# Thematically Focused Search
# in Web 2.0 Folksonomies

Sergej Sizov, Olaf Görlitz

University of Koblenz-Landau

Germany

UNIVERSITÄT
KOBLENZ · LANDAU

can this huge mass of users give us new insights that would not be possible by considering individual contributions ?!

# The mass makes the difference?

# Flickr example – what you get ..

## Query: "Athos fire"



У Светој српској царској лаври Хиландару на Светој Гори Атонској, у ноћи између 3. и 4. марта 2004. године, избио је пожар великих размера.

# Flickr: queries with low recall



flickr LOVES YOU™

Sie sind nicht angemeldet     Anmelden     Hilfe

**Startseite     Die Tour     Registrieren     Entdecken** ▾                    **Suchen** ▾

## Suchen      **Fotos**   Gruppen   Personen

bamberg aula                          **SUCHEN**   Erweiterte Suche
                                                   Nach Kamera suchen

● Volltext        ○ Nur Tags

⊘ Wir konnten keine Ergebnisse zu  für **bamberg** und **aula** finden.

Möchten Sie stattdessen nach germany, deutschland, bavaria, bayern or franken suchen?

Sie          Anmelden | Kostenlosen Account einrichten

Entdecken    Places | Letzte 7 Tage | Dieser Monat | Beliebte Tags | Creative Commons | Suchen

Hilfe        Community-Richtlinien | Hilfeforum | FAQ | Sitemap | E-Mail-Hilfe

Flickr Blog | Über Flickr | Nutzungsbedingungen | Datenschutz | Copyright-Richtlinien | Missbrauch melden         a YAHOO! company

繁體中文 | Deutsch | English | Español | Français | 한글 | Italiano | Português

Copyright © 2008 Yahoo! Alle Rechte vorbehalten.

# Outline

Problem formalization

Thematically focused search and ranking

Distributed setting – pro & contra

Evaluation

# Formalizing the problem

Collaborative content sharing framework:

users $u \in U$   tags $t \in T$   resources $r \in R$

Information cloud:   $T := (Y^*, f), f(t) : Y^* \subseteq Y, Y^* \to [0..1]$

- user-centric:   $T_u := (Y_u, f), Y_u \subseteq u \times T \times R,$
- resource-centric:   $T_r := (Y_r, f), Y_u \subseteq U \times T \times r,$
- community-specific: $T_{U^*} := (Y_{U^*}, f),$   $Y_{U^*} \subseteq U^* \times T \times R$
- collection-specific   $T_{R^*} := (Y_{R^*}, f),$   $Y_{R^*} \subseteq U \times T \times R^*$
- arbitrary   $T_{U^*R^*} := (Y_{U^*R^*}, f),$   $Y_{U^*R^*} \subseteq U^* \times T \times R^*$

   .. e.g. obtained by traversing the hypergraph up to certain depth

Common recommender scenarios:

- ◆ Given a user, recommend photos which may be of interest.
- ◆ Given a user, recommend users they may like to contact.
- ◆ Given a user, recommend groups they may like to join.

# The IR background – constructing feature vectors



$$if(i) = (a_i, b_i)$$
$$iif(i) = \left(\log \frac{|J|}{|J^*|}, \log \frac{|K|}{|K^*|}\right),$$
$$j \in J^*/k \in K^*:$$
$$(i, j, k) \in Y^*$$
$$weight(i) = iif(i) \cdot if(i)^T$$

.. defined analogously to tf·idf

## tag features:

$$iif \cdot if(p2p) = (\log \frac{5}{2}, \log \frac{5}{5}) \cdot (6, 6)^T$$

$$iif \cdot if(talk) = (\log \frac{5}{1}, \log \frac{5}{1}) \cdot (1, 1)^T$$

$$iif \cdot if(slide) = (\log \frac{5}{1}, \log \frac{5}{1}) \cdot (1, 1)^T$$

$$iif \cdot if(pdf) = (\log \frac{5}{1}, \log \frac{5}{1}) \cdot (1, 1)^T$$

$$iif \cdot if(routing) = (\log \frac{5}{5}, \log \frac{5}{1}) \cdot (6, 6)^T$$

Further dimensions of interest:

