15 Genome comparison

(*Multiple*) genome comparison (or genome alignment) quite obviously deals with the comparison of several genomic sized sequences.

The aim is to better reconstruct genomic architecture (i.e. detect shared motifs) as well as phylogeny if the input consists of related genomes. The term can also be used when one compares for example different assemblies of the same genome or individual genomes of the same species. Depending on the input, adequate methods can be different.

15.1 Genome comparison

In general we can distinguish two main approaches:

- 1. *collinear alignment*: Here we are in the same model as in standard sequence alignment, except that the size of the input can be hundreds of millions of base pairs.
- 2. *non-collinear alignment*: Here we assume, that the genomes under inverstigation have undergone some evolutionary changes like duplications, rearrangement, or transposition of genomic sequence in addition to nucleotide changes.

We will address in this unit both approaches starting with the collinear alignment.

15.2 Genome comparison

A first step in collinear (and non-collinear) genome comparison is to use fast methods to compute local correspondences or *segment matches*. Using such correspondences, one can devise hierachical schemes to save on computation time as depicted below.



To do this, one has to address the problem of overlapping segment matches.

15.3 Segment match refinement and applications

This exposition is based on:

1. Aaron Halpern, Daniel Huson, Knut Reinert: Segment Match Refinement and Applications, WABI 2002.

We will present the algorithms used for refining a set of segment matches in such a way that the resulting set of refined matches does not contain any two matches that overlap, but the union of all matches in the refined set is the same as in the original set.

15.4 Motivation

Comparison of large, unfinished genomic sequences requires fast methods that are robust to misordering, misorientation, and duplications. A number of fast methods exist that can compute local similarities between such sequences, from which an optimal one-to-one correspondence might be desired. However, existing methods for computing such a correspondence are either too costly to run or are inappropriate for unfinished sequence.

We discuss an efficient method for refining a set of segment matches such that the resulting segments are of maximal size without non-identity overlaps. This resolved set of segments can be used in various ways to compute a similarity measure between any two large sequences, and hence can be used in alignment, matching, or tree construction algorithms for two or more sequences.

Assume the following two local overlaps are given. It can be easily seen that a greedy algorithm which takes the largest match first and does not allow overlaps can be suboptimal by a factor of almost 2 (ore more if more matches are involved) in the sense that we could also choose a large subsegment of the second match.



We assume the existence of a pre-computed set of local *segment matches* between two sequences. Such sets can be obtained efficiently by various techniques, e.g. BLAST, using suffix arrays, or hashing based schemes. However, a set of local segment matches is not disjoint: more than one match may cover a given interval, whereas the two analyses described above require either a one-to-one matching or a many-to-one matching.

In the following, we introduce an efficient and provably optimal approach to the problem of constructing one-to-one or many-to-one match sets by minimally "refining" (subdividing) matches until all overlaps between the refined matches are "resolved": the projections of any two refined matches onto each sequence are either disjoint or are identical.

Given such a resolved set of matches, several analyses that are otherwise difficult at best become quite straightforward, including not only the one-to-one matching problem and the many-to-one matching problems introduced above but also sequence alignment based on conserved segments.

For instance the one-to-one correspondence on the refined set is reduced to the classical matching problem. In addition, it is straightforward to compute an alignment between the sequences if desired (How?).

15.5 Definitions

Let $A = a_1 a_2 \dots a_p$ and $B = b_1 b_2 \dots b_q$ be two sequences. We call $A_{ij} = a_{i+1} a_{i+2} \dots a_j$ an *A*-segment, with $0 \le i < j \le p$. A *B*-segment is defined similarly. In the following, we will assume without specific mention that any definitions relative to *A* carry over to analogous definitions relative to *B*.

A segment match $S = (A_{ij}, B_{kl})$ consists of an A-segment A_{ij} and a B-segment B_{kl} . We call i, j, k and l the *left-A*, *right-A*, *left-B* and *right-B* positions of S, respectively. Given a set of segment matches S, let $supp_A(S)$ denote the A-support of S, i. e. the set of left-A and right-A positions of all matches $S \in S$.

Assume that every segment match $S = (A_{ij}, B_{kl})$ comes with two *projection maps* $\alpha_S : [i, j] \rightarrow [k, l]$ and $\beta_S : [k, l] \rightarrow [i, j]$ that reflect e.g. the alignment associated with *S*, as depicted in the below Figure.



In the case of a direct (i. e. orientation preserving) match, we require that α_S and β_S are *non-crossing*, i. e. that the two maps are monotonically increasing and that we have $\alpha_S(h) \le h' \Leftrightarrow h \le \beta_S(h')$ for all $h \in [i, j]$ and $h' \in [k, l]$.

In the case of a reversal, we require that the projections are (*completely*) *crossing*, i. e. that they are monotonically decreasing and that we have $\alpha_S(h) \le h' \Leftrightarrow h \ge \beta(h')$ for all $h \in [i, j]$ and $h' \in [k, l]$.

