1  **Comparison of the *h*-index scores among pathogens identified as emerging hazards in**
2  **North America**
3
4  **Ruth Cox[1], K. Marie McIntyre[2], Javier Sanchez[1], Christian Setzkorn[2], Matthew Baylis[2],**
5  **Crawford W. Revie[1]**
6
7  Author affiliations
8  [1]Centre for Veterinary Epidemiological Research, Atlantic Veterinary College, University of
9  Prince Edward Island, 550 University Avenue, Charlottetown, Canada
10
11  [2]Department of Epidemiology and Population Health, Institute of Infection and Global Health,
12  University of Liverpool, Liverpool, UK
13
14  Corresponding author:
15  Ruth Cox
16  Centre for Veterinary Epidemiological Research, Atlantic Veterinary College, University of
17  Prince Edward Island, 550 University Avenue, Charlottetown, Canada, C1A 4P3.
18  Email: rucox@upei.ca
19  Telephone: +1 902 566 0815
20  Fax: +1 902 620 5053
21
22  Short title: The Hirsch index as a novel method to rank pathogens

1

23  **Summary**

24  Disease surveillance must assess the relative importance of pathogen hazards. Here we use the

25  Hirsch index (*h*-index) as a novel method to identify and rank infectious pathogens that are likely

26  to be a hazard to human health in the North American region. This bibliometric index was

27  developed to quantify an individual's scientific research output and was recently used as a proxy

28  measure for pathogen impact. Analysis of more than 3000 infectious organisms indicated that

29  651 were human pathogen species that had been recorded in the North American region. The *h*-

30  index of these pathogens ranged from 0 to 584. The *h*-index of emerging pathogens was greater

31  than non-emerging pathogens as was the *h*-index of frequently pathogenic pathogens when

32  compared to non-pathogenic pathogens.  As expected the *h*-index of pathogens varied over time

33  between 1960 and 2011. We discuss how the *h*-index can contribute to pathogen prioritisation

34  and as an indicator of pathogen emergence.

35

36  **Keywords**: *h*-index, bibliometric, pathogen hazard, pathogen prioritisation, pathogen emergence

37

38  **Introduction**

39  Effective disease surveillance and control rely on an ability to assess the relative importance of

40  diseases and pathogens. Such prioritisation often involves the use of decision support tools to

41  identify which diseases to target and where to focus resources and funding. These can be biased

42  by the quality of evidence utilised, time taken for its collection and therefore the timeliness of

43  results, or by the opinion of experts employed to make judgements on topics (McIntyre et al,

44  2011). In this work we use a quantitative method to identify and compare pathogens that are

45  hazardous to human health in the North America region, which is quick and relatively simple to

2

46  calculate, and we consider whether it might be used as a method to rank pathogens according to

47  their impact. This novel method involves the use of the Hirsch index ($h$-index); a bibliometric

48  index which was originally developed to quantify an individual's scientific research output

49  (Hirsch, 2005), by accounting for the number of publications produced and the number of

50  citations of those publications. "A scientist has an index $h$ if $h$ of his or her $Np$ publications have

51  at least $h$ citations each and the other (Np-$h$) papers have $\leq h$ citations each" (Hirsch, 2005). It

52  can be calculated using a range of bibliometric services such as those available from the Institute

53  for Scientific Information's Web of Science (WOS) (Thomson Reuters, 2011), and is given as a

54  standard metric in output generated by Google Scholar (Google, 2013). While the $h$-index was

55  initially devised to assess the output of individual scholars, it has been extended to measure the

56  productivity of research groups (Van Raan, 2006), and some services now provide $h$-index

57  values associated with groups of keywords or phrases in a given bibliometric database (Thomson

58  Reuters, 2011).

59

60  Recently the $h$-indices of a number of human pathogens (n=27) were shown to be significantly

61  positively correlated with their impact as measured by disability adjusted life year (DALY)

62  estimates (McIntyre et al, 2011). DALYs provide a combined measure of the years of healthy

63  life lost as a result of poor health or disability in combination with an estimation of the potential

64  years of life lost due to premature death. They were developed by the World Health Organization

65  (Murray, 1994) and although they have only been estimated for a small number of diseases, they

66  have become the most widely-used measure of the true burden of disease (Mathers et al, 2004).

67  Thus, although the $h$-index is a measure that reflects global scientific interest in a pathogen (and

3

68    inherently reflects research trends and funding), it is a reasonable proxy indicator of high impact

69    human pathogens (McIntyre et al, 2011).

70

71    The aims of this work are to investigate:

72    1) whether the *h*-index might be used to identify the pathogens that are likely to have a high

73    impact upon human health in the North American region.

74    2) how the *h*-index might be used to apply a relative ranking to a set of pathogens identified as

75    emerging hazards. In this second aim we focus on examples of pathogens of interest to Canada

76    because our research institution and funding agency are based in Canada.

77    3) how the *h*-index of a pathogen changes over time.

78

79

80    **Materials and Methods**

81    *Identification of pathogen species*

82    The 'ENHanCEd Infectious Diseases' (EID2) database (University of Liverpool, 2011) provided

83    the raw data for this study. The purpose of this database is to provide a method of studying the

84    main pathogens and hosts involved in disease transmission (McIntyre et al, 2013). It contains

85    information about more than 740,000 organisms (such as vectors, hosts and pathogens) and their

86    structure in the phylogenetic tree. This includes details regarding all pathogens that are known to

87    infect humans and some known to infect domestic or companion animal species. Information

88    about pathogens is assigned using data-mining of meta-data and semi-automated literature

89    searches, for further details see McIntyre et al (2013).

90

4

91   Information was extracted about all organisms classified as 'pathogen species' that were known

92   to infect humans. Analysis only included human pathogen species, since the database did not

93   include all known North American animal host species. All data searches were undertaken in

94   October 2011.

95

96   Information about the taxonomic division and pathogenic status of each pathogen was extracted

97   from the database. Pathogenic status was defined as one of 'frequently pathogenic' (frequently

98   causes morbidity and/or mortality in the general population), 'non-pathogenic' (does not

99   frequently cause clinical signs within the general population, but may affect immune-

100   compromised individuals) or 'unknown' (there was insufficient evidence to determine

101   pathogenic effects).

