## Slide 1

# Visualization and Biology:
# Fertile Ground for Collaboration

Tamara Munzner
Department of Computer Science
University of British Columbia
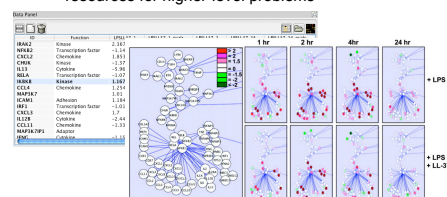
June 2009

http://www.cs.ubc.ca/~tmm/talks.html#eindhoven09

## Slide 2

# Outline

- visualization ideas and background

- combining interaction networks, microarray data
  – Cerebral system

- comparing phylogenetic trees
  – TreeJuxtaposer system

2

## Slide 3

# Why do visualization?

- pictures help us think
  – substitute perception for cognition
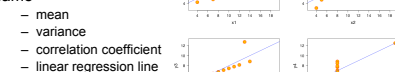  – external memory: free up limited cognitive/memory resources for higher-level problems



3

## Slide 4

# When should we bother doing vis?

- need a human in the loop
  – augment, not replace, human cognition
  – for problems that cannot be (completely) automated
- simple summary not adequate
  – statistics may not adequately characterize complexity of dataset distribution

Anscombe's quartet: same
  – mean
  – variance
  – correlation coefficient
  – linear regression line

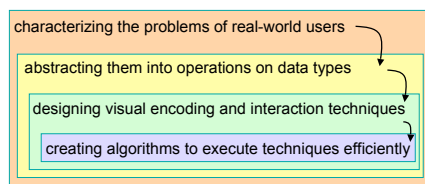http://upload.wikimedia.org/wikipedia/commons/b/b6/Anscombe.svg

4

## Slide 5

# What does visualization allow?

- discovering new things
  – hypothesis discovery, "eureka moment"
- confirming conjectured things
  – hypothesis confirmation
- contradicting conjectured things
  – especially (inevitably?) data cleansing

- novel capabilities
  – tool supports fundamentally new operations
- **speedup**
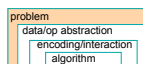  – tool accelerates workflow (most common!)

5

## Slide 6

# Multiple levels of problem-driven vis

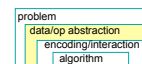- cascading levels: output above is input below

characterizing the problems of real-world users

abstracting them into operations on data types

designing visual encoding and interaction techniques

creating algorithms to execute techniques efficiently

6

## Slide 7

# Characterizing problems

problem
  data/op abstraction
    encoding/interaction
      algorithm

- understanding domain concepts and current workflow
- finding gaps, breakdowns, slowdowns
  – where conjecture that vis would help

7

## Slide 8
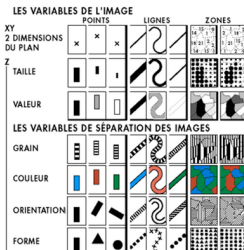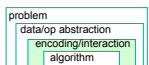
# Abstracting into operations on data types

problem
  data/op abstraction
    encoding/interaction
      algorithm

- data types
  – tables of numbers
  – relations: networks/graphs, hierarchies/trees
  – spatial data: geographic, positions in space
- operations
  – sorting, filtering, browsing, comparison, characterizing trends and distributions, finding anomalies and outliers, finding correlation...
  – relations: following path through network...
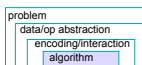
8

## Slide 9

# Designing encoding and interaction

problem
  data/op abstraction
    encoding/interaction
      algorithm

- visual encoding
  – marks: points, lines, areas
  – attributes: position, color, shape, size, orientation, ...
- interaction
  – selecting, navigating, ordering,...

LES VARIABLES DE L'IMAGE
LES VARIABLES DE SÉPARATION DES IMAGES

Semiology of Graphics. Jacques Bertin, Gauthier-Villars 1967, EHESS 1998    9

## Slide 10

# Creating efficient algorithms

problem
  data/op abstraction
    encoding/interaction
      algorithm

- classic computer science problem
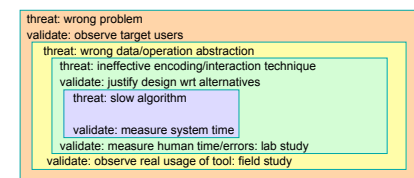  – create algorithm given clear specification

10

## Slide 11

# Design decisions

- huge space of design alternatives

- many/most choices are ineffective
  – wrong visual encoding can mislead, confuse
  – principled reasons to make choices usually not obvious to untrained people

