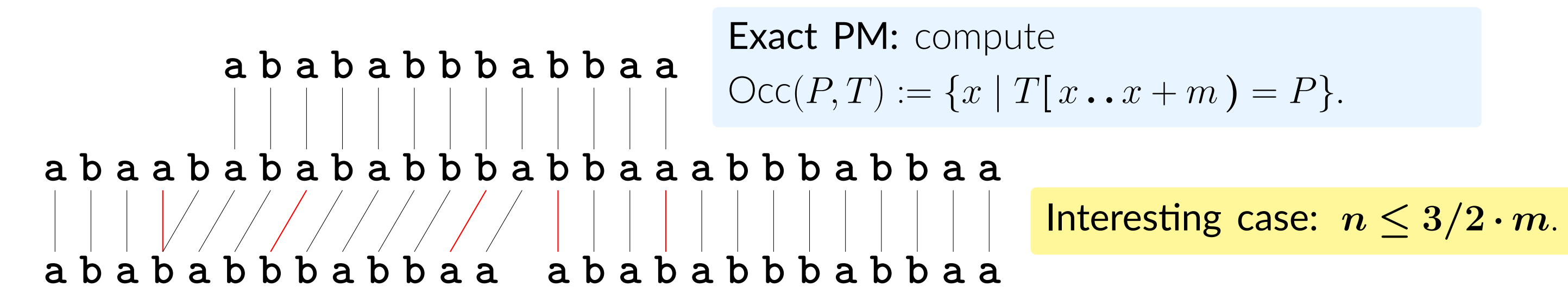# On the Communcation Complexity of Approximate Pattern Matching

Tomasz Kociumaka[1], **Jakob Nogler**[2], and Philip Wellnitz[3]

[1]Max Planck Institute for Informatics, SIC, Germany     [2]ETH Zurich, Switzerland     [3]National Institute of Informatics, SOKENDAI, Japan

## Pattern Matching (PM)

Given a *text* $T$ and a *pattern* $P$ of length $|T| = n$ and $|P| = m$, find all substrings of $T$ where $P$ occurs exactly (exact PM) or with few errors (approximate PM).



**Exact PM:** compute
$\text{Occ}(P,T) := \{x \mid T[x\mathbin{.\,.}x+m) = P\}.$

Interesting case: $n \leq 3/2 \cdot m$.

**PM with edits:** compute
$\text{Occ}_k^E(P,T) := \{x \mid \exists y\ \text{ED}(T[x\mathbin{.\,.}y),P) \leq k\}.$

**PM with mismatches:** compute
$\text{Occ}_k^H(P,T) := \{x \mid \text{HD}(T[x\mathbin{.\,.}x+m),P) \leq k\}.$

## Communication Complexity

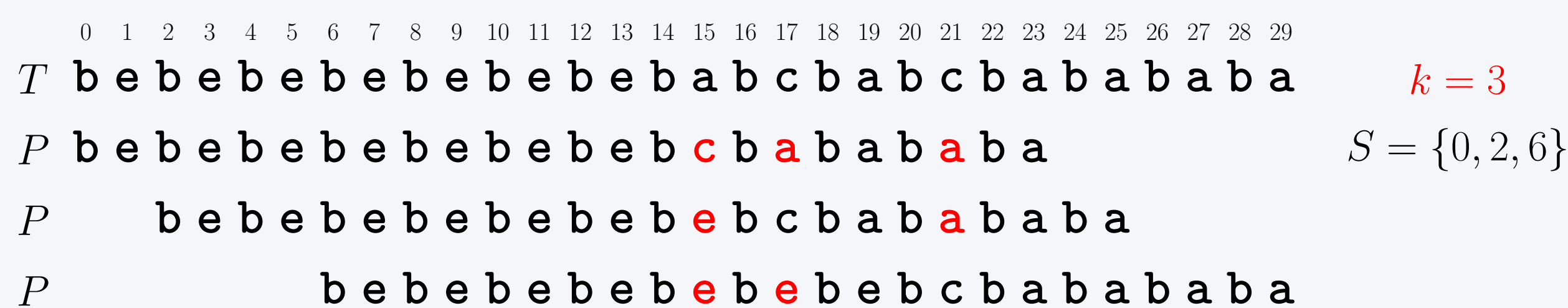Alice ———— One-way protocol ————→ Bob

① Alice receives a PM instance.
*Text $T$, Pattern $P$, Threshold $k$*

② Alice compresses the input.

