

# ***Searching Nucleotide Databases***

*{MATRIX}  
{SCIENCE}*

K P I R L T A D L L A E T L Q A R R E W G P I F N I  
 A S P S D # Q Q I S W Q K L Y K P E E S G G Q Y S T F  
 Q A H Q T N S R S L G R N S T S Q K R V G A N I Q H  
 CAAGCCCATCAGACTAACAGCAGATCTCTGGCAGAAACTCTACAAGCCAGAAGAGAGTGGGGCCAAATTCAACAT  
 299200 299210 299220 299230 299240 299250 299260 299270  
 TTGGGGTAGTCTGATTGTCGTCTAGAGAACCGTCTTTGAGATGTTCTGGTCTTCTCACCCTTATAAGTTGTAA  
 C A W \* V L L L D R P L F E V L W F L T P A L I \* C E  
 L G D S + C C I E Q C F S + L G S S L P P W Y E V N  
 L G M L S V A S R K A S V R C A L L S H P G I N L M

*(MATRIX)  
(SCIENCE)*

When we search a nucleic acid databases, Mascot always performs a 6 frame translation on the fly. That is, 3 reading frames from the forward strand and 3 reading frames from the complementary strand.



## ***NA Translation***

- **Always translate in all 6 reading frames**
- **Translation starts from the beginning of the sequence, not from a start codon**
- **When a stop codon is encountered, insert a gap and re-start translation.**
- **No attempt to resolve codon ambiguity.**
- **All translation uses the NCBI standard genetic code.**

*{MATRIX}*  
*{SCIENCE}*

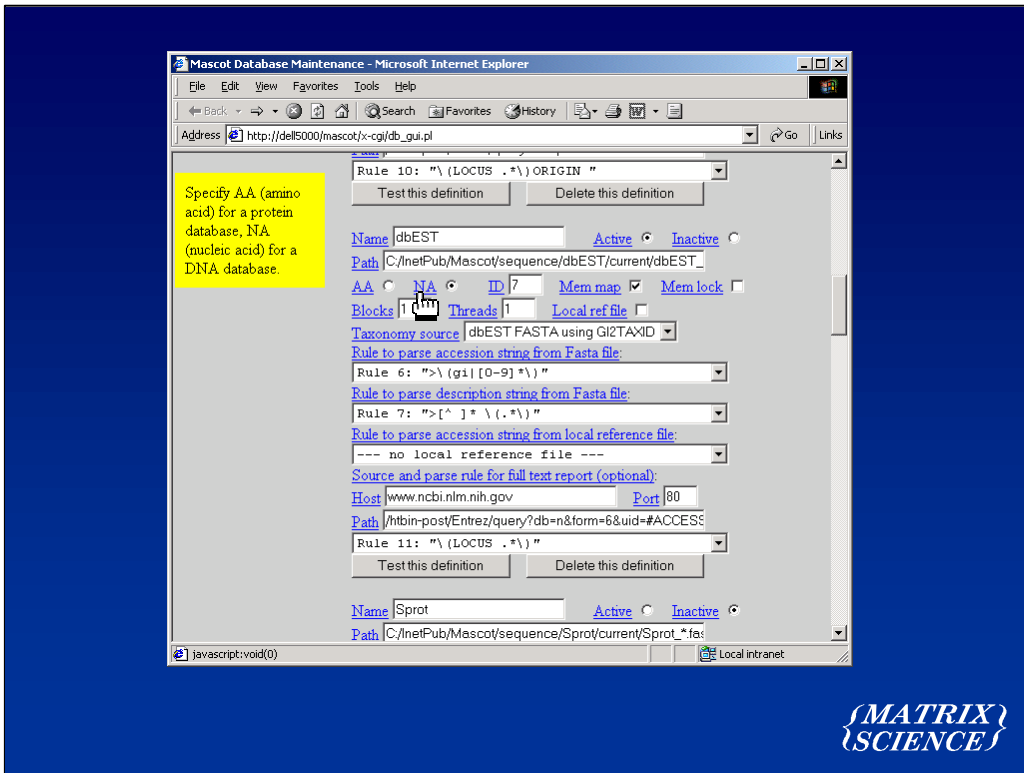
The rules for Nucleic Acid translation in Mascot are:

we translate the entire sequence, we don't look for a start codon.

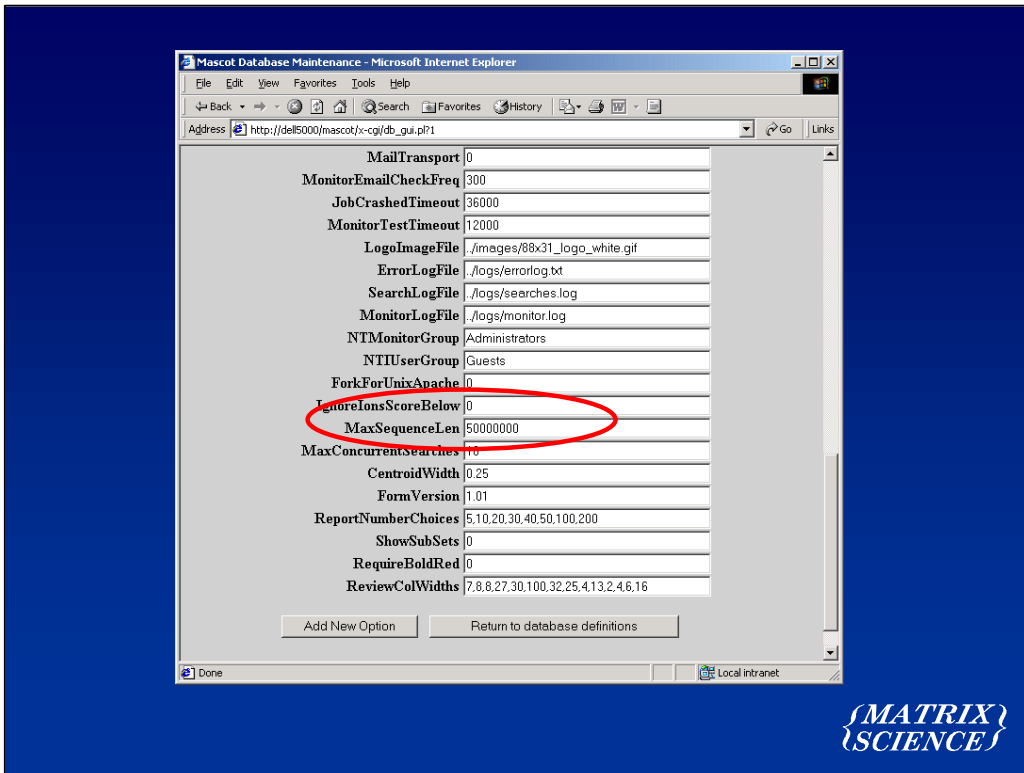
When a stop codon is encountered, we leave a gap, and immediately re-start translation.

There is no attempt to resolve ambiguous codons. For example, ACX can be translated as Threonine, because the identity of the last base is a don't care. However, this is not done in the current code.

Finally, all translations use the standard genetic code. Ideally, we would use species specific code where a sequence has a known taxonomy. But, again, this is not done at present.



Setting up a nucleic acid database in Mascot is no different from setting up a protein database. The only two things to watch are that the database type is specified as NA



*(MATRIX)  
(SCIENCE)*

And, if the sequences are very long, you may need to increase the upper limit on the sequence length of individual entries.

## **Standard data set**

- **Collaboration with Walter Blackstock and Jyoti Choudhary, Cell Map Project, GlaxoSmithKline R&D, Stevenage, UK**
- **Human embryonic kidney cell lysate; digested with trypsin; analysed by LC-MS/MS using a Micromass Q-TOF**
- **169 MS/MS spectra after data reduction (.pkl file)**

*{MATRIX}*  
*{SCIENCE}*

Here are examples of searching the same data set against protein, EST and genomic sequence databases. This data set was generated by Jyoti Choudhary and Walter Blackstock of GlaxoSmithKline.

They generated a high quality LC-MS/MS data from a tryptic digest of whole cell lysate from human embryonic kidney cells.

