

## 5 Fast filtering algorithms

This exposition is based on

1. Flexible Pattern Matching in Strings, Navarro, Raffinot, 2002, chapter 6.5, pages 162ff.

We present the **hierarchical filtering** approach called **PEX** of Navarro and Baeza-Yates.

### 5.1 Filtering algorithms

The **idea** behind **filtering algorithms** is that it might be easier to **check** that a text **position does not match** a pattern string than to verify that it does.

Filtering algorithms **filter out** portions of the text that **cannot possibly contain a match**, leaving positions that could match.

The **potential match positions** then need to be **verified** with another algorithm like for example the bit-parallel algorithm of Myers (BPM).

**Filtering algorithms** are very sensitive to the **error level**  $\alpha := k/m$  since this normally affects the **amount of text** that can be **discarded** from **further consideration**. ( $m$  = pattern length,  $k$  = errors.)

**If most of the text has to be verified, the additional filtering steps are an overhead** compared to the strategy of just verifying the pattern in the first place.

On the other hand, if large portions of the text can be discarded quickly, then the filtering results in a faster search.

**Filtering algorithms** can **improve** the **average-case** performance (sometimes dramatically), but **not** the **worst-case** performance.

Assume that we want to **find all occurrences** of a pattern  $P = p_1, \dots, p_m$  in a text  $T = t_1, \dots, t_n$  that have an edit distance of at most  $k$ .

**If we divide** the **pattern** into  $k+1$  pieces  $P = p^1, \dots, p^{k+1}$ , then at least **one** of the pattern pieces **match without error**.

There is a more general version of this principle first formalized by **Myers** in 1994:

**Lemma 1.** Let **Occ match**  $P$  with  $k$  errors,  $P = p^1, \dots, p^j$  be a concatenation of subpatterns, and  $a_1, \dots, a_j$  be nonnegative integers such that  $A = \sum_{i=1}^j a_i$ . Then, for some  $i \in \{1, \dots, j\}$ , **Occ includes a substring that matches**  $p^i$  **with at most**  $\lfloor a_i k / A \rfloor$  errors.

**Proof:** Exercise.

So the **basic procedure** is:

1. **Divide:** **Divide** the **pattern** into  $k+1$  pieces of approximately the same length.
2. **Search:** **Search** all the **pieces simultaneously** with a **multi-pattern string matching algorithm**. According to the above lemma, **each possible occurrence will match at least one of the pattern pieces**.
3. **Verify:** For **each found pattern piece**, check the neighborhood with a **verification algorithm** that is able to detect an occurrence of the **whole pattern** with edit distance at most  $k$ . Since we allow indels, if  $p_{i_1} \dots p_{i_2}$  matches the text  $t_j \dots t_{j+i_2-i_1}$ , then the **verification** has to **consider** the text area  $t_{j-(i_1-1)-k} \dots t_{j+(i_2-i_1)+k}$ , which is of **length**  $m+2k$ .

### 5.2 An example

Say we want to find the pattern *annual* in the texts

$t_1 = \text{any\_annealing}$  and

$t_2 = \text{an\_unusual\_example\_with\_numerous\_verifications}$

with at most **2 errors**.

1. *Divide*: We divide the pattern *annual* into  $p^1 = an$ ,  $p^2 = nu$ , and  $p^3 = al$ . One of these subpattern has to match with 0 errors.
2. *Search*: We search for all subpatterns:
  - 1: searching for *an*: in  $t_1$ : find positions 1, 5  
in  $t_2$ : find position 1
  - 2: searching for *nu*: in  $t_1$ : find no positions  
in  $t_2$ : find positions 5, 25
  - 3: searching for *al*: in  $t_1$ : find position 9  
in  $t_2$ : find position 9
3. *Verification*: We have to verify 3 positions in  $t_1$ , and 4 positions in  $t_2$ , to find 3 occurrences in  $t_1$  and none in  $t_2$ .

## 5.3 Hierarchical verification

The toy example makes clear that many verifications can be triggered that are unsuccessful and that many subpatterns can trigger the same verification. Repeated verifications can be avoided by carefully sorting the occurrences of the pattern.

It was shown by Baeza-Yates and Navarro that the running time is dominated by the multipattern search for error levels  $\alpha = k/m$  below  $1/(3 \log_{|\Sigma|} m)$ . In this region, the search cost is about  $O(kn \frac{\log_{|\Sigma|} m}{m})$ . For higher error levels, the cost for verifications starts to dominate, and the filter efficiency deteriorates abruptly.

