

## Digitally Assessing Protein Properties

Biochemistry Boot Camp 2019  
 Session #2  
 Nick Fitzkee  
 nfitzkee@chemistry.msstate.edu

1

## 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?

2

## Sequence of GB3

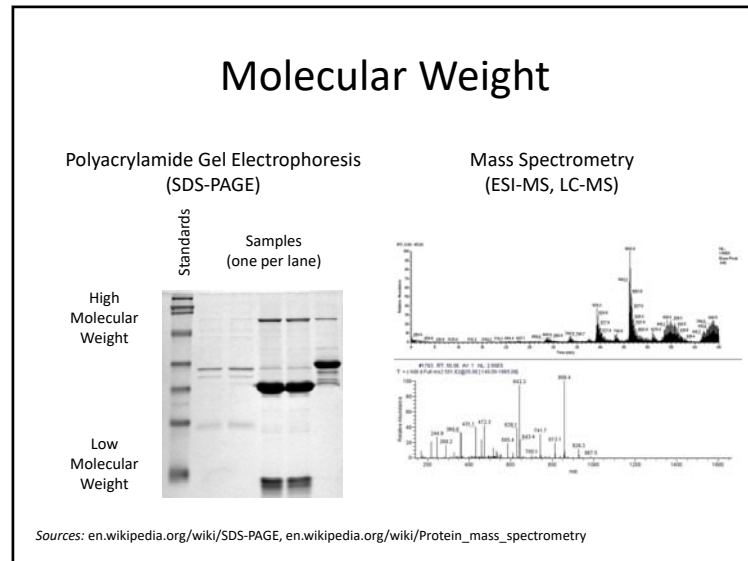
- Primary Structure:  
**NT**-Met-Gln-Tyr-Lys-...-Thr-Glu-**CT**
- More convenient:  
 MQYKLVINGK TLKGETTTKA VDAETAEKAF  
 KQYANDNGVD GVWTYDDATK TFTVTE
- Can we search this (think Google)?

