

# Sparse Suffix Sorting

Tomohiro I

# Suffix sorting

- Sorting suffixes of a string  $T$  is a fundamental task having a lot of applications such as text indexing.

# Suffix sorting

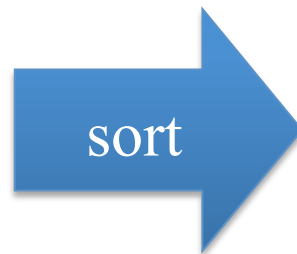
- Sorting suffixes of a string  $T$  is a fundamental task having a lot of applications such as text indexing.

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
 $T =$  l i o n l y i n g o n l y o n l y o n \$

# Suffix sorting

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
 $T = \text{l i o n l y i n g o n l y o n l y o n \$}$

lonlyingonlyonlyon\$  
ionlyingonlyonlyon\$  
onlyingonlyonlyon\$  
nlyingonlyonlyon\$  
lyingonlyonlyon\$  
yingonlyonlyon\$  
ingonlyonlyon\$  
ngonlyonlyon\$  
gonlyonlyon\$  
onlyonlyon\$  
nlyonlyon\$  
lyonlyon\$  
yonlyon\$  
onlyon\$  
nlyon\$  
lyon\$  
yon\$  
on\$  
n\$  
\$



20 \$  
9 gonlyonlyon\$  
7 ingonlyonlyon\$  
2 ionlyingonlyonlyon\$  
1 lionlyingonlyonlyon\$  
5 lyingonlyonlyon\$  
16 lyon\$  
12 lyonlyon\$  
19 n\$  
8 ngonlyonlyon\$  
4 nlyingonlyonlyon\$  
15 nlyon\$  
11 nlyonlyon\$  
18 on\$  
3 onlyingonlyonlyon\$  
14 onlyon\$  
10 onlyonlyon\$  
6 yingonlyonlyon\$  
17 yon\$  
13 yonlyon\$

all suffixes of  $T$

# Suffix sorting

Suffix Array of  $T$

1 2 3 4 5 6 7 8 9 10 11 12 13  
 $T = \text{lionlyingonlyyon\$}$

lionlyingonlyyon\$  
ionlyingonlyyon\$  
onlyingonlyyon\$  
nlyingonlyyon\$  
lyingonlyyon\$  
yingonlyyon\$  
ingonlyyon\$  
ngonlyyon\$  
gonlyonlyyon\$  
onlyonlyyon\$  
nlyonlyyon\$  
lyonlyyon\$  
yonlyyon\$  
onlyon\$  
nlyon\$  
lyon\$  
yon\$  
on\$  
n\$  
\$

sort

20 \$  
9 gonlyonlyyon\$  
7 ingonlyonlyyon\$  
2 ionlyingonlyonlyyon\$  
1 lionlyingonlyonlyyon\$  
5 lyingonlyonlyyon\$  
16 lyon\$  
12 lyonlyyon\$  
19 n\$  
8 ngonlyonlyyon\$  
4 nlyingonlyonlyyon\$  
15 nlyon\$  
11 nlyonlyyon\$  
18 on\$  
3 onlyingonlyonlyyon\$  
14 onlyon\$  
10 onlyonlyyon\$  
6 yingonlyonlyyon\$  
17 yon\$  
13 yonlyyon\$

all suffixes of  $T$

# Suffix sorting

Suffix Array of  $T$

1 2 3 4 5 6 7 8 9 10 11 12 13  
 $T = \text{l i o n l y i n g o n l y o n l y o n } \$$

l ionly ingonly yonly on\$  
 i only ingonly yonly on\$  
 o nly ingonly yonly on\$  
 n l y ingonly yonly on\$  
 l y ingonly yonly on\$  
 y ingonly yonly on\$  
 i ngonly yonly on\$  
 n gonly yonly on\$  
 g only yonly on\$  
 o nly yonly on\$  
 n l y only on\$  
 l y only on\$  
 y only on\$  
 o nly on\$  
 n l y on\$  
 l y on\$  
 y on\$  
 o n\$  
 n\$  
 \$

20 \$  
 9 gonly yonly on\$  
 7 ingonly yonly on\$  
 2 ionly ingonly yonly on\$  
 1 lionly ingonly yonly on\$  
 5 l y ingonly yonly on\$  
 16 l y on\$  
 12 l y only on\$  
 19 n\$  
 8 n gonly yonly on\$  
 4 n l y ingonly yonly on\$  
 15 n l y on\$  
 11 n l y only on\$  
 18 o n\$  
 3 o nly ingonly yonly on\$  
 14 o nly on\$  
 10 o nly yonly on\$  
 6 y ingonly yonly on\$  
 17 y on\$  
 13 y only on\$



# Suffix sorting

Suffix Array of  $T$

1 2 3 4 5 6 7 8 9 10 11 12 13  
 $T = \text{l i o n l y i n g o n l y o n l y o n } \$$

l ionlyingonlylyon\$  
 ionlyingonlylyon\$  
 onlyingonlylyon\$  
 nlyingonlylyon\$  
 lyingonlylyon\$  
 yingonlylyon\$  
 ingonlylyon\$  
 ngonlylyon\$  
 gonlylyon\$  
 onlyonlylyon\$  
 nlyonlylyon\$  
 lyonlylyon\$  
 yonlylyon\$  
 onlyon\$  
 nlyon\$  
 lyon\$  
 yon\$  
 on\$  
 n\$  
 \$



20	\$
9	gonlyonlyon\$
7	ingonlyonlyon\$
2	ionlyingonlyonlyon\$
→ 1	lionlyingonlyonlyon\$
5	lyingonlyonlyon\$
→ 16	lyon\$
12	lyonlyon\$
19	n\$
→ 8	ngonlyonlyon\$
4	nlyingonlyonlyon\$
15	nlyon\$
11	nlyonlyon\$
18	on\$
3	onlyingonlyonlyon\$
14	onlyon\$
10	onlyonlyon\$
6	yingonlyonlyon\$
17	yon\$
13	yonlyon\$





# Sparse suffix sorting

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
 $T =$  l i o n l y i n g o n l y o n l y o n \$

l i o n l y i n g o n l y o n l y o n \$

i o n l y i n g o n l y o n l y o n \$

o n l y i n g o n l y o n l y o n \$

n l y i n g o n l y o n l y o n \$

l y i n g o n l y o n l y o n \$

y i n g o n l y o n l y o n \$

i n g o n l y o n l y o n \$

n g o n l y o n l y o n \$

g o n l y o n l y o n \$

o n l y o n l y o n \$

n l y o n l y o n \$

l y o n l y o n \$

y o n l y o n \$

o n l y o n \$

n l y o n \$

l y o n \$

y o n \$

o n \$

n \$

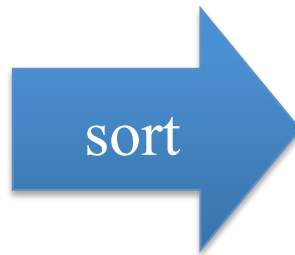
\$

If we are interested in some (small) set of suffixes of  $T$ .

