

Graph Based Reference Genomes

Knut Dagestad Rand

April 23, 2019

Thesis submitted for the degree of Philosophiæ Doctor

Contents

1	List of papers	1
2	Introduction	3
3	Background	5
3.1	Molecular Biology	5
3.1.1	DNA	5
3.1.2	Replication and Mutation	6
3.1.3	Epigenomics	6
3.2	Sequencing	7
3.2.1	Illumina Dye Sequencing	7
3.2.2	ChIP-seq	9
3.2.3	Third Generation Sequencing	9
3.3	Reference Genomes	9
3.4	Mathematical Framework	11
3.4.1	Strings	11
3.4.2	Notation	11
3.4.3	Graphs	11
3.5	Alignment	12
3.5.1	Edit Distance	12
3.5.2	Pairwise Sequence Alignment	12
3.5.3	Sequence Graph Alignment	16
3.6	Mapping	18
3.6.1	Linear Mapping	19
3.6.2	Graph Mapping	22
3.7	Peak Calling	23

4	Summary of Papers	27
4.0.1	Paper I: Coordinates and Intervals in Graph Based Reference Genomes	27
4.0.2	Paper II: Graph Peak Caller: Calling ChIP-seq peaks on graph-based reference genomes	27
4.0.3	Paper III: Assessing graph-based read mappers against a novel baseline approach highlights strengths and weaknesses of the current generation of methods	27
5	Discussion	29
5.0.4	Computational Complexity	29
5.0.5	Invalid Sequences	30

Chapter 1

List of papers

Chapter 2

Introduction

The broad topic of this thesis graph based reference genomes, and specifically of mapping and peak calling on such references. The introduction here is meant to provide an introduction to the topic as well as a coherent brief of the content of the provided papers. As is common in bioinformatics, full understanding of the field of graph based references requires at least a rudimentary understanding of the underlying biology, mathematics and informatics. An brief introduction to the specific subtopics

The first section will be an introduction to two main biological topics, reference genomes and genomic variation, in order to familiarize the reader with topics in addition to establish a perspective that establishes the benefits of graph based references. Section 2 introduces the concept of sequence graphs and their interpretation in different settings. Section 3 describes the mathematical formalism used in this thesis as well as the data structures embodying this formalism. These sections might be persived as too rigorous and opaque for a brief introduction, but this highlights one of the main drawbacks of graphical models, namely their complexity. Section 4 and 5 gets to the main focus of this thesis, covering the topics of read mapping and peak calling on graph based reference genomes.

Chapter 3

Background

3.1 Molecular Biology

3.1.1 DNA

At the heart of molecular biology is *deoxyribonucleic acid*, *DNA*. DNA-molecules consists of pairs of four different *nucleotides* linked together to form a double stranded helix. The four different nucleotides are often represented by the letters A (*adenin*), C (*cytosin*), T (*thymine*) and G (*guanin*), forming an alphabet of nucleotides $\Gamma_{DNA} = \{A, C, T, G\}$. Each nucleotide can only be paired with one of the others, such that each base-pair in the helix is either (A-T) or (G-C) or opposite. (*chromosomes*)

The most important function of DNA is to serve as a template for the synthesis of proteins, which in turn are responsible for a wide range of bio-molecular functions. This is a two-step process where DNA is first *transcribed* to *RNA*, molecules similar to DNA but single stranded and with *thymine* replaced by *uracil* (U). The resulting RNA-sequence can in turn be translated to proteins by mapping triplets of ribonucleotides to amino acids, which are the building blocks of proteins.

Sequences of DNA that are transcribed to RNA that either has a function in itself, or is in turn translated to protein, are called *genes*. Genes constitute only a minor proportion of human DNA. The remaining DNA can be important due to the biochemical properties of the DNA itself.

3.1.2 Replication and Mutation

In addition to serving as template for RNA molecules, DNA can also be replicated. This leads to inheritance both on cellular and organismal level. In normal cell division, the chromosomes are replicated such that each of the two resulting cells have a copy of the DNA from the original cell. Furthermore, germ line cells are copied and transferes half of the organisms chromosomes to an offspring.

However, this DNA-replication is not always accurate. Errors can occur leading to a copy of the DNA molecule which is not identical with the original. Most common are the substitution of a single nucleotide, insertion of a small dna sequence, or deletion of a small subsequence. Such errors lead to difference between cells within an organism, or difference between individuals.

Difference in the DNA-sequence between two individuals can lead to a change in the transcribed RNA sequence and further in the sequence, and thereby the form and function, of the expressed protein. Or it can lead to less drastic changes such as changes in the RNA-structure or the shape of the DNA molecule itself. In this way genomic variation can determine differences between specimens of the same species and also differences between the species themselves.

3.1.3 Epigenomics

Difference in DNA-sequence can explain phenotypic (*explain geno/pheno*) difference between individuals, but fails to explain the difference between the different cells within an organism. Within an organism, the cells contain mostly the same DNA-sequences, but exhibit vastly different phenotypes.

The reason for this difference is fundamentally mostly attributed to differences in expression levels of genes. Much attention has been devoted in recent years to the mechanisms determining gene expression levels. Two high-level consortia, ENCODE ?? and Roadmap Epigenomics ??, have systematically conducted experiments geared to understanding some of these factors, including: 3D-configuration of the chromosomes, accessibility of the DNA in different regions, transcription factor binding sites, and methylation patterns across the genome. Of these, this thesis is concentrated on transcription factor binding sites, but for an introduction to epigenomics as a whole see ??.

Transcription factors are proteins that bind to DNA, with an affinity for certain DNA-sequences. These proteins functions as signals to other proteins that a gene in it's vicinity should be transcribed, thus affecting expression levels of genes. As such the precence or absence of transcripion factors binding to a binding site can

affect the phenotype of a cell. Also since the binding of transcription factors are dependent on the DNA-sequence, a mutation in the DNA-sequence in a binding site can affect the binding of the protein and thus the expression of a nearby gene.

