

# Consistent Layout for Thematic Software Maps

Adrian Kuhn, Peter Loretan, **Oscar Nierstrasz**

SCG, University of Bern

Published and presented at WCRE 2008.





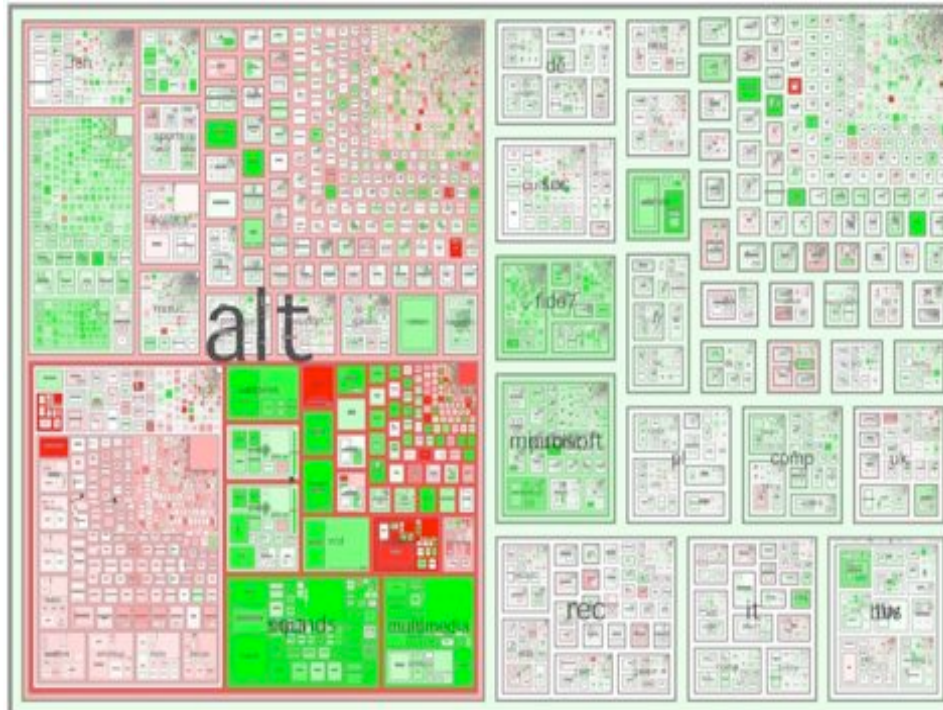
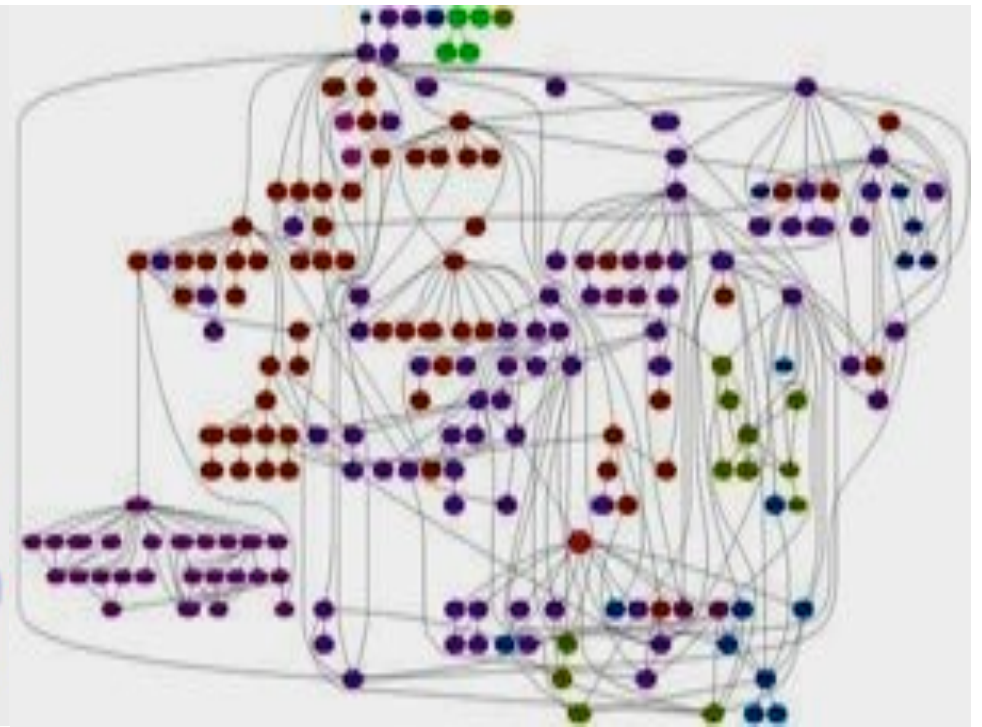
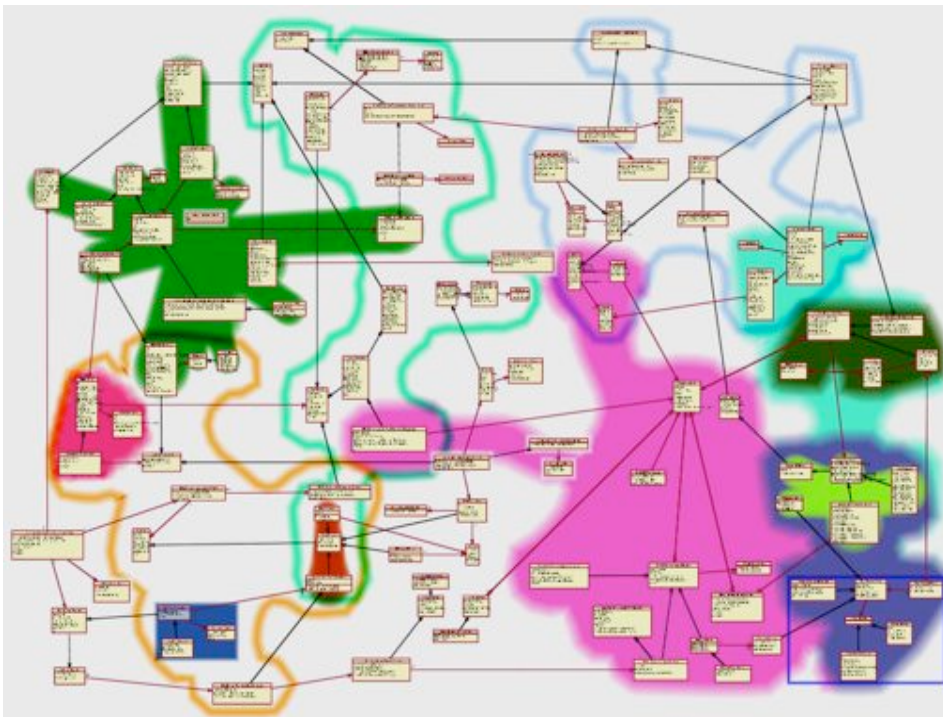




