# Indexing for complex queries

Indexing: preprocess a text or a collection of texts into a data structure that allows locating occurrences of a query pattern in the texts.

- If the query is a <u>string</u>, multiple space- and time-efficient solutions exist

However, *it is desirable to allow for **more general queries!***

- For <u>regular expression</u> patterns, there cannot be a  data structure with polynomial-time preprocessing and sublinear query time, conditioned on the online matrix-vector multiplication conjecture [Thankanchan and Gibney, 2021]

- What about simpler models with just two patterns?

# Indexing for complex queries

**Indexing:** preprocess a text or a collection of texts into a data structure that allows locating occurrences of a query pattern in the texts.

- If we are looking for all texts in a collection that contains two string patterns $P_1, P_2$, or all texts containing $P_1$ but not $P_2$, the asymptotically fastest linear-space solutions use $O(\sqrt{N})$ query time, where $N$ is the total length of the texts [Hon et al. 2010, Hon et al. 2012]

- This was shown to be optimal conditioned on Boolean Matrix Multiplication [Larsen et al. 2014] and on the 3SUM conjecture [Kopelowitz et al. 2016]

# Indexing for complex queries

**Indexing:** preprocess a text or a collection of texts into a data structure that allows locating occurrences of a query pattern in the texts.

- [Kopelowitz and Krauthgamer 2016] considered the problem of retrieving the pair of closest occurrences of two patterns $P_1, P_2$ in a text $T$.

- For a text of length $N$, they showed an index using space $\tilde{O}(N^{1.5})$ with $\tilde{O}(N\sqrt{N})$ preprocessing time and $\tilde{O}(|P_1| + |P_2| + \sqrt{N})$ query time.

- By establishing a connection with Boolean Matrix Multiplication, they highlighted a difficulty in removing the $\sqrt{N}$ factor both from the preprocessing and the query time.

$\tilde{O}$ hides the logarithmic factors.

# Gapped consecutive indexing

**Gapped indexing for consecutive occurrences:** preprocess an text $T$ of length $N$ into a data structure, which allows, given a range $[a, b]$ and two patterns $P_1, P_2$, to retrieve all pairs of consecutive occurrences of $P_1, P_2$ separated by distance $d \in [a, b]$.

# Gapped consecutive indexing

**Gapped indexing for consecutive occurrences:** preprocess an $N$-length text $T$ into a data structure, which allows, given a range $[a, b]$ and two string patterns $P_1, P_2$, to retrieve all pairs of consecutive occurrences of $P_1, P_2$ separated by distance $d \in [a, b]$.

- [Navarro and Thankanchan 2016] For the case $P_1 = P_2$, $O(N \log N)$-space index with $O(|P_1| + |P_2| + \mathrm{occ})$ query time.

- [Bille et al. 2021] For the general case $P_1 \neq P_2$, no $\tilde{O}(N)$-space index can achieve $\tilde{O}(|P_1| + |P_2| + \sqrt{N})$ query time conditioned on the Set Disjointness conjecture.

**For highly compressible texts, can we design an efficient index for this problem?**

$\tilde{O}$ hides the logarithmic factors.

# Choice of compression method

**The answer, of course, depends on the chosen compression method…**

- We assume that the text is represented by a straight-line program (SLP), which is a context-free grammar describing exactly one string.

**Example:** The SLP $\{A \to BC, B \to ba, C \to DD, D \to na\}$ generates $banana$.

- SLPs are **capable of describing strings of exponential length** (in the size of the representation).

- Capture the popular **Lempel-Ziv compression method** up to a log factor.

- On the other hand, SLPs provide a convenient interface, allowing e.g. for efficient random access [Bille et al. 2015].

# Indexing in compressed space

Assuming that a string $T$ of length $N$ is described by an SLP with $g$ productions, there are multiple $\tilde{O}(g)$-space indexes for **classic pattern matching**:

| Space | Query time | Reference |
|---|---|---|
| $O(g)$ | $O(m \log \log N + \mathrm{occ} \log g)$ | Claude and Navarro 2012 |
| $O(g \log N)$ | $O((m + \mathrm{occ})\log g)$ | Claude et al. 2021 |
| $O(g \log N)$ | $O(m + \mathrm{occ} \log^{\varepsilon} N)$ | Christiansen et al. 2021 |
| $O(g \log N)$ | $O((m \log m + \mathrm{occ})\log g)$ | Díaz-Domínguez et al. 2021 |

# Lower bounds for compressed data

Some problems cannot avoid a high dependency on the size of an uncompressed string:

- Pattern matching with wildcards [Aboud et al. 2017]

- Longest common subsequence [Aboud et al. 2017]

- Median edit distance [Kociumaka et al. 2022]

- Center edit distance [Kociumaka et al. 2022]

**What about consecutive pattern matching?**

# Our results

**Grammar compressed indexing for gapped consecutive occurrences:** preprocess an $N$-length text $T$ given as a grammar of size $g$ into a data structure, which allows, given a range $[a, b]$ and two string patterns $P_1, P_2$, of length $m$ to retrieve all pairs of consecutive occurrences of $P_1, P_2$ separated by distance $d \in [a, b]$.

| Cases | Space | Query time |
|---|---|---|
| unbounded (a = 0, b = N) | $O(g^2 \log^4 N)$ | $O(m \log N + (1 + occ) \cdot \log^3 N \log \log N)$ |
| a = 0 | $O(g^5 \log^5 N)$ | $O(m \log N + (1 + occ) \cdot \log^4 N \log \log N)$ |

We obtain those results by using **locally consistent grammars** !

# Our results

**Gapped indexing for consecutive occurrences:**

preprocess an $N$-length text $T$ into a data structure, which allows, given a range $[a, b]$ and two string patterns $P_1, P_2$, of length $m$ to retrieve all pairs of consecutive occurrences of $P_1, P_2$ separated by distance $d \in [a, b]$.

| Cases | Space | time |
|---|---|---|
| unbounded (a = 0, b = N) | $O(g^2 \log^4 N)$ | $(1 + \text{occ}) \cdot \log^3 N \log \log N)$ |
| a = 0 | $O(g^5 \log^5 N)$ | $O(m \log N + (1 + \text{occ}) \cdot \log^4 N \log \log N)$ |

Achieving $\tilde{O}(g)$ space and $\tilde{O}(m + \text{occ})$ query time would contradict the lower bound of Bille et al. 2021

We obtain those results by using **locally consistent grammars** !

