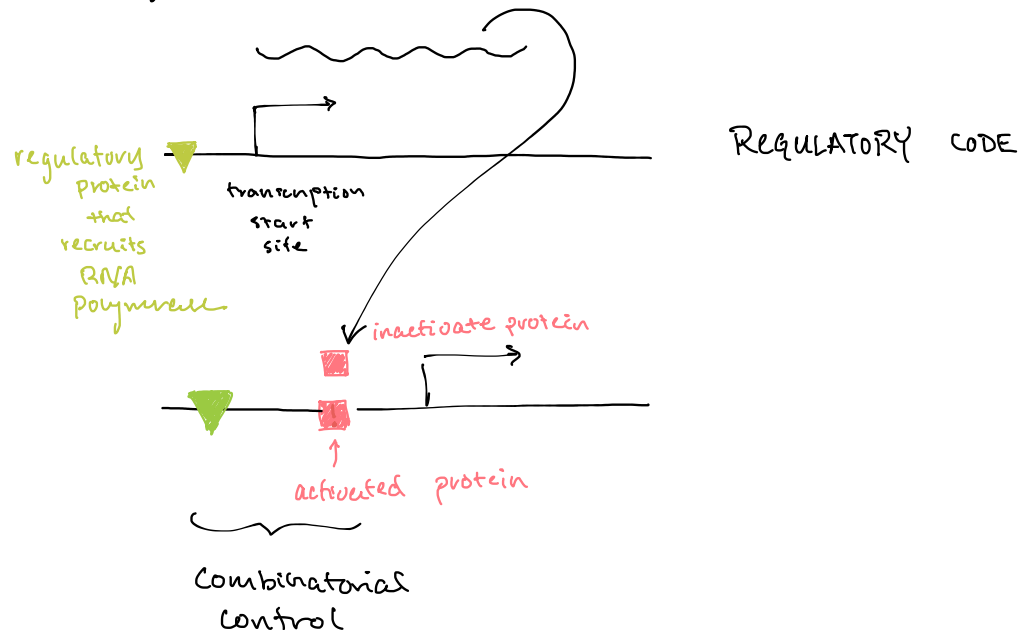


7 ChIP-seq Analysis; DNA-protein interactions

1. Overview

- * How do we study gene regulation?
- * Transcription factors regulate gene expression.

