

# Using Tree Edit Distance as an RNA Secondary Structure Similarity Metric

Jon-Michael Deldin

Dept. of Computer Science  
University of Montana  
`jon-michael.deldin@umontana.edu`

2011-05-03

# Outline

- 1 Project
- 2 Trees
- 3 Current Work
- 4 Future Work

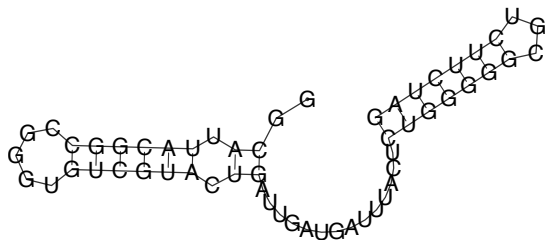
1 Project

2 Trees

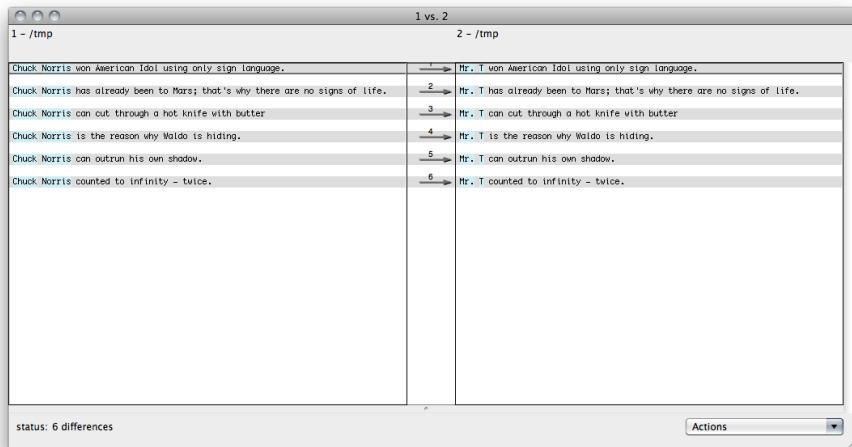
3 Current Work

4 Future Work

- 1 Represent an RNA secondary structure as a tree data structure
- 2 Determine the similarity between two structures using tree edit distance



# Similarity



The screenshot displays a diff window titled "1 vs. 2". The left pane, labeled "1 - /tmp", contains the following text:  
Chuck Norris won American Idol using only sign language.  
Chuck Norris has already been to Mars; that's why there are no signs of life.  
Chuck Norris can cut through a hot knife with butter  
Chuck Norris is the reason why Waldo is hiding.  
Chuck Norris can outrun his own shadow.  
Chuck Norris counted to infinity - twice.

The right pane, labeled "2 - /tmp", contains the following text:  
Mr. T won American Idol using only sign language.  
Mr. T has already been to Mars; that's why there are no signs of life.  
Mr. T can cut through a hot knife with butter  
Mr. T is the reason why Waldo is hiding.  
Mr. T can outrun his own shadow.  
Mr. T counted to infinity - twice.

Arrows labeled 1 through 6 indicate the line-by-line comparison between the two files. The status bar at the bottom left shows "status: 6 differences", and the bottom right has an "Actions" dropdown menu.

# Pipeline



- 100 SELEX aptamers  
(M. Ellenbecker, J.M. Lanchy, & J.S. Lodmell)
- Random RNA sequences

```
>MBE2A  
GGCATTACGGCCGGG TGTCGTACTGATTGATGATTACTCTGGGG GCGTCTTCTAG
```

