# Modeling the structure and evolution of online discussion cascades

Vicenç Gómez[1]    Hilbert J Kappen[1]    Nelly Litvak[2]
Andreas Kaltenbrunner[3]

[1]Donders Institute for Brain, Cognition and Behaviour,
Radboud University, Nijmegen, The Netherlands

[2]Faculty of Electrical Engineering, Mathematics and Computer Sciences,
University of Twente, Enschede, The Netherlands

[3]Social Media Research Group, Barcelona Media,
Barcelona, Spain

ETC* July 26th, 2012

# Outline

# Agenda

## Structure and evolution of online discussion cascades

Gómez V., Kappen H. J., Litvak N., and Kaltenbrunner, A. (2012).
A likelihood-based framework for the analysis of discussion threads.
*World Wide Web Journal*, 2012, pp 1-31.

Gómez V., Kappen H. J., and Kaltenbrunner, A. (2011).
Modelling the Structure and Evolution of Discussion Cascades.
In *HT2011 22nd ACM Conference on Hypertext and Hypermedia*, , Eindhoven, The Netherlands.

## But first ...

- a brief presentation of Fundació Barcelona Media.

# Outline

# Barcelona Media

## What is Barcelona Media?

- A private non-profit organisation for the research and the innovation in the communication industry.
- A Technology Centre collaborating with companies and institutions to foster the competitiveness of the sector.
- Close relations with Universitat Pompeu Fabra.
- Participation in 41 European projects (12 as a coordinator).

## Research Lines

- Audio
- Image
- Voice and Language
- Perception and Cognition
- Information retrieval (Yahoo! Labs Barcelona)
- Social Media

# The Social Media Research Group at Barcelona Media

## Mission

- **Combine qualitative and quantitative methods** to generate knowledge about the interplay of social behaviour and "Social Media".

## Main research lines

- Survey Research
- Sentiment analysis
- Social Network Analysis
- Study of online conversation

## Main data sources

- Twitter
- Wikipedia
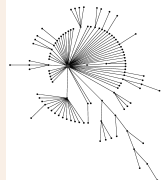- Online Forums
- Online Social Networks

# Outline

# Motivation

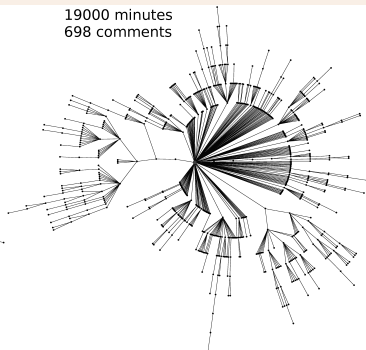## Example of online discussion (from Slashdot)



2000 minutes
109 comments

5000 minutes
314 comments

19000 minutes
698 comments

Title: *"Can Ordinary PC Users Ditch Windows for Linux?*.

- Online conversations as networks: **nodes** correspond to comments, **edges** represent a reply action.

# Motivation - Online discussion threads

## Scientific questions

- What are the structural patterns governing these responses?
- What determines the growth of a conversation?
- Is there a generative model that captures their statistical properties?
- Can we use the model parameters to characterize websites, user behaviour, discussions?

## Implications / Applications

- Understanding communication in large webspaces that comprise many-to-many interaction.
- Understanding diffusion of news and opinion in social networks.
- Community management, forum design/maintenance, ...

# Outline

# Online discussion threads
Datasets

We collected data from the following sources:

Slashdot (SL) : Technological news aggregator.
$473,065$ discussions, $2 \cdot 10^6$ comments,
$93 \cdot 10^3$ users

Barrapunto (BP) : Spanish version of Slashdot.
$44,208$ discussions, $4 \cdot 10^5$ comments,
$50 \cdot 10^3$ users

Meneame (MN) : Spanish Digg clone (general news
aggregator)
$58,613$ discussions, $2.1 \cdot 10^6$ comments,
$5,4 \cdot 10^4$ users.

