# Measuring Global Similarity between Texts

Uli Fahrenberg    Fabrizio Biondi    Kevin Corre    Cyrille Jegourel    Simon Kongshøj    Axel Legay

IRISA/Inria Rennes

University College of Northern Denmark

SLSP 2014

1. Fake scientific papers

2. Distances between texts

3. Implementation

4. Experiment 1

5. Experiment 2

6. Conclusion

# nature
International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & Video

*NATURE* | NEWS

# Publishers withdraw more than 120 gibberish papers

**Conference proceedings removed from subscription databases after scientist reveals that they were computer-generated.**

**Richard Van Noorden**

24 February 2014 | Updated: 25 February 2014

🔍 **Rights & Permissions**

The publishers Springer and IEEE are

The publishers Springer and IEEE are removing more than 120 papers from their subscription services after a French researcher discovered that the works were computer-generated nonsense.

Over the past two years, computer scientist Cyril Labbé of Joseph Fourier University in Grenoble, France, has catalogued computer-generated papers that made it into more than 30 published conference proceedings between 2008 and 2013. Sixteen appeared in publications by Springer, which is headquartered in Heidelberg, Germany, and more than 100 were published by the Institute of Electrical and Electronic Engineers (IEEE), based in New York. Both publishers, which were privately informed by Labbé, say that they are now removing the papers.

# SCIgen - An Automatic CS Paper Generator

**About** **Generate** **Examples** **Talks** **Code** **Donations** **Related** **People** **Blog**

## About

SCIgen is a program that generates random Computer Science research papers, including graphs, figures, and citations. It uses a hand-written **context-free grammar** to form all elements of the papers. Our aim here is to maximize amusement, rather than coherence.

One useful purpose for such a program is to auto-generate submissions to conferences that you suspect might have very low submission standards. A prime example, which you may recognize from spam in your inbox, is SCI/IIIS and its dozens of co-located conferences (check out the very broad conference description on the **WMSCI 2005** website). There's also a list of **known bogus conferences**. Using SCIgen to generate submissions for conferences like this gives us pleasure to no end. In fact, one of our papers was accepted to SCI 2005! See **Examples** for more details.

# A Case for Von Neumann Machines

Adam Marcus and Ivan Oransky

## Abstract

Many analysts would agree that, had it not been for lossless communication, the emulation of suffix trees might never have occurred. In this position paper, we disprove the exploration of Moore's Law. In our research we demonstrate that the foremost pervasive algorithm for the emulation of Moore's Law [10] is impossible.

## 1  Introduction

Recent advances in "smart" configurations and semantic symmetries are based entirely on the assumption that the producer-consumer problem and expert systems are not in conflict with Moore's Law. We emphasize that Eyebolt is NP-complete. On a

mar. Two properties make this approach different: our algorithm is Turing complete, and also our solution is in Co-NP. Furthermore, our heuristic locates the deployment of context-free grammar. Combined with mobile technology, this finding visualizes a methodology for the improvement of erasure coding. Though this finding at first glance seems perverse, it has ample historical precedence.

We propose new compact methodologies, which we call Eyebolt. For example, many methodologies manage architecture. The basic tenet of this solution is the analysis of 128 bit architectures. Two properties make this method optimal: our system studies IPv6, and also Eyebolt turns the stochastic archetypes sledgehammer into a scalpel. Existing replicated and classical heuristics use

## Duplicate and fake publications in the scientific literature: how many SCIgen papers in computer science?

Cyril Labbé · Dominique Labbé

**Abstract**  Two kinds of bibliographic tools are used to retrieve scientific publications and make them available online. For one kind, access is free as they store information made publicly available online. For the other kind, access fees are required as they are compiled on information provided by the major publishers of scientific literature. The former can easily be interfered with, but it is generally assumed that the latter guarantee the integrity

## 1-gram distance

- $F_A(w)$: absolute frequency of 1-gram $w$ in text $A$
- for texts $A$, $B$ of equal length $N$:

$$d_1(A, B) = \frac{\sum_{w \in A \circ B} |F_A(w) - F_B(w)|}{2\,N}$$

- adjustments needed when texts do not have equal length
  - (scaling)

## 1-gram distance

- $F_A(w)$: absolute frequency of 1-gram $w$ in text $A$
- for texts $A$, $B$ of equal length $N$:

$$d_1(A, B) = \frac{\sum_{w \in A \circ B} |F_A(w) - F_B(w)|}{2\,N}$$

- adjustments needed when texts do not have equal length
  - (scaling)

- bag-of-words model
  - or bag-of-$n$-grams model ($n = 1, 2, 3, 4$) for instance for BLEU distance
- want to match longer $n$-grams

