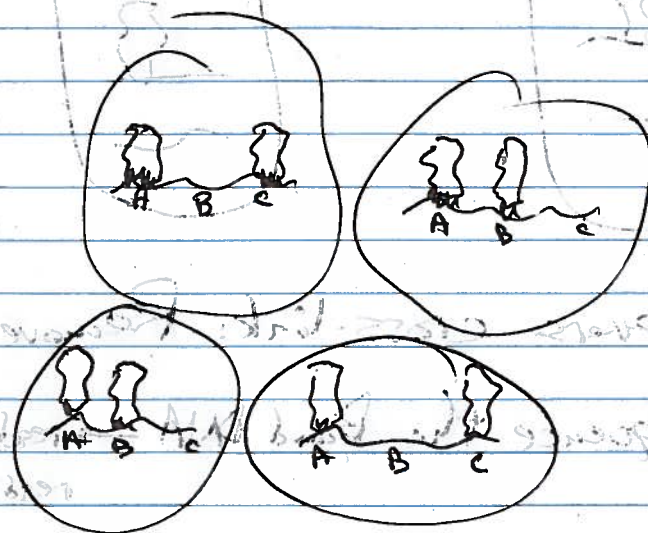


ChIP-Seq

ChIP-Seq + Motif discovery

Goal of ChIP-Seq: Identify regions of genome bound by a given TF in a given type of cells.

ChIP-Seq: Chromatin Immunoprecipitation followed by Seq.



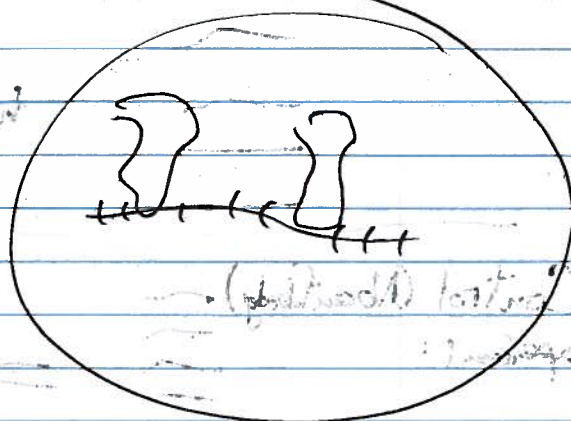
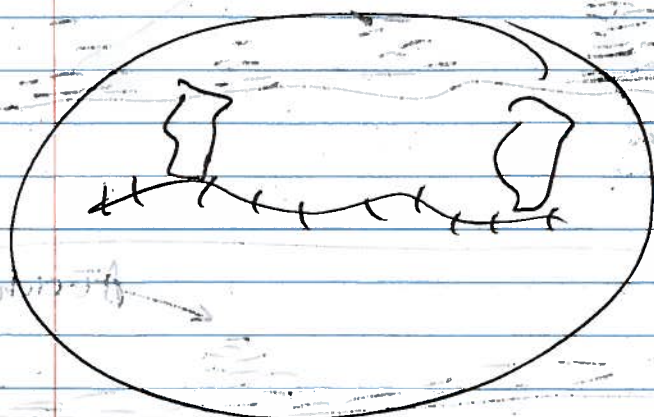
(10 Million cells)

