

# Gaming the H-index

Hugo Jonker, Sjouke Mauw

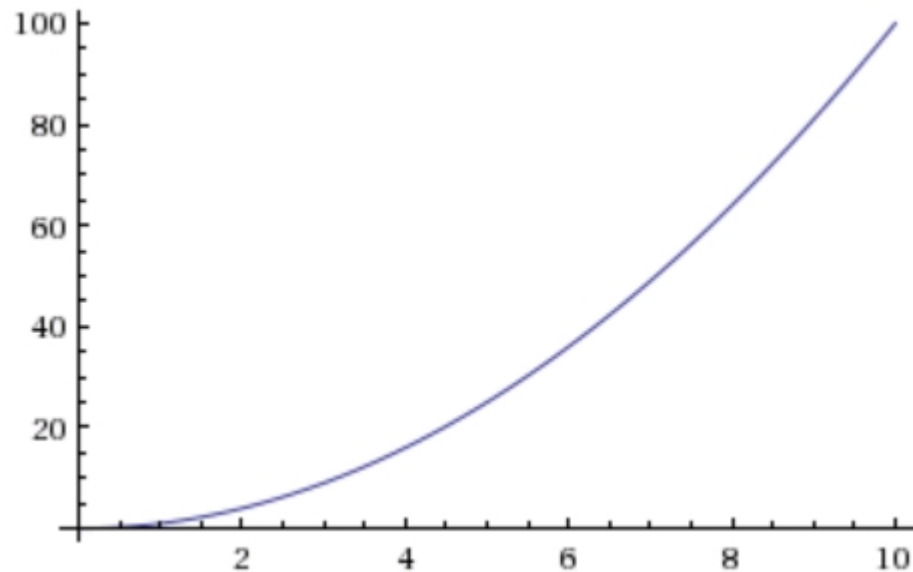
**Open Universiteit**

[www.ou.nl](http://www.ou.nl)



# H-index

The h-index of  $A$  is the maximum number  $N$  such that  $A$  (co-)authored at least  $N$  papers, each of which was cited at least  $N$  times.



# Goodhart's Law [Goodhart75]

When a measure becomes a target,  
it ceases to be a good measure.

*As soon as the government attempts to regulate any particular set of financial assets,  
these become unreliable as indicators of economic trends.*

# Attacking the publication process

Various attacks, both real-life<sup>1</sup> and scientific<sup>2</sup>

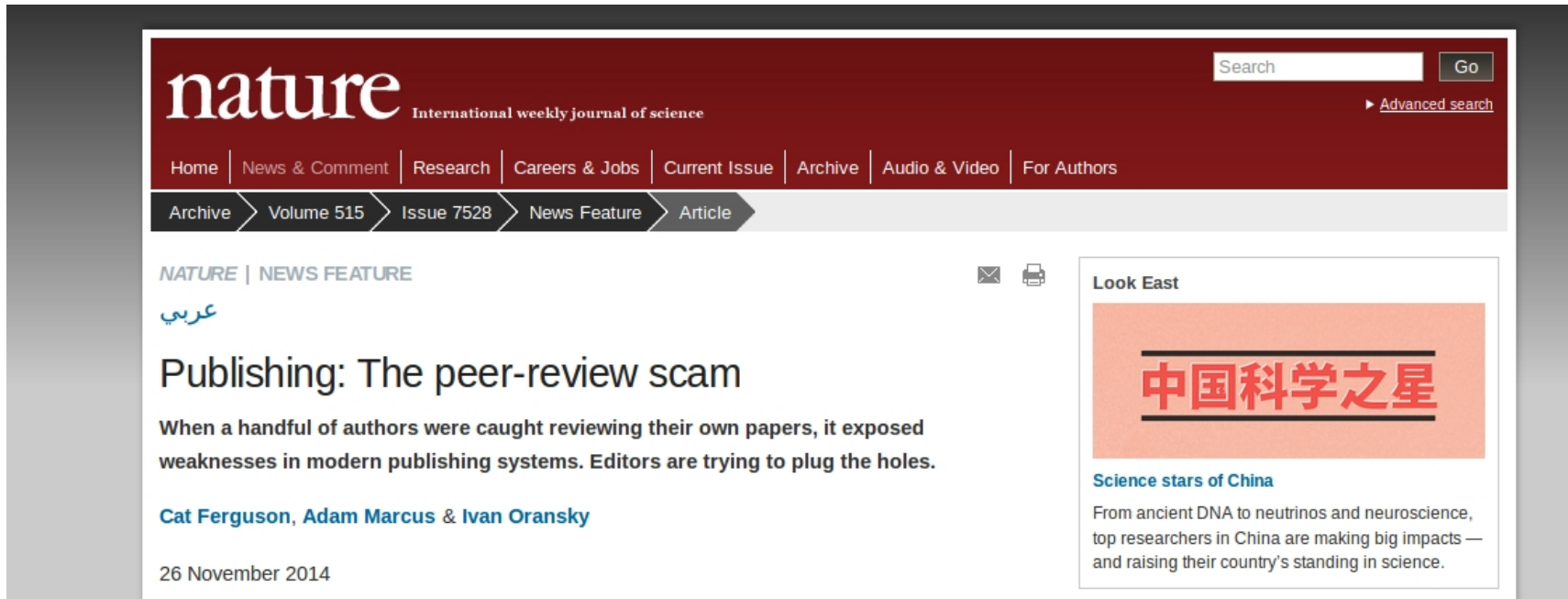
- Faking experimental data
- Manipulating peer reviews
- Evading peer reviews
- Randomly generated papers

