# Growing a Tree in the Forest: Constructing Folksonomies by Integrating Structured Metadata

**Anon Plangprasopchok (USC/ISI)**
**Kristina Lerman (USC/ISI)**
**Lise Getoor (UMD)**

- **Explosion of user-generated content**
  - *Images: Flickr, Picasa,..*
  - *Videos: YouTube, Vimeo,..*
  - *Maps: WikiMapia,..*
  - *Story: Blogs, Twitter,..*
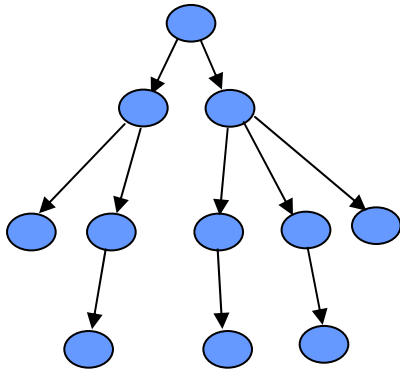  - *Relational Data: Metaweb, Google Base,..*


- **User-generated semantics: annotation/metadata**
  - *Tags, Geotags*
  - *Personal Hierarchies*


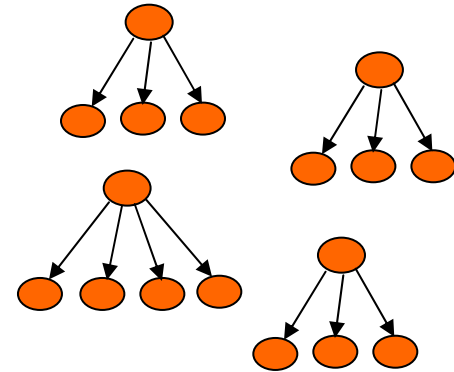*Goal: extract users' knowledge (folk knowledge) from this metadata*

# Folksonomy (communal taxonomy/hierarchy)

Folksonomy that users commonly have in their *mind (hidden)*

*Users select a portion of the hierarchy to organize their content.*

Personal hierarchies from various users (*observed*) such as users' folder-sub folders



• • •

[deep & bushy]

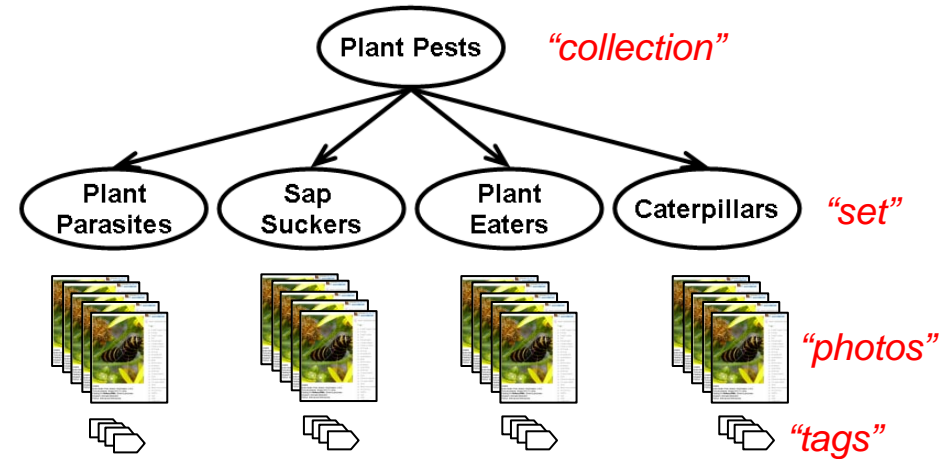[shallow, noisy, sparse(incomplete) & inconsistent]

*Can we recover the folksonomy back from many observed hierarchies? → folksonomy learning*

USC

- **Personal Hierarchies in Social Web**
- **Challenges in Folksonomy Learning**
- **Integrate Personal Hierarchies to Folksonomy**
  - *Relational Clustering for Learning Folksonomy*
- **Evaluations**
  - *Metrics*
  - *Results*
- **Related Work**
- **Discussion & Conclusions**