3.2 Sequencing

The fact that the sequence of nucleotides in the DNA has direct links to biological function have made determining such sequences as good as possible an important goal in molecular biology. The process of determining the sequence of nucleotides is referred to as *DNA-sequencing*. The first widely used method for DNA-sequencing was Sanger Sequencing [?] in (year). The throughput from Sanger sequencing was pretty low since it required measuring the weight of many nucleotides. In the beginning of the 2000s, several technologies were developed that was able to parallelize the process of sequencing and thus increased the throughput. Illumina(trademark?) dye sequencing is the most commonly used, and will be described here.

3.2.1 Illumina Dye Sequencing

Illumina dye sequencing uses sequencing by synthesis, which is the same base methodology as Sanger sequencing. This uses the property that given a single strand of DNA and the presence of DNA polymerase, nucleotides will sequentially bind to the strand according to the complementary pairings (A-T) and (C-G). This usually happens so fast that it is difficult to measure any point in this process. However nucleotides can be adapted so that the synthesis is terminated after including the adapted nucleotide (chain termination). Sanger sequencing uses this by combining a small amount of adapted nucleotides with a larger amount of normal nucleotides, and then using this combined set for DNA synthesis. Illumina sequencing evolves this concept by also being able to reverse these adaptations so that synthesis can continue. Thus it can step for step synthesize one adapted nucleotide to the single DNA-strand, find out which nucleotide was included, then reverse the adaption and synthesise a new adapted nucleotide(see figure 3.1.

The big breakthrough for this sequencing technology is that it could be performed massively in parallel and thereby generate much more sequencing data. The downside is that while Sanger sequencing could yield up to 700 base pairs long sequences, Illumina sequencing typically yields shorter reads (25-200)(accurate numbers). This massive amount of relatively short reads have are able to give great

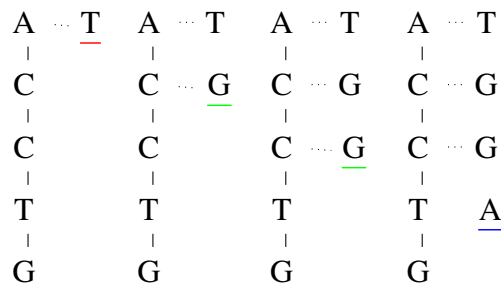


Figure 3.1: Sequencing of one DNA-chain in Illumina sequencing. For each step, an adapted nucleotide is added with a phosphorescent color, a picture is taken - capturing the color, and the adaption is reversed allowing for the next nucleotide to be bound.

insights into biology by mapping them to reference genomes. Reference genomes are covered in section ??, and the mapping process is covered in 3.6.

Illumina sequencing main error source is that a non-adapted nucleotide can get included in one of the steps, whereupon it is possible for a new nucleotide to get included without the previous one being captured. To alleviate this, and to obtain a stronger signal, Illumina performs the sequencing steps concurrently on many identical DNA templates. Even so, the number of out-of-sync strands grows with the length of the sequence and puts a practical limit on how many base pairs can be read with certainty.

Depending on which input DNA is provided, this technology can be used to answer a range of biological questions. The simplest(*other word*) application, is to input unfiltered DNA for sequencing and using the resulting reads to deduce as much of the whole genomic sequence as possible. One can also filter the DNA to contain mostly fragments from protein coding parts of the genome, in order to more efficiently find the DNA-sequence of these regions, from which variations can cause differences in protein structure and function.

In addition to deducing the DNA-sequence of a sample, NGS can also be used to deduce where in the genome biochemical processes occur, for instance the binding of proteins to the genomes. Of importance to this thesis are experiments used to find binding sites for transcription factors, described below.

3.2.2 ChIP-seq

Chromatin Immunoprecipitation followed by *sequencing* (ChIP-seq) is a sequencing experiment that does not aim to determine the genetic variation of the sample, but rather to give insight to where a protein of interest binds to the DNA.

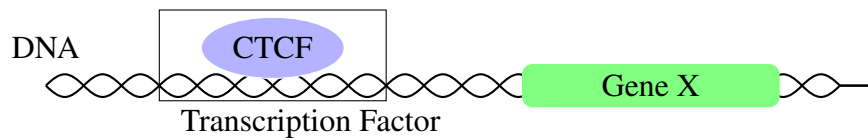
In order to do this, the binding of the protein to DNA is strengthened in order to keep the bindings after the DNA is cut into fragments. Antibodies for the protein of interest is then used to retrieve DNA-fragments with the protein attached. These DNA-fragments can then be sequenced, yielding reads that tend to come from the vicinity of a binding site (figure 3.2). The process of analyzing these reads is covered in section 3.7.

3.2.3 Third Generation Sequencing

New technologies are continuously being developed, and a set of more recent technologies allows for sequencing of long DNA-molecules. This does however come at the price of much higher error rates than Illumina sequencing. Most notable is Oxford Nanopore [?], passing DNA-molecules through a molecule which emits a different electric current for each type of nucleotide, and PacBio [?], using nucleotide-dependent light waves emitted continuously while synthesising. These techniques offer possibilities for answering biological questions where short read lengths are insufficient, and also new computational challenges for dealing with the high error rates and long read lengths.


3.3 Reference Genomes

Since it can be assumed that the genomic DNA-sequences from individuals of the same population is highly similar, it is possible to use one DNA-sequence as a *reference genome* for the population. Such reference genomes serves two main purposes. They make the deduction of the DNA-sequence from a sample a simpler problem through mapping, and they make it possible to represent the outcome of a range of biological experiments as intervals or coordinates on the coordinate system induced by the reference genome. The process of mapping is covered in section 3.6, while the representation of different experiment outputs are illustrated in figure ??.





