

# CS 240 – Data Structures and Data Management

## Module 10: Data Compression

O. Veksler

Based on lecture notes by many previous cs240 instructors

David R. Cheriton School of Computer Science, University of Waterloo

Winter 2024

# Outline

- Data Compression
  - Background
  - Single-Character Encodings
  - Huffman Codes
  - Run-Length Encoding
  - Lempel-Ziv-Welch
  - Combining Compression Schemes: bzip2
  - Burrows-Wheeler Transform

# Outline

- **Data Compression**
  - **Background**
  - Single-Character Encodings
  - Huffman Codes
  - Run-Length Encoding
  - Lempel-Ziv-Welch
  - Combining Compression Schemes: bzip2
  - Burrows-Wheeler Transform

# Data Compression Introduction

- **The problem:** How to store and transmit data efficiently?
- **Source text:**
  - original data, string  $S$  of characters from **source alphabet**  $\Sigma_S$
- **Coded text**
  - encoded data, string  $C$  of characters from **coded alphabet**  $\Sigma_C$
- **Encoding [scheme]**
  - algorithm mapping source text to coded text
- **Decoding [scheme]**
  - algorithm mapping coded text back to original source text



- Source “text” can be any sort of data (not always text)
- Usually the coded alphabet is binary  $\Sigma_C = \{0,1\}$
- Consider lossless compression: exact recovery of  $S$  from  $C$

# Judging Encoding Schemes

- **Main objective:** for data compression, want to minimize the size of the coded text
- Measure the **compression ratio**

$$\frac{|C| \cdot \log|\Sigma_C|}{|S| \cdot \log|\Sigma_S|}$$

- **Examples:**

$(73)_{10} \rightarrow (1001001)_2$       compression ratio  $\frac{7 \cdot \log 2}{2 \cdot \log 10} \approx 1.05$       X

$(127)_{10} \rightarrow (7F)_{16}$       compression ratio  $\frac{2 \cdot \log 16}{3 \cdot \log 10} \approx 0.8$       ✓

- Want to achieve compression ration better than 1
  - smaller than 1
- Can achieve compression ratio of 1 by sending  $S$  without changes

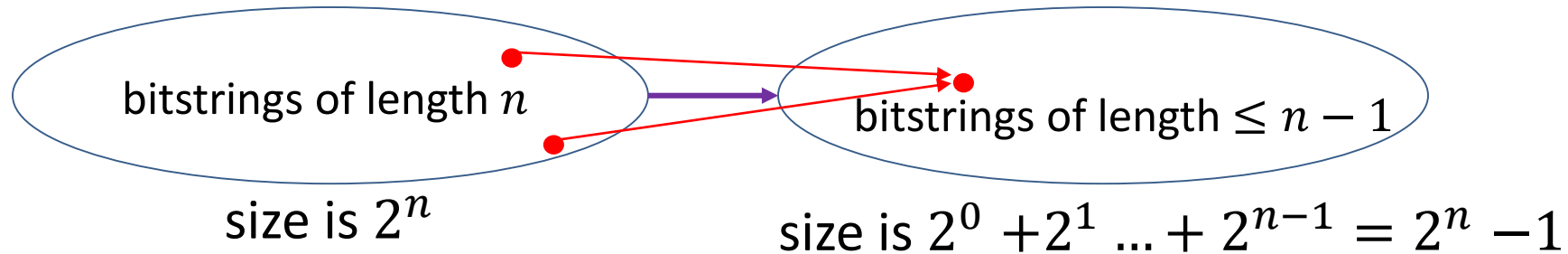
# Judging Encoding Schemes

- Also measure efficiency of encoding/decoding algorithms, as for any usual algorithm
  - always need time  $\Omega(|S| + |C|)$ 
    - sometimes need more time
- Other possible goals, not studied in this course
  - reliability (e.g. error-correcting codes)
  - security (e.g. encryption)

# Impossibility of Compressing

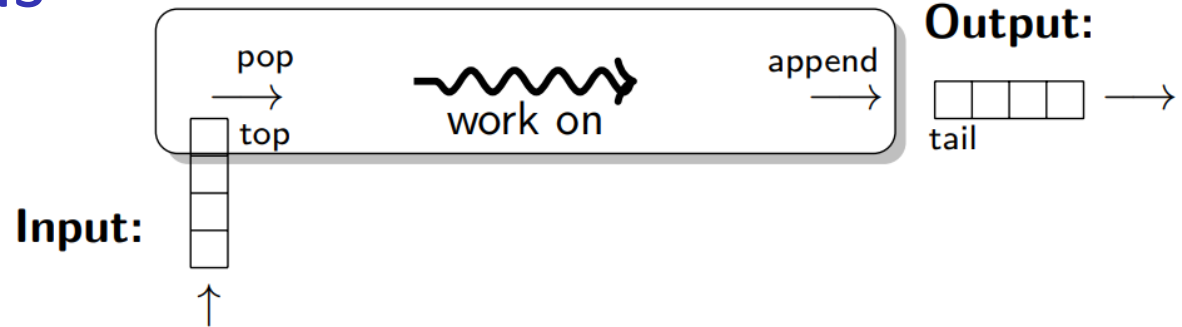
- **Observation:** No lossless encoding scheme can have compression ratio  $< 1$  for **all** input strings
- **Proof:** (for  $\Sigma_S = \Sigma_C = \{0,1\}$ , by contradiction)

Fix  $n$ , and assume all length  $n$  strings get shorter



- So impossible to provide good worst-case compression bounds
- However real-life data is usually far from random, it has some patterns that occur more frequently
  - can design compression schemes that work well for frequently occurring patterns

# Detour: Streams



- Usually texts are huge and do not fit into computer memory
- Therefore usually store *S* and *C* as *streams*
  - input stream
    - read one character at a time
      - `pop( ), top( ), isEmpty()`
      - sometimes need `reset()` to start processing from the start
    - output stream
      - write one character at a time
        - `append( ), isEmpty()`
- Advantage of streams
  - can start processing text while it is still being loaded
  - avoids needing to hold the entire text in memory at once



# Outline

- **Data Compression**
  - Background
  - **Single-Character Encodings**
  - Huffman Codes
  - Run-Length Encoding
  - Lempel-Ziv-Welch
  - Combining Compression Schemes: bzip2
  - Burrows-Wheeler Transform

# Character Encodings

- **Definition:** character encoding  $E$  (or single-character encoding) maps each *character* in the source alphabet to a *string* in coded alphabet

$$E : \Sigma_S \rightarrow \Sigma_C^*$$

- for  $c \in \Sigma_S$ ,  $E(c)$  is called the *codeword* (or *code*) of  $c$
- Two possibilities
  1. **Fixed-length code:** all codewords have the same length
    - compression ratio  $\geq 1$  (not good)
  2. **Variable-length code:** codewords may have different lengths

# Fixed Length Character Encoding

- Example: ASCII (American Standard Code for Information Interchange), 1963

char in $\Sigma_S$	null	start of heading	start of text	...	0	1	...	A	B	...	~	delete
code	0	1	2	...	48	49	...	65	66	...	126	127
code in binary	0000000	0000001	0000010		0110000	0110001		01000001	01000010		1111110	1111111

- Each codeword  $E(c)$  has length 7 bits
- Encoding/Decoding is easy: just concatenate/decode the next 7 bits
  - A P P L E  $\leftrightarrow$  (65, 80, 80, 76, 69)  $\leftrightarrow$  01000001 1010000 1010000 1001100 1000101
    - here  $|S| = 5, |C| = 5 \cdot 7, |\Sigma_S| = 128$
- Standard in all computers and often our source alphabet
- Other (earlier) fixed-length codes: Baudot code, Murray code
- Fixed-length codes do not compress
  - let  $|E(c)| = b$  and assume binary coded alphabet

$$\frac{|C| \cdot \log |\Sigma_C|}{|S| \cdot \log |\Sigma_S|} = \frac{|C|}{b \cdot |C| \cdot \log 2^b} = 1$$

# Better Idea: Variable-Length Codes

- **Observation:** Some alphabet letters occur more often than others
  - example: frequency of letters in typical English text

<b>e</b>	12.70%	<b>d</b>	4.25%	<b>p</b>	1.93%
<b>t</b>	9.06%	<b>l</b>	4.03%	<b>b</b>	1.49%
<b>a</b>	8.17%	<b>c</b>	2.78%	<b>v</b>	0.98%
<b>o</b>	7.51%	<b>u</b>	2.76%	<b>k</b>	0.77%
<b>i</b>	6.97%	<b>m</b>	2.41%	<b>j</b>	0.15%
<b>n</b>	6.75%	<b>w</b>	2.36%	<b>x</b>	0.15%
<b>s</b>	6.33%	<b>f</b>	2.23%	<b>q</b>	0.10%
<b>h</b>	6.09%	<b>g</b>	2.02%	<b>z</b>	0.07%
<b>r</b>	5.99%	<b>y</b>	1.97%		

- **Idea:** use shorter codes for more frequent characters
  - as before, map source alphabet to codewords:  $E : \Sigma_S \rightarrow \Sigma_C^*$
  - but not all codewords have the same length
  - this could make the coded text shorter



# Encoding

- Assume we have some character encoding  $E: \Sigma_S \rightarrow \Sigma_C^*$
- $E$  is a dictionary with keys in  $\Sigma_S$

*singleChar::Encoding*( $E, S, C$ )

$E$ : encoding dictionary,  $S$ : input stream with characters in  $\Sigma_S$

$C$ : output stream

**while**  $S$  is non-empty

$w \leftarrow E.\text{search}(S.\text{pop}())$

append each bit of  $w$  to  $C$

- Using dictionary below, encode **ANENT**  $\rightarrow$  **01 001 000 0100111**

$c \in \Sigma_S$		U	A	E	N	O	T
$E(c)$		000	01	101	001	100	11

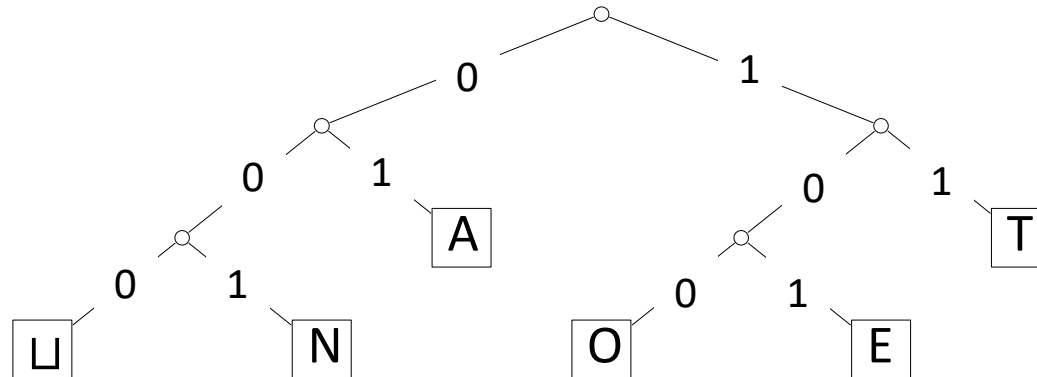
# Decoding

- The **decoding algorithm** must map  $\Sigma_C^*$  to  $\Sigma_S$
- The code must be *uniquely decodable*
  - false for Morse code as described

A	● ■■	U	● ● ■■
B	■■ ● ● ●	V	● ● ● ■■
C	■■ ● ■■ ●	W	● ■■ ■■
D	■■ ● ●	X	■■ ● ● ■■
E	●	Y	■■ ● ■■ ■■



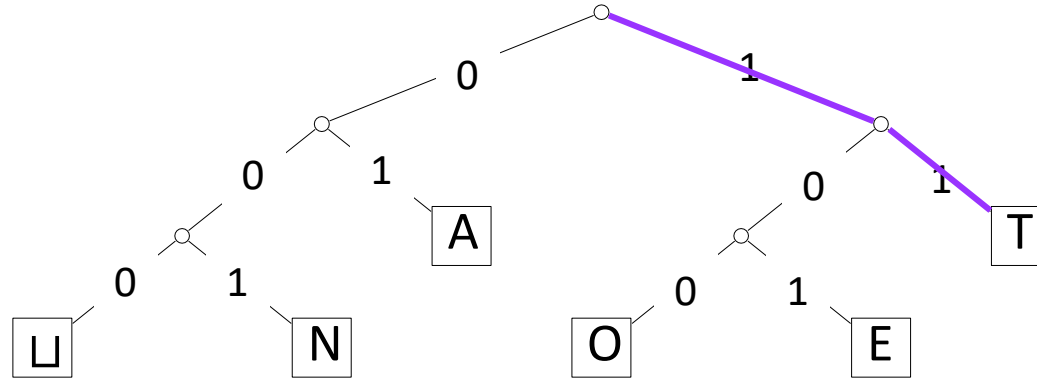
- Morse code uses 'end of character' pause to avoid ambiguity
- From now on only consider **prefix-free** codes  $E$ 
  - no codeword is a prefix of another codeword
- Store codes in a **trie** with characters of  $\Sigma_S$  at the leaves



- Do not need symbol \$, codewords are prefix-free by definition

# Example: Prefix-free Decoding

- Decode from a trie

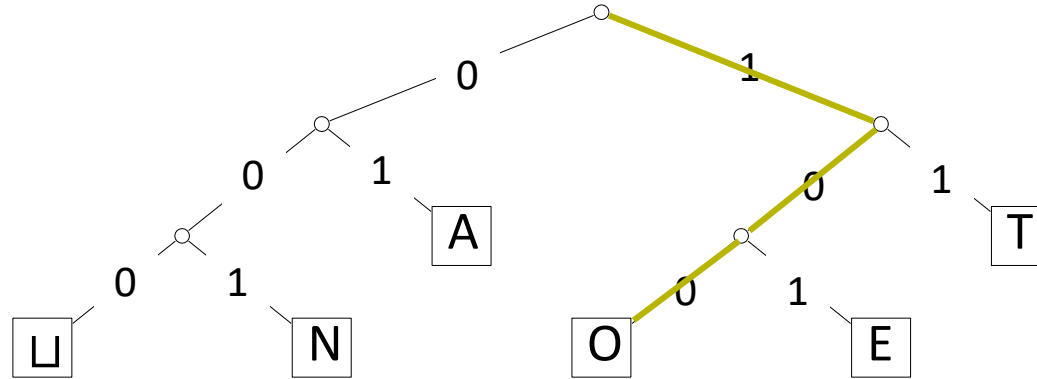


- Decode **111000001010111** → **T**



# Example: Prefix-free Decoding

- Decode from a trie

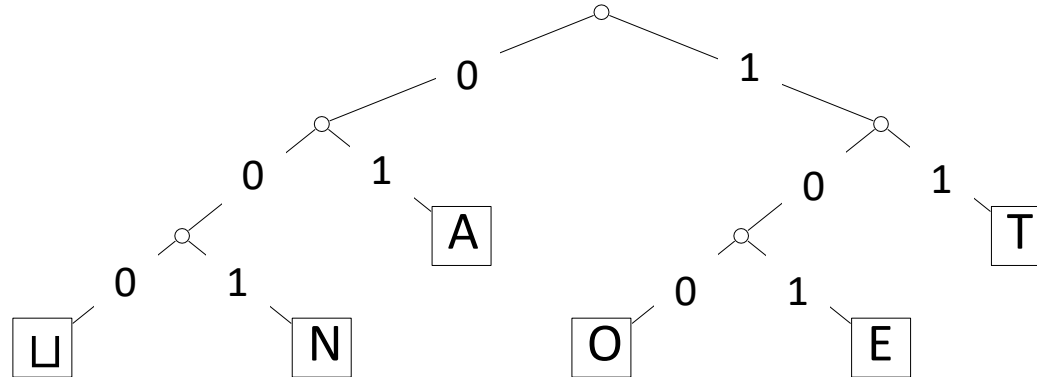


- Decode **111000001010111** → **T****O**



# Example: Prefix-free Decoding

- Decode from a trie



- Decode **11100000**1010111  $\rightarrow$  **T****O****L****E****A****T**
- Run-time:  $O(|C|)$

