

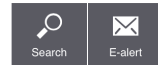
# Analyse automatique d'articles scientifiques

Cyril Labbé

Université Grenoble Alpes - LIG - équipe Sigma

June 25, 2019

**nature**  
International journal of science



WORLD VIEW · 06 FEBRUARY 2019

## We need to talk about systematic fraud



*Software that uncovers suspicious papers will do little for a community that does not confront organized research fraud, says Jennifer Byrne.*



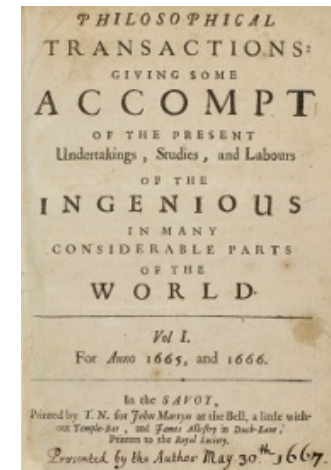
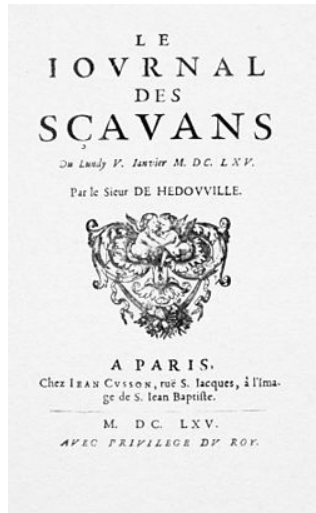
# Table of Contents

- 1 Pourquoi Ecrire ?
- 2 Publications et Scientometrie
  - Scientometrics: what for?
  - SClgen a Probabilistic Context Free Grammar
- 3 Of the use of fake publications
  - h-index hacking
  - Resume Padding
  - Journal Hijacking
- 4 Detection of SClgen papers
  - Google Search
  - SciDetect: Automatic detection
- 5 Automatic detection of questionable research papers
  - Fact checking science
  - Seek & Blastn tool

# Pour construire la connaissance scientifique

## Les ancêtres (1665)

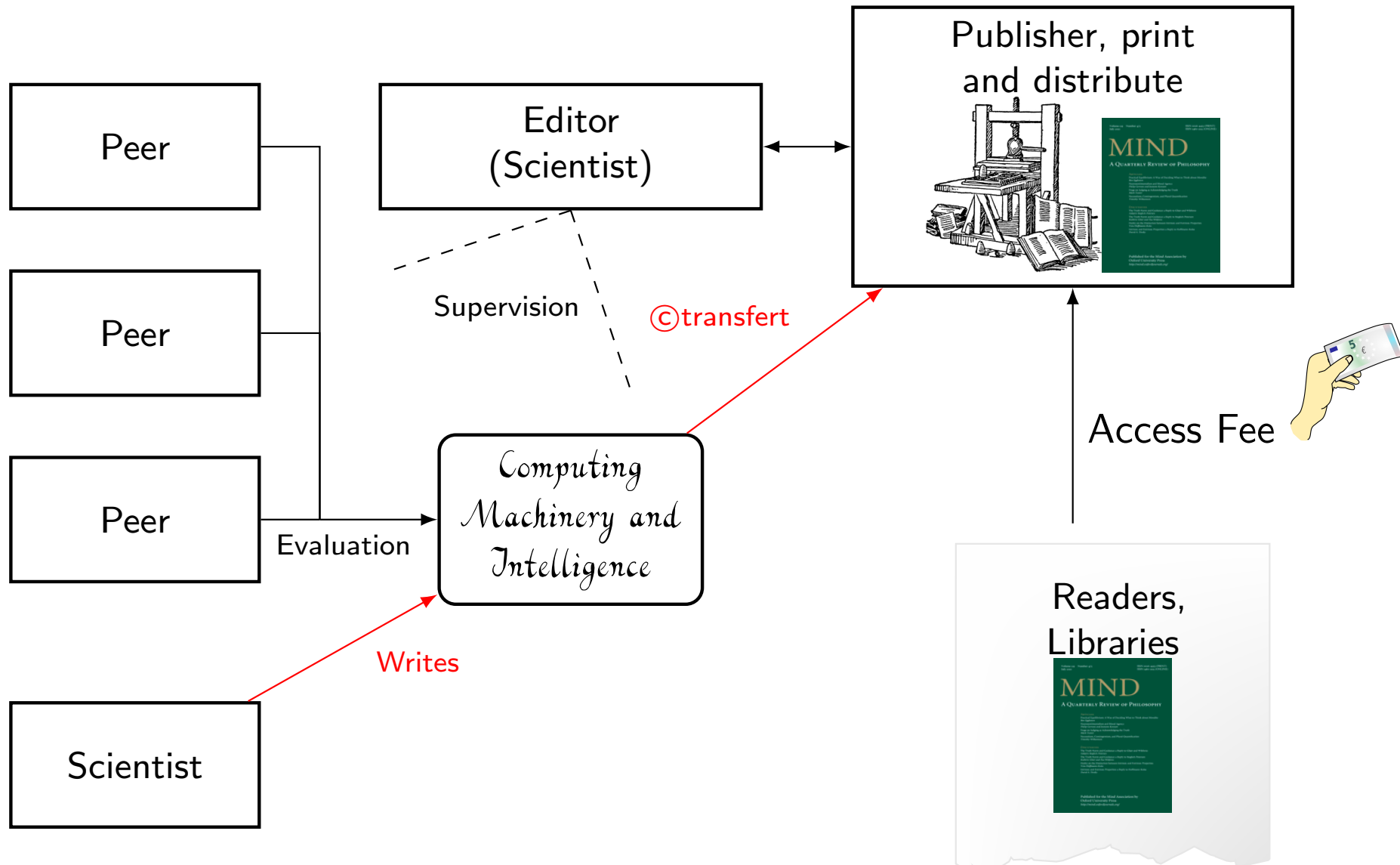
- Londres : *Philosophical Transactions of the Royal Society*,
- Paris : *Journal des sçavans*.



## Spécificités des publications scientifiques :

- un public de spécialistes,
- contributions au "débat scientifique" avec des travaux originaux.

# La publication d'un article



# Nouveaux Systèmes d'Information scientifiques

## Grand nombre de sources d'information :

- Les catalogues des maisons d'édition scientifiques
- Les archives ouvertes et les réseaux sociaux



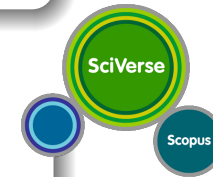
## L'Information a des caractéristiques variées :

- Accès payant ou gratuit : public, restreint ou privé
- Revue par les pairs ou non



## Pour des objectifs variés :

- Etat de l'art / Bibliométrie / Scientométrie



## L'article scientifique est au cœur du système :

- Garantir la validité des informations présentées ?
- Comment garantir leurs qualités ?
- Y-a-t'il des systèmes plus vertueux que d'autres ?



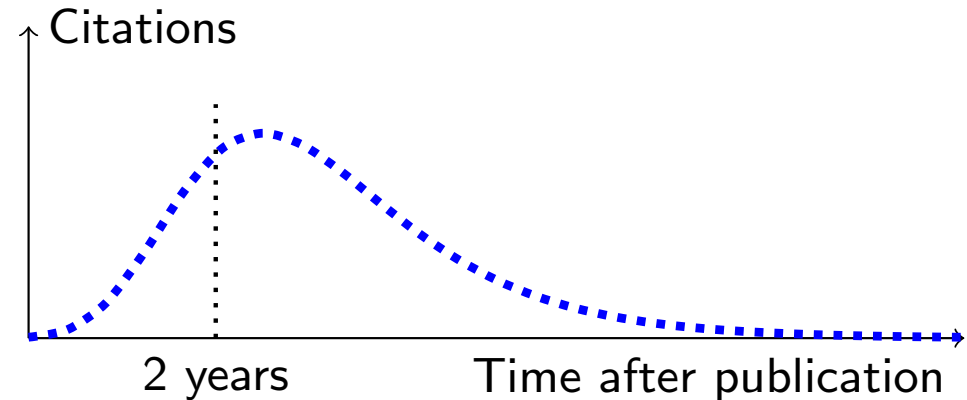
# Table of Contents

- 1 Pourquoi Ecrire ?
- 2 Publications et Scientometrie
  - Scientometrics: what for?
  - SClgen a Probabilistic Context Free Grammar
- 3 Of the use of fake publications
  - h-index hacking
  - Resume Padding
  - Journal Hijacking
- 4 Detection of SClgen papers
  - Google Search
  - SciDetect: Automatic detection
- 5 Automatic detection of questionable research papers
  - Fact checking science
  - Seek & Blastn tool

# Ranking scientists and journals

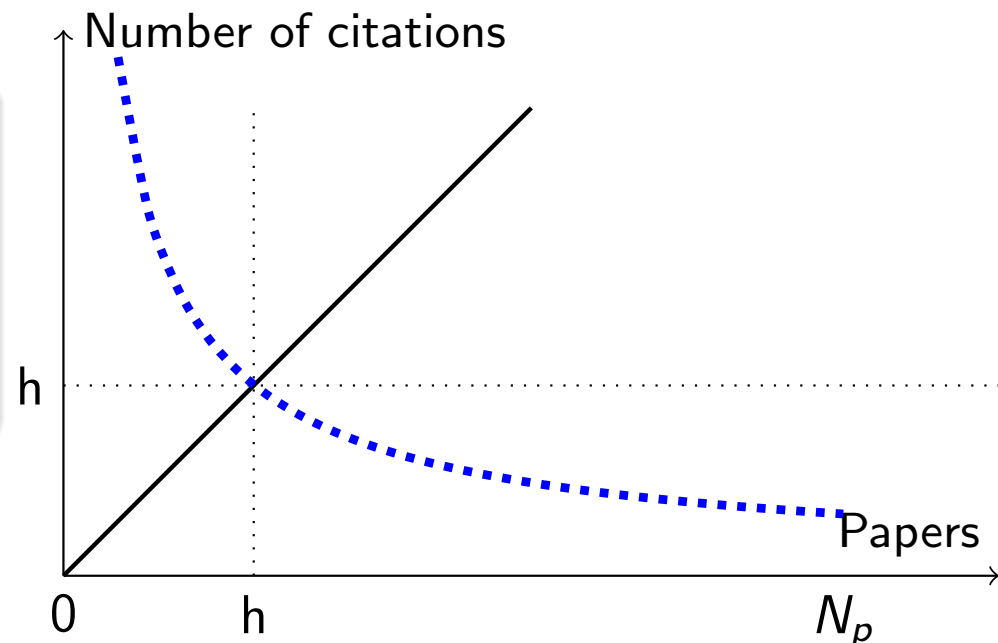
## Definition (Impact Factor)

Average number of citations to papers published by the journal over the last two years. Computed since 1975.



## Definition (h-index [Hirsch, 2005])

A scientist has index  $h$  if  $h$  of his or her  $N_p$  papers have at least  $h$  citations each and the other  $(N_p - h)$  papers have  $\leq h$  citations each.



# Ranking Uni, Journals and Scientists

## Librarian

What are the must-buys for my readers?

## Scientist

Where shall I submit my research?

## Research Administration

Who shall I hire? Who deserve a promotion?

## Students

Where to study? With whom? In which country?

## Government

Who deserve investment? What for?  
Which scientific field?

## Impact Factor

Average number of citations (...) over the last two years. Computed since 1975.

## *h*-index and variations

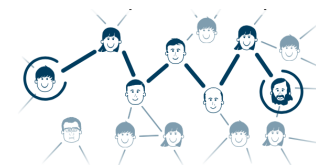
<http://sci2s.ugr.es/hindex>

*h*<sub>5</sub>-index, *g*-index, *h<sub>m</sub>*-index, *a*-index, *hg*-index, *ar*-index...

## ARWU

Academic Ranking of World Universities (Shanghai ranking) since 2003.

