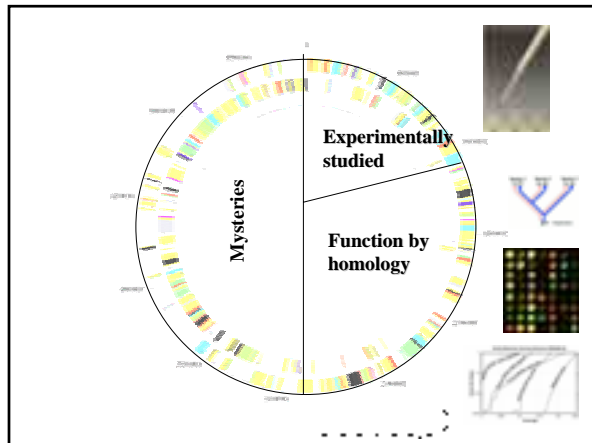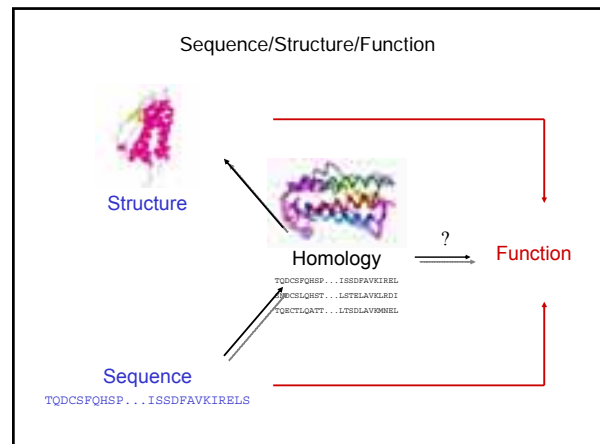# Ancient Protein
# Evolutionary Relationships
# Inferred from Structure

Emma E. Hill &
Steven E. Brenner

University of California, Berkeley

IPAM - Structural Proteomics
10 May 2004

---

Sequence/Structure/Function



Structure

Homology    ? → Function

TQDCSFQHSP...ISSDFAVKIREL
TQDCSLQHST...LSTELAVKLRDI
TQECTLQATT...LTSDLAVKMNEL

Sequence
TQDCSFQHSP...ISSDFAVKIRELS

---



Experimentally studied

Mysteries

Function by homology

---

# Genome Annotation Quality

- What is the quality of genome annotation?
- Quality of sequence well known
- Quality of gene prediction at least roughly understood
- Functional accuracy of 99.5% claimed…
      … but not tested experimentally
- *We rely upon functional assignments for biological interpretation*

---

# The Annotation of *M. genitalium*

1. TIGR sequences genome and makes initial annotation
2. GeneQuiz consortium automatically annotates
3. Eugene Koonin et al (NCBI) manually make annotations
4. GeneQuiz consortium automatically re-annotates
5. Updates
   - Several groups make automated structural annotations
   - TIGR makes updates to annotation, including new genefinding

Different groups use similar methods and operated sequentially, reviewing each others' results

---

# Compatible Annotations

**mg463**
**TIGR:** ● high level kasgamycin resistance (ksgA)
**NCBI:** ● rRNA (adenosine-N6, N6-)-dimethyltransferase (ksgA)
**GeneQuiz:** ● Dimethyladenosine transfe [sic]

**mg010**
**TIGR:** ● DNA primase (dnaE)
**NCBI:** ● DNA primase (truncated version) (DnaGp)
**GeneQuiz:** ● DNA primase (EC 2.7.7.-)

**mg225**
**TIGR:** ● hypothetical protein
**NCBI:** ● amino acid permease
**GeneQuiz:** ● histidine permease

## Slide 1: Incompatible Annotations

**mg302**
| | | |
|---|---|---|
| **TIGR:** | ● | **no database match** |
| **NCBI:** | ● | **(glycerol-3-phosphate?) permease** |
| **GeneQuiz:** | ● | **mitochondrial 60S ribosomal protein L2** |

**mg448**
| | | |
|---|---|---|
| **TIGR:** | ● | **pilin repressor (pilB)** |
| **NCBI:** | ● | **putative chaperone-like protein** |
| **GeneQuiz:** | ● | **pilB protein** |

