



# H-index manipulation by merging articles: Models, theory, and experiments ☆



René van Bevern<sup>a,b,d,\*</sup>, Christian Komusiewicz<sup>c,d</sup>, Rolf Niedermeier<sup>d</sup>,  
Manuel Sorge<sup>d</sup>, Toby Walsh<sup>d,e,f</sup>

<sup>a</sup> Novosibirsk State University, Novosibirsk, Russian Federation

<sup>b</sup> Sobolev Institute of Mathematics, Siberian Branch of the Russian Academy of Sciences, Novosibirsk, Russian Federation

<sup>c</sup> Institut für Informatik, Friedrich-Schiller-Universität Jena, Germany

<sup>d</sup> Institut für Softwaretechnik und Theoretische Informatik, TU Berlin, Germany

<sup>e</sup> University of New South Wales, Sydney, Australia

<sup>f</sup> Data61, Sydney, Australia

## ARTICLE INFO

### Article history:

Received 9 March 2016

Received in revised form 26 July 2016

Accepted 5 August 2016

Available online 10 August 2016

### Keywords:

Citation index

Hirsch index

Parameterized complexity

Exact algorithms

AI's 10 to watch

## ABSTRACT

An author's profile on Google Scholar consists of indexed articles and associated data, such as the number of citations and the H-index. The author is allowed to merge articles; this may affect the H-index. We analyze the (parameterized) computational complexity of maximizing the H-index using article merges. Herein, to model realistic manipulation scenarios, we define a compatibility graph whose edges correspond to plausible merges. Moreover, we consider several different measures for computing the citation count of a merged article. For the measure used by Google Scholar, we give an algorithm that maximizes the H-index in linear time if the compatibility graph has constant-size connected components. In contrast, if we allow to merge arbitrary articles (that is, for compatibility graphs that are cliques), then already increasing the H-index by one is NP-hard. Experiments on Google Scholar profiles of AI researchers show that the H-index can be manipulated substantially only if one merges articles with highly dissimilar titles.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

The H-index is a widely used measure for estimating the productivity and impact of researchers, journals, and institutions. Hirsch [22] defined the index as follows: a researcher has H-index  $h$  if  $h$  of the researcher's articles have at least  $h$  citations and all other articles have at most  $h$  citations. Several publicly accessible databases such as AMiner, Google Scholar, Scopus, and Web of Science compute the H-index of researchers. Such metrics are therefore visible to hiring committees and funding agencies when comparing researchers and proposals.<sup>1</sup>

☆ An extended abstract of this article appeared at IJCAI 2015 [4]. This version provides full proof details, new kernelization results, as well as additional experiments.

\* Corresponding author at: Novosibirsk State University, ul. Pirogova 2, 630090 Novosibirsk, Russian Federation.

E-mail addresses: [rvb@nsu.ru](mailto:rvb@nsu.ru) (R. van Bevern), [christian.komusiewicz@uni-jena.de](mailto:christian.komusiewicz@uni-jena.de) (C. Komusiewicz), [rolf.niedermeier@tu-berlin.de](mailto:rolf.niedermeier@tu-berlin.de) (R. Niedermeier), [manuel.sorge@tu-berlin.de](mailto:manuel.sorge@tu-berlin.de) (M. Sorge), [toby.walsh@nicta.com.au](mailto:toby.walsh@nicta.com.au) (T. Walsh).

<sup>1</sup> Our study on H-index manipulation is not meant to endorse or discourage the use of the H-index as an evaluation tool. In this regard, we merely aim to raise awareness for the various possibilities for manipulation.

Although the H-index of Google Scholar profiles is computed automatically, profile owners can still affect their H-index by merging articles in their profile. The intention of providing the option to merge articles is to enable researchers to identify different versions of the same article. For example, a researcher may want to merge a journal version and a version on arXiv.org, which are found as two different articles by Google's web crawlers. This may decrease a researcher's H-index if both articles counted towards it before merging, or increase the H-index since the merged article may have more citations than each of the individual articles. Since the Google Scholar interface permits to merge arbitrary pairs of articles, this leaves the H-index of Google Scholar profiles vulnerable to manipulation by insincere authors.

