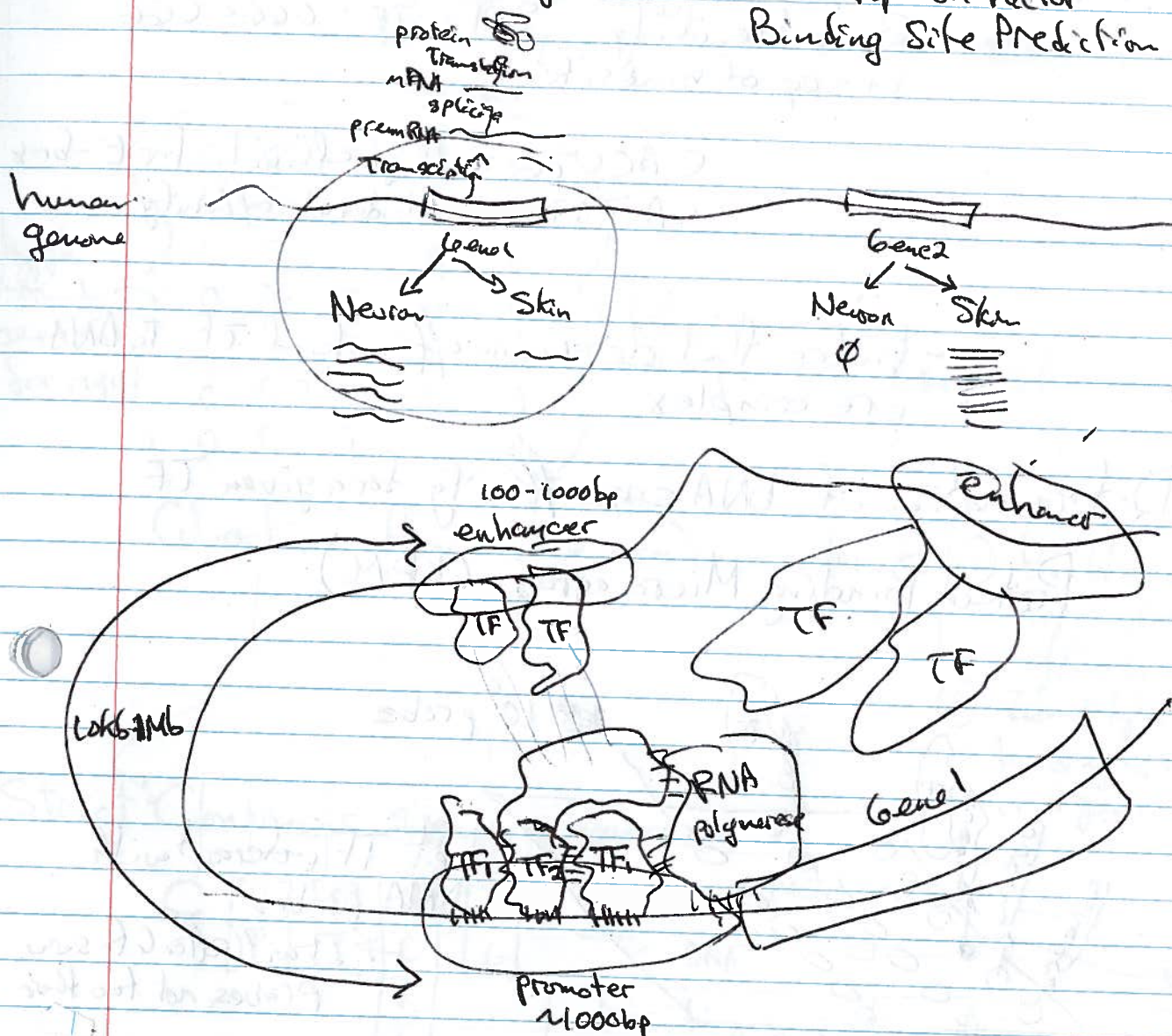


Gene Regulation + Transcription Factor Binding Site Prediction



- Transcription factors: - Proteins that
- Bind specific DNA seq.
 - Alter expression of nearby gene(s)
 - activation
 - repression
- In humans: 2000 different TFs
- ↳ each binds to different DNA sequences
- Transcription Factor Binding sites

