

Fragment Assembly.

⑦ DNA sequencing & assembly:

Shotgun sequencing starts with a large sample of genomic DNA. The sample is sonicated, a process which randomly partitions each piece of DNA in the sample into inserts; the inserts that are smaller than 500 nucleotides are removed from further consideration. Before the inserts can be read, each one must be multiplied billions of times, so that it is possible to read them through gel electrophoresis.

⑧ Shortest superstring problem:

Given a set of strings find the shortest string that contains all of them.

→ graph Set of strings $\{s_1, s_2, \dots, s_n\}$

→ graphical representation

$s_i \rightarrow s_j$ if s_i overlaps s_j
(Prefaces of s_i = prefixes of s_j).

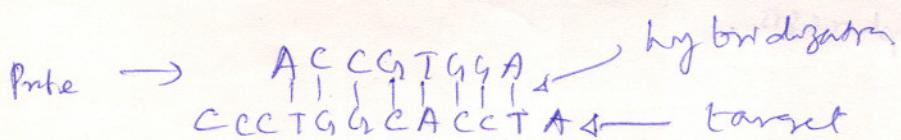
→ equivalent to Travelling Salesman Problem

(or a Hamiltonian Path finding problem) in a directed graph with n vertices (correspond to s_i 's) & edges of length ($\text{overlap}(s_u, s_v)$).

→ greedy strategy: repeatedly merge a pair of strings with max. overlap until only one string remains.

② Sequencing by hybridization (SBH)

SBH involves in building a miniature DNA array (known as also ^aDNA chip) that contains thousands of short DNA fragments called probes. Given a short probe (an 8- to 30 nucleotide single-stranded synthetic DNA fragment), the target will hybridize with the probe if the probe is substring of the target's ~~Watson~~ complement. When the



SBH relies on the hybridization of the target DNA fragment against a very large array of short probes. In this manner, probes can be used to test the unknown target DNA to determine its l-mer composition. The universal DNA array ~~contains~~ contains all 4^l probes of length l.
~~and so on~~

- Spectrum (s, l): For a string s of length n , the l-mer composition or spectrum of s , is the multiset of $n-l+1$ l-mers in s ,

e.g. Spectrum (TATGGGTGC, 3)

$$= \{TAT, ATG, TG, GGT, GTG, TG, GC\}$$

• SBH Problem:

Reconstruct a string from its l -mer composition.

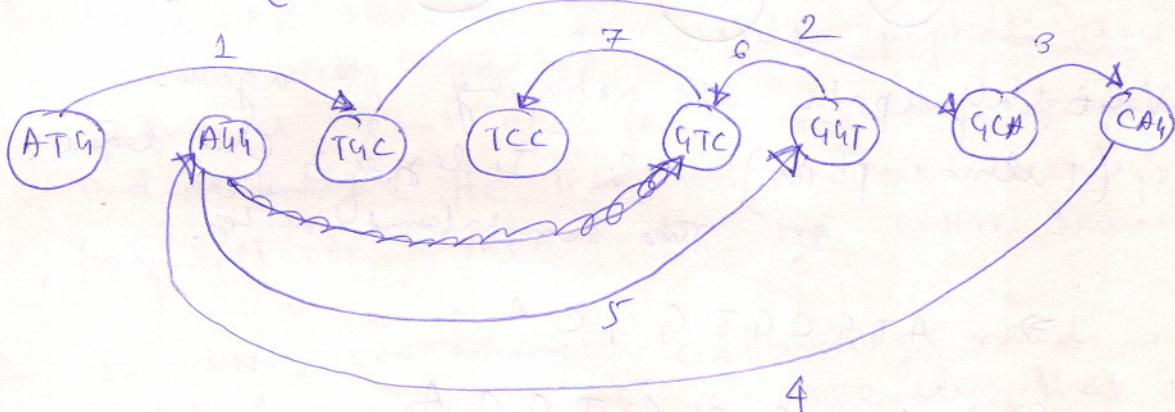
Input: Spectrum = (s_1, s_2, \dots, s_n) ; $O/p = s$.

