

Mining the Human Genome Using Protein Structure Homology

Randal R. Ketchum, Ph.D.
Immunex Corporation

Need for gene mining

Scale of problem

Protein structure

Structure prediction

Mining the genome

Some results

Some problems

Future work

```
TCTCGAGGGCCACGCGTTTAAACGTTCGAGGTACCTATCCCGGGCCGCCAC
CATGGCTACAGGCTCCCGGACGTCCCTGCTCCTGGCTTTTGGCCTGCTCT
GCCTGCCCTGGCTTCAAGAGGGCAGTGCAACTAGTTCTGACCGTATGAAA
CAGATAGAGGATAAGATCGAAGAGATCCTAAGTAAGATTTATCATATAGA
GAATGAAATCGCCCGTATCAAAAAGCTGATTGGCGAGCGGACTAGATCTA
GTTTGGGGAGCCGGGCATCGCTGTCCGCCAGGAGCCTGCCAGGAGGAG
CTGGTGGCAGAGGAGGACCAGGACCCGTCCGAACTGAATCCCAGACAGA
AGAAAGCCAGGATCCTGCGCCTTTCCTGAACCGACTAGTTCCGGCCTCGCA
GAAGTGCACCTAAAGGCCGAAAACACGGGCTCGAAGAGCGATCGCAGCC
CATTATGAAGTTCATCCACGACCTGGACAGGACGGAGCGCAGGCAGGTGT
GGACGGGACAGTGAGTGGCTGGGAGGAAGCCAGAATCAACAGCTCCAGCC
CTCTGCGCTACAACCGCCAGATCGGGGAGTTTATAGTCACCCGGGCTGGG
CTCTACTACCTGTACTGTCAGGTGCACTTTGATGAGGGGAAGGCTGTCTA
CCTGAAGCTGGACTTGCTGGTGGATGGTGTGCTGGCCCTGCGCTGCCTGG
AGGAATTCTCAGCCACTGCGGCCAGTTCCCTCGGGCCCCAGCTCCGCCTC
TGCCAGGTGTCTGGGCTGTTGGCCCTGCGGCCAGGGTCCCTCCCTGCGGAT
CCGCACCCTCCCCTGGGCCATCTCAAGGCTGCCCCCTTCTCACCTACT
TCGGACTCTTCCAGGTTCACTGAGCGGCCGCGGATCTGTTTAAACTAG
```

```
MATGSRTSLLLAFGLLCLPWLQEGSATSSDRMKQIEDKIEEILSKIYHIE
NEIARIKKLIGERTRSSLGSRASLSAQEPAQEELVAEEDQDPSELNPQTE
ESQDPAPFLNRLVRRRSAPKGRKTRARRAIAAHYEVHPRPGDGAQAGV
DGTVSGWEEARINSSPLRYNRQIGEFIVTRAGLYYLYCQVHFDEGKAVY
LKLDDLVDGVLALRCLLEEFSAATAASSLGPQLRLCQVSGLLALRPGSSLRI
RTL PWAHLKAAPFLTYFGLFQVH
```

Need For Gene Mining

Human Genome contains approximately 30-60 thousand genes

Only 30-40% of these are classified into known function families

Function of proteins needed to enable development of therapeutics



Need For Gene Mining

Experimental methods too slow for complete classification

Computational methods for elucidating function needed

Weeks or months, around \$100K, to solve single, globular structure



NIGMS Structural Genomics Initiative

Proteins fold into a limited number of shapes

Estimates of ~10K protein folds - ~700 currently in the PDB

Solve key structures within families - homology can be used for rest

Around 10 years to solve 10K unique structures

Problem - many proteins have same fold with little or no sequence homology

Scale of the Problem

~15K structures in the Protein Data Bank

Around 4K are unique (< 90% identical)

This represents ~1500 families and ~700 folds

Less than 10% of all chains discovered in 2001 were new folds

So - many genes are for unknown function with no hope of change in the near future

Family: Short-chain cytokines

Lineage:

1. Root: scop

2. Class: All alpha proteins

3. Fold: 4-helical cytokines

core: 4 helices; bundle, closed; left-handed twist; 2 crossover connections

4. Superfamily: 4-helical cytokines

there are two different topoisomers of this fold with different entanglements of the two crossover connections

5. Family: Short-chain cytokines

Protein Domains:

1. Erythropoietin

long chain cytokine with a short-chain cytokine topology

1. Human (Homo sapiens) (3)

2. Granulocyte-macrophage colony-stimulating factor (GM-CSF)

1. Human (Homo sapiens) (2)

3. Interleukin-4 (IL-4)

1. Human (Homo sapiens) (13)

4. Interleukin-5

intertwined dimer

1. Human (Homo sapiens) (1)

5. Macrophage colony-stimulating factor (M-CSF)

forms dimer similar to the Flt3 ligand and SCF dimers

1. Human (Homo sapiens) (1)

Etc.

