

# Scientific Text Mining and Knowledge Graphs

## Chapter 2

### Part 1: Taxonomy Construction

**Presenter: Jingbo Shang**

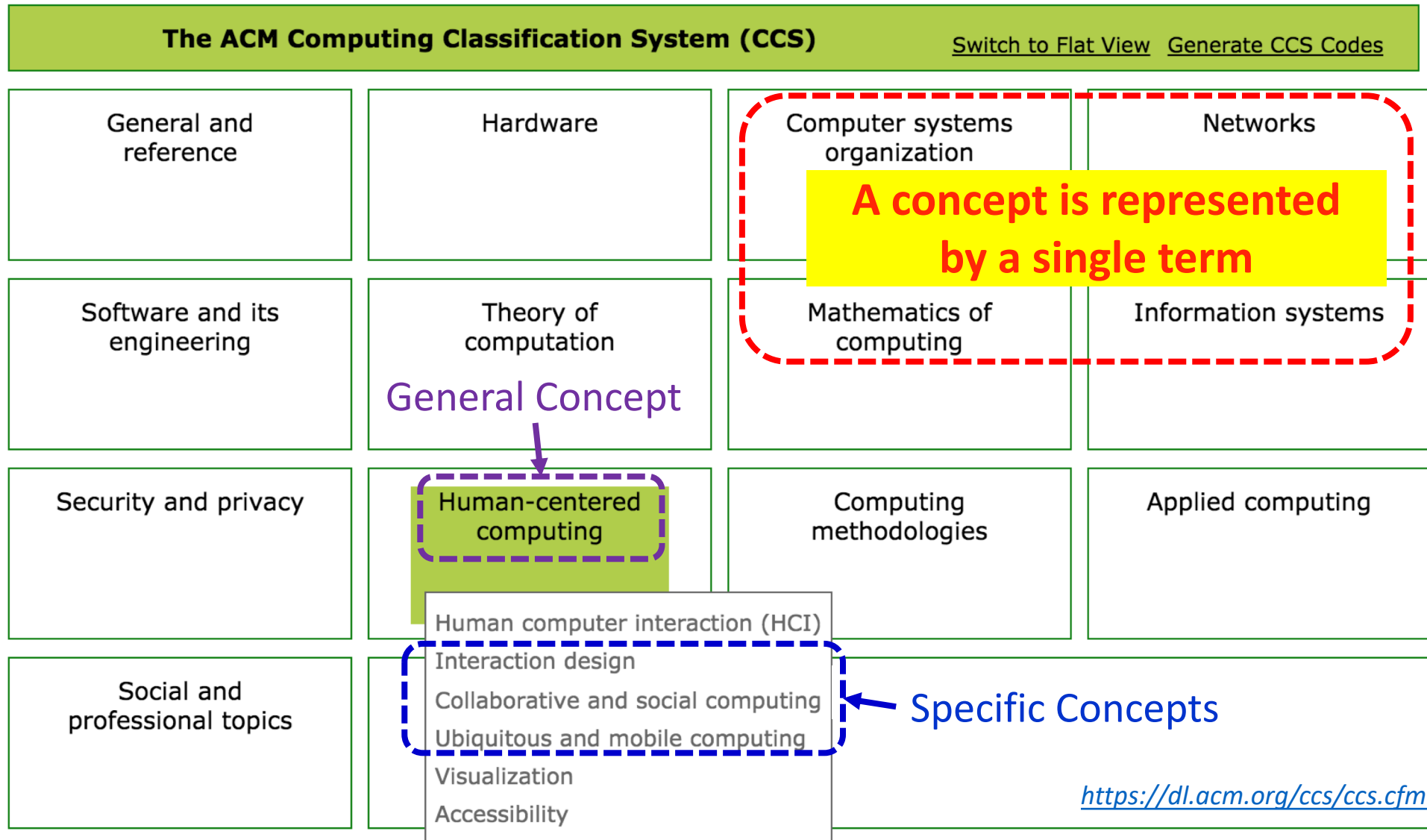
University of California, San Diego

[jshang@ucsd.edu](mailto:jshang@ucsd.edu)

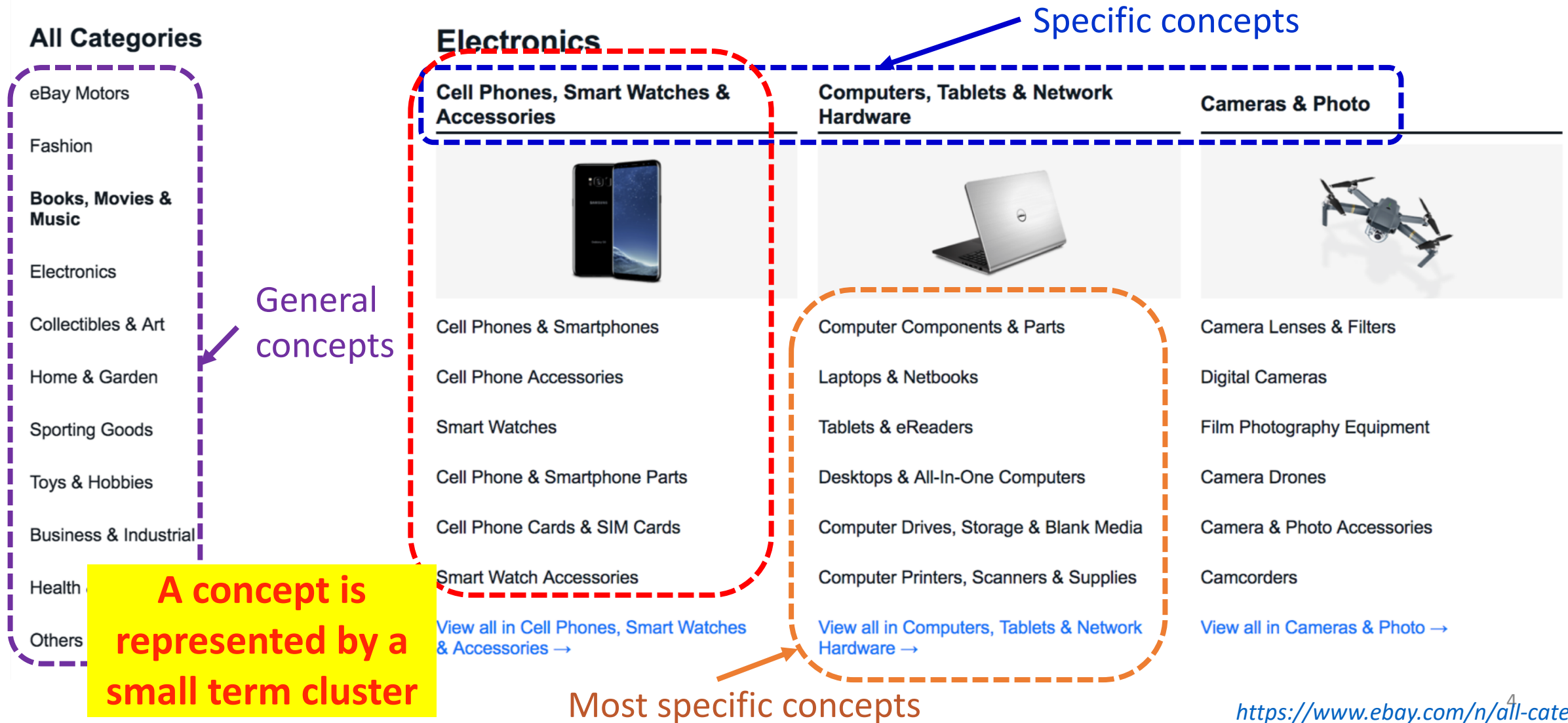
# What is a taxonomy?

