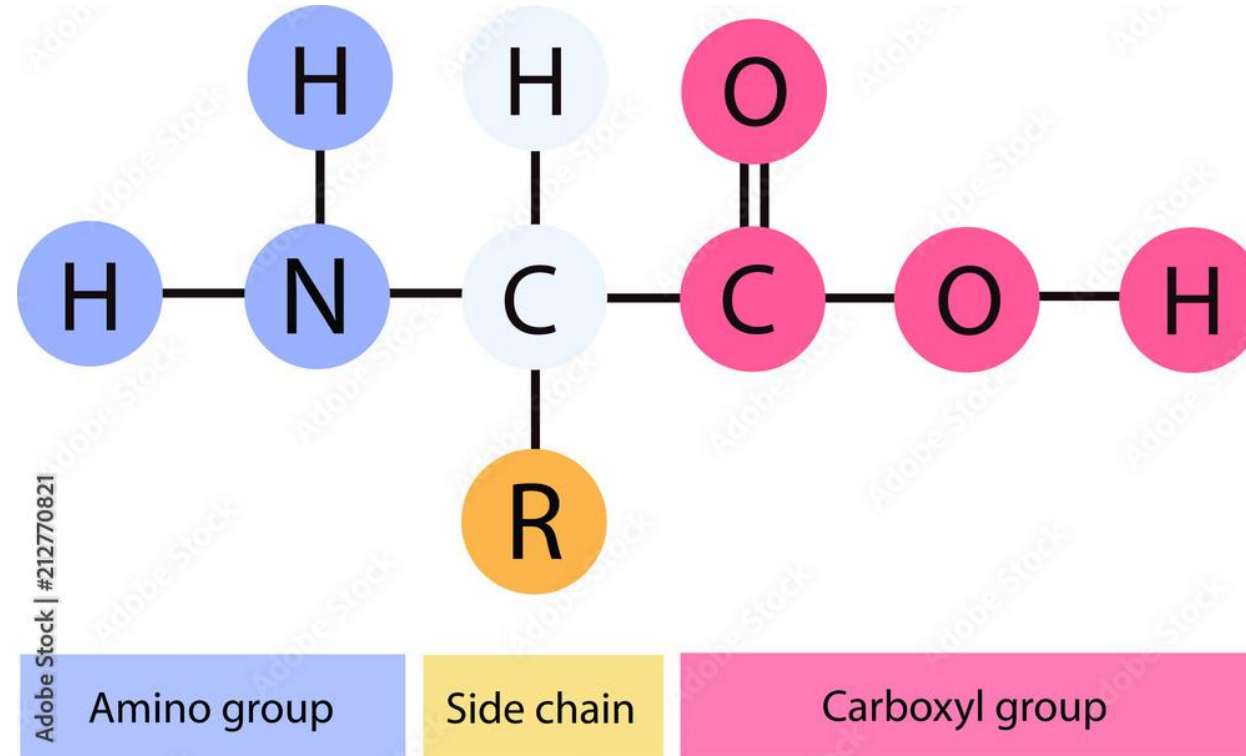# Motif and Domain Discovery
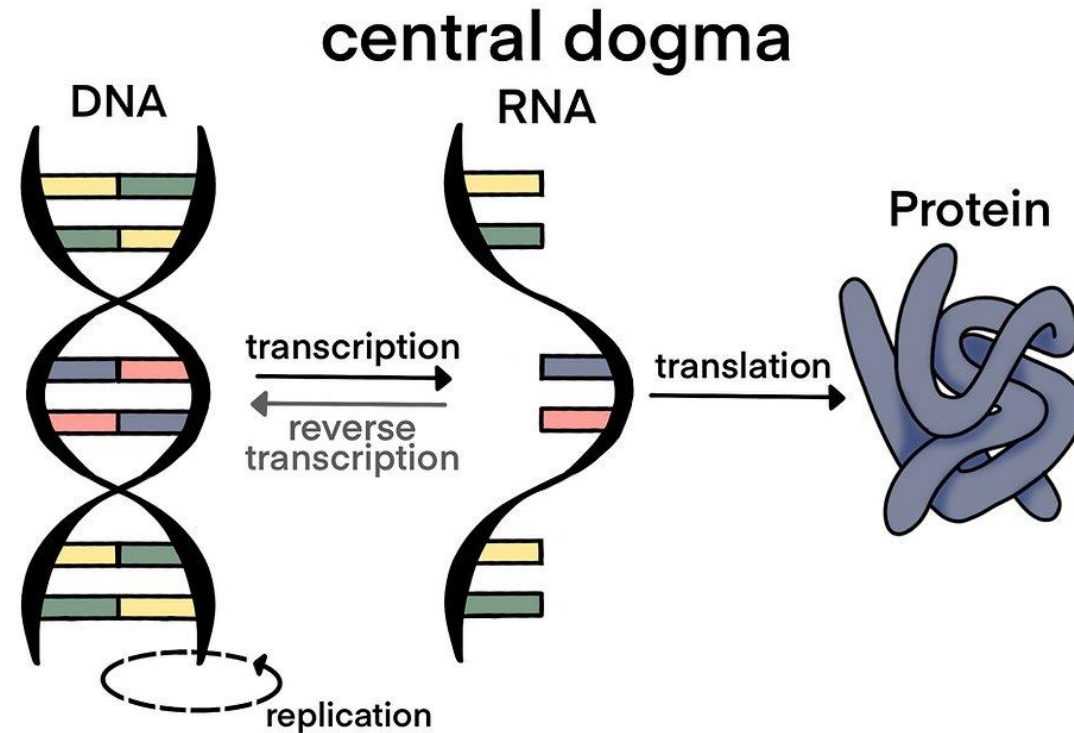
**Rohan Babaria and Joely Swanson**
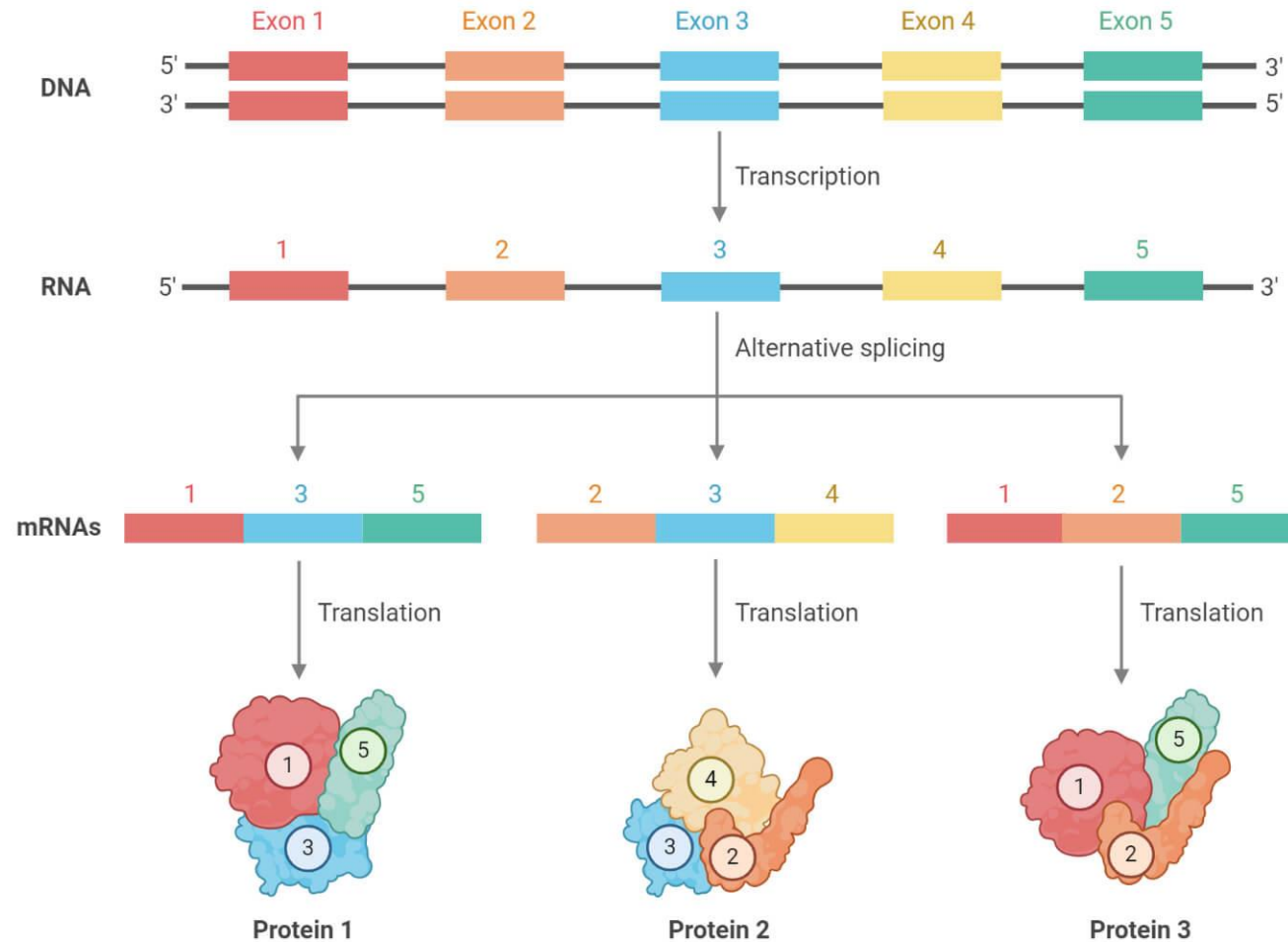
# What are Proteins?