- favorites
- group membership
- contact lists
- comments on other's resources

# Generalization: the tensor model

idea: using multi-dimensional arrays for representing relationships



user-centric clouds

resource-centric clouds

tag-centric clouds

matricize:

$$X_{(users)} = \boxed{\text{community cloud}}$$

$$X_{(resources)} = \boxed{\text{cloud for resources}}$$

$$X_{(tags)} = \boxed{\text{cloud for tags}}$$

slide shows 3rd order tensor for common Web 2.0 dimensions, but can (and should) be extended by other relationships (favorites, comments, groups..)

$$Y = X \times_1 A$$
$$Y_{(::k)} = X_{(::k)} \times A^T$$

$$Y = X \times_3 B$$
$$Y_{(i::)} = X_{(i::)} \times C^T$$

$$Y = X \times_2 B$$
$$Y_{(:j:)} = X_{(:j:)} \times B^T$$

general idea: decompose tensor in order to identify significant "factors" along each dimension  (multi-dimensional methods analogously to LSI, PCA)

our current approach: input-tensor R is decomposed into R=UxDxV, V contains the orthogonal mapping of R into space of real tensors with target dimension (i.e. V is also a matrix). Restrict the resulting feature vectors to 5-10 most significant dimensions, analogously to LSI.

## User-centered focusing

1. Compute characteristic feature vectors for resources, tags, contacts, or favorites of the given user

2. Construct appropriate decision model (centroid, naive bayes, SVM, etc.)

3. Explore the tagging cloud around the user, order matches acording wrt estimated utility function (cosine similarity, classification confidence, etc.)

4. Return the top-k result set (e.g. top-10, top-20) to the user

## Datasets

- Flickr dataset (2004-2005)
  - *319,686 users,*
  - *1,607,879 tags,*
  - *28,153,045 resources,*
  - *112,900,000 tag assignm.*

- Del.icio.us dataset (2003-2006)
  - *532,924 users,*
  - *2,481,698 tags,*
  - *17,262,480 resources,*
  - *140,126,586 tag assignm.*

Evaluation: apriori method

- remove a certain fraction of relationships (e.g. group participation, comments, ..) from the test cloud

- test the ability of the recommender to reconstruct missing relationships (i.e. to place them within top-k of the result set)

# Results: user-focused recommendations

## recommending favorites

| User representation | Training:10 prec@10 | Training:10 prec@20 |
|---|---|---|
| *Random* | 0.167 | 0.167 |
| *User items* | 0.259 | 0.268 |
| *Commented items* | 0.236 | 0.221 |
| *Favorites* | 0.872 | 0.727 |
| *Combined* | 0.854 | 0.713 |
| | Training:20 prec@10 | Training:20 prec@20 |
| *Random* | 0.167 | 0.167 |
| *Commented items* | 0.255 | 0.248 |
| *Favorites* | 0.918 | 0.851 |
| *Combined* | 0.899 | 0.828 |
| | Training:40 prec@10 | Training:40 prec@20 |
| *Random* | 0.167 | 0.167 |
| *Commented items* | 0.265 | 0.266 |
| *Favorites* | 0.933 | 0.903 |
| *Combined* | 0.914 | 0.876 |

## recommending contacts

| User representation | prec@5 | prec@10 |
|---|---|---|
| *Random* | 0.167 | 0.167 |
| *User Items* | 0.574 | 0,472 |
| *Commented Items* | 0.576 | 0.473 |
| *Favorites* | 0.535 | 0.455 |
| *Contacts (training 10)* | 0.604 | 0.497 |
| *Contacts (training 20)* | 0.611 | 0.498 |

tensor based recommendation: consistently better accuracy
in preliminary experiments, now under evaluation

# Decentralized setting: pro & contra argumentation

☹ multiple accounts for different resource types

☹ space limitations
(e.g. max 200 photos in Flickr)

