

CSE 490 GZ Introduction to Data Compression Winter 2004

Sequitur

Sequitur

- Nevill-Manning and Witten, 1996.
- Uses a context-free grammar (without recursion) to represent a string.
- The grammar is inferred from the string.
- If there is structure and repetition in the string then the grammar may be very small compared to the original string.
- Clever encoding of the grammar yields impressive compression ratios.
- Compression plus structure!

CSE 490gz - Lecture 9 - Winter 2004

2

Context-Free Grammars

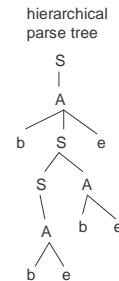
- Invented by Chomsky in 1959 to explain the grammar of natural languages.
- Also invented by Backus in 1959 to generate and parse Fortran.
- Example:
 - terminals: b, e
 - non-terminals: S, A
 - Production Rules:
 - $S \rightarrow SA$, $S \rightarrow A$, $A \rightarrow bSe$, $A \rightarrow be$
 - S is the start symbol

CSE 490gz - Lecture 9 - Winter 2004

3

Context-Free Grammar Example

- $S \rightarrow SA$
 - $S \rightarrow A$
 - $A \rightarrow bSe$
 - $A \rightarrow be$
- derivation of bbebee
- Example: b and e matched as parentheses



CSE 490gz - Lecture 9 - Winter 2004

4

Arithmetic Expressions

- $S \rightarrow S + T$
 - $S \rightarrow T$
 - $T \rightarrow T^*F$
 - $T \rightarrow F$
 - $F \rightarrow a$
 - $F \rightarrow (S)$
- derivation of $a^* (a + a) + a$
- parse tree
-
- ```

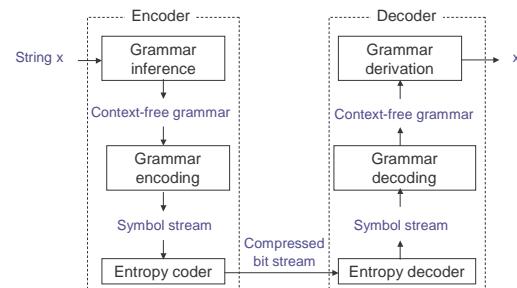
graph TD
 S1[S] --- T1[T]
 T1 --- T2[T^*]
 T2 --- F1[F]
 F1 --- a1[a]
 T2 --- T3[T]
 T3 --- F2[F]
 F2 --- a2[a]
 T3 --- T4[T]
 T4 --- F3[F]
 F3 --- a3[a]
 T4 --- T5[T]
 T5 --- F4[F]
 F4 --- a4[a]

```

CSE 490gz - Lecture 9 - Winter 2004

5

### Overview of Grammar Compression



CSE 490gz - Lecture 9 - Winter 2004

6

## Sequitur Principles

- Digram Uniqueness:
  - no pair of adjacent symbols (digram) appears more than once in the grammar.
- Rule Utility:
  - Every production rule is used more than once.
- These two principles are maintained as an invariant while inferring a grammar for the input string.

CSE 490gz - Lecture 9 - Winter 2004

7

## Sequitur Example (1)

bbebeebebebbeebee

$S \rightarrow b$

CSE 490gz - Lecture 9 - Winter 2004

8

## Sequitur Example (2)

bbebeebebebbeebee

$S \rightarrow bb$

CSE 490gz - Lecture 9 - Winter 2004

9

## Sequitur Example (3)

bbebeebebebbeebee

$S \rightarrow bbe$

CSE 490gz - Lecture 9 - Winter 2004

10

## Sequitur Example (4)

bbebbeebebebbee

$S \rightarrow bb eb$

CSE 490gz - Lecture 9 - Winter 2004

11

## Sequitur Example (5)

bbebebebebbee

$S \rightarrow bbebe$

Enforce digram uniqueness.  
be occurs twice.  
Create new rule A  $\rightarrow$  be.

