

Digitally Assessing Protein Properties

Biochemistry Boot Camp 2017
Session #2
Nick Fitzkee
nfitzkee@chemistry.msstate.edu

Protein as Chemicals

- Molecular weight
- Chemical formula (e.g. $C_{274}H_{427}N_{69}O_{93}S_1$)
- Isoelectric point
- Sequence & Residue composition
- Solubility
- Structure
- Concentration/extinction coefficient

→ How do we access this information?

Sequence of GB3

- Primary Structure:

NT-Met-Gln-Tyr-Lys-...-Thr-Glu-**CT**

- More convenient:

```
MQYKLVINGK TLKGETTTKA VDAETAEKAF  
KQYANDNGVD GVWTYDDATK TFTVTE
```

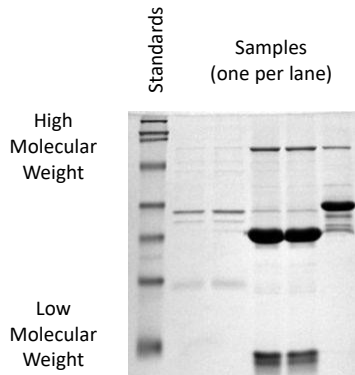
- Can we search this (think Google)?

Website #1: Protparam

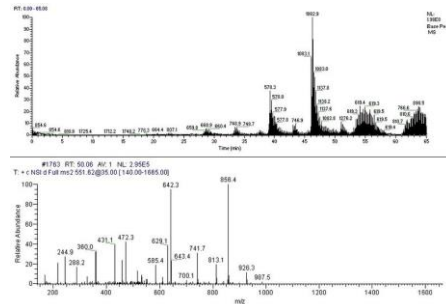
- <http://web.expasy.org/protparam/>
- **Input:** Protein sequence (one-letter codes)
- **Output:** Basic chemical properties
 - Molecular weight
 - Isoelectric point (pI)
 - Extinction coefficient

Molecular Weight

Polyacrylamide Gel Electrophoresis
(SDS-PAGE)



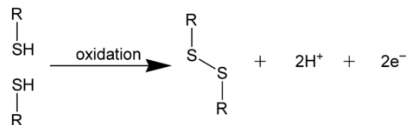
Mass Spectrometry
(ESI-MS, LC-MS)



Sources: en.wikipedia.org/wiki/SDS-PAGE, en.wikipedia.org/wiki/Protein_mass_spectrometry

Residue Composition

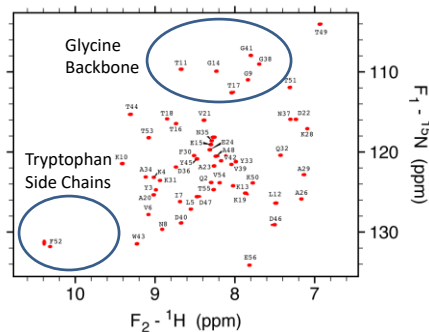
Disulfide Formation
(Cysteine Content)



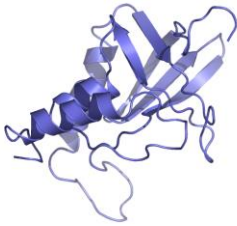
Reducing Agents:

- 2-Mercaptoethanol (BME, 5-10 mM)
- Dithiothreitol (DTT, 1-5 mM)
- Tris-(2 carboxyethyl) phosphine (TCEP, < 1 mM)

Protein ^{15}N HSQC (NMR)

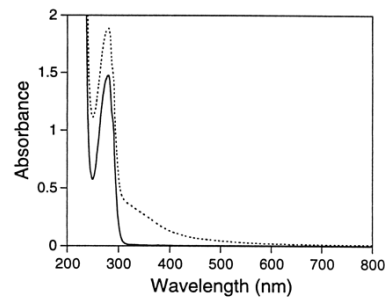


Extinction Coefficient



Tryptophan side chain absorbs light at 280 nm

More absorbance → More protein



If we know the extinction coefficient, we can *estimate* the concentration.

Calculating Protein Concentration

(Beer's Law)

- **UV-Vis:** Absorbance at 280 nm is 0.348 in a 0.3 cm quartz cuvette
– Most cuvettes are 1 cm
- **Protparam:** Extinction coefficient at 280 nm is $9970 \text{ M}^{-1} \text{ cm}^{-1}$
- **Beer's Law:** $A = \epsilon Cl$



Source: www.malvernstore.com

What If My Protein Doesn't Have Trp?

