



# Visual Analysis of Unstructured Data Sets



**Michael Granitzer**

**Know-Center GmbH & Graz University of  
Technology**

<http://www.know-center.at>

# Outline

## • Motivation

- Facetted Retrieval + Scatter/Gather + Some Visual stuff
- Why visual stuff?

## • Clustering Approach (TIR 10)

- Scalable Top-Down recursive Clustering approach with Model Selection
- Experiments

## • Labelling (SIGIR 2010)

- Effects of structural relationships: Parent Child and Sibling Relationships
- Experiments

## • Feedback mechanisms (for discussion)

## • Experiments

- Visual Analysis
- Inex

# Know-Center ?!?

- The Know-Center is Austria's Competence Center for Knowledge-Based Applications and Systems, funded in the COMET program
- Application oriented research  
Bridge the gap between science and industry
- 21 Industry partners, 5 scientific partners (e.g. APA, Bertelsmann, Infonova ...)

Area 1: Knowledge Services – Technology enhanced learning, Context Detection

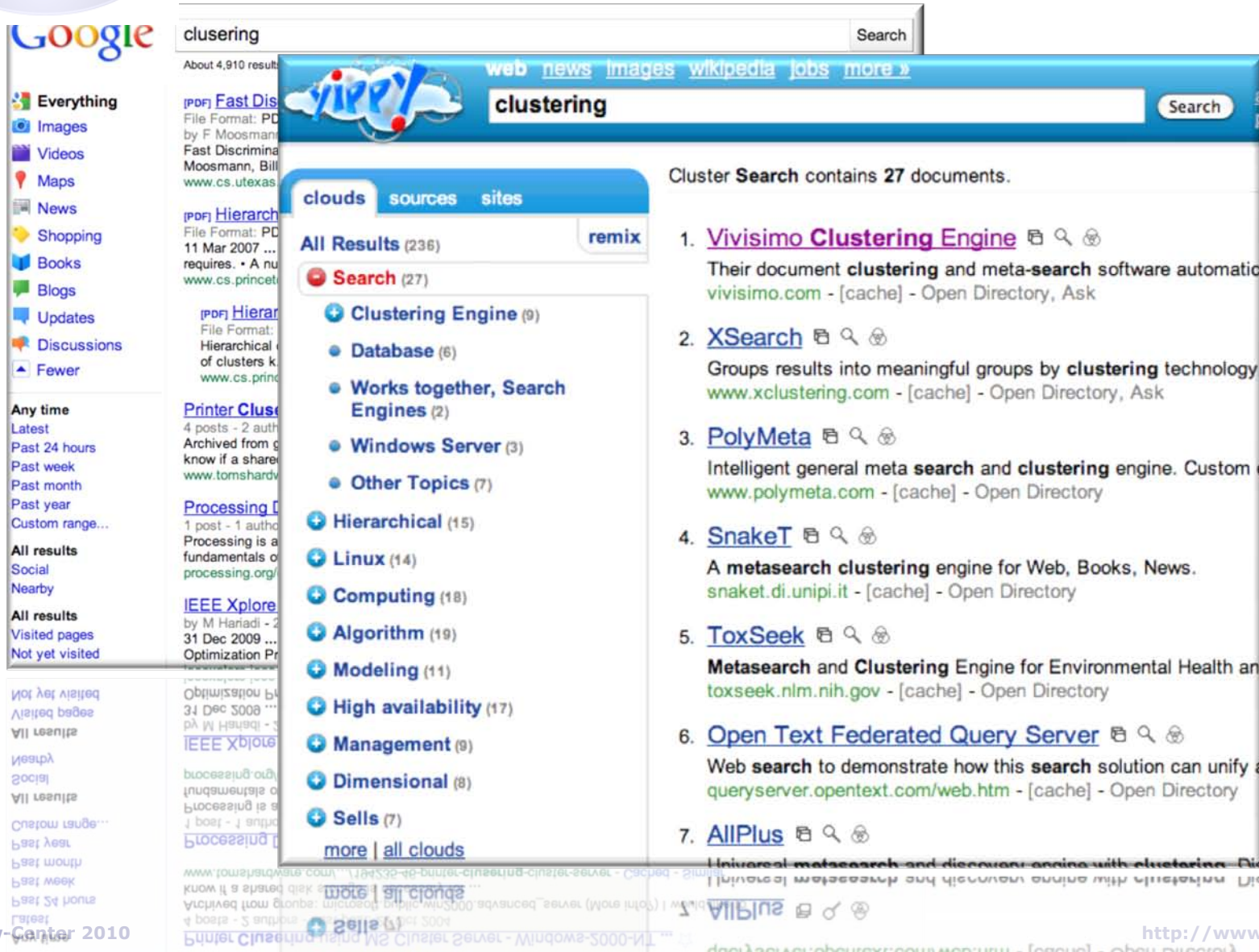
Area 2: Knowledge Relationship Discovery – Text Analysis, Visualisation, Retrieval, Plagiarisma Analysis, Social Media (PAN, CLEF, TREC etc.

Roman Kern, Elisabeth Lex, (Wolfgang Kienreich, Markus Muhr, Vedran Sabol, Christin Seifert, Christopher Horn, Mari Zechner, Werner Klieber)

- Applying Basic Research results in different application scenarios  
Plagiarism Analysis == Media Diffusion Analysis (E.g. „Nike, just Sports“)  
Enterprise Search: not solved

Patent Analysis

# Motivation Facetted Retrieval



The screenshot shows a Google search for "clustering" with approximately 4,910 results. A faceted search interface is overlaid, displaying various categories and their counts:

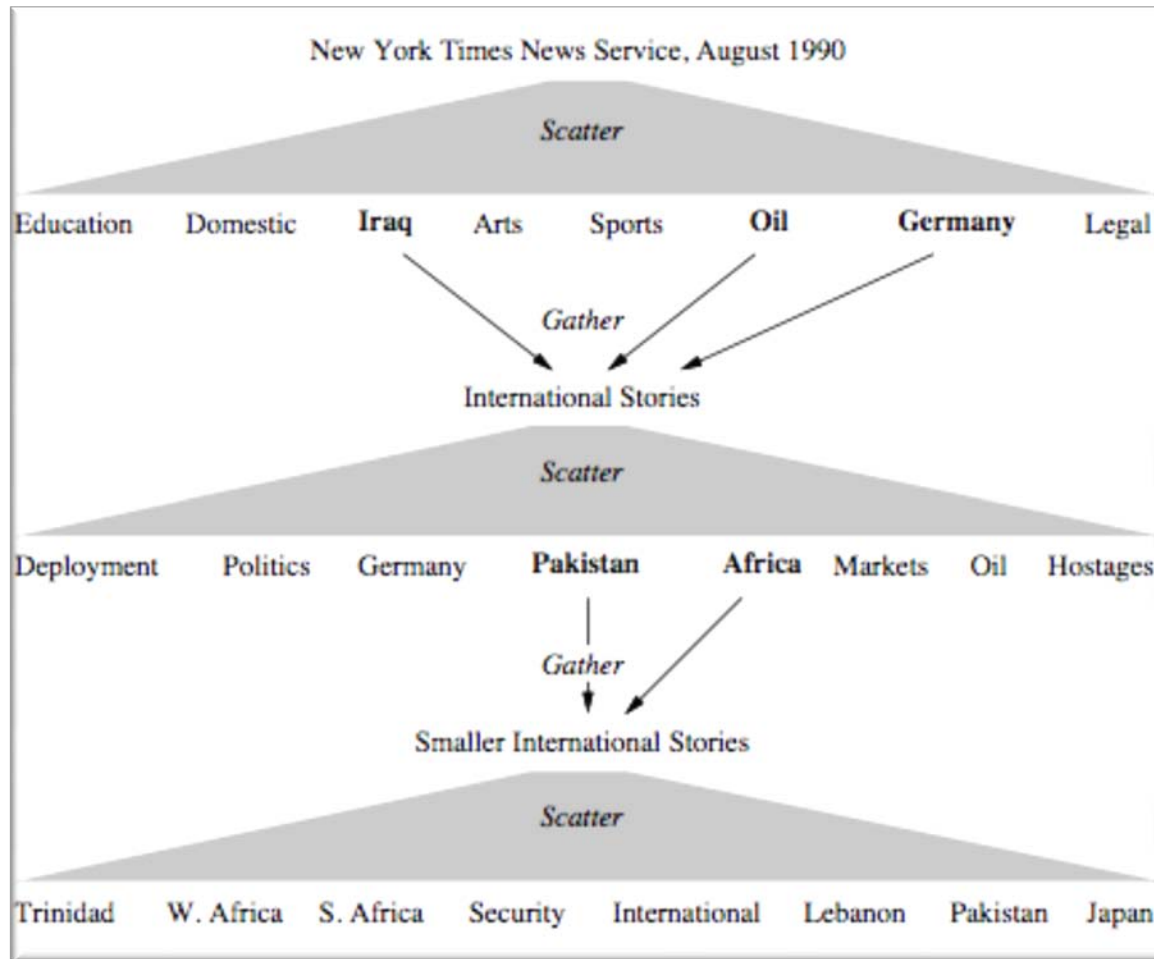
- All Results (236)
- Search (27)
- Clustering Engine (9)
- Database (6)
- Works together, Search Engines (2)
- Windows Server (3)
- Other Topics (7)
- Hierarchical (15)
- Linux (14)
- Computing (18)
- Algorithm (19)
- Modeling (11)
- High availability (17)
- Management (9)
- Dimensional (8)
- Sells (7)

The main search results list includes:

- Vivisimo Clustering Engine**: Their document clustering and meta-search software automatic vivisimo.com - [cache] - Open Directory, Ask
- XSearch**: Groups results into meaningful groups by clustering technology. www.xclustering.com - [cache] - Open Directory, Ask
- PolyMeta**: Intelligent general meta search and clustering engine. Custom e www.polymeta.com - [cache] - Open Directory
- SnakeT**: A metasearch clustering engine for Web, Books, News. snaket.di.unipi.it - [cache] - Open Directory
- ToxSeek**: Metasearch and Clustering Engine for Environmental Health and toxseek.nlm.nih.gov - [cache] - Open Directory
- Open Text Federated Query Server**: Web search to demonstrate how this search solution can unify a queryserver.opentext.com/web.htm - [cache] - Open Directory
- AllPlus**: Universal metasearch and discovery engine with clustering. Dis

# Motivation

## Scatter/Gather [Cutting et. al. 1992]



# Motivation

## InfoSky: Visual Exploration [Andrews et. al. 2002]

The screenshot displays the InfoSky application interface. On the left, a file explorer tree shows a hierarchy starting from 'Root' down to various folders like 'Artificial Intelligence', 'Home\_Automation', etc. The main area is a circular sunburst visualization where each sector represents a folder and its sub-items, with labels such as 'Artificial Life', 'Hacking History', 'Supercomputing', and 'Programming'. The bottom pane shows details for the selected folder, including a table of files and their metadata.

Name	Size	Modified	Keywords
Artificial_Intelligence	1721 documents	Thu Jan 01 01:00:00 CET 1970	Artificial, Intelligence, artificial, intelligence, resources, links, Resources
Home_Automation	97 documents	Thu Jan 01 01:00:00 CET 1970	home, Home, automation, systems, networking, Internet, control
Organizations	308 documents	Thu Jan 01 01:00:00 CET 1970	Computing, technology, FOCUS, Circle, PenDragon, Humour, programming
Robotics	1105 documents	Thu Jan 01 01:00:00 CET 1970	robotics, Robotics, news, discussion, robot, site, Robots

# Motivation

## Visualization, why?

- Exploit the capacity of the visual cortex to immediately recognize certain circumstances
- Example: Preattentive Processing

A restricted set of visual properties can be recognized immediately

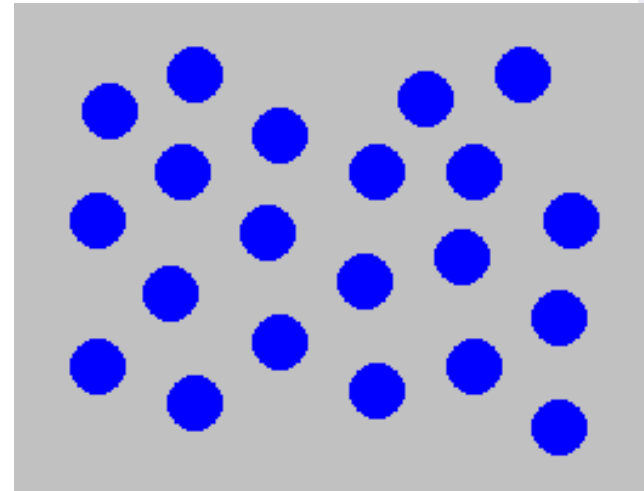
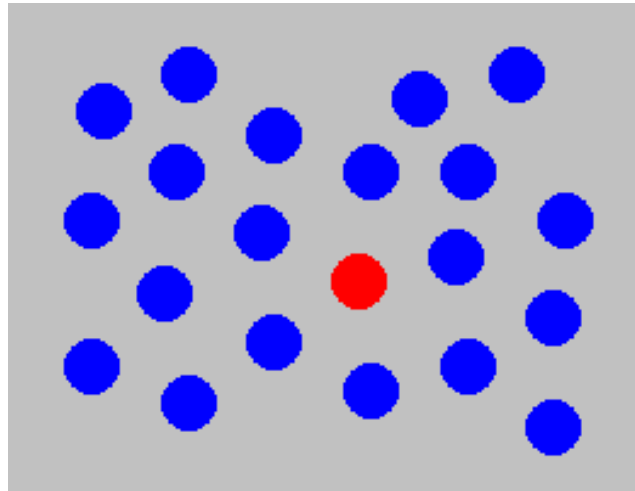
Criteria 1: Processing time below  $<200 - 250\text{ms}$  (within the blink of an eye =  $200\text{ms}$ )

Criteria 2: fixed time period independent of the number of noise

Where is the red circle?

