

in silico protein recombination: a genetic algorithm
applied to template and alignment selection in
comparative modelling

Bruno Contreras-Moreira
Biomolecular Modelling Laboratory
London Research Institute
September 2003

comparative modelling

Predictive technique to build a molecular model for a sequence based on homologous proteins whose structure is known.

query sequence



search for templates
and selection



alignments
to template(s)



query inherits
backbone from
template(s) **+FR**



loops are
modelled

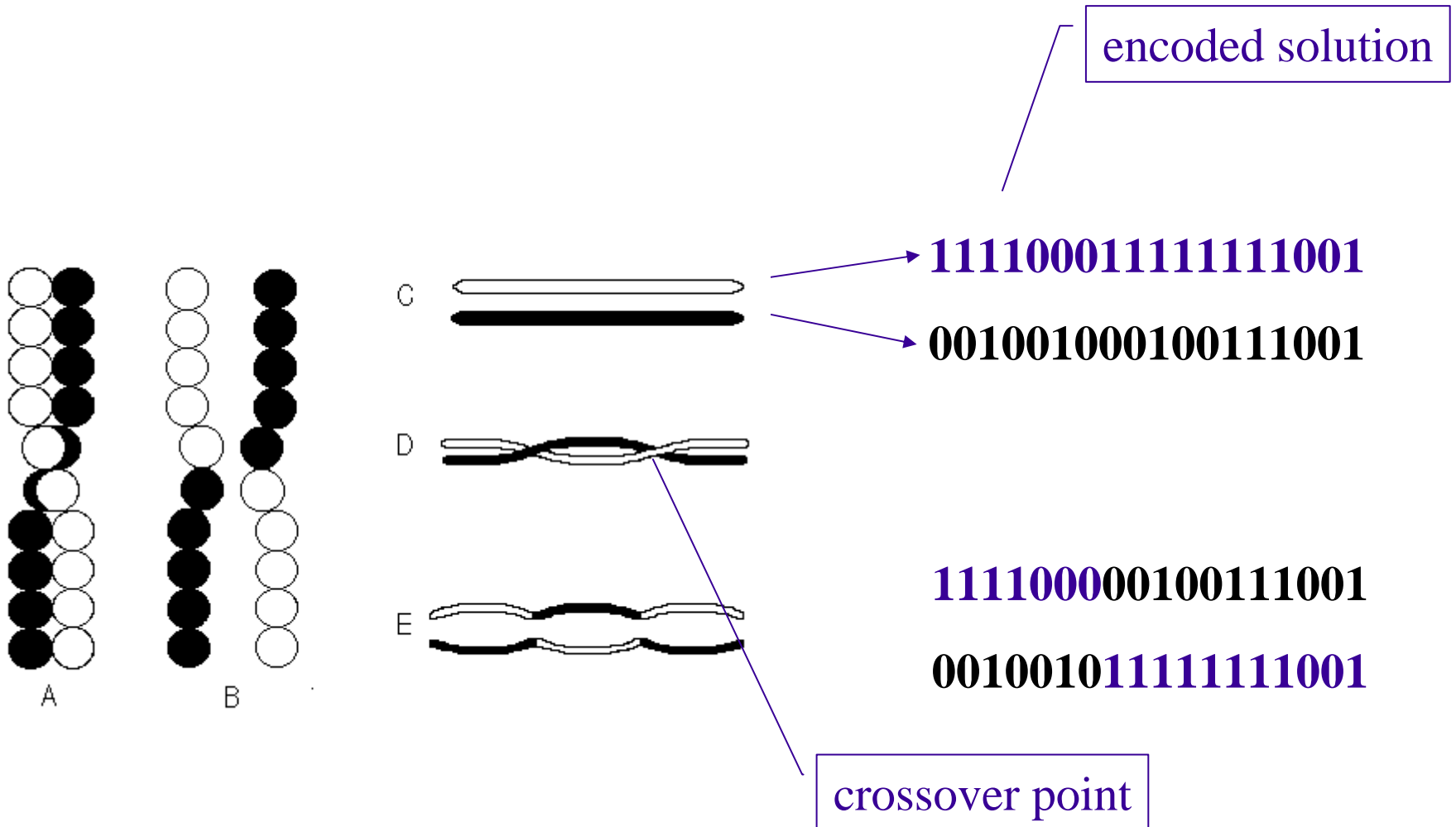


model refinement



error estimates
on final model

chromosome evolution & computational analogy: genetic algorithms



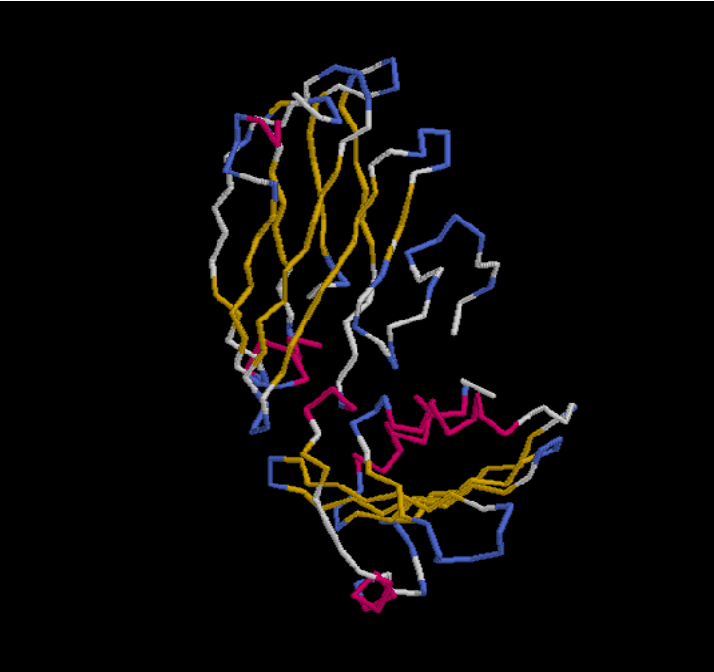
a genetic algorithm applied to Comparative Modelling