102

103   Zoonotic potential and emerging status of the pathogens were taken from (Taylor et al, 2001;

104   Woolhouse and Gowtage-Sequeria, 2005) where available. Definitions were described in those

105   publications as follows. Zoonotic potential was classified as either 'non-zoonotic' or 'zoonotic'.

106   Zoonotic pathogens were those that are naturally transmitted between vertebrate animals and

107   humans. Pathogens previously but no longer transmitted from animals, such as HIV, were not

108   regarded as zoonotic. Pathogens were classified as 'non-emerging' or 'emerging'. Emerging

109   pathogens were those that have appeared in a human population for the first time, or have

110   occurred previously but are increasing in incidence or expanding into areas where they had not

111   previously been reported over the last 20 years. Once the $h$-index of pathogens had been

112   calculated (see below), the indices of pathogenic versus non-pathogenic, zoonotic versus

5

113    non-zoonotic and emerging versus non-emerging pathogens were compared using the

114    non-parametric Mann Whitney U test.

115

116    *Calculation of the h-index of pathogens*

117    The *h*-index scores were obtained for all pathogen species named in the database using WOS

118    (Thomson Reuters, 2011). The following specific search protocol was followed in order to

119    identify all scientific papers relating to each pathogen, and thus to calculate the *h*-index of the

120    pathogen.

121

122    For each pathogen, literature searches were undertaken using search phrases specified using

123    quotation marks ("”), the 'topic' search field and with lemmatization turned off. Search phrases

124    were compiled which included the scientific name and any alternative names, synonyms or

125    alternative spellings according to the National Center for Biotechnology Information (NCBI)

126    taxonomy website (National Center for Biotechnology Information, 2011); searches for

127    organisms contained 'exclusion terms' when necessary. Searches for viruses were more complex

128    because of the frequent existence of synonyms and acronyms. Synonyms and acronyms were

129    obtained from the NCBI taxonomy website (National Center for Biotechnology Information,

130    2011) or the International Committee on Taxonomy of Viruses (ICTV) website (International

131    Committee on Taxonomy of Viruses, 2011) and were included as additional search terms. Since

132    some acronyms were used for more than one virus, or occurred in a non-viral context, searches

133    also included the term 'virus' if they had 'virus' within their pathogen name or if they were

134    within the 'virus' division of the NCBI Taxonomy database and excluded any other entities (viral

135    or non-viral) which shared the acronym. The Boolean operators 'AND', 'OR' and 'NOT' were

6

136  used to link multiple search phrases. For example the query for *Sin nombre virus* contained the

137  following search terms: ('sin nombre virus' OR 'sin nombre hantavirus' OR ('snv' AND

138  'virus')) AND NOT ('spleen' OR 'sindbis').  All searches were restricted to the years from 1900

139  to 2011, inclusive. Search terms are available on request from the authors.

140

141  *Identification of pathogens that occur in the North American region*

142  Having calculated the *h*-index for all pathogens in the EID2 database, we firstly, identified which

143  pathogens are able to establish in the North American region, and secondly, ranked them

144  according to their *h*-index. Thus we use previous occurrence in the North American region as an

145  indicator of the pathogens that are more likely to emerge again in the same area, either because

146  they are endemic and have the potential to re-emerge or because, in the past, they have had the

147  opportunity to establish in the region. Clearly this is a simple indicator, however it provides a

148  method of identifying pathogens that are able to occur in a specific geographical region. Our

149  ranking of pathogens 'of interest' to Canada (see below) takes into account pathogens that are

150  exotic to Canada.

151

152  Two methods were used to identify which of the pathogens had been recorded in at least one of

153  the following North American countries: Canada, United States, Mexico or Greenland. These

154  countries (which we defined as the North American region and now refer to as 'North America')

155  were selected to comprise the North American land mass, while excluding the countries of

156  Central America for simplicity.

157

7

158    The first method involved searching for the pathogens within the NCBI Nucleotide database

159    (National Center for Biotechnology Information, 2011a). This database is a collection of genome

160    sequences from sequencing projects around the world. The metadata for nucleotide sequences in

161    some cases contains information about the location of pathogen isolation. In order to identify

162    location, searches established where the pathogen and at least one of the geographical 'Medical

163    Subject Headings'(MeSH) terms for Canada, United States, Mexico or Greenland co-occurred.

164    MeSH terms act as a controlled thesaurus and are used for indexing articles by the US Library of

165    Medicine (US National Library of Medicine, 2012a). If one sequence from a pathogen had been

166    recorded in North America within the Nucleotide database, then this was used as confirmation of

167    pathogen presence. A second method was also used to identify pathogen location because the

168    NCBI nucleotide database did not include location information about all pathogens in our study.

169    This second method used the PubMed database, a database that contains more than 21 million

170    citations of biomedical and life sciences literature (US National Library of Medicine, 2012b).

171    The database was searched for all publications where the pathogen search terms (described

172    above) and at least one of the geographical MeSH terms for Canada, United States, Mexico or

173    Greenland co-occurred. The search terms had to be recorded in the title or abstract of the

174    publication. There was a degree of inaccuracy associated with this method, since co-occurrence

175    of a pathogen and a North American search term does not necessarily indicate that the pathogen

176    has occurred in that region. Co-occurrence could also arise in publications that describe pathogen

177    absence, animal models or simulation models for example. In order to account for this

178    inaccuracy, only searches for pathogens which generated at least five references in the same

179    country were used as confirmation of pathogen presence in North America. The threshold of five

180    was chosen following sensitivity testing of the results from searches conducted for 21 randomly

181    selected pathogens. In brief, this involved stratifying the scientific publications according to the

182    pathogen and the continent to which they were linked via a MeSH term for a country. The

183    association was checked to substantiate that the pathogen was found in hosts (including humans)

184    within a MeSH term country. This indicated that on average 95% of the associations in single

185    papers were accurate. Therefore setting the threshold at five papers would result in a positive

186    predictive value PPV (i.e. proportion of predicted interactions for which papers provide

187    supporting evidence) exceeding 99.9%. A high enough threshold to avoid false positives was

188    balanced with the need to avoid causing any major bias against 'newer' pathogens that have

189    fewer publications. For detailed description see McIntyre et al, (2013).

