# Automated Homology-based Approach for the Identification of Transposable Elements: TESeeker

Ryan C. Kennedy
November 2010

# Transposable Elements (TEs)

- First found and analyzed by Barbara McClintock in 1948
  - Won Nobel Prize in 1983
- TEs are mobile pieces of DNA
- Typically divided into Class I and Class II elements
  - Class I elements are RNA-mediated
  - Class II elements are DNA-mediated
- Example *mariner* Class II TE:

| TA | ACGC...GTAA | GTATCAGCCA...CAAATTACG | TTAC...GCGT | TA |
|---|---|---|---|---|
| Target Site Duplication | Inverted Repeat | Transposase | Inverted Repeat | Target Site Duplication |
| 2 bp | 20-30 bp | ~900 bp | 20-30 bp | 2 bp |

UNIVERSITY OF NOTRE DAME

# Motivation

- ☐ Why study transposable elements?
  - ■ Have been found in all eukaryotic genomes
  - ■ Occupy large portions of genomes
    - ☐ 50% of human genome
    - ☐ 47% of *Aedes aegypti* mosquito genome
  - ■ Can influence genome evolution and gene expression

  - ■ Mouse Genome Sequencing Consortium:
  
    "The single most prevalent feature of mammalian genomes is their repetitive sequences, most of which are interspersed repeats representing 'fossils' of transposable elements. *Transposable elements are a principal force in reshaping the genome, and their fossils thus provide powerful reporters for measuring evolutionary forces acting on the genome.*"

R.H. Waterston, et al. Initial sequencing and comparative analysis of the mouse genome. Nature, 420:520-562, December 2002.

UNIVERSITY OF NOTRE DAME

# TE Discovery Techniques

- Bergman and Quesneville categorize TE discovery into four categories:
  1) Comparative Genomic Methods
     - Perform multiple sequence alignment of related genomes and look for large changes amongst them
     - Good for finding new TE families, but relies on readily available, properly sequenced, related genomes
  2) *De novo*
     - Detect similar sequences found throughout the genome and cluster
     - Can discover new TE families, but often difficult to distinguish closely related TEs

C. M. Bergman and H. Quesneville. Discovering and detecting transposable elements in genome sequences. Briefings in Bioinformatics, 8(6):382-392, 2007.

