



A Machine Learning Approach for Prediction of Protein Structure and Function

Tom Wilson

HCI 5903 Capstone

2019.04.15

Objectives



Become familiar with traditional methods of protein structure and function determination



Recognize the utility of protein characterization for potential medical therapies



Understand the computational challenge of brute force protein structural prediction



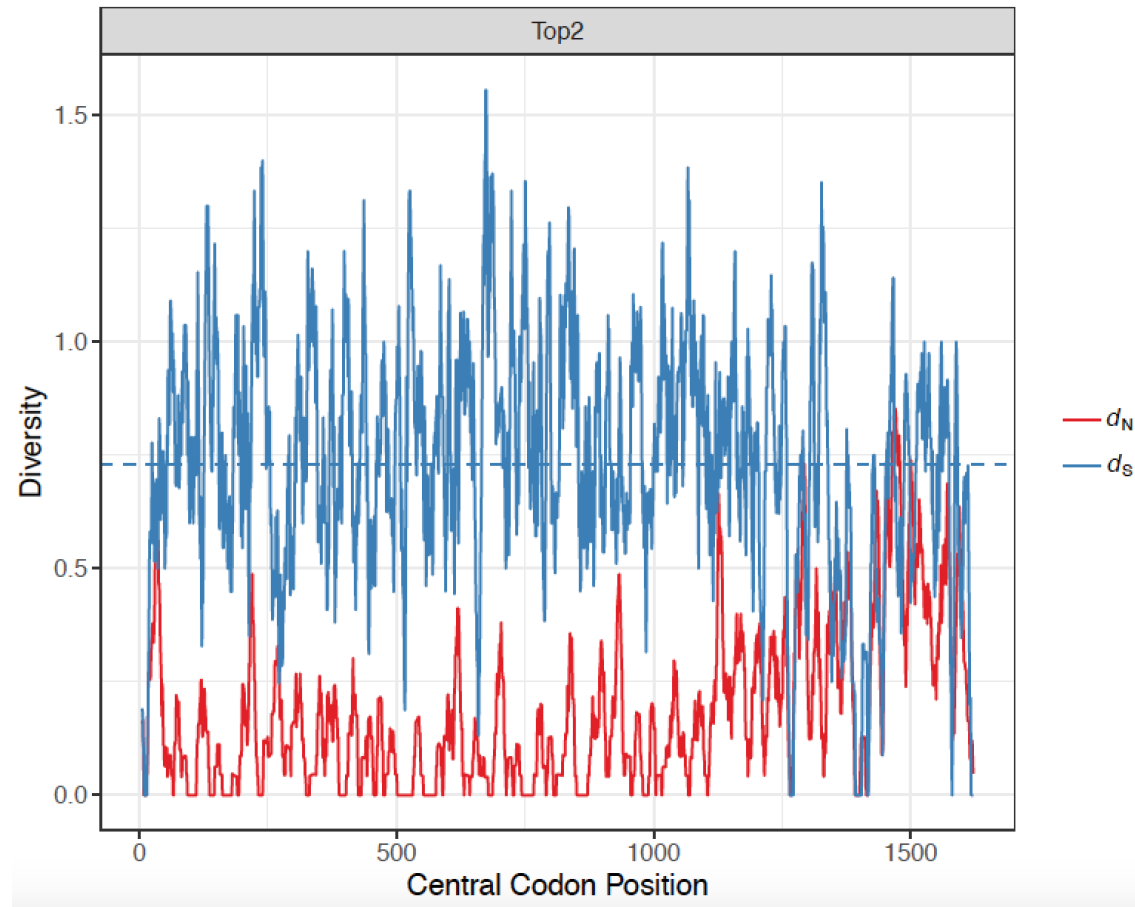
Understand the role of machine learning in protein structure and function prediction

Project Inspiration

- Deweese Laboratory
- Topoisomerase II isoform comparison
- C-terminus is not conserved between the two isoforms¹
- C-terminus has been resistant to structural determination¹
- Structure and function determination of the C-terminus for each isoform may allow targeted drug therapies with reduced adverse reactions
 - Secondary leukemias

Differences in C-termini of Topoisomerase II¹

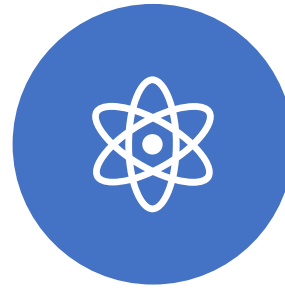
1. Human_TOP2A_P11388 R A A T K T K F T M D L D S D E D F - - - - - S D F D E K T - - - - - D D E D F V P S D A S - - - - - P P K T K T S P K L S N K E L K P Q - - - - - K
13. Human_TOP2B_Q02880 A A A E R P K Y T F D F S E E E D D - D A D D D D D D N N D L E E L K V K A S P I T N D - G E D E F V P S D G L D K D E Y T F S P G K S K A T P E K S L H D K K S Q D F G N L F S F P S Y