After data reduction, we were left with 169 MS/MS spectra.

http://dell5000/mascot/cg/master\_results.pl?file=../data/20001016/F003980.dat - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address http://dell5000/mascot/cg/master\_results.pl?file=../data/20001016/F003980.dat#H#1

## MASCOT (SCIENCE) Mascot Search Results

User : JSC  
 Email : JSC@work  
 Search title : Annexin mix  
 MS data file : U:\Mascot test data\Glaxo\gtof10348.pk1  
 Database : MSDB 20000621 (508120 sequences; 156794043 residues)  
 Timestamp : 16 Oct 2000 at 13:38:06 GMT

Significant hits:

<a href="#">L10HU</a>	annexin I - human
<a href="#">AAC52068</a>	HSTALDR3 NID: - Homo sapiens
<a href="#">AAB19866</a>	S57440S13 NID: - Rattus sp.
<a href="#">AAC78495</a>	OCU24656 NID: - Oryctolagus cuniculus
<a href="#">1T6S2</a>	trypsin (EC 3.4.21.4) precursor (with pancreatic secretory trypsin inhibitor), chain Z - bov
<a href="#">1NTP</a>	trypsin (EC 3.4.21.4) (isopropylphosphorylated) - bovine
<a href="#">AAA36574</a>	HUMRNP2A NID: - Homo sapiens
<a href="#">A32915</a>	nucleophosmin - human
<a href="#">Q91TV2</a>	VITAMIN D RESPONSE ELEMENT BINDING PROTEIN.- Saguinus oedipus (Cotton-top tamarin).
<a href="#">Q9XSY6</a>	HHRNP A/B RELATED PROTEIN (FRAGMENT).- Felis silvestris catus (Cat).
<a href="#">S60335</a>	TGF-beta receptor interacting protein 1 - human
<a href="#">Q9QZD9</a>	TGF-BETA RECEPTOR BINDING PROTEIN.- Mus musculus (Mouse).
<a href="#">B38611</a>	casein kinase II (EC 2.7.1.-) alpha' chain - chicken
<a href="#">B35838</a>	casein kinase II (EC 2.7.1.-) alpha' chain - human
<a href="#">S55282</a>	isocitrate dehydrogenase (NAD+) (EC 1.1.1.41) alpha chain precursor - human
<a href="#">AAB47721</a>	HUMKRT1X NID: - Homo sapiens
<a href="#">S41754</a>	CRKL protein - human
<a href="#">KRHU2</a>	keratin 1, type II, cytoskeletal - human
<a href="#">LUGP1</a>	annexin I - guinea pig
<a href="#">BAA37117</a>	AB001915 NID: - Homo sapiens
<a href="#">S40776</a>	ribonucleoprotein - African clawed frog
<a href="#">CAA64477</a>	SSANNEXHI NID: - Sus scrofa
<a href="#">AAA59468</a>	HUMKRT10A NID: - Homo sapiens
<a href="#">PC4375</a>	telomeric and tetraplex DNA binding protein gTBP42 V - rat (fragment)
<a href="#">JCS660</a>	hepatoma-derived growth factor - mouse
<a href="#">1HA11</a>	hnrnp a1 hnrnp a1 (rbd1, rbd2) hnrnp a1 1-184, fragment 1 - human
<a href="#">I52962</a>	FBRNP - human
<a href="#">093446</a>	ANNEXIN MAX3.- Oryzias latipes (Medaka fish).
<a href="#">G3P2_HUMAN</a>	GLYCERALDEHYDE 3-PHOSPHATE DEHYDROGENASE, LIVER (EC 1.2.1.12).- Homo sapiens (Human).

Local intranet

First of all, we searched the data against a comprehensive non-identical protein database, MSDB. We found significant matches to 22 human proteins ... and one non-human, our frequent flyer, bovine trypsin.



Mascot Search Results - Microsoft Internet Explorer

Address [http://dell5000/mascot/cgi/master\\_results.pl?file=../data/20001016/F289840.dat](http://dell5000/mascot/cgi/master_results.pl?file=../data/20001016/F289840.dat)

**MASCOT**  
*(SCIENCE)* **Mascot Search Results**

User : JSC  
 Email : jcottrell@matrixscience.com  
 Search title : Annexin  
 MS data file : U:\Mascot test data\Glaxo\gtof10348.pk1  
 Database : dbEST 20001001 (35277150 sequences; 4704334140 residues)  
 Taxonomy : Homo sapiens (human) (14927850 sequences)  
 Timestamp : 16 Oct 2000 at 21:15:15 GMT

Significant hits:

<a href="#">gi 10348033</a>	601512345F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3913811 5'
<a href="#">gi 10347940</a>	601512293F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3913822 5'
<a href="#">gi 10340616</a>	601509784F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3911108 5'
<a href="#">gi 6571609</a>	xb37e04.y1 NCI_CGAP_Lu31 Homo sapiens cDNA clone IMAGE:2578494 5' similar to SW:ANX1_HUMAN PO
<a href="#">gi 10330826</a>	601505336F2 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3906872 5'
<a href="#">gi 9145520</a>	601140912F1 NIH_MGC_9 Homo sapiens cDNA clone IMAGE:3140763 5'
<a href="#">gi 9891351</a>	601483620F1 NIH_MGC_69 Homo sapiens cDNA clone IMAGE:3886133 5'
<a href="#">gi 10345301</a>	601513625F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3914791 5'
<a href="#">gi 10198665</a>	601348510F1 NIH_MGC_54 Homo sapiens cDNA clone IMAGE:3686738 5'
<a href="#">gi 10347200</a>	601512680F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3914346 5'
<a href="#">gi 10319788</a>	601449048F1 NIH_MGC_65 Homo sapiens cDNA clone IMAGE:3853282 5'
<a href="#">gi 10332342</a>	601508038F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3909657 5'
<a href="#">gi 10099330</a>	RC3-HT0649-150600-023-e08 HT0649 Homo sapiens cDNA
<a href="#">gi 10145359</a>	601565022F1 NIH_MGC_21 Homo sapiens cDNA clone IMAGE:3840124 5'
<a href="#">gi 9155988</a>	601156847F1 NIH_MGC_21 Homo sapiens cDNA clone IMAGE:3140530 5'
<a href="#">gi 9890400</a>	601473296F1 NIH_MGC_68 Homo sapiens cDNA clone IMAGE:3876359 5'
<a href="#">gi 9323971</a>	601236972F1 NIH_MGC_44 Homo sapiens cDNA clone IMAGE:3609088 5'
<a href="#">gi 7948606</a>	RC1-CT0249-030300-026-b03 CT0249 Homo sapiens cDNA
<a href="#">gi 10153647</a>	601556364T1 NIH_MGC_58 Homo sapiens cDNA clone IMAGE:3826407 3'
<a href="#">gi 10319146</a>	601447547F1 NIH_MGC_65 Homo sapiens cDNA clone IMAGE:3851699 5'
<a href="#">gi 6798648</a>	dr04f01.x1 NIH_MGC_3 Homo sapiens cDNA clone IMAGE:2847121 5'
<a href="#">gi 10327945</a>	601486909F1 NIH_MGC_69 Homo sapiens cDNA clone IMAGE:3889058 5'
<a href="#">gi 9898273</a>	601279337F1 NIH_MGC_39 Homo sapiens cDNA clone IMAGE:3611240 5'
<a href="#">gi 10401105</a>	601497445F1 NIH_MGC_70 Homo sapiens cDNA clone IMAGE:3899602 5'
<a href="#">gi 8906307</a>	hw47b05.x1 NCI_CGAP_Lu24 Homo sapiens cDNA clone IMAGE:3176529 3' similar to gb:M55268 CASEIN
<a href="#">gi 8627305</a>	RC3-HT0470-120200-013-g02 HT0470 Homo sapiens cDNA
<a href="#">gi 10146282</a>	601570283F1 NIH_MGC_21 Homo sapiens cDNA clone IMAGE:3844988 5'
<a href="#">gi 8114235</a>	PM3-DT0064-260300-002-g12 DT0064 Homo sapiens cDNA

Done Local intranet

With dbEST, we obtained almost the same results, just a couple of additional peptide matches. However, unlike the protein database search, it doesn't immediately communicate which proteins have been found.