• Hamiltonian approach:

$$V = \{s_1, s_2, \dots, s_n\}$$

$E = \{s_i \rightarrow s_j\}$, if s_i is $l-1$ mer suffix of s_j
 $= l-1$ mer prefix of s_j .

e.g. $s = (\text{ATG } \text{AGG } \text{TAC } \text{TCC } \text{GTC } \text{GGT } \text{GCA})$



(nos. denotes sequence of edges to be traversed).

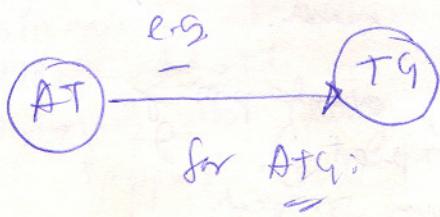
\therefore String: ATGCAAGGTCC

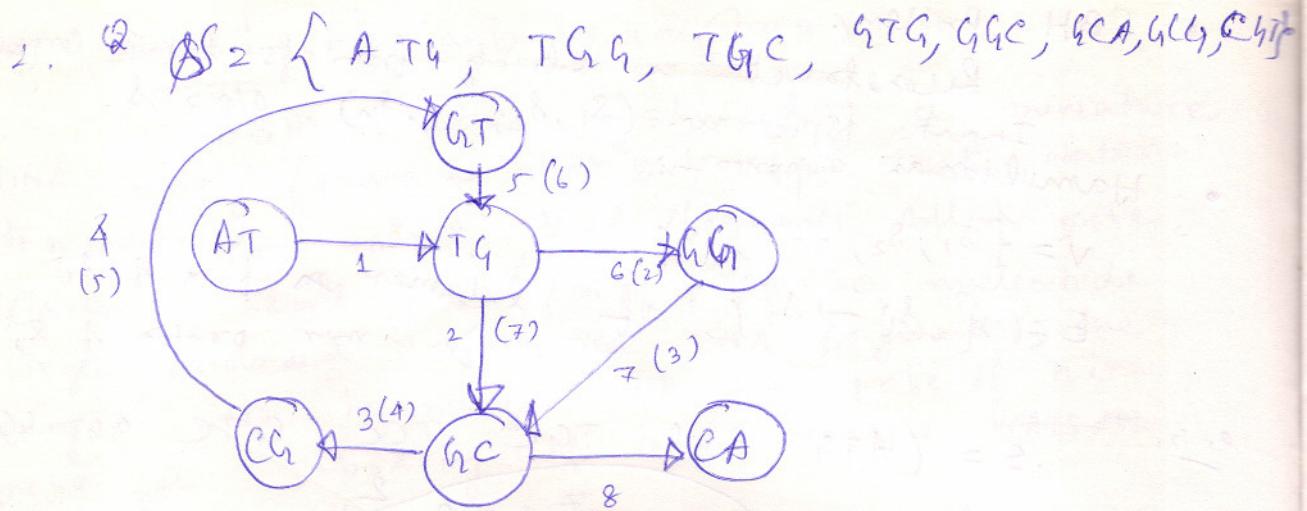
belongs to
However the problem is \sim NP complete set.

• Eulerian Path Approach:

Build a graph with set of all $l-1$ -mers.

as vertex & an edge is connected between two vertices, if a string has its prefix l suffixes as those vertices.





Visit all path w/o retracing it again.

(Euler path) i.e. Indegree = Out degree,
& Two semi balanced vertex.

⇒ ATGCGTGGCA

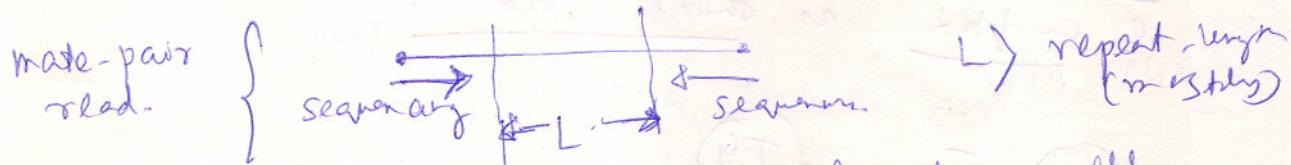
or, or, ATGCGCTGCA

→ linear time traversal & reconstruction

④ Problems in Fragment Assembly in DNA sequencing:

- 1% to 3% error rate in reading insert.
- Double stranded DNA → which strand?
- Repeats ⇒ major problem.