# Our results

**Gapped indexing for consecutive occurrences:**

preprocess an $N$-length text $T$ into a data structure, which allows, given a range $[a, b]$ and two string patterns $P_1, P_2$, of length ~~$m$~~ ~~all~~ pairs of consecutive occurrences of $P_1, P_2$ separate~~d~~ $b]$.

| Cases | Space | | Query time |
|:---:|:---:|:---:|:---:|
| unbounded (a = 0, b = N) | $O(g^2 \log^4 N)$ | | $O(m \log N + (1 + \mathrm{occ}) \cdot \log^3 N \log \log N)$ |
| a = 0 | $O(g^5 \log^5 N)$ | | $O(m \log N + (1 + \mathrm{occ}) \cdot \log^4 N \log \log N)$ |

In the unbounded case, it might be possible to achieve $\tilde{O}(g)$ space and $\tilde{O}(m + \mathrm{occ})$ query time

We obtain those results by using **locally consistent grammars** !

# Run-length SLP

**Run-length SLP** a set of non-terminals, a set of terminals, an initial symbol, and a set of productions, such that:
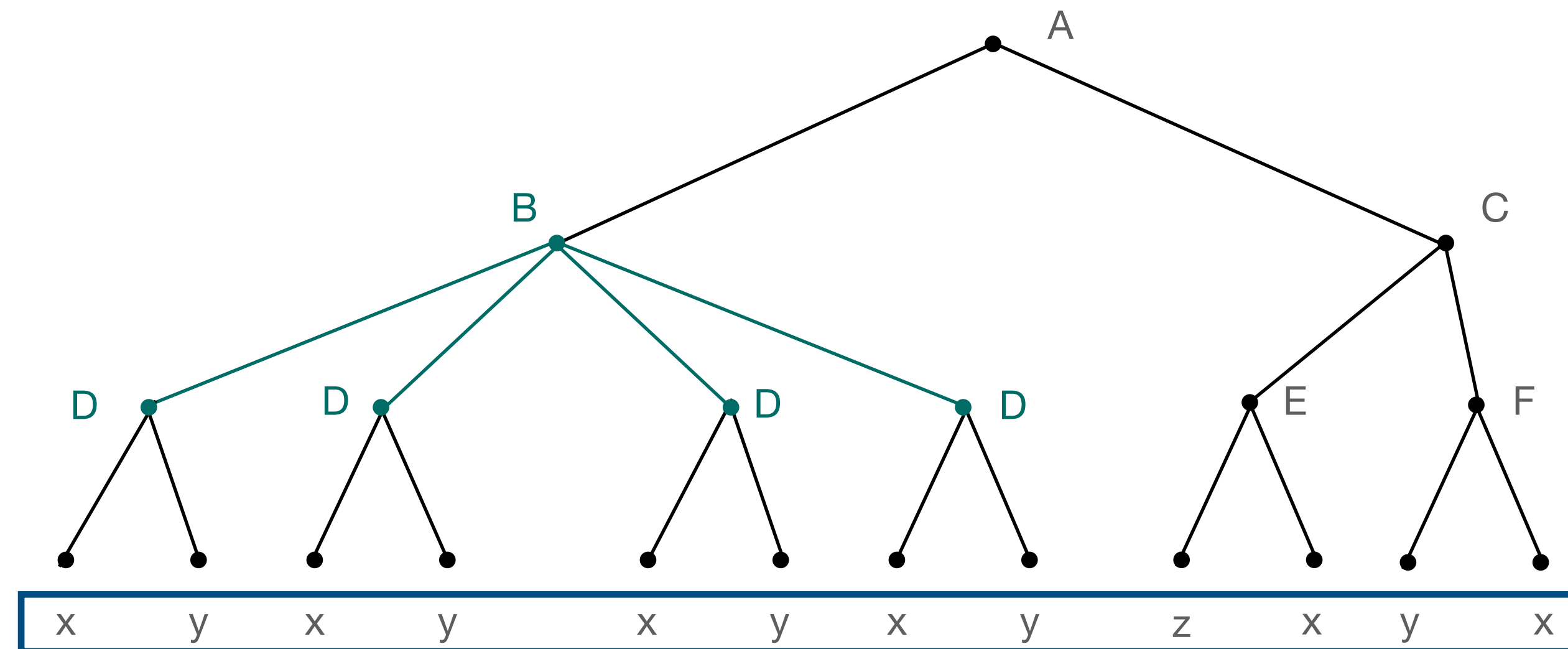
- Each production has form $A \to BC$ or $A \to B^k$, where $A$ is a non terminal and $B, C$ can be either terminals or non-terminals.

- Every non-terminal is on the left-hand side of exactly one production (=> it generates exactly one string).

**Expansion(S)**, is the string "generated" by the non-terminal S.

The string obtained by iterative replacement of non-terminals by the right-hand sides of the production rules, until only terminals remain.

