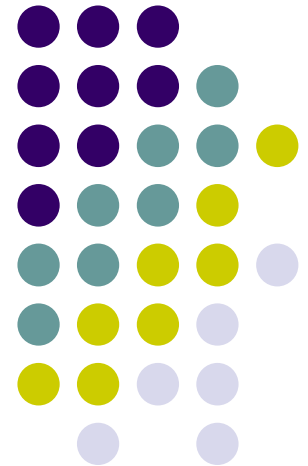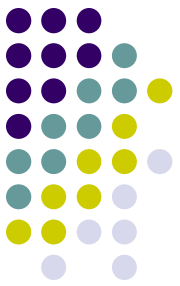# Applications of suffix trees

Lecture 3.2
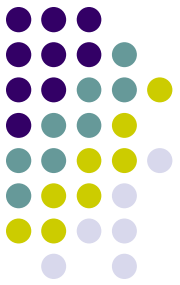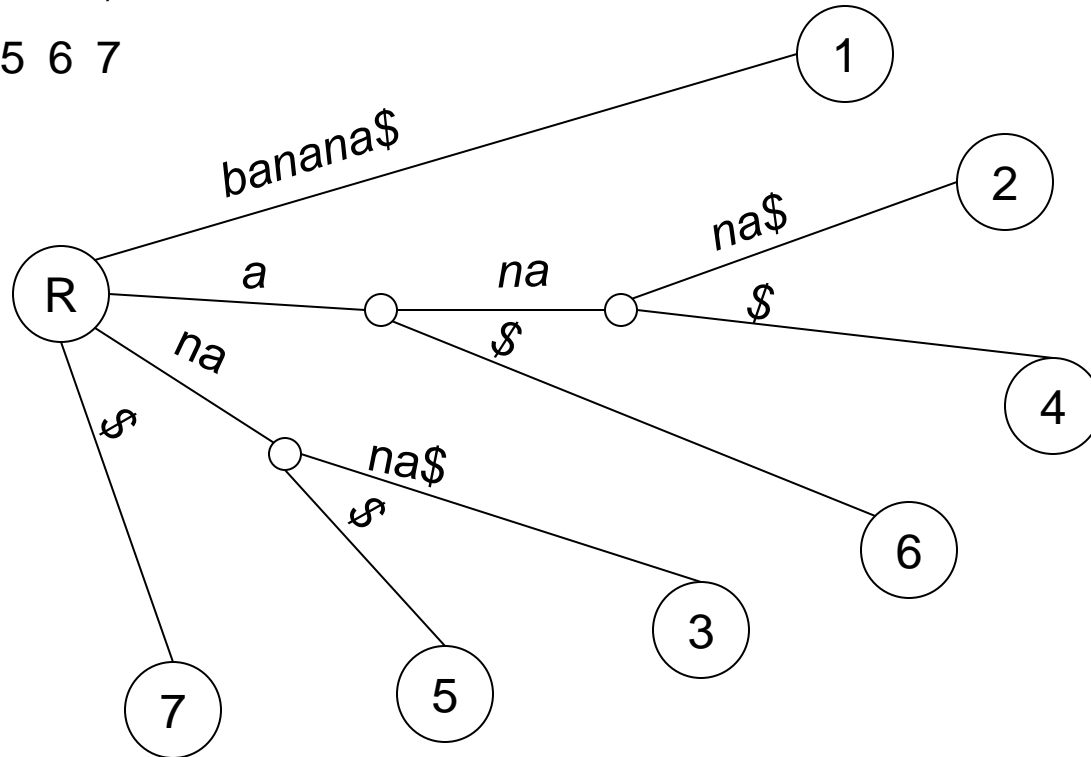*by Marina Barsky*

# Suffix tree - recap

- 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 add at the end 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
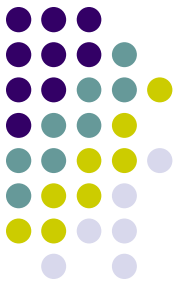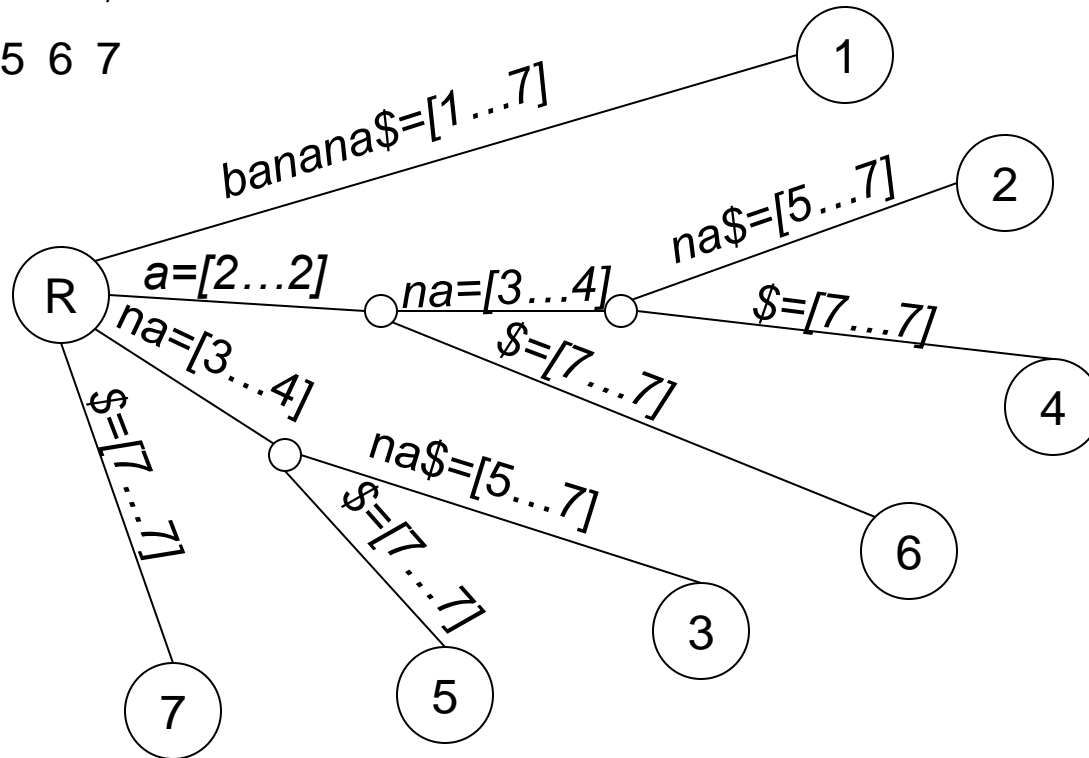
# Example: *banana$*