Traditional Methods of Protein Characterization^{2,3}



X-ray crystallography



NMR Spectroscopy



Activity assays



Agonists vs antagonists

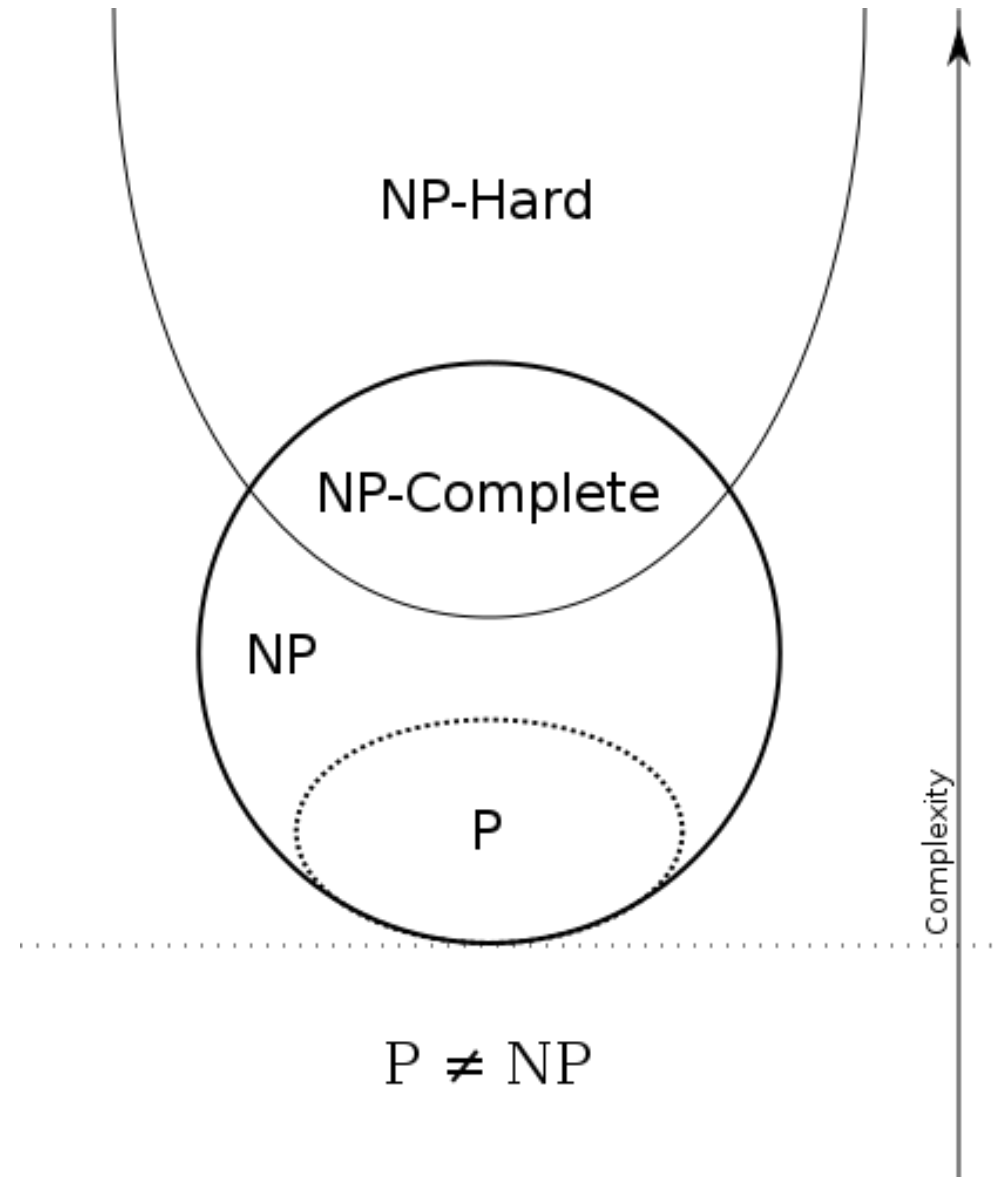
Computational Analysis⁴

- Computational protein structure prediction (PSP) typically utilizes enumeration or searching strategies with optimization which involve enormous probability spaces
- $$U = \sum_b K_b^{bond} (l_b - l_b^0)^2 + \sum_a K_a^{angle} (\theta_a - \theta_a^0)^2 + \sum_t K_t^{torsion} (1 - \cos[n_t(\varphi_t - \varphi_t^0)]) + \sum_{i>j} K_{ij}^{nonlocal} f(r_{ij} / r_{ij}^0)$$
- related optimization algorithm for the current method of PSP
- Algorithm is NP-Hard

P, NP, and Computational Complexity⁵

- $P = NP$ vs $P \neq NP$
- P is the set of all problems that can be solved in polynomial time
- NP is the set of all problems that must be solved in non-deterministic polynomial time
- NP -Hard is as difficult as the most difficult NP problem in the set
- N -Hard problems become computationally intractable if not optimized or abstracted

P, NP, and Computational Complexity^{5,9}



Levinthal Paradox

- “In theory a protein is expected to require exponential time to fold, given an arbitrary starting configuration, whereas in practice proteins are observed to fold within seconds to minutes, independent of size... [where] exponential-time folding is expected because of the exponential size of the protein’s conformational space”⁴
- “Exhaustive search of a protein's conformational space is clearly not a feasible algorithmic strategy. The number of possible conformations is exponential in the length of the protein sequence, and powerful computational hardware would not be capable of searching this space for even moderately large proteins”⁶

Machine Learning

- K-modes Attribute Clustering Algorithm⁷
 - Learning algorithm for non-numeric data
 - take protein sequences and predict a structure in polynomial time based on probability and the mutual similarity of sequences between different groups of species
- Multiple Sequence Alignment
- Complexity⁸
 - $O(knp^2t)$



