

Visualization Analysis and Design for Business Intelligence

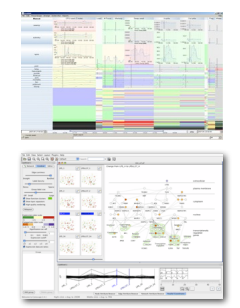
Tamara Munzner
 Department of Computer Science
 University of British Columbia

Disney Research
 20 July 2015, Glendale CA
<http://www.cs.ubc.ca/~tmm/talks.html#disney15>

@tamaramunzner

Outline

- **interactive visual analysis**
 - role and advantages
- LiveRAC
 - time-series data: managed web hosting (with AT&T)
- Cerebral
 - network of relationships: genes (with Agilent and UBC Immunology)
- wrapup



2

Defining visualization (vis)

Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

- what counts as effective?
 - novel: enable entirely new kinds of analysis
 - faster: speed up existing workflows
 - most common case!

Why have a human in the loop?

Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

Visualization is suitable when there is a need to augment human capabilities rather than replace people with computational decision-making methods.

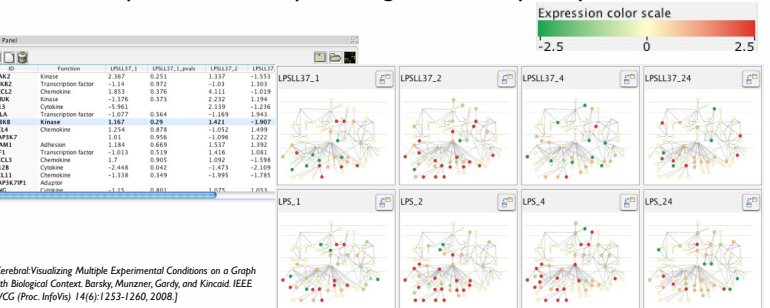
- don't need vis when fully automatic solution exists and is trusted
- many analysis problems ill-specified
 - don't know exactly what questions to ask in advance
- possibilities
 - long-term use for end users (e.g. exploratory analysis of scientific data)
 - presentation of known results
 - stepping stone to better understanding of requirements before developing models
 - help developers of automatic solution refine/debug, determine parameters
 - help end users of automatic solutions verify, build trust

4

Why use an external representation?

Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

- external representation: replace cognition with perception



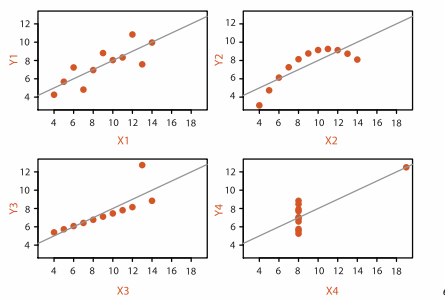
[Cerebral Visualizing Multiple Experimental Conditions on a Graph with Biological Context. Bersky, Munzner, Gandy, and Kincaid. IEEE TVCG (Proc. InfoVis) 14(6):1253-1260, 2008.]

Why show the data in detail?

- summaries lose information
 - confirm expected and find unexpected patterns
 - assess validity of statistical model

Anscombe's Quartet

Identical statistics	
x mean	9
x variance	10
y mean	8
y variance	4
x/y correlation	1



6

What are the resource limitations?

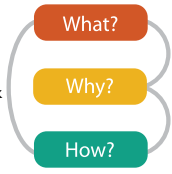
Vis designers must take into account three very different kinds of resource limitations: those of computers, of humans, and of displays.

- computational limits
 - processing time
 - system memory
- human limits
 - human attention and memory
- display limits
 - pixels are precious resource, the most constrained resource
 - **information density**: ratio of space used to encode info vs unused whitespace
 - tradeoff between clutter and wasting space, find sweet spot between dense and sparse

How to analyze vis design?

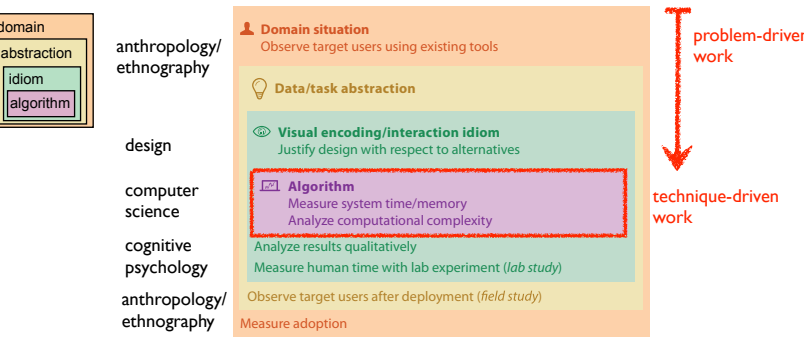
Vis usage can be analyzed in terms of what data is shown, why the user needs it, and how the idiom is designed.