# Motivation

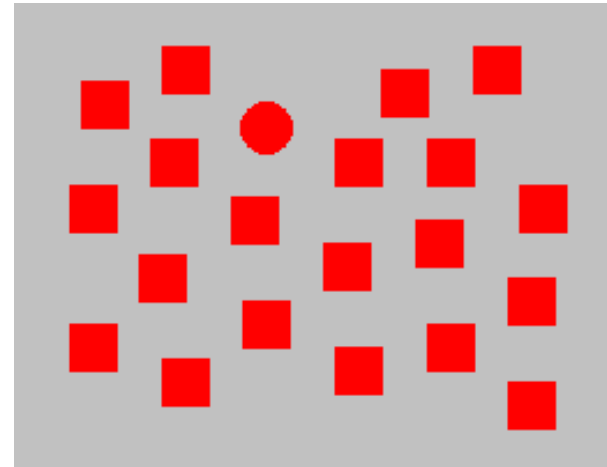
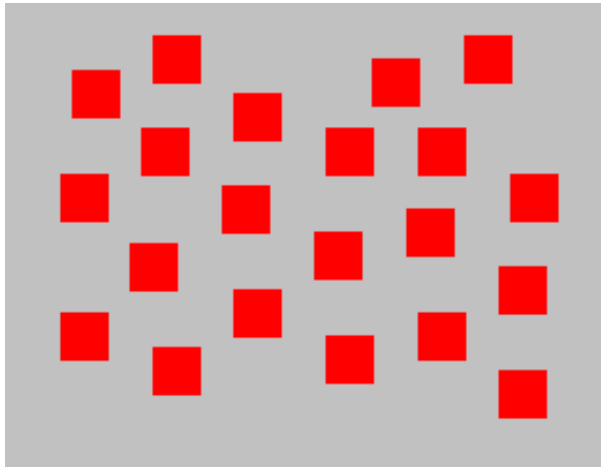
## Visualization, why?





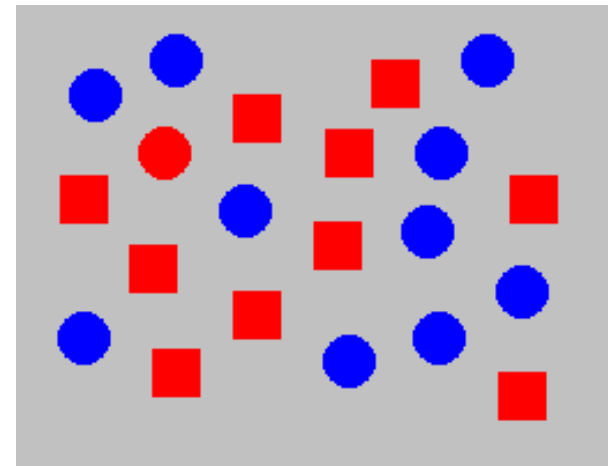
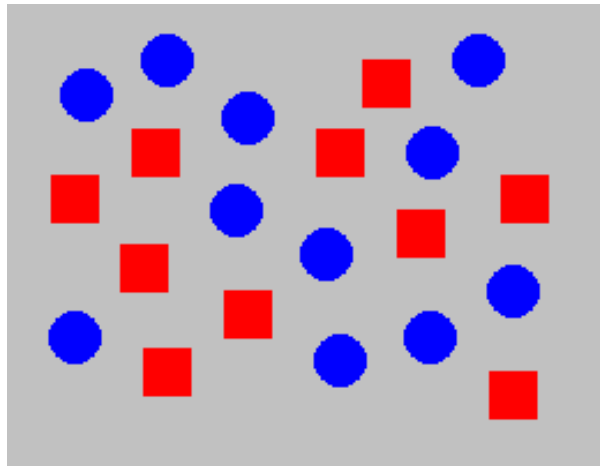
# Motivation

## Visualization, why?



# Motivation

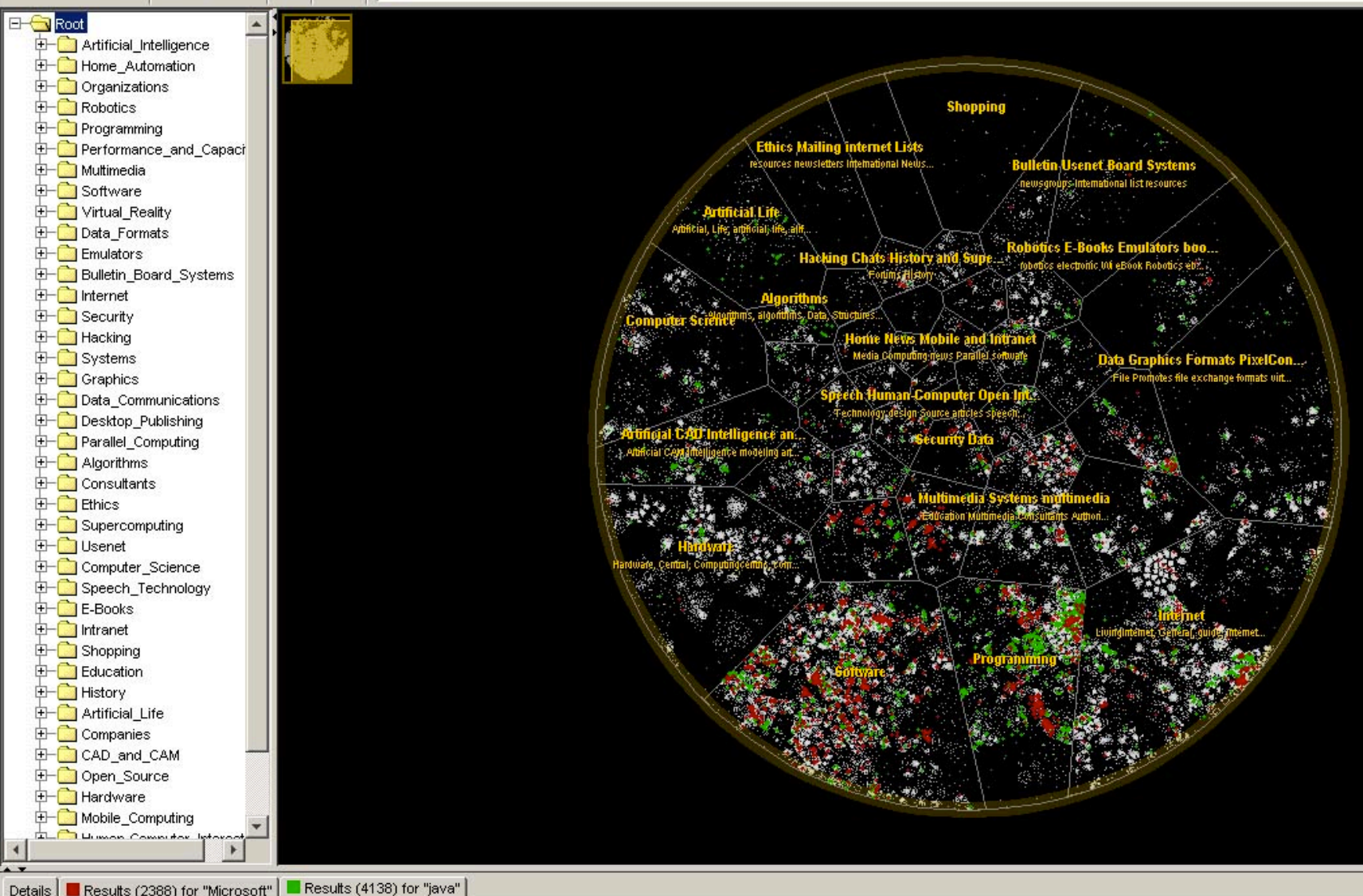
## Visualization, why?