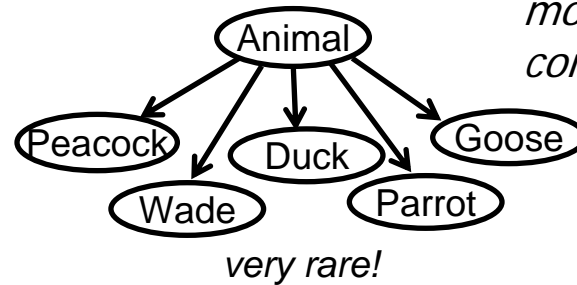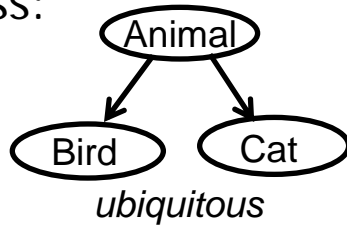**Shallow Hierarchy by user "maxi millipede"**

*"collection"*

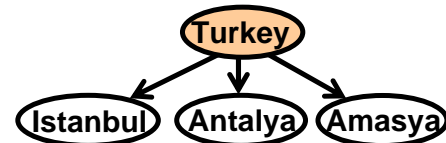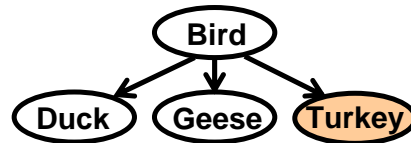*"set"*

*"photos"*

*"tags"*

**Tags on each photo**

Assume:
1) A set of tags on a set is an aggregation of all tags of all photos in the set
2) A set of tags on a collection is an aggregation of all tags of all sets in the collection
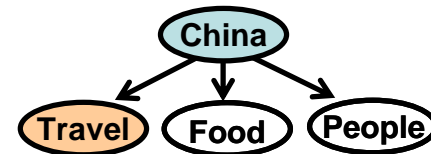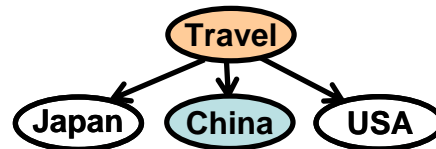
1.) Sparseness:

Animal → Bird, Cat

*ubiquitous*

Animal → Peacock, Duck, Goose, Wade, Parrot

*mostly personal hierarchies contain very few child nodes*

*very rare!*

2.) Ambiguity:

Bird → Duck, Geese, Turkey

Turkey → Istanbul, Antalya, Amasya

3.) Conflict:

Travel → Japan, China, USA

China → Travel, Food, People

4.) Varying Granularity:

UK → Scotland, London

UK → Glasgow, Edinburgh, London

Scotland → Glasgow, Shetland
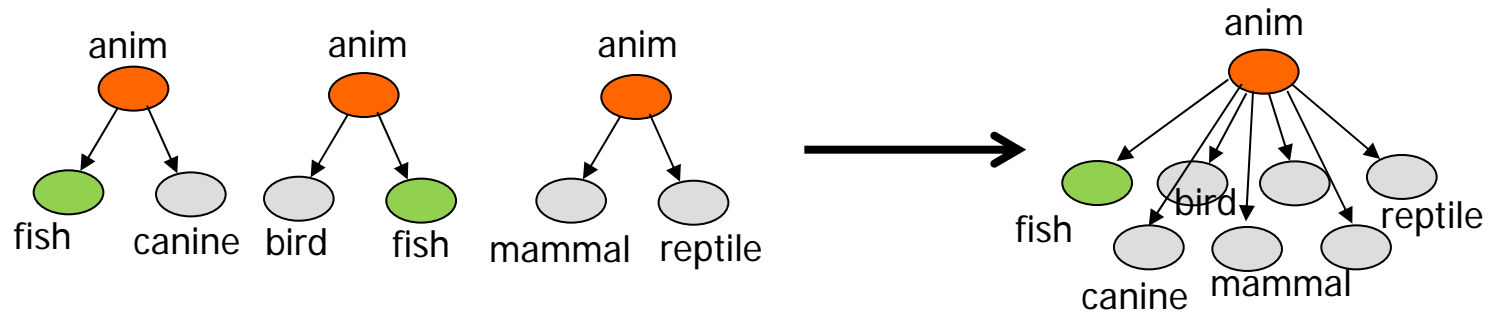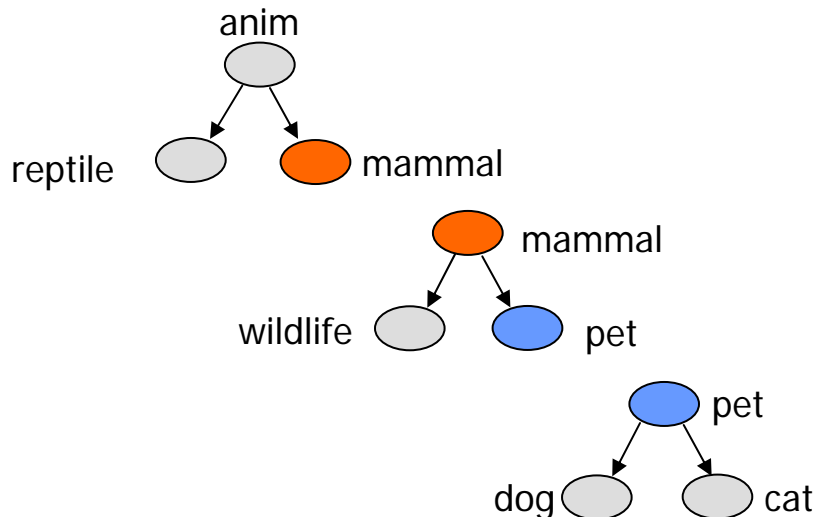
**Sketched idea:** *combine/aggregate personal hierarchies together in both <u>horizontal</u> and <u>vertical</u> directions.*
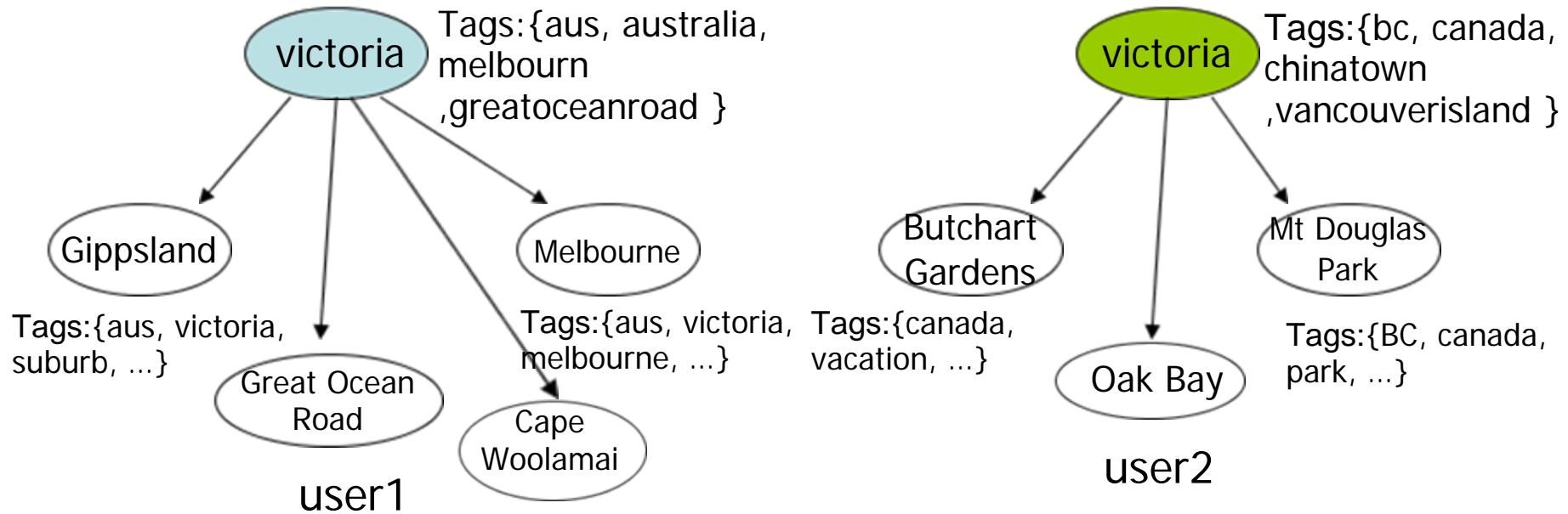