**Proteins are polypeptides, chains of amino acids, that act as the functional unit of life.**
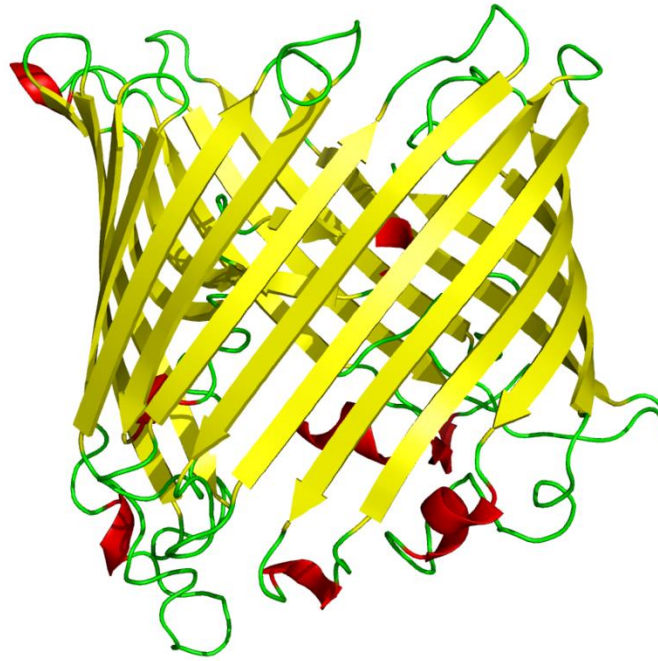
# How are proteins made?



Proteins are translated from mRNA transcripts derived from DNA.

