

T-Coffee: A Novel Method for Fast and Accurate Multiple Sequence Alignment

Beth Crompton



Multiple Alignment is Hard

- General empirical models fail at <30% sequence identity
- Progressive Alignment is too greedy
- Simultaneous alignment is very resource intensive

T-Coffee is the solution to all your problems



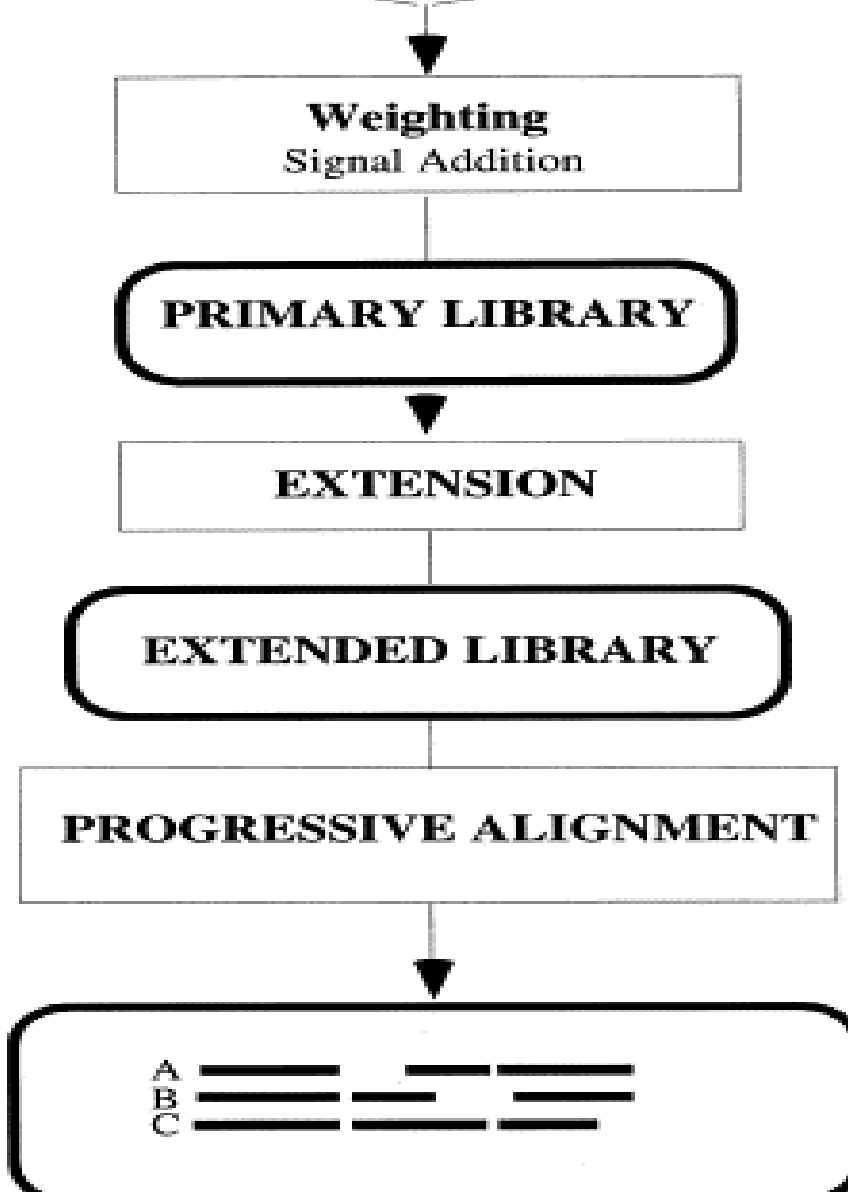
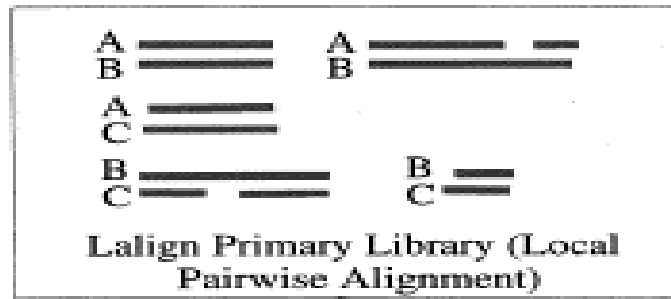
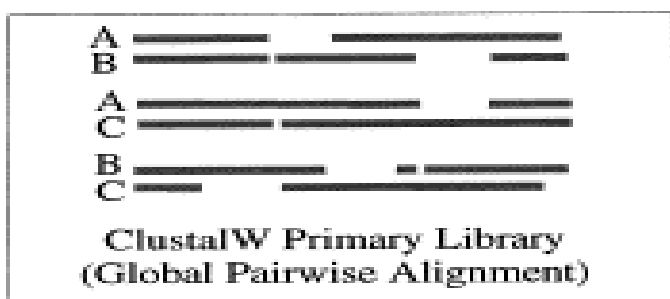
<http://www.saynotocrack.com/wp-content/uploads/2007/07/evil-latte.jpg>