Weber & Mayer suggested reading from two ends of a long insert (with a fixed gap).



(This is unlikely, overlapping will
be ambiguous due to repeat). If it is
unlikely both end will contain repeat

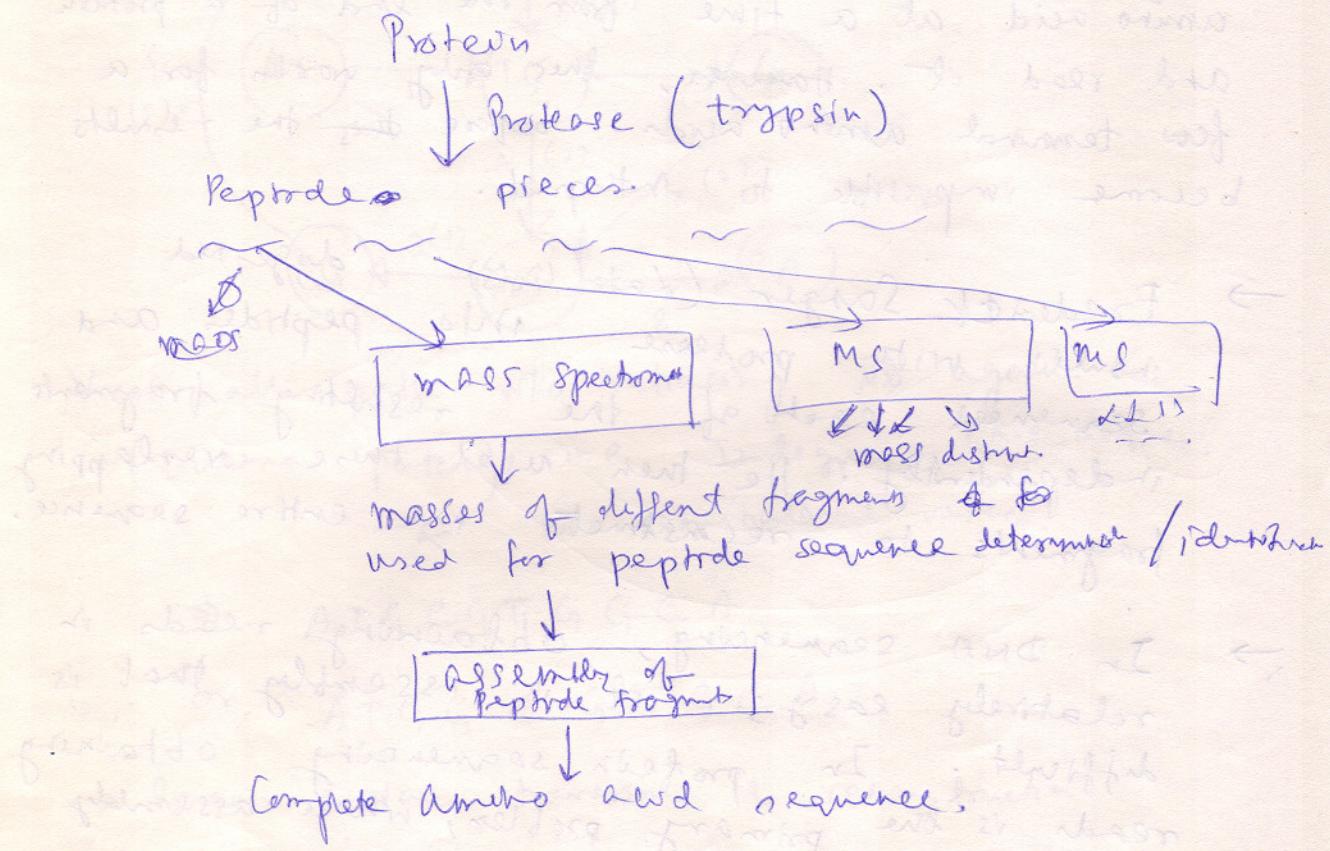
① Protein Sequencing and Identification:

- Edman degradation reaction to chop off an amino acid at a time from the end of a protein and read it. However, this only works for a few terminal amino acids before the results become impossible to interpret.
- Frederick Sanger (late 40's) digested insulin with protease into peptides and sequenced each of the resulting fragments independently. He then used these overlapping fragments to reconstruct the entire sequence.
- In DNA sequencing, obtaining reads is relatively easy; it is assembly that is difficult. In protein sequencing obtaining reads is the primary problem, while assembly is easy.

* Protein analysis by mass spectrometry:

A mass spectrometer works like a charged sieve. A large molecule (peptide) gets broken into smaller fragments that have an electrical charge. These fragments are then spun around and accelerated in a magnetic field until they hit a detector. Because large fragments are harder to spin than small ones, one can distinguish between fragments with different masses based on the amount of energy required to flip the different fragments around. It happens that most molecules can be broken in several places, generating several dif.

ion types, the problem is to reconstruct the amino acid sequences of the peptide from the masses of these broken pieces.



- Peptide Sequencing problem:

$A = \{a_1, a_2, \dots, a_{20}\}$ \leftarrow Set of amino acids,

$P = p_1 \dots p_n$ \leftarrow a peptide (Set of amino acids $p_i \in A$)

$m(a_i)$ \leftarrow mass of amino acid a_i .

$$\therefore m(P) = \sum_{i=1}^n m(p_i)$$

N-terminal peptide : $p_1 \dots p_i \Rightarrow P_1 = m_1 = \sum_{j=1}^{i-1} m(p_j)$

C-terminal peptide : $p_i \dots p_n \Rightarrow P_2 = m(P) - m_i$

$$i \leq n$$

Mass spectra obtained by tandem mass spectrometry (MS/MS) consists predominantly of partial N-terminal peptides & C-terminal peptides.

Partial peptide may also loose a water H_2O (18 dalton) or Ammonia (NH_3) (17) or both (18 + 17).

Let $\Delta = \{\delta_1, \delta_2, \dots, \delta_K\}$ be the set of such numbers, which may be lost from the molecular weight of a peptide fragment P_i with mass m_i so that it is converted to any one of the elements of $\{m_i - \delta_1, m_i - \delta_2, \dots, m_i - \delta_K\}$

δ -ion of N-terminal partial peptide P_i modifies its mass m_i to $m_i - \delta$ & the ion is called b-ions, similarly δ -ions of C-terminal partial peptide P_i^- , are called y-ions.

e.g. $b-\text{H}_2\text{O}$, $y-\text{NH}_3$ etc., $b-\text{H}_2\text{O}-\text{NH}_3$ etc.

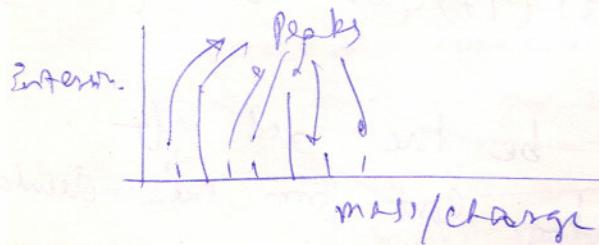
Theoretical spectrum = $T(P) = \left\{ \begin{array}{l} \text{masses of} \\ \text{partial peptides} \\ \text{and masses of all possible derivatives due to} \\ \text{subtraction of } \delta\text{-ions} \end{array} \right\}$

$$= \left\{ \begin{array}{l} \{m(P), m(P) - \delta_1, \dots, m(P) - \delta_K\} \\ \{m(P) - m_i, m(P) - m_i - \delta_1, \dots, m(P) - m_i - \delta_K\} \end{array} \right\}$$

Experimental spectrum = $S = \{s_1, \dots, s_N\} =$
experimentally observed masses from tandem mass spectrometry.

$$\text{Shared peak counts} = \text{No. of } |SNT(P)|$$

Mass spectrometry provides the intensity of different mass distribution. Hence, maxima of which are denoted by discrete peaks in the interval.



e.g.

