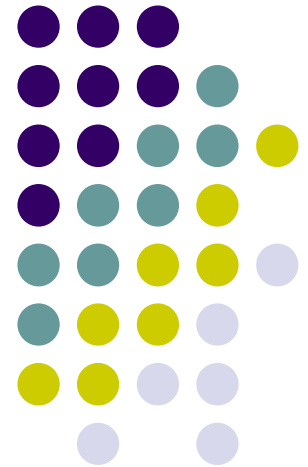
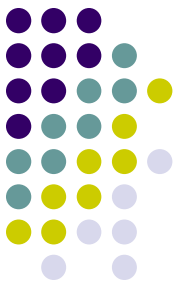


# Applications of suffix trees

## Lecture 4





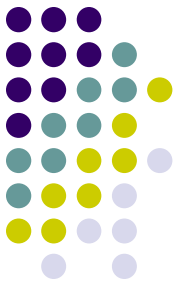
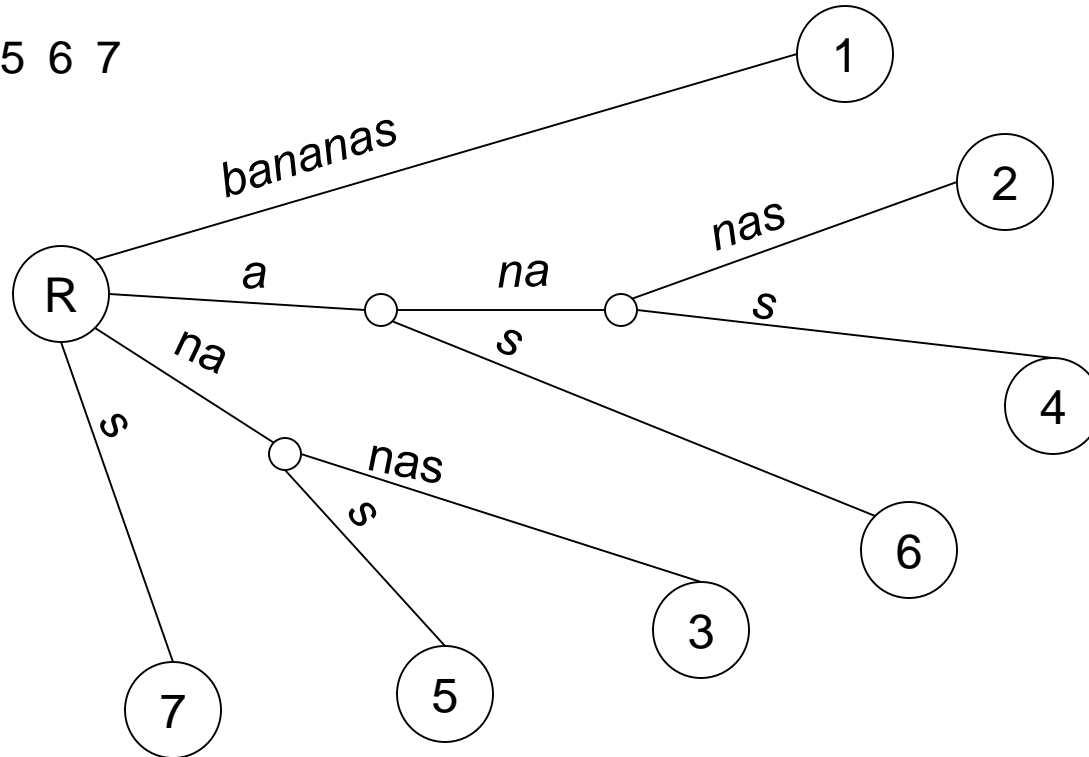
# Suffix tree - summary

- Suffix tree is a digital tree of all suffixes of text  $T$  (of length  $N$ )
- The suffixes are inserted by following a path of characters from the root, and a new branch of a tree is created only if the next character in a current suffix does not match the existing path
- Suffix tree has  $N$  leaves (1 for each suffix), where we store the starting position of a suffix in  $T$
- In order for each suffix to have its own leaf, we need to have at the end of  $T$  a special character which can not be found anywhere else in  $T$  – this ensures that a special branch will be created for each suffix which is also a prefix of another suffix

# Example: *bananas*

*b a n a n a s*

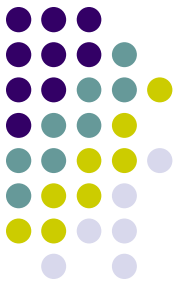
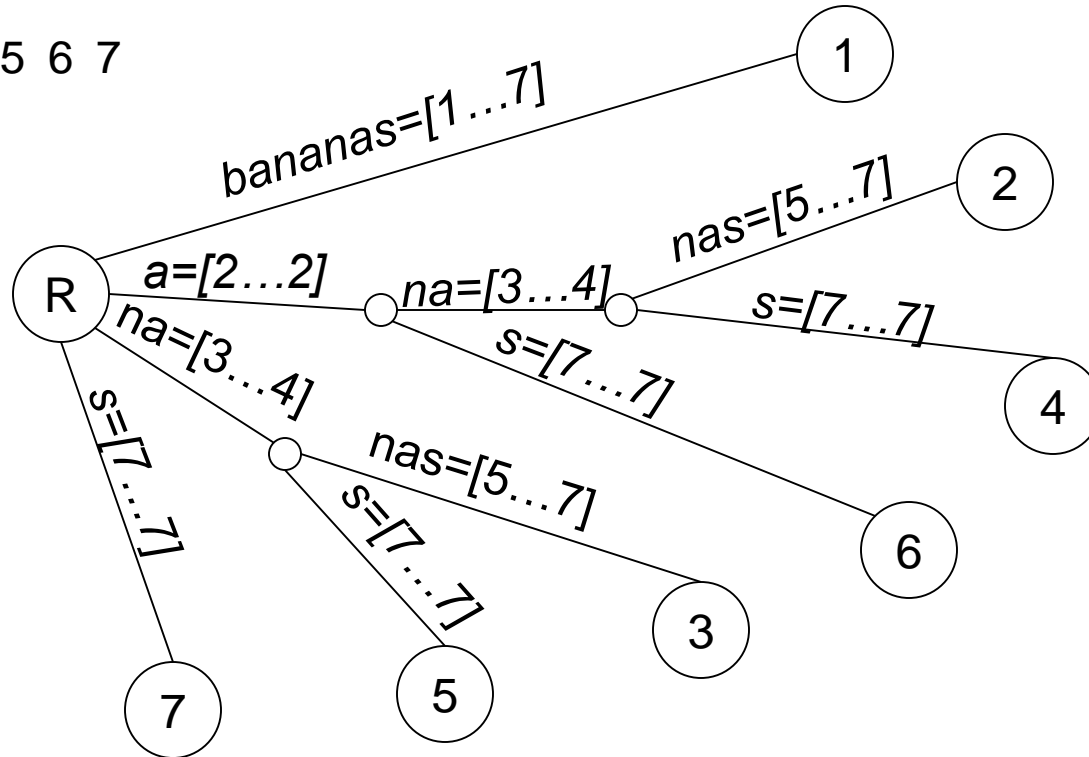
1 2 3 4 5 6 7

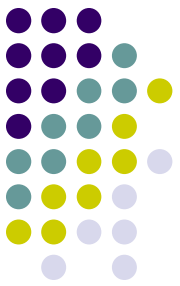


# Space reduction

*b a n a n a s*

1 2 3 4 5 6 7





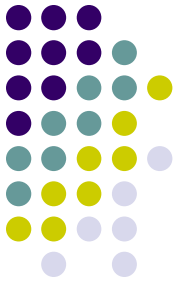
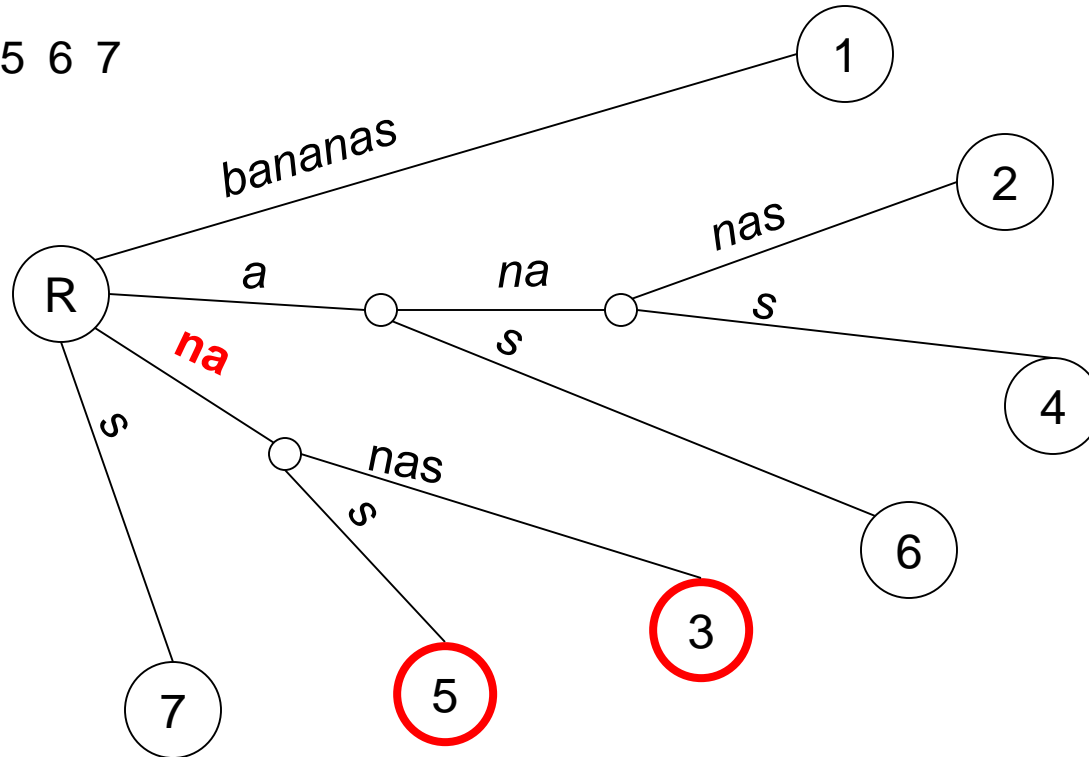
# Suffix tree - search

- In order to find all the occurrences of pattern  $P$  (of length  $M$ ), follow the path of symbols from the root. If there is a path corresponding to all  $M$  symbols of  $P$ , the positions where  $P$  occurs in  $T$  can be collected by the depth-first traversal of the subtree below the end of this path.
- If there is no path in  $T$  for all the characters of  $P$ , then pattern  $P$  does not occur in  $T$

# Example: $P=na$

*b a n a n a s*

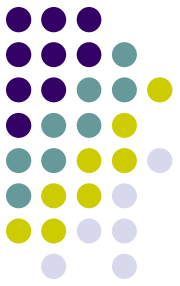
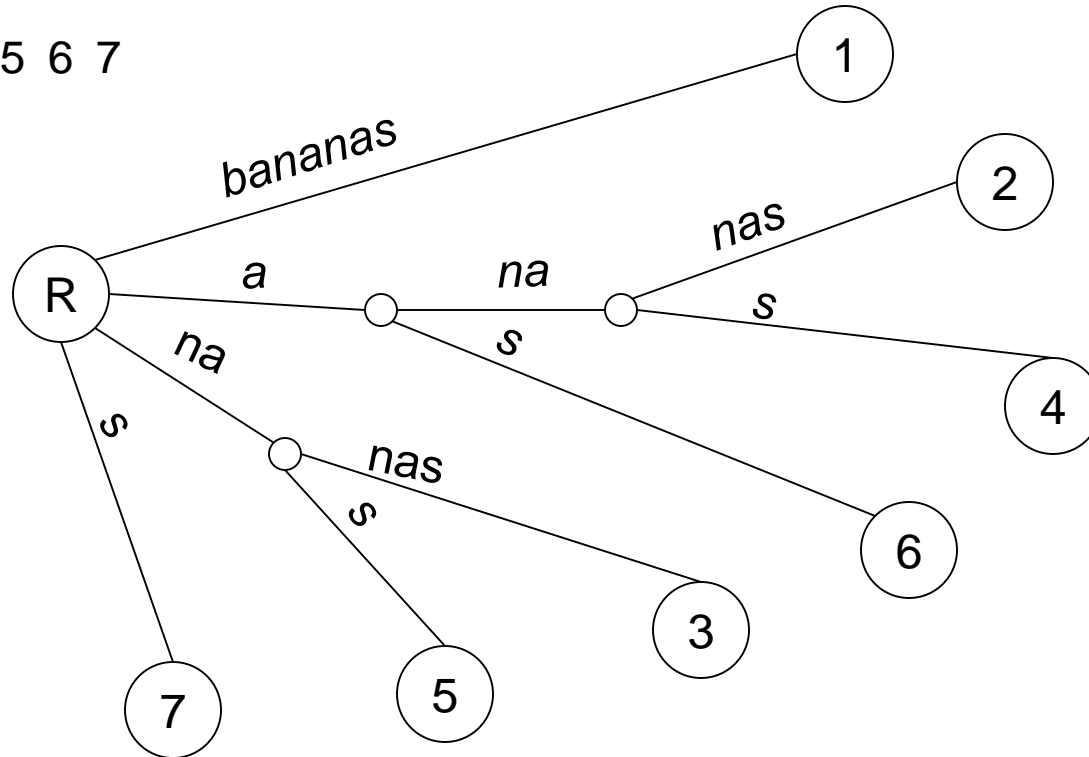
1 2 3 4 5 6 7



# Example: $P=an$

*b a n a n a s*

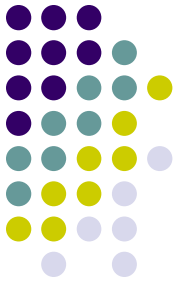
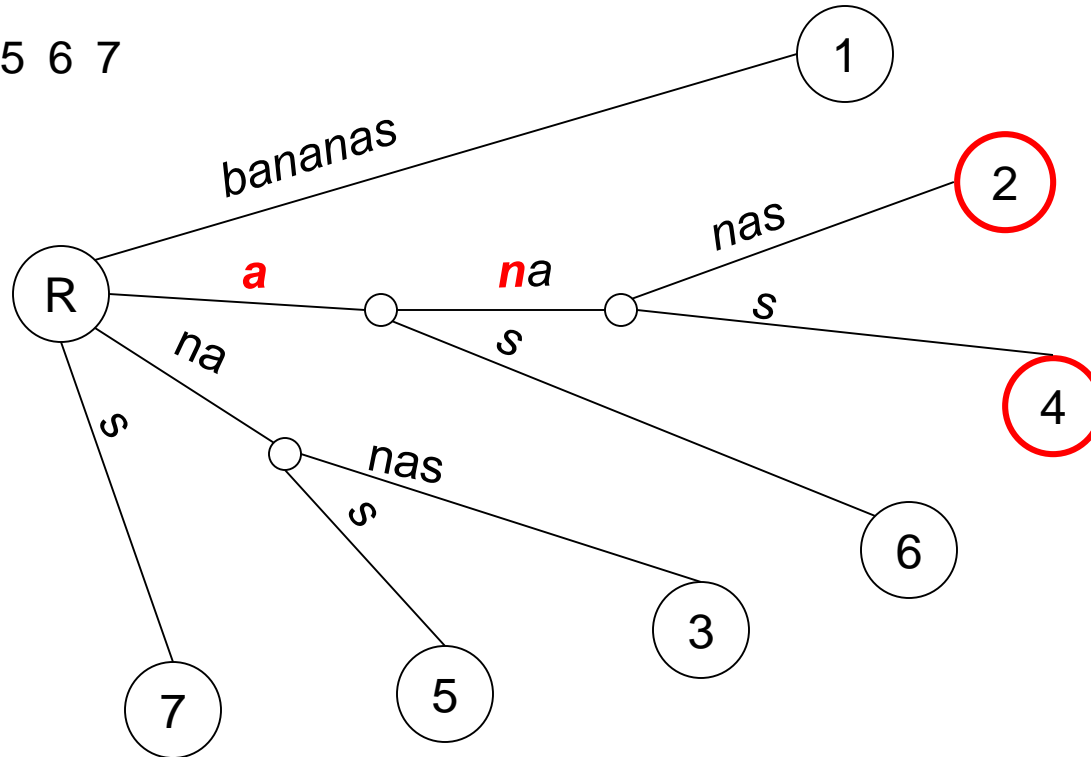
1 2 3 4 5 6 7



# Example: $P=an$

*b a n a n a s*

1 2 3 4 5 6 7

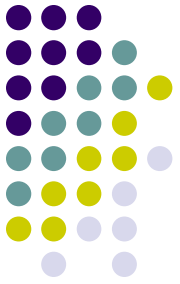
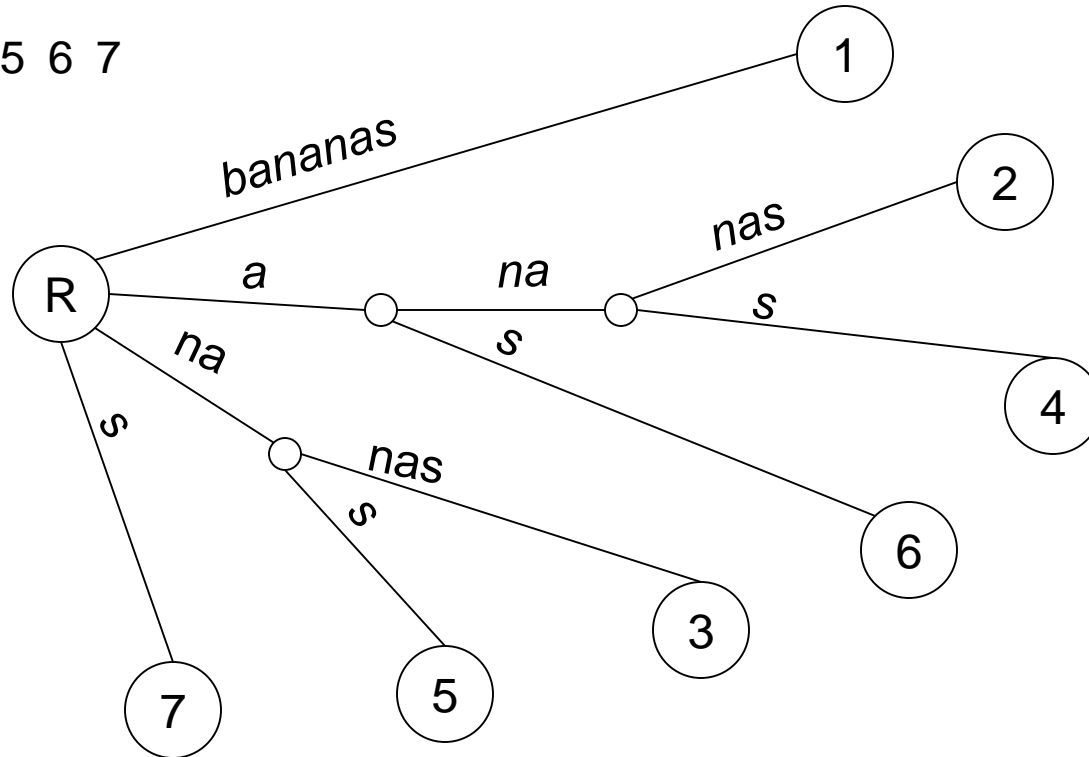


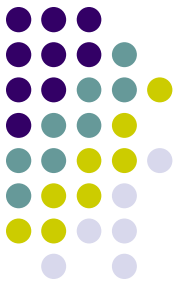


# Example: $P=naa$

*b a n a n a s*

1 2 3 4 5 6 7

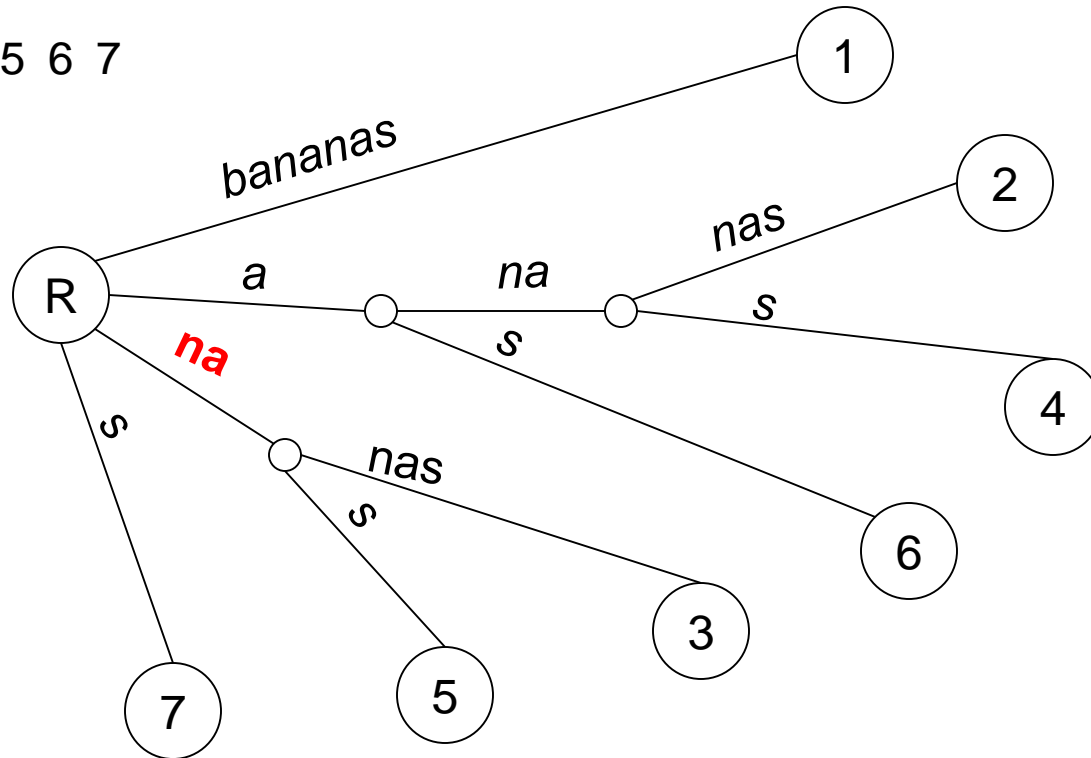




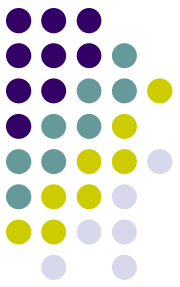
# Example: $P=naa$

*b a n a n a s*

1 2 3 4 5 6 7



Not found



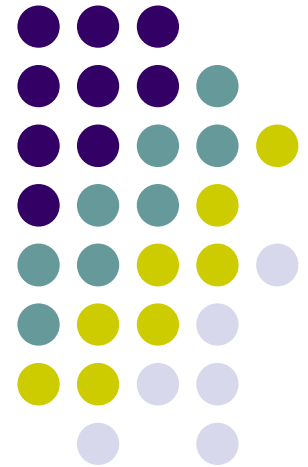
# Search efficiency

Input	Sub-tree size	Total number of sub-trees	Query time for 25 random patterns of length 10		Max occurrences
			Ave	Max	
8 GB (100 chromosomes) of eukaryotic genomes	131 MB	3125	1.3 sec	1.6 sec	5,598,876
	13 MB	31253	0.2 sec	0.3 sec	
3 GB (23 chromosomes of HG)	131 MB	1107	1.2 sec	1.4 sec	2,559,998
	13 MB	11063	0.2 sec	0.3 sec	
	1.3 MB	110635	<b>0.01 sec</b>	<b>0.03 sec</b>	
113 MB (6,300 viral genomes)	131 MB	75	1.2 sec	1.4 sec	15,534
	13 MB	754	0.2 sec	0.3 sec	