 ACGTTCGTATATCGTAGCTACTCGAGCTGTAGTTTGATAGATAT

(a) Transcription factor binding to DNA, regulating the expression of Gene X



Cell1:	ACGTTCGTATATCGTAGCTACTCGAGCTGTAGTTTGATAGATAT
Cell2:	ACGTTCGTATATCGTAGCTACTCG.....
Cell3:TATATCGTAGCTACTCGAGCTGTA.....
TCGTAGCTACTCGAGCTGTAGTTT.....
	..GTTCGTATATCGTAGCTACTCGAG.....
	...TTCGTATATCGTAGCTACTCGAGC.....
CTACTCGAGCTGTAGTTTGATAGA...
ACTCGAGCTGTAGTTTGATAGATA.
GTAGCTACTCGAGCTGTAGTTTGA.....

(b) DNA Fragments obtained from ChIP different cells



ACGTTCGTATATCGTAGCTACTCGAGCTGTAGTTTGATAGATAT
 ACGTTCGTATATCGTAGCTACTCG.....
TATATCGTAGCTACTCGAGCTGTA.....
TCGTAGCTACTCGAGCTGTAGTTT.....
 ..GTTCGTATATCGTAGCTACTCGAG.....
 ...TTCGTATATCGTAGCTACTCGAGC.....
CTACTCGAGCTGTAGTTTGATAGA...
ACTCGAGCTGTAGTTTGATAGATA..
GTAGCTACTCGAGCTGTAGTTTGA.....

↓

ACGTTCGTAT, TACAGCTCGA, TCGTAGCTAC, GTTCGTATAT
 GCTCGAGTAG, CTACTCGAGC, TATCTATCAA, TCAAACATA

(c) One end of each fragment is sequenced

Figure 3.2: Figure showing the main steps of a ChIP-seq experiment. A protein of interest (CTCF) binds to the DNA, the DNA is cut up into fragments, from which the ones with proteins bound to it are kept. These fragments are then sequenced, yielding reads that are from the vicinity of protein binding sites.

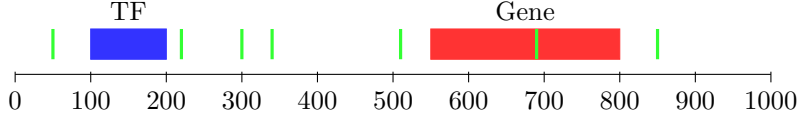


Figure 3.3: Figure showing the locations of several genomic elements on the same coordinate system. Locations of genes (red), transcription factor binding sites (blue), and methylation sites (green) from different sequencing experiments can all be compared and analyzed together.

3.4 Mathematical Framework

3.4.1 Strings

3.4.2 Notation

Strings are important in this thesis as they are the means to represent DNA-sequences. We will here formally work with an alphabet $\Gamma_{DNA} = \{A, C, G, T\}$ and say that a *sequence* of length n over that alphabet are a tuple in the set Γ_{DNA}^n . The set of all sequences over $\Gamma_{DNA} = \{A, C, G, T\}$ is represented by Γ_{DNA}^* . The length of a sequence $s \in \Gamma_{DNA}^*$ is notated $|s|$. The i th symbol in a sequence s is notated $s[i]$ and a subsequence of s starting at the i th symbol (inclusive) and ending at the j th symbol (exclusive) is denoted $s[i : j]$. A *prefix* of s , i.e. a substring of s starting from the first symbol, of length k is denoted $s[: k]$, while a suffix of s starting at the k th symbol is denoted $s[k :]$. For convenience a string of length n is sometimes represented as $s[: n]$ to convey the size of the string. Concatenation of two strings S, T are represented as $S * T$, while the concatenation of a symbol a to a string S is denoted Sa , and a string to a symbol as aS .

3.4.3 Graphs

A graph is a tuple $G = (V, E)$ of vertices and edges, where the edges are pairs of vertices $E \subset V^2$. We will here deal with directed graphs where we say that edge (v_1, v_2) is an edge from v_1 to v_2 . An *path* is an alternating sequence of vertices of edges $(v_1, e_1, v_2, e_2, \dots, e_{n-1}, v_n)$ where $e_i = (v_i, v_{i+1})$. A cycle is a path $(v_1, e_1, v_2, e_2, \dots, e_{n-1}, v_n)$ where $v_1 = v_n$ and $n > 1$, i.e. a path that starts and ends at the same vertex. A directed acyclic graph (DAG) is a directed graph in which there are no cycles. We call the set of all paths starting at v_s and ending at v_e paths_{v_s, v_e} .

Sequence Graphs

We define a *simple sequence graph* $SG = \{G, s\}$ over an alphabet Γ as a graph $G = \{V, E\}$ and a label function $\text{label} : V \rightarrow \Gamma$ labeling each vertex with a letter from Γ . We refer to the label of a path $\text{label}((v_1, e_1, v_2, \dots, e_{n-1}, v_n))$ as the concatenation of the labels of its vertices $\text{label}(v_1) * \text{label}(v_2) \dots \text{label}(v_n)$. The *language* recognized by a sequence graph $L((G, \text{label}))$ is the label of each path in the set $\text{paths}(v_s, v_e)$.

3.5 Alignment

3.5.1 Edit Distance

An edit distance is a measure of How many mutations are required to convert one sequence into another. Different edit distance measures exists, varying in the set of allowed mutations. The most common is the *Levenshtein distance* [?], allowing single character substitutions, insertions and deletions.

Finding the edit distance between two sequences, and which set of edits this edit distance corresponds to, is of central importance in bioinformatics since it can give an estimate of how related the two sequences are and which mutations have separated them. The process of finding these estimates is called *sequence alignment* and is important for this thesis in two respects. Firstly, it is one of the earliest and most intuitive applications of sequence graphs, and secondly it is an important part of *mapping* which will be covered in the next section.

In the following we will refer to the edit distance between two sequences, S and T , as $D(S, T)$, meaning the Levenshtein distance unless specifically mentioned. The concepts discussed are however generalisable to other (weighted) edit distances by small changes.

The set of edits between two sequences can be represented with an *alignment block* where “-” symbols are inserted into the sequences to represent indels, (see figure ??).