Four levels of protein structure

Primary - amino acid sequence

>1csgA

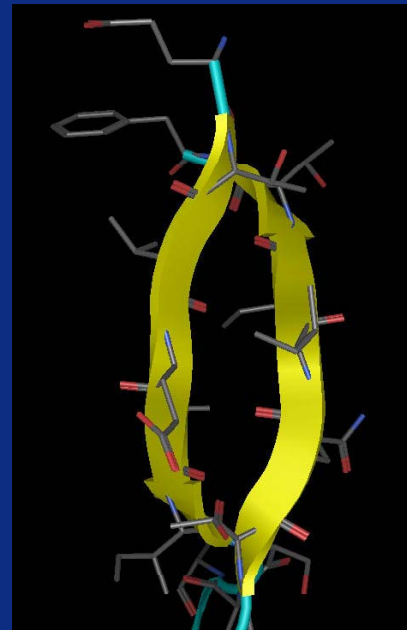
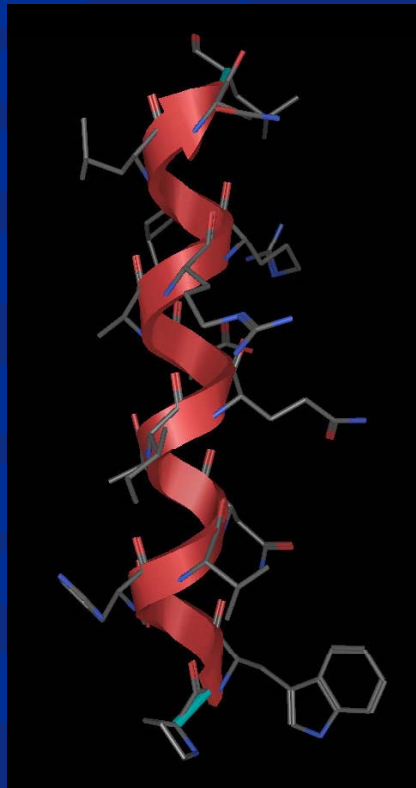
SPSPSTQPWEHVNAIQEARRLLNLSRD TAAEMNETVEVISEMFDLQEPTC
LQTRLELYKQGLRGS LTKLKGPLTMMASHYKQHCPPTPETS CATQIITFE
SFKENLKDFLLV I PFDCWEP

GRANULOCYTE-MACROPHAGE COLONY-STIMULATING FACTOR (GM-CSF) HUMAN
(HOMO SAPIENS) RECOMBINANT FORM EXPRESSED IN (ESCHERICHIA COLI)
M.R.WALTER, W.J.COOK, S.E.EALICK

Protein Structure

Four levels of protein structure

Secondary - local structure such as α helices and β strands



Protein Structure

Four levels of protein structure

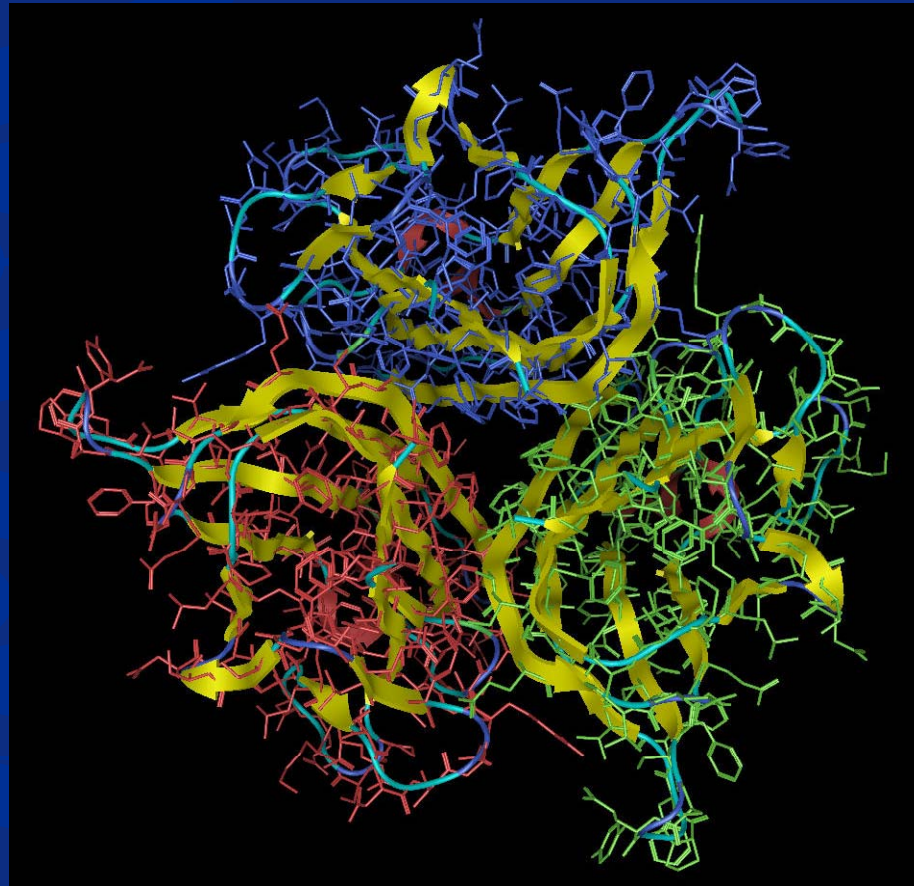
Tertiary - packing secondary structure elements into domains