Text is abstract and hardly pre-attentive in contrast to images

# Motivation

# InfoSky: Visual Exploration [Andrews et. al. 2002]





# Approach & Contribution



- Preprocessing
  - Nothingnewwhere.....
- Clustering
  - Combinewell-knownntechniques (Growingk-means, Model Selection....)
- Projection
  - Clustering + Force DirectedPlacement:  $O(n^3) \rightarrow O(n \cdot \log(n))$
- Labelling
  - Label qualitydepends on thehierarchystructure
  - Ad-hocsolution, yet no well foundedtheoreticalapproach
  - Clustering + Force DirectedPlacement:  $O(n^3) \rightarrow O(n \cdot \log(n))$
- Metric Feedback: Just fordiscussions...

- Hierarchical, top-down, polythetic, document clustering approach
- Dynamic cluster structure on each level of the hierarchy supporting splitting and merging of clusters.
- Constraints on the maximum and minimum number of elements per hierarchy level
- Resulting reduced computational costs of the layout algorithm
- Scalable to datasets consisting of millions of documents with a reasonable trade-off between runtime and accuracy



**Top-Down, scalable clustering algorithm for creating a topical hierarchy**



# Clustering Overview

Divide and conquer: decompose into tasks starting at the root node

For every task

- Step 1: Preprocess documents to be clustered
  - Bag-of-Words, BM 25, cosine inner product
- Step 2: Cluster documents using a flat clustering algorithm
- Step 3: Split and merge clusters till constraints are met
- Step 4: Recursion: Evaluate the stopping criterion for dividing into further sub-tasks
- Step 5: Cluster Labeling
- Step 6: Project clusters into a 2 dimensional space



# Clustering

## Step 2: Clustering Algorithm (1/4)

Given a set of documents  $X$ , find a set of  $K$  groups of similar documents (clusters)

- Utilize existing clustering methods

HAC, DBScan or Chameleon  $> O(n^2)$

GNG, BIRCH fast and storage efficient, but order dependent

- Growing k-means

Online Competitive Learning with Winner-takes it all approach

trade-off between runtime and accuracy [Zhao and Karypis 02]

Allows for efficient model selection (determine  $k$ )



# Clustering

## Step 2: Clustering Algorithm (2/4)

---

### Algorithm 1 Growing Spherical K-Means

---

**input:**

$\mathcal{X} = \{x_1, \dots, x_N\}$  with  $x_i \in \mathbb{R}^d$ ,  $K, l, \eta, \nu$

**output:**

$\mathcal{C} = \{c_1, \dots, c_K\}, \mathcal{Y} = \{y_1, \dots, y_N\} \forall y_n \in \{1, \dots, K\}$

**steps:**

initialize centroids  $c_1$  and  $c_2$  by a seeding mechanism

**for**  $m = 2$  to  $K$  **do**

**for**  $n = 1$  to  $N$  **do**

$y_p = y_n$

$y_n = \arg \max_{1 \leq k \leq m} x_n^T c_k$

$c_{y_n} = c_{y_n} + \eta x_n$

$c_{y_p} = c_{y_p} - \nu x_n$

**if**  $\|c_{y_n}\| - 1.0 > l$  **then**

$c_{y_n} = \frac{c_{y_n}}{\|c_{y_n}\|}$

**for**  $n = 1$  to  $N$  **do**

$y_n = \arg \max_{1 \leq k \leq m} x_n^T c_k$

$s_k = s_k + \max_{1 \leq k \leq m} x_n^T c_k$

**if**  $m < K$  **then**

$c_i = \arg \min_{1 \leq k \leq m} S(c_k)$

$x_j = \arg \min_{x \in \mathcal{X}_i} x^T c_i$  with  $\mathcal{X}_i = \{x_n | y_n = i\}$

$c_t = \frac{c_i - x_j}{2}, \mathcal{C} = \mathcal{C} \cup \{c_t\}$

Init and loop for maximum  $k$ -clusters

Update cluster hypothesis

Runtime improvement of centroid update

Assign documents and average similarity

Create  $m$ -th centroid

# Clustering

## Step 2: Clustering Algorithm (3/4)

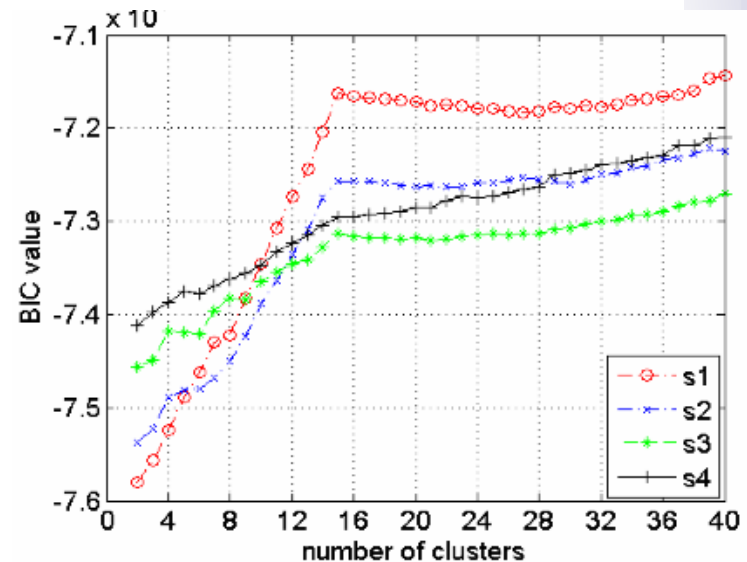
### Model Selection methods

Obtain fitness criterion for different number of clusters (Bayesian Information Criterion (BIC), Stability based approaches)

Monotonical increasing/decreasing

Overtraining on the data

Determine the „best cluster number“  
using knee-point detection  
[Zhao et. al. 2008]



Efficient calculation for the growing k-means by simply calculating the fitness criterion for each new centroid

# Clustering

## Step 2: Clustering Algorithm (4/4)

### Heuristics

- Efficient update rules [Zhong 2005]