*b a n a n a $*

1 2 3 4 5 6 7

# Space reduction
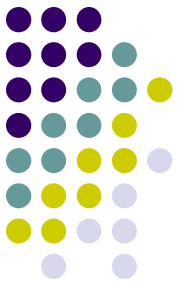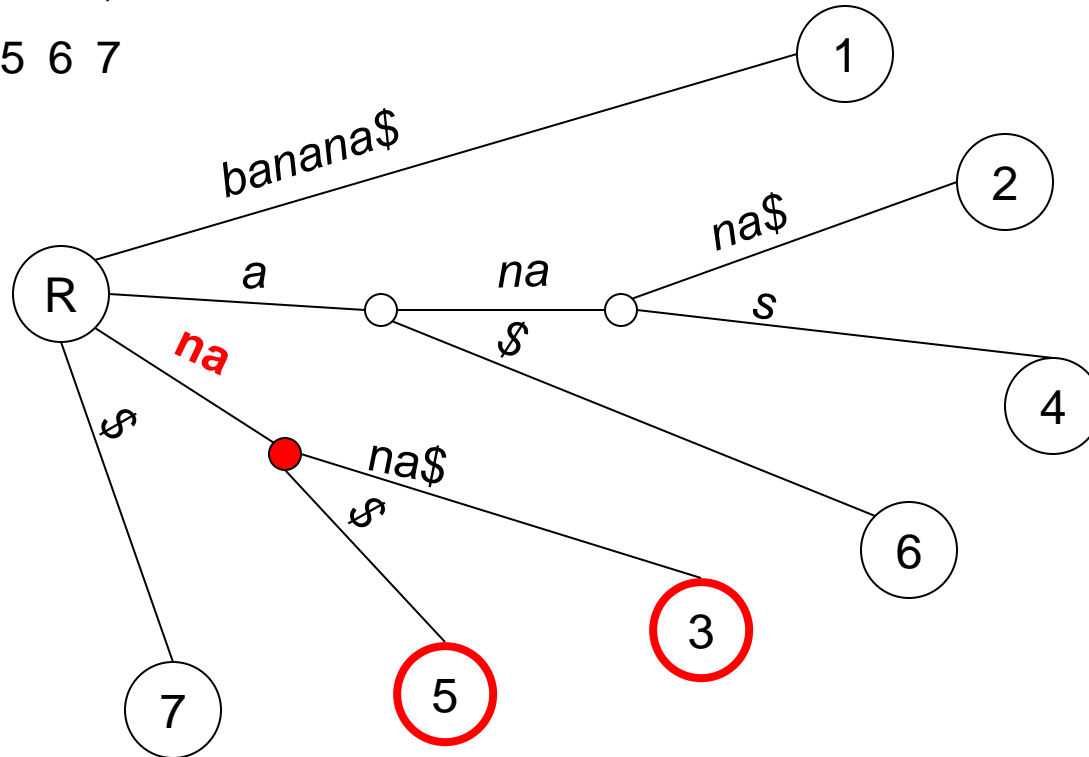
*b a n a n a $*

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 rooted at the node located 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*
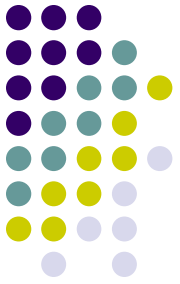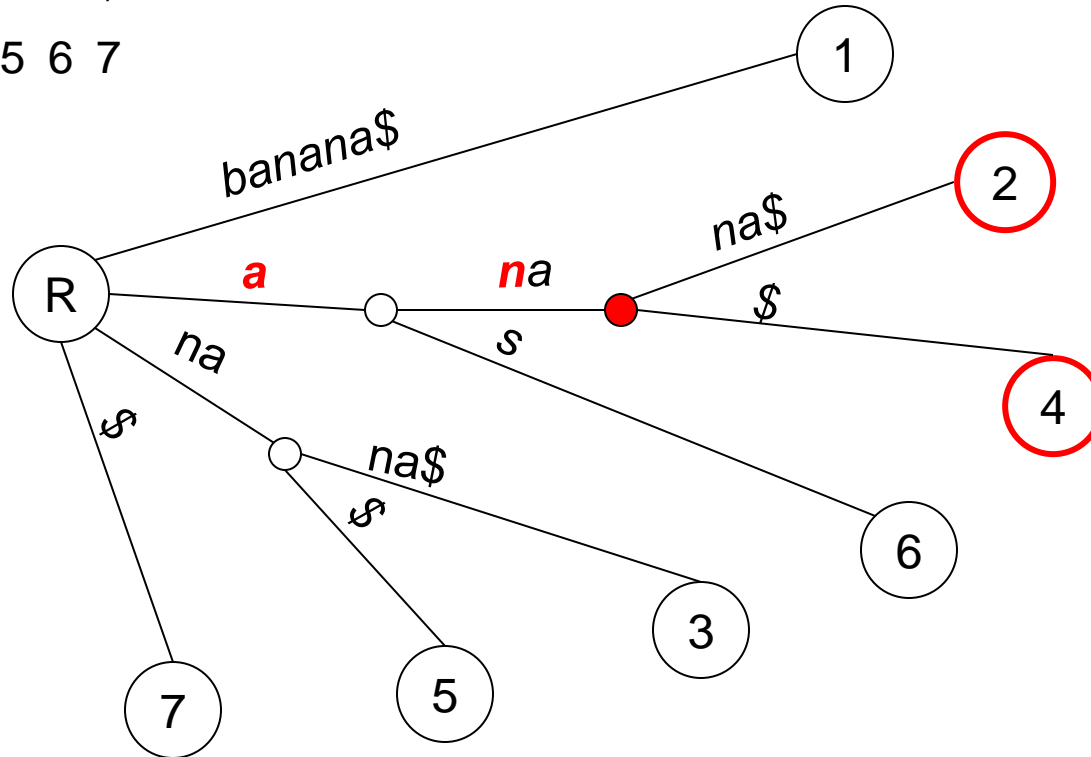
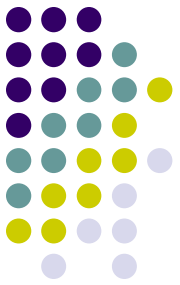# Example: *P=na*

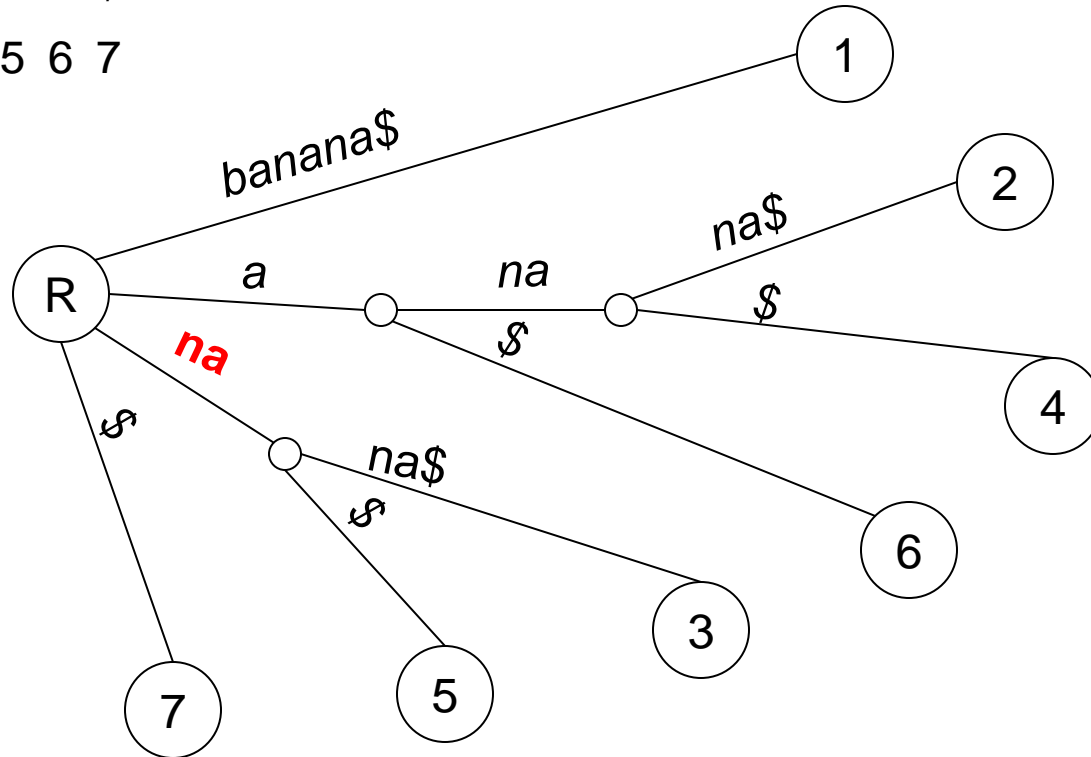*b a n a n a $*

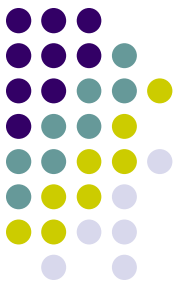1 2 3 4 5 6 7

# Example: *P=an*

*b a n a n a $*

1 2 3 4 5 6 7

# Example: *P=naa*

*b a n a n a $*

1 2 3 4 5 6 7



banana$
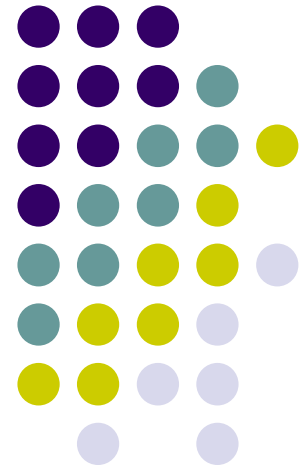
R

a

na

na$

$

na

$

$

na$

$

1
2
4
6
3
5
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 |
|---|---|---|---|---|---|
| | | | Average | 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 | **0.01 sec** | **0.03 sec** | |
| 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

These are my experiments with large disk-based suffix trees
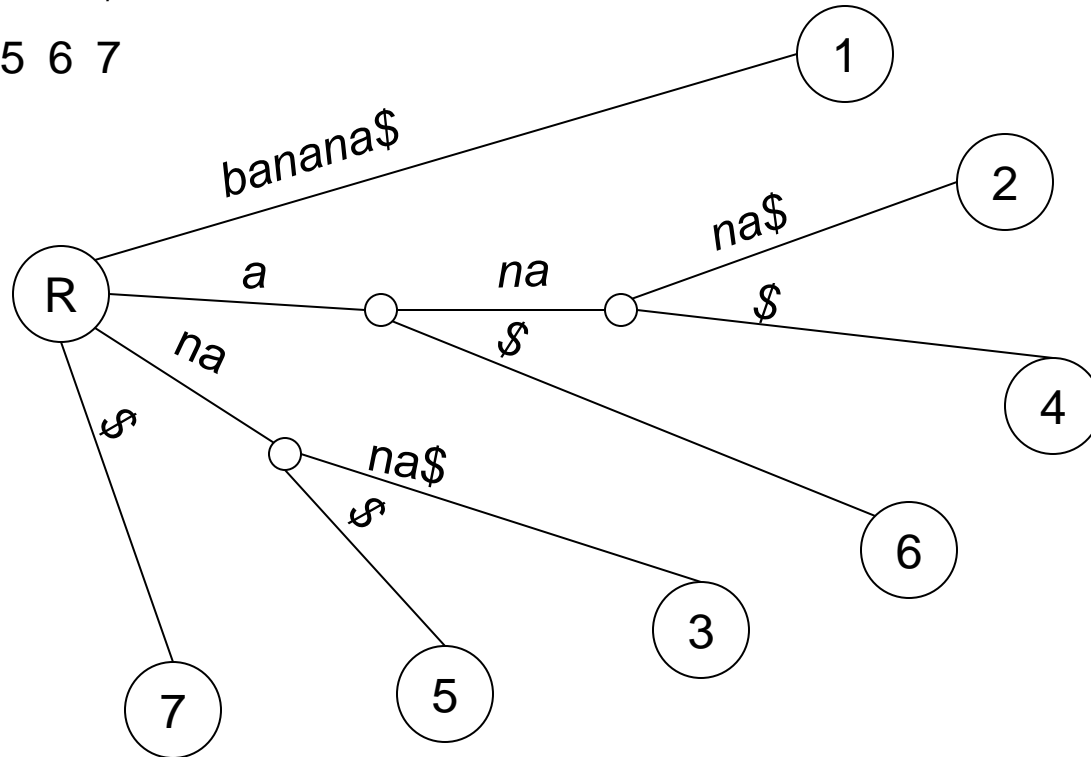
9

