

Room E3607  
Protein bioinformatics  
260.841 Protein Bioinformatics

Computer lab  
Tuesday, May 17, 2005  
Sean Prigge  
Jonathan Pevsner  
Ingo Ruczinski

Outline of today's lab

	Topic	Suggested time
1	Find a protein's primary amino acid sequence	15 minutes
2	Assess the secondary structure of the protein	15 minutes
3	Find out whether three-dimensional structure data are available for this protein (or a region of the protein)	15 minutes
4	Analyze the 3D structure of a protein	30 minutes
5	2D gels, mass spec data, and protein interaction networks	as available

We will focus on six proteins that you told us you are currently researching:

	Protein
1	Yeast Esa1
2	Human IL-13
3	Rat tech: neuronal rho gef
4	Human huntingtin
5	Nuclear receptor NR2E3
6	CD4-like Lag3

Part 1 (suggested time: 15 minutes)  
 Topic: find the right protein sequence  
 Approach 1: NCBI  
 Approach 2: ExPASy

	Protein (by name)	Entrez Gene result (NCBI)		
		#	Name	RefSeq or Gb
1	Yeast Esa1	7	ESA1 Esa1p [Saccharomyces cerevisiae] Other Aliases: YOR244W, TAS1	NP_014887
2	Human IL-13	83	IL13 interleukin 13 [Homo sapiens]	NP_998968
3	Rat tech: neuronal rho gef	2	RGD:1303132 Interim Symbol: Tech and Name: neuronal RhoA GEF protein [Rattus norvegicus]	NP_958429
4	Human huntingtin	214 (172 H.s.)	huntingtin [Homo sapiens]	NP_002102
5	Nuclear receptor NR2E3	7	NR2E3 nuclear receptor subfamily 2, group E, member 3 [Homo sapiens]	NP_055064 (410 aa) NP_057430 (367 aa)
6	CD4-like Lag3	3	LAG3 lymphocyte- activation gene 3 [Homo sapiens]	NP_002277

Key points:

[1] NCBI Entrez is a good starting point for finding protein and DNA sequences. From the main page of NCBI, enter a search. Then try “Entrez Gene.”

[2] Names can be ambiguous, so be careful. For Esa1 there are matches to non-yeast proteins.

[3] For huntingtin there are 214 entries in Entrez Gene. You can reduce the list of results by selecting *Homo sapiens*, and by using the limit option to select “from RefSeq only.”

[4] RefSeq (reference sequence) accession numbers have the format NP\_XXXXXX (for proteins) or NM\_XXXXXX (for DNA). RefSeq is preferred because these are curated (agree-upon by experts).