In extreme cases, the merging operation may yield an arbitrarily large H-index even if each single article is cited only a few times: If the author has, for example,  $h^2$  articles that are cited once, each by a distinct article from another author, then the H-index of the profile is 1. Creating  $h$  merged articles, each consisting of  $h$  original articles, gives a profile with H-index  $h$ . This is the maximum H-index achievable with  $h^2$  citations.

Increasing the H-index even by small values could be tempting in particular for young researchers, who are scrutinized more often than established researchers.<sup>2</sup> Hirsch [22] estimates that, for the field of physics, the H-index of a successful researcher increases by roughly one per year of activity. Hence, an insincere author might try to save years of research work with the push of a few buttons.

H-index manipulation by article merging has been studied by de Keijzer and Apt [9]. In their model, each article in a profile comes with a number of citations. Merging two articles, one with  $x$  and one with  $y$  citations, replaces these articles by a new article with  $x + y$  citations. The obtained article may then be merged with further articles to obtain articles with even higher citation numbers. In this model, one can determine in polynomial time whether it is possible to improve the H-index by merging, but maximizing the H-index by merging is strongly NP-hard [9]. We extend the results of de Keijzer and Apt [9] as follows.

1. We propose two further ways of measuring the number of citations of a merged article. One of them seems to be the measure used by Google Scholar.
2. We propose a model for restricting the set of allowed merge operations. Although Google Scholar allows merges between arbitrary articles, such a restriction is well motivated: An insincere author may try to merge only similar articles in order to conceal the manipulation.
3. We consider the variant of H-index manipulation in which only a limited number of merges may be applied in order to achieve a desired H-index. This is again motivated by the fact that an insincere author may try to conceal the manipulation by performing only few changes to her or his own profile.
4. We analyze each problem variant presented here within the framework of parameterized computational complexity [8,12,19,25]. That is, we identify parameters  $p$ —properties of the input measured in integers—and aim to design fixed-parameter algorithms, which have running time  $f(p) \cdot n^{O(1)}$  for a computable function  $f$  independent of the input size  $n$ . In some cases, this allows us to give efficient algorithms for realistic problem instances despite the NP-hardness of the problems in general. We also show parameters that presumably cannot lead to fixed-parameter algorithms by showing some problem variants to be  $W[1]$ -hard for these parameters.
5. We evaluate our theoretical findings by performing experiments with real-world data based on the publication profiles of AI researchers. In particular, we use profiles of some young and up-and-coming researchers from the 2011 and 2013 editions of the IEEE “AI's 10 to watch” list [1,33].

**Related work** Using the models introduced here, Elkind and Pavlou [27] recently studied manipulation for two alternatives to the H-index: the  $i10$ -index, the number of articles with at least ten citations, and the  $g$ -index [14], which is the largest number  $g$  such that the  $g$  most-cited articles are cited at least  $g$  times *on average*. They also considered the scenario where merging articles can influence the profiles of *other* authors. In a follow-up work to our findings, we analyzed the complexity of *unmerging* already merged articles so to manipulate the H-index with respect to the citation measures introduced here [5]. Notably, in the model corresponding to Google Scholar, the complexity is much lower for unmerging rather than for merging articles.

