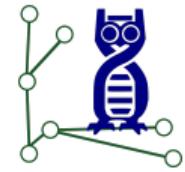




UNIVERSITÄT
DES
SAARLANDES



Design of Worst-Case-Optimal Spaced Seeds

Jens Zentgraf, Sven Rahmann

WABI 2025

Motivation

k-mers

A *k*-mer is a substring of length *k*.

Motivation

k-mers

A *k*-mer is a substring of length *k*.

Alignment-free applications

- Taxonomic classification
- Contamination removal
- Gene expression
- Phylogenetic tree reconstruction
- Pangenomics
- ...

Motivation

k-mers

A *k*-mer is a substring of length *k*.

Alignment-free applications

- Taxonomic classification
- Contamination removal
- Gene expression
- Phylogenetic tree reconstruction
- Pangenomics
- ...

Advantages

- Fast information retrieval
(hashing, sorting, ...)
- Subsampling
 - Minimizers
 - Open/Closed Syncmers
 - ...

Motivation

k-mers

A *k*-mer is a substring of length *k*.

Alignment-free applications

- Taxonomic classification
- Contamination removal
- Gene expression
- Phylogenetic tree reconstruction
- Pangenomics
- ...

Advantages

- Fast information retrieval
(hashing, sorting, ...)
- Subsampling
 - Minimizers
 - Open/Closed Syncmers
 - ...

Disadvantages

- Intolerant against errors
- Reference bias

Disadvantages of k -mers

An error changes k **consecutive** k -mers
in a read.

Disadvantages of k -mers

An error changes k **consecutive** k -mers in a read.

TACGAGCTA**G**CATCAGC
CGAGCT**A**
GAGCT**A**G
AGCT**A**GC
GCT**A**GCA
CT**A**GCAT
T**A**GCATC
AGCATCA

Disadvantages of k -mers

An error changes k **consecutive** k -mers in a read.

Large k

- **High** specificity
- **Lower** number of k -mers
- **Few (No)** unchanged k -mers if an error occurs

TACGAGCT**A**GCATCAGC
TACGAGCT**A**GCATCA
ACGAGCT**A**GCATCAG
CGAGCT**A**GCATCAGC

Disadvantages of k -mers

An error changes k **consecutive** k -mers in a read.

Large k

- **High** specificity
- **Lower** number of k -mers
- **Few (No)** unchanged k -mers if an error occurs

Small k

- **Low** specificity
- **High** number of k -mers
- **Many** unchanged k -mers if an error occurs

TACGAGCT**A**GCATCAGC
TAC AGC **A**GC TCA
ACG GCT GCA CAG
CGA CTA**A** CAT AGC
GAG T**A**G ATC

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)

$k = 5, w = 7$

TACGAGCTAGCATCAGC
###_#_##

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)
- Mask, or tuple of offsets of #
 - Mask: #_#_#
 - Tuple: (0,2,4)

$k = 5, w = 7$

TACGAGCTAGCATCAGC
###_ _##

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)
- Mask, or tuple of offsets of #
 - Mask: #_#_#
 - Tuple: (0,2,4)
- Constraints for masks:
 - Symmetric
 - Start and end with #

$k = 5, w = 7$

TACGAGCTAGCATCAGC
###_#_##

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)
- Mask, or tuple of offsets of #
 - Mask: #_#_#
 - Tuple: (0,2,4)
- Constraints for masks:
 - Symmetric
 - Start and end with #

$k = 5, w = 7$

TACGAGCTAGCATCAGC
#_###_#

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)
- Mask, or tuple of offsets of #
 - Mask: #_#_#
 - Tuple: (0,2,4)
- Constraints for masks:
 - Symmetric
 - Start and end with #

$$k = 5, w = 7$$

TACGAGCTAGCATCAGC
T C G A C
A G A G T
C A G C A
G G C T G
A C T A C
G T A G A
C A G C T
T G C A C
A C A T A
G A T C G
C A T C A G C

Spaced seeds (gapped k -mers)