<sup>1</sup> cf. e.g. <http://www.the-scientist.com/?articles.view/articleNo/41777/title/The-Top-10-Retractions-of-2014/>

<sup>2</sup> e.g. C. Labbé and D. Labbé, “Duplicate and fake publications in the scientific literature: how many SciGen papers in computer science?” *Scientometrics*, vol. 94, no. 1, pp. 379–396, 2013.

# Not incidental

- SAGE: “peer review and **citation ring**”
- Thomson Reuters: **citation stacking**.
- Hyung-in Moon:



The screenshot shows the Nature journal website interface. At the top, the 'nature' logo is displayed in white on a dark red background, with the tagline 'International weekly journal of science' below it. A search bar with a 'Go' button and a link to 'Advanced search' is located in the top right corner. A navigation menu below the logo includes links for Home, News & Comment, Research, Careers & Jobs, Current Issue, Archive, Audio & Video, and For Authors. Below the navigation menu, a breadcrumb trail shows the path: Archive > Volume 515 > Issue 7528 > News Feature > Article. The main content area features the text 'NATURE | NEWS FEATURE' and a blue Arabic script icon. The article title is 'Publishing: The peer-review scam', followed by a summary: 'When a handful of authors were caught reviewing their own papers, it exposed weaknesses in modern publishing systems. Editors are trying to plug the holes.' The authors listed are 'Cat Ferguson, Adam Marcus & Ivan Oransky', and the date is '26 November 2014'. On the right side, there is a 'Look East' sidebar with a red background and the Chinese characters '中国科学之星' (China Science Stars) in white. Below this, the text reads 'Science stars of China' and 'From ancient DNA to neutrinos and neuroscience, top researchers in China are making big impacts — and raising their country's standing in science.'

# Gaming the H-index

- artificial manipulations to improve one's scientific standing
  - Bona fide effort increases quadratically
- focus on author metrics (like h-index)
- applies equally to venue metrics (like impact factor)

# What is gaming?

- gaming: exploiting weaknesses
- BUT: gaming != hacking.
  - **Hacking**: exploiting weaknesses in implementation
  - **Gaming**: exploiting weaknesses in specification

