

- **Modern Information Retrieval (1999)**
Ricardo-Baeza Yates and Berthier Ribeiro-Neto
- **Flexible Pattern Matching in Strings (2002)**
Gonzalo Navarro and Mathieu Raffinot
- **Algorithms on strings (2001)**
M. Crochemore, C. Hancart and T. Lecroq
- <http://www-igm.univ-mlv.fr/~lecroq/string/index.html>

String matching: definition of the problem (text,pattern)

- **Exact matching:** depends on what we have: text or patterns
 - **The patterns** ---> Data structures for the patterns
 - 1 pattern ---> The algorithm depends on $|p|$ and $|\Sigma|$
 - k patterns ---> The algorithm depends on k, $|p|$ and $|\Sigma|$
 - Extensions
 - Regular Expressions
 - **The text** ----> Data structure for the text (suffix tree, ...)
- **Approximate matching:**
 - Dynamic programming
 - Sequence alignment (pairwise and multiple)
 - Sequence assembly: hash algorithm
- **Probabilistic search:** Hidden Markov Models

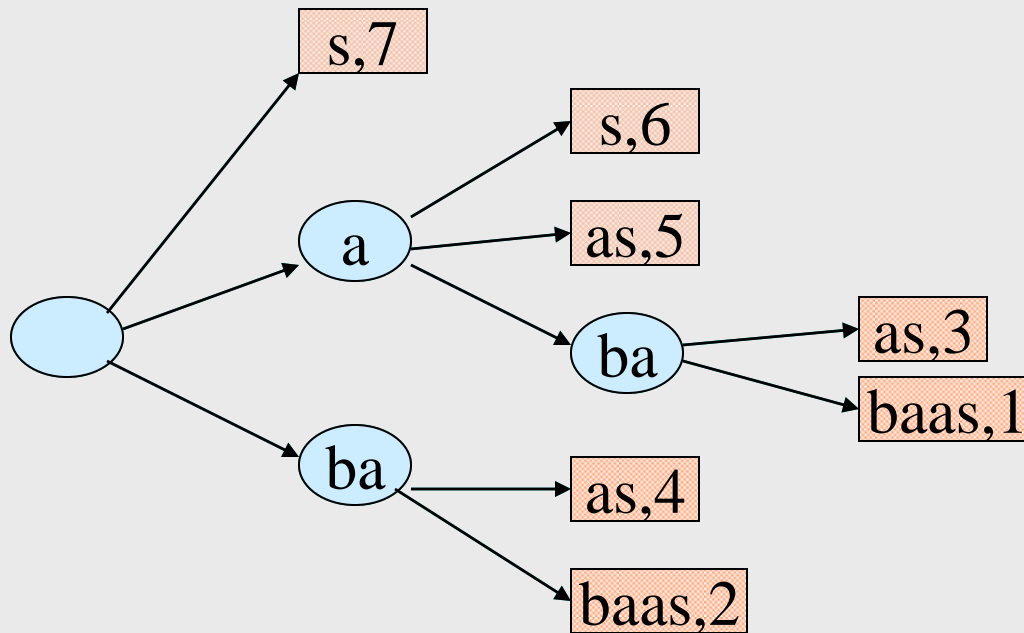
1a. Part: Suffix trees

Algorithms on strings, trees and sequences,
Dan Gusfield
Cambridge University Press

2a. Part: Suffix arrays

Suffix-arrays: a new method for on-line
string searches,
G. Myers, U. Manber

Given string **ababaas**:



Suffixes:

7: s

6: as

5: aas

3: abaas

1: ababaas

4: baas

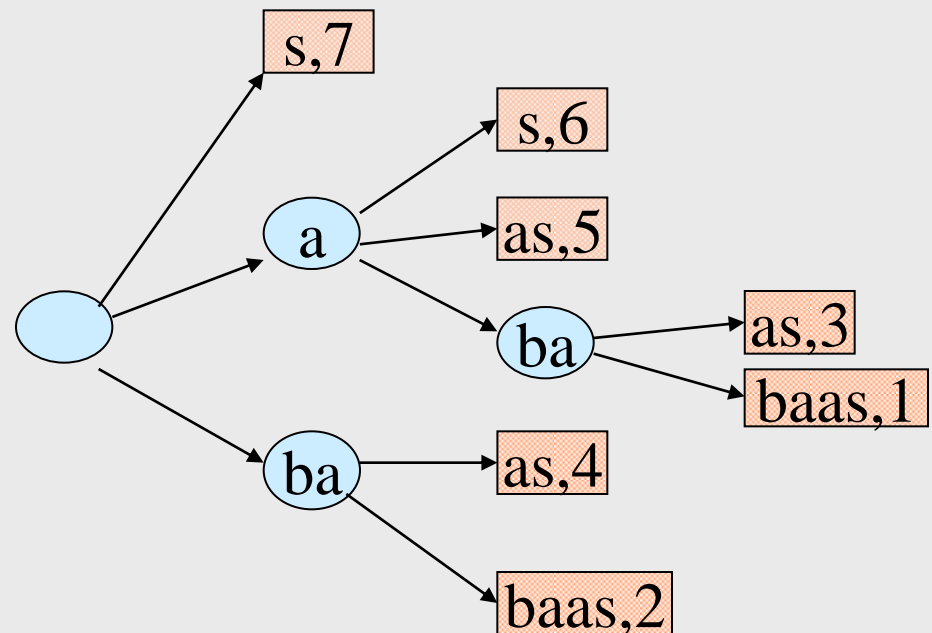
2: babaas

What kind of queries?

Applications of Suffix trees

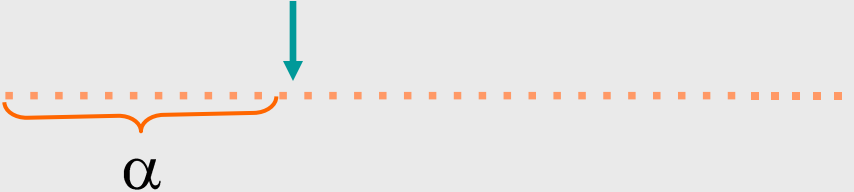
1. Exact string matching

- Does the sequence **ababaas** contain any occurrence of patterns **abab**, **aab**, and **ab**?

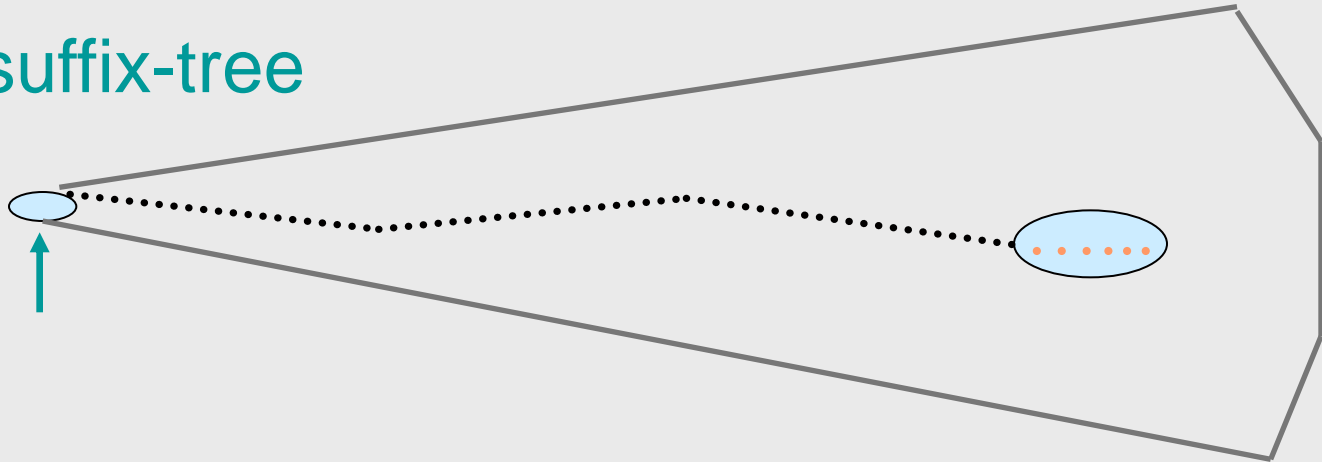


Quadratic insertion algorithm

Invariant Properties:

Given the string 

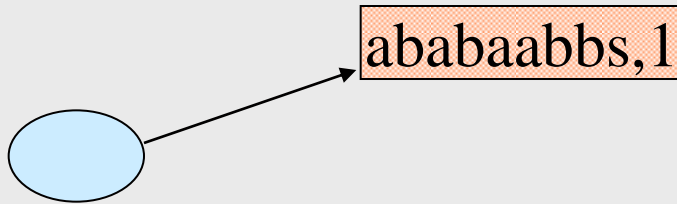
and the suffix-tree



P1: the leaves of suffixes from α have been inserted

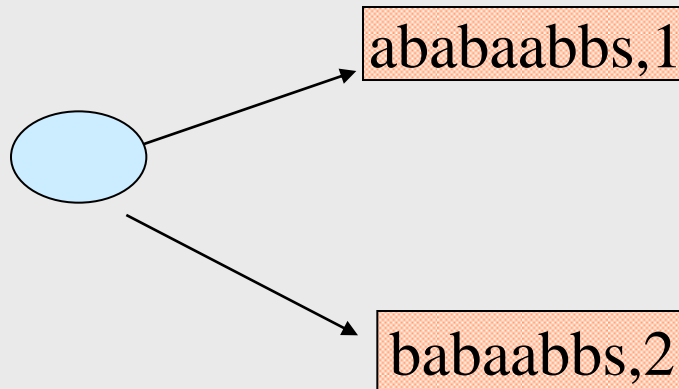
Quadratic insertion algorithm

Given the string ababaabbs



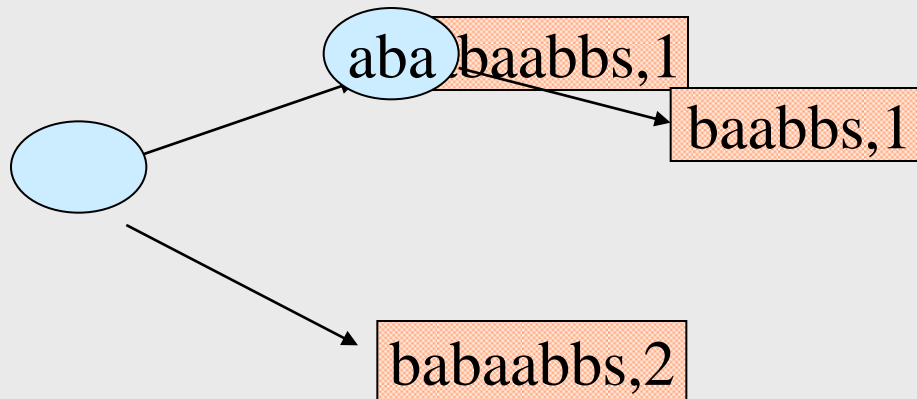
Quadratic insertion algorithm

Given the string ababaabbs



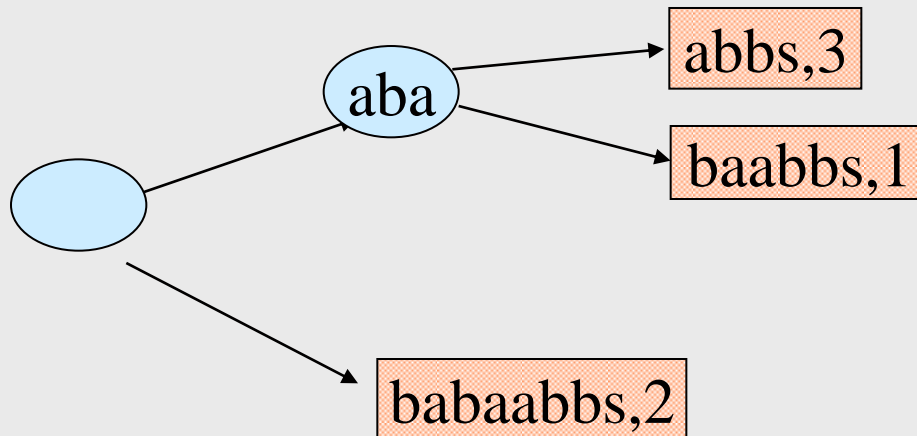
Quadratic insertion algorithm

Given the string **ababaabbs**



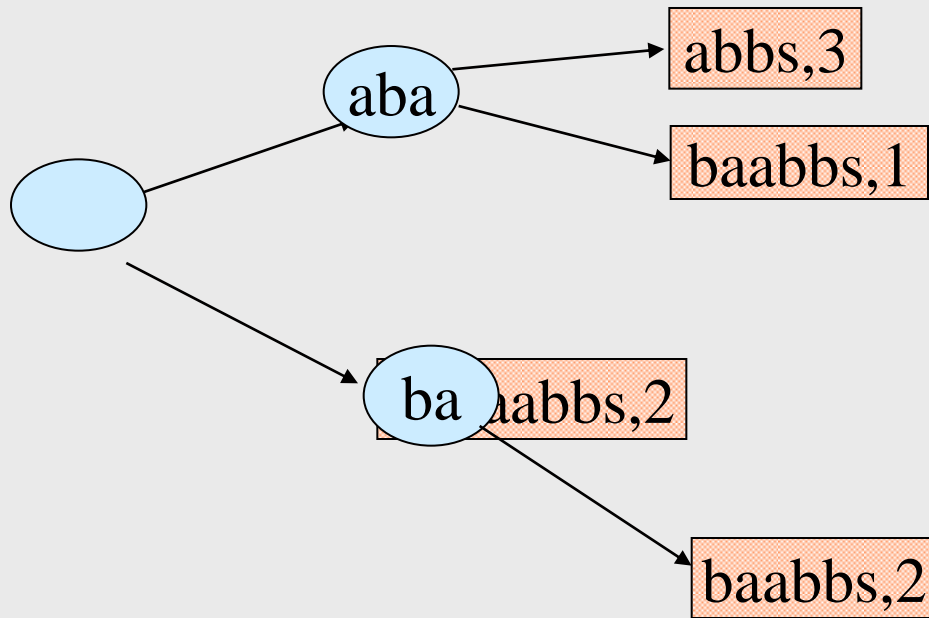
Quadratic insertion algorithm

Given the string **ababaabbs**



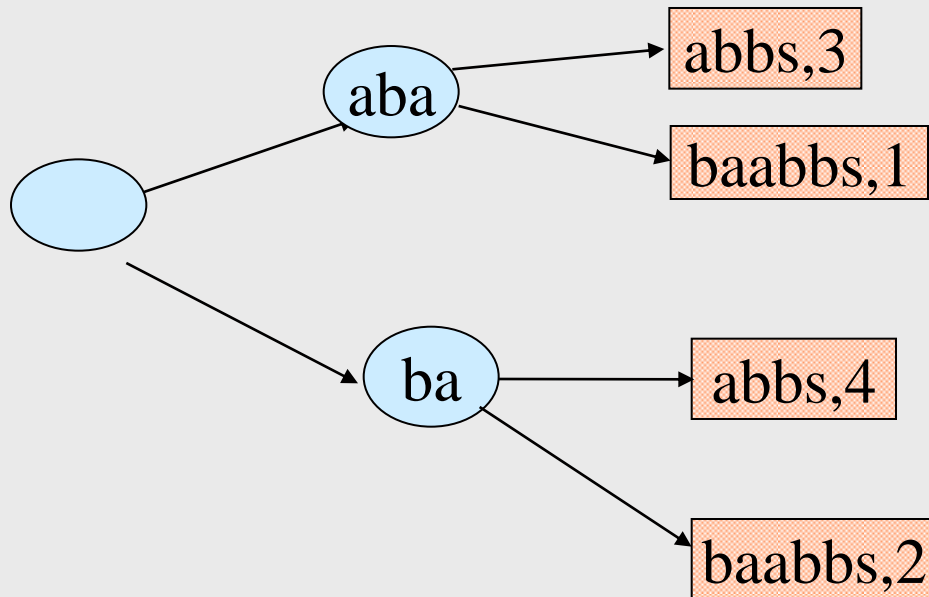
Quadratic insertion algorithm

Given the string **ababaabbs**



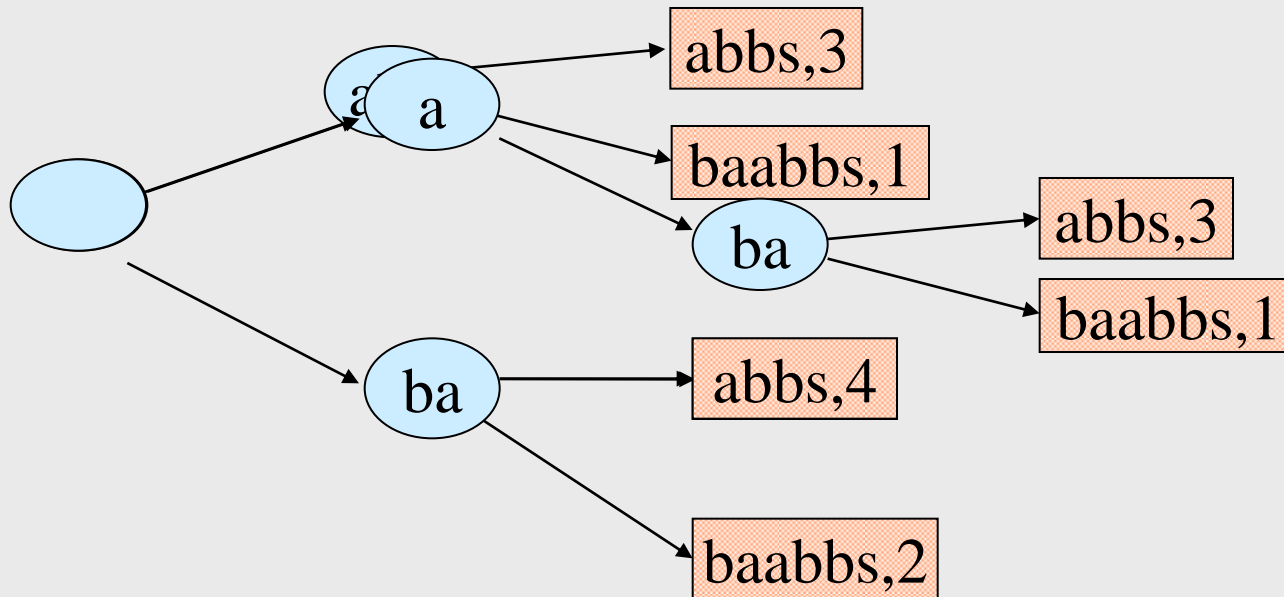
Quadratic insertion algorithm

Given the string **ababaabbs**



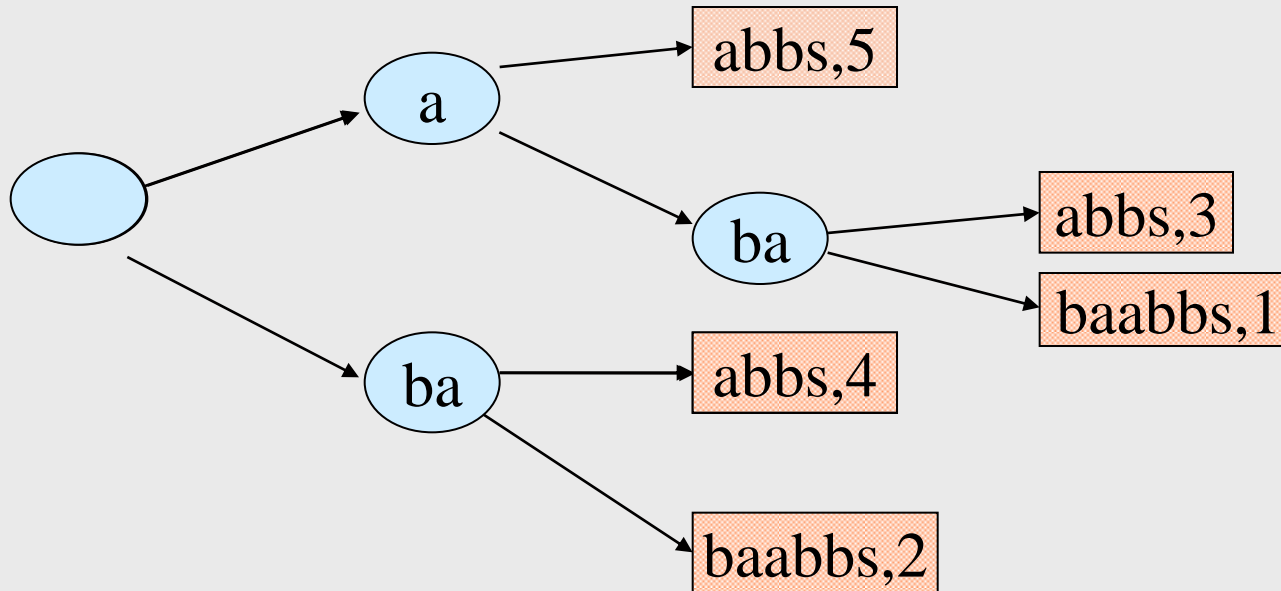
Quadratic insertion algorithm

Given the string **ababaabbs**



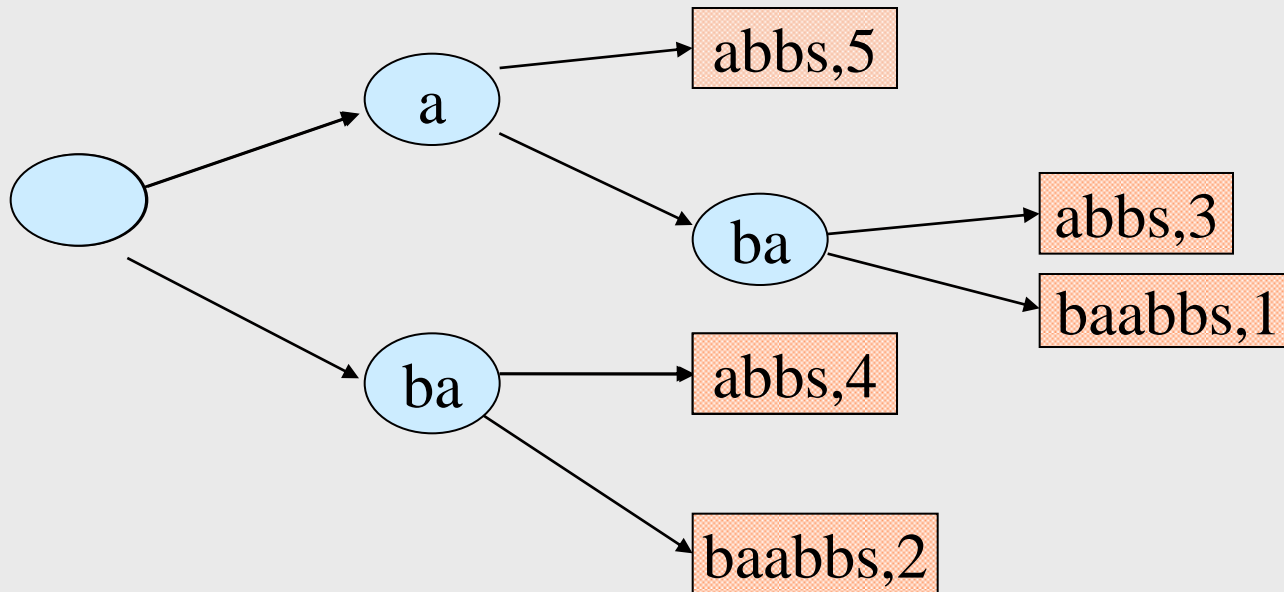
Quadratic insertion algorithm

Given the string **ababaabbs**



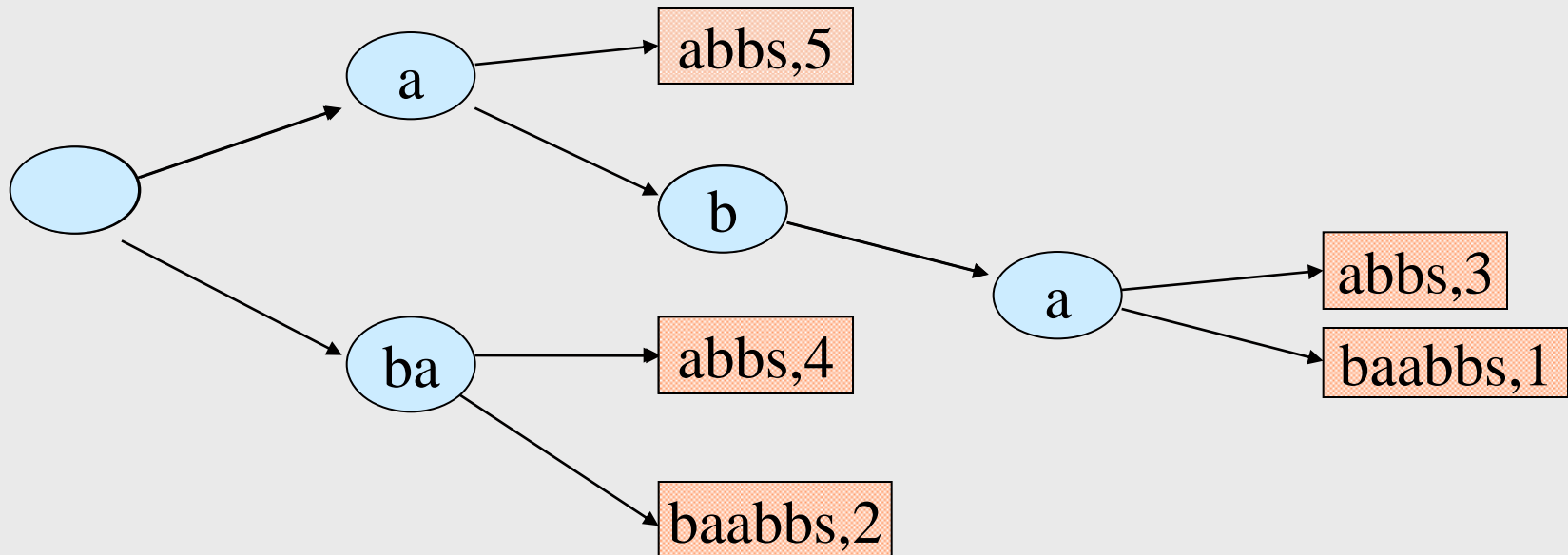
Quadratic insertion algorithm

Given the string **ababaabbs**



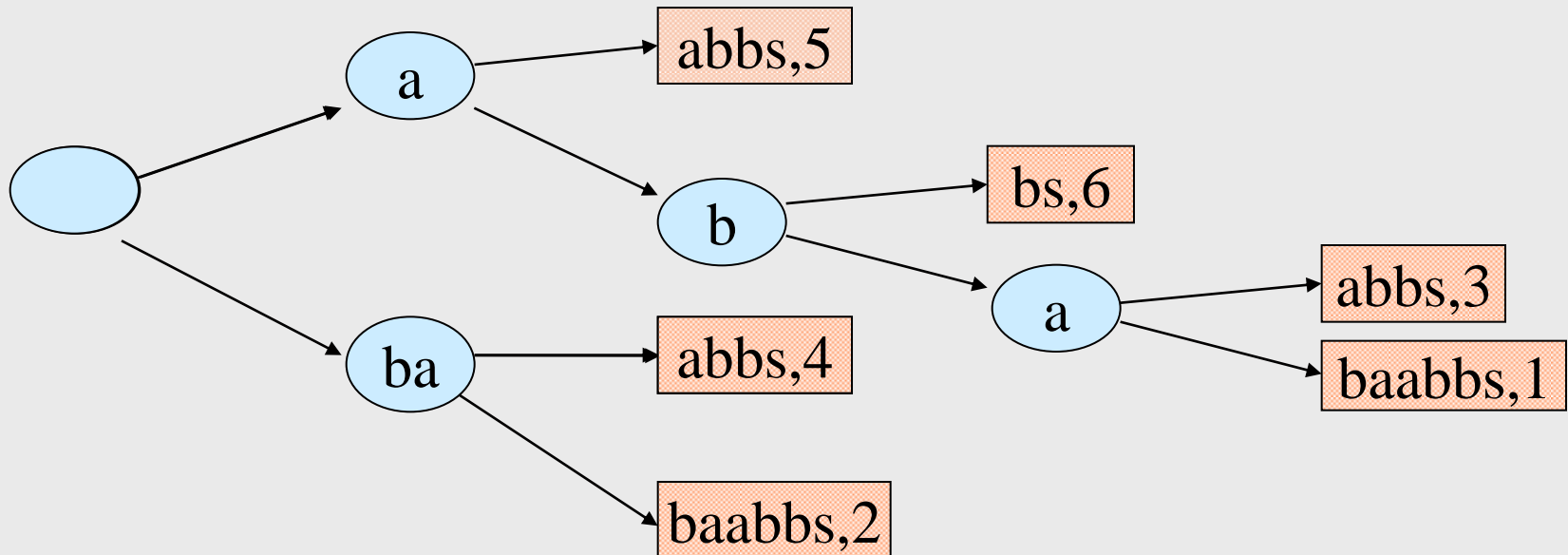
Quadratic insertion algorithm

Given the string **ababaabbs**



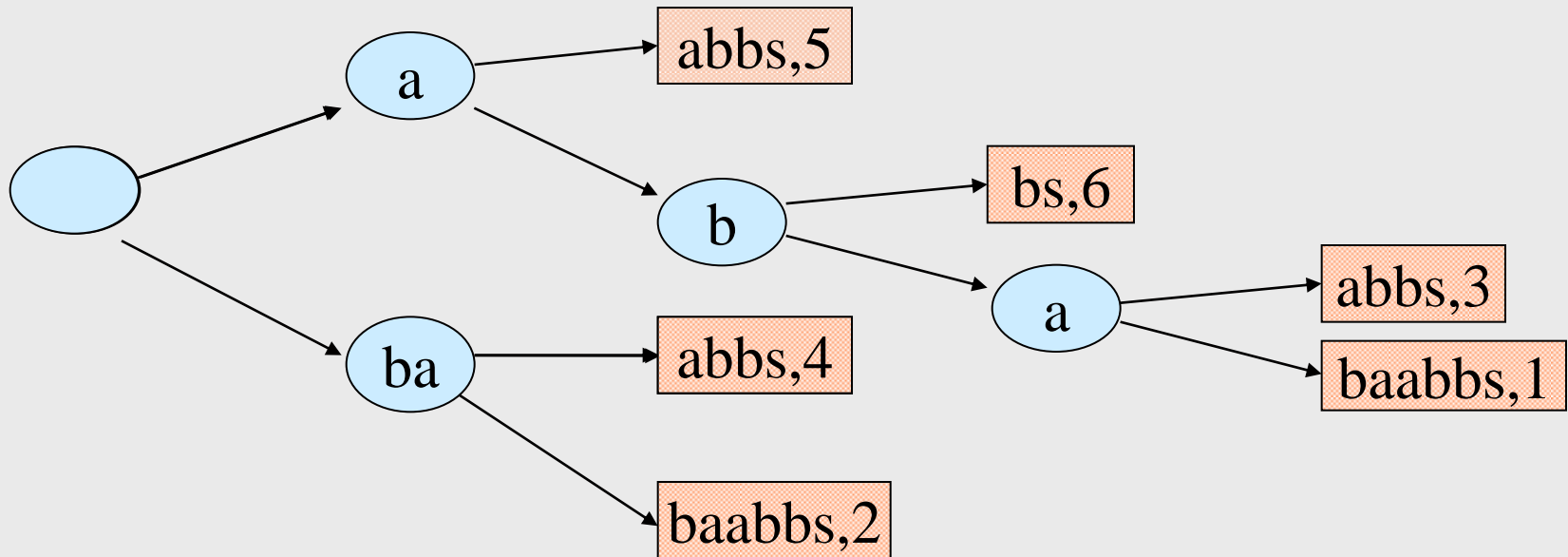
Quadratic insertion algorithm

Given the string **ababaabbs**



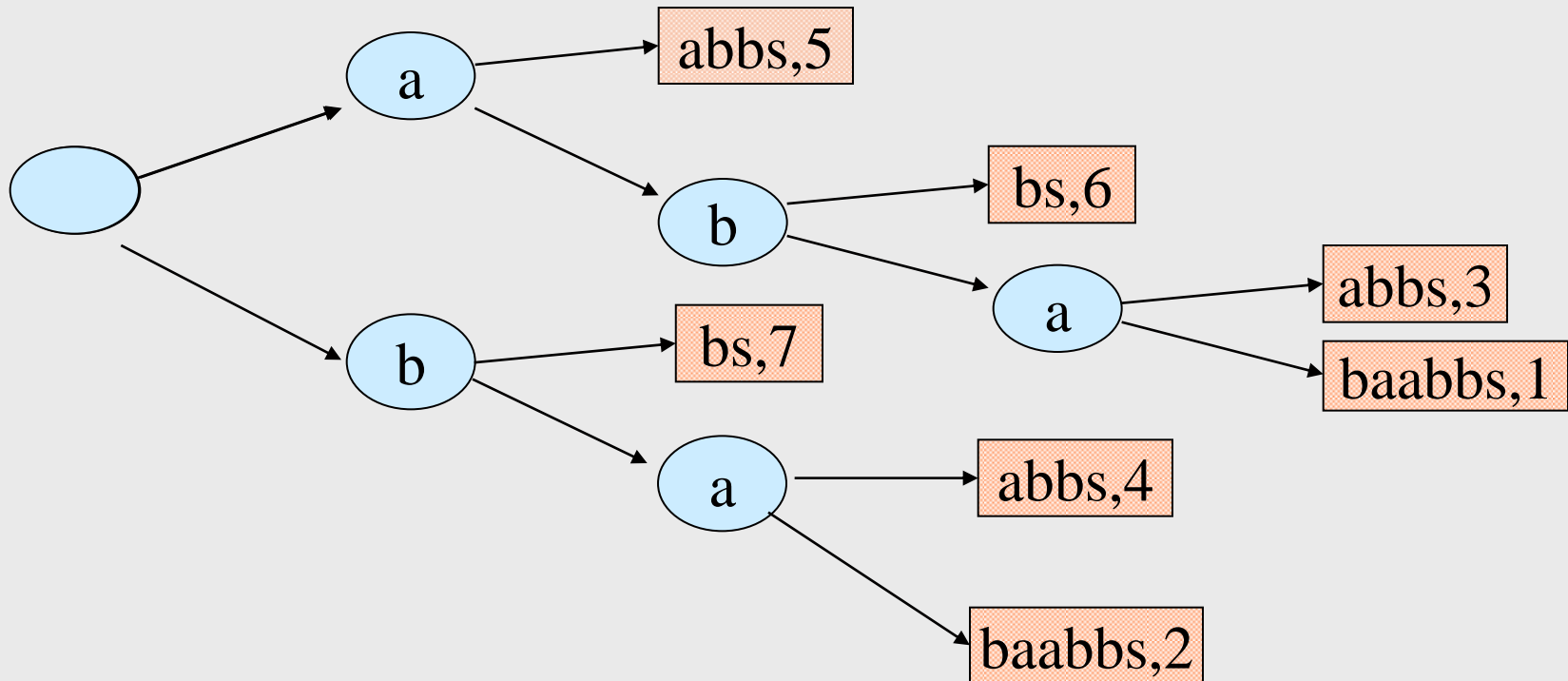
Quadratic insertion algorithm

Given the string **ababaabbs**



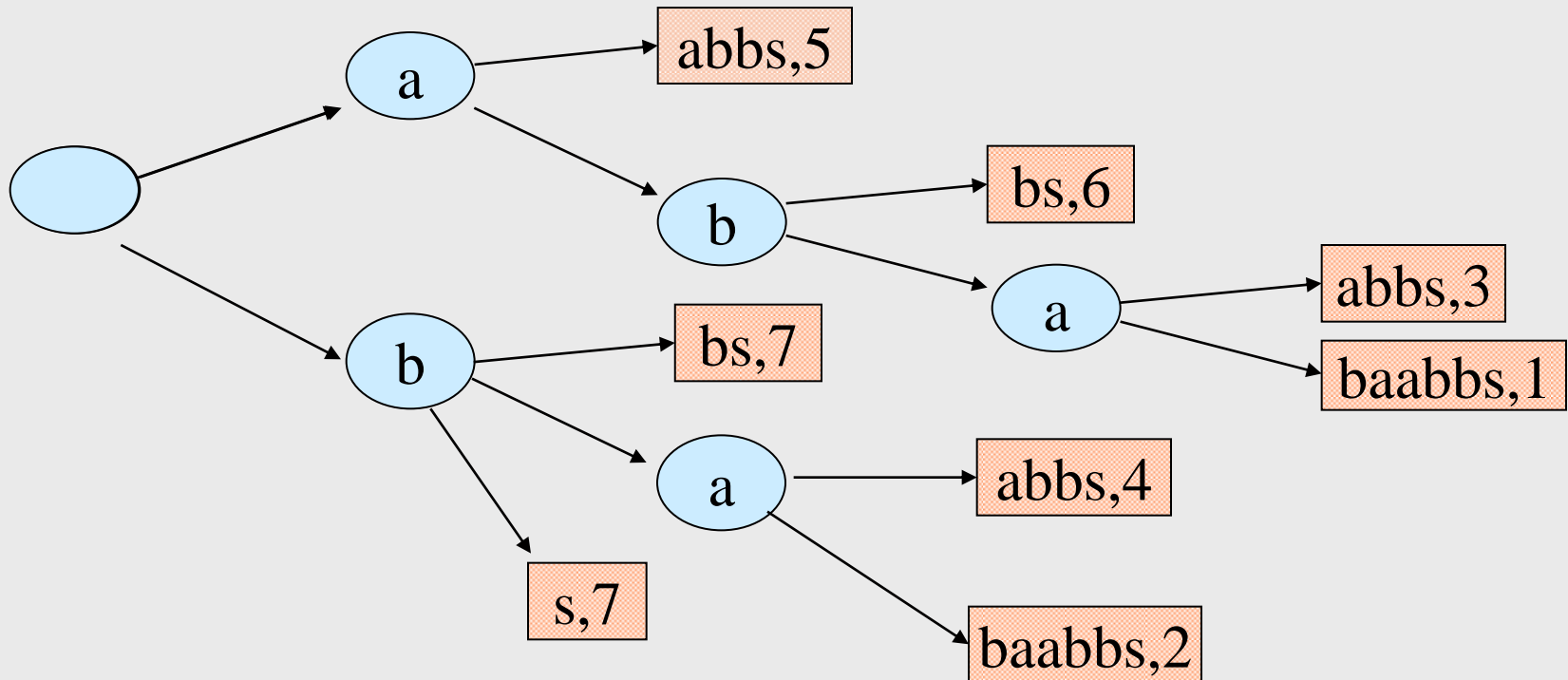
Quadratic insertion algorithm

Given the string **ababaabbs**



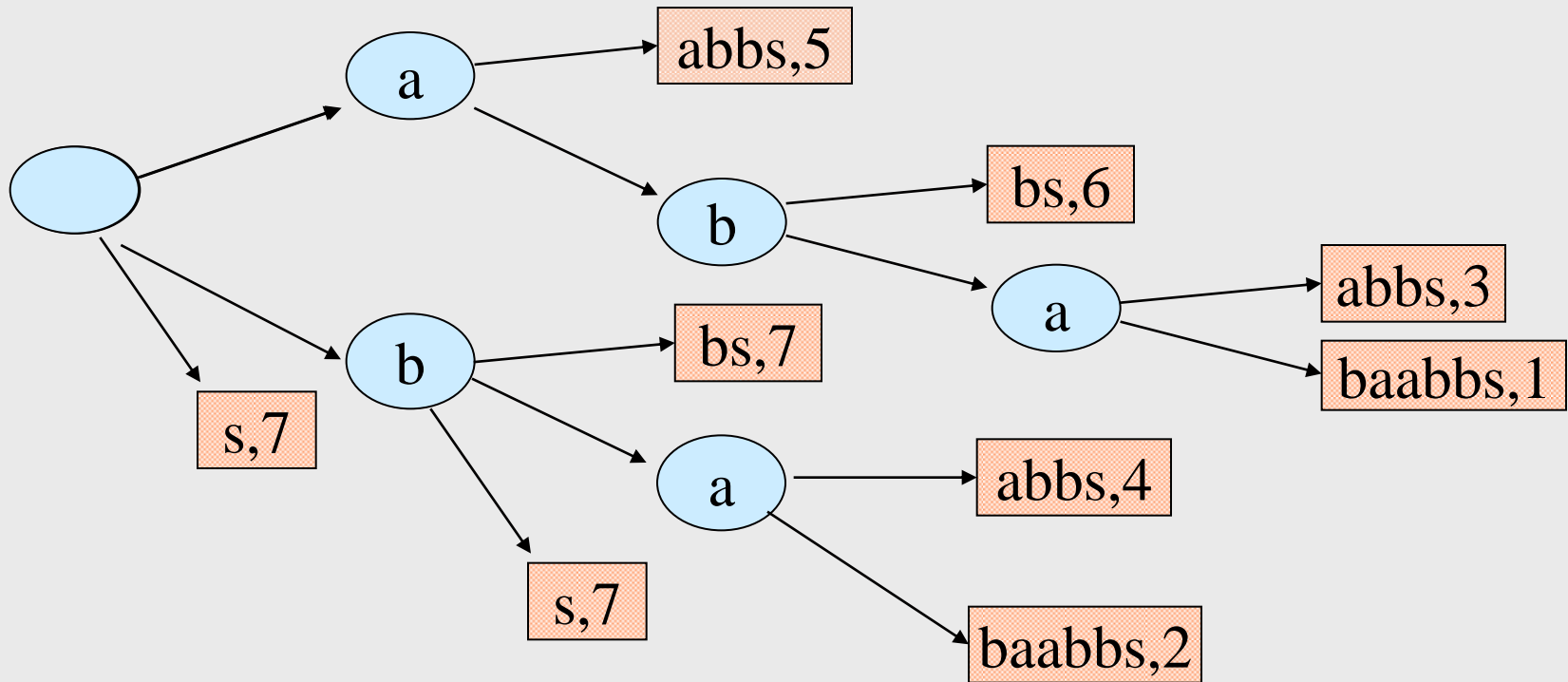
Quadratic insertion algorithm

Given the string **ababaabbs**



Quadratic insertion algorithm

Given the string **ababaabbs**



The suffix tree of many strings ...

is called the generalized suffix tree ...

and it is the suffix tree of the concatenation
of strings.

For instance,

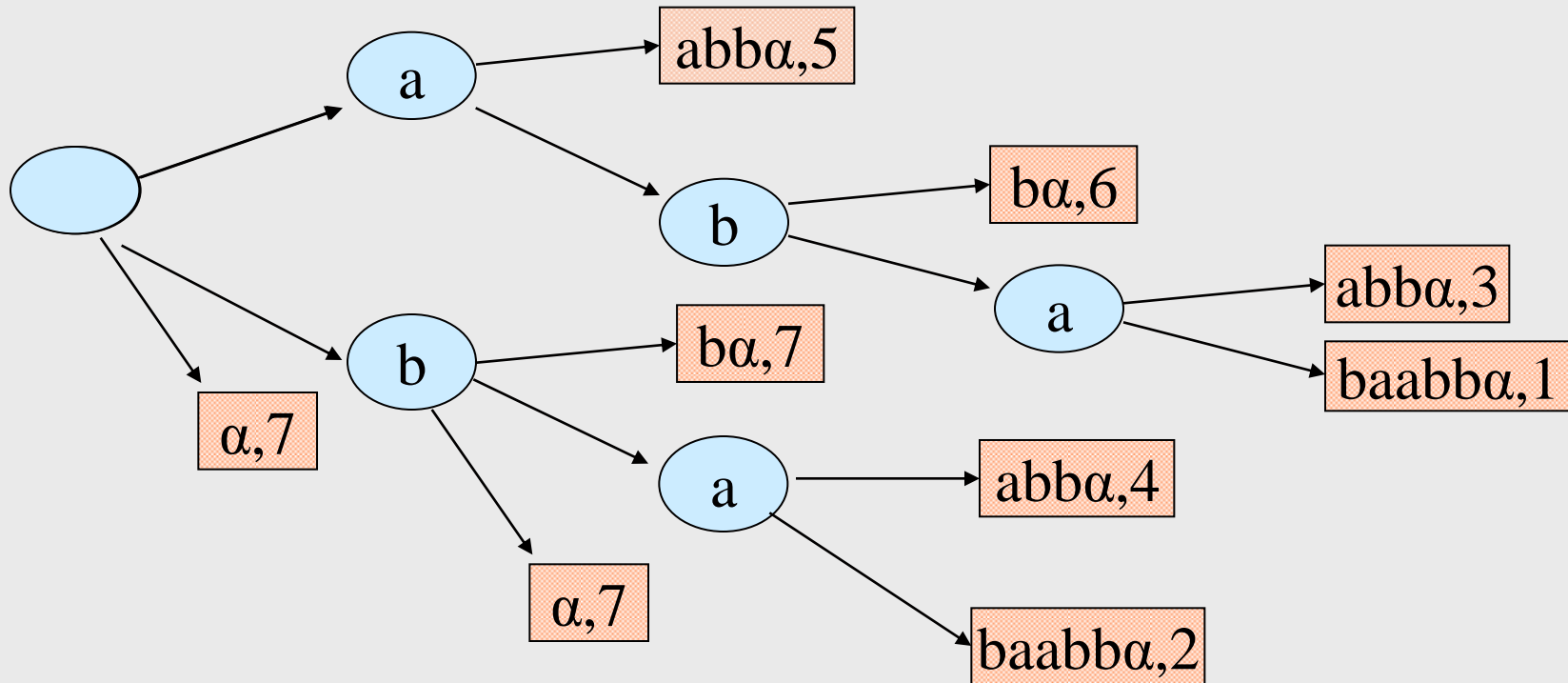
the generalized suffix tree of **ababaabb** and **aabaat** ...

is the suffix tree of **ababaab** α **aabaat** β , :

Generalized suffix tree

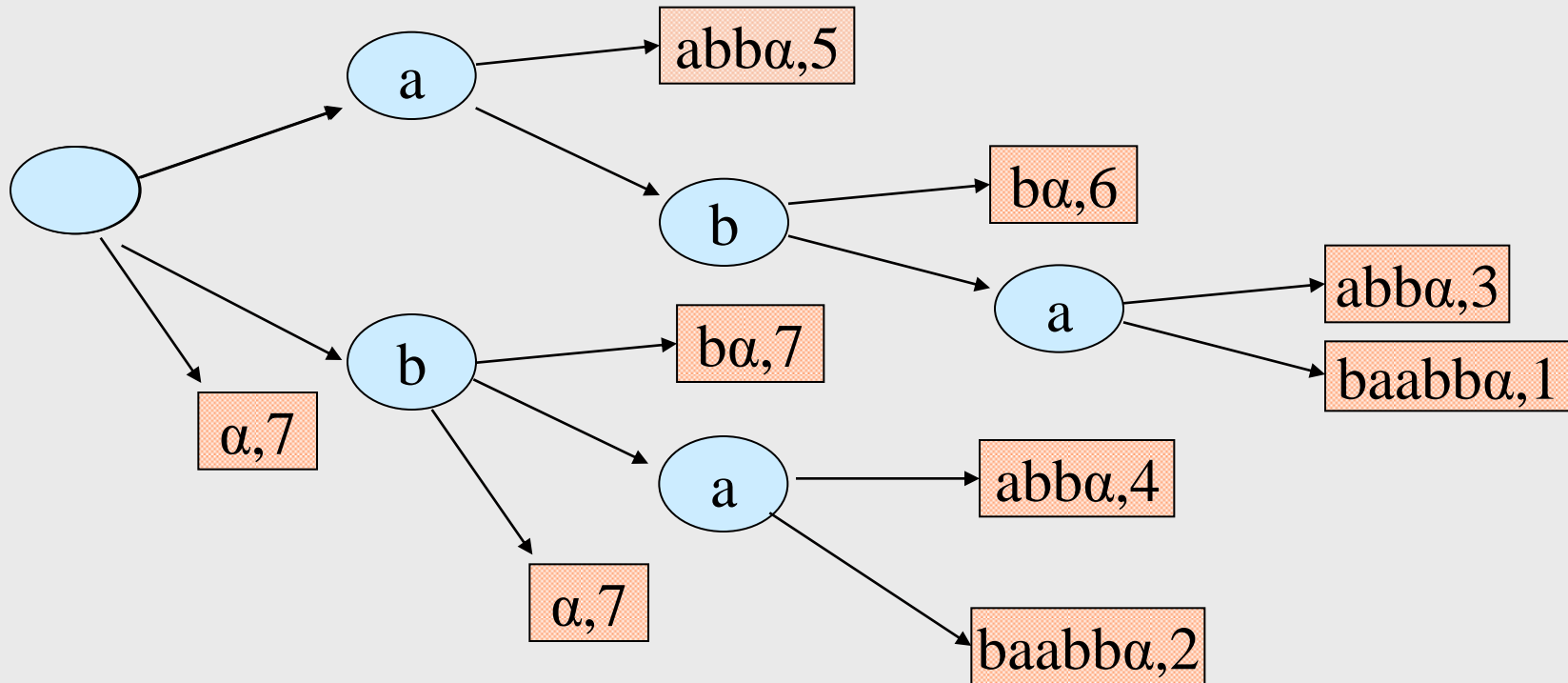
Construction of the suffix tree of $ababaabb\alpha aaba\beta$:

Given the suffix tree of $ababaab\alpha$:



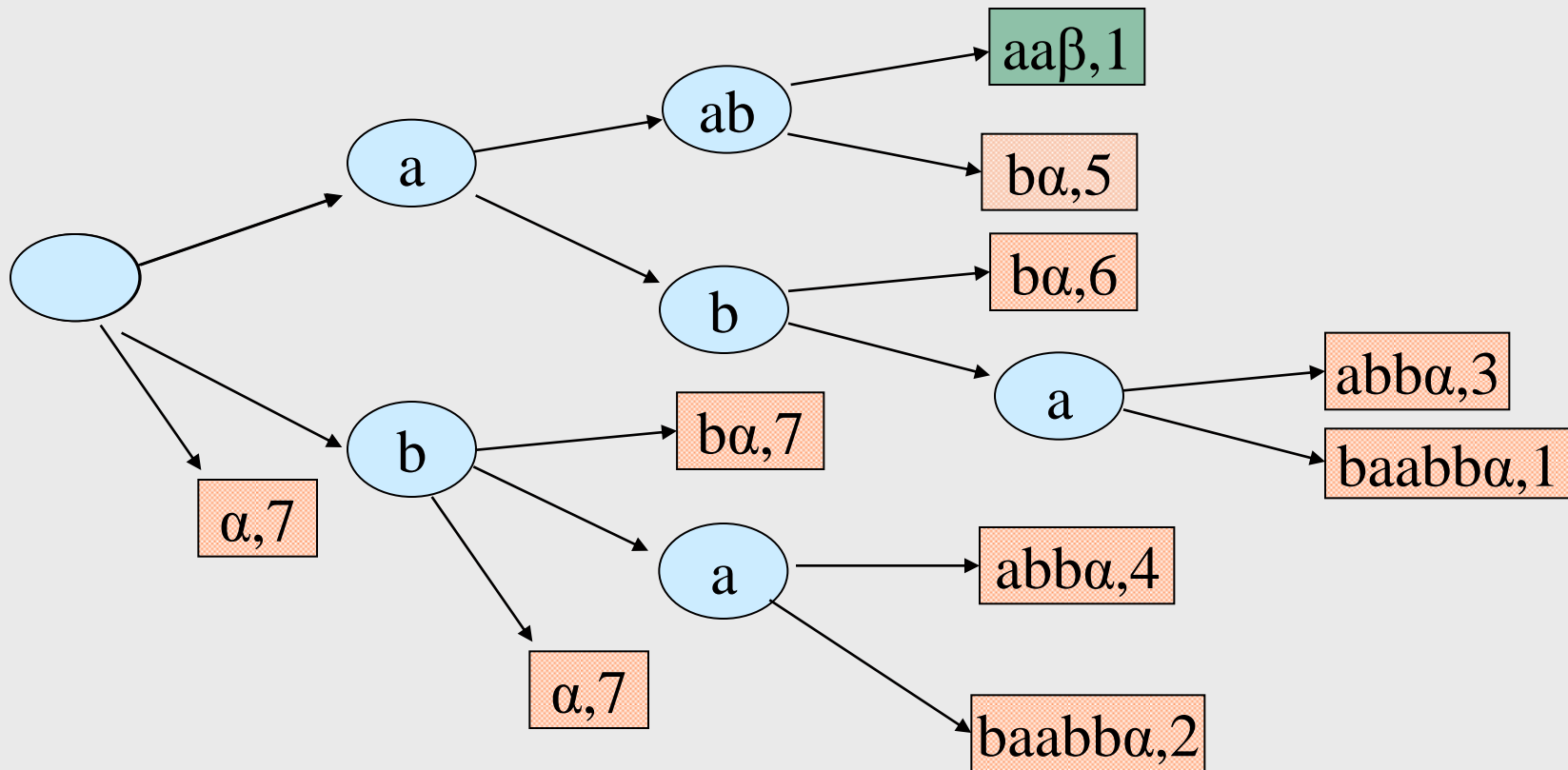
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\alpha\beta$:



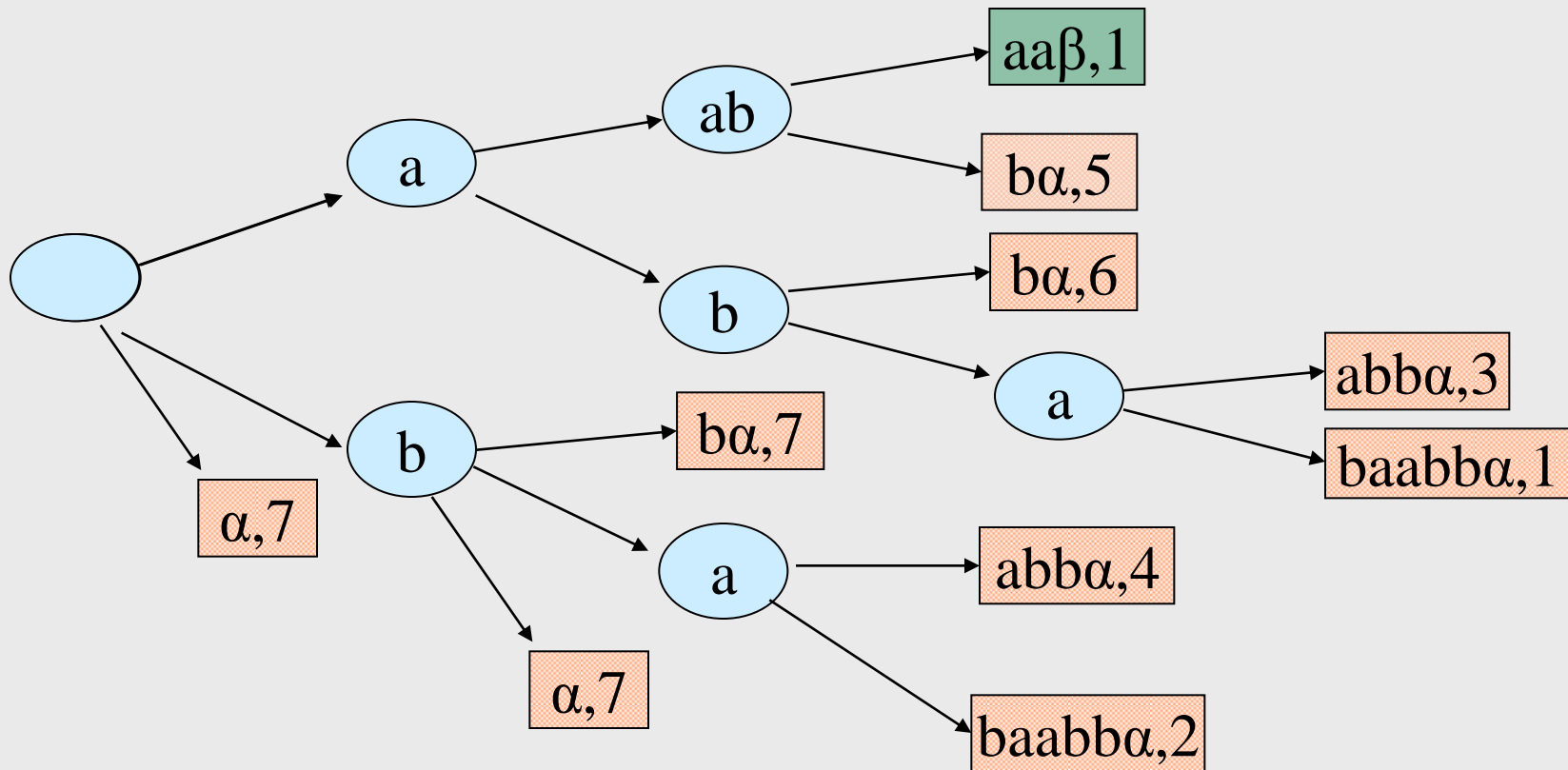
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\alpha\beta$:



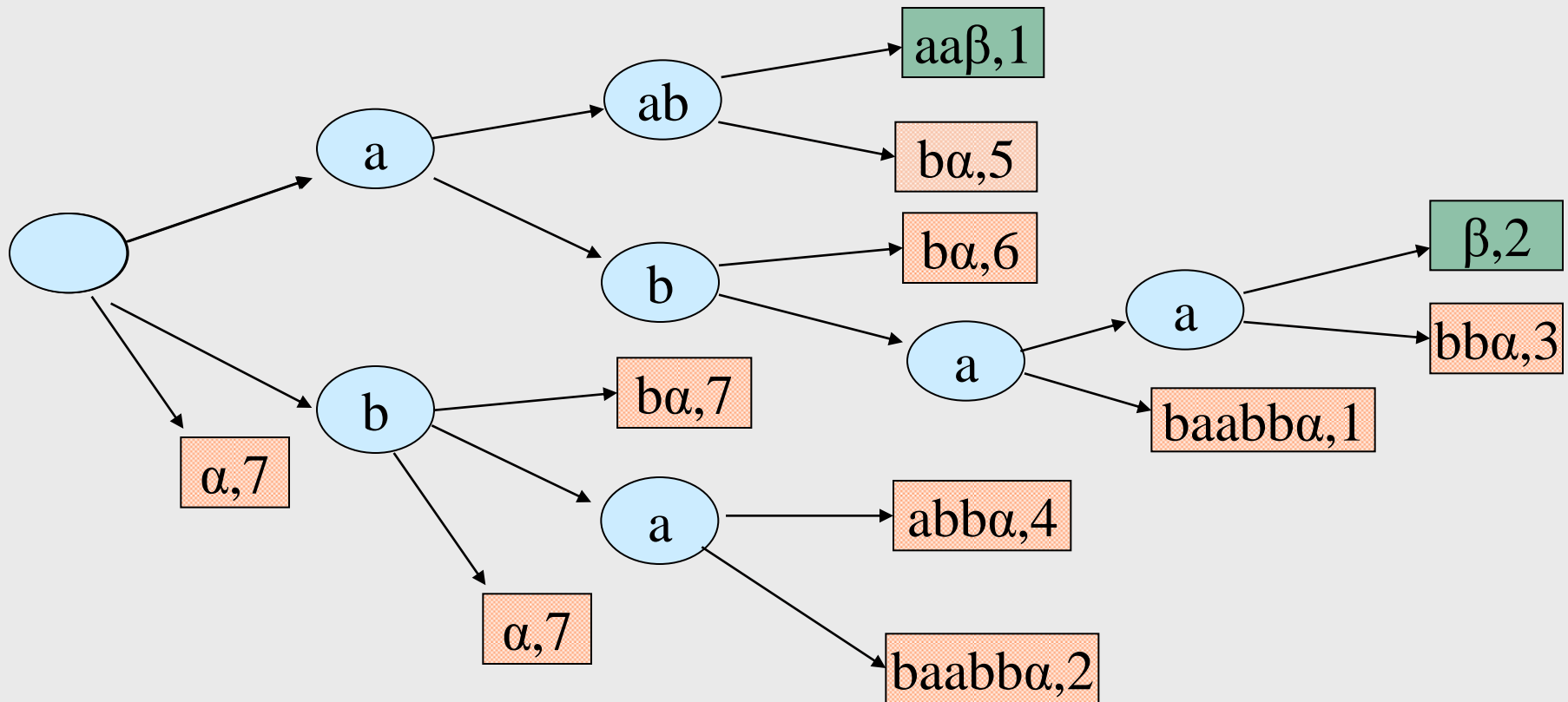
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\alpha\beta$:



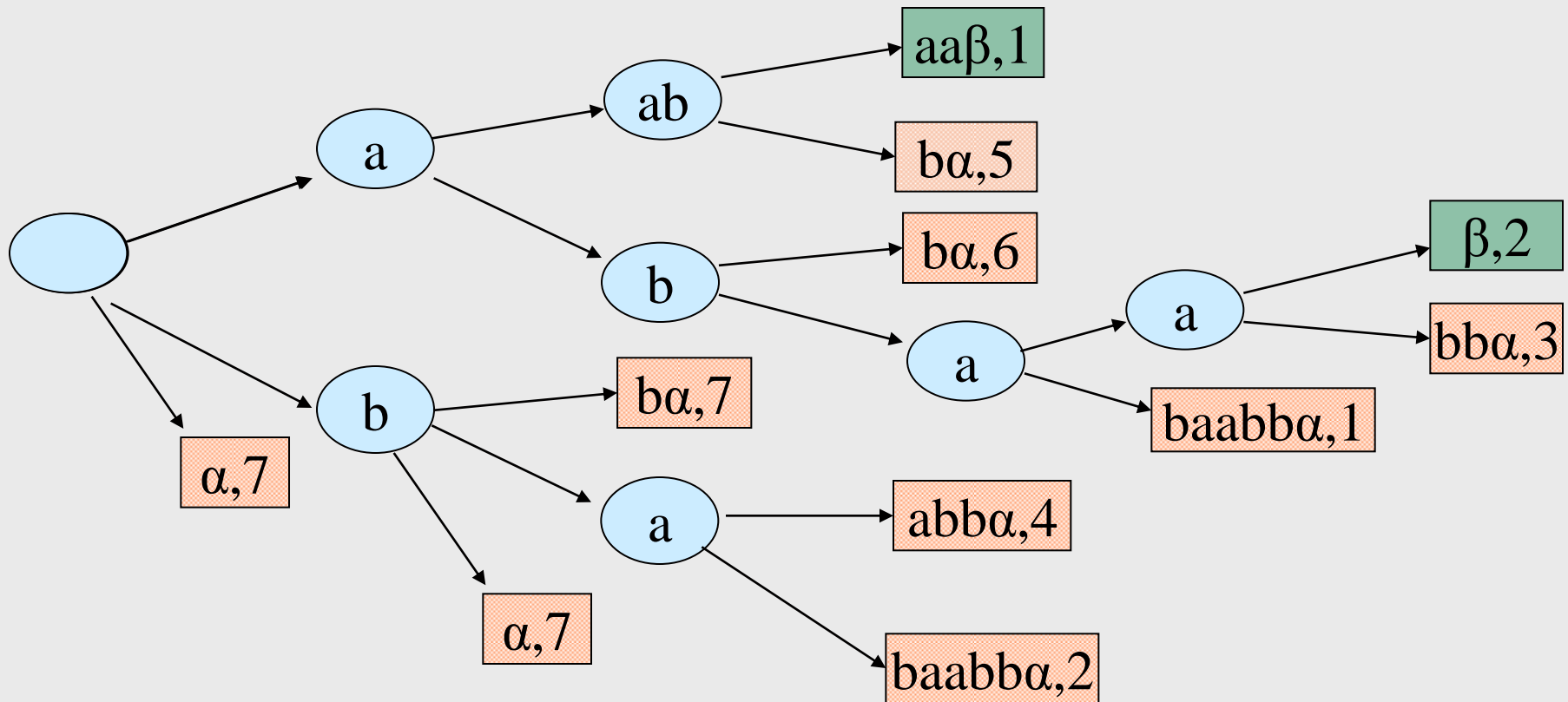
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



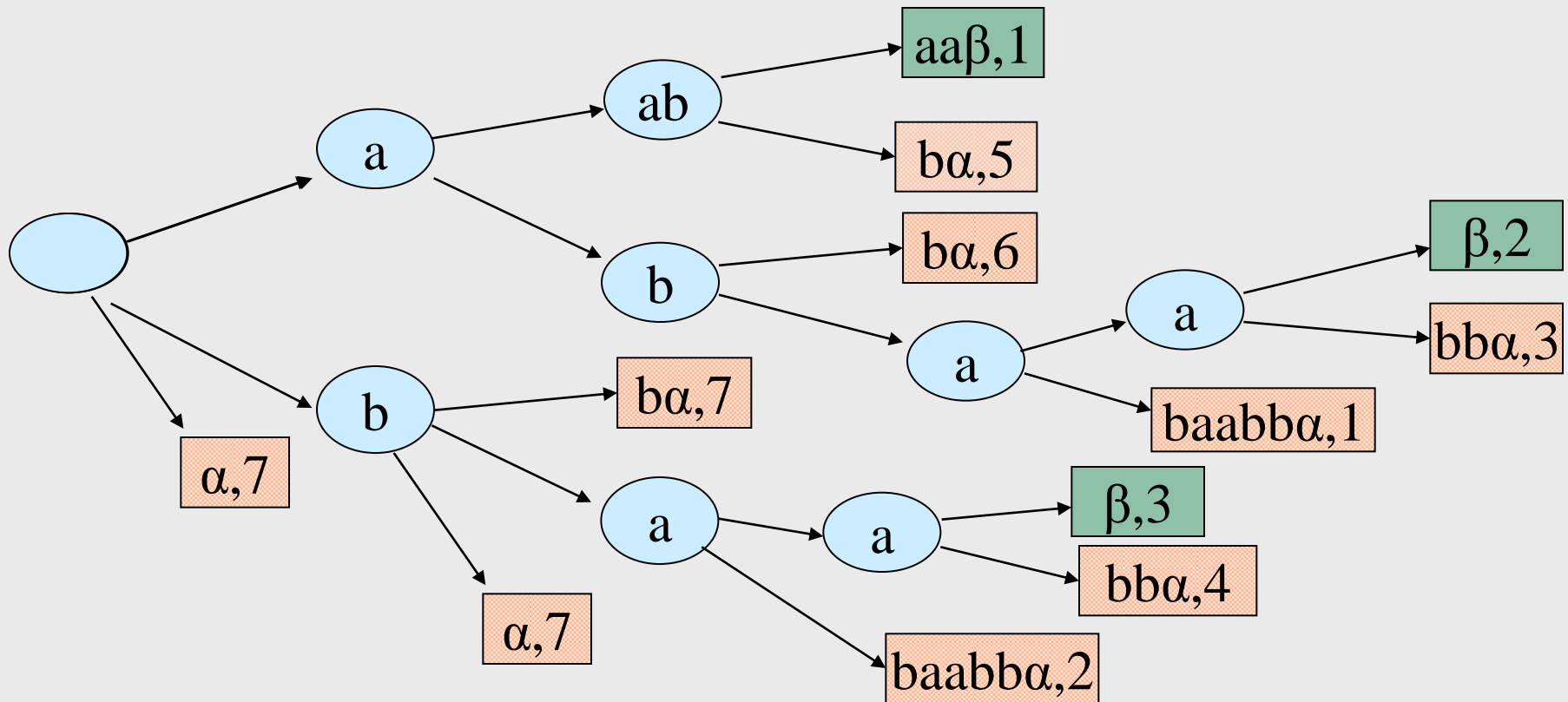
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:

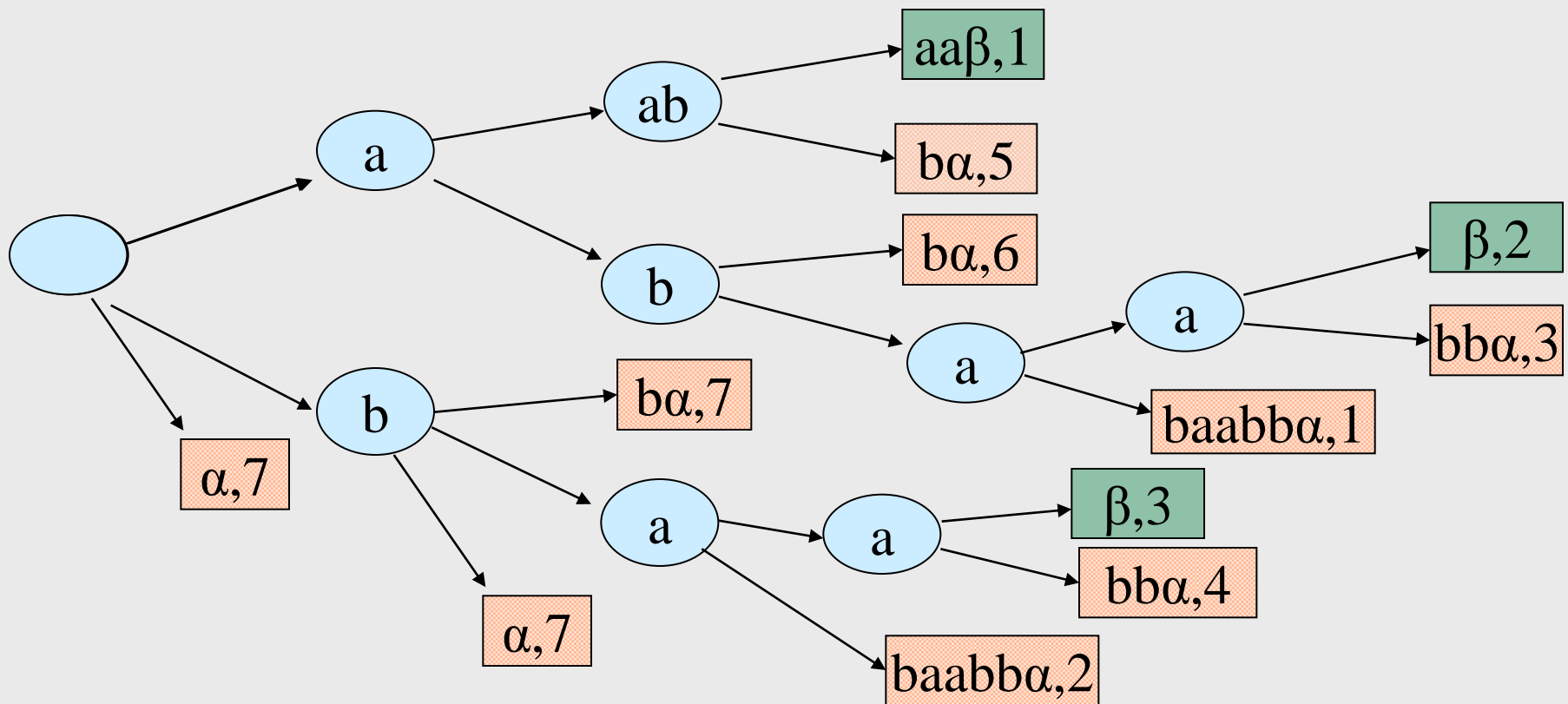


Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:

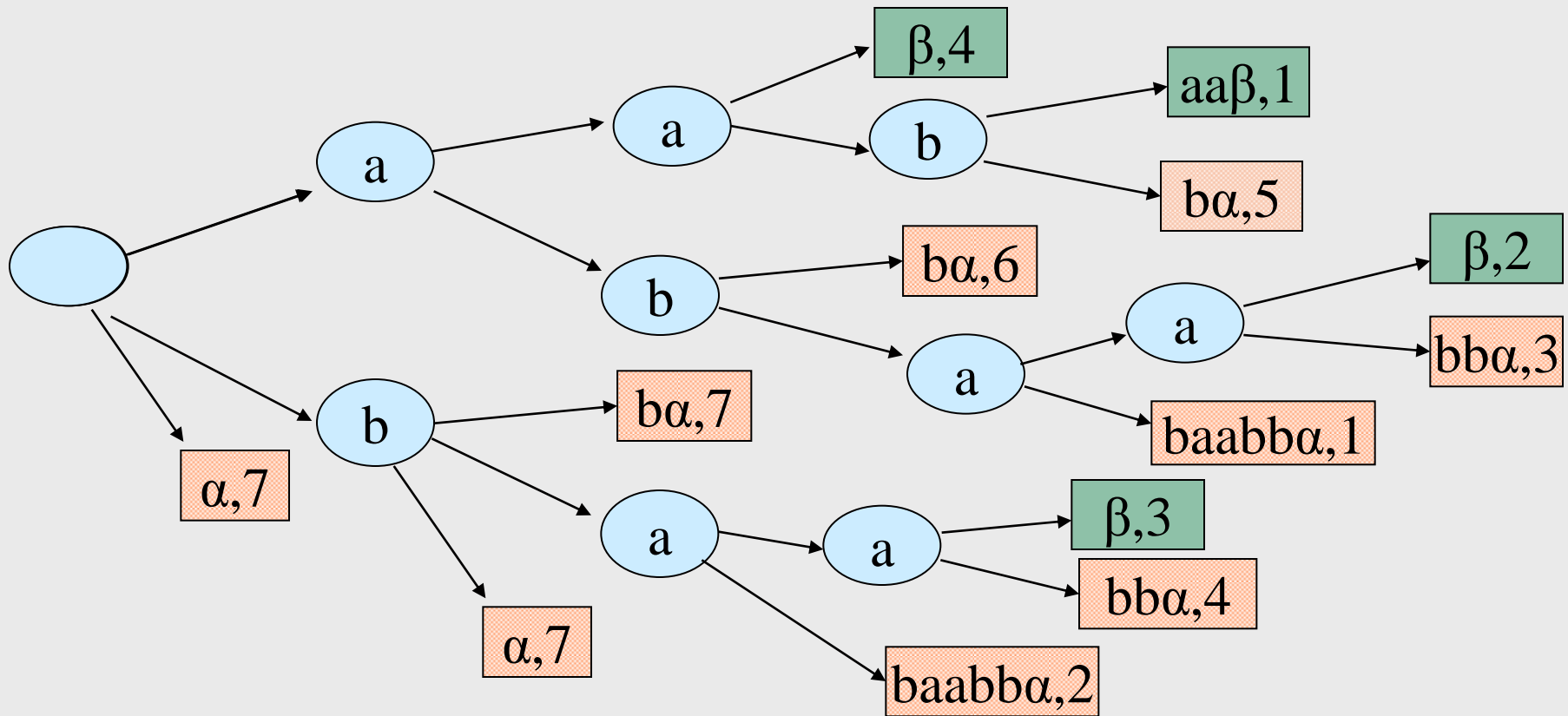


Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



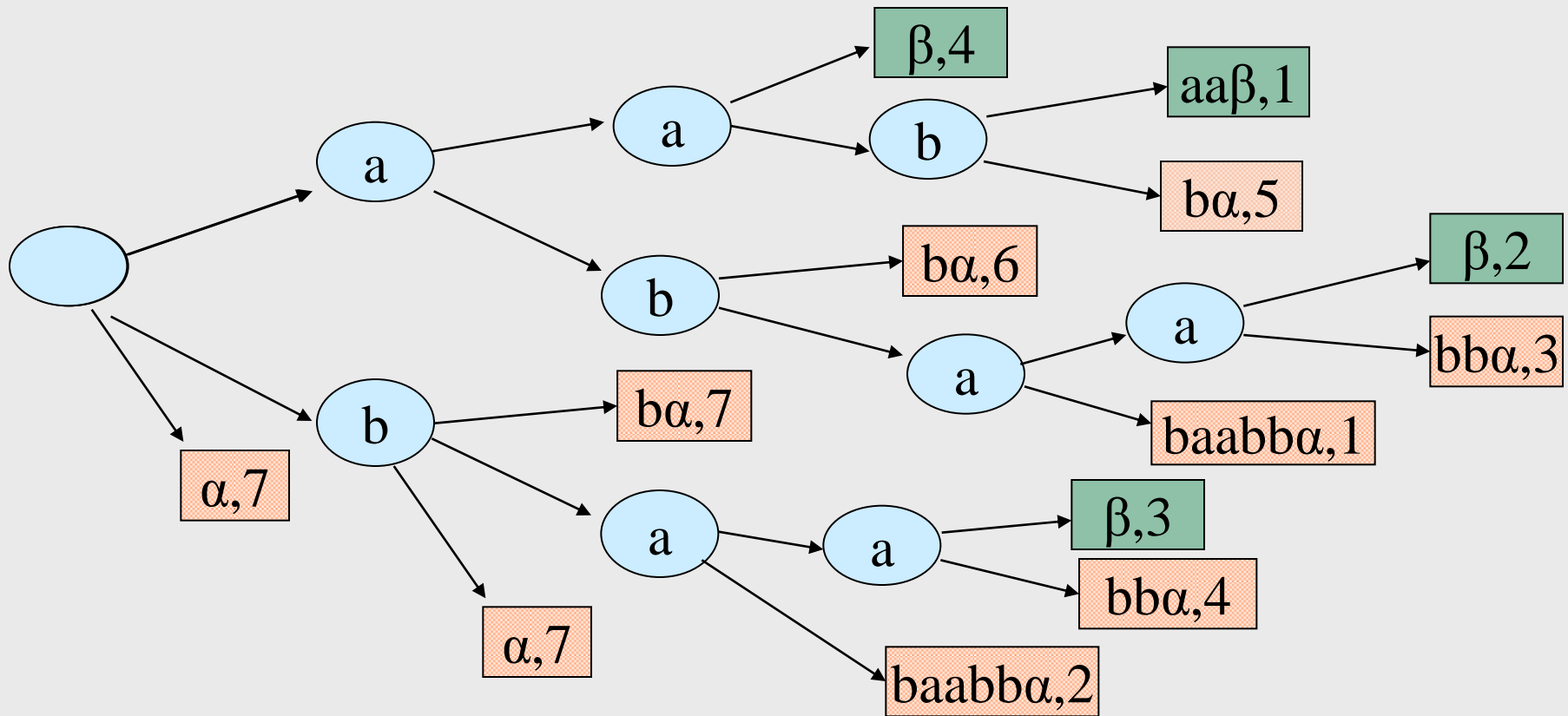
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



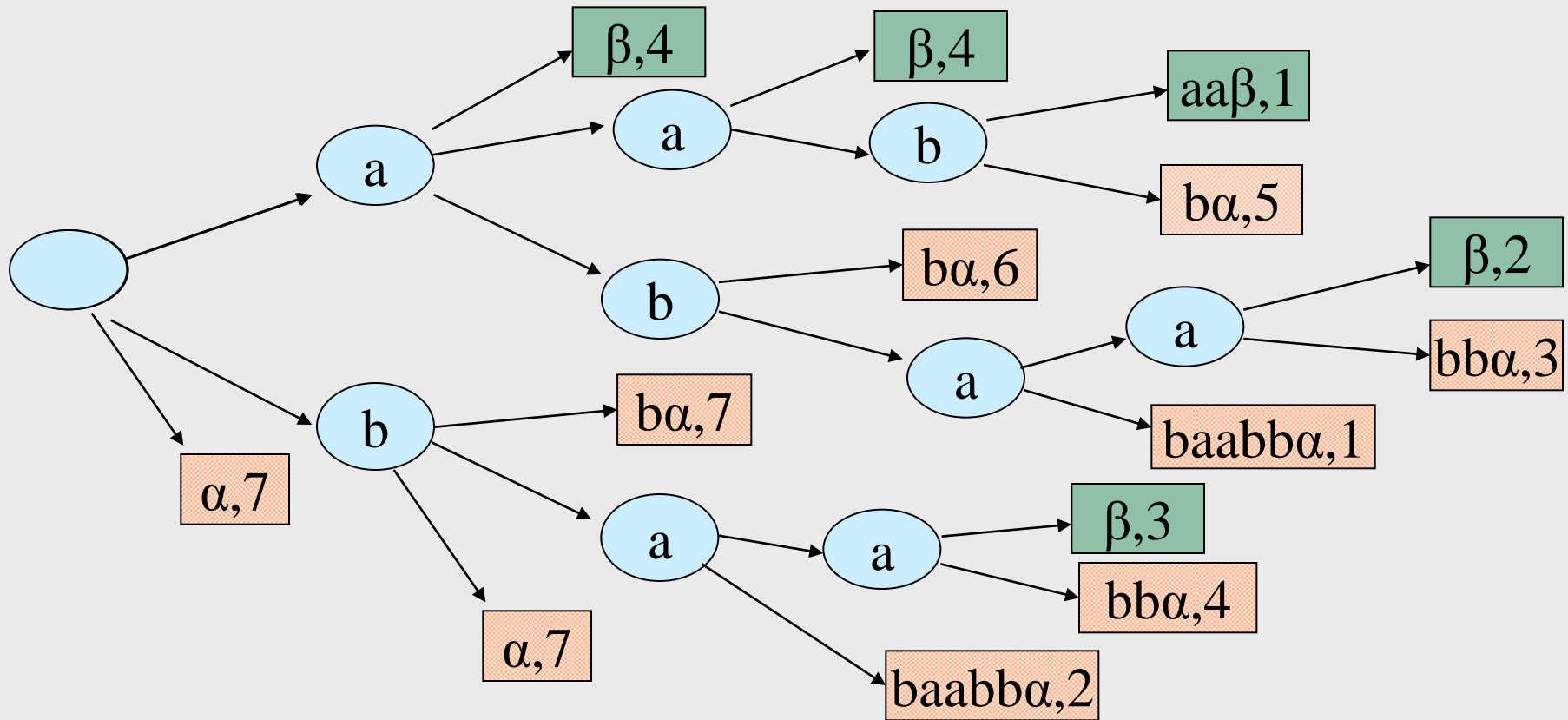
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



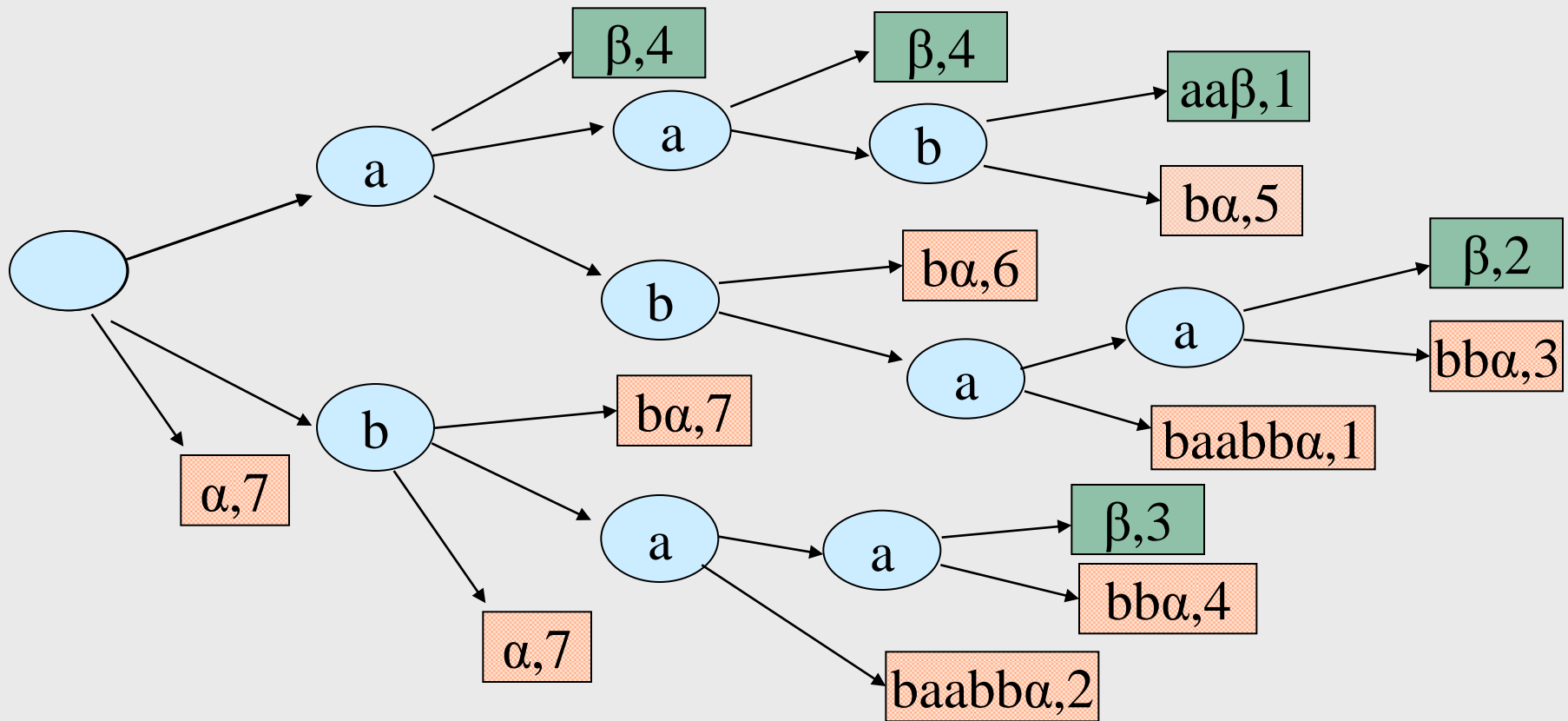
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



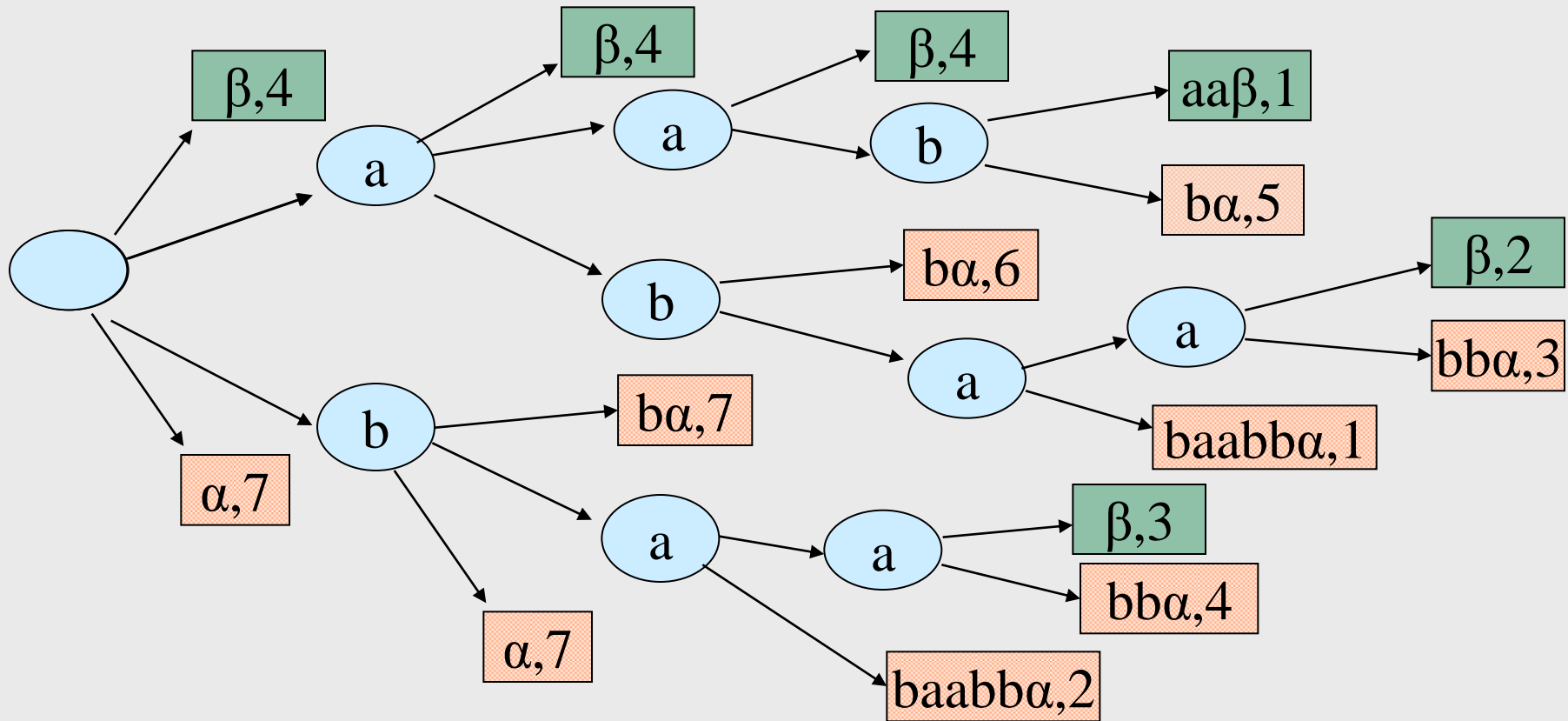
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



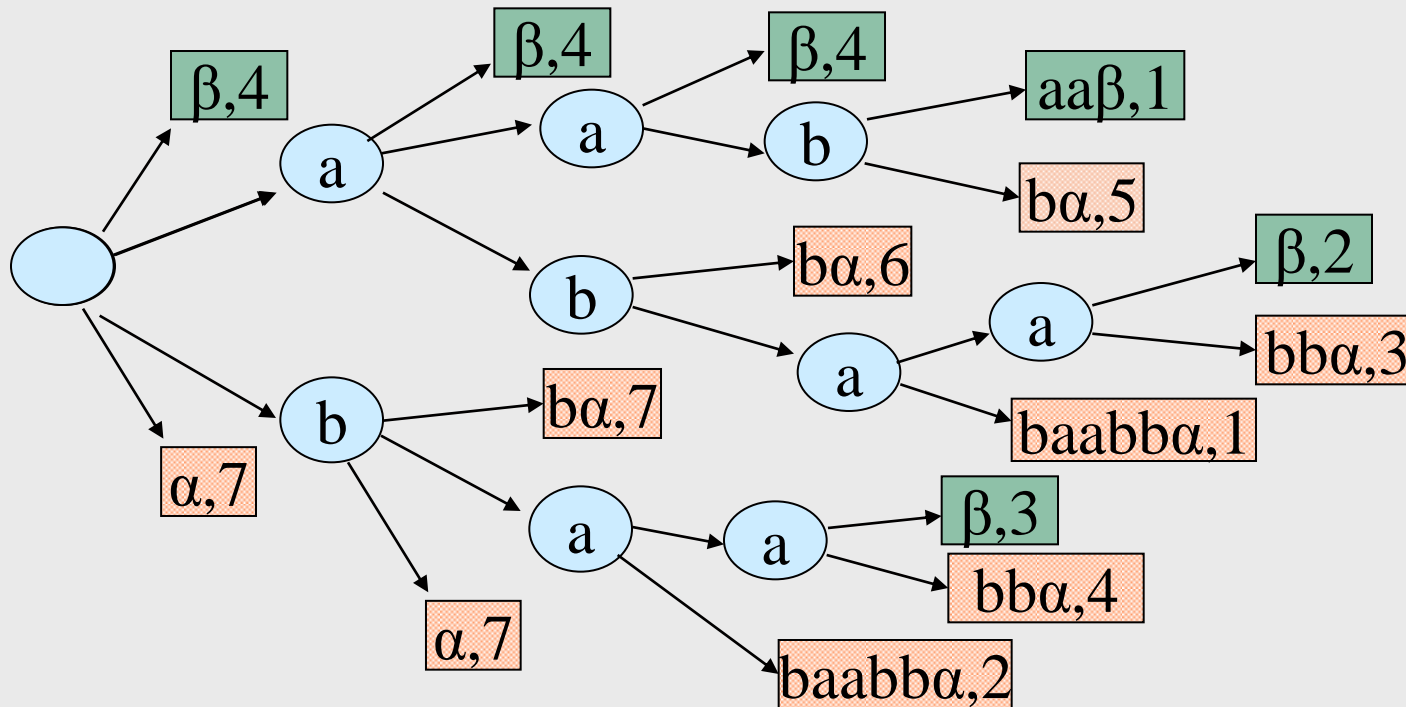
Generalized suffix tree

Construction of the suffix tree of $ababaabb\alpha aaba\beta$:



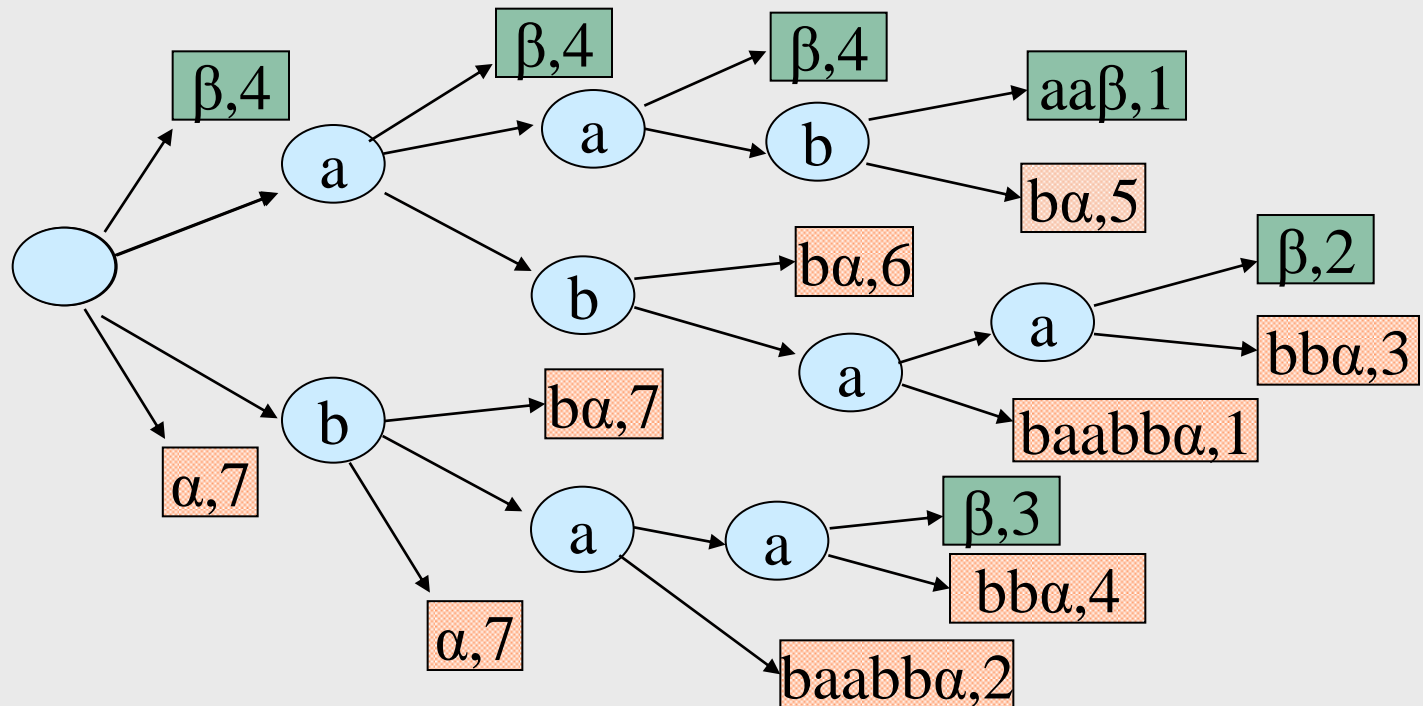
Generalized suffix tree

Generalized suffix tree of $ababaabb\alpha aaba\beta$:

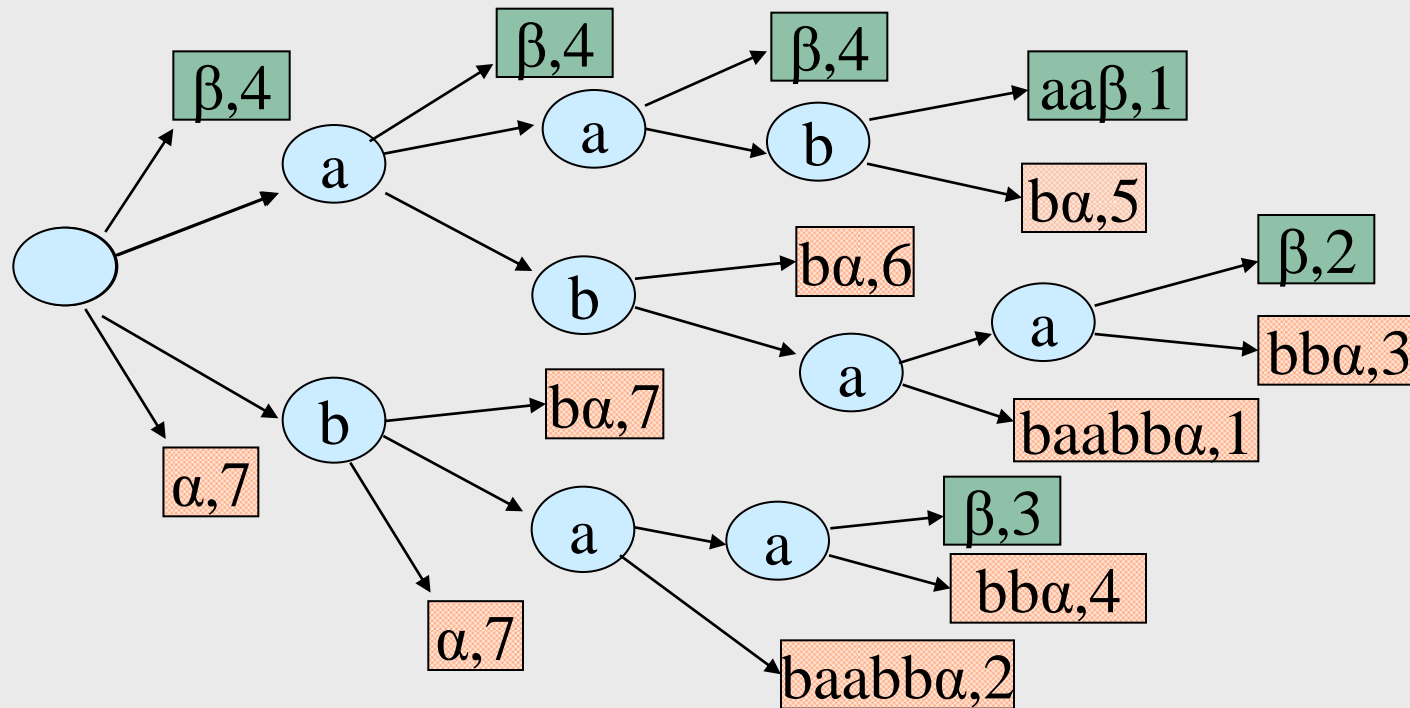


1. The substring problem for a database of strings DB

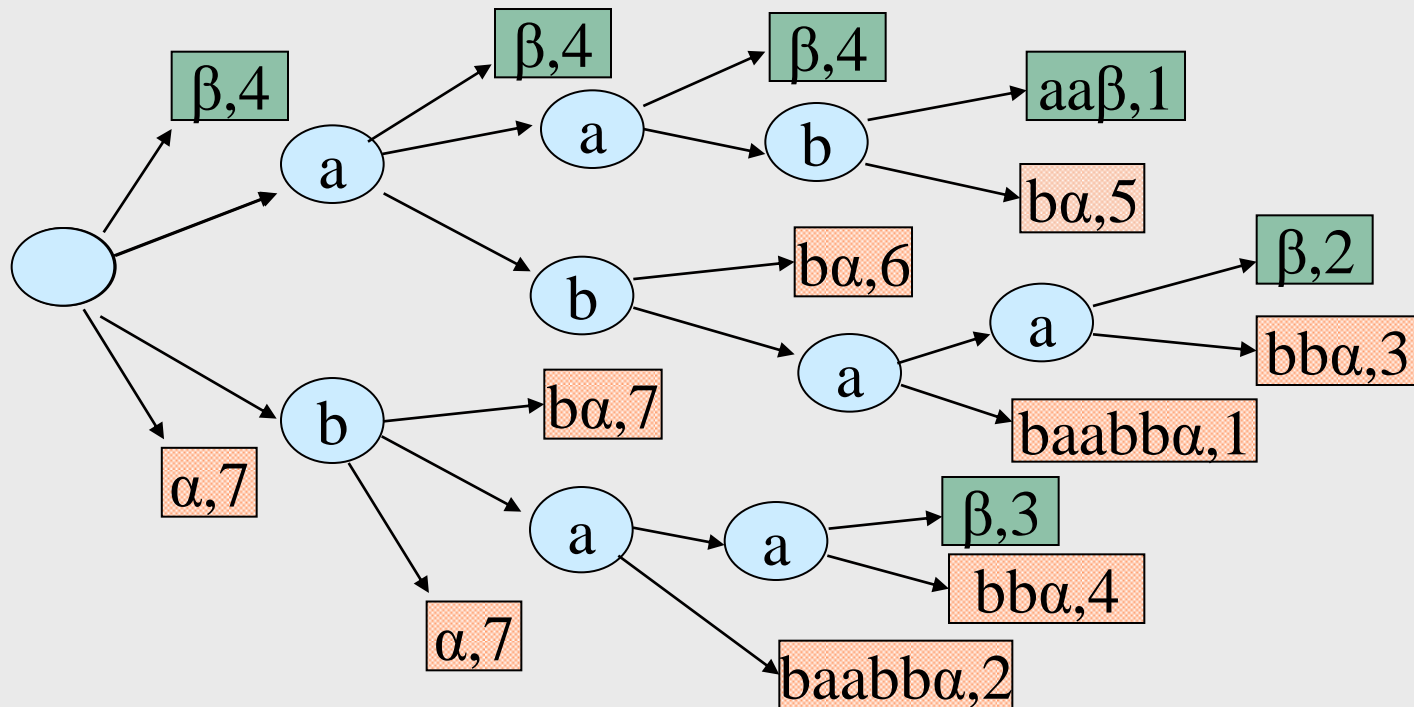
- Does the DB contain any occurrence of patterns **abab**, **aab**, and **ab**?