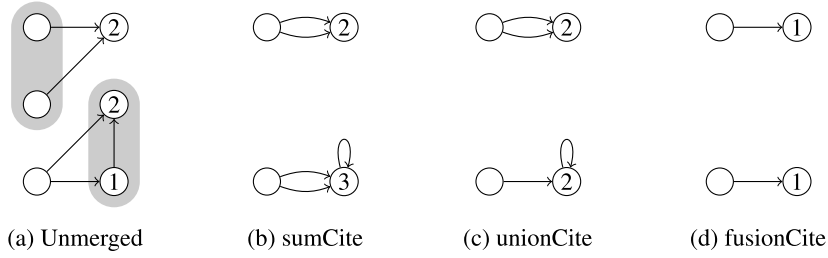
A different way of manipulating the H-index is by strategic self-citations [10,28]; Bartneck and Kokkelmans [3] consider approaches to detect these. Strategic self-citations take some effort and are irreversible. Thus, they can permanently damage an author's reputation. In comparison, article merging is easy, reversible and usually justified.

Bodlaender and van Kreveld [6] showed that, in a previous version of the Google Scholar interface, which only allowed merges of articles displayed together on one page, it was NP-hard to decide whether a given set of articles can be merged at all.

The problem of maximizing the H-index in the model of de Keijzer and Apt [9] is essentially a special case of the scheduling problems BIN COVERING [2,7] and MACHINE COVERING [20,29].

A considerable body of work on manipulation can be found in the computational social choice literature [15,16]. If we view citations as articles voting on other articles, then the problem we consider here is somewhat analogous to strategic candidacy [13].

<sup>2</sup> In fact, for senior researchers with many citations, the H-index is barely more expressive than the total citation count [32].



**Fig. 1.** Vertices represent articles in our profile  $W$ , arrows represent citations, numbers are citation counts (note that, in general, there may be articles in  $V \setminus W$ , which are not in our profile and not displayed here). The articles on a gray background in (a) have been merged in (b)–(d), and citation counts are given according to the measures sumCite, unionCite, and fusionCite, respectively. The arrows represent the citations counted by the corresponding measure.

### 1.1. Our models

We propose two new models for the merging of articles. These models take into consideration two aspects that are not captured by the model of de Keijzer and Apt [9]:

1. The number of citations of an article resulting from a merge is not necessarily the sum of the citations of the merged articles. This is in particular the case for Google Scholar.
2. In order to hide manipulation, it would be desirable to only merge related articles instead of arbitrary ones. For example, one could only merge articles with similar titles.

To capture the second aspect, our model allows for constraints on the compatibility of articles. To capture the first aspect, we represent citations not by mere citation counts, but using a directed *citation graph*  $D = (V, A)$ . The vertices of  $D$  are the articles of the author's profile plus the articles that cite them, and there is an arc  $(u, v)$  in  $D$  if article  $u$  cites article  $v$ .

To simplify notation, we assume from now on that we are an author who wants to maximize her or his H-index by merging articles. Let  $W \subseteq V$  denote the articles in our profile. In the following, these articles are called *atomic articles* and we aim to maximize our H-index by merging some articles in  $W$ . The result of a sequence of article merges is a partition  $\mathcal{P}$  of  $W$ . We call each part  $P \in \mathcal{P}$  with  $|P| \geq 2$  a *merged article*. Note that having a merged article  $P$  corresponds to performing  $|P| - 1$  successive merges on the articles contained in  $P$ . It is sometimes convenient to alternate between the partitioning and merging interpretations.

The aim is to find a partition  $\mathcal{P}$  of  $W$  with a large H-index, where the *H-index of a partition*  $\mathcal{P}$  is the largest number  $h$  such that there are at least  $h$  parts  $P \in \mathcal{P}$  whose number  $\mu(P)$  of citations is at least  $h$ . Herein, we have multiple possibilities of defining the measure  $\mu(P)$  of citations of an article in  $\mathcal{P}$ . Before describing these possibilities, we introduce some further notation.