*Horizontal aggregation: expanding folksonomy's width*



*Vertical aggregation: extending folksonomy's depth*
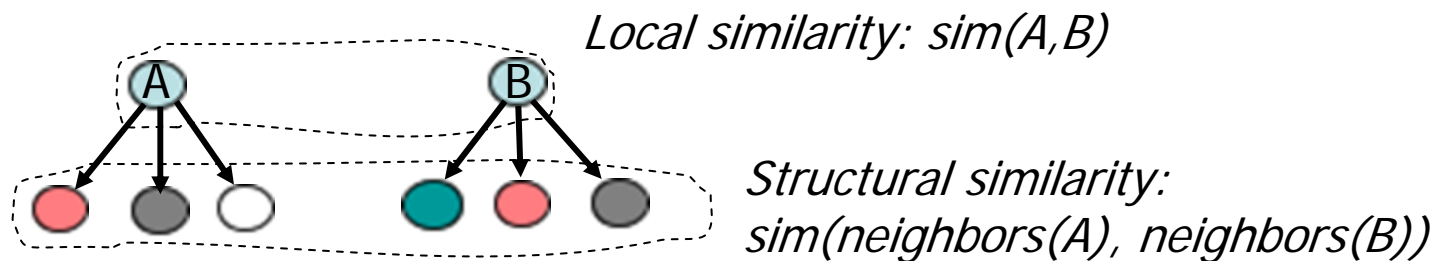
**Sketched idea**:  2 nodes should be clustered if they are <u>similar</u> enough
– similarity is computed using *contextual & relational information*



*Check common tags & child nodes*

USC

Formally, two nodes are considered similar if:

(1) their features are similar, i.e., have similar names, have many
    common tags – *local similarity*
(2) their neighbors are similar – *structural similarity*