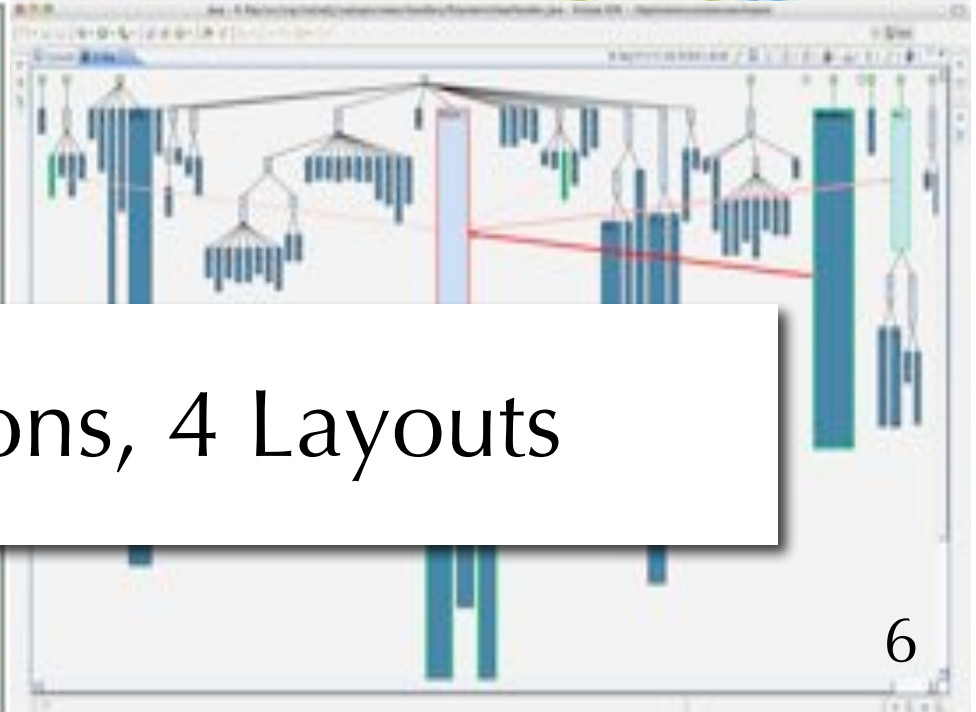
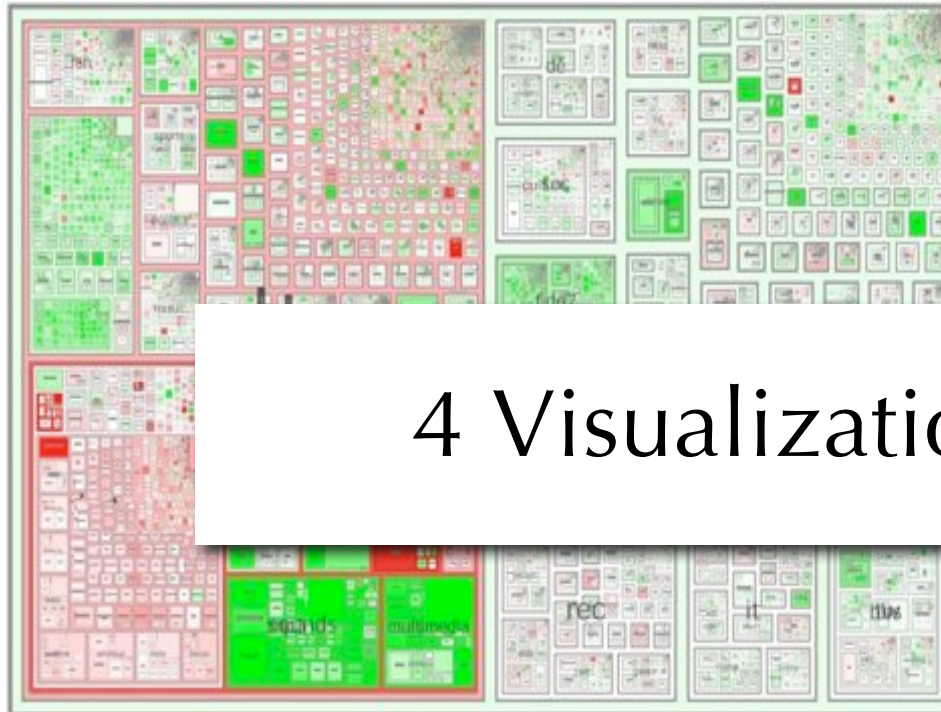
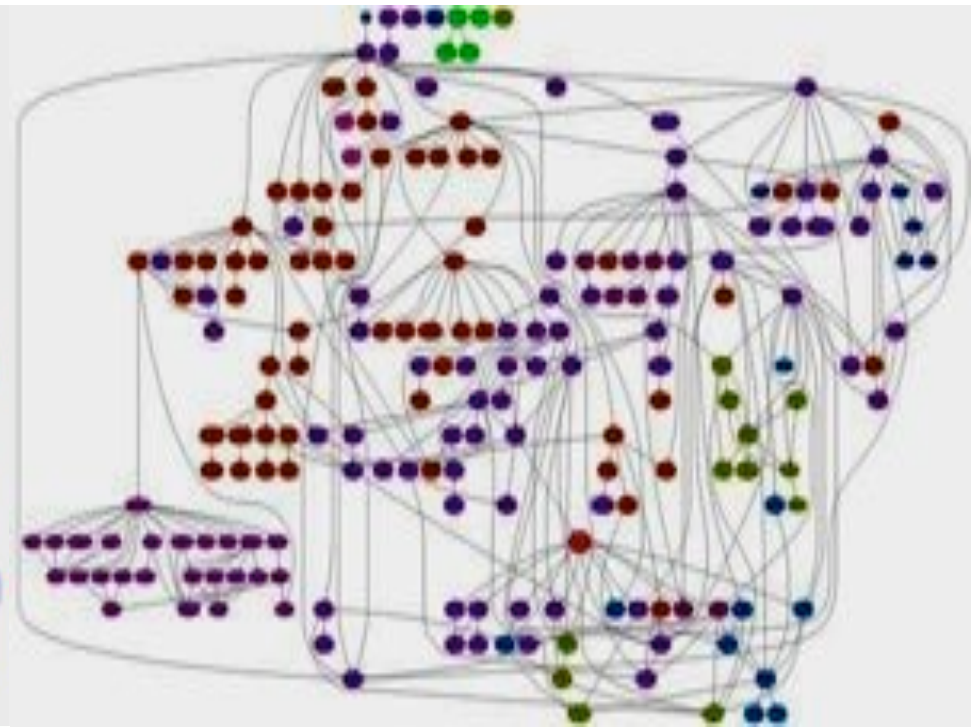
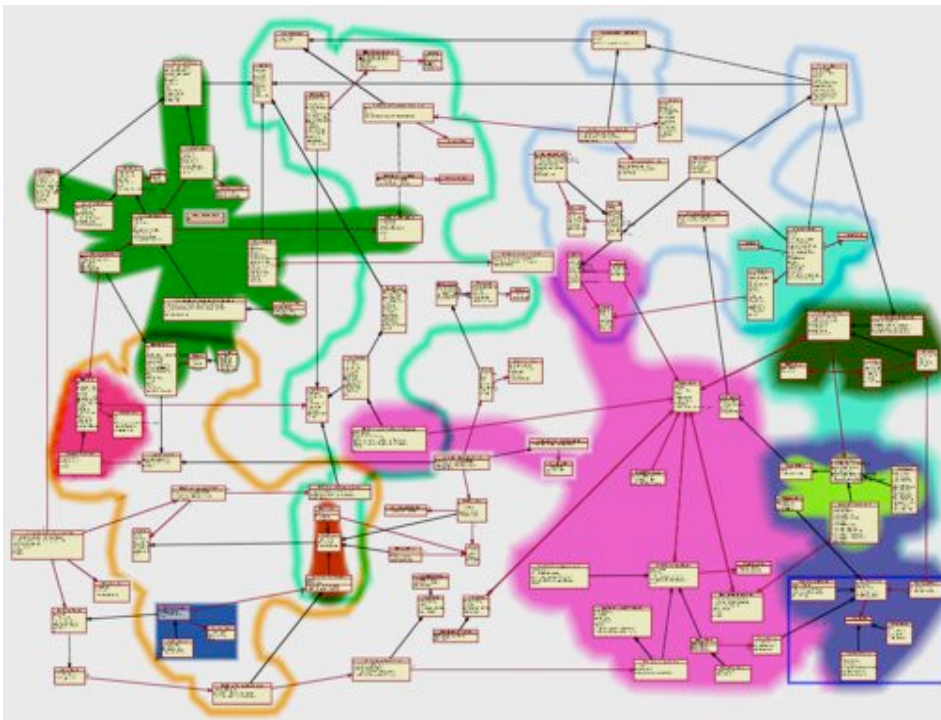


