

# A Comparative Evaluation of XML Difference Algorithms with Genomic Data

*Cornelia Hedeler (Connie)* and Norman W. Paton School of Computer Science, The University of Manchester, UK

Combining the strengths of UMIST and The Victoria University of Manchester

## Overview

- Motivation
- Problem description
- XML difference algorithms
- Experiments involving controlled modifications
- Evaluation involving real genomic data
- Conclusions

#### Motivation – Genome sequence data

- Genome sequence data for increasing number of genomes
- Undergoes regular reannotation (and re-assembly)
- Annotation is done automatically using computational analysis
- Very little information on the changes between releases
- But previous versions are available (Ensembl, EMBL)

Number of species sequenced (source: NCBI)

Organism	Number of	
group	genomes	
Viruses	1952	
Eukaryota	1565	
Bacteria	956	

#### Motivation – From a biologist's perspective



## Motivation – Information on changes (1)

#### Saccharomyces Genome Database (SGD)

2006-01-11	YCL058W-A
	YCLD58W-A was originally added per Brachat et al., based on conservation with Ashbya gossypii. Kellis et al. also predicted a protein in the same frame, but with a 19 amino acid N-terminal deletion relative to the ORF predicted by Brachat et al. SGD has changed the annotation to that predicted by Kellis et al., based on conservation of that start site in sensu stricto strains of Saccharomyces.
	Kellis M, et al. (2003) Sequencing and comparison of yeast species to identify genes and regulatory elements. <i>Nature</i> 423(6937):241-54          Sap <sup>©</sup> grated       Access       Access       Web         Full Text       Full Text       Supplement
	Brachat S, et al. (2003) Reinvestigation of the Saccharomyces cerevisiae genome annotation by comparison to the genome of a related fungus: Ashbya gossypii. Genome Biol 4(7):R45
2005-11-14	YCLD48W-A
	Based on genome sequence comparisons among six Saccharomyces species, Cliften et al. 2003 suggested that a new ORF, YCL048W-A, be added to the S. cerevisiae genome annotation. This suggestion was later confirmed by Zhang and Dietrich 2005 using high-throughput identification of transcription start sites.
	Cliften P, et al. (2003) Finding functional features in Saccharomyces genomes by phylogenetic footprinting. <i>Science</i> 301(5629):71-6 September Publicat Septement
	Zhang Z and Dietrich FS (2005) Mapping of transcription start sites in Saccharomyces cerevisiae using 5' SAGE. Nucleic Acids Res 33(9):2838-51         SGD <sup>C</sup> urated Full Text       Access Full Text       Access Full Text       Access Full Text       Access Full Text       Supplement
1998-07-24	YCR055C, YCR057C, YCR058C
	ORFs YCR055C and YCR058C were merged with existing ORF YCR057C. The coding region of YCR057C had been 1320 nt long, but is now 2772 nt in length.
1998-07-24	YCR029C

ORF YCR029C has been deleted from the genome annotation due to sequence correction.

#### http://www.yeastgenome.org/cache/genomeSnapshot.html#ChrSeqAnnotUpdates

## Motivation – Information on changes (2)

#### Ensembl

#### What's New in Ensembl 49

#### Mus musculus News

Mouse transcripts updated

Several mouse transcripts that were "culled" in Ensembl 47 have been reinstated. Also, a bug in the pipeline that caused ditags to be ignored has been fixed, which has resulted in some shorter UTRs.

#### Minor updates

A number of minor changes have been made to the core databases. Read more...

#### Variation updates

New variation databases have been released for human, mouse, chicken and zebrafish (dbSNP 128) and rat (dbSNP 126). Transcript variation has been rerun for dog, platypus, chimp, and tetraodon due to gene set changes.

#### Vega updates

Ensembl Vega has a new Mouse database, and Vega Human no longer shown external Vega annotation.

For the convenience of our users, we have created a web-based Wizard' which enables you to convert a GFF file from Mouse assembly m36 to m37.