  – conflicting tradeoffs
    • iterative refinement often necessary

11

## Slide 12

# Validation: Is problem solved?

- humans in the loop for outer three levels

threat: wrong problem
validate: observe target users
  threat: wrong data/operation abstraction
    threat: ineffective encoding/interaction technique
    validate: justify design wrt alternatives
      threat: slow algorithm
      validate: measure system time
    validate: measure human time/errors: lab study
  validate: observe real usage of tool: field study

12

## Slide 13

# Collaboration: Complementary expertise

- vis researchers
  – vis design alternatives
  – human perceptual capabilities
  – scalable graphics algorithms
  – validation methodology
- domain scientists
  – deep knowledge of driving problems, data

- both benefit from new tools
  – scientist: you get something helpful
  – vis researcher: we get to watch you use it
    • see if problem actually solved
    • feed new knowledge back into our design principles

13

## Slide 14

# Good driving problems for vis research

- big data
- reasonably clear questions
- need for humans in the loop

- many areas of science are a great match
  – biology particularly appealing

14

## Slide 15

# Outline

- visualization ideas and background

- combining interaction networks, microarray data
  – Cerebral system

- comparing phylogenetic trees
  – TreeJuxtaposer system

15

## Slide 16

# Cerebral

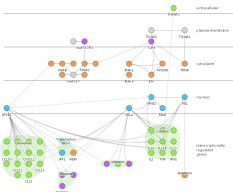collaboration with researchers at UBC Hancock Lab studying innate immunity

open-source software download (Cytoscape plugin)
http://www.pathogenomics.ca/cerebral/
deployed in InnateDB (mammalian innate immunity database)
http://www.innatedb.ca

16

## Systems biology model

- graph G = {V, E}
  - V: proteins, genes, DNA, RNA, tRNA, etc.
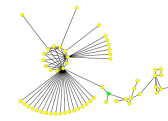  - E: interacting molecules



17

## Model - Experiment cycle

- conduct experiments on cells
- interpret results in current graph model
- propose modifications to refine model

- vis tool to accelerate workflow?

18

## Goal: Integrate model with measurements

- system model
  - interaction graph G = {V, E}
  - meta-data for each v in V
    - labels, biological attributes
- experimental measurements
  - multiple floats for each v in V
    - microarray data



19

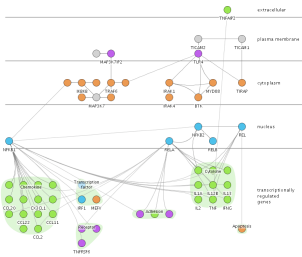## Model summarizes extensive lab work

- graphs come from hand-curated databases
  - dynamic, change with each new publication

- each edge has provenance from experimental evidence
  - TIRAP: an adapter molecule in the Toll signaling pathway. *Horng T, Barton GM, Medzhitov R.*
  - Mal (MyD88-adapter-like) is required for Toll-like receptor-4 signal transduction. *Fitzgerald KA, Palsson-McDermott EM, Bowie AG, Jefferies CA, Mansell AS, Brady G, Brint E, Dunne A, Gray P, Harte MT, McMurray D, Smith DE, Sims JE, Bird TA, O'Neill LA.*

- choose scope to manage complexity
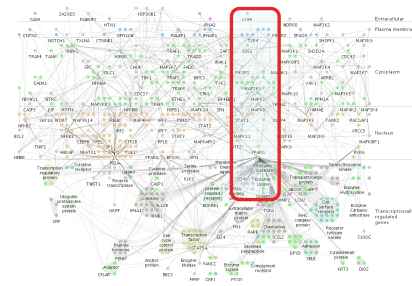
20

## TLR4 biomolecule:  E=74, V=54

- very local view
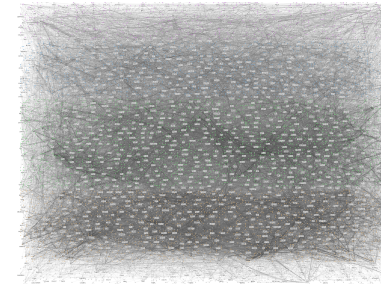


21

## Immune system: E=1263, V=760
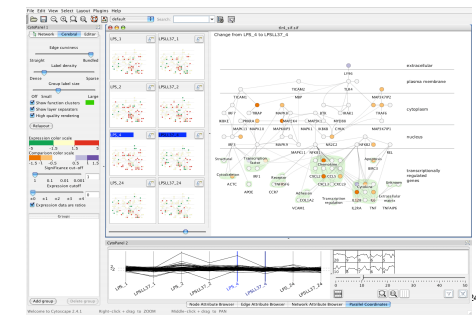
- bigger picture, target size for Cerebral



22

## Human interactome: E~50,000, V~10,000

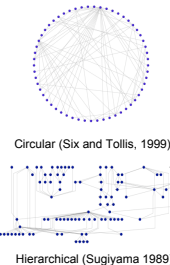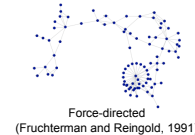- too complex, beyond scope of tool



