

# Data Visualization in Genomics and In-Car Network Engineering

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

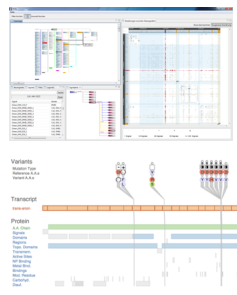
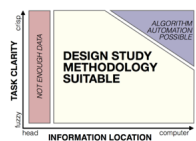
City University of London, Computer Science Department Seminar

1 July 2014, London UK

<http://www.cs.ubc.ca/~tmm/talks.html#london14>

## Outline

- Design Study Methodology  
 – meta-paper: how to do design studies
- RelEx  
 – overlay network optimization for in-car networks
- Variant View  
 – sequence variant analysis in gene context



## Defining Visualization

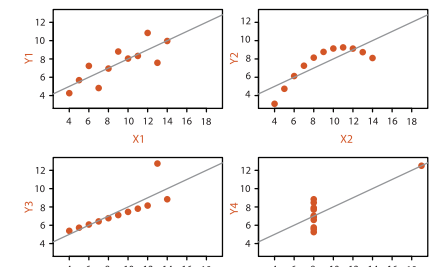
Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

## Defining Visualization

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

- human in the loop needs the details  
 – doesn't know exactly what questions to ask in advance

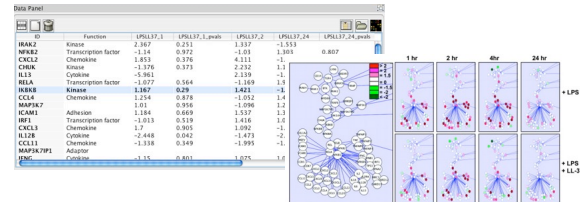
Identical statistics	
x mean	9.0
x variance	10.0
y mean	7.50
y variance	3.75
x/y correlation	0.816



## Defining Visualization

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

- human in the loop needs the details  
 – doesn't know exactly what questions to ask in advance
- external representation: replace cognition with perception



## Defining Visualization

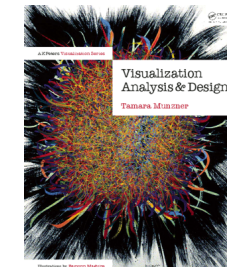
Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

- human in the loop needs the details  
 – doesn't know exactly what questions to ask in advance
- external representation: perception vs cognition
- intended task

## Defining Visualization

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

- human in the loop needs the details  
 – doesn't know exactly what questions to ask in advance
- external representation: perception vs cognition
- intended task
- measurable definitions of effectiveness



more at:  
 Visualization Analysis and Design, Chapter 1.  
 Munzner, AK Peters, 2014, to appear.

# 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 TVCG 18(12):2431-2440, 2012 (Proc. InfoVis 2012).

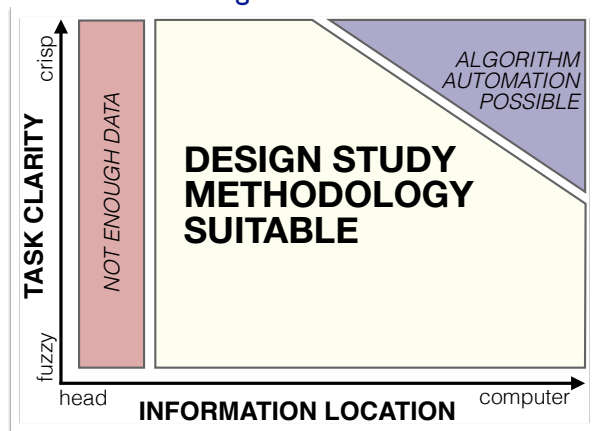
## Defining Design Study

- a specific **real-world** problem  
 – real users and real data,  
 – collaboration is (often) fundamental
- **design** a visualization system  
 – implications: requirements, multiple ideas
- **validate** the design  
 – at appropriate levels
- **reflect** about lessons learned  
 – transferable research: improve design guidelines for vis in general  
 • confirm, refine, reject, propose

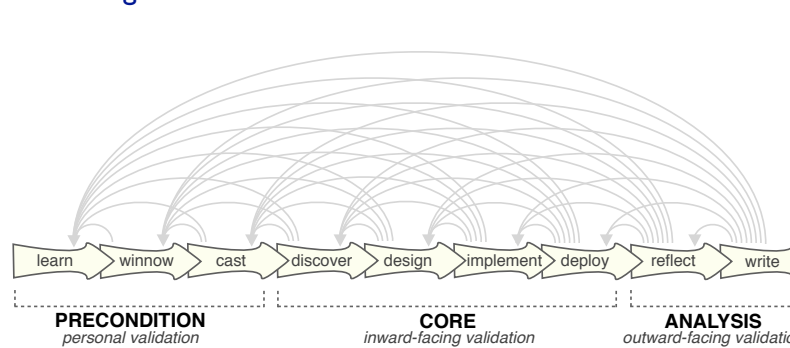
more at:  
 A Nested Model of Visualization Design and Validation.  
 Munzner. IEEE TVCG 15(6):921-928, 2009 (Proc. InfoVis 2009).

more at:  
 The Nested Blocks and Guidelines Model.  
 Meyer, Sedlmair, Quinan, Munzner. Information Visualization Journal, 2014,  
 to appear.