Mascot Search Results - Microsoft Internet Explorer

Address: http://dell5000/mascot/cgi/master\_results.pl?file=../data/20001016/F289840.dat

[Switch to Protein Summary Report](#)

To create a bookmark for this report, right click this link: [Peptide Summary Report \(Annexin\)](#)

Cluster matches using UniGene index for [human](#) [bovine](#)

Select All Select None Search Selected Archive Report

1. [gi|10348033](#) Mass: 35874 Total score: 700 Peptides matched: 14  
601512345F1 NIH\_MGC\_71 Homo sapiens cDNA clone IMAGE:3913811 5'

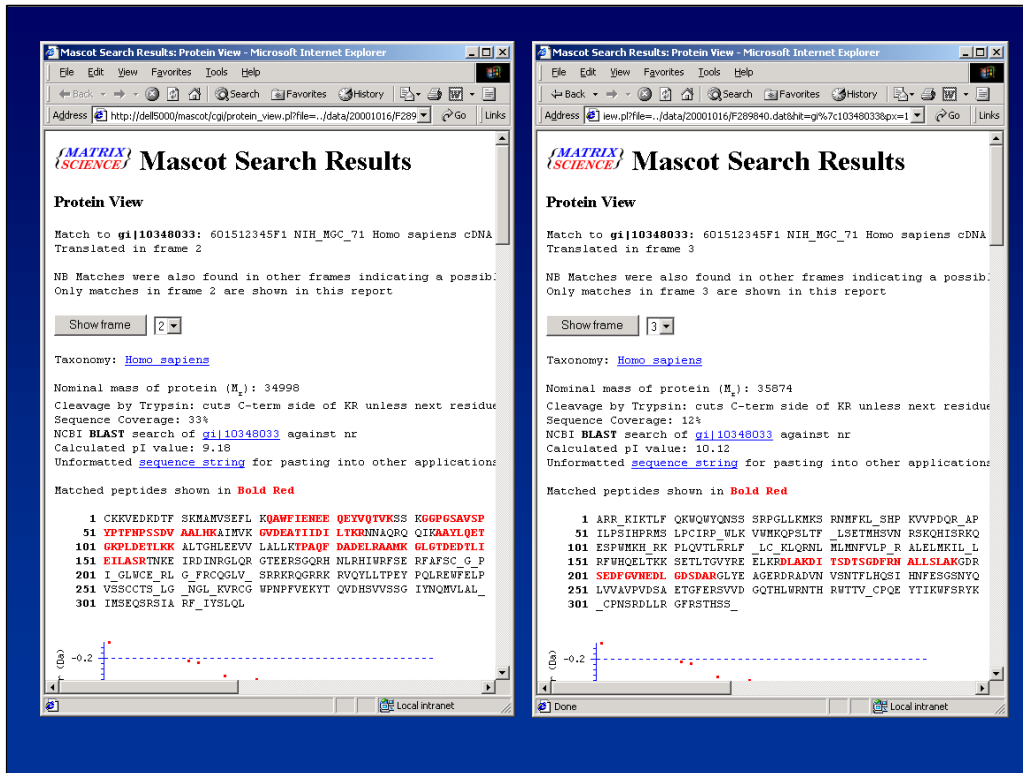
Check to include this hit in archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Rank	Peptide
<input checked="" type="checkbox"/> <a href="#">12</a>	415.19	828.36	828.51	-0.14	0	33	1	NALLSLAK
<input checked="" type="checkbox"/> <a href="#">45</a>	607.16	1212.31	1212.53	-0.21	0	70	1	DITSDTSGDFR
<input checked="" type="checkbox"/> <a href="#">53</a>	631.70	1261.38	1261.59	-0.22	0	69	1	TPAQFDADELK
<input checked="" type="checkbox"/> <a href="#">69</a>	694.25	1386.49	1386.76	-0.27	0	73	1	GVDEATIIDILTK
<input checked="" type="checkbox"/> <a href="#">91</a>	515.20	1542.58	1542.86	-0.28	1	46	1	GVDEATIIDILTKR
<input checked="" type="checkbox"/> <a href="#">98</a>	547.49	1639.45	1639.77	-0.32	1	(41)	1	DLAKDITSDTSGDFR
<input checked="" type="checkbox"/> <a href="#">99</a>	820.75	1639.48	1639.77	-0.29	1	52	1	DLAKDITSDTSGDFR
<input checked="" type="checkbox"/> <a href="#">103</a>	851.77	1701.52	1701.88	-0.36	0	82	1	GLGTDEDTLIEILASR
<input checked="" type="checkbox"/> <a href="#">105</a>	870.21	1738.41	1738.73	-0.32	0	82	2	SEDFGVNEDLGSDAR + 1 Methyl ester (DE)
<input checked="" type="checkbox"/> <a href="#">123</a>	476.92	1903.67	1904.03	-0.36	1	22	1	AAYLQETGKPLEDTLKK
<input checked="" type="checkbox"/> <a href="#">131</a>	707.22	2118.63	2119.08	-0.45	1	35	2	AAMKGLGTDEDTLIEILASR + 1 Oxidation (M)
<input checked="" type="checkbox"/> <a href="#">132</a>	1062.33	2122.64	2122.98	-0.35	0	(72)	1	QAWFIENEQEYVQTVK + 1 Pyro-glu (N-term Q)
<input checked="" type="checkbox"/> <a href="#">133</a>	1070.83	2139.64	2140.01	-0.37	0	84	1	QAWFIENEQEYVQTVK
<input checked="" type="checkbox"/> <a href="#">149</a>	785.91	2354.72	2355.15	-0.43	0	66	1	GGPGSAVSPYPTFNPSSDVAALHK

2. [gi|10347940](#) Mass: 29372 Total score: 694 Peptides matched: 13  
601512293F1 NIH\_MGC\_71 Homo sapiens cDNA clone IMAGE:3913822 5'

Local intranet

The master results report from the EST search looks pretty similar to the MSDB search, except that the EST sequences are mostly shorter than full length proteins, so the peptide matches are more scattered. If we click on the protein accession number link...



we get a protein view. This is similar to the protein view for a protein database entry, except we have drop down list for the different translation frames. For this particular entry, most of the matches have been found in reading frame 2.

But, as so often happens, there is a frame shift in this entry, and there are additional matches in frame 3.

Mascot Search Results - Microsoft Internet Explorer

Address [http://dell5000/mascot/cgi/master\\_results.pl?file=.../data/20001016/F289840.dat](http://dell5000/mascot/cgi/master_results.pl?file=.../data/20001016/F289840.dat)

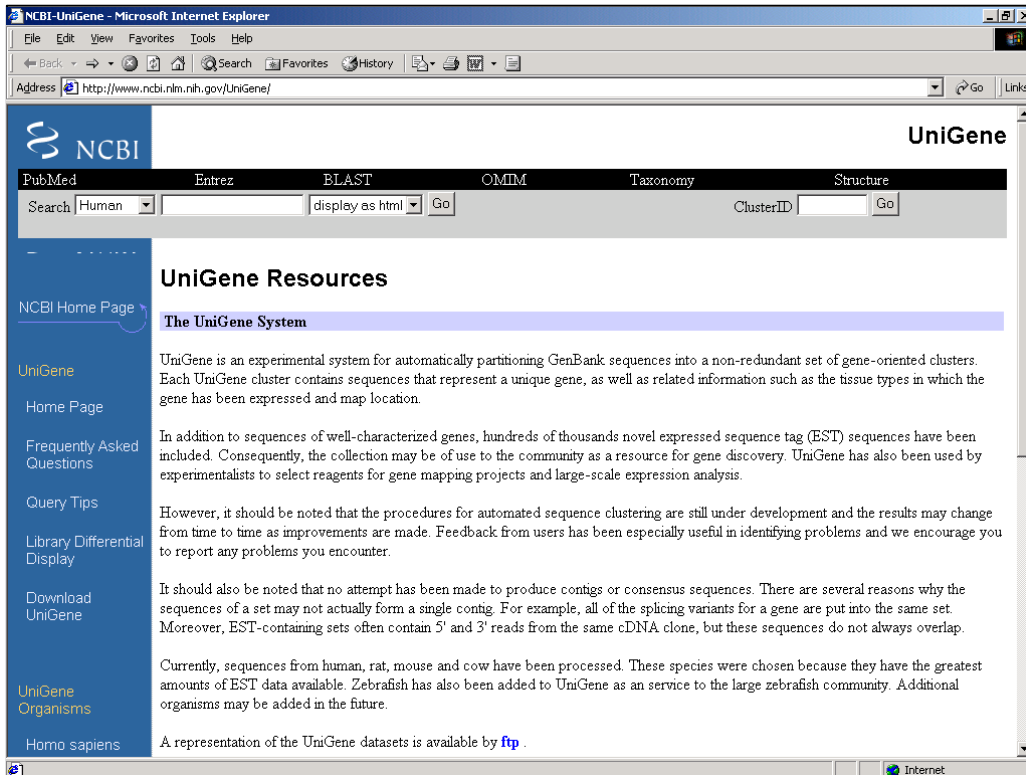
**Mascot Search Results**

User : JSC  
 Email : jcottrell@matrixscience.com  
 Search title : Annexin  
 MS data file : U:\Mascot test data\Glaxo\gtof10348.pk1  
 Database : dbEST 20001001 (35277150 sequences; 4704334140 residues)  
 Taxonomy : Homo sapiens (human) (14927850 sequences)  
 Timestamp : 16 Oct 2000 at 21:15:15 GMT

