

# Author Name Disambiguation for the Inspire Project

Mateusz Susik

14 marca 2018



## Hawking, Stephen W.

[View Profile](#) [Manage Profile](#) [Manage Publications](#) [Help](#)Profile Name  

© 2015-11-09 14:55:29

## PERSONAL INFORMATION

## Personal Details (HepNames)

**Name** Stephen William Hawking

**Current Institution** Cambridge U., DAMTP

**E-mails** [sw11@damtp.cam.ac.uk](mailto:sw11@damtp.cam.ac.uk)  
[Prof.Hawking@damtp.cam.ac.uk](mailto:Prof.Hawking@damtp.cam.ac.uk)

**Links** <http://www.hawking.org.uk>  
[https://twitter.com/Prof\\_S\\_Haw...](https://twitter.com/Prof_S_Haw...)

**Fields** HEP-TH  
GR-QC

**Identifiers** BA: [S.W.Hawking.1](#)  
INSPIRE: [INSPIRE-00140145](#)  
GoogleScholar: [q74uWAAAAJ](#)

Period	Rank	Institution
1959 - 1962	UG	Oxford U.
1962 - 1965	PHD	Cambridge U.
1965	SENIOR	Cambridge U., DAMTP

[Update Details](#)

## Name Variants

Hawking, Stephen W. (11)  
Hawking, Stephen (10)  
Hawking, S.W. (18)  
Hawking, S. W. (1)  
Hawking, S. (1-1)

## Affiliations

Cambridge U. (125)  
Cambridge U., DAMTP (90)  
Caltech (9)  
Cambridge U., Inst. of Astron. (4)  
Newton Inst. Math. Sci., Cambridge (2)  
Santa Barbara, KITP (2)  
Murch. Max Planck Inst. (1)  
UC, Santa Barbara (1)  
CAMBRIDGE U. (1)  
Tufts U. (1)

## Collaborations

No Collaborations

## PUBLICATIONS AND OUTPUT

## Publications Datasets External

- The Information Paradox for Black Holes
  - Information Preservation and Weather Forecasting for Black Holes
  - Vector Fields in Holographic Cosmology
  - Quantum Probabilities for Inflation from Holography
  - Accelerated Expansion from Negative  $\Lambda$
  - The dreams that stuff is made of: The most astounding papers of quantum physics - and how they shook the scientific world
  - Local Observation in Eternal Inflation
  - The No-Boundary Measure in the Regime of Eternal Inflation
  - General Relativity
  - General Relativity
- [Click here to see all](#)

## Co-Authors

G.W.Gibbons (14)  
T.Horng (13)  
J.B.Harrie (11)  
R.Bousso (11)  
D.N.Page (8)  
C.A.Pope (5)  
H.S.Reall (5)  
C.Horne (4)  
N.Turok (4)  
S.F.Ross (4)  
# more

## Papers

	All papers	Single authored
All papers	224	114
Book	17	7
Conference Paper	52	44
Introductory	12	9
Lectures	10	10
Published	142	50
Review	5	3
Thesis	0	0
Proceedings	3	0

## Subject Categories

Gravitation and Cosmology (128)  
Theory-HEP (64)  
Astrophysics (26)  
General Physics (22)  
Phenomenology-HEP (6)  
Math and Math Physics (1)  
Other (1)

## Frequent Keywords

black hole (43)  
general relativity (38)  
boundary condition (28)  
inflation (27)  
quantum gravity (27)  
cosmological model (19)  
path integral (19)  
wormhole (19)  
field theory, scalar (18)  
BOUNDARY CONDITION (16)  
# more

## STATS

## Citations Summary

224 papers found, 162 of them citable (published or arXiv)

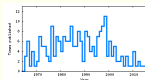
	Citable papers	Published only
<b>Number of papers analyzed:</b>	162	142
<b>Number of citations:</b>	37715	37498
<b>Citations per paper (average):</b>	232.8	264.1
<b><math>h_{2015}</math> index [7]</b>	78	78

