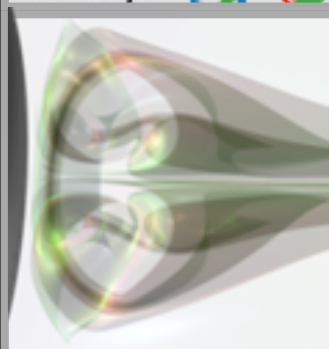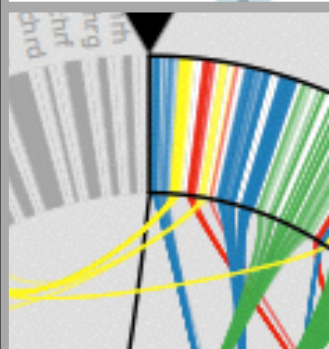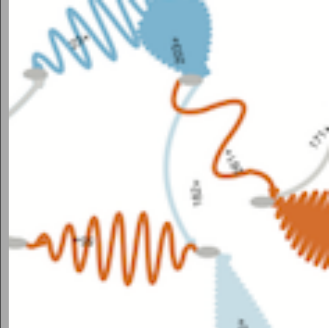# TEXT & SETS

Miriah Meyer
*University of Utah*

administrivia . . .

-parallel coordinates due next Thursday

last time . . .

# dataset types

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists |
|---|---|---|---|---|
| Items | Items (nodes) | Grids | Items | Items |
| Attributes | Links | Positions | Positions | |
| | Attributes | Attributes | | |

Attributes (columns)

Items (rows)

Cell containing value

→ *Multidimensional Table*

Key 1

Key 2

Value in cell

Attributes

Link

Node (item)

Cell

→ *Trees*

Grid of positions

Attributes (columns)

Value in cell

Position

# GRAPHS & TREES

**-graphs**
  -model relations amount data
  - *nodes* and *edges*

**-trees**
  -graphs with hierarchical structure
  -nodes as *parents* and *children*

flare

animate
- interpolate
  - ArrayInterpolator
  - ColorInterpolator
  - DateInterpolator
  - Interpolator
  - MatrixInterpolator
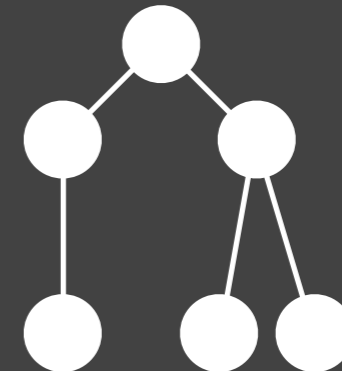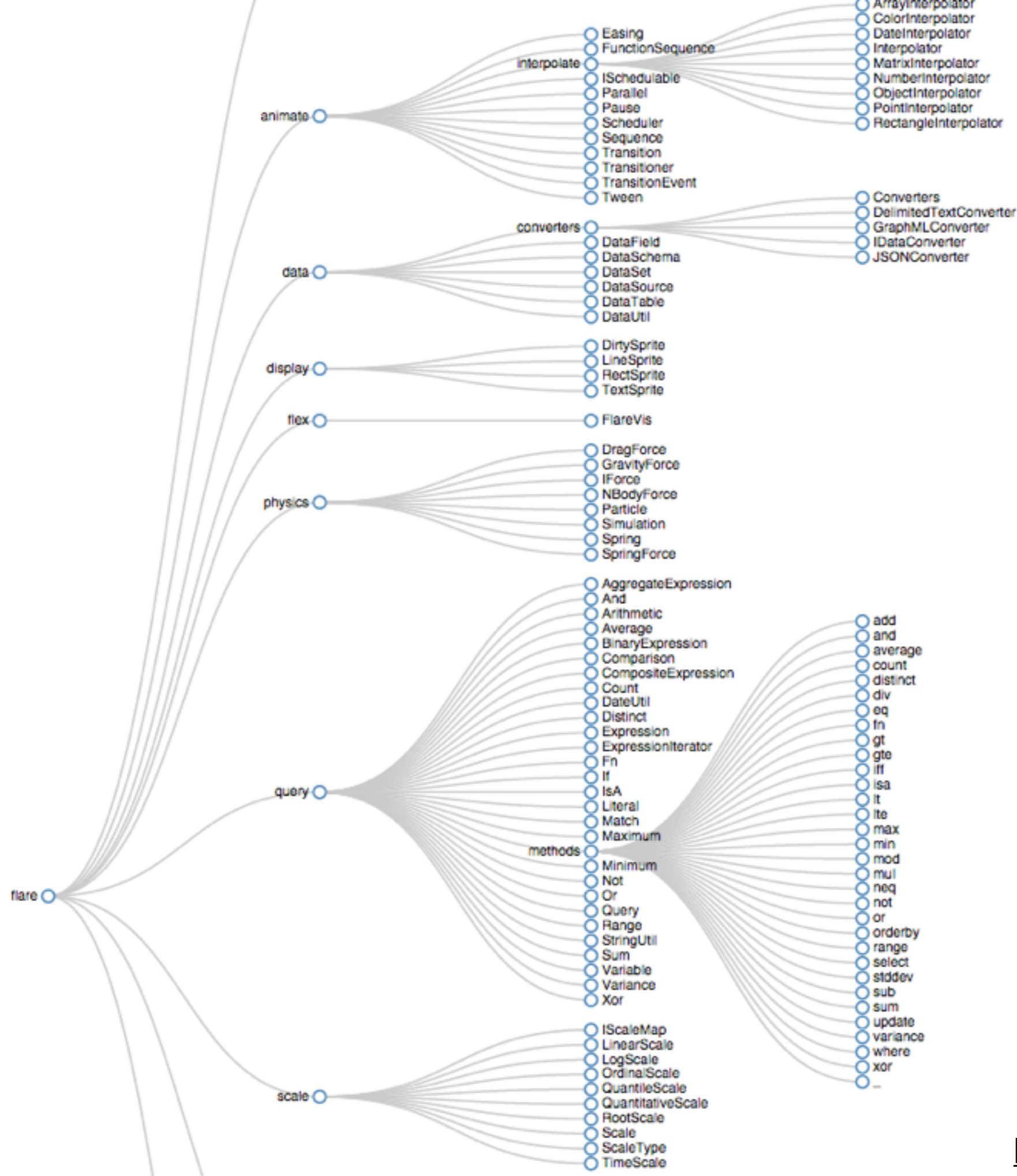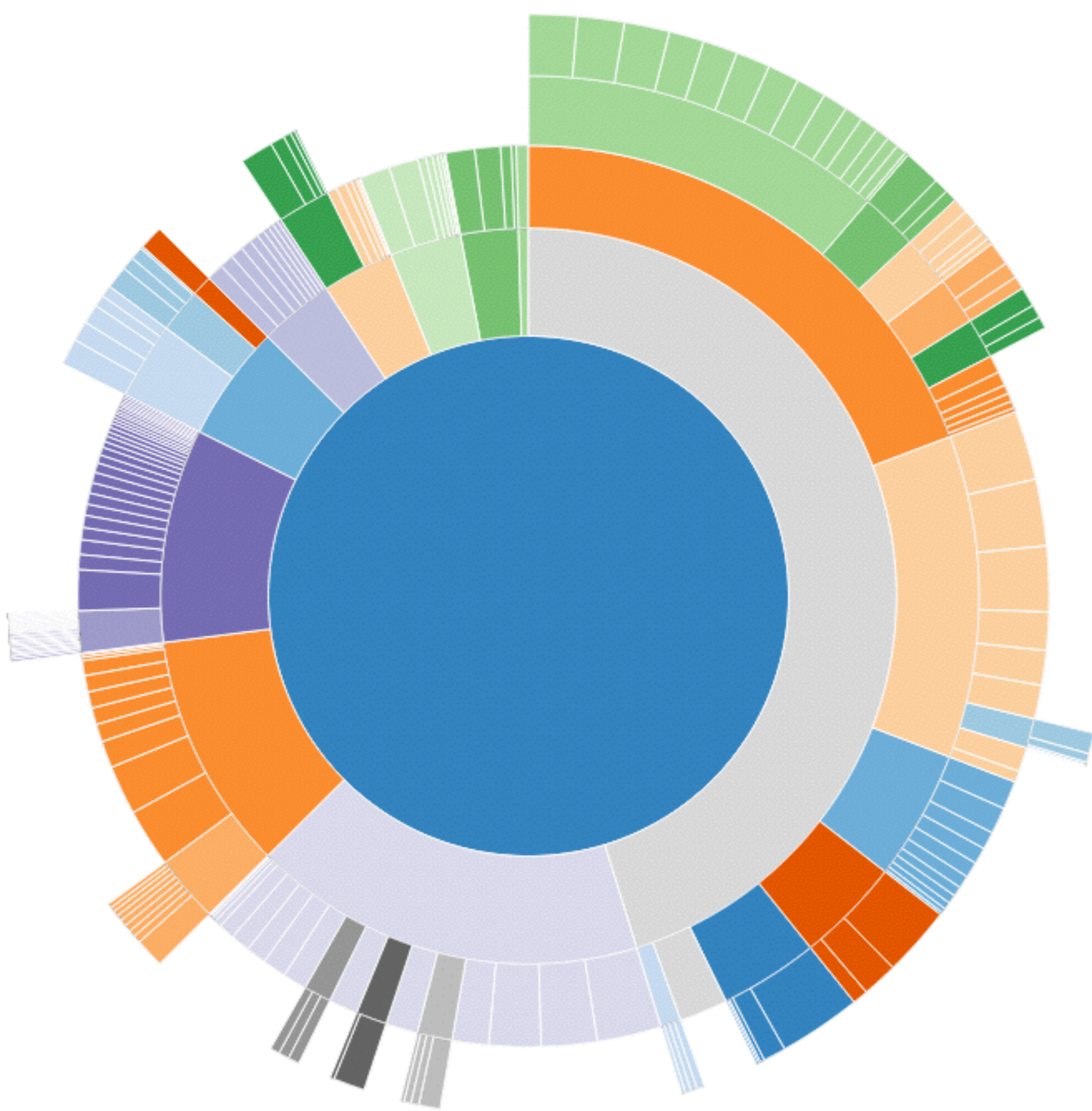  - NumberInterpolator
  - ObjectInterpolator
  - PointInterpolator
  - RectangleInterpolator
- Easing
- FunctionSequence
- ISchedulable
- Parallel
- Pause
- Scheduler
- Sequence
- Transition
- Transitioner
- TransitionEvent
- Tween