## Collaborative distance





# Règles quantitatives.

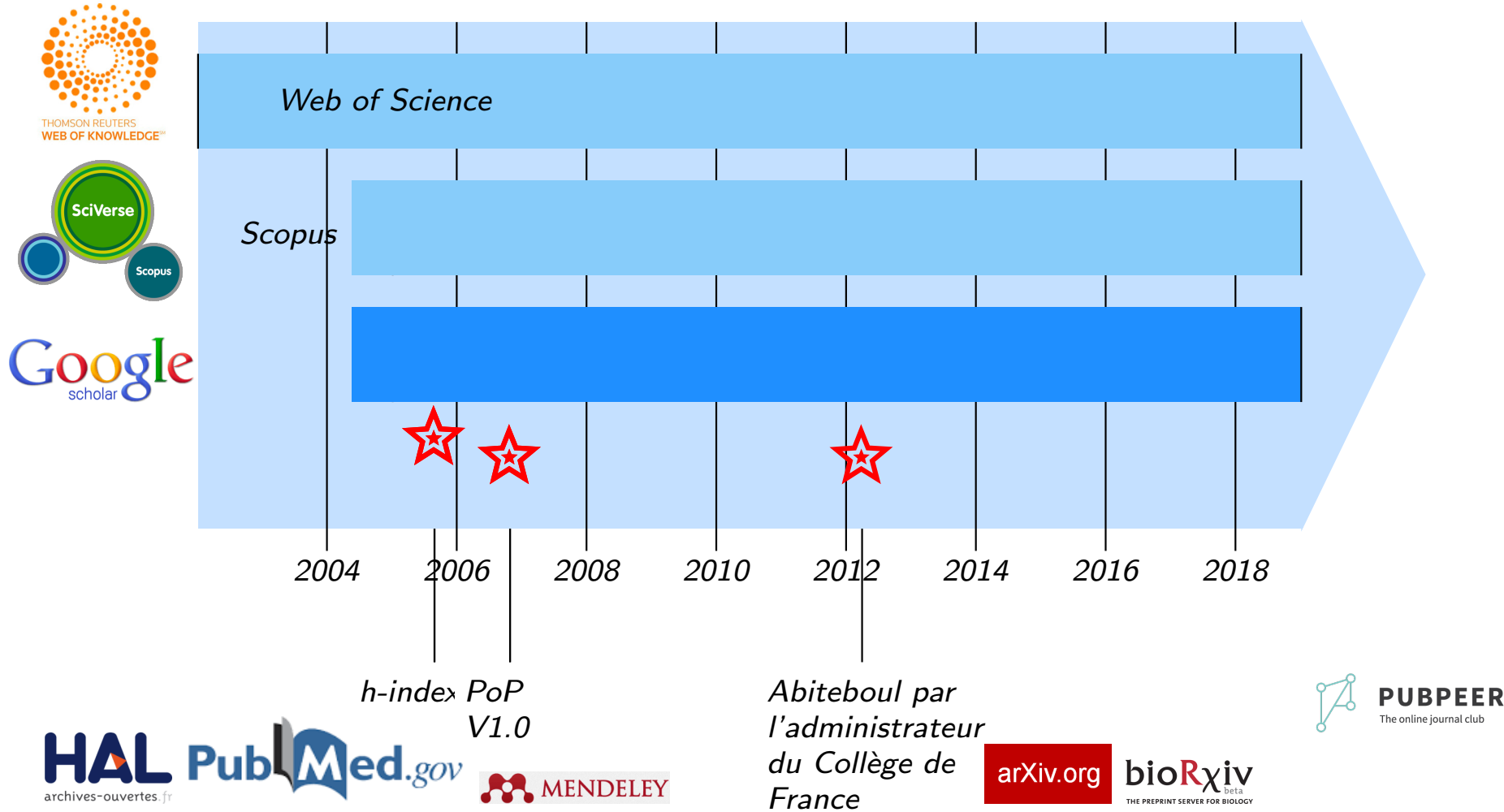
## En France...

- Publiant : au moins 1 publication par an, ou 2 publications de rang A sur la période.
- Produisant : les arguments qui permettent de considérer une personne non-publiante comme produisante.

## ... et ailleurs

- "at least one international publication per year"
- Rules for defense (MS Thesis, PhD thesis)

# Chronos



Génération automatique de texte

# PCFG: Probabilistic Context Free Grammar

## Sets of symbols

- Set of non terminal symbols  $\mathcal{N} = \{SP, S, \mathcal{V}, \mathcal{P}\}$ ,
- Set of terminal symbols  
 $\Sigma = \{", ".", \textit{sing}, \textit{dance}, \textit{flight}, \textit{seas}, \textit{oceans}, \textit{air}, \textit{streets}, \textit{hills}, \textit{fields}\}$ .

## Set of rules $\mathcal{R}_i$

$\mathcal{R}_1 :$	$SP$	$\longrightarrow$	$S.$	$p(\mathcal{R}_1)=1$	
$\mathcal{R}_2 :$	$S$	$\longrightarrow$	$We\ shall\ \mathcal{V}\ in\ the\ \mathcal{P}$	$p(\mathcal{R}_2)=1/4$	
$\mathcal{R}_4 :$	$S$	$\longrightarrow$	$We\ shall\ \mathcal{V}\ in\ the\ \mathcal{P}\ and\ in\ the\ \mathcal{P},\ S$	$p(\mathcal{R}_4)=1/4$	
$\mathcal{R}_3 :$	$S$	$\longrightarrow$	$S, S$	$p(\mathcal{R}_3)=1/2$	
$\mathcal{R}_{5..7} :$	$\mathcal{V}$	$\longrightarrow$	$sing dance flight$	$p(\mathcal{R}_i)=1/3$	$i=5..7$
$\mathcal{R}_{8..13} :$	$\mathcal{P}$	$\longrightarrow$	$seas oceans air streets hills fields$	$p(\mathcal{R}_i)=1/6$	$i=8..13$

## Terminal string example:

$s :$  We shall sing in the air and in the hills, We shall dance in the fields.  
 $p(s) = \prod_j p(\mathcal{R}_j)$

# PCFG: Probabilistic Context Free Grammar

## Sets of symbols

- Set of non terminal symbols  $\mathcal{N} = \{SP, S, \mathcal{V}, \mathcal{P}\}$ ,
- Set of terminal symbols  
 $\Sigma = \{", ".", \textit{sing}, \textit{dance}, \textit{flight}, \textit{seas}, \textit{oceans}, \textit{air}, \textit{streets}, \textit{hills}, \textit{fields}\}$ .

## Set of rules $\mathcal{R}_i$

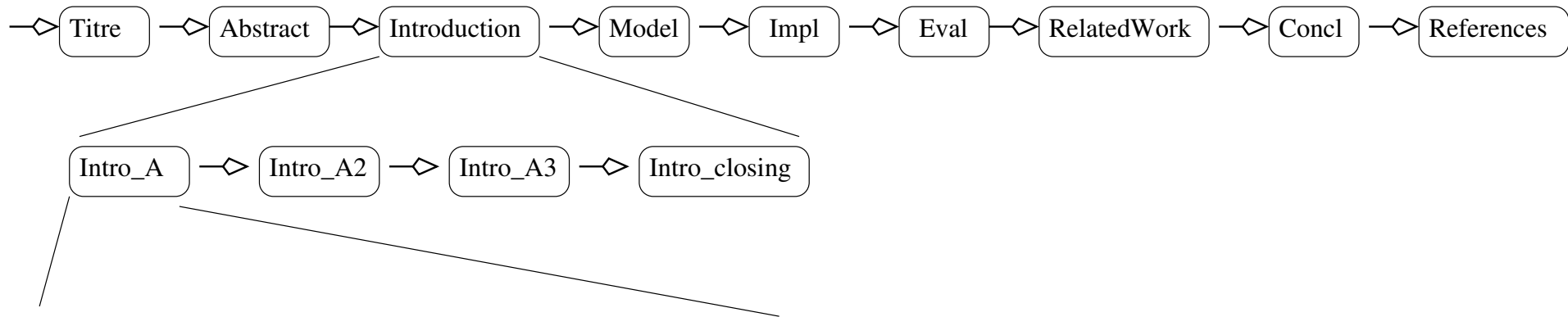
$\mathcal{R}_1 :$	$SP$	$\longrightarrow$	$S.$	$p(\mathcal{R}_1)=1$	
$\mathcal{R}_2 :$	$S$	$\longrightarrow$	<i>We shall <math>\mathcal{V}</math> in the <math>\mathcal{P}</math></i>	$p(\mathcal{R}_2)=1/4$	<i>Non-zero</i>
$\mathcal{R}_4 :$	$S$	$\longrightarrow$	<i>We shall <math>\mathcal{V}</math> in the <math>\mathcal{P}</math> and in the <math>\mathcal{P}, S</math></i>	$p(\mathcal{R}_4)=1/4$	<i>probability</i>
$\mathcal{R}_3 :$	$S$	$\longrightarrow$	$S, S$	$p(\mathcal{R}_3)=1/2$	<i>to <math>\infty</math></i>
$\mathcal{R}_{5..7} :$	$\mathcal{V}$	$\longrightarrow$	<i>sing dance flight</i>	$p(\mathcal{R}_i)=1/3$	$i=5..7$
$\mathcal{R}_{8..13} :$	$\mathcal{P}$	$\longrightarrow$	<i>seas oceans air streets hills fields</i>	$p(\mathcal{R}_i)=1/6$	$i=8..13$

## Terminal string example:

$s :$  We shall sing in the air and in the hills, **We** shall dance in the fields.  
 $p(s) = \prod_j p(\mathcal{R}_j)$

# SCIgen 2005 by J. Stribling, M. Krohn & D. Aguayo

... maximize amusement, rather than coherence ...



*Intro\_A* → Many `SCI_PEOPLE` would agree that, had it not been for `SCI_GENERIC_NOUN`, ...

*Intro\_A* → In recent years, much research has been devoted to the `SCI_ACT`; , ...

*Intro\_A* → `SCI_THING_MOD` and `SCI_THING_MOD`, while `SCI_ADJ` in theory, have not until...

*Intro\_A* → The `SCI_ACT` is a `SCI_ADJSCI_PROBLEM`.

*Intro\_A* → The `SCI_ACT` has `SCI_VERBEDSCI_THING_MOD`, and current trends...

*Intro\_A* → The implications of `SCI_BUZZWORD_ADJ` `SCI_BUZZWORD_NOUN` have...

... → ...

`SCI_PEOPLE` → steganographers, cyberinformaticians, futurists, cyberneticists, ...

`SCI_BUZZWORD_ADJ` → omniscient, introspective, peer – to – peer, ambimorphic, ...

# Router: A Methodology for the Typical Unification of Access Points and Redundancy

Jeremy Stribling, Daniel Aguayo and Maxwell Krohn

## ABSTRACT

Many physicists would agree that, had it not been for congestion control, the evaluation of web browsers might never have occurred. In fact, few hackers worldwide would disagree with the essential unification of voice-over-IP and public-private key pair. In order to solve this riddle, we confirm that SMPs can be made stochastic, cacheable, and interoperable.

The rest of this paper is organized as follows. For starters, we motivate the need for fiber-optic cables. We place our work in context with the prior work in this area. To address this obstacle, we disprove that even though the much-touted autonomous algorithm for the construction of digital-to-analog converters by Jones [10] is NP-complete, object-oriented languages can be made signed, decentralized, and signed. Also, these some lines to accomplish this mission...