Move a fraction of the distance  
between sample and centroid

$$c_{y_n} = \frac{c_{y_n} + \eta(x_n - c_{y_n})}{\|c_{y_n} + \eta(x_n - c_{y_n})\|}$$

Simply update the angle and ignore  
non unit length

Track norm changes and rescale after  
norm exceeds numerical boundaries

$$c_{y_n} = c_{y_n} + \eta x_n$$

$$c_{y_p} = c_{y_p} - \nu x_n$$

**if**  $\|c_{y_n}\| - 1.0 > l$  **then**

$$c_{y_n} = \frac{c_{y_n}}{\|c_{y_n}\|}$$

- Decreasing learning rate with the size of the cluster for balancing

$$\eta = 1 / |\sqrt{\mathcal{X}_{k(x)}}|$$

# Clustering

## Step 3: Split and Merge

Split and Merge Clusters to fulfill the following constraints

- # Cluster at one level

Merge the most similar cluster if  $\#cluster > \text{maximum number of clusters}$

Split the least coherent or biggest cluster if  $\#cluster < \text{minimum number of clusters}$

- # documents in a cluster

Below the Maximum number of documents for a cluster →  
`clusterokforbrowsing`

More than 1.5 times the upper limit to ensure meaningful clustering at next hierarchical level

If all clusters fulfill this constraint, cluster recursively  
(Step 4)



# Clustering Experiments

## INEX Clustering

- Initiativ for Evaluation of XML Retrieval
- XML Mining Track – Cluster the English Wikipedia
  - Small data set 54k documents
  - Large data set 2.6 Million Documents
  - Preprocessed document vectors (uni and bi-grams)
- Ground truth provided by YAGO ontology, but no hierarchical structure
- Document assigned to each cluster on the path to facilitat multi cluster assignment as it is the case in Wikipedia



# Clustering Experiments

## INEX Clustering

- 10,467 Clusters for the small data set

4 Minutes to compute on a 16GB Quad Core including I/O

MacroPurity	BIC	Stability
73k Categories	0.4959	0.4945
12k Categories	0.5473	0.5303

- 133,704 Clusters on the large data set

Runtime 2 hours

348 k Categories: Macro Purity of 0.4457

12k Categories: Macro Purity of 0.5359

- Clusters appear to be reasonable, but good evaluation strategy remains an open issue

High level clusters are more important

Accurate ground truth reflecting good browsing strategies



# Clustering

## Step 5: Labeling - Overview

- Labeling via Jensen Shannon Divergence

  - How to achieve good labeling quality for browsing?

  - Does level of the hierarchy have an impact on the labeling quality?

- Intuition

  - Take structural relationships into account to improve labeling quality

  - Siblings - labels should help to separate neighbor clusters

  - Hierarchies -

    - labels should become more generic the higher the cluster is within the hierarchy

- Open Issues here

  - Most state-of-the-art labeling approaches do not exploit structural relationships

  - No standardized test dataset

  - No evaluation for browsing purpose



# Clustering Labeling - Approach

- Extend existing well-known labeling techniques by structural relationships

- Maximum term weight based methods

- Reference collection based methods

- Types of structural relationships

- Sibling relationships

- Parent-child relationship

Assumption: All labeling algorithms are based on a bag of word model.  
Extension possible with bi-grams, tri-grams etc.



# Clustering Labeling – Maximum Term Weight Labeling

- ▶ Pick the top  $k$  terms according to a weighting scheme by summing over all cluster documents
  - ▶ Local weights
  - ▶ Global weights (IDF, BM25)
  - ▶ Named in the evaluation:  $MTWL_{raw}$

$$L_j \leftarrow best_k \left( \sum_{d_i \in \mathbf{D}_{c_j \rightarrow *}} idf_{global} \cdot tfWeight(d_i) \right)$$



# Clustering Labeling – Reference Collection based

- ▶ Compare the distribution of terms within a cluster with a reference collection
  - ▶  $\chi^2$  Popescul and Ungar [2000]
  - ▶ Information Gain Geraci et al. [2007]
  - ▶ Jensen-Shannon Divergence Carmel et al. [2006]
  - ▶ Named in the evaluation: *JSD*

$$L_j \leftarrow \text{best}_k \left( JSD(\mathbf{D}_{ref}, \mathbf{D}_{c_j \rightarrow *}) \right)$$



# Clustering Labeling – Inverse Cluster Weight Labelling (ICWL)

How to exploit the sibling relationship?

- ▶ Follow the approach of the CFC classification algorithm Guan et al. [2009]
- ▶ *Intuition*: If one term occurs often in one sibling cluster only, this term should be preferred over terms occurring in all sibling clusters
- ▶ Integrate sibling weighting into the maximum term weight labeling
- ▶ Named in the evaluation:  $ICWL_{raw}$

$$icf_{j,k} = \exp \left( \frac{\#(t_k, \mathbf{D}_{c_j \rightarrow *})}{|\mathbf{D}_{c_j \rightarrow *}|} \right) \log \left( \frac{\#(c_p)}{\#(t_k, c_p)} + 1 \right)$$

# Clustering Labeling – Hierarchical Labelling

## How to exploit the parent-child relationship?

- ▶ *Intuition*: Integrate the path length (distance between cluster to label and document a term occurs in) into the label calculation and promote terms occurring in a higher number of child clusters.
- ▶ Hierarchical labeling extends all introduced labeling approaches
- ▶ Added prefix *hier* in the evaluation

$$L_j \leftarrow \text{best}_k \left( \sum_{c_i \in \mathbf{C}_{c_j \rightarrow *}} \frac{1}{l(j, i)} * cf_{l(j, i)} \cdot v_{j, i} \right)$$



# Clustering Labeling - Evaluation

- Open Directory Project (ODP)

  - Top categories: arts, business, games, health, home, news, society, sports

  - Ignored soft links, ignored single letter categories

- Wikipedia

  - Top categories: arts, computing, health, sports

  - Restricted to 10 sub-categories and 80 articles (drawn randomly)

  - Ignored internal categories, ignored "authors by year" categories, limited number of documents per category, ignored cycles