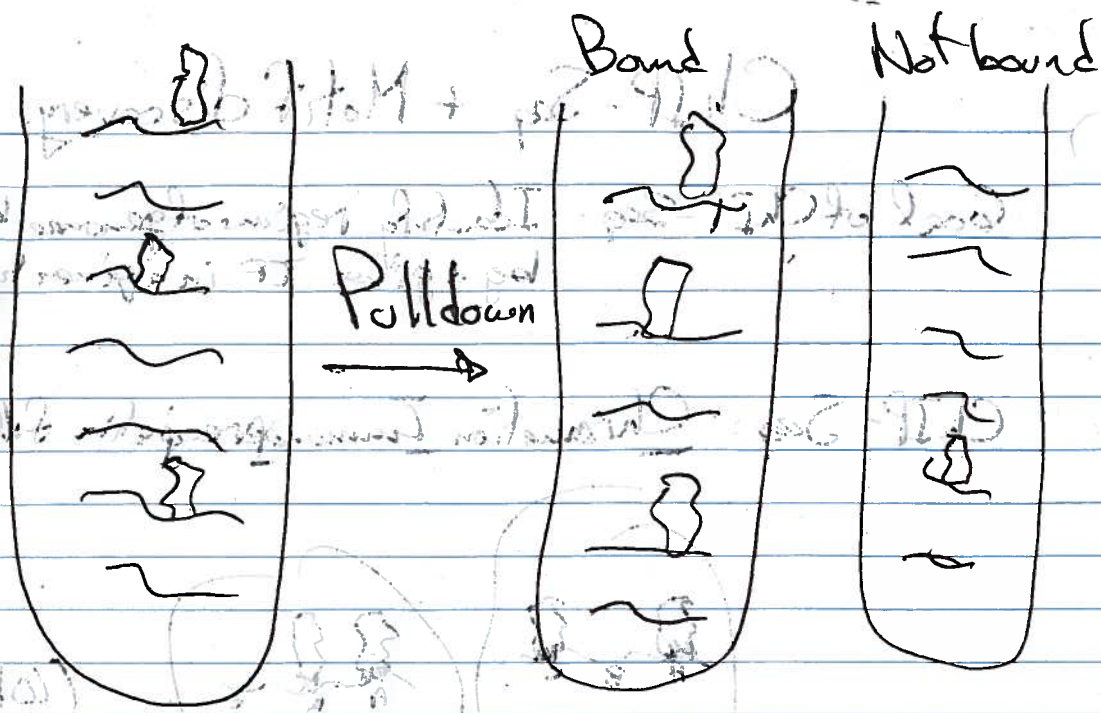
(Typically:
100-100,000 sites
occupied by TF)

① Cross-linking proteins to DNA: Strengthen bonds btw prot. and DNA

② Extract DNA with proteins attached to it

③ Fragment DNA in pieces of ~200bp (~~cut~~ Sonication)





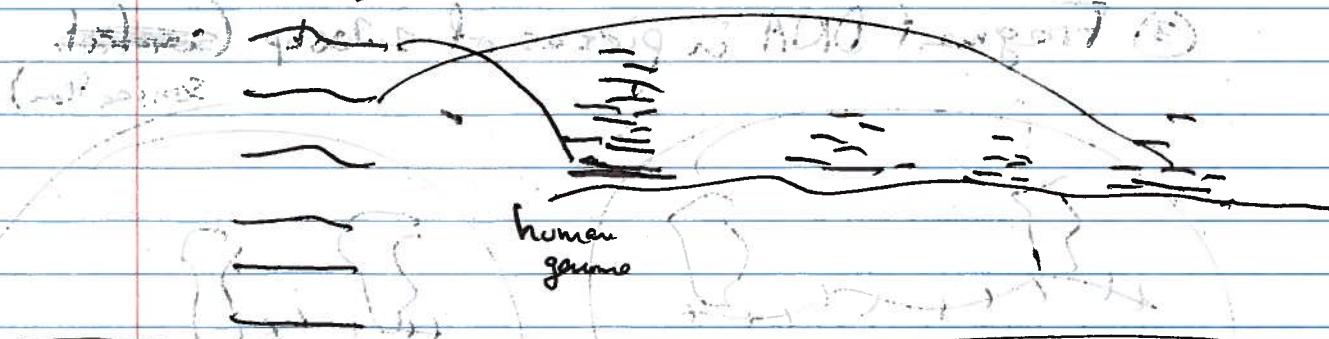
⑤: Reverse cross-link: Remove the TF from DNA

⑥: Sequence the Bound DNA → read1: ACTAGGTC
 read2: TCGTCCTA..
 ...
 read 10⁷: TAACTGAT..