③ Alice sends compressed data to Bob.

④ Bob needs to reconstructs the output of the instance.
*Set $\text{Occ}_k^E(P,T)$*

Communication Complexity = "minimum # of machine words to send to Bob"

|  | UB | LB | Refererence |
|---|---|---|---|
| Exact PM | $\mathcal{O}(1)$ | $\Omega(1)$ | Periodicity Lemma, [FW65] |
| PM with mismatches | $\mathcal{O}(k)$ | $\Omega(k)$ | [CKP19] |
| PM with edits | $\mathcal{O}(k^3)$ | | [CKW20] |
| PM with edits | $\mathcal{O}(k \log m)$ | $\Omega(k)$ | This Work |

## How to Achieve $\mathcal{O}(k \log m)$ for PM with Mismatches
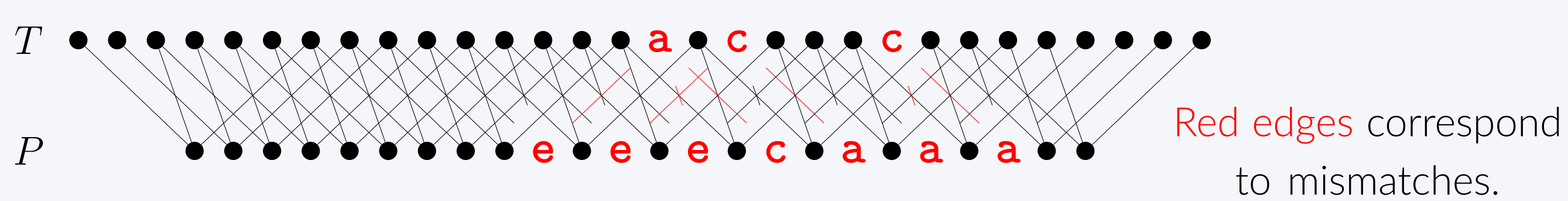
How does Alice encode $\text{Occ}_k^H(P,T)$?



Alice sends $S$, $\{(15,\mathsf{c},\mathsf{a}),(17,\mathsf{a},\mathsf{c}),(21,\mathsf{a},\mathsf{c})\}$, $\{(13,\mathsf{e},\mathsf{a}),(19,\mathsf{a},\mathsf{c})\}$, and $\{(9,\mathsf{e},\mathsf{a}),(11,\mathsf{e},\mathsf{c})\}$.

1. Crops $T$ s.t. $\{0, n-m\} \subseteq \text{Occ}_k^H(P,T)$.
2. Selects set $\{0, n-m\} \subseteq S \subseteq \text{Occ}_k^H(P,T)$ s.t. $\gcd(S) = \gcd(\text{Occ}_k^H(P,T)) =: g$.
3. Sends $S$ and the *mismatch information* for every $i \in S$ to Bob, i.e.,
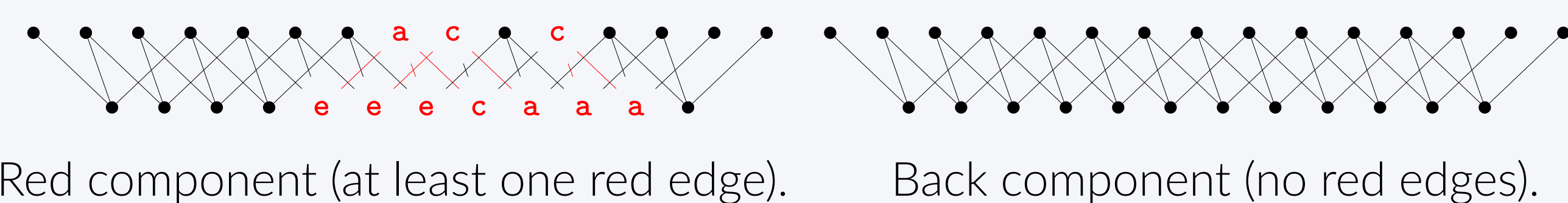
$$\{(j, P[j], T[i+j]) \mid P[j] \neq T[i+j]\}.$$

**Claim:** Alice can choose $S$ s.t. $|S| = \mathcal{O}(\log m)$.

How does Bob decode Alice's message?