2. The longest common substring of two strings



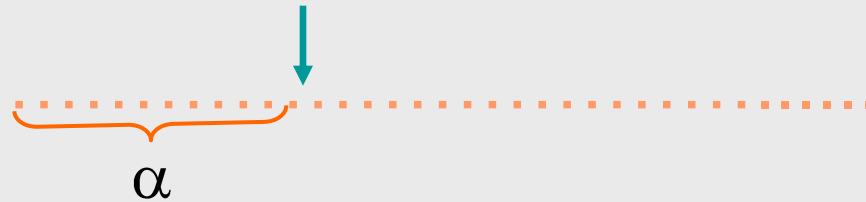
3. Finding MUMs.



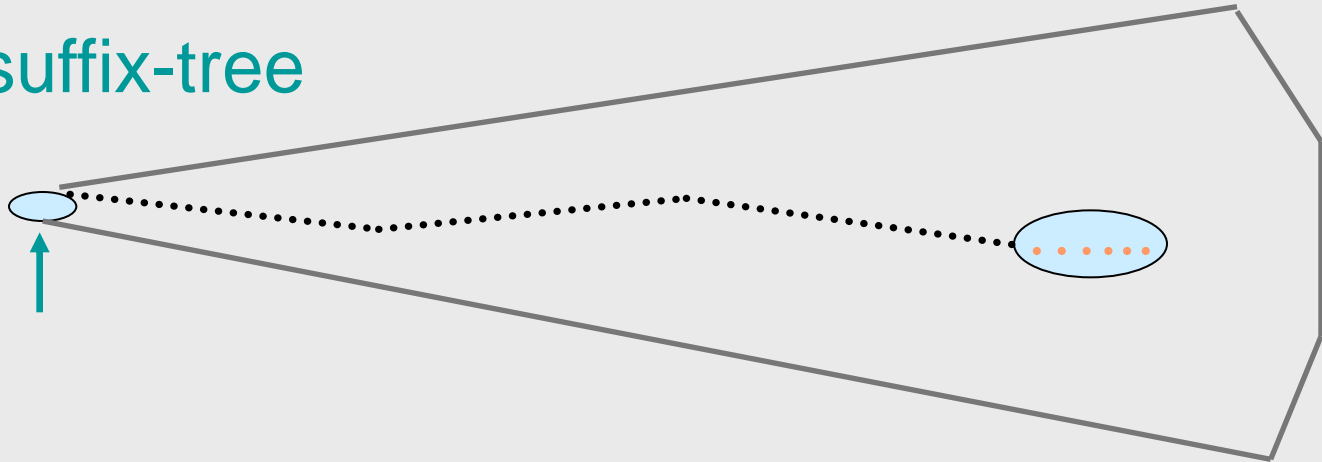
Quadratic insertion algorithm

Invariant Properties:

Given the string



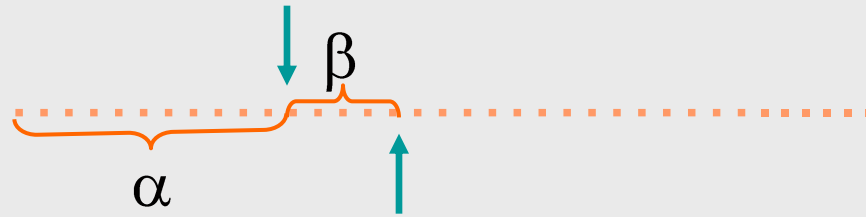
and the suffix-tree



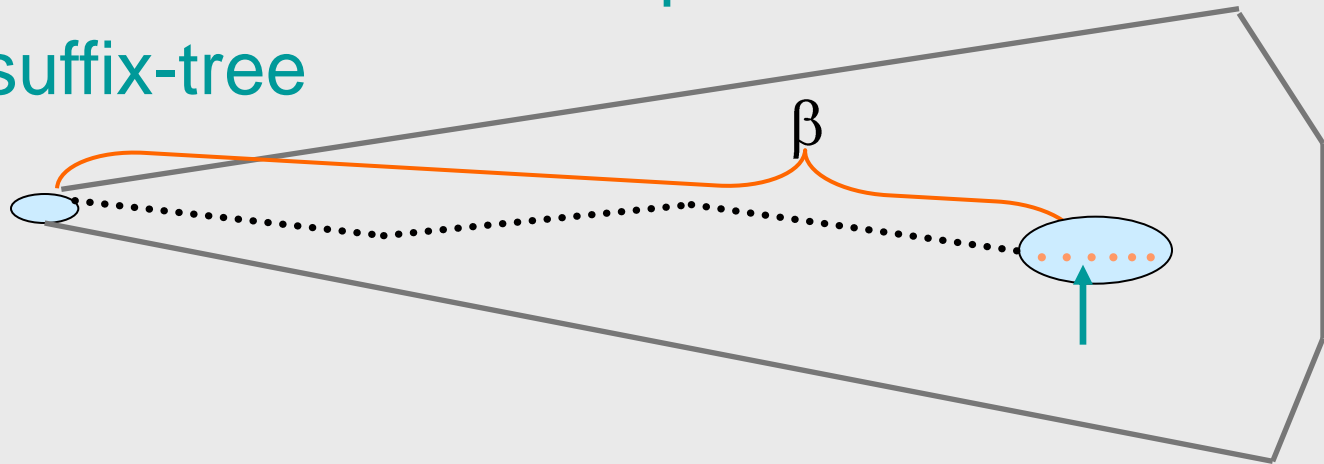
P1: the leaves of suffixes from α have been inserted

Invariant Properties:

Given the string



and the suffix-tree



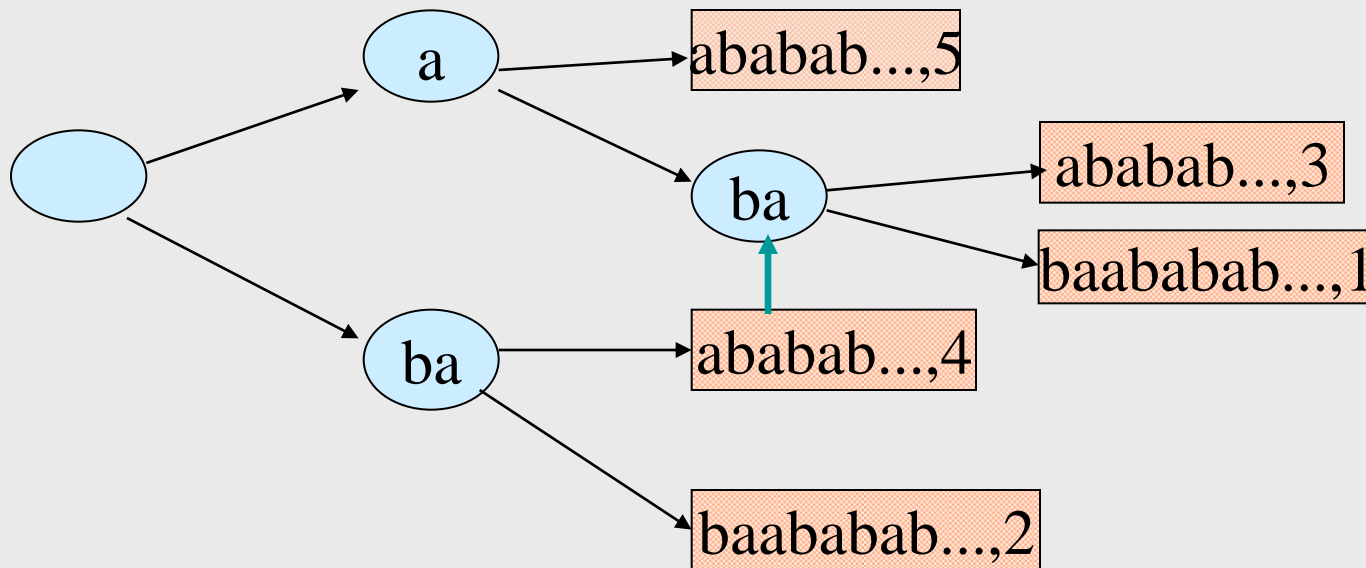
P1: the leaves of suffixes from α have been inserted

P2: the string β is the longest string that can be spelt through the tree.

