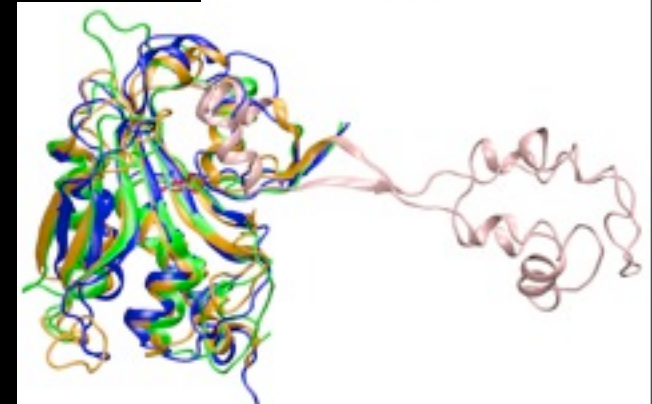


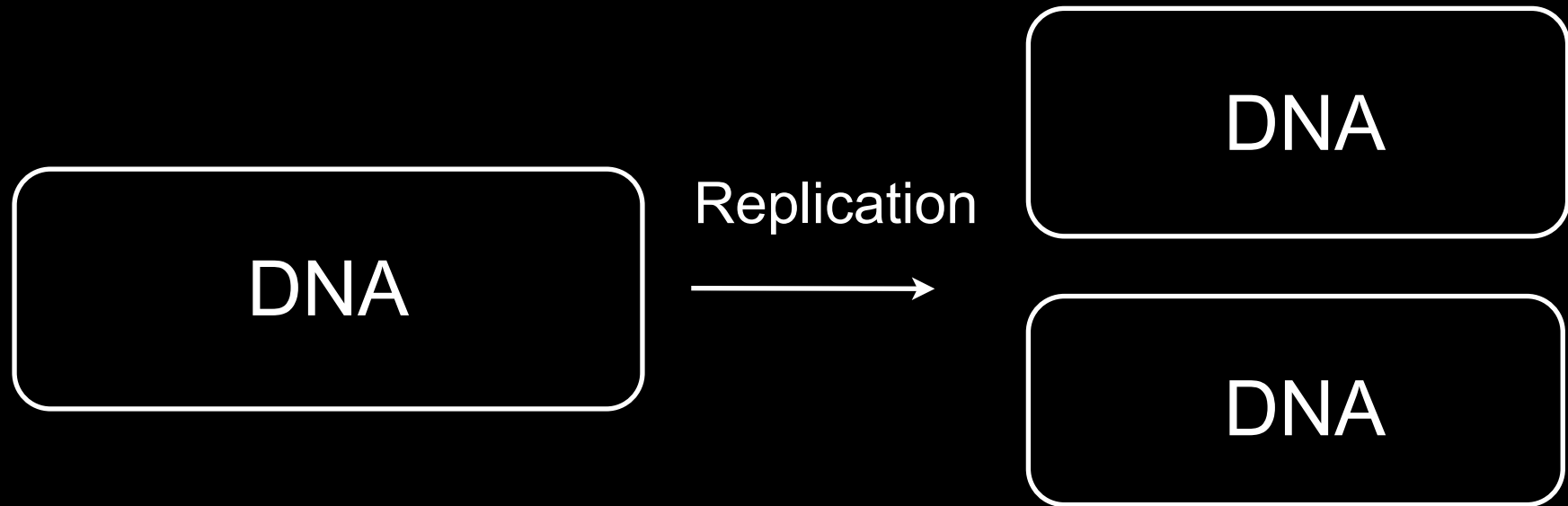
# Modeling Biological Information Processing Pathways

Anurag Sethi  
LANL  
Yale University  
June 2013

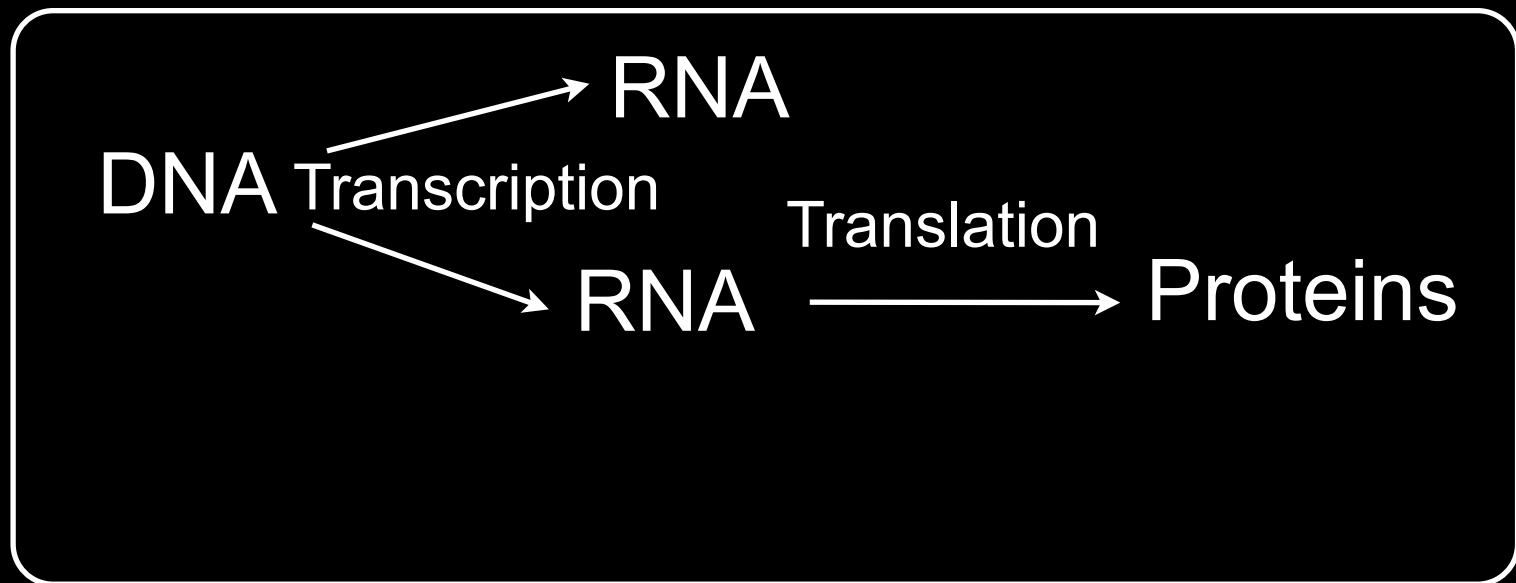
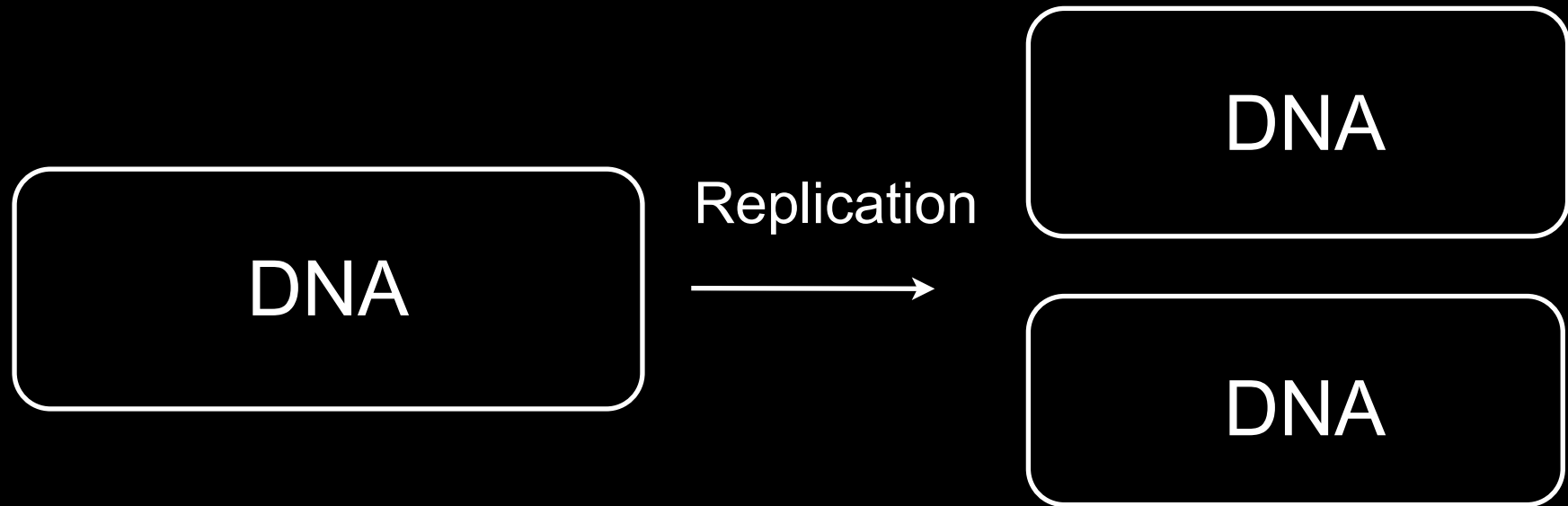


# Life is Information Processing

# Life is Information Processing

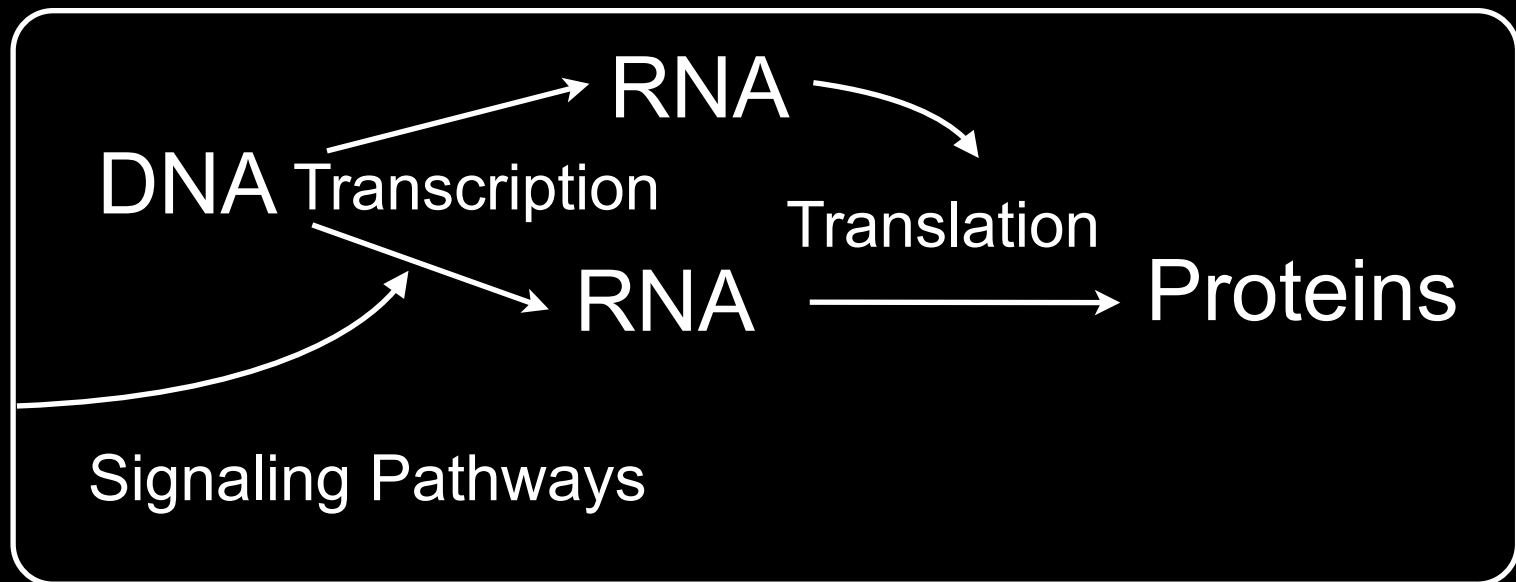
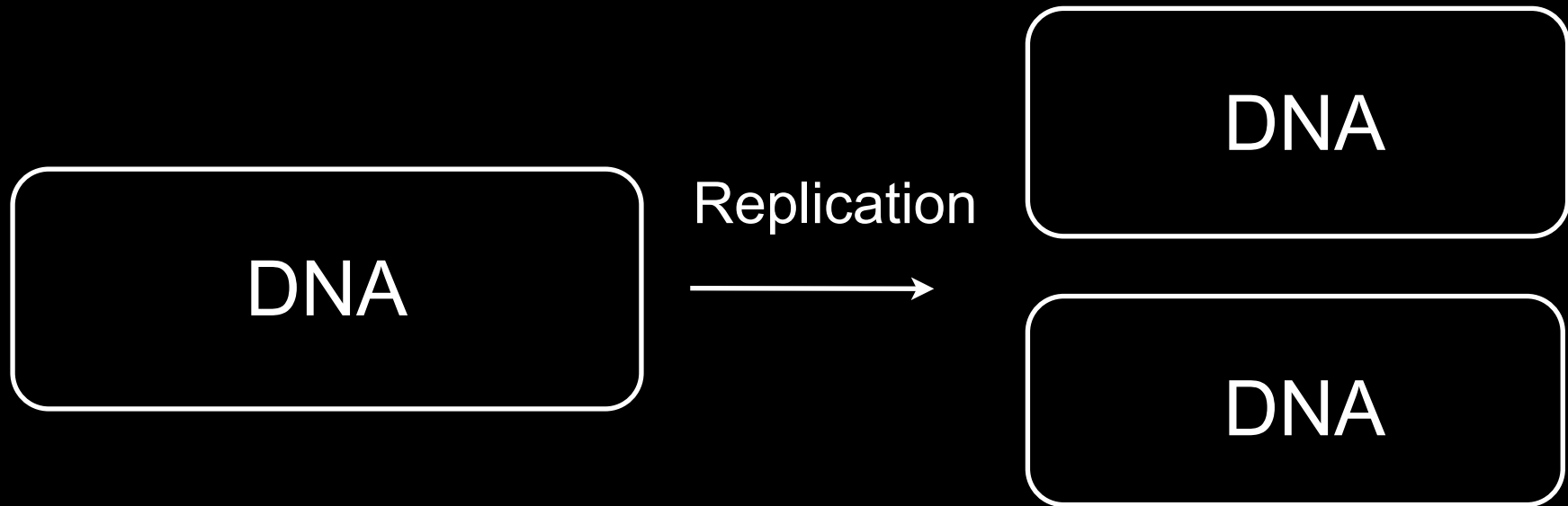


# Life is Information Processing



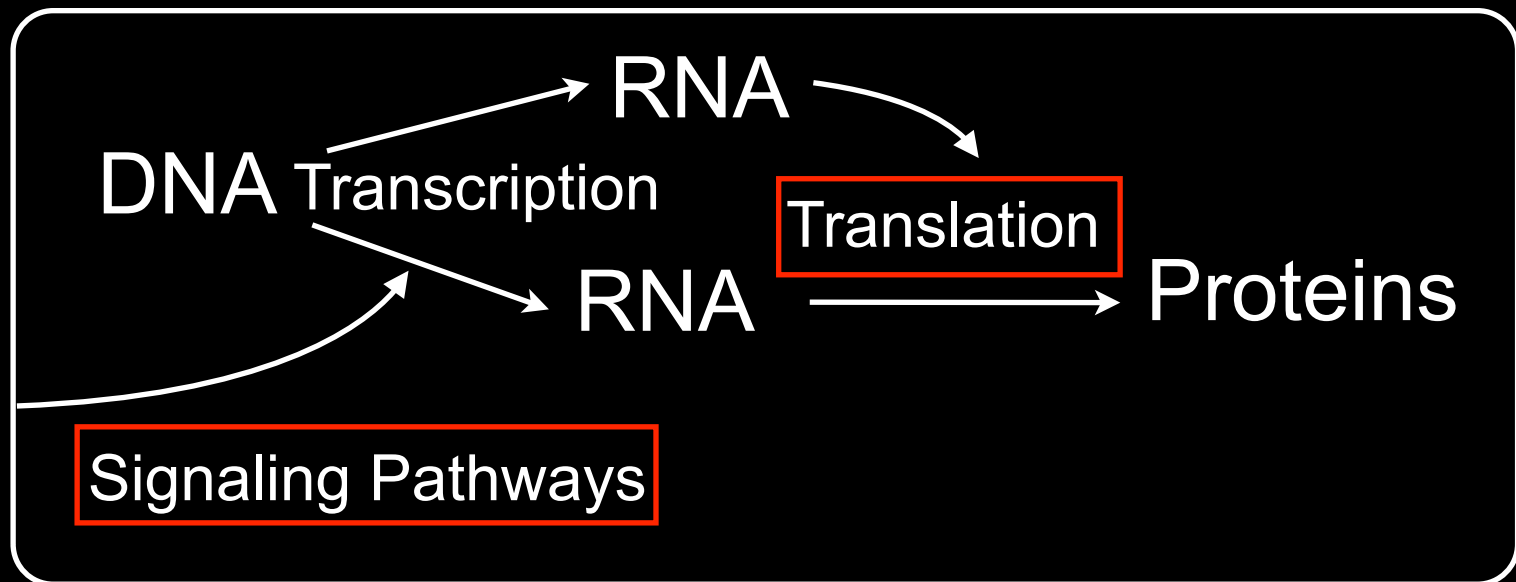
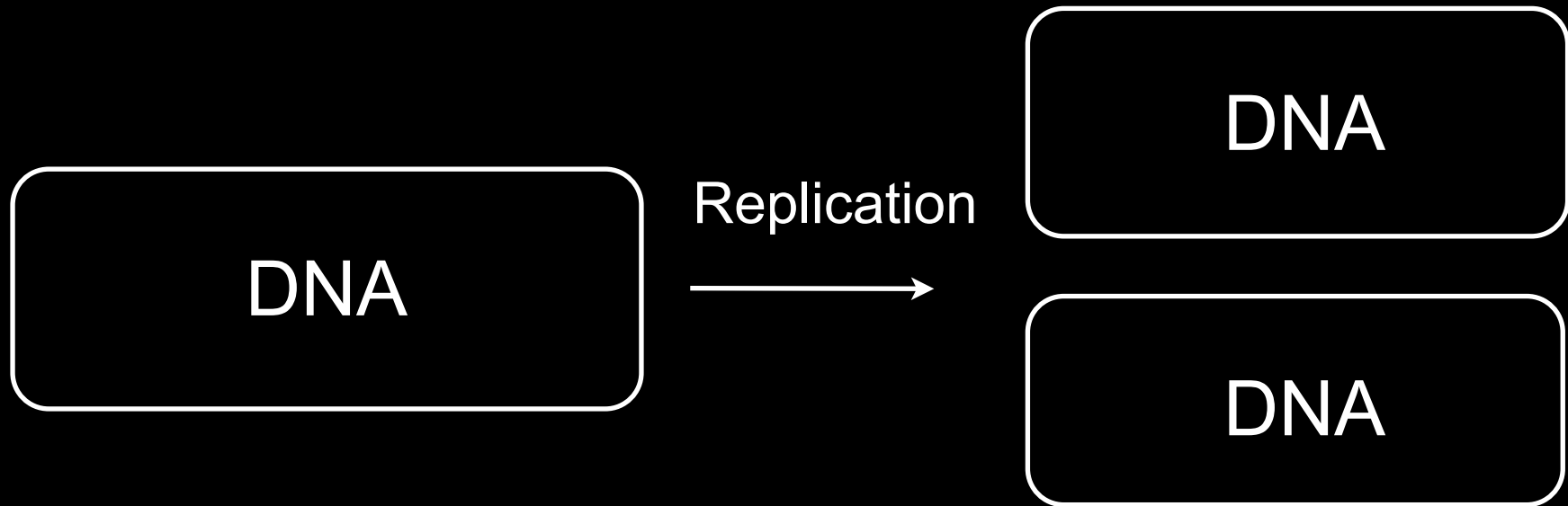
Genes are translated into proteins and proteins perform functions subject to cues from the environment.

# Life is Information Processing



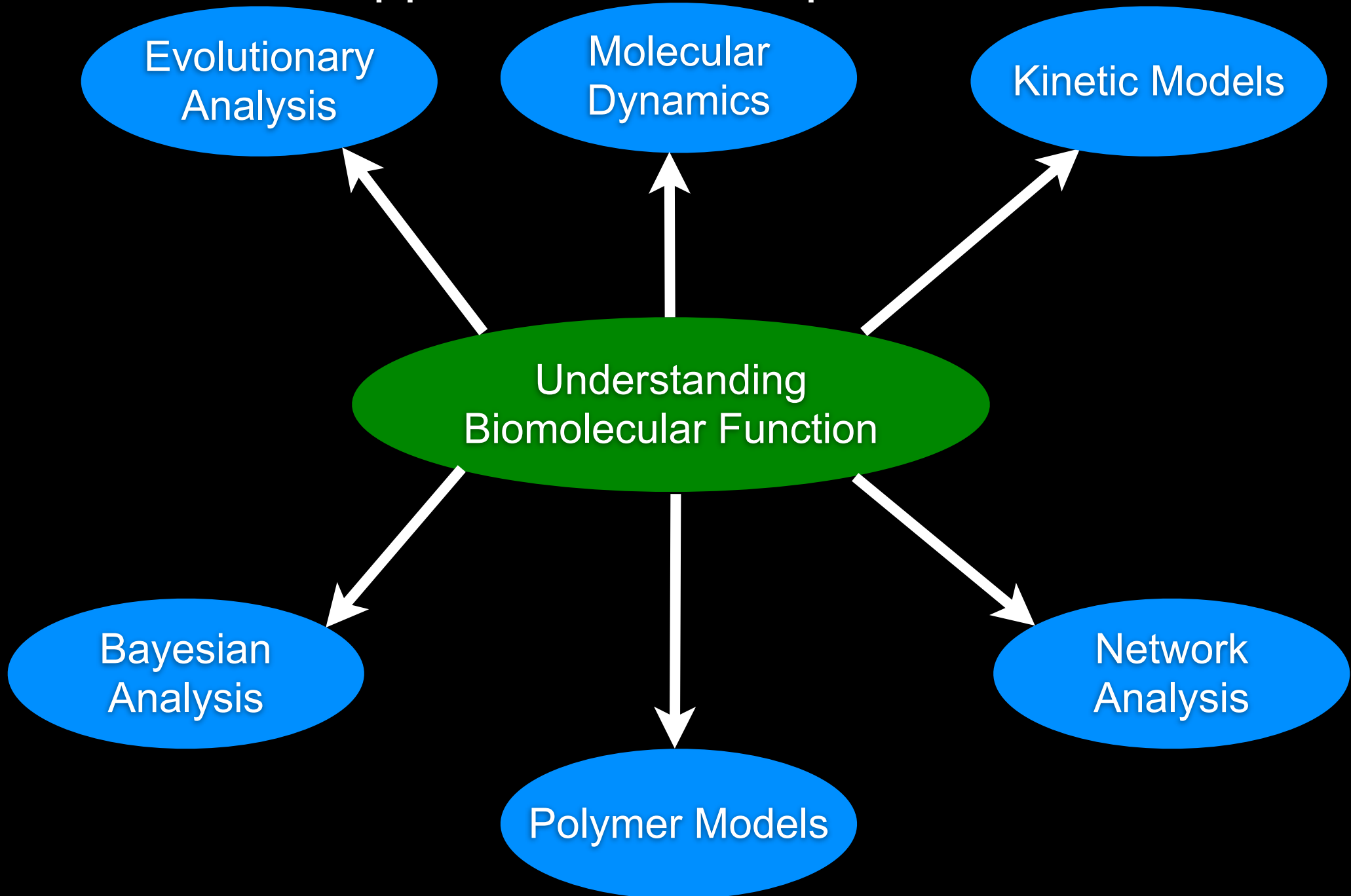
Genes are translated into proteins and proteins perform functions subject to cues from the environment.

# Life is Information Processing



Genes are translated into proteins and proteins perform functions subject to cues from the environment.

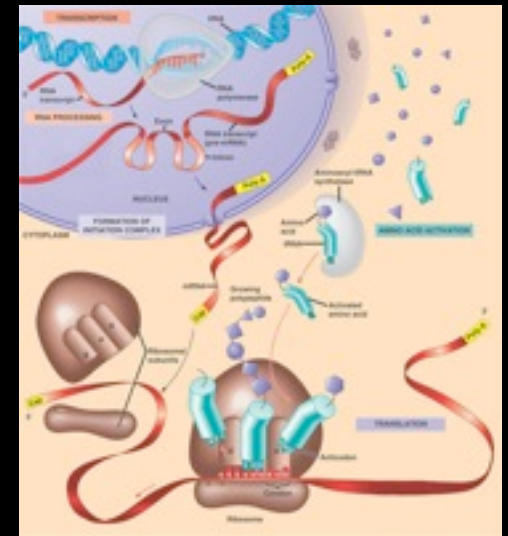
# Investigation of complex systems requires the application of multiple tools



# Organization of Talk

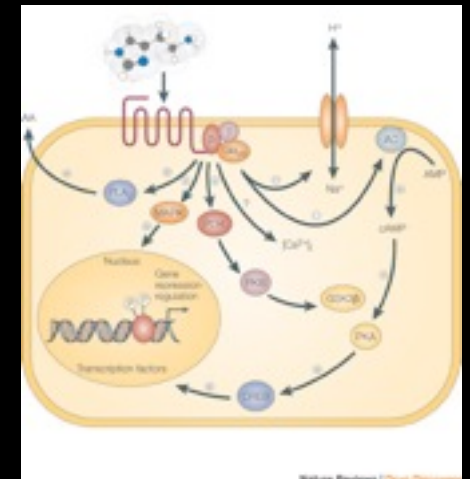
## Translating Information in RNA

Evolutionary analysis of biomolecules  
Allosteric signaling pathways



## Disordered Regions in Signaling Cascades

Multivalent Proteins  
Modeling Disordered Proteins

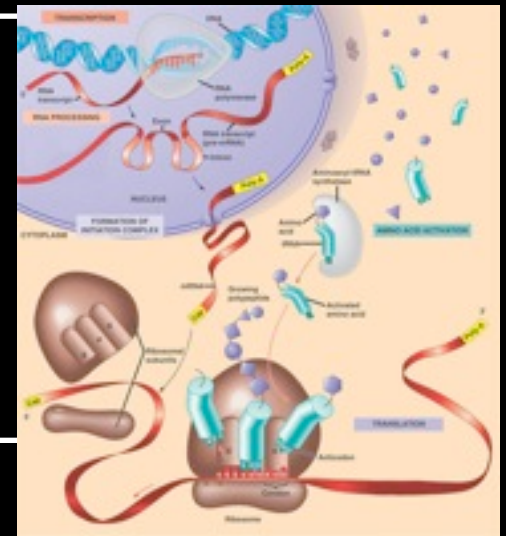




# Organization of Talk

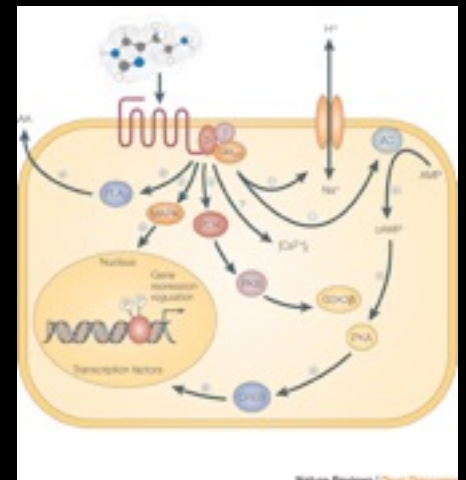
## Translating Information in RNA

Evolutionary analysis of biomolecules  
Allosteric signaling pathways



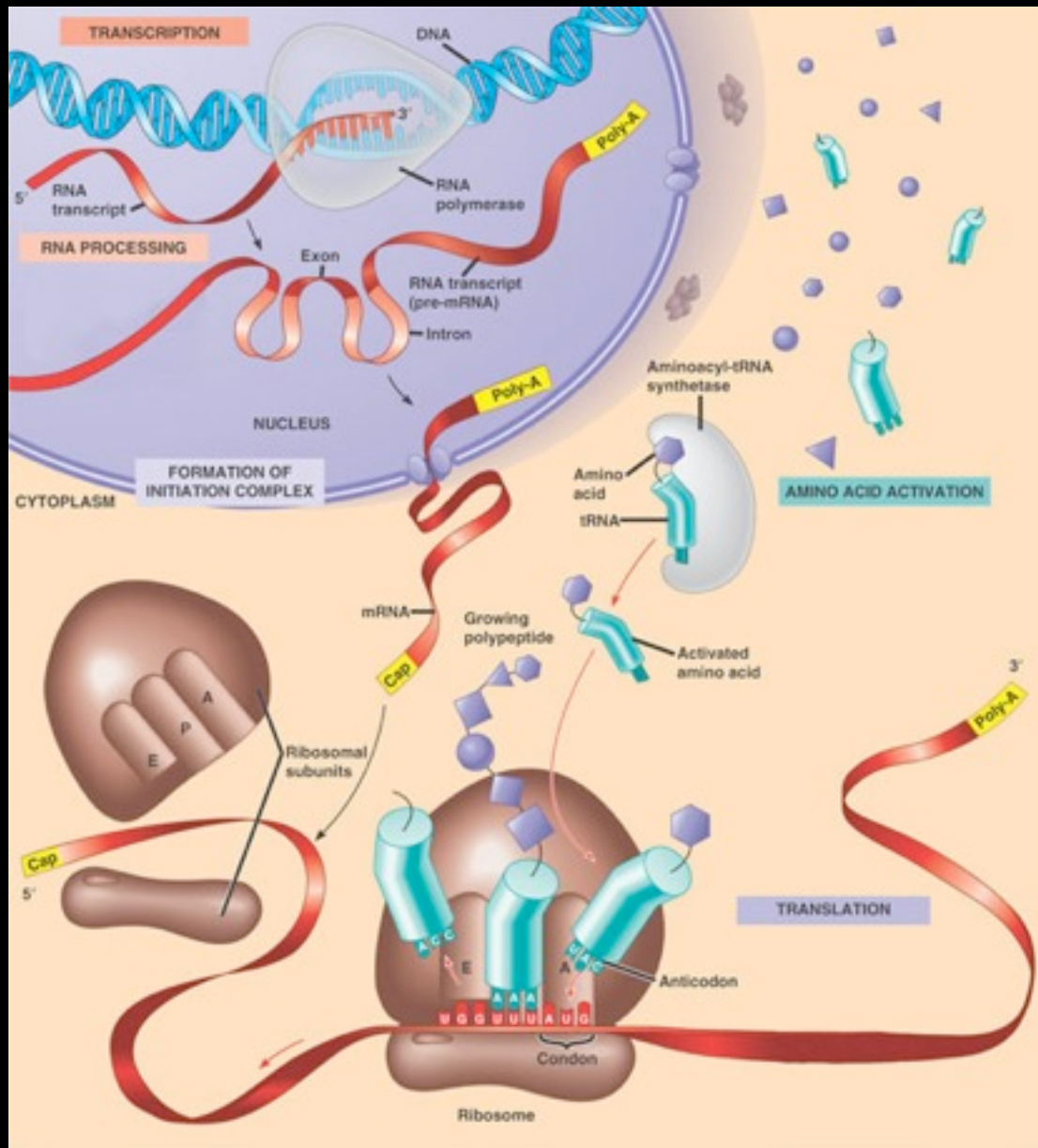
## Disordered Regions in Signaling Cascades

Multivalent Proteins  
Modeling Disordered Proteins



## Setting the Genetic Code: Evolutionary Analysis

The central dogma of molecular biology describes how DNA is translated into proteins



The universal genetic code is used to translate the information in the DNA into functional proteins.