Protein Structure

Four levels of protein structure

Quaternary - multiple chains



IMMUNEX®

Experimental Structure

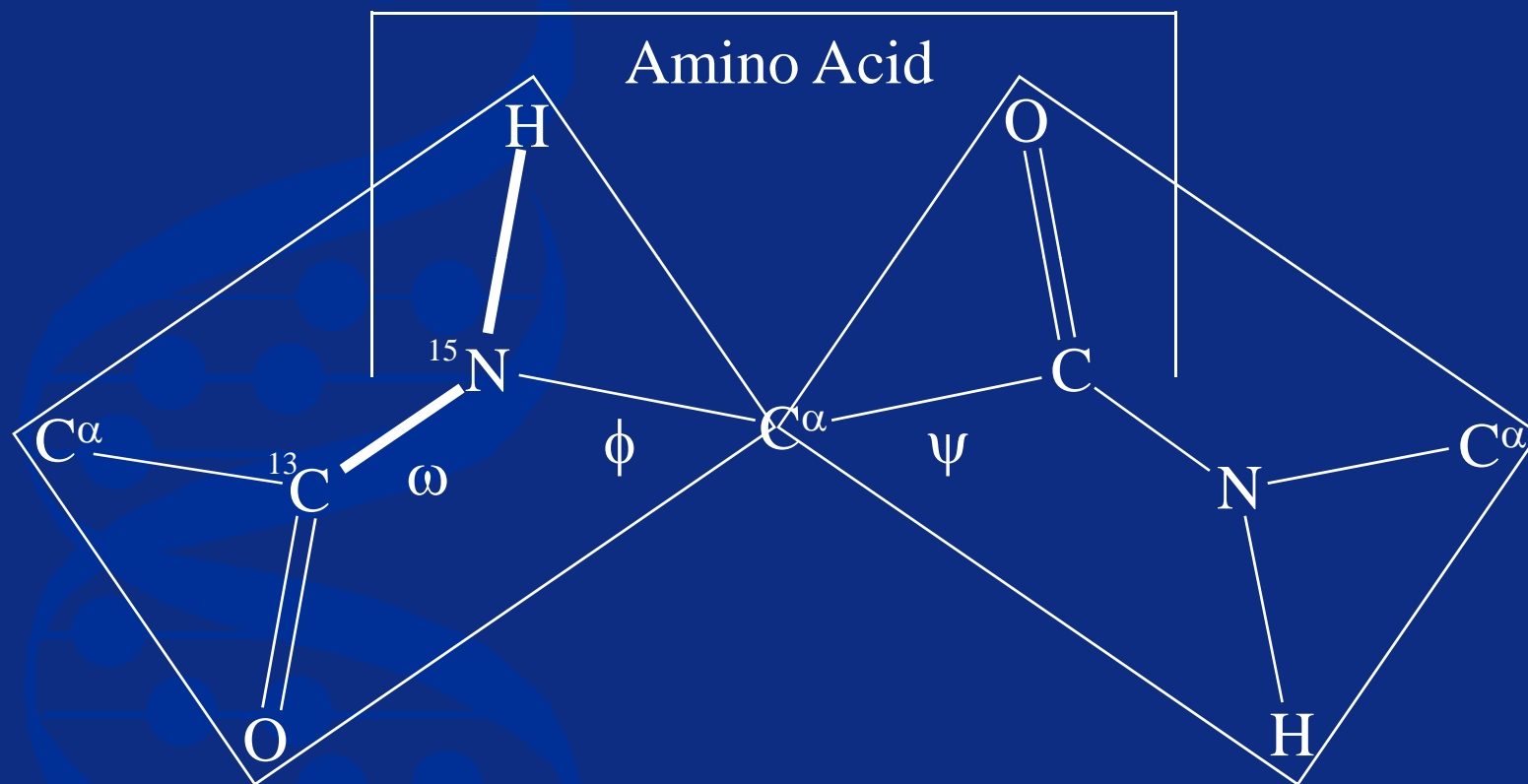
Proteins too small to see

Solid State NMR

Solution NMR

X-Ray Crystallography

Backbone consists of diplanes



Bond angles measurable to external magnetic field

Two intersecting vectors defines plane orientation

Join planes to determine dihedrals

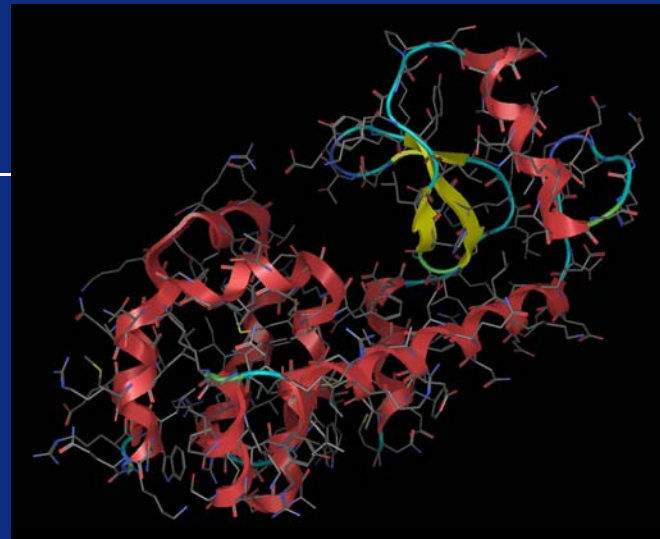
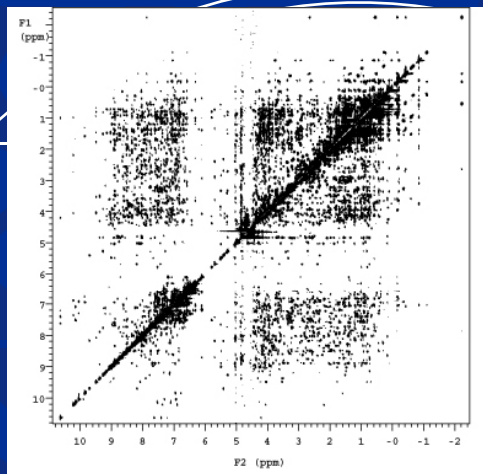