data
- converters
  - Converters
  - DelimitedTextConverter
  - GraphMLConverter
  - IDataConverter
  - JSONConverter
- DataField
- DataSchema
- DataSet
- DataSource
- DataTable
- DataUtil

display
- DirtySprite
- LineSprite
- RectSprite
- TextSprite

flex
- FlareVis

physics
- DragForce
- GravityForce
- IForce
- NBodyForce
- Particle
- Simulation
- Spring
- SpringForce

query
- AggregateExpression
- And
- Arithmetic
- Average
- BinaryExpression
- Comparison
- CompositeExpression
- Count
- DateUtil
- Distinct
- Expression
- ExpressionIterator
- Fn
- If
- IsA
- Literal
- Match
- Maximum
- methods
  - add
  - and
  - average
  - count
  - distinct
  - div
  - eq
  - fn
  - gt
  - gte
  - iff
  - isa
  - lt
  - lte
  - max
  - min
  - mod
  - mul
  - neq
  - not
  - or
  - orderby
  - range
  - select
  - stddev
  - sub
  - sum
  - update
  - variance
  - where
  - xor
  - _
- Minimum
- Not
- Or
- Query
- Range
- StringUtil
- Sum
- Variable
- Variance
- Xor

scale
- IScaleMap
- LinearScale
- LogScale
- OrdinalScale
- QuantileScale
- QuantitativeScale
- RootScale
- Scale
- ScaleType
- TimeScale

http://bl.ocks.org/mbostock/4339184

http://bl.ocks.org/mbostock/raw/4348373/

# VISUALIZING GRAPHS

- **node link layouts**
  - Reingold-Tilford (trees only)
  - Sugiyama (directed acyclic graphs)
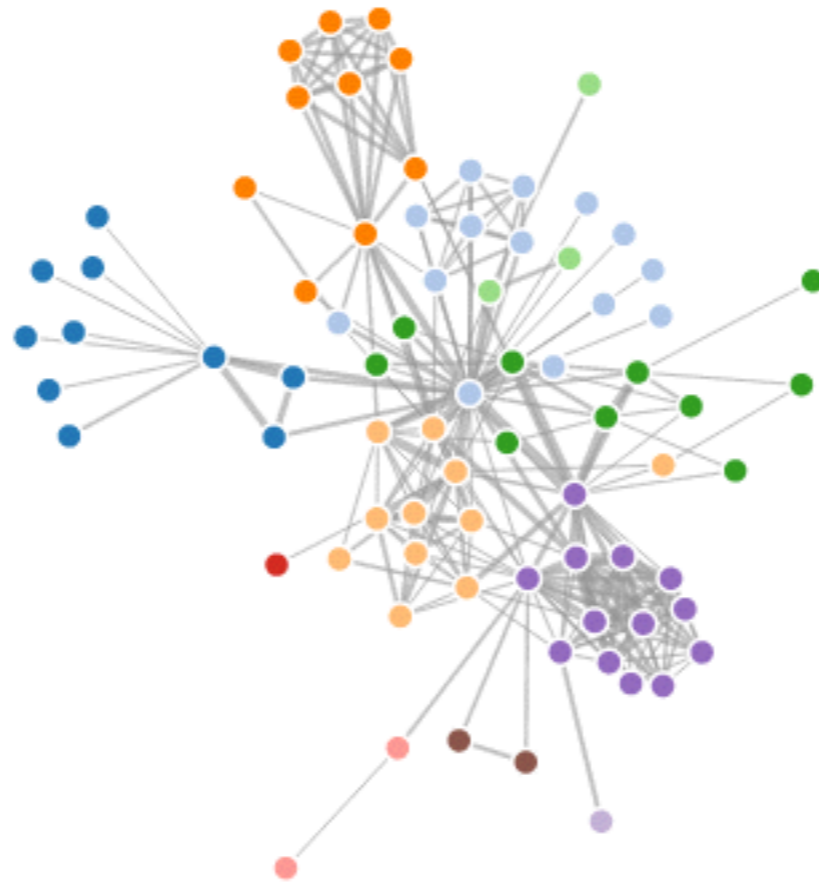  - Force directed
  - Attribute-based