- 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  GEF-VQNGAPEEE--QLPPSSYSLLAENSYVKMTCDIRGSLQEDSQVTVAIVLENRSS
lqts_A  GSPGIRLGSSSEDNFARFVCKNMGVLF-ENQLLQI--GLKSEFRQNLG-RMFIFYGNKTS
SS      CCCCCCCCCCCCCCHHHHCCCCCEEEE-ECCCEEE--EEEEEEECCEE-EEEEEEEECCC
```

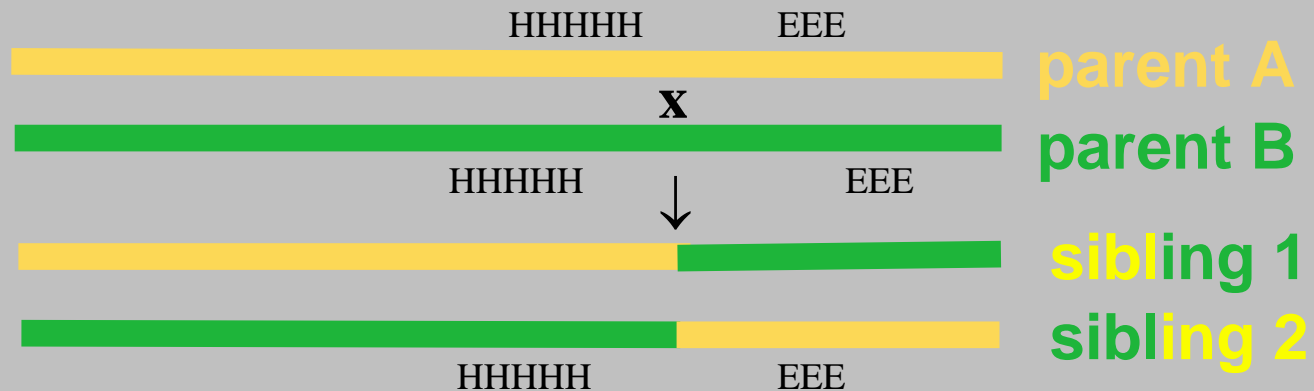


(1model = 1PDB template + 1alignment)



recombination

```
model recombination( model A , model B )  
{  
  do sequence_alignment( A , B );  
  do sequence_superimposition_Cβ( A , B );  
  do refine_superimposition_close_Cβ( A , B );  
  do draw_crossover_point( A , B ); /* out of SS? */  
  return create_model(A , B , crosspoint );  
}
```



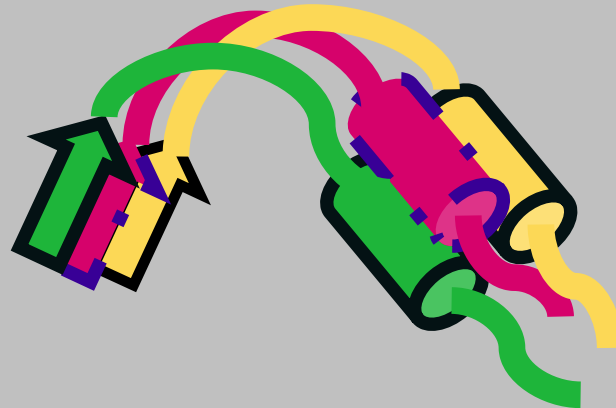
mutation