⑦: ~~Read Mapping~~ Local alignment btw

reads and the human genome

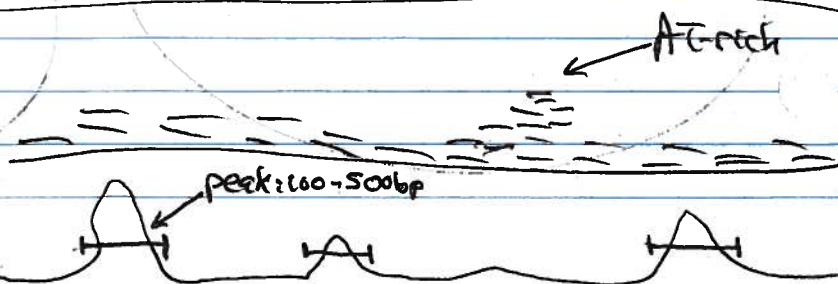
Reads



Control (No antibody)

Experiment:

Normalize: $\frac{\text{readcount (exp)}}{\text{readcount (control)}}$



ChIP-seq outcome: List of \checkmark genomic regions 100-500bp
that are believed to be bound by TF

Problem: Doesn't tell us which 6-15bp regions bound
 \Rightarrow Can't build consensus ~~or~~ PWM for TF

Motif Discovery Problem

Given: Set of sequences S_1, S_2, \dots, S_n each \rightarrow of length 100-500
believed to contain binding site for TF

Find: Consensus sequence for TF
OR
PWM for TF

Consensus sequence: Regular expression

$$w = \{A\} \{ \begin{matrix} C \\ T \end{matrix} \} \{ \begin{matrix} A \\ C \\ G \\ T \end{matrix} \} \{T\} \{ \begin{matrix} A \\ C \end{matrix} \}$$

How should we score a candidate cons. seq $w = w_1 \dots w_k$
Criteria:

- w should occur in each of S_1, \dots, S_n
 - \rightarrow too strict: Some S_i might be false positives
 - $\rightarrow \{ \begin{matrix} A \\ C \end{matrix} \} \{ \begin{matrix} A \\ C \end{matrix} \} \dots \{ \begin{matrix} A \\ C \end{matrix} \}$ matches every where

Motif enrichment ^{consensus sequence} approach

Let $N_w = \# \text{ of matches for } w \text{ in } S_1, S_2, \dots, S_n$
 $w \in \Sigma^k$ ~~is a motif~~ ^{is a motif}

$E_w = \text{Expected \# of matches of } w \text{ in}$
 $\text{a set of random sequences } R_1, R_2, \dots, R_n,$
 $\text{where } R_i \text{ has same length } S_i$

→ each nucleotide is chosen indep. with $P_A: 0.3$
 $P_C: 0.2$
 $P_G: 0.2$
 $P_T: 0.3$

How to compute E_w ? $w = w_1 w_2 \dots w_k$
 $\text{where } w_i \in \{A, C, G, T\}$

$w = \text{ACGTC}$
 $w_1 = A, w_2 = C, w_3 = G, w_4 = T, w_5 = C$

$P_R[w \text{ has a match starting at position } p \text{ in random seq}] = ?$
 $R = \dots \dots \dots$

$$\Rightarrow \prod_{i=1}^k P_R[w_i \text{ has match at position } p+i-1]$$

$$P_{\text{match}}(w) = \prod_{i=1}^k \left(\sum_{a \in \Sigma} P_a \right)$$

$$E_w = (\# \text{ positions eligible for match}) \cdot P_{\text{match}}(w)$$

$$= \sum_{i=1}^n (\text{length of } S_i - k + 1) \cdot P_{\text{match}}(w)$$

Finally, from N_w , E_w and Z_w

2. We want to find w where N_w and E_w are the most different, i.e. $N_w \gg E_w$

Z-score approach

$$Z_w = \frac{N_w - E_w}{\sqrt{E_w}}$$

Complete algo. ?

For each possible consensus seq. w
Calculate N_w , E_w , Z_w

Report word w with highest Z_w

$$P(w) = \prod_{i=1}^L \frac{1}{\sum_{a \in \Sigma} P(a)}$$

$$P(w) = \prod_{i=1}^L \frac{1}{\sum_{a \in \Sigma} P(a)}$$

$$P(w) = \prod_{i=1}^L \frac{1}{\sum_{a \in \Sigma} P(a)}$$