# Indexing a large set of reads

Nicolas Philippe[1,2]    Mikaël Salson[3]    Thierry Lecroq[4]
Martine Léonard[4]    Thérèse Commes[2]    Éric Rivals[1]

[1]  LIRMM, CNRS and Université de Montpellier 2
[2]  IGH, CNRS, Montpellier
[3]  LIFL, CNRS and Université de Lille I – INRIA Lille-Nord Europe
[4]  LITIS, Université de Rouen

10 January 2011

# Introduction

**Context**

Next generation sequencers produce gigabytes of reads in a single run

# Introduction

**Context**

Next generation sequencers produce gigabytes of reads in a single run

**Problem**

How to search efficiently any relevant information?

# Introduction

**Context**

Next generation sequencers produce gigabytes of reads in a single run

**Problem**

How to search efficiently any relevant information?

Read ———————————

# Introduction

**Context**

Next generation sequencers produce gigabytes of reads in a single run

**Problem**

How to search efficiently any relevant information?

Factor $f$

Read ────────────

# Introduction

**Context**

Next generation sequencers produce gigabytes of reads in a single run

**Problem**

How to search efficiently any relevant information?

Factor $f$

Read ————

**Interesting questions**

- How many reads share this factor $f$?
- Which reads share this factor? At which positions?

# Introduction

**Context**

Next generation sequencers produce gigabytes of reads in a single run

**Problem**

How to search efficiently any relevant information?



Factor $f$

Read

**Interesting questions**

- How many reads share this factor $f$?
- Which reads share this factor? At which positions?

**Why is it interesting?**

- Genome assembly
- Read mapping
- ...

# Fixed-length factors

**Question**

Should we consider every factor?

# Fixed-length factors

**Question**

Should we consider every factor?

**Remarks**

- Factors of length 2 are quite uninformative
- At a certain point, increasing factor lengths does not help in identifying unique genome location ([Philippe *et al.*, 2009])

# Fixed-length factors

**Question**

Should we consider every factor?

**Remarks**

- Factors of length 2 are quite uninformative
- At a certain point, increasing factor lengths does not help in identifying unique genome location ([Philippe *et al.*, 2009])

**Conclusion**

We only consider $k$-length factors ($k$-factors or $k$-mers), $k$ being fixed

# Queries

**Queries for $k$-factors of a given read**

Given a read, and a $k$-factor in that read, we would like to know:

- Q1 the number of times this $k$-factor appears in the whole set of reads
- Q2 the reads and the positions in the reads in which it occurs
- Q3 the number of distinct reads in which it occurs

    . . .

# An immediate solution

> **Remark**
>
> We need to search patterns in a text

# An immediate solution

**Remark**

We need to search patterns in a text

**Classical solution**

Use a text index

# An immediate solution

**Remark**

We need to search patterns in a text

**Classical solution**

Use a text index

- Suffix tree
- Suffix array
- Compressed text index (FM-index, LZ-index, . . . )

# An immediate solution

**Remark**

We need to search patterns in a text

**Classical solution**

Use a text index

- Suffix tree
- Suffix array
- Compressed text index (FM-index, LZ-index, . . . )

Ok, let's try a suffix array!

# Using a Suffix Array for querying reads

Reads:
$r_0 = $ A T A A C G          $r_1 = $ A T A G T C          $r_2 = $ G A T A A C

# Using a Suffix Array for querying reads

Reads:
$r_0 = $ A T A A C G        $r_1 = $ A T A G T C        $r_2 = $ G A T A A C

(positions: 0 1 2 3 4 5    6 7 8 9 10 11    12 13 14 15 16 17)

$R = r_0 \cdot r_1 \cdot r_2 \cdot \$$

# Using a Suffix Array for querying reads

Reads:  $r_0 = $ A T A A C G $\qquad$ $r_1 = $ A T A G T C $\qquad$ $r_2 = $ G A T A A C

$R = r_0 \cdot r_1 \cdot r_2 \cdot \$ = $ A T A A C G A T A G T C G A T A A C $

$r_0$ $\qquad$ $r_1$ $\qquad$ $r_2$

# Using a Suffix Array for querying reads

Reads:
$r_0 = $ A T A A C G     $r_1 = $ A T A G T C     $r_2 = $ G A T A A C

$R = r_0 \cdot r_1 \cdot r_2 \cdot \$ = $ A T A A C G A T A G T C G A T A A C $

                 $r_0$            $r_1$            $r_2$

Let's build the suffix array (sort suffixes in lexicographic ascending order)

| SA | Suffixes |
|----|----------|
| 18 | $ |
| 15 | AAC$ |
| 2 | AACGATAGTCGATAAC$ |
| 16 | AC$ |
| 3 | ACGATAGTCGATAAC$ |
| 8 | AGTCGATAAC$ |
| 13 | ATAAC$ |
| 0 | ATAACGATAGTCGATAAC$ |
| 17 | C$ |
| 11 | CGATAAC$ |
| ⋮ | ⋮ |

# Using a Suffix Array for querying reads

Reads:      $r_0 = $ A T A A C G      $r_1 = $ A T A G T C      $r_2 = $ G A T A A C

$R = r_0 \cdot r_1 \cdot r_2 \cdot \$ = $ A T A A C G A T A G T C G A T A A C \$

$r_0$          $r_1$          $r_2$

Let's build the suffix array (sort suffixes in lexicographic ascending order)

| SA | Suffixes |
|---|---|
| 18 | \$ |
| 15 | AAC\$ |
| 2 | AACGATAGTCGATAAC\$ |
| 16 | AC\$ |
| 3 | ACGATAGTCGATAAC\$ |
| 8 | AGTCGATAAC\$ |
| 13 | ATAAC\$ |
| 0 | ATAACGATAGTCGATAAC\$ |
| 17 | C\$ |
| 11 | CGATAAC\$ |
| ⋮ | ⋮ |

**Remark**

Only the $k$ first letters of each suffix are interesting ($k = 3$)

# Using a Suffix Array for querying reads

Reads:

$r_0 = $ A T A A C G     $r_1 = $ A T A G T C     $r_2 = $ G A T A A C

$R = r_0 \cdot r_1 \cdot r_2 \cdot \$ = $ A T A A C G A T A G T C G A T A A C \$

$r_0$      $r_1$      $r_2$

Let's build the suffix array (sort suffixes in lexicographic ascending order)

| SA | Suffixes |
|---|---|
| 18 | \$ |
| 15 | AAC\$ |
| 2 | AACGATAGTCGATAAC\$ |
| 16 | AC\$ |
| 3 | ACGATAGTCGATAAC\$ |
| 8 | AGTCGATAAC\$ |
| 13 | ATAAC\$ |
| 0 | ATAACGATAGTCGATAAC\$ |
| 17 | C\$ |
| 11 | CGATAAC\$ |
| ⋮ | ⋮ |

**Remark**

Only the $k$ first letters of each suffix are interesting ($k = 3$)

5

# Using a Suffix Array for querying reads

Reads:

$r_0 = $ A T A A C G     $r_1 = $ A T A G T C     $r_2 = $ G A T A A C

$R = r_0 \cdot r_1 \cdot r_2 \cdot \$ = $ A T A A C G A T A G T C G A T A A C $

$r_0$          $r_1$          $r_2$

Let's build the suffix array (sort suffixes in lexicographic ascending order)

| SA | Suffixes |
|----|----------|
| 18 | $ |
| 15 | AAC$ |
| 2 | AACGATAGTCGATAAC$ |
| 16 | AC$ |
| 3 | ACGATAGTCGATAAC$ |
| 8 | AGTCGATAAC$ |
| 13 | ATAAC$ |
| 0 | ATAACGATAGTCGATAAC$ |
| 17 | C$ |
| 11 | CGATAAC$ |
| ⋮ | ⋮ |

**Remark**

Only the $k$ first letters of each suffix are interesting ($k = 3$)