Wikipedia (WK) : discussion pages related to every article.
$871,485$ discussions, $\approx 10^7$ comments,
$3.5 \cdot 10^5$ users.

# Motivation

Example of discussion in Slashdot (post):

# Motivation
Example of discussion in Slashdot (comments):

# Motivation
Example of discussion in Barrapunto (comments):

Pues claro que es un lujo... (Puntos:2)
por coricpablo (20771) <(coricpablo) (at) (gmail.com)> el Lunes, 23 julio de 2012, 15:20h (#1335551)
( Última bitácora: jueves, 24 febrero de 2011, 14:21h )

Deberían invertir todo ese presupuesto de ciencia que desperdiciamos en intentar curar el cáncer o hacer más eficientes las energías renovables en algo útil y más acorde al país en que vivimos como formar a cocteleros y dj's para nuestros chiringuitos playeros.

Somos un país de servicios, tenemos que damos cuenta de ello. Los Vascos, como siempre, son los que se han dado cuenta y van un paso por delante [bculinary.com].

Asumámoslo, nunca vamos a ser los médicos ni los científicos ni los ingenieros de Europa, somos los que ponen las copas cuando se vienen de vacaciones.
--
All my life, I have enjoyed the reputation of being someone who disrupted prevailing ideas.

(Benoit Mandelbrot)

[ Responder ]

Re:Pues claro que es un lujo... (Puntos:0)
por porocito hablador el Lunes, 23 julio de 2012, 16:48h (#1335566)

"Deberían invertir todo ese presupuesto de ciencia que desperdiciamos en intentar curar el cáncer o hacer más eficientes las energías renovables en algo útil y más acorde al país en que vivimos como formar a cocteleros y dj's para nuestros chiringuitos playeros. "

Pues no se si lo has dicho con ironía, pero...

http://www.cefe.gva.es/agenda.asp?id=613 [cefe.gva.es]
La conselleria de Educación, Formación y Empleo implanta 26 nuevos ciclos formativos

(blablablá) ...

Como novedad para este curso 2012-2013 se han implantado las nuevas titulaciones de Técnico Superior de Artista Fallero y Construcción de Escenografías y el de Técnico en Vídeo Disc-jockey y Sonido.

(blablablá)...

Ya ves, los vascos no son los únicos "visionarios" :D

[ Responder | Padre ]

y mientras tanto se suben sus sueldos... (Puntos:0)
por porocito hablador el Martes, 24 julio de 2012, 06:21h (#1335608)

El Gobierno de Castilla-La Mancha ha aumentado la dotación económica para alta dirección un 157% en los presupuestos de la Junta para 2012, con incrementos del 8,6% y del 10,9 % en algunos sueldos, según ha dicho hoy el secretario regional de Organización del PSOE, Jesús Fernández Vaquero. En rueda de prensa, Fernández Vaquero ha explicado que esta subida se recoge en el tomo 1 de los presupuestos generales de la Junta, en el que se constata que la dotación económica para alta dirección del Gobierno de la comunidad autónoma ha pasado de 22,3 millones de euros en 2011 a 56,64% en 2012. Fernández Vaquero ha indicado que los sueldos de los directores generales se incrementan un 10,9%, pasando de 53.450 euros brutos anuales a 59.270. A su vez, el sueldo de los viceconsejeros pasan de 58.230 euros brutos a 63.320, según ha precisado. Esta subida contrasta, según Fernández Vaquero, con el hecho de que a todos los funcionarios se les baje un 3% su salario. El secretario regional de organización del PSOE ha reconocido que los sueldos de los altos cargos de Castilla-La Mancha son más bajos que los de otras comunidades autónomas pero ha criticado la doble vara de medir: "Al señor que gana 1.000 euros o gana 900 euros se le baja o se le reduce el sueldo". "No podemos estar hablando de contener el gasto y reducir, cuando si hay algo que le afecta a ella lo que hace es aumentar el gasto", ha añadido Fernández Vaquero en alusión a la presidenta de Castilla-La Mancha, María Dolores de Cospedal. http://politica.elpais.com/politica/2012/06/01/act ualidad/1338573325_921691.html [elpais.com]

