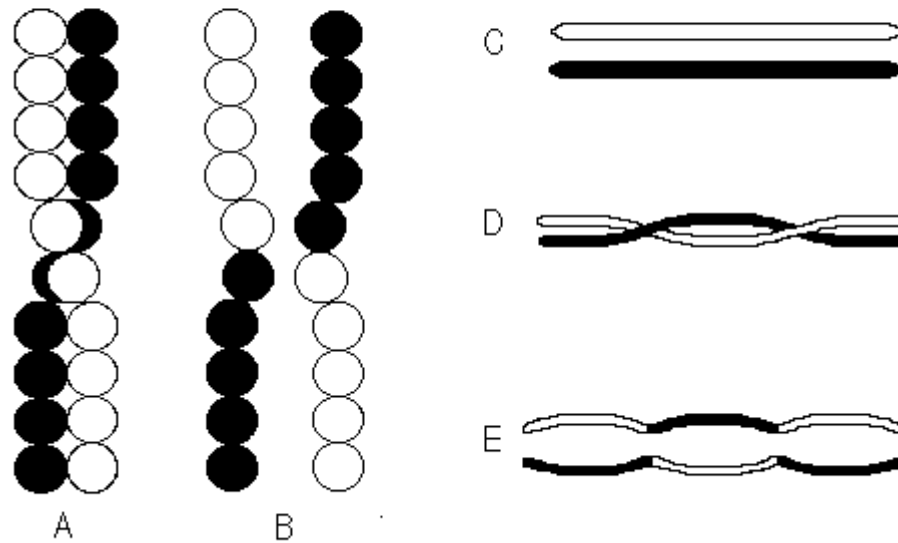


in silico protein recombination
applied to Comparative Modelling

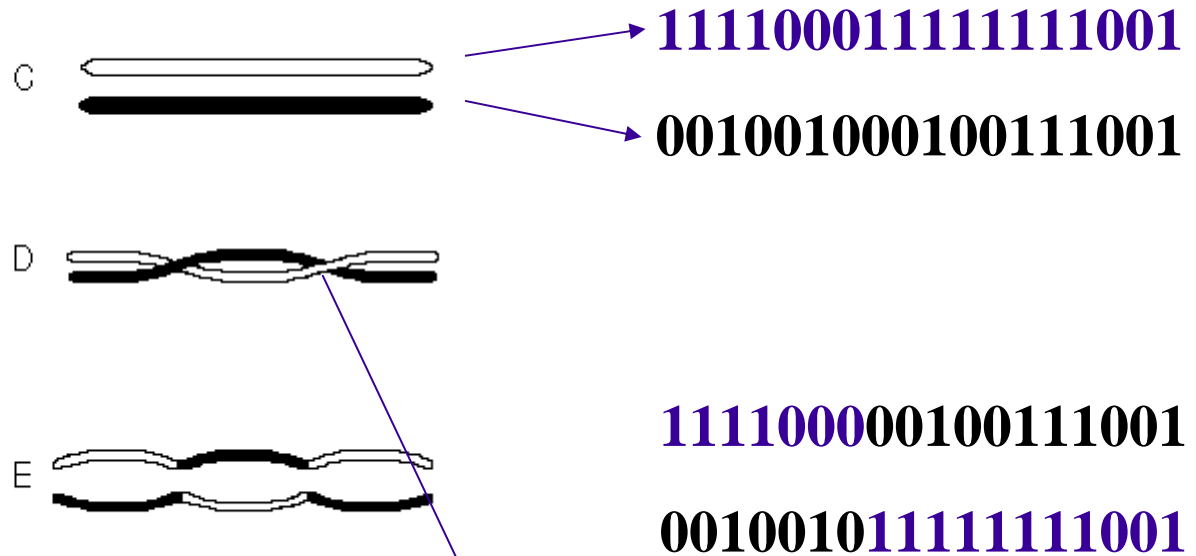
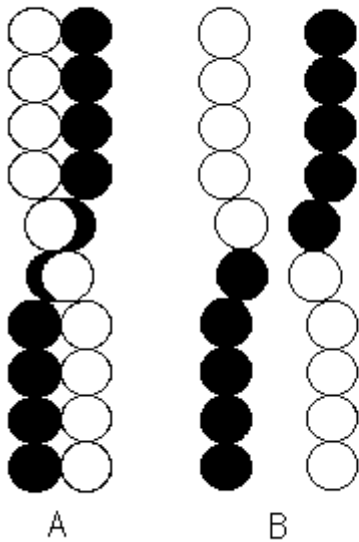
Bruno Contreras-Moreira,
Paul W. Fitzjohn and Paul A. Bates
Biomolecular Modelling Laboratory
London Research Institute
SAC-CASP5, December 2002

the biological inspiration

Chromosome Crossing-over



the computational analogy: genetic algorithms



encoded solution

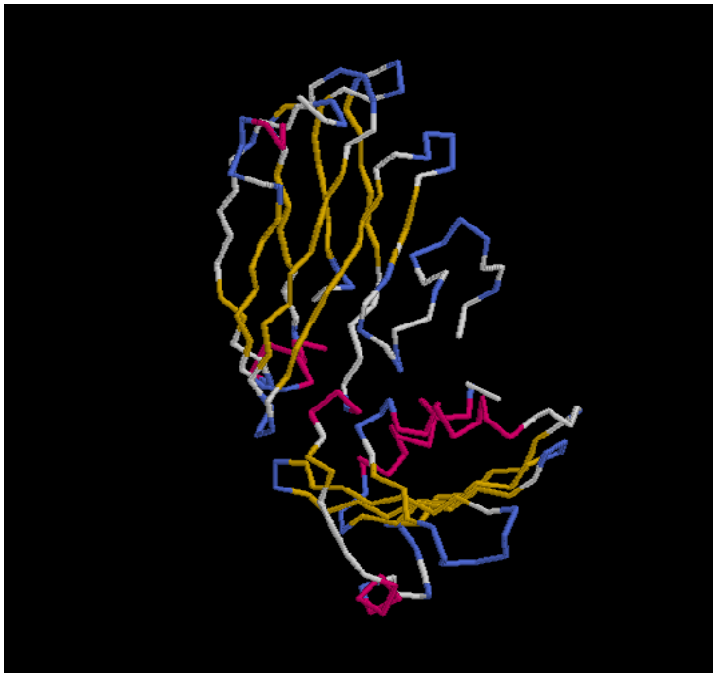
crossover point

a genetic algorithm applied to Comparative Modelling

- how are solutions coded?
 - genetic operators
 - definition of fitness
- design of the algorithm

proteins models are implicitly coded solutions

- **linear molecules:** arrays of residues connected by peptide bonds
 - **fitness** = likelihood of its fold