*Local similarity: sim(A,B)*

*Structural similarity:*
*sim(neighbors(A), neighbors(B))*

$$Sim(A,B) = (1-\alpha)*localsim(A,B) + \alpha*structuralSim(A,B)$$

$\alpha$ is a weight
on how much we
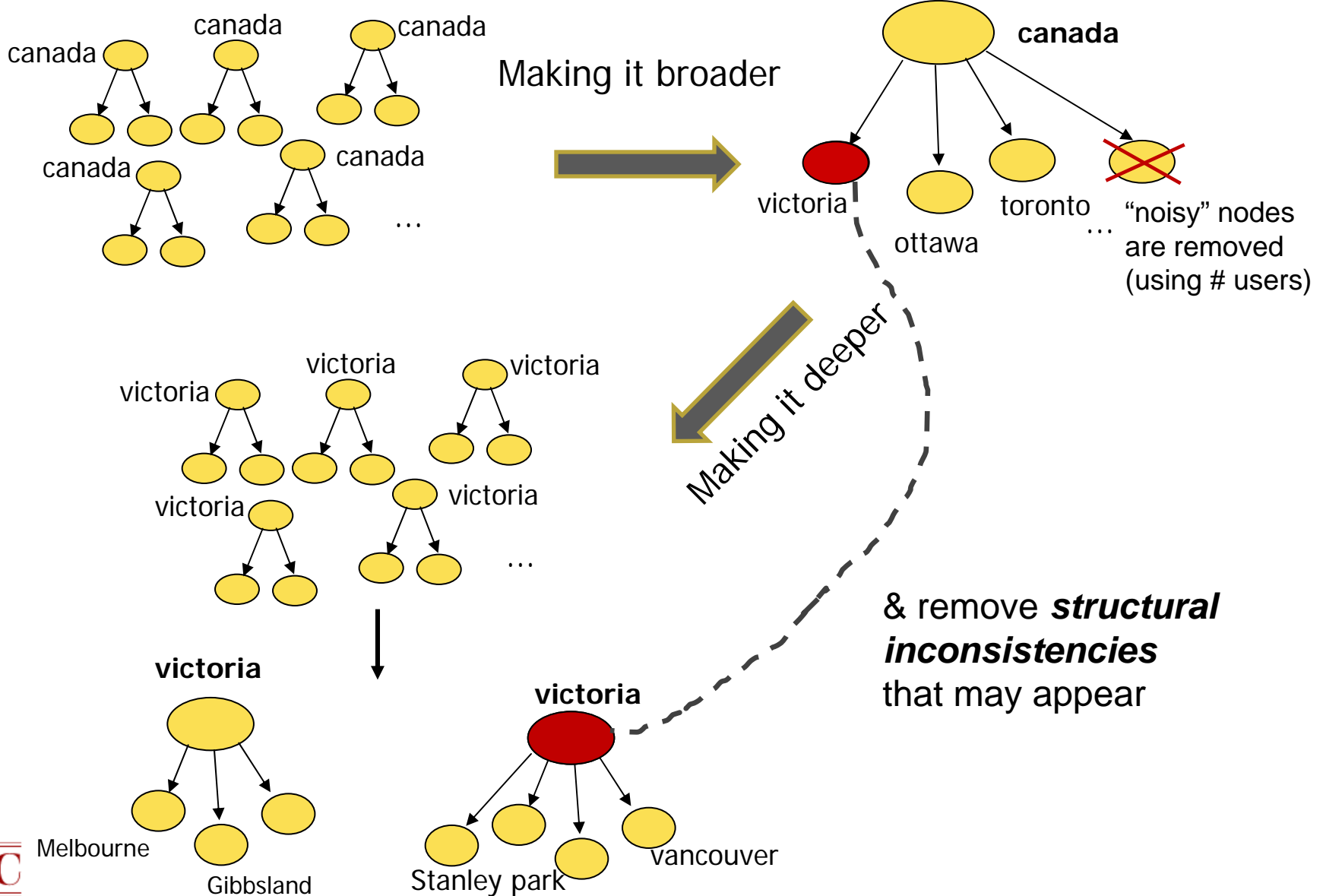rely on structural
information

We then merge nodes together if they are similar enough.

