

# Variant View

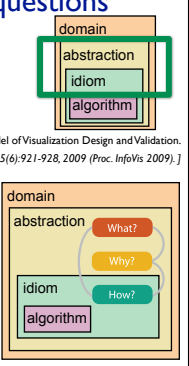
Genomics, Big Data, and Patient Privacy Implications

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

**THINK Conference**  
 6 Nov 2015, Santa Cruz CA  
<http://www.cs.ubc.ca/~tmm/talks.html#think15>  
 @tamaramunzner

## Visualization 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**
    - often don't just draw what you're given: transform to new form
  - **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



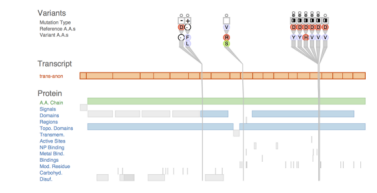
[A Nested Model of Visualization Design and Validation. Munzner. IEEE TVCG 15(6):921-928, 2009 (Proc. InfoVis 2009).]

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

# Variant View

Visualizing Sequence Variants in their Gene Context

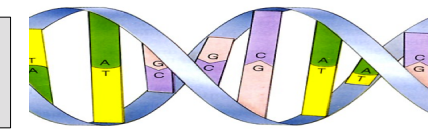
joint work with:  
 Joel Ferstay, Cydney Nielsen  
<http://www.cs.ubc.ca/labs/imager/tr/2012/VariantView/>



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 person's genome



Reference Genome DNA: ATA TGA TCA ACA CTT

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

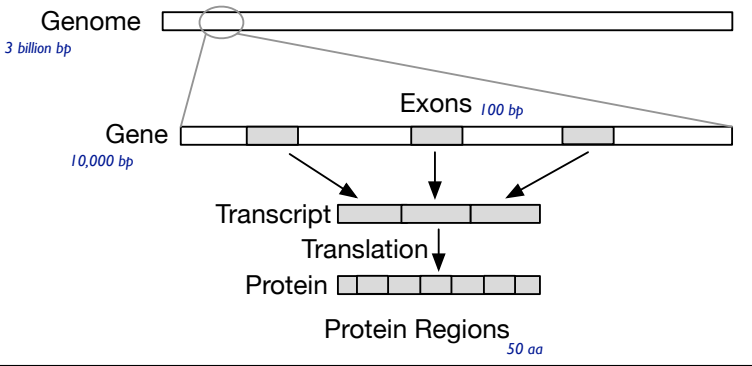
Sample 2 Genome DNA: ATA TGA TGA ACA CCT **Harmless?**

## Cancer Research

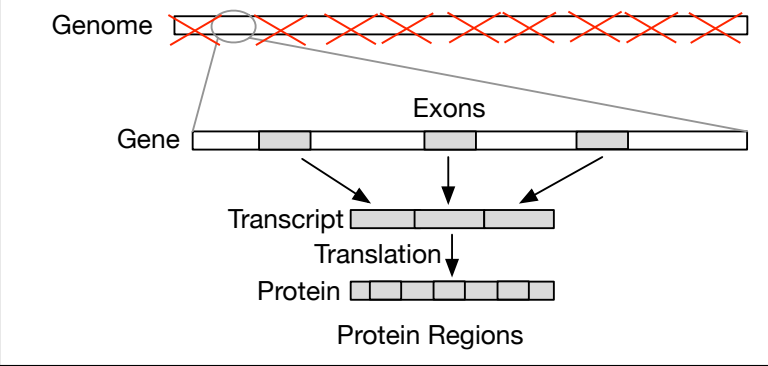
- collaboration with analysts at BC Genome Sciences Center
  - studying genetic basis of leukemia
- driving task
  - discover new candidate genes with harmful variants
- two big questions
  - what to show
    - data abstraction
    - challenge: enormous range of scales in the data
  - how to show it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