## Setting the Genetic Code: Evolutionary Analysis

The genetic code is set by enzymes called the aminoacyl-tRNA synthetases



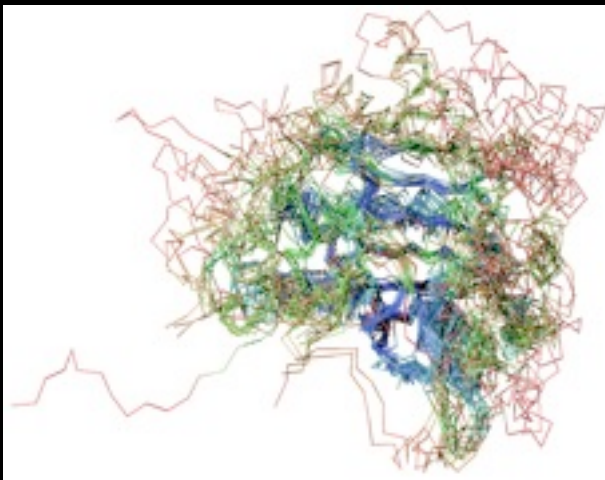
		Second Letter				
		T	C	A	G	
First Letter	T	TTT } Phe TTC } TTA } Leu TTG }	TCT } Ser TCC } TCA } TCG }	TAT } Tyr TAC } TAA Stop TAG Stop	TGT } Cys TGC } TGA Stop TGG Trp	T C A G
	C	CTT } Leu CTC } CTA } CTG }	CCT } Pro CCC } CCA } CCG }	CAT } His CAC } CAA } Gln CAG }	CGT } Arg CGC } CGA } CGG }	T C A G
	A	ATT } Ile ATC } ATA } ATG Met	ACT } Thr ACC } ACA } ACG }	AAT } Asn AAC } AAA } Lys AAG }	AGT } Ser AGC } AGA } Arg AGG }	T C A G
	G	GTT } Val GTC } GTA } GTG }	GCT } Ala GCC } GCA } GCG }	GAT } Asp GAC } GAA } Glu GAG }	GGT } Gly GGC } GGA } GGG }	T C A G

How do complex biomolecules perform their function?

# Setting the Genetic Code: Evolutionary Analysis

## Goal of Evolutionary Analysis

Nothing in biology makes sense except in the light of evolution. - Theodosius Dobzhansky



Structural alignment

Sos1		P1	P2	P3	RP	P4
<i>H. sapiens</i>	1149	VPPVPPRRR	20 SPPAIPPRQP	22 SPPLLPREP	48 GTRRHLPSPP	11 AGPPVPPRQS
<i>M. musculus</i>	1135	VPPVPPRRR	20 SPPAIPPRQP	22 SPPLLPREP	48 GTRRHLPSPP	11 AGPPVPPRQS
<i>E. caballus</i>	1288	VPPVPPRRR	20 SPPAIPPRQP	22 SPPLLPREP	48 GTRRHLPSPP	11 AGPPVPPRQS
Sos2		M1	M2	M3	M4	M5
<i>H. sapiens</i>	1144	IPPLPPRKK	18 DPPAIPPRQP	35 TPPPVPLRPP	48 PSPRVPRRCY	12 PAPPVPPRQN
<i>M. musculus</i>	1145	IPPLPPRKK	18 DPPAIPPRQP	35 TPPPVPLRPP	48 PSPRIPRSCH	12 PAPPVPPRQN
<i>G. gallus</i>	1280	NPPPLPPRKK	18 DPPAIPPRQP	35 TPPPVHRPP	48 PSPRIPRRCH	12 PAPPVPPRQN
Sos		S1	S2	S3		
<i>D. melanogaster</i>	1342	VPPVPPRRR	18 DAPTLPPRDG	3 SPPPIPPRLN		
<i>D. mojavensis</i>	1325	VPPVPPRRR	18 DAPILPPRDG	3 SPPPIPPRLN		
<i>D. erecta</i>	1335	VPPVPPRRR	18 DAPTLPPRDG	3 SPPPIPPRLN		

Sequence alignment

Conservation analysis provides information on the constraints in the evolution of biomolecular families.

## Setting the Genetic Code: Evolutionary Analysis

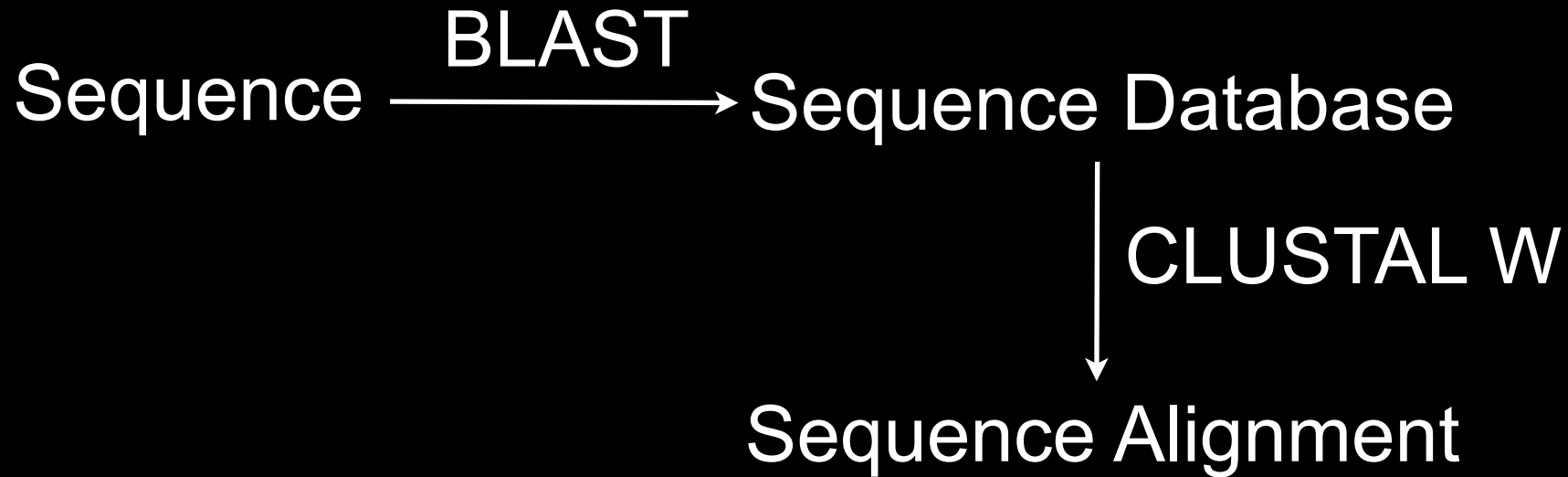
# Evolutionary Analysis

Sequence  $\xrightarrow{\text{BLAST}}$  Sequence Database

The sequences used to represent the evolutionary history of a biomolecule are the sequences found in a database.

## Setting the Genetic Code: Evolutionary Analysis

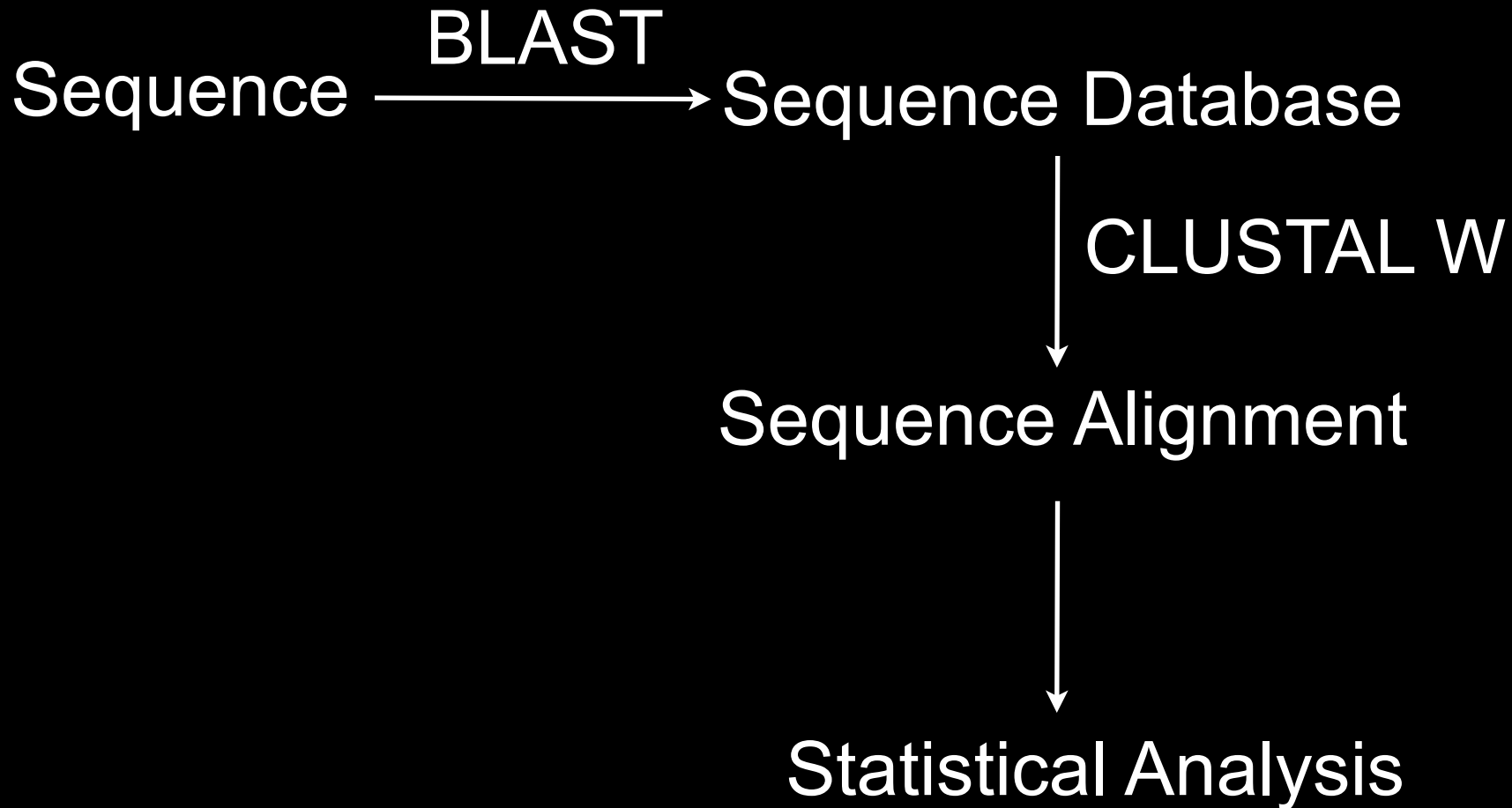
# Evolutionary Analysis



The sequences used to represent the evolutionary history of a biomolecule are the sequences found in a database.

## Setting the Genetic Code: Evolutionary Analysis

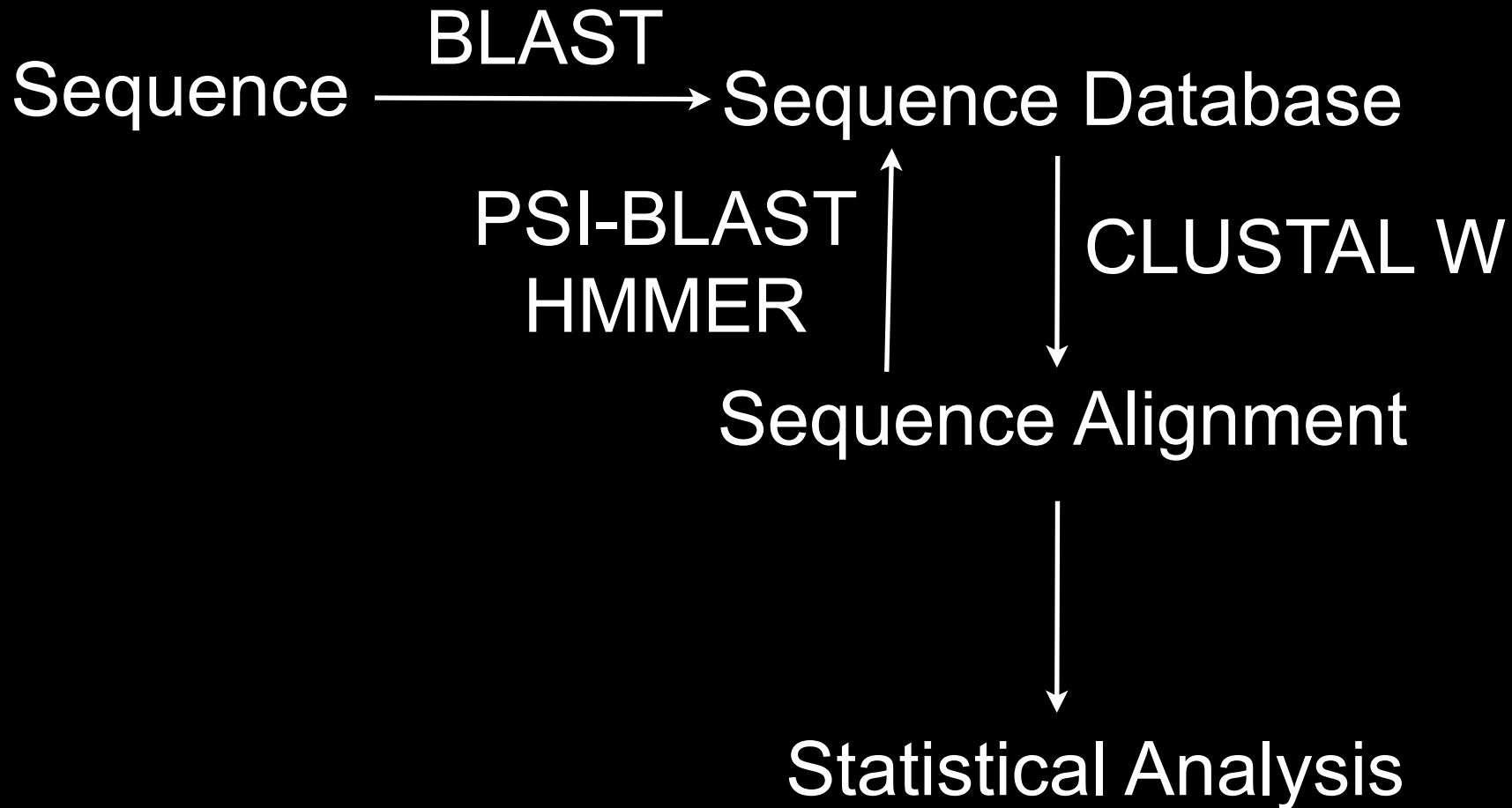
# Evolutionary Analysis



The sequences used to represent the evolutionary history of a biomolecule are the sequences found in a database.

# Setting the Genetic Code: Evolutionary Analysis

## Evolutionary Analysis

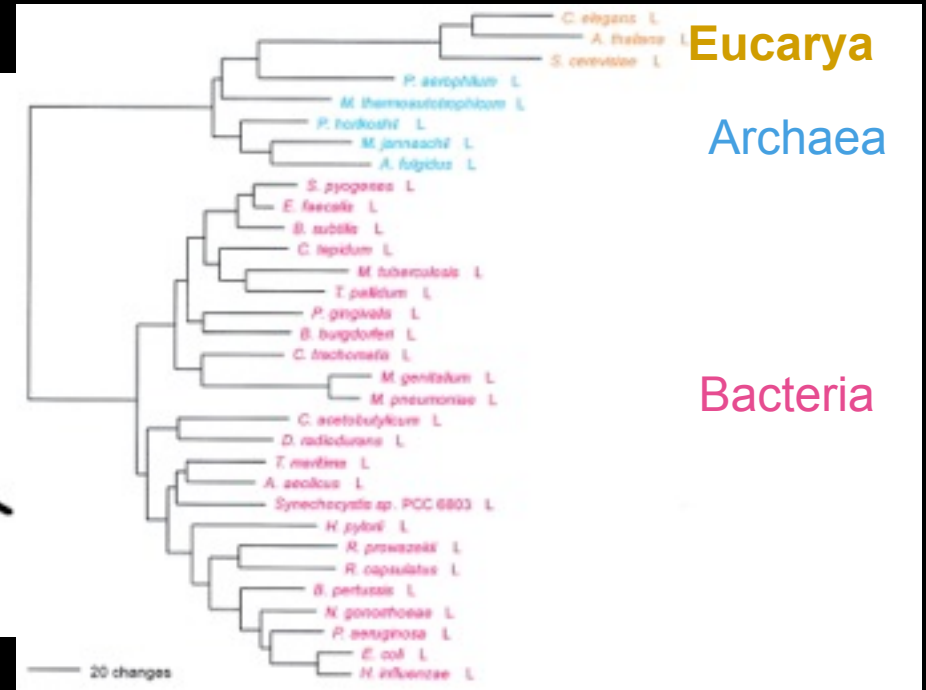
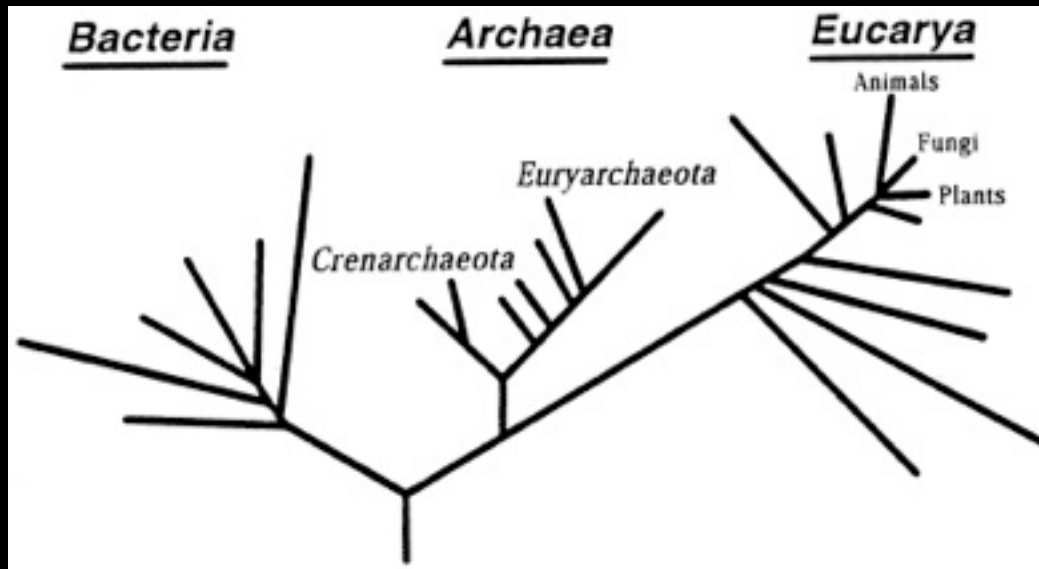


The sequences used to represent the evolutionary history of a biomolecule are the sequences found in a database.



# Setting the Genetic Code: Evolutionary Analysis

The sequence and structure databases are biased towards bacterial domain of life



Universal Phylogenetic Tree

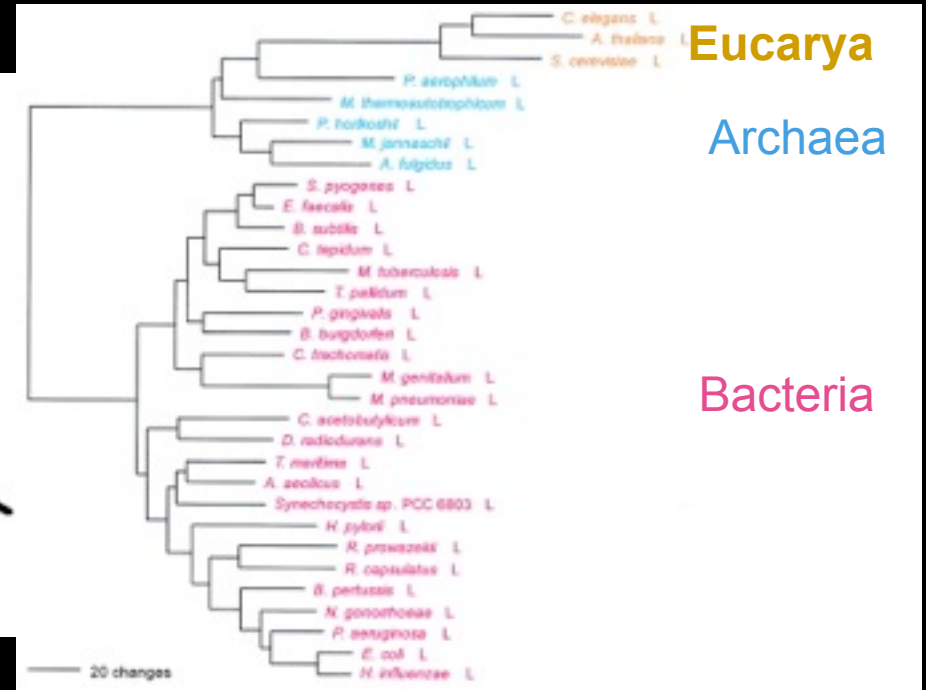
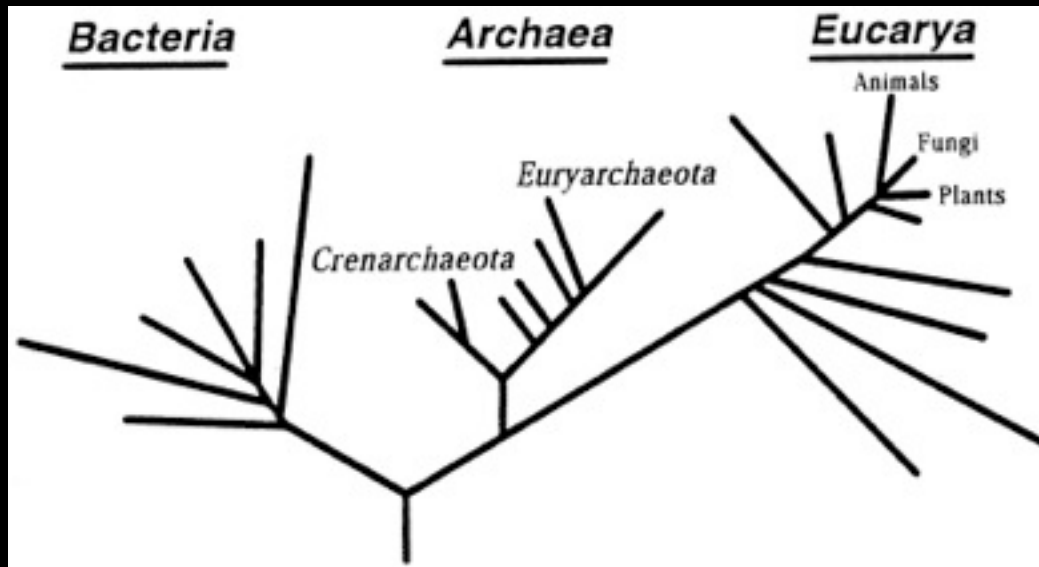
**Three domains of life**

Based on rRNA

Leucyl-tRNA synthetase displays the full canonical phylogenetic distribution.

# Setting the Genetic Code: Evolutionary Analysis

The sequence and structure databases are biased towards bacterial domain of life



Universal Phylogenetic Tree

**Three domains of life**

Based on rRNA

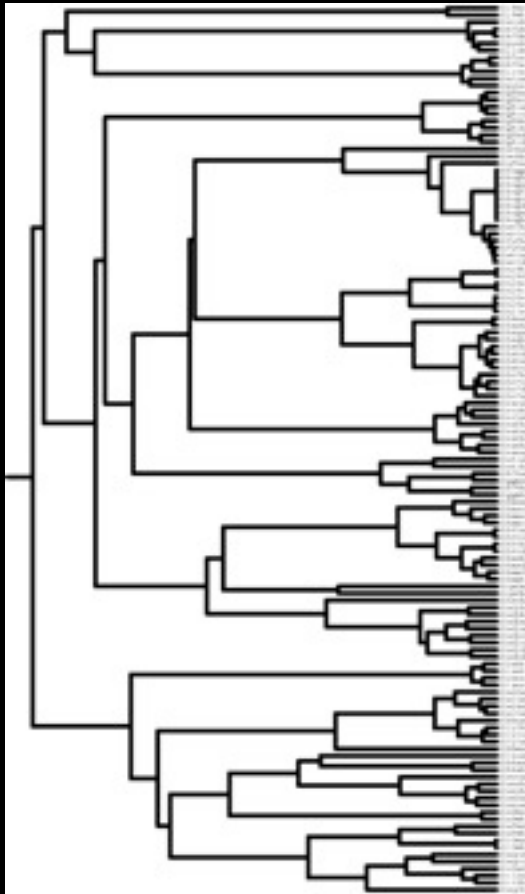
Leucyl-tRNA synthetase displays the full canonical phylogenetic distribution

The databases are biased and statistical analysis of sequence and structure profiles implement ad-hoc sequence weighting methods.

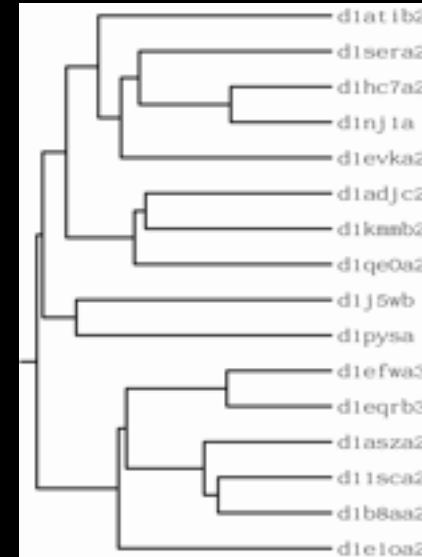
# Setting the Genetic Code: Evolutionary Analysis

Non-redundant sets of sequences and structures can be used for statistical analysis.

