Problem: Given an alphabet Σ , a *text* $T \in \Sigma^*$ and a *pattern* $P \in \Sigma^*$, find some occurrence/all occurrences of P in T.



 $P = \operatorname{art}$

$$\Sigma = \{ \texttt{A}, \texttt{B}, \dots, \texttt{Z}, \texttt{a}, \texttt{b}, \dots, \texttt{z}, \lrcorner \}$$

 $T = \texttt{Bart_played_darts_at_the_party}$



Problem: Given an alphabet Σ , a *text* $T \in \Sigma^*$ and a *pattern* $P \in \Sigma^*$, find some occurrence/all occurrences of P in T.



$$\Sigma = \{A, B, \dots, Z, a, b, \dots, z, _\}$$

 $T = Bart_played_darts_at_the_party$
 $P = art$



Problem: Given an alphabet Σ , a *text* $T \in \Sigma^*$ and a *pattern* $P \in \Sigma^*$, find some occurrence/all occurrences of P in T.



$$\Sigma = \{ A, B, \dots, Z, a, b, \dots, z, _ \}$$

 $T = Bart_played_darts_at_the_party$
 $P = art$





- $\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{C}, \mathtt{G}, \mathtt{T}\}$
- T = ACGTGCTTGCAGTGTGCATTACCTGAGTGC...

 $P={\tt GTG}$

Problem: Given an alphabet Σ , a *text* $T \in \Sigma^*$ and a *pattern* $P \in \Sigma^*$, find some occurrence/all occurrences of P in T.



$$\Sigma = \{A, B, \dots, Z, a, b, \dots, z, _\}$$

 $T = Bart_played_darts_at_the_party$
 $P = art$





 $\Sigma = \{A, C, G, T\}$ T = ACGTGCTTGCAGTGTGCATTACCTGAGTGC...P = GTG

One-shot:

- Both the text and the pattern are **part of the input**
- Algorithm design problem

One-shot:

- Both the text and the pattern are **part of the input**
- Algorithm design problem

Repeated:

- The text is **static** and known beforehand (can be preprocessed)
- Patterns are revealed on-demand
- We want to answer each *query* as **quickly** as possible
- Data structure design problem

One-shot:

- Both the text and the pattern are **part of the input**
- Algorithm design problem

Repeated:

- The text is **static** and known beforehand (can be preprocessed)
- Patterns are revealed on-demand
- We want to answer each *query* as **quickly** as possible
- Data structure design problem



Data structure to store a dynamic collection of k strings over an alphabet Σ

 $\Sigma = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$

 $\{ RAD, RADAR, RAG, RAGE, RAGS, RATE \}$

- Insert(T): add T to the collection of strings
- **Delete**(T): remove T from the collection of strings
- **Find**(P): return whether P is in the collection

Data structure to store a dynamic collection of k strings over an alphabet Σ

 $\Sigma = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$

 $\{ RAD, RADAR, RAG, RAGE, RAGS, RATE \}$

- **Insert**(T): add T to the collection of strings
- **Delete**(T): remove T from the collection of strings
- **Find**(P): return whether P is in the collection

Obs: A string comparison requires time O(string length). Binary searching requires time $O(\max \text{ string length} \cdot \log k)$

Data structure to store a dynamic collection of k strings over an alphabet Σ

 $\Sigma = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$

 $\{ RAD, RADAR, RAG, RAGE, RAGS, RATE \}$

- Insert(T): add T to the collection of strings
- **Delete**(T): remove T from the collection of strings
- **Find**(P): return whether P is in the collection
- Count/return the strings in the collection that start with ${\cal P}$

Data structure to store a dynamic collection of k strings over an alphabet Σ

 $\Sigma = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$

 $\{ RAD, RADAR, RAG, RAGE, RAGS, RATE \}$

- Insert(T): add T to the collection of strings
- **Delete**(T): remove T from the collection of strings
- **Find**(P): return whether P is in the collection
- **Count/return** the strings in the collection that start with P
- **Predecessor**(T): return the largest string in the collection that is "not smaller than" T (w.r.t. the lexicopraphic order)

Data structure to store a dynamic collection of k strings over an alphabet Σ

 $\Sigma = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$

 $\{ RAD, RADAR, RAG, RAGE, RAGS, RATE \}$

- **insert**(T). add T to the collection of strings
- **Delete**(T): remove T from the collection of strings
- **Find**(P): return whether P is in the collection
- **Count/return** the strings in the collection that start with P
- **Predecessor**(T): return the largest string in the collection that is "not smaller than" T (w.r.t. the lexicopraphic order)

We will only focus on the static case

Pretend that each string ends with a special "end marker" symbol \$

RAD RADAR RAG RAGE RAGS RATE

Pretend that each string ends with a special "end marker" symbol \$

RAD \$ RADAR \$ RAG \$ RAGE \$ RAGS \$ RATE \$

Pretend that each string ends with a special "end marker" symbol \$

RAD \$ RADAR \$ RAG \$ RAGE \$ RAGS \$ RATE \$

Build a tree in which:

• Edges are labelled with a symbol in $\Sigma \cup \{\$\}$ and are sorted



Pretend that each string ends with a special "end marker" symbol \$

RAD\$ RADAR\$ RAG\$ RAGE\$ RAGS\$ RATE\$

- Edges are labelled with a symbol in $\Sigma \cup \{\$\}$ and are sorted
- Each string T_i corresponds to a root-to-leaf path and vice-versa



Pretend that each string ends with a special "end marker" symbol \$

RAD $\$ RADAR $\$ RAG $\$ RAGE $\$ RAGE $\$ RAGS $\$ RATE $\$

- Edges are labelled with a symbol in $\Sigma \cup \{\$\}$ and are sorted
- Each string T_i corresponds to a root-to-leaf path and vice-versa
- Satellite data is often useful, e.g.:
 Number of leaves in each subtree



Pretend that each string ends with a special "end marker" symbol \$

RAD \$ RADAR \$ RAG \$ RAGE \$ RAGS \$ RATE \$

- Edges are labelled with a symbol in $\Sigma \cup \{\$\}$ and are sorted
- Each string T_i corresponds to a root-to-leaf path and vice-versa
- Satellite data is often useful, e.g.:
- Number of leaves in each subtree
- Pointers to the first/last leaf in the subtree



Pretend that each string ends with a special "end marker" symbol \$

RAD \$ RADAR \$ RAG \$ RAGE \$ RAGS \$ RATE \$

- Edges are labelled with a symbol in $\Sigma \cup \{\$\}$ and are sorted
- Each string T_i corresponds to a root-to-leaf path and vice-versa
- Satellite data is often useful, e.g.:
- Number of leaves in each subtree
- Pointers to the first/last leaf in the subtree
- Pointers from leaves to strings



Pretend that each string ends with a special "end marker" symbol \$

RAD \$ RADAR \$ RAG \$ RAGE \$ RAGS \$ RATE \$

- Edges are labelled with a symbol in $\Sigma \cup \{\$\}$ and are sorted
- Each string T_i corresponds to a root-to-leaf path and vice-versa
- Satellite data is often useful, e.g.:
- Number of leaves in each subtree
- Pointers to the first/last leaf in the subtree'
- Pointers from leaves to strings
- Leaves arranged in a (doubly) linked list



Find(P):

• Walk down the tree matching the characters in P^{\$} with the edge labels





Find(P):

• Walk down the tree matching the characters in P^{\$} with the edge labels



P = RAG

Find(P):

• Walk down the tree matching the characters in P^{\$} with the edge labels

To count the number of strings that start with P:

 $\bullet\,$ Find the node corresponding to P



P = RAG

Find(P):

• Walk down the tree matching the characters in P^{\$} with the edge labels

To count the number of strings that start with P:

- $\bullet\,$ Find the node corresponding to P
- Return the number of leaves in the subtree (stored in the node)



P = RAG

Find(P):

• Walk down the tree matching the characters in P^{\$} with the edge labels

To count the number of strings that start with P:

- $\bullet\,$ Find the node corresponding to P
- Return the number of leaves in the subtree (stored in the node)
- The actual matches can be listed in O(1) additional time per match by following pointers



 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

• Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T^{\$} with the edge labels



 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

• Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T^{\$} with the edge labels



 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

- Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T^{\$} with the edge labels
- If T\$ is found we are done



 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

- Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T\$ with the edge labels
- If T^{\$} is found we are done
- Otherwise, stop at the node v_i matching the longest prefix $T_1T_2...T_i$



 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

- Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T^{\$} with the edge labels
- If T^{\$} is found we are done
- Otherwise, stop at the node v_i matching the longest prefix $T_1T_2...T_i$
- Find the deepest ancestor of v_j of v_i (possibly v_i itself) such that T_j has a strict predecessor u w.r.t. v_j .



The strict predecessor of $\sigma \in \Sigma$ w.r.t. a node v, if it exists, is the child u of v such that (v, u) has the largest label that is smaller than σ

 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

- Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T\$ with the edge labels
- If T\$ is found we are done
- Otherwise, stop at the node v_i matching the longest prefix $T_1T_2...T_i$
- Find the deepest ancestor of v_j of v_i (possibly v_i itself) such that T_j has a strict predecessor u w.r.t. v_j .
- Follow the pointers from u to the maximum string in its subtree



 $T\$ = T_1 T_2 T_3 \dots$

Predecessor(T):

Time?

- Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T^{\$} with the edge labels
- If T^{\$} is found we are done
- Otherwise, stop at the node v_i matching the longest prefix $T_1T_2...T_i$
- Find the deepest ancestor of v_j of v_i (possibly v_i itself) such that T_j has a strict predecessor u w.r.t. v_j .
- Follow the pointers from u to the maximum string in its subtree

T = RATA R A C D G

Α

R

\$

\$

RAGS

F

\$

\$

E

\$

T^{\$} = $T_1 T_2 T_3 \dots$

Predecessor(T):

Time?

- Walk down a path $\langle v_0, v_1, v_2 \dots \rangle$ of the tree matching the characters in T^{\$} with the edge labels
- If T is found we are done
- Otherwise, stop at the node v_i matching the longest prefix $T_1T_2\ldots T_i$
- Find the deepest ancestor of v_i of v_i (possibly v_i itself) such that T_i has a strict predecessor u w.r.t. v_i .
- Follow the pointers from u to the maximum string in its subtree

T = RATA



Depends on how the tree is stored

Representing Tries



- $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$
- $\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$






 $\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$



Space: $O(|\Sigma|)$

Time to find a symbol's edge: O(1)



 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$



Space: $O(|\Sigma|)$

Time to find a symbol's edge: O(1)

Time to find predecessor: $O(|\Sigma|)$



 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$



Space: $O(|\Sigma|)$

Time to find a symbol's edge: O(1)

Time to find predecessor: $O(|\Sigma|)$



 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$



Space: $O(|\Sigma|)$

Time to find a symbol's edge: O(1)

Time to find predecessor: $O(|\Sigma|) O(1)$



 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$



Space: $O(|\Sigma|)$

Time to find a symbol's edge: O(1)

Time to find predecessor: $O(|\Sigma|) O(1)$

Overall space: $O(|\Sigma| \cdot n)$ **Overall time:** O(|P|)



 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$



 $\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$





Balanced Binary Search Tree

$$n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$







 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

 $\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$



Balanced Binary Search Tree





Space: O(#children)



Array (sparse)

 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$







Space: O(#children)

Time to find a symbol's edge/predecessor: $O(\log \# \text{children}) = O(\log |\Sigma|)$



Array (sparse)

Balanced Binary Search Tree



Overall space: O(n)**Overall time:** $O(|P| \log |\Sigma|)$

Space: O(#children)

Time to find a symbol's edge/predecessor: $O(\log \# \text{children}) = O(\log |\Sigma|)$

$$\begin{array}{c|c} D & G & T \\ \hline D & C & d \end{array}$$

$$\boldsymbol{\Sigma} = \{\mathtt{A}, \mathtt{D}, \mathtt{E}, \mathtt{G}, \mathtt{R}, \mathtt{S}, \mathtt{T}\}$$

 $n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$

Weight-Balanced BSTs

(a)

b

9

 $(\tilde{\underline{C}})$

$$n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$$

Each vertex of the trie has a weight equal to the number of leaves in its subtree

Recursively construct a binary search tree by splitting the children in the trie so that the sum of their weights is as balanced as possible



Weight-Balanced BSTs

$$n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$$

Each vertex of the trie has a weight equal to the number of leaves in its subtree

Recursively construct a binary search tree by splitting the children in the trie so that the sum of their weights is as balanced as possible





Weight-Balanced BSTs

$$n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$$

Each vertex of the trie has a weight equal to the number of leaves in its subtree

Recursively construct a binary search tree by splitting the children in the trie so that the sum of their weights is as balanced as possible





Weight-Balanced BSTs

$$n = \# \mathsf{nodes} = O\left(\sum_i |T_i|\right)$$

Each vertex of the trie has a weight equal to the number of leaves in its subtree

Recursively construct a binary search tree by splitting the children in the trie so that the sum of their weights is as balanced as possible





Space: O(#children) Overall space: O(n)

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves.

Imagine the leaves in the subtree of v as consecutive segments with lengths equal to their weights



Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves.

Imagine the leaves in the subtree of v as consecutive segments with lengths equal to their weights

If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains more than one segment:



Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves.

Imagine the leaves in the subtree of v as consecutive segments with lengths equal to their weights

If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains more than one segment:



Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves.

Imagine the leaves in the subtree of v as consecutive segments with lengths equal to their weights

If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains more than one segment:



• The weight of each children of v is at most $\frac{2}{3}w(v)$

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves. If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains a single segment, let x be the corresponding leaf



Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves. If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains a single segment, let x be the corresponding leaf



 $\frac{1}{3}w(v) \qquad \frac{2}{3}w(v)$



Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves. If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains a single segment, let x be the corresponding leaf



Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves. If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains a single segment, let x be the corresponding leaf



•
$$w(v'') \leq \frac{1}{3}w(v)$$
.

Weight-Balanced BSTs

 $\frac{1}{2}w(v')$

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves. If the interval $[\frac{1}{3}w(v), \frac{2}{3}w(v)]$ contains a single segment, let x be the corresponding leaf

• v splits the segments immediately before/after x.