$$\Delta\nu = \nu_{\parallel} (3\cos^2\theta - 1)$$

Magnetization transfers between nuclei

Distance dependent

Assign measured NOE's to atoms

Fold structure using Distance Geometry



X-Ray Crystallography

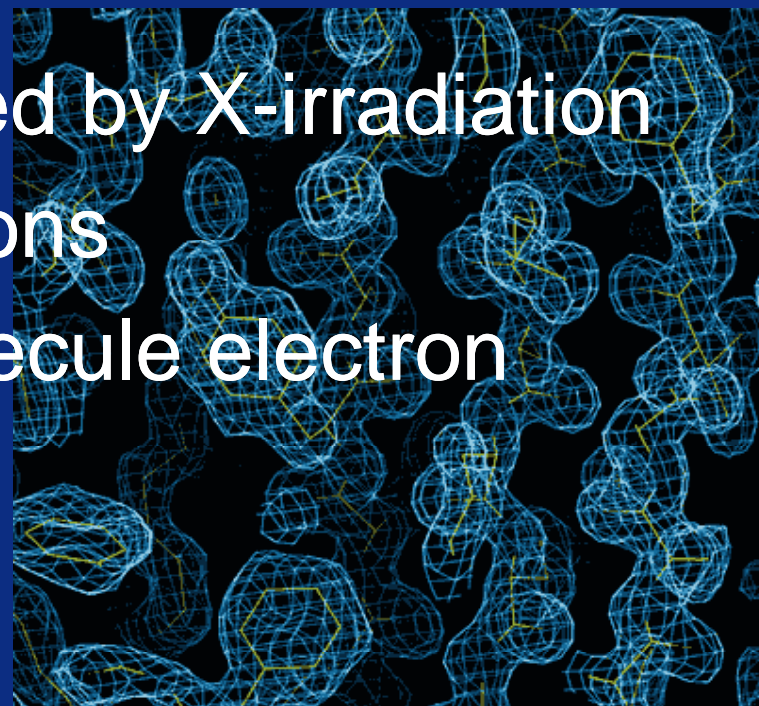
Molecule crystallized, crystals singular,

perfect quality

Diffraction pattern produced by X-irradiation

X-rays diffracted by electrons

Result is 3D image of molecule electron
clouds



Homology Modeling

Align sequence with unknown structure to
sequence with known structure

Extract structural parameters from known
and apply to unknown

Evaluate, modify alignment, and repeat

Higher homology produces more accurate
homology model



Structure Prediction

Homology modeling is routine with
sequence identity $> 30\%$

Less than 25% homology is termed the
twilight zone and requires other methods

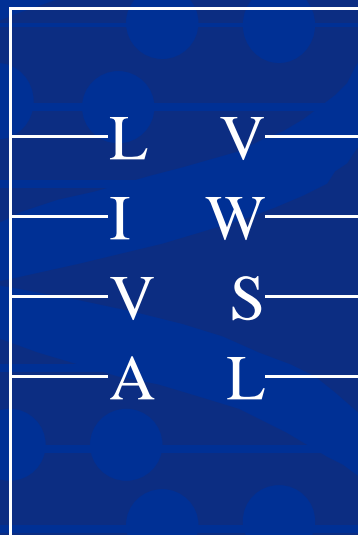
Protein Structure Prediction Using Inverse
Folding (Threading)