From the Entrez Protein page, select “FASTA” from the display options for each protein. The FASTA format consists of a carat > followed by an unbroken string of text (it can wrap across several lines, as shown in the first example below). The FASTA format is extremely useful for pasting your protein sequence into a variety of on-line tools, such as secondary structure prediction programs.

```
>gi|6324818|ref|NP_014887.1| Histone acetyltransferase catalytic subunit of the native multisubunit complex (NuA4) that acetylates four, conserved internal lysines of histone H4 N-terminal tail; required for cell cycle progression; Esalp [Saccharomyces cerevisiae]
MSHDGKEEPGIACKKINSVDDII IKCQCWVQKNDEERLAEILSINTRKAPPKFYVHYVNYNKRLEDEWITTD
RINLDKEVLYPKLKATDEDNKKQKKKATNTSETPQDSLQDGVDFGFSRENTDVMMLDNLNVQGIKDENIS
HEDEIKKLRITSGSMTQNPHEVARVRNLNRIIMGKYEIEPWFYFSPYPIELTDEDFIYIDDFTLQYFGSKKQ
YERYRKKCTLRHPPGNEIYRDDYVSFFEIDGRKQRTWCRNLCLLSKFLFDHKTLYYDVPFLFYCMTRRD
ELGHHLVGYFSKEKESADGYNVACILTLPPQYQRMGYGKLLIEFSYELSKKENKVGSPKPLSDLGLLSYR
AYWSDTLITLLVEHQKEITIDEISSMTSMTTDLHTAKTLNLRYYKQHIIFLNEDILDYRNLKAKK
RRTIDPNRLIWKPPVFTASQLRFAW
```

```
>gi|47523238|ref|NP_998968.1| interleukin-13 [Sus scrofa]
MALWLTLVIALTCFGLASPGVPVPHSTALKEIEELVNIQTQKTPCNGSMVWSVNLTTSMQYCAALE
SLINISDCSAIQKTQRMLSALCSHKPPSEQVPGKHIRDTKIEVAQFVKDLLKHLRMIFRHG
```

```
>gi|41152499|ref|NP_958429.1| neuronal RhoA GEF protein [Rattus norvegicus]
MDKGRAAKVCHHADCCQLHHRGPNLCEICDSKHFHNTTHYDGHVRFDLPPQGSVLARNVSTRSCPPRTSP
AGDLEEEDEGYTNGKDRKSAGLKI SKKKARRRHTDDPSKECFTLKFIDLNDIETEIVPAMKKKSLGEVL
LPVFERKGIAGLVDIYLDQSNTPLSLTFEAYRFGGHYLRVKAKPGDEGKVEQGVKDSKLSL PALRPSG
AGTPVLERVDPQSRRESSLDILAPGRRRKNMSEFLGDT SIPGQES PAPSSCSLPVGS SVGSSGSSES WKN
RAASRFGSFFSSPSTGAFGREVDKMEQLESK L HAYS L FGLPRMPRRLRFDHDSWEEEEEDDEEEDNSGL
RLED SWRELIDGHEKLRQCHQQA V WEL LHTEVSYIRKLRVITNLFCLLN LQESGLLCEVEAERLF
SNIPELARLHRGLWSSVMVPVLEKARRTRALLQPSDFLKGFKMFGSLFKPYIRYCMEEEGCMEYMRSLLR
DNDLFRAYVTWAEKHQCCQLKLS DMLAKPHQRLTKYPLLLKSVLRKTDEPRAKEA IITMISSVERFIHH
VNTCMRQRQERQRLAGVVSRI DAYEVVEGSNDEVDKFLKEFLHLDLTAPMPGTSPEETRQLLLEGLSRMK
EGKDSKMDVYCFLFTDLLLVTKAVKKAERTKVI RPELLVDKIVCRELRDPGSFLLIHLNEFHSAVGAYTF
QASSQALCRSWVDTLYNAQNQLQQLRAQLLCAQEHPGTQHLSLEEEEEEDEQEEEGEESGTS AASSPTILR
KSSNSLDSEHCASDGSTETLAMVVVEPGETLSSPEFDRGPFSSQSDEASLSNTTSSITPTSELLPLGPVD
GRSCSMDSAYGTLSPSLQDFAAHPVVEPVVPQTLSPPQSPRLRRRTPVQLLPRLLPHLLKSKSEASLL
QLLSGTTT SVSPPAPSRSLSELCLITMAPGVRTQSSLQEGGPGWNC PGACGPGCQGPPLSESENRP SHKAG
GPADSARRKCREMPCGTVPRVQPEPSPGISAQHRKLT LAQLYRIRTTLLLNSTLTASEV
```

```
>gi|4753163|ref|NP_002102.2| huntingtin [Homo sapiens]
MATLEKLMKAFESLKS FQQQQQQQQQQQQQQQQQQQQQQQQPPPPPPPPPPQLPQPPPPQAQPLLQPQPP
PPPPPPPPGPAVAEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPEFQKLLGIAMELFLLCSDD
AESDVRMVADECLNKVIKALMDSNLPRLQLELYKEIKKNGAPRSLRAALWRFAELAH LVRPQKCRPYLVN
LLPCLTRTSKRPEESVQETLAAAVPKIMASFGNFANDNEIKVLLKAFIANLKSSSPTIRRTAAGSAVSIC
QHSRRTQYFYSWLLNVLLGLLVPVEDEHSTLLILGVLLTLRYLVPLLQQQVKDTS LKGSFGVTRKEMEVS
PSAEQLVQVYELTLHHTQH QDHNVVTGALELLQQLFRTPPPPELLQTLTAVGGIGQLTAAKEESGGRSRSG
SIVELIAGGGSSCSPVLSRKQKGVLLGEEEALEDDSESRSDVSSALTASVKDEISGELAASSGVSTPG
SAGHDIITEQPRSQHTLQADSVDLASC DLTSSATDGDEEDILSHSSQVSAVPSDPAMD LNDGTQASSPI
SDSSQTTTEGPDSAVTPSDSSEIVLDGTDNQYLG LQIGQPQDEDEEATGILPDEASEAFRNSSMALQQA H
LLKNMSHCRQPSDSSVDKFLRDEATEPGDQENKPCRIKGDIGQSTDDDSAPLVHCVRLLSASFLLTGGK
NVLVPRDRVRSVKALALSCVGAVALHPESFFSKLYKVP LDTTEYPEEQYVSDILNYIDHGD PQVRGAT
AILCGTLICISILSRSRFHVGDWMTIRTLTGNTFSLADCIPLLRKTLKDESSVTCKLACTAVRNCVMSLC
SSSYSELGLQLIIDVLT LRNSSYWLVRTELETLAEIDFRLVSFLEAKAENLHRGAHHTGLLK LQERVL
NNVVIHLLGDEDPVRVHVAASLIRLVPKLFYKCDQGOADPVVAVARDQSSVYLKLLMHETQPPSHFSVS
```



