

The LZW Approach

Key idea: Each ***codeword*** represents multiple symbols.

- Start with ‘trivial’ codeword table where each codeword corresponds to one ASCII symbol.
- Every time a codeword X is used, record a new codeword Y corresponding to X concatenated with the next symbol.

Example begins on next slide for

B="aababcabcdabcdeabcdefabcdefgabcdefgh"

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefgabcdefgh”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣

$C(B) =$

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefabcdeffgabcdeffgh”

- Best prefix match in codeword table is a, so output 0x61.
- Also add aa to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣

$$C(B) = 0x61$$

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefabcdeffgabcdeffgh”

- Best prefix match in codeword table is a, so output 0x61.
- Also add aa to table.

$$C(B) = 0x61$$

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcabcdabcdeabcdefgabcdefgabcdefg”

- Best prefix match in codeword table is a, so output 0x61.
- Also add ab to table.

$$C(B) = 0x6161$$

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcabcdabcdeabcdefgabcdefgabcdefg”

- Best prefix match in codeword table is a, so output 0x61.
- Also add ab to table.

$$C(B) = 0x6161$$

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbabcdabcdeabcdefabcdefgabcdefgh”

- Best prefix match in codeword table is b, so output 0x62.
- Also add ba to table.

$$C(B) = 0x616162$$

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefabcdeffgabcdeffgh”

- Best prefix match in codeword table is b, so output 0x62.
- Also add ba to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
------	----

$C(B) = 0x616162$

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcabcdabcdeabcdefgabcdefggh”

- Best prefix match in codeword table is ab, so output 0x81.
- Also add abc to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc

$C(B) = 0x61616281$

Only half as many bits!

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcabcdabcdeabcdefgabcdefg”

- Best prefix match in codeword table is c, so output 0x63.
- Also add ca to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca

$C(B) = 0x6161628163$

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcabcdabcdeabcdefabcdeffgabcdeffgh”

- Best prefix match in codeword table is ???, so output ???.
- Also add ??? to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
???	???

$C(B) = 0x6161628163$

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcabcdabcdeabcdefabcdegh”

- Best prefix match in codeword table is abc, so output 0x83.
- Also add abcd to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd

$C(B) = 0x616162816383$

LZW Prefix Matching

How should we store our codeword table to easily support fast prefix matching?

- B = “aababcabcdabcdeabcdefgabcdefggh”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd

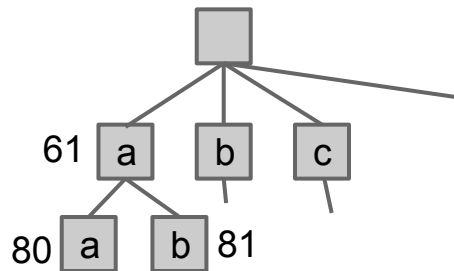
LZW Prefix Matching

How should we store our codeword table to easily support fast prefix matching? A trie mapping strings to codewords.

- B = “aababcabcdabcdeabcdefgabcdefggh”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd



The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefabcdefgabcdefgh”

- Best prefix match in codeword table is d, so output 0x64.
- Also add da to table.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd
0x86	da

$C(B) = 0x61616281638364$

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefabcdefgabcdefgh”

- What will be the codeword for abcdefg (to think about when reviewing this material later)? What will be C(B)?

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd
0x86	da

C(B) = 0x61616281638364

The LZW Approach: Codewords for Multiple Symbols

Example: B = “aababcbcdabcdeabcdefabcdefgabcdefgh”

- abcdefg will be 0x8b.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd
0x86	da
0x87	abcde
0x88	ea
0x89	abcdef

0x8a	fa
0x8b	abcdefg
0x8c	ga
0x8d	abcdefgh

C(B) =

0x616162816383648565

876689678b68

288 bits.

120 bits.

Decoding

The Decompression Challenge

Suppose we're given a bitstream that we need to decompress.
 $C(B) = 0x616162816383$, but don't have table used for encoding.

- B clearly starts with aab, but what is 0x81?

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣

Decompressing LZW

Key idea: After processing each codeword, add the codeword that would have been added by the *previous* character.

- Start with 'trivial' codeword table where each codeword corresponds to one ASCII symbol.
- When codeword X is encountered, return appropriate symbols $S(X)$ given in codeword table.
- Given consecutive codewords X_1 and X_2 add codeword corresponding to $[S(X_1) + \text{firstCharacterOf}(S(X_2))]$

Example on following slides for : $C(B) = 0x61616281638364$

LZW Decompression

$C(B) = 0x\underline{61}616281638364$

- $0x\underline{61}$ is a in the codeword table.
- No previous codeword so add nothing.

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣

$B = \text{"a"}$

LZW Decompression

$C(B) = 0x61\underline{61}6281638364$

- $0x\underline{61}$ is a in the codeword table, so output “a”
- $S(0x61) = \text{“a”}$ and $S(0x61) = \text{“a”}$, so add “aa”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa

$B = \text{“a}\underline{\text{a}}\text{”}$

LZW Decompression

$C(B) = 0x616\underline{2}81638364$

- $0x\underline{62}$ is b in the codeword table, so output “b”.
- $S(0x61)=\text{“a”}$ and $S(0x62)=\text{“b”}$, so add “ab”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

$B = \text{“aab\underline{b}”}$

LZW Decompression

$C(B) = 0x616162\underline{81}638364$

- $0x\underline{81}$ is ab in the codeword table, so output “ab”
- $S(0x62) = \text{“b”}$ and $S(0x81) = \text{“ab”}$, so add “ba”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
------	----

