# Scalable Visual Comparison of Biological Trees and Sequences

**Tamara Munzner**

**University of British Columbia**

Stanford Computer Systems Lab Colloqium (EE380)
5 May 2004

---

## Collaborators

TreeJuxtaposer joint work with
- Francois Guimbretiere, Maryland
- Serdar Tasiran, Compaq SRC
- Li Zhang, Compaq SRC
- Yunhong Zhou, Compaq SRC
- James Slack, UBC

SequenceJuxtaposer joint work with
- James Slack, UBC
- Kristian Hildebrand, UBC
- Katherine St. John, CUNY/Lehman

TJC, TJC–Q joint work with
- Dale Beerman, Virginia
- Greg Humphreys, Virginia

2

---

## Outline

Comparing big phylogenetic trees
- TreeJuxtaposer
    phylogeny background
    structural difference computation
    guaranteed visibility

Browsing huge trees
- TJC, TJC–Q

Comparing many large gene sequences
- SequenceJuxtaposer

3

---

## Tree comparison

active area: hierarchy browsing

- previous work: browsing

- comparison still open problem
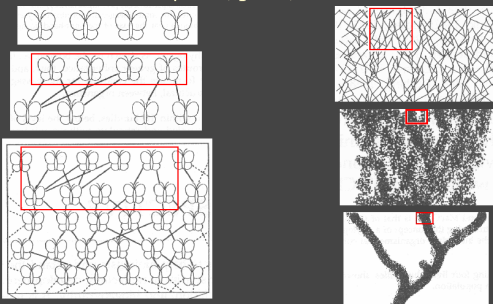
bioinformatics applicationn

- phylogenetic trees reconstructed from DNA

4

---

## Phylogeny background

tree describing evolutionary relationships
- leaves (taxa): species, genes, disease strains



[Maddison and Maddison, MacClade, 1992, p 25–26]

5

---

## Phylogenetic reconstruction

know leaves, infer interior nodes
- similarity:
  parallel evolution or common ancestor?

old: morphology
- observable similarities

new: molecular
- DNA sequences – nucleotides
- protein sequences – amino acids



[research.amnh.org/programs/genomelab]
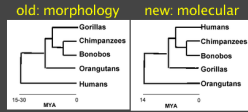
[gwis2.circ.gwu.edu/~atkins]

```
horse: ...CCTGAACCG...
tapir: ...ACTCTACCG...
rhino: ...GCTCTACCG...
```

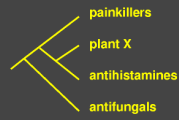6

## Phylogeny uses

establish relationships

· understand species evolution

old: morphology | new: molecular



· track diseases
  genes evolve 1M x faster

predict characteristics

· design drugs

painkillers
plant X
antihistamines
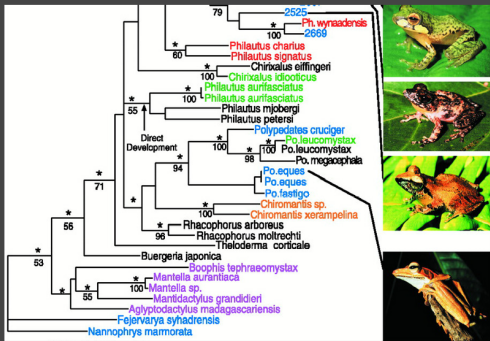antifungals

· reveal gene function

---

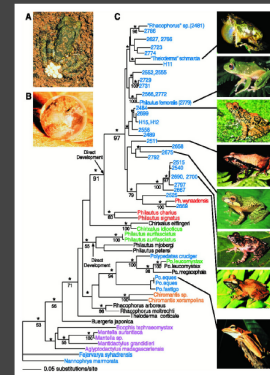## Inferring species relationships

---

## Phylogenetic/Evolutionary tree



[M Meegaskumbura et al., Science, 298:379 (2002)]

---

## Common tree size now

---

## Tree of Life: 10M species



[David Hillis, Science, 300:1687, 2003]