Grep (Boyer-Moore) –  
44 sec

# 1. Finding repeats

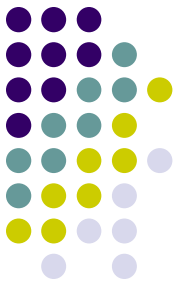
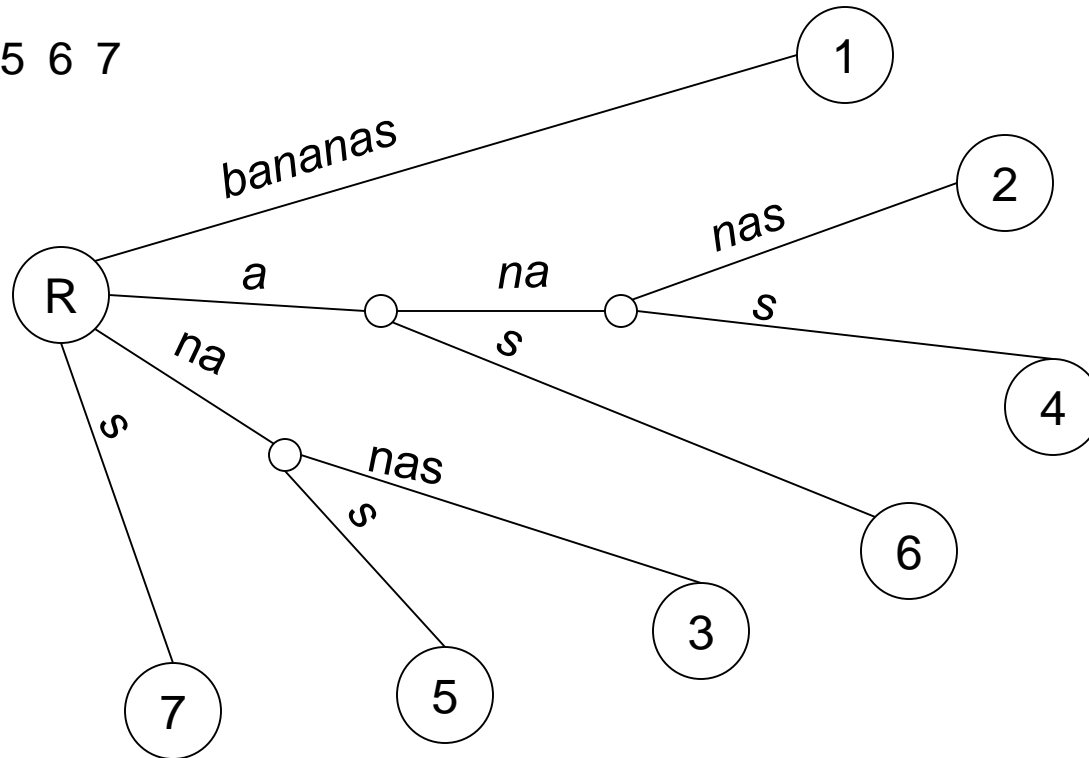
---



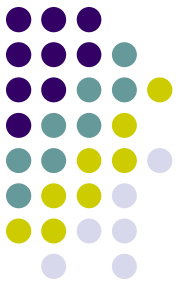
# Example: *bananas*

*b a n a n a s*

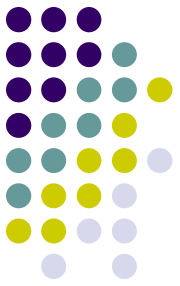
1 2 3 4 5 6 7



# Finding repeating substrings



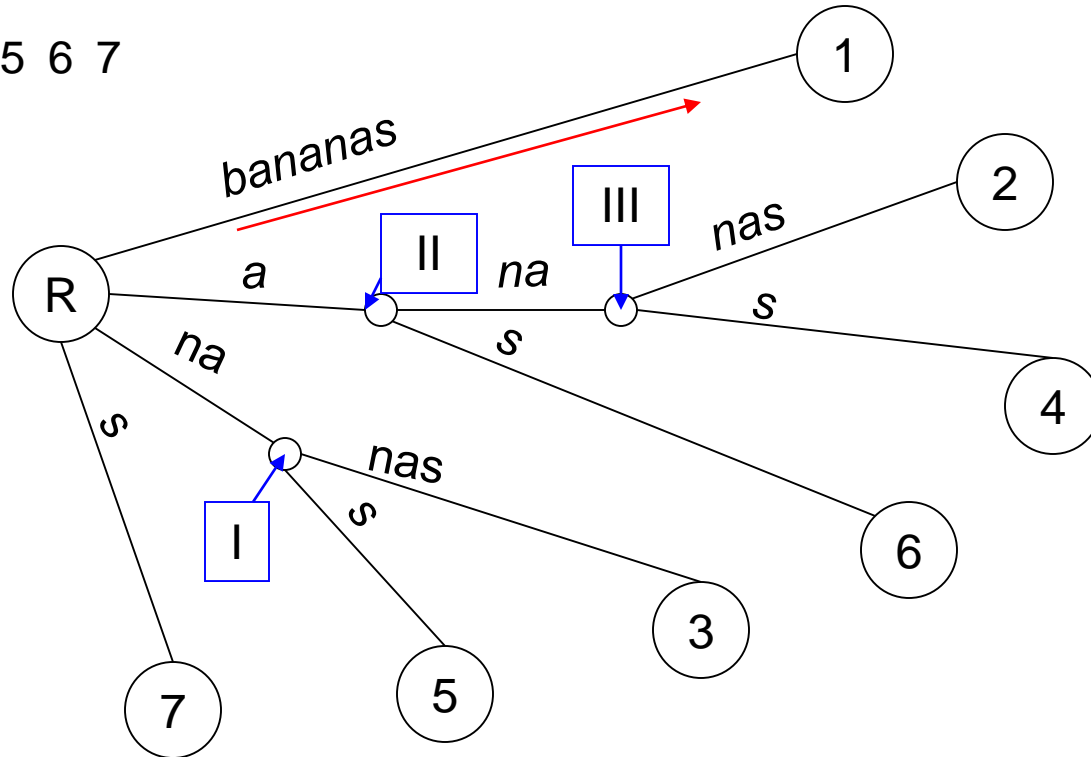
- The path from the root to any internal node of the suffix tree represents a substring of  $T$  which occurs at least twice in  $T$ , since it corresponds to a common prefix of at least 2 different suffixes
- Thus, all repeating substrings can be found by collecting the internal nodes of the suffix tree during the depth-first traversal



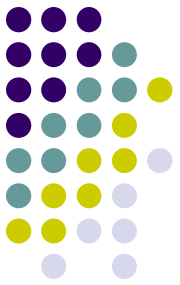
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



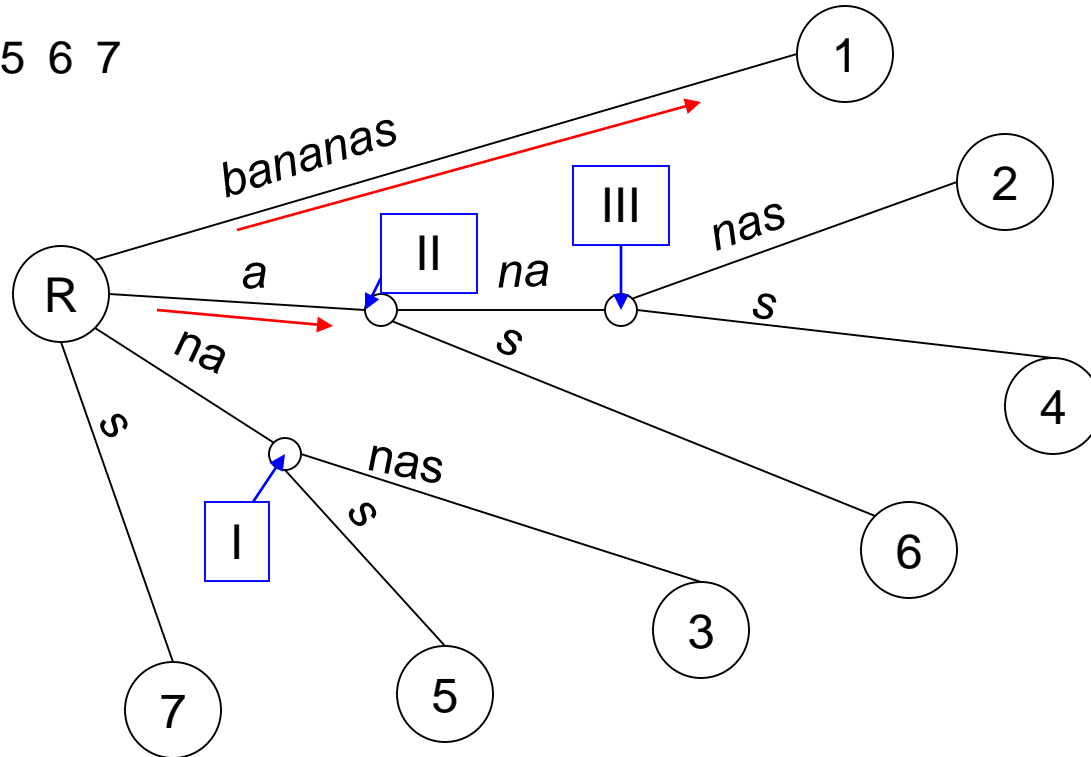
Sequence: R 1 R



# The depth-first traversal

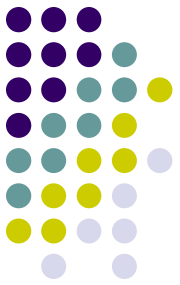
*b a n a n a s*

1 2 3 4 5 6 7



Sequence: R 1 R II

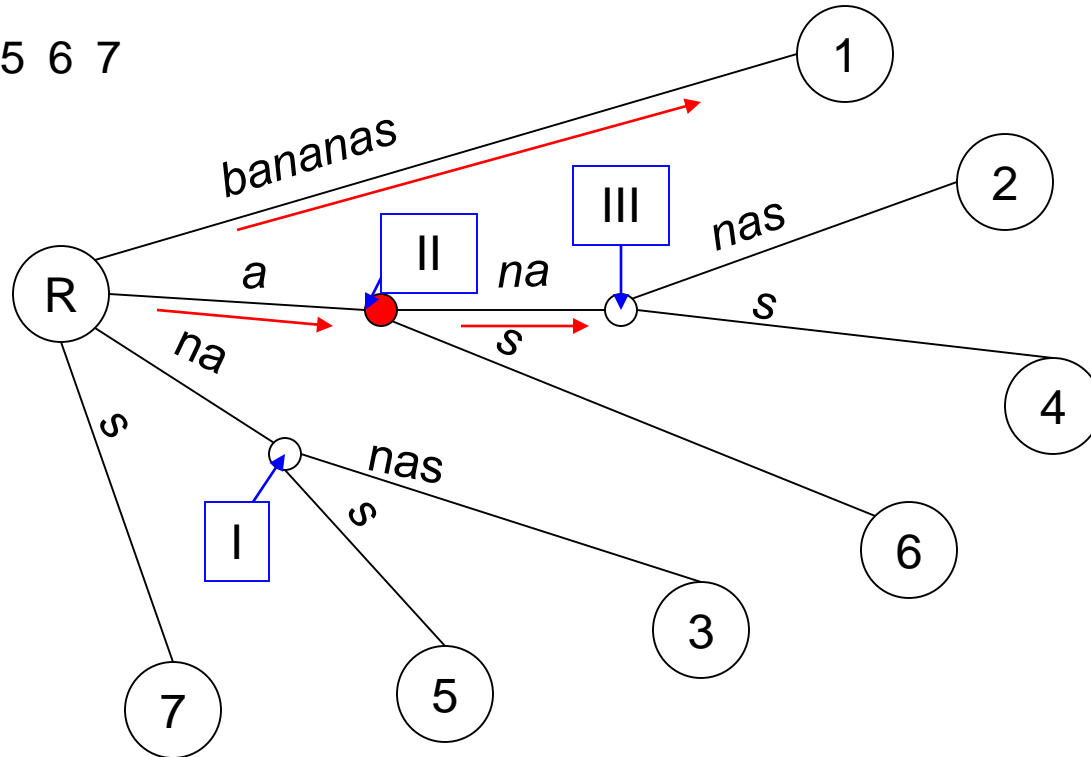




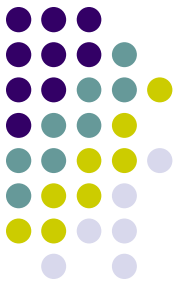
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



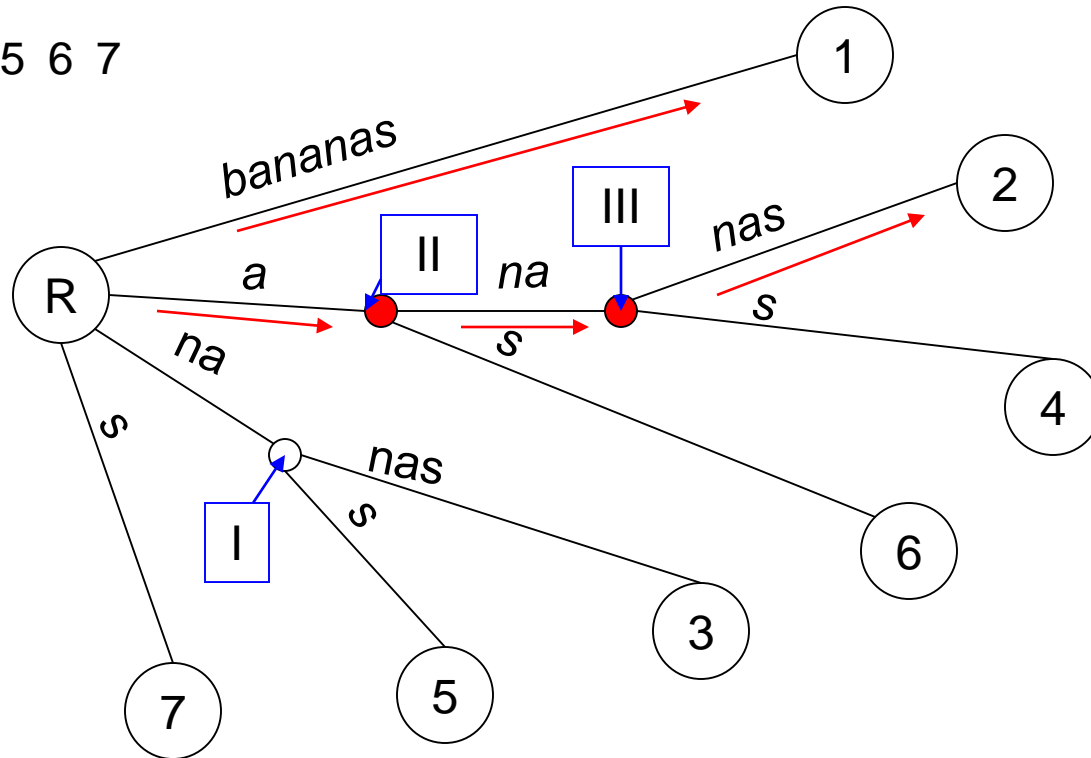
Sequence: R 1 R II III



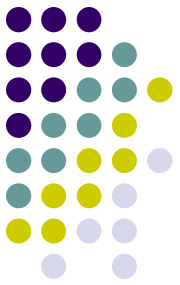
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



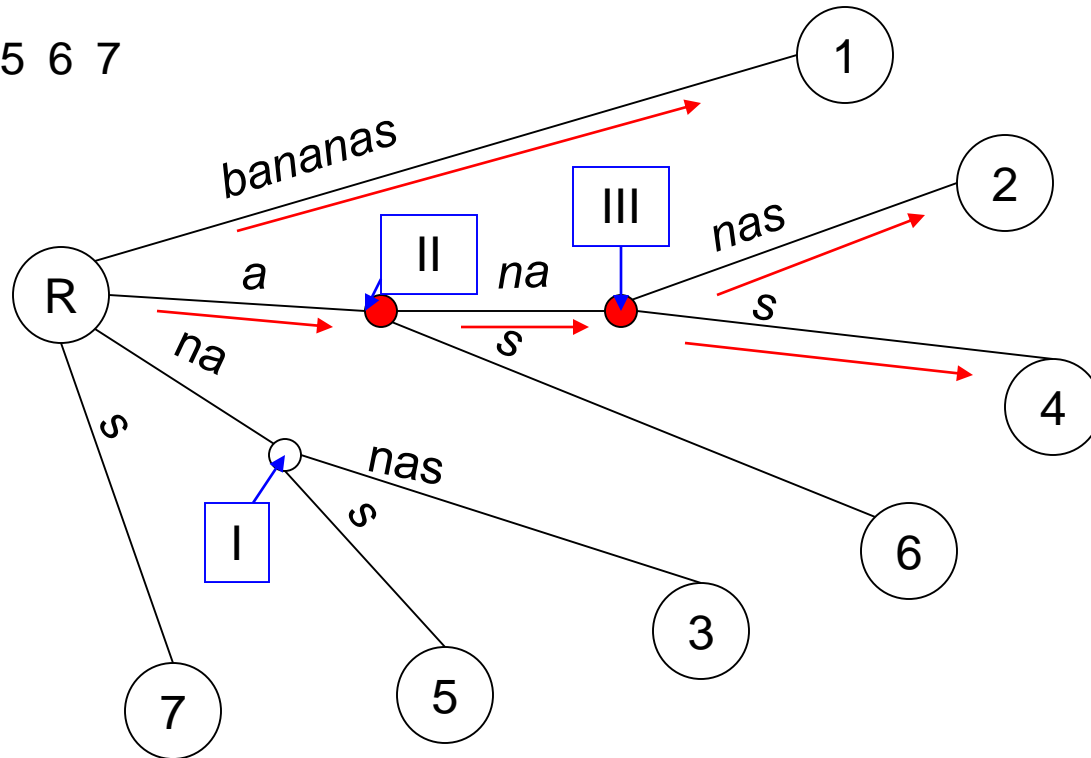
Sequence: R 1 R II III 2 III



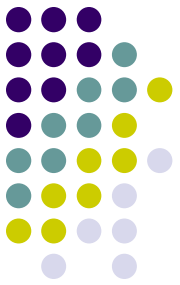
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



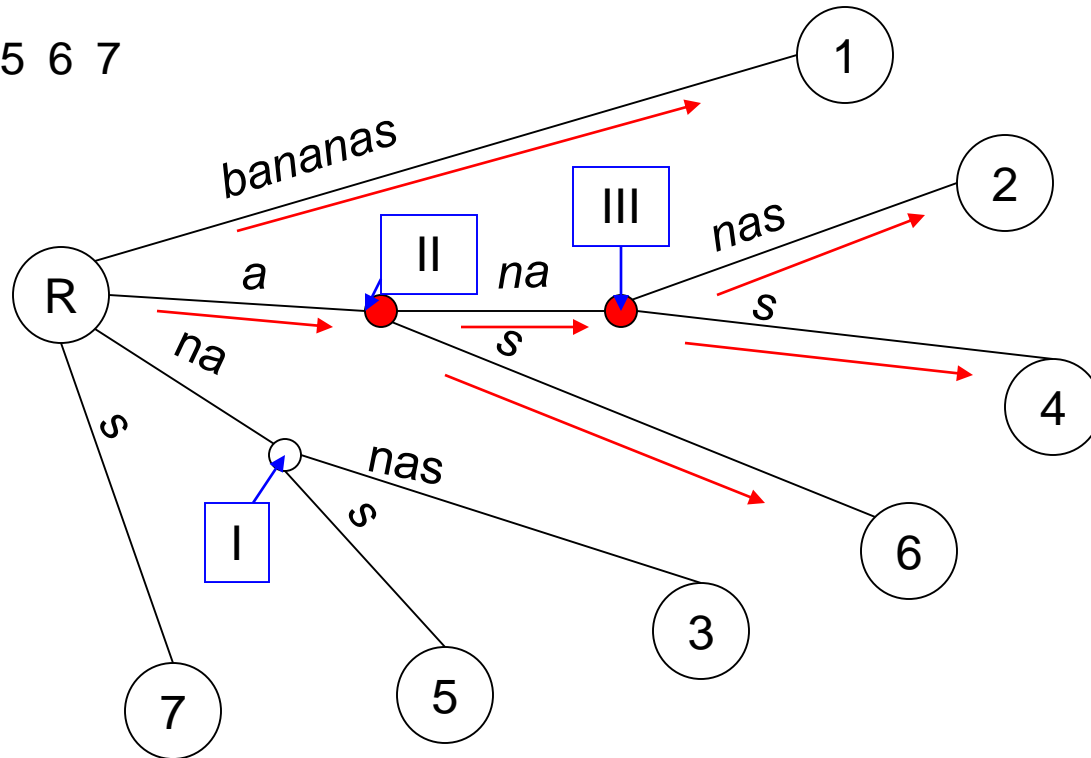
Sequence: R 1 R II III 2 III 4 III



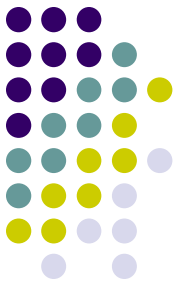
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



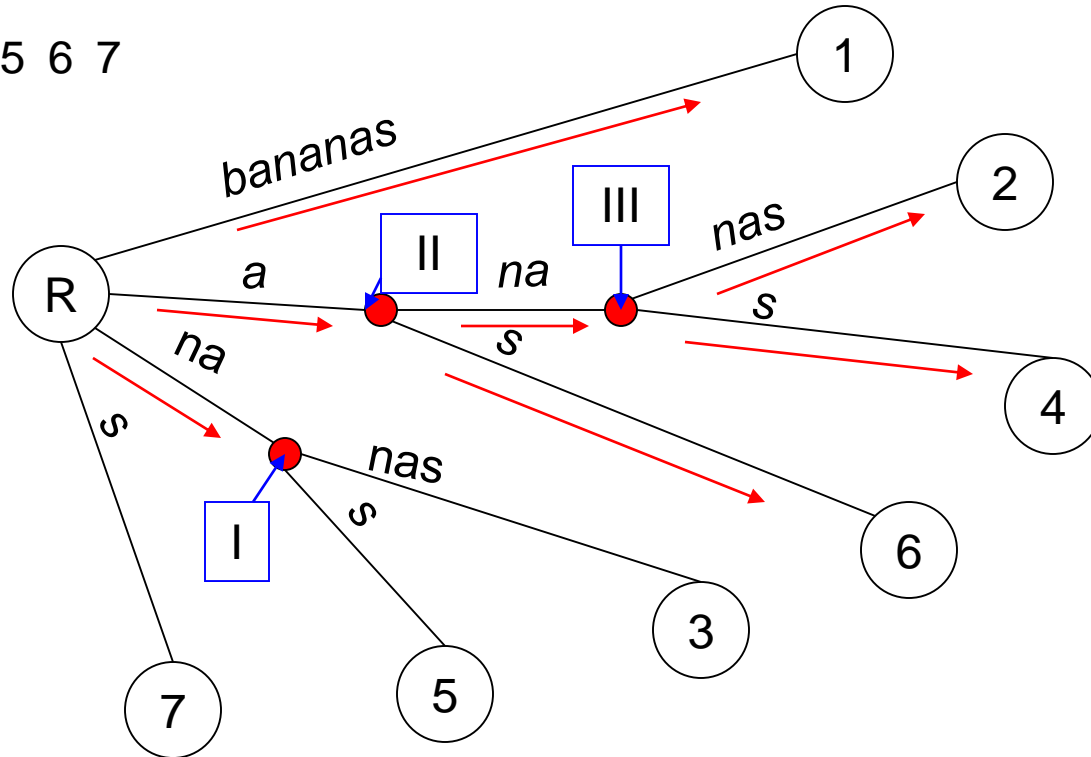
Sequence: R 1 R II III 2 III 4 III II 8 II R