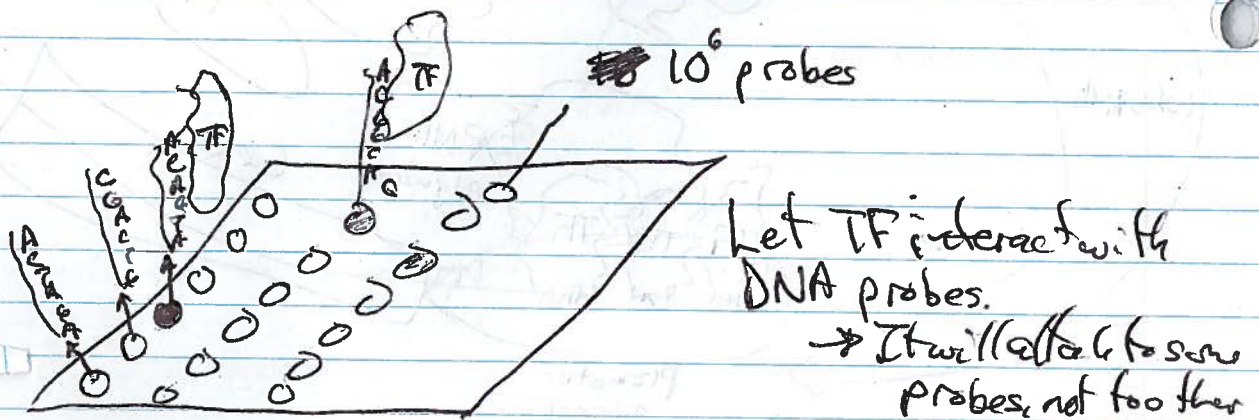
TFBS : - 6-15 bp long
 - Some flexibility in seq. of bind site

E-box TF bind : CACGTC
 SPI TF : GGGCGGG

CACGTC : High affinity for E-box
 CATGTC : Moderate affinity

- Rules that determine affinity of TF to DNA are complex

Determination of DNA seq. affinity for a given TF
 Protein Binding Micro-array (PBM)



Take picture of PBM ⇒ infer affinity of TF to each probe

⇒ intensity

ACAGTA	: 1000
ACGATA	: 250

Example: ~~Binding site~~

DNA sequences bound by TF ~~myc~~ myc

not independent

high affinity sequences for myc

C	C	T	A	A
C	T	A	C	G
C	T	T	A	G
C	T	T	G	G
C	C	T	T	A
C	C	T	C	A
C	C	T	T	A

Representative sample of all possible sites TF

Question: How can we use in order to

- 1 Build model for that TF's affinity
- 2 Identify new binding sites in a genome

Strict Consensus sequence approach

Match consensus

DNA

C	[C]	[A]	[A]	[A]
	[T]	[T]	C	G
			G	
			T	

or

C	[C]	[T]	[A]	[A]
	[T]		C	G
			C	
			T	

Position weight matrices

	1	2	3	4	5
A	0	0	1/7	2/7	4/7
C	1	4/7	0	2/7	0
G	0	0	0	1/7	3/7
T	0	3/7	6/7	2/7	0

Candidate sequence
score

$$1 \times \frac{3}{7} \times \frac{4}{7} \times \frac{1}{7} \times \frac{2}{7} = \frac{9}{7^4} = 0.00...$$

How to identify matches for a given PWM in a given sequence

$S =$ A C A G T C A C T T

For each starting position i in S

Calculate score of $S[i, i+1, \dots, i+5-i]$ on PWM

If score > Threshold : then predict Binding
other no binding