## Abstractions

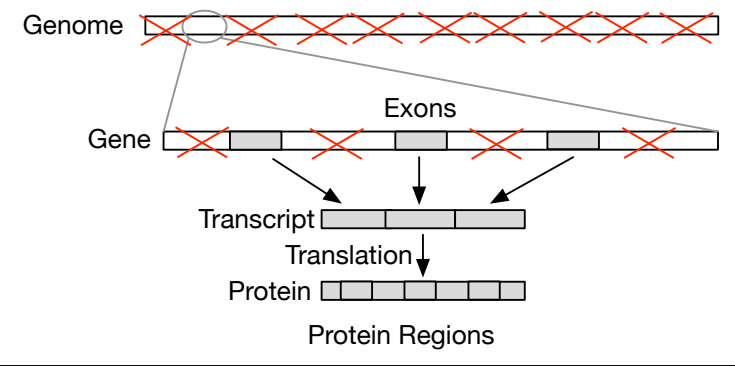
## Data: Filtering to relevant biological levels and scales



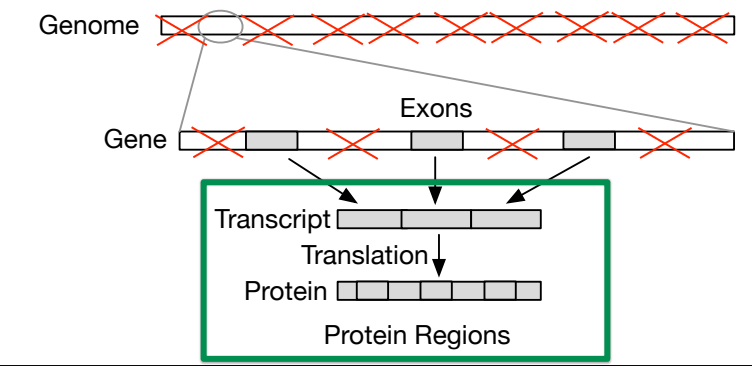
## Filter out whole genome; keep genes



## Filter out non-exon regions

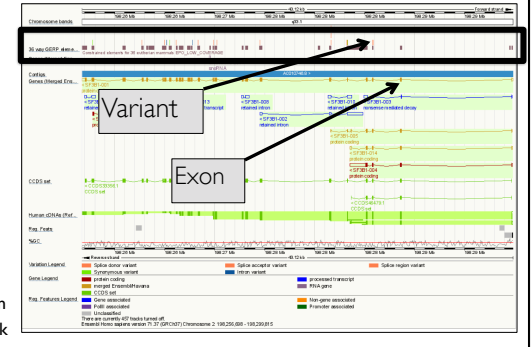


## Data abstraction: highly filtered scope



## 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

1st Screen

- Exon regions small
- Color coding difficult to see
- Protein regions overlap on same track

Ensembl Variant Image  
 Chen et al, BMC Bioinformatics 2010.

## Idioms

## Variant View

Gene Search:  Submit

Sort By: Gene, Alpha, Cluster Score

Alternative Transcripts: gene-annot (trans-annot)

Variants: Mutation Type, Reference A.A., Variant A.A.

Transcript: gene-annot (trans-annot)

Protein: A.A. Chain, Domains, Regions, Active Sites, Binding Sites, Most Residue

Variant Data: Patient ID, Chr, Coord, Ref Base, Var Base, dbSNP128, dbSNP136, dbSNP137, COSMIC, A.A. Chng, Gene, Ref

## Variant View

Previous: Table, one row per variant

Information-dense single gene view

Gene Search:  Submit

Sort By: Gene, Alpha, Cluster Score

Alternative Transcripts: gene-annot (trans-annot)

Variants: Mutation Type, Reference A.A., Variant A.A.

Transcript: gene-annot (trans-annot)

Protein: A.A. Chain, Domains, Regions, Active Sites, Binding Sites, Most Residue

Variant Data: Patient ID, Chr, Coord, Ref Base, Var Base, dbSNP128, dbSNP136, dbSNP137, COSMIC, A.A. Chng, Gene, Ref

## Variant View

Information-dense single gene view

Gene Search:  Submit

Sort By: Gene, Alpha, Cluster Score

Alternative Transcripts: gene-annot (trans-annot)

Variants: Mutation Type, Reference A.A., Variant A.A.