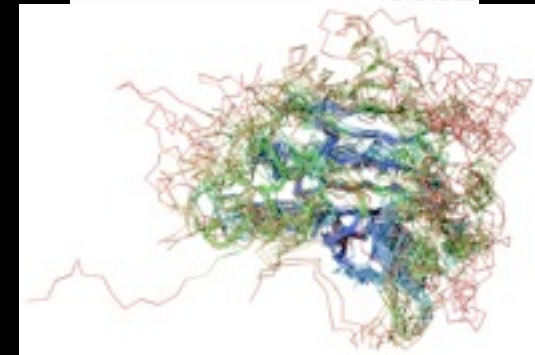
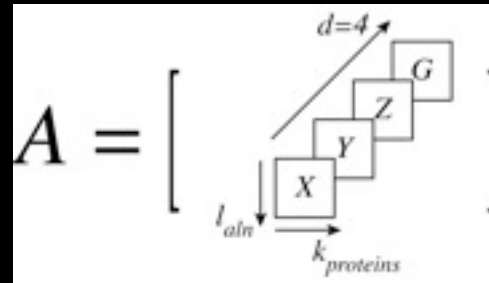
Too much information  
129 Structures



Economy of information  
16 representatives



Multidimensional QR factorization of alignment matrix,  $A$ .





## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS

GlnRS

LysRS

CysRS

## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS }  
GlnRS } Indirect Mechanism  
LysRS  
CysRS

## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS }  
GlnRS } Indirect Mechanism  
LysRS Class I aaRS  
CysRS

## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS }  
GlnRS } Indirect Mechanism  
LysRS Class I aaRS  
CysRS

CysteinyI-tRNA(Cys) formation in *Methanocaldococcus jannaschii*: the mechanism is still unknown. *J. Bacteriology*, Jan. 2004, **186**:8-14. Ruan B, Nakano H, Tanaka M, Mills JA, DeVito JA, Min B, Low KB, Battista JR, and Söll D.



## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS }  
GlnRS } Indirect Mechanism  
LysRS Class I aaRS  
CysRS

CysteinyI-tRNA(Cys) formation in *Methanocaldococcus jannaschii*: the mechanism is still unknown. *J. Bacteriology*, Jan. 2004, **186**:8-14. Ruan B, Nakano H, Tanaka M, Mills JA, DeVito JA, Min B, Low KB, Battista JR, and Söll D.

*M. jannaschii*  
genome database  
search using EP of  
class II aaRS with  
HMMER

## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS }  
GlnRS } Indirect Mechanism  
LysRS Class I aaRS  
CysRS

CysteinyI-tRNA(Cys) formation in *Methanocaldococcus jannaschii*: the mechanism is still unknown. *J. Bacteriology*, Jan. 2004, **186**:8-14. Ruan B, Nakano H, Tanaka M, Mills JA, DeVito JA, Min B, Low KB, Battista JR, and Söll D.

*M. jannaschii*  
genome database  
search using EP of  
class II aaRS with  
HMMER

Protein	E-value
HisRS	1.1e-10
AspRS	1.9e-10
PheRS $\alpha$ -chain	9.5e-10
ThrRS	6.6e-04
ProRS	9.1e-03
SerRS	9.2e-03
putative CysRS	1.6e-02
AlaRS	5.1e-02
GlyRS	0.12
PheRS $\beta$ -chain	0.15
DNA repair protein	7.5

A Sethi, et al., PNAS, 2005.

## Setting the Genetic Code: Evolutionary Analysis

# The story of the missing archaeal CysRS

*M.jannaschii* genome was completely sequenced in 1996.

Genome had four missing aaRSs:

AsnRS } Indirect Mechanism  
GlnRS }  
LysRS Class I aaRS  
CysRS

CysteinyI-tRNA(Cys) formation in *Methanocaldococcus jannaschii*: the mechanism is still unknown. *J. Bacteriology*, Jan. 2004, **186**:8-14. Ruan B, Nakano H, Tanaka M, Mills JA, DeVito JA, Min B, Low KB, Battista JR, and Söll D.

*M. jannaschii*  
genome database  
search using EP of  
class II aaRS with  
HMMER

Protein	E-value
HisRS	1.1e-10
AspRS	1.9e-10
PheRS $\alpha$ -chain	9.5e-10
ThrRS	6.6e-04
ProRS	9.1e-03
SerRS	9.2e-03
putative CysRS	1.6e-02
AlaRS	5.1e-02
GlyRS	0.12
PheRS $\beta$ -chain	0.15
DNA repair protein	7.5

← MJ1660

A Sethi, et al., PNAS, 2005.

# Setting the Genetic Code: Evolutionary Analysis

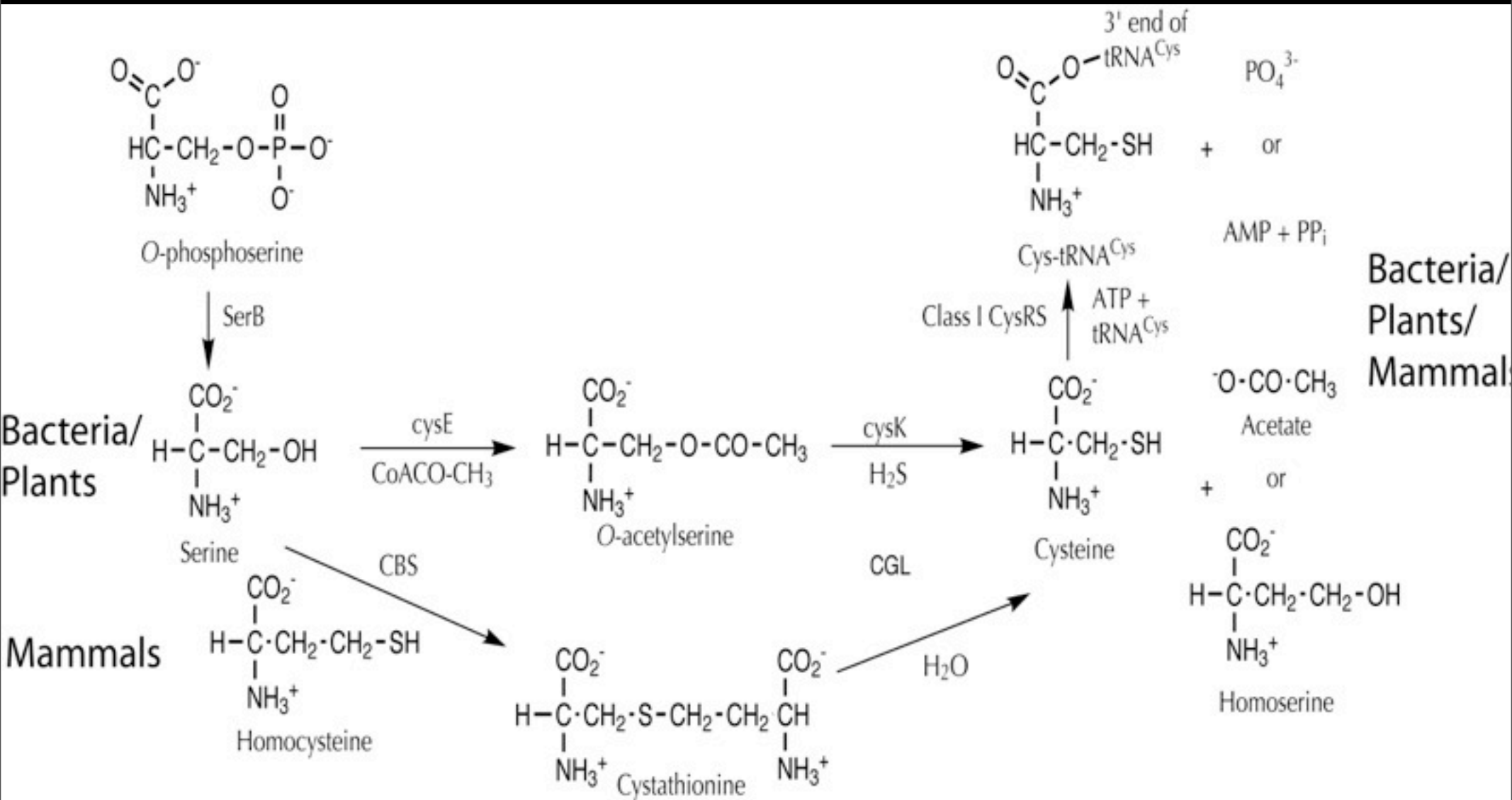
## The Complete Story

Sauerwald, et al., Science, 2005.

# Setting the Genetic Code: Evolutionary Analysis

## The Complete Story

Direct pathway for cysteine aminoacylation



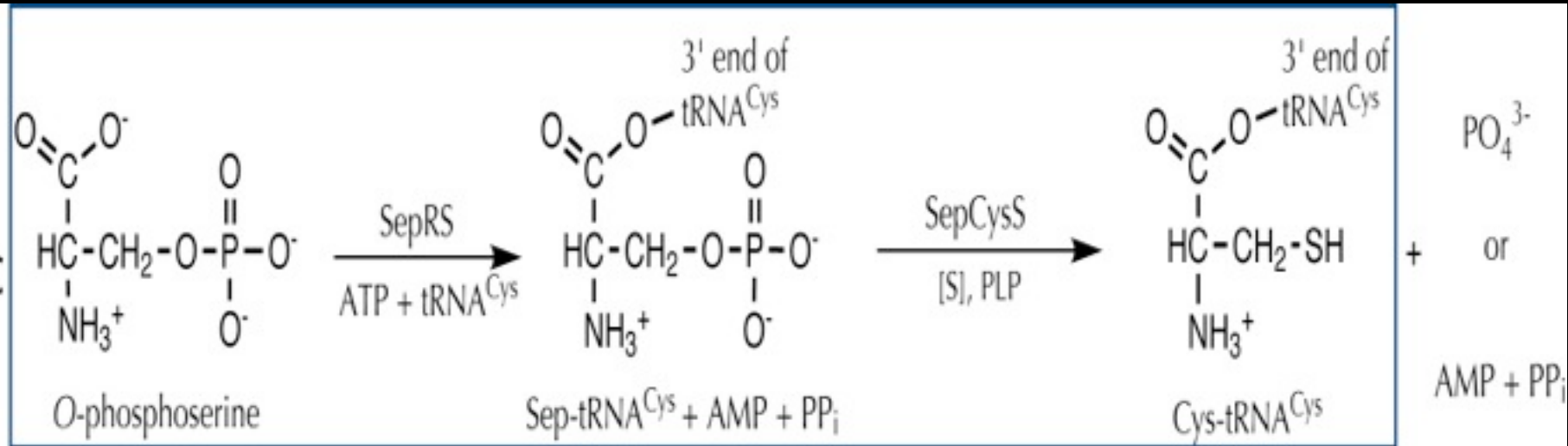
Sauerwald, et al., Science, 2005.

# Setting the Genetic Code: Evolutionary Analysis

## The Complete Story

Indirect pathway for cysteine aminoacylation

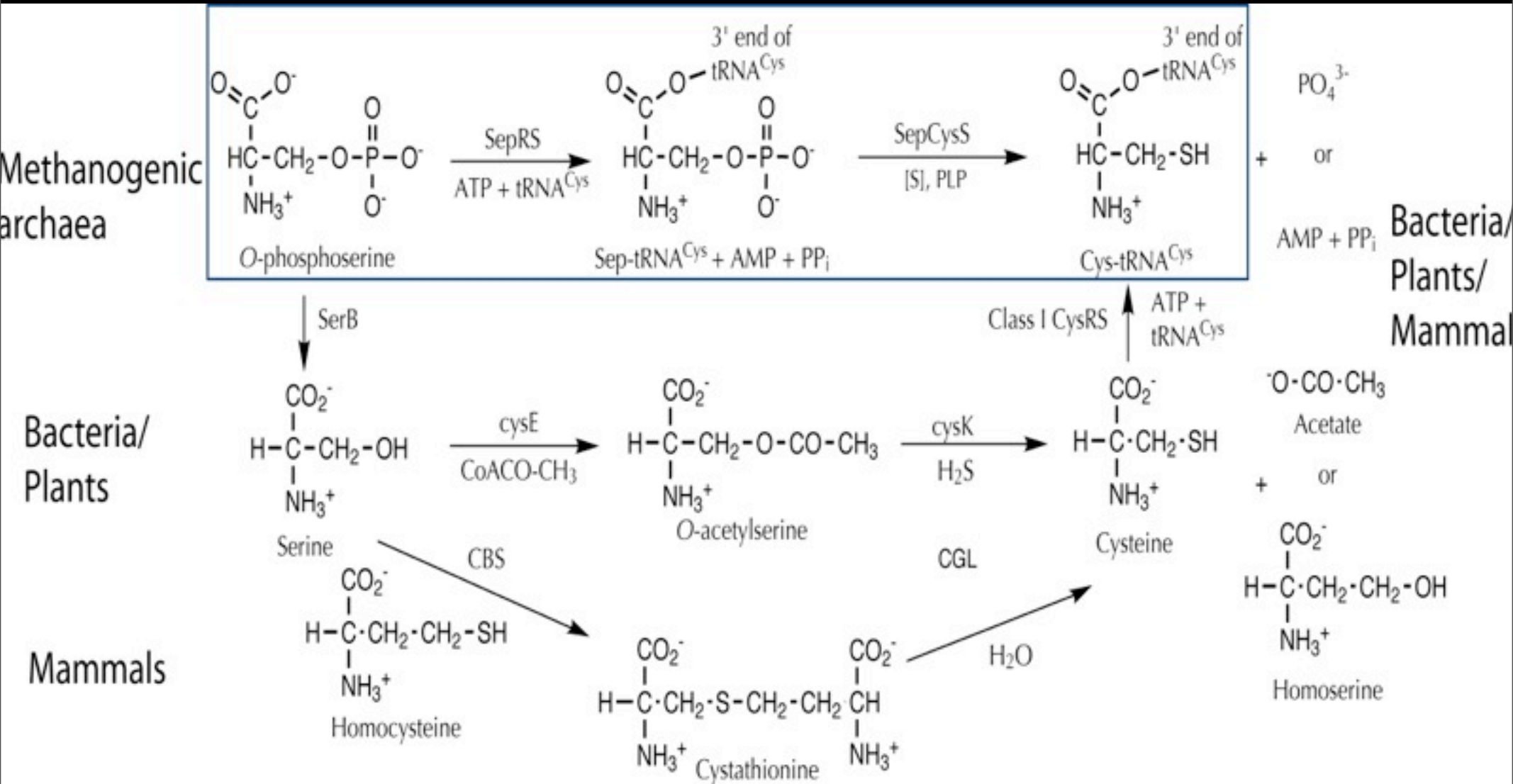
Methanogenic  
archaea



Sauerwald, et al., Science, 2005.

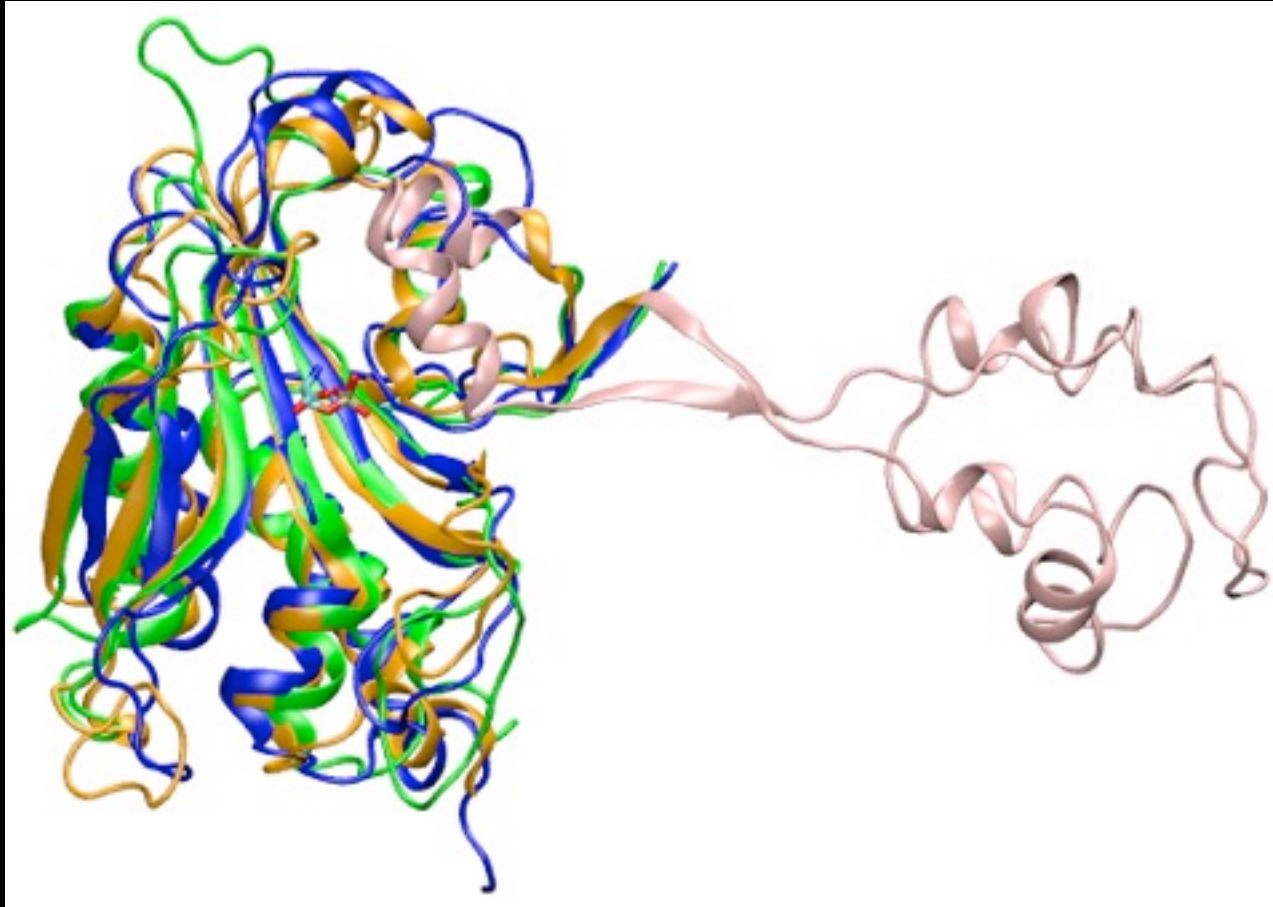
# Setting the Genetic Code: Evolutionary Analysis

## The Complete Story



Sauerwald, et al., Science, 2005.

# Success from Structure Prediction



RMSD = 2.72 Å

Orientation of O-phosphoserine in SepRS is different from that of all other amino acid substrates in the catalytic site of class II aaRS.

Fukunaga and Yokoyama, *Nature Str. Mol. Biol.*, 2007.

Sethi, et al., *PNAS*, 2005.



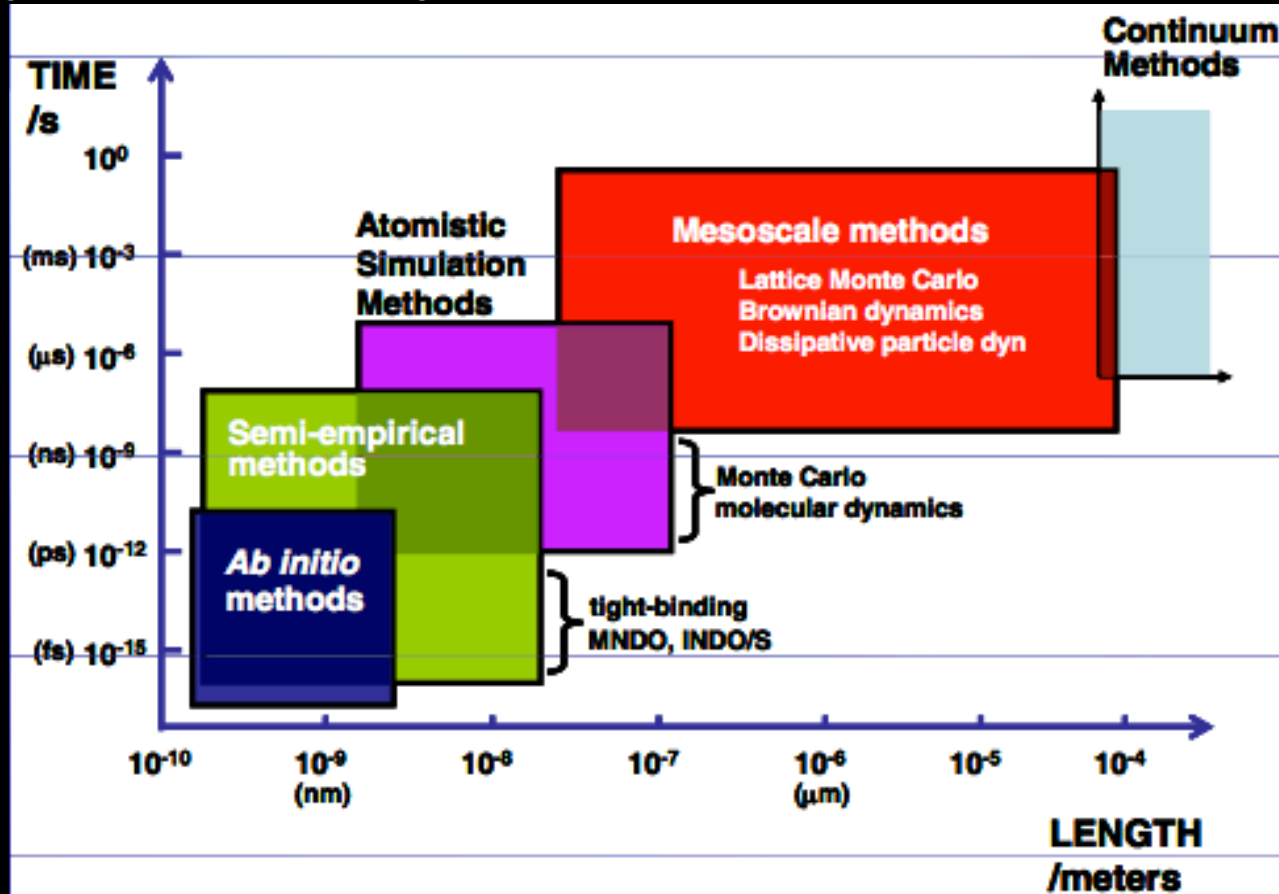
SeqQR was used to study all the steps of translation.

Evolution gives valuable clues to understand how complex systems function.

However, conservation can be due to a variety of reasons and teasing out the details requires physical models.

# Setting the Genetic Code: Signaling Within Biomolecules

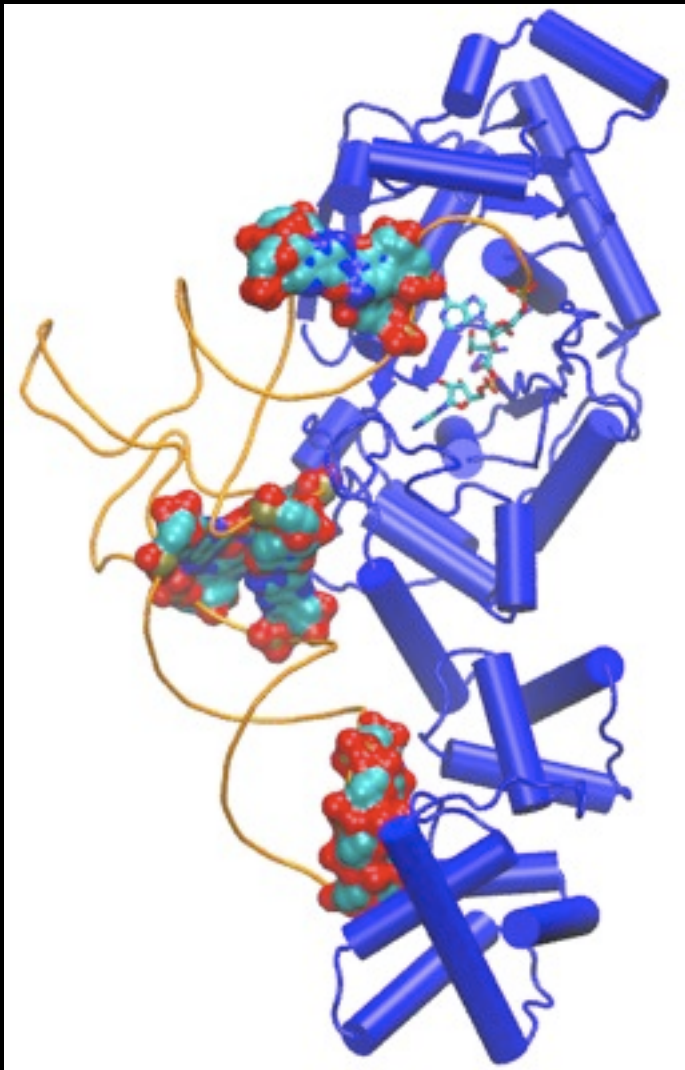
Molecular dynamics simulations can be used to analyze the dynamics within biomolecules



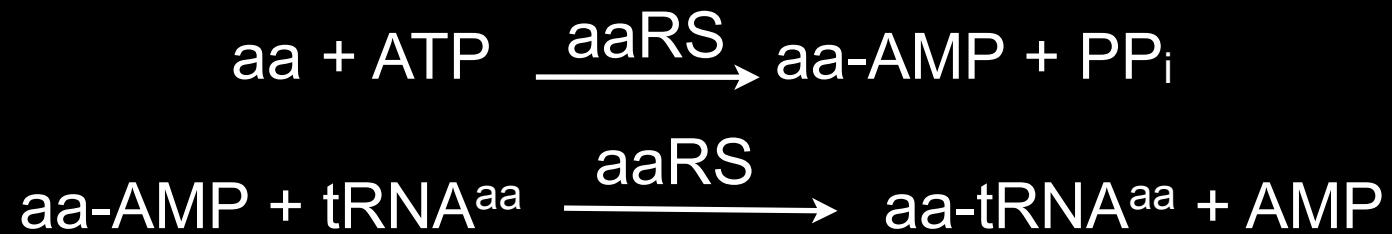
$$\begin{aligned}
 U(\vec{R}) = & \sum_{\text{bonds}} K_b (b - b_0)^2 + \sum_{\text{UB}} K_{UB} (S - S_0)^2 + \sum_{\text{angle}} K_\theta (\theta - \theta_0)^2 \\
 & + \sum_{\text{dihedrals}} K_\chi (1 + \cos(n\chi - \delta)) + \sum_{\text{impropers}} K_{\text{imp}} (\phi - \phi_0)^2 \\
 & + \sum_{\text{nonbond}} \epsilon \left[ \left( \frac{R_{\text{min}_{ij}}}{r} \right)^{12} - \left( \frac{R_{\text{min}_{ij}}}{r} \right)^6 \right] + \frac{q_i q_j}{\epsilon_1 r_{ij}}
 \end{aligned}$$

## Setting the Genetic Code: Signaling Within Biomolecules

Long range communication is necessary for setting the genetic code



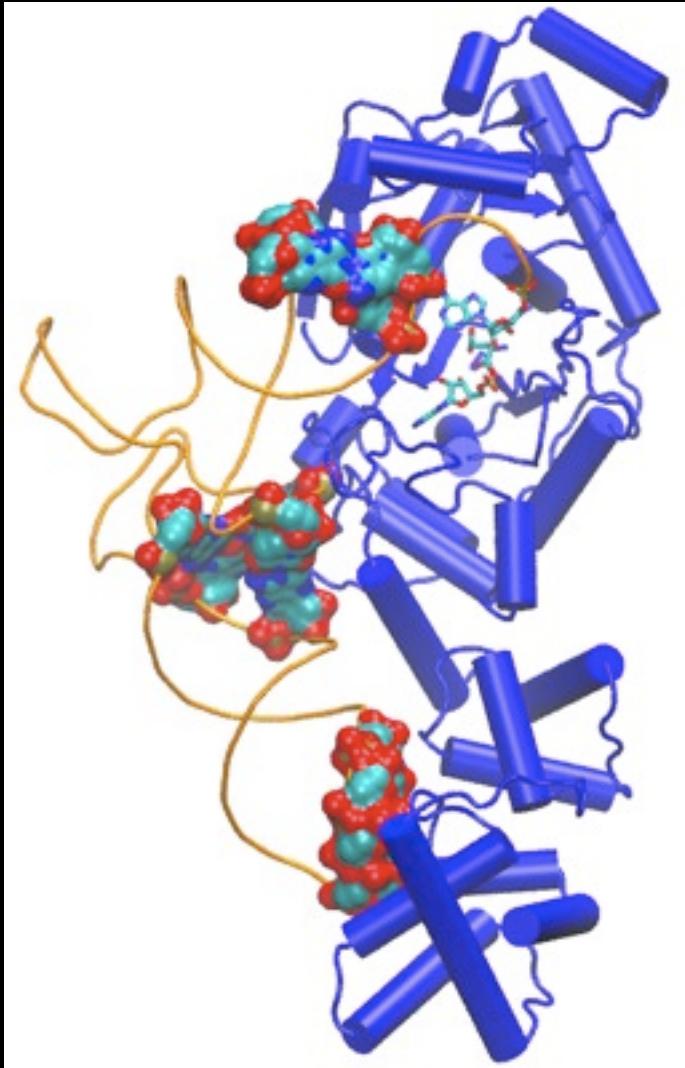
PDB ID: 1N78



Sekine, et. al., JMB, 1996. Sekine, et. al., Eur. J. Biochem, 1999.