- **adjacency matrices**

- **aggregate views**
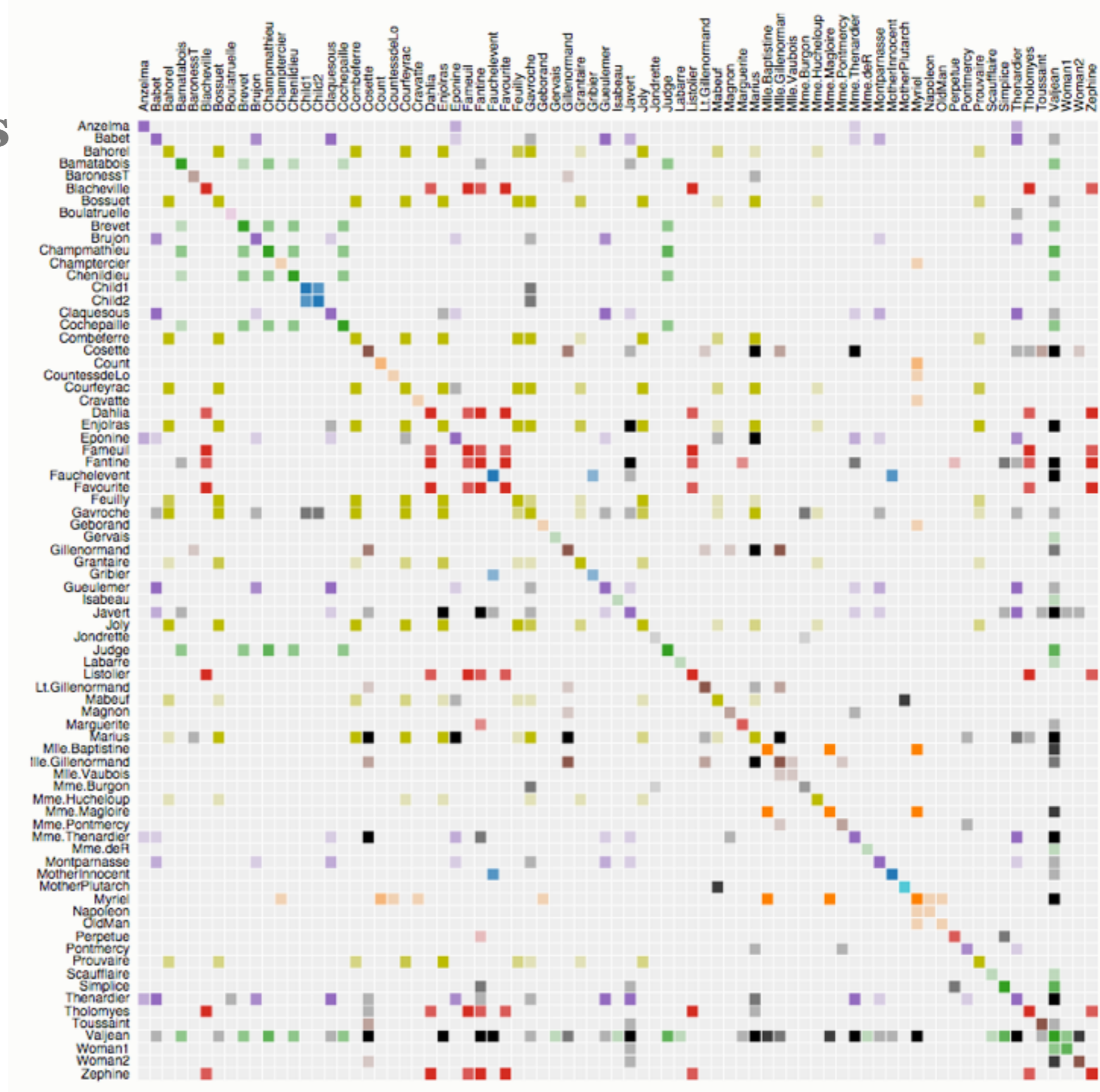  - Motif Glyphs
  - PivotGraph
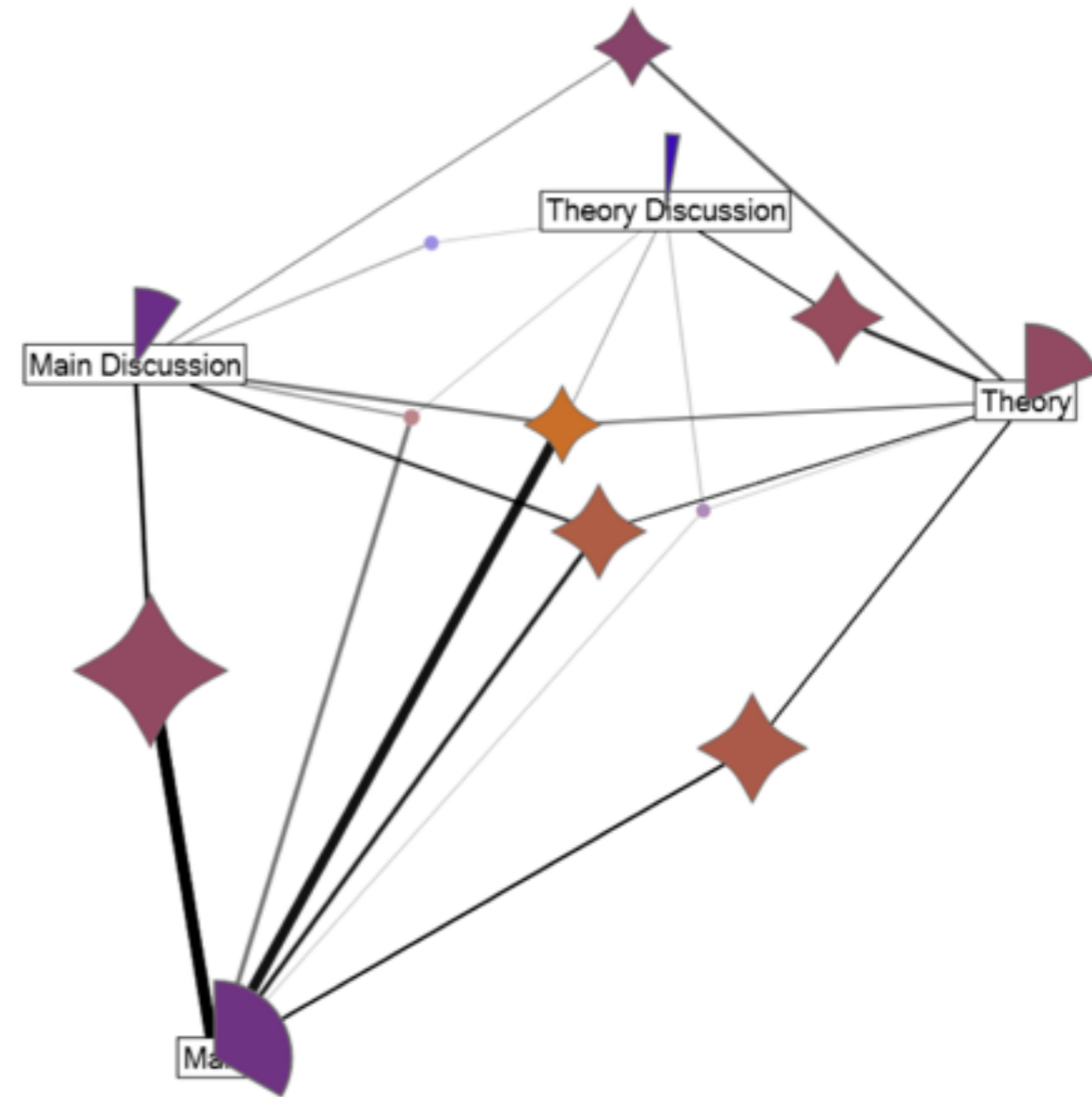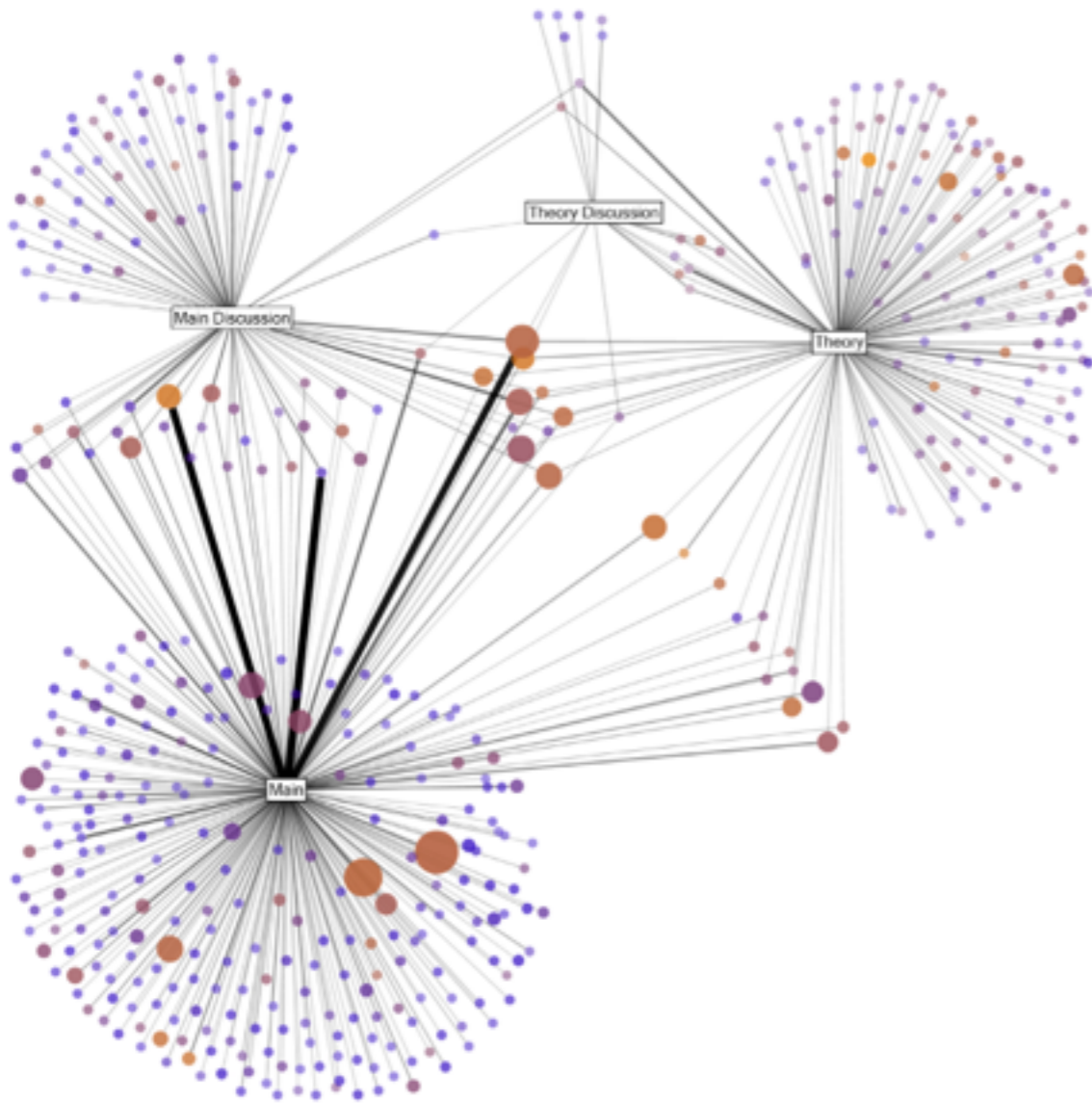
# Les Misérables
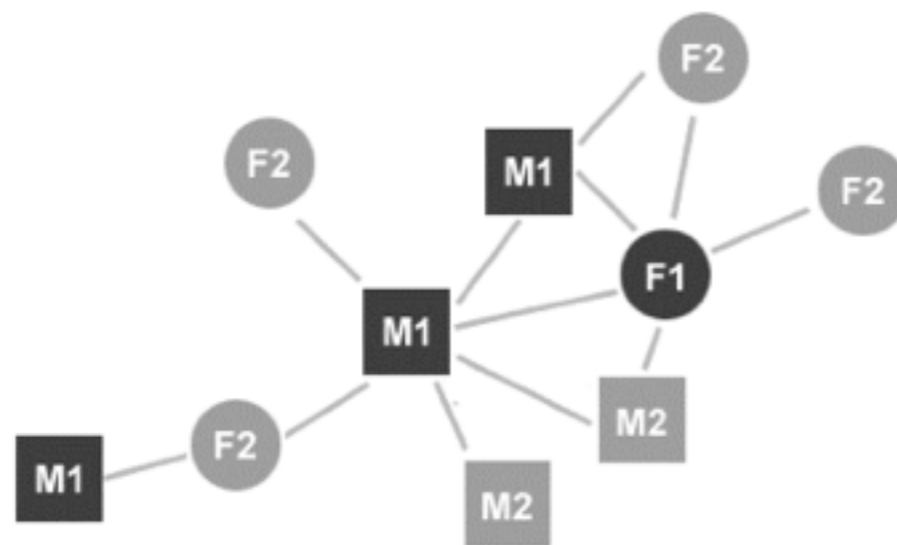character co-occurrence

# Les Misérables
character
co-occurrence
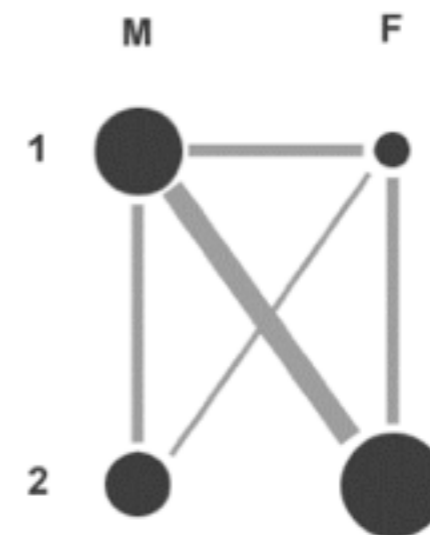
# MOTIF GLYPHS



Dunne 2013

# PIVOT GRAPHS

- new graph, derived from categorical node attributes

- 1D or 2D layouts possible

- size of nodes and edges related to number of aggregated original nodes and edges

- scalability through abstraction, not layout algorithm



Node and Link Diagram

PivotGraph Roll-up

today . . .

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists |
|--------|------------------|--------|----------|------------------------|
| Items | Items (nodes) | Grids | Items | Items |
| Attributes | Links | Positions | Positions | |
| | Attributes | Attributes | | |

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists | Text |
|---|---|---|---|---|---|
| Items | Items (nodes) | Grids | Items | Items | |
| Attributes | Links | Positions | Positions | | |
| | Attributes | Attributes | | | |

# TEXT

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists | Text |
|---|---|---|---|---|---|
| Items | Items (nodes) | Grids | Items | Items | ? |
| Attributes | Links | Positions | Positions | | |
| | Attributes | Attributes | | | |

WHAT DOES IT MEAN TO BE AN "ITEM"?