```
model mutation( model A , model B )  
{  
  do sequence_alignment( A , B );  
  do sequence_superimposition_C $\beta$ ( A , B );  
  return create_Cartesian_average_model(A , B);  
  /* quality checks, minimization? */  
}
```

parent A

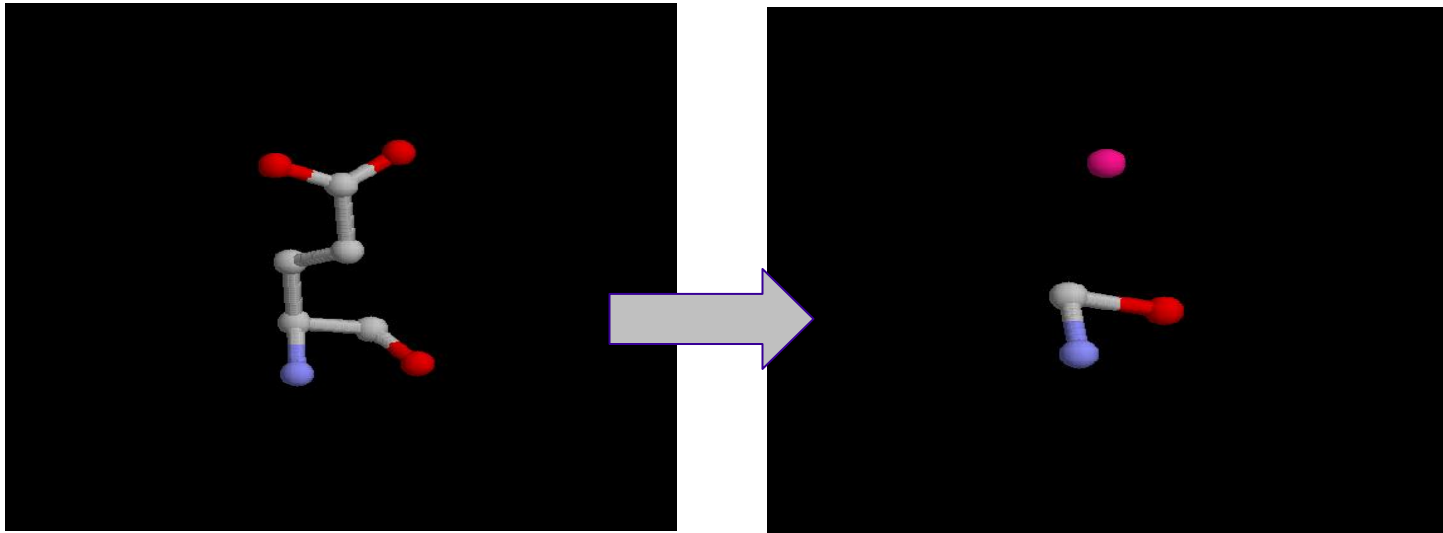
parent B

sibling



protein fitness

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



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

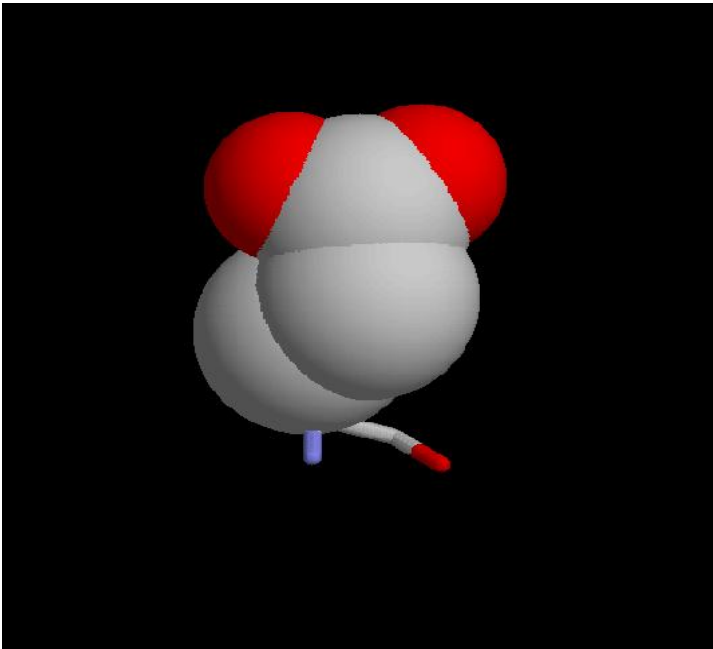
where i, j are pairs of pseudoatoms in protein p

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

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