Taxonomy is the practice and science of **describing** and **organizing concepts**

# Example: ACM CCS Taxonomy



# Example: eBay Product Taxonomy

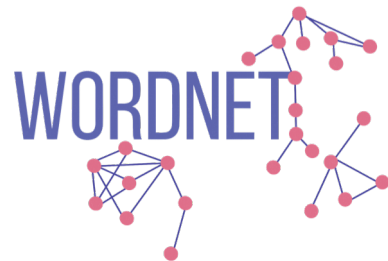


# Two General Types of Taxonomy

□ Each concept in the taxonomy is called a ***taxon*** and based on how we represent a taxon, we categorize taxonomies into two types:

□ **Instance-based** Taxonomy – each taxon is a single term (+ strict synonyms)

ACM  
CCS



ProBase

□ **Clustering-based** Taxonomy – each taxon is a topically related term cluster

ebay

amazon

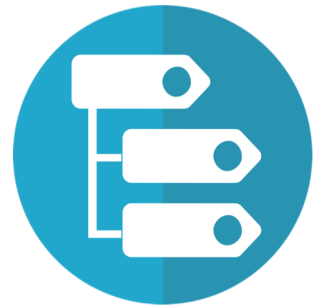


wayfair

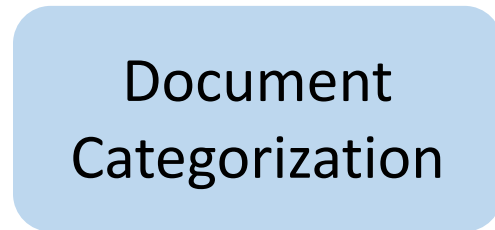
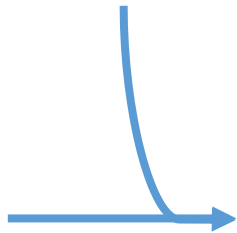
# Taxonomy Underpins Digital Library



**Scientific Papers**



**Taxonomy**



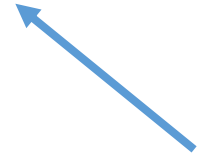
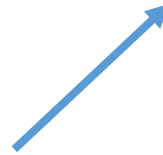
*Data Mining Papers*



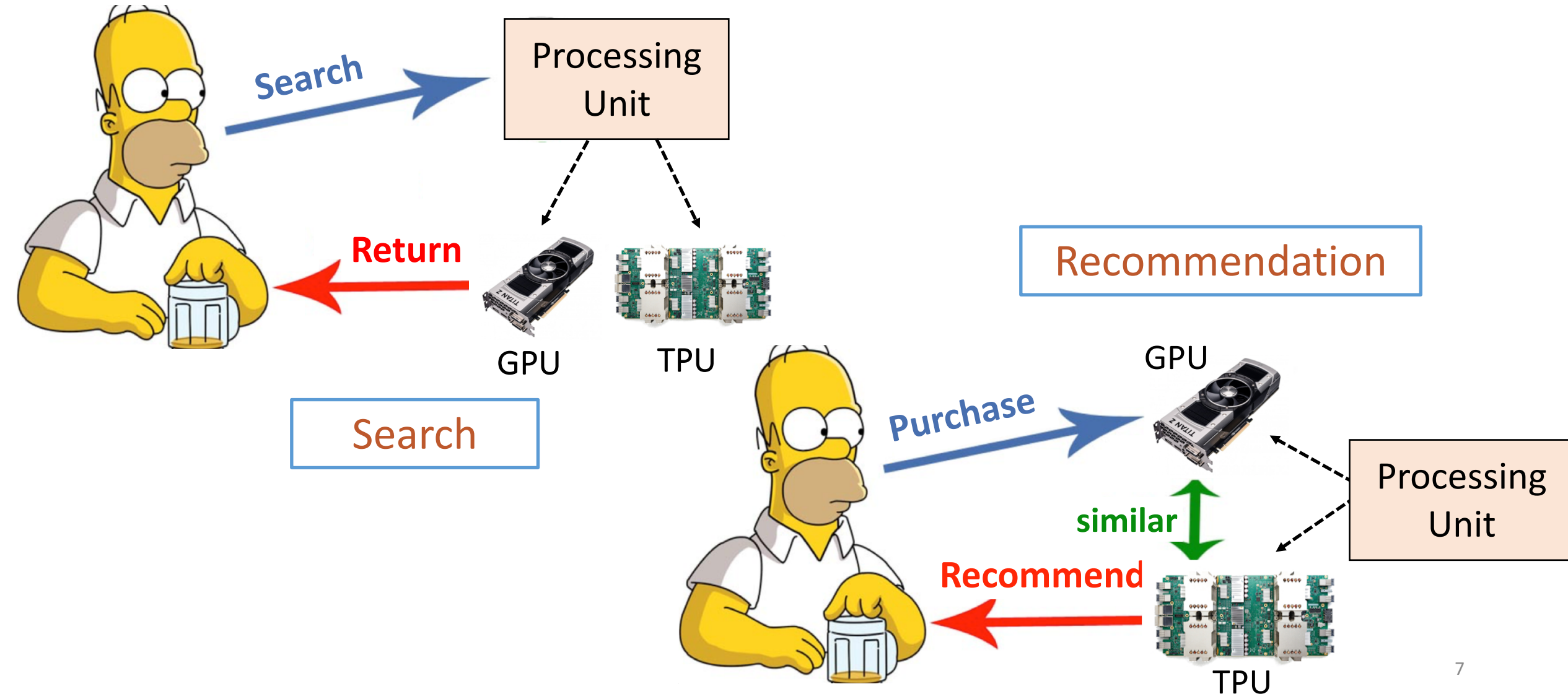
*Theory Papers*



*HCI<sub>6</sub> Papers*

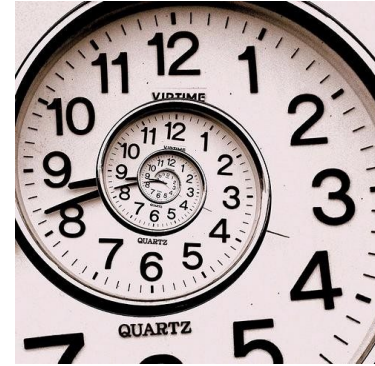


# Taxonomy Helps Search & Recommendation



# How to Build a Taxonomy?

- ❑ Manual curation
  - ❑ Time-consuming and expensive
  - ❑ Human (expert) labor-intensive
- ❑ Automated construction
  - ❑ Scalable and extensible