Note that we use naïve version of relational clustering by simply using neighbors' local
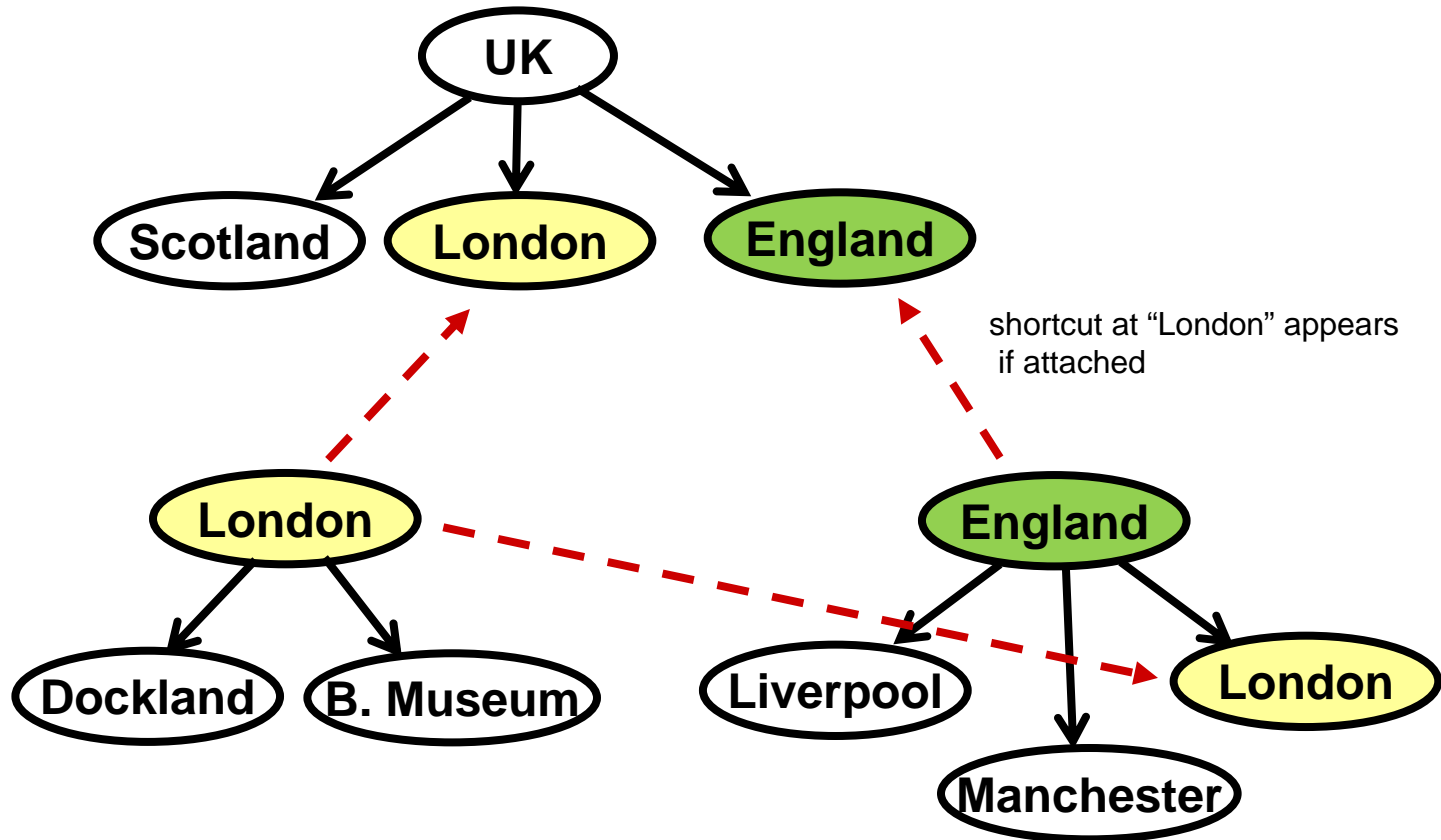features to compute structural similarity, instead of neighbors' class labels*.

*see Bhattacharya & Getoor, 2007, Collective Entity Resolution in Relational Data, TKDD for more detail

- **Pick a seed(root) term, e.g., "canada"**



Making it broader

Making it deeper

"noisy" nodes are removed (using # users)

& remove *structural inconsistencies* that may appear

Suppose we have the following clusters of hierarchies:



shortcut at "London" appears
if attached

- Shortcuts have to be removed to make the learned hierarchy consistent

- Keep the longer path since it captures more specific knowledge

- Growing trees from 32 seed terms & uses personal hierarchies from Flickr as in the previous work.*

**Evaluation Methodologies:**

1) *Against the reference hierarchy* (DMOZ)
2) *Structural evaluation*
3) *Manual evaluation*

**Baseline Approach***

- Assume nodes having the same name refer to the same concept
- Keep the relations between node pairs if they are not generated at random (using significance test)
- Then, combine all relations into a tree

* A. Plangprasopchok and K. Lerman, 2009, Constructing folksonomies from user-specified relations on flickr, WWW

## 1.) an automatic comparison to the reference hierarchy

-**Taxonomic Overlap** [adapted from Maedche & Staab] measuring <u>structure similarity</u> between two trees. For each node, determining how many ancestor and descendant nodes overlap to those in the reference tree.