# Decoding of Prefix-Free Codes

- Any prefix-free code is uniquely decodable

```
prefixFree::decoding(T, C, S)
```

*T*: trie of a prefix-free code, *C*: input-stream with characters in  $\Sigma_C$

*S*: output-stream

```
while C is non-empty // iterate over all codewords
```

```
  z  $\leftarrow$  T.root
```

```
  while z is not a leaf // read next codeword
```

```
    if C is empty or z has no child labelled C.top()
```

```
      return "invalid encoding"
```

```
    z  $\leftarrow$  child of z that is labelled with C.pop()
```

```
  S.append(character stored at z)
```

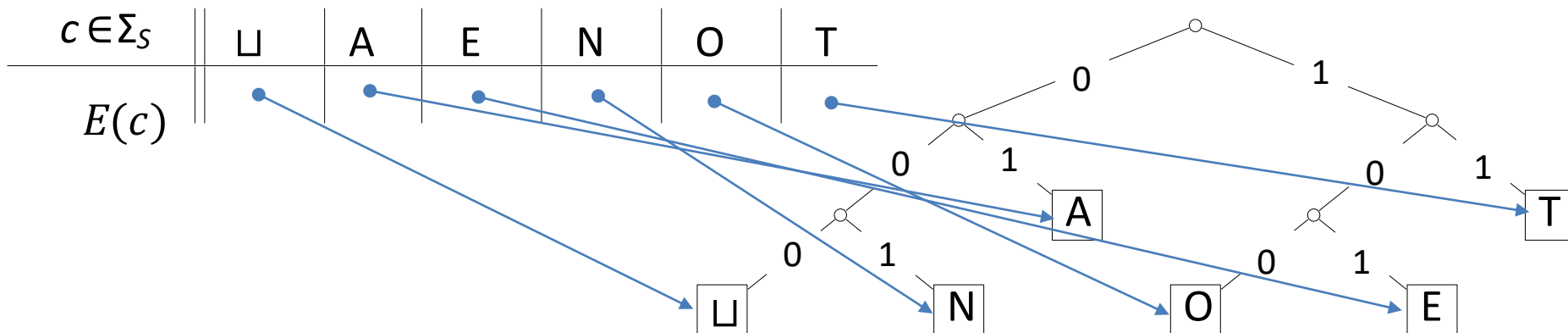
- Run-time:  $O(|C|)$
- Detects if the encoding is invalid

# Encoding from the Trie

- Explained previously how to encode from a table

$c \in \Sigma_S$	␣	A	E	N	O	T
$E(c)$	000	01	101	001	100	11

- Table wastes space, codewords can be quite long
- Better idea: store codewords via links to the trie leaves



# Encoding from the Trie

- Can encode directly from the trie  $T$

$prefixFree::encoding(T, S, C)$

$T$  : prefix-free code trie,  $S$ : input-stream with characters in  $\Sigma_S$

$E \leftarrow$  array of nodes in  $T$  indexed by  $\Sigma_S$

for all leaves  $l$  in  $T$

$E[\text{character at } l] \leftarrow l$

while  $S$  is non-empty

$w \leftarrow$  empty bitstring;  $v \leftarrow E[S.pop()]$

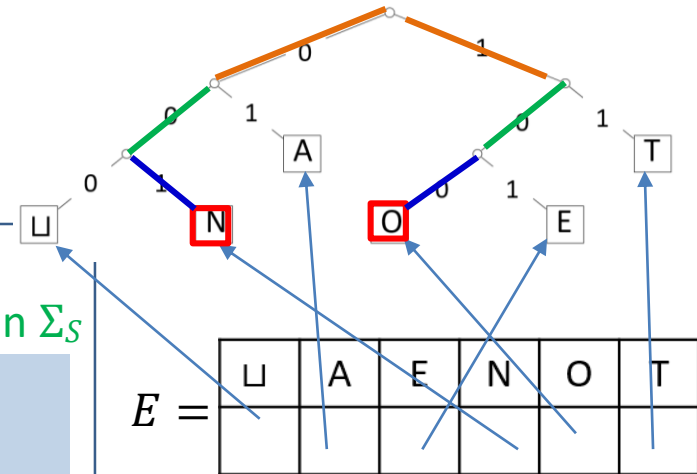
while  $v$  is not the root

$w.prepend$  (character from  $v$  to its parent)

$v \leftarrow$  parent( $v$ )

// now  $w$  is the encoding of  $S$

append each bit  $w$  of to  $C$



$E =$

□	A	E	N	O	T

$S = ON$

$i = 0$  (letter  $O$ )

$w = \Lambda$

$w = 0$

$w = 00$

$w = 100$

$C = 100$

$i = 1$  (letter  $N$ )

$w = \Lambda$

$w = 1$

$w = 01$

$w = 001$

$C = 100\ 001$

- Run-time:  $O(|T| + |C|)$

- have to visit all trie nodes, and insert leaves into  $E$

- if  $T$  has no nodes with one child

- #leaves  $- 1 \geq$  #internal nodes

- $|\Sigma_S| \cdot 2 - 1 \geq$  #internal nodes + #leaves =  $|T|$

- $O(|\Sigma_S| + |C|)$

# Outline

- **Data Compression**
  - Background
  - Single-Character Encodings
  - **Huffman Codes**
  - Run-Length Encoding
  - Lempel-Ziv-Welch
  - Combining Compression Schemes: bzip2
  - Burrows-Wheeler Transform

# Huffman's Algorithm: Building the Best Trie

- How to determine the best trie for a given source text  $S$  ?
  - i.e. trie giving shortest  $|C|$
- Idea: infrequent characters should be far down in the trie



# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies  
 $G: 2, R: 2, E: 4, N: 2, Y: 1$
- Put each character into its own trie (single node, height 0)
  - each trie has a frequency
  - initially, frequency is equal to its character frequency

2  
G

2  
R

4  
E

2  
N

1  
Y

# Example: Huffman Tree Construction

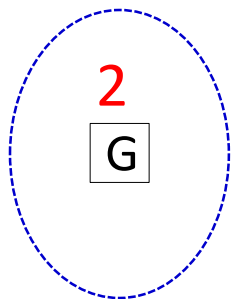
- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$

- Calculate character frequencies

$G: 2, R: 2, E: 4, N: 2, Y: 1$

- Join two least frequent tries into a new trie

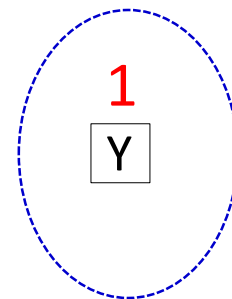
- frequency of the new trie = sum of old trie frequencies



2  
R

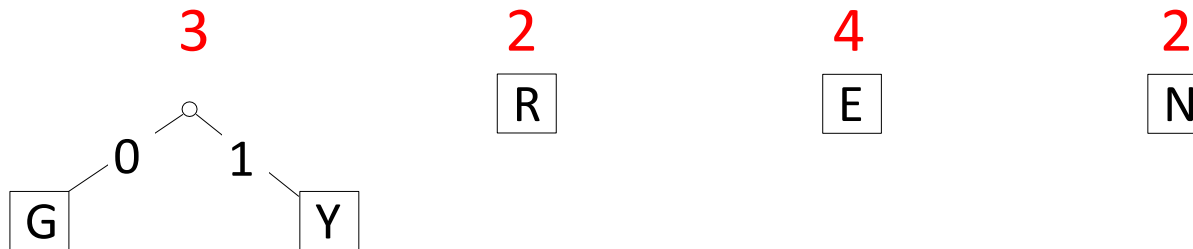
4  
E

2  
N



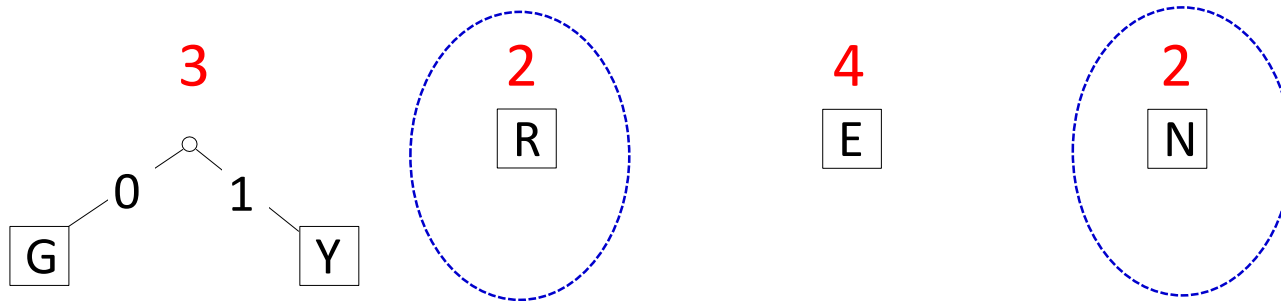
# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies  
 $G: 2, R: 2, E: 4, N: 2, Y: 1$
- Join two least frequent tries into a new trie
  - frequency of the new trie = sum of old trie frequencies



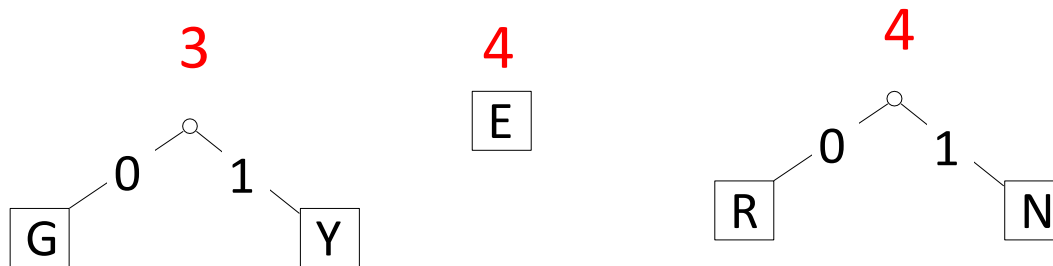
# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies  
 $G: 2, R: 2, E: 4, N: 2, Y: 1$
- Join two least frequent tries into a new trie
  - frequency of the new trie = sum of old trie frequencies



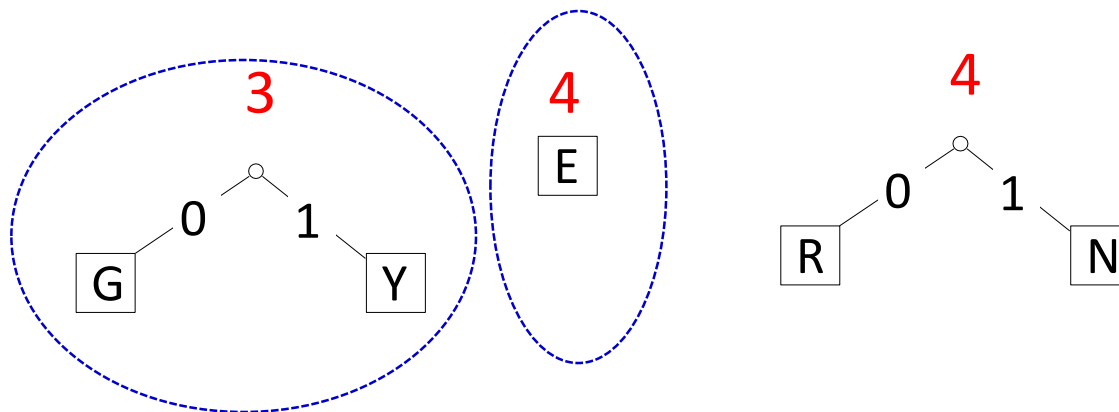
# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies  
 $G: 2, R: 2, E: 4, N: 2, Y: 1$
- Join two least frequent tries into a new trie
  - frequency of the new trie = sum of old trie frequencies



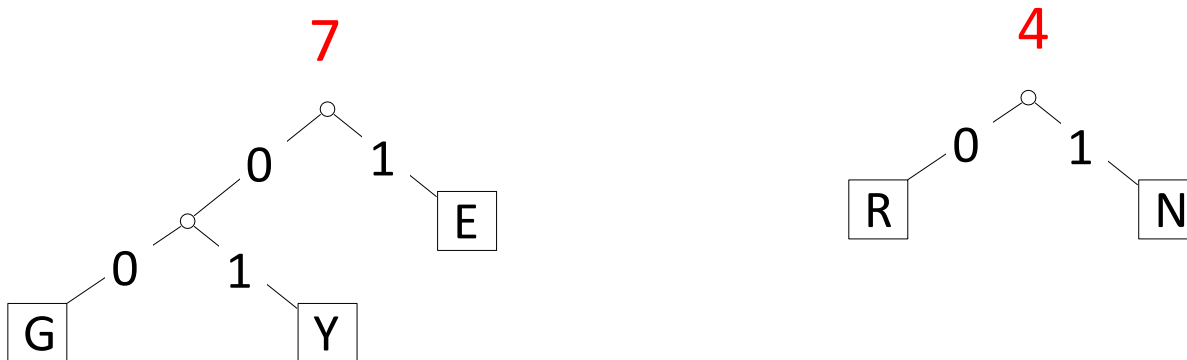
# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies  
 $G: 2, R: 2, E: 4, N: 2, Y: 1$
- Join two least frequent tries into a new trie
  - frequency of the new trie = sum of old trie frequencies



# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies  
 $G: 2, R: 2, E: 4, N: 2, Y: 1$
- Join two least frequent tries into a new trie
  - frequency of the new trie = sum of old trie frequencies