UNIVERSITY OF
NOTRE DAME

# TE Discovery Techniques

3) Structure-based

   ☐ Use TE structural data, such as inverted repeats, to find TEs

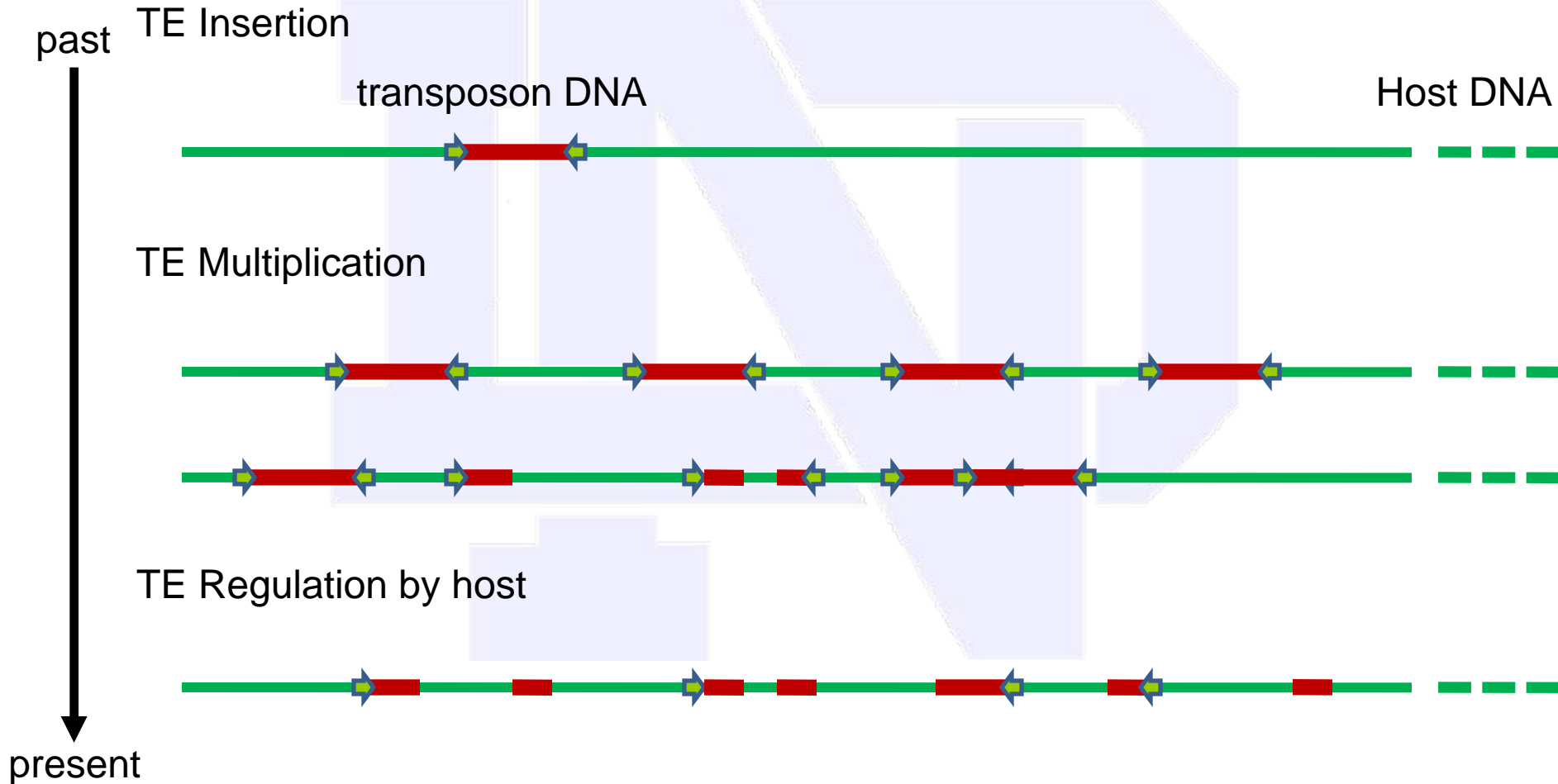   ☐ Works well for characterized TEs, but does not locate degraded TEs or TEs with non-distinct structures

4) Homology-based (our approach)

   ☐ Use known TEs as seeds to search in novel genomes

   ☐ Can discover new TE families, but requires additional verification

UNIVERSITY OF
NOTRE DAME

# Challenges in Locating TEs

☐ Although present in all eukaryotic genomes, difficult to annotate

   ■ Varying structural characteristics

   ■ Mobile nature often leads to copies within copies

   ■ TEs often are very degraded

UNIVERSITY OF
NOTRE DAME

# Class II TE Evolution

past

TE Insertion

transposon DNA                                                    Host DNA

TE Multiplication

TE Regulation by host

present

UNIVERSITY OF
NOTRE DAME

# Manual Approach

- Developed and utilized during TE search on very different genome projects:
  - *Pediculus humanus humanus* (body louse)
    - Comprehensive search for all TE families
  - *Culex quinquefasciatus* (mosquito)
    - Search for non-LTR TEs
- Homology-based
  - Assembled representative TE library of high-quality TEs
    - Intact open reading frames
  - Results appear in TE sections of respective genome papers

UNIVERSITY OF
NOTRE DAME

# Manual Approach



Representative transposases

Genome

tblastn

Combine, add flanks, extract

BLAST hits

Assemble in DNASTAR SeqMan II

9

Consensus TE

# *P. humanus humanus* Results

| Class I | Family | Element | Length (bp) | Full-length Copies | Partial Hits | Density |
|---|---|---|---|---|---|---|
| Non-LTR | SART | *Hope-like* | 4655 | 1 | 522 | 0.18% |
| | R4 | *Dong-like* | 5266 | 4 | 1739 | 0.45% |
| LTR | Ty3/gypsy | *Mdg1* | 5395 | 2 | 976 | 0.28% |
| Class II | Family | Element | Length (bp) | Full-length Copies | Copies | Density |
| | Mariner/Tc1 | *mariner* | 1276 | 24 | 216 | 0.09% |
| TOTAL | | | | | | 1.0% |

E.F. Kirkness et al., "Genome sequences of the human body louse and its primary endosymbiont provide insights into the permanent parasitic lifestyle." *Proceedings of the National Academy of Sciences*, 107(27):12168-12173, July 2010.