# Sparse suffix sorting

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
 $T = l i o n l y i n g o n l y o n l y o n \$$

l ionly ingonly only on \$  
i only ingonly only on \$  
o nly ingonly only on \$  
n l y ingonly only on \$  
l y ingonly only on \$  
y ingonly only on \$  
i ngonly only on \$  
n gonly only on \$  
g only only on \$  
o nly only on \$  
n l y only on \$  
l y only on \$  
y only on \$  
o nly on \$  
n l y on \$  
l y on \$  
y on \$  
o n \$  
n \$  
\$



Sparse Suffix Array (SSA)

1	l ionly ingonly only on \$
12	l y only on \$
18	o n \$
14	o nly on \$
10	o nly only on \$

For sparse indexing, the sparse suffix array is enough, which is much smaller than full suffix array.

If we are interested in some (small) set of suffixes of  $T$ .

# Problem

Given  $T$  of length  $n$  and set of  $b$  positions, sort the designated suffixes.

- It can be solved in  $O(n)$  time by computing full suffix array and make it sparse, but it needs  $O(n)$  space in addition to the text.
- An alternative way is to use any sorting algorithm for a set of strings (not for suffixes) that uses  $O(b)$  additional space, but it takes  $\Omega(bn)$  time.

## Question

Using  $O(s)$  additional space for some  $s \in [b, n]$ , how fast can we sort sparse suffixes?

# Monte-Carlo Algorithm

Randomized algorithms that may return a wrong output with low probability.

We can win the gamble  
with high probability!!

# Monte-Carlo Algorithm

Randomized algorithms that may return a wrong output with low probability.

We can win the gamble with high probability!!

- Monte-Carlo algorithm using Karp-Rabin fingerprints.
  - ◆ Time-space tradeoffs: for any  $s \in [b, n]$   
 $O(n + (bn/s) \log s)$  time using  $O(s)$  additional space.
  - ◆ When  $s = b$ , the time complexity is  $O(n \log b)$ .
  - ◆ When  $s = b \log b$ , the time complexity is  $O(n)$ .

# Karp-Rabin fingerprints (rolling hash)

[Karp and Rabin, '87]

$T$  

$i$    $j$



$$FP[i..j] = \sum_{k=i}^j T[k] \cdot r^{j-k} \pmod{q}$$

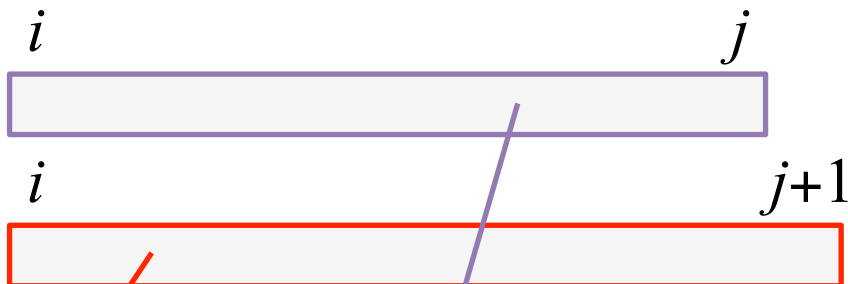
- $q$  is a prime number.
- $r$  is chosen from  $[0, q-1]$  **uniformly at random**.

For sufficiently large  $q$ , the hash function becomes perfect with high probability, i.e.,  
 $T[i..i+l] = T[i'..i'+l]$  iff  $FP[i..i+l] = FP[i'..i'+l]$ .

# Karp-Rabin fingerprints (rolling hash)

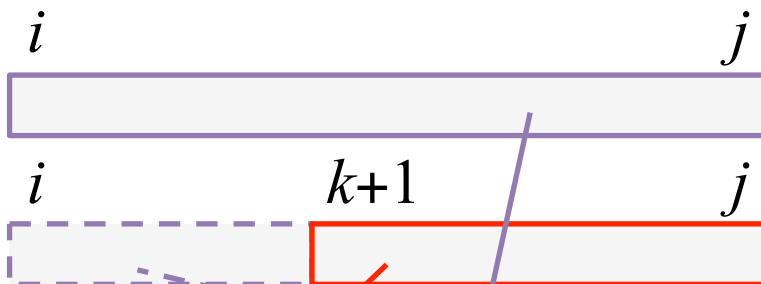
[Karp and Rabin, '87]

$T$  



We can compute  
FPs incrementally.

$$FP[i..j+1] = FP[i..j] \cdot r + T[j+1] \bmod q$$

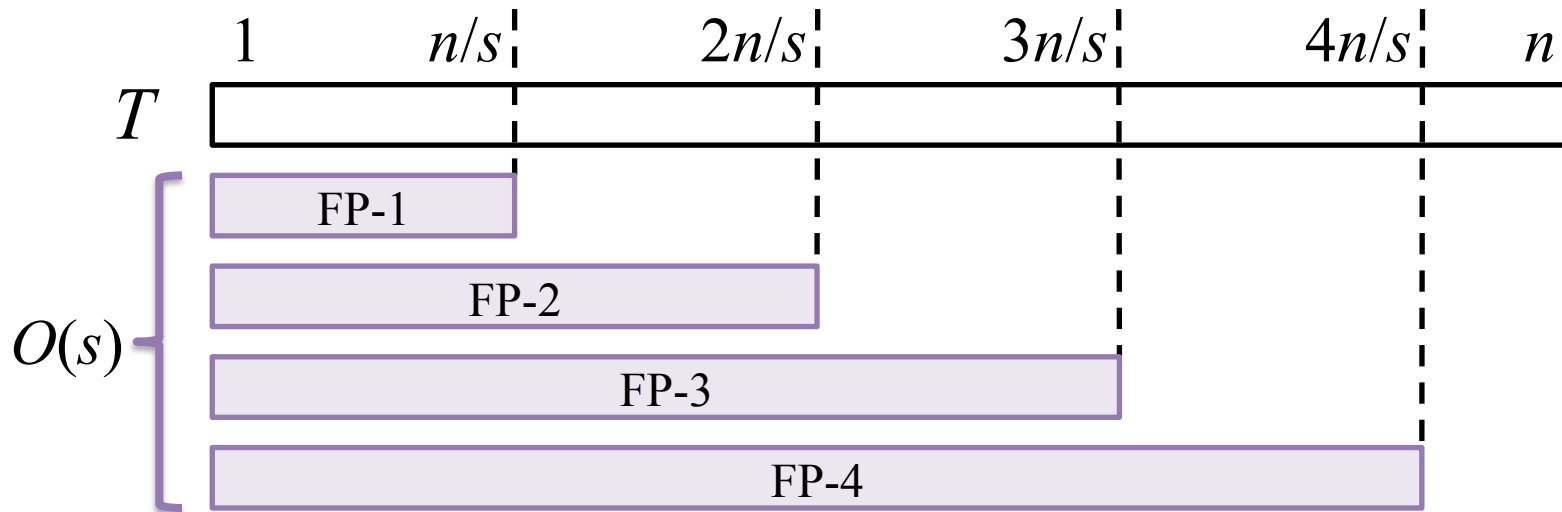


We can “subtract”  
prefix FP’s influence.