# Example: Huffman Tree Construction

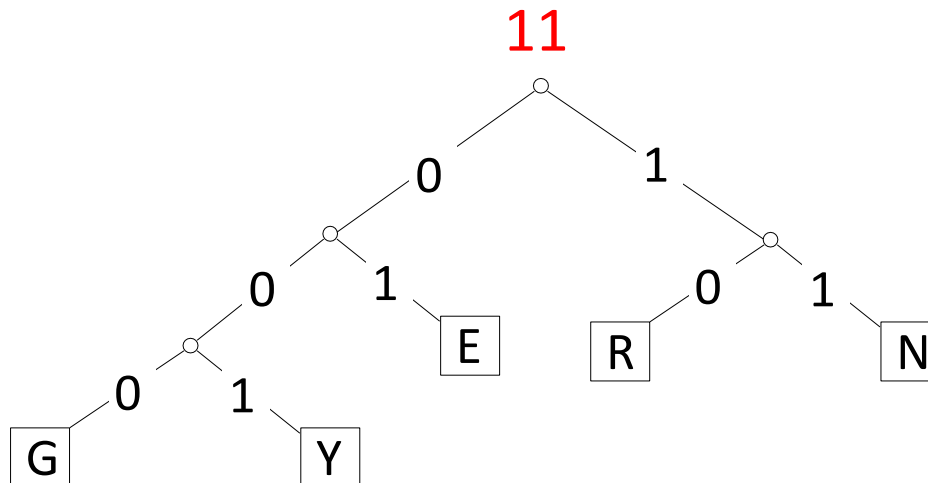
- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$

- Calculate character frequencies

$G: 2, R: 2, E: 4, N: 2, Y: 1$

- Join two least frequent tries into a new trie

- frequency of the new trie = sum of old trie frequencies



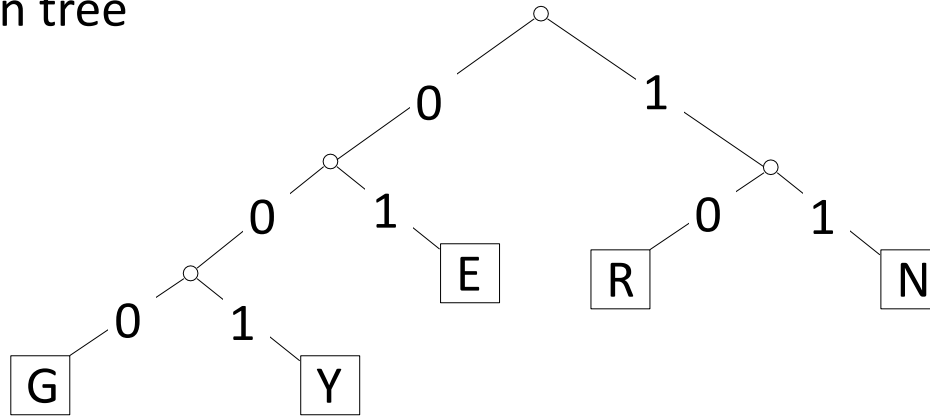


# Example: Huffman Tree Construction

- Example text: *GREENENERGY*,  $\Sigma_S = \{G, R, E, N, Y\}$
- Calculate character frequencies

$G: 2, R: 2, E: 4, N: 2, Y: 1$

- Final Huffman tree



- **GREENENERGY** → 000 10 01 01 11 01 11 01 10 000 001

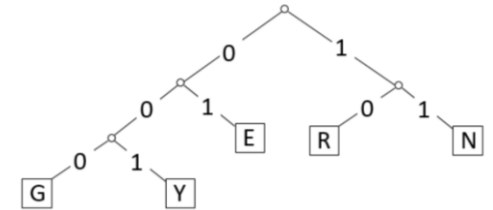
- Compression ratio

$$\frac{25}{11 \cdot \log 5} \approx 97\%$$

- Frequencies are not skewed enough to lead to good compression

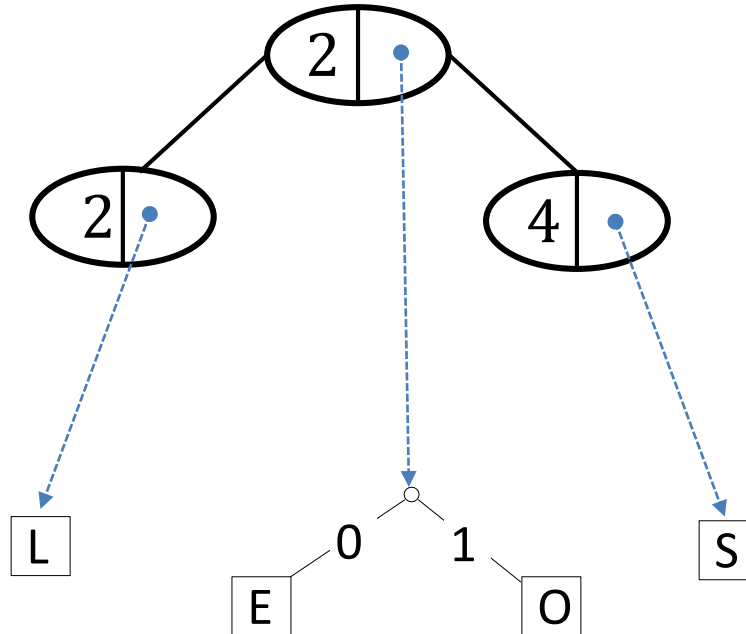
# Huffman Algorithm Summary

- **Greedy algorithm:** always pair up most frequent characters
  - 1) determine frequency of each character  $c \in \Sigma$  in  $S$
  - 2) for each  $c \in \Sigma$ , create trie of height 0 holding only  $c$ 
    - call it  $c$ -trie
  - 3) assign weight to each trie
    - weight trie character
  - 4) find and merge two tries with the minimum weight
    - new interior node added
    - the new weight is the sum of merged tries weights
    - corresponds to adding one bit to encoding of each character
  - 5) repeat Steps 4 until there is only 1 trie left
    - this is  $D$ , the final decoder
- Min-heap for efficient implementation: step 4 is two *delete-min* one *insert*



# Heap Storing Tries during Huffman Tree Construction

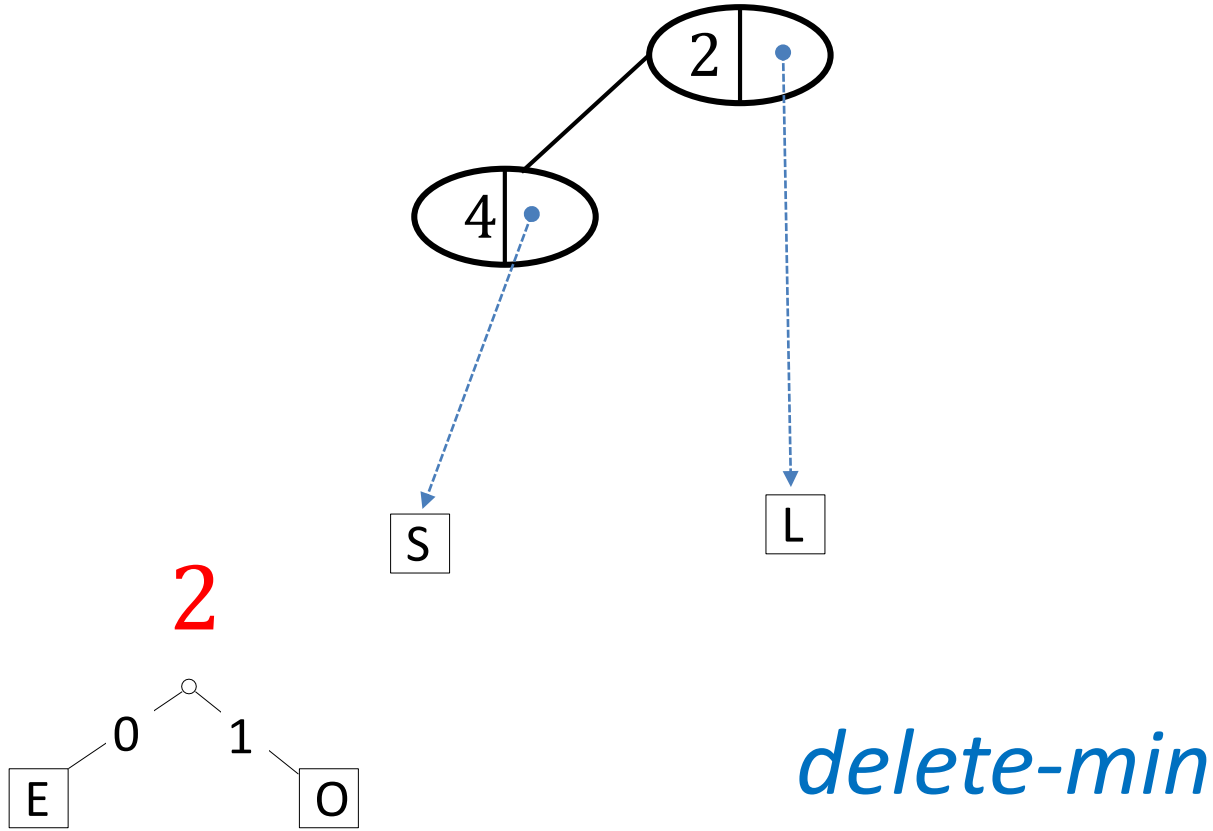
- (key,value) = (trie weight, link to trie)
- step 4 is two *delete-mins*, one *insert*



*delete-min*

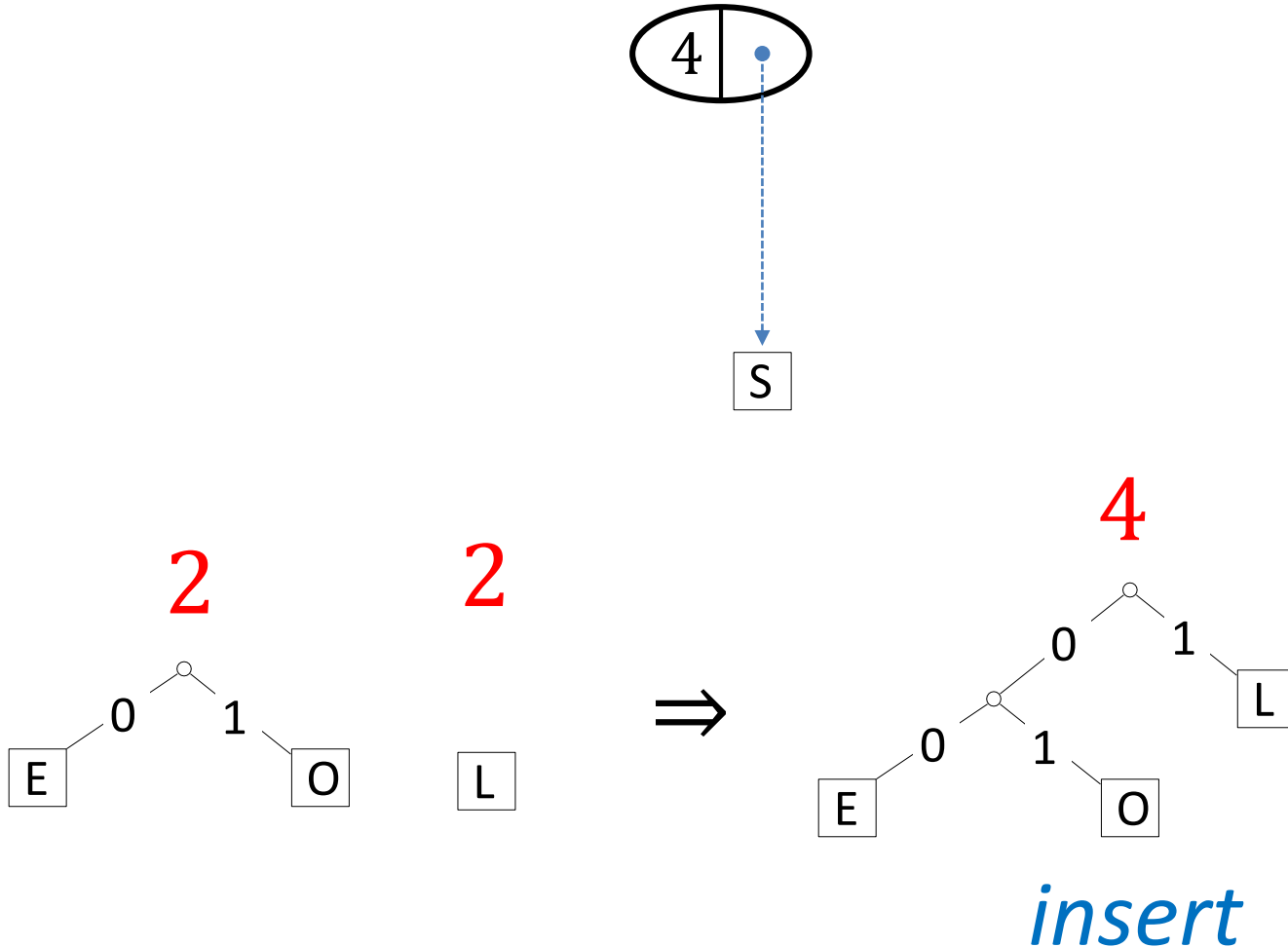
# Heap Storing Tries during Huffman Tree Construction

- (key,value) = (trie weight, link to trie)
- step 4 is two *delete-mins*, one *insert*



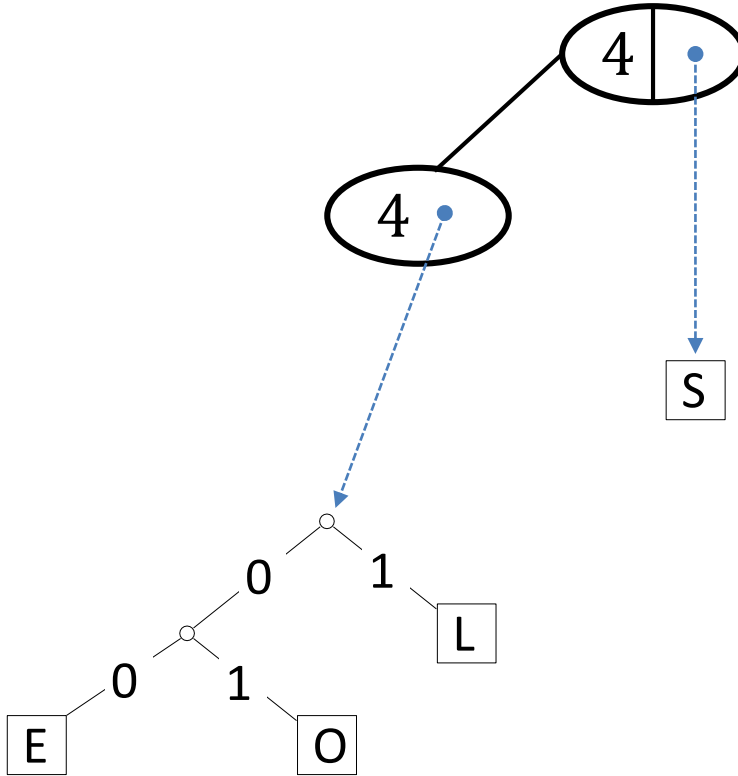
# Heap Storing Tries during Huffman Tree Construction

- (key,value) = (trie weight, link to trie)
- step 4 is two *delete-mins*, one *insert*