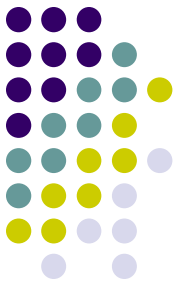# 1. Finding repeats

# Example: *bananas*

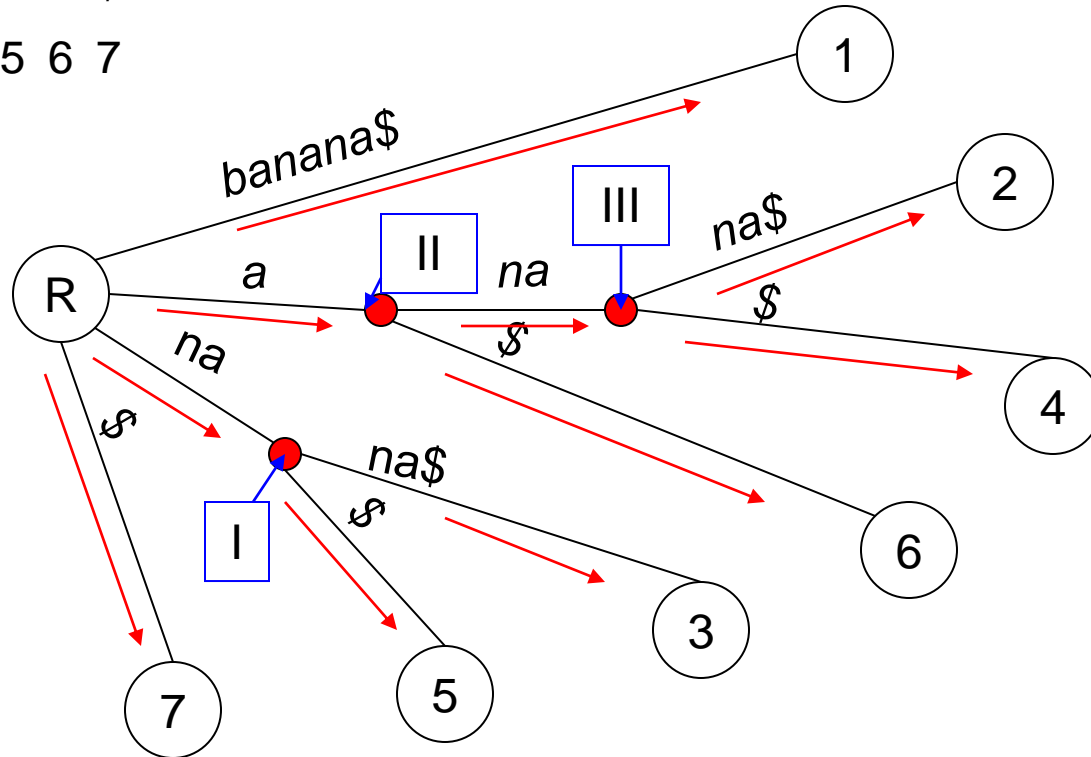b a n a n a $
1 2 3 4 5 6 7

# 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
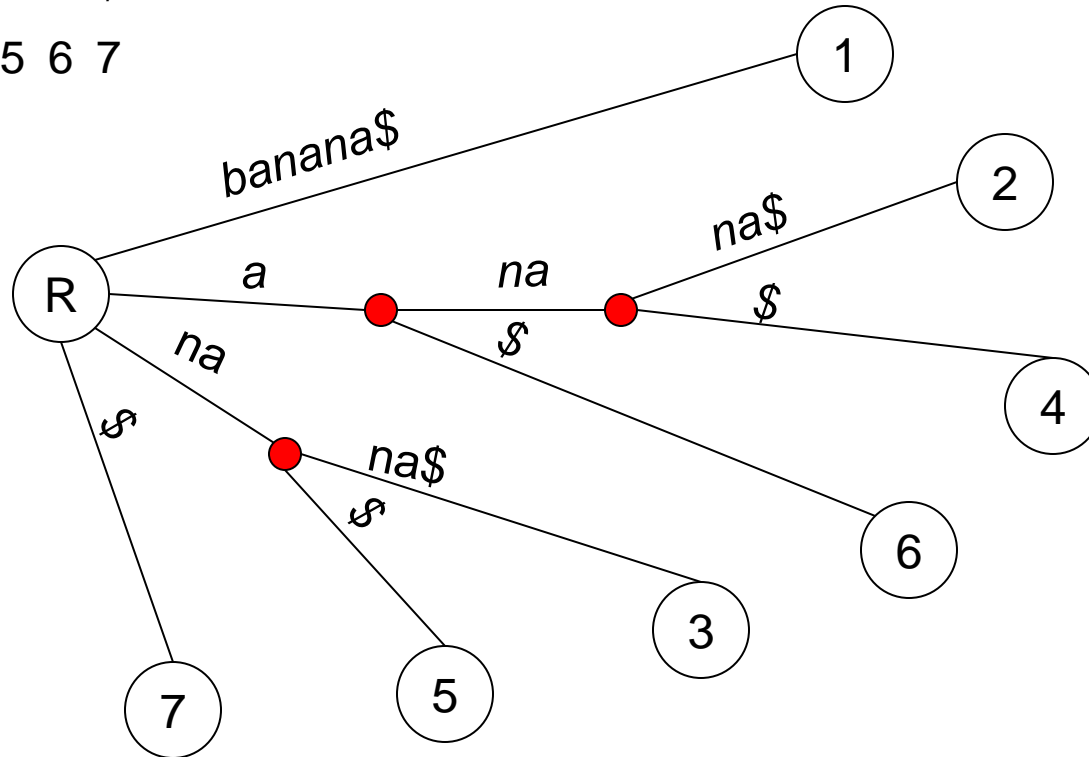
*b a n a n a $*

1 2 3 4 5 6 7



Sequence of nodes visited during traversal:

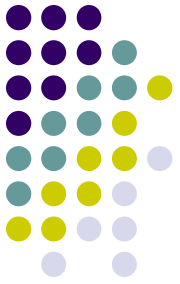R 1 R II III 2 III 4 III II 8 II R I 3 I 5 I R 7 R

# All repetitions

*b a n a n a $*

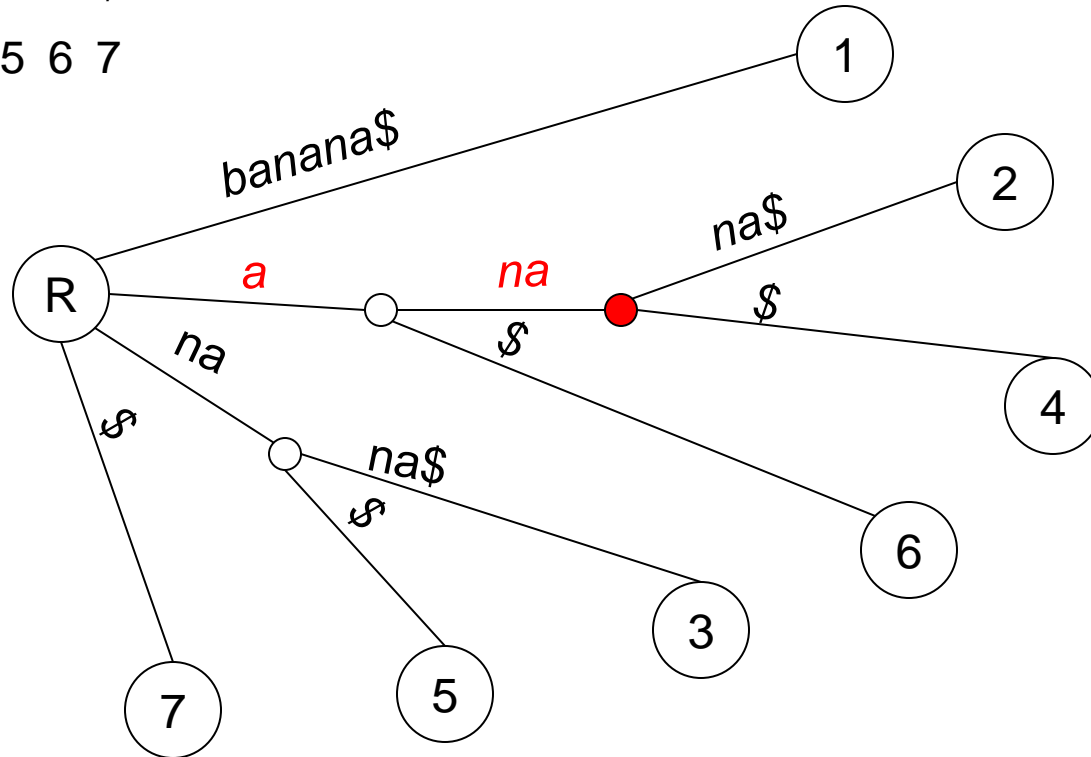1 2 3 4 5 6 7



*n, na; a, an, ana;*
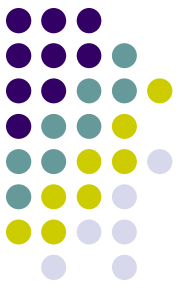
# The longest repeating substring in linear time

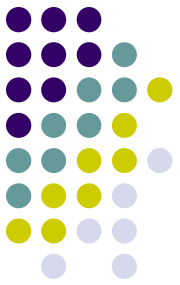*b a n a n a $*

1 2 3 4 5 6 7



*ana*

# 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).
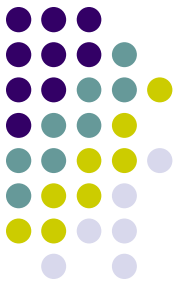
# Maximal repeats example
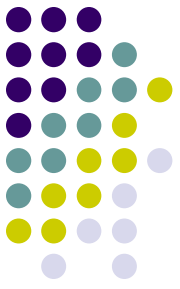
2          10          14
↓          ↓           ↓

- *T=xabcyiiizabcqabcyrxar*

- Which of the following repeated pairs of length 3 are maximal repeats?

  A.  *(2,10)*
  B.  *(2,14)*
  C.  *(10,14)*

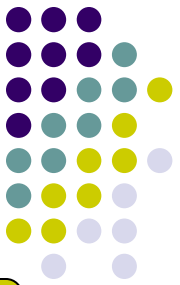# An efficient algorithm for finding all maximal repeats

- The substring labeling the path to any internal node of the suffix tree always represents a right-maximal pair (Why?)