# How are multiple protein variants produced?
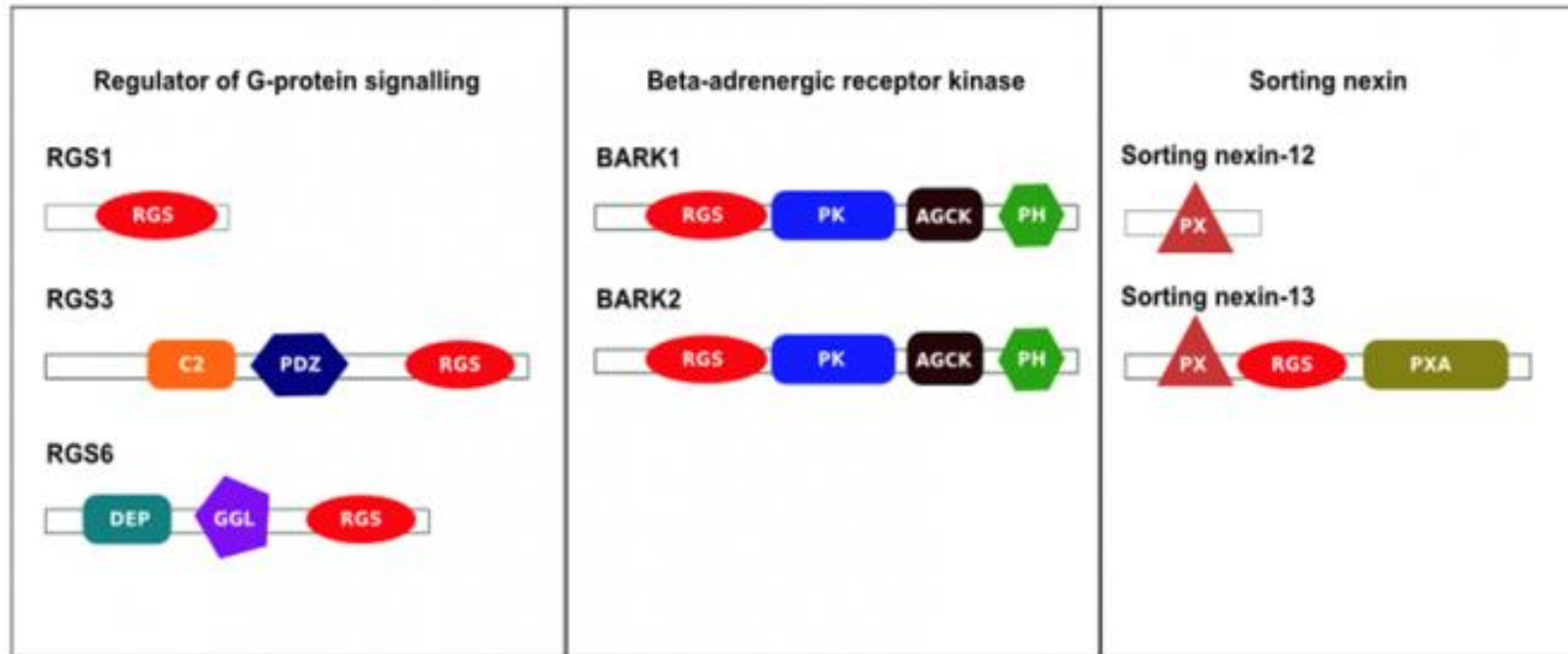
# What are the 4 structure types?

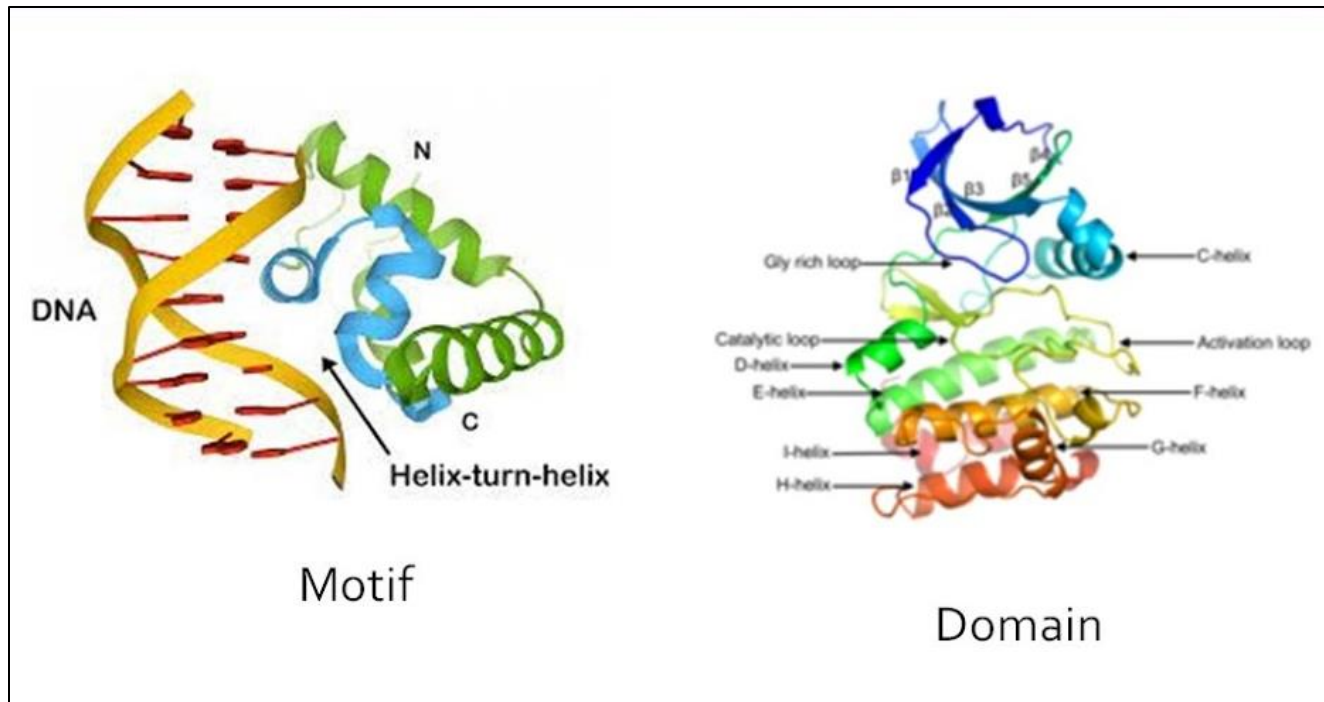# What are Protein Motifs?



**Small region of protein with a common three-dimensional structure/sequence shared among proteins.**

# What are Protein Domains?



**A conserved sequence pattern that acts as an independent functional and structural unit.**
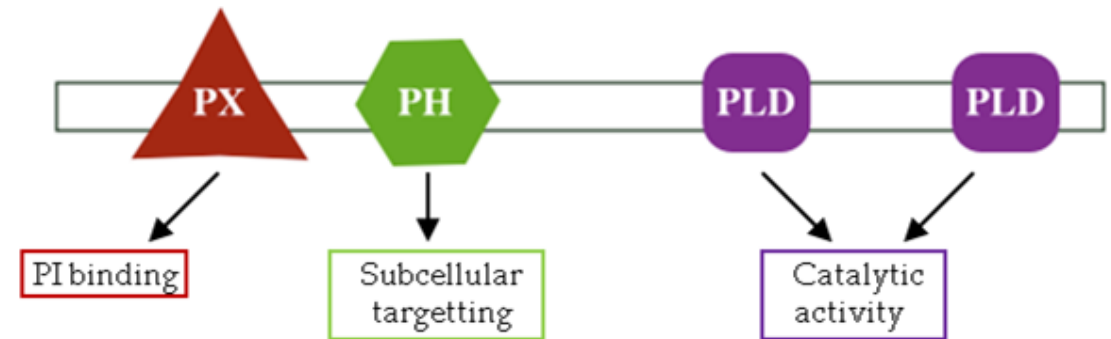
# What are the differences?