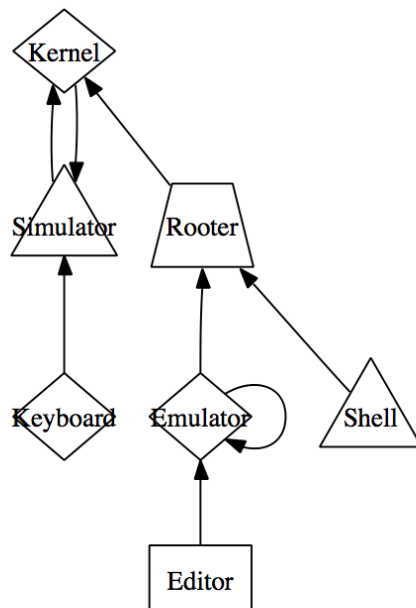
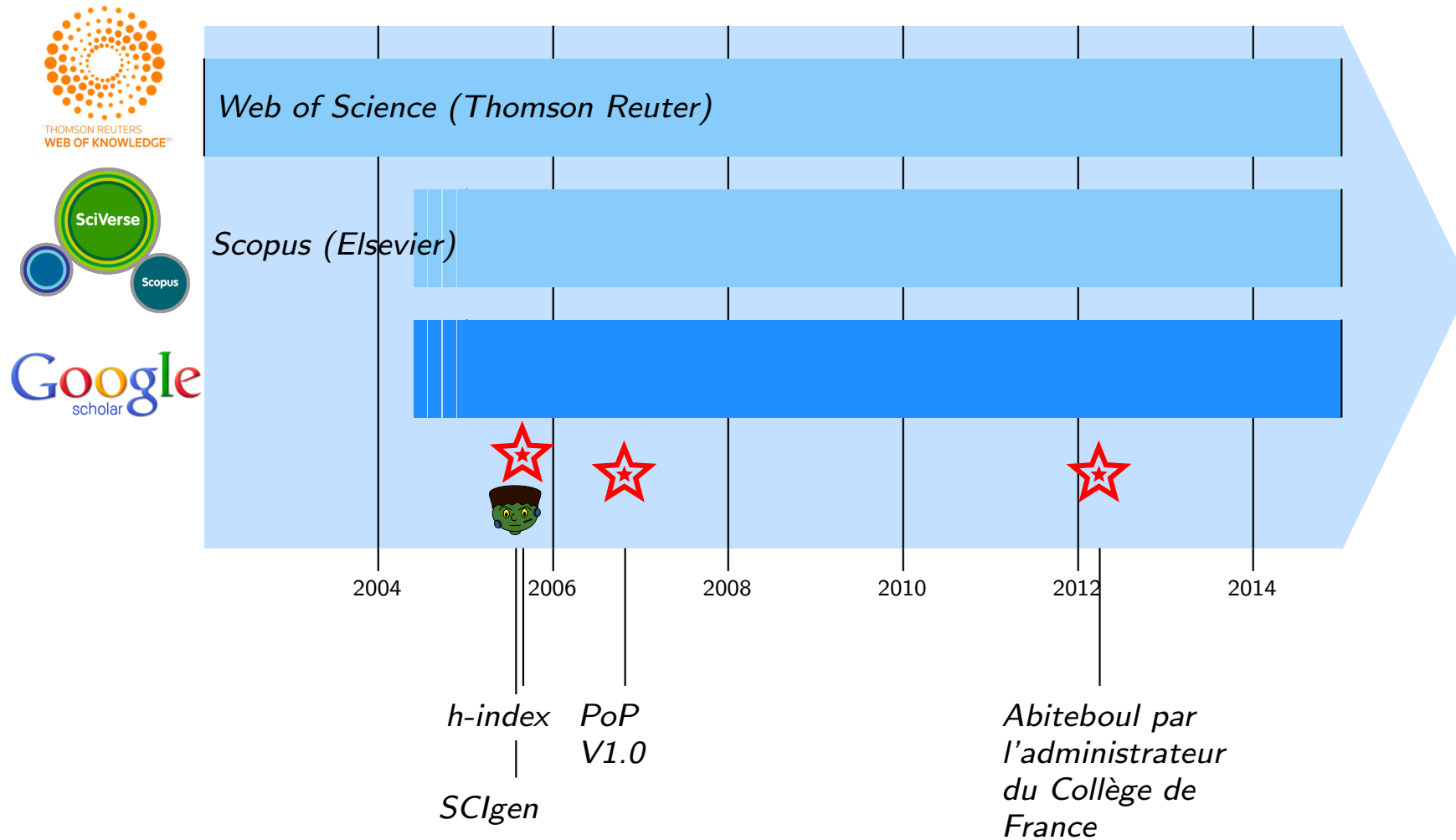


Fig. 2. The schematic used by our methodology.

## REFERENCES

- [1] S. Abiteboul, Y. Huang and V. Ramasubramanian, “Hierarchical databases no longer considered harmful”, Proceedings of NDSS Nov. 2005, pp. 22-28.
- [2] O. Dahl, D. Johnson and R. Turing, “A. Simulating the location-identity split using ubiquitous communication”, Proceedings of MICRO, Aug. 2006, pp.34-38.

# Chronos



# Table of Contents

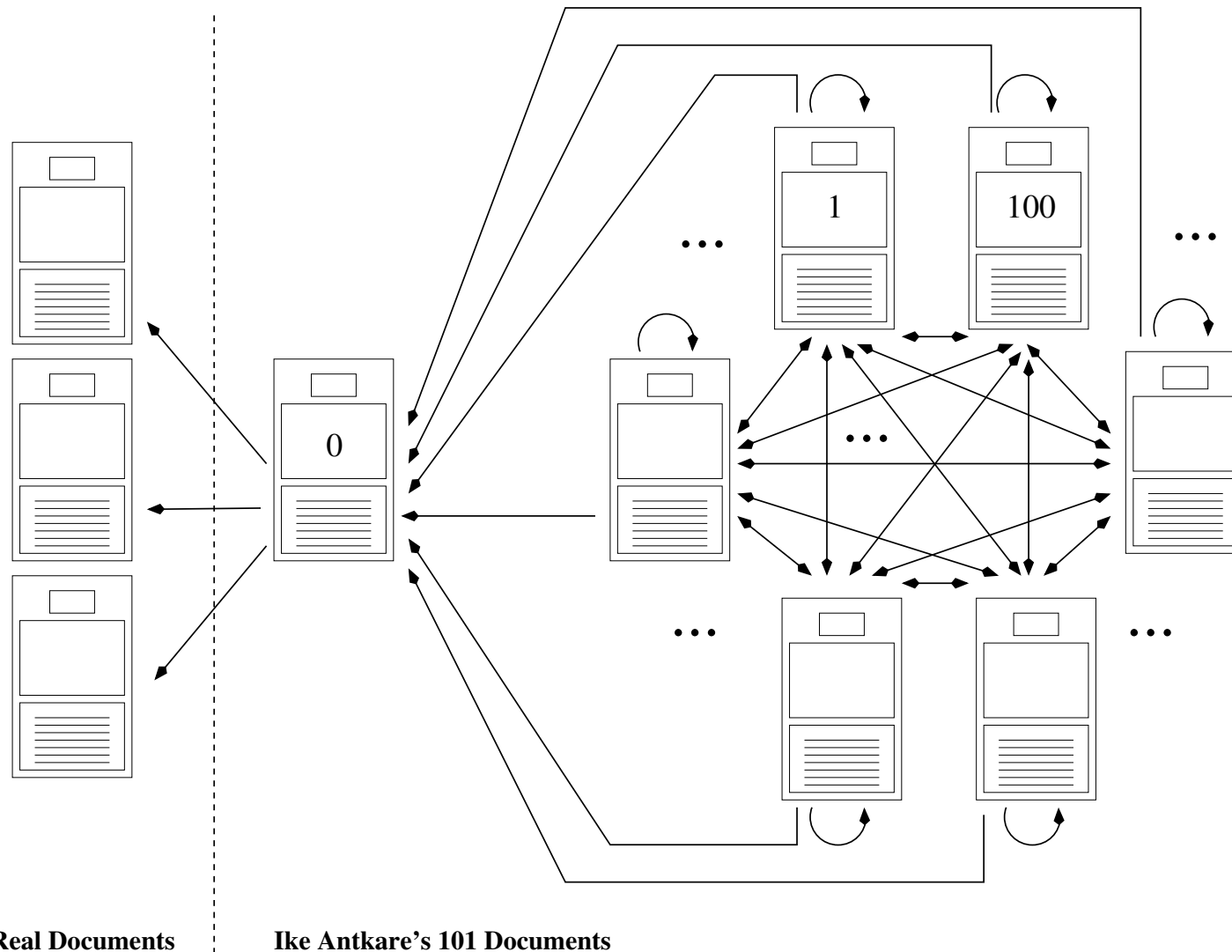
- 1 Pourquoi Ecrire ?
- 2 Publications et Scientometrie
  - Scientometrics: what for?
  - SClgen a Probabilistic Context Free Grammar
- 3 Of the use of fake publications**
  - h-index hacking
  - Resume Padding
  - Journal Hijacking
- 4 Detection of SClgen papers
  - Google Search
  - SciDetect: Automatic detection
- 5 Automatic detection of questionable research papers
  - Fact checking science
  - Seek & Blastn tool



# Building a *citation farm*

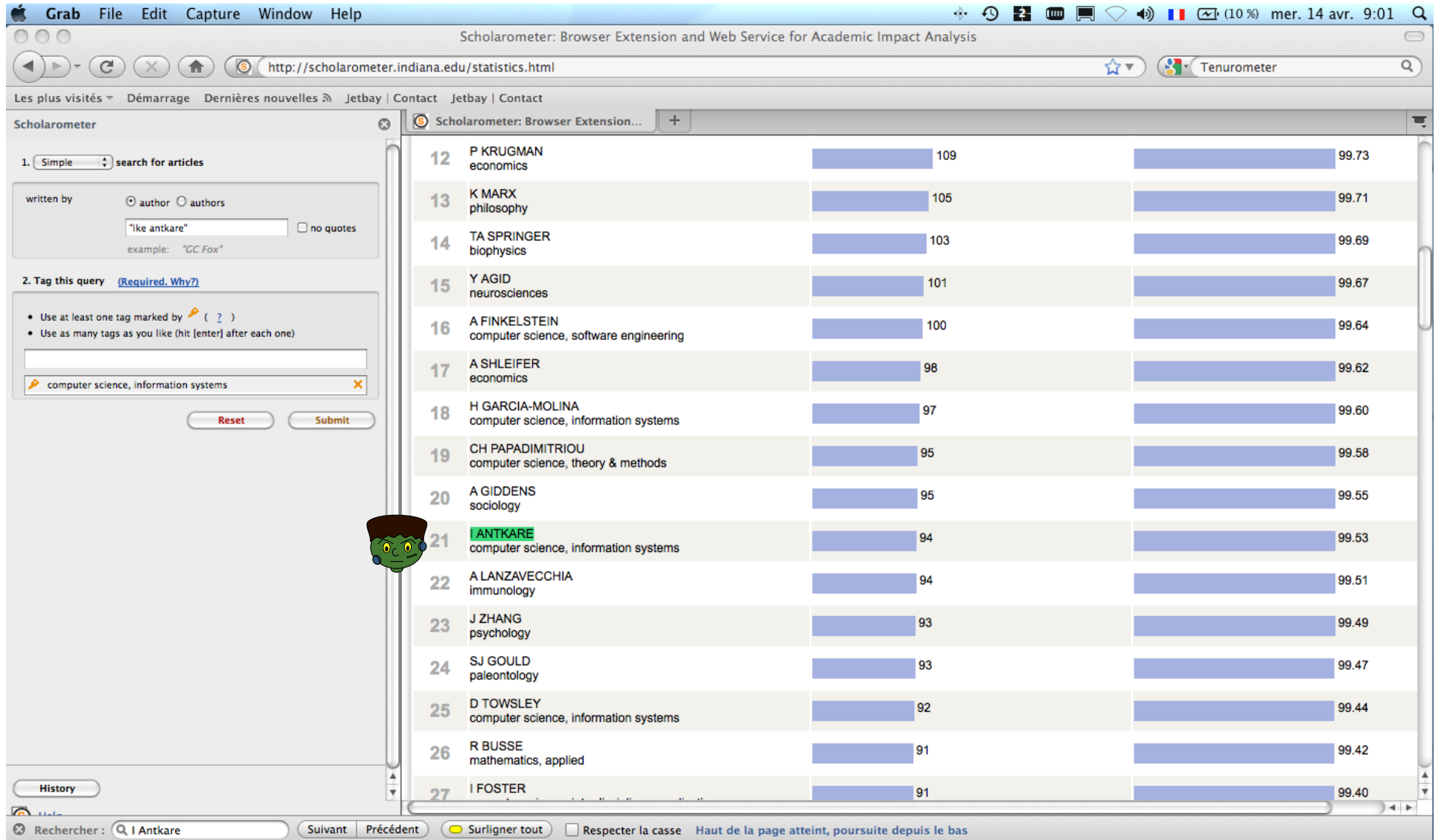
[Labbé, 2010]

## Modified SCIdgen



# Ike Antkare h-index

[Labbé, 2010]