- Oshumed

  - Mesh tree hierarchy

  - Documents only at leaf categories



- European Patents

  - Years 1991-2000

  - IPC classification hierarchy

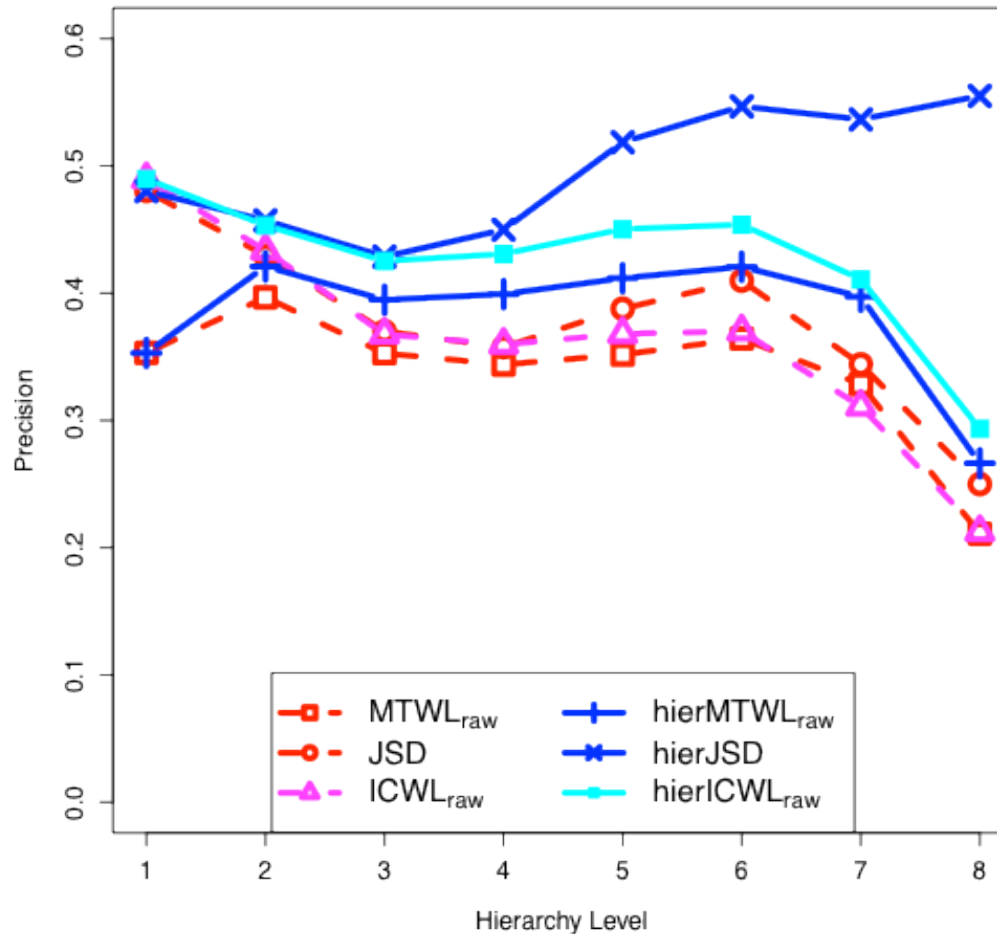
# Clustering Labeling - Evaluation

	Categories	Documents
ODP	150,000	800,000
Wikipedia	50,000	400,000
Oshumed	7,724	348,564
Patents	60,000	265,409

