

# An index to quantify an individual's scientific research output

J. E. Hirsch\*

Department of Physics, University of California at San Diego, La Jolla, CA 92093-0319

Communicated by Manuel Cardona, Max Planck Institute for Solid State Research, Stuttgart, Germany, September 1, 2005 (received for review August 15, 2005)

**I propose the index  $h$ , defined as the number of papers with citation number  $\geq h$ , as a useful index to characterize the scientific output of a researcher.**

citations | impact | unbiased

For the few scientists who earn a Nobel prize, the impact and relevance of their research is unquestionable. Among the rest of us, how does one quantify the cumulative impact and relevance of an individual's scientific research output? In a world of limited resources, such quantification (even if potentially distasteful) is often needed for evaluation and comparison purposes (e.g., for university faculty recruitment and advancement, award of grants, etc.).

The publication record of an individual and the citation record clearly are data that contain useful information. That information includes the number ( $N_p$ ) of papers published over  $n$  years, the number of citations ( $N_c$ ) for each paper ( $j$ ), the journals where the papers were published, their impact parameter, etc. This large amount of information will be evaluated with different criteria by different people. Here, I would like to propose a single number, the " $h$  index," as a particularly simple and useful way to characterize the scientific output of a researcher.

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.

The research reported here concentrated on physicists; however, I suggest that the  $h$  index should be useful for other scientific disciplines as well. (At the end of the paper I discuss some observations for the  $h$  index in biological sciences.) The highest  $h$  among physicists appears to be E. Witten's  $h$ , which is 110. That is, Witten has written 110 papers with at least 110 citations each. That gives a lower bound on the total number of citations to Witten's papers at  $h^2 = 12,100$ . Of course, the total number of citations ( $N_{c,tot}$ ) will usually be much larger than  $h^2$ , because  $h^2$  both underestimates the total number of citations of the  $h$  most-cited papers and ignores the papers with  $< h$  citations. The relation between  $N_{c,tot}$  and  $h$  will depend on the detailed form of the particular distribution (1), and it is useful to define the proportionality constant  $a$  as

$$N_{c,tot} = ah^2. \quad [1]$$

I find empirically that  $a$  ranges between 3 and 5.

Other prominent physicists with high  $h$ s are A. J. Heeger ( $h = 107$ ), M. L. Cohen ( $h = 94$ ), A. C. Gossard ( $h = 94$ ), P. W. Anderson ( $h = 91$ ), S. Weinberg ( $h = 88$ ), M. E. Fisher ( $h = 88$ ), M. Cardona ( $h = 86$ ), P. G. deGennes ( $h = 79$ ), J. N. Bahcall ( $h = 77$ ), Z. Fisk ( $h = 75$ ), D. J. Scalapino ( $h = 75$ ), G. Parisi ( $h = 73$ ), S. G. Louie ( $h = 70$ ), R. Jackiw ( $h = 69$ ), F. Wilczek ( $h = 68$ ), C. Vafa ( $h = 66$ ), M. B. Maple ( $h = 66$ ), D. J. Gross ( $h = 66$ ), M. S. Dresselhaus ( $h = 62$ ), and S. W. Hawking ( $h = 62$ ). I argue that  $h$  is preferable to other single-number criteria commonly used to evaluate scientific output of a researcher, as follows:

- (i) Total number of papers ( $N_p$ ). Advantage: measures productivity. Disadvantage: does not measure importance or impact of papers.
- (ii) Total number of citations ( $N_{c,tot}$ ). Advantage: measures total impact. Disadvantage: hard to find and may be inflated by a small number of "big hits," which may not be representative of the individual if he or she is a coauthor with many others on those papers. In such cases, the relation in Eq. 1 will imply a very atypical value of  $a$ ,  $> 5$ . Another disadvantage is that  $N_{c,tot}$  gives undue weight to highly cited review articles versus original research contributions.
- (iii) Citations per paper (i.e., ratio of  $N_{c,tot}$  to  $N_p$ ). Advantage: allows comparison of scientists of different ages. Disadvantage: hard to find, rewards low productivity, and penalizes high productivity.
- (iv) Number of "significant papers," defined as the number of papers with  $> y$  citations (for example,  $y = 50$ ). Advantage: eliminates the disadvantages of criteria *i*, *ii*, and *iii* and gives an idea of broad and sustained impact. Disadvantage:  $y$  is arbitrary and will randomly favor or disfavor individuals, and  $y$  needs to be adjusted for different levels of seniority.
- (v) Number of citations to each of the  $q$  most-cited papers (for example,  $q = 5$ ). Advantage: overcomes many of the disadvantages of the criteria above. Disadvantage: It is not a single number, making it more difficult to obtain and compare. Also,  $q$  is arbitrary and will randomly favor and disfavor individuals.