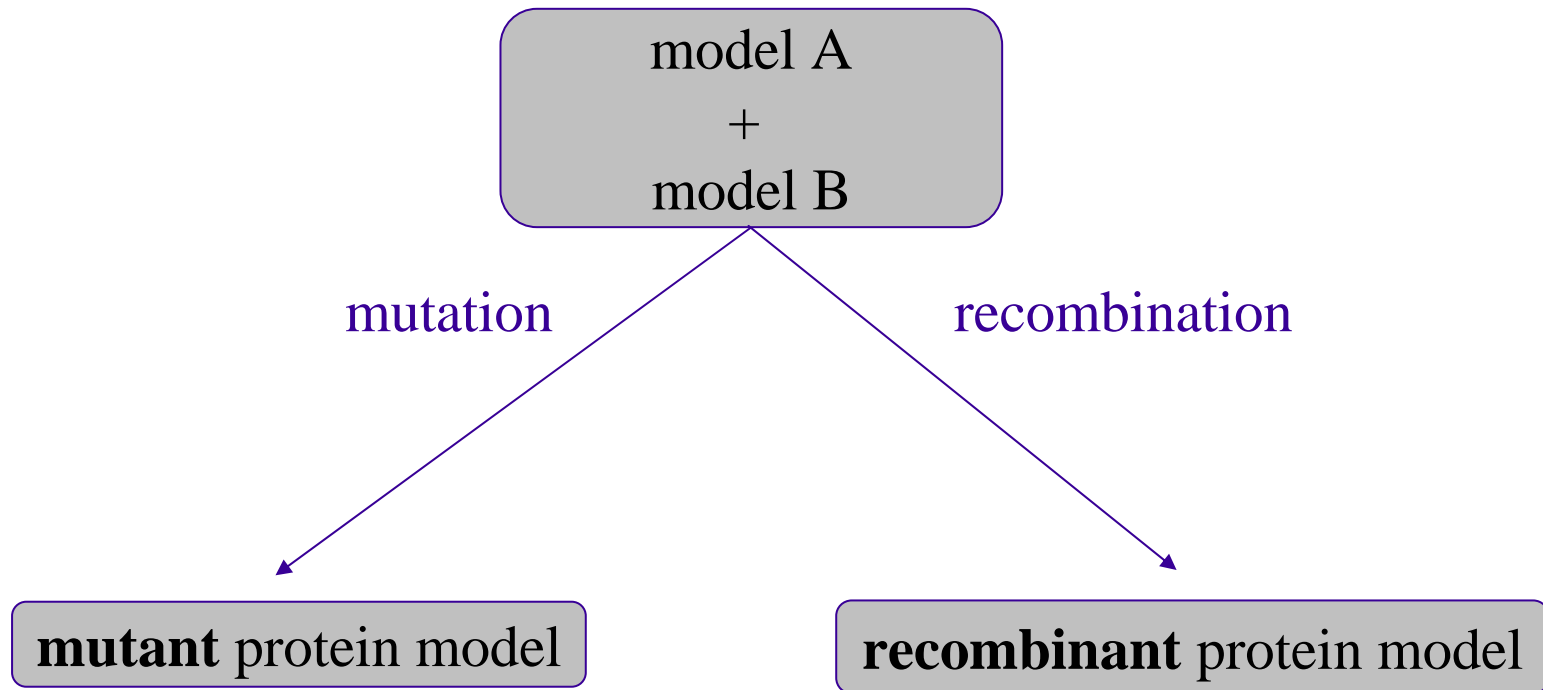


```
T0134  GEP-VQNGAPEEE--QLPPSSYSLLAENSYVMTCDIRGSLQEDSQVTVAIVLENRSS
lqt8_A  GSPGIRLGSSSEDNFARFVCKNNGVLF-ENQLIQI--GLKSEFRQNLG-RMFI FVGNKTS
SS      CCCCCCCCCCCHHHHOCCEEEE-ECCCEE--EEEEEEECCEE-EEEEEEECCE
```

(1model = 1PDB template + 1alignment)



genetic operators



recombination

model A

+

model B

- sequence alignment
 - superimpose on C β of equivalent residues
 - refine fit on close equivalent residues (2 x C α -C β)
- draw crossover point (!helix && !strand, after STICK)

mutation

model A

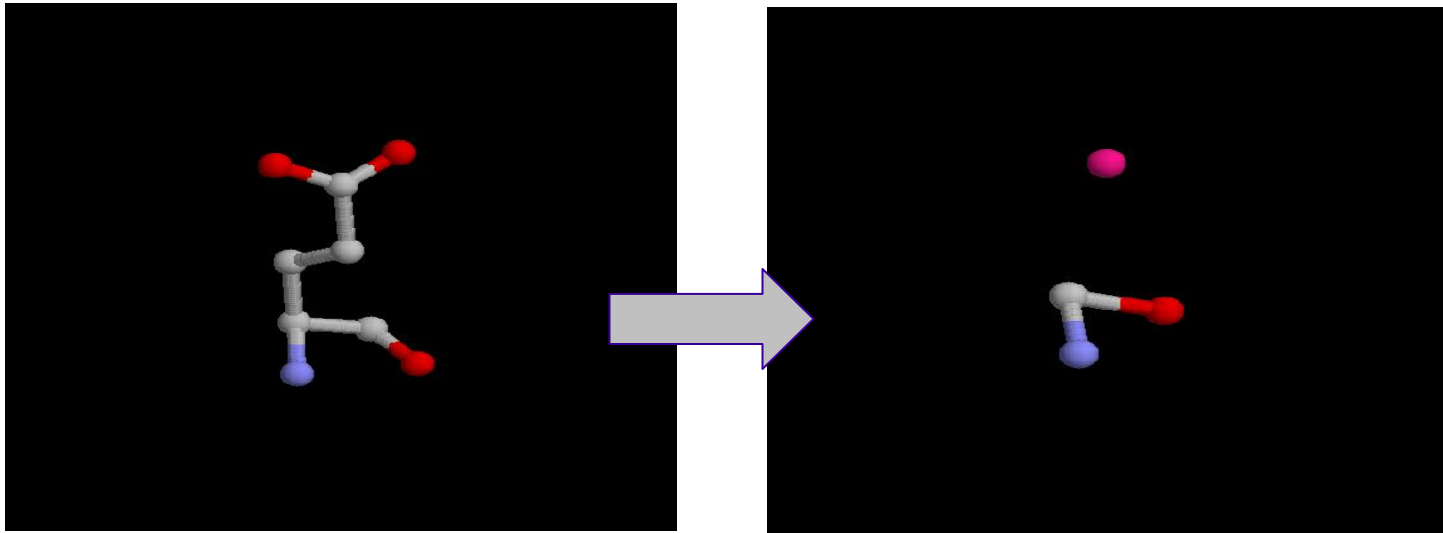
+

model B

- sequence alignment
- superimpose on C β of equivalent residues
 - all-atom Cartesian average (no checks)

protein fitness

$$\text{fitness}(p) = \text{internal_contacts}(p) + \text{solvation}(p)$$



$$\sum_i \sum_j (A_{ij}/r_{ij}^9) - (B_{ij}/r_{ij}^6) \quad (\text{in Kcal/mol})$$