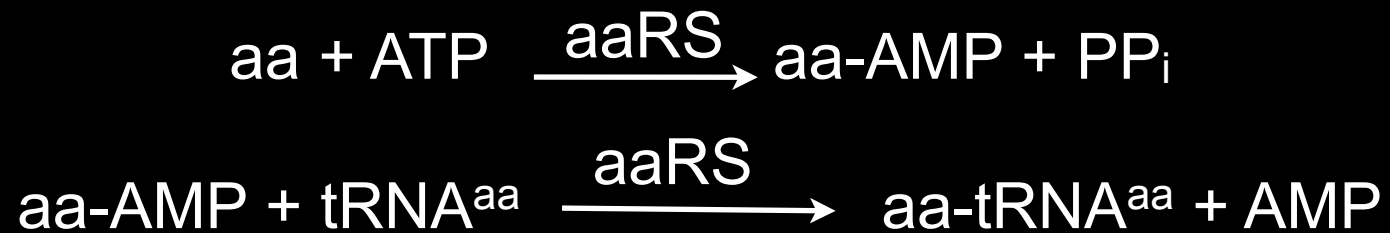
## Setting the Genetic Code: Signaling Within Biomolecules

Long range communication is necessary for setting the genetic code



PDB ID: 1N78

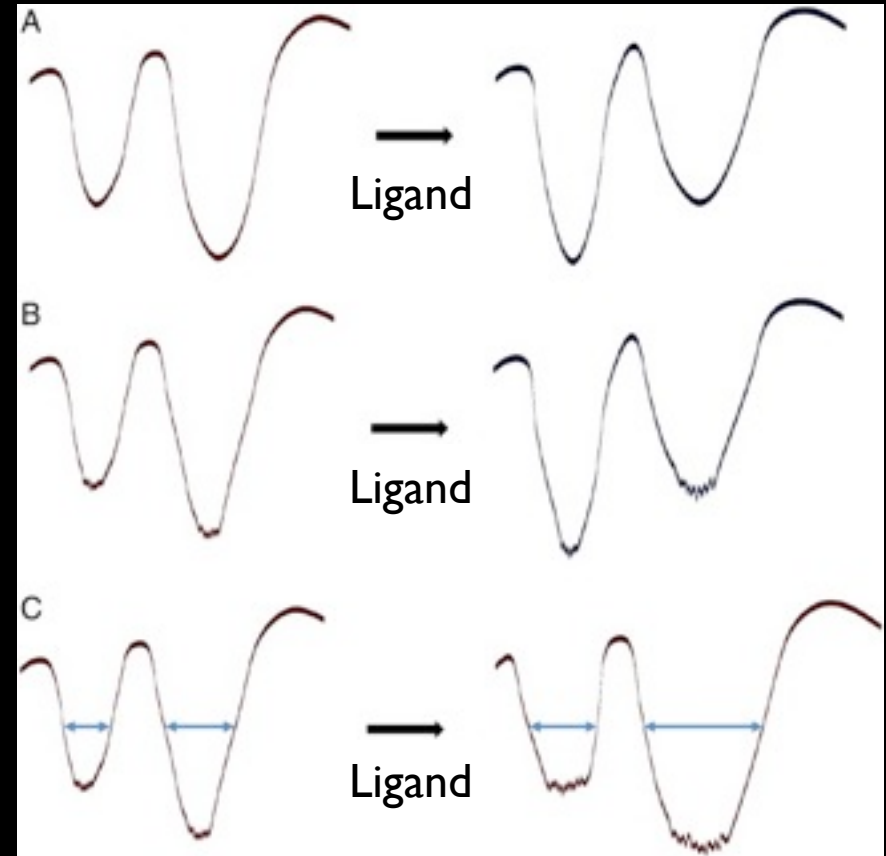
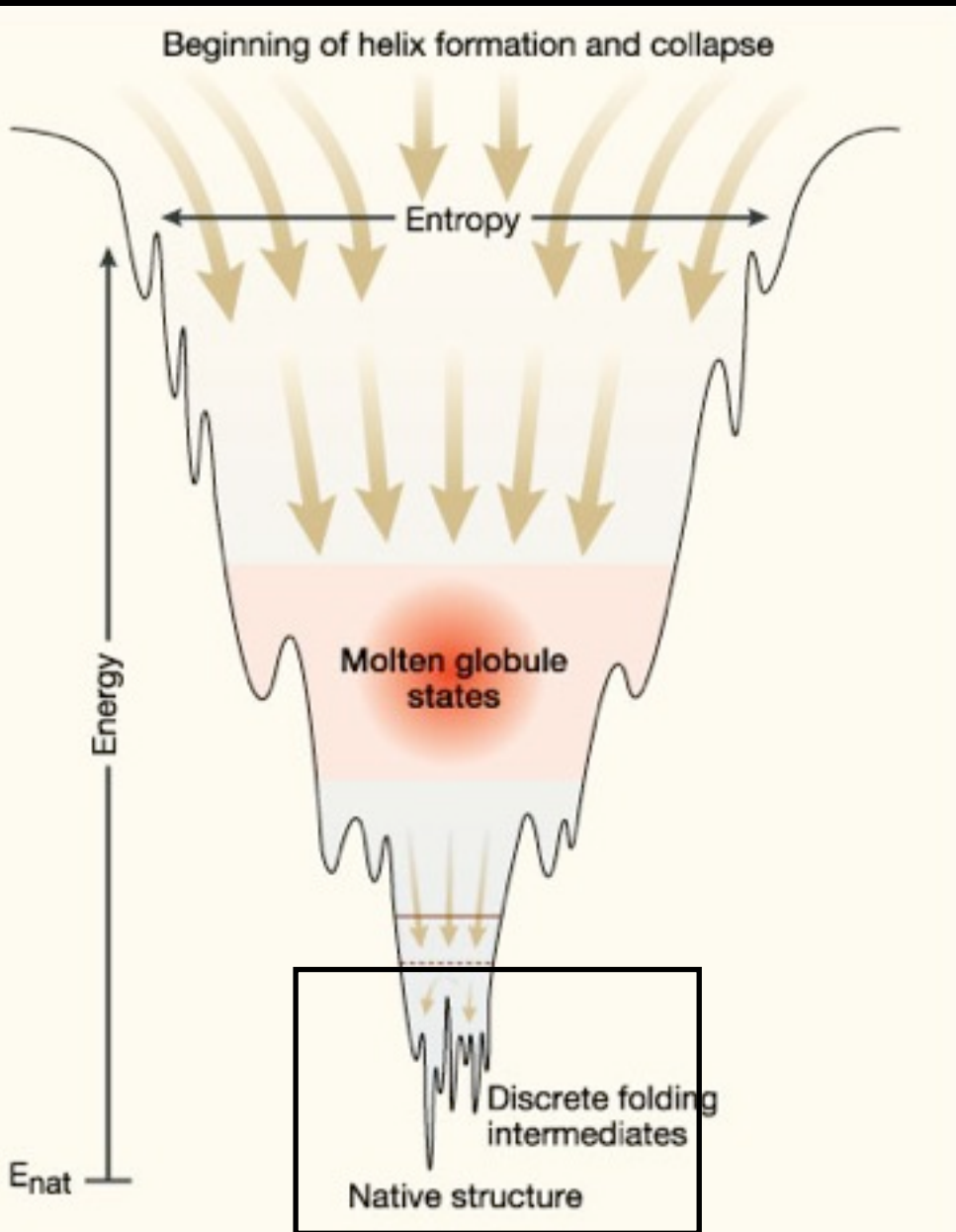
Nucleotides that affect the efficiency of the reaction (catalytic rate) can occur up to 50-70 Angstroms away from the catalytic site.



Sekine, et. al., JMB, 1996. Sekine, et. al., Eur. J. Biochem, 1999.

# Setting the Genetic Code: Signaling Within Biomolecules

## Energy landscape theory explains how macromolecules fold and function

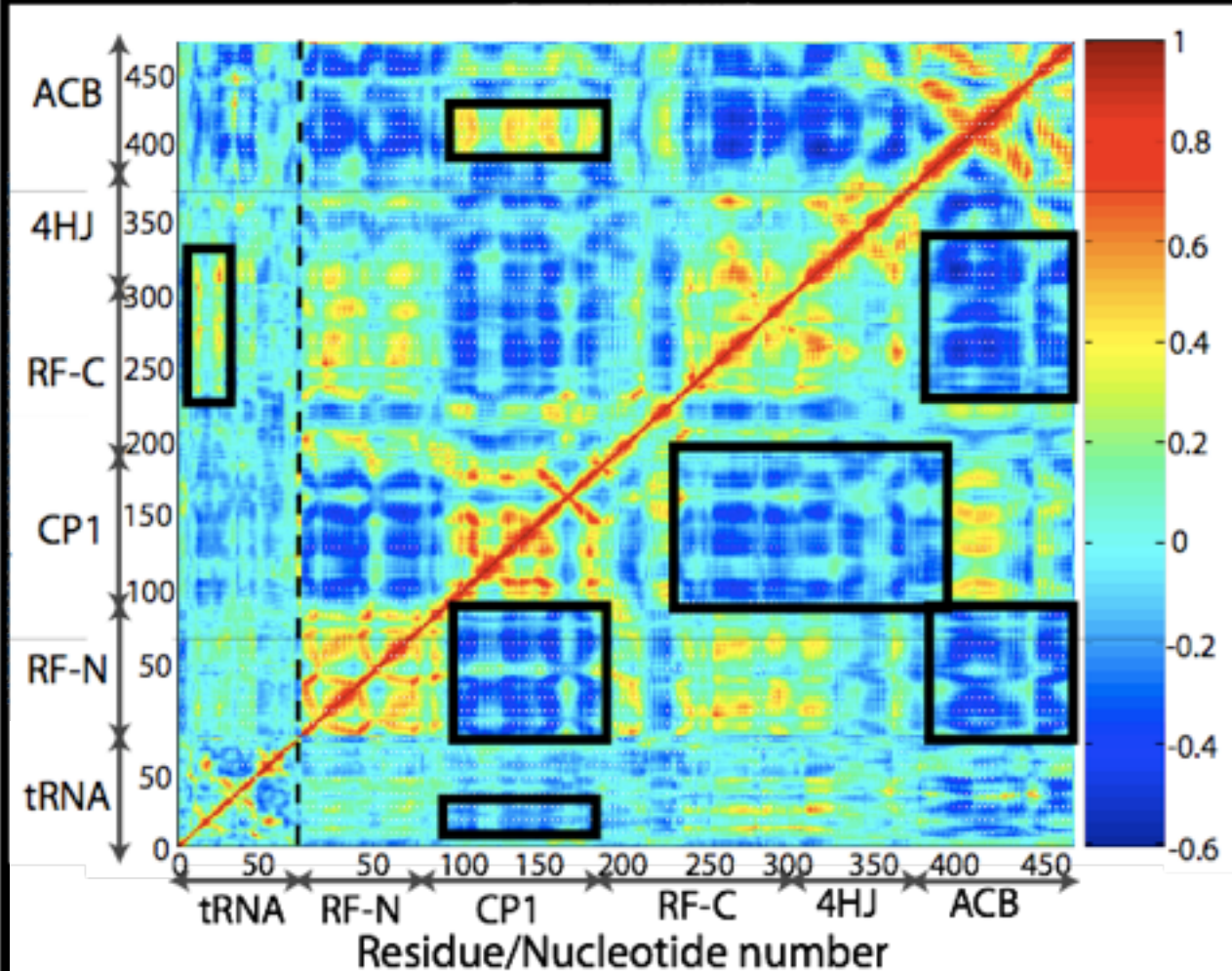


Allostery involves a change in conformation or dynamics. Structural changes might occur through a network of local changes.

Onuchic, et al., Ann Rev Phys Chem, 1997  
Lila Gierarsch, Curr Opin Str Biol, 2006

# Setting the Genetic Code: Signaling Within Biomolecules

## Observation of Allostery in MD Simulations



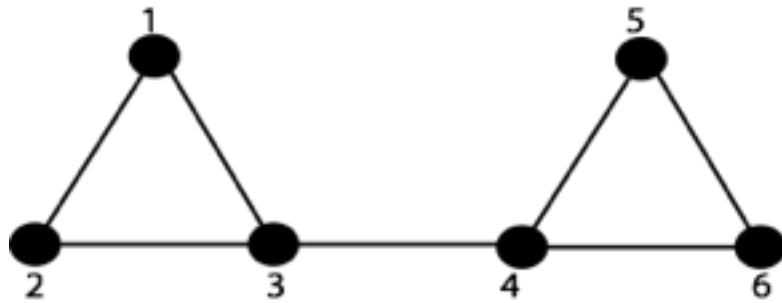
$$C(i, j) = \frac{\langle \Delta r_i(t) \cdot \Delta r_j(t) \rangle}{(\langle \Delta r_i(t)^2 \rangle \langle \Delta r_j(t)^2 \rangle)^{0.5}}$$

Correlation in motion between residues and nucleotides in protein:RNA complex.

A Sethi, et al., PNAS, 2009.

# Setting the Genetic Code: Signaling Within Biomolecules

## Ideas borrowed from network theory



	1	2	3	4	5	6
1		X	X			
2	X		X			
3	X	X		X		
4			X		X	X
5				X		X
6				X	X	

Nodes represent the residues/nucleotides in the complex.

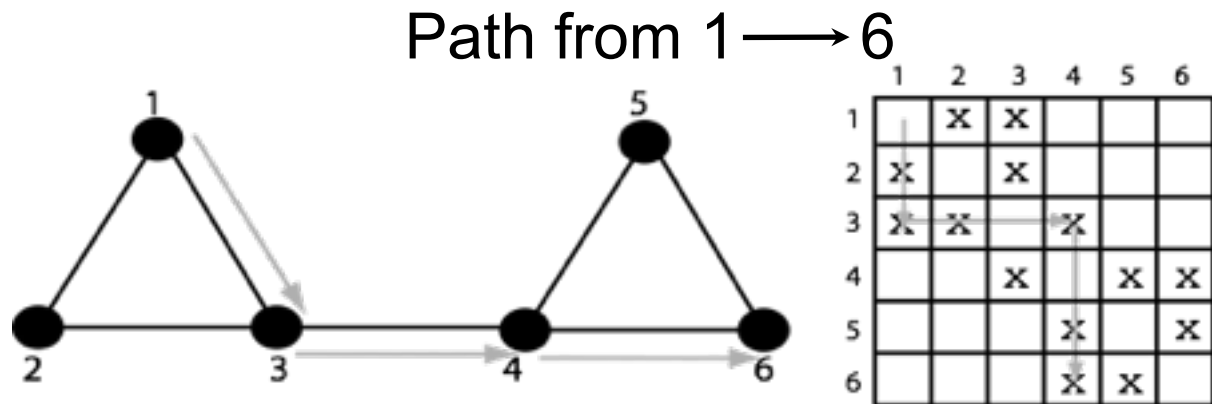
Edges represent contact between monomers in the complex.

The edges can either be unweighted or weighted.  
The edges are weighted by correlation ( $C_{ij}$ )  
between contacts in the simulation:

$$w_{ij} = -\ln(|C_{ij}|)$$

# Setting the Genetic Code: Signaling Within Biomolecules

## Ideas borrowed from network theory



Nodes represent the residues/nucleotides in the complex.

Edges represent contact between monomers in the complex.

The edges can either be unweighted or weighted. The edges are weighted by correlation ( $C_{ij}$ ) between contacts in the simulation:

$$w_{ij} = -\ln(|C_{ij}|)$$

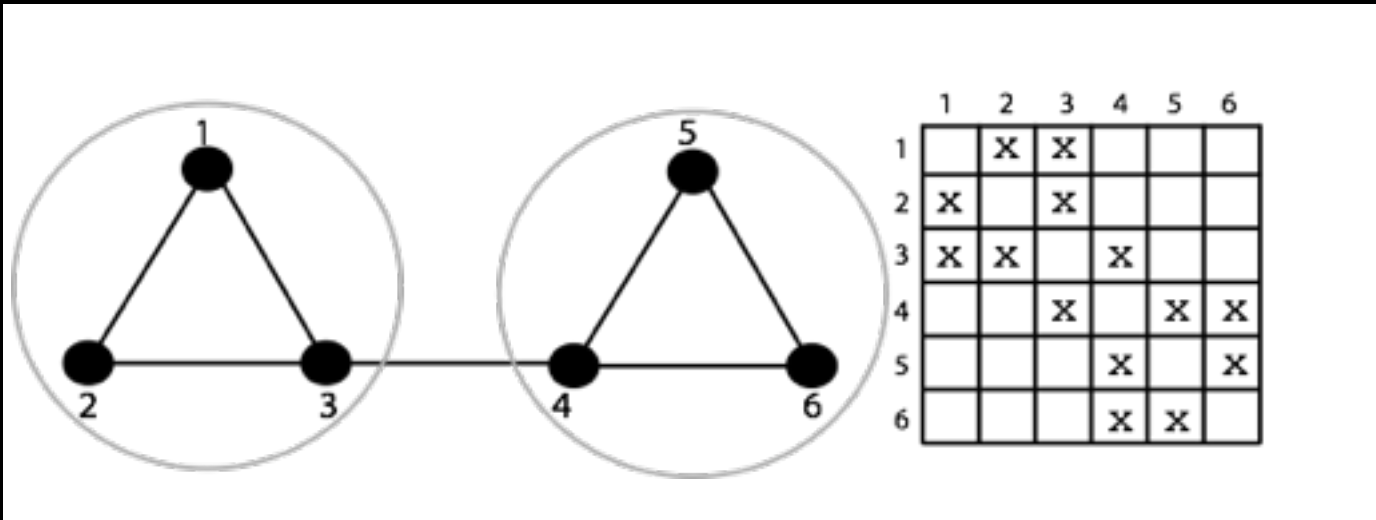
The path distance is the sum of the weights or distance of each edge:

$$D_{ij} = \sum_{k,l} w_{kl}$$



# Setting the Genetic Code: Signaling Within Biomolecules

## Ideas borrowed from network theory



Nodes represent the residues/nucleotides in the complex.

Edges represent contact between monomers in the complex.

The edges can either be unweighted or weighted. The edges are weighted by correlation ( $C_{ij}$ ) between contacts in the simulation:

$$w_{ij} = -\ln(|C_{ij}|)$$

The path distance is the sum of the weights or distance of each edge:

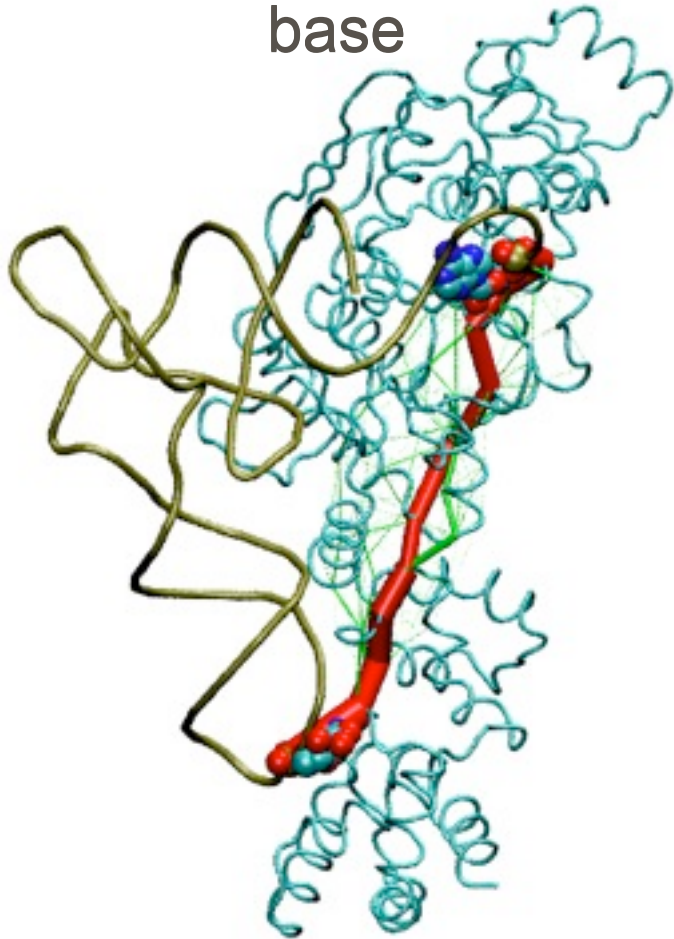
$$D_{ij} = \sum_{k,l} w_{kl}$$

Modules have fewer connections between them

# Setting the Genetic Code: Signaling Within Biomolecules

## The Importance of Suboptimal Paths

Suboptimal paths from anticodon  
base



There are a number of suboptimal paths for communication between identity elements and active site.

A Sethi, et al., PNAS, 2009.

# Setting the Genetic Code: Signaling Within Biomolecules

## The Importance of Suboptimal Paths

Suboptimal paths from anticodon  
base



There are a number of suboptimal paths for communication between identity elements and active site.

A Sethi, et al., PNAS, 2009.

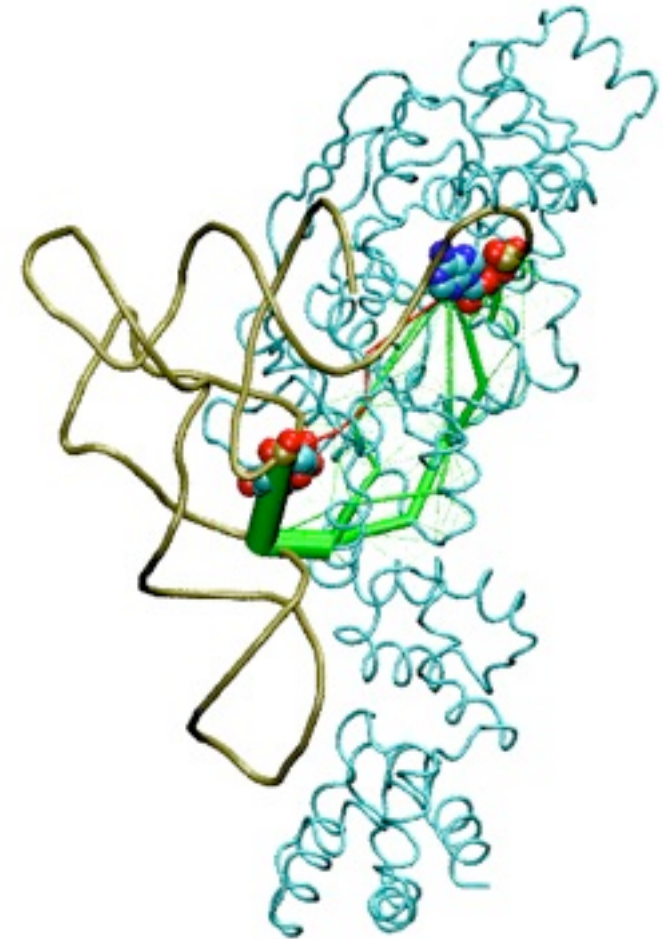
# Setting the Genetic Code: Signaling Within Biomolecules

## The Importance of Suboptimal Paths

Suboptimal paths from anticodon base



Suboptimal paths from Ura11



There are a number of suboptimal paths for communication between identity elements and active site.

A Sethi, et al., PNAS, 2009.

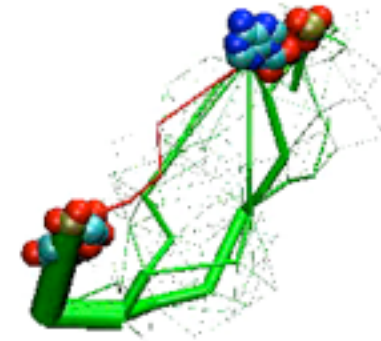
# Setting the Genetic Code: Signaling Within Biomolecules

## The Importance of Suboptimal Paths

Suboptimal paths from anticodon base



Suboptimal paths from Ura11



There are a number of suboptimal paths for communication between identity elements and active site.

A Sethi, et al., PNAS, 2009.

## Setting the Genetic Code: Signaling Within Biomolecules

Regions connecting modules form hotspots for communication in the network

Communities are modules within the network that move in a correlated fashion during the MD simulation.

Residues connecting modules are critical for communication in the biomolecular network

They are conserved in evolution

They affect network properties

They occur in majority of suboptimal pathways



## Setting the Genetic Code: Signaling Within Biomolecules

Regions connecting modules form hotspots for communication in the network

Communities are modules within the network that move in a correlated fashion during the MD simulation.

Residues connecting modules are critical for communication in the biomolecular network

They are conserved in evolution

They affect network properties

They occur in majority of suboptimal pathways

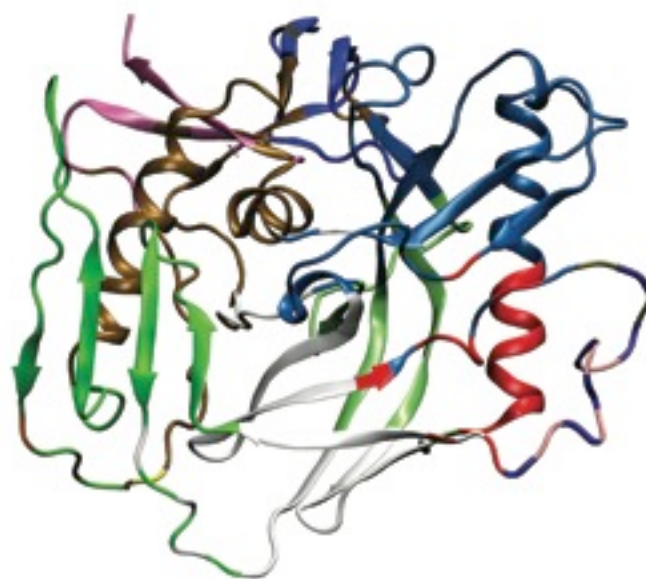


## Signaling Within Biomolecules

The modules in the protein are more highly conserved



CAP210



YU2



HXBC2

The communities in the network are highly conserved between different sequences of gp120.

The intermodular contacts are under high immune pressure to evolve.



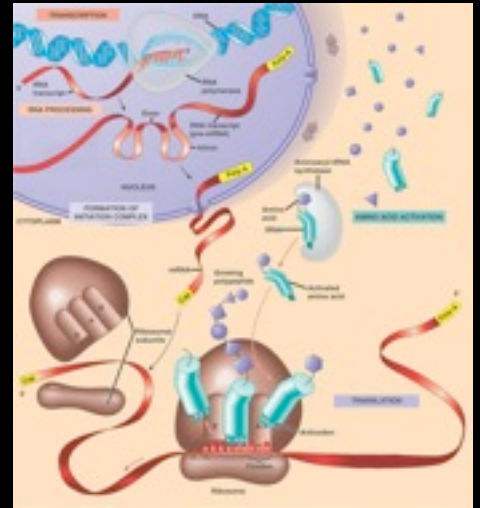
# Conclusions

- Balanced evolutionary profiles provide an economy of information that can be used for gene annotation.
- Evolutionary profiles were successful at identifying the protein responsible for cysteine aminoacylation in methanogens.
- The suboptimal paths should be considered while studying communication between distant sites in biomolecular complexes.
- The residues involved in communication between modules in the dynamical network are highly conserved and form hot spots for communication in biomolecular complexes.