1. Bob constructs the graph $\mathbf{G}_S$.
   $V =$ characters of $P$ & $T$
   $E = \{\{P[j], T[i+j]\} \mid i \in S, j \in [0\mathbin{.\,.}m)\}$
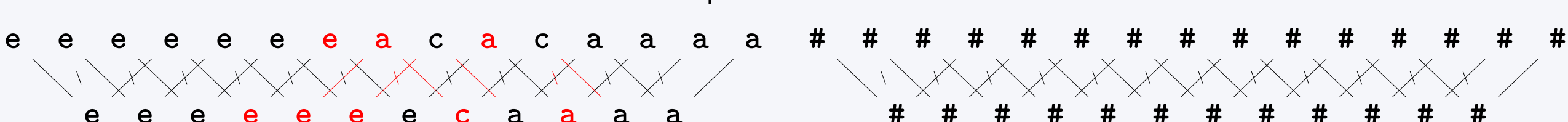


Red edges correspond to mismatches.

2. Bob divides $\mathbf{G}_S$ into *black and red connected components*.



Red component (at least one red edge).     Back component (no red edges).

**Claim:** For every remainder $c \in [0\mathbin{.\,.}g)$ modulo $g$ there is a connected components consisting of all vertices $P[i]$ and $T[i]$ with $i \equiv_g c$.

3. Bob propagates characters in red connected components, and replaces characters in black connected components with a sentinel character #.
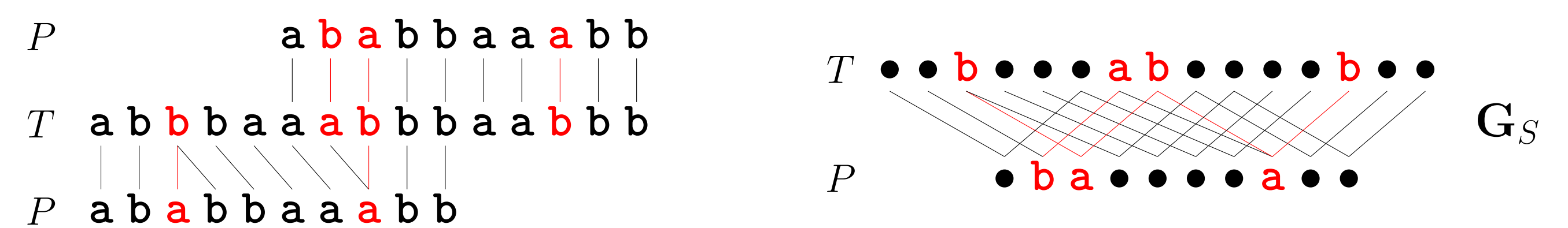


4. Bob thereby obtains $T^\#$ and $P^\#$.

$T^\#$ #e#e#e#e#e#e#e#a#c#a#c#a#a#a#a

$P^\#$ #e#e#e#e#e#e#e#c#a#a#a#a

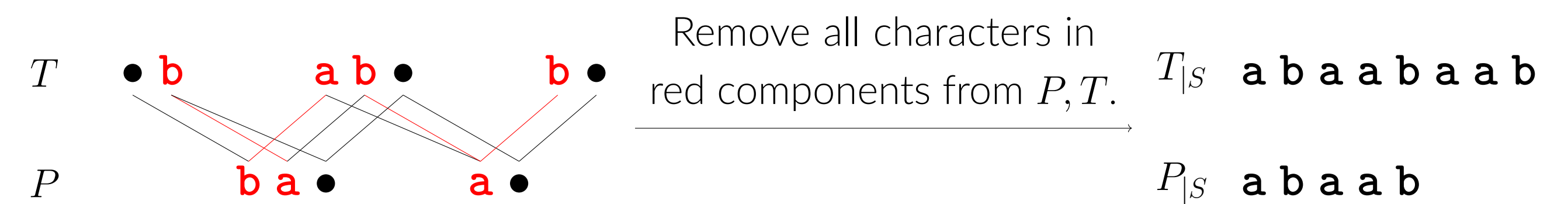5. Bob returns $\text{Occ}_k^H(P^\#, T^\#)$.