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)
- Mask, or tuple of offsets of #
 - Mask: #_#_#
 - Tuple: (0,2,4)
- Constraints for masks:
 - Symmetric
 - Start and end with #

$$k = 5, w = 7$$

TACGAGCTAGCATCAGC
T CGA C
A GAG T
C AGC A
G GCT G
A CTA C
G TAG A
C AGC T
T GCA C
A CAT A
G ATC G
CATCAGC

Spaced seeds (gapped k -mers)

$k = 5, w = 7$

Spaced seeds

- k significant positions (#)
- Window of size w
- $w - k$ ignored positions, gaps, ...(_)
- Mask, or tuple of offsets of #
 - Mask: #_#_#
 - Tuple: (0,2,4)
- Constraints for masks:
 - Symmetric
 - Start and end with #

An error changes k distributed k -mers.

Designing spaced seeds

Lossless setting ([NR08, Bri14, FCLST07])

- Given k , w , n (length of the sequence), c number of changes.
- Find a seed that guarantees at least one hit independently of the error positions.
- Decision problem already NP-complete.

Designing spaced seeds

Lossless setting ([NR08, Bri14, FCLST07])

- Given k, w, n (length of the sequence), c number of changes.
- Find a seed that guarantees at least one hit independently of the error positions.
- Decision problem already NP-complete.

Probabilistic setting ([BM08, CP10, No 17])

- No fixed number of changes.
- Compare two sequences that differ at each position with a probability.
- Compute probability distribution of the random number of hits (at least one hit).
- Computing sensitivity of a mask is NP-hard.

Goal

Goal

- Compute the following guarantees for a mask:

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.

CTATTTGTGTAATAGCCCATCGAGACGGCCCTCATGGCGCTCAAGTGAACCTCTATTGGGAGTCGTATGCATGTCAGCGGTACGGTATAGATG

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.

$k = 27$ $w = 27$ #####

CTATTGCTAATAGCCCCATCGAGACGGCCCTCATGGCGCTCAAGTGAACCTCTATTGGGGAGTCGGTCATATGCATGTCAGCGGTGACGGTATAGATG

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.

$k = 27 \ w = 27 \ #####\#####\#####\#####\#####\#####\#$

CTATTTGTGTAATAGCCCATCGAGACGCCCCCTCATGGCGCTCAAGTGAACCTCTATTGGGAGTCGTATGCATGTCAGCGGTACGGTATAGATG

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.

$k = 27 \ w = 29 \ #####_#####_#####$

CTATTTGTGTAATAGCCCATCGAGACGCCCCCTATGGCGCTCAAGTGAACCTCTATTGGGAGTCGTATGCATGTCAGCGGTGACGGTATAGATG

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.

$k = 27$ $w = 29$ ##### ##### ###### ##### #######

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.

$k = 27 \ w = 29 \ #####_#####_#####$

CTATTTGTGTAATAGCCCCATCGAGACGGCCCTCATGGCGCTCAAGTGAACCTCTATTGGGAGTCGTATATGCATGTCAGCGGTACGGTATAGATG

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.

$k = 27$ $w = 29$ #####_#####_#####_#####_#####

CTATTTGTGTAATAGCCCATCGAG A CGGCCCTCATGGCGCTC A GTGAACCTCTATTGGGAGTCCGTCA TGCATGTCAGCGGTGACGGTATAAGATG
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

$k = 27$ $w = 29$ #####_#####_#####_#####_#####

CTATTTGTGTAATAGCCCATCGAGA CGGGCCCTCATGGCGCTA AGTGAACCTCTATTGGGAGTCCGTCA TAGCATGTCAGCGGTGACGGTATAAGATG
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####
#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####_#####

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

$k = 27$ $w = 29$ #####_#####_#####_#####_#####

CTATTGTGTAATAGCCCATCGAG  CGGGCCCTCATGGCGCTC  AGTGAACCTCTATTGGGAGTCGTCA  TGCAATGTCAGCGGTACGGTATAGATG
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