**mg085**
| | | |
|---|---|---|
| **TIGR:** | ● | **hydroxymethylglutaryl-CoA reductase (NADPH)** |
| **NCBI:** | ● | **ATP(GTP?)-utilizing enzyme** |
| **GeneQuiz:** | ● | **NADH-ubiquinone oxidoredu [sic]** |

## Slide 2



## Slide 3: Genome Annotation Quality

- **Average error rate at least 8%**
  - Actual error rate likely to be 2-3 times higher

- **Where do errors come from?**
  - Poor sequence comparison: not homology at all
  - Incorrect inferences of function from homology
  - Propagation of erroneous data

- **Solutions?**
  - Careful sequence comparison
  - Avoidance of over-annotation
  - Complete description of method in database
  - New methods for functional characterization
    - ***…Phylogenomics***

Legend:
- 2 OK
- 2 Wrong
- 3 OK
- 3 Wrong
- 0 or 1

## Slide 4: Inferring Homology

```
DR1776      EVPAELPHGAFSVLDNTDTGFEWVRLDELGARPVYPLLVRDLLSVPVGEVRHLVIRS--
DR2272      LT-GELPA---TVLDNPHVFFRWLAVDALDDHTLYPRCVPQLLRLPAGEIGHFVTDERA
            . .***   :****... *.*: :* *. :.:**  * :** :*.**: *:*  .
```

**Percentage Identity scale**

| 0 | 25 | 30 | 50 | 100 |
|---|---|---|---|---|

**0-25%**
- Homology possible but often unclear
- Detection and Alignment difficult

**25-30%**
- Highly likely homology
- Detection and Alignment becoming tricky

**30-50%**
- Good homology
- Relatively easy detection and alignment

**>50%**
- Close homology
- Trivial detection and alignment

## Slide 5

ATGTTGCAT

→ Transcription

AUGUUGCAU

→ Translation

MLH

## Slide 6

2

/footer_navigation

## Moore's Law



**Abstract:**
"With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65,000 components on a single silicon chip"

Actually: 5,000 in 1974
Intel 8080

**Further:**
by 2003, there would be $2^{32}$ (~17 billion; ~$10^{10}$) components on each integrated circuit
Actually: $4 \times 10^8$ ($2^{28.6}$)
Intel Madison

MOORE, G. E. *Cramming More Components Onto Integrated Circuits.* **Electronics** 38, 8 (April 19 1965), 114--117.

## Moore's Law



Number of Transistors (vertical axis: $10^2$ to $10^{12}$)

Transistors/Processor (Intel)

Horizontal axis: 1972 1976 1980 1984 1988 1992 1996 2000

## DNA Sequence Data Growth



Number of Monomers or Transistors (vertical axis: $10^2$ to $10^{12}$)

GenBank

1982

Transistors/Processor (Intel)

Horizontal axis: 1972 1976 1980 1984 1988 1992 1996 2000



$n = \exp(0.19\, y)$

n = new structures/year
y = years since 1960

**Prediction:**
At the end of 2001, there would be **13,941** crystal structure entries available in the PDB

**14,000 crystal structures were actually available!**

Arthur Arnone
http://www.rcsb.org/pdb/pdb_news2002.html

## Database growth

Number of Monomers or Transistors

$10^{12}$
$10^{10}$
$10^8$
$10^6$
$10^4$
$10^2$

GenBank
SP+TrEMBL
1996
SwissProt
1982   1986
Protein Data Bank
Transistors/Processor
Open-access articles

1972  1976  1980  1984  1988  1992  1996  2000

## Structural Genomics

- Provide a 3D structure or quality model for every tractable biomacromolecule
- LBNL Berkeley Structural Genomics Center is one of 9 NIH-funded pilot centers
- Comparable worldwide efforts, esp Japan
- Experimental & computational effort

## Practical Target Selection

Proteins in Sequence Space

## Practical Target Selection

Proteins in Sequence Space

## Practical Target Selection

Proteins in Sequence Space

## Practical Target Selection

Proteins in Sequence Space

## Practical Target Selection



Proteins in Sequence Space

## Practical Target Selection



Proteins in Sequence Space

## Protein Selection

- All proteins in *Mycoplasma* are potential targets
- Potentially easy to characterize members
  - Cloning (accessible organisms)
  - Expression (UGA/Trp)
  - Purification (thermo) & Crystallization