Tree-based Consistency Objective Function For AlignmEnt Evaluation



- ♦Generates multiple alignments from heterogeneous data source in reasonable runtime
- ♦Optimizes the multiple alignments to fit the pairwise alignments

http://farm1.static.flickr.com/120/284971917_d60dcce9a0.jpg?v=0



Weighted Primary Library

Align pairs using ClustalW & Lalign

Each aligned pair is a constraint

Assign weights to each constraint using Sequence Identity method

Weight represents correctness of constraints



<http://www.emmitsburg.net/humor/pictures/2007/Fw%20Coffee%20art.jpg>

Extending the library

Problem: fitting a set of weighted constraints into a multiple alignment is NP-complete

Solution: heuristics!



Weights

Already indicate similarity between two sequences

Should indicate consistency with all other sequences involving those residues

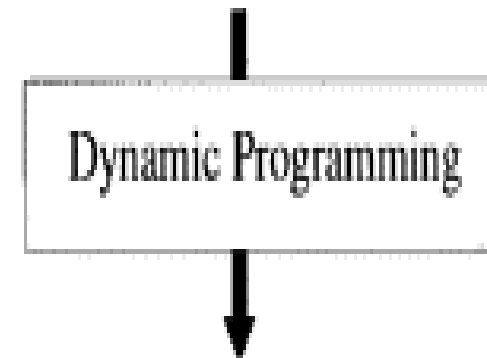
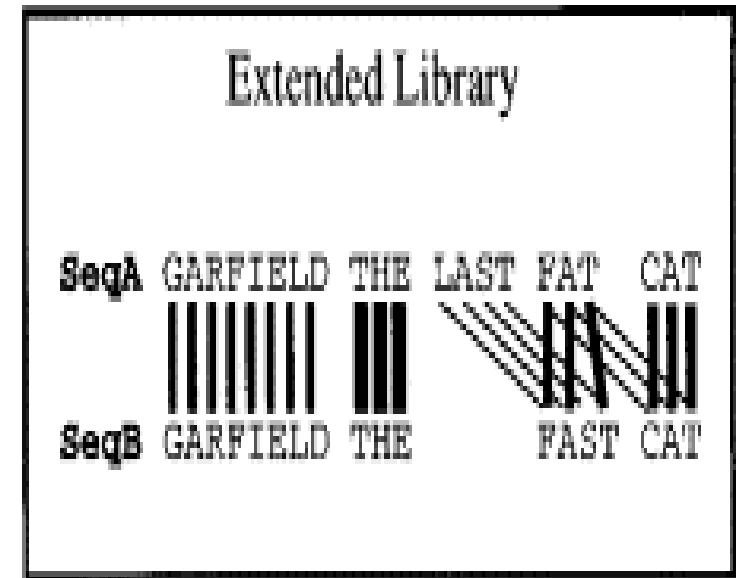
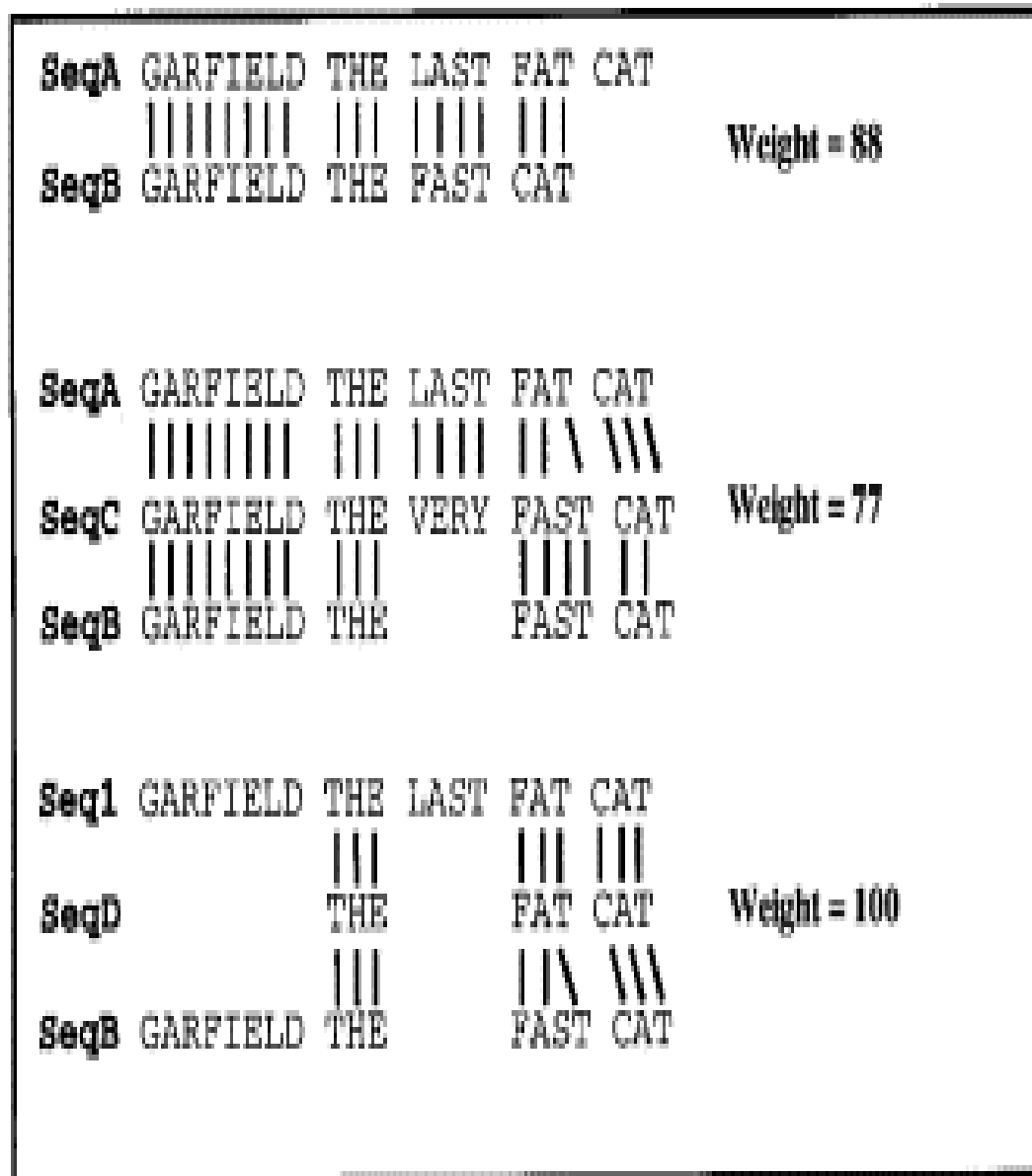
Weight of a pair becomes sum of all weights found examining all triplets involving that pair