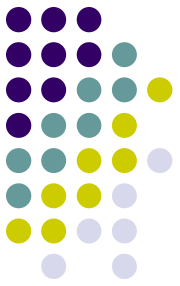
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



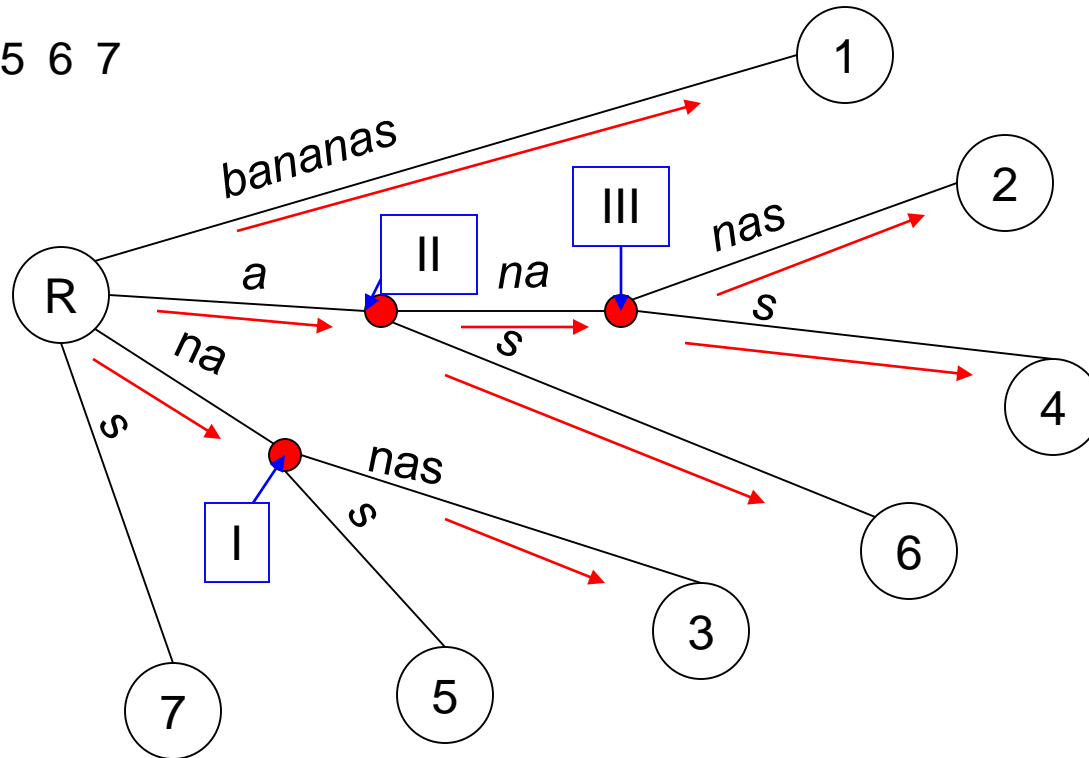
Sequence: R 1 R II III 2 III 4 III II 8 II R I



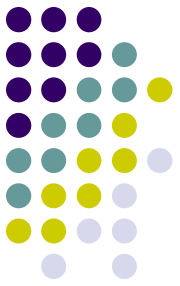
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



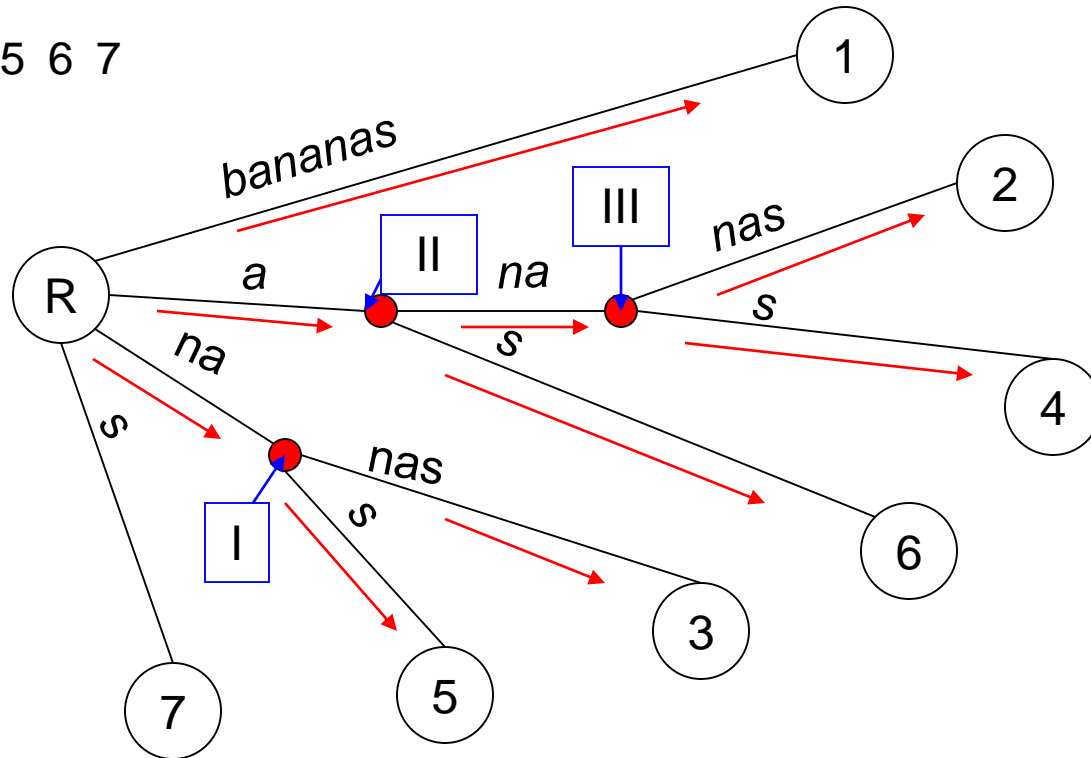
Sequence: R 1 R II III 2 III 4 III II 8 II R I 3 I



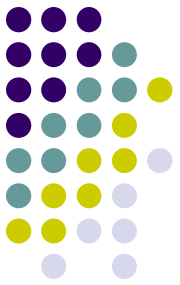
# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



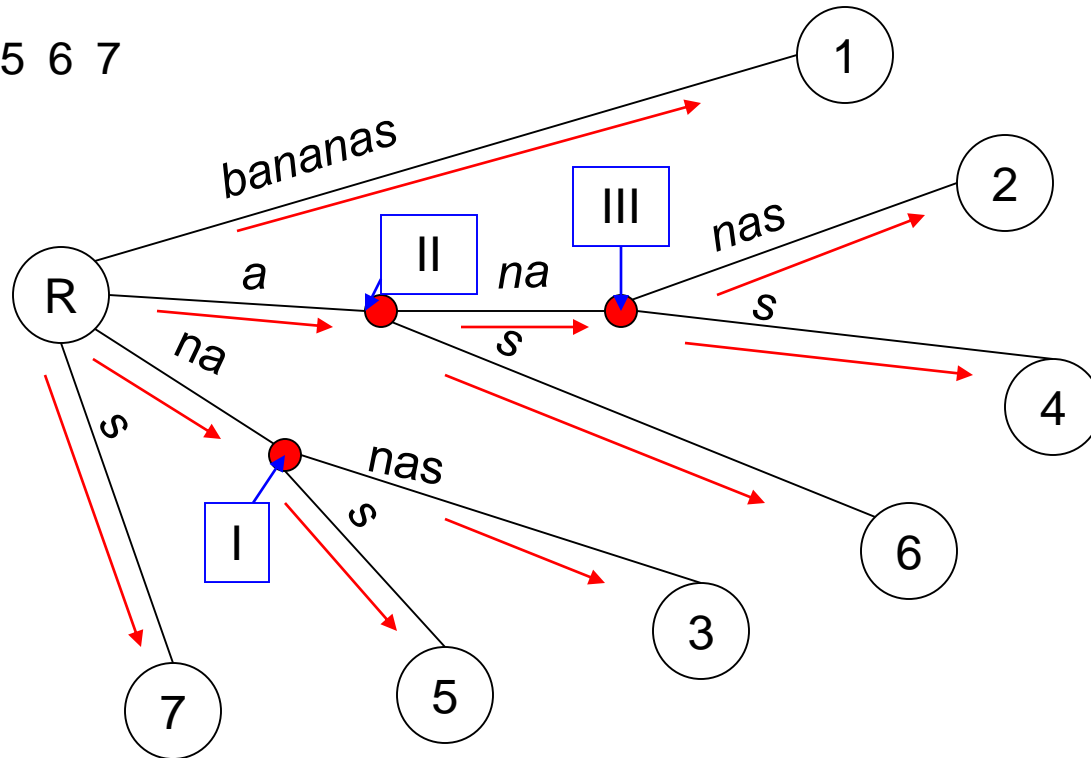
Sequence: R 1 R II III 2 III 4 III II 8 II R I 3 I 5 I R



# The depth-first traversal

*b a n a n a s*

1 2 3 4 5 6 7



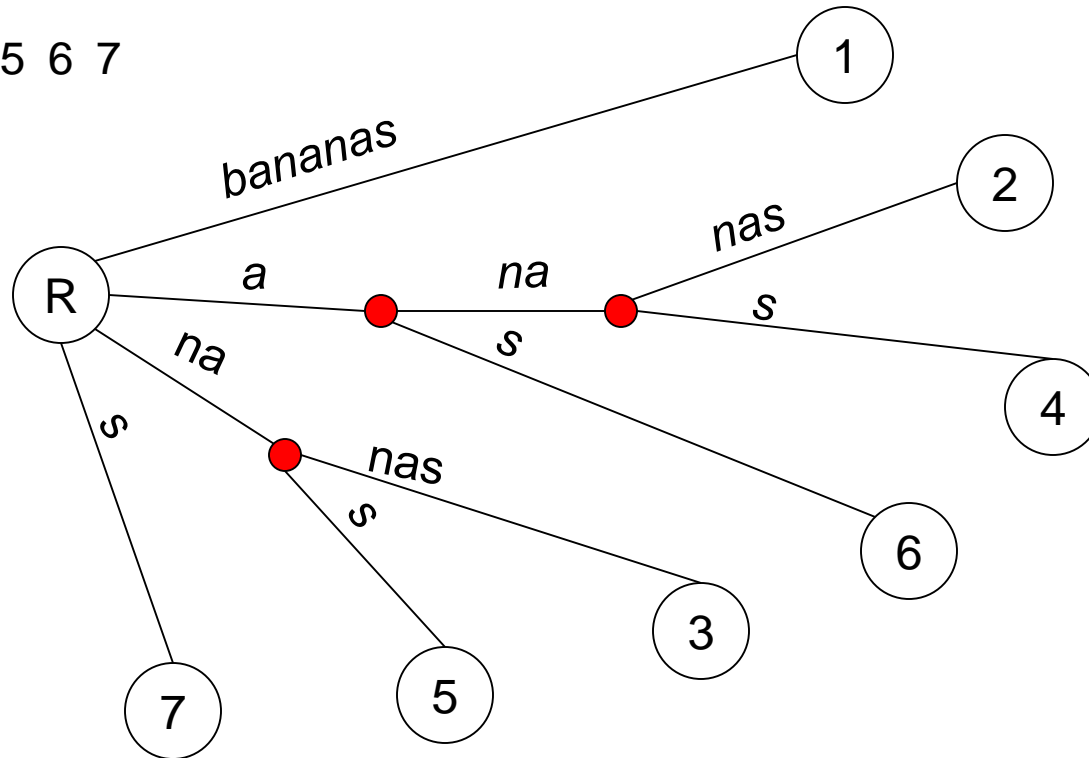
Sequence: R 1 R II III 2 III 4 III II 8 II R I 3 I 5 I R 7 R



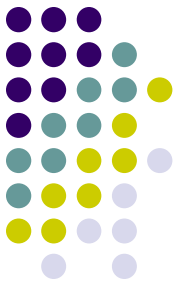
# All repetitions

*b a n a n a s*

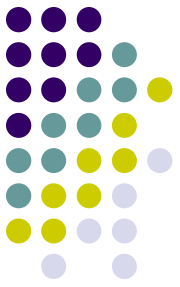
1 2 3 4 5 6 7



*n, na; a, an, ana;*

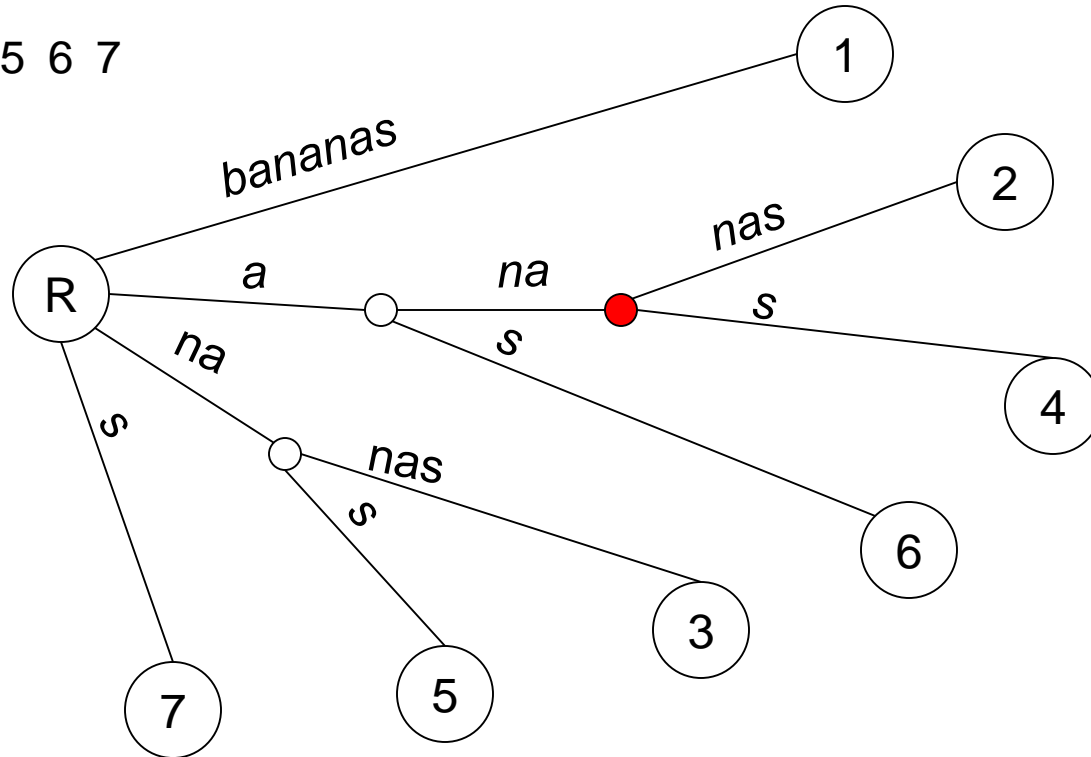


# The longest repeating substring in linear time

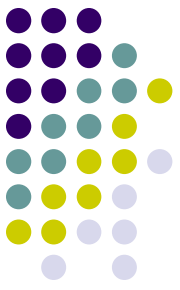


*b a n a n a s*

1 2 3 4 5 6 7

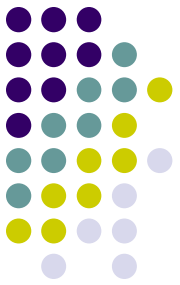


*ana*



# All maximal repeats

- Definition: A *maximal repeated pair* (MR) in a string  $T$  is a pair of identical substrings  $t_1$  and  $t_2$  such that the character to the immediate right (left) of  $t_1$  is different from the character to the immediate right (left) of  $t_2$ . Each MR pair can be represented by a tuple  $(i, j, k)$ , where  $i$  and  $j$  are start positions of the corresponding substrings, and  $k$  is the substring length
- If the characters to the right of  $t_1$  and  $t_2$  are different, we will call such repeat *right-maximal* (cannot be extended to the right).
- If the characters to the left of  $t_1$  and  $t_2$  are different, we will call such repeat *left-maximal* (cannot be extended to the left).

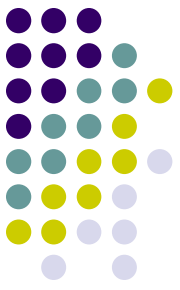


# Maximal repeats example

2            10            14  
↓            ↓            ↓

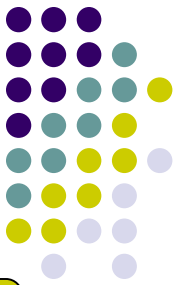
- $T = x$ *abc**yiiiz**abcq**abc**yrxar*
- Which of the following repeated pairs of length 3 are maximal repeats?
  - $(2, 10)$
  - $(2, 14)$
  - $(10, 14)$

# An efficient algorithm for finding all maximal repeats

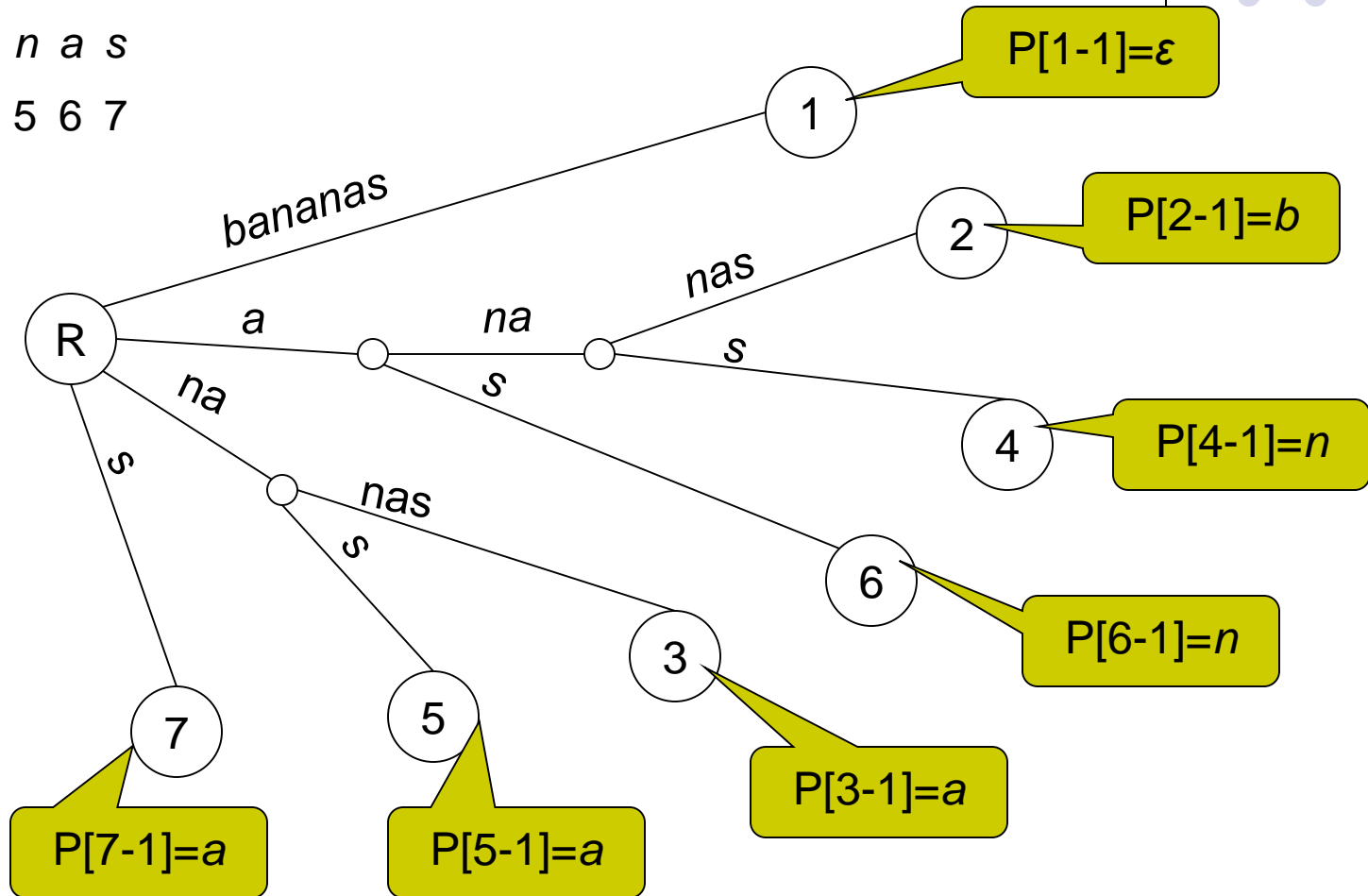


- The substring labeling the path to any internal node of the suffix tree always represents a right-maximal pair (Why?)
- Each such substring represents a prefix of some pair of suffixes  $T[i...N]$  and  $T[j...N]$ . In order to check for each such substring if it is also a left-maximal repeat, we need only to check if the characters at positions  $T[i-1]$  and  $T[j-1]$  are different.
- This can be done in a linear time.

# Step 1. Mark leaves with the left character



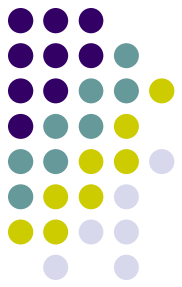
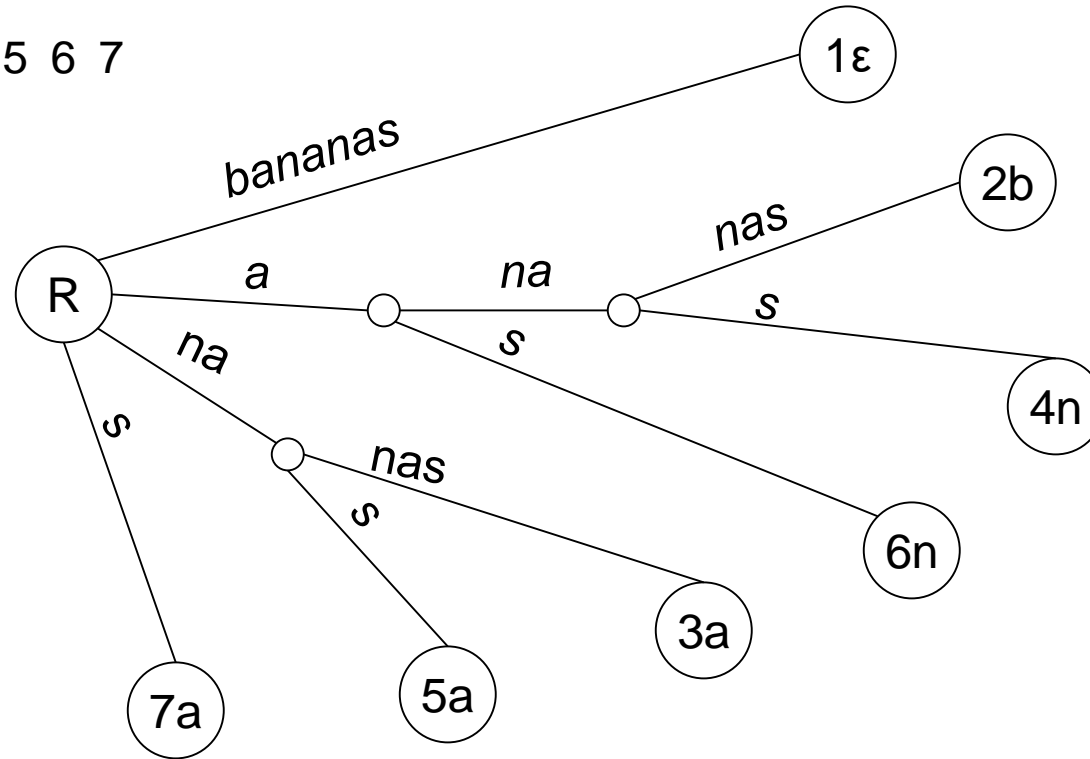
*b a n a n a s*  
1 2 3 4 5 6 7



# Step 2. Traverse

*b a n a n a s*

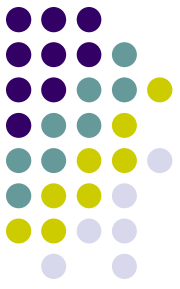
1 2 3 4 5 6 7



If both children of an internal node have the same character to the left of the suffix, then mark this internal node with this character.

