

Validation & Evaluation

CS 7250

SPRING 2020

Prof. Cody Dunne

NORTHEASTERN UNIVERSITY

Slides and inspiration from Michelle Borkin, Krzysztof Gajos, Hanspeter Pfister, Miriah Meyer, Jonathan Schwabish, and David Sprague

BURNING QUESTIONS?

Home

Syllabus

Pages

Announcements

Assignments

Quizzes

Discussions

Grades

People

Files

Search for Quiz

▼ Assignment Quizzes



Quiz – Validation & Evaluation

Not available until Mar 26 | Due Mar 26 at 12pm | 3 pts | 3 Questions

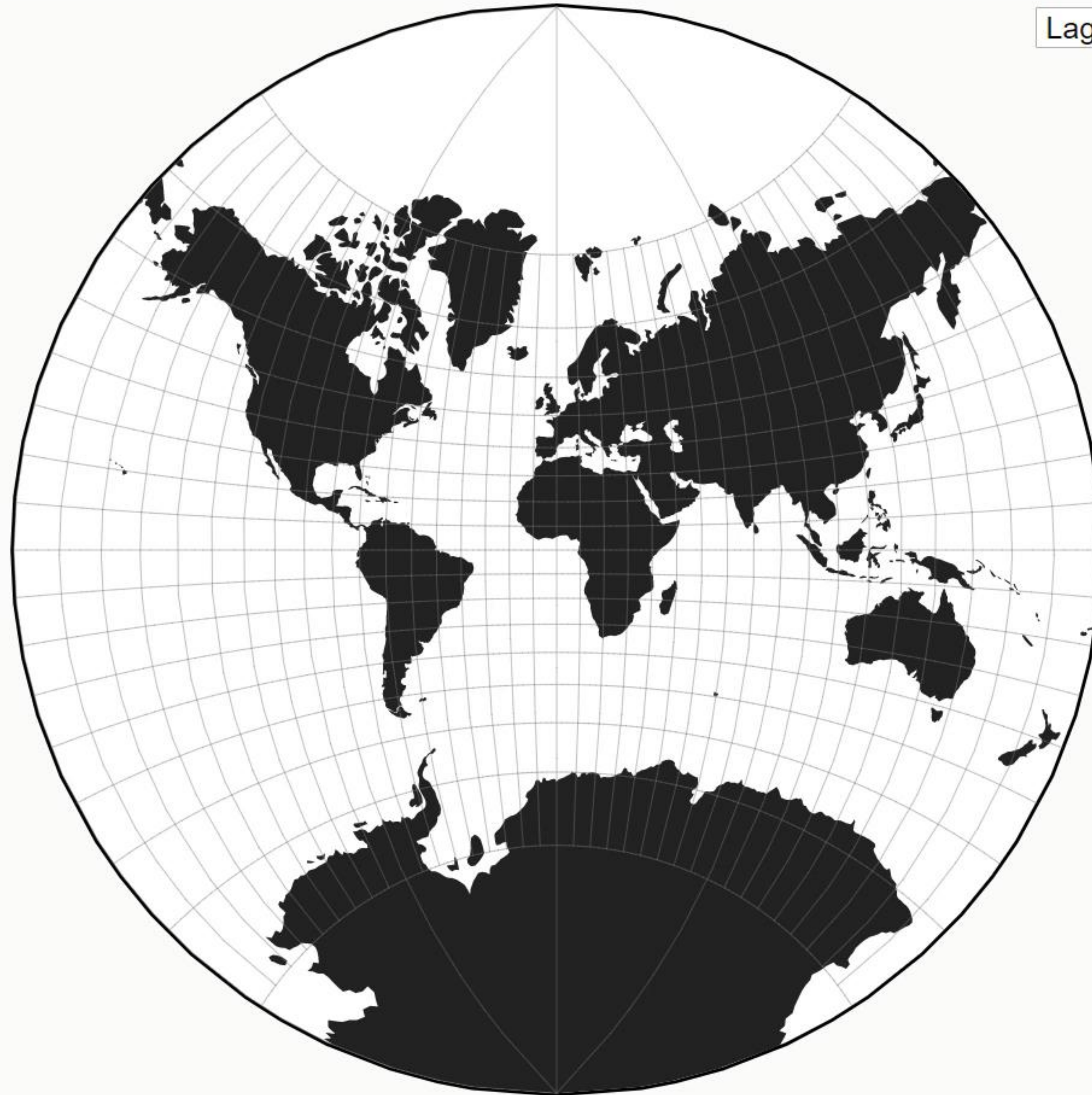
READING QUIZ

Quiz – Validation & Evaluation

~6 min

PREVIOUSLY, ON CS 7250...

Projection Transitions



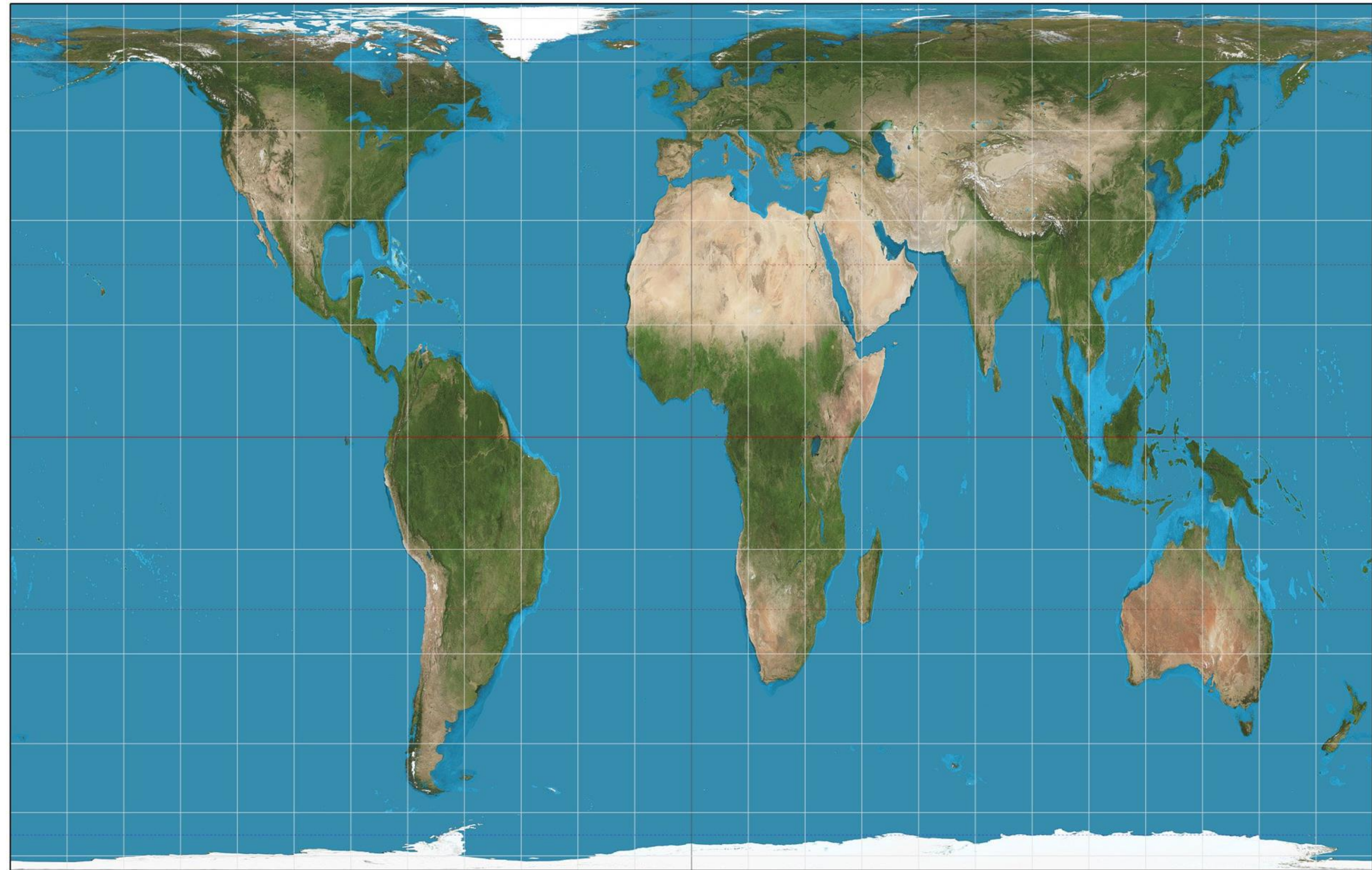
Lagrange ▼

Mercator Projection



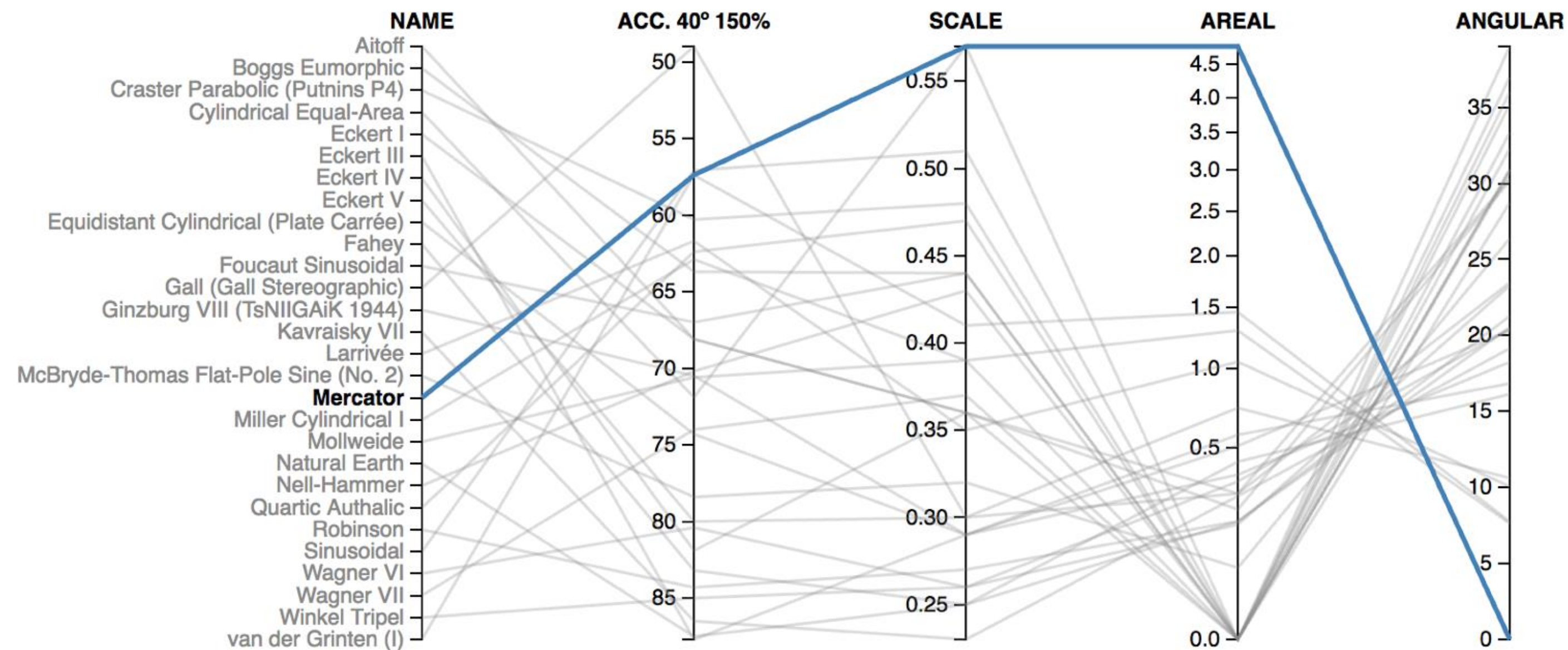
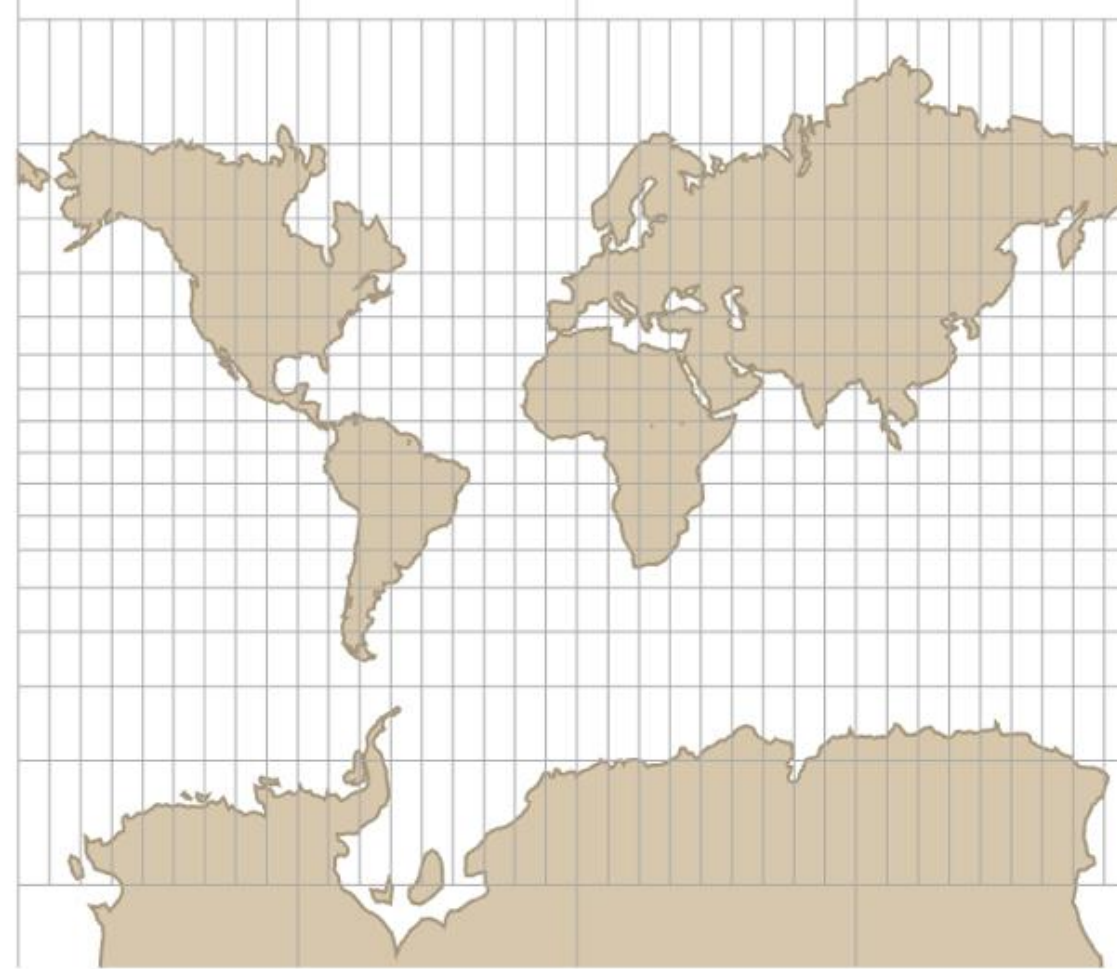
Great for ocean navigation,
but dramatically exaggerates poles.

Gall-Peters Projection



More accurate land areas.
(Officially endorsed by the UN.)