CSE 490gz - Lecture 9 - Winter 2004

12

### Sequitur Example (6)

bbebeebebebbebee

$S \rightarrow bAA$   
 $A \rightarrow be$

CSE 490gz - Lecture 9 - Winter 2004

13

### Sequitur Example (7)

bbebeebebbebee

$S \rightarrow bAAe$   
 $A \rightarrow be$

CSE 490gz - Lecture 9 - Winter 2004

14

### Sequitur Example (8)

bbebeebebebbebee

$S \rightarrow bAAeb$   
 $A \rightarrow be$

CSE 490gz - Lecture 9 - Winter 2004

15

### Sequitur Example (9)

bbebebebebbee

$S \rightarrow bAAebe$   
 $A \rightarrow be$

Enforce digram uniqueness.  
be occurs twice.  
Use existing rule  $A \rightarrow be$ .

CSE 490gz - Lecture 9 - Winter 2004

16

### Sequitur Example (10)

bbebebebebbee

$S \rightarrow bAAeA$   
 $A \rightarrow be$

CSE 490gz - Lecture 9 - Winter 2004

17

### Sequitur Example (11)

bbebebebebbee

$S \rightarrow bAAeAb$   
 $A \rightarrow be$

CSE 490gz - Lecture 9 - Winter 2004

18

### Sequitur Example (12)

bbebeebebebbebee

$S \rightarrow bAAeAbe$   
 $A \rightarrow be$

Enforce digram uniqueness.  
be occurs twice.  
Use existing rule  $A \rightarrow be$ .

CSE 490gz - Lecture 9 - Winter 2004

19

### Sequitur Example (13)

bbebeebebebbebee

$S \rightarrow bAAeAA$   
 $A \rightarrow be$

Enforce digram uniqueness  
AA occurs twice.  
Create new rule  $B \rightarrow AA$ .

CSE 490gz - Lecture 9 - Winter 2004

20

### Sequitur Example (14)

bbebeebebebbebee

$S \rightarrow bBeB$   
 $A \rightarrow be$   
 $B \rightarrow AA$

CSE 490gz - Lecture 9 - Winter 2004

21

### Sequitur Example (15)

bbebeebebebbebee

$S \rightarrow bBeBb$   
 $A \rightarrow be$   
 $B \rightarrow AA$

CSE 490gz - Lecture 9 - Winter 2004

22

### Sequitur Example (16)

bbebeebebebbebeee

$S \rightarrow bBeBbb$   
 $A \rightarrow be$   
 $B \rightarrow AA$

CSE 490gz - Lecture 9 - Winter 2004

23

### Sequitur Example (17)

bbebeebebebbebeee

$S \rightarrow bBeBbbe$   
 $A \rightarrow be$   
 $B \rightarrow AA$

Enforce digram uniqueness.  
be occurs twice.  
Use existing rule  $A \rightarrow be$ .

CSE 490gz - Lecture 9 - Winter 2004

24

### Sequitur Example (18)

bbebeebebebbebee

S → bBeBbA  
A → be  
B → AA

CSE 490gz - Lecture 9 - Winter 2004

25

### Sequitur Example (19)

bbebeebebebbeee

S → bBeBbAb  
A → be  
B → AA

CSE 490gz - Lecture 9 - Winter 2004

26

### Sequitur Example (20)

bbebeebebebbeee

S → bBeBbAbe  
A → **be**  
B → AA

Enforce digram uniqueness.  
be occurs twice.  
Use existing rule A → be.

CSE 490gz - Lecture 9 - Winter 2004

27

### Sequitur Example (21)

bbebeebebebbeee

S → bBeBbAA  
A → be  
B → **AA**

Enforce digram uniqueness.  
AA occurs twice.  
Use existing rule B → AA.

CSE 490gz - Lecture 9 - Winter 2004

28

### Sequitur Example (22)

bbebeebebebbeee

S → **bBeBbB**  
A → be  
B → AA

Enforce digram uniqueness.  
bB occurs twice.  
Create new rule C → bB.