Linear insertion algorithm: example

Given the string $ababaababb\dots$

α β

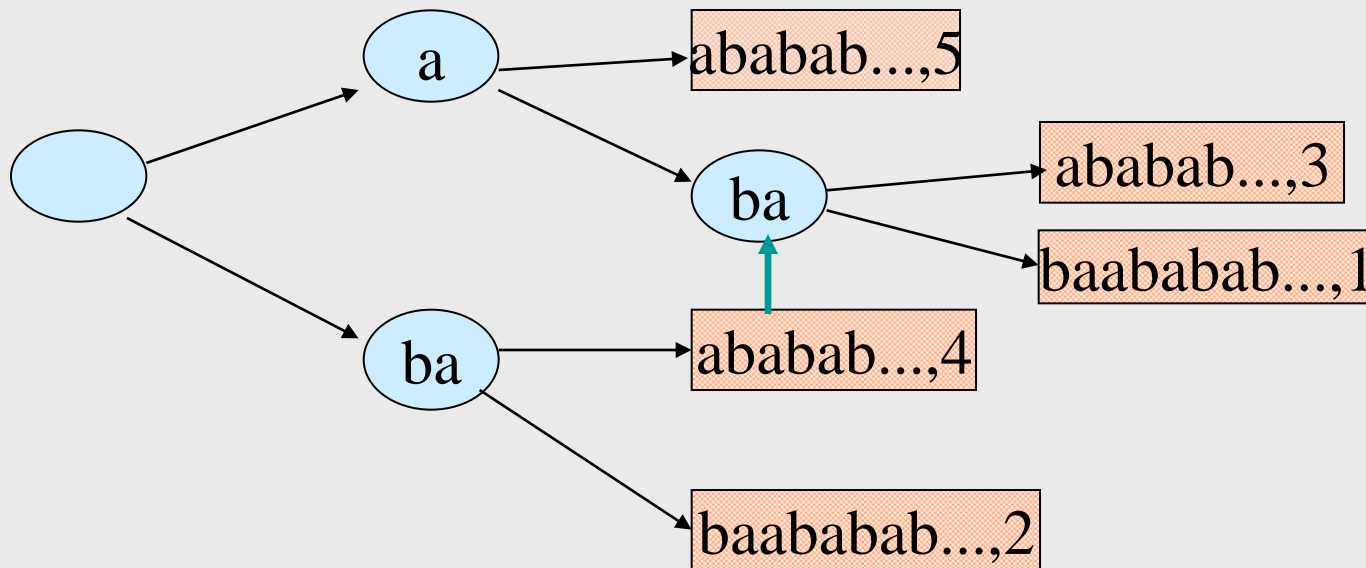


Linear insertion algorithm: example

Given the string $ababaababb\dots$

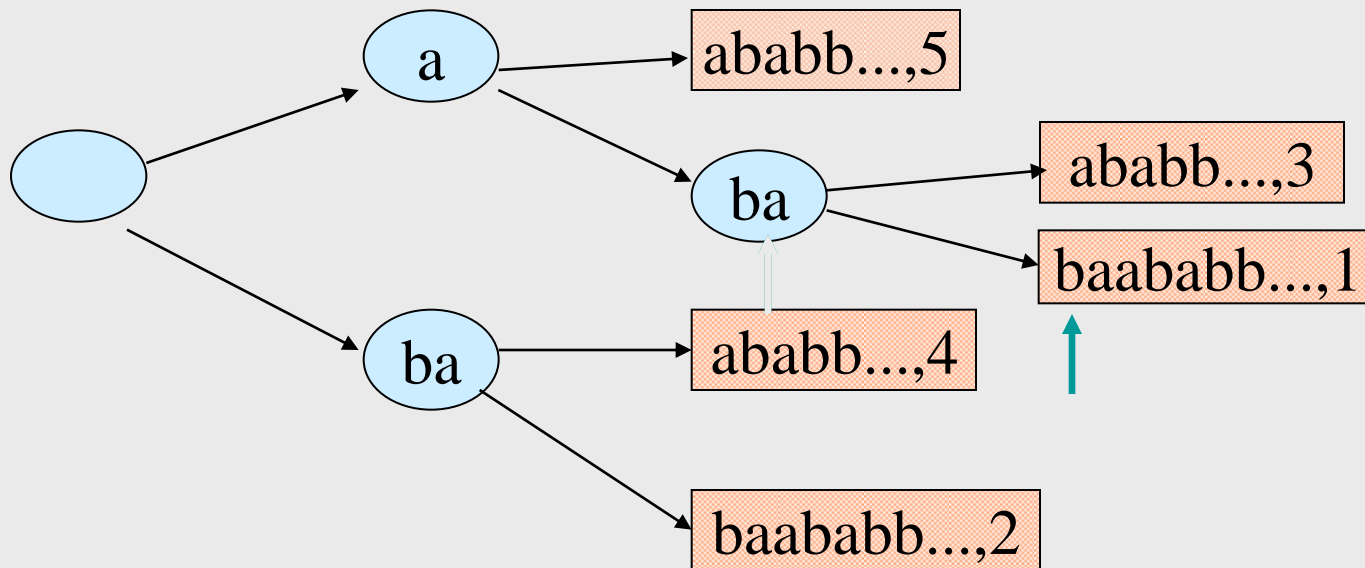
α β

6 7 8



Linear insertion algorithm: example

Given the string $ababaababb\dots$

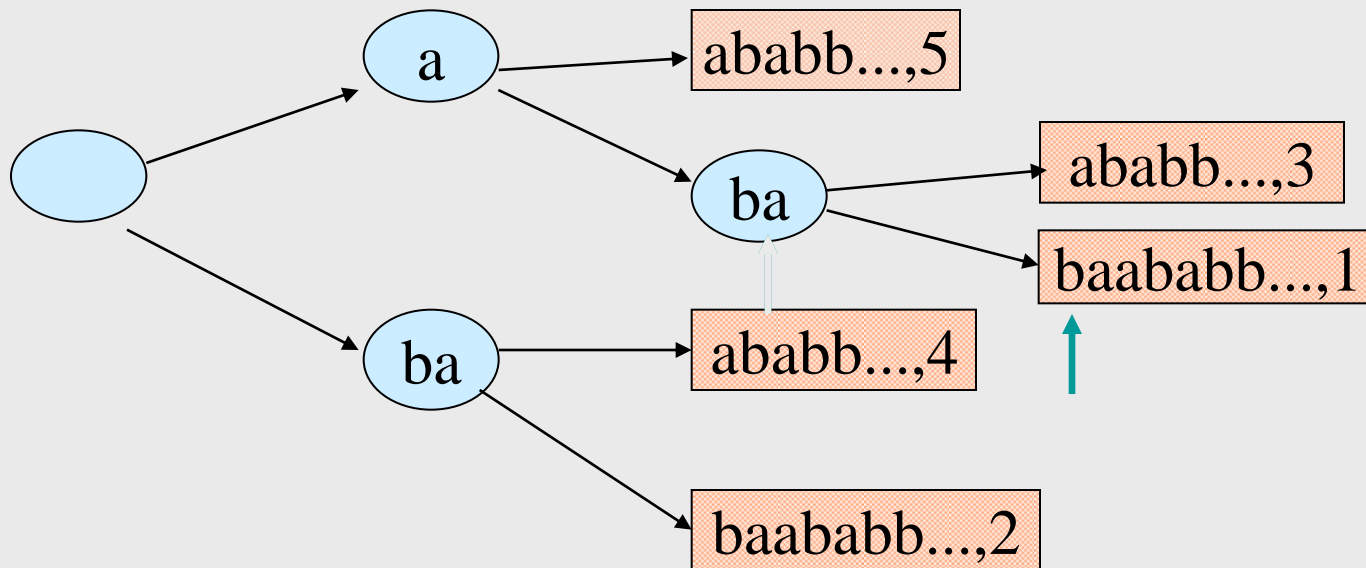


Linear insertion algorithm: example

Given the string $ababaababb\dots$

α β

6 7 8 9

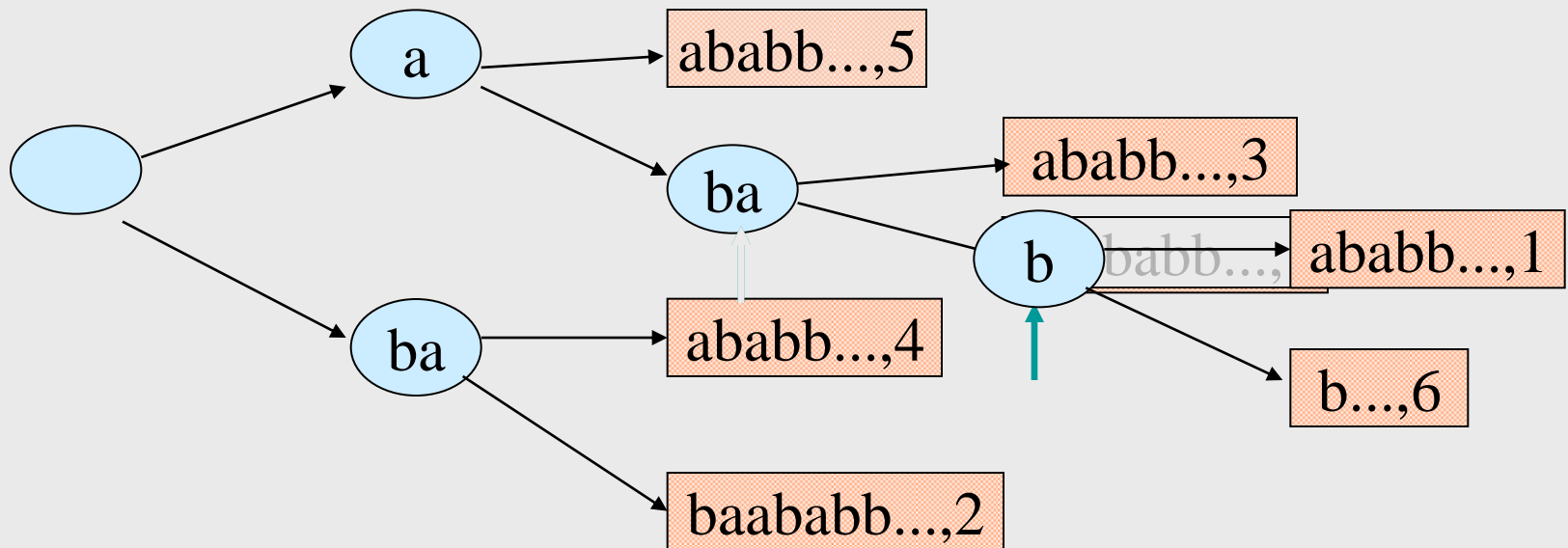


Linear insertion algorithm: example

Given the string $ababaababb\dots$

α β

6 7 8 9

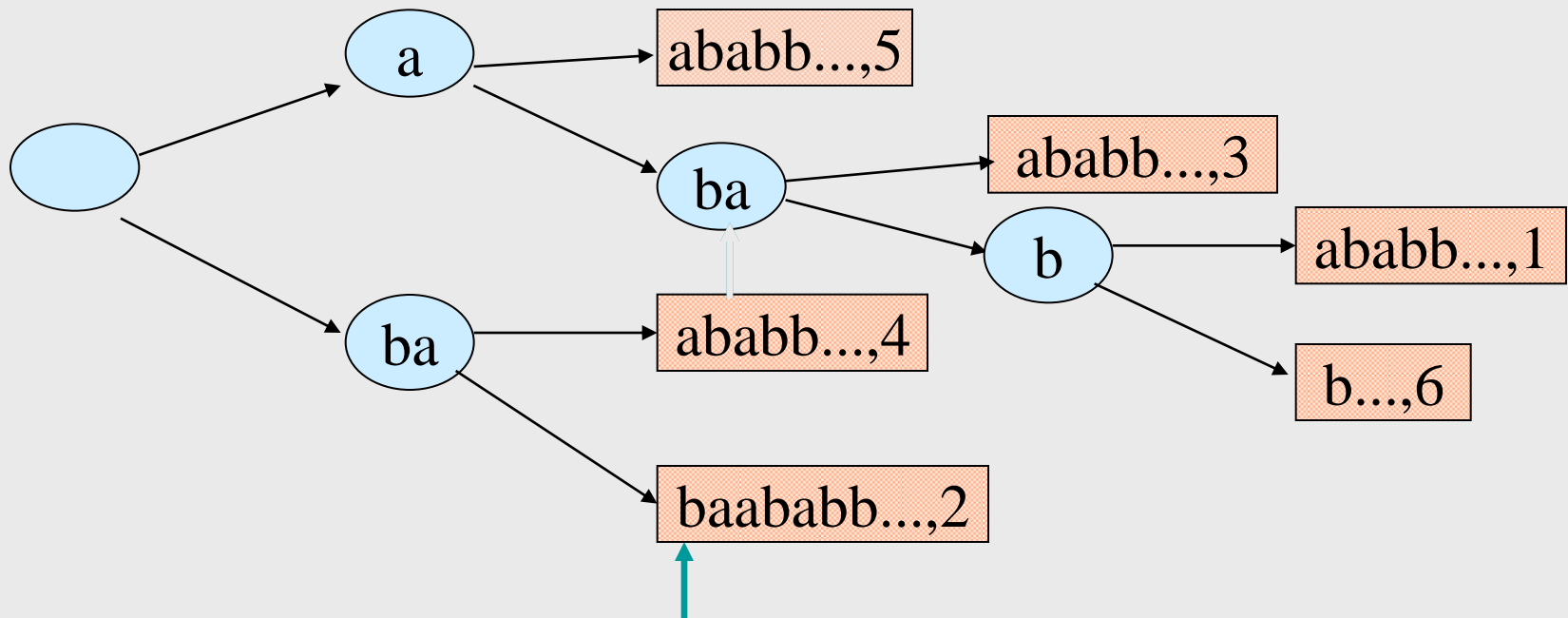


Linear insertion algorithm: example

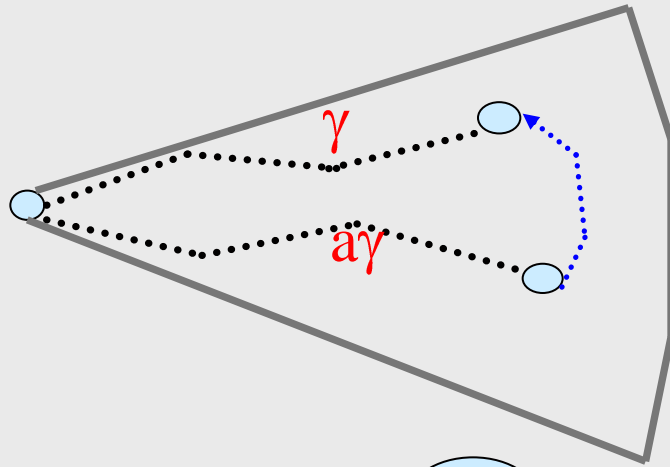
Given the string $ababaababb\dots$

α β

7 8 9



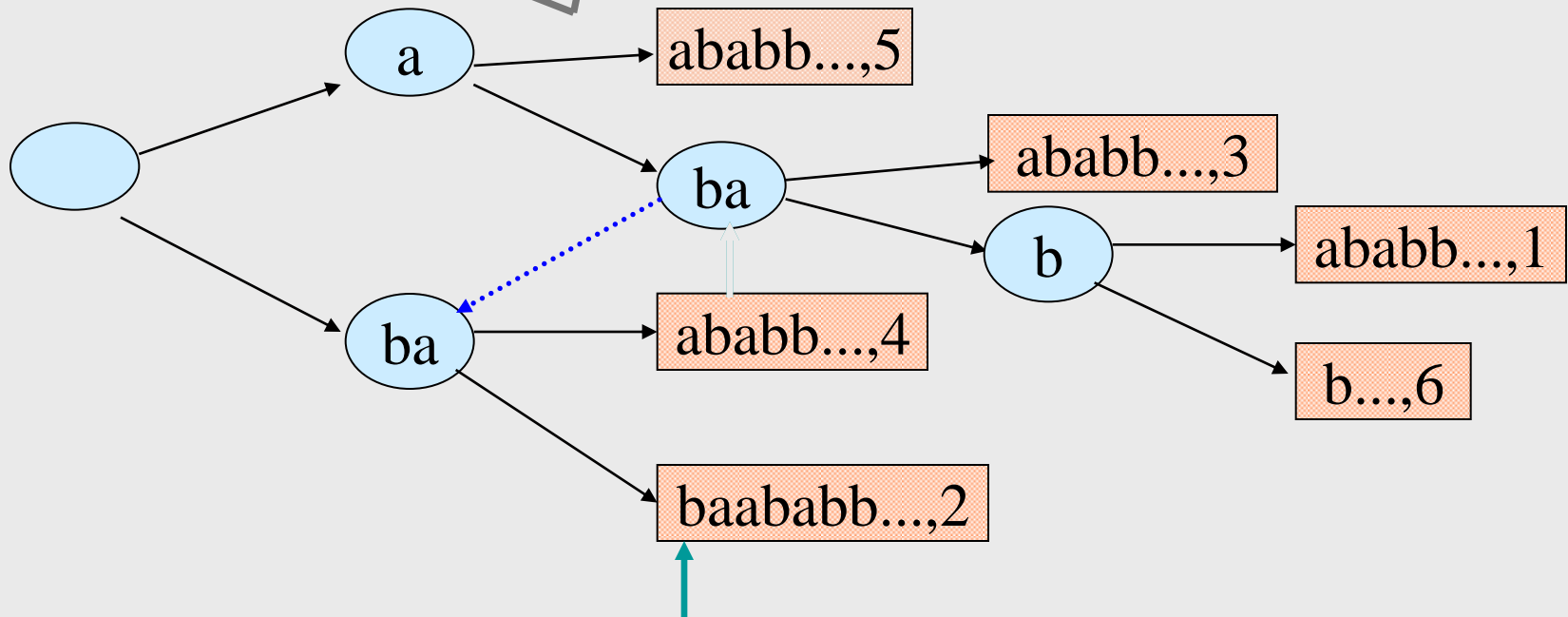
Linear insertion algorithm: example



Given the string **ababaababb...**

α β

7 8 9

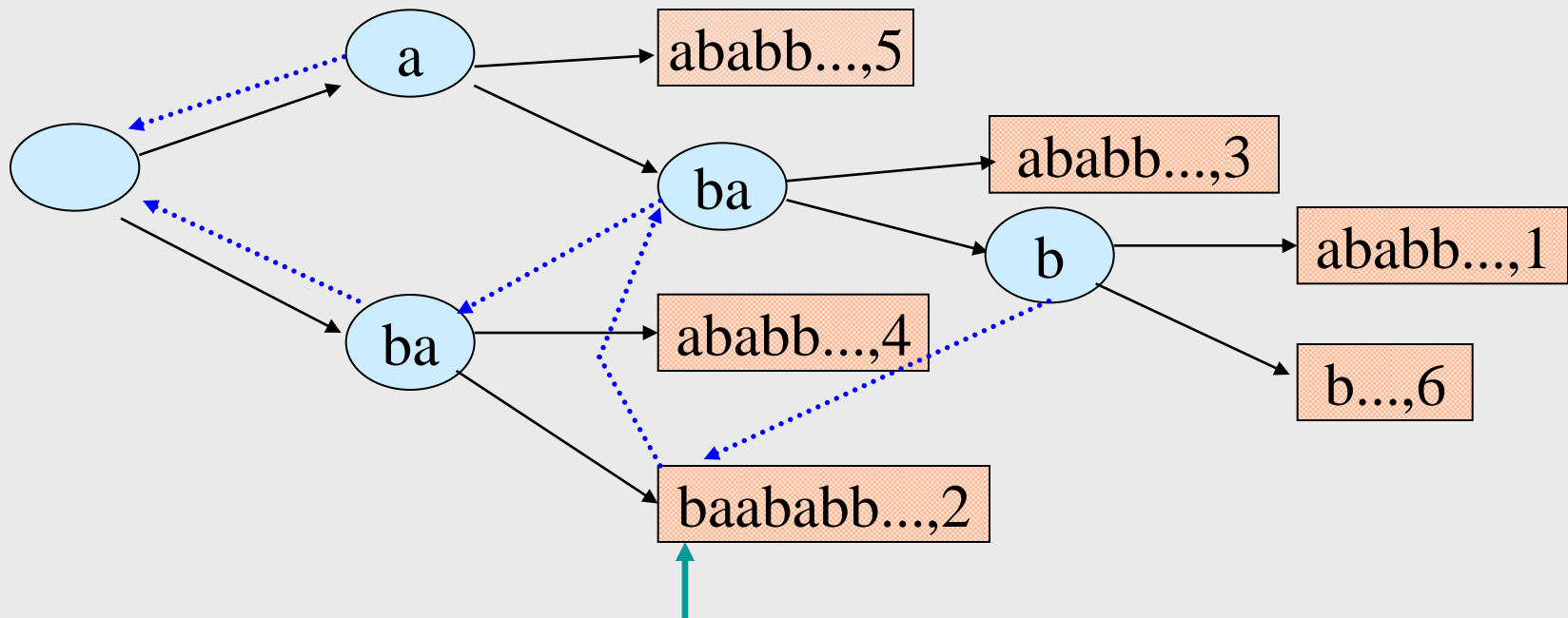


Linear insertion algorithm: example

Given the string $ababaababb\dots$

α β

7 8 9

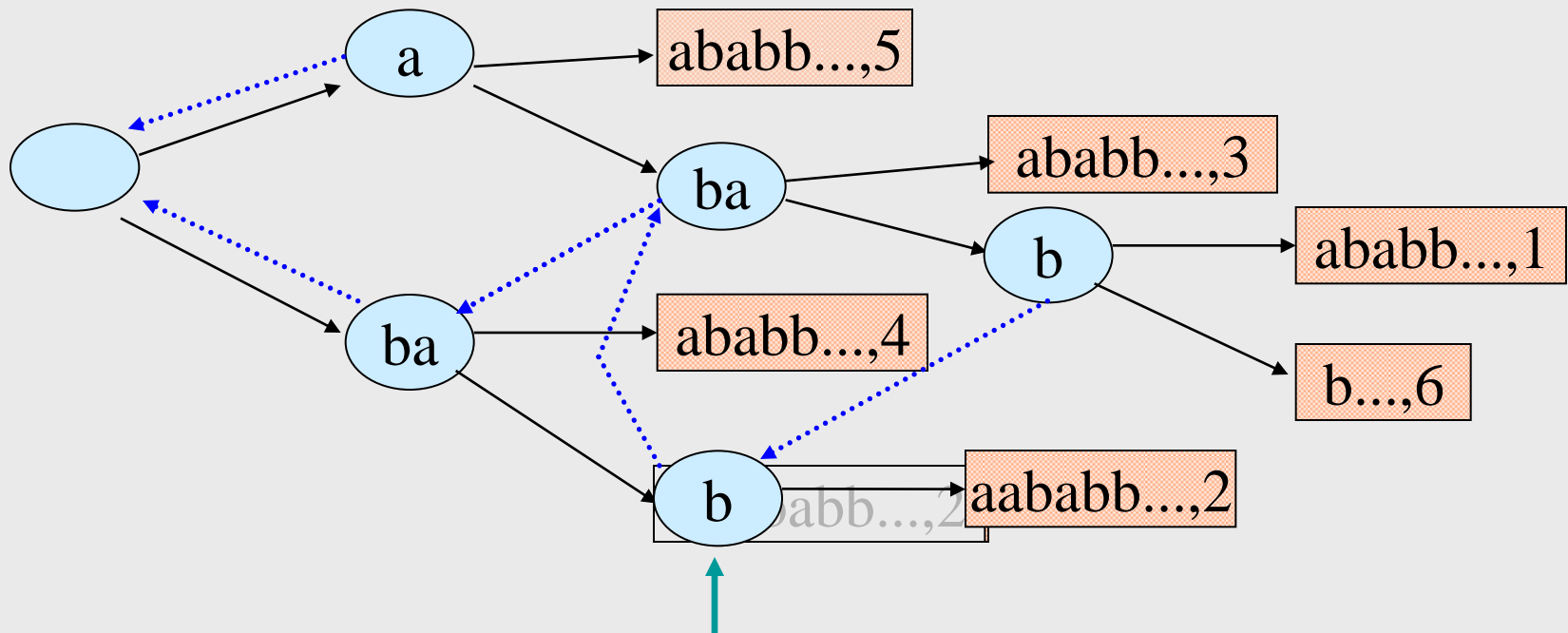


Linear insertion algorithm: example

Given the string $ababaababb\dots$

α β

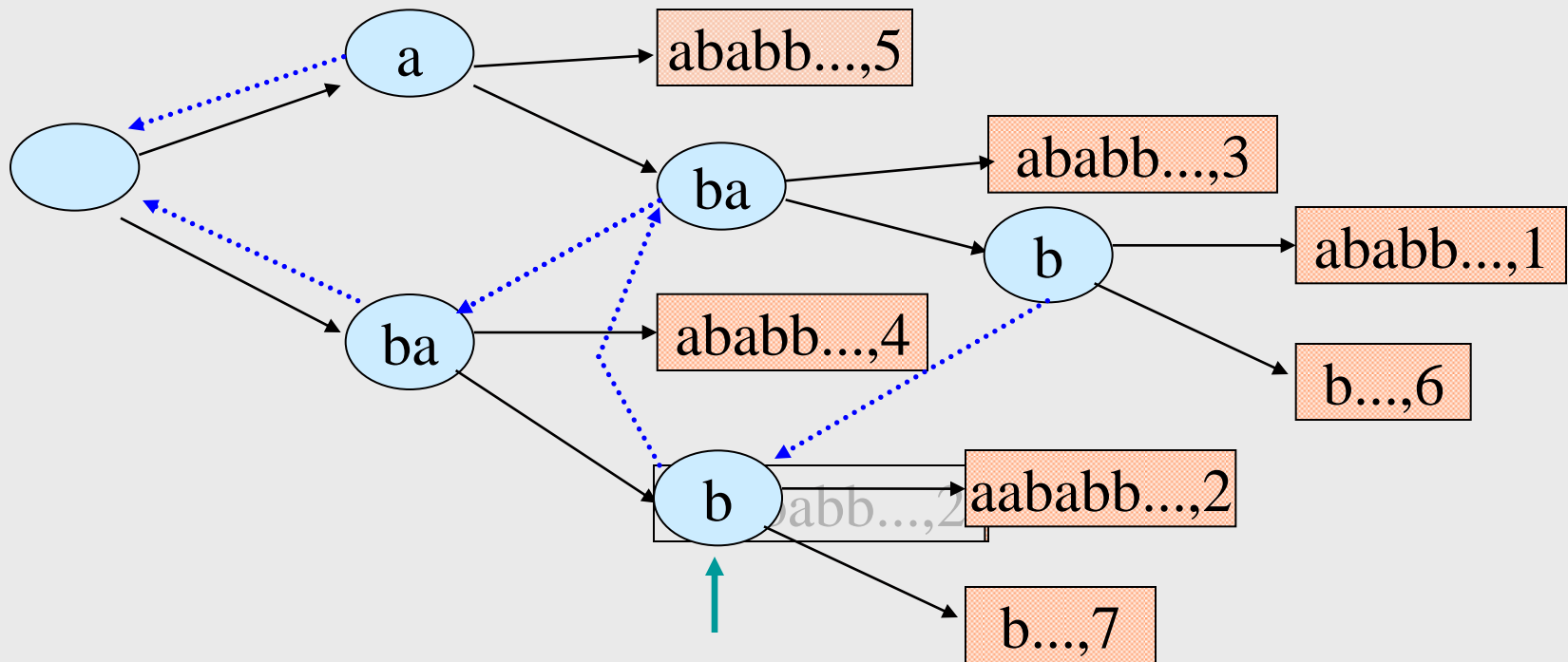
7 8 9



Linear insertion algorithm: example

Given the string $ababaababb\dots$

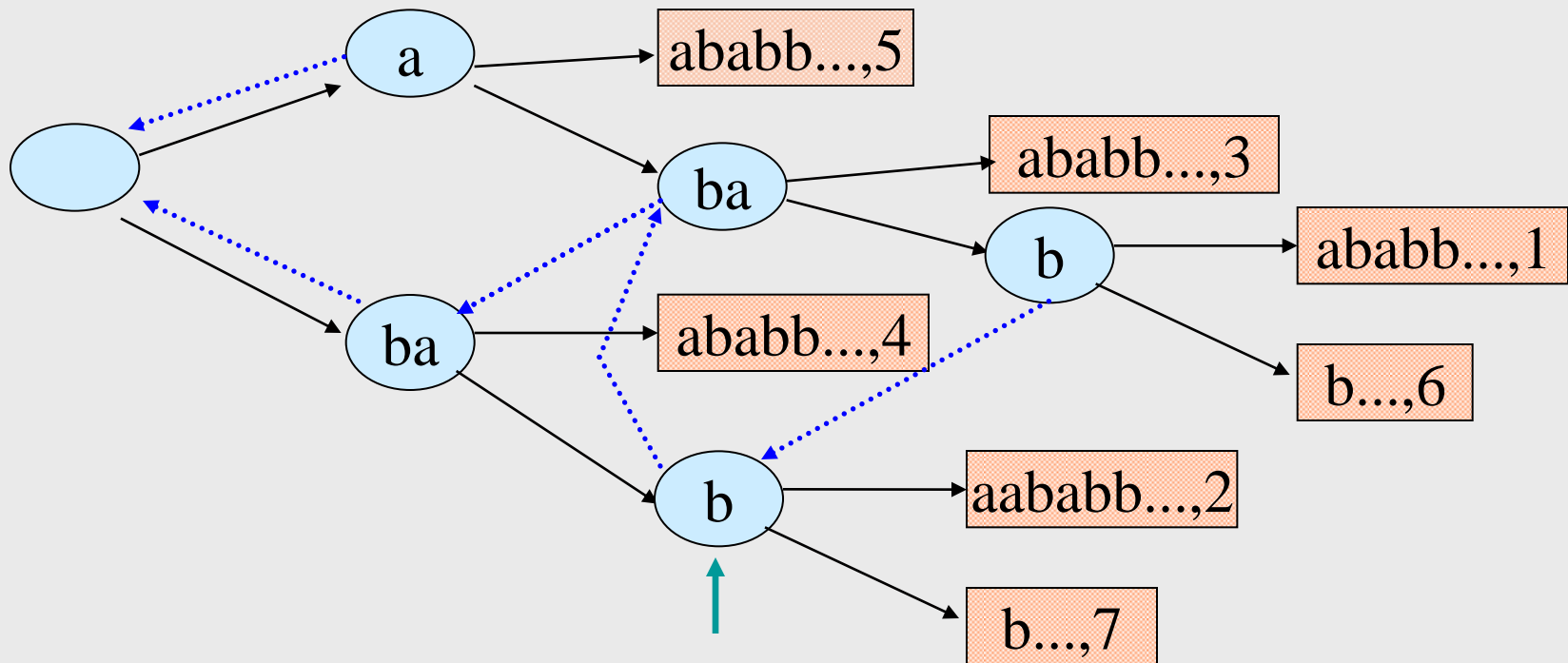
α β
7 8... \uparrow



Linear insertion algorithm: example

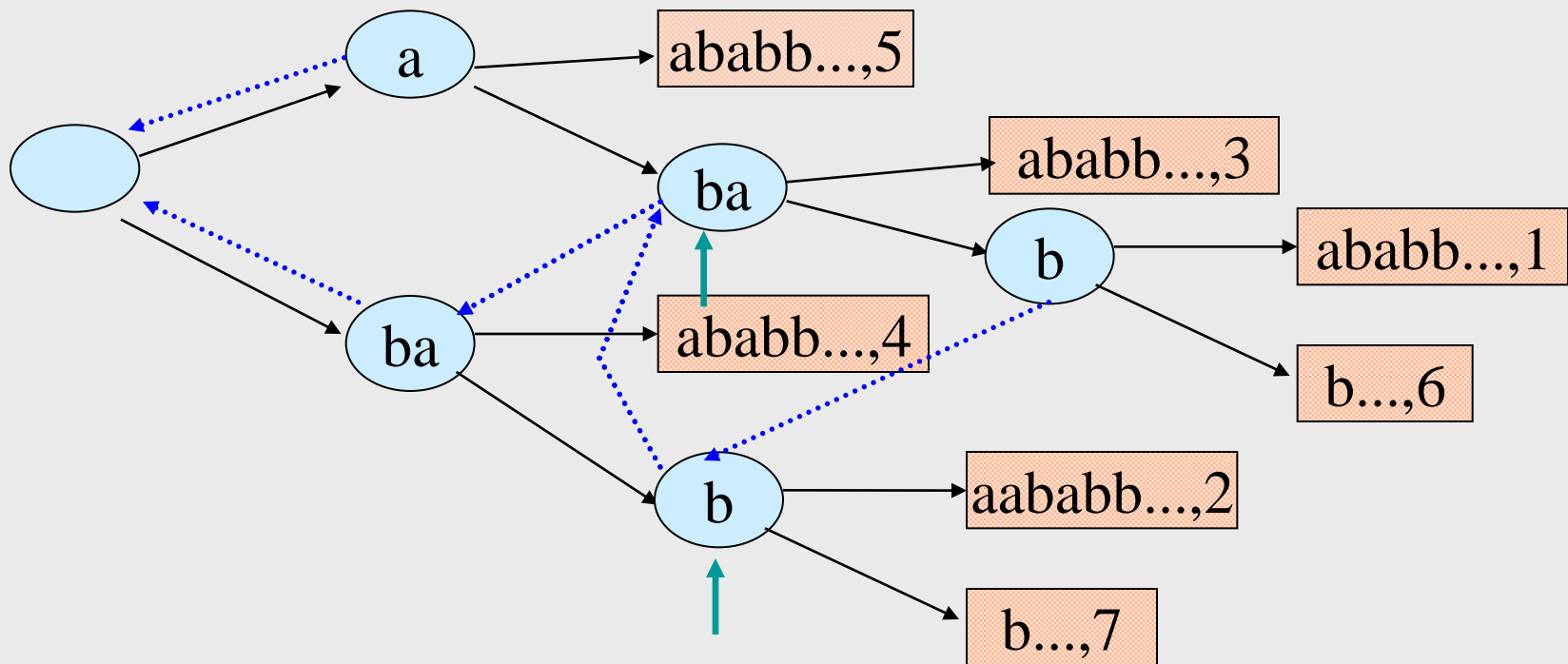
Given the string $ababaababb\dots$

α β 89



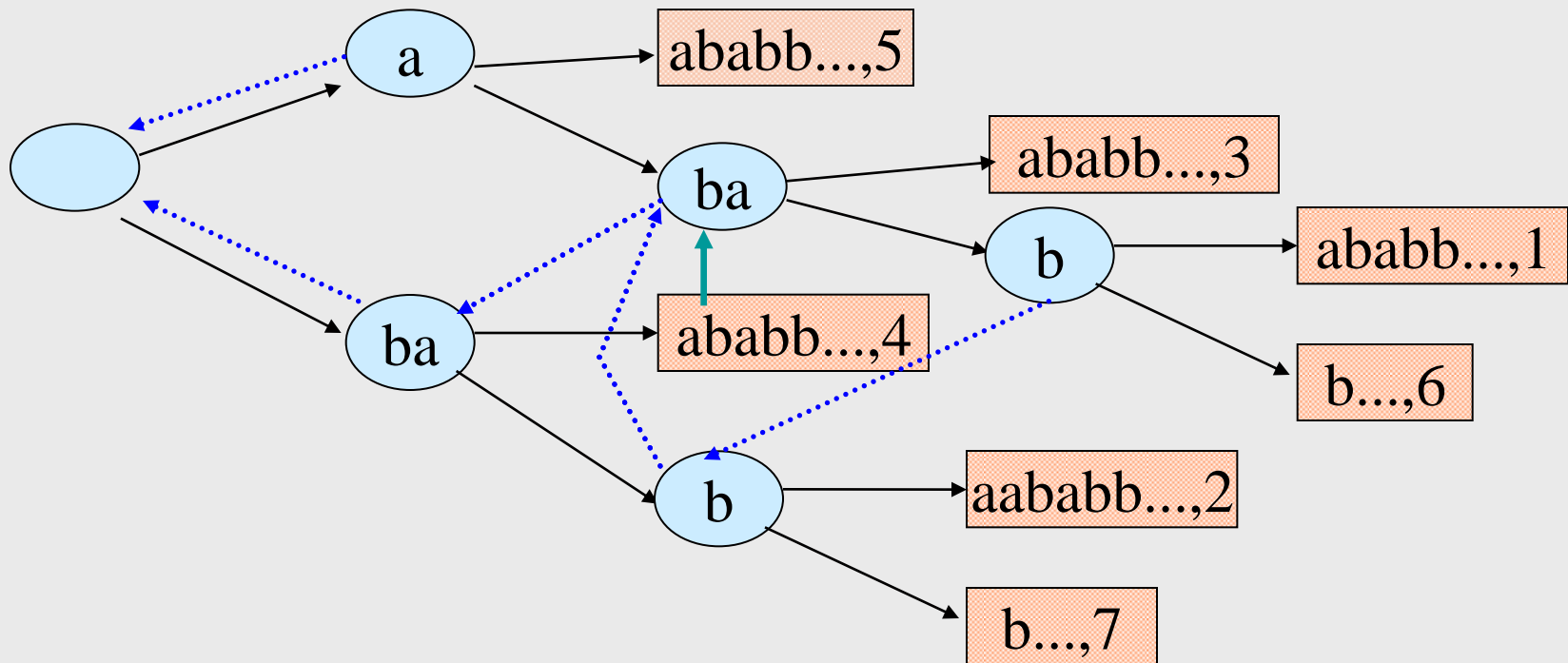
Linear insertion algorithm: example

Given the string $ababaabb\alpha\beta\ldots$



Linear insertion algorithm: example

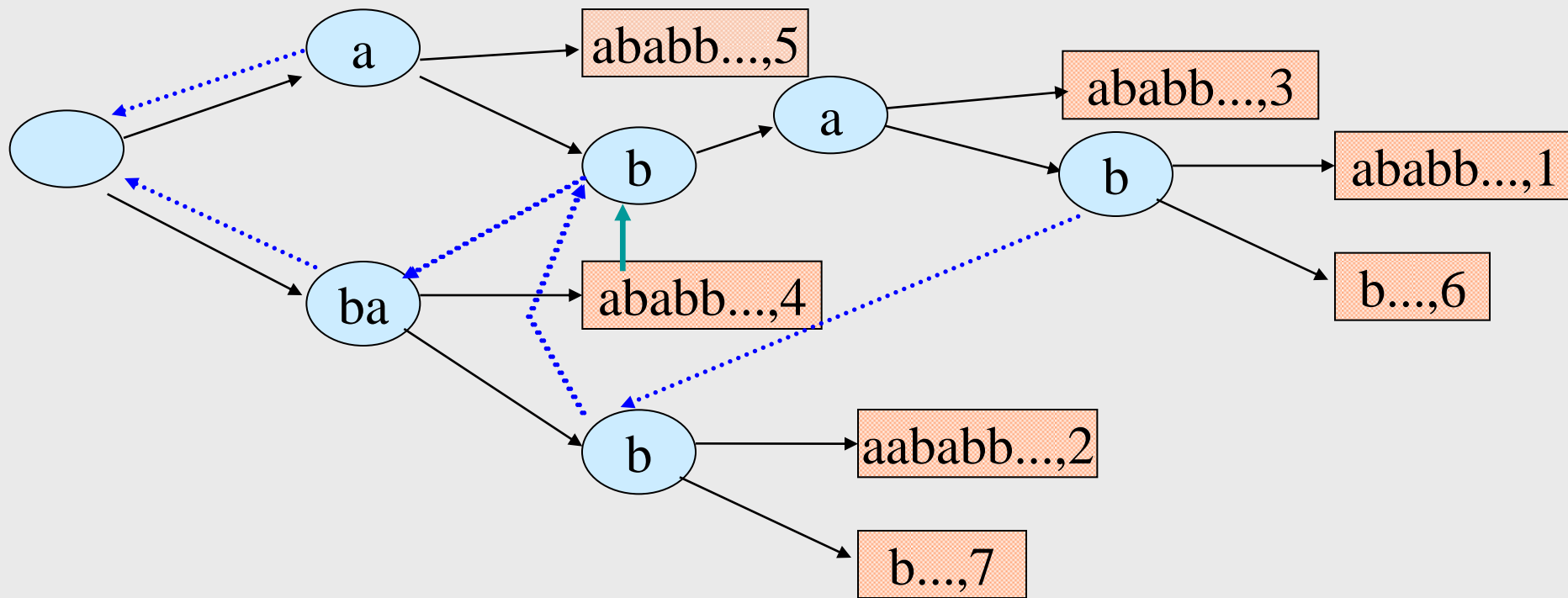
Given the string $ababaababb\dots$



Linear insertion algorithm: example

Given the string $ababaababb\dots$

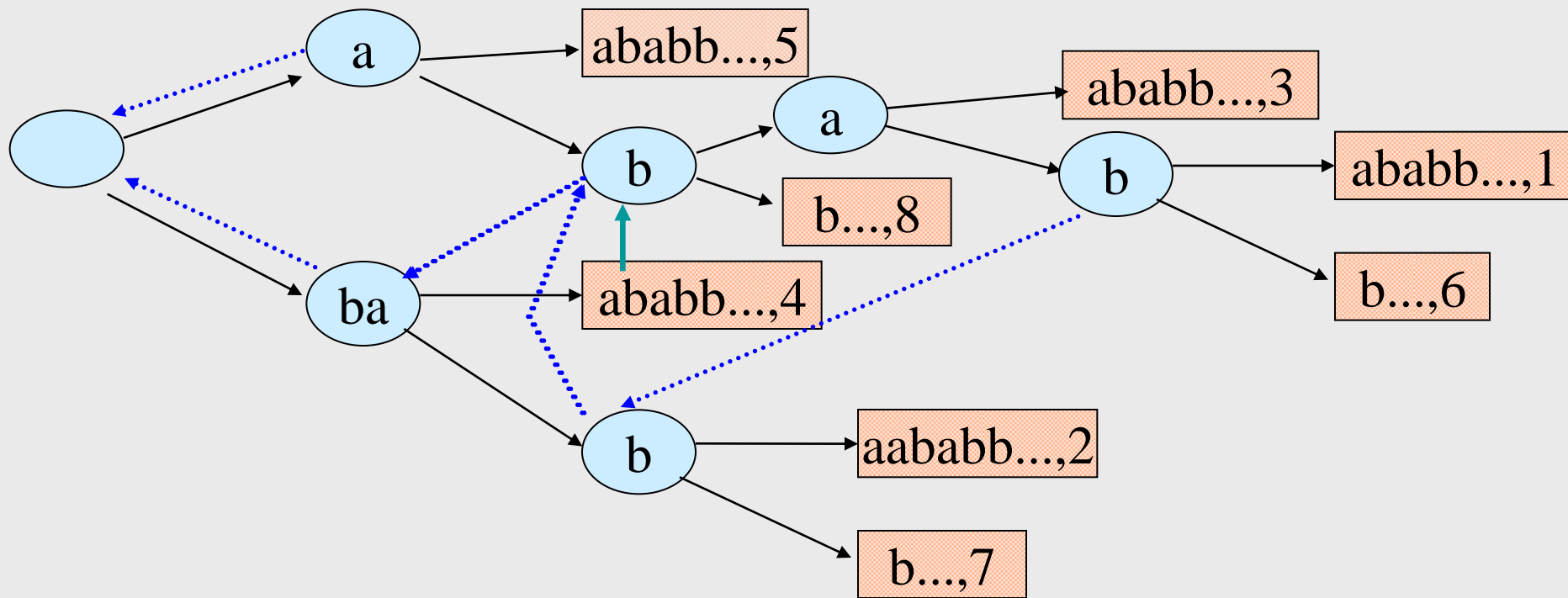
α β 89



