



# Information Theory: One-Minute Lesson

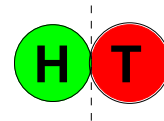
number of symbols	number of bits	example
-------------------	----------------	---------

**M**

**B**

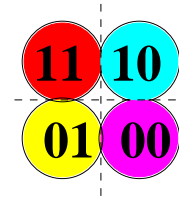
**2**

**1**



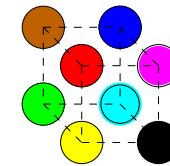
**4**

**2**



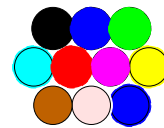
**8**

**3**



$$M=2^B$$

$$B=\log_2 M$$



# Information Theory: One-Minute Lesson

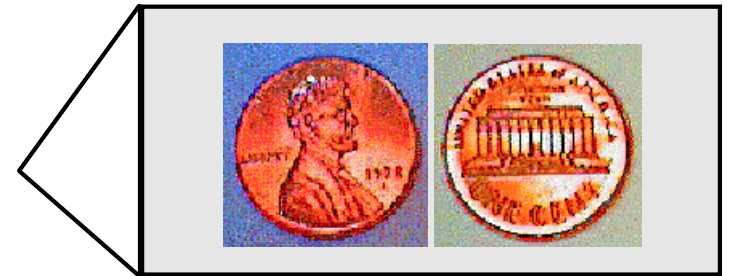
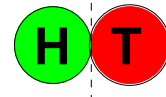
number of symbols	number of bits	example
-------------------	----------------	---------

**M**

**B**

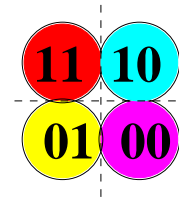
**2**

**1**



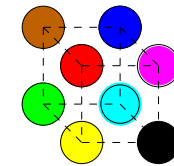
**4**

**2**



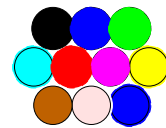
**8**

**3**



$$M=2^B$$

$$B=\log_2 M$$



# Information Theory: One-Minute Lesson

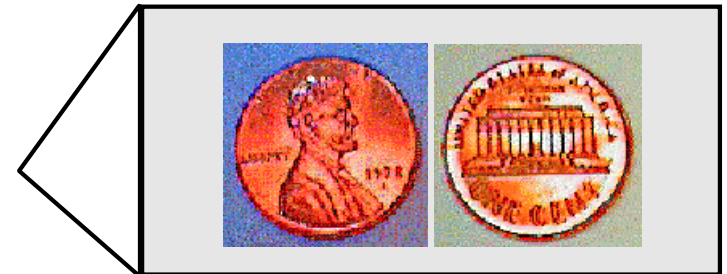
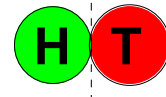
number of symbols	number of bits	example
-------------------	----------------	---------

**M**

**B**

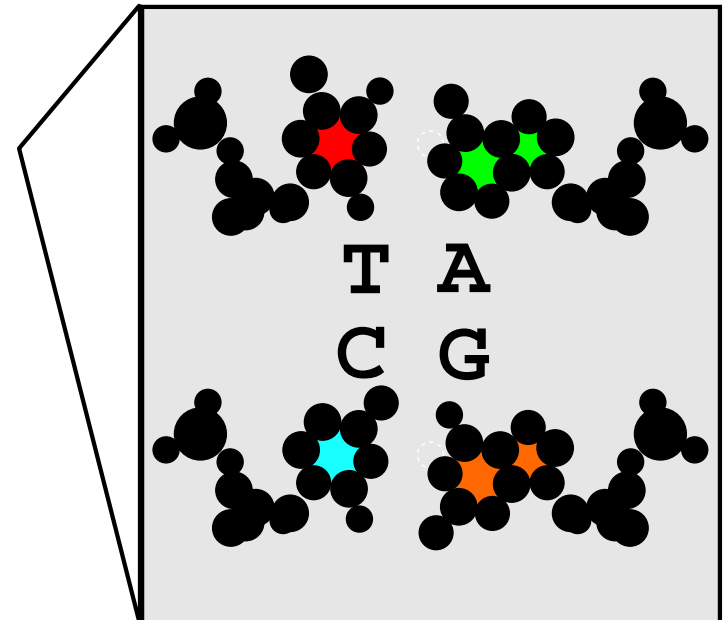
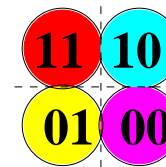
**2**

**1**



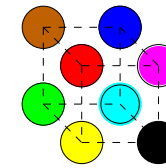
**4**

**2**



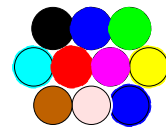
**8**

**3**



$$M=2^B$$

$$B=\log_2 M$$



# Information Theory: One-Minute Lesson

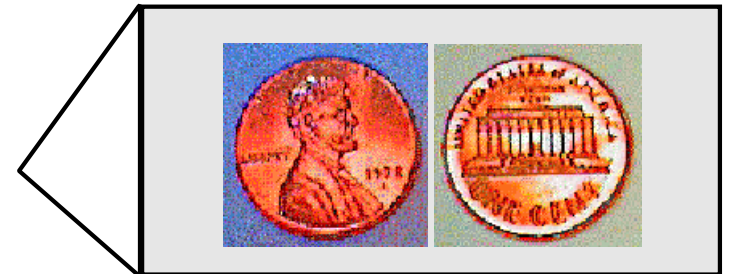
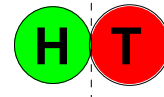
number of symbols	number of bits	example
-------------------	----------------	---------

**M**

**B**

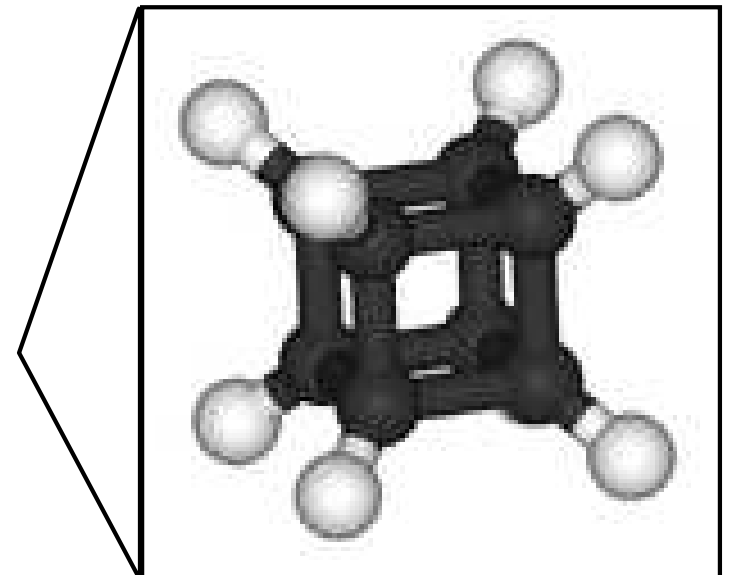
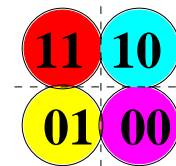
**2**

**1**



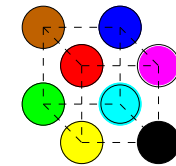
**4**

**2**



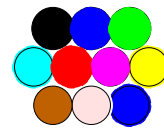
**8**

**3**



$$M=2^B$$

$$B=\log_2 M$$



# Information Theory: One-Minute Lesson

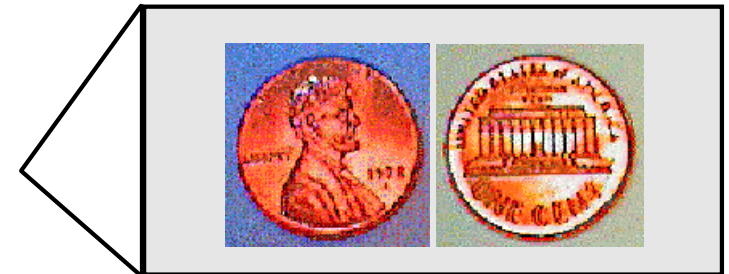
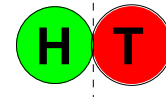
number of symbols	number of bits	example
-------------------	----------------	---------

**M**

**B**

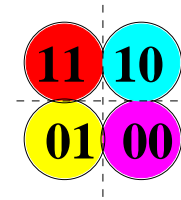
**2**

**1**



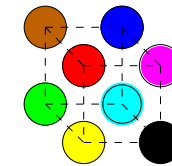
**4**

**2**



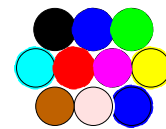
**8**

**3**



$$M=2^B$$

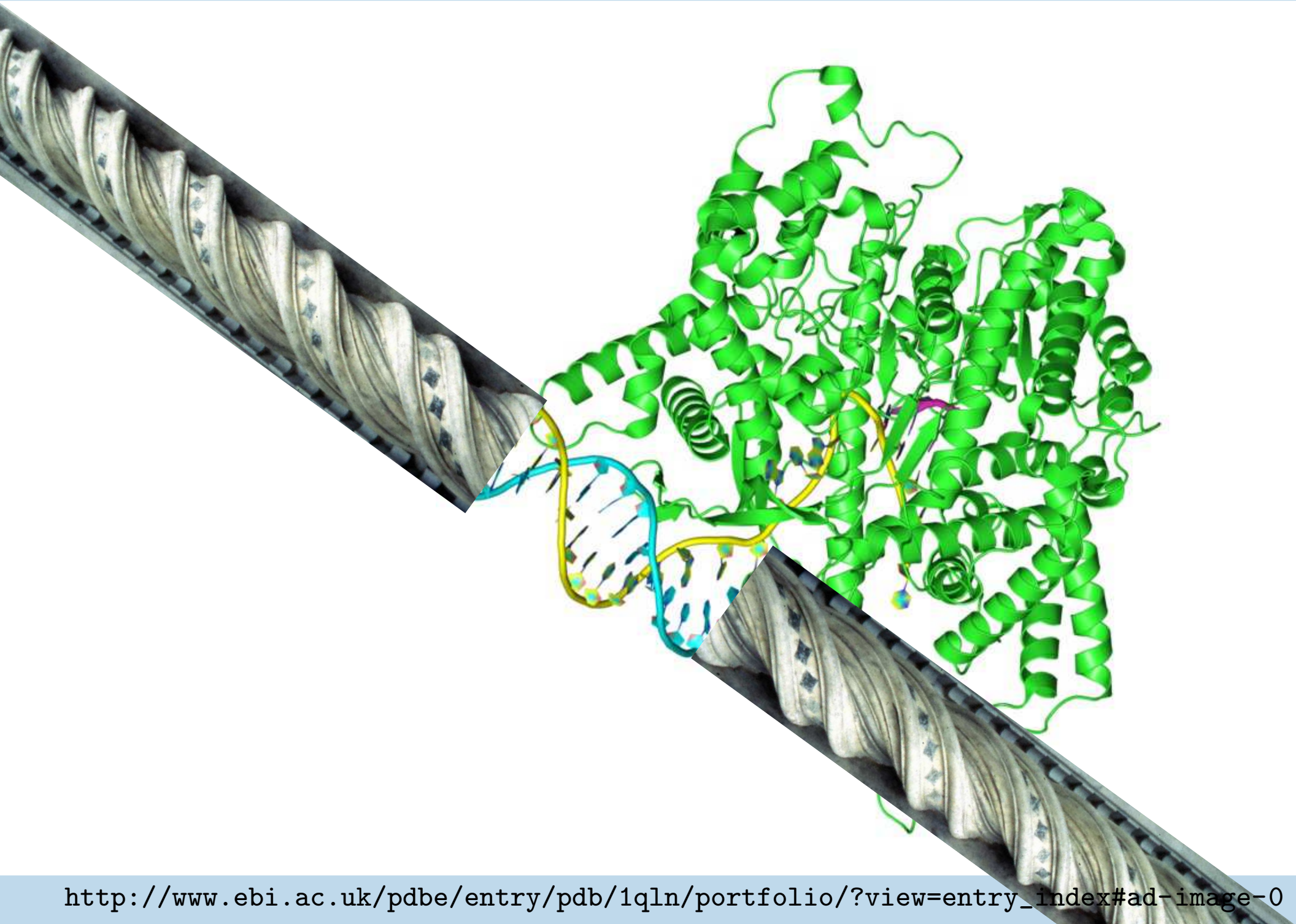
$$B=\log_2 M$$



# El Duomo, Florence, Italy

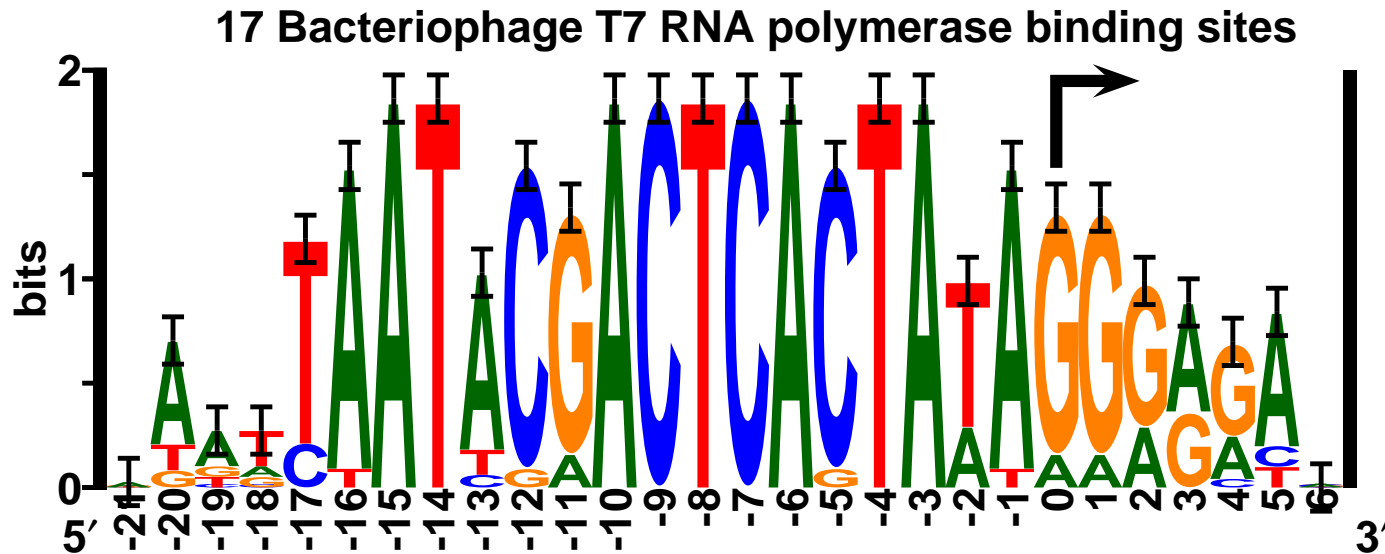


# T7 RNA polymerase + DNA





# Sequence Logo

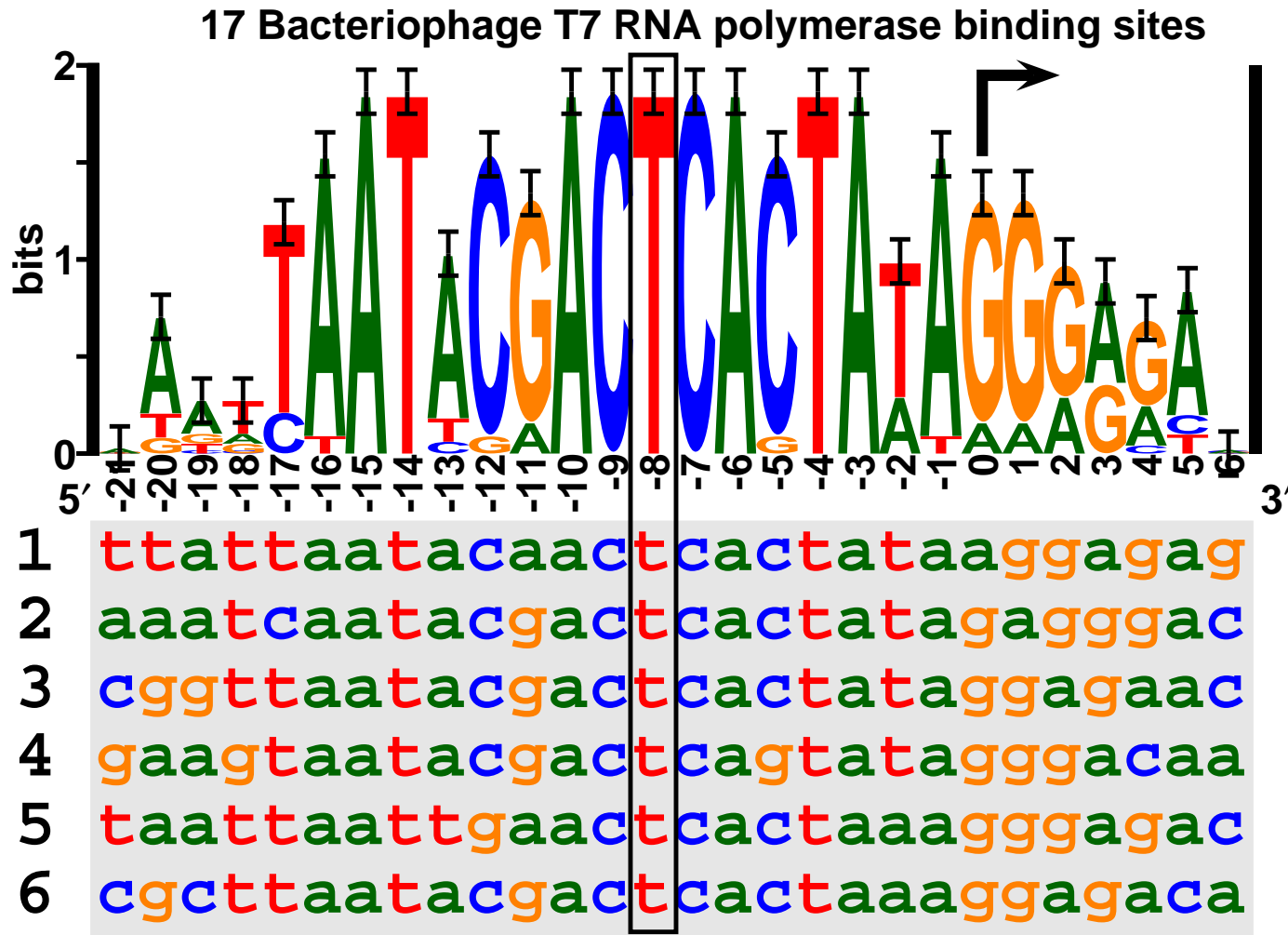


Schneider &  
Stephens  
*Nucl. Acids Res.*  
18: 6097-6100  
1990

```
1 ttattaatacaactcactataaggagag
2 aaatcaatacgaactcactatagaggac
3 cggttaatacgaactcactataggagaac
4 gaagtaatacgaactcagtatagggacaa
5 taattaattgaactcactaaaggggagac
6 cgcttaatacgaactcactaaaggagaca
```

6 of 17 sites

# Sequence Logo

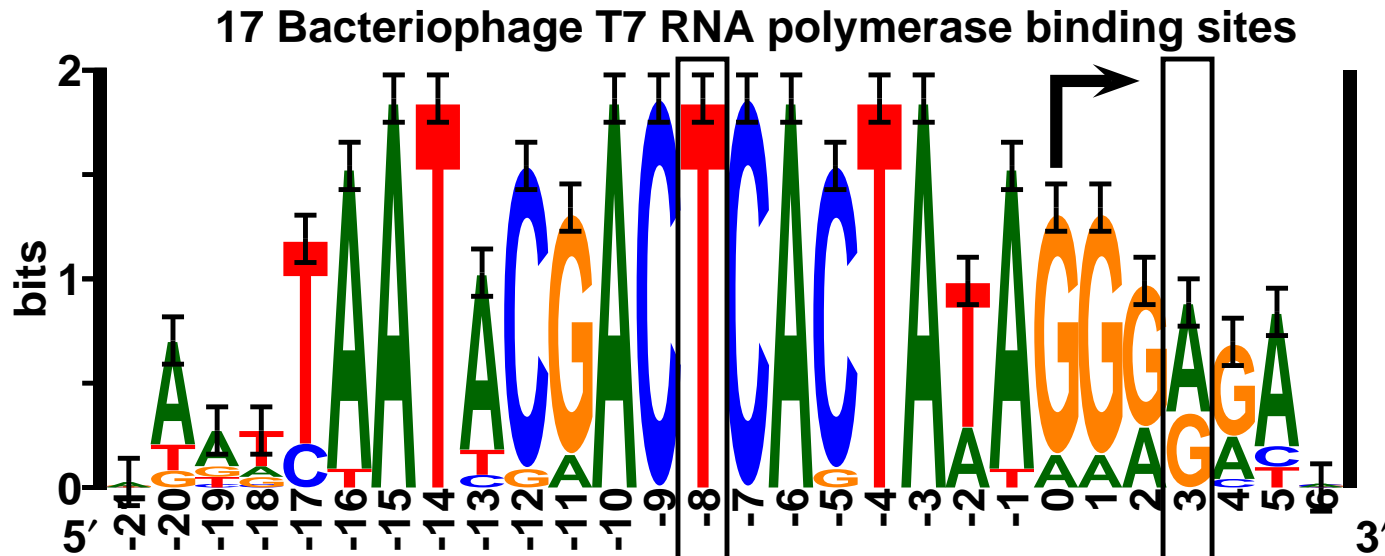


Schneider &  
Stephens  
*Nucl. Acids Res.*  
18: 6097-6100  
1990

6 of 17 sites

2 bits/base

# Sequence Logo



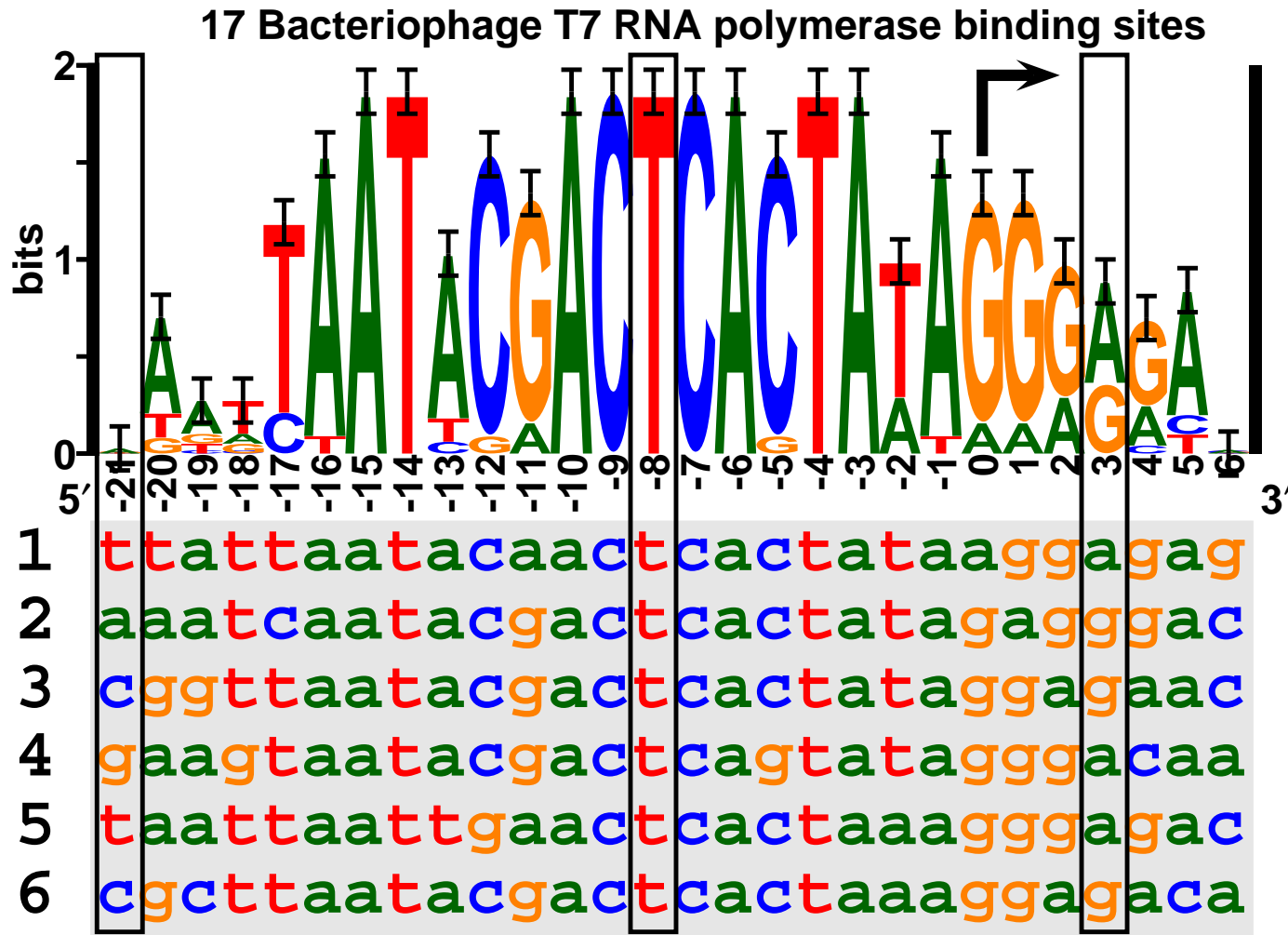
Schneider &  
Stephens  
*Nucl. Acids Res.*  
18: 6097-6100  
1990

1 ttattaatacaactcactataaggagag  
2 aaatcaatacgaactcactatagaggac  
3 cggttaatacgaactcactataggagaac  
4 gaagtaatacgaactcagtatagggacaa  
5 taattaattgaactcactaaaggaggac  
6 cgcttaatacgaactcactaaaggagaca

6 of 17 sites

1 bit/base

# Sequence Logo

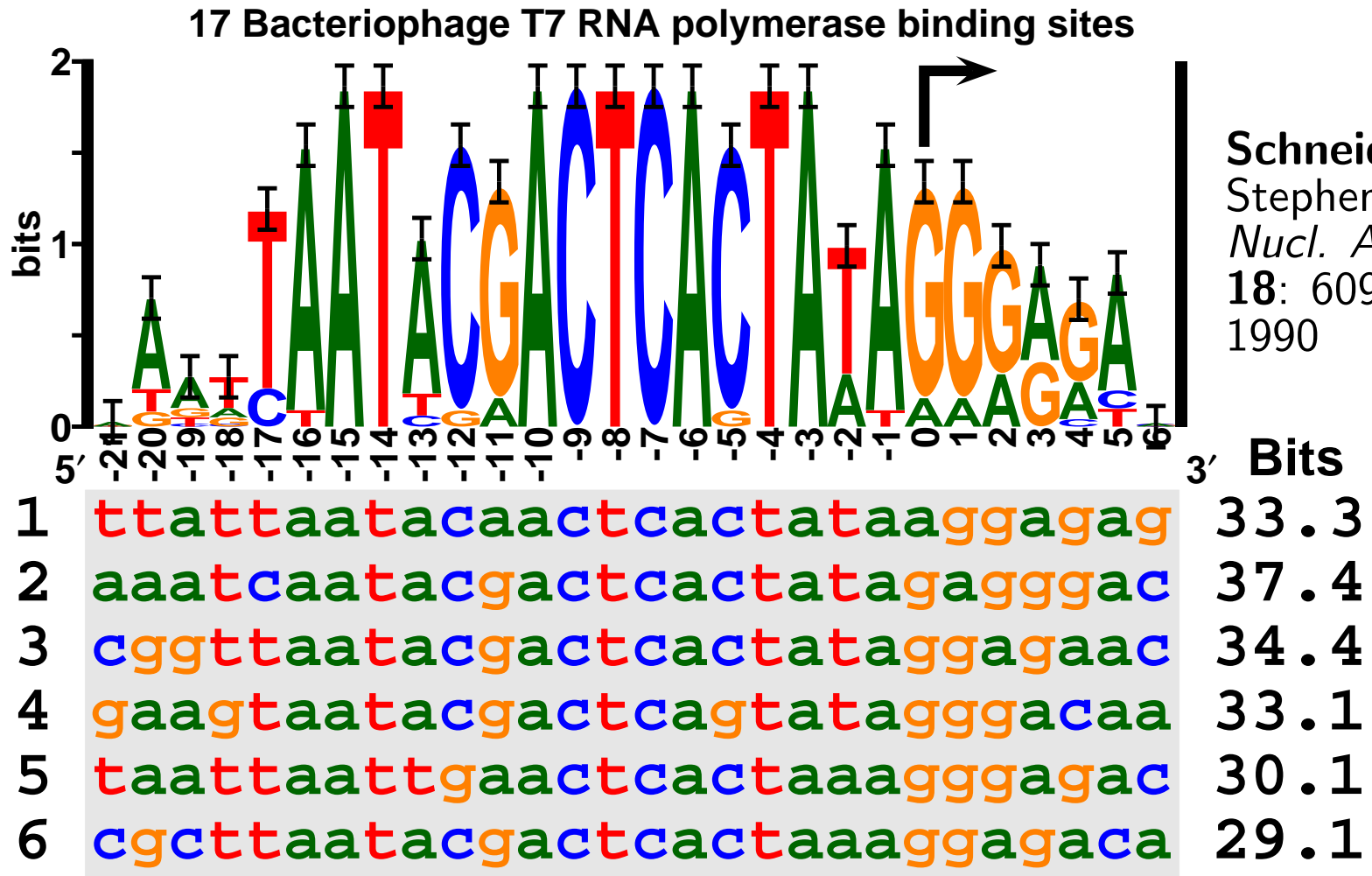


Schneider &  
Stephens  
*Nucl. Acids Res.*  
18: 6097-6100  
1990

6 of 17 sites

0 bits/base

# Sequence Logo

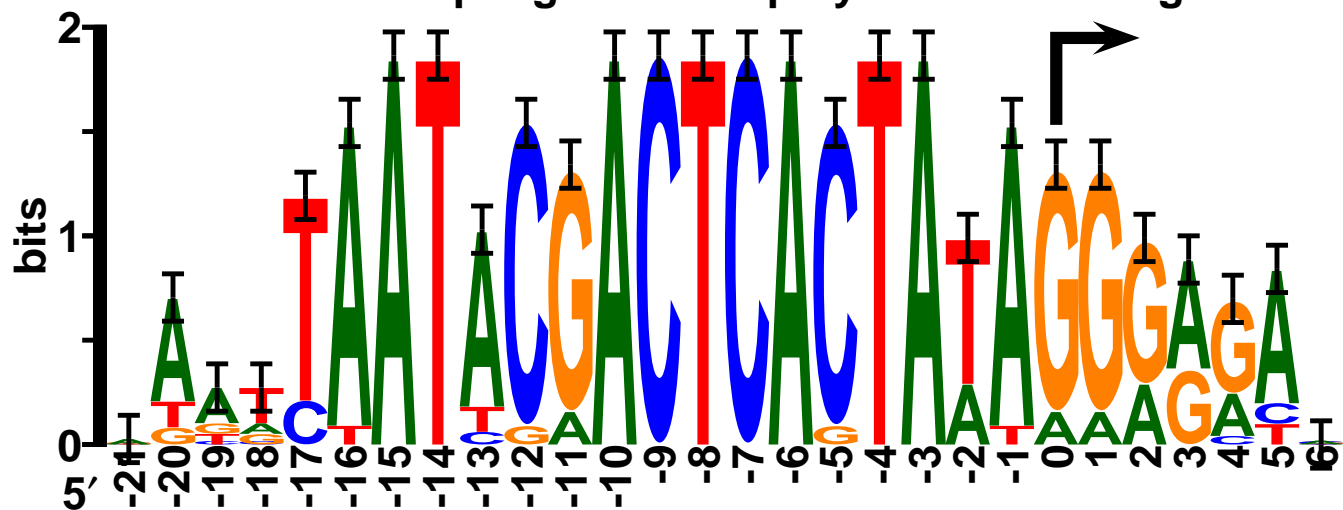


Schneider &  
Stephens  
*Nucl. Acids Res.*  
18: 6097-6100  
1990

Individual Information

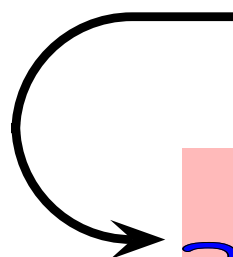
# Sequence Logo and Sequence Walker

17 Bacteriophage T7 RNA polymerase binding sites



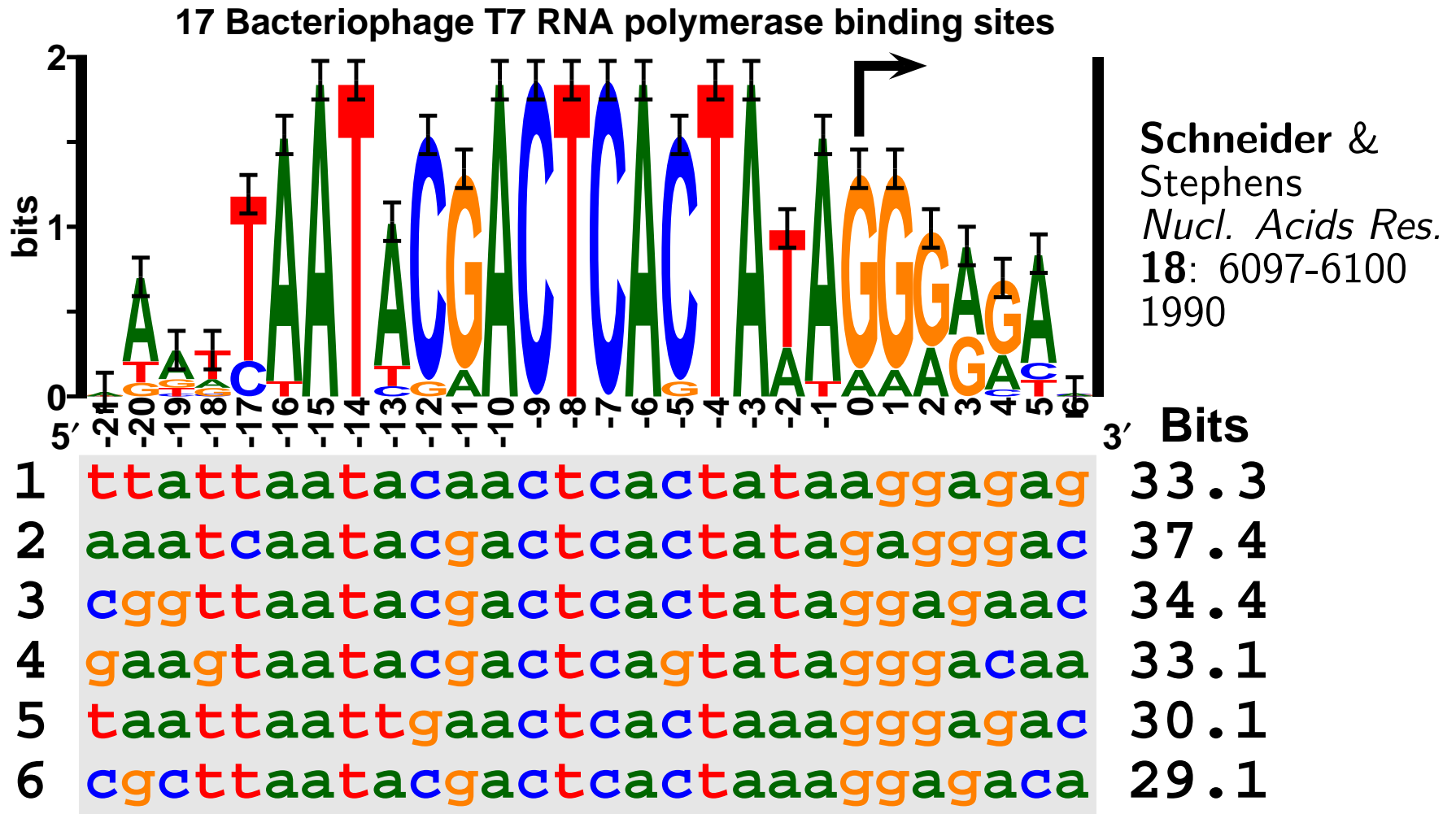
Schneider &  
Stephens  
*Nucl. Acids Res.*  
18: 6097-6100  
1990

	Sequence	Bits
1	ttattaatacaactcactataaggagag	33.3
2	aatcaatacgactcactatagagggac	37.4
3	ggttaatacgactcactataggagaac	34.4
4	gaagtaatacgactcagtatagggacaa	33.1
5	taattaattgaactcactaaaggggagac	30.1
6	cgcttaatacgactcactaaaggagaca	29.1



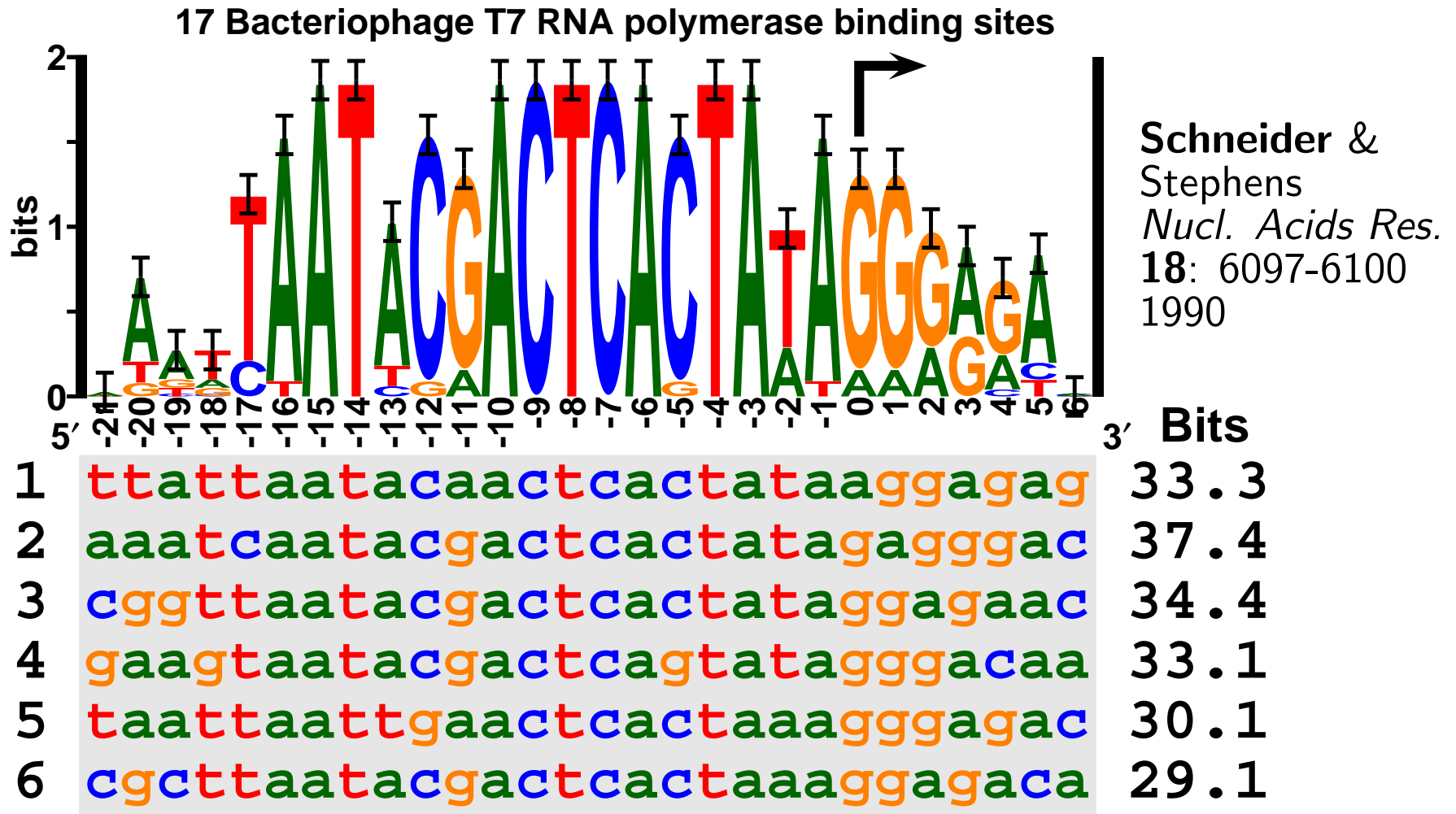
Sequence  
Walker  
Patent  
5,867,402

# Sequence Logo and Sequence Walker and Rsequence



Rsequence is the average:  $35.0 \pm 0.6$  bits

# Sequence Logo and Sequence Walker and Rsequence



Rsequence is the average:  $35.0 \pm 0.6$  bits  
= "area under the logo"



## Information required to find a set of binding sites

$G = \#$  of potential binding sites

## Information required to find a set of binding sites

$G$  = # of potential binding sites  
= genome size in some cases

## Information required to find a set of binding sites

$G$  = # of potential binding sites  
= genome size in some cases

$\gamma$  = number of binding sites on genome

## Information required to find a set of binding sites

$G$  = # of potential binding sites  
= genome size in some cases

$\gamma$  = number of binding sites on genome

$$R_{frequency} = H_{before} - H_{after}$$

## Information required to find a set of binding sites

$$\begin{aligned} G &= \# \text{ of potential binding sites} \\ &= \text{genome size in some cases} \end{aligned}$$

$\gamma$  = number of binding sites on genome

$$\begin{aligned} R_{\text{frequency}} &= H_{\text{before}} - H_{\text{after}} \\ &= \log_2 G - \log_2 \gamma \end{aligned}$$