$$FP[k+1..j] = FP[i..j] - FP[i..k] \cdot r^{j-k} \bmod q$$

## Lemma

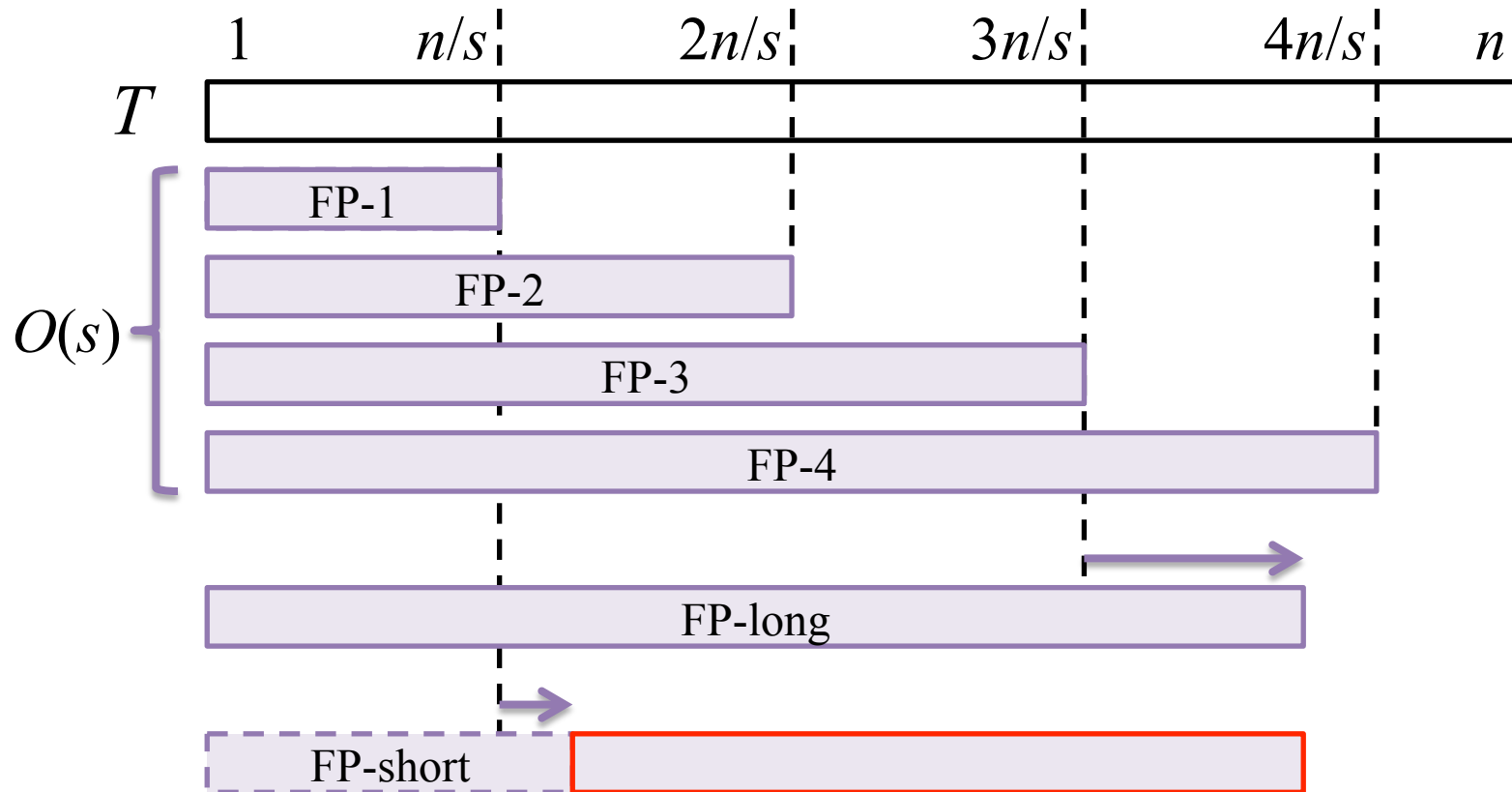
We can preprocess  $T$  in  $O(n)$  time and  $O(s)$  space so that the fingerprint for any substring of length  $l$  can be computed in  $O(\min\{l, n/s\})$  time.





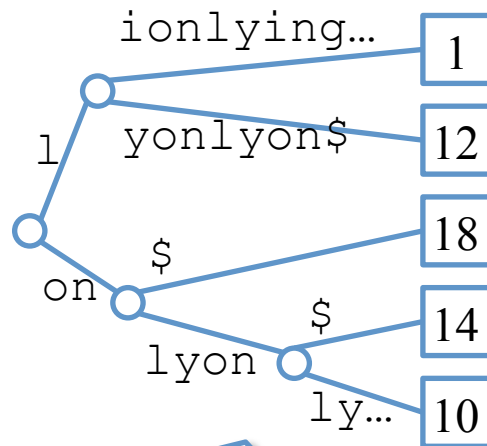
## Lemma

We can preprocess  $T$  in  $O(n)$  time and  $O(s)$  space so that the fingerprint for any substring of length  $l$  can be computed in  $O(\min\{l, n/s\})$  time.



# Monte-Carlo algorithm

- The algorithm constructs sparse suffix tree.