Linear insertion algorithm: example

Given the string $ababaababb\dots$

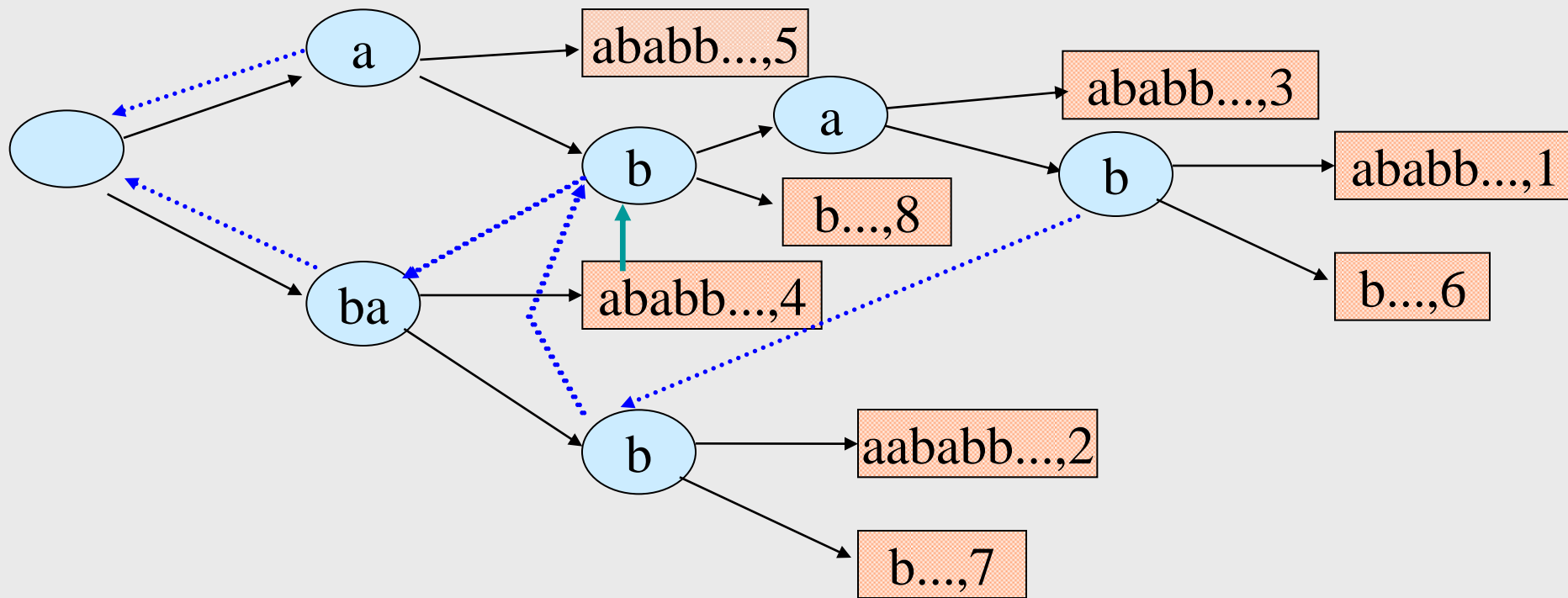
α β 89



Linear insertion algorithm: example

Given the string $ababaababb\dots$

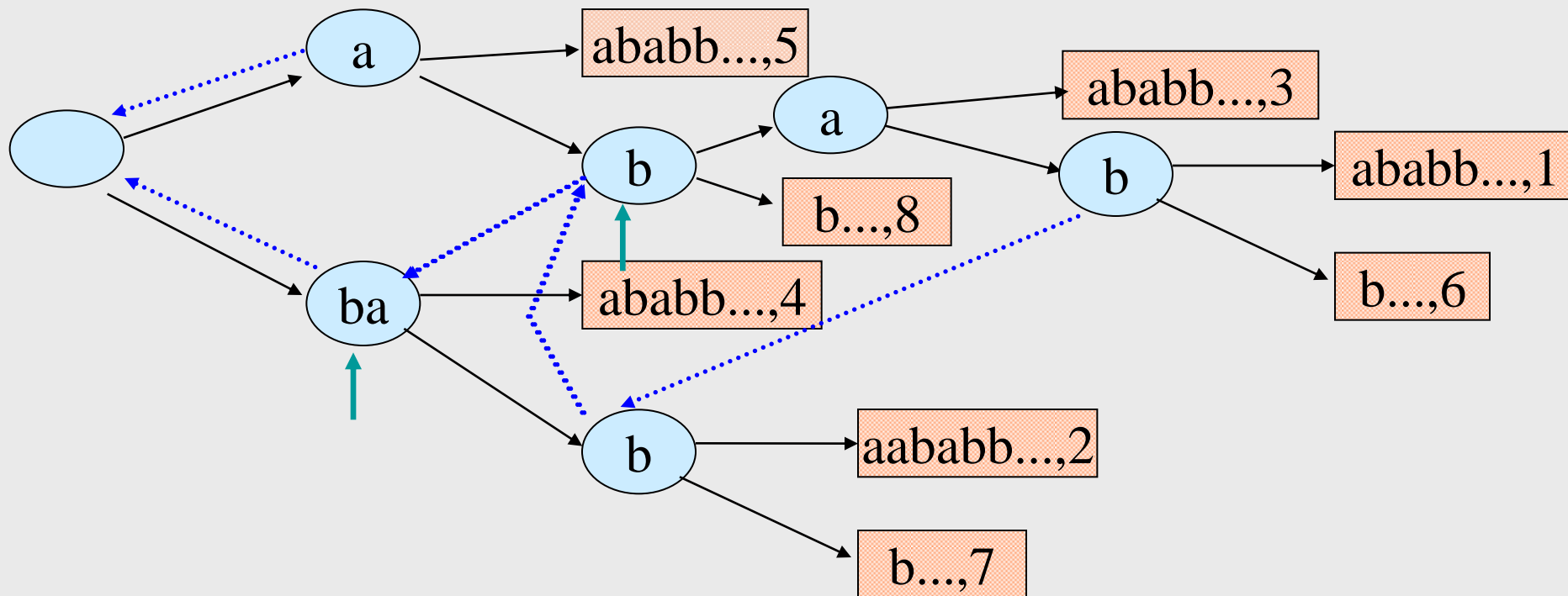
α β 9



Linear insertion algorithm: example

Given the string $ababaababb\dots$

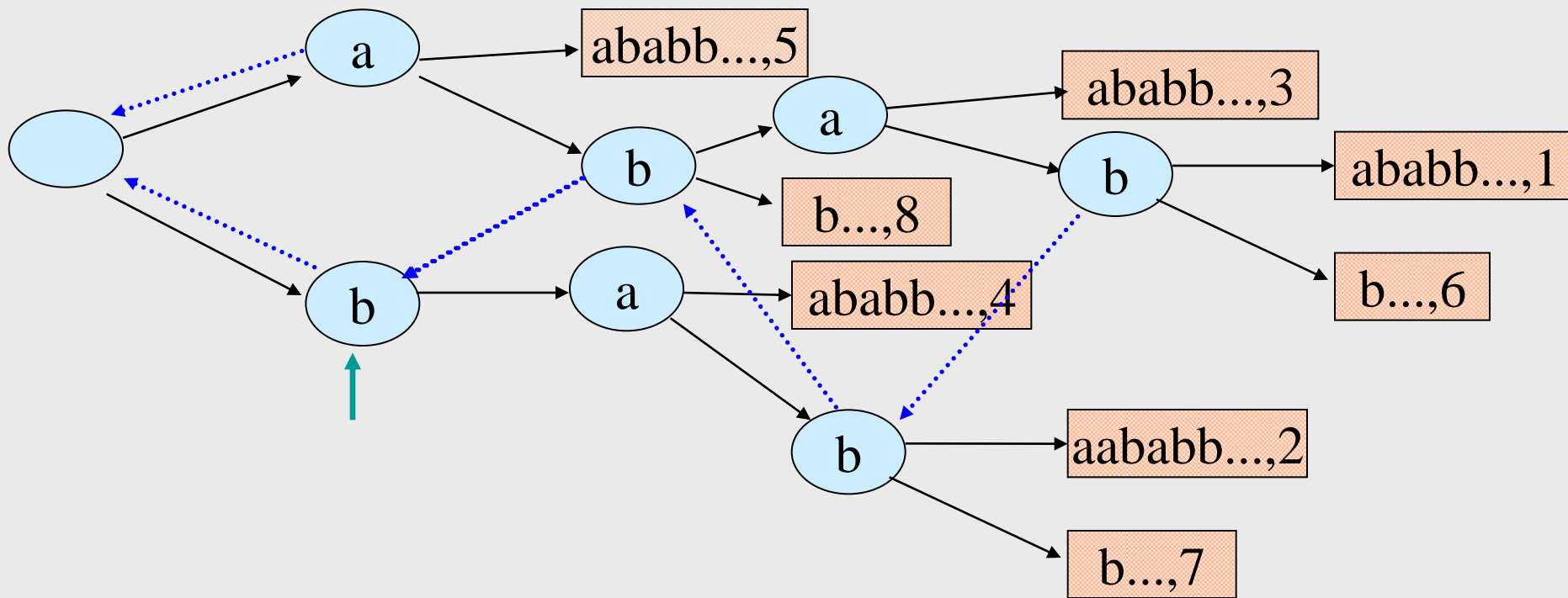
α β 9



Linear insertion algorithm: example

Given the string $ababaababb\dots$

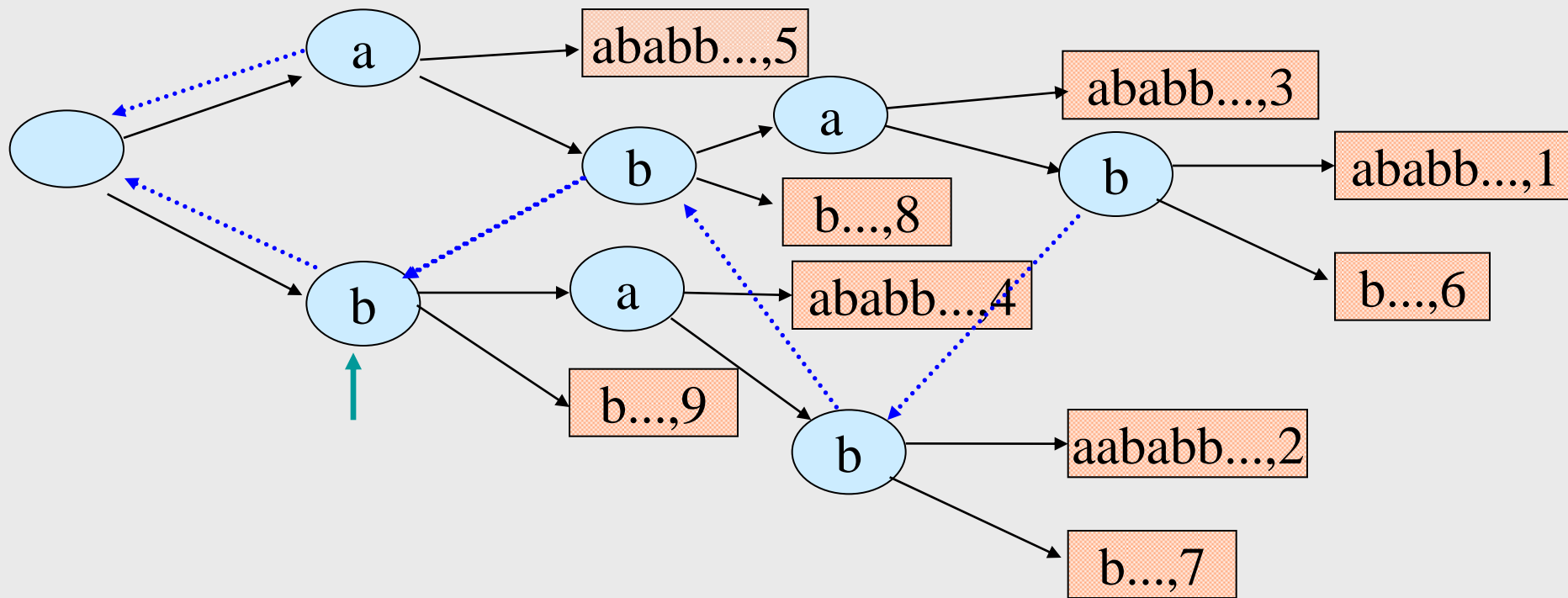
α β 9



Linear insertion algorithm: example

Given the string $\underbrace{ababaababb}_{\alpha} \dots$

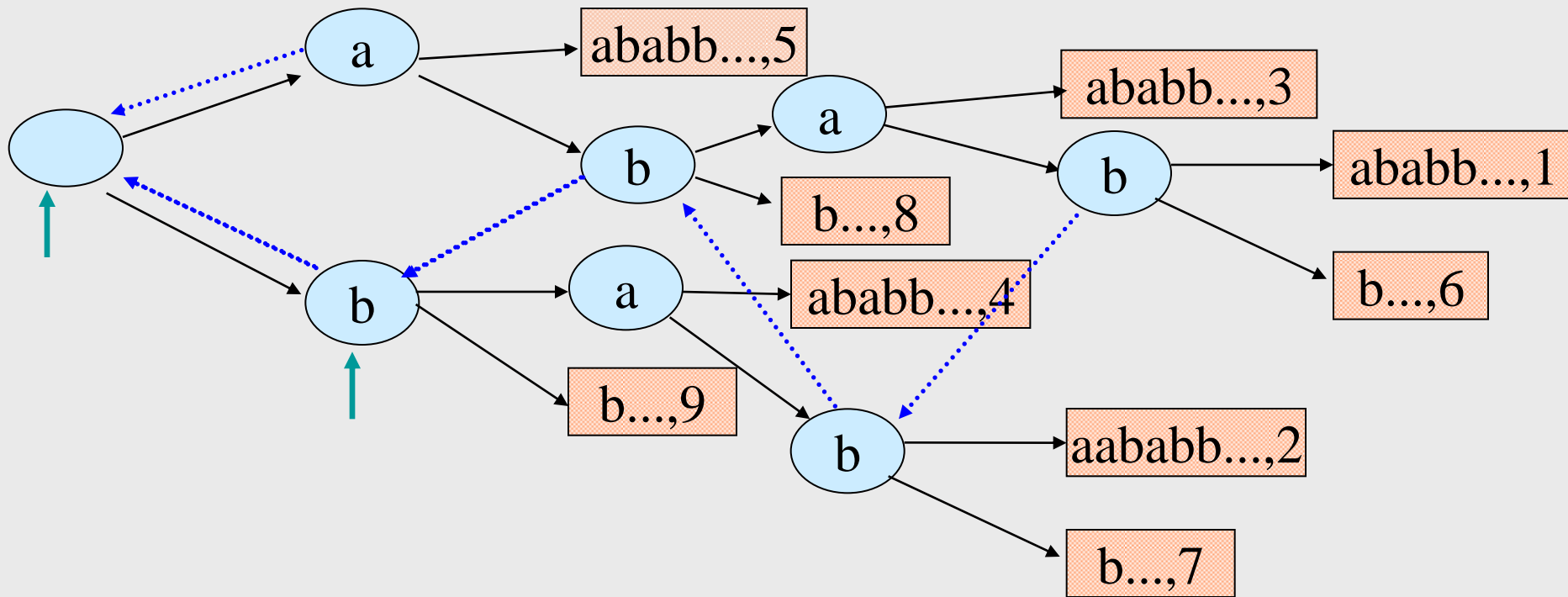
$\downarrow \beta$
 $\uparrow 9$



Linear insertion algorithm: example

Given the string $ababaababb\dots$

$\underbrace{\hspace{10em}}_{\alpha}$ \uparrow
9

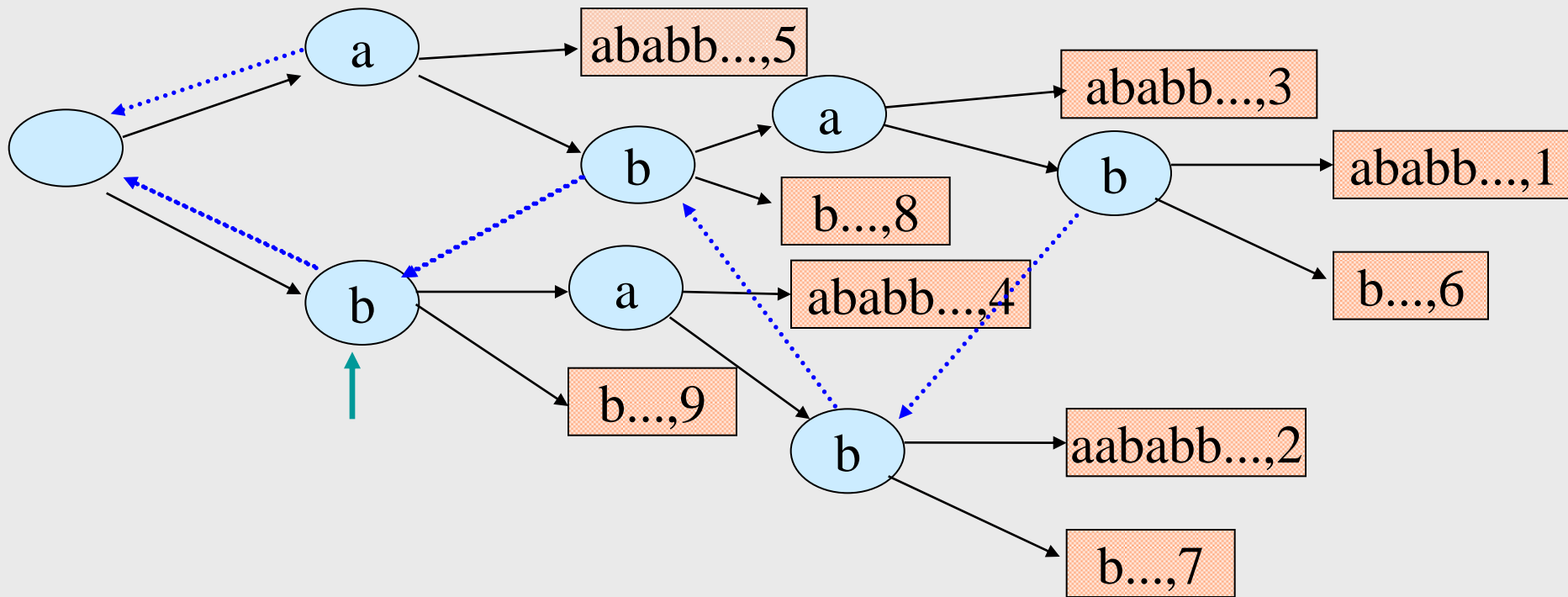


Linear insertion algorithm: example

Given the string $ababaababb\dots$

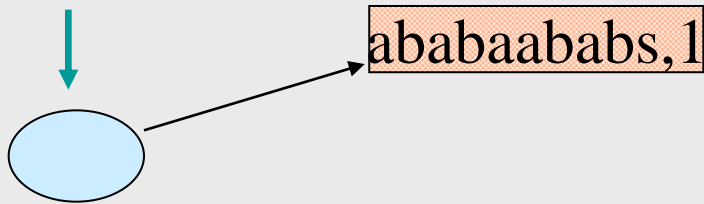
α β

9



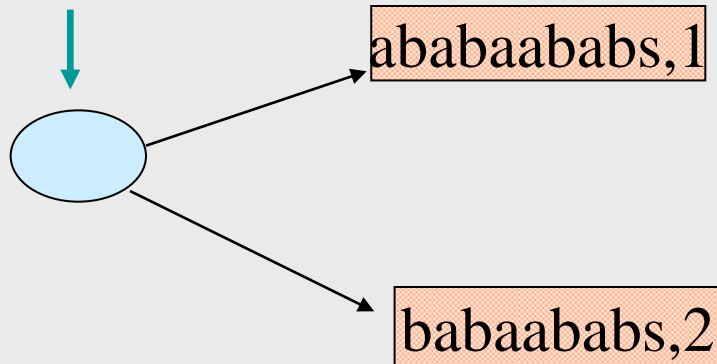
Linear insertion algorithm

Given the string ababaababs



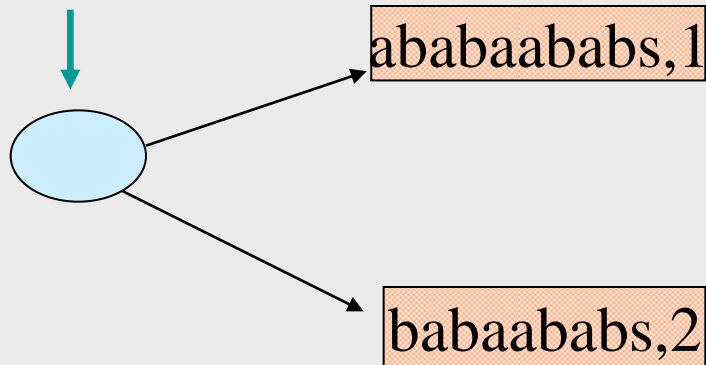
Linear insertion algorithm

Given the string **ababaababs**



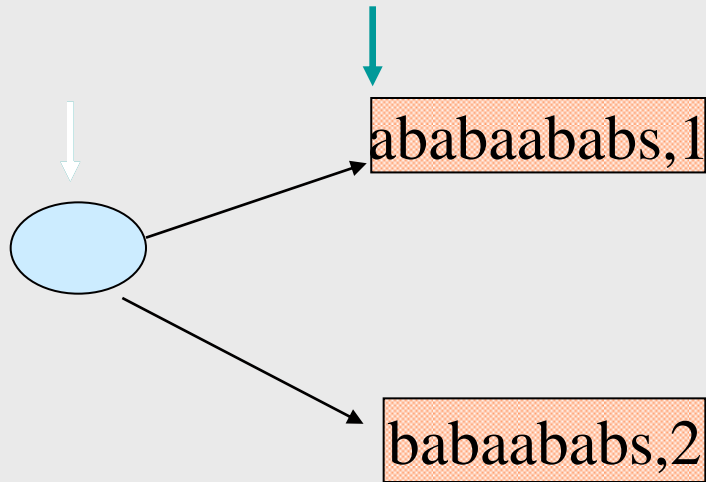
Linear insertion algorithm

Given the string **ababaababs**



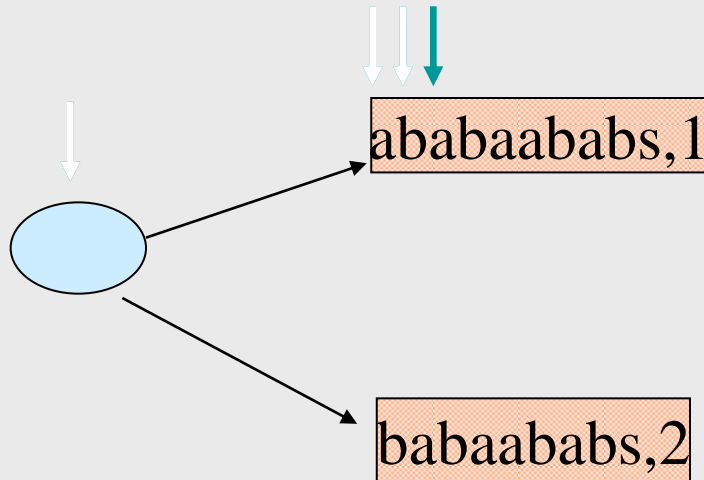
Linear insertion algorithm

Given the string ababaababs



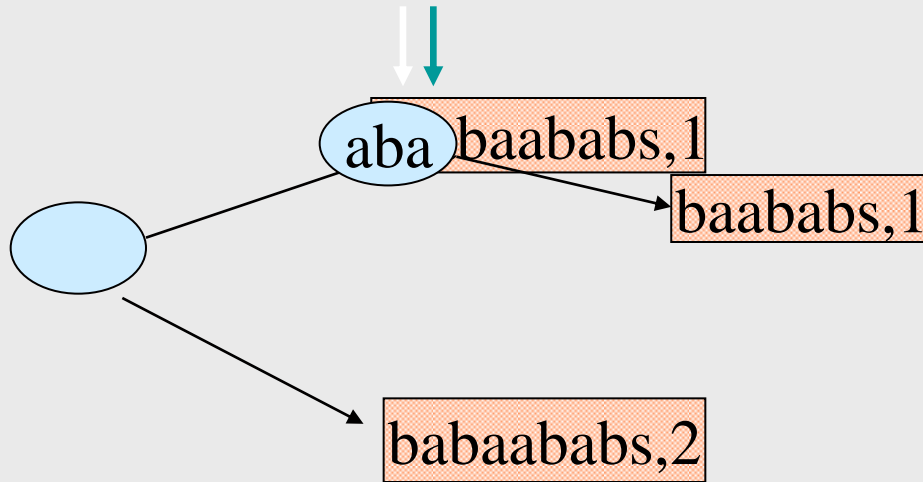
Linear insertion algorithm

Given the string **ababaababs**



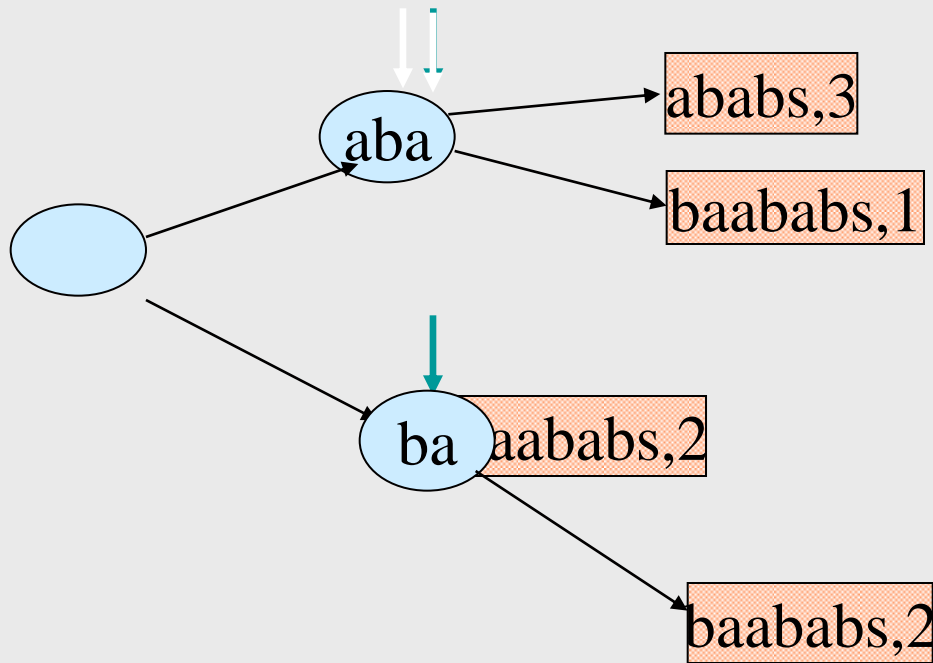
Linear insertion algorithm

Given the string **ababaababs**



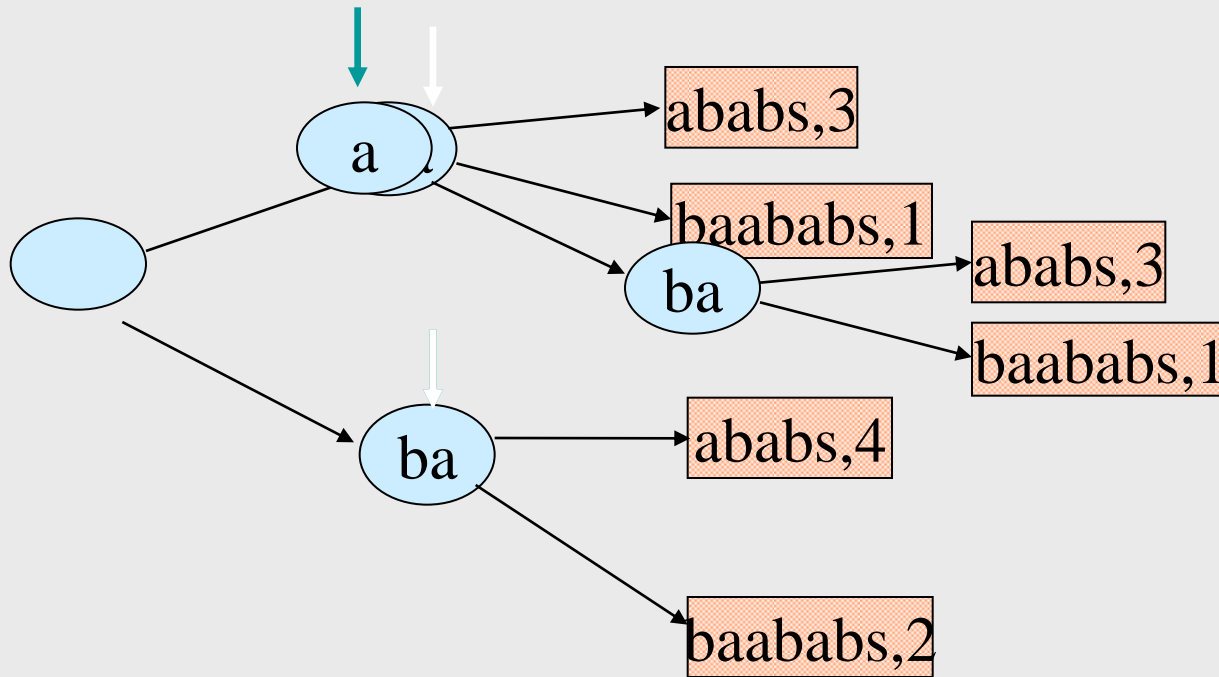
Linear insertion algorithm

Given the string **ababaababs**



Linear insertion algorithm

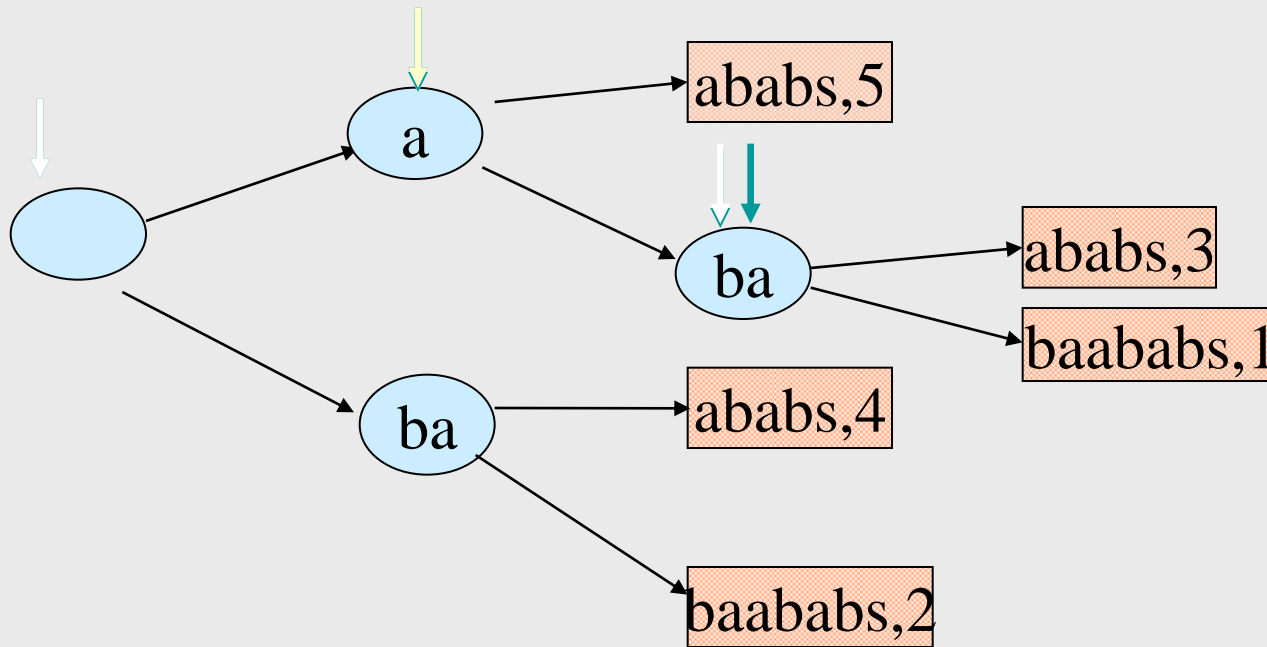
Given the string **ababaababs**



Linear insertion algorithm

Given the string **ababaababs**

