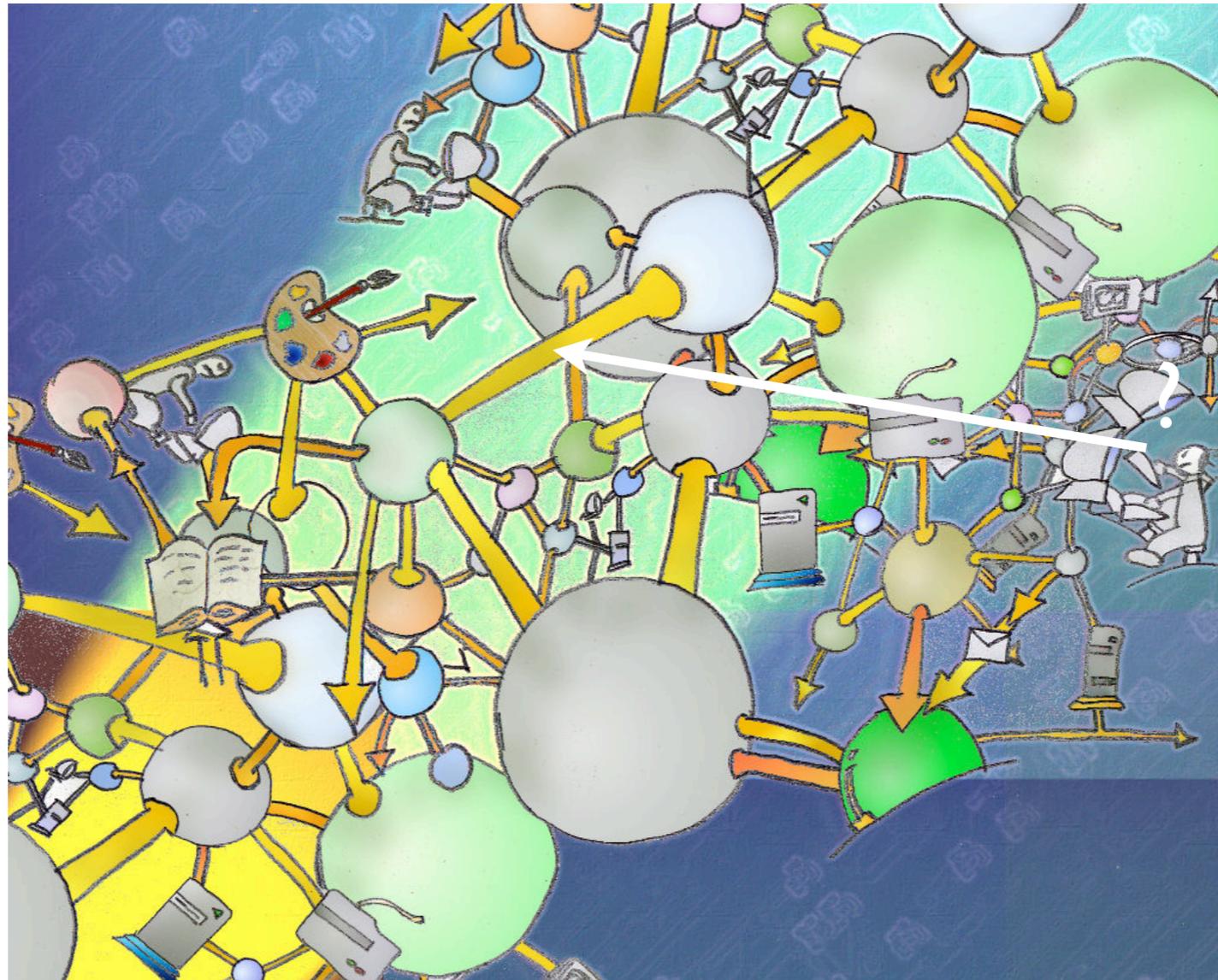


# Learning in hyperlinked environments

“There are finer fish in the sea that have ever been caught,” Irish proverb



*Marco Gori*  
University of Siena (Italy)

ECIR-2007 Rome, 5 April 2007

# Outline

- Beyond link analysis
- Diffusion learning
- Impact on information retrieval
- Conclusions

# **Beyond Link Analysis**

ECIR-2007 Rome, 5 April 2007

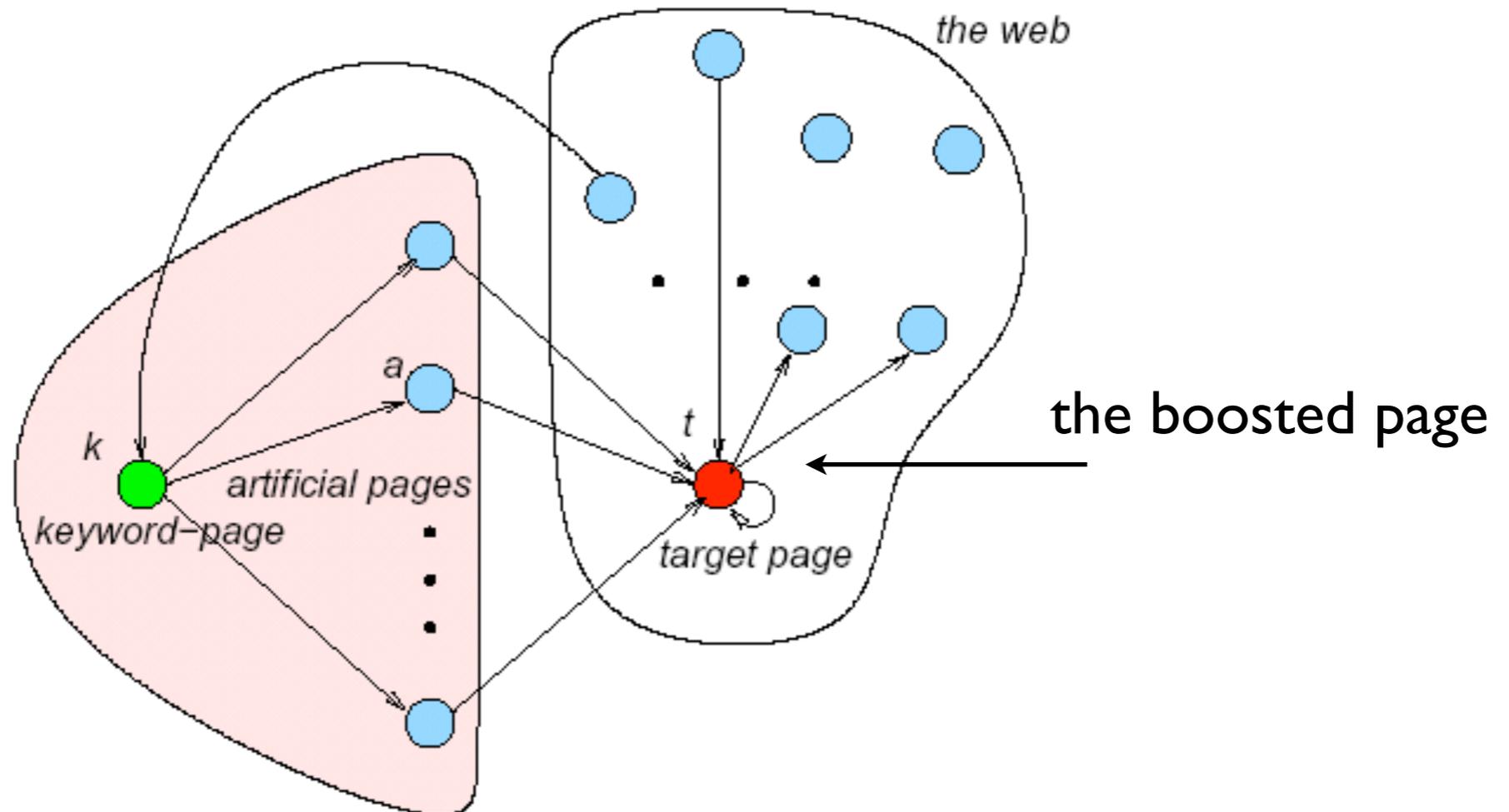
# Link Analysis

By analysing how pages link to each other, determine what a page is about and whether that page is deemed to be “relevant” ( does it deserve a ranking boost?)