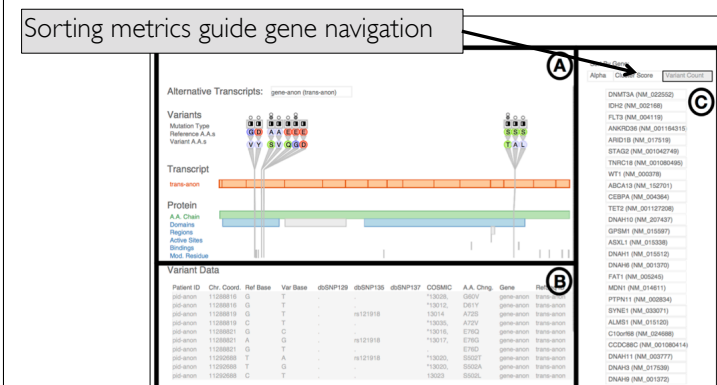
Transcript: gene-annot (trans-annot)

Protein: A.A. Chain, Domains, Regions, Active Sites, Binding Sites, Most Residue

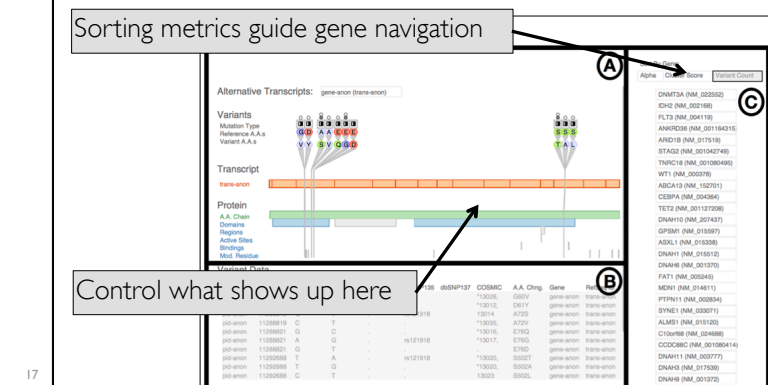
Variant Data: Patient ID, Chr, Coord, Ref Base, Var Base, dbSNP128, dbSNP136, dbSNP137, COSMIC, A.A. Chng, Gene, Ref