190

191    *Comparison of the h-index calculated from WOS with PubMed*

192    The *h*-index was calculated from the WOS, which differs in its bibliographic content to other

193    bibliographic databases. In order to compare the output from WOS with the PubMed database,

194    non-parametric Spearman rank was used to correlate the WOS *h*-index of pathogens that

195    occurred in North America with the number of publications for that pathogen in PubMed.

196

197    *Ranking 'pathogens of concern' in Canada.*

198    Additional descriptive analysis focused on 'pathogens of interest' in Canada. These were

199    identified from three different sources. The first source was a recent publication which

200    highlighted pathogens that are likely to emerge in North America in response to climate change

201    (Greer et al, 2008). The second source was researchers from the Zoonotics Division of the Public

202    Health Agency of Canada (PHAC) (N. Ogden *pers comm*.). PHAC provided funding support to

203    the project and the researchers provided a list of pathogens that are of interest due to their

204  potential to become emerging hazards in Canada. Details about the characteristics of these

205  pathogens, including whether they have occurred in North America, were collated and they were

206  ranked according to their *h*-index. The third source was the Ontario burden of infectious disease

207  study (Kwong et al, 2010), a study that describes the mortality and morbidity of infectious

208  diseases in Ontario. It lists three measures of disease burden for infectious diseases that have

209  occurred in Ontario. The measures are YLL: Years of Life Lost due to premature mortality,

210  YERF: Year-Equivalents of Reduced Function as a result of disease or condition, and HALY:

211  Health Adjusted Life Years, which are calculated by adding the YLL and YERF for each

212  pathogen. Spearman rank correlation was used to compare each of these measures with the *h*-

213  index of the pathogen.  Kwong et al, (2010) calculated the burden for a total of 69 diseases,

214  however we only included those for which the pathogen was specified and could therefore be

215  matched with an *h*-index. Thus, general terms describing a disease or condition such as

216  'Septicaemia' were excluded from this analysis.

217

218  *Calculation of change in h-index over time*

219  Time-bounded *h*-index scores were obtained for a selected set of pathogens using the same

220  phrase searches as described above. However, here the cumulative *h*-index was calculated every

221  year from 1960 to 2011 inclusive to assess how the index changed over time. The pathogens

222  chosen were *Chikungunya virus, Hendra virus, Monkeypox virus, Nipah virus, Rift Valley Fever*

223  *virus, Trypanasoma cruzi* (the cause of Chagas disease) and *West Nile Virus.* These pathogens

224  were chosen as example pathogens that have been classified as either 'emerging' or 'non-

225  emerging'; with the intent to compare the change in *h*-index of both types of pathogens.

226  Furthermore, they are examples of pathogens that were deemed to be of particular interest to the

10

227     PHAC (N. Ogden *pers comm*.) due to their potential for emergence in Canada. In addition we

228     calculated the cumulative *h*-index for *Plasmodium falciparum*, because it is a pathogen of

229     worldwide concern and because preliminary calculations showed that it has one of the highest *h*-

230     indices.

231

232     In order to assess the rate of change in *h*-index for these pathogens, two negative binomial

233     models were evaluated. The first model assessed the cumulative *h*-index as the outcome where

234     the time since first publication was used as an offset. Since the rate of change will be largely

235     influenced by the number of publications for pathogens of major importance (e.g. *P. falciparum*),

236     the second model assessed the rate of change of *h*-index by year. The outcome for this model

237     was calculated by subtracting the cumulative *h*-index for a particular year from the previous

238     year, except for the first year of the series. Similarly, the number of years since first publication

239     was used as an offset. Both models included a categorical variable indicating the pathogen, a

240     variable indicating the calendar year when *h*-index was computed and the interaction between

241     these two variables. These models were assessed using the Deviance and X2 goodness of fit tests

242     (Dohoo et al, 2009). The predicted rates from these models were calculated and plotted against

243     time for each pathogen.

244

245     **Results**

246     *Pathogens likely to have a high impact in North America*

247     A total of 3627 pathogen species were recorded in EID2 and of these 1827 were classified as

248     human pathogens species. Of these, 651 were human pathogens that have been recorded in North

249     America. These consisted of 474 pathogen species that have occurred in North America

11

250     according to the Nucleotide database, and an additional 177 pathogens that were identified when

251     the pathogen search terms occurred in at least five publications in conjunction with the North

252     American search terms entered into the PubMed database. A total of 258 occurred in both the

253     Nucleotide database and the PubMed searches.

254

255     The *h*-index of the human pathogen species ranged from 0 to 584 and was highly over dispersed

256     (Figure 1). Only a limited number of pathogens had an *h*-index over 100, with most pathogens

257     scoring a relatively low value (median=37). Although, the *h*-index was calculated from WOS,

258     which differs to some extent in its bibliographic content from PubMed, the *h*-index was

259     significantly correlated with the number of publications recorded in the PubMed database

260     (Spearman rank correlation $r_s$=0.736, p<0.001, n=651), (Figure 2).

261

262     The largest proportion (42.2%) of pathogen species were bacteria, followed by fungi (21.2%)

263     and viruses/prions (16.0%) (Table 1). The 10 pathogens with the highest *h*-index included one

264     yeast (*Saccharomyces cerevisiae*), five viruses and four bacteria species (Table 2).