## Cerebral video



## Encoding and interaction design decisions

- create custom graph layout
  - guided by biological metadata
- use small multiple views
  - one view per experimental condition
- show measured data in graph context
  - not in isolation

25

## Traditional graph layout

- given graph G={V,E}
- create layout in 2D plane
- heavily studied
  - hundreds of papers
  - annual Graph Drawing conf

Circular (Six and Tollis, 1999)

Force-directed (Fruchterman and Reingold, 1991)
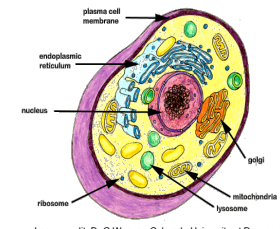
Hierarchical (Sugiyama 1989)

26

## Existing layouts did not suit immunologists

- graph drawing goals
  - visualize graph structure
- biologist goals
  - visualize biological knowledge
  - some relationships happen to form a graph
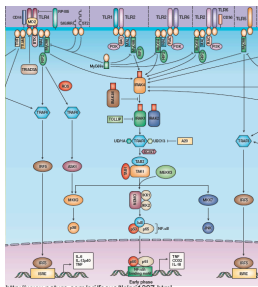  - cell location also relevant

27

## Biological cells divided by membranes

- interactions generally occur within a compartment
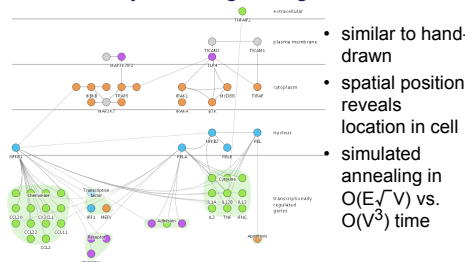- crossing membranes is interesting



Image credit: Dr.G Weaver, Colorado University at Denver

28

## Hand-drawn diagrams



- cellular location encoded spatially
- infeasible to create by hand in era of big data

http://www.nature.com/nri/focus/tlr/nri1397.html

29

## Cerebral layout using biological metadata
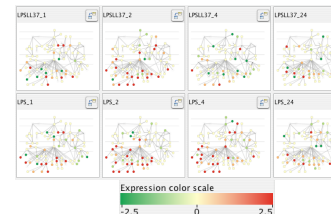


- similar to hand-drawn
- spatial position reveals location in cell
- simulated annealing in $O(E\sqrt{V})$ vs. $O(V^3)$ time
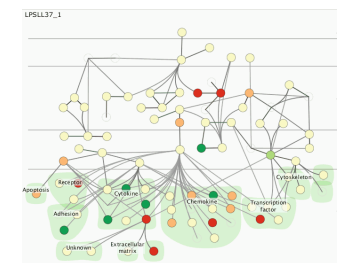
30

## Use small multiple views

- one graph instance per experimental condition
  - same spatial layout
  - color differently, by condition



Expression color scale
-2.5   0   2.5

31

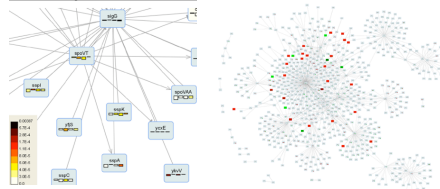## Why not animation?

- global comparison difficult



32

# Why not animation?

- limits of human visual memory
  - compared to side by side visual comparison

- Matthew Plumlee and Colin Ware. Zooming versus multiple window interfaces: Cognitive costs of visual comparisons. *ACM Trans. Computer-Human Interaction (ToCHI)*,13(2):179-209, 2006.

- Barbara Tversky, Julie Bauer Morrison, and Mireille Betrancourt. Animation: can it facilitate? *International Journal of Human-Computer Studies,* 57(4):247-262, 2002.