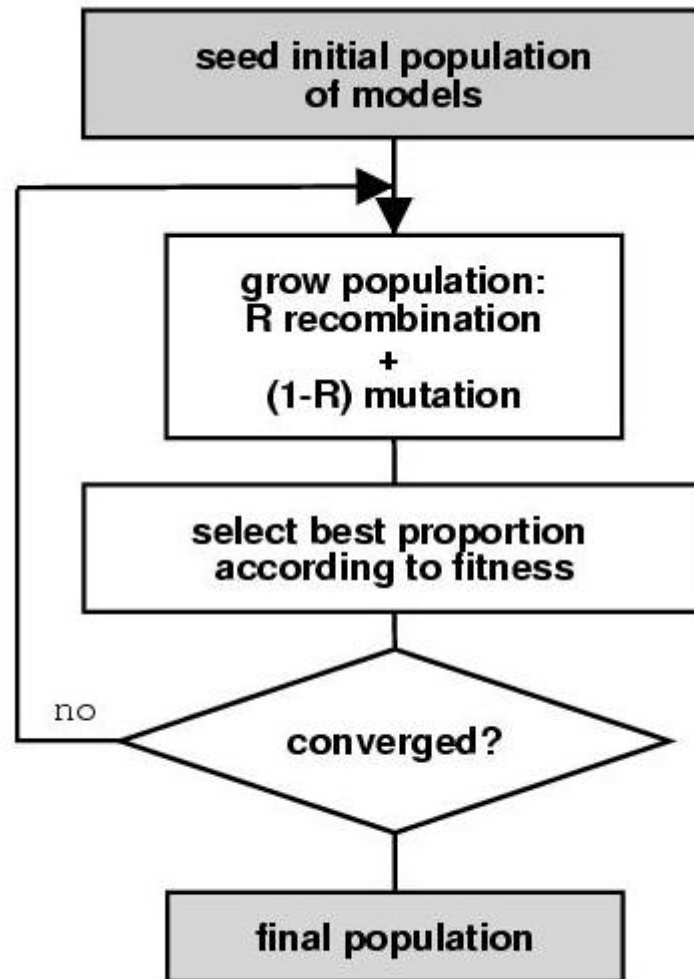
$$\sum_i (\text{SA}_i \cdot \Delta G_{\text{solv}_i}) \quad (\text{in Kcal/mol})$$

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

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

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

in silico protein recombination algorithm

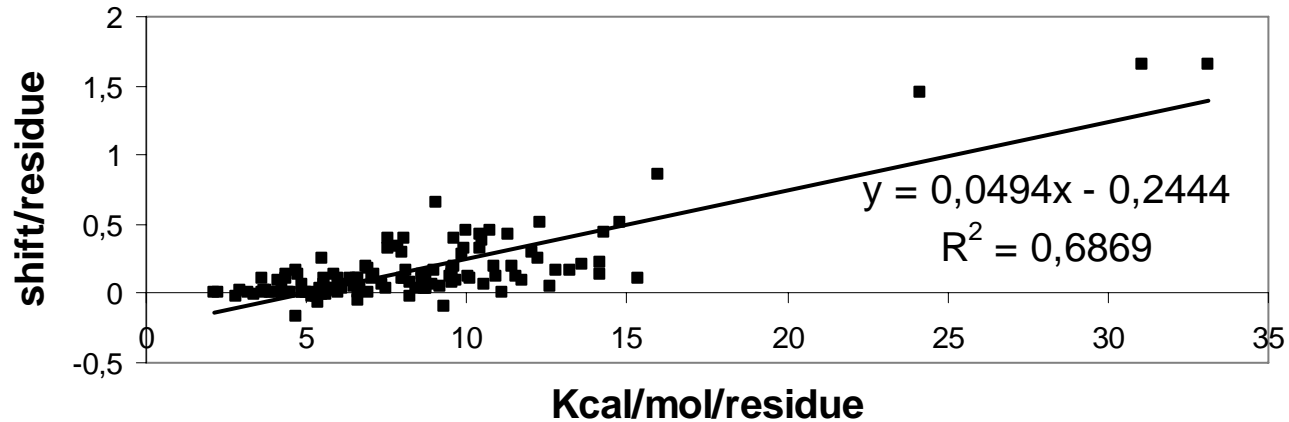


Contreras-Moreira, Fitzjohn and Bates (2003) *J Mol Biol*, **328**: 593-608.

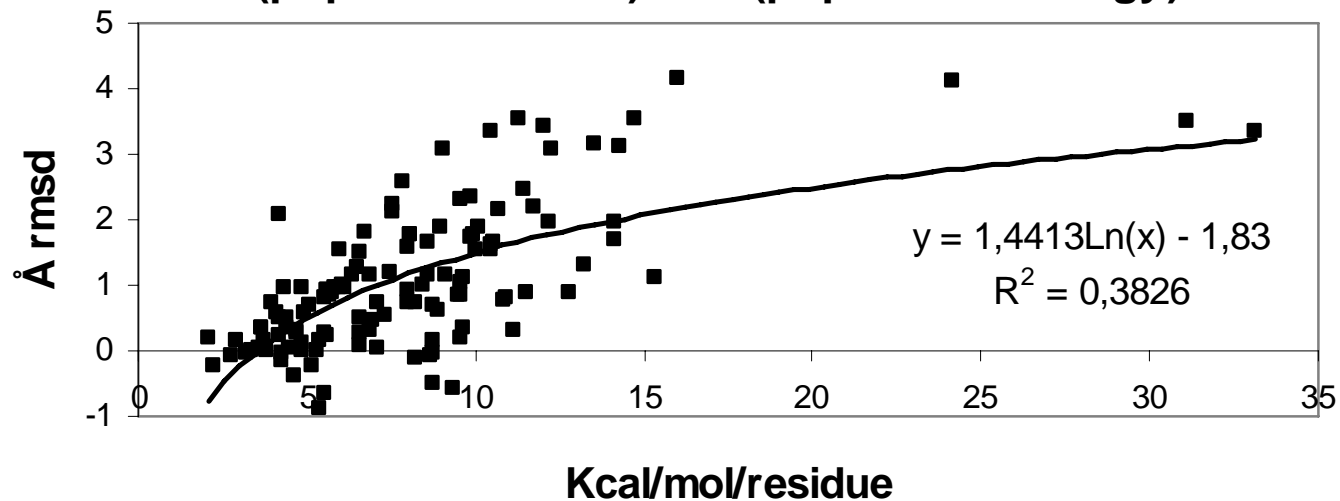
in silico protein recombination: performance

(benchmark on 130 SCOP families)

d(population energy) vs d(alignment shift)