265

266     Information about emergence status (emerging or non-emerging) and zoonotic potential

267     (zoonotic or non-zoonotic) was obtained from two publications. These publications assigned an

268     emergence status to 462 (71%) and a zoonotic status to 464 (71%) of the 651 pathogens included

269     in our analysis. Of the 462 pathogens that had an emergence status, 26.2% were classified as

270     emerging. Pathogen species with the highest *h*-index recorded in WOS that were classified as

271     emerging included *Escherichia coli*, *Human immunodeficiency virus 1* and *2* and *Hepatitis C*

272     *virus* (Table 2). Emerging pathogens had a significantly higher *h*-index than non-emerging

12

273    pathogens (Mann Whitney U, p<0.001) (Table 3). A total of 464 of the pathogens had been

274    assigned a zoonotic potential status and 67.9% of these were zoonotic (Table 3). The *h*-index

275    values of zoonotic and non-zoonotic pathogens were not significantly different (Mann Whitney

276    U, p = 0.718). Pathogens that were frequently pathogenic had significantly higher *h*-index scores

277    than pathogens that were non-pathogenic (Mann Whitney U, p<0.001) (Table 3). There were 13

278    pathogens of 'unknown' pathogenicity, which were excluded from this analysis.

279

280    *Using the h-index to apply a relative ranking to pathogens of interest*

281    Additional analysis focused on pathogens that had been identified as potential emerging hazards

282    within Canada in the literature or by PHAC. These pathogens were both endemic and exotic to

283    Canada. Of the pathogens of interest to PHAC, two (*Plasmodium falciparum*) and Verotoxic *E.*

284    *coli*) cause notifiable diseases in Canada (Public Health Agency of Canada, 2010) and three

285    (*Nipah virus*, *Hendra virus* and *Rift Valley Fever virus*) had not previously been recorded in

286    North America (and therefore did not feature in our list of North American pathogens).

287

288    All of the pathogens of interest from both sources were classed as frequently pathogenic. Of the

289    pathogens that were highlighted by Greer *et al.* (2008), those with the highest *h*-index were *E.*

290    *coli*, *P. falciparum* and *Streptococcus pneumoniae* (Table 4). All had previously been recorded in

291    North America. Additional pathogens of concern to PHAC with the highest *h*-indices were

292    *Trypanosoma cruzi* (the cause of Chagas Disease), *Nipah Virus* and *Hendra Virus* (Table 5).

293    Only *T. cruzi* has been recorded in North America and had a much higher *h*-index (130) than any

294    of the others deemed to be of interest by PHAC. Overall, the median *h*-index of the pathogens

295    listed in Tables 4 and 5 is 82, which is considerably greater than the median value of 37 for all

13

296 North American pathogens analysed. The *h*-index of 31 of these 33 pathogens were ranked in the

297 top 50% of the North American pathogens (figure 1). The only exceptions were the food and

298 water borne pathogens *Cryptosporidium hominis* and *Shigella boydii.*

299

300 The *h*-index of pathogens was positively correlated with the HALY measure of pathogen impact

301 in Ontario (Spearman rank correlation $r_s$=0.627, p<0.001, n=41), (Figure 3). The *h*-index was

302 also positively correlated with the two measures that make up the HALY score, namely the YLL

303 ($r_s$=0.676, p<0.001, n=41) and the YERF ($r_s$=0.448, p<0.003, n=41). Of the 20 pathogens with

304 the highest HALY score in Ontario, a total of 8 also feature in the top 20 pathogens with the

305 highest *h*-index, while 15 feature in the top 50 and 16 have an *h*-index of greater than 100. The

306 strength of this correlation is likely influenced by a few very high impact pathogens and we

307 highlight that there are also a few pathogens that have a relatively high *h*-index score, although a

308 relatively low HALY measure, e.g. malaria.

309

310 *Change in h-index over time*

311 The *h*-index of pathogens varied considerably over time. Figure 4 shows the time series for seven

312 pathogens; *P. falciparum* was excluded from this figure because it has a high *h*-index that tends

313 to obscure the series of the other pathogens. The *h*-index of *Rift Valley fever virus* and

314 *Monkeypox virus* increased gradually from 1960 onwards. This was also the case for *T. cruzi,*

315 although the *h*-index value for this pathogen was much greater than for any of the other six. The

316 *h*-index of *Chikungunya virus* increased gradually from 1960 onwards, showing a steep increase

317 in 2005 until 2007. *West Nile virus* showed a steady increase in *h*-index from 1960, until around

318 1998 at which point it was associated with a sharp increase which only recently appeared to have

14

319 reached a plateau. The *h*-index scores for *Hendra virus* and *Nipah virus* were zero until the mid-

320 1990s, but then increased relatively rapidly until around 2005 when both appear to have

321 plateaued to some extent.

322

323 When the *h*-index was adjusted for the number of years since the first record of the pathogen (or

324 in the case of 'older' pathogens the record in 1960, when our dataset began), the pathogens with

325 the highest *h*-index were *P. falciparum* and *T. cruzi* (Figure 5a) throughout this time. In the

326 2000s the *h*-index of the pathogens *Hendra virus* and *Nipah virus* increased more rapidly than

327 the other pathogens tested. When the yearly rate of change of the *h*-index was measured

328 (adjusted by discounting the *h*-index from previous years), the *h*-index of *P. falciparum* and of

329 *West Nile virus* increased at a higher rate than any of the other pathogens (Figure 5b). In

330 comparison the rate of change of the *h*-index for *T. cruzi* gradually decreased from 1960. Finally,

331 both *Hendra virus* and *Nipah virus* showed a rapid increase in the 1990s, although in more recent

332 years (since the early 2000s) the rate of change of the *h*-index of these viruses has decreased.

333

334 **Discussion**

335 *Pathogens likely to have a high impact in North America*

336 The *h*-index of a pathogen can be viewed as an indicator of the relative scientific interest in that

337 pathogen. Although it likely reflects trends in research interest, research funding and regional

338 bias, the *h*-index of a limited number of pathogens has been correlated with their DALY measure

339 which suggests that it might be used as a measure of impact (McIntyre et al, 2011). We focused

340 on human pathogens that have been recorded in North America. We used previous occurrence in

341 North America as an indicator of the pathogens that are more likely to emerge again in the same

15

342　area, because geographic proximity is a characteristic that has been deemed a risk for emergence,

343　for example in Canada (Cox et al, 2012). Clearly, this is a simple indicator and other non-

344　endemic pathogens that have not previously been recorded in the region could still emerge.

