① Architecture Team: Discuss TSP interface.
② Class: Discussion on TSP. { Go over LP-relaxation }.
③ Optical Mapping:
④ | Quiz. |
⑤ Optical Mapping Architecture Team.

## Optical Mapping:

◊ Restriction Map Model:

SMRM (Single Molecule Restriction Map)

A vector with ordered set of rational numbers on the open interval $(0,1)$:

$$D_j = (s_{1j}, s_{2j}, \cdots, s_{M_j j}),$$

$$0 < s_{1j} < s_{2j} < \cdots < s_{M_j j} < 1. \qquad s_{ij} \in \mathbb{Q}$$

◊ Problem

Data: A collection of SMRM vectors:
$$D_1, D_2, \cdots, D_m$$

Desiderata: Compute a consensus vector
$$H = (h_1, h_2, \cdots, h_N)$$

such that $H$ is "consistent" with each $D_j$.

$$H^* = \arg\min_{H, j} \overset{dist}{}(D_j, H).$$

Consensus:

$$H^* = \underset{H, j}{\text{argmin}} \; \text{dist}(D_j, H)$$

$$D_j = (s_{1j}, s_{2j}, \dots, s_{M_j j})$$

$$\Rightarrow \quad D_j + c = (s_{1j} + c, s_{2j} + c, \dots, s_{M_j j} + c)$$

$$c \in [0, 1) \quad c \in \mathbb{Q} \quad -s_{1j} < c < 1 - s_{M_j j}$$

$$\text{dist}(D_j, H) = \text{dist}(D_j + c, H)$$

$$D_j^R = (1 - s_{M_j j}, \dots, 1 - s_{2j}, 1 - s_{1j})$$

$$\text{dist}(D_j^R, H) = \text{dist}(D_j, H).$$

Consensus:

$$H^* = \underset{H, j}{\text{argmin}} \; \{ \text{dist}(D_j, H), \text{dist}(D_j^R, H) \}$$

or

$$H^* = \underset{H, j}{\text{argmin}} \; \{ \text{dist}(D_j + c, H) \mid -s_{1j} < c < 1 - s_{M_j j} \}$$

or

$$H^* = \underset{H, j}{\text{argmin}} \; \{ \text{dist}(D_j + c, H), \text{dist}(D_j^R + c, H) \mid -s_{1j} < c < 1 - s_{M_j j} \}$$

Assume some distribution generating $D_j$'s

$$\Rightarrow \quad \text{Maximum Likelihood formulation} \Rightarrow$$

$$\langle H \rangle = \underset{H}{\text{argmin}} \sum_j \min \{ \text{dist}(D_j + c, H), \text{dist}(D_j^R + c, H) \mid 1 - s_{1j} < c < 1 - s_{M_j j} \}$$

(Unknown Orientation:)

Data: A set of ordered vectors with rational entries in the open interval $(0,1)$:

$$D_1, D_2, \cdots, D_\ell, D_{\ell+1}, \cdots, D_m$$

A rational number $p_c \in (0,1)$ and an integer $N$.

An admissible alignment of the data can be represented as

$$D_1', D_2', \cdots, D_\ell', D_{\ell+1}', \cdots, D_m'$$

where

$$D_j' \in \{ D_j, D_j^R \} \qquad (1 \le j \le \ell)$$

and

$$D_j' = D_j \qquad (j > \ell)$$

$$\left.\begin{array}{c}\\\\\end{array}\right\} = \text{An Alignment } (A_k)$$

For any rational number $h_i \in [0,1]$, define an indicator variable

$$m_{ijk} = \begin{cases} 1 & \text{if } h_i \in D_j' \\ 0 & \text{otherwise.} \end{cases}$$

Define a characteristic function

$$\chi_k : [0,1] \longrightarrow \{0,1\}$$

$$: h_i \longmapsto \begin{cases} 1 & \text{iff } \sum_j m_{ijk} > p_c m. \end{cases}$$

<u>Desiderata</u>: Find an admissible alignment $A_k$ such that

$$\left| \left\{ h \in [0,1] \mid \chi_k(h) = 1 \right\} \right| \geq N.$$

## NP-Completeness

Consider an instance of a 3-SAT problem:

With $\ell$ variables:

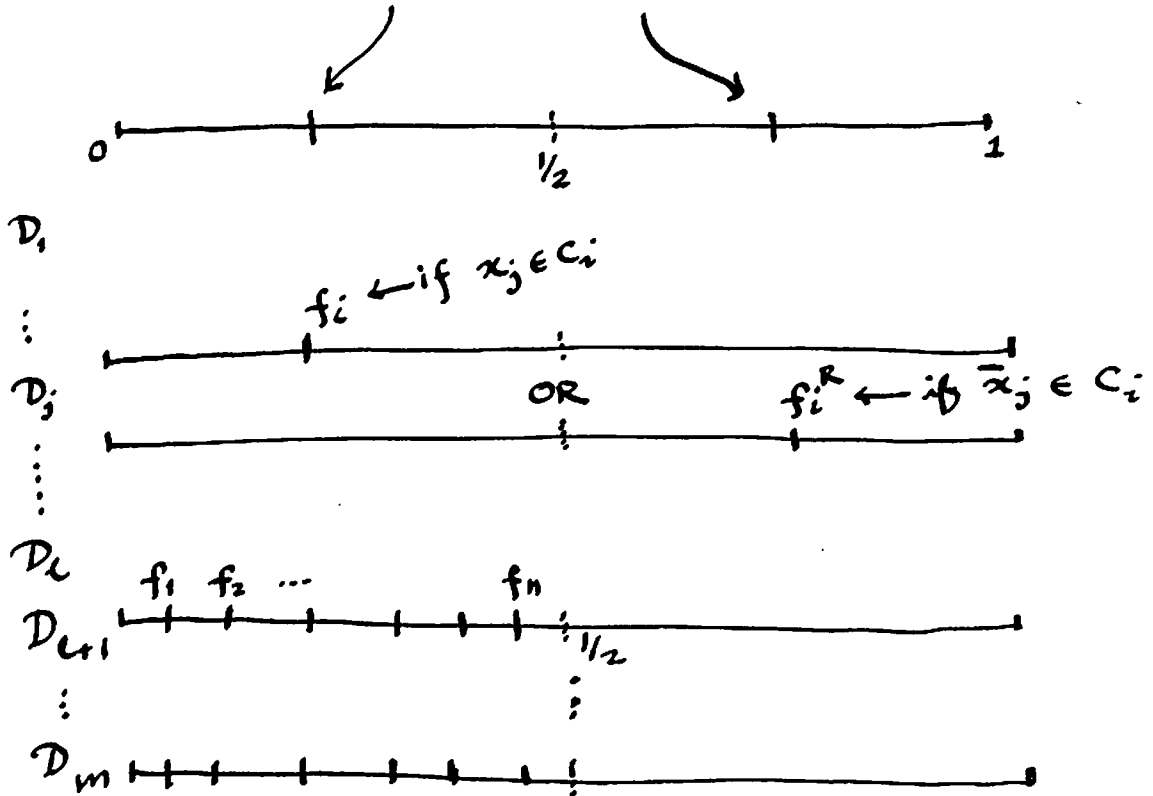$$x_1, x_2, \ldots, x_\ell$$

And $n$ clauses:

$$C_1, C_2, \ldots C_n \qquad (n \geq \ell)$$

◇ Assume that no clause contains a variable and its negation: $x_j$ and $\bar{x}_j$ (The clause is a tautology $\equiv T$)

◇ Restriction site associated with a ~~cut~~ clause $C_i$.

$$f_i = \frac{i}{2(n+1)} \qquad f_i^R = 1 - f_i = \frac{2n-i+2}{2(n+1)}$$



$f_i \leftarrow$ if $x_j \in C_i$

OR

$f_i^R \leftarrow$ if $\bar{x}_j \in C_i$

$D_1$

$D_j$

$D_\ell$

$D_{\ell+1}$

$f_1 \; f_2 \; \cdots \qquad f_n$

$D_m$

● Create a dataset $D_1, D_2, \ldots, D_\ell, D_{\ell+1}, \ldots D_m$ (22)

with $m = 2\ell - 1$ as follows:

$D_j$ has a cuts at $f_i$ or $f_i^R$, only:

$$f_i \in D_j \quad \text{iff} \quad x_j \in C_i$$

$$\left( f_i^R \in D_j \quad \text{iff} \quad \bar{x}_j \in C_i \right)$$

$$N \equiv n, \quad p_c = \tfrac{1}{2}$$

CNF has a satisfying assignment

● $\Rightarrow$ Choose an admissible alignment in which

$$D_j' = \begin{cases} D_j & \text{if} \quad x_j = \text{true} \\ D_j^R & \text{if} \quad x_j = \text{false} \end{cases} \Bigg\} \; 1 \leq j \leq \ell$$

$$D_j' = D_j, \quad \ell < j \leq m.$$

∴ For every $f_i$, $(1 \leq i \leq n)$ there are $(\ell-1)$ matches
from $D_{\ell+1}, \ldots, D_m$

& at least one more from $D_1', \ldots, D_\ell'$
(since each clause must be satisfied)

● ∴ $\forall_{1 \leq i \leq n}$

$$\sum_j m_{ijk} \geq \ell > \frac{2\ell - 1}{2} = p_c m.$$

$$\Rightarrow \{ h \in [0,1] \mid x_k(h) = 1 \} = \{ f_1, f_2, \ldots, f_n \}$$

$$\Rightarrow | h \in [0,1] \mid x_k(h) = 1 \} = n \geq N.$$

Conversely, if the CNF has no satisfying assignment, then for every admissible alignment there exists an

$$1 \le i \le n$$

$$\forall_k \exists_i \qquad \sum_j m_{ijk} = (\ell - 1) < p_c m \qquad \text{and}$$

$$\left| \left\{ h \in [0,1] \mid \chi_k(h) = 1 \right\} \right| < n.$$

$$\square.$$

## Problem Generation:
### Statistical Model:

◇ A model or hypothesis $H$.

$$= \{ h_1, h_2, \ldots, h_N \}$$

$N \approx 40.$
Distribution for $h_i$'s
Exponential gaps
or uniform gaps.

◇ $Pr[D_j \mid H]$

$$D_j \sim H. \qquad \left\{ \begin{array}{l} \text{Pairwise Conditional Indep.} \\ Pr[D_j \mid D_{j_1}, \ldots, D_{j_m}, H] \\ \qquad = Pr[D_j, H] \end{array} \right.$$

⬦ $Pr[bad]$, $Pr[good] = 1 - Pr[bad]$

$$Pr[D_j \mid H] = \frac{1}{2} \sum Pr[D_j^{(k)} \mid H, good] \, Pr[good]$$

$$+ \frac{1}{2} \sum Pr[D_j^{(k)} \mid H, bad] \, Pr[bad]$$

$$(k) \to \text{Alignment}.$$

$$D_j^{(k)} = D_j \text{ or } D_j^R \quad \text{with equal probability:}$$

$D_j = $ Good $\Rightarrow$

Choose parameters $p_c, \sigma, f$.

$h_i \in H \quad \Rightarrow \quad s_i \sim N(h_i, \sigma) \quad \text{with } pr = p_c$.

$$s_i = \text{absent} \quad \text{with } pr = 1 - p_c.$$

spurious cuts $\Rightarrow$ ~~Expon~~ Poisson. $\quad e^{-\lambda_f} \dfrac{\lambda_f^{F_{jk}}}{F_{jk}!}$

$D_j = $ Bad $\Rightarrow$

Poisson: $\quad e^{-\lambda_n} \dfrac{\lambda_n^{M_j}}{M_j!}$

$$Pr\left[ D_j^{(k)} \Big| \; _{good}^{H,} \right]$$

$$= \prod_{i=1}^{N} \left\{ \left[ \left( p_c \; \frac{e^{-(s_{ij} - h_i)^2/2\sigma^2}}{\sqrt{2\pi}\,\sigma_i} \right)^{m_{ijk}} (1-p_c)^{(1-m_{ijk})} \right] \right.$$

$$\times \, e^{-\lambda_f} \lambda_f^{F_{jk}}.$$

$$Pr\left[ D_j^{(k)} \Big| \; _{bad}^{H,} \right] = e^{-\lambda_n} \lambda_n^{M_j}$$