☹ censorship, rank manipulations

☹ single point of failure

## Distributed Tagging System ?

◆ tag any kind of personal data on the Desktop
◆ share and browse tagged data in a P2P network

Our implementation: Tagster
open source, available at http://isweb.uni-koblenz.de

given: set of methods $V = \{v_1,...,v_k\}$, confidence grades $res(v_i,d)$ for document $d$

Meta result (restrictivity by thresholds $t_1$ and $t_2$, *tuning by weights* $w(v_i)$ ):

$$Meta(d) = \begin{cases} +1 & \text{if } \sum_i res_i(d) \cdot w(v_i) > t_1 \\ -1 & \text{if } \sum_i res_i(d) \cdot w(v_i) < t_2 \\ 0 & \text{otherwise} \end{cases}$$

Special cases:
- "Unanimous Decision"
- "Voting"
- "Weighted Average" (e.g., weighted by some quality estimator)

# Collaborative organization of document collections

**Given**: set of methods $V=\{v_1,...,v_L\}$, „unanimous decision"

$$X_i = \begin{cases} 1 & \text{if } v_i \text{ assigns document correctly} \\ 0 & \text{otherwise} \end{cases}$$

$$P(X_1 = 1, ..., X_L = 1) = P(X_1 = 1) \cdot \prod_{i=1}^{L-1} \frac{P(X_i = 1)P(X_{i+1} = 1) + cov(X_i, X_{i+1})}{P(X_i = 1)}$$

$$error(meta) = P(X_1 = 0, ..., X_L = 0 | X_1 = .. = X_L)$$

$$loss(meta) = 1 - P(X_1 = ... = X_L)$$



$P(X_1{=}1), P(X_2{=}1), P(X_3{=}1), cov$
$loss(meta), error(meta), L$

$\mathbf{V_2}\ V_3\ V_1$
*meta* $V_1 .. V_6$

$\mathbf{V_1} V_2 V_3$ *meta*

$\mathbf{V_3}\ V_2\ V_1$
*meta* $V_1 .. V_6$

$\mathbf{V_4}$ *meta* $V_1 .. V_6$    $\mathbf{V_5}$ *meta* $V_1 .. V_6$    $\mathbf{V_6}$ *meta* $V_1 .. V_6$

# Decentralized Collaboration: Results



Accuracy: del.icio.us, Junk=1/2

- 1 peer
- 2 peers
- 4 peers
- 8 peers
- 16 peers

Junk Reduction and Document Loss: del.icio.us, Junk=1/2

- 1 peer (docs)
- 16 peers (docs)
- 1 peer (junk)
- 16 peers (junk)

# Distributed scenario: an example

bamberg
iuf: log(5/2) = **0.4**

BOB

| bamberg | 65·**0.4** |
|---------|-----------|
| koblenz | 114·0.1 |
| belin | 49·0.3 |
| cologne | 68·0.2 |

local ◄    ► global

TOM

| bamberg | 37·**0.4** |
|---------|-----------|
| berlin | 71·0.3 |
| paris | 54·0.1 |
| usa | 69·0.5 |

bamberg
iuf: **0.4 → 0.22**

$\delta = 0.98$

BOB

update?

local    global

$$sim(v_{BOB}^{\star}, v_{BOB}) = sim \left( \left[ \begin{array}{c|c} 65 \cdot 0.4 & 65 \cdot 0.22 \\ 114 \cdot 0.1 & 114 \cdot 0.1 \\ 49 \cdot 0.3 & 49 \cdot 0.3 \\ 68 \cdot 0.2 & 68 \cdot 0.2 \end{array} \right] \right) = 0.958$$

**update required!**

# The PINTS approach

- Index peers monitor feature vector accuracy for their tags
- Compare feature approximation with the tag's true *iuf* value