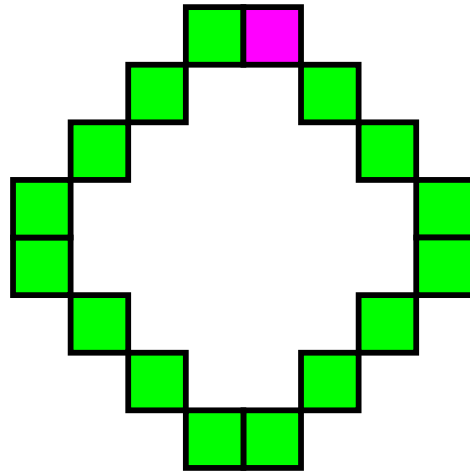
## Information required to find a set of binding sites

$$\begin{aligned} G &= \# \text{ of potential binding sites} \\ &= \text{genome size in some cases} \end{aligned}$$

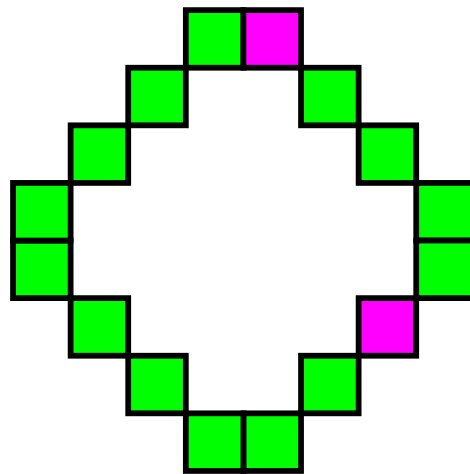
$\gamma$  = number of binding sites on genome

$$\begin{aligned} R_{\text{frequency}} &= H_{\text{before}} - H_{\text{after}} \\ &= \log_2 G - \log_2 \gamma \\ &= -\log_2 \gamma/G \end{aligned}$$

# Information required to find a set of binding sites in a genome



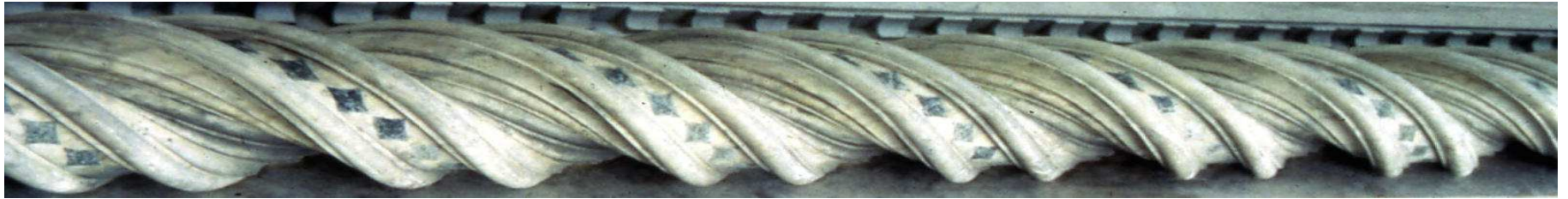
16 positions  
1 site  
 $\log_2 16/1 = 4$  bits



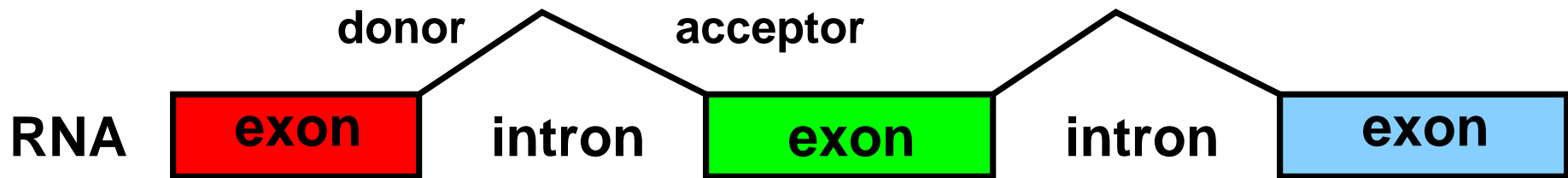
16 positions  
2 sites  
 $\log_2 16/2 = 3$  bits

# RNA Splicing

DNA



↓ Copy DNA (transcription)



↓ RNA Splicing



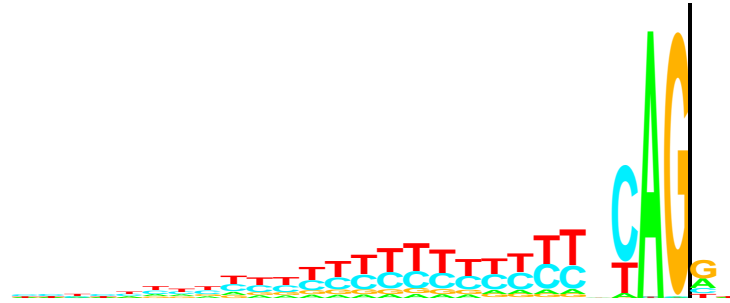






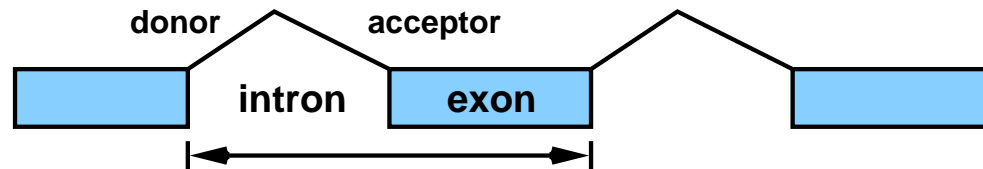
# Rsequence and Rfrequency for Splice Acceptors

$R_{sequence}$



- Information at binding site sequences (area under sequence logo)
- from: binding site sequences
- 9.4 bits per site

$R_{frequency}$



- Information needed to locate the sites
- from: size of genome and number of sites (length of intron+exon)
- 9.7 bits per site

$$R_{frequency} / R_{sequence} = 0.97$$

# Hypothesis:

The information in  
binding site patterns  
is just sufficient  
for the sites to be found  
in the genome

# Rsequence versus Rfrequency

Binding Site Recognizer <sup>1</sup>	Total Pattern Information = $R_{\text{sequence}}$ (bits)	Information needed to Locate Site in Genome = $R_{\text{frequency}}$ (bits)	$\frac{\text{Pattern Info}}{\text{Location Info}}$ = $\frac{R_{\text{sequence}}}{R_{\text{frequency}}}$
Spliceosome acceptor <sup>2</sup>	<b>9.35 ± 0.12</b>	<b>9.66</b>	<b>0.97 ± 0.01</b>
Spliceosome donor	<b>7.92 ± 0.09</b>	<b>9.66</b>	<b>0.82 ± 0.01</b>
Ribosome	<b>11.0</b>	<b>10.6</b>	<b>1.0</b>
$\lambda$ cl/cro	<b>17.7 ± 1.6</b>	<b>19.3</b>	<b>0.9 ± 0.1</b>
LexA	<b>21.5 ± 1.7</b>	<b>18.4</b>	<b>1.2 ± 0.1</b>
TrpR	<b>23.4 ± 1.9</b>	<b>20.3</b>	<b>1.2 ± 0.1</b>
LacI	<b>19.2 ± 2.8</b>	<b>21.9</b>	<b>0.9 ± 0.1</b>
ArgR	<b>16.4</b>	<b>18.4</b>	<b>0.9</b>
O ( $\lambda$ Origin)	<b>20.9</b>	<b>19.9</b>	<b>1.0</b>
Ara C	<b>19.3</b>	<b>19.3</b>	<b>1.0</b>
Transcription at TATA <sup>3</sup>	<b>3.3</b>	<b>~ 3</b>	<b>~ 1</b>
T7 Promoter	<b>35.4</b>	<b>16.5</b>	<b>2.1</b>

<sup>1</sup>T. D. Schneider, G. D. Stormo, L. Gold, and A. Ehrenfeucht. J. Mol. Biol., 188:415-431, 1986.

<sup>2</sup>R. M. Stephens and T. D. Schneider. J. Mol. Biol., 228:1124-1136, 1992.

<sup>3</sup>F. E. Penotti. J Mol Biol, 213:37-52, 1990.

## $R_{sequence}$ versus $R_{frequency}$ - meaning

The information in the binding site pattern ( $R_{sequence}$ )  
is close to  
The information needed to find the binding sites ( $R_{frequency}$ )

# $R_{sequence}$ versus $R_{frequency}$ - meaning

The information in the binding site pattern ( $R_{sequence}$ )  
is close to

The information needed to find the binding sites ( $R_{frequency}$ )

But for a species in a stable environment:

- size of genome ( $G$ ) is fixed (e. g. *E. coli* has  $4.7 \times 10^6$  bp)
- number of binding sites ( $\gamma$ ) is fixed (e. g. there are  $\sim 50$  *E. coli* LexA sites)

so  $R_{frequency} = \log_2 G/\gamma$  is fixed

# $R_{sequence}$ versus $R_{frequency}$ - meaning

The information in the binding site pattern ( $R_{sequence}$ )  
is close to

The information needed to find the binding sites ( $R_{frequency}$ )

But for a species in a stable environment:

- size of genome ( $G$ ) is fixed (e. g. *E. coli* has  $4.7 \times 10^6$  bp)
- number of binding sites ( $\gamma$ ) is fixed (e. g. there are  $\sim 50$  *E. coli* LexA sites)

so  $R_{frequency} = \log_2 G/\gamma$  is fixed

$R_{sequence}$  must evolve towards  $R_{frequency}$ !



## Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$

## Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

## Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

Setup a Computer Model, 'Ev':

A population of "creatures" with

# Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

Setup a Computer Model, 'Ev':

A population of "creatures" with

- genomes containing 4 bases (A, C, G, T)

# Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

Setup a Computer Model, 'Ev':

A population of "creatures" with

- genomes containing 4 bases (A, C, G, T)
- a defined genome size ( $G$ )

# Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

Setup a Computer Model, 'Ev':

A population of "creatures" with

- genomes containing 4 bases (A, C, G, T)
- a defined genome size ( $G$ )
- predetermined binding site locations ( $\gamma$ )  
(to fix the frequency of sites)

# Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

Setup a Computer Model, 'Ev':

A population of "creatures" with

- genomes containing 4 bases (A, C, G, T)
  - a defined genome size ( $G$ )
  - predetermined binding site locations ( $\gamma$ )  
(to fix the frequency of sites)
- }  $R_{frequency}$   
is fixed

# Evolution of Binding Sites

- $R_{frequency}$  is fixed relative to  $R_{sequence}$
- Does  $R_{sequence}$  evolve toward  $R_{frequency}$ ?

Setup a Computer Model, 'Ev':

A population of "creatures" with

- genomes containing 4 bases (A, C, G, T)
  - a defined genome size ( $G$ )
  - predetermined binding site locations ( $\gamma$ )  
(to fix the frequency of sites)
  - a recognizer gene encoded in the sequence:  
use a weight matrix
- }  $R_{frequency}$   
is fixed



# How A Weight Matrix Works

Sequence matrix,  $s(b, l, j)$  for sequence  $j$

base b	position l									
	C	A	G	G	T	C	T	G	C	A
	-3	-2	-1	0	1	2	3	4	5	6
A	0	1	0	0	0	0	0	0	0	1
C	1	0	0	0	0	1	0	0	1	0
G	0	0	1	1	0	0	0	1	0	0
T	0	0	0	0	1	0	1	0	0	0

Individual information weight matrix,  $R_{iw}(b, l)$

base b	position l									
	-3	-2	-1	0	1	2	3	4	5	6
A	+0.4	+1.3	-1.4	-8.8	-5.8	+1.1	+1.5	-1.8	-0.7	+0.0
C	+0.6	-0.8	-2.4	-7.8	-5.5	-3.7	-1.6	-2.2	-0.5	-0.2
G	-0.6	-1.0	+1.6	+2.0	-6.2	+0.7	-1.1	+1.7	-0.3	+0.4
T	-1.0	-0.9	-1.7	-5.8	+2.0	-3.4	-1.6	-2.2	+0.9	-0.5

# How A Weight Matrix Works

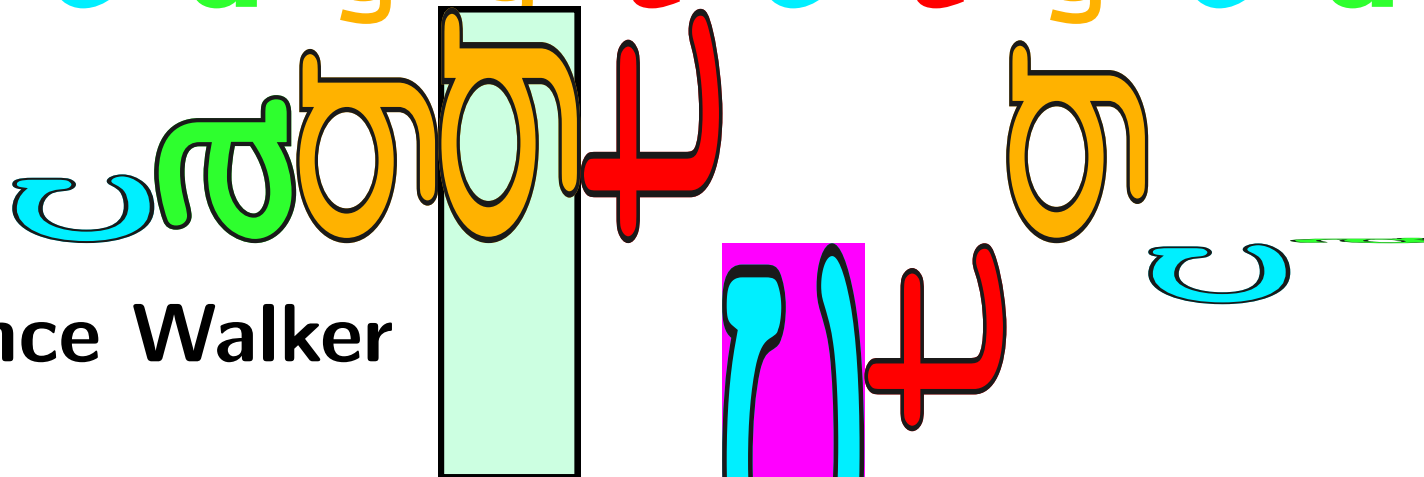
Sequence matrix,  $s(b, l, j)$  for sequence  $j$

base b	position l									
	C	A	G	G	T	C	T	G	C	A
	-3	-2	-1	0	1	2	3	4	5	6
A	0	1	0	0	0	0	0	0	0	1
C	1	0	0	0	0	1	0	0	1	0
G	0	0	1	1	0	0	0	1	0	0
T	0	0	0	0	1	0	1	0	0	0

Individual information weight matrix,  $R_{iw}(b, l)$

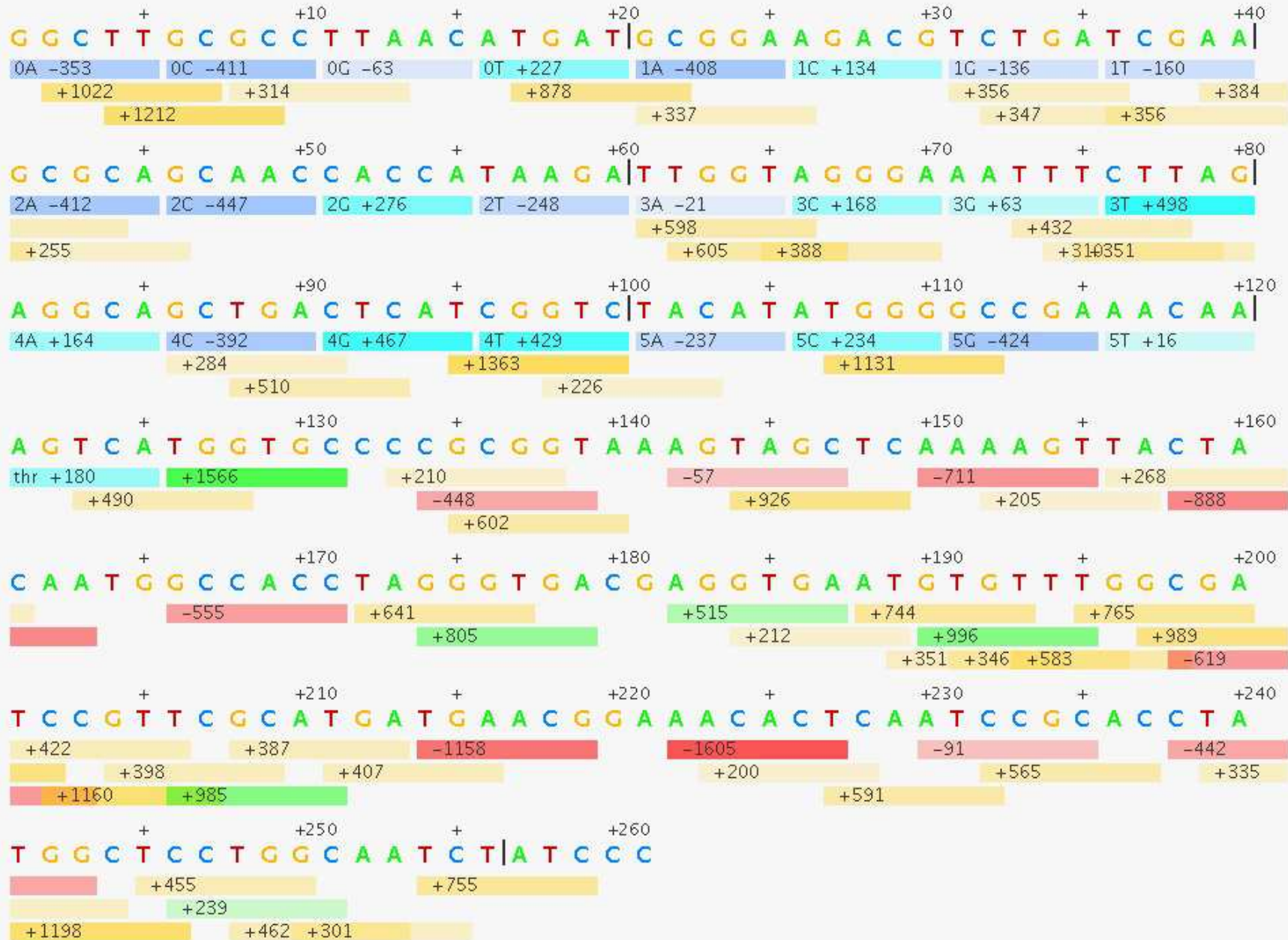
base b	position l									
	-3	-2	-1	0	1	2	3	4	5	6
A	+0.4	+1.3	-1.4	-8.8	-5.8	+1.1	+1.5	-1.8	-0.7	+0.0
C	+0.6	-0.8	-2.4	-7.8	-5.5	-3.7	-1.6	-2.2	-0.5	-0.2
G	-0.6	-1.0	+1.6	+2.0	-6.2	+0.7	-1.1	+1.7	-0.3	+0.4
T	-1.0	-0.9	-1.7	-5.8	+2.0	-3.4	-1.6	-2.2	+0.9	-0.5

5' **c** **a** **g** **g** **t** **c** **t** **g** **c** **a** 3'

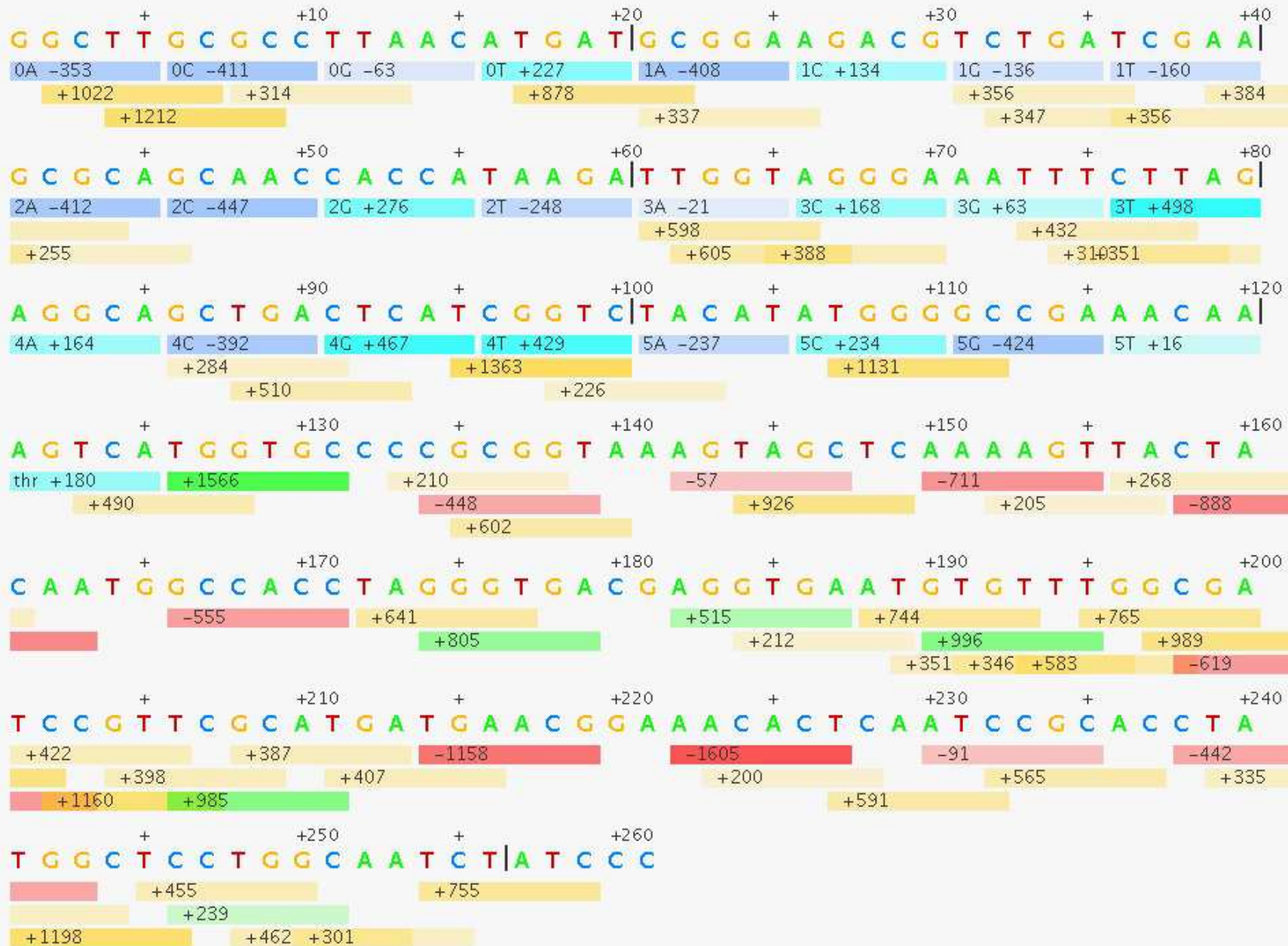


Sequence Walker

# Unevolved Ev Creature

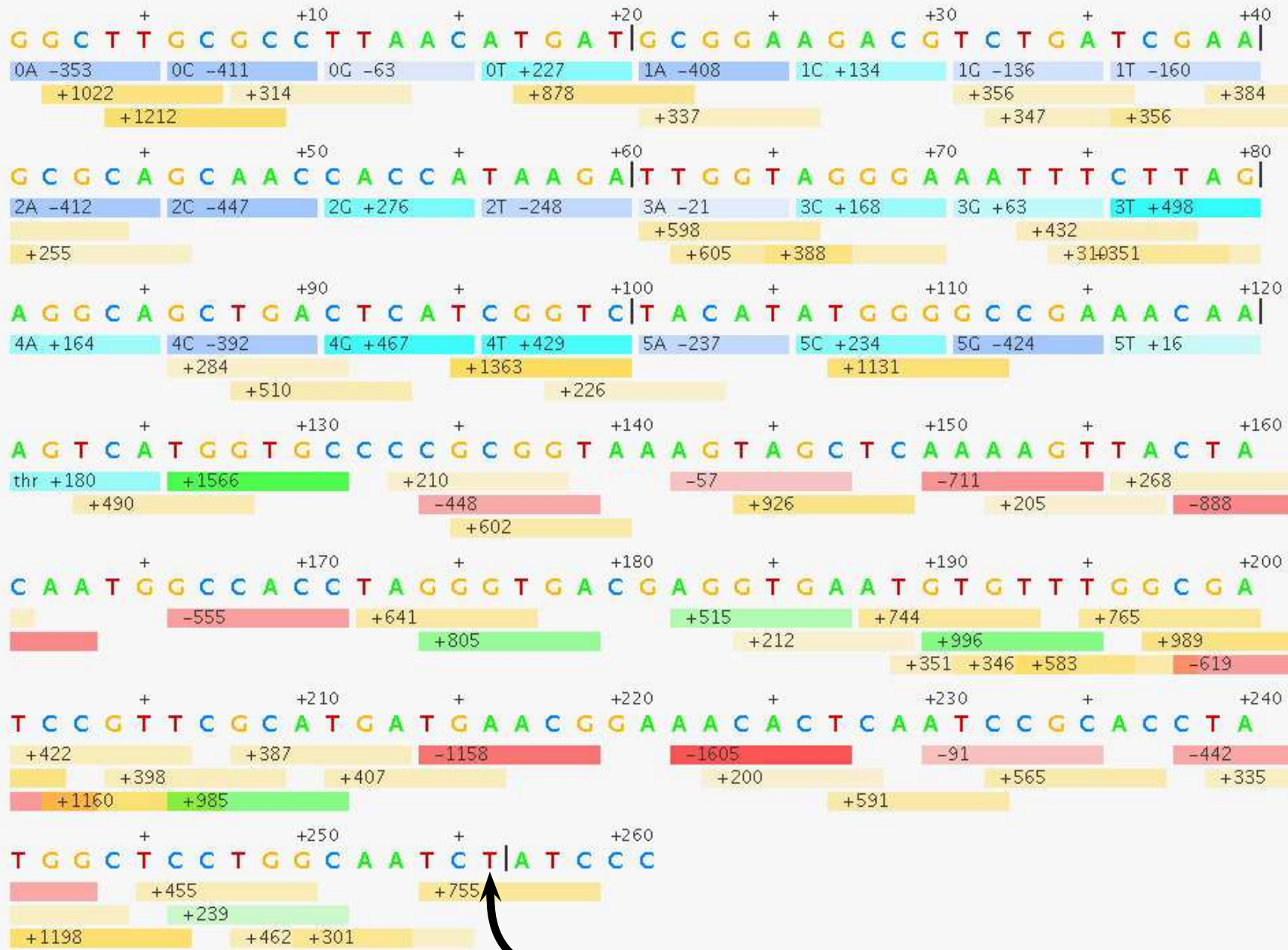


# Unevolved Ev Creature



“blue”  
gene  
weight  
matrix:  
6 bp  
wide

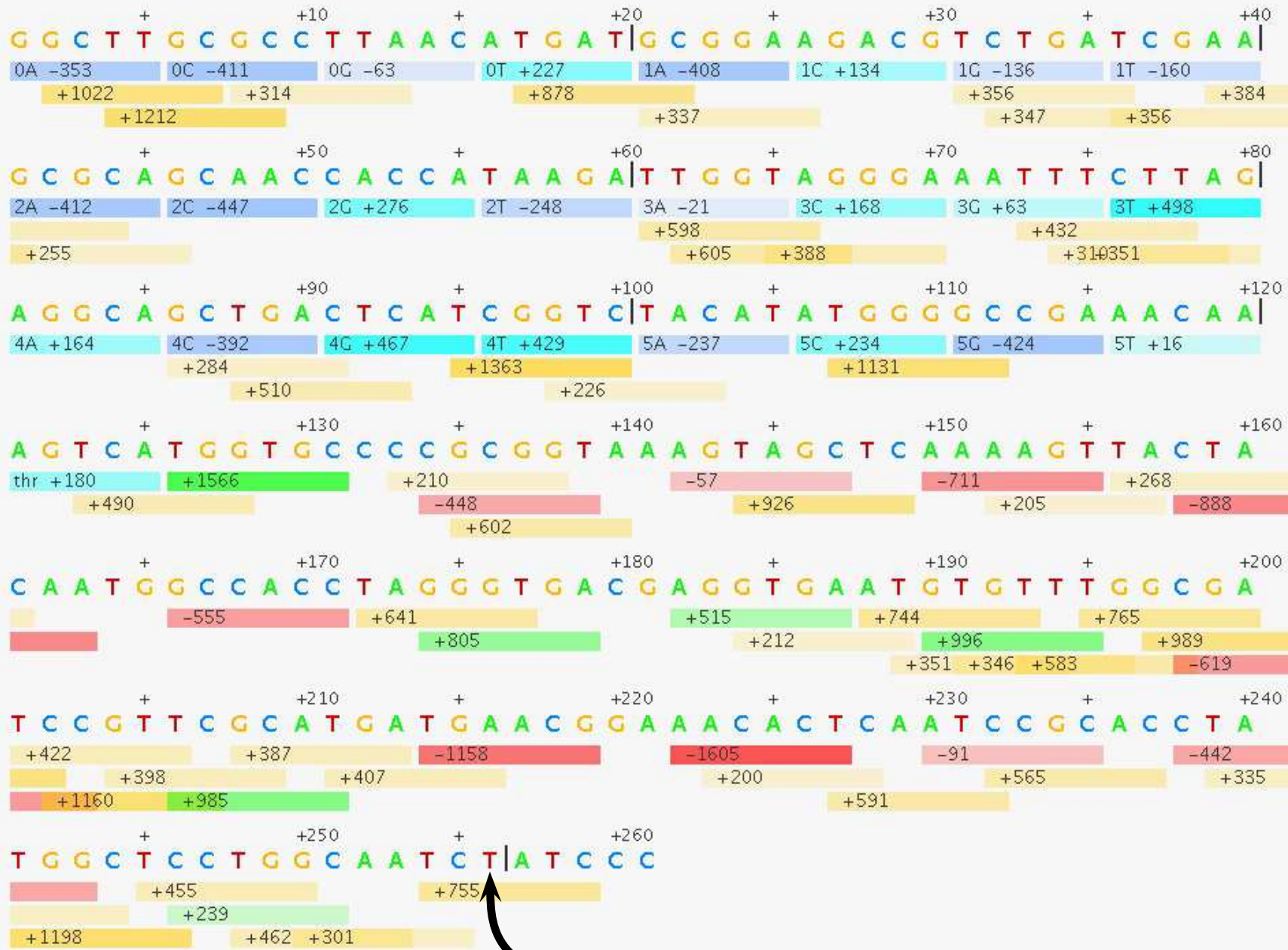
# Unevolved Ev Creature



“blue”  
gene  
weight  
matrix:  
6 bp  
wide

Genome positions available  $G = 256$  bases

# Unevolved Ev Creature



“blue”  
gene  
weight  
matrix:  
6 bp  
wide

$\gamma = 16$   
binding  
sites

Genome positions available  $G = 256$  bases  
 $R_{frequency} = \log_2 256/16 = 4$  bits

# Unevolved Ev Creature



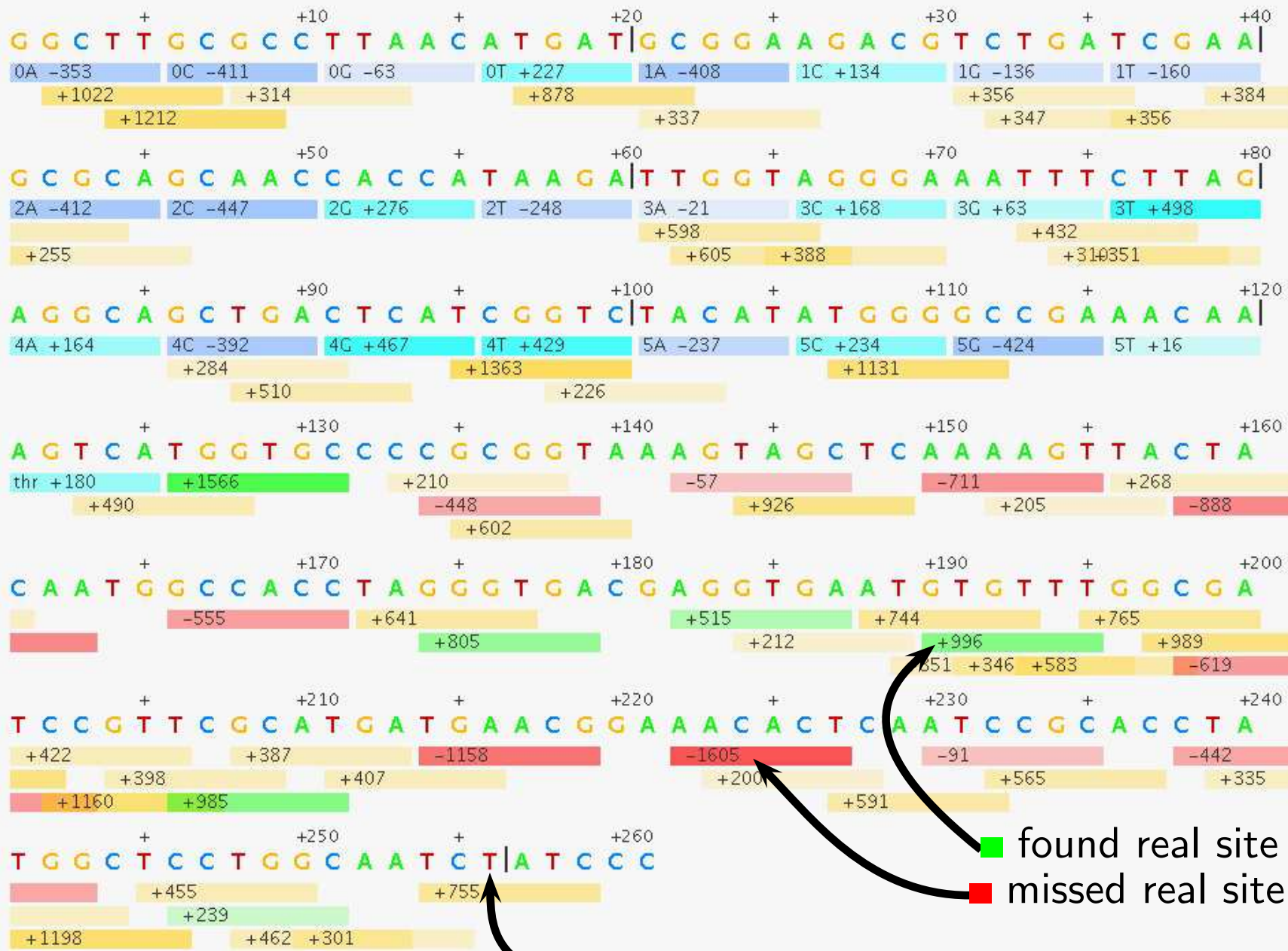
“blue”  
gene  
weight  
matrix:  
6 bp  
wide

$\gamma = 16$   
binding  
sites

found real site

Genome positions available  $G = 256$  bases  
 $R_{frequency} = \log_2 256/16 = 4$  bits

# Unevolved Ev Creature



“blue”  
gene  
weight  
matrix:  
6 bp  
wide

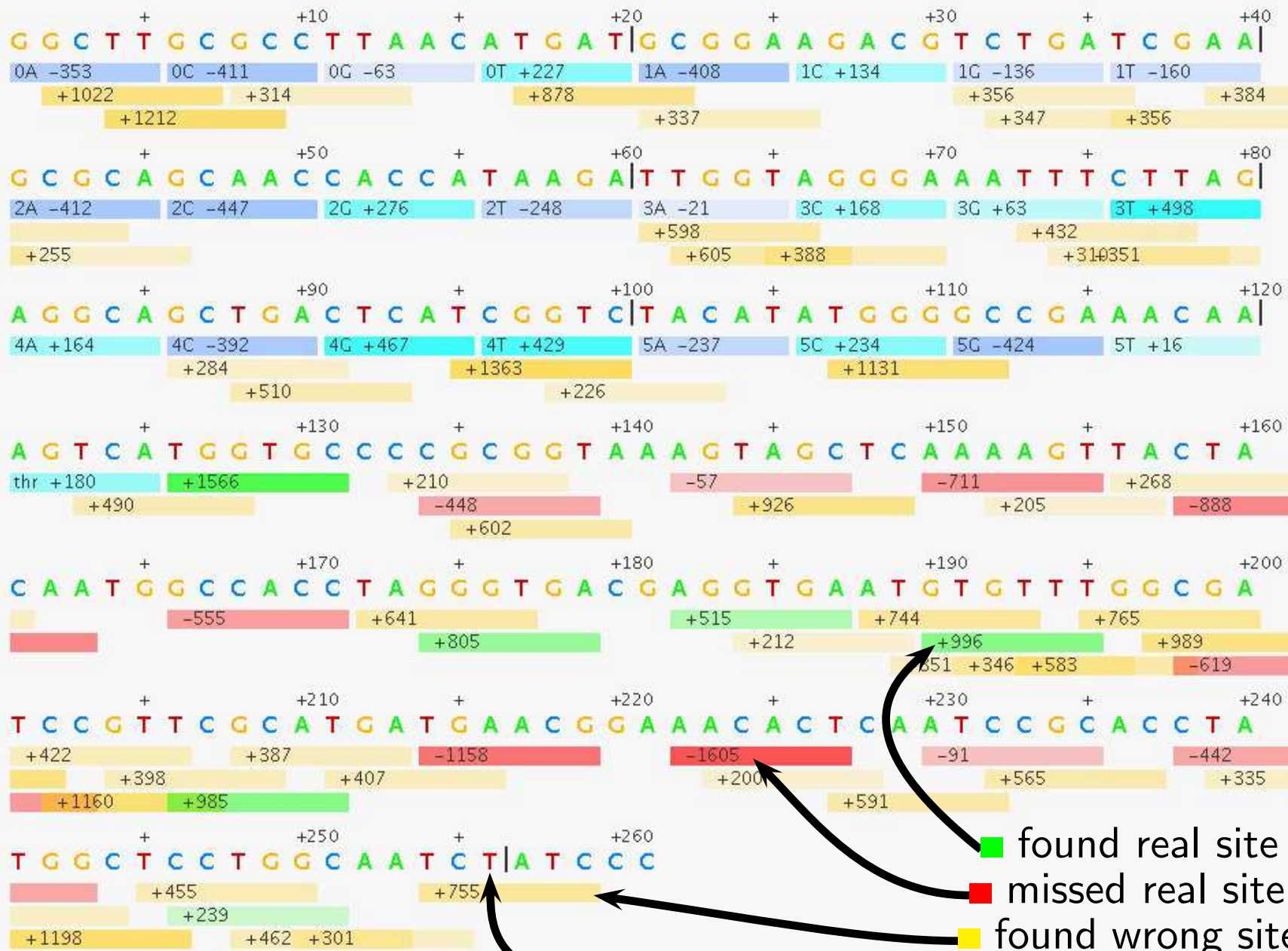
$\gamma = 16$   
binding  
sites

■ found real site  
■ missed real site

Genome positions available  $G = 256$  bases  
 $R_{frequency} = \log_2 256/16 = 4$  bits



# Unevolved Ev Creature



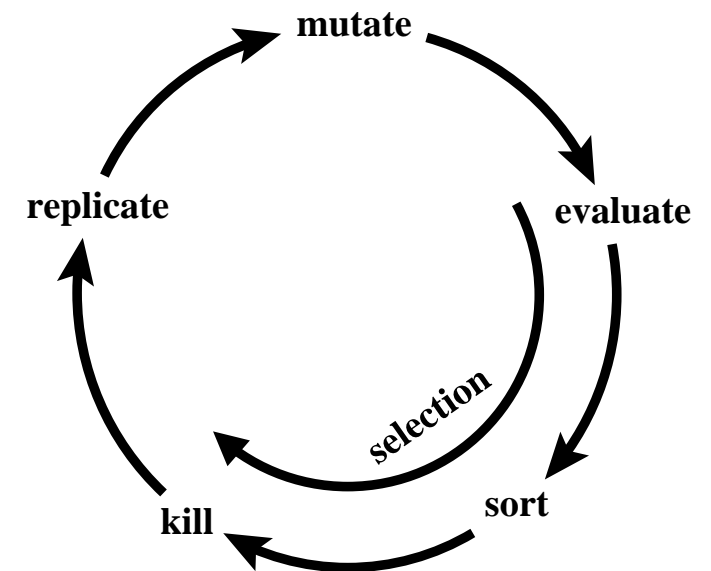
“blue”  
gene  
weight  
matrix:  
6 bp  
wide

$\gamma = 16$   
binding  
sites

Genome positions available  $G = 256$  bases  
 $R_{frequency} = \log_2 256/16 = 4$  bits