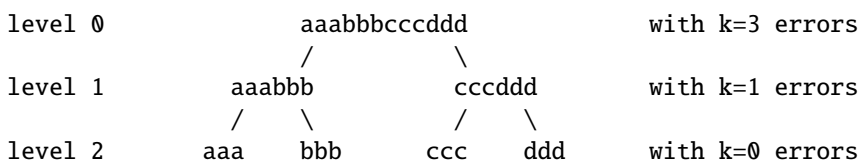
Baeza-Yates and Navarro introduced the idea of hierarchical verification to reduce the verification costs, which we will explain next. Then we will work out more details of the three steps.

Navarro and Baeza-Yates use Lemma 1 for a hierarchical verification. The idea is that, since the verification cost is high, we pay too much for verifying the whole pattern each time a small piece matches. We could possibly reject the occurrence with a cheaper test for a shorter pattern.

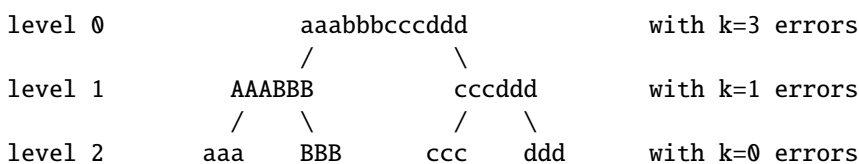
So, instead of directly dividing the pattern into  $k+1$  pieces, we do it hierarchically. We split the pattern first in two pieces and search for each piece with  $\lfloor k/2 \rfloor$  errors, following Lemma 1. The halves are then recursively split and searched until the error rate reaches zero, i. e. we can search for exact matches.

With hierarchical verification the area of applicability of the filtering algorithm grows to  $\alpha < 1/\log_{|\Sigma|} m$ , an error level three times as high as for the naive partitioning and verification. In practice, the filtering algorithm pays off for  $\alpha < 1/3$  for medium long patterns.

**Example.** Say we want to find the pattern  $P = aaabbbcccd$  in the text  $T = xxxbbbx$  with at most  $k = 3$  differences. The pattern is split into four pieces  $p^1 = aaa$ ,  $p^2 = bbb$ ,  $p^3 = ccc$ ,  $p^4 = ddd$ . We search with  $k = 0$  errors in level 2 and find *bbb*.



Now instead of verifying the complete pattern in the complete text (at level 0) with  $k = 3$  errors, we only have to check a slightly bigger pattern (*aaabbb*) at level 1 with one error. This is much cheaper. In this example we can decide that the occurrence *bbb* cannot be extended to a match.



## 5.4 The PEX algorithm

**Divide:** Split pattern into  $k + 1$  pieces, such that each piece has equal probability of occurring in the text. If no other information is available, the uniform distribution is assumed and hence the pattern is divided in pieces of equal length.

**Build Tree:** Build a tree of the pattern for the hierarchical verification. If  $k + 1$  is not a power of 2, we try to keep the binary tree as balanced as possible.

Each node has two members *from* and *to* indicating the first and the last position of the pattern piece represented by it. The member *err* holds the number of allowed errors. A pointer *myParent* leads to its parent in the tree. (There are no child pointers, since we traverse the tree only from the leafs to the root.) An internal variable *left* holds the number of pattern pieces in the left subtree. *idx* is the next leaf index to assign. *plen* is the length of a pattern piece.

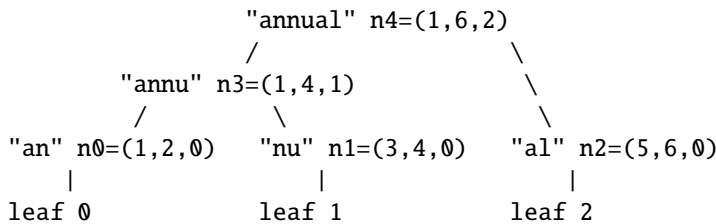
Algorithm CreateTree generates a hierarchical verification tree for a single pattern. (Lines 12 and 14 are justified by Lemma 1.)

```

(1) CreateTree(  $p = p_i p_{i+1} \dots p_j$ ,  $k$ , myParent, idx, plen )
(2) // Note: the initial call is: CreateTree(  $p$ ,  $k$ , nil, 0,  $\lfloor m/(k+1) \rfloor$  )
(3) Create new node node
(4) from(node) =  $i$ 
(5) to(node) =  $j$ 
(6) left =  $\lceil (k+1)/2 \rceil$ 
(7) parent(node) = myParent
(8) err(node) =  $k$ 
(9) if  $k = 0$ 
(10) then leafidx = node
(11) else
(12)    $lk = \lfloor (left \cdot k) / (k+1) \rfloor$ 
(13)   CreateTree(  $p_i \dots p_{i+left-plen-1}$ ,  $lk$ , node, idx, plen )
(14)    $rk = \lfloor ((k+1-left) \cdot k) / (k+1) \rfloor$ 
(15)   CreateTree(  $p_{i+left-plen} \dots p_j$ ,  $rk$ , node,  $idx + left$ , plen )
(16) fi