345

346　The species with the highest *h*-index values included yeast (*Saccharomyces cerevisiae*), which

347　can cause opportunistic infection, food-borne pathogens (*E. coli*), person to person transmitted

348　viruses (*Hepatitis* B and C virus, *Human Immunodeficiency virus*, human herpesvirus), bacteria

349　that cause multiple clinical infections (*Staphylococcus aureus*) and person to person transmitted

350　bacteria (*Helicobacter pylori*).  While some of these pathogens have a high impact on the human

351　population, others are likely to have generated a high *h*-index for other reasons. For example, the

352　vast majority of publications about *S. cerevisiae* are related to its industrial use in brewing and

353　baking, rather than opportunistic infection. Similarly, the high *h*-index of *S. aureus* is likely to be

354　associated with non-zoonotic infections in multiple species, rather than simply human illness.

355

356　There are two implications of these findings. Firstly that there may be a need to refine our search

357　terms, as we increase our understanding of the biases of the *h*-index. Secondly, that the *h*-index

358　may be most useful for ranking selected pathogens of concern. We suggest, therefore that it

359　might be most reliably used as one complementary component of a pathogen prioritisation risk

360　assessment particularly since such studies often rely on qualitative data or expert opinion (Cox et

361　al, 2012, Krause and Working Group on Prioritization at Robert Koch Institute, 2008). Indeed

362　the first publication on the *h*-index notes (when assessing the *h*-index of the individual

363　researcher) that, 'a single number can never give more than a rough approximation to an

16

364 individual's multifaceted profile, and many other factors should be considered in combination in

365 evaluating an individual' (Hirsch, 2005).

366

367 *Using the h-index to apply a relative ranking to pathogens of interest*

368 An additional part of our work focused on pathogens that have been identified as 'pathogens of

369 interest' in Canada. All of the pathogens that were identified by PHAC or by Greer et al, (2008)

370 had a relatively high ranking *h*-index. Those with the highest *h*-index were *E. coli, P. falciparum*

371 and *S. pneumonia*. These pathogens are likely to have a high *h*-index either because they tend to

372 be virulent and/or because they spread relatively easily in the human population, either via

373 vectors, food and water or from person to person. It is important to note that our analysis only

374 included species level pathogens and that we did not differentiate between strains of pathogens.

375 This may be a useful distinction in future analyses. *E coli*, for example, are a large and diverse

376 group of bacteria, which includes both virulent and non-virulent strains as well as zoonotic and

377 non-zoonotic strains. In our analysis, *E. coli* has been classified as zoonotic, because at least

378 some strains are zoonotic. This group is likely to score a high *h*-index not only due to the impact

379 of virulent strains such as zoonotic Verotoxic *E. coli* O157, but also due to the prevalence of

380 non-zoonotic illness such as urinary tract infections and neonatal meningitis.

381

382 Within our list of pathogens 'of interest' there are some that score a relatively high *h*-index but

383 that do not cause especially severe disease. Examples, with an *h*-index greater than 100, include

384 *Salmonella enterica*, *Respiratory synctial virus* and influenza virus. These tend to cause mild

385 symptoms in the general population, (although they can be severe in individuals who are

386 immunocompromised). They are likely to have generated scientific interest (and therefore a high

17

387   *h*-index) due to their morbidity and their ease of transmission. *Salmonella enterica* for example,

388   which has an *h*-index of 107, causes a diarrheal infection and occurs worldwide. In Canada there

389   are an estimated 6,000 to 12,000 cases per year (Health Canada, 2012), although it is likely that

390   cases are under-reported and that the actual number of infections is much higher.

391

392   A positive correlation between the *h*-index and the HALY score indicated that the *h*-index is a

393   proxy for this measure of pathogen impact in Ontario and it could therefore be used to rank

394   pathogens that are known to occur in a specific region. While we found a positive relationship

395   between the two measurements, we also show that the ranking needs to be interpreted in the

396   correct context. For example, *P. falciparum* scores a high *h*-index, but a relatively low HALY in

397   Ontario. This shows that it has a high impact on a global scale, but that its impact within the

398   cooler climate of Ontario is relatively low because it does not commonly occur.

399

400   Overall, our analysis of pathogens 'of interest' from three different sources, supports the idea

401   that the *h*-index could be a practical method to compare potential pathogen hazards. There are

402   two particular ways that it could be best used. Firstly as a method to separate high and low

403   priority pathogens and therefore act as a rapid screening method for pathogens that require

404   further risk analysis. Secondly, to rank pathogens that are 'of interest' in a specific region. For

405   example to rank pathogens that are exotic to a region, but are of concern due to their global

406   impact, or to rank pathogens that are endemic in a region and that occur frequently enough to

407   have become 'of interest'.

408

409   *Change in h-index over time*

18

410    Analysis of time series data demonstrated that the *h*-index of a pathogen changes over time, even

411    after accounting for the increasing trend in total number of publications. We hypothesise that the

412    rate of change of the *h*-index might be used as a crude indicator of a pathogen's emergence

413    and/or the spread of infection. *Hendra virus,* for example, was discovered in horses in Australia

414    in 1994 and its *h*-index began increasing from 0 in 1995. Similarly the *h*-index of *Nipah virus*

415    increased from the time that it was discovered in a pig population in Malaysia in 1999. The

416    *h*-index of both of these recently emerging pathogens has increased rapidly since their

417    identification compared to the other pathogens studied here. It is also notable that the *h*-index of

418    *West Nile virus*, which increased steadily from 1960 showed a relatively rapid increase from

419    1999 onwards and we hypothesise that this increase coincides with the emergence of the disease

420    in the Eastern United States in 1999 (Soverow et al, 2009). Finally, the increase in the *h*-index of

421    Chikungunya virus from 2005 to 2007 coincides with its outbreak in the Indian Ocean territories

422    in 2005 (Schuffenecker et al, 2006).

