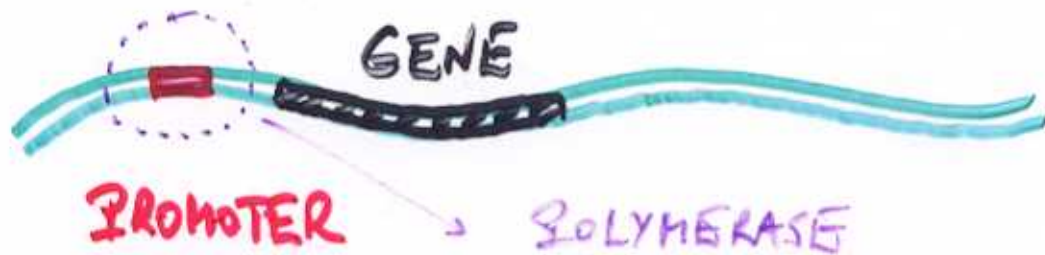


TRANSCRIPTION SIGNALS AND PROMOTER RECOGNITION

Q1: WHAT'S THE PROBLEM?



Q2: WHY CARE?



1. THEORETICAL:

- UNDERSTAND GENE EXPRESSION
- HOW DOES THE CELL FIND GENES
- GENE REGULATION, CANCER, ...

2. PRACTICAL:

- GENE FINDING
- PREDICT GENE EXPRESSION
- ...

Q3: WHAT DO WE DO?

OUTLINE

I. MOLECULAR BIOLOGY (short)

II. METHODS (BIO INFORMATICS)

- CONSENSUS SEQUENCES
- REGULAR EXPRESSIONS
- POSITION WEIGHT MATRICES (PWM)
- HOW TO CONSTRUCT PWMs FROM SET OF ALIGNED SEQUENCES
- HOW TO ALIGN SEQUENCES TO PWM
- HOW TO FIND PWM AND ALIGNMENT FROM SET OF UN-ALIGNED SEQUENCES
 - DEFINITION OF GIBBS SAMPLER
 - PERFORMANCE OF GIBBS SAMPLER
- OUTLOOK

Q4: WHERE DO PROTEINS COME FROM?

→ TRANSCRIPTION

→ TRANSLATION

• PROMOTER = DNA BINDING SITE

FOR DOCKING / ASSEMBLING TRANSCRIPT.

MACHINERY

↳ RNA POLYMERASE AND
OTHER FACTORS / COFACTORS

• -55 ... +20

• UPSTREAM ... -1 +1 ... DOWNSTREAM

• -10 SIGNAL ("TATA")

• -35 SIGNAL

• 7000 RNA POL. COPIES / CELL

• 40 nt/sec (20 × SLOWER THAN REPLIC.)

~ 15 aa/sec (TRANSLATION)

Q5: HOW DOES THE TRANSCRIPTION MACHINERY
RECOGNIZE PROMOTERS ?

AS: ?

Q6: HOW CAN WE RECOGNIZE PROMOTERS ?

AG: 1. EXPERIMENTALLY

- TIME, \$\$\$
- DO NOT REALLY UNDERSTAND
BINDING OF σ TO PROMOTER

2. THEORETICALLY (BIOPHYSICS)

- MODEL ALL MOLECULES (50 - 150)
- COMPUTE BINDING ENERGIES
- SOLVE EQUATIONS OF MOTION

→ CONFORMATION / BINDING

- WRONG POTENTIALS

- TOO COMPLEX !

3. THEORETICALLY (BIOINFORMATICS)

- SIMPLIFY MODELS TREMENDOUSLY
- NOT BIOPHYS. MODELS ANYMORE
- "SUPER" ABSTRACT
- ADVANTAGE: THEY WORK!

Q7: WHAT ARE CONSENSUS SEQUENCES ?

... A A T C A T ...

... T A C A T T ...

... T A C A T T ...

... T A T G A T ...

... T A T T A T ...

