









#### Double Digest Problem

Double Digest Problem DDP

given

- 1. multiset A of lengths from enzyme A
- 2. multiset B of lengths from enzyme B
- 3. multiset C of lengths from enzymes A  $\land$  B

determine an ordering of the fragments that is consistent with the three experiments  $\mathbf{k} \in \mathbb{R}^{2}$ 

Saad Mnein





- The solution for a DDP might not be unique.
- The number of solutions grows exponentially





## Equivalence of Solutions

- Some different solutions might be equivalent.
- For instance, if (a<sub>1</sub>, a<sub>2</sub>, ..., a<sub>m</sub>) (b<sub>1</sub>, b<sub>2</sub>, ..., b<sub>n</sub>) is a solution, then (a<sub>m</sub>, a<sub>m-1</sub>, ..., a<sub>1</sub>) (b<sub>n</sub>, b<sub>n-1</sub>, ..., b<sub>1</sub>) is also a solution.
- This is a Reflection. No fragment length data could possibly distinguish between the two, they only differ by orientation.



## Overlap Equivalence

- Let's define a more general type of equivalence called overlap equivalence.
- Let {A<sub>i</sub>} be the set of fragments from A and {B<sub>j</sub>} be the set of fragments from B.
- A solution defines an overlap matrix *O*, s.t. *O*<sub>ij</sub> = 1 if *A*<sub>i</sub> overlaps with *B*<sub>j</sub>.
- Two solutions are overlap equivalent if they define the same overlap matrix O.
- Reflections are overlap equivalent.



#### Equivalence class

- A solution with all its overlap equivalent solutions form an equivalence class (this is an equivalence relation).
- Given a solution,
  - $-\operatorname{What}$  is the size of its equivalence class?
  - Can we generate all solutions in the class?





- If a solution has *t* 1 coincident cuts sites, then it has *t* components.
- They can be permuted in *t*! ways without changing the overlap data.
- Each component can also be reflected without changing the overlap.
- Therefore, we can generate 2<sup>t</sup>t! solutions.





#### Another observation

- Given a solution, let  $-\mathcal{A}_{j} = \{ A_{k}: A_{k} \subset B_{j} \}$  $-\mathcal{B}_{i} = \{ B_{k}: B_{k} \subset A_{i} \}$
- Permuting  $\mathcal{A}_i$  and  $\mathcal{B}_i$  does not change the overlap data





### Size of equivalence class

- Is it  $2^{t} ! \Pi |\mathcal{A}_{j}|! \Pi |\mathcal{B}_{j}|!$
- Not quite!
- If a component has only one fragment in either *A* or *B*, then a reflection is also a permutation.
- Let s be the number of such components, then the size of the equivalence class is:  $2^{(t\cdot s)} t! \prod |\mathcal{A}_{j}|! \prod |\mathcal{B}_{j}|!$

Other Equivalences?

- We can define other kinds of equivalences.
- Consider overlap size equivalence, i.e. two solutions are equivalent if they produce the same overlap sizes.
- Overlap equivalence => overlap size equivalence, but not the other way around.







- For 1 ≤ i ≤ j ≤ l, I<sub>c</sub> = { C<sub>k</sub>: i ≤ k ≤ j } is the set of fragments from C<sub>i</sub> to C<sub>j</sub>.
- The cassette defined by  $I_C$  is the set of fragments ( $I_A$ ,  $I_B$ ) that contain a fragment from  $I_C$ .















#### Cassette Equivalence

- Two solutions are *cassette equivalent* if there exists a series of cassette transformations (exchanges and reflections) that take on to the other.
- What is the size of an equivalence class?



Saad M

Soud Mr

#### Alternating Euler Paths

- Consider a graph with colored edges
- An Euler path (cycle) is a path (cycle) that goes through every edge
  once
- An alternating Euler path (cycle) is an Euler path (cycle) such that consecutive edges on the path (cycle) have different colors
- Pevzner 1995 showed that given a solution, we can construct a special bi-colored graph called the border block graph.
- Each cassette equivalent solution corresponds to an alternating Euler path (cycle) in the graph and vice-versa.

#### Fact

- Let d<sub>i</sub>(v) be the number of edges of color i incident to v.
- An edge bi-colored connected graph with  $d_A(v) = d_B(v)$  has an alternating Euler cycle.
- Proof:
  - Every vertex has even degree; therefore, the graph contains an Euler cycle.
  - Construct the Euler cycle the usual way, but by using only distinct color edges when traversing a vertex.