No need for pan and zoom

# Variant View



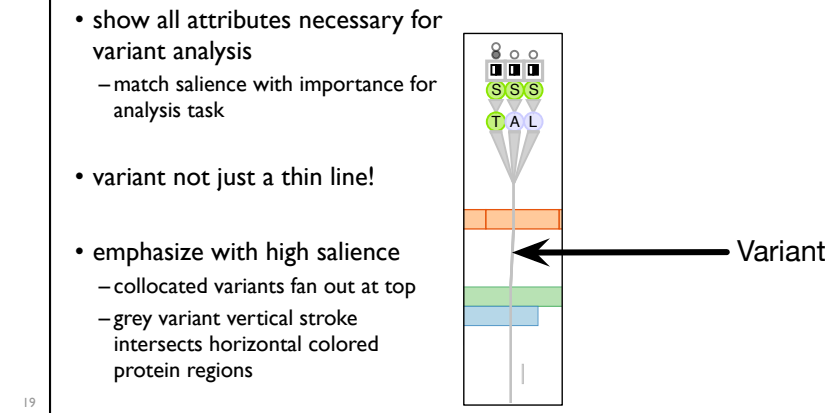
# Variant View



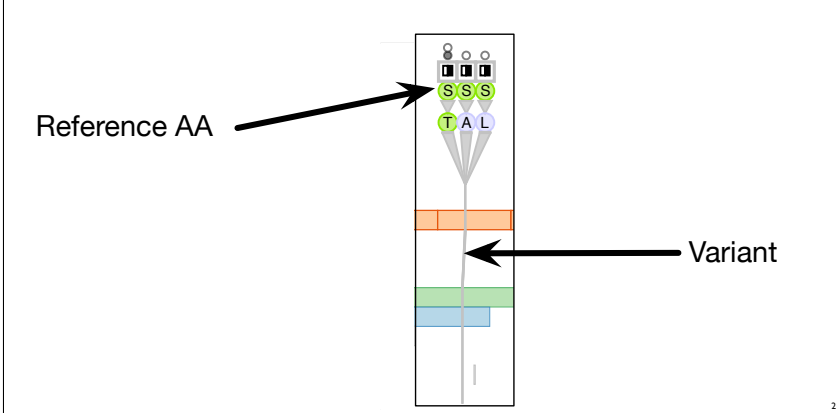
# Variant View



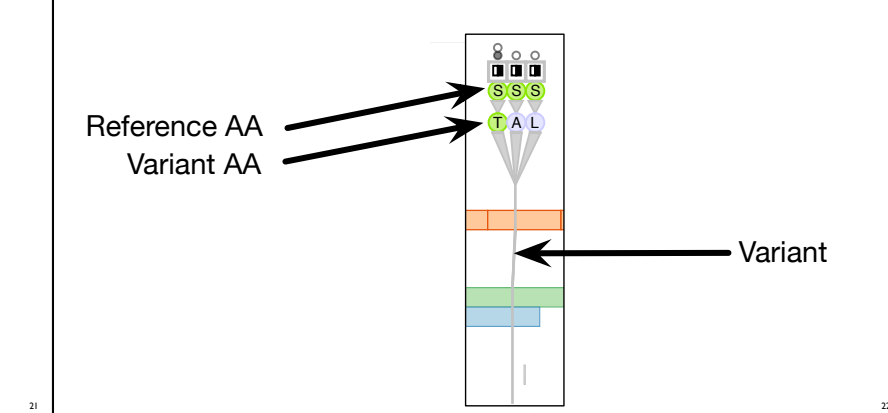
# Design information-dense visual encoding



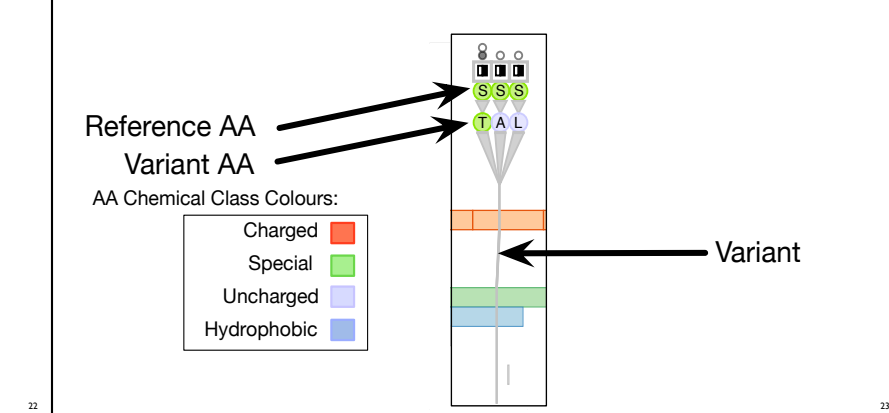
# Design information-dense visual encoding



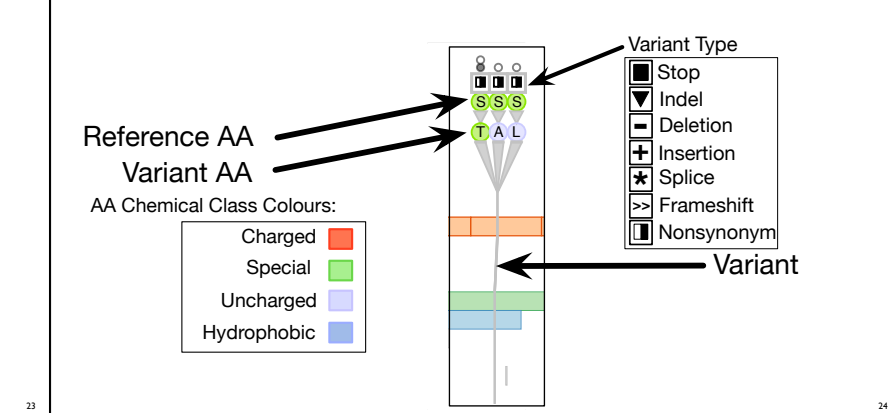
# Design information-dense visual encoding