Let  $\deg_D^{\text{in}}(v)$  denote the indegree of an article  $v$  in the citation graph  $D$ , that is, its number of citations. Analogously, we will use  $\deg_D^{\text{out}}(v)$  to denote the outdegree of an article  $v$  in the citation graph  $D$ . Moreover, let  $N_D^{\text{in}}(v) := \{u \mid (u, v) \in A\}$  denote the set of articles that cite  $v$  and let  $N_{D-W}^{\text{in}}(v) := \{u \mid (u, v) \in A \wedge u \notin W\}$  denote the set of articles that cite  $v$  and are not contained in  $W$  (thus, the articles in  $N_{D-W}^{\text{in}}(v)$  cannot be merged). For each part  $P \in \mathcal{P}$ , we consider the following three citation measures for defining the number  $\mu(P)$  of citations of  $P$ . They are illustrated in Fig. 1. The measure

$$\text{sumCite}(P) := \sum_{v \in P} \deg_D^{\text{in}}(v)$$

defines the number of citations of a merged article  $P$  to be the sum of the citations of the atomic articles it contains. This is the measure proposed by de Keijzer and Apt [9]. In contrast, the measure

$$\text{unionCite}(P) := \left| \bigcup_{v \in P} N_D^{\text{in}}(v) \right|$$

defines the number of citations of a merged article  $P$  as the number of distinct atomic articles citing at least one atomic article in  $P$ . We verified empirically that, at the time of writing, Google Scholar used the unionCite measure. The measure

$$\text{fusionCite}_{\mathcal{P}}(P) := \left| \bigcup_{v \in P} N_{D-W}^{\text{in}}(v) \right| + \sum_{P' \in \mathcal{P} \setminus \{P\}} \begin{cases} 1 & \text{if } \exists v \in P' \exists w \in P : (v, w) \in A, \\ 0 & \text{otherwise} \end{cases}$$

is, in our opinion, the most natural one: a set of merged articles is indeed considered to be one article, that is, at most one citation of a part  $P' \in \mathcal{P}$  to a part  $P \in \mathcal{P}$  is counted. In contrast to the two other measures, merging two articles under the fusionCite measure may lower the number of citations of the resulting article and of other articles. Note that, in contrast to

**Table 1**

The complexity of H-INDEX MANIPULATION for the citation measures sumCite, unionCite, fusionCite, and the parameters “H-index  $h$  to achieve”, “size  $c$  of the largest connected component of the compatibility graph  $G$ ”, and “number  $k$  of allowed article merges”. The last row shows the complexity of the variant where we only aim to improve the H-index compared to the profile without merges.

	sumCite	unionCite	fusionCite
$h$	W[1]-hard (Corollary 1), but linear-time for constant $h$ if $G$ is a clique (Corollary 2)		
$c$	Solvable in $O(3^c \cdot (n+m))$ time (Theorem 1)		NP-hard even for $c=2$ (Theorem 2)
$k$	W[1]-hard (Theorem 3), but solvable in $O(9^k k \cdot (k+n+m))$ time if $G$ is a clique (Theorem 4)	W[1]-hard even if $G$ is a clique (Theorem 5)	
	Improving H-index by one is NP-hard (Theorem 3), but polynomial-time solvable if $G$ is a clique [9]	Improving H-index by one is NP-hard even if $G$ is a clique (Theorem 6)	

unionCite and sumCite, the number of citations of an article according to fusionCite $_{\mathcal{P}}$  may depend on the partition  $\mathcal{P}$ . We omit the index  $\mathcal{P}$  where it is clear from the context.

To model constraints on permitted article merges, we consider an undirected *compatibility graph*  $G = (V, E)$ . We call two articles *compatible* if they are adjacent in  $G$ . We say that a partition  $\mathcal{P}$  of the articles  $W$  *complies* with  $G$  if, for each part  $P \in \mathcal{P}$ , all articles in  $P$  are pairwise compatible, that is, if the subgraph  $G[P]$  of  $G$  induced by  $P$  is a clique. Thus, if the compatibility graph  $G$  is a clique, then there are no constraints: all partitions of  $W$  comply with  $G$  in this case.

Formally, for each measure  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$ , we are interested in the following problem:

H-INDEX MANIPULATION( $\mu$ )

**Input:** A citation graph  $D = (V, A)$ , a compatibility graph  $G = (V, E)$ , a set  $W \subseteq V$  of articles, and a non-negative integer  $h$ .