DNA

Helix-turn-helix

Motif

Domain

N

C

Gly rich loop
Catalytic loop
D-helix
E-helix
I-helix
H-helix
C-helix
Activation loop
F-helix
G-helix

β1 β2 β3 β4 β5

## MOTIF VERSUS DOMAIN IN PROTEIN STRUCTURE

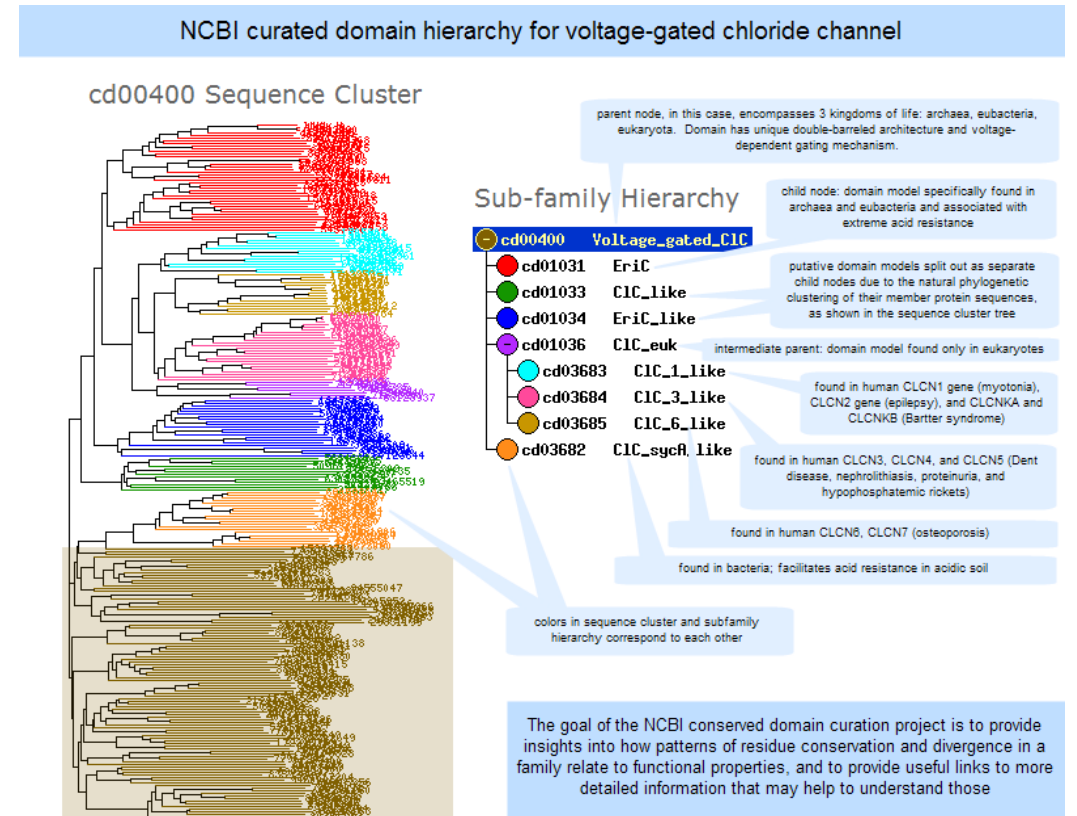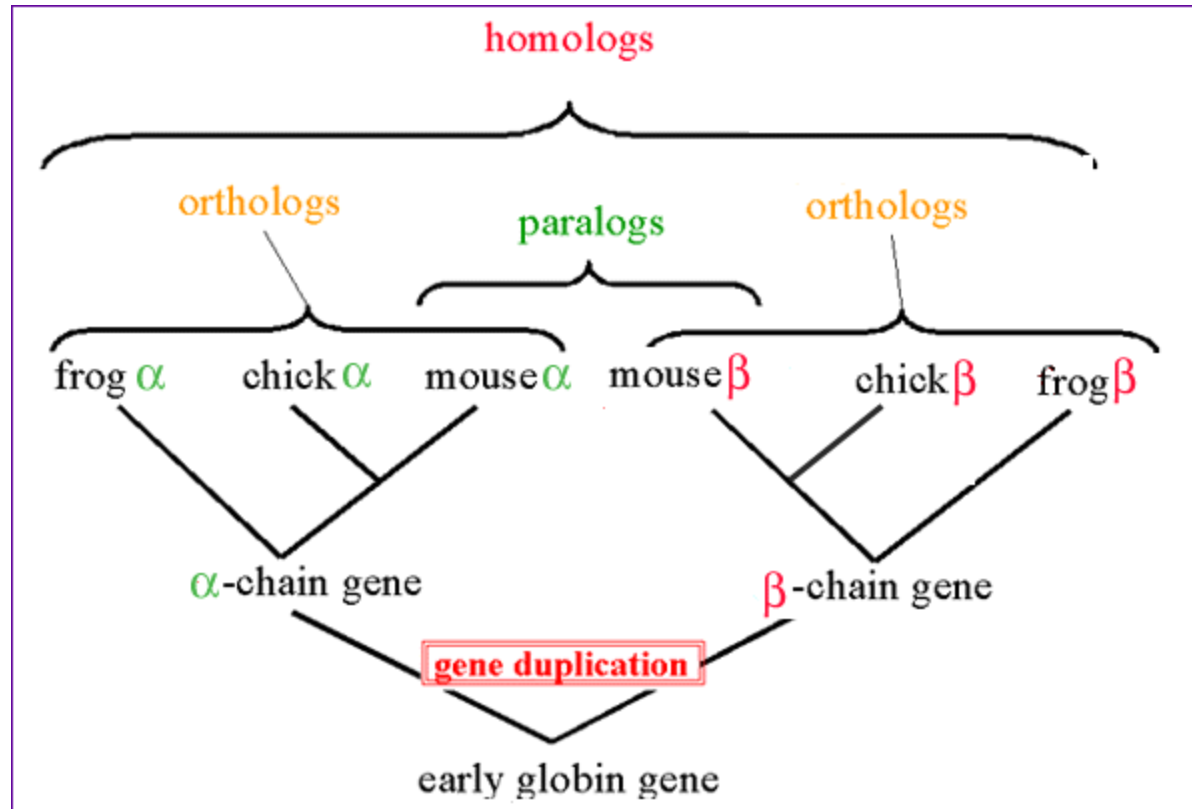| MOTIF | DOMAIN |
|---|---|
| A chain-like biological structure made up of connectivity between secondary structural elements | An independent folding unit of the three-dimensional protein structure |
| A supersecondary structure of a protein | A tertiary structure of the protein |
| Formed by the connected alpha-helices and beta-sheets through loops | Formed by the formation of disulfide bridges, ionic bonds, and hydrogen bonds between amino acid side chains |
| Mainly have a structural function in the protein structure | Mainly have functional importance |
| Have similar functions through protein families | Have unique functions |
| Are not stable independently | Are independently stable |