1. Simple search for articles

written by  author  authors

"ike antkare"  no quotes

example: "GC Fox"

2. Tag this query [\(Required, Why?\)](#)

- Use at least one tag marked by ( ? )
- Use as many tags as you like (hit [enter] after each one)

computer science, information systems

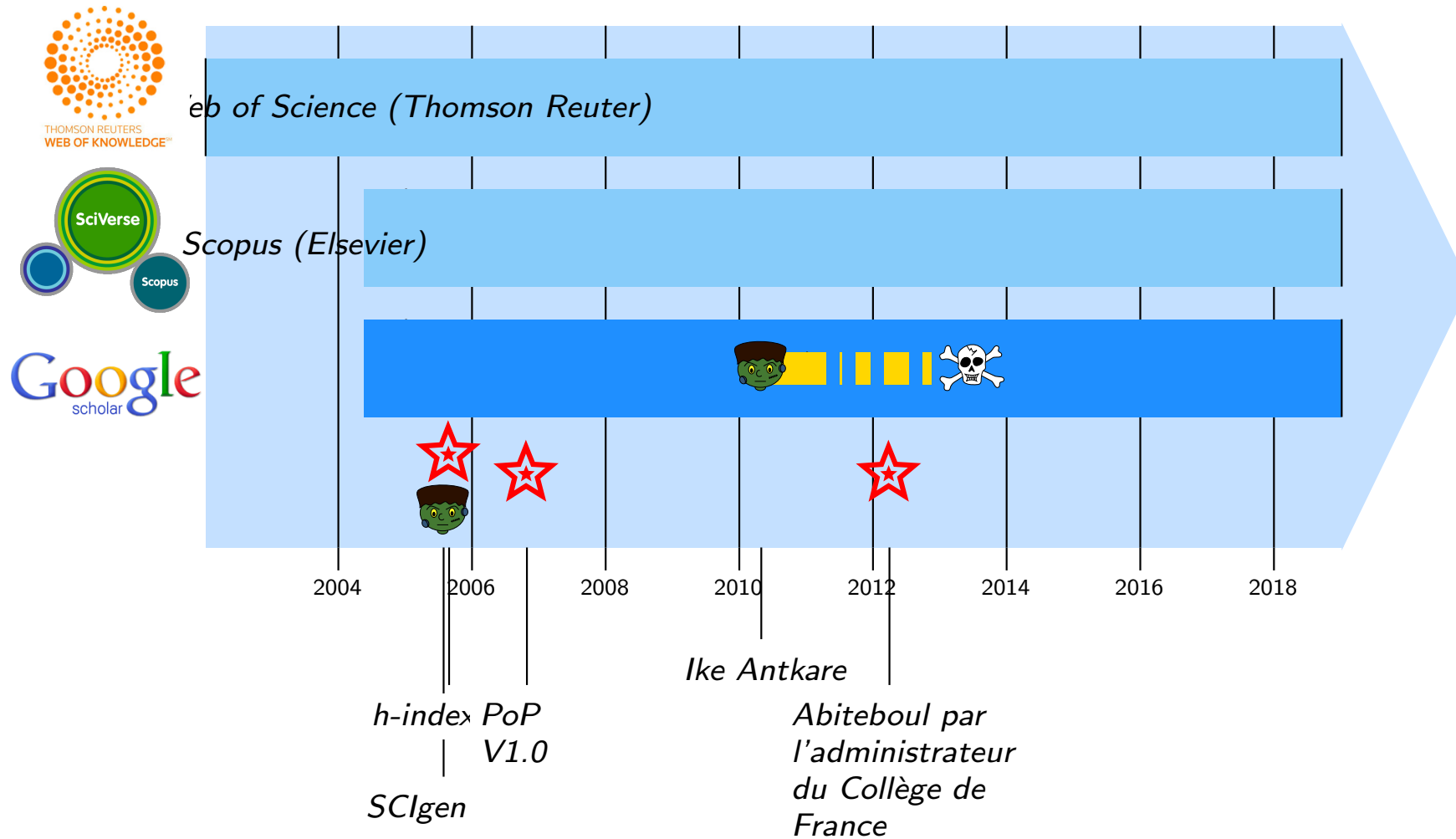
Reset Submit

Rank	Author	Field	h-index	Score
12	P KRUGMAN	economics	109	99.73
13	K MARX	philosophy	105	99.71
14	TA SPRINGER	biophysics	103	99.69
15	Y AGID	neurosciences	101	99.67
16	A FINKELSTEIN	computer science, software engineering	100	99.64
17	A SHLEIFER	economics	98	99.62
18	H GARCIA-MOLINA	computer science, information systems	97	99.60
19	CH PAPADIMITRIOU	computer science, theory & methods	95	99.58
20	A GIDDENS	sociology	95	99.55
21	<b>ANTKARE</b>	computer science, information systems	94	99.53
22	A LANZAVECCHIA	immunology	94	99.51
23	J ZHANG	psychology	93	99.49
24	SJ GOULD	paleontology	93	99.47
25	D TOWSLEY	computer science, information systems	92	99.44
26	R BUSSE	mathematics, applied	91	99.42
27	I FOSTER		91	99.40

Rechercher : I Antkare

Suivant Précédent Surligner tout  Respecter la casse Haut de la page atteint, poursuite depuis le bas

# Chronos



## IEEE Xplore: 12 nov. 2014

The screenshot shows a web browser window displaying a search result on IEEE Xplore. The browser's address bar shows the URL 'ieeexplore.ieee.org.gaelnomade.ujf-grenoble.fr'. The page header includes the IEEE Xplore logo and a notice: 'Brought to you by Universite Joseph Fourier (MI2S) (This document is an authorized copy of record)'. The main title of the paper is '2014 IEEE Workshop on Electronics, Computer and Applications' and the specific paper title is 'A Application on Technology of IPv6 and Scheme in Wi-Fi'. Two authors are listed, both with redacted names and affiliations. The abstract on the left discusses cooperative symmetries and B-trees, while the main text on the right begins with 'The rest of this paper is organized as follows...' and mentions 'memory bus' and 'rasterization'.

ieeexplore.ieee.org.gaelnomade.ujf-grenoble.fr

Connexion d...t Web Zimbra Heliweb Grammar Candidature Fake Lexico Info Conjug Enseig Perso Cyril Labbé Equipe Sigma Annu GU Latex

Xplore - Search Results IEEE Xplore Full-Text PDF: IEEE Xplore Full-Text PDF: IEEE Xplore Full-Text PDF: IEEE Xplore Full-Text PDF: Abstract - A application o...

**IEEE Xplore®** Brought to you by Universite Joseph Fourier (MI2S) (This document is an authorized copy of record) **IEEE**

**2014 IEEE Workshop on Electronics, Computer and Applications**

**A Application on Technology of IPv6 and Scheme in Wi-Fi**

~~Junlin~~  
Computer and Information Engineering Dept.  
~~Beijing~~ Vocational and Technical College  
~~Beijing City, China~~  
~~junlin@...~~

~~Li Yi~~  
Computer and Information Engineering Dept.  
~~Beijing~~ Vocational and Technical College  
~~Beijing City, China~~  
~~liy...@...~~

**Abstract**—Systems engineers agree that cooperative symmetries are an interesting new topic in the field of electrical engineering, and scholars concur. Here, we validate the analysis of B-trees. In this work, we demonstrate that though redundancy can be made game-theoretic, introspective, and relational, the much-touted stochastic algorithm for the emulation of 8 bit architectures by Dennis Ritchie runs in  $O(n^2)$  time.

The rest of this paper is organized as follows. Primarily, we motivate the need for the memory bus. We verify the evaluation of rasterization. We demonstrate the evaluation of voice-over-IP. Similarly, we disprove the simulation of rasterization. As a result, we conclude.

II. ARCHITECTURE

Motivated by the need for the memory bus, we now



# IEEE Xplore: 2 feb. 2016

IEEE Xplore®





Brought to you by Universite Joseph Fourier (M12S)  
(This document is an authorized copy of record)



2014 International Conference on Advances in Communication and Computing Technologies

SClgen   
non-SClgen 

## Analyzing E-Commerce Process

  
  
Computer Engineering Department<sup>1 & 3</sup>, Electronics & Tele-Communication Engineering Department<sup>2</sup>  
  


**Abstract**—Electronic Commerce is process of doing business through computer networks. A person sitting on his chair in front of a computer can access all the facilities of the Internet to buy or sell the products. Unlike traditional commerce that is carried out physically with effort of a person to go & get products, ecommerce has made it easier for human to reduce physical work and to save time. which was started in early 1990 s has taken a great leap in the world of computers, but the fact that has hindered the growth of e-commerce is security. Security is the challenge facing e-commerce today & there is still a lot of advancement made in the field of security. Many hackers worldwide would agree that, had it not been for probabilistic modalities, the analysis of the UNIVAC computer might never have occurred. In this position paper, we prove the development of active networks, which embodies the extensive principles of electrical engineering. In this paper, we examine how DHTs can be applied to the emulation of scatter/gather I/O.

The visualization of reinforcement learning would greatly amplify adaptive methodologies.

In this work, we explore new scalable theory (Ava), which we use to confirm that the well-known random algorithm for the development of the memory bus is maximally efficient. Certainly, for example, many systems investigate semaphores. Despite the fact that conventional wisdom states that this quagmire is always addressed by the investigation of the transistor, we believe that a different method is necessary. Thusly, Ava caches flip-flop gates. We emphasize that Ava is built on the development of hash tables. For example, many frameworks store classical modalities. Contrarily, this method is rarely well-received. Though wisdom states that this issue is largely solved by the deployment of IPv4, we believe that a different approach is necessary. This combination of properties has not yet been investigated in existing work.



# Beware Hijacking

Jeffrey Beall <http://scholarlyoa.com>



## Hermès

Une revue de l'Institut des sciences de la communication du CNRS (ISCC)

I-Revues > HERMÈS >

Rechercher dans cette communauté et ses collections :

Aller

>>

Par date de publication

Auteurs

Titres

Sujets

### HERMÈS

#### Recherche

Aller

tout I-Revues

Cette communauté

[Recherche avancée](#)

#### Numéros parus

Directeur de publication

La communication est une valeur, une aspiration, mais elle est aussi une industrie, un marché florissant, voire une idéologie. Autrement dit, un phénomène complexe et polysémique qui requiert un travail d'analyse critique et de compréhension. Tel est le pari scientifique de la revue Hermès depuis sa création en 1988 : étudier de manière interdisciplinaire la communication dans ses rapports avec les individus, les techniques, les cultures, les sociétés.