Comparing Map Projections



Overlays

Place Locations

Transportation

MBTA Subway Lines

BLUE

GREEN

ORANGE

RED

SILVER

Bike Trails

Evacuation Routes

Investment & Growth (Building Permits)

Assessed Value (Tax Assessments)

Trends in Assessed Value

Annual Changes in Assessed Value

Building Age

1941 - 1953

1953 - 1959

1959 - 1967

1967 - 1976

1976 - 1987

1987 - 2007

Medical Emergencies (911 reports) (2015)

Social Disorder and Crime (911 Report) (2015)

Physical Disorder (311 reports) (2015)

Usage of 311 System (2015)

Gentrification (2000 - 2014)

American Community Survey (2011-2015)

Basic Characteristics, ACS 2011-2015

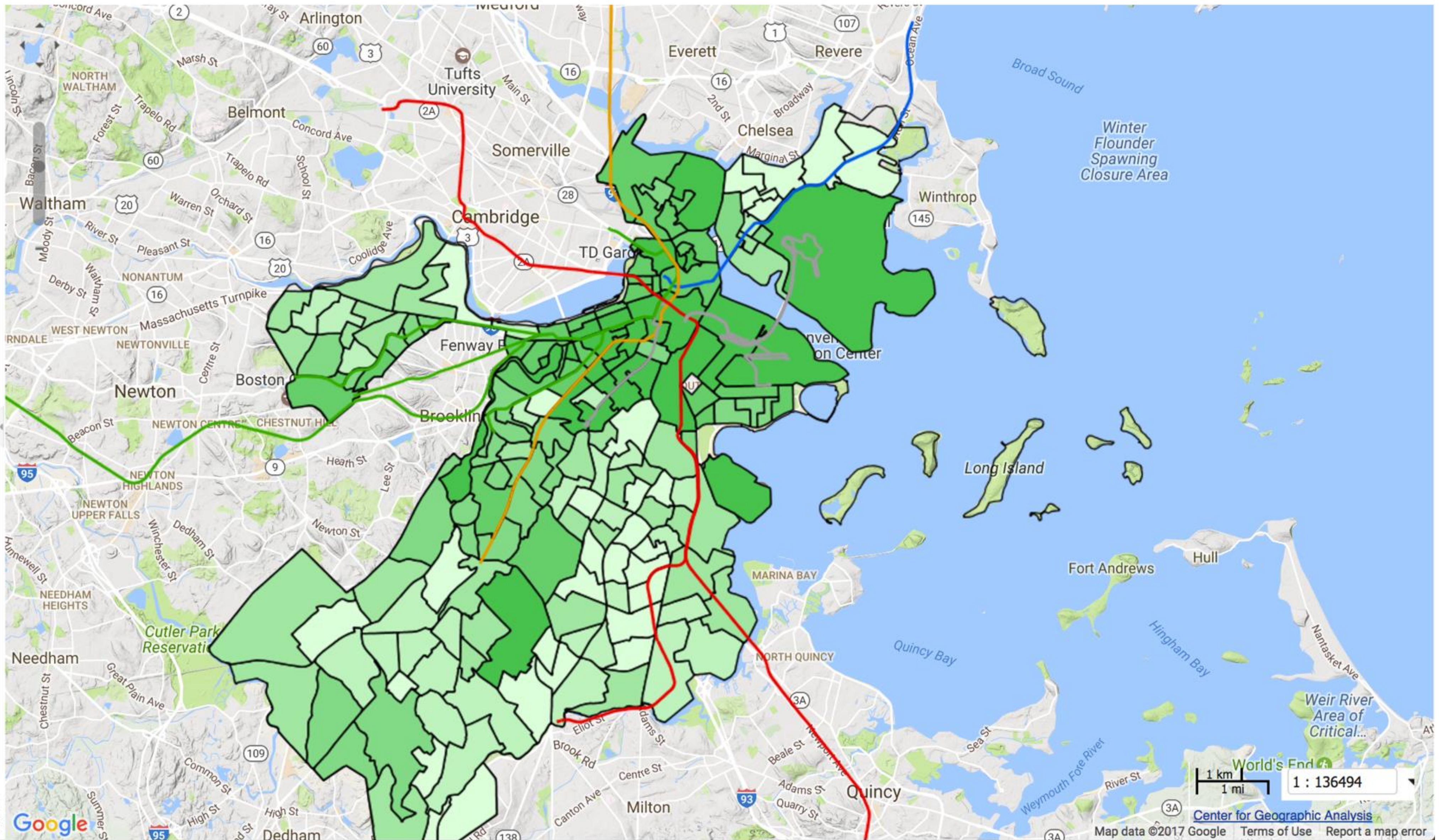
Racial and Ethnic Composition, ACS 2011-2015

Economic Characteristics, ACS 2011-2015

Education Levels, ACS 2011-2015

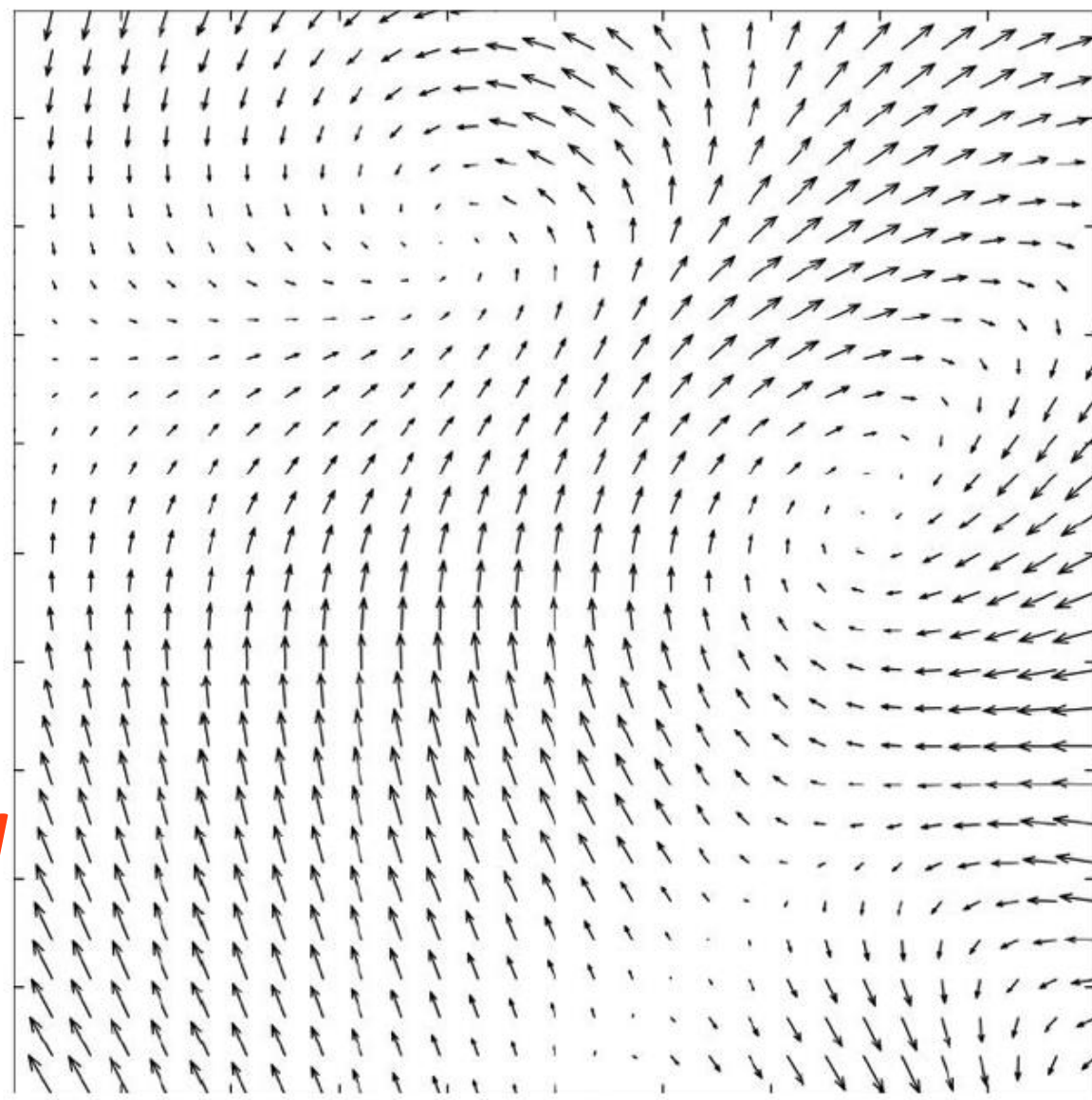
Family and Household Characteristics, ACS 2011-2015

Transportation to Work, ACS 2011-2015

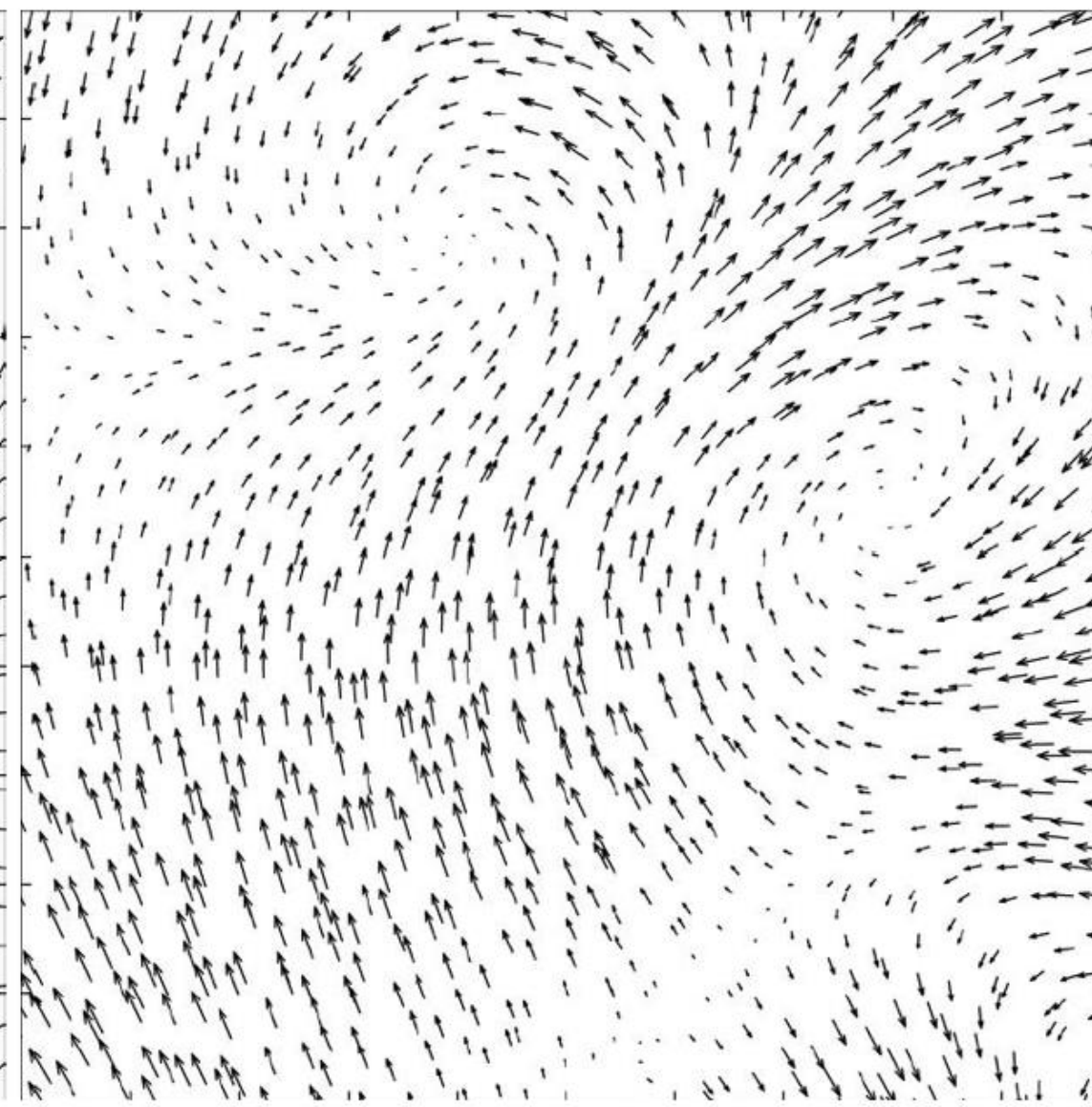


Vector Field Encoding Examples:

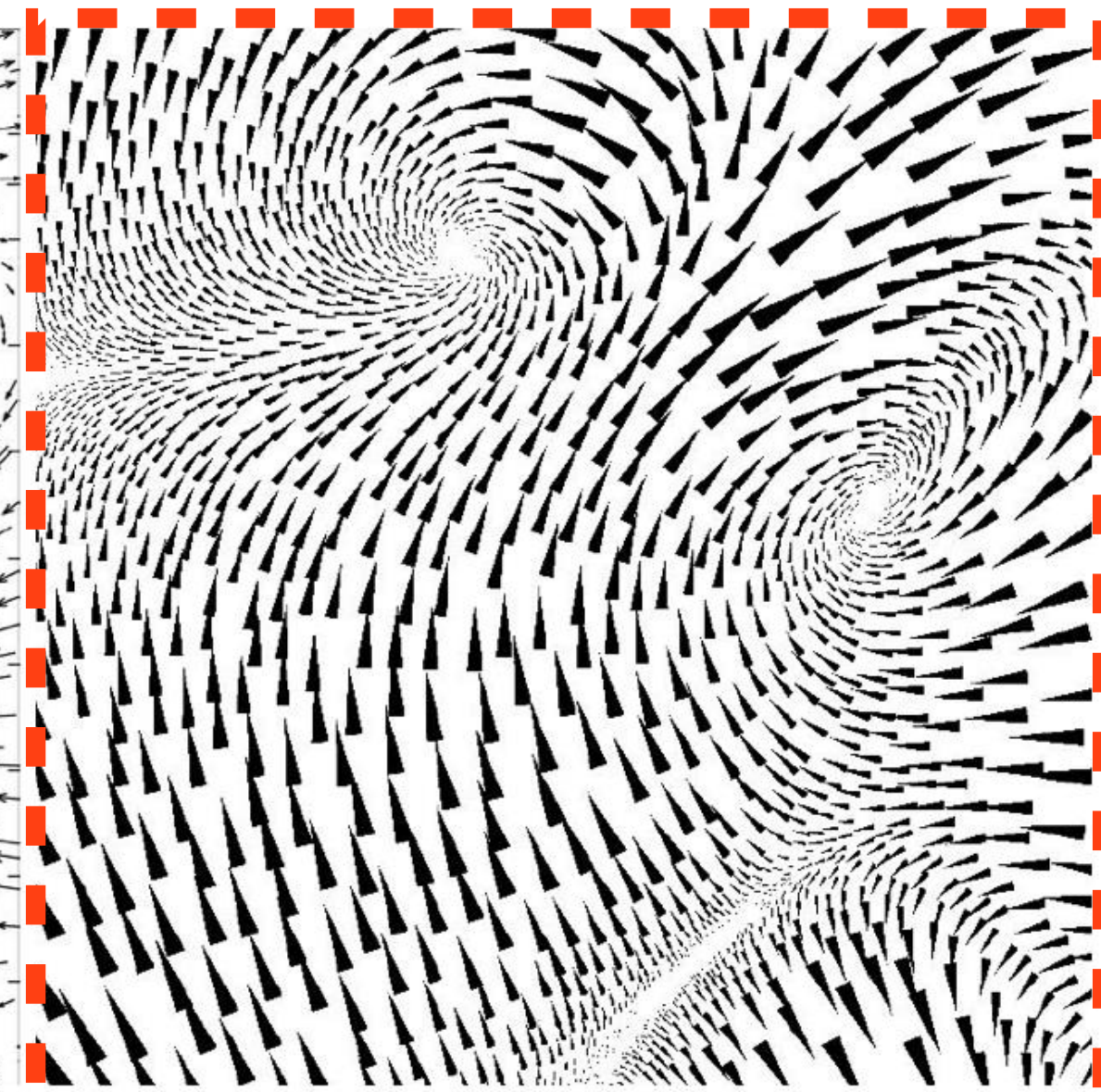
Most accurate and efficient for certain spatial tasks