```

**Example:** Find the pattern  $P = \text{annual}$  in the text  $T = \text{annual.CPM.anniversary}$  with at most  $k = 2$  errors. First we build the tree with  $k + 1 = 3$  leaves. Below we write at each node  $n_i$  the variables (*from*, *to*, *error*).



**Search:** After constructing the tree, we have  $k + 1$  leafs *leaf<sub>i</sub>*. The  $k + 1$  subpatterns

$$\{ p_{from(n)}, \dots, p_{to(n)}, n = leaf_i, i \in \{0, \dots, k\} \}$$

are sent as input to a multi-pattern search algorithm (e. g. Aho-Corasick, Wu-Manbers, or SBOM). This algorithm gives as output a list of pairs (*pos*, *i*) where *pos* is the text position that matched and *i* is the number of the piece that matched.

The PEX algorithm performs verifications on its way upward in the tree, checking the presence of longer and longer pieces of the pattern, as specified by the nodes.

```

(1) Search phase of algorithm PEX
(2) for (pos, i) ∈ output of multi-pattern search do
(3)    $n = leaf_i$ ;  $in = from(n)$ ;  $n = parent(n)$ ;
(4)   cand = true;

```

```

(5)   while cand = true and n ≠ nil do
(6)       p1 = pos - (in - from(n)) - err(n);
(7)       p2 = pos + (to(n) - in) + err(n);
(8)       verify text tp1 ... tp2 for pattern piece pfrom(n) ... pto(n)
(9)       allowing err(n) errors;
(10)      if pattern piece was not found
(11)          then cand = false;
(12)      else n = parent(n);
(13)      fi
(14)  od
(15)  if cand = true
(16)      then report the positions where the whole p was found;
(17)  fi
(18) od

```

We search for *annual* in *annual.CPM.anniversary*. We constructed the tree for *annual*. A multi-pattern search algorithm finds: (1, 1), (12, 1), (3, 2), (5, 3). (Note that leaf *i* corresponds to pattern  $p^{i+1}$ ). For each of these positions we do the hierarchical verification:

```

Initialization for (1,1);
n=n0; in=1; n=n3; cand=true;
While loop;
  a) p1=1-(1-1)-1=0; p2=1+(4-1)+1=5;
  verify pattern annu in text annua with 1 error => found !
  b) p1=1-(1-1)-2=-1; p2=1+(6-1)+2=8;
  verify pattern annual in text annual_C => found !
  c) report end positions (6,7,8)

```

```

Initialization for (3,2);
n=n1; in=3; n=n3; cand=true;
While loop;
  a) p1=3-(3-1)-1=0; p2=3+(4-3)+1=5;
  verify pattern annu in text annua with 1 error => found !
  b) p1=3-(3-1)-2=-1; p2=3+(6-3)+2=8;
  verify pattern annual in text annual_C => found !
  c) report end positions (6,7,8)

```

```

Initialization for (12,1);
n=n0; in=1; n=n3; cand=true;
While loop;
  a) p1=12-(1-1)-1=11; p2=12+(4-1)+1=16;
  verify pattern annu in text _anniv with 1 error => found !
  b) p1=12-(1-1)-2=10; p2=12+(6-1)+2=19;
  verify pattern annual in text M_annivers => NOT found !

```

## 5.5 Summary

- Filtering algorithms prevent a large portion of the text from being looked at.
- The larger  $\alpha = k/m$ , the less efficient filtering algorithms become.
- Filtering algorithms based on the pigeonhole principle need an exact, multi-pattern search algorithm and a verification capable approximate string matching algorithm.
- The PEX algorithm starts verification from short exact matches and considers longer and longer substrings of the pattern as the verification proceeds upward in the tree.

## 6 QUASAR - q-gram based database searching

This exposition has been developed by Knut Reinert and Clemens Gröpl. It is based on the following sources, which are all recommended reading:

1. Burkhardt et al. (1999) *q-gram Based Database Searching Using a Suffix Array (QUASAR)*, Proc. RECOMB 99.
2. Burkhardt and Kärkkäinen (2001) *Better Filtering with Gapped q-grams*, Proc. CPM 01.