**A run-length SLP describes the expansion of its initial symbol.**
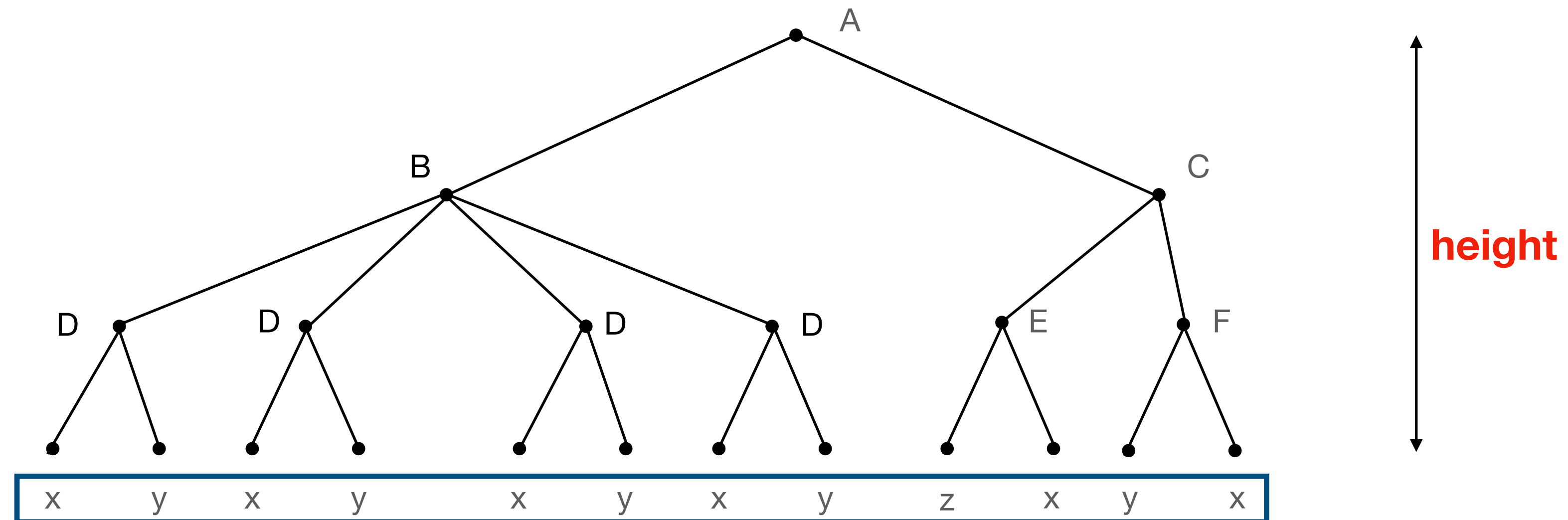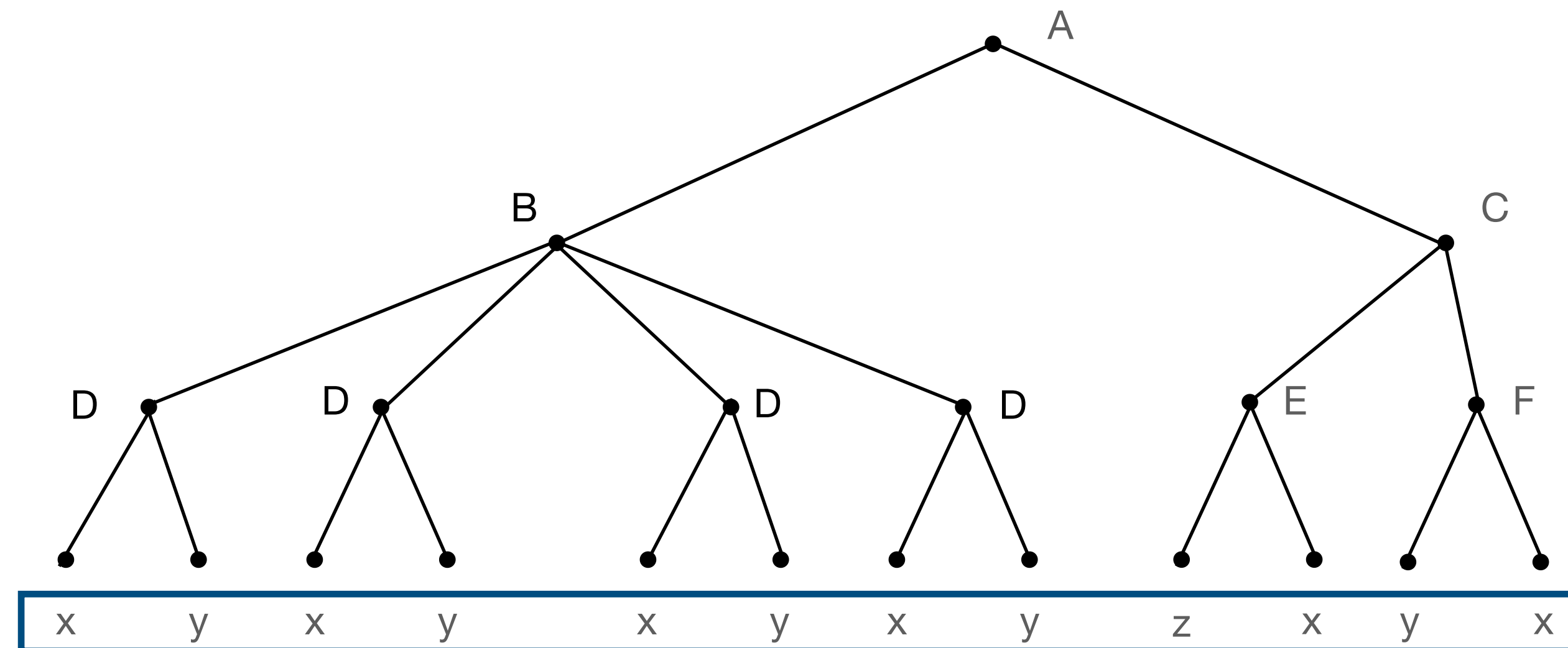
# Run-length SLP



**The parse tree of a run-length SLP**

non-terminals $= \{A, B, C, D, E, F\}$ and terminals $= \{x, y, z\}$

productions: $A \rightarrow BC, B \rightarrow D^4, D \rightarrow xy, C \rightarrow EF, E \rightarrow zx, F \rightarrow yx$