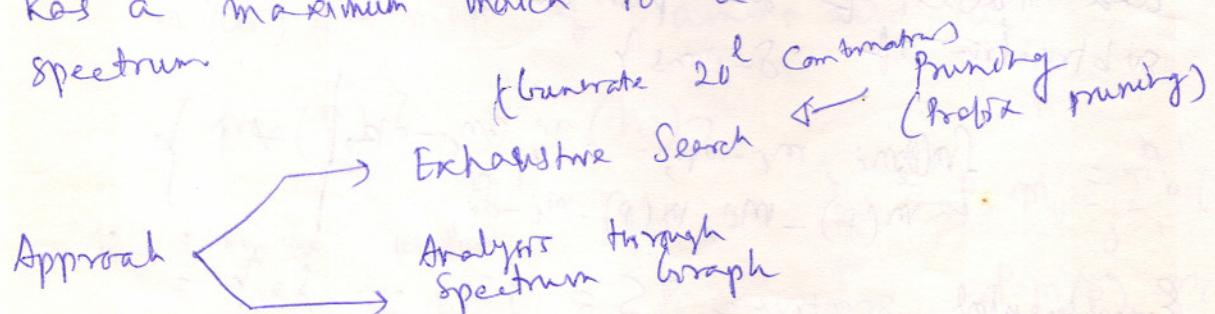
$$\phi = GPFNA \quad d_1=0 \quad \delta_2=18 \quad \delta_3=17 \quad \delta_4=35$$

	Sequence	mass	m/z H_2O	m/z NH_2	m/z H_2O+NH_3
	GPFNA	498	480	481	463
b_1	G	58	40	41	23
y_1	PFNA	442	424	425	405
b_2	GP	119	131	132	119
y_2	FNA	351	133	334	316

; and so on $(b_3, y_3), (b_4, y_4), \dots$

Prob. Statement:

Find a peptide whose theoretical spectrum has a maximum match to a measured experimental spectrum.



- Spectrum Graph: Let us assume Spectrum contains masses of b-ions only (N-terminal partial peptide)

$$\text{wt. } S = \{s_1, s_2, \dots, s_n\}$$

$$\text{Let } \Delta = \{\delta_1, \delta_2, \dots, \delta_k\}$$

For each s_i , there may be any of $s_i + \delta_j, 1 \leq j \leq k$ masses. There are $(k+1)$ vertices for s_i (including itself). We define edge between two vertices if it differs by a single amino acid mol. wt. (and $\neq 0$).



$$\text{if } (s_i + \delta_m) - (s_i + \delta_l) = m(a_k)$$

$$a_k \in \{a_1, a_2, \dots, a_{20}\}$$

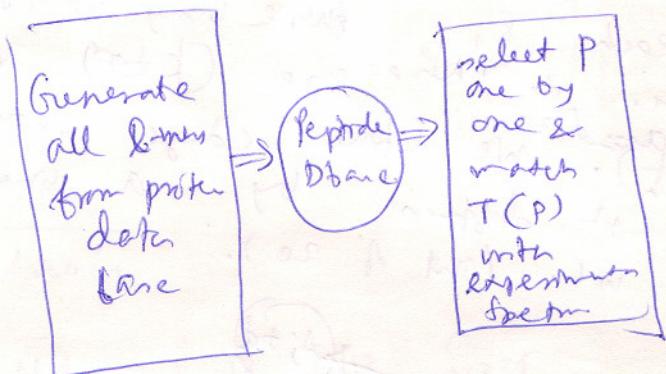
Add vertex O & m. Then peptide sequencing problem can be cast as finding a path from O to m in the resulting DAG.