## When To Do Design Studies

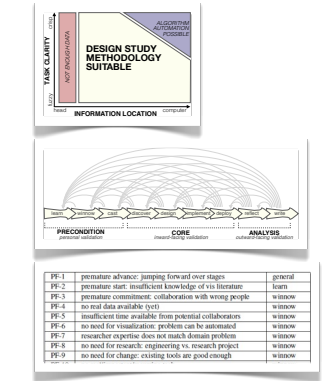


## Nine-Stage Framework



## How To Do Design Studies

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



## Pitfall Example: Premature Publishing

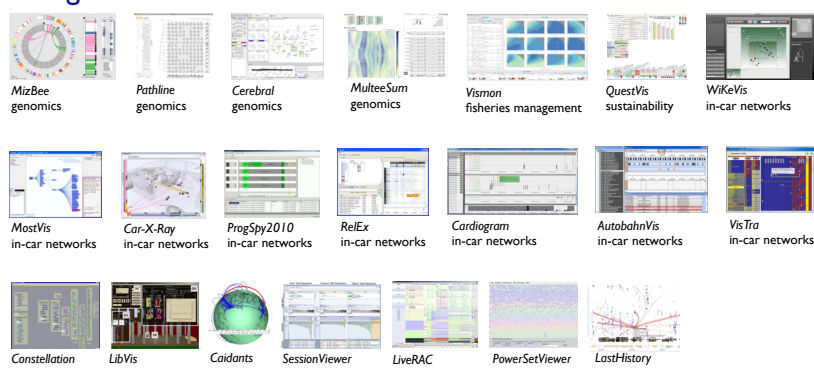
algorithm innovation      design studies

Must be first!

Am I ready?



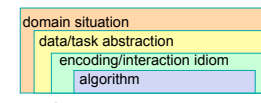
## Design Studies: Lessons learned after 21 of them



• commonality of representations cross-cuts domains!

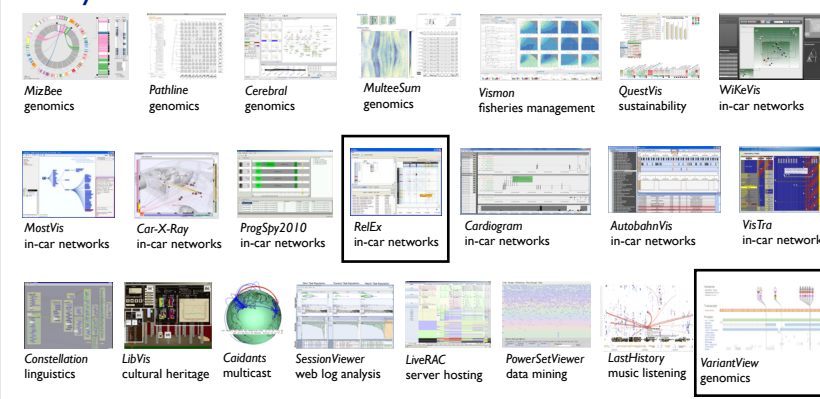
## Abstractions and Idioms

- abstractions  
 – **translate** from specifics of domain to vocabulary of vis  
 • task abstraction: **why** they're looking at it  
 • data abstraction: **what** to draw  
 – **transform** data into form useful for task at hand  
 • don't just draw what you're given; decide what is the right thing!
- idioms  
 – visual encoding idiom: **how** to draw  
 – interaction idiom: **how** to manipulate
- focus today: two mappings  
 – from domain to abstraction  
 – from abstraction to idiom



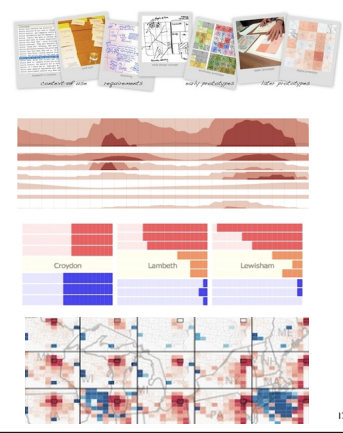
A Nested Model of Visualization Design and Validation.  
 Munzner. IEEE TVCG 15(6):921-928, 2009 (Proc. InfoVis 2009).