T A T A A T

CONSENSUS SEQUENCE = SEQUENCE
OF THE MOST FREQUENT NUCLEOTIDE
AT EACH POSITION.

CONSENSUS SEQ \neq MOST FREQUENT SEQ.

CONSENSUS SEQ DOES OFTEN NOT OCCUR
AT ALL.

SHORTCOMINGS OF CONSENSUS SEQ.

1. ONLY INFORMATION ABOUT **MOST-FREQUENT**, NOT **LEAST-FREQUENT** NT.

2. NOT : **HOW** FREQUENT

→ **REGULAR EXPRESSIONS** (solve 1, not 2)

[AT]A[CT]*[AT]

→ **GENERALIZED CONSENSUS SEQ** (solve 2, not 1)

T₂₀ A₉₅ T₄₅ A₆₀ A₅₀ T₉₆ → **FREQUENCY OF NT IN %**

→ **POSITION WEIGHT MATRICES (PWM)**

SOLVE BOTH 1 AND 2

Q8: WHAT ARE POSITION WEIGHT MATRICES?

P_{ie} = PROBABILITY OF FINDING
NUCLEOTIDE i IN
POSITION e

i	{	1	A	0.2	1.0	0	0.4	0.6	0
		2	C	0	0	0.4	0.2	0	0
		3	G	0	0	0	0.2	0	0
		I=4	T	0.2	0	0.6	0.2	0.4	1.0
				1	2	.	.	.	L=6
				}					e

IWM CONTAINS ALL INFORMATION ABOUT
ONSENSUS SEQ, GENERALIZED CONS. SEQ,
REGULAR EXPRESS.

Q9: DO PWMs HAVE AN INTUITIVE MEANING?

PWM $\hat{=}$ "SUPER" SIMPLE MODEL
(ISING MODEL, POTTS MODEL)

$E_{i,l}$ = ENERGY IF POLYMERASE BINDS TO
NUCLEOTIDE i AT POSITION l

$\tilde{E}_{i,l}$ = ENERGY IF ... DOES NOT BIND TO ...

$$\tilde{E} = \tilde{E}_{i_1 l} + \tilde{E}_{i_2 l} + \dots + \tilde{E}_{i_l l}$$

$$\Delta E_{i,l} = E_{i,l} - \tilde{E}_{i,l} \quad \text{BINDING ENERGY}$$

$$\Delta E = E - \tilde{E} \quad \text{TOTAL BIND. ENERGY}$$

Q10 : HOW TO OBTAIN BWH FROM SET
OF **ALIGNED SEQUENCES** ?

$$\hat{p}_i = \frac{\sum x_i}{N}$$

MAXIMUM LIKELIHOOD ESTIMATOR

$$\hat{p}_i = \frac{x_i + 1}{N + 1}$$

MINIMUM VARIANCE ESTIMATOR
(BAYES)

PSEUDO COUNTS

Q12: HOW TO ALIGN A SEQUENCE /
A SET OF SEQUENCES TO A PWM?

GIVEN: $P_{i,j}$, Q_j , SEQUENCE

$$\text{SCORE} = \frac{P_{i,1}}{Q_1} \cdot \frac{P_{i,2}}{Q_2} \cdot \dots \cdot \frac{P_{i,L}}{Q_L}$$

... ACGTACCGTTTACCGT...



SCORE

- COMPUTE SCORE FOR EACH POSITION OF THE SEQUENCE
- CHOOSE THAT POSITION FOR WHICH SCORE = MAX.

Q13 : WHAT IF NEITHER PWH
NOR ALIGNMENT IS KNOWN ?

- SEARCH FOR PATTERN, BUT DON'T KNOW WHICH PATTERN AND WHERE TO LOOK !
- WE ARE GIVEN A SET OF UN-ALIGNED SEQUENCES AND ARE TO SEARCH FOR PATTERNS UNKNOWN
- CONSTRUCT PWH AND ALIGNMENT SIMULTANEOUSLY

Q: CAN WE SOLVE THIS PROBLEM ?