# Typical Taxonomy Construction Process

Extract Terms

Decide Taxon  
Representation

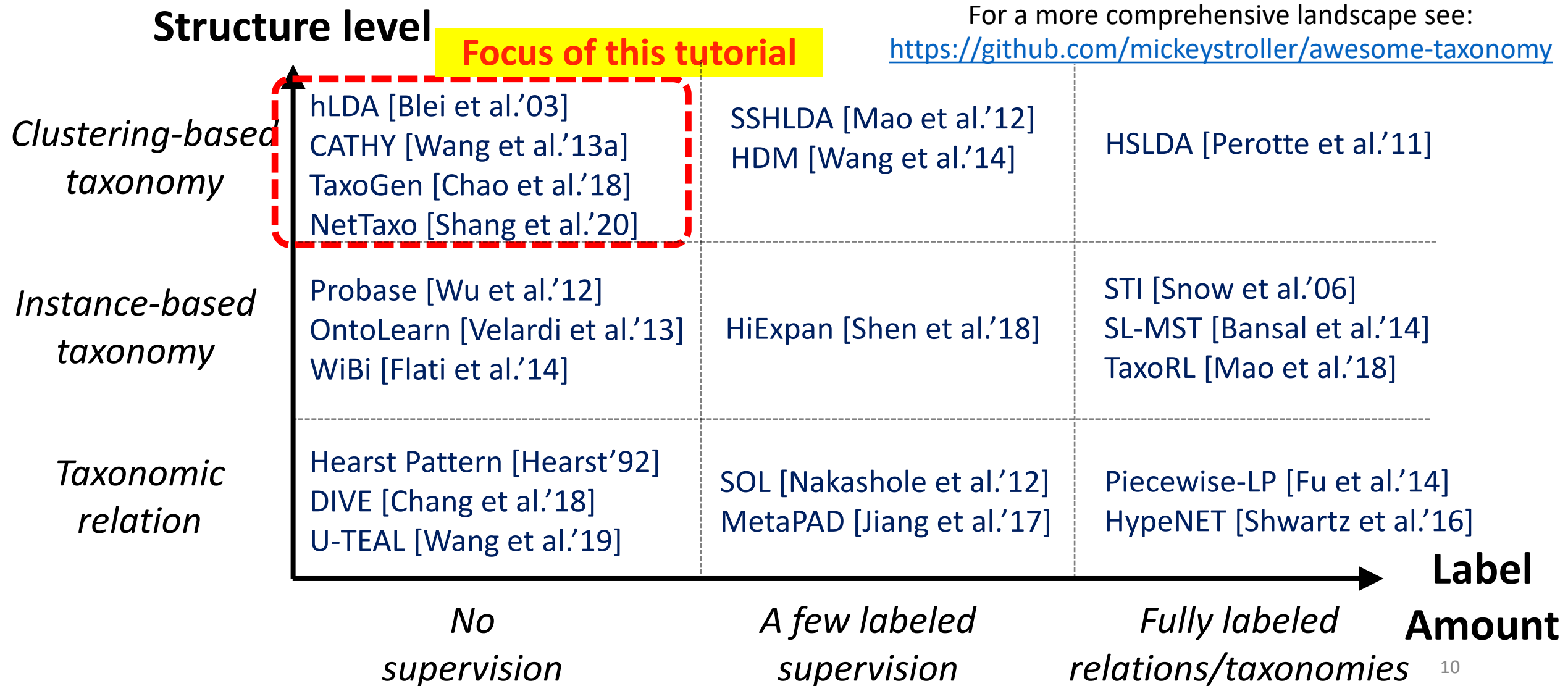
Determine Taxonomic  
Relations

Construct Taxonomy  
Global Structure

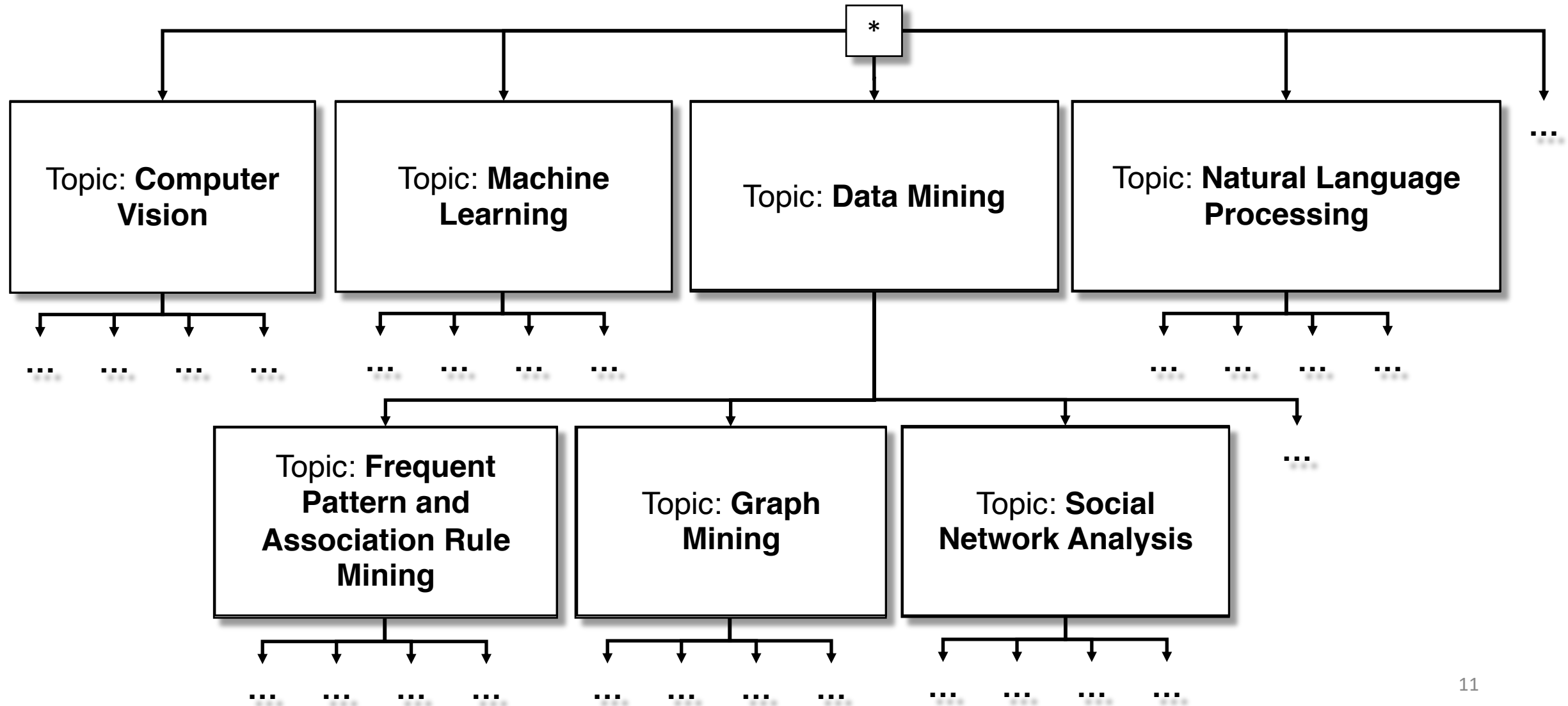
- What data sources you have?
  - Raw text? Query logs? Clicks? Networks?
- ☐ What type of taxonomy you want?
  - ☐ Instance-based? Clustering-based?
- ☐ What type of relation you want?
  - ☐ Is-A relation? Part-of relation? Other relation?
- ☐ How many supervision you have?
  - ☐ Many labeled taxonomies? Several seed taxons?