## Today's Focus



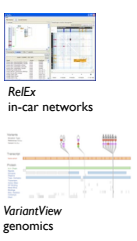
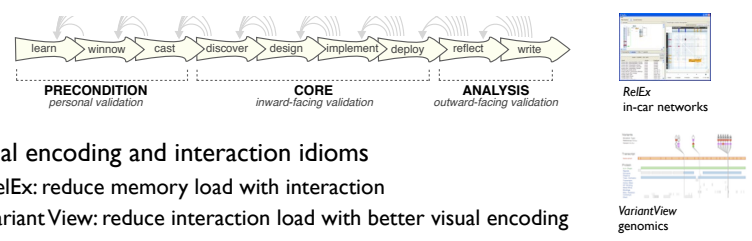
# Design Studies: giCentre Context

- methodology  
Human-centered approaches in geovisualization design: investigating multiple methods through a long-term case study. Lloyd and Dykes. IEEE Transactions on Visualization and Computer Graphics, 17(12):2498-2507, 2011.
- energy analysis  
Creative user-centered visualization design for energy analysts and modelers. Goodwin, Dykes, Jones, Dillingham, Dove, Duffy, Kachkoey, Slingsby, Wood. IEEE Transactions on Visualization and Computer Graphics, 19(12), pp. 2516-2525, 2013.
- BallotMaps  
BallotMaps: Detecting name bias in alphabetically ordered ballot papers. Wood, Badawood, Dykes, Slingsby. IEEE Transactions on Visualization and Computer Graphics, 17(12), pp. 2384-2391, 2011
- ODMaps  
Visualisation of origins, destinations and flows with OD maps. Wood, Dykes, Slingsby. The Cartographic Journal, 47(2), pp. 117-129, 2010.



# Themes

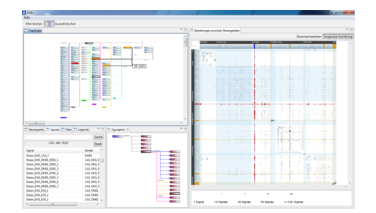
- task and data abstraction
  - both cases: complex and tricky
  - clear description in final talk/paper is end of a long, long road
  - writing as research: refine during reflection even after vis tool is finalized...
- visual encoding and interaction idioms
  - RelEx: reduce memory load with interaction
  - Variant View: reduce interaction load with better visual encoding



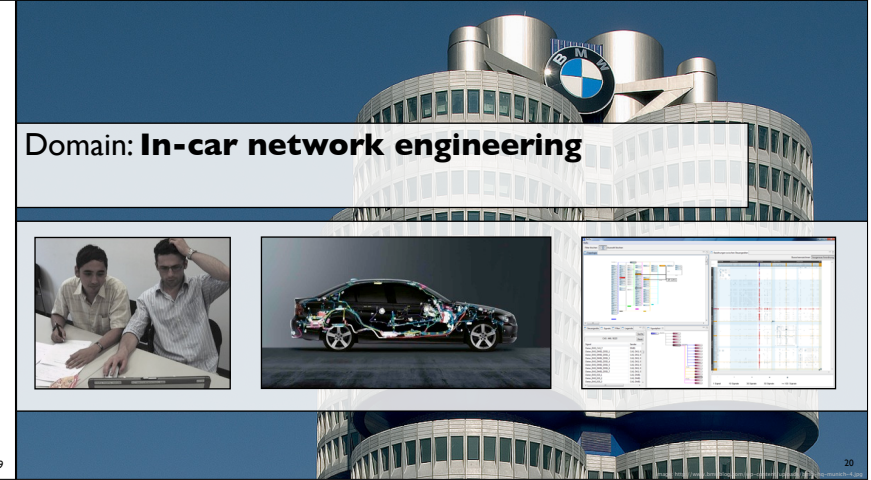
# RelEx

Visualization for Actively Changing Overlay Network Specifications

joint work with:  
Michael Sedlmair, Annika Frank, Andreas Butz  
<http://www.cs.ubc.ca/labs/imager/tr/2012/relex/>



RelEx: Visualization for Actively Changing Overlay Network Specifications. Sedlmair, Frank, Butz, Munzner. IEEE TVCG 18(12): 2729-2738, 2012 (Proc. InfoVis 2012).



# Domain: In-car network engineering

## Abstractions

21

## DATA In-car Electronics

22

## Data Abstraction: 3 Networks

- physical network
  - 100 nodes: Electronic Control Units
  - 10-15 hyperedges: bus systems
  - hardware engineers
- logical network
  - same nodes
  - 10,000 multigraph edges: signals
  - 1,000 weighted edges: signal counts
  - software engineers
- overlay network
  - maps logical onto physical
  - 30,000 edges: signal paths
  - target engineers

23

## Task Abstraction: Mapping

- specify overlay network that maps logical onto physical

24

## Task Abstraction: Optimizing

- traffic optimization

Many constraints  
bandwidth ... delay/real time ...  
path length ... load balance ...  
reliability ... money ...