# Heap Storing Tries during Huffman Tree Construction

- (key,value) = (trie weight, link to trie)
- step 4 is two *delete-mins*, one *insert*



# Huffman's Algorithm Pseudocode

*Huffman::encoding*( $S, C$ )

$S$ : input-stream (length  $n$ ) with characters in  $\Sigma_S$ ,  $C$ : output-stream, initially empty

$f \leftarrow$  array indexed by  $\Sigma_S$ , initialized to 0

**while**  $S$  is non-empty **do increase**  $f[S.pop()]$  by 1 // get frequencies  $O(n)$

$Q \leftarrow$  min-oriented priority queue to store tries

**for all**  $c \in \Sigma_S$  with  $f[c] > 0$

$Q.insert(\text{single-node trie for } c, f[c])$   $O(|\Sigma_S| \log |\Sigma_S|)$

**while**  $Q.size() > 1$

$(T_1, f_1) \leftarrow Q.deleteMin()$

$(T_2, f_2) \leftarrow Q.deleteMin()$   $O(|\Sigma_S| \log |\Sigma_S|)$

$Q.insert(\text{trie with } T_1, T_2 \text{ as subtrees, } f_1 + f_2)$

$T \leftarrow Q.deleteMin()$  // trie for decoding

reset input-stream  $S$  // read all of  $S$ , need to read again for encoding

*prefixFree::encoding*( $T, S, C$ ) // perform actual encoding  $O(|\Sigma_S| + |C|)$

- Total time is  $O(|\Sigma_S| \log |\Sigma_S| + |C|)$ 
  - $n < |C|$

# Huffman Coding Discussion

- We require  $|\Sigma_S| \geq 2$
- The constructed trie is *optimal* in the sense that the coded text  $C$  is shortest among all prefix-free character encodings with  $\Sigma_C = \{0, 1\}$ 
  - proof is in the course notes
- Constructed trie is **not unique**
  - so decoding trie must be transmitted along with the coded text
  - this may make encoding bigger than source text!
- Many variations
  - tie-breaking rules, estimate frequencies, adaptively change encoding, etc.
- Encoding must pass through stream twice
  1. to compute frequencies and to encode
  2. cannot use stream unless it can be reset
- Encoding runtime:  $O(|\Sigma_S| \log |\Sigma_S| + |C|)$
- Decoding run-time:  $O(|C|)$



# Outline

- **Data Compression**
  - Background
  - Single-Character Encodings
  - Huffman Codes
  - **Run-Length Encoding**
  - Lempel-Ziv-Welch
  - Combining Compression Schemes: bzip2
  - Burrows-Wheeler Transform

# Single-Character vs Multi-Character Encoding

- **Single character encoding:** each source-text character receives one codeword

$S = \text{b a n a n a}$

01 1 11 1 11 1

- **Multi-character encoding:** multiple source-text characters can receive one codeword

$S = \text{b a n a n a}$

01 11 101

# Run-Length Encoding

- RLE is an example of multi-character encoding
- Source alphabet and coded alphabet are both binary:  $\Sigma = \{0, 1\}$ 
  - can be extended to non-binary alphabets
- Useful  $S$  has long runs of the same character: 00000 111 0000
- Dictionary is uniquely defined by *algorithm*
  - no need to store it explicitly

# Run-Length Encoding

- **Encoding idea**
  - give the first bit of  $S$  (either 0 or 1)
  - then give a sequence of integers indicating run lengths
  - do not have to give the bit for runs since they alternate
- Example 00000 111 0000
  - becomes: 0 5 3 4
- Need to encode run length in binary, how?
  - cannot use variable length binary encoding 10111100
    - do not know how to parse in individual run lengths
  - fixed length binary encoding (say 16 bits) wastes space, bad compression
    - 000000000000010100000000000011000000000000100

# Towards Prefix-free Encoding for Positive Integers

- To know where each number begins/ends, need to know number length
  - first say how many digits the number has
    - by printing as many zeros as the number of digits
  - then print the actual number

<i>number</i>		<i>encoding</i>
1	→	<del>0</del> 1
10	→	<del>00</del> 10
11	→	<del>00</del> 11
100	→	<del>000</del> 100
101	→	<del>000</del> 101

- The first zero is actually not necessary
  - # digits = #zeros + 1
  - get shorter encoding if remove the first zero

# Prefix-free Encoding for Positive Integers

- Use *Elias gamma code* to encode  $k$ 
  - $\lfloor \log k \rfloor$  copies of 0, followed by
  - binary representation of  $k$  (always starts with 1)

$k$	$\lfloor \log k \rfloor$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110

- Easy to decode
  - (number of zeros+1) tells you the length of  $k$  (in binary representation)
  - after zeros, read binary representation of  $k$  (it starts with 1)



# RLE Example: Encoding

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
7	2	111	<b>00111</b>

- Encoding

$S =$  **1111111**00100000000000000000000011111111111

$k = 7$

$C =$  1**00111**



# RLE Example: Encoding

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	<b>010</b>
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
7	2	111	00111

- Encoding

$S = \cancel{1111111} \mathbf{00} 100000000000000000000000011111111111$

$k = 2$

$C = 100111 \mathbf{010}$



# RLE Example: Encoding

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
7	2	111	00111
20	4	10100	<b>000010100</b>

- Encoding

$S = 1111111001$ **000000000000000000000000** $11111111111$

$k = 20$

$C = 1001110101$  **000010100**



# RLE Encoding

*RLE::encoding*( $S, C$ )

$S$ : input-stream of bits,  $C$ : output-stream

$b \leftarrow S.top()$

$C.append(b)$  //  $C$  is initialized to the first bit of  $S$

**while**  $S$  is non-empty **do**

$k \leftarrow 1$  // initialize run length

**while** ( $S$  is non-empty and  $S.top() = b$ ) //compute run length

$k++$ ;  $S.pop()$

    // compute Elias gamma code  $K$  (binary string) for  $k$

$w \leftarrow$  empty string

**while** ( $k > 1$ )

$C.append(0)$  // 0 appended to output  $C$  directly

$w.prepend(k \bmod 2)$  //  $K$  is built from last digit forwards

$k \leftarrow \lfloor k/2 \rfloor$

$w.prepend(1)$  // the very first digit of  $K$  was not computed

    append each bit of  $w$  to  $C$

$b \leftarrow 1 - b$

# RLE Example: Decoding

- Recall that  $(\# \text{ zeros} + 1)$  tells you the length of  $k$  in binary representation

- Decoding

$$C = 00001101001001010$$

$$b = 0$$

$$l = 4$$

$$k = 13$$

$$S = 00000000000000$$

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
<b>13</b>	3	<b>1101</b>	0001101

# RLE Example: Decoding

- Decoding

$C = 00001101001001010$

$b = 1$

$l = 3$

$k = 4$

$S = 000000000000001111$

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
<b>4</b>	2	<b>100</b>	00100
5	2	101	00101
6	2	110	00110
13	3	1101	0001101

# RLE Example: Decoding

- Decoding

$C = 00001101001001010$

$b = 0$

$l = 1$

$k = 1$

$S = 000000000000011110$

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
13	3	1101	0001101



# RLE Example: Decoding

- Decoding

$C = \text{00001101001001010}$

$b = 1$

$l = 2$

$k = 2$

$S = \text{00000000000001111011}$

$k$	$\lceil \log k \rceil$	$k$ in binary	encoding
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
13	3	1101	0001101

# RLE Decoding

*RLE::decoding*( $C, S$ )

$C$ : input stream of bits,  $S$ : output-stream

$b \leftarrow C.pop()$  // bit-value for the first run

**while**  $C$  is non-empty

$l \leftarrow 0$  // length of base-2 number - 1

**while**  $C.pop() = 0$

$l++$

$k \leftarrow 1$  // base-2 number converted

**for** ( $j = 1$  to  $l$ ) // translate  $k$  from binary string to integer

$k \leftarrow k * 2 + C.pop()$

    // if  $C$  runs out of bits then encoding was invalid

**for** ( $j = 1$  to  $k$ )

$S.append(b)$

$b \leftarrow 1 - b$  // alternate bit-value

# RLE Discussion

- Best compression (for most  $n$ ) is for  $S = 000 \dots 000$  of length  $n$ 
  - compressed to  $2\lceil \log n \rceil + 2 \in o(n)$  bits
    - 1 for the initial bit
    - $\lceil \log n \rceil$  zeros to encode the length of binary representation of integer  $n$
    - $\lceil \log n \rceil + 1$  digits that represent  $n$  itself in binary
- Usually not that lucky
  - no compression until run-length  $k \geq 6$
  - **expansion** when run-length  $k = 2$  or  $4$
- RLE was popular for fax-machines
- Can be adapted to larger alphabet sizes
  - but then the encoding for each run must also store the character
- Can be adapted to encode only runs of 0
  - used inside bzip2 (will see later)

# Outline

- **Compression**

- Encoding Basics
- Huffman Codes
- Run-Length Encoding
- **Lempel-Ziv-Welch**
- bzip2
- Burrows-Wheeler Transform

# Longer Patterns in Input

- Huffman and RLE take advantage of frequent or repeated *single characters*
- **Observation:** certain *substrings* are much more frequent than others
- Examples
  - English text
    - most frequent digraphs: TH, ER, ON, AN, RE, HE, IN, ED, ND, HA
    - most frequent trigraphs: THE, AND, THA, ENT, ION, TIO, FOR, NDE
  - HTML
    - “<a href”, “<img src”, “<br>”
  - Video
    - repeated background between frames, shifted sub-image
- Could find the most frequent substrings of length up to  $k$  and store them in a dictionary (in addition to characters, i.e. strings of length 1)

	null	start of heading	start of text	...	A	...	delete	er	in	...	ed	the
code	0	1	2	...	65	...	127	128	129	...		255
code in binary	00000000	00000001	00000010		001000001		01111111	11000001	11000010	...	11111110	11111111

- however, each text has its own set of most frequently occurring substrings

# Lempel-Ziv-Welch Compression

- **Ingredient 1** for Lempel-Ziv-Welch compression
  - encode characters and frequent substrings
    - discover and encode frequent substring as we process text
      - no need to know frequent substrings beforehand

# Adaptive Dictionaries

- ASCII and RLE use *fixed* dictionaries
  - same dictionary for any text encoded
  - no need to pass dictionary to the decoder
- Huffman's dictionary is not *fixed* but it is *static*
  - dictionary is not fixed: each text has its own dictionary
  - dictionary is *static*: dictionary *does not change* for entire encoding/decoding
  - need to pass dictionary to the decoder
- **Ingredient 2** for LZW: *adaptive dictionary*
  - dictionary constructed during encoding/decoding
    - no need to send dictionary with the encoding,
  - start with some initial fixed dictionary  $D_0$ 
    - usually ASCII
  - at iteration  $i \geq 0$ ,  $D_i$  is used to determine the  $i$ th output
  - after  $i$ th output (iteration  $i$ ), update  $D_i$  to  $D_{i+1}$ 
    - $D_{i+1} \leftarrow D_i.\text{insert}(\text{new character combination})$
  - decoder must know (i.e. be able to reconstruct from the coded text) how encoder changed the dictionary

# LZW Encoding: Main Idea

- Iteration  $i$  of encoding
- Current  $D_i = \{a:65, b: 66, ab:140, \mathbf{bb:145}, bbc:146\}$

S = abbbc**bb**ad

C = 78 95 **145**



- find longest substring that starts at current pointer and is in the dictionary
- encode 'bb' with 145
- $D_{i+1} = D_i.\mathit{insert}(\mathbf{bba}, \mathit{nextAvailableCodenum})$
- 'bba' would have been useful at iteration  $i$ , so likely useful in the future

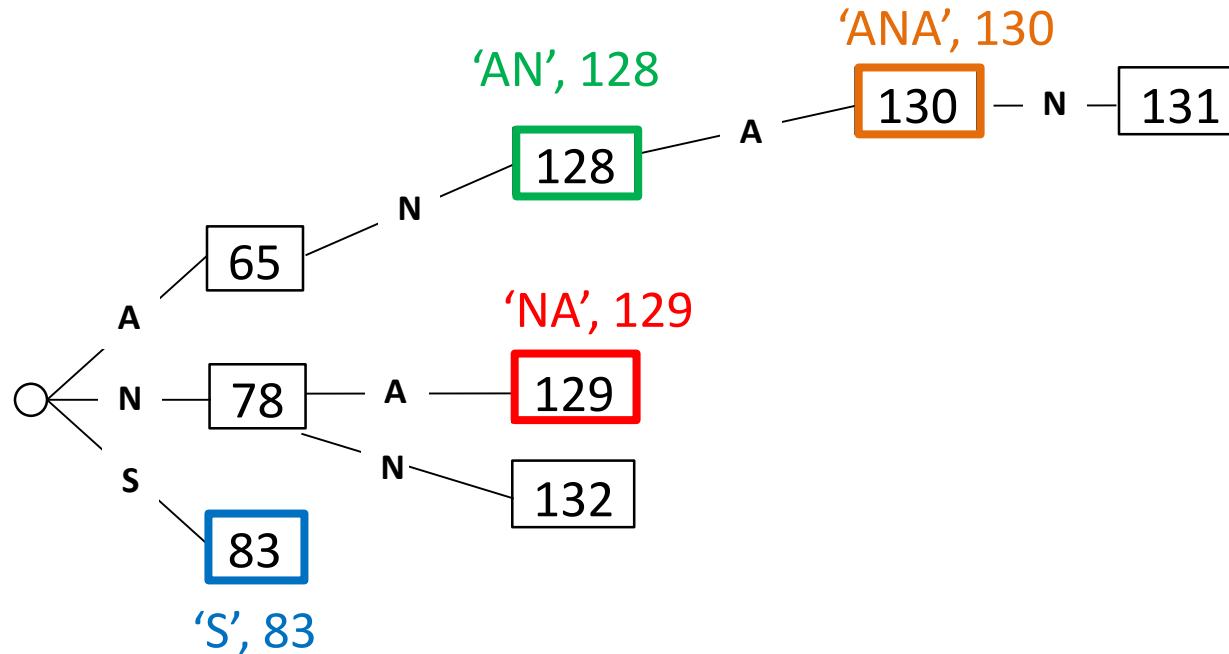
- After iteration  $i$

$D_{i+1} = \{a:65, b: 66, ab:140, bb:145, bbc:146, bba:147\}$

- $\mathit{codenumber} = \mathit{codeword} = \mathit{code}$



# Tries for LZW Encoding

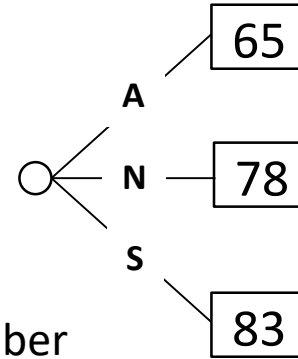


- Store (string, codeword) pairs, with string being the key
- Variation of tries different from what we have seen before
- Trie stores codenumbers at **all** nodes (external and internal) except the root
  - works because a string is inserted only after all its prefixes are inserted
- Do not store the string key explicitly, store only the codenumber
  - read the string key corresponding to each codenumber from the edges