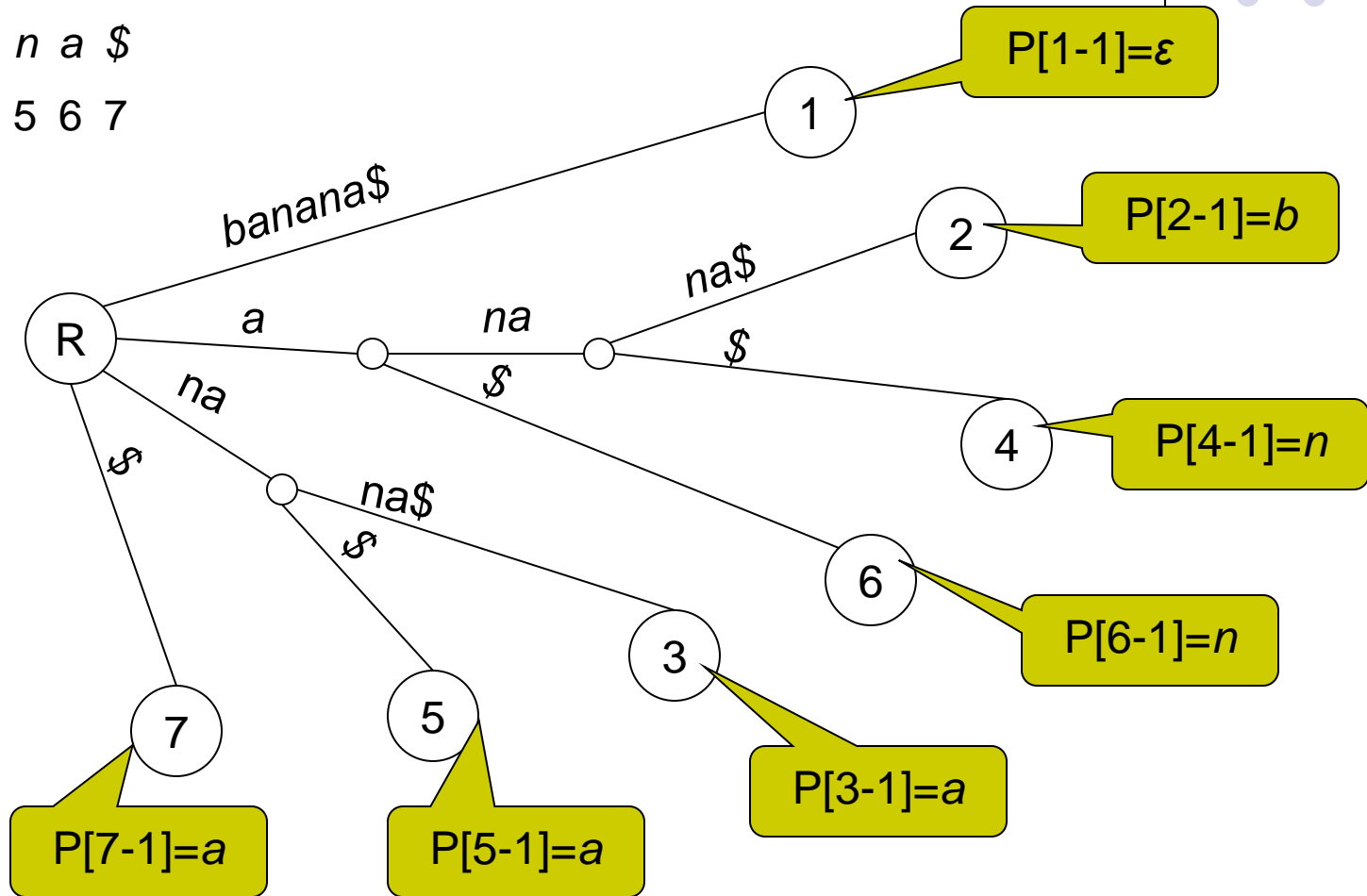# An efficient algorithm for finding all maximal repeats

- 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 if such a substring 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.

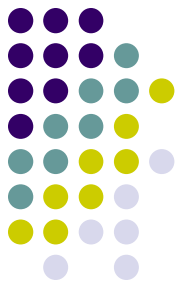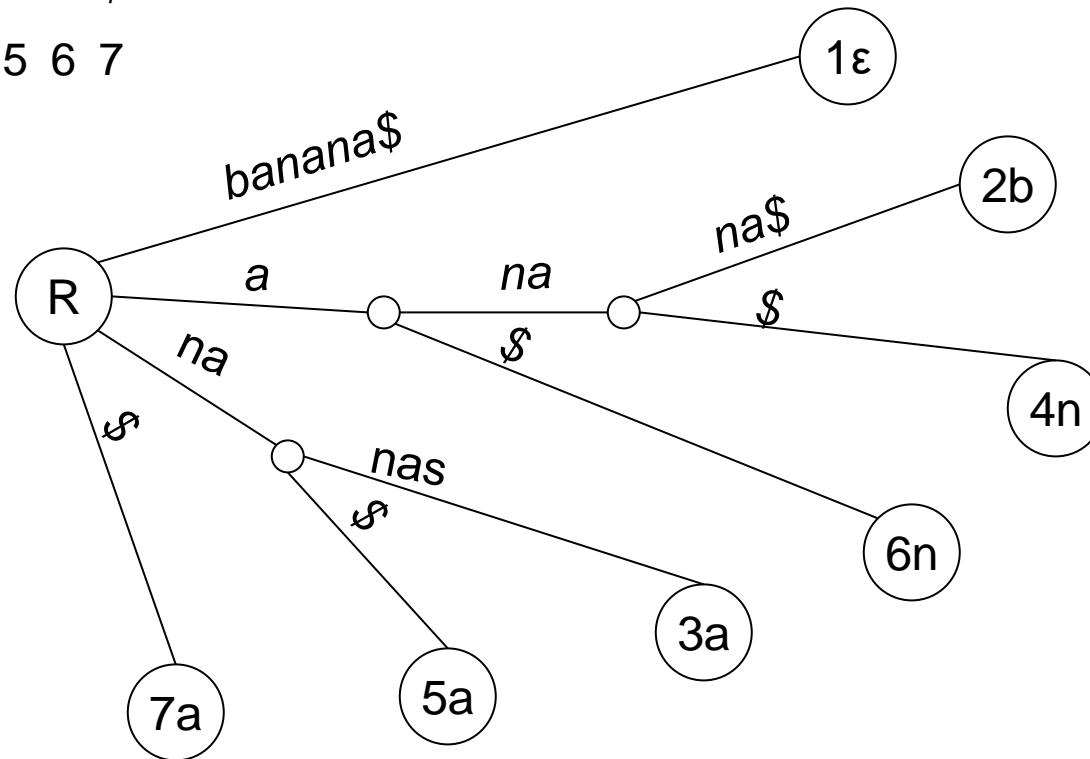# Step 1. Mark leaves with the left character

*b a n a n a $*

1 2 3 4 5 6 7

banana$

a

na

$

na

na$

$

na

$

na$

$

R

1   P[1-1]=ε

2   P[2-1]=b

4   P[4-1]=n

6   P[6-1]=n

3   P[3-1]=a

5   P[5-1]=a

7   P[7-1]=a

# Step 2. Traverse

*b  a  n  a  n  a  $*

1  2  3  4  5  6  7

banana$ — 1ε

a

na — na$ — 2b
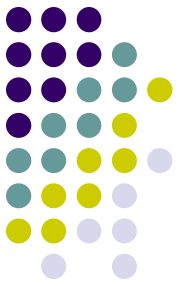
$

$ — 4n

na

$

6n

nas

$

3a

7a

5a

R

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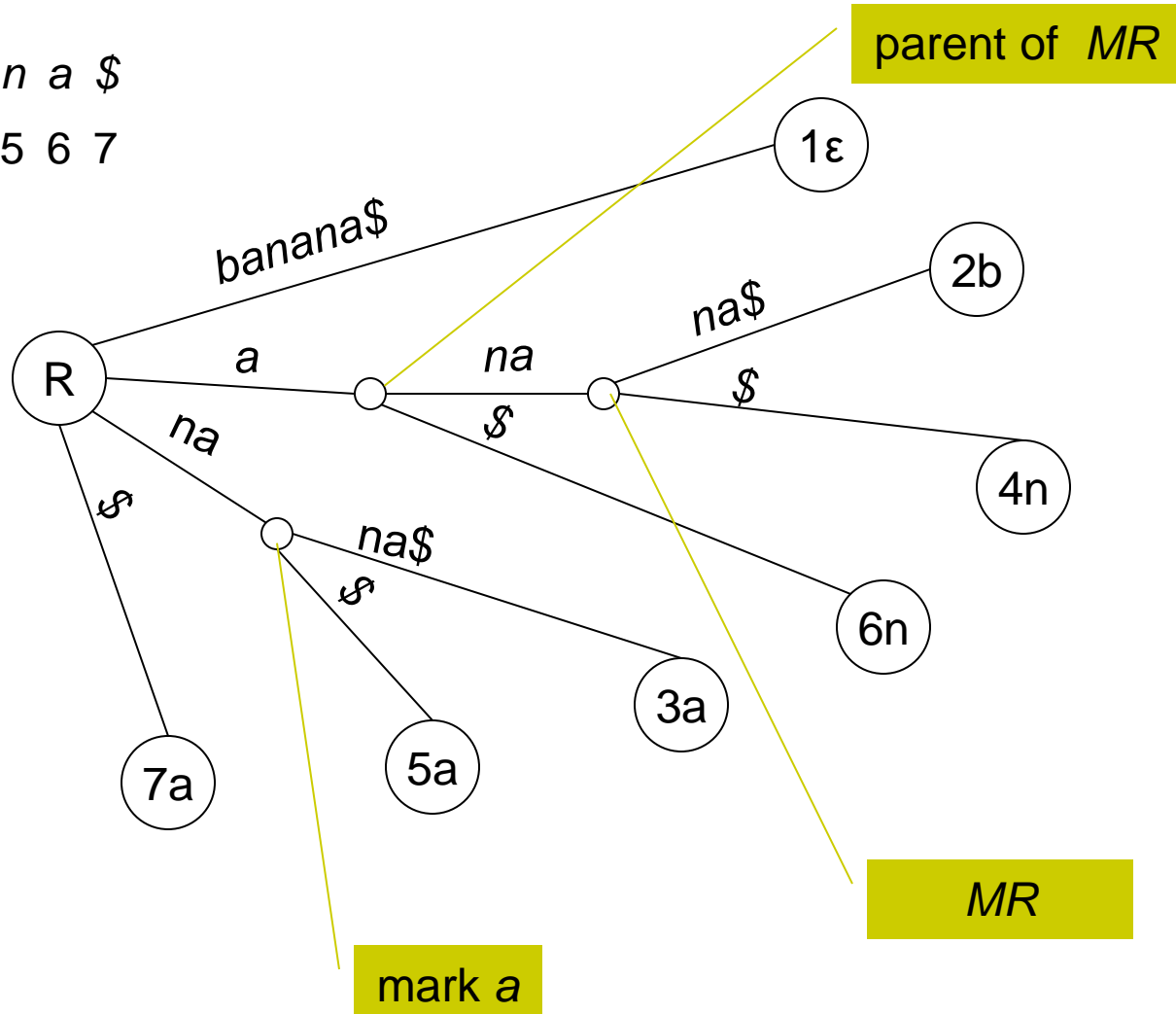
If the left characters are different, then the path from the root to this node represents a maximal repeat, so mark the node as maximal repeat

