# Improving an Exact Solution to the (l,d) Planted Motif Problem

Maria Clara Isabel Sia

October 12, 2015

Introduction
EMS-GT
Methodology
Results
Conclusion

The ($l$,$d$) planted motif problem
Definitions of key concepts

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Introduction

Introduction
EMS-GT
Methodology
Results
Conclusion

The ($l,d$) planted motif problem
Definitions of key concepts

## Introduction

- ▶ motifs: repeated sub-sequences in DNA that have some biological significance

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

## Introduction

- ▶ motifs: repeated sub-sequences in DNA that have some biological significance

- ▶ DNA motif finding: search for motifs over a set of DNA sequences, allowing for mismatches due to mutation

Introduction
EMS-GT
Methodology
Results
Conclusion

The ($l,d$) planted motif problem
Definitions of key concepts

# Introduction

- ▶ motifs: repeated sub-sequences in DNA that have some biological significance

- ▶ DNA motif finding: search for motifs over a set of DNA sequences, allowing for mismatches due to mutation

- ▶ known as a difficult problem in computational biology and CS (proven NP-complete)

Introduction
EMS-GT
Methodology
Results
Conclusion

The (*l,d*) planted motif problem
Definitions of key concepts

# The (*l,d*) planted motif problem

*Find a motif of length l=8 across 5 DNA sequences,*
*each containing the motif with at most d=2 mismatches.*

at<mark>cactcgtt</mark>ctcctctaatgtgtaaagacgtactaccgacctta
acgccgaccggtc<mark>cgatcctt</mark>gtatagctcctaacgggcatcagc
tcctgactgcatcgcgatctcggtagtttcctgt<mark>tcatcatt</mark>tt
ggccctca<mark>gcatcgtg</mark>cgtcctgctaacacattcccatgcagctt
tgaaaagaatttacggtaaaggatccacatc<mark>caatcgtg</mark>tgaaag

*Motif:* ccatcgtt

Introduction
EMS-GT
Methodology
Results
Conclusion

The (*l,d*) planted motif problem
Definitions of key concepts

# The (*l,d*) planted motif problem

*Find a motif of length l=8 across 5 DNA sequences,*
*each containing the motif with at most d=2 mismatches.*

at`cactcgtt`ctcctctaatgtgtaaagacgtactaccgacctta
acgccgaccggtc`cgatcctt`gtatagctcctaacgggcatcagc
tcctgactgcatcgcgatctcggtagtttcctgt`tcatcatt`tt
ggccctca`gcatcgtg`cgtcctgctaacacattcccatgcagctt
tgaaaagaatttacggtaaaggatccacatc`caatcgtg`tgaaag

*Motif:* ccatcgtt

Given a set $\mathcal{S} = \{S_1, ... S_n\}$ of $n$ DNA sequences of length $L$,
find $M$, the set of sub-sequences (motifs) of length $l < L$
which occur with at most $d$ mismatches in each sequence in $\mathcal{S}$.

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Definitions of key concepts

- $l$-mer
- Hamming distance $dH(x_1, x_2)$
- $d$-neighborhood $N(x, d)$ of $l$-mer $x$
- $d$-neighborhood $\mathcal{N}(S, d)$ of sequence $S$

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Definitions of key concepts

- $l$-mer
  - sequence of length $l$

    For $l = 7$,
    $S = $ acgcc**gattaca**tccgatccttgtatagctcctaacgggcatcac
    $\hookrightarrow$ **gattaca** is the $6^{th}$ $l$-mer in $S$.

- Hamming distance $dH(x_1, x_2)$
- $d$-neighborhood $N(x, d)$ of $l$-mer $x$
- $d$-neighborhood $\mathcal{N}(S, d)$ of sequence $S$

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Definitions of key concepts

- $l$-mer
- Hamming distance $dH(x_1, x_2)$
  - number of mismatches between $l$-mers $x_1$ and $x_2$

  $x_1 = $ **ga**tta**c**a
  $x_2 = $ **cg**tta**g**a
  $\hookrightarrow x_1$ and $x_2$ differ in their first, second and sixth characters.
  Thus, $dH(x_1, x_2) = 3$.

- $d$-neighborhood $N(x, d)$ of $l$-mer $x$
- $d$-neighborhood $\mathcal{N}(S, d)$ of sequence $S$

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Definitions of key concepts

- $l$-mer

- Hamming distance $dH(x_1, x_2)$

- $d$-neighborhood $N(x, d)$ of $l$-mer $x$
  - set of all $l$-mers having at most $d$ mismatches with $x$

$$N(\text{gattaca}, 2) = \{ \text{gattaca},$$
$$\text{aattaca, cattaca, tattaca,}$$
$$\text{gcttaca, ggttaca, gtttaca,}$$
$$...,$$
$$\text{acttaca, agttaca, atttaca, ..., tcttaca, tgttaca, ttttaca,}$$
$$\text{aaataca, aactaca, aagtaca, ..., taataca, tactaca, tagtaca,}$$
$$\text{aataaca, aatcaca, aatgaca, ..., tataaca, tatcaca, tatgaca,}$$
$$...,$$
$$\}$$

- $d$-neighborhood $\mathcal{N}(S, d)$ of sequence $S$

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Definitions of key concepts

- $l$-mer
- Hamming distance $dH(x_1, x_2)$
- $d$-neighborhood $N(x, d)$ of $l$-mer $x$
- $d$-neighborhood $\mathcal{N}(S, d)$ of sequence $S$
  - union of $d$-neighborhoods of all $l$-mers in $S$

$S = \text{acgccga}\text{ttacatccgatccttgtatagctcctaacgg}\text{gcatcac}$

$\mathcal{N}(S, 2) = N(\text{acgccga}, 2) \cup N(\text{cgccgat}, 2) \cup ... \cup N(\text{gcatcac}, 2)$
$\hookrightarrow d\text{-neighborhood of first } l\text{-mer in } S$