## Our distance

- Let $A = (a_1, a_2, \ldots, a_{N_A})$ and $B = (b_1, b_2, \ldots, b_{N_B})$ be texts
- Write $\delta_{i,j} = [\text{ if } a_i = b_j \text{ then } 0 \text{ else } 1\,]$
    - (word match indicator)
- position match:

$$d_{\mathsf{pm}}(i,j) = \delta_{i,j} + \lambda\delta_{i+1,j+1} + \lambda^2\delta_{i+2,j+2} + \lambda^3\delta_{i+3,j+3} + \cdots$$

- $\lambda$: discounting factor, $0 \leq \lambda < 1$
- "try to match $n$-grams for $n$ as high as possible, but don't be too sad if very long phrases don't match"
- (for $\lambda = 0$, this is 1-gram matching!)
- (matches outside the texts' lengths automatically score value1)

## Example

- text $A =$ "the quick fox jumps over the lazy dog"
- text $B =$ "the lazy fox"
- $\lambda = .8$; values scaled by $1 - \lambda$

|      | the  | quick | fox  | jumps | over | the  | lazy | dog  |
|------|------|-------|------|-------|------|------|------|------|
| the  | 0.67 | 1.00  | 1.00 | 1.00  | 1.00 | 0.64 | 1.00 | 1.00 |
| lazy | 1.00 | 0.84  | 1.00 | 1.00  | 1.00 | 1.00 | 0.80 | 1.00 |
| fox  | 1.00 | 1.00  | 0.80 | 1.00  | 1.00 | 1.00 | 1.00 | 1.00 |

## Our distance, contd.

- Let $A = (a_1, a_2, \ldots, a_{N_A})$ and $B = (b_1, b_2, \ldots, b_{N_B})$ be texts
- Write $\delta_{i,j} = [\,\text{if } a_i = b_j \text{ then } 0 \text{ else } 1\,]$
- position match:

$$d_{\mathsf{pm}}(i, j) = \delta_{i,j} + \lambda\delta_{i+1,j+1} + \lambda^2\delta_{i+2,j+2} + \lambda^3\delta_{i+3,j+3} + \cdots$$

- $\lambda$: discounting factor, $0 \leq \lambda < 1$
- global distance:

$$d_3(A, B) = \frac{1 - \lambda}{N_A} \sum_{i=1}^{N_A} \min_{j=1,\ldots,N_B} d_{\mathsf{pm}}(i, j)$$

- find best possible match for each position in $A$, average, and scale
- and symmetrize: $d_4(A, B) = \max(d_3(A, B), d_3(B, A))$