# Run-length SLP



**The parse tree of a run-length SLP**

non-terminals $= \{A, B, C, D, E, F\}$ and terminals $= \{x, y, z\}$

productions: $A \to BC, B \to D^4, D \to xy, C \to EF, E \to zx, F \to yx$

**size**

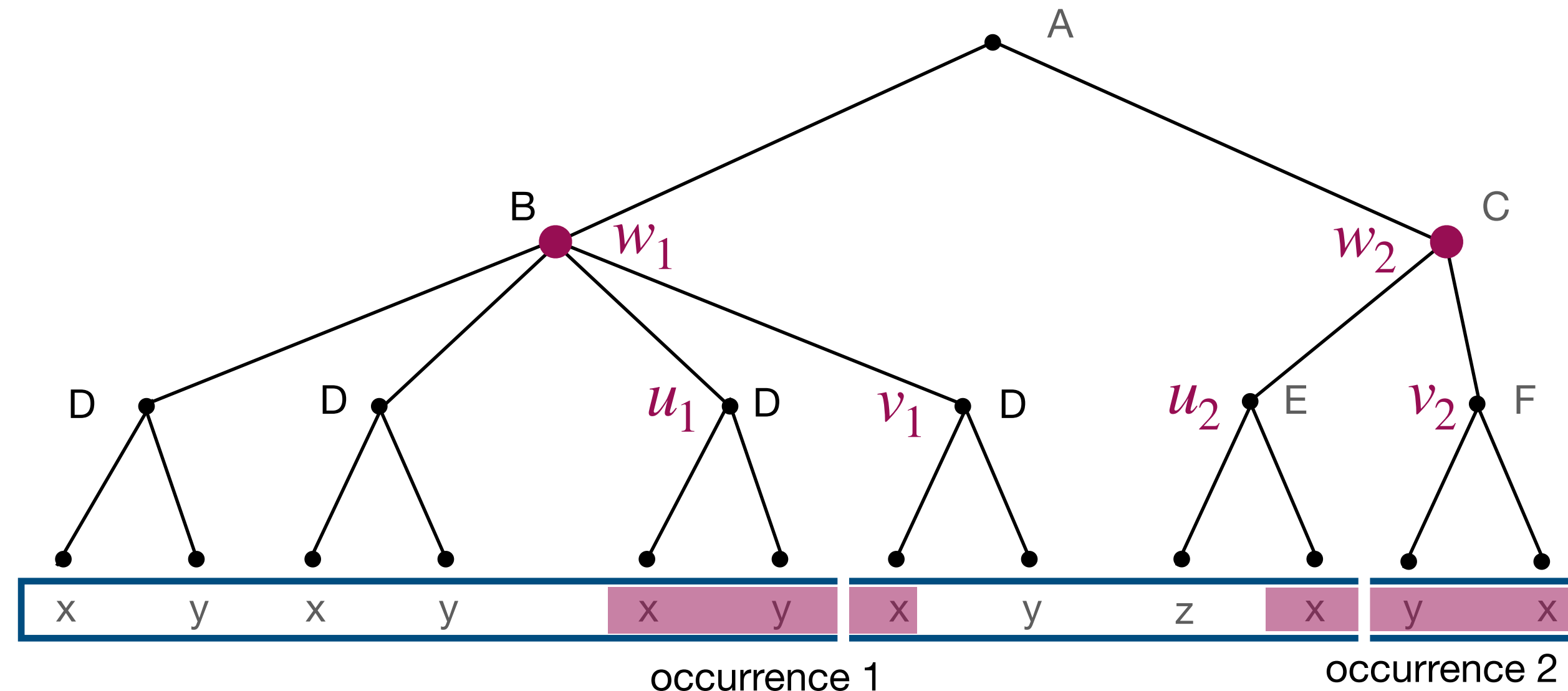# SLP to a locally-consistent run-length SLP



**Corollary of [Gawrychowski et al. 2018]**

There is a Las-Vegas algorithm that converts an SLP of size $g$ describing a string $T$ of length $N$ into a **locally-consistent** run-length SLP of size $O(g \log N)$ and height $O(\log N)$ describing the same string $T$ in $O(g \log N)$ time.