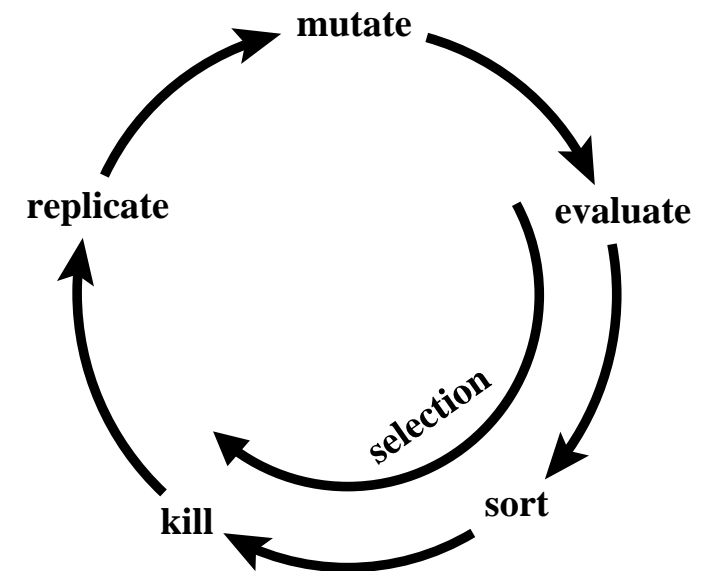
# Evolution Cycle

- EVALUATE each creature



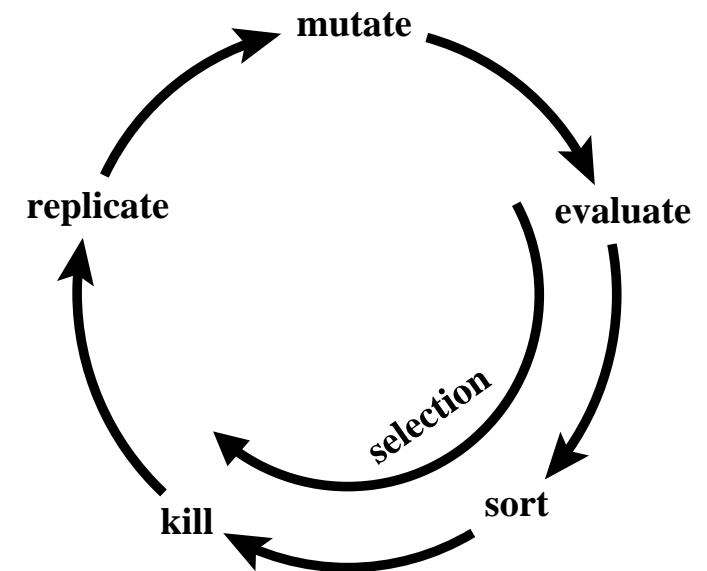
# Evolution Cycle

- EVALUATE each creature
  - translate the recognizer gene into a weight matrix



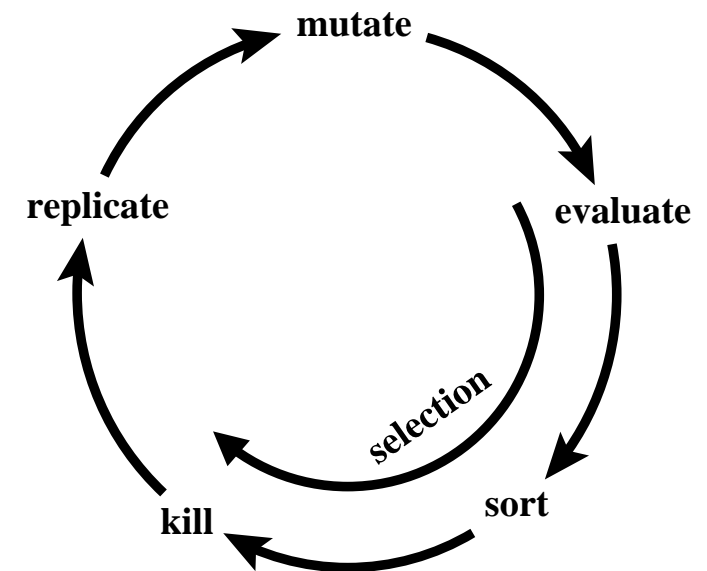
# Evolution Cycle

- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome



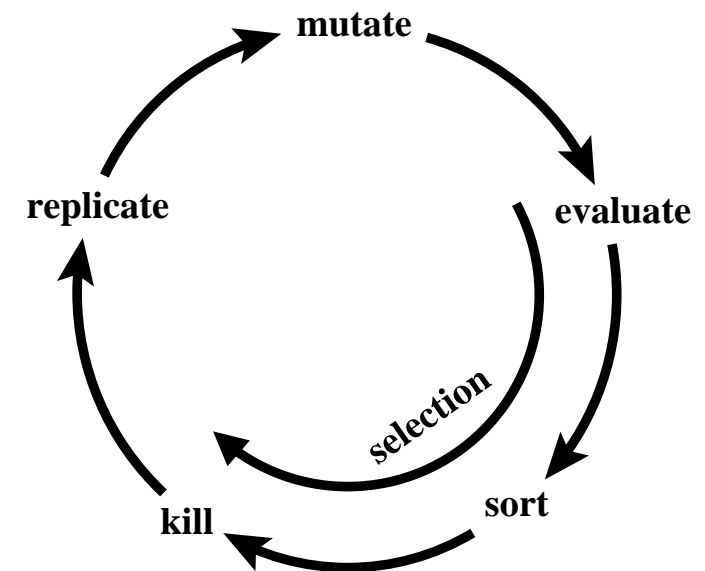
# Evolution Cycle

- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome
  - count the number of mistakes:



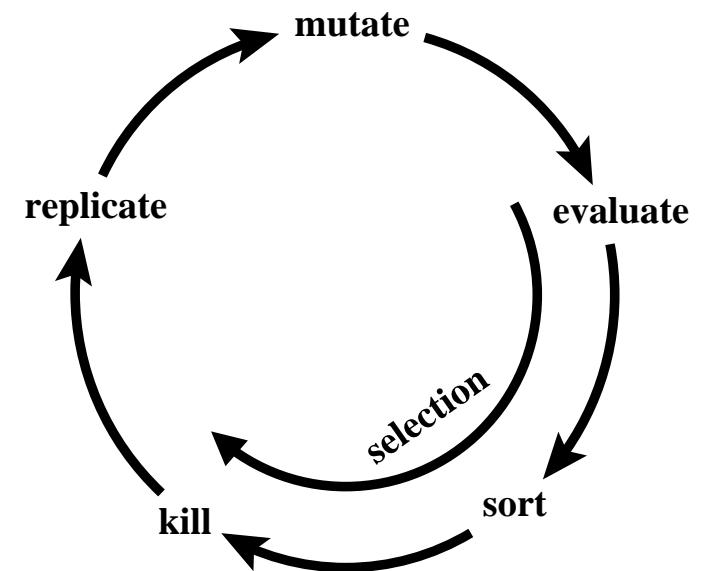
# Evolution Cycle

- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome
  - count the number of mistakes:
    - missing a site at a right place



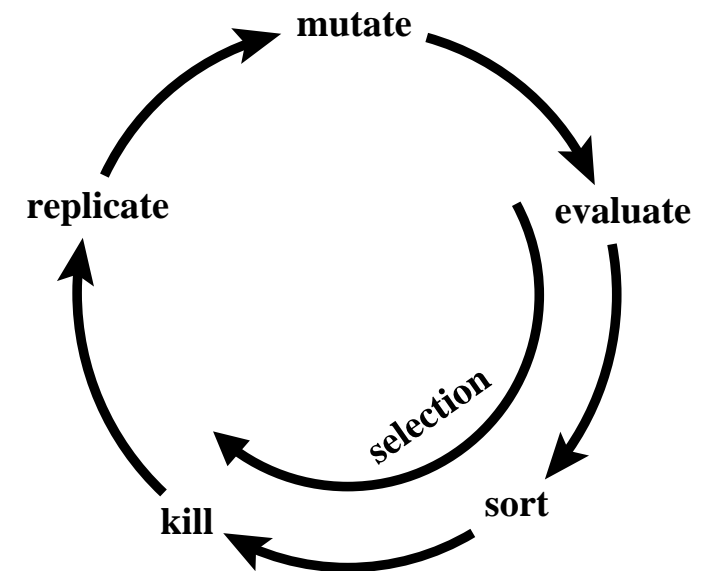
# Evolution Cycle

- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome
  - count the number of mistakes:
    - missing a site at a right place
    - finding a site at a wrong place



# Evolution Cycle

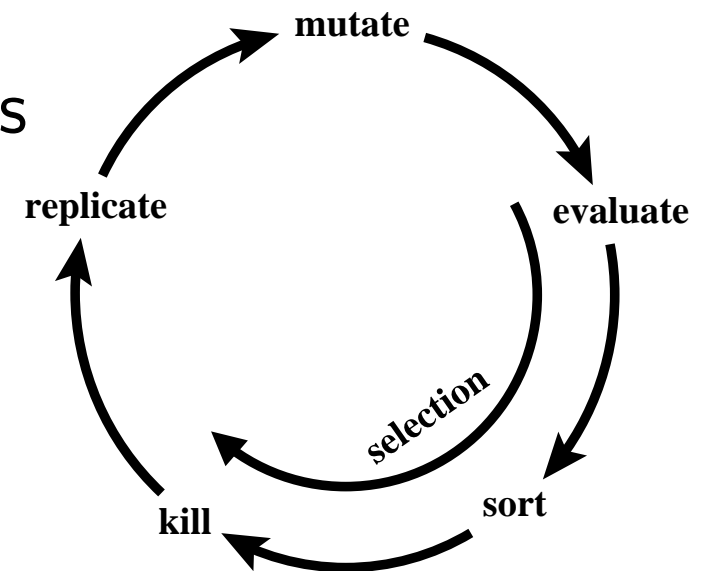
- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome
  - count the number of mistakes:
    - missing a site at a right place
    - finding a site at a wrong place
  - Sort the creatures by their mistakes





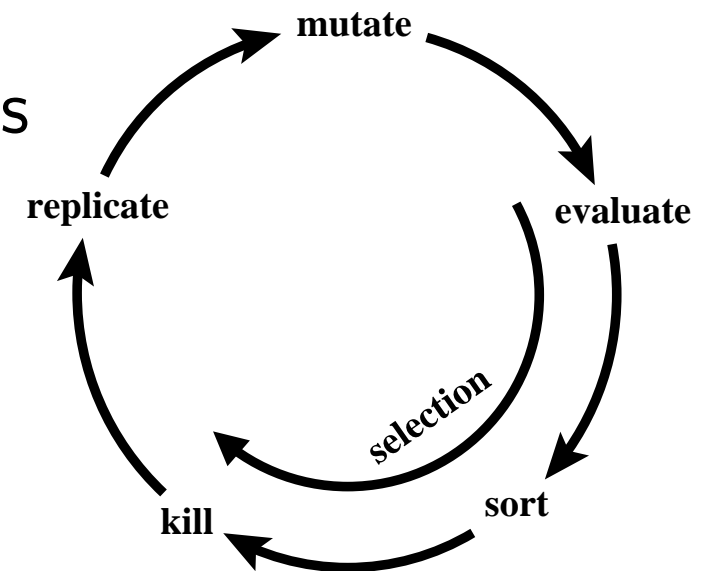
# Evolution Cycle

- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome
  - count the number of mistakes:
    - missing a site at a right place
    - finding a site at a wrong place
  - Sort the creatures by their mistakes
- REPLICATE: the best creatures are duplicated and replace the worst ones

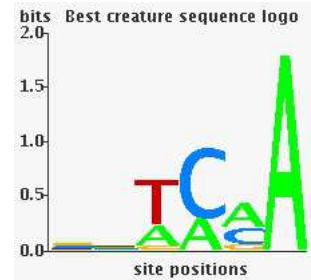


# Evolution Cycle

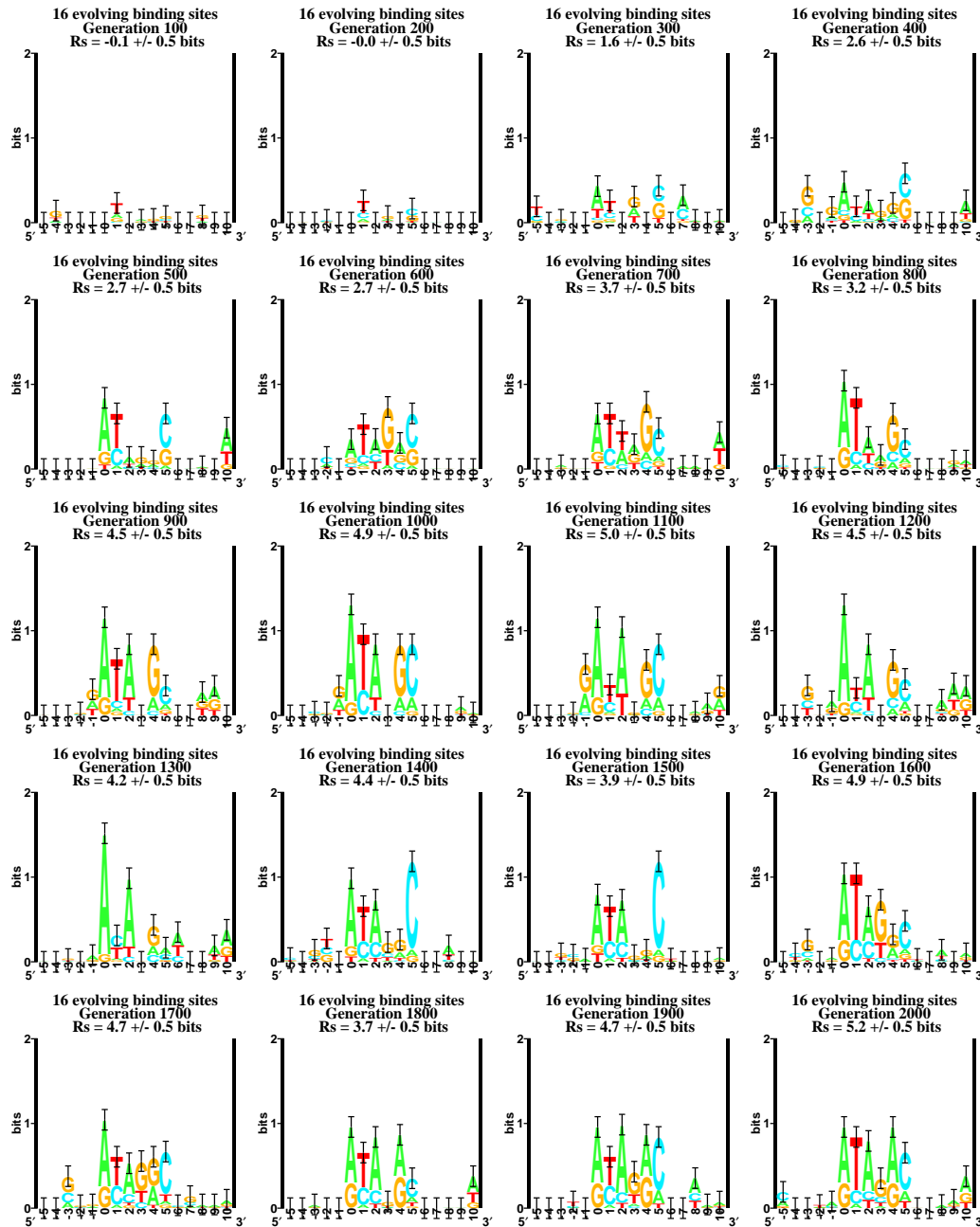
- EVALUATE each creature
  - translate the recognizer gene into a weight matrix
  - scan the weight matrix across the genome
  - count the number of mistakes:
    - missing a site at a right place
    - finding a site at a wrong place
  - Sort the creatures by their mistakes
- REPLICATE: the best creatures are duplicated and replace the worst ones
- MUTATE all genomes randomly



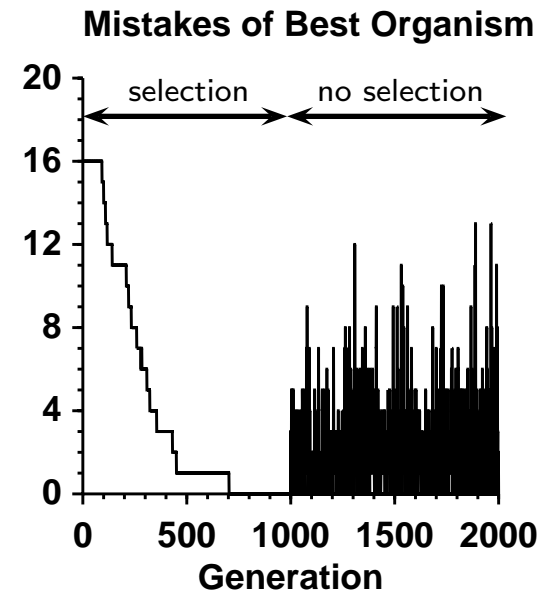
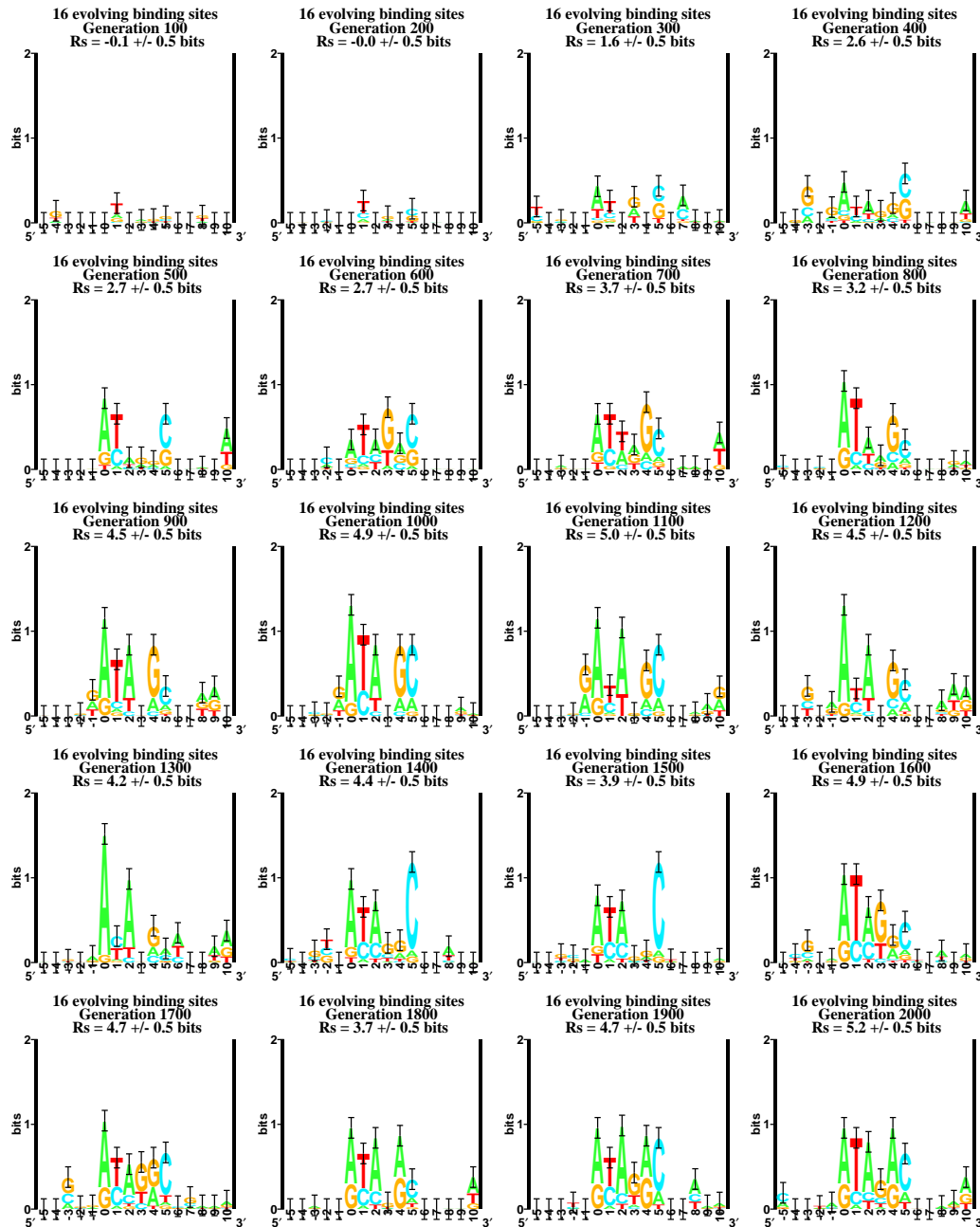
# Evolved Ev Creature



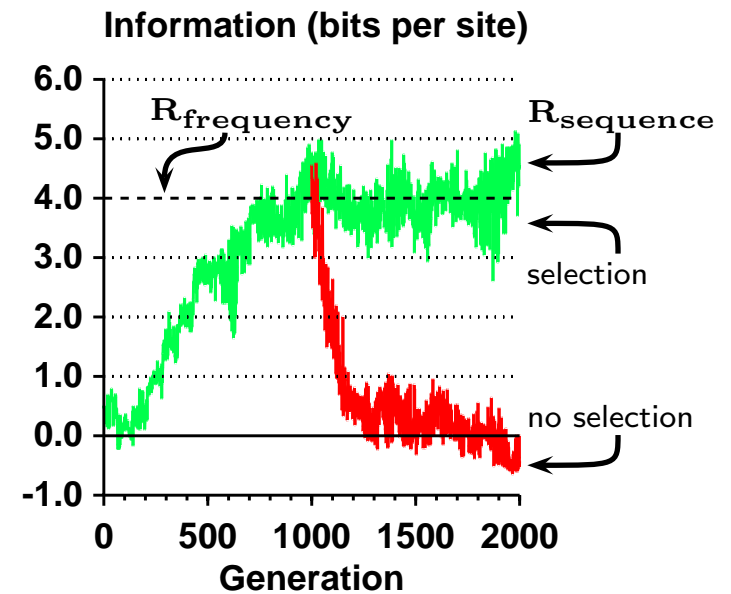
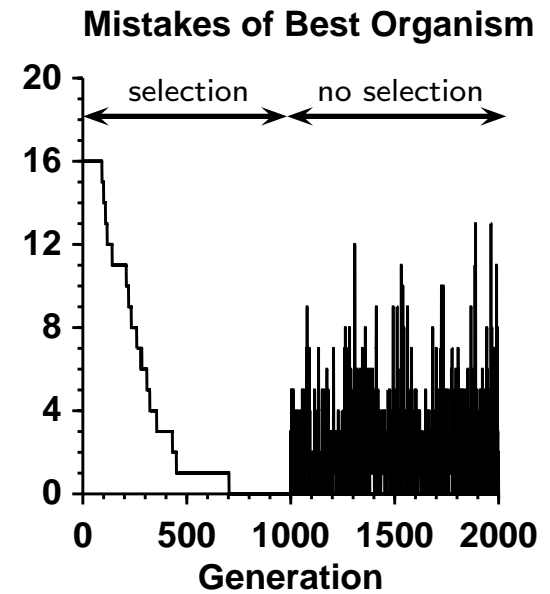
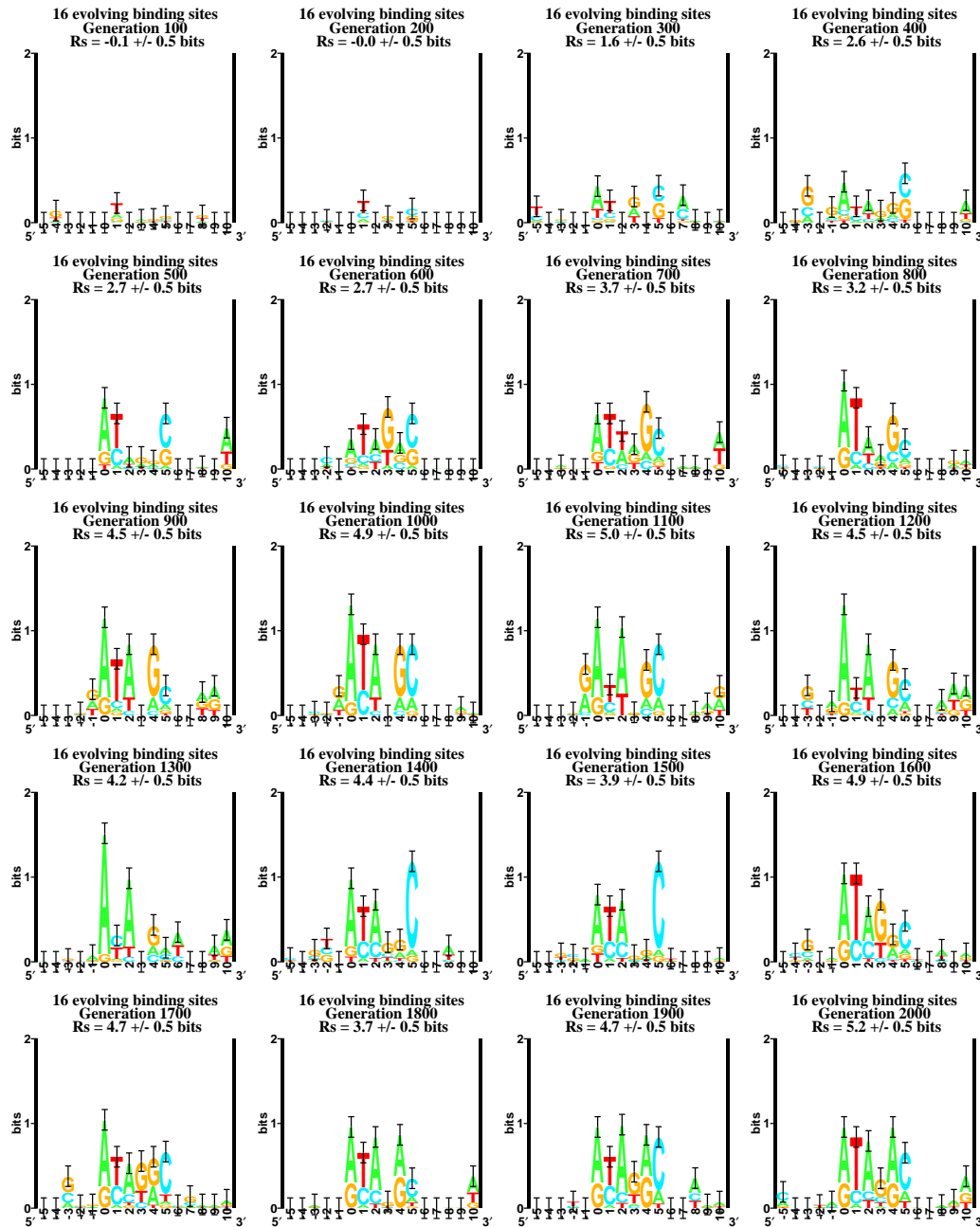
# Evolution of Binding Sites



# Evolution of Binding Sites



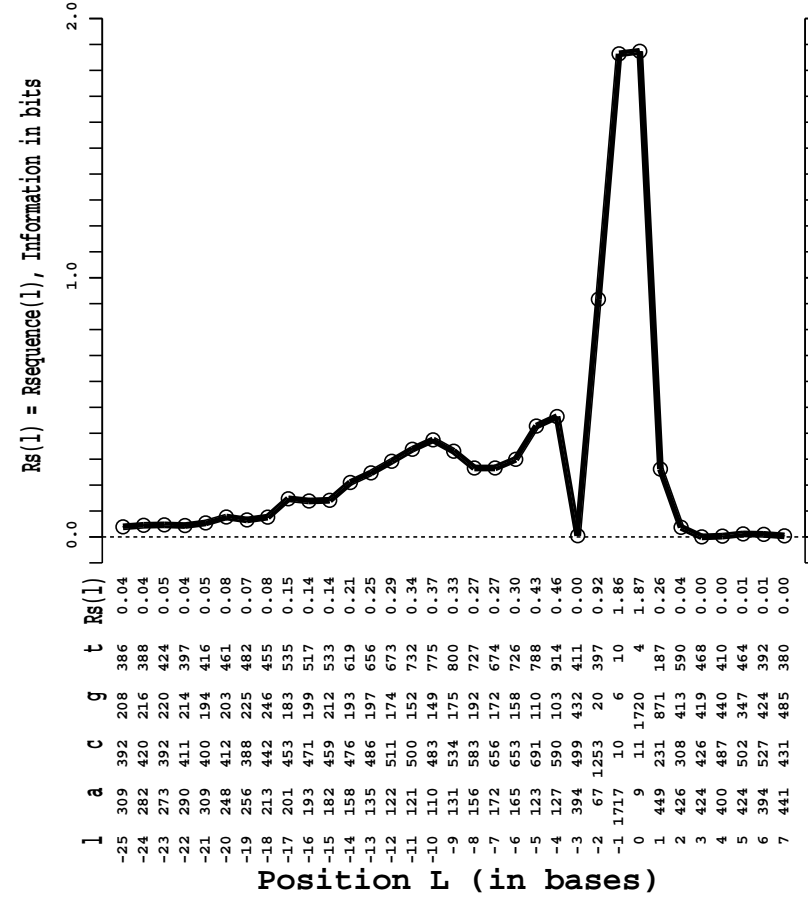
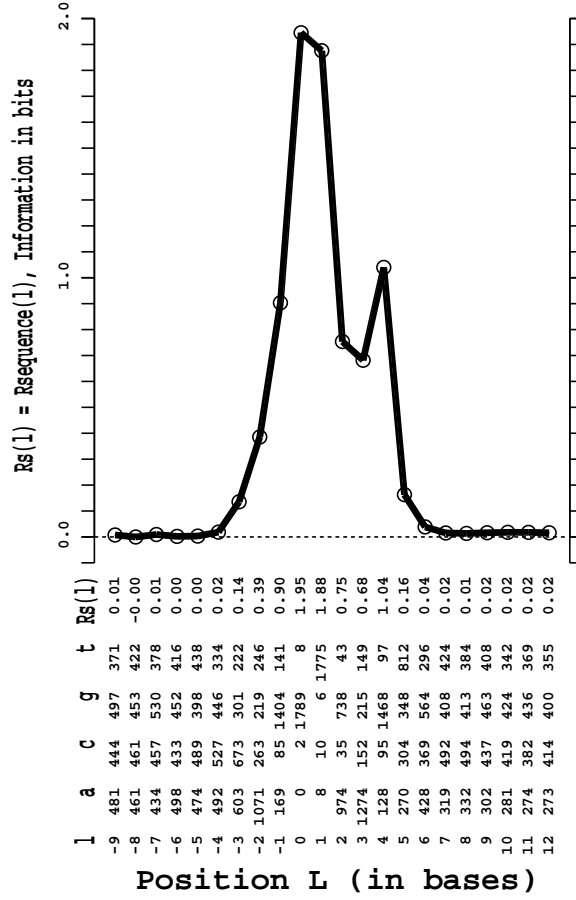
# Evolution of Binding Sites





# Human Splice Junction Information Curves

Sequence Conservation  $\rightarrow$   
in bits per base **Donor**



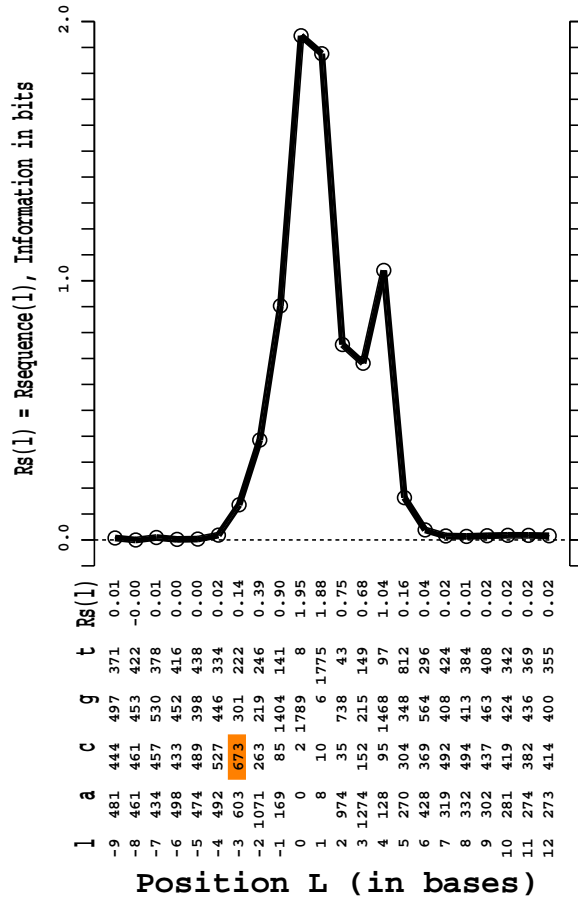
**Acceptor**

- The consensus sequences match ...



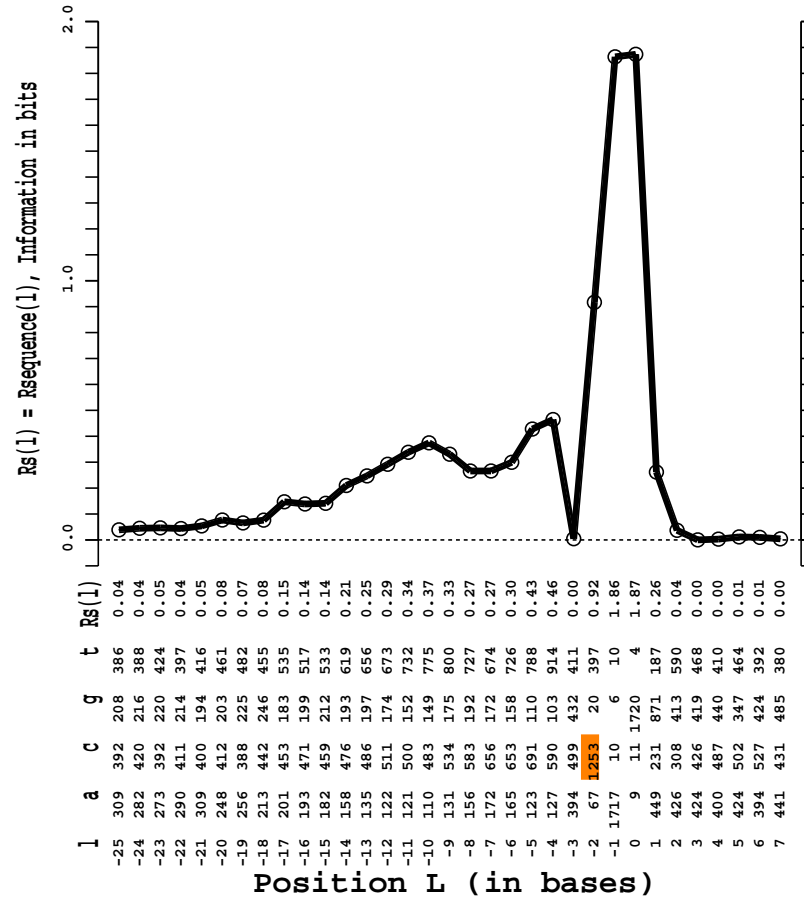
# Human Splice Junction Information Curves

Sequence Conservation →  
in bits per base **Donor**



C

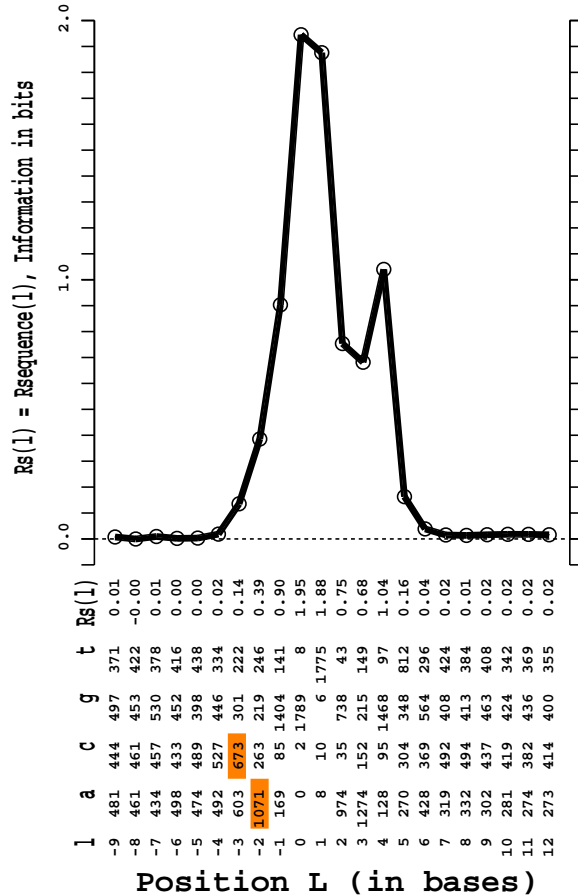
**Acceptor**



- The consensus sequences match ...

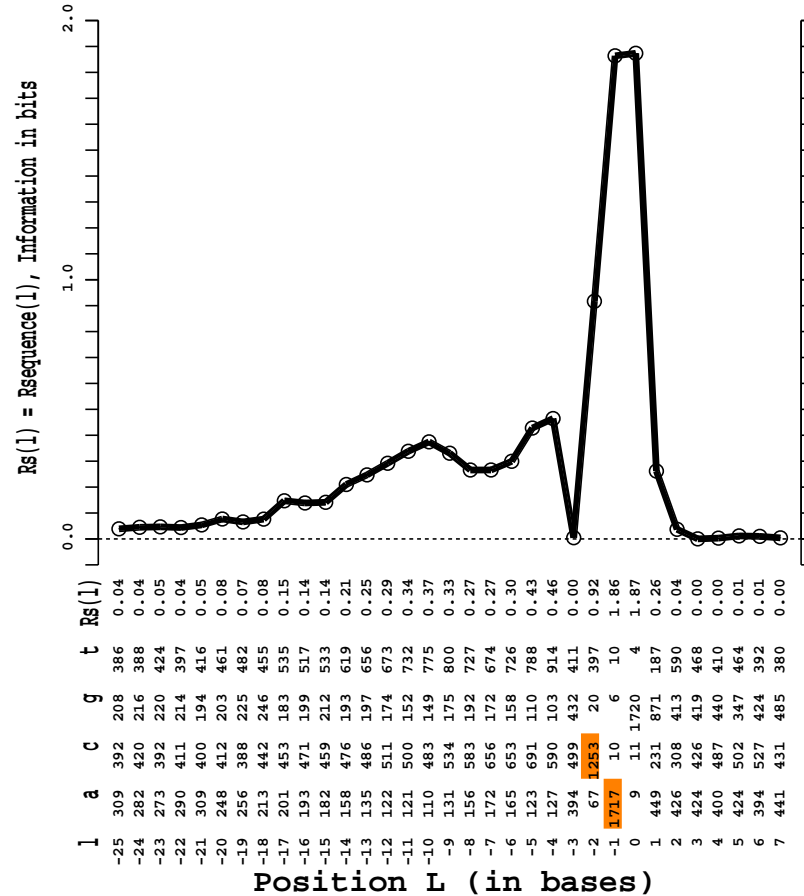
# Human Splice Junction Information Curves

Sequence Conservation →  
in bits per base **Donor**



**C A**

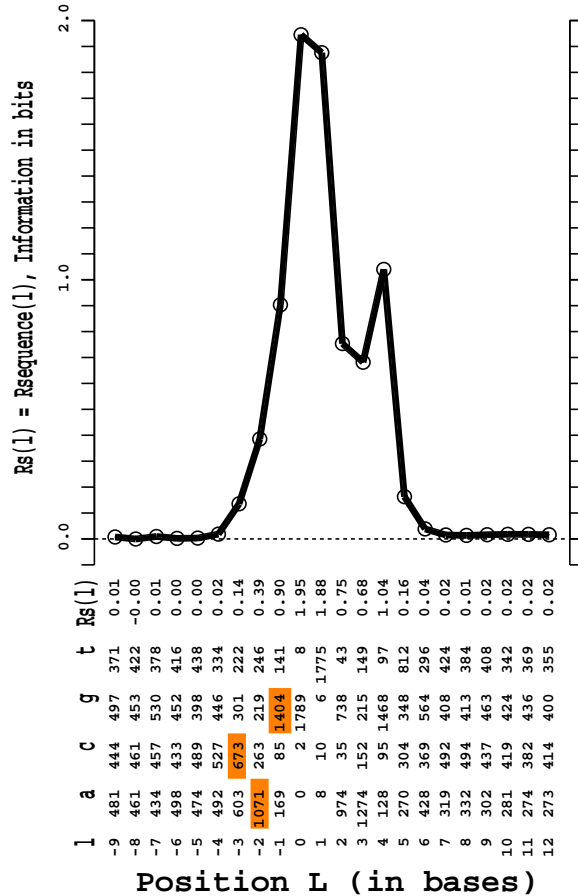
**Acceptor**



- The consensus sequences match ...

# Human Splice Junction Information Curves

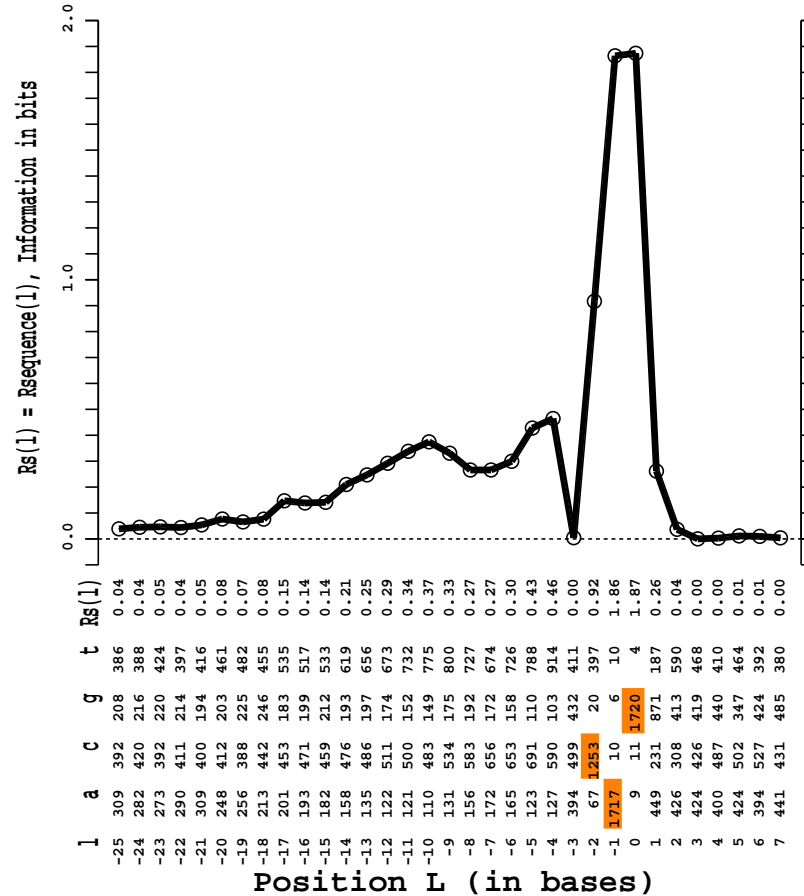
Sequence Conservation →  
in bits per base **Donor**



C A G

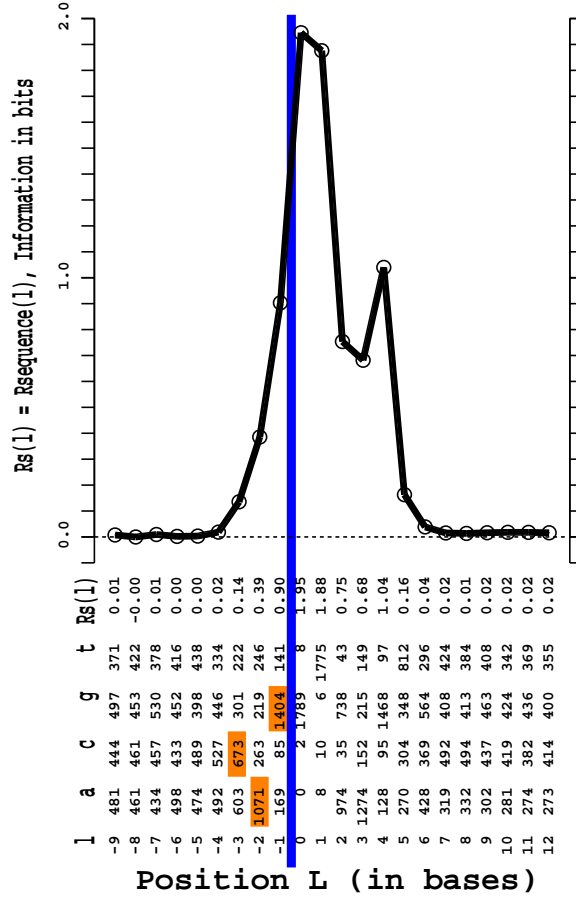
- The consensus sequences match ...