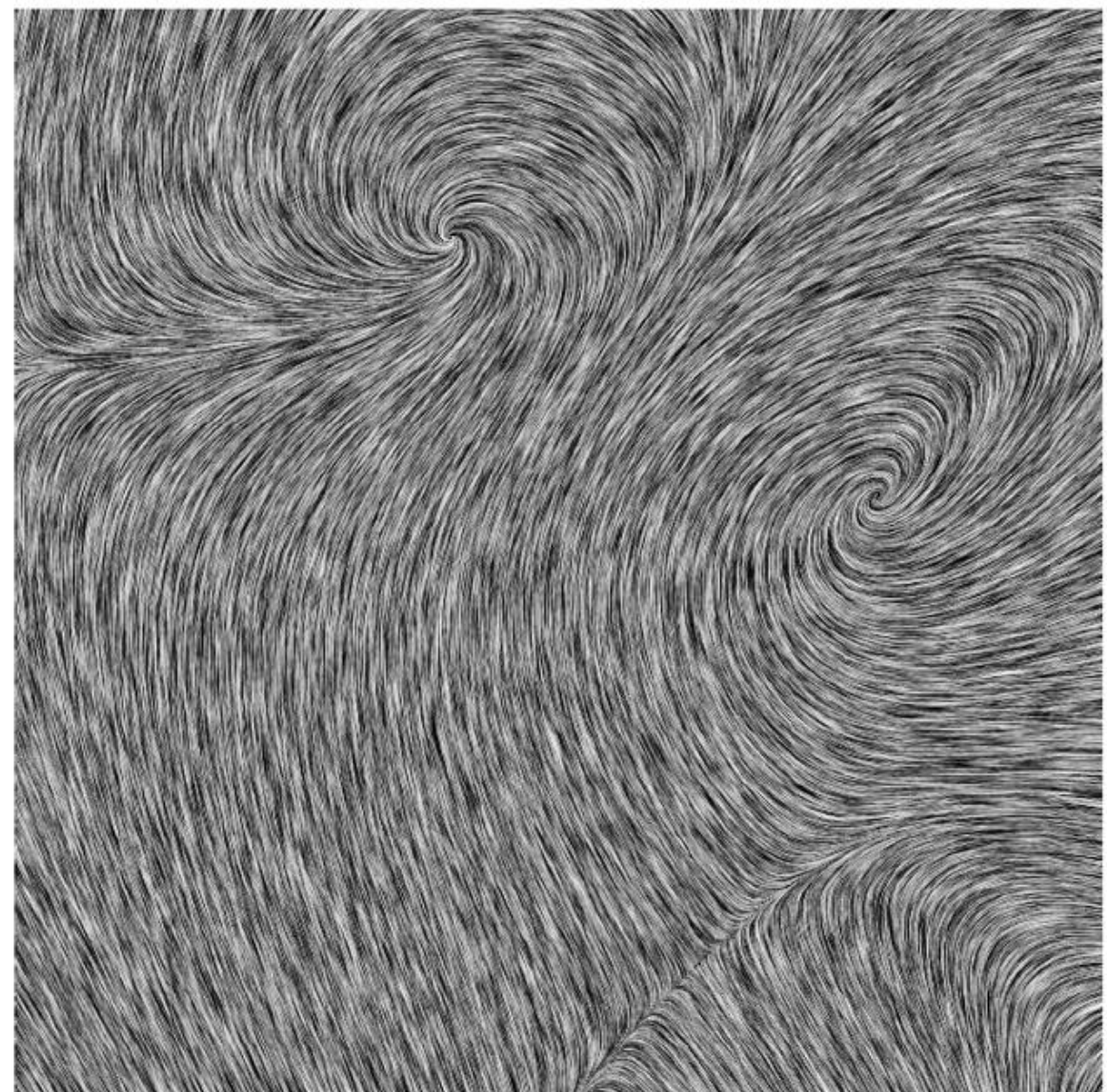
GRID



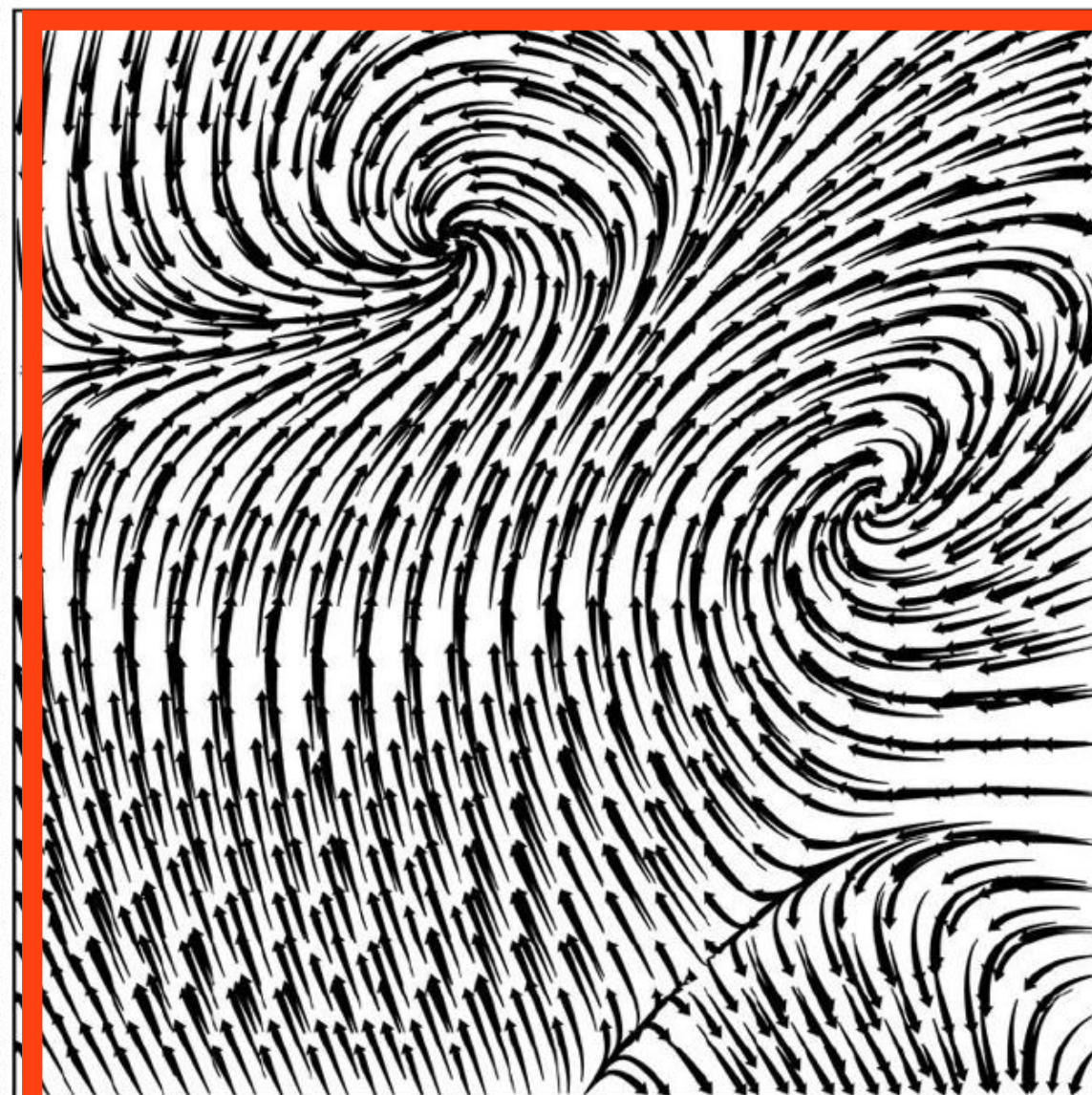
JIT



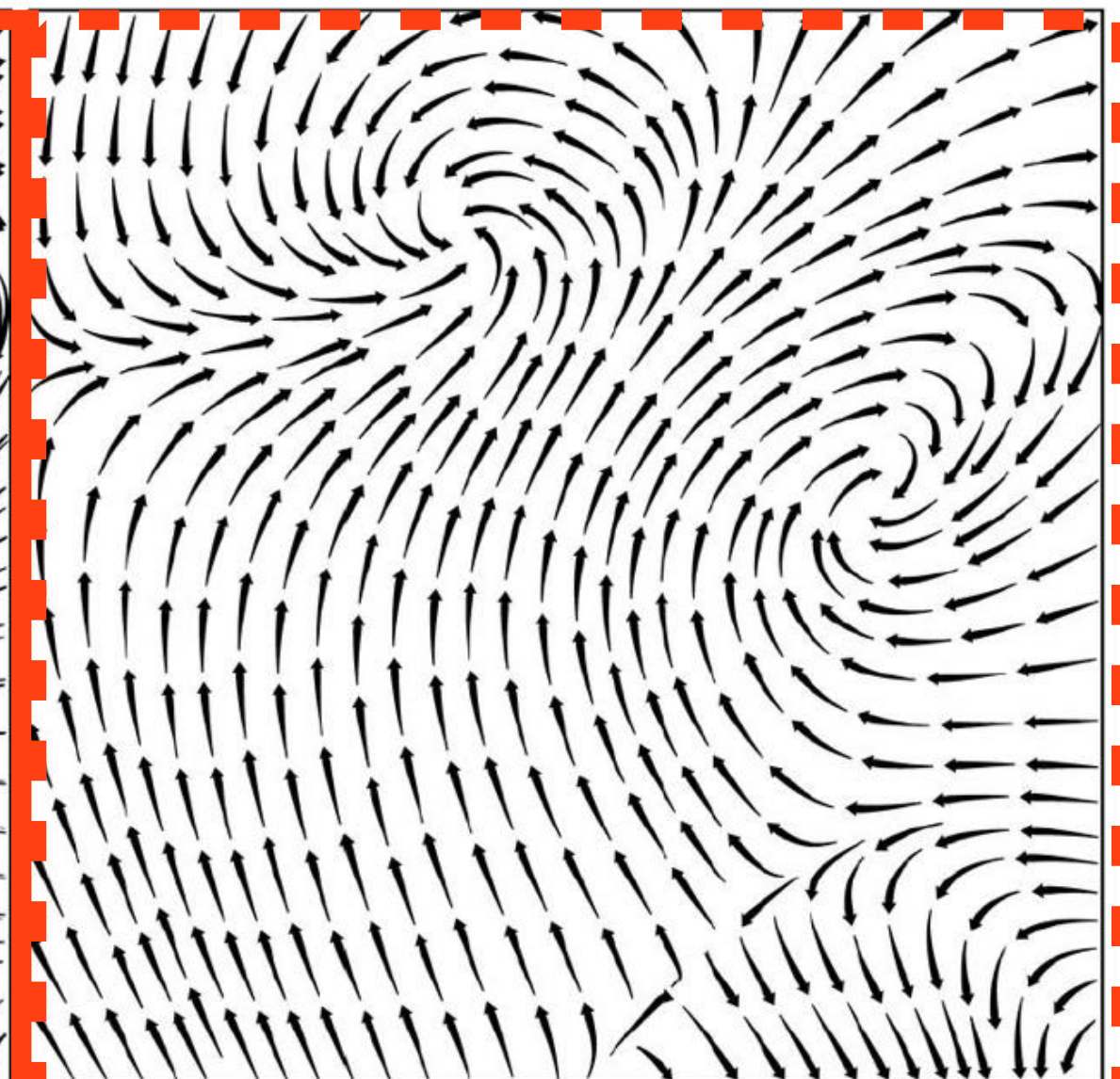
LIT



LIC



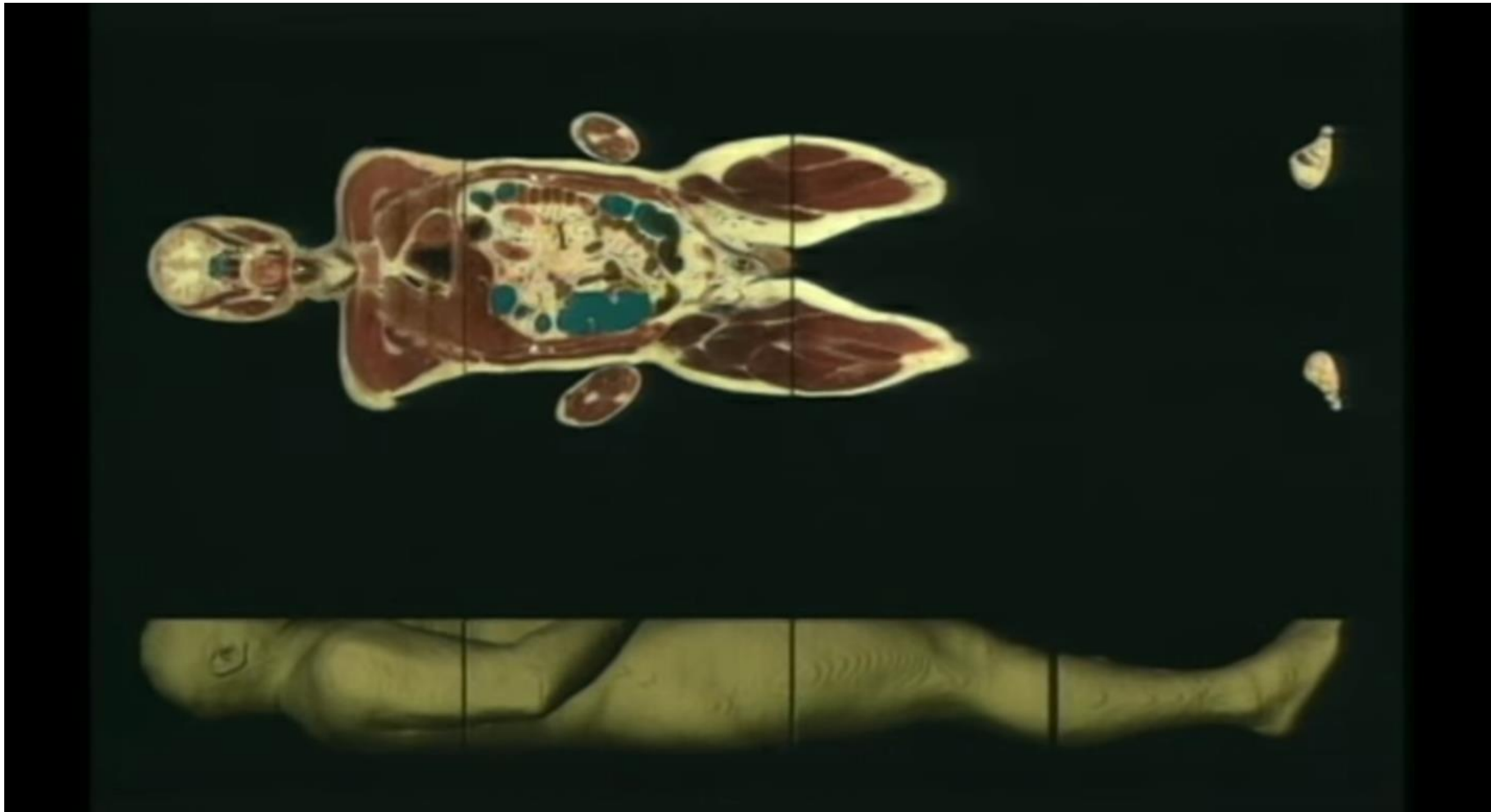
GSTR



OSTR

Isosurfaces & Volume Rendering

[Visible Human Project](#)

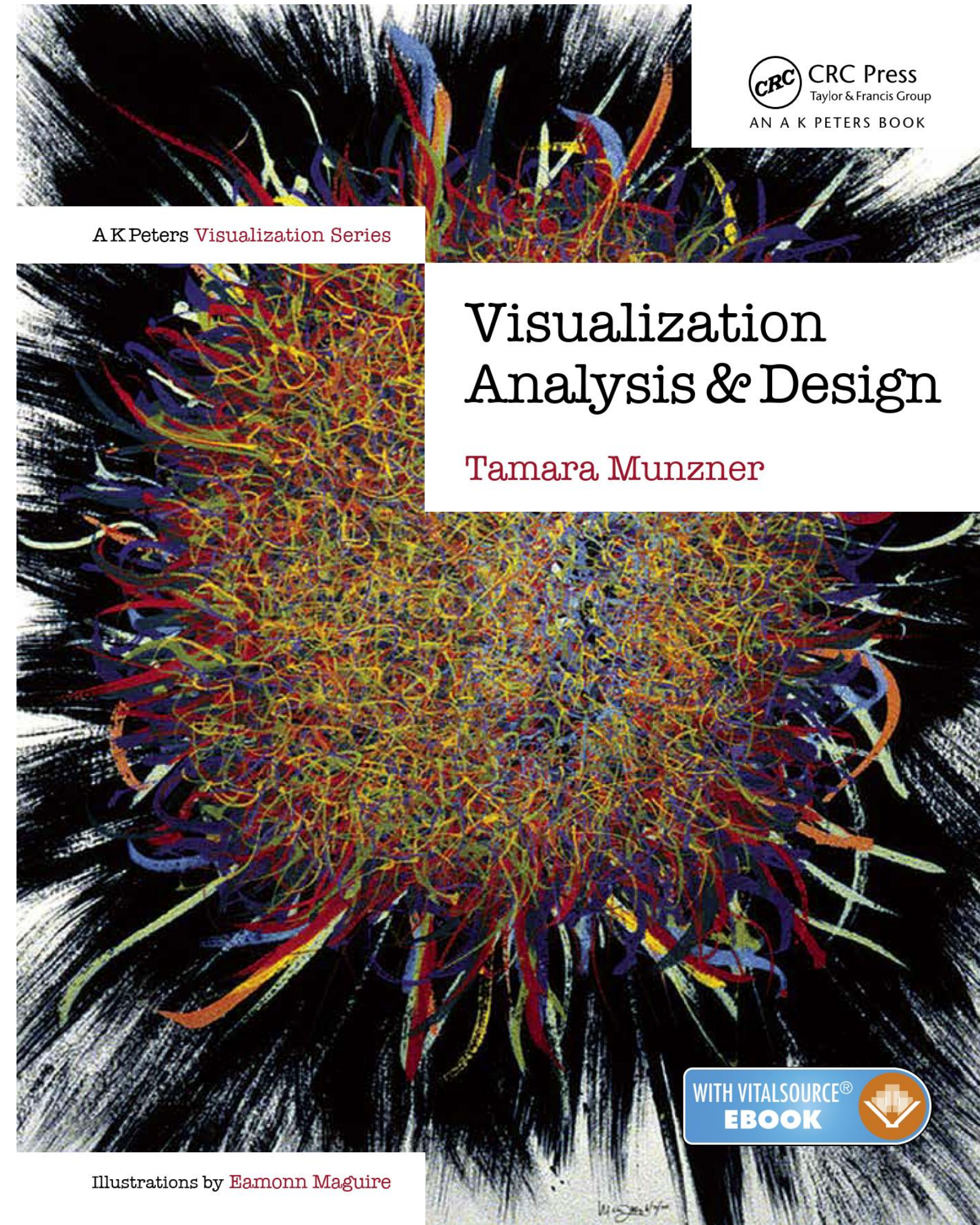


<https://www.youtube.com/watch?v=7GPB1sjEhIQ>

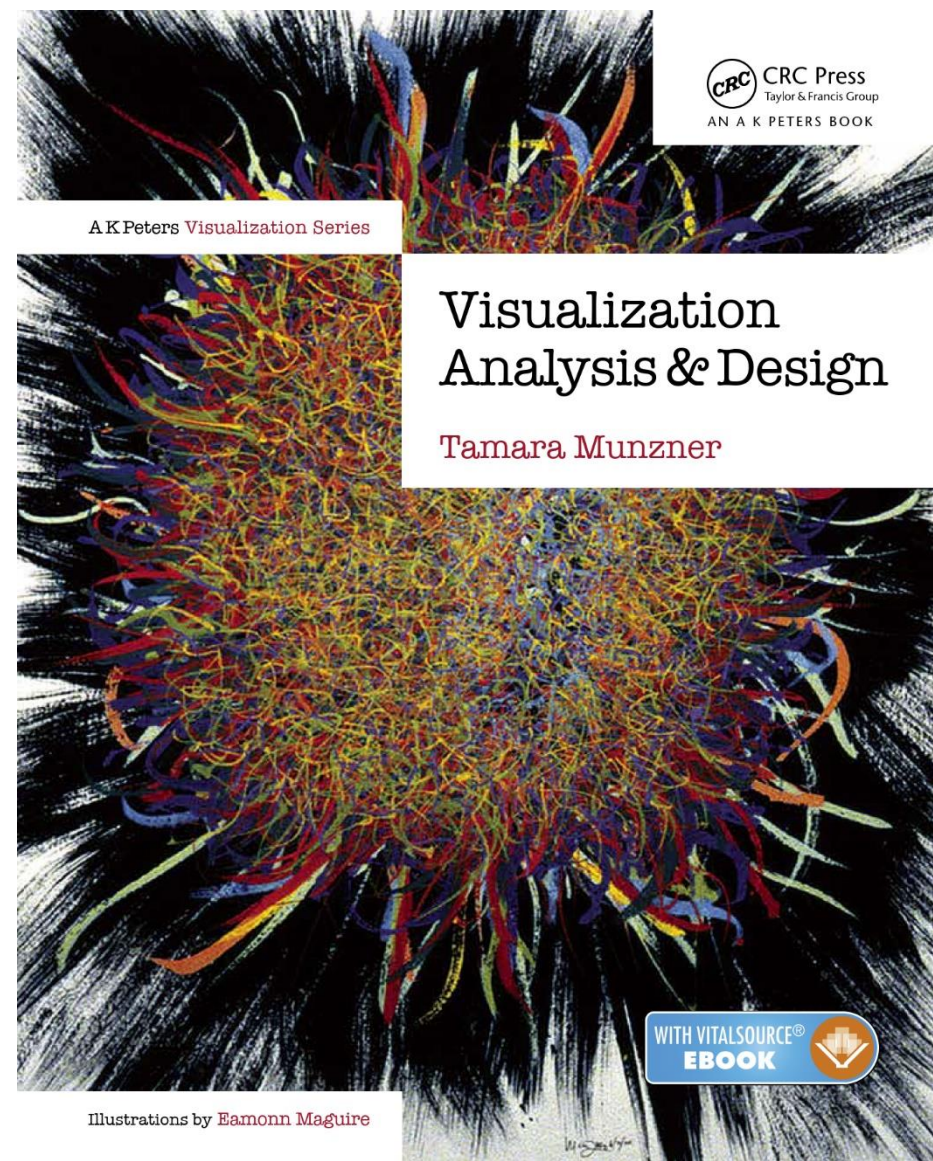
Now, ON CS 7250...