$k = 27$ $w = 29$ #####_#####_#####_#####

CTATTTGTGTAATAGCCCATCGAGACGGC_↓CCCTCATG_↓GCGCTCAAGTGAACCTCTATTGGAGTCGTCA_↓TATGCATGTCAGCGGTACGGTATAAGATG

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

$k = 27$ $w = 29$ ##### ##### ###### ##### #######

CTATTGTAAAGCCATCGAGACGGCCCCTCATG_GCGCTCAAGTGAACCTCTATTGGAGTCGTCA_TATGCATGTCAGCGGTACGGTATAAGATG

- ##### - #####
- ##### - #####
- ##### - #####

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

$k = 27$ $w = 29$ #####_#####_#####_#####_#####

CTATTGTGTAATAGCCCATCGAGACGGCCCTCATGCGCTCAAGTGAACCTCTATTGGAGTCGTATGCATGTCAGCGGTACGGTATAGATG
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####
#####_#####_#####_#####_#####
#####_#####_#####_#####_##### (51)

Goal

Goal

- Compute the following guarantees for a mask:
 - How many errors c are tolerated, until we cannot find at least one k -mer.
 - For a given c , what is the minimal number of k -mers we can find.
 - For a given c , what is the minimal number of bases covered by k -mers.

Solution

We provide 3 ILPs to compute these values for a given mask.

Number of Tolerated Errors

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

Number of Tolerated Errors

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

$k=9$
####-####-###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Number of Tolerated Errors

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$\overbrace{\text{### - ##### - ####}}^{k=9}$
 $\overbrace{\quad\quad\quad\quad}^{w=11}$
 $(0, 1, 2, 4, 5, 6, 8, 9, 10)$

Number of Tolerated Errors

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$n=28$ CTATTTGTGTAATAGCCCCATCGAGACGG
x 00000000000000000000000000000000
Algorithmic Bioinformatics



$k=9$
 $\overbrace{\#\#\# - \#\#\# - \#\#\#}^{w=11}$
(0, 1, 2, 4, 5, 6, 8, 9, 10)

Number of Tolerated Errors

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$n=28$ CTATTTGTGTAATAGCCCCATCGAGACGG
x 00000000000010000000000000000000
Algorithmic Bioinformatics



$k=9$
 $\underbrace{\#\#\# - \#\#\# - \#\#\#}$
 $w=11$
 $(0, 1, 2, 4, 5, 6, 8, 9, 10)$

Number of Tolerated Errors

Variables

- k : weight of the mask
 - n : sequence length
 - w : window length
 - κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$$\underbrace{\# \# \# - \# \# \# - \# \# \#}_{w=11}^{k=9}$$

Number of Tolerated Errors

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$n=28$ CTATTTGTGTAATAGCCCCATCGAGACGG
x 00000000000010000000000000000000
Algorithmic Bioinformatics



$k=9$
 $\underbrace{\#\#\# - \#\#\# - \#\#\#}$
 $w=11$
 $(0, 1, 2, 4, 5, 6, 8, 9, 10)$

Number of Tolerated Errors

Minimize $\sum_{i \in [n]} x_i$,

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$n=28$ CTATTTGTGTAATAGCCCCATCGAGACGG
x 00000000000010000000000000000000
Algorithmic Bioinformatics



$k=9$
 $\overbrace{\text{#### - ##### - ####}}$
 $w=11$
(0, 1, 2, 4, 5, 6, 8, 9, 10)

Number of Tolerated Errors

Minimize $\sum_{i \in [n]} x_i$,

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$n=28$ CTATTTGTGTAA**T**AGCCCCATCGAGACGG
 x 00000000000010000000000000000000
Algorithmic Bioinformatics



$k=9$
####_####_###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Number of Tolerated Errors

Minimize $\sum_{i \in [n]} x_i$,

such that $\sum_{j \in \kappa} x_{p+j} \geq 1 \quad p \in [n - w + 1]$.