423

424    There is likely to be a bias in the *h*-index towards 'old' pathogens compared to newly emerging

425    pathogens, for which papers have not yet had time to accumulate citations. Indeed, it has been

426    suggested that the *h*-index can only provide a realistic assessment of the achievement of

427    academics (and therefore in our work – the impact of a pathogen) who have been publishing for

428    at least ten years (Harzing, 2008). One way to compare between pathogens with different lengths

429    of 'academic publishing' is to divide the *h*-index by the number of years since the first

430    publication, a measure referred to as the 'm-quotient' (Hirsch, 2005). Our analysis, which

431    controlled for the number of years of publication, revealed how the rate of change of *P.*

432    *falciparum* and *West Nile virus* was higher than the other pathogens tested. The high rate of

19

433    increase in the *h*-index of the malaria pathogen reflects the impact of the disease for which there

434    were approximately 219 million cases worldwide in 2010 (World Health Organisation, 2012).

435    Although the impact in terms of mortality rates has fallen by 26% since 2000, the increasing

436    *h*-index also accounts for the fact that malaria is a risk to over half of the world's population, and

437    that international disbursements and government funding for malaria control rose steeply during

438    this time (World Health Organisation, 2012). We suggest that the *h*-index of *West Nile virus* has

439    increased at a high rate because this reflects the impact of the pathogen as it has spread across the

440    USA since it emerged in 1999 and because its emergence has been attributed to climate warming

441    (Soverow et al, 2009). In contrast, the rate of change in *h*-indices of other pathogens such as *T.*

442    *cruzi* and *Chikingunya virus* have decreased yearly. These pathogens have both been described

443    as 'neglected tropical diseases', which tend to be endemic in low income, developing regions and

444    typically have a high morbidity, but low mortality (Hotez, 2011). The rate of change in the *h*-

445    indices of newly emerging pathogens (*Nipah virus* and *Hendra virus*) showed a different pattern

446    to that of the 'older' pathogens, with a rapid increase following their emergence which then

447    slowed in more recent years. This trajectory is likely to reflect the increasing scientific interest in

448    a newly emerging pathogen, which then levels off as knowledge is established.

449

450    Identifying patterns of disease emergence using bibliometric measures or electronic resources

451    has proven a valuable tool to augment disease monitoring and surveillance. For example,

452    patterns of disease reporting in ProMED (the Internet-based 'Program for Monitoring Emerging

453    Diseases' (International Society for Infectious Diseases, 2013)) have been used as an early-

454    warning of disease emergence (Cowen et al, 2006), while records of Internet queries have been

455    used to track the spread of influenza infections (Google Flu Trends, 2012; Ginsberg et al, 2009)

456 and Methicillin Resistant *S. aureus* (Dukic et al, 2011). Similarly, social networking tools such

457 as Twitter, have proven to be real-time indicators of public health concerns, since the number of

458 Twitter posts relating to 'swine flu' and/or 'H1N1' in 2009 correlated well with H1N1 incidence

459 data (Chew and Eysenbach, 2010). Twitter has also been used to measure the uptake of research

460 findings, with the number of tweets generated within the first 3 days of an article's publication

461 being a good predictor of highly cited articles. A proposed 'twimpact factor' has therefore been

462 suggested as a timely metric to gauge research impact and influence (Eysenbach, 2011). A

463 comparison of 'twimpact factor' with *h*-index may provide some predictive value in the case of

464 disease monitoring.

465

466 In comparison to the *h*-index, the indicators described above are more instantaneous measures

467 and it is unlikely that the change in *h*-index could be used for real-time surveillance purposes due

468 to the time lag in the measure of the *h*-index and the relative impact. In addition, newly emerging

469 pathogens are likely to be under-represented. However, the trajectory of the *h*-index may be

470 relatively predictable if combined with other measures. Work has shown that it is possible to

471 predict the future *h*-index of scientists as far as five to ten years into the future, on the basis of

472 additional publicly available information, including years since publishing their first article,

473 number of distinct journals published in and the number of articles in five prestigious journals

474 (Acuna et al, 2012).

475

476 Assessment of a wider range of pathogens would be beneficial, with a particular focus on

477 emerging pathogens. Specific incidences of emerging diseases and global emerging disease

478 hotspots have been identified in the past (Jones et al, 2008). Comparison of the *h*-index of

21

479 pathogens with their global emergence may reveal the typical time delay between disease

480 emergence and changes in associated *h*-indices, as well as whether there is a level of increase in

481 *h*-index that can be reliably interpreted as an early warning of future disease emergence.

482

483 *Comparison of the h-index with other bibliometric sources*

484 The *h*-index scores in the present study were generated from one bibliometric source and

485 comparison was not made with other sources. Other bibliometric services, such as SCOPUS or

486 Google Scholar, search different literature sources over differing temporal periods. Although

487 alternative sources produce slightly different *h*-index values, these tend to be comparable across

488 platforms (McIntyre et al, 2011). Our work demonstrated a clear correlation between the *h*-index

489 calculated in WOS and the number of publications recorded in PubMed.

490

491 Overall the *h*-index combines an assessment of both the quantity of publications and the quantity

492 of citations. A pathogen cannot have a high *h*-index without having a substantial number of

493 papers published about it. However the number of papers is not enough – a reasonable number of

494 these papers need to have been cited in order to increase the *h*-index value. The *h*-index thus

495 corrects for pathogens that might have a limited number of highly cited papers, or many that

496 have not been cited. It therefore tends to highlight pathogens that generate a continuous stream of

497 publications with above average publication impact. While the *h*-index is the most commonly

498 cited metric, alternative methods of assessing research output have been suggested (Harzing,

499 2007; Alonso et al, 2009) and might be considered in future assessments of pathogen impact. For

500 example, the *g*-index could be used to give more weight to highly-cited articles (Egghe, 2006)

501 and has been suggested as a useful complement to the *h*-index (Harzing, 2008). We also suggest

502    an evaluation of the measure '*cf*', which takes into account the differences in number of citations

503    received by all articles in a given year, so that scientific impact can be compared across years

504    (Radicchi et al, 2008).