# text data type

- no numbers   *(implicitly)*

- characters:   ASCII

- strings

**USASCII code chart**

| b4 b3 b2 b1 | Column / Row | 0 NUL | 1 DLE | 2 SP | 3 0 | 4 @ | 5 P | 6 ` | 7 p |
|---|---|---|---|---|---|---|---|---|---|
| 0 0 0 0 | 0 | NUL | DLE | SP | 0 | @ | P | ` | p |
| 0 0 0 1 | 1 | SOH | DC1 | ! | 1 | A | Q | a | q |
| 0 0 1 0 | 2 | STX | DC2 | " | 2 | B | R | b | r |
| 0 0 1 1 | 3 | ETX | DC3 | # | 3 | C | S | c | s |
| 0 1 0 0 | 4 | EOT | DC4 | $ | 4 | D | T | d | t |
| 0 1 0 1 | 5 | ENQ | NAK | % | 5 | E | U | e | u |
| 0 1 1 0 | 6 | ACK | SYN | & | 6 | F | V | f | v |
| 0 1 1 1 | 7 | BEL | ETB | ' | 7 | G | W | g | w |
| 1 0 0 0 | 8 | BS | CAN | ( | 8 | H | X | h | x |
| 1 0 0 1 | 9 | HT | EM | ) | 9 | I | Y | i | y |
| 1 0 1 0 | 10 | LF | SUB | * | : | J | Z | j | z |
| 1 0 1 1 | 11 | VT | ESC | + | ; | K | [ | k | { |
| 1 1 0 0 | 12 | FF | FS | , | < | L | \ | l | \| |
| 1 1 0 1 | 13 | CR | GS | - | = | M | ] | m | } |
| 1 1 1 0 | 14 | SO | RS | . | > | N | ^ | n | ~ |
| 1 1 1 1 | 15 | SI | US | / | ? | O | — | o | DEL |

# text data type

- no numbers  *(implicitly)*

- characters:   ASCII

- strings



USASCII code chart

# text data semantics

love

visualization

I

.
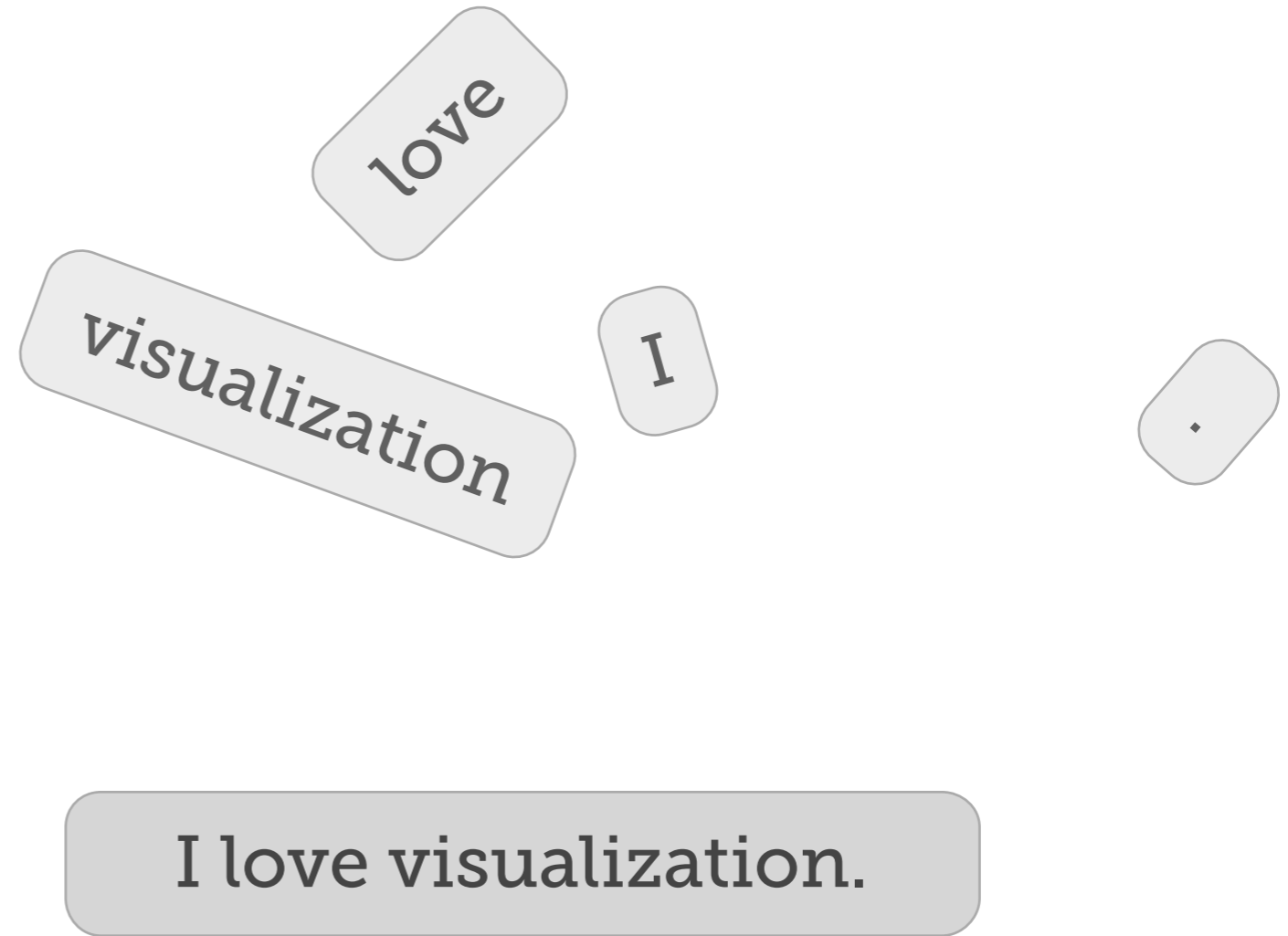
- words


- sentences
  - paragraphs
  - chapters


- lines

# text data semantics

-words
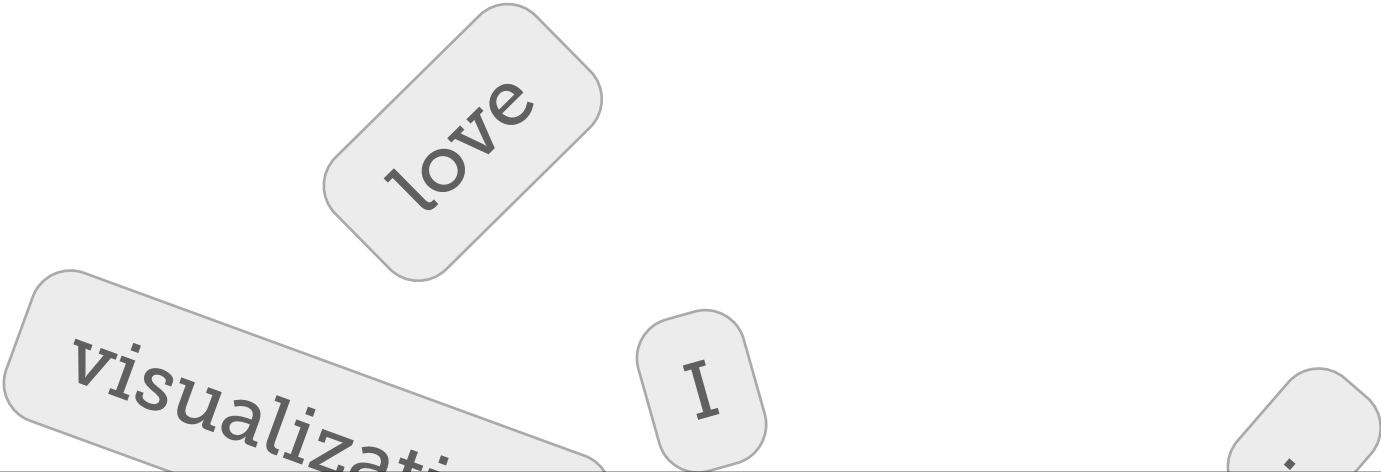
-sentences
  -paragraphs
  -chapters

-lines

love

visualization

I

.

I love visualization.

# text data semantics

- words

- sentences
  - paragraphs
  - chapters

- lines

love

visualizati

I

```
16  // displays a data set using parallel coordinates
17
18  // dataset info
19  String dataSet = "cars";
20  String fileName = dataSet + ".csv";
21  boolean cluster = true;
22  FloatTable table;
23  float[][] data;
24
25  // row, column info
26  String[] colNames;
27  int col = 0;
28  int colTot;
29  String[] rowNames;
30  int row = 0;
31  int rowTot;
```

# text data semantics

- **documents**
  - books
  - papers
  - webpages
  - emails
  - twitter post

- **corpus:** *collection of documents*

# text data semantics

- **documents**
  - books
  - papers
  - webpages
  - emails
  - twitter post

- **corpus:** *collection of documents*

single document

# Tag Clouds / Word Clouds



http://www.tagcrowd.com

http://www.wordle.com

# Visualizations : definitions of visualization word tree

# Text Arc



Wattenberg, Viegas 2008

# DocuBurst

# Arc Diagrams



Analysis of the Characters from Les Misérables: http://mbostock.github.io/protovis/ex/arc.html

# Rule-Based: Poetry

collection of documents

Parallel Tag Clouds to Explore and Analyze Faceted Text Corpora

Christopher Collins

# Document Cards
## (small multiples)

# Showing Temporal Relationships: ThemeRiver (Stream Graph)



Havre, Hetzler, Nowell 2000

# Jigsaw: Many Linked Views

**Visual Analytics Support for Intelligence Analysis**
**Case Study: The 9/11 Report**

Carsten Görg
Youn-ah Kang
Zhicheng Liu
John Stasko

Information Interfaces Group
Georgia Institute of Technology

# Jigsaw: Many Linked Views

# SETS

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists | Text |
|--------|------------------|--------|----------|----------------------|------|
| Items | Items (nodes) | Grids | Items | Items | ? |
| Attributes | Links | Positions | Positions | | |
| | Attributes | Attributes | | | |

36

# thought experiment...

- **item:**   Lego

- **attributes**

# thought experiment...

- **item:**   Lego

- **attributes**
  - color
  - height
  - width
  - length
  - shape

# dataset: option 1

# dataset: option 2

# dataset: more realistic

# dataset

·**where do we start?** we need to organize! but, how?

# dataset

- sort by color

# dataset

-sort by size, shape

# dataset

-**task:** organization

-drawbacks?

# dataset

- organization leads us to a set problem

- so what are sets?