BOB's approximation:

$$v^{\star}_{BOB}(\theta) = \begin{vmatrix} tf(t_1) \cdot (a_{t_1} \cdot \theta + b_{t_1}) \\ \vdots \\ tf(t_m) \cdot (a_{t_m} \cdot \theta + b_{t_m}) \\ \vdots \\ tf(t_N) \cdot (a_{t_N} \cdot \theta + b_{t_N}) \end{vmatrix}$$

index peer's view:

$$v^{\circ}_{BOB,t_m}(\theta) = \begin{vmatrix} tf(t_1) \cdot (a_{t_1} \cdot \theta + b_{t_1}) \\ \vdots \\ \boxed{tf(t_m) \cdot iuf^{true}_{t_m}} \\ \vdots \\ tf(t_N) \cdot (a_{t_N} \cdot \theta + b_{t_N}) \end{vmatrix}$$

- Index peer needs to know the other approximations
- Vector similarity must be above threshold $\delta$

$$sim(v^{\star}, v^{\circ}_{t_m}) > \delta \qquad sim(v^{\star}, v^{\circ}_{t_m}) = \frac{v^{\star} \cdot v^{\circ}_{t_m}}{\|v^{\star}\| \, \|v^{\circ}_{t_m}\|}$$

$$sim(v^\star, v^\circ_{t_m}) = \frac{v^\star \cdot v^\circ_{t_m}}{\|v^\star\| \, \|v^\circ_{t_m}\|}$$

$$v^\star \cdot v^\circ_{t_m} = \sum_{t_i \neq t_m} \left( tf(t_i)^2 \cdot (a_{t_i} \cdot \theta + b_{t_i})^2 \right) + tf(t_m)^2 \cdot (a_{t_m} \cdot \theta + b_{t_m}) \cdot iuf^{true}_{t_m}$$

$$\|v^\star\| = \sqrt{\sum_{t_i \neq t_m} \left( tf(t_i)^2 \cdot (a_{t_i} \cdot \theta + b_{t_i})^2 \right) + tf(t_m)^2 \cdot (a_{t_m} \cdot \theta + b_{t_m})^2}$$

$$\|v^\circ_{t_m}\| = \sqrt{\sum_{t_i \neq t_m} \left( tf(t_i)^2 \cdot (a_{t_i} \cdot \theta + b_{t_i})^2 \right) + tf(t_m)^2 \cdot (iuf^{true}_{t_m})^2}$$

$$A_{t_m} = \sum_{t_i / t_m} \left( tf(t_i)^2 \cdot a_{t_i}^2 \right) \qquad B_{t_m} = \sum_{t_i / t_m} \left( tf(t_i)^2 \cdot a_{t_i} \cdot b_{t_i} \right) \qquad C_{t_m} = \sum_{t_i / t_m} \left( tf(t_i)^2 \cdot b_{t_i}^2 \right)$$

$$v^\star \cdot v^\circ_{t_m} = A_{t_m} \theta^2 + 2 B_{t_m} \theta + C_{t_m} + tf(t_m)^2 \cdot (a_{t_m} \cdot \theta + b_{t_m}) \cdot iuf^{true}_{t_m}$$

$$\|v^\star\| = \sqrt{A_{t_m} \theta^2 + 2 B_{t_m} \theta + C_{t_m} + tf(t_m)^2 \cdot (a_{t_m} \cdot \theta + b_{t_m})^2}$$

$$\|v^\circ_{t_m}\| = \sqrt{A_{t_m} \theta^2 + 2 B_{t_m} \theta + C_{t_m} + tf(t_m)^2 \cdot (iuf^{true}_{t_m})^2}$$