[ Responder | Padre ]

Re:y mientras tanto se suben sus sueldos... (Puntos:0)
por porocito hablador el Martes, 24 julio de 2012, 06:36h (#1335609)

Ya se ve donde va a parar el dinero que no pagan a los farmacéuticos de Castilla la Mancha, a los enchufados de la señora Cospedal

[ Responder | Padre ]

Re:Pues claro que es un lujo... (Puntos:0)
por porocito hablador el Martes, 24 julio de 2012, 08:25h (#1335613)

Con idiotas como tú diciendo esas sopiapolleces es como me doy cuenta de que España es un país que tiene exactamente aquello que merece.

# Motivation
## Example of discussion in Meneame:

# Motivation
## Example of discussion in Wikipedia (I)

# Motivation
## Example of discussion in Wikipedia (II)

### I think the bot is archiving this talk page too soon and too frequently.                    [edit]

What do you think? Grundle2600 (talk) 22:53, 23 June 2009 (UTC)

> I think this page is very short, so no harm in keeping the discussions a bit longer. As an incremental step I'll increased the archive time from 7 days to 10 days, per your comment. Wikidemon (talk) 01:19, 24 June 2009 (UTC)

>> Thank you. I didn't even know you could change it like that! Grundle2600 (talk) 02:40, 24 June 2009 (UTC)

### Coup d'etat in Honduras                                                                      [edit]

Why does this article not mention Obama's first coup d'etat. Just go to http://www.globalresearch.ca ⧉ to find out more information about this current event. —Preceding unsigned comment added by 99.255.173.155 (talk) 03:26, 5 July 2009 (UTC)

> Because there needs to be better sourcing for something as controversial as this. QueenofBattle (talk) 19:39, 5 July 2009 (UTC)

>> Better sourcing ⧉ Grundle2600 (talk) 01:22, 8 July 2009 (UTC)

>>> Typing in Obama's name, consitution, and the guys name does not mean better sourcing. Simply put, there has to be a very reliable source that says exactly that this was Obama's first coup d'etat. There can be no WP:OR or WP:SYNTH. Brothejr (talk) 01:30, 8 July 2009 (UTC)

>>>> OK. That's a good point. ABC News ⧉ says that Honduras' President Manuel Zelaya tried to give himself another term, despite the fact that Honduras' constitution prohibits such a thing, and that Obama said he supports Zelaya in this action. Grundle2600 (talk) 01:37, 8 July 2009 (UTC)

>>>>> And of course it says much more than that - you leave out a great deal of context, not to mention leaving to the side more in-depth coverage in any number of other sources. Also you should explain exactly how you envision the Honduran coup being mentioned in this article. Personally I think it's rather important in a global history sense, but I'm not sure it deserves mention in the article on the Presidency of Barack Obama. He has talked about it, of course, but then again so has most every major leader and foreign policy NGO in the Western Hemisphere. --Bigtimepeace | talk | contribs 07:35, 8 July 2009 (UTC)

>>>>>> I was just helping out with sources. I don't necessarily think that it should, or shouldn't, be cited in this article. Perhaps it could be cited in the article about the event itself, instead. Grundle2600 (talk) 08:54, 8 July 2009 (UTC)

>>>>>>> 2009 Honduran coup d'état already says, "President Barack Obama of the United States said "We believe that the coup was not legal and that President Zelaya remains the President of Honduras."[55][157]" I think that's enough, and it doesn't necessarily have to be added to this article. It could be cited in this article to, but I don't give it high priority in this article. Grundle2600 (talk) 09:04, 8 July 2009 (UTC)

>>>>>>>> If you "help out with sources" on an article talk page the assumption is going to be that you want that

