# ParsEvalMPI: comparison of gene structure annotations in parallel

Daniel S. Standage

Bioinformatics and Computational Biology
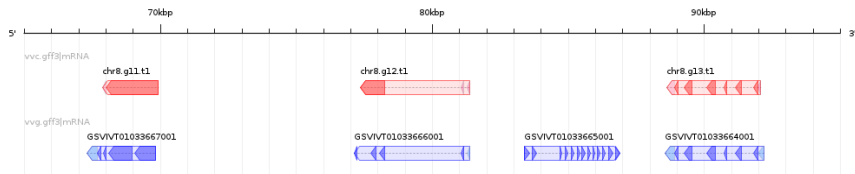
May 3, 2011

# Gene prediction (annotation)

## Input

String representing DNA sequence

## Output

Regions (coordinates) of sequence that correspond to genes and their
structural components

# Comparing annotations

# Previous implementation

## ParsEval

- written in Perl

# Previous implementation

## ParsEval

- written in Perl
- external dependencies (platform-specific)

# Previous implementation

## ParsEval

- written in Perl
- external dependencies (platform-specific)
- significant memory demands

# Previous implementation

## ParsEval

- written in Perl
- external dependencies (platform-specific)
- significant memory demands
- serial execution: run time in minutes to hours

# New implementation

## ParsEvalMPI

- written in C

# New implementation

## ParsEvalMPI

- written in C
- external dependencies (cross-platform)

# New implementation

## ParsEvalMPI

- written in C
- external dependencies (cross-platform)
- reduced memory demands

# New implementation

## ParsEvalMPI

- written in C
- external dependencies (cross-platform)
- reduced memory demands
- parallel execution: run time in ...

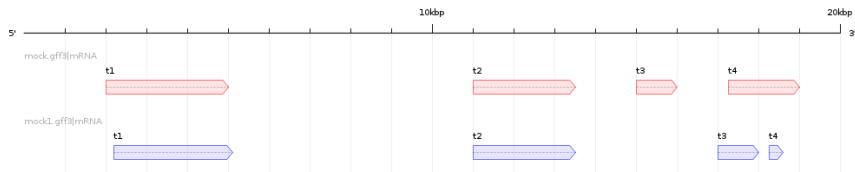# Program overview

# Program overview

- delegation

# Program overview
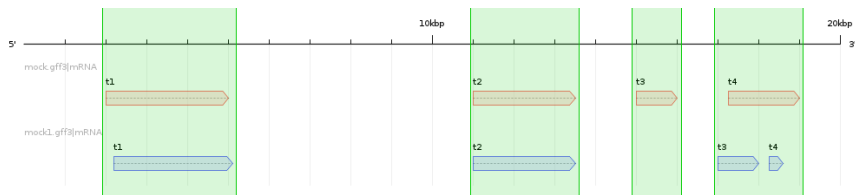
- delegation
- local analysis

# Program overview

- delegation
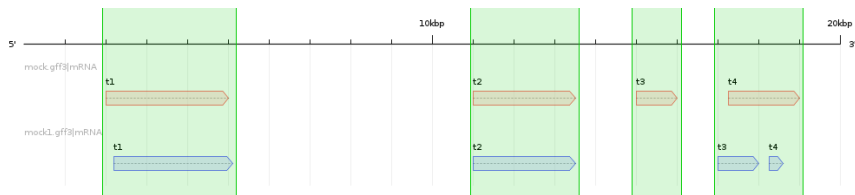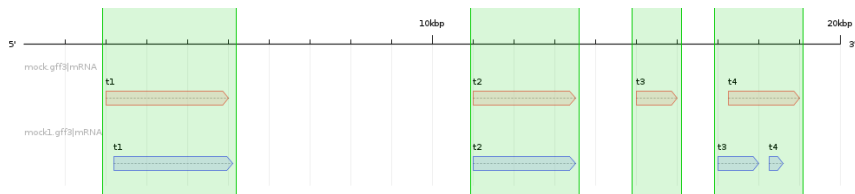- local analysis
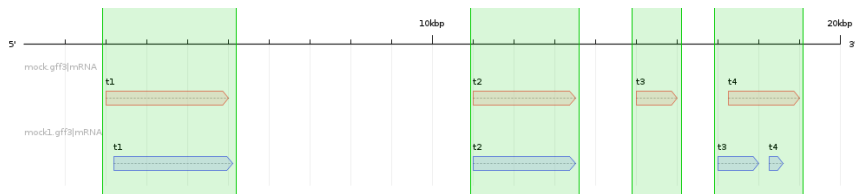- global analysis