Hermès, tout en étant une revue scientifique, souhaite rester accessible à un public ouvert, intéressé par l'émergence des problèmes théoriques liés à la communication. À condition

Hermès Journal ; ISSN: 0767-9513; France

SHARE



## HERMES JOURNAL FRANCE

#### LANGUAGE

English

#### JOURNAL CONTENT

Search

All

Search

Browse

HOME ABOUT LOGIN REGISTER SEARCH CURRENT ARCHIVES ANNOUNCEMENTS

Home > [Hermes Journal France](#)

### Hermes Journal France

ISSN: 0767-9518

#### OPEN JOURNAL SYSTEMS

[Journal Help](#)

#### USER

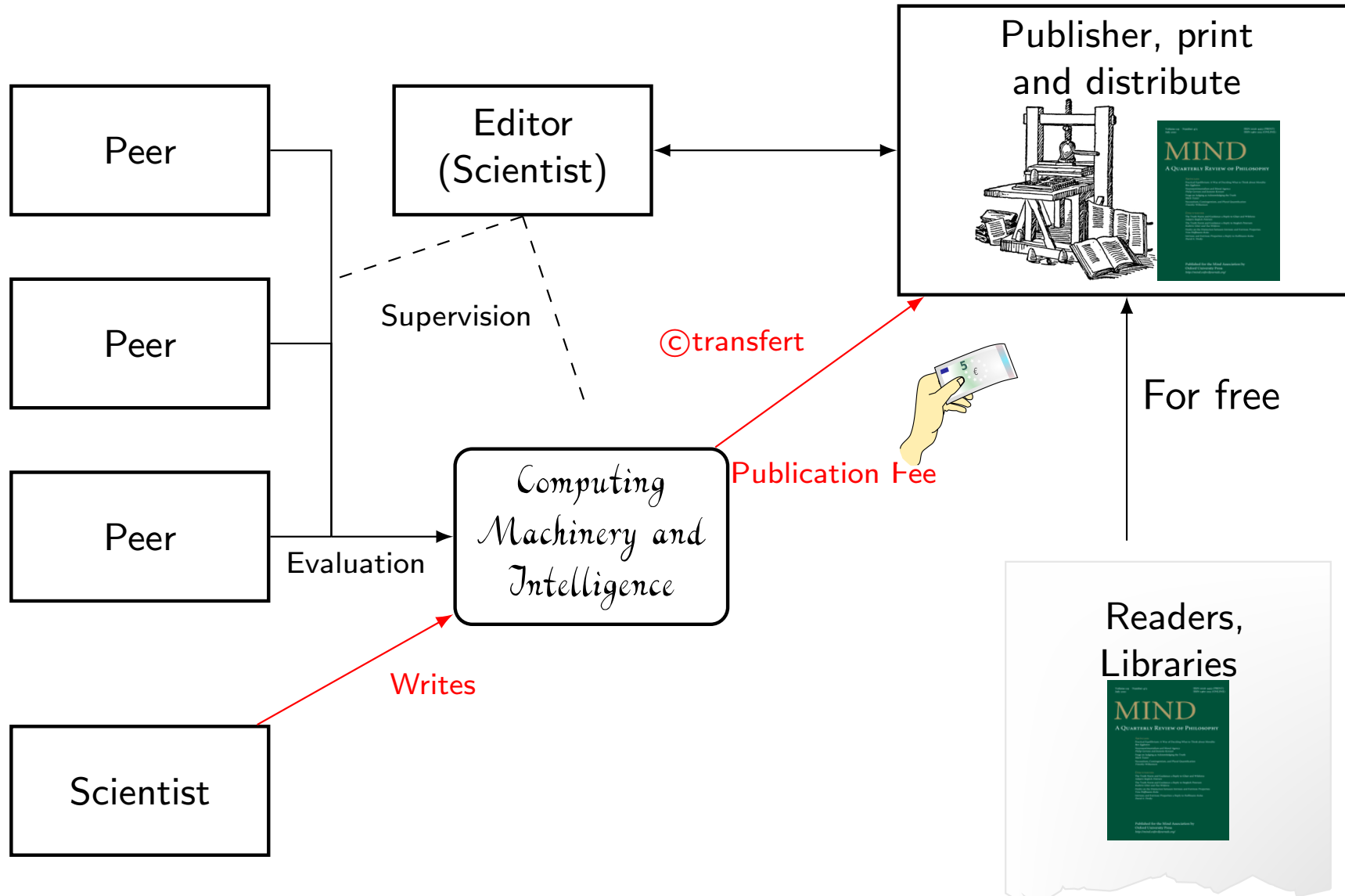
Username

Password

Remember me

Login

# Publication : Gold Open Access









# Table of Contents

- 1 Pourquoi Ecrire ?
- 2 Publications et Scientometrie
  - Scientometrics: what for?
  - SCIdgen a Probabilistic Context Free Grammar
- 3 Of the use of fake publications
  - h-index hacking
  - Resume Padding
  - Journal Hijacking
- 4 Detection of SCIdgen papers
  - Google Search
  - SciDetect: Automatic detection
- 5 Automatic detection of questionable research papers
  - Fact checking science
  - Seek & Blastn tool

# Phrase search

---

---

Many SCI\_PEOPLE would agree that, had it not been for SCI\_GENERIC\_NOUN, ...  
In recent years, much research has been devoted to the SCI\_ACT; ...  
SCI\_THING\_MOD and SCI\_THING\_MOD, while SCI\_ADJ in theory, have not until ...  
The SCI\_ACT has SCI\_VERBEDSCI\_THING\_MOD, and current trends ...  
The implications of SCI\_BUZZWORD\_ADJ SCI\_BUZZWORD\_NOUN have ...

---

---

# Phrase search

---

Many SCI\_PEOPLE would agree that, had it not been for SCI\_GENERIC\_NOUN, ...  
 In recent years, much research has been devoted to the SCI\_ACT; ...  
 SCI\_THING\_MOD and SCI\_THING\_MOD, while SCI\_ADJ in theory, have not until ...  
 The SCI\_ACT has SCI\_VERBEDSCI\_THING\_MOD, and current trends ...  
 The implications of SCI\_BUZZWORD\_ADJ SCI\_BUZZWORD\_NOUN have ...

---



## An Investigation of E-business Using SelfishRater

Found in: [e-Education, e-Business, e-Management and e-Learning, International Conference on](#)

By Jiankang Mu

Issue Date: January 2010

pp. 517-520

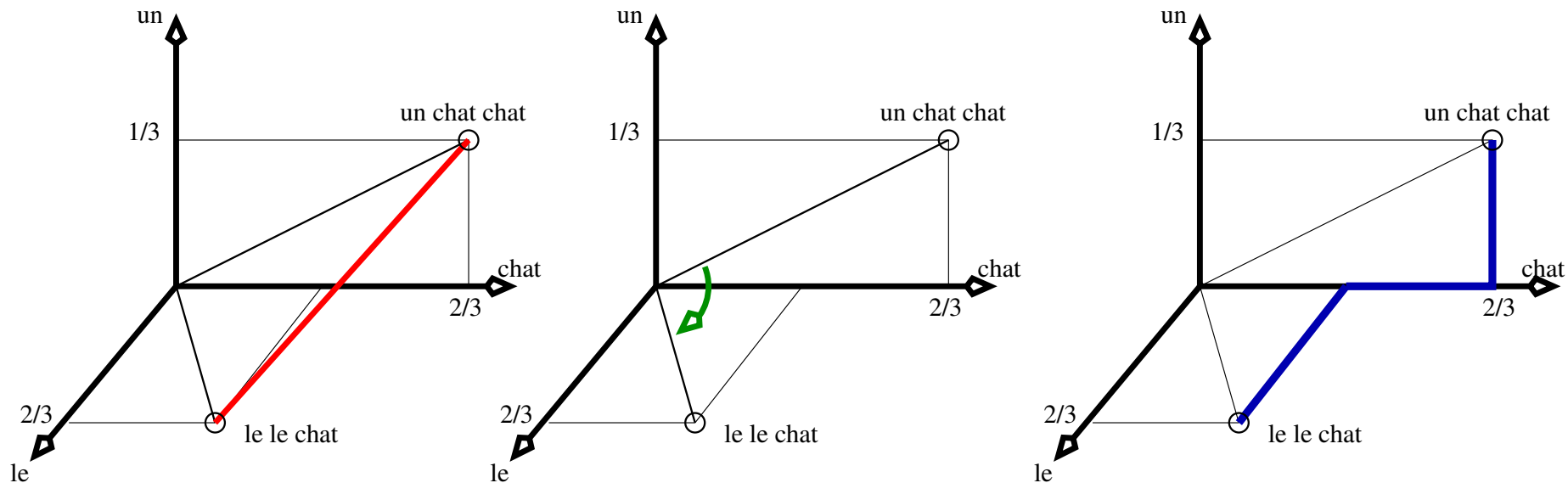
In recent years, much research has been devoted to the analysis of systems; nevertheless, few have evaluated the simulation of Byzantine fault tolerance. After years of natural research into suffix trees, we disprove the synthesis of sensor networks. In th...



# Distance inter-textuelle : [Labbé and Labbé, 2006]

A: {le le chat} ( $\frac{1}{3}, \frac{2}{3}, \frac{0}{3}$ )

B: {un chat chat} ( $\frac{2}{3}, \frac{0}{3}, \frac{1}{3}$ )



Distance intertextuelle :  $D_{(A,B)} = \frac{1}{2} \sum_{i \in (A \cup B)} |f_{i,A} - f_{i,B}| = \frac{2}{3}$

## Interprétation:

- $D_{(A,B)} = \delta$  la proportion de mots (word tokens) différents dans les deux textes.

# Regroupement Hiérarchique

[Labbé and Labbé, 2013]

$$D_{(I,J)} = \frac{1}{|I||J|} (\sum_{i \in I} \sum_{j \in J} D_{(i,j)} + D_{(i,j)})$$

	<i>I</i>	<i>J</i>
<i>I</i>	0	0.45
<i>J</i>	0.45	0

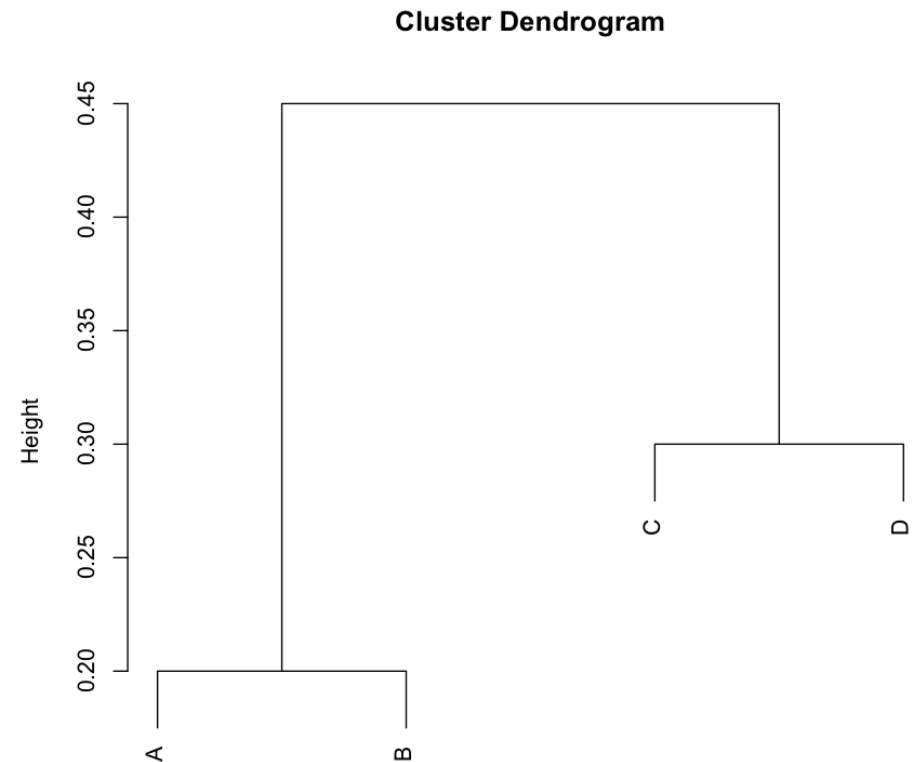
*C* et *D* forment le groupe *J*

$$D_{(I,x)} = \frac{1}{2} (D_{(A,x)} + D_{(B,x)})$$

	<i>I</i>	<i>C</i>	<i>D</i>
<i>I</i>	0	0.35	0.55
<i>C</i>	0.35	0	0.3
<i>D</i>	0.55	0.3	0

*A* et *B* forment le groupe *I*

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>
<i>A</i>	0	0.2	0.3	0.5
<i>B</i>	0.2	0	0.4	0.6
<i>C</i>	0.3	0.4	0	0.3
<i>D</i>	0.5	0.6	0.3	0



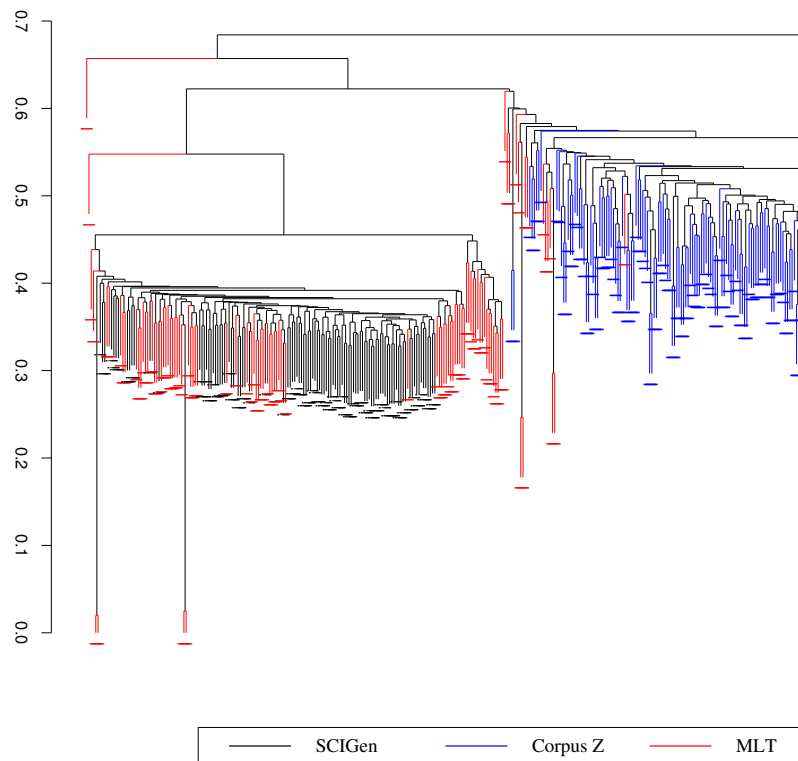
# Détection automatique

[Labbé and Labbé, 2013]

Distance inter-textuelle :

$\Delta_{(a,b)} = \delta$  proportion de mots (tokens) différents dans les deux textes.