CSE 490gz - Lecture 9 - Winter 2004

29

### Sequitur Example (23)

bbebeebebebbeee

S → CeBC  
A → be  
B → AA  
C → bB

CSE 490gz - Lecture 9 - Winter 2004

30

## Sequitur Example (24)

bbebeebebebbeebee

$S \rightarrow \text{CeBCe}$  Enforce digram uniqueness.  
 $A \rightarrow \text{be}$  Ce occurs twice.  
 $B \rightarrow \text{AA}$  Create new rule  $D \rightarrow \text{Ce}$ .  
 $C \rightarrow \text{bB}$

CSE 490gz - Lecture 9 - Winter 2004

31

## Sequitur Example (25)

bbebeebebebbeebee

$S \rightarrow \text{DBD}$  Enforce rule utility.  
 $A \rightarrow \text{be}$  C occurs only once.  
 $B \rightarrow \text{AA}$  Remove  $C \rightarrow \text{bB}$ .  
 $C \rightarrow \text{bB}$   
 $D \rightarrow \text{Ce}$

CSE 490gz - Lecture 9 - Winter 2004

32

## Sequitur Example (26)

bbebeebebebbeebee

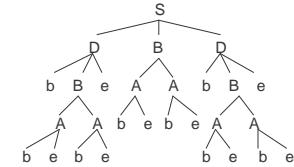
$S \rightarrow \text{DBD}$   
 $A \rightarrow \text{be}$   
 $B \rightarrow \text{AA}$   
 $D \rightarrow \text{bBe}$

CSE 490gz - Lecture 9 - Winter 2004

33

## The Hierarchy

bbebeebebebbeebee



Is there compression? In this small example, probably not.

CSE 490gz - Lecture 9 - Winter 2004

34

## Sequitur Algorithm

Input the first symbol  $s$  to create the production  $S \rightarrow s$ ;  
 repeat  
 match an existing rule:  
 $A \rightarrow \dots XY\dots \quad \rightarrow \quad A \rightarrow \dots B\dots$   
 $B \rightarrow XY \quad \rightarrow \quad B \rightarrow XY$   
 create a new rule:  
 $A \rightarrow \dots XY\dots \quad \rightarrow \quad A \rightarrow \dots C\dots$   
 $B \rightarrow \dots XY\dots \quad \rightarrow \quad B \rightarrow \dots C\dots$   
 remove a rule:  
 $A \rightarrow \dots B\dots \quad \rightarrow \quad$   
 $B \rightarrow X_1X_2\dots X_k \quad \rightarrow \quad A \rightarrow \dots X_1X_2\dots X_k \dots$   
 input a new symbol:  
 $S \rightarrow X_1X_2\dots X_k \quad \rightarrow \quad S \rightarrow X_1X_2\dots X_k s$   
 until no symbols left

CSE 490gz - Lecture 9 - Winter 2004

35

## Exercise

Use Sequitur to construct a grammar for  $aaaaaaaaaa = a^{10}$

CSE 490gz - Lecture 9 - Winter 2004

36

## Complexity

- The number of non-input sequitur operations applied  $< 2n$  where  $n$  is the input length.
- Since each operation takes constant time, sequitur is a linear time algorithm

CSE 490gz - Lecture 9 - Winter 2004

37

## Amortized Complexity Argument

- Let  $m = \#$  of non-input sequitur operations.  
Let  $n =$  input length. Show  $m \leq 2n$ .
- Let  $s =$  the sum of the right hand sides of all the production rules. Let  $r =$  the number of rules.
- We evaluate  $2s - r$ .
- Initially  $2s - r = 1$  because  $s = 1$  and  $r = 1$ .
- $2s - r > 0$  at all times because each rule has at least 1 symbol on the right hand side.