## BSGC Target Selection (Current)



Known 3D Structure (BLAST, PSI-BLAST, Pfam)

Too Long (>700 AA)

Transmembrane (TMHMM, PHDhtm)

Coiled Coil ≥ 20% (CCP)

Low Complexity ≥ 20% (SEG)

Too Many Nonstandard Codons

Select up to 10 homologues per *MP* gene, prioritize thermophiles

## Practical Structural Genomics



Proteins in Sequence Space

## Why Classify

"Physics or Stamp Collecting"
--Ernest Rutherford

General trends provide insight into underlying principles

Unusual features only become apparent with knowledge of the principles

*The Jenny*

Assist predictions

## Slide 1: SCOP Sample Hierarchy



SCOP Sample Hierarchy

| | | | | | |
|---|---|---|---|---|---|
| | scop | | | **Root** | |
| α | β | α/β | α+β | **Class** | Determined by structure |
| Rossmann fold | Flavodoxin-like | α/β-Barrel | | **Fold** | |
| TIM | Trp biosynthesis | Glycosyltransferase | RuBisCo (C) | **Superfamily** | |
| β-Galactosidase (3) | β-Glucanase | α-Amylase (N) | β-Amylase | **Family** | Related by homology |
| Acid α-amylase | Cyclodextrin glycosyltransferase | Oligo-1,6 glucosidase | | **Protein** | |
| A. niger | B. circulans | B. stearothermophilus | B. cereus | **Species** | |
| 2aaa:1-353 | 1cdg:1-382 1cgt:1-382 | 1cyg:1-378 | J. Biochem 113:646-649 | **PDB/Ref** | |

## Slide 2 (image)



## Slide 3: SCOP Sample Hierarchy

SCOP Sample Hierarchy

| | | | | | |
|---|---|---|---|---|---|
| | scop | | | **Root** | |
| α | β | α/β | α+β | **Class** | Determined by structure |
| Rossmann fold | Flavodoxin-like | α/β-Barrel | | **Fold** | |
| TIM | Trp biosynthesis | Glycosyltransferase | RuBisCo (C) | **Superfamily** | |
| β-Galactosidase (3) | β-Glucanase | α-Amylase (N) | β-Amylase | **Family** | Related by homology |
| Acid α-amylase | Cyclodextrin glycosyltransferase | Oligo-1,6 glucosidase | | **Protein** | |
| A. niger | B. circulans | B. stearothermophilus | B. cereus | **Species** | |
| 2aaa:1-353 | 1cdg:1-382 1cgt:1-382 | 1cyg:1-378 | J. Biochem 113:646-649 | **PDB/Ref** | |

## Slide 4: Making Structure Classification Consistent and Automated

Making Structure Classification Consistent and Automated

1. Automatically determine homology from these features for proteins of known structure
2. Calculate phylogenies for protein superfamilies
3. Apply phylogenomic techniques to predict protein function

Identifying & learning on protein structure features



START → New structure → Structural similarity? → No STOP / Yes → **1** Homology? → No STOP / Yes → **2** Phylogeny → **3** Function

CAPER (Classification of Ancient Protein Evolution)

## Slide 5: The Plan

The Plan



Data Structure Features → **1. Kernel function** K(x,y) → Kernel matrices → **2. Learning Algorithms** LA → Homology prediction No/Yes

Individual mini-kernels and joining of multiple kernels → Feature selection

Phylogeny reconstruction

**Representation in CAPER**

## Slide 6: Validation on the SCOP database

Validation on the SCOP database

| | | | | | |
|---|---|---|---|---|---|
| | scop | | | **Root** | |
| α | β | α/β | α+β | **Class** | Determined by structure |
| Rossmann fold | Flavodoxin-like | α/β-Barrel | | **Fold** | |
| TIM | Trp biosynthesis | Glycosyltransferase | RuBisCo (C) | **Superfamily** | |
| β-Galactosidase (3) | β-Glucanase | α-Amylase (N) | β-Amylase | **Family** | Related by homology |
| Acid α-amylase | Cyclodextrin glycosyltransferase | Oligo-1,6 glucosidase | | **Protein** | |
| A. niger | B. circulans | B. stearothermophilus | B. cereus | **Species** | |
| 2aaa:1-353 | 1cdg:1-382 1cgt:1-382 | 1cyg:1-378 | J. Biochem 113:646-649 | **PDB/Ref** | |