- abstractions
 - translate from specifics of domain to vocabulary of vis
 - data abstraction: **what** to show
 - might not draw what you're given: transform data into form useful for task
 - task abstraction: **why** they're looking at it
- idioms
 - visual encoding idiom: **how** to draw
 - interaction idiom: **how** to manipulate
- analysis framework: scaffold to think systematically about design space
 - huge, and most possibilities ineffective for particular task/data combination



8

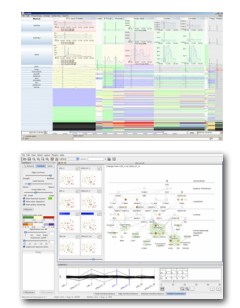
How to validate design?



[A Nested Model of Visualization Design and Validation. Munzner. IEEE TVCG 15(6):921-928, 2009 (Proc. InfoVis 2009).]

Outline

- interactive visual analysis
 - role and advantages
- LiveRAC
 - time-series data: managed web hosting (with AT&T)
- Cerebral
 - network of relationships: genes (with Agilent and UBC Immunology)
- wrapup



10

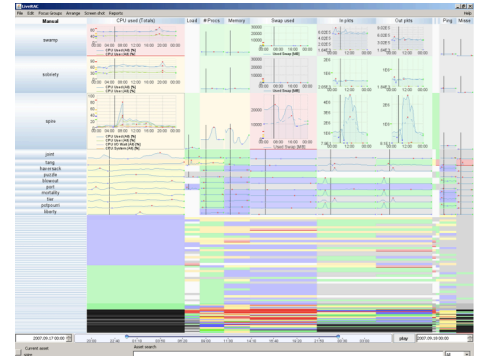
LiveRAC

Interactive Visual Exploration of System Management Time-Series Data

joint work with:
 Peter McLachlan, Eleftherios Koutsosifos, Stephen North.
<http://www.cs.ubc.ca/labs/imager/tr/2008/liverac>

LiveRAC - Interactive Visual Exploration of System Management Time-Series Data. McLachlan, Munzner, Koutsosifos, North. Proc. SIGCHI Conference on Human Factors in Computing Systems (CHI'08), p. 1483-1492, 2008.

LiveRAC video

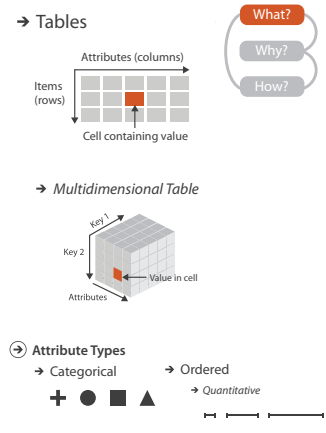


<http://youtu.be/id0c3H0Y5kw>

12

What: Data abstraction

- multidimensional table: time series data
 - key attributes
 - time
 - 50,000: 5-minute intervals over 6 months
 - multiscale levels of interest
 - devices
 - 4000
 - parameters
 - 20
 - ex: CPU usage, memory load, network traffic, alarms, ...
 - value attributes
 - parameter value for device at time point
 - quantitative
 - device groups
 - categorical



Why: Tasks in domain language

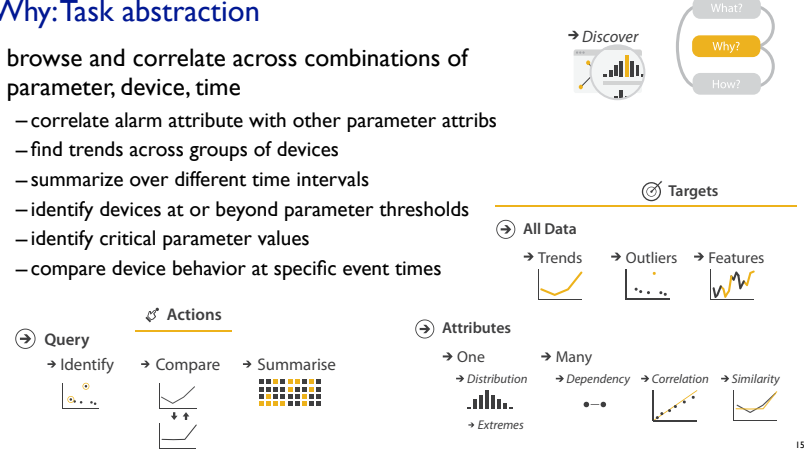
- interpret network environment status
- report generation
- capacity planning
- event investigation/forensics
- coordination
 - between customers, engineering, ops