**Remark**

Factors overlapping two reads are undesirable

# Using a Suffix Array for querying reads

Reads:
$$r_0 = \overset{0\ 1\ 2\ 3\ 4\ 5}{\text{A T A A C G}} \qquad r_1 = \overset{6\ 7\ 8\ 9\ 10\ 11}{\text{A T A G T C}} \qquad r_2 = \overset{12\ 13\ 14\ 15\ 16\ 17}{\text{G A T A A C}}$$

$$R = r_0 \cdot r_1 \cdot r_2 \cdot \$ = \overset{0\ \ 1\ \ 2\ \ 3\ \ 4\ \ 5\ \ 6\ \ 7\ \ 8\ \ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17\ 18}{\text{A T A A C G A T A G T C G A T A A C \$}}$$

$$\underbrace{\phantom{\text{ATAACG}}}_{r_0} \qquad \underbrace{\phantom{\text{ATAGTC}}}_{r_1} \qquad \underbrace{\phantom{\text{GATAAC}}}_{r_2}$$

Let's build the suffix array (sort suffixes in lexicographic ascending order)

| SA | Suffixes |
|----|----------|
| 18 | $\$$ |
| 15 | AAC$ |
| 2 | AAC GATAGTCGATAAC$ |
| 16 | AC$ |
| 3 | ACG ATAGTCGATAAC$ |
| 8 | AGT CGATAAC$ |
| 13 | ATA AC$ |
| 0 | ATA ACGATAGTCGATAAC$ |
| 17 | C$ |
| 11 | CGA TAAC$ |
| ⋮ | ⋮ |

**Remark**

Only the $k$ first letters of each suffix are interesting ($k = 3$)

**Remark**

Factors overlapping two reads are undesirable

# Discarding useless positions

# Discarding useless positions

# Discarding useless positions

# Discarding useless positions

# Discarding useless positions

$P$-positions

0  1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16 17 18

$R =$ A T A A C G A T A G T C G A T A A C $

---

**$P$-positions**

- Set of positions where a $k$-factor belonging to a single read starts.
- This set is not a permutation

A $k$-factor starting at a $P$-position is called a $P$-$k$-factor

---

# Discarding useless positions



$P$-positions

$R = \text{ATAACGATAGTCGATAAC}\$$

$Q$-positions

### $P$-positions

- Set of positions where a $k$-factor belonging to a single read starts.
- This set is not a permutation

A $k$-factor starting at a $P$-position is called a $P$-$k$-factor

### $Q$-positions

Renumbered $P$-positions so that the set of $Q$-positions is a permutation

# Generalized $k$-Factor Array (GkFA)



$P$-positions

$$R = \text{A T A A C G A T A G T C G A T A A C } \$$$

$Q$-positions

**Generalized $k$-factor array**

Index suffixes starting at $P$-positions. Positions are renumbered to $Q$-positions.

# Generalized $k$-Factor Array (GkFA)

| SA | Suffixes |
|---|---|
| 18 | $ |
| 15 | AAC$ |
| 2 | AACGATAGTCGATAAC$ |
| 16 | AC$ |
| 3 | ACGATAGTCGATAAC$ |
| 8 | AGTCGATAAC$ |
| 13 | ATAAC$ |
| 0 | ATAACGATAGTCGATAAC$ |
| 6 | ATAGTCGATAAC$ |
| 17 | C$ |
| 11 | CGATAAC$ |
| 4 | CGATAGTCGATAAC$ |
| 12 | GATAAC$ |
| 5 | GATAGTCGATAAC$ |
| 9 | GTCGATAAC$ |
| 14 | TAAC$ |
| 1 | TAACGATAGTCGATAAC$ |
| 7 | TAGTCGATAAC$ |
| 10 | TCGATAAC$ |

$P$-positions

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

$R = $ A T A A C G A T A G T C G A T A A C $

0 1 2 3 4 5 6 7 8 9 10 11

$Q$-positions

**Generalized $k$-factor array**

Index suffixes starting at $P$-positions. Positions are renumbered to $Q$-positions.

7

# Generalized $k$-Factor Array (GkFA)



| SA | Suffixes |
|----|----------|
| ~~18~~ | ~~$~~ |
| 15 | AAC$ |
| 2 | AACGATAGTCGATAAC$ |
| ~~16~~ | ~~AC$~~ |
| 3 | ACGATAGTCGATAAC$ |
| 8 | AGTCGATAAC$ |
| 13 | ATAAC$ |
| 0 | ATAACGATAGTCGATAAC$ |
| 6 | ATAGTCGATAAC$ |
| ~~17~~ | ~~C$~~ |
| ~~11~~ | ~~CGATAAC$~~ |
| ~~4~~ | ~~CGATAGTCGATAAC$~~ |
| 12 | GATAAC$ |
| ~~5~~ | ~~GATAGTCGATAAC$~~ |
| 9 | GTCGATAAC$ |
| 14 | TAAC$ |
| 1 | TAACGATAGTCGATAAC$ |
| 7 | TAGTCGATAAC$ |
| ~~10~~ | ~~TCGATAAC$~~ |

**Generalized $k$-factor array**

Index suffixes starting at $P$-positions. Positions are renumbered to $Q$-positions.

# Generalized $k$-Factor Array (GkFA)



| SA | Suffixes | GkFA ($P$-positions) |
|---|---|---|
| ~~18~~ | ~~$~~ | |
| 15 | AAC$ | 15 |
| 2 | AACGATAGTCGATAAC$ | 2 |
| ~~16~~ | ~~AC$~~ | |
| 3 | ACGATAGTCGATAAC$ | 3 |
| 8 | AGTCGATAAC$ | 8 |
| 13 | ATAAC$ | 13 |
| 0 | ATAACGATAGTCGATAAC$ | 0 |
| 6 | ATAGTCGATAAC$ | 6 |
| ~~17~~ | ~~C$~~ | |
| ~~11~~ | ~~CGATAAC$~~ | |
| ~~4~~ | ~~CGATAGTCGATAAC$~~ | |
| 12 | GATAAC$ | 12 |
| ~~5~~ | ~~GATAGTCGATAAC$~~ | |
| 9 | GTCGATAAC$ | 9 |
| 14 | TAAC$ | 14 |
| 1 | TAACGATAGTCGATAAC$ | 1 |
| 7 | TAGTCGATAAC$ | 7 |
| ~~10~~ | ~~TCGATAAC$~~ | |

**Generalized $k$-factor array**

Index suffixes starting at $P$-positions. Positions are renumbered to $Q$-positions.

# Generalized $k$-Factor Array (GkFA)



| SA | Suffixes | GkFA ($P$-positions) | GkFA ($Q$-positions) |
|---|---|---|---|
| ~~18~~ | ~~$~~ | | |
| 15 | AAC$ | 15 | 11 |
| 2 | AACGATAGTCGATAAC$ | 2 | 2 |
| ~~16~~ | ~~AC$~~ | | |
| 3 | ACGATAGTCGATAAC$ | 3 | 3 |
| 8 | AGTCGATAAC$ | 8 | 6 |
| 13 | ATAAC$ | 13 | 9 |
| 0 | ATAACGATAGTCGATAAC$ | 0 | 0 |
| 6 | ATAGTCGATAAC$ | 6 | 4 |
| ~~17~~ | ~~C$~~ | | |
| ~~11~~ | ~~CGATAAC$~~ | | |
| ~~4~~ | ~~CGATAGTCGATAAC$~~ | | |
| 12 | GATAAC$ | 12 | 8 |
| ~~5~~ | ~~GATAGTCGATAAC$~~ | | |
| 9 | GTCGATAAC$ | 9 | 7 |
| 14 | TAAC$ | 14 | 10 |
| 1 | TAACGATAGTCGATAAC$ | 1 | 1 |
| 7 | TAGTCGATAAC$ | 7 | 5 |
| ~~10~~ | ~~TCGATAAC$~~ | | |

**Generalized $k$-factor array**

Index suffixes starting at $P$-positions. Positions are renumbered to $Q$-positions.