# Organization of Talk

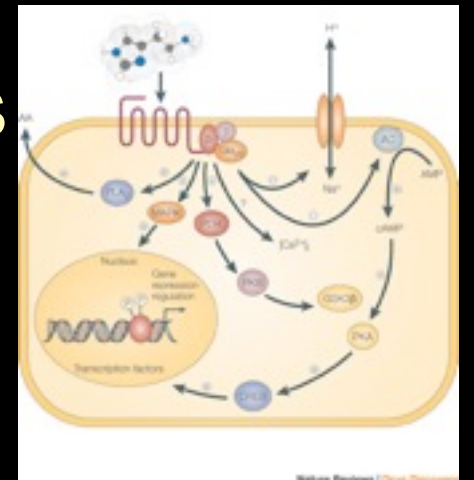
## Translating Information in RNA

Evolutionary analysis of biomolecules  
Allosteric signaling pathways

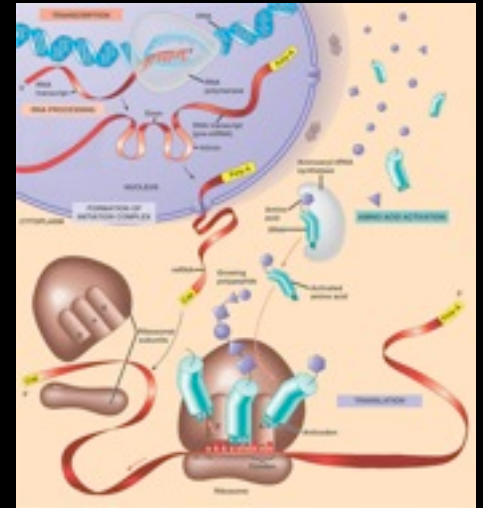


## Disordered regions in Signaling Proteins

Multivalent Proteins  
Modeling Disordered Proteins



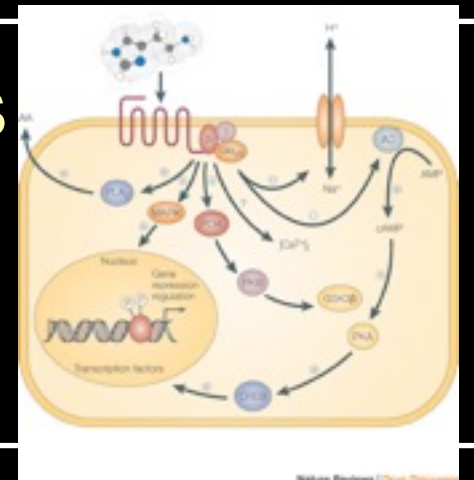
# Organization of Talk



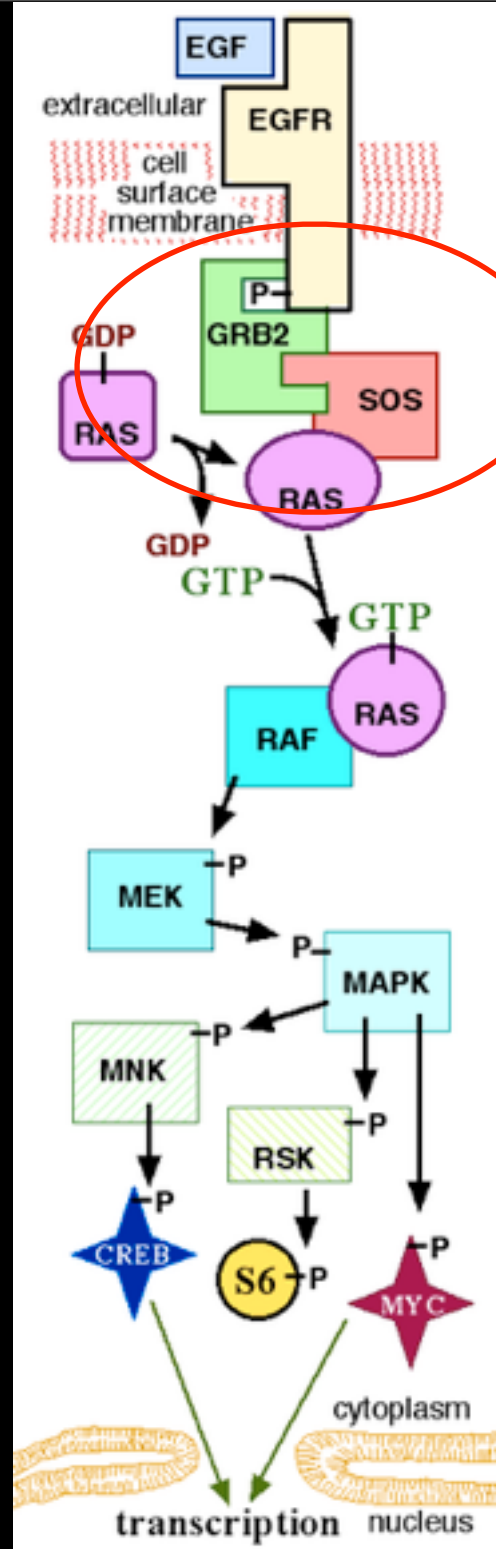
## Disordered regions in Signaling Proteins

Multivalent Proteins

Modeling Disordered Proteins



Part 2:  
Structurally  
Modeling  
Signaling  
Cascades



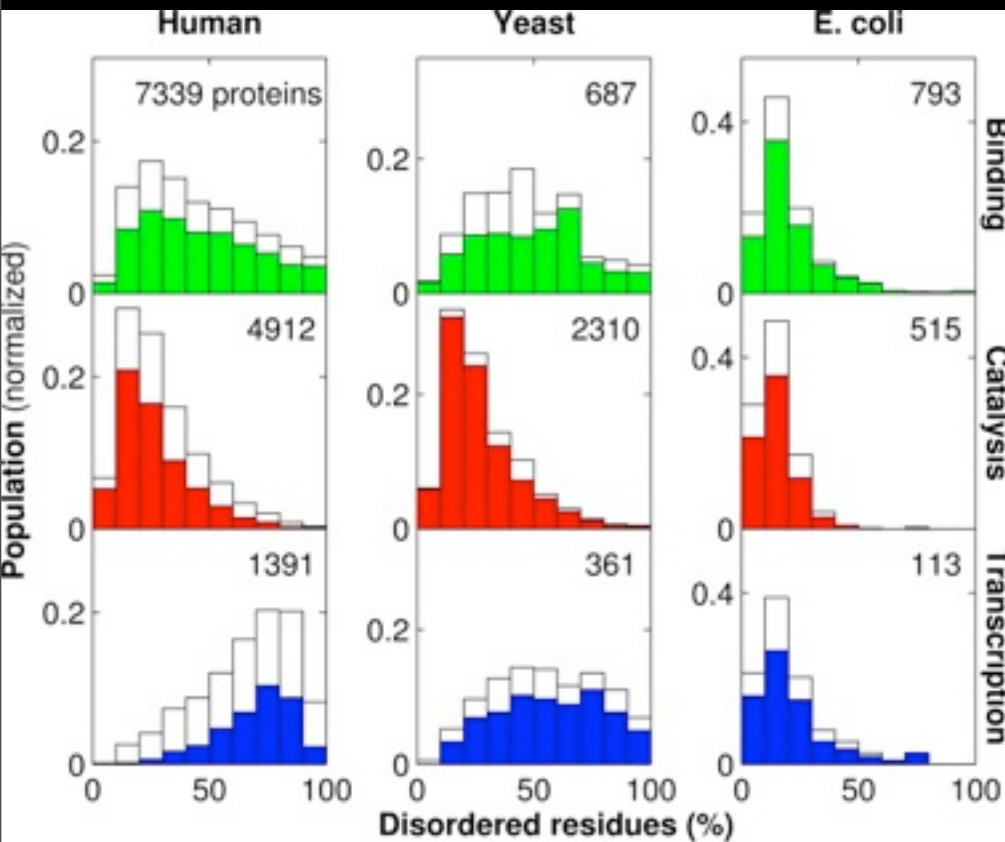
[http://en.wikipedia.org/wiki/Cell\\_signaling](http://en.wikipedia.org/wiki/Cell_signaling)

Signaling cascades  
regulate information  
transfer inside the  
cell.

# Eukaryotes have a significant proportion of their proteins that are disordered.

Classical interpretation:

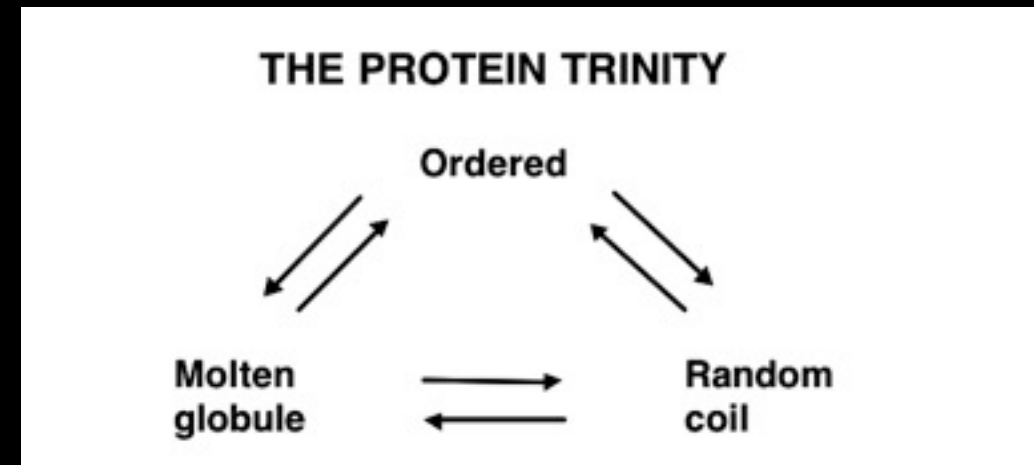
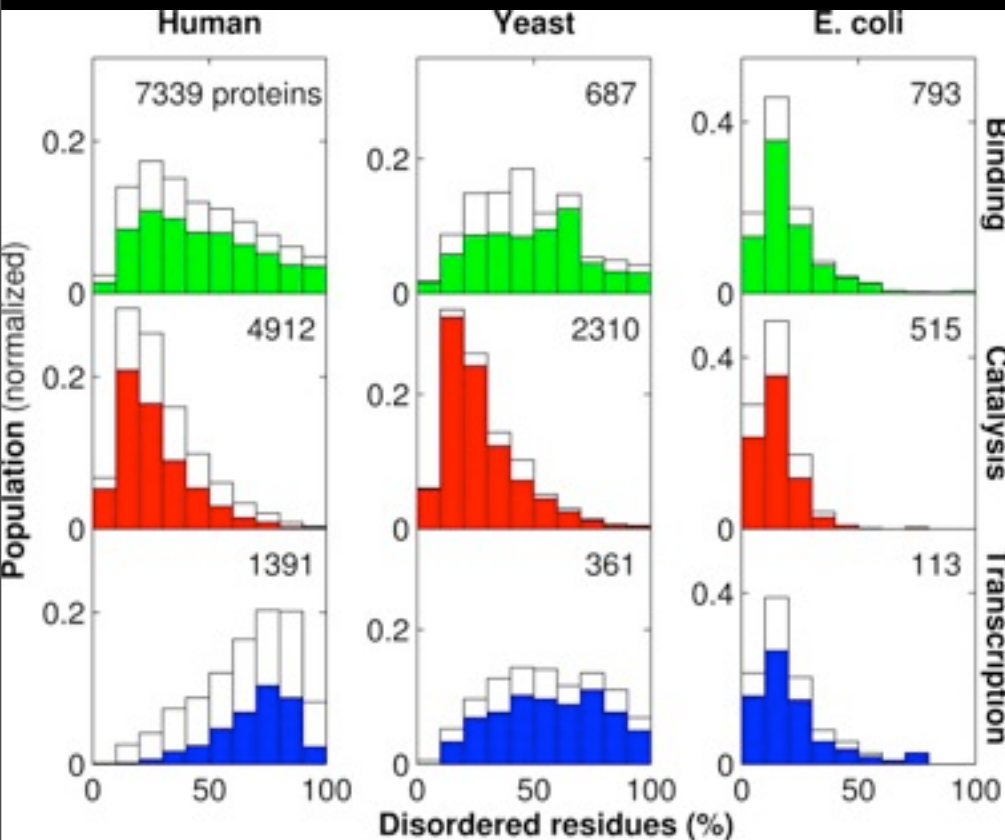
Protein's ordered structure is related to its function.



# Eukaryotes have a significant proportion of their proteins that are disordered.

Classical interpretation:

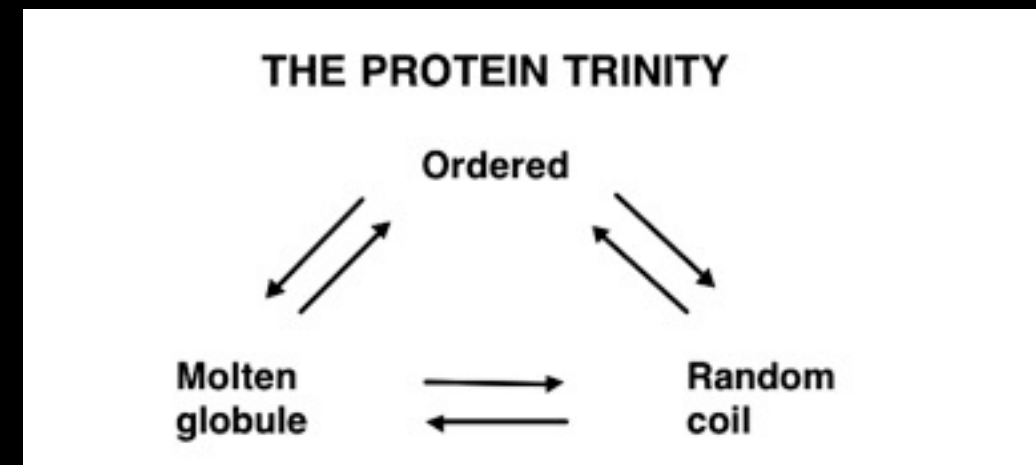
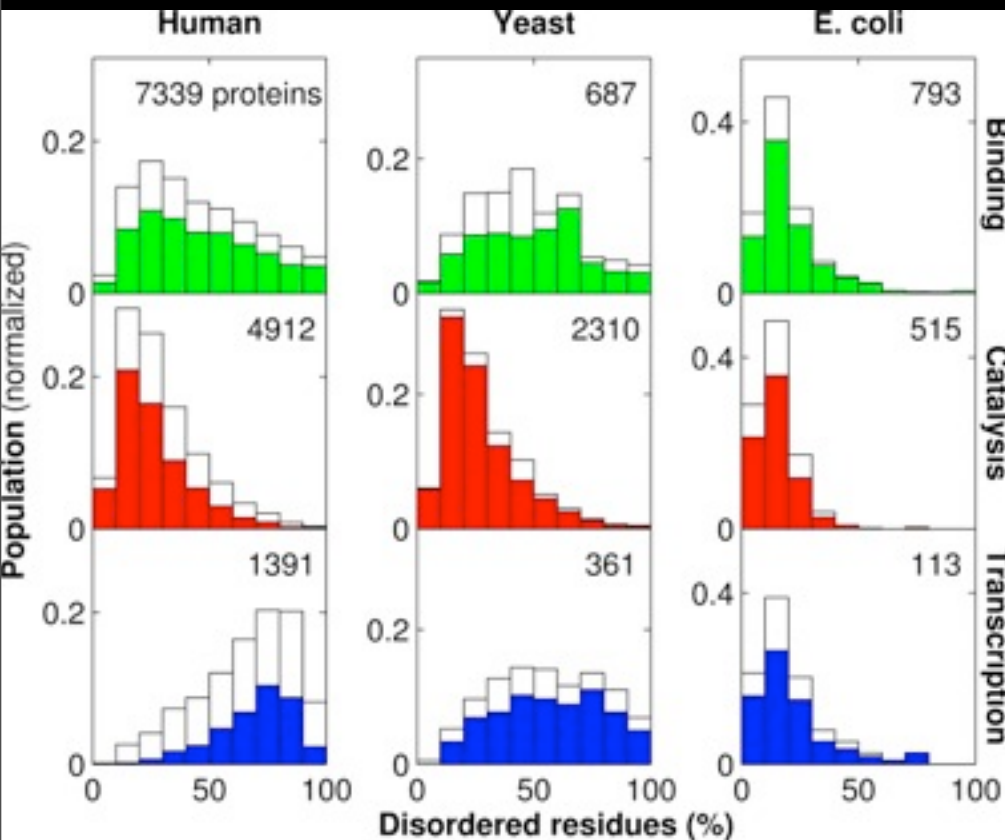
Protein's ordered structure is related to its function.



# Eukaryotes have a significant proportion of their proteins that are disordered.

Classical interpretation:

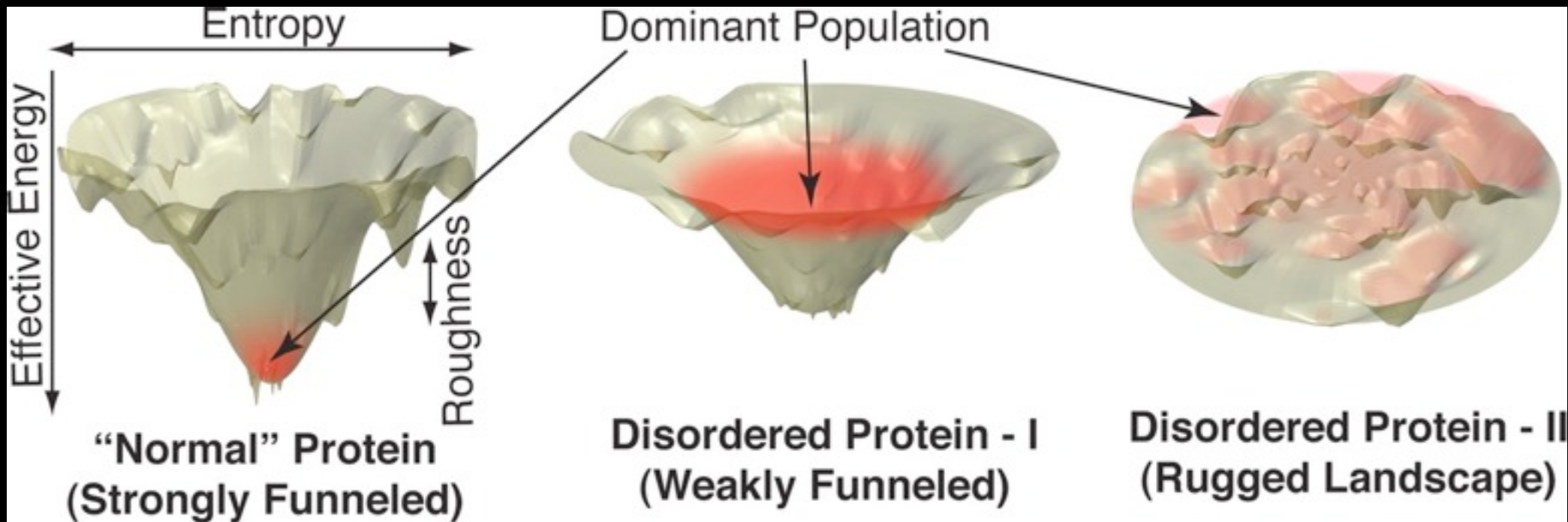
Protein's ordered structure is related to its function.



Folding upon binding  
Fuzzy complexes  
Entropic chains

## Signaling Cascades: Disordered Regions

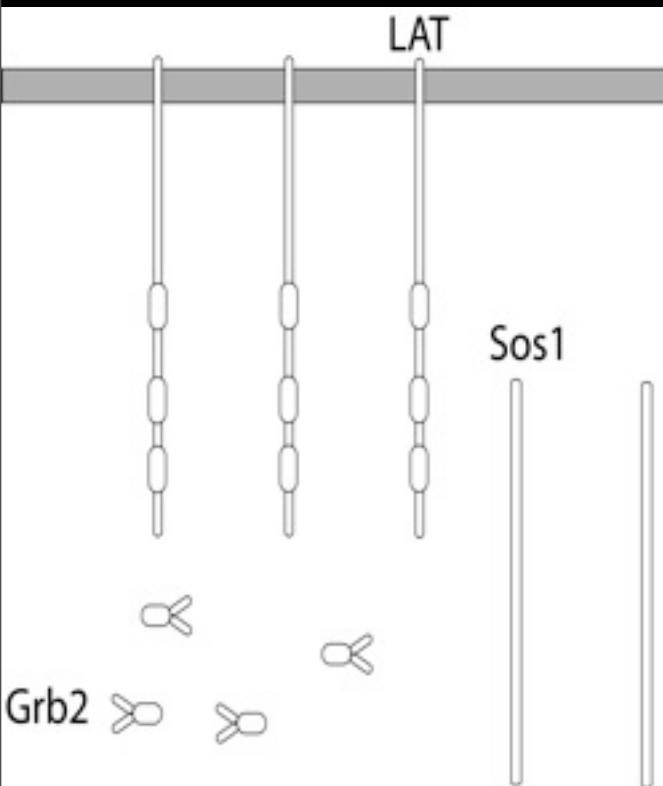
Disordered proteins are considered to either have a weakly funneled or a rugged energy landscape





# Signaling Cascades: Quantifying Multivalent Binding

## Signaling proteins utilize multivalent interactions

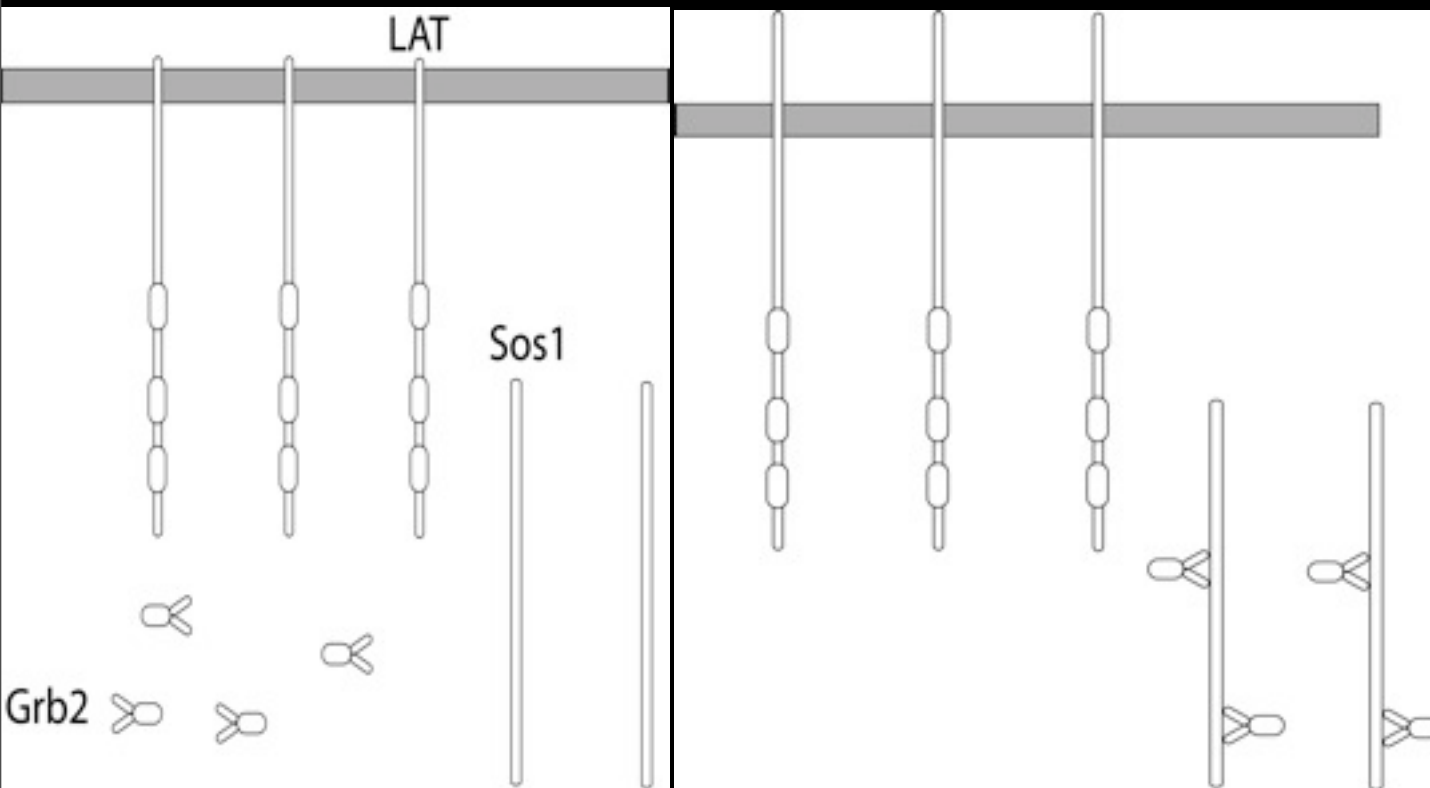


The nucleotide exchange factor Sos1 has to be localized near the plasma membrane.

Houtman, et al., NSB, 2006.  
Nag, et al., Biophys J., 2009.

# Signaling Cascades: Quantifying Multivalent Binding

## Signaling proteins utilize multivalent interactions



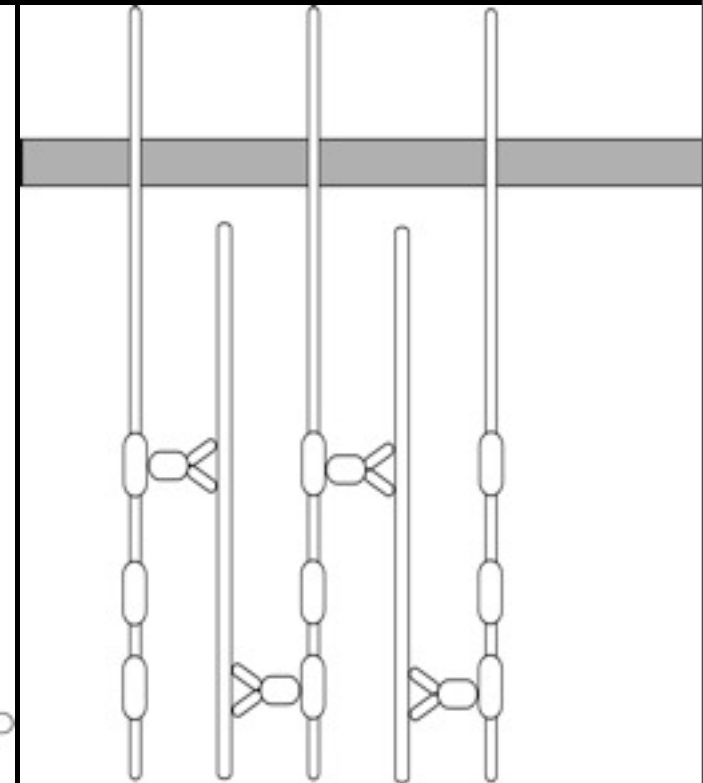
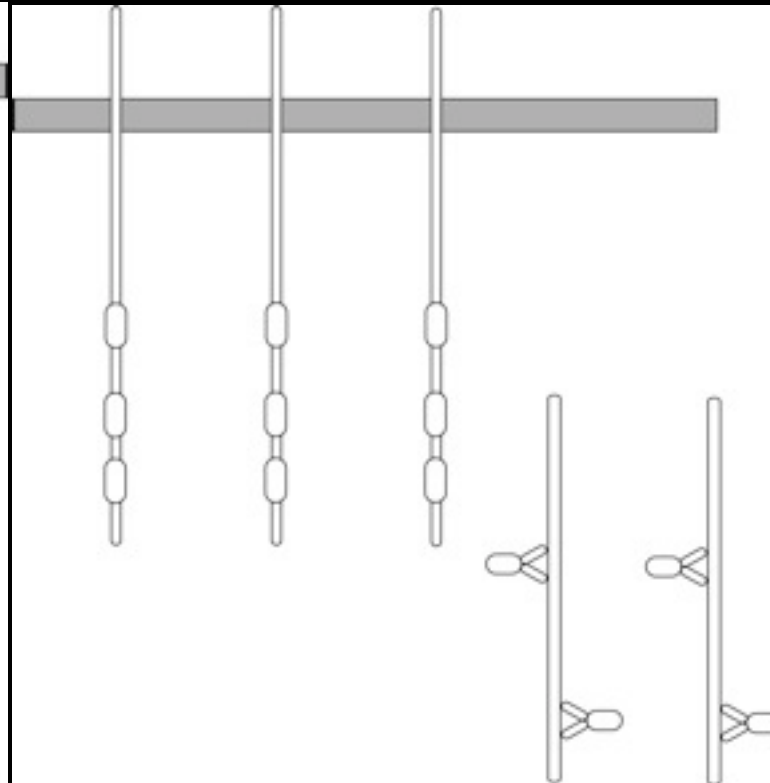
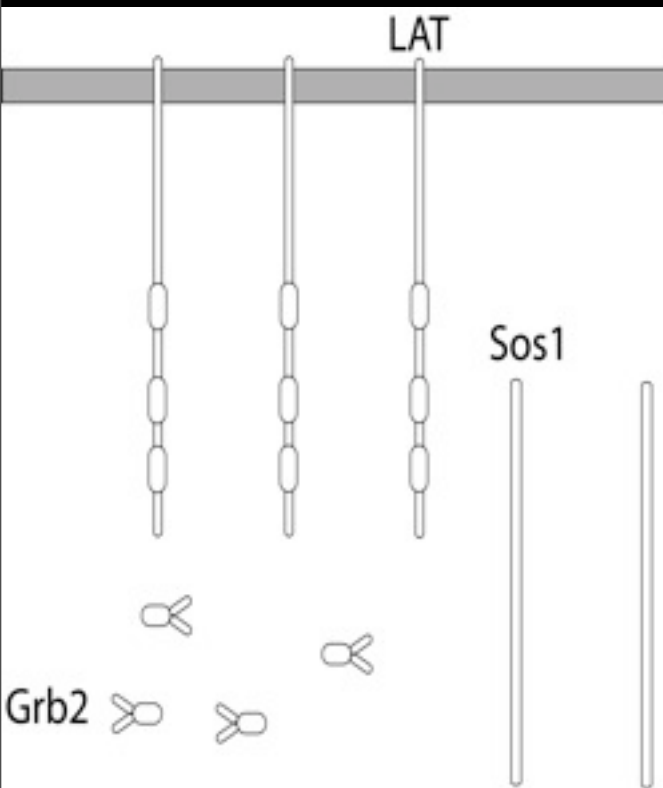
The nucleotide exchange factor Sos1 has to be localized near the plasma membrane.