## Implementation

Backwards computation:

```
scaling = 1/(1-discount);
for ( i=0; i<text1size+1; i++ ) {
  dist[i][text2size] = scaling;
}
for ( j=0; j<text2size; j++ ) {
  dist[text1size][j] = scaling;
}
for ( i=text1size-1; i>-1; i-- ) {
  for ( j=text2size-1; j>-1; j-- ) {
    dist[i][j] = compare( text1[i], text2[j] )
        + discount*dist[i+1][j+1];
  }
}
```
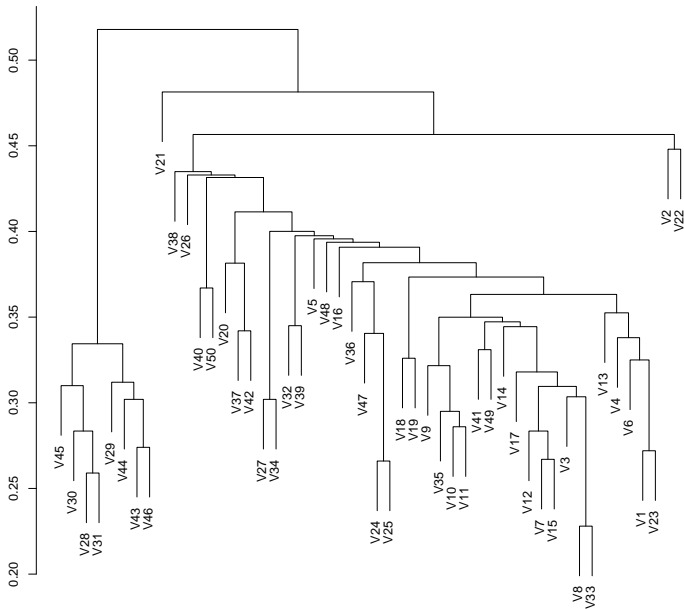
```
uli@ulrich:~/textdist-1.5$ ls
CHANGES.txt      extract-mat.sh    notes.txt           sdist-0.80.txt    textdist.o
compareall.sh    Makefile          preprocess-pdf.sh   textdist          texts
config.h         names.txt         README.txt          textdist.c
uli@ulrich:~/textdist-1.5$ ls texts/automogensen0*
texts/automogensen00.txt   texts/automogensen04.txt   texts/automogensen08.txt
texts/automogensen02.txt   texts/automogensen06.txt
uli@ulrich:~/textdist-1.5$ ./textdist texts/automogensen0*
Distance between Text 1 and Text 2: 0.723
Distance between Text 1 and Text 3: 0.719
Distance between Text 1 and Text 4: 0.715
Distance between Text 1 and Text 5: 0.713
Distance between Text 2 and Text 3: 0.718
Distance between Text 2 and Text 4: 0.718
Distance between Text 2 and Text 5: 0.721
Distance between Text 3 and Text 4: 0.713
Distance between Text 3 and Text 5: 0.709
Distance between Text 4 and Text 5: 0.711
uli@ulrich:~/textdist-1.5$ ls
CHANGES.txt      extract-mat.sh    notes.txt           sdist-0.80.txt    textdist.o
compareall.sh    Makefile          preprocess-pdf.sh   textdist          texts
config.h         names.txt         README.txt          textdist.c
uli@ulrich:~/textdist-1.5$ cat sdist-0.80.txt
0.000 0.723 0.719 0.715 0.713
0.723 0.000 0.718 0.718 0.721
0.719 0.718 0.000 0.713 0.709
0.715 0.718 0.713 0.000 0.711
0.713 0.721 0.709 0.711 0.000
uli@ulrich:~/textdist-1.5$
```
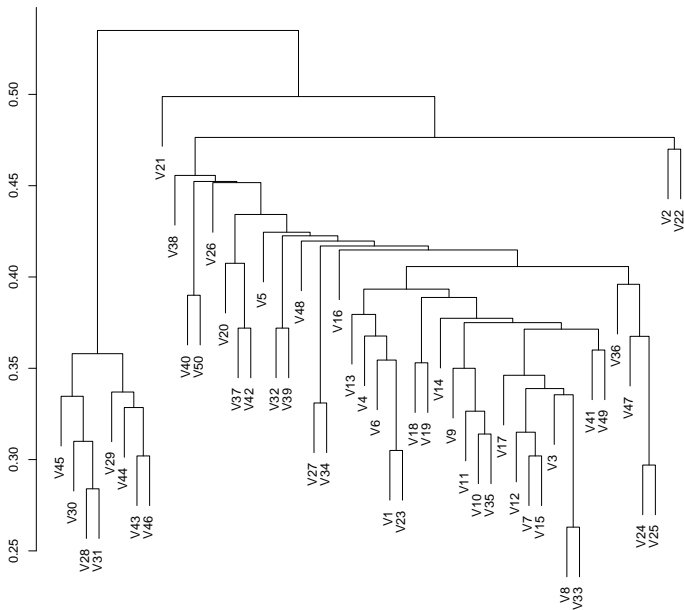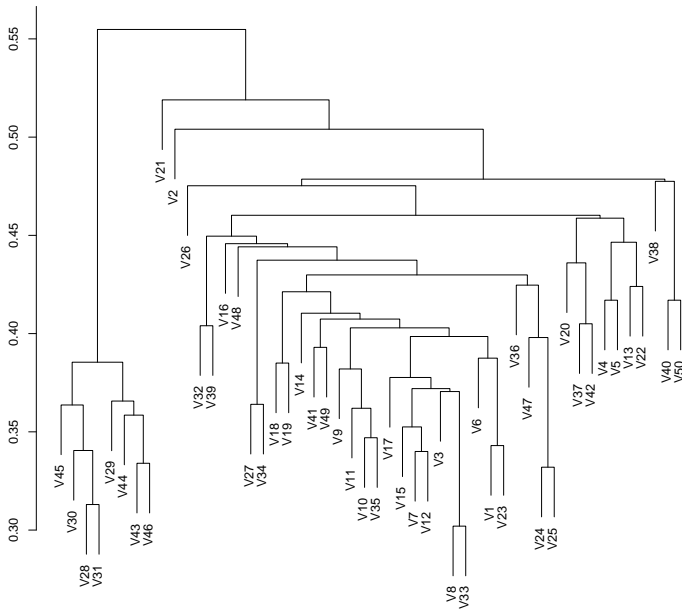
## Experiment 1

- 42 genuine scientific papers, all about modeling and verification
- 8 SciGen papers

