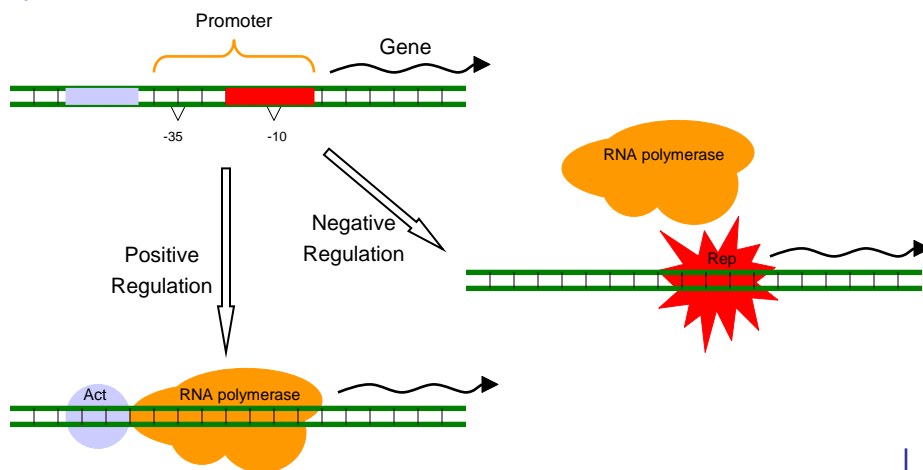


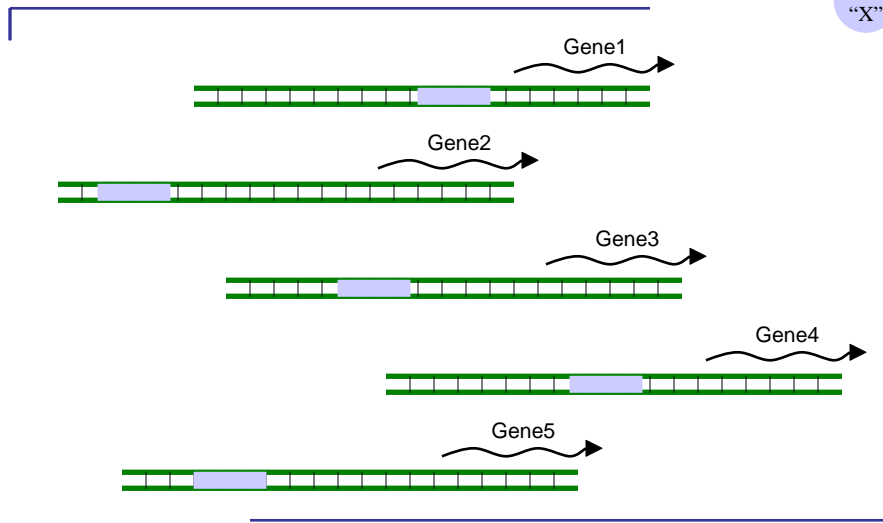
Regulatory Motifs

Gene Regulation

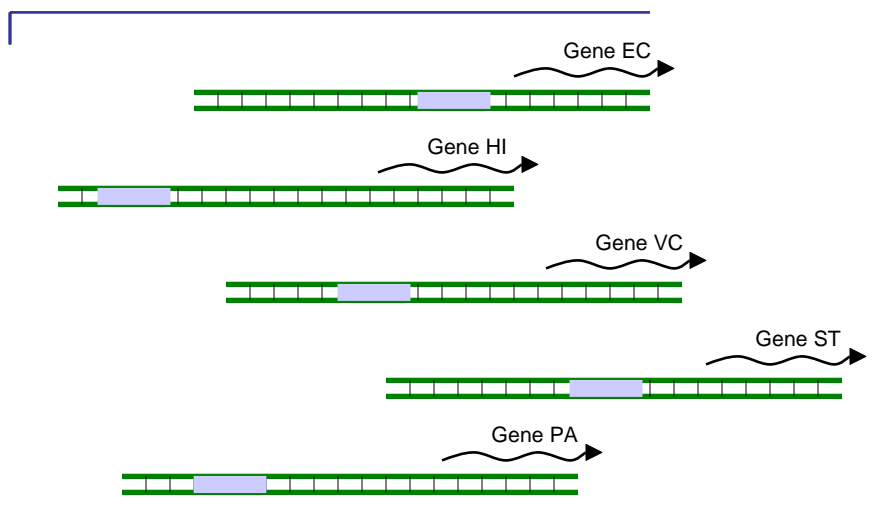


What if we believed that a number of genes were regulated by the same transcription factor?