14

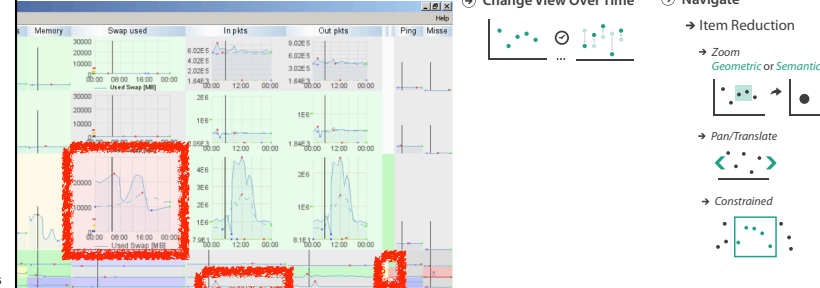
Why: Task abstraction

- browse and correlate across combinations of parameter, device, time
 - correlate alarm attribute with other parameter attribs
 - find trends across groups of devices
 - summarize over different time intervals
 - identify devices at or beyond parameter thresholds
 - identify critical parameter values
 - compare device behavior at specific event times



How: Navigate

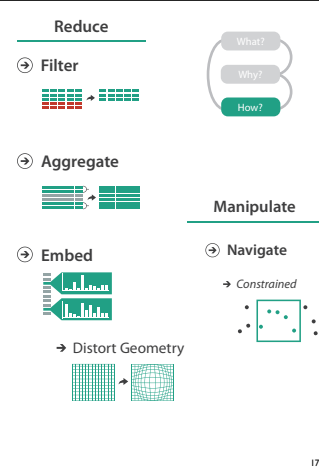
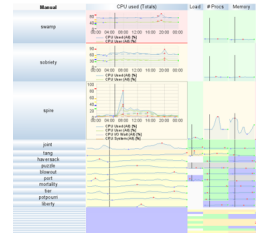
- semantic zooming: adapts to pixels available
 - many: superimposed line charts with full labeling
 - some: iconic line chart (sparkline)
 - few: color-coded box (heatmap)



16

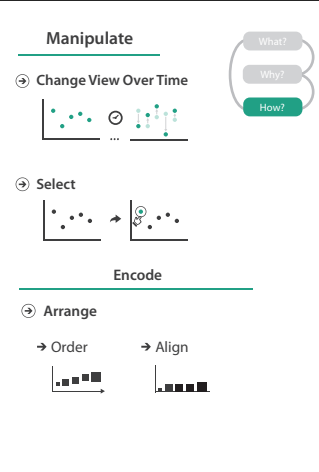
How: Reduce

- reduce data shown with complex combination of filtering and aggregation
- embed focus+context in single view
- distort geometry
 - metaphor: stretch and squish navigation
 - shape: rectilinear
 - foci: multiple
 - impact: global



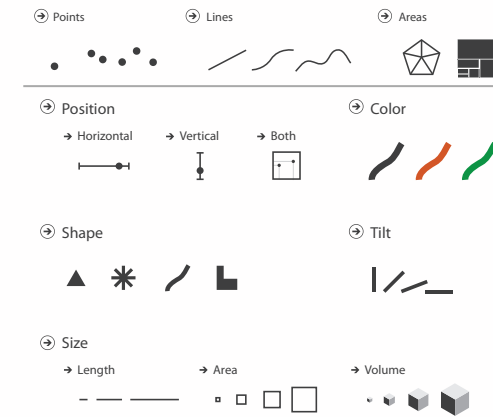
How: Reordering

- change spatial arrangement
 - resort by selected attribute
 - check for correlations between aligned attribute columns
 - ex: high load without high CPU, maybe I/O bound