Significant hits:

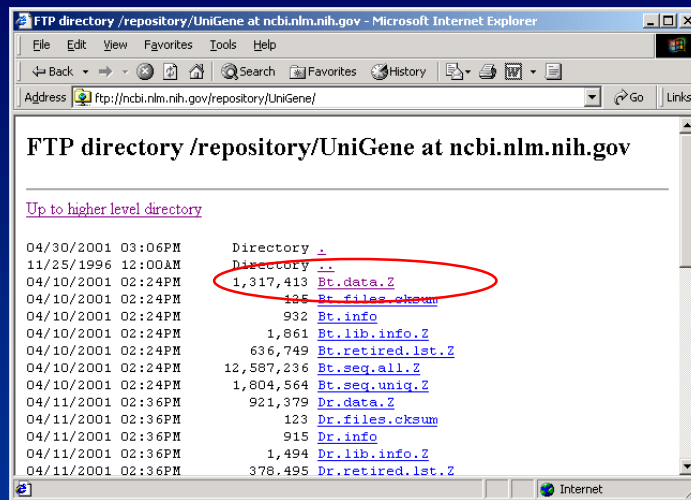
<a href="#">gi 10348033</a>	601512345F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3913811 5'
<a href="#">gi 10347940</a>	601512293F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3913822 5'
<a href="#">gi 10340616</a>	601509784F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3911108 5'
<a href="#">gi 6571609</a>	xb37e04.y1 NCI_CGAP_Lu31 Homo sapiens cDNA clone IMAGE:2578494 5' similar to SW:ANX1_HUMAN PO
<a href="#">gi 10330826</a>	601505336F2 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3906872 5'
<a href="#">gi 9145520</a>	601140912F1 NIH_MGC_9 Homo sapiens cDNA clone IMAGE:3140763 5'
<a href="#">gi 9891351</a>	601483620F1 NIH_MGC_69 Homo sapiens cDNA clone IMAGE:3886133 5'
<a href="#">gi 10345301</a>	601513625F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3914791 5'
<a href="#">gi 10198665</a>	601348510F1 NIH_MGC_54 Homo sapiens cDNA clone IMAGE:3686738 5'
<a href="#">gi 10347200</a>	601512680F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3914346 5'
<a href="#">gi 10319788</a>	601449048F1 NIH_MGC_65 Homo sapiens cDNA clone IMAGE:3853282 5'
<a href="#">gi 10332342</a>	601508038F1 NIH_MGC_71 Homo sapiens cDNA clone IMAGE:3909657 5'
<a href="#">gi 10099330</a>	RC3-HT0649-150600-023-e08 HT0649 Homo sapiens cDNA
<a href="#">gi 10145359</a>	601565022F1 NIH_MGC_21 Homo sapiens cDNA clone IMAGE:3840124 5'
<a href="#">gi 9155988</a>	601156847F1 NIH_MGC_21 Homo sapiens cDNA clone IMAGE:3140530 5'
<a href="#">gi 9890400</a>	601473296F1 NIH_MGC_68 Homo sapiens cDNA clone IMAGE:3876359 5'
<a href="#">gi 9323971</a>	601236972F1 NIH_MGC_44 Homo sapiens cDNA clone IMAGE:3609088 5'
<a href="#">gi 7948606</a>	RC1-CT0249-030300-026-b03 CT0249 Homo sapiens cDNA
<a href="#">gi 10153647</a>	601556364F1 NIH_MGC_58 Homo sapiens cDNA clone IMAGE:3826407 3'
<a href="#">gi 10319146</a>	601447547F1 NIH_MGC_65 Homo sapiens cDNA clone IMAGE:3851699 5'
<a href="#">gi 6798648</a>	dr04f01.x1 NIH_MGC_3 Homo sapiens cDNA clone IMAGE:2847121 5'
<a href="#">gi 10327945</a>	601486909F1 NIH_MGC_69 Homo sapiens cDNA clone IMAGE:3889058 5'
<a href="#">gi 9898273</a>	601279337F1 NIH_MGC_39 Homo sapiens cDNA clone IMAGE:3611240 5'
<a href="#">gi 10401105</a>	601497445F1 NIH_MGC_70 Homo sapiens cDNA clone IMAGE:3899602 5'
<a href="#">gi 8906307</a>	hv47b05.x1 NCI_CGAP_Lu24 Homo sapiens cDNA clone IMAGE:3176529 3' similar to gb:M55268 CASEIN
<a href="#">gi 8627305</a>	RC3-HT0470-120200-013-g02 HT0470 Homo sapiens cDNA
<a href="#">gi 10146282</a>	601570283F1 NIH_MGC_21 Homo sapiens cDNA clone IMAGE:3844988 5'
<a href="#">gi 8114235</a>	PM3-DT0064-260300-002-g12 DT0064 Homo sapiens cDNA

Going back to the issue of the hit list and the descriptions not saying very much. There are several problems here. One is that EST databases usually have a huge amount of redundancy, which can make for very long reports. Another problem is that the sequences tend to be short, so we don't get much grouping of peptide matches into protein matches.



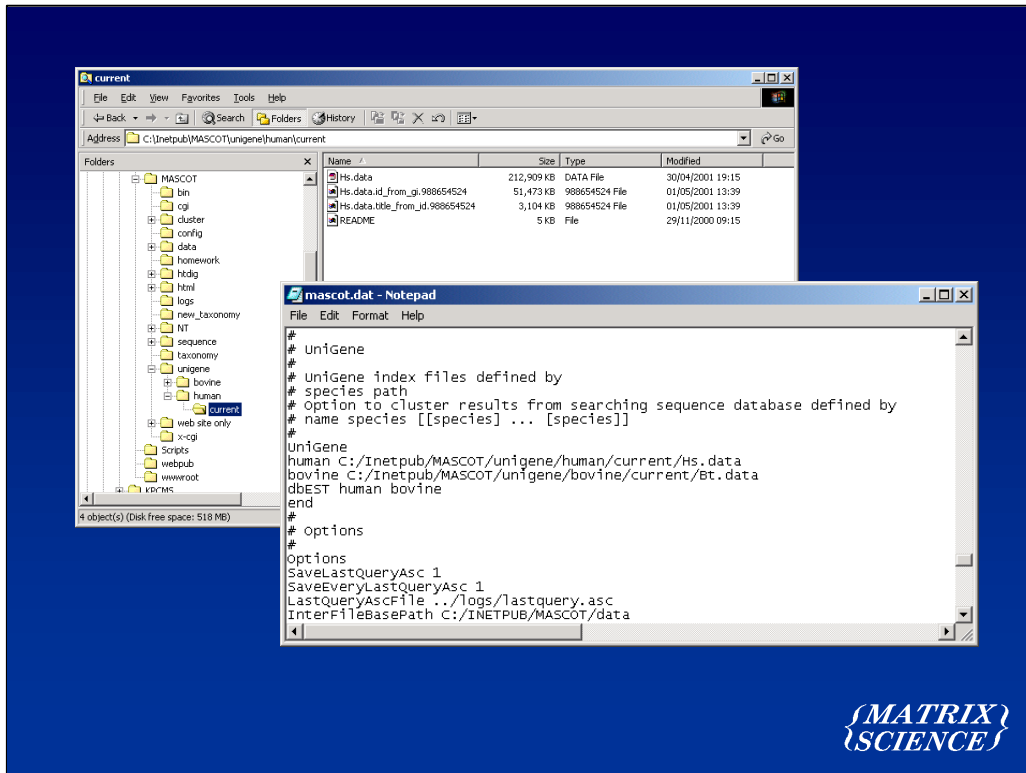
To address this problem, we have recently started using the UniGene index from the National Center for Biotechnology Information to simplify the search results.

UniGene is not a consensus sequence, it is an index which is created by BLASTing GenBank sequences against themselves to cluster them into gene families.



*MATRIX*  
*SCIENCE*

Unigene can be downloaded from the NCBI FTP site. Several important species are available: human, mouse, rat, cow and zebra fish.



To use a unigene index in Mascot, the data file for the species is downloaded and unpacked into a suitable directory structure...  
 Then a few lines are added to the Mascot configuration file, mascot.dat.