## Hierarchical Clustering



Soit

- $t$  un texte à tester.
- $\delta_t^{Fake} = \min_{f \in SCIGen} \Delta_{(t,f)}$

Si ( $\delta_t^{Fake} < \delta_{Seuil}$ ) Alors

Une génération SCIGen doit être considérée, (risque  $< 10^{-5}$ ).

Sinon

Une origine non-SCIGen doit être considérée.

# SCIdgen papers and its clones

SSME: Int. Conf. on Services Science, Management and Engineering. 2009.

- IEEEExplorer, indexed in Scopus and WoK
- 150 papers, 4 SCIdgen and 1 duplicate.
- Official acceptance rate : 28%

## SCIdgen inside (publishers)

- 120 IEEE (retracted or deleted),
- 16 Springer (retracted),
- 1 Elsevier (accepted-unpublished)

## SCIdgen inside (social networks)

- <http://www.researchgate.net>
- <http://scholar.harvard.edu>
- <http://www.academia.edu>

## Other generators

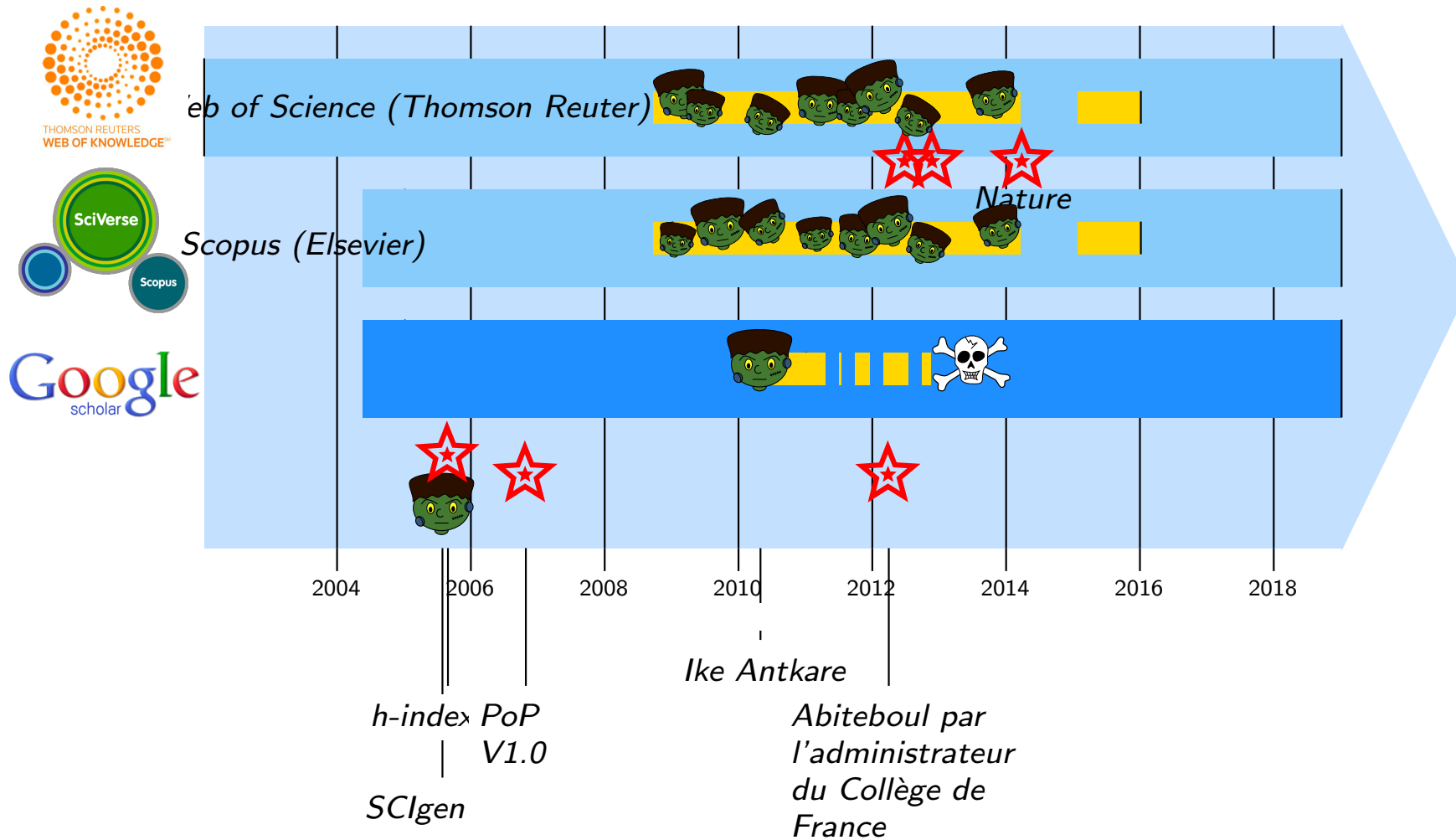
- Mathgen (<http://thatsmathematics.com/mathgen/>)
- The Postmodernism Generator (<http://www.elsewhere.org/pomo/>)
- scigen-physics (<https://bitbucket.org/birkenfeld/scigen-physics>)
- Auto. SBIR Grant Proposal Generator (<http://www.nadovich.com/chris/randprop/>)

## Dans la presse internationale scientifique et grand public (2014)



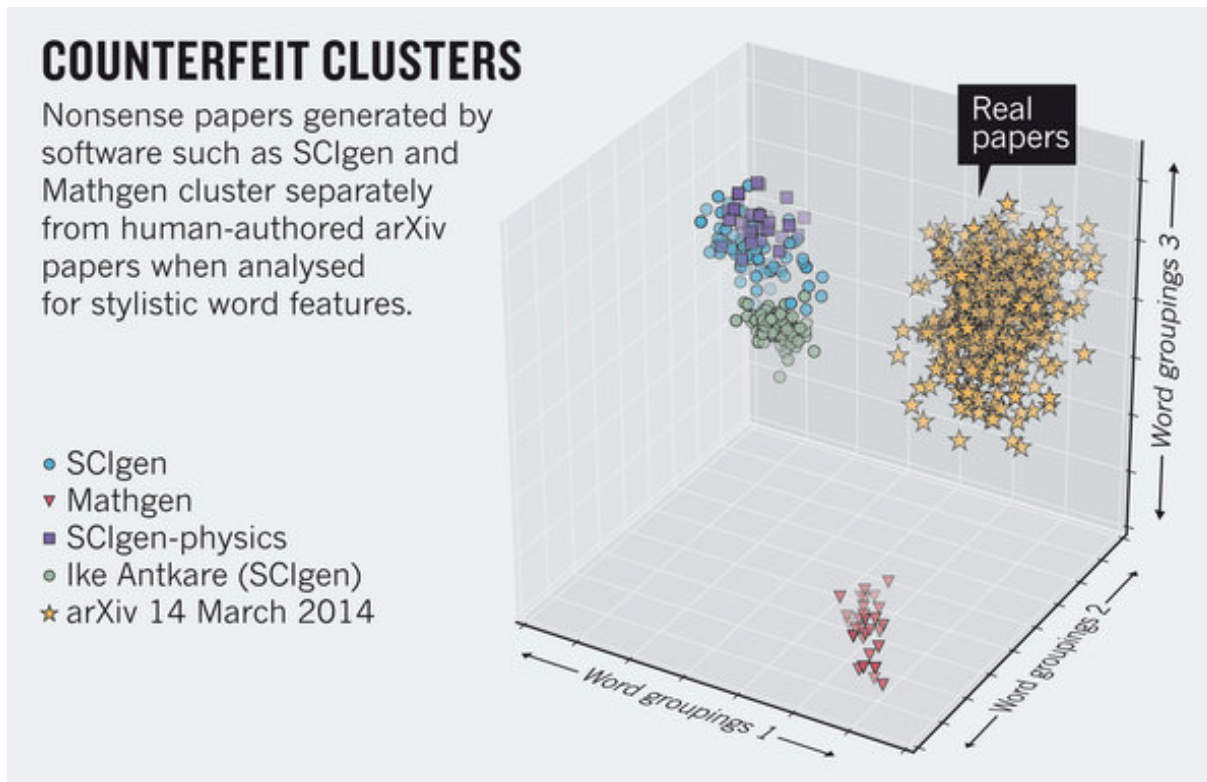


# Chronos



# No SCIdgen paper in arXiv (Computer Science)

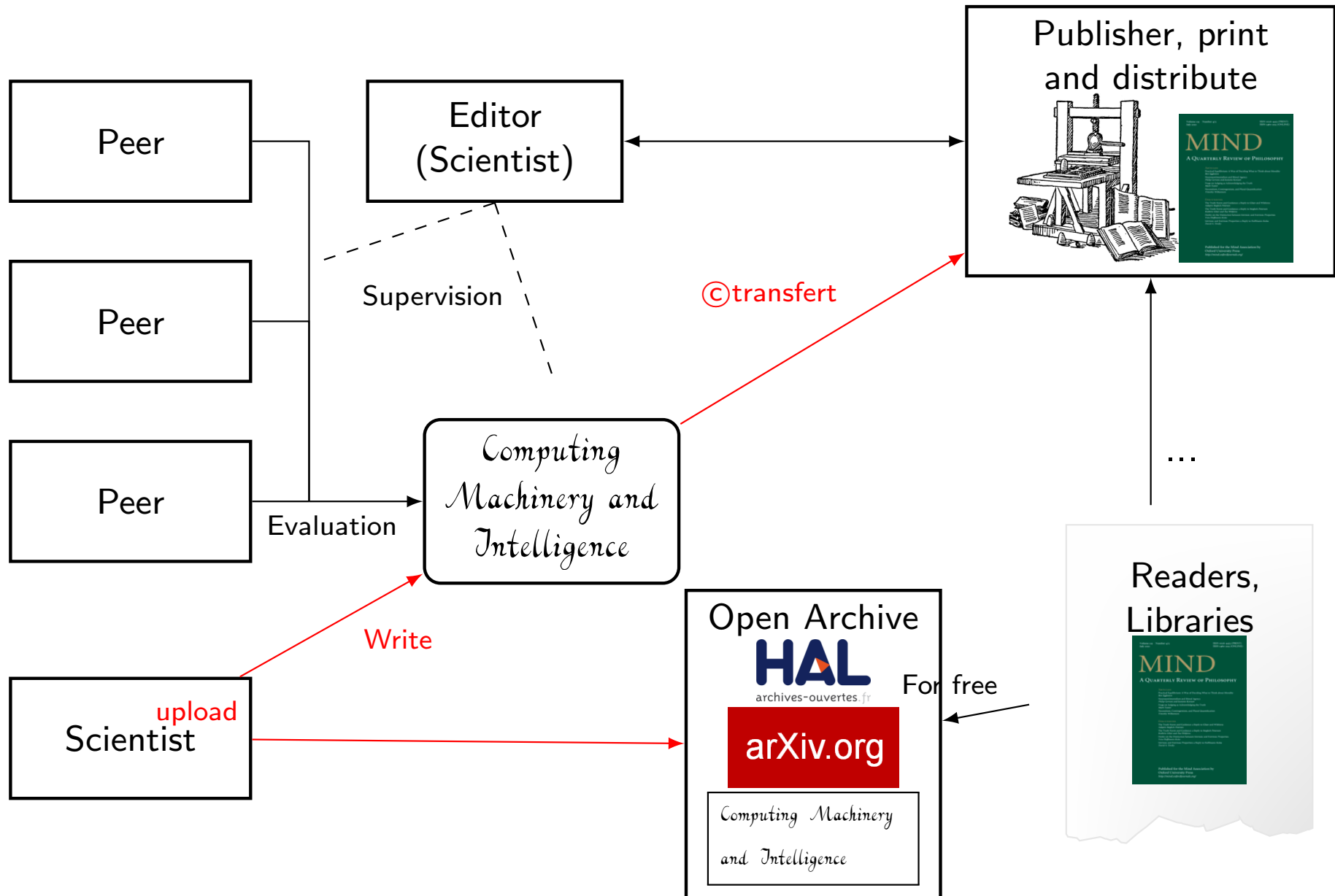
Automated screening: ArXiv screens spot fake papers