Grb2 consists of two SH3 domains that interact with Sos1.

Houtman, et al., NSB, 2006.  
Nag, et al., Biophys J., 2009.

# Signaling Cascades: Quantifying Multivalent Binding

## Signaling proteins utilize multivalent interactions



The nucleotide exchange factor Sos1 has to be localized near the plasma membrane.

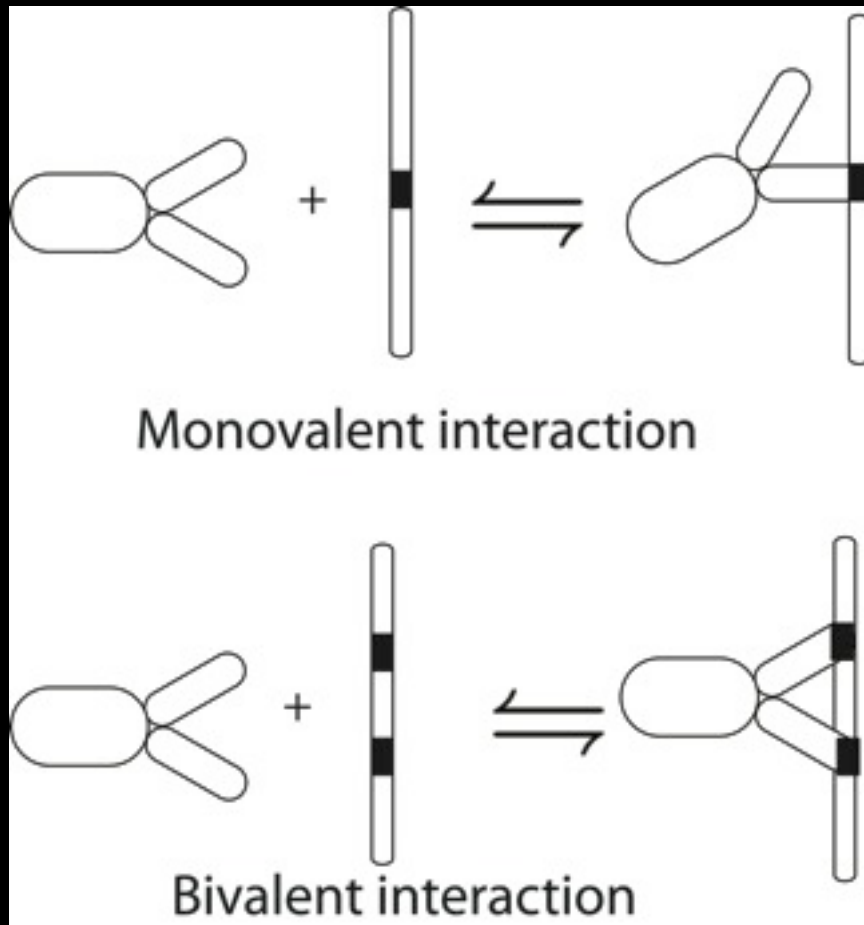
Grb2 consists of two SH3 domains that interact with Sos1.

Grb2 consists of a SH2 domain that interacts with LAT.

Houtman, et al., NSB, 2006.  
Nag, et al., Biophys J., 2009.

## Signaling Cascades: Quantifying Multivalent Binding

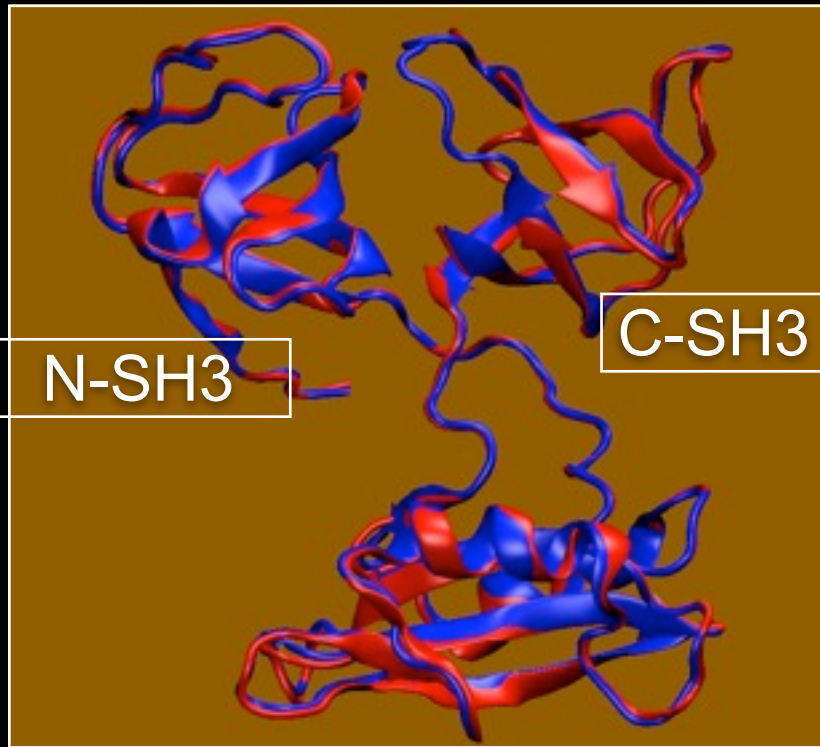
Multivalent interactions are ubiquitous in biology.



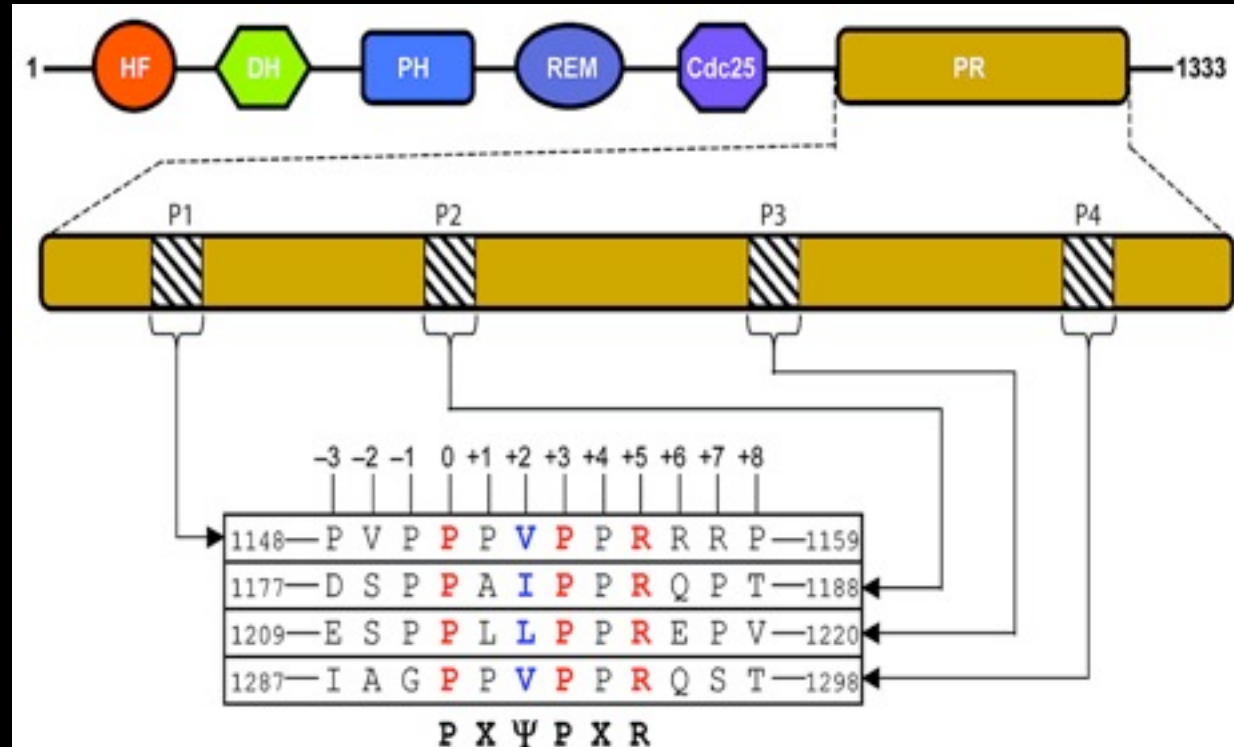
Signaling proteins often use multiple binding sites to increase the overall strength and specificity of the complexes formed.

# Signaling Cascades: Quantifying Multivalent Binding

## Controversy on stoichiometry of complexes formed under physiological conditions.



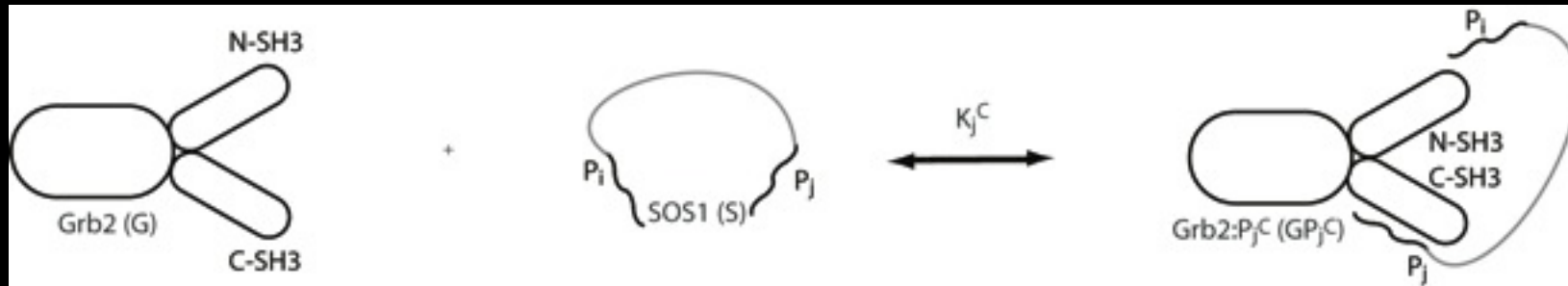
Grb2



Sos1

# Signaling Cascades: Quantifying Multivalent Binding

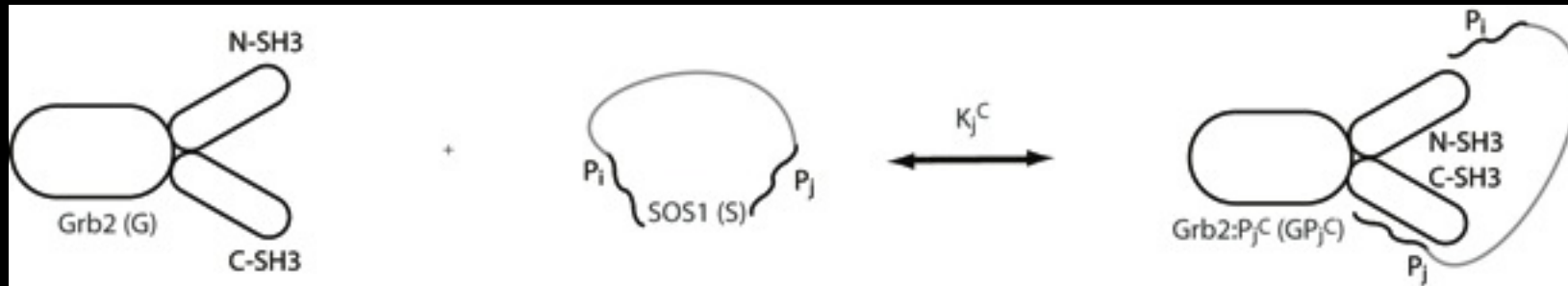
Simultaneous binding of both SH3 domains to two motifs in Sos1



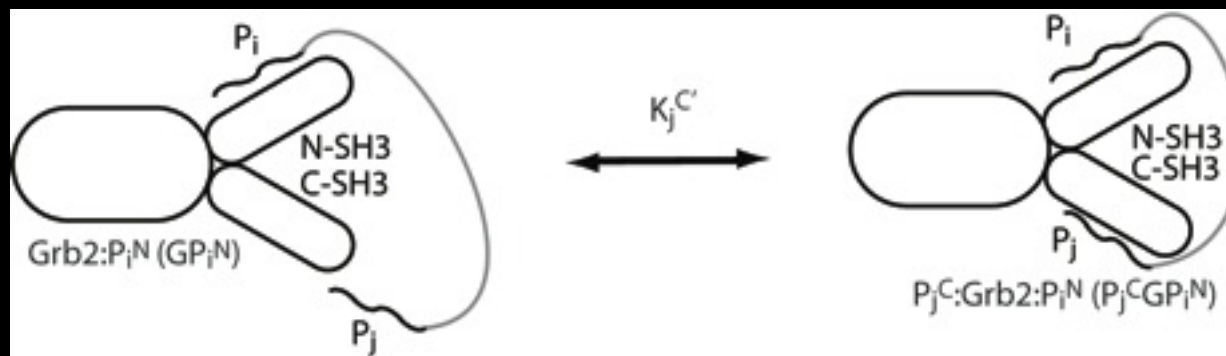
$$[GP_j^C] = K_j^C [G] [P_j]$$

# Signaling Cascades: Quantifying Multivalent Binding

Simultaneous binding of both SH3 domains to two motifs in Sos1



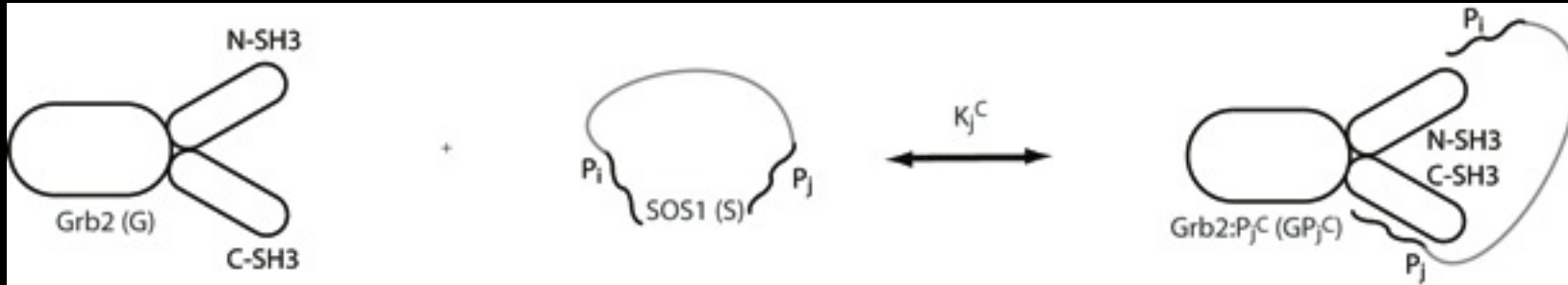
$$[GP_j^C] = K_j^C [G] [P_j]$$



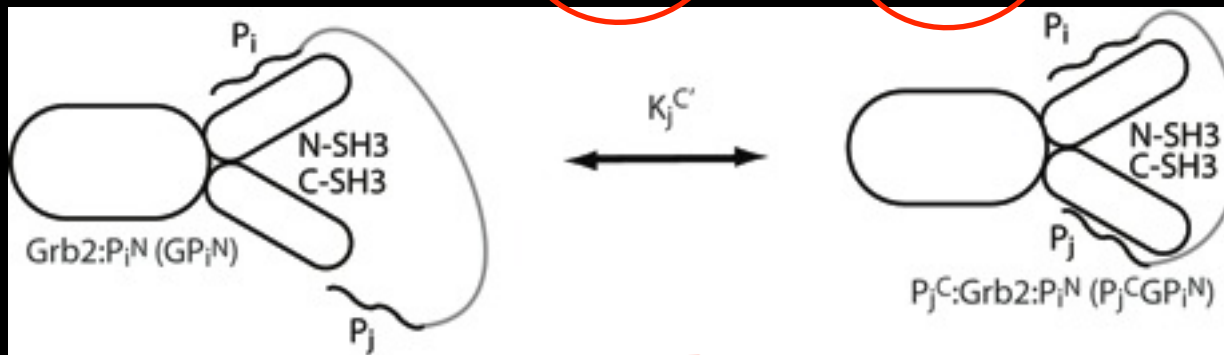
$$[P_j^C GP_i^N] = K_j^{C'} [GP_j^N]$$

# Signaling Cascades: Quantifying Multivalent Binding

Simultaneous binding of both SH3 domains to two motifs in Sos1



$$[GP_j^C] = K_j^C [G] [P_j]$$



$$[P_j^C GP_i^N] = K_j^{C'} [GP_i^N]$$

$$K_j^{C'} = K_j^C \times C_{eff}(P_i^N, P_j^C)$$



# Signaling Cascades: Quantifying Multivalent Binding

## Modeling multivalent interactions

Motifs of Sos1 that bind to Grb2

- Evolutionary analysis
- Binding Energy Calculations

# Signaling Cascades: Quantifying Multivalent Binding

## Modeling multivalent interactions

### Motifs of Sos1 that bind to Grb2

- Evolutionary analysis
- Binding Energy Calculations

### Simultaneous binding of two motifs in Sos1 to Grb2

- Hybrid model from polymer theory and MD Simulations.

# Signaling Cascades: Quantifying Multivalent Binding

## Modeling multivalent interactions

### Motifs of Sos1 that bind to Grb2

- Evolutionary analysis
- Binding Energy Calculations

### Simultaneous binding of two motifs in Sos1 to Grb2

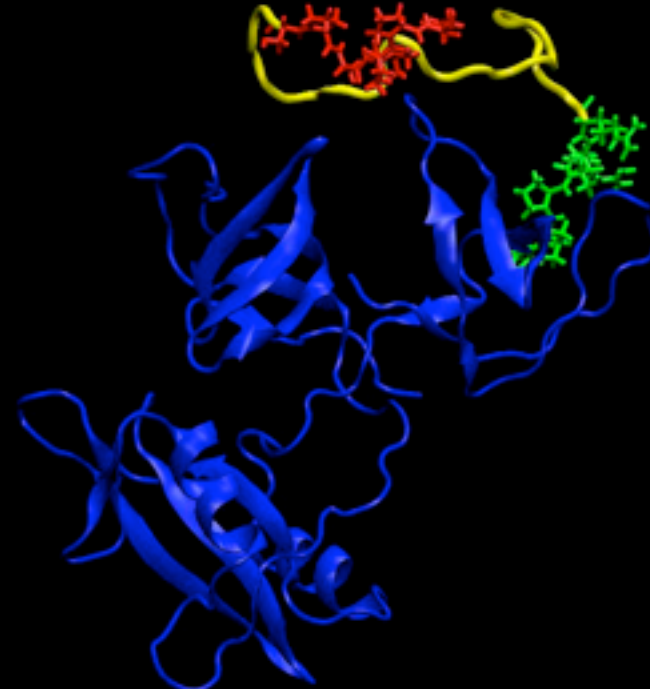
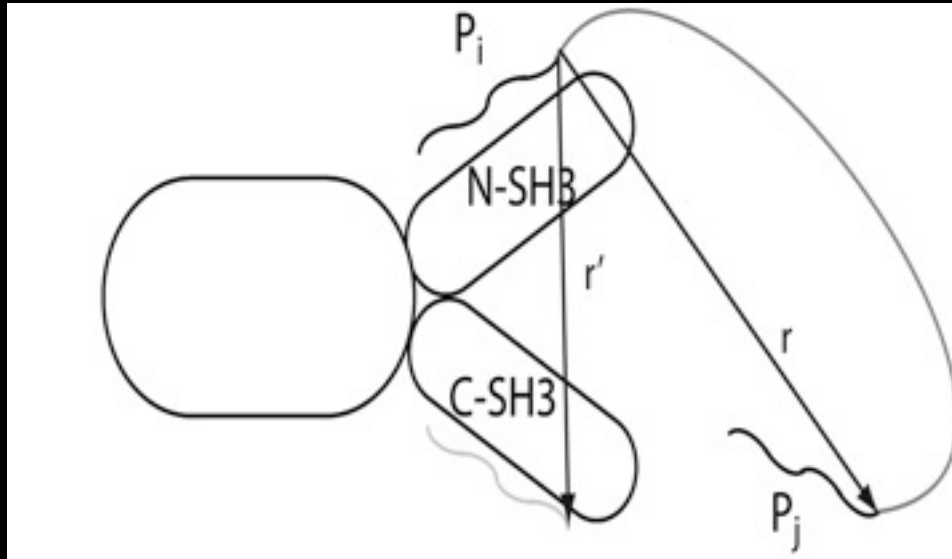
- Hybrid model from polymer theory and MD Simulations.

### Binding of Grb2 to Sos1

- Thermodynamic modeling of Grb2-Sos1 complexes.

# Signaling Cascades: Quantifying Multivalent Binding

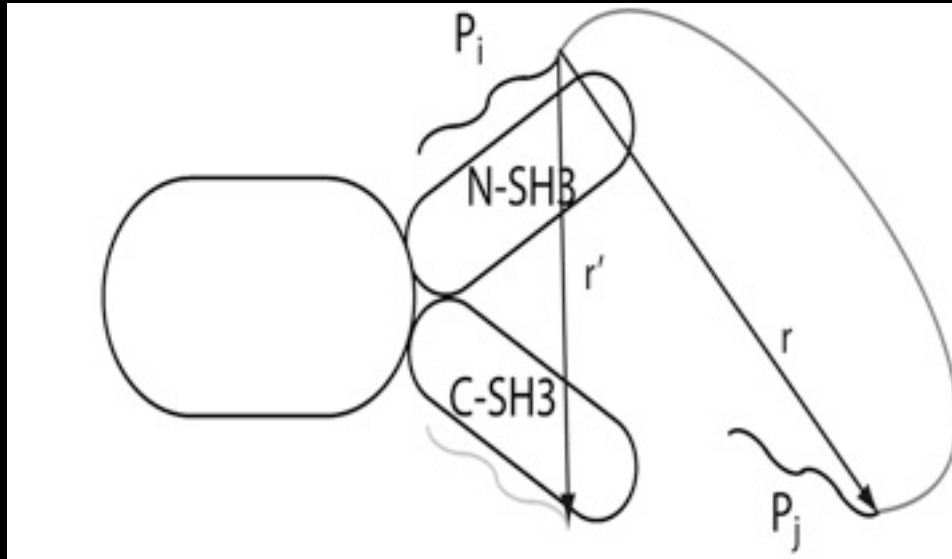
The flexibility of both binding partners determine  $C_{eff}$ .



$$C_{eff} (P_i^N, P_j^C) = \int_{r=0}^{\infty} \int_{r'=0}^{\infty} p_{bs}(\vec{r}') p_{pep}(P_i^N, P_j^C, \vec{r}) \delta(\vec{r} - \vec{r}') d^3 r d^3 r'$$
$$= \int_{r=0}^{\infty} p_{bs}(\vec{r}) p_{pep}(P_i^N, P_j^C, \vec{r}) d^3 r$$

# Signaling Cascades: Quantifying Multivalent Binding

The flexibility of both binding partners determine  $C_{eff}$ .



$$C_{eff} (P_i^N, P_j^C) = \int_{r=0}^{\infty} \int_{r'=0}^{\infty} p_{bs}(\vec{r}') p_{pep}(P_i^N, P_j^C, \vec{r}) \delta(\vec{r} - \vec{r}') d^3 r d^3 r'$$

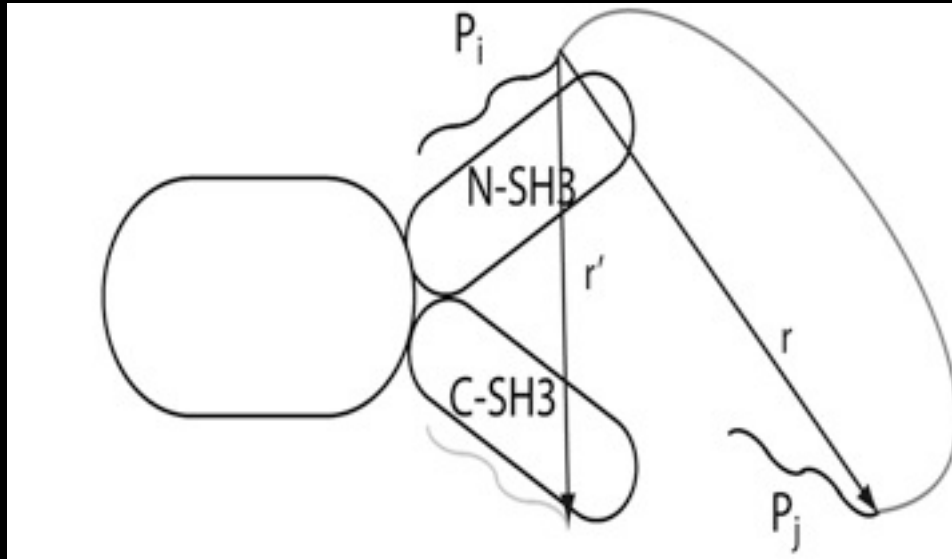
$$= \int_{r=0}^{\infty} p_{bs}(\vec{r}) p_{pep}(P_i^N, P_j^C, \vec{r}) d^3 r$$

Worm-like chain model for linker

$$p_{pep}(P_i^N, P_j^C, r) = \left( \frac{3}{4\pi l_p l_c} \right)^{3/2} \times \exp\left( \frac{-3r^2}{4l_p l_c} \right)$$

# Signaling Cascades: Quantifying Multivalent Binding

The flexibility of both binding partners determine  $C_{eff}$ .



$$C_{eff} (P_i^N, P_j^C) = \int_{r=0}^{\infty} \int_{r'=0}^{\infty} p_{bs}(\vec{r}') p_{pep}(P_i^N, P_j^C, \vec{r}) \delta(\vec{r} - \vec{r}') d^3 r d^3 r'$$