3.5.2 Pairwise Sequence Alignment

Finding the edit distance between two strings is easy if indels are not considered. It is merely the process of counting the number of mismatches between them. Indels introduce a dependency since an insertion at position i affects the pairings

of all the symbols after position i . The number of meaningful combinations of insertions and deletions grow exponentially with sequence length, so exploring all of them is not a viable solution even for short sequences. The problem is however tractable since the best alignment of two sequences $S[: m], T[: n]$ is either the best alignment of $S[: m - 1], T[: n]$ with an insertion at the end, or of $S[: m], T[: n - 1]$ with a deletion at the end, or of $T[: n - 1], T[: m - 1]$ with a match or substitution at the end. Letting $d_{i,j} = D(S[: i], T[: j])$ and $m_{i,j}$ be zero if $S[i] = T[j]$ else 1, we can write this as the recurrence relation:

$$d_{k,l} = \min \begin{cases} d_{k-1,l-1} + m_{k-1,l-1} & \text{(match/substitution)} \\ d_{k-1,l} + 1 & \text{(insertion)} \\ d_{k,l-1} + 1 & \text{(deletion)} \end{cases}$$

The fact that the edit distance to a empty string is just the sequence length gives $d_{k,0} = k, d_{0,l} = l$. The Needleman-Wunch algorithm uses dynamic programming to calculate the matrix of edit distances d_{kl} , where d_{mn} is the edit distance between S and T . In order to find the specific edits, a backtracking algorithm is used that starts at the $(i, j) = (m, n)$ corner and finds out which of the possible predecessors $(i - 1, j), (i - 1, j - 1), (i, j - 1)$ contributed to the (i, j) edit distance. Then repeating this until the $(0, 0)$ corner is reached. Each of these steps correspond to a column in the alignment block. Figure ?? details the Needleman Wunch algorithm. This algorithm can be adopted in a number of ways in order to solve related problems. Notably, by using an affine gap penalty one can reduce the cost of long indels compared to many small indels [?], or one can use positive scores for matches to be able to find subsequences that align well to each other [?].

An adaption relevant for this thesis, is to find the a subsequence of one sequence R that minimizes the edit distance to another Q . I.e find $\min_{k,l}(D(R[k, l], Q))$. This can be done by changing just the initial conditions of the Needleman-Wunch algorithm, to remove the cost of gaps at the beginning and end of R . We set $d_{k0} = 0$ and start the backtracking algorithm at (m, l) where $j = \operatorname{argmin}_i d_{mi}$. Figure?? shows an example of this algorithm. This algorithm provides three results: the coordinate k, l of the subsequence of R most similar to Q , the edit distance from that subsequence to Q , and the set of edits contributing to the edit distance. In this way it solves in a exact, but slow way the problem of the next chapter, namely mapping a read Q to a reference genome R .

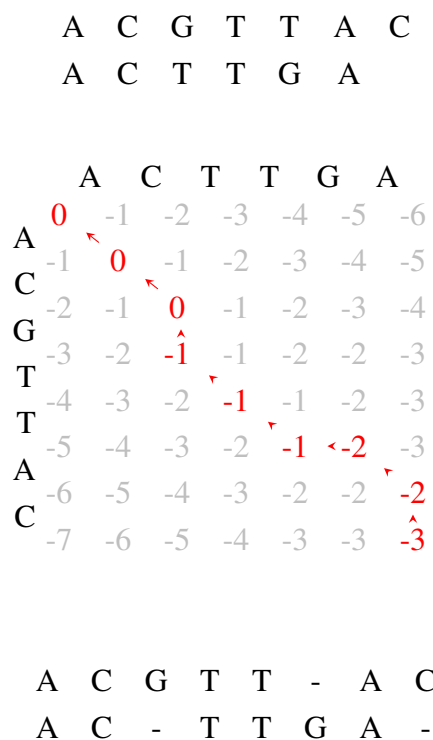


Figure 3.4: Figure showing the alignment of two sequences using Needleman-Wunch. Each cell in the matrix corresponds to the edit distance between prefixes of S and T . The red path shows the backtracking, resulting in an alignment block where each diagonal arrow gives a symbol from both sequences, while horizontal or vertical arrows give an insertion or deletion.

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

	A	C	T	T	G	A
G	0	-1	-2	-3	-4	-5
T	0	-1	-2	-3	-4	-5
G	0	-1	-2	-2	-3	-4
G	0	-1	-2	-3	-3	-4
A	0	-1	-2	-3	-4	-3
A	0	0	-1	-2	-3	-4
C	0	-1	0	-1	-2	-3
G	0	-1	-1	-1	-2	-3
T	0	-1	-2	-1	-2	-3
T	0	-1	-2	-2	-1	-2
A	0	-1	-2	-2	-2	-3
C	0	0	-1	-2	-2	-2
C	0	-1	0	-1	-2	-3
G	0	-1	-1	-1	-2	-3
C	0	-1	-1	-2	-2	-3
G	0	-1	-2	-2	-3	-2
G	0	-1	-2	-3	-3	-3
T	0	-1	-2	-2	-3	-4

A C G T T - A
A C - T T G A

Figure 3.5: Figure showing the alignment of a query sequence Q to a reference sequence R . Each cell in the matrix corresponds to the edit distance between a prefix of S and a subsequence of R . The red path shows the backtracking, starting at the lowest value of the last column, resulting in an alignment block where each diagonal arrow gives a symbol from both sequences, while horizontal or vertical arrows give an insertion or deletion.

3.5.3 Sequence Graph Alignment

The framework used to align sequences extends naturally to acyclic sequence graphs [?, ?]. We define the alignment of a sequence S to a sequence graph G as the alignment of S to a sequence $T \in L(G)$ in the language of G which yields the lowest edit distance. Similarly, the alignment of two sequence graphs G, H , is the alignment of a sequence $S \in L(G)$ to a sequence $T \in L(H)$ that yields the lowest edit distance.