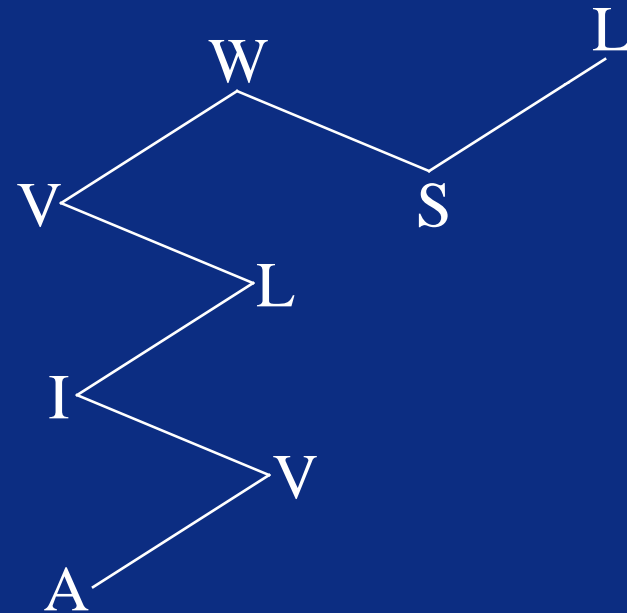
Threading

“Thread” a protein sequence onto a known structure

Score the threaded fold



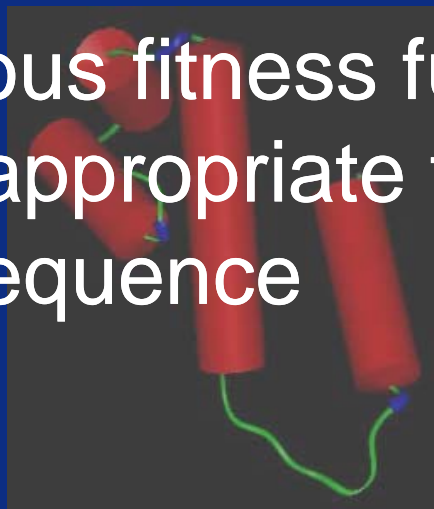
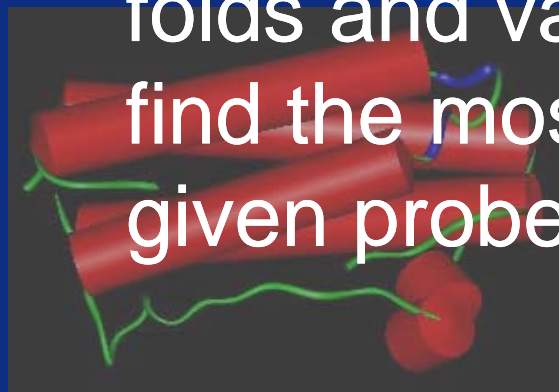
Happy



Sad

GeneFold Threading

Uses a representative library of protein folds and various fitness functions to find the most appropriate fold for a given probe sequence



```
KPAAHLIGDPSKQNSLLWRANTDRAFLQDGFSLSNNSLLVPTSGIYFVYSQVVFSGKAYS  
PKATSSPLYLAHEVQLFSSQYPFHVPLLSSQKMVYPGLQEPWLHSMYHGAAFQLTQGDQL  
STHTDGIPHLVLS PSTVFFGAFAL
```

GeneFold Threading

Describes each template protein in terms of:
of:

Sequence

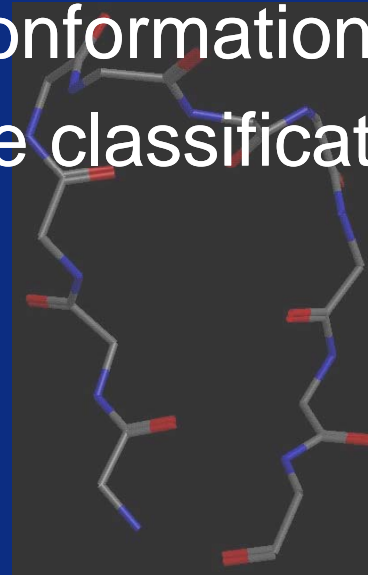
Burial pattern of residues

Local main chain conformation

Secondary structure classification



KPAAHLIGDPSKQNSLLWRANT
DRAFLQDGFSLSNLLVPTSG
IYFVYSQVVFSGKAYS



GeneFold Threading

Structure database based on PDB

Clustered by 50% sequence identity

Theoretical, long (>900) and short (<40) structures removed

1500 Clusters - highest resolution

structure chosen as representative (if no x-ray, choose NMR - grr)

GeneFold Threading

Scores a target sequence using:

Sequence-sequence: No structural information

Sequence-structure: Pseudo-energy of a single residue mounted in the template structural environment

Structure-structure: Comparison between predicted and actual secondary structure

GeneFold Threading

Three scoring methods

Sequence similarity: sequence term only

Hybrid sequence/structure similarity:
sequence, local conformation and burial

Full hybrid: Sequence, secondary structure,
local conformation and burial

GeneFold Threading

No one method produces a reliable prediction, but different methods give consistently correct answers

Jury Prediction

Two methods agree or

One of the three has a high reliability

GeneFold Scores

A given probe is aligned with every template and scored

P-value is calculated for alignment ensemble using distribution of scores

The inverse of the P-value is reported

This process is repeated independently for the three methods

Mining the Genome

Database of all gene predictions
translated to protein sequences

Calculate GeneFold scores for each
sequence

Relate interesting families using known
proteins

Search by family

An example: Mining the Family of Interleukins

Celera Genefold Data

Celera human r26b and mouse r12 Otto predictions and GeneFold 6.7
[instructions](#)

Enter a Celera ID (HCP...):

Human: Mouse:

Sort by:

Or, Select a GeneFold family as related to ProtBase (ProtBase Category: Possible GeneFold family):

•Bear in mind that this is merely an alternate method of choosing the GeneFold family. As such, several ProtBase categories map to the same GeneFold family, and therefore provide an identical list of Celera id's, regardless of belonging to different ProtBase categories. For example, 4BHC:1lki_CYTOKINE is identical to RTK-CSF:1lki_CYTOKINE.

•This list is prepared by running all ProtBase proteins classified as known through GeneFold and selecting the strong hits from those runs. The hits are then sorted and the known assignment is associated with its possible GeneFold families. This is merely a help in choosing GeneFold families to mine. The comprehensive list of possible GeneFold families is available below.

•The listed families contain PDB ID's. The first four characters are the PDB ID. The last character is the chain. An underscore indicates that there is a single chain for this ID. You can get details for a PDB ID at <http://www.rcsb.org/pdb/>

FIL:1itn_ BINDING PROTEIN
FIL:2frtE COMPLEX (RECEPTOR/IMMUNOGLOBULIN)
FIL:1ita_ CYTOKINE
FIL:7i1b_ CYTOKINE
FILR:1ic1B CELL ADHESION
FILR:1vscB CELL ADHESION PROTEIN

Browse the hits for the selected PDB chain

Celera IDs for which family 'lita_ CYTOKINE' is possible:

GeneBase info color code:

Known (and source not celera or sanger)
 Known and Categorized
 Unknown

Contig numbers are relative to the chromosome. Protein numbers are relative to the contig, with the exons ordered, begin being the begin of the first exon, end being the end of the last exon.