- PageRank (Page and Brin), HITS (Kleimberg), TrustRank (Gyongy, Garcia-Molina, and Pedersen)

**topological analysis only!**

# Boosting PageRank: Link farms and the SEO/Search Engine “war”



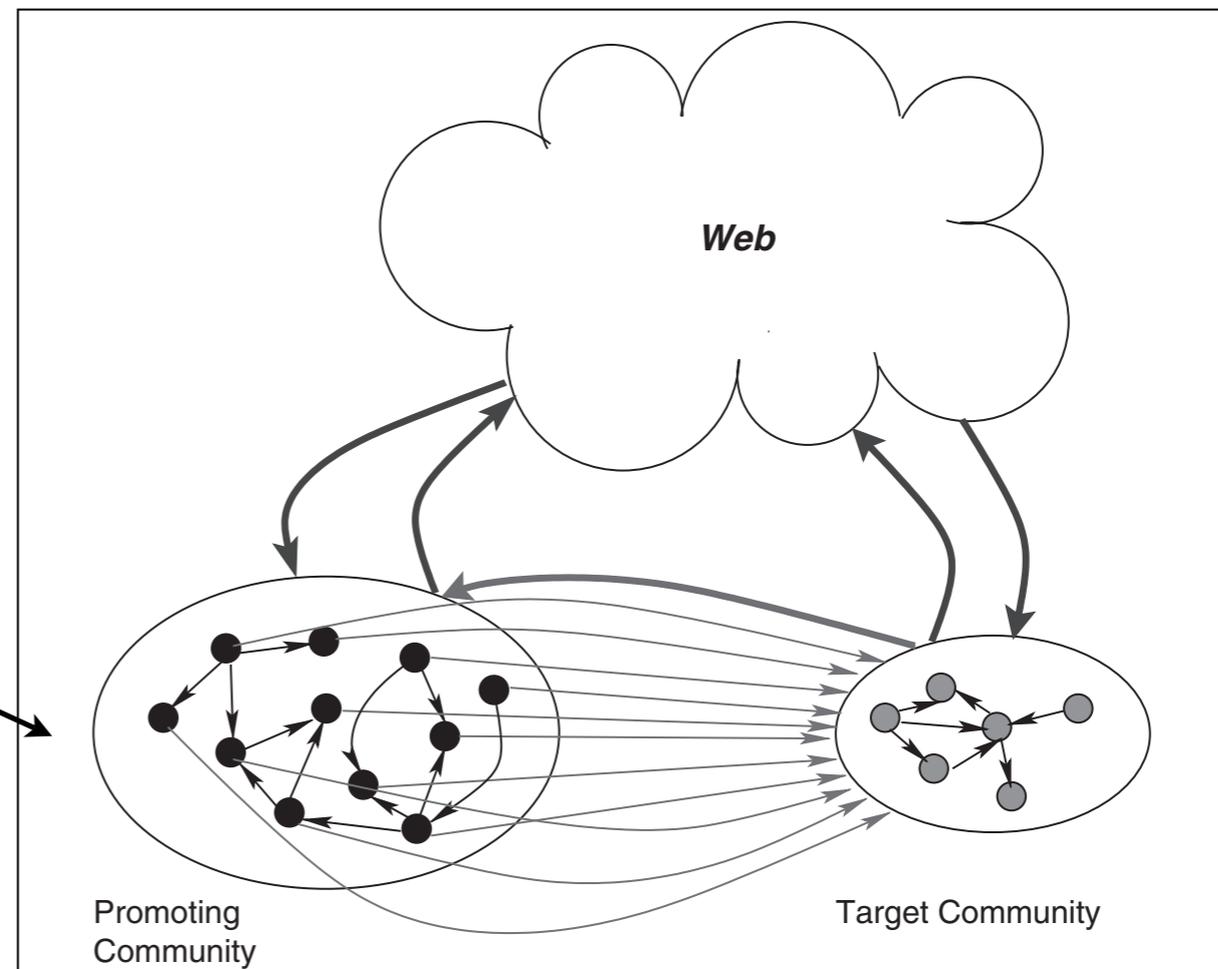
Can I combat spam  
by detecting the topology of the “promoting community”?

Gori&Witten,  
“The bubble of Web visibility”, CACM (2005)

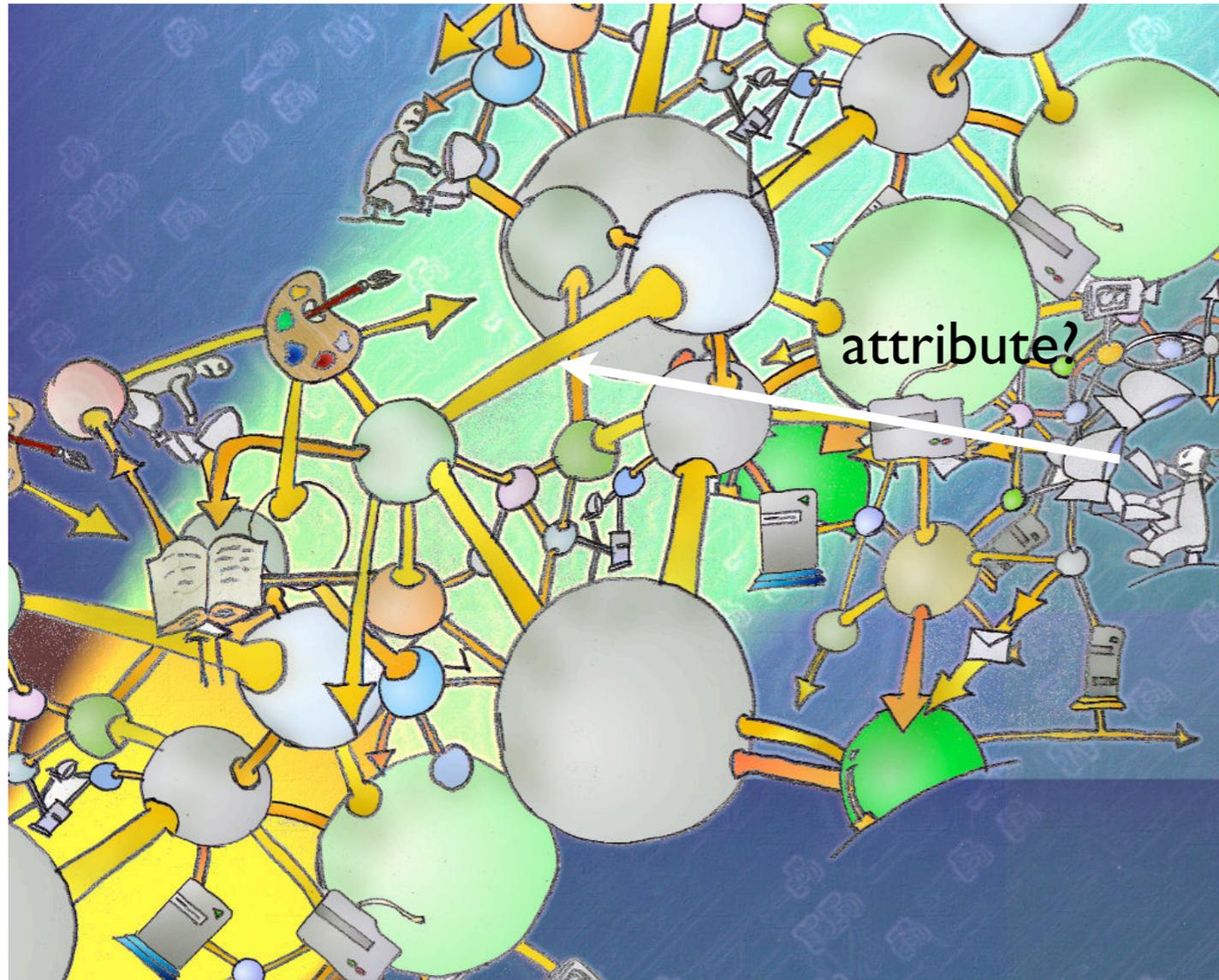
# Link farms cannot be detected

Bianchini, Gori, Scarselli, ACM-TOIT, 2005 “Inside PageRank”  
spam is independent of the topology of the promoting community!

Link analysis  
is not enough!



# Do I trust links?



No, I don't trust it, I need an attribute!