If we let $d_{ij} = \min_{p_g \in \text{paths}(v_0, v_i), p_h \in \text{paths}(w_0, w_k)} (D(\text{label}(p_h), \text{label}(p_g)))$ we get a recurrence relation:

$$d_{ij} = \min_{v_k \in e^- v_i, w_l \in e^- w_j} d^*(i, j, k, l)$$

$$d^*(i, j, k, l) = \min \begin{cases} d_{il} + 1 \\ d_{kj} + 1 \\ d_{kl} + m_{kl} \end{cases}$$

$$m_{ij} = 1 \text{ if } \text{label}(v_i) \neq \text{label}(w_j) \text{ else } 0$$

This is the same as for ordinary sequence alignment, except that all predecessor nodes of v_i and w_k have to be considered, not only $i - 1$ and $j - 1$ as in the linear case. If the graph is acyclic, then all the d_{ij} can be calculated using dynamic programming, without incurring any infinite loops.

An alignment between two sequences S, T can be represented by a sequence graph in a meaningful manner in that one can construct a sequence graph $AG(S, T)$ from the alignment of the sequences in such a way that

$$\forall (R \in L(AG(S, T))) [D(S, R) + D(R, T) = D(S, T)]$$

. For an optimal alignment, this means that all sequences recognized by the sequence graph have the property that the sum of the distance to the original sequences is as low as it can be. These sequences thus represents combinations of S and T that are natural estimates of an ancestor of the two sequences. Thus aligning a sequence S to an alignment graph $AG(T, R)$ can be seen as aligning S to the best fitting ancestor sequence of T and R . Similarly, aligning two alignment graphs can be seen as finding the ancestors of two pairs of sequences that fits best together.[?]. This is illustrated in figure 3.7.

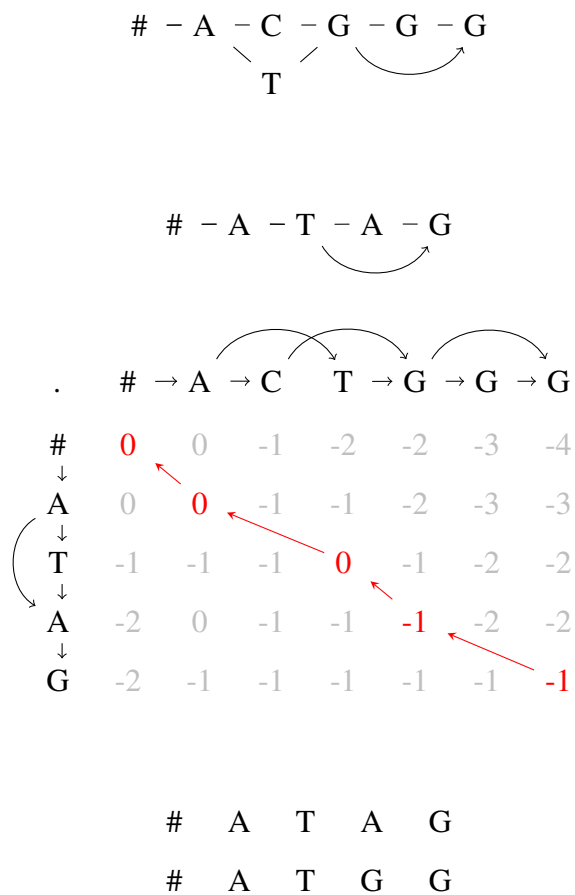


Figure 3.6: Figure showing the alignment of two sequence graphs. Marked in red is the path taken during backtracking. The result is the best alignment of a sequence from the language of each sequence graph. This alignment can again be represented as a sequence graph

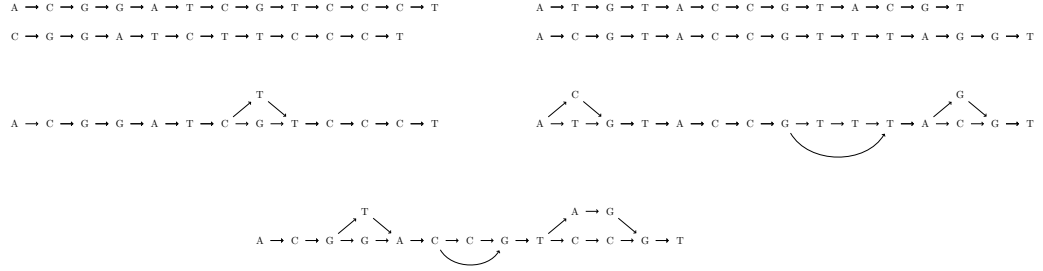


Figure 3.7: Iterative sequence graph alignments of four sequences evolved in two generations from *ACGTACGTACGT*. The sequences are represented as linear sequence graphs (a) and pairwise alignment is performed on the closest pairs yielding two sequence graphs (b). These two sequence graphs are then aligned to each other, yielding a sequence graph representing an alignment of the two closest paths in the graphs (c). As seen the original sequence is in the language of the final sequence graph

The sequence graph alignment algorithm can also be adapted to find the subsequence of a sequence graph G_R that minimizes the edit distance to a sequence Q . This solves the problem of mapping a read to a graph, which will be discussed further in section ??

3.6 Mapping

The alignment methods described in the previous section do not scale well to large sequences. Aligning a read $Q[: m]$ to a reference genome $R[: N]$ will require computing mN values. A typical NGS sequencing experiment typically yields hundreds of thousands of reads of length 100 which needs to be located in the human reference genome of length 3 billions. This would need to calculate 10^{16} values, which is impossible? even on modern computers.

This has led to much development the last 13 years in developing alignment methods that avoid the complexity of the exact dynamic programming methods. This has been achieved mainly by two measures, creating searchable indexes of the reference genome and using heuristics to limit the search space of possible alignments.

In the following is a brief description of the developments in this field, with focus on the aligner used in this thesis' projects: BWA-mem [?]. After this we will look at the methods used in for graph alignment with focus on *vg*.