— engineer, BMW —

25

## Task Abstraction: Changing

- external change requests

Change  
(trivial requests might lead to complex changes)

26

## Idioms

27

## RELEX: Relation Explorer

28

## Vis Guideline [Ghoniem 2005] Matrix for dense network data

**SIGNAL COUNT NETWORK**

visual encoding: size-coded matrix

**Logical Network View: Overview**

29

## Vis Guideline [Ghoniem 2005] Node-link for path following tasks

**SIGNAL PATH NETWORK**

filtered by signal

**Signal Path View: Selected Signal**

30

## INTERACTION IDIOM: Cross-Network Relations

linked highlighting

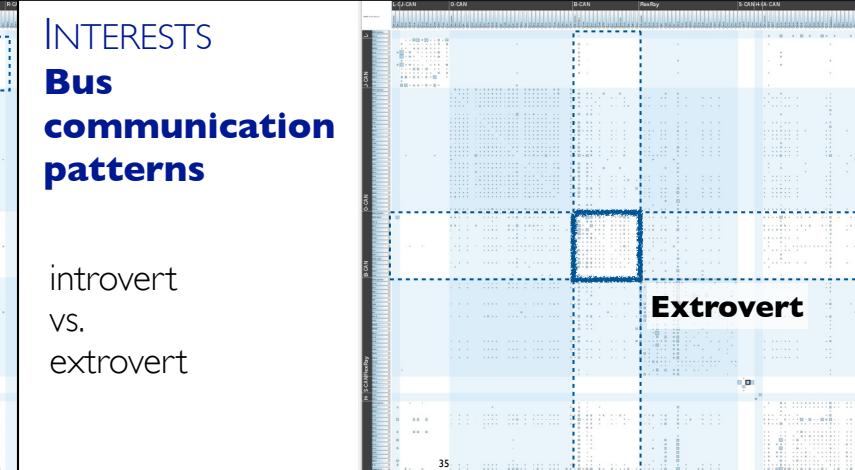
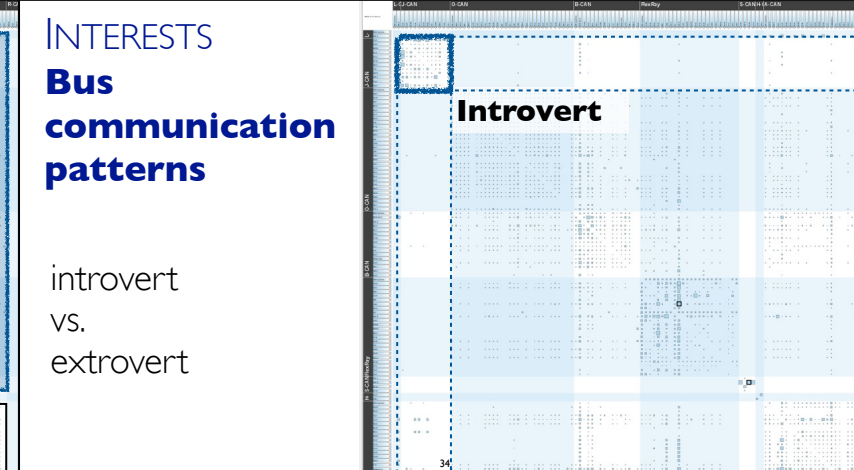
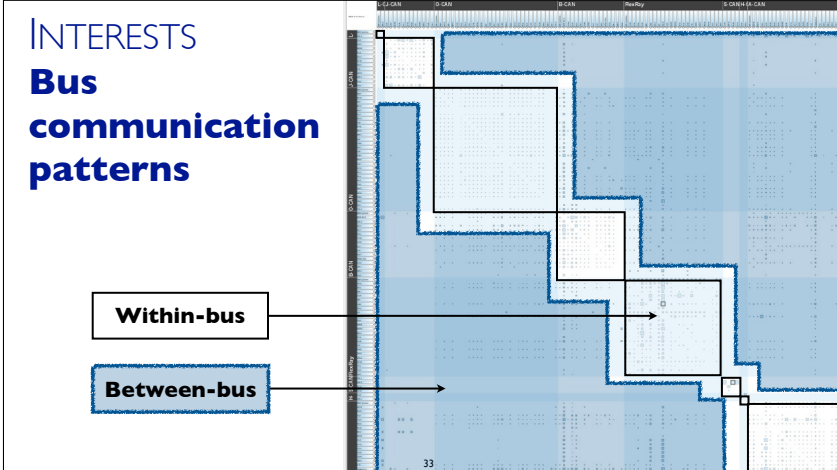
VIDEO

31

## INTERESTS Bus communication patterns

CAN Bus

32



## Methods

### Phase 1: Discover

3 months

- embedded within BMW
  - phases 1, 2, 3
- contextual inquiry
- abstracting
- deriving design requirements

37

### Phase 2: Design, implement, deploy

4 months

- iterative paper prototyping
- agile software development
  - 3 lead users (engineers)
  - 6 deployed releases
- usability engineering
  - domain experts
  - HCI students

38

### Phase 3: Summative evaluation

2 months

- field study
  - 7 engineers
  - 5 weeks
- think aloud study
  - 10 engineers
  - ~1 hour each session
- adoption
  - 15+ users, 3 months post-study