The tool **QUASAR** aims at aligning a query  $S = s_1, \dots, s_m$  in a text, also called database  $D = d_1, \dots, d_n$ . It can be seen as an efficient filter that uses exact matches. In contrast to online filtering algorithms, QUASAR uses a suffix array as indexing structure for the database.

### 6.1 Quasar

QUASAR, or “Q-gram Alignment based on Suffix ARrays”, is a filtering approach. QUASAR finds all local approximate matches of a query sequence  $S$  in a database  $D = \{d, \dots\}$ . The verification is performed by other means.

**Definition.** A sequence  $d$  is locally similar to  $S$ , if there exists at least one pair  $(S_{i,i+w-1}, d')$  of substrings such that:

1.  $S_{i,i+w-1}$  is a substring of length  $w$  and  $d'$  is a substring of  $D$ , and
2. the substrings  $d'$  and  $S_{i,i+w-1}$  have edit distance at most  $k$ .

We call this the *approximate matching problem with  $k$  differences and window length  $w$* .

For simplicity, we assume that the database consists of only one sequence, i. e.  $D = \{d\}$ .

### 6.2 The q-gram lemma

A short subsequence of length  $q$  is called a *q-gram*. In the following we start by considering the first  $w$  letters of  $S$ . The algorithm uses the following lemma:

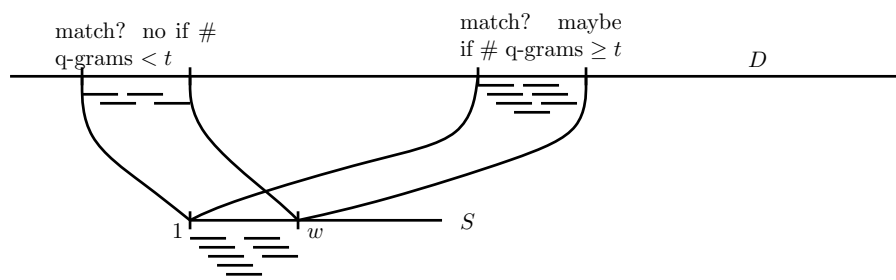
**Lemma 2.** Let  $P$  and  $S$  be strings of length  $w$  with at most  $k$  differences. Then  $P$  and  $S$  share at least  $w + 1 - (k + 1)q$  common  $q$ -grams.

In our case, this means:

**Lemma 3.** Let an occurrence of  $S_{1,w}$  with at most  $k$  differences end at position  $j$  in  $D$ . Then at least  $w + 1 - (k + 1)q$  of the  $q$ -grams in  $S_{1,w}$  occur in the substring  $D_{j-w+1,j}$ .

**Proof:** Exercise. . .

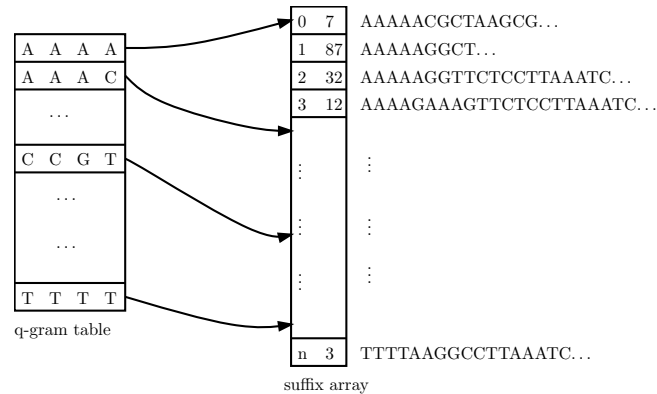
That means that as a necessary condition for an approximate match, at least  $t = w + 1 - (k + 1)q$  of the  $q$ -grams contained in  $S_{1,w}$  occur in a substring of  $D$  with length  $w$ . For example the strings ACAGCTTA and ACACCTTA have  $8 + 1 - (1 + 1)3 = 3$  common 3-grams, namely ACA, CTT and TTA.



## 6.3 q-gram index

The algorithm builds in a **first step** an **indexing structure** as follows:

1. Build a **suffix array**  $A$  over  $D$ .
2. Given  $q$ , **compute** for all possible  $|\Sigma|^q$   $q$ -grams the **start position** of the hitlist. This allows to **lookup** a  $q$ -gram in constant time.
3. If **another**  $q$  is specified,  $A$  is used to **recompute** the above **table**.