1 Project

2 Trees

3 Current Work

4 Future Work

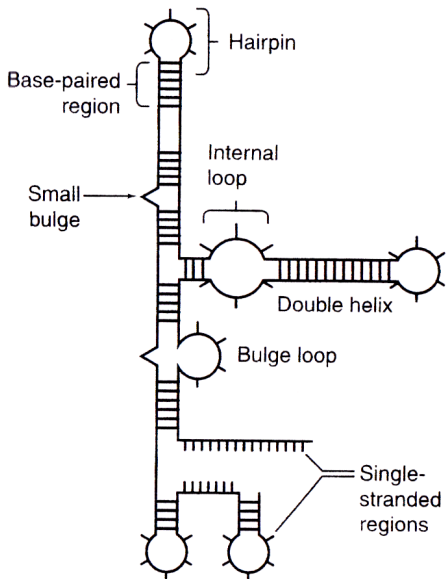


## Definition (Tree Data Structure)

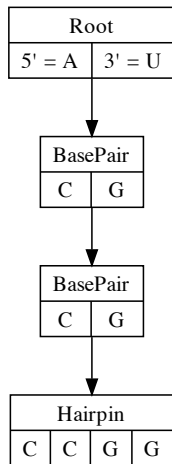
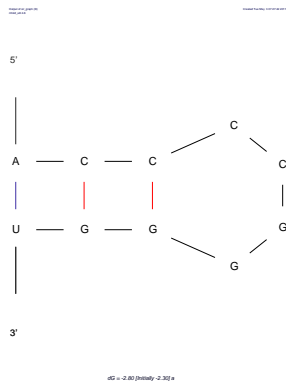
An abstract datatype that represents a hierarchy of objects (e.g., family tree). It is an ordered, directed, acyclic graph.

- Root node
- Traversal order matters (e.g.,  $5' \rightarrow 3'$ )

# Secondary Structure



# Representation



- Pruning
  - Removing nodes from a tree
  - Rules
    - If deleting a nt in a BP connected to a loop, move a nt from out of the loop
    - Cannot delete a single BP in a series of BPs (gap?)
- Grafting
  - Inserting nodes into a tree

- Two different trees  $T_0$  and  $T_1$
- How many operations (insertions/deletions) are needed until  $T_0 = T_1$ ?
- Not easy

1 Project

2 Trees

3 Current Work

4 Future Work

## Implementing a genetic algorithm

- Take an RNA secondary structure
- Permute it (rearrange/insert/delete nodes)
- fitness(): BP score to original structure

## Implementing a genetic algorithm

- Take an RNA secondary structure
- Permute it (rearrange/insert/delete nodes)
- fitness(): BP score to original structure

Chromosome Secondary structure  
Genes Hairpins, bulges, base pairs



- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$

# Genetic Algorithm

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:
  - 1 Get the fitness of all chromosomes

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:
  - 1 Get the fitness of all chromosomes
  - 2 Select 2 parent chromosomes,  $x_0, x_1$

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:
  - 1 Get the fitness of all chromosomes
  - 2 Select 2 parent chromosomes,  $x_0, x_1$
  - 3 offspring  $\leftarrow$  crossover( $x_0, x_1$ ) according to a crossover rate

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:
  - 1 Get the fitness of all chromosomes
  - 2 Select 2 parent chromosomes,  $x_0, x_1$
  - 3 offspring  $\leftarrow$  crossover( $x_0, x_1$ ) according to a crossover rate
  - 4 Mutate the offspring according to a mutation rate

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:
  - 1 Get the fitness of all chromosomes
  - 2 Select 2 parent chromosomes,  $x_0, x_1$
  - 3 offspring  $\leftarrow$  crossover( $x_0, x_1$ ) according to a crossover rate
  - 4 Mutate the offspring according to a mutation rate
  - 5 Add offspring to the population

- 1 Permute ancestor  $n$  times or randomly generate a population of size  $n$
- 2 Until a solution is found:
  - 1 Get the fitness of all chromosomes
  - 2 Select 2 parent chromosomes,  $x_0, x_1$
  - 3 offspring  $\leftarrow$  crossover( $x_0, x_1$ ) according to a crossover rate
  - 4 Mutate the offspring according to a mutation rate
  - 5 Add offspring to the population
- 3 Print solution



- 1 Project
- 2 Trees
- 3 Current Work
- 4 Future Work**

Tree edit distance