Instead, the proposed  $h$  index measures the broad impact of an individual's work, avoids all of the disadvantages of the criteria listed above, usually can be found very easily by ordering papers by "times cited" in the Thomson ISI Web of Science database (<http://isiknowledge.com>),<sup>†</sup> and gives a ballpark estimate of the total number of citations (Eq. 1).

Thus, I argue that two individuals with similar  $h$ s are comparable in terms of their overall scientific impact, even if their total number of papers or their total number of citations is very different. Conversely, comparing two individuals (of the same scientific age) with a similar number of total papers or of total citation count and very different  $h$  values, the one with the higher  $h$  is likely to be the more accomplished scientist.

For a given individual, one expects that  $h$  should increase approximately linearly with time. In the simplest possible model, assume that the researcher publishes  $p$  papers per year and that each published paper earns  $c$  new citations per year every subsequent year. The total number of citations after  $n + 1$  years is then

$$N_{c,tot} = \sum_{j=1}^n pcj = \frac{pcn(n+1)}{2}. \quad [2]$$

\*E-mail: [jhirsch@ucsd.edu](mailto:jhirsch@ucsd.edu).

<sup>†</sup>Of course, the database used must be complete enough to cover the full period spanned by the individual's publications.

© 2005 by The National Academy of Sciences of the USA

Assuming all papers up to year  $y$  contribute to the index  $h$ , we have

$$(n - y)c = h \quad [3a]$$

$$py = h. \quad [3b]$$

The left side of Eq. 3a is the number of citations to the most recent of the papers contributing to  $h$ ; the left side of Eq. 3b is the total number of papers contributing to  $h$ . Hence, from Eq. 3,

$$h = \frac{c}{1 + c/p} n. \quad [4]$$

The total number of citations (for not-too-small  $n$ ) is then approximately

$$N_{c,tot} \sim \frac{(1 + c/p)^2}{2c/p} h^2 \quad [5]$$

of the form Eq. 1. The coefficient  $a$  depends on the number of papers and the number of citations per paper earned per year as given by Eq. 5. As stated earlier, we find empirically that  $a \approx 3$ –5 is a typical value. The linear relation

$$h \sim mn \quad [6]$$

should hold quite generally for scientists who produce papers of similar quality at a steady rate over the course of their careers; of course,  $m$  will vary widely among different researchers. In the simple linear model,  $m$  is related to  $c$  and  $p$  as given by Eq. 4. Quite generally, the slope of  $h$  versus  $n$ , the parameter  $m$ , should provide a useful yardstick to compare scientists of different seniority.

In the linear model, the minimum value of  $a$  in Eq. 1 is  $a = 2$ , for the case  $c = p$ , where the papers with  $>h$  citations and those with  $<h$  citations contribute equally to the total  $N_{c,tot}$ . The value of  $a$  will be larger for both  $c > p$  and  $c < p$ . For  $c > p$ , most contributions to the total number of citations arise from the “highly cited papers” (the  $h$  papers that have  $N_c > h$ ), whereas for  $c < p$ , it is the sparsely cited papers (the  $N_p - h$  papers that have  $<h$  citations each) that give the largest contribution to  $N_{c,tot}$ . We find that the first situation holds in the vast majority of, if not all, cases. For the linear model defined in this example,  $a = 4$  corresponds to  $c/p = 5.83$  (the other value that yields  $a = 4$ ,  $c/p = 0.17$ , is unrealistic).

The linear model defined above corresponds to the distribution

$$N_c(y) = N_0 - \left(\frac{N_0}{h} - 1\right)y, \quad [7]$$

where  $N_c(y)$  is the number of citations to the  $y$ th paper (ordered from most cited to least cited) and  $N_0$  is the number of citations of the most highly cited paper ( $N_0 = cn$  in the example above). The total number of papers  $y_m$  is given by  $N_c(y_m) = 0$ ; hence,

$$y_m = \frac{N_0 h}{N_0 - h}. \quad [8]$$

We can write  $N_0$  and  $y_m$  in terms of  $a$  defined in Eq. 1 as

$$N_0 = h \left[ a \pm \sqrt{a^2 - 2a} \right] \quad [9a]$$

$$y_m = h \left[ a \mp \sqrt{a^2 - 2a} \right]. \quad [9b]$$

For  $a = 2$ ,  $N_0 = y_m = 2h$ . For larger  $a$ , the upper sign in Eq. 9 corresponds to the case where the highly cited papers dominate