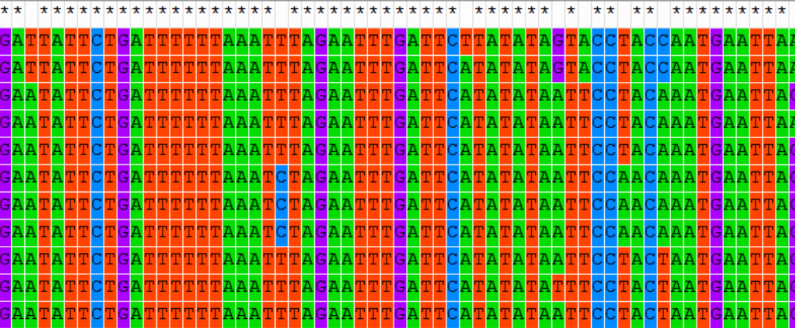Visit www.PEDIAA.com

# Why do we care?



**Conserved domain sequences give us insights into protein function and history.**

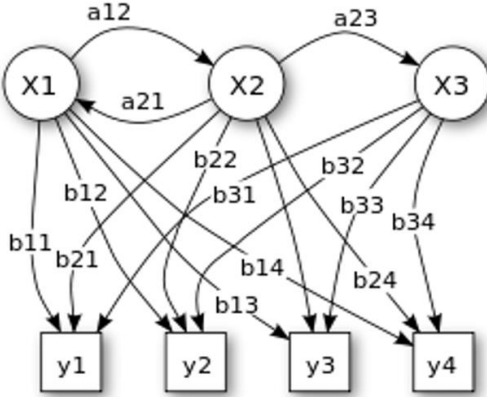# How do we characterize proteins through homology?

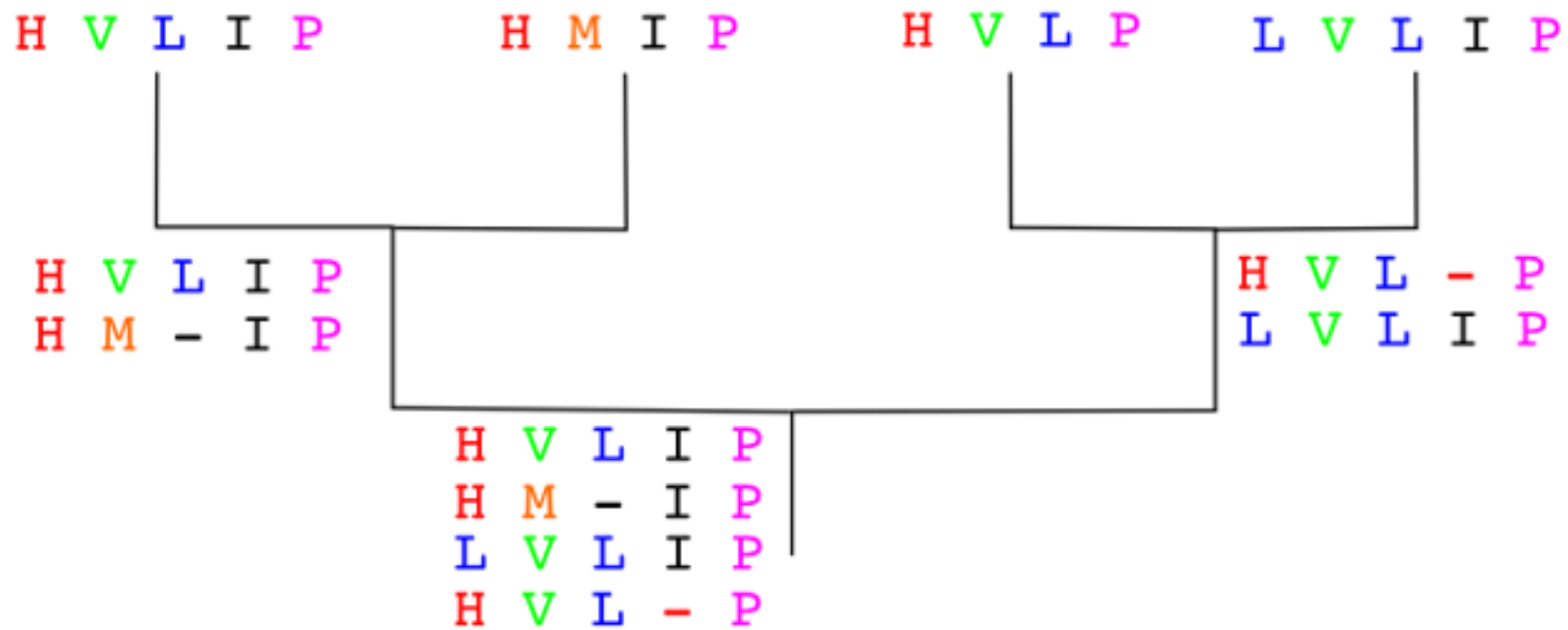# How can we identify conserved domains?