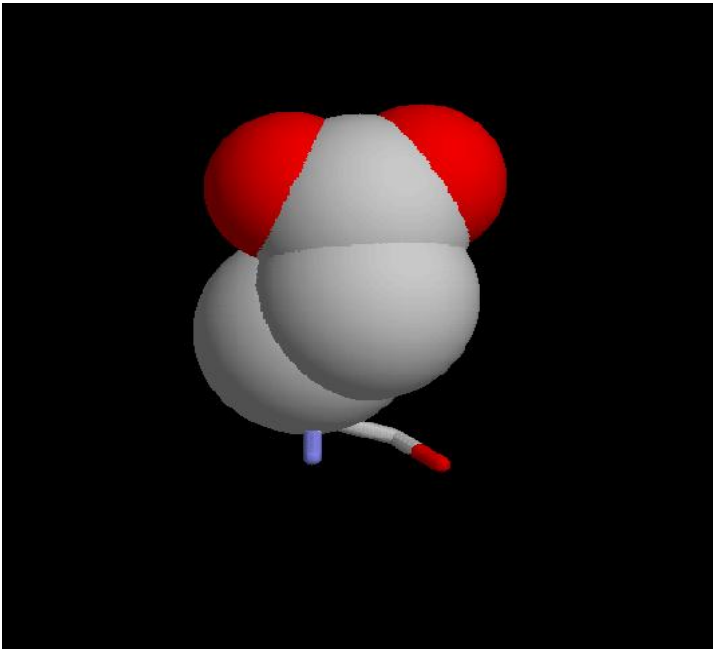
where i, j are pairs of pseudoatoms in protein p

and A and B are statistical potentials

(taken from Robson and Osguthorpe (1979) *J.Mol.Biol.* **132**(1):19-51, code by Paul Fitzjohn)

protein fitness

$$\text{fitness}(p) = \text{internal_contacts}(p) + \text{solvation}(p)$$



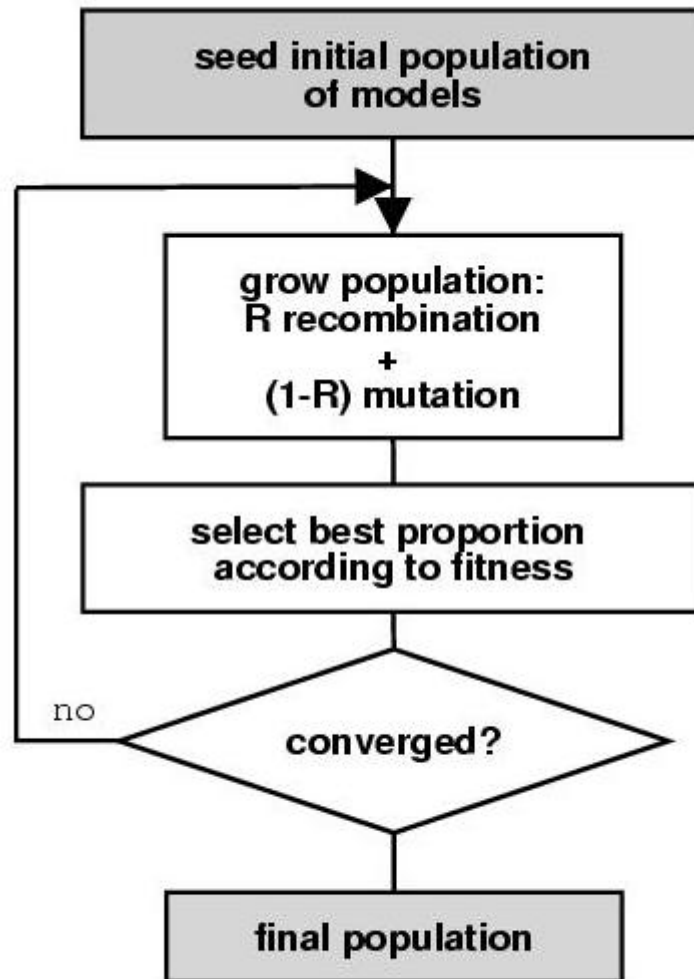
$$\sum_i (SA_i \cdot \Delta G_{\text{solv}_i}) \quad (\text{in Kcal/mol})$$

where i is a residue in protein p ,
 SA is the side-chain solvent
accessible area calculated by
NACCESS* and $\Delta G_{\text{solv}}^{\dagger}$ is the
experimental solvation free
energy change for each residue
type

* NACCESS (Hubbard and Thornton see <http://wolf.bms.umist.ac.uk/naccess>)

\dagger Eisenberg and MacLachlan (1986) *Nature*, **319**: 199-203.

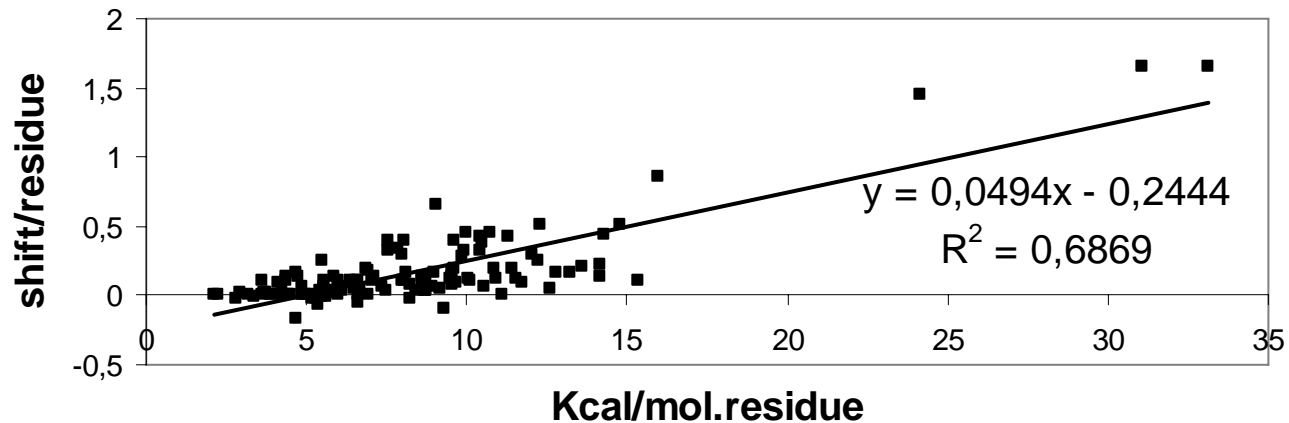
in silico protein recombination algorithm