Distance has a meaning



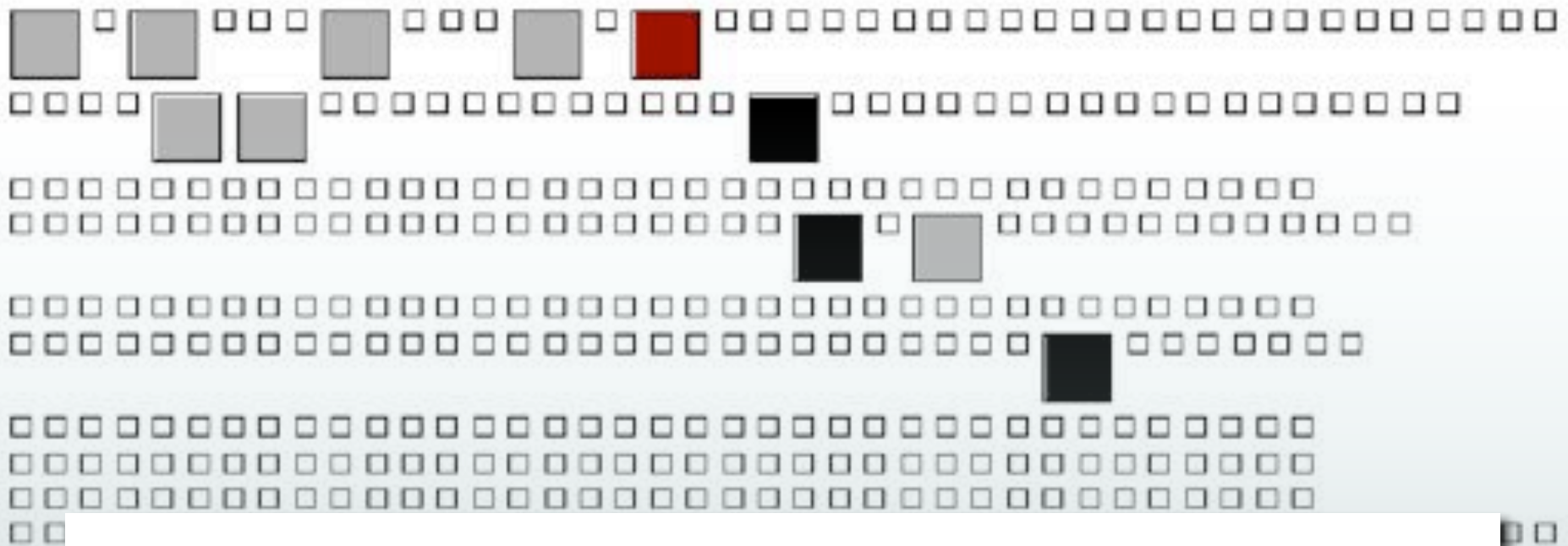






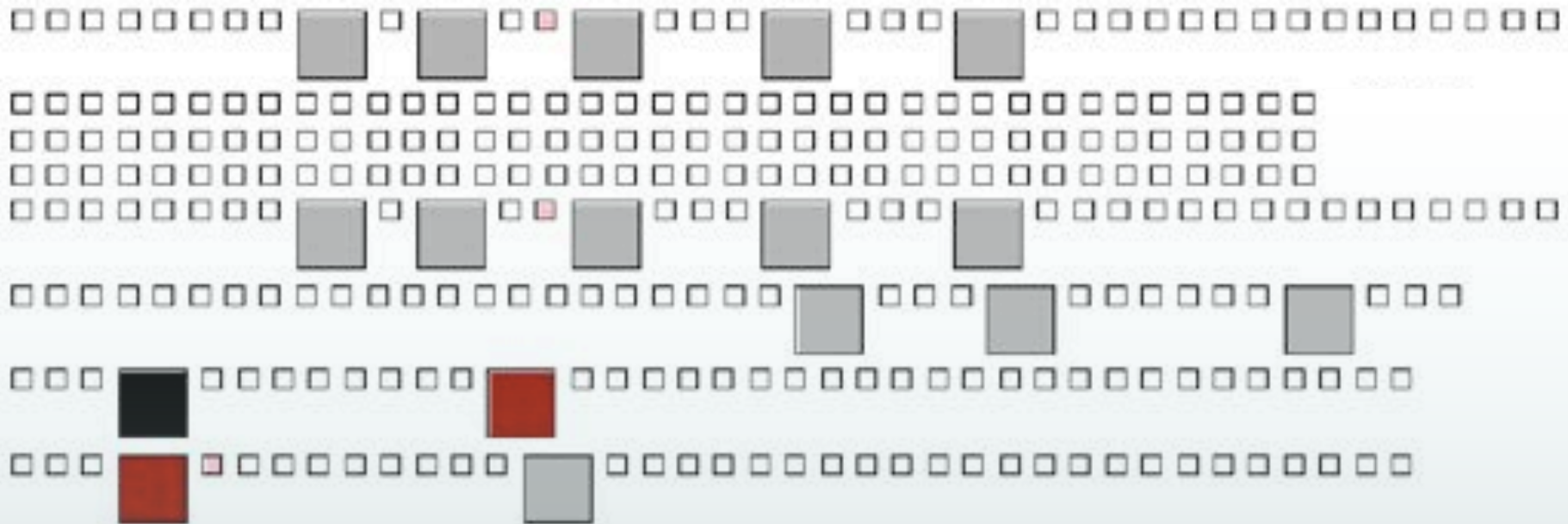
4 Visualizations, 4 Layouts

A ← F → Z



Often alphabetic order is used,  
where distance has little meaning.

0x017 ↔ 0x02A ↔ 0x7FF



Or worse "hash key" order, where distance has no meaning at all.



0x017 ↔ 0x02A ↔ 0x7FF

```
for (each : Set all) {  
    // non-deterministic order!  
}
```