**Question:** Is there a partition of  $W$  that complies with  $G$  and that has H-index at least  $h$  with respect to measure  $\mu$ ?

Throughout this work, we use  $n := |V|$  to denote the number of input articles and  $m := |E| + |A|$  to denote the overall number of edges and arcs in the two input graphs. Moreover, we use  $G[S]$  to denote the subgraph of  $G$  induced by a subset  $S$  of its vertices.

## 1.2. Our results

We study the complexity of H-INDEX MANIPULATION with respect to several structural features of the input instances. In particular, we consider the following three parameters:

- The size  $c$  of the largest connected component in the compatibility graph  $G$ . We expect this size to be small if only reasonable merges are allowed or if all merges have to appear reasonable, that is, if compatible articles should superficially look similar.
- The number  $k$  of merges. An insincere author would hide manipulations using a small number of merges.
- The H-index to be achieved. Although one is interested in maximizing the H-index, we expect this number also to be relatively small, since even experienced researchers seldom have an H-index of greater than 50.<sup>3</sup>

Table 1 summarizes our theoretical results. For example, we find that, with respect to the unionCite measure used by Google Scholar, it is easier to manipulate the H-index if only a small number of articles can be merged into one (small  $c$ ). The unionCite measure is complex enough to make increasing the H-index by one an NP-hard problem even if the compatibility graph  $G$  is a clique. In contrast, for the sumCite measure and the compatibility graph being a clique, it can be decided in polynomial time whether the H-index can be increased by one [9].

We implemented the manipulation algorithms exploiting small  $k$  and small  $c$ . Experimental results show that all of our sample AI authors can increase their H-index by only three merges but that usually merging articles with highly dissimilar titles is required to obtain a substantial improvement.

## 1.3. Preliminaries

We analyze H-INDEX MANIPULATION with respect to its classic and its parameterized complexity. The aim of parameterized complexity theory is to analyze problem difficulty not only in terms of the input size, but also with respect to an additional parameter, typically an integer  $p$  [8,12,19,25]. Thus, formally, an instance of a parameterized problem is a pair  $(I, p)$  consisting of the input  $I$  and the parameter  $p$ . A parameterized problem with parameter  $p$  is *fixed-parameter tractable* (FPT) if there is an algorithm that decides an instance  $(I, p)$  in  $f(p) \cdot |I|^{O(1)}$  time, where  $f$  is an arbitrary computable function depending only on  $p$ . Clearly, if the problem is NP-hard, then we expect  $f$  to grow superpolynomially.

<sup>3</sup> More than 99.99% of the authors listed at [aminer.org](http://aminer.org) (accession date 2/27/2016) have an H-index of at most 50.

There are parameterized problems for which there is good evidence that they are not fixed-parameter tractable: Analogously to the concept of NP-hardness, the concept of W[1]-hardness was developed. It is widely assumed that a W[1]-hard parameterized problem cannot be fixed-parameter tractable. To show that a parameterized problem with parameter  $p'$  is W[1]-hard, a *parameterized reduction* from a known W[1]-hard parameterized problem with parameter  $p$  can be used. This is a reduction that runs in  $f(p) \cdot |I|^{O(1)}$  time and produces instances such that the parameter  $p'$  is upper-bounded by some function  $g(p)$ . Determining whether an undirected graph  $G$  has a clique of order  $p$  is well known to be W[1]-hard with respect to  $p$ .

The notion of a *problem kernel* tries to capture the existence of efficient and provably effective preprocessing rules [21, 24]. More precisely, we say that a parameterized problem has a problem kernel if every instance can be reduced in polynomial time to an equivalent instance whose size depends only on the parameter. The algorithm computing the problem kernel is called *kernelization* and is often presented as a series of data reduction rules. A data reduction rule transforms an instance  $(I, p)$  of a parameterized problem into an instance  $(I', p')$  of the same problem; a data reduction rule is *correct* if  $(I, p)$  is a yes-instance if and only if  $(I', p')$  is a yes-instance.