K-Modes Algorithm^{7,8}

- Data-driven, machine learning algorithm for clustering similar data
- Clusters based on how dependent one attribute is on another
- Machine learning has immense potential for use in the biological and medical sciences, especially for process optimization and probabilistic prediction

Multiple Sequence Alignments¹⁰

- Alignment of a protein's sequence across multiple organisms
 - Takes into consideration substitution, insertion, and deletion events
- Necessary for probabilistic purposes
 - More sequences allows for better resolution of results

```
Q9VEC8_DROME/3-75      ITIKV..LKGK....DCTI.EVAPTSTILEVKHQIEAEL.Q...ISATNQ.KLLLLGRPL..NNEQTIASYPNIKEG.TKLNLVVVKP
UBL4A_MOUSE/3-74      LTVKA..LQGR....ECSL.QVAEDELVSTLKHLSVDKL.N...VPVRQQ.RLLFKGKAL..ADEKRLSDY.NIGPN.SKLNLVVVKP
R7SL1_ARATH/3-74      VYIDT..ETGS....SFSI.TIDFGETVLEIEKEIEKESQ.G...IPVSKQ.ILYLDGKAL..EDDLHKIDY.MILFE.SRLLRLSPD
Q9ZQZ4_ARATH/3-74      MTVEN..ESGS....TFSI.DIGLQDTVLTFKKKIEMTQ.R...IPVSRQ.TIFFQKGKL..EDHLDFEFW.DILQN.PLLHLSISPD
Q9SYF2_ARATH/71-142    IFIKT..LTGR....TNTY.EVKGSDTIRELKAKAHEEKE.G...IPVEQQ.RLIFQGRVL..EDSKAISDY.NIKHE.STLHITHQPC
RUB1_YEAST/3-74        VKVKT..LTGK....EISV.ELKESDTIYVHIKELLEKEE.G...IPPSQQ.RLIFQKGKL..DDKLTVTDA.HLVGE.MQLHLVLTIR
NEDD8_HUMAN/3-74      IKVKT..LTGK....EIEI.DIEPTDKVERIKERVEEKE.G...IPQQQ.RLIYSQKQM..NDEKTAADY.KILGG.SVLHLVLIALR
RUB3_ARATH/3-74        IKVKT..LTEK....QIDI.EIETDITIERIKERIEEKE.G...IPPVHQ.RIVYTGKQL..ADDLTAKHY.NLREG.SVLHLVLIALR
B0BLT8_XENTR/30-101    LFJET..LTGT....CFEL.RVSPYETVASVSKIQRLG.G...IPVAQQ.HLIWNMEL..EDECLSLDY.NISEG.CTLKRMVLMAR
Q9T284_CAEL/185-255    FNV1..LHGS....RIPM.ELDSLDITTYTGMALLKHG.G...LALHRQ.RMLNDKEL..QYNQTLSEY.KIKDG.SVIRKQHTDEK
P4KG4_ARATH/37-108     IYLT..LPGS....VIPM.RVLESDESIESVKLRIQSYR.G...FVVRNQ.KLVFSGREL.ARSNSNMRYD.GVSEG.NILHLVLKLS
P4KG2_ARATH/37-108     VFLS..VSGS....TMPM.LILESDDISAEVKLRIQTCN.G...FVVRNQ.KLVFSGREL.ARNASRVKDY.GVTGG.SVLHLVLKLY
UBQ8_ARATH/81-152      IFVQT..LTGK....TITL.EVKSSDTIDNVKAKIQDKE.G...ILPRQQ.RLIFAGKQL..EDGRTLADY.NIQKE.STLHLVLRLC
UBQ8_ARATH/471-549     IFVKTFSFSGSETPTCKTITL.EVESSDTIDNVKVIQHKV.G...IPLDRQ.RLIFGGRVL..VGSRTLDDY.NIQKG.STIHQFLQQR
UBQ8_ARATH/240-316     IFVKN..LPYNSFTGENFIL.EVESSDTIDNVKAKIQDKE.R...IPMDLH.RLIFAGKPL..EGGRTLADY.NIQKG.STLYLVTRFR
UBQ8_ARATH/157-235     IFVST..FSGKNFTSDTLTL.KVESSDTIDNVKAKIQDKE.G...LRPDHQ.RLIFHGEELFTEDNRTLADY.GIRNR.STLCLALRLR
UBQ8_ARATH/395-466     IFVKL..FGGK....IITL.EVLSDDTIKSVKAKIQDKV.G...SPPDQQ.ILLFRGGQL..QDGRTLGDY.NIRNE.STLHLFFHTR
R5RDZ3_9PROT/23-94     IFVKT..LTGK....HITL.EVEPTDRIEDVNAKIQDKE.G...IPPDQQ.RLIFAGKQL..EDGNTLQDY.SIQKD.STLHLVLRLR
UPL5_ARATH/97-169      IFVRM..MSGG....KTIVIAEKYDTVEKLHQRIEWT.K...IPALEQ.RVIYKGKQL..QRENSLTY.SIEQD.ASLQLVARMQ
P91050_CAEL/144-214    LCIAV.SMPGR....LFSI.GANKMESVEQLKMKIECQT.G...IPRTKF.WLRHLHGKPL..YDDKLADY...KWD.STVELLVAS
Q9SV28_ARATH/11-79     FFVRL..LDGK....SLTFSFSSPLAYGEQIKQRIEQT.K...IPTHLQ.RLISGGYQI..SDGSAISQP....D.ATVNLVLISLR
UHRF1_MOUSE/3-76       IQVRT..MDGK....ETHVTNSLSRLTKVQELRKKIEEVF.H...VEPQLQ.RLFYRGKQM..EDGHTLFDY.DVRLN.DTIQLLVRSQ
OASL_HUMAN/436-507     VFVKN..PDGG....SYAY.AINPNSFILGLKQQIEDQ.G...LPKKQQ.QLEFQGGVL..QDNLGLGY.GIQDS.DTLILSKKKG
YKA4_CAEL/343-415      ILIKLFMNDT.....EKNTYASLEDTVAKFKVDHFTNLANQVIRLIYQQQL.REDHRTLEEY.GLQPG.SVICHISTT
H9L0X6_CHICK/393-463   VLVK..DSNK....TTYV.TVRPTDTVKQLKQQIYACQ.H...VPVEQQ.RLTYETKEL..ENHHTLEHY.HVQPR.STIYLLRLR
ISG15_HUMAN/84-155     ILVRN..NKGR....SSTY.EVRLTQTVAHLKQQVSGLE.G...VQDDLF.WLTFEGKPL..EDQLPLGEY.GLPL.STVFMNLRQ
ISG15_BOVIN/81-152     ILVRN..DKGR....SSPY.EVQLKQTVAKLQKQVCQKE.R...VQADQF.WLSFEGRPM..DDEHPLEY.GLMKG.CTVFMNLRQ
ISG15_MOUSE/82-153     ILVRN..ERGH....SNYI.EVFLTQTVDTLKKKVSQRE.Q...VHEDQF.WLSFEGRPM..EDKELGEY.GLKPQ.CTVIKHLRLR