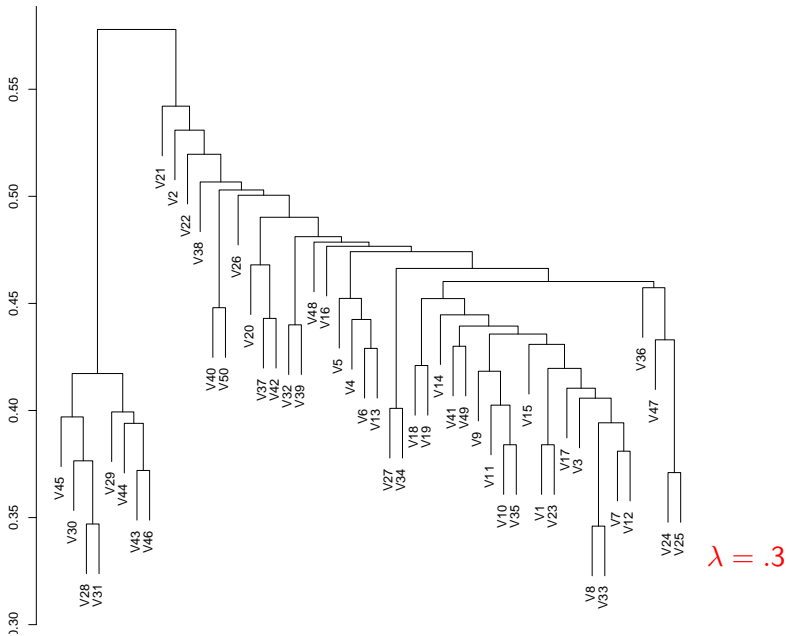- nice and clear separation; gets better with higher discounting
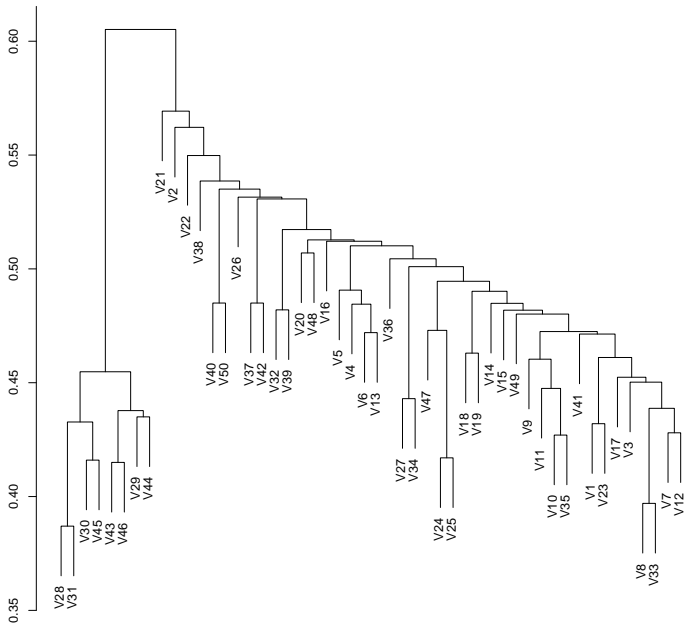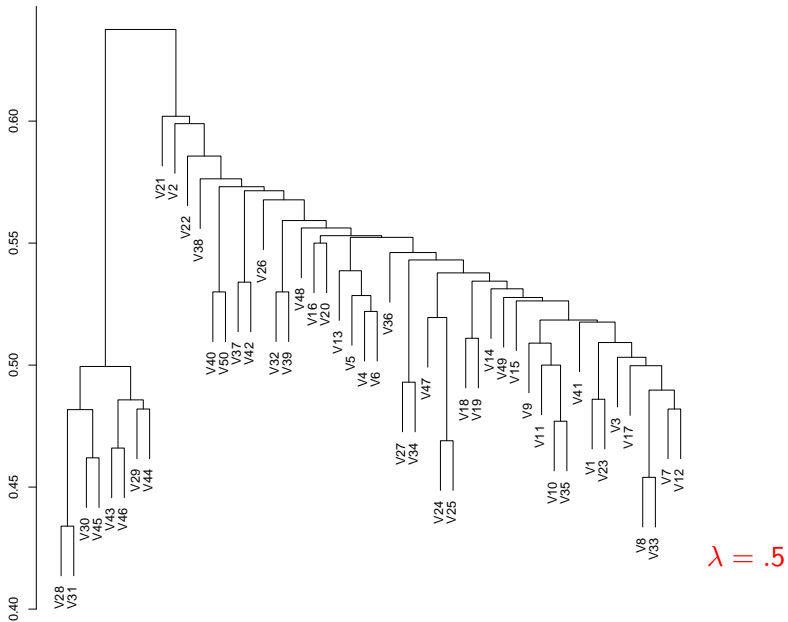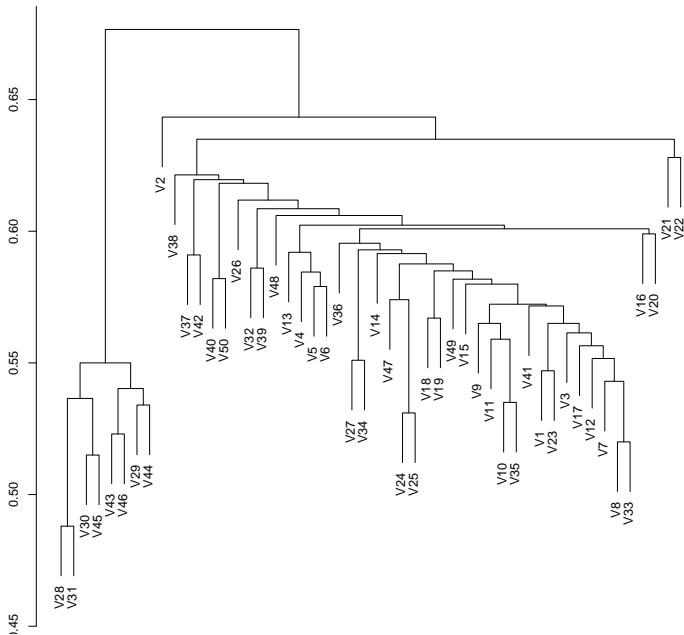
$\lambda = 0$