TF
"X"



What if we believed that a number of genes were orthologous?



How do we search upstream sequences for instances of a motif?

> Escherichia coli
TTGATTCCCTGAATGCCCGCTTAGTGTAACTACTGTAACCGGCATTTTCTGCTTTTCC
TGCCGATATTTTTCTTATCTACCTCACAAAGTTAGCAATAACTGCTGGGAAAATTCCG
AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
ATTTGTGGCGTTTTCAAATGCTTGGAGAATGGGTACATGATCCGCTTGGCATTGAAGGT
GAGGCTGGCAGCAGCGAAGGCTGGGGCTGTTGAACGTTACACGAGTGAACCGCGAA
CCATGTTGACACGAATTCTG

> Salmonella typhi
GGTCGGCTTAGACTAGTGTGACCAAAAAGCTTTTGCTGAAGTTTCAGGGTAAGAAGAACC
AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCCGCATGCCGGCGCCGC
TGCTGGAACGCTCGACCCAGGGCTACACTAGTTTAAACCGGAACGCCGCAAGTGGATCG
GCCTGCCCCAGCTATTGCTC

If we knew *where* the motif instances were located in each sequence...

> Escherichia coli
TTGATTCCCTGAATGCCCGCTTAGT**GTAACACTACTGTAAC**CGGCATTTTCTGCTTTTCC
TGCCGATATTTTTCTTATCTACCTCACAAAGTTAGCAATAACTGCTGGGAAAATTCCG
AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
ATTTGTGGCGTTTTCAAATGCTTGGAGAATGGGTACATGATCCGCTTGGCATTGAAGGT
GAGGCTGGCAGCAGCGAAGGCTGGGGCTGTTGAAC**GTTACACGAGTGAAC**CGCGAA
CCATGTTGACACGAATTCTG

> Salmonella typhi
GGTCGG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGGTAAGAAGAACC
AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCCGCATGCCGGCGCCGC
TGCTGGAACGCTCGACCCAG**GCTACACTAGTTTAAAC**CGGAACGCCGCAAGTGGATCG
GCCTGCCCCAGCTATTGCTC

Then we could determine a motif model!

GTAACACTACTGTAAC

GTTACACTAGTGGGAC

GTTACACGAGTGTAAAC

CTTAGACTAGTGTGAC

GCTACACTAGTTTAAAC

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0	0.0

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

Consensus Sequence

But we don't know the locations of the motif instances...

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAACCGGCATTTTCTGCTTTTCC
 TGCCGATATTTTTTCTTATCTACCTCACAAAGTTAGCAATAACTGCTGGGAAAATTCG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGCCTATGGAATCTTGTAAGAATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTCCTCCGCCATT
 AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCGTTTTCAAATGCTTGGAGAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGTTGGGGCTGTTTGAACGTTACACGAGTGAACCGCGGAA
 CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCCGCTTAGACTAGTGTGACCAAAAAGCTTTTGTGTAAGTTTCAGGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGCTCAGCGCAAAGCCGACCCGACA
 AAACGCAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGTCCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGCGCGCC
 TGCTGGAACGCCCTCGACCCAGGGCTACACTAGTTTAAACCGAAACGCCGCAAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

What if we knew the motif model...

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0	0.0

We could determine the location of the motif instance which best matches the model...

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0	0.0

$$\text{Score} = 0.0 * .80 * 0.0 * 1.0 * 0.0 * 0.0 * 1.0 * 0.0 * 0.0 * 0.0 * 0.0 * 0.0 * 0.0 * 0.0 * 0.0 * 1.0$$

$$\text{Score} = 0.01 * .80 * 0.01 * 1.0 * 0.01 * 0.01 * 1.0 * 0.01 * 0.01 * 0.01 * 0.01 * 0.01 * 0.01 * 0.01 * 0.01 * 1.0$$

$$\text{Score} = 8.0 * 10^{-27}$$

TTGATTCCCTGAATGCCCGCTTAGTGTAACTACTGTAA

We could determine the location of the motif instance which best matches the model...