Multiple Sequence Alignment
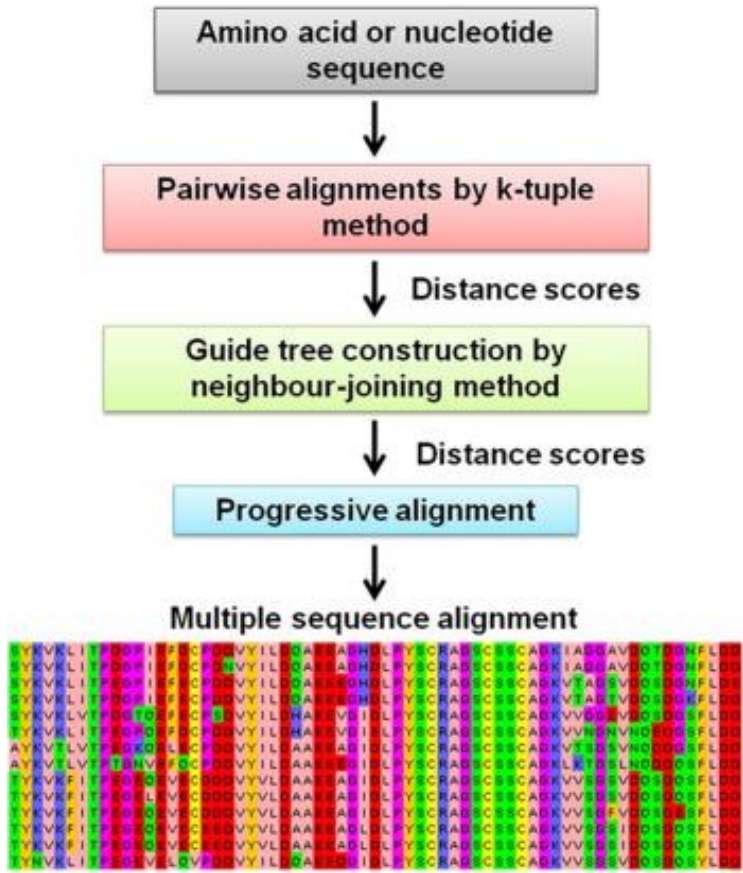Theory and Practice - Step-by-Step
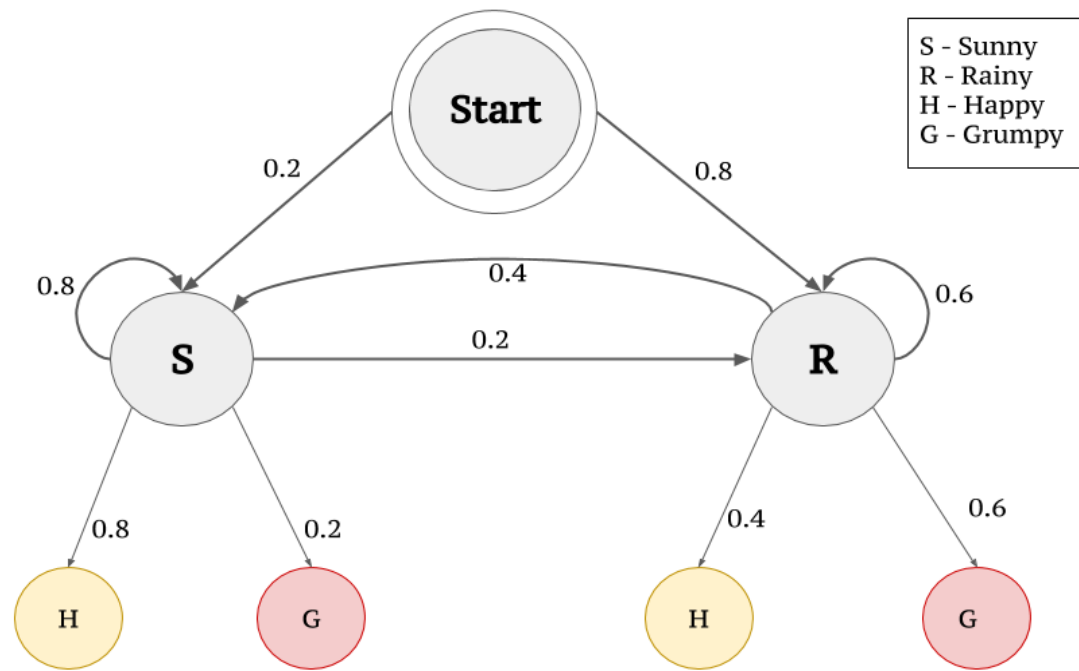


Hidden Markov Model

# What is MUSCLE?



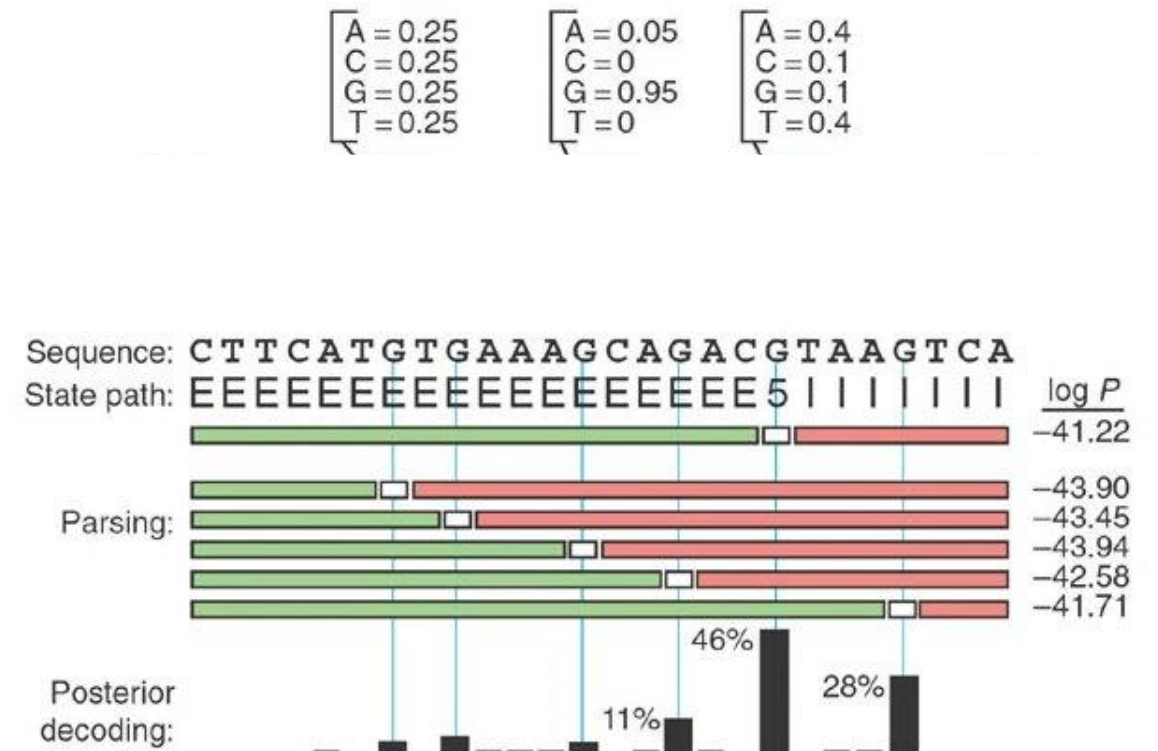**Multiple Sequence Alignment using Multiple Sequence Comparison Log-Expectation**
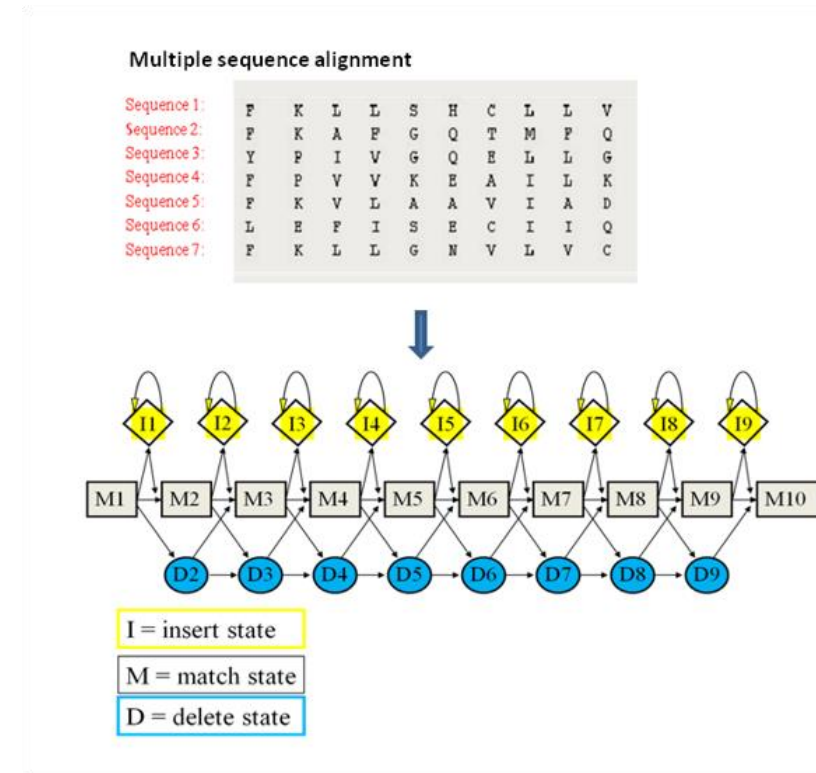
# How does MUSCLE work?