## **Learning in hyperlinked environments**

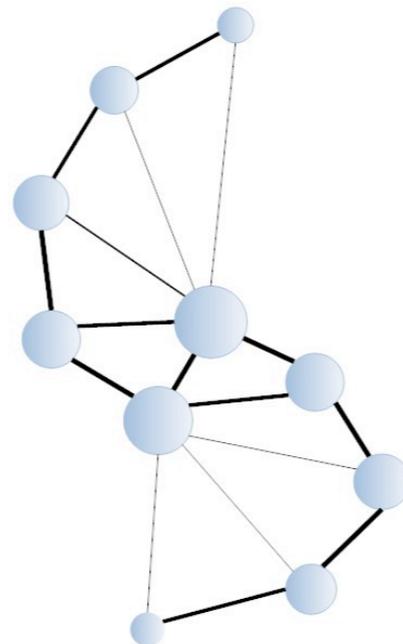
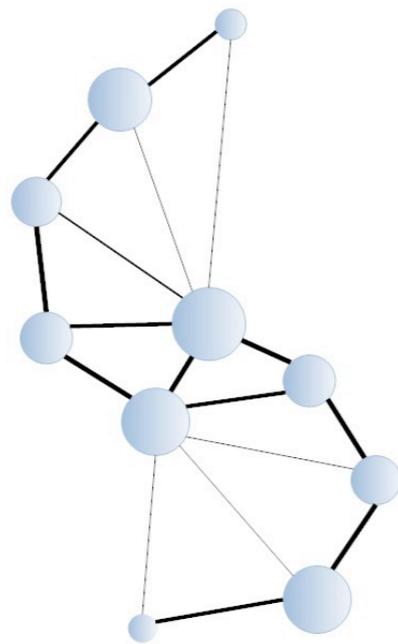
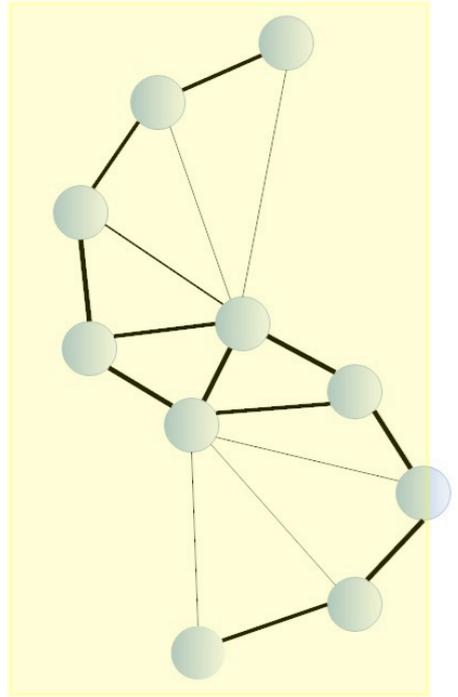
looks for methods to discover appropriate attributes for links (categorical or numerical)

# **Diffusion Learning**

ECIR-2007 Rome, 5 April 2007

# Diffusion machines

init: same rank

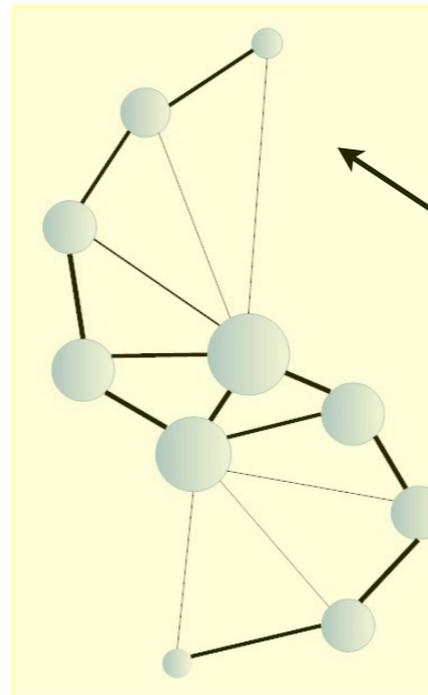
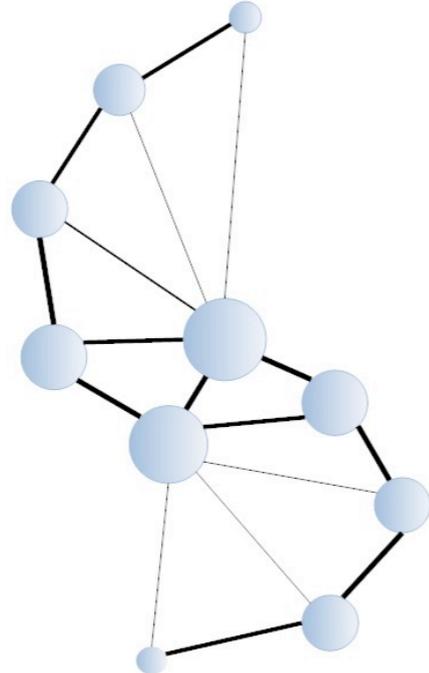
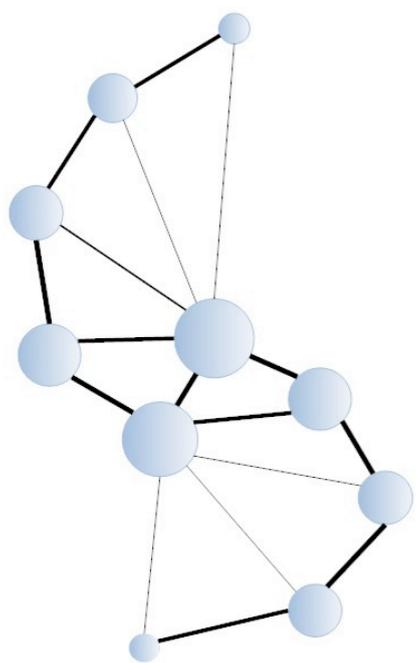


PageRank: A noticeable example

$$x_p = d \sum_{q \in pa[p]} \frac{x_q}{h_q} + (1 - d)$$

$$\mathbf{x} = d\mathbf{W}\mathbf{x} + (1 - d)\mathbf{1}_N$$

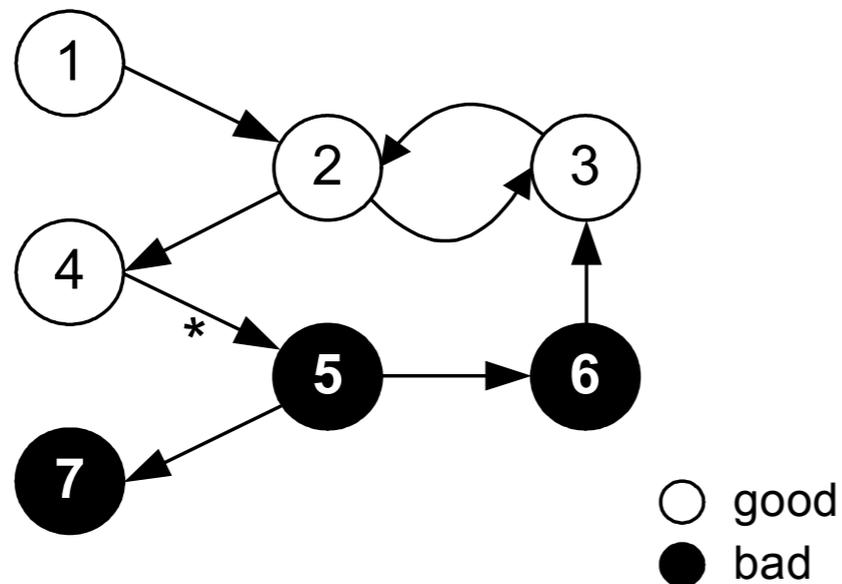
$$\mathbf{x}(t) = d\mathbf{W}\mathbf{x}(t - 1) + (1 - d)\mathbf{1}_N$$



Social nets  
random walk  
the efficiency issue

stationary distribution

# TrustRank



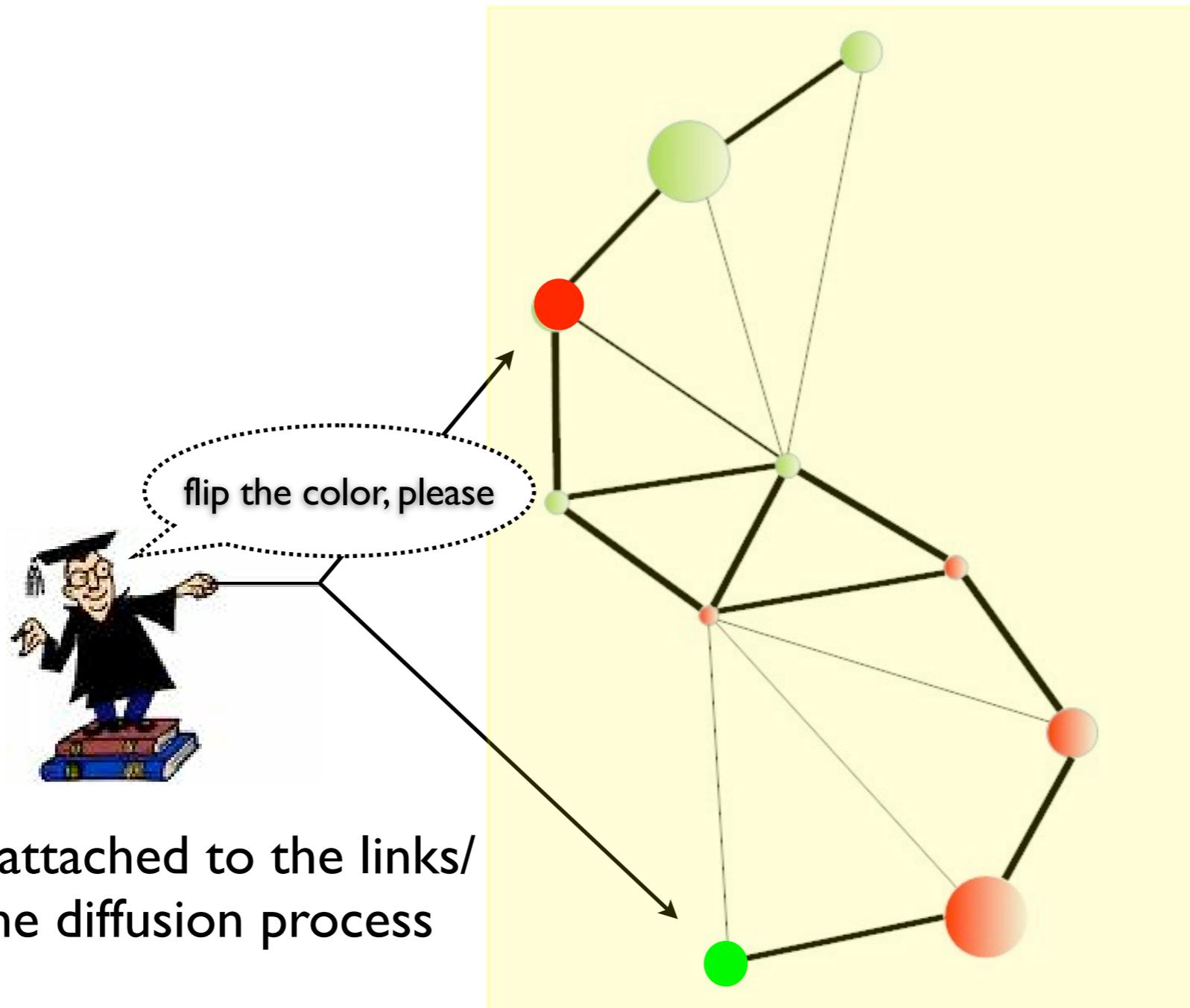
Invoke the “oracle”