# LZW Example

- Start dictionary  $D$

- ASCII characters
- codes from 0 to 127
- next inserted code will be 128
- variable  $idx$  keeps track of next available codenumber
- initialize  $idx = 128$

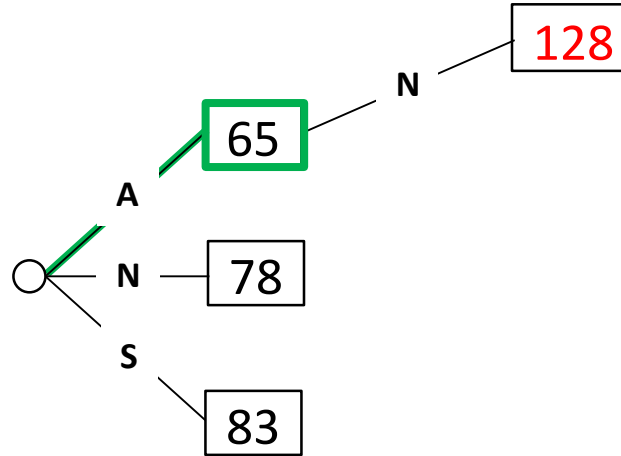


- Text

A N A N A N A N N A

# LZW Example

- Dictionary  $D$ 
  - $idx = 129$

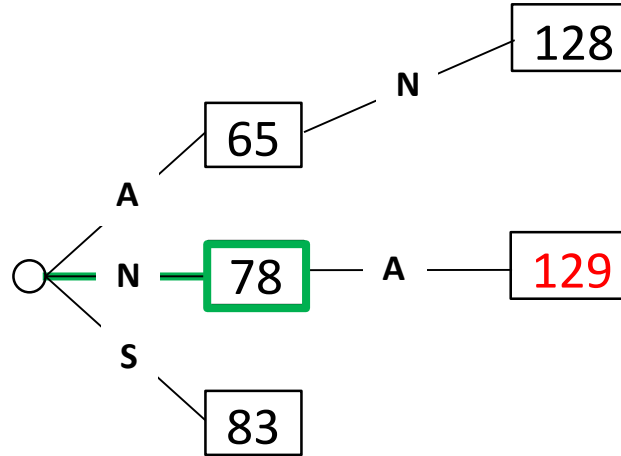


- Text      A    N    A    N    A    N    N    A
- Encoding    65

- Add to dictionary “string just encoded” + “first character of next string to be encoded”
- Inserting new item into  $D$  is  $O(1)$  since we stopped at the right node in the trie when we searched for ‘A’

# LZW Example

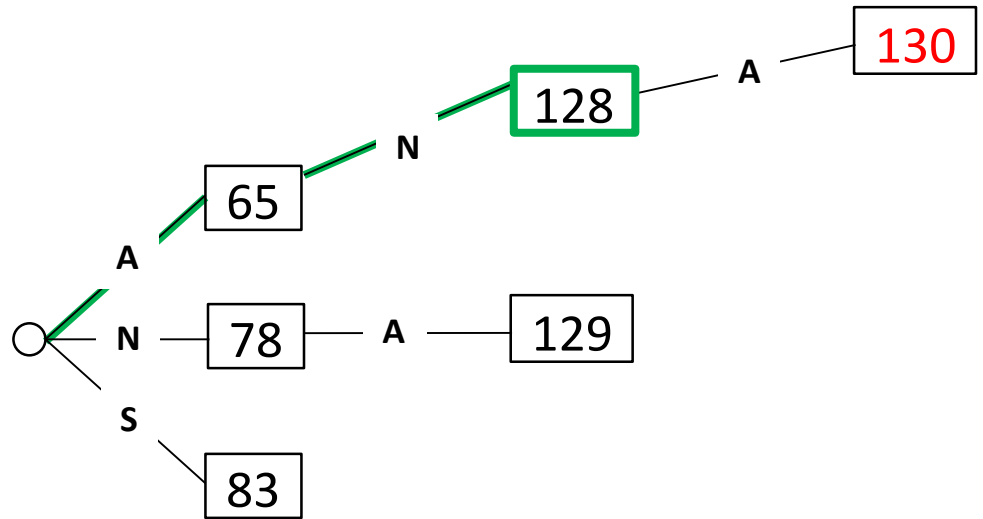
- Dictionary  $D$ 
  - $idx = 130$



- Text      A    **N** A    N    A    N    A    N    N    A
- Encoding    65    **78**

# LZW Example

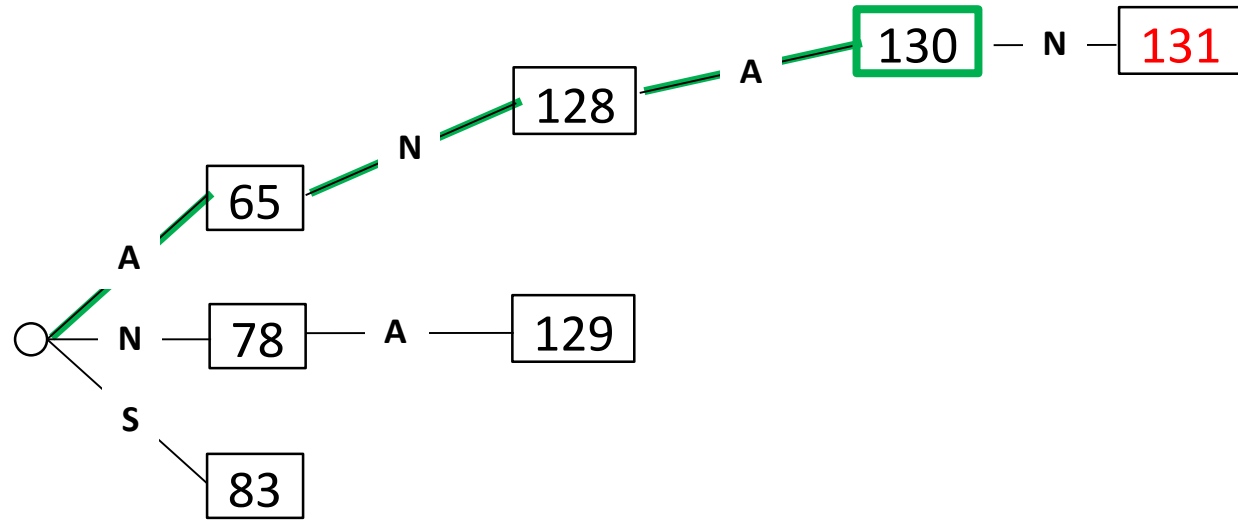
- Dictionary  $D$ 
  - $idx = 131$



Text	A	N	A N A	N	A	N	N	A
Encoding	65	78	128					

*Note: The sequence 'A N A' in the Text row and the value '128' in the Encoding row are highlighted in green and enclosed in a red dashed box, with the text 'add to dictionary' written above the box.*

# LZW Example



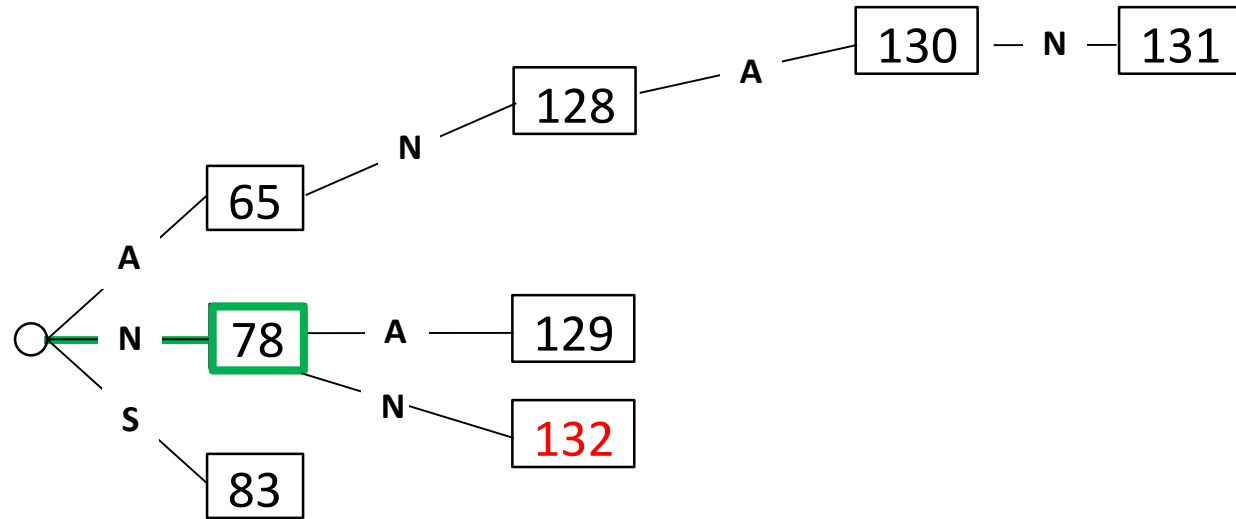
- Dictionary D
  - $idx = 132$

Text	A	N	A	N	A	N	A	N	N	A
Encoding	65	78	128		130					

add to dictionary

# LZW Example

- Dictionary D
  - $idx = 133$



Text

A | N | A N | A N | A | N N A

Encoding

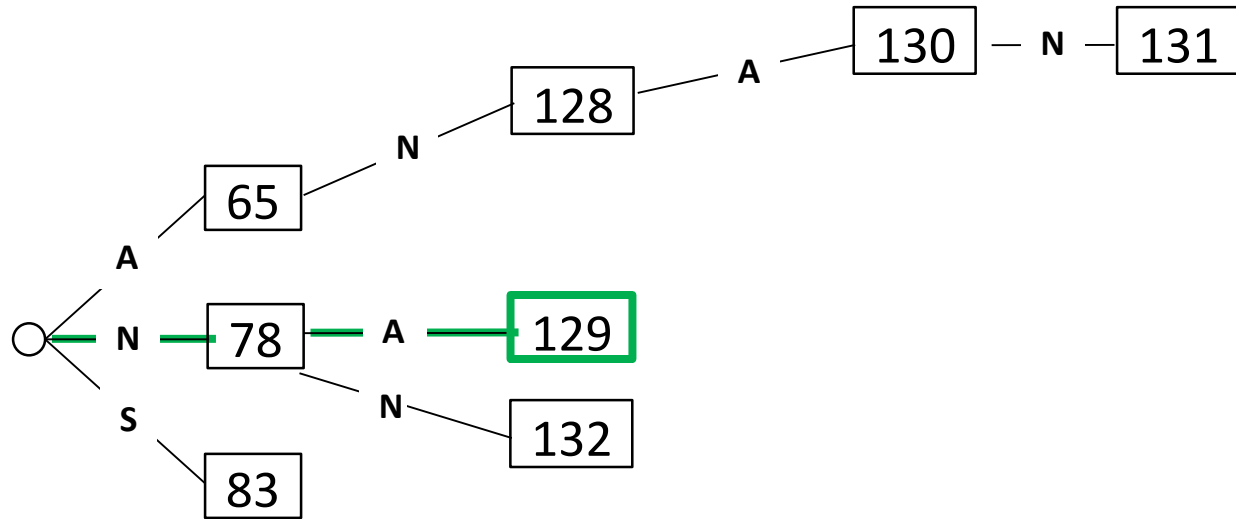
65 | 78 | 128 | 130

add to dictionary



78

# LZW Example

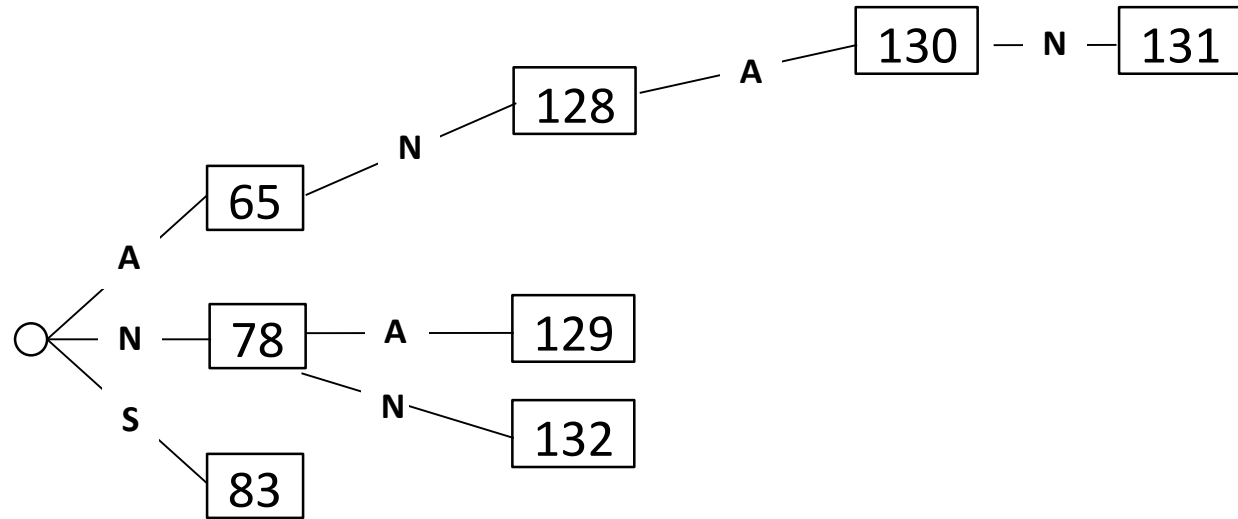


- Dictionary D
  - $idx = 133$

Text	A	N	A	N	A	N	A	N	N	A
Encoding	65	78	128		130		78	129		



# LZW Example



- Dictionary D
  - $idx = 133$

Text	A	N	A	N	A	N	A	N	N	A
Encoding	65	78	128		130		78	129		

Final output 00001000001 00001001110 00010000000 0001000010 00001001110 00010000001

- Use fixed length (12 bits) per codenumber
  - 12 bit binary string representation for each code
  - total of  $2^{12} = 4096$  codesnumbers available during encoding
    - if you run out of codenumbers, stop inserting new elements in the dictionary

# LZW encoding pseudocode

*LZW::encoding*( $S, C$ )

$S$  : input stream of characters,  $C$ : output-stream

initialize dictionary  $D$  with ASCII in a trie

$idx \leftarrow 128$

**while**  $S$  is not empty **do**

$v \leftarrow$  root of trie  $D$

**while**  $S$  is non-empty and  $v$  has a child  $c$  labelled  $S.top()$

$v \leftarrow c$

$S.pop()$

$C.append$ (codenumber stored at  $v$ )

**if**  $S$  is non-empty

create child of  $v$  labelled  $S.top()$  with code  $idx$

$idx ++$

trie  
search

new  
dictionary  
entry

- Running time is  $O(|S|)$ 
  - assuming can look up child labeled with  $c$  in  $O(1)$  time
    - i.e. trie node stores children in an array