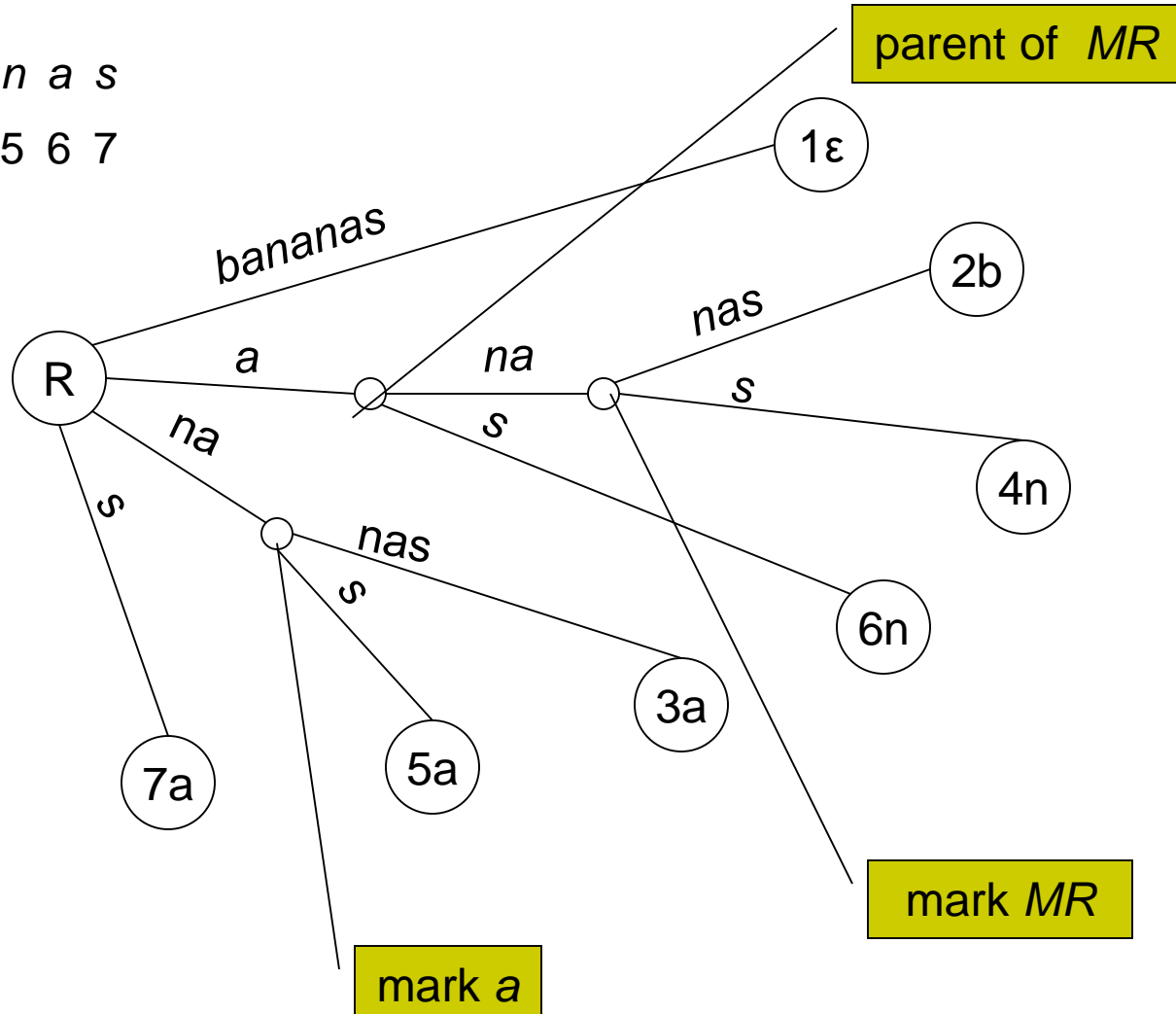
If the left characters are different, then the path from the root to this node represents a maximal repeat, so mark the node as maximal repeat

All parents of MR node are maximal repeats too



# Step 2. Mark internal nodes

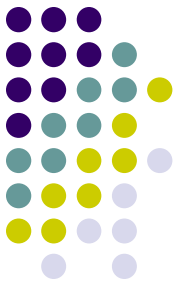
*b a n a n a s*  
 1 2 3 4 5 6 7



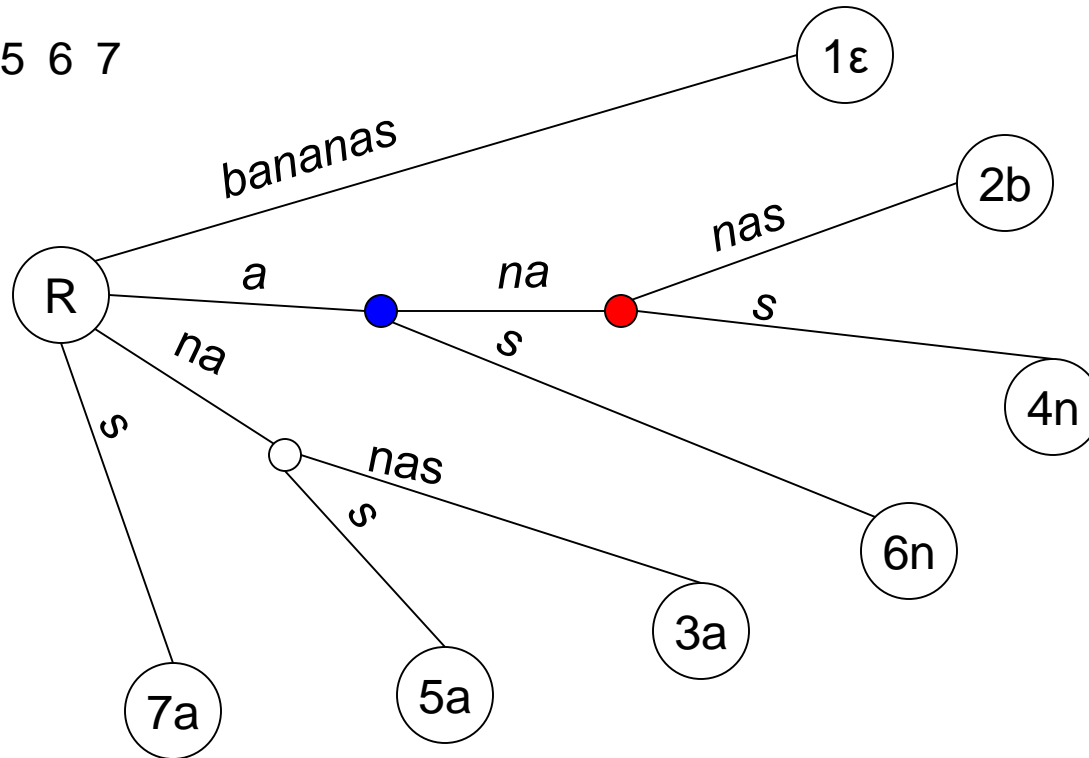
The parent of maximal repeat is a maximal repeat too (why?)



# Step 3. Output

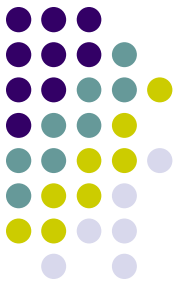


*b a n a n a s*  
1 2 3 4 5 6 7



Maximal repeat is *ana* (2,4)

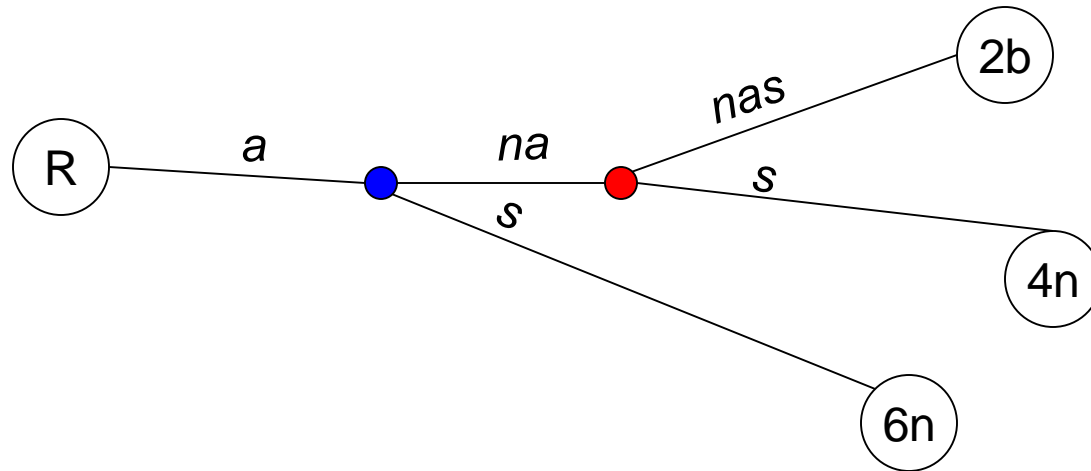
Maximal repeat is also *a* (2,6)



# Step 3. Output

*b a n a n a s*

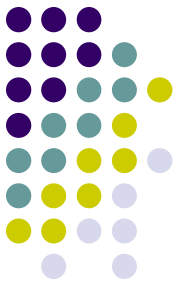
1 2 3 4 5 6 7



There can be up to  $N^2$  maximal repeats in any string (why?)

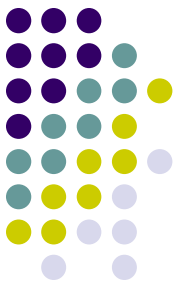
These maximal repeats can be efficiently represented in a linear space using the same suffix tree with the nodes corresponding to repeats only

# Repetitions observed in genome sequences



- Families of reiterated sequences account for about one third of the human genome
- For  $3.6 \times 10^6$  nucleotides of the C. Elegans genome, 7000 families of repetitive sequences were discovered
- Prokaryotes<sup>1)</sup> have in total little repetitive DNA
- The mechanism – unequal crossing-over

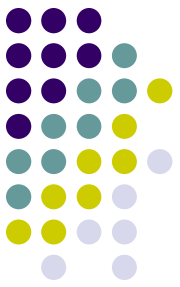
<sup>1)</sup> Prokaryotes (for example, Bacteria) have a circular DNA not enclosed into a nucleus



# Repetitions in genome sequences I

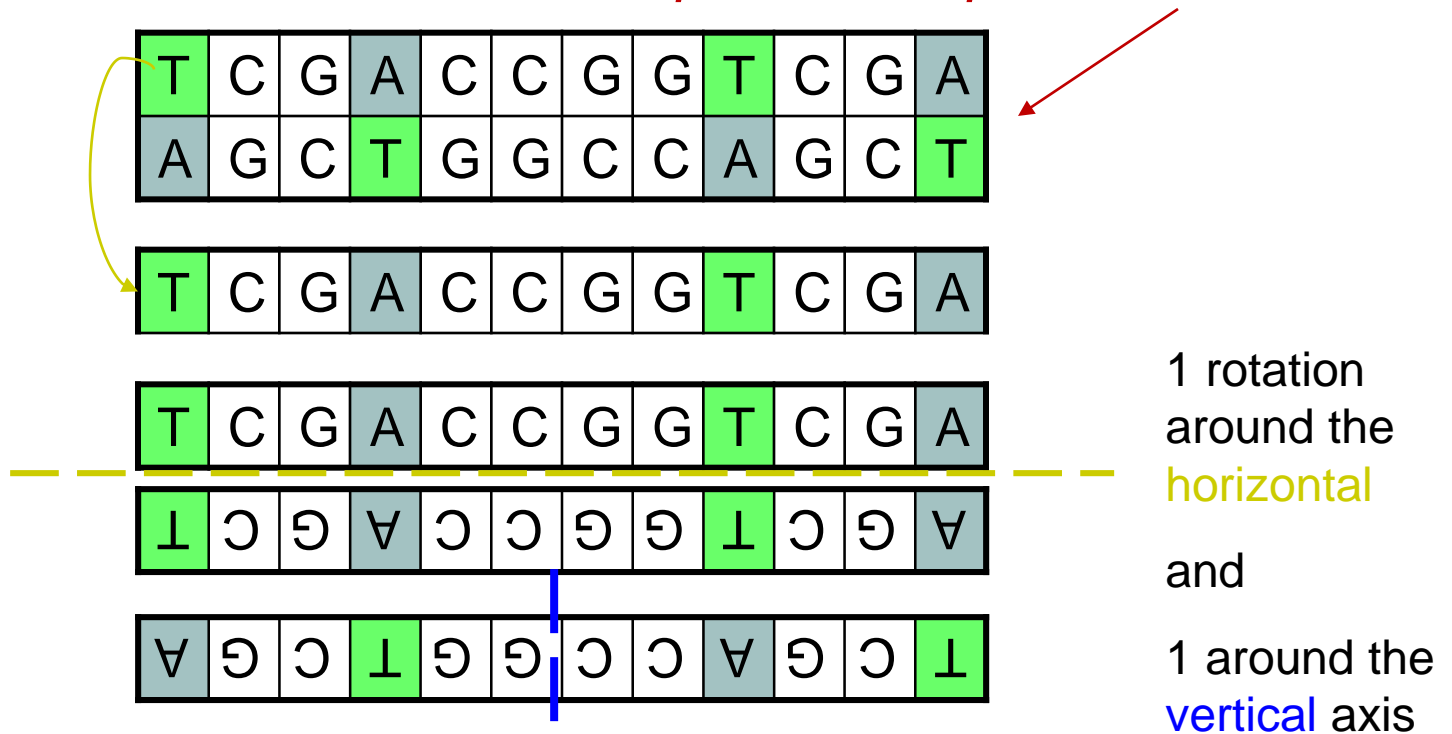
- Gene families
  - Many genes occur in multiple copies in the genome
  - They may be identical copies (r-RNA\* code) or just similar sequences of the same gene, modified by mutations
  - Some contain only a short similar motif – **homeobox** (~160 Bp)– which defines the shape of the protein-binding site
  - The copies may occur in **tandem** (one after another) or are dispersed through different areas of the genome
- Some of these multiple copies serve the purpose of an enhanced gene expression (r-RNA), others are redundant and are used interchangeably when one copy is damaged
- The copy of the original gene can mutate and acquire a new function
- This is believed to be a main mechanism of evolution

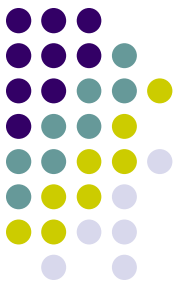
\*R-RNA is an RNA component of the ribosomes



# Repetitions in genome sequences II

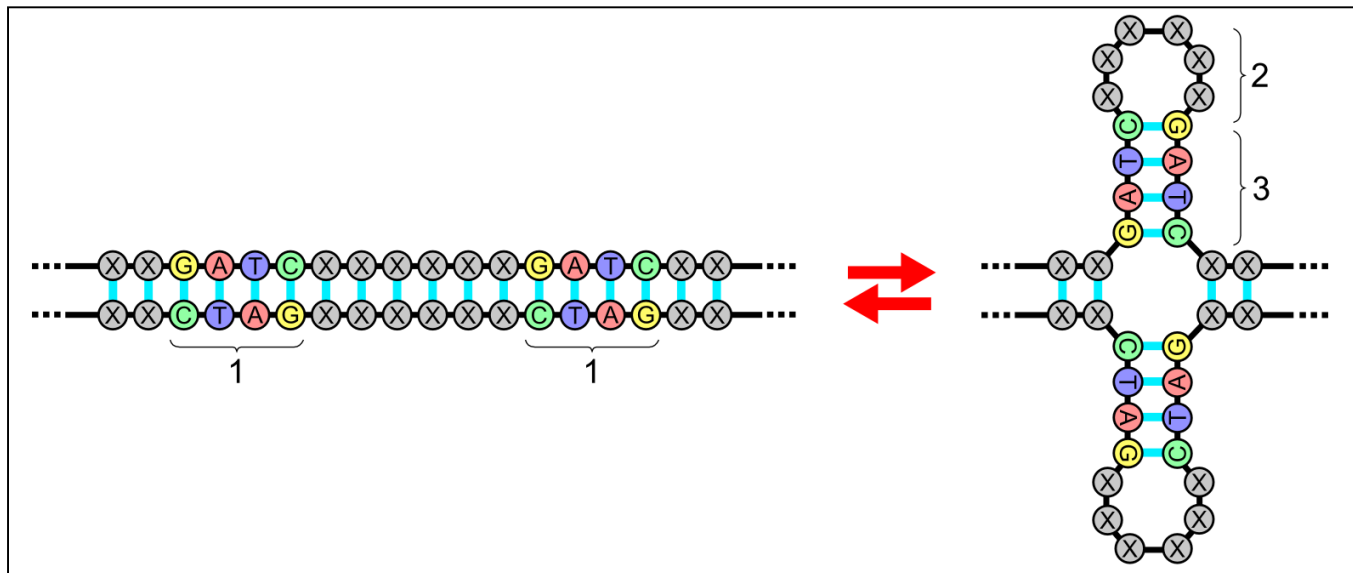
- Functional repeats – short repeats encoding the same functional sites (transcription sites and protein-binding sites on DNA)
- They often have a form of a *complemented palindrome*



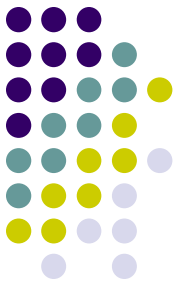


# Repetitions in genome sequences II

- These complemented palindromic repeats have a potential to form secondary structures such as hairpins and stem loops reflecting the dimeric nature of proteins. They serve for recognition of the transcription sites by enzymes

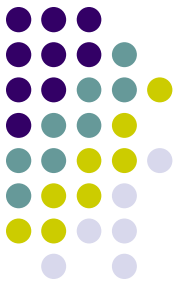


# Repetitions in genome sequences III



- *Transposons* – dispersed repetitive elements – the remains of viral DNA which were incorporated into genome and lost their functionality
  - SINE – Short Interspersed Nuclear Elements – *Alu* element (in human but not in mouse) 300 bp flanked by direct repeats –  $10^6$  copies, 1 such element per each 4 kBp sequence
  - LINE – Long Interspersed Nuclear Element – L1 element – 6 kbp long and  $10^5$  copies. They are not in the protein-coding regions, but often in introns, or at the ends of the transcribed region, so they are transcribed as a part of a gene

# Repetitions in genome sequences IV

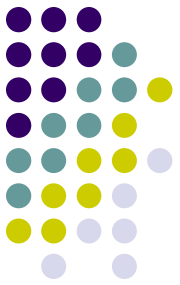


- *Satellite sequences*

Microsatellites – distributed through the entire genome  
1-4 bp repeats in clusters of ~200 Bp. They are highly polymorphic (in the number of copies) and make an ideal genetic marker. *VNTR*, variable number of tandem repeats, is used for personal identification

- When these repeats are inside a protein-coding region, they cause severe diseases (for example, Huntington's disease, if more than 20 *CAG* repeats are present inside the coding region for the *huntingtin* protein)



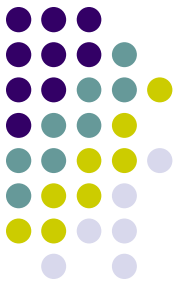


# Repetitions in genome sequences IV

- *Satellite sequences*

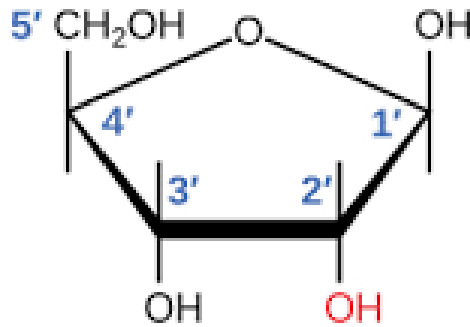
Minisatellites – occur as tandem repeats at the end of chromosomes

They have an important function discussed below

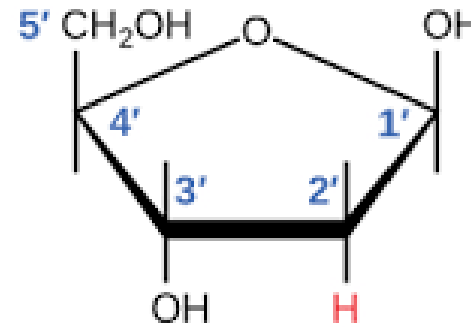


# Chemistry of nucleic acids

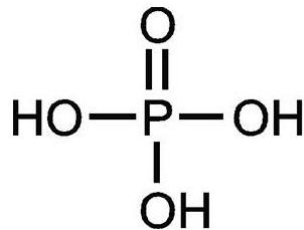
- DeoxyriboNucleic Acid (DNA)
- RiboNucleic Acid (RNA)