All parents of MR node are maximal repeats too (do you see why?)
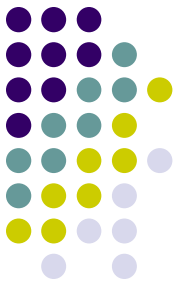
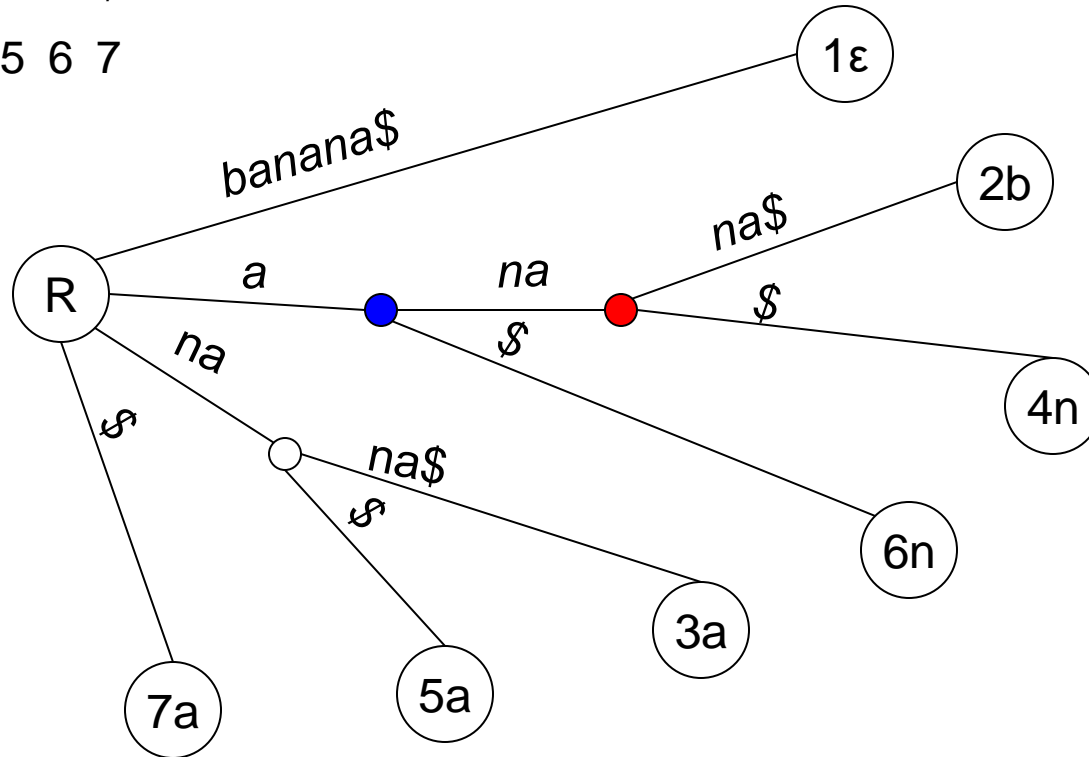# Step 2. Mark internal nodes

*b a n a n a $*

1 2 3 4 5 6 7

parent of *MR*

banana$

1ε

a

na$

2b

na

R

$

na

4n

$

$

6n

na$

$

3a

7a

5a

MR

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

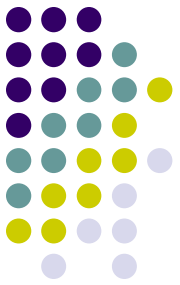mark *a*

# Step 3. Output

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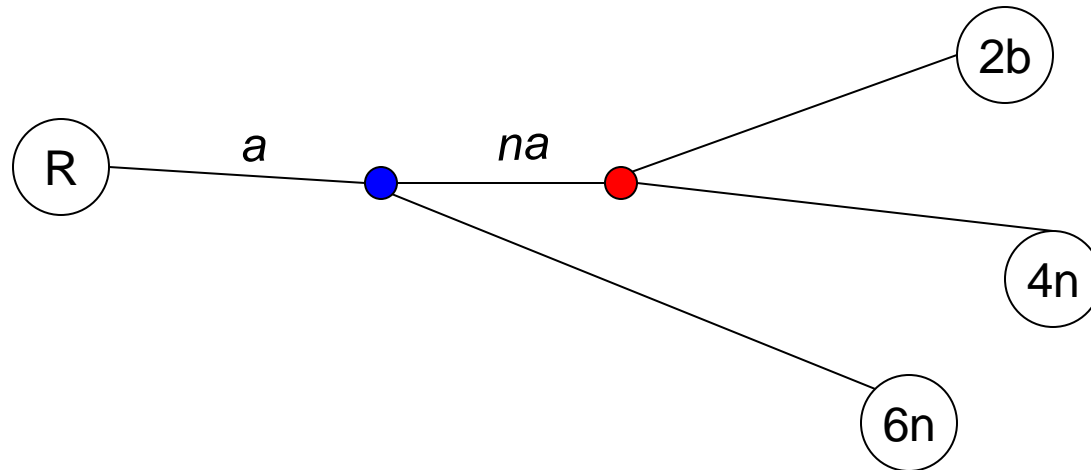
1 2 3 4 5 6 7



Maximal repeat is *ana* (2,4)

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

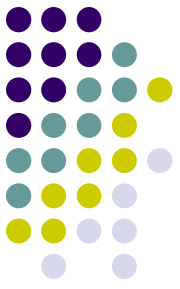# Step 3. Output

*b a n a n a $*

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
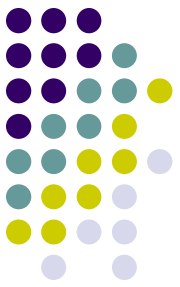
# Significance of repetitions in genome sequences

- Families of reiterated sequences account for about <span style="color:blue">one third</span> of the human genome

- For $3.6*10^6$ nucleotides of the C. Elegans genome, 7000 families of repetitive sequences were discovered

- Mechanism of creating repetitions: error during crossing-over

- Prokaryotes[1] have little repetitive DNA

[1] Prokaryotes (for example, Bacteria) have a circular DNA not enclosed into a nucleus
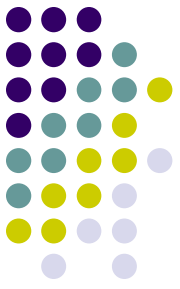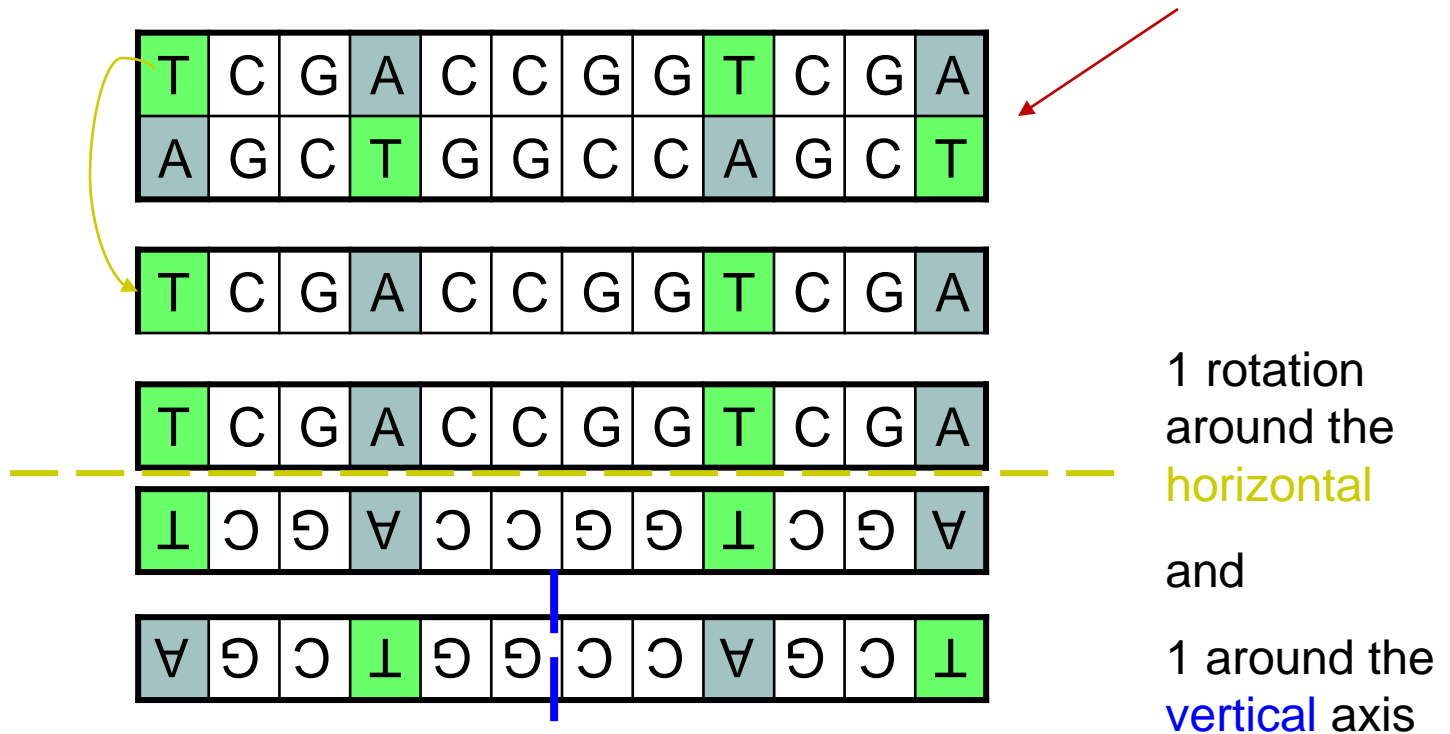
# Repetitions in genome sequences I

- ## <u>Gene families</u>
  - Many genes occur in multiple copies
  - 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 main mechanism of evolution