## Sparse Suffix Array

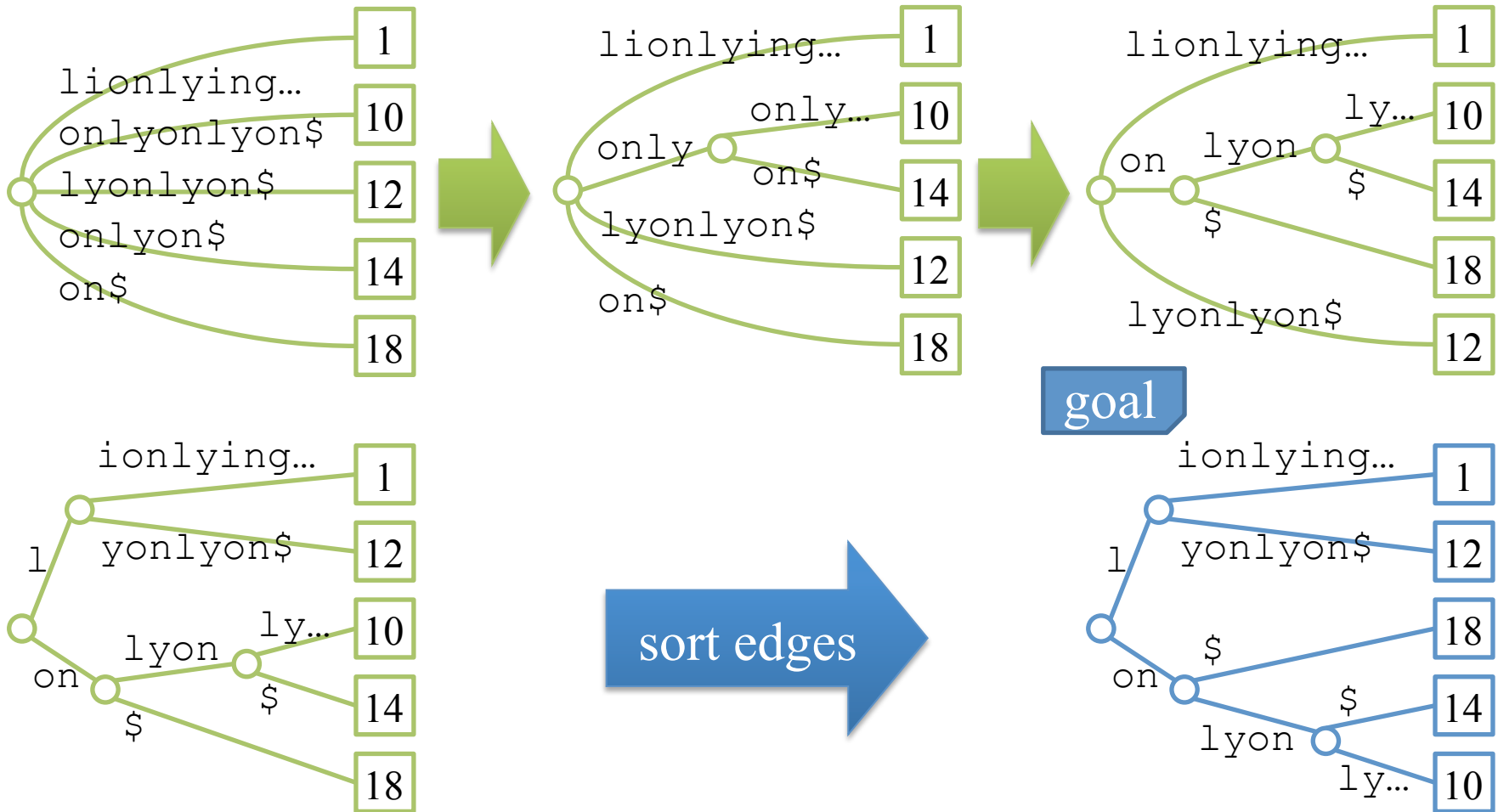
1	lionlyingonlyon\$
12	lyonlyon\$
18	on\$
14	onlyon\$
10	onlyonlyon\$

## Sparse Suffix Tree:

- Compacted trie representing all designated suffixes.
- For each node, child edges are sorted.
- Since edge labels can be encoded by pointers to  $T$ , sparse suffix tree can be represented in  $O(b)$  space.

# Overview: gradual refinement

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
 $T = \text{lionlyingonlyononlyon}\$$

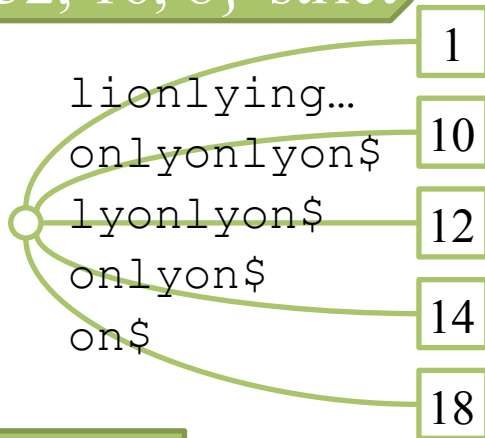


# *l*-strict (unsorted) sparse suffix trees

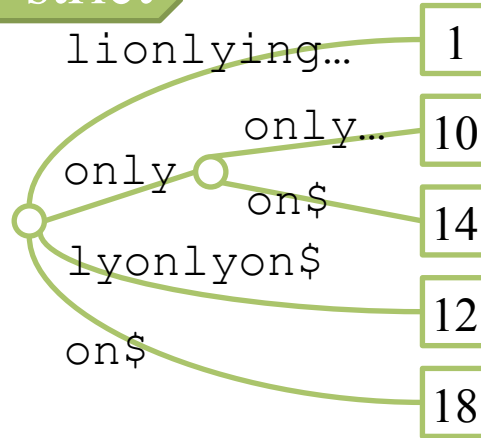
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