# SLP to a locally-consistent run-length SLP



**Locally-consistent run-length SLP**

There are only $O(\log N)$ possible splits of $P$, and we can compute them in $O(|P|\log N)$ time
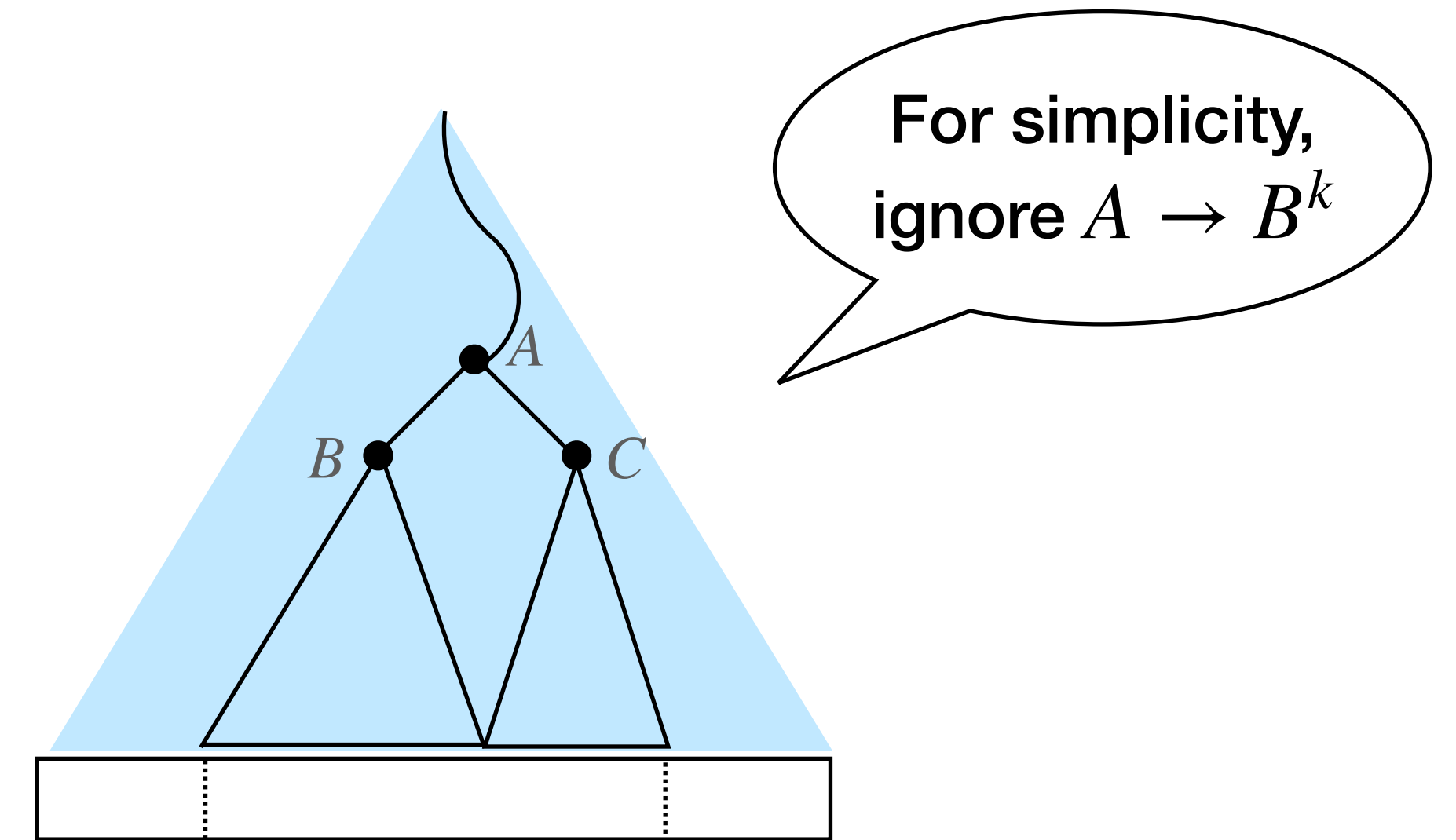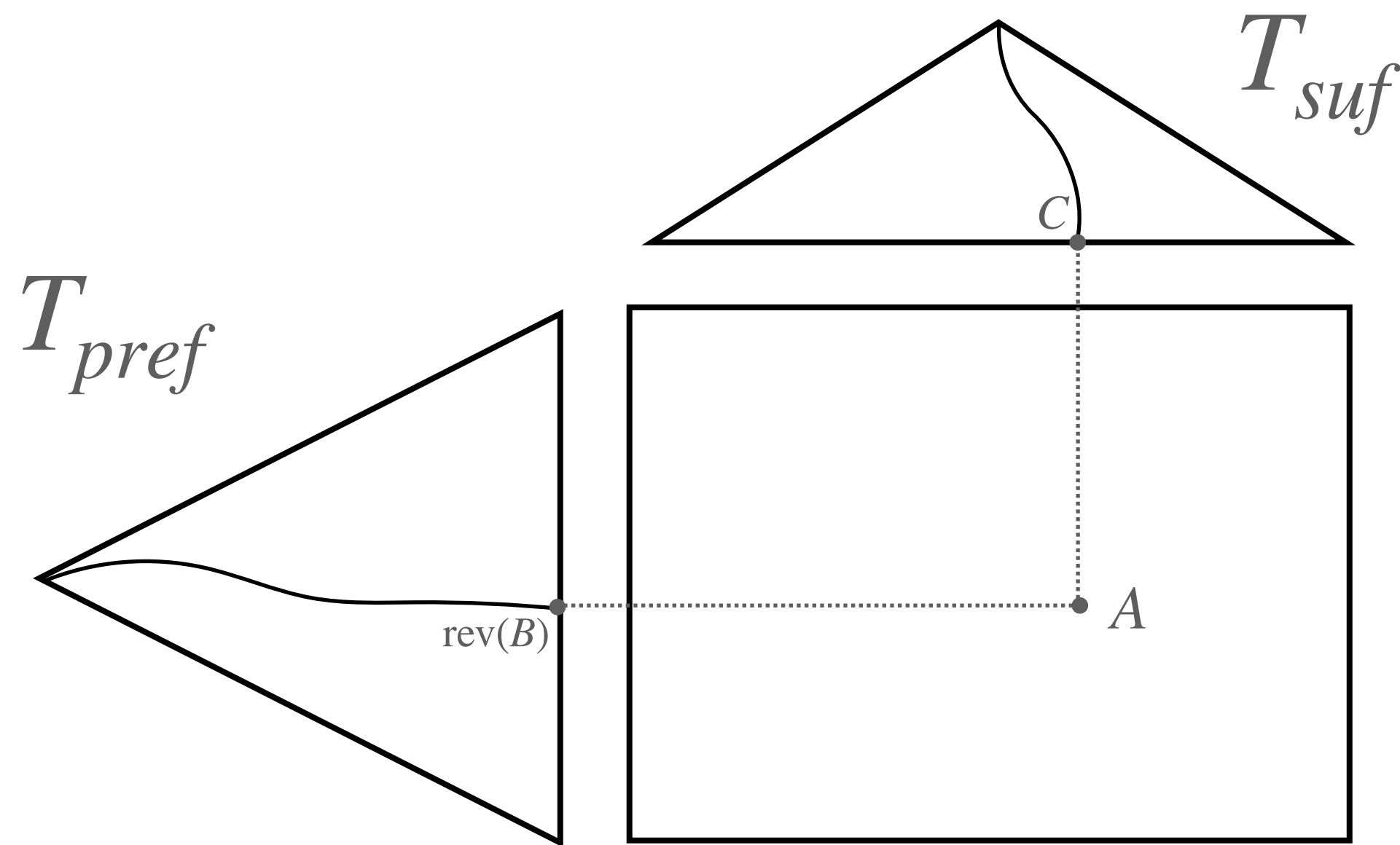
**Definition: Split**

For an occurrence $T[\ell, r]$ of a pattern $P$, let $w$ be the lowest node of a parse tree containing it, we say that $T[\ell, r]$ is **relevant** for the label of $w$ (a non-terminal).

$T[\ell, r]$ is **split at position** $i$ if there exist children $u, v$ of $w$ such that $T[\ell, \ell + i]$ is contained in $u$ and $T[\ell + i + 1, r]$ in $v$.