---

## Phylogenetic reconstruction

multiple trees
· reconstruction algorithm returns many possibilities
· different biological assumptions or data

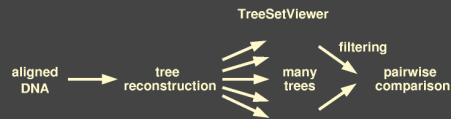aligned DNA → tree reconstruction → many trees → filtering → pairwise comparison

## Phylogenetic reconstruction

multiple trees
- · reconstruction algorithm returns many possibilities
- · different biological assumptions or data

**TreeSetViewer**



aligned DNA → tree reconstruction → many trees → filtering → pairwise comparison

visually filtering large sets of trees
- · TreeSet Viewer
  - [Amenta and Klingner, InfoVis 2002]
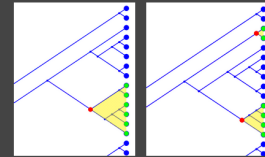
visual pairwise comparison
- · open problem

13

---

## Clades

comparing contiguous groups
- · clade: ancestor + all descendants
- · is a clade in one tree also a clade in other?
- · is some group a clade?


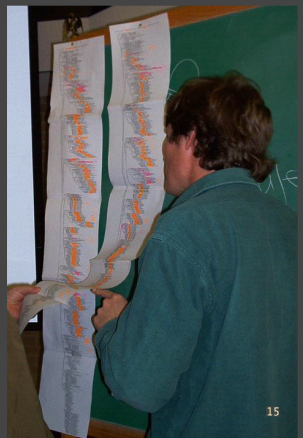
[Tree Juxtaposer first prototype, unpublished]

14

---

## Paper comparison

focus



context

Will Fischer, UT-Austin, May 2003

15

---

## Biologists' requirements

reliable detection of structural differences
- · rapid identification of interesting spots

analyses of differences in context
- · mostly side by side comparison

manipulation of increasingly larger trees

support for multiple platforms

16

---

## TreeJuxtaposer contributions

interactive tree comparison system
- · automatic detection of structural differences
  - sub-quadratic preprocessing

- · efficient Focus+Context navigation and layout
  - merge overview and detail in single view

- · guaranteed visibility under extreme distortion

scalable
- · dataset size: handles 280K–500K nodes
- · display size: handles 3800x2400 display

17

---

## TreeJuxtaposer video

platforms shown
- · java 1.4, GL4Java 2.7 bindings for OpenGL

Windows
- · 2.4 GHz P3, nVidia Quadro4 700XGL
- · 1.1GB java heap
- · window sizes 1280x1024, 3800x2400

Linux
- · 3.1 GHz P4, nVidia GeForce FX 5800 Ultra
- · 1.7GB java heap
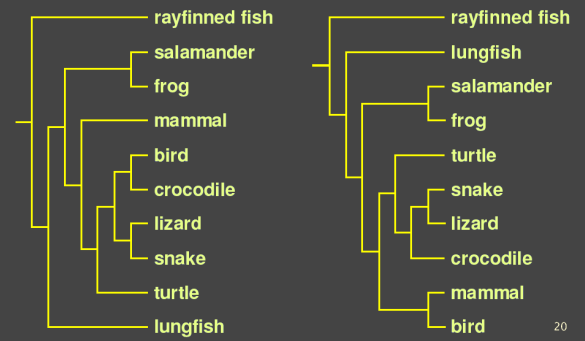- · window size 800x600

18

## Outline

Comparing big phylogenetic trees
· TreeJuxtaposer
phylogeny background
structural difference computation
guaranteed visibility