DSK2A_ARATH/20-91      VNVRC..SNGT....KFSV.TTSLDSTVESFKELIAQNS.D...VPANQQ.RLIYKGRIL..KDDQTLLSY.GLQAD.HTVHMRGVF
DSK2_SCHPO/6-77        LTICA..ANDQ....KYAV.TVDSSESVLAKELIAPVA.D...IEKERQ.RLIYAGRVL..KDEESLRTY.KIQDG.HSIHLVKTG
DSK2_YEAST/5-75        IHIKS....GQD....KWEV.NVAPESTVLQFKEAANKAN.G...IPVANQ.RLIYSGKIL..KDDQTVESY.KIQDG.HSVHLVKSQ
G5GEG66_CAEL/10-81     VHVKS...PSN....KYDV.EIAADASVSELKDKVLVFP.P...TANKEQVCIIYTGKIL..KDEETLTQH.KIADG.HTVHLVIRNQ
UBQL1_HUMAN/39-109     VTVKT...PKE....KEEF.AVPENSSVQKKEEISKR.F.K...SHTDQL.VLIFAGKIL..KDDQTLQSGH.GIHDG.LTVHLVKTQ
Q9VWD9_DROME/11-81     VVVKT...PKD....KKTV.EVDEDSGGLKDKILVAQKF.E...AEPEQL.VLIFAGKIM..KDDTDLQMH.NIKDN.LTVHLVKAIP
HERP1_HUMAN/12-89      LLVKSPPNQRRH....DLEL.SGDRGWSVGHLLKAHLRSVYPE...RPRPEDQLIYSGKLL..LDHQCRLDLPQKRRVHLVLCNVK
YB92_SCHPO/6-79        IRVTTVDQ.....KVGIFQVPRTKTIVLEKELIAVTF.E...APADRL.KLIHAGRVL..RNETPLEEIHLDATDLVTFHLVIAVF
Q9VNB2_DROME/165-235   LRISSTMTDVK.....L.PVYSKDTVGQCKKKLQAAE.G...VDACQ.RWFYSGKLL..GDKVPIDEC.SIHQG.YVQVIVNTE
Q9VS82_DROME/5-75      LKVKT..LDAR....IHEF.SIDNELTIRQFKDQIAEKT.N...IAAENQ.RIYQGRVL..VDDQKVKEY.DV.DG.KVLHVAERPP
UBIM_HUMAN/3-72        LFRVA...QE....LHTF.EVTGQETVAQIKAHVASLE.G...IAPEDQ.VVLLAGAPL..EDEATLGQC.GVEAL.TTLEVAGRML
Q18231_CAEL/3-69       IFLLG..LDNT....THTL.DVDASTTLSAIGVIG.....AGEEF.SISYSGKVL..SEELTLGEC.QIESL.STLSVNGRLL
PRKN_MOUSE/3-74        VFVRF..NSSY....GFPV.EVDSDTSLQLKEVVAKRQ.G...VPADQL.RVIFAGKEL..PNHLTQVNC.DLEQQ.SVIVHVQRR
Q9VXF9_DROME/3-75      ITVTT..SDDK....VFCL.DVAQDLELENLALCAMEI.G...AEVSQI.AVIFNGREL.SSDKQTLQCC.GVGDG.DFIMLERRR
SUMO1_MOUSE/22-95      IKIKVVGQDSN....EVHF.RVKYGTSMALKKSYADRT.G...VAVNSL.RLFDGRRIL..NDDTPTKTL.EMEDD.DVIEVYQEQ
SUMT3_YEAST/24-96      IKLVIGQDSS....EIH.F.KVKMTTHLKKLKECYCQRQ.G...VPMNSL.RLFEQGRI..ADNHTPEKL.GMEEI.DVIEVYQEQ
PMT3_SCHPO/36-109      INLVK.SDGSS....EIFF.KIKKTTPLRRMLDEAFAKRQ.G...KEMDSL.RFLYDGIRI..QADQTPEDL.DMEDD.DQIEAHREI
SUMO1_ARATH/18-91      INLVKVGQDNN....EVFF.KIKKTTESFKLMKIYCARQ.G...KSMNSL.RFLVDGERI..RPDQTPAEL.DMEDG.DQIEAVLEQL
D7LEX6_ARALL/78-148    INLVKVGQDGN....EVFF.RIKRSTQKLKLMNAYCDRQ.S...VDMNSI.AFLFDGRRIL..RAEQTPEDL.DMEDG.DEIDAMLHQT
Q20899_CAEL/22-96      .TVK...FPSK....QFTV.EVDRTETVSSLSKDKHIVE.N...TPIKRM.QLYYSGLIEL.ADDYRNLINEY.GISEF.SEIVVFLKSI
SF3A1_HUMAN/709-788    ITVSSVMQGVK....QIVV.EMDKKETVSLKNRIEQT.E...VLTNRQ.VLLFKGMEL.KDNKNRITDC.GINSN.AKITMNVKMS
UBP6_SCHPO/5-75        VQVPMQDKTEWKLNGQVLVLTPLTDQVSVIKVKIHEAT.G...MPAGKQ.KLYEGIFI..KDSNSLAYY.NMANG.AVHILALKER
UBFD1_HUMAN/90-157     IAIR...WQGK....KYDL.EIEPNETGSTLKHQLYLSLT.Q...VPPERQ.KVIVKGQGL..KDDVLLGSV.GIKPN.ATLLMMGTAG
Y2010_ARATH/15-85      I....WNKT....KHDV.KFPLDSTGSSELKQKHSIT.G...LPPAMQ.KVMYKGLV..PEDKTLREI.KVTSG.AKIMVVGSTI
MDY2_YEAST/76-151      LTVK...FGGK....SIPL.SVSPDCTVKDLKSQLPIT.N...VLPRGQ.KLIFKGKVL..VETSTLKQS.DVSGG.AKLMLMASQG
YQ77_SCHPO/3-71        LTLKKIQAPKF....SIEH.DFSPSDTTLQIKQHLSIEKA..SHISEI.KLLKKGKVL..HDNLFLSDL.KVTPANSTITVMIKPN
RAD23_YEAST/4-75       LKFS...CRGN....VIAL.SFNENDTVLDAKEKLGQEI.D...VSPSLI.KLLYKGNL..SDDSHLQDV..VKNE.SKIMCLIRQD
RAD23_DROME/3-76       LTFKN..FKKE....KVPL.DLEPNTIETKTTLKQSI....SCEESQIKLIYSKVL..QDSKTVCSE.GLKDQ.DQVFMVSQK
Q9V3W9_DROME/3-76      LSIRM..LDQR....TITL.EMNESQEVRAKQKLGNL.P.EV.AMPAENL.QLIYSGRIM..EDAMPLSEY.RIAED.KIIVLMGKKK
Q23451_CAEL/5-78       ITIKN..LQQQ....TFTI.EFAPEKTVLELKKKIFEEER.GP.EYVAEQ.KLIYAGVIL..TDORTVGSY.NVDEK.KFIVVMLTRD
RD23A_MOUSE/5-79       VTFRT..LTQV....NFNL.ELNEDQTAIEVKALVASEK.GD.DYAPELQ.KLIYNGKIL..DDSVKVGVE.GFDS.SKIVVMLSKR
RHP23_SCHPO/3-75       ITLKT..LQQQ....TFKI.RMEPDDETVKVLKEKEIAEK.GRDAFPVAGQ.KLIYAGKIL..SDDVPIKEY.HIDEK.NFVVMVTKRA
RAD23_DICDI/3-74       LTFKN..LQQQ....KFVISDVADTKISELKEKIQTOQ.N...YEVERQ.KLIYSGRIL..ADDKTVGEY.NIKEQ.DFIVCMVSRP
                        VTIKN..INKE....IYVF.EVNGDLTVAELKNLISEKH.N...QTPSWQ.TLIYSGRIL..EDKRTLESY.NITDS.GFIVMMIKKP
```