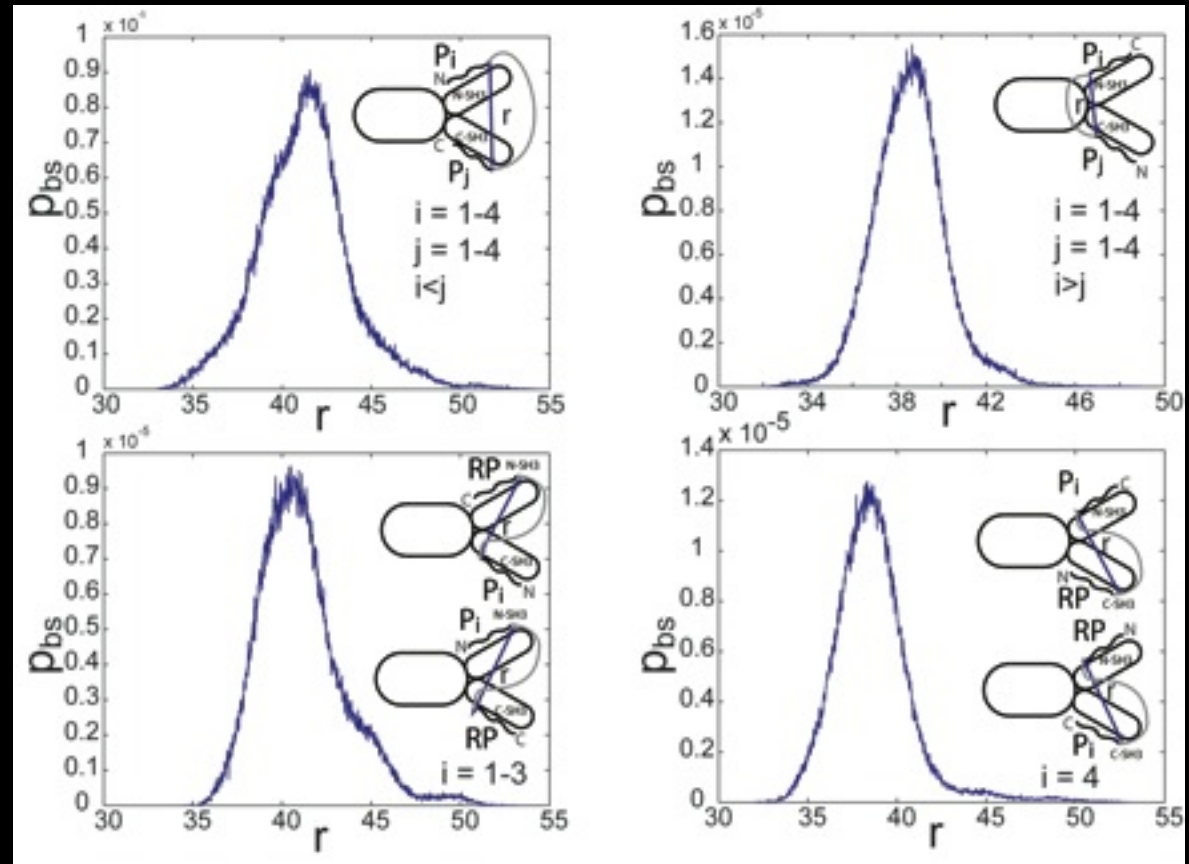
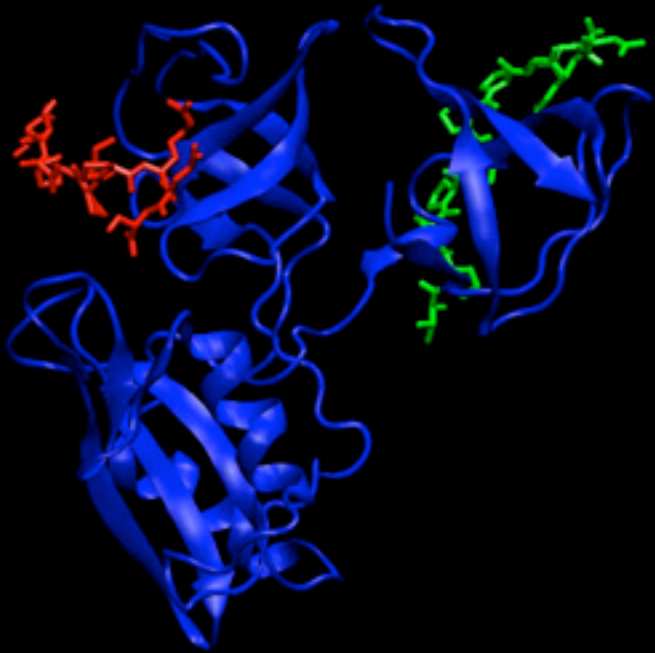
$$= \int_{r=0}^{\infty} p_{bs}(\vec{r}) p_{pep}(P_i^N, P_j^C, \vec{r}) d^3 r$$

Worm-like chain model for linker

$$p_{pep}(P_i^N, P_j^C, r) = \left( \frac{3}{4\pi l_p l_c} \right)^{3/2} \times \exp\left( \frac{-3r^2}{4l_p l_c} \right)$$

# Signaling Cascades: Quantifying Multivalent Binding

The flexibility of the linker and the motion of the two domains with respect to each other influence  $C_{\text{eff}}$



## Signaling Cascades: Quantifying Multivalent Binding

The local concentration of the other motifs near the second binding site of Grb2 is in mM range

Motif bound to N-SH3 domain

	P1	P2	P3	RP	P4
P1	-	0.6	2.1	1.6	1.6
P2	0.3	-	0.7	1.7	1.9
P3	1.6	0.4	-	1.5	2.1
RP	1.6	1.7	1.5	-	0.07
P4	1.4	1.6	1.7	0.07	-

Motif bound to C-SH3 domain

$C_{\text{eff}}$  (mM)



# Signaling Cascades: Quantifying Multivalent Binding

The local concentration of the other motifs near the second binding site of Grb2 is in mM range

Motif bound to N-SH3 domain

	P1	P2	P3	RP	P4
P1	-	0.6	2.1	1.6	1.6
P2	0.3	-	0.7	1.7	1.9
P3	1.6	0.4	-	1.5	2.1
RP	1.6	1.7	1.5	-	0.07
P4	1.4	1.6	1.7	0.07	-

Motif bound to C-SH3 domain

$C_{\text{eff}}$  (mM)

Motif bound to N-SH3 domain

	P1	P2	P3	RP	P4
P1	-	12	7	9	7
P2	164	-	220	94	62
P3	42	220	-	133	67
RP	43	56	130	-	2100
P4	37	46	90	2300	-

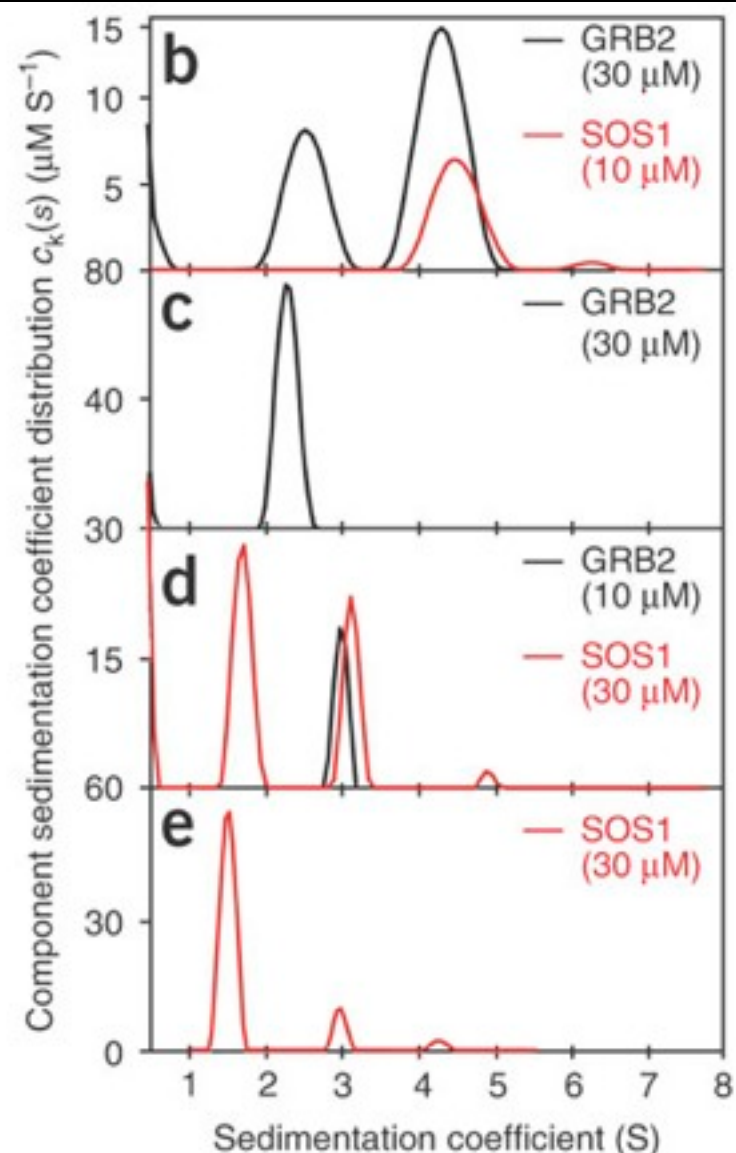
Motif bound to C-SH3 domain

$1/\bar{K}_{ij}^{NC}$  ( $\mu M$ )

A Sethi, et al., PLoS Comp. Biol., 2011

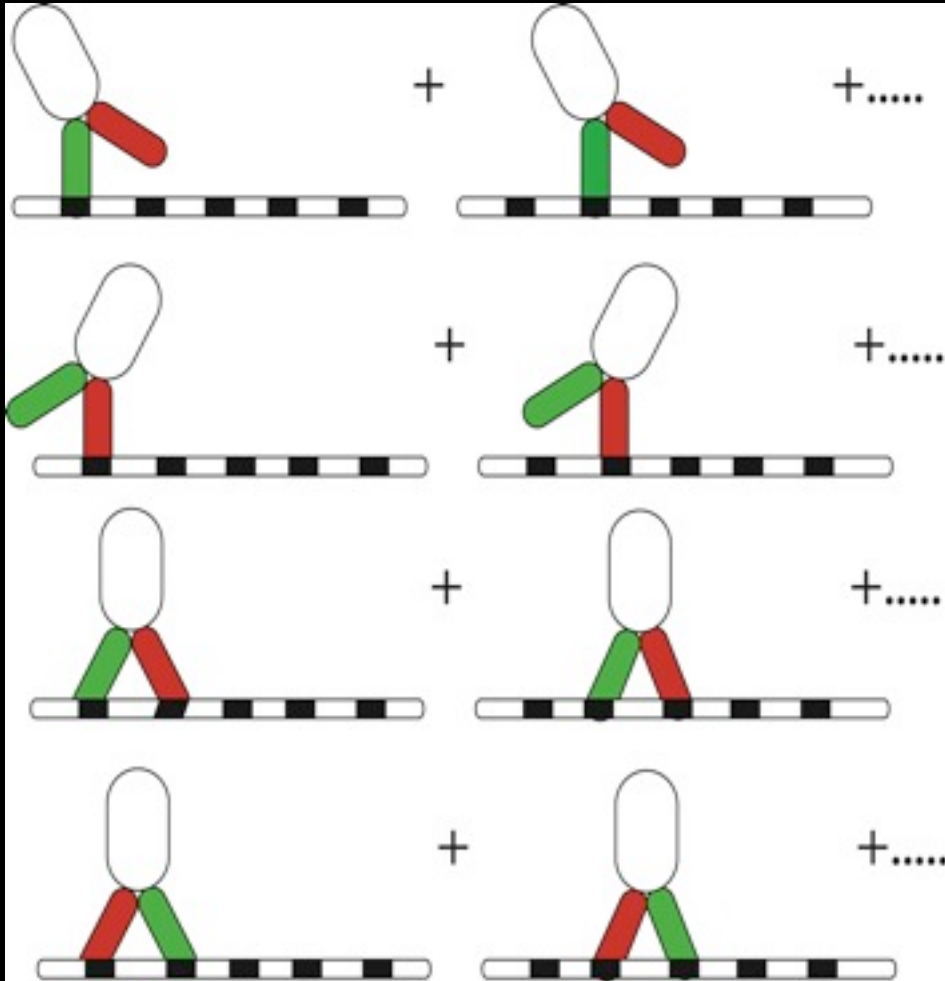
# Signaling Cascades: Quantifying Multivalent Binding

It is very difficult to distinguish different complexes formed between Sos1 and Grb2 in experiments



# Signaling Cascades: Quantifying Multivalent Binding

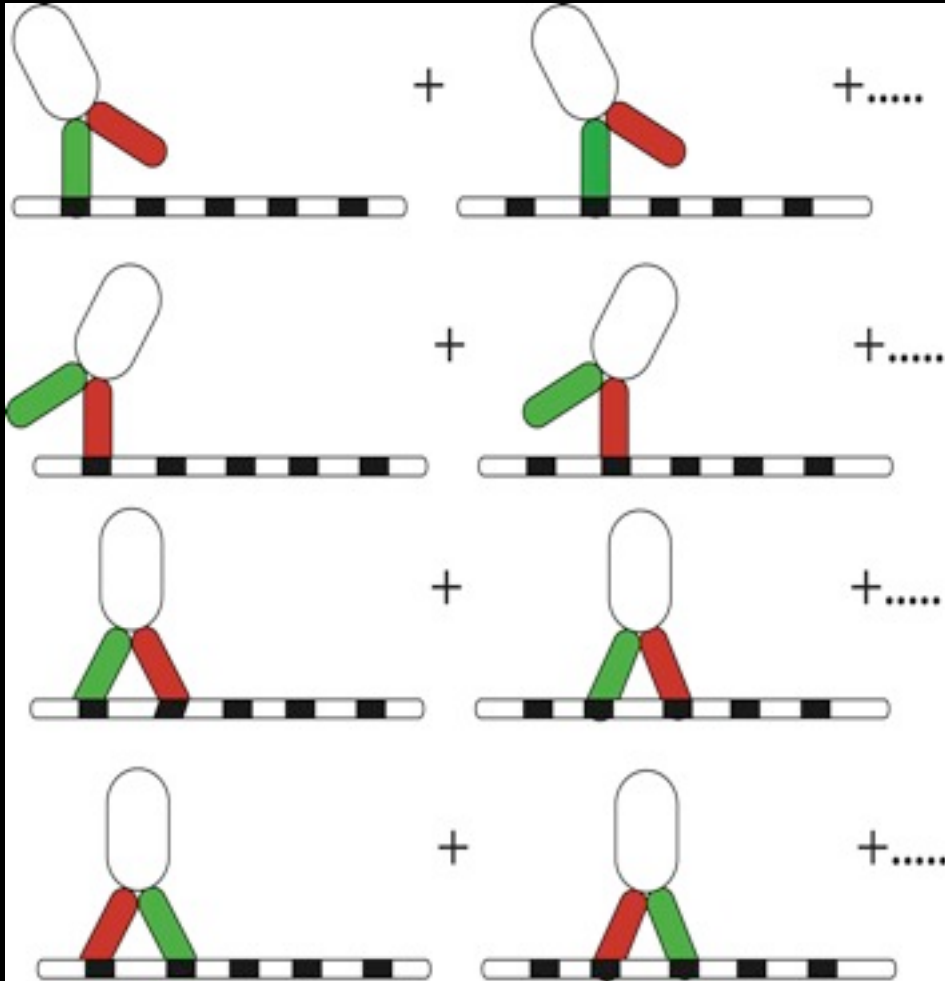
The 1:1 complex is actually a large combination of complexes that should be taken into account.



$$\bar{K}_1 = \sum_{X=N,C} \sum_{i=1}^5 K_i^X + \sum_{i=1}^5 \sum_{j=1, j \neq i}^5 \bar{K}_{ij}^{NC}$$

# Signaling Cascades: Quantifying Multivalent Binding

The 1:1 complex is actually a large combination of complexes that should be taken into account.

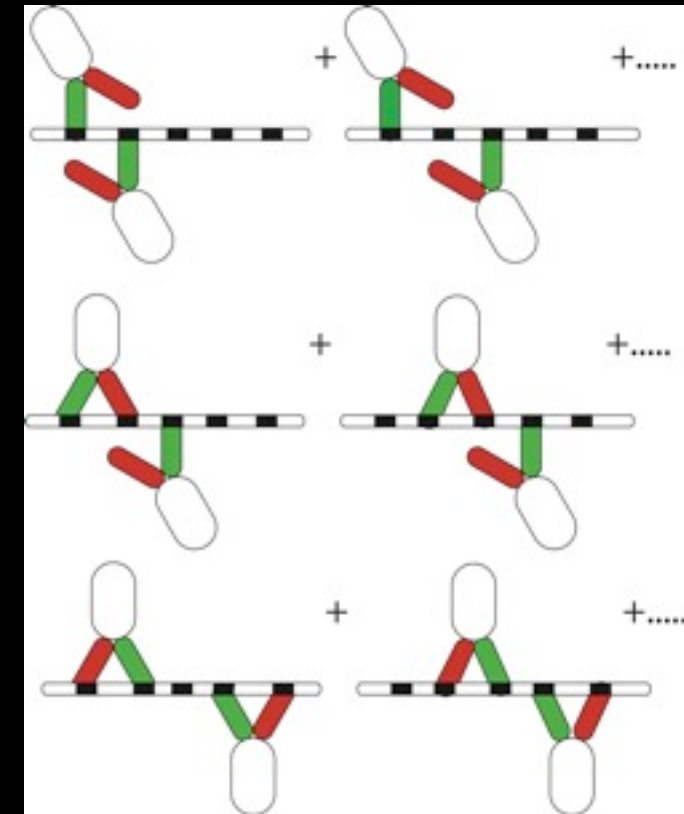


$$\bar{K}_1 = \sum_{X=N,C} \sum_{i=1}^5 K_i^X + \sum_{i=1}^5 \sum_{j=1, j \neq i}^5 \bar{K}_{ij}^{NC}$$

Predicted value = 1.2  $\mu\text{M}$ .  
 Experimental value = 0.3  $\mu\text{M}$ .

# Signaling Cascades: Quantifying Multivalent Binding

The 2:1 complex contains an even larger combination of complexes.



Predicted value = 14  $\mu\text{M}$ .  
 Experimental value = 1  $\mu\text{M}$ .

$$\bar{K}_2 = \frac{1}{2\bar{K}_1} \left( \sum_{X,Y=N,C} \sum_{i=1}^5 \sum_{j=1, j \neq i}^5 K_i^X K_j^Y + \sum_{X=N,C} \sum_{i=1}^5 \sum_{j=1}^5 \sum_{k=1, k \neq j \neq i}^5 K_i^X \bar{K}_{jk}^{NC} \right. \\ \left. + \sum_{i=1}^5 \sum_{j=1}^5 \sum_{k=1}^5 \sum_{l=1, l \neq k \neq j \neq i}^5 \bar{K}_{ij}^{NC} \bar{K}_{kl}^{NC} \right)$$

A Sethi, et al., PLoS Comp. Biol., 2011

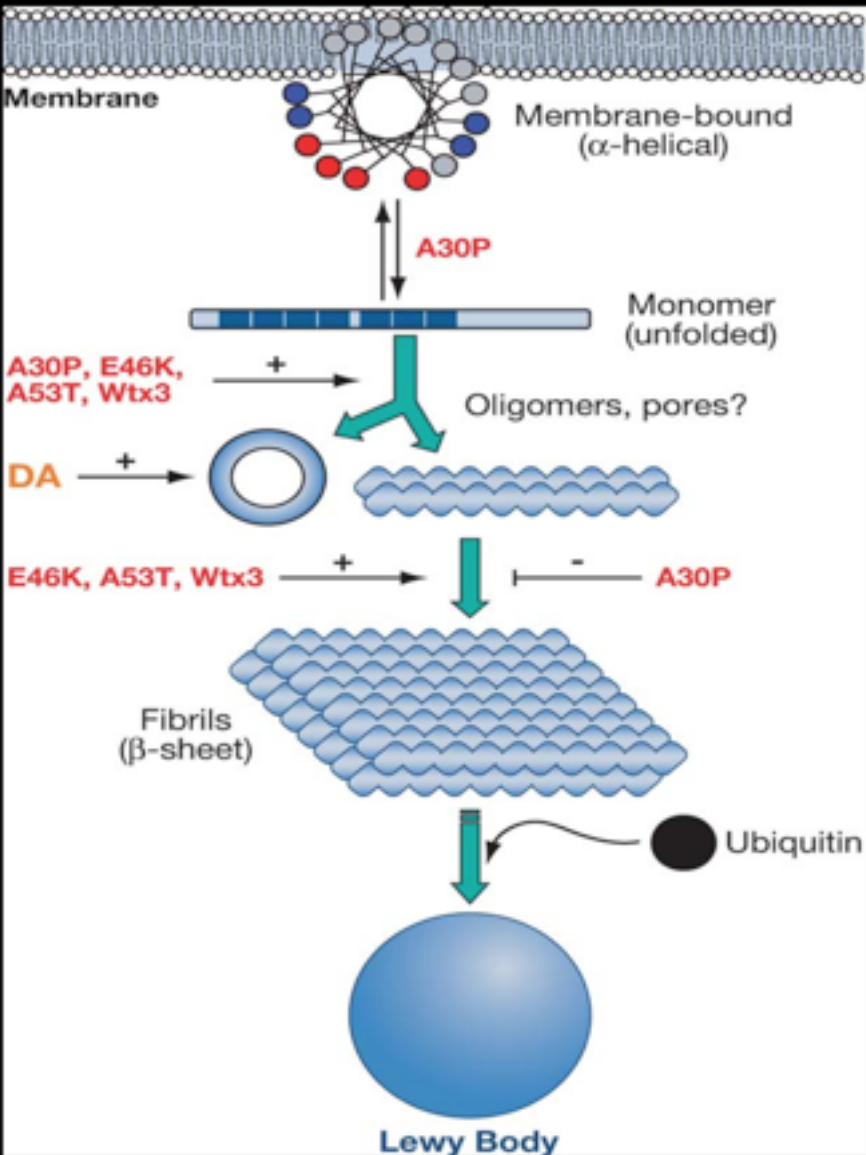


Is a WLC model a good model to get the probability of the distance between the two ends of a linker?

Can we use MD-based methods to figure this probability density?

## Signaling Cascades: Disordered Regions

Fibrils or oligomers of disordered proteins are often implicated in neurological diseases





## Signaling Cascades: Disordered Regions

Fibrils or oligomers of disordered proteins are often implicated in neurological diseases

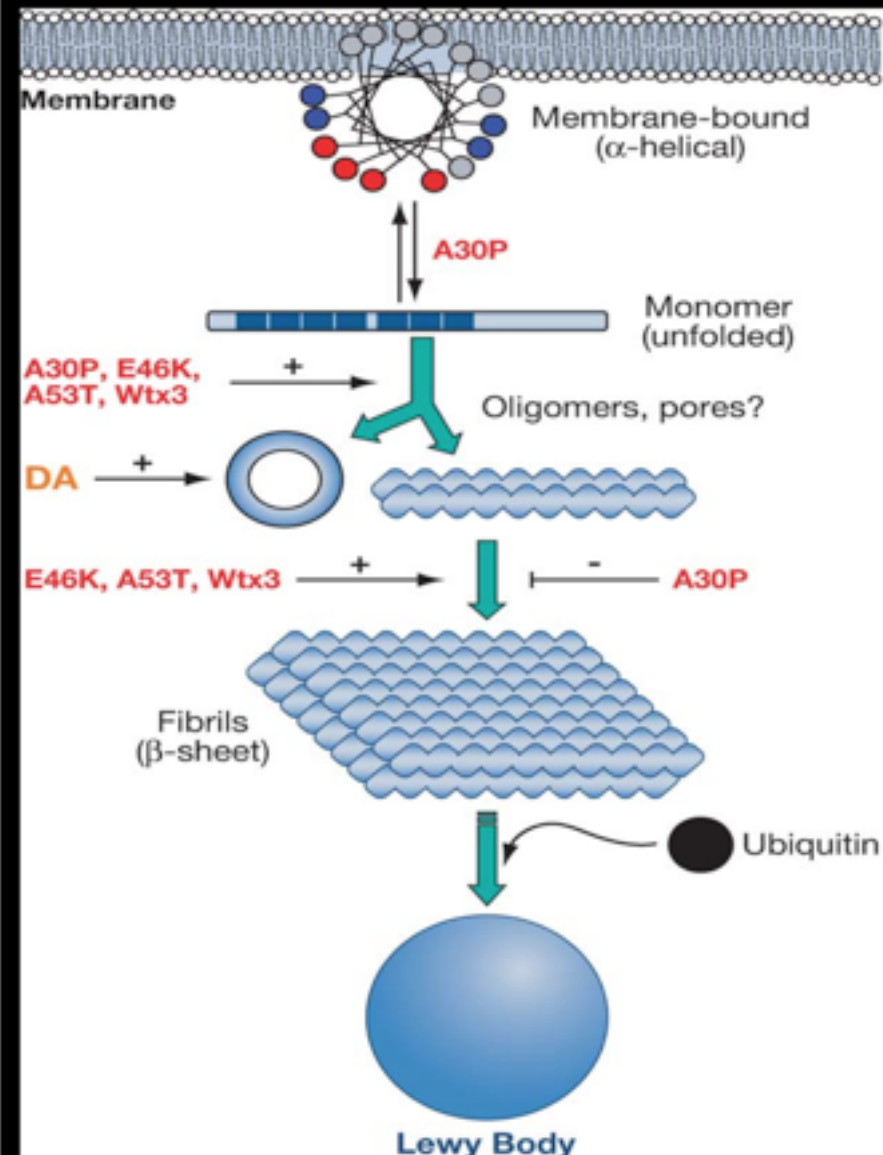


Image Source:

<http://www.genome.gov/pressDisplay.cfm?photoID=10004>

Cookson, Ann. Rev. Biochem., 2005

## Signaling Cascades: Disordered Regions

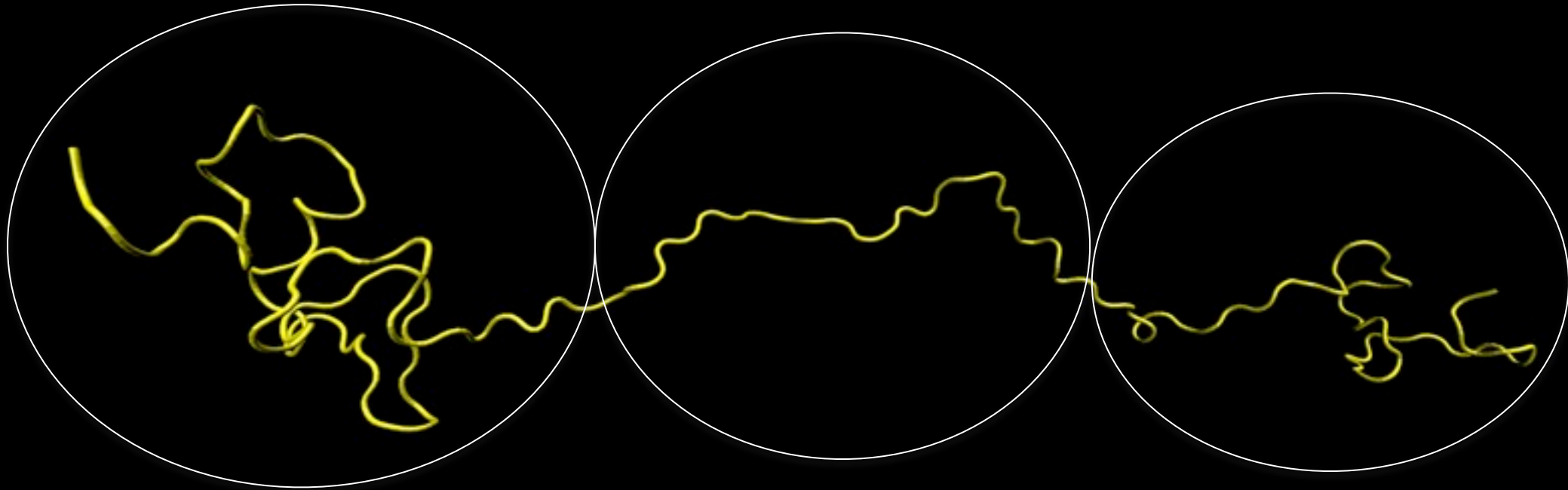
Is divide and conquer technique possible for IDPs?



Can we break an IDP into smaller, more manageable, pieces in order to calculate its conformational network (divide and conquer approach)?