Ribose

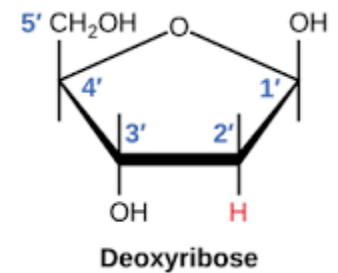
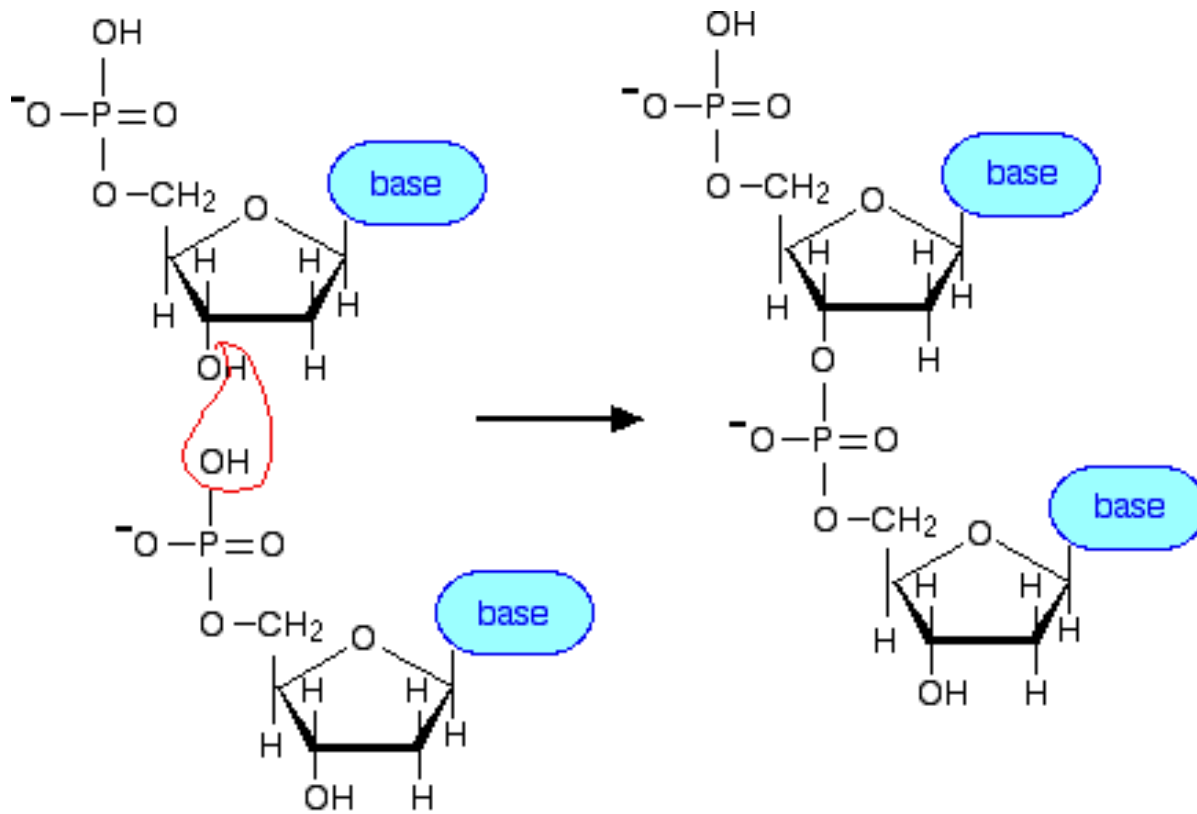
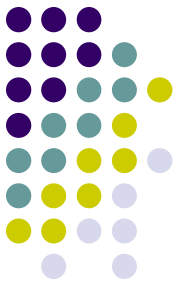


Deoxyribose

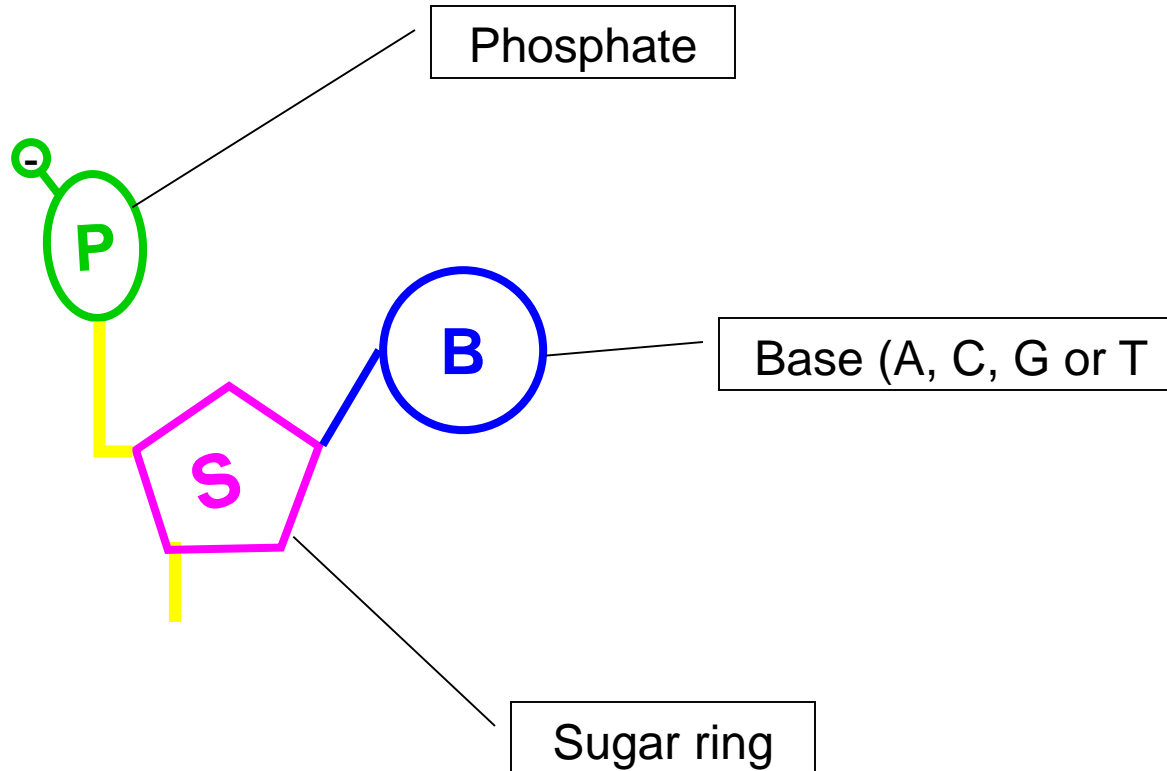
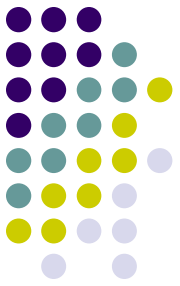


Phosphoric acid

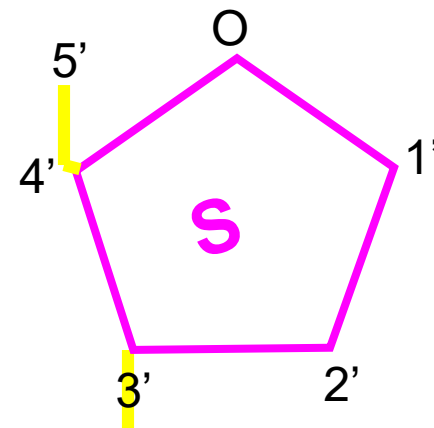
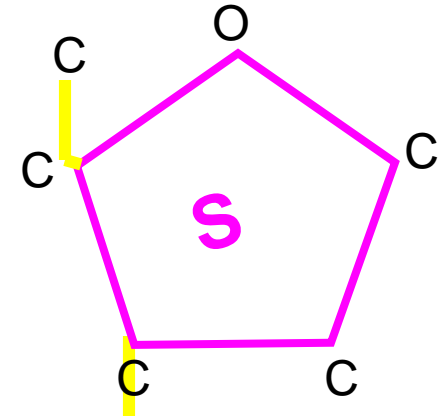
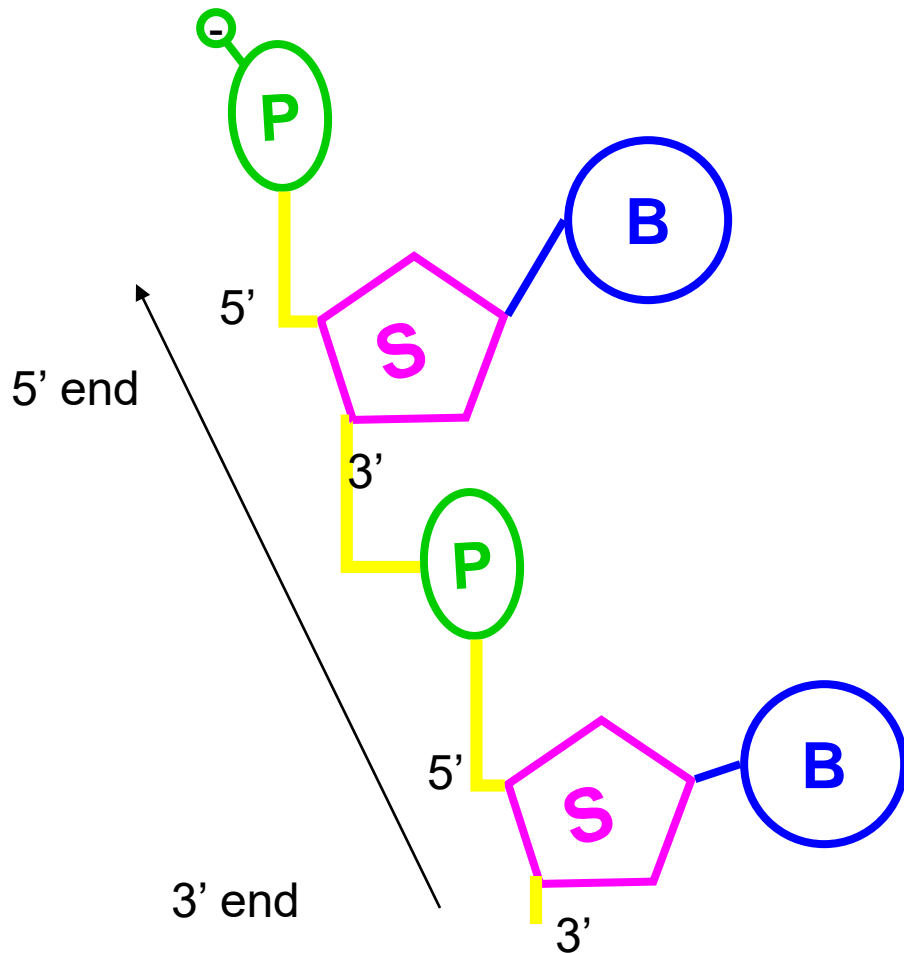
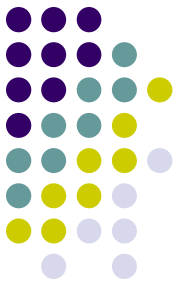
# Chemistry of DNA



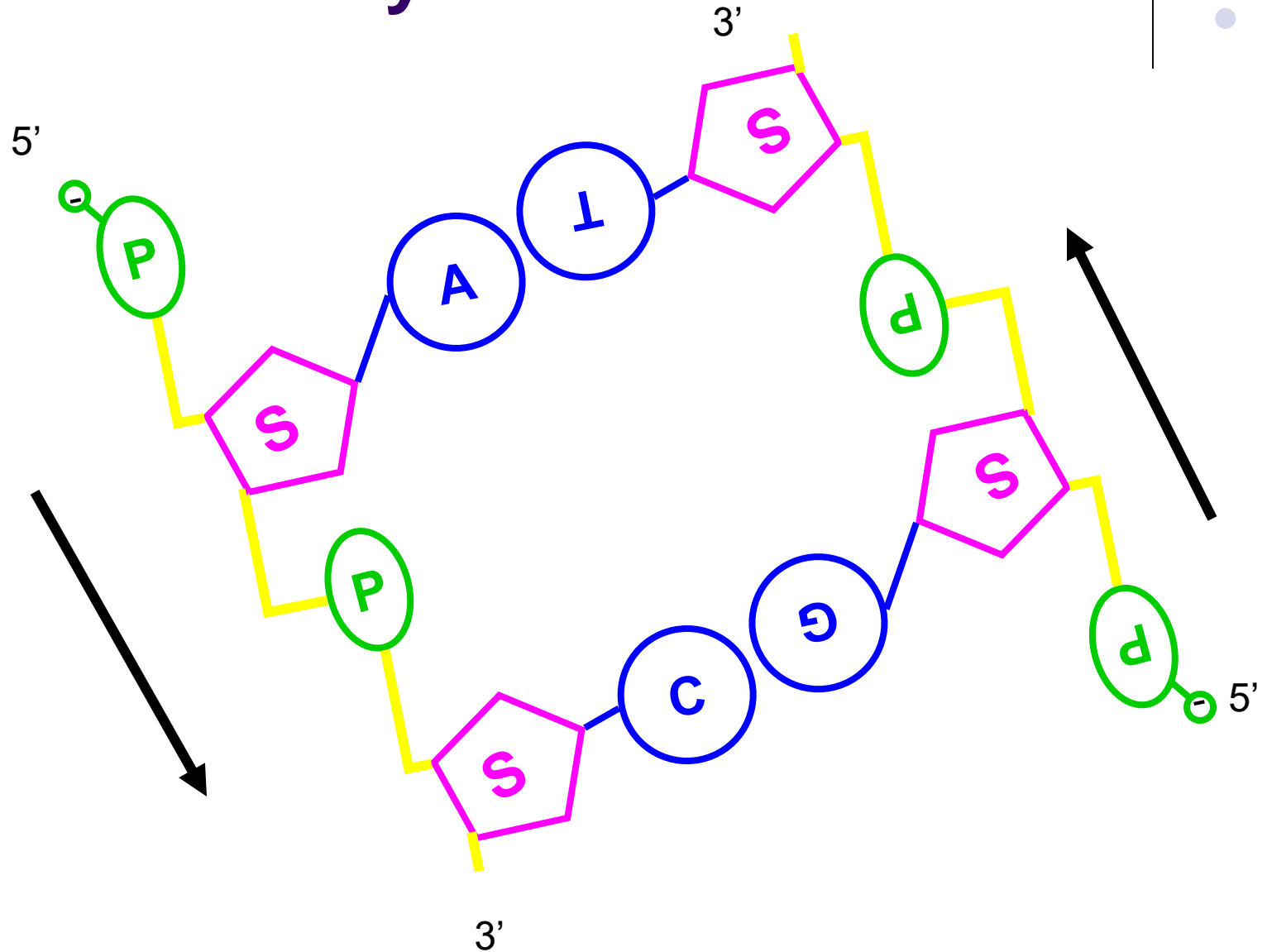
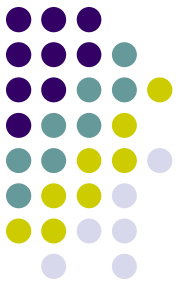
# Nucleotide – DNA polymer building block



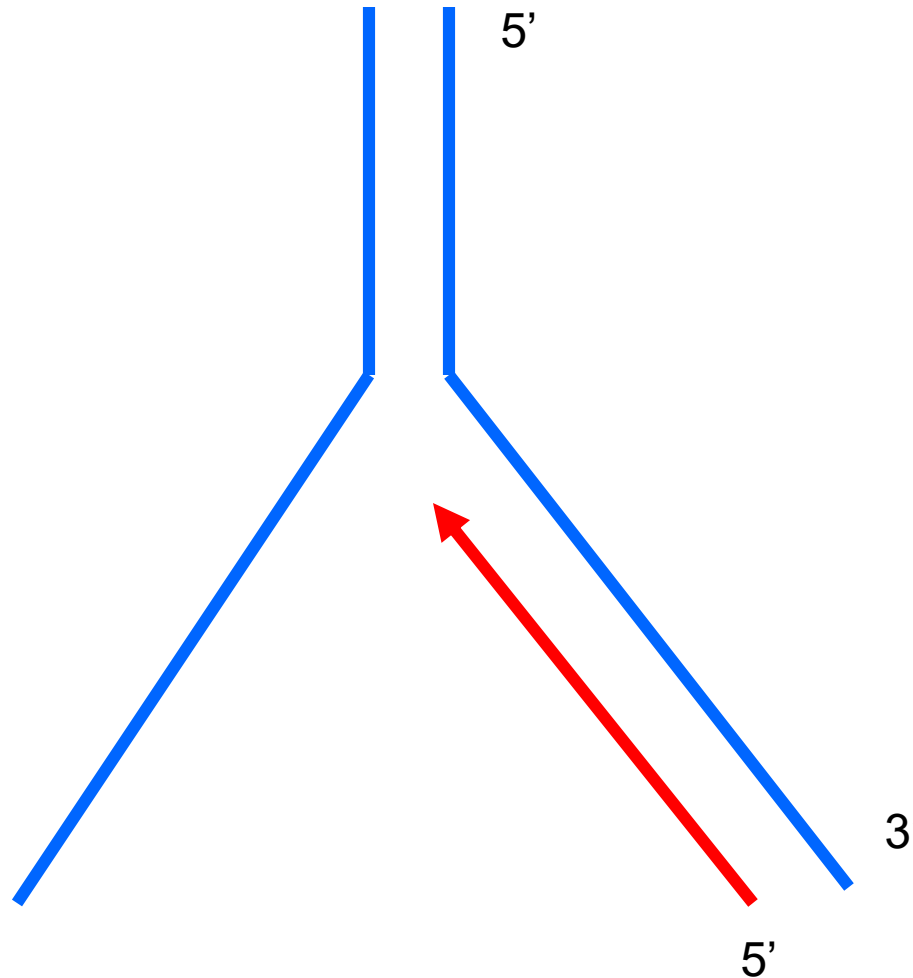
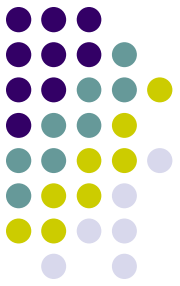
# Chaining nucleotides



# Synthesis of the chain can be performed only in one direction

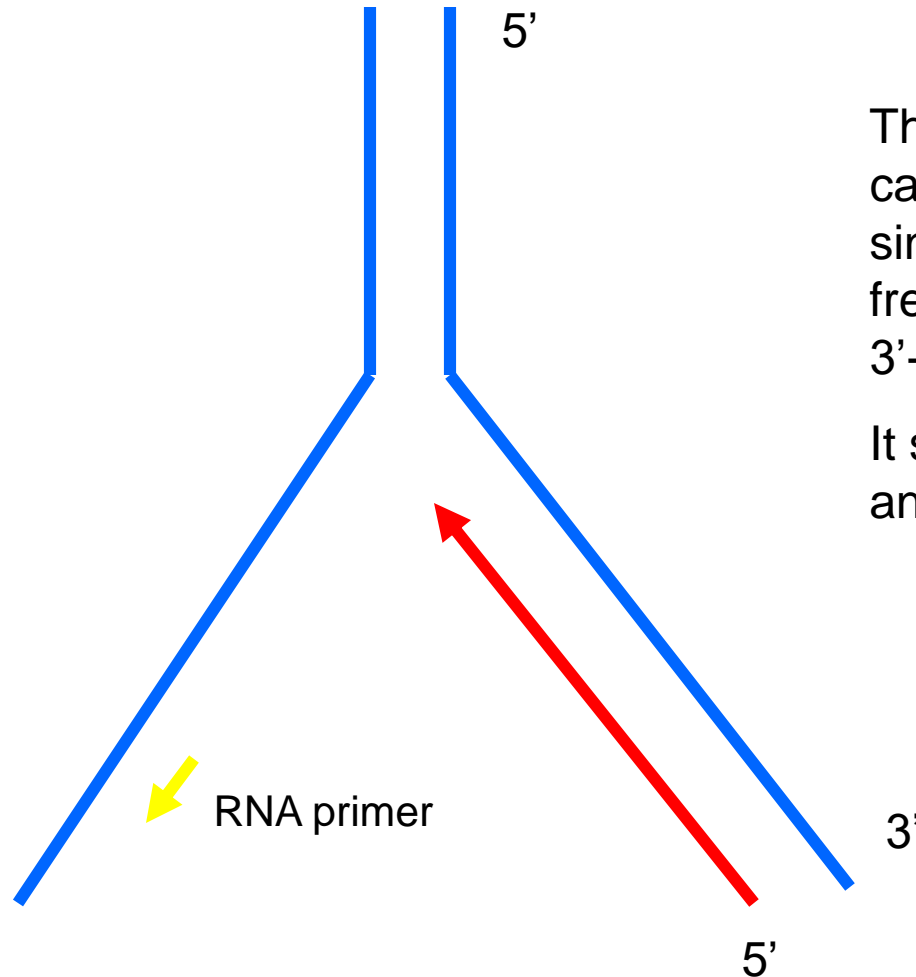
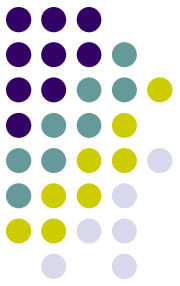


# DNA Replication



The leading strand  
is replicating  
without problems

# DNA Replication

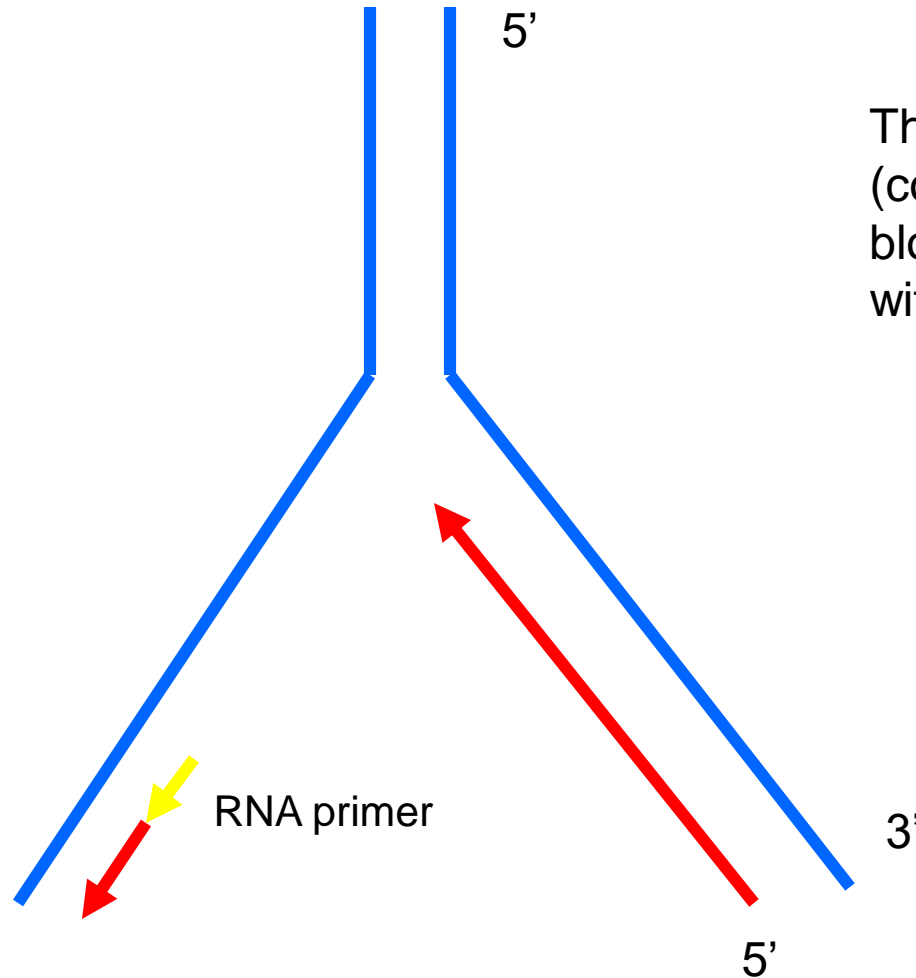
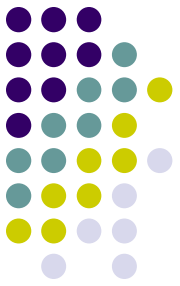


The lagging strand cannot even start, since there is no free complementary 3'-end