# PINTS: update strategy

$$sim(v^{\star}, v_{t_m}^{\circ}) = \frac{v^{\star} \cdot v_{t_m}^{\circ}}{\|v^{\star}\| \, \|v_{t_m}^{\circ}\|}$$
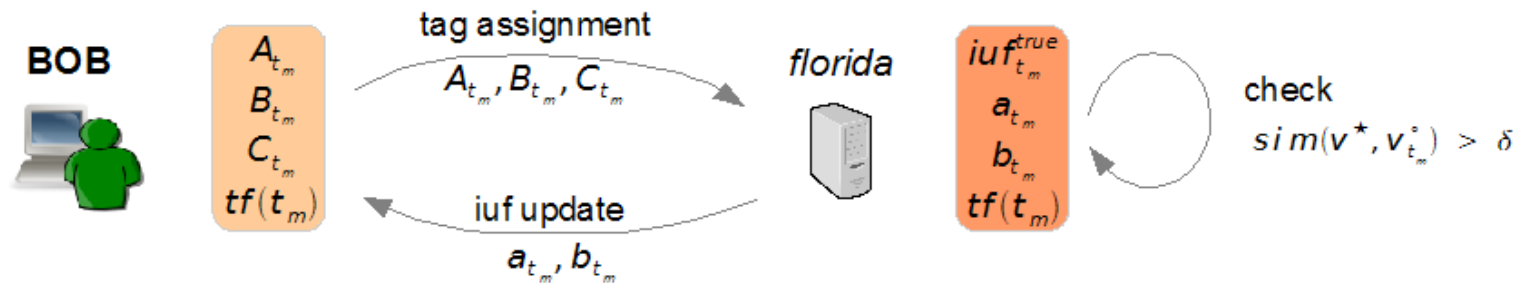


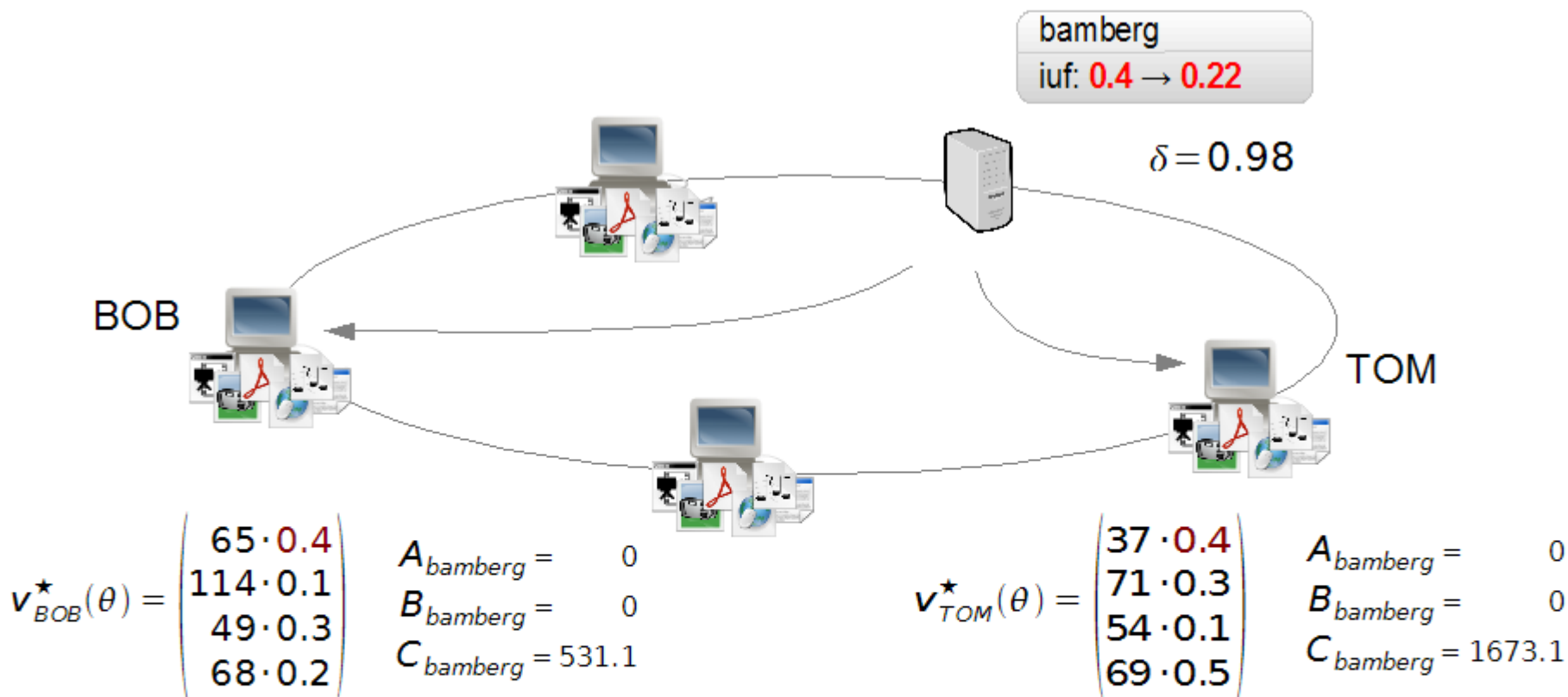$$v^{\star} \cdot v_{t_m}^{\circ} = A_{t_m}\theta^2 + 2B_{t_m}\theta + C_{t_m} + tf(t_m)^2 \cdot (a_{t_m} \cdot \theta + b_{t_m}) \cdot iuf_{t_m}^{true}$$