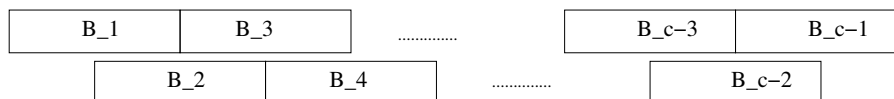


## 6.4 Counting q-grams

Now we have to **find** all **approximate matches** between  $S_{1,w}$  and  $D$ , that means we have to **find all substrings** in  $D$  that **share at least  $t$  q-grams** with  $S_{1,w}$ . The algorithm proceeds in the following basic steps on which we will elaborate:

1. Define **two arrays** of **non-overlapping blocks** of size  $b \geq 2w$ . The **first array** is **shifted by  $b/2$**  against the **other**.
2. **Process all q-grams** in  $S_{1,w}$  and **increment the counters** of the corresponding blocks.
3. All **blocks containing approximate matches** will have a **counter of at least  $t$** . (The reverse is not true).
4. **Shift** the **search window by one**. Now we **consider  $S_{2,w+1}$** .

## 6.5 Blocking



Since we want to count the  $q$ -grams that are in common between the query and the database, we use counters. Ideally we would use a counter of size  $w$  for each substring of this size. Since this uses too much memory, we build larger, non-overlapping blocks. While this decreases the memory usage, it also decreases the specificity.

Since the **blocks are not overlapping** we might **miss q-grams** that cross the block boundary. As a **remedy**, we use a **second, shifted array of blocks**.

## 6.6 Window Shifting

We started the search for approximate matches of window length  $w$  with the first  $w$ -mer in  $S$ , namely  $S_{1,w}$ . In order to determine the approximate matches for the next window  $S_{2,w+1}$ , we only have to discard the old  $q$ -gram  $S_{1,q}$  and consider the new  $q$ -gram  $S_{w-q+2,w+1}$ .

To do that we decrement the counters of all blocks that contain  $S_{1,q}$  that have not reached the threshold  $t$ . However, if the counter has already reached  $t$  it stays at this value to indicate a match for the extension phase.

For the new block we use the precomputed index and the suffix array to find the occurrences of the new  $q$ -gram and increment the corresponding block counters (at most two).

## 6.7 Alignment

After having computed the list of blocks, QUASAR uses BLAST to actually search the blocks. Here are some results from the initial implementation. QUASAR was run with  $w = 50$ ,  $q = 11$ , and  $t$  such that windows with at most 6% differences are found. Reasonable values for the block size are 512 to 4096.

DB size	query	id. res.	filtr. ratio	QUASAR	BLAST
73.5 Mb	368	91.4%	0.24%	0.123 s	3.27 s
280 Mb	393	97.1%	0.17%	0.38 s	13.27 s

“A database in BLAST format is built in main memory which is then passed to the BLAST search engine. The construction of this database requires a significant amount of time and introduces unnecessary overhead.”

## 6.8 Gapped $q$ -grams

In order to achieve a high filtration rate, we would like to choose  $q$  as large as possible, since the number of hits decreases exponentially in  $q$ . On the other hand, the threshold  $t = w - q - qk + 1$  also decreases with increasing  $q$  thereby reducing the filtering efficiency. The question is whether we could increase the length of the  $q$ -grams somehow, such that the threshold  $t$  stays high.

This can indeed be achieved by using *gapped  $q$ -grams*. For example the 3-grams with the *shape*  $\#\#.\#$  in the string ACAGCT are AC.G, CA.C, and AG.T:

ACAGCT
AC G
CA C
AG T

Next we define the concept formally.

**Definition 4.**

- A *shape*  $Q$  is a set of non-negative integers containing 0.
- The *size* of  $Q$ , denoted by  $|Q|$ , is the cardinality of the set.
- The *span* of  $Q$  is  $s(Q) = \max Q + 1$ .
- A shape of size  $q$  and span  $s$  is called  $(q, s)$ -*shape*.
- For any integer  $i$  and shape  $Q$ , the *positioned shape*  $Q_i$  is the set  $\{i + j \mid j \in Q\}$ .
- Let  $Q_i = \{i_1, i_2, \dots, i_q\}$ , where  $i = i_1 < i_2 < i_3 < \dots < i_q$ , and let  $S = s_1 s_2 \dots s_m$  be a string. For  $1 \leq i \leq m - s(Q) + 1$ , the  $Q$ -gram at position  $i$  in  $S$ , denoted by  $S[Q_i]$ , is the string  $s_{i_1} s_{i_2} \dots s_{i_q}$ .
- Two strings  $P$  and  $S$  have a common  $Q$ -gram at position  $i$  if  $P[Q_i] = S[Q_i]$ .