UNIVERSITY OF
NOTRE DAME

# *C. quinquefasciatus* Results

| Class I | Family | Full-length Copies | Partial Hits | Density |
|---------|--------|--------------------|--------------|---------|
| Non-LTR | CR1 | 31 | 973 | 0.28% |
| | I | 11 | 63 | 0.02% |
| | Jockey | 14 | 5028 | 1.77% |
| | L1 | 57 | 662 | 0.15% |
| | L2 | 9 | 1416 | 0.61% |
| | Loa | 9 | 184 | 0.09% |
| | Loner | 2 | 127 | 0.12% |
| | Outcast | 4 | 15 | 0.00% |
| | R1 | 32 | 250 | 0.14% |
| | RTE | 8 | 892 | 0.38% |
| | Unclassified LINE | 32 | 11,117 | 0.88% |
| **TOTAL** | | | | **4.44%** |

P. Arensburger et al., "Sequence of *Culex quinquefasciatus* Establishes a Platform for vector Mosquito Comparative Genomics." *Science*, 330(6000):86-88, October 2010.

UNIVERSITY OF
NOTRE DAME

# Manual Approach



Representative transposases

Genome

tblastn

BLAST hits

Combine, add flanks, extract

major burden

Assemble in DNASTAR SeqMan II

12

Consensus TE

# DNASTAR SeqMan II

☐ Manually trimming hits and additional processing is time consuming

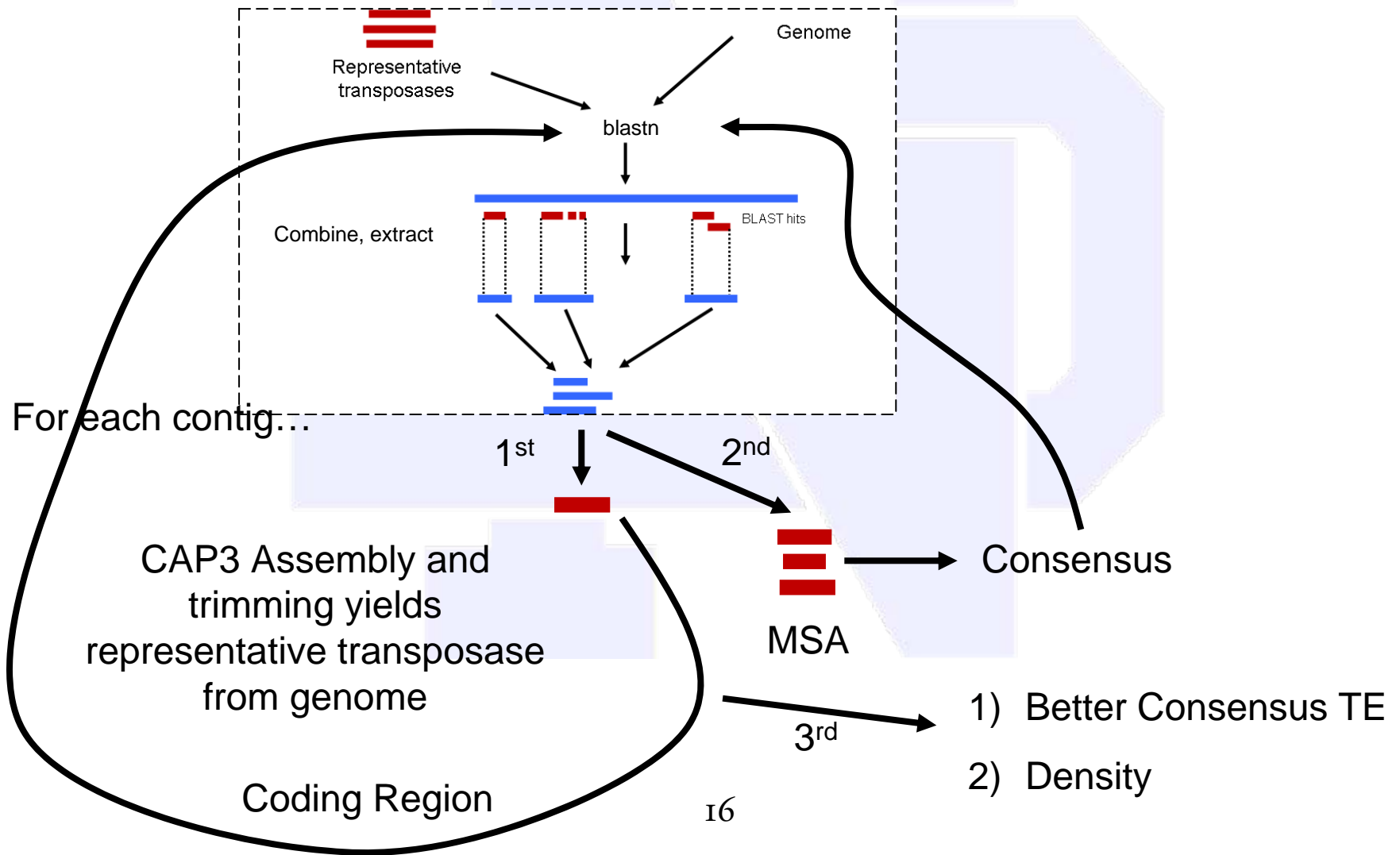☐ Can only assemble limited number of sequences at a time