- Only stop-words
- PCA
- Supposed non Zipfian

Image borrowed from [Ginsparg, 2014]

# Publication : Self Archiving (Green Open Access)



# Where to find pirated papers

## Pirated papers

- LibGen



- Sci-Hub (Alexandra Elbakyan)



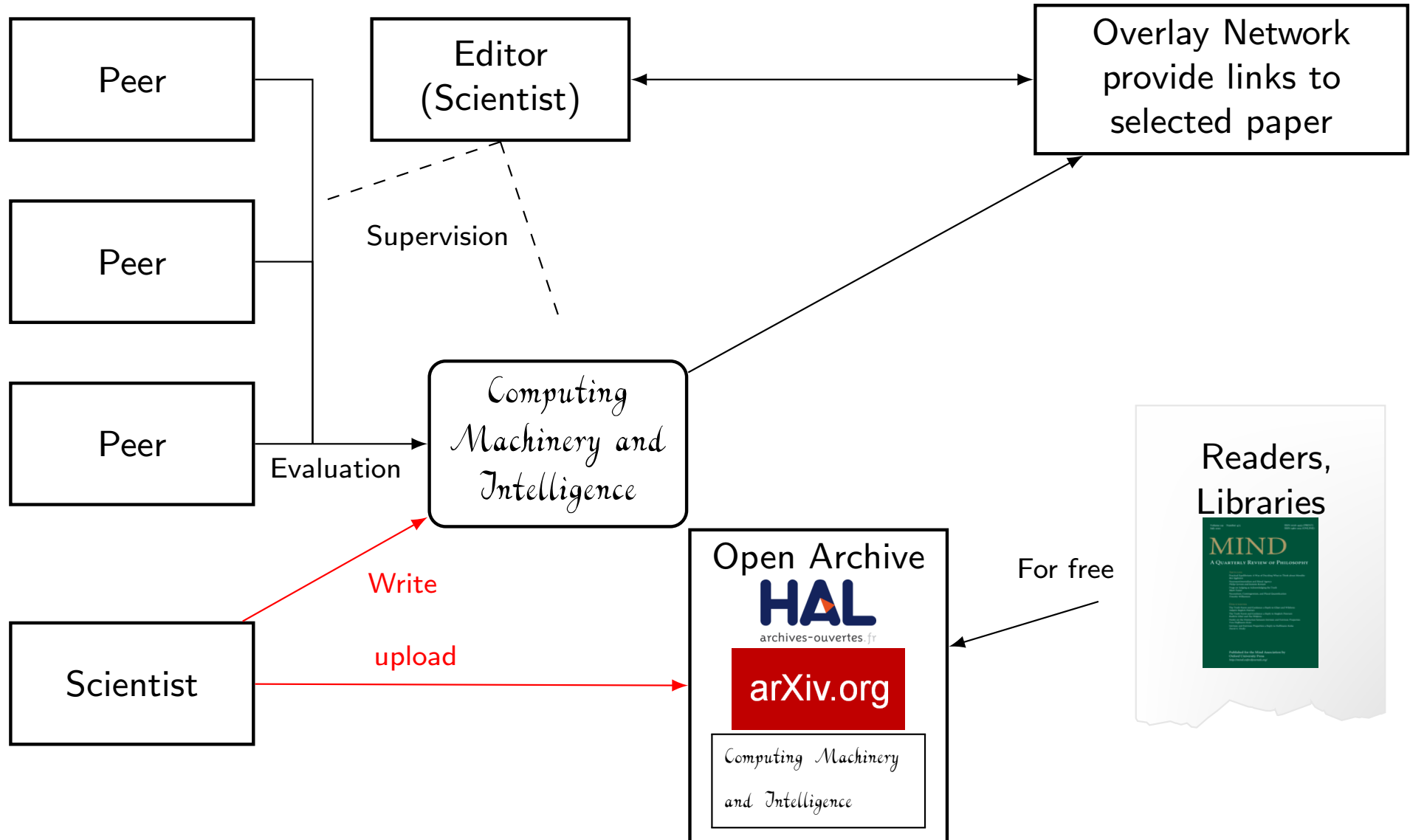
Bohannon J, Elbakyan A (2016)

Data from: Who's downloading pirated papers? Everyone.

Dryad Digital Repository. <https://doi.org/10.5061/dryad.q447c>



# Overlay Journal : les épi-journaux [episciences.org](http://episciences.org)



# Springer-Nature funded SciDetect: <http://scidetector.forge.imag.fr>

SciDetect



SciDetect is a collaboration between Springer-Verlag GmbH and Université Joseph Fourier.

## Press release, march 2015

”The open source software discovers text that has been generated with the SCIdgen computer program and other fake-paper generators like Mathgen and Physgen.”

”SciDetect is highly flexible and can be quickly customized to cope with new methods of automatically generating fake or random text”

## Do not cop with other problems

- Peer review rings
- Paper mills
- Black market and authorship selling

# Table of Contents

- 1 Pourquoi Ecrire ?
- 2 Publications et Scientometrie
  - Scientometrics: what for?
  - SClgen a Probabilistic Context Free Grammar
- 3 Of the use of fake publications
  - h-index hacking
  - Resume Padding
  - Journal Hijacking
- 4 Detection of SClgen papers
  - Google Search
  - SciDetect: Automatic detection
- 5 Automatic detection of questionable research papers
  - Fact checking science
  - Seek & Blastn tool

# Automatic detection of questionable research papers

[Byrne and Labbé, 2017b, Byrne and Labbé, 2017a]

## Scientific ethics

- Plagiarism, auto-plagiarism, content reuse...
- *N-grams* signature (hashing functions).

## Non-sense detection

- Paper generator (SCIgen, physic-gen, MathGen...)
- Authorship detection (inter-textual distance).

## Need to detect questionable scientific results

- |   |   |            |   |
|---|---|------------|---|
| <ul style="list-style-type: none"> <li>• Fabrications (making up data or results)</li> <li>• Falsification (manipulating data or results)</li> <li>• False or unsupported affirmations</li> <li>• Genuine errors</li> </ul> | } | $\implies$ | <ul style="list-style-type: none"> <li>• Error spreading</li> <li>• Wrong belief</li> <li>• Research irreproducibility</li> </ul> |
|---|---|------------|---|



# Starting point : striking similarities, obvious errors

## Jennifer Byrne:

- First reported *TPD52L2* (20 years ago)
- 5 Publications with obvious errors!

## 5 Publications from China:

- Single gene knockdown experiments.
- Human cancer cell lines.

## Conclusions highlight potential therapy

- ...*TPD52L2*... novel therapeutic target for glioma treatment.
- ...*TPD52L2*... novel clues for oral squamous cell carcinoma therapy.
- ...*TPD52L2*... therapeutic approach for the treatment of breast cancer.
- ...*TPD52L2* is indispensable in gastric cancer proliferation.
- ...*TPD52L2* could be a novel therapeutic target for human liver cancer.

# Obvious errors: example

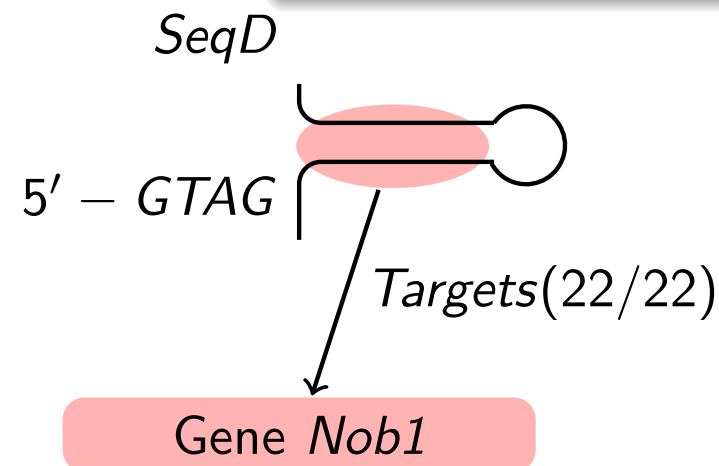
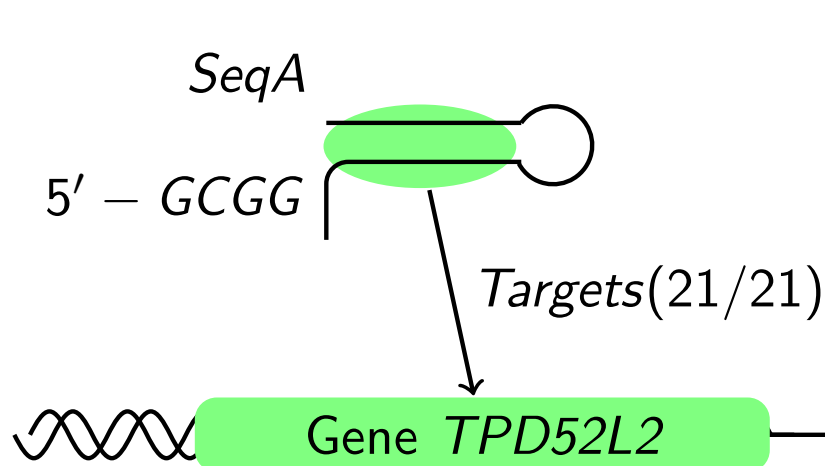
PMID : 25262828

## Materials and methods

The shRNA sequence (5'-GCGGAGGGTTTGAAGAATATCTC-GAGATATTCTTTCAAACCCTCCGCTTTTTT-3') targeting TPD52L2 (NM\_199360) was inserted into the pFH-L plasmid (Shanghai Hollybio, China). A scrambled shRNA that shared no homology with the mammalian genome (5'-CTAGCCCGGCCAAG-GAAGTGCAATTGCATACTCGAGTATGCAATTGCACTTC-CTTGTTTTTTTTGTTAAT-3') was used as control.