**Focus of this part of tutorial**

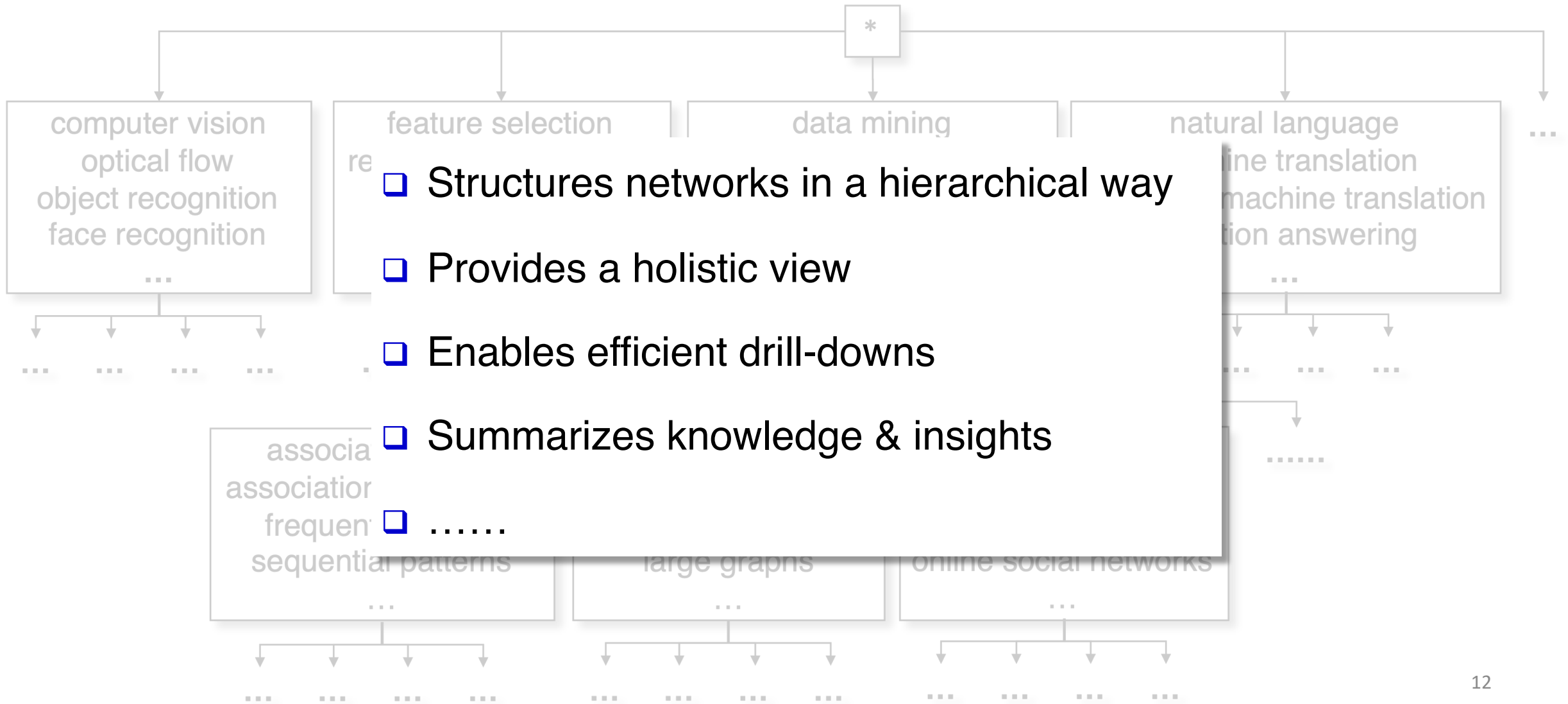
# Taxonomy Construction Methods: A Landscape



# Constructed Topic Taxonomy: Example



# Why Topic Taxonomy?



# Hierarchical Topic Model

- Use a cluster of terms (i.e., a topic) to represent a concept and organize topics in a hierarchical way
- Pose different statistical assumptions on the data generation process
  - Nested Chinese Restaurant Process:
    - hLDA [Blei et al.'03], hLDA-nCRP [Blei et al.' 10]
  - Pachinko Allocation Model:
    - PAM [Li and McCallum'06], hPAM [Mimno et al.'07]
  - Dirichlet Forest Model :
    - DF [Andrzejewski et al.'09], Guided HTM [Shin and Moon'17]

# Example: hLDA

## Document generation from Chinese Restaurant Process

1. Let  $c_1$  be the root restaurant.
2. For each level  $\ell \in \{2, \dots, L\}$ :
  - (a) Draw a table from restaurant  $c_{\ell-1}$  using Eq. (1). Set  $c_\ell$  to be the restaurant referred to by that table.
3. Draw an  $L$ -dimensional topic proportion vector  $\theta$  from  $\text{Dir}(\alpha)$ .
4. For each word  $n \in \{1, \dots, N\}$ :
  - (a) Draw  $z \in \{1, \dots, L\}$  from  $\text{Mult}(\theta)$ .
  - (b) Draw  $w_n$  from the topic associated with restaurant  $c_z$ .

Generates



We develop an approach to risk minimization and stochastic optimization that provides a convex surrogate for variance, allowing near-optimal and computationally efficient trading between approximation and estimation error.

**“Observed” documents**

Inference

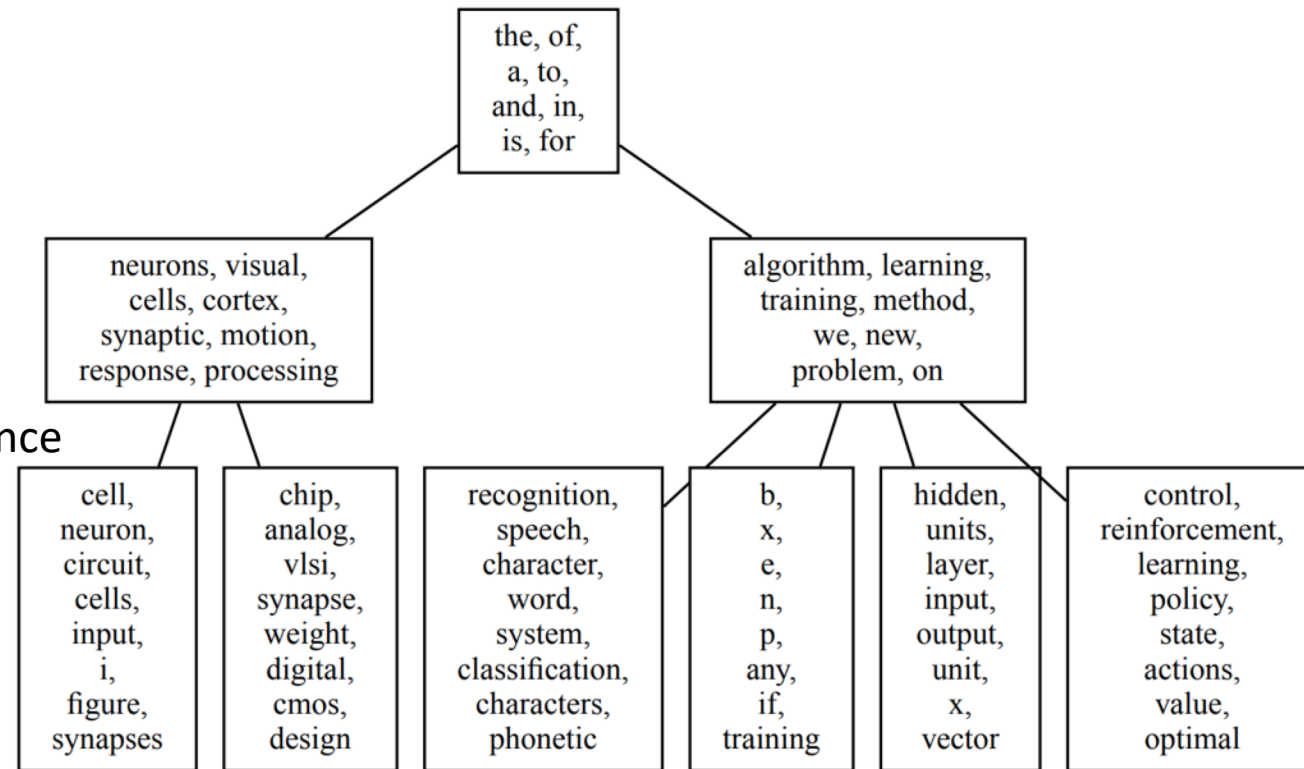
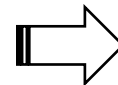