# set theory

- **set**
  - a collection of objects
  - some set: A

- **object**
  - some object: z
  - z ∈ A

A

z

# set theory

- **set**
  - a collection of objects
  - some set: A

A

# set theory

-multiple sets:   A & B

B

A

# set theory

-**union:**   A ∪ B

# set theory

-intersection:   A ∩ B



B    A

# set theory

-set difference:  A \ B



B          A
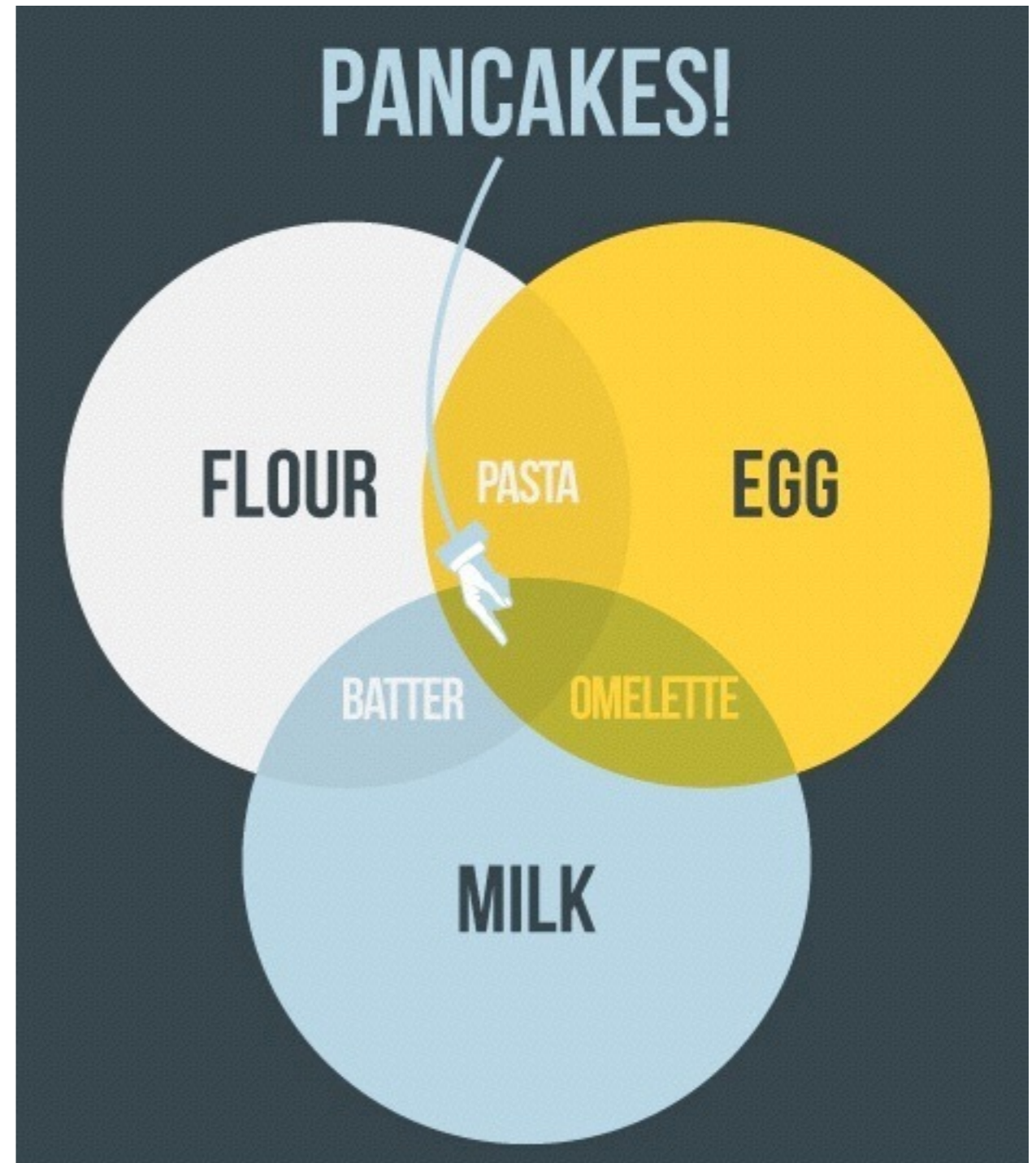
# set theory

- symmetric difference:  $A \ominus B$