It starts by creating an RNA primer

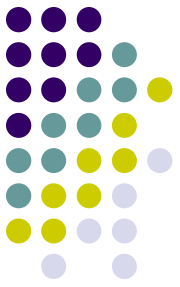


# DNA Replication

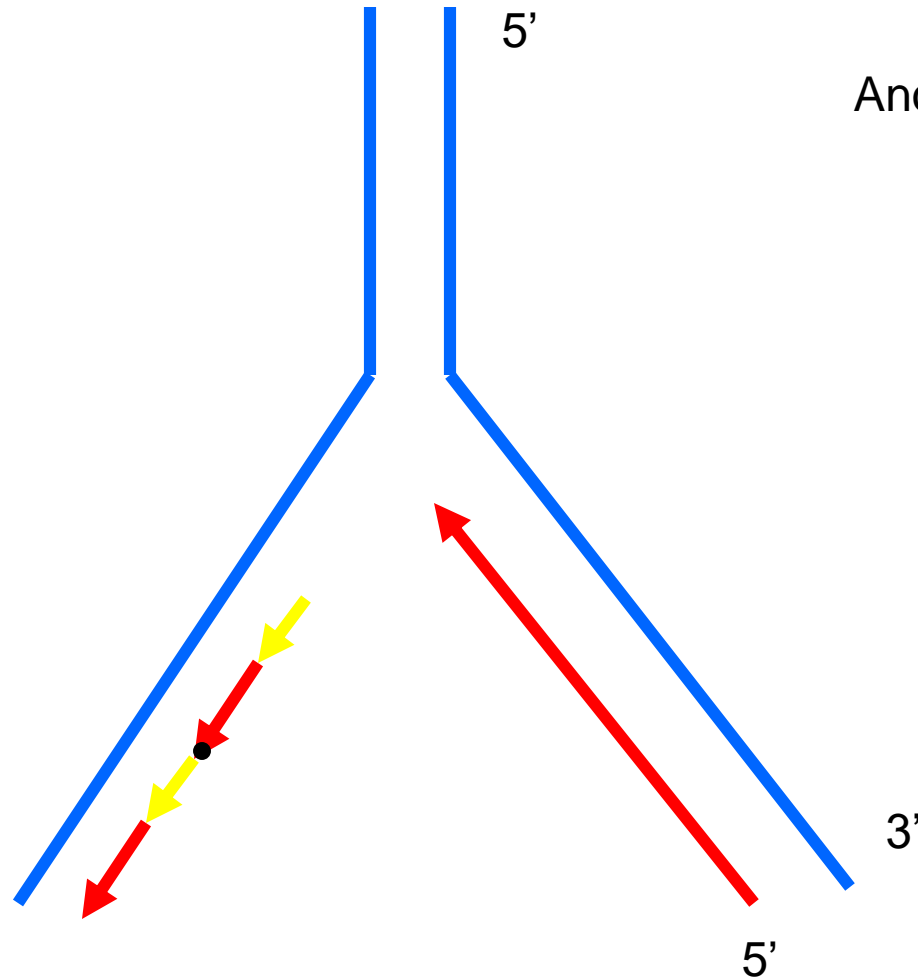


The primer  
(consisting of RNA  
blocks) is extended  
with DNA blocks

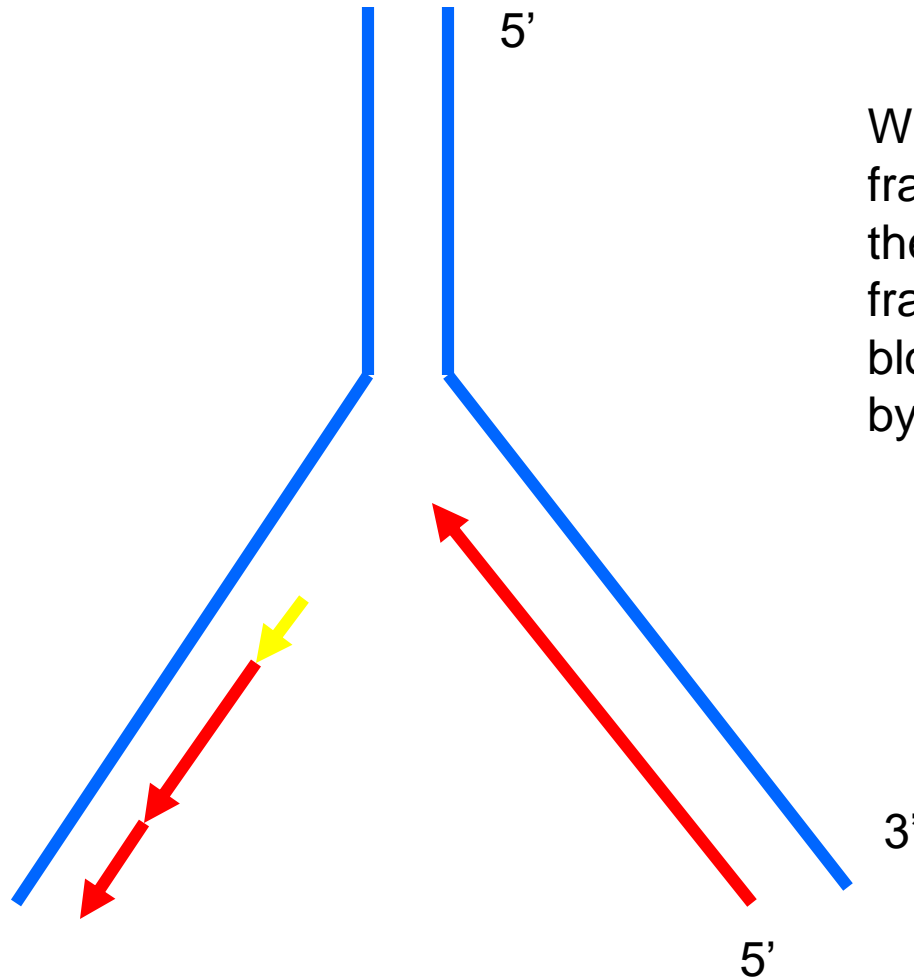
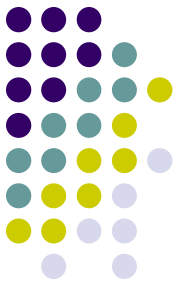
# DNA Replication



Another fragment

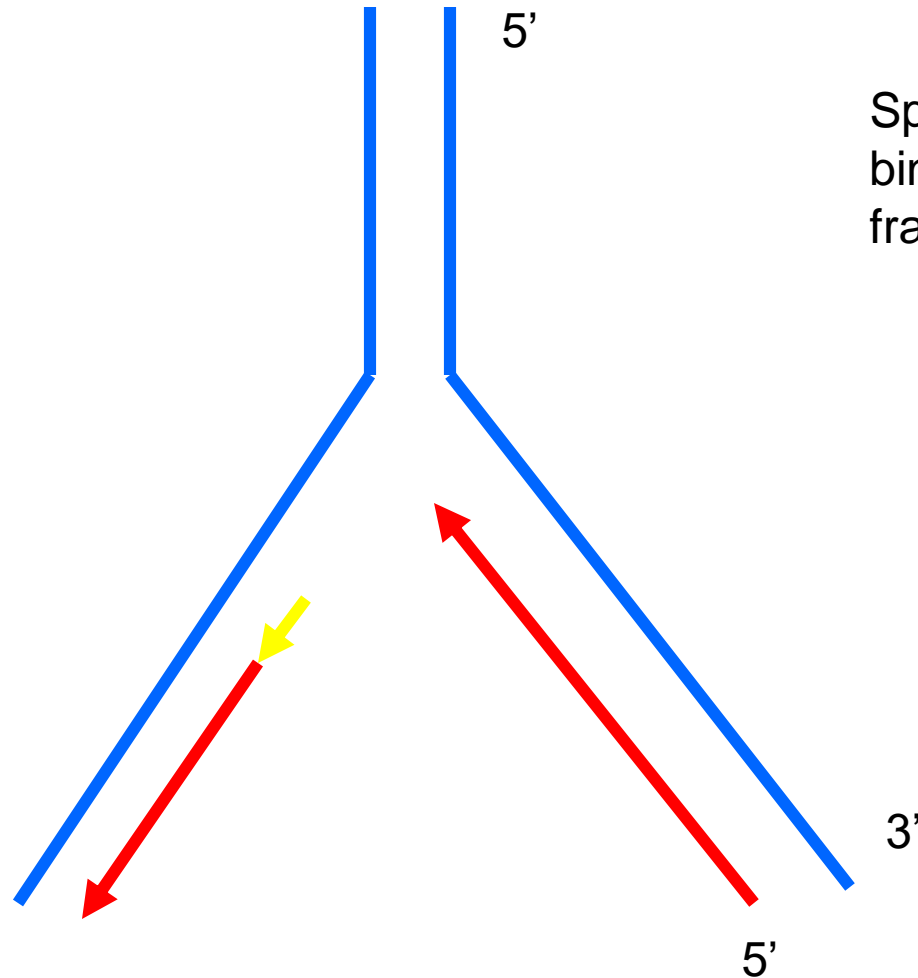
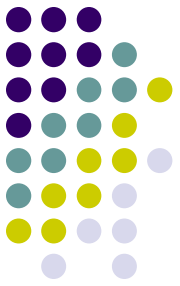


# DNA Replication



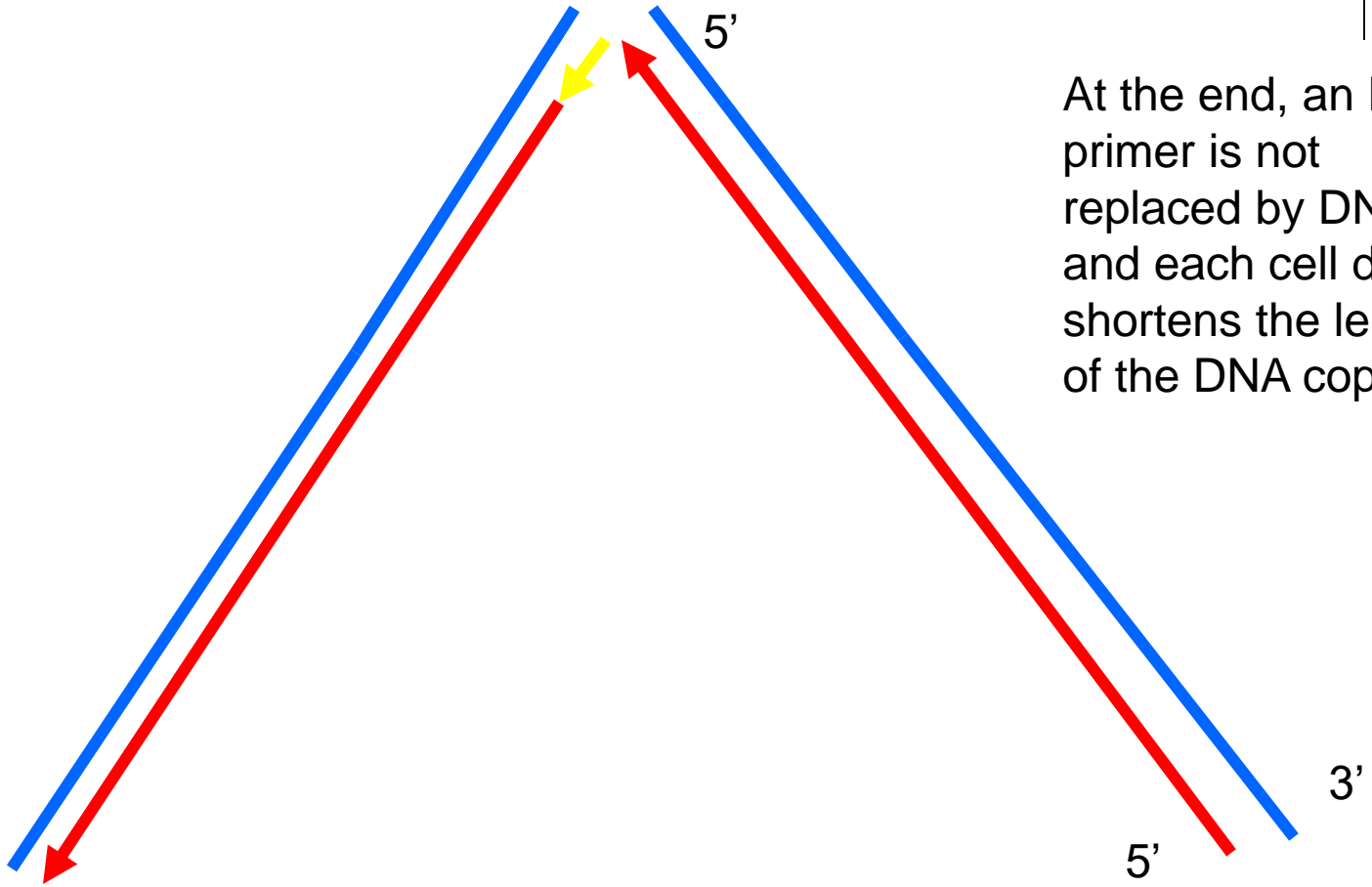
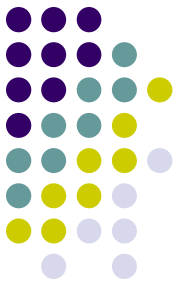
When the second fragment reaches the start of the first fragment, RNA blocks are replaced by DNA blocks

# DNA Replication



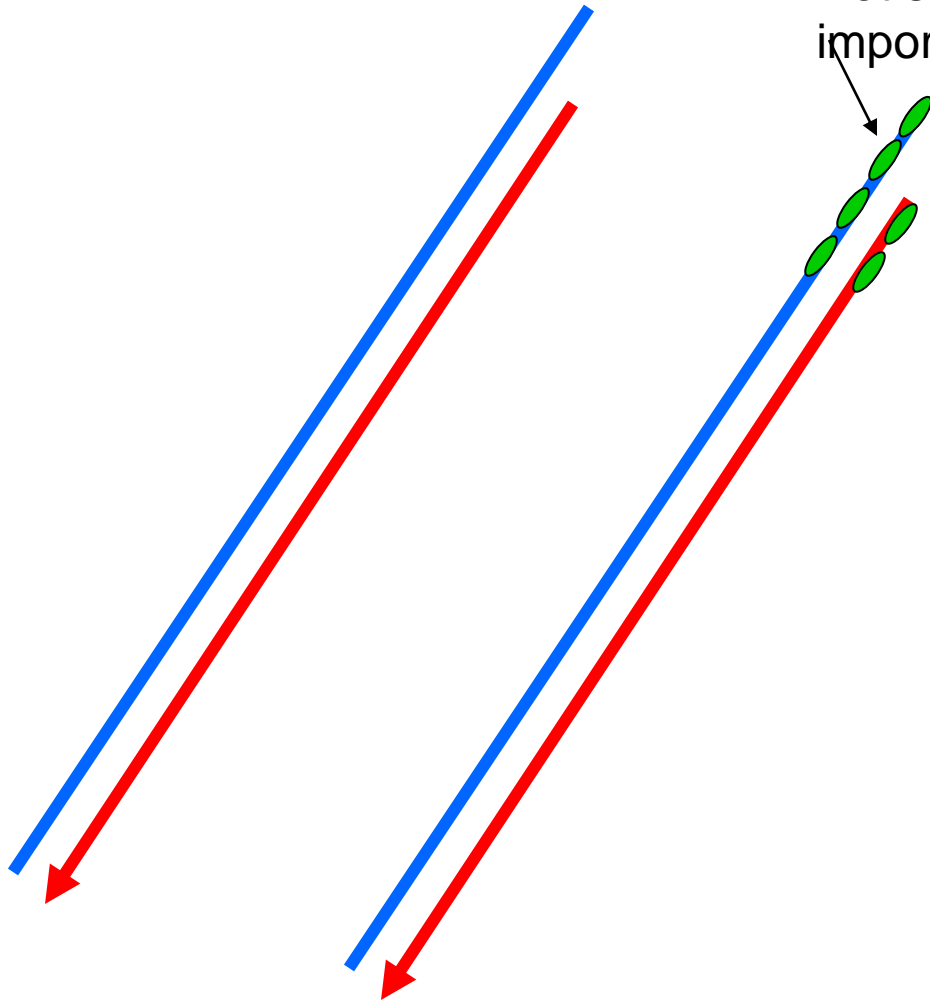
Special enzyme  
binds the DNA  
fragments together

# DNA Replication



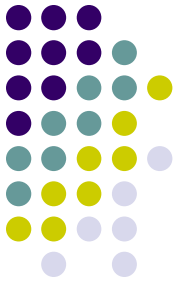
At the end, an RNA primer is not replaced by DNA, and each cell division shortens the length of the DNA copy

# Telomerase

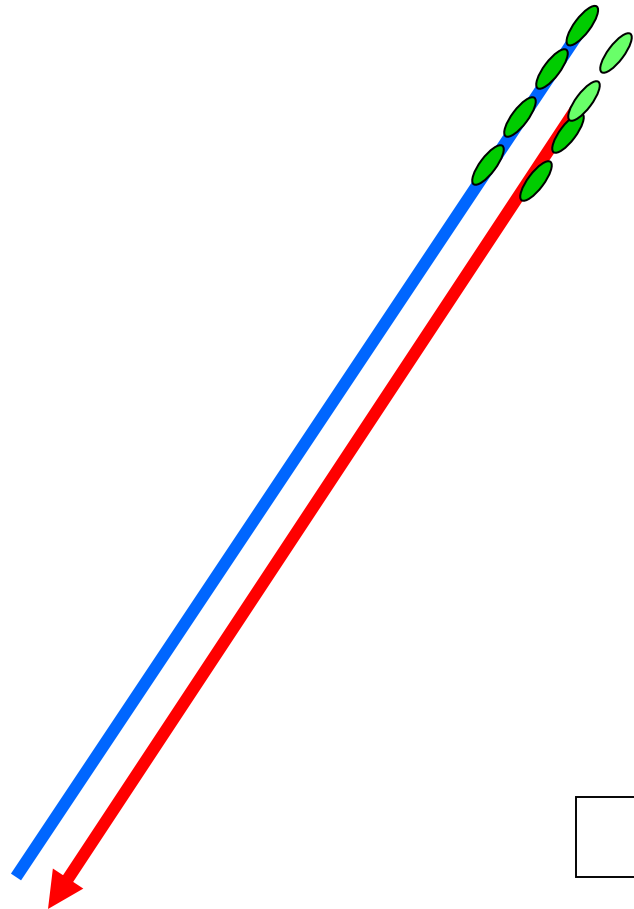
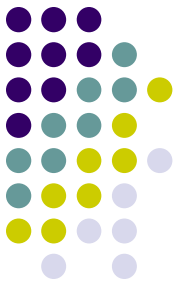


Not something  
important

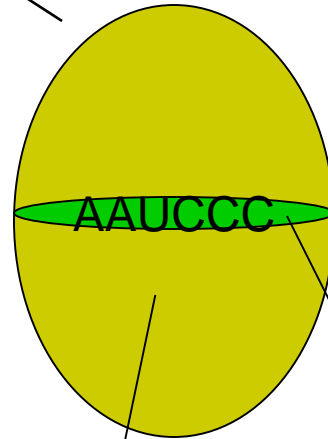
At the end of each  
chromosome there  
are no genes, but  
tandem repeats. For  
example, *TTAGGG*  
(Mammals)



# Telomerase

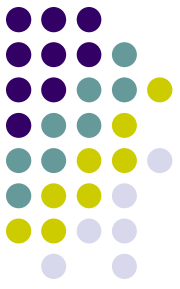


Special molecule –  
telomerase –  
rebuilds the identical  
repeats after each  
replication



Protein

RNA, which  
encodes a repeat

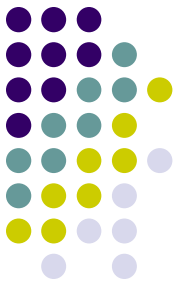


# Telomerase

- Hypothesis: the activity of telomerase is decreasing with aging, and the chromosomes start shortening
- When the shortage reaches the encoding zone, organism dies
- To prevent aging we cannot just add telomerase, since it also promotes cancer
- The knock-out mice without telomerase within several generation become early-aging



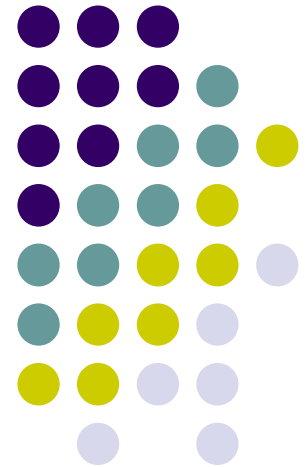
# With an efficient algorithm for repeats



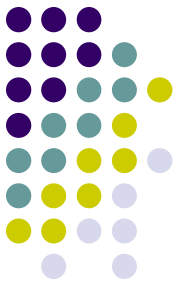
- We can discover new repeating sequences in genomes
  - New disease markers
  - New personal identifiers
  - New viral insertions

## 2. Longest common substrings

---

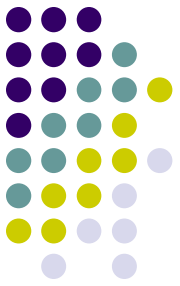


# The longest common substring of several strings



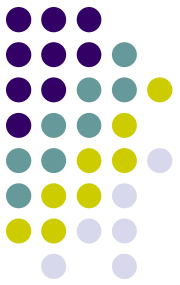
- The problem: find the longest substring common to two given strings I and II.
  - For example, if I=*superiorcalifornialives* and II=*sealiver*, then the longest common substring of I and II is *alive*.
- 1970 – Knuth conjectured that the linear-time solution to the longest common substring problem would be impossible

# The longest common substring for 2 strings in linear time

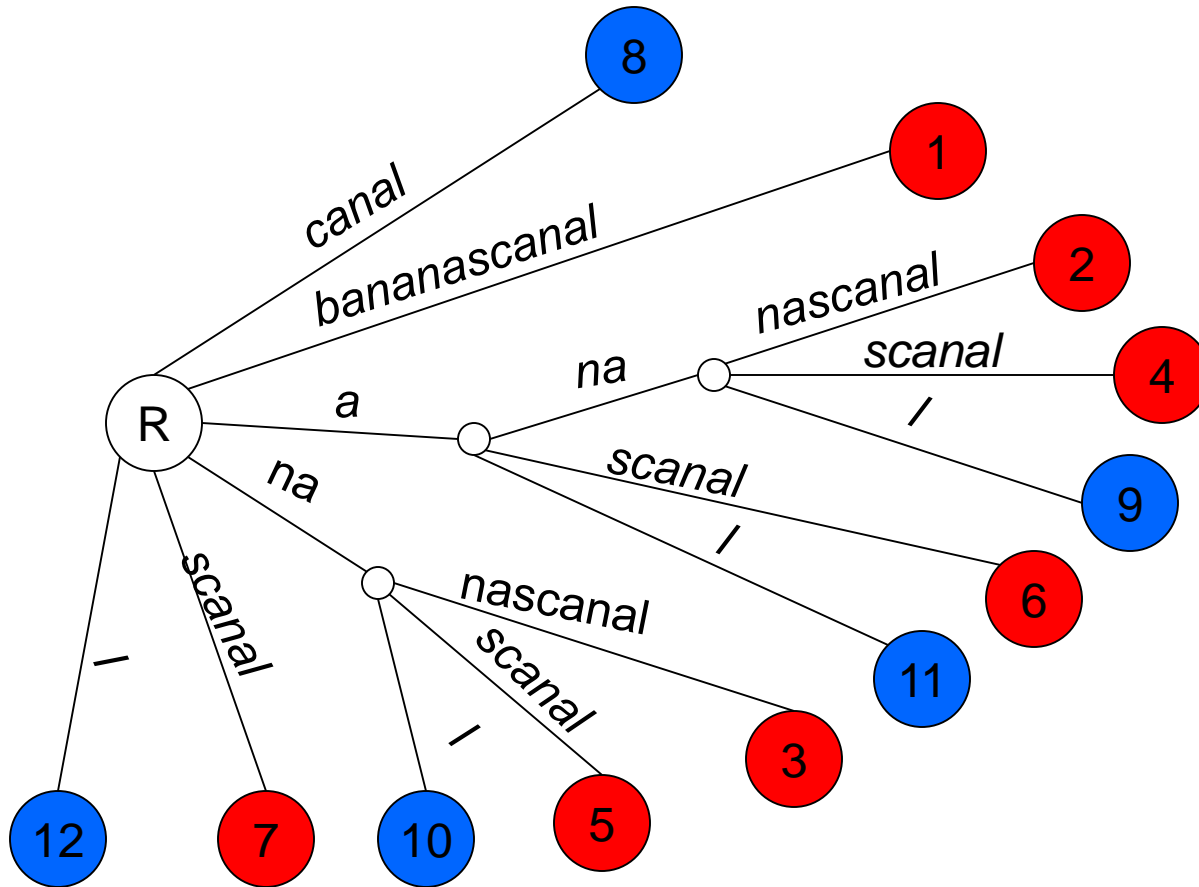


- Concatenate 2 strings and build the suffix tree for the concatenated string
- Label each leaf with the corresponding suffix start position, plus the ID of the string (I or II)
- Perform the depth-first traversal and mark each internal node by I, II or both, depending what suffixes are found in the subtree for this node
- Find the deepest internal node which is marked by both I and II