Human: Mouse:

Sort by: [Method for Sorting This Table](#)

Show only hits where:

family is at least number and score is at least (zero ignores these cutoffs).

HCP34318.1 Sequence Info Known: IL-1 family GeneBase (IMX189)	Family is number 1 of 15 possible score 999.9, length 277	contig: GA_X54KRE9YM0J chrom: 2 begin: 108313232 end: 109165726 Protein begin: 350388 end: 340348
HCP34322.1 Sequence Info Known: IL-1 family GeneBase (IMX115)	Family is number 1 of 15 possible score 999.9, length 298	contig: GA_X54KRE9YM0J chrom: 2 begin: 108313232 end: 109165726 Protein begin: 402705 end: 395685
HCP1628454.1 Sequence Info Unknown GeneBase (IMX181783)	Family is number 1 of 15 possible score 163.1, length 251	contig: GA_X54KREBBWRK chrom: 8 begin: 121558207 end: 136784473 Protein begin: 7220385 end: 7250443

Drill in on possible hits

<p>Possible GeneFold Families for ID: HCP1628454.1 Description: /len=251 /protein_uid=101000087703514 /ga_name=GA_x54KREBBWRK /ga_uid=181000067218227 /transcript_na Length: 251 Sequence Info Unknown GeneBase (IMX181783)</p>	
<p>Run GeneFold</p>	<p>FASTA Sequence:</p> <pre>>HCP1628454.1-PC1 /len=251 /protein_uid=101000087703514 /ga_name=GA_x54KREBBWRK /g MLIKINQOKVETEWNLTEHITQSAYVRVEEMDTVSLPKAKENKNGIECYAPELSKFTFCIHPSIKGPKLLMYAIV AGVFQDKNTLIFQKC SKMGTSARANSREI PSTFLWLKKS VHRRLLHLAVVDSYWCMSAGPRSGI SAGGKHLPLVP SGAASYLRSVTLI WDL SGI LERVHHPKAEAAWGPKYLHAQPAAGTPCCQEAQGLAHPD PPLAPKYDAQGLKQ KAGPLPTPLVPEDHPSGVLFP RKDSP</pre>
<p>Blast vs. GeneBase</p>	
<p>GeneFold Hits:</p>	
<p><u>lita_</u> CYTOKINE Journal Title: STRUCTURE AND FUNCTION OF INTERLEUKIN-1, BASED ON CRYSTALLOGRAPHIC AND MODELING STUDIES score: 163.1, score type: br, hit name: lita_</p>	
<p><u>lvolA</u> COMPLEX (TRANSCRIPTION FACTOR/REGN/DNA) Journal Title: CRYSTAL STRUCTURE OF A TFIIB-TBP-TATA-ELEMENT TERNARY COMPLEX score: 134.1, score type: br, hit name: lvolA</p>	
<p><u>4rhv3</u> RHINOVIRUS COAT PROTEIN Journal Title: THE USE OF MOLECULAR-*REPLACEMENT PHASES FOR THE REFINEMENT OF THE HUMAN RHINOVIRUS 14 STRUCTURE score: 132.8, score type: sq, hit name: 4rhv3</p>	
<p><u>ludiE</u> COMPLEX (HYDROLASE/INHIBITOR) Journal Title: NUCLEOTIDE MIMICRY IN THE CRYSTAL STRUCTURE OF THE URACIL-DNA GLYCOSYLASE - URACIL GLYCOSYLASE INHIBITOR PROTEIN COMPLEX score: 130.9, score type: br, hit name: ludiE</p>	
<p><u>lita_</u> CYTOKINE Journal Title: STRUCTURE AND FUNCTION OF INTERLEUKIN-1, BASED ON CRYSTALLOGRAPHIC AND MODELING STUDIES score: 129.6, score type: sq, hit name: lita_</p>	

Mining the Genome

Verify by viewing full GeneFold run

GeneFold: Display Run ketchemr_20020520103949.txt
 (Instructions) (Return to the [main GeneFold page](#)) (Display Printable Version)