-**Lexical Recall** measuring how well an approach can discover concepts, existing in the reference hierarchy (coverage)
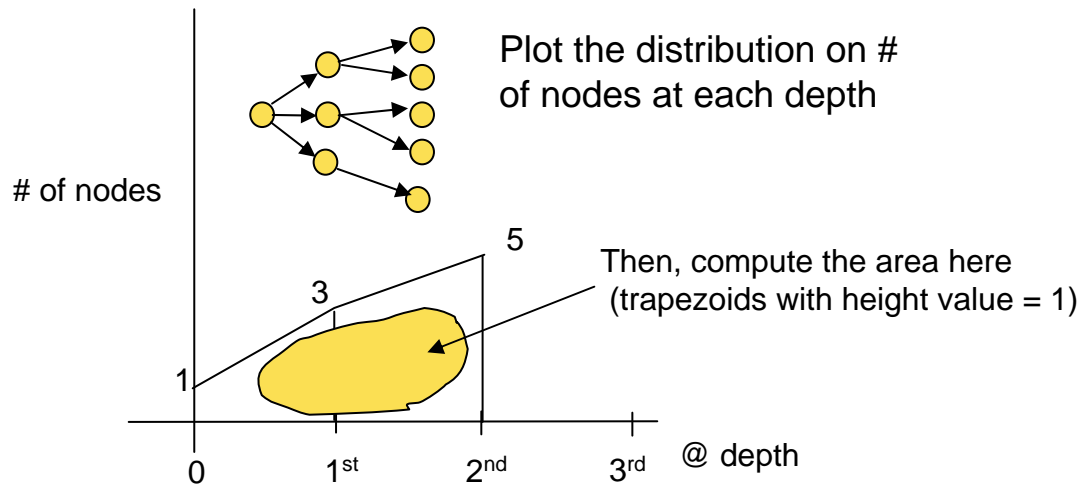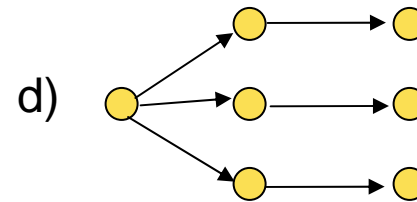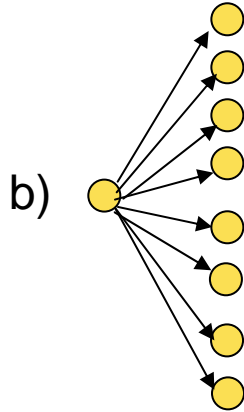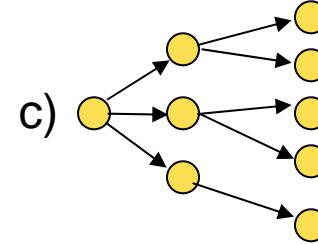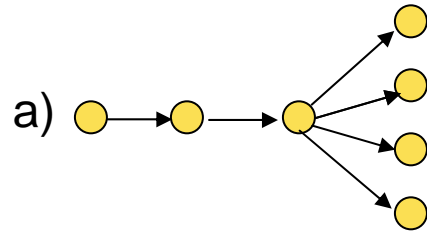
## 2.) Structural evaluation