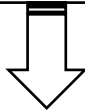
Figure credits to [Blei et al.'03]

# Example: hPAM

## Document generation from Pachinko Allocation Model

1. For each document  $d$ , sample a distribution  $\theta_0$  over super-topics and a distribution  $\theta_T$  over sub-topics for each super-topic.
2. For each word  $w$ ,
  - (a) Sample a super-topic  $z_T$  from  $\theta_0$ .
  - (b) Sample a sub-topic  $z_t$  from  $\theta_{z_T}$ .
  - (c) Sample a level  $\ell$  from  $\zeta_{z_T z_t}$ .
  - (d) Sample a word from  $\phi_0$  if  $\ell = 1$ ,  $\phi_{z_T}$  if  $\ell = 2$ , or  $\phi_{z_t}$  if  $\ell = 3$ .

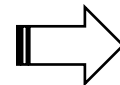
Generates



We develop an approach to risk minimization and stochastic optimization that provides a convex surrogate for variance, allowing near-optimal and computationally efficient trading between approximation and estimation error.

“Observed” documents

Inference



writes article don time apr

super-topic

god jesus christ people christian

faith wrong read spiritual passage

agree reason matter statement means

history support community house involved

key government encryption president clipper

sub-topic

agree reason matter statement means

power arms president home vote

history support community house involved

israel jews israeli jewish arab

history support community house involved

side left happened committee region

agree reason matter statement means

turkish armenian armenians people turkey

side left happened committee region

history support community house involved

hundred clothes tyre bosnians origin

file ftp windows window image

bit fax manager lib uk

site dec sources key public

release size function appreciated box

# Hierarchical Clustering

- ❑ Group terms into hierarchical clusters and each cluster represents an interested concept
- ❑ Top-down approaches:
  - ❑ CATHY [Wang et al.'13a]
  - ❑ CATHYHIN [Wang et al.'13b]
- ❑ Bottom-up approaches:
  - ❑ BRT [Liu et al.'12] [Song et al.' 15]

# Example: CATHY [Wang et al.'13a]

- ❑ Step 1: Construct term co-occurrence network
- ❑ Step 2: Cluster co-occurrence network into subtopic's sub-networks and estimate each sub-topical phrase's frequency
- ❑ Step 3: Extract candidate phrases using topical frequency
- ❑ Step 4: Rank topical phrases based on topical frequency
- ❑ Step 5: Apply steps 2-5 to each subtopic recursively and construct the hierarchy in a top-down fashion

# Example: BRT [Liu et al.'12]

- Agglomerative multi-branch clustering using Bayesian Rose Tree

---

**Algorithm 1** Bayesian Rose Tree (BRT).

---

**Input:** A set of documents  $\mathcal{D}$ .

$T_i \leftarrow \mathbf{x}_i$  for  $i = 1, 2, \dots, n$

$c \leftarrow n$

**while**  $c > 1$  **do**

1. Select  $T_i$  and  $T_j$  and merge them into  $T_m$  which maximizes

**Join:**  $T_m = \{T_i, T_j\}$

$$L(T_m) = \frac{p(\mathcal{D}_m | T_m)}{p(\mathcal{D}_i | T_i)p(\mathcal{D}_j | T_j)}, \quad \textbf{Absorb: } T_m = \{\text{children}(T_i) \cup T_j\}$$

where the merge operation is join, absorb, or collapse.

2. Replace  $T_i$  and  $T_j$  with  $T_m$  in the tree.

3.  $c \leftarrow c - 1$

**Collapse:**  $T_m = \{\text{children}(T_i) \cup \text{children}(T_j)\}$

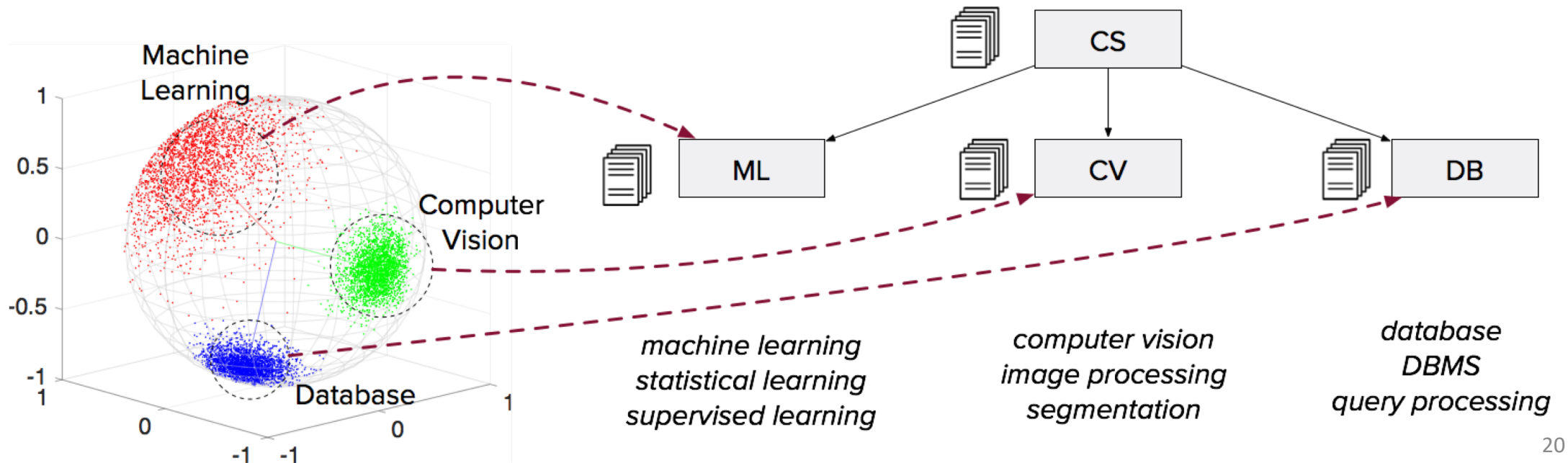
**end while**

# Limitations of Previous Methods

- ❑ Too strong assumptions on document generation process
  - ❑ Bag-of-word document representation ignores word order information
  - ❑ Real-world data may not follow these statistical distributions/processes
- ❑ Computationally slow
  - ❑ Slow inference restricts their applications to large-scale data

# Recent Methods: Uses Term Embedding

- Most existing work follows this idea: using term embedding to construct topic taxonomy based on hierarchical clustering
  - Learns term embedding to capture their semantic correlations
  - Constructs topic taxonomy in a recursive, top-down fashion