- Let v' be the child of v that contains x and let v'' be the other child
- $w(v'') \leq \frac{1}{3}w(v)$.
- x is the first or last leaf in the subtree of v' and $w(x) \ge \frac{1}{2}w(v')$

Weight-Balanced BSTs

 $\frac{1}{2}w(v')$

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}w(v)$ or are leaves. If the interval $\left[\frac{1}{3}w(v), \frac{2}{3}w(v)\right]$ contains a single segment, let x be the corresponding leaf

• v splits the segments immediately before/after x.



- $\bullet\,$ Let v' be the child of v that contains x and let v'' be the other child
- $w(v'') \leq \frac{1}{3}w(v)$.
- x is the first or last leaf in the subtree of v' and $w(x) \ge \frac{1}{2}w(v')$
- One child of v' is x and the other child weighs $\leq \frac{1}{2}w(v') \leq \frac{1}{2}w(v)$

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}(v)$ or are leaves.

Traversing two edges of a weight-balanced BST either:

• Brings us to the next node in the trie, i.e., we advance one character into P; or

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}(v)$ or are leaves.

Traversing two edges of a weight-balanced BST either:

• Brings us to the next node in the trie, i.e., we advance one character into P; or

• Reduces the weight (i.e, the number of leaves in the trie reachable from the current node) by at least a 2/3 factor

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}(v)$ or are leaves.

Traversing two edges of a weight-balanced BST either:

- Brings us to the next node in the trie, i.e., we advance one character into P; or
 Can only happen O(|P|) times
- Reduces the weight (i.e, the number of leaves in the trie reachable from the current node) by at least a 2/3 factor

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}(v)$ or are leaves.

Traversing two edges of a weight-balanced BST either:

Brings us to the next node in the trie, i.e., we advance one character into P; or
 Can only happen O(|P|) times

Can only happen O(|P|) times

Reduces the weight (i.e, the number of leaves in the trie reachable from the current node) by at least a 2/3 factor
 Can only happen O(log_{3/2} #leaves) = O(log k) times

Weight-Balanced BSTs

Claim: All the grand-children u of v satisfy $w(u) \leq \frac{2}{3}(v)$ or are leaves.

Traversing two edges of a weight-balanced BST either:

- Brings us to the next node in the trie, i.e., we advance one character into P; or
 Can only happen O(|P|) times
- Reduces the weight (i.e, the number of leaves in the trie reachable from the current node) by at least a 2/3 factor
 Can only happen O(log_{3/2} #leaves) = O(log k) times

Overall space: O(n) **Overall time:** $O(|P| + \log k)$

Representing Tries: Recap				
	Space	Query Time		
Array (dense)	$O(\Sigma \cdot n)$	O(P)		
Array (sparse) / BST	O(n)	$O(P \log \Sigma)$		
Weight-balanced BST	O(n)	$O(P + \log k)$		

Representing Tries: Recap			
	Space	Query Time	
Array (dense)	$O(\Sigma \cdot n)$	O(P)	
Array (sparse) / BST	O(n)	$O(P \log \Sigma)$	
Weight-balanced BST	O(n)	$O(P + \log k)$	
Optimal			

Representing Tries: Recap				
	Space	Query Time		
Array (dense)	$O(\Sigma \cdot n)$	O(P)		
Array (sparse) / BST	O(n)	$O(P \log \Sigma)$		
Weight-balanced BST	O(n)	$O(P + \log k)$		
		Can we get rid of this term?		
Optimal				

Representing Tries: Recap				
	Space	Query Time		
Array (dense)	$O(\Sigma \cdot n)$	O(P)		
Array (sparse) / BST	O(n)	$O(P \log \Sigma)$		
Weight-balanced BST	O(n)	$O(P + \log k)$		
		Can we get rid of this term?		
Optimal		Almost		

Indirection

We can use a similar technique to the one we encountered while designing level ancestor oracles



Indirection

We can use a similar technique to the one we encountered while designing level ancestor oracles



Find the set M of all maximally deep vertices with at least $|\Sigma|$ descendants

Indirection

We can use a similar technique to the one we encountered while designing level ancestor oracles



Find the set M of all maximally deep vertices with at least $|\Sigma|$ descendants
We can use a similar technique to the one we encountered while designing level ancestor oracles



Find the set M of all maximally deep vertices with at least $|\Sigma|$ descendants

Split the trie into a top-tree T' containing all the ancestors of the vertices in M and several bottom-trees in $T \setminus T'$.

We can use a similar technique to the one we encountered while designing level ancestor oracles



Find the set M of all maximally deep vertices with at least $|\Sigma|$ descendants

Split the trie into a top-tree T' containing all the ancestors of the vertices in M and several bottom-trees in $T \setminus T'$.

