)    | 9   | (d, 2)       |
| 10  | (b, 10)   | 10  | (r, 1)       |
| 11  | (b, 11)   | 11  | (r, 4)       |

52

## 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!*

|  | $D$ | | sorted $D$ | |
|---|---|---|---|---|
|  | | | | char   next |
| **▶ "Magic" solution:** | 0 | (a, 0) | 0 | ($, 3) |
| **1.** Create array $D[0..n]$ of pairs: | 1 | (r, 1) | 1 | (a, 0) |
| $D[r] = (B[r], r)$. | 2 | (d, 2) | 2 | (a, 6) |
| **2.** Sort $D$ *stably* with | 3 | ($, 3) | 3 | (a, 7) |
| respect to *first entry*. | 4 | (r, 4) | 4 | (a, 8) |
| **3.** Use $D$ as linked list with | 5 | (c, 5) | 5 | (a, 9) |
| (char, next entry) | 6 | (a, 6) | 6 | (b, 10) |
|  | 7 | (a, 7) | 7 | (b, 11) |
| **Example:** | 8 | (a, 8) | 8 | (c, 5) |
| $B = \text{ard\$rcaaaabb}$ | 9 | (a, 9) | 9 | (d, 2) |
| $S = \text{abracadabr}$ | 10 | (b, 10) | 10 | (r, 1) |
|  | 11 | (b, 11) | 11 | (r, 4) |

52

## 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!

▶ **"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 = $ abracadabra

|   | $D$ |
|---|---|
| 0 | (a, 0) |
| 1 | (r, 1) |
| 2 | (d, 2) |
| 3 | (\$, 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) |

|   | sorted $D$ |
|---|---|
|   | char  next |
| 0 | (\$, 3) |
| 1 | (a, 0) |
| 2 | (a, 6) |
| 3 | (a, 7) |
| 4 | (a, 8) |
| 5 | (a, 9) |
| 6 | (b, 10) |
| 7 | (b, 11) |
| 8 | (c, 5) |
| 9 | (d, 2) |
| 10 | (r, 1) |
| 11 | (r, 4) |

52

# 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!*

|     | $D$ |   |     | sorted $D$ |       |
|-----|-----|---|-----|------------|-------|
|     |     |   |     | char       | next  |
| 0   | (a, | 0) | 0   | ($,        | 3)    |
| 1   | (r, | 1) | 1   | (a,        | 0)    |
| 2   | (d, | 2) | 2   | (a,        | 6)    |
| 3   | ($, | 3) | 3   | (a,        | 7)    |
| 4   | (r, | 4) | 4   | (a,        | 8)    |
| 5   | (c, | 5) | 5   | (a,        | 9)    |
| 6   | (a, | 6) | 6   | (b,        | 10)   |
| 7   | (a, | 7) | 7   | (b,        | 11)   |
| 8   | (a, | 8) | 8   | (c,        | 5)    |
| 9   | (a, | 9) | 9   | (d,        | 2)    |
| 10  | (b, | 10) | 10 | (r,        | 1)    |
| 11  | (b, | 11) | 11 | (r,        | 4)    |

▶ **"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 = $ abracadabra\$

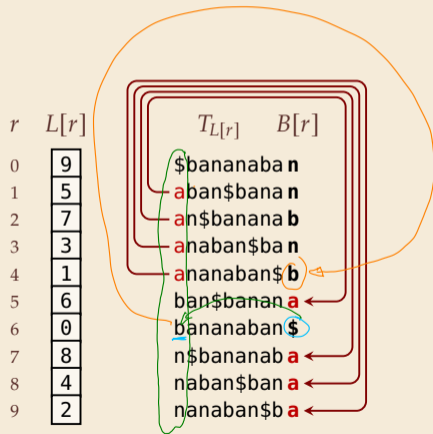# Inverse BWT – The magic revealed

- ▶ Inverse BWT very easy to compute:
    - ▶ only sort individual characters in $B$   (not suffixes)
    - ⤳ $O(n)$ 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)
  ⤳ $O(n)$ with counting sort

▶ *but why does this work!?*

▶ decode char by char
  ▶ can find unique \$  ⤳ 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

# Inverse BWT – The magic revealed

▶ Inverse BWT very easy to compute:
  ▶ only sort individual characters in $B$ (not suffixes)
  ⤳ $O(n)$ with counting sort

▶ *but why does this work!?*

▶ decode char by char
  ▶ can find unique $ ⤳ 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

- 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 $ $\rightsquigarrow$ 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
- for (i): first column = characters of $B$ in sorted order ✓
- for (ii): relative order of same character stays same:
  $i$th a in first column = $i$th a in BWT
  - $\rightsquigarrow$ stably sorting $(B[r], r)$ by first entry enough ✓

| $r$ | $L[r]$ | $T_{L[r]}$ | $B[r]$ |
|---|---|---|---|
| 0 | 9 | \$bananaba | **n** |
| 1 | 5 | aban\$bana | **n** |
| 2 | 7 | an\$banana | **b** |
| 3 | 3 | anaban\$ba | **n** |
| 4 | 1 | ananaban\$ | **b** |
| 5 | 6 | ban\$banan | **a** |
| 6 | 0 | bananaban | **\$** |
| 7 | 8 | n\$bananab | **a** |
| 8 | 4 | naban\$ban | **a** |
| 9 | 2 | nanaban\$b | **a** |

## BWT – Discussion

▶ Running time: $\Theta(n)$
  ▶ **encoding** uses suffix sorting
  ▶ decoding only needs counting sort
  ⇝ 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 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