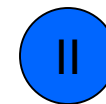
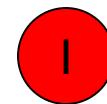




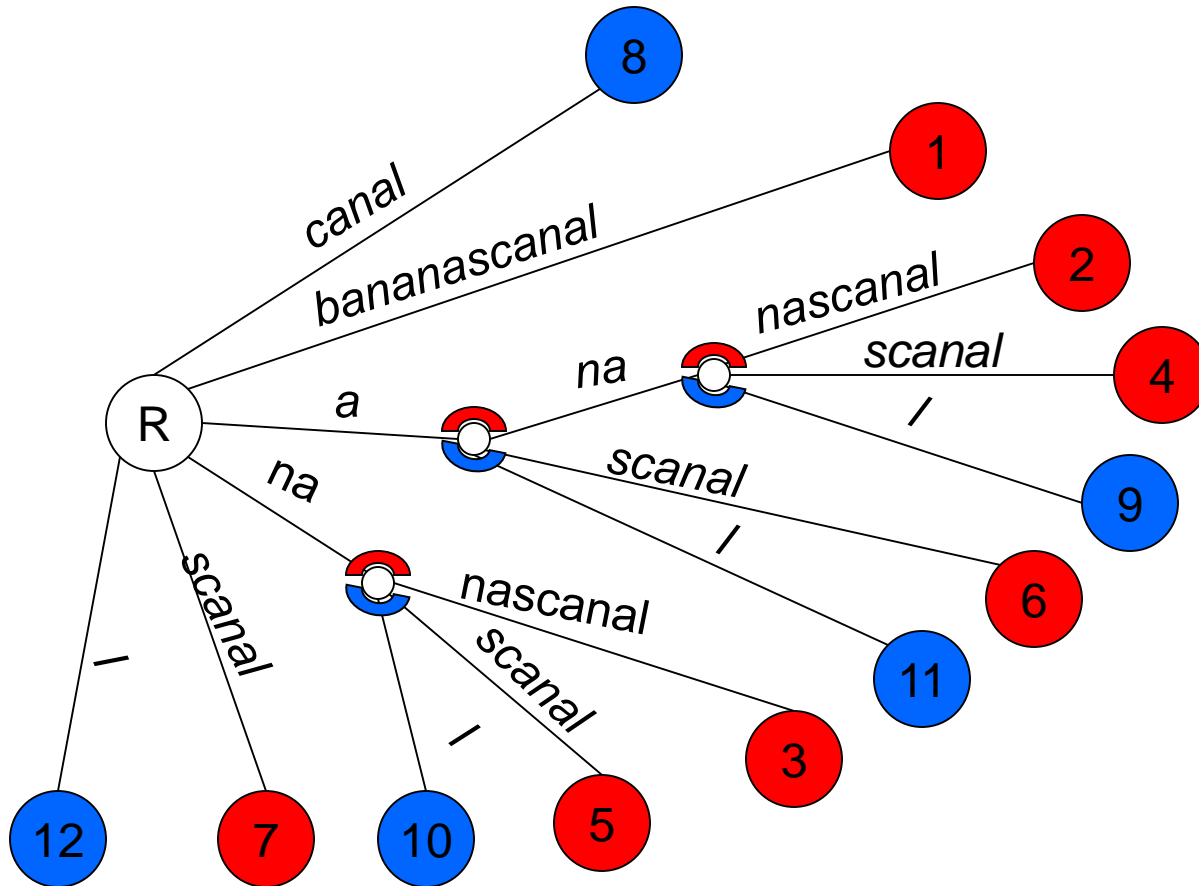
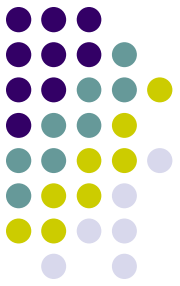
# Example: $I=bananas$ $II=canal$



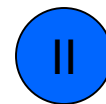
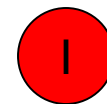
<i>b</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>s</i>	<i>c</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>l</i>
1	2	3	4	5	6	7	8	9	1	1	1
									0	1	2



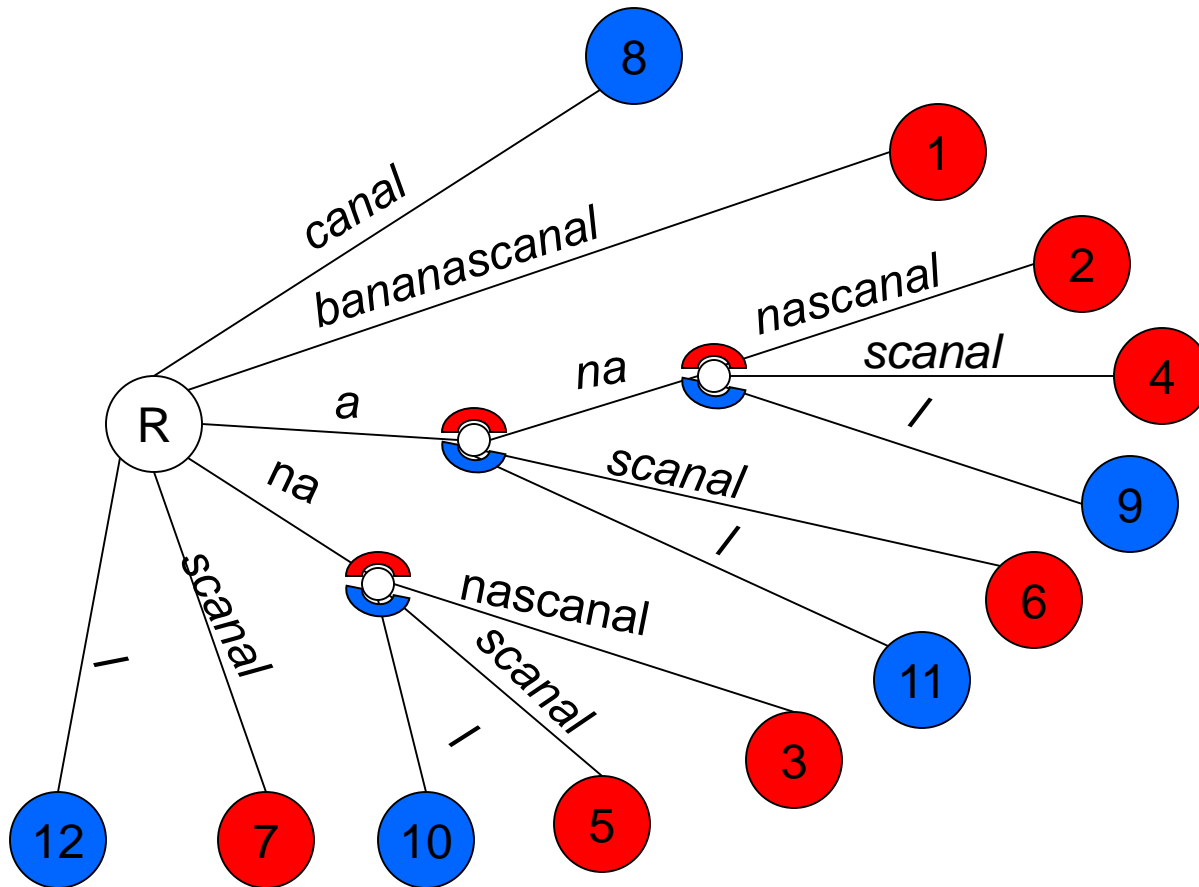
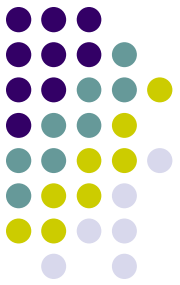
# Example: Marking internal nodes



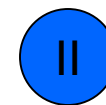
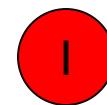
<i>b</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>s</i>	<i>c</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>l</i>
1	2	3	4	5	6	7	8	9	1	1	1
									0	1	2



# Example: What is the longest common substring?

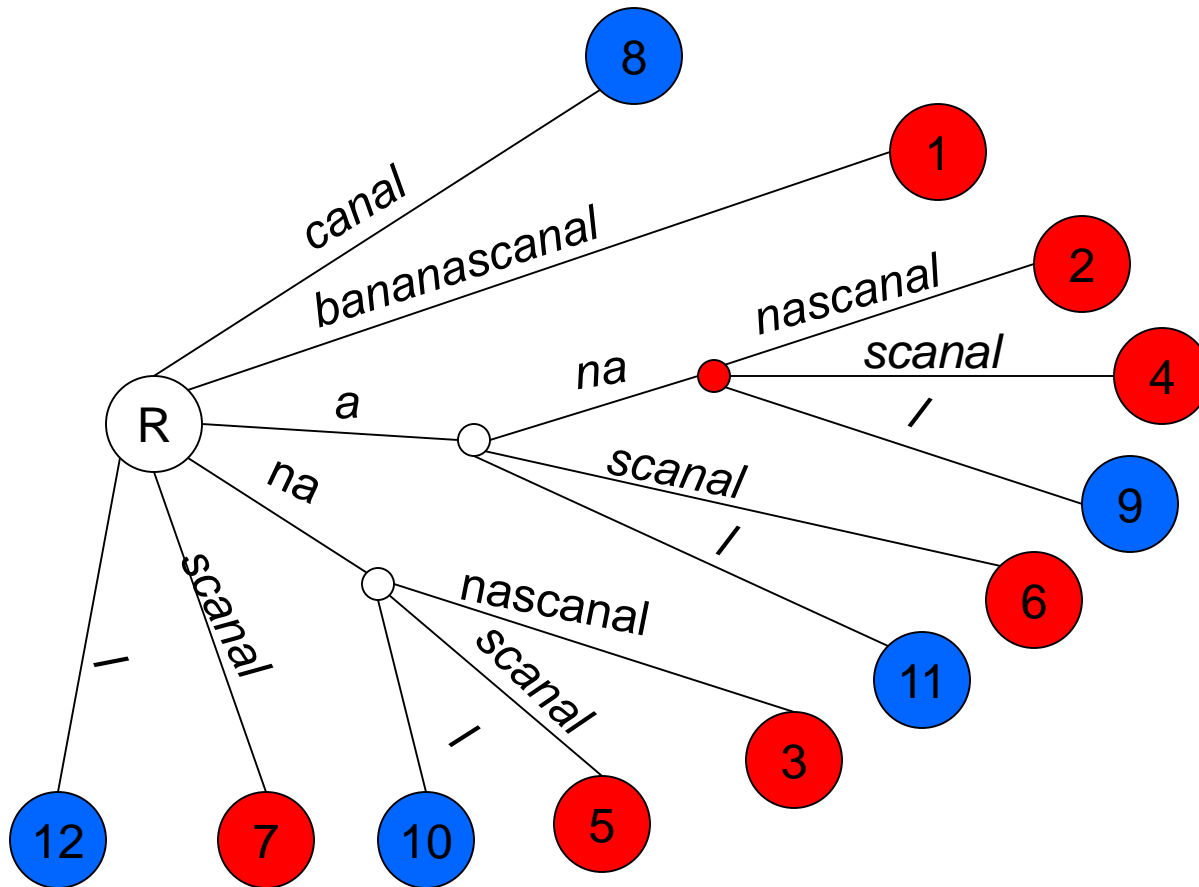
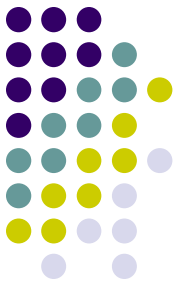


<i>b</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>s</i>	<i>c</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>l</i>
1	2	3	4	5	6	7	8	9	1	1	1
								0	1	2	

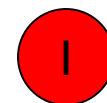




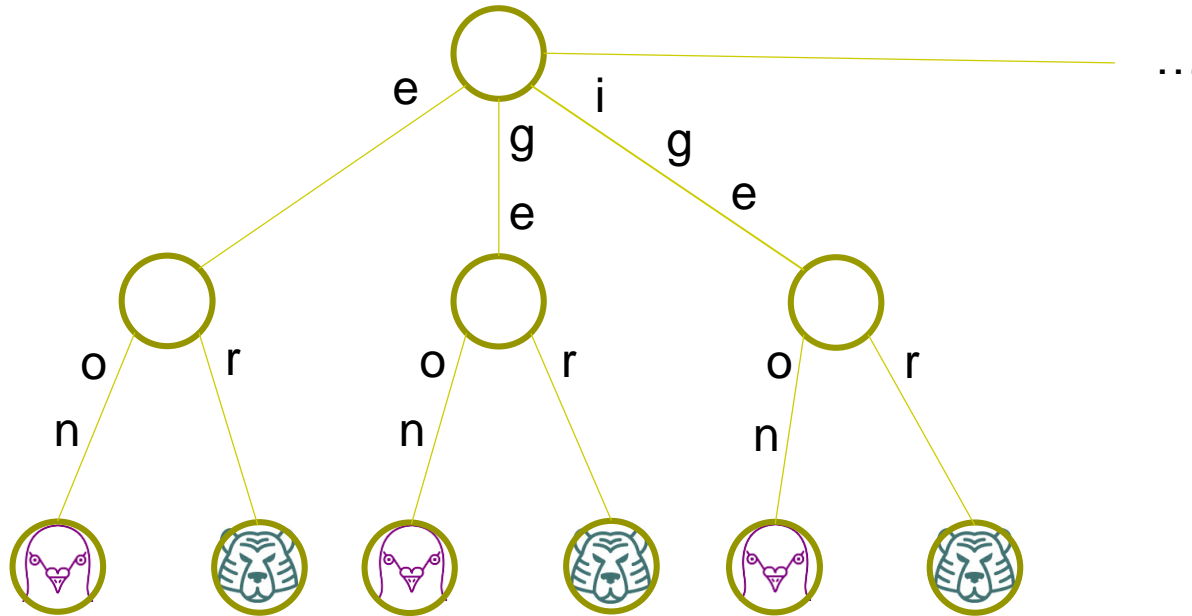
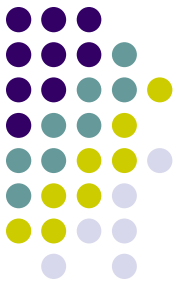
# Example: LCS=*ana*



<i>b</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>s</i>	<i>c</i>	<i>a</i>	<i>n</i>	<i>a</i>	<i>l</i>
1	2	3	4	5	6	7	8	9	1	1	1
									0	1	2

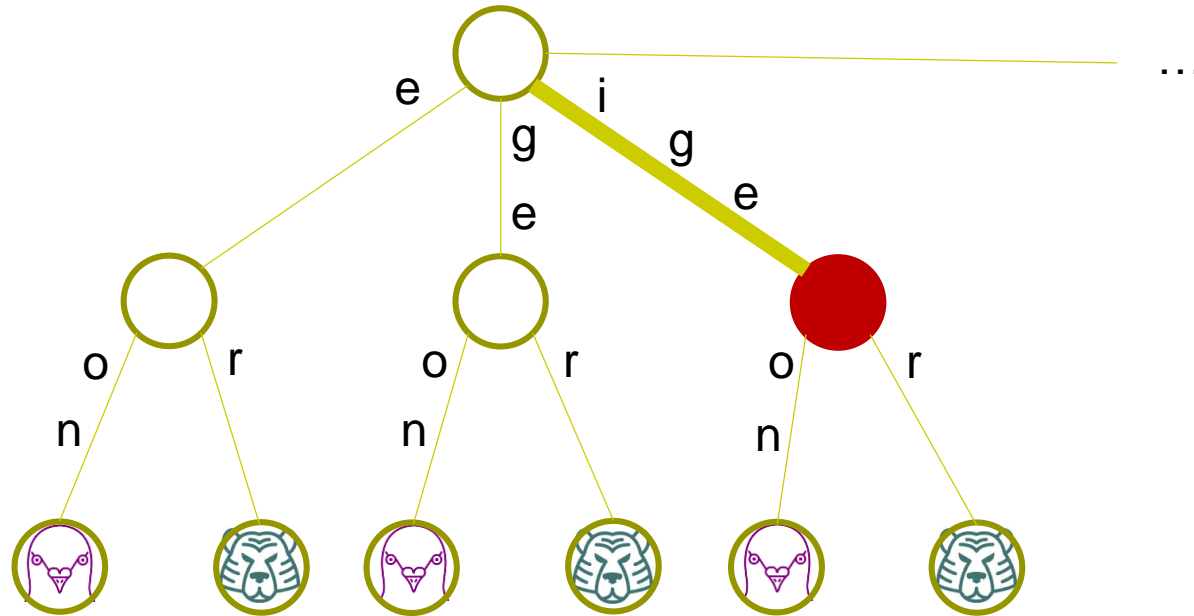
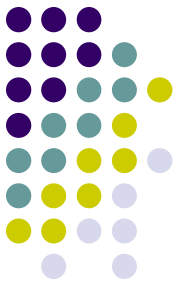


# Longest common substrings: example



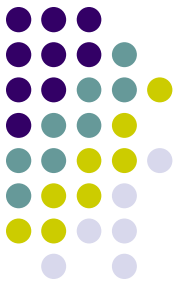
Query: what do *tiger* and *pigeon* have in common?

# Longest common substrings: example



Query: what do *tiger* and *pigeon* have in common?

# Common substrings for a set of DNA sequences



Insert suffixes of multiple strings into one tree

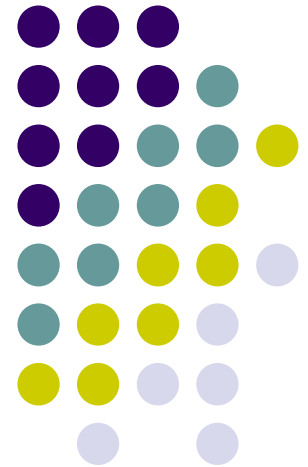
- Discover substrings common to viruses and humans
- Discover substrings unique to cancer

Used in the identification of the remains of US military personnel

- Mitochondrial DNA from live person is collected, sequenced and the sequences are stored in the database (I)
- Later, the DNA is extracted from the remains (II), and the longest common substring of I and II helps to narrow down the search

---

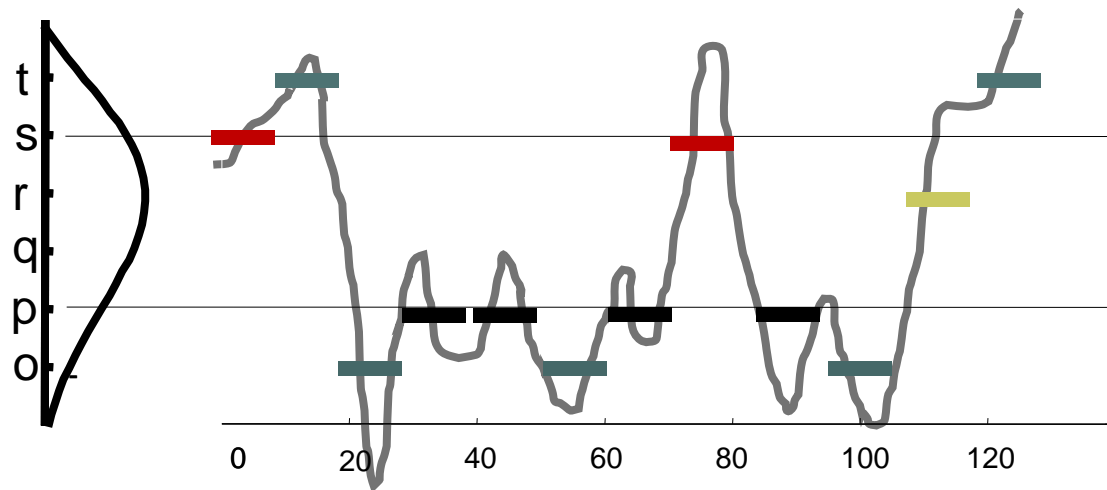
# Potential applications of suffix trees for any sequential data



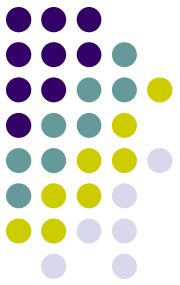


# Time series as strings

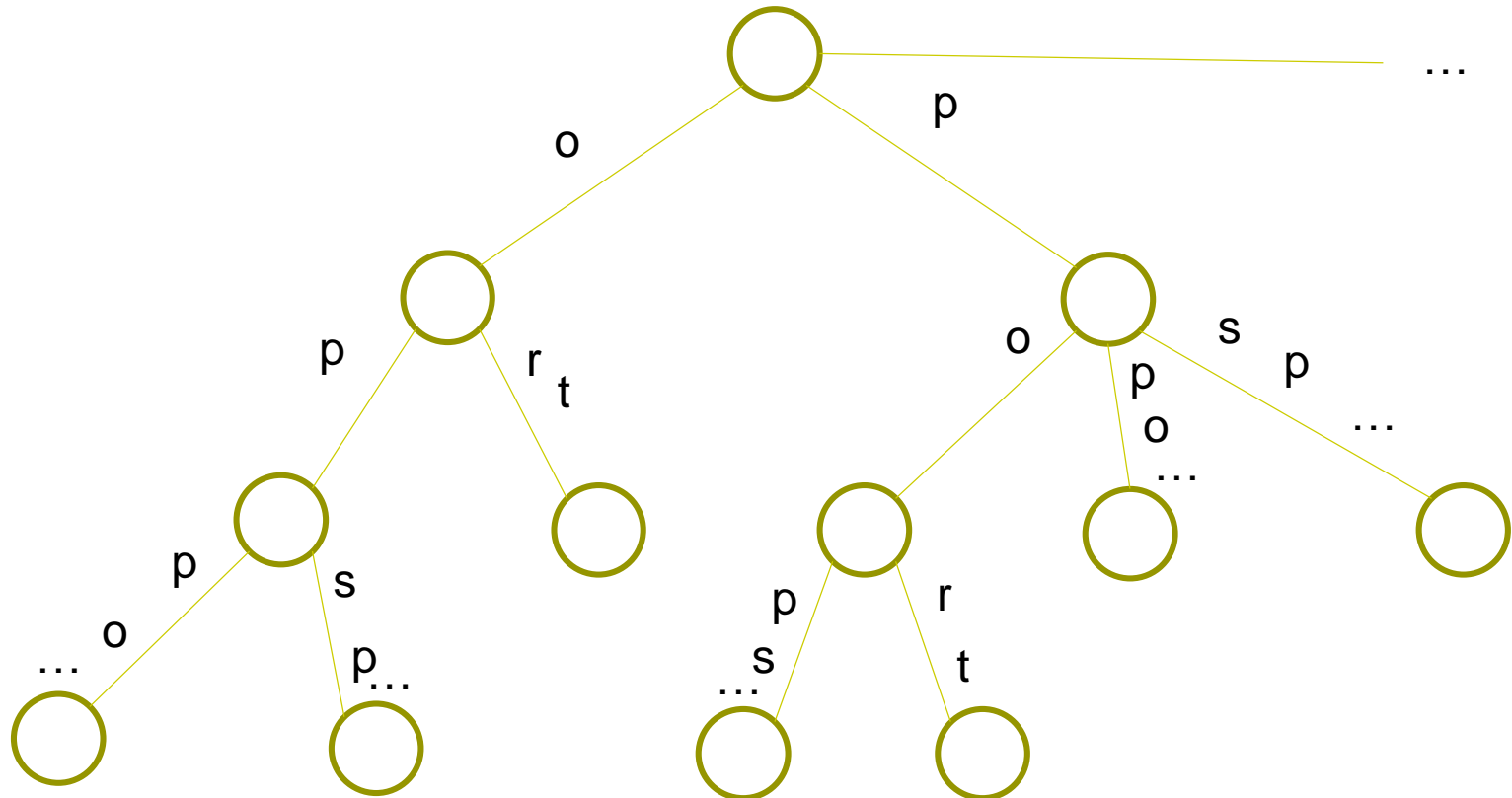
SAX - Symbolic Aggregate approximation (by Eamon Keough, 2001)