**Acceptor**



# Human Splice Junction Information Curves

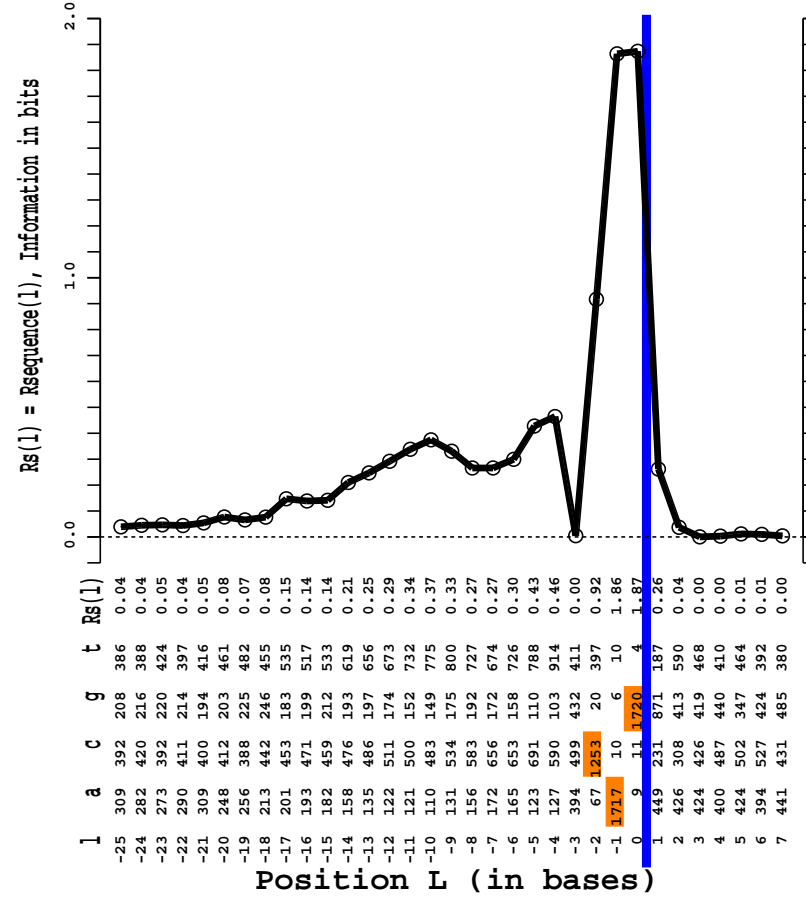
Sequence Conservation →  
in bits per base **Donor**



C A G —

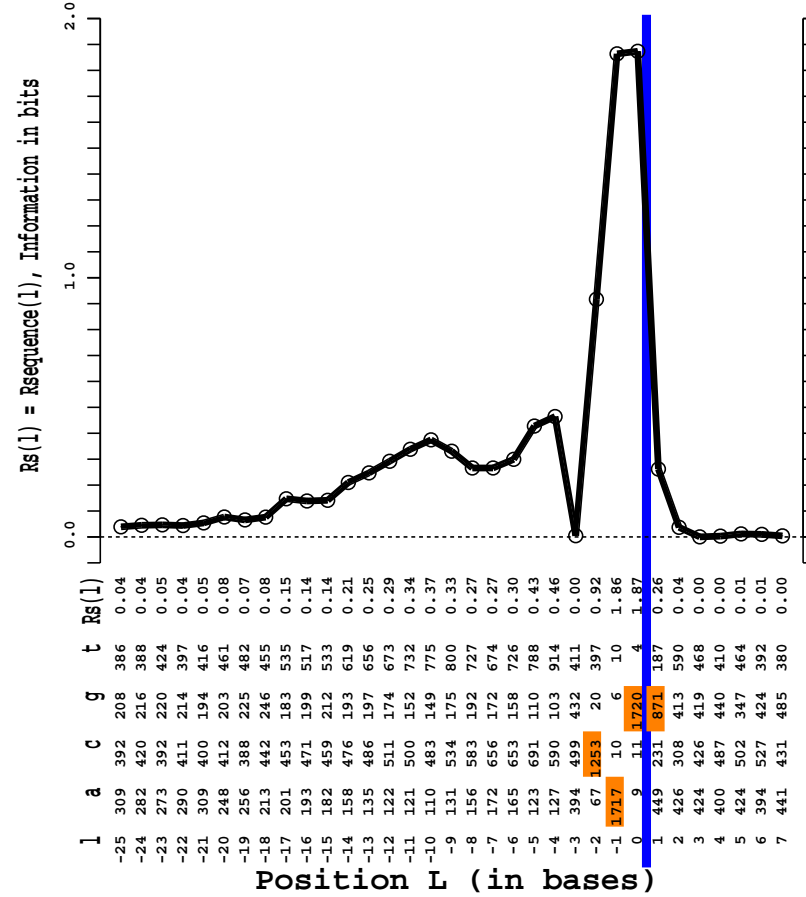
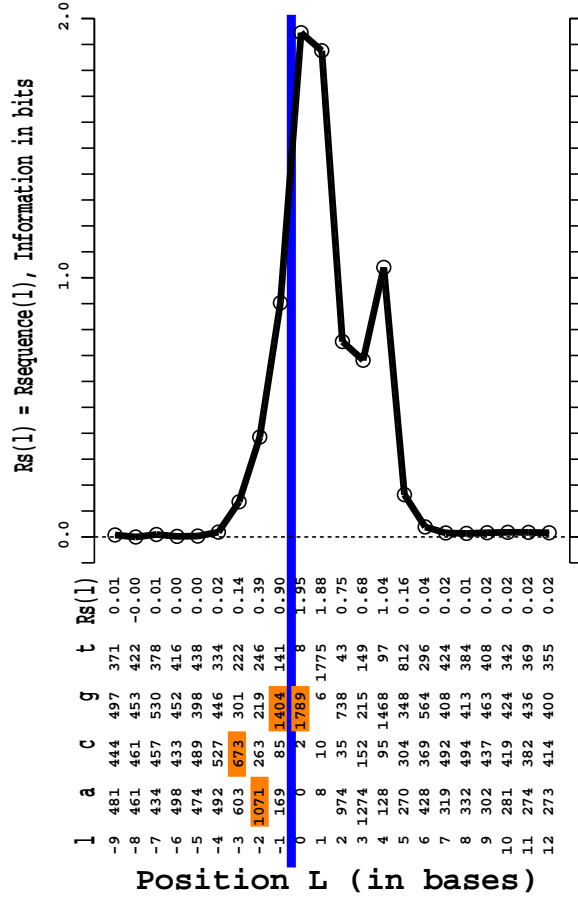
- The consensus sequences match ...

**Acceptor**



# Human Splice Junction Information Curves

Sequence Conservation →  
in bits per base **Donor**



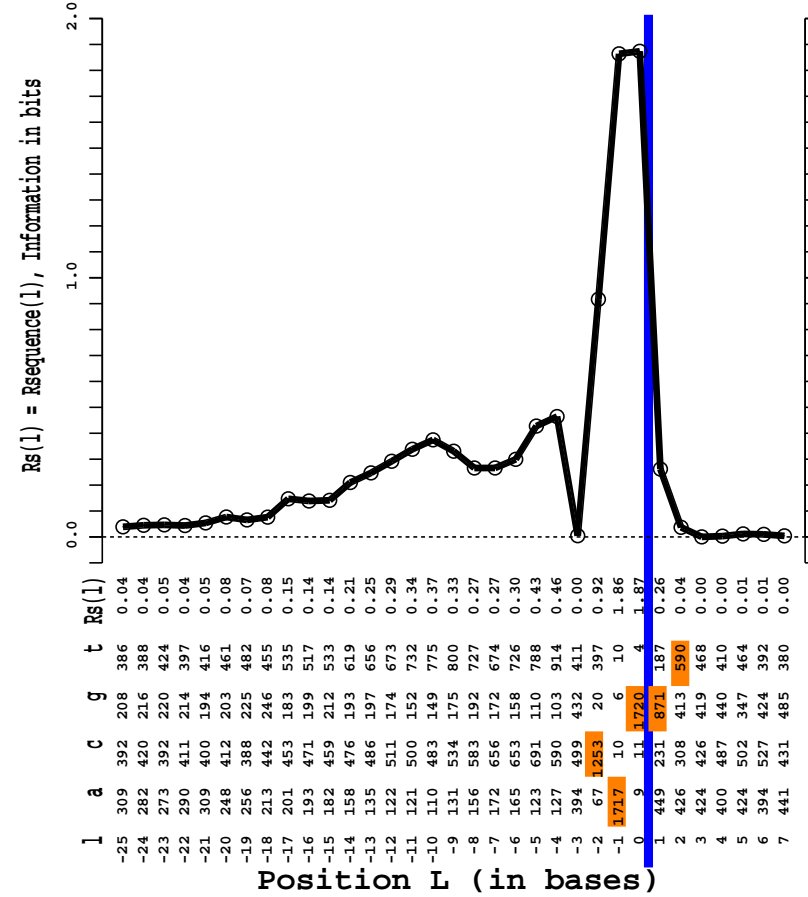
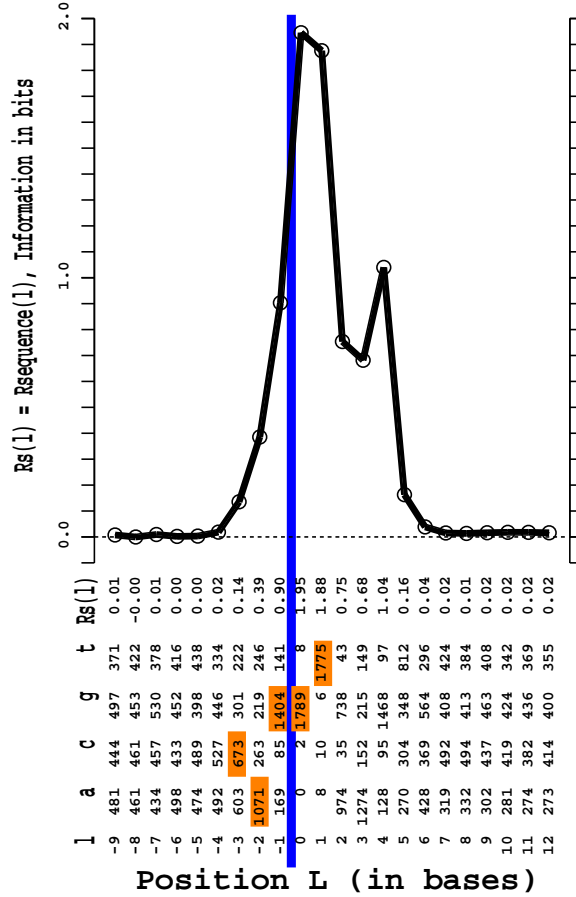
**Acceptor**

C A G — G

- The consensus sequences match ...

# Human Splice Junction Information Curves

Sequence Conservation →  
in bits per base **Donor**



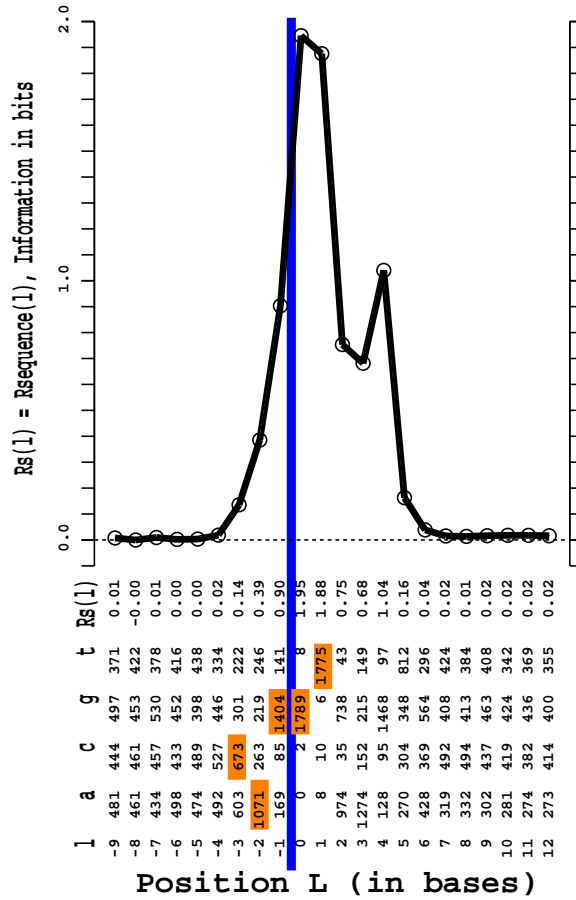
**Acceptor**

C A G — G T

- The consensus sequences match ...

# Human Splice Junction Information Curves

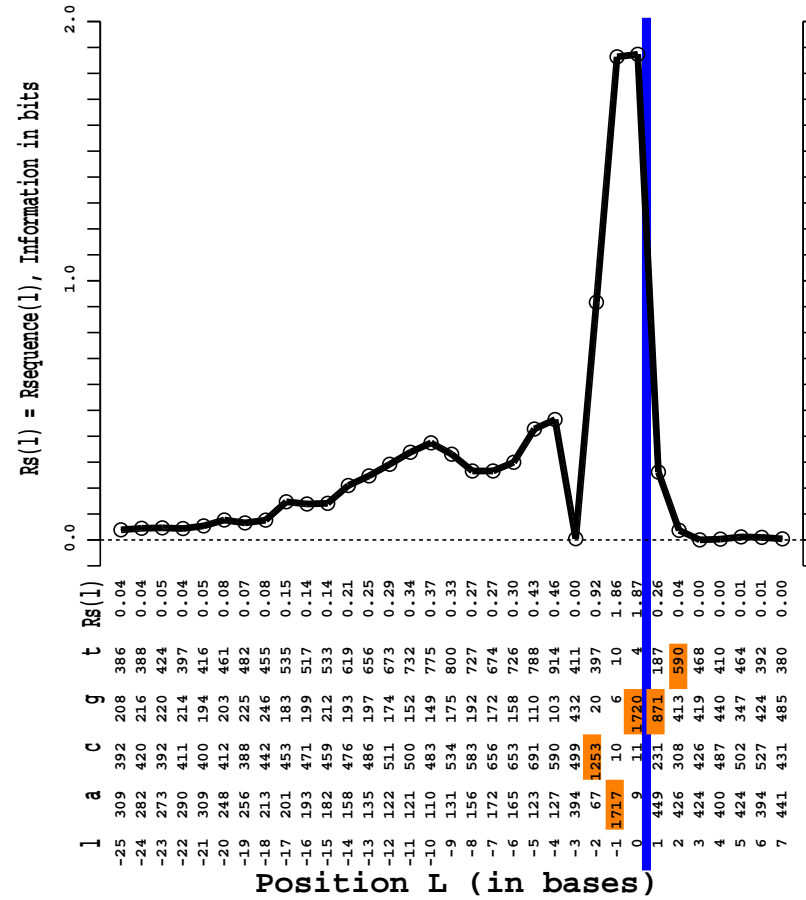
Sequence Conservation →  
in bits per base **Donor**



C A G — G T

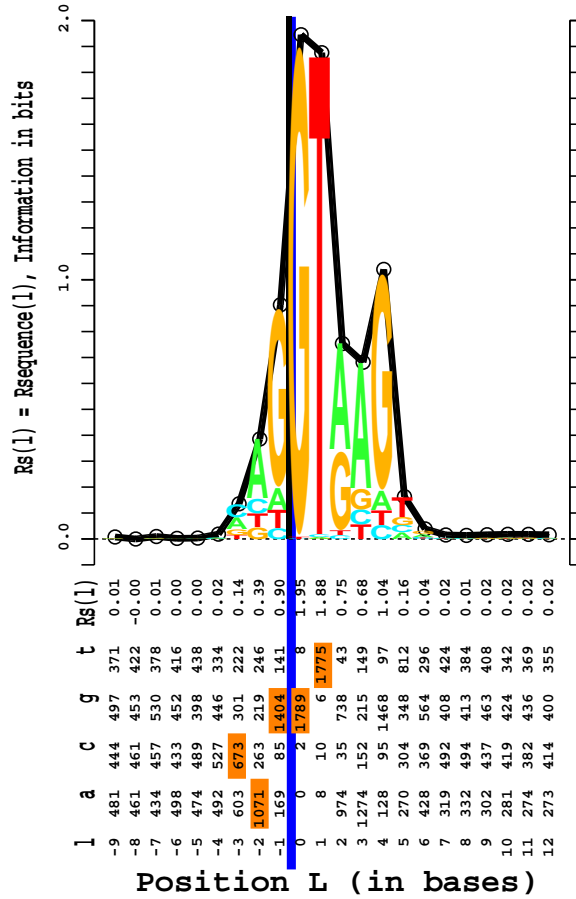
- The consensus sequences match ...
- BUT the information curves (sequence conservation) differ!

**Acceptor**



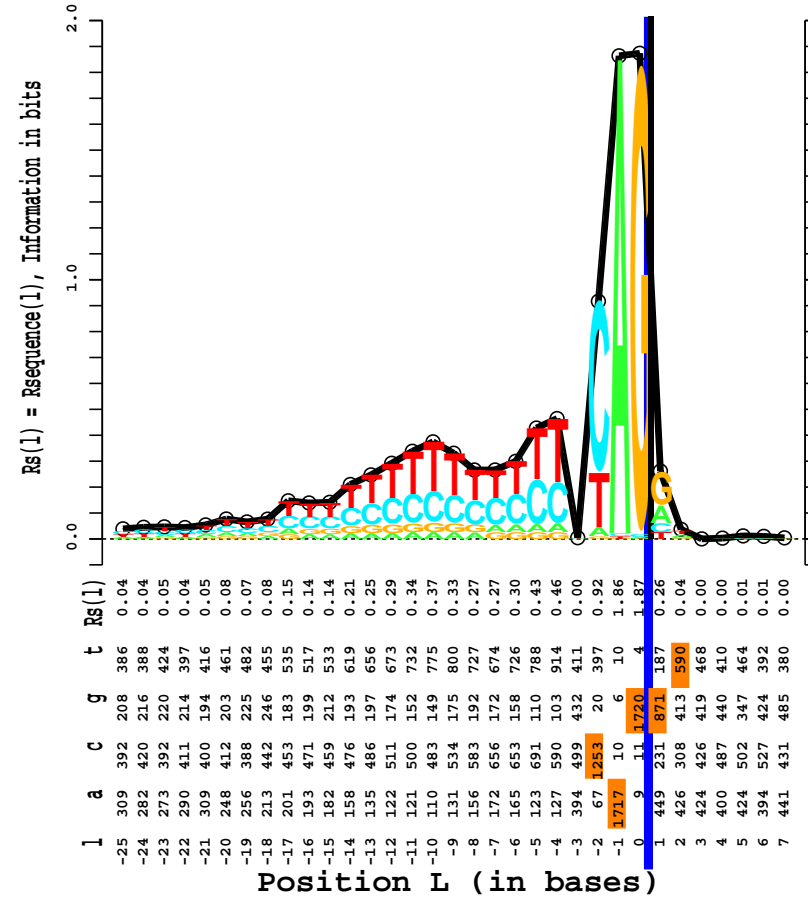
# Human Splice Junction Information Curves

Sequence Conservation →  
in bits per base Donor



C A G — G T

- The consensus sequences match ...
- BUT the information curves (sequence conservation) differ!
- Put letters into the graph proportional to their frequency!

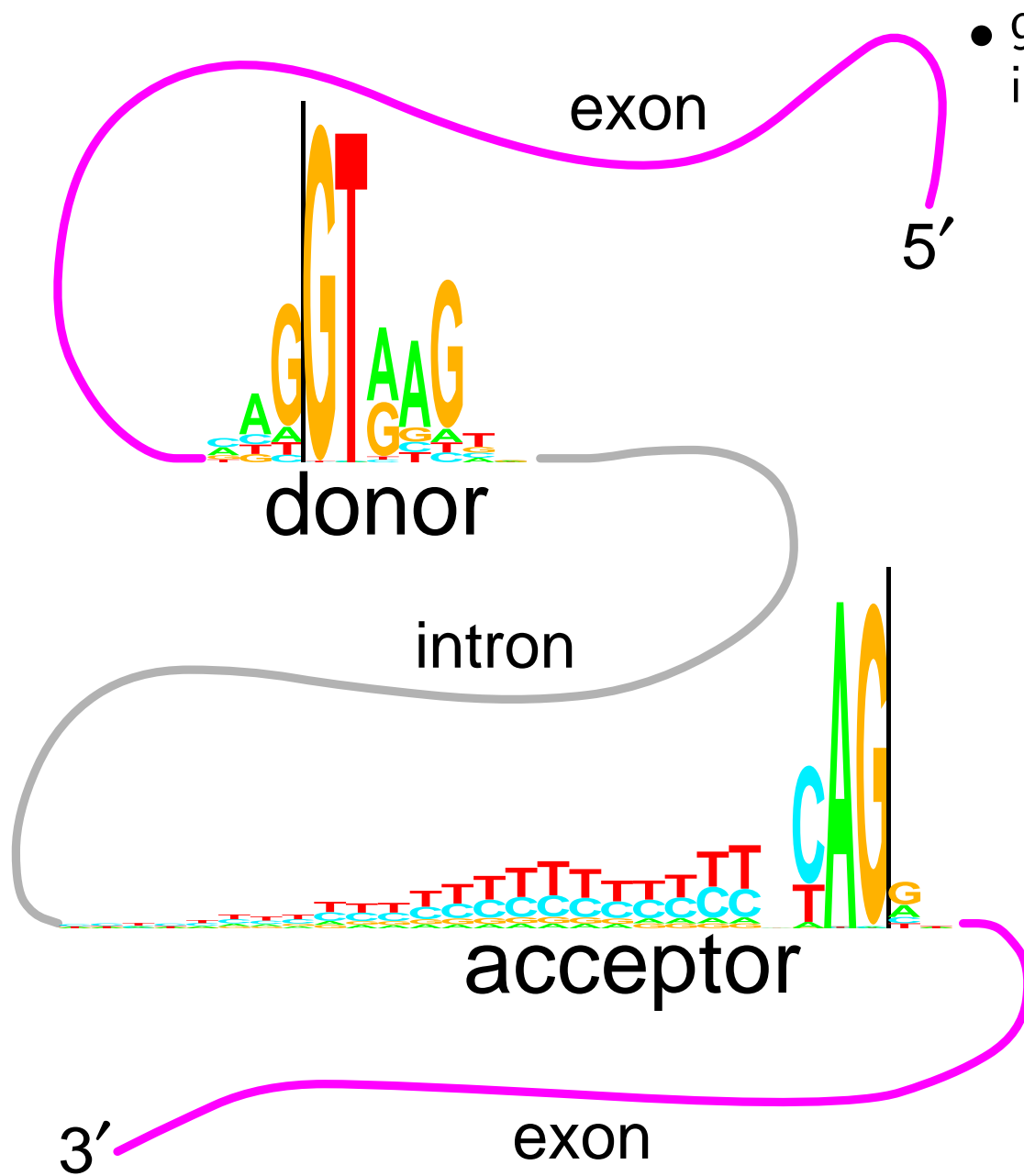


Acceptor



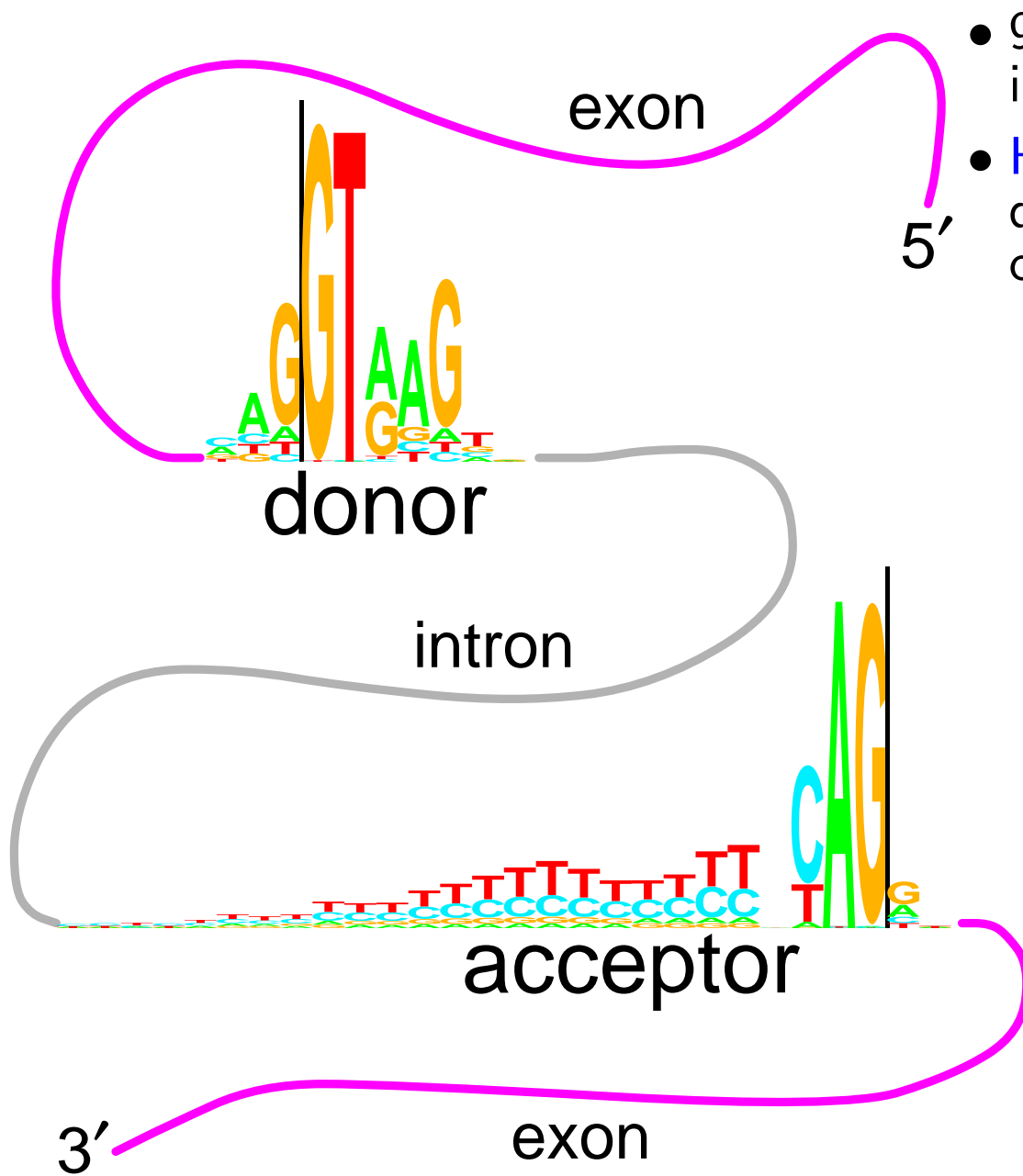


# Splice Junction Sequence Logos

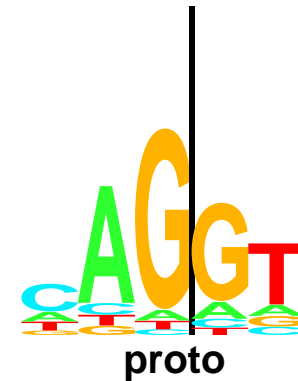


- 90% of the splice junction information is on the intron side

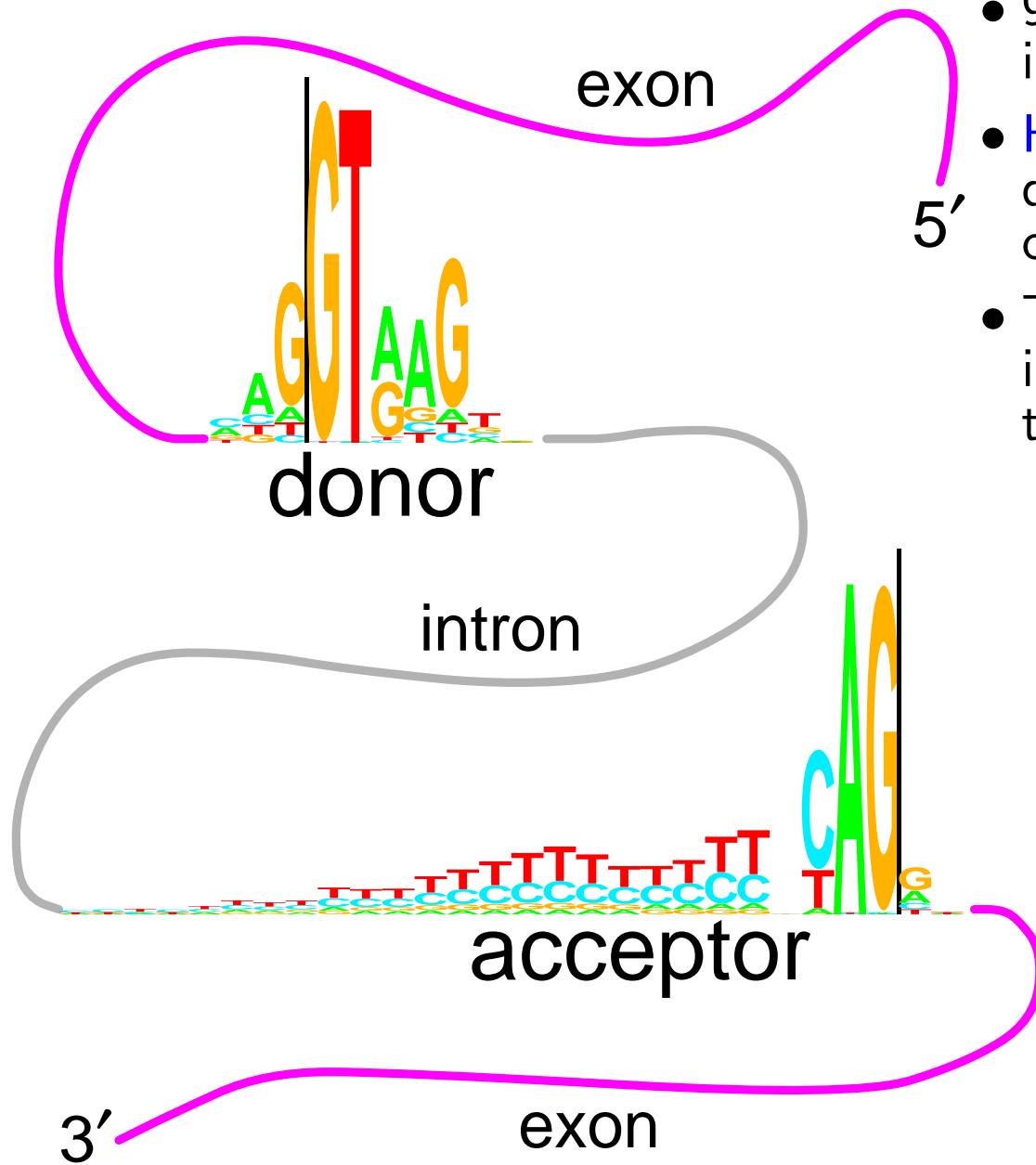
# Splice Junction Sequence Logos



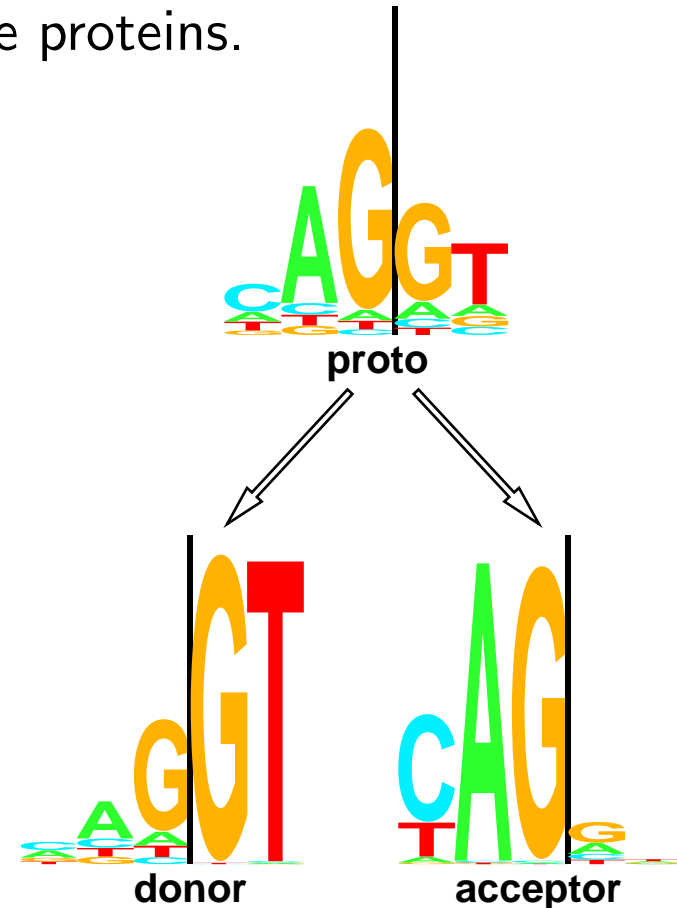
- 90% of the splice junction information is on the intron side
- **Hypothesis:** donor and acceptor sites had a common ancestor that duplicated



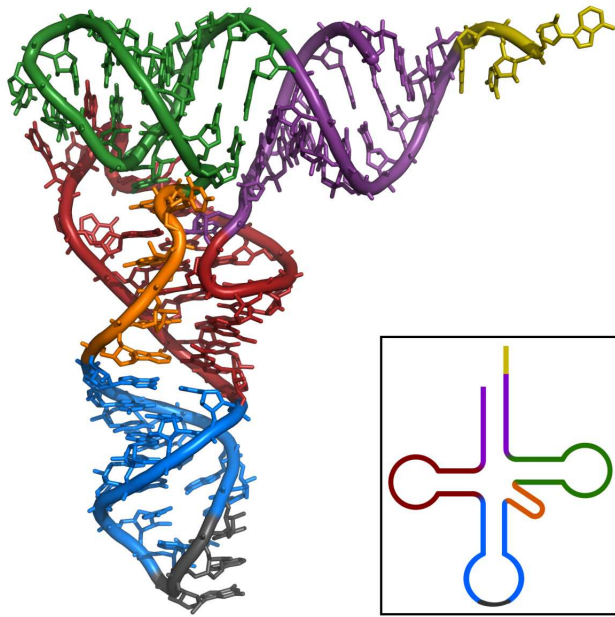
# Splice Junction Sequence Logos



- 90% of the splice junction information is on the intron side
- **Hypothesis:** donor and acceptor sites had a common ancestor that duplicated
- They evolved to put the information into the intron. This avoids affecting the proteins.

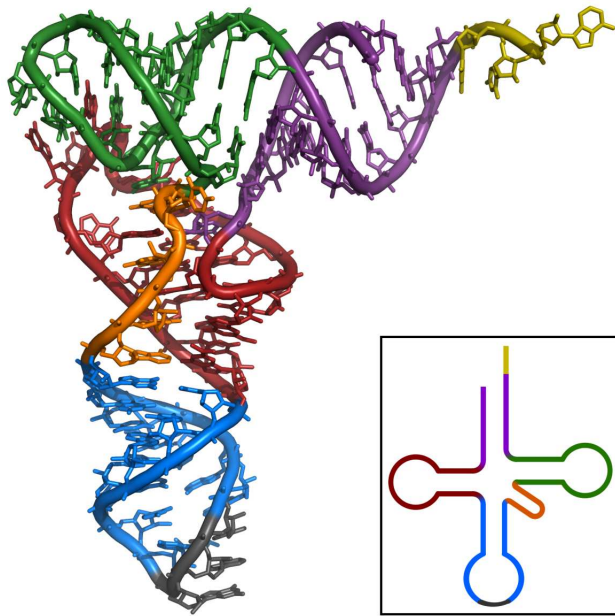


# 3D Sequence Logos for tRNA Correlations



- tRNA reads RNA to make protein

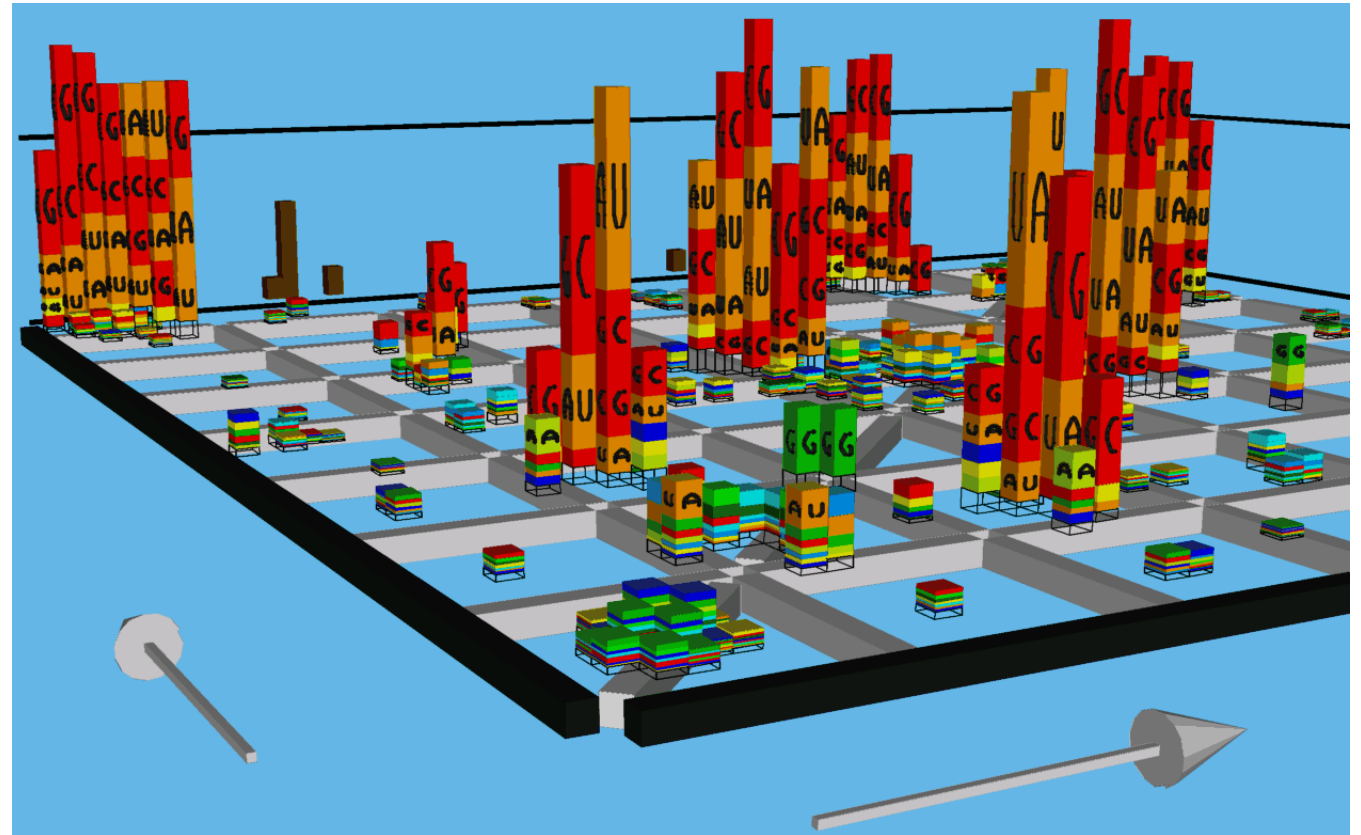
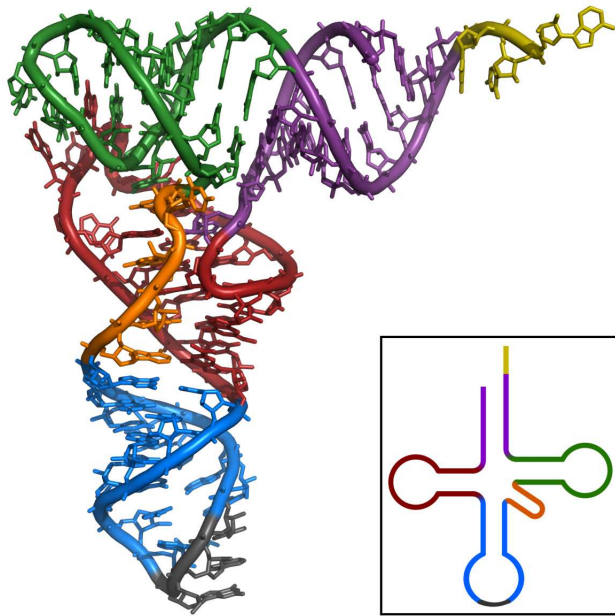
# 3D Sequence Logos for tRNA Correlations



- tRNA reads RNA to make protein
- Correlations can be measured in bits!