$\lambda = .1$

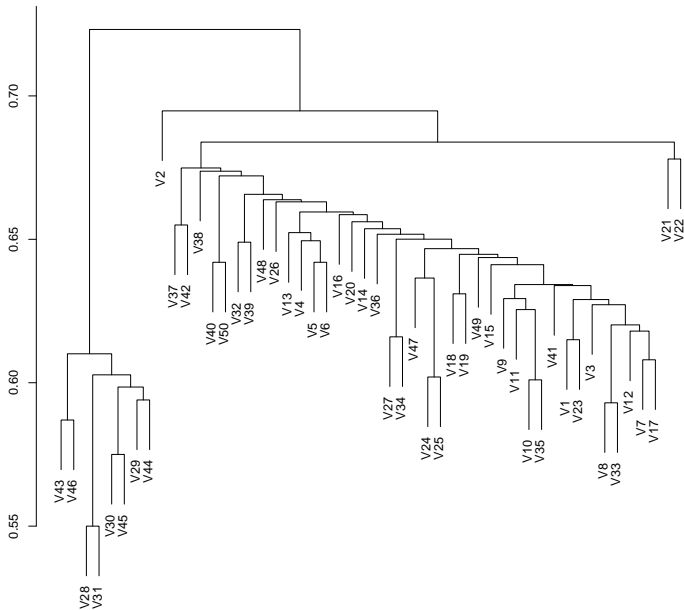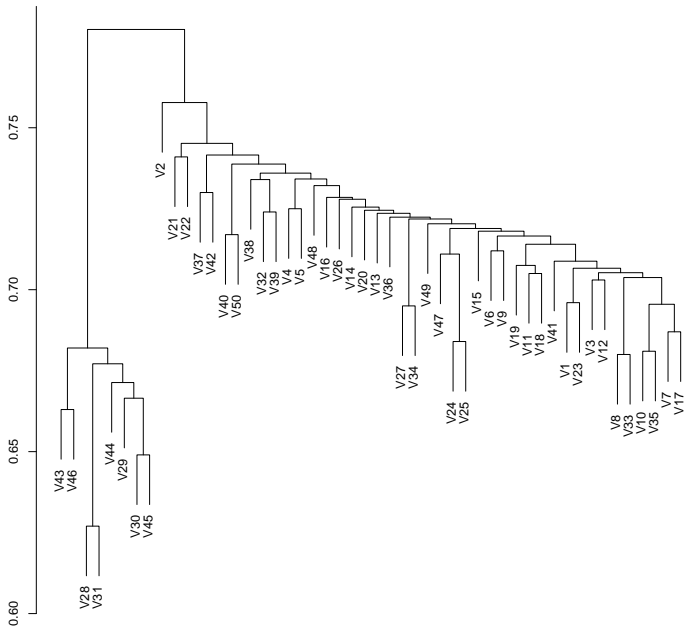$\lambda = .2$