**Example:** Occurrence 1  xyx is split at position  2 and relevant for B, occurrence 2 is split at position 1 and relevant for C.
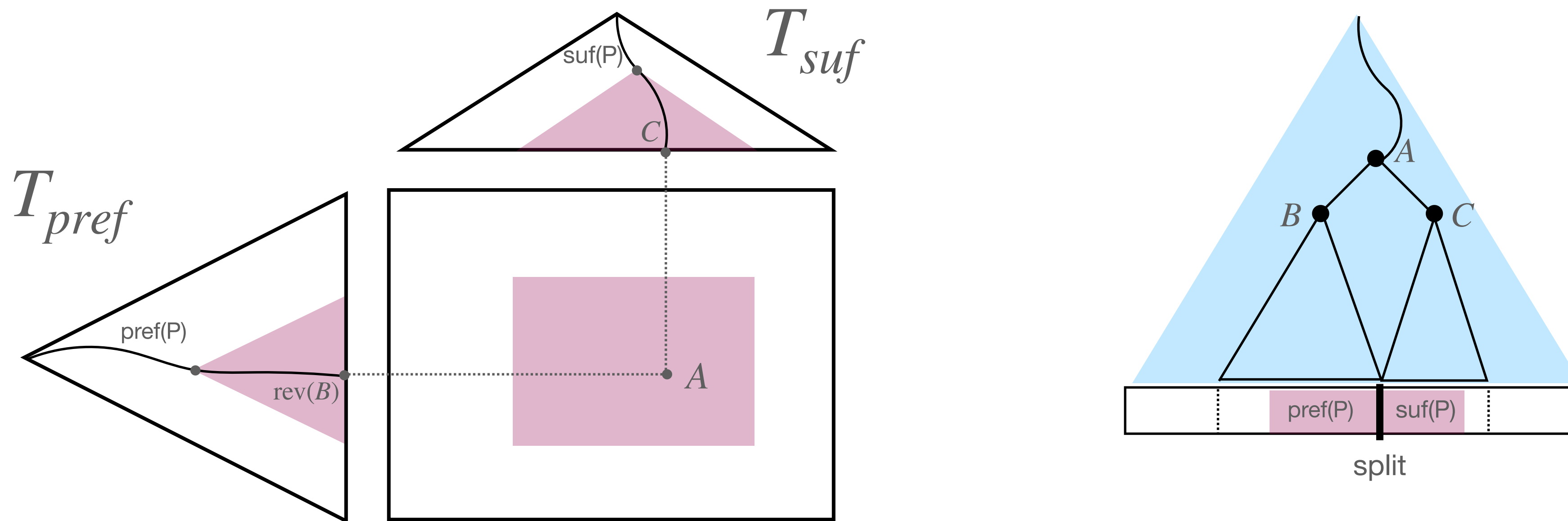
# Why is local consistency interesting?



…because we can search for relevant occurrences quickly using the following structure:

- For every $A \to BC$, add $\mathrm{rev}(\mathrm{expansion}(B))$ to $T_{pref}$ and $\mathrm{expansion}(C)$ to $T_{suf}$

- Create a point $(r_B, r_C)$ (the lexicographic rank of the expensions) for every $A \to BC$

- Build an orthogonal range data structure on the points

# Why is local consistency interesting?



To find relevant occurrences of a pattern $P$ in non-terminals:

- For each split s, search for $\text{pref}(P) = \text{rev}(P[1..s])$ in $T_{pref}$ to obtain an interval $I_{pref}$ of leaves starting with it, and for $\text{suf}(P) = P[s+1..]$ in $T_{suf}$ to obtain an interval $I_{suf}$

- Report all non-terminals in $I_{pref} \times I_{suf}$

# We show an even stronger result

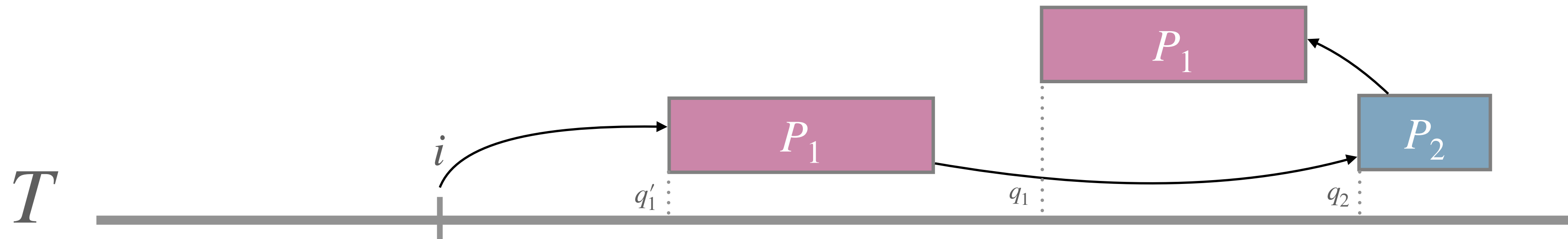For a run-length SLP representing a string $T$ of length $N$, with size $g$ and height $O(\log n)$,

There is **a** $O(g^2 \log^2 N)$**-space data structure** that preprocesses an $m$-length pattern $P$ in $O(m \log N + \log^2 N)$ time and can, answer the following queries in $\text{polylog } N$ time:

For a given non terminal $A$, in $\text{expansion}(A)$,

- Report relevant occurrences of $P$;

- Decide whether there is an occurrence of $P$;

- Report the leftmost/rightmost occurrence of $P$ ;

- Find a predecessor/successor occurrence of $P$ given a position $q$.

# Corollary: unbounded case

**Task:** report all consecutive occurrences of $P_1, P_2$ in a $N$-length string $T$ described by a run-length SLP of size $g$ and height $\log N$.