- No Trp means low (no) absorbance at 280 nm
- Protein backbone has intrinsic absorbance at 205 nm
 - See Anthis, N.J. and Clore, G.M. (2013) *Protein Science*.
<http://www.ncbi.nlm.nih.gov/pubmed/?term=23526461>
 - Website: <http://nickanthis.com/tools/a205.html>
- Complications:
 - Protein concentration will need to be quite low, which may introduce dilution errors
 - Many buffers absorb at 205 nm, these can overwhelm the protein signal (even when using a blank)
 - **Solution:** Careful dilution, use water as a blank if possible

Caveats: Extinction Coefficient

- Uncertainty can be as much as 10%
 - Can be worse if your technique is poor!
- Absorbance values need to be between 0.1-1.0 for highest accuracy
 - Estimate your expected A_{280} and dilute if necessary
- **Scattering of aggregates:** If the baseline is not zero at 600 nm, you are probably not getting an accurate value!
- DNA, other impurities or other compounds may artificially increase absorbance at 280 nm

Think and Discuss

The extinction coefficient can be calculated from primary structure alone. Why is this important?

Website #2: NCBI Databases

- <http://www.ncbi.nlm.nih.gov/sites/gquery>
- **Input:** Gene names, organisms, authors, etc.
- **Output:** Curated summary of research
 - Accepted DNA and protein sequences
 - Summaries of associated diseases
 - Recent research papers

NCBI Tricks #1

- Database restriction

srcdb refseq [prop]	Only search reference sequences
srcdb pdb [prop]	Only search the PDB

- Journal restriction

1998:2003 [dp]	Dates from 1998-2003
fitzkee_nc [auth]	Author name is Fitzkee, N. C.
j am chem soc [jour]	Journal name is JACS (need to know abbreviation)

NCBI Tricks #2

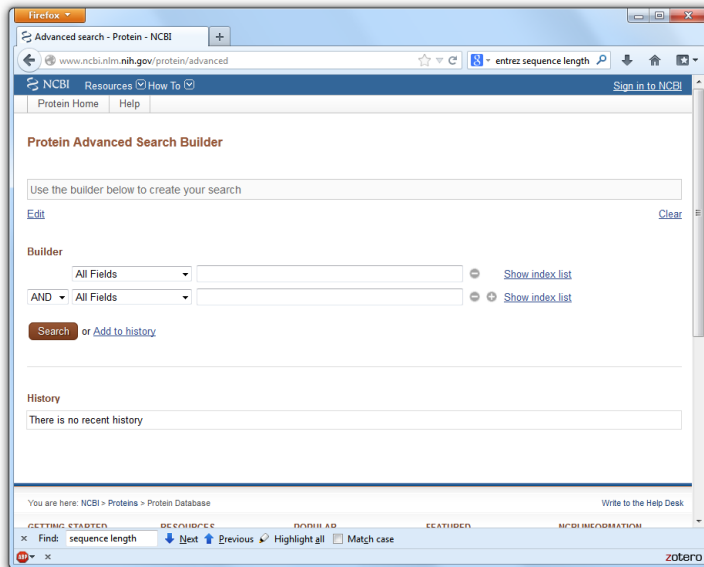
- Combining Terms

xx AND yy	Must have xx and yy
xx OR yy	Must have either xx or yy
NOT zz	Without term zz
xx AND (yy OR zz)	Complex example

- Chemical Properties

75:100 [sequence length]
3500:6000 [molecular weight]

Advanced Searches



Practice

- What's the sequence of your favorite protein?
- What's the extinction coefficient of human heart fatty acid binding protein?
- What human disease is associated with phenylalanine hydroxylase?

Website #3: Protein Data Bank

- <http://rcsb.org/>
- **Input:** Protein name, PDB ID, authors, etc.
- **Output:** 3D coordinates of protein structures
 - Author information on methods
 - Cofactors and other information

What is a PDB file?