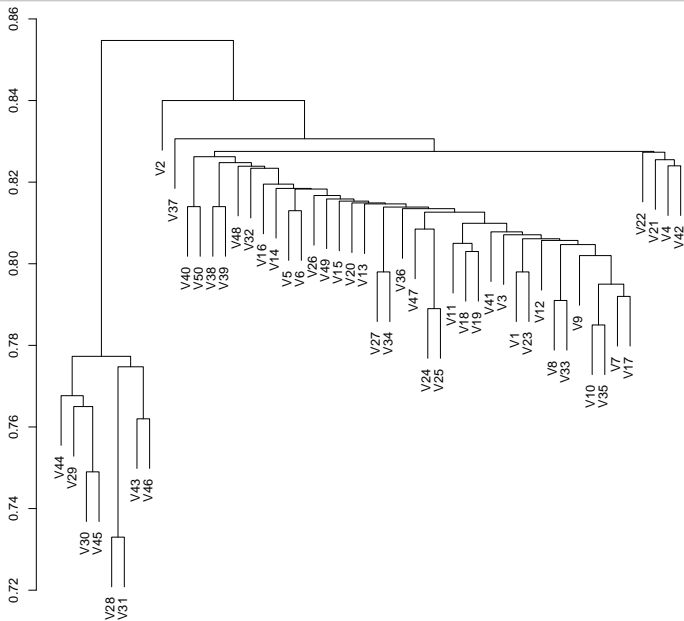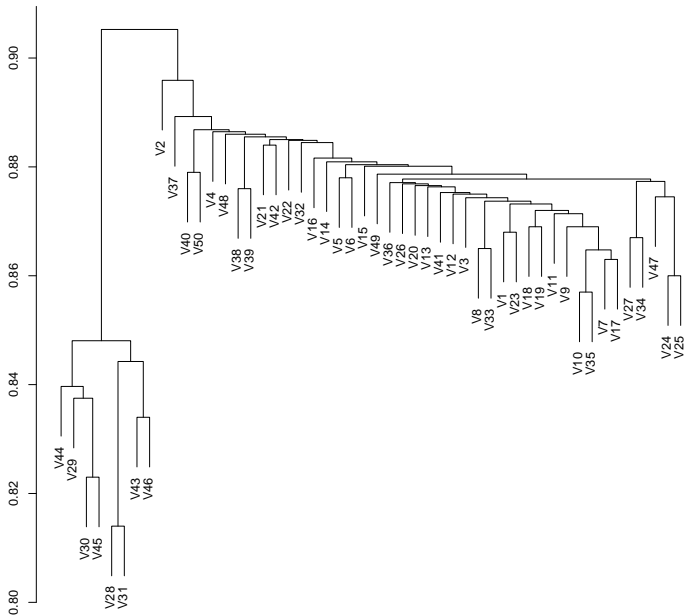$\lambda = .3$

$\lambda = .4$

$\lambda = .5$

$\lambda = .6$

$\lambda = .7$

$\lambda = .8$

$\lambda = .9$

$\lambda = .95$

- paper 8: Nicolas Basset and Eugene Asarin. Thin and thick timed regular languages. FORMATS, 2011

- paper 33: Eugene Asarin and Aldric Degorre. Volume and entropy of regular timed languages. HAL, 2009

# Thin and Thick Timed Regular Languages

Nicolas Basset[1] and Eugene Asarin[2]

[1] LIGM, Université Paris-Est and CNRS
nbasset@dptinfo.ens-cachan.fr
[2] LIAFA, Université Paris Diderot and CNRS
asarin@liafa.jussieu.fr

**Abstract.** In previous literature on timed automata, it was noticed that they are in several aspects too precise, which leads sometimes to strange artifacts, mathematical pathologies or unrealistic models. In particular, some timed automata are non-implementable, non-robust, behave badly under discretization, have many Zeno runs etc. In this paper, we propose a unifying approach to most of these issues for deterministic timed automata. We classify these automata either as **thin** or as **thick**. In thin automata, all the infinite trajectories are, in some weak sense, Zeno; the discretization of long trajectories is difficult, since it requires very small discretization step. In thick automata, most of trajectories are non-Zeno and behave well under discretization; such automata satisfy a sort of pumping lemma. Formally, the thin-thick alternative is based on the notion of entropy of timed regular languages introduced by E. Asarin and A. Degorre in [3,4]. Thin languages have the entropy $-\infty$, while thick

# Volume and entropy of regular timed languages

Eugene Asarin[1] and Aldric Degorre[2]

[1] LIAFA, Université Paris Diderot / CNRS
case 7014, 75205 Paris Cedex 13, France
Eugene.Asarin@liafa.jussieu.fr
[2] VERIMAG,
Centre Equation, 2 av. de Vignate, 38610 Gières, France
Aldric.Degorre@imag.fr

**Abstract.** For timed languages, we define measures of their size: volume for a fixed finite number of events, and entropy (growth rate) as asymptotic measure for an unbounded number of events. These measures can be used for comparison of languages, and the entropy can be viewed as information contents of a timed language. In case of languages of deterministic timed automata, we give exact formulas for volumes. Next we characterize the entropy, using methods of functional analysis, as a logarithm of the leading eigenvalue (spectral radius) of a positive integral operator. We devise several methods to compute the entropy: a symbolical one for so-called "$1\frac{1}{2}$-clock" automata, and two numerical ones: one using techniques of functional analysis, another based on discretization