Weight of non-existent residues becomes 0

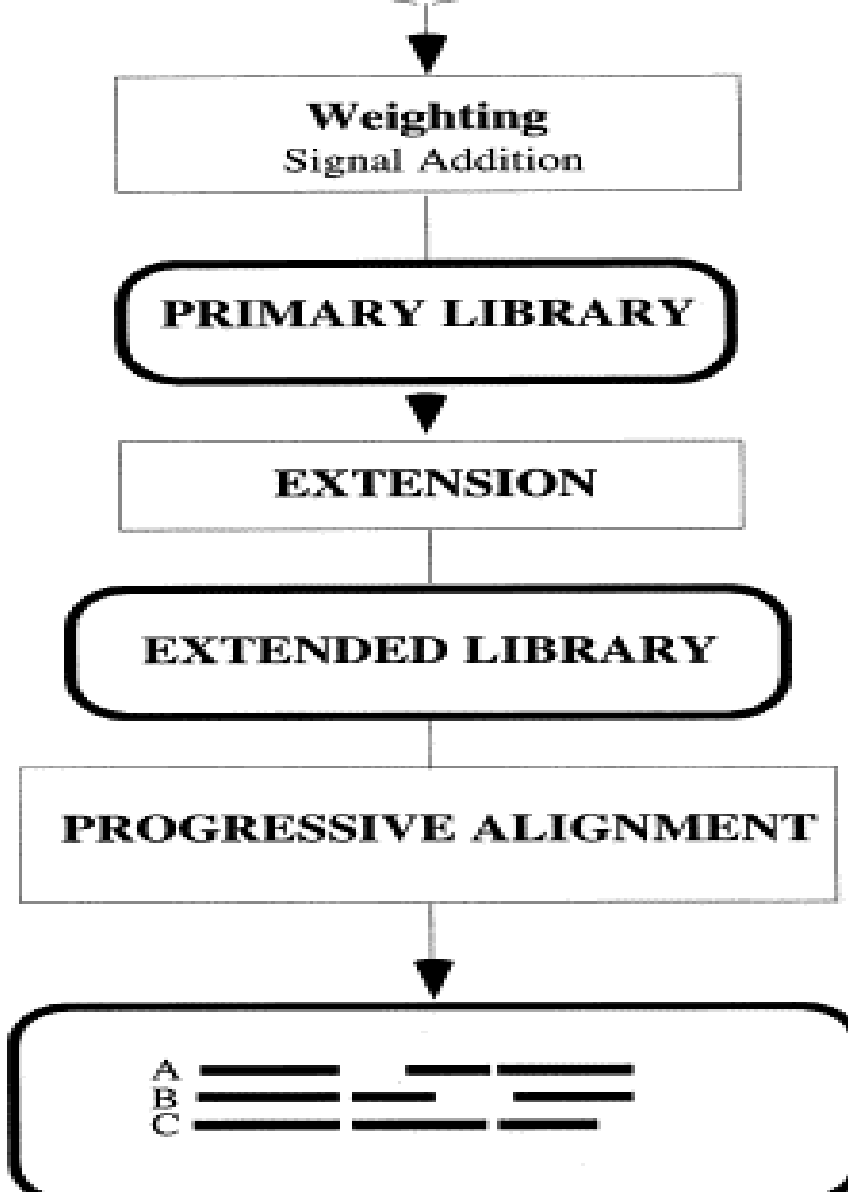
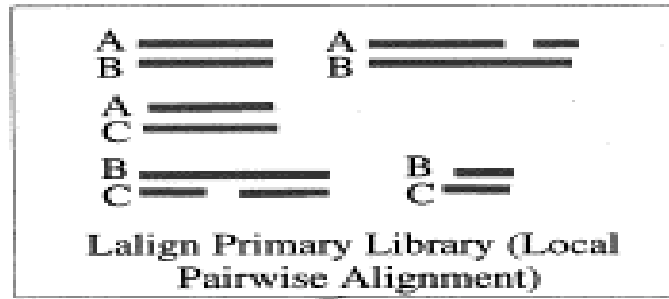
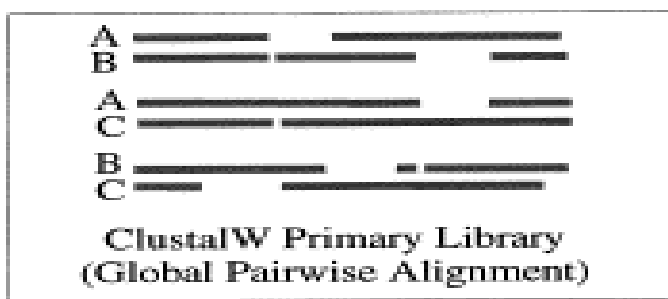