$T = \text{lionlyingonlyononlyon}\$$

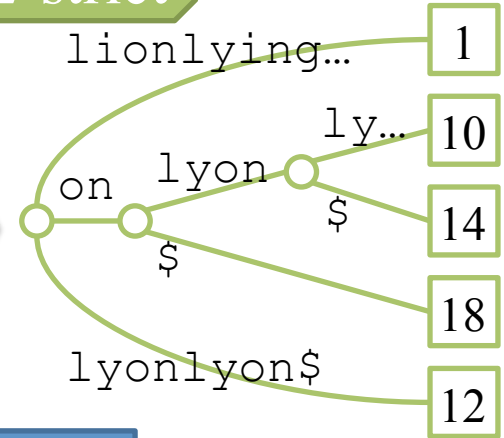
**{32, 16, 8}-strict**



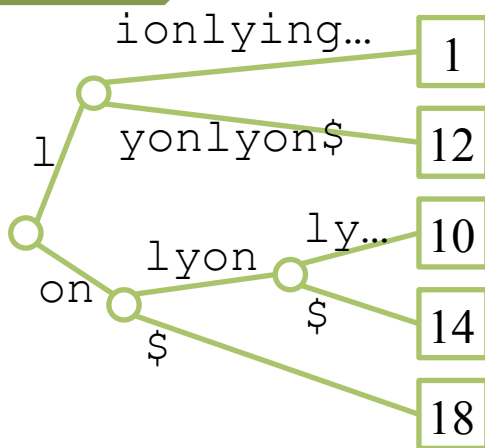
**4-strict**



**2-strict**

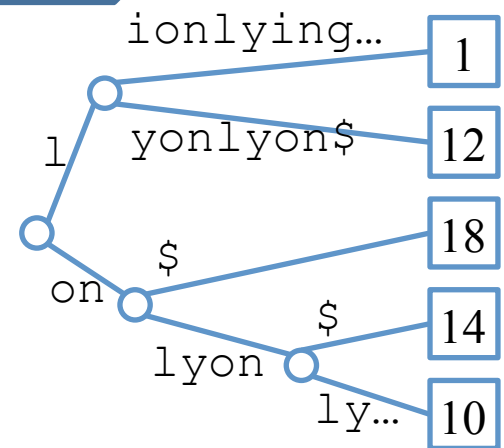


**1-strict**



sort edges

**goal**

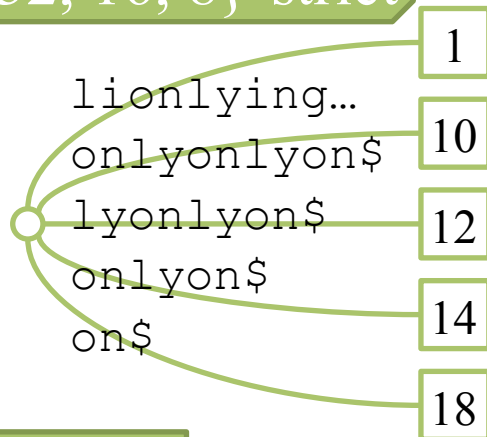


# $l$ -strict (unsorted) sparse suffix trees

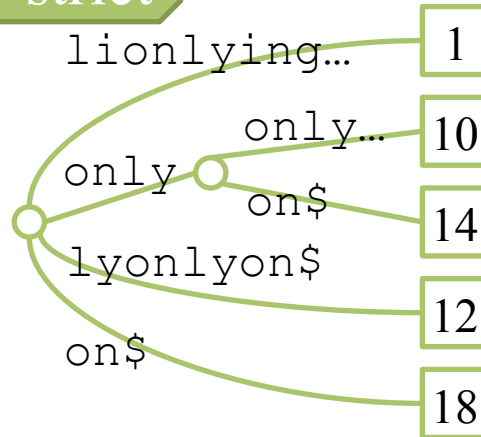
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

$T = \text{lionlyingonlyonlyon\$}$

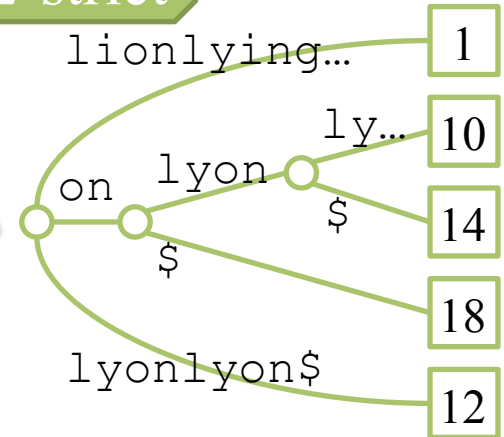
**{32, 16, 8}-strict**



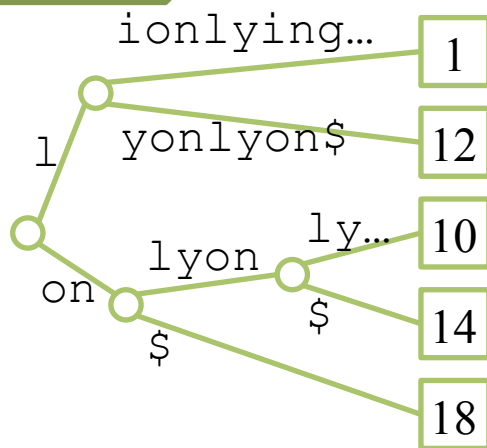
**4-strict**



**2-strict**



**1-strict**



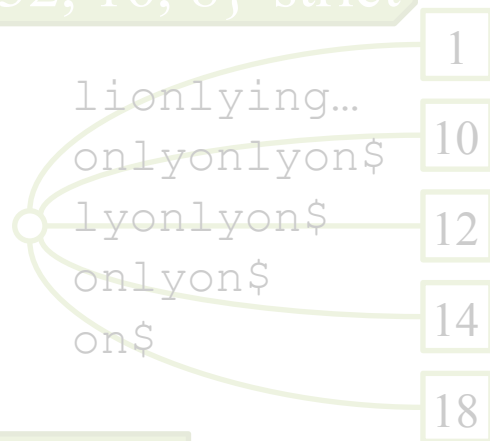
## Condition of $l$ -strict tree

For any internal node of  $l$ -strict tree, any pair of child edge labels has common prefix of length less than  $l$ .

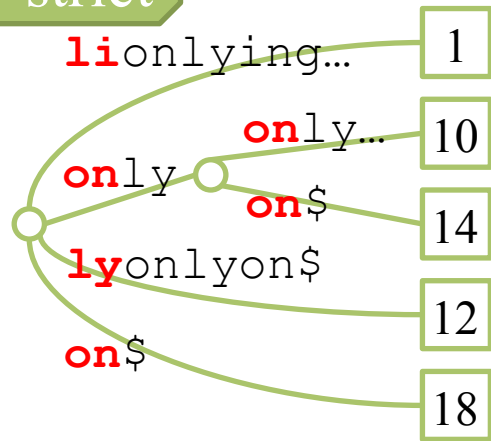
# From 4-strict tree to 2-strict tree

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
*T* = l i o n l y i n g o n l y o n l y o n \$