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0

Score = $0.0 * 0.0 * .20 * 0.0 * 0.0 * 0.0 * 1.0 * 0.0 * 0.0 * .80 * 0.0 * 0.0 * .80 * .40 * 0.0 * 1.0$
 Score = $0.01 * 0.01 * .20 * 0.01 * 0.01 * 0.01 * 1.0 * 0.01 * 0.01 * .80 * 0.01 * 0.01 * .80 * .40 * 0.01 * 1.0$
 Score = $5.12 * 10^{-22}$

TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAA

We could determine the location of the motif instance which best matches the model...

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0

Score = $7.16 * 10^{-28}$

TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAA

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0	0.0

```

> Escherichia coli
TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAACCGGCATTTTCTGCTTTTCC
TGCCGATATTTTTTCTTATCTACCTCACAAAGTTAGCAATAACTGCTGGGAAAATTCG
AGTTAGTCGTTATATCTAT

> Haemophilus influenzae
ATCTAACGGTACGGATTCTCCAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
ATTTGTGGCGTTTTCAAATGCTTGAGAAATGGGTACATGATCCGCTTGGCATTGAAGGT
GAGGCTGGCAGCAGCGAAGGTCTGGGGCTGTTGAACGTTACACGAGTGTAAACCGCCGAA
CCATGTTGACACGAATTCTG

> Salmonella typhi
GGTCGGCTTAGACTAGTGTGACCAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
TGCTGGAACGCCTCGACCCAGGGCTACACTAGTTTAAACCGGAACGCCAGTGGATCG
GCCTGCCCCAGCTATTGCTC

```

Expectation-Maximization (EM)

- Randomly guess the locations of each motif instance
- Repeat until convergence
 - Calculate a new motif model from the motif instances
 - Calculate new locations for the motif instances from the motif model

EM - Randomly guess the locations of each motif instance

```

> Escherichia coli
TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAACCGGCATTTTCTGCTTTTCC
TGCCGATATTTTTTCTTATCTACCTCACAAAGGTTAGCAATAACTGCTGGGAAAATTCGG
AGTTAGTCGTTTATATTCTAT

> Haemophilus influenzae
ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
ATTTGTGGCGTTTTCAAATGCTTGGAGAATGGGTACATGATCCGCTTGGCATTGAAGGT
GAGGCTGGCAGCAGCGAAGGTCTGGGGCTGTTGAACGTTACACGAGTGAACCGCCGAA
CCATGTTGACACGAATTCTG

> Salmonella typhi
GGTCGGCTTAGACTAGTGTGACCAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTACAGCGCAAAGCCGACCCGACA
AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
TGCTGGAACGCTCGACCCAGGGCTACACTAGTTTAAACCGAAACGCCGCAAGTGATCG
GCCTGCCCCAGCTATTGCTC
  
```

A	.40	.20	0.0	.20	.40	0.0	.20	.20	.40	.40	.20	.60	.20	.20	0.0	0.0
C	.20	.40	0.0	0.0	.40	.60	.20	.40	0.0	.40	.20	0.0	.20	0.0	.20	.40
G	0.0	.20	.40	.40	0.0	.40	.40	.40	.40	0.0	.20	.40	.60	.20	.40	.40
T	.40	.20	.60	.40	.20	0.0	.20	0.0	.20	.20	.40	0.0	0.0	.60	.40	.20

```

> Escherichia coli
TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAACCGGCATTTTCTGCTTTTCC
TGCCGATATTTTTTCTTATCTACCTCACAAAGGTTAGCAATAACTGCTGGGAAAATTCGG
AGTTAGTCGTTTATATTCTAT

> Haemophilus influenzae
ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
ATTTGTGGCGTTTTCAAATGCTTGGAGAATGGGTACATGATCCGCTTGGCATTGAAGGT
GAGGCTGGCAGCAGCGAAGGTCTGGGGCTGTTGAACGTTACACGAGTGAACCGCCGAA
CCATGTTGACACGAATTCTG

> Salmonella typhi
GGTCGGCTTAGACTAGTGTGACCAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTACAGCGCAAAGCCGACCCGACA
AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
TGCTGGAACGCTCGACCCAGGGCTACACTAGTTTAAACCGAAACGCCGCAAGTGATCG
GCCTGCCCCAGCTATTGCTC
  