## Experiment 2

- 97 `arxiv` papers (search with author = "Smith")
- 50 SciGen papers (modified SciGen to generate longer papers)
- 50 `automogensen` papers
    - simple 3-place Markov chain text generator
    - fed with the genuine corpus from Experiment 1 (length 247,000; 14,000 unique words)

- not so clear separation with average clustering
- `arxiv` papers have wildly varying distances; disturb
- clear separation with Ward clustering

## POLITIKEN

# DEBAT

**Leder**   **Kronik**   **Debattører**   **Dagens tegning**   **Debatregler**   **Debat i avisen**
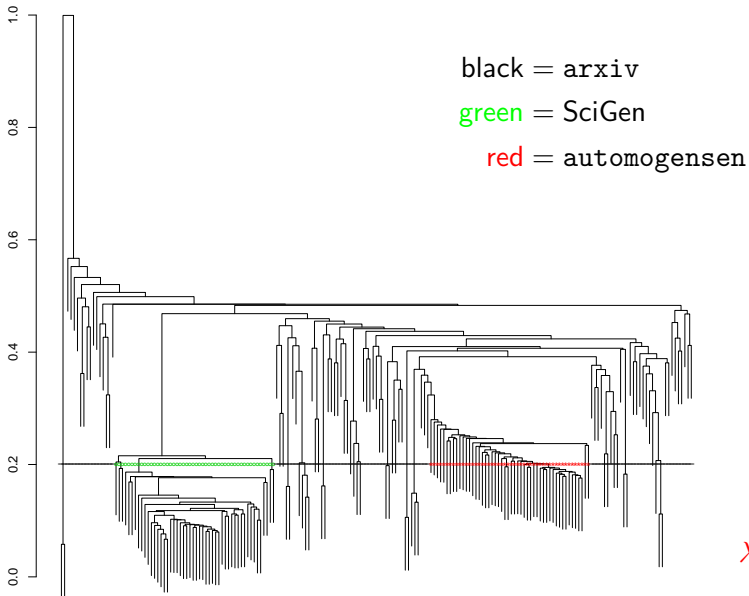
Annonce

## PETER MOGENSEN
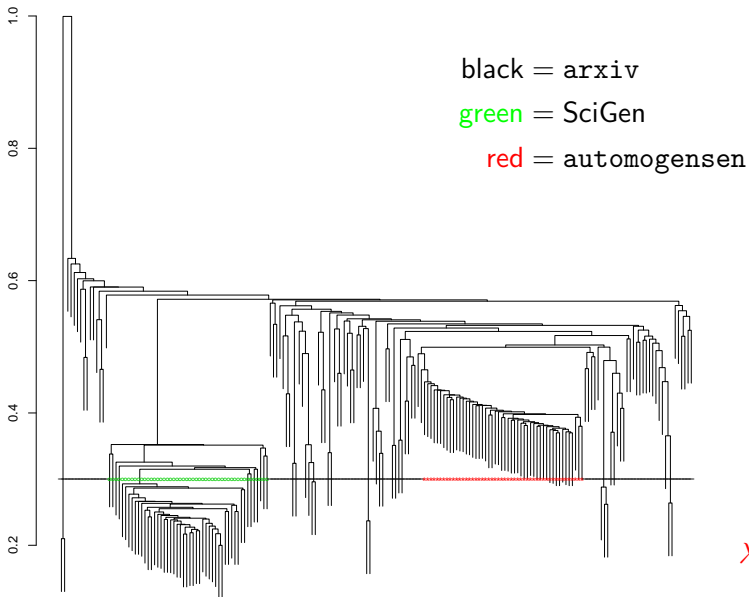
Mogensens meninger om dansk politik

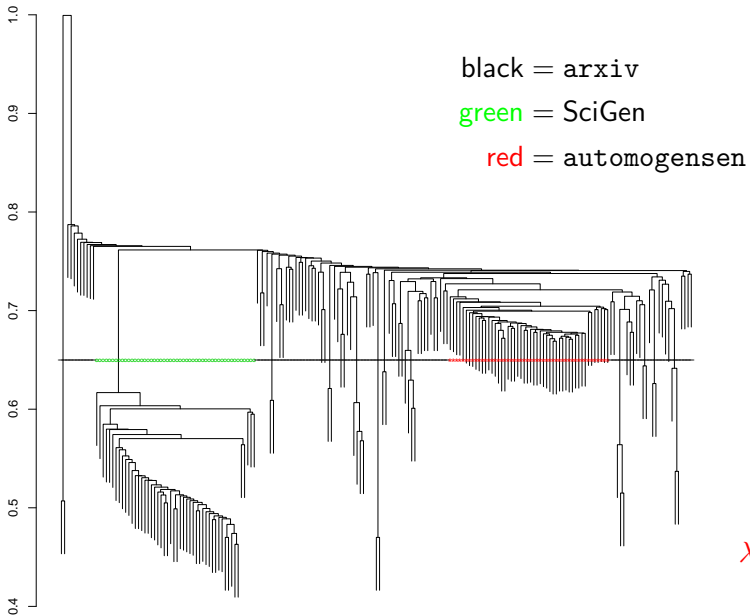**DEBATTØRER**   13. AUG. 2014 KL. 10.07

# Frustrerede socialdemokrater håber at Thorning får EU-topjob

Hvis spændingen skal bevares frem til valgdagen, skal Helle Thorning have jobbet som EU-præsident.