# Generalized $k$-Factor Array (GkFA)



| SA | Suffixes | GkFA ($Q$-positions) |
|---|---|---|
| ~~18~~ | ~~$~~ | |
| 15 | AAC$ | 11 |
| 2 | AACGATAGTCGATAAC$ | 2 |
| ~~16~~ | ~~AC$~~ | |
| 3 | ACGATAGTCGATAAC$ | 3 |
| 8 | AGTCGATAAC$ | 6 |
| 13 | ATAAC$ | 9 |
| 0 | ATAACGATAGTCGATAAC$ | 0 |
| 6 | ATAGTCGATAAC$ | 4 |
| ~~17~~ | ~~C$~~ | |
| ~~11~~ | ~~CGATAAC$~~ | |
| ~~4~~ | ~~CGATAGTCGATAAC$~~ | |
| 12 | GATAAC$ | 8 |
| ~~5~~ | ~~GATAGTCGATAAC$~~ | |
| 9 | GTCGATAAC$ | 7 |
| 14 | TAAC$ | 10 |
| 1 | TAACGATAGTCGATAAC$ | 1 |
| 7 | TAGTCGATAAC$ | 5 |
| ~~10~~ | ~~TCGATAAC$~~ | |

**Generalized $k$-factor array**

Index suffixes starting at $P$-positions. Positions are renumbered to $Q$-positions.

7

# Generalized $k$ Count Factor Array

$$R = \underset{\underset{Q\text{-positions}}{0\ \ 1\ \ 2\ \ 3\qquad 4\ \ 5\ \ 6\ \ 7\qquad 8\ \ 9\ 1011}}{\text{A T A A C G A T A G T C G A T A A C}}\,\$$$

# Generalized $k$ Count Factor Array

$$R = \underset{\underset{0\ \ \ 1\ \ \ 2\ \ \ 3\quad\quad\ 4\ \ \ 5\ \ \ 6\ \ \ 7\quad\quad\ 8\ \ \ 9\ \ 10\,11}{Q\text{-positions}}}{\text{A T A A C G A T A G T C G A T A A C}}\ \$}$$

**GkCFA**

Count the number of occurrences of a
$P$-$k$-factor

# Generalized $k$ Count Factor Array

$$R = \underset{0\ \ 1\ \ 2\ \ 3}{\text{A T A A}} \text{C G} \underset{4\ \ 5\ \ 6\ \ 7}{\text{A T A G}} \text{T C} \underset{8\ \ 9\ \ 10\ 11}{\text{G A T A}} \text{A C \$}$$

$Q$-positions

**GkCFA**

Count the number of occurrences of a
$P$-$k$-factor

**Purpose**

Compute the read coverage of a given
region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$
$$\underset{0 \ 1 \ 2 \ 3 \qquad 4 \ 5 \ 6 \ 7 \qquad 8 \ 9 \ 10 \ 11}{}$$

$Q$-positions

| GkFA | $k$-factor |
|---|---|
| 11 | AAC |
| 2 | AAC |
| 3 | ACG |
| 6 | AGT |
| 9 | ATA |
| 0 | ATA |
| 4 | ATA |
| 8 | GAT |
| 7 | GTC |
| 10 | TAA |
| 1 | TAA |
| 5 | TAG |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0 1 2 3    4 5 6 7    8 9 10 11

$Q$-positions

| GkFA | | $k$-factor |
|---|---|---|
| 0 | 11 | AAC |
| | 2 | AAC |
| | 3 | ACG |
| | 6 | AGT |
| | 9 | ATA |
| | 0 | ATA |
| | 4 | ATA |
| | 8 | GAT |
| | 7 | GTC |
| | 10 | TAA |
| | 1 | TAA |
| | 5 | TAG |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0\ 1\ 2\ 3}}{\text{A T A A}}\ \text{C}\ \underset{\substack{4\ 5\ 6\ 7}}{\text{G A T A}}\ \text{G T C}\ \underset{\substack{8\ 9\ 10\ 11}}{\text{G A T A}}\ \text{A C \$}$$

$Q$-positions

| GkFA | $k$-factor | |
| --- | --- | --- |
| 11 | AAC | |
| 2 | AAC | |
| 3 | ACG | |
| 6 | AGT | |
| 9 | ATA | |
| 0 | ATA | |
| 4 | ATA | |
| 8 | GAT | |
| 7 | GTC | |
| 10 | TAA | |
| 1 | TAA | |
| 5 | TAG | |

0    2

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0 \ 1 \ 2 \ 3}}{ATAA}\underset{\substack{4 \ 5 \ 6 \ 7}}{CGATAGTC}\underset{\substack{8 \ 9 \ 1011}}{GATAAC}\$$$

$Q$-positions

| GkFA | | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | |
| | 6 | AGT | |
| | 9 | ATA | |
| | 0 | ATA | |
| | 4 | ATA | |
| | 8 | GAT | |
| | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0  1  2  3       4  5  6  7       8  9  10  11

$Q$-positions

| GkFA |  | $k$-factor |  |
|------|------|------------|------|
| 0 | 11 | AAC | 2 |
|  | 2 | AAC |  |
| 1 | 3 | ACG | 1 |
|  | 6 | AGT |  |
|  | 9 | ATA |  |
|  | 0 | ATA |  |
|  | 4 | ATA |  |
|  | 8 | GAT |  |
|  | 7 | GTC |  |
|  | 10 | TAA |  |
|  | 1 | TAA |  |
|  | 5 | TAG |  |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0 1 2 3   4 5 6 7   8 9 10 11

$Q$-positions

| GkFA | $k$-factor | |
|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | |
| | 9 | ATA | |
| | 0 | ATA | |
| | 4 | ATA | |
| | 8 | GAT | |
| | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{ATAACGATAGTCGATAAC}\$$$

0 1 2 3    4 5 6 7    8 9 10 11

$Q$-positions

| GkFA | | $k$-factor | |
|------|------|-----------|------|
| 0 | 11 | AAC | 2 |
|   | 2  | AAC |   |
| 1 | 3  | ACG | 1 |
| 2 | 6  | AGT | 1 |
|   | 9  | ATA |   |
|   | 0  | ATA |   |
|   | 4  | ATA |   |
|   | 8  | GAT |   |
|   | 7  | GTC |   |
|   | 10 | TAA |   |
|   | 1  | TAA |   |
|   | 5  | TAG |   |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0  1  2  3        4  5  6  7        8  9  10 11

$Q$-positions

| GkFA | | $k$-factor | |
|---|---|---|---|
| | 11 | AAC | |
| 0 | 2 | AAC | 2 |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| | 9 | ATA | |
| 3 | 0 | ATA | |
| | 4 | ATA | |
| | 8 | GAT | |
| | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0  1  2  3          4  5  6  7          8  9  10 11

$Q$-positions

| GkFA | | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| | 8 | GAT | |
| | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$
$$\quad\; 0\; 1\; 2\; 3 \qquad\quad 4\; 5\; 6\; 7 \qquad\quad 8\; 9\; 10\; 11$$
$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | |
| | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0\ 1\ 2\ 3}}{A\,T\,A\,A}\,C\,\underset{\substack{4\ 5\ 6\ 7}}{G\,A\,T\,A}\,G\,T\,C\,\underset{\substack{8\ 9\ 10\ 11}}{G\,A\,T\,A}\,A\,C\,\$$$

$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

8

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0\ 1\ 2\ 3}}{\text{A T A A}} \text{C G} \underset{\substack{4\ 5\ 6\ 7}}{\text{A T A G}} \text{T C} \underset{\substack{8\ 9\ 10\ 11}}{\text{G A T A}} \text{A C \$}$$

$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$R = \text{A T A A C G A T A G T C G A T A A C \$}$

0 1 2 3     4 5 6 7     8 9 10 11

$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
| | 10 | TAA | |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0 \ 1 \ 2 \ 3 \qquad 4 \ 5 \ 6 \ 7 \qquad 8 \ 9 \ 10 \ 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}$$

$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
| 6 | 10 | TAA | 2 |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0\ 1\ 2\ 3}}{\text{A T A A}}\ \text{C G}\ \underset{\substack{4\ 5\ 6\ 7}}{\text{A T A G}}\ \text{T C}\ \underset{\substack{8\ 9\ 10\ 11}}{\text{G A T A}}\ \text{A C \$}$$