THE NESTED MODEL FOR VISUALIZATION VALIDATION


TEXTBOOK



Additional “recommended” books as resources in syllabus



“Nested Model”


 **Domain situation**
Observe target users using existing tools

Example

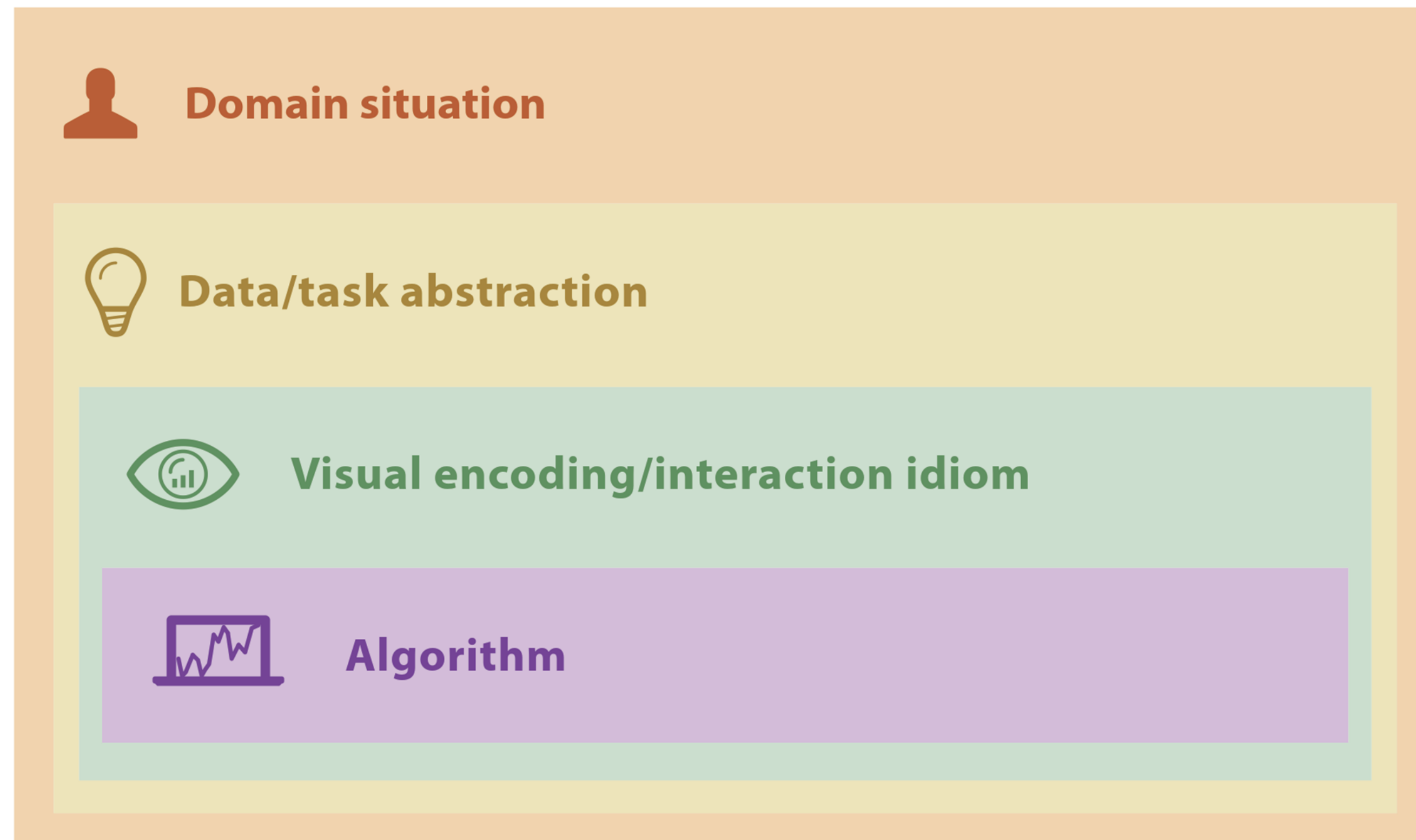
FAA (aviation)

What is the busiest time of day at Logan Airport?

Map vs. Scatter Plot vs. Bar

 Tamara
Munzner

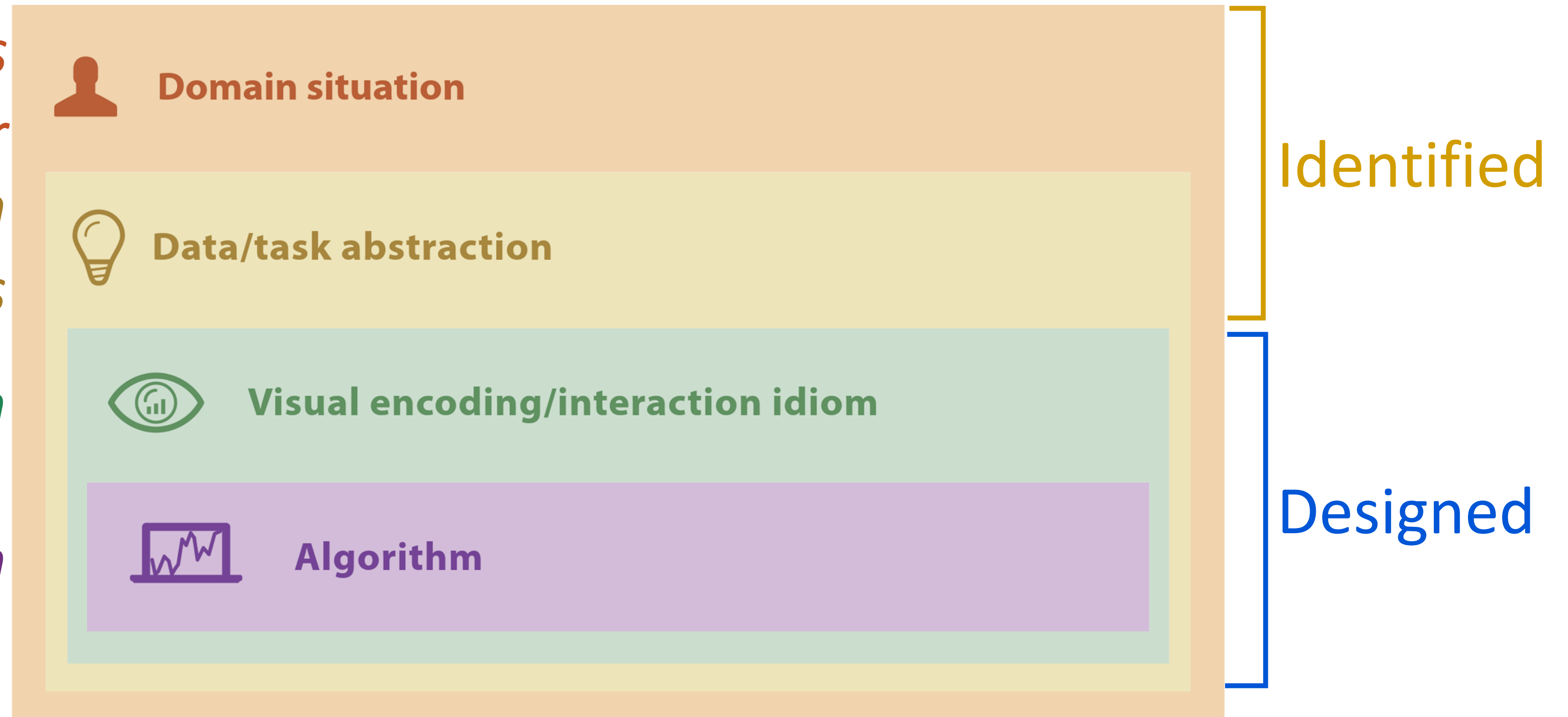
Nested Model



Nested Model

Human-centered design

Designer understands user
Abstract domain tasks
Visualization design
Implementation



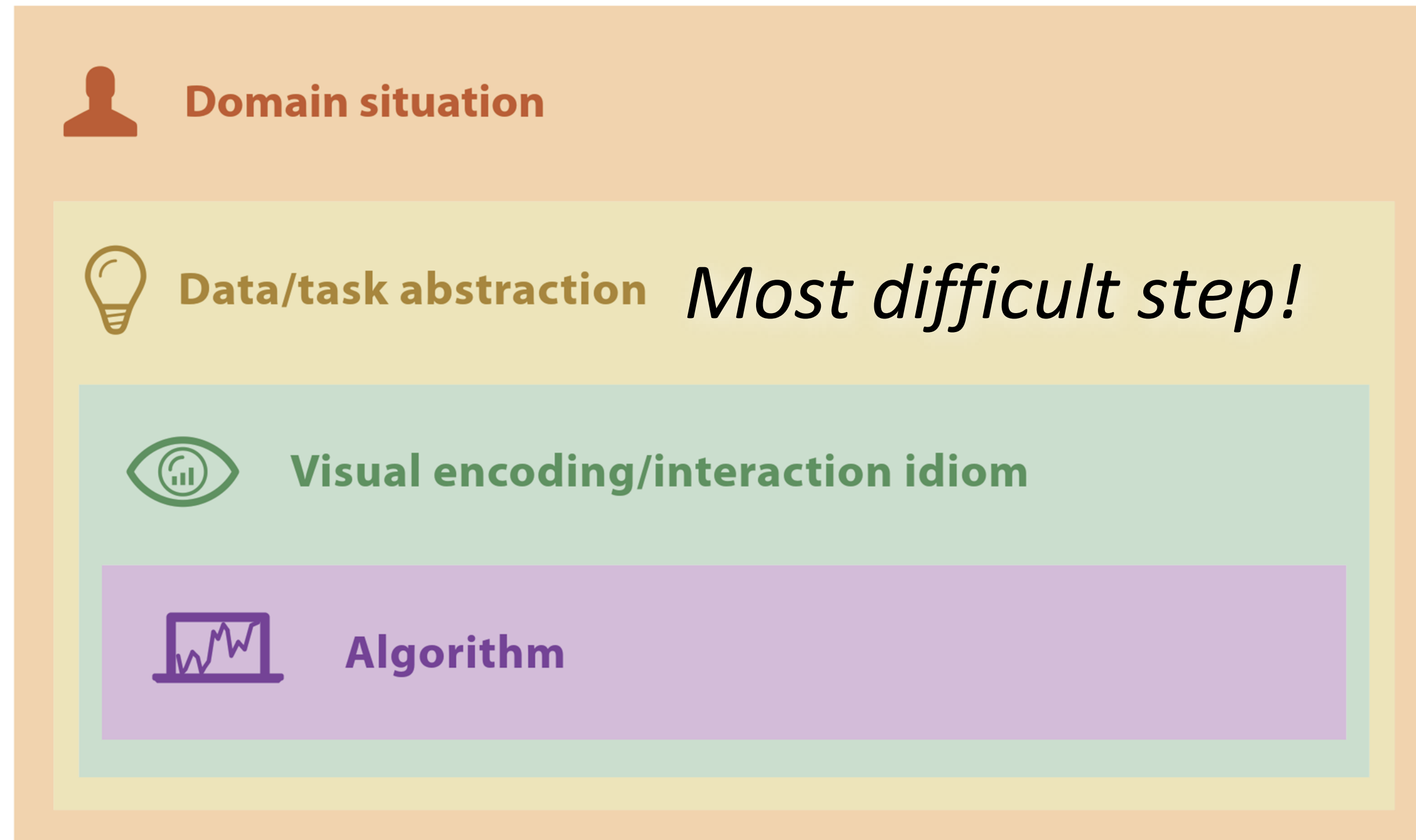
Nested Model

Design Study

Technique

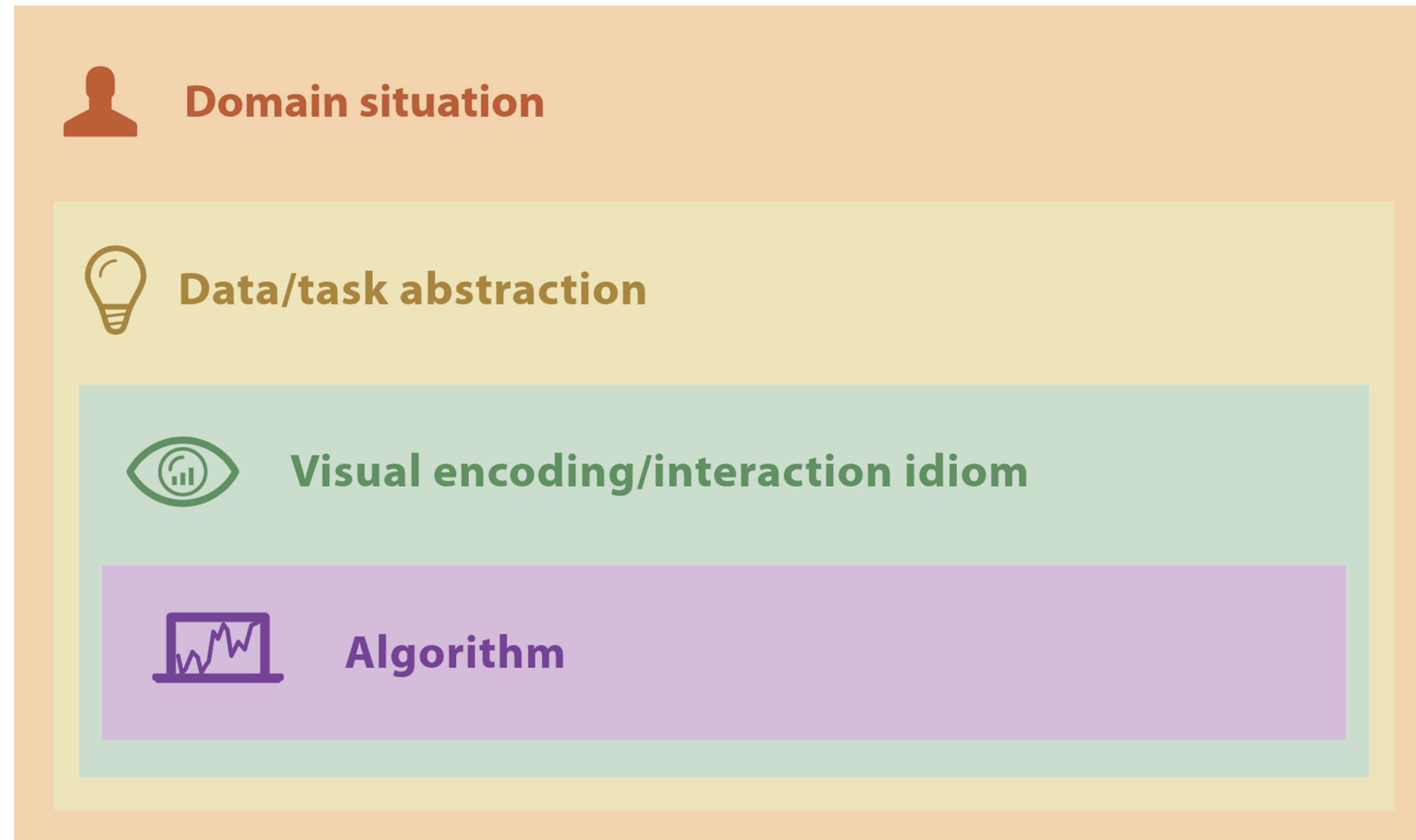
TOP-DOWN
“problem-
driven”

BOTTOM-UP
“technique
-driven”