**Mascot Search Results**

User : JSC  
 Email : jcottrell@matrixscience.com  
 Search title : Annexin  
 MS data file : U:\Mascot test data\Glaxo\gtof10348.pk1  
 Database : dbEST 20001001 (35277150 sequences; 4704334140 residues)  
 Taxonomy : Homo sapiens (human) (14927850 sequences)  
 Timestamp : 16 Oct 2000 at 21:15:15 GMT

Significant hits:

- [Hs.78225](#) ANXA1 annexin A1
- [Hs.77290](#) TALD01 transaldolase 1
- [Hs.173205](#) NPM1 nucleophosmin (nucleolar phosphoprotein B23, numatrin)
- [Hs.75598](#) HNRPA2B1 heterogeneous nuclear ribonucleoprotein A2/B1
- [Hs.192023](#) EIF3S2 eukaryotic translation initiation factor 3, subunit 2 (beta, 36kD)
- [Hs.81361](#) HNRPA8 heterogeneous nuclear ribonucleoprotein A/B
- [Hs.82201](#) CSNK2A2 casein kinase 2, alpha prime polypeptide
- [Hs.250616](#) IDH3A isocitrate dehydrogenase 3 (NAD+) alpha
- [Hs.80828](#) KRT1 keratin 1 (epidermolytic hyperkeratosis)
- [Hs.278572](#) ALK anaplastic lymphoma kinase (Ki-1)
- [Hs.169476](#) GAPD glyceraldehyde-3-phosphate dehydrogenase
- [gi|9770088](#) 601066289F1 NIH MGC 10 Homo sapiens cDNA clone IMAGE:3452447 5'
- [Hs.156110](#) IGKC immunoglobulin kappa constant
- [Hs.37078](#) CRKL v-crk avian sarcoma virus CT10 oncogene homolog-like
- [gi|7309507](#) CM0-CT0341-260100-160-d10 CT0341 Homo sapiens cDNA
- [Hs.256309](#) DXS1357E accessory proteins BAP31/BAP29
- [Hs.99936](#) KRT10 keratin 10 (epidermolytic hyperkeratosis; keratosis palmaris et plantaris)
- [Hs.181165](#) EEF1A1 eukaryotic translation elongation factor 1 alpha 1
- [gi|7306319](#) RC0-BT0387-170100-011-a08 BT0387 Homo sapiens cDNA
- [Hs.89525](#) HDGF hepatoma-derived growth factor (high-mobility group protein 1-like)
- [Hs.289109](#) DDAH1 dimethylarginine dimethylaminohydrolase 1
- [Hs.217493](#) ANXA2 annexin A2
- [Hs.249495](#) HNRPA1 heterogeneous nuclear ribonucleoprotein A1
- [Hs.249247](#) FBRNP heterogeneous nuclear protein similar to rat helix destabilizing protein
- [Hs.183704](#) UBC ubiquitin C
- [Hs.65114](#) KRT18 keratin 18
- [Hs.278242](#) K-ALPHA-1 tubulin, alpha, ubiquitous
- [gi|10244362](#) RC0-AN0040-200700-022-f03 AN0040 Homo sapiens cDNA

Now, using the UniGene index as a lookup table, we can transform the results of a dbEST search.

This is now a much clearer picture, very similar to the protein database result. Please remember that we are not clustering the database sequences into consensus sequences prior to searching. This could lead to matches being missed. UniGene is being used after the search, to simplify the results.



**Mascot Search Results - Microsoft Internet Explorer**

34. [q116890858](#) Mass: 21281 Total score: 169 Peptides matched: 4  
 RCD-PT0006-271199-011-G06 PT0006 Homo sapiens cDNA  
 Check to include this hit in archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss Score	Rank	Peptide
<input checked="" type="checkbox"/> 17	466.16	930.31	930.47	-0.16	0	68	1 GPSSVEDIK
<input checked="" type="checkbox"/> 50	414.17	1239.48	1239.69	-0.21	2	18	4 SAPGG6SKVPQKK
<input checked="" type="checkbox"/> 118	910.23	1818.44	1818.84	-0.39	0	(60)	1 MTDQEAIQDLWQWR
<input checked="" type="checkbox"/> 120	910.24	1834.47	1834.83	-0.36	0	83	1 MTDQEAIQDLWQWR + 1 Oxidation (M)

**Proteins matching the same set of peptides:**  
[q115876282](#) Mass: 20076 NP1 nucleophosmin (nucleolar phosphoprotein B23, numatrin)  
[q1143808.v1](#) Norton Fer  
[q116890852](#) Mass: 21515 RCD-PT0006-271199-011

**Mascot Search Results - Microsoft Internet Explorer**

3. [Hs\\_173205](#) Total score: 372 Peptides matched: 10  
 Check to include this hit in archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss Score	Rank	Peptide
<input checked="" type="checkbox"/> 1	430.11	429.10	428.19	0.91	0	6	6 DAAP + 1 Acetyl (N-term); 1 Methyl ester (DE)
<input checked="" type="checkbox"/> 5	402.17	802.33	802.45	-0.12	0	33	1 TVSLGAGAK
<input checked="" type="checkbox"/> 17	466.16	930.31	930.47	-0.16	0	68	1 GPSSVEDIK
<input checked="" type="checkbox"/> 50	414.17	1239.48	1239.69	-0.21	2	18	4 SAPGG6SKVPQKK
<input checked="" type="checkbox"/> 118	910.23	1818.44	1818.84	-0.39	0	(60)	1 MTDQEAIQDLWQWR
<input checked="" type="checkbox"/> 120	910.24	1834.47	1834.83	-0.36	0	83	1 MTDQEAIQDLWQWR + 1 Oxidation (M)
<input checked="" type="checkbox"/> 140	1121.88	2241.75	2242.20	-0.46	0	70	1 HSNVQPTVSLGGFEITPPVWLR + 1 Oxidation (M)
<input checked="" type="checkbox"/> 141	748.26	2241.75	2242.20	-0.45	0	(52)	1 HSNVQPTVSLGGFEITPPVWLR + 1 Oxidation (M)
<input checked="" type="checkbox"/> 160	982.29	2943.84	2944.45	-0.61	1	45	1 TVSLGAGAKDELHIVEAEAMHYEGSPIK + 1 Oxidation (M)
<input checked="" type="checkbox"/> 161	736.97	2943.85	2944.41	-0.56	0	49	1 TVSLGAGAKDELHIVEAEAMHYEGSPIK + 1 Oxidation (M)

**Proteins matching the same set of peptides:**  
[q118115409](#) Mass: 20037 Total score: 166 Peptides matched: 6  
 Q01-PT0071-230200-083-c03 PT0071 Homo sapiens cDNA  
[q118115410](#) Mass: 19216 Total score: 166 Peptides matched: 6

**Mascot Search Results - Microsoft Internet Explorer**

36. [q1110315208](#) Mass: 2  
 601678979F1 NIH\_MGC\_S

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss Score	Rank	Peptide
<input checked="" type="checkbox"/> 5	402.17	802.33	802.45	-0.12	0	33	1 TVSLGAGAK
<input checked="" type="checkbox"/> 50	414.17	1239.48	1239.69	-0.21	2	18	4 SAPGG6SKVPQKK
<input checked="" type="checkbox"/> 140	1121.88	2241.75	2242.20	-0.46	0	70	1 HSNVQPTVSLGGFEITPPVWLR + 1 Oxidation (M)
<input checked="" type="checkbox"/> 141	748.26	2241.75	2242.20	-0.45	0	(52)	1 HSNVQPTVSLGGFEITPPVWLR + 1 Oxidation (M)
<input checked="" type="checkbox"/> 160	982.29	2943.84	2944.45	-0.61	1	(45)	1 TVSLGAGAKDELHIVEAEAMHYEGSPIK + 1 Oxidation (M)
<input checked="" type="checkbox"/> 161	736.97	2943.85	2944.45	-0.60	1	46	2 TVSLGAGAKDELHIVEAEAMHYEGSPIK + 1 Oxidation (M)

**Proteins matching the same set of peptides:**  
[q118115409](#) Mass: 20037 Total score: 166 Peptides matched: 6  
 Q01-PT0071-230200-083-c03 PT0071 Homo sapiens cDNA  
[q118115410](#) Mass: 19216 Total score: 166 Peptides matched: 6

**MATRIX SCIENCE**

When we look further down the report, at details of individual matches, we see the benefits of clustering the ESTs. Here we have two groups of matches from the dbEST search. These matches are well down the report, and appear to have little in common apart from a very weak match to query 50. There is no particular reason to connect these two hits.