# Delegation

# Delegation

# Delegation



1. all data on all processors

# Delegation



1. all data on all processors
2. even distribution of DNA

# Delegation



1. all data on all processors
2. even distribution of DNA
3. even distribution of genes

# Local and global analysis

- Local analysis

- Global analysis

# Local and global analysis

- Local analysis
    - generate local model vector

- Global analysis

# Local and global analysis

- Local analysis
    - generate local model vector
    - send local vector to global vector on root processor
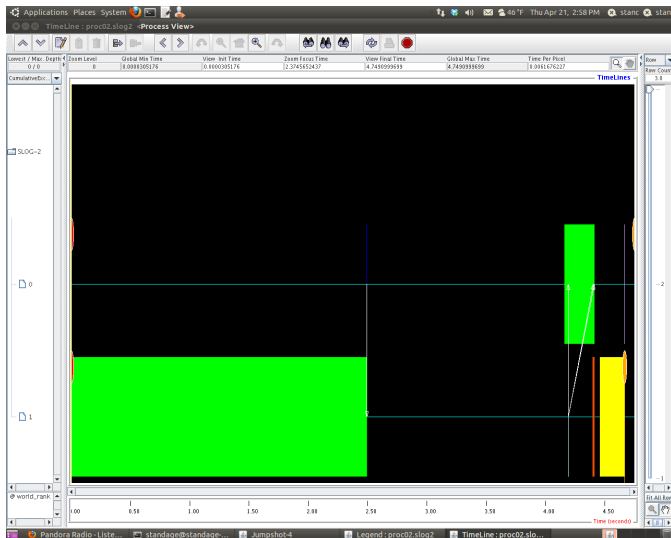
- Global analysis

# Local and global analysis

- Local analysis
  - generate local model vector
  - send local vector to global vector on root processor
  - analyze local vector, print scores
- Global analysis
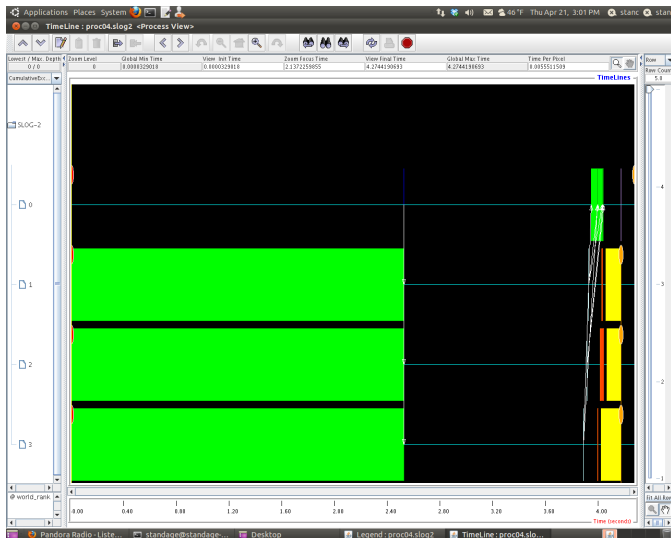
# Local and global analysis

- Local analysis
  - generate local model vector
  - send local vector to global vector on root processor
  - analyze local vector, print scores
- Global analysis
  - receive local vectors from each processor