propagate from “good” pages

```

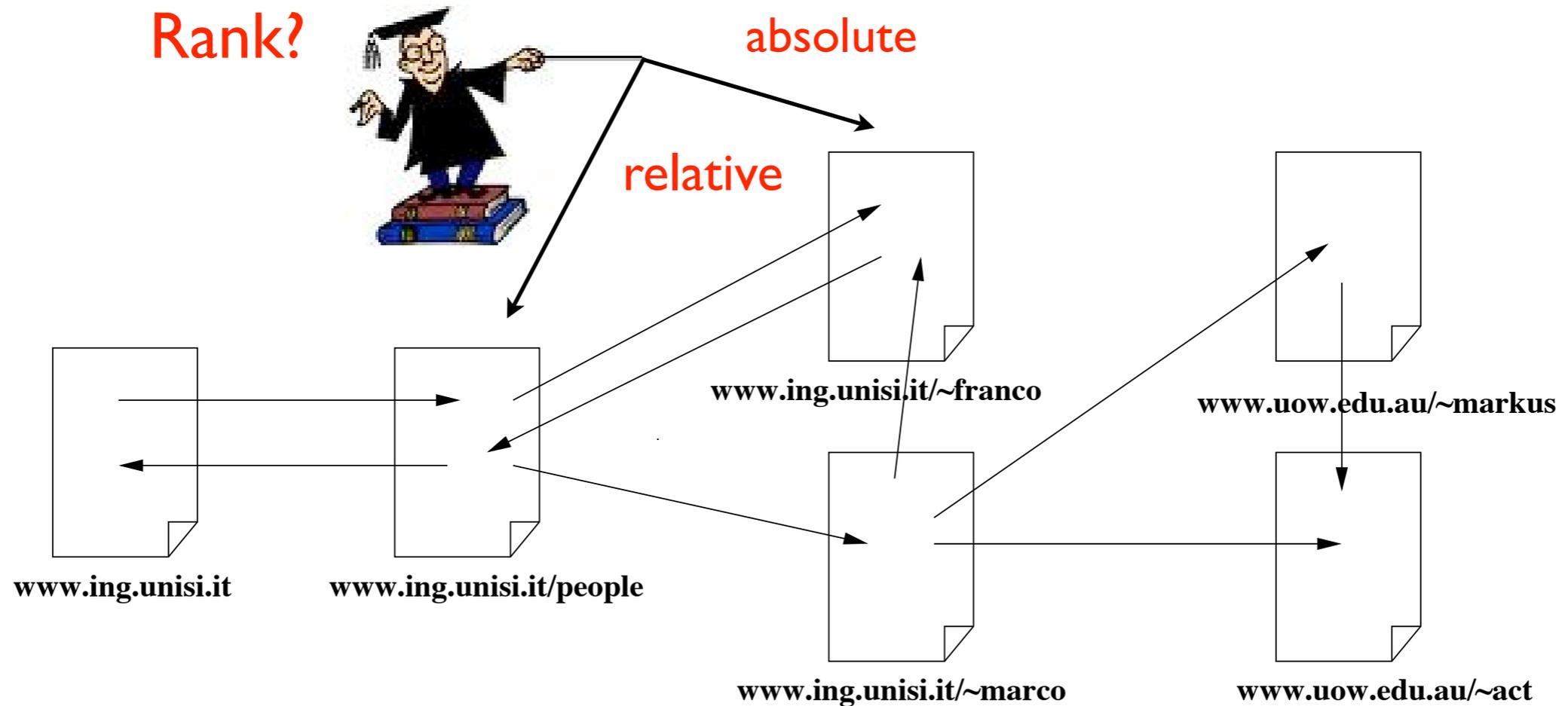
begin
  // evaluate seed-desirability of pages
  (1) s = SelectSeed(...)
  // generate corresponding ordering
  (2)  $\sigma = \text{Rank}(\{1, \dots, N\}, \mathbf{s})$ 
  // select good seeds
  (3) d =  $\mathbf{0}_N$ 
  for i = 1 to L do
    if  $O(\sigma(i)) == 1$  then
      d( $\sigma(i)$ ) = 1
  // normalize static score distribution vector
  (4) d = d/|d|
  // compute TrustRank scores
  (5) t* = d
  for i = 1 to MB do
    t* =  $\alpha_B \cdot \mathbf{T} \cdot \mathbf{t}^* + (1 - \alpha_B) \cdot \mathbf{d}$ 
  return t*
end
  
```

# Diffusion Learning machines



the parameters attached to the links/  
nodes affect the diffusion process

# Learning to Rank Web Pages



Set a *supervision value* (topology & content)

we put supervision on “some pages” and ask inference on the rest ...  
this is based on a learning machine that attaches attributes to links (trust)

# Deeper Inside Diffusion ... and Learning

- Diffusion machines to compute functions on graphs
- Variational approach: machines induced by target tracking and regularization

# Re-foundation of Calculus

(see e.g. A. Bensoussan & J.-L. Menaldi)

[Difference Equations on Weighted Graphs](#) (with A. Bensoussan). (*Journal of Convex Analysis (Special issue in honor of Claude LeMarechal)*, 12 (2005), pp. 13-44.)

- Notions of boundary, Hilbert space of functions, norms & semi-norms (Sobolev spaces)
- Extensions of differential operators ... Green's formula
- Variational problems, Harmonic functions, Dirichlet & Neumann' problems
- Relationships with random walk