3.6.1 Linear Mapping

As in section ?? (*Define this explicitly in alignmetn*), we define the mapping problem as finding the interval i_R for which the edit distance $D(R[i_r], Q)$ is smallest, and the alignment $A(R[i_R], Q)$.

Indexing

An essential tool for quick alignment is to have a searchable index of the reference genome which is capable of returning sets of interval-pairs from the reference and query sequence for which the reference sequence is identical or near identical to the query sequence. The indexes used by different tools vary, but can mainly be divided into fixed-length(kmer) indexes [?, ?, ?, ?] and variable-length indexes [?, ?, ?, ?]. Fixed length indexes typically uses hash tables to store the location of all, or subsets of all, kmers (*define kmer*) in the reference sequence. Variable length indexes usually uses a variation of the *full text minute index* (FM-index) [?], described below. The main idea covered here is the seed-and-extend paradigm. This idea consists of finding a set of exact matches to the reference sequence, and using these as anchors for using dynamic programmic based alignment of the query. We let $EM(Q, R)$ be the set of all exact matches, represented by the tuples $\{(q, r, l) \mid Q[q : q + l] = R[r : r + 1]\}$. Since the dynamic programming methods can still be computationally expensive, the goal is to make the seeds as small a subset of exact matches as possible, while still making sure that the true alignment covers one of the seeds.

FM-index

The FM-index uses a succinct representation of suffix array [?] and Burrows-Wheeler transform [?] in order to find exact matches of a query string S in $O(|S|)$ time (see figure ??). Thus if a read has one or more exact matches in the reference sequence it can be found directly using the FM-index. For inexact matches, it can also be used to find an exact alignment of Q and R in $O(|R|^{0.628} |Q|)$ time [?, ?]. Some algorithms also use the FM-index to find exact matches for permutations of the query sequence [?, ?], but for longer reads the search space gets too big when allowing for indels.

In order to handle indels in longer reads, the main methodology is to use the FM-index to find exact matches between substrings of the query and reference sequence, called seeds. And then aligning the reads using dynmaic programing to

C	T	A	G	C	T	G	C	A	T	G
0	1	2	3	4	5	6	7	8	9	10

11	\$	CTAGCTGCAT	G_0	0	0	0	0
2	A_0	GCTGCATG\$C	T_0	0	0	1	0
8	A_1	TG\$CTAGCTG	C_0	0	0	1	1
7	C_0	ATG\$CTAGCT	G_1	0	1	1	1
0	C_1	TAGCTGCATG	\$	0	1	2	1
4	C_2	TGCATG\$CTA	G_2	0	1	2	1
10	G_0	\$CTAGCTGCA	T_1	0	1	3	1
6	G_1	CATG\$CTAGC	T_2	0	1	3	2
3	G_2	CTGCA TG\$CT	A_0	0	1	3	3
1	T_0	AGCTGCATG\$	C_1	1	1	3	3
9	T_1	G\$CTAGCTGC	A_1	1	2	3	3
5	T_2	GCATG\$CTAG	C_2	2	2	3	3

C T G

Figure 3.8: Illustration of backward extension using the last-first (LF) property of the FM index. The rows represent sorted suffixes of the reference. The SA column holds the indices in the reference sequence for each suffix, the F column holds the first character of each suffix, while the L column holds the preceding character of each suffix. Subfixes gives the occurrences of each character in each column. I.e. T_i is the i th occurrence of T in that column. Finding the string CTG is done by (1) starting with the range of all G 's in the F column; (2) finding all T 's in this range in the L column; (3) mapping by occurrence number those T 's to the F column; (4) mapping the C 's in the current range in the L folder by occurrence number. The end result is in this case a single row which represents position 4 in the reference sequences.

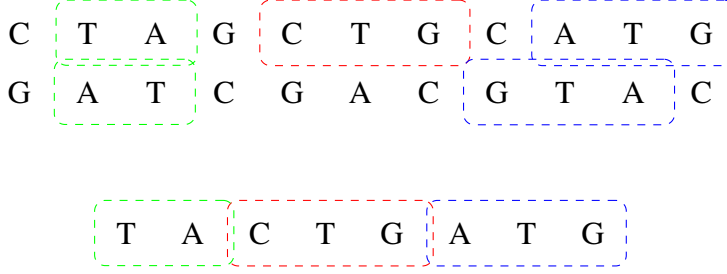


Figure 3.9: SMEMs found between query Q and both strands of reference R . Note that the MEM $(q, r, l) = (2, 0, 2)$ (CT) is not a SMEM, since it is contained in the SMEM $(q, r, l) = (2, 4, 3)$

intervals surrounding the seed matches on the reference sequence [?], often called the seed-and-extend paradigm.

One way of finding seeds is to find Maximal Exact Matches (MEM) [?, ?]. Maximal Exact Matches are Exact Matches that cannot be extended in either direction, ie.

$$MEM(Q, R) = \{(q, r, l) \in EM \mid (q, r, l+1) \notin EM \wedge (q-1, r-1, l+1) \notin EM\}$$

BWA-mem furthers this concept to SuperMaximal Exact Matches (SMEM) [?]. A SMEM is a MEM where the query interval cannot be extended either side and still yield an Exact Match ??.

$$SMEM(Q, R) = \{(q, r, l) \in EM \mid \nexists r^* [(q, r^*, l+1) \notin EM \vee (q-1, r^*-1, l+1) \notin EM]\}$$

SMEMS are natural to use as seeds as they cover for each subsequence in Q the longest exact match in R . They are however vulnerable for spurious long matches hiding shorter exact matches. To account for this BWA-mem allows an option to split long SMEMs into shorter MEMs if they are longer than a certain threshold. Splitting SMEMs like this increases accuracy, since it increases the number of seeds, but can negatively affect performance. In order to find SMEMS, an adaption of the FM index is used, the FMD index, where $FMD(R) = FM(R * \bar{R})$.