$$\|v^{\star}\| = \sqrt{A_{t_m}\theta^2 + 2B_{t_m}\theta + C_{t_m} + tf(t_m)^2 \cdot (a_{t_m} \cdot \theta + b_{t_m})^2}$$

$$\|v_{t_m}^{\circ}\| = \sqrt{A_{t_m}\theta^2 + 2B_{t_m}\theta + C_{t_m} + tf(t_m)^2 \cdot (iuf_{t_m}^{true})^2}$$

# PINTS updates: Beispiel



bamberg
iuf: **0.4 → 0.22**

$\delta = 0.98$

BOB

TOM

$$v^{\star}_{BOB}(\theta) = \begin{pmatrix} 65 \cdot 0.4 \\ 114 \cdot 0.1 \\ 49 \cdot 0.3 \\ 68 \cdot 0.2 \end{pmatrix} \qquad \begin{matrix} A_{bamberg} = & 0 \\ B_{bamberg} = & 0 \\ C_{bamberg} = 531.1 \end{matrix}$$

$$v^{\star}_{TOM}(\theta) = \begin{pmatrix} 37 \cdot 0.4 \\ 71 \cdot 0.3 \\ 54 \cdot 0.1 \\ 69 \cdot 0.5 \end{pmatrix} \qquad \begin{matrix} A_{bamberg} = & 0 \\ B_{bamberg} = & 0 \\ C_{bamberg} = 1673.1 \end{matrix}$$