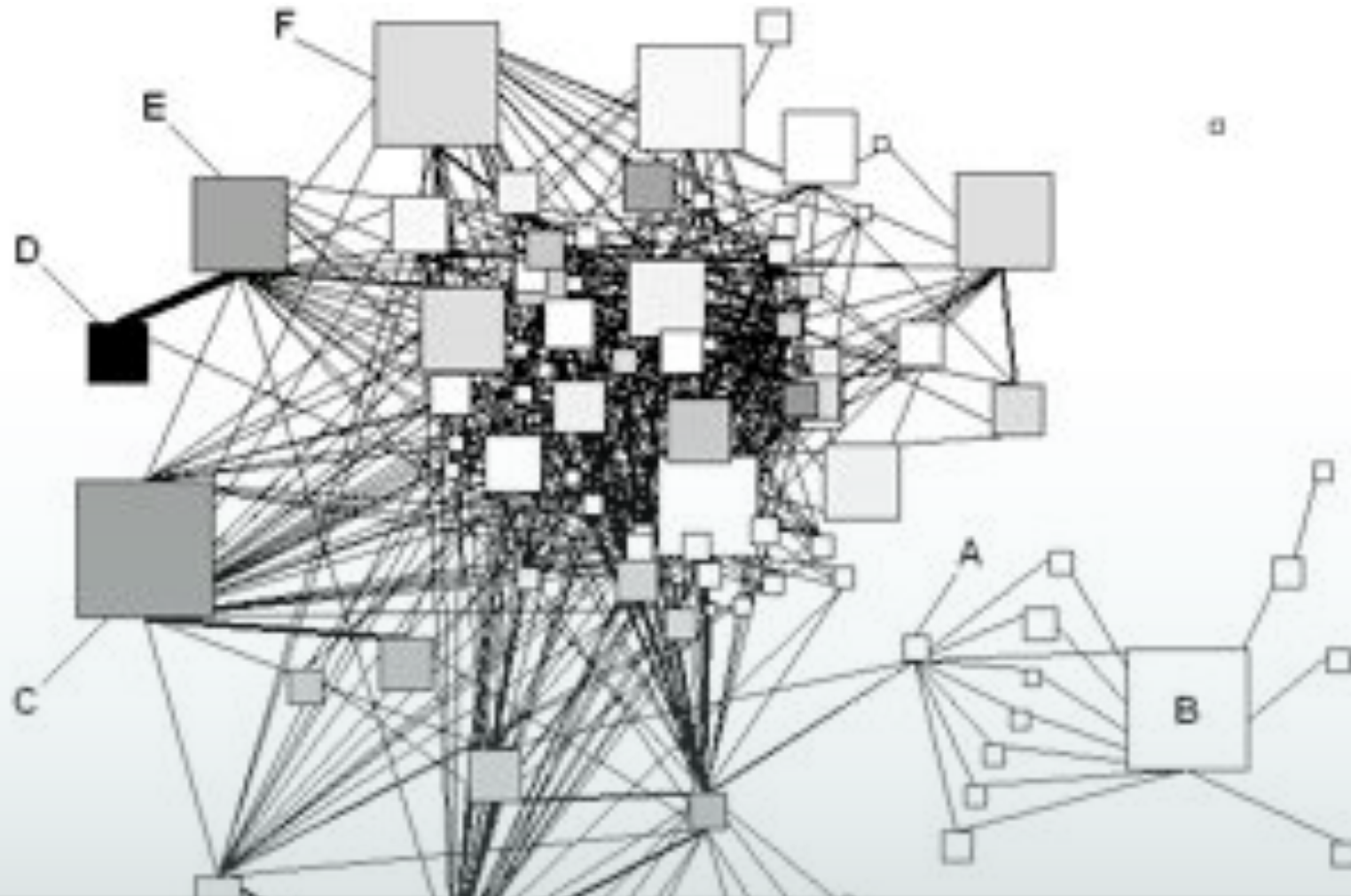
Why? Set is internally realized as a dictionary. Iteration order is given by hash keys, which are arbitrary and thus meaningless.

Or worse "hash key" order, where distance has no meaning at all.

We need: a consistent and meaningful  
layout for software.

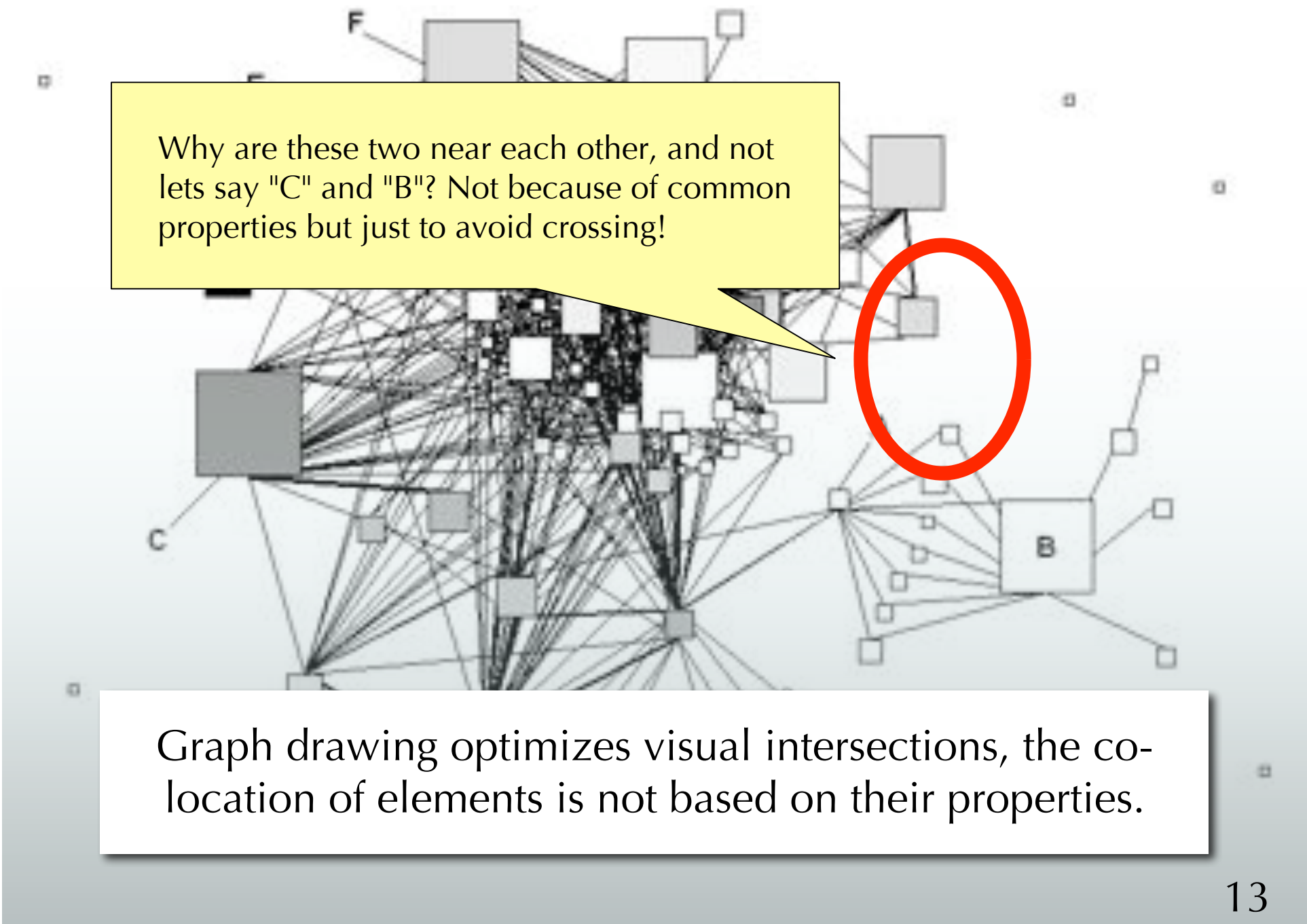


Has it been done?