39

### Phase 4: Reflect and write

3 months

- revisit abstractions
- relate to other design studies
- write up

40

## Abstraction Innovation

### Previous Work

#### Focus on social network analysis

- radically different task and data abstractions

41

### Task Abstraction

#### Social Network Analysis Domain

- find clusters

42

### Task Abstraction

#### Social Network Analysis Domain

- find clusters
- find high-degree nodes

43

### Task Abstraction

#### Social Network Analysis Domain

- find clusters
- find high-degree nodes
- find bridge nodes

44

### Task Abstraction

#### Social Network Analysis Domain

- find clusters
- find high-degree nodes
- find bridge nodes
- understand temporal dynamics
  - passively notice changes

45

### Data Abstraction

#### Social Network Analysis Domain

- single graph

46

### Data Abstraction

#### Social Network Analysis

- single graph
- scalability challenge: nodes

47

### Social Network Analysis vs Overlay Network Optimization

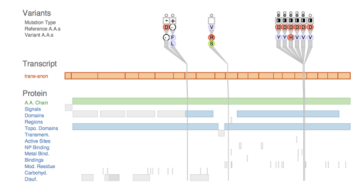


- data
    - single network
    - node scalability
      - sparse edges
  - task
    - find clusters, high-degree nodes, bridge nodes
    - passive changes
- data
    - three related networks
      - physical, logical, overlay
    - path scalability
    - dense edges, few nodes
  - task
    - traffic optimization
    - active changes

## Variant View

Visualizing Sequence Variants in their Gene Context

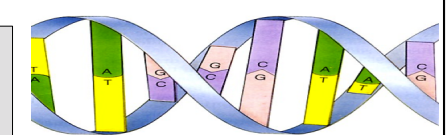
joint work with: Joel Ferstay, Cydney Nielsen  
<http://www.cs.ubc.ca/labs/imager/tr/2012/VariantView/>



Variant View: Visualizing Sequence Variants in their Gene Context. Ferstay, Nielsen, Munzner. IEEE TVCG 19(12): 2546-2555, 2013 (Proc. InfoVis 2013).

### Sequence Variant Definition

- Sequence variants
  - Difference between reference and given genome



Reference Genome DNA: ATA TGA TCA ACA CTT

Sample 1 Genome DNA: ATA TGG TCA ATA CTT **Harmful?**

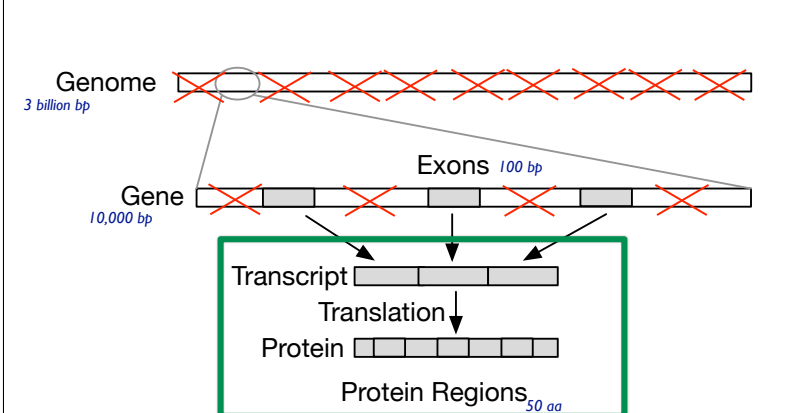
Sample 2 Genome DNA: ATA TGA TGA ACA CCT **Harmless?**

### Cancer Research

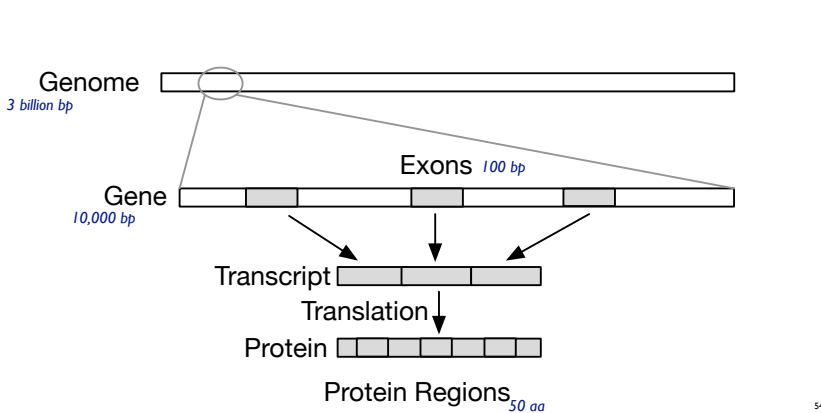
- collaboration with analysts at BC Genome Sciences Center
  - studying genetic basis of leukemia
- driving task
  - discover new candidate genes with harmful variants
- two big questions
  - what to show
    - data abstraction
    - challenge: enormous range of scales in the data
  - how to show it
    - visual encoding idiom