No. of vertices: $9R+2$.

Complete Spectrum graph: If every N-type partial peptide is present, for a complete graph. If a $(n+1)$ length path from O to m is $m - R + (N+1)$. Often it is the longest path (max edges) of the DAG. There may be ambiguous paths. The graph may be also incomplete.

Protein Identification:

Prob. Statement: Find a protein from a database that best matches the experimental spectrum.



Problem needs modification
get modified due to many amino acids
phosphorylation etc

e.g. A chemical modification at Nth amino acid of a peptide sequence P_1, P_2, \dots, P_n results in increased mass of the N-terminal peptides $P_1^-, P_2^-, \dots, P_n^-$ & increased C-terminal peptides $P_1^-, P_2^-, \dots, P_{n-1}^-$

Modified prob. statement: Find a peptide from the database that best matches the experimental spectrum with up to k modifications.

Spectral Convolution:

$$\text{Multiset } S_2 \ominus S_1 = \{ s_2 - s_1 : s_1 \in S_1, s_2 \in S_2 \}$$

Let $S_2 \oplus S_1(x) = \text{multiplicity of } x$

If $S_2 \oplus S_1$ match (for, &
w/o any modifications $S_2 \ominus S_1(0)$ should
be high (maximum). It gives the count of
shared peaks.

For k -modifications : $S_2 \ominus S_1(x)$ will have peaks at different values of x

e.g. for $k=1$, if the mutation takes place at t th amino acid,

$$S_2 \ominus S_1(0) = t - \text{No. of N-terminal peptide}$$

for \oplus

$$S_2 \ominus S_1(m(P_{1t}) - m(P_{2t})) = n - t.$$

It is therefore reasonable to define spectral similarity as overall height of k highest peaks of $S_2 \ominus S_1$.

However, different combinations may result same $S_2 \ominus S_1$ (though some of the combinations are not valid set of partial peptide masses), as they violate const monotonicity of masses).

Spectral Alignment:

Let $S = \{s_1, s_2, \dots, s_n\}$ be an ordered set

of integers e.v.e. $s_1 < s_2 < \dots < s_n$.

A shift Δ_i transforms S into $\{s_1 + \Delta_i, s_2 + \Delta_i, \dots, s_{i-1} + \Delta_i, s_i + \Delta_i, \dots, s_n + \Delta_i\}$. (Ordering remains intact).

k similarity between A and B

$\equiv D(k) = \text{max. no. of elements in common}$
between these sets after k -shift.

Ex:

$$A = \{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$$

$$B = \{10, 20, 30, 40, 50, \underset{+5}{55}, \underset{-5}{65}, 75, 85, 95\}$$

$$\therefore D(CD) = 10, \text{ due to } -5 \text{ shift at } 65_{B_{10}}$$

$$C = \{10, 15, 30, 35, 50, 55, 70, 75, \underset{-5}{90}, \underset{+5}{95}\}$$

$$D(CD) = 6, \quad 5 \text{ element match}$$

Spectral Alignment Problem:

Find the k -similarity b/w two sets.

$$A = \{a_1, \dots, a_n\} \equiv \{00 \dots 0 \overset{a_1-1}{\overbrace{1}} 0 \dots 0 \overset{a_2-1}{\overbrace{1}} \dots\}^T$$

\downarrow a_{n-m} zeroes

$$B = \{b_1, b_2, \dots, b_m\} \equiv \{0 \dots 0 \overset{b_1-1}{\overbrace{1}} 0 \dots 0 \overset{b_2-1}{\overbrace{1}} \dots\}^T$$

\downarrow b_{m-n} zeroes

Any shift δ by δ in A means insertion or deletion of zeroes. Dynamic programming for the no. of matches with k -shift. \Rightarrow longest path problem.

$$\text{Spectral product: } A \otimes B = \{(a_i, b_j) \mid a_i \in A, b_j \in B\}$$