We can use a similar technique to the one we encountered while designing level ancestor oracles



Find the set M of all maximally deep vertices with at least $|\Sigma|$ descendants

Split the trie into a top-tree T' containing all the ancestors of the vertices in M and several bottom-trees in $T \setminus T'$.

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children)

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children)

• Store leaves using dense arrays

Space $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children)

- Store leaves using dense arrays
- Store branching nodes using dense arrays

Space $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$ $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children) **Space**

- Store leaves using dense arrays
- Store branching nodes using dense arrays
- Store the unique child of each non-branching node explicitly O(n)

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children) **Space**

- Store leaves using dense arrays
- Store branching nodes using dense arrays
- Store the unique child of each non-branching node explicitly O(n)Time to find the next node O(1)

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children) **Space**

- Store leaves using dense arrays
- Store branching nodes using dense arrays
- Store the unique child of each non-branching node explicitly O(n)Time to find the next node O(1)

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

Storing the bottom trees:

• Store each bottom tree using a weight-balanced BST Total space of all bottom trees: O(n)

Storing the top tree:

The number of leaves of T' is at most $\frac{n}{|\Sigma|}$

Fact: A tree with ℓ leaves has at most $\ell - 1$ branching nodes (i.e., nodes with at least 2 children) **Space**

- Store leaves using dense arrays
- Store branching nodes using dense arrays
- Store the unique child of each non-branching node explicitly O(n)Time to find the next node O(1)

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

 $O(|\Sigma| \cdot \frac{n}{|\Sigma|}) = O(n)$

Storing the bottom trees:

- Store each bottom tree using a weight-balanced BST Total space of all bottom trees: O(n)
- Each bottom tree has at most $|\Sigma|$ leaves Time to navigate a bottom tree: $O(|P| + \log |\Sigma|)$

Representing Tries: Recap		
	Space	Query Time
Array (dense)	$O(\Sigma \cdot n)$	O(P)
Array (sparse) / BST	O(n)	$O(P \log \Sigma)$
Weight-balanced BST	O(n)	$O(P + \log k)$

Representing Tries: Recap		
Space	Query Time	
$O(\Sigma \cdot n)$	O(P)	
O(n)	$O(P \log \Sigma)$	
O(n)	$O(P + \log k)$	
O(n)	$O(P + \log \Sigma)$	
	ing Tries: Space $O(\Sigma \cdot n)$ O(n) O(n)	

Representing Tries: Recap		
	Space	Query Time
Array (dense)	$O(\Sigma \cdot n)$	O(P)
Array (sparse) / BST	O(n)	$O(P \log \Sigma)$
Weight-balanced BST	O(n)	$O(P + \log k)$
Indirection	O(n)	$O(P + \log \Sigma)$

Can be made dynamic with a time complexity of $O(|T| + \log |\Sigma|)$ per insertion/deletion of T

Sort a collection of k strings T_1, T_2, \ldots, T_k over Σ

 $L = \max_{i=1,\dots,k} |T_i|$

Obs: A string comparison requires time O(L). Naive sorting algorithms take time $O(Lk \log k)$ or $O(L(k + |\Sigma|))$

Sort a collection of k strings T_1, T_2, \ldots, T_k over Σ

 $L = \max_{i=1,\dots,k} |T_i|$

Obs: A string comparison requires time O(L). Naive sorting algorithms take time $O(Lk \log k)$ or $O(L(k + |\Sigma|))$

- Create an empty trie
- For i = 1, ..., k:
 - Insert T_i into the trie
- An in-order visit of the trie returns the strings in lexicographic order

Sort a collection of k strings T_1, T_2, \ldots, T_k over Σ $L = \max_{i=1,\ldots,k} |T_i|$

Obs: A string comparison requires time O(L). Naive sorting algorithms take time $O(Lk \log k)$ or $O(L(k + |\Sigma|))$

Time

 $\left\{ O\left(\sum_{i=1}^{k} (|T_i| + \log |\Sigma|)\right) \right.$

- Create an empty trie
- For i = 1, ..., k:
 - Insert T_i into the trie
- An in-order visit of the trie returns the strings in lexicographic order

Sort a collection of k strings T_1, T_2, \ldots, T_k over Σ $L = \max_{i=1,\ldots,k} |T_i|$

Obs: A string comparison requires time O(L). Naive sorting algorithms take time $O(Lk \log k)$ or $O(L(k + |\Sigma|))$

Time

- Create an empty trie
- For i = 1, ..., k:

• Insert T_i into the trie

- $O\left(n+k\log|\Sigma|\right)$
- An in-order visit of the trie returns the strings in lexicographic order

Sort a collection of k strings T_1, T_2, \ldots, T_k over Σ $L = \max_{i=1,\dots,k} |T_i|$

Obs: A string comparison requires time O(L). Naive sorting algorithms take time $O(Lk \log k)$ or $O(L(k + |\Sigma|))$