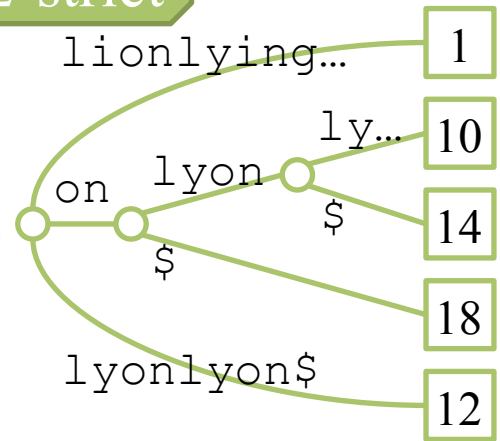
{32, 16, 8}-strict



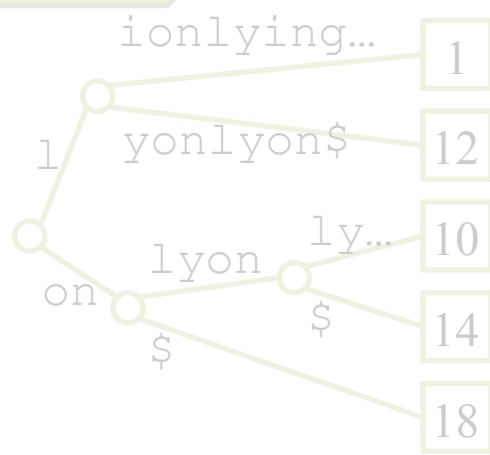
4-strict



2-strict



1-strict

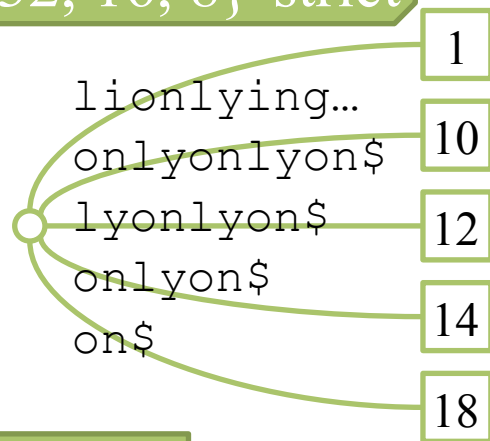


For every edge label of 4-strict tree, we compute FP for the prefix of length 2, and merge edge labels having the same value.

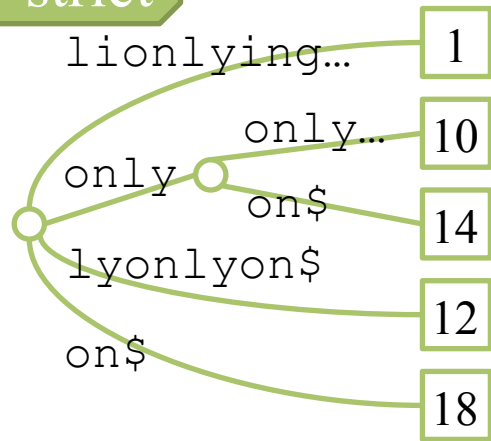
# Gradual refinement

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20  
 $T = \text{lionlyingonlyonllyon\$}$

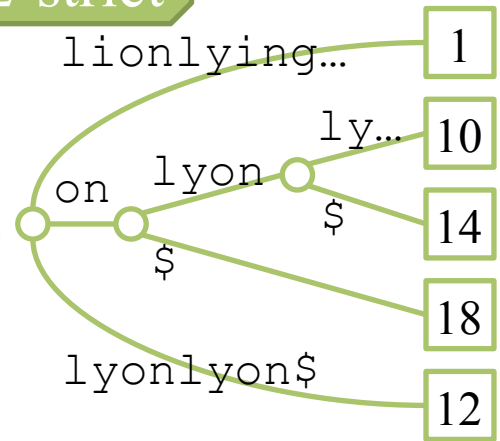
**{32, 16, 8}-strict**



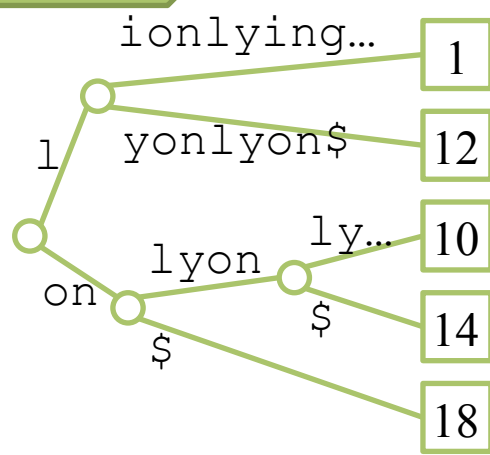
**4-strict**



**2-strict**



**1-strict**



# The total cost for gradual refinement

The total cost for FP computation is  $O((bn/s)\log s)$ .

- Proof. ■ For each refinement step,  
the FPs can be computed in  $O(b \min\{l, n/s\})$  time.
- For the first  $\log s$  steps, the cost is  $O((bn/s)\log s)$ .
  - After  $\log s$  steps, since  $l < n/2^{\log s} = n/s$ , the cost is  $bn/s + bn/2s + bn/4s + \dots + b = O(bn/s)$



# The total cost for gradual refinement

The total cost for merging edge labels is  
 $O(b \log n \log_s n) = O(b \log^2 n / \log s) = O((bn/s) \log s)$ .

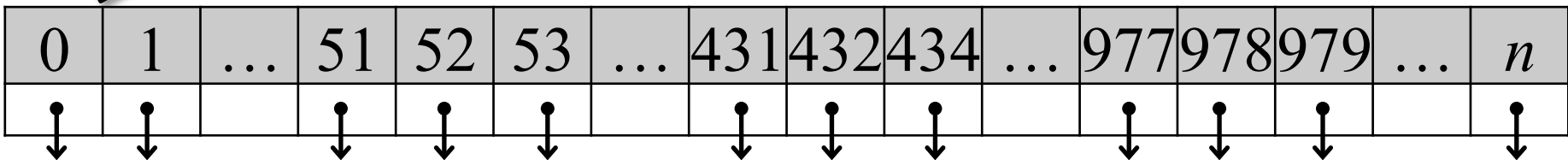
Proof. ■ For each refinement step, we conduct **radix grouping** of  $O(b)$  integer keys with radix  $s$ , which can be done in  $O(b \log_s n)$  time.

# Grouping by distribute-and-collect

$k$  elements associated  
with keys (integers in  $[0, n]$ )

- ★ 978
- ★ 52
- ★ 432
- ★ 52
- ★ 978
- ★ 52

Empty-initialized  
buckets of size  $n$

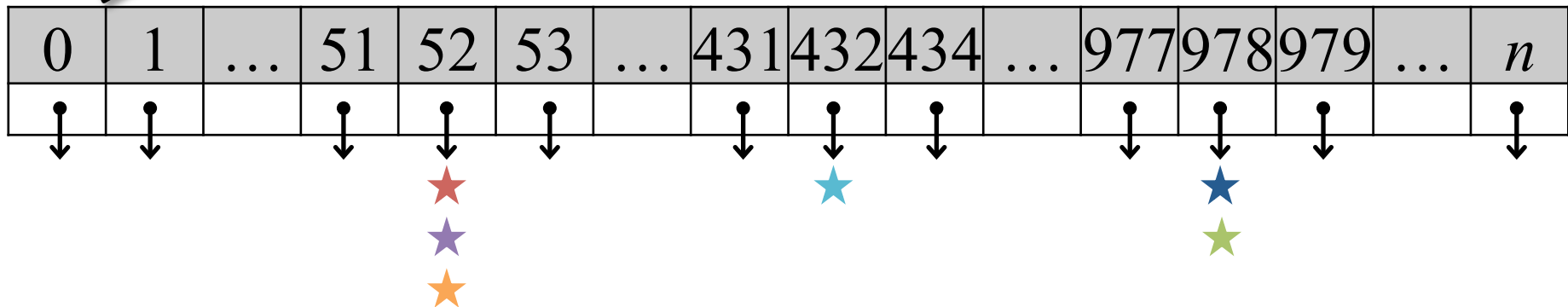


# Grouping by distribute-and-collect

$k$  elements associated  
with keys (integers in  $[0, n]$ )