Attributes

- Once the multiple sequence alignment data has been obtained, it needs to be parsed through and manipulated for prediction calculations
 - Array manipulation
- Attributes within the data array
 - An attribute in this case is a column of amino acids within the aligned sequence array

```
data = [  
    ['C', 'A', 'R', 'C', 'A', 'W', 'A', 'A'],  
    ['C', 'G', 'K', 'C', 'G', 'Y', 'G', 'G'],  
    ['C', 'N', 'M', 'C', 'N', 'F', 'N', 'N'],  
    ['C', 'D', 'I', 'C', 'D', 'V', 'D', 'D'],  
    ['C', 'A', 'L', 'C', 'A', 'H', 'A', 'A'],  
    ['C', 'G', 'R', 'C', 'G', 'Q', 'G', 'G'],  
    ['C', 'N', 'K', 'C', 'N', 'E', 'N', 'N'],  
    ['T', 'D', 'M', 'T', 'D', 'P', 'D', 'D'],  
    ['T', 'A', 'I', 'T', 'A', 'W', 'A', 'A'],  
    ['T', 'G', 'L', 'T', 'G', 'Y', 'G', 'G'],  
    ['T', 'N', 'R', 'T', 'N', 'F', 'N', 'N'],  
    ['T', 'D', 'K', 'T', 'D', 'V', 'D', 'D'],  
    ['S', 'A', 'M', 'S', 'A', 'H', 'A', 'A'],  
    ['S', 'G', 'I', 'S', 'G', 'Q', 'G', 'G'],  
    ['S', 'N', 'L', 'S', 'N', 'E', 'N', 'N']  
]
```