- Create an empty trie
- $O\left(n+k\log|\Sigma|\right)$ • For i = 1, ..., k:
 - Insert T_i into the trie
- An in-order visit of the trie returns the strings in O(n)lexicographic order

Time

Sort a collection of k strings T_1, T_2, \ldots, T_k over Σ $L = \max_{i=1,\ldots,k} |T_i|$

Obs: A string comparison requires time O(L). Naive sorting algorithms take time $O(Lk \log k)$ or $O(L(k + |\Sigma|))$

- Create an empty trie
- For i = 1, ..., k:
 - Insert T_i into the trie
- An in-order visit of the trie returns the strings in lexicographic order

Overall time: $O(n + k \log |\Sigma|)$

Time

O(n)

 $O\left(n+k\log|\Sigma|\right)$

Among all the destinations that match, a packet gets routed to the one with the most specific rule

Packet	Routing Table	
Src: 192.168.42.10	Destination	Interface
Dst: 101.167.200.15	169.0.0.0/11	eth1
	169.48.0.0/12	ppp0
'	169.128.0.0/10	eth1
	169.160.0.0/11	eth0
	96.0.0.0/3	tun1
	96.0.0.0/6	tun0
	100.0.0/8	eth0
	127.0.0.0/8	lo
	default	wlan0

Among all the destinations that match, a packet gets routed to the one with the most specific rule

Packet	Routing Table	
Src: 192.168.42.10	Destination	Interface
Dst: 0110010110100111	10101001000 <mark>\$</mark>	eth1
	101010010011 <mark>\$</mark>	ppp0
'	1010100110 <mark>\$</mark>	eth1
	10101001101 <mark>\$</mark>	eth0
	011 <mark>\$</mark>	tun1
	011000 <mark>\$</mark>	tun0
	01100100 <mark>\$</mark>	eth0
	01111111 <mark>\$</mark>	lo
	\$	wlan0

Among all the destinations that match, a packet gets routed to the one with the most specific rule

Packet	Routing Table	
Src: 192.168.42.10	Destination	Interface
Dst: 0110010110100111	10101001000 <mark>\$</mark>	eth1
	101010010011 <mark>\$</mark>	ppp0
'	1010100110 <mark>\$</mark>	eth1
	10101001101 <mark>\$</mark>	eth0
	011	tun1
	011000 <mark>\$</mark>	tun0
	01100100 <mark>\$</mark>	eth0
	0111111 <mark>\$</mark>	lo
	\$	wlan0

Among all the destinations that match, a packet gets routed to the one with the most specific rule

Packet	Routing Table	
Src: 192.168.42.10	Destination	Interface
Dst: 0110010110100111	10101001000\$	eth1
\tilde{P}	101010010011 <mark>\$</mark>	ppp0
·'	1010100110 <mark>\$</mark>	eth1
	10101001101 <mark>\$</mark>	eth0
	011\$	tun1
	011000\$	tun0
	01100100\$	eth0
	01111111 <mark>\$</mark>	lo
	\$	wlan0

Given a pattern ${\cal P}$ we want the longest string in our collection that appears as a prefix of ${\cal P}$

Build a trie T with all the addresses in the routing table.



 $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P



 $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ
- Route the packet towards the interface stored in ℓ



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ
- Route the packet towards the interface stored in ℓ



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ
- Route the packet towards the interface stored in ℓ



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ
- Route the packet towards the interface stored in ℓ



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ
- Route the packet towards the interface stored in ℓ

Build a trie T with all the addresses in the routing table.



- $\bullet\,$ Find the node v corresponding to the maximal prefix that matches P
- Walk up the tree searching for the deepest ancestor u of v incident to a "§" edge towards a leaf ℓ
- Route the packet towards the interface stored in ℓ

Time: O(address length)








Back to String Matching

Problem: Given an alphabet Σ , a *text* $T \in \Sigma^*$ and a *pattern* $P \in \Sigma^*$, find some occurrence/all occurrences of P in T.



$$\Sigma = \{A, B, \dots, Z, a, b, \dots, z, \bot\}$$
$$T = Bart_played_darts_at_the_party$$
$$P = art$$



Want: A data structure that can preprocesses T and answer string matching queries

The suffix tree of T is the compressed trie of all the suffixes of T

 $\Sigma = \{ \mathtt{A}, \mathtt{B}, \mathtt{N}, \mathtt{S} \}$ $T = \mathtt{BANANAS}$

The suffix tree of T is the compressed trie of all the suffixes of T^{\$}

 $\Sigma = \{A, B, N, S\}$ T = BANANAS

- 7 \$
- 6 S\$
- 5 AS\$
- 4 NAS\$
- 3 ANAS\$
- 2 NANAS\$
- 1 ANANAS\$
- 0 BANANAS\$

The suffix tree of T is the compressed trie of all the suffixes of T

01234567 $\Sigma = \{A, B, N, S\}$ T = BANANAS7 \$ BANANAS\$ 6 S\$ S\$ \$ NA Α 5 AS\$ 4 NAS\$ S\$ NA NAS<mark>\$</mark> **S**\$ 3 ANAS\$ 2 NANAS^{\$} S\$ NAS\$ 1 ANANAS\$