# visualizing sets
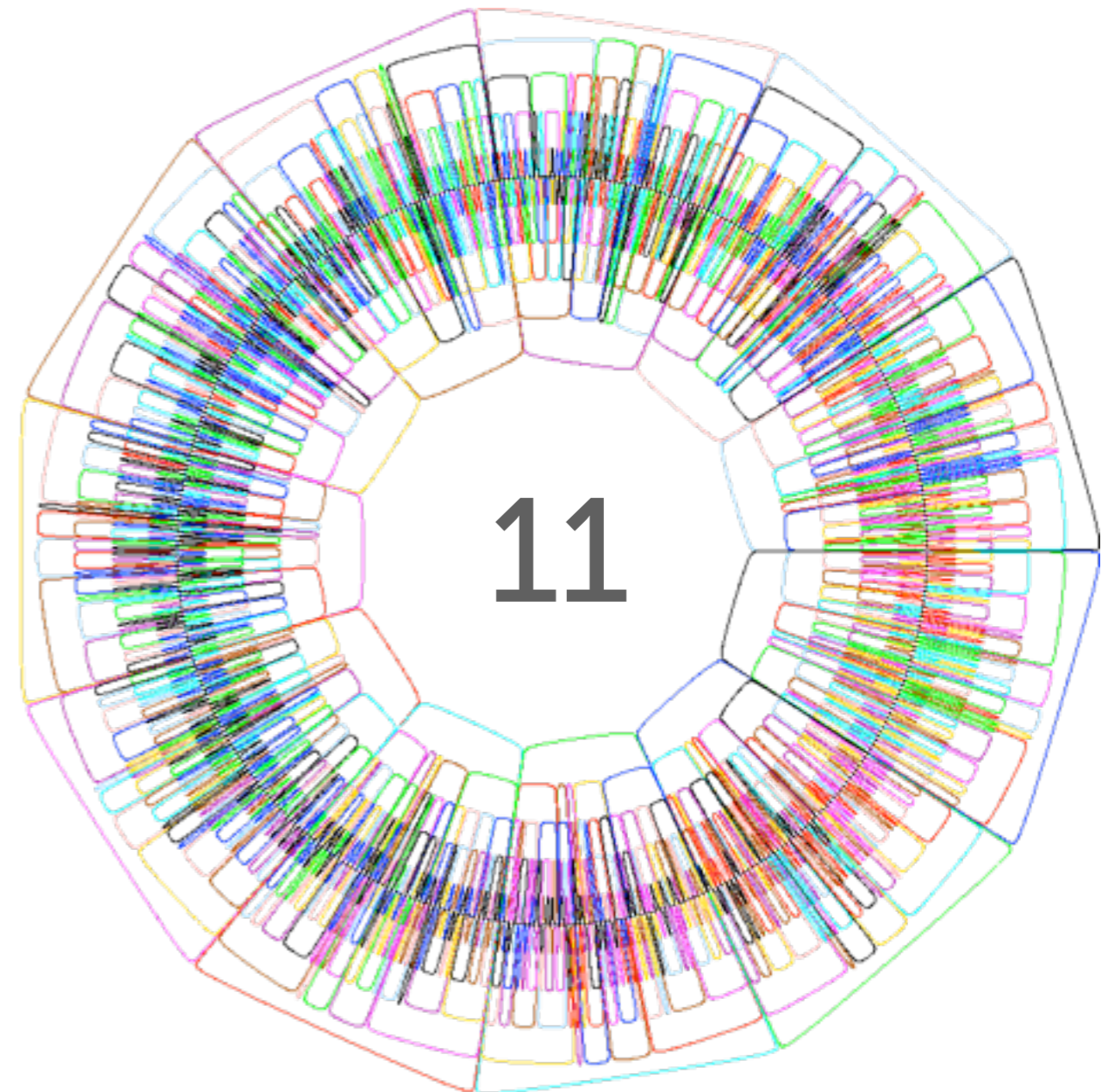
# venn diagrams

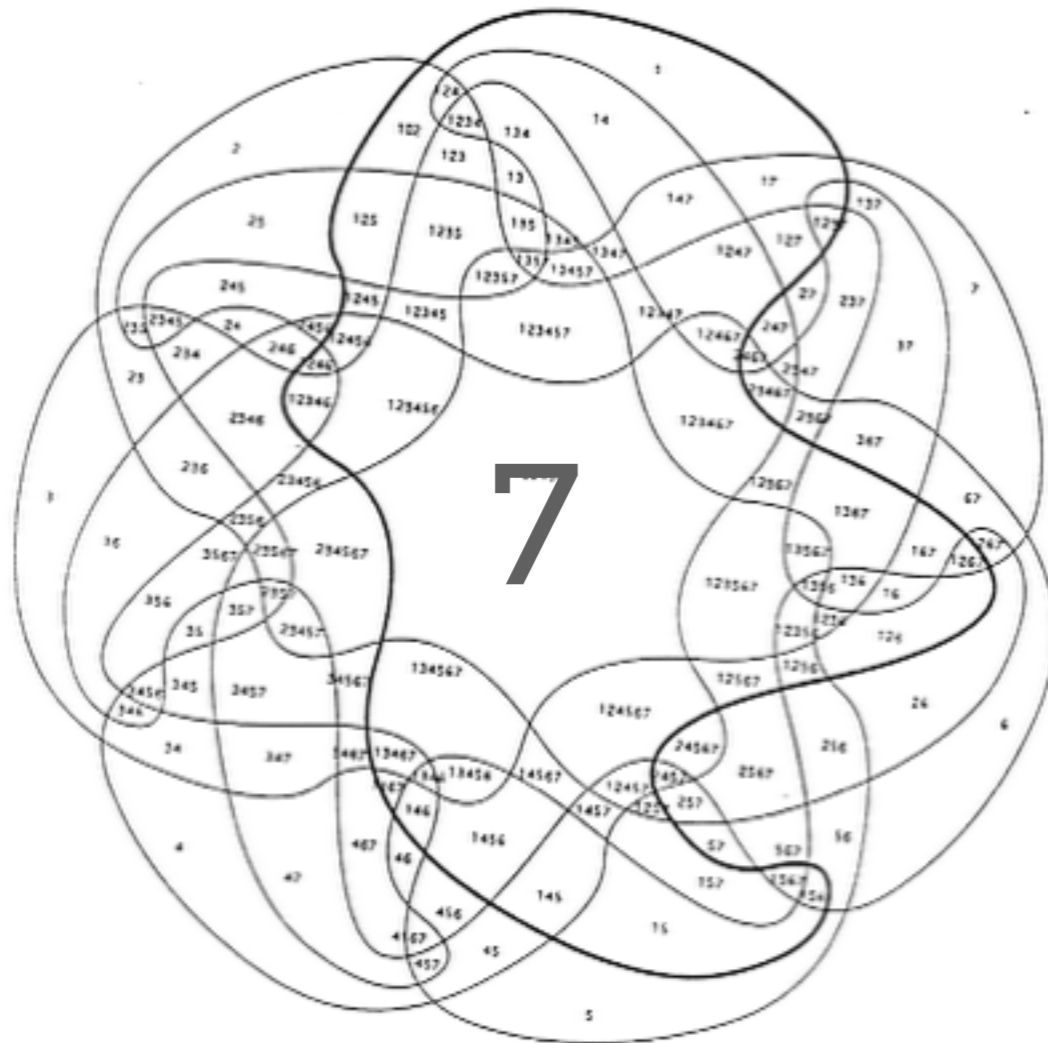-show all possible relationships
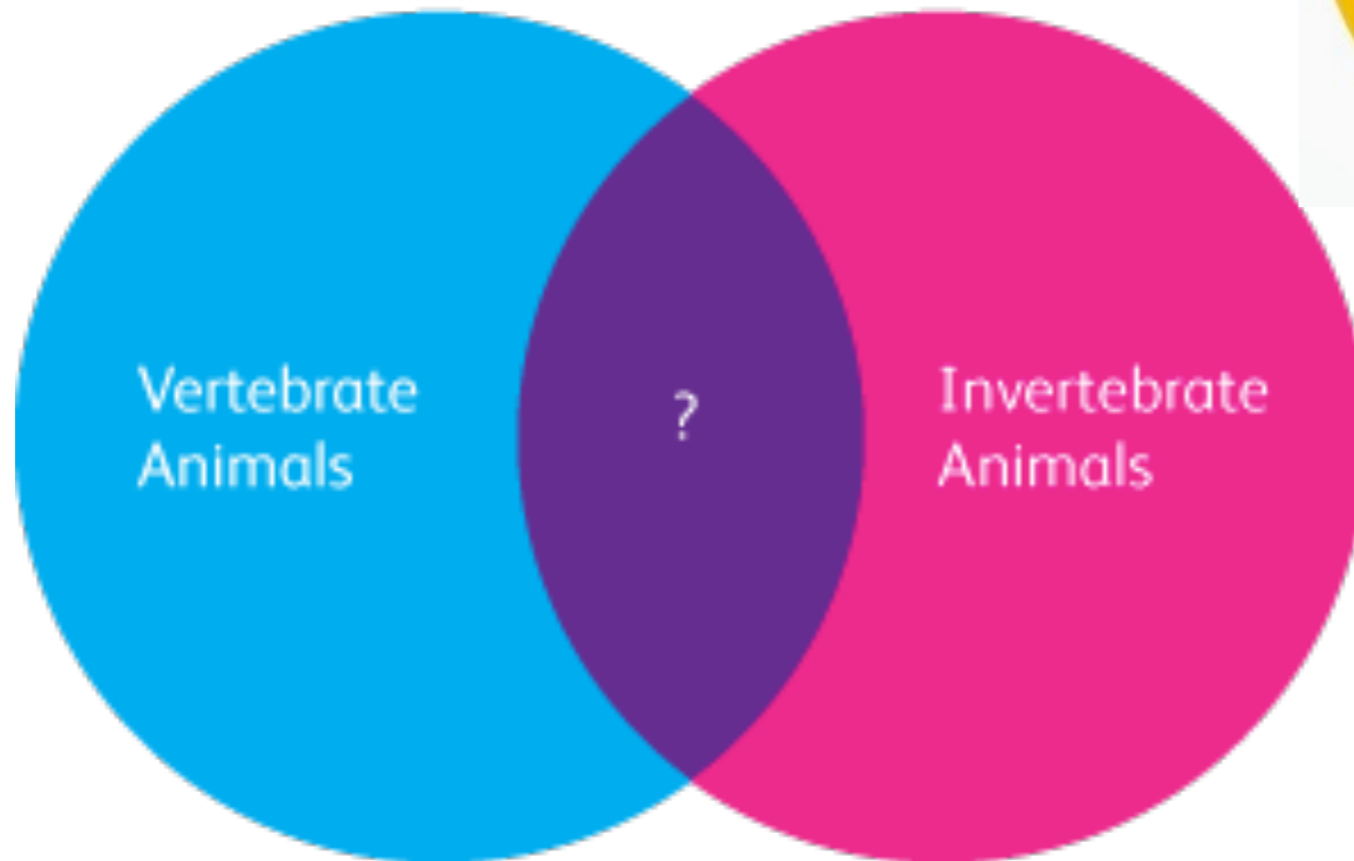
# venn diagrams

-casual infovis
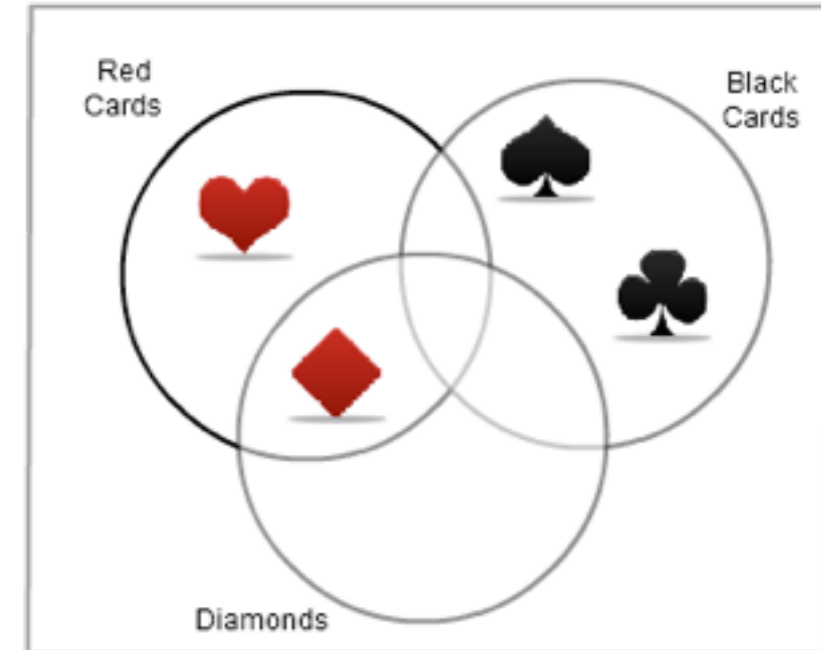
# venn diagrams

-get messy fast

# venn diagrams

-non-sensical

# euler diagrams

-show only existing relationships

# euler diagrams

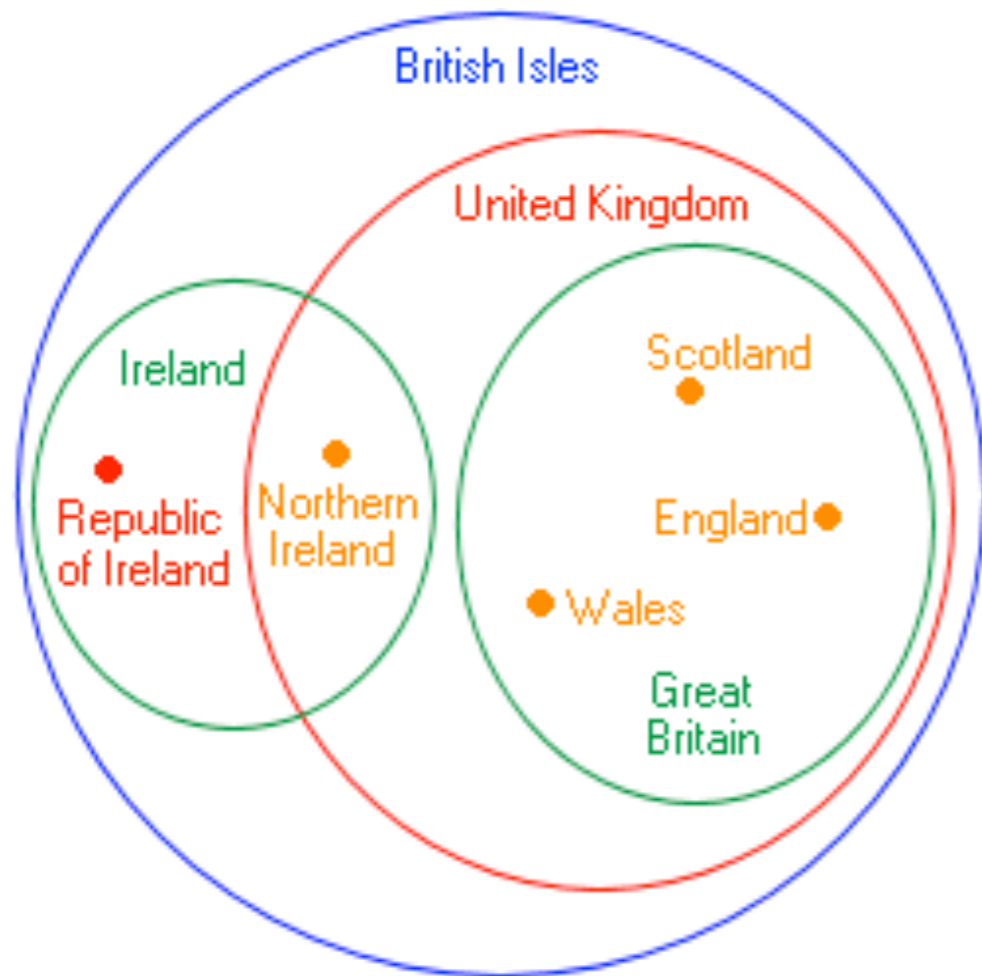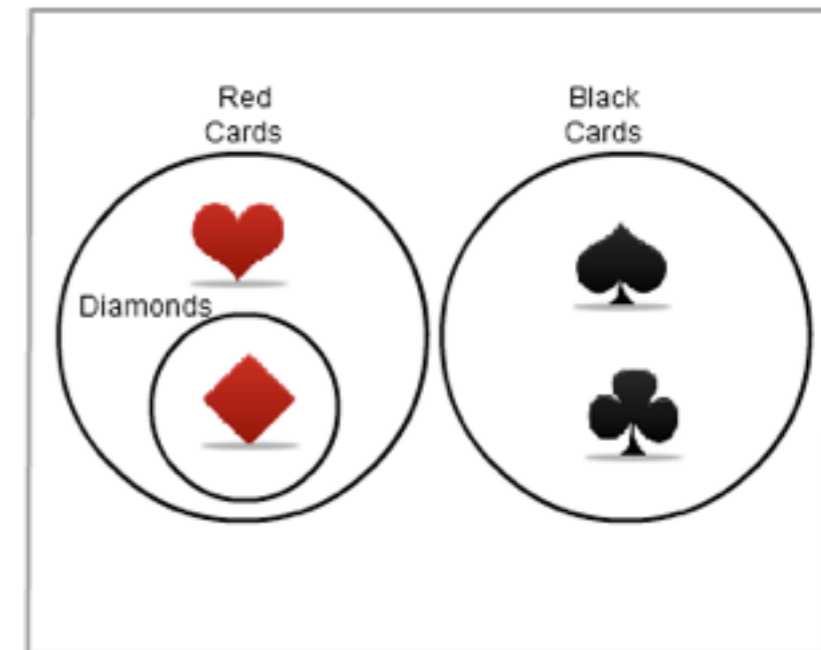-show only existing relationships