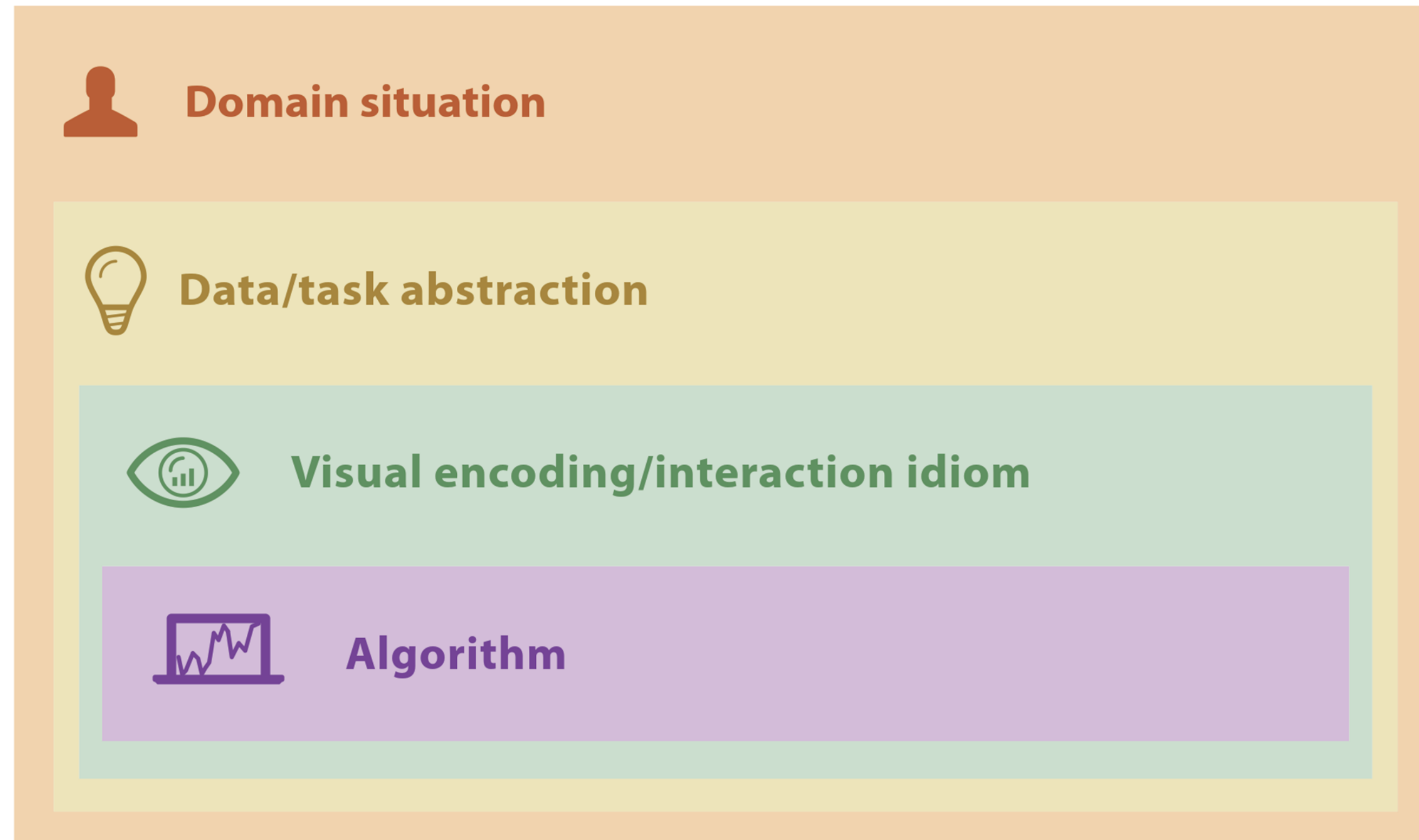


Nested Model

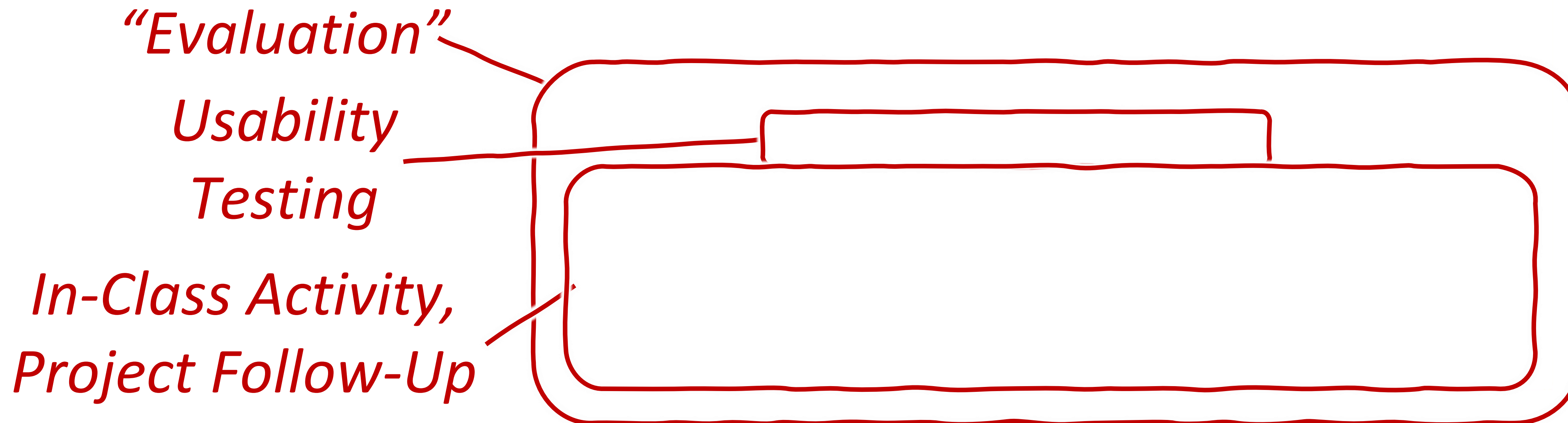
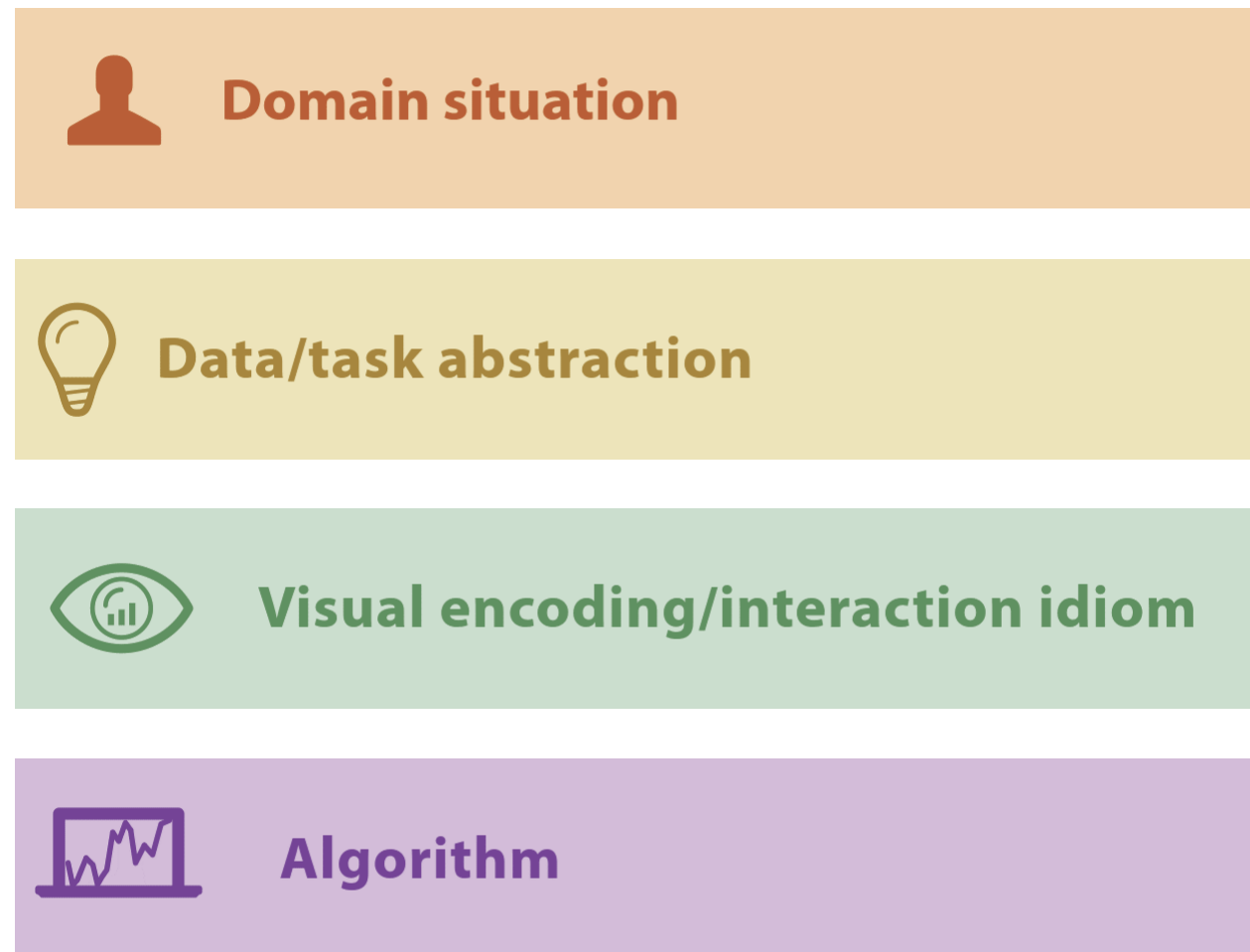
*Mistakes
propagate
through model!*



Threats to Validity

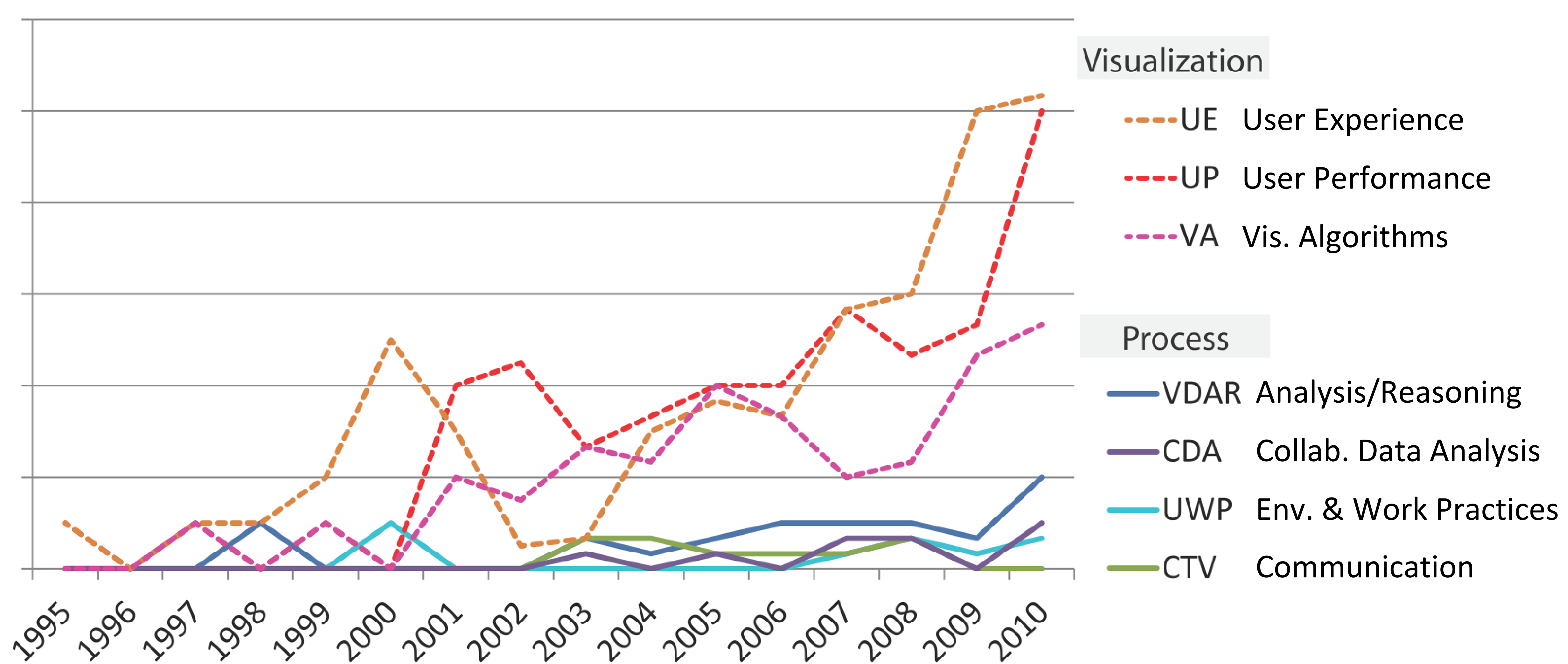


Threats to Validity *✓ Final Project validation*



EMPIRICAL STUDIES IN
INFORMATION VISUALIZATION:
SEVEN SCENARIOS

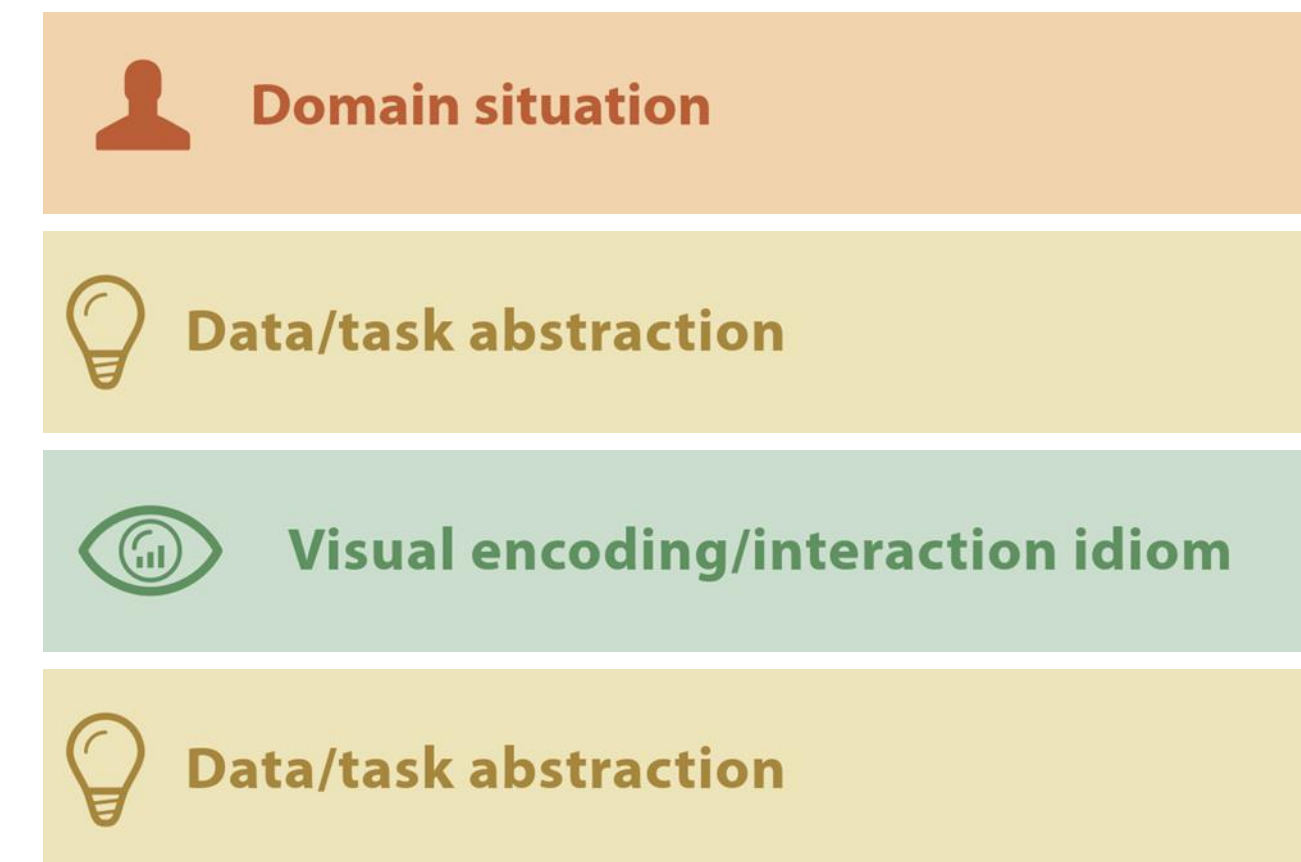
Empirical Studies in Information Visualization: Seven Scenarios



7 Evaluation Scenarios

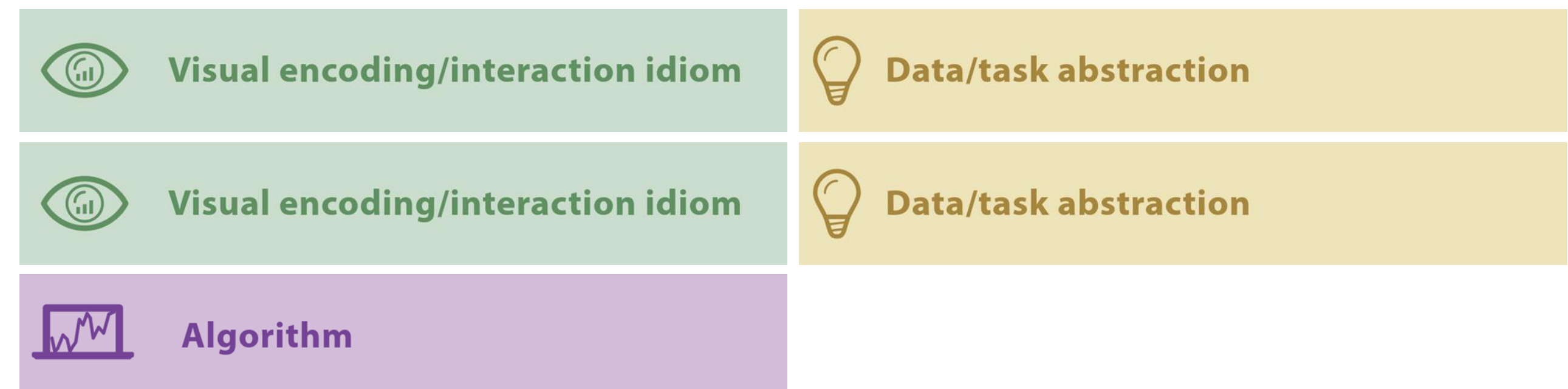
How to understand your data:

- Understanding Environments and Work Practices
- Evaluating Visual Data Analysis and Reasoning
- Evaluating Communication Through Visualization
- Evaluating Collaborative Data Analysis



How to understand your visualization:

- Evaluating User Performance
- Evaluating User Experience
- Evaluating Visualization Algorithms



Understanding environments and work practices



Domain situation

- Goals & outputs
 - Understand work, analysis, or information processing practices of people
 - Without software in use: inform design
 - With software in use: assess factors for adoption, how appropriated for future design
- Evaluation Questions
 - Context of use?
 - Integrate into which daily activities?
 - Supported analyses?
 - Characteristics of user group and environment?
 - What data & tasks?
 - What visualizations/tools used?
 - How current tools solve tasks?
 - Challenges and usage barrier?