## Breakdown of papers by citations:

	Citable papers	Published only
Renowned papers (500+)	16	16
Famous papers (250-499)	18	18
Very well-known papers (100-249)	32	32
Well-known papers (50-99)	27	26
Known papers (10-49)	41	30
Less known papers (1-9)	20	13
Unknown papers (0)	8	2

[Click here to view statistics without self-citations or RPP](#)**Warning:** The citations count should be interpreted with great care. [Read the fine print](#)

## Publication Graph



## Effects of Limited Calorimeter Coverage on ET

Frank E. Paige, A.V. Vanyashin (SSCI)

Mar 1992 · 9 pages

## Search for single $b^*$ -quark production with the ATLAS detector at $\sqrt{s} = 7$ TeV

Planck Inst.), Wainer Vandelli (CERN), Alexandre Vaniachine (Argonne), Peter Vankov (DESY), Francois Vannucci (Paris U., VI-VII), Riccardo Vari (INFN, Rome), Erich Varnes (Arizona U.),

Jan 2013 · 11 pages

Phys.Lett. B721 (2013) 171-189  
(2013-04-25)

DOI: [10.1016/j.physletb.2013.03.016](https://doi.org/10.1016/j.physletb.2013.03.016)  
CERN-PH-EP-2012-344

e-Print: [arXiv:1301.1583](https://arxiv.org/abs/1301.1583) [hep-ex] | PDF  
Experiment: [CERN-LHC-ATLAS](#)

# Jak sobie poradzisz?

## Effects of Limited Calorimeter Coverage on ET

Frank E. Paige, A.V. Vanyashin (SSCI)

Mar 1992 · 9 pages

## Search for single $b^*$ -quark production with the ATLAS detector at $\sqrt{s} = 7$ TeV

Planck Inst.) , Wainer Vandelli (CERN) , Alexandre Vaniachine (Argonne) , Peter Vankov (DESY) , Francois Vannucci (Paris U., VI-VII) , Riccardo Vari (INFN, Rome) , Erich Varnes (Arizona U.) ,

Jan 2013 · 11 pages

Phys.Lett. B721 (2013) 171-189  
(2013-04-25)

DOI: [10.1016/j.physletb.2013.03.016](https://doi.org/10.1016/j.physletb.2013.03.016)  
CERN-PH-EP-2012-344

e-Print: [arXiv:1301.1583](https://arxiv.org/abs/1301.1583) [hep-ex] | PDF  
Experiment: [CERN-LHC-ATLAS](#)

✓ Ten sam autor

# Homonimy wśród chińskich nazwisk pisanych w języku angielskim

Poznajcie Yanga Wanga, Wanga Yanga i Yanga Wanga!

杨阳



杨洋



杨旻



# Definicije

## Publications



## Signatures



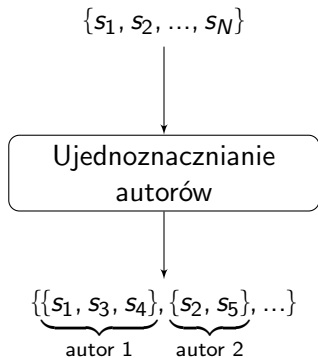
## Signature for Doe, John

<b>Title</b>	Lorem ipsum dolor sit amet, consectetur adipiscing elit
<b>Author</b>	Doe, John
<b>Affiliation</b>	University of Foo
<b>Co-authors</b>	Smith, John; Chen, Wang
<b>Year</b>	2015

Zgrupuj publikacje każdego autora, i tylko jego.

Inspirehep.net jest biblioteką cyfrową, która zawiera:

- 1M publikacji zawierających 10M sygnatur
- 1.2M sygnatur jest oznaczonych:
  - Osobiście przez autorów (podobnie jak w Google Scholar).
  - Za pomocą identyfikatorów (ORCID).
  - Przez profesjonalnych kustoszy.