*stoppopsport*



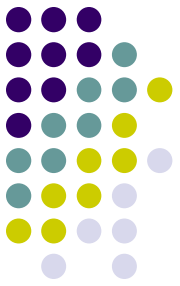
# Suffix trees for time series



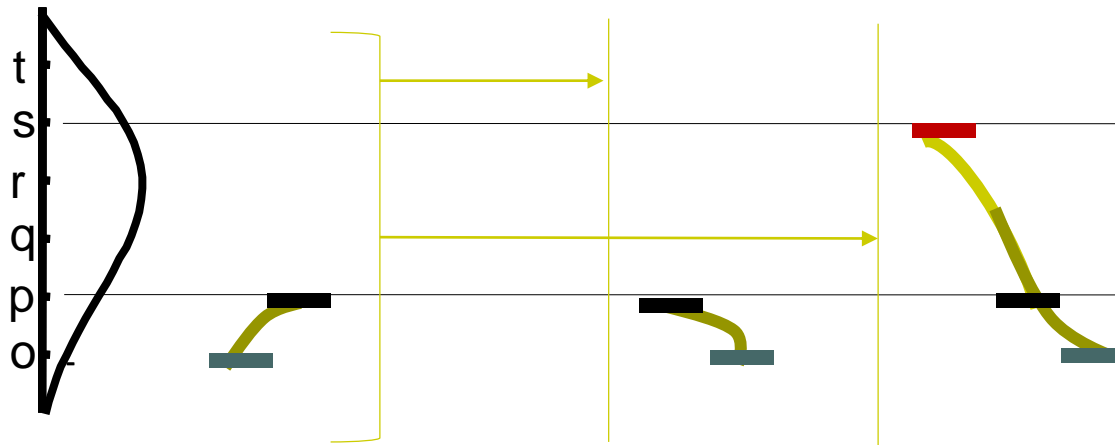




# Suffix trees for time series: rise and fall of stocks - prediction



50% *po*, 50% *spo*



Query: what happened after *op*?

Web

**Images**

Videos

News

More ▾

Search tools

Size ▾

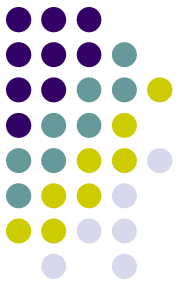
Color ▾

**Clip art** ▾

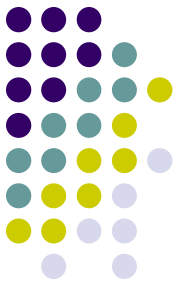
Time ▾

Usage rights ▾

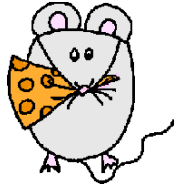
More tool



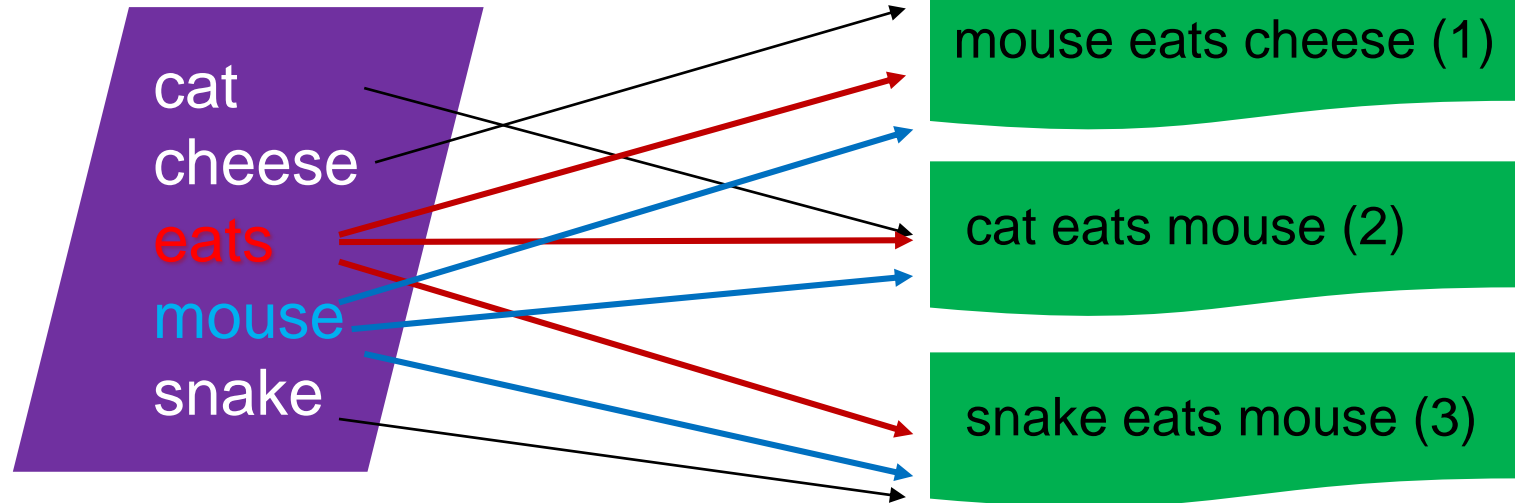
# Meaningful search: not with inverted index



Query: What animal “**eats mouse**”

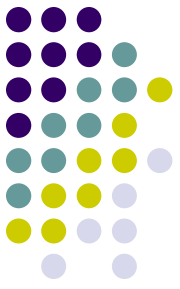


WORD-BASED INDEX  
(inverted index):



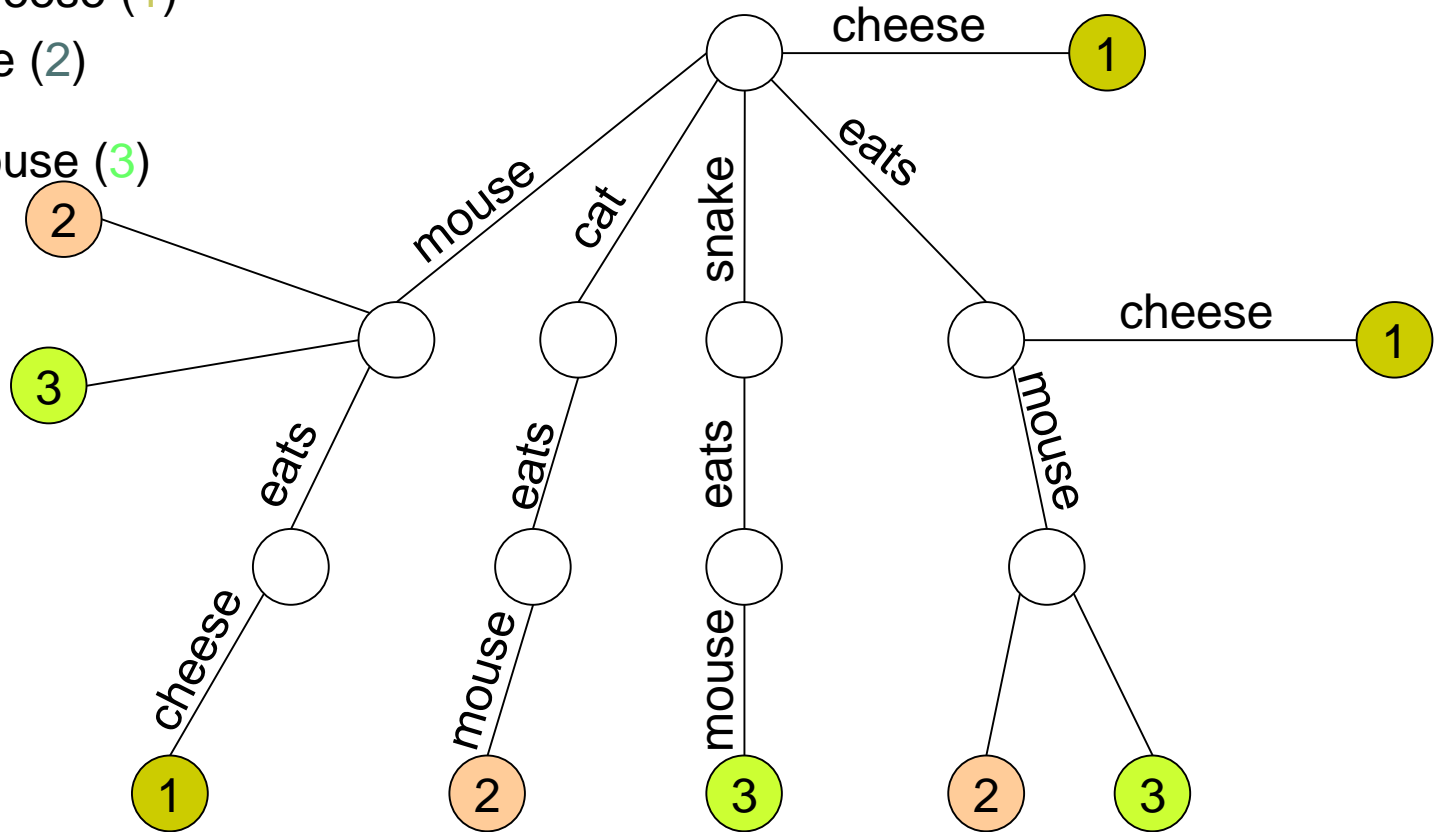
Answer is in documents 1,2,3

# Meaningful search: example

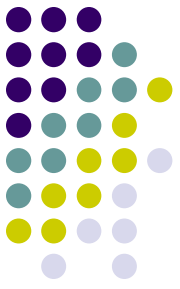


Collection of 1-sentence documents

- ❑ mouse eats cheese (1)
- ❑ cat eats mouse (2)
- ❑ snake eats mouse (3)



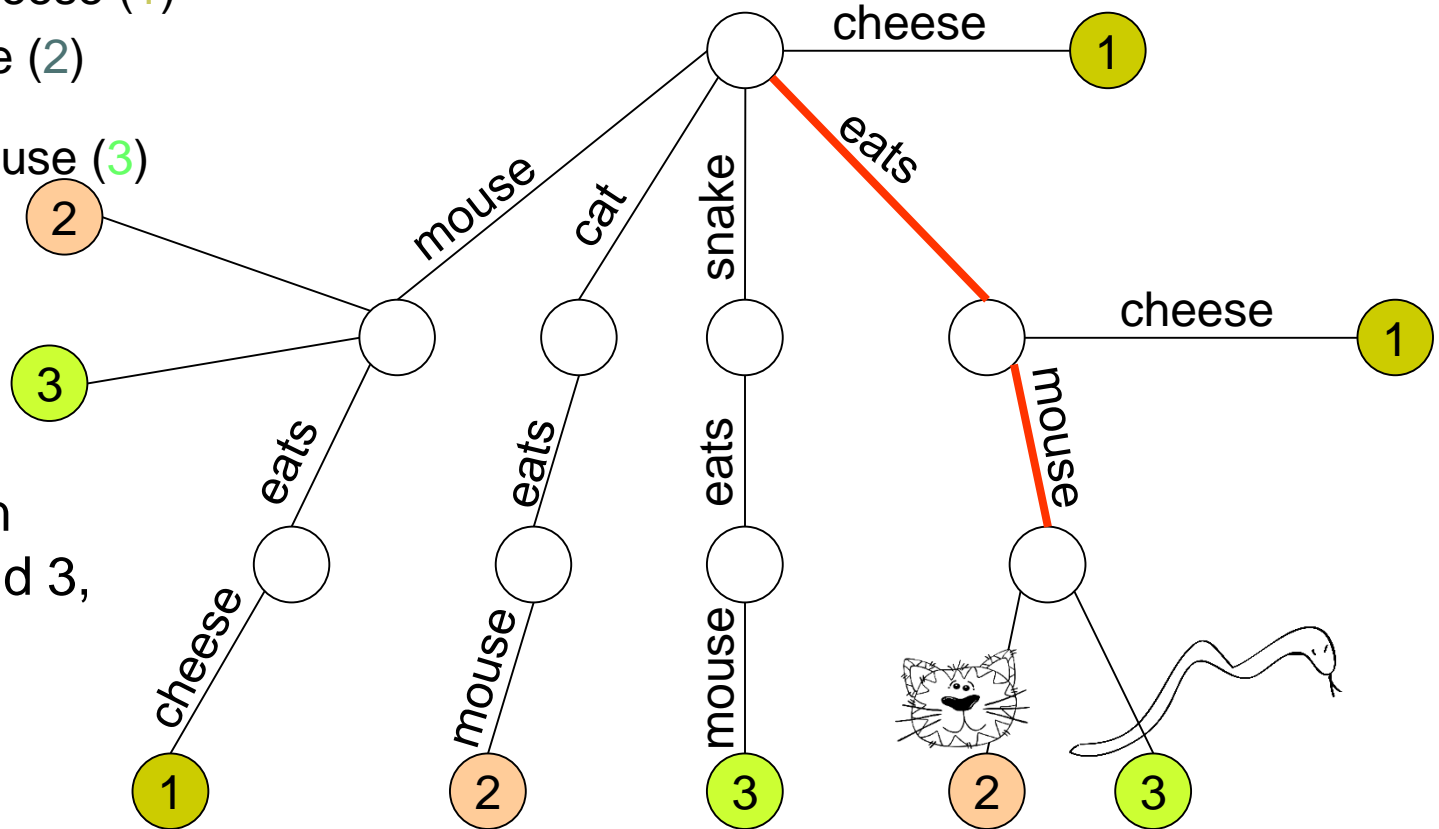
# Meaningful search: example



Query: What animal “**eats mouse**”

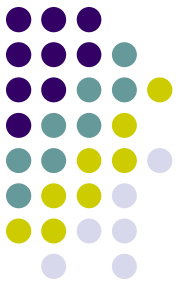
Collection of 1-sentence documents

- ❑ mouse eats cheese (1)
- ❑ cat eats mouse (2)
- ❑ snake eats mouse (3)



The answer is in documents 2 and 3, but not in 1

# Suffix tree for melodies ...



Saint-Saëns, Camille (1835-1921), Carnival des Animaux, Orch. & 2 Pfts., Aquarium

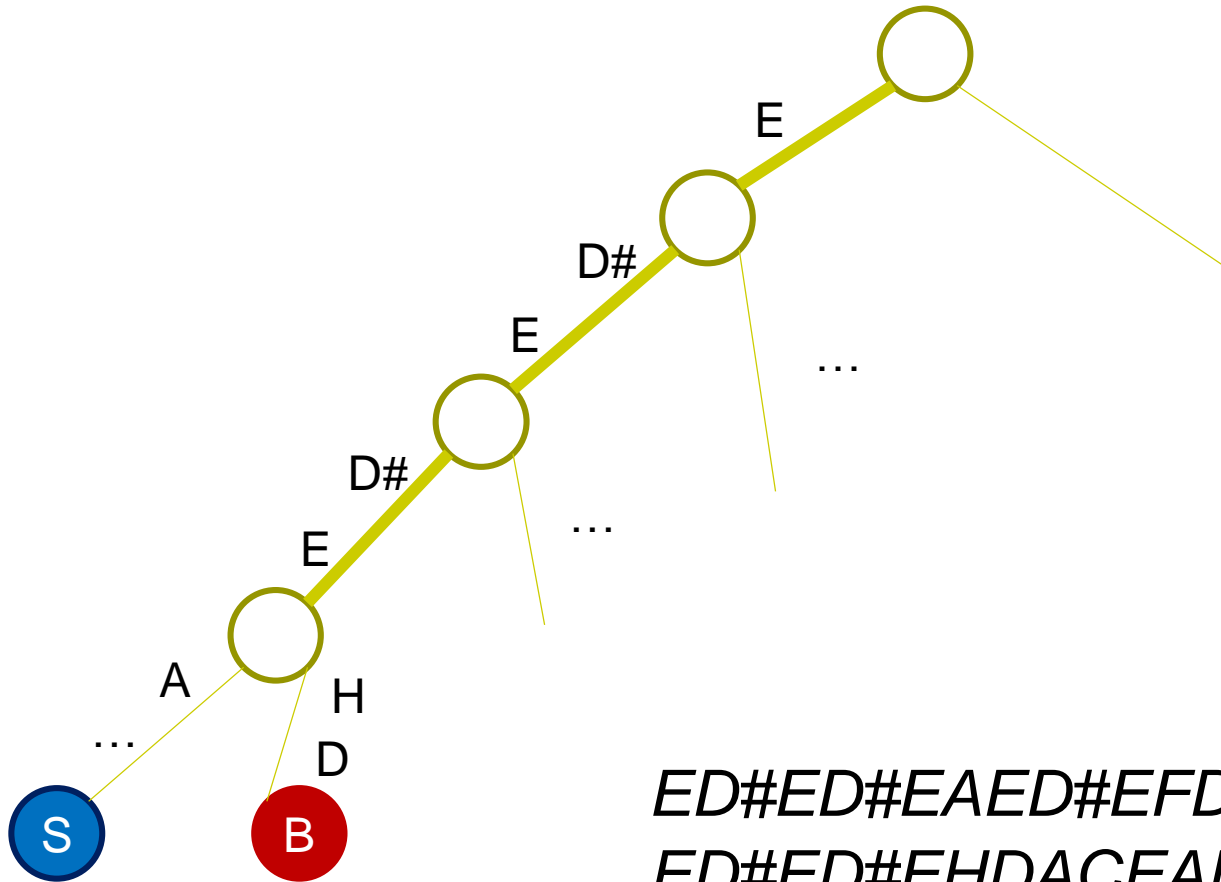
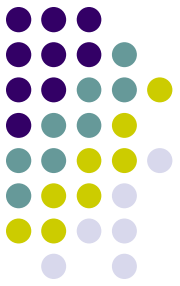


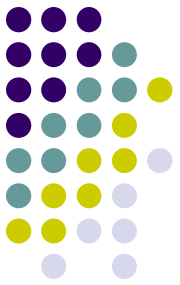
Beethoven, Ludwig Van (1770-1827), Für Elise, Pft.

*ED#ED#EAED#EFDC#DECHCDH* (**S-S**)

*ED#ED#EHDACEAHEG#C* (**B**)

# Suffix tree for melodies and plagiarism detection





# Indexing melodies: example

## Song 1

F [BALLADE]  
[Zu Strassburg steht ein hohes Haus]  
REG[Deutschland / Frankreich, Lothringen]  
MEL[-5\_ 1\_.23\_4\_ 2\_.31\_  
-5\_ 1\_.23\_4\_ 2\_231\_  
3\_ 5\_5\_5\_66 5\_2\_  
2\_ 5\_4\_3\_2\_ 1\_-6\_-5\_  
-5\_ 1\_2\_3\_4\_ 2\_\_1\_ //] >>  
FCT[Ballade, Braut - Werbung, Erpressung]



## Song 2

F[KRIEGS]  
[In Boehmen liegt ein staedtchen]  
REG[Deutschland, Hessen, Marburg]  
MEL[-5\_ -5\_.33\_3\_ 3\_\_1\_  
3\_ 5\_.55\_6\_ 5\_0\_  
5\_ 7\_.67\_6\_ 6\_5\_  
4\_ 3\_5\_2\_5\_ 1\_\_0\_ //] >>  
FCT[Staende -, Soldaten -, Kriegs – Lied]



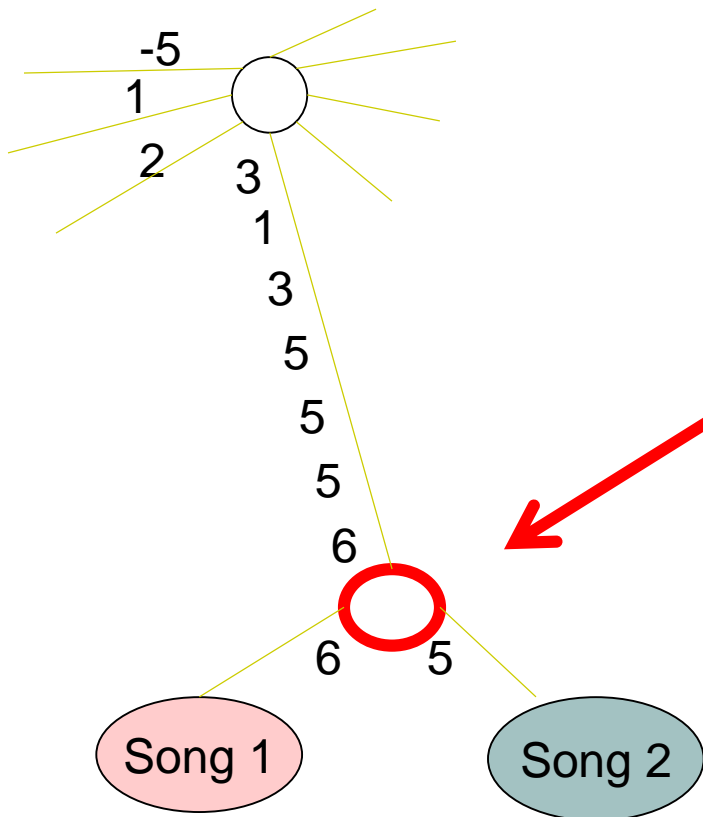
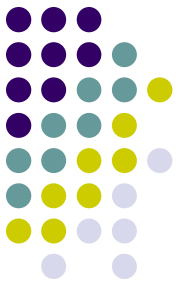
2 folk songs from the Essen Associative Code (EsAC) database

<http://www.esac-data.org/data/>



# ...and plagiarism detection

## Generalized suffix tree for two songs



The longest  
common  
substring