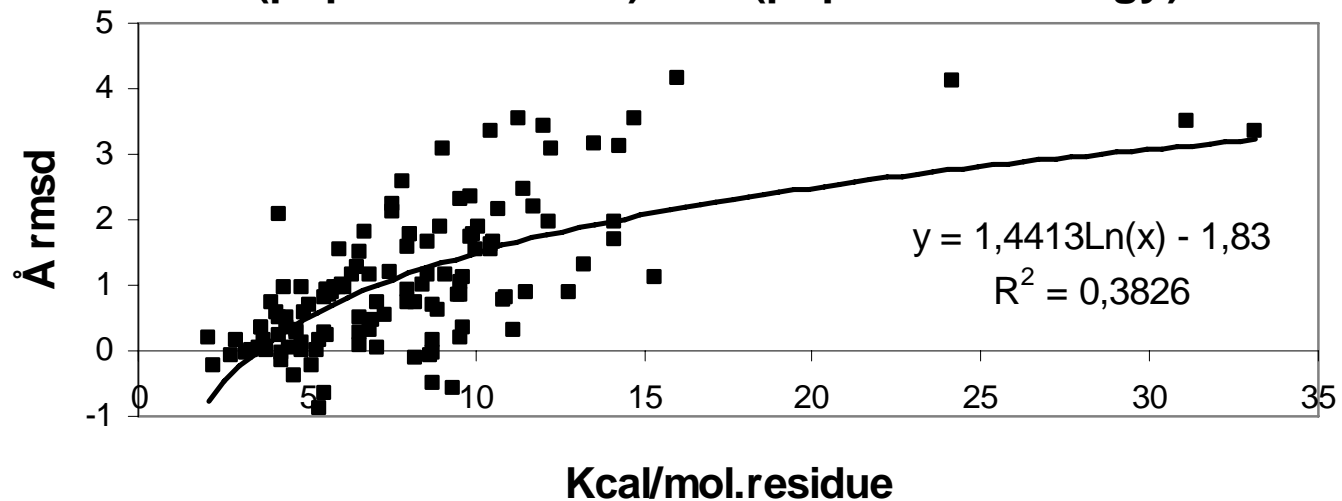
in silico protein recombination: performance

(benchmark on 130 SCOP families)

d(population energy) vs d(alignment shift)



d(population rmsd) vs d(population energy)



in silico protein recombination: evaluation

PROBLEMS

- models in the last population have sometimes **broken loops**
 - models need often to be **minimized** after the simulation
 - longer **computing time** than traditional methods
 - current **mutation** implementation does not help much

ADVANTAGES

- converges close to the best initial model
 - is able to recover alignment errors
- often last population contains different conformations

in silico protein recombination: example T0134

INITIAL POPULATION

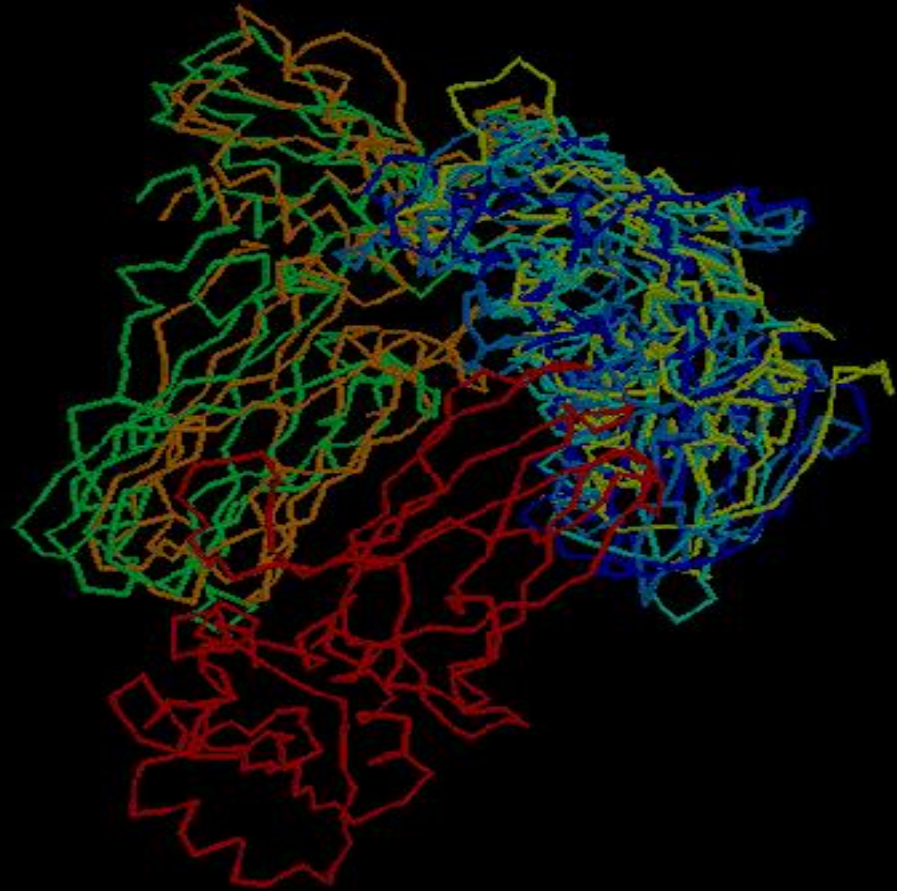
- 2 different templates

1QTS & 1E42

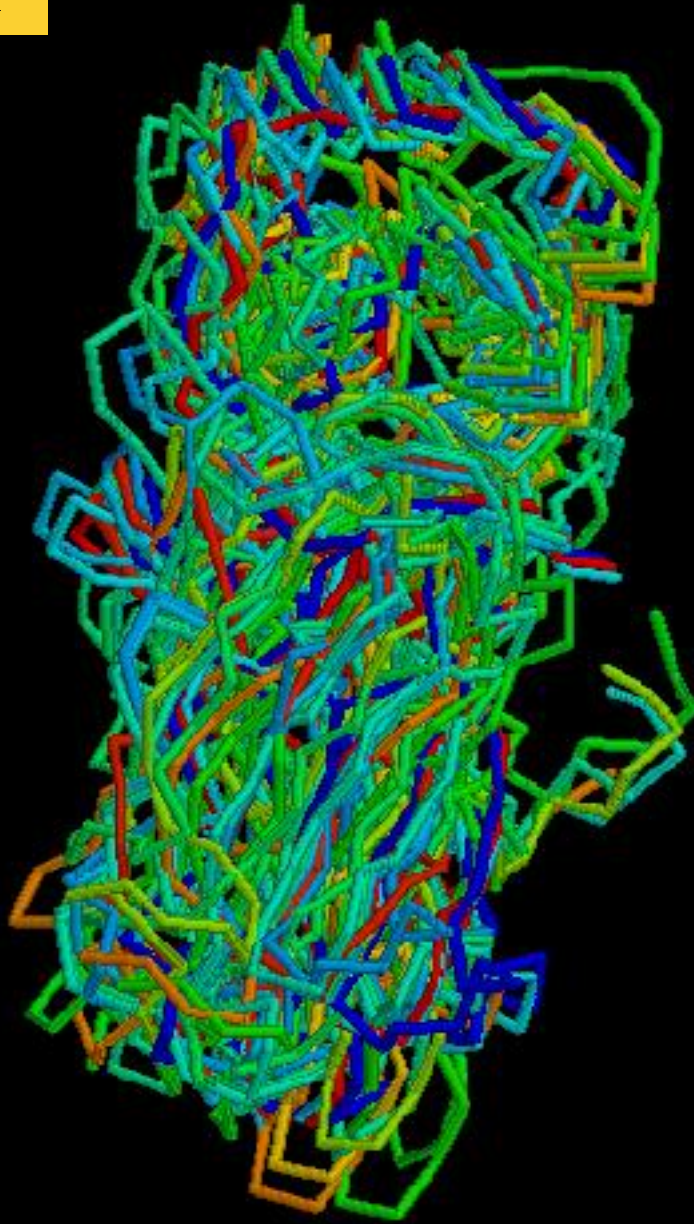
- 8 different alignments

- 2 different programs:

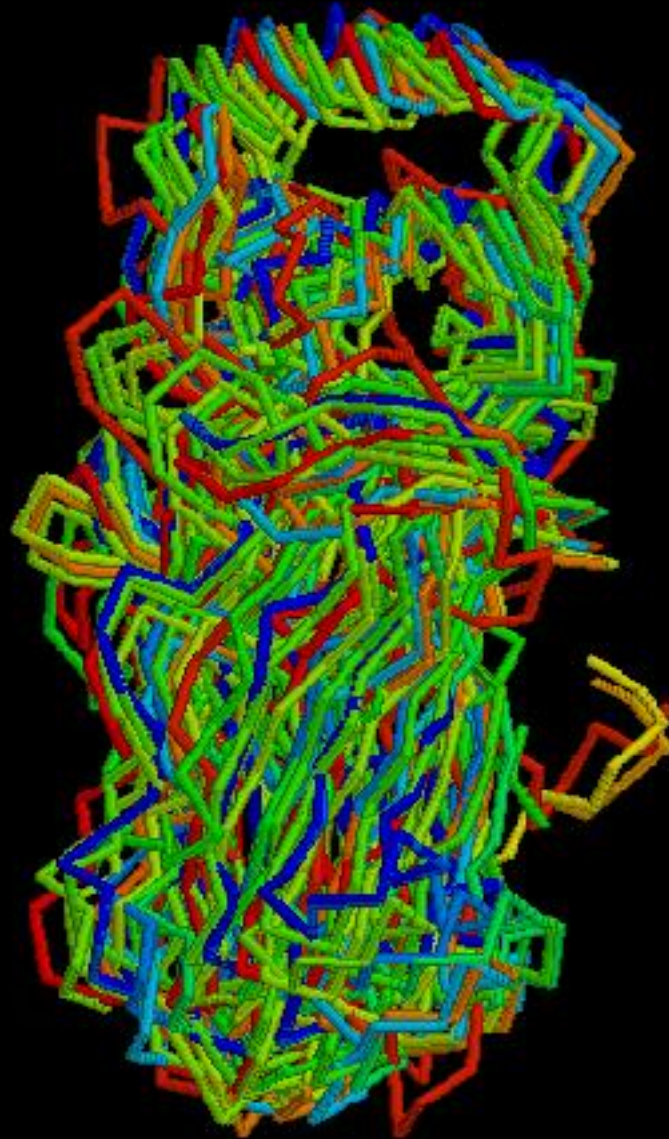
3D-JIGSAW & *Pmodeller*



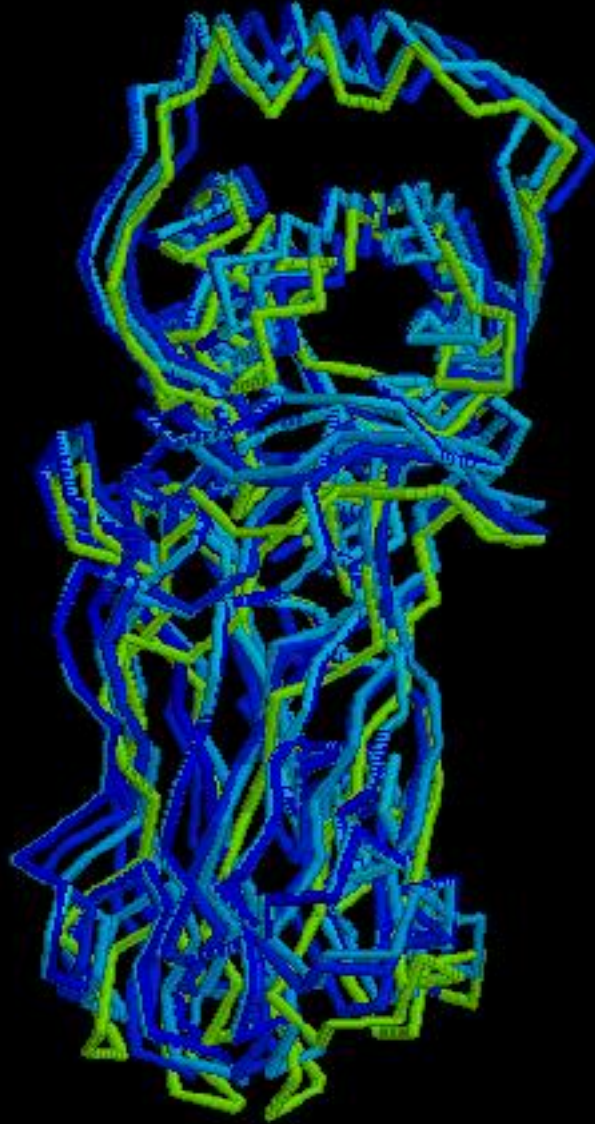
generation 1



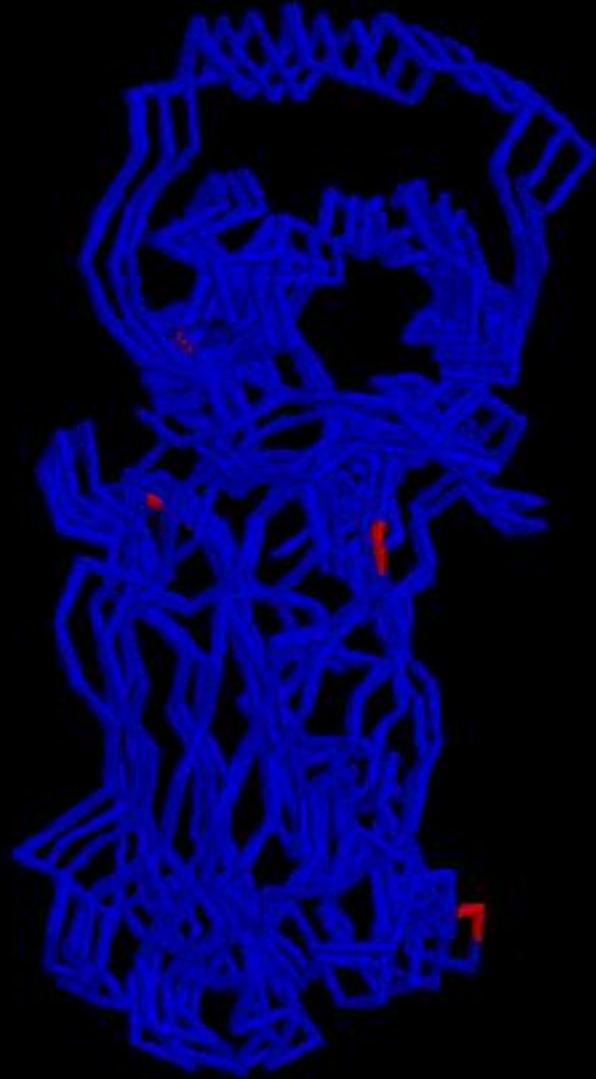
generation 4



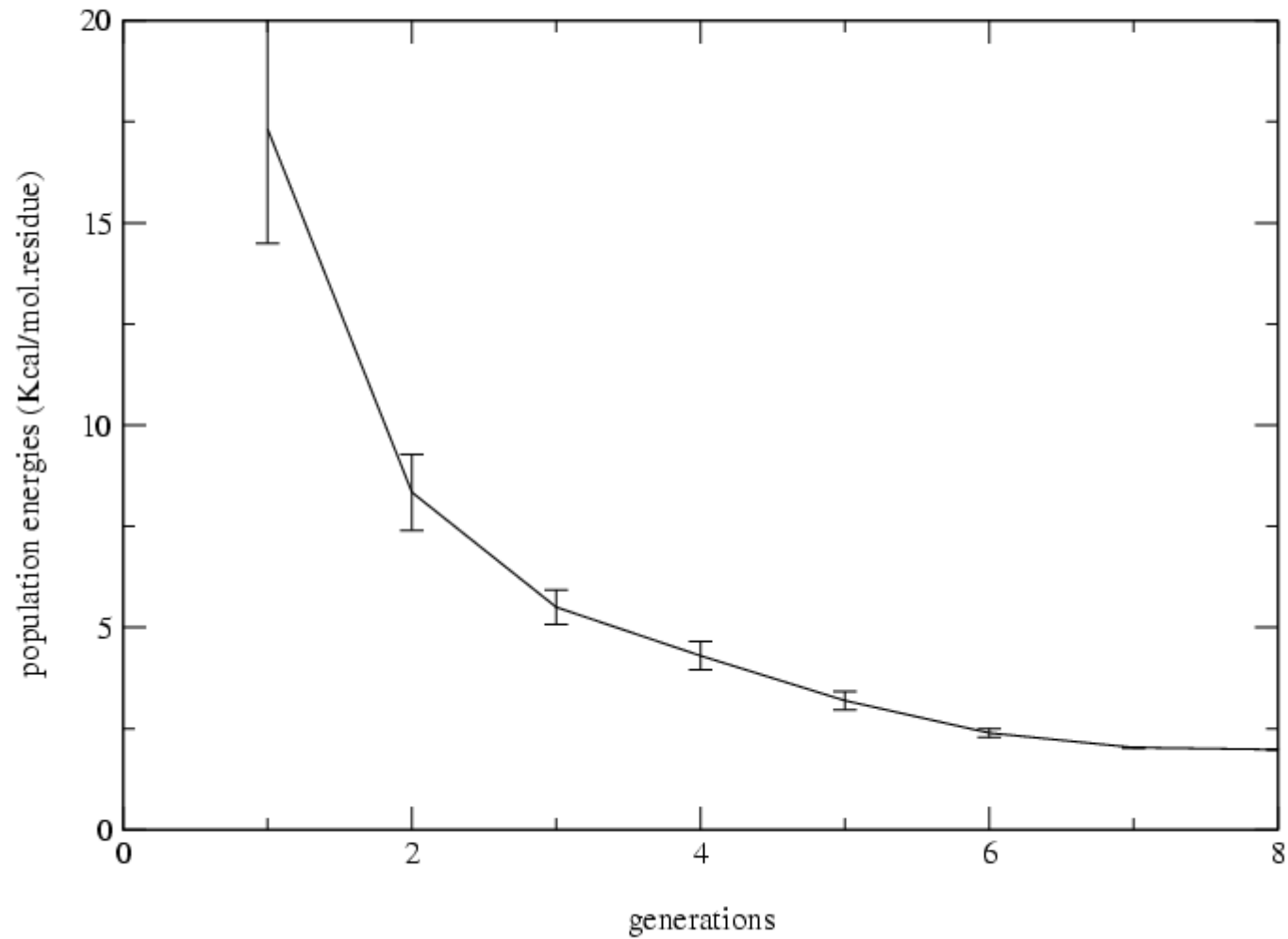
generation 7



last population



in silico Protein Recombination experiment: T0134



in silico protein^xrecombination (test version)

<http://www.bmm.icnet.uk/3djigsaw/recomb>

Description

This program performs artificial selection (through recombination + mutation) over a population of protein atomic models seeded by the user, with the aim of obtaining a more accurate and energetically favourable atomic conformation than any starting model but based on all. So please make sure that your input file contains only models for the same protein. Enjoy yourself.

job identifier

your e-mail address

Please specify a file ([PDB format](#), [TER-minated chains](#)):

or

... paste your PDB coordinates here:

pop size

selected prop

forced improvement

last gen

mutation rate

[help](#)