# Harmolodic Functions

jazz saxophonist  
Ornette Coleman

## Local computation

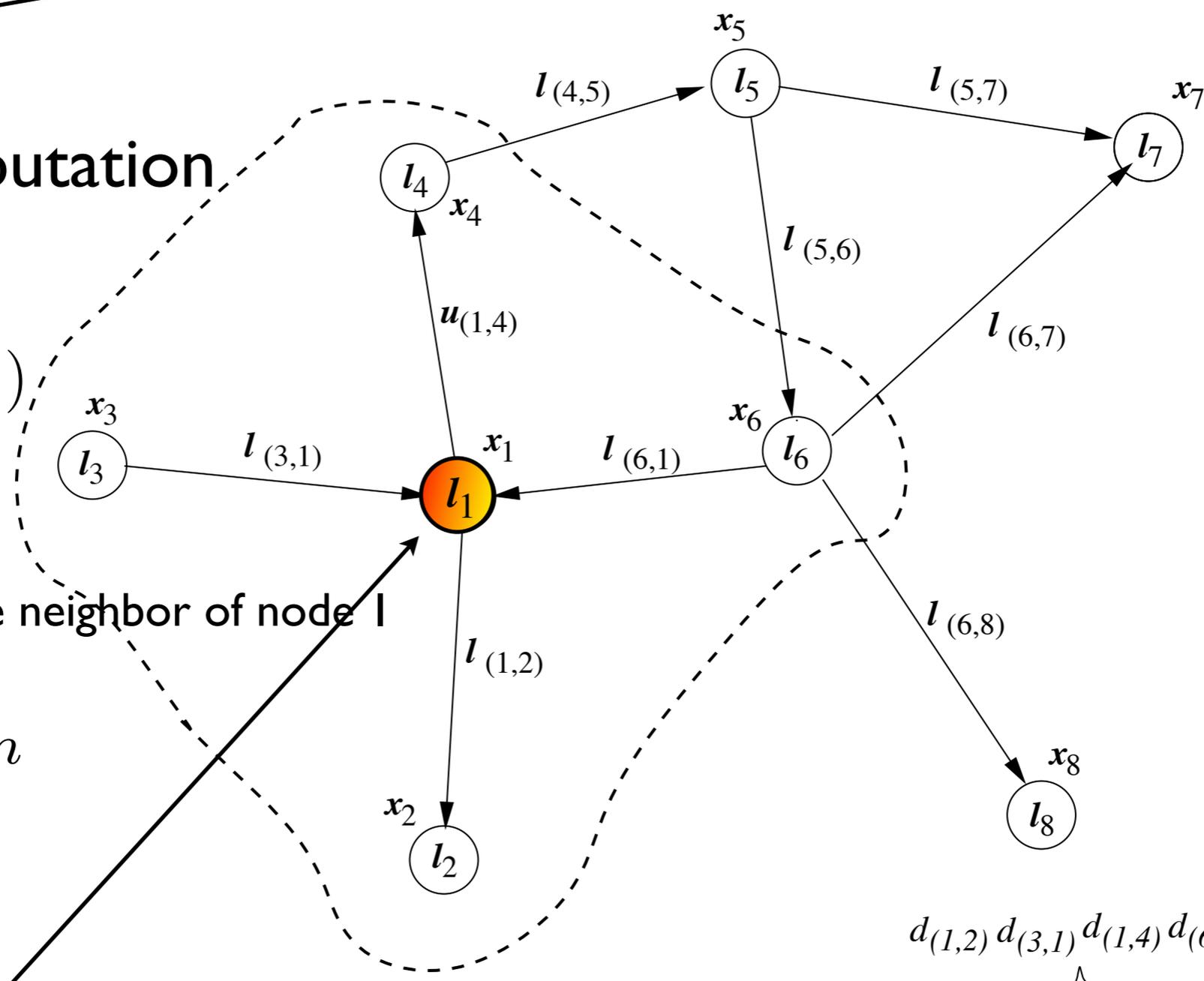
$$\mathbf{x} = F_w(\mathbf{x}, \mathbf{l})$$

$$\mathbf{o} = G_w(\mathbf{x}, \mathbf{l}_N)$$

contraction

$$\varphi_w(\mathbf{G}, n) = \mathbf{o}_n$$

harmolodic function



$$x_1 = f_w(l_1, \underbrace{x_2, x_3, x_4, x_6}_{x_{ne[1]}}, \underbrace{l_{(1,2)}, l_{(3,1)}, l_{(1,4)}, l_{(6,1)}}_{l_{col[1]}}, \underbrace{l_2, l_3, l_4, l_6}_{l_{ne[n]}}, \underbrace{0, 1, 0, 1}_{d_{(1,2)} d_{(3,1)} d_{(1,4)} d_{(6,1)}})$$

parametric dependence

ECIR-2007 Rome, 5 April 2007

coding the direction

# Diffusion Machines

$$\mathbf{x}_n = \sum_{u \in \text{ne}[n]} h_w(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{x}_u, \mathbf{l}_u), \quad n \in N$$

Something familiar ...  $h_w(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{x}_u, \mathbf{l}_u) = \mathbf{A}_{n,u} \mathbf{x}_u + \mathbf{b}_n$

$$\frac{\partial \mathbf{x}(t)}{\partial t} = \chi \mathcal{L} \mathbf{x}(t) + \beta$$

$$\mathcal{L} \mathbf{x} = \mathbf{b}$$

Poisson equation

$$\mathcal{L}(h_w) \mathbf{x} = \mathbf{x} - \sum_{u \in \text{ne}[n]} h_w(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{x}_u, \mathbf{l}_u)$$

$$\frac{\partial \mathbf{x}(t)}{\partial t} = \mathcal{L}(h_w) \mathbf{x}(t)$$

convergence to a  
stationary configuration

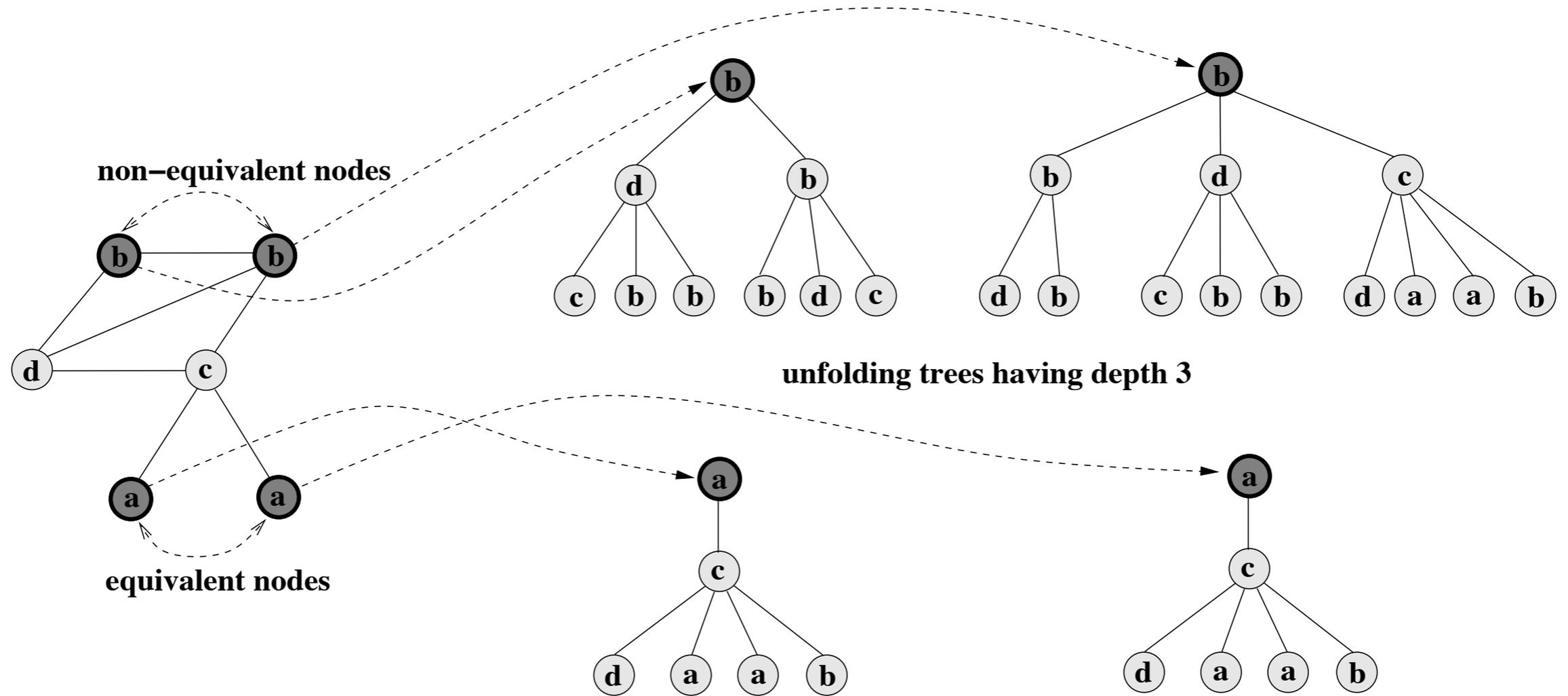
ECIR-2007 Rome, 5 April 2007

# Representation Power of Harmonic Functions

*Crucial question*

What are we missing with DLM  
(computing harmonic functions)?

# Unfolding Equivance



“Graph fibrations, graph isomorphism, and PageRank,” Boldi et al  
unfolding equivalence is equivalent to fibration prime graph  
... crucial for the isomorphism problem!

# Unfolding-Equivalence

**Definition 3** A function  $l : \mathcal{G} \times \mathcal{N} \rightarrow \mathbb{R}^m$  is said to preserve the unfolding equivalence on  $\mathcal{G}$ , if  $n \sim u$  implies  $l(\mathbf{G}, n) = l(\mathbf{G}, u)$ , for any nodes  $n, u$  of  $n\mathcal{G}$ . The set of functions that preserve the unfolding equivalence on  $\mathcal{G}$  will be denoted by  $\mathcal{F}(\mathcal{G})$ .