$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
| 6 | 10 | TAA | 2 |
| | 1 | TAA | |
| | 5 | TAG | |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0\ 1\ 2\ 3\quad\ 4\ 5\ 6\ 7\quad\ 8\ 9\ 10\ 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}$$

$Q$-positions

|   | GkFA |  $k$-factor |   |
|---|------|-------------|---|
| 0 | 11 | AAC | 2 |
|   | 2 | AAC |   |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
|   | 9 | ATA |   |
| 3 | 0 | ATA | 3 |
|   | 4 | ATA |   |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
|   | 10 | TAA | |
| 6 | 1 | TAA | 2 |
| 7 | 5 | TAG | 1 |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

8

# Generalized $k$ Count Factor Array

$$R = \underset{\substack{0 \ 1 \ 2 \ 3 \quad\quad 4 \ 5 \ 6 \ 7 \quad\quad 8 \ 9 \ 10 11}}{\text{A T A A C G A T A G T C G A T A A C}} \$$$

$Q$-positions

| | GkFA | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
| 6 | 10 | TAA | 2 |
| | 1 | TAA | |
| 7 | 5 | TAG | 1 |

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{ATAACGATAGTCGATAAC\$}$$

positions: 0 1 2 3   4 5 6 7   8 9 10 11

$Q$-positions

| GkFA | | $k$-factor | |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
| | 2 | AAC | |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
| 3 | 9 | ATA | 3 |
| | 0 | ATA | |
| | 4 | ATA | |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
| 6 | 10 | TAA | 2 |
| | 1 | TAA | |
| 7 | 5 | TAG | 1 |

GkCFA

**GkCFA**
Count the number of occurrences of a $P$-$k$-factor

**Purpose**
Compute the read coverage of a given region *inside* a read

# Generalized $k$ Count Factor Array

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

positions: 0 1 2 3    4 5 6 7    8 9 10 11

$Q$-positions

|  | GkFA | $k$-factor | GkCFA |
|---|---|---|---|
| 0 | 11 | AAC | 2 |
|  | 2 | AAC |  |
| 1 | 3 | ACG | 1 |
| 2 | 6 | AGT | 1 |
|  | 9 | ATA |  |
| 3 | 0 | ATA | 3 |
|  | 4 | ATA |  |
| 4 | 8 | GAT | 1 |
| 5 | 7 | GTC | 1 |
| 6 | 10 | TAA | 2 |
|  | 1 | TAA |  |
| 7 | 5 | TAG | 1 |

IDs            GkCFA

**GkCFA**

Count the number of occurrences of a $P$-$k$-factor

**Purpose**

Compute the read coverage of a given region *inside* a read

8

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0  1  2  3      4  5  6  7      8  9 10 11

$Q$-positions

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0 1 2 3    4 5 6 7    8 9 10 11

$Q$-positions

## GkIFA

- ▶ "Inverse" of GkFA
- ▶ Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

9

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0  1  2  3    4  5  6  7    8  9 10 11

$Q$-positions

**GkIFA**

- ▶ "Inverse" of GkFA
- ▶ Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

**Purpose**

Quickly find the id associated to a $k$-factor coming from a read

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

Q-positions: 0 1 2 3   4 5 6 7   8 9 10 11

**GkIFA**

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

**Purpose**

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | |
| | 2 | AAC | 1 | |
| 1 | 3 | ACG | 2 | |
| 2 | 6 | AGT | 3 | |
| 3 | 9 | ATA | 4 | |
| | 0 | ATA | 5 | |
| | 4 | ATA | 6 | |
| 4 | 8 | GAT | 7 | |
| 5 | 7 | GTC | 8 | |
| 6 | 10 | TAA | 9 | |
| | 1 | TAA | 10 | |
| 7 | 5 | TAG | 11 | |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \underset{\substack{0 \ 1 \ 2 \ 3}}{\text{A T A A}} \text{C G} \underset{\substack{4 \ 5 \ 6 \ 7}}{\text{A T A G}} \text{T C} \underset{\substack{8 \ 9 \ 10 \ 11}}{\text{G A T A}} \text{A C \$}$$

$Q$-positions

**GkIFA**

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

**Purpose**

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | |
| | 2 | AAC | 1 | |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | |
| | 9 | ATA | 4 | |
| 3 | 0 | ATA | 5 | |
| | 4 | ATA | 6 | |
| 4 | 8 | GAT | 7 | |
| 5 | 7 | GTC | 8 | |
| | 10 | TAA | 9 | |
| 6 | 1 | TAA | 10 | |
| 7 | 5 | TAG | 11 | 0 |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C } \$$$

$Q$-positions at 0 1 2 3, 4 5 6 7, 8 9 10 11

**GkIFA**

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

**Purpose**

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | $k$-factor | $i$ | GkIFA |
|---|---|---|---|
| 0 | 11 AAC | 0 | |
| | 2 AAC | 1 | |
| 1 | 3 ACG | 2 | 0 |
| 2 | 6 AGT | 3 | 1 |
| | 9 ATA | 4 | |
| 3 | 0 ATA | 5 | |
| | 4 ATA | 6 | |
| 4 | 8 GAT | 7 | |
| 5 | 7 GTC | 8 | |
| 6 | 10 TAA | 9 | |
| | 1 TAA | 10 | |
| 7 | 5 TAG | 11 | 0 |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$
0 1 2 3   4 5 6 7   8 9 10 11

$Q$-positions

### GkIFA

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

### Purpose

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | |
| | 2 | AAC | 1 | |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
| | 9 | ATA | 4 | |
| 3 | 0 | ATA | 5 | |
| | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | |
| 5 | 7 | GTC | 8 | |
| | 10 | TAA | 9 | |
| 6 | 1 | TAA | 10 | |
| 7 | 5 | TAG | 11 | 0 |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \underset{\underset{Q\text{-positions}}{0\ 1\ 2\ 3\quad\quad 4\ 5\ 6\ 7\quad\ 8\ 9\ 10\ 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}$$

**GkIFA**

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

**Purpose**

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | 3 |
| | 2 | AAC | 1 | |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
| | 9 | ATA | 4 | 3 |
| 3 | 0 | ATA | 5 | |
| | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | |
| 5 | 7 | GTC | 8 | |
| | 10 | TAA | 9 | 3 |
| 6 | 1 | TAA | 10 | |
| 7 | 5 | TAG | 11 | 0 |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$
0 1 2 3    4 5 6 7    8 9 10 11
$Q$-positions

## GkIFA

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

## Purpose

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | 3 |
| | 2 | AAC | 1 | |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
| | 9 | ATA | 4 | 3 |
| 3 | 0 | ATA | 5 | |
| | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | |
| 5 | 7 | GTC | 8 | 4 |
| | 10 | TAA | 9 | 3 |
| 6 | 1 | TAA | 10 | |
| 7 | 5 | TAG | 11 | 0 |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

$Q$-positions

| GkFA | $k$-factor | $i$ | GkIFA |
|------|-----------|-----|-------|
| 0 | 11 | AAC | 0 | 3 |
|   | 2 | AAC | 1 | |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
| 3 | 9 | ATA | 4 | 3 |
|   | 0 | ATA | 5 | |
|   | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | 5 |
| 5 | 7 | GTC | 8 | 4 |
| 6 | 10 | TAA | 9 | 3 |
|   | 1 | TAA | 10 | |
| 7 | 5 | TAG | 11 | 0 |

## GkIFA

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

## Purpose

Quickly find the id associated to a $k$-factor coming from a read

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$
$$\phantom{R = \text{A T}} 0\ 1\ 2\ 3 \quad\ \ 4\ 5\ 6\ 7 \quad\ \ 8\ 9\ 10\ 11$$
$$Q\text{-positions}$$

## GkIFA

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

## Purpose

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | 3 |
|   | 2 | AAC | 1 | 6 |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
| 3 | 9 | ATA | 4 | 3 |
|   | 0 | ATA | 5 | |
|   | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | 5 |
| 5 | 7 | GTC | 8 | 4 |
| 6 | 10 | TAA | 9 | 3 |
|   | 1 | TAA | 10 | 6 |
| 7 | 5 | TAG | 11 | 0 |