Browsing huge trees
· TJC, TJC-Q
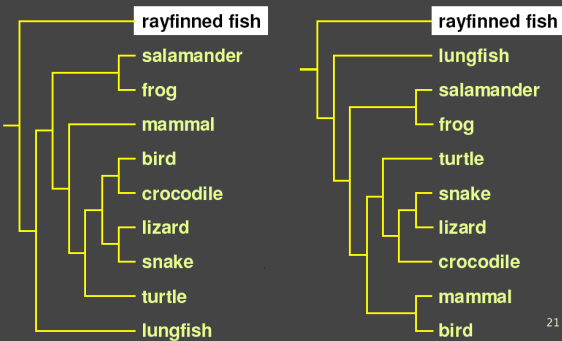
Comparing many large gene sequences
· SequenceJuxtaposer
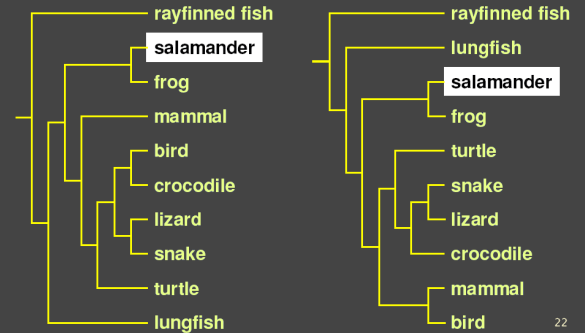
19

## Computing structural differences
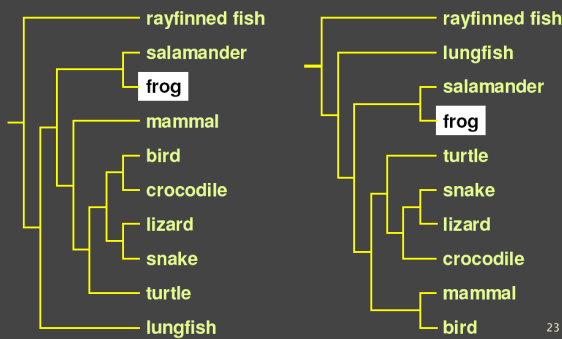


20

## Computing structural differences



21

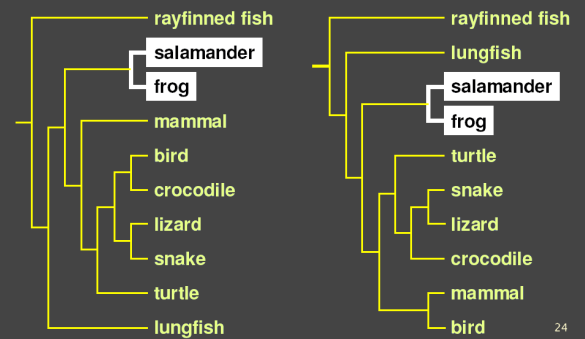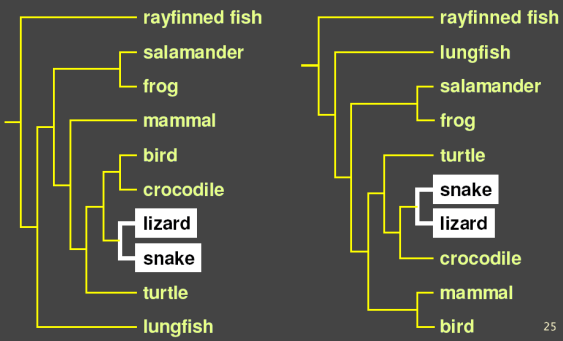## Computing structural differences



22

## Computing structural differences



23

## Computing structural differences



24

## Computing structural differences

rayfinned fish
salamander
frog
mammal
bird
crocodile
**lizard**
**snake**
turtle
lungfish

rayfinned fish
lungfish
salamander
frog
turtle
**snake**
**lizard**
crocodile
mammal
bird

25

## Computing structural differences

rayfinned fish
salamander
frog
mammal
**bird**
**crocodile**
lizard
snake
turtle
lungfish

rayfinned fish
lungfish
salamander
frog
turtle
snake
lizard
**crocodile**
**bird**
mammal