1. Find the leftmost occurrence $q_1'$ of $P_1$ in $T[i\,..\,]$ (successor)

2. Find the leftmost occurrence $q_2 \geq q_1'$ of $P_2$ (successor)

3. Find the rightmost occurrence $q_1 \leq q_2$ of $P_1$ (predecessor)
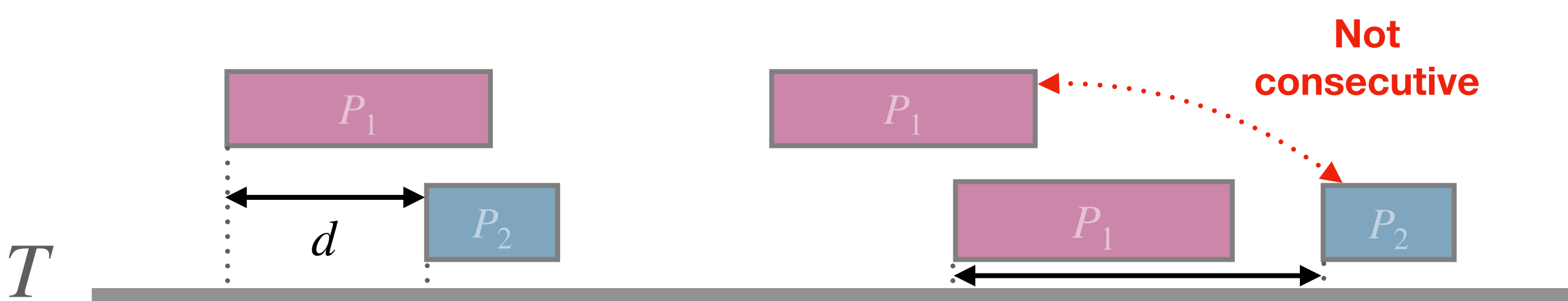
4. Report $(q_1, q_2)$ and set $i = q_2 + 1$

**Time** $\tilde{O}(m + (\text{occ} + 1)\text{polylog } N)$,

**space** $\tilde{O}(g^2)$**.**

# Idea of our index for the case $a = 0$

<div style="border:1px solid black; padding:10px;">

**Task:** given an $N$-length string $T$ described by a run-length SLP of size $g$ and height $\log N$, report all consecutive occurrences of $P_1, P_2$ in $T$ separated by distance in $[0,b]$.

</div>

- For each non-terminal of the grammar, retrieve relevant consecutive occurrences separated by distance in $[0,b]$.

- Generate all consecutive occurrences separated by distance in $[0,b]$ by traversing a pruned parse tree of the grammar (using a standard technique borrowed from classic pattern matching compressed-space indexes).

# Summary



**Complex matching: Gapped consecutive occurrences:** given a range $[a, b]$ and two string patterns $P_1, P_2$, retrieve all pairs of consecutive occurrences of $P_1, P_2$ separated by distance $d \in [a, b]$.

**Sketch as input:** Grammar compressed input (local consistency preserves splits)

**Indexing:** preprocess a text $T$ of length $N$ given as grammar $g$ into a data structure.

| Case | Space | Query time |
|------|-------|-----------|
| unbounded (a = 0, b = N) | $O(g^2 \log^4 N)$ | $O(m \log N + (1 + \mathrm{occ}) \cdot \log^3 N \log \log N)$ |
| a = 0 | $O(g^5 \log^5 N)$ | $O(m \log N + (1 + \mathrm{occ}) \cdot \log^4 N \log \log N)$ |

Can better space be achieved?  Solution for the general case?

**Is the dual problem of consecutive compressed pattern matching easier?**

# Grammar compressed consecutive pattern matching

**Dual problem:** Given $T$ of length N as grammar of size $g$, a range $[a, b]$, and two string patterns $P_1, P_2$, retrieve all pairs of consecutive occurrences of $P_1, P_2$ separated by distance $d \in [a, b]$. **Process the text and the patterns at the same time !**



**If $T$ is given uncompressed**: we can just go from left to right, keeping track of the most recent occurrences of $P_1$ and $P_2$.

$$\implies O(|T| + |P_1| + |P_2| + occ) \text{ time algorithm.}$$

# Grammar compressed pattern matching

For **a single pattern** $P$, matching in a text $T$ given as a grammar of size $g$, **[Ganardi & Gawrychowski, SODA'22]** showed that we can detect whether $P$ occurs in $T$ in $O(g + |P|)$ **time.**

**Can we extend this result to two patterns consecutive (and reporting)? Yes!**

**Gawrychowski, Gourdel, Starikovskaya, Steiner (Unpublished)**

Given $T$ of length N as grammar of size $g$, and two string patterns $P_1, P_2$, we can report all consecutive occurrences of $P_1, P_2$ in $T$ in $O(g + |P_1| + |P_2| + occ)$ time.
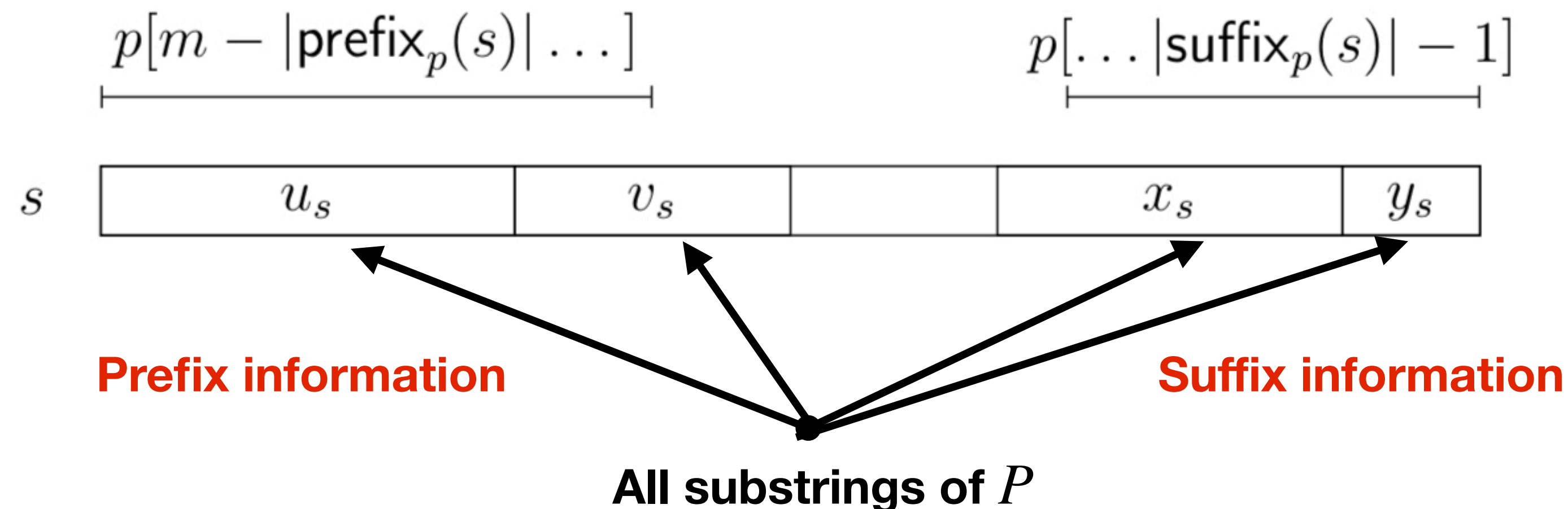
# Boundary information