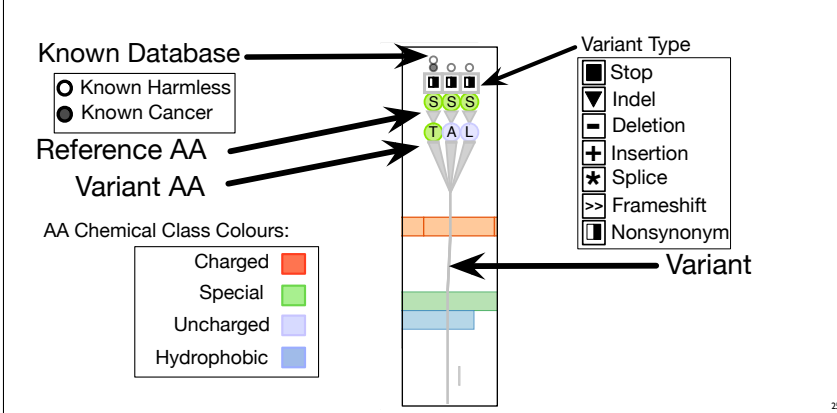
# Design information-dense visual encoding



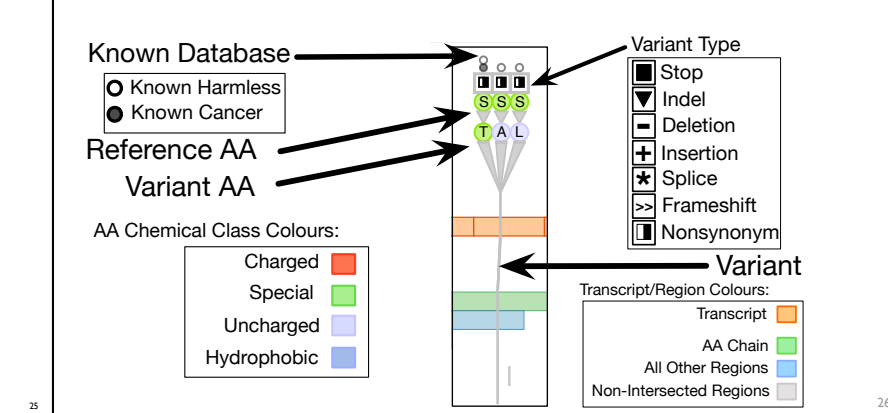
# Design information-dense visual encoding



# Design information-dense visual encoding



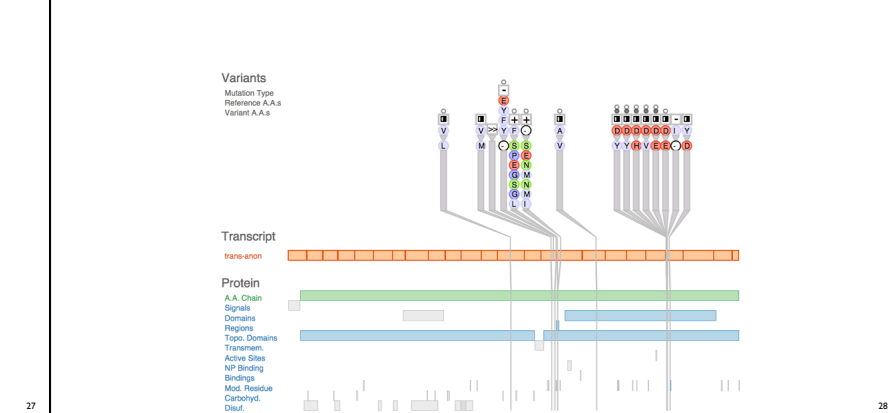
# Design information-dense visual encoding