seq. (sq)	seq. +loc. +bur. (br)	seq. +loc. +bur. +ss (tt)
* 1. 4rhv3 P F 132.8 4rhv3 RHINOVIRUS C	* 1. lita P F 163.1 lita CYTOKINE (ST	* 1. ludiE P F 39.7 ludiE COMPLEX (HYDR
* 2. lita P F 129.6 lita CYTOKINE (ST	* 2. lvolA P F 134.1 lvolA COMPLEX(TRAN	* 2. lita P F 33.8 lita CYTOKINE (STRU
* 3. ludiE P F 111.3 ludiE COMPLEX (HYD	* 3. ludiE P F 130.9 ludiE COMPLEX (HYD	* 3. lntyH P F 25.9 lntyH MONOOXYGENASE
* 4. lurk P F 109.5 lurk PLASMINOGEN	* 4. 4eng P F 77.2 4eng GLYCOSYL HYD	* 4. lurk P F 22.9 lurk PLASMINOGEN AC
* 5. 9pap P F 87.5 9pap HYDROLASE (S	* 5. 2cab P F 67.3 2cab HYDRO-LYASE	* 5. lgnbJ P F 18.8 lgnbJ ACUTE-PHASE PR
* 6. 2cab P F 76.3 2cab HYDRO-LYASE	* 6. 1thw P F 61.0 1thw SWEET TASTIN	* 6. 4rhv3 P F 17.6 4rhv3 RHINOVIRUS COA
* 7. lntyH P F 69.4 lntyH MONOOXYGENAS	* 7. lgnbJ P F 55.5 lgnbJ ACUTE-PHASE	* 7. 1c52 P F 15.6 1c52 ELECTRON TRANS
* 8. lvolA P F 62.0 lvolA COMPLEX(TRAN	* 8. 7timB P F 49.2 7timB INTRAMOLECUL	* 8. 1fxrB P F 14.8 1fxrB ELECTRON TRANS
* 9. 2hlpH P F 60.5 2hlpH COMPLEX (ANT	* 9. 1svcP P F 46.1 1svcP COMPLEX (TRA	* 9. 1kst P F 14.4 1kst AGGREGATION IM
* 10. 1ppo P F 59.6 1ppo HYDROLASE(TH	* 10. 1mri P F 44.2 1mri RIBOSOME-TNA	* 10. 1sebF P F 14.1 1sebF COMPLEX (MHC I
* 11. lgnbJ P F 57.7 lgnbJ ACUTE-PHASE	* 11. 1sebF P F 42.2 1sebF COMPLEX (MHC	* 11. 7timB P F 13.5 7timB INTRAMOLECULAR
* 12. 1sebF P F 55.4 1sebF COMPLEX (MHC	* 12. 1jud P F 42.0 1jud DEHALOGENASE	* 12. 1iceA P F 12.7 1iceA CYTOKINE (STRU
* 13. 6fabH P F 53.3 6fabH IMMUNOGLOBUL	* 13. 1poiD P F 38.7 1poiD TRANSFERASE	* 13. lvolA P F 12.4 lvolA COMPLEX(TRANSC
* 14. 3drcB P F 52.4 3drcB OXIDOREDUCTA	* 14. 4jdwA P F 38.6 4jdwA TRANSFERASE	* 14. 1melB P F 12.3 1melB COMPLEX (ANTIB
* 15. 6fabL P F 52.0 6fabL IMMUNOGLOBUL	* 15. 1aw9 P F 35.4 1aw9 TRANSFERASE	* 15. 3dfr P F 12.1 3dfr OXIDO-REDUCTAS
* 16. 3sdpB P F 44.7 3sdpB OXIDOREDUCTA	* 16. 2fvwL P F 35.2 2fvwL IMMUNOGLOBUL	* 16. 1jvr P F 12.0 1jvr MATRIX PROTEIN
* 17. 1an3C P F 40.6 1an3C COMPLEX (HORM	* 17. 8fabD P F 34.9 8fabD IMMUNOGLOBUL	* 17. 1vcqB P F 11.7 1vcqB NUCLEOCAPSID P
* 18. 1mpaH P F 36.6 1mpaH COMPLEX (IMM	* 18. 1arc P F 34.7 1arc HYDROLASE (SE	* 18. 1thw P F 11.6 1thw SWEET TASTING
* 19. 1cv8 P F 34.9 1cv8 CYSTEINE PRO	* 19. 1ecrA P F 34.3 1ecrA COMPLEX (DNA	* 19. 1b5m P F 11.6 1b5m ELECTRON TRANS
* 20. 1jud P F 34.1 1jud DEHALOGENASE	* 20. 1p38 P F 34.2 1p38 TRANSFERASE	* 20. 6fabH P F 11.6 6fabH IMMUNOGLOBULIN
* 21. 1hnf P F 33.8 1hnf T LYMPHOCYTE	* 21. 2bbvB P F 34.1 2bbvB COMPLEX (VIRU	* 21. 1nfa P F 11.5 1nfa TRANSCRIPTION
* 22. 4eng P F 33.3 4eng GLYCOSYL HYD	* 22. 1eaf P F 33.8 1eaf DIHYDROLIPOA	* 22. 2cab P F 11.5 2cab HYDRO-LYASE
* 23. 1wdcB P F 31.5 1wdcB MUSCLE PROTEI	* 23. 1vcqB P F 32.9 1vcqB NUCLEOCAPSID	* 23. 1cd8 P F 11.1 1cd8 SURFACE GLYCOP
* 24. 1cnv P F 31.3 1cnv SEED PROTEIN	* 24. 1mhcD P F 32.4 1mhcD HISTOCOMPATI	* 24. 1npoC P F 11.1 1npoC COMPLEX (HORMO
* 25. 1gc1H P F 31.0 1gc1H COMPLEX (HIV	* 25. 1ako P F 32.4 1ako NUCLEASE	* 25. 1yaiC P F 11.0 1yaiC OXIDOREDUCTASE
* 26. 3hfmH P F 27.3 3hfmH COMPLEX (ANTI	* 26. 2hvm P F 31.4 2hvm HYDROLASE (T	* 26. 1tiv P F 10.9 1tiv TRANSCRIPTION
* 27. 1lx1 P F 26.1 1lx1 APOPTOSIS (X	* 27. 2acu P F 31.3 2acu OXIDOREDUCTA	* 27. 1hunB P F 10.5 1hunB CYTOKINE (CHEMO
* 28. 8fabD P F 25.8 8fabD IMMUNOGLOBUL	* 28. 1nbaD P F 30.8 1nbaD HYDROLASE (IN	* 28. 1ar1D P F 10.3 1ar1D COMPLEX (OXIDO
* 29. 1tam P F 25.6 1tam MATRIX PROTEI	* 29. 1bbmB P F 29.0 1bbmB COMPLEX (EMU	* 29. 4vgcC P F 10.3 4vgcC SERINE PROTEAS
* 30. 1vcqB P F 23.9 1vcqB NUCLEOCAPSID	* 30. 6fabL P F 28.7 6fabL IMMUNOGLOBUL	* 30. 1rtoB P F 10.1 1rtoB CHEMOKINE (PRO
* 31. 1innB P F 23.0 1innB COMPLEX (HOM	* 31. 1frfS P F 28.7 1frfS NI-FE HYDROG	* 31. 6rhn P F 10.1 6rhn NUCLEOTIDE-RTM