[REDACTED]

gtctcggatcaccctccctagacgataagggcccgcacttgtagtggcgctctagcagat
cctccctagacgataagggcccgcacttgtactcgctgcttcaccgggtcttatcgattt
gtcacatgatactccctagacgataagggcccgcacttgtatactttcagcccaaacga
accaatcaagatacactccctagacgataagggcccgcacttgtattctgtcctgtgggg
tgacgacgctccctagacgataagggcccgcacttgtaatagctccagcactacgctcgt
gatcgttctatgtcgcttatagacgcctccctagacgataagggcccgcacttgtaggtc
ctccctagacgataagggcccgcacttgtacccaatgagaatgttccccgcctaccgctc
ctctccctagacgataagggcccgcacttgtatctactctcgtccaaccagagccgatgt
ggtgctccctagacgataagggcccgcacttgtatcgagagctcctgacctgtctaagtt
ctcagtatgccctccctagacgataagggcccgcacttgtaccaggagttacttctcaga

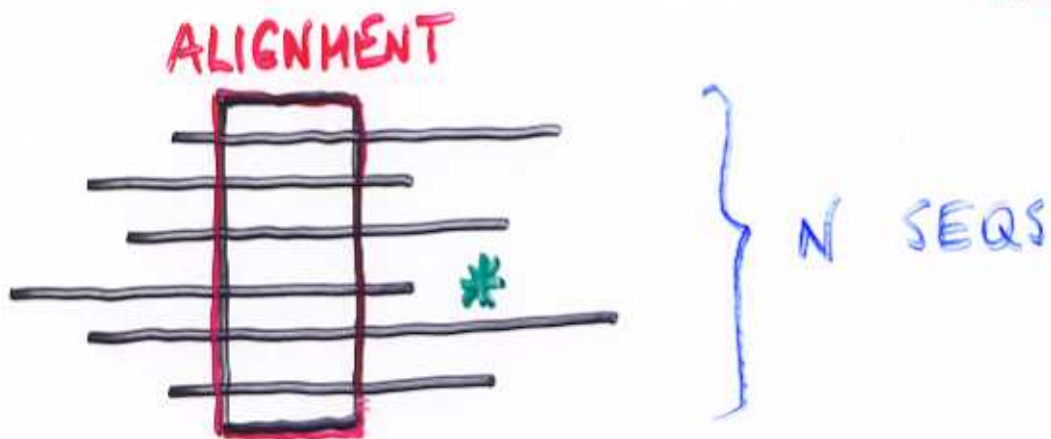
gtctcggatcaccaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaagtggcgctctagcagat
caaaaaaaaaaaaaaaaaaaaaaaaaaaaaaactcgctgcttcaccgggtcttatcgattt
gtcacatgataaaaaaaaaaaaaaaaaaaaaaaaaaaaaaatactttcagcccaaacga
accaatcaagatacaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaattctgtcctgtgggg
tgacgacgaaaaaaaaaaaaaaaaaaaaaaaaaaaaaataagctccagcactacgctcgt
gatcgttctatgtcgcttatagacgcaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaggtc
aaaaaaaaaaaaaaaaaaaaaaaaaaaaaacccaatgagaatgttccccgcctaccgctc
ctaaaaaaaaaaaaaaaaaaaaaaaaaaaaaatactactctcgtccaaccagagccgatgt
ggtgaaaaaaaaaaaaaaaaaaaaaaaaaaaaaatacgagagctcctgacctgtctaagtt
ctcagtatgcccaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaccaggagttacttctcaga

gtctcggatcaccCTCCCTAGACGATAAGGGCCCGCACTTGTAggtggcgctctagcagat
cCTCCCTAGACGATAAGGGCCCGCACTTGTActcgctgcttcaccgggtcttatcgattt
gtcacatgataCTCCCTAGACGATAAGGGCCCGCACTTGTAtactttcagcccaaacga
accaatcaagatacaCTCCCTAGACGATAAGGGCCCGCACTTGTAttctgtcctgtgggg
tgacgacgCTCCCTAGACGATAAGGGCCCGCACTTGTAAatagctccagcactacgctcgt
gatcgttctatgtcgcttatagacgcCTCCCTAGACGATAAGGGCCCGCACTTGTAggtc
CTCCCTAGACGATAAGGGCCCGCACTTGTAcccaatgagaatgttccccgcctaccgctc
ctCTCCCTAGACGATAAGGGCCCGCACTTGTAtctactctcgtccaaccagagccgatgt
ggtgCTCCCTAGACGATAAGGGCCCGCACTTGTAtcgagagctcctgacctgtctaagtt
ctcagtatgccCTCCCTAGACGATAAGGGCCCGCACTTGTAccaggagttacttctcaga

agactcttcattgaggaccatggtcacgcccgggaatagcagatccctcccgtatgactc
ttgaatctgtccagtcctcgagagctatggtcacgcccgggaatagcagatccctcgtgg
tgtagccgagaccaacctgctctcatctatggtcacgcccgggaatagcagatccctctc
ctcgccatccgccccttgaagagtaacctatggtcacgcccgggaatagcagatccctc
ctggatggtcacgcccgggaatagcagatccctccgcagatattcgctgtatcttgctag
tgctcgcacacgacctcgataatggtcacgcccgggaatagcagatccctcgcagcagt
ggggtgctcctgtcttatggtcacgcccgggaatagcagatccctcacatagaactaacca
agcaaaacccgactttcatatggtcacgcccgggaatagcagatccctcatagtacactg
tttagctattctgggcccacttcgtcgctcatggtcacgcccgggaatagcagatccctcc
tagacgatctcggggtgatggtcacgcccgggaatagcagatccctcccactaggtctcgt

Q14: HOW TO FIND PWM AND ALIGNMENT
SIMULTANEOUSLY?

IDEA: RANDOM ALIGNMENT PWM
ITERATE!



1. CHOOSE 1 SEQ RANDOMLY (*)
TAKE OUT *
2. COMPUTE PWM FROM N-1 SEQS
3. ALIGN SEQ * TO PWM
4. GOTO 1.

LAWRENCE, C.E., et al. SCIENCE 266: 208-214
(1993)

COMMENT:

(3.) IS NOT DETERMINISTIC

BUT STOCHASTIC

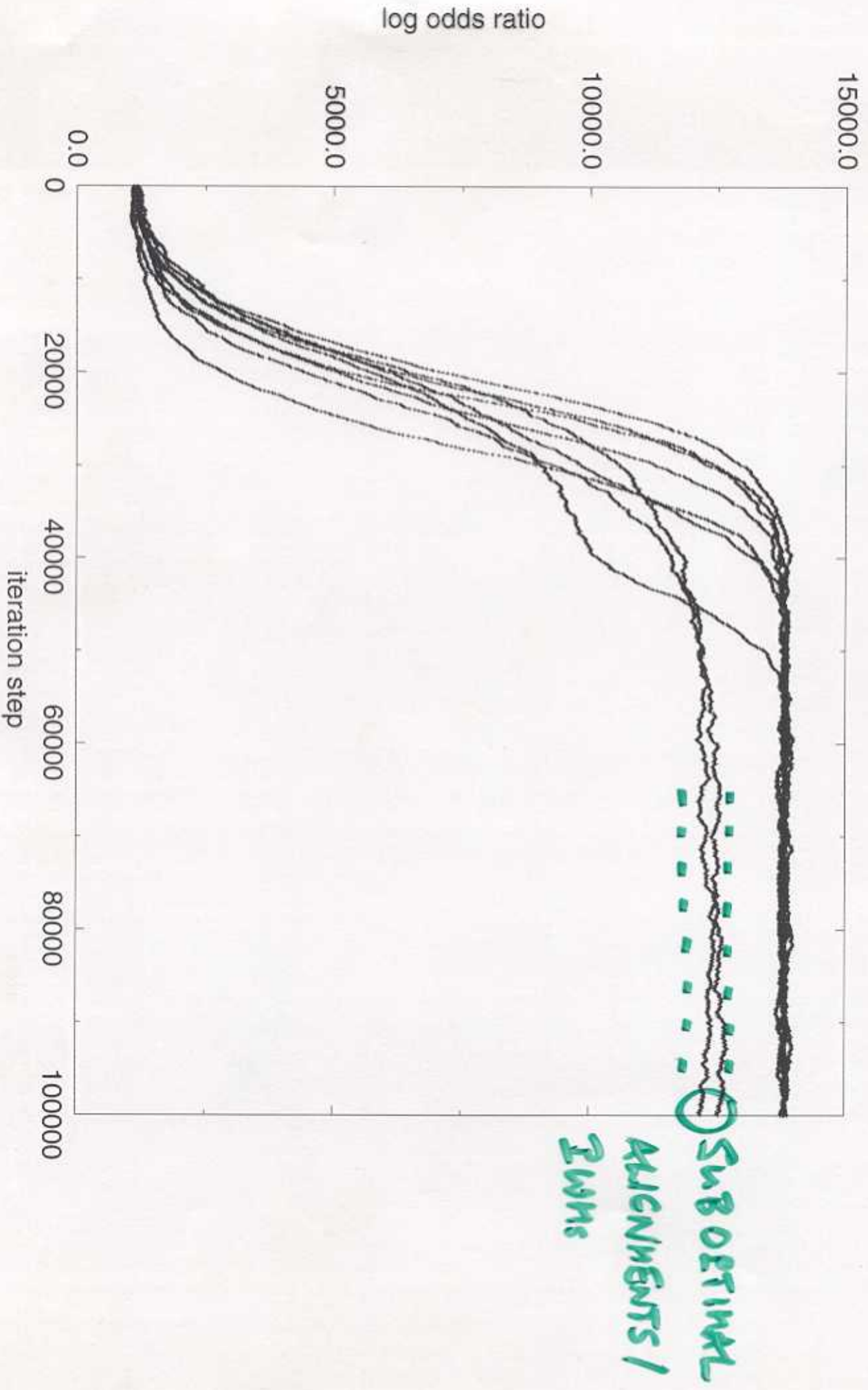
3.A: COMPUTE AT EACH POSITION

$$\text{SCORE} = \prod_{l=1}^L \frac{P_{i_l}}{Q_{i_l}}$$

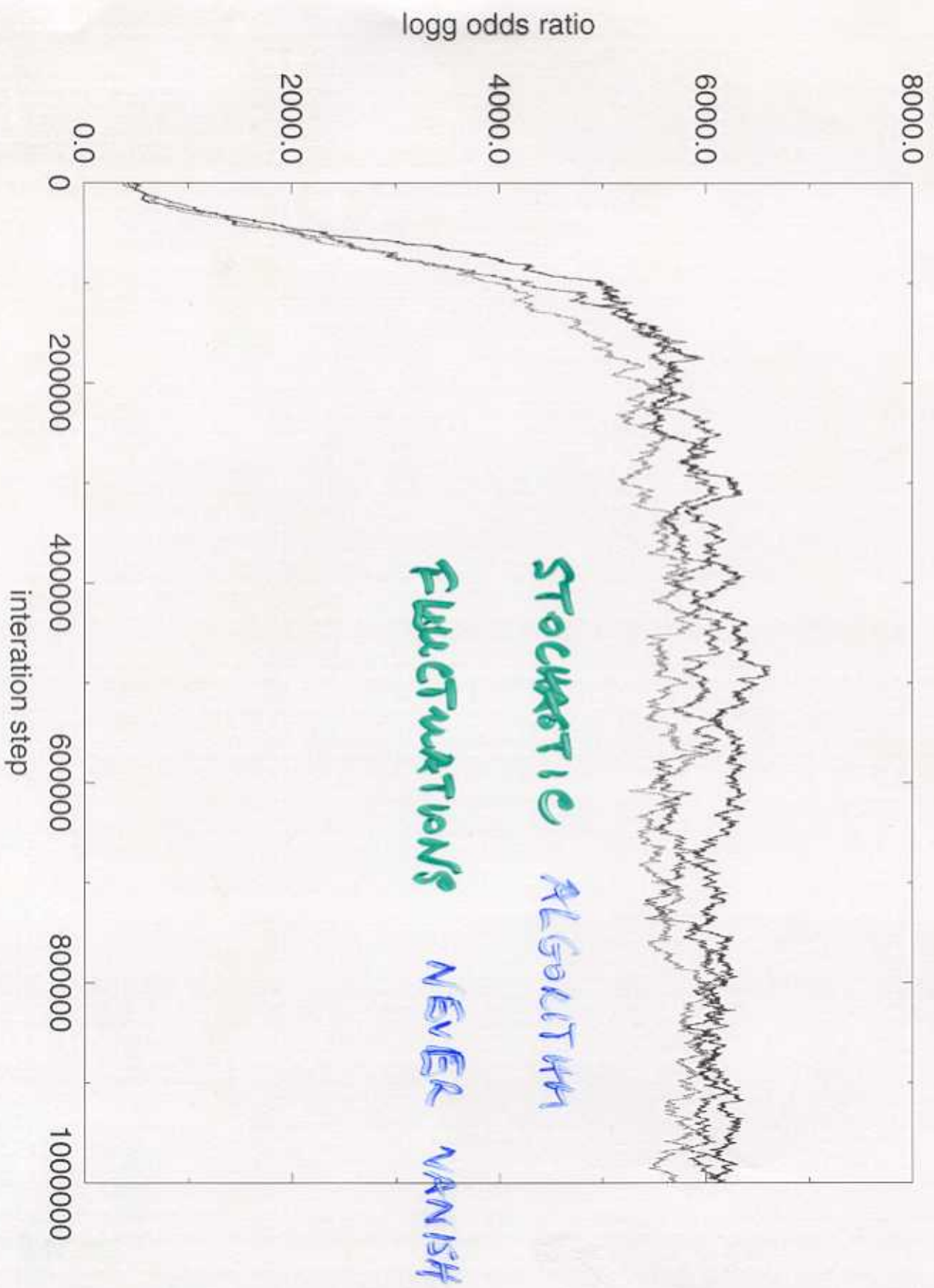
3.B: CHOOSE POSITION RANDOMLY
FROM DISTRIBUTION $\{\text{SCORE}(\text{POS})\}$

(DO NOT CHOOSE POSITION
AT WHICH $\text{SCORE} = \text{MAX.}$)

Convergence of Gibbs Sampler



Convergence of Gibbs Sampler



Q15: LIMITATIONS OF PWMs ?

1. HETEROGENEITY

- RNA POLYMERASE VERY COMPLEX MOLECULE (NOT ISING)
- RNA POL. BINDS TO MANY FACTORS

2. CORRELATIONS

- NEAREST NEIGHBORS
- LONG-RANGE CORR.

Q16: WHAT CAN WE DO ?

- PWM WITH DIMERS, TRIMERS, ...
- BETTER MODELS !!!
- GENE EXPRESSION DATA ANALYSIS !!!

Q₀₀ : WHAT DO WE UNDERSTAND NOW
THAT WE DID NOT PREVIOUSLY UNDERSTAND?

- MOL ~~BIO~~L., BIO ~~PHYS.~~, BIO INFO !!!
- PWM (NT ARE INDEPENDENT $\hat{=}$
E IS ADDITIVE)
- PWM FROM UNALIGNED SEQs
(GIBBS SAMPLER)
- PWM SIMPLE, BUT AS GOOD AS MOST-
SOPHISTICATED METHODS
- NEVERTHELESS: LIMITATIONS! (WE ARE
MISSING ESSENTIAL INFORMATION)
- HOPE FOR BETTER MODELS IN THE FUTURE
- GENE FINDING CANNOT YET RELY
ON PROMOTER RECOGNITION