Part 2 (suggested time: 15 minutes)

Topic: assess the secondary structure of the protein

Approach 1: POLE

Approach 2: ExPASy.

Part 3 (suggested time: 15 minutes)

Topic: find out whether three-dimensional structure data are available for this protein (or a region of the protein)

Approach 1: Entrez Structure at NCBI

Approach 2: blastp search at NCBI

Approach 3: text search at PDB

Approach 4: FASTA search at PDB

Example 1: huntingtin

	Entrez Structure	PDB ID	NCBI blast versus PDB
1	4 (ESA1)	1FY7, 1MJB 1MJA, 1MJ9	
	1 (Esa1p)	1FY7	
2	1 (nmr structure)	1GA3	
3	0 matches	---	matches found
4	5 (all HIPs)	---	no sig. matches
5	0 matches	---	matches found
6	0 matches	---	matches found

protein 3: tech NP\_958429  
 blastp versus PDB  
 matches found e.g. RhoA, 1XCG

NCBI *formatting* **BLAST**  
 Nucleotide Protein Translations Retrieve results for an RID

Your request has been successfully submitted and put into the Blast Queue.

Query = gi|41152499 (1039 letters)

Putative conserved domains have been detected, click on the image below for detailed results.

The diagram shows a horizontal line representing the query sequence from position 1 to 1039. Two domains are highlighted: a red box labeled 'RhoGEF' spanning approximately from position 400 to 600, and a blue box labeled 'PH' spanning approximately from position 650 to 750.