# Automated Approach

☐ Homology-based

☐ Replace DNASTAR SeqMan II and manual analysis with other tools

■ CAP3, Clustal, various scripts

☐ Iterative- repeat steps if necessary

UNIVERSITY OF
NOTRE DAME

# Automated Approach Steps

1. Identify transposase in target genome
2. Find copies in target genome with flanks
3. Generate consensus from multiple sequence alignment of copies
4. Use consensus to identify TE

☐ Output: putative high-quality consensus TE which can in turn be used locate instances within the genome

☐ Runs in a matter of minutes/hours
  ■ Dependent on genome size, size of representative TEs, and richness of TEs in the genome

☐ Runs via web interface or via automated scripts

15

UNIVERSITY OF NOTRE DAME

# Automated Approach



Representative transposases

Genome

blastn

Combine, extract

BLAST hits

For each contig…

1st    2nd

CAP3 Assembly and trimming yields representative transposase from genome

Consensus

MSA

Coding Region

3rd    1) Better Consensus TE

2) Density

16

# Step 1

- Identify transposase(s) in target genome
  - *tblastn* representative transposases against genome
    - Parse BLAST file with the following parameters:
      - combine threshold: maximum distance sequences can be apart to join as a single hit
      - minimum length percentage: must be at least this percentage of query sequence to be considered
      - e-value cutoff: ignore everything worse than this value, typically 1E-20
      - flank size: amount of extra sequence to add to each end of hit (0)
  - Extract genomic sequences from above and iteratively assemble with CAP3
    - With CAP3, specify quality window size and threshold, as well as combine threshold

transposase(s) within genome

UNIVERSITY OF
NOTRE DAME

# Step 2

- Find copies in target genome with flanks
  - *blastn* transposase(s) against genome
    - Parse BLAST file with the following parameters:
      - combine threshold: maximum distance sequences can be apart to join as a single hit
      - minimum length percentage: must be at least this percentage of query sequence to be considered
      - e-value cutoff: ignore everything worse than this value, typically 1E-20
      - flank size: amount of extra sequence to add to each end of hit
  - Extract genomic sequences from above

Copies within genome with flanks

UNIVERSITY OF NOTRE DAME

# Step 3

☐ Obtain consensus from multiple sequence alignment (MSA) of copies

- Perform MSA on sequences
- Generate consensus from MSA
  - ☐ Can specify percentage of nucleotides that must be common amongst sequences to count in consensus