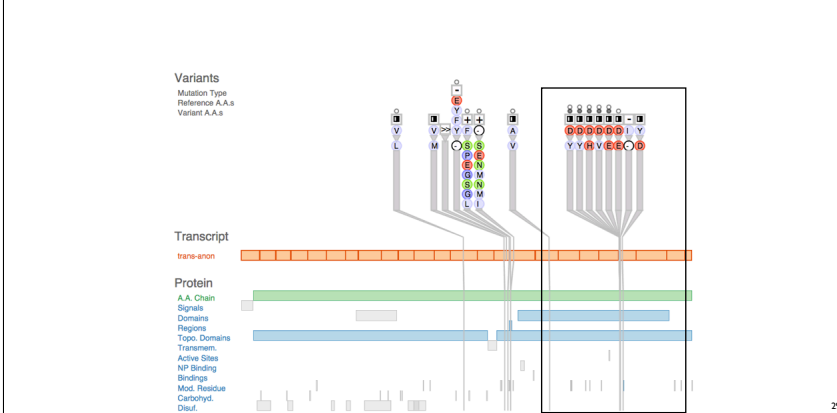
# Results



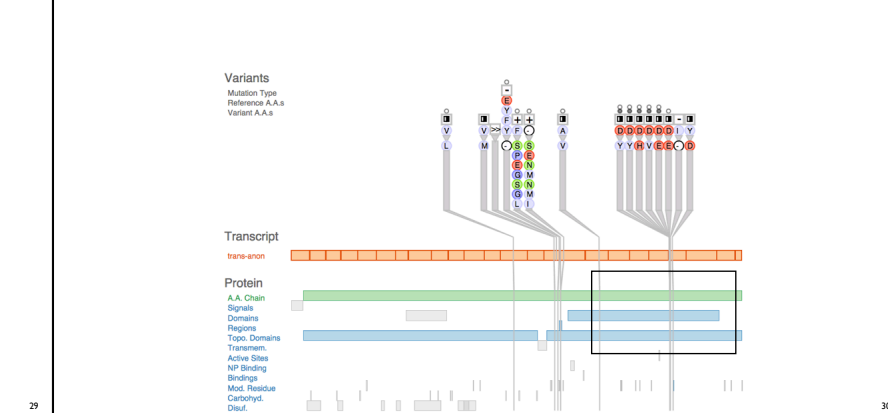
# Verify known leukemia gene: Highly scored by sorting metric



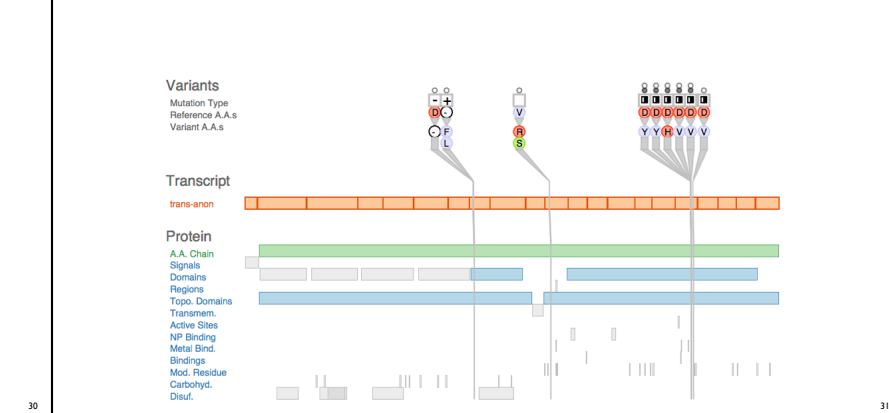
# Visual inspection reveals collocation of variants



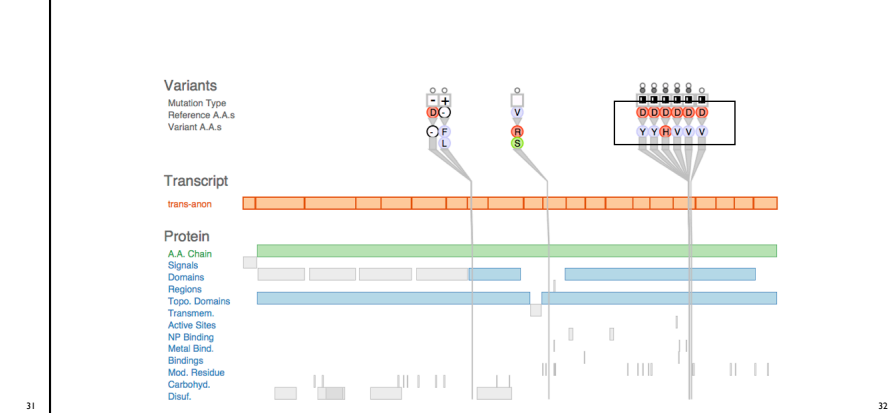
# Several functional protein regions affected