-**Area Under Tree (AUT)** combining bushiness and depth of the tree into a single number: the higher value, the bushier and deeper tree.
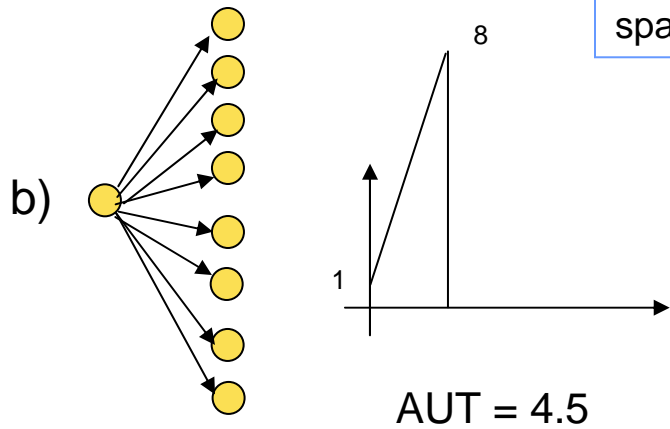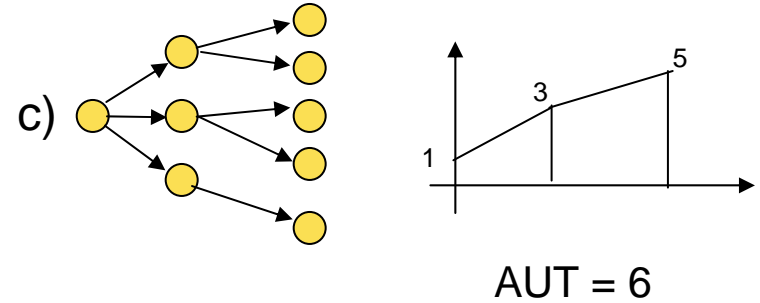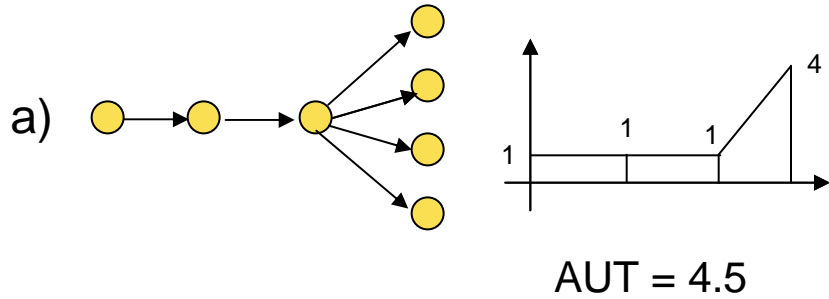
## 3.) Manual evaluation

- **Accuracy**: simply asking users whether a path from root to leaf of is correct: if there are some nodes misplaced in the wrong order, users will judge the whole path incorrect

*A. Maedche & S. Staab, 2002, Measuring Similarity between Ontologies, in EKAW

Which structures are the best in term of "bushiness" and "depth"?



Plot the distribution on # of nodes at each depth

Then, compute the area here (trapezoids with height value = 1)

# of nodes

5

3

1

0    1st    2nd    3rd    @ depth

USC

a)

AUT = 4.5

c)

AUT = 6

The largest area we can get is from the tree that keeps spanning at each level

b)

AUT = 4.5

d)

AUT = 5

Evaluate on 32 cases

| Metrics | # of cases that are superior to the other approach | |
| --- | --- | --- |
| | Baseline | The present work |
| Taxonimic Overlap | 7 | 15 |
| Lexical Recall | 6 | 19 |
| AUT | 3 | 18 |
| **TO+LR+AUT** | 0 | 11 |
| Accuracy (Manual) | 5 | 5 |

Terms are stemmed

- **Learning concept hierarchy from text data**
  - **Syntactic based [Hearst92, Caraballo99, Pasca04, Cimiano+05, Snow+06]**
  - **Word clustering [e.g., Segal+02, Blei+03]**

- **Induce concept hierarchy from tags**
  - **Graph-based & clustering based [Mika05, Brooks+06, Heymann+06, Zhou07+]**
  - **Probabilistic subsumption [Schmitz06]**

- **Ontology alignment [e.g., Udrea+07]**

- **Exploit user-specified hierarchy for recommendation**
  - **GiveALink [Markines06+]**

- The present work can create **more accurate** and **more detailed** folksonomies than the current state-of-the-art approach, since it exploits structural information during the merging process

- The present work is **more scalable**: incrementally growing the folksonomies rather than using on an exhaustive search

- Future work:
  - *Automatically separate broader/narrower from related-to relations (facets)*
  - *combining more sources of evidence such as geographical information*
  - *Apply on different data sets: e.g., personal workspaces, semantic network*