d(population rmsd) vs d(population energy)

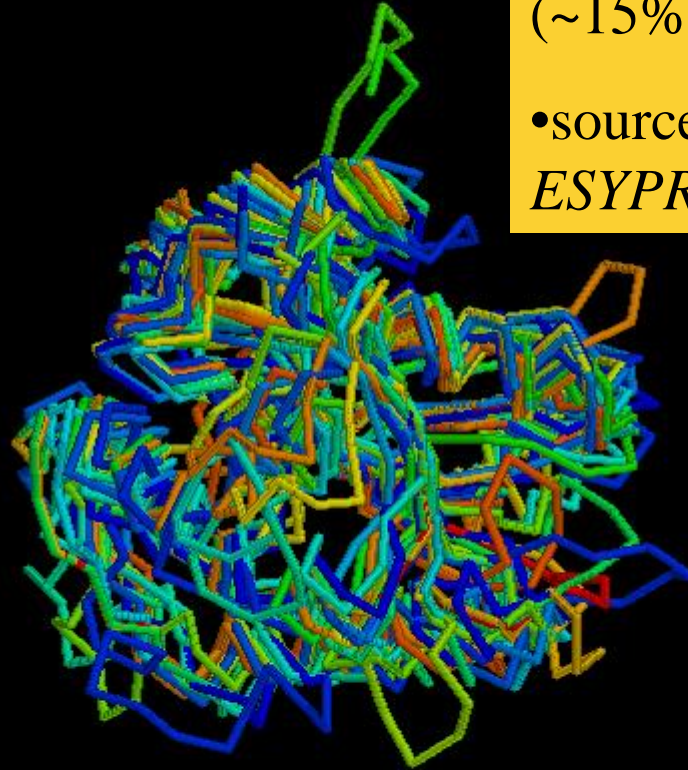


CASP5 example: T0192

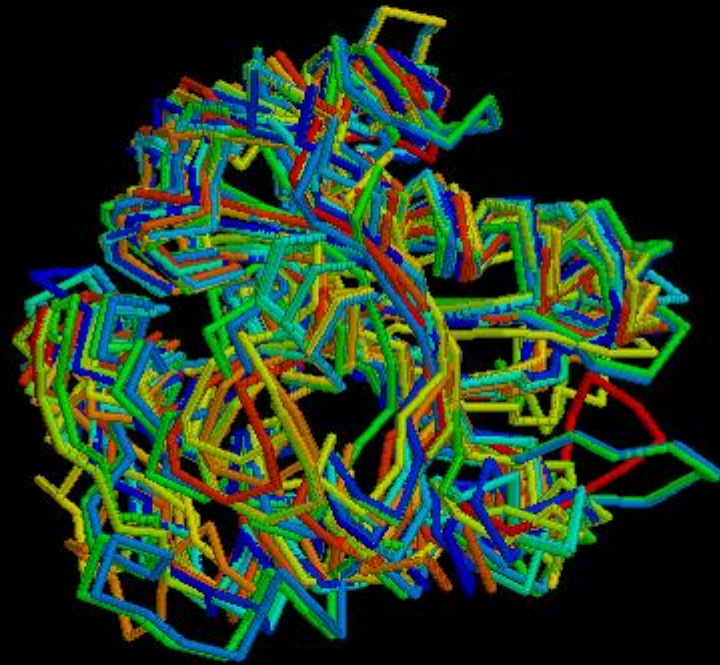
Human acetyltransferase

- 2 templates: 1QSM & 1QSO (~15% SeqID), 12 alignments
- sources: *3D-JIGSAW*, *FAMS*, *ESYPRED* & *Pmodeller*

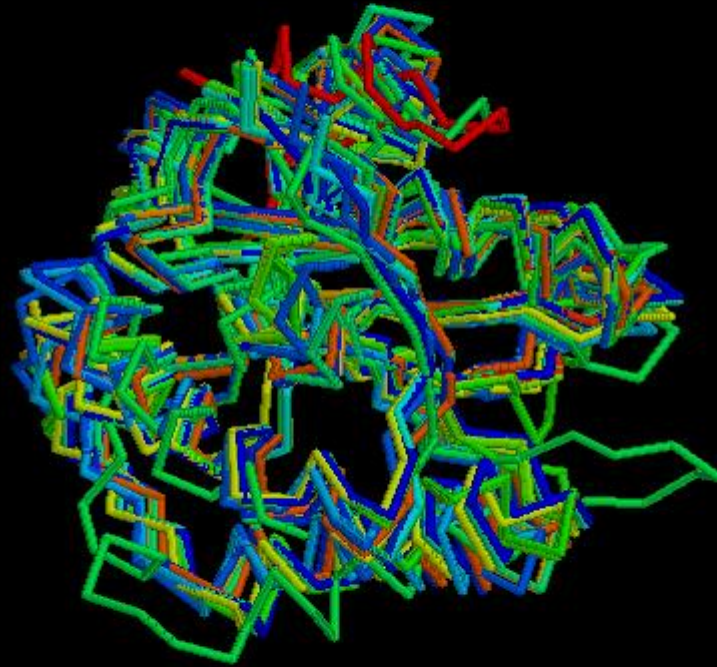
generation 0



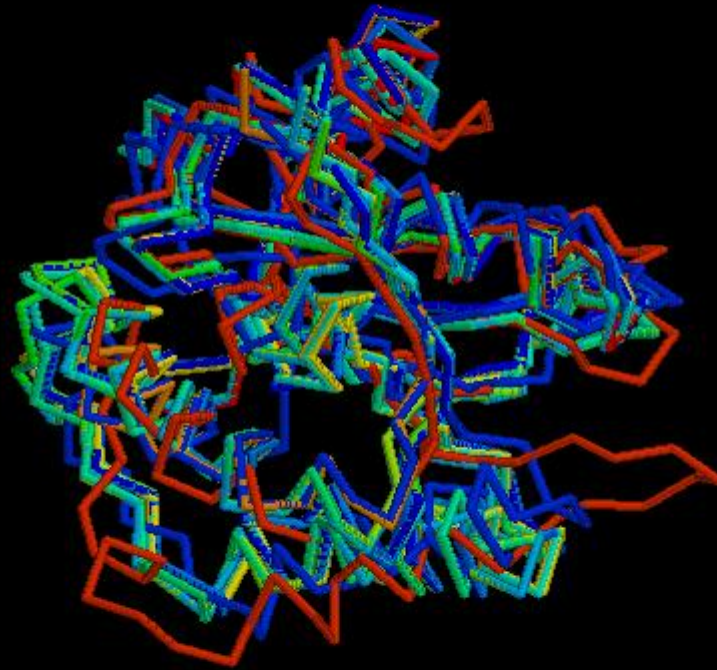
generation 2



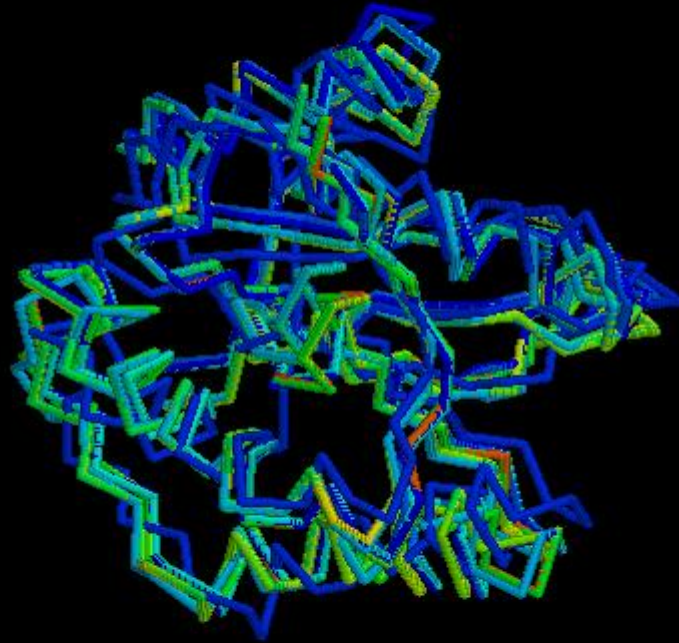
generation 4



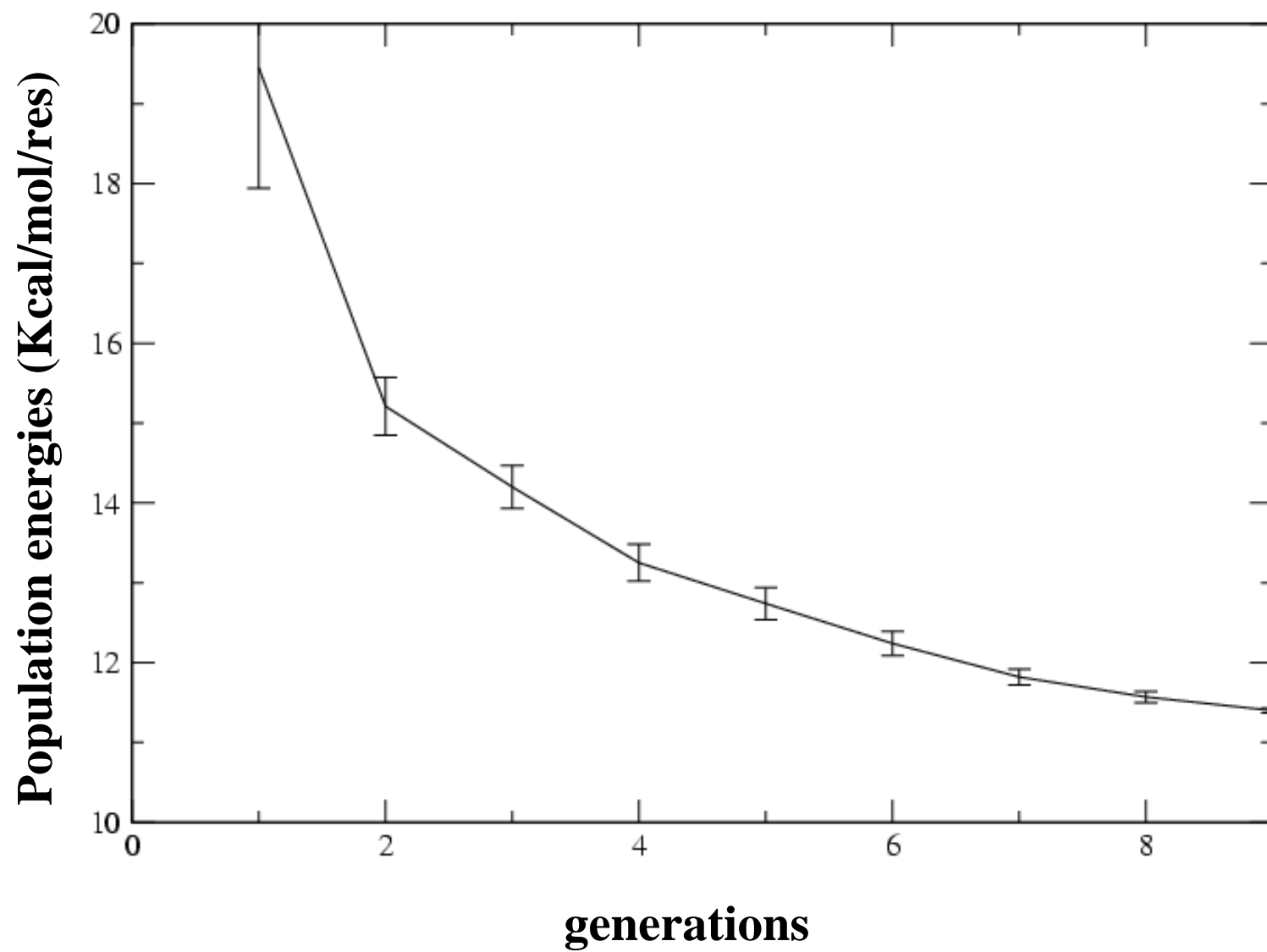
generation 6



generation 8(last)

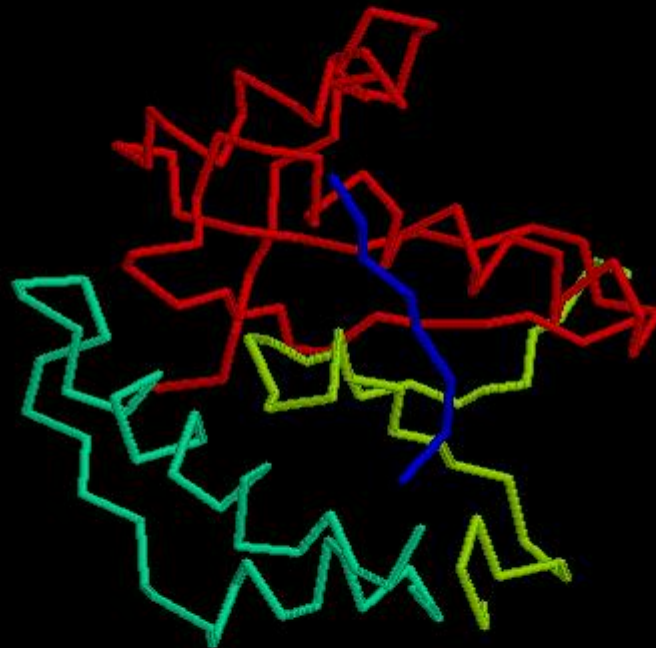


in silico Protein Recombination experiment: T0192_2