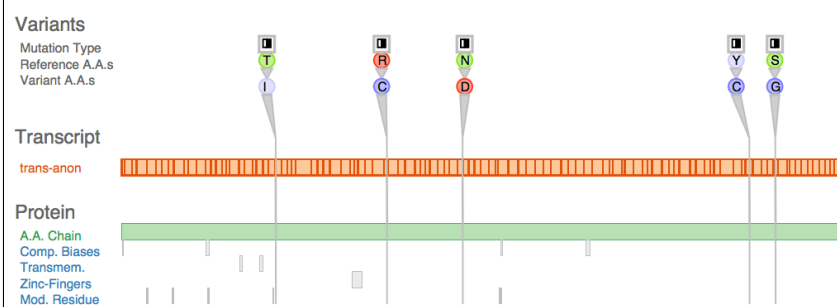
# Highly scored by metric: not previously known, good candidate



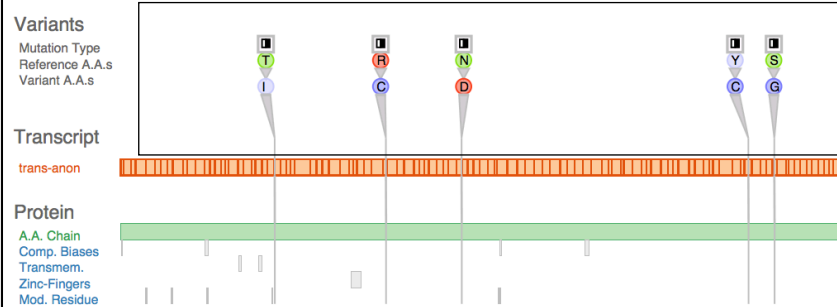
# Protein chemical class change evident



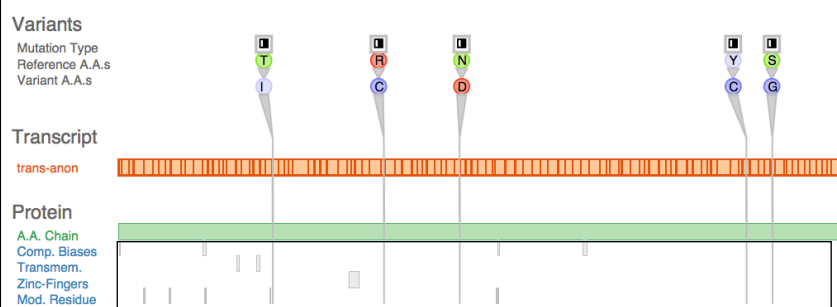
## In contrast, low scoring gene



## No collocation of variants



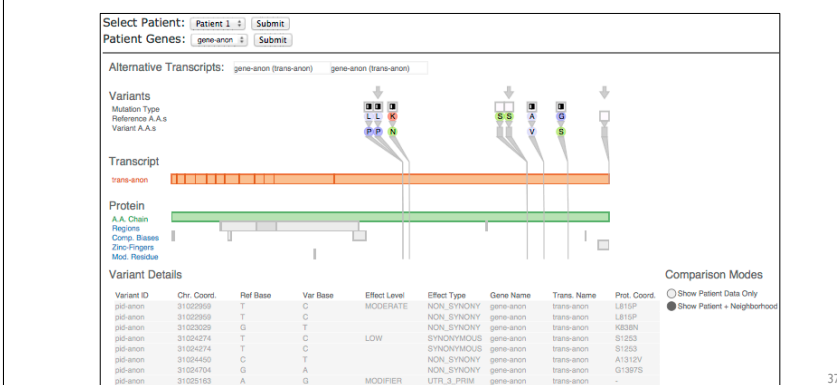
## Mostly unaffected protein regions



## Additional tasks

- task 2: compare patients
  - clinical setting application
  - compare patient data to known harmful variants
  - challenge
    - similarity is loosely understood rather than fully characterized
    - visual inspection for what constitutes a match