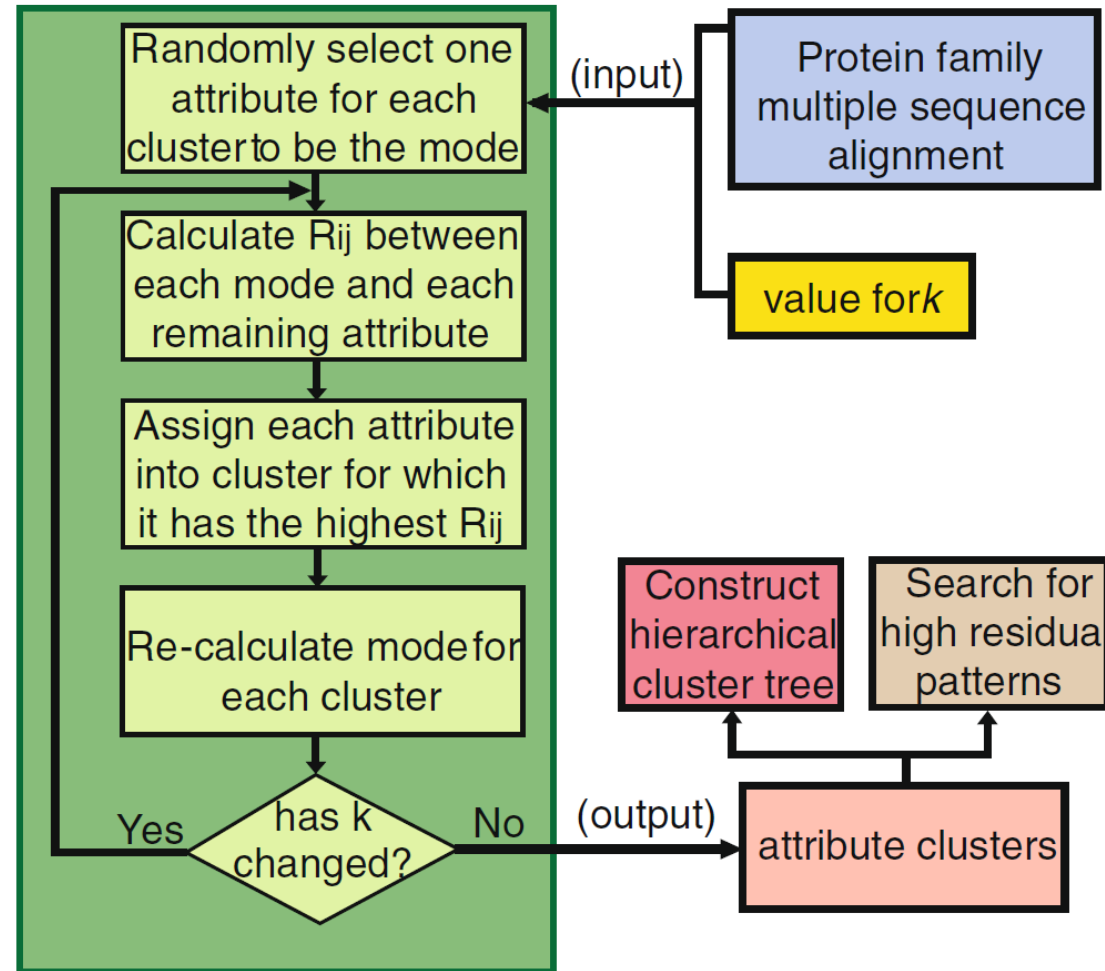
Interdependence and Mutual Information^{7,11,12}

- Mutual information is calculated to determine the interdependency relationship between two attributes
- Mutual information is normalized by the joint entropy (Shannon entropy) of the sequences as well to yield a more interpretable result
- $I(X_i, X_j) = \sum_{x \in X_i} \sum_{y \in X_j} p(x, y) \log\left(\frac{p(x, y)}{p(x)p(y)}\right)$
 - Mutual Information
- $H(X_i, X_j) = - \sum_{x \in X_i} \sum_{y \in X_j} p(x, y) \log(p(x, y))$
 - Shannon Entropy
- $R_{ij} = \frac{I(X_i, X_j)}{H(X_i, X_j)}$
 - Normalized Mutual Information
 - Interdependency Redundancy

Algorithmic Approach¹³

- Python, Scikit-learn, and Numpy
- Once a multiple sequence alignment has been input into the program, a number of steps take place in a loop-wise fashion until the prescribed value for k has been reached
 - algorithm will run for a number of iterations from $k - 1$ to $k = 2$
- Output from each iteration will be stored in a file which is then used to construct a cluster tree graph