RID=1115908591-29811-209543723098.BLASTQ2, gi|41152499[ref|NP\_958429.1|neuron RhoA GEF protein - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi

If you have any problems or questions with the results of this search please refer to the [BLAST FAQs](#)

[Taxonomy reports](#)

**Distribution of 13 Blast Hits on the Query Sequence**

Mouse-over to show define and scores. Click to show alignments

Color Key for Alignment Scores

<40	40-50	50-80	80-200	>=200
-----	-------	-------	--------	-------

41152499

0 250 500 750 1000

**Related Structures**

Sequences producing significant alignments:

Accession	Description	Score (bits)	E Value
<a href="#">gi 58177305 pdb 1XCG E</a>	Chain E, Crystal Structure Of Human ...	120	2e-27
<a href="#">gi 55670771 pdb 1X86 G</a>	Chain G, Crystal Structure Of The Dh...	106	2e-23
<a href="#">gi 21465839 pdb 1KI D</a>	Chain D, Guanine Nucleotide Exchange...	83	3e-16
<a href="#">gi 13096553 pdb 1FOE G</a>	Chain G, Crystal Structure Of Rac1 I...	50	2e-06
<a href="#">gi 50513392 pdb 1RJ2 J</a>	Chain J, Crystal Structure Of The Dh...	40	0.002
<a href="#">gi 21466030 pdb 1LB1 G</a>	Chain G, Crystal Structure Of The Dh...	40	0.002
<a href="#">gi 20151150 pdb 1K29 C</a>	Chain C, Dbscdc42 (Y889f) >gi 2015114...	40	0.002
<a href="#">gi 50513258 pdb 1NTY A</a>	Chain A, Crystal Structure Of The Fi...	39	0.006
<a href="#">gi 20151146 pdb 1K27 C</a>	Chain C, Crystal Structure Of The Dh...	38	0.008
<a href="#">gi 10835779 pdb 1F5X A</a>	Chain A, Nmr Structure Of The Y174 A...	34	0.15
<a href="#">gi 6435590 pdb 1BY1 A</a>	Chain A, Dbl Homology Domain From Bet...	29	4.9
<a href="#">gi 51247863 pdb 1UZ7 A</a>	Chain A, Complex Of The Anti-Hyperte...	29	4.9
<a href="#">gi 31615824 pdb 1ON9 F</a>	Chain F, Transcarboxylase 12s Crysta...	28	6.4

Alignments

Get selected sequences Select all Deselect all

protein 4: huntingtin NP\_002102  
 blastp versus PDB  
 result: 1 match, not significant

RID=1115908110-11810-61489232233.BLASTQ2, gi|4753163|ref|NP\_002102.2|huntingtin (Homo sapiens) - Microsoft Internet Explorer

Address <http://www.ncbi.nlm.nih.gov/BLAST/Blast.cgi>

**Distribution of 1 Blast Hits on the Query Sequence**

Mouse-over to show define and scores. Click to show alignments

**Related Structures**

Sequences producing significant alignments:	Score (bits)	E Value
<a href="#">gi 4558259 pdb 1B3U E</a> Chain B, Crystal Structure Of Constan...	34	0.36

**Alignments**

Get selected sequences    Select all    Deselect all

[gi|4558259|pdb|1B3U|E](#) Chain B, Crystal Structure Of Constant Regulatory Domain Of Human Pp2a, Pr65alpha

[gi|4558258|pdb|1B3U|A](#) Chain A, Crystal Structure Of Constant Regulatory Domain Of Human Pp2a, Pr65alpha  
 Length = 588

Score = 34.3 bits (77), Expect = 0.36  
 Identities = 39/194 (20%), Positives = 70/194 (36%), Gaps = 49/194 (25%)

```

Query: 122 EFQKLLGI-----AMELFLLCSDDAESDVRMVADECLNKVIKALMDSNLPRLQLELYKE 175
      EF K+L + + +F + D + VR++A E + + L +L L + ++
Sbjct: 190 EFAKVLLELDNVKSEIIPNFSNLASDEQDSVRLLAVEACVNIQLLPQEDLEALVMPTRLRQ 249

Query: 176 IKKNGAPRSLRAALWRF AELAHVLRPQ----- 202
      ++ + R +F EL V P+
Sbjct: 250 AAEDKSWRVRYMVDKFTLQKAVGPEITKTDLVPAFQNLMKDCEAEVRAAASHKVKEFC 309

Query: 203 -----KCRPYLV--NLLPCLTRTSKRPEESVQETLAAAVPKIMASFGNFAANDNEIKVLLK 255
      CR ++ +LPC+ + V+ LA+ + + G DN I+ LL
Sbjct: 310 ENLSADCRENVIMSQILPCIKELVSDANQHVKSALASVIMGLSPILGK---DNTIEHLLP 366

Query: 256 AFIANLKSSSP TIR 269
      F+A LK P +R
Sbjct: 367 LFLAQLKDECEPEVR 380
  
