

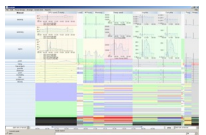
# Visualization for Hackers: Why It's Tricky, and Where to Start

**Tamara Munzner**  
 Department of Computer Science  
 University of British Columbia

Hackers on Planet Earth (HOPE) X  
 19 July 2014, New York NY  
<http://www.cs.ubc.ca/~tmm/talks.html#hope14>

## Outline

- introduction
  - what's vis anyway?
- LiveRAC
  - server logs: managed web hosting (with AT&T)
- Overview
  - text: visual document mining for journalists (with Associated Press)
- big picture and wrapup



2

## Defining visualization (vis)

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

Why?...

## Why have a human in the loop?

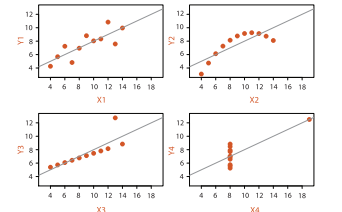
Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

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

- many analysis problems ill-specified, not clear what questions to ask in advance
  - don't need vis when fully automatic solution exists and is trusted

### Anscombe's Quartet

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

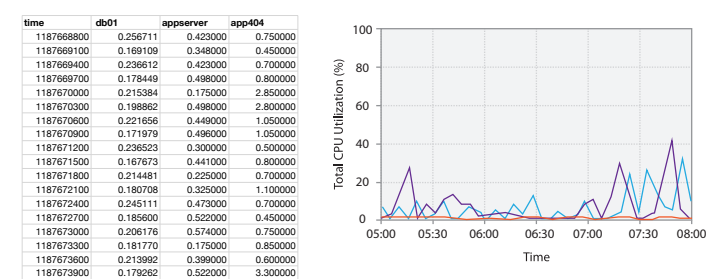


4

## Why use an external representation?

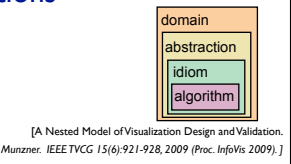
Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

- external representation: replace cognition with perception

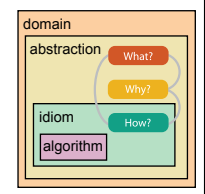


## Analysis framework: Four levels, three questions

- domain situation
  - who are the target users?
- abstraction
  - translate from specifics of domain to vocabulary of vis
  - what is shown? **data abstraction**
  - why is the user looking at it? **task abstraction**
- idiom
  - how is it shown?
    - visual encoding idiom: how to draw
    - interaction idiom: how to manipulate
- algorithm
  - efficient computation



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



[A Multi-Level Typology of Abstract Visualization Tasks. Brehmer and Munzner. IEEE TVCG 19(12):2376-2385, 2013 (Proc. InfoVis 2013).]

6

## Why analyze?

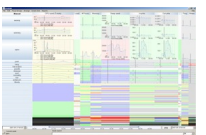
- huge design space
  - visual encoding: combinatorial explosion of choices
  - add interaction: even bigger
  - add data abstraction transformation: truly enormous
- most possibilities ineffective for particular task/data combination
  - implication: avoid random walk, be guided by principles
- analysis framework: scaffold to think systematically about design space
  - ensure that consideration space encompasses full scope of possibilities
  - improve chances that selected solution is good not mediocre
  - today's focus: abstractions and idioms, what-why-how



7

## Outline

- introduction
  - what's vis anyway?
- LiveRAC
  - server logs: managed web hosting (with AT&T)
- Overview
  - text: visual document mining for journalists (with Associated Press)
- big picture and wrapup