# Publication metrics weaknesses

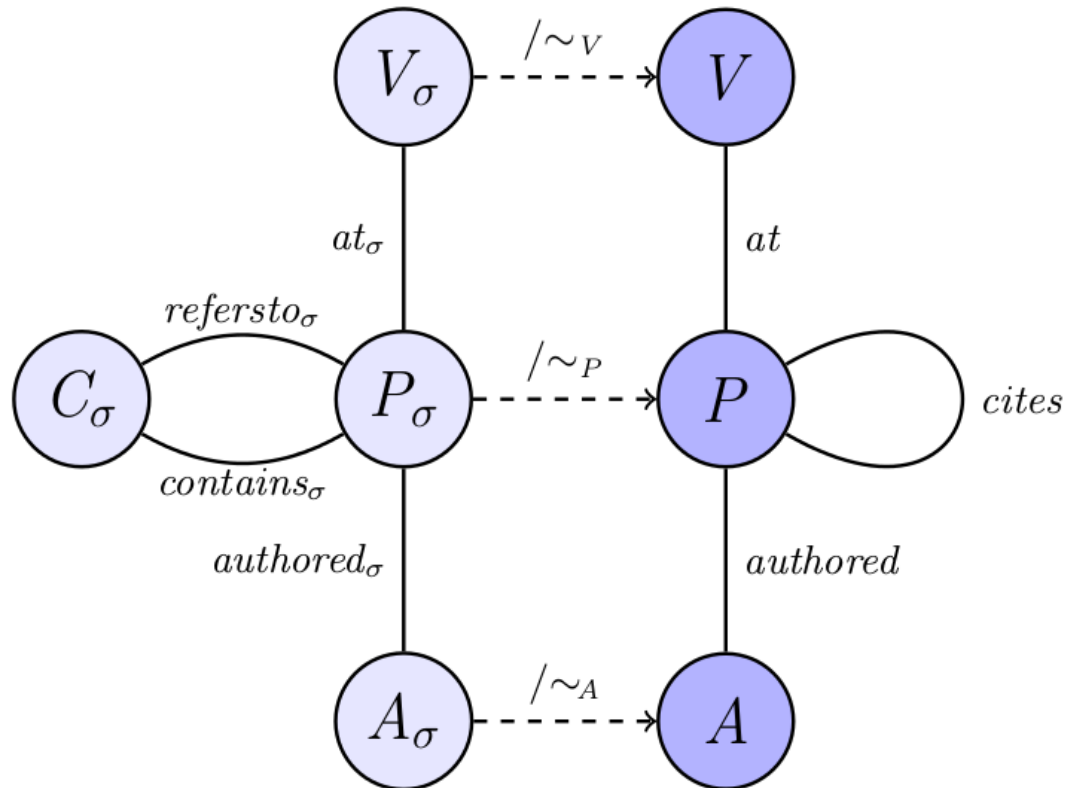
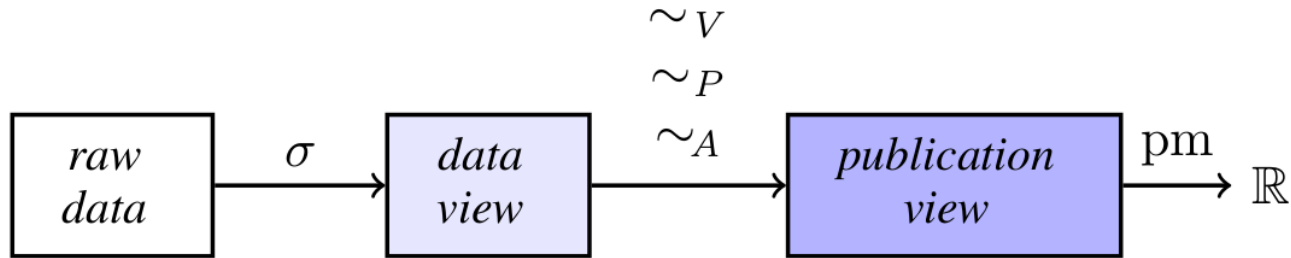
	<i>Implementation</i>	<i>Design</i>
<i>Errors</i>	Implementation errors	Methodological drawbacks
<i>Attacks</i>	Hacking	Gaming



# Possible gaming attacks?

- Model publication process
- Formal definition of publication metrics
- Derive attack surface from formal definition