## Four-helical cytokines



http://scop.berkeley.edu

---

## Low Sequence Identity but Definite Homology

Long Chain
Short Chain

| seed | Erythropoietin d1eera_ 2.1 | GM-CSF d2gmfa_ 2.2 | IL-4 d1iara_ 2.3 | IL-5 d1hula_ 2.4 | M-CSF d1hmca_ 2.5 | Flt3 d1etea_ 2.6 | IL-2 d3inkc_ 2.7 | IL-3 d1jli__ 2.8 |
|---|---|---|---|---|---|---|---|---|
| d1eera_ 2.1 | 2.80E-71 166 | None | None | None | None | None | None | None |
| d2gmfa_ 2.2 | None | 7.90E-65 121 | None | None | None | None | None | None |
| d1iara_ 2.3 | None | None | 2.20E-54 129 | None | None | None | None | None |
| d1hula_ 2.4 | None | None | None | 1.00E-51 108 | None | None | None | None |
| d1hmca_ 2.5 | None | None | None | None | 7.40E-82 145 | None | None | None |
| d1etea_ 2.6 | None | None | None | None | None | 1.40E-80 134 | None | None |
| d3inkc_ 2.7 | None | None | None | None | None | None | 7.60E-64 128 | None |
| d1jli__ 2.8 | None | None | None | None | None | None | None | 2.50E-43 111 |

• Long Chain Family   - 3 inter-protein hits.

• Short Chain Family  - 0 inter-protein matches.

---

## Why are we certain they are homologous?



Unique topology (up-up-down-down)

Restricted range of functions

Homologous receptors

Conserved exon/intron boundaries

---

**Family-level Phylogeny of the four-helical cytokines.**



| | Interferons/IL-10 | Long chain | Short chain |
|---|---|---|---|
| **Protein (PDBID):** | Ifn-β (1au1) | LIF (1lki) | IL-4 (1rcb) |
| **Reference:** | Karpusas et al., 1997 | Robinson et al., 1994 | Wlodawer et al., 1992 |
| **Av. Chain Length:** | 180 | 180 | 140 |
| **Av. Helix Length:** | 20-30 | 20-30 | 10-20 |
| **AB Helix Linker:** (w.r.t Helix D) | Outside | Outside | Inside |
| **Other:** | CD linker forms 5th helix | Small Helix in AB linker | Small 2 stranded β-sheet |

---



---

## The Globin-like superfamily



**Families:**   Globins   Nerve tissue mini-hemoglobin   Truncated hemoglobin   Phycocyanin-like

prototype

lacks first helix (but more similar to globins than truncated hemoglobins.)

lacks first helix

Oligomers of two different homologous subunits each subunit has two additional N-terminus helices binds a chromophore

## Slide 1

**Phylogenomic Application - MOP-like Superfamily**

?                    ?

**Given a new structure with no known sequence similar proteins & unknown function, what can we do???**

**Structural similarity?**

**Homology?** MOP(1gut)     BiMOP (1h9r)     MalK(1g29)

**Phylogenetic relationship?**

**Function?**

duplication          duplication

fusion with DNA binding domain

## Slide 2

**Structure characters for homology detection and phylogeny estimation**

CE
DALI
SSM

1. **Structural alignment**
2. **Shared physical structure features**
3. **Functional residues and interactions**
4. **Structural domain context**
5. **Species distributions**
6. **Function**

EC
GO

Superfamily
Pfam
COGs

β −bulges.
Helix capping.
Left handed β-α-β crossover.
Loop length and non functional surfaces.
Side chain conformations (rotamers).
Disulfide topology.
Secondary structure.
Hydrogen bonding patterns.
Cis-Proline conformation.
Torsion angles.
Center of mass.
Average density.
Angles of secondary structure elements with respect to one another.
Surface area.

Energetically unfavorable residues (e.g. buried charge).
Active sites.
Cofactor binding sites.
Peripheral binding sites (e.g. metal).
Residue-residue contacts.
Co-variation of interacting residues.
Functional residue clusters.
Ligands/co-factors.
Domain interface surface area.
Electrostatic potential.

Domain architecture.
Sub-domain architecture.
Interacting partners.
Oligomeric state.
Circular permutations.
Intron/exon boundaries.