$n=28$ CTATTTGTGTAA**T**AGCCCCATCGAGACGG
x 00000000000010000000000000000000
Algorithmic Bioinformatics



Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$k=9$
 $\overbrace{\#\# _ \#\# _ \#\#}$
 $w=11$
 $(0, 1, 2, 4, 5, 6, 8, 9, 10)$

Number of Tolerated Errors

Minimize $\sum_{i \in [n]} x_i,$

such that $\sum_{j \in \kappa} x_{p+j} \geq 1 \quad p \in [n - w + 1].$

Tolerated Errors

$$c = \sum_{i \in [n]} x_i - 1$$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$k=9$
####_####_###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Minimal Number of Hits

Variables

- k : weight of the mask
- n : sequence length
- w : window length
- c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions

$k=9$
####_####_###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Minimal Number of Hits

Variables

- k : weight of the mask
- n : sequence length
- w : window length
- c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
####_####_###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Minimal Number of Hits

$n=28$ CTATTGTGTAATAGCCCCATCGAGACGG
 x 00000000000010000000000000000000

Variables

- k : weight of the mask
- n : sequence length
- w : window length
- c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
####-####-###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Minimal Number of Hits

$n=28$ $\text{CTATTGTGTAATAGCCCATCGAGACGG}$
 x 000000000001000000000000000
 y 110001000100011111

Variables

- k : weight of the mask
- n : sequence length
- w : window length
- c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
####_####_###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Minimal Number of Hits

Minimize $\sum_{p \in [n-w+1]} y_p,$

$n=28$ CTATTGTGTAATAGCCCCATCGAGACGG
 x 000000000000100000000000000000
 y 110001000100011111

Variables

- k : weight of the mask
- n : sequence length
- w : window length
- c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
####_####_###
 $w=11$
(0,1,2,4,5,6,8,9,10)

Minimal Number of Hits

Minimize $\sum_{p \in [n-w+1]} y_p,$

such that $\sum_{i \in [n]} x_i = c,$

$n=28$ $\text{CTATTGTGTAATAGCCCATCGAGACGG}$
x 0000000000001000000000000000
y 110001000100011111

Variables

- k : weight of the mask
- n : sequence length
- w : window length
- c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
 $\overbrace{\text{####_####_####}}$
 $w=11$
 $(0,1,2,4,5,6,8,9,10)$

Minimal Number of Hits

Minimize $\sum_{p \in [n-w+1]} y_p,$

such that $\sum_{i \in [n]} x_i = c,$

$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n-w+1],$$

$n=28$ $\text{CTATTGTGTAATAGCCCATCGAGACGG}$
x 0000000000001000000000000000
y 110001000100011111

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
 $\overbrace{\text{####_####_####}}$
 $w=11$
 $(0,1,2,4,5,6,8,9,10)$

Minimal Number of Hits

Minimize $\sum_{p \in [n-w+1]} y_p,$

such that $\sum_{i \in [n]} x_i = c,$

$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n-w+1],$$

$$y_p \leq 1 - x_{p+j} \quad p \in [n-w+1], \\ j \in \kappa.$$

$n=28$ $\text{CTATTGTGTAATAGCCCATCGAGACGG}$
x 0000000000001000000000000000
y 110001000100011111

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
 $\overbrace{\text{CTATTGTGTAATAGCCCATCGAGACGG}}$
 $w=11$
 $(0, 1, 2, 4, 5, 6, 8, 9, 10)$

Minimal number of Covered Bases

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers

$k=9$
####_####_###
 $w=11$
(0, 1, 2, 4, 5, 6, 8, 9, 10)

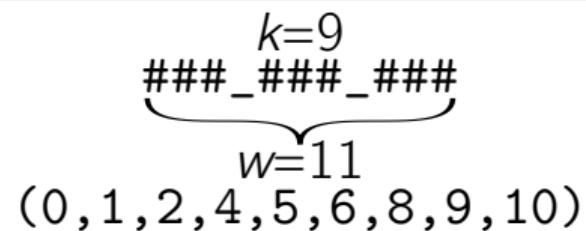
Minimal number of Covered Bases

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Minimal number of Covered Bases