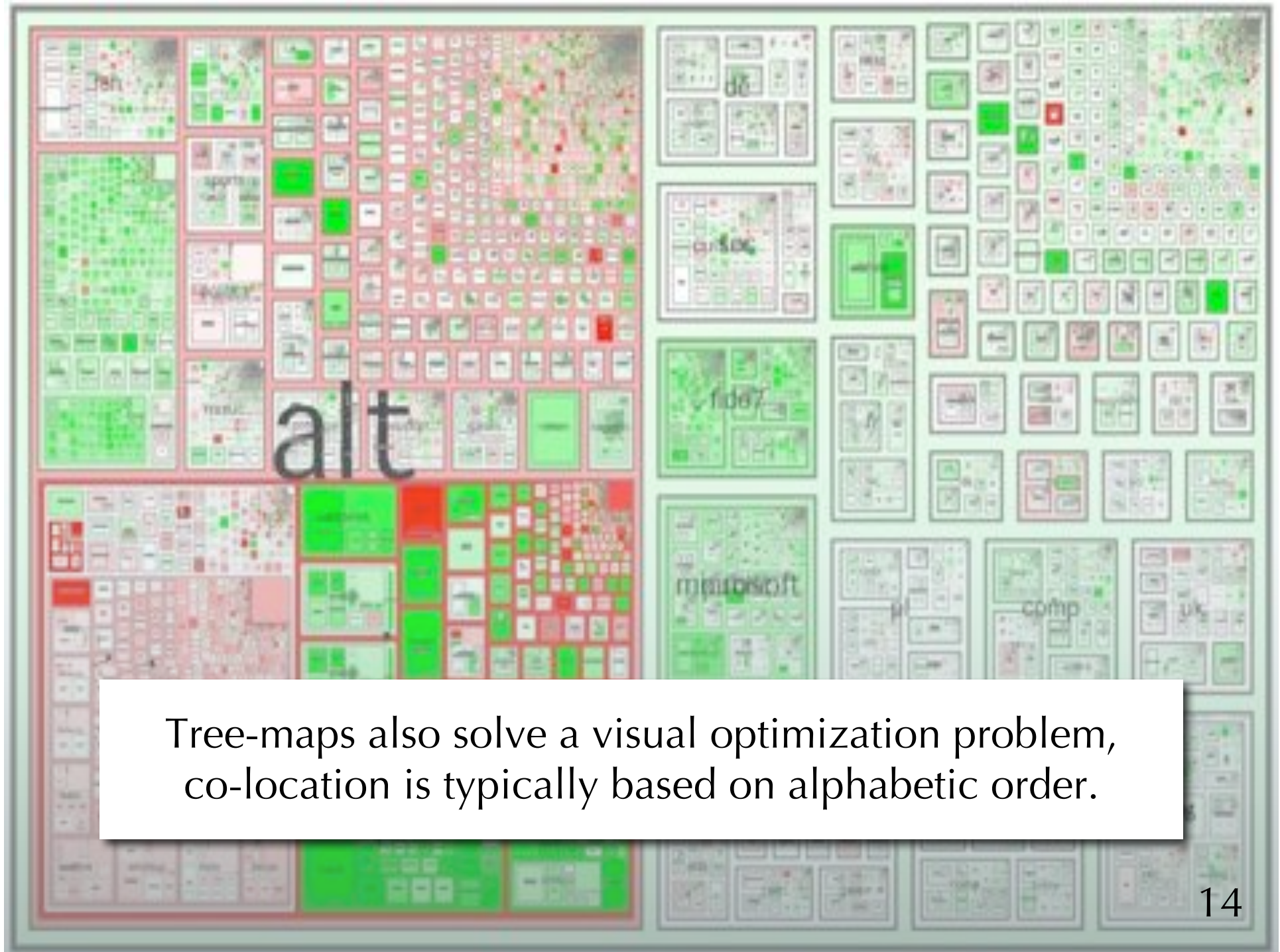
Graph drawing optimizes visual intersections, the co-location of elements is not based on their properties.





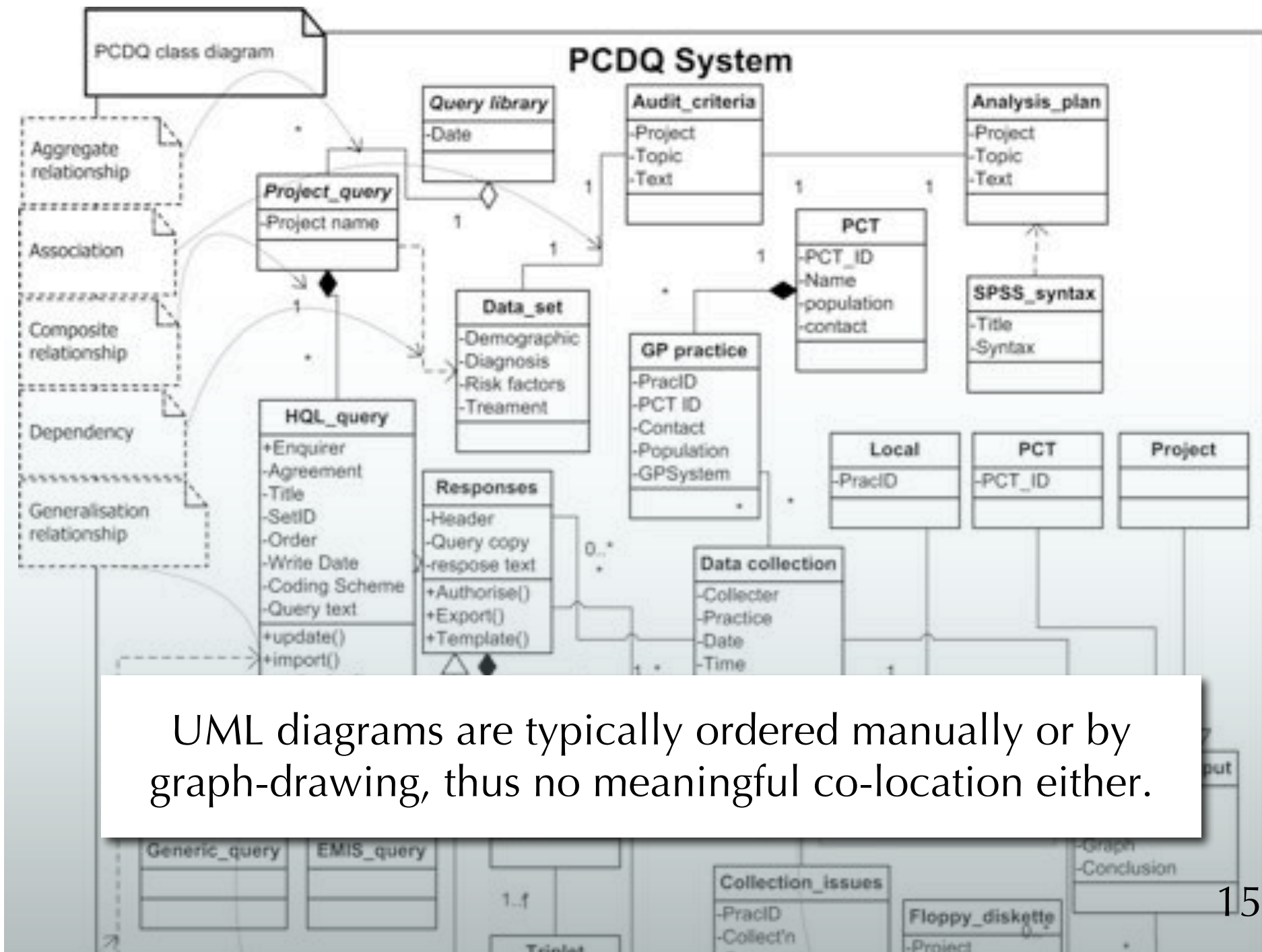
Why are these two near each other, and not lets say "C" and "B"? Not because of common properties but just to avoid crossing!

Graph drawing optimizes visual intersections, the co-location of elements is not based on their properties.



Tree-maps also solve a visual optimization problem, co-location is typically based on alphabetic order.