2-D matrix representation can \Rightarrow

$$A \otimes B = \sum_{i=1}^n \sum_{j=1}^m \begin{bmatrix} & \\ & \\ & 1 \\ & \\ & \end{bmatrix}$$

$$A \otimes B = \begin{bmatrix} a_{00} & a_{01} & \dots & a_{0m} \\ a_{10} & a_{11} & \dots & a_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k0} & a_{k1} & \dots & a_{km} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m0} & a_{m1} & \dots & a_{mm} \end{bmatrix}$$

a_{kl}

$m \times n$
main Diagonal

$$A \otimes B(k, l) = 1, \quad \text{if } k \in A \text{ & } l \in B$$

$D(k) \Rightarrow k \stackrel{=0}{\sim} 0$ otherwise.

$D(0) = \text{No. of } 1's \text{ in main Diagonal.}$

$D(1) = D(0) + \text{Another diagonal containing max. no. of } 1's$
(out of the set of diagonals).

$$\therefore D(k) = D(0) + \sum_{i=k}^{n-1} \text{No. of } 1's \text{ in } i-th \text{ diagonal.}$$

$D(k) \geq k$ similarity between spectra

$\Rightarrow \text{max. no. of } 1's \text{ on a path through}$
the spectral matrix that uses at ~~the~~
most $(k+1)$ diagonals ~~and~~

k -spectral alignment \Rightarrow the path that uses
 $(k+1)$ diagonals leading to $D(k)$.

Algorithm: : Let $A_i = \{a_1, a_2, \dots, a_i\}$
 $B_j = \{b_1, b_2, \dots, b_j\}$

\rightarrow Compute $D_{ij}(k)$ for A_i & B_j .

\rightarrow Using Dynamic Prog. & solve for $D_{mn}(k)$.

$\rightarrow (i', j')$ & $(0, j)$ are codiagonal if
 $a_i - a_{i'} = b_j - b_{j'}$.

$\rightarrow (i', j') < (0, j)$ if $i' < i$ & $j' < j$

\rightarrow Every (i, j) is codiagonal with $(0, 0)$.
 $(0, 0) \in \text{Null set}$ $A_0 = \emptyset$, $B_0 = \emptyset$

$\rightarrow D_{00}(R) = 0, \forall R.$

$$\therefore D_{ij}(R) = \max_{(i', j') < (i, j)} \begin{cases} D_{i'j'}(k) + 1, & \text{if } (i', j') \text{ & } (i, j) \text{ are codiagonal.} \\ D_{i'j'}(k-1) + 1, & \text{else.} \end{cases}$$

Time Complexity:

For computing at state (i, j) by scanning the previously visited states (s.t. $i' < i \leq j' \leq j$)
no. of such states: D_{ij} .

$$\therefore \text{Total cost: } R \cdot \sum_{j=1}^m \sum_{i=1}^n D_{ij} = R \cdot \sum_j \sum_i = \frac{m(n+1)^2}{4} \cdot R$$

Compute D_{ij}
upto R^{th} level.
 $i = 0, 1, \dots$

* An efficient algo.
for reducing the time complexity make
compute the similarity from the closest diagonal
state s.t. (1) maximal diagonal pair $(0, j')$
s.t. $i' < i$ & $j' < j$ or $(0, 0)$ if
no s.t. pair does not exist)

Define also :

$$M_{ij}(k) = \max_{(i',j') < (i,j)} D_{i'j'}(k)$$

$$\therefore D_{ij}(k) = \max \left\{ \begin{array}{l} D_{\text{diag}(i,j)}(k) + 1 \\ M_{i-1,j-1}(k-1) + 1 \end{array} \right.$$

$$M_{ij}(k) = \max \left\{ \begin{array}{l} D_{ij}(k) \\ M_{i-1,j}(k) \\ M_{i,j-1}(k) \end{array} \right.$$

$$\Rightarrow \text{Time} \geq O(n^2 k)$$

* Complexities:

→ presence of
Both N-terminal & C-terminal ions
complicate the problem.