↓

Putative Consensus

UNIVERSITY OF
NOTRE DAME

# Step 4

- ☐ Use consensus to identify proper TE
  - ■ *blastn* representative transposases against genome
    - ☐ Parse BLAST file with the following parameters:
      - ▪ combine threshold: maximum distance sequences can be apart to join as a single hit
      - ▪ minimum length percentage: must be at least this percentage of query sequence to be considered
      - ▪ e-value cutoff: ignore everything worse than this value, typically 1E-20
      - ▪ flank size: amount of extra sequence to add to each end of hit
  - ■ Extract genomic sequences from above and iteratively assemble with CAP3
    - ☐ With CAP3, specify quality window size and threshold, as well as combine threshold

Consensus TE → Density

UNIVERSITY OF
NOTRE DAME

# Automated Approach Schematic

# Validation Strategy

1. Initially evaluated automated approach on *P. humanus humanus* and *C. quinquefasciatus*
   - Validate against high-quality manually verified annotation
   - Identify default starting parameters
2. Check automated results versus published results
3. Genomes in General:
   - Translate consensus TE sequences
     - Identify open reading frame
       - *blastp* against *non-redundant protein (nr)* database at NCBI and check for conserved domains/hits
   - Can check for structural signatures

UNIVERSITY OF
NOTRE DAME

# Validation (1): *P. humanus humanus mariner*

☐ Full *mariner* element identified following Step 4

☐ Validated against manual effort

  ■ TSDs; 14 bp terminal inverted repeats (TIRs); well-trimmed

UNIVERSITY OF NOTRE DAME

# Validation (2): *Anopheles gambiae* PEST

- ❏ *P* elements (Class II)
  - ■ Sarkar et al. (2003) identified 6 distinct elements
  - ■ Oliveira de Carvalho et al. (2004) identified 4 additional elements
  - ■ Quesneville et al. (2006) identified 9 elements at least 30% divergent at nucleotide level
  - ■ **Total:** 12 elements at least 30% divergent at nucleotide level

  - ■ Automated Approach
    - ❏ Identified 11/12 elements + 2 partial hits
    - ❏ Captured TIRs where previously described

UNIVERSITY OF
NOTRE DAME

# Validation (3)

- Searched for *mariner* in a number of genomes
  - In agreement where previously reported
    - Human, frog, chicken
  - In agreement where not reported
    - Dog, cat, horse

  - Possible discovery
    - *Drosophila melanogaster* putative *mariner*
      - 1061 bp element has TIRs
      - 26 bp TIRs
      - no apparent TSDs
      - Single full-length copy, as well as several partial hits
      - Transposase is most similar to that of *Chymomyza amoena*, 77% identical at the amino acid level
      - Searches for this element in existing TE annotations for D. melanogaster produced no hits

UNIVERSITY OF
NOTRE DAME

# Implementation

- Approach implemented as TESeeker
  - VirtualBox virtual appliance
    - Cross-platform
    - Completely configured, no need to install scripts
      - Provide only genome FASTA file
      - Optionally provide additional library files
    - Local web interface

- http://www.nd.edu/~teseeker
  - Virtual appliance
  - Documentation
  - TE Library

UNIVERSITY OF
NOTRE DAME

# TESeeker Desktop

# TESeeker Desktop

# TESeeker Desktop

# TESeeker Desktop

# TESeeker Desktop

# TESeeker Desktop

# TESeeker Walkthrough

☐ Identify *mariner* element in *P. humanus humanus:*

- ■ Start with default parameters
- ■ Make sure genome file and library file are present
- ■ Start search