# Local and global analysis

- Local analysis
  - generate local model vector
  - send local vector to global vector on root processor
  - analyze local vector, print scores
- Global analysis
  - receive local vectors from each processor
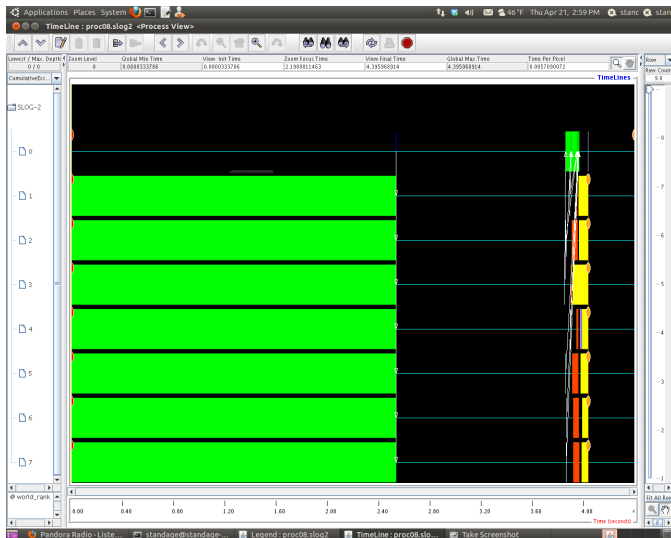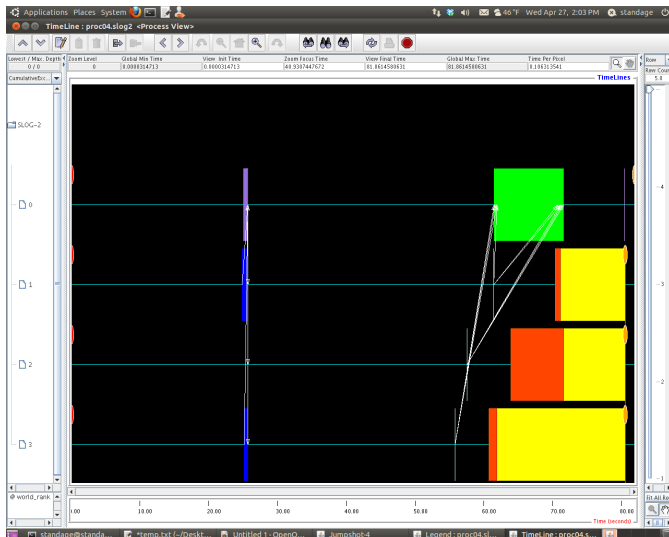  - analyze combined global vector, print scores

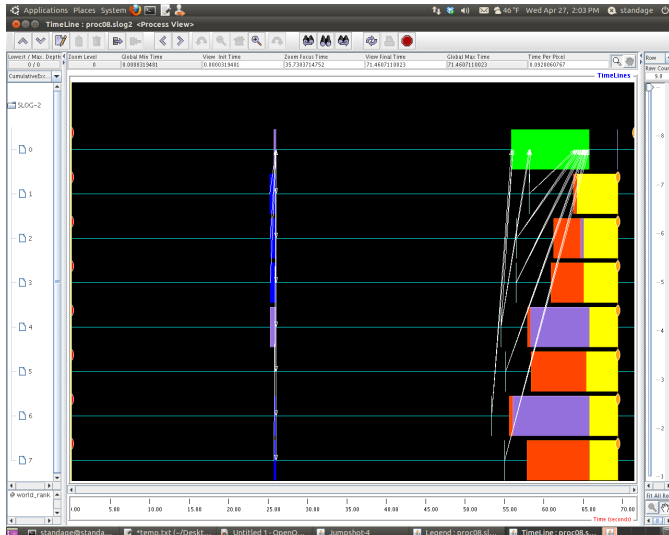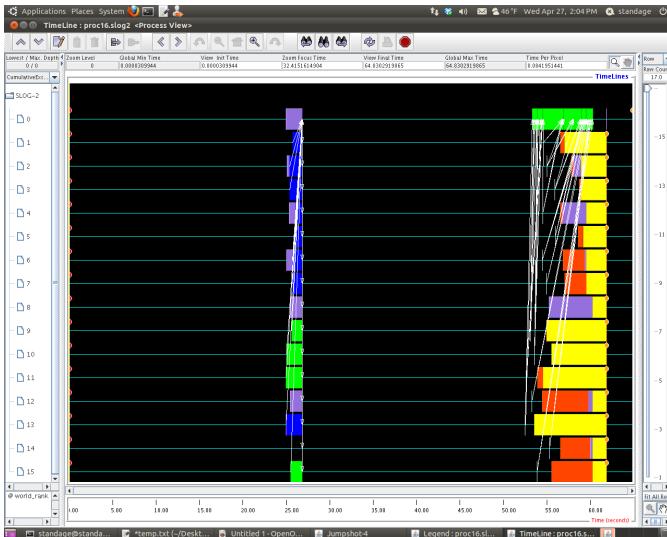# Load balancing