## Abstractions

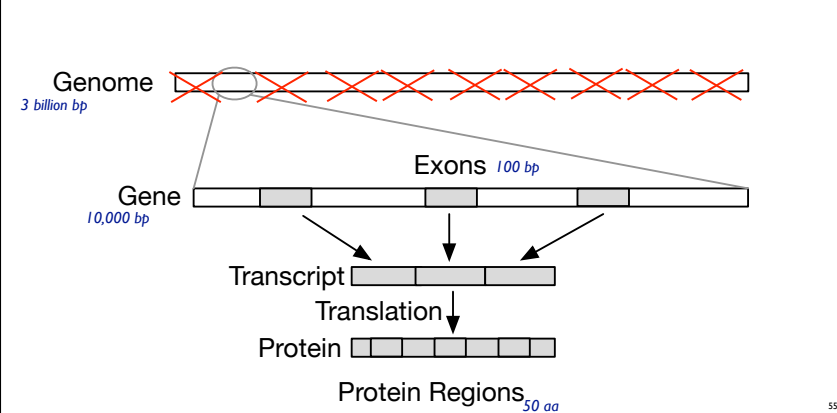
### Data abstraction: highly filtered scope of transcript coordinates



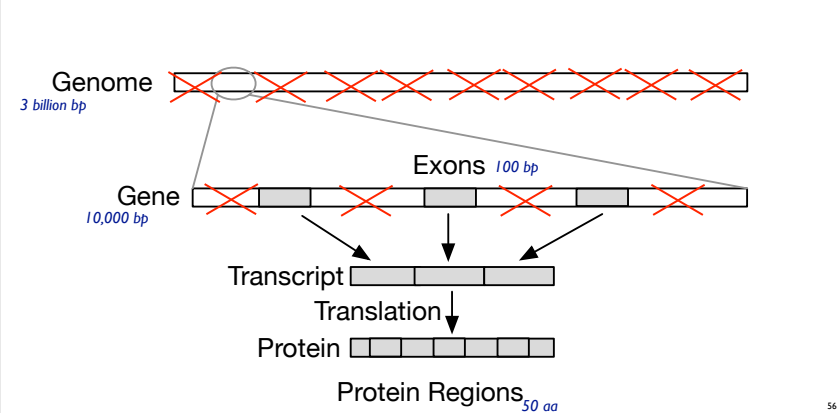
### Data: Filtering to relevant biological levels and scales



### Filter out whole genome; keep genes

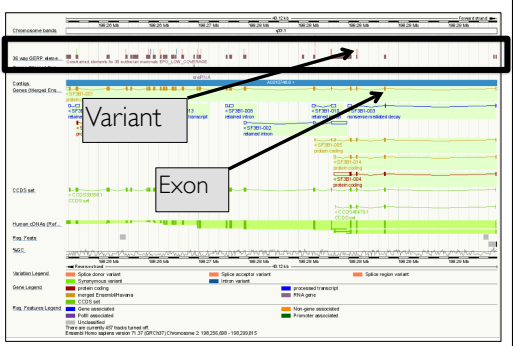


### Filter out non-exon regions



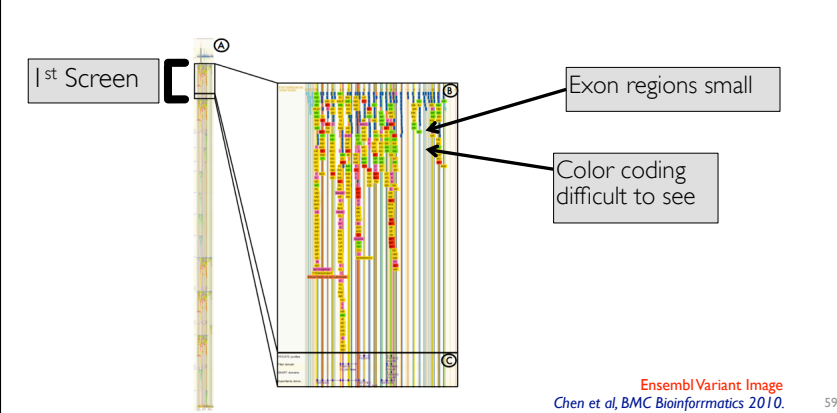
### Dominant paradigm: genome browsers

- strengths: flexible and powerful
  - horizontal tracks: user data
  - shared coordinate system: genome coordinates (bp)
- problems
  - tiny features of interest spread out across large extent
    - must zoom far in to inspect known feature, then zoom out and pan to locate next
  - high cognitive load for interaction
  - must already know where to look



representative example: Ensembl  
Chen et al, BMC Bioinformatics 2010.

### Features of interest small even in variant-specific view



Ensembl Variant Image  
Chen et al, BMC Bioinformatics 2010.