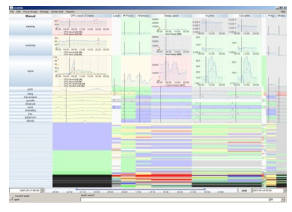


8

## LiveRAC

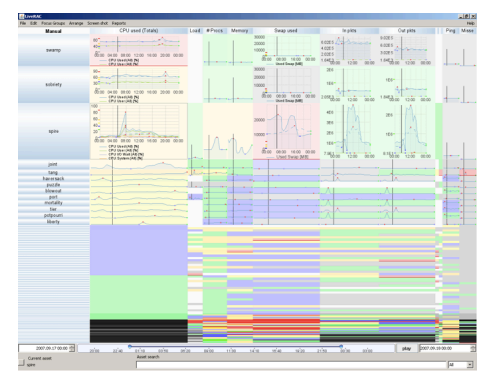
Interactive Visual Exploration of System Management Time-Series Data

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



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

## LiveRAC video

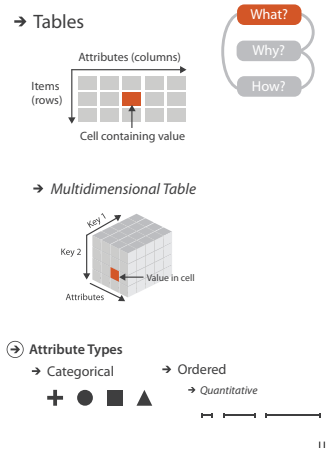


<http://youtu.be/Id0c3H0VSkw>

10

## 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



11

## Why: Tasks in domain language

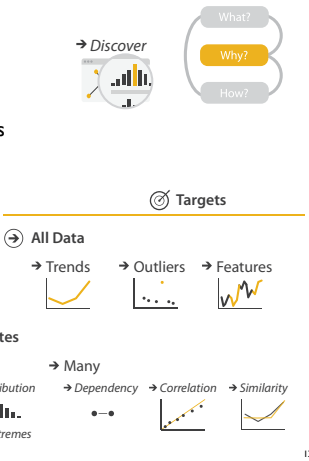
- interpret network environment status
- report generation
- capacity planning
- event investigation/forensics
- coordination
  - between customers, engineering, ops



12

## Why: Task abstraction

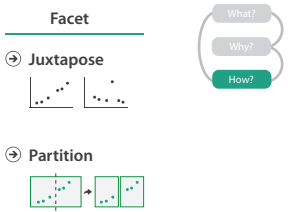
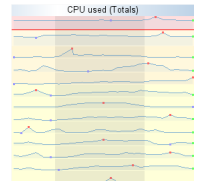
- browse and correlate across combinations of parameter, device, time
  - correlate alarm attribute with other parameter attribs
  - find trends across groups of devices
  - summarize over different time intervals
  - identify devices at or beyond parameter thresholds
  - identify critical parameter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



13

## How: Facet

- facet: partition data into multiple views
  - juxtapose views side by side
    - same encoding, different data: small multiples

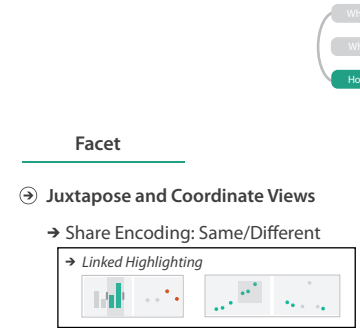
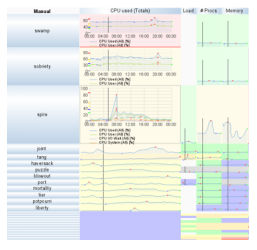


		Data		
		All	Subset	None
Encoding	Same	Redundant	Overview/Detail	Small Multiples
	Different	Multiform	Multiform, Overview/Detail	No Linkage

14

## How: Juxtapose

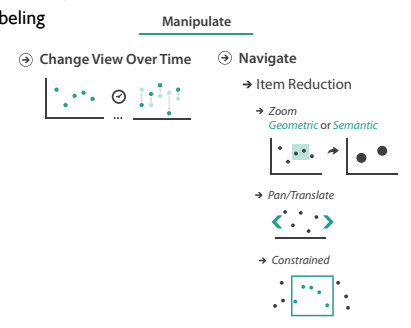
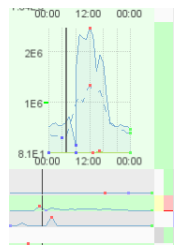
- juxtapose linked views
  - linked highlighting
    - marker line tracks across views