*Preprocessing:* Tokenized *openNLP*, Stemmed , Stop-word removal 

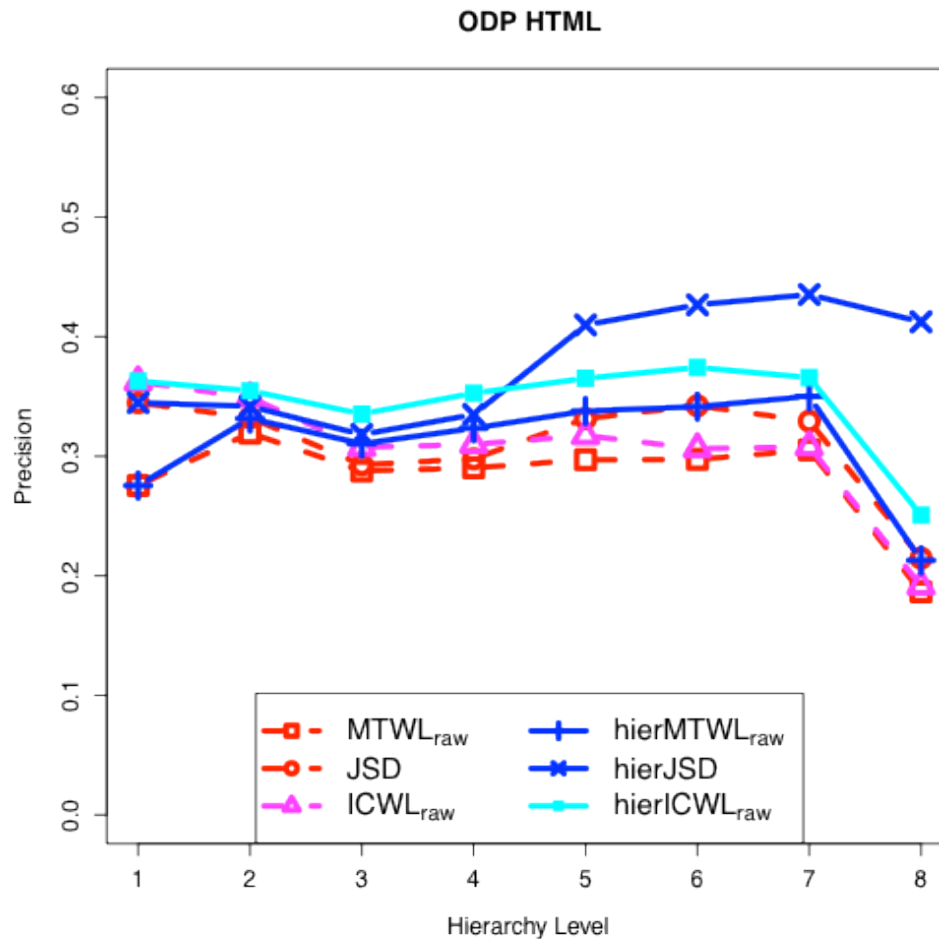
# Clustering Labeling - Evaluation

- Precision over hierarchies with different depths  
ODP Title & Description



# Clustering Labeling - Evaluation

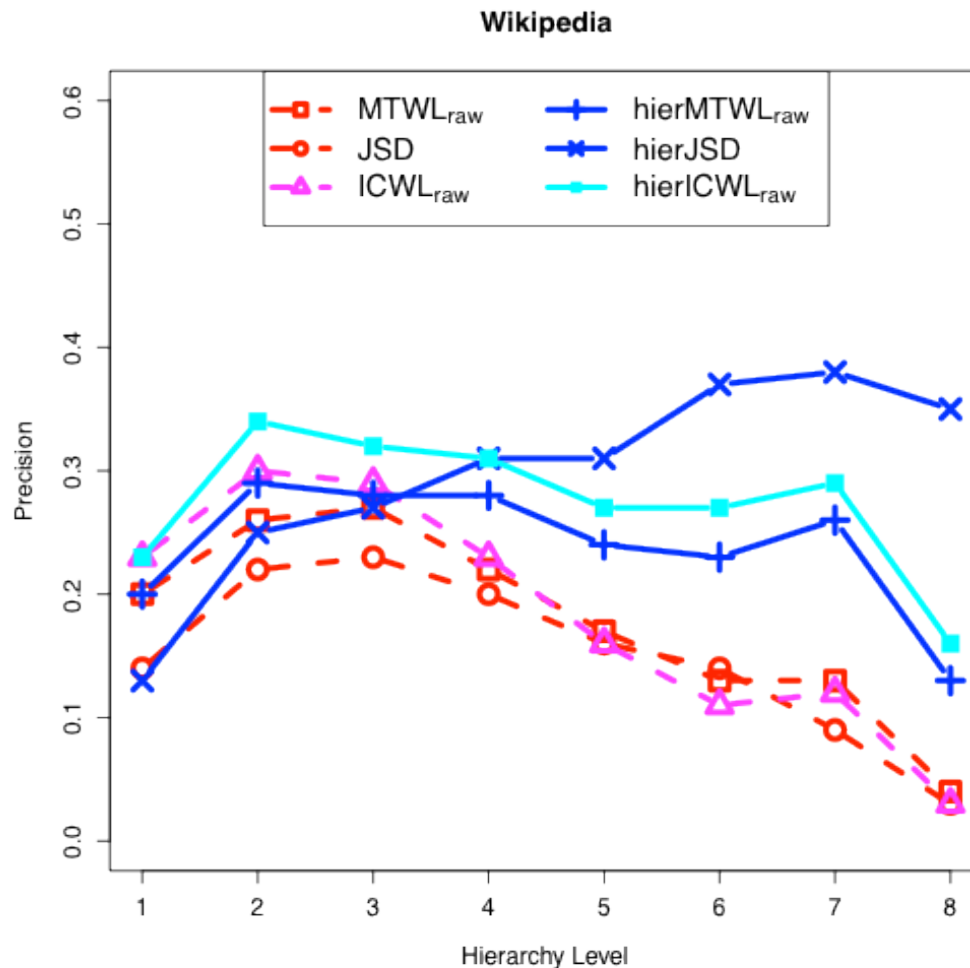
- Precision over hierarchies with different depths





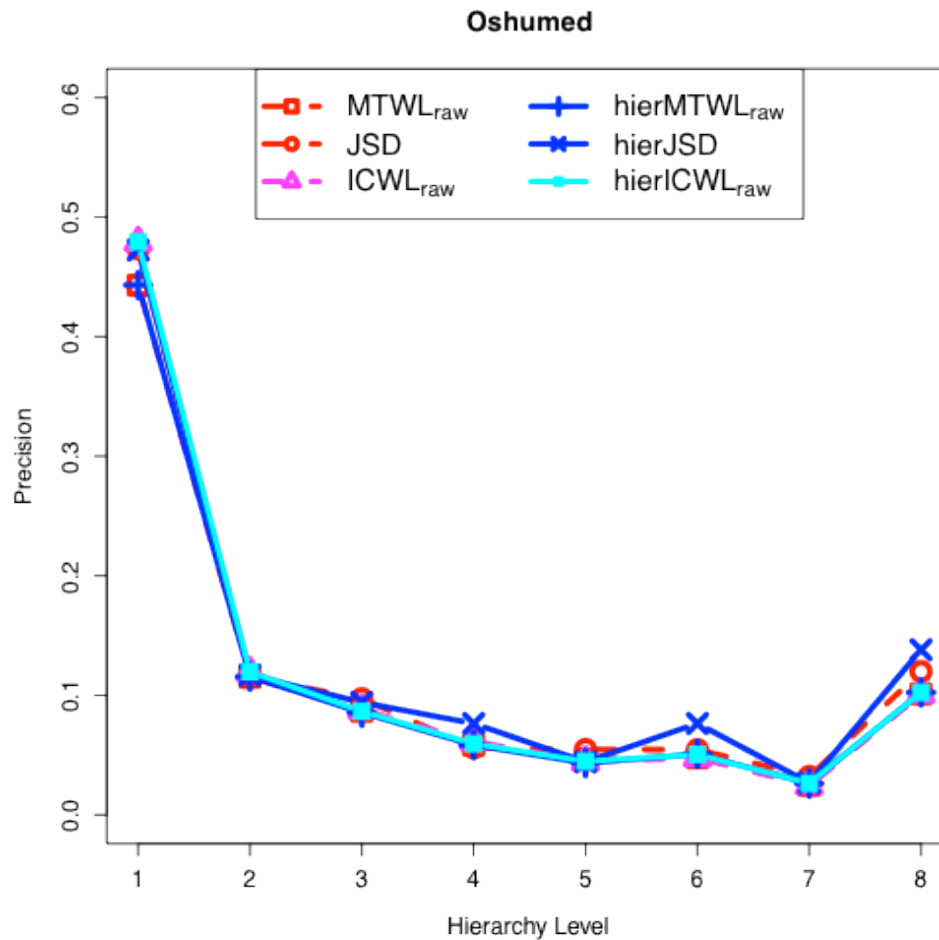
# Clustering Labeling - Evaluation

- Precision over hierarchies with different depths



# Clustering Labeling - Evaluation

- Precision over hierarchies with different depths



# Clustering Labeling – Evaluation - Summary

## ● Sibling Relations

No impact on ODP and Ohsumed

Slight improvement over the respective MTWL methods for the Wikipedia dataset

## ● Summary Parent Child Relations

	<i>MTWL<sub>raw</sub></i>	<i>JSD</i>	<i>ICWL<sub>raw</sub></i>
ODP - T&D	0.06	0.15	0.08
ODP - HTML	0.04	0.09	0.05
Wikipedia	0.08	0.19	0.12
Oshumed	0.00	0.01	0.00

**Table:** Average relative difference of the precision for all hierarchy levels greater than 2 for all datasets between the different methods either with and without exploitation of hierarchical information.



# Clustering Labeling – Evaluation - Conclusio

## Interpretation of the Results

- Flat labeling approachessupportthebrowsing of leafnodesratherthanthebrowsing of high levelnodes → a resultquitecontradictory to theusersneed
- Usingsiblinginformationincreaseslabelingaccuracy in somedatasets
- Integratinghierarchicalinformationproducesbetterlabelin resultsfor all datasets
- Labelingaccuracyisstronglydomainindependent

# Clustering

## Step 6: Projection

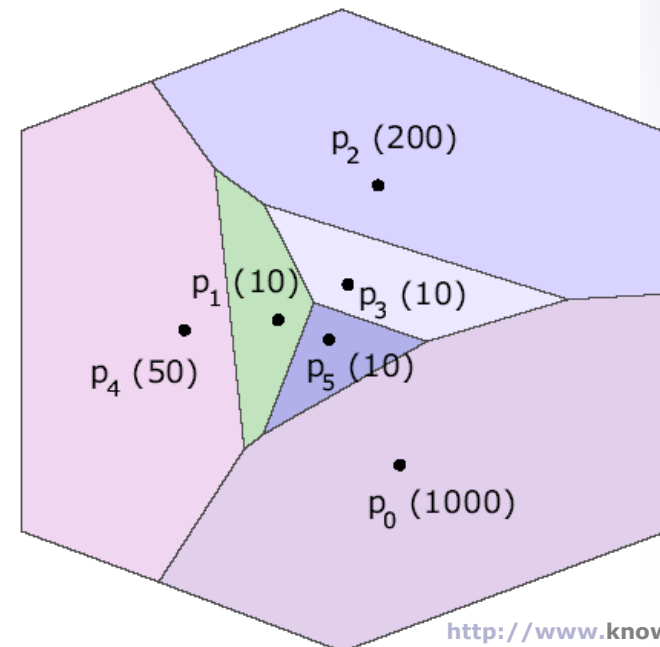
- Projection [Andrews et. Al. 2004]