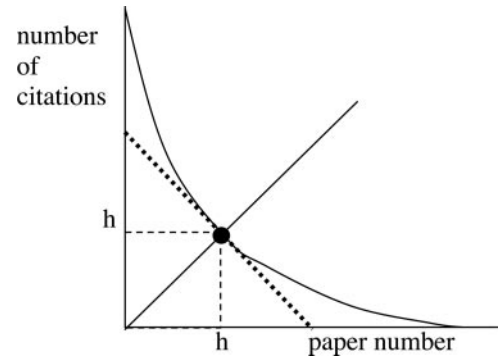


Fig. 1. Schematic curve of number of citations versus paper number, with papers numbered in order of decreasing citations. The intersection of the 45° line with the curve gives  $h$ . The total number of citations is the area under the curve. Assuming the second derivative is nonnegative everywhere, the minimum area is given by the distribution indicated by the dotted line, yielding  $a = 2$  in Eq. 1.

(the more realistic case), and the lower sign corresponds to the case where the less frequently cited papers dominate the total citation count.

In a more realistic model,  $N_c(y)$  will not be a linear function of  $y$ . Note that  $a = 2$  can safely be assumed to be a lower bound quite generally, because a smaller value of  $a$  would require the second derivative  $\partial^2 N_c / \partial y^2$  to be negative over large regions of  $y$ , which is not realistic. The total number of citations is given by the area under the  $N_c(y)$  curve that passes through the point  $N_c(h) = h$ . In the linear model, the lowest  $a = 2$  corresponds to the line of slope  $-1$ , as shown in Fig. 1.

A more realistic model would be a stretched exponential of the form

$$N_c(y) = N_0 e^{-\left(\frac{y}{y_0}\right)^\beta}. \quad [10]$$

Note that for  $\beta \leq 1$ ,  $N_c''(y) > 0$  for all  $y$ ; hence,  $a > 2$  is true. We can write the distribution in terms of  $h$  and  $a$  as

$$N_c(y) = \frac{a}{\alpha I(\beta)} h e^{-\left(\frac{y}{h\alpha}\right)^\beta} \quad [11]$$

with  $I(\beta)$  the integral

$$I(\beta) = \int_0^\infty dz e^{-z^\beta} \quad [12]$$

and  $\alpha$  determined by the equation

$$\alpha e^{\alpha^{-\beta}} = \frac{a}{I(\beta)}. \quad [13]$$

The maximally cited paper has citations

$$N_0 = \frac{a}{\alpha I(\beta)} h, \quad [14]$$

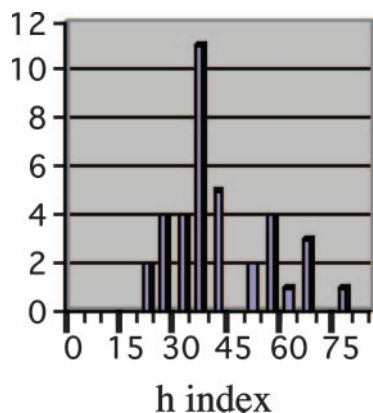
and the total number of papers (with at least one citation) is determined by  $N(y_m) = 1$  as

$$y_m = h[1 + \alpha^\beta \ln(h)]^{1/\beta}. \quad [15]$$

A given researcher's distribution can be modeled by choosing the most appropriate  $\beta$  and  $a$  for that case. For example, for  $\beta = 1$ , if  $a = 3$ ,  $\alpha = 0.661$ ,  $N_0 = 4.54h$ , and  $y_m = h[1 + .66 \ln h]$ . With  $a = 4$ ,  $\alpha = 0.4644$ ,  $N_0 = 8.61h$ , and  $y_m = h[1 + 0.46 \ln(h)]$ . For







**Fig. 2.** Histogram giving the number of Nobel prize recipients in physics in the last 20 years versus their  $h$  index. The peak is at the  $h$  index between 35 and 39.

35, lower than the mean due to the tail for high  $h$  values. It is interesting that Nobel prize winners have substantial  $h$  indices (84% had an  $h$  of at least 30), indicating that Nobel prizes do not originate in one stroke of luck but in a body of scientific work. Notably, the values of  $m$  found are often not high compared with other successful scientists (49% of our sample had  $m < 1$ ), clearly because Nobel prizes are often awarded long after the period of maximum productivity of the researchers.