15

## How: Navigate

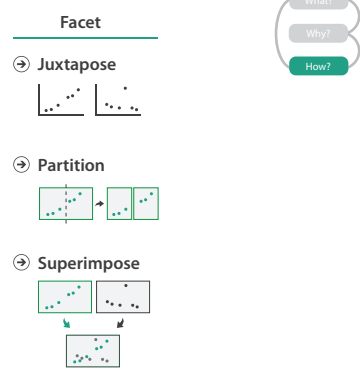
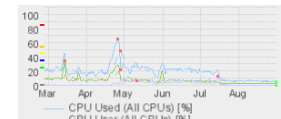
- semantic zooming
  - representation adapts to pixels available for object
    - many: superimposed line charts with full labeling
    - some: iconic line chart (sparkline)
    - few: color-coded box (heatmap)



16

## How: Superimpose

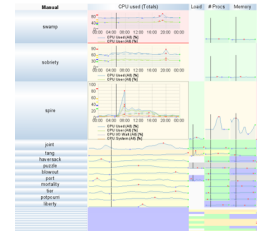
- superimpose layers
- vs juxtapose side by side



17

## How: Reduce

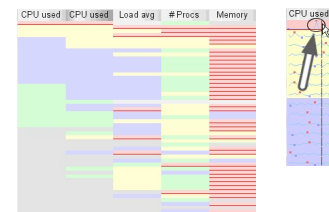
- reduce data shown with complex combination of filtering and aggregation
- embed focus+context in single view
- distort geometry
  - metaphor: stretch and squish navigation
  - shape: rectilinear
  - foci: multiple
  - impact: global



18

## How: Reordering

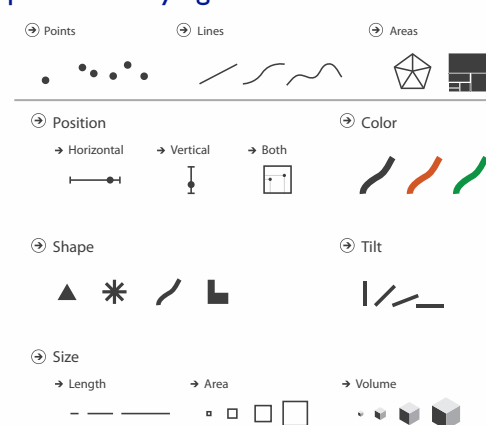
- change spatial arrangement
  - resort by selected attribute
  - check for correlations between aligned attribute columns
  - ex: high load without high CPU, maybe I/O bound



19

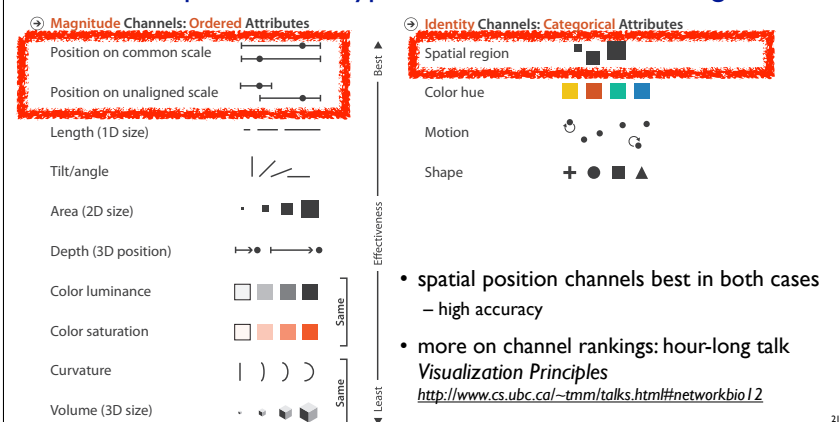
## Importance of arranging space: Underlying definitions