- Example: Ricin (2AAI)
- Text file contains a summary of information used in structure determination
- Most important: ATOM records contain X, Y, Z in *Ångströms* (1×10^{-10} m)
 - Most atoms have a radius of 0.5-2 Å

Properties of PDB Files

- Experimental methodology:
 - X-Ray: Typically more precise
 - NMR: Need lots of “restraints;” sometimes hard to assess quality
- “Good” Structures (for X-Ray)
 - Low resolution ($< 2\text{\AA}$)
 - Low R-value ($< 20\%$)
 - Low R_{free} -value ($< 25\%$)

Searching the PDB

Query Refinements: [Select an item or pie chart](#) Hide

Organism <ul style="list-style-type: none"> Homo sapiens (458) Mus musculus (15) Haemophilus influenzae (14) Methanosarcina thermophila (12) Drosophila melanogaster (6) Thalassiosira weissflogii (6) Coccomyxa sp. PA (5) Other (36) 	Taxonomy <ul style="list-style-type: none"> Eukaryota (502) Bacteria (31) Archaea (18) Unassigned (5) Viruses (1) Other (1) 	Experimental Method <ul style="list-style-type: none"> X-ray (554) Neutron Diffraction (2) Hybrid (1) 	X-ray Resolution <ul style="list-style-type: none"> less than 1.5 Å (86) 1.5 - 2.0 Å (273) 2.0 - 2.5 Å (165) 2.5 - 3.0 Å (30) 3.0 and more Å (1) more choices... 	Release Date <ul style="list-style-type: none"> before 2000 (128) 2000 - 2005 (68) 2005 - 2010 (168) 2010 - today (193) this year (11) more choices...
Polymer Type <ul style="list-style-type: none"> Protein (557) 	Enzyme Classification <ul style="list-style-type: none"> 4: Lyases (538) 3: Hydrolases (9) 	SCOP Classification <ul style="list-style-type: none"> All beta proteins (261) Alpha and beta proteins (a/b) (12) 	Protein Symmetry <ul style="list-style-type: none"> Cyclic (528) Dihedral (29) 	
Protein Stoichiometry <ul style="list-style-type: none"> Monomer (480) Homomer (71) Heteromer (6) 				

[Refine Query with Advanced Search](#) Remove Similar: [Select Percent Similarity](#)

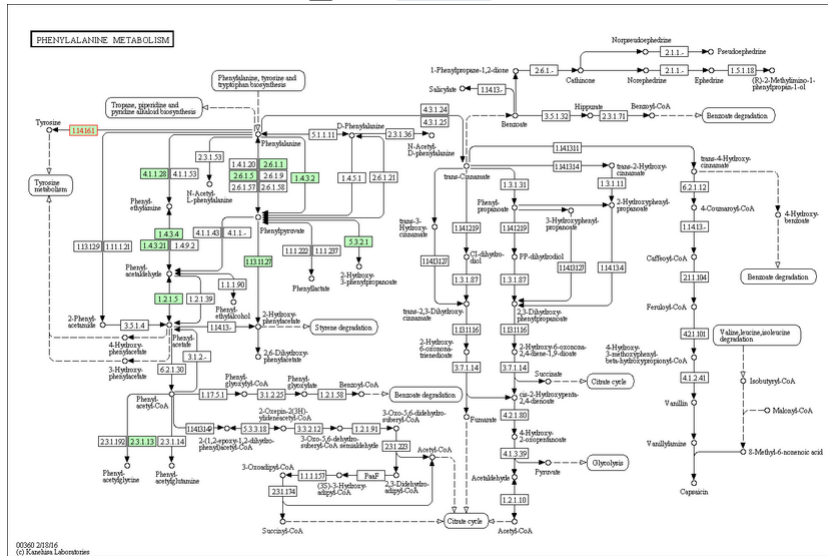
Advanced Searching

The screenshot displays the 'Advanced Search Interface' with a search bar containing 'carbonic anhydrase'. To the right of the search bar, it shows 'Result Count: 557 PDB Entries (Structures)'. Below the search bar, there is a section labeled 'AND' with a 'Choose a Query Type:' dropdown menu. At the bottom of the interface, there are checkboxes for 'Remove Similar Sequences at: 90% Identity' and 'Match: all of the above conditions'. There are also buttons for 'Clear All Parameters' and 'Submit Query'.

Website #4: KEGG

- <http://www.genome.jp/kegg/>
(Kyoto Encyclopedia of Genes and Genomes)
- **Input:** Protein name, PDB ID, authors, etc.
- **Output:** What reactions does an enzyme catalyze?
 - Metabolic pathways
 - The “big picture”

Pathway for Phenylalanine Hydroxylase



Think and Discuss

What are the advantages to large, public databases of scientific information? Are there any disadvantages?

Summary

- Protein properties depend on their primary, secondary, tertiary, and quaternary structure
- Computer databases can organize huge amounts of data on biomolecular systems
- Entrez and the PDB are curated from published research worldwide