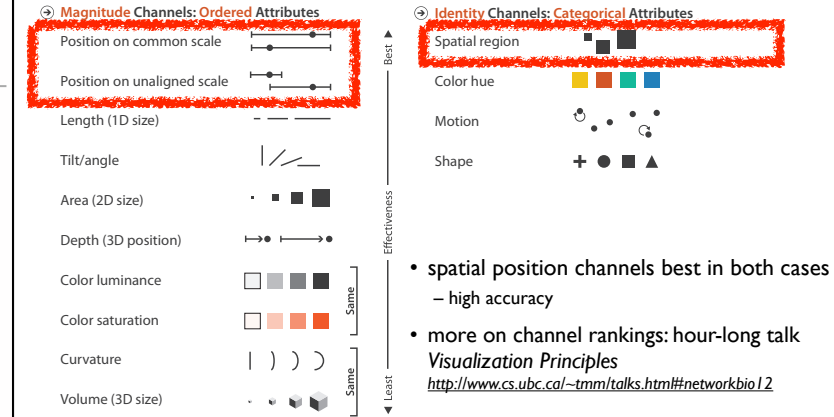


Importance of arranging space: Underlying definitions

- marks
 - geometric primitives
- channels
 - control appearance of marks

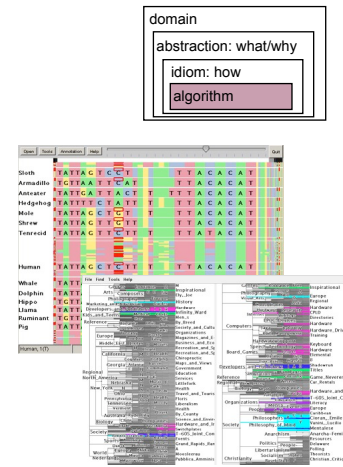


Channels: Expressiveness types and effectiveness rankings



Algorithms

- back end: SWIFT server
- front end: PRISAD rendering
 - separate threads for render vs server update
 - guaranteed visibility of semantically important marks even when squished small
 - sublinear rendering: $O(p)$ where p = pixel count
 - scalable for n of millions
 - generic framework
 - » time series charts, gene sequences, trees



[Partitioned Rendering Infrastructure for Scalable Accordion Drawing (Extended Version). Slack, Hildebrand, and Munzner. Information Visualization, 5(2), p. 137-151, 2006.]

Outline

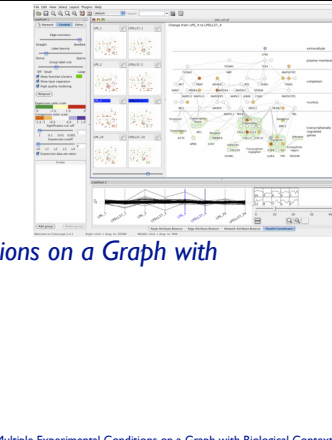
- interactive visual analysis
 - role and advantages
- LiveRAC
 - time-series data: managed web hosting (with AT&T)
- Cerebral
 - network of relationships: genes (with Agilent and UBC Immunology)
- wrapup



Cerebral

Visualizing Multiple Experimental Conditions on a Graph with Biological Context

joint work with:
Aaron Barsky, Jennifer Gardy, Robert Kincaid
<http://www.pathogenomics.ca/cerebral/>



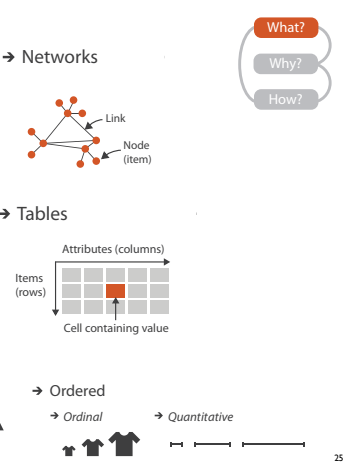
Cerebral: Visualizing Multiple Experimental Conditions on a Graph with Biological Context. Barsky, Munzner, Gardy, Kincaid. IEEE Trans. Visualization and Computer Graphics 14(6):1253-1260 2008. (Proc. InfoVis 2008).