# euler diagrams

-misunderstood



Maybe this Venn Diagram will explain this better :

What you just said

Things I care about

# euler diagrams

1: People who know what a Venn Diagram is.
2: People who know what an Euler Diagram is.
3: People who know the difference.

# venn & euler diagrams

-adjust for area

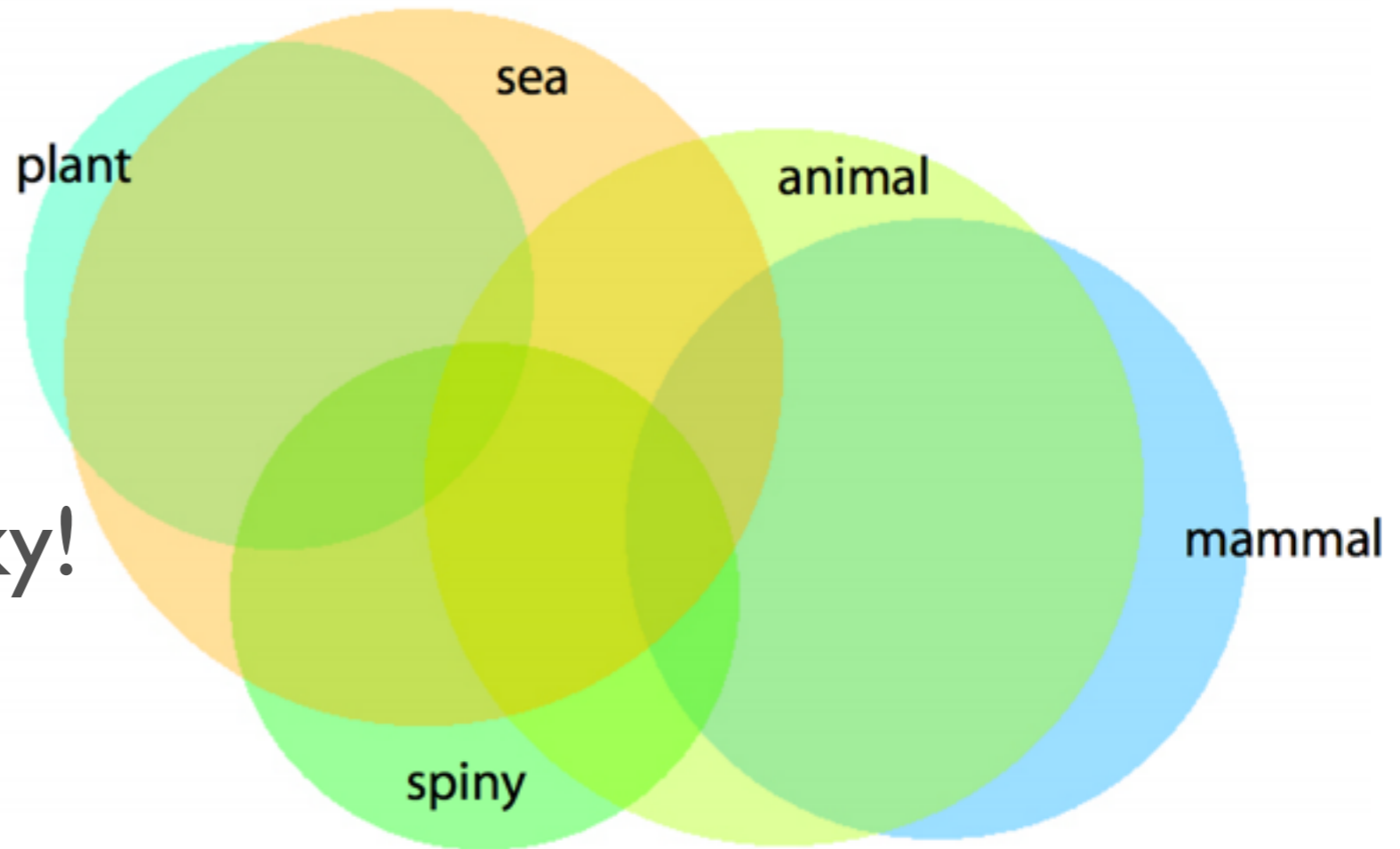-starts getting tricky!
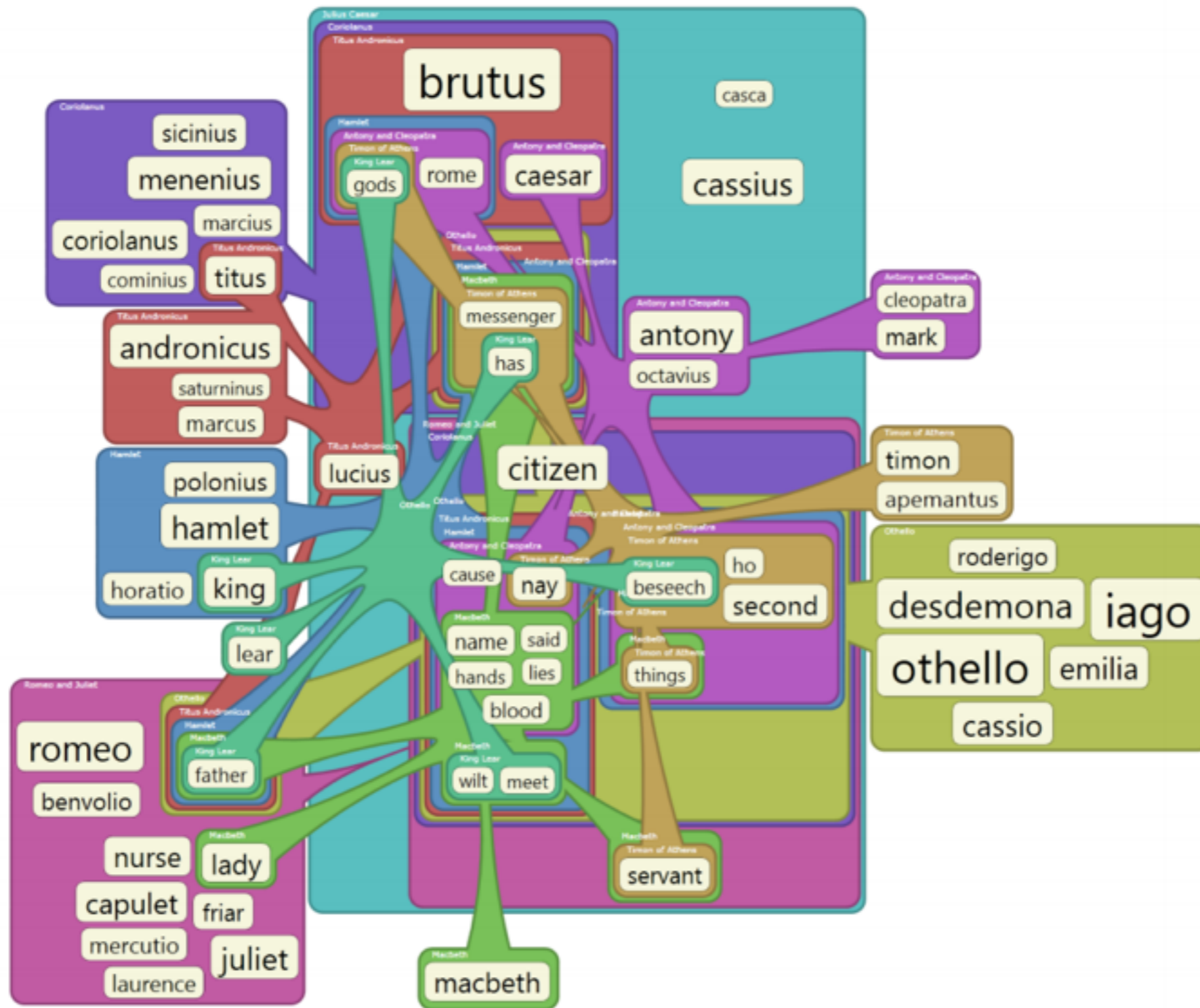
# venn & euler diagrams

-adjust for area

-starts getting tricky!