# 3D Sequence Logos for tRNA Correlations



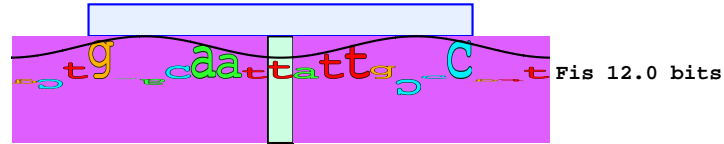
- tRNA reads RNA to make protein
- Correlations can be measured in bits!
- 3D Sequence logo
- **OBSERVED:** tRNA stems



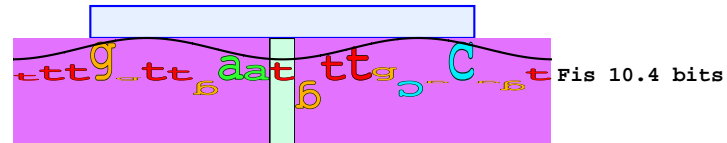
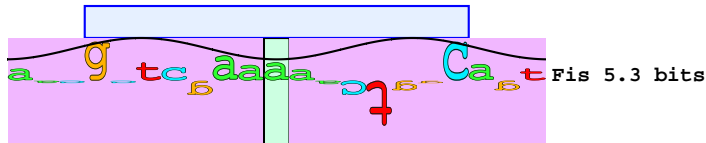
# Sequence Walker example: *rrnB* P1

*rrnB* P1

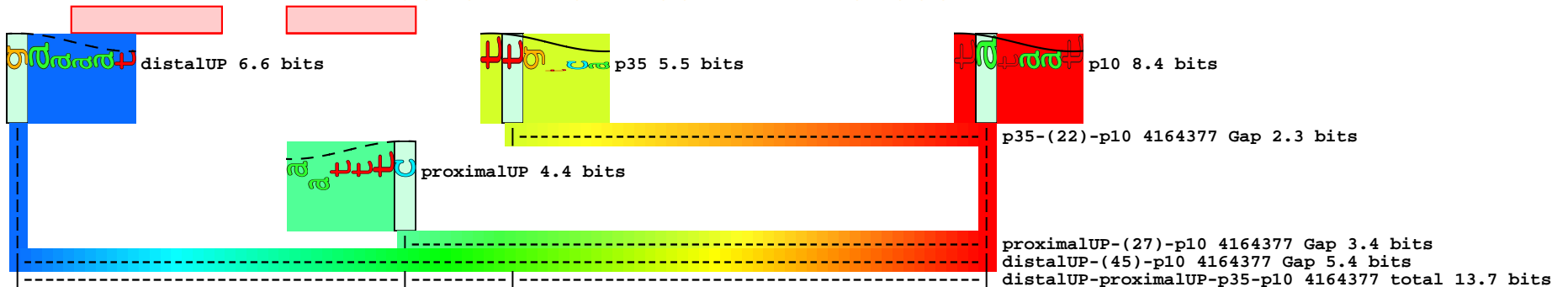
5' *g g a g c t g a a c a a t t a t t g c c c g t t t t a c a g c g t t a c g g c t t c g a* 3'  
 3' *c c t c g a c t t g t t a a t a a c g g g c a a a a t g t c g c a a t g c c g a a g c t* 5'



\*4164280 . \*4164290 . \*4164300 . \*4164310 . \*4164320 . \*4164330  
 5' *a a c g c t c g a a a a a c t g g c a g t t t t a g g c t g a t t t g g t t g a a t g t t g c g c g g t c a* 3'  
 3' *t t g c g a g c t t t t t g a c c g t c a a a a t c c g a c t a a a c c a a c t t a c a a c g c g c c a g t* 5'



\*4164340 . \*4164350 . \*4164360 . \*4164370 . \*4164380  
 5' *g a a a a t t a t t t t a a a t t t c c t c t t g t c a g g c c g g a a t a a c t c c c t a t a a t g* 3'  
 3' *c t t t t a a t a a a a t t t a a a g g a g a a c a g t c c g g c c t t a t t g a g g a t a t t a c* 5'



# Complex Sequence Walker Example

- $\sigma^{70}$  promoters have a  $-35$  and a  $-10$

# Complex Sequence Walker Example

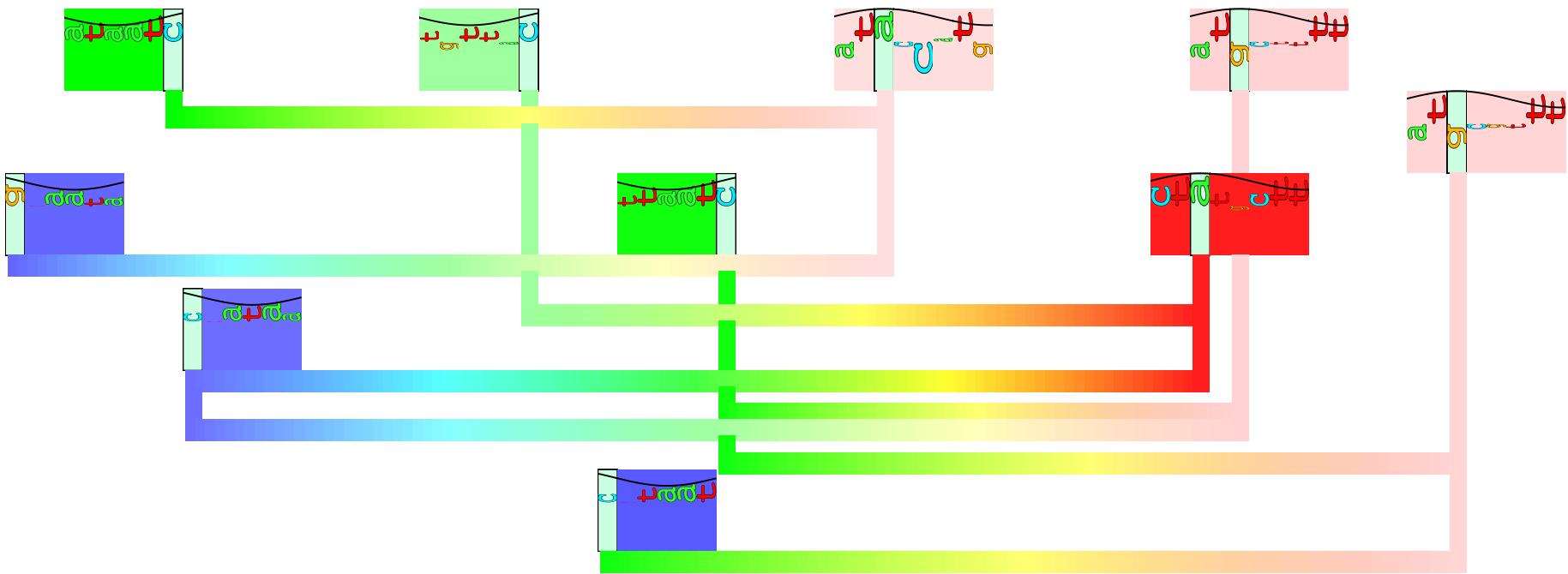
- $\sigma^{70}$  promoters have a  $-35$  and a  $-10$
- Using information theory we discovered that stress-response  $\sigma^{38}$  promoters do not have a  $-35$

# Complex Sequence Walker Example

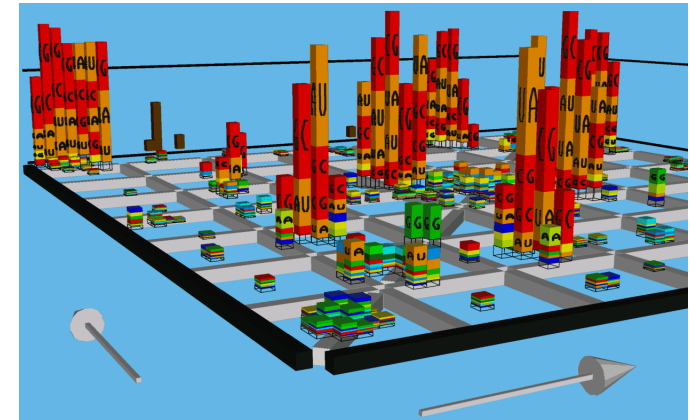
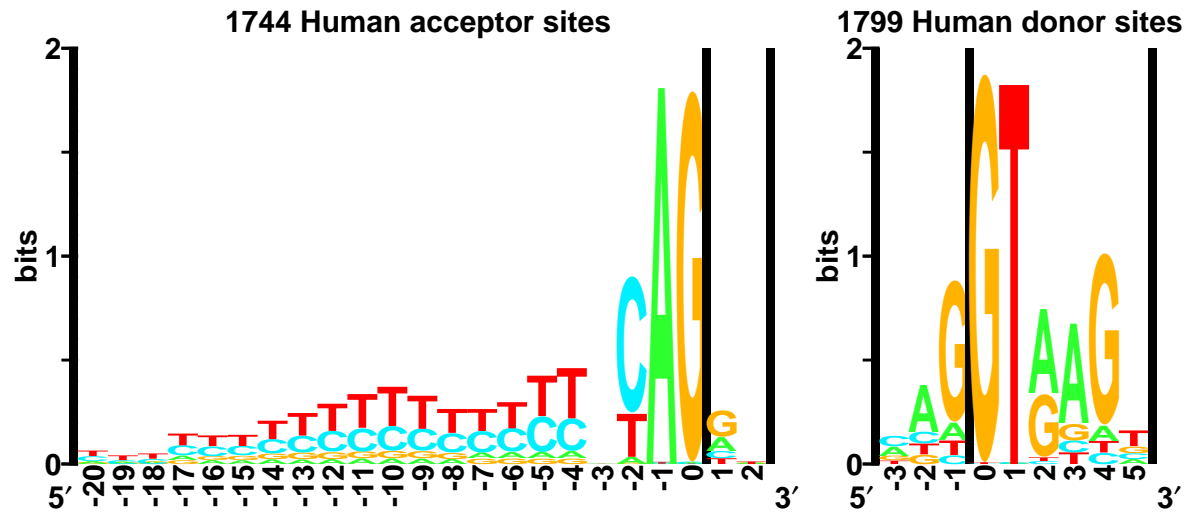
- $\sigma^{70}$  promoters have a  $-35$  and a  $-10$
- Using information theory we discovered that stress-response  $\sigma^{38}$  promoters do not have a  $-35$
- Instead, they have a  $-10$  and two UP elements

# Complex Sequence Walker Example

- $\sigma^{70}$  promoters have a  $-35$  and a  $-10$
- Using information theory we discovered that stress-response  $\sigma^{38}$  promoters do not have a  $-35$
- Instead, they have a  $-10$  and two UP elements
- $\sigma^{38}$  promoter *talA* P1 is complex!

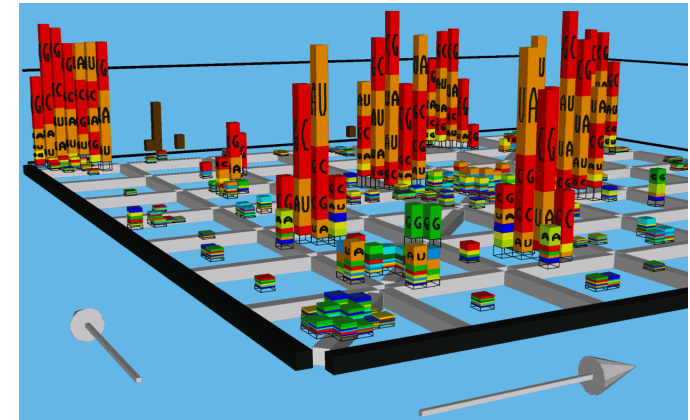
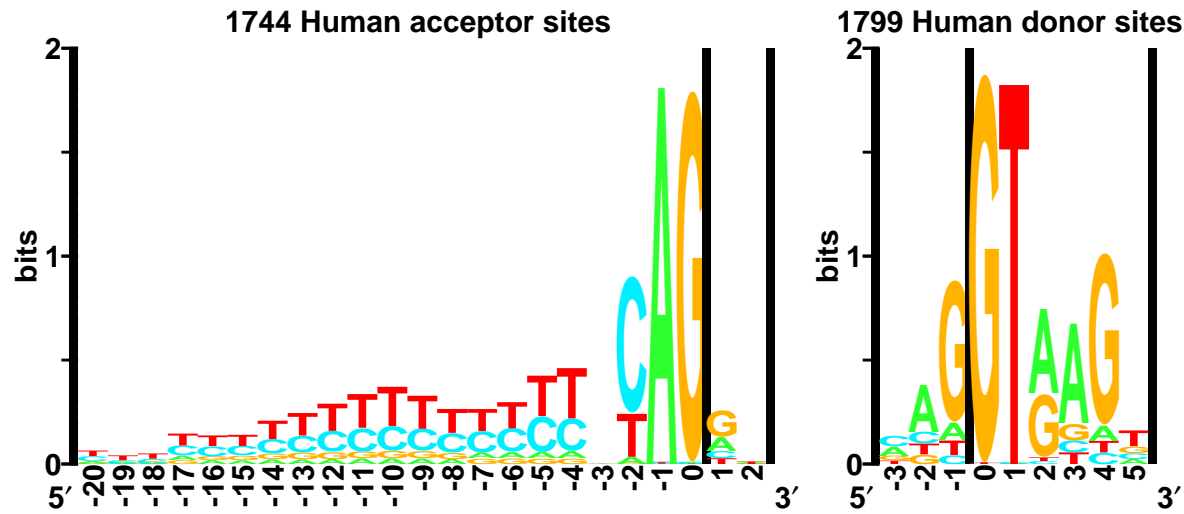


# Important Discovery



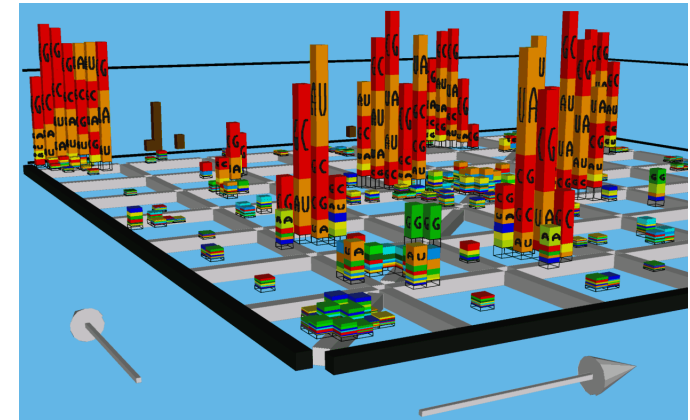
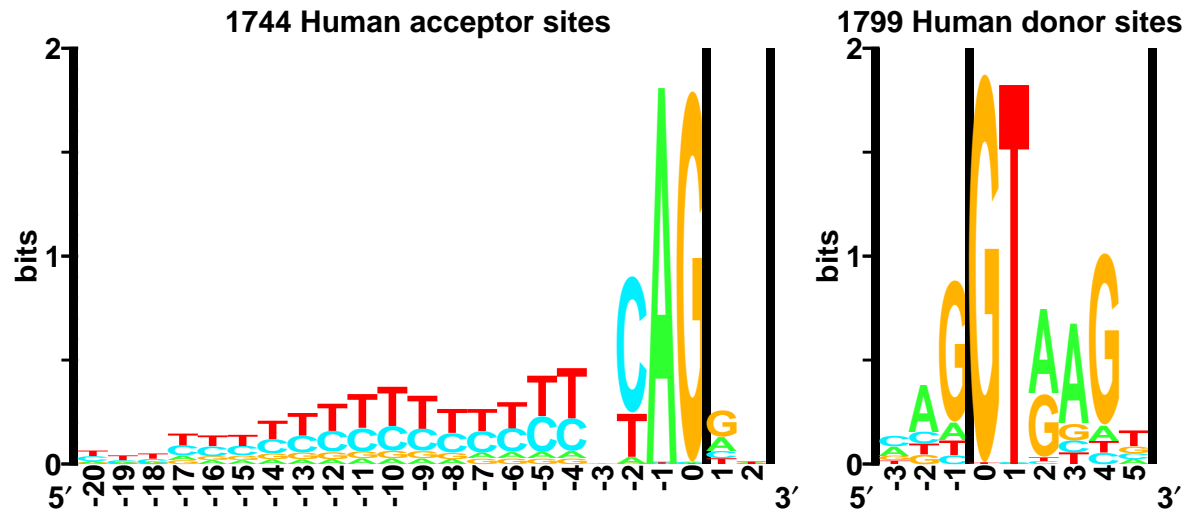
- Area under a sequence logo is the total information.

# Important Discovery



- Area under a sequence logo is the total information.
- How is that related to the binding energy?

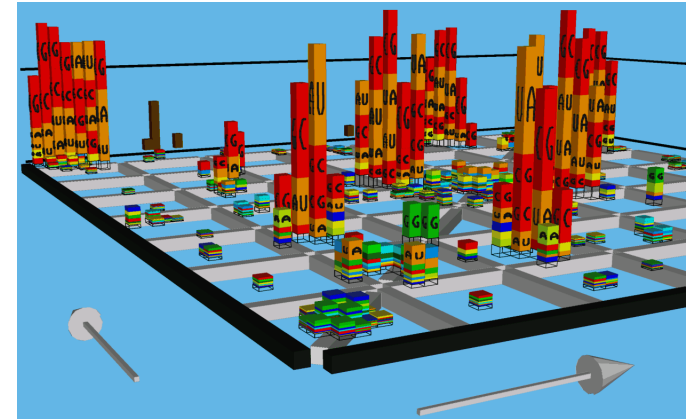
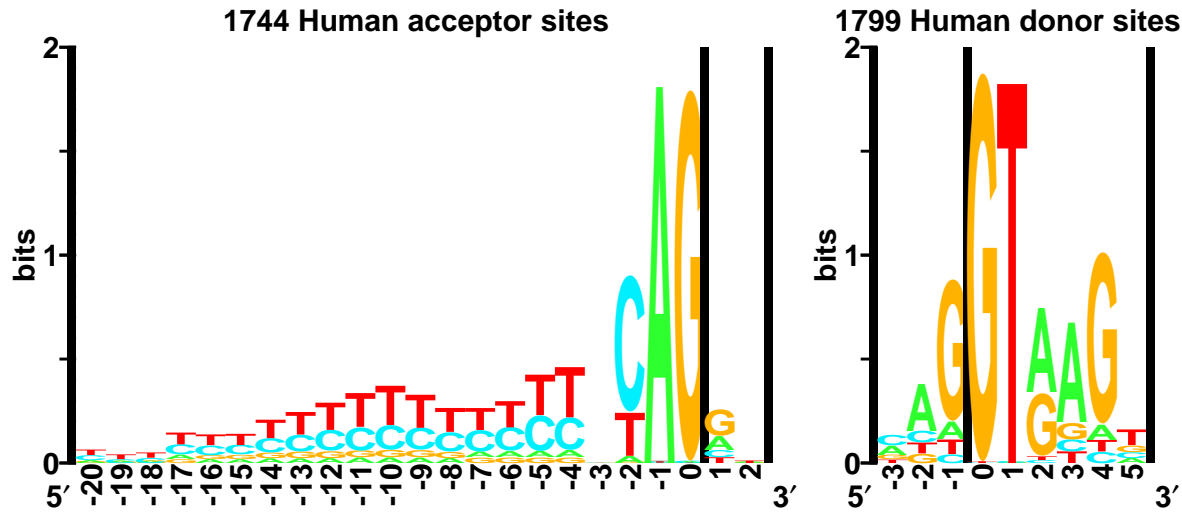
# Important Discovery



- Area under a sequence logo is the total information.
- How is that related to the binding energy?
- Information gained for energy dissipated

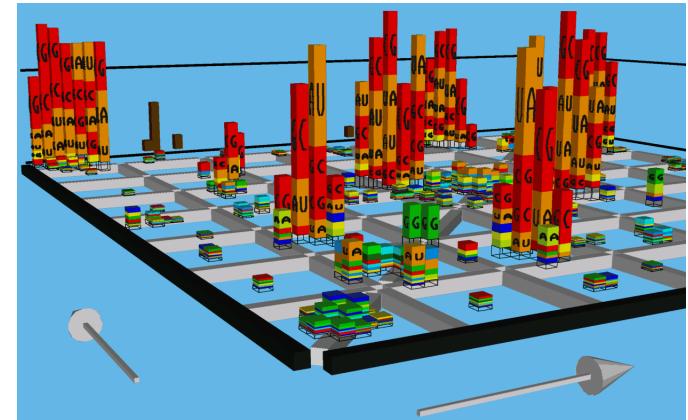
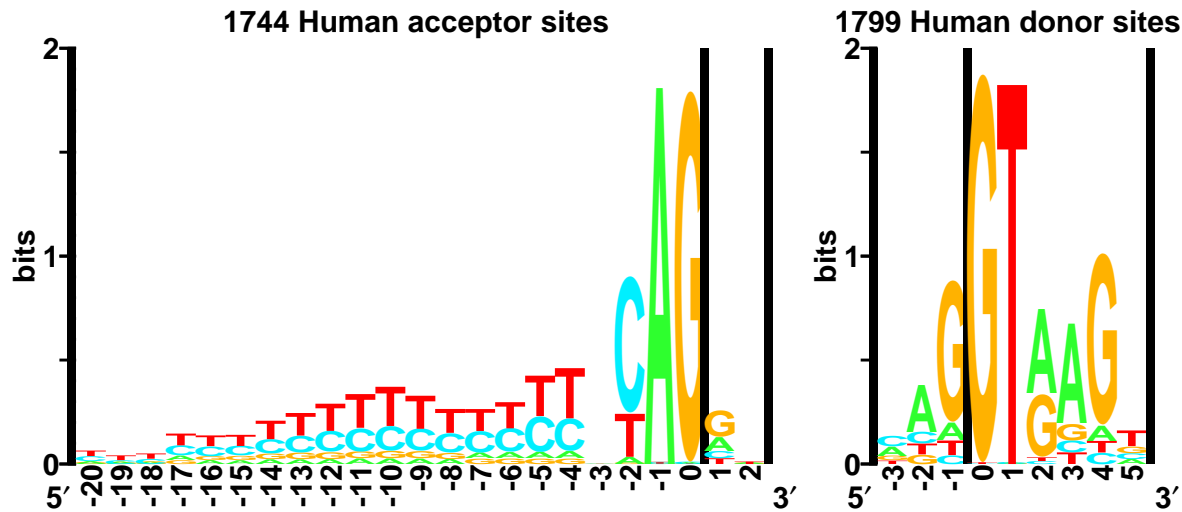


# Important Discovery



- Area under a sequence logo is the total information.
- How is that related to the binding energy?
- Information gained for energy dissipated
- Isothermal efficiency

# Important Discovery

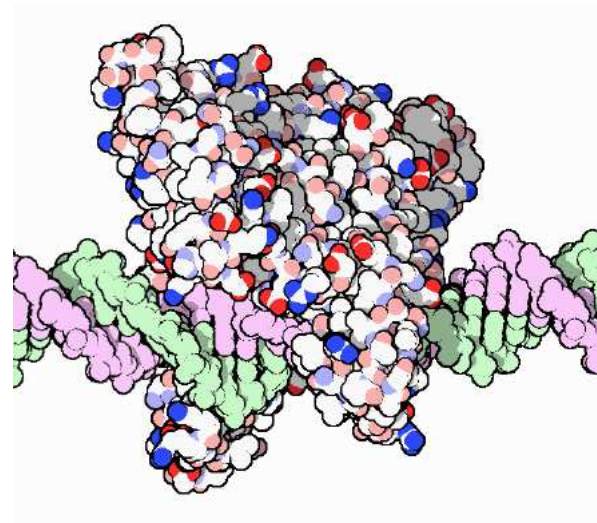


- Area under a sequence logo is the total information.
- How is that related to the binding energy?
- Information gained for energy dissipated
- Isothermal efficiency
- My most important discovery:

Molecules are often 70% efficient

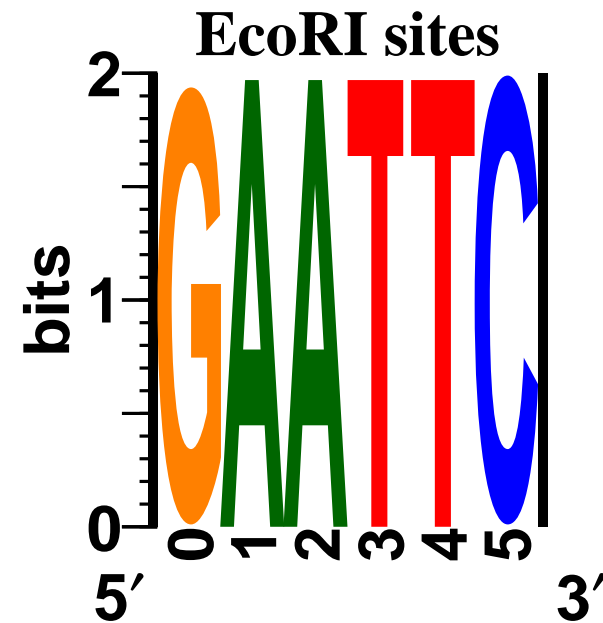
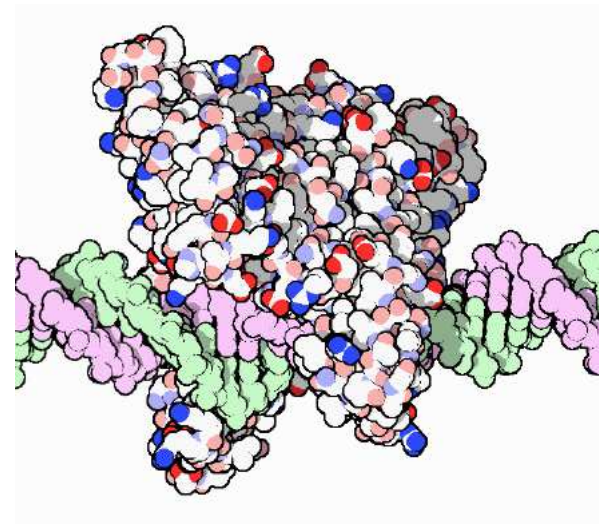
# Information of EcoRI DNA Binding

- EcoRI - restriction enzyme



# Information of EcoRI DNA Binding

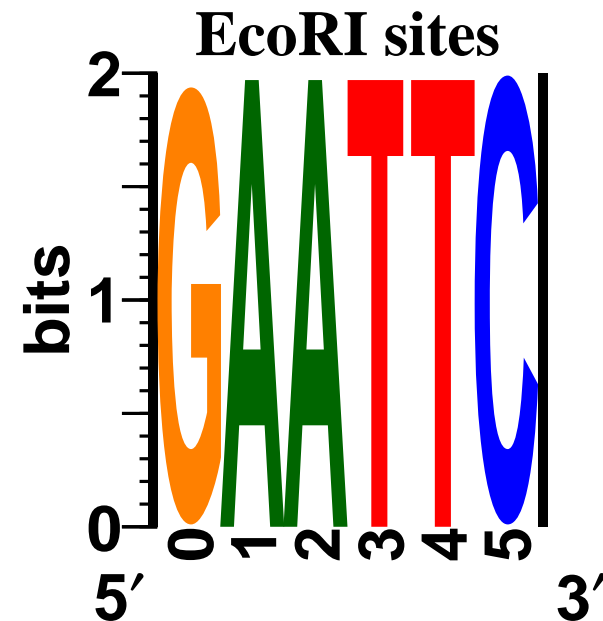
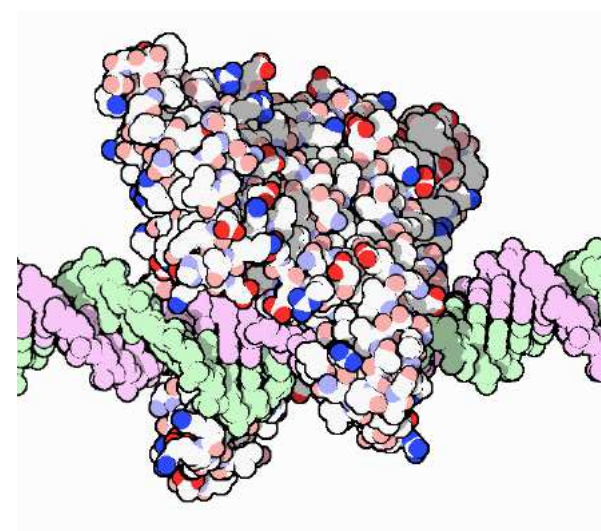
- EcoRI - restriction enzyme
- EcoRI binds DNA at 5' GAATTC 3'



# Information of EcoRI DNA Binding

- EcoRI - restriction enzyme
- EcoRI binds DNA at 5' GAATTC 3'
- information required:

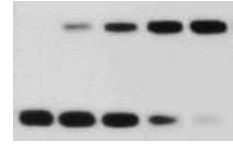
$$6 \text{ bases} \times 2 \text{ bits per base} = \boxed{12 \text{ bits}}$$



# Energy Dissipation by EcoRI

- Measured specific binding constant:

$$K_{spec} = 1.6 \times 10^5$$



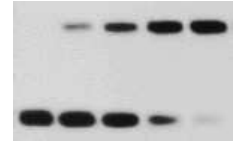
# Energy Dissipation by EcoRI

- Measured specific binding constant:

$$K_{spec} = 1.6 \times 10^5$$

- Average energy dissipated by one molecule as it binds:

$$\Delta G_{spec}^{\circ} = -k_B T \ln K_{spec} \quad (\text{joules per binding})$$



# Energy Dissipation by EcoRI

- Measured specific binding constant:

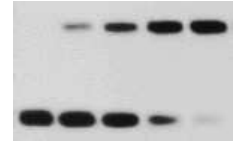
$$K_{spec} = 1.6 \times 10^5$$

- Average energy dissipated by one molecule as it binds:

$$\Delta G_{spec}^{\circ} = -k_B T \ln K_{spec} \quad (\text{joules per binding})$$

- The Second Law of Thermodynamics as a conversion factor:

$$\mathcal{E}_{min} = k_B T \ln 2 \quad (\text{joules per bit})$$

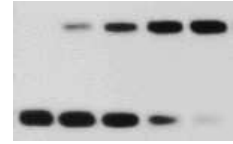




# Energy Dissipation by EcoRI

- Measured specific binding constant:

$$K_{spec} = 1.6 \times 10^5$$



- Average energy dissipated by one molecule as it binds:

$$\Delta G_{spec}^{\circ} = -k_B T \ln K_{spec} \quad (\text{joules per binding})$$

- The Second Law of Thermodynamics as a conversion factor:

$$\mathcal{E}_{min} = k_B T \ln 2 \quad (\text{joules per bit})$$

- Number of bits that could have been selected:

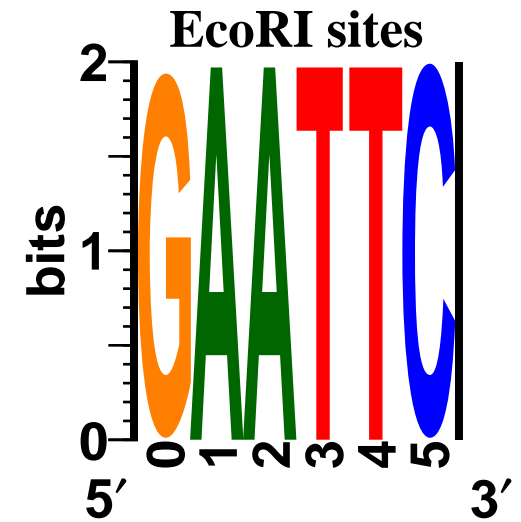
$$\begin{aligned} R_{energy} &= -\Delta G^{\circ} / \mathcal{E}_{min} \\ &= k_B T \ln K_{spec} / k_B T \ln 2 \\ &= \log_2 K_{spec} \quad \Leftarrow \text{SO SIMPLE!} \\ &= \boxed{17.3 \text{ bits per binding}} \end{aligned}$$

# Information/Energy = Efficiency of EcoRI

EcoRI could have made 17.3 binary choices

# Information/Energy = Efficiency of EcoRI

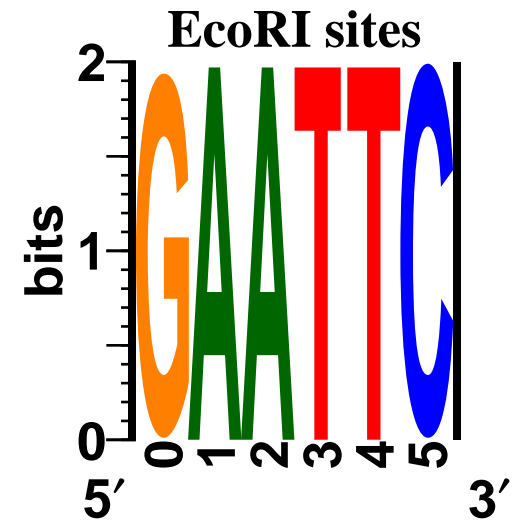
EcoRI could have made 17.3 binary choices  
... but it only made 12 choices.



# Information/Energy = Efficiency of EcoRI

EcoRI could have made 17.3 binary choices  
...but it only made 12 choices.

Efficiency is  
'WORK' DONE / ENERGY DISSIPATED

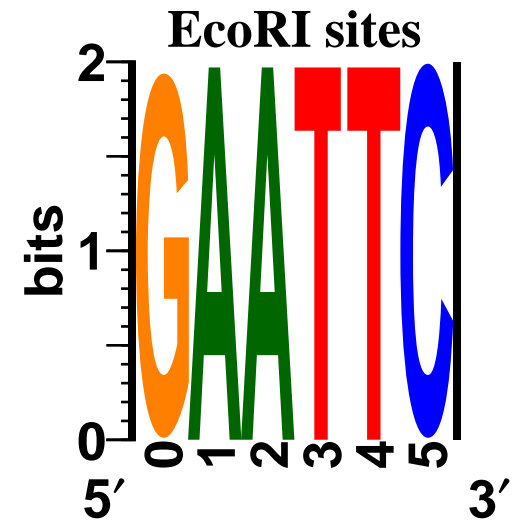


# Information/Energy = Efficiency of EcoRI

EcoRI could have made 17.3 binary choices  
...but it only made 12 choices.

Efficiency is  
'WORK' DONE / ENERGY DISSIPATED

$$\frac{12 \text{ bits per binding}}{17.3 \text{ bits per binding}} = 0.7$$



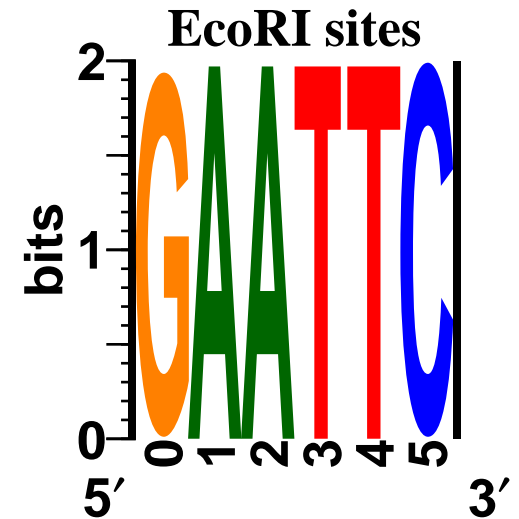
# Information/Energy = Efficiency of EcoRI = 70%

EcoRI could have made 17.3 binary choices  
...but it only made 12 choices.

Efficiency is  
'WORK' DONE / ENERGY DISSIPATED

$$\frac{12 \text{ bits per binding}}{17.3 \text{ bits per binding}} = 0.7$$

**The efficiency is 70%.**



# Information/Energy = Efficiency of EcoRI = 70%

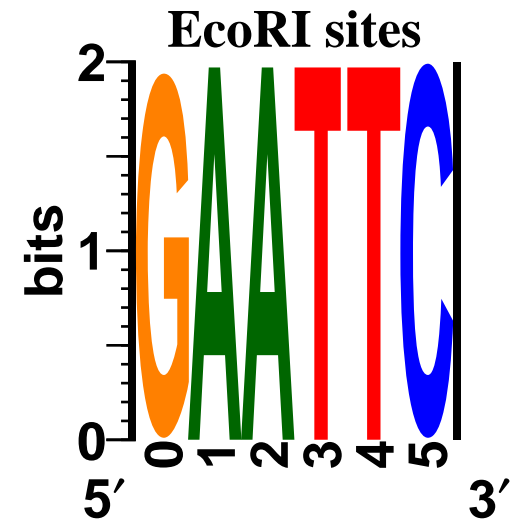
EcoRI could have made 17.3 binary choices  
...but it only made 12 choices.