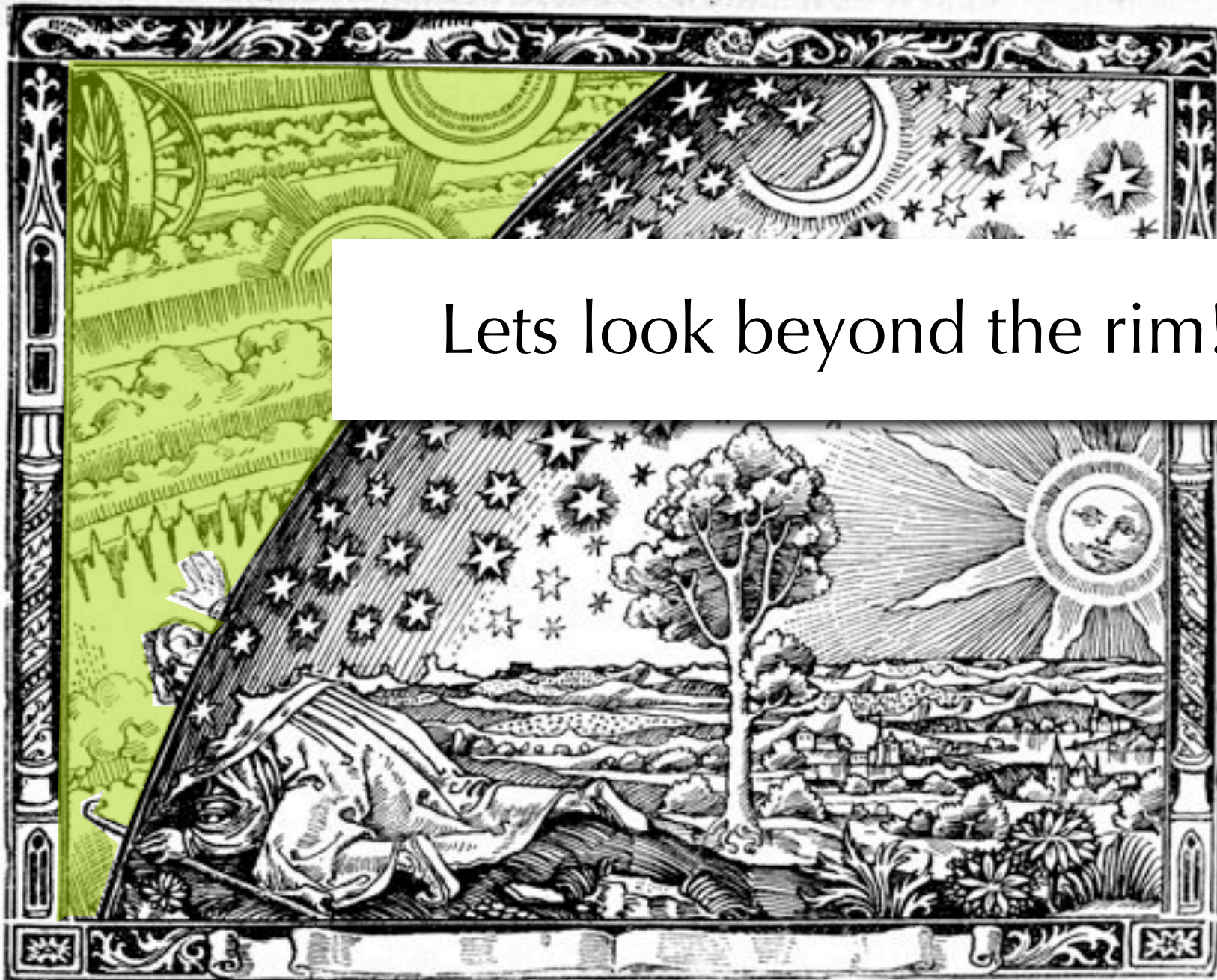


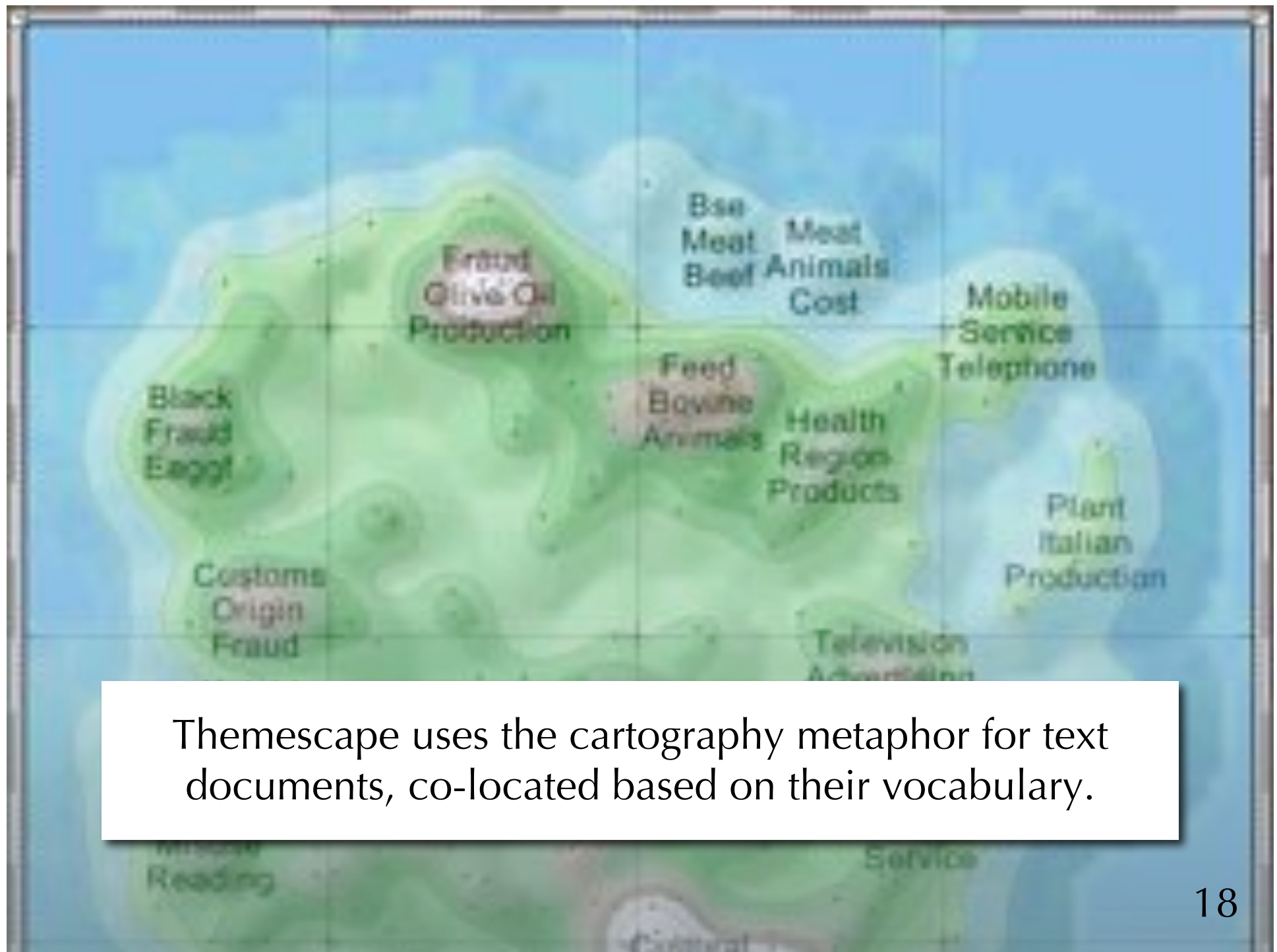


DistributionMap was our first attempt to achieve meaningful co-location, but it's limited to boxology. [Ducasse 2006]



Lets look beyond the rim!