```

A	.40	.20	0.0	.20	.40	0.0	.20	.20	.40	.40	.20	.60	.20	.20	0.0	0.0
C	.20	.40	0.0	0.0	.40	.60	.20	.40	0.0	.40	.20	0.0	.20	0.0	.20	.40
G	0.0	.20	.40	.40	0.0	.40	.40	.40	.40	0.0	.20	.40	.60	.20	.40	.40
T	.40	.20	.60	.40	.20	0.0	.20	0.0	.20	.20	.40	0.0	0.0	.60	.40	.20

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGTGTAACTACTGTAACCGCATTTCCTGCTTTCC
 TGCCGATATTTTTCTTATCTACCTCAC**AAAGGTAGCAATAAC**TGCTGGGAAAATTCG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAATATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCG**GCATAAGTAGTAAGT**GCCGTGCCCTCCGCCATT
 AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCGTTTTCAAATGCTTGAGAAATGGGTACATGATCCGCTTGGCATTGAAGT
 GAGGCTGGCAGCAGCGAAGGCTGGGGCTGTTGAACGTTACACG**AGTGTAAACCGCGAA**
CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCCGCTTAGACTAGTGTGACCAAAAAGCTTTTGCTGAA**GTTTCAGGGTAAGAAG**AACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGGCGCGCCGCGCATGCCGGCGCCG
 TGCTGGAACGCTCGACCCAGGGCTACACTA**GTTTAAACCGAACGCC**GCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	.40	.20	.40	0.0	.40	.80	.20	.20	.20	0.0	.60	.80	.20	.60	.60	0.0
C	0.0	.20	0.0	0.0	.20	0.0	.20	.40	.20	.20	.20	.20	.20	0.0	.20	.60
G	.60	.20	0.0	.40	.20	0.0	.40	.20	.60	.60	0.0	0.0	.40	.20	0.0	.40
T	0.0	.40	.60	.60	.20	.20	.20	.20	0.0	.20	.20	0.0	.20	0.0	.20	0.0

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGTGTAACTACTGTAACCGCATTTCCTGCTTTCC
 TGCCGATATTTTTCTTATCTACCTCAC**AAAGGTAGCAATAAC**TGCTGGGAAAATTCG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAATATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCG**GCATAAGTAGTAAGT**GCCGTGCCCTCCGCCATT
 AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCGTTTTCAAATGCTTGAGAAATGGGTACATGATCCGCTTGGCATTGAAGT
 GAGGCTGGCAGCAGCGAAGGCTGGGGCTGTTGAACGTTACACG**AGTGTAAACCGCGAA**
CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCCGCTTAGACTAGTGTGACCAAAAAGCTTTTGCTGAA**GTTTCAGGGTAAGAAG**AACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGGCGCGCCGCGCATGCCGGCGCCG
 TGCTGGAACGCTCGACCCAGGGCTACACTA**GTTTAAACCGAACGCC**GCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	.40	.20	.40	0.0	.40	.80	.20	.20	.20	0.0	.60	.80	.20	.60	.60	0.0
C	0.0	.20	0.0	0.0	.20	0.0	.20	.40	.20	.20	.20	.20	.20	0.0	.20	.60
G	.60	.20	0.0	.40	.20	0.0	.40	.20	.60	.60	0.0	0.0	.40	.20	0.0	.40
T	0.0	.40	.60	.60	.20	.20	.20	.20	0.0	.20	.20	0.0	.20	0.0	.20	0.0

> *Escherichia coli*
 TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAACCGGCATTTTCTGCTTTTCC
 TGCCGAT**ATTTTTCTTATCTAC**CTCACAAAGGTTAGCAATAACTGCTGGGAAAATTCGG
 AGTTAGTCGTTATATTCTAT