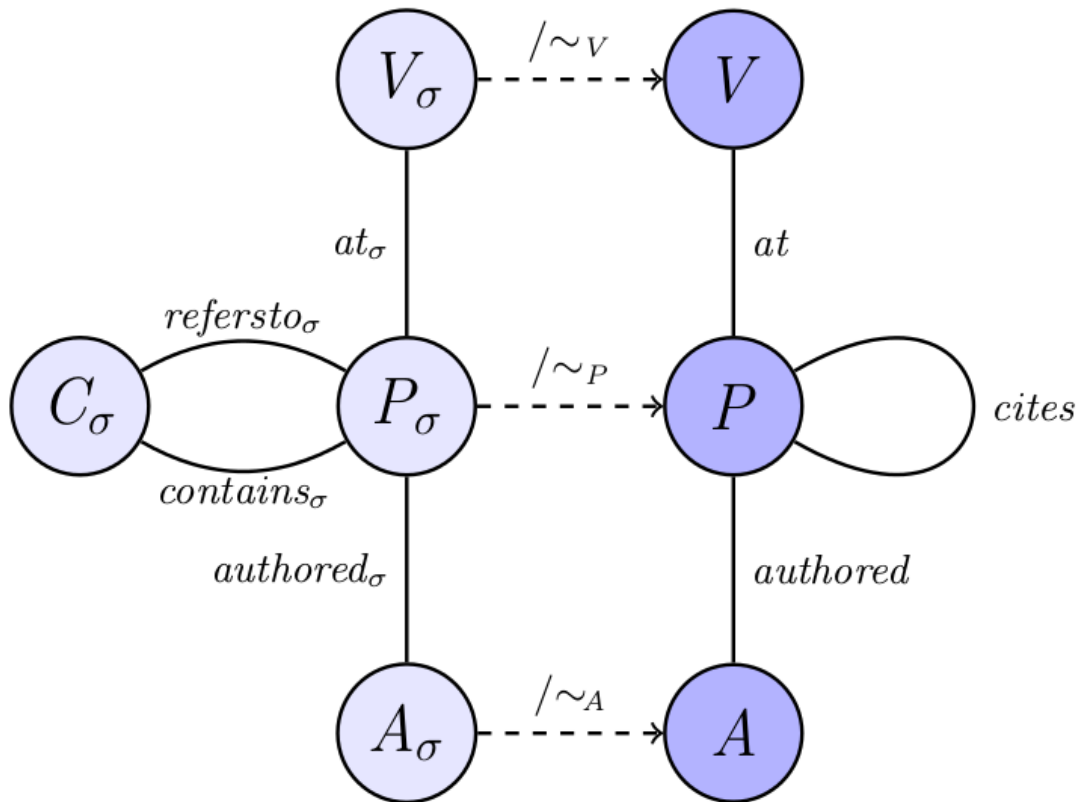
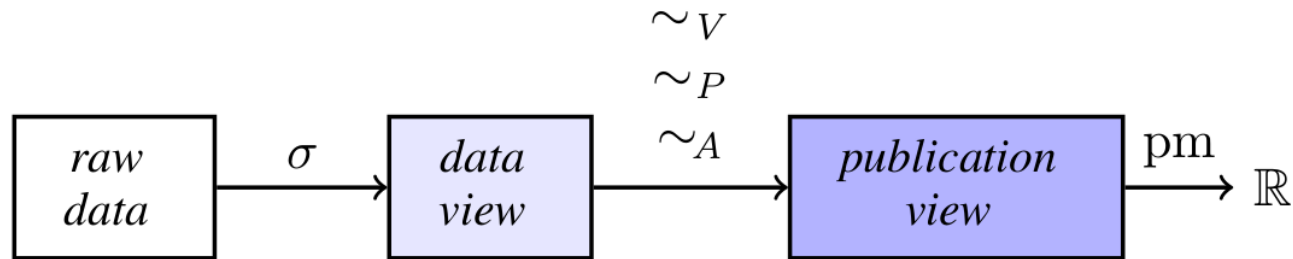
# Publication structure



# Examples

- $h\text{-index}(a) = \max\{ i \in \mathbb{N} \mid \exists T \subseteq \text{pubs}(a) |T| = i \wedge \forall p \in T \#citing(p) \geq i \}$ .
- $i10\text{-index}(a) = |\{ p \in P \mid \text{authored}(a, p) \wedge \#citing(p) \geq 10 \}|$ .
- $g\text{-index}(a) = \max\{ i \in \mathbb{N} \mid \exists T \subseteq \text{pubs}(a) |T| = i \wedge \sum_{p \in T} \#citing(p) \geq i^2 \}$ .
- $\text{AR}(v) = \frac{|\{ p \in P \mid \text{at}(p, v) \}|}{|\{ p \in P \mid \text{submitted-to}(p, v) \}|}$ .