26

## Computing structural differences

rayfinned fish
salamander
frog
mammal
**?**
bird
crocodile
lizard
snake
turtle
lungfish

rayfinned fish
lungfish
salamander
frog
turtle
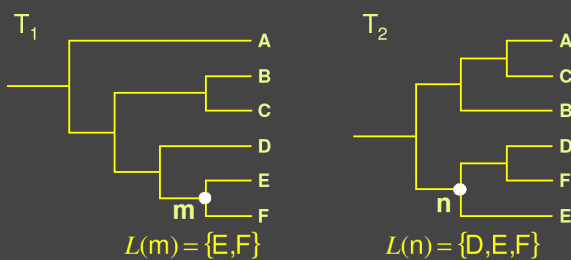snake
lizard
crocodile
mammal
bird

27

## Previous work

tree comparison

· RF distance [Robinson and Foulds 81]

· perfect node matching [Day 85]

· creation/deletion [Chi and Card 99]
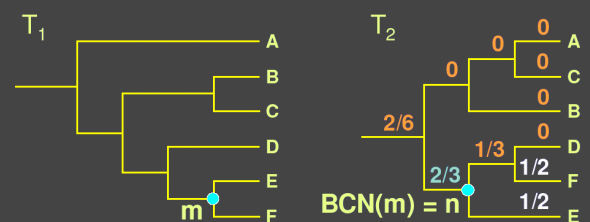
· leaves only [Graham and Kennedy 01]

28

## Similarity score

$T_1$

A
B
C
D
E
**m** F

$L(m) = \{E,F\}$

$T_2$

A
C
B
D
F
**n** E

$L(n) = \{D,E,F\}$

$$S(m,n) = \frac{|L(m) \cap L(n)|}{|L(m) \cup L(n)|} = \frac{|\{E,F\}|}{|\{D,E,F\}|} = \frac{2}{3}$$

29

## Best corresponding node

$T_1$

A
B
C
D
E
**m** F

$T_2$

0 A
0 0 C
0 B
2/6 0 D
1/3 1/2 F
2/3 1/2 E

**BCN(m) = n**

• $BCN(m) = \text{argmax}_{v \in T_2} (S(m,v))$

– computable in $O(n \log^2 n)$

– linked highlighting

30

## Marking structural differences



$T_1$    A B C D E F    **m**

$T_2$    A C B D F E    **n**

- Nodes for which $S(v, \mathrm{BCN}(v)) \neq 1$

## Structural difference algorithm

powerful and totally automatic

matches intuition
- · UT–Austin biology lab
- · other biologists
- · other domains

leads users to important locations

efficient algorithms: 7s for 2 x 140K nodes

## Outline

Comparing big phylogenetic trees
- · TreeJuxtaposer
  - phylogeny background
  - structural difference computation
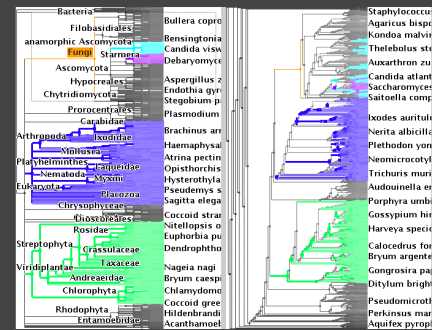  - guaranteed visibility

Browsing huge trees
- · TJC, TJC–Q

Comparing many large gene sequences
- · SequenceJuxtaposer
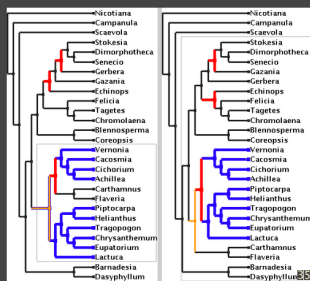
## Guaranteed mark visibility

## Marks

regions of interest shown with color highlight
- · structural difference
- · search results
- · user–specified

purpose
- · guide navigation
- · provide landmarks
- · contiguity check

## How can a mark disappear?

moving outside viewport
- · choose global Focus+Context navigation
  - "tacked–down" borders