Efficiency is  
'WORK' DONE / ENERGY DISSIPATED

$$\frac{12 \text{ bits per binding}}{17.3 \text{ bits per binding}} = 0.7$$

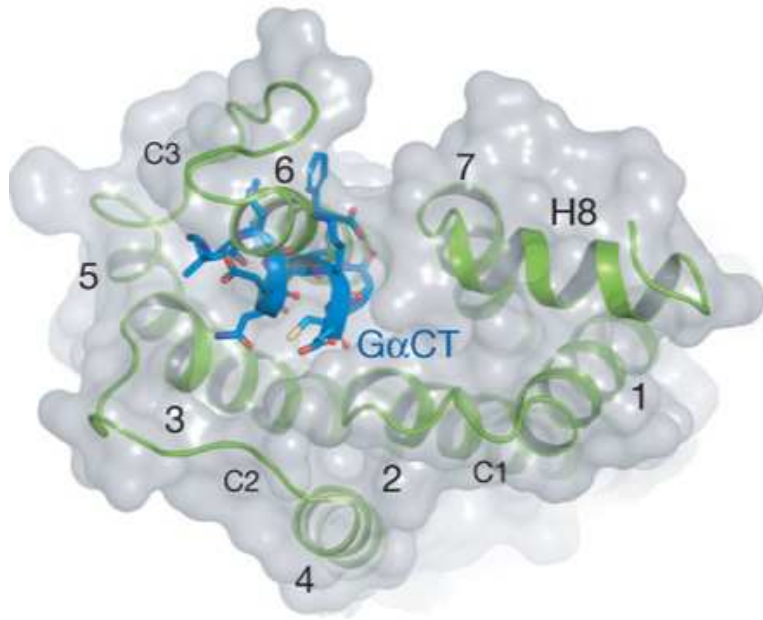
**The efficiency is 70%.**

**18 out of 19 DNA binding proteins give ~70% efficiency.**



# Rhodopsin Shape Change

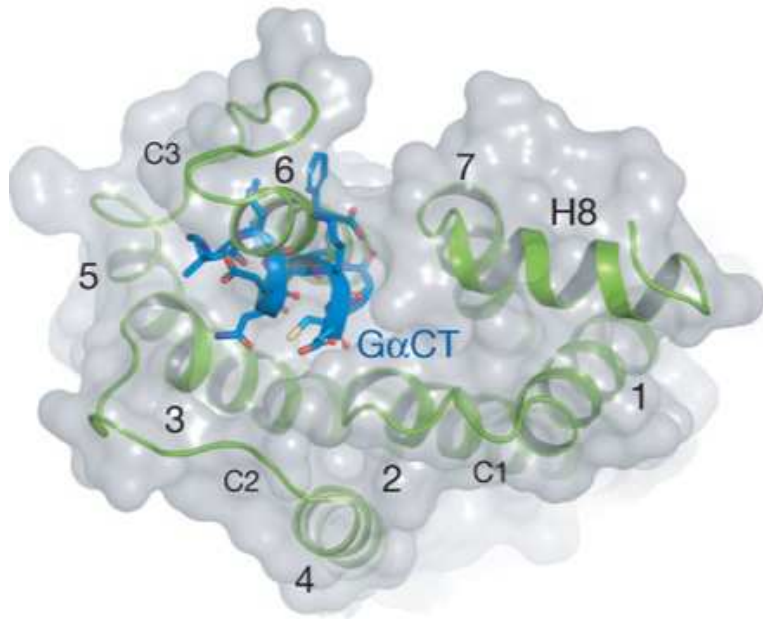
Dark State





# Rhodopsin Shape Change

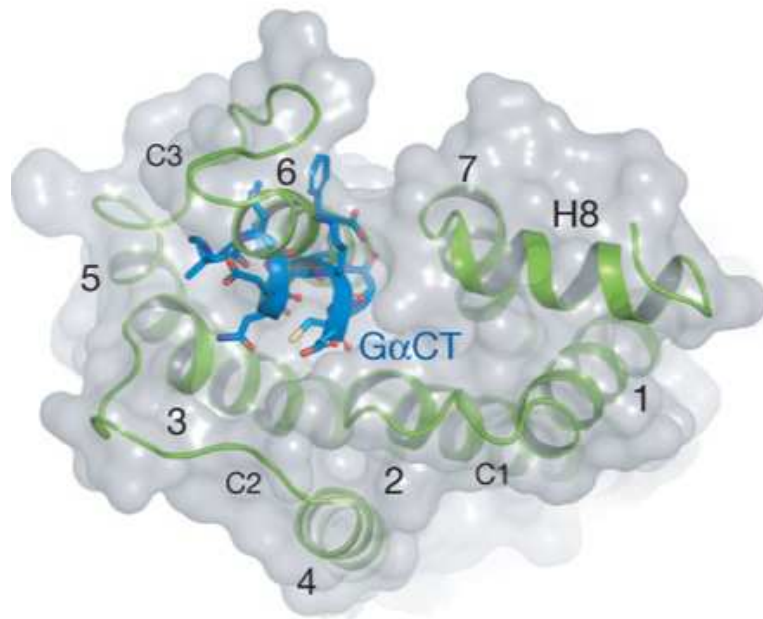
Dark State



$h\nu$

# Rhodopsin Shape Change

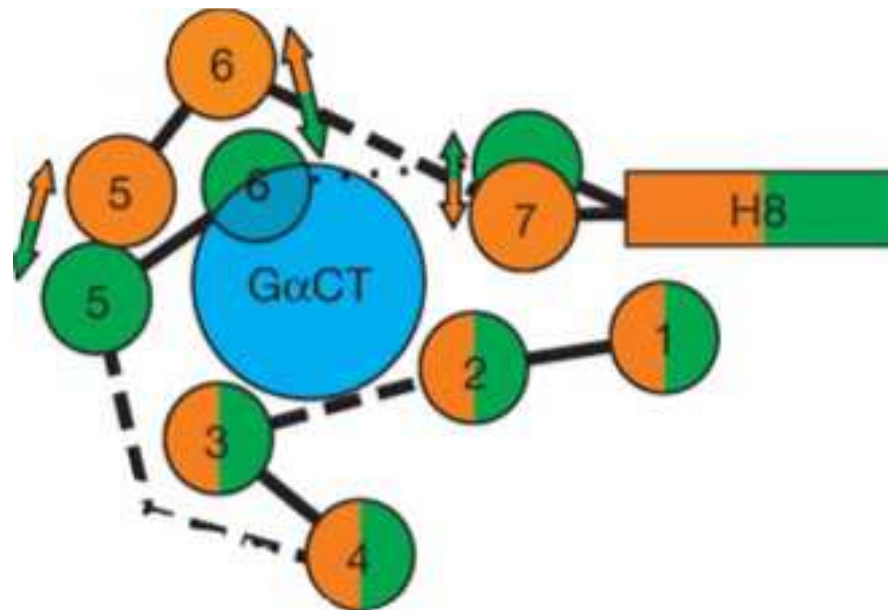
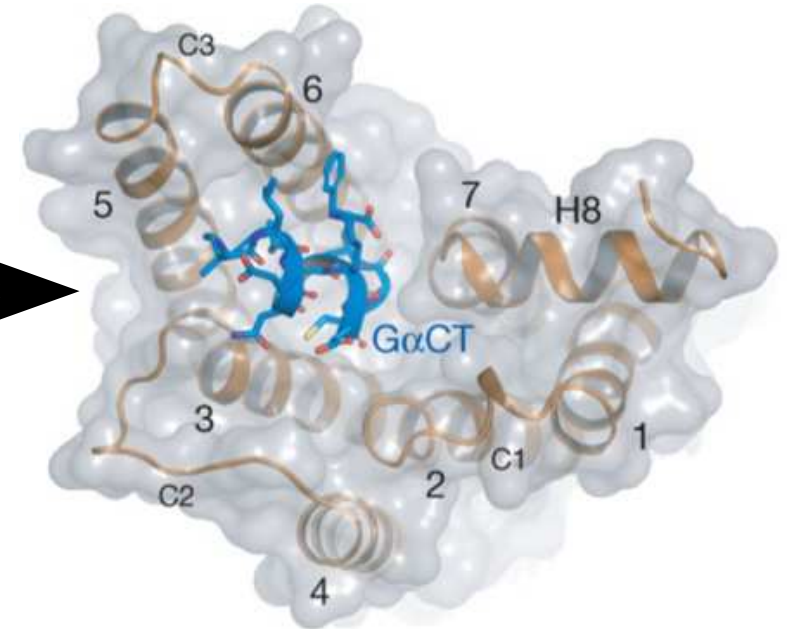
Dark State



$h\nu$

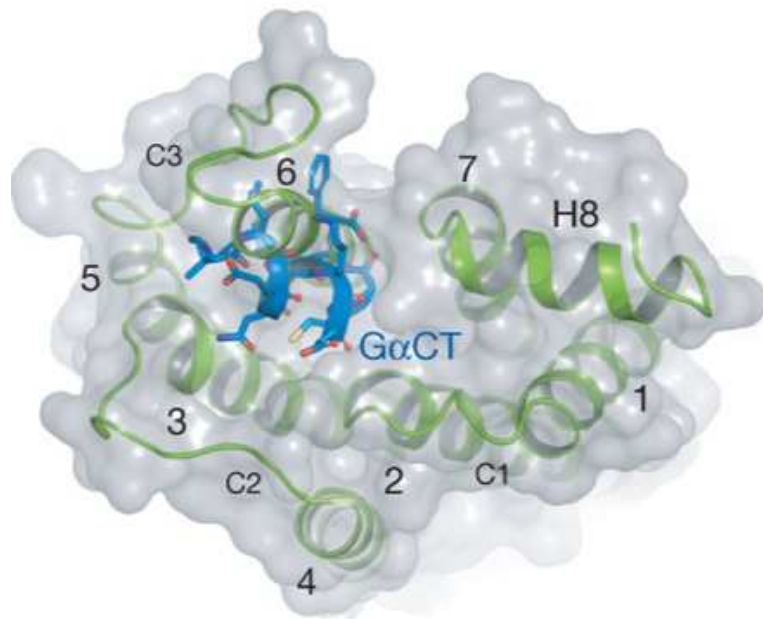


After Photon - Light State

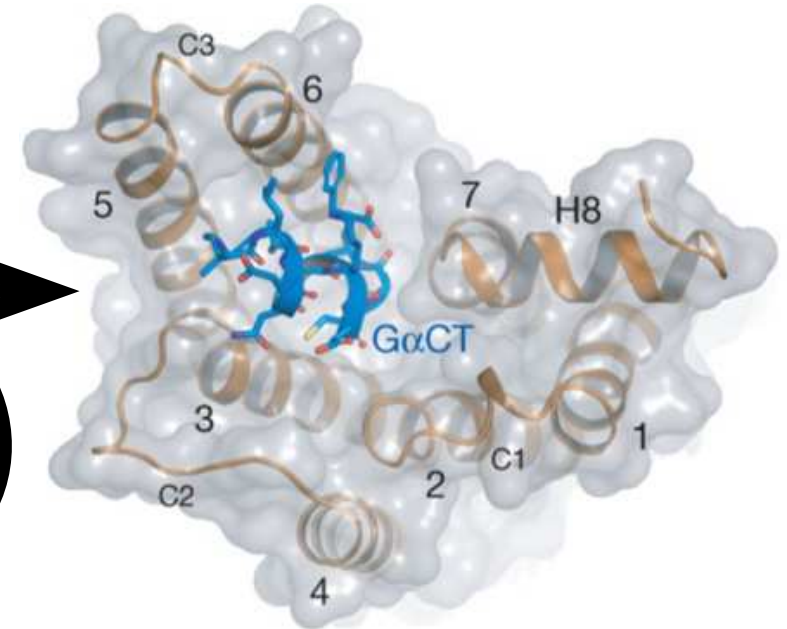


# Rhodopsin Shape Change

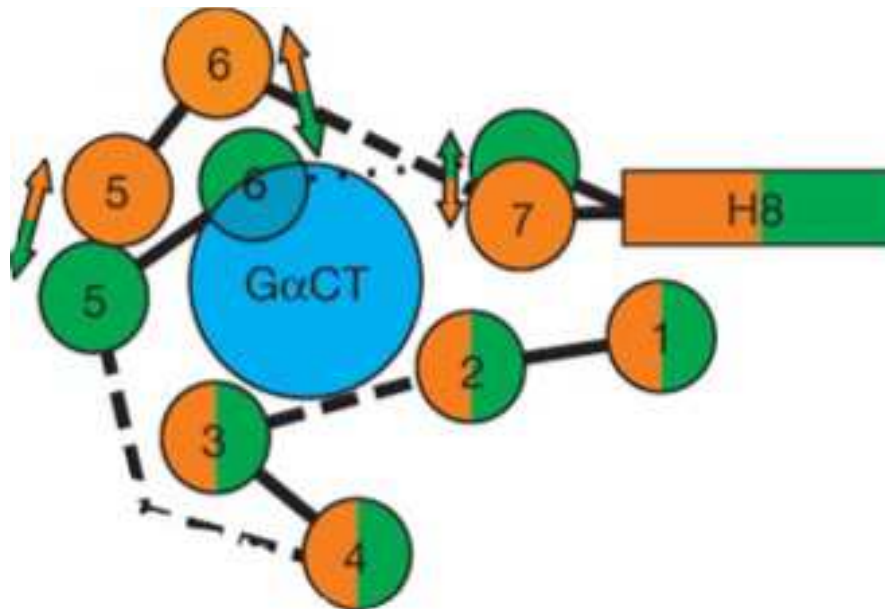
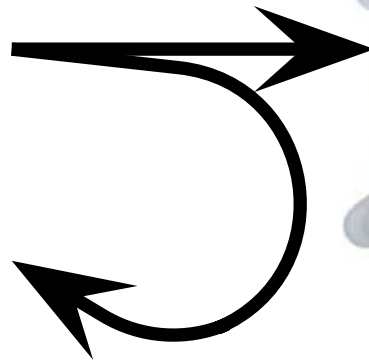
Dark State



After Photon - Light State



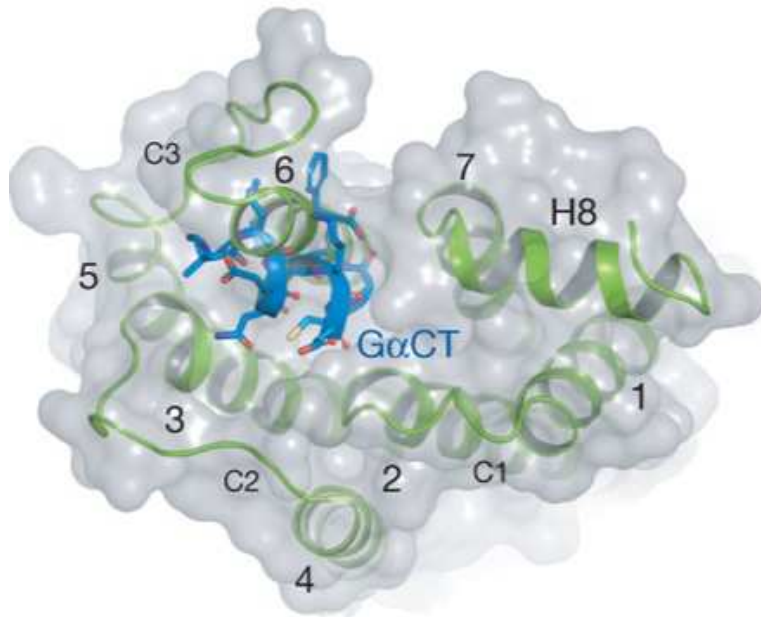
$h\nu$



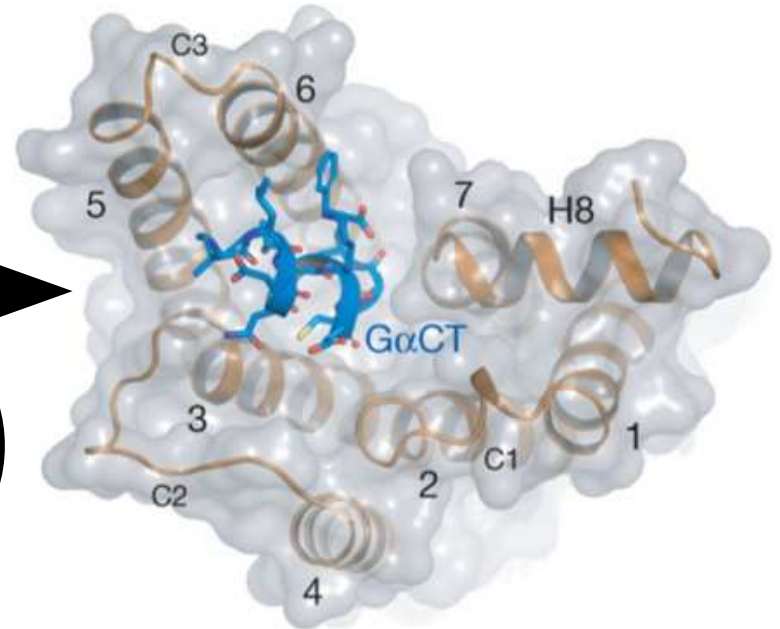


# Rhodopsin Shape Change

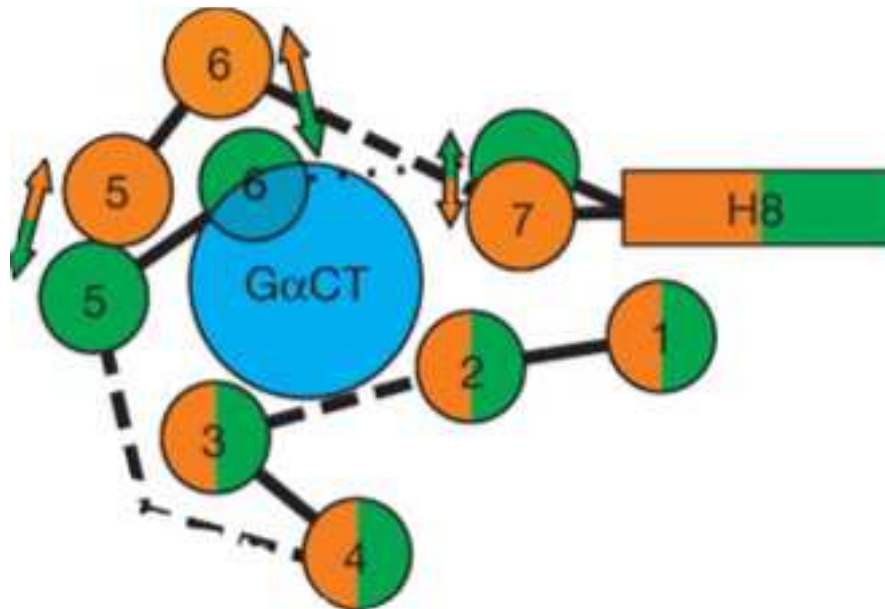
Dark State



After Photon - Light State



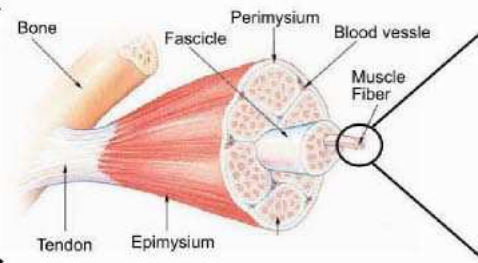
$h\nu$   
70%  
30%



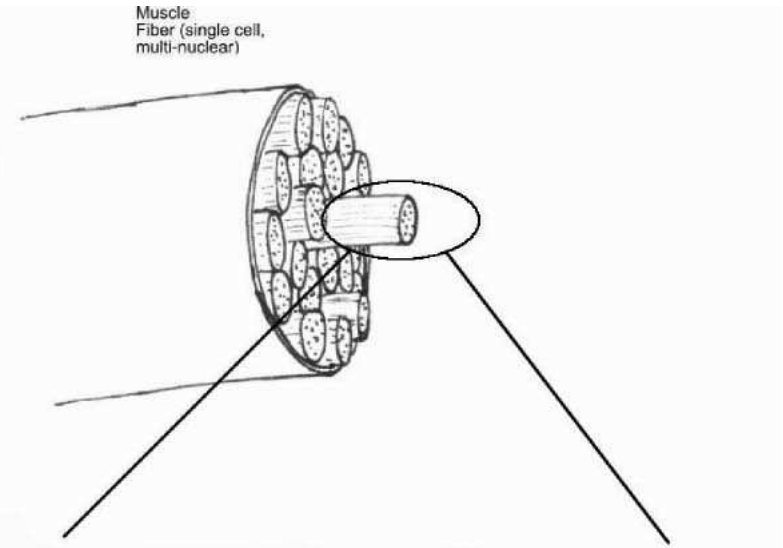
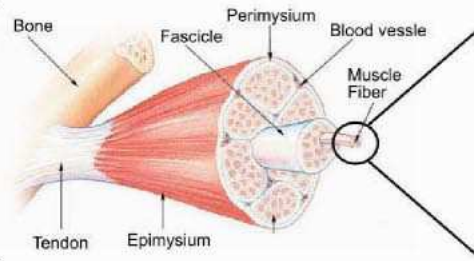
# Muscle Structure



# Muscle Structure

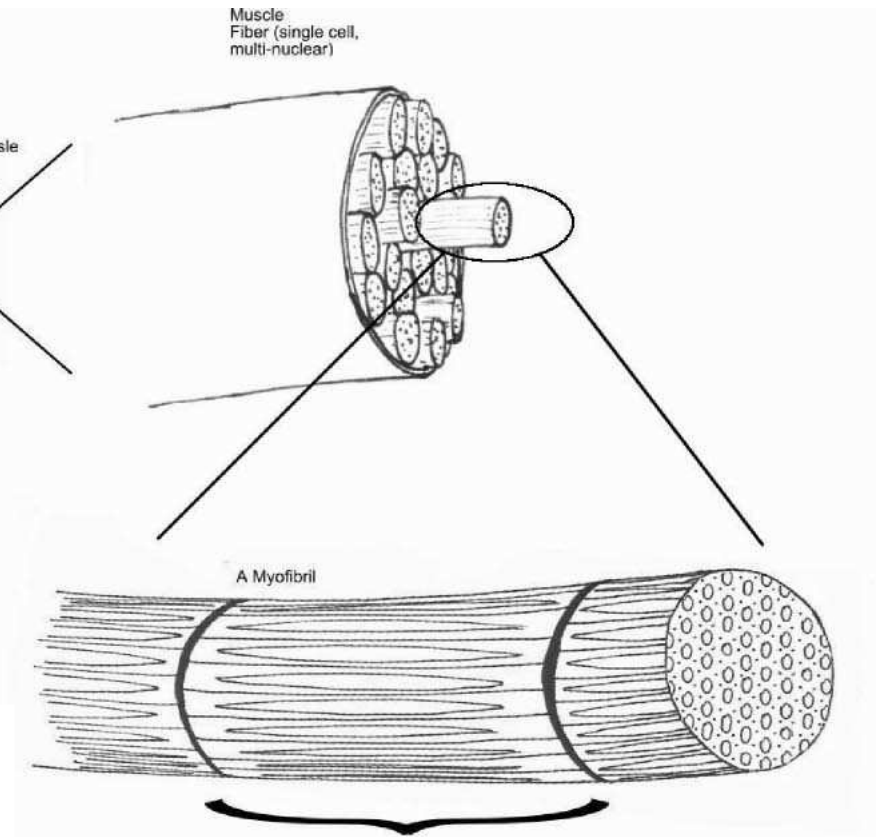
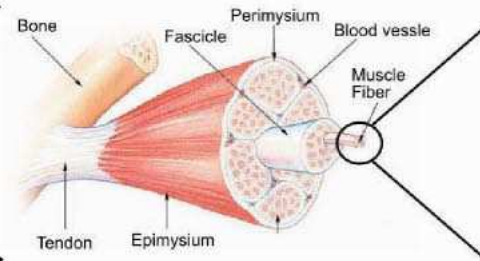


# Muscle Structure

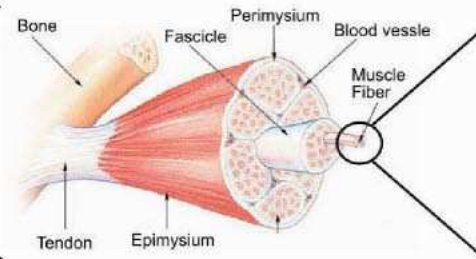




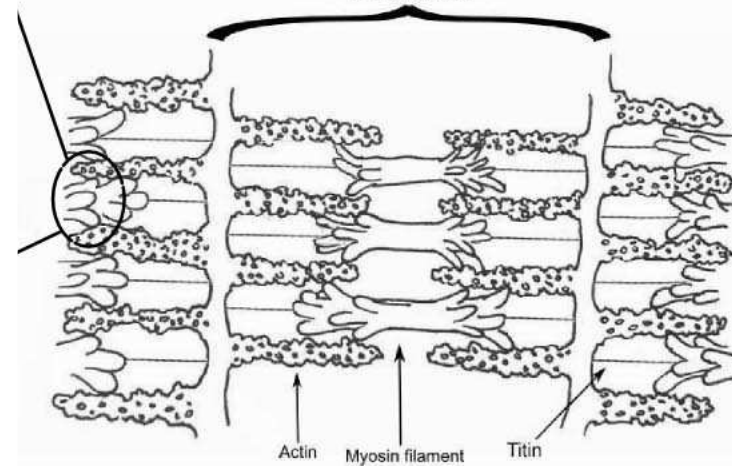
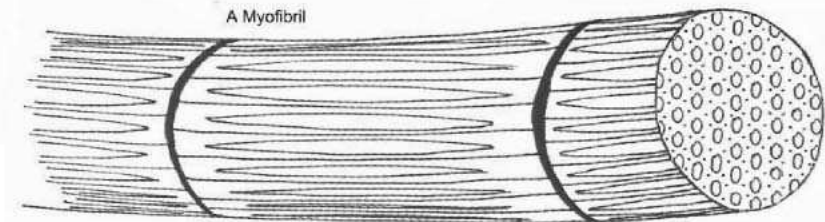
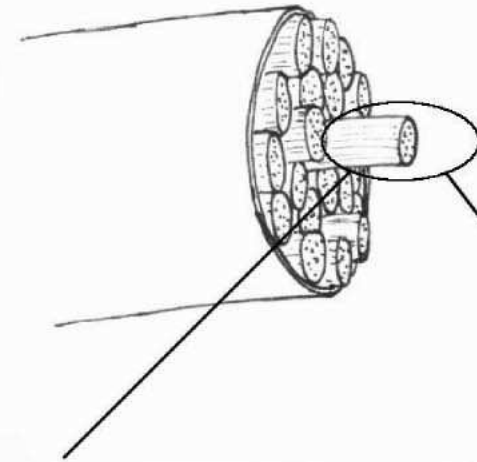
# Muscle Structure



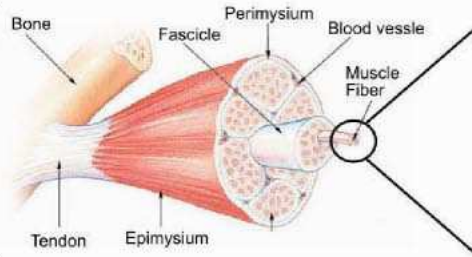
# Muscle Structure



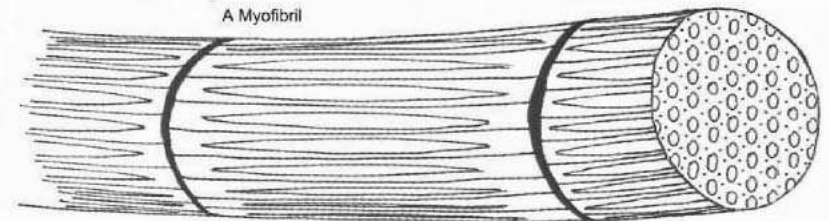
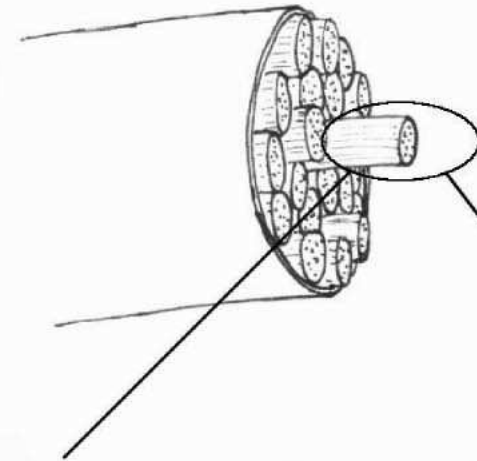
Muscle Fiber (single cell, multi-nuclear)



# Muscle Structure

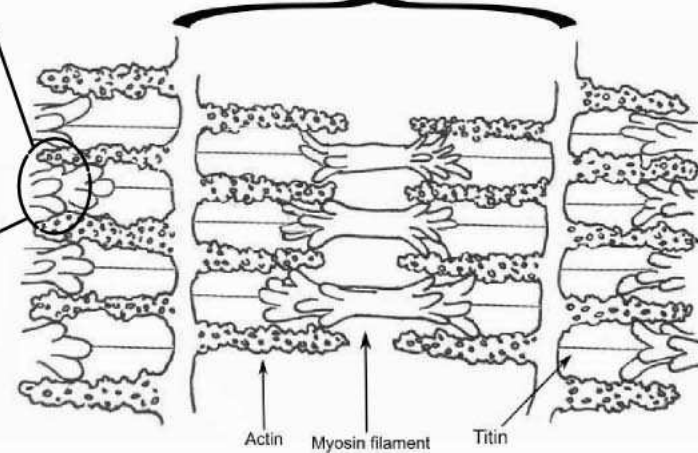
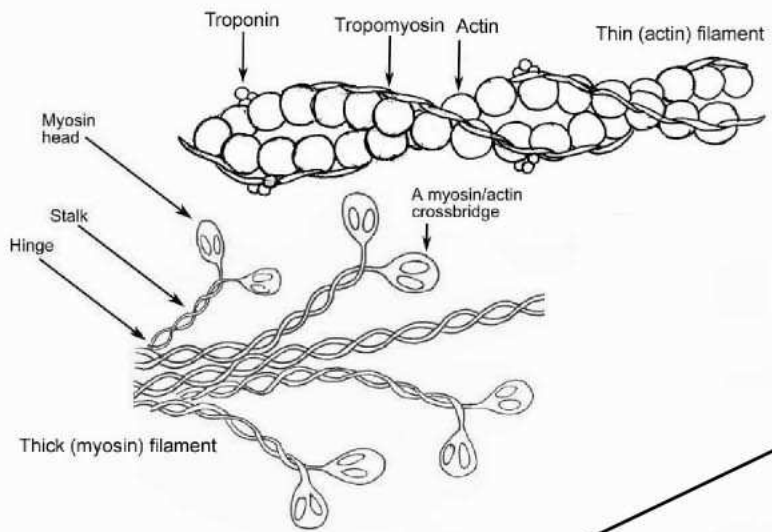


Muscle Fiber (single cell, multi-nuclear)

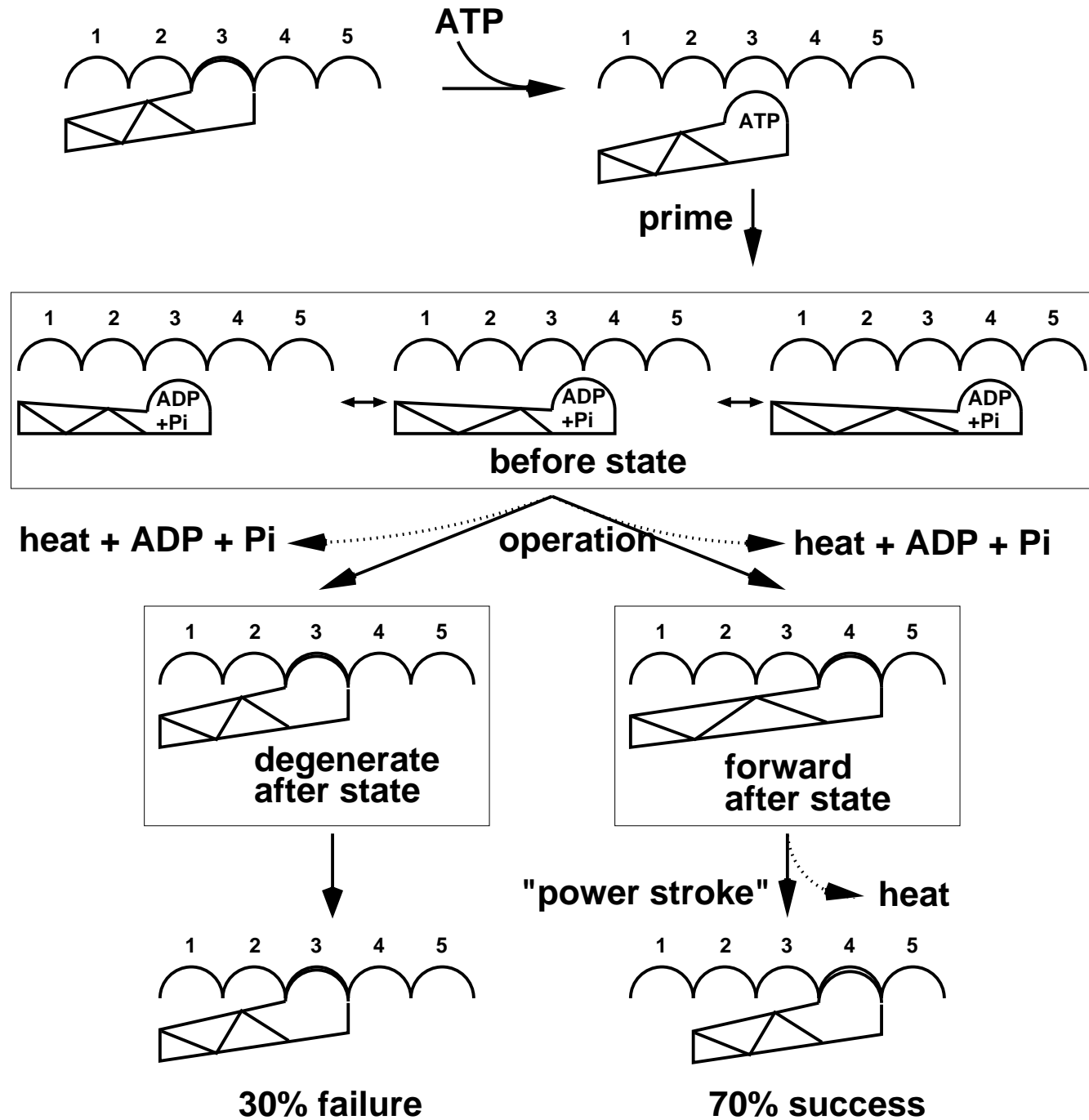


A Myofibril

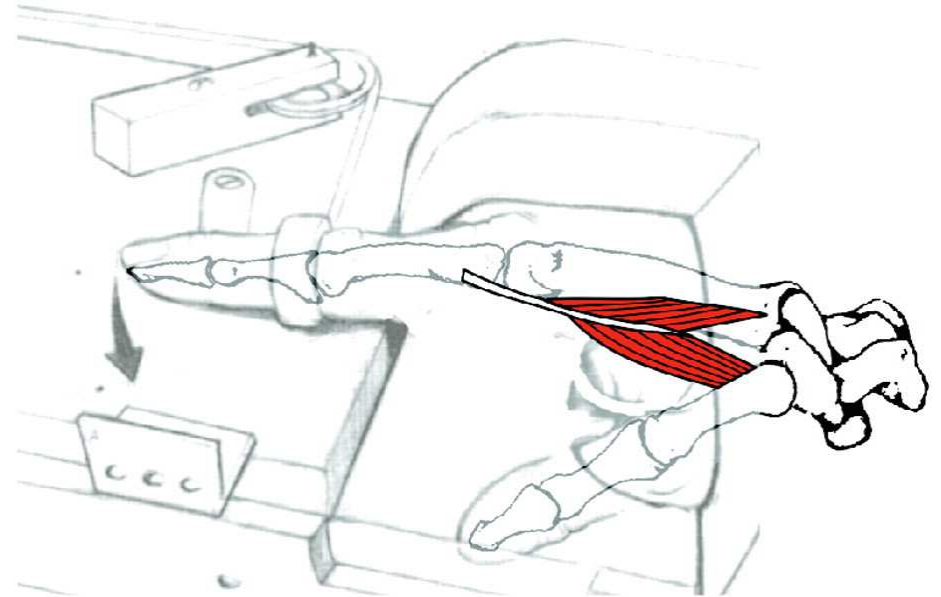
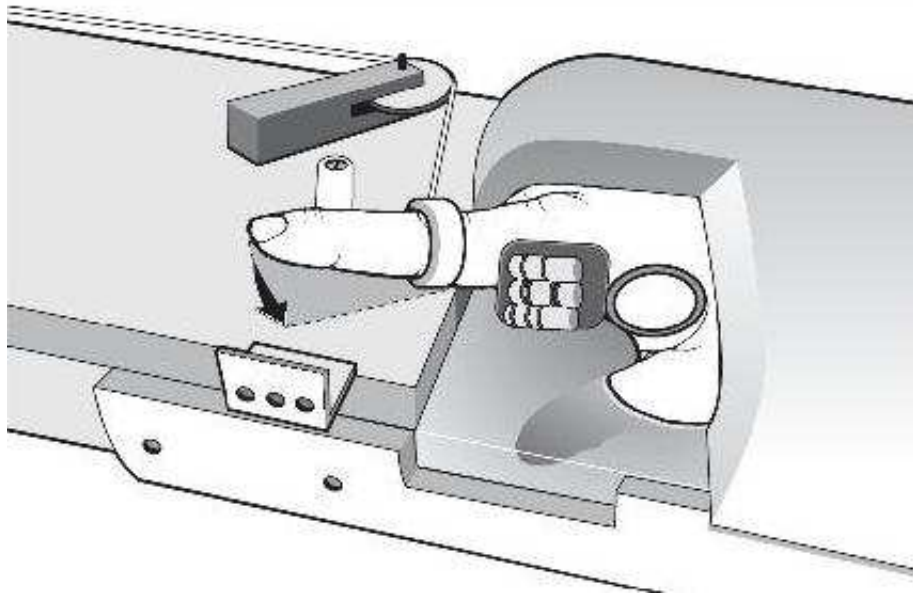
One sarcomere



# Tom's Model of Muscle Mechanism

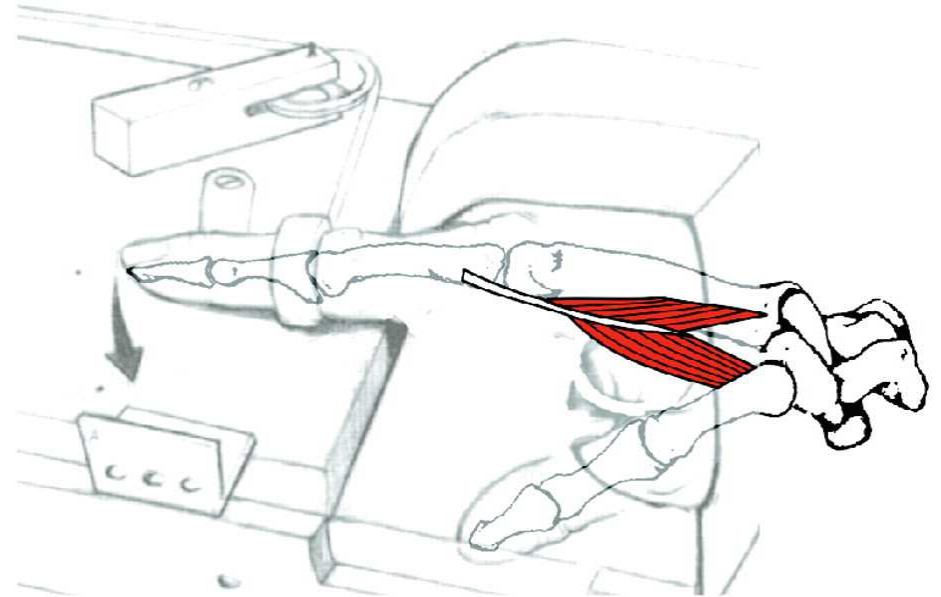
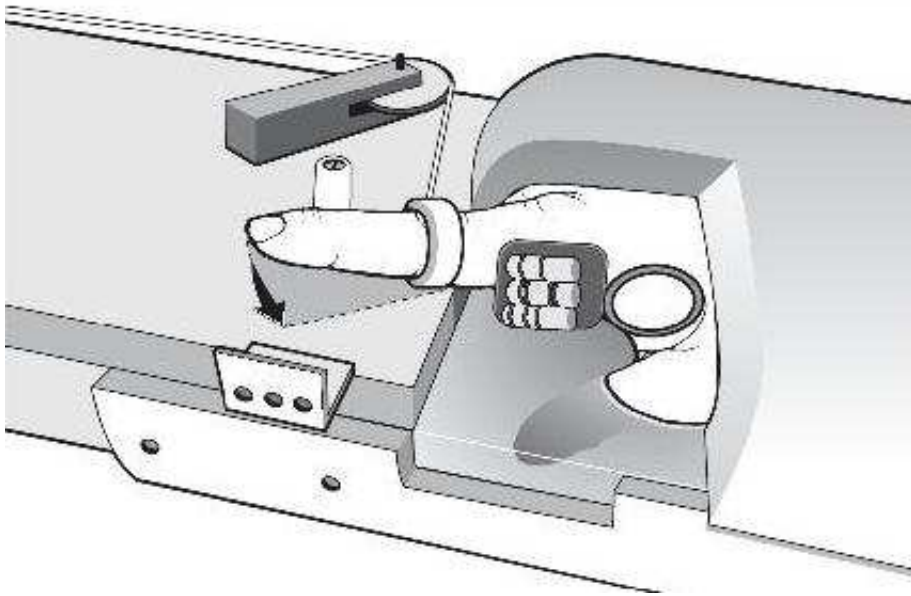


# Efficiency of Muscle



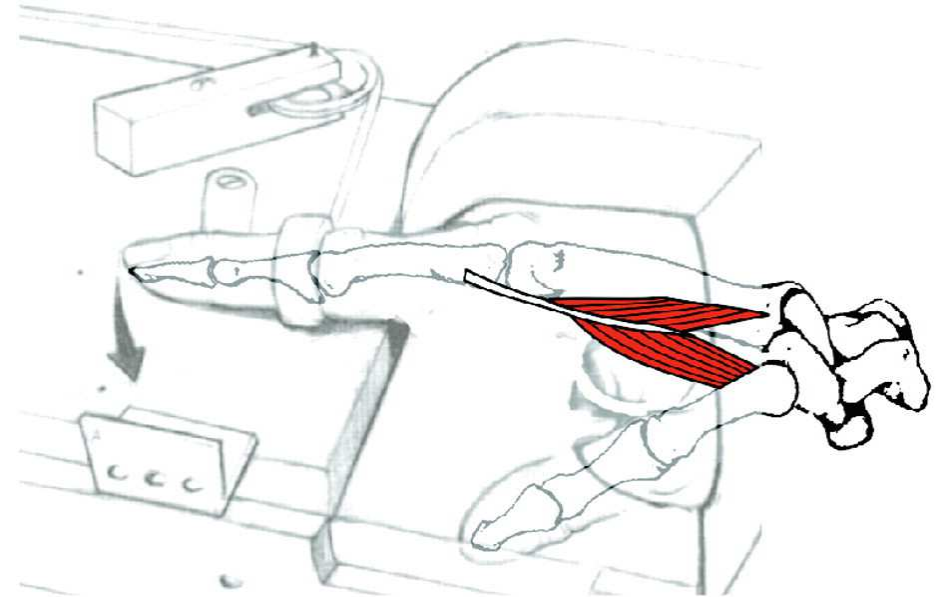
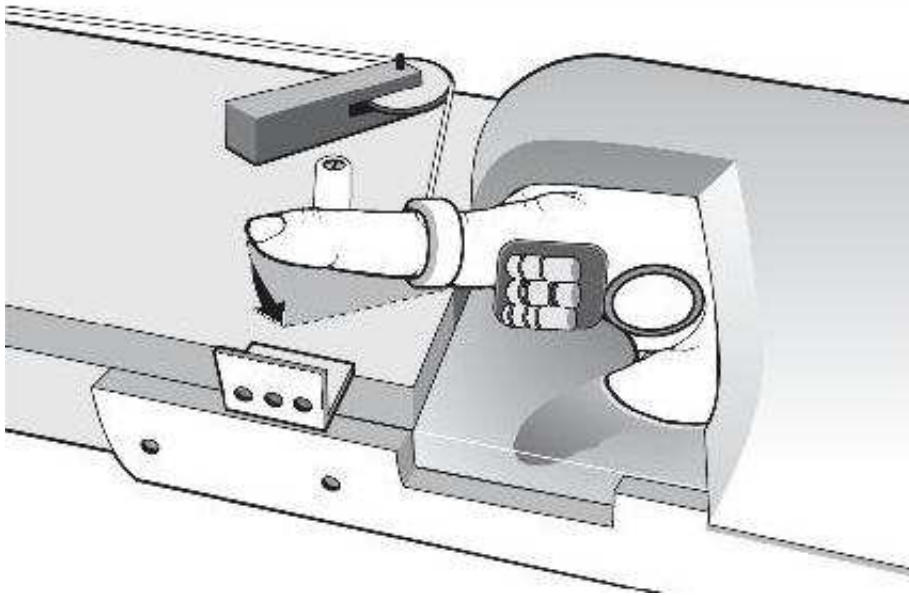
- Experiments by Kushmerick's lab since (at least) 1969