$B = \text{“aabab”}$

LZW Decompression

$C(B) = 0x61616281\underline{63}8364$

- $0x\underline{63}$ is c in the codeword table, so output “c”
- $S(0x81) = \text{“ab”}$ and $S(0x63) = \text{“c”}$, so add “abc”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc

$B = \text{“aabab}\underline{c}\text{”}$

LZW Decompression

$C(B) = 0x6161628163\underline{83}64$

- $0x\underline{83}$ is ??? in the codeword table, so output “???”
- $S(???) = \text{“}???\text{”}$ and $S(???) = \text{“}???\text{”}$, so add “???”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc

$B = \text{“}aababc\underline{???}\text{”}$

LZW Decompression

$C(B) = 0x6161628163\underline{83}64$

- $0x\underline{83}$ is abc in the codeword table, so output “abc”
- $S(0x63) = \text{“c”}$ and $S(0x83) = \text{“abc”}$, so add “ca”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca

$B = \text{“aababca}\underline{\text{bc}}\text{”}$

LZW Decompression

$C(B) = 0x616162816383\underline{64}$

- $0x\underline{64}$ is d in the codeword table, so output “d”
- $S(0x83) = \text{“abc”}$ and $S(0x64) = \text{“d”}$, so add “abcd”

0x61	a
0x62	b
0x63	c
...	
0x7e	~
0x7f	␣
0x80	aa
0x81	ab

0x82	ba
0x83	abc
0x84	ca
0x85	abcd

$B = \text{“aababcbcd\underline{d}”}$

LZW Decompression: Special Case

LZW Compression (exactly like before)

B="OXOXOXO"

- Best prefix match is O, so output 0x4F.
- Add codeword "OX"

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX

$C(B) = 0x4F$

Compressing this is exactly as we did before, but this one will turn out to be slightly trickier to decompress (as we'll soon see).

LZW Compression (exactly like before)

B="OXOXOXO"

- Best prefix match is X, so output 0x58.
- Add codeword "XO"

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

$C(B) = 0x4F58$

LZW Compression (exactly like before)

B="OXOXOXO"

- Best prefix match is OX, so output 0x80.
- Add codeword "OXO"

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

0x82	OXO
------	-----

$C(B) = 0x4F5880$

LZW Compression (exactly like before)

B="OXOXOXO"

- Best prefix match is OXO, so output 0x82.
- Don't add anything.

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

0x82	OXO
------	-----

$C(B) = 0x4F588082$

LZW Compression (exactly like before)

B="OXOXOXO"

- Best prefix match is OXO, so output 0x82.
- Don't add anything.

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

0x82	OXO
------	-----

$C(B) = 0x4F588082$

LZW Decompression

$C(B) = 0x\underline{4F}588082$

- $0x\underline{4F}$ is O in the codeword table, so output “O”
- No previous codeword, so output nothing.

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	⌘

$B = \text{“}\underline{O}\text{”}$

LZW Decompression

$C(B) = 0x4F\underline{58}8082$

- $0x\underline{58}$ is X in the codeword table, so output “X”
- $S(0x4F) = \text{“O”}$ and $S(0x58) = \text{“X”}$, so add “OX”

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX

$B = \text{“O}\underline{\text{X}}\text{”}$

LZW Decompression

$C(B) = 0x4F58\underline{8}082$

- $0x\underline{8}0$ is OX in the codeword table, so output “OX”
- $S(0x58) = "X"$ and $S(0x80) = "OX"$, so add “XO”

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

$B = "OX\underline{OX}"$

LZW Decompression

$C(B) = 0x4F588082$

- $0x82$ is ... uh oh.

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

The problem is that the compression algorithm looked ahead one character to create new codewords. When decompressing we look BACK to create new codewords!

Solution: Figure out what 82 is going to be by looking back.

$B = \text{“OXOX?????”}$

LZW Decompression

$C(B) = 0x4F588082$

- $0x82$ is ****.
- $S(0x80)$ is "OX" and $S(0x82)$ is "****", so add "****"

$B = "OXOX\underline{****}"$

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

All the same thing!

LZW Decompression

$C(B) = 0x4F588082$

- $0x82$ is ****.

- $S(0x80)$ is "OX" and $S(0x82)$ is "OX**", so add "OX**"

We know that **** must start with OX.



$B = \text{"OXOX****"}$

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

LZW Decompression

$C(B) = 0x4F5880\underline{82}$

- $0x\underline{82}$ is ****.
- $S(0x80)$ is "OX" and $S(0x82)$ is "OX**", so add "OX**"

We know that ** is just the leftmost character of $0x82$!



$B = \text{"OXOX"}\underline{\text{****}}$

...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

LZW Decompression

$C(B) = 0x4F5880\underline{82}$

- $0x\underline{82}$ is ****.
- $S(0x80)$ is “OX” and $S(0x82)$ is “OXO”, so add “OXO”

We also know that ** is just the leftmost character of $0x82$!



...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

0x82	OXO
------	-----

$B = \text{“OXOX}\underline{\text{****}}\text{”}$

LZW Decompression

$C(B) = 0x4F5880\underline{82}$

- $0x\underline{82}$ is OXO.
- $S(0x80)$ is “OX” and $S(0x82)$ is “OXO”, so add “OXO”

We also know that ** is just the leftmost character of $0x82$!



...	...
0x4F	O
...	...
0x58	X
...	...
0x7f	␣
0x80	OX
0x81	XO

0x82	OXO
------	-----

$B = \text{“OXOXOXO”}$