A match $S' = (A_{i'j'}, B_{k'l'})$ is called a *submatch of* S, if $i \le i' \le j' \le j$, $k \le k' \le l' \le l$, and we have $\alpha_{S'}(i') = k'$ or $\beta_{S'}(k') = i'$, and $\alpha_{S'}(j') = l'$ or $\beta_{S'}(l') = k'$. Here, we assume that α'_S , and β'_S , equals the restriction of α_S to [i', j'], and β_S to [k', l'], respectively.

All the previous definitions just formalize trivialities. Now we define the key concept, a *resolved refinement*. Consider a match $S = (A_{ij}, B_{kl}) \in S$. A set of matches S' is called a *refinement of* S, if each $S' \in S'$ is a submatch of S and S' *tiles* S, i. e.:

$$(i, j] = \bigcup_{(A_{i'j'}, B_{k'l'}) \in S'} (i', j']$$
 and $(k, l] = \bigcup_{(A_{i'j'}, B_{k'l'}) \in S'} (k', l'].$

Definition 1. A set S' of matches is called a *refinement of* S, if there exists a partitioning $S' = S'_1 \cup S'_2 \cup \ldots \cup S_n$ such that S'_i is a refinement of S_i for every $S_i \in S$.

We are particularly interested in refinements of S that have the property that the segments of any two matches are either disjoint or are identical and capture this property as follows:

Definition 2. A set of matches S is called *resolved*, if for any $S = (A_{ij}, B_{kl}) \in S$ we have $[i, j] \cap supp_A(S) = \{i, j\}$ and $[k, l] \cap supp_B(S) = \{k, l\}$. For technical reasons, we also require $\{\alpha_S(i), \alpha_S(j)\} \subseteq supp_B(S)$ and $\{\beta_S(k), \beta_S(l)\} \subseteq supp_A(S)$.

15.6 Example



The preceding example shows (a) A valid match $S = (A_{ij}, B_{kl})$, together with arrows depicting α_S (pointing down) and β_S (pointing up). (b) The set $S = \{S, S'\}$ is unresolved. (c) The set $\{S_1, S_2, S_3, S'\}$ is a minimal resolved refinement of S.

Any set of matches S possesses a (trivial) resolved refinement obtained by simply refining every $S \in S$ into a set of single position matches. We are interested in resolved refinements of minimal cardinality and claim:

Lemma 3. There exists a unique resolved refinement \overline{S} of S of minimal cardinality.

Proof: (Exercise). Sketch: Consider two different resolved refinements S_1 and S_2 of S, both of minimal cardinality. Divide proof into two cases. 1) ($supp_A(S_1) \neq supp_A(S_2)$ 2) $supp_A(S_1) = supp_A(S_2)$, $supp_B(S_1) = supp_B(S_2)$

We now show how to efficiently compute the minimal resolved refinement.

15.7 The algorithm

The input consists of set of segment matches S and projections α_S , β_S for all $S \in S$. The output is a resolved refinement \overline{S} of S of minimal cardinality.

We maintain a bipartite graph $G = (V_A \cup V_B, E \subseteq S \times V_A \times V_B)$ with nodes sets $V_A \subseteq \{a_0, \ldots, a_p\}$ and $V_B \subseteq \{b_0, \ldots, b_q\}$ representing *A*- and *B*-positions, respectively, and a set *E* of (labeled) edges that keep track of pairs of positions projected onto each other by a given sequence match.

(1) Refine($S, \{\alpha_S, \beta_S\}$); (2) $V_A = V_B = E = \emptyset$; (3) $W_A = supp_A(S)$; $W_B = supp_B(S)$; (4) while ($W_A \neq \emptyset \land W_B \neq \emptyset$) do (5) $A_refine(W_A)$; (6) $B_refine(W_B)$; (7) od (8) lexicographically order edges in E using (9) (S, i, k) or (S, i, -k) depending on whether α_S and β_S (10) are non-crossing or crossing; (11) $\overline{S} := \{(A_{ij}, B_{kl}) \mid (S, i, j) \text{ and } (S', k, l)$ (12) are consecutive and $S = S'\}$;

```
(1) A_Refine(W_A);
 (2) if this is not the first execution of A_refine;
       then W_B = \emptyset;
 (3)
 (4) fi
 (5) for each h \in W_A \setminus V_A do
         for each match S = (A_{ij}, B_{kl}) \in S with i \le h \le j do
 (6)
 (7)
             h' = \alpha_S(h);
             if \nexists edge (S, h, h')
 (8)
               then if \nexists h' \in W_B
 (9)
                        then create node h';
(10)
                      fi
(11)
                      G.new_edge(S, h, h'); insert h' into W_B;
(12)
             fi
(13)
(14)
         od
(15) od
(16) V_A = V_A \cup W_A;
```

(1) $B_{-}Refine(W_B)$; (2) $W_A = \emptyset;$ (3) for each $h \in W_B \setminus V_B$ do for each match $S = (A_{ij}, B_{kl}) \in S$ with $i \le h \le j$ do (4) $h' = \beta_S(h);$ (5) if \nexists edge (S, h, h')(6) then if $\nexists h' \in W_A$ (7) then create node *h*'; (8) fi (9) *G.new_edge*(S, h, h'); insert h' into W_A ; (10)fi (11) (12) od (13) Od (14) $V_B = V_B \cup W_B$;

15.8 Example



We go through the algorithm using the above example.