3

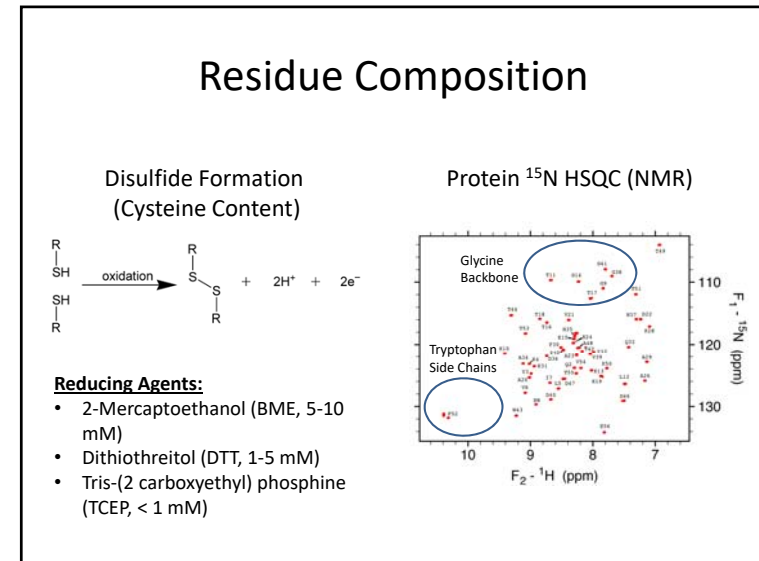
## Website #1: Protparam

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

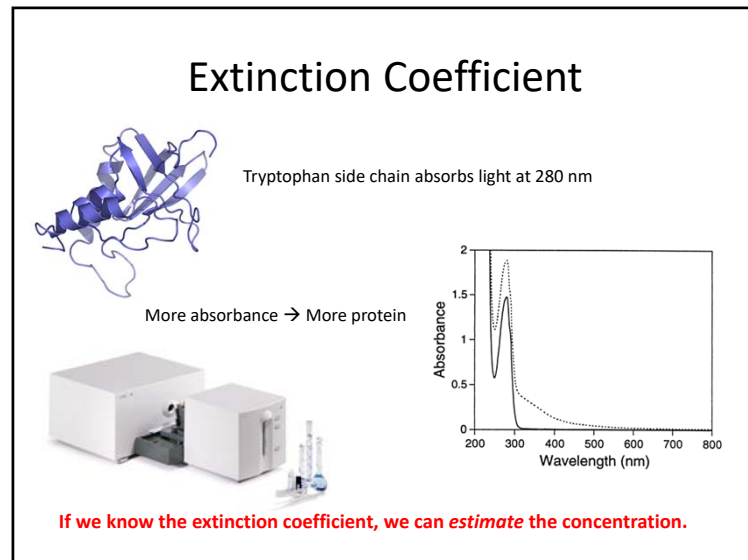
4



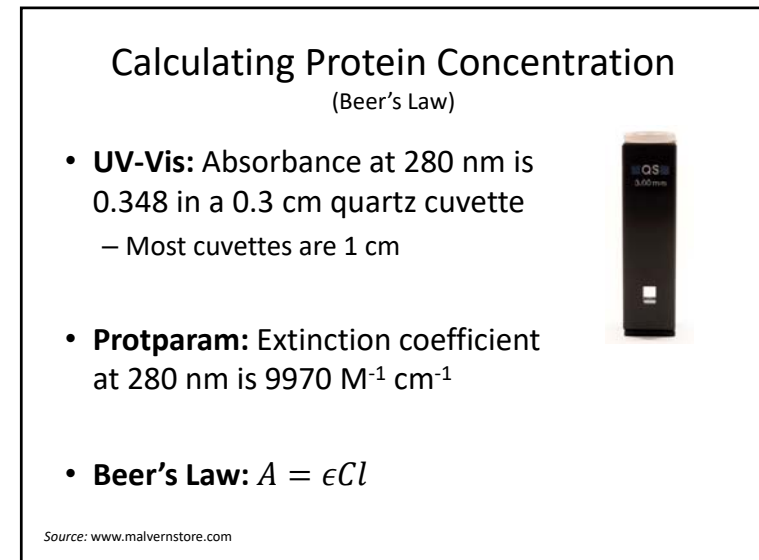
5



6



7



8

## 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

9

## 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

10

## Think and Discuss

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

11

## 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

12

## 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)

13

## NCBI Tricks #2

- Combining Terms

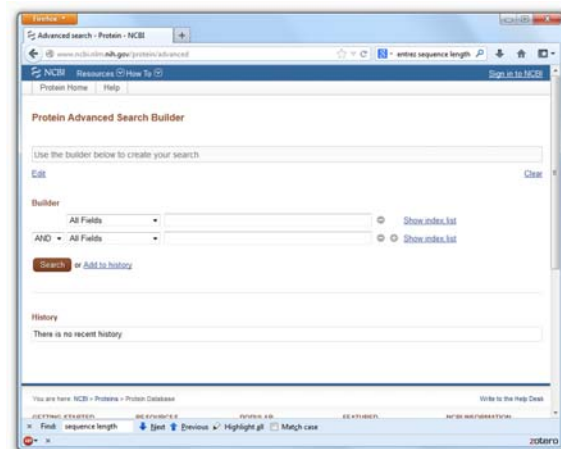
xx AND yy                Must have xx and yy  
 xx OR yy                Must have either xx or yy  
 NOT zz                  Without term zz  
 xx AND (yy OR zz)      Complex example

- Chemical Properties

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

14

## Advanced Searches



15

## *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?

16

## 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

17

## 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* ( $1 \times 10^{-10}$  m)
  - Most atoms have a radius of 0.5-2 Å

18

## 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 ( $< 2\text{Å}$ )
  - Low R-value ( $< 20\%$ )
  - Low  $R_{\text{free}}$ -value ( $< 25\%$ )

19

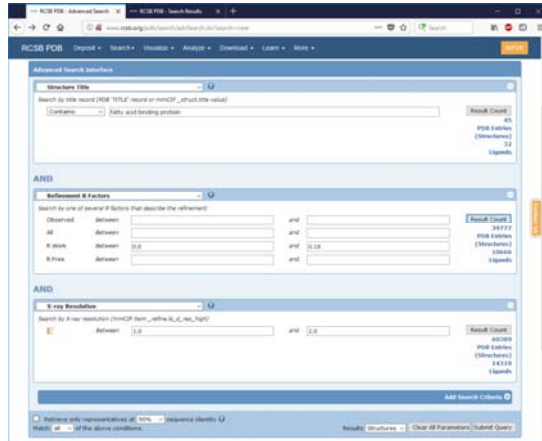
## Searching the PDB

The screenshot shows the RCSB PDB search results page. The search parameter is 'Fatty acid binding protein'. The results list several structures, with '3PL5' selected. The 'Refinements' section is circled in red, showing a list of structures with their respective R-values and R-free values. The details for '3PL5' are also visible, including the method (X-ray diffraction) and resolution (2.04 Å).

