# Use of parallelism on MSA tools

Miquel Orobitg Cortada

September 9, 2011

# Index

# Index

# Sequence Alignment

## Definition

A sequence alignment is a way of arranging the sequences of DNA, RNA, or protein to identify regions of similarity that may be a consequence of functional, structural, or evolutionary relationships between the sequences.

# Sequence Alignment
Types

## Types

- Pairwise Alignments.
- Multiple Sequence Alignments (MSA)

# Multiple Sequence Alignment
Global Optimization Methods

## Dynamic programming

- Is a technique to identify the globally optimal alignment solution.
- Exists different algorithms (global, local, glocal).

## Problems

- Computationally difficult to produce the alignment (NP-complete problem).

# Sequence Alignment
Heuristics (1/2)

## Progressive Alignment (PA)

- The alignment is produced by a successive construction of pair-wise alignments.
- Advantages:
  - Good compromise between time spend and accuracy.
- Disadvantages:
  - Heavy dependence on the initial alignment.
  - It is not guaranteed to converge to a global optimum.
- Common methods: T-Coffee and ClustalW

# Sequence Alignment
Heuristics (2/2)

## Iterative methods

- Tries to reduce the errors made in progressive methods.
- Works similarly to progressive methods but repeatedly realign the initial sequences as well as adding new sequences to the growing MSA.
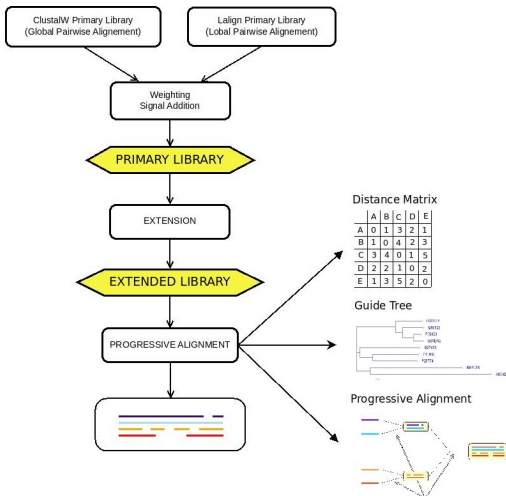- Common methods: Dialign and Muscle

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

T-Coffee
Parallel-TCoffee

# Index

1. Multiple Sequence Alignment

2. MSA Tools
   - T-Coffee
   - Parallel-TCoffee

3. Proposed Solutions
   - Balanced Guide Tree
   - Multiple Trees

4. Future Work

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

T-Coffee
Parallel-TCoffee

# T-Coffee

## T-Coffee

- Is a MSA method that combines the consistency based scoring function COFFEE with the progressive alignment algorithm.
- Advantages:
  - Improvement in the accuracy compared with progressive methods.
  - Reduce the dependency on the initial alignment.

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

T-Coffee
Parallel-TCoffee

# T-Coffee
## Structure

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

T-Coffee
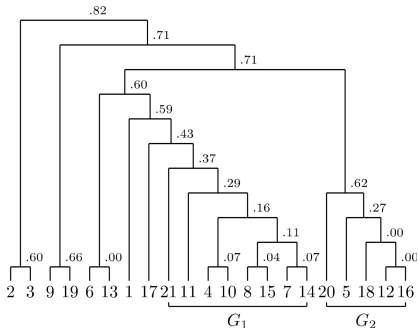Parallel-TCoffee

# T-Coffee
## 1. Library

### Library

- List of pairs of alignments evaluated by a weight that is given by a percentage of identity.
- Generated using different resources.
- Can be extended by transitive properties.
- Used in the progressive alignment.

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

**T-Coffee**
Parallel-TCoffee

# Structure
## 2. Distance Matrix (DM) & Guide Tree (GT)

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

T-Coffee
Parallel-TCoffee

# Structure
## 3. Progressive Alignment

### Progressive Alignment (PA)

- Align the two input sequences using the information of the library.
- The order is determined by the alignment guide tree.

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

**T-Coffee**
Parallel-TCoffee

# Disadvantages

### Library

- $Size = N^2 * L$
- Primary library complexity: $O(N^2 L^2)$
- Extended library complexity: $O(N^3 L^2)$

### Progressive aligment

- Requires n-1 partial multiple alignments using the library. Each alignment can be computation intensive.
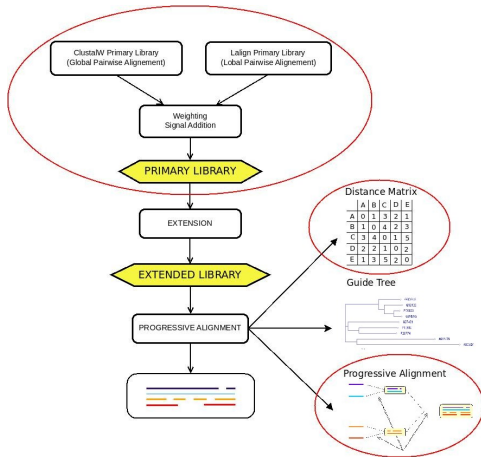- Complexity: $O(NL^2)$