contrera@cancer.org.uk BMM disclaimer

Thanks to
the Biomolecular Modelling Laboratory

Paul Bates

Paul Fitzjohn

Pall Jonsson

Graham Smith

Chris Page

Marc Offman

www.bmm.icnet.uk

Thanks to Cancer Research UK and
to SAC-CASP5 organizers

Possible applications of comparative modelling*

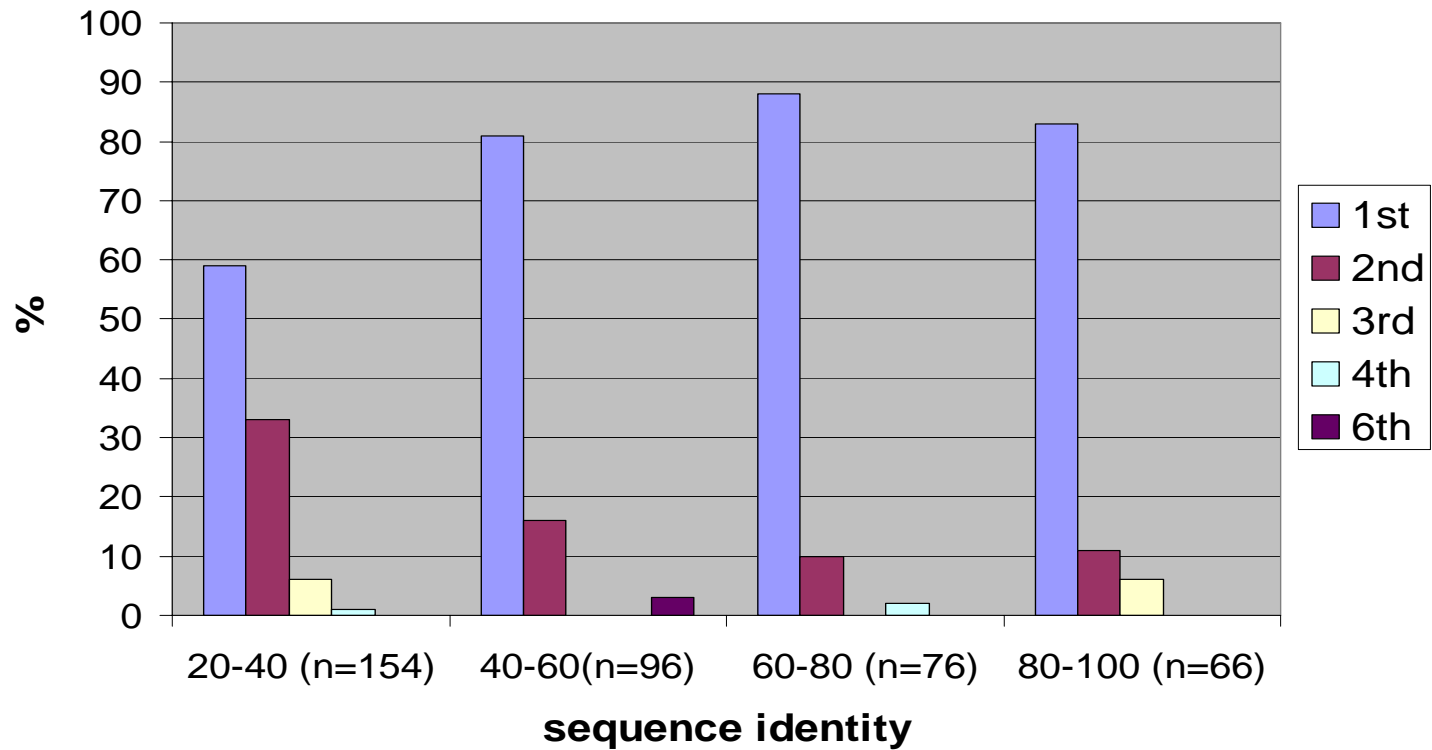
Depending on the sequence identity between query and template:

- >90% virtual ligand screening
- >40% defining antibody epitopes
- >40% molecular replacement in X-ray crystallography
- >20% support site directed mutagenesis
- >20% fitting into low resolution electron density maps

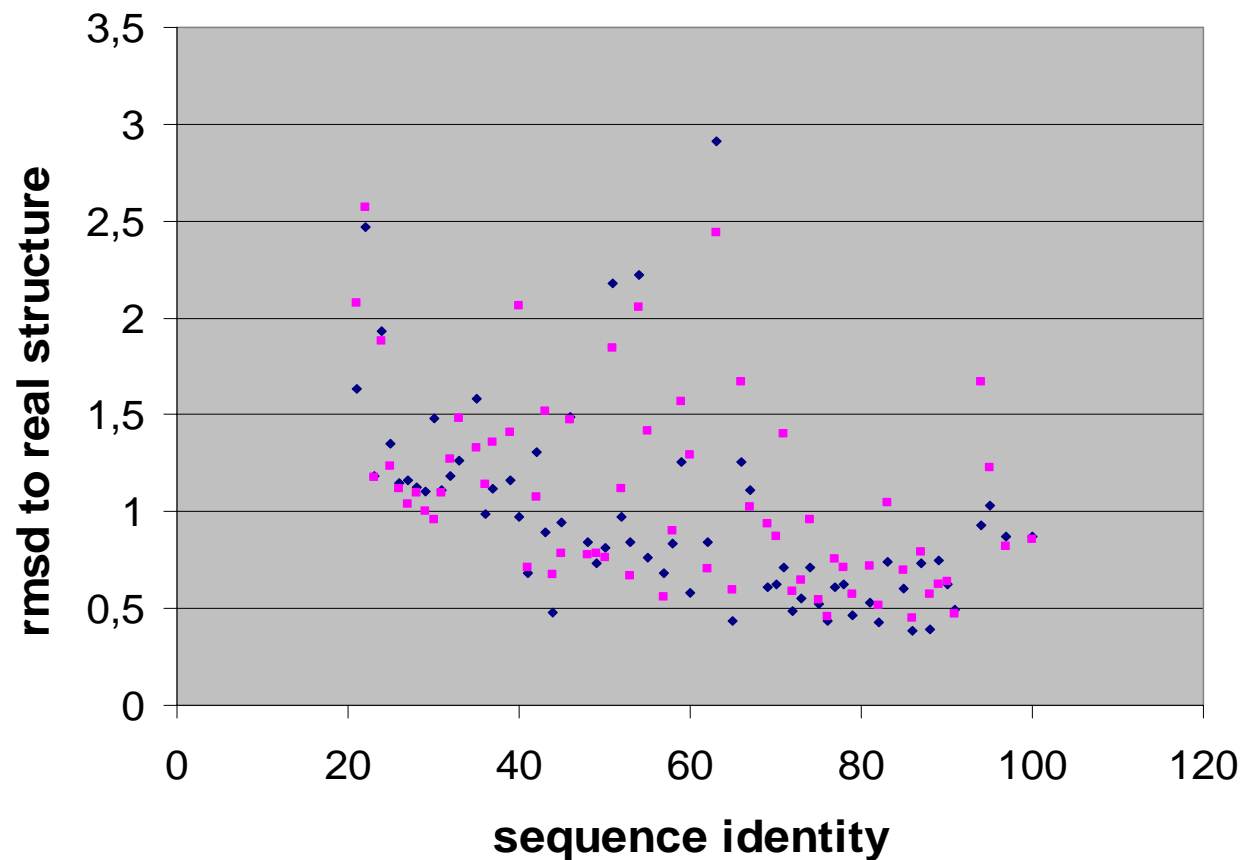
* Baker & Sali (2001) Science 294: 93-96