Cerebral video



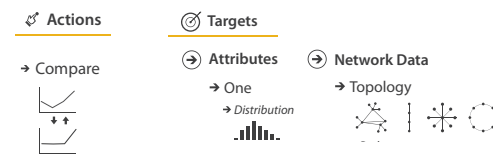
What: Data abstraction

- dataset types
 - network
 - nodes: genes
 - links: known interactions between genes
 - table
 - quantitative attributes
 - gene expression levels for nodes across different experimental conditions
 - categorical attributes
 - subcellular location of interaction
 - functional groups



Why: Task abstraction

- task: interpret experiment results with respect to gene network
 - goal: accelerate existing discovery workflow
 - compare distributions between attributes
 - different experiments
 - interpret attributes in context of known network structure



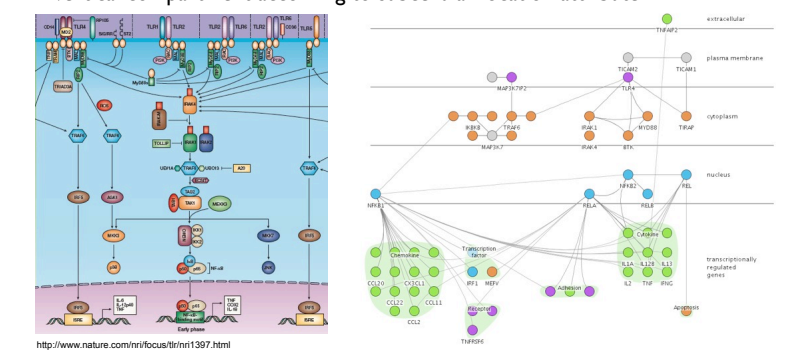
How: Idiom design decisions

- arrange space for networks
 - custom node-link diagram layout
 - points for nodes
 - connection marks for links
 - vertical compartment according to subcellular location attribute
 - cluster according to functional grouping



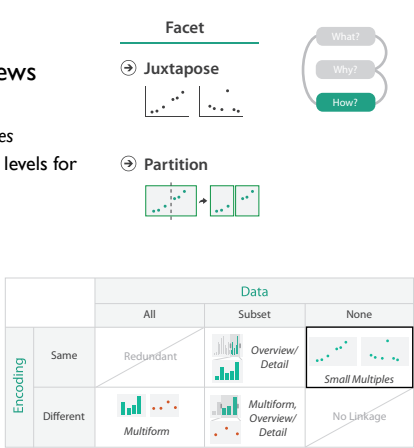
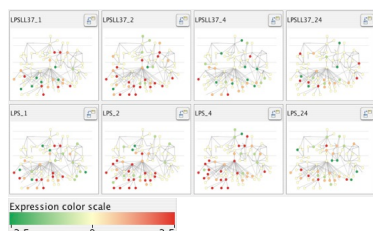
How: Arrange space

- automatic layout similar to hand-drawn diagrams
 - vertical compartment according to subcellular location attribute



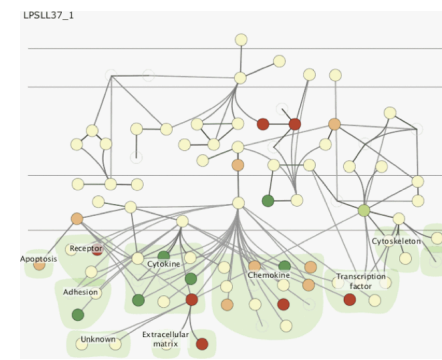
How: Idiom design decisions

- facet: partition data into multiple views
 - juxtapose views side by side
 - same encoding, different data: small multiples
 - nodes in each view colored by expression levels for experimental condition



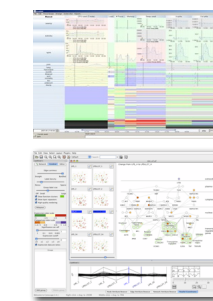
How: Juxtapose vs. animate

- comparison difficult across many frames with many changes everywhere
- rule of thumb: eyes beat memory
 - principle: external cognition vs. internal memory
 - easy to compare by moving eyes between side-by-side views
 - harder to compare memory of what you saw to visible view



Outline

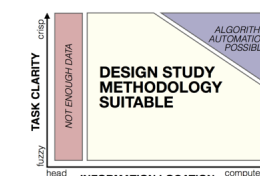
- interactive visual analysis
 - role and advantages
- LiveRAC
 - time-series data: managed web hosting (with AT&T)
- Cerebral
 - network of relationships: genes (with Agilent and UBC Immunology)
- wrapup