Note refinements!

20

## Advanced Searching



21

## 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”

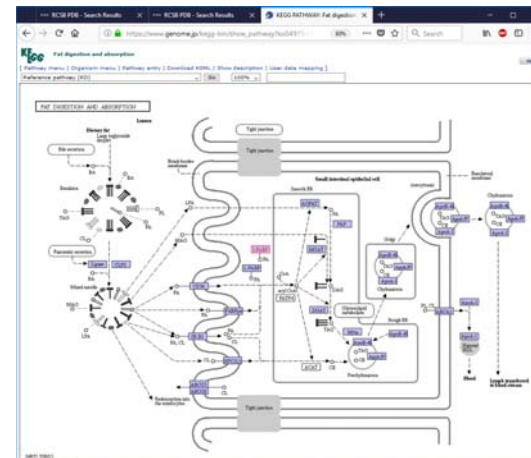
22

## Search Result: Intestinal FABP

KEGG ORTHOLOGY: K08751	
<b>Entry</b>	K08751 KO
<b>Name</b>	FABP2
<b>Definition</b>	fatty acid-binding protein 2, intestinal
<b>Pathway</b>	ko03328 PPAR signaling pathway ko04075 Fat digestion and absorption
<b>Write</b>	KEGG Orthology (KO) [88,100000] Organismal System Endocrine system 03328 PPAR signaling pathway K08751 FABP2; fatty acid-binding protein 2, intestinal Digestion system 04075 Fat digestion and absorption K08751 FABP2; fatty acid-binding protein 2, intestinal
<b>Genes</b>	<a href="#">MGI: 2159(FABP2)</a> <a href="#">PPI: 74821(FABP2)</a> <a href="#">PM: 1009512(FABP2)</a> <a href="#">GGC: 1815128(FABP2)</a> <a href="#">PM: 10044937(FABP2)</a> <a href="#">MGI: 1004515(FABP2)</a> <a href="#">MGI: 785475(FABP2)</a> <a href="#">MGI: 16144095(FABP2)</a> <a href="#">CAG: 10121078(FABP2)</a> <a href="#">MGI: 18466108(FABP2)</a> <a href="#">MGI: show all</a>
<b>Reference</b>	<a href="#">PMID: 2071617</a> <b>Authors</b> Storch J, Thumser AE <b>Title</b> Issue-specific functions in the fatty acid-binding protein family. <b>Journal</b> J Biol Chem 205:12679-83 (2030) <a href="#">DOI:10.1074/jbc.210.13.12679</a>

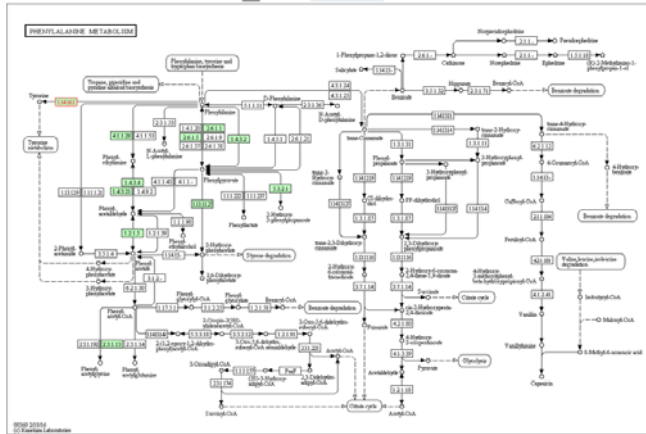
23

## Search Result: Fat Digestion and Absorption



24

## Pathway for Phenylalanine Hydroxylase



25

## Think and Discuss

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

26

## 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

27