#### http://www.ensembl.org

#### **Fungal Genome Initiative**

OLD LOCUS, NEW LOCUS, UPDATE TYPE, UPDATE DETAIL NCU00001.2,NCU00001.3,UNCHANGED, NCU00002.2,,DELETE, NCU00003.2, NCU00003.3, UNCHANGED, NCU00004.2, NCU00004.3, UNCHANGED, NCU00005.2, NCU00005.3, CHANGE, UTR CHANGED NCU00006.2, NCU00006.3, CHANGE, EVIDENCE CHANGED NCU00007.2,NCU00007.3,CHANGE,EVIDENCE CHANGED NCU00165.2,NCU11172.3,SPLIT, NCU00165.2,NCU11173.3,SPLIT, NCU00166.2, NCU00166.3, CHANGE, EVIDENCE CHANGED NCU11163.2,,DELETE, NCU11165.2,,DELETE, NCU11166.2, NCU11166.3, UNCHANGED, NCU11167.2, NCU11167.3, UNCHANGED, ,NCU11174.3,ADD, ,NCU11182.3,ADD, NCU03818.2, NCU03818.3, UNCHANGED, NCU03819.2, NCU03819.3, CHANGE, SPLICING CHANGED NCU03820.2, NCU11262.3, MERGE, NCU03821.2, NCU11262.3, MERGE,

# Neurospora crassa – changes between version 2 and version 3

http://www.broad.mit.edu/annotation/fgi/

# Motivation – Tools (1)

#### Converter – Ensembl

#### Convert your data

Assembly Converter enables you to convert y any other species or assembly.	our data from Mouse assembly m36 to m37. Please note that this facility is not currently co	ompatible with
Supported formats: Only GFF files with chrom Ensembl archive	osomal coordinates are currently supported. Data in this format can be exported from the	August 2007
Paste file content		
or upload file	Browse	
or use file URL		]
Upload data >		

Combining the strengths of UMIST and The Victoria University of Manchester http://www.ensembl.org/Mus\_musculus/assemblyconverter

# Motivation – Tools (2)

#### Compare versions - EMBL

Accession Number or Sequence Version: X59720 Go! Go! Go!				
Snap	Snapshot at day-month-year ( <i>e.g.</i> 30-11-1998 or 30-NOV-1998)			
Diffe	rences fo	or X59720 24-AUG-2004 / 29-MAY-2008 Back to List		
		Lines unchanged Lines removed Lines inserted		
FT	CDS	complement(1150312285)		
FT		∕product="hypothetical protein"		
FT	FT /db_xref="SGD:S0000573"			
FT	FT /db_xref="UniProt/Swiss-Prot:P25593"			
FT	FT /note="ORF YCL068c"			
FT	FT /db_xref="GOA:P25593"			
FT	T /db_xref="InterPro:IPR008937"			
FT	FT /db_xref="SGD:S00000573"			
FT	FT /db_xref="UniProtKB/Swiss-Prot:P25593"			
FT		∕product="hypothetical protein"		
FT		<pre>/protein_id="CAC42951.1"</pre>		
FT	$\prime$ translation="MFVLIDNVLAYLLEQDDLFVTARFAIQGQIVSRRVNKIHISNITD			
FT	VLLQQFISHTLPYNDNIVPKKILDSMRTAVRQLLEATACVSRECPLVKRSQDIKRARKR			
FT	LLSDWYRLGADANMDAVLLVVNSAWRFLAVWRPFVNSIQHATQELYQNIAHYLLHGNVN			
FT	I QRVTALLQLVMGQDDLLFSMDDVLQEVFRIQLYLNKMLPHNSHKWQKPSPFDSANLLL			
FT		NFRDWITDNALLQELLLSYPTINKNKHKNHSVPRLIQV"		

http://www.ebi.ac.uk/cgi-bin/sva/sva.pl

Combining the strengths of UMIST and The Victoria University of Manchester

# Motivation – Genomic data

#### Representation

- Flat file formats
  - EMBL, Genbank: information rich, well established
  - GFF (General Feature Format): less information than EMBL or Genbank flat files
- XML
  - EMBL
  - INSD, based on Genbank flat file format
  - BIODAS, based on GFF flat file format

# Overview

- Motivation
- Problem description
- XML difference algorithms
- Experiments involving controlled modifications
- Evaluation involving real genomic data
- Conclusions

#### Problem description – EMBL

gene	335649 /gene=YAL069W /locus_tag="O13588_YEAST /note="Yal069wp (YAL069W)"
mRNA	335649 /gene="YAL069W
	/note="transcript_id=YAL069W"
CDS	335649
	/gene="YAL069W"
	/protein id="YAL069W"
	/note="transcript_id=YAL069W"
	/db_xref="SGD:YAL069W"
	/db <sup>_</sup> xref="EMBL:U73805
	/translation="MIVNNTHVLTL"
exon	33564

335..64 /note="exon\_id=YAL069W.1«

. . .