# LZW Encoder vs Decoder

- For decoding, need a dictionary
- Construct a dictionary during decoding, imitating what encoder does
- But will be forced to be one step behind
  - at iteration  $i$  of decoding can reconstruct substring which encoder inserted into dictionary at iteration  $i - 1$ 
    - delay is due to not having access to the original text

# LZW Decoding Example

- Given encoding to decode back to the source text

65            78            128            130            78            129

- Build dictionary adaptively, while decoding
- Decoding starts with the same initial dictionary as encoding
  - use array instead of trie, need  $D$  that allows efficient search by code
- We will show the original text during decoding in this example, but just for reference
  - do not need original text to decode

initial  $D$

65	A
78	N
83	S

$idx = 128$

# LZW Decoding Example

	$i=0$									
Text	A	N	A	N	A	N	A	N	N	A
Encoding	65	78	128		130		78	129		
Decoding	A									

iter  $i = 0$

$D =$

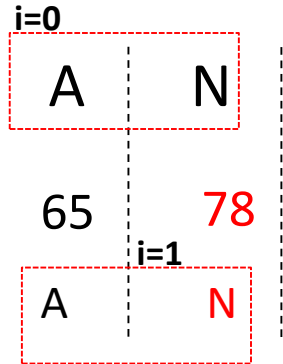
65	A
78	N
83	S

$idx = 128$

- First step:  $s = D(65) = A$
- Encoding iteration  $i = 0$ 
  - looked ahead in the text, saw N, and added AN to the dictionary
- Decoding iteration  $i = 0$ 
  - know text starts with A, but cannot look ahead as the text is not available
  - no new word added at iteration  $i = 0$
  - keep track of  $s_{prev}$  = string decoded at previous iteration
    - $s_{prev}$  is also string encoder encoded at previous iteration

# LZW Decoding Example

- Text
- Encoding
- Decoding  
iter  $i = 1$



A N A N A N N A  
128 130 78 129

$D =$

65	A
78	N
83	S
128	AN

$idx = 129$

- $s_{prev} = A$ 
  - string encoded/decoded at previous iteration
- First step:  $s = D(78) = N$
- The first letter of  $s$  is exactly what the encoder looked ahead at during previous iteration!
- So at previous iteration, encoder added to the dictionary  $s_{prev} + s[0]$ 
  - A N
- Starting at iteration  $i = 1$  of decoding
  - add  $s_{prev} + s[0]$  to dictionary

# LZW Decoding Example Continued

Text	A	N	A	N	A	N	A	N	N	A
Encoding	65	78	128		130		78		129	
Decoding	A	N	AN							

iter  $i = 2$

$D =$

65	A
78	N
83	S
128	AN
129	NA

$idx = 130$

- $s_{prev} = N$ 
  - string encoded/decoded at previous iteration
- First step:  $s = D(128) = AN$
- Next step: add to dictionary  $s_{prev} + s[0]$ 

$$N + A = NA$$
  - encoder added this string at previous iteration

# LZW Decoding Example

▪ Text	A	N	A	N	A	N	A	N	N	A
▪ Encoding	65	78	128		130		78	129		
▪ Decoding	A	N	AN		$s = ???$					

iter  $i = 3$

$D =$

65	A
78	N
83	S
128	AN
129	NA

$idx = 130$

- $s_{prev} = AN$ 
  - string encoded/decoded at previous iteration
- First step:  $s = D(130) = ???$  (code 130 is not in  $D$ )
- Dictionary is exactly one step behind at decoding
- Encoder added  $(s, 130)$  to  $D$  at previous iteration
- What did the encoder add at the previous iteration?

$$\begin{array}{l}
 \text{known} \quad \text{unknown} \quad \text{unknown} \\
 s_{prev} + s[0] = s \\
 \text{AN} + s[0] = s \\
 \text{AN} \xrightarrow{\text{A}} \text{ANA} = s \\
 s[0] = s_{prev}[0] = \text{A} \\
 \text{ANA} = s
 \end{array}$$



# LZW Decoding Example

- Text            A    N    A    N    A    N    N    A
  - Encoding       65   78   128            130            78            129
  - Decoding       A    N    AN            ANA
- iter  $i = 3$

$D =$

65	A
78	N
83	S
128	AN
129	NA
130	ANA

$idx = 131$

- General rule: if code  $C$  is not in  $D$ 
  - $s = s_{prev} + s_{prev}[0]$
- in our example,  $s_{prev} = AN$ 
  - $s = AN + A = ANA$
- Now that we recovered  $s$ , continue as usual
- Add to dictionary  $s_{prev} + s[0]$

# LZW Decoding Example

▪ Text	A	N	A	N	A	N	A	N	N	A
▪ Encoding	65	78	128		130		78	129		
▪ Decoding	A	N	AN		ANA		N			

iter  $i = 4$

$D =$

65	A
78	N
83	S
128	AN
129	NA
130	ANA
131	ANAN

$idx = 132$

- $s_{prev} = ANA$
- If code  $C$  is not in  $D$ 

$$s = s_{prev} + s_{prev}[0]$$
- Add to dictionary  $s_{prev} + s[0]$

# LZW Decoding Example

Text	A	N	A	N	A	N	A	N	N	A
Encoding	65	78	128		130		78	129		
Decoding	A	N	AN		ANA		N	NA		

iter  $i = 5$

$D =$

65	A
78	N
83	S
128	AN
129	NA
130	ANA
131	ANAN

$idx = 132$

- $s_{prev} = N$
- If code  $C$  is not in  $D$ 

$$s = s_{prev} + s_{prev}[0]$$
- Add to dictionary  $s_{prev} + s[0]$

# LZW decoding

- To save space, store new codes using its prefix code + one character
  - given a codenumber, can find corresponding string  $s$  in  $O(|s|)$  time

$D =$

65	A
78	N
83	S
128	AN
129	NA
130	ANA
131	ANAN

wasteful storage

65	A
78	N
83	S
128	65, N
129	78, A
130	128, A
131	130, N

ANAN

65, NAN

128, AN

The diagram illustrates the compact LZW dictionary. It shows a table with 7 rows. The first three rows are the same as in the 'wasteful storage' table. The last four rows contain the new codes and their corresponding strings: 128 maps to '65, N', 129 to '78, A', 130 to '128, A', and 131 to '130, N'. To the right of the table, the string 'ANAN' is shown. Three curved arrows point from the string to the dictionary entries: one from 'ANAN' to the entry for code 128, one from 'ANAN' to the entry for code 129, and one from 'ANAN' to the entry for code 130. This demonstrates how the compact dictionary stores the strings used in the previous step, allowing for efficient decoding.

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

$D =$

next  
available  
code

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128		

# LZW decoding, Another Example

Encoding:            98  97 114 128 114 97 131 134 129 101 135

Decoding:            B

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128		

$s = B$

nothing added to dictionary at iteration 0

# LZW decoding, Another Example

Encoding:            98  97  114  128  114  97  131  134  129  101  135

Decoding:            B    A

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A

$s_{prev} = B, code_{prev} = 98$

$s = A$

add to dictionary  $s_{prev} + s[0] = BA$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R

$s_{prev} = A, code_{prev} = 97$   
 $s = R$

add to dictionary  $s_{prev} + s[0] = AR$



# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B

$s_{prev} = R, code_{prev} = 114$   
 $s = BA$

add to dictionary  $s_{prev} + s[0] = RB$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R

$s_{prev} = BA, code_{prev} = 128$

$s = R$

add to dictionary  $s_{prev} + s[0] = BAR$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R    A

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R
132	RA	114, A

$s_{prev} = R, code_{prev} = 114$

$s = A$

add to dictionary  $s_{prev} + s[0] = RA$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R    A   BAR

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R
132	RA	114, A
133	AB	97, B

$s_{prev} = A, code_{prev} = 97$

$s = BAR$

add to dictionary  $s_{prev} + s[0] = AB$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R    A   BAR   BARB

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R
132	RA	114, A
133	AB	97, B
134	<b>BARB</b>	<b>131, B</b>

$$s_{prev} = \text{BAR}, code_{prev} = 131$$

$$s = ?$$

if code is not in dictionary

$$s = s_{prev} + s_{prev}[0]$$

$$s = \text{BAR} + \text{B} = \text{BARB}$$

add to dictionary  $s_{prev} + s[0] = \text{BARB}$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R    A   BAR   BARB   AR

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R
132	RA	114, A
133	AB	97, B
134	BARB	131, B
135	BARBA	134, A

$$s_{prev} = \text{BARB}, code_{prev} = 134$$

$$s = \text{AR}$$

add to dictionary  $s_{prev} + s[0] = \text{BARBA}$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R    A   BAR   BARB   AR   E

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R
132	RA	114, A
133	AB	97, B
134	BARB	131, B
135	BARBA	134, A
136	ARE	129, E

$$s_{prev} = AR, code_{prev} = 129$$

$$s = E$$

add to dictionary  $s_{prev} + s[0] = ARE$

# LZW decoding, Another Example

Encoding:            98   97   114   128   114   97   131   134   129   101   135

Decoding:            B    A    R    BA    R    A   BAR   BARB   AR   E   BARBA

$D =$

code	string (human)	string (implementation)
97	A	
98	B	
101	E	
110	N	
114	R	
128	BA	98, A
129	AR	97, R
130	RB	114, B
131	BAR	128, R
132	RA	114, A
133	AB	97, B
134	BARB	131, B
135	BARBA	134, A
136	ARE	129, E

$s_{prev} = E$   
 $s = BARBA$



# LZW Decoding Pseudocode

*LZW::decoding*( $C, S$ )

$C$  : input-stream of integers,  $S$ : output-stream

$D \leftarrow$  dictionary that maps  $\{0, \dots, 127\}$  to ASCII

$idx \leftarrow 128$  // next available code

$code \leftarrow C.pop()$ ;  $s \leftarrow D.search(code)$ ;  $S.append(s)$

**while** there are more codes in  $C$  **do**

$s_{prev} \leftarrow s$ ;  $code \leftarrow C.pop()$

**if**  $code < idx$  **then**

$s \leftarrow D.search(code)$  //code in  $D$ , look up string  $s$

**if**  $code = idx$  // code not in  $D$  yet, reconstruct string

$s \leftarrow s_{prev} + s_{prev}[0]$

**else** **Fail** // invalid encoding

append each character of  $s$  to  $S$

$D.insert(idx, s_{prev} + s[0])$

$idx ++$

- Running time is  $O(|S|)$

# LZW Discussion

- Encoding is  $O(|S|)$  time, uses a trie of encoded substrings to store the dictionary
- Decoding is  $O(|S|)$  time, uses an array indexed by code numbers to store the dictionary
- Encoding and decoding need to go through the string only one time and do not need to see the whole string
  - can do compression while streaming the text
- Works badly if no repeated substrings
  - dictionary gets bigger, but no new useful substrings inserted
- In practice, compression rate is around 45% on English text

# Lempel-Ziv Family

- Lempel-Ziv is a family of *adaptive* compression algorithms
  - **LZ77** Original version (“sliding window”)
    - Derivatives: LZSS, LZFG, LZRW, LZIP, DEFLATE, . . .
      - DEFLATE used in (pk)zip, gzip, PNG
  - **LZ78** Second (slightly improved) version
    - Derivatives LZW, LZMW, LZAP, LZJ, . . .
    - LZW used in compress, GIF
      - patent issues

# Outline

- **Data Compression**
  - Background
  - Single-Character Encodings
  - Huffman Codes
  - Run-Length Encoding
  - Lempel-Ziv-Welch
  - **Combining Compression Schemes: bzip2**
  - Burrows-Wheeler Transform

# Overview of bzip2

- **Idea**: Combine multiple compression schemes and *text transforms*
  - *text transform*: change input text into a *different text*
    - output is not shorter, but likely to compresses better

$T_0 = \text{alfeatsalfalfa}$

Burrows-Wheeler transform

if  $T_0$  has repeated longer substrings, then  $T_1$  has long runs of characters

$T_1 = \text{affs\$eflllaaata}$

Move-to-front transform

if  $T_1$  has long runs of characters, then  $T_2$  has long runs of 0 and skewed frequencies

text  $T_2 = 13053435006006$

Modified RLE

if  $T_2$  has long runs of zeroes, then  $T_3$  is shorter; skewed frequencies remain

text  $T_3$

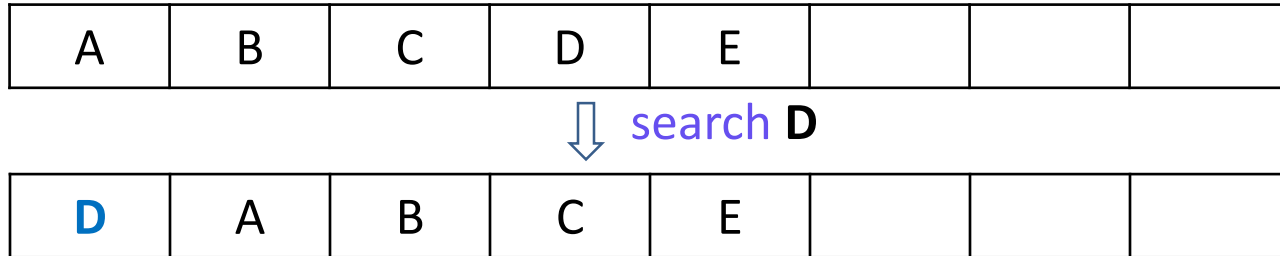
Huffman encoding

compresses well since frequencies are skewed

text  $T_4$

# Move-to-Front transform

- Recall the MTF heuristic
  - after an element is accessed, move it to array front



- Use this idea for **MTF** (move to front) text transformation
  - transformed text is likely to have text with repeated zeros and skewed frequencies

# MTF Encoding Example

- Source alphabet  $\Sigma_S$  with size  $|\Sigma_S| = m$
- Put alphabet in array  $L$ , initially in sorted order, but allow  $L$  to get unsorted

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

- This gives encoding dictionary  $L$ 
  - single character encoding  $E$
- Code of any character = index of array where character stored in dictionary  $L$ 
  - $E('B') = 1$
  - $E('H') = 7$
- After each encoding, update  $L$  with Move-To-Front heuristic
- Coded alphabet is  $\Sigma_C = \{0, 1, \dots, m - 1\}$
- Change dictionary  $D$  dynamically (like LZW)
  - unlike LZW
    - no new items added to dictionary
    - codeword for one or more letters can change at each iteration

# MTF Encoding 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

$S = \text{MISSISSIPPI}$

$C =$



# MTF Encoding 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

$S = \text{MISSISSIPPI}$

$C = 12$

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

S = MISSISSIPPI

C = 12 9

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

S = ~~M~~ISSISSIPPI

C = 12 9 **18**

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

S = MISSISSIPPI

C = 12 9 18 0

# MTF Encoding Example

0	<b>1</b>	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	<b>I</b>	M	A	B	C	D	E	F	G	H	J	K	L	N	O	P	Q	R	T	U	V	W	X	Y	Z

S = MISSISSIPPI

C = 12 9 18 0 **1**

# MTF Encoding Example

0	<b>1</b>	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	<b>S</b>	M	A	B	C	D	E	F	G	H	J	K	L	N	O	P	Q	R	T	U	V	W	X	Y	Z