**Claim:** $\text{Occ}_k^H(P^\#, T^\#) = \text{Occ}_k^H(P,T)$.

## Extending the Techinques to PM with Edits

We store a set $S$ of $\mathcal{O}(\log m)$ $k$-edits occurences (including a prefix and a suffix), along with the information for their edits. We use them to construct a graph $\mathbf{G}_S$.
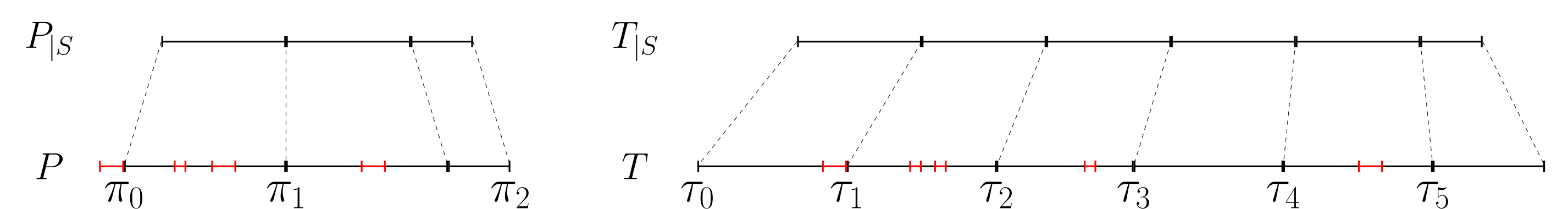


What changes?



Remove all characters in red components from $P, T$.

$T_{|S}$  abaabaab

$P_{|S}$  abaab

Red connected components do not have a periodic structure anymore.

Black connected components have a periodic structure in $P_{|S}$ and $T_{|S}$.

**Claim:** $P_{|S}, T_{|S}$ are periodic with period $\text{bc}(\mathbf{G}_S)$ (# number of black components). Moreover, for every remainder $c \in [0\mathbin{.\,.}\text{bc}(\mathbf{G}_S))$, there is a connected components consisting of all vertices $P[i]$ and $T[i]$ with $i \equiv_{\text{bc}(\mathbf{G}_S)} c$.
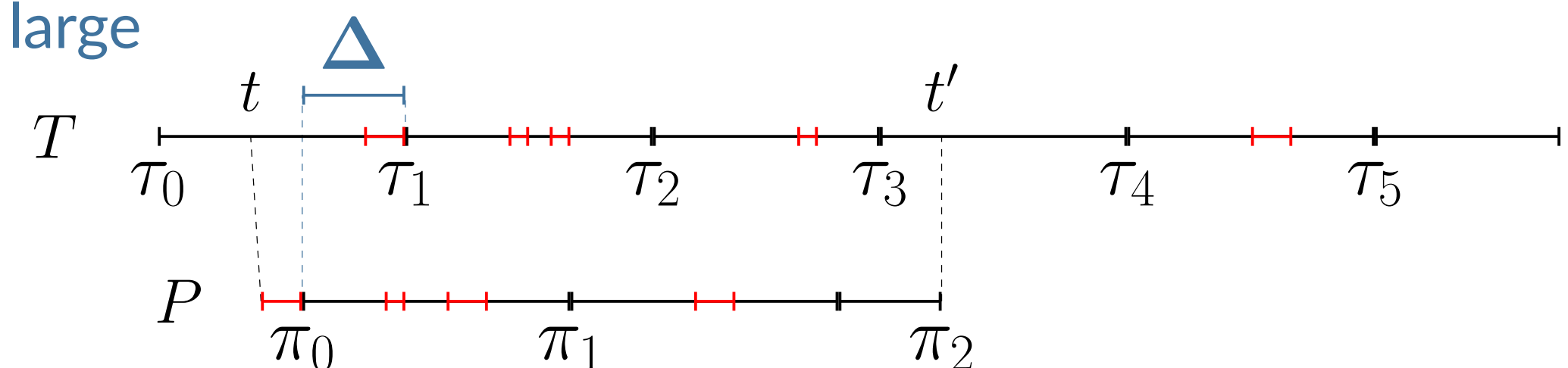
How do we keep $S$ small? We add alignments iteratively to $S$. For each added alignment we aim to at least halve the period of $P_{|S}$ and $T_{|S}$.