# Efficiency of Muscle



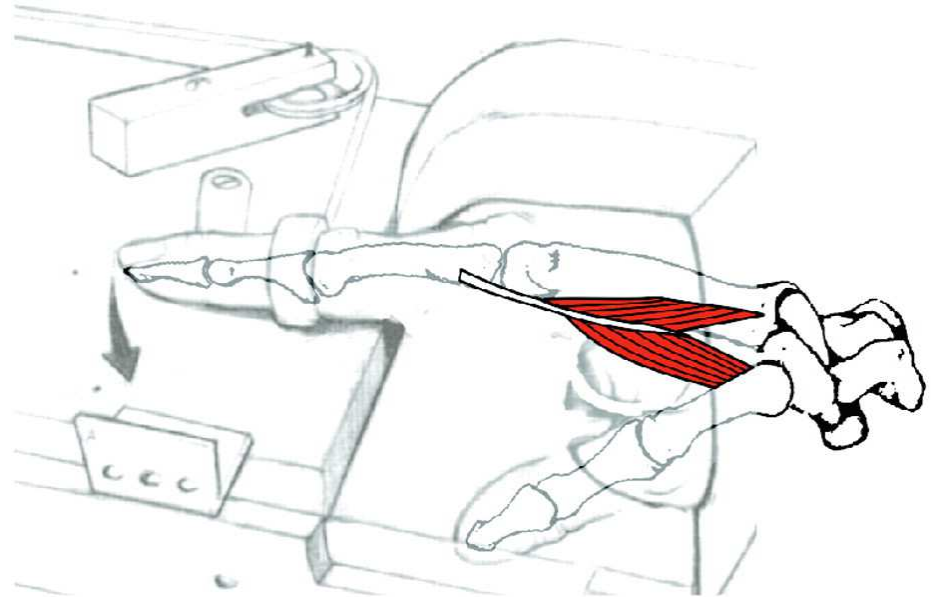
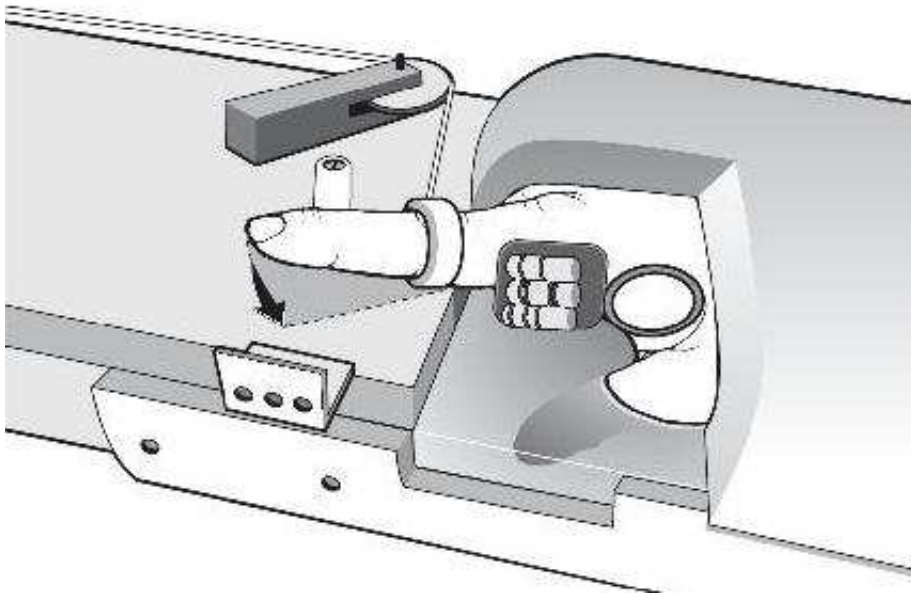
- Experiments by Kushmerick's lab since (at least) 1969
- new work: 2008, 2011

# Efficiency of Muscle



- Experiments by Kushmerick's lab since (at least) 1969
- new work: 2008, 2011
- Weight lifting gives work done

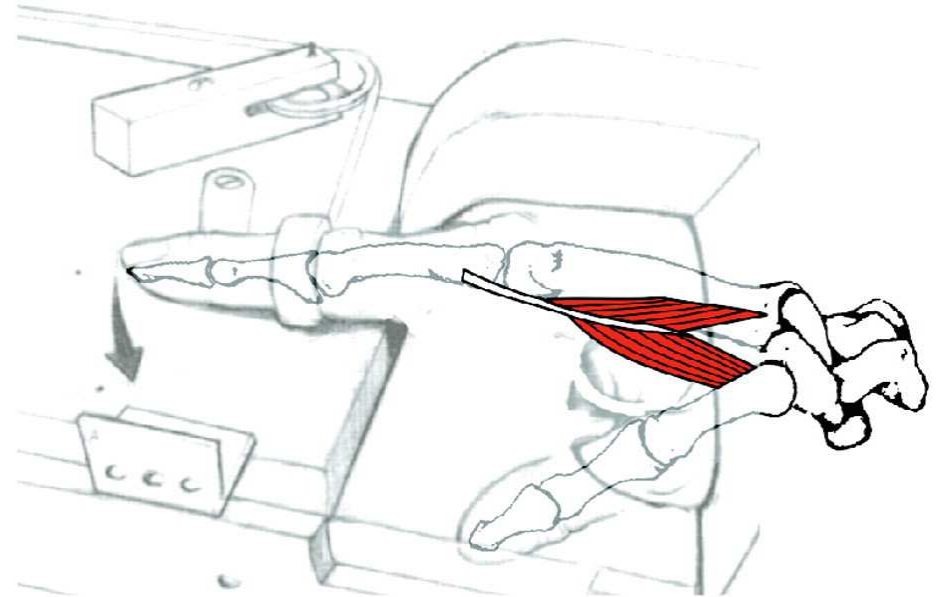
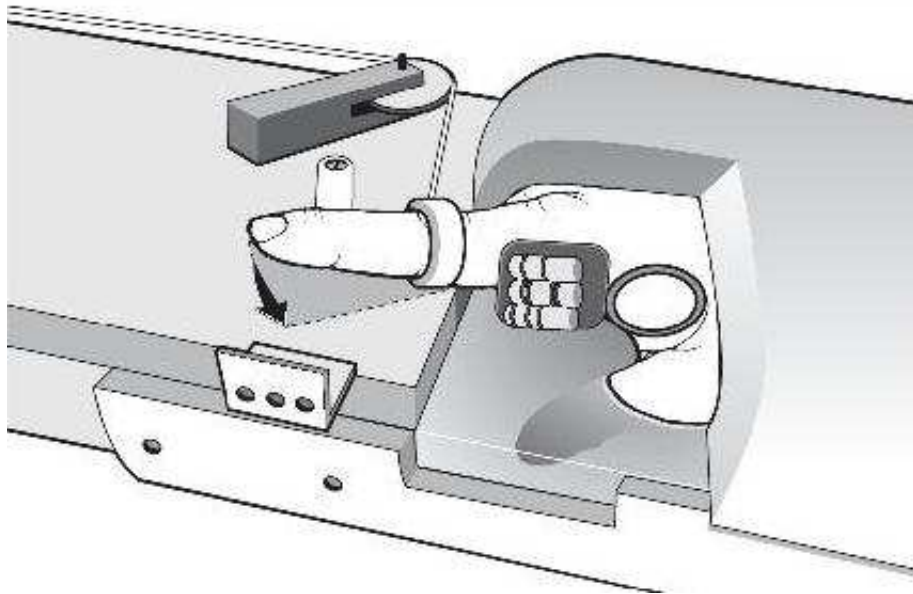
# Efficiency of Muscle



- Experiments by Kushmerick's lab since (at least) 1969
- new work: 2008, 2011
- Weight lifting gives work done
- NMR coil gives ATP = energy used



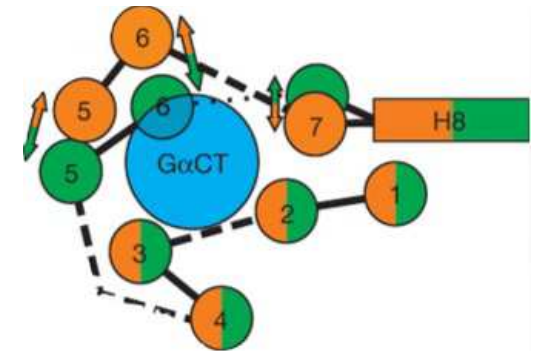
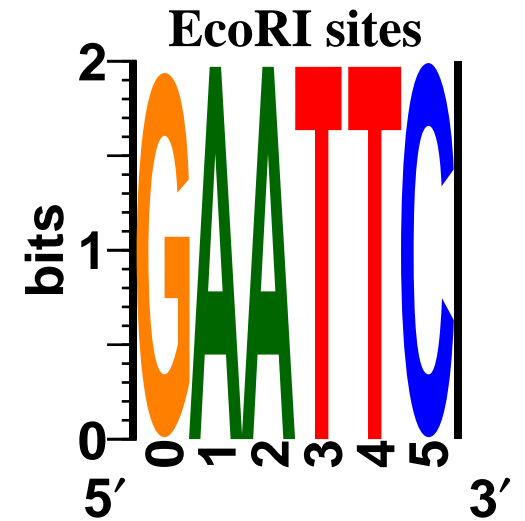
# Efficiency of Muscle



- Experiments by Kushmerick's lab since (at least) 1969
- new work: 2008, 2011
- Weight lifting gives work done
- NMR coil gives ATP = energy used
- **Efficiency:  $0.68 \pm 0.09$**

# Why are molecular machines 70% efficient?

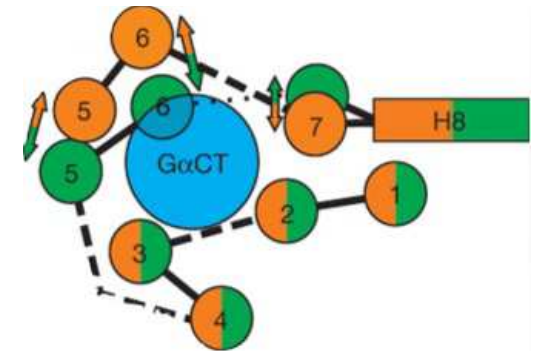
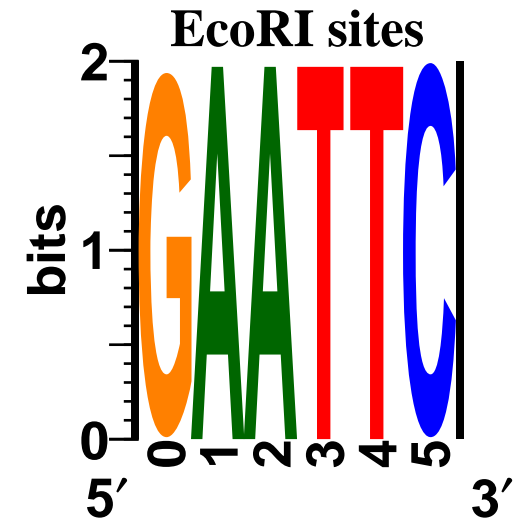
70% efficiency appears widely in biology:



# Why are molecular machines 70% efficient?

70% efficiency appears widely in biology:

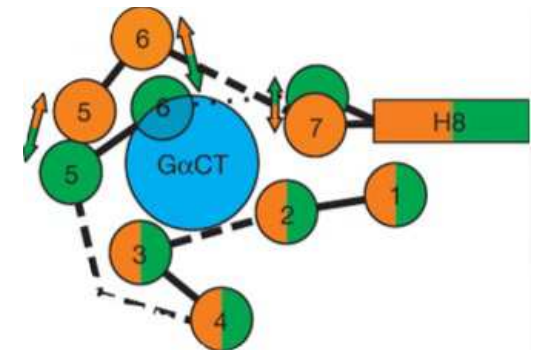
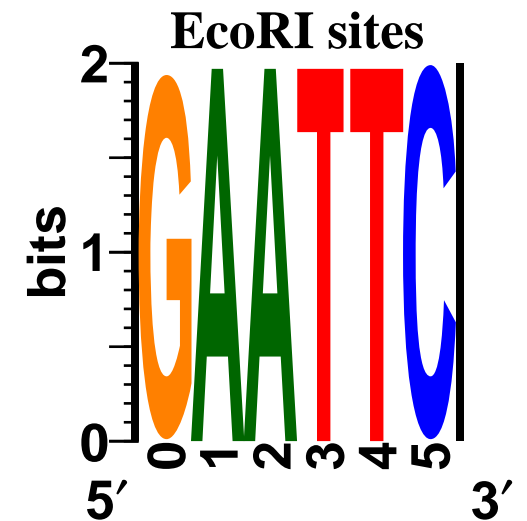
- DNA - protein binding



# Why are molecular machines 70% efficient?

70% efficiency appears widely in biology:

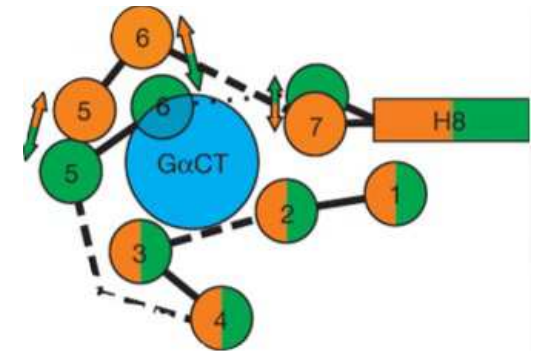
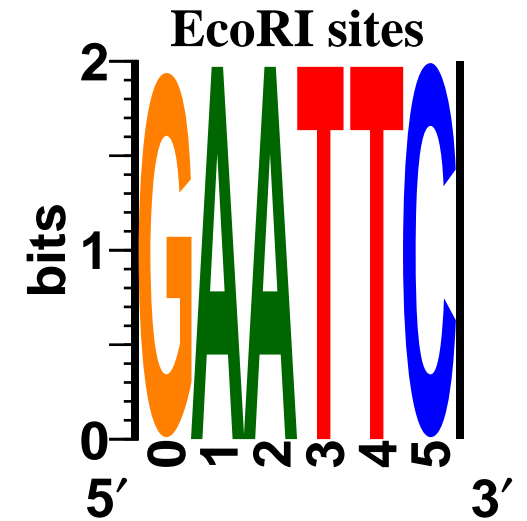
- DNA - protein binding
- rhodopsin



# Why are molecular machines 70% efficient?

70% efficiency appears widely in biology:

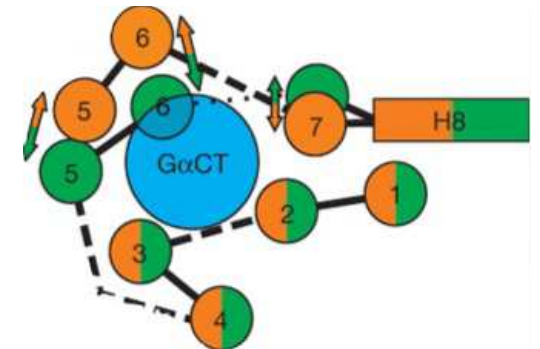
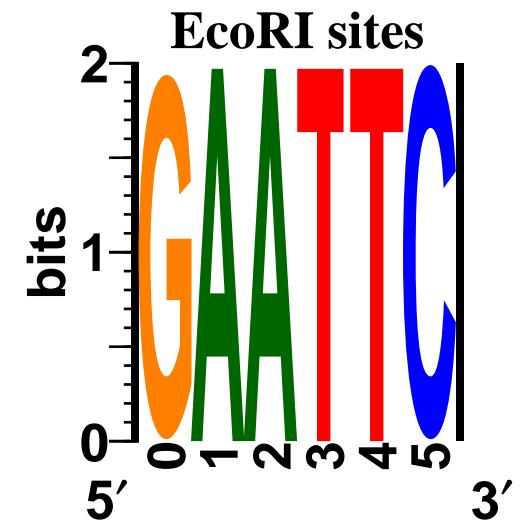
- DNA - protein binding
- rhodopsin
- muscle



# Why are molecular machines 70% efficient?

70% efficiency appears widely in biology:

- DNA - protein binding
- rhodopsin
- muscle
- other systems

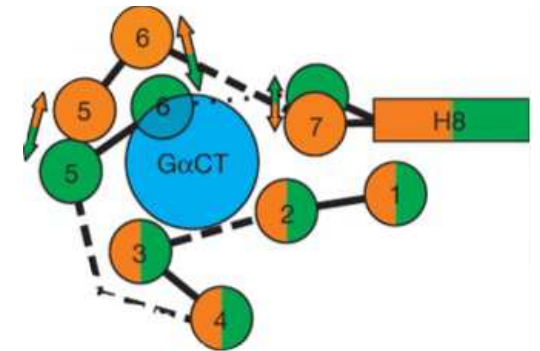
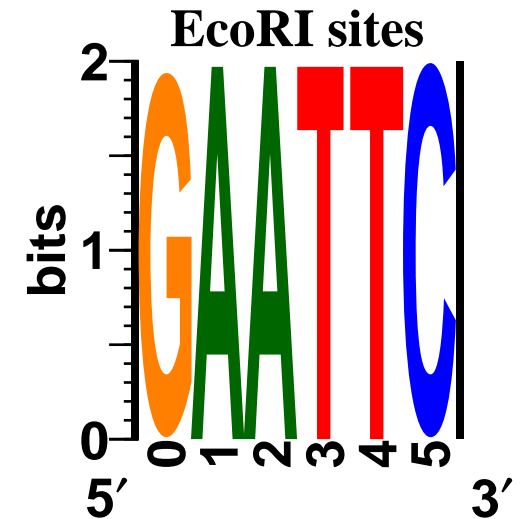


# Why are molecular machines 70% efficient?

70% efficiency appears widely in biology:

- DNA - protein binding
- rhodopsin
- muscle
- other systems

Why 70% efficiency?



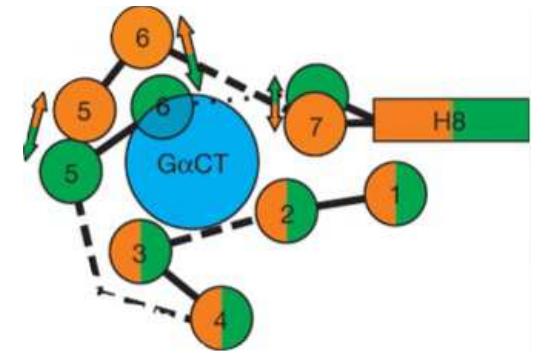
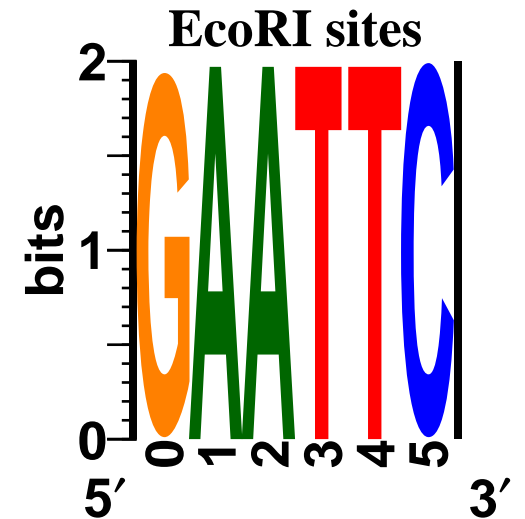
# Why are molecular machines 70% efficient?

70% efficiency appears widely in biology:

- DNA - protein binding
- rhodopsin
- muscle
- other systems

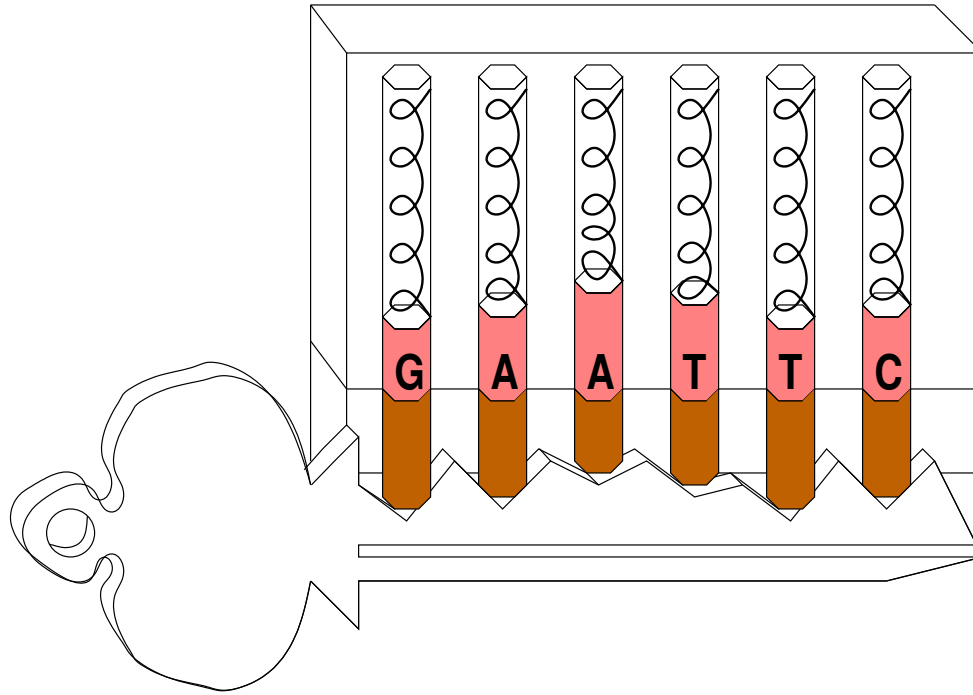
Why 70% efficiency?

Information theory explanation





# Lock and Key

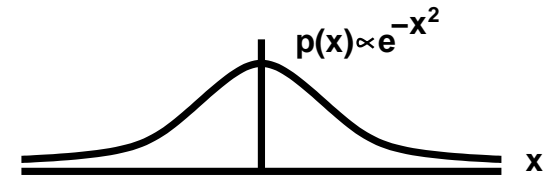


**Like a key in a lock  
which has many independent pins,  
it takes many numbers  
to describe the vibrational state  
of a molecular machine**

# Gaussians

- Pin motion  $x$  has a Gaussian distribution:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

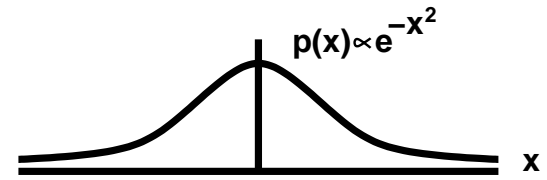


$\mu$  = mean,  $\sigma$  = standard deviation

# Gaussians

- Pin motion  $x$  has a Gaussian distribution:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



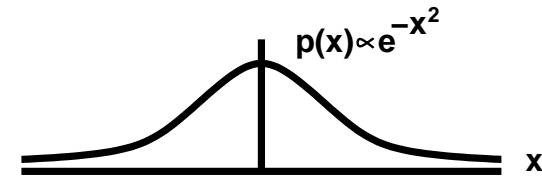
$\mu$  = mean,  $\sigma$  = standard deviation

- Gaussian distributions are generated by the sum of many small random variables

# Gaussians

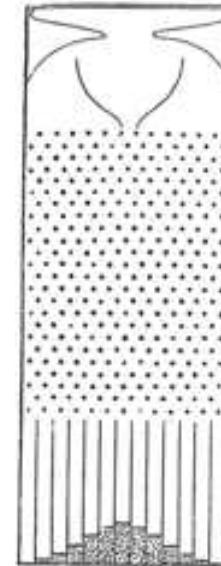
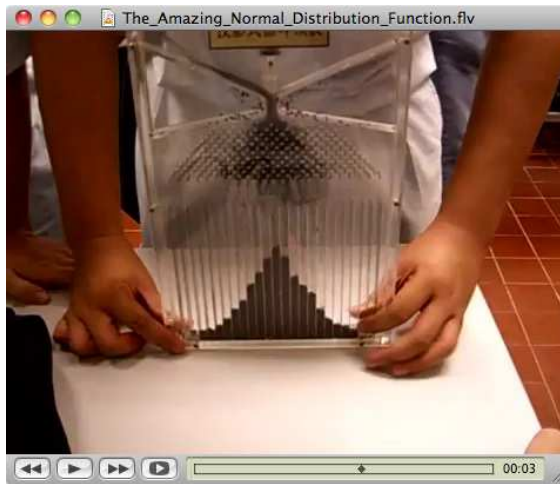
- Pin motion  $x$  has a Gaussian distribution:

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



$\mu$  = mean,  $\sigma$  = standard deviation

- Gaussian distributions are generated by the sum of many small random variables
- Drunkard's walk: Galton's quincunx device!



# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \quad (2)$$

# Two Gaussians

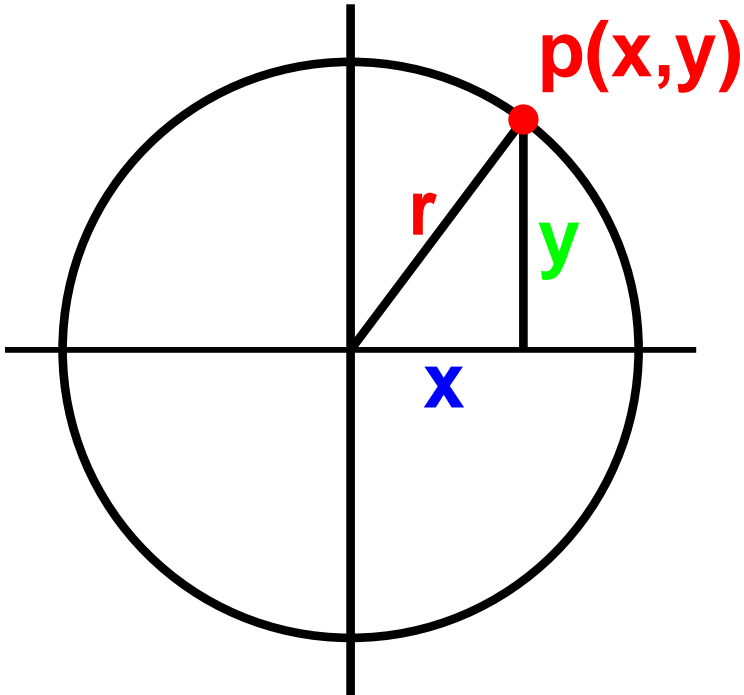
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$

# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$

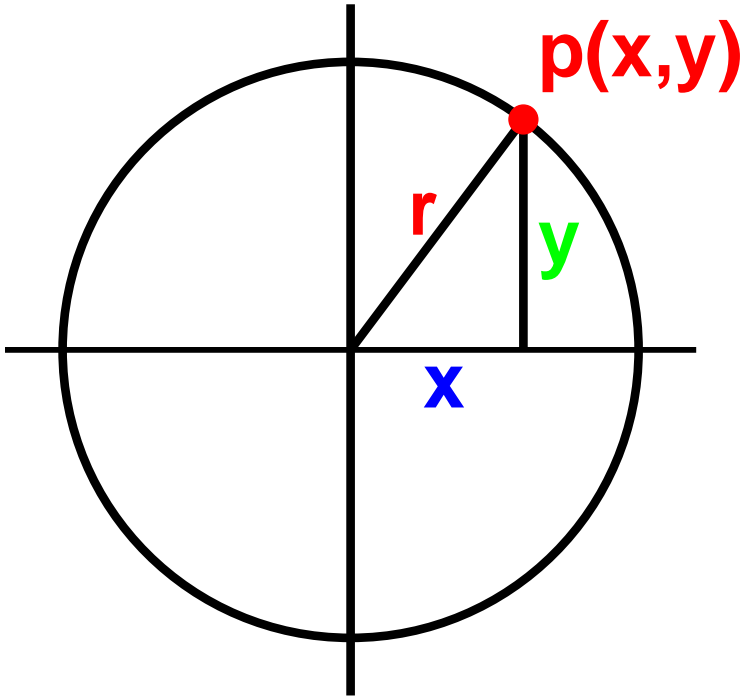


# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$

$$p(x, y) = p(x) \times p(y) \quad (3)$$

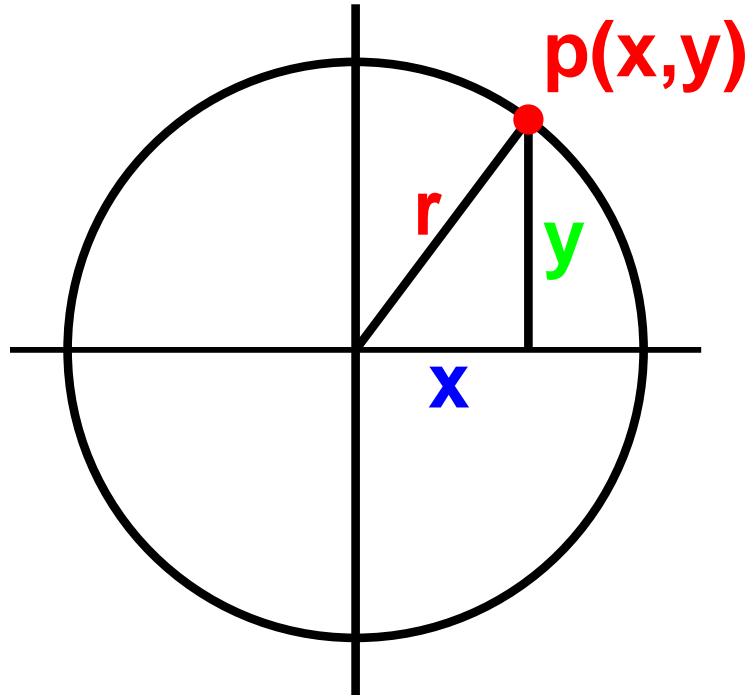




# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$



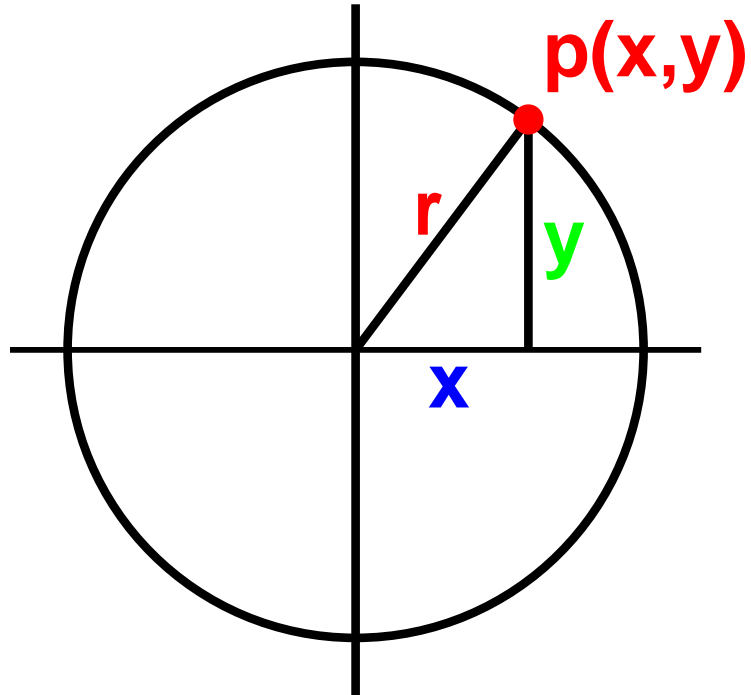
$$p(x, y) = p(x) \times p(y) \quad (3)$$

$$\propto e^{-x^2} \times e^{-y^2} \quad (4)$$

# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$



$$p(x, y) = p(x) \times p(y) \quad (3)$$

$$\propto e^{-x^2} \times e^{-y^2} \quad (4)$$

$$\propto e^{-(x^2+y^2)} \quad (5)$$

# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

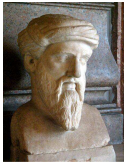
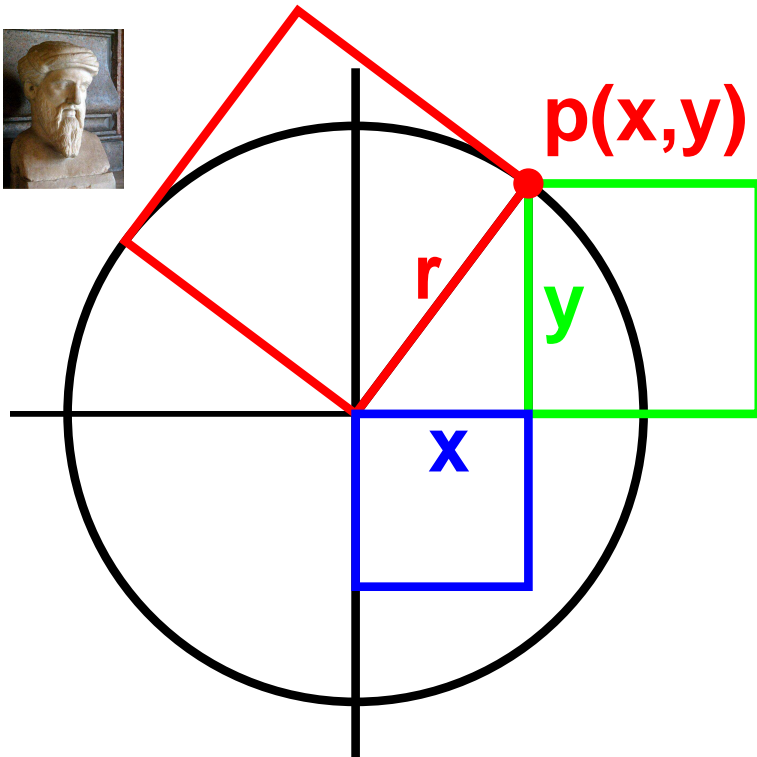
$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$

$$p(x, y) = p(x) \times p(y) \quad (3)$$

$$\propto e^{-x^2} \times e^{-y^2} \quad (4)$$

$$\propto e^{-(x^2+y^2)} \quad (5)$$

$$\propto e^{-r^2} \quad (6)$$



# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$

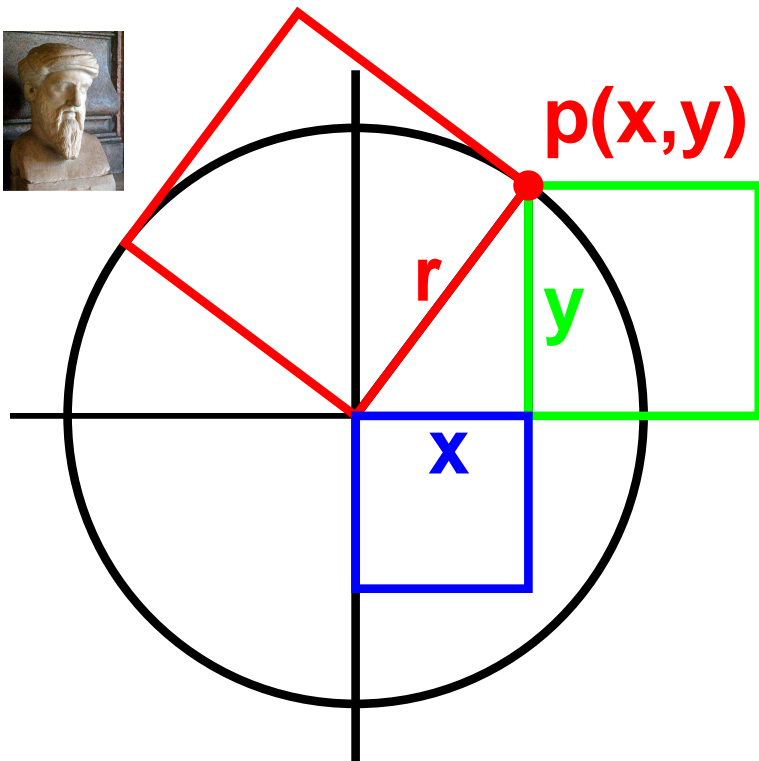
$$p(x, y) = p(x) \times p(y) \quad (3)$$

$$\propto e^{-x^2} \times e^{-y^2} \quad (4)$$

$$\propto e^{-(x^2+y^2)} \quad (5)$$

$$\propto e^{-r^2} \quad (6)$$

If  $p(x, y)$  is a constant,  
then  $r$  is a constant.



# Two Gaussians

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-x^2} \quad (1)$$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \propto e^{-y^2} \quad (2)$$

$$p(x, y) = p(x) \times p(y) \quad (3)$$

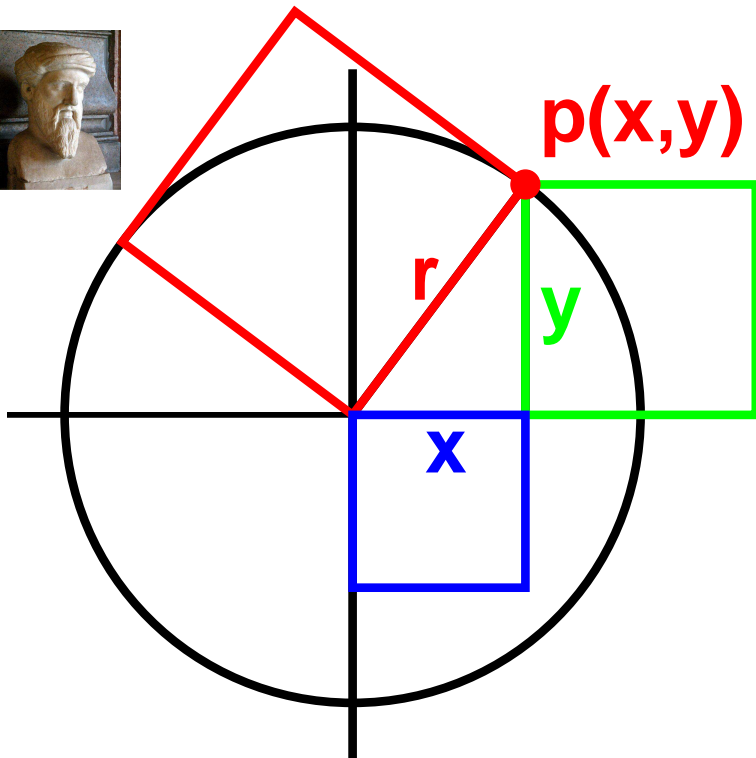
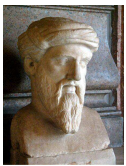
$$\propto e^{-x^2} \times e^{-y^2} \quad (4)$$

$$\propto e^{-(x^2+y^2)} \quad (5)$$

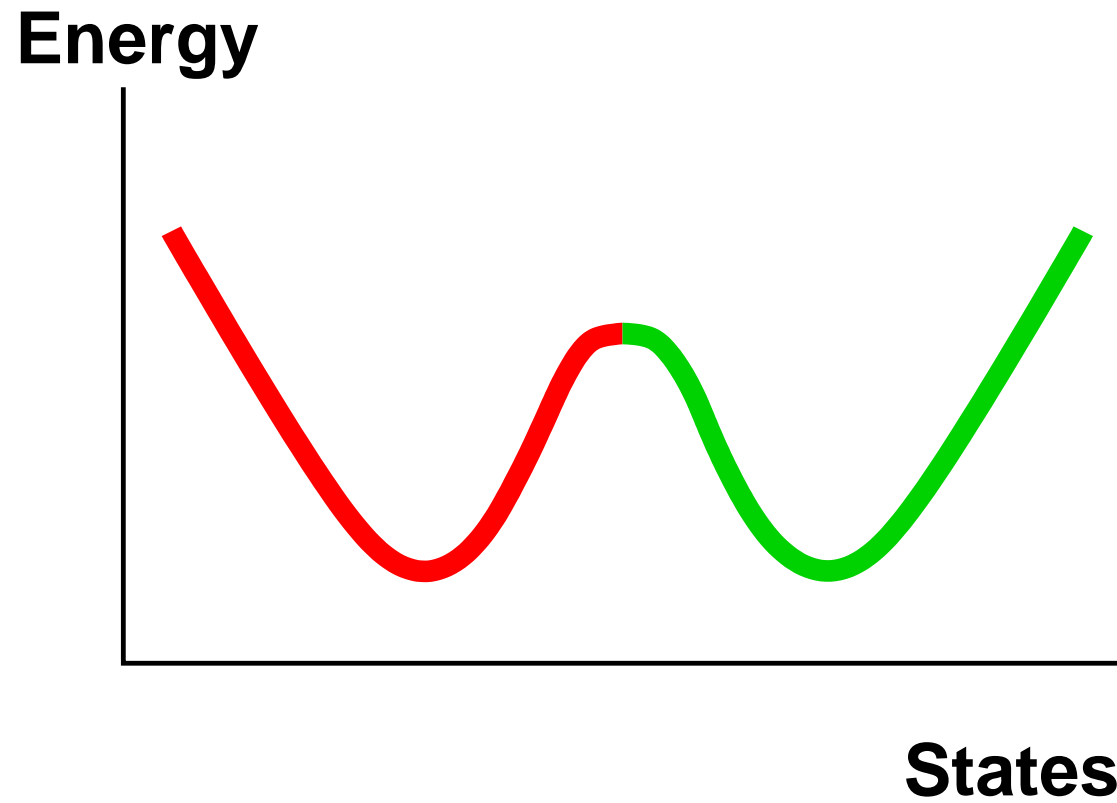
$$\propto e^{-r^2} \quad (6)$$

If  $p(x, y)$  is a constant,  
then  $r$  is a constant.