## **Theorem 4** APPROXIMATION BY POSITIONAL DLMS

Let  $\mathcal{D}$  be a domain that contains positional graphs. For any measurable function  $\tau \in \mathcal{F}(\mathcal{D})$  preserving the unfolding equivalence, any norm  $\|\cdot\|$  on  $\mathbb{R}^m$ , any probability measure  $P$  on  $\mathcal{D}$ , and any reals  $\varepsilon, \mu, \lambda$ , where  $\varepsilon > 0$ ,  $0 < \lambda < 1$ ,  $0 < \mu < 1$ , there exist two continuously differentiable functions  $f$  and  $g$  such that

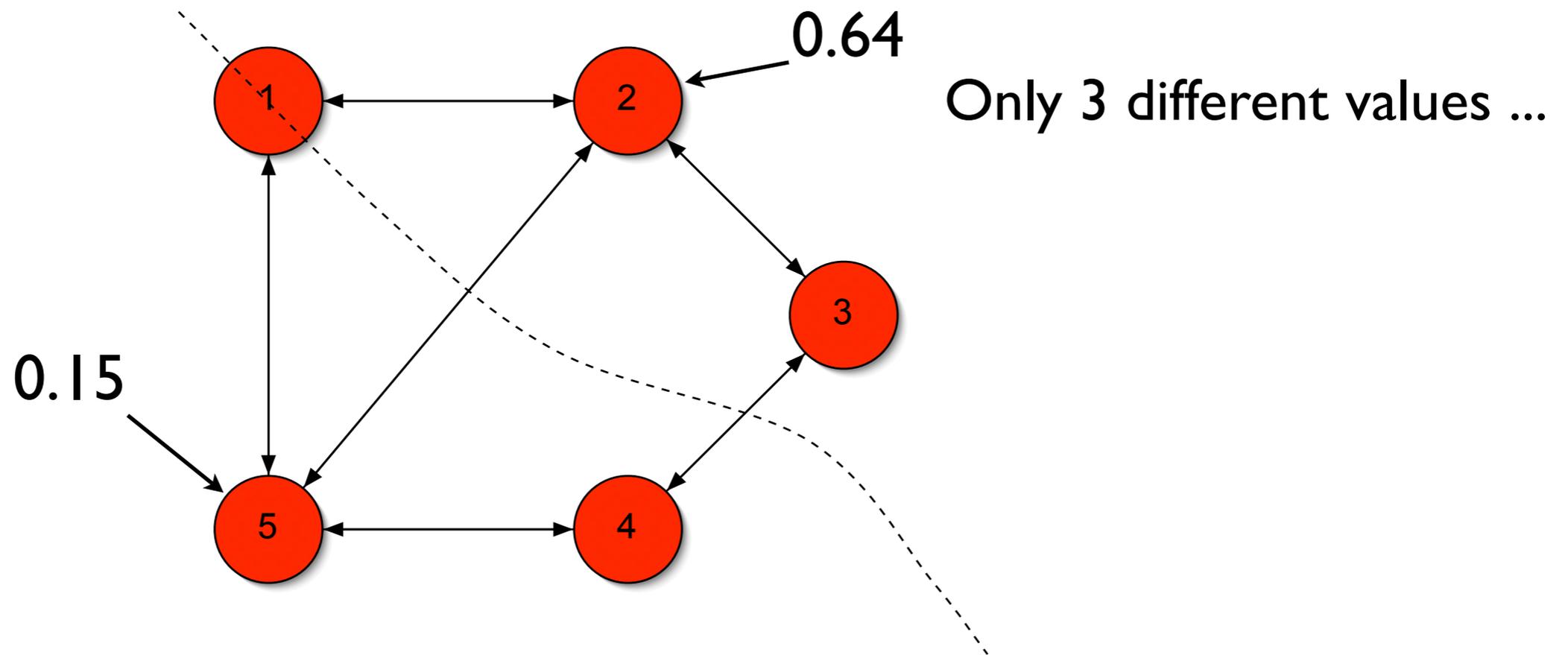
$$\begin{aligned}\mathbf{x}_n &= f(\mathbf{l}_n, \mathbf{l}_{\text{co}[n]}, \mathbf{x}_{\text{ne}[n]}, \mathbf{l}_{\text{ne}[n]}) \\ \mathbf{o}_n &= g(\mathbf{x}_n, \mathbf{l}_n), \quad n \in \mathbf{N},\end{aligned}$$

the global transition function  $F$  is a contraction with contracting constant  $\mu$ , the state dimension is  $s = 1$ , and the corresponding harmonic function defined by  $\varphi(\mathbf{G}, n) \doteq \mathbf{o}_n$  satisfies the condition

$$P(\|\tau(\mathbf{G}, n) - \varphi(\mathbf{G}, n)\| \geq \varepsilon) \leq 1 - \lambda$$

# What Cannot Be Calculated?

Warning!  
The symmetry constraints



Impossible!

The function doesn't preserve unfolding equivalence

# Diffusion Learning Machines (DLM)

$$\begin{aligned} \mathbf{x}_n(t+1) &= f_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{\text{co}[n]}, \mathbf{x}_{\text{ne}[n]}(t), \mathbf{l}_{\text{ne}[n]}) \\ \mathbf{o}_n(t) &= g_{\mathbf{w}}(\mathbf{x}_n(t), \mathbf{l}_n), \end{aligned}$$

discrete-time version

learning parameters

$$\frac{\partial \mathbf{x}(t)}{\partial t} = \mathcal{L}(h_{\mathbf{w}}) \mathbf{x}(t)$$

The dynamic evolution of DLM yields harmonic functions ...  
convergence to the fixed point of the map

# Semi-Supervised Framework

quadratic error w.r.t. the target

$$e_w = \sum_{i=1}^q (t_i - \varphi_w(\mathbf{G}, n_i))^2$$

entropy-based measure for ranking:  
Burges et al - RankNet, ICML05

$$C_{ij} \equiv C(o_{ij}) = -\bar{P}_{ij} \log P_{ij} - (1 - \bar{P}_{ij}) \log (1 - P_{ij})$$

$$P_{ij} \equiv \frac{e^{o_{ij}}}{1 + e^{o_{ij}}}$$

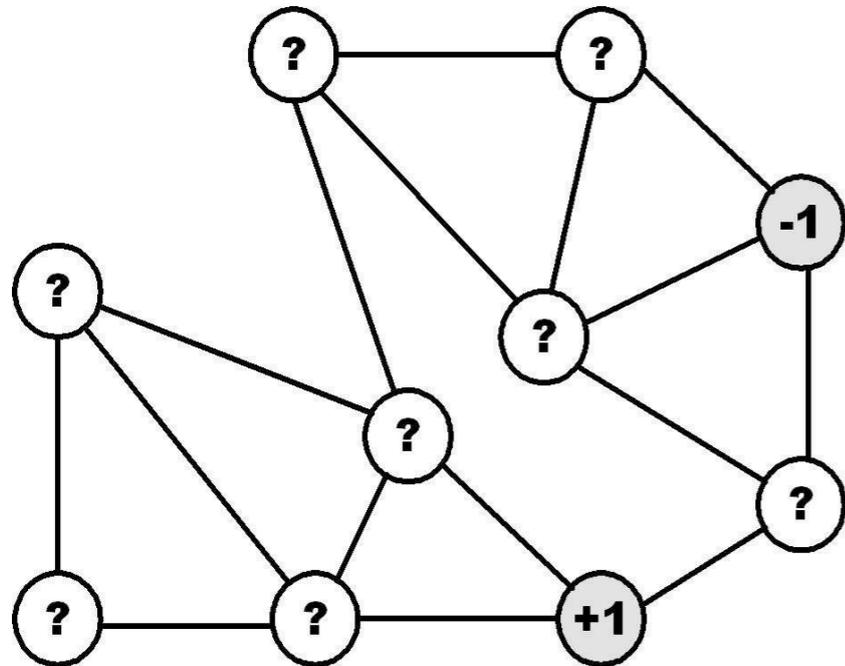
# From “B&W to color TV”

traditional machine learning framework

towards graphical domains

```
1  MAIN
2  initialize  $w$ ;
3   $x$ =FORWARD( $w$ );
4  repeat
5      $\frac{\partial e_w}{\partial w}$ =BACKWARD( $x, w$ );
6      $w=w - \lambda \cdot (\frac{\partial e_w}{\partial w})'$ ;
7      $x$ =FORWARD( $w$ );
8  until the stopping criterion is met;
9  return  $w$ ;
10 end
```

# A variational approach (Zhou & Schoelkopf, 2004)



$$\arg \min_{f \in L^2(V)} \left\{ \mathcal{S}(f) + \frac{\mu}{2} \|f - y\|^2 \right\}$$