$sim(v^{\star}_{BOB}, v^{\circ}_{BOB,bamberg}) = 0.958$    **update required**

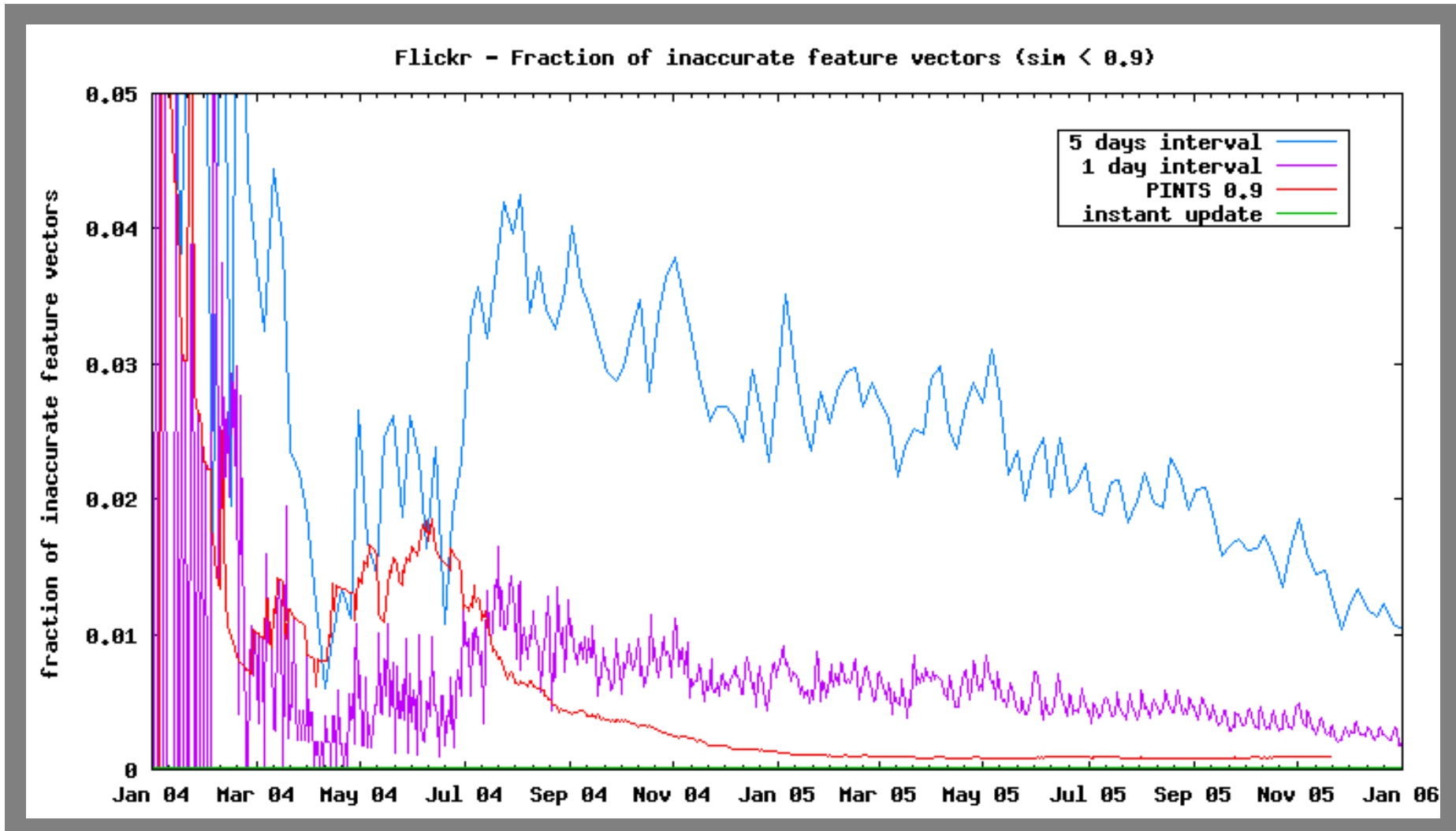$sim(v^{\star}_{TOM}, v^{\circ}_{TOM,bamberg}) = 0.989$    **no update required**

## Objectives

- check if (and how frequent) specified thresholds violated
- compare message complexity for various methods