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



20

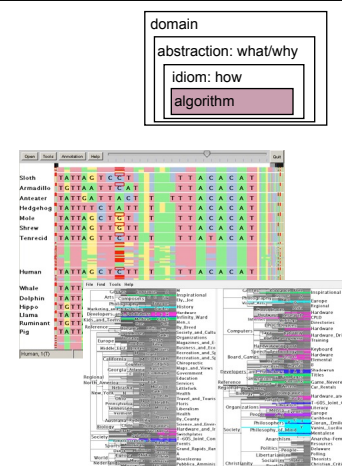
## Channels: Expressiveness types and effectiveness rankings



21

## 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

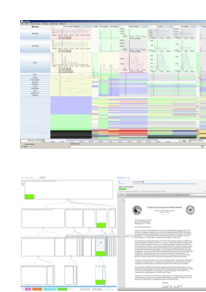


[Partitioned Rendering Infrastructure for Scalable Accordion Drawing (Extended Version). Slack, Hildebrand, and Munzner. Information Visualization, 5(2), p. 137-151, 2006.]

22

## Outline

- introduction
  - what's vis anyway?
- LiveRAC
  - server logs: managed web hosting (with AT&T)
- Overview
  - text: visual document mining for journalists (with Associated Press)
- big picture and wrapup



23

## Overview

The Design, Adoption, and Analysis of a Visual Document Mining Tool For Investigative Journalists

joint work with:  
Matt Brehmer, Stephen Ingram, Jonathan Stray  
<http://www.cs.ubc.ca/labs/imager/tr/2014/Overview/>  
<https://www.overviewproject.org/>

Overview: The Design, Adoption, and Analysis of a Visual Document Mining Tool For Investigative Journalists. Brehmer, Ingram, Stray, and Munzner. IEEE Trans. Visualization and Computer Graphics (Proc. InfoVis 2014), to appear.

24

## Origin story: WikiLeaks meets Glimmer

- WikiLeaks: hacker-journalist Jonathan Stray analyzing Iraq warlogs
  - conjecture that existing label classification falls short of showing all meaningful structure in data
  - friendly action, criminal incident, ...
  - had some NLP, needed better vis tools



- Glimmer: multilevel dimensionality reduction algorithm
- scalability to 30K documents and terms

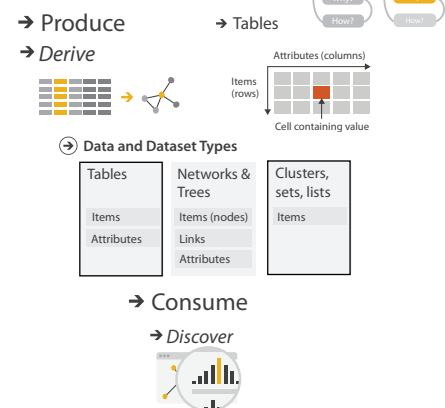
[Glimmer: Multilevel MDS on the GPU. Ingram, Munzner, Olano. IEEE TVCG 15(2):249-261, 2009.]



25

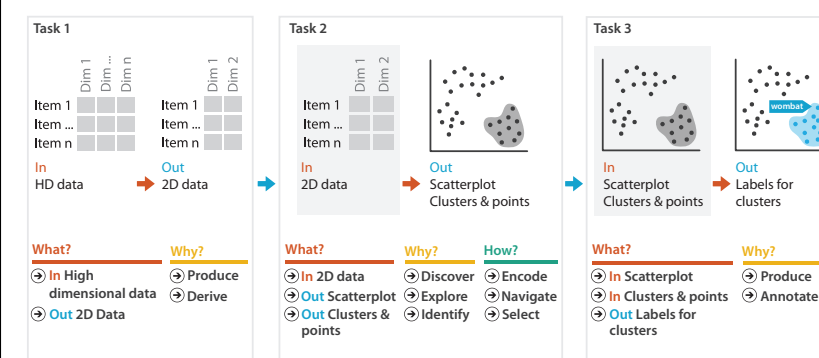
## What: Data and task abstraction

- derive data to transform text into visualizable dataset
  - from documents to high-dimensional table
  - bag of words model
    - attribute: any word that appears across entire collection
    - document/item: word counts (sparse)
  - from high-dimensional table to low-dimensional table
  - synthesize new dimensions that capture most of high-dim proximity structure
  - find clusters of items in lowD space
    - discover: generate or verify