505

506    In conclusion, the *h*-index is a quantitative measure that can be used to estimate the potential

507    impact of a pathogen and that can be calculated quickly and easily. It can be used to identify and

508    to rank individual pathogens or types of pathogens (e.g. zoonotic, emerging and pathogenic) and

509    to measure changes over time. It could provide a rapid method of screening for pathogens that

510    are likely to be important and therefore it would be particularly useful if incorporated into a

511    prioritisation tool to complement a set of more qualitative criteria.

512

520

521    **Author contributions**

522    Conceived and designed the study: RC, KMM, JS, CS, MB, CWR

523    Provided the *h-indices* for analysis: KMM, CS, MB

524    Performed the data collection: RC, KMM, CS

23

525    Analyzed the data: RC, KMM, JS, CS, CWR

526    Wrote and commented on the manuscript: RC, KMM, JS, CS, MB, CWR

527
528    **Figure captions**

529
530    Figure 1 A scatter-plot showing the *h*-index value (y axis) of 651 pathogen species that are
531    infectious to humans and have been recorded in the North American region against the rank
532    position of each of those pathogens (x axis).
533    Points that have been coloured black indicate pathogens that were identified as potential
534    emerging hazards and therefore of interest in Canada.

535
536    Figure 2 Correlation of *h*-index with number of publications reported in the PubMed database for
537    651 human pathogen species.
538    (The x-axis has been truncated at 300 to better demonstrate the association of the *h*-index and the
539    total number of publications. There were only four pathogens with an *h*-index greater than 300.
540    These were *Saccharomyces cerevisiae*, *Escherichia coli*, *Human Immunodeficiency virus 1* and
541    *2*).

542
543    Figure 3 Correlation of *h*-index with Health Adjusted Life Year measurement of 36 infectious
544    diseases that occurred in Ontario in 2010.
545    Data labels show the pathogen or disease named in the study by (Kwong et al, 2010).

546
547    Figure 4 The *h*-index score by year from 1960 to 2009 for seven selected pathogens.

548
549    Figure 5 The modelled rate of change of the *h*-index from 1960 to 2009 for eight pathogens.
550    A. The *h*-index has been adjusted according to the number of years since the first record.
551    B. The yearly rate of change of the *h*-index. This model has been adjusted according to the
552    number of years since the first record and it also discounts the *h*-index from previous years.

553
554    **Table captions**

555
556    Table 1 Taxonomic classification of 651 human pathogen species that have been recorded in the
557    North American region.

558
559    Table 2 Pathogen species with the highest *h*-index recorded in Web of Science from those human
560    pathogen species recorded in the North American region. All were classed as frequently
561    pathogenic.

562
563    Table 3 Summary of *h*-index values for human pathogen species that have been recorded in the
564    North American region, grouped according to emerging, zoonotic and pathogenic status
565    A total of 651 human pathogen species were recorded in the North American region, however
566    not all had been assigned a status for each characteristic.

567

24

568     Table 4 The *h*-index of 33 pathogens that have been identified as an emergence risk in Canada
569     by Greer et al (2008).
570     Pathogen names in grey indicate a pathogen that can cause the associated disease, but that is not
571     commonly the main cause of the disease in the North American region.
572
573     Table 5 The *h*-index of six pathogens that have been identified as pathogens of emergence
574     concern in Canada by the Public Health Agency of Canada (N. Ogden *pers comm*.).
575
576

577     **References**

578     Acuna, D. E., S. Allesina, and K. P. Kording, 2012: Future impact: Predicting scientific success.
579     Nature 489, 201-202.

580     Alonso, S., F. J. Cabrerizo, E. Herrera-Viedma, and F. Herrera, 2009: h-index: A review focused
581     on its variants, computation and standardization for different scientific fields. J. Informetr. 3,
582     273-289. Chew, C., and G. Eysenbach, 2010: Pandemics in the age of Twitter: content analysis
583     of Tweets during the 2009 H1N1 outbreak. PLoS One 5, e14118.

584     Cowen, P., T. Garland, M. E. Hugh-Jones, A. Shimshony, S. Handysides, D. Kaye, L. C.
585     Madoff, M. P. Pollack, and J. Woodall, 2006: Evaluation of ProMED-mail as an electronic early
586     warning system for emerging animal diseases: 1996 to 2004. J. Am. Vet. Med. Assoc. 229, 1090-
587     1099.

588     Cox, R., C. W. Revie, and J. S. Sanchez, 2012: The use of expert opinion to assess the risk of
589     emergence or re-emergence of infectious diseases in Canada associated with climate change.
590     PLoS ONE 7: e41590.

591     Dohoo, I. R., W. Martin, and H. Stryhn, 2009: Veterinary Epidemiological Research, Veterinary
592     Epidemiological Research, 2nd edn , pp. 865-866. VER Inc.

593     Dukic, V., M. David, and D. S. Lauderdale, 2011: Internet Queries and Methicillin-Resistant
594     *Staphylococcus aureus* Surveillance. Emerg. Infect. Diseases 17, 1068-1070.

595     Egghe, L., 2006: Theory and practice of the g-index. Scientometrics 69, 131-152.

596     Eysenbach, G., 2011: Can tweets predict citations? Metrics of social impact based on Twitter and
597     correlation with traditional metrics of scientific impact. J. Med. Internet. Res. 13, e123.

598     Ginsberg, J., M. H. Mohebbi, R. S. Patel, L. Brammer, M. S. Smolinski, and L. Brilliant, 2009:
599     Detecting influenza epidemics using search engine query data. Nature 457, 1012-1014.

600     Google, 2013: Google Scholar. Available at: http://scholar.google.ca/ (Accessed May 2012).

601     Google Flu Trends, 2012: Google flu trends around the world. Available at:
602     http://www.google.org/flutrends/ (Accessed May 2012).

25

603     Greer, A., V. Ng, and D. Fisman, 2008: Climate change and infectious diseases in North
604     America: the road ahead. CMAJ 178, 715-722.

605     Harzing, A., 2008: Reflections on the *h*-index. Available at:
606     http://www.harzing.com/pop_hindex.htm (Accessed May 2012).