S = ~~MISSI~~**S**SIPPI

C = 12 9 18 0 1 **1**

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

S = MISSISSIPPI

C = 12 9 18 0 1 1 0

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

$S = \text{MISSISSIPPI}$

$C = 12\ 9\ 18\ 0\ 1\ 1\ 0\ 1\ 16\ 0\ 1$

- What does a run in  $C$  mean about the source  $S$ ?
  - zeros tell us about consecutive character runs



# MTF Decoding 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

$S =$

$C = 12\ 9\ 18\ 0\ 1\ 1\ 0\ 1\ 16\ 0\ 1$

- Decoding is similar
- Start with the same dictionary  $D$  as encoding
- Apply the same MTF transformation at each iteration
  - dictionary  $D$  undergoes exactly the transformations when decoding
  - no delays, identical dictionary at encoding and decoding iteration  $i$
  - can always decode original letter

# MTF Decoding 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

$S = M$

$C = 12\ 9\ 18\ 0\ 1\ 1\ 0\ 1\ 16\ 0\ 1$

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

$S = M \quad |$

$C = 12 \quad 9 \quad 18 \quad 0 \quad 1 \quad 1 \quad 0 \quad 1 \quad 16 \quad 0 \quad 1$

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

$S = M \ I \ S$

$C = 12 \ 9 \ 18 \ 0 \ 1 \ 1 \ 0 \ 1 \ 16 \ 0 \ 1$

# Move-to-Front Transform: Properties



- If a character in  $S$  repeats  $k$  times, then  $C$  has a run of  $k - 1$  zeros
- $C$  contains a lot of small numbers and a few big ones
- $C$  has the same length as  $S$ , but better properties for encoding

# Move-to-Front Encoding/Decoding Pseudocode

*MTF::encoding*( $S, C$ )

$L \leftarrow$  array with  $\Sigma_S$  in some pre-agreed, fixed order (i.e. ASCII)

**while**  $S$  is non-empty **do**

$c \leftarrow S.\text{pop}()$

$i \leftarrow$  index such that  $L[i] = c$

**for**  $j = i - 1$  down to 0

    swap  $L[j]$  and  $L[j + 1]$

*MTF::decoding*( $C, S$ )

$L \leftarrow$  array with  $\Sigma_S$  in some pre-agreed, fixed order (i.e. ASCII)

**while**  $C$  is non-empty **do**

$i \leftarrow$  next integer of  $C$

$S.\text{append}(L[i])$

**for**  $j = i - 1$  down to 0

    swap  $L[j]$  and  $L[j + 1]$

# Move-to-Front Transform Summary

- **MTF text transform**
  - source alphabet is  $\Sigma_S$  with size  $|\Sigma_S| = m$
  - store alphabet in an array
    - code of any character = index of array where character stored
    - coded alphabet is  $\Sigma_C = \{0, 1, \dots, m - 1\}$
  - Dictionary is adaptive
    - nothing new is added, but meaning of codewords are changed
  - MTF is an *adaptive* text-transform algorithm
    - it does not compress input
    - the output has the same length as input
    - but output has better properties for compression

# Outline

- **Data Compression**
  - Background
  - Single-Character Encodings
  - Huffman Codes
  - Run-Length Encoding
  - Lempel-Ziv-Welch
  - Combining Compression Schemes: bzip2
  - **Burrows-Wheeler Transform**



# Burrows-Wheeler Transform

- Transformation (not compression) algorithm
  - transforms source text to coded text with same letters but in different order
    - source and coded alphabets are the same
  - if original text had frequently occurring substrings, transformed text should have many runs of the same character
    - more suitable for MTF transformation*



- Required: the source text  $S$  ends with *end-of-word character*  $\$$ 
  - $\$$  occurs nowhere else in  $S$
  - count  $\$$  towards length of  $S$
- Based on *cyclic* shifts for a string
  - example
    - string  
 $\text{abcde}$
    - cyclic shift  
 $\text{cdeab}$
- Formal definition
  - a *cyclic shift* of string  $X$  of length  $n$  is the concatenation of  $X[i + 1 \dots n - 1]$  and  $X[0 \dots i]$ , for  $0 \leq i < n$



# BWT Algorithm and Example

$S = \text{alfeatsalfalfa\$}$

- Write all consecutive cyclic shifts
  - forms *an array of shifts*
  - last letter in any row is the first letter of the previous row

```
alfeatsalfalfa$
lfeatsalfalfa$a
featsalfalfa$al
eatsalfalfa$aalf
atsalfalfa$aalfe
tsalfalfa$aalfea
salfalfa$aalfeats
alfalfa$aalfeatsa
falffa$aalfeatsal
alfa$aalfeatsalf
lfa$aalfeatsalfal
fa$aalfeatsalfal
a$aalfeatsalfalf
$aalfeatsalfalfa
```

# BWT Algorithm and Example

$S = \text{alfeatsalfalfa\$}$

- Array of cyclic shifts
  - the first column is the original  $S$

```
a l f e a t s a l f a l f a $
l f e a t s a l f a l f a $ a
f e a t s a l f a l f a $ a l
e a t s a l f a l f a $ a l f
a t s a l f a l f a $ a l f e
t s a l f a l f a $ a l f e a
s a l f a l f a $ a l f e a t
a l f a l f a $ a l f e a t s
l f a l f a $ a l f e a t s a
f a l f a $ a l f e a t s a l
a l f a $ a l f e a t s a l f
l f a $ a l f e a t s a l f a
f a $ a l f e a t s a l f a l
a $ a l f e a t s a l f a l f
$ a l f e a t s a l f a l f a
```

# BWT Algorithm and Example

$S = \mathbf{a} | \mathbf{lf} eats | \mathbf{a} | \mathbf{lf} a | \mathbf{lf} a \$$

- Array of cyclic shifts
- $S$  has **alf** repeated 3 times
  - 3 different shifts start with **lf** and end with **a**

```
alf eatsalfalfa$  
lfeatsalfalfa$aa  
featsalfalfa$alf  
eatsalfalfa$alf  
atsalfalfa$alf  
tsalfalfa$alf  
salfalfa$alf  
alfalfa$alf  
lfalfa$alf  
alfalfa$alf  
alfalfa$alf  
lfa$alf  
fa$alf  
a$alf  
$alf
```

# BWT Algorithm and Example

$S = \mathbf{a} | \mathbf{l f} e a t s \mathbf{a} | \mathbf{l f a} | \mathbf{l f} a \$$

- Array of cyclic shifts
- Sort (lexographically) cyclic shifts
  - strict sorting order due to \$
- First column (of course) has many consecutive character runs
- But also the last column has many consecutive character runs
  - 3 different shifts start with **lf** and end with **a**
  - sort groups **lf** lines together, and they all end with **a**

sorted shifts array

```
$ a l f e a t s a l f a l f a  
a $ a l f e a t s a l f a l f  
a l f a $ a l f e a t s a l f  
a l f a l f a $ a l f e a t s  
a l f e a t s a l f a l f a $  
a t s a l f a l f a $ a l f e  
e a t s a l f a l f a $ a l f  
f a $ a l f e a t s a l f a l  
f a l f a $ a l f e a t s a l  
f e a t s a l f a l f a $ a l  
l f a $ a l f e a t s a l f a  
l f a l f a $ a l f e a t s a  
l f e a t s a l f a l f a $ a  
s a l f a l f a $ a l f e a t  
t s a l f a l f a $ a l f e a
```

# BWT Algorithm and Example

$S = \mathbf{a} | \mathbf{l f} e a t s \mathbf{a} | \mathbf{l f a} | \mathbf{l f} a \$$

- Array of cyclic shifts
- Sort (lexographically) cyclic shifts
  - strict sorting order due to '\$'
- First column (of course) has many consecutive character runs
- But also the last column has many consecutive character runs
  - 3 different shifts start with **lf** and end with **a**
  - sort groups **lf** lines together, and they all end with **a**
  - could happen that another pattern will interfere
    - **hlf**d broken into **h** and **lfd**
  - chance of interference is small

sorted shifts array

```

$ a l f e a t s a l f a l f a
a $ a l f e a t s a l f a l f
a l f a $ a l f e a t s a l f
a l f a l f a $ a l f e a t s
a l f e a t s a l f a l f a $
a t s a l f a l f a $ a l f e
e a t s a l f a l f a $ a l f
f a $ a l f e a t s a l f a l
f a l f a $ a l f e a t s a l
f e a t s a l f a l f a $ a l
l f a $ a l f e a t s a l f a
l f a l f a $ a l f e a t s a
l f d          ... .. h
l f e a t s a l f a l f a $ a
s a l f a l f a $ a l f e a t
t s a l f a l f a $ a l f e a
    
```

# BWT Algorithm and Example

$S = \text{alfeatsalfalfa}\$$

- Sorted array of cyclic shifts
- First column is useless for encoding
  - cannot decode it
- Last column can be decoded
- BWT Encoding
  - last characters from sorted shifts
    - i.e. the last column

$C = \text{affs}\$ \text{eflllaaata}$

## sorted shifts array

```
$alfeatsalfalfa
a$alfeatsalfalff
alf$aalfeatsalff
alfalf$aalfeatss
alfeatsalfalfa
atsalfalf$aalfe
eatsalfalf$aalff
fa$alfeatsalfl
falf$aalfeatsalfl
featsalfalf$al
lfa$alfeatsalfa
lalf$aalfeatsa
lfeatsalfalf$aa
salfalf$aalfeatst
tsalfalf$aalfea
```

# BWT Fast Encoding: Efficient Sorting

$S = \text{a l f e a t s a l f a l f a } \$$

i	cyclic shift
0	a l f e a t s a l f a l f a \$
1	l f e a t s a l f a l f a \$ a
2	f e a t s a l f a l f a \$ a l
3	e a t s a l f a l f a \$ a l f
4	a t s a l f a l f a \$ a l f e
5	t s a l f a l f a \$ a l f e a
6	s a l f a l f a \$ a l f e a t
7	a l f a l f a \$ a l f e a t s
8	l f a l f a \$ a l f e a t s a
9	f a l f a \$ a l f e a t s a l
10	a l f a \$ a l f e a t s a l f
11	l f a \$ a l f e a t s a l f a
12	f a \$ a l f e a t s a l f a l
13	a \$ a l f e a t s a l f a l f
14	\$ a l f e a t s a l f a l f a

- Refer to a cyclic shift by the start index in the text, no need to write it out explicitly
- For sorting, letters after \$ do not matter

**a l f a l f a \$ a l f e a t s**  
<  
**l f a \$ a l f e a t s a l f a**



# BWT Fast Encoding: Efficient Sorting

$S = \text{a l f e a t s a l f a l f a } \$$

i	cyclic shift
0	a l f e a t s a l f a l f a \$
1	l f e a t s a l f a l f a \$ a
2	f e a t s a l f a l f a \$ a l
3	e a t s a l f a l f a \$ a l f
4	a t s a l f a l f a \$ a l f e
5	t s a l f a l f a \$ a l f e a
6	s a l f a l f a \$ a l f e a t
7	a l f a l f a \$ a l f e a t s
8	l f a l f a \$ a l f e a t s a
9	f a l f a \$ a l f e a t s a l
10	a l f a \$ a l f e a t s a l f
11	l f a \$ a l f e a t s a l f a
12	f a \$ a l f e a t s a l f a l
13	a \$ a l f e a t s a l f a l f
14	\$ a l f e a t s a l f a l f a

- Refer to a cyclic shift by the start index in the text, no need to write it out explicitly
- For sorting, letters after \$ do not matter

l f a \$ a l f e a t s a l f a  
<  
l f a l f a \$ a l f e a t s a

# BWT Fast Encoding: Efficient Sorting

$S = \text{a l f e a t s a l f a l f a } \$$

i	cyclic shift
0	<b>a l f e a t s a l f a l f a</b> \$
1	<b>l f e a t s a l f a l f a</b> \$ a
2	<b>f e a t s a l f a l f a</b> \$ a l
3	<b>e a t s a l f a l f a</b> \$ a l f
4	<b>a t s a l f a l f a</b> \$ a l f e
5	<b>t s a l f a l f a</b> \$ a l f e a
6	<b>s a l f a l f a</b> \$ a l f e a t
7	<b>a l f a l f a</b> \$ a l f e a t s
8	<b>l f a l f a</b> \$ a l f e a t s a
9	<b>f a l f a</b> \$ a l f e a t s a l
10	<b>a l f a</b> \$ a l f e a t s a l f
11	<b>l f a</b> \$ a l f e a t s a l f a
12	<b>f a</b> \$ a l f e a t s a l f a l
13	<b>a</b> \$ a l f e a t s a l f a l f
14	<b>\$</b> a l f e a t s a l f a l f a

- Refer to a cyclic shift by the start index in the text, no need to write it out explicitly
- For sorting, letters after \$ do not matter
- This is the same as sorting suffixes of  $S$
- We already know how to do it
  - exactly as for suffix arrays, with MSD-Radix-Sort
  - $O(n \log n)$  running time

# BWT Fast Encoding: Efficient Sorting

$S = \text{alfeatsalfalfa\$}$

i	cyclic shift
0	alfeatsalfalfa\$
1	lfeatsalfalfa\$a
2	featsalfalfa\$al
3	eatsalfalfa\$alf
4	atsalfalfa\$alfe
5	tsalfalfa\$alfea
6	salfalfa\$alfeat
7	alfalfa\$alfeats
8	lfa\$alfeatsalfa
9	fa\$alfeatsalf
10	alfa\$alfeatsalf
11	lfa\$alfeatsalfa
12	fa\$alfeatsalfal
13	a\$alfeatsalfalf
14	\$alfeatsalfalfa

j	$A^S[j]$	sorted cyclic shifts
0	14	\$alfeatsalfalfa
1	13	a\$alfeatsalfalf
2	10	alfa\$alfeatsalf
3	7	alfalfa\$alfeats
4	0	alfeatsalfalfa\$
5	4	atsalfalfa\$alfe
6	3	eatsalfalfa\$alf
7	12	fa\$alfeatsalfal
8	9	falfa\$alfeatsalf
9	2	featsalfalfa\$al
10	11	lfa\$alfeatsalf
11	8	lfa\$alfeatsalf
12	1	lfeatsalfalfa\$a
13	6	salfalfa\$alfeat
14	5	tsalfalfa\$alfea

# BWT Fast Encoding: Efficient Sorting

- Can read BWT encoding from suffix array in  $O(n)$  time

$S =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
a	l	f	e	a	t	s	a	l	f	a	l	f	a	\$

$A^s =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
14	13	10	7	0	4	3	12	9	2	11	8	1	6	5



cyclic shift starts at  $S[14]$

need last letter of that cyclic shift, it is at  $S[13]$

**a**

# BWT Fast Encoding: Efficient Sorting

- Can read BWT encoding from suffix array in  $O(n)$  time

$S =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
a	l	f	e	a	t	s	a	l	f	a	l	f	a	\$