## Focus+Context previous work

combine overview and detail into single view

Focus+Context
- large tree browsing
    - Cone Trees [Robertson et al 91]
    - Hyperbolic Trees [Lamping et al 95, Munzner 97]
    - Space Tree [Plaisant et al 03]
    - DOI Tree [Card and Nation 02]
- global
    - Document Lens [Robertson and Mackinlay 93]
    - Rubber Sheets [Sarker et al 93]

our contribution
- scalability, guaranteed visibility

## How can a mark disappear?

moving outside viewport
- choose global Focus+Context navigation "tacked-down" borders
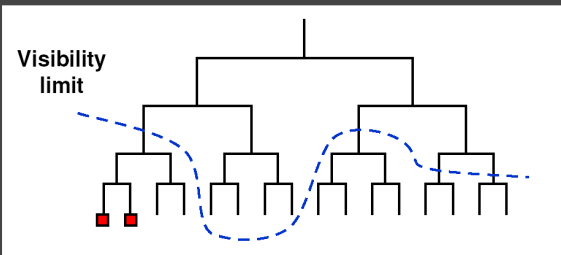
occlusion
- choose 2D++ layout

culling at subpixel sizes
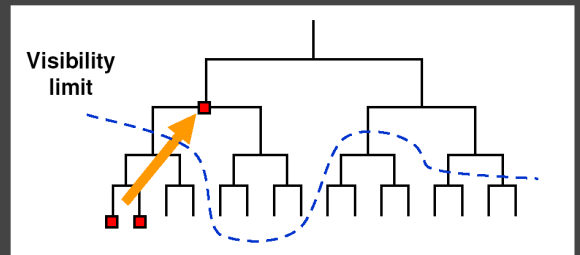- develop efficient check for marks when culling

## Preserving marks while culling
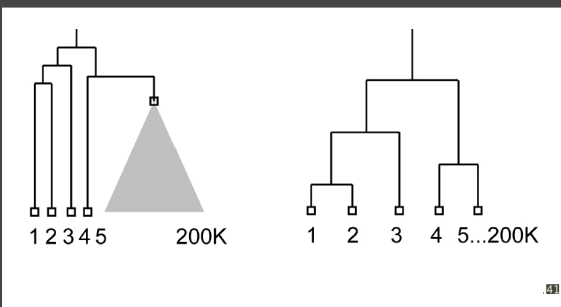
show mark at unculled node

## Preserving marks while culling

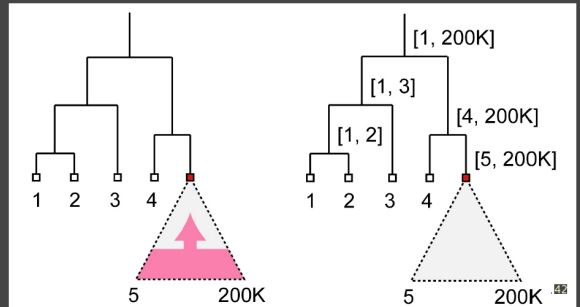show mark at unculled node

## Mark preservation strategies

compress large subtree to small spatial area

## Precompute subtree ranges

propagation: cost depends on total nodes
precomputation: cost depends on visible nodes

## Marks and linked highlighting

also check for linked marks from other tree

check if best match for node is marked
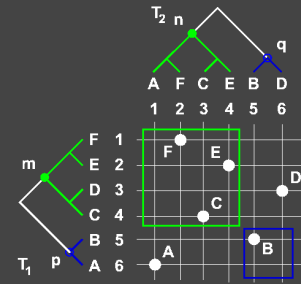- up to O(n) to look up each node in range

intersect node ranges between trees
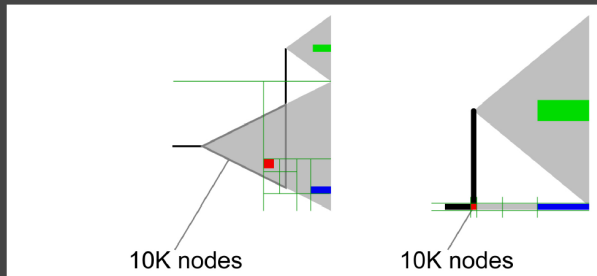- reduces to point in polygon test, $O(\log^2 n)$