26

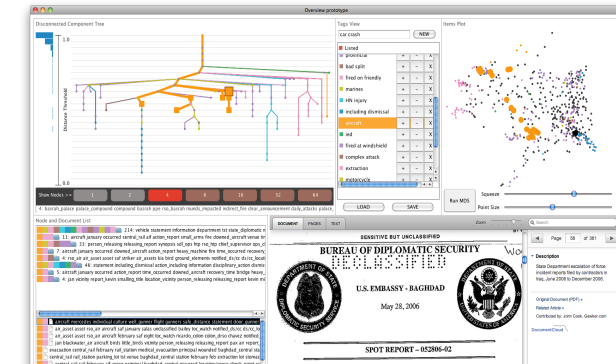
## Dimensionality reduction for document datasets



- more on DR: hour-long talk Dimensionality Reduction from Several Angles  
<http://www.cs.ubc.ca/~tmm/talks.htm#linz14>

27

## Overview video (version 1)

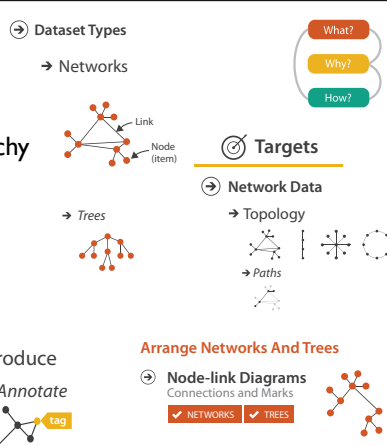


<http://www.cs.ubc.ca/labs/imager/tr/2012/modiscotag>

28

## What/Why/How interplay

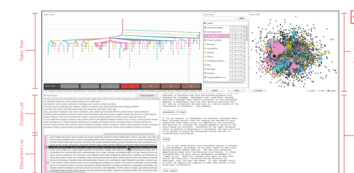
- why: understand clusters
- what: derive data of full cluster hierarchy
  - explore space of possible clusterings
- how: show cluster hierarchy
  - arrange space: node-link
- how: support tagging clusters/docs
  - following or cross-cutting hierarchy!
  - simple annotation
  - progress tracking
  - user-defined semantics



29

## How: Idiom design decisions

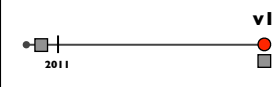
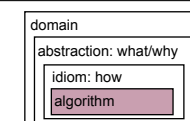
- facet: juxtapose linked views
  - linked color coding
    - cluster hierarchy tree
    - DR scatterplot
    - tags
  - reading text/keywords
    - cluster list
    - doc reader



30

## Algorithm

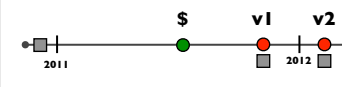
- version 1
  - fast cluster hierarchy construction for sparse data
  - research prototype by PhD student
  - positive initial assessment from AP Caracas bureau chief
  - barrier to adoption: difficult install/load process



31

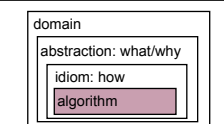
## Algorithm

- version 1
  - fast cluster hierarchy construction for sparse data
  - research prototype by PhD student
  - positive initial assessment from AP Caracas bureau chief
  - barrier to adoption: difficult install/load process
- version 2
  - web deployment, DocumentCloud integration, usability
  - many months of engineering
    - Knight Foundation funding to the rescue!
  - published story by unaffiliated reporter: police corruption in Tulsa