$n=28$ CTATTGTGTAATAGCCCCATCGAGACGG
 x 000000000000100000000000000000
 y 110001000100011111

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions

$k=9$
####_####_###
 $w=11$
(0, 1, 2, 4, 5, 6, 8, 9, 10)

Minimal number of Covered Bases

$n=28$ CTATTGTGTAATAGCCCCATCGAGACGG
 x 000000000000100000000000000000
 y 110001000100011111
 z 111111111110111111111111111111

###_###_###
↓

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions

$k=9$
###_###_###
 $w=11$
(0, 1, 2, 4, 5, 6, 8, 9, 10)

Minimal number of Covered Bases

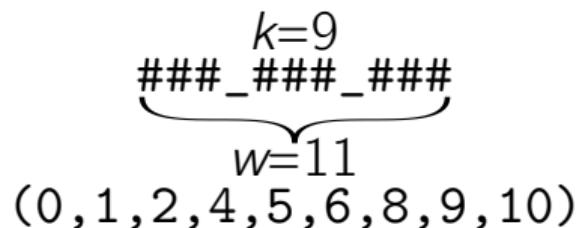
$$\text{Minimize} \quad \sum_{i \in [n]} z_i,$$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Minimal number of Covered Bases

Minimize $\sum_{i \in [n]} z_i,$

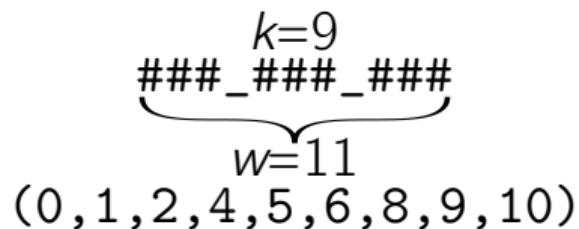
such that $\sum_{i \in [n]} x_i = c,$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Minimal number of Covered Bases

Minimize $\sum_{i \in [n]} z_i,$

such that $\sum_{i \in [n]} x_i = c,$

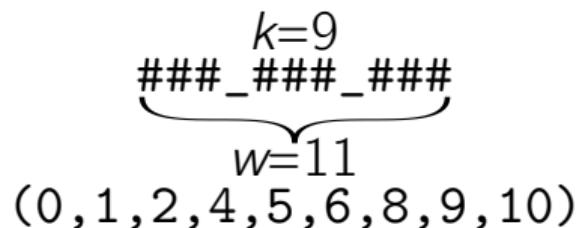
$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n - w + 1],$$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n - w + 1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Minimal number of Covered Bases

Minimize $\sum_{i \in [n]} z_i,$

such that $\sum_{i \in [n]} x_i = c,$

$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n - w + 1],$$

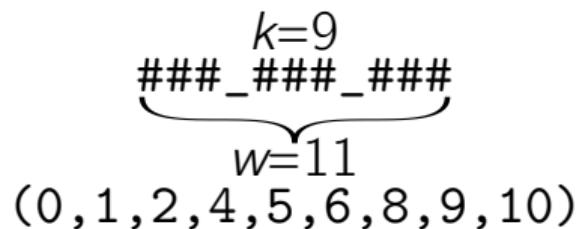
$$y_p \leq 1 - x_{p+j} \quad p \in [n - w + 1], \\ j \in \kappa,$$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n - w + 1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Minimal number of Covered Bases

$$\text{Minimize} \quad \sum_{i \in [n]} z_i,$$

$$\text{such that} \quad \sum_{i \in [n]} x_i = c,$$

$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n - w + 1],$$

$$y_p \leq 1 - x_{p+j} \quad p \in [n - w + 1], \\ j \in \kappa,$$

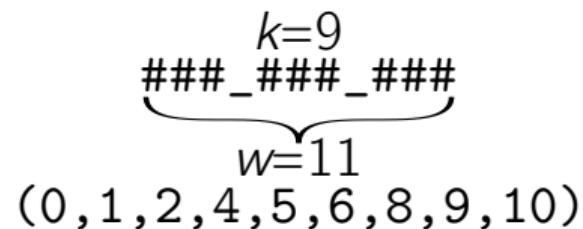
$$z_{p+j} \geq y_p \quad p \in [n - w + 1], \\ j \in \kappa,$$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n - w + 1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Minimal number of Covered Bases