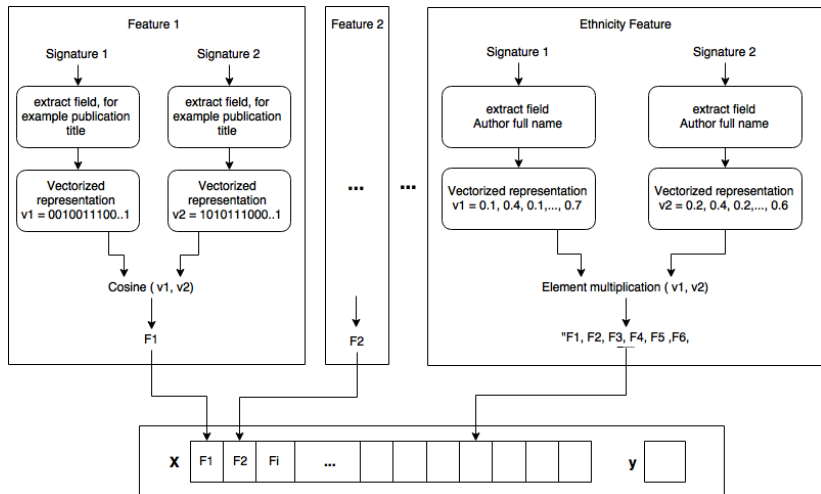


- Ręczne ujednoznacznianie jest **źmudne i trudne**, nawet dla doświadczonych kustoszy.
- Czy nie moglibyśmy **automatycznie znaleźć zbioru reguł**, które ujednoznaczą parę sygnatur?

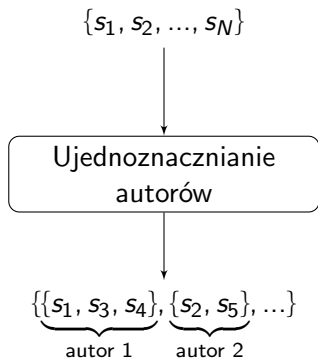
$$\varphi(s_1, s_2) = \begin{cases} 0 & \text{jeśli } s_1 \text{ i } s_2 \text{ należą do tego samego autora,} \\ 1 & \text{wpp.} \end{cases}$$

- Jest to przykład **uczenia nadzorowanego**.

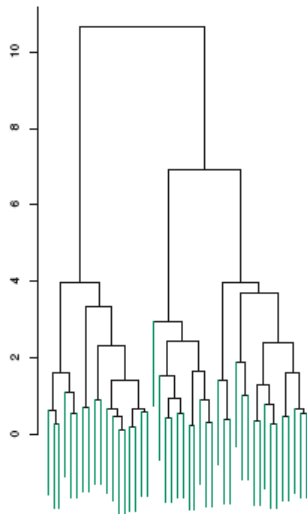
# Tworzenie cech na podstawie par sygnatur



# Ujednoznacznianie jako problem klastrowania



# Aglomeracyjny klastrowanie hierarchiczne



- Rodzina algorytmów klastrujących, która **buduje klastry iteracyjnie je łącząc**.
- Hierarchia klastrow może być reprezentowana jako dendrogram.
- Korzeń drzewa jest klastrem zawierającym wszystkie sygnatury, liście to pojedyncze sygnatury

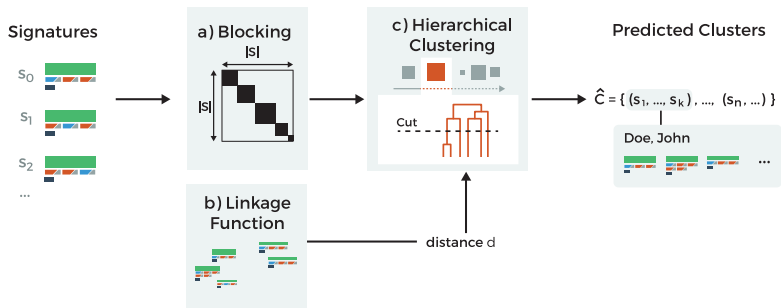
- Złożoność tego klastrowania to  $O(N^2)$ . Mamy  $N = 10^7$  sygnatur.