Force directed placement  $O(n^3)$

Recursive application on cluster hierarchy using document and cluster centroids as points to layout

Due to the constraints we achieve a runtime of roughly  $O(n \cdot \log(n))$

Voronoi inscription of rectangular Layout



# Clustering

## Step 7: Metric Feedback

Not implemented/analysed yet

- High dimensional distances → Low Dimensional Distances
- User moves points on the plain
- New Low Dimensional Distances → Update high dimensional similarity
- Metric Learning: There are some approaches

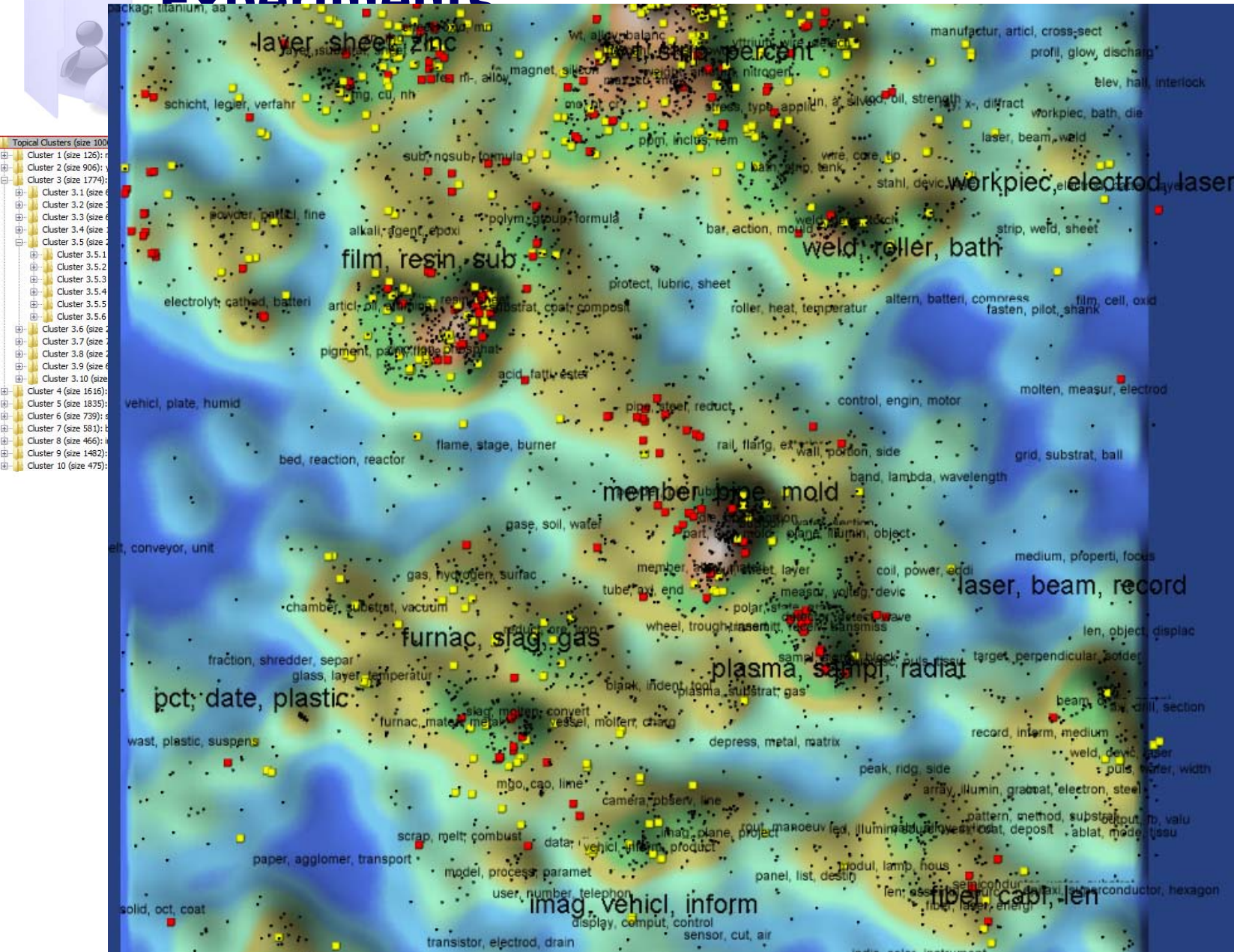
Donald Metzler and Hugo Zaragoza. Semi-parametric and non-parametric term weighting for information retrieval. In Proceedings of the 2nd International Conference on the Theory of Information Retrieval (ICTIR 2009), 2009.

Marco Ernandes, Giovanni Angelini, Marco Gori, Leonardo Rigutini, and Franco Scarselli. Adaptive context-based term (re)weighting: an experiment on single-word question answering. *Frontiers in Artificial Intelligence and Applications*; Vol. 141, page 1, 2006.

Shai S. Shwartz, Yoram Singer, and Andrew Y. Ng. Online and batch learning of pseudo-metrics. In *ICML '04: Proceedings of the twenty-first international conference on Machine learning*, pages 94+, New York, NY, USA, 2004. ACM

Granitzer M., Adaptive Term Weighting through Stochastic Optimization, *CICLING 2010*, Springer

# Experiments



isio



# Experiments

## Clustering based Visualisation

- Not for search, but for analysis of unstructured text documents
- Preliminary user evaluation
  - Combination of visualisation and standard components helpful for explorative tasks [Andrews et. Al. 2002]
  - Improved interaction and navigation paradigms to support explorative search tasks
  - Patent analysis tasks improved in real world use case
  - Suitable for high recall search tasks
- Detailed evaluation still missing
- Similarity biases results
- ?? Could the user be utilized to learn similarity metrics via such visualisations??





# Summary & Conclusio

- Support explorative search and analysis tasks, not standard retrieval
- Top-down, recursive algorithm with different model selection strategy to scale

K-Means based approaches simply work well, invest in features in stead of algorithms

- Labeling exploiting hierarchical relationships improves labeling accuracy

External resources + hierarchical relationships + !bag-of-words=  
??



**Evaluation for Browsing behaviour hard to conduct:  
Missing measures & datasets; no comparison to  
literature**

Thank for your attention  
**Questions?**



**Michael Granitzer**

Scientific Director  
Know-Center Graz  
Inffeldgasse 21a  
8020 Graz

+43 316 873 9263  
**mgrani@know-center.at**  
[www.know-center.at](http://www.know-center.at)