However, when we look at the unigene report, we find that these matches all belong to the same gene, nucleophosmin. And, of course, because this gene now has 10 matches, it is listed near the top of the report.

Mascot Search Results: Protein View - Microsoft Internet Explorer

Address: http://dell500/mascot/cgi/protein\_view.pl?file=../data/20001016/F289840.dat&hit=Hs%2e78225&px=18&UNIGENE=human

## Mascot Search Results

### UniGene View

**ID** Hs.78225  
**TITLE** annexin A1  
**GENE** ANXA1  
**CYTOBAND** 9q12-q21.2  
**LOCUSLINK** 301  
**EXPRESS** ; Aorta; Bone; Brain; Breast; CNS; Colon; Ear; Esophagus; Eye; Foreskin; Gall bladder; Heart; Kidney; Lung; Lymph; Nose; Omentum;  
**CHROMOSOME** 9  
**STS** ACC=GO6372 NAME=SHGC-12349 UNISTS=78440  
**STS** ACC=- NAME=A009X30 UNISTS=5599  
**STS** ACC=G32952 NAME=A009X30 UNISTS=117530  
**STS** ACC=- NAME=sts-H67867 UNISTS=29146  
**STS** ACC=- NAME=H29761 UNISTS=3341  
**TXMAP** D9S166-D9S1876; MARKER=SHGC-12349; RHPANEL=G3  
**TXMAP** D9S1876-D9S175; MARKER=WI-7046; RHPANEL=GB4  
**TXMAP** D9S1876-D9S175; MARKER=A009X30; RHPANEL=GB4  
**TXMAP** D9S1876-D9S175; MARKER=sts-H67867; RHPANEL=GB4  
**TXMAP** D9S1876-D9S175; MARKER=H29761; RHPANEL=GB4  
**PROTSIM** ORG=Caenorhabditis elegans; PROTI=482227; PROTIID=pir:S41022; PCT=42; ALN=311  
**PROTSIM** ORG=Homo sapiens; PROTI=6729710; PROTIID=pcdb:1B09; PCT=100; ALN=345  
**PROTSIM** ORG=Mus musculus; PROTI=71759; PROTIID=pir:LUMS1; PCT=87; ALN=345  
**PROTSIM** ORG=Rattus norvegicus; PROTI=71758; PROTIID=pir:LURT1; PCT=89; ALN=345  
**SCOUNT** 722  
**SEQUENCE** ACC=BC001275; NID=g12654862; PID=g12654863  
**SEQUENCE** ACC=XM\_005665; NID=g11429703; PID=g11429704  
**SEQUENCE** ACC=AW950757; NID=g8140416; LID=2983  
**SEQUENCE** ACC=AX004164; NID=g9927714; PID=g9927715  
**SEQUENCE** ACC=AV734845; NID=g10852390; CLONE=cdAAHD12; END=5'; LID=4713  
**SEQUENCE** ACC=BF218154; NID=g11111740; CLONE=IMAGE:4104209; END=5'; LID=3915; MGC=4502100  
**SEQUENCE** ACC=XD5908; NID=g34387; PID=g34388  
**SEQUENCE** ACC=BF241022; NID=g11154947; CLONE=IMAGE:4109288; END=5'; LID=3917; MGC=4502100  
**SEQUENCE** ACC=BE171982; NID=g8634708; LID=3484  
**SEQUENCE** ACC=AV717869; NID=g10815021; CLONE=DCBARH06; END=5'; LID=4704  
**SEQUENCE** ACC=AB973706; NID=g8164890; LID=3051

When you click on the accession number link of a unigene filtered report, you get full details for that particular gene family.

## ***Human genome***

- ~  $3 \times 10^9$  bases
- (dbEST was  $3.4 \times 10^9$  bases on 1 May 2001)
- ~  $6 \times 10^9$  residues in 6 frame translation
- 99.75% of translated sequence is non-coding
- $1.5 \times 10^5$  tryptic limit peptides of 1500 Da  $\pm$  0.5
- $6 \times 10^7$  no-enzyme peptides of 1500 Da  $\pm$  0.5

*{MATRIX}*  
*{SCIENCE}*

We can also perform MS/MS searches on continuous raw genomic sequence data. The recent availability of a draft assembly of the human genome has made this a focus of great interest. Let's just look at some numbers.

The human genome assembly is approximately 3 billion bases which makes it similar in size to dbEST.

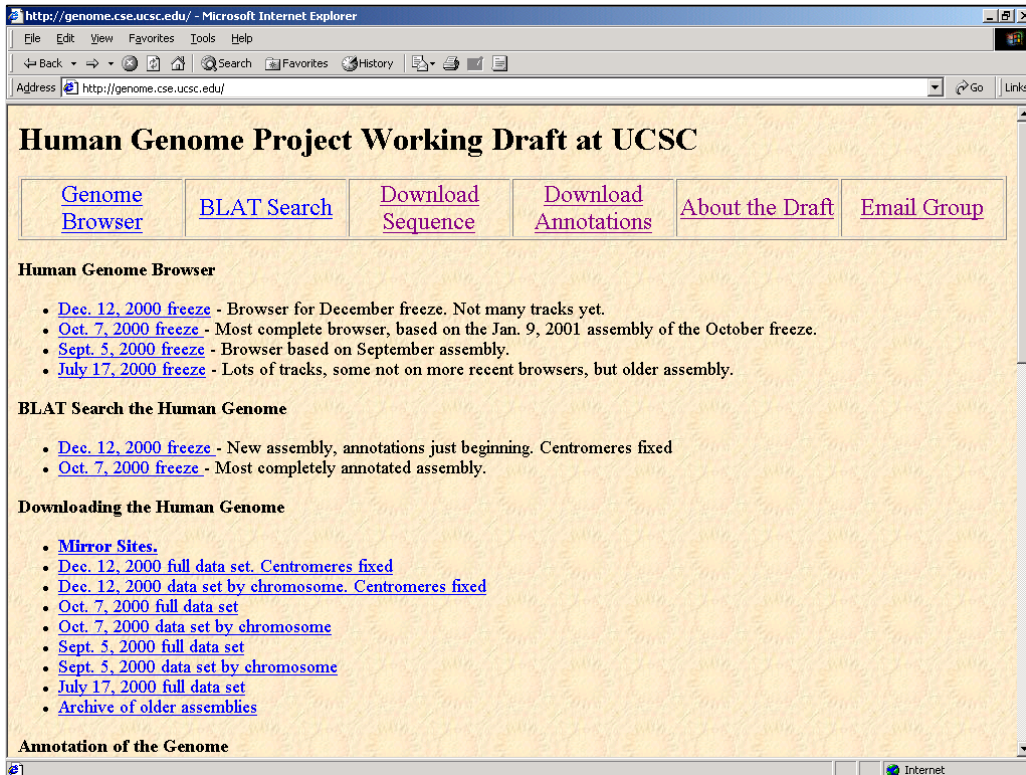
Since we must translate in all 6 reading frames, this corresponds to 6 billion amino acid residues.

In the human genome, only 1.5% of the sequence codes for proteins. Conversely, 99.75% of the translated sequence is non-coding and simply contributes to the background of random matches. This is a severe test of the discrimination of the scoring scheme.

This is only slightly worse than dbEST. Although dbEST is essentially all coding sequences, after translation in 6 frames, we know that 83% (5/6) must be junk.

If we are matching MS/MS data from a tryptic peptide of nominal mass 1500 Da against the human genome, we are going to have to test 150 thousand peptides. Which sounds bad...

but is not nearly as bad as the no-enzyme case where we have to test 60 million.



The draft assembly can be downloaded from UC Santa Cruz GoldenPath web site

**Index of /goldenPath/12dec2000/bigZips**

Name	Last modified	Size	Descr
<a href="#">Parent Directory</a>	29-Mar-2001 03:38	-	
<a href="#">chromApp.zip</a>	05-Apr-2001 16:18	4.7M	
<a href="#">chromFa.zip</a>	05-Apr-2001 16:58	801M	
<a href="#">chromFaMasked.zip</a>	05-Apr-2001 16:22	508M	
<a href="#">chromOut.zip</a>	05-Apr-2001 17:45	71.7M	
<a href="#">contigApp.zip</a>	29-Mar-2001 03:51	4.7M	
<a href="#">contigFa.zip</a>	31-Mar-2001 21:26	799M	
<a href="#">contigFaMasked.zip</a>	31-Mar-2001 21:37	506M	
<a href="#">liftAll.zip</a>	11-Apr-2001 10:32	12k	

This directory contains the Dec. 12, 2000 Genbank freeze ordered and oriented according to the corresponding fingerprint map and genome layout from Wash. U., taking into account overlap between fragments and bridging mRNA, EST, plasmid and BAC end pairs. Chromosomes 21 and 22 are the finished versions as obtained from NCBI. This directory includes the following files:

contigFa.zip - The working draft sequence one file per fingerprint contig (clone layout) in Fasta format. Unpacks with one directory for each chromosome and one subdirectory for each fingerprint contig.

contigFaMasked.zip - RepeatMasked version of contigFa.zip (with N's in place of repeating elements.)