# Most discussed Wikipedia articles
Top 20 articles ordered by number of chains in the discussion [Laniado et al. 2011]

| # | Title | chains | comments | users | h-index | max. depth | edits |
|---|-------|--------|----------|-------|---------|------------|-------|
| 1 | Intelligent design | 2413 | 22454 (3) | 954 (13) | 16 (20) | 20 (358) | 9179 (53) |
| 2 | Gaza War | 2358 | 17961 (6) | 607 (47) | 19 (2) | 27 (28) | 11499 (29) |
| 3 | Barack Obama | 2301 | 22756 (2) | 2360 (2) | 18 (6) | 21 (245) | 17453 (6) |
| 4 | Sarah Palin | 2182 | 19634 (4) | 1221 (9) | 17 (10) | 25 (56) | 12093 (24) |
| 5 | Global warming | 2178 | 19138 (5) | 1382 (5) | 17 (10) | 20 (358) | 14074 (15) |
| 6 | Main Page | 2065 | 32664 (1) | 5969 (1) | 15 (34) | 22 (169) | 4003 (674) |
| 7 | Chiropractic | 1772 | 13684 (13) | 243 (389) | 18 (6) | 29 (17) | 6190 (204) |
| 8 | Race and intelligence | 1764 | 13790 (12) | 410 (126) | 17 (10) | 24 (74) | 7615 (100) |
| 9 | Anarchism | 1589 | 14385 (9) | 496 (76) | 20 (1) | 28 (22) | 12589 (19) |
| 10 | British Isles | 1556 | 12044 (16) | 576 (56) | 17 (10) | 23 (113) | 4047 (658) |
| 11 | CRU[1] hacking incident | 1551 | 11536 (17) | 474 (88) | 17 (10) | 20 (358) | 2346 (2364) |
| 12 | Jesus | 1397 | 17916 (7) | 1239 (7) | 13 (119) | 16 (1383) | 17081 (7) |
| 13 | Circumcision | 1356 | 10469 (21) | 436 (113) | 17 (10) | 26 (42) | 7354 (117) |
| 14 | Homeopathy | 1323 | 13509 (14) | 516 (68) | 17 (10) | 25 (56) | 6902 (151) |
| 15 | George W. Bush | 1281 | 15257 (8) | 1969 (3) | 14 (65) | 18 (676) | 32314 (1) |
| 16 | September 11 attacks | 1250 | 13830 (11) | 1244 (6) | 16 (20) | 26 (42) | 11086 (30) |
| 17 | Evolution | 1165 | 13404 (15) | 942 (16) | 13 (119) | 23 (113) | 9780 (44) |
| 18 | Catholic Church | 1162 | 14104 (10) | 620 (43) | 15 (34) | 18 (676) | 14082 (14) |
| 19 | Cold fusion | 1098 | 8354 (29) | 359 (174) | 15 (34) | 20 (358) | 4320 (557) |
| 20 | 2008 South Ossetia war | 1075 | 10596 (20) | 853 (20) | 17 (10) | 23 (113) | 9930 (43) |

In parenthesis: rank according to the corresponding variable

[1] Climatic Research Unit

# Temporal patterns. From Kaltenbrunner et al (LAWEB 2007)
Time series of total number of comments



1. "Sustained" activity coupled with the circadian rhythm.

# Temporal patterns. From Kaltenbrunner et al (LAWEB 2007)
## Single post level analysis



- Posts create cascades of comments which propagate over the network.
- All posts show a stereotyped behaviour.
- Response times can be described using a log-normal distribution.

# Online discussion threads
## Examples of real discussions

Typical cascades for each website:



(a) Slashdot   (b) Barrapunto   (c) Meneame   (d) Wikipedia

Degrees Slashdot:

# Online discussion threads
## Global analysis



- SL, BP and MN present a distribution with a defined scale.
- Cascade sizes in Wikipedia "seem to be" scale-free.

# Outline

# Model definition

## Our approach:

- The model must reproduce:
    - The statistical structure of threads.
    - Their evolution.
- No content involved.
- No authorship.
- Essentially *"Which comment is going to be replied next?"*

## Empirical facts:

- Popular comments receive more replies: *preferential attachment*.
- New comments are more *attractive* than old ones.
- Replies to the post behave different than replies to comments.