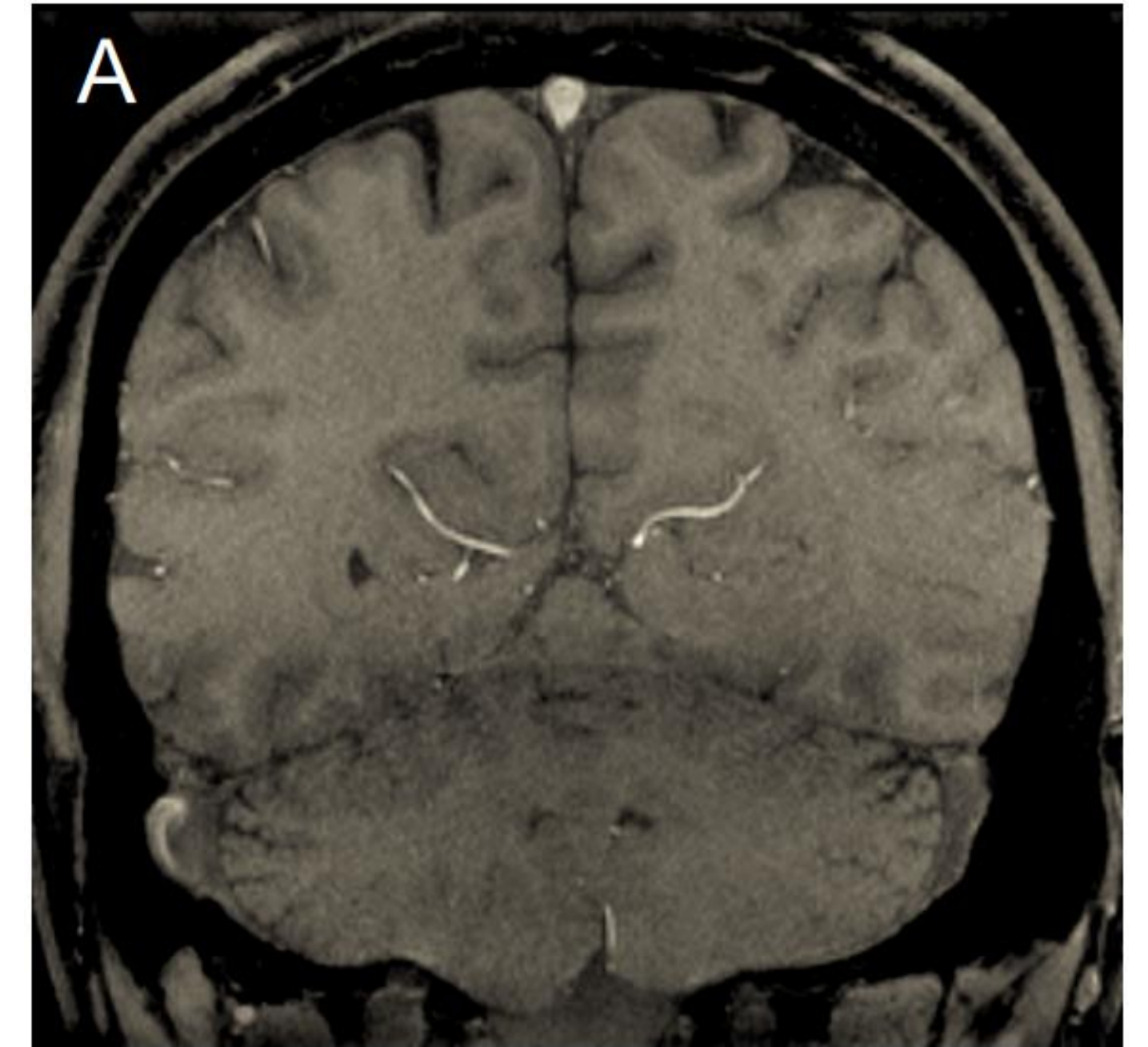
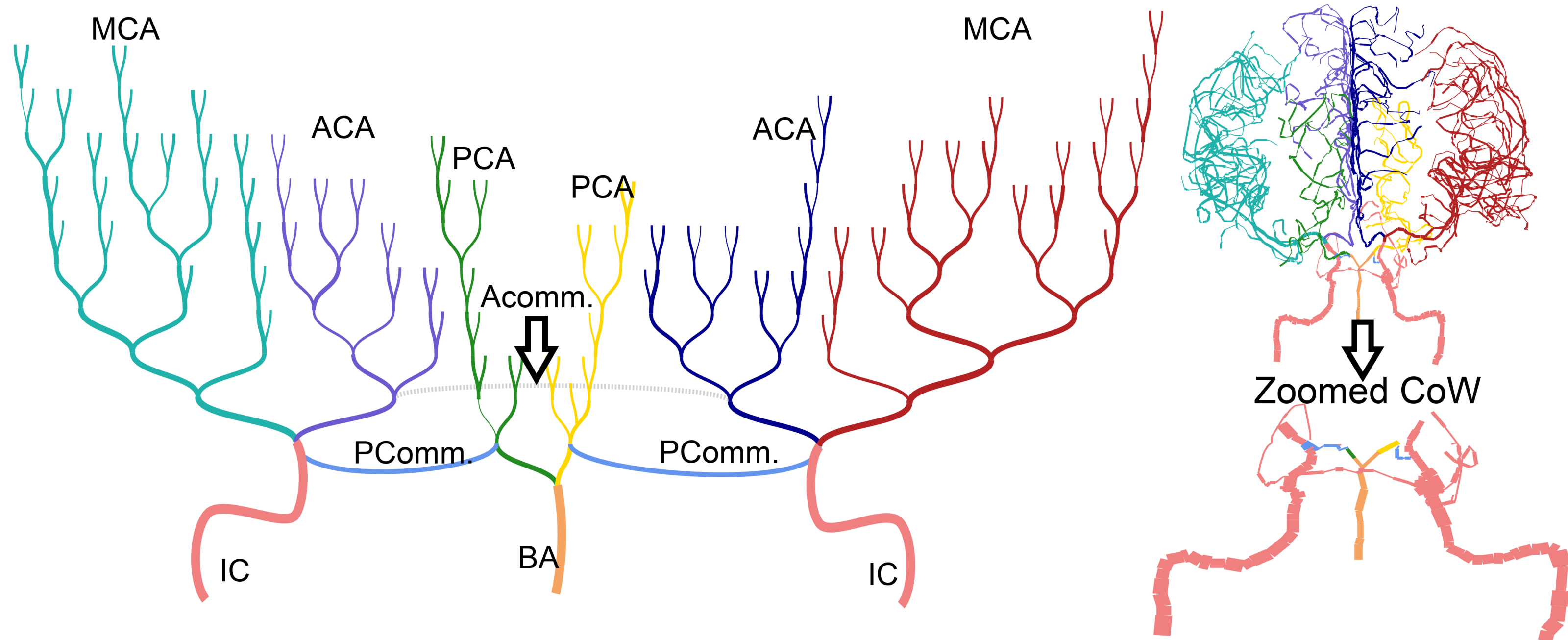
Understanding environments and work practices



Domain situation

- Methods
 - Field Observation
 - Real world, free use of tool
 - Derive requirements
 - Interviews
 - Contextual inquiry: interview then observe in routines, with little interference
 - Pick the right person
 - Laboratory context w/domain expert
 - Laboratory Observation
 - How people interact with each other, tools
 - More control of situation

Understanding environments and work practices: Example



Evaluating visual data analysis and reasoning



Data/task abstraction

- Goals & outputs
 - Assess visualization tool's ability to support visual analysis and reasoning
 - As a whole! Not just a technique
 - Quantifiable metrics or subjective feedback
- Evaluation Questions: Does it support...
 - Data exploration?
 - Knowledge discovery?
 - Hypothesis generation?
 - Decision making?

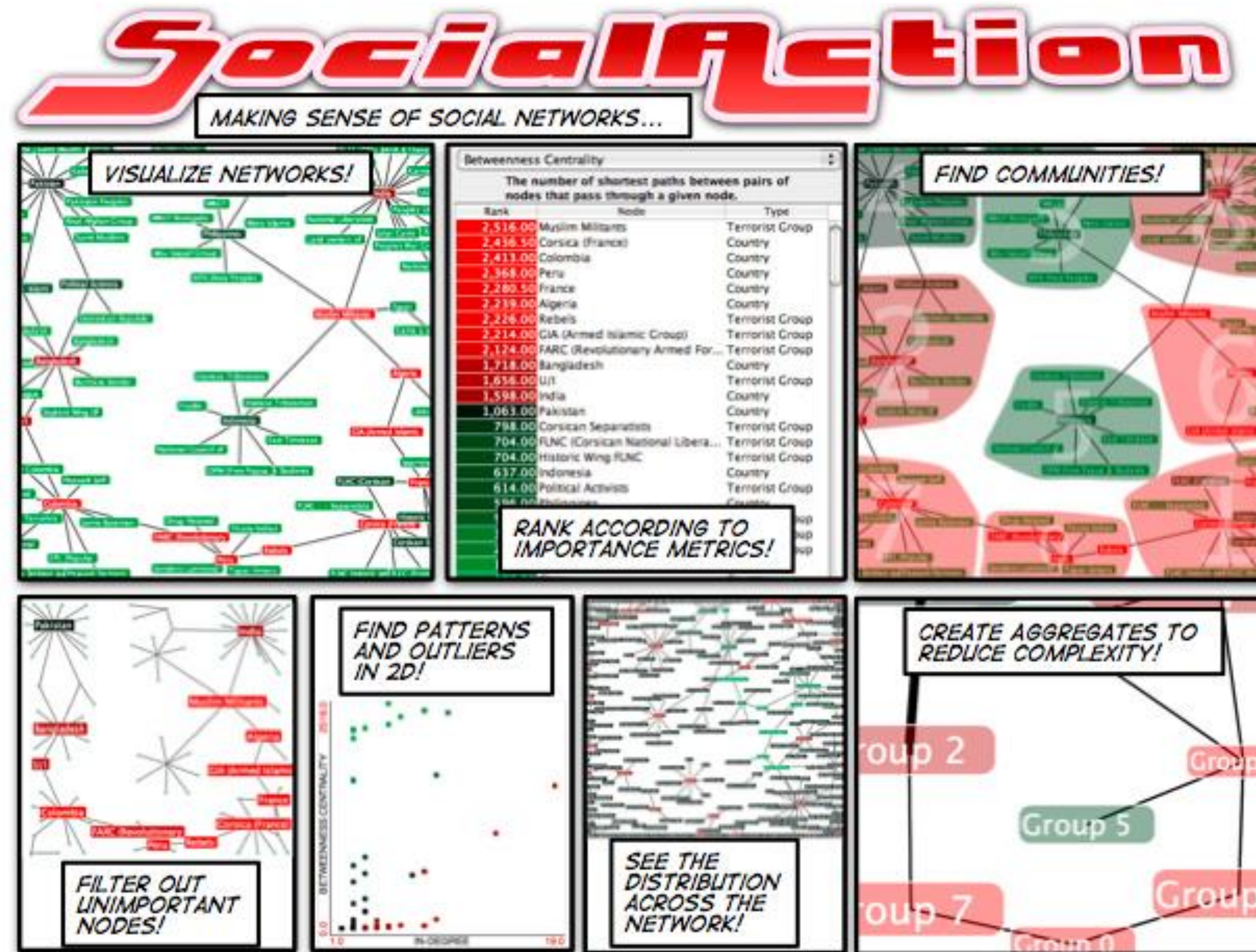
Evaluating visual data analysis and reasoning



Data/task abstraction

- Methods
 - Case studies
 - Motivated experts with own data in own environment
 - Can be longitudinal
 - Insight-Based ([Saraiya et al., 2004](#))
 - Unguided, diary, debriefing meetings
 - MILCS: Multidimensional In-depth Long-term Case studies (Shneiderman & Plaisant, 2006)
 - Guided, observations, interviews, surveys, automated logging
 - Assess interface efficacy, user performance, interface utility
 - Improve system during
 - Lab observations and interviews
 - Code results
 - Think aloud
 - Controlled Experiment
 - Isolate important factors

Evaluating visual data analysis and reasoning



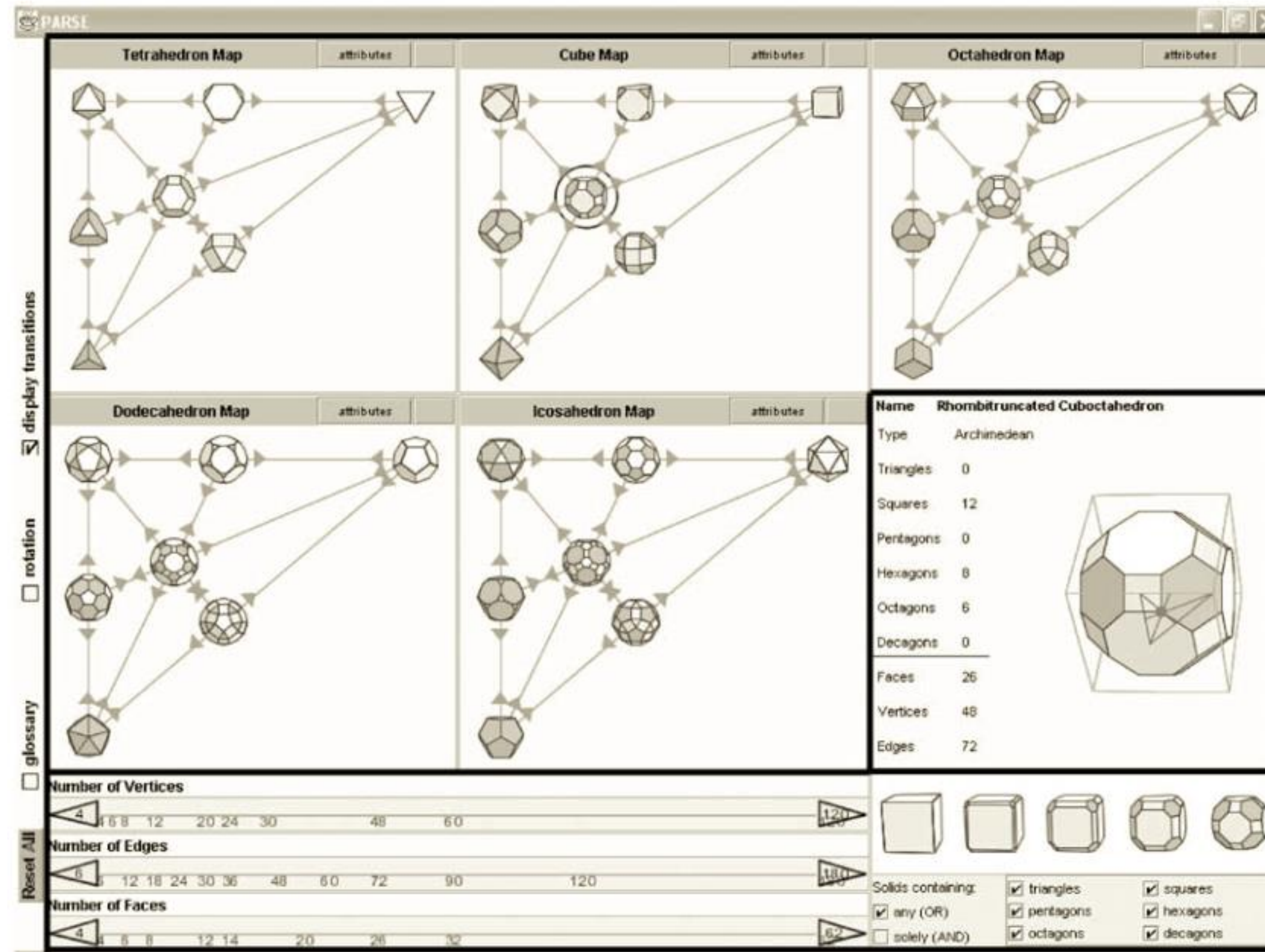
Evaluating communication through visualization



Visual encoding/interaction idiom

- Goals & outputs
 - How effectively is a message delivered and acquired
- Evaluation Questions
 - Quantitative: learning rate, information retention and accuracy
 - Qualitative: interaction patterns
- Methods
 - Controlled experiments
 - Field observation & interviews

Evaluating communication through visualization: Example



Evaluating Collaborative Data Analysis



Data/task abstraction

- Goals & outputs
 - Evaluate support for taskwork and teamwork
 - Holistic understanding of group work processes or tool use
 - Derive design implications
- Evaluation Questions
 - Effective and efficient?
 - Satisfactorily support or stimulate group sensemaking?
 - Support group insight?
 - Is social exchange and communication facilitated?
 - How is the tool used? Features, patterns...
 - What is the process? User requirements?

Evaluating Collaborative Data Analysis



Data/task abstraction

- Methods
 - Context critical, but early formative studies less dependant
 - Heuristic evaluation
 - Heuristics: actions, mechanics, interactions, locales needed
 - Log analysis
 - Distributed or web-based tools
 - Combine with questionnaire or interview
 - Hard to evaluate unlogged & qualitative aspects
 - Field or laboratory observation
 - Involve group interactions and harmony/disharmony
 - Combine with insight-based?