best model (after 8 generations)

model	GDT_TS	AL_4
mod1	45	61
mod2	63	81
mod3	57	72
mod4	54	64
mod5	54	64
mod6	61	80
mod7	61	76
mod8	61	80
mod9	62	78
mod10	65	77
mod11	62	78
mod12	60	71
<i>average</i>	58	74
rec_8gen	61	81
<i>bestCASP5</i>	66	85



$$\text{GDT} = (\%<1\text{\AA} + \%<2\text{\AA} + \%<4\text{\AA} + \%<8\text{\AA}) / 4$$

$$\text{AL}_4 = \%(<4\text{\AA} \text{ AND shift}\pm 4)$$

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

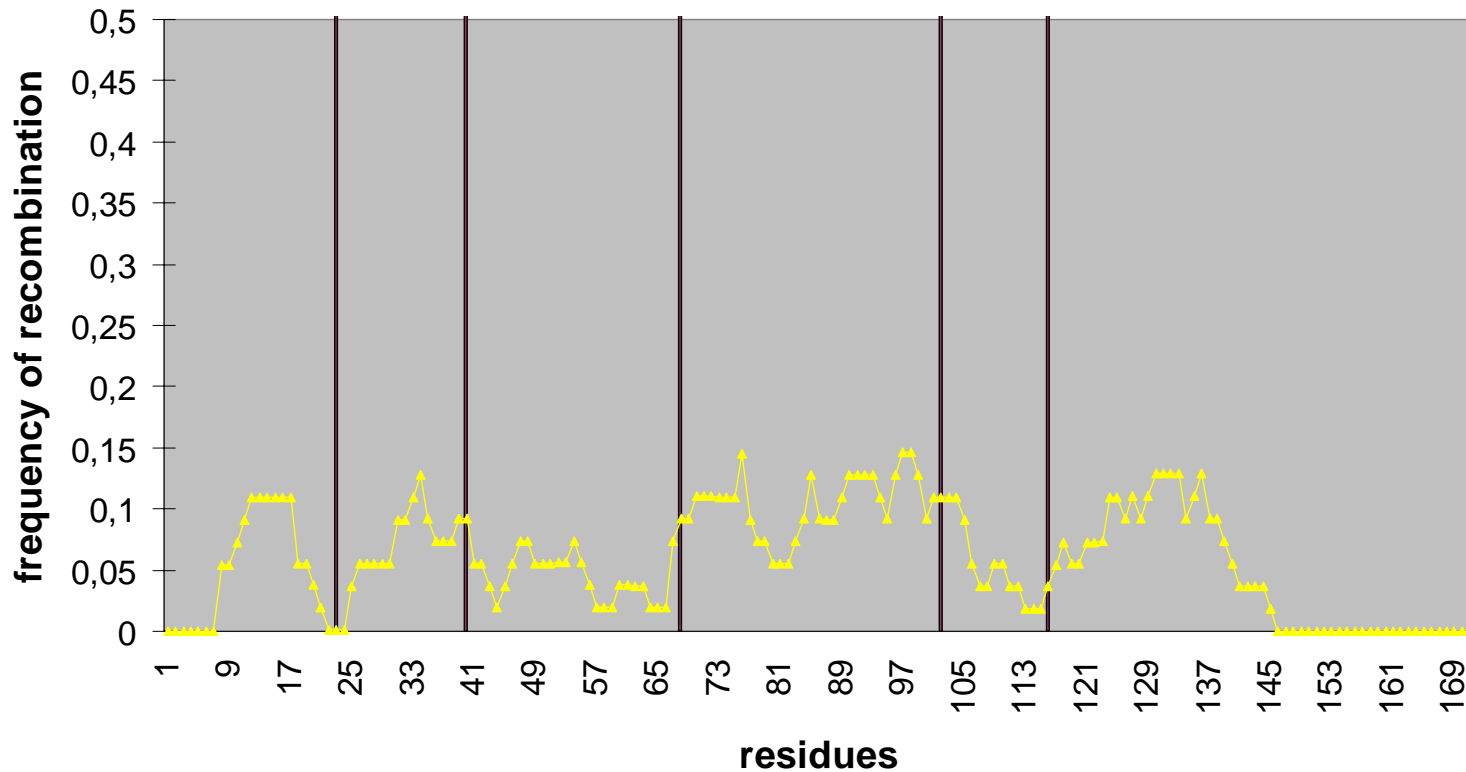
- converges close to the best initial model
 - is able to recover alignment errors