Prioritizing and Extending

The seeds found from the index are next used as seeds for DP based alignment. Q is then aligned against an area around $(r, r + l)$ for each seed (q, r, l) . Since this is a computationally expensive step, further limiting the set of seeds is advantageous. BWA-mem does this by *chaining* the seeds. This is done by grouping approximately colinear, nearby seeds into chains, and removing small chains that overlap with larger chains. Approximately colinear means that $|(q_1 - q_2) - (r_1 - r_2)| < w$ for some set threshold.

3.6.2 Graph Mapping

Mapping to a graph based reference is similar to the linear case, except that instead of finding a linear interval, the goal is to find a graph interval i_r such that $D(\text{label}(G_R(i_r), Q))$ is minimized.

It is however deceptively more complicated. Firstly because the number of subsequences in a sequence graph grows exponentially with the complexity of the graph. And secondly since chaining subsequence-matches is more complicated due to the possible existence of multiple paths between two matches. *vg* [?] has been at the forefront of mapping to a graph reference, showing that it can lead to better mapping accuracy than BWA-mem. Below is a brief description of the methodology used by *vg*.

vg

vg uses much the same methodology as BWA-mem to align reads. It uses the GCSA2-index to find SMEMs, uses chain filtering to filter the SMEMs, and uses the graph adaption of Smith-Waterman to extend the seeds. *vg* is able to align reads to more complicated graph structures than the simple directed sequence graphs considered in this thesis. For simplicity the descriptions below will be contained to simple graphs, which entails that the GCSA index described is GCSA1[?] which is only able to index directed sequence graphs.

GCSA

The GCSA-index [?, ?] is a generalization of the FM-index, where arbitrary-length sequences can be looked up in a sequence graph. Originally constructed to work on acyclic sequence graphs, *gcsa2* extends the functionality to general

variation graphs. For simplicity we will here constrain the discussion to acyclic graphs.

GCSA uses the same concept of LF mapping as the FM-index does. Here the F column holds the label of each node in the graph, sorted by the suffix starting from that node. The L column holds the label of the corresponding node's predecessor nodes. The problem with this setup is that there are many suffixes starting from each node, depending on which path is taken in the graph. To resolve this problem, the graph needs to be expanded so that for each node, all suffixes starting from that node shares a prefix that is not found from any other node in the graph(see figure ??).

For complicated regions in the graph, this expansion procedure gets too costly, and the GCSA index therefore needs to prune edges in such areas in order to be able to index it. This means that not all possible sequences in the graph gets indexed. Even after pruning, this procedure is costly and makes the GCSA index significantly slower to create, and more memory consuming than the FM-index.

Since the gcsa index provides the same functionality as an FM-index, it can be used to find SMEMs in much the same manner. These SMEMs are in *vg* used as seeds.

GCSA2 solves the problem in a slightly different way that allows for indexing general variation graphs. This is based on succinct de-Bruijn graph [?, ?] structure which sets a limit to the length of the unique prefixes by allowing for false positive edges to occur in the index. This means that results from an index lookup needs to be validated by traversing the graph.

Filtering and Extending

vg employs a similar method as BWA-mem for filtering out seeds: chaining the seeds and filtering out small chains that overlap bigger chains. The chaining procedure is however more complicated on a graph, and involves clustering the seeds by position in addition to a Markov Model based method to find approximately colinear seeds.

3.7 Peak Calling

Peak Calling is an important step in the ChIP-seq experiment pipeline. The sequencing in a ChIP-seq experiments yields reads that tend to come from DNA-fragments with the transcription factor bound to it. After mapping these reads to

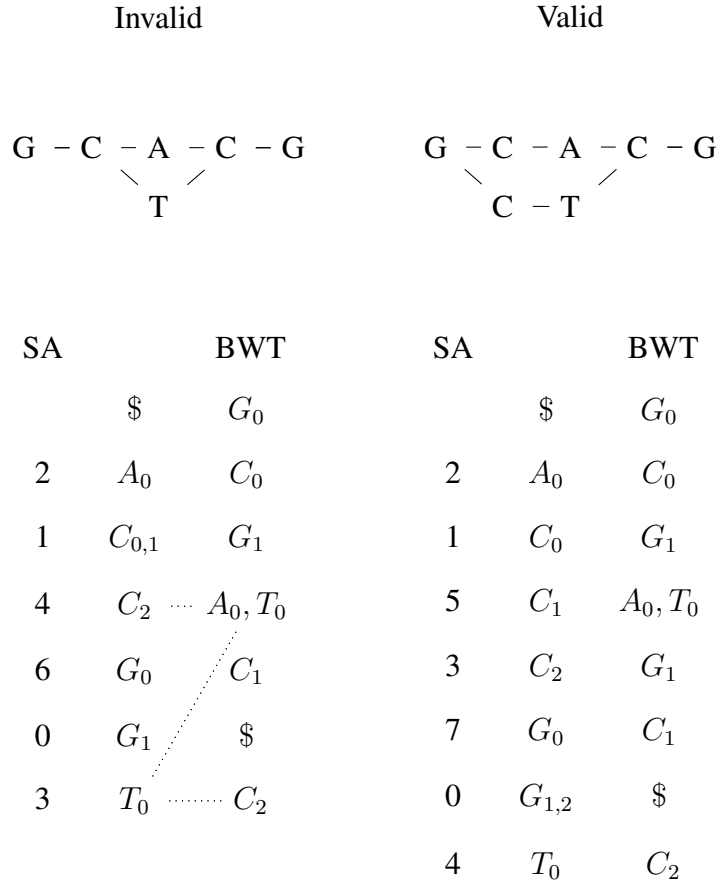


Figure 3.10: Valid and invalid GCSA index. In the first graph, the two edges from C to A and T makes a direct GCSA index impossible due the other C with an edge to a G . Both T and G needs to point to the same row in the index. In the valid example, the C node have been duplicated. The bifurcation now happens at an earlier point, and the GC now present is unique so that the the LF mapping is valid.

the reference genome, the result is a set of genomic intervals which should be in the vicinity of the binding site. Peak calling is the process of using these intervals to predict the actual binding sites.