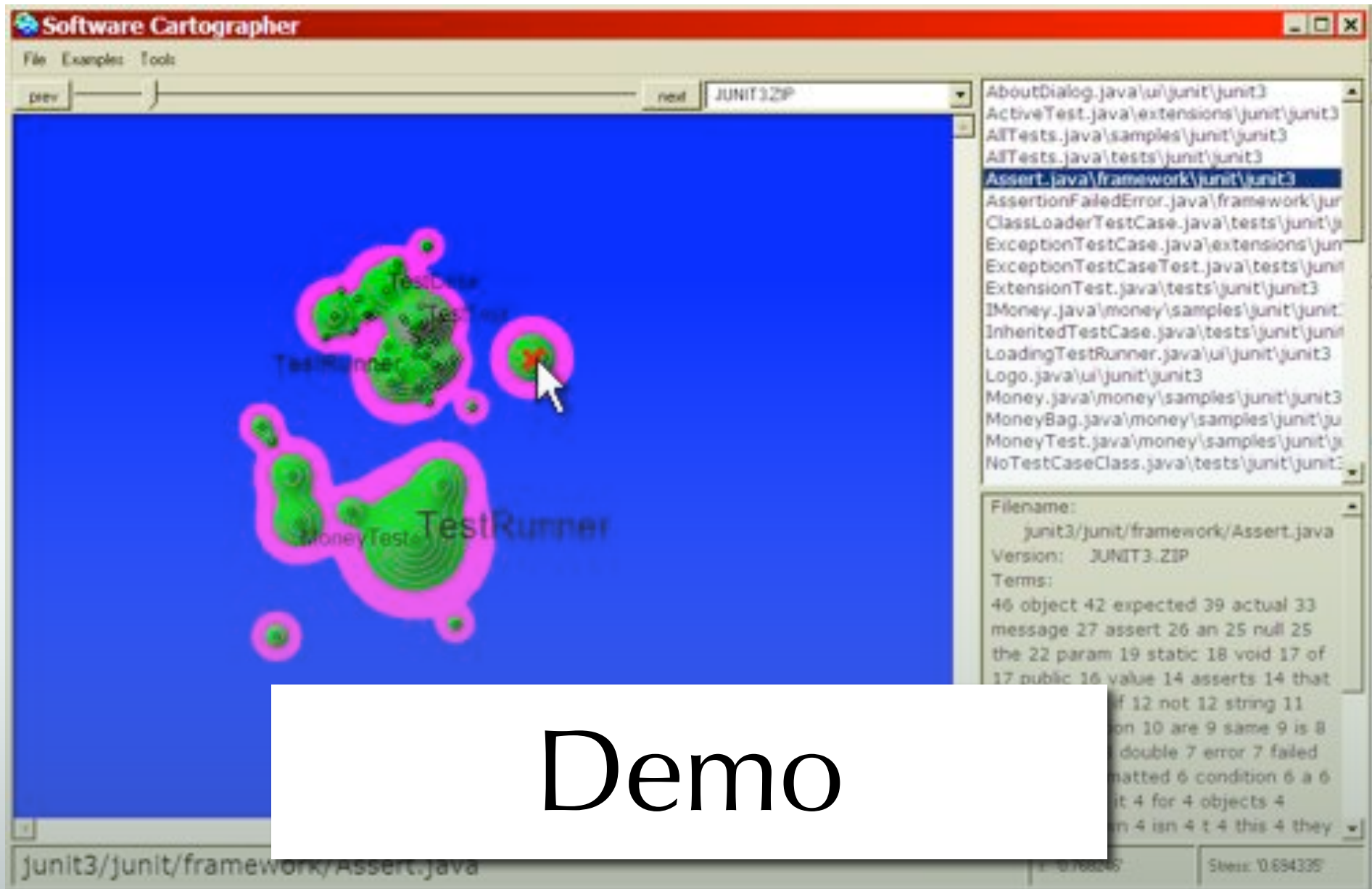




Themescape uses the cartography metaphor for text documents, co-located based on their vocabulary.

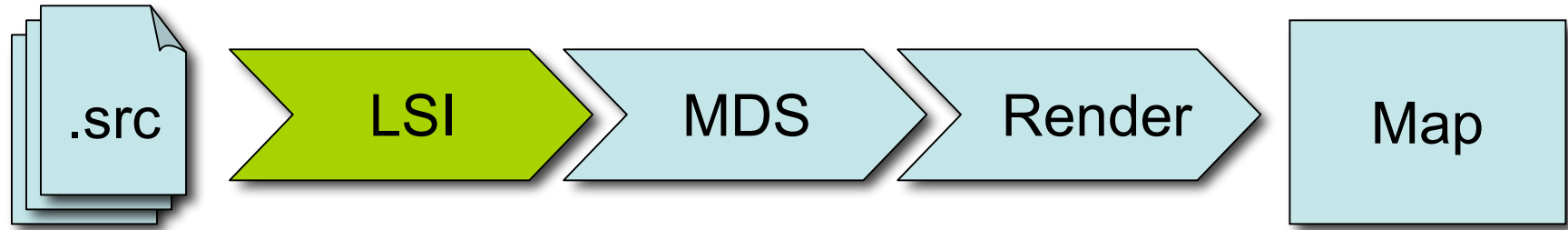
Can this be done for code?





# Why Vocabulary?

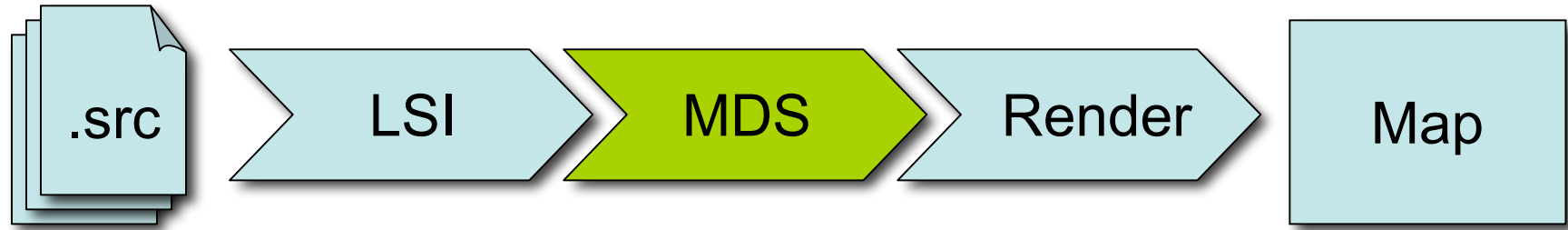
- Abstracts away from structure [Maletic 2001]
- Useful for software analysis [Marcus 2004]
- Metric distance between code [Poshyvanyk 2006]
- Clusters software by topic [Kuhn 2007]
- Over time, software grows... [Vasa 2007]
- ...but vocabulary remains stable [Antoniol 2007]



## Semantic Clustering [Kuhn 2007]

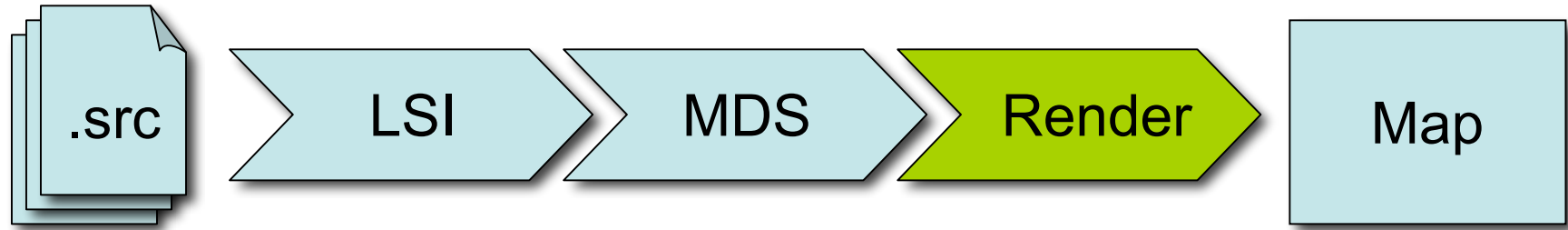
- Analyses the **vocabulary** of source files.
- Latent Semantic Indexing (LSI) [Deerwester 1990]
- Maps files into n-dim vector space.
- Files with similar vocabulary are close.