# Limitation of Term Embedding: Example






- ❑ Two terms in the Computer Science publications

**SIGKDD & Data Mining Researchers** — “*Frequent Pattern*” vs. “*Transaction Database*” — **SIGMOD & Database Researchers**

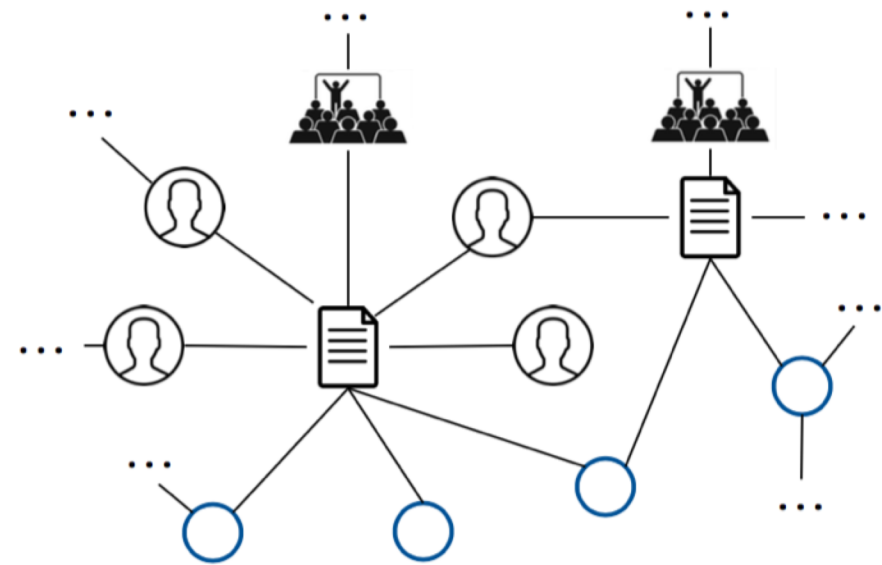
- ❑ From the taxonomy view
  - ❑ We should separate them into “*Data Mining*” and “*Database*”
- ❑ From the term embedding view
  - ❑ They are very similar due to similar contexts

# Text-Rich Network: Text & Meta-Data

- Terms are extracted by AutoPhrase from raw texts (e.g., paper & review)

Text Data	Typed Meta-Data		
docs (raw text)	 venues	 authors	 terms
	SIGKDD	C. Aggarwal Yan Li ...	freq. pattern uncertain data ...
	WWW	J. Leskovec J. Kleinberg ...	social networks info. cascade ...
...		...	

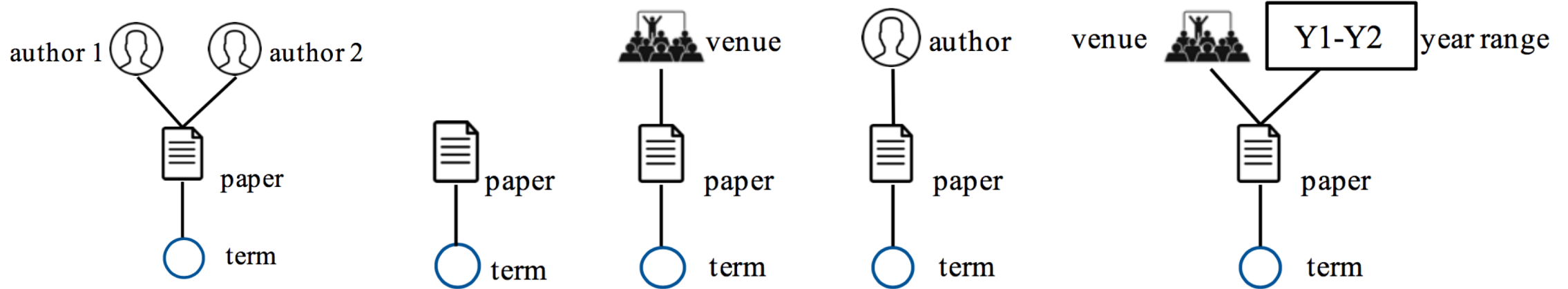
(a) An example digital collection of massive scientific papers.



(b) An text-rich network view of the example digital collection.

# Network Motifs: Contexts from Network

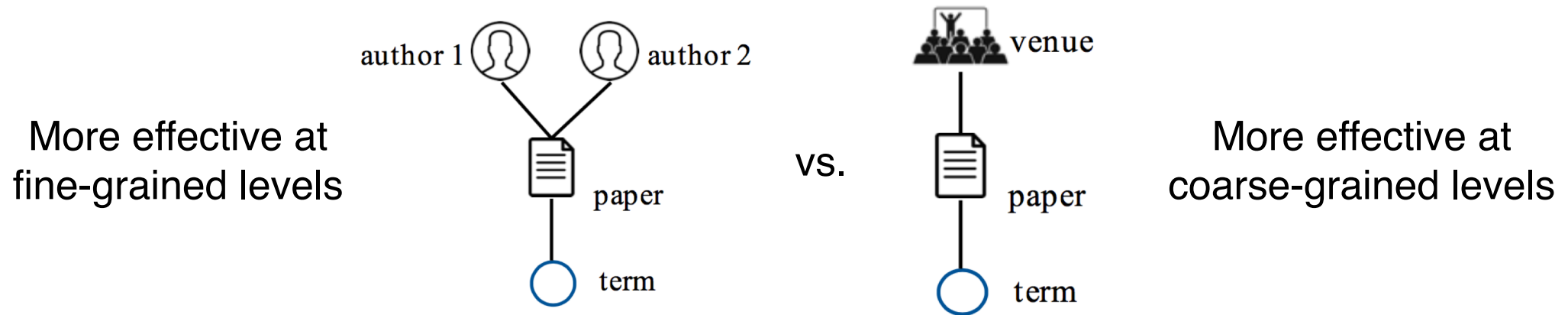
- Motif patterns capture subgraph contexts
  - Those nodes “connecting” two terms describes contexts
- Meta-path can be viewed as a special case of motif
  - T-P-T, T-P-V-P-T, T-P-A-P-T, ...