CSE 490gz - Lecture 9 - Winter 2004

38

## Sequitur Rule Complexity

- Digram Uniqueness - match an existing rule.
- $$\begin{array}{ccc} A \rightarrow \dots XY\dots & \longrightarrow & A \rightarrow \dots B\dots \\ B \rightarrow XY & & B \rightarrow XY \end{array} \quad \begin{array}{ccc} s & r & 2s - r \\ -1 & 0 & -2 \end{array}$$
- Digram Uniqueness - create a new rule.
- $$\begin{array}{ccc} A \rightarrow \dots XY\dots & \longrightarrow & A \rightarrow \dots C\dots \\ B \rightarrow \dots XY\dots & & B \rightarrow \dots C\dots \\ & & C \rightarrow XY \end{array} \quad \begin{array}{ccc} s & r & 2s - r \\ 0 & 1 & -1 \end{array}$$
- Rule Utility - Remove a rule.
- $$\begin{array}{ccc} A \rightarrow \dots B\dots & \longrightarrow & A \rightarrow \dots X_1X_2\dots X_k\dots \\ B \rightarrow X_1X_2\dots X_k & & B \rightarrow X_1X_2\dots X_k \end{array} \quad \begin{array}{ccc} s & r & 2s - r \\ -1 & -1 & -1 \end{array}$$

CSE 490gz - Lecture 9 - Winter 2004

39

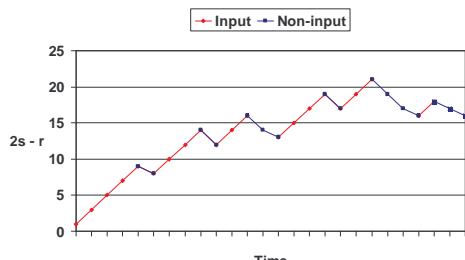
## Amortized Complexity Argument

- $2s - r \geq 0$  at all times because each rule has at least 1 symbol on the right hand side.
- $2s - r$  increases by 2 for every input operation.
- $2s - r$  decreases by at least 1 for each non-input sequitur rule applied.
- $n =$  number of input symbols  
 $m =$  number of non-input operations
- $2n - m \geq 0$ .  $m \leq 2n$ .

CSE 490gz - Lecture 9 - Winter 2004

40

## Amortized Complexity Argument



CSE 490gz - Lecture 9 - Winter 2004

41

## Linear Time Algorithm

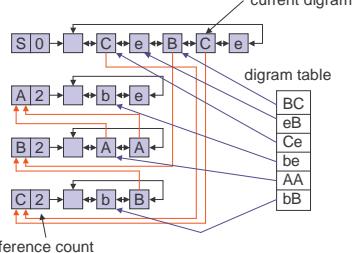
- There is a data structure to implement all the sequitur operations in constant time.
  - Production rules in an array of doubly linked lists.
  - Each production rule has reference count of the number of times used.
  - Each nonterminal points to its production rule.
  - Digrams stored in a hash table for quick lookup.

CSE 490gz - Lecture 9 - Winter 2004

42

## Data Structure Example

$S \rightarrow CeBCe$   
 $A \rightarrow be$   
 $B \rightarrow AA$   
 $C \rightarrow bB$



CSE 490gz - Lecture 9 - Winter 2004

43

## Basic Encoding a Grammar

|         |                                                                                        |             |                                                              |                                        |
|---------|----------------------------------------------------------------------------------------|-------------|--------------------------------------------------------------|----------------------------------------|
| Grammar | $S \rightarrow DBD$<br>$A \rightarrow be$<br>$B \rightarrow AA$<br>$D \rightarrow bBe$ | Symbol Code | <b>A</b> 010<br><b>B</b> 011<br><b>D</b> 100<br><b>#</b> 101 | b 000<br>e 001<br>No code for S needed |
|---------|----------------------------------------------------------------------------------------|-------------|--------------------------------------------------------------|----------------------------------------|