## Slide 3

**Structure features examples**

**Domain organization:** who are my neighbors?

**Secondary structure:** what am I made up of?

**Size:** How big am I? (Volume, # amino acids, Mol. wt., surface area)

**Oligomeric state:** What is my biologically active form? (monomer, dimer etc.)

**Species distribution:** In which organisms am I found?

## Slide 4

**A detailed example: Disulfide Bridges**

An oxidative environment (extracellular) trigger formation of a disulfide bond between two sulfhydryl groups of cysteine residues.

The S-S bond disappears when the protein is reduced.

R-SH + R'-SH $\xrightarrow{\text{oxidation}}$ R-S-S-R'

## Slide 5

**What do disulfide bridges do?**

**Inter-chain**
Hold together different polypeptide chains (e.g. chains A & B of insulin)

**Intra-chain**
Stabilize folding of a single chain making it less susceptible to degradation (e.g. snake venom toxins & protease inhibitors)

Bovine insulin chains A & B (1cph)
Gursky *et al.*, 1993
Chains A & B are held together by three disulfide bonds

Bovine pancreatic trypsin inhibitor (1bpi),
Parkin *et al.*, 1996
When the three disulfide bonds are reduced this small protein unfolds

## Slide 6

**Structure Feature → Feature vectors → Kernels**

**Disulfide bonds in any protein structure…**

1. **Presence/absence** – how many?

2. If more than 1 – what is their **relative ordering** on the chain? (which bonds with which)

3. What is their **secondary structural context** (where are the cysteines located)?

4. What is the **secondary structure context for the entire domain**?

5. **Distances** between cysteines, & lengths of secondary structure elements.

6. Are they exactly **equivalent in position**?

7. **Distances & three-dimensional orientations** from one another &/or from centre of mass of protein?
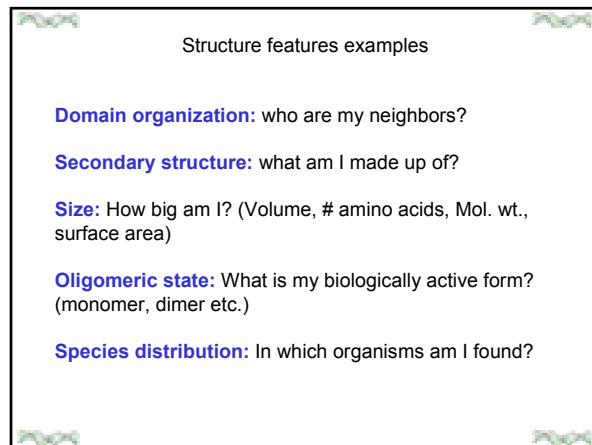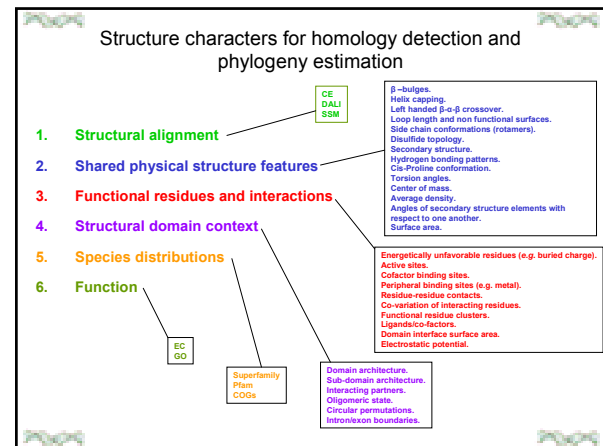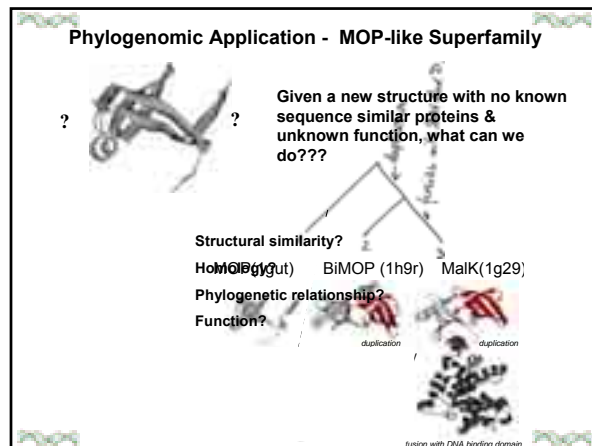
## Relative orientations (# and their ordering)