- often last population contains different conformations

EVOLUTIONARY INSIGHT

- high-frequency crossover spots seem to occur away from intron boundaries (work with Pall Jonsson)

Crossover points and introns boundaries: T0192

average of 5 simulations
7 homologues < 20%SeqID
origin: yeast , *B.subtilis* , *M.tuberculosis*



Biomolecular Modelling Laboratory

Paul Bates

Paul Fitzjohn

José Jimenez

Pall Jonsson

Chris Page

Graham Smith

&

Marc Offman

Fabien Birzele

(students from LMU München)

www.bmm.icnet.uk

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

the basic principle

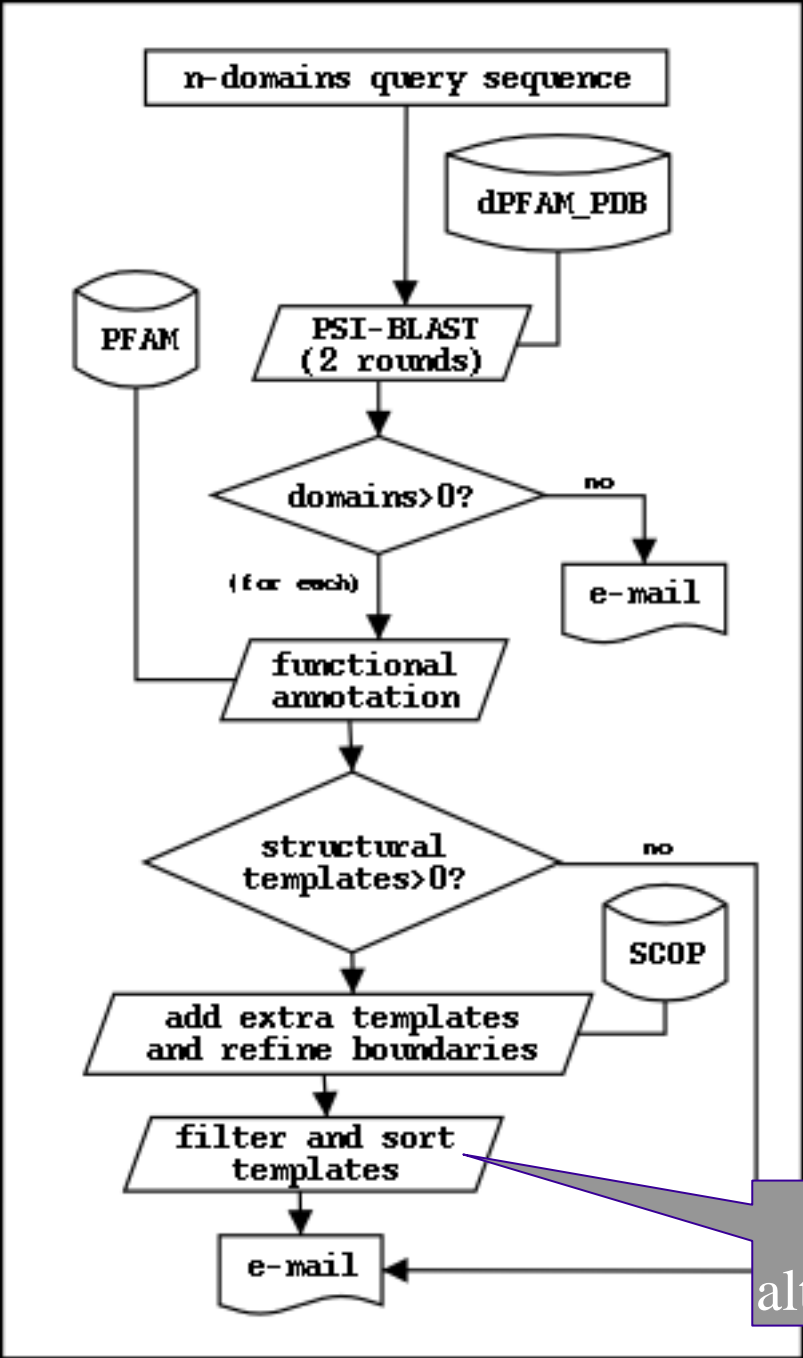