0 BANANAS\$

The suffix tree of T is the compressed trie of all the suffixes of T



Label edges with indices into \boldsymbol{T}

Label leaves with the index of the start of the corresponding suffix

The suffix tree of T is the compressed trie of all the suffixes of T

01234567 $\Sigma = \{A, B, N, S\}$ T = BANANAS7 \$ BANANAS\$ 6 S\$ S\$ \$ NA Α 5 AS\$ 6 4 NASS\$ NA **S**\$ NAS\$ 3 ANAS\$ 5 2 NANAS^{\$} S**\$** NAS<mark>\$</mark> 1 ANANAS^{\$} BANANAS\$ 0

Label edges with indices into ${\boldsymbol{T}}$

Label leaves with the index of the start of the corresponding suffix Space: O(# nodes) = O(# leaves) = O(|T|)



- $\bullet\,$ Find the node v corresponding to P
- The occurrences of ${\cal P}$ are all and only the leaves in the subtree of v



- $\bullet\,$ Find the node v corresponding to P
- The occurrences of ${\cal P}$ are all and only the leaves in the subtree of v



- $\bullet\,$ Find the node v corresponding to P
- $\bullet\,$ The occurrences of P are all and only the leaves in the subtree of v
- Arrange leaves in a linked list to find the next match in O(1) time



- $\bullet\,$ Find the node v corresponding to P
- $\bullet\,$ The occurrences of P are all and only the leaves in the subtree of v
- Arrange leaves in a linked list to find the next match in ${\cal O}(1)$ time

Time: $O(|P| + \log |\Sigma| + \# \text{desired matches})$



- $\bullet\,$ Find the node v corresponding to P
- $\bullet\,$ The occurrences of P are all and only the leaves in the subtree of v
- Arrange leaves in a linked list to find the next match in ${\cal O}(1)$ time

Time: $O(|P| + \log |\Sigma| + \# \text{desired matches})$ Number of matches in time $O(|P| + \log |\Sigma|)$

Each vertex stores the # of leaves in its subtrees



• Look at the leaves u_i , u_j corresponding to T[i:] and T[j:]



- Look at the leaves u_i , u_j corresponding to T[i:] and T[j:]
- Find the common prefix of the paths from the root to u_i and u_j



- Look at the leaves u_i , u_j corresponding to T[i:] and T[j:]
- Find the common prefix of the paths from the root to u_i and u_j
- This is the path from the root to the lowest common ancestor of \boldsymbol{u}_i and \boldsymbol{u}_j



- Look at the leaves u_i , u_j corresponding to T[i:] and T[j:]
- Find the common prefix of the paths from the root to u_i and u_j
- This is the path from the root to the lowest common ancestor of \boldsymbol{u}_i and \boldsymbol{u}_j

We already know how to answer LCA queries in constant time!





Find the longest string that appears at least twice in T as a substring:

Applications: Longest Repeated Substring



Find the longest string that appears at least twice in T as a substring:

• Assign a length to each edge equal to the number of symbols in its label

Applications: Longest Repeated Substring



Find the longest string that appears at least twice in T as a substring:

- Assign a length to each edge equal to the number of symbols in its label
- Find the deepest (w.r.t. edge lengths) node with at least two descendants

Applications: Longest Repeated Substring



Find the longest string that appears at least twice in T as a substring:

- Assign a length to each edge equal to the number of symbols in its label
- Find the deepest (w.r.t. edge lengths) node with at least two descendants

Time: O(|T|)



Given an occurrence T[i:j] of P in T, find all other occurrences of P:

 $\bullet\,$ We want to quickly find the node that corresponds to P



Given an occurrence T[i:j] of P in T, find all other occurrences of P:

- $\bullet\,$ We want to quickly find the node that corresponds to P
- Start from the leaf corresponding to T[i:]



Given an occurrence T[i:j] of P in T, find all other occurrences of P:

- $\bullet\,$ We want to quickly find the node that corresponds to P
- Start from the leaf corresponding to T[i:]
- Walk up the tree for "|T| j" characters



Given an occurrence T[i:j] of P in T, find all other occurrences of P:

- $\bullet\,$ We want to quickly find the node that corresponds to P
- Start from the leaf corresponding to T[i:]
- Walk up the tree for "|T| j" characters
- This is a **weighted** level ancestor query!



Given an occurrence T[i:j] of P in T, find all other occurrences of P:

- $\bullet\,$ We want to quickly find the node that corresponds to P
- Start from the leaf corresponding to T[i:]
- Walk up the tree for "|T| j" characters
- This is a **weighted** level ancestor query!