# Model definition

- Thread **representation**: vector of parent nodes $\boldsymbol{\pi}$, where $\pi_t$ denotes the parent of the node with id $t + 1$ added at time-step $t$.

$$\boldsymbol{\pi}_0 = ()$$
$$\boldsymbol{\pi}_1 = (1)$$
$$\cdots$$



$$\boldsymbol{\pi} = \boxed{1\;1\;2\;1\;5\;2\;1\;6\;?}$$

---

### Parameters of the model

- At time $t$, the **popularity** of node $k$ is its degree:

$$d_{k,t}(\boldsymbol{\pi}_{(1:t-1)}) = \begin{cases} 1 + \sum_{m=2}^{t-1} \delta_{k\pi_m} & \text{for } k \in \{1, \ldots, t\} \\ 0 & \text{otherwise,} \end{cases} \quad (d_{k,t} \text{ is weighted by } \alpha)$$

- At time $t$, the **novelty** of node $k$ is $n_{k,t} = \tau^{t-k+1}, \quad \tau \in [0, 1]$.

- **Root bias**: The bias of a node $k$ is is either zero or $\beta$ for the root:

$$b_k = \beta, \qquad \text{for } k = 1, \text{ and } 0 \text{ otherwise.}$$

## Model definition

- We define a model by means of its associated *attractiveness* function $\phi(\cdot)$, which is defined for each of the nodes.

- At time $t + 1$, a new node is linked to node $k$ with probability:

$$p(\pi_t = k | \boldsymbol{\pi}_{(1:t-1)}) = \frac{\phi(k)}{Z_t}, \qquad Z_t = \sum_{l=1}^{t} \phi(l),$$

- Different model variants:

| Model | Attractiveness funct. $\phi(\cdot)$ | Parameters $\boldsymbol{\theta}$ | Constraint |
|-------|-------------------------------------|----------------------------------|------------|
| Full model (**FM**) | $\alpha d_{k,t} + b_k + \tau^{t-k+1}$ | $\{\alpha, \tau, \beta\}$ | |
| Model without popularity (**NO-$\alpha$**) | $b_k + \tau^{t-k+1}$ | $\{\tau, \beta\}$ | $\alpha = 0$ |
| Model without novelty (**NO-$\tau$**) | $\alpha d_{k,t} + b_k + 1$ | $\{\alpha, \beta\}$ | $\tau = $ |
| Model without bias (**NO-bias**) | $\alpha d_{k,t} + \tau^{t-k+1}$ | $\{\alpha, \tau\}$ | $\beta = 0$ |

# Outline

# Parameter estimation
Maximum likelihood

## We can compute the likelihood of the full model

- The likelihood of a set $\Pi := \{\boldsymbol{\pi}_1, \ldots \boldsymbol{\pi}_N\}$ of $N$ trees with respective sizes $|\boldsymbol{\pi}_i|$, $i \in \{1, \ldots N\}$, given the values of $\boldsymbol{\theta}$ can be written as:

$$\mathcal{L}(\boldsymbol{\Pi}|\boldsymbol{\theta}) = \prod_{i=1}^{N} p(\boldsymbol{\pi}_i|\boldsymbol{\theta}) = \prod_{i=1}^{N}\prod_{t=2}^{|\boldsymbol{\pi}_i|} p(\pi_{t,i}|\boldsymbol{\pi}_{(1:t-1),i}, \boldsymbol{\theta}) = \prod_{i=1}^{N}\prod_{t=2}^{|\boldsymbol{\pi}_i|} \frac{\phi(\pi_{t,i})}{Z_{t,i}}$$

- We minimise the negative of the log-likelihood function:

$$-\log\mathcal{L}(\Pi|\boldsymbol{\theta}) = -\sum_{i=1}^{N}\sum_{t=2}^{|\boldsymbol{\pi}_i|} \phi(\pi_{t,i}) - \log Z_{t,i}.$$

# Parameter estimation
## Validation

For each model:

- Choose $\theta^*$ randomly.
- Generate $N$ threads.
- Find estimates $\hat{\theta}$.
- Compute residuals $\theta^* - \hat{\theta}$.
- Repeat for $100$ times.