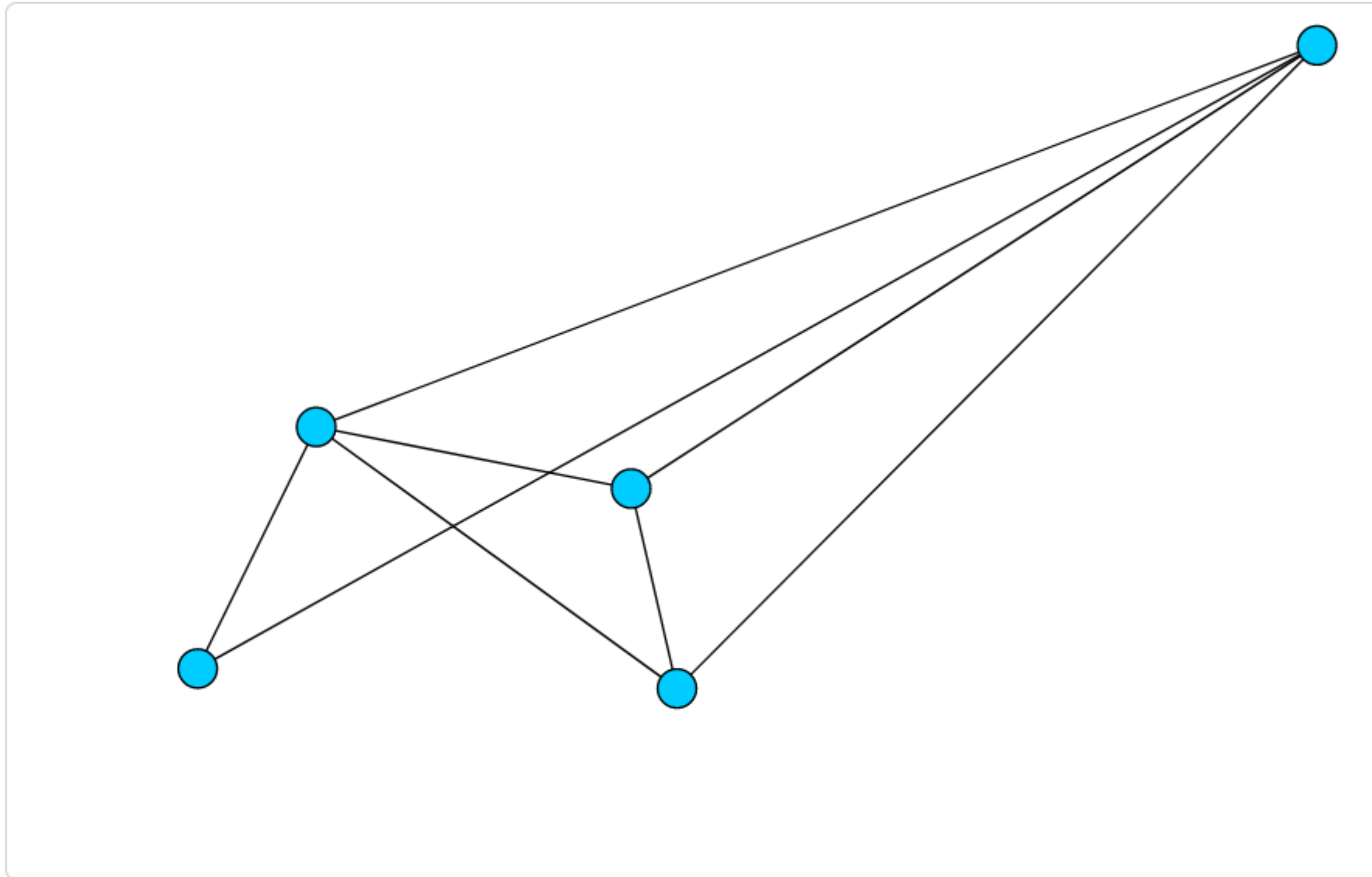
Evaluating Collaborative Data Analysis: Examples

Planarity Party

Can you untangle the graph? See if you can position the vertices so that no two lines cross.

Level 1. Number of line crossings detected: 2.

0 moves. [Next Level](#)



Evaluating User Performance

- Goals & outputs
 - Measure specific features
 - Time, accuracy, and error; work quality (if quantifiable); memorability
 - Descriptive statistics results
- Evaluation Questions
 - What are the limits of human perception and cognition?
 - How do techniques compare?
- Methods
 - Controlled experiment → design guideline, model, head-to-head
 - Few variables
 - Simple tasks
 - Individual differences matter
 - Field logs
 - Suggest improvements, recommendation systems



Data/task abstraction



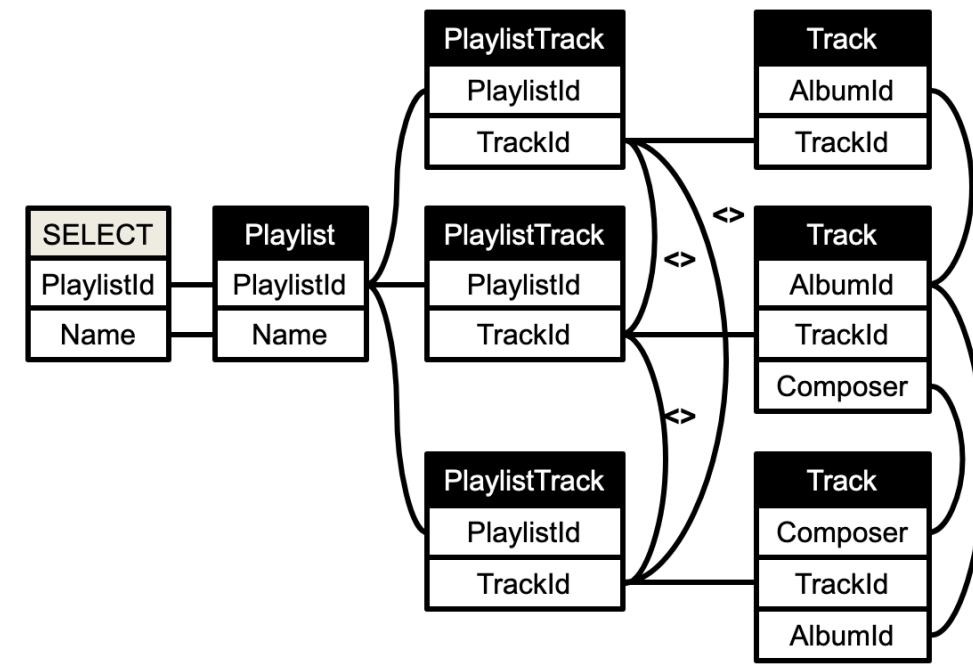
Visual encoding/interaction idiom

Evaluating User Performance: Examples

Question 6 / 12

Time remaining: 48:39 minutes

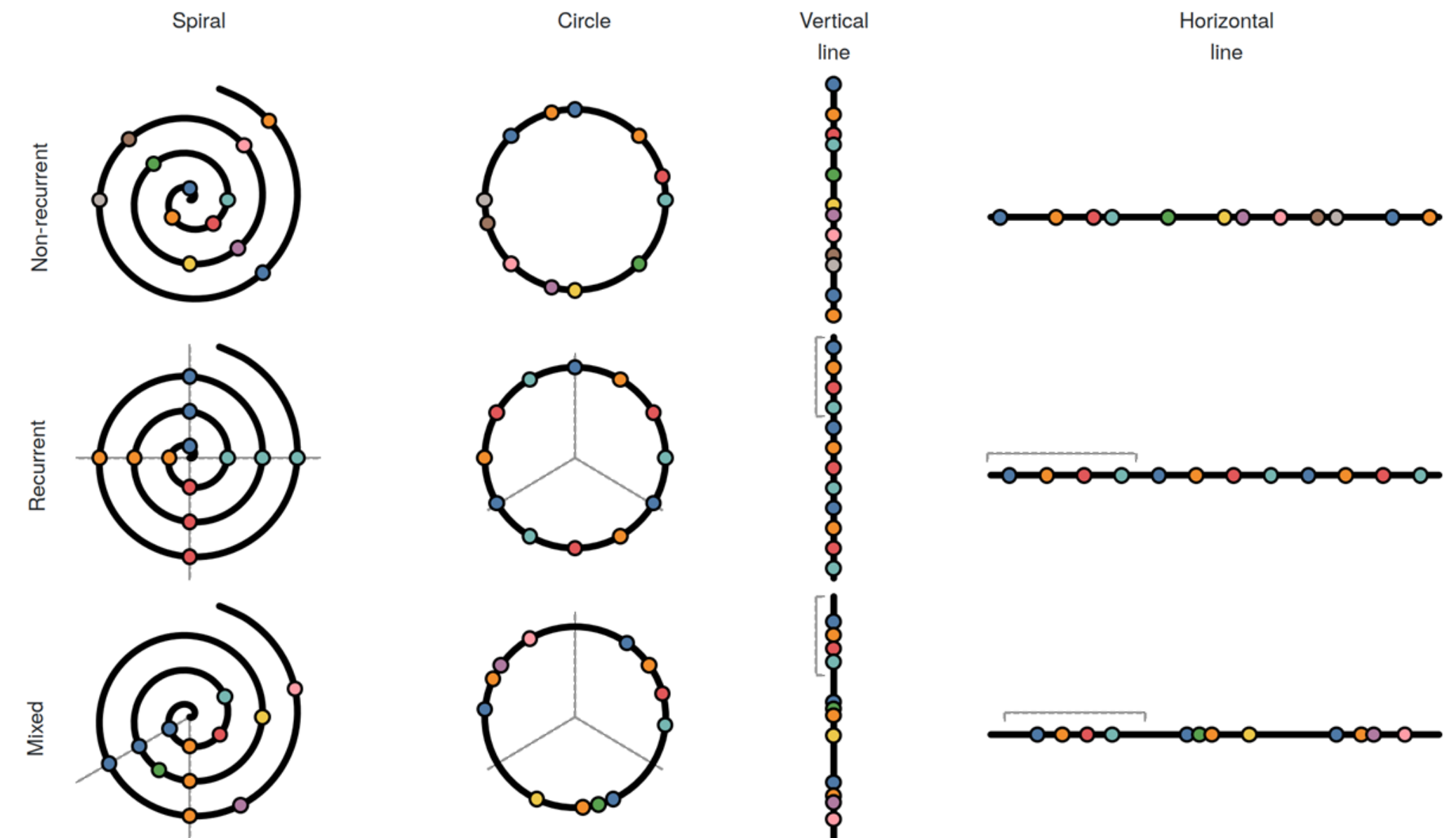
```
SELECT P.PlaylistId, P.Name
FROM Playlist P, PlaylistTrack PT1,
PlaylistTrack PT2, PlaylistTrack PT3,
Track T1, Track T2, Track T3
WHERE P.PlaylistId = PT1.PlaylistId
AND P.PlaylistId = PT2.PlaylistId
AND P.PlaylistId = PT3.PlaylistId
AND PT1.TrackId <> PT2.TrackId
AND PT2.TrackId <> PT3.TrackId
AND PT1.TrackId <> PT3.TrackId
AND PT1.TrackId = T1.TrackId
AND PT2.TrackId = T2.TrackId
AND PT3.TrackId = T3.TrackId
AND T1.AlbumId = T2.AlbumId
AND T2.AlbumId = T3.AlbumId
AND T2.Composer = T3.Composer;
```



- Find playlists that have at least 3 different tracks that are in the same album and they are all made by the same composer.
- Find playlists that have at least 3 different tracks so that at least 2 of them are in the same album but all 3 tracks are made by the same composer.
- Find playlists that have at least 3 different tracks so that at least 2 of them are in the same album and made by the same composer.
- Find playlists that have at least 3 different tracks that are in the same album and at least 2 of them are made by the same composer.

Submit

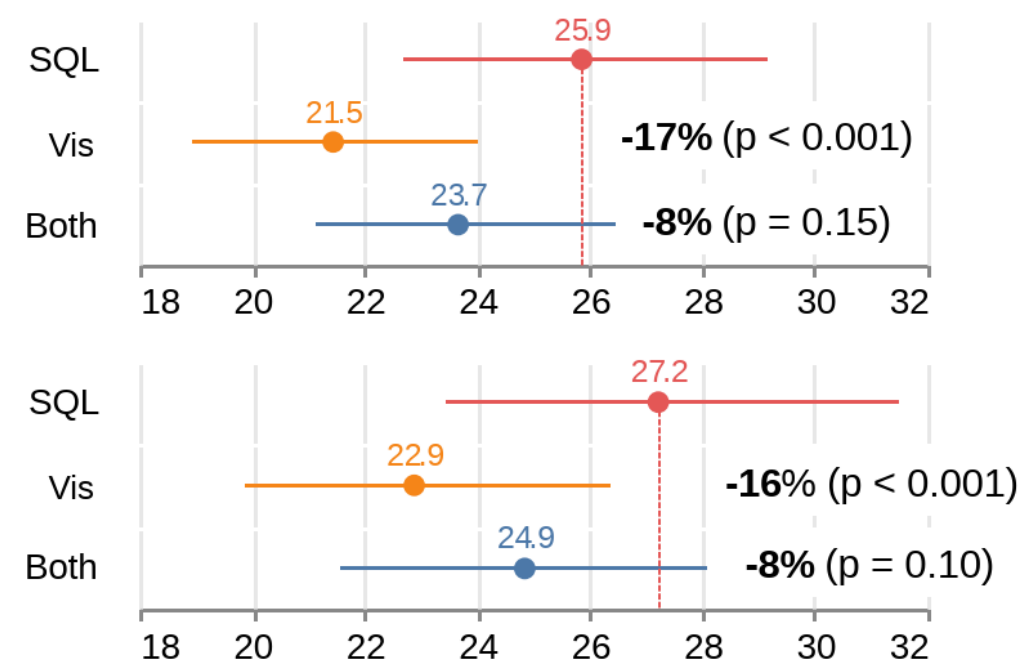
[Tutorial \(PDF\)](#)



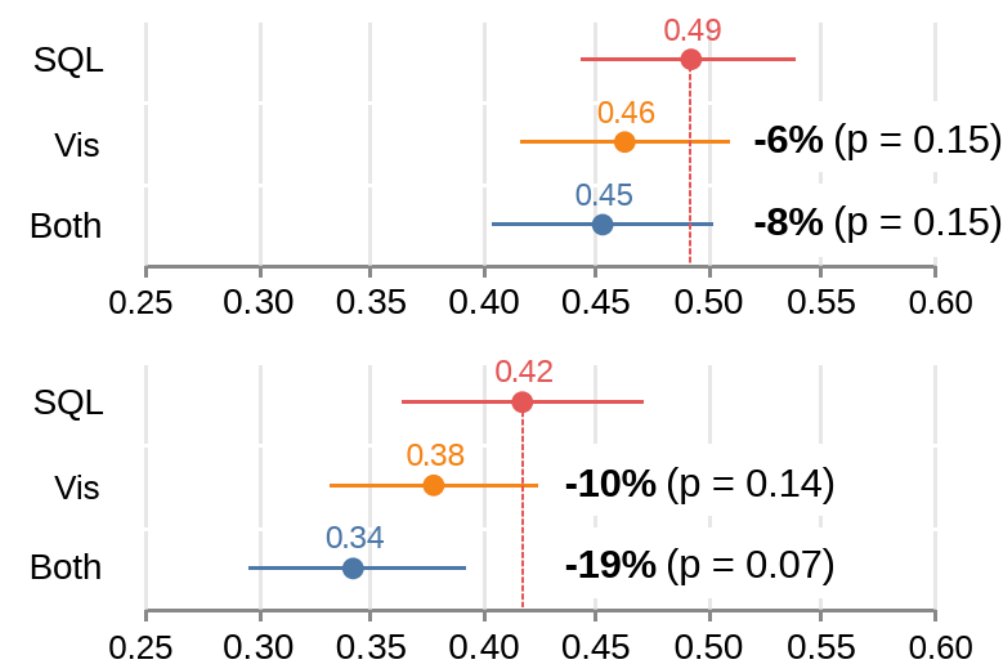
All partici- pants (n=102)

Accepted partici- pants (n=70)

Mean Time per question (seconds)



Mean Error



Dataset Sample question Mean Completion Time Mean Per-Worker Log Change in Completion Time $\ln(\text{Shape} / \text{LH})$ Mean Proportion Correct

Task	Dataset	Sample question	Mean Completion Time	Mean Per-Worker Log Change in Completion Time $\ln(\text{Shape} / \text{LH})$	Mean Proportion Correct
Task 1: When	Mixed	At what time was the Group Meeting on Wednesday?	[Plot]	[Plot]	[Plot]
	Non-recurrent	In what year was writing invented?	[Plot]	[Plot]	[Plot]
	Recurrent	In which season do you plant Puffapod?	[Plot]	[Plot]	[Plot]

Evaluating User Experience

- Goals & outputs

- Inform design: uncover gaps in functionality, limitations, directions for improvement
- Subjective: user responses
 - Effectiveness, efficiency, correctness, satisfaction, trust, features liked/disliked
- Objective: body sensors, eye tracking

- Evaluation Questions

- Features: useful, missing, to rework?
- Are there limitations that hinder adoption?
- Is the tool understandable/learnable?