```

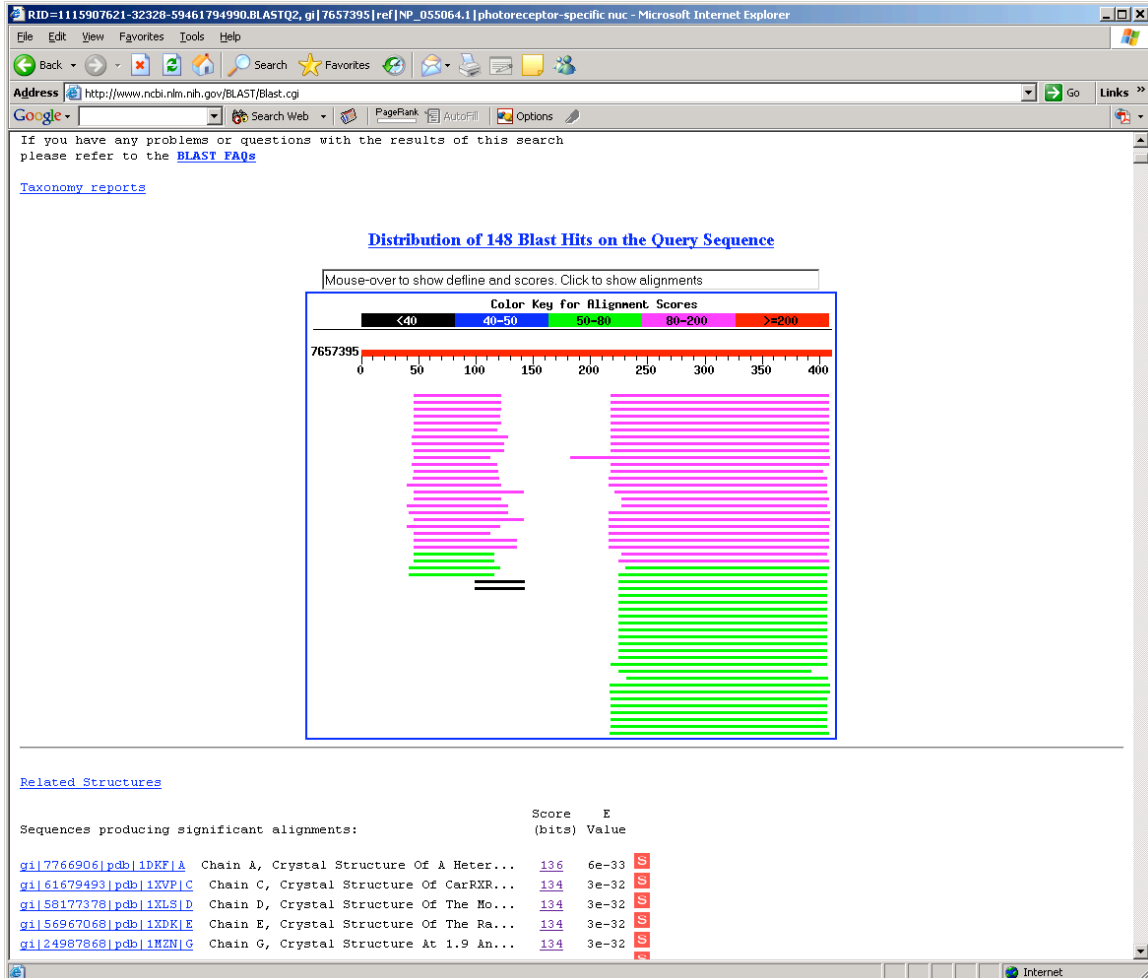
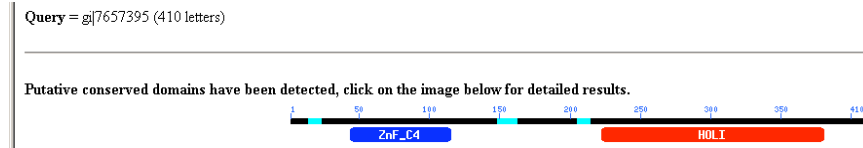
Done Internet