- Estimation is unbiased.
- Good estimates can be obtained using $N = 500$.

# Parameter estimation
## Model Comparison

For each dataset:

- Select $N = 5 \cdot 10^4$ threads randomly with replacement.
- Find estimates $\hat{\boldsymbol{\theta}}$.
- Compute likelihoods.
- Repeat for $100$ times.

- Model comparison based on likelihoods for each dataset.

# Parameter estimation
Parameter estimates for the different datasets

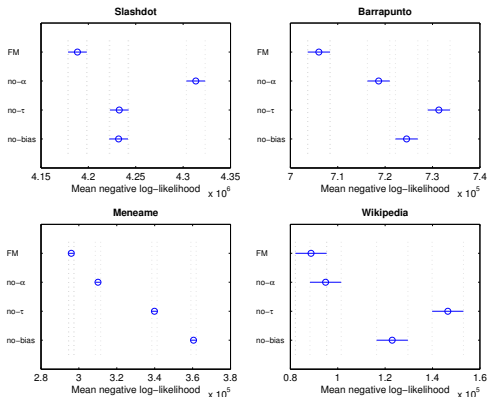| Dataset | $\log \beta$ | | $\alpha$ | | $\tau$ | |
|---------|--------------|--|----------|--|--------|--|
| $N = 50$ | | | | | | |
| SL | 2.39 | (0.17) | 0.31 | (0.02) | 0.98 | (0.02) |
| BP | 0.93 | (0.12) | 0.08 | (0.04) | 0.92 | (0.00) |
| MN | 1.66 | (0.16) | 0.03 | (0.01) | 0.72 | (0.04) |
| WK | $-0.21$ | (0.81) | 0.00 | (0.00) | 0.40 | (0.19) |
| $N = 5000$ | | | | | | |
| SL | **2.39** | (0.01) | **0.31** | (0.01) | **0.98** | (0.00) |
| BP | **0.96** | (0.02) | **0.08** | (0.00) | **0.92** | (0.00) |
| MN | **1.69** | (0.03) | **0.02** | (0.00) | **0.74** | (0.01) |
| WK | **0.39** | (0.22) | **0.00** | (0.00) | **0.60** | (0.01) |

- Bootstrap with $N = 50$ threads already gives good estimates.



Dataset parameters

# Outline

# Growing tree model for discussion threads

### Validation of the model

We calculate the following quantities from the empirical data and from the synthetic threads produced by the model:

- Degrees distribution.
- Subtree sizes distribution.
- Mean node depth versus size.
- Node depths distribution.

- Size of the post $N$ is drawn from the empirical distribution.
- We use model **NO-BIAS** for comparison [Kumar et al. 2010].

# Barrapunto dataset

# Slashdot dataset

# Meneame dataset

# Wikipedia dataset

# Growing tree model for discussion threads

Real cascades:



(a) Slashdot  (b) Barrapunto  (c) Meneame  (d) Wikipedia

Synthetic cascades:



(a) Slashdot  (b) Barrapunto  (c) Meneame  (d) Wikipedia

# Evolution of mean depths and mean widths

FULL MODEL:



NO-BIAS model:

# Theoretical Result
Asymptotics for degree distribution in FM

## It can be shown that ...

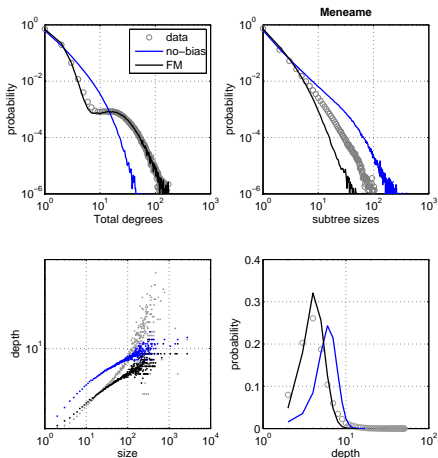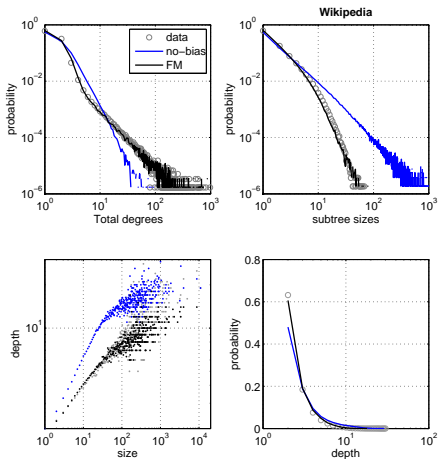- the degree distribution follows (asymptotically) a power-law with exponent 3.
- The parameter $\tau$ does not affect the power-law exponent

## but ...

- formally

$$c_1 x^{-2} \leq P(degree \geq x) \leq c_2 x^{-2}, 0 \leq c_1 \leq c_2$$

- with $\tau$ affecting $c_2$ which is bounded by $\exp(\frac{\tau}{1-\tau})$.
- Thus $\tau$ can affect the fraction of nodes with a degree larger than $x$ by several orders of magnitude.

# Related work
## Galton-Watson branching process

### Idea

- Tree grows level by level.
- Nodes at level $i$ receive a random number of child-nodes at level $i + 1$ (according to a probability distribution).

### Pros

- The model is simple.
- Explains chain letter trees (combined with a selection bias)
  [Golub & Jackson 2010].

### Cons

- Not a generative model.
- Does not capture the order of message creation.

# Conclusions and future work

## Conclusions

- Framework which allows to re-create discussions with similar structural features as real instances.
- Likelihood-based optimisation on the entire cascade evolution.
- Large datasets are not necessary.
- Parameters allow to characterize audience and platform:
  - Same platform : differences between SL and BP.
  - Influence of the interface: MN (flat) characterised by bias.
  - Main difference between news media and WK: popularity.

## Future work

- Include prior authorship structure in model and analysis
- Application to other types of information cascades.

# Bibliography I

V. Gómez, H. J. Kappen, N. Litvak & A. Kaltenbrunner.
*A likelihood-based framework for the analysis of discussion threads.*
World Wide Web Journal

V. Gómez, H. J. Kappen & A. Kaltenbrunner.
*Modeling the structure and evolution of discussion cascades.*
In HT '11, Eindhoven, The Netherlands, 2011. ACM.

B. Golub & M.O. Jackson.
*Using selection bias to explain the observed structure of Internet diffusions.*
Proceedings of the National Academy of Sciences, vol. 107, no. 24, page 10833, 2010.

A. Kaltenbrunner, V. Gómez & V. López.
*Description and Prediction of Slashdot Activity.*
In LA-WEB '07, Santiago de Chile, 2007. IEEE.

R. Kumar, M. Mahdian & M. McGlohon.
*Dynamics of conversations.*
In SIGKDD '10, pages 553–562, New York, USA, 2010. ACM.

D. Laniado, R. Tasso, Y. Volkovich & A. Kaltenbrunner.
*When the Wikipedians talk: Network and tree structure of Wikipedia discussion pages.*
In ICWSM-11 - 5th International AAAI Conference on Weblogs and Social Media. The AAAI Press, 2011.

# Related work II
## T-MODEL [Kumar et al. 2010]

### Features

- Equivalent to model **NO-bias** with an extra parameter to model the death of a discussion.
- Model is illustrated on USENET.
- Authorship model (TI-model).
- Both are independent of the structure.
- Could be build on top of other structural models as well.
- T-model re-creates a power-law relation in the data between size and depth of the discussions, but is not the best model.