UNIVERSITY OF
NOTRE DAME

# TESeeker Desktop

# TESeeker Local Web Interface

# TESeeker Status

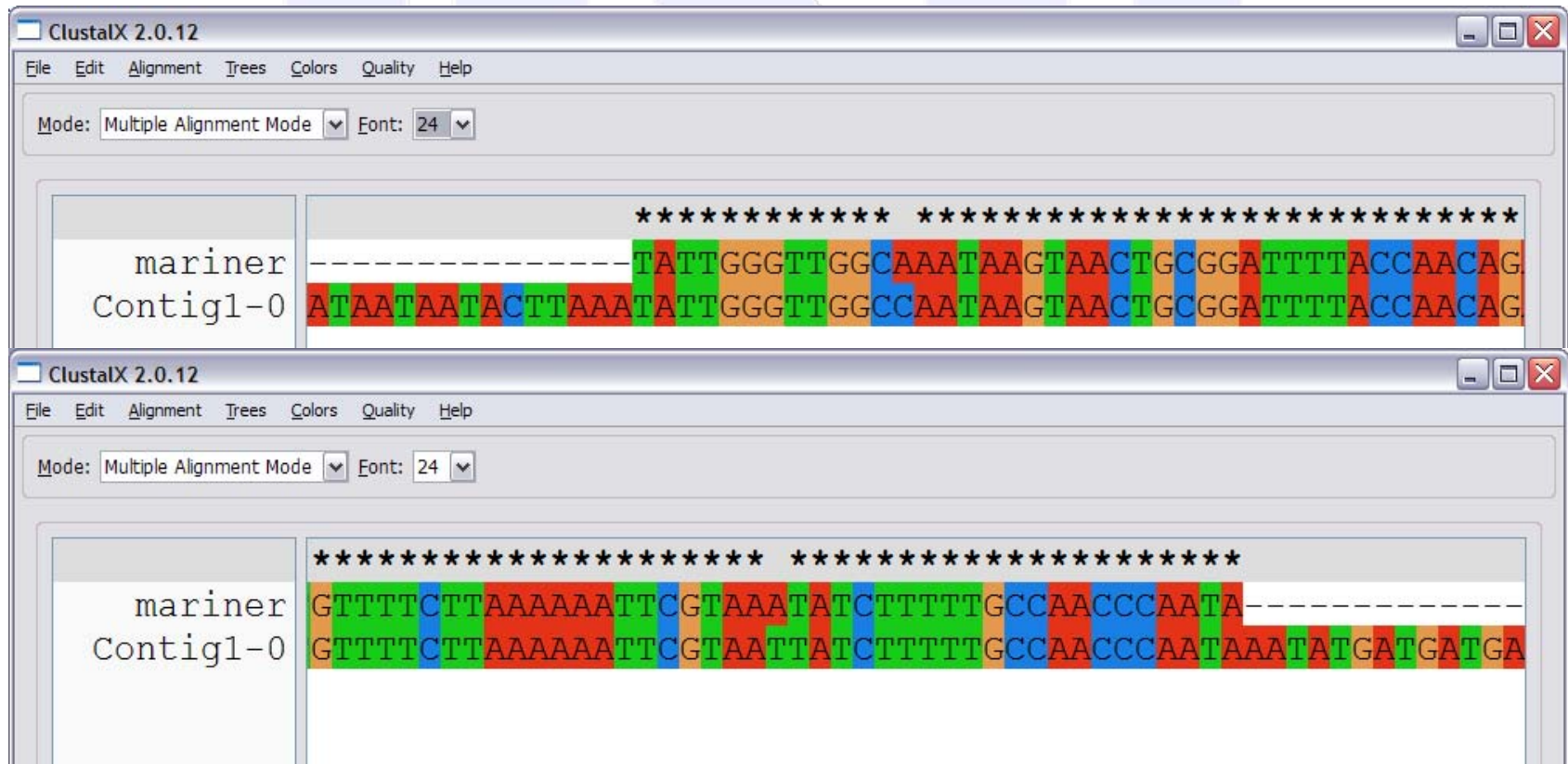# TESeeker Status

# TESeeker Results

# TESeeker Results

# TESeeker Results

☐ Alignment
  ■ TESeeker top result with default parameters
  ■ 99% identity with manually annotated *mariner*

# TE Identification Summary

- Developed and automated a homology-based approach to identify TEs
  - Tedious and time-consuming task now automated
    - From months to hours or days
  - Output: high-quality consensus TEs
    - Can be used to determine instances in genome (density)
- Implemented as TESeeker
  - Distributed as a virtual appliance
    - All tools and scripts
  - Web interface
  - Distributed with high-quality library of representative coding regions from major TE families
- Approach contributed to multiple genome annotation projects
  - Sequences available in TEfam database
  - Most rigorously tested in arthropod genomes

UNIVERSITY OF
NOTRE DAME

# Acknowledgments

UNIVERSITY OF
NOTRE DAME

# Questions or Comments?