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

T-Coffee
**Parallel-TCoffee**

# Index

1. Multiple Sequence Alignment

2. MSA Tools
   - T-Coffee
   - Parallel-TCoffee

3. Proposed Solutions
   - Balanced Guide Tree
   - Multiple Trees

4. Future Work

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

T-Coffee
**Parallel-TCoffee**

# Parallel-TCoffee
## Parallelization analysis

### 3 parallelization steps

1. Library generation
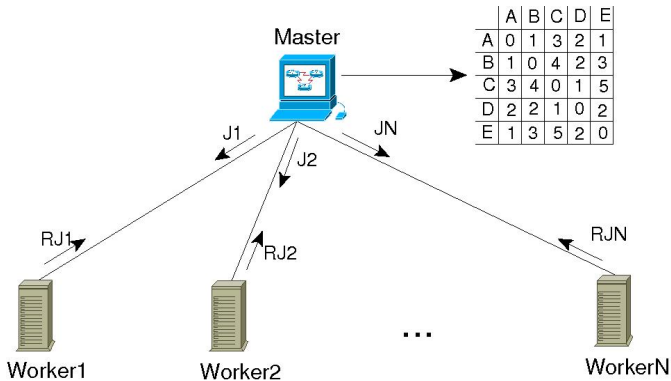2. Distance Matrix computation
3. Progressive alignment

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

T-Coffee
**Parallel-TCoffee**

# Parallel-TCoffee
1. Library Generation

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

T-Coffee
Parallel-TCoffee

# Parallel-TCoffee
## 2. Distance Matrix

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

T-Coffee
**Parallel-TCoffee**

# Parallel-TCoffee
## 3. Progressive Alignment

Multiple Sequence Alignment
**MSA Tools**
Proposed Solutions
Future Work

T-Coffee
**Parallel-TCoffee**

# Parallel-TCoffee Performance
## PF00231 Execution times (554 sequences)



600_Sequences

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# Index

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# NJ Guide Tree analysis
## Critical Path (CP)



**Critical Path Length = 7**

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# NJ Guide Tree analysis
## Maximum Parallelism Degree (MPD)



**Maximum Parallelism Degree = 3**

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# NJ Guide Tree analysis
Examples

| Sequence Set | Nseqs | CP/OCP | MPD |
|---|---|---|---|
| PF00859 | 105 | 37/7 | 19 |
| PF00074 | 442 | 24/9 | 137 |
| PF00349 | 515 | 21/10 | 144 |
| PF01057 | 563 | 84/10 | 87 |
| PF00007 | 731 | 54/10 | 186 |

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# NJ Guide Tree analysis

## Guide Tree problems

- Trees generated with T-Coffee are unbalanced.
- Dependence between iterations.
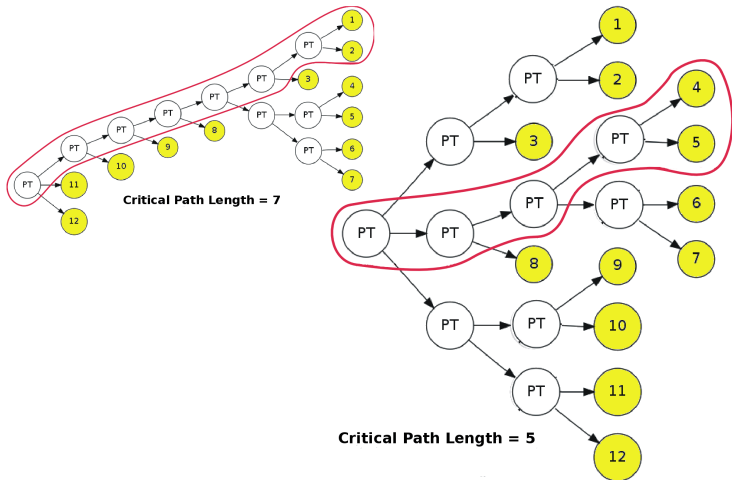- Low degree of parallelism.
- Limited scalability.

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
Multiple Trees

# Balanced Guide Tree heuristic

## Balanced Guide Tree (BGT)

- BGT: Heuristic to balance the nj guide tree maintaining the alignment accuracy.
- Goals:
  - Reduce the number of precedence relations.
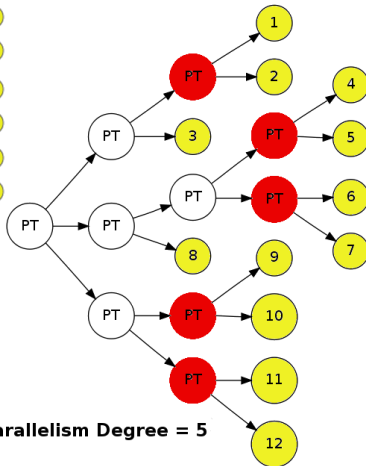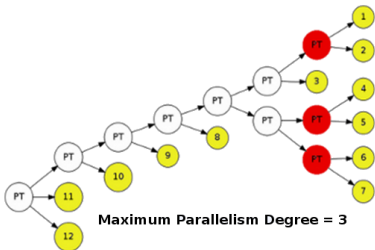  - Decrease the critical path.
  - Increase the parallelism degree.