- ★ 978
- ★ 52
- ★ 432
- ★ 52
- ★ 978
- ★ 52

Empty-initialized  
buckets of size  $n$





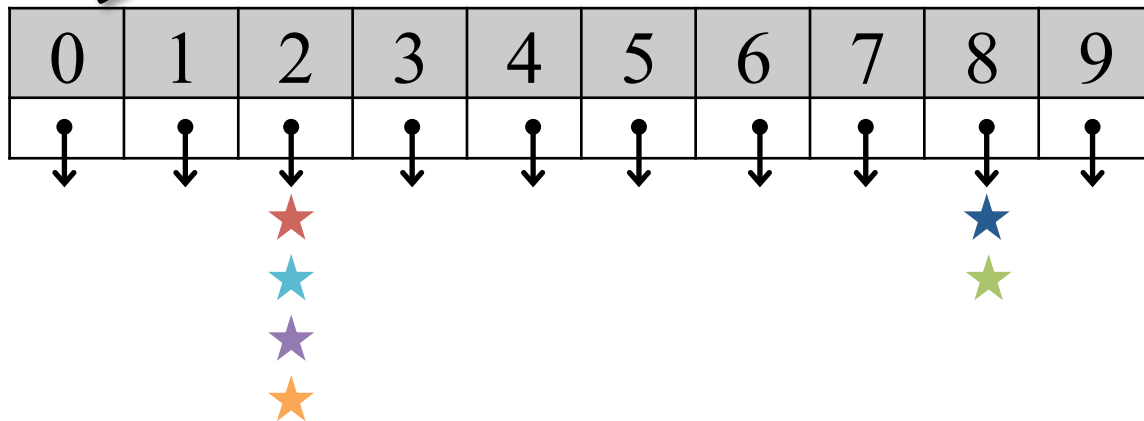


# Grouping by distribute-and-collect

$k$  elements associated  
with keys (integers in  $[0, n]$ )

- ★ 978
- ★ 52
- ★ 432
- ★ 52
- ★ 978
- ★ 52

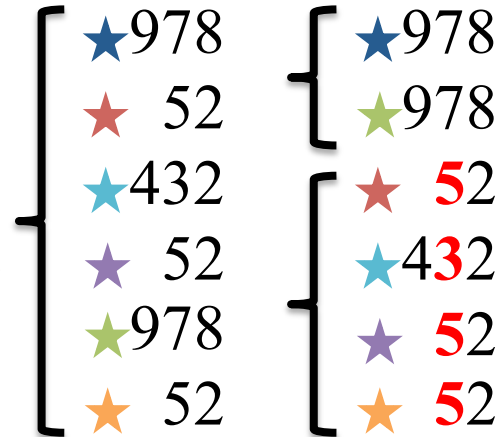
Empty-initialized  
buckets of size 10



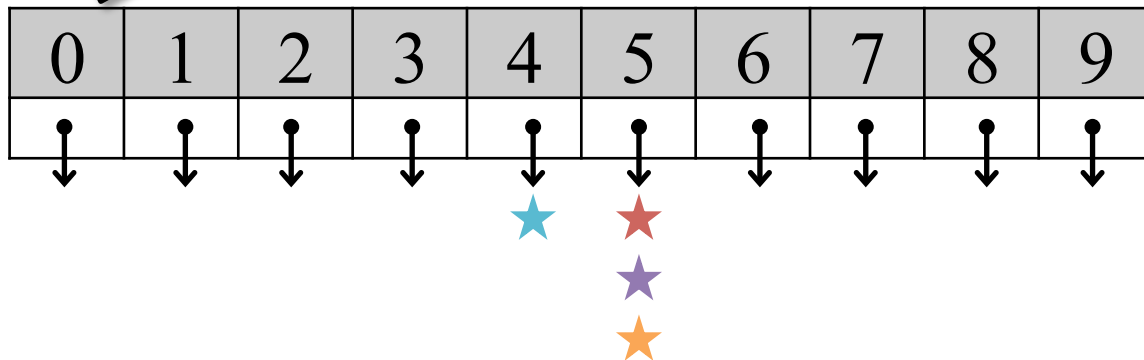


# Grouping by distribute-and-collect

$k$  elements associated  
with keys (integers in  $[0, n]$ )



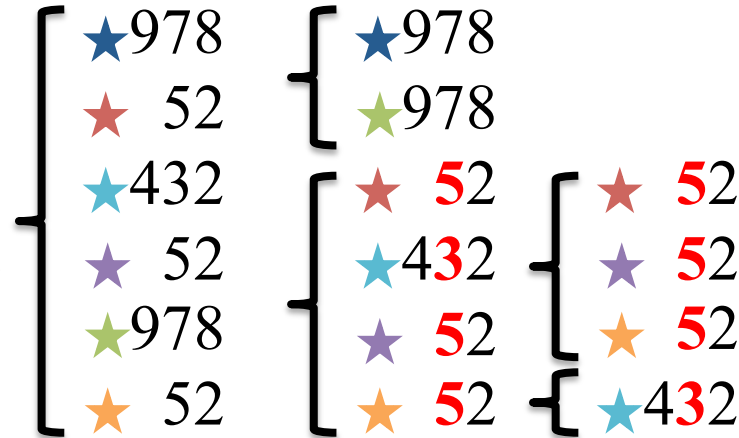
Empty-initialized  
buckets of size 10



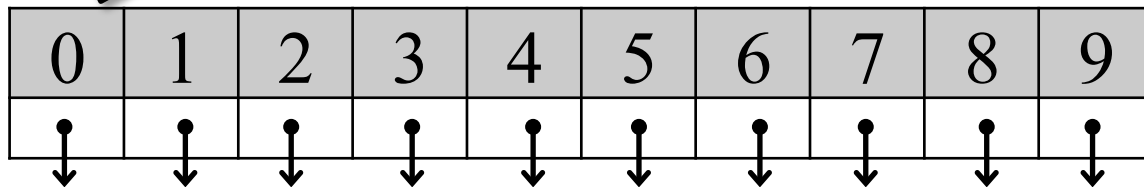


# Grouping by distribute-and-collect

$k$  elements associated  
with keys (integers in  $[0, n]$ )