protein 5: NR2E3 NP\_055064

blastp versus PDB

result: many matches, e.g. 1DKF (residues 219-409 of NR2E3), 1YNW (vitamin D receptor) (residues 47-123)

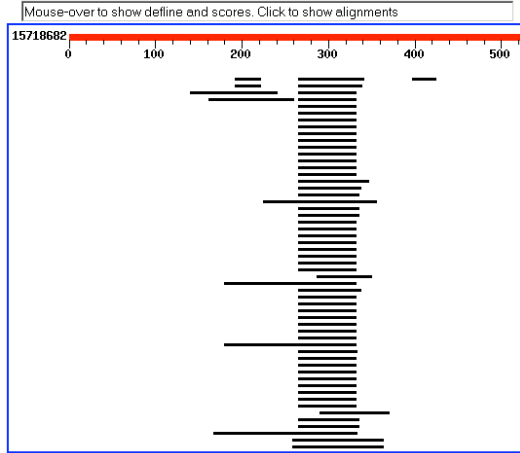


protein 6: lag3 NP\_002277  
 blastp versus pdb  
 match found, residues 266-342 (approximate)

Putative conserved domains have been detected, click on the image below for detailed results.



**Distribution of 149 Blast Hits on the Query Sequence**



Related Structures

Sequences producing significant alignments:	Score (bits)	E Value
<a href="#">gi 7766934 pdb 32C2 A</a> Chain A, Structure Of An Activity Sup...	<a href="#">.37</a>	0.009
<a href="#">gi 1064953 pdb 1IBG L</a> Chain L, Igg Fab (Igg2b, Kappa) Fram...	<a href="#">.37</a>	0.009
<a href="#">gi 1127091 pdb 1N5N L</a> Chain L, Immunoglobulin, Staphylococc...	<a href="#">.37</a>	0.011
<a href="#">gi 15825780 pdb 1I72 C</a> Chain C, Antibody Gnc92h2 Bound To L...	<a href="#">.36</a>	0.019
<a href="#">gi 24158825 pdb 1L7T L</a> Chain L, Crystal Structure Analysis ...	<a href="#">.33</a>	0.096
<a href="#">gi 10835827 pdb 1F3D J</a> Chain J, Catalytic Antibody 4b2 In C...	<a href="#">.33</a>	0.16
<a href="#">gi 20150130 pdb 1I9R Y</a> Chain Y, Structure Of Cd401 In Compl...	<a href="#">.32</a>	0.21
<a href="#">gi 4558340 pdb 3F58 L</a> Chain L, Igg1 Fab Fragment (58.2) Com...	<a href="#">.32</a>	0.21
<a href="#">gi 4558336 pdb 2F58 L</a> Chain L, Igg1 Fab Fragment (58.2) Com...	<a href="#">.32</a>	0.21

Part 4 (suggested time: 30 minutes)  
Topic: analyze the 3D structure of a protein  
Approach: CN3D

Part 5 (as time permits)

Topic: 2D gels, mass spec data, and protein interaction networks

Approach: visit ExPASy