## Efficient marking detection

intersecting ranges between trees

## Storing spatial ranges

in each box, store range of objects inside



10K nodes          10K nodes

## Spatial range solution

recursive spatial subdivision
- quadtree
- store range of objects enclosed for each cell
- quick check: spatial range vs. selection range

extending quadtrees to Focus+Context
- quadtree cells also "painted on rubber sheet"
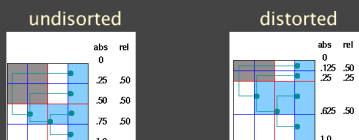- efficient O(log n) update when stretch/shrink

## Focus+Context quadtree

sparse cell instantiation
map from cell boundary to object location
store absolute location of cell boundaries?
- lookup: O(1), update: O(n)
instead, store boundaries hierarchically
- relative "split" between parent cell boundaries



undisorted          distorted

## Guaranteed visibility

infrastructure needed for efficient computation

relief from exhaustive exploration
- missed marks lead to false conclusions
- hard to determine completion
- tedious, error-prone

compelling reason for Focus+Context
- controversy: does distortion help or hurt?
- strong rationale for comparison

constraint to fit everything in viewport
- instead could show indirectly
- ideas: Halo [Baudisch 03]

## Guaranteed visibility previous work

visibility of abstract information

- · effective view navigation [Furnas 97]

- · critical zones [Jul and Furnas 98]

49

## TreeJuxtaposer contributions

first interactive tree comparison system
- · automatic structural difference computation
- · guaranteed visibility of landmark areas

scalable to large datasets
- · 250,000 to 500,000 total nodes
- · all preprecessing subquadratic
- · all realtime rendering sublinear

techniques broadly applicable
- · not limited to biological trees

overall winner: InfoVis Contest 2003

50

## Outline

Comparing big phylogenetic trees
- · TreeJuxtaposer
    phylogeny background
    structural difference computation
    guaranteed visibility

Browsing huge trees
- · TJC, TJC-Q

Comparing many large gene sequences
- · SequenceJuxtaposer

51

## Scaling up

TreeJuxtaposer limits
- · memory footprint
- · rendering CPU bound, want graphics bound

goal: browse huge trees
- · concentrate on browsing
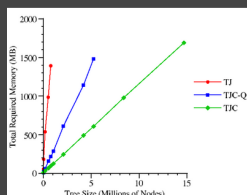
TJC-Q: 5M nodes
- · commodity platforms

TJC: 15M nodes
- · leading-edge graphics hardware

[video]

52

## Memory footprint reduction



TJ quadtrees
- · navigating, culling, drawing, picking

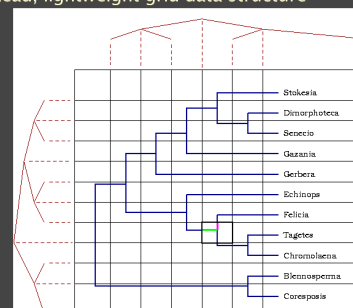new algorithms for drawing/culling
new data structures
- · TJC-Q: low-memory quadtrees
- · TJC: no quadtrees, picking with hardware

53

## Quadtree: navigating

navigating with stretch/shrink
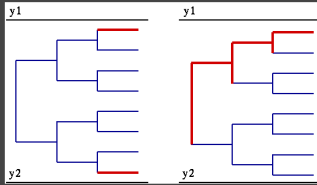- · instead, lightweight grid data structure



54

## Quadtree: culling and drawing

culling subpixel objects
· leaf overlap test, not gridcell size test



drawing in order of importance
· new alg fast enough to ignore order

## Quadtree: picking

TJ: picking with spatial subdivision

TJC: multiple render target buffer
· encode object ID into offscreen buffer
· supported in hardware on latest ATI cards