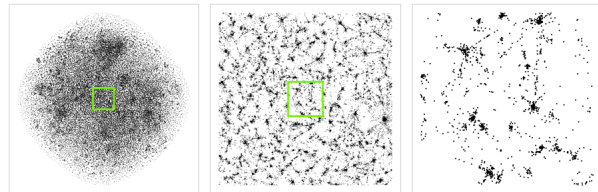


32

# Algorithm: Spinoff series

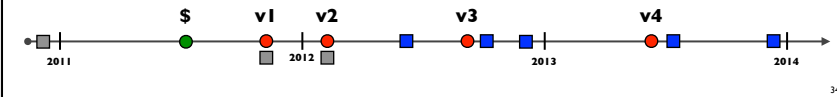


- dimensionality reduction for huge text collections
  - great algorithm problem in its own right!
  - QSNE: fast and high-quality DR for millions of documents
    - key feature: handle sparseness appropriately

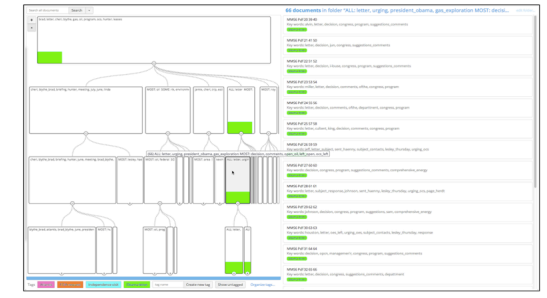


# Path to adoption

- even more rounds of what/why/how interplay
  - which views needed? what should they show? how should they show it?
  - usability and utility
- version 3
  - published story: VP candidate Ryan asked for federal help even as championed cuts
  - published story: gun control debate
- version 4
  - followup investigation: government corruption in Texas
  - published story: police corruption in New York (*Pulitzer prize finalist!*)



# Overview v4 video



- versions 3 and 4
  - no DR scatterplot
  - tree arrangement emphasizing nodes not links
  - combined doc/cluster viewer

# Why: Task abstractions revisited

- what's in this collection? (of leaked docs)
  - generate hypothesis
  - summarize clusters
  - explore clusters
- locate evidence (within FOIA dump)
  - verify hypothesis
  - identify clusters/documents
  - locate clusters/documents
- prove non-existence of evidence
  - even harder!
  - exhaustive reading vs filtering out irrelevant

Discover

Query

Identify

Compare

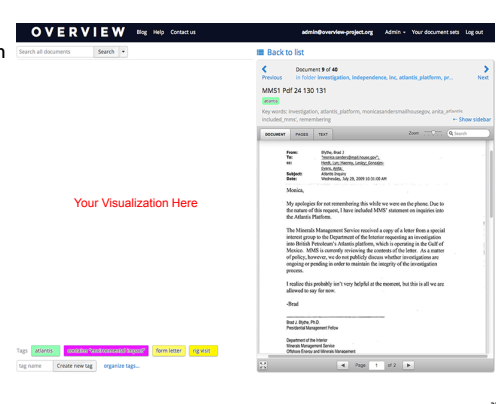
Summarise

	Target known	Target unknown
Location known	Lookup	Browse
Location unknown	Locate	Explore

[A Multi-Level Typology of Abstract Visualization Tasks. Brehmer and Munzner. IEEE TVCG 19(12):2376-2385, 2013 (Proc. InfoVis 2013).]

# Now what?

- continuing adoption
  - food stamp distribution delays in North Carolina
  - credit card agreements allow repossession
  - this week
    - The Brilliance of Louis C.K.'s Emails: He Writes Like a Politician
- continuing development
  - Knight Foundation funds v5
    - named entity recognition
    - plugin API

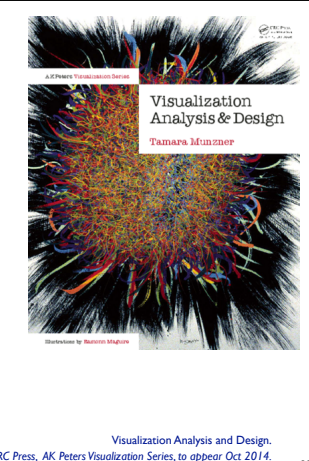


# Outline

- introduction
  - what's vis anyway?
- LiveRAC
  - server logs: managed web hosting (with AT&T)
- Overview
  - text: visual document mining for journalists (with Associated Press)
- big picture and wrapup



# Visualization Analysis & Design



<http://www.cs.ubc.ca/~tmm/vadbook>

Munzner. Taylor and Francis / CRC Press, AK Peters Visualization Series, to appear Oct 2014.