Library extension-continued

c) Extended Library for seq1 and seq2



SeqA GARFIELD THE LAST FA-T CAT
SeqB GARFIELD THE ---- FAST CAT



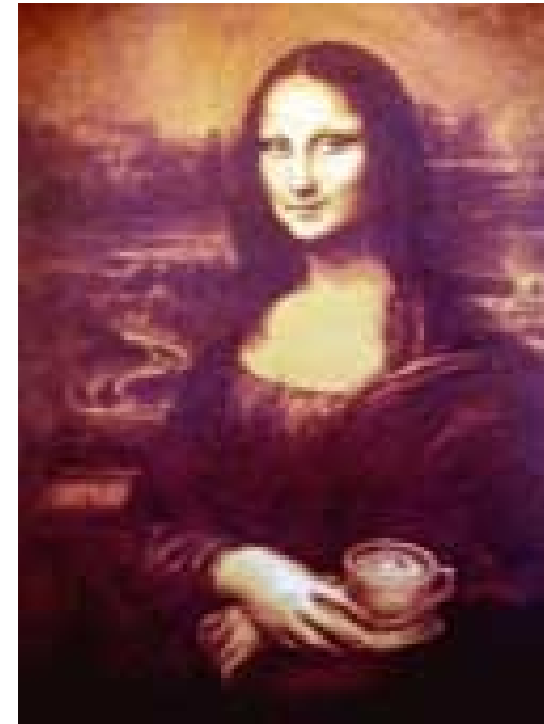
Progressive Alignment

Pair-wise alignments

distance matrices

guide tree

multiple alignment



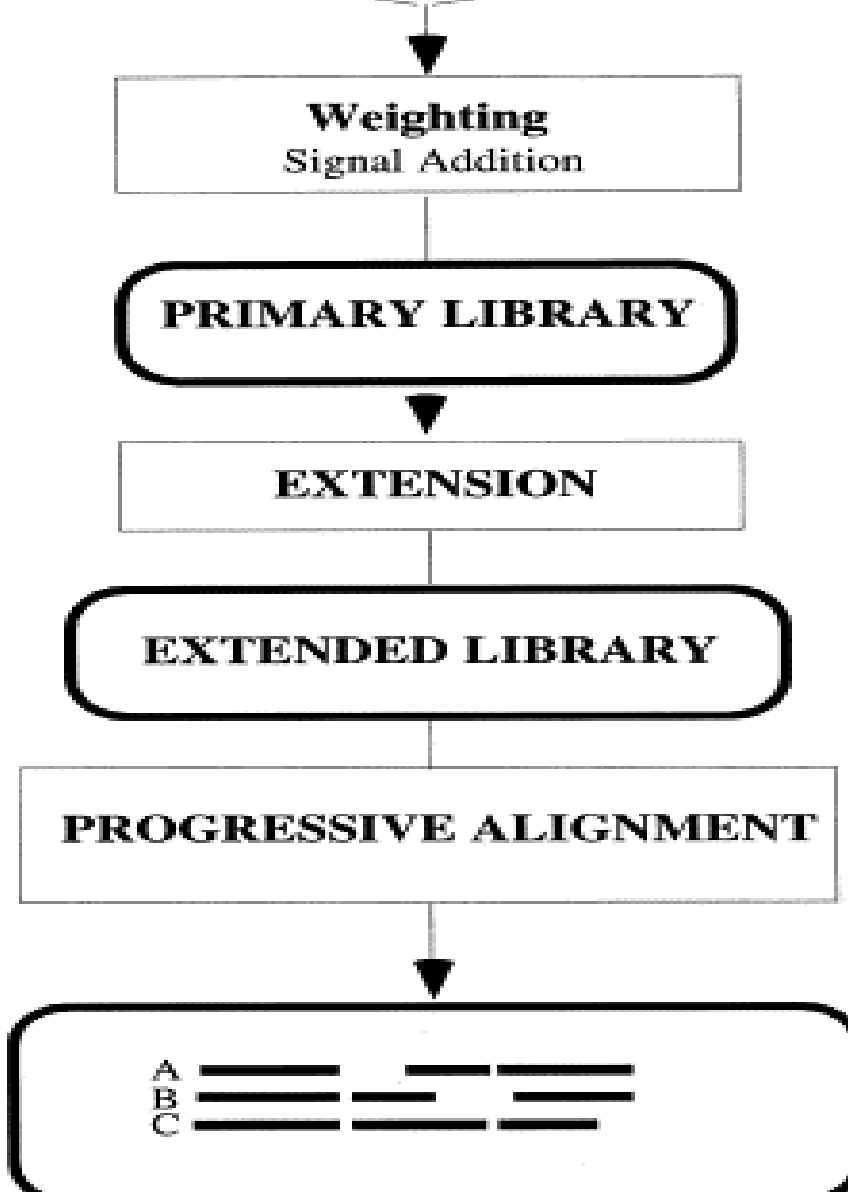
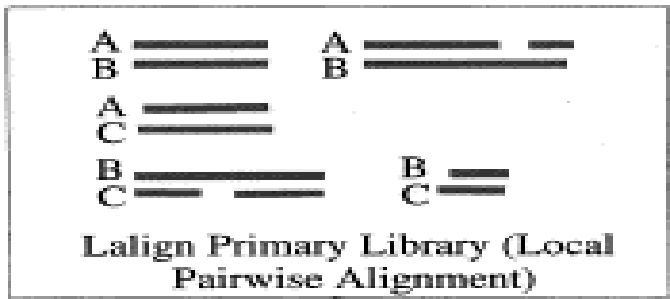
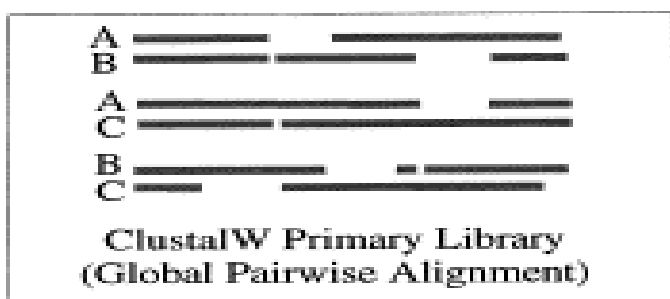
<http://www.monalisamania.com/graphics/art/monalatte.jpg>