↓ ↓ ↓
↑ ↑ ↑ ↑ ↑



1a. Part: Suffix trees

Algorithms on strings, trees and sequences,
Dan Gusfield
Cambridge University Press

2a. Part: Suffix arrays

Suffix-arrays: a new method for on-line
string searches,
G. Myers, U. Manber

Given string **ababaa#**:

Suffixes: **1: ababaa# ...** but lexicographically sorted

2: babaa#

3: abaa#

4: baa#

5: aa#

6: a#

7: #

1	1: #
2	6: a#
3	5: aa#
4	3: abaa#
5	1: ababaa#
6	4: baa#
7	2: babaa#

Which is the cost? $O(n \log(n))$

1. Exact string matching

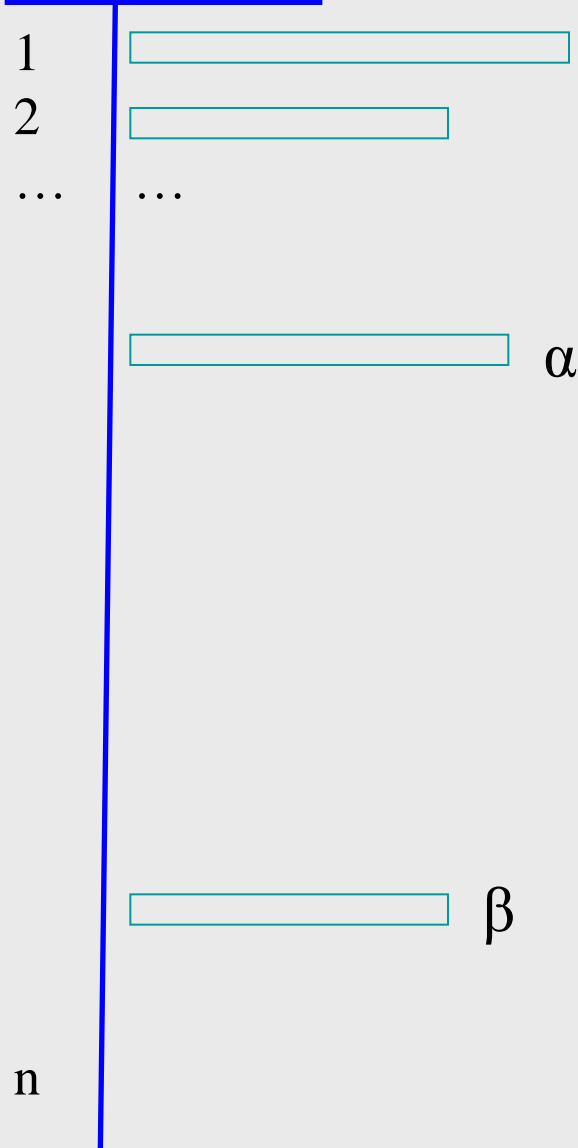
- Does the sequence **ababaas** contain any occurrence of patterns **abab**, **aab**, and **ab**?

1	1: #
2	6: a#
3	5: aa#
4	3: abaa#
5	1: ababaa#
6	4: baa#
7	2: babaa#

Binary search

Search with cost $O(\log(n) |P|)$

Suffix array



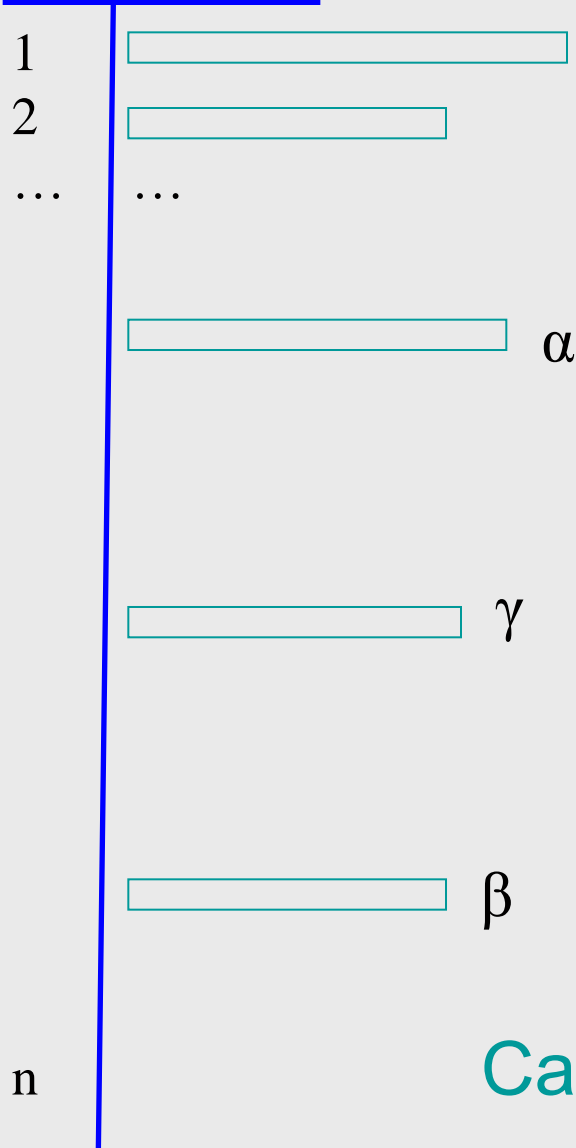
Query: ████████████████████

Invariant Properties:

$$P1: \alpha < \text{query} \leq \beta$$

Search with cost $O(\log(n) |P|)$

Suffix array



Query: ████████████████████

Invariant Properties:

P1: $\alpha < \text{query} \leq \beta$

Algorithm:

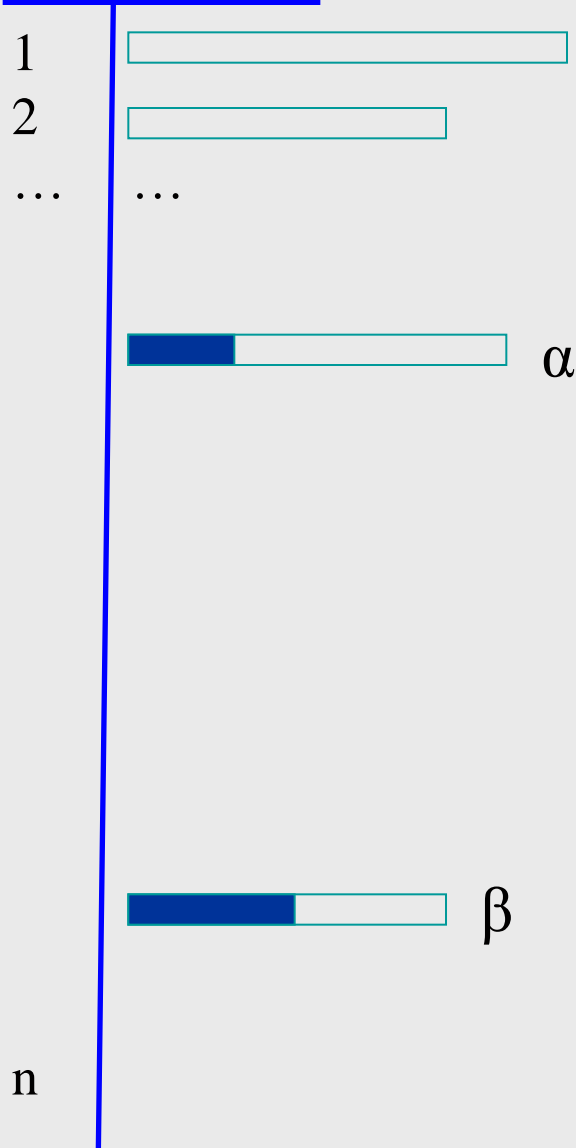
If $\gamma < \text{query}$ **then** $\alpha = \gamma$
else $\beta = \gamma$

Cost: $O(\log(n) |P|)$

Can it be improved to ... $O(\log(n) + |P|)$?

Fast search with cost $O(\log(n)+|P|)$

Suffix array



Query: 

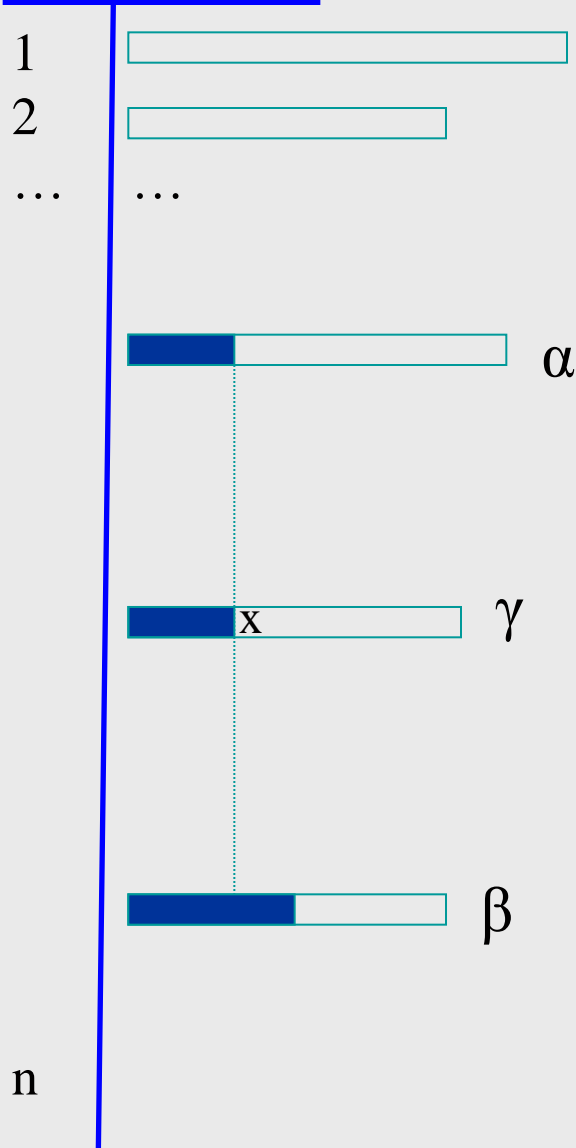
Invariant Properties:

P1: $\alpha < \text{query} \leq \beta$

P2:  matches $\text{pref}(\text{query})$

Fast search with cost $O(\log(n)+|P|)$

Suffix array



Query: $\blacksquare y \blacksquare$

Invariant Properties:

P1: $\alpha < \text{query} \leq \beta$

P2: \blacksquare matches $\text{pref}(\text{query})$

Algorithm:

if

$\mathbf{x < y}$ **then** $\alpha = \gamma$

$\mathbf{x > y}$ **then** $\beta = \gamma$

$\mathbf{x = y}$ **then** ...

fi