• Initially $W_A = \{a_0, \ldots, a_5\}$, $W_B = \{b_0, \ldots, b_5\}$ and $V_A = \emptyset$, $V_B = \emptyset$. We first *A_Refine*. a_3, a_4 lie in the match S_2 . Hence we create nodes v_1, v_2 and edges (S_2, a_3, v_1) and (S_2, a_4, v_2) . $W_B = \{b_0, \ldots, b_5, v_1, v_2\}$. $V_A = W_A$, since V_A was empty.

15.9 Example

- Now we call *B_Refine* for the first time. b_1, b_2, v_1, v_2 lie in the match S_2 . Hence we create a node u_0 (a_3 exists already) and edges (S_2, b_1, u_1) and (S_2, b_2, a_3). $W_A = \{u_0\}$. $V_B = W_B$ since V_B was empty.
- Now we *A*_*Refine*. u_0, a_3, a_4 lie in the match S_2 . But only u_0 is new. No new edge needs to be created since we assume the projections in the example to point to each other. Hence $W_B = \emptyset$ and we can stop the refinement.
- Finally we have to sort the edges lexicographically and determine the new segment matches.

15.10 Running time

Lemma 4. Given a set of n segment matches S between two sequences $A = a_1a_2...a_p$ and $B = b_1b_2...b_q$, with $p \ge q$. The graph G in algorithm Refine has at most O(kp) edges and can be computed in at most $O(p \log n + kp \log p)$ steps, where k is the maximal number of segment matches containing any given position.

Proof: The maximum number of different *A*-positions is *p*. To insert an *A*-position *h* into V_A , we first determine which segment matches contain it in $O(\log n + k)$ steps, using an interval tree.

For each of the at most *k* matches $S \in S$ that contain *h*, determine whether $\alpha_S(h)$ is contained in V_B in $O(\log p)$ steps. Putting this together, we obtain a bound of $O(p(\log n + k + k \log p)) = O(p \log n + kp \log p)$. Note that each of O(p) positions gives rise to at most *k* edges in *G*.

Lemma 5. Given a set of n segment matches S between two sequences $A = a_1a_2...a_p$ and $B = b_1b_2...b_q$, with $p \ge q$. The minimal resolved refinement \overline{S} of S can be computed in at most $O(p \log n + kp \log kp)$ steps, where k is the maximal number of segment matches containing any given position.

Proof: The preceding result implies that the ordering step in algorithm Refine will take at most *O*(*kp* log *kp*) computations.

Note that the size of a minimal resolved refinement depends on the length of the input sequences. For example, a set of only four segment matches between segments of length *d*, say, can give resize to 4*d* matches, as demonstrated in the below Figure.



15.11 Experiments

The method was tested with two (partial) versions of the drosophila genome. The first was obtained from the BDGP (http://www.fruitfly.org/sequence/dlMfasta.shtml), corresponding to 50,107,720 bp of two finished chromosome arms (2L and 3R). The second was a unpublished assembly of whole genome shotgun sequencing data conducted at Celera in late 2001.

A suffix-tree based method (Mummer) yielded 292,769 alignments seeded with a perfect match of at least 50 bp. These matches covered 49,953,627bp of the finished sequence (i.e., they covered all but 154,093 bp).

43,740,070bp were covered by a single match, but the remainder was covered two or more times; since the total length of all matches was 125,863,932, a relatively small fraction of the sequence is clearly involved in a large number of matches.

The refinement method described in this paper transformed the match set into 17,167,891 distinct matches. This took a total of 880 seconds on a Compaq Unix workstation.

The initial match set, with mean length 429 bp, was cut up into much smaller pieces (mean length 7bp) by the refinement. However, since most of the original sequence was covered by a single match, there should still be matches of considerable size; indeed, the maximum length refined match was 466,877bp, and of the finished sequence covered by matches, half is covered by a refined match of length at least 97,008 bp.

An optimal matching was determined on the refined match set, as described above. The resulting set of matches involved 262,885 refined matches and covered 49,875,599bp, leaving 232,121bp of the finished sequence uncovered.

The greedy method selected 1042 matches totaling 49,754,205 bp, leaving 353,515 bp of the finished sequence uncovered. The mean length is 47,749bp and the N50 is 136,332.

Hence it is questionable whether the refinment is suitable for this analysis. It has probably more merits when there are more overlapping matches and the refined matches are used to compute an alignment.

15.12 Summary

- In many applications it is important to have a 1-to-1 correspondence of subsequences, given a set of sequences.
- Normally, one has a given set of overlapping local alignments.
- One can transform this given set in a set of minimal cardinality such that no two segments intersect properly (they are disjoiunt or identical).