Selecting templates

Best template for comparative modelling

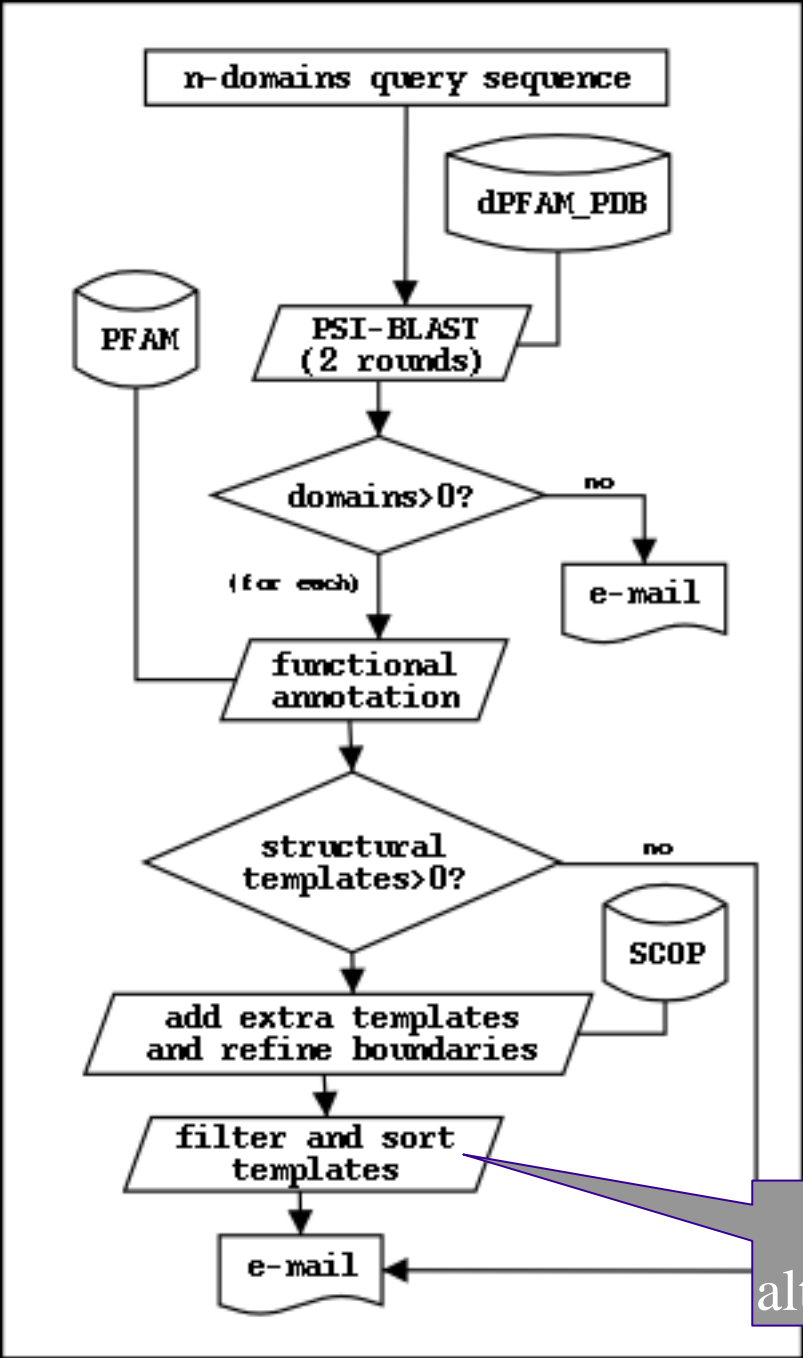


Single vs Multiple template modelling



Overall:
54/97 S
43/97 M

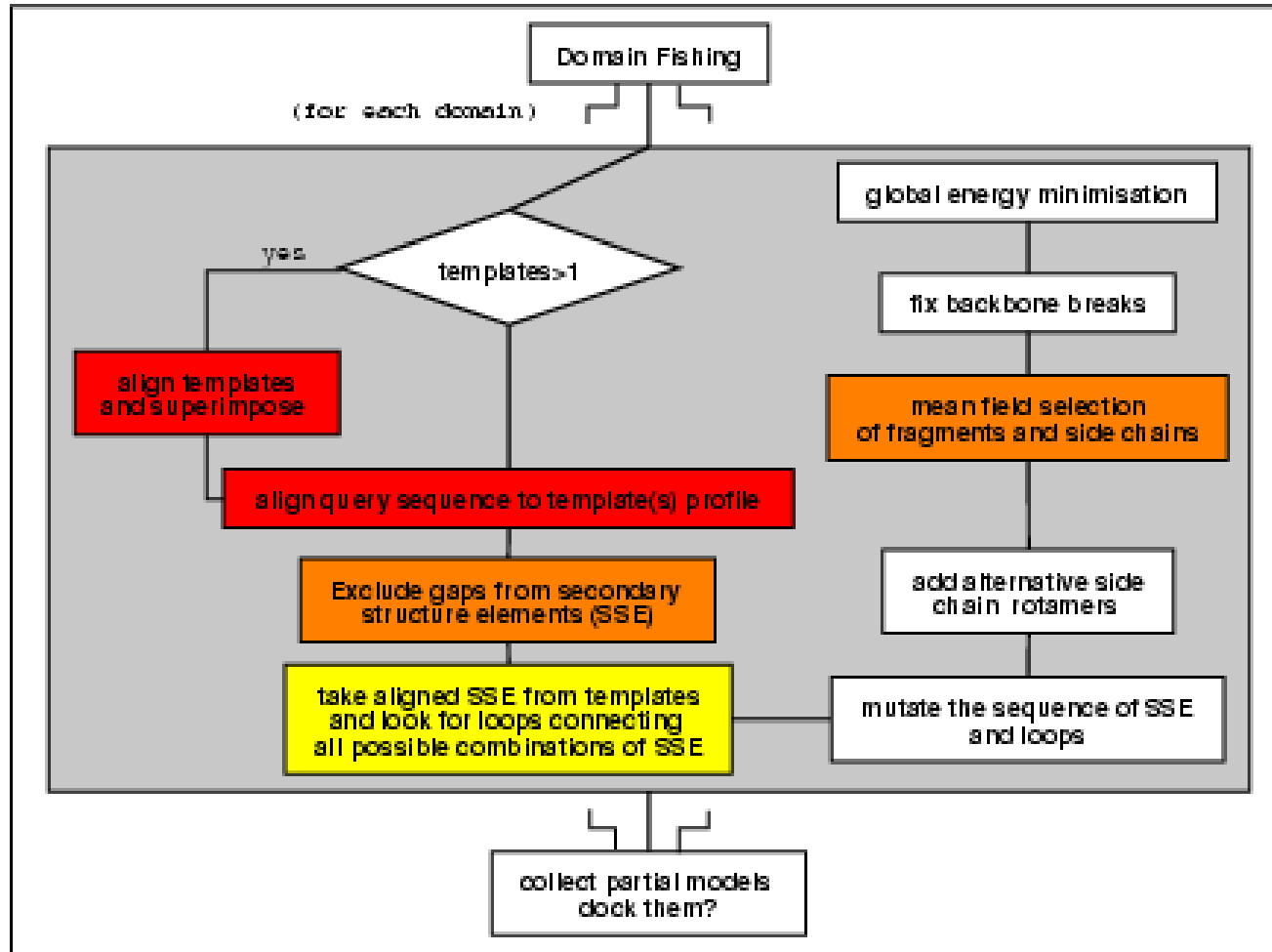
Domain Fishing



up to 7
alternative alignments

3D-JIGSAW

Example



Conclusion

We have done:

- automatic domain identification
- improved alignments
- multidomain modelling

We want to do next:

- better template selection (energies)
- connecting domains
- different multi-template strategies