- Consider an alternating path ....*x*...*y*...*x*...*y*...
- It consists of 5 parts  $F_1F_2F_3F_4F_5$
- $F_1F_2F_3F_4F_5 \rightarrow F_1F_4F_3F_2F_5$  is called an exchange if  $F_1F_4F_3F_2F_5$  is an alternating path





# Reflection• Consider an alternating path<br/> $\dots x \dots x \dots$ • It consists of 3 parts $F_1F_2F_3$ • It consists of 3 parts $F_1F_2F_3$ • $F_1F_2F_3 \rightarrow F_1F_2^-F_3$ is called a reflection if<br/> $F_1F_2^-F_3$ is an alternating path, where $F_2^-$ is<br/>the reverse of $F_2$ .



#### Fact

- Every two alternating Euler cycles in a bi-colored graph can be transformed into each other by a series of exchanges and reflections.
- Proof: Pevzner p. 29



#### The border blocks

- Let  $I(A_i) = \{ C_k: C_k \subset A_i \}$
- Let  $I(B_j) = \{ C_k: C_k \subset B_j \}$
- If |*I*(*X*)| > 1, define the border blocks of *X* to be the left most and right most block in *I*(*X*).
- *C<sub>i</sub>* is a border block is it is a border block for some fragment *X*.





#### Border graph

- Let  $\mathcal{B} = \{ C_k : C_k \text{ is a border block } \}$
- $V = \{ |C_k| : C_k \in \mathcal{B} \}$  vertices
- $\mathsf{E} = \{ (|C_i|, |C_j|) : C_i \text{ and } C_j \in \mathcal{B} \cap I(X) \text{ for some } X \}$
- Each edge labeled by its X and colored A if  $X \in A$  and B if  $X \in B$ .

Saad M



# Alternating Euler path in border block graph

- Each vertex has equal number of edges of each color, except possibly for  $|C_1|$  and  $|C_l|$ .
- By adding one or two edges (depending on the colors) we can fix this. Therefore, the graph has an alternating Euler path or cycle.
- Let  $C_1...C_m$  be the ordered set of border blocks, then  $P = |C_1|...|C_m|$  is an alternating Euler path (cycle).

Saad Mne

#### Result

Cassette transformations do not change the border graph.

Let P be the alternating Euler path (cycle) corresponding a solution [A, B].

- If a solution [A', B'] is obtained from [A, B] by cassette exchange (reflection), it will have a path P that can be obtained from P by an exchange (reflection).
- Let P' be an alternating Euler path (cycle) obtained from P by exchange (reflection). Then there is a solution [A', B] that can be obtained from [A, B] by cassette exchange (reflection), where P corresponds to [A', B]













#### DDP is NP-complete

#### Proof:

- $-DPP \in NP$ . A solution for DDP can be verified in polynomial time.
- Set Partition problem (classical one), which is NP-complete, reduces to *DDP* in polynomial time.



Saad Mne

#### $DDP \in NP$

#### Given

- multiset A of lengths from enzyme A - multiset B of lengths from enzyme B
- multiset C of lengths from enzymes A + B
- Solution
  - two sets of restriction sites, a and b.

#### • Verification:

- Sort  $g = a \cup b \cup \{0, L\}$ , L =sum of all lengths in A
- Compute multiset  $c = \{ c_i, c_i = g_{i+1} g_{i}, 0 < i < |g| \text{ and } g_{i+1} \neq g_i \}$  Sort *c* and *C* and compare them





## Partial Digest Problem (turnpike problem)

Partial Digest Problem PDP

given

multiset  $\Delta x$  of distances between every pair of points on the line  $|\Delta x| = \binom{n}{2} = \frac{n(n-1)}{2}$ 

reconstruct the points on the line

#### PDP

- No polynomial time algorithm know
- Not known to be NP-complete
- Practical Backtracking algorithm due to Skiena et al. 1990



#### The algorithm

- Find longest distance in  $\Delta X$  this decides the two outermost points, delete that distance from  $\Delta X$
- Repeatedly position the longest remaining distance of  $\Delta X$ .
- Since the longer distance must be realized from one of the two outermost points, we have two possible positions (left and right) for the point.
- For each of these two positions, check whether all the distances from the position to the points already positioned are in  $\Delta X.$
- If they are, delete all those distances from  $\Delta X$  and proceed.
- Backtrack if they are not for both of the two positions.

