# What is HMM and how is it used?



Hidden Markov Model

# How is HMM utilized?



**Searching through sequences and identifying protein families and generating profiles from MUSCLE.**

# What are Family Based Resource Groups?

Techniques that group together protein sequences or protein domains into evolutionary families.

Some utilize techniques like sequence clustering, while others use significant manual curation.

Ex. PFAM, ProDom, PANTHER, SMART, PhyloFacts

# What is **Pfam**



**PFAM is a database that allows you to analyze proteins and find families using HMM to detect homology.**

# What can you do with PFAM?

# What can you do with PFAM continued?

# What is InterPro?

Interpro is a large secondary protein database that incorporates many different databases

# How to use Interpro?



**Search by sequence or Ascension -> Interpro entry**

# How can SMART be used?



**SMART is a database that is used in the analysis of protein domains within protein sequences.**

**Also uses Hidden Markov models to best sort through the database.**

# SMART vs PFAM- when to use what?

**SMART**

- Specializes in extracellular, signaling, and chromatin associated domains

- Significant manual curation (>1200 manually curated models)

- Exclusive annotation

**PFAM**

- Many more proteins in the database

- Novel sequences are classified into families

- No longer exists (incorporated into InterPro)

# What is Gene Ontology?

Ontology: A set of well-defined terms with well-defined relationships

Gene ontology is a way of categorizing and organizing data to facilitate data retrieval.

# Categories of Gene ontology

**Biological process:** the objective that the gene/protein contributes to.

**Molecular function:** The role/biochemical activity of a gene product.

**Cellular component:** Place in the cell where a gene/protein is active.

# Summary

1) Domains are conserved, functional units of protein

2) Domains can be uncovered and analyzed by various homology and non-homology directed methods.

3) Domain analysis allows us to infer protein function and better classify protein families.

# Dr. Jason Shepherd & the Arc protein

**Molecular function of Arc protein in long-term memory formation**

# Questions?

The Neuronal Gene *Arc* Encodes a Repurposed Retrotransposon Gag Protein that Mediates Intercellular RNA Transfer

Pastuzyn et al., 2018

# How is information stored in the brain?



Retrotransposons store long-lasting information along with genetic memory.

# What is a retrovirus?



Retroviral Genome Structure

| LTR | Gag | Pol | Env | LTR |

Gag → Viral capsid (outer shell)
Pol → Reverse transcriptase & other enzymes
Env → Protective lipid envelope

Gag and Env: form capsid shell
Pol: encodes enzymes

# What is a synapse?



A synapse is the gap between neurons that allow movement of information in the brain.

# What is Arc?



MA = Matrix

CA-NTD = Capsid N-terminal domain

CA-CTD = Capsid C-terminal domain

Arc is a neuronal gene that is important for memory and synaptic plasticity regulation.

# What are the evolutionary origins of Arc?



Retroviral Gag domain

Ty3/gypsy transposons

# Does the Arc protein form capsids?



Arc proteins self-assemble into virus-like capsids.

# Is CTD required for capsid formation?



Capsids are not formed without CTD.

# Does Arc bind and encapsulate mRNA?



Arc binds and encapsulates mRNA, protecting it from degradation.

# Is RNA required for proper capsid assembly?



Removing RNA bases decreased proper capsid formation.

Addition of mRNA increased proper capsid formation.

# Is Arc mRNA found in extracellular vesicles (EVs)?



Arc mRNA is found in EVs.

# Can Arc transfer mRNA with capsids or EVs?

# Are capsids required for neuronal uptake?

# Is mRNA transferred by Arc available for translation?

# Is mRNA transferred by Arc available for translation?



Translation is occurring after being transferred to neurons.

# How does Arc work in the brain?



Like a retrovirus!

# Summary



Arc shares properties of the retroviral Gag protein



Arc can form stable virus-like capsids



These capsid structures allow for mRNA transfer from neuron to neuron.

# Questions?

# Citations

Xiong J. Protein Motifs and Domain Prediction. In: *Essential Bioinformatics*. Cambridge University Press; 2006:85-94.