Differences

No need for scoring matrices or gap penalties!

Pairwise alignments are guided toward consistency with multiple alignment anyway.

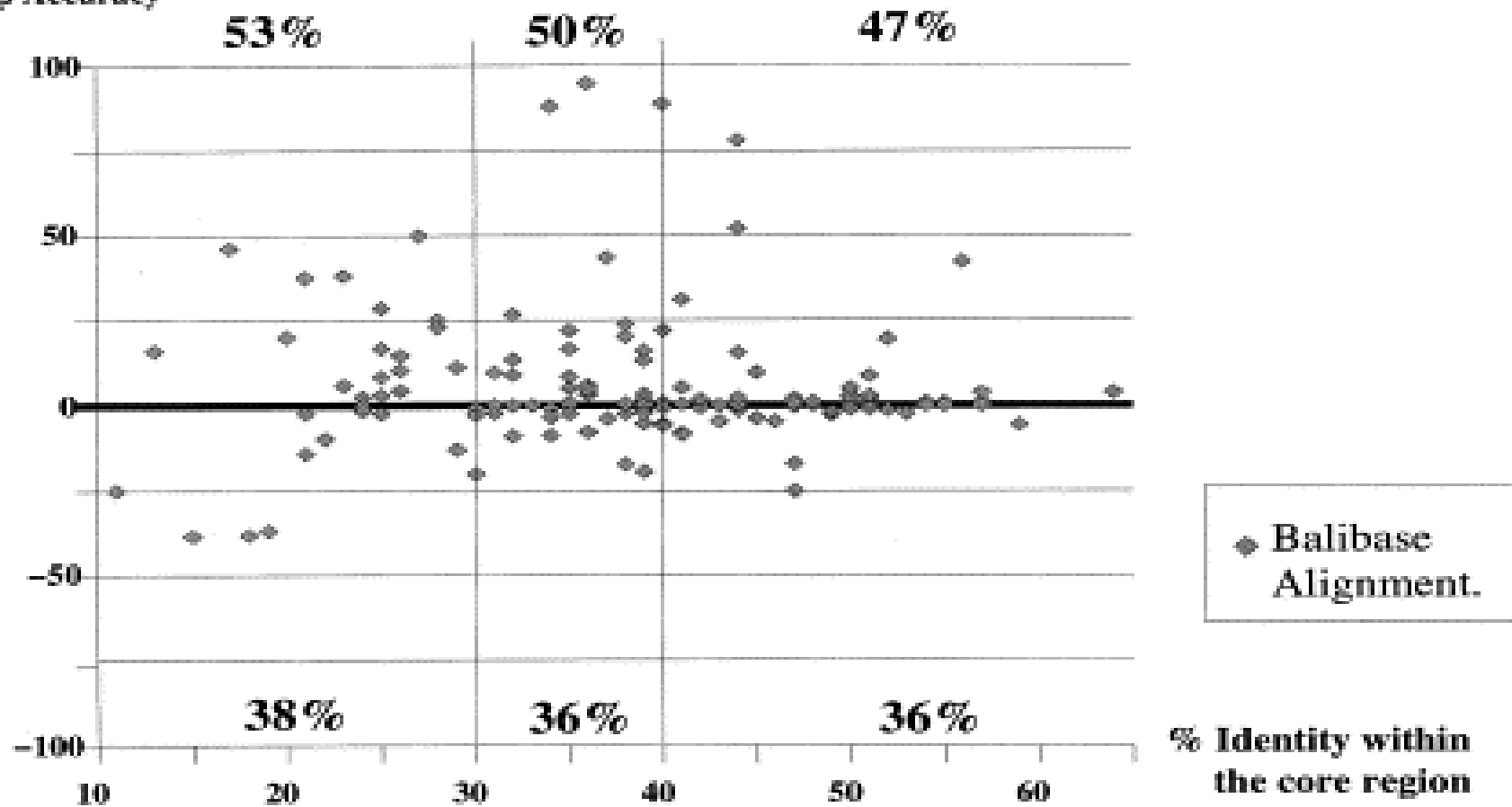


But is it biologically accurate?

Comparison of T-Coffee and Prnp

T-Coffee Accuracy

Prnp Accuracy



Why T-Coffee is worth the effort

Method	Cat1 (81)	Cat2 (23)	Cat3 (4)	Cat4 (12)	Cat5 (11)	Total1 (141)	Total2 (141)	Significance
Dialign	71.0	25.2	35.1	74.7	80.4	61.5	57.3	11.3 ^a
ClustalW	78.5	32.2	42.5	65.7	74.3	66.4	58.6	26.2 ^a
Prnp	78.6	32.5	50.2	51.1	82.7	66.4	59.0	36.9 ^a
T-Coffee	<u>80.7</u>	<u>37.3</u>	<u>52.9</u>	<u>83.2</u>	<u>88.7</u>	<u>72.1</u>	<u>68.7</u>	

Library extension makes all the difference

Particularly useful for Kinases

NBS

```
g11a_orysa  lshfkllkklgsgdigsvylsels---gtesyfamKVMDKas-----
kp68_human  gmdfkeieligsgggfgqvfkakhr---idgktyviKRVKYnn-----
gcn2        --tlkrlnfsgggafgqvvkarna---ldsryyaiKKIRNte-----
st11_yeast  pknwlkgacigsgsfgsvylgmna---htgelnavKQVEIknnnigvpt