For a pattern $P$, the $P$-boundary information of a string $S$ is substrings occurring both in $P$ and $S$:

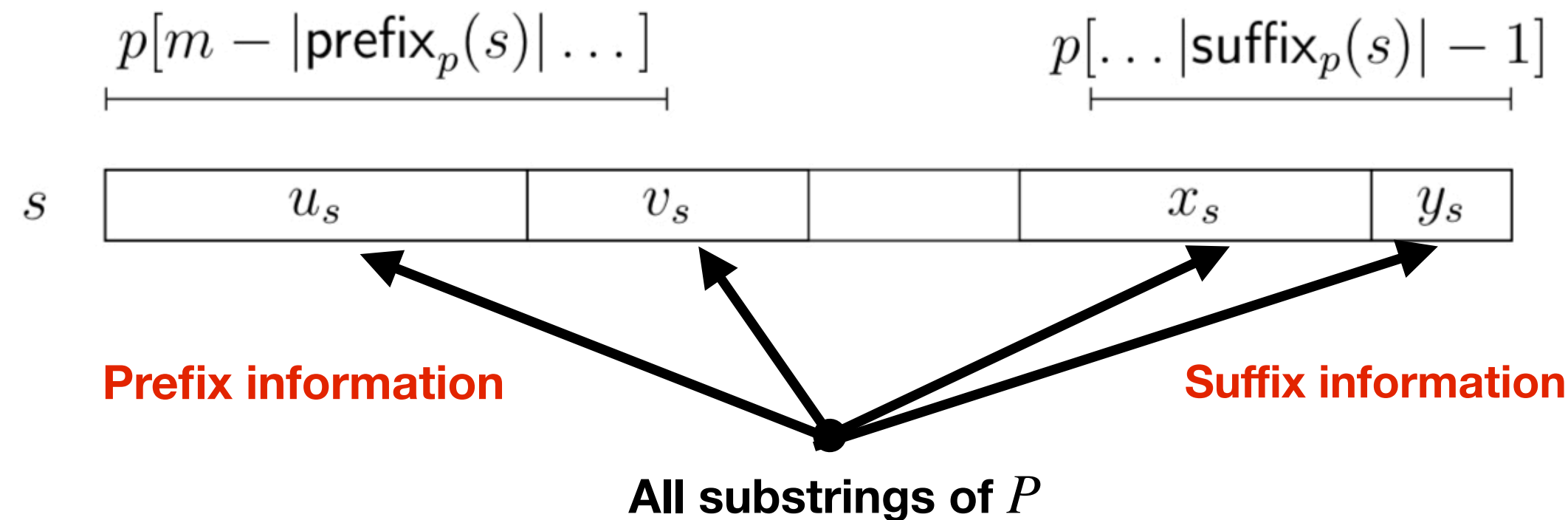If $S$ occurs in $P$, then the position where it occurs is the $P$-substring information.

Else, let $\mathrm{prefix}_P(S)$ be the **longest prefix** of $S$ which is a suffix of $P$.

let $\mathrm{suffix}_P(S)$ be the **longest suffix** of $S$ which is a prefix of $P$.

$$p[m - |\mathsf{prefix}_p(s)| \ldots] \qquad p[\ldots |\mathsf{suffix}_p(s)| - 1]$$

$s$  | $u_s$ | $v_s$ | | $x_s$ | $y_s$ |

**Prefix information**     **Suffix information**

**All substrings of $P$**

$P$-boundary information of two strings $S$ and $T$ allows to efficiently report new occurrences of $P$ appearing in $ST$.
+ the boundary information for $ST$ can be computed quickly.

# (Secondary) Boundary information



$p[m - |\text{prefix}_p(s)| \ldots]$

$p[\ldots |\text{suffix}_p(s)| - 1]$

$s$ | $u_s$ | $v_s$ | | $x_s$ | $y_s$

**Prefix information**

**Suffix information**

**All substrings of $P$**

$P$-boundary information of two strings $S$ and $T$ allows to efficiently report new occurrences of $P$ appearing in $ST$.
+ the boundary information for $ST$ can be computed quickly.

---

For an SLP rule $A \rightarrow BC$, Compute bottom to top:

- $P_1$-boundary information and $P_2$-boundary information for $\overline{A}$ (from the boundary informations for $\overline{B}$ and $\overline{C}$).

- All crossing occurrences of $P_1, P_2$ in $A$.

- The rightmost occurrences of $P_1, P_2$ in $\overline{A}$.

**Enough to detect pirmary co-occ in $O(g + |P_1| + |P_2|)$ time !**

- If the $P_2$-suffix information for $\overline{A}$ is $(x_A, y_A)$: $P_1$-boundary information for $x_A$ and $y_A$, all crossing occurrences of $P_1$ and leftmost rightmost.

# Summary

**Complex matching:** **Consecutive occurrences:** given two string patterns $P_1, P_2$, retrieve all pairs of consecutive occurrences of $P_1, P_2$ .

**Sketch as input:** Grammar compressed input (handled efficiently through boundary information)

**Pattern matching:** process $P$ and $T$ at the same time.

**Gawrychowski, Gourdel, Starikovskaya, Steiner (Unpublished)**

Given $T$ of length N as grammar of size $g$, and two string patterns $P_1, P_2$, we can report all consecutive occurrences of $P_1, P_2$ in $T$ in $O(g + |P_1| + |P_2| + occ)$ time.

**Corollary:** Gapped consecutive matching in $O(g + |P_1| + |P_2| + occ)$ time and Top-$k$ closest occurrences $O(g + |P_1| + |P_2| + k)$ time.