You can download the assembly as chromosome length sequences or as collections of contigs. Searching the chromosome length sequences in Mascot is possible, but not advisable.

## ***Chromosome length sequences***

- **Chromosome 1 is 285 Mbp**
- **Mascot workspace scales as the length of the longest sequence. Hence need lots of physical RAM to avoid disk thrashing**
- **More efficient to work with smaller pieces, e.g. 600 kbp with 600 bp overlaps**

*{MATRIX}*  
*{SCIENCE}*

The longest human chromosome is chromosome 1, 285 million bp. Mascot requires a significant memory overhead to manipulate such long sequences, which means that unless you have a very large amount of RAM, the search is going to be using virtual memory ... i.e. swapping out to disk ... and run relatively slowly. So, we recommend working with contigs or just chopping the chromosomes into more manageable lengths with small overlaps. In any case, we don't know of any tools for reviewing the results which can handle 250 Mbp sequences.

Mascot Search Results - Microsoft Internet Explorer

Address [http://dell5000/mascot/cgi/master\\_results.pl?file=../data/20001207/F001327.dat](http://dell5000/mascot/cgi/master_results.pl?file=../data/20001207/F001327.dat)

**MASCOT**  
*(SCIENCE)* **Mascot Search Results**

User : JSC  
 Email : JSC@work  
 Search title : Annexin mix: Oct 7th  
 MS data file : U:\Mascot test data\Glaxo\gtof10348.pk1  
 Database : HG (34668 sequences: 6913762562 residues)  
 Timestamp : 8 Dec 2000 at 09:50:44 GMT

Significant hits:

<a href="#">chr9_123</a>	73200001-73800598
<a href="#">chr7_46</a>	27000001-27600598
<a href="#">chr12_94</a>	55800001-56400598
<a href="#">chr16_117</a>	69600001-70200599
<a href="#">chr15_127</a>	75600001-76200599
<a href="#">chr5_175</a>	104400001-105000598
<a href="#">chr1_60</a>	35400001-36000599
<a href="#">chr8_115</a>	68400001-69000597
<a href="#">chr10_132</a>	78600001-79200599
<a href="#">chr2_354</a>	211800001-212400597
<a href="#">chr22_30</a>	17400001-18000597
<a href="#">chr7_66</a>	39000001-39600598
<a href="#">chr10_177</a>	105600001-106200599
<a href="#">chr2_312</a>	186600001-187200597
<a href="#">chr6_70</a>	41400001-42000599
<a href="#">chr15_164</a>	97800001-98400599
<a href="#">chr12_98</a>	58200001-58800598
<a href="#">chr9_152</a>	90600001-91200598
<a href="#">chr11_137</a>	81600001-82200599
<a href="#">chr13_129</a>	76800001-77400599
<a href="#">chrX_208</a>	124200001-124800597
<a href="#">chr1_215</a>	128400001-129000599

**Probability Based Mowse Score**

Score is  $-10 \cdot \log(P)$ , where P is the probability that the observed match is a random event.  
 Individual ions scores > 68 indicate identity or extensive homology ( $p < 0.05$ ).

Done Local intranet

This is the result of searching our standard data set against the unmasked human genome assembly.

Mascot Search Results - Microsoft Internet Explorer

Address: http://dell5000/mascot/cgi/master\_results.pl?file=../data/20001207/F001327.dat

1. [chr9\\_123](#) Total score: 491 Peptides matched: 13  
 Could not retrieve title string  
 Check to include this hit in archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Rank	Peptide
<a href="#">1</a>	430.11	429.10	429.22	-0.12	0	7	4	DAPK
<a href="#">11</a>	415.19	828.36	828.50	-0.13	0	32	6	VLDELEK
<a href="#">12</a>	415.19	828.36	828.51	-0.14	0	33	2	NALLSLAK
<input checked="" type="checkbox"/> <a href="#">45</a>	607.16	1212.31	1212.53	-0.21	0	70	1	DITSDTSGDFR
<input checked="" type="checkbox"/> <a href="#">53</a>	631.70	1261.38	1261.59	-0.22	0	69	1	TPAQFDADELK
<input checked="" type="checkbox"/> <a href="#">65</a>	457.85	1370.54	1370.77	-0.23	1	46	1	VLDELEKGDIEK
<input checked="" type="checkbox"/> <a href="#">93</a>	775.76	1549.50	1549.81	-0.31	0	69	1	GTDVVVFNTILTR
<input checked="" type="checkbox"/> <a href="#">98</a>	547.49	1639.45	1639.77	-0.32	1	(41)	1	DLAKDITSDTSGDFR
<input checked="" type="checkbox"/> <a href="#">99</a>	820.75	1639.48	1639.77	-0.29	1	52	1	DLAKDITSDTSGDFR
<a href="#">101</a>	840.24	1678.47	1677.90	0.57	1	19	7	RGTDVVVFNTILTR
<input checked="" type="checkbox"/> <a href="#">103</a>	851.77	1701.52	1701.88	-0.36	0	82	1	GLGTDEDLIEILASR
<input checked="" type="checkbox"/> <a href="#">105</a>	870.21	1738.41	1738.73	-0.32	0	99	1	SEDFGVNEDLADSDAR
<input checked="" type="checkbox"/> <a href="#">149</a>	785.91	2354.72	2355.15	-0.43	0	66	1	GGPGSAVSPYPTFNPSSDVAALHK

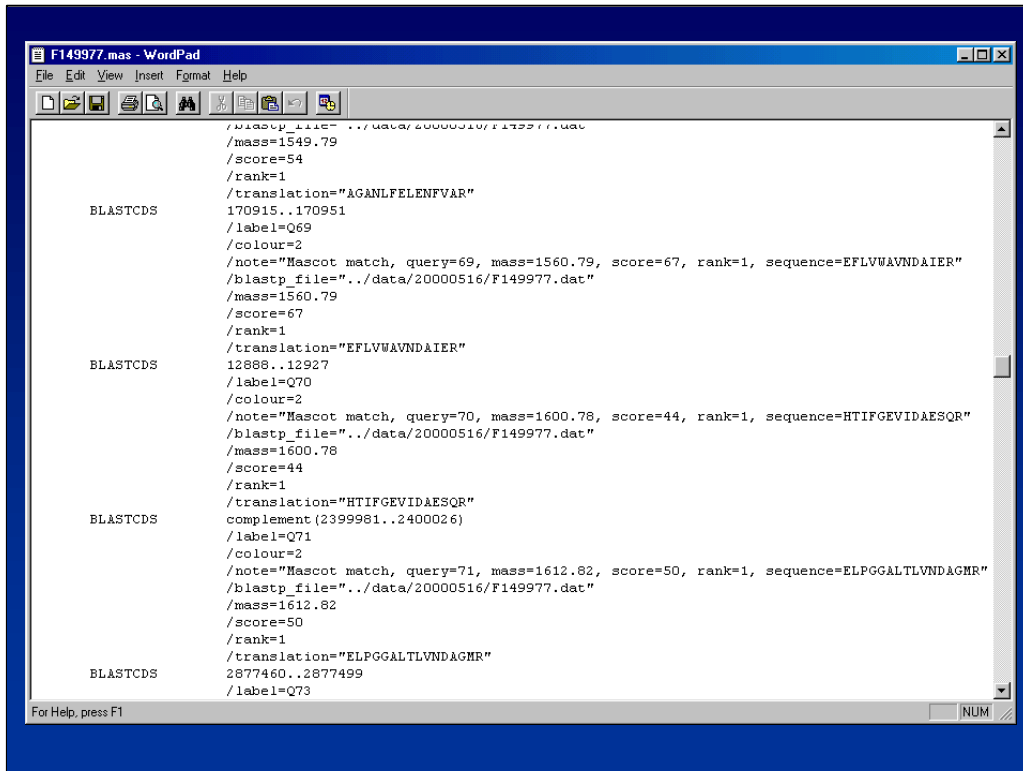
2. [chr7\\_46](#) Total score: 204 Peptides matched: 7  
 Could not retrieve title string  
 Check to include this hit in archive report