$A^s =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
14	13	10	7	0	4	3	12	9	2	11	8	1	6	5



cyclic shift starts at  $S[13]$

need last letter of that cyclic shift, it is at  $S[12]$

**a f**

# BWT Fast Encoding: Efficient Sorting

- Can read BWT encoding from suffix array in  $O(n)$  time

$S =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
a	l	f	e	a	t	s	a	l	f	a	l	f	a	\$

$A^s =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
14	13	10	7	0	4	3	12	9	2	11	8	1	6	5



cyclic shift starts at  $S[10]$

need last letter of that cyclic shift, it is at  $S[9]$

**a f f**

# BWT Fast Encoding: Efficient Sorting

- Can read BWT encoding from suffix array in  $O(n)$  time

$S =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
a	l	f	e	a	t	s	a	l	f	a	l	f	a	\$

$A^s =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
14	13	10	7	0	4	3	12	9	2	11	8	1	6	5



cyclic shift starts at  $S[5]$

need last letter of that cyclic shift, it is at  $S[4]$

a f f s \$ e f l l l a a a t a

# BWT Fast Encoding: Efficient Sorting

- Can read BWT encoding from suffix array in  $O(n)$

$S =$

0	1	2	3	4	5	6	7	8	9
a	l	f	e	a	t	s	a	l	f

$A^s =$

0	1	2	3	4	5	6	7	8	9
14	13	10	7	0	4	3	12	9	2

cyclic shift starts at  $S[5]$

need last letter of that cyclic shift, it is at  $S[4]$

$j$	$A^s[j]$	
0	14	\$alfeatsalfalfa
1	13	a\$alfeatsalfalf
2	10	alf\$a\$alfeatsalf
3	7	alfalfa\$alfeats
4	0	alfeatsalfalfa\$
5	4	atsalfalfa\$alfe
6	3	eatsalfalfa\$alf
7	12	fa\$alfeatsalfal
8	9	falfa\$alfeatsal
9	2	featsalfalfa\$a
10	11	lfa\$alfeatsalf
11	8	lalfa\$alfeatsa
12	1	lfeatsalfalfa\$a
13	6	salfalfa\$alfeate
14	5	tsalfalfa\$alfeaa

a f f s \$ e f l l l a a a t a



# BWT Fast Encoding: Efficient Sorting

- Can read BWT encoding from suffix array in  $O(n)$  time

$S =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
a	l	f	e	a	t	s	a	l	f	a	l	f	a	\$

$A^s =$

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
14	13	10	7	0	4	3	12	9	2	11	8	1	6	5



cyclic shift starts at  $S[5]$

need last letter of that cyclic shift, it is at  $S[4]$

a f f s \$ e f l l l a a a t a

- Formula:  $C[i] = S[A^s[i] - 1]$ 
  - array is circular, i.e.  $S[-1] = S[n - 1]$

# BWT Decoding

$C = \text{affs}\$ \text{eflll} \text{aaata}$

- In unsorted shifts array, first column is  $S$
- So decoding = determining the first letter of each row in unsorted shifts array
  - when decoding, do not have unsorted shifts array

## unsorted shifts array

```
alfeatsalfalfa$  
lfeatsalfalfa$a  
ffeatsalfalfa$al  
eatsalfalfa$alf  
atsalfalfa$alfe  
ttsalfalfa$alfea  
salfalfa$alfeat  
alalfalfa$alfeats  
lfalfalfa$alfeatsa  
falfalfa$alfeatsal  
alfa$alfeatsalf  
lfa$alfeatsalf  
fa$alfeatsalfal  
a$alfeatsalfalf  
$alfeatsalfalfa
```



# BWT Decoding

$C = \text{affs}\$eflll\text{laaata}$

- Now have first and last columns of **sorted** shifts array
- Need the first column of **unsorted** shifts array

## unsorted shifts array

$S[0]$  **a** l f e a t s a l f a l f a \$  
 $S[1]$  **l** f e a t s a l f a l f a \$ a  
 $S[2]$  **f** e a t s a l f a l f a \$ a l  
 $S[3]$  **e** a t s a l f a l f a \$ a l f  
 .....

- Where in **sorted** shifts array are rows  $0, 1, \dots, n - 1$  of **unsorted** shifts array?
- Where is row 0 of **unsorted** shifts array?

## sorted shifts array

\$ . . . . . a  
 a . . . . . f  
 a . . . . . f  
 a . . . . . s  
 a . . . . . \$  
 a . . . . . e  
 e . . . . . f  
 f . . . . . l  
 f . . . . . l  
 f . . . . . l  
 l . . . . . a  
 l . . . . . a  
 l . . . . . a  
 s . . . . . t  
 t . . . . . a

# BWT Decoding

$C = \text{affs}\$ \text{eflll} \text{aaata}$

- Now have first and last columns of **sorted** shifts array
- Need the first column of **unsorted** shifts array

## unsorted shifts array

```

S[0] a l f e a t s a l f a l f a $
S[1] l f e a t s a l f a l f a $ a
S[2] f e a t s a l f a l f a $ a l
S[3] e a t s a l f a l f a $ a l f
.....

```

- Where in **sorted** shifts array are rows  $0, 1, \dots, n - 1$  of **unsorted** shifts array?
- Where is row 0 of **unsorted** shifts array?
  - must end with \$

## sorted shifts array

```

$ . . . . . a
a . . . . . f
a . . . . . f
a . . . . . s
a . . . . . $ row 0
a . . . . . e
e . . . . . f
f . . . . . l
f . . . . . l
f . . . . . l
l . . . . . a
l . . . . . a
l . . . . . a
s . . . . . t
t . . . . . a

```





# BWT Decoding

- Row 1 of unsorted shifts array ends with **a**
- Multiple rows end with **a**, which one is row 1 of unsorted shifts?
  - row 1 is a cyclic shift by one of row 0

## sorted shifts array

\$	.	.	.	.	.	.	.	<b>a</b>	<b>?</b>
a	.	.	.	.	.	.	.	f	
a	.	.	.	.	.	.	.	f	
a	.	.	.	.	.	.	.	s	
<b>a</b>	.	.	.	.	.	.	.	<b>\$</b>	<b>row 0</b>
a	.	.	.	.	.	.	.	e	
e	.	.	.	.	.	.	.	f	
f	.	.	.	.	.	.	.	l	
f	.	.	.	.	.	.	.	l	
f	.	.	.	.	.	.	.	l	
l	.	.	.	.	.	.	.	<b>a</b>	<b>?</b>
l	.	.	.	.	.	.	.	<b>a</b>	<b>?</b>
l	.	.	.	.	.	.	.	<b>a</b>	<b>?</b>
s	.	.	.	.	.	.	.	t	
t	.	.	.	.	.	.	.	<b>a</b>	<b>?</b>



# BWT Algorithm and Example

- Multiple rows end with **a**, which one is row 1 of unsorted shifts?
  - row 1 is a cyclic shift by one of row 0
- Rows that end with **a** are cyclic shifts by one of rows that start with **a**
- Rows that start with **a** appear in exactly the same order as their cyclic shifts by 1 (i.e. rows that end with **a**)

sorted shifts array

```
$alfeatsalffalffa  
a$alfeatsalffalff  
alfa$alfeatsalff  
alfalfa$alfeats  
alfeatsalffalffa$row 0  
atsalffalffa$alfe  
eatsalffalffa$alff  
fa$alfeatsalffal  
falffa$alfeatsal  
featsalffalffa$al  
lfa$alfeatsalffa  
lfalfa$alfeatsa  
lfeatsalffalffa$a  
salffalffa$alfeate  
tsalffalffa$alfea
```

# BWT Algorithm and Example

- for both group of patterns, sorting does not depend on **a**, and all other letters are the same between these two groups

rows starting with **a**

a	\$	a	l	f	e	a	t	s	a	l	f	a	l	f
a	l	f	a	\$	a	l	f	e	a	t	s	a	l	f
a	l	f	a	l	f	a	\$	a	l	f	e	a	t	s
a	l	f	e	a	t	s	a	l	f	a	l	f	a	\$
a	t	s	a	l	f	a	l	f	a	\$	a	l	f	e

row 0 of unsorted shifts is #4

their cyclic shifts by 1

\$	a	l	f	e	a	t	s	a	l	f	a	l	f	a
l	f	a	\$	a	l	f	e	a	t	s	a	l	f	a
l	f	a	l	f	a	\$	a	l	f	e	a	t	s	a
l	f	e	a	t	s	a	l	f	a	l	f	a	\$	a
t	s	a	l	f	a	l	f	a	\$	a	l	f	e	a

its cyclic shift by one is **also** #4

# BWT Algorithm and Example

- Multiple rows end with **a**, which one is row 1 of unsorted shifts?
  - row 1 is a cyclic shift by one of row 0
- Rows that end with **a** are cyclic shifts by one of rows that start with **a**
- Rows that start with **a** appear in exactly the same order as their cyclic shifts by 1 (i.e. rows that end with **a**)
- But direct 'counting' takes  $O(n)$  to find row 1

sorted shifts array

```

$alf eatsal falf a 1
1 a$alf eatsal falf
2 a lfa$alf eatsal f
3 a l falf a$alf eats
4 a l featsal falf a$ row 0
a t sal falf a$alf e
eatsal falf a$alf
fa$alf eatsal falf
falf a$alf eatsal
featsal falf a$alf
lfa$alf eatsal f a 2
l falf a$alf eats a 3
l featsal falf a$ a 4 row 1
sal falf a$alf eat
tsal falf a$alf e a
    
```









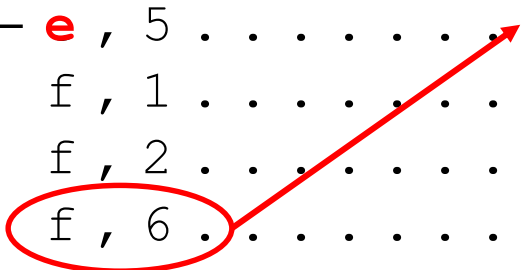
# BWT Decoding Continued

$C = \text{affs}\$eflll\text{aaata}$

$S = \text{alf}\mathbf{e}$

sorted shifts array

	\$ , 4	. . . . .	a , 0	
	a , 0	. . . . .	f , 1	
	a , 1 0	. . . . .	f , 2	
	a , 1 1	. . . . .	s , 3	
	a , 1 2	. . . . .	\$ , 4	row 0
	a , 1 4	. . . . .	e , 5	
$S[3] = \mathbf{e} \leftarrow$	$\mathbf{e} , 5$	. . . . .	f , 6	row 3
	f , 1	. . . . .	l , 7	
	f , 2	. . . . .	l , 8	
	f , 6	. . . . .	l , 9	row 2
	l , 7	. . . . .	a , 1 0	
	l , 8	. . . . .	a , 1 1	
	l , 9	. . . . .	a , 1 2	row 1
	s , 3	. . . . .	t , 1 3	
	t , 1 3	. . . . .	a , 1 4	





# BWT Decoding Continued

$C = \text{affs}\$eflll\text{aaata}$

$S = \text{alfe}\mathbf{a}$

sorted shifts array

	\$ , 4 . . . . . a , 0	
	a , 0 . . . . . f , 1	
	a , 1 0 . . . . . f , 2	
	a , 1 1 . . . . . s , 3	
	a , 1 2 . . . . . \$ , 4	<b>row 0</b>
$S[4] = \mathbf{a} \leftarrow$	<b>a</b> , 1 4 . . . . . e , 5	<b>row 4</b>
	e , 5 . . . . . f , 6	<b>row 3</b>
	f , 1 . . . . . l , 7	
	f , 2 . . . . . l , 8	
	f , 6 . . . . . l , 9	<b>row 2</b>
	l , 7 . . . . . a , 1 0	
	l , 8 . . . . . a , 1 1	
	l , 9 . . . . . a , 1 2	<b>row 1</b>
	s , 3 . . . . . t , 1 3	
	t , 1 3 . . . . . a , 1 4	

# BWT Decoding Continued

$C = \text{affs}\$eflll\text{aaata}$

$S = \text{alfeatsalfalfa}\$$

sorted shifts array

\$ , 4	. . . . .	a , 0	<b>row 14</b>
a , 0	. . . . .	f , 1	<b>row 13</b>
a , 1 0	. . . . .	f , 2	<b>row 10</b>
a , 1 1	. . . . .	s , 3	<b>row 7</b>
a , 1 2	. . . . .	\$ , 4	<b>row 0</b>
a , 1 4	. . . . .	e , 5	<b>row 4</b>
e , 5	. . . . .	f , 6	<b>row 3</b>
f , 1	. . . . .	l , 7	<b>row 12</b>
f , 2	. . . . .	l , 8	<b>row 9</b>
f , 6	. . . . .	l , 9	<b>row 2</b>
l , 7	. . . . .	a , 1 0	<b>row 11</b>
l , 8	. . . . .	a , 1 1	<b>row 8</b>
l , 9	. . . . .	a , 1 2	<b>row 1</b>
s , 3	. . . . .	t , 1 3	<b>row 6</b>
t , 1 3	. . . . .	a , 1 4	<b>row 5</b>

# BWT Decoding Pseudocode

*BWT::decoding*( $C, S$ )

$C$ : string of characters over alphabet  $\Sigma_C$ , one of which is \$

$S$ : output stream

initialize array  $A$  // leftmost column

**for** all indices  $i$  of  $C$

$A[i] \leftarrow (C[i], i)$  // store character and index

stably sort  $A$  by character (the first aspect)

**for** all indices  $j$  of  $C$  // find \$

**if**  $C[j] = \$$  **break**

**do**

$S.append(\text{character stored in } A[j])$

$j \leftarrow \text{index stored in } A[j]$

**while** appended character is not \$

- What sorting algorithm would you use?

# BWT and bzip2 Discussion

## ■ BWT

- encoding cost
  - $O(n \log n)$  with special sorting algorithm
  - read encoding from the suffix array
- decoding cost
  - $O(n + |\Sigma_S|)$ 
    - faster than encoding
- encoding and decoding both use  $O(n)$  space
- they need all of the text (no streaming possible)
- can use on blocks of text (block compression method)

## ■ bzip2

- encoding cost:  $O(n [\log n + |\Sigma|])$  with a big multiplicative constant
- decoding cost:  $O(n|\Sigma|)$  with a big multiplicative constant
- tends to be slower than other methods but gives better compression

# Compression Summary

<b>Huffman</b>	<b>Run-length encoding</b>	<b>Lempel-Ziv-Welch</b>	<b>bzip2</b> (uses Burrows-Wheeler)
variable-length	variable-length	fixed-length	multi-step
single-character	multi-character	multi-character	multi-step
2-pass, must send dictionary	1-pass	1-pass	not streamable
optimal 01-prefix-code	good on long runs (e.g., pictures)	good on English text	better on English text
requires uneven frequencies	requires runs	requires repeated substrings	requires repeated substrings
rarely used directly	rarely used directly	frequently used	used but slow
part of pzip, JPEG, MP3	fax machines, TIFF	GIF, some variants of PDF, compress	bzip2 and variants