```
-QTSVSPS-KVILPRGGSVLVTCSTS-CDQPKLLGIET---P-LPKKELLPGNNRKVYE  
FKIETTPESRYLAQIGDSVSLTCSTTGCESP-FFSWRTQIDSPLNGKVTN--EGTTSTLT
```

```
LS--NVQE-DSQPMCYSNCPDGGSTAKTFLTV--  
MNPVS-FGNEHSYLCTATCESRKLEKGIQVEIYS
```

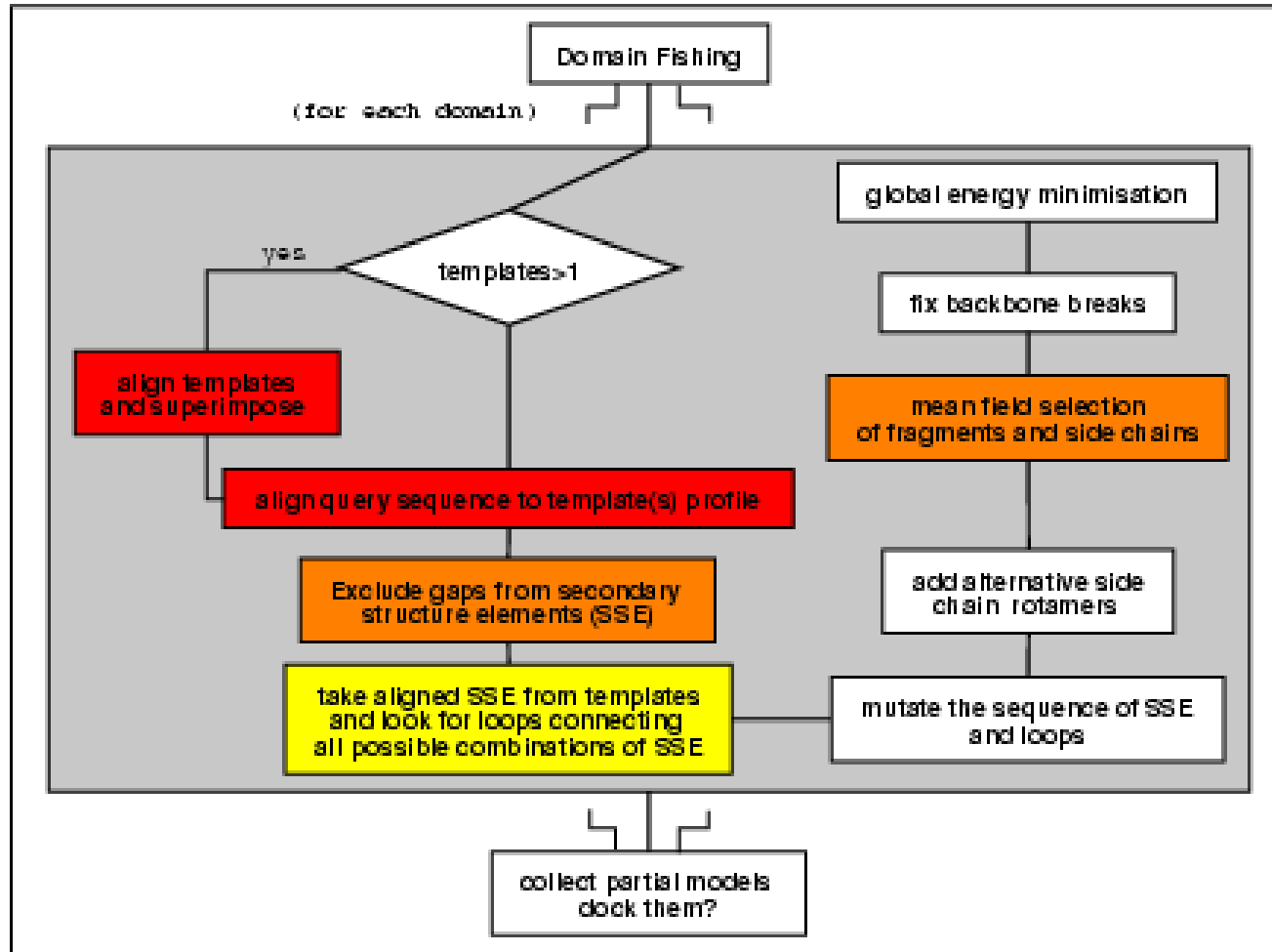
structural agreement = f(sequence similarity)

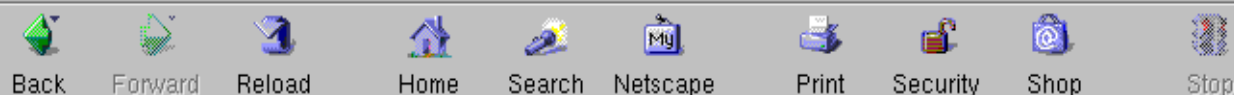
Domain Fishing



up to 7
alternative alignments

3D-JIGSAW





[Interactive 3D-JIGSAW](#) [legend](#) [home](#) [disclaimer](#) [contact us](#) **HHCCEE: predicted Helix, Coil or Strand**

Possible structural templates in PDB

name	from	to
1bza_# Model!?	28	287
1shv_A Model!?	26	292
1g56_A Model!?	26	292
1ck3_A Model!?	26	290
1jtd_A Model!?	27	288
1fgg_A Model!?	26	288
1bt1_# Model!?	26	290
1bt5_A Model!?	26	290
1erq_A Model!?	26	288

truncated alignments?
wrong templates? PDB code chain first residue last