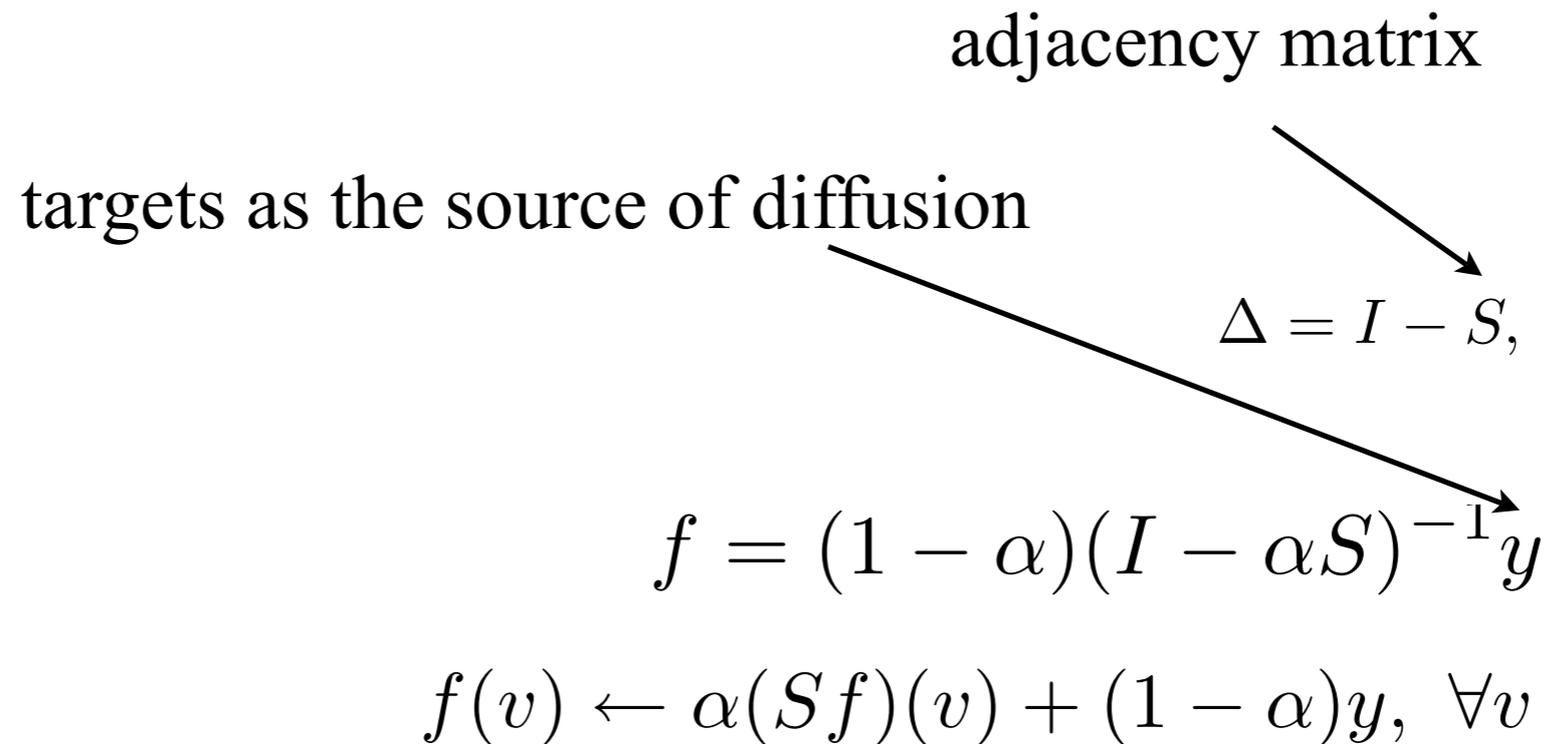
regularization term  $\nearrow$   $\mathcal{S}(f)$   
 $\searrow$  error w.r.t. target  $\frac{\mu}{2} \|f - y\|^2$

$$\mathcal{S}(f) = \frac{1}{2} \sum_v \|\nabla_v f\|^2.$$

$$\frac{\partial f}{\partial e} \Big|_u = \sqrt{\frac{w(u, v)}{d(u)}} f(u) - \sqrt{\frac{w(u, v)}{d(v)}} f(v).$$

# An elegant solution

**Theorem 1.** *The solution of the optimization problem (2.9) satisfies  $\Delta f + \mu(f - y) = 0$ .*



**A PageRank-like solution!**

# One step further ... at airgroup.unisi.it

initial heuristic value

$$\left\{ \begin{array}{l} \varphi^e = \arg \min_{\varphi \in \mathcal{H}(V)} \left\{ S_2(\varphi, w^e) + \frac{\mu \|\varphi - y\|^2}{2} \right\} \\ w^{e+1} = \arg \min_{w \in \mathcal{H}(E)} \left\{ S_2(\varphi^e, w) + \frac{\mu \|w - \bar{w}\|^2}{2} \right\} \end{array} \right.$$

weight refinement

# **Impact on information retrieval**

# What is *diffusion learning* for?

- pattern classification: Hancock, Bekin, ...
- link prediction: Zhou&Schoelkopf, SRL04
- network evolution: Barabasi (Linked) et al ...
- community detection (M.B. Hasting, Physical Review 06) ...
- ...

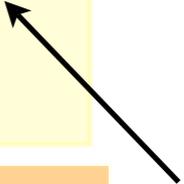
# Ranking

- Tsoi et al ACM-TOIT 2007
- Diligenti et al IJCAI05
- Scarselli et al WVIC05
- Chakabarti et al KDD06

# airgroup.dii.unisi.it

- Learning to rank papers (ACM portal)
- HITS (Kleimberg, 1998)
- PageRank-like functions

links only

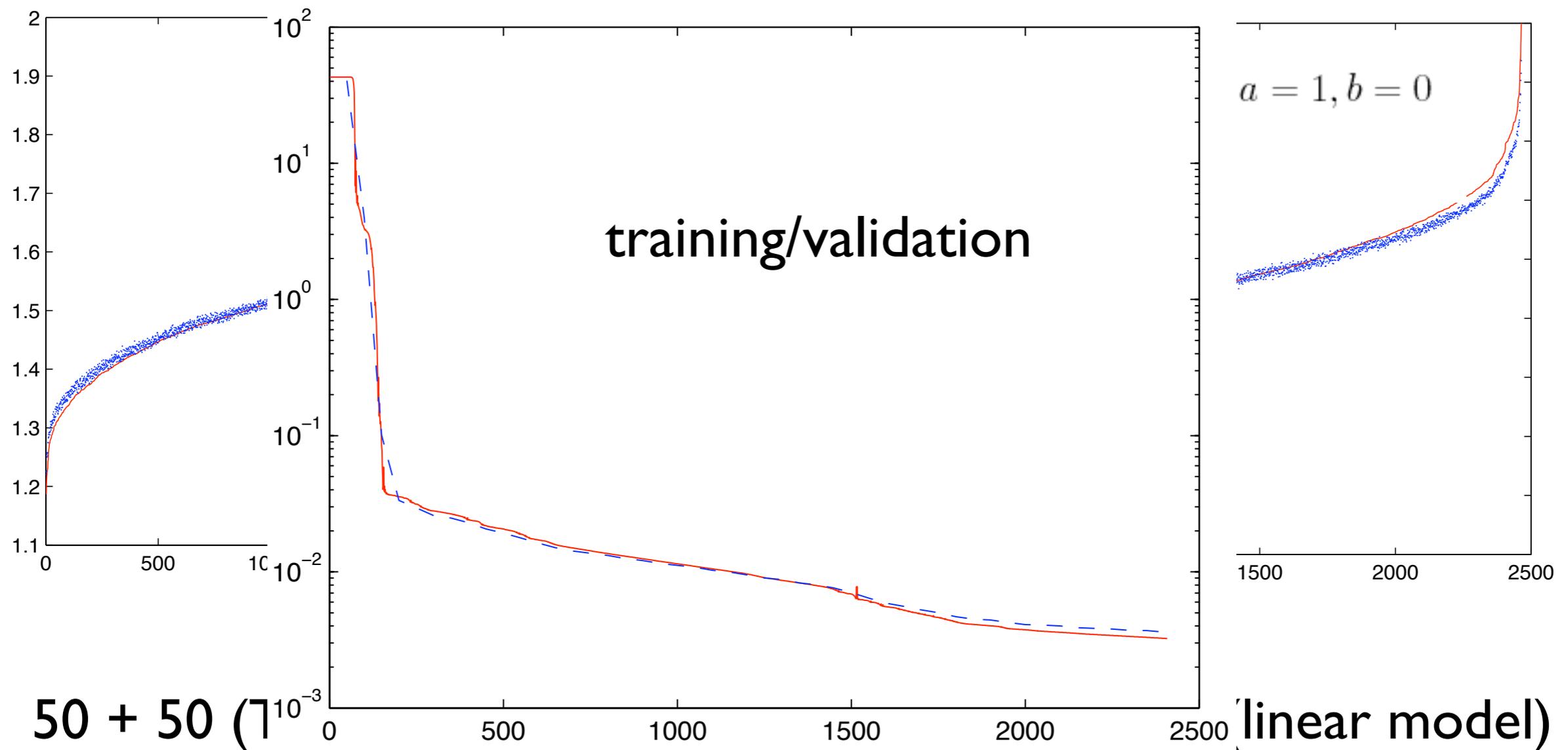


links and content!