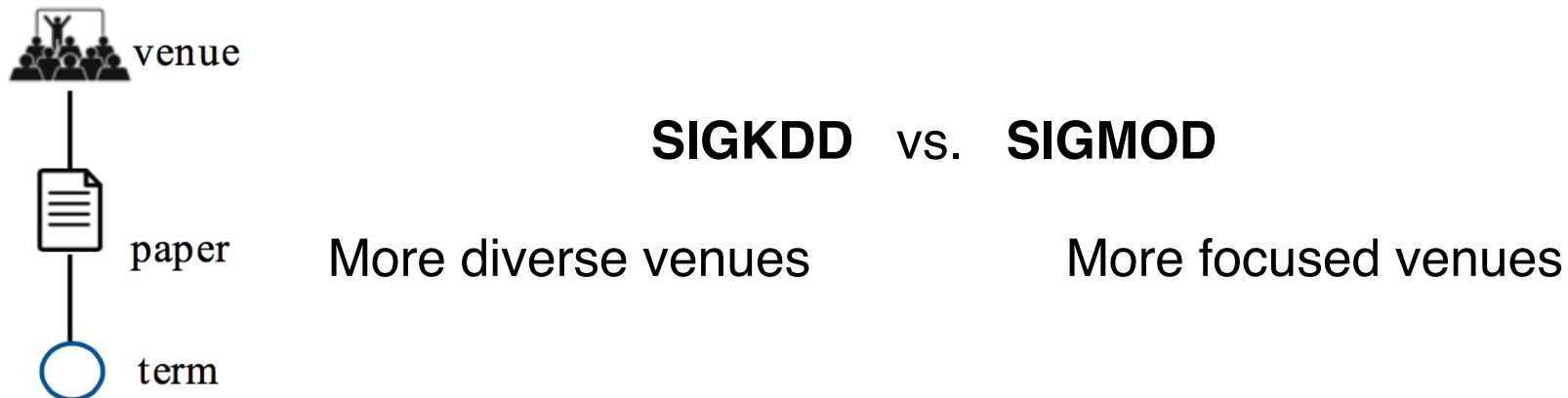
Motif Instance More Motif Patterns

# Text-Network Collaboration: Challenges

- ❑ Motif patterns are not created equally, even given by human experts

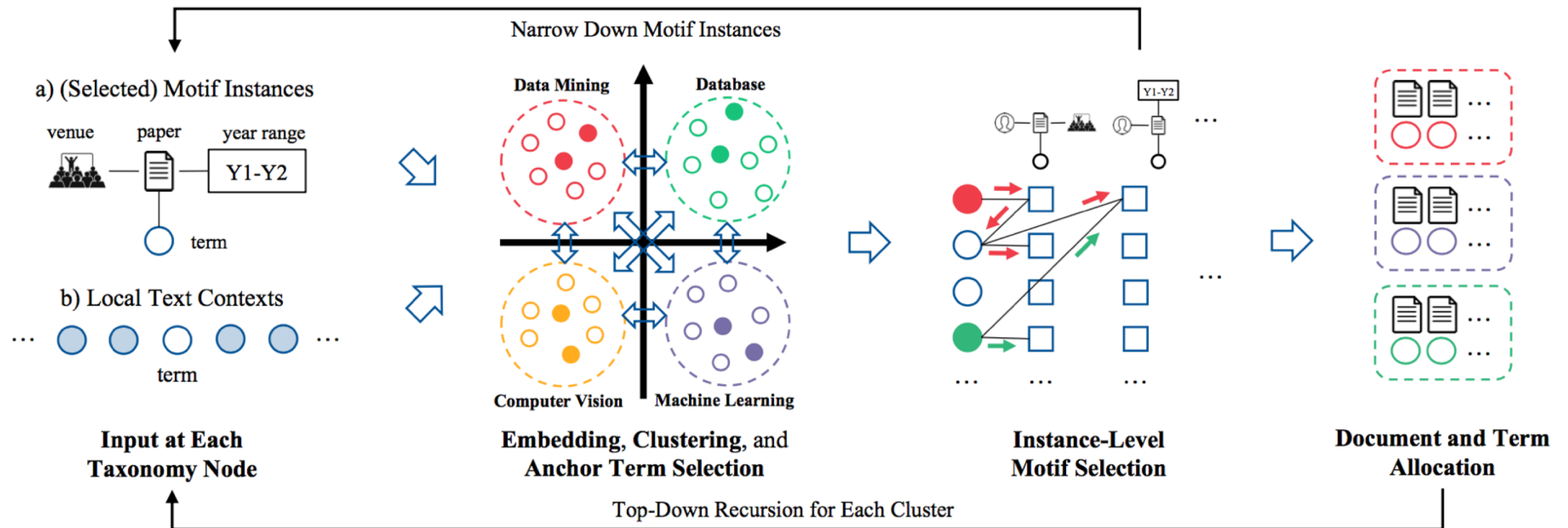


- ❑ Motif instances of the same motif patterns are not equally informative

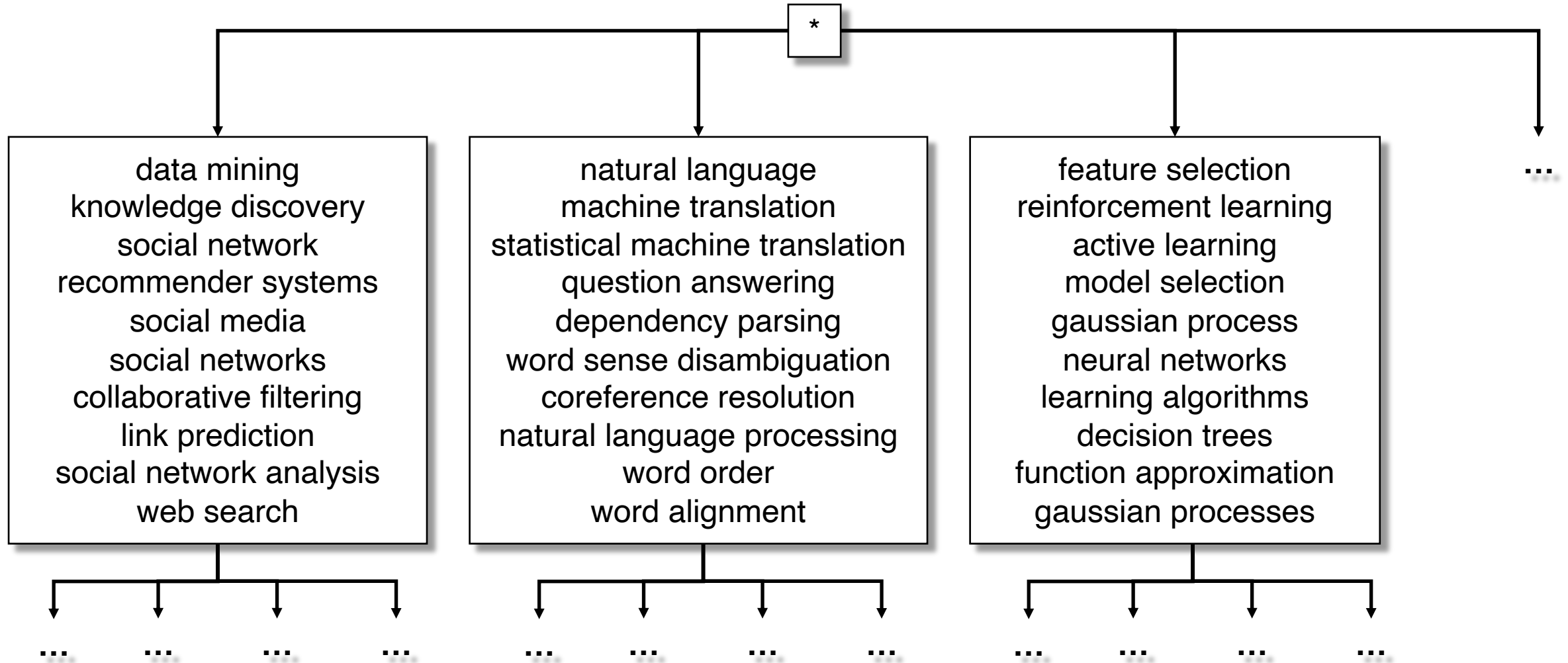


# NetTaxo: Instance-Level Motif Selection

- Starts from term embedding using textual contexts only
- Anchor terms makes the initial clustering results more robust
- Joint term embedding based on selected motif instances & textual contexts