$$\text{Minimize} \quad \sum_{i \in [n]} z_i,$$

$$\text{such that} \quad \sum_{i \in [n]} x_i = c,$$

$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n - w + 1],$$

$$y_p \leq 1 - x_{p+j} \quad p \in [n - w + 1], \quad j \in \kappa,$$

$$z_{p+j} \geq y_p \quad p \in [n - w + 1], \quad j \in \kappa,$$

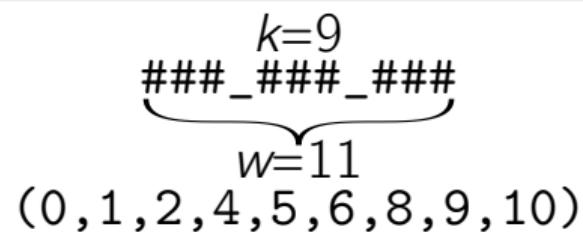
$$z_i \leq \sum_{j \in \kappa, i-j \geq 0} y_{i-j} \quad i \in [n].$$

Variables

- k : weight of the mask ■ n : sequence length
- w : window length ■ c : allowed changes
- κ : k -tuple of offsets

We use the following **binary** variables:

- $x = (x_i)_{i \in [n]}$, indicators of changed positions
- $y = (y_p)_{p \in [n-w+1]}$, indicators of (starting positions of) unchanged κ -mers
- $z = (z_i)_{i \in [n]}$, indicators of covered positions



Evaluation

Read mapper using spaced seeds

- Index which contains the position of each k -mer in the reference.
- Query all k -mers in a read.
- Subtract relative position in the read to get a start position.
- Combine near positions to one.

Evaluation

Read mapper using spaced seeds

- Index which contains the position of each k -mer in the reference.
- Query all k -mers in a read.
- Subtract relative position in the read to get a start position.
- Combine near positions to one.

Challenge data set

- Restrict to area without repeats in the genome.
- All k -mers in these areas occur only once.
- Also all k -mers do not have a Hamming-distance one neighbor anywhere in the genome.

Evaluation

Read mapper using spaced seeds

- Index which contains the position of each k -mer in the reference.
- Query all k -mers in a read.
- Subtract relative position in the read to get a start position.
- Combine near positions to one.

Challenge data set

- Restrict to area without repeats in the genome.
- All k -mers in these areas occur only once.
- Also all k -mers do not have a Hamming-distance one neighbor anywhere in the genome.
- Create reads of length 100.
- Insert errors in equidistant positions.
- Add a random offset uniformly chosen from $\{-3, \dots, 0, \dots, 3\}$