Data/task abstraction



Visual encoding/interaction idiom

Evaluating User Experience

- Methods
 - Informal evaluation
 - Demo for domain experts (usually) and collect feedback
 - Usability test
 - Watch (video) how participants perform set of tasks to perfect design
 - Take note of behaviors, remarks, problems
 - Carefully prepare tasks, interview script, questionnaires
 - Field observation
 - Understand interaction in real setting
 - Laboratory questionnaire
 - Likert scale
 - Open ended

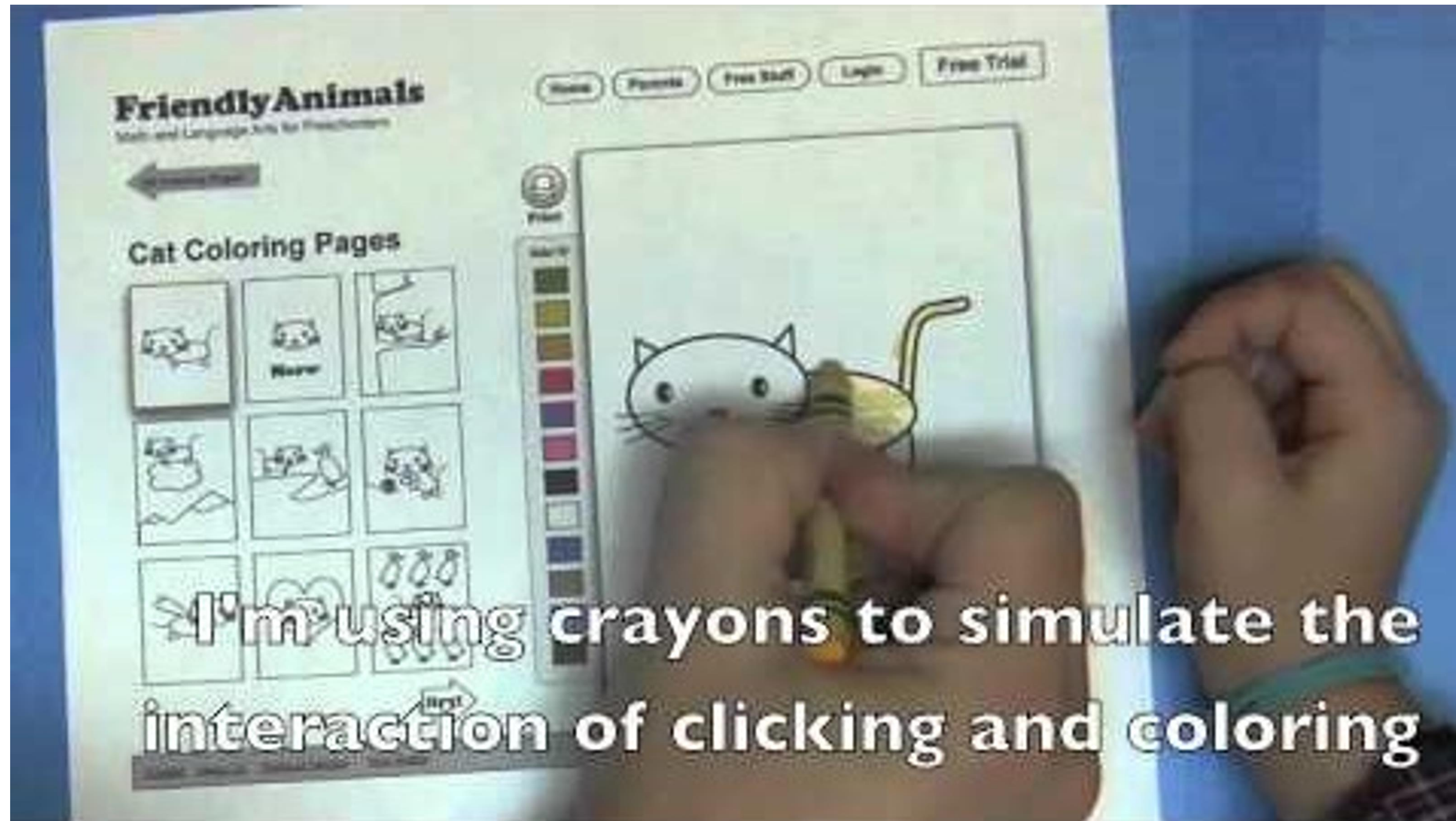


Data/task abstraction



Visual encoding/interaction idiom

Evaluating User Experience: Example



Evaluating Visualization Algorithms

- Goals & outputs
 - Quantitatively or qualitatively judge generated output quality (metrics) & performance
 - How scores vs. alternatives
 - Explore limits & behavior
- Evaluation Questions
 - Which shows interesting patterns best?
 - Which is more truthful?
 - Which is less cluttered?
 - Faster, less memory, less money?
 - How does it scale?
 - Extreme cases?



Data/task abstraction



Visual encoding/interaction idiom

Evaluating Visualization Algorithms

- Methods
 - Visualization quality assessment
 - Readability metrics, image quality measures
 - Algorithmic performance
 - Varied data, size, complexity, corner cases
 - Benchmark data sets



Data/task abstraction



Visual encoding/interaction idiom

Evaluating Visualization Algorithms: Example

Type	Name	$recn_{\Gamma}$				
		GVA	FM ³	FMS	ACE	HDE
Kind Artificial	rnd_grid_032	<u>3.82</u>	0	0	0	<0.01
	rnd_grid_100	<u>14.75</u>	0	0	0	<0.01
	rnd_grid_320	<u>181.51</u>	0	(N)	<0.01	<0.01
	sierpinski_06	<u>2.00</u>	0.05	<0.01	0	0.02
	sierpinski_08	<u>9.49</u>	0.07	0.01	0.02	0.08
	sierpinski_10	<u>99.97</u>	0.09	(N)	0.27	0.01
Kind Real World	crack	<u>30.82</u>	<0.01	(N)	0	0.07
	fe_pwt	<u>150.70</u>	2.45	(N)	(N)	1.61
	finan_512	<u>301.25</u>	18.81	(N)	12.27	21.27
	fe_ocean	<u>622.48</u>	7.13	(N)	9.07	8.24
Challenging Artificial	tree_06_04	<u>2.21</u>	1.16	7.89	0.01	0
	tree_06_05	9.33	1.89	11.48	0	<u>22.92</u>
	tree_06_06	70.68	3.31	(N)	4.16	<u>128.82</u>
	snowflake_A	<u>0.63</u>	0	0.10	<0.01	0.62
	snowflake_B	1.46	0	<u>8.18</u>	(N)	6.92
	snowflake_C	15.53	0	(N)	(N)	<u>195.87</u>
	spider_A	15.62	<u>16.55</u>	1.17	6.60	1.25
	spider_B	<u>154.70</u>	132.96	1.64	0	0
	spider_C	<u>2522.89</u>	1029.64	(N)	0	0
	flower_A	46.71	<u>49.08</u>	5.63	0.26	0.55
flower_B	<u>64.90</u>	51.57	1.90	0.06	0.34	
flower_C	<u>578.22</u>	53.39	(N)	(N)	0.30	
Challenging Real World	ug_380	<u>22.93</u>	19.55	13.67	20.99	1.35
	esslingen	<u>47.52</u>	23.71	28.42	20.81	3.89
	add_32	<u>8.65</u>	1.69	5.75	0.89	5.80
	dg_1087	1.74	< 0.01	<u>37.07</u>	5.92	6.49
	bcsstk_33	720.94	376.18	<u>4171.05</u>	413.56	113.86
bcsstk_31	<u>708.69</u>	94.26	(N)	63.00	611.21	

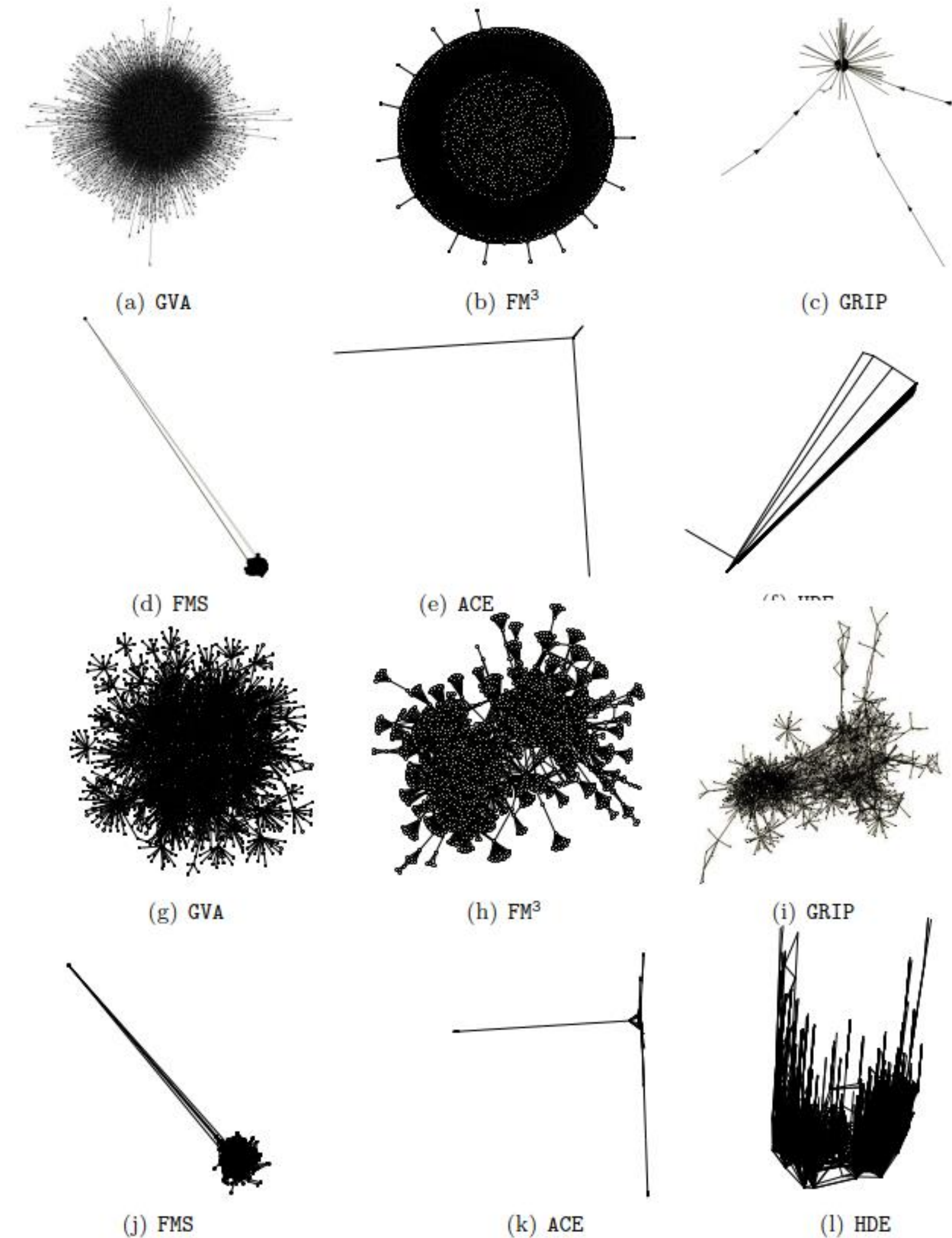


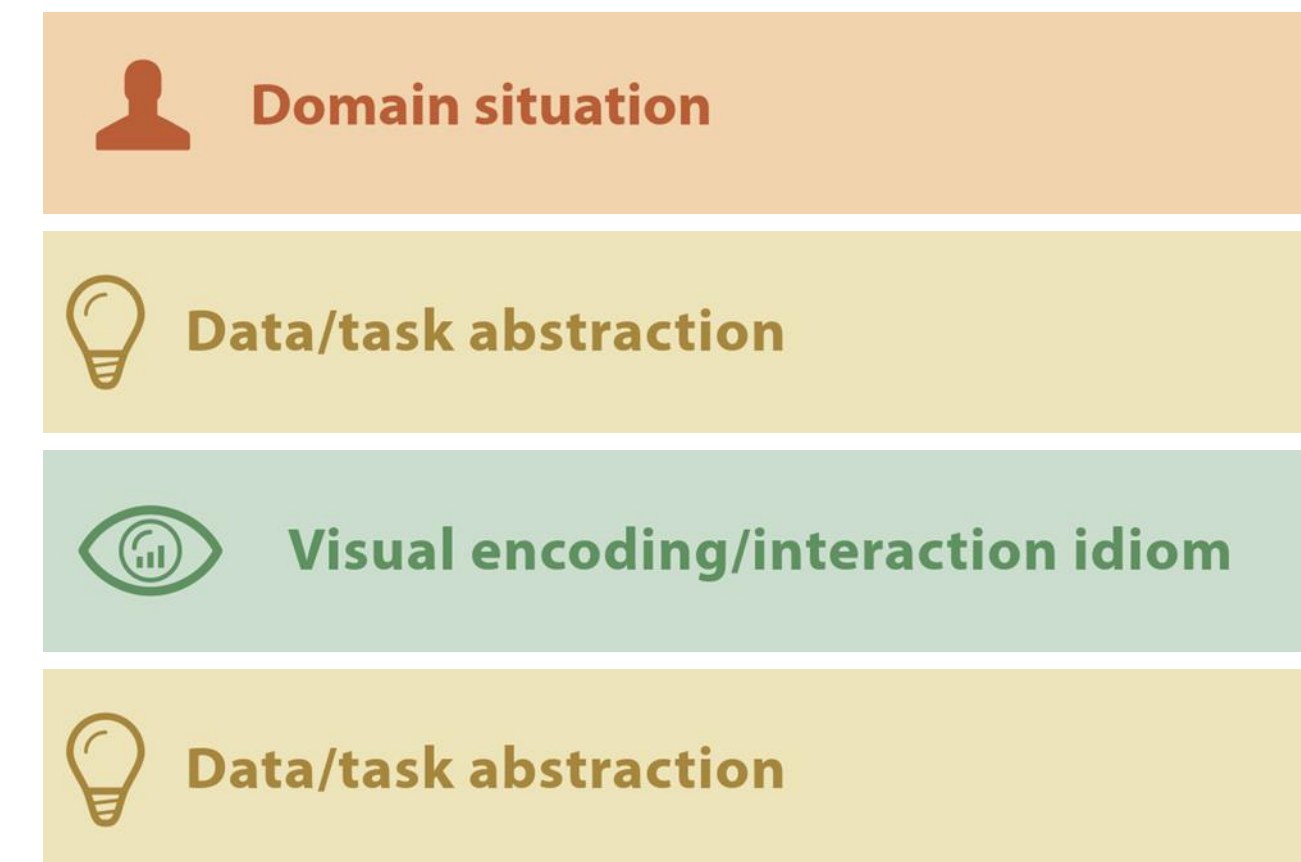
Figure 8: (a)-(f) Drawings of dg_1087 and (g)-(l) esslingen generated by different algorithms

Table 3: The relative edge-crossing numbers ($recn_{\Gamma}$) of the drawings Γ computed by the tested algorithms. The entry (N) indicates that no drawing was computed. Best values are printed bold. Worst values are underlined.

7 Evaluation Scenarios

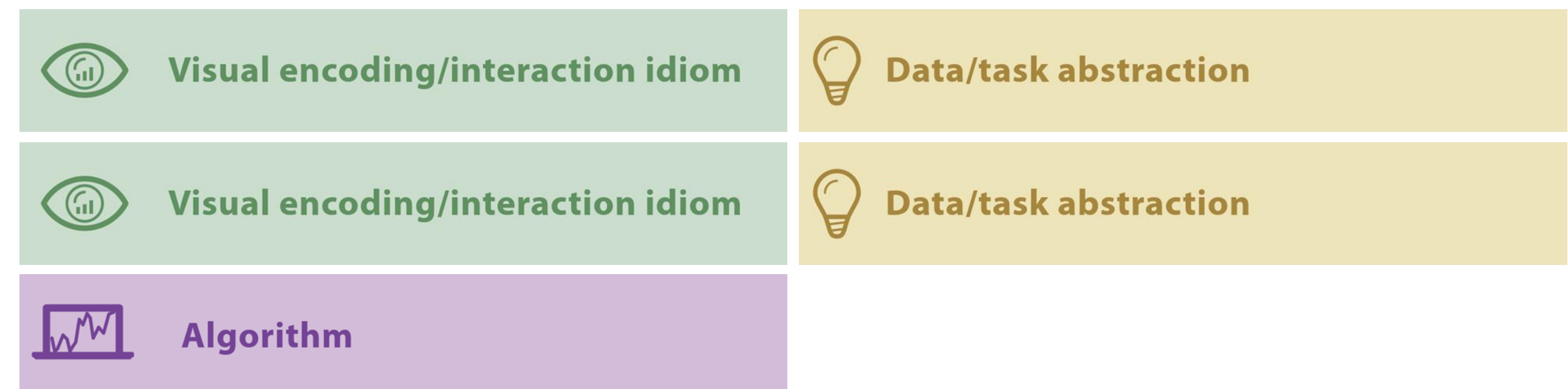
How to understand your data:

- Understanding Environments and Work Practices
- Evaluating Visual Data Analysis and Reasoning
- Evaluating Communication Through Visualization
- Evaluating Collaborative Data Analysis



How to understand your visualization:

- Evaluating User Performance
- Evaluating User Experience
- Evaluating Visualization Algorithms



7 Evaluation Scenarios

How to understand your data:

- Understanding Environments and Work Practices Field Observations, Interviews
- Evaluating Visual Data Analysis and Reasoning Case Studies, Controlled Experiment
- Evaluating Communication Through Visualization Field Observation, Controlled Experiment
- Evaluating Collaborative Data Analysis Field Observation, Heuristic Evaluation, Log Analysis

How to understand your visualization:

- Evaluating User Performance Controlled Experiment, Log Analysis
- Evaluating User Experience Informal Evaluation, Usability Test, Field Observation
- Evaluating Visualization Algorithms Visualization Quality Assessment, Algorithm Performance

In-Class Validation — Final Project Evaluation

~35 min