For a protein with two S-S bonds there are 3 possible orientations



| # DS | # C | # poss orders |
|---|---|---|
| 1 | 2 | 1 |
| 2 | 4 | 3 |
| 3 | 6 | 15 |
| n | 2n | $X_n = (2n-1)X_{n-1}$ |

Consideration of non-bonded Cysteines

| n | 2n + c | $X_n = [(2n-1)X_{n-1}](2n+1)^c$ |
|---|---|---|

(0) 1(0)2(0)2 (0)1(0)

## Ordering & Secondary structure context



Beta strand    Loop    Helix

**Representative alphabet for secondary structure**

E – extended (beta-strand)
H – Helix
L – Loop
…

1E    1L    2H  2H
1E    2L    1H  2H
1E    2L    2H  1H

## Ordering & More secondary structure context



Beta strand    Loop    Helix

$H_1, 1E_1, E_2, 1L_1, E_2, H_2, 22H_3, H_4$

Next step is to add…
- # residues in each secondary structure element
- actual position of C within the secondary structure element

## Structural equivalence???

If two proteins have the same signature:
$H_1, 1E_1, E_2, 1L_1, E_2, H_2, 22H_3, H_4$

Are their disulfide bridges at equivalent positions???
Do they have to be the same to be homologous?



## Mini-kernels of increasing complexity

| Feature representation | Protein A  SCOP classification | Protein B  SCOP classification | … |
|---|---|---|---|
| # DS bonds | n | n | |
| # & orientation | 11022 | 1212 | |
| #, orientation & SSE | 1H1E0L2H2H | 1H2H1E2L | |
| #, orientation & all SSE | $H_1E_2E_21L_1E_2H_222H_3$ | $H_1E_2E_21L_1E_2H_222H_3$ | |
| Above + specific lengths | $(8)H_1(6)1E_1(9)E_2(3)1L_1$ | $(8)H_1(6)1E_1(9)E_2(3)1L_1$ | |
| Alignment | Dynamic programming | Dynamic programming | |

Will we need alignment information?

## A more three-dimensional approach…
(this could remove the need to make alignments?)

Fix first Cα of first SS bond (relative to sequence) as planar, (x=0 for both Cα atoms) with x,y,z of first Cys as (0,0,0)

Give relative positions of other Cα atoms from other disulfide Cysteines

Calculate euclidean distances

Calculate RMSDs based on fitting the disulfide bonds

### Slide 1: Structural similarity → Function prediction

**BMC Structural Biology**

Crystal structure of the YfiB protein from Pseudomonas aeruginosa suggests a glutathione-dependent thiol reductase function.

**Background:** ...

**Results:** ...

**Conclusion:** ...



### Slide 2: Phylogeny → Understanding protein structure evolution



Fig. 1 from Newlove *et al.,* Structure 2004

### Slide 3: Overview of our research



START → New structure → Structural similarity? (No → STOP / Yes) → **1** Homology? (No → STOP / Yes → **2** Phylogeny) → **3** Function

**CAPER (Classification of Ancient Protein Evolution)**

1. A prediction facility to calculate homology from structure
2. Phylogenetic reconstructions for homologous proteins of known structure
3. A method to predict function based on phylogenetic location

### Slide 4: Applications

This has applications for…

**Function prediction**

**Reconstructing phylogenetic relationships**

**Relating phylogenetic lineage to protein structure evolution**
→ Understanding how changes have occurred in protein structures
→ Resolving questions about the ancestral form of proteins

Elucidating **which structure features are important** in which superfamilies
→ Relating these features back to the proteins in question

### Slide 5: Acknowledgements

**Steven E. Brenner** (UCB)
John-Marc Chandonia
Marcin Joachimiak
Donna Hendrix
Gavin Crooks
Ed Green

**Michael Jordan** (UCB)
Barbara Engelhardt
Gert Lanckriet
Neil Lawrence
Matthias Seeger
Guillaume Obozinski

**Alexey G. Murzin** (CPE Cambridge, UK)
**Cyrus Chothia** (MRC-LMB, UK)
**Michael Levitt** (Stanford)
**Tim Hubbard** (Sanger Centre, UK)
**Loredana Lo Conte** (MRC-LMB, UK)
**Veronica Morea** (La Sapienza, Rome)
**Dick Karp** (UCB)