607     Harzing, A., 2007: Publish or perish. Available at: http://www.harzing.com/pop.htm. (Accessed
608     May 2012)

609     Health Canada, 2012: Salmonella prevention. Available at: http://www.hc-sc.gc.ca/hl-vs/iyh-
610     vsv/food-aliment/salmonella-eng.php. (Accessed Dec 2012).

611     Hirsch, J. E., 2005: An index to quantify an individual's scientific research output. Proc. Natl.
612     Acad. Sci. USA. 102, 16569-16572.

613     Hotez, P. J., 2011: The neglected tropical diseases and the neglected infections of poverty:
614     Overview of their common features, global disease burden and distribution, new control tools,
615     and prospects for disease elimination , National Academies Press, Washington DC.

616     International Committee on Taxonomy of Viruses, 2011: Taxonomy of Viruses. Available at:
617     http://www.ictvdb.org/Ictv/index.htm (Accessed October 2011).

618     International Society for Infectious Diseases, 2013: ProMED. Available at:
619     http://www.promedmail.org/.

620     Jones, K. E., N. G. Patel, M. A. Levy, A. Storeygard, D. Balk, J. L. Gittleman, and P. Daszak,
621     2008: Global trends in emerging infectious diseases. Nature 451, 990-993.

622     Krause, G., and Working Group on Prioritization at Robert Koch Institute, 2008: How can
623     infectious diseases be prioritized in public health? A standardized prioritization scheme for
624     discussion. EMBO Rep. 9 Suppl. 1, S22-7.

625     Kwong, J.C., N. S. Crowcroft, M. A. Campitelli, S. Ratnasingham, N. Daneman, S. L. Deeks, D.
626     G. Manuel, and Ontario Burden of Infectious Disease Study Advisory Group 2010:  Ontario
627     Burden of Infectious Disease Study (ONBOIDS). An OAHPP/ICES Report. Ontario Agency for
628     Health Protection and Promotion, Institute for Clinical Evaluative Sciences. Available at:
629     http://www.publichealthontario.ca/

630     Mathers, C. D., D. Ma Fat, J. T. Boerma and World Health Organization, 2004: The global
631     burden of disease: 2004 update. Available at:
632     http://www.who.int/healthinfo/global_burden_disease/2004_report_update/en/.

633     McIntyre, K. M., I. Hawkes, A. Waret-Szkuta, S. Morand, and M. Baylis, 2011: The H-index as
634     a quantitative indicator of the relative impact of human diseases. PLoS ONE 6: e19558.

635    McIntyre, K. M., C. Setzkorn, M. Wardeh, P. J. Hepworth, A. D. Radford, and M. Baylis, 2013:
636    Using open-access comprehensive database for the study of Mammalian and avian livestosk and
637    pet infections. Prev. Vet. Med. http://dx.doi.org/10.1016/j.prevetmed.2013.07.002

638    Murray, C.J.L, 1994: Quantifying the burden of disease: the technical basis for disability
639    adjusted life years. Bulletin of the World Health Organisation 72, 429-445.

640    National Center for Biotechnology Information, 2011a: Nucleotide database. Available at:
641    http://www.ncbi.nlm.nih.gov/nuccore (Accessed October 2011).

642    National Center for Biotechnology Information, 2011b: Taxonomy browser. Available at:
643    http://www.ncbi.nlm.nih.gov/Taxonomy/taxonomyhome.html (Accessed October 2011).

644    Public Health Agency of Canada, 2010: National notifiable diseases: notifiable diseases on line.
645    Available at: http://dsol-smed.phac-aspc.gc.ca/dsol-smed/ndis/list-eng.php (Accessed May
646    2010).

647    Radicchi, F., S. Fortunato, and C. Castellano, 2008: Universality of citation distributions: toward
648    an objective measure of scientific impact. PNAS 105, 17268-17272.

649    Schuffenecker, I., I. Iteman, A. Michault, S. Murri, L. Frangeul, M. C. Vaney, R. Lavenir, N.
650    Pardigon, J. M. Reynes, F. Pettinelli, L. Biscornet, L. Diancourt, S. Michel, S. Duquerroy, G.
651    Guigon, M. P. Frenkiel, A. C. Bréhin, N. Cubito, P. Després, F. Kunst, F. A. Rey, H. Zeller, S.
652    Brisse, 2006: Microevolution of Chikingunya viruses causing the Indian Ocean outbreak. PLoS.
653    Med. 3(7): e263.

654    Soverow, J. E., G. A. Wellenius, D. N. Fisman, and M. A. Mittleman, 2009: Infectious disease in
655    a warming world: how weather influenced West Nile virus in the United States (2001-2005).
656    Environ. Health. Perspect. 117, 1049-1052.

657    Taylor, L. H., S. M. Latham, and M. E. Woolhouse, 2001: Risk factors for human disease
658    emergence. Philos. Trans. R. Soc. Lond. B. Biol. Sci. 356, 983-989.

659    Thomson Reuters, 2011: Overview - Web of Science. Available at:
660    http://thomsonreuters.com/products_services/science/science_products/a-
661    z/web_of_science/(Accessed October 2011).

662    University of Liverpool, 2011: ENHanCEd Infectious Diseases database (EID2). Available at:
663    www.zoonosis.ac.uk/eid2 (Accessed October 2011).

664    US National Library of Medicine, 2012a: MeSH. Available at:
665    http://www.ncbi.nlm.nih.gov/mesh (Accessed October 2011).

666    US National Library of Medicine, 2012b: PubMed. Available at:
667    http://www.ncbi.nlm.nih.gov/sites/entrez?db=pubmed (Accessed October 2011).

668 Van Raan, A., 2006: Comparison of the Hirsch-index with standard bibliometric indicators and
669 with peer judgment for 147 chemistry research groups. Scientometrics 67, 491-502.

670 Woolhouse, M. E., and S. Gowtage-Sequeria, 2005: Host range and emerging and reemerging
671 pathogens. Emerg. Infect. Dis. 11, 1842-1847.

672 World Health Organisation, 2012: World Malaria Report 2012, World Health Organisation,
673 Geneva.

674
675