## Methodology

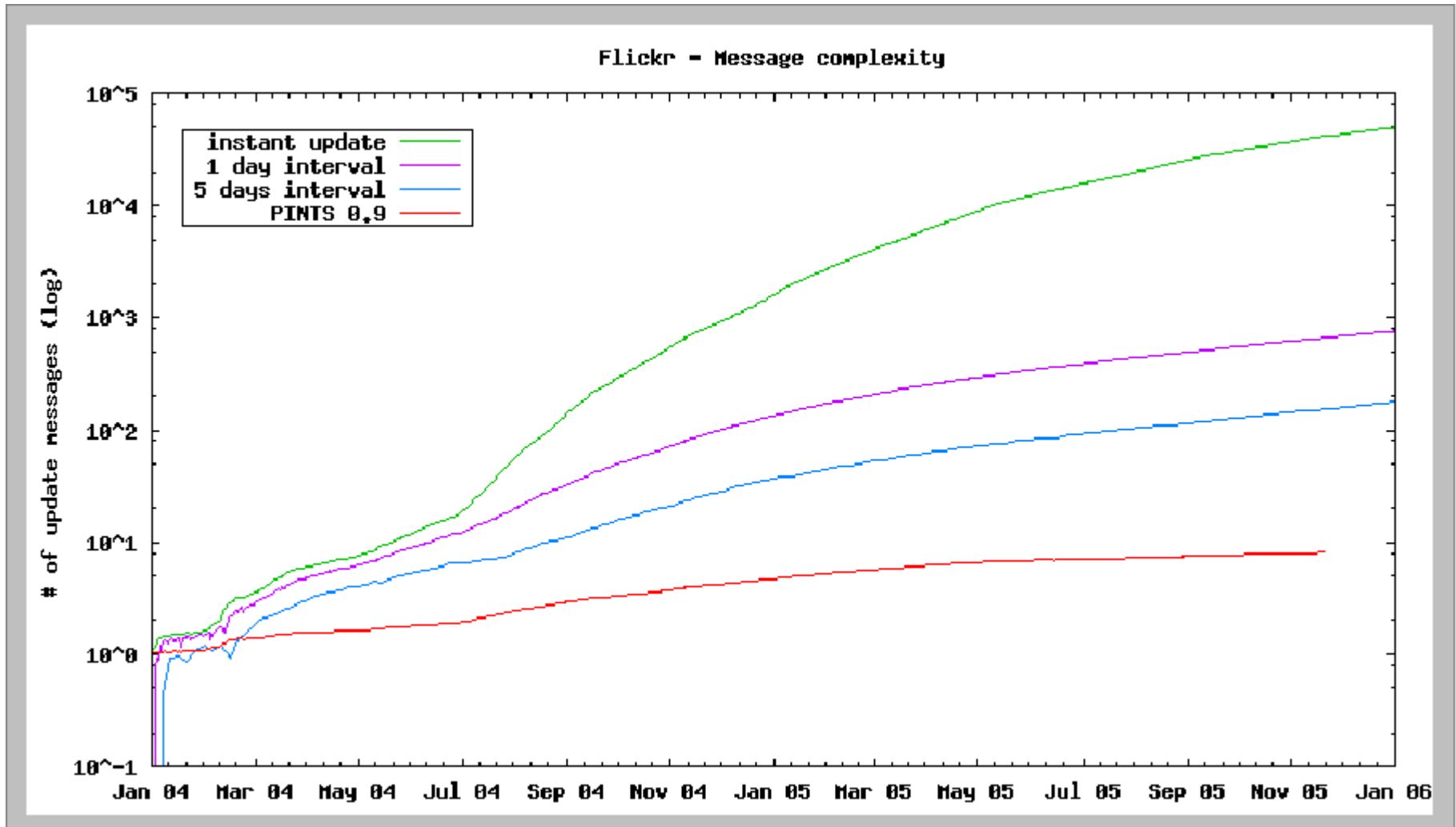- use real world tagging traces (time-ordered tas assignm.)
  - flickr.com  (~320k users, ~1.6m tags, ~28.2m resources)
  - del.icio.us (~533k users, ~2.3m tags, ~17.3m resources)
- replay tagging traces in P2P simulation
- measure inaccurate feature vectors, message complexity
- evaluate against interval-based update

Flickr – Fraction of inaccurate feature vectors (sim < 0.9)

PINTS: higher accuracy than interval updates

Flickr - Message complexity

PINTS: high accuracy at low message complexity

## Conclusions

Focused search & recommendation in Web 2.0 folksonomies:

- ◆ IR-like problem formalization
- ◆ Personal and social aspects/dimensions are important
- ◆ Multi-dimensional setting helps to improve accuracy
- ◆ Can be realized for centralized and decentralized architectures

## Future work

- ◆ bridging the semantic gap between low-level and high-level features
- ◆ decentralized computations on large sparse matrices (e.g P2P based PageRank or HITS estimation)
- ◆ evaluation methodology for Web 2.0 applications
- ◆ better understanding of Web 2.0 evolution patterns

thank you