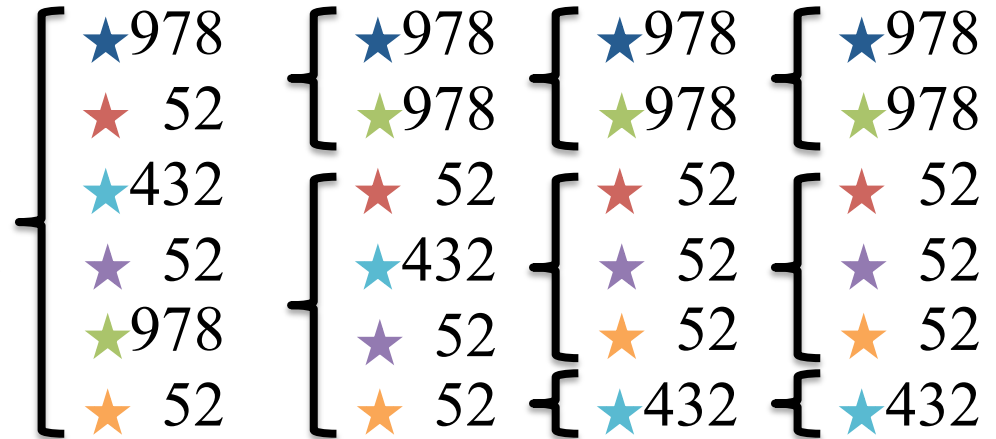


Empty-initialized  
buckets of size 10

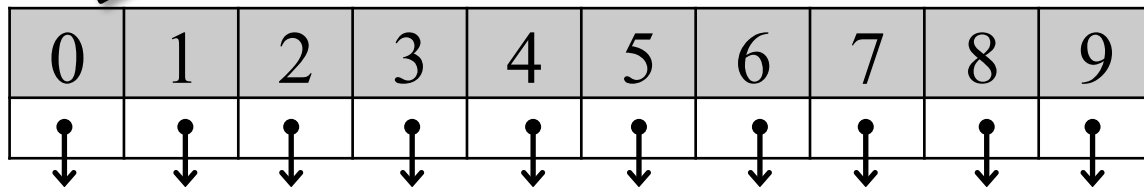


# Grouping by distribute-and-collect

$k$  elements associated  
with keys (integers in  $[0, n]$ )



Empty-initialized  
buckets of size 10

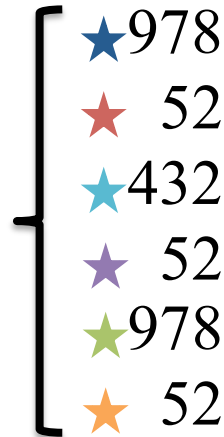


$O(k \log_{10} n)$  time

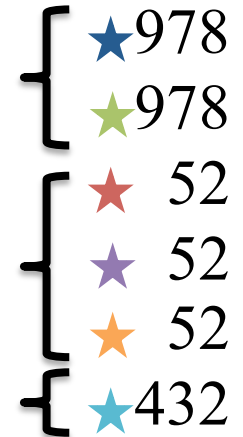
# Grouping by distribute-and-collect

## Radix grouping with radix $s$

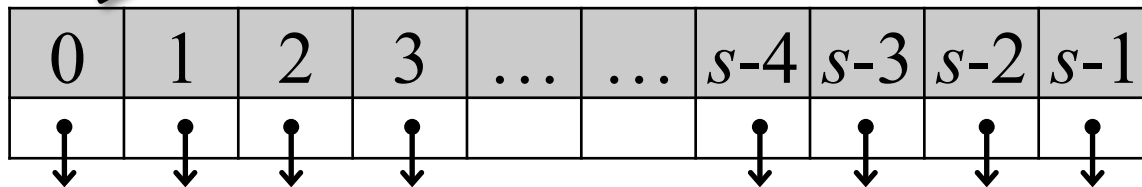
$k$  elements associated  
with keys (integers in  $[0, n]$ )



$\log_s n$  steps



Empty-initialized  
buckets of size  $s$



$O(k \log_s n)$  time

# The total cost for gradual refinement

The total cost for merging edge labels is  
 $O(b \log n \log_s n) = O(b \log^2 n / \log s) = O((bn/s) \log s)$ .

Proof. ■ For each refinement step, we conduct **radix grouping** of  $O(b)$  integer keys with radix  $s$ , which can be done in  $O(b \log_s n)$  time.

# Monte-Carlo algorithm

The total cost for FP computation is  $O((bn/s)\log s)$ .

The total cost for merging edge labels is  
 $O(b \log n \log_s n) = O(b \log^2 n / \log s) = O((bn/s)\log s)$ .

## Theorem

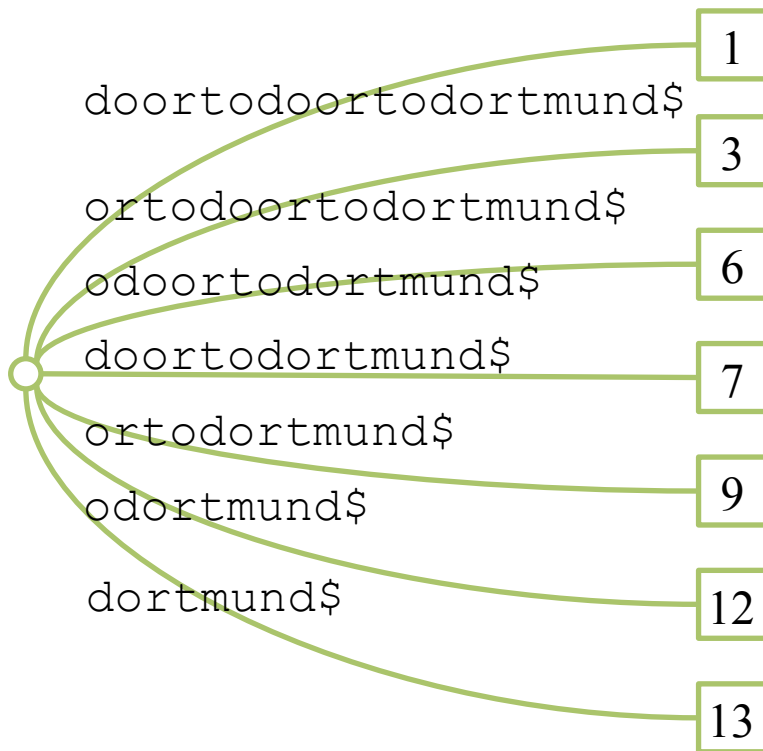
For any  $s \in [b, n]$ , we can construct sparse suffix tree correctly with high probability in  $O(n + (bn/s)\log s)$  time using  $O(s)$  additional space.

- When  $s = b$ , the time complexity is  $O(n \log b)$ .
- When  $s = b \log b$ , the time complexity is  $O(n)$ .

# Exercise

- Follow the gradual refinement steps for this instance.

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21  
 $T = d o o r t o d o o r t o d o r t m u n d \$$



?