## Signaling Cascades: Disordered Regions

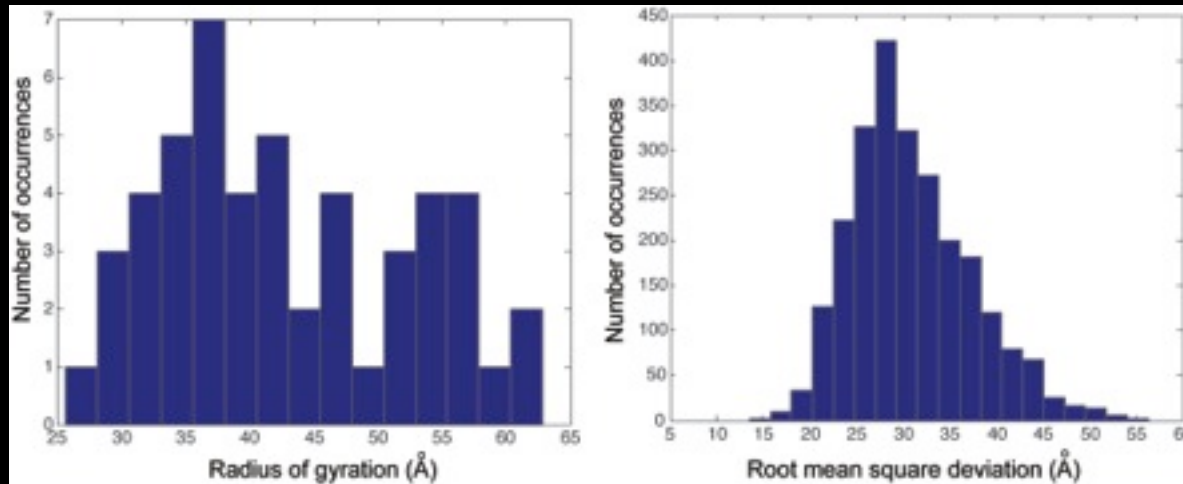
Is divide and conquer technique possible for IDPs?



Can we break an IDP into smaller, more manageable, pieces in order to calculate its conformational network (divide and conquer approach)?

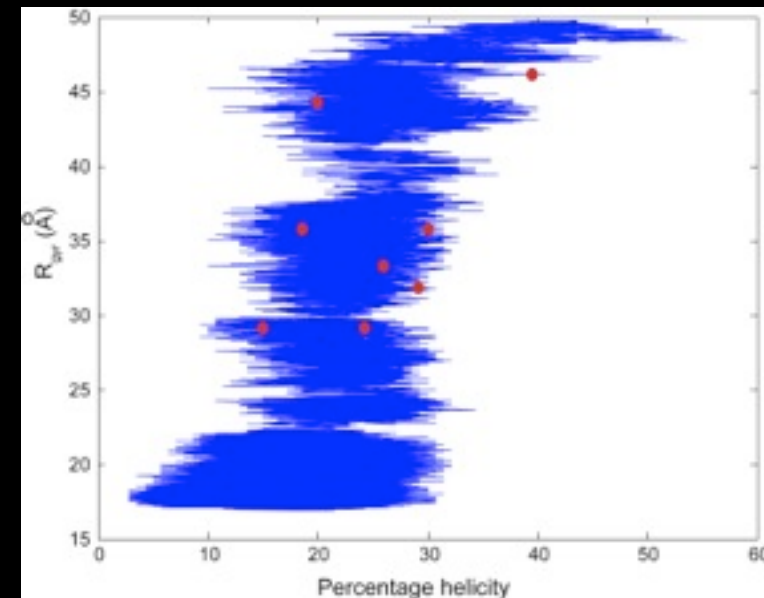
## Signaling Cascades: Disordered Regions

We performed 100 simulations to sample the entire phase space of  $\alpha$ -synuclein



50 conformations generated using random coil generator (conformations of random coils in pdb database)

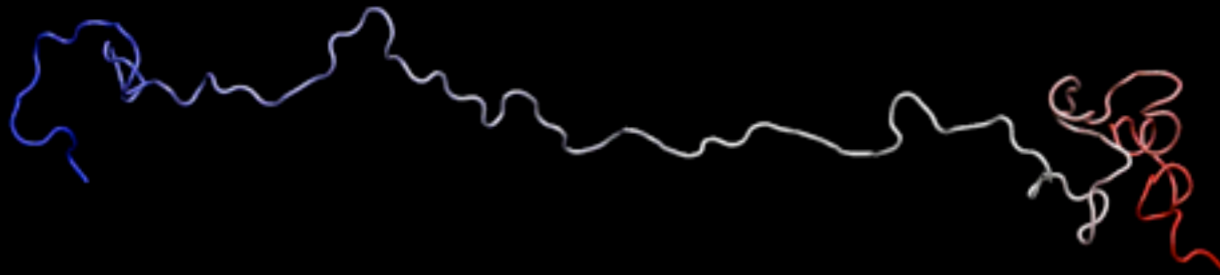
8 partially helical conformations chosen from one simulation of the helix turn helix conformation + 2 native conformations of  $\alpha$ -synuclein = 50 partially helical simulations (five simulations starting from 10 different conformations).



A Sethi, et al., Biophys. J., 2012

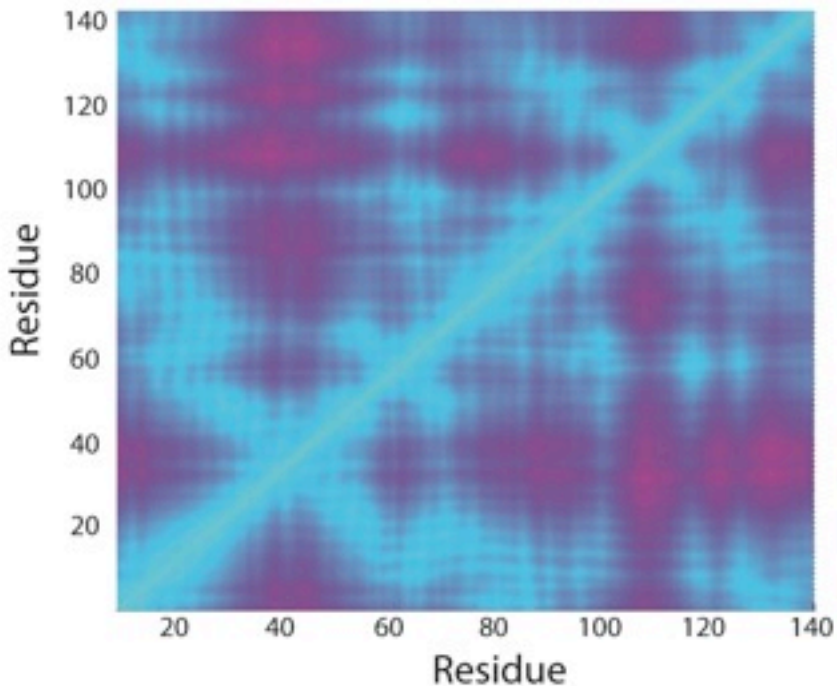
## Signaling Cascades: Disordered Regions

The protein collapses and remains collapsed during the timescale of these simulations.

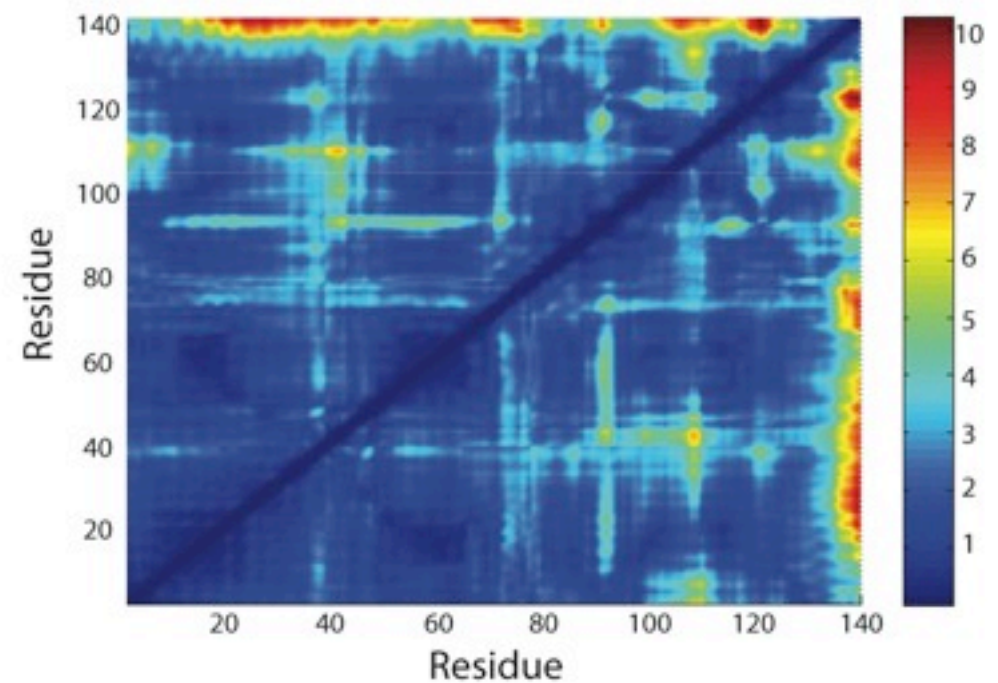


## Signaling Cascades: Disordered Regions

The protein stays trapped in a minima upto hundreds of nanoseconds.



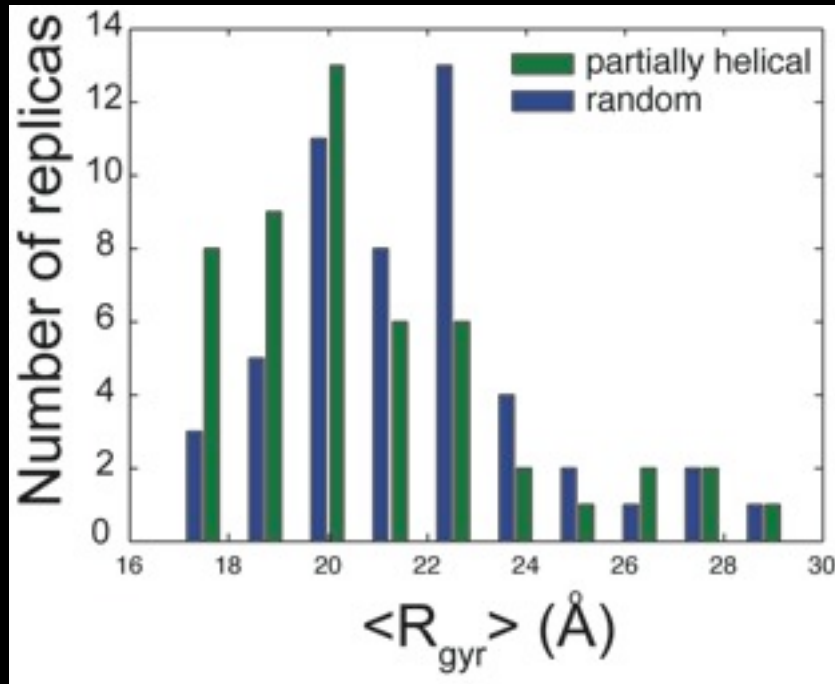
$\langle r(i,j) \rangle$  (in Angstroms) during  
300ns of 298K (NMR)



Standard deviation in  $r(i,j)$   
(in Angstroms) during  
300ns of the 298K (NMR)  
simulation.

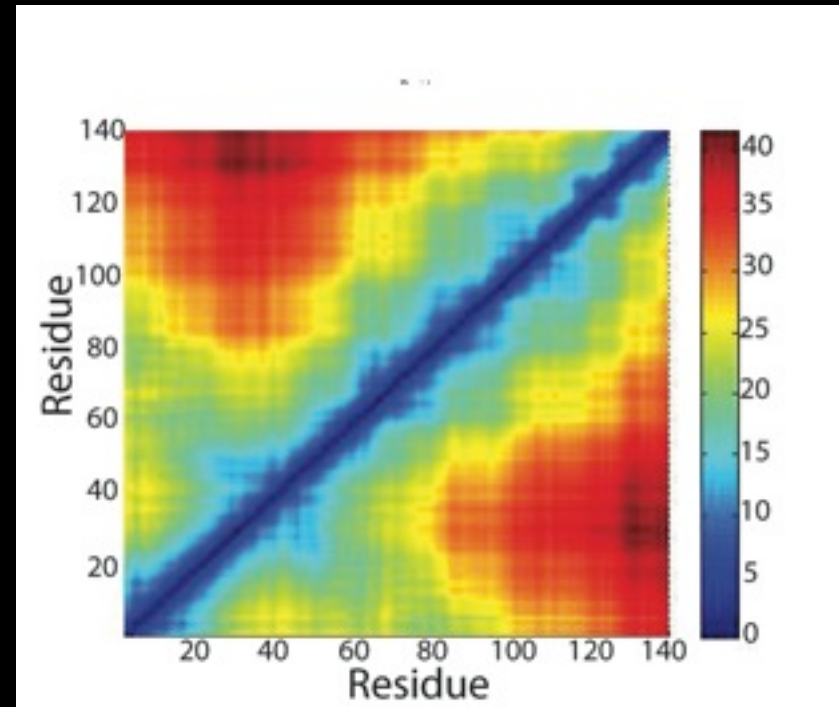
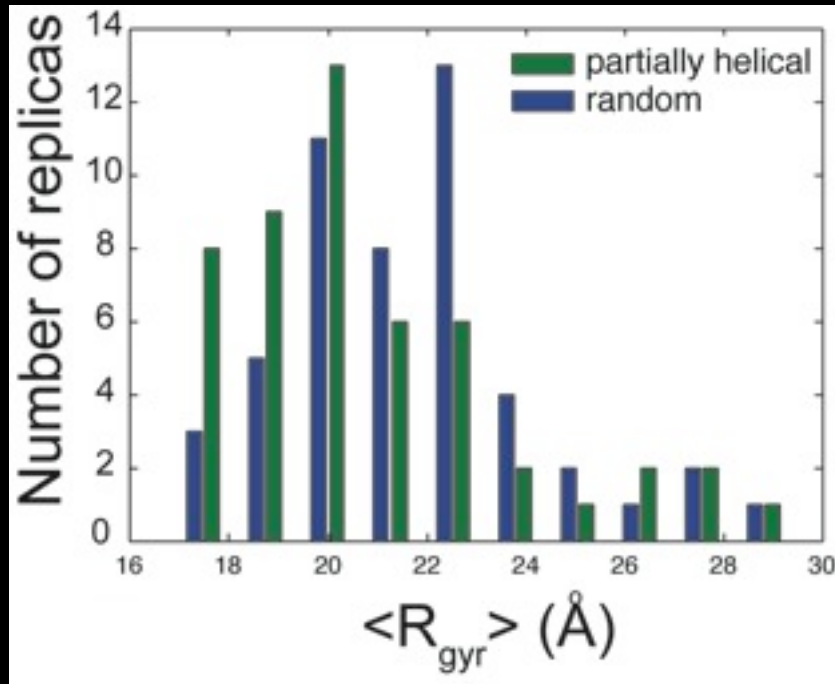
## Signaling Cascades: Disordered Regions

The protein forms a heterogeneous collapsed state.



## Signaling Cascades: Disordered Regions

The protein forms a heterogeneous collapsed state.

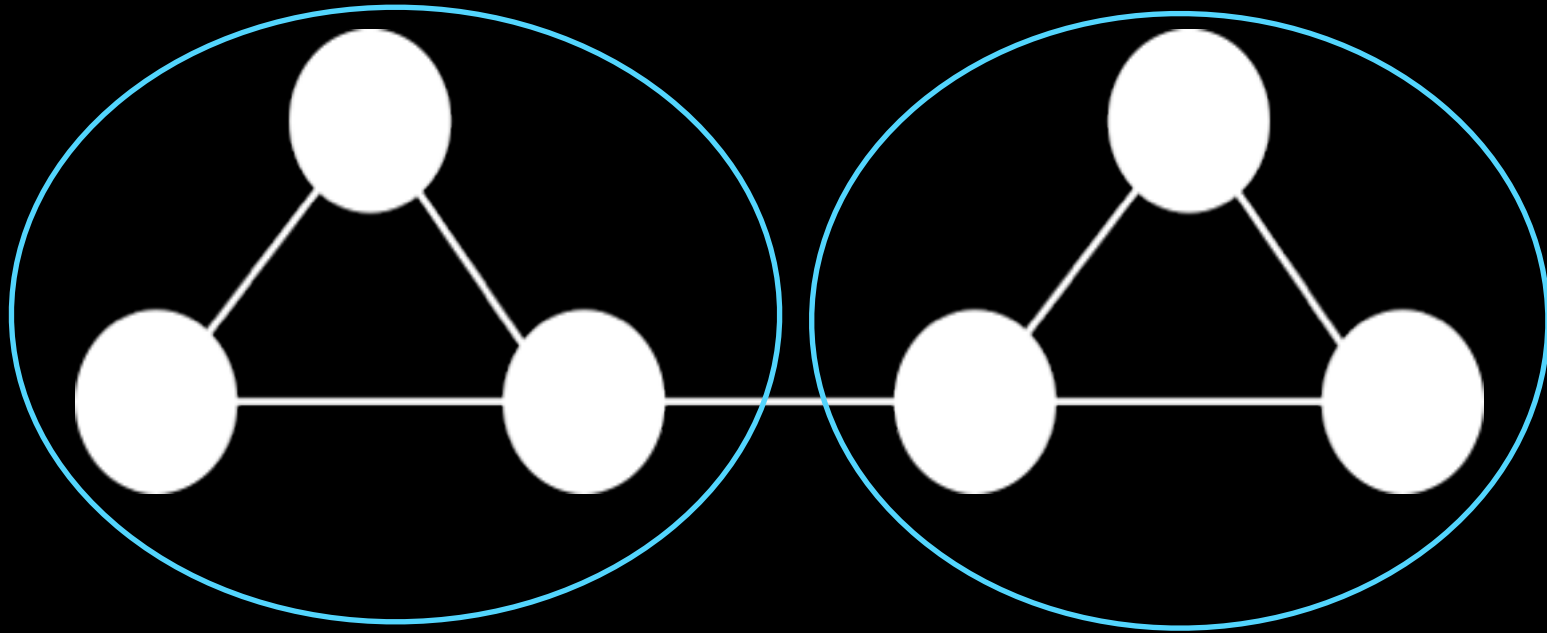


The protein behaves like an ideal chain in poor solvent.  
There are no persistent long-range contacts in the protein.



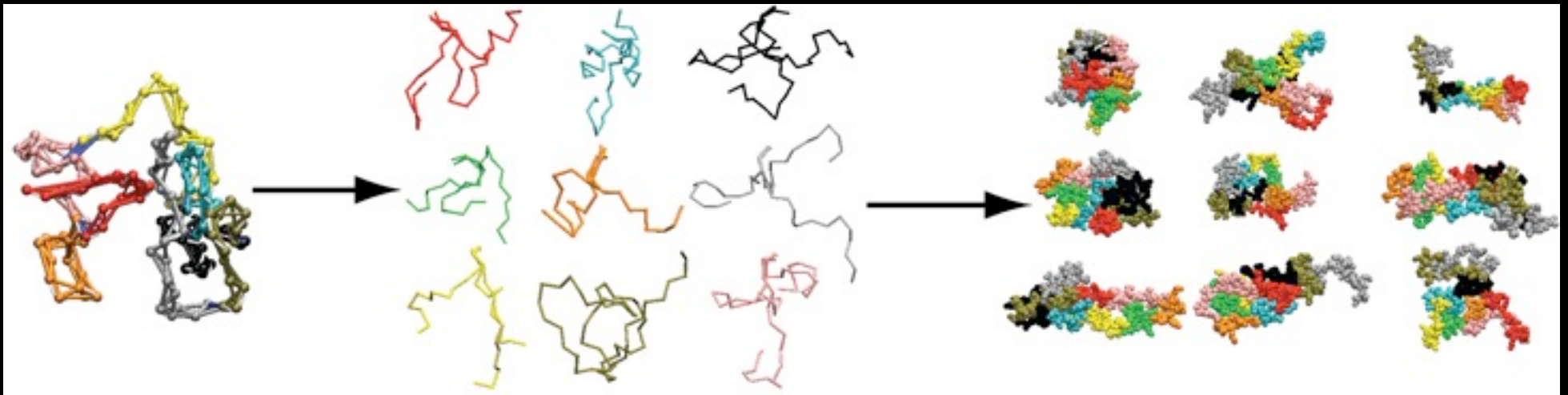
## Signaling Cascades: Disordered Regions

Communities are modules within the network.



## Signaling Cascades: Disordered Regions

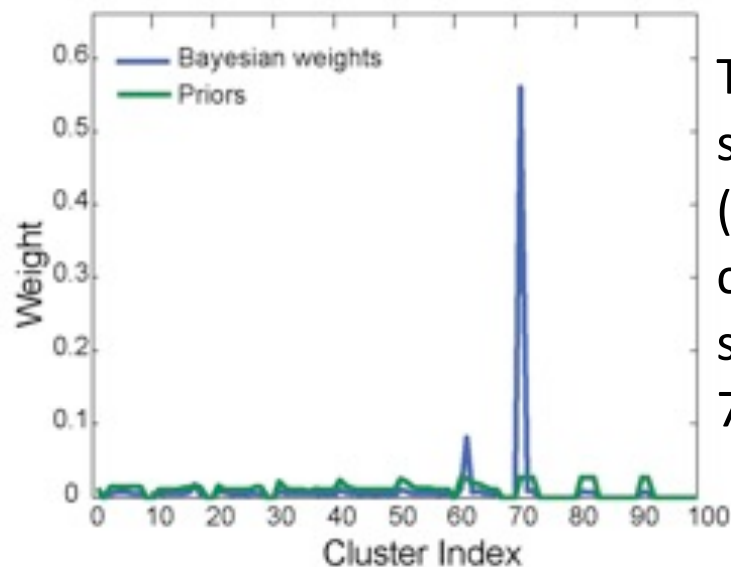
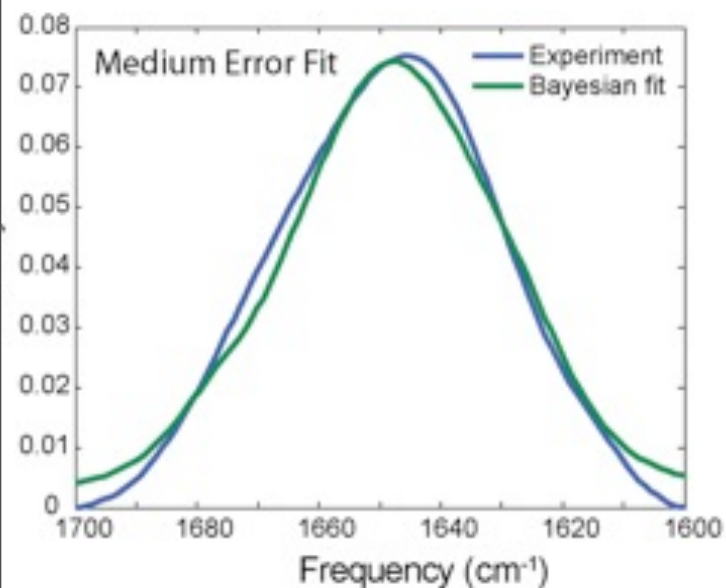
The heterogeneous nature of the compact states may make it possible to split  $\alpha$ -synuclein into smaller peptides



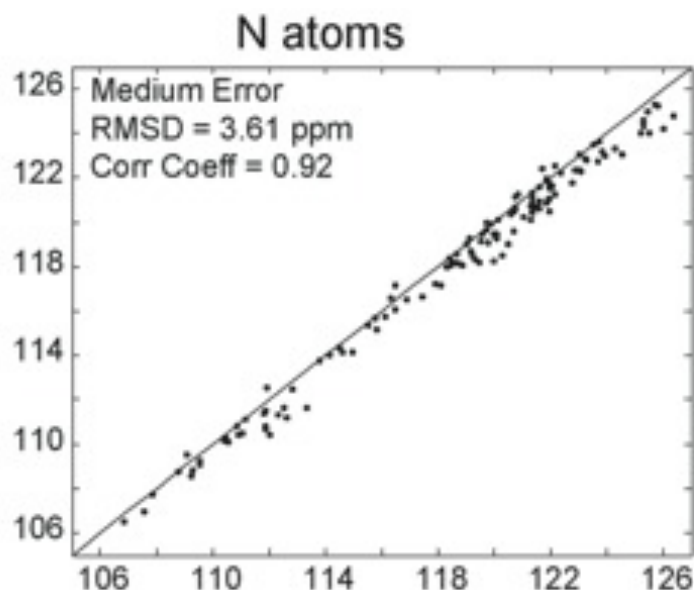
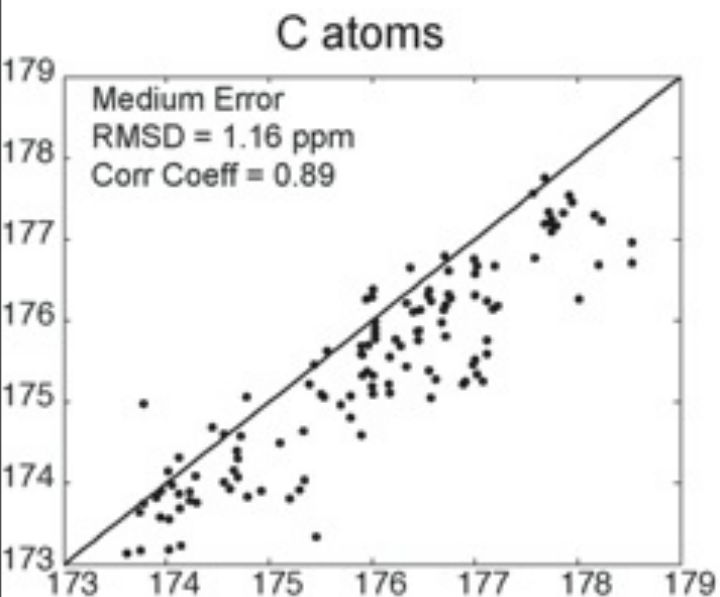
Nine communities in combined trajectory.

The communities are made of sequentially contiguous residues.

# Bayesian Statistics of Intrinsically Disordered Proteins



The Bayesian ensemble of  $\alpha$ -synuclein is semi-collapsed (mean of 35 Angstroms) and contains very low residual secondary structure (mean of 7% helicity)



The predicted chemical shifts for the Bayesian ensemble are consistent with experimental observations

Sethi et al., Chem. Phys., In Press.

# Signaling Cascades

## Conclusions

## Signaling Cascades

# Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

# Signaling Cascades

## Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

The distribution of the multivalent complexes further increases the effective equilibrium constants for complexes of different stoichiometry to form.

## Signaling Cascades

# Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

The distribution of the multivalent complexes further increases the effective equilibrium constants for complexes of different stoichiometry to form.

IDPs can exist in a heterogeneous collapsed state in aqueous environment.

## Signaling Cascades

# Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

The distribution of the multivalent complexes further increases the effective equilibrium constants for complexes of different stoichiometry to form.

IDPs can exist in a heterogeneous collapsed state in aqueous environment.



## Signaling Cascades

# Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

The distribution of the multivalent complexes further increases the effective equilibrium constants for complexes of different stoichiometry to form.

IDPs can exist in a heterogeneous collapsed state in aqueous environment.

The collapsed states are stabilized by contacts that remain for hundreds of nanoseconds.

## Signaling Cascades

# Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

The distribution of the multivalent complexes further increases the effective equilibrium constants for complexes of different stoichiometry to form.

IDPs can exist in a heterogeneous collapsed state in aqueous environment.

The collapsed states are stabilized by contacts that remain for hundreds of nanoseconds.

## Signaling Cascades

# Conclusions

The local concentration in multivalent complexes can be quantified using a combination of polymer models and MD simulations.

The distribution of the multivalent complexes further increases the effective equilibrium constants for complexes of different stoichiometry to form.

IDPs can exist in a heterogeneous collapsed state in aqueous environment.

The collapsed states are stabilized by contacts that remain for hundreds of nanoseconds.

The heterogeneity in the collapsed states might make it possible to divide  $\alpha$ -Synuclein into smaller fragments.

# Acknowledgements

## Mentors:

S. Gnanakaran (LANL)  
Zan Schulten (UIUC)  
Byron Goldstein (LANL)



## Collaborators:

Carl Woese (UIUC) - aaRS and ribosome  
Dung Vu (LANL) - IDP  
Martin Gruebele (UIUC) - RNA unfolding  
V Jo Davisson (Purdue) - Allostery  
Cynthia Derdeyn (Emory) - HIV  
Bette Korber (LANL) - HIV  
Shishir Chundawat (GLBRC) - Biofuels  
Peter Goodwin (LANL) - Biofuels  
Bridget Wilson (UNM) - Signaling  
Chang-Shung Tung (LANL) - Signaling



Patrick O'Donoghue (Yale) - aaRS  
Rommie Amaro (UCSD) - Allostery  
Elijah Roberts (JHU) - Ribosome  
John Eargle (UIUC) - aaRS  
Jianhui Tian (ORNL) - IDP/HIV