For which new alignments added to $S$ does the period of $P_{|S}$ and $T_{|S}$ decrease?
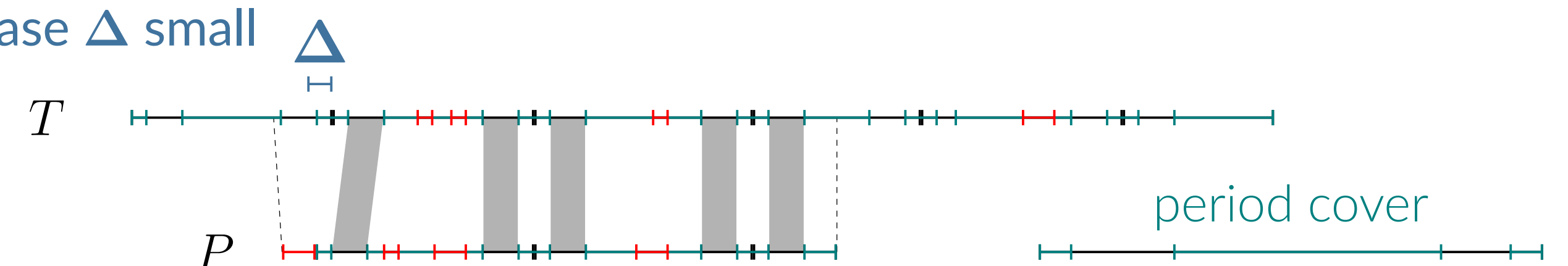


For a new $k$-edit alignment $P \rightsquigarrow T[t\mathbin{.\,.}t')$ we consider $\Delta = \min_i |\tau_i - t - \pi_0|$.

① Case $\Delta$ large



Then, any node in a black components is far away from any another node in the same component. If the alignment is added to $S$, every black component (if it not becomes red) is merged with another one

② Case $\Delta$ small



period cover

Then, we learn a set of character contained in black components around characters in red components (*period cover*). Characters outside from the period cover always match with characters belonging to the same black component.

**Why does this work?** New results linking "close alignments" and compressibility (see for example [CKW23], [GJKT24], [GK24]).

## Application: PM with Edits in the Quantum Setting

There is a quantum algorithm that, given a PM with edits instance:

- computes $\text{Occ}_k^E(P,T)$ using $\widehat{\mathcal{O}}(n/m \cdot \sqrt{km})$ queries and $\widehat{\mathcal{O}}(n/m \cdot (\sqrt{km} + k^{3.5}))$ time;
- decides whether $\text{Occ}_k^E(P,T) \neq \emptyset$ using $\widehat{\mathcal{O}}(\sqrt{n/m} \cdot \sqrt{km})$ queries and $\widehat{\mathcal{O}}(\sqrt{n/m} \cdot (\sqrt{km} + k^{3.5}))$ time.

The number of queries is optimal for small $k$ (up to a subpolynomial factor).

Main ingredients of the algorithm:

- Structural insights for approximate pattern matching [CKW20].
- Quantum algorithm for computing bounded edit distance [GJKT24].
- Quantum *Gap Edit Distance* algorithm (adapted from [GKKS22]).

Where do we use the previous results?
Suppose that we have a set $\text{Occ}_k^E(P,T) \subseteq H \subseteq [0\mathbin{.\,.}n)$ of *candidate positions*. We can retrieve $\text{Occ}_k^E(P,T)$ by only testing $\mathcal{O}(\log n)$ of the positions in $H$:

① Keep a set of $K$-edit alignments for $K > k$ (but not too large).

② Discard positions in $H$ "covered" by $S$. For the others run a $(k,K)$-Gap Edit Distance oracle (faster than computing the edit distance).

③ If there is one with dist. $\leq K$, retrieve the edits and repeat ② until $H = \emptyset$.