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$
$$\phantom{R = }\text{0 1 2 3 \quad 4 5 6 7 \quad 8 9 10 11}$$

$Q$-positions

## GkIFA

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

## Purpose

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | 3 |
|   | 2 | AAC | 1 | 6 |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
|   | 9 | ATA | 4 | 3 |
| 3 | 0 | ATA | 5 | 7 |
|   | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | 5 |
| 5 | 7 | GTC | 8 | 4 |
|   | 10 | TAA | 9 | 3 |
| 6 | 1 | TAA | 10 | 6 |
| 7 | 5 | TAG | 11 | 0 |

9

# Generalized $k$ Inverse Factor Array (GkIFA)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0 1 2 3   4 5 6 7   8 9 10 11

$Q$-positions

**GkIFA**

- "Inverse" of GkFA
- Given a $Q$-position in $R$, stores the id associated to the corresponding $k$-factor

**Purpose**

Quickly find the id associated to a $k$-factor coming from a read

| GkFA | | $k$-factor | $i$ | GkIFA |
|---|---|---|---|---|
| 0 | 11 | AAC | 0 | 3 |
|   | 2 | AAC | 1 | 6 |
| 1 | 3 | ACG | 2 | 0 |
| 2 | 6 | AGT | 3 | 1 |
|   | 9 | ATA | 4 | 3 |
| 3 | 0 | ATA | 5 | 7 |
|   | 4 | ATA | 6 | 2 |
| 4 | 8 | GAT | 7 | 5 |
| 5 | 7 | GTC | 8 | 4 |
|   | 10 | TAA | 9 | 3 |
| 6 | 1 | TAA | 10 | 6 |
| 7 | 5 | TAG | 11 | 0 |

# Using Gk arrays (I)

$$R = \begin{array}{cccccccccccccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{array}$$

0  1  2  3        4  5  6  7        8  9  10 11

$Q$-positions

# Using Gk arrays (I)

$$R = \text{ATAACGATAGTCGATAAC\$}$$

positions above: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

$Q$-positions below: 0 1 2 3 4 5 6 7 8 9 10 11

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

# Using Gk arrays (I)

$$R = \overset{\text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18}}{\underset{\text{0 1 2 3 \quad\quad 4 5 6 7 \quad\quad 8 9 10 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

| $r_2$ at position 0 |
|:---:|

# Using Gk arrays (I)

$$R = \overset{\substack{0\ 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17\ 18}}{\text{A T A A C G A T A G T C G A T A A C \$}}$$

$Q$-positions

(positions: 0 1 2 3 — 4 5 6 7 — 8 9 10 11)

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

| $r_2$ at position 0 |
|---|

Read length: 6 $\longrightarrow$

# Using Gk arrays (I)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

positions (top): 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

$Q$-positions (bottom): 0 1 2 3    4 5 6 7    8 9 10 11

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

| $r_2$ at position 0 | Read length: 6 $\rightarrow$ | $P$-position $2 \times 6 + 0 = 12$ in $R$ |
|---|---|---|

# Using Gk arrays (I)

$$R = \begin{array}{c} {\scriptstyle 0\ 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17\ 18} \\ \texttt{A T A A C G A T A G T C G A T A A C \$} \\ {\scriptstyle 0\ 1\ 2\ 3\quad\ \ 4\ 5\ 6\ 7\quad\ 8\ 9\ 10\ 11} \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

| $r_2$ at position 0 | Read length: 6 | $P$-position $2 \times 6 + 0 = 12$ in $R$ | $k = 3$ |

# Using Gk arrays (I)

$$R = \begin{array}{c} \text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18} \\ \text{A T A A C G A T A G T C G A T A A C \$} \\ \text{0 1 2 3 \quad 4 5 6 7 \quad 8 9 10 11} \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

Position 0 in read $r_2$ corresponds to $Q$-position 8

| $r_2$ at position 0 | Read length: 6 $\rightarrow$ | $P$-position $2 \times 6 + 0 = 12$ in $R$ | $k = 3$ $\rightarrow$ | $Q$-position $12 - (3-1) * 2 = 8$ in $R$ |
|---|---|---|---|---|

# Using Gk arrays (I)

$$R = \begin{array}{cccccccccccccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{array}$$

0 1 2 3     4 5 6 7     8 9 10 11

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

Position 0 in read $r_2$ corresponds to $Q$-position 8

# Using Gk arrays (I)

$$R = \begin{array}{cccccccccccccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{array}$$

$Q$-positions: 0 1 2 3 · 4 5 6 7 · 8 9 10 11

---

### Searching a $k$-factor of a given read

Read $r_2$ at position 0 → GAT

Position 0 in read $r_2$ corresponds to $Q$-position 8

---

| $i$ | GkIFA |
|---|---|
| 0 | 3 |
| 1 | 6 |
| 2 | 0 |
| 3 | 1 |
| 4 | 3 |
| 5 | 7 |
| 6 | 2 |
| 7 | 5 |
| 8 | 4 |
| 9 | 3 |
| 10 | 6 |
| 11 | 0 |

| $b$ | GkCFA |
|---|---|
| 0 | 2 |
| 1 | 1 |
| 2 | 1 |
| 3 | 3 |
| 4 | 1 |
| 5 | 1 |
| 6 | 2 |
| 7 | 1 |

| $j$ | GkFA |
|---|---|
| 0 | 11 |
| 1 | 2 |
| 2 | 3 |
| 3 | 6 |
| 4 | 9 |
| 5 | 0 |
| 6 | 4 |
| 7 | 8 |
| 8 | 7 |
| 9 | 10 |
| 10 | 1 |
| 11 | 5 |

# Using Gk arrays (I)

$$R = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{matrix}$$

$Q$-positions

## Searching a $k$-factor of a given read

Read $r_2$ at position 0 → GAT

Position 0 in read $r_2$ corresponds to $Q$-position 8

| $i$ | GkIFA | | $b$ | GkCFA | | $j$ | GkFA |
|-----|-------|---|-----|-------|---|-----|------|
| 0 | 3 | | 0 | 2 | | 0 | 11 |
| 1 | 6 | | | | | 1 | 2 |
| 2 | 0 | | 1 | 1 | | 2 | 3 |
| 3 | 1 | | 2 | 1 | | 3 | 6 |
| 4 | 3 | | | | | 4 | 9 |
| 5 | 7 | | 3 | 3 | | 5 | 0 |
| 6 | 2 | | | | | 6 | 4 |
| 7 | 5 | | 4 | 1 | | 7 | 8 |
| 8 | 4 | | 5 | 1 | | 8 | 7 |
| 9 | 3 | | 6 | 2 | | 9 | 10 |
| 10 | 6 | | | | | 10 | 1 |
| 11 | 0 | | 7 | 1 | | 11 | 5 |

# Using Gk arrays (I)

$$R = \begin{array}{cccccccccccccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{array}$$

$Q$-positions with positions: 0 1 2 3   4 5 6 7   8 9 10 11

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

Position 0 in read $r_2$ corresponds to $Q$-position **8**

| $i$ | GkIFA | | $b$ | GkCFA | | $j$ | GkFA |
|---|---|---|---|---|---|---|---|
| 0 | 3 | | 0 | 2 | | 0 | 11 |
| 1 | 6 | | | | | 1 | 2 |
| 2 | 0 | | 1 | 1 | | 2 | 3 |
| 3 | 1 | | 2 | 1 | | 3 | 6 |
| 4 | 3 | | | | | 4 | 9 |
| 5 | 7 | | 3 | 3 | | 5 | 0 |
| 6 | 2 | | | | | 6 | 4 |
| 7 | 5 | | 4 | 1 | | 7 | 8 |
| **8** | 4 | | 5 | 1 | | 8 | 7 |
| 9 | 3 | | 6 | 2 | | 9 | 10 |
| 10 | 6 | | | | | 10 | 1 |
| 11 | 0 | | 7 | 1 | | 11 | 5 |