Alignment of [lita_](#), method br, run gf
 50.6% of the query is aligned, 84.1% of the template is aligned

Primary sequence:
 query: **ELIKIN**QK**YET**EWNLT**EH**I**EQSAYV**RV**EE**MDTV**SL**PKAKEN**KNG**IECY**AP**ELSK-----**FTFC**I**HP**S-----**IKGPK****EL****IV****AT****AG**V**FQ**DK**MT****LI****FQ**CK**SK**MG**T**SAR**AN**SRE**IP****ST****FL****WL**
 lita_ :-----**NV**K**Y****N****F****M****R**-----**LI****K****Y****E****F****I****L****N****D****A****L****N****Q**-----**SI****I****R****A****N****D****Q****Y**-----**LT****A****A****A****L****N****H****L****D****E****A****V****K****F****D****M****G****A****Y****K****S****K****D****A****K****I****T****V****I****L****R****I****S****K****T****Q****L****Y****V****T****A****Q**-----**DE****D****Q****P****V****L****L****K****E****M****P****E****I****P****K****T****T****G****S****E****T****N****L**-----**LF****F****W****E**

Secondary structure:
 query: **EEEE** **E** ----- **RRRRRR** ----- **EEEE** ----- **EEEE** ----- **RRRRRR** ----- **EE** ----- **RRRRRR**
 lita_ :-----**EEEE** **EEEEEEEE** ----- **EEEE** ----- **EEEE** ----- **TT** ----- **EEEEEE** ----- **EEEEEE** ----- **EEEE** ----- **EEEE** ----- **EE** ----- **GGGGGG** ----- **EEEE**

Some Results

Gene mining by remote homology
detection has been very successful

Identified several novel human cytokines

Verification lead to discovery of further
novel cytokines

Some Problems

MANY false positives - many hits to wade through

Requires expert in particular family to identify true positives

Genes must be verified experimentally

Some folds hard to score - TNFR's

Not all folds represented

Utilize advances in remote homology detection

As structure representatives grow, so will ability of remote homology detection

Utilize fast, automated methods for assigning structure family

Fast Threading Data Analysis

To mine threading data, we need programs that:

Repeat interpretation of threading output consistently and quickly

We can train to recognize different folds in the output

Aid protein structure experts by applying similar logic

Fast Threading Data Analysis

We chose:

The support vector machine algorithm

Quick to train, quick to give answers

Generates score which can be used as
measure of confidence in answer
generated

How SVM's work

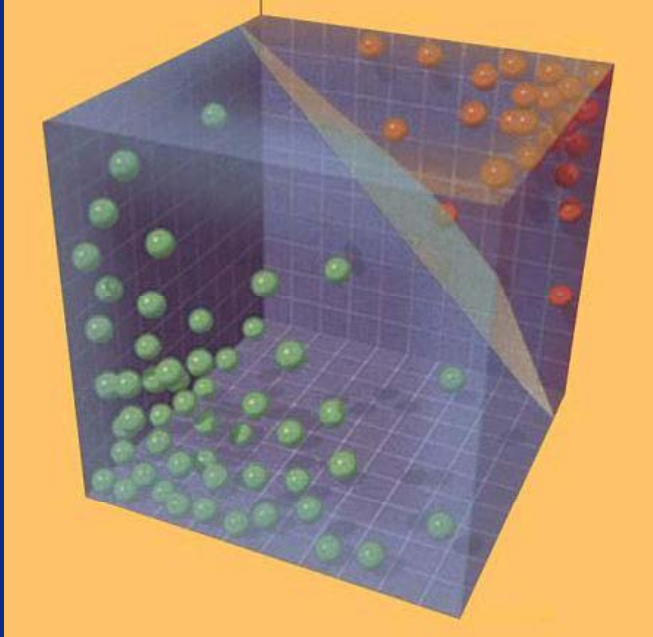


Figure from cover of book: Introduction to Support Vector Machines. Christianni and Shaw-Taylor

SVM work by Paul McDonagh,
Immunex Corporation

SVM takes 'positive' and 'negative' fold examples

Red = positive, Green = negative

Uses function (kernel function) to plot data to different type of space

Graphic shows 3-D linear kernel space

Threading uses 1823-D spherical space

Regression techniques fit a plane

Vectors from points 'support' the plane

Term coined - support vector machine

If fall on red side of plane - new member of the fold

Distance from plane gives measure of confidence in prediction

Support Vector Machines

Trained by scientists

Success primarily depends on scientific input to training set

Scientist finds members of fold - positive training set

Scientist identifies all other folds - negative training set

Support Vector Machines

Threading algorithm run on unknowns

1823*3 data points for each protein in a set

Support vector machines find which of the 1823*3 points and values carry the most predictive power

Early results have been very promising

Genome sequenced, but still a LONG way to go for function

Structure homology methods valuable in identifying unknown sequences

Many structure families not represented

Need better remote homology detection methods

Need fast, automated methods