33

---

# Why not glyphs?

- embed multiple conditions as a chart inside node
- clearly visible when zoomed in
- but cannot see from global view
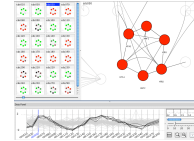  - only one value shown in overview



[M. A. Westenberg, S. A. F. T. van Hijum, O. P. Kuipers, J. B. T. M. Roerdink. Visualizing Genome Expression and Regulatory 34 Network Dynamics in Genomic and Metabolic Context. Computer Graphics Forum, 27(3):887-894, 2008.]
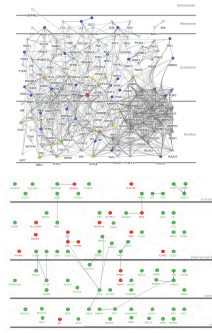
---

# Show measured data in graph context

- data driven hypothesis
  - clusters indicate similar function?
  - same pattern of gene expression → same role in cell?
- clusters are often untrustworthy artifacts!
  - noisy data: different clustering alg. → different results
  - measured data alone potentially misleading
  - **show in context of graph model**

35

---

# Adoption by biologists

- Matthew D Dyer, T. M Murali, and Bruno W Sobral. The landscape of human proteins interacting with viruses and other pathogens. PLoS Pathogens, 4(2):e32, 2008.

- Liqun He et al. The glomerular transcriptome and a predicted protein-protein interaction network. Journal of the American Society of  Nephrology, 19(2):260-268, 2008.
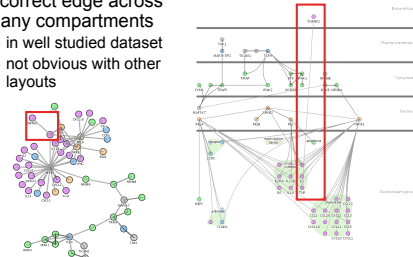
---

# InnateDB links to Cerebral

- InnateDB: facilitating systems-level analyses of the mammalian innate immune response
  - David J Lynn, Geoffrey L Winsor, Calvin Chan, Nicolas Richard, Matthew R Laird, Aaron Barsky, Jennifer L Gardy, Fiona M Roche, Timothy H W Chan, Naisha Shah, Raymond Lo, Misbah Naseer, Jaimmie Que, Melissa Yau, Michael Acab, Dan Tulpan, Matthew D Whiteside, Avinash Chikatamarla, Bernadette Mah, Tamara Munzner, Karsten Hokamp, Robert E W Hancock, Fiona S L Brinkman. Molecular Systems Biology 2008; 4:218
  - http://innatedb.ca

37

---

# Data cleansing example

- incorrect edge across many compartments
  - in well studied dataset
  - not obvious with other layouts

38

---

# Cerebral summary

- supports interactive exploration of multiple experimental conditions in graph context

- provides familiar representation by using biological metadata to guide graph layout

39

---

# Outline

- visualization ideas and background

- combining interaction networks, microarray data
  - Cerebral system

- comparing phylogenetic trees
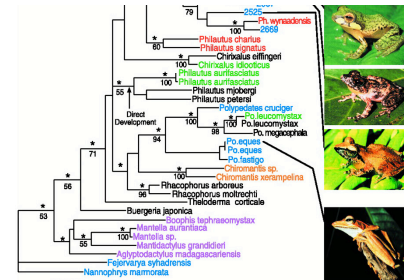  - TreeJuxtaposer system

40

---

# TreeJuxtaposer

collaboration with biologists at UT-Austin Hillis Lab

TreeJuxtaposer: Scalable Tree Comparison using Focus+Context with Guaranteed Visibility.
Tamara Munzner, François Guimbretière, Serdar Tasiran, Li Zhang, Yunhong Zhou.
ACM Trans. Graphics 22(3): 453-462, 2003 (Proc. SIGGRAPH 2003).
http://www.cs.ubc.ca/labs/imager/tr/2003/tj

open-source software download
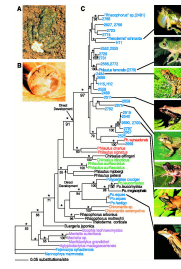http://olduvai.sourceforge.net/tj

41

---

# Phylogenetic (evolutionary) tree
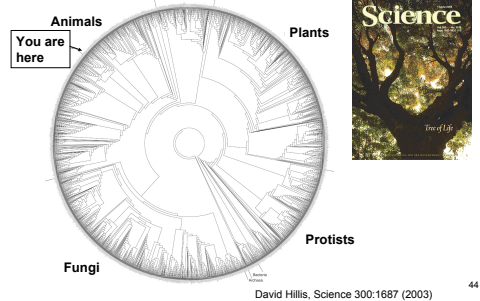
M Meegaskumbura et al., Science 298:379 (2002)