## Idioms

### Variant View

Gene Search:  Submit

Sort By: Gene

Alternative Transcripts: gene-exon (trans-annot)

Variants

Transcript

Protein

Variant Data

### Variant View

Gene Search:  Submit

Sort By: Gene

Alternative Transcripts: gene-exon (trans-annot)

Variants

Transcript

Protein

Variant Data

Information-dense single gene view

### Variant View

Gene Search:  Submit

Sort By: Gene

Alternative Transcripts: gene-exon (trans-annot)

Variants

Transcript

Protein

Variant Data

Information-dense single gene view

No need for pan and zoom

### Variant View

Gene Search:  Submit

Sort By: Gene

Alternative Transcripts: gene-exon (trans-annot)

Variants

Transcript

Protein

Variant Data

Sorting metrics guide gene navigation

### Variant View

Sorting metrics guide gene navigation

Control what shows up here

65

### Variant View

Peripheral supporting data

66

### Design information-dense visual encoding

- show all attributes necessary for variant analysis
  - match salience with importance for analysis task
- variant not just a thin line!
- emphasize with high salience
  - collocated variants fan out at top
  - grey variant vertical stroke intersects horizontal colored protein regions

67

### Design information-dense visual encoding

Reference AA

Variant

68

### Design information-dense visual encoding

Reference AA

Variant AA

Variant

69

### Design information-dense visual encoding

Reference AA

Variant AA

AA Chemical Class Colours:

- Charged
- Special
- Uncharged
- Hydrophobic

Variant

70

### Design information-dense visual encoding

Reference AA

Variant AA

AA Chemical Class Colours:

- Charged
- Special
- Uncharged
- Hydrophobic

Variant Type

- Stop
- Indel
- Deletion
- Insertion
- Splice
- Frameshift
- Nonsynonym

Variant

71

### Design information-dense visual encoding

Known Database

- Known Harmless
- Known Cancer

Reference AA

Variant AA

AA Chemical Class Colours:

- Charged
- Special
- Uncharged
- Hydrophobic

Variant Type

- Stop
- Indel
- Deletion
- Insertion
- Splice
- Frameshift
- Nonsynonym

Variant

72

### Design information-dense visual encoding

Known Database

- Known Harmless
- Known Cancer

Reference AA

Variant AA

AA Chemical Class Colours:

- Charged
- Special
- Uncharged
- Hydrophobic

Variant Type

- Stop
- Indel
- Deletion
- Insertion
- Splice
- Frameshift
- Nonsynonym

Transcript/Region Colours:

- Transcript
- AA Chain
- All Other Regions
- Non-Intersected Regions

Variant

73

### Results

74

### Highly scored gene by sorting metric: known leukemia gene

Variants

Transcript

Protein

75

### Visual inspection reveals collocation of variants

Variants

Transcript

Protein

76

### Several functional protein regions affected

Variants

Transcript

Protein

77

### Highly scored by metric: not previously known, good candidate

Variants

Transcript

Protein

78

### Protein chemical class change evident

Variants

Transcript

Protein

79

### In contrast, low scoring gene

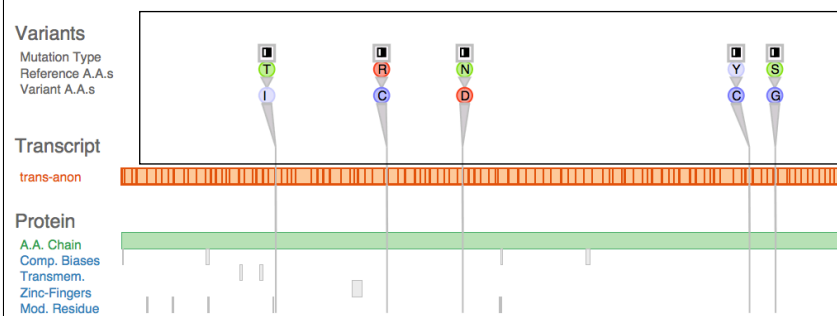
Variants

Transcript

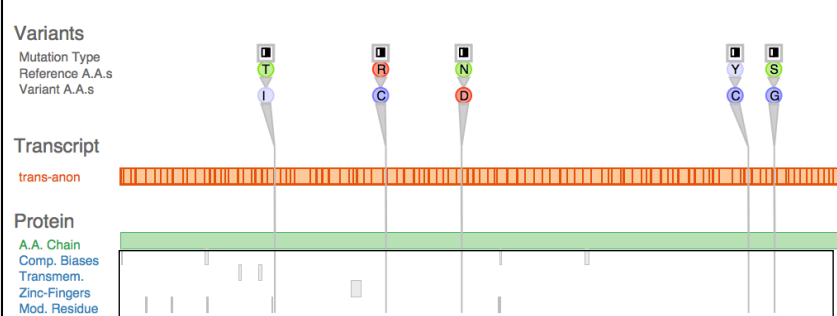
Protein

80

## No collocation of variants



## Mostly unaffected protein regions



## Methods

## Phase I: Winnow and Cast



- embedded within GSC for all stages
- winnow stage
  - considered and ruled out many potential collaborators
- cast stage
  - gatekeeper (PI)
  - two front-line analysts (postdocs)