> *Haemophilus influenzae*
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGAATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTCCGTGCCCTCCGCCATT
AGTTACACTAGTGGGACACC

> *Vibrio cholerae*
 ATTTGTGGCG**TTTTCAAATGCTTGG**AGAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGTTCTGGGGCTGTTTGAACGTTACACGAGTGAACCGCGGAA
 CCATGTTGACACGAATTCTG

> *Salmonella typhi*
 GGTCCG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> *Pseudomonas aeruginosa*
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGGCGCGCCGCGCATGCCGGCGCCGC
 TGCTGGAACGCTCGACCCAGGGCTACACTA**GTTTAACCGGAACGCC**GCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	.20	0.0	0.0	.40	.20	.60	.20	.20	.60	0.0	.40	0.0	.20	0.0	.60	0.0
C	.20	0.0	0.0	0.0	.20	.20	.60	.40	0.0	0.0	0.0	.40	.20	.40	.20	.80
G	.60	0.0	0.0	0.0	.20	0.0	0.0	0.0	.20	.60	.20	.20	.40	.20	.20	.20
T	0.0	1.0	1.0	.60	.40	.20	.20	.40	.20	.40	.40	.40	.20	.40	0.0	0.0

> *Escherichia coli*
 TTGATTCCCTGAATGCCCGCTTAGTGTAACACTACTGTAACCGGCATTTTCTGCTTTTCC
 TGCCGAT**ATTTTTCTTATCTAC**CTCACAAAGGTTAGCAATAACTGCTGGGAAAATTCGG
 AGTTAGTCGTTATATTCTAT

> *Haemophilus influenzae*
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGAATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTCCGTGCCCTCCGCCATT
AGTTACACTAGTGGGACACC

> *Vibrio cholerae*
 ATTTGTGGCG**TTTTCAAATGCTTGG**AGAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGTTCTGGGGCTGTTTGAACGTTACACGAGTGAACCGCGGAA
 CCATGTTGACACGAATTCTG

> *Salmonella typhi*
 GGTCCG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> *Pseudomonas aeruginosa*
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGGCGCGCCGCGCATGCCGGCGCCGC
 TGCTGGAACGCTCGACCCAGGGCTACACTA**GTTTAACCGGAACGCC**GCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	.20	0.0	0.0	.40	.20	.60	.20	.20	.60	0.0	.40	0.0	.20	0.0	.60	0.0
C	.20	0.0	0.0	0.0	.20	.20	.60	.40	0.0	0.0	0.0	.40	.20	.40	.20	.80
G	.60	0.0	0.0	0.0	.20	0.0	0.0	0.0	.20	.60	.20	.20	.40	.20	.20	.20
T	0.0	1.0	1.0	.60	.40	.20	.20	.40	.20	.40	.40	.40	.20	.40	0.0	0.0

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGT**GTAACACTACTGTAAC**CGGCATTTTCTGCTTTTCC
 TGCCGATATTTTTTCTTATCTACCTCACAAAGGTTAGCAATAACTGCTGGGAAAATTCGG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCG**GTTTTCAAATGCTTGG**AGAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGCTTGGGGCTGTTGAACGTTACACGAGTGAACCGCGGAA
 CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCGG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
 TGCTGGAACGCTCGACCCAGG**GCTACACTAGTTTAAC**CGGAACGCCGCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	0.0	0.0	.20	.80	0.0	.80	.20	.20	1.0	0.0	0.0	0.0	0.0	.40	.80	0.0
C	.20	0.0	0.0	0.0	.60	.20	.80	0.0	0.0	.20	0.0	.20	0.0	0.0	0.0	.80
G	.80	0.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	.60	.20	.60	.20	.40	.20	.20
T	0.0	1.0	.80	.20	.20	0.0	0.0	.80	0.0	.20	.80	.20	.80	.20	0.0	0.0