## 2. Compatibility graphs with small connected components

In this section, we analyze the parameterized complexity of H-INDEX MANIPULATION parameterized by the size  $c$  of the largest connected component of the compatibility graph. This parameterization is motivated by the fact that one would merge only similar articles and that usually each article is similar to only few other articles.

The following theorem shows that H-INDEX MANIPULATION is solvable in linear time for the citation measures sumCite and unionCite if  $c$  is constant. The algorithm exploits that, for these two measures, merging articles does not affect other articles. Thus, we can solve each connected component independently of the others.

**Theorem 1.** H-INDEX MANIPULATION( $\mu$ ) is solvable in  $O(3^c \cdot (n + m))$  time for  $\mu \in \{\text{sumCite}, \text{unionCite}\}$  if the connected components of the compatibility graph  $G$  have size at most  $c$ .

**Proof.** Clearly, articles from different connected components of  $G$  cannot be together in a part of any partition complying with  $G$ . Thus, independently for each connected component  $C$  of  $G$ , we compute a partition of the articles of  $C$  that complies with  $G$  and has the maximum number of parts  $P$  with  $\mu(P) \geq h$ .

We first show that this approach is correct and then show how to execute it efficiently. Obviously, if an algorithm creates a partition  $\mathcal{P}$  of the set  $W$  of our own articles that complies with  $G$  and has at least  $h$  parts  $P$  with  $\mu(P) \geq h$ , then we face a yes-instance. Conversely, if the input is a yes-instance, then there is a partition  $\mathcal{P}$  of  $W$  complying with  $G$  and having at least  $h$  parts  $P$  with  $\mu(P) \geq h$ . Consider any connected component  $C$  of  $G$  and the restriction  $\mathcal{P}_C = \{P \in \mathcal{P} \mid P \subseteq V(C)\}$  of  $\mathcal{P}$  to  $C$ , where  $V(C)$  is the vertex set of  $C$ . Note that each part in  $\mathcal{P}$  is either contained in  $V(C)$  or disjoint from it and, thus,  $\mathcal{P}_C$  is a partition of  $V(C)$ . Moreover, merging articles of one connected component does not affect the number of citations of articles in other connected components with respect to sumCite or unionCite. Thus, if we replace the sets of  $\mathcal{P}_C$  in  $\mathcal{P}$  by a partition of  $C$  that has a maximum number of parts  $P$  with  $\mu(P) \geq h$ , then we obtain a partition that still has H-index at least  $h$ . Thus, our algorithm indeed finds a partition with H-index at least  $h$ .

We now show how to compute, for each connected component  $C$  of  $G$ , a partition that maximizes the number of parts with at least  $h$  citations. In order to achieve a running time of  $O(3^c \cdot (n + m))$ , we employ dynamic programming. First, for each connected component  $C$  of  $G$  and every  $V' \subseteq V(C)$ , we initialize a table

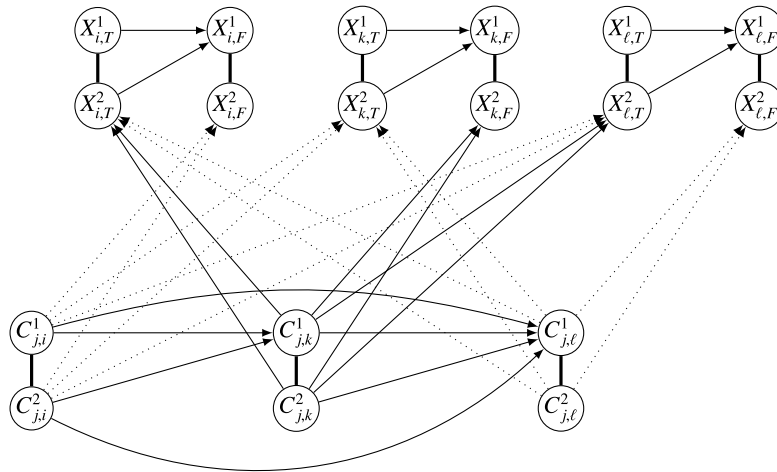
$$Q[V'] := \begin{cases} 1 & \text{if } G[V'] \text{ is a clique and } \mu(V') \geq h, \\ 0 & \text{if } G[V'] \text{ is a clique and } \mu(V') < h, \text{ and} \\ -\infty & \text{otherwise.} \end{cases}$$