**Example 5.** Let  $Q = \{0, 1, 3, 6\}$  be a shape. Using the graphical representation it is the shape  $\#\#.\#.\#$ . Its size is  $|Q| = 4$  and its span is  $s(Q) = 7$ . The string ACGGATTAC has three  $Q$ -grams:  $S[Q_1] = s_1 s_2 s_4 s_7 = ACGT$ ,  $S[Q_2] = CGAA$ , and  $S[Q_3] = GGTC$ .

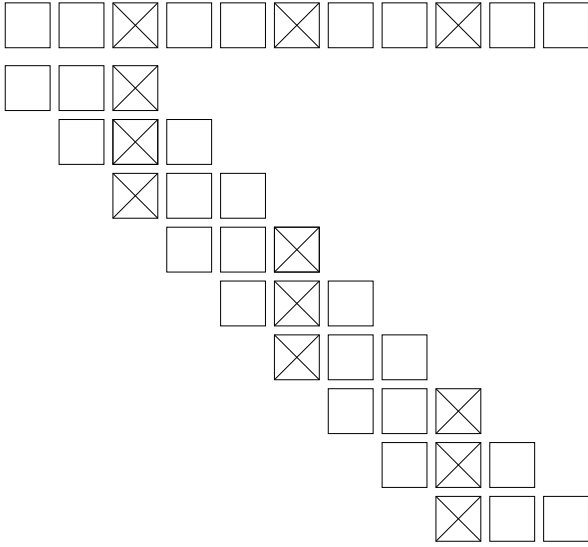
The  $q$ -gram lemma can be extended for gapped  $q$ -grams. A generalization gives

$$t = w - s(Q) - |Q|k + 1.$$

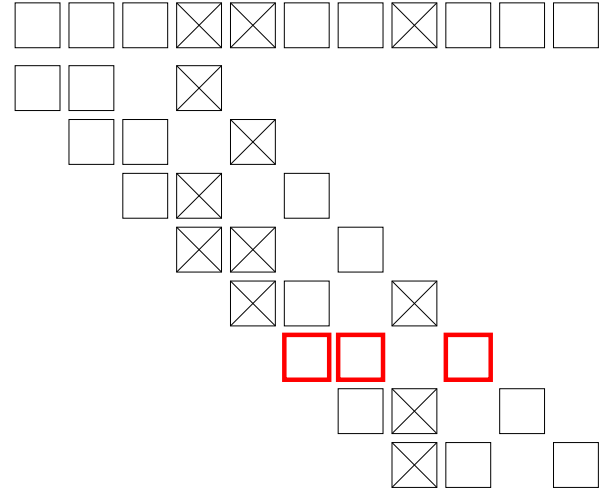
However it is not tight anymore (we will prove this).

**Example 6.** Let  $w = 11$  and  $k = 3$  and consider the 3-shapes  $###$  and  $##. \#$ . The above threshold for the two shapes is  $0 = 11 - 3 \cdot 4 + 1$  and  $-1 = 11 - 4 - 3 \cdot 3 + 1$  respectively. Thus neither shape would be useful for filtering. However, the real threshold for  $##. \#$  is 1. This can be checked by a full enumeration of all combinations of 3 mismatches.

shape:  $###$



shape:  $##. \#$



Worst-case mismatch positions

## 6.9 New threshold

What is the (tight) threshold for arbitrary  $Q$ -shapes?

Let  $P = p_1, \dots, p_w$  and  $S = s_1, \dots, s_w$  be two strings of length  $w$ . Let  $R(P, S)$  be the set of positions where  $S$  and  $P$  do not match. Then  $|R(P, S)|$  is the Hamming distance of  $P$  and  $S$ .

To determine the common  $Q$ -grams of  $P$  and  $S$  only the mismatch set is needed: It holds that

$$P[Q_i] = S[Q_i] \quad \text{if and only if} \quad Q_i \cap R(P, S) = \emptyset.$$

Using this notation we can define the threshold of a shape  $Q$  for a pattern of length  $w$  and Hamming distance  $k$  as:

$$t(Q, w, k) := \min_{R \subseteq \{1, \dots, w\}, |R|=k} \left| \{i \in \{1, \dots, w - s(Q) + 1\} \mid Q_i \cap R = \emptyset\} \right|$$

From the above discussion we get the following tight form of the  $q$ -gram lemma for arbitrary shapes:

**Lemma 7.** Let  $Q$  be a shape. For any two strings  $P$  and  $S$  of length  $w$  with Hamming distance  $k$ , the number of common  $Q$ -grams of  $P$  and  $S$  is at least  $t(Q, w, k)$ . Furthermore, there exist two strings  $P$  and  $S$  of length  $w$  and Hamming distance  $k$ , for which the number of common  $Q$ -grams is exactly  $t(Q, w, k)$ .