# Using Gk arrays (I)

$$R = \overset{\text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18}}{\underset{\text{0 1 2 3 \quad 4 5 6 7 \quad 8 9 10 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_2$ at position 0 $\rightarrow$ GAT

Position 0 in read $r_2$ corresponds to $Q$-position **8**

| $i$ | GkIFA |
|-----|-------|
| 0 | 3 |
| 1 | 6 |
| 2 | 0 |
| 3 | 1 |
| 4 | 3 |
| 5 | 7 |
| 6 | 2 |
| 7 | 5 |
| **8** | 4 |
| 9 | 3 |
| 10 | 6 |
| 11 | 0 |

| $b$ | GkCFA |
|-----|-------|
| 0 | 2 |
| 1 | 1 |
| 2 | 1 |
| 3 | 3 |
| 4 | 1 |
| 5 | 1 |
| 6 | 2 |
| 7 | 1 |

| $j$ | GkFA |
|-----|------|
| 0 | 11 |
| 1 | 2 |
| 2 | 3 |
| 3 | 6 |
| 4 | 9 |
| 5 | 0 |
| 6 | 4 |
| 7 | 8 |
| 8 | 7 |
| 9 | 10 |
| 10 | 1 |
| 11 | 5 |

**Information**

There is only 1
$P$-$k$-factor GAT
in $R$.

# Using Gk arrays (II)

$$R = \begin{array}{c} \phantom{.} \\ \mathtt{A\,T\,A\,A\,C\,G\,A\,T\,A\,G\,T\,C\,G\,A\,T\,A\,A\,C\,\$} \end{array}$$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

0 1 2 3    4 5 6 7    8 9 10 11

$Q$-positions

# Using Gk arrays (II)

$$R = \text{A T A A C G A T A G T C G A T A A C \$}$$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

0 1 2 3    4 5 6 7    8 9 10 11

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 $\rightarrow$ ATA

# Using Gk arrays (II)

$$R = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{matrix}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA

| $r_1$ at position 0 |
|:---:|

# Using Gk arrays (II)

$$R = \overset{\text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18}}{\underset{\text{0 1 2 3 \quad 4 5 6 7 \quad 8 9 10 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}}$$

$Q$-positions

## Searching a $k$-factor of a given read

Read $r_1$ at position 0 → ATA

| $r_1$ at position 0 | → | Read length: 6 → |
|---|---|---|

# Using Gk arrays (II)

$$R = \overset{\substack{0\ 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17\ 18}}{\underset{\substack{0\ 1\ 2\ 3\qquad\ 4\ 5\ 6\ 7\qquad\ 8\ 9\ 10\ 11}}{\text{A T A A C G A T A G T C G A T A A C \$}}}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 $\rightarrow$ ATA

| $r_1$ at position 0 | Read length: 6 | $P$-position $1 \times 6 + 0 = 6$ in $R$ |
|---|---|---|

# Using Gk arrays (II)

$$R = \begin{array}{c} \text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18} \\ \texttt{A T A A C G A T A G T C G A T A A C \$} \\ \text{0 1 2 3 ~~~~ 4 5 6 7 ~~~~ 8 9 10 11} \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 $\rightarrow$ ATA

| $r_1$ at position 0 | Read length: 6 → | $P$-position $1 \times 6 + 0 = 6$ in $R$ | $k = 3$ → |

# Using Gk arrays (II)

$$R = \begin{array}{c} 0\ 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\ 10\ 11\ 12\ 13\ 14\ 15\ 16\ 17\ 18 \\ \text{A T A A C G A T A G T C G A T A A C \$} \\ 0\ 1\ 2\ 3\quad\quad 4\ 5\ 6\ 7\quad\quad 8\ 9\ 10\ 11 \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA

Position 0 in read $r_1$ corresponds to $Q$-position 4

| $r_1$ at position 0 | | |
|---|---|---|

Read length: 6 →

| $P$-position $1 \times 6 + 0 = 6$ in $R$ |
|---|

$k = 3$ →

| $Q$-position $6 - (3 - 1) * 1 = 4$ in $R$ |
|---|

# Using Gk arrays (II)

$$R = \begin{array}{c} \text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18} \\ \text{A T A A C G A T A G T C G A T A A C \$} \\ \text{0 1 2 3 \quad 4 5 6 7 \quad 8 9 10 11} \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 $\rightarrow$ ATA

Position 0 in read $r_1$ corresponds to $Q$-position 4

# Using Gk arrays (II)

$$R = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \\ 0 & 1 & 2 & 3 & & & 4 & 5 & 6 & 7 & & & 8 & 9 & 10 & 11 \end{matrix}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 $\rightarrow$ ATA

Position 0 in read $r_1$ corresponds to $Q$-position **4**

| $i$ | GkIFA |   | $b$ | GkCFA |   | $j$ | GkFA |
|-----|-------|---|-----|-------|---|-----|------|
| 0 | 3 | | 0 | 2 | | 0 | 11 |
| 1 | 6 | | | | | 1 | 2 |
| 2 | 0 | | 1 | 1 | | 2 | 3 |
| 3 | 1 | | 2 | 1 | | 3 | 6 |
| 4 | 3 | | | | | 4 | 9 |
| 5 | 7 | | 3 | 3 | | 5 | 0 |
| 6 | 2 | | | | | 6 | 4 |
| 7 | 5 | | 4 | 1 | | 7 | 8 |
| 8 | 4 | | 5 | 1 | | 8 | 7 |
| 9 | 3 | | 6 | 2 | | 9 | 10 |
| 10 | 6 | | | | | 10 | 1 |
| 11 | 0 | | 7 | 1 | | 11 | 5 |

# Using Gk arrays (II)

$$R = \overset{\text{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18}}{\text{A T A A C G A T A G T C G A T A A C \$}}$$
$$\underset{\text{0 1 2 3 ~~~ 4 5 6 7 ~~~~ 8 9 10 11}}{}$$
$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA
Position 0 in read $r_1$ corresponds to $Q$-position 4

| $i$ | GkIFA | | $b$ | GkCFA | | $j$ | GkFA |
|---|---|---|---|---|---|---|---|
| 0 | 3 | | 0 | 2 | | 0 | 11 |
| 1 | 6 | | | | | 1 | 2 |
| 2 | 0 | | 1 | 1 | | 2 | 3 |
| 3 | 1 | | 2 | 1 | | 3 | 6 |
| 4 | 3 | | | | | 4 | 9 |
| 5 | 7 | | 3 | 3 | | 5 | 0 |
| 6 | 2 | | | | | 6 | 4 |
| 7 | 5 | | 4 | 1 | | 7 | 8 |
| 8 | 4 | | 5 | 1 | | 8 | 7 |
| 9 | 3 | | 6 | 2 | | 9 | 10 |
| 10 | 6 | | | | | 10 | 1 |
| 11 | 0 | | 7 | 1 | | 11 | 5 |

11

# Using Gk arrays (II)

$$R = \begin{array}{cccccccccccccccccccc} \scriptstyle 0 & \scriptstyle 1 & \scriptstyle 2 & \scriptstyle 3 & \scriptstyle 4 & \scriptstyle 5 & \scriptstyle 6 & \scriptstyle 7 & \scriptstyle 8 & \scriptstyle 9 & \scriptstyle 10 & \scriptstyle 11 & \scriptstyle 12 & \scriptstyle 13 & \scriptstyle 14 & \scriptstyle 15 & \scriptstyle 16 & \scriptstyle 17 & \scriptstyle 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA

Position 0 in read $r_1$ corresponds to $Q$-position **4**

| $i$ | GkIFA |
|---|---|
| 0 | 3 |
| 1 | 6 |
| 2 | 0 |
| 3 | 1 |
| 4 | 3 |
| 5 | 7 |
| 6 | 2 |
| 7 | 5 |
| 8 | 4 |
| 9 | 3 |
| 10 | 6 |
| 11 | 0 |

| $b$ | GkCFA |
|---|---|
| 0 | 2 |
| 1 | 1 |
| 2 | 1 |
| 3 | 3 |
| 4 | 1 |
| 5 | 1 |
| 6 | 2 |
| 7 | 1 |

| $j$ | GkFA |
|---|---|
| 0 | 11 |
| 1 | 2 |
| 2 | 3 |
| 3 | 6 |
| 4 | 9 |
| 5 | 0 |
| 6 | 4 |
| 7 | 8 |
| 8 | 7 |
| 9 | 10 |
| 10 | 1 |
| 11 | 5 |

# Using Gk arrays (II)

$$R = \begin{array}{c} \texttt{0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18} \\ \texttt{A T A A C G A T A G T C G A T A A C \$} \\ \texttt{0 1 2 3 \quad 4 5 6 7 \quad 8 9 10 11} \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 $\rightarrow$ ATA

Position 0 in read $r_1$ corresponds to $Q$-position **4**

| $i$ | GkIFA |
|-----|-------|
| 0 | 3 |
| 1 | 6 |
| 2 | 0 |
| 3 | 1 |
| **4** | 3 |
| 5 | 7 |
| 6 | 2 |
| 7 | 5 |
| 8 | 4 |
| 9 | 3 |
| 10 | 6 |
| 11 | 0 |

| $b$ | GkCFA |
|-----|-------|
| 0 | 2 |
| 1 | 1 |
| 2 | 1 |
| 3 | **3** |
| 4 | 1 |
| 5 | 1 |
| 6 | 2 |
| 7 | 1 |

| $j$ | GkFA |
|-----|------|
| 0 | 11 |
| 1 | 2 |
| 2 | 3 |
| 3 | 6 |
| 4 | 9 |
| 5 | 0 |
| 6 | 4 |
| 7 | 8 |
| 8 | 7 |
| 9 | 10 |
| 10 | 1 |
| 11 | 5 |

**Information**

There are 3
$P$-$k$-factors ATA in $R$.
But...where are they?

# Using Gk arrays (II)



**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA

Position 0 in read $r_1$ corresponds to $Q$-position **4**

**Information**

There are 3
$P$-$k$-factors ATA in $R$.
But... where are they?

# Using Gk arrays (II)

$R =$ A T A A C G A T A G T C G A T A A C $

positions: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

$Q$-positions: 0 1 2 3 (gap) 4 5 6 7 (gap) 8 9 10 11

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA

Position 0 in read $r_1$ corresponds to $Q$-position **4**

| $i$ | GkIFA |
|-----|-------|
| 0 | 3 |
| 1 | 6 |
| 2 | 0 |
| 3 | 1 |
| 4 | 3 |
| 5 | 7 |
| 6 | 2 |
| 7 | 5 |
| 8 | 4 |
| 9 | 3 |
| 10 | 6 |
| 11 | 0 |

| $b$ | GkCFA |
|-----|-------|
| 0 | 2 |
| 1 | 1 |
| 2 | 1 |
| 3 | 3 |
| 4 | 1 |
| 5 | 1 |
| 6 | 2 |
| 7 | 1 |

| $j$ | GkFA |
|-----|------|
| 0 | 11 |
| 1 | 2 |
| 2 | 3 |
| 3 | 6 |
| 4 | 9 |
| 5 | 0 |
| 6 | 4 |
| 7 | 8 |
| 8 | 7 |
| 9 | 10 |
| 10 | 1 |
| 11 | 5 |

**Information**

There are 3 $P$-$k$-factors ATA in $R$. But...where are they?

**$Q$-positions of ATA**

9  $r_2$, position 1

0  $r_0$, position 0

4  $r_1$, position 0

# Using Gk arrays (II)

$$R = \begin{array}{ccccccccccccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 & 17 & 18 \\ A & T & A & A & C & G & A & T & A & G & T & C & G & A & T & A & A & C & \$ \end{array}$$

$Q$-positions

**Searching a $k$-factor of a given read**

Read $r_1$ at position 0 → ATA

Position 0 in read $r_1$ corresponds to $Q$-position **4**

| $i$ | GkIFA |
|---|---|
| 0 | 3 |
| 1 | 6 |
| 2 | 0 |
| 3 | 1 |
| 4 | 3 |
| 5 | 7 |
| 6 | 2 |
| 7 | 5 |
| 8 | 4 |
| 9 | 3 |
| 10 | 6 |
| 11 | 0 |

| $b$ | GkCFA | GkCFPS |
|---|---|---|
| 0 | 2 | 2 |
| 1 | 1 | 3 |
| 2 | 1 | 4 |
| 3 | 3 | 7 |
| 4 | 1 | 8 |
| 5 | 1 | 9 |
| 6 | 2 | 11 |
| 7 | 1 | 12 |

| $j$ | GkFA |
|---|---|
| 0 | 11 |
| 1 | 2 |
| 2 | 3 |
| 3 | 6 |
| 4 | 9 |
| 5 | 0 |
| 6 | 4 |
| 7 | 8 |
| 8 | 7 |
| 9 | 10 |
| 10 | 1 |
| 11 | 5 |

**Information**

There are 3
$P$-$k$-factors ATA in $R$.
But... where are they?

**$Q$-positions of ATA**

9 $r_2$, position 1

0 $r_0$, position 0

4 $r_1$, position 0

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

$\Rightarrow$ The number of reads in which occurs a $P$-$k$-factor is not its total number of occurrences in $R$

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

$\Rightarrow$ The number of reads in which occurs a $P$-$k$-factor is not its total number of occurrences in $R$

**(Modified) Example**

| GkCFA | $i$ | GkFA |
|-------|-----|------|
|  | $\vdots$ | $\vdots$ |
|  | 4 | 11 |
|  | 5 | 0 |
| 3 | 6 | 4 |
|  | 7 | 9 |
|  | $\vdots$ | $\vdots$ |

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

$\Rightarrow$ The number of reads in which occurs a $P$-$k$-factor is not its total number of occurrences in $R$

**(Modified) Example**

| GkCFA | $i$ | GkFA | |
|---|---|---|---|
| | $\vdots$ | $\vdots$ | |
| | 4 | 11 | $\longrightarrow$ read 2 |
| | 5 | 0 | $\longrightarrow$ read 0 |
| 3 | 6 | 4 | $\longrightarrow$ read 1 |
| | 7 | 9 | $\longrightarrow$ read 2 |
| | $\vdots$ | $\vdots$ | |

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

$\Rightarrow$ The number of reads in which occurs a $P$-$k$-factor is not its total number of occurrences in $R$

**(Modified) Example**

| GkCFA | $i$ | GkFA | |
|---|---|---|---|
| | $\vdots$ | $\vdots$ | |
| | 4 | 11 | $\longrightarrow$ read 2 |
| | 5 | 0 | $\longrightarrow$ read 0 |
| 3 | 6 | 4 | $\longrightarrow$ read 1 |
| | 7 | 9 | $\longrightarrow$ read 2 |
| | $\vdots$ | $\vdots$ | |

**Solutions for counting reads**

► Use a mask to known which read have already been counted

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

$\Rightarrow$ The number of reads in which occurs a $P$-$k$-factor is not its total number of occurrences in $R$

**(Modified) Example**

| GkCFA | $i$ | GkFA | |
|-------|-----|------|---|
| | $\vdots$ | $\vdots$ | |
| | 4 | 11 | $\longrightarrow$ read 2 |
| | 5 | 0 | $\longrightarrow$ read 0 |
| 3 | 6 | 4 | $\longrightarrow$ read 1 |
| | 7 | 9 | $\longrightarrow$ read 2 |
| | $\vdots$ | $\vdots$ | |

**Solutions for counting reads**

- Use a mask to known which read have already been counted
- Sort the entries (when querying or at construction).

# Multiplicity of $P$-$k$-factors

**Problem**

What if a $P$-$k$-factor occurs many times in the same read?