What?

Why?

How?

**What?**

Attribute Types

- Categorical
- Ordered
- Ordinal
- Quantitative

Ordering Direction

- Sequential
- Diverging
- Cyclic

**How?**

Dataset Availability

- Static
- Dynamic

What?

Why?

How?

**Why?**

Actions

- Analyze
  - Consume
  - Discover
  - Present
  - Enjoy
  - Produce
  - Annotate
  - Record
  - Derive
- Search
  - Target known: Lookup, Browse
  - Target unknown: Locate, Explore
- Query
  - Identify
  - Compare
  - Summarise

**Targets**

- All Data
  - Trends
  - Outliers
  - Features
- Attributes
  - One
  - Many
  - Distribution
  - Dependency
  - Correlation
  - Similarity
  - Extremes
- Network Data
  - Topology
- Spatial Data
  - Shape

What?

Why?

How?

**How?**

Encode

- Arrange
  - Express
  - Separate
  - Order
  - Align
  - Use
- Map
  - from categorical and ordered attributes
  - Color
    - Hue
    - Saturation
    - Luminance
  - Size, Angle, Curvature, ...
  - Shape
    - Direction, Rate, Frequency, ...

**Manipulate**

- Change
- Select
- Navigate

**Facet**

- Juxtapose
- Partition
- Superimpose

**Reduce**

- Filter
- Aggregate
- Embed

# Channels: Expressiveness types and effectiveness rankings

**Magnitude Channels: Ordered Attributes**

- Position on common scale
- Position on unaligned scale
- Length (1D size)
- Tilt/angle
- Area (2D size)
- Depth (3D position)
- Color luminance
- Color saturation
- Curvature
- Volume (3D size)

**Identity Channels: Categorical Attributes**

- Spatial region
- Color hue
- Motion
- Shape

Effectiveness: Best (top) to Least (bottom)

# Four levels of design

inverse cases: technique-driven vs. problem-driven work

– both useful, but learning curve to switch between

**Domain situation** → **problem-driven work**

↓

**Data/task abstraction**

↓

**Visual encoding/interaction idiom**

↕

**Algorithm** → **technique-driven work**

# 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 Studies: Lessons learned after 21 of them (+more)

# Methodology for Problem-Driven Work

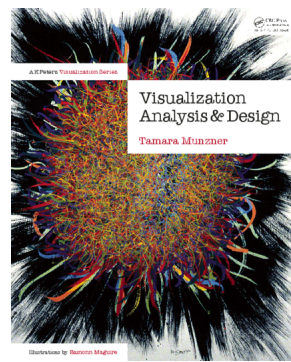
- definitions
- 9-stage framework
- 32 pitfalls and how to avoid them

# Wrapup

- two systems analyzed
  - LiveRAC, Overview
- analysis framework big ideas
  - what: data abstraction
    - characterize and derive data
  - why: task abstraction
    - translate from domain-specific to generic
  - how: visual encoding and interaction idioms
    - separate from questions of algorithm design
- scaffolding for thinking systematically about full design space
  - describing existing systems helps with generating new ones

## More Information

- this talk  
<http://www.cs.ubc.ca/~tmm/talks.html#hope14>
- papers, videos, software, talks, courses  
<http://www.cs.ubc.ca/group/infovis>  
<http://www.cs.ubc.ca/~tmm>
- book (to appear Oct 2014)  
<http://www.cs.ubc.ca/~tmm/vadbook>
- acknowledgements  
– funding: AT&T, Knight Foundation, NSERC  
– talk feedback: Matt Brehmer



Visualization Analysis and Design.  
Munzner. Taylor and Francis / CRC Press, AK Peters Visualization Series, to appear Oct 2014.