Algorithm⁷



Code

- Hard-coded test sequence
 - In the future, would allow for upload of code via file parsing

```
data1 = np.array([
    ['C', 'A', 'R', 'C', 'A', 'W', 'A', 'A'],
    ['C', 'G', 'K', 'C', 'G', 'Y', 'G', 'G'],
    ['C', 'N', 'M', 'C', 'N', 'F', 'N', 'N'],
    ['C', 'D', 'I', 'C', 'D', 'V', 'D', 'D'],
    ['C', 'A', 'L', 'C', 'A', 'H', 'A', 'A'],
    ['C', 'G', 'R', 'C', 'G', 'Q', 'G', 'G'],
    ['C', 'N', 'K', 'C', 'N', 'E', 'N', 'N'],
    ['T', 'D', 'M', 'T', 'D', 'P', 'D', 'D'],
    ['T', 'A', 'I', 'T', 'A', 'W', 'A', 'A'],
    ['T', 'G', 'L', 'T', 'G', 'Y', 'G', 'G'],
    ['T', 'N', 'R', 'T', 'N', 'F', 'N', 'N'],
    ['T', 'D', 'K', 'T', 'D', 'V', 'D', 'D'],
    ['S', 'A', 'M', 'S', 'A', 'H', 'A', 'A'],
    ['S', 'G', 'I', 'S', 'G', 'Q', 'G', 'G'],
    ['S', 'N', 'L', 'S', 'N', 'E', 'N', 'N']
])
```

```
num_rows = len(data)
num_cols = len(data[0])
```

```
# Random attribute selection for first pass
rand = np.random.randint(0, num_cols)
```

Code

- Computing the length of the data array to assign to k
- Random attribute selection for first iteration

Code

- Looping through the data array to calculate the interdependency redundancy between attributes

```
# First, random clustering
for i in range(num_cols):
    if(rand != i):
        # A[:,i] returns column at index i of array A for
        numpy arrays
        rii = nmis(data[:, rand], data[:, i])
        RAND_OUTPUT_DICT[rand, i] = rii
    if(rii > max_R):
        max_R = rii
        # Index of the attribute with the highest Rii
score
        starting_attribute = i
```

Code

- Calculating the Normalized Interdependency Redundancy (NIR) and Average Normalized Interdependency Redundancy (ANIR)
- Performing clustering again without random selection for remaining iterations

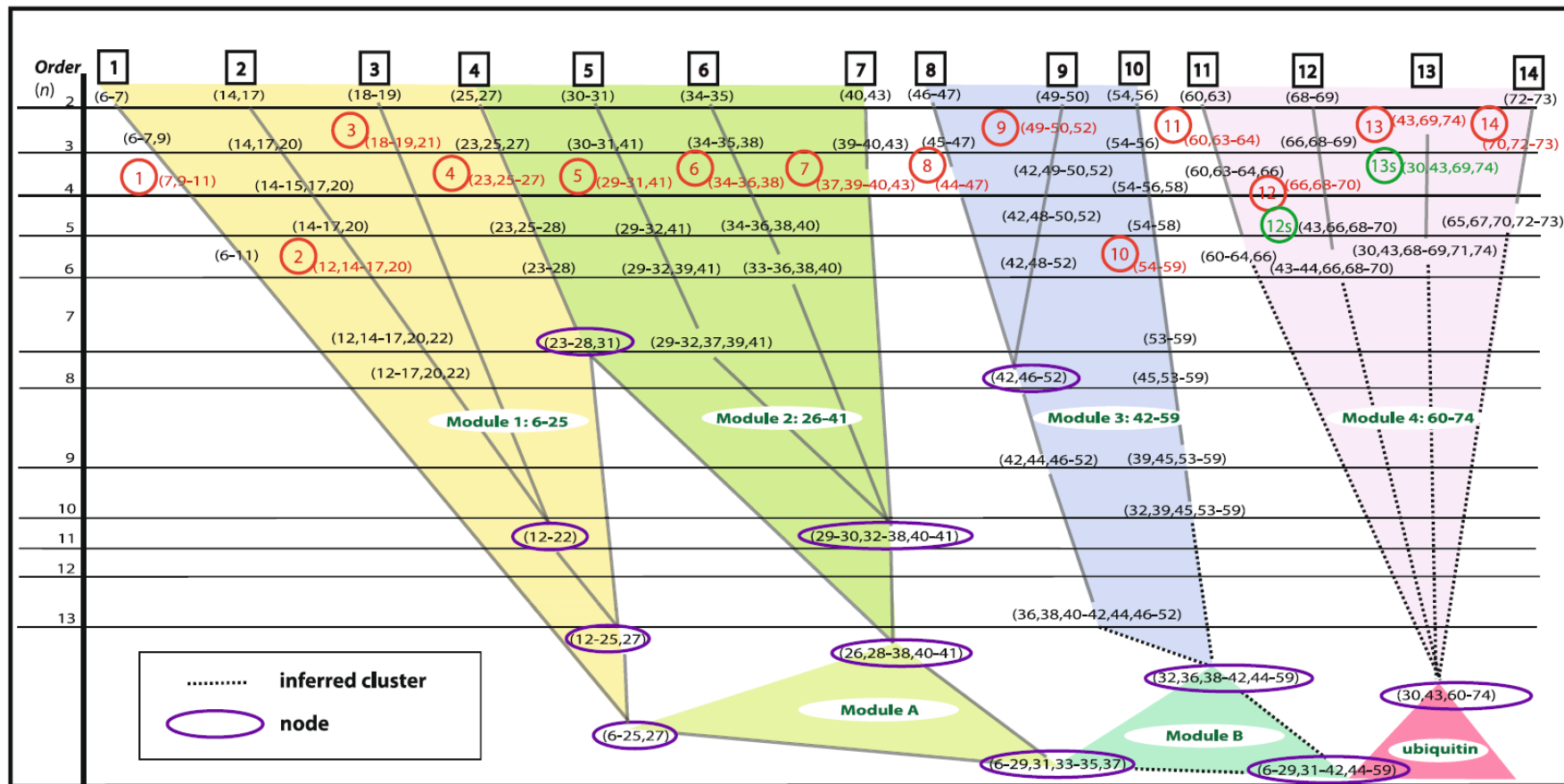
```
for cluster in clusters:  
    NIR += rii[cluster]  
  
ANIR = NIR / len(clusters)
```

```
# Taking output of random clustering to form a true clustering  
max_R = 0  
for i in range(num_cols):  
    if(starting_attribute != i):  
        rii = nmis(data[:, starting_attribute], data[:, i])  
        FIRST_RUN_DICT[starting_attribute, i] = rii  
        if(rii > max_R):  
            max_R = rii
```