A table entry  $Q[V']$  thus stores whether merging  $V'$  results in an article with at least  $h$  citations. Obviously, if  $G[V']$  is not a clique, then  $V'$  cannot be a part in any partition complying with  $G$ . Therefore, we set  $Q[V'] := -\infty$  in this case. All table entries  $Q[V']$  for all vertex subsets  $V' \subseteq V(C)$  of all connected components of  $G$  can be computed in  $O(2^c \cdot (n + m))$  time.

Now, for every vertex subset  $V' \subseteq V(C)$  of a connected component  $C$ , we define  $T[V']$  to be the maximum number of parts  $P$  with  $\mu(P) \geq h$  in any partition of  $V'$  complying with  $G$ . Obviously, we have the recurrence relation

$$T[V'] = \begin{cases} 0 & \text{if } V' = \emptyset, \text{ and} \\ \max_{V'' \subsetneq V'} (T[V' \setminus V''] + Q[V'']) & \text{otherwise.} \end{cases}$$

After computing the table  $Q$ , we can compute  $T[V(C)]$  for each connected component  $C$  using dynamic programming: compute  $T[V']$  for each subset  $V' \subseteq V(C)$  in the order of increasing cardinality. To this end, for each such subset  $V' \subseteq V(C)$ , we iterate over all subsets  $V'' \subsetneq V'$ . Thus, the computation of  $T$  works in  $O(3^c)$  time since there are at most  $3^c$  partitions of  $V(C)$  into  $V(C) \setminus (V' \cup V'')$ ,  $V' \setminus V''$ , and  $V''$ . Herein, the factor  $c$  accounts for operations with sets of cardinality at most  $c$ . Thus, the total running time is  $O(2^c \cdot (n + m) + 3^c c) \subseteq O(3^c \cdot (n + m))$ .  $\square$



**Fig. 2.** Construction for a clause  $c_j = (x_i \vee x_k \vee x_l)$ . Undirected bold edges belong to the compatibility graph  $G$ , directed arcs are citations in the citation graph  $D$ . Only vertices in  $W$  and citations between vertices in  $W$  are shown. Some arcs are dotted to keep the picture clean. Observe that all citations point from the bottom to the top or from the left to the right, and thus form a directed acyclic graph.

We have seen that H-INDEX MANIPULATION is solvable in linear time for the citation measures sumCite and unionCite if the compatibility graph has constant-size connected components. In contrast, constant-size components of the compatibility graph do not help when the fusionCite measure is used. This we show by a reduction from the NP-hard 3-BOUNDED POSITIVE 1-IN-3-SAT problem [11].

**Theorem 2.** H-INDEX MANIPULATION(fusionCite) is NP-hard even if all of the following conditions hold:

- i) the largest connected component of the compatibility graph has size two,
- ii) the citation graph is acyclic, and
- iii) the input instance has H-index  $h - 1$ .

Regarding (ii), note that citation graphs are often acyclic in practice as papers tend to cite only earlier papers. Thus, it is important that Theorem 2 does not require cycles in the citation graph.

**Proof.** We prove Theorem 2 using a polynomial-time many-one reduction from the NP-hard 3-BOUNDED POSITIVE 1-IN-3-SAT problem [11].

#### 3-BOUNDED POSITIVE 1-IN-3-SAT

**Input:** A formula  $\phi$  in 3-conjunctive normal form containing only positive literals and with each literal contained in at most three clauses.