### Grammar Code

D B D # b e # A A # b B e  
100 011 100 101 000 001 101 010 010 101 000 011 001 39 bits

$$|\text{Grammar Code}| = (s + r - 1) \lceil \log_2(r + a) \rceil$$

r = number of rules

s = sum of right hand sides

a = number in original symbol alphabet

CSE 490gz - Lecture 9 - Winter 2004

44

## Better Encoding of the Grammar

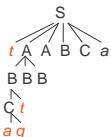
- Nevill-Manning and Witten suggest a more efficient encoding of the grammar that uses LZ77 ideas.
  - Send the right hand side of the S production.
  - The first time a nonterminal is sent, its right hand side is transmitted instead.
  - The second time a nonterminal is sent as a triple  $[i, j, d]$ , which says the right hand side starts at position  $j$  in production rule  $i$  and is  $d$  long. A new production rule is then added to a dictionary.
  - Subsequently, the nonterminal is represented by the index of the production rule.

CSE 490gz - Lecture 9 - Winter 2004

45

## Transmission Example

|                        |                      |
|------------------------|----------------------|
| $S \rightarrow tAABCa$ | T = Transmitted      |
| $A \rightarrow BBB$    |                      |
| $B \rightarrow Ct$     | T tagt               |
| $C \rightarrow ag$     | $X_0 tagt$ $l_0 = 4$ |



CSE 490gz - Lecture 9 - Winter 2004

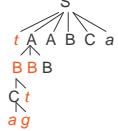
46

## Transmission Example

|                        |                 |
|------------------------|-----------------|
| $S \rightarrow tAABCa$ | T = Transmitted |
| $A \rightarrow BBB$    |                 |

T tagt [0, 1, 3]

$X_0 t X_1 X_1$   
 $X_1 agt$        $l_0 = 3$   
 $l_1 = 3$

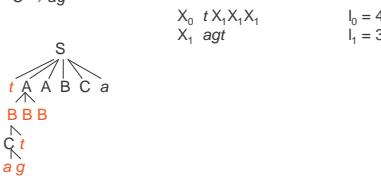


CSE 490gz - Lecture 9 - Winter 2004

47

## Transmission Example

|                        |                    |
|------------------------|--------------------|
| $S \rightarrow tAABCa$ | T = Transmitted    |
| $A \rightarrow BBB$    |                    |
| $B \rightarrow Ct$     | T tagt [0, 1, 3] 1 |

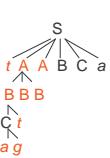


CSE 490gz - Lecture 9 - Winter 2004

48

## Transmission Example

$S \rightarrow tAABCa$   
 $A \rightarrow BBB$   
 $B \rightarrow Ct$   
 $C \rightarrow ag$



$T = \text{Transmitted}$

$T \quad tagt[0, 1, 3] 1 [0, 1, 3]$

$X_0 \quad t X_2 X_2$   
 $X_1 \quad agt$   
 $X_2 \quad X_1 X_1 X_1$

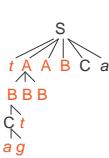
$l_0 = 3$   
 $l_1 = 3$   
 $l_2 = 3$

CSE 490gz - Lecture 9 - Winter 2004

49

## Transmission Example

$S \rightarrow tAABCa$   
 $A \rightarrow BBB$   
 $B \rightarrow Ct$   
 $C \rightarrow ag$



$T = \text{Transmitted}$

$T \quad tagt[0, 1, 3] 1 [0, 1, 3] 1$

$X_0 \quad t X_2 X_2 X_1$   
 $X_1 \quad agt$   
 $X_2 \quad X_1 X_1 X_1$

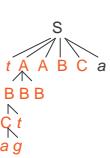
$l_0 = 4$   
 $l_1 = 3$   
 $l_2 = 3$

CSE 490gz - Lecture 9 - Winter 2004

50

## Transmission Example

$S \rightarrow tAABCa$   
 $A \rightarrow BBB$   
 $B \rightarrow Ct$   
 $C \rightarrow ag$