CCACACCACA CCCACACACC CACACACCAC ACCACACCACC ACACCACCACC ...

Combining the strengths of UMIST and The Victoria University of Manchester

## Problem description – EMBL XML



Combining the strengths of UMIST and The Victoria University of Manchester

# Problem description – EMBL XML

<feature name="mRNA"> mRNA</feature>	335649
<qualifier name="gene"></qualifier>	
YAL069W -	/gene="YAL069W
<qualifier name="note"></qualifier>	
transcript_id=YAL069W <	<pre>- /note="transcript_id=YAL069W"</pre>
location	
<feature name="CDS">  CDS</feature>	335649
<pre><dbreference db="EMBL" primary="U73805"></dbreference> </pre>	/db_xref="EMBL:U73805
<pre><dbreference db="SGD" primary="YAL069W"></dbreference> </pre>	- /db_xref="SGD:YAL069W"
<qualifier name="gene">YAL069W </qualifier> -	/gene="YAL069W"
<qualifier name="protein_id">YAL069W </qualifier>	<pre>- /protein_id="YAL069W"</pre>
<qualifier name="note"></qualifier>	
transcript_id=YAL069W <	<pre>- /note="transcript_id=YAL069W"</pre>
<qualifier name="translation"></qualifier>	
MIVNNTHVLTL	<pre>/translation="MIVNNTHVLTL"</pre>
location	
Combining the strengths of UMIST and The Victoria University of Manchester	

### **Problem description - Changes**

- Update sequence
- Identification of a new gene
- Removal of previously predicted gene
- Merging of two neighbouring genes into one gene
- Splitting of a gene into two neighbouring genes
- ...

```
<feature name="gene">
             <qualifier name="gene">YAL069W </qualifier>
             <qualifier name="locus tag">
                                           O13588 YEAST
             </gualifier>
             <qualifier name="note">
                                            Yal069wp (YAL069W)
             </gualifier>
             <location type="single" complement="false">
                           <locationElement complement="false"
                                              type="range">
                                            <basePosition type="simple">
                                                          335
                                            </basePosition
                                            <basePosition type="simple">
                                                          649
                                            </basePosition>
                           </locationElement>
             </location>
</feature>
<feature name="mRNA">
             <qualifier name="gene">YAL069W</qualifier>
             <qualifier name="note">
                                           transcript_id=YAL069W
             </gualifier>
             location
</r></feature</td>Image: Second Control of Co
```

<feature name="CDS"> <dbreference db="EMBL" primary="U73805"/> <dbreference db="SGD" primary="YAL069W"/> <qualifier name="gene">YAL069W</qualifier> <qualifier name="protein\_id"> YAL069W </gualifier> <qualifier name="note"> transcript\_id=YAL069W </gualifier> <qualifier name="translation"> MIVNNTHVLTL... </gualifier> location </feature> <sequence length="987" type="DNA" version="0.0"> gtccagttaaggcctatc... </sequence> <feature name="exon"> <qualifier name="note"> exon id=YAL069W.1

```
</qualifier>
```

# Overview

- Motivation
- Problem description
- XML difference algorithms
- Experiments involving controlled modifications
- Evaluation involving real genomic data
- Conclusions

Summary of key properties

Algorithm	Ordered/ Unordered tree	Changes detected
X-Diff	Unordered	Insert, Delete of leaf nodes; Update of values of text- and attribute nodes
JXyDiff (XyDiff)	Ordered	Insert, Delete, Move of leaf nodes; Update of values of text- or attribute nodes
3-Way Merge and Diff (3DM)	Ordered	Insert, Delete, Move, Copy of leaf nodes; Update of values of text- or attribute nodes