*Rozwiązanie:* podzielenie sygnatur na bloki i następnie potraktowanie tych bloków jako oddzielne zagadnienia ujednoznaczniania. Można np. grupować razem sygnatury z tym samym nazwiskiem.

- Jak zdecydować jak wysoko **przeciąć dendogram?**

*Rozwiązanie:* używając oznaczonych sygnatur, wybieramy wysokość, która maksymalizuje wynik  $B^3$  F-measure.

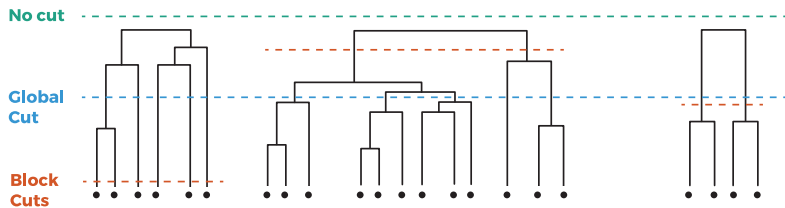
# Cały algorytm



Przeprowadzone analizy pokazują, że jest parę przypadków, gdzie grupowanie na podstawie nazwiska nie działa:

- "Mueller, R." i "Muller, R."
- "Martinez Torres, A." i "Torres, A. Martinez"
- "Smith-Jones, A." i "Smith, A."
- Nazwisko autora zmieniło się (np. przez ślub).

# Jak wysoko ciąć dendrogramy





Opis	$B^3$		
	Prec.	Recall	F1
Prosty algorytm	0.9024	0.9828	0.9409
<u>Grupowanie = Nazwisko i pierwszy inicjał</u>	0.9901	0.9760	0.9830
Grupowanie = Double metaphone	0.9856	0.9827	0.9841
Grupowanie = NYSIIS	0.9875	0.9826	<b>0.9850</b>
Grupowanie = Soundex	0.9886	0.9745	0.9815
<u>Klasyfikator = Gradient Boosting DT</u>	0.9901	0.9760	0.9830
Klasyfikator = Lasy losowe	0.9909	0.9783	<b>0.9846</b>
Klasyfikator = Regresja Liniowa	0.9749	0.9584	0.9666

Opis	$B^3$		
	<i>Prec.</i>	<i>Recall</i>	<i>F1</i>
Prosty algorytm	0.9024	0.9828	0.9409
<u>Klastrowanie = Average linkage</u>	0.9901	0.9760	<b>0.9830</b>
Klastrowanie = Single linkage	0.9741	0.9603	0.9671
Klastrowanie = Complete linkage	0.9862	0.9709	0.9785
No cut (prosty algorytm)	0.9024	0.9828	0.9409
Global cut	0.9892	0.9737	0.9814
<u>Block cut</u>	0.9901	0.9760	<b>0.9830</b>
<b>Najlepsze ustawienia razem</b>	0.9888	0.9848	<b>0.9868</b>

Rozwiązanie jest obecnie używane w CERNie przez projekty INSPIRE i INVENIO.

Czas wykonania: 20 godzin dla 10M sygnatur, na 16 rdzeniach i 32GB RAMu.

**Ale** ujednoznacznianie inkrementalne zajmuje tylko parę minut!

Nasze rozwiązanie jest open-source<sup>1</sup> i udostępniliśmy nasz zbiór danych<sup>2</sup>.

---

<sup>1</sup>[github.com/inspirehep/beard](https://github.com/inspirehep/beard)

<sup>2</sup>[github.com/gloupppe/paper-author-disambiguation/data](https://github.com/gloupppe/paper-author-disambiguation/data)

Dziękuję za uwagę!