42

---

# Common dataset size today

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

43

---

# Future goal: Full Tree of Life, ~10M nodes

Animals

You are here

Plants

Protists

Fungi

David Hillis, Science 300:1687 (2003)

44

---

# Operation: Comparing multiple trees

- presentation: single tree shown as final result

- exploration: determine true tree from many possibilities
  - different biological conjectures or data
  - different phylogenetic reconstruction algorithms
  - multiple alternatives from same reconstruction algorithm

- most previous work on browsing
  - necessary but not sufficient for comparison

45

---

# Limitations of paper: Scale and speed

- literal: actual paper
- figurative: interfaces with same semantics as paper

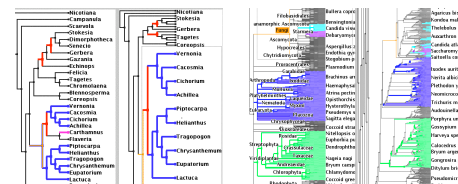need to focus on details                yet maintain context

46

---

# TreeJuxtaposer video

- stretch and squish navigation
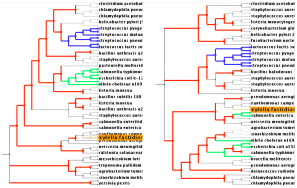- linked side by side comparison

47

---

# Encoding and interaction design decisions

- guaranteed visibility of small marks
  - scaling up to millions of nodes

48

## Guaranteed visibility

- marks are always visible
  - structural differences, search results, user selections
- easy with small datasets
  - regions of interest shown with color highlights



49

## Guaranteed visibility challenges

- hard with larger datasets
- reasons a mark could be invisible
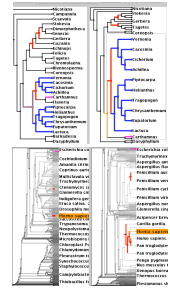
50

## Guaranteed visibility challenges

- hard with larger datasets
- reasons a mark could be invisible
  - mark outside the window
    - solution: constrained navigation

51

## Constrained navigation for visibility

- stretch and squish navigation
  - stretch out part of surface, the rest squishes
  - borders nailed down
  - integrated focus and context
- items never fall outside camera
  - but squished regions can have many items per pixel

52

## Guaranteed visibility challenges

- hard with larger datasets
- reasons a mark could be invisible
  - mark outside the window
    - solution: constrained navigation
  - mark underneath other marks
    - solution: use 2D not 3D layout

53

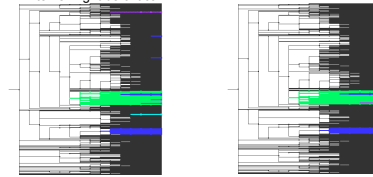## Guaranteed visibility challenges

- hard with larger datasets
- reasons a mark could be invisible
  - mark outside the window
    - solution: constrained navigation
  - mark underneath other marks
    - solution: use 2D not 3D layout
  - mark smaller than a pixel
    - solution: smart culling

54

## Smart culling for small item visibility

- naïve culling does not draw all marked items
  - graphics cards optimized for realism: small items far away and thus not important
  - rendering infrastructure for visualization semantics: small items might be critical!

guaranteed mark visibility    no guaranteed visibility    55

## Guaranteed visibility benefits

- with GV
  - no mark is visible means no need to explore area further
- without GV
  - risk of false negative conclusions, or
  - user must do tedious exhaustive search to ensure nothing missed
- algorithm scalability challenge
  - rendering complexity based on number of onscreen pixels
    - not total number of items in dataset
    - Partitioned Rendering Infrastructure for Scalable Accordion Drawing (Extended Version). James Slack, Kristian Hildebrand, and Tamara Munzner. Information Visualization, 5(2), p. 137-151, 2006
    - Composite Rectilinear Deformation for Stretch and Squish Navigation. James Slack and Tamara Munzner. Proc. Visualization 2006, published as Transactions on Visualization and Computer Graphics 12(5), September 2006.

56

## TJ summary

- first interactive tree comparison system
  - automatic structural difference computation
  - guaranteed visibility of small marks
- scalable to large datasets
  - 250K to 500K total nodes: original
  - up to 4M nodes: later, with PRISAD
  - subquadratic preprocessing
  - sublinear realtime rendering
    - depends on number of pixels, not number of nodes

57

## More information

- this talk
  http://www.cs.ubc.ca/~tmm/talks.html#eindhoven09

- papers, videos
  http://www.cs.ubc.ca/~tmm

- software
  http://olduvai.sourceforge.net/tj

  http://www.pathogenomics.ca/cerebral

  http://www.innatedb.ca

58