It is easy to see that this bound is as least as tight as the lower bound we already introduced:

**Lemma 8.**

$$t(Q, w, k) \geq \max\{0, w - s(Q) - |Q|k + 1\}$$

**Proof:** Let  $R$  be the set minimizing the expression in the definition of  $t(Q, w, k)$ . For each  $j \in R$  there are exactly  $|Q|$  integers  $i$  such that  $j \in Q_i$ . Therefore, at most  $k|Q|$  of the positioned shapes  $Q_i$ ,  $i \in \{1, \dots, w - s(Q) + 1\}$ , intersect with  $R$ , and at least  $w - s(Q) - k|Q| + 1$  do not intersect with  $R$ .

T The above lemma gives indeed the exact threshold for ungapped  $q$ -grams.



**Lemma 9.** Let  $Q$  be a contiguous shape, i. e.,  $Q = \{0, \dots, q-1\}$ . Then

$$t(Q, w, k) = \max\{0, w - s(Q) - |Q|k + 1\} = \max\{0, w - q(k+1) + 1\}.$$

**Proof:** The lower bound is shown by Lemma 8. For the upper bound we choose  $R = \{q, 2q, \dots, kq\}$ . Then  $Q_i$  intersects with  $R$  if and only if  $i \in \{1, \dots, kq\}$ , and thus does not intersect with  $R$  if  $i \in \{kq+1, \dots, w-q+1\}$ . Hence for this  $R$  we have only  $w - q + 1 - kq - 1 + 1 = w - (k+1)q + 1$  common  $q$ -grams.

The following table gives the exact thresholds for all shapes for  $w = 50$  and  $k = 5$ . One can see that in many cases, especially for higher values of  $q$ , best gapped shapes have higher thresholds than contiguous shapes of the same or even smaller size.

$s \downarrow : q \rightarrow$	4	5	6	7	8	9	10
5	26	21	–	–	–	–	–
6	25	20	15	–	–	–	–
7	24	19	14	9	–	–	–
8	23	18	13	8	3	–	–
9	22	<b>18</b> > 17	<b>14</b> > 12	<b>9</b> > 7	<b>5</b> > 2	0	–
10	21	<b>18</b> > 16	<b>13</b> > 11	<b>10</b> > 6	<b>6</b> > 1	<b>3</b> > 0	0
11	20	<b>16</b> > 15	<b>13</b> > 10	<b>10</b> > 5	<b>7</b> > 0	<b>4</b> > 0	<b>2</b> > 0
12	19	<b>16</b> > 14	<b>12</b> > 9	<b>9</b> > 4	<b>7</b> > 0	<b>4</b> > 0	<b>2</b> > 0

It has to be noted that it does not suffice to put in gaps somewhere; the gaps have to be chosen carefully. For example in the above table ( $w = 50$ ,  $k = 5$ , and  $q = 12$ ) there are only two shapes with a positive threshold, namely **###.#...###.#...###.#** and **#.#.#...#...#.#.#...#...#.#.#...#** and their mirror images.

## 6.10 Minimum coverage

The filtering efficiency of a  $Q$ -gram clearly depends on the threshold  $t(Q, w, k)$ . However there is also another factor that influences it. This factor is called *minimum coverage*.

Before we define it formally let's have a look at an example.

**Example 10.** Let  $w = 13$  and  $k = 3$ . Then both shapes **###** and **##.#** have a threshold of two. If two strings have four consecutive characters then they have two common 3-grams of shape **###**. In contrast, in order to have two common 3-grams of shape **##.#**, two strings need at least 5 matching characters.

This means, that the gapped 3-gram would have a lower count of common  $q$ -grams on strings that have only four consecutively matching characters although it has the same threshold.