# Multidimensional Scaling

- Maps vectors from  $n$ -dim to 2D.
- MDS attempts to preserve the relative distance between vectors.
- Iterative, non-deterministic approximation.



# Map Rendering

- Each file is rendered as a hill.
- Hill height equals file size.
- Standard algorithms from cartography:
  - Digital elevation model (DEM)
  - Hill Shading
  - Contour lines

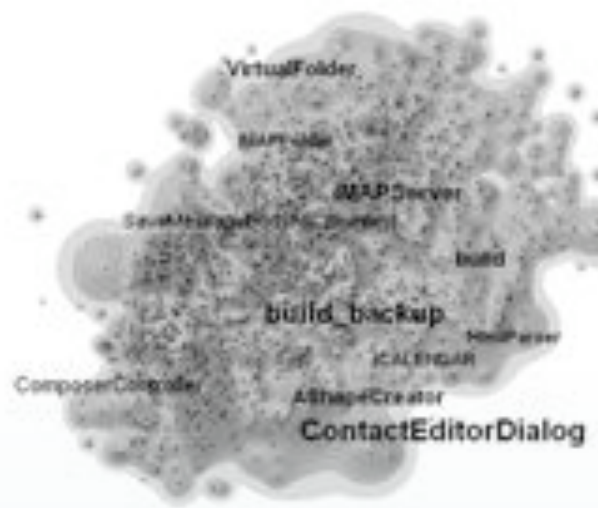


Consistent Layout over Time





Apache Tomcat



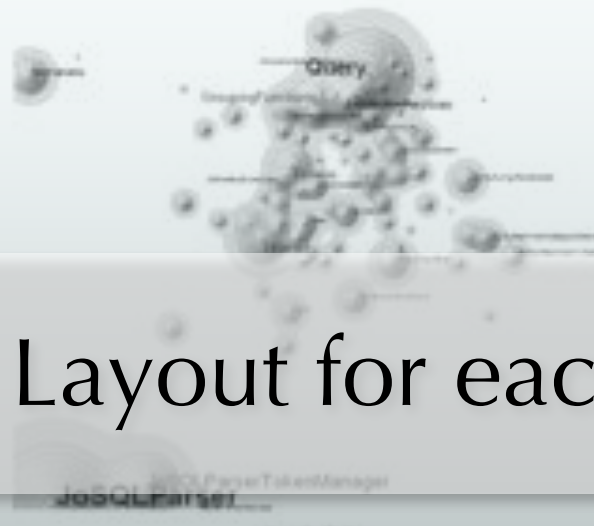
Columba



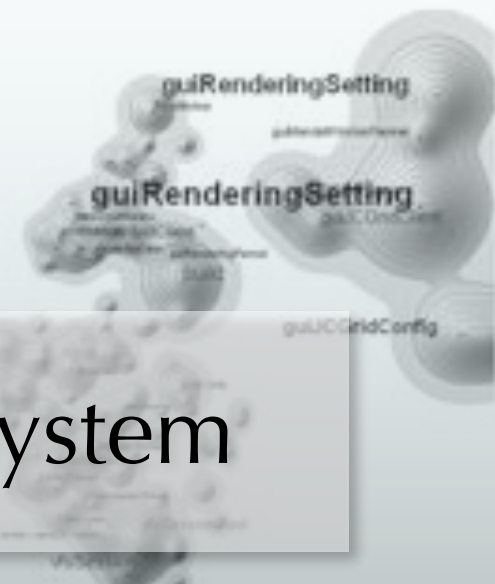
Google Taglib



JFtp

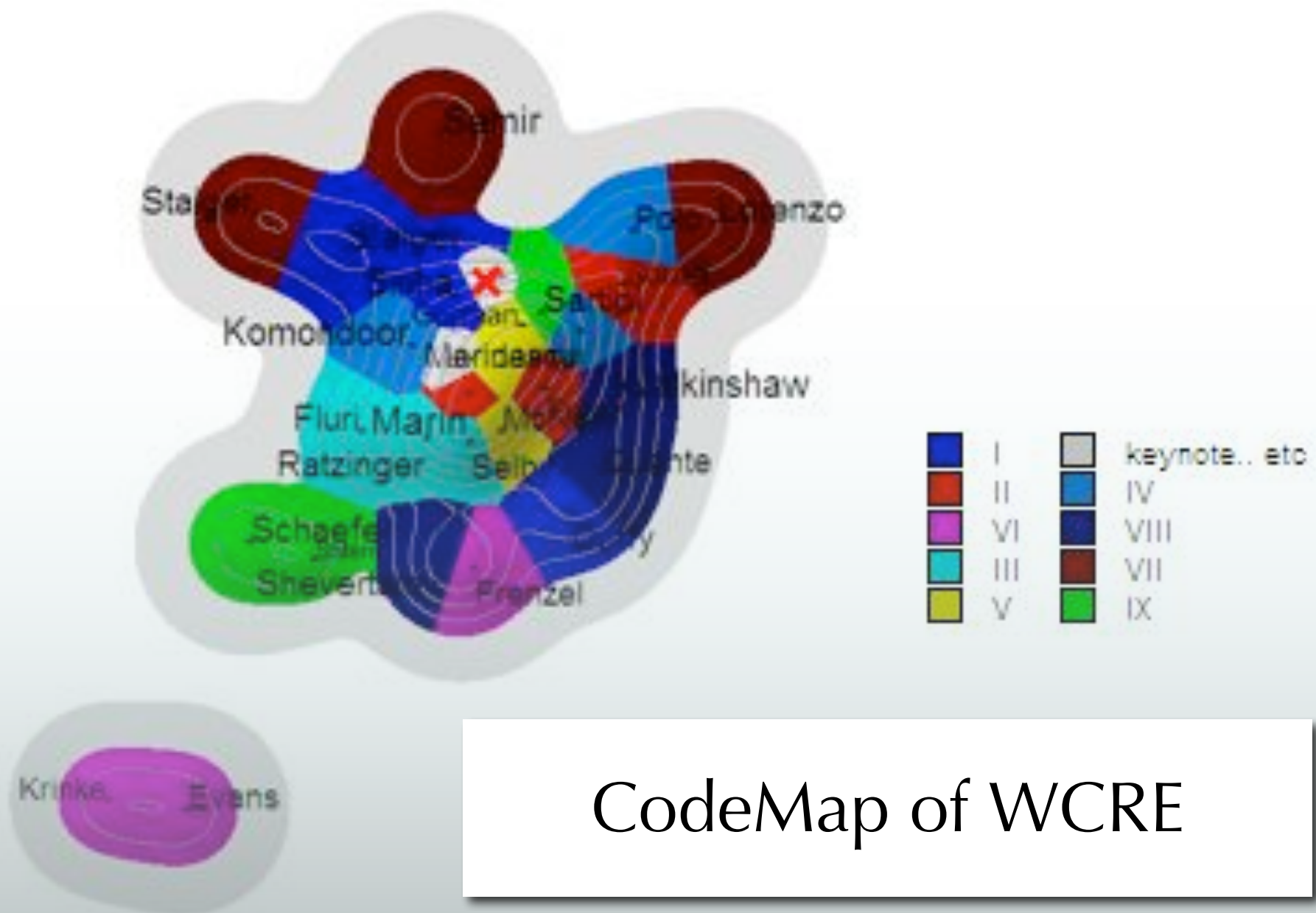


JoSQL



JCGrid

Unique Layout for each System



## CodeMap of WCRE