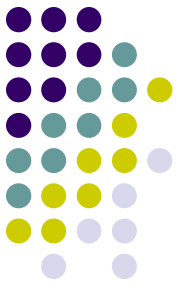*R-RNA is an RNA component of the ribosomes
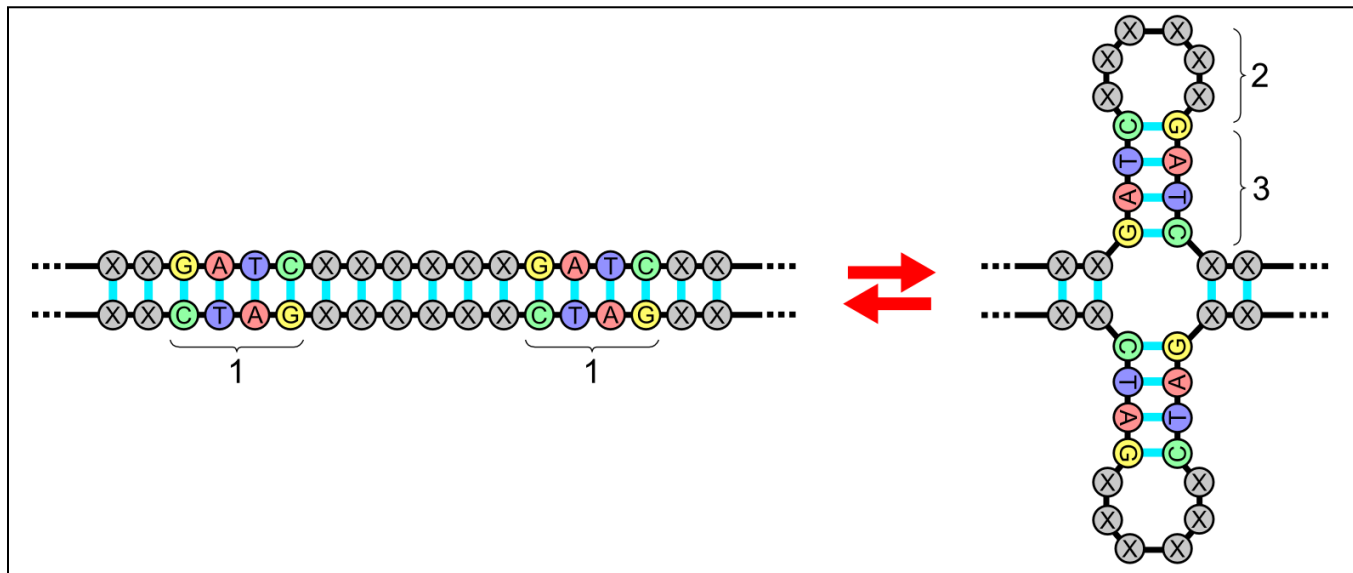
# Repetitions in genome sequences II

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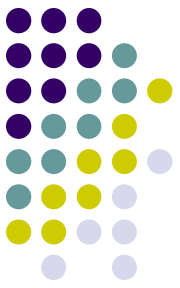- They often have a form of a *complemented palindrome*

| T | C | G | A | C | C | G | G | T | C | G | A |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | G | C | T | G | G | C | C | A | G | C | T |

| T | C | G | A | C | C | G | G | T | C | G | A |
|---|---|---|---|---|---|---|---|---|---|---|---|

| T | C | G | A | C | C | G | G | T | C | G | A |
|---|---|---|---|---|---|---|---|---|---|---|---|

1 rotation around the **horizontal**

and

1 around the **vertical** axis

# 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 aminoacids.
- 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*

<u>Microsatellites</u> – 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)
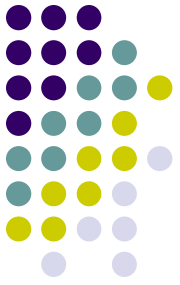


Ribose



Deoxyribose



Phosphoric acid

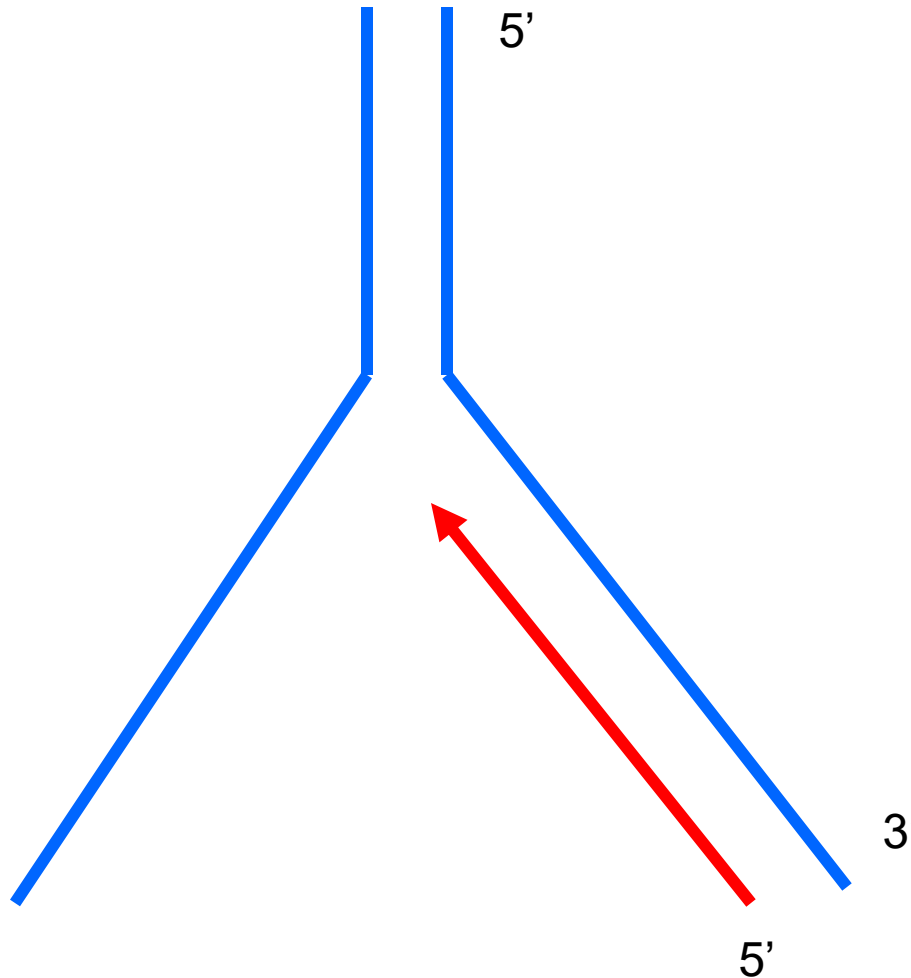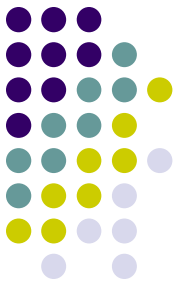# Chemistry of DNA



Deoxyribose

# Nucleotide – building block of the DNA polymer

Phosphate

Base (A, C, G or T

Sugar ring

P

B

S

# Chaining nucleotides

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

# DNA Replication

5'

3'

5'

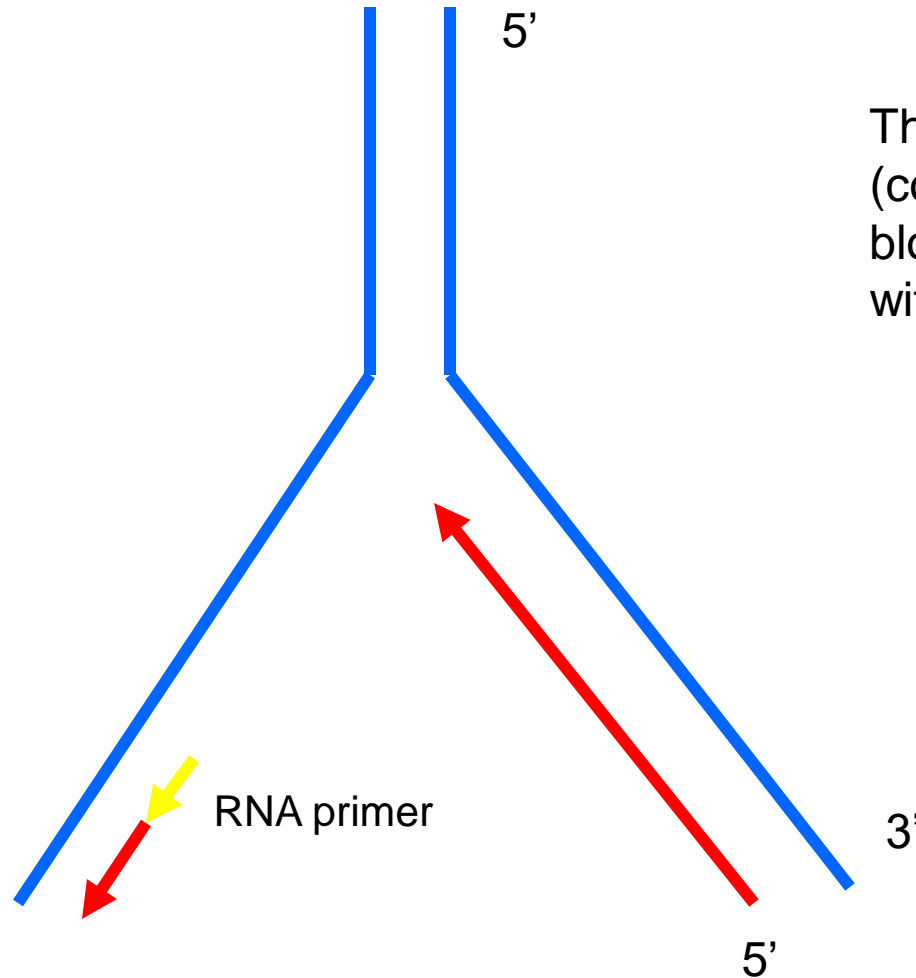The leading strand is replicating without problems

# DNA Replication

5'