$\Rightarrow$ The number of reads in which occurs a $P$-$k$-factor is not its total number of occurrences in $R$

**(Modified) Example**

| GkCFA | $i$ | GkFA | |
|-------|-----|------|--|
| | $\vdots$ | $\vdots$ | |
| | 4 | 0 | $\longrightarrow$ read 0 |
| | 5 | 4 | $\longrightarrow$ read 1 |
| 3 | 6 | 9 | $\longrightarrow$ read 2 |
| | 7 | 11 | $\longrightarrow$ read 2 |
| | $\vdots$ | $\vdots$ | |

**Solutions for counting reads**

- Use a mask to known which read have already been counted
- Sort the entries (when querying or at construction).

# Complexities

**Space complexities**

GkFA, GkIFA Number of entries: Number of reads $\times$ (Read length $-k+1$)

GkCFA Number of entries: Number of distinct $P$-$k$-factors

# Complexities

## Space complexities

GkFA, GkIFA  Number of entries: Number of reads $\times$ (Read length $-k + 1$)

GkCFA  Number of entries: Number of distinct $P$-$k$-factors

## Time complexities

**Q1**  (counting $P$-$k$-factors)          $O(1)$
**Q2**  (retrieving positions in reads)     $O(occ)$
**Q3**  (counting reads)                    $O(occ)$

where $occ$ is the number of occurrences of the $P$-$k$-factor in the reads.

# Complexities with the classical solution

**SA-based solution**

Build the suffix array of $R$, the inverse suffix array and the LCP array.

# Complexities with the classical solution

**SA-based solution**

Build the suffix array of $R$, the inverse suffix array and the LCP array.

**Space complexities**

Three arrays containing (number of reads $\times$ length of the reads) elements each

# Complexities with the classical solution

**SA-based solution**

Build the suffix array of $R$, the inverse suffix array and the LCP array.

**Space complexities**

Three arrays containing (number of reads $\times$ length of the reads) elements each

**Time complexities**

| **Q1** | (counting $P$-$k$-factors) | $O(occ_R)$ |
| **Q2** | (retrieving positions in reads) | $O(occ_R)$ |
| **Q3** | (counting reads) | $O(occ_R+\text{number of reads})$ |

where $occ_R$ is the number of occurrences of the $k$-factors in $R$.

# Complexities with the classical solution

**SA-based solution**

Build the suffix array of $R$, the inverse suffix array and the LCP array.

**Space complexities**

Three arrays containing (number of reads $\times$ length of the reads) elements each

**Time complexities**

| **Q1** | (counting $P$-$k$-factors) | $O(occ_R)$ |
| **Q2** | (retrieving positions in reads) | $O(occ_R)$ |
| **Q3** | (counting reads) | $O(occ_R+\text{number of reads})$ |

where $occ_R$ is the number of occurrences of the $k$-factors in $R$.

**Improvements over a SA-based solution**

Space At least $(3 \times (k-1) \times$ number of reads) elements

Time No dependency on the number of reads, no dependency on the number of occurrences in $R$

# Time and space construction in practice

**Data**

- Fruit fly sequences from a Genome Analyzer II
- 7,000,000 reads
- read length: 75

# Time and space construction in practice



Construction time
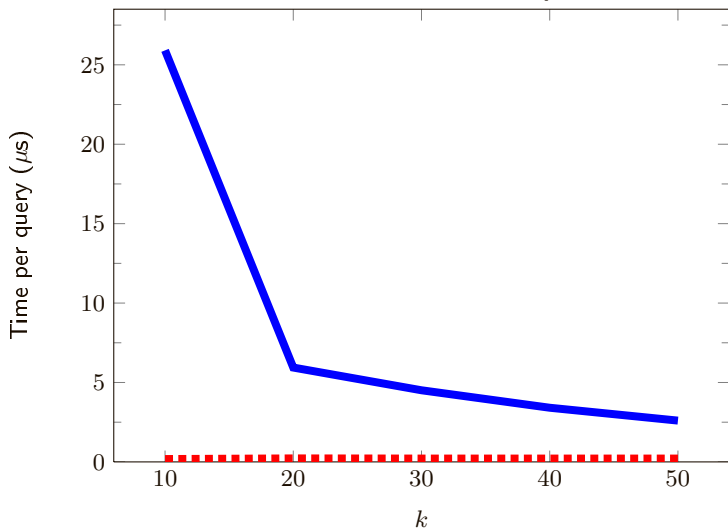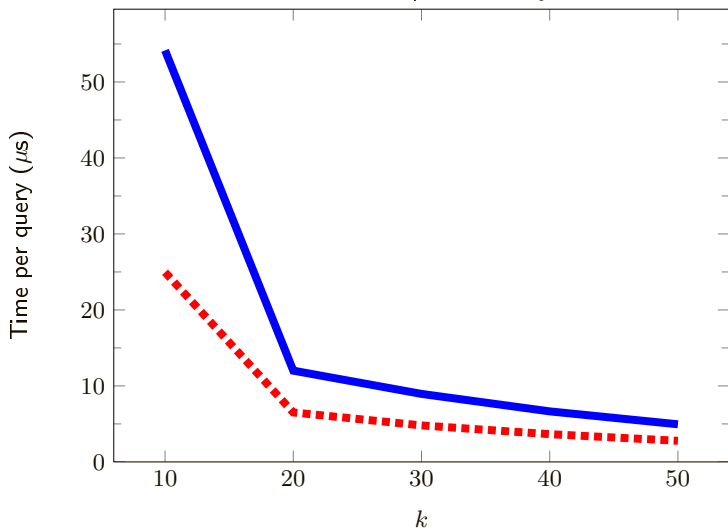
# Time and space construction in practice



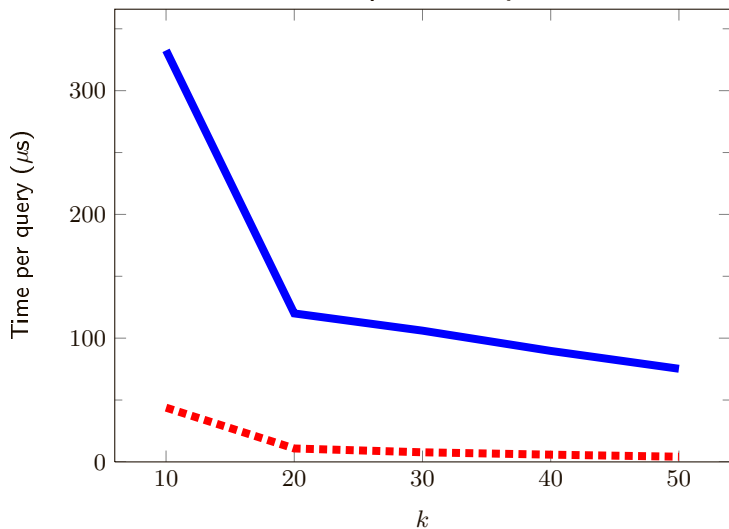Maximal space usage

# Query time

What is the number of occurrences of $f$ in the reads?

# Query time

## What are the occurrence positions of $f$ in the reads?

# Query time

In how many reads does $f$ occur?

# Conclusions and Perspectives

**Efficiency**

Gk arrays allow to query more reads in less time than a suffix array-based method

# Conclusions and Perspectives

**Efficiency**

Gk arrays allow to query more reads in less time than a suffix array-based method

**Variable-length reads**

We can deal with variable-length reads by adding a bit vector for identifying the end of reads in $R$

# Conclusions and Perspectives

**Efficiency**

Gk arrays allow to query more reads in less time than a suffix array-based method

**Variable-length reads**

We can deal with variable-length reads by adding a bit vector for identifying the end of reads in $R$

**Compressing Gk arrays**

Can we adapt compression techniques to Gk arrays?
$\rightarrow$ new space/time tradeoff

# Conclusions and Perspectives

**Efficiency**

Gk arrays allow to query more reads in less time than a suffix array-based method

**Variable-length reads**

We can deal with variable-length reads by adding a bit vector for identifying the end of reads in $R$

**Compressing Gk arrays**

Can we adapt compression techniques to Gk arrays?
$\rightarrow$ new space/time tradeoff

**Updating Gk arrays**

Can we efficiently update Gk arrays?
$\rightarrow$ read correction