2. Gene Regulation

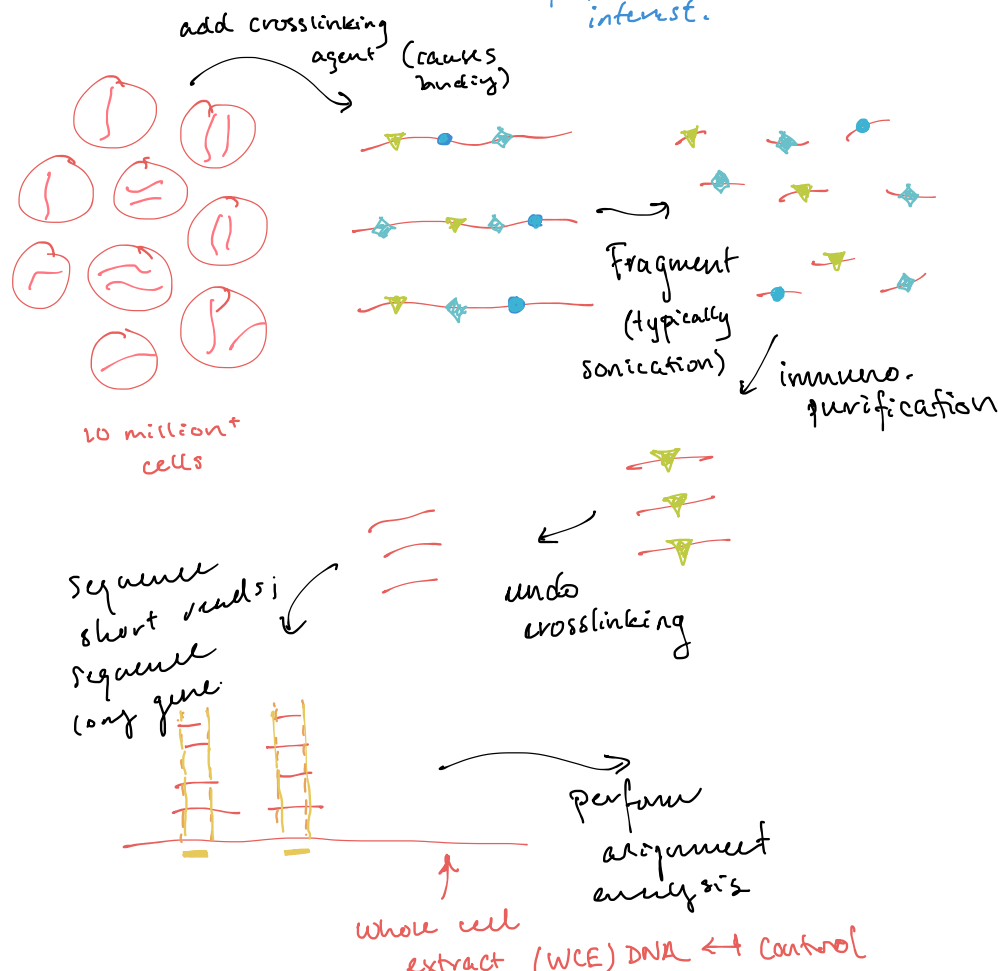


- * currently, we can only consider 1 protein at a time.
- * we can also look at more general epigenetic marks.

3. ChIP-seq [Assumes we know what the regulator is]

- * chromatin immunoprecipitation, sequencing.
- * we can see where proteins bind within 10 bp.
- * needs

- good antibody or, } for immunoprecipitation of proteins of interest.
- epitopic tag



- * We need to get to the point where we can discover regulatory proteins de novo.

- * constructed mixture model, map of spatial distribution of reads.

- * GPS [genome positioning system]

$$P(r_i | b_j) = P(r_i | z_{ij} = 1) = \exp(-l^{z_i}(r_i - b_j))$$

binding event at position b_j

$$P(R | \pi) = \prod_{i=1}^N \sum_{m=1}^M P(r_i | m) \pi_m \quad \text{subject to} \quad \sum_{m=1}^M \pi_m = 1$$

- * Goal: we want to find

$$\pi^* = \operatorname{argmax}_{\pi} P(R | \pi).$$

- * Insight: $g(z_n = m) = \begin{cases} 1 & \text{if read } n \text{ came from } m \\ 0 & \text{otherwise} \end{cases}$

$$N_m = \sum_{i=1}^N g(z_i = m) \Rightarrow \pi_m = \frac{N_m}{\sum_{m=1}^M N_m}$$

(g is a latent variable)

- * estimate g \rightarrow estimate π_m \rightarrow update estimate for g

E step

$$\gamma(z_n = m) = \frac{\pi_m P(r_n | m)}{\sum_{m'=1}^M \pi_{m'} P(r_n | m')}$$

fraction of read n assigned to event m .

$$\hat{\pi}_m^{(1)} = \frac{N_m}{\sum_{m'=1}^M N_{m'}}$$

$$N_m = \sum_{n=1}^N \gamma(z_n = m).$$

effective # reads assigned to event m

Assuming punctate (single point) binding proteins.

finds MLE for assumed $\gamma \approx g$

Add a sparse prior on π , the binding events.

$$P(\pi) \propto \prod_{m=1}^M \frac{1}{(\pi_m)^\alpha}, \quad \alpha > 0. \quad \rightarrow \text{now use MAP, } \hat{\pi}_i \text{ now}$$

negative Dirichlet prior

$$\hat{\pi}_m^{(1)} = \frac{\max(0, N_m - \alpha)}{\sum_{m'=1}^M \max(0, N_{m'} - \alpha)}$$

need α reads or eliminated (component elimination)

- * Benjamin-Hochberg Correction

$$P_i = \frac{i}{N} \cdot \alpha$$

desired FALSE DISCOVERY RATE

\sim p-values.

(replicates)

- * Run experiment twice (at least) — are the results concordant?

- * $\psi_n(t) \equiv$ fraction of events paired in the top $n \times t$ events

- rank in order of significance
- use rank correlation (gets rid of numerical dependencies, like # of reads).

fraction of events we are considering

IDR (irreproducibility discovery rate)