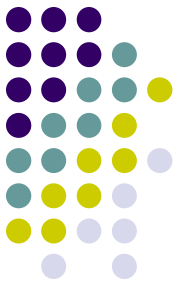The lagging strand cannot even start, since there is no free complementary 3'-end
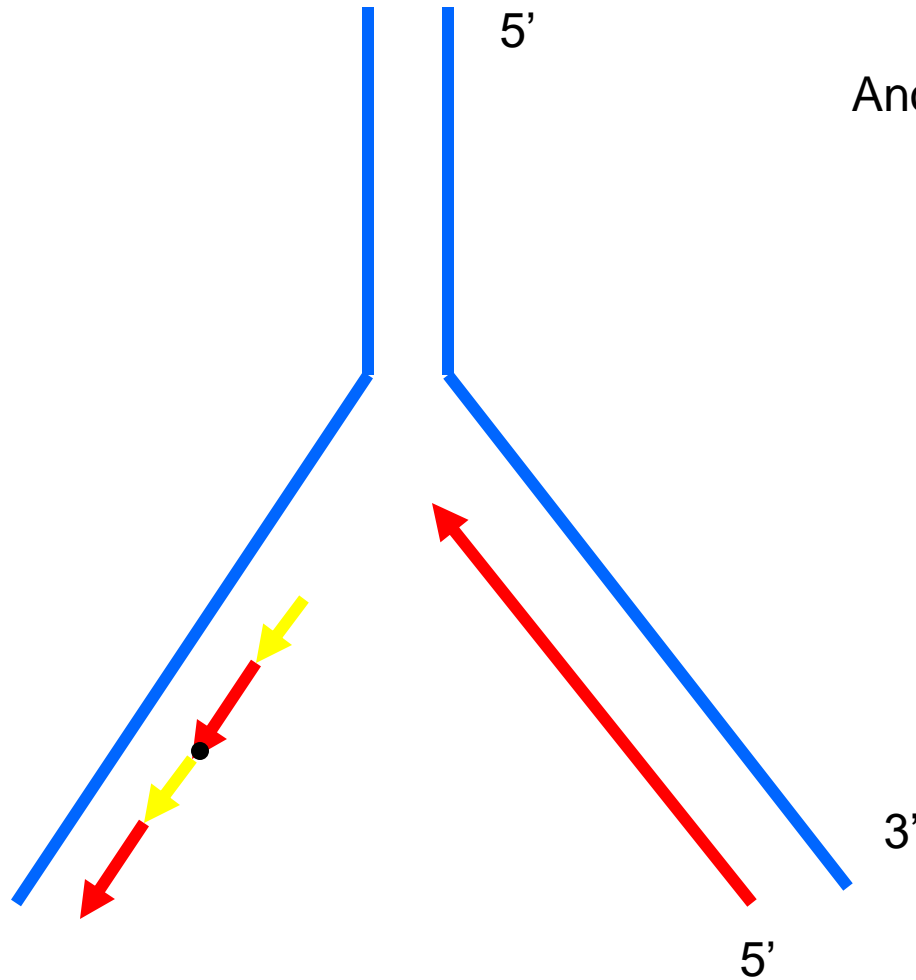
It starts by creating an RNA primer

RNA primer

3'

5'

# DNA Replication
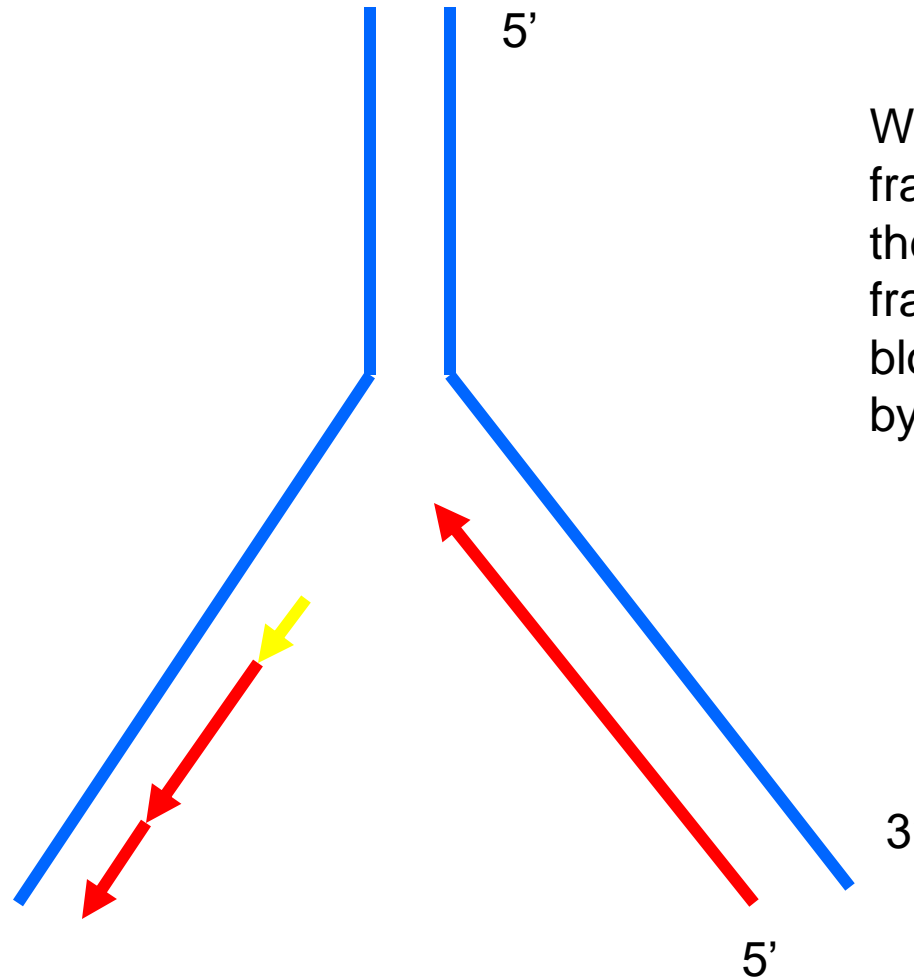
5'

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

RNA primer
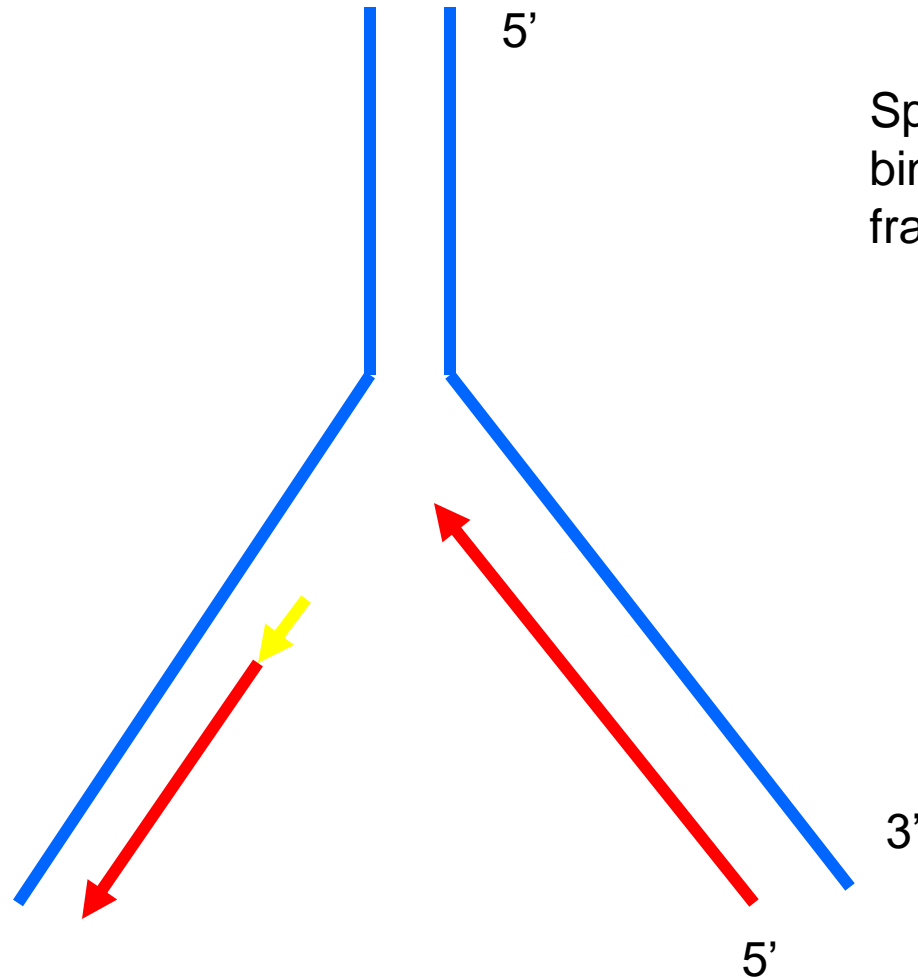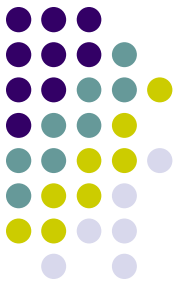
3'

5'

# DNA Replication



5'