Code Output

```
1 RUN = 1
2 K = 9
3 Clusters: (1, 3);
4 SR(1, 3) = 0.9;
5 SR(mode(1, 3)) = 0.9;
6 ---
7 RUN = 2
8 K = 8
9 Clusters: (1, 3); (5, 6);
10 SR(1, 3) = 0.9; SR(5, 6) = 0.899;
11 SR(mode(1, 3)) = 0.9; SR(mode(5, 6)) = 0.899;
12 ---
```

Cluster Trees⁷



Next Steps

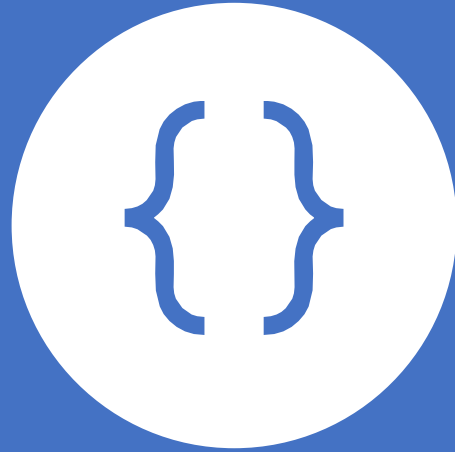
- Verification and repeatability are of the utmost importance if this tool is to be used by other researchers
- Next steps should focus on validating this algorithm against protein sequences with known and verified crystallographic structures
- Algorithmic scalability is also an area ripe for optimization
 - Protein sequence and probability space grows exponentially with the length of the sequence
 - Intractable program runtime
 - Computational Clusters

Reiteration

- The importance of the scope of this project lies in the utility of proteins as medicinal targets for various therapeutic strategies
- Having a tool that can accurately and efficiently predict protein structure and function from a sequence would be useful for new drug design and targeting techniques as well as understanding the pathophysiology of disease states

References

1. Deweese, J., et al. (2019). The variable C-terminal domain of human type II topoisomerases as a functionally relevant therapeutic target. American Society for Biochemistry and Molecular Biology, Orlando, FL.
2. Carpenter, E. P., et al. (2008). "Overcoming the challenges of membrane protein crystallography." Curr Opin Struct Biol **18**(5): 581-586.
3. Lacapere, J. J., et al. (2007). "Determining membrane protein structures: still a challenge!" Trends Biochem Sci **32**(6): 259-270.
4. Ngo, J. T., et al. (1994). "Computational Complexity, Protein Structure Prediction, and the Levinthal Paradox." The Protein Folding Problem and Tertiary Structure Prediction: 433-506.
5. Knuth, D. E. (1974). "Postscript about NP-hard problems." ACM SIGACT News **6**(2): 15-16.
6. Newman, A. and W. E. Hart (2001). "The Computational Complexity of Protein Structure Prediction in Simple Lattice Models." CRC Press.
7. Durston, K. K., et al. (2012). "Statistical discovery of site inter-dependencies in sub-molecular hierarchical protein structuring." EURASIP J Bioinform Syst Biol **2012**(1): 8.
8. Au, W. H., et al. (2005). "Attribute clustering for grouping, selection, and classification of gene expression data." IEEE/ACM Trans Comput Biol Bioinform **2**(2): 83-101.
9. Esfahbod, B. (2007). N. Euler diagram for P, NP-Complete, and NP-Hard set of problems. Wikimedia Commons.
10. Laboratory, E. M. B. (2019). "Family: ubiquitin (PF00240)." 2019, from <https://pfam.xfam.org/family/PF00240#tabview=tab3>.
11. Wong, A. K. C. and G. C. L. Li (2008). "Simultaneous Pattern and Data Clustering for Pattern Cluster Analysis." IEEE Transactions on Knowledge and Data Engineering **20**(7): 911-923.
12. Wong, A. K. C., et al. (1976). "Statistical analysis of residue variability in cytochrome c." Journal of Molecular Biology **102**(2): 287-295.
13. Pedregosa, F., et al. (2013). "Scikit-learn: Machine Learning in Python." Journal of Machine Learning Research **12**: 2825-2830.



Questions

A Machine Learning Approach for Prediction of Protein Structure
and Function