# compact euler diagrams

# parallel sets



Titanic Survivors

# set o'gram



Titanic

# visualizing sets with constraints

# bubble sets

- connect points

# line sets

- restaurants                     social communities



Animation
Set Visualizations and Clutering
High-Dimentional Data Visualizations
Matrix Visualizations
Text Visualizations
Social Visualizations

# kelp diagrams



- cities on a map

- metabolic network

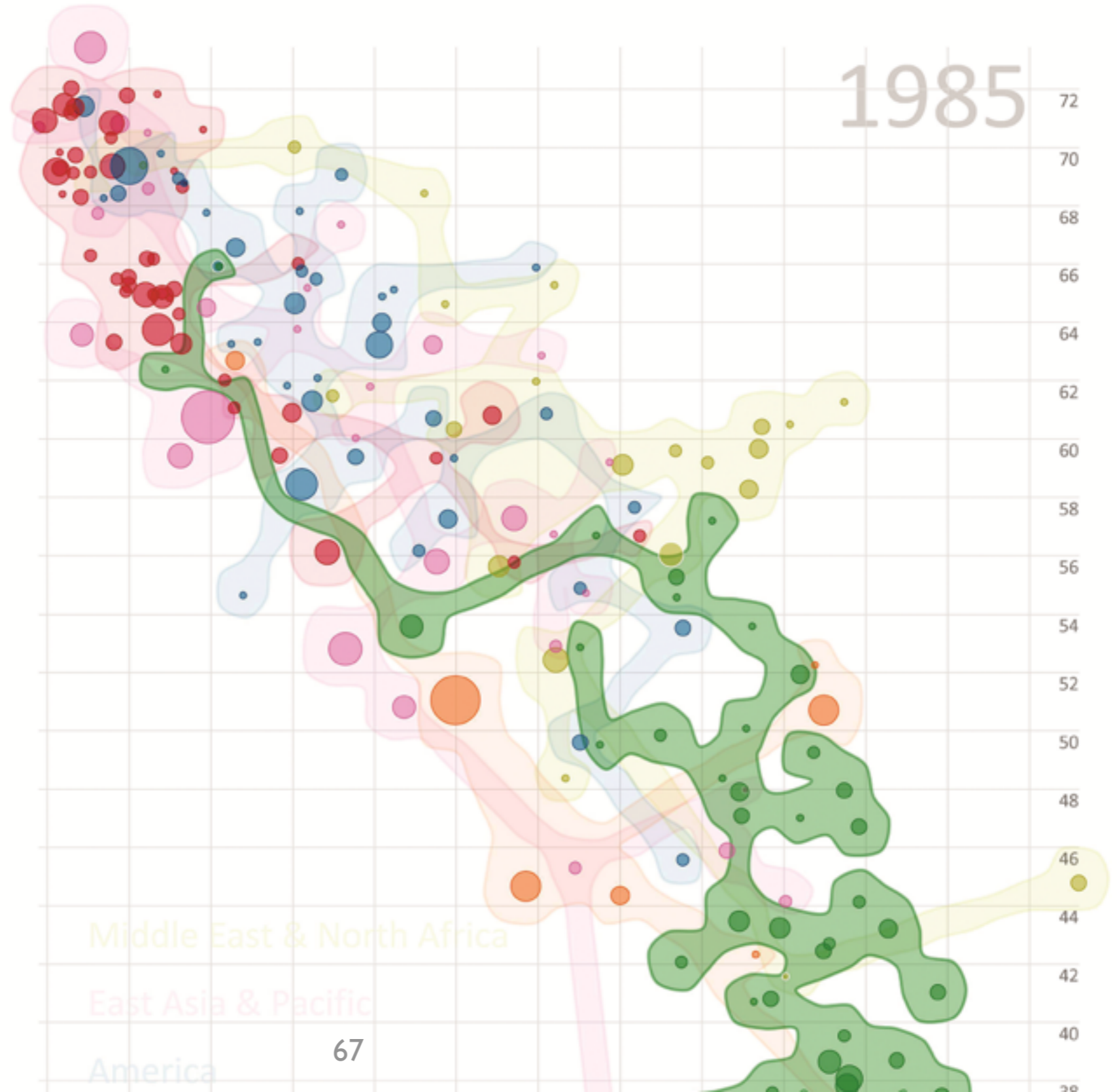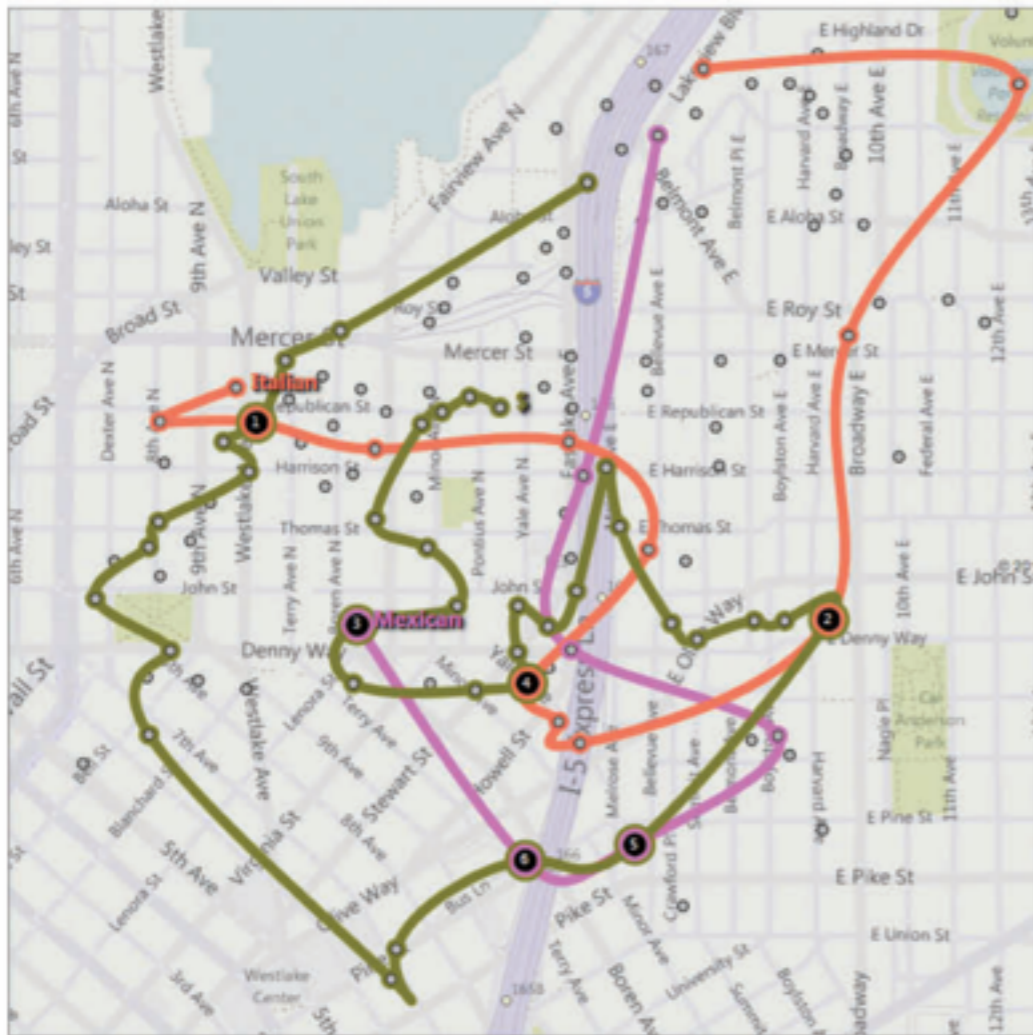# kelp fusion

- cities on map

- lines & areas

# sets

- applies to many datasets

- many combinations may be interesting

- limited numbers of sets

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists | Text |
|---|---|---|---|---|---|
| Items | Items (nodes) | Grids | Items | Items | ? |
| Attributes | Links | Positions | Positions | | |
| | Attributes | Attributes | | | |

L15: Maps
# REQUIRED READING

# Chapter 8

## Arrange Spatial Data

### 8.1 The Big Picture

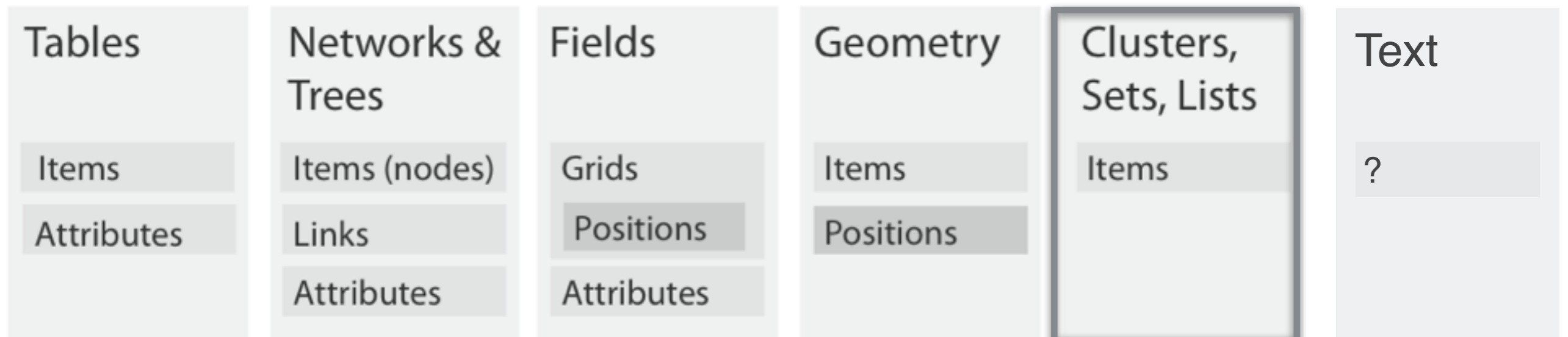For datasets with spatial semantics, the usual choice for *arrange* is to *use* the given spatial information to guide the layout. In this case, the choices of *express*, *separate*, *order*, and *align* do not apply because the position channel is not available for directly encoding attributes. The two main spatial data types are geometry, where shape information is directly conveyed by spatial elements that do not necessarily have associated attributes, and spatial fields, where attributes are associated with each cell in the field. (See Figure 8.1.) For scalar fields with one attribute at each field cell, the two main visual encoding idiom families are isocontours and direct volume rendering. For both vector and tensor fields, with multiple attributes at each cell, there are four families of encoding idioms: flow glyphs that show local information, geometric approaches that compute derived geometry from a sparse set of seed points, texture approaches that use a dense set of seeds, and feature approaches where data is derived with global computations using information from the entire spatial field.

### 8.2 Why Use Given?