Introduction
EMS-GT
Methodology
Results
Conclusion

The $(l,d)$ planted motif problem
Definitions of key concepts

# Definitions of key concepts

- $l$-mer
- Hamming distance $dH(x_1, x_2)$
- $d$-neighborhood $N(x, d)$ of $l$-mer $x$
- $d$-neighborhood $\mathcal{N}(S, d)$ of sequence $S$

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# The EMS-GT algorithm

The EMS-GT algorithm proceeds in two main steps:

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# The EMS-GT algorithm

The EMS-GT algorithm proceeds in two main steps:

1. *Generate candidates*

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# The EMS-GT algorithm

The EMS-GT algorithm proceeds in two main steps:

1. *Generate candidates*
   Take the intersection of the $d$-neighborhoods of the first $n'$ sequences $S_1, S_2, ..., S_{n'}$. Every $l$-mer in the resulting set $C$ is a candidate motif.

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# The EMS-GT algorithm

The EMS-GT algorithm proceeds in two main steps:

1. *Generate candidates*
   Take the intersection of the $d$-neighborhoods of the first $n'$
   sequences $S_1, S_2, ..., S_{n'}$. Every $l$-mer in the resulting set $C$ is
   a candidate motif.

2. *Test candidates*

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# The EMS-GT algorithm

The EMS-GT algorithm proceeds in two main steps:

1.  *Generate candidates*
    Take the intersection of the $d$-neighborhoods of the first $n'$
    sequences $S_1, S_2, ..., S_{n'}$. Every $l$-mer in the resulting set $C$ is
    a candidate motif.

2.  *Test candidates*
    For every candidate $c$ in $C$, check whether a $d$-neighbor of $c$
    appears in each of the remaining sequences $S_{n'+1}, S_{n'+2}, ... S_n$.
    If this is the case, accept $c$ as a motif.

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# EMS-GT efficiency strategies

- $l$-mer enumeration
- Bit-based set representation
- Recursive neighborhood generation

# EMS-GT efficiency strategies

- $l$-mer enumeration

- Bit-based set representation
- Recursive neighborhood generation

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# EMS-GT efficiency strategies

- $l$-mer enumeration
- Bit-based set representation

- Recursive neighborhood generation

Introduction
EMS-GT
Methodology
Results
Conclusion

The EMS-GT algorithm
Efficiency strategies

# EMS-GT efficiency strategies

- $l$-mer enumeration
- Bit-based set representation
- Recursive neighborhood generation

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

The main objectives of this research are:

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

The main objectives of this research are:

1. To develop a speedup technique for EMS-GT that takes advantage of distance-related patterns in the search space;

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

The main objectives of this research are:

1. To develop a speedup technique for EMS-GT that takes advantage of distance-related patterns in the search space;

2. To evaluate the speedup technique with regard to improvement in runtime; and

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

The main objectives of this research are:

1. To develop a speedup technique for EMS-GT that takes advantage of distance-related patterns in the search space;

2. To evaluate the speedup technique with regard to improvement in runtime; and

3. To evaluate the improved version of EMS-GT against state-of-the-art motif search algorithms.

# Methodology

- *Developing an EMS-GT speedup technique*

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

- *Developing an EMS-GT speedup technique*
  In EMS-GT, *l*-mer neighborhoods are represented with
  repeating patterns of bits. We exploit these patterns in a
  speedup technique that sets bits quickly and in blocks.

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

- *Developing an EMS-GT speedup technique*
  In EMS-GT, $l$-mer neighborhoods are represented with
  repeating patterns of bits. We exploit these patterns in a
  speedup technique that sets bits quickly and in blocks.

- *Evaluating performance*

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

▶ *Developing an EMS-GT speedup technique*
In EMS-GT, *l*-mer neighborhoods are represented with
repeating patterns of bits. We exploit these patterns in a
speedup technique that sets bits quickly and in blocks.

▶ *Evaluating performance*
Performance is benchmarked on "challenging" $(l,d)$ instances:
(9,2), (11,3), (13,4), (15,5) and (17,6).

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

- *Developing an EMS-GT speedup technique*
  In EMS-GT, *l*-mer neighborhoods are represented with
  repeating patterns of bits. We exploit these patterns in a
  speedup technique that sets bits quickly and in blocks.

- *Evaluating performance*
  Performance is benchmarked on "challenging" $(l,d)$ instances:
  (9,2), (11,3), (13,4), (15,5) and (17,6).

  - *Synthetic datasets*

Introduction
EMS-GT
Methodology
Results
Conclusion

Research objectives
Development and evaluation

# Methodology

- *Developing an EMS-GT speedup technique*
  In EMS-GT, $l$-mer neighborhoods are represented with
  repeating patterns of bits. We exploit these patterns in a
  speedup technique that sets bits quickly and in blocks.

- *Evaluating performance*
  Performance is benchmarked on "challenging" $(l,d)$ instances:
  (9,2), (11,3), (13,4), (15,5) and (17,6).

  - *Synthetic datasets*
    - sets of 20 randomly-generated DNA sequences of length 600,
    with the same $(l,d)$ motif planted once in each sequence.

Introduction
EMS-GT
Methodology
**Results**
Conclusion

Pattern-based speedup technique
Performance

Introduction
EMS-GT
Methodology
**Results**
Conclusion

Pattern-based speedup technique
Performance

# Pattern-based speedup technique

Introduction
EMS-GT
Methodology
**Results**
Conclusion

Pattern-based speedup technique
**Performance**

# Performance improvement with speedup

# Performance vs PMS8, qPMS9

Thank you!