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# Balanced Guide Tree heuristic
## BGT Tree Features - Critical Path (CP)



Critical Path Length = 7

Critical Path Length = 5

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# Balanced Guide Tree heuristic
## BGT Tree Features - Parallelism Degree (PD)



Maximum Parallelism Degree = 3

Maximum Parallelism Degree = 5

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# BGT Results
Tree features comparison

| Sequence Set | Standard Tree | | BGT Tree | |
|---|---|---|---|---|
| | CP/OCP | MPD | CP/OCP | MPD |
| PF00859 (105) | 37/7 | 19 | 8/7 | 51 |
| PF00074 (442) | 24/9 | 137 | 14/9 | 216 |
| PF00349 (515) | 21/10 | 144 | 15/10 | 249 |
| PF01057 (563) | 84/10 | 87 | 17/10 | 274 |
| PF00007 (731) | 54/10 | 186 | 24/10 | 355 |

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
Multiple Trees

# BGT Results
## Alignment Accuracy

| | T-Coffee | | BGT | |
|---|---|---|---|---|
| **Balibase** | **SP** | **TC** | **SP** | **TC** |
| **Ref 1** | 0.764 | 0.579 | 0.763 | 0.577 |
| **Ref 2** | 0.877 | 0.362 | 0.877 | 0.363 |
| **Ref 3** | 0.785 | 0.393 | 0.783 | 0.390 |
| **Ref 4** | 0.804 | 0.419 | 0,805 | 0.426 |
| **Ref 5** | 0.788 | 0.424 | 0.786 | 0.426 |
| **Ref 6** | 0,807 | 0,393 | 0,807 | 0,402 |
| **Ref 7** | 0,804 | 0.360 | 0.809 | 0.353 |
| **Ref 8** | 0.700 | 0.180 | 0.700 | 0.180 |
| **Ref 9** | 0.742 | 0.481 | 0.742 | 0.482 |
| **Total** | **0.783** | **0.457** | **0.783** | **0.458** |

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# BGT Results
## Alignment Accuracy

|         | T-Coffee | BGT     |
|---------|----------|---------|
| **Prefab** | **Q**    | **Q**   |
| **0 - 15** | 0.421    | 0.422   |
| **15 - 25** | 0.724   | 0.725   |
| **25 - 35** | 0.877   | 0.875   |
| **35 -100** | 0.955   | 0.954   |
| **Total** | **0.711** | **0.711** |

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
Multiple Trees

# BGT Results
## Parallel-TCoffee Performance



749 Sequences - 128 Processors



1318 Sequences - 128 Processors

Total: Total execution time
PA: Progressive alignment execution time

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# BGT Results
## ClustalW-MPI Performance



Total: Total execution time
PA: Progressive alignment execution time

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

**Balanced Guide Tree**
Multiple Trees

# BGT Results
## Parallel-TCoffee Performance - 600 Sequences



600_Sequences

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
**Multiple Trees**

# Index

1. Multiple Sequence Alignment

2. MSA Tools
   - T-Coffee
   - Parallel-TCoffee

3. Proposed Solutions
   - Balanced Guide Tree
   - **Multiple Trees**

4. Future Work

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# Proposal

## Proposal

- Create multiple different subtrees of a guide tree.
- Calculate the alignment of each subtree and its score.
- Use the alignment which gets the best score.

## Objective

- Improve the alignment accuracy.

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
**Multiple Trees**

# Proposal
## Algorithm proposal

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# First Implementation
## Algorithm

Multiple Sequence Alignment
MSA Tools
Proposed Solutions
Future Work

Balanced Guide Tree
Multiple Trees

# Experimentation results
## Prefab - SP Score

- **T-Coffee**:



- **ClustalW**:

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
**Multiple Trees**

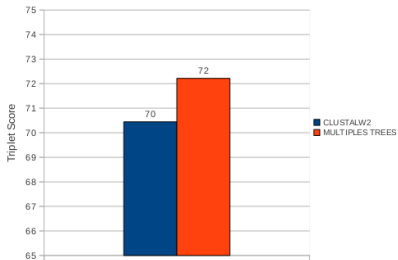# Experimentation results
## Prefab - NORMD Score

● **T-Coffee:**



● **ClustalW:**

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work

Balanced Guide Tree
**Multiple Trees**

# Experimentation results
## Prefab - Triplet Score

- **ClustalW:**

Multiple Sequence Alignment
MSA Tools
**Proposed Solutions**
Future Work
Balanced Guide Tree
**Multiple Trees**

# Disadvantages

## Disadvantages

- To find an evaluation score that defines the best tree.
- The MSA with the best score is not always the best MSA using the benchmark scores.

# Index

# Future Work

## Future Work

- Finish the implementation of the Multiple Trees algorithms.
- Test the performance of the Parallel Multiple Trees solution.
- Publish the Multiple Trees solution.
- Study new parallel algorithms for MSA.

# Questions?