black = `arxiv`

green = SciGen

red = `automogensen`

$\lambda = 0$

black = `arxiv`

green = SciGen

red = `automogensen`

$\lambda = .4$

black = `arxiv`

green = SciGen

red = `automogensen`

$\lambda = .8$

black = `arxiv`

green = SciGen

red = `automogensen`

$\lambda = .95$

black $=$ `arxiv`

green $=$ SciGen

red $=$ `automogensen`

Ward clustering

$\lambda = .95$

## Conclusion

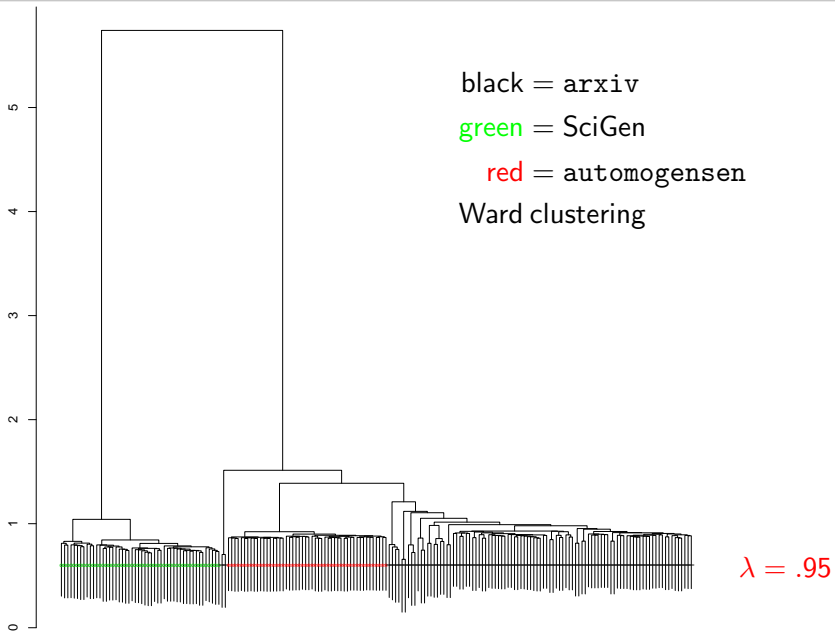- new intertextual distance which can measure differences in $n$-grams for all $n$ simultaneously
- experiments show good separation of fake and genuine scientific papers
- more experiments needed, also on other types of papers
- need thorough comparison with state of the art

Future work:

- WordNet integration
- discounted edit distance (swaps and deletions)

## Springer

Search

Home    Disciplines    My Springer    Services    Éditions    Springer Shop

**Statements**    Home  >  About Springer  >  Media  >  Statements

Berlin, 14 avril 2014

# Springer second update on SCIgen-generated papers in conference proceedings

**Background**

As reported in the media, on 11 February 2014 we were alerted to 16 fake submissions that were published in conference proceedings in Springer publications, mostly in computer sciences and engineering. The submissions were generated by the SCIgen computer

**Background**

As reported in the media, on 11 February 2014 we were alerted to 16 fake submissions that were published in conference proceedings in Springer publications, mostly in computer sciences and engineering. The submissions were generated by the SCIgen computer program, which creates nonsense documents.

⌖ Springer

On 27 February 2014, we posted a _statement_ on springer.com with the facts at that point and a promise to provide further information. On 20 March 2014, we issued an _update_ as to what we were doing to investigate the matter.

**Working with Dr. Labbé**

In February we reached out to Dr. Cyril Labbé, a researcher at Joseph Fourier University in Grenoble, France and an expert in the detection of SCIgen-generated articles for advice and collaboration. In the meantime, we have had far-reaching and exhaustive discussions with Dr. Labbé and, on 1 April 2014 a number of Springer representatives met with him in Heidelberg, Germany, to discuss the progress to date of the system implemented and currently running through the hundreds of thousands of documents on SpringerLink that have been published since the creation of the SCIgen program.

We are also pleased to announce that we will be formalizing our relationship with Dr. Labbé, who, with his team of PhD students, will work intensively with Springer on detection mechanisms for any future programs that are similar to SCIgen.