As another example, among newly elected members of the National Academy of Sciences in physics and astronomy in 2005, I find  $\langle h \rangle = 44$ ,  $\sigma_h = 14$ , highest  $h = 71$ , lowest  $h = 20$ , and median  $h_m = 46$ . Among the total membership in the National Academy of Sciences in physics, the subgroup of last names starting with “A” and “B” has  $\langle h \rangle = 38$ ,  $\sigma_h = 10$ , and  $h_m = 37$ . These examples further indicate that the index  $h$  is a stable and consistent estimator of scientific achievement.

An intriguing idea is the extension of the  $h$ -index concept to groups of individuals.<sup>‡</sup> The SPIRES high-energy physics literature database ([www.slac.stanford.edu/spires/hep](http://www.slac.stanford.edu/spires/hep)) recently implemented the  $h$  index in their citation summaries, and it also allows the computation of  $h$  for groups of scientists. The overall  $h$  index of a group will generally be larger than that of each of the members of the group but smaller than the sum of the individual  $h$  indices, because some of the papers that contribute to each individual’s  $h$  will no longer contribute to the group’s  $h$ . For example, the overall  $h$  index of the condensed matter group at the University of California at San Diego physics department

is  $h = 118$ , of which the largest individual contribution is 25; the highest individual  $h$  is 66, and the sum of individual  $h$ s is  $>300$ . The contribution of each individual to the group’s  $h$  is not necessarily proportional to the individual’s  $h$ , and the highest contributor to the group’s  $h$  will not necessarily be the individual with highest  $h$ . In fact, in principle (although rarely in practice), the lowest- $h$  individual in a group could be the largest contributor to the group’s  $h$ . For a prospective graduate student considering different graduate programs, a ranking of groups or departments in his or her chosen area according to their overall  $h$  index would likely be of interest, and for administrators concerned with these issues, the ranking of their departments or entire institution according to the overall  $h$  could also be of interest.

To conclude, I discuss some observations in the fields of biological and biomedical sciences. From the list compiled by Christopher King of Thomson ISI of the most highly cited scientists in the period 1983–2002 (5), I found the  $h$  indices for the top 10 on that list, all in the life sciences, which are, in order of decreasing  $h$ : S. H. Snyder,  $h = 191$ ; D. Baltimore,  $h = 160$ ; R. C. Gallo,  $h = 154$ ; P. Chambon,  $h = 153$ ; B. Vogelstein,  $h = 151$ ; S. Moncada,  $h = 143$ ; C. A. Dinarello,  $h = 138$ ; T. Kishimoto,  $h = 134$ ; R. Evans,  $h = 127$ ; and A. Ullrich,  $h = 120$ . It can be seen that, not surprisingly, all of these highly cited researchers also have high  $h$  indices and that high  $h$  indices in the life sciences are much higher than in physics. Among 36 new inductees in the National Academy of Sciences in biological and biomedical sciences in 2005, I find  $\langle h \rangle = 57$ ,  $\sigma_h = 22$ , highest  $h = 135$ , lowest  $h = 18$ , and median  $h_m = 57$ . These latter results confirm that  $h$  indices in biological sciences tend to be higher than in physics; however, they also indicate that the difference appears to be much higher at the high end than on average. Clearly, more research in understanding similarities and differences of  $h$  index distributions in different fields of science would be of interest.

In summary, I have proposed an easily computable index,  $h$ , which gives an estimate of the importance, significance, and broad impact of a scientist’s cumulative research contributions. I suggest that this index may provide a useful yardstick with which to compare, in an unbiased way, different individuals competing for the same resource when an important evaluation criterion is scientific achievement.

I am grateful to many colleagues in the University of California at San Diego Condensed Matter group and especially Ivan Schuller for stimulating discussions on these topics and encouragement to publish these ideas. I also thank the many readers who wrote with interesting comments since this paper was first posted at arXiv.org (6); the referees who made constructive suggestions, all of which led to improvements in the paper; and Travis Brooks and the SPIRES database administration for rapidly implementing the  $h$  index in their database.

<sup>‡</sup>This was first introduced in the SPIRES database.

1. Laherrere, J. & Sornette, D. (1998) *Eur. Phys. J. E Soft Matter* **B2**, 525–539.
2. Redner, S. (1998) *Eur. Phys. J. E Soft Matter* **B4**, 131–134.
3. Redner, S. (2005) *Phys. Today* **58**, 49–54.

4. van Raan, A. F. J. (2004) *Scientometrics* **59**, 467–472.
5. King, C. (2003) *Sci. Watch* **14**, no. 5, 1.
6. Hirsch, J. E. (2005) *arXiv.org E-Print Archive* (Aug. 3, 2005). Available at <http://arxiv.org/abs/physics/0508025>.