Another fragment
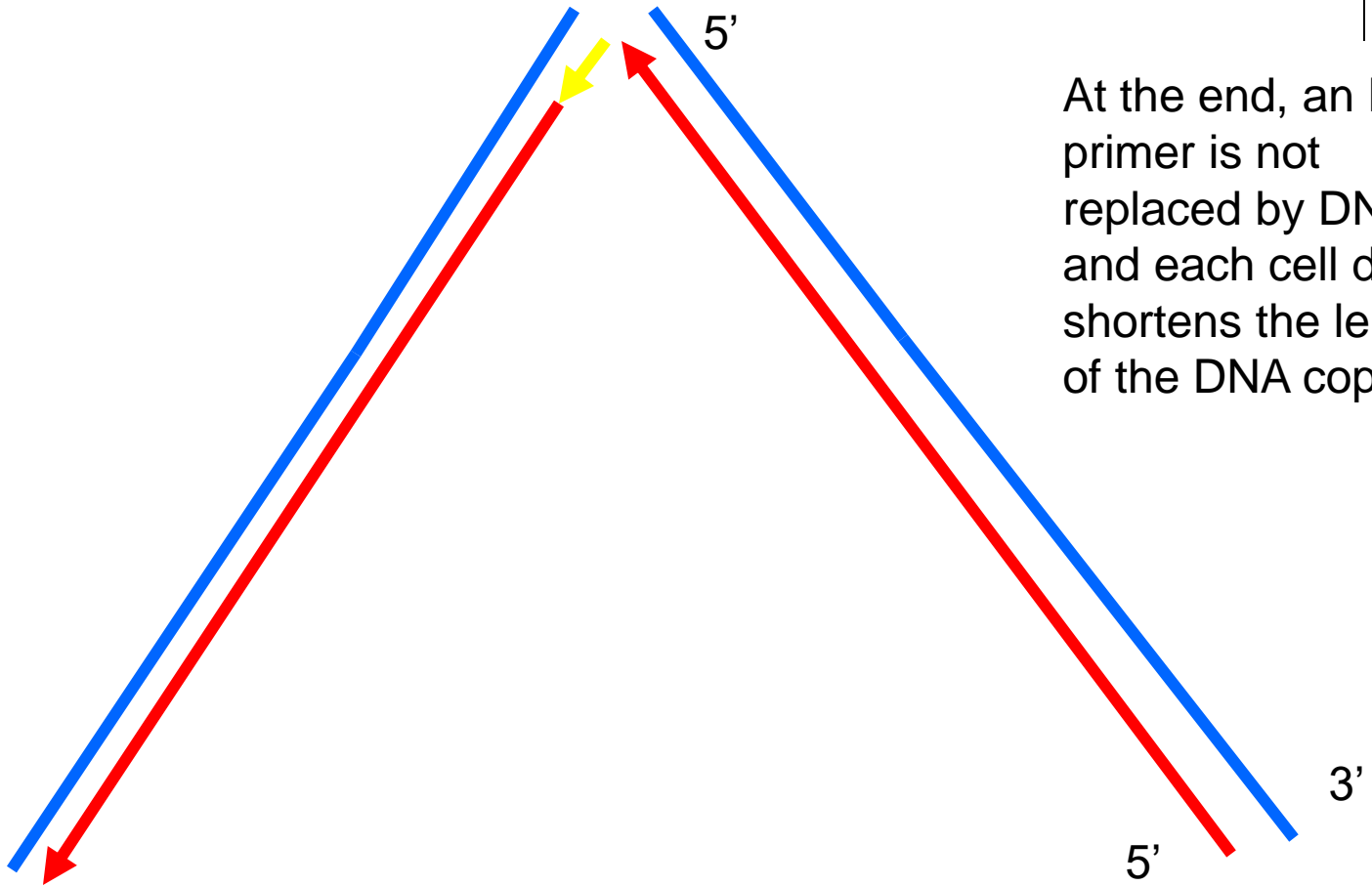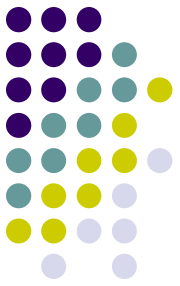
3'

5'

# DNA Replication

5'

When the second fragment reaches the start of the first fragment, RNA blocks are replaced by DNA blocks
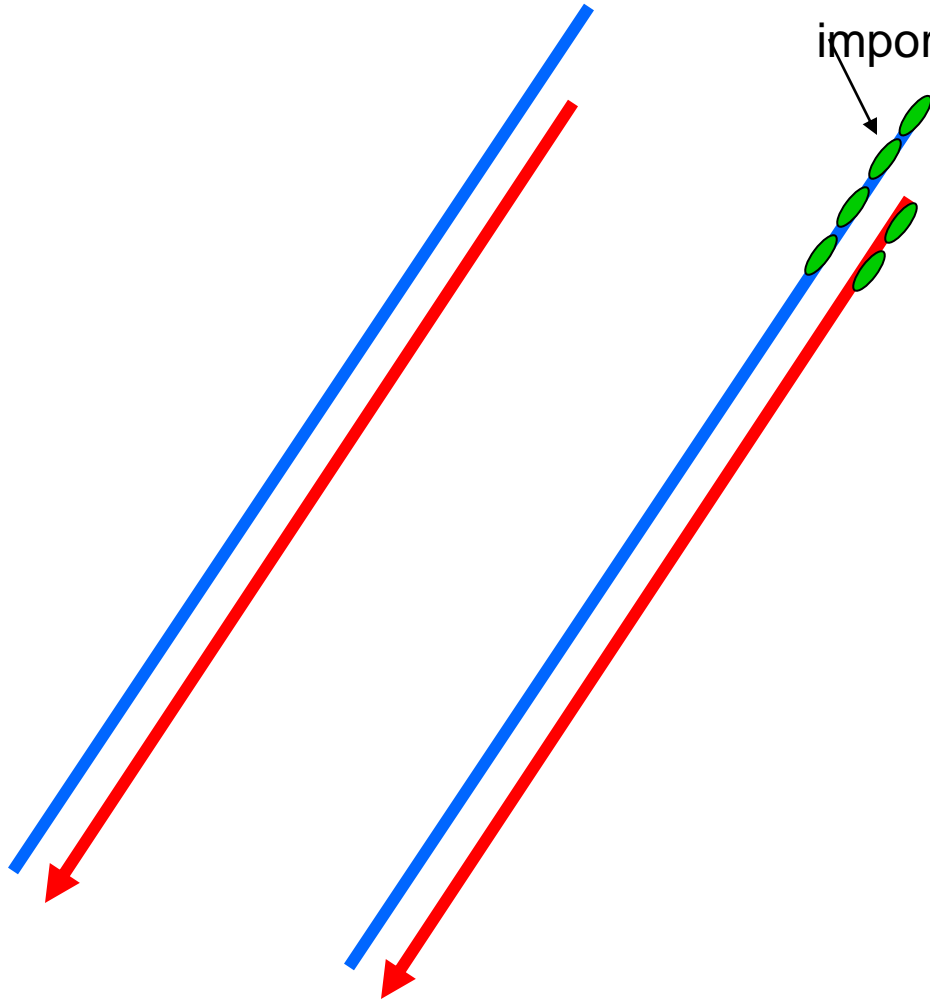
3'

5'

# **DNA Replication**

5'

Special enzyme
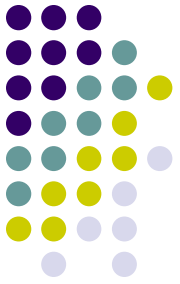binds the DNA
fragments together

3'

5'

# DNA Replication

5'

At the end, an RNA primer is not replaced by DNA, and each cell division shortens the length of the DNA copy
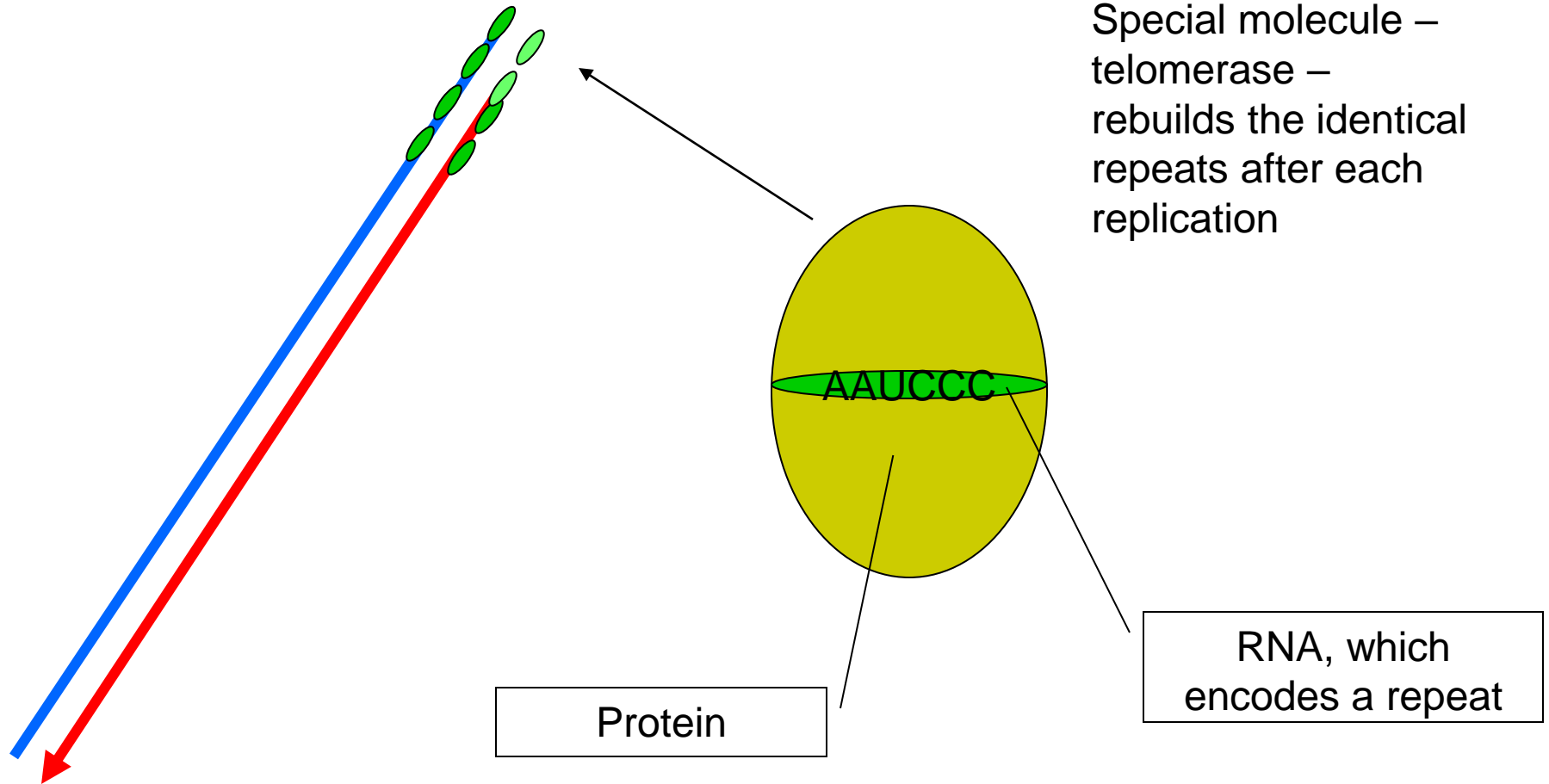
3'

5'

# Telomerase
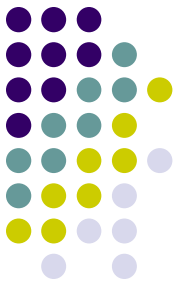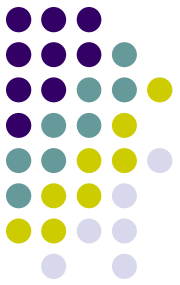
Not something important

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

AAUCCC
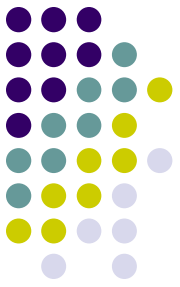
Protein

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 efficient algorithm for finding repeats

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