$T = \text{Transmitted}$

$T \quad tagt[0, 1, 3] 1 [0, 1, 3] 1 [1, 0, 2]$

$X_0 \quad t X_2 X_2 X_1 X_3$   
 $X_1 \quad X_3 t$   
 $X_2 \quad X_1 X_1 X_1$   
 $X_3 \quad ag$

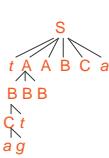
$l_0 = 5$   
 $l_1 = 2$   
 $l_2 = 3$   
 $l_3 = 2$

CSE 490gz - Lecture 9 - Winter 2004

51

## Transmission Example

$S \rightarrow tAABCa$   
 $A \rightarrow BBB$   
 $B \rightarrow Ct$   
 $C \rightarrow ag$



$T = \text{Transmitted}$

$T \quad tagt[0, 1, 3] 1 [0, 1, 3] 1 [1, 0, 2] a$

$X_0 \quad t X_2 X_2 X_1 X_3 a$   
 $X_1 \quad X_3 t$   
 $X_2 \quad X_1 X_1 X_1$   
 $X_3 \quad ag$

$l_0 = 6$   
 $l_1 = 2$   
 $l_2 = 3$   
 $l_3 = 2$

CSE 490gz - Lecture 9 - Winter 2004

52

## Kieffer-Yang Improvement

- Kieffer and Yang developed a theoretical framework for studying these types of grammars in 2000.
  - KY is universal; it achieves entropy in the limit
- Add to sequitur Reduction Rule 5:

$$\begin{array}{ll}
 S \rightarrow AB & S \rightarrow AA \\
 A \rightarrow CD & A \rightarrow CD \\
 B \rightarrow aE & \Rightarrow \quad B \rightarrow aE \quad \text{Adding this} \\
 C \rightarrow ab & C \rightarrow ab \quad \text{constraint} \\
 D \rightarrow cd & D \rightarrow cd \quad \text{makes sequitur} \\
 E \rightarrow bD & E \rightarrow bD \quad \text{universal.}
 \end{array}$$

$\langle A \rangle = \langle B \rangle = abcd$

CSE 490gz - Lecture 9 - Winter 2004

53

## Compression Quality

- Neville-Manning and Witten 1997

|       | size   | comp | gzip        | sequitur    | PPMC        | bzip2       |
|-------|--------|------|-------------|-------------|-------------|-------------|
| bib   | 111261 | 3.35 | 2.51        | <b>2.48</b> | <b>2.12</b> | 1.98        |
| book  | 768771 | 3.46 | 3.35        | <b>2.82</b> | <b>2.52</b> | 2.42        |
| geo   | 102400 | 6.08 | 5.34        | <b>4.74</b> | <b>5.01</b> | 4.45        |
| obj2  | 246814 | 4.17 | <b>2.63</b> | <b>2.68</b> | 2.77        | 2.48        |
| pic   | 513216 | 0.97 | <b>0.82</b> | <b>0.90</b> | 0.98        | 0.78        |
| prog2 | 38611  | 3.87 | <b>2.68</b> | 2.83        | <b>2.49</b> | <b>2.53</b> |

■ First; ■ Second; ■ Third.

Files from the Calgary Corpus

Units in bits per character (8 bits)

compress - based on LZ77

gzip - based on LZ77

PPMC - adaptive arithmetic coding with context

bzip2 - Burrows-Wheeler block sorting

CSE 490gz - Lecture 9 - Winter 2004

54

### Notes on Sequitur

- Yields compression and hierarchical structure simultaneously.
- With clever encoding is competitive with the best of the standards.
- The grammar size is not close to approximation algorithms
  - Upper =  $O((n/\log n)^{3/4})$ ; Lower =  $\Omega(n^{1/3})$ . (Lehman, 2002)
- *But!* Practical linear time encoding and decoding.

CSE 490gz - Lecture 9 - Winter 2004

55

### Other Grammar Based Methods

- Longest Match
- Most frequent digram
- Match producing the best compression

CSE 490gz - Lecture 9 - Winter 2004

56