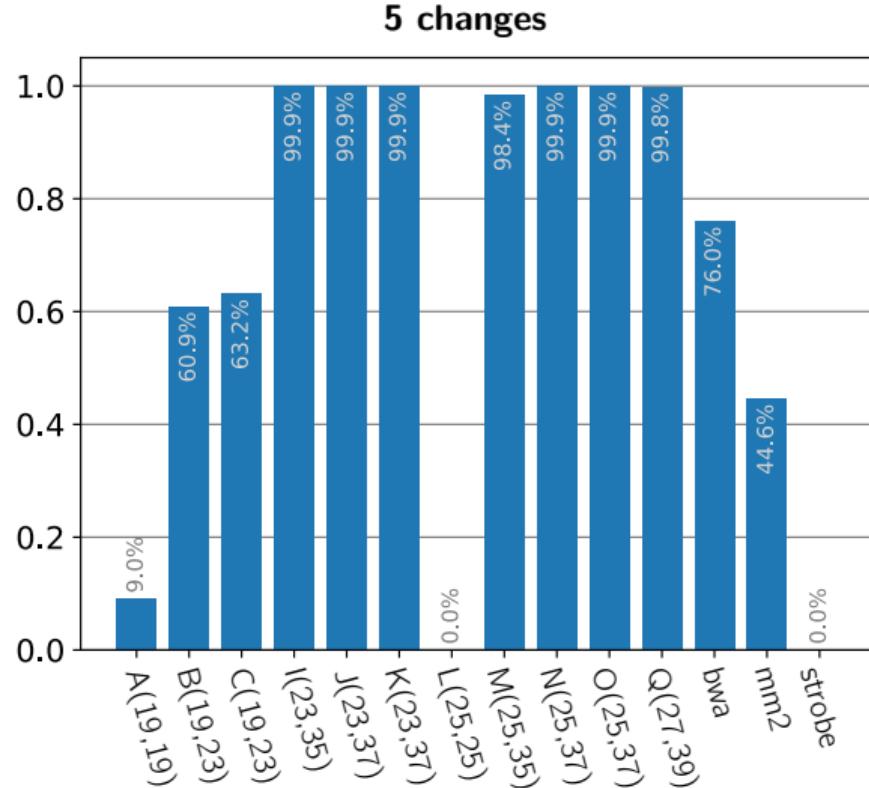
Evaluation

(k, w) shape	mask	$c = 3$		$c = 4$		$c = 5$	
		MH	MC	MH	MC	MH	MC
A (19, 19)	#####_#####_#####_#####_#####_#####	25	43	6	24	0	0
B (19, 23)	####_###_###_###_###_###_###	21	68	11	48	6	42
C (19, 23)	####_###_###_###_###_###_###	21	68	11	48	6	42
I (23, 35)	##_###_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	12	59	5	48	2	34
J (23, 37)	####_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	11	55	4	45	2	34
K (23, 37)	##_###_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	12	58	6	47	2	34
L (25, 25)	#####_#####_#####_#####_#####_#####	1	25	0	0	0	0
M (25, 35)	####_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	9	55	4	42	0	0
N (25, 37)	####_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	9	53	4	45	0	0
O (25, 37)	####_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	8	54	2	47	0	0
Q (27, 39)	#_###_#_#_#_#_#_#_#_#_#_#_#_#_#_#_#	7	52	2	42	0	0

Evaluation

Benchmark setup

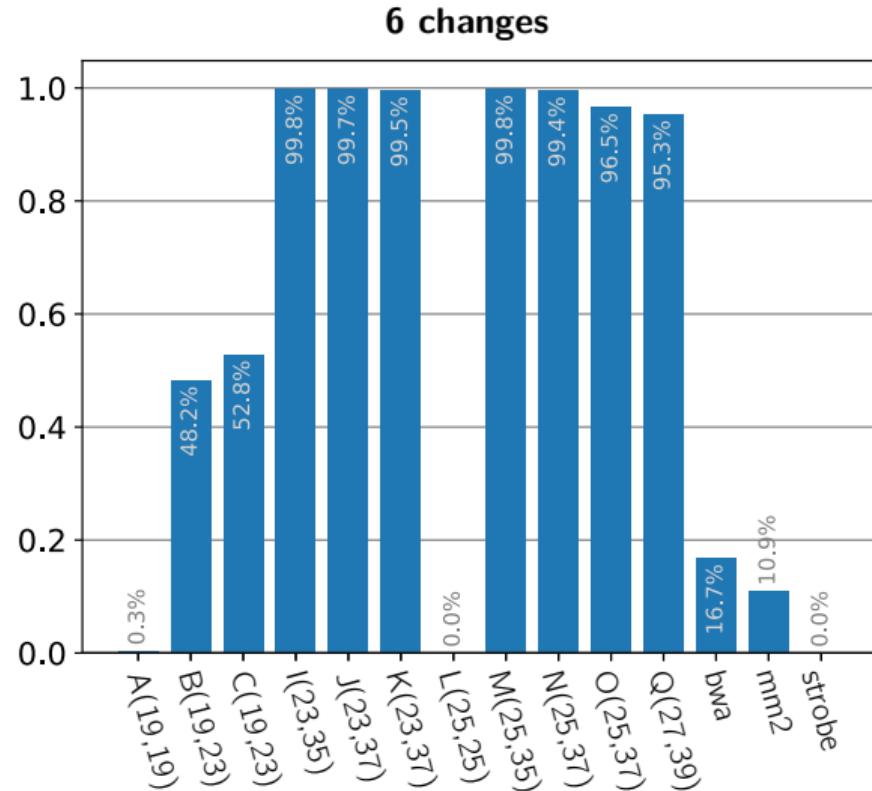
- AMD Ryzen 9 5950X
16-core CPU
- 128 GB RAM
- Gurobi (10.0.3)
- bwa-mem2 (v2.2.1)
- minimap2 (v2.26)
- strobealign (v0.11.0)