#### X-Diff

- Pre-processing:
  - Parse, calculate XHash
- Matching:
  - Reduce search space
  - Bottom-up: Match nodes of same type and with matching ancestor names and calculate edit distance
  - Top-down: Create minimum-cost matching using matches and edit distances
    - One-to-one matches
    - Match child nodes of parents that are matched
- Edit script

#### JXyDiff

- Pre-processing:
  - Bottom-up: Calculate hash values and weight
  - Place sub-trees into priority queue ordered by weight
- Matching:
  - Starting with heaviest sub-tree, match sub-trees
  - If no match found, add sub-trees rooted in child nodes into priority queue
  - Propagate matches further to unmatched nodes with same labels (up - controlled by weight of sub-tree)
- Edit script

#### 3DM

- Matching:
  - For each node in base document
    - Find exact or close matching nodes in updated document (using q-gram string distance measure)
    - For pairs of matched nodes
      - Match sub-trees by depth-first traversal as long as child nodes are matched
- Post-processing:
  - Propagate matches
- Edit script

# Overview

- Motivation
- Problem description
- XML difference algorithms
- Experiments involving controlled modifications
- Evaluation involving real genomic data
- Conclusions

#### Experiments involving controlled modifications

#### Experimental setup

- Genomic data of chromosome 10 of yeast with 400 genes
- Introduce n (n = 4, 20, 40, 60, 80) changes and pair wise combinations of changes and produce new documents
  - Update sequence
  - Insert gene
  - Delete gene
  - Merge genes
  - Split gene
- Compare base document with each of new documents using X-Diff, JXy-Diff, 3DM
- Process edit scripts

#### Experiment: Insert gene



- Correct identification by X-Diff, JXyDiff, 3DM
- Additional copies of inserted sequence elements detected by 3DM

### Experiment: Insert gene & Delete gene



- 3DM:
  - + Inserted and deleted elements
  - Additional copies of sequence
- JXyDiff:
  - + Inserts and deletes of mRNA and gene
  - Incorrect matching of CDS elements
- X-Diff
  - Elements are reported as updates

# Experiment: Merge neighbouring genes



- X-Diff:
  - + Deleted elements
  - Updated elements (mRNAs - CDS)
- JXyDiff:
  - Fraction of deleted elements
  - Very few updated elements (CDS)
- 3DM:
  - + Majority of deleted and updated elements
  - Updates as inserts & deletes

### Experiment: Split gene



- X-Diff:
  - + Inserted, majority of updated
  - Updated mRNA matched with gene
- JXyDiff:
  - + Inserted
  - Few updated
- 3DM:
  - + Inserted
  - Updated

# Overview

- Motivation
- Problem description
- XML difference algorithms
- Experiments involving controlled modifications
- Evaluation involving real genomic data
- Conclusions

# Evaluation involving real genomic data

Experimental setup and results

- Ensembl releases 41 and 42 of yeast chromosomes 3 and 5
- X-Diff
  - + Singular change within an element
  - Large number of change within an element
- JXyDiff
  - Matching elements correctly, in particular CDS
- 3DM
  - Detects majority of changes
  - Insert, delete of dbreference, qualifier in CDS
  - Update of sequence
  - Change of location of genes

Chromosome	X-Diff	JXyDiff	3DM
3	1805	13223	1864
5	1130	18475	679

# Overview

- Motivation
- Problem description
- XML difference algorithms
- Experiments involving controlled modifications
- Evaluation involving real genomic data
- Conclusions

# Conclusions

- Algorithms effective for detecting subset of changes
- All require significant post-processing of edit scripts
- Properties of XML representation of genomic data
  - Large number of siblings
  - Related elements not easily identifiable as such
  - Very similar contents of nodes and sub-trees
- Properties of algorithms
  - Treating XML documents as unordered trees
  - Fuzzy matching using q-grams
  - Propagating matches bottom-up

# Future work

- Alterations to XML representation of genomic data to address:
  - Large number of siblings
  - Related elements not easily identifiable as such
  - Very similar contents of nodes and sub-trees
- New XML difference algorithm



# Thanks!

# Questions?

Combining the strengths of UMIST and The Victoria University of Manchester