# Learning PR-like Functions

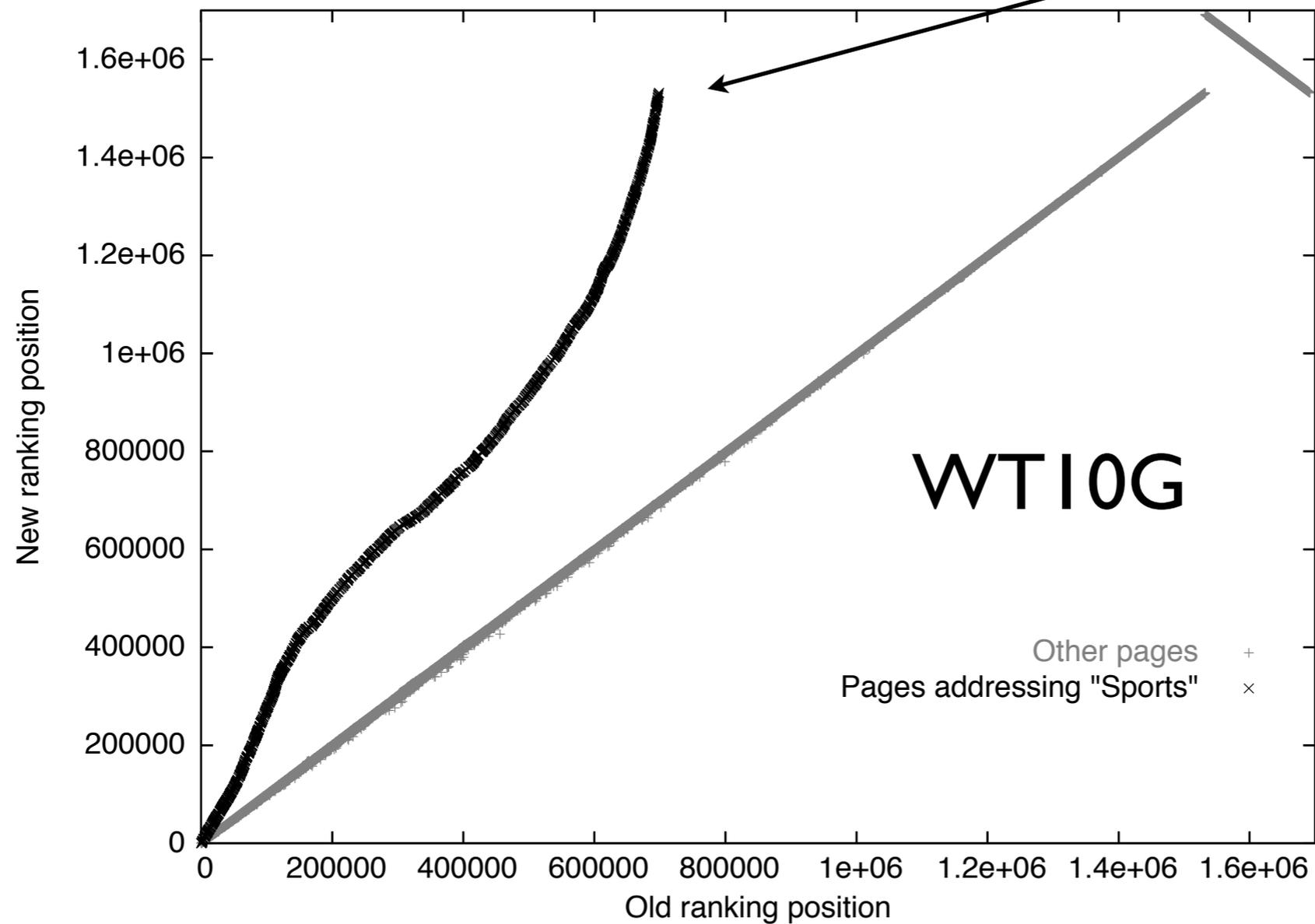
Learn the PageRank (PR)  
Each node has a label with 2 Booleans



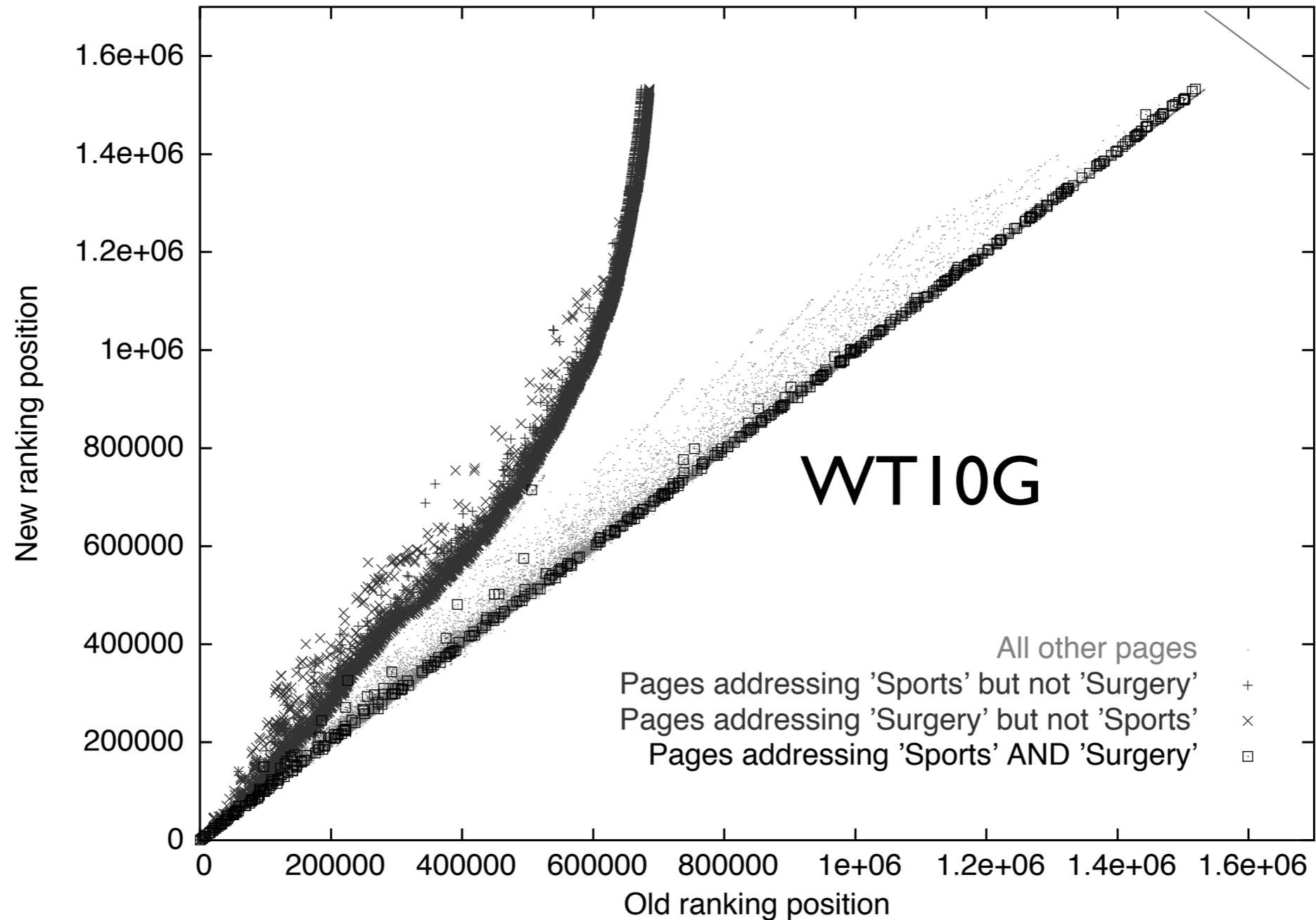
ECIR-2007 Rome, 5 April 2007

# more focus on a selective topic

double on sport



# The "XOR"



# Conclusions

- Beyond link analysis: A general framework to build functions on the Web (content and hyperlinks)
- Links with relational learning, Dagstuhl seminar (15-20 April) on “Probabilistic, Logical, and Relational Learning”
- An interesting perspective: Facing the link trust problem (Web spam, ... harder time for spammers)
- A lot of work for choosing the “right trade-off” ... see e.g. complexity issues in Domingo et al (QD-PageRank)

# Acknowledgements

- Research at Univ. of Siena (Franco Scarselli, Gabriele Monfardini, Augusto Pucci, Vincenzo Di Massa)
- Markus Hagenbuchner, Linus Yong (Univ. of Wollongong, Australia)
- Ah Chung Tsoi (Monash Univ., Melbourne, Australia)
- Hendrik Blockeel & Werner Uwents Université catholique de Louvain (Belgium)
- Paolo Frasconi (Univ. of Florence, Italy)
- Alessandro Sperduti (Univ. of Padua, Italy)
- Yoshua Bengio (Université de Montréal, Canada)

**Thanks for your attention!**