Design Study Methodology

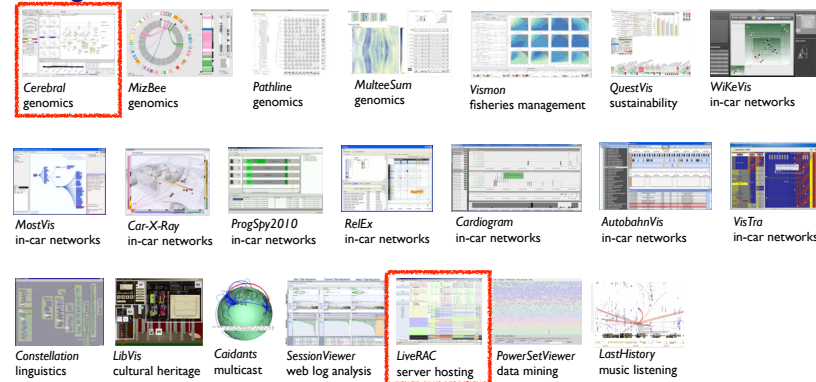
Reflections from the Trenches and from the Stacks

joint work with:
Michael Sedlmair, Miriah Meyer
<http://www.cs.ubc.ca/labs/imager/tr/2012/dsm/>



Design Study Methodology: Reflections from the Trenches and from the Stacks. Sedlmair, Meyer, Munzner. IEEE Trans. Visualization and Computer Graphics 18(12): 2431-2440, 2012. (Proc. InfoVis 2012).

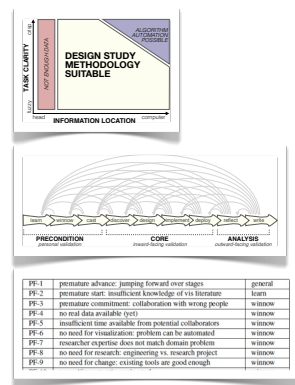
Design Studies: Lessons learned after 21 of them



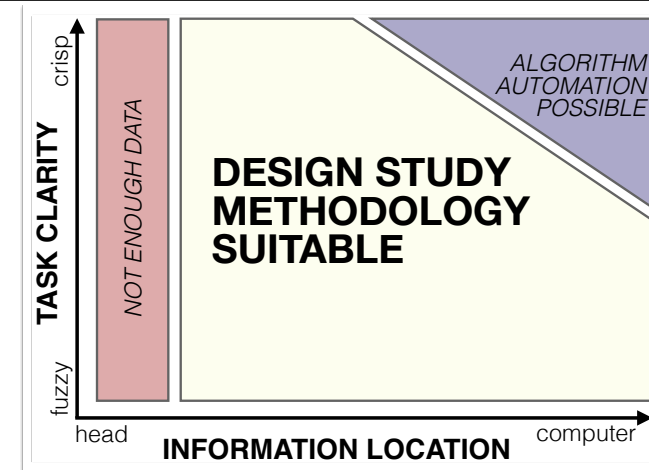
33

Methodology for Problem-Driven Work

- definitions
- 9-stage framework
- 32 pitfalls
– and how to avoid them



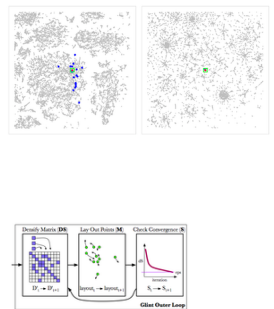
34



35

Techniques: Dimensionality Reduction

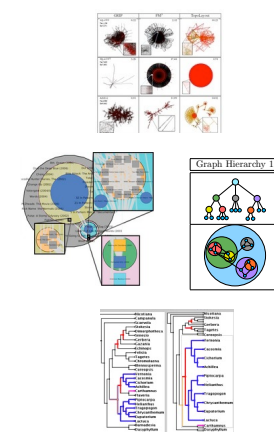
- reducing high-dimensional data to tractable low-dimensional form
- Q-SNE: high-quality clusters for millions of documents
- Glint: costly distance functions
 - incl. preferences elicited from people



36

Techniques: Networks & Trees

- large multi-level networks
 - layout
 - TopoLayout
 - interaction
 - Grouse
 - GrouseFlocks
 - TugGraph
- large tree comparison
 - TreeJuxtaposer



37

More Information

- this talk
<http://www.cs.ubc.ca/~tmm/talks.html#disney15>
- papers, videos, software, talks, courses
<http://www.cs.ubc.ca/group/infovis>
<http://www.cs.ubc.ca/~tmm>
- book
<http://www.cs.ubc.ca/~tmm/vadbook>
- acknowledgements
– funding: Agilent, AT&T, NSERC, NSF



@tamaramunzner

38