Eddy, S. What is a hidden Markov model?. *Nat Biotechnol* **22**, 1315–1316 (2004). https://doi.org/10.1038/nbt1004-1315

*Finn, R.D., Mistry, J., Tate, J.G., Coggill, P.C., Heger, A., Pollington, J.E., Gavin, O.L., Gunasekaran, P., Ceric, G., Forslund, K., Holm, L., Sonnhammer, E.L., Eddy, S.R., & Bateman, A. (2007). The Pfam protein families database. Nucleic Acids Research, 38, D211 - D222.*

*Punta, Marco et al. "The Pfam protein families database." Nucleic Acids Research 40 (2011): D290 - D301.*

UniProt Consortium. The universal protein resource (UniProt). *Nucleic Acids Res*. 2008;36(Database issue):D190-D195. doi:10.1093/nar/gkm895

Blum M, Chang HY, Chuguransky S, et al. The InterPro protein families and domains database: 20 years on. *Nucleic Acids Res*. 2021;49(D1):D344-D354. doi:10.1093/nar/gkaa977
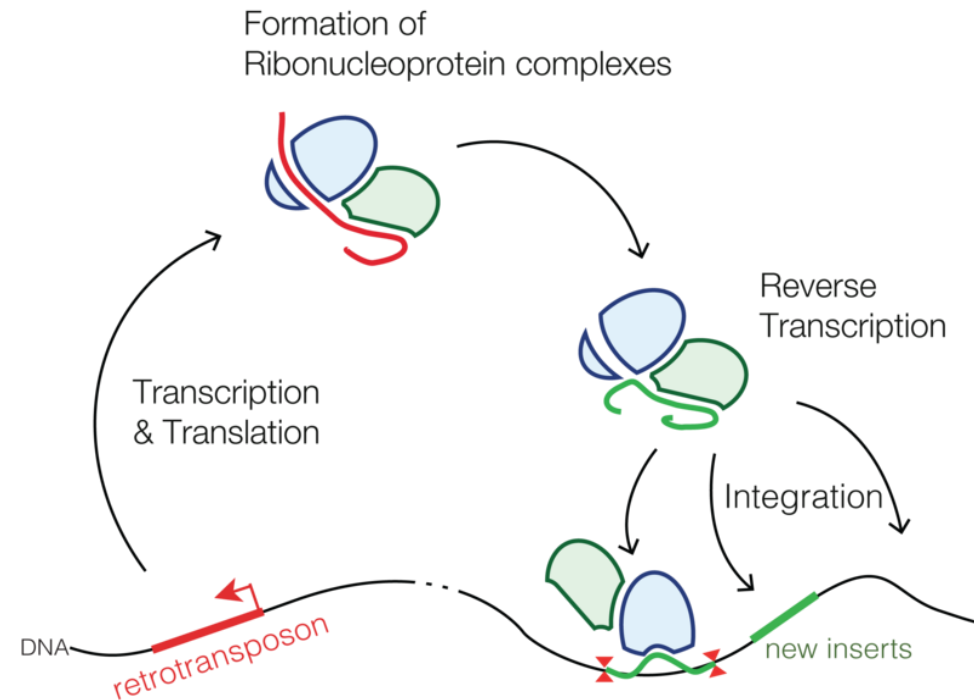
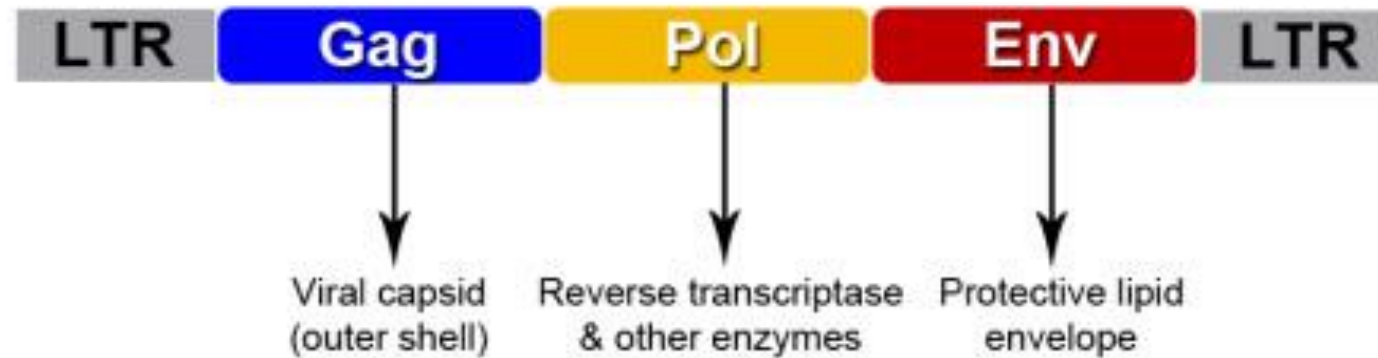Pastuzyn ED, et al., The Neuronal Gene Arc Encodes a Repurposed Retrotransposon Gag Protein that Mediates Intercellular RNA Transfer. Cell. 2018 Jan 11;172(1-2):275-288.e18. doi: 10.1016/j.cell.2017.12.024.

# Image links

- https://static.wixstatic.com/media/309e2a_1fccf0f30277416b8ed317c7548b3401~mv2.png/v1/fill/w_1000,h_667,al_c,q_90,usm_0.66_1.00_0.01/309e2a_1fccf0f30277416b8ed317c7548b3401~mv2.png

- https://microbenotes.com/wp-content/uploads/2020/10/Alternative-Splicing.jpeg

- https://www.ebi.ac.uk/training/online/courses/protein-classification-intro-ebi-resources/wp-content/uploads/sites/96/2020/07/figure7.png

- https://www.ncbi.nlm.nih.gov/Structure/cdd/docs/images/cd00400_hierarchy_tree.png

- https://i.ytimg.com/vi/6OMo9M-nnuk/maxresdefault.jpg

- https://i0.wp.com/pediaa.com/wp-content/uploads/2019/06/Difference-Between-Motif-and-Domain-in-Protein-Structure-Comparison-Summary.jpg?resize=475%2C600&ssl=1

- https://www.ebi.ac.uk/training/online/courses/interpro-functional-and-structural-analysis/wp-content/uploads/sites/32/h5p/content/5/images/image-65b9130f47105.png

- https://www.sciencedirect.com/science/article/pii/S0092867417315040

- https://clarkesworldmagazine.com/koboldt_02_16/

- https://en.wikipedia.org/wiki/Retrotransposon