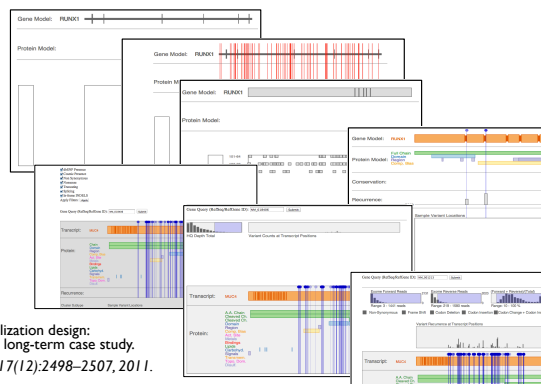


more at:  
 Design Study Methodology: Reflections from the Trenches and from the Stacks.  
 Sedlmair, Meyer, Munzner. IEEE TVCG 18(12): 2431-2440, 2012 (Proc. InfoVis 2012).

## Phase 2: Core Design



- main task abstraction
  - discover gene
- semi-structured interviews
  - every week for 1 hr
- iterative refinement
  - 8 data sketches deployed



Human-centered approaches in geovisualization design:  
 investigating multiple methods through a long-term case study.  
 Lloyd and Dykes. IEEE TVCG (Proc. InfoVis), 17(12):2498–2507, 2011.

## Phase 3: Two More Tasks



- two new analysts
  - connected by enthusiastic gatekeeper
- new task abstractions
  - compare patients
  - debug pipeline
- transferrable with minimal changes



## Phase 4: Reflect and write



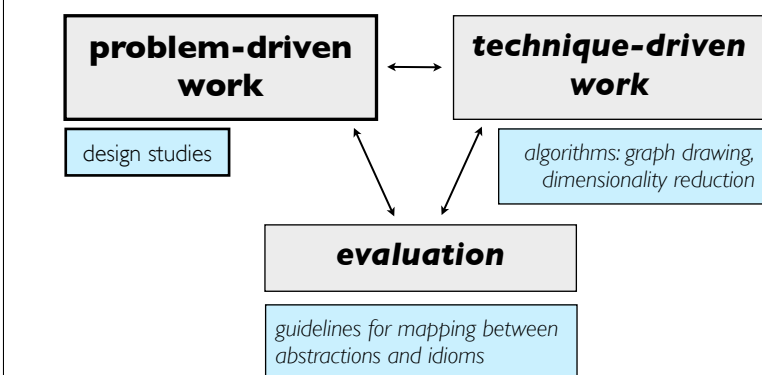
- abstraction innovation
  - data abstraction: highly filtered *transcript coordinates* (vs genome coordinates)
- guidelines
  - specialize first, generalize later
    - good for domains with complex data
  - high-level considerations
    - identifying scales of interest
    - what to visually encode directly vs what to support through interaction
    - when (and how) to eliminate navigation

## Themes, Revisited

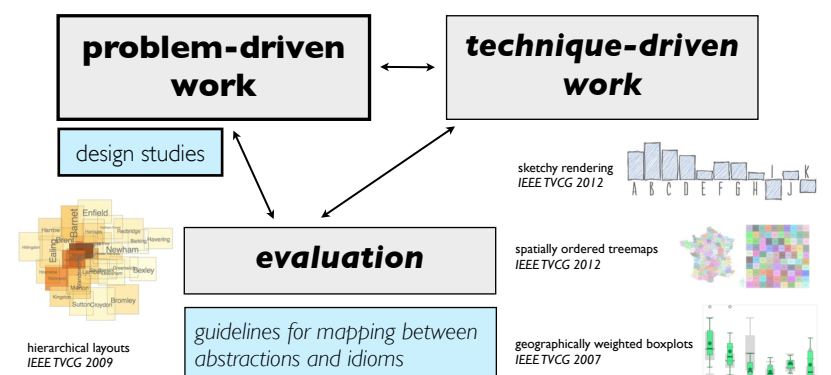
- what and why to show: task and data abstraction
  - task and data commonalities cross-cut domains
- how to show: visual encoding and interaction idioms
  - RelEx: reduce memory load with interaction
  - VariantView: reduce interaction load with better visual encoding
- transferability from design studies
  - DSM: reflection to confirm/refute/refine/propose guidelines



## Research Interests



## Research Interests: giCentre Context



## Further Information

- further info
  - <http://www.cs.ubc.ca/~tmm/talks.html#london14> (this talk, and many others)
  - <http://www.cs.ubc.ca/group/infovis> (papers, software, videos)
  - <http://www.cs.ubc.ca/~tmm/courses/infovis/book> (book: to appear)
    - Visualization Analysis and Design. Munzner. AK Peters 2014
- open source software downloads
  - <http://www.cs.ubc.ca/labs/imager/tr/2013/VariantView/VariantViewSoftware/>
- acknowledgements
  - funding: NSERC, NSF
  - joint work: all co-authors
    - Andreas Butz, Annika Frank, Joel Ferstay, Miriah Meyer, Cydney Nielsen, Michael Sedlmair
  - feedback on this talk
    - Matthew Brehmer, Stephen Ingram