TJC-Q: low-memory quadtrees

## Outline

Comparing big phylogenetic trees
· TreeJuxtaposer
    phylogeny background
    structural difference computation
    guaranteed visibility

Browsing huge trees
· TJC, TJC-Q

Comparing many large gene sequences
· SequenceJuxtaposer

## Accordion drawing

general scalable visualization infrastructure

· "rubber sheet" navigation

· guaranteed visibility of marked areas

modular package

· layer below TreeJuxtaposer

· not just for trees

## SequenceJuxtaposer

accordion drawing for DNA/RNA

previous work: web-based sequence browsers
· Ensembl, UCSC Genome Browser, NCBI MapViewer
· heavily used, huge server-side databases

· zoom or pan in jumps
· can't see context

fluid Focus+Context navigation
guaranteed visibility
· establish when these features useful
· proof of concept prototype, eventually merge

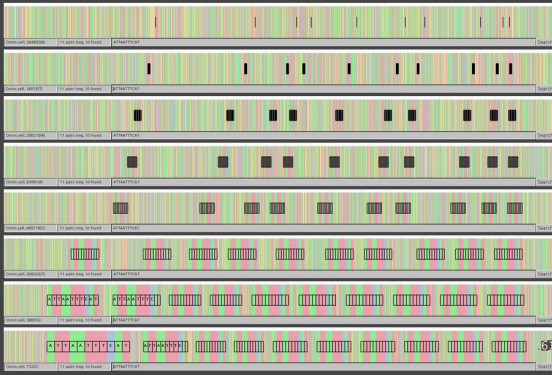## SJ in action

shown on publicly available data

· onion yellows phytoplasma: whole genome
    860 Kbp

· Murphy: 22 genes
    44 mammals x 17000 bp each = 748 Kbp

· Treezilla: single gene
    500 plants x  1428 bp each = 714 Kbp

scales to 1.7 Mbp with 1.7GB heap

[videos]

## Expanding search results



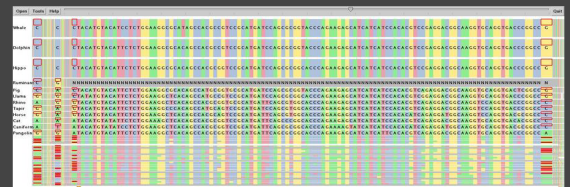## Changing difference thresholds



25%

## Changing difference thresholds



50%

## Changing difference thresholds



60%

## Changing difference thresholds



67%

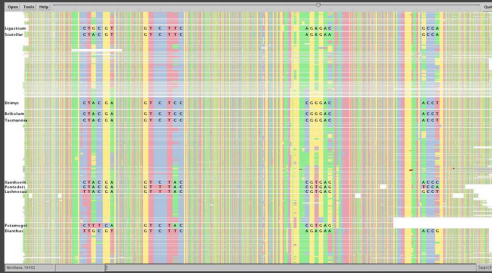phylogenetic signal visible
inspecting 1 of 22 genes

## Codon bias shown with visual patterns

# Codon bias shown with visual patterns



67

# Work in progress

trees with weighted edges

protein sequences

linking tree and sequence navigation

accordion drawing for sets
· data mining: transaction processing
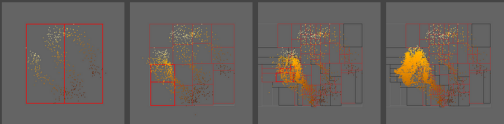
open-source release
· stay tuned!

68

# Other projects in progress

dimensionality reduction
· steerable MDS (multidimensional scaling)
· (with Matt Williams)



perception experiments
· quantifying cost of Focus+Context fisheye distortions
· no-cost and low-cost regions for visual search task
· (with Keith Lau, Ron Rensink)

69

# More information

www.cs.ubc.ca/~tmm/papers.html
www.cs.ubc.ca/~tmm/talks.html

papers, slides, images, movies

software: beta now, public release very soon