We can answer weighted LA queries in $O(\log \log |T|)$ time!



01234567T = BANANAS

Given an occurrence T[i:j] of P in T, find all other occurrences of P:

- $\bullet\,$ We want to quickly find the node that corresponds to P
- Start from the leaf corresponding to T[i:]
- Walk up the tree for "|T| j" characters
- This is a **weighted** level ancestor query!
- Link leaves to find the other occurrences in O(1) additional time each

We can answer weighted LA queries in $O(\log \log |T|)$ time!

Preprocess collection of documents T_1, T_2, \ldots, T_k to quickly find all documents that contain a pattern P

Preprocess collection of documents T_1, T_2, \ldots, T_k to quickly find all documents that contain a pattern P

Use the end symbol \mathbf{s}_i for document T_i and build a suffix-tree with the suffixes of all the strings $T_i\mathbf{s}_i$

Preprocess collection of documents T_1, T_2, \ldots, T_k to quickly find all documents that contain a pattern P

Use the end symbol \mathbf{s}_i for document T_i and build a suffix-tree with the suffixes of all the strings $T_i\mathbf{s}_i$



Applications: Document Retrieval HOSASS 1525 05833 es/ g $\$_1$ $\$_2$ \$₃ $\$_4$. <mark>૬</mark>૪ ิล S Sage<mark>%</mark> ugar St gar\$2 nosas \$3 \$4 \$ \$ \$ ˈsa S а \$₃ Index the leaves from left to right

Applications: Document Retrieval HOSASS Isas 05283 es/ g $\$_1$ $\$_2$ **\$**3 $\$_4$ ં ૪ઙુ а S Sage 👌 ugars losa -sa (ଦୁ ତ ତ ତ a 3_3 \$3 Index the leaves from left to right







Find all distinct documents (colors) in A[i:j]



Constructing Suffix Trees & & Suffix Arrays
01234567

 $T = {\tt BANANAS}$

Sort all suffixes along with their start index

- 0 BANANAS\$
- 1 ANANAS\$
- 2 NANAS\$
- 3 ANAS\$
- 4 NAS\$
- 5 AS\$
- 6 S\$
- 7 \$

01234567

 $T = {\tt BANANAS}$

Sort all suffixes along with their start index

7 \$

- 1 ANANAS\$
- 3 ANAS\$
- 5 AS\$
- 0 BANANAS\$
- 2 NANAS\$
- 4 NAS\$
- 6 S\$

01234567

 $T = {\tt BANANAS}$

Sort all suffixes along with their start index



_ Suffix array

01234567

 $T = {\tt BANANAS}$

Length of the longest common prefix between adjacent suffixes in the sorted order \sim

7	\$	
1	ANANAS <mark>\$</mark>	0
3	ANAS <mark>\$</mark>	う 1
5	AS\$	
0	BANANAS <mark>\$</mark>	0
2	NANAS <mark>\$</mark>	0
4	NAS <mark>\$</mark>	
6	S\$	U



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order \sim



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays

0 3 1	0	$0 \mid 2 \mid 0$	
-------	---	-------------------	--

01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays



01234567

T = BANANAS

Length of the longest common prefix between adjacent suffixes in the sorted order



We can construct a suffix tree from the suffix and LCP arrays



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit



Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit

01234567 T = BANANAS()()U BANANAS\$ \$ NA **S**\$ А 6 \$ BANANAS^{\$} **S**\$ S\$ NA **S**\$ NAS^{\$} 3 AS\$ NANAS^{\$} NAS^{\$} S\$ NAS<mark>\$</mark> ANANAS^{\$} ANAS<mark>\$</mark>

Construction time (given Suffix + LCP Arrays): O(|T|)

Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit

 $\begin{array}{c} 01234567\\ T = \texttt{BANANAS} \end{array}$



Construction time (given Suffix + LCP Arrays): O(|T|)

Suffix + LCP Arrays can be built in O(|T|) time [J. Kärkkäinen, P. Sanders, ICALP'03]

Branching vertices are labelled with their *letter depth*, i.e., the number of letters of the prefix encoded in the path from the root to the vertex

Edge labels are easy to reconstruct with a post-order visit

01234567 Construction time T = BANANAS(given Suffix + LCP Arrays): O(|T|)()()U BANANAS\$ \$ NA S\$ Α Suffix + LCP Arrays can be 6 \$ built in O(|T|) time BANANAS^{\$} **S**\$ S\$ NA [J. Kärkkäinen, P. Sanders, ICALP'03] S\$ NAS^{\$} 3 AS\$ NANAS^{\$} NAS^{\$} S\$ NAS<mark>\$</mark> Suffix trees can be built in O) time! ANANAS^{\$} ANAS^{\$}