# Load balancing

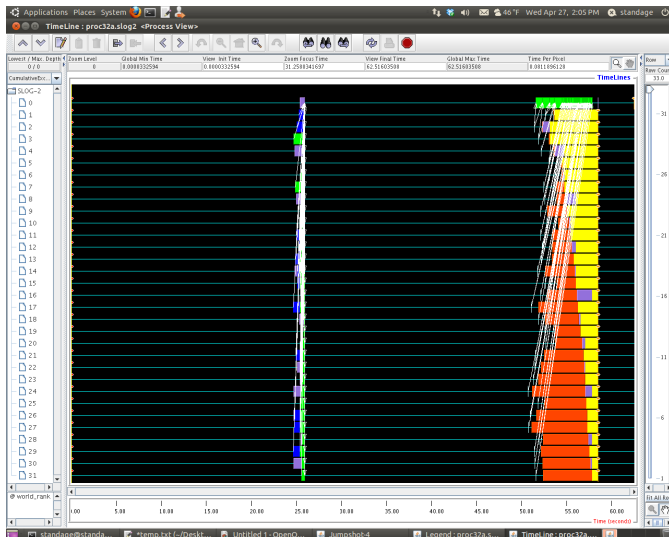# Load balancing

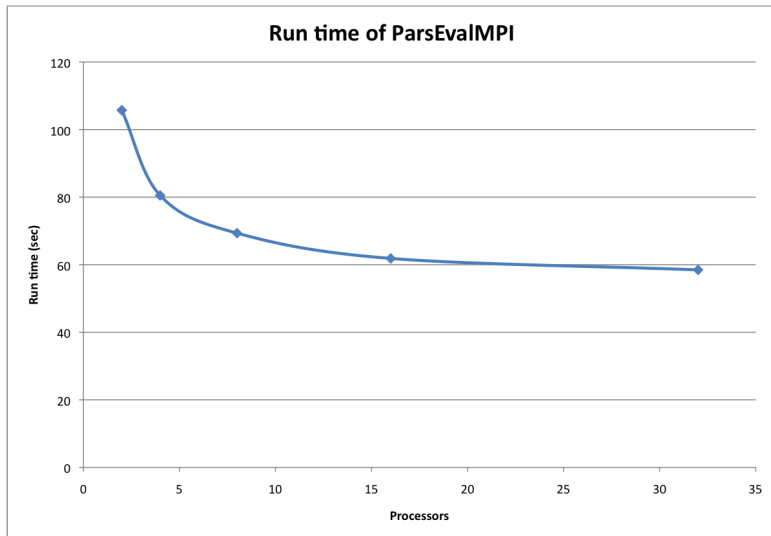# Load balancing

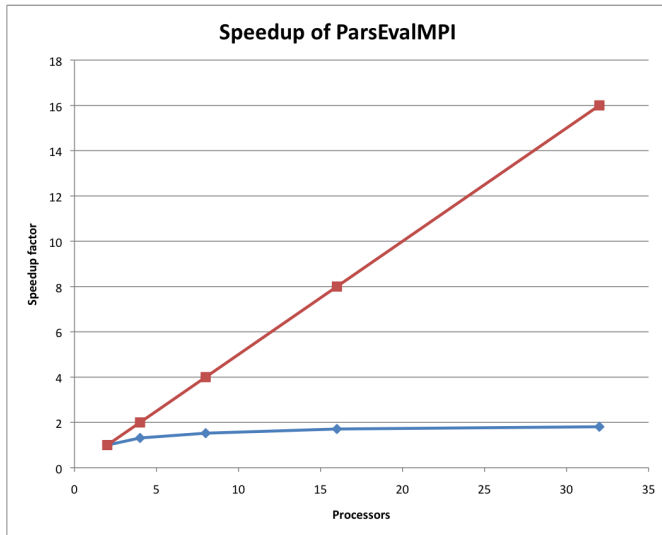# Load balancing

# Load balancing

# Load balancing

# Scalability

# Scalability

# Serial optimization

# Serial optimization

- Good
  - native data types
  - static arrays

# Serial optimization

- Good
  - native data types
  - static arrays
- Bad
  - pointers, dynamic arrays
  - dynamic data structures

# Serial optimization

- Good
  - native data types
  - static arrays
- Bad
  - pointers, dynamic arrays
  - dynamic data structures
- Ugly
  - copying data

# Conclusions

# Conclusions

- Excellent data distribution, load balancing

# Conclusions

- Excellent data distribution, load balancing
- Very poor scaling properties

# Conclusions

- Excellent data distribution, load balancing
- Very poor scaling properties
    - maximum scaling factor of 2?!?!

# Conclusions

- Excellent data distribution, load balancing
- Very poor scaling properties
  - maximum scaling factor of 2?!?!
  - perhaps try OpenMP

# Conclusions

- Excellent data distribution, load balancing
- Very poor scaling properties
  - maximum scaling factor of 2?!?!
  - perhaps try OpenMP
- Significant improvement