# Theoretical “gaming surface”



# In practice

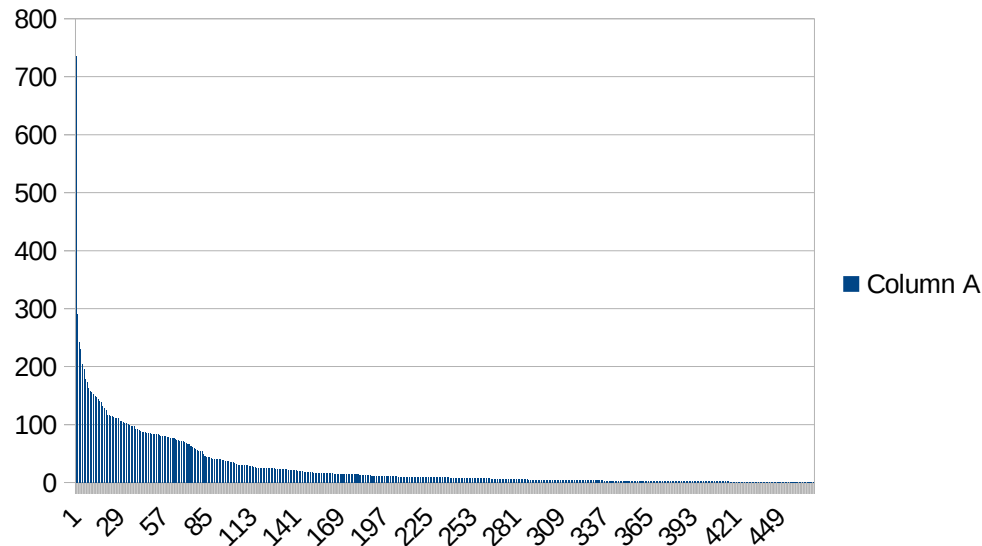
- Collect data
- Understand typical change rates
- Find outliers

# Suspect #1



## Scientific production and originality

More than 3000 citations, h-index 30 (due to quality and quantity of recent works these figures will probably be doubled in the next 2 years)



# Initial findings

- 2 extremely similar papers found  
Difference: ~10 words per page
- 1 paper from 2016 cited in 2015... hmms.
  - ~50 times cited in 2015
  - 104 times cited by June 2016.
  - 80+ from one journal
  - 10+ from one other journal
  - Guess who is editor there?

# Suspect #2

- Self-citations
- Book editor → chapter editor
  - More self-citations
- Publishing almost identical paper twice





# Detection algorithm

- Find strange h-curve
- Find journals where suspect is editor
  - Many publications in own journal?
  - Many self-citations?
  - Many citations from own journal?

# Next steps

- Msc thesis:
  - Niels van Tielenburg
  - Automating detection of fraud
- Paper on h-index+manual fraud detection
- Paper on automating fraud detection
- ...?

Bedankt voor de aandacht!



- We worden geregeerd door uiterlijkheden: uiterlijk goed = trust is goed.
- De wetenschappelijke methode is prima
- Maar de implementatie is niet optimaal: zoveel mogelijkheden voor ruis, dat mensen van deze ruimte misbruik kunnen maken om wetenschappelijker te lijken dan ze zijn.
- Wegens enorme grootte wetenschappelijke wereld hebben we een trustmodel en methode nodig om een oordeel te vellen over kwaliteit van wetenschap.
  - cf. climate change debate
- Op lange termijn is wetenschappelijke methode okee(?), op korte termijn niet. Vergelijk met discussie over wetenschappelijke vooruitgang: helemaal niet smooth, maar met horten en stoten en bruised ego's.