**Circular distribution!**



# 1 Dimension



**1 dimension is too simple!**

# Bowls in 2 Dimensions

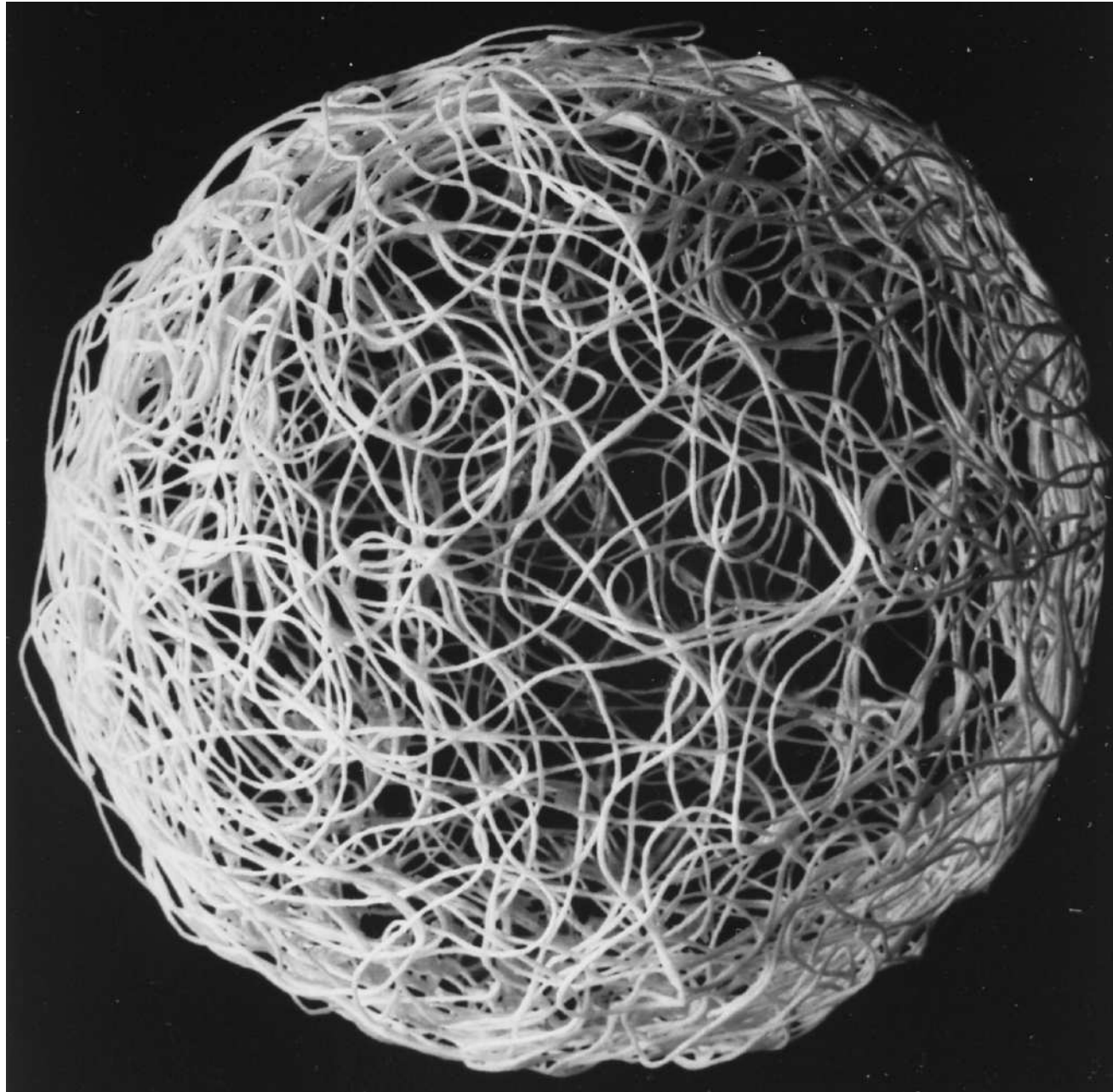


# Spheres in 3 Dimensions

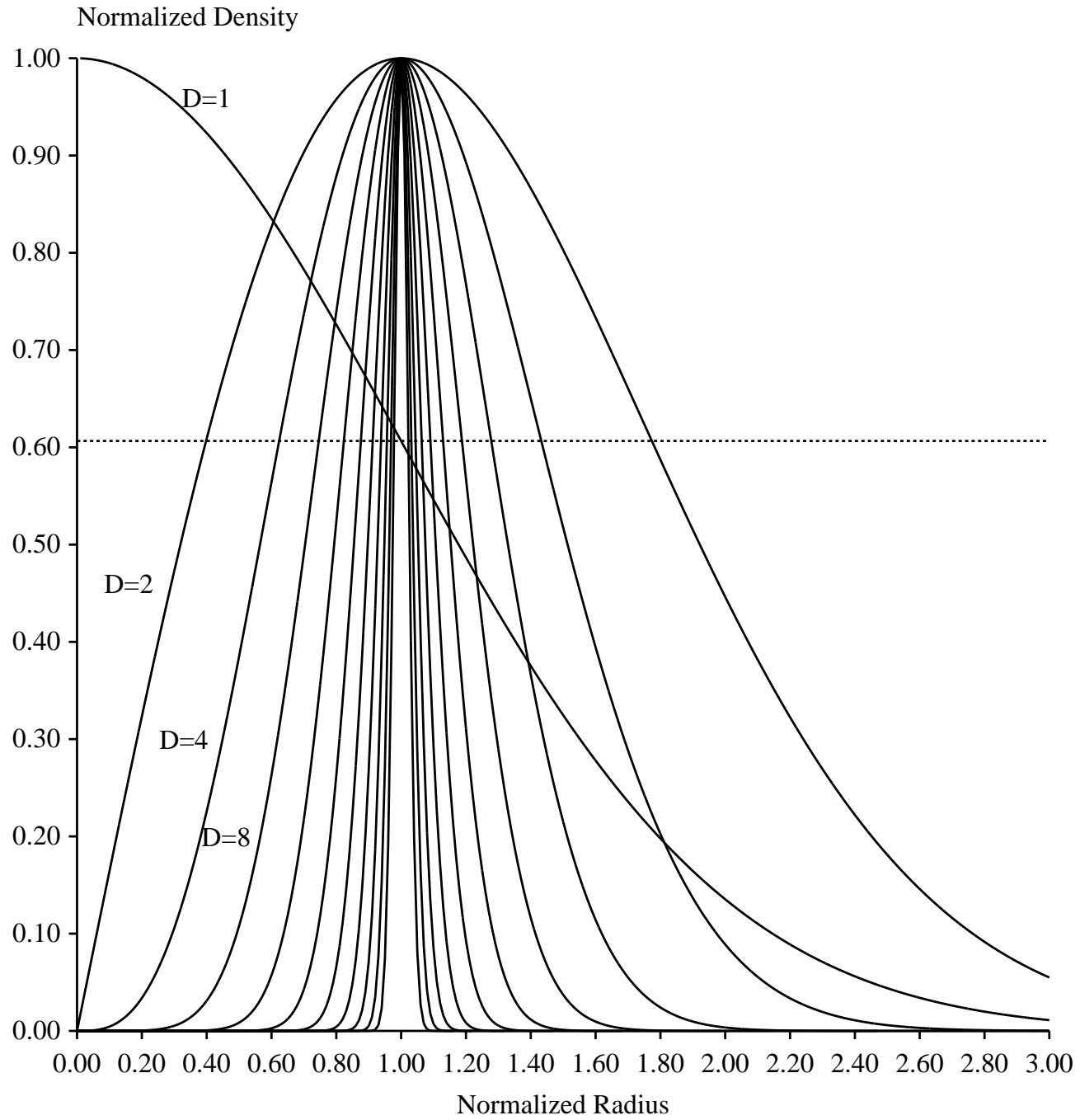




# N Dimensional Sphere



# Spheres tighten in high dimensions



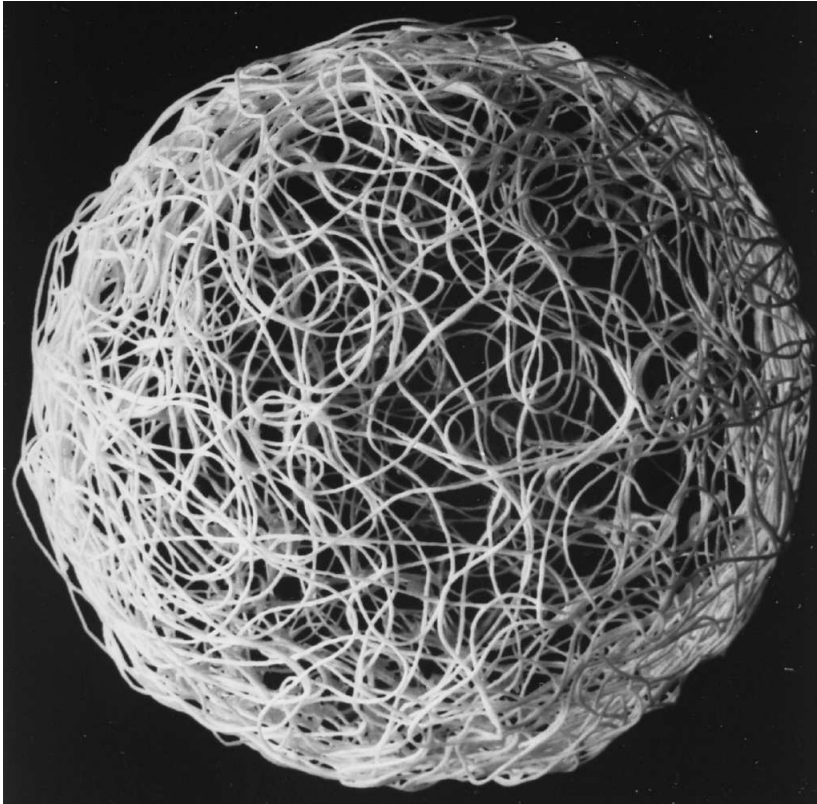
# Good Sphere Packing



Good packing of spheres  
gives a molecule  
the capacity  
to make selections efficiently

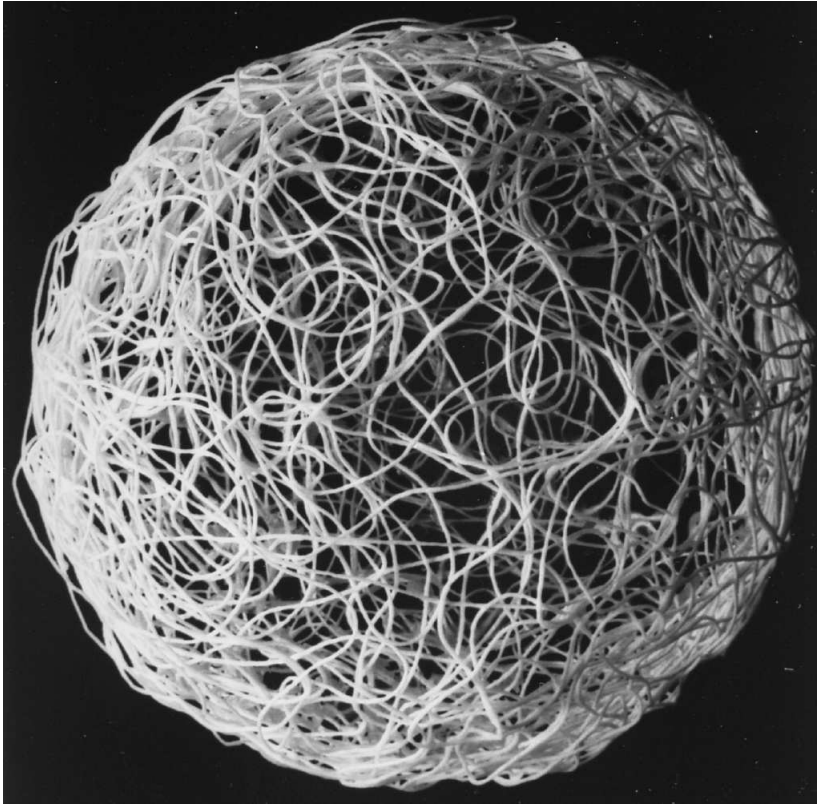
# N Dimensional Sphere Separation

Degenerate Sphere

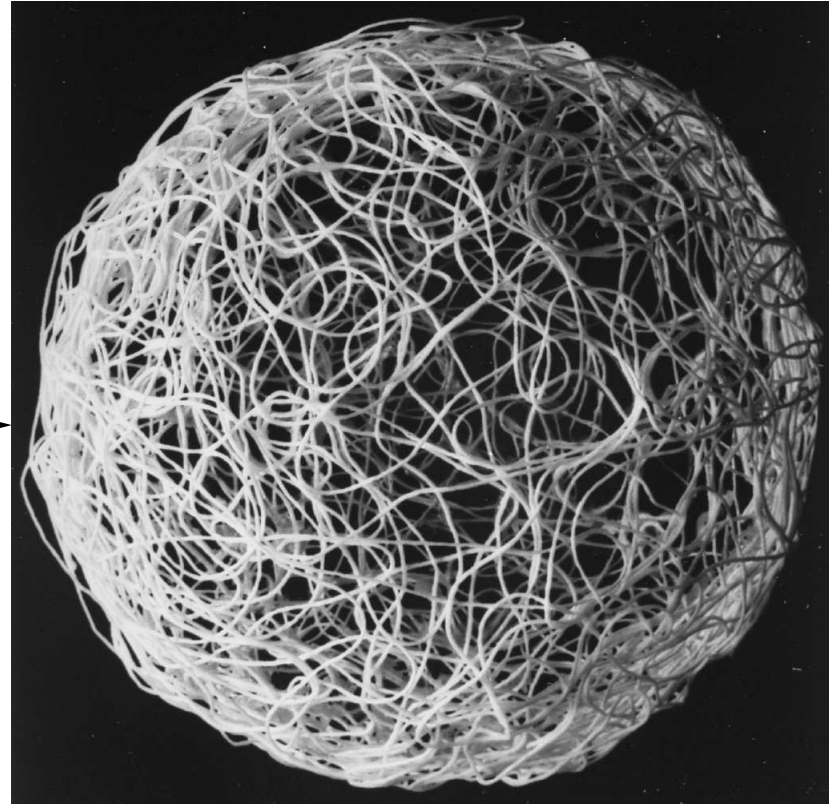


# N Dimensional Sphere Separation

Degenerate Sphere

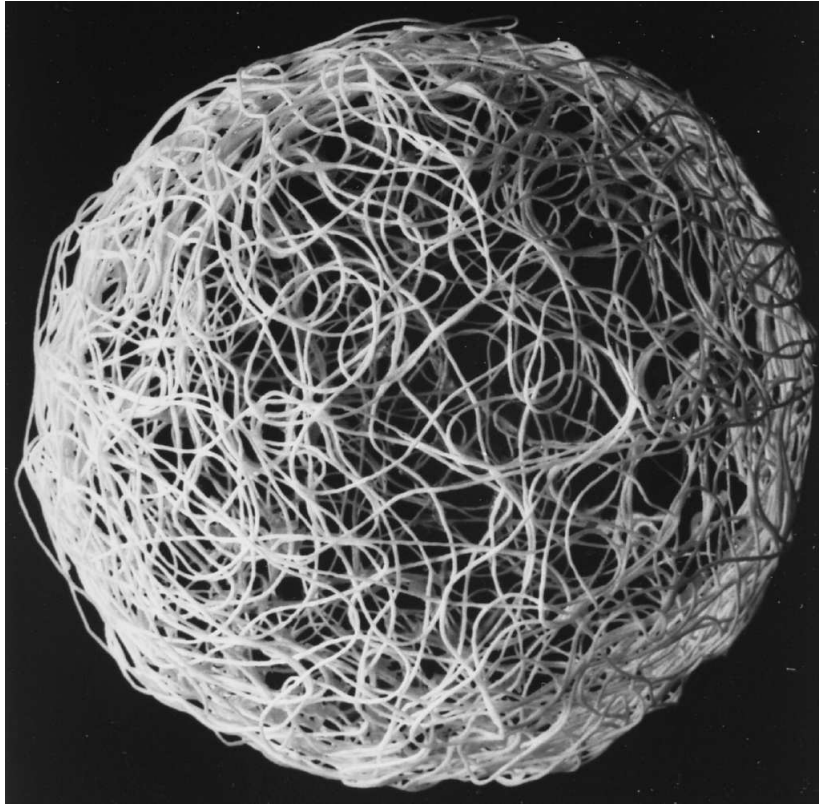


Forward Sphere

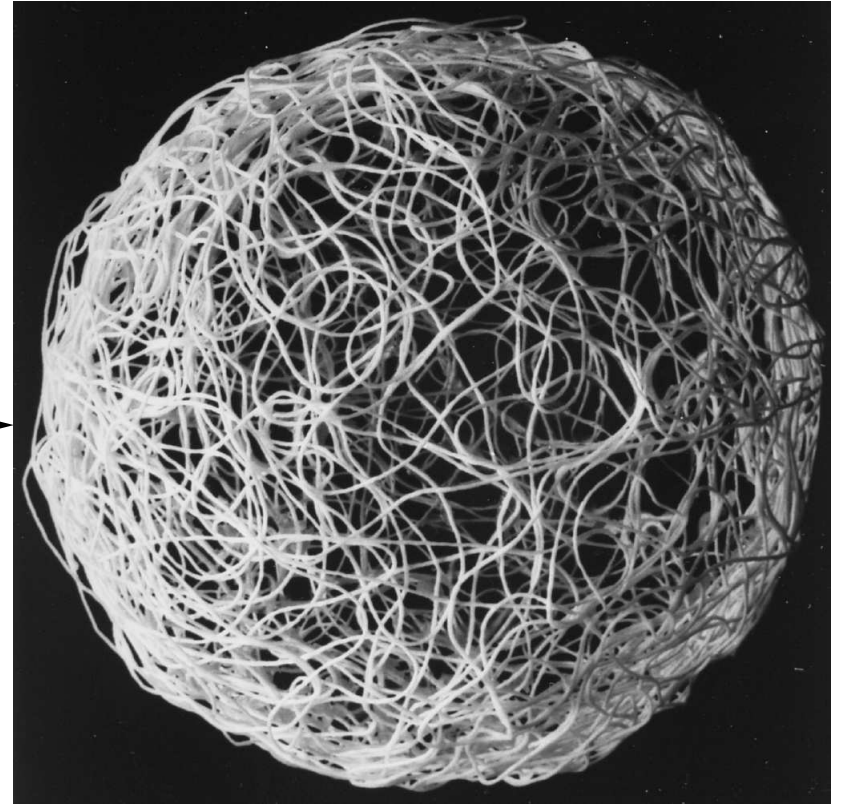


# N Dimensional Sphere Separation

Degenerate Sphere



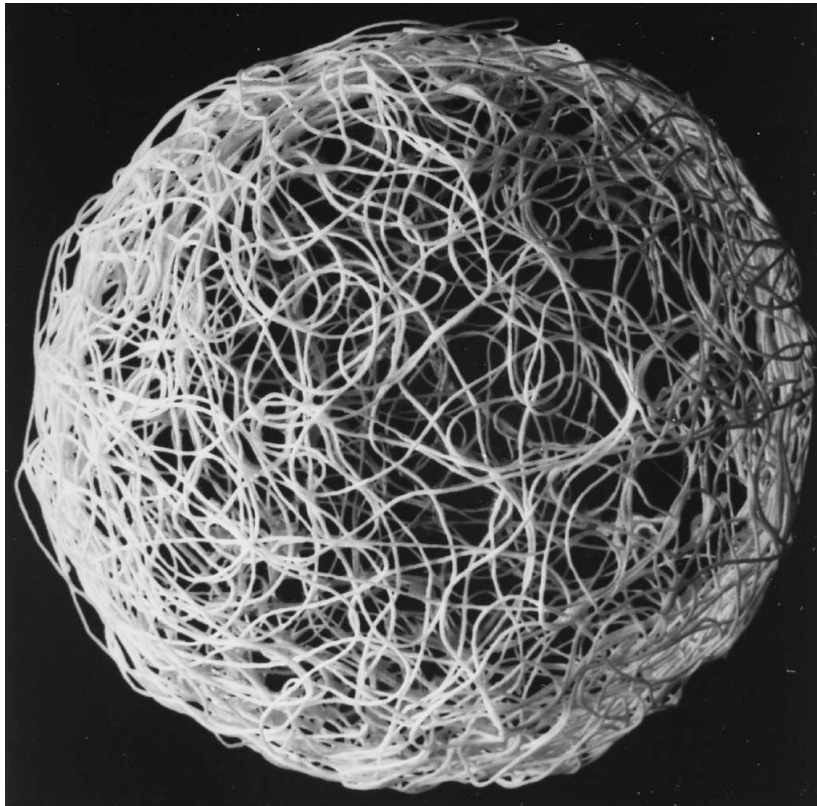
Forward Sphere



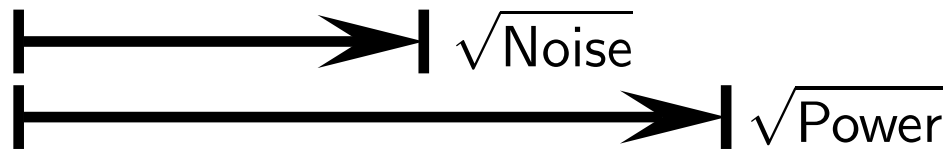
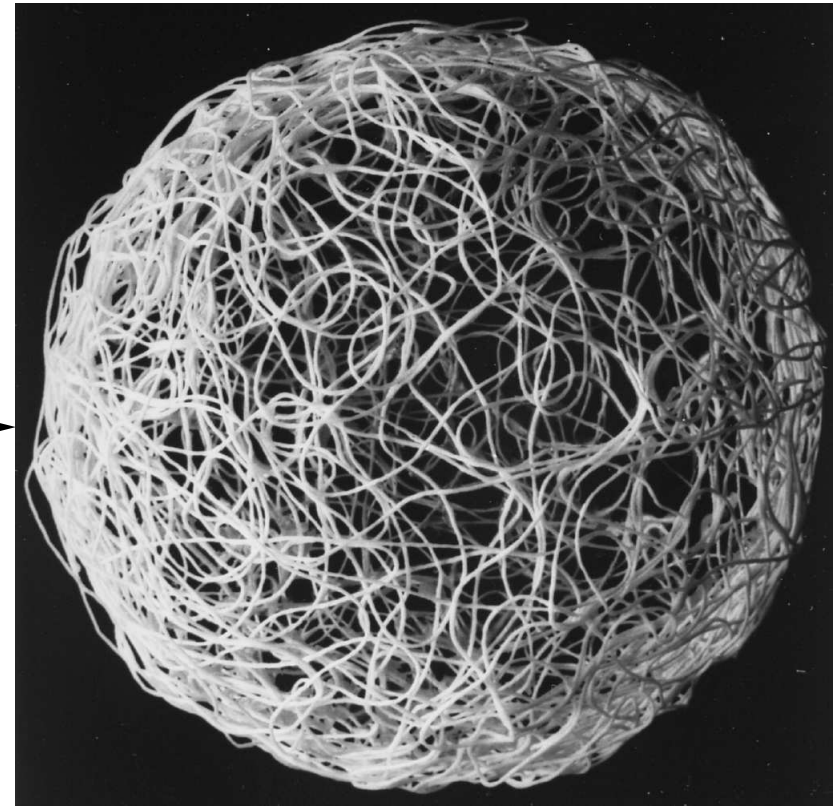
$\sqrt{\text{Noise}}$

# N Dimensional Sphere Separation

Degenerate Sphere

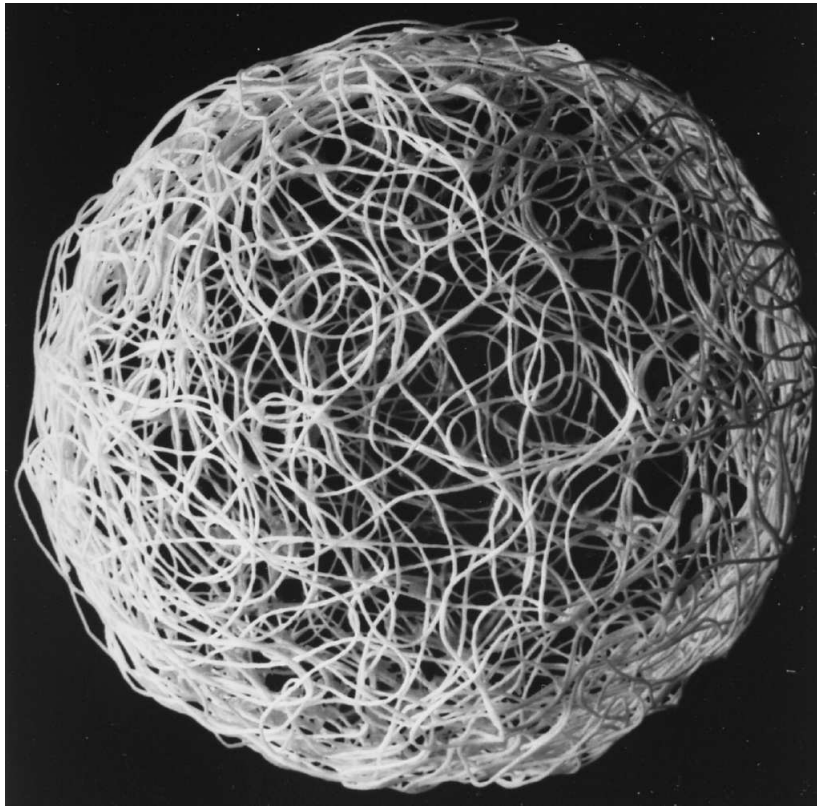


Forward Sphere

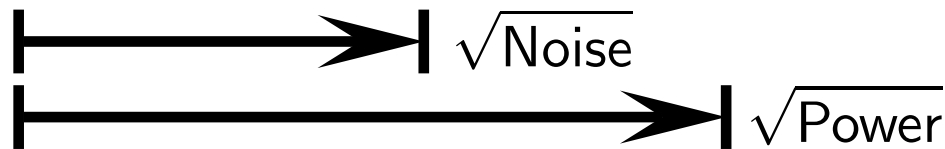
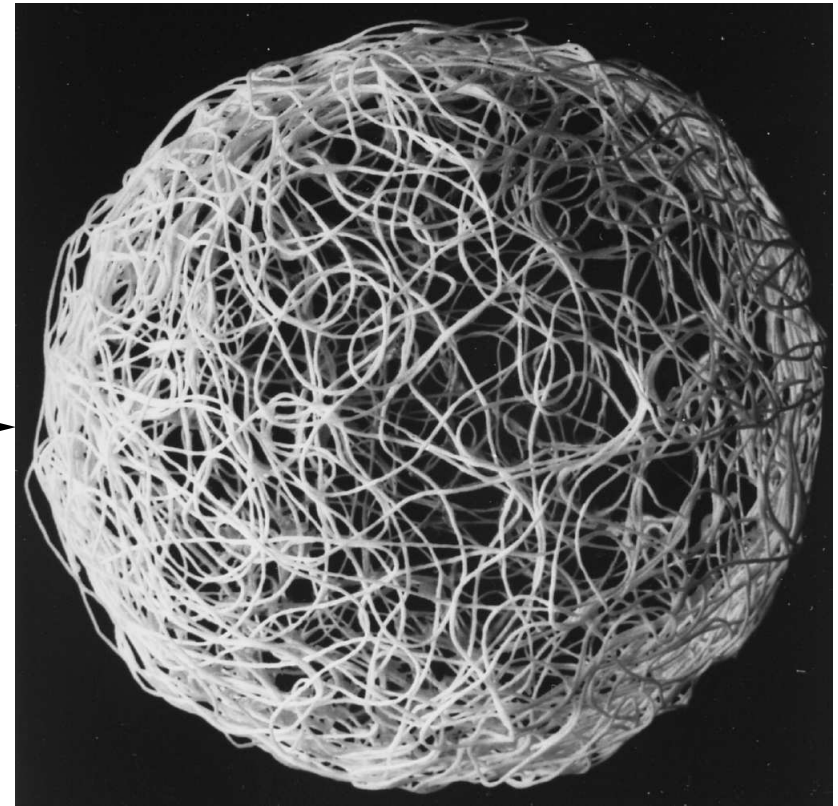


# N Dimensional Sphere Separation

Degenerate Sphere



Forward Sphere

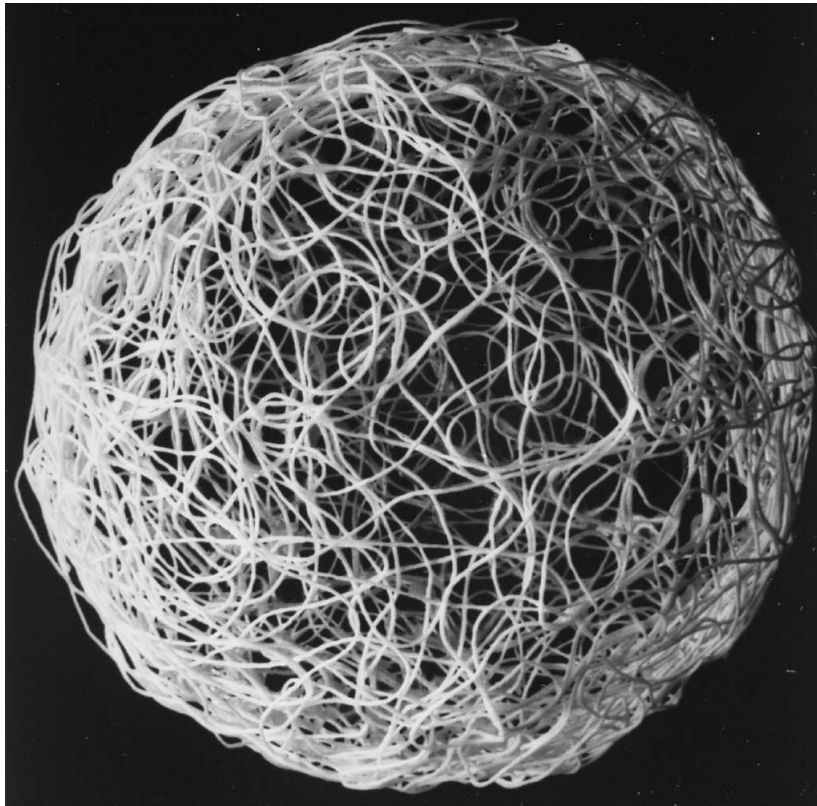


Energy dissipated to escape the Degenerate Sphere must exceed the Noise

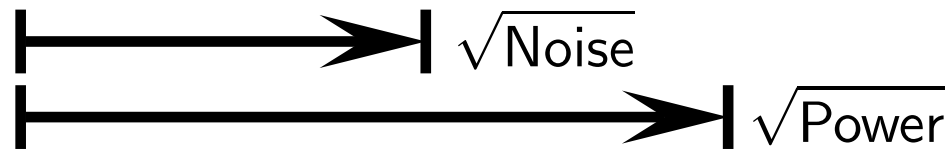
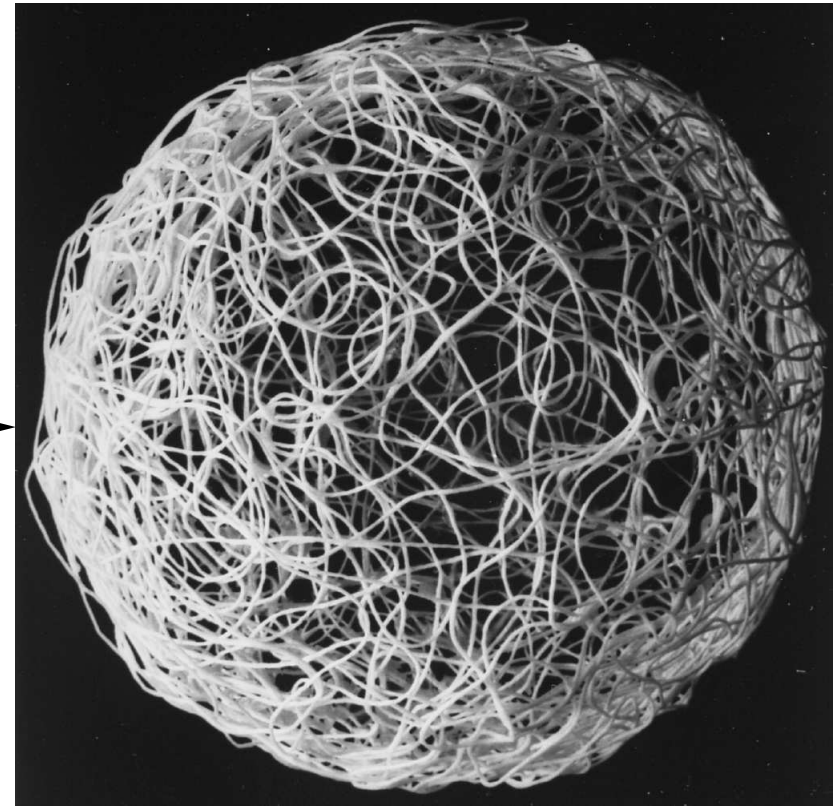


# N Dimensional Sphere Separation

Degenerate Sphere



Forward Sphere



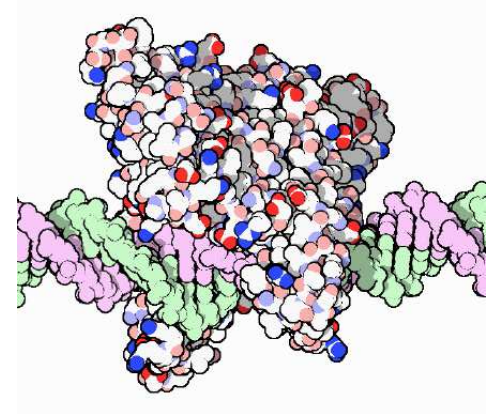
Energy dissipated to escape the Degenerate Sphere must exceed the Noise

$$\sqrt{\text{Power}} > \sqrt{\text{Noise}}$$

# Theoretical Isothermal Efficiency

- For molecular states of molecules with  $d_{space}$  'parts'  $P_y$  energy is dissipated for noise  $N_y$  and

$$C_y = d_{space} \log_2(P_y/N_y + 1) \leftarrow \text{machine capacity}$$

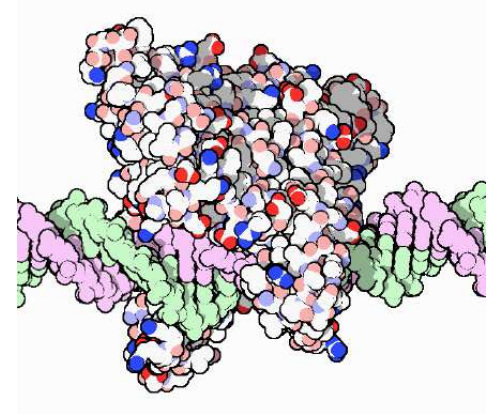


# Theoretical Isothermal Efficiency

- For molecular states of molecules with  $d_{space}$  'parts'  $P_y$  energy is dissipated for noise  $N_y$  and

$$C_y = d_{space} \log_2(P_y/N_y + 1) \leftarrow \text{machine capacity}$$

$$\epsilon_t = \frac{\ln\left(\frac{P_y}{N_y} + 1\right)}{\frac{P_y}{N_y}} \leftarrow \text{isothermal efficiency}$$

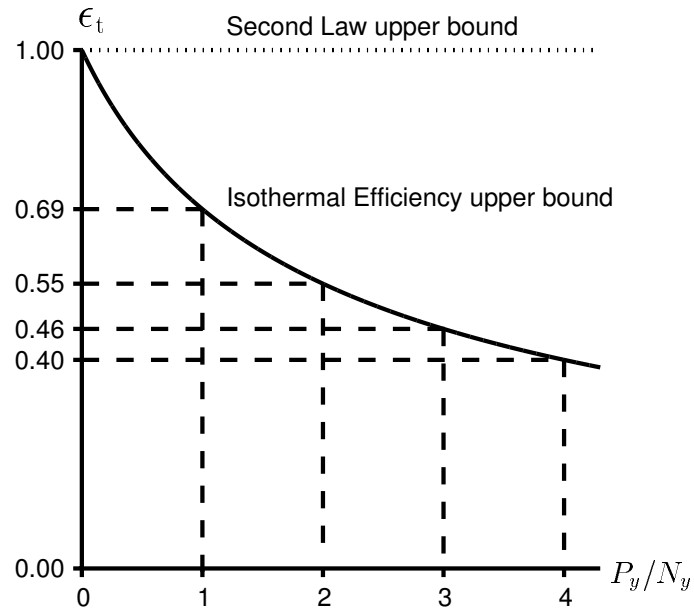
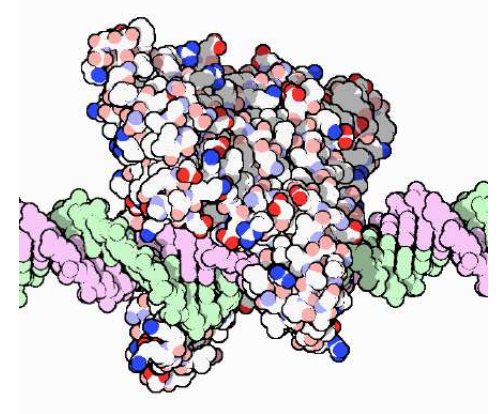


# Theoretical Isothermal Efficiency

- For molecular states of molecules with  $d_{space}$  'parts'  $P_y$  energy is dissipated for noise  $N_y$  and

$$C_y = d_{space} \log_2(P_y/N_y + 1) \leftarrow \text{machine capacity}$$

$$\epsilon_t = \frac{\ln\left(\frac{P_y}{N_y} + 1\right)}{\frac{P_y}{N_y}} \leftarrow \text{isothermal efficiency}$$



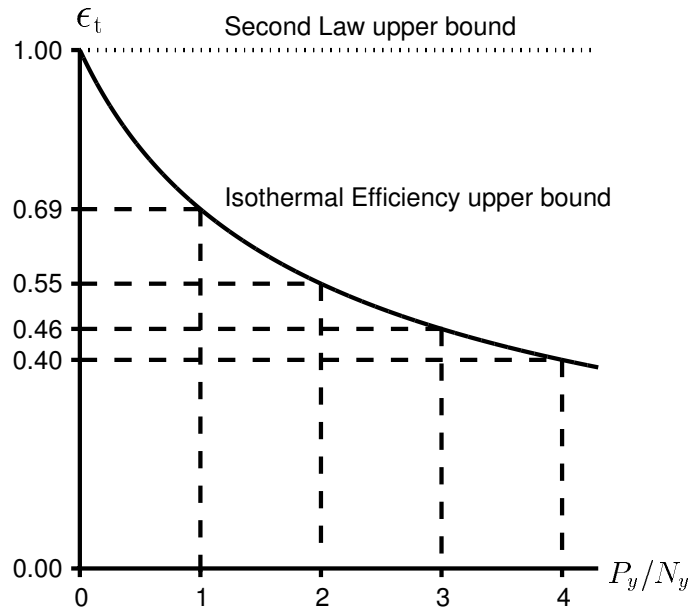
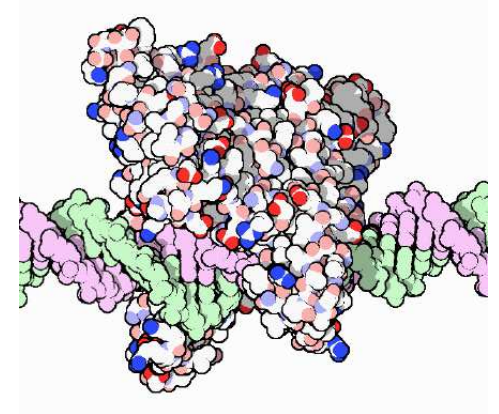
The curve is an upper bound

# Theoretical Isothermal Efficiency

- For molecular states of molecules with  $d_{space}$  'parts'  $P_y$  energy is dissipated for noise  $N_y$  and

$$C_y = d_{space} \log_2(P_y/N_y + 1) \leftarrow \text{machine capacity}$$

$$\epsilon_t = \frac{\ln\left(\frac{P_y}{N_y} + 1\right)}{\frac{P_y}{N_y}} \leftarrow \text{isothermal efficiency}$$

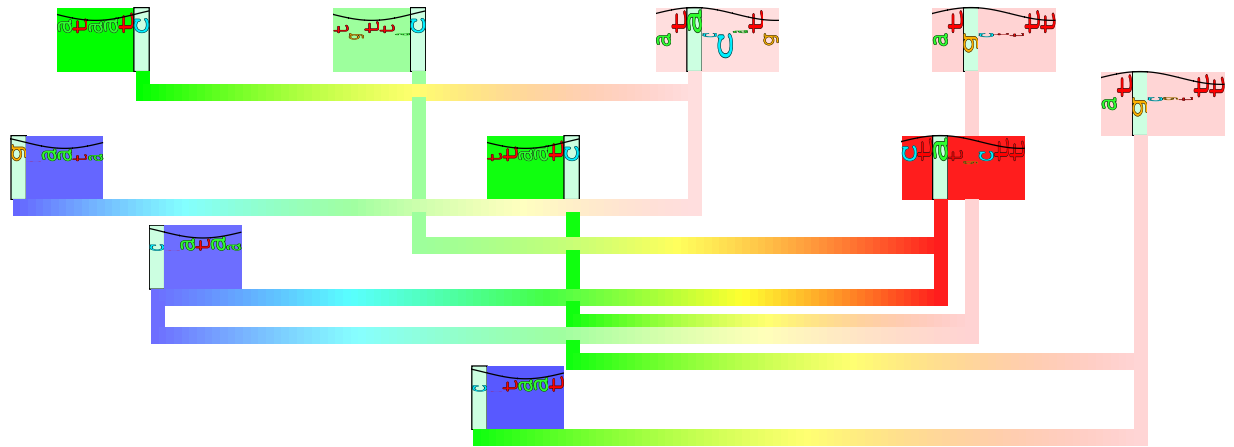


The curve is an upper bound

- If  $P_y/N_y = 1$  the efficiency is 70%!

# Acknowledgements

- Mentors:
  - Larry Gold (graduate school mentor)
  - Andrej Ehrenfeucht (information theory)
- Students:
  - Mike Stephens (logos, splicing)
  - John Spouge (mathematics)
  - Paul C. Anagnostopoulos (Evolution model)
  - Bruce Shapiro and Eckart Bindewald (3D logo)
  - Kevin Franco ( $\sigma^{38}$ )
  - Ding Jin ( $\sigma^{38}$ )
- Useful discussions with
  - Jeff Strathern
  - Amar Klar
  - Kemi Abolude
  - Susan Lauffer
  - Cedric Cagliero
  - Amar Klar
  - Zhi-Ming Zheng
  - Mark Lewandoski
- This research was supported by the Intramural Research Program of the NIH, National Cancer Institute, Center for Cancer Research.









# Version

version = 1.89 of kanpurtalk.tex 2016 Oct 16