Query	Observed	Mr(expt)	Mr(calc)	Delta	Miss	Score	Rank	Peptide
<a href="#">1</a>	430.11	429.10	429.22	-0.12	0	7	4	DAPK
<input checked="" type="checkbox"/> <a href="#">33</a>	544.15	1086.29	1086.48	-0.19	0	34	1	NYYEQWCK
<input checked="" type="checkbox"/> <a href="#">43</a>	594.71	1187.41	1187.64	-0.23	0	66	1	IDTIEIITDR
<input checked="" type="checkbox"/> <a href="#">66</a>	689.19	1376.36	1376.62	-0.26	0	44	1	GGGGNFPGPGSNFR
<input checked="" type="checkbox"/> <a href="#">102</a>	565.82	1694.44	1694.76	-0.32	0	42	1	GFGFVTFDDHDPVDK
<input checked="" type="checkbox"/> <a href="#">117</a>	899.82	1797.63	1797.91	-0.29	0	63	1	LFIGGLSFETTEESLR

Done Local intranet

The tabular reports which we and others have developed for reporting results from protein and dbEST databases are just not suitable for a sequence the length and complexity of a human chromosome.



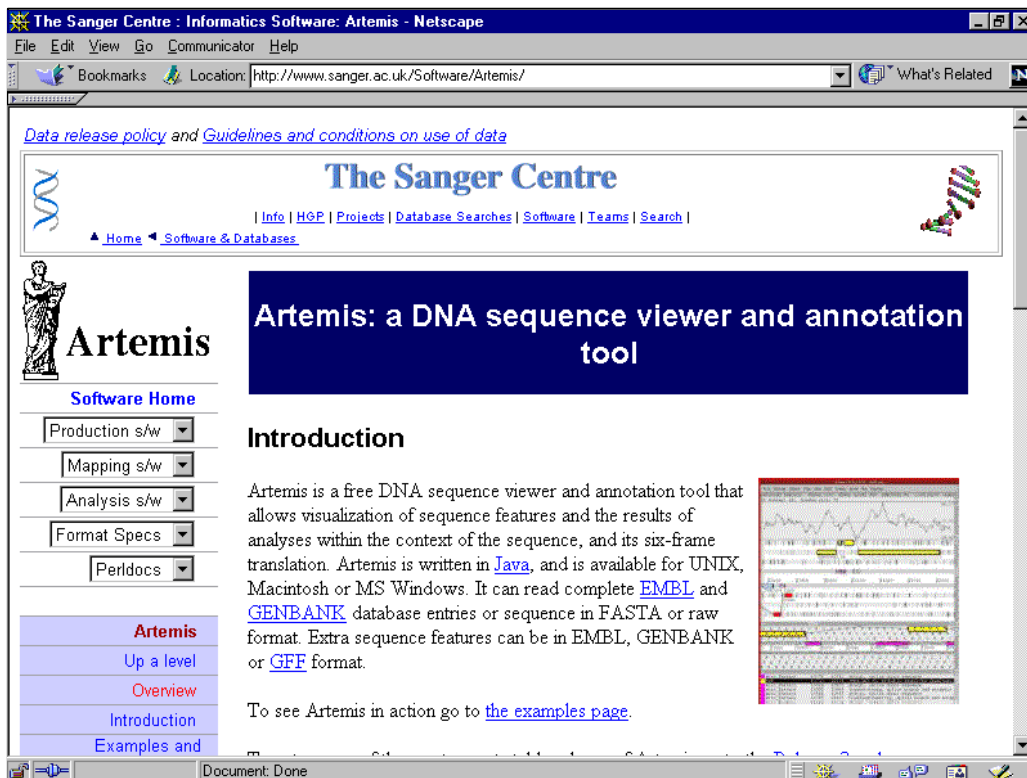


The screenshot shows a WordPad window titled "F149977.mas - WordPad" with a menu bar (File, Edit, View, Insert, Format, Help) and a toolbar. The main text area contains BLAST results in EMBL/GenBank format. The results are organized into four entries, each starting with "BLASTCDS".

```
/blastp_file=../data/20000516/F149977.dat
/mass=1549.79
/score=54
/rank=1
/translation="AGANLFELENFVAR"
BLASTCDS
170915..170951
/label=Q69
/colour=2
/note="Mascot match, query=69, mass=1560.79, score=67, rank=1, sequence=EFLVWAVNDAIER"
/blastp_file=../data/20000516/F149977.dat"
/mass=1560.79
/score=67
/rank=1
/translation="EFLVWAVNDAIER"
BLASTCDS
12888..12927
/label=Q70
/colour=2
/note="Mascot match, query=70, mass=1600.78, score=44, rank=1, sequence=HTIFGEVIDAESQR"
/blastp_file=../data/20000516/F149977.dat"
/mass=1600.78
/score=44
/rank=1
/translation="HTIFGEVIDAESQR"
BLASTCDS
complement(2399981..2400026)
/label=Q71
/colour=2
/note="Mascot match, query=71, mass=1612.82, score=50, rank=1, sequence=ELPGGALTLVNDAGMR"
/blastp_file=../data/20000516/F149977.dat"
/mass=1612.82
/score=50
/rank=1
/translation="ELPGGALTLVNDAGMR"
BLASTCDS
2877460..2877499
/label=Q73
```

At the bottom left of the window, it says "For Help, press F1". At the bottom right, there is a "NUM" button.

For very long DNA sequences, such as this, what we have done is to switch from our standard protein view report to outputting the peptide match results as an EMBL / GenBank format feature table. This may not look very friendly, but the advantage is that this report can now be read into a standard genome browser.



Having generated a feature table, we can now use a genome browser to view it. One which we find works well is Artemis, a Java based genome browser developed and distributed by the Sanger Centre.



Here is an Artemis screenshot showing three views of a portion of the genome. In the upper third, we have a low resolution view. This can be zoomed out to show an entire sequence as a single strip. We have the forward and complementary DNA strands, and the 6 frame translation. The vertical bars are stop codons. The yellow blocks are exons, while the blue blocks here are coding sequences. Individual Mascot peptide matches are shown in red. This particular gene has 8 peptide matches.

The middle third is a similar arrangement, but at high enough resolution to see individual bases and residues.

Finally, the lower third shows a tabular view of the feature table. When a match is selected, it is highlighted in all three views, and we can see the spectrum number, sequence, molecular weight, Mascot score, etc.

Not only does this allow us to zoom and pan around these extremely long sequences, it also allows us to view the peptide matches found by Mascot in the context of all the existing annotations. This gives us a powerful way to present the results of MS based searching complete genomes.

Key	Category	MSDB	dbEST	HG
a	Top match with significant ions score	74	56	33
b	Top match, but ions score not significant	26	37	13
c	Not top match and ions score not significant	10	11	11
d	No match because of higher scoring non-significant matches	-	6	11
e	No match because peptide sequence not found in MSDB	4	-	-
f	No match because peptide sequence not found in dbEST	-	2	-
g	No match because coding sequence substantially missing from HG	-	-	15
h	No match because coding sequence poorly aligned in HG	-	-	10
i	No match because peptide spans exon / intron boundary in HG	-	-	19
k	No match because peptide results from non-tryptic post-translational processing	-	2	2

*{MATRIX}*  
*{SCIENCE}*

When we make a detailed comparison of the results from searching the same data against the three different types of database, the major differences are caused by two factors.

First, the human genome assembly is only a draft assembly and, at the time we did this study, there were complete mRNA's which were either poorly aligned or even missing. This accounted for 25 "missing" peptide matches. Obviously, this situation will change over the coming months as the assembly is refined.

The second factor will not change, if we choose to search the raw genomic sequence. Approximately one quarter of peptide matches are missed because they span exon / intron boundaries. This is not a severe problem if we have multiple peptides from the protein, but is clearly a limitation.

## ***HGP Draft Assembly***

- **Searching complete chromosome entries is possible, but unwieldy.**
- **Scoring statistics very similar to dbEST**
- **Some well characterised proteins are 'missing' (e.g. transaldolase)**
- **When error rate is high, better off searching redundant EST's than a single consensus sequence**
- **Still too early to use public draft assembly for routine protein ID?**

*{MATRIX}*  
*{SCIENCE}*

The main conclusions of our database comparison study are listed.

***MATRIX***  
***SCIENCE***

<http://www.matrixscience.com>