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGT**GTAACACTACTGTAAC**CGGCATTTTCTGCTTTTCC
 TGCCGATATTTTTTCTTATCTACCTCACAAAGGTTAGCAATAACTGCTGGGAAAATTCGG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCG**GTTTTCAAATGCTTGG**AGAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGCTTGGGGCTGTTGAACGTTACACGAGTGAACCGCGGAA
 CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCGG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCGTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
 TGCTGGAACGCTCGACCCAGG**GCTACACTAGTTTAAC**CGGAACGCCGCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	0.0	0.0	.20	.80	0.0	.80	.20	.20	1.0	0.0	0.0	0.0	0.0	.40	.80	0.0
C	.20	0.0	0.0	0.0	.60	.20	.80	0.0	0.0	.20	0.0	.20	0.0	0.0	0.0	.80
G	.80	0.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	.60	.20	.60	.20	.40	.20	.20
T	0.0	1.0	.80	.20	.20	0.0	0.0	.80	0.0	.20	.80	.20	.80	.20	0.0	0.0

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGT**GTAACACTACTGTAAC**CGGCATTTTCTGCTTTTCC
 TGCCGATATTTTTTCTTATCTACCTCACAAAGTTAGCAATAACTGCTGGGAAAATTCG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCGTTTTCAAATGCTTGAGAAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGCTGGGGCTGTTGAAC**GTTACACGAGTGTAAAC**CGCCGAA
 CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCCG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
 TGCTGGAACGCCTCGACCCAG**GCTACACTAGTTTAAAC**CGGAACGCCGCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

A	0.0	0.0	.20	1.0	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	.60	1.0	0.0
C	.20	.20	0.0	0.0	.80	0.0	1.0	0.0	0.0	.20	0.0	0.0	0.0	0.0	0.0	1.0
G	.80	0.0	0.0	0.0	.20	0.0	0.0	.20	0.0	.80	0.0	.80	.20	.40	0.0	0.0
T	0.0	.80	.80	0.0	0.0	0.0	0.0	.80	0.0	0.0	1.0	.20	.80	0.0	0.0	0.0

> Escherichia coli
 TTGATTCCCTGAATGCCCGCTTAGT**GTAACACTACTGTAAC**CGGCATTTTCTGCTTTTCC
 TGCCGATATTTTTTCTTATCTACCTCACAAAGTTAGCAATAACTGCTGGGAAAATTCG
 AGTTAGTCGTTATATTCTAT

> Haemophilus influenzae
 ATCTAACGGTACGGATTCTCCAAAGGCCTATGGAATCTTGTAAGATATGAAACGTTCTAA
 TAAATCATAAAGTTGGAGCAAACGCTCGGCATAAGTAGTAAGTGCCGTGCCTCCGCCATT
AGTTACACTAGTGGGACACC

> Vibrio cholerae
 ATTTGTGGCGTTTTCAAATGCTTGAGAAATGGGTACATGATCCGCTTGGCATTGAAGGT
 GAGGCTGGCAGCAGCGAAGGCTGGGGCTGTTGAAC**GTTACACGAGTGTAAAC**CGCCGAA
 CCATGTTGACACGAATTCTG

> Salmonella typhi
 GGTCCG**CTTAGACTAGTGTGAC**CAAAAAGCTTTTGCTGAAGTTTCAGGTAAGAAGAACC
 AGCTCCTAGTAAAAGACTATTGTGACTGAAAAGCGCTCAGCGCAAAGCCGACCCGACA
 AAACGCACAAGGAGTTACAG

> Pseudomonas aeruginosa
 ACGCGCCAGGGTCTTCTCCTGCGAGATCATGCGCGCGCGCCGCGCATGCCGGCGCCGC
 TGCTGGAACGCCTCGACCCAG**GCTACACTAGTTTAAAC**CGGAACGCCGCCAGTGGATCG
 GCCTGCCCCAGCTATTGCTC

Expectation-Maximization (EM)

- Randomly guess the locations of each motif instance
- Repeat until convergence
 - Calculate a new motif model from the motif instances
 - Calculate new locations for the motif instances from the motif model

Each motif instance is *best scoring* match to motif model

Gibbs Sampling

- Randomly guess the locations of each motif instance
- Repeat until convergence
 - Calculate a new motif model from the motif instances
 - Calculate new locations for the motif instances from the motif model

Each motif instance is *sampled* from scores of matches to motif model