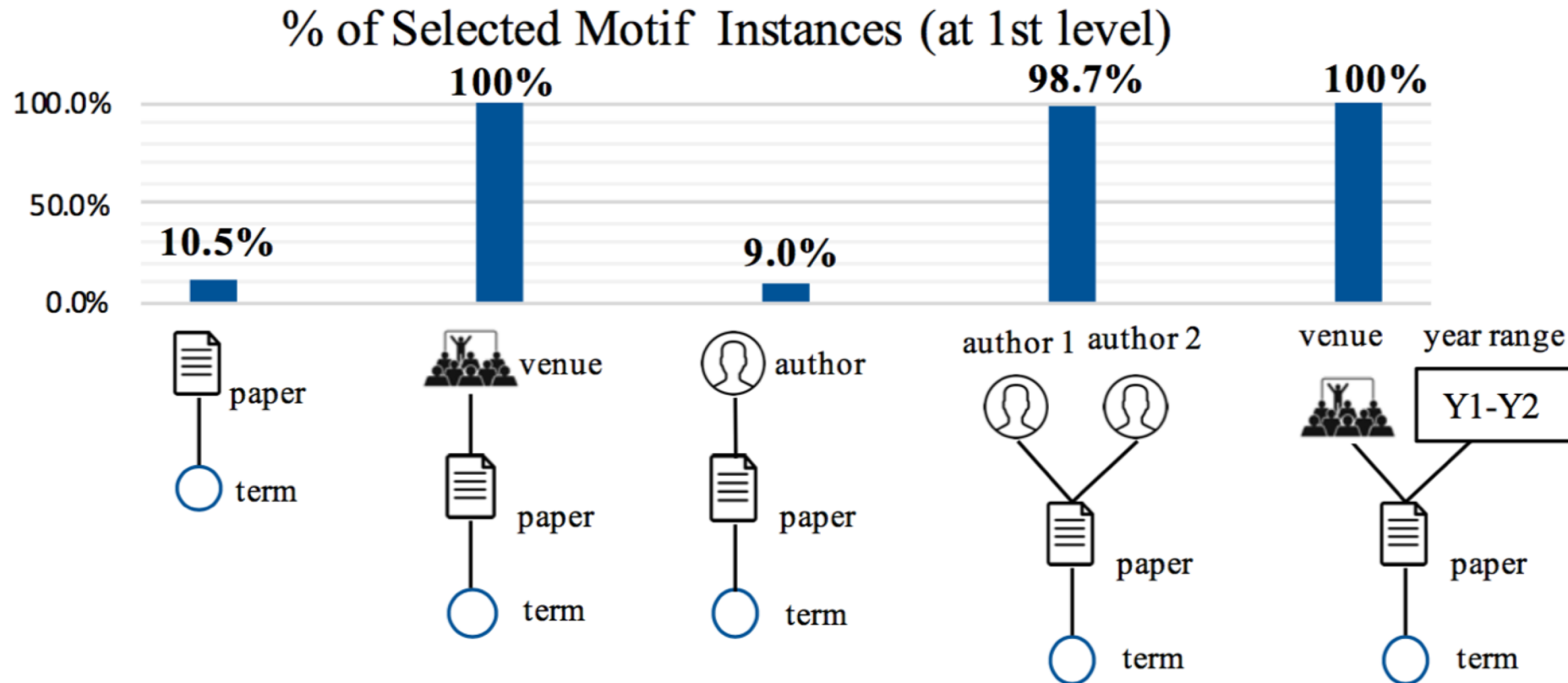


# NetTaxo: Joint Clustering Results (CS domain, 1<sup>st</sup> level)



# NetTaxo: Implicit Motif Pattern Selection

- Instance-level selection implicitly filters useless motif patterns too



# Case Study: Selected Motif Instances I

- ❑ Interesting motif instances selected at different levels
- ❑ Author pairs with more focused research topics are selected at the 2<sup>nd</sup> level

## Hua Wu - Zhanyi Liu

sentiment analysis  
semantic features  
semantic relations  
textual similarity  
sentiment words  
sentiment classification  
...

## Roland Kuhn - George F. Foster

source language  
bilingual corpora  
bilingual word  
machine translation  
statistical machine translation  
bleu score  
...

## Level 2: NLP -> Sub-Areas

Hua Wu - Zhanyi Liu  
Omar F. Zaidan - Chris Callison-Burch  
Boxing Chen - Roland Kuhn  
Hua Wu - Haifeng Wang  
Roland Kuhn - George F. Foster  
Yoan Gutiérrez - Andrés Montoyo  
John Makhoul - Richard M. Schwartz  
...

# Case Study: Selected Motif Instances II

- ❑ Interesting motif instances selected at different levels
- ❑ Recent, diverse & early, focused venues may help at the 2<sup>nd</sup> level

## CIKM 2010-2014

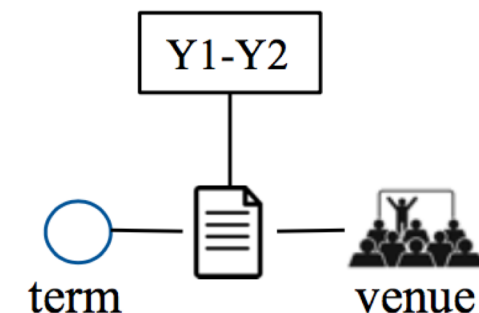
question answering  
information extraction  
language models  
sentiment analysis  
sentiment classification  
knowledge base

...

## Level 2: NLP -> Sub-Areas

ACL 1985-1989  
COLING 1985-1989  
EACL 1985-1989  
COLING 1980-1984  
ACL 1980-1984  
CIKM 2010-2014  
COLING 1990-1994

...



# NetTaxo: Evaluation Metrics

- ❑ **Coherence Measure**

- ❑ Are the terms at the same taxonomy node coherent?

- ❑ **Sibling Exclusiveness**

- ❑ Are the terms at a taxonomy node more similar compared to the terms at its sibling nodes?

- ❑ **Parent-Child Relationship**

- ❑ Are the relationships between the taxonomy nodes correct?

- ❑ All metrics are between 0 and 1; The bigger, the better.

# NetTaxo: Experimental Results

- Two domains: CS papers & Yelp reviews
- Compared with (“++” means “enhanced by our phrase mining results”)
  - TaxoGen (KDD’18) & HPAM++: Using text data only
  - CATHYHIN++ (ICDM’13, TKDE’18): Using network data only
  - HClusEmbed (NeurIPS’13, WWW’15): combines both term and node embeddings

	DBLP-5					Yelp-5				
	Coherence	Sibling	Parent-Child Relations			Coherence	Sibling	Parent-Child Relations		
	Measure	Exclusiveness	Precision	Recall	F <sub>1</sub>	Measure	Exclusiveness	Precision	Recall	F <sub>1</sub>
HPAM++	0.796	0.680	0.348	0.451	0.393	0.832	0.740	0.171	0.247	0.202
TaxoGen	0.840	0.740	0.780	0.713	0.745	0.920	0.800	0.650	0.618	0.633
CATHYHIN++	0.880	0.533	0.850	0.744	0.793	0.742	0.420	0.705	0.638	0.670
HClusEmbed	0.624	0.420	0.525	0.409	0.460	0.744	0.560	0.655	0.610	0.632
NetTaxo w/o Selection	0.908	0.680	0.895	0.808	0.849	0.816	0.540	0.668	0.681	0.674
NetTaxo	<b>0.912</b>	<b>0.880</b>	<b>0.898</b>	<b>0.810</b>	<b>0.852</b>	<b>0.928</b>	<b>0.854</b>	<b>0.790</b>	<b>0.825</b>	<b>0.807</b>