Calling peaks faces two main problems, noise and bias. Noise can come from either of the previous steps of the experiment: isolating DNA fragments with the TF bound to it, sequencing the fragments, or mapping the resulting reads to the reference genome. The bias can be either biological, in that certain regions of the genome are over/under-represented in the sample independently of the TF, or from the mapping, in that certain regions of the reference genome are more or less likely to be mapped to. In addition to this, the sequenced reads do not always contain the actual binding site, but rather comes from the beginning of a fragment containing the binding site.

In light of this, peak callers seek to recognize the pattern expected from a true binding site: that the mapped intervals should generally be close to and point towards the true binding site. Several algorithms exist for this [?, ?, ?, ?], of which MACS2 is one of the most widely used. As it forms the basis for Graph Peak Caller (paper 2), the algorithm is described in some detail in figure 3.11:

```

Reads:      ACGTTCGTAT, TACAGCTCGA, TCGTAGCTAC, GTTCGTATAT
            GCTCGAGTAG , CTACTCGAGC, TATCTATCAA, TCAAACATACA

Mapped:      ACGTTCGTATATCGTAGCTACTCGAGCTGTAGTTTGATAGATAT
            [----->.....<-----] .....
            .....[----->.....<-----] .....
            ..[----->.....<-----] .....
            .....[----->.....<-----] .....
            .....[----->.....<-----] .....
            .....[----->.....<-----] .....
            .....[----->.....<-----] .....
            .....[----->.....<-----] .....

            (0, 10, +), (21, 31, -), (7, 17, +), (2, 12, +)
            (17, 27, -), (17, 27, +), (33, 43, -), (27, 37, -)

Extended:    [-----+++++-----] .....
            .....[+++++-----] .....
            .....[-----+++++-----] .....
            ..[-----+++++-----] .....
            ...[+++++-----] .....
            .....[-----+++++-----] .....
            .....[+++++-----] .....
            .....[-----+++++-----] .....
            .....[+++++-----] .....

            (0, 24), (7, 31), (7, 31), (2, 26)
            (3, 27), (17, 41), (19, 43), (13, 37)

Pileup:      .....-----.....
            .....-----.....
            .....-----.....
            .....-----.....
            .....-----.....
            .....-----.....
            .....-----.....
            .....-----.....

Peak:        .....-----.....

```

Figure 3.11: Illustration of the main MACS2 algorithm. After reads are mapped to the reference genome, all mapped intervals are extended to match the estimated fragment length. A pileup is then created based on the coverage of each position, and positions with a high enough coverage are marked as reads

Chapter 4

Summary of Papers

4.0.1 Paper I: Coordinates and Intervals in Graph Based Reference Genomes

In paper I, we discussed the implications of using a graphical reference structure for representing and comparing genomic locations and intervals. The work focused on how to obtain succinct and interpretable representations that were robust against changes to the reference graph topology.

4.0.2 Paper II: Graph Peak Caller: Calling ChIP-seq peaks on graph-based reference genomes

Paper III developed a new method of calling ChIP-seq peaks using a sequence graph as reference. The method is an adaptation of MACS2 that uses the benefits of read-mapping to a variation graph to obtain more accurate peak calling of potential transcription factor binding sites.

4.0.3 Paper III: Assessing graph-based read mappers against a novel baseline approach highlights strengths and weaknesses of the current generation of methods

Paper III looked into the potential for using a two step approach for read mapping using the benefits of mapping to a graph for estimating a genotype for the sample, before using a linear mapper to obtain a final mapping that is not suffer from the too big search space of graph mapping.

Chapter 5

Discussion

This thesis comprises three projects that investigate different facets of graph based reference genomes. In this work, and in the work of others, some problems have become clear which relate to the potentially large number of sequences in the language of a sequence graph. Foremost of these is the computational complexity of dealing with this large set of potential sequences, and the large number of sequences that does not exist in real.

5.0.4 Computational Complexity

The work in paper II and III involved using mapping reads using *vg*. Although the method used there is in principle very similar to what is used by BWA-mem, the memory requirements and running times are significantly higher. Graph Peak Caller itself also suffers from requiring more memory and having longer running times than MACS2. Although algorithmic and data structural advances might alleviate these problems, some complexities might not be resolvable without turning to further approximations, or changing the interpretation of sequence graphs. Indeed, (*name*) et al showed that the indexing of sequence graphs is a hard problem with theoretical lower bound on complexity [?].

When the field of bioinformatics in recent years have been dominated by the sequencing capacity growing faster than Moore's law, it is worth to wonder how much accuracy gains is needed to justify the increase in computational complexity.

5.0.5 Invalid Sequences

Not only performance is affected by the large number of sequences in a graphs language. It can also lead to a large number of sequences that one would not expect to observe in nature. This is because all combinations of variants are represented in the language, but in nature variants can be highly positively or negatively correlated with each other. Such sequences makes graph mappers prone to mapping reads to sub sequences in the graph that matches the query string, but is unlikely to be the true origin of the read. This can lead to a loss of accuracy when mapping, but can also introduce biases in downstream analysis: Regions of the graph with many variants will be able to match many different query sequences and can thus become over represented when mapping reads to a graph. This potential bias is discussed in briefly in Paper II, where such over represented regions could be interpreted as peaks by the peak caller. The loss of accuracy in general can be seen in paper II in the relatively poor performance of graph mappers on reads not containing any variants. The two step approach introduced there can remedy such mapping effects, and it would be interesting in the future to use this approach also when calling ChIP-seq peaks, as this could remove the bias.

Recent work has introduced indexes for sequence graph that only index sub strings that are present in one of the sequences which was used to create the graph [?]. This approach has the potential of both reducing the computational complexity, and improve the accuracy of the read mapping. But in doing this, the simple interpretation of the language of a graph is not used, which makes this approach not as much mapping to a sequence graph as mapping to a set of similar sequences.