kin3_yeast  rseyqvleeigrsgsfgsvrkvihi---ptkklivrKDIKYgh-----
nima_emeni  adkyevlekigcgsfgiirkvkrk---sdgfilcrKEINYik-----
kin1_yeast  lgdwefvetvgagsmgkvklakhr---ytnevcavKIVNRat----kaf
kcc1_yeast  kkkyvfgktlgagtfgvvrqaknt---etgedvavKILIKka-----
ks62_human  psqfellkvlggsfgkvflvkkisgsdarqlyamKVLKKat-----
kpc1_yeast  ldnfvllkvlkggnfgkvilsksk---ntdrlcaiKVLKKdn-----
ypk2_yeast  iddfdllkvigkgsfgkvmqvrkk---dtgkiyaIKALRKay-----
krac_diedi  vadfellnlvgkgsfgkviqvrkk---dtgevyamKVLSKkh-----
kcp2_drome  ltdlrviatlgvggfgrvelvqtn---gdssrsfalkqmkksg-----
kapa_mouse  ldgfdriktlgtgsfgrvmlvkhk---esgnhyamKILDKqk-----
kdca_drome  lenyitravlgngsfgtvmlvrek---sgknyyaaKMMSKed-----
ark1_human  mndfsvhriigrsggfgevygcrkr---dtgkmyamKCLDKkr-----
dmk_human   rddfeilkvigrgafsevavvkmk---qtgqvvyamKIMNKwd-----
dbf2_yeast  nrdfemitqvsggggygqvyllarkk---dtkevcalKILNKKl-----
pim1_human  esqyqvgp1lgsgggfgsvysgirv---sdnlpvaiKHVEKdr-----
```

Efficiency

Runs in quadratic time
 $O(N^2L^2) + O(N^3L) + O(N^3) + O(NL^2)$

2X Slower than ClustalW
(more overhead)