Evaluation

Benchmark setup

- AMD Ryzen 9 5950X
16-core CPU
- 128 GB RAM
- Gurobi (10.0.3)
- bwa-mem2 (v2.2.1)
- minimap2 (v2.26)
- strobealign (v0.11.0)



Summary

- Solved ILPs for

- $n = 100$
- $c \in \{3, 4, 5, 6, 7\}$
- $k \in \{15, \dots, 31\}$
- $w \in \{15, \dots, 47\}$

- Running time for one mask is small

- **Good** masks are rare



zentgraf@cs.uni-saarland.de
sven.rahmann@uni-saarland.de

<https://doi.org/10.5281/zenodo.15690315>

$$\text{Minimize} \sum_{i \in [n]} z_i,$$

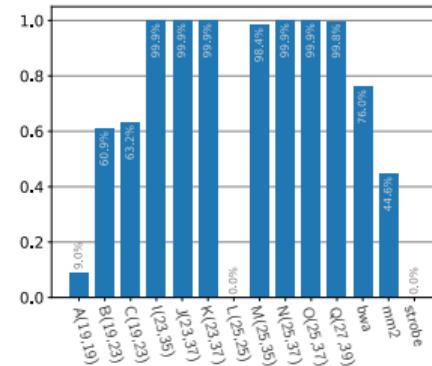
$$\text{such that} \sum_{i \in [n]} x_i = c,$$

$$y_p \geq 1 - \sum_{j \in \kappa} x_{p+j} \quad p \in [n-w+1],$$

$$y_p \leq 1 - x_{p+j} \quad p \in [n-w+1],$$

$$z_{p+j} \geq y_p \quad p \in [n-w+1],$$

$$z_i \leq \sum_{j \in \kappa, i-j \geq 0} y_{i-j} \quad i \in [n].$$



References I

-  [Gary Benson and Denise Y. F. Mak.](#)
Exact distribution of a spaced seed statistic for DNA homology detection.
In Amihood Amir, Andrew Turpin, and Alistair Moffat, editors, *String Processing and Information Retrieval, 15th International Symposium, SPIRE 2008, Melbourne, Australia, November 10-12, 2008. Proceedings*, volume 5280 of *Lecture Notes in Computer Science*, pages 282–293. Springer, 2008.
-  [Karel Brinda.](#)
Languages of lossless seeds.
In Zoltán Ésik and Zoltán Fülöp, editors, *Proceedings 14th International Conference on Automata and Formal Languages, AFL 2014, Szeged, Hungary, May 27-29, 2014*, volume 151 of *EPTCS*, pages 139–150, 2014.
-  [W. H. Chung and S. B. Park.](#)
Hit integration for identifying optimal spaced seeds.
BMC Bioinformatics, 11 Suppl 1(Suppl 1):S37, Jan 2010.
-  [Martin Farach-Colton, Gad M. Landau, Süleyman Cenk Sahinalp, and Dekel Tsur.](#)
Optimal spaced seeds for faster approximate string matching.
J. Comput. Syst. Sci., 73(7):1035–1044, 2007.

References II



L. Noé.

Best hits of 11110110111: model-free selection and parameter-free sensitivity calculation of spaced seeds.
Algorithms Mol Biol, 12:1, 2017.



François Nicolas and Eric Rivals.

Hardness of optimal spaced seed design.
J. Comput. Syst. Sci., 74(5):831–849, 2008.