

# Influence of alignment uncertainty on homology and phylogenetic modeling

Jia-Ming Chang<sup>1</sup>, Cedric Notredame<sup>2</sup>

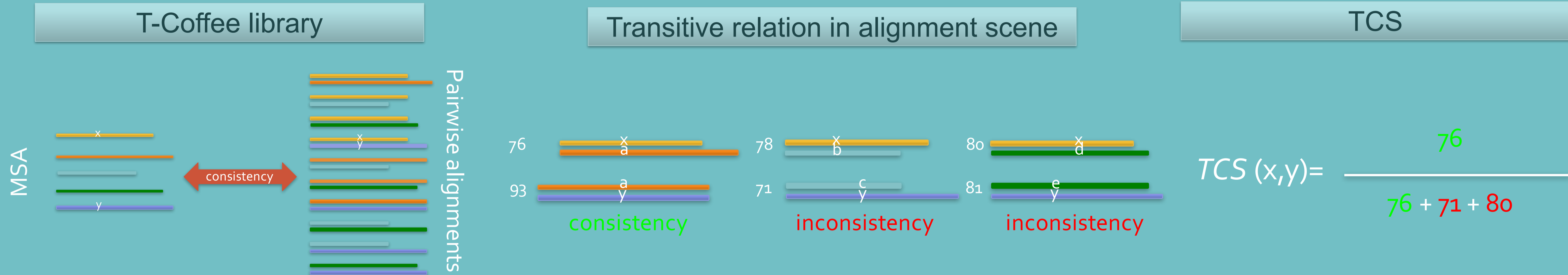
1.Computer Science, National Chengchi University, Taipei, Taiwan; [chang.jiaming@gmail.com](mailto:chang.jiaming@gmail.com)

2.Bioinformatics and Genomics, Centre for Genomic Regulation (CRG), Barcelona, Spain

## Homology extension and sampling approach

Most evolutionary analyses or structure modeling are based upon pre-estimated multiple sequence alignment (MSA) models. From a computational point of view, it is too complex to estimate a correct alignment. Hence, increasing or identifying signal inside sequence alignment has intensified over the last few years. I would like to share two approaches, homology extension and sampling, on this topic.

## Transitive Consistency Score (TCS)



## Local evaluation function on homology and evolutionary modeling

**Residue level**

| Col | row | row | TCS   |
|-----|-----|-----|-------|
| 1   | 1   | 2   | 0.762 |
| 1   | 1   | 3   | 0.748 |
| 1   | 1   | 4   | 0.741 |
| 1   | 2   | 3   | 0.651 |
| 1   | 2   | 4   | 0.677 |
| 1   | 3   | 4   | 0.693 |
| 2   | 1   | 3   | 0.562 |
| 2   | 1   | 4   | 0.632 |
| 2   | 3   | 4   | 0.526 |

**Column level**

| Seq1                              | Seq2                               | Seqn                               | confidence1 | confidence2 |
|-----------------------------------|------------------------------------|------------------------------------|-------------|-------------|
| Seq1 ...SALMLWLSARESIREN...YPD... | Seq2 ...SAYNIYVSPQ----RESA...KD... | Seqn ...SAYNIYVSAQ----RENA...KD... | confidence1 | confidence2 |
| Seq1 ...SALMLWLSARESIREN...YPD... | Seq2 ...SAYNIYVSF----QRESA...KD... | Seqn ...SAYNIYVSA----QRENA...KD... | confidence1 | confidence2 |

**Simulation**

- 16 tips
- 32 tips
- 64 tips
- Yeasts: 853

**RF: average Robinson-Foulds distance respect to Yeast ToL.**  
**TPs: the number of genes whose tree topology is identical with yeast ToL.**

|          | Original | Gblocks relaxed | Gblocks stringent | trimAl gapout | trimAl strictplus | TCS replicate |      |     |      |     |      |     |
|----------|----------|-----------------|-------------------|---------------|-------------------|---------------|------|-----|------|-----|------|-----|
|          | RF       | TPs             | RF                | TPs           | RF                | TPs           | RF   | TPs | RF   | TPs |      |     |
| ClustalW | 0.90     | 643             | 0.99              | 629           | 1.24              | 584           | 0.95 | 628 | 1.31 | 561 | 0.91 | 649 |
| MAFFT    | 0.80     | 665             | 0.83              | 653           | 1.26              | 573           | 0.83 | 657 | 1.28 | 562 | 0.76 | 669 |
| Muscle   | 0.95     | 639             | 0.91              | 646           | 1.26              | 578           | 0.96 | 633 | 1.29 | 559 | 0.84 | 662 |
| PRANK    | 0.79     | 665             | 0.88              | 642           | 1.28              | 565           | 0.84 | 648 | 1.19 | 575 | 0.81 | 662 |
| SATe     | 0.86     | 660             | 0.87              | 650           | 1.28              | 578           | 0.85 | 655 | 1.25 | 567 | 0.79 | 666 |
| AVE      | 0.86     | 654             | 0.896             | 644           | 1.26              | 575           | 0.88 | 644 | 1.26 | 565 | 0.82 | 661 |

## Web site

We show that one can identify the most reliable portions of an MSA, as judged from BALiBASE and PREFAB structure-based reference alignments. We also show how this measure can be used to improve phylogenetic tree reconstruction using both an established simulated data set and a novel empirical yeast data set. For this purpose, we describe a novel lossless alternative to site filtering that involves overweighting the trustworthy columns. We compared TCS with Heads-or-Tails, GUIDANCE, Gblocks, and trimAl and found it to lead to significantly better estimates of structural accuracy and more accurate phylogenetic trees.

TCS <http://tcoffee.org.cat/tcs>

PSI/TM-Coffee <http://tcoffee.org.cat/tmcoffee>

The screenshot shows the TCOFFEE web interface. The 'TCS' section is highlighted, showing the 'Alignment input' and 'Filter options' sections. The 'Filter options' section includes a 'Filter' dropdown set to 'columns', and 'Min' and 'Max' sliders. The 'Output options' section includes 'Filtered' (clustalw, fasta, phylip) and 'Weighted' (score, tcs\_weighted, tcs\_replicate) options. The 'Library Computation' section includes 'Pairwise Methods' (Myriads, MSA, MSA, MSA, MSA) and 'Library Methods' (Msa, Msa, Msa, Msa) options.

The screenshot shows the TCOFFEE web interface. The 'MSA' section is highlighted, showing the 'Result files' and 'Phylogenetic reconstruction' sections. The 'Result files' section includes 'ClustalW', 'MAFFT', and 'Muscle' options. The 'Phylogenetic reconstruction' section includes 'Filter', 'Weighted', and 'Bootstrap' options.

The screenshot shows the TCOFFEE web interface. The 'TM-Coffee alignment result' section is highlighted, showing the 'MSA' section. The 'MSA' section includes a 'TCS' column and a 'Score' column, with a color scale from 0.0 to 1.0.

1. PSI/TM-Coffee: a web server for fast and accurate multiple sequence alignments of regular and transmembrane proteins using homology extension on reduced databases. *Nucleic acids research* 44, W339–343(2016).
2. TCS: a web server for multiple sequence alignment evaluation and phylogenetic reconstruction. *Nucleic acids research* 43, W3–6 (2015).
3. TCS: a new multiple sequence alignment reliability measure to estimate alignment accuracy and improve phylogenetic tree reconstruction. *Molecular biology and evolution* 31, 1625–37 (2014).
4. Accurate multiple sequence alignment of transmembrane proteins with PSI-Coffee. *BMC Bioinformatics* 13, S1 (2012).