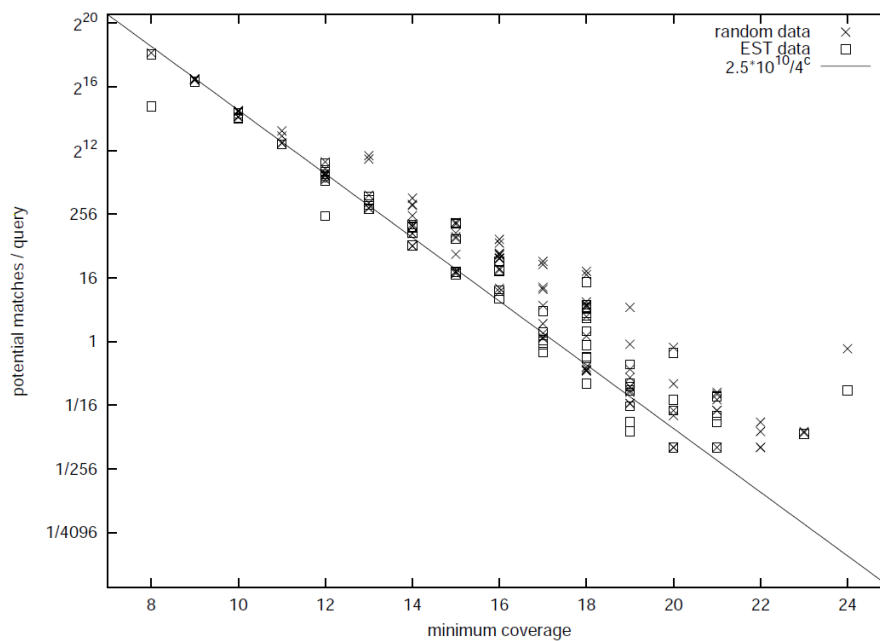
**Definition 11.** Let  $Q$  be a shape and  $t$  be a non-negative integer. The *minimum coverage* of  $Q$  for threshold  $t$  is:

$$c(Q, t) = \min_{C \subseteq \mathbb{N}, |C|=t} |\cup_{i \in C} Q_i|.$$

Hence the minimum coverage is the minimum number of characters that need to match between a pattern and a text substring for there to be  $t$  matching  $Q$ -grams.

Whenever possible, gapped Quasar chooses the highest minimum coverage, since it makes it more unlikely that a random string matches  $t$   $Q$ -grams. This improves the filter efficiency.

Computational experiments indicate that there is a strong correlation between the minimum coverage  $c(Q, t(Q, m, k))$  and the filter efficiency.



Correlation between expected and actual number of potential matches.

The following table shows different shapes for  $k = 5$ . The column *best* shows the shape with the highest minimum coverage (ties are broken using the threshold). The column *median* shows the median shape ordered by minimum coverage. If one chooses a random shape, the chance is 50% to be better (or worse) than this one. The last column show the best *one-gapped* shape. (The details of the tie breaking used here can be read in the paper.)

q	best	median	1-gapped
6	##.....#..#..#.#	#.###.....#.#	#####...#
7	#.###.....##..#.#	##..#..#...##	#####...##
8	##..#.#.....#.#..#.#	#.##.#####...#.#	#####...##
9	###..#..#.#...#..##	#####...#.#.#	#####...##
10	###..#..#.#...###.#	##.###..#.#.###.#	#####

## 6.11 Index structure

It is not necessary to use a suffix array for ungapped  $q$ -grams, and it is not possible anymore to use a suffix array for the gapped  $Q$ -grams. Instead, the database is scanned twice. The first time the number of occurrences of all  $Q$ -grams is counted.

In the second scan, the positions at which a  $q$ -gram starts are recorded in an array of size  $n$ . During that scan, the index points to the start of the respective list.

The detail shall be worked out as an *exercise*.

## 6.12 Extension to Levenshtein distance

Note that the  $q$ -gram method presented so far can only be used to find local approximate matches with the *Hamming* distance.

The  $q$ -gram method can be generalized to the *Levenshtein* distance. Burkhardt and Kärkkäinen have described an extension that uses 'one-gapped  $q$ -grams'.

The idea is to model insertions and deletions by additional  $Q$ -grams. For example, with the basic shape  $##-#$  applied the text, we would use  $##-#$ ,  $##--#$ , and  $###$  for the pattern.

The filter then compares all three shapes in the pattern to the  $q$ -grams of the basic shape in the text. Thus matching  $q$ -grams are even found in the presence of indels.

Otherwise the algorithm stays essentially unchanged, except that the threshold computation is slightly different.

## 6.13 Summary

- Filtering based on  $q$ -grams using a suffix array with an index is an efficient filtering method.
- In the gapless case, filtering efficiencies of  $\approx 0.2\%$  were observed for genomic sequences.
- Gapped  $Q$ -grams improve the filtering efficiency further (by orders of magnitude).
- The threshold  $t$  and the minimum coverage both influence the filter efficiency.
- No closed formula is known for computing  $t$  for gapped  $Q$ -grams.