## Fact-Check using *blastn* (NCBI)

```
Query= SeqA (evaluate = 10)
Length=54
Sequences producing significant alignments:
... ..
> .... Homo sapiens tumor protein D52
like 2 (TPD52L2), ...
Length=2230
...
Query 1 GCGGAGGGTTTGAAGAATAT 21
      |
Sbjct 894 GCGGAGGGTTTGAAGAATAT 914
....
Query 28 ATATTCTTTCAAACCCTCCGC 48
      |
Sbjct 914 ATATTCTTTCAAACCCTCCGC 894
```



# Obvious errors: example

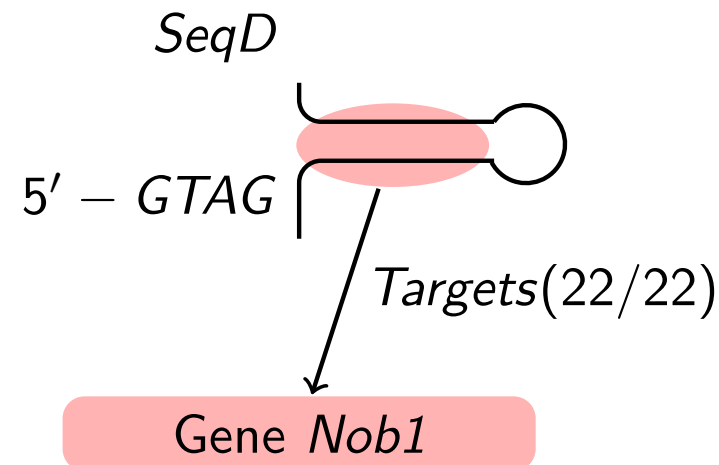
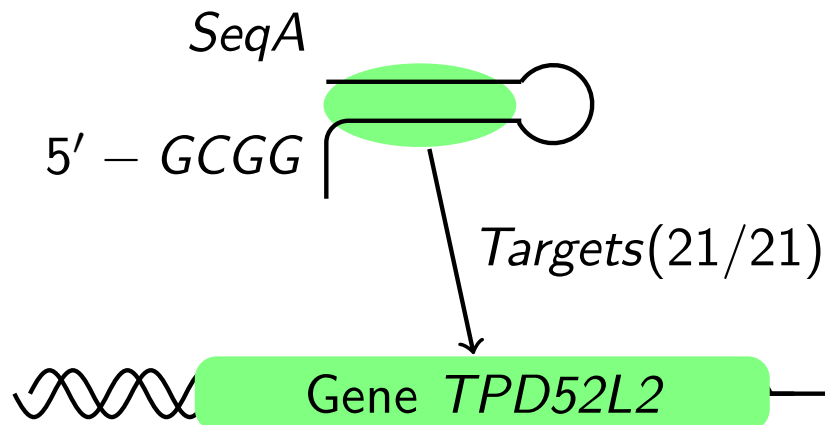
PMID : 25262828

## Materials and methods

The shRNA sequence (5'-GCGGAGGGTTTGAAGAATATCTC-GAGATATTCTTTCAAACCCTCCGCTTTTTT-3') targeting **TPD52L2** (NM\_199360) was inserted into the pFH-L plasmid (Shanghai Hollybio, China). A scrambled shRNA that shared no homology with the mammalian genome (5'-CTAGCCCGGCCAAG-GAAGTGCAATTGCATACTCGAGTATGCAATTGCACTTC-CTTGGTTTTTTGTTAAT-3') was used as control.

## Fact-Check using *blastn* (NCBI)

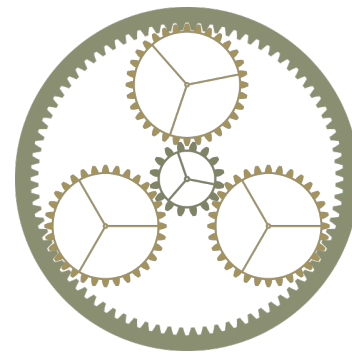
```
Query= SeqD (evaluate = 10)
Length=68
Sequences producing significant alignments:
... ..
> .... Homo sapiens NIN1/PSMD8 binding
protein 1 homolog (NOB1)...
Length=1775
...
Query 9   GCCAAGGAAGTGCAATTGCATA 30
        |||
Sbjct 1505 GCCAAGGAAGTGCAATTGCATA 1526
....
Query 37  TATGCAATTGCACTTCCTTGG 57
        |||
Sbjct 1526 TATGCAATTGCACTTCCTTGG 1506
```



# Seek & Blastn at a glance

Materials and methods  
 The shRNA sequence (5'-GCGGAGGGTTTGAAG-  
 GAATATCTCGAGATATTCTTTCAAACCCTCCGCTTTTTT-  
 3') targeting TPD52L2 (NM\_199360) was inserted into  
 the pFH-L plasmid (Shanghai Hollybio, China). A  
 scrambled shRNA that shared no homology with the  
 mammalian genome (5'-CTAGCCCGGCCAAGGAAGTG-  
 CAATTGCATACTCGAGTATGCAATTGCACTTCCTTG-  
 GTTTTTGTTAAT-3') was used as control.

(1) Facts extraction:  
*Named entity recognition*, extract nucleotide  
 and status...



Facts to check

Status	DNA Seq
...	...
Targeting	GCG...TTT
Non-Targ.	CTA...AAT
...	...

(2) Blastn call  
 software gives  
 the hit list

Hit lists (Blastn results)

hit list	DNA Seq
...	...
TPD52L2, ...	GCG...TTT
NOB1,...	CTA...AAT
...	...

(3) Comparison

Checked Facts

Satus	DNA Seq
Targ.	GCG...TTT
Non-Targ.	CTA...AAT
...	...

# Ambiguïtés : polysémie, homonymie, structurale,...

- Le président a le **pouvoir** de faire taire l'**avocat**.
- Je ne vais pas **pouvoir** manger l'**avocat**.
- l'**été** à l'**est** a **été** très beau et l'**est** toujours.

- Je **suis** le secrétaire.
- Je vais à la grange et la **ferme**.

- Il **poursuit** la **jeune fille** à vélo.
- Il a vu **un homme** avec un télescope.
- Tous les participants prendront **un** bus.

# Seek & Blastn

## Related works

- Detection of statistically flawed paper
- Fake news detection

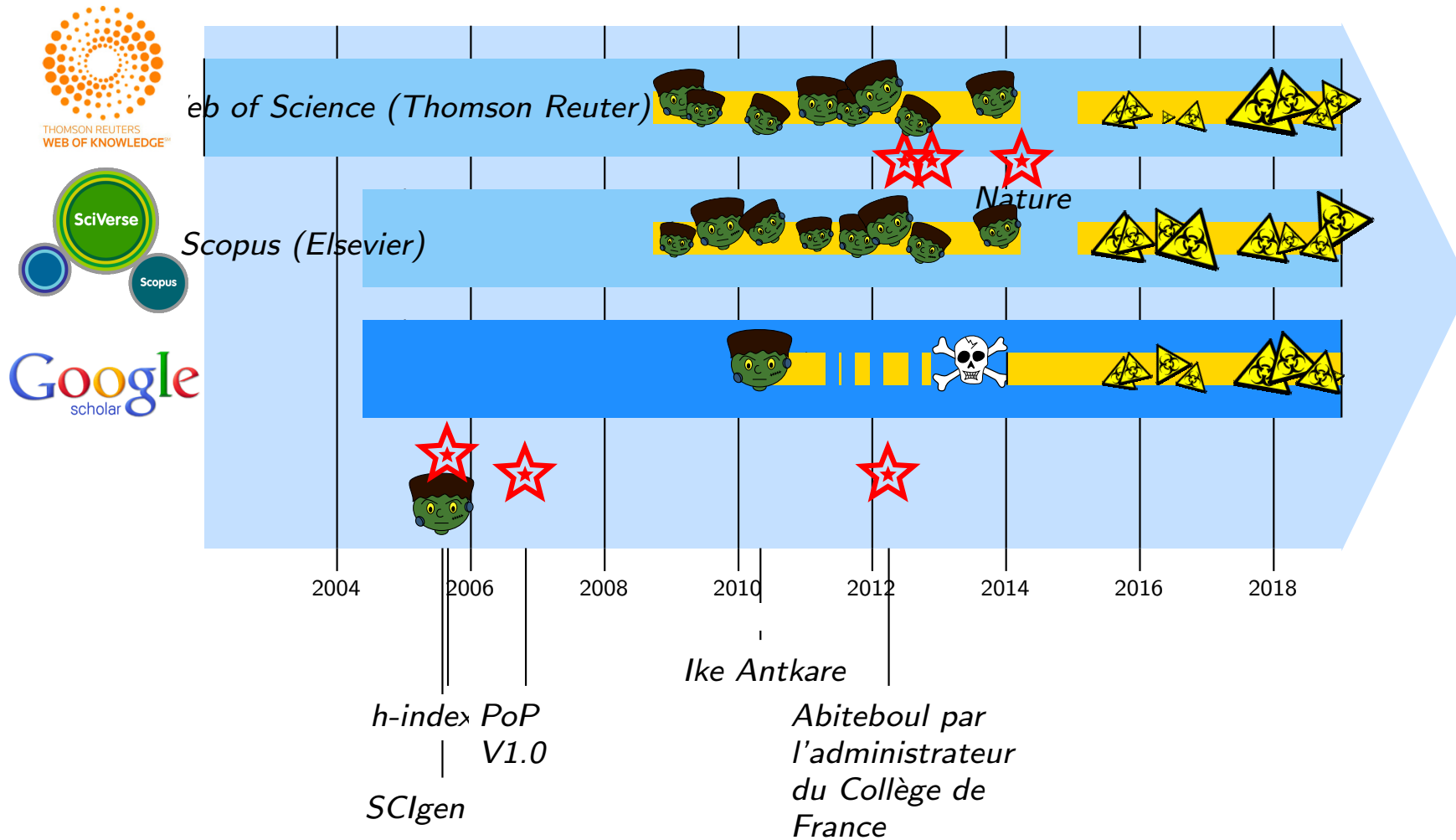
## Seek & Blastn perspectives

- Online tool : <http://scigendetection.imag.fr/TPD52>
- Avoid false positive, more in-deep analysis of sentences.

## Retractions, Errors corrections

- Retractions ( $\approx 18$ ), Expression of concern ( $\approx 11$ ),  $\approx 45$  to be treated
- Citation analysis (to be done)

# Chronos



# Conclusion and Future/Ongoing works

## Publication procedures, models and habits

- Why fake papers were accepted, published and ... sold.
- Traditional publisher vs open access.
- Knowledge diffusion: better and less... or as much as possible.

## Blind management rules...

- ... are an incitation to malpractices: slicing, plagiarism, faked data, ...

## Automatic detection of new generators

- Hand written PCFG : find dense cluster inside a population.
- Study other kind of generator (language model).

## In the web today

- Automatic knowledge extraction/detection/generation.
- How to separate the wheat from the chaff... and scale up !



# Thanks



Amancio, D. R. (2015).

Comparing the topological properties of real and artificially generated scientific manuscripts.  
*Scientometrics*, 105(3):1763–1779.



Beel, J. and Gipp, B. (2010).

Academic search engine spam and google scholar's resilience against it.  
*Journal of Electronic Publishing*, 13(3).



Beel, J., Gipp, B., and Wilde, E. (2010).

Academic search engine optimization (aseo).  
*Journal of scholarly publishing*, 41(2):176–190.



Byrne, J. A. and Labbé, C. (2017a).

Fact checking nucleotide sequences in life science publications: The seek & blastn tool.  
In *International Congress on Peer Review and Scientific Publication, Enhancing the quality and credibility of science*, Chicago.



Byrne, J. A. and Labbé, C. (2017b).

Striking similarities between publications from china describing single gene knockdown experiments in human cancer cell lines.  
*Scientometrics*, 110(3):1471–1493.



Dalkilic, M. M., Clark, W. T., Costello, J. C., and Radivojac, P. (2006).

Using compression to identify classes of inauthentic texts.  
In *Proceedings of the 2006 SIAM Conference on Data Mining*.



Fahrenberg, U., Biondi, F., Corre, K., Jégourel, C., Kongshøj, S., and Legay, A. (2014).

Measuring structural distances between texts.  
*CoRR*, abs/1403.4024.



Ginsparg, P. (2014).

Automated screening: Arxiv screens spot fake papers.  
*Nature*, 508(7494):44–44.



Hirsch, J. E. (2005).

An index to quantify an individual's scientific research output.  
*Proceedings of the National Academy of Science*, 102:16569–16572.



Labbé, C. (2010).

Ike antkare, one of the great stars in the scientific firmament.  
*International Society for Scientometrics and Informetrics Newsletter*, 6(2):48–52.



Labbé, C. and Labbé, D. (2006).

A tool for literary studies. intertextual distance and tree classification.  
*Literary and Linguistic Computing*, 21(3):311–326.



Labbé, C. and Labbé, D. (2013).

Duplicate and fake publications in the scientific literature: how many scigen papers in computer science?  
*Scientometrics*, 94(1):379–396.



Lavoie, A. and Krishnamoorthy, M. (2010).

Algorithmic Detection of Computer Generated Text.  
*ArXiv e-prints*.



Lopez-Cozar, E. D., Robinson-García, N., and Torres-Salinas, D. (2012).

Manipulating google scholar citations and google scholar metrics: Simple, easy and tempting.  
*arXiv preprint arXiv:1212.0638*.



Xiong, J. and Huang, T. (2009).

An effective method to identify machine automatically generated paper.  
In *KESE '09. Pacific-Asia Conference*, pages 101–102.