## Adapted Variant View with minimal changes



## Navigate through patient data with list



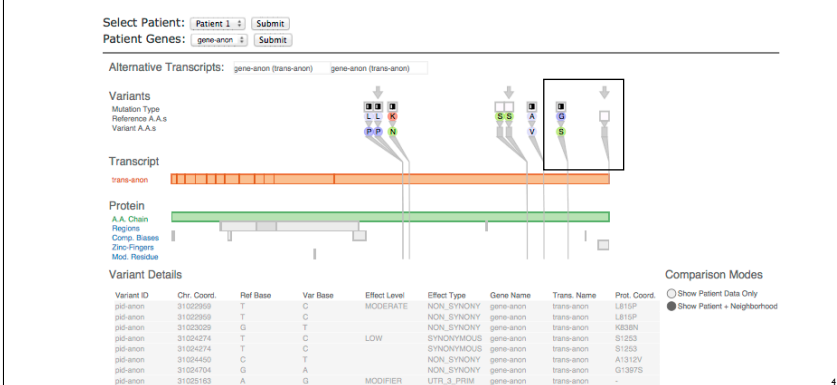
## Patient data emphasized with arrows



## Patient has same harmful L to P mutation



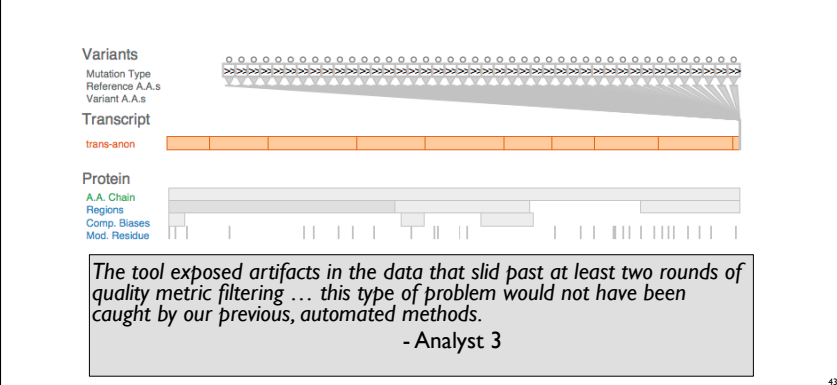
## Nonmatching variants



## Additional tasks

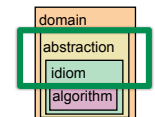
- task 3: debug pipeline
  - data cleansing before analysis
  - analysts originally thought pipeline fully debugged
    - no perceived need for vis support

## Tool revealed errors in the data



## Reduce from big data to manageable data

- at abstraction level, not algorithm level
  - filter away huge amounts of detail
  - transform from variant centric to gene centric



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

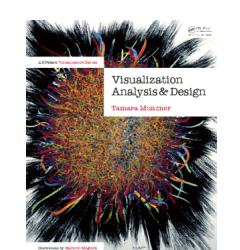
## Privacy implications

- anonymizing
  - protect patients: strip/change patient names
  - protect researchers: strip/change gene names
  - deanonymization threats: minor in this case
- future projects
  - genomics meets clinical meets administrative meets demographic
    - enable data re-use: link individuals across datasets to create research study cohorts
    - health research opportunity: major
    - deanonymization threat: major
  - ferociously privacy-protected data silos/vaults
    - can't see data without permission; what to ask for?
    - upcoming project: privacy preserving visual metadata browsing



## More Information

- paper & open source download  
<http://www.cs.ubc.ca/labs/imager/tr/2013/VariantView/>
- this talk  
<http://www.cs.ubc.ca/~tmm/talks.html#think15>
- papers, videos, software, talks, courses  
<http://www.cs.ubc.ca/group/infovis>  
<http://www.cs.ubc.ca/~tmm>
- book: Visualization Analysis and Design. CRC Press 2014.  
<http://www.cs.ubc.ca/~tmm/vadbook>
  - 20% promo code for book+ebook combo: HVN17
  - <http://www.crcpress.com/product/isbn/9781466508910>



- acknowledgements
  - funding: VIVA, AeroInfo/Boeing, MITACS

@tamaramunzner