**Question:** Is there a truth assignment to the variables of  $\phi$  that sets exactly one variable per clause to “true”?

Let  $n$  be the number of variables of  $\phi$  and let  $m$  be the number of clauses. If  $m + n$  is odd, then we simply duplicate the instance. If  $(m + n)/2 < 18$ , then we solve  $\phi$  using brute force and output a corresponding trivial yes- or no-instance of H-INDEX MANIPULATION(fusionCite). Otherwise, we now create an instance of H-INDEX MANIPULATION with  $h := m + n$ . The construction is illustrated in Fig. 2.

For each variable  $x_i$  of  $\phi$ , we introduce a variable gadget consisting of

- four articles  $X_{i,T}^1, X_{i,T}^2, X_{i,F}^1$ , and  $X_{i,F}^2$ ,
- two edges  $\{X_{i,T}^1, X_{i,T}^2\}$  and  $\{X_{i,F}^1, X_{i,F}^2\}$  in the compatibility graph  $G$ , and
- two arcs  $(X_{i,T}^1, X_{i,F}^1)$  and  $(X_{i,T}^2, X_{i,F}^2)$  in the citation graph  $D$ .

Merging the pair  $\{X_{i,T}^1, X_{i,T}^2\}$  will correspond to setting  $x_i$  to true, merging the pair  $\{X_{i,F}^1, X_{i,F}^2\}$  will correspond to setting  $x_i$  to false. For each clause  $c_j = (x_i \vee x_k \vee x_l)$ , we add a clause gadget consisting of

- six articles  $C_{j,z}^1, C_{j,z}^2$  for  $z \in \{i, k, l\}$ ,
- three edges  $\{C_{j,z}^1, C_{j,z}^2\}$  for  $z \in \{i, k, l\}$  in the compatibility graph  $G$ , and

- six arcs  $(C_{j,i}^1, C_{j,k}^1), (C_{j,i}^2, C_{j,k}^1), (C_{j,i}^1, C_{j,\ell}^1), (C_{j,i}^2, C_{j,\ell}^1), (C_{j,k}^1, C_{j,\ell}^1),$  and  $(C_{j,k}^2, C_{j,\ell}^1)$  in the citation graph  $D$ .

Merging a pair  $\{C_{j,z}^1, C_{j,z}^2\}$  for  $z \in \{i, k, \ell\}$  will correspond to setting the literal  $x_z$  of  $c_j$  to true.

To connect the clause gadget for the clause  $c_j = (x_i \vee x_k \vee x_\ell)$  to the corresponding variable gadgets, for all  $z \in \{i, k, \ell\}$  and all  $y \in \{i, k, \ell\} \setminus \{z\}$ , we add the arcs  $(C_{j,z}^1, X_{z,F}^2), (C_{j,z}^2, X_{z,F}^2), (C_{j,z}^1, X_{y,T}^2),$  and  $(C_{j,z}^2, X_{y,T}^2)$ .

Observe that the constructed citation graph is acyclic since each variable gadget and each clause gadget is acyclic and all other arcs are from clause gadgets to variable gadgets. Moreover, since each variable occurs in at most three clauses of  $\phi$  and each clause has only three variables, every created article has at most  $3 \cdot 3 \cdot 2 = 18$  incoming citations. Since  $(m+n)/2 = h/2 \geq 18$ , we can finally add, for each of the created articles, a distinct set of articles to  $D$  such that each pair  $\{X_{i,T}^1, X_{i,T}^2\}, \{X_{i,F}^1, X_{i,F}^2\},$  or  $\{C_{j,z}^1, C_{j,z}^2\}$  is cited exactly  $h$  times in total and each single article is cited less than  $h$  times. In particular, because articles of the form  $X_{i,T}^1$  and  $C_{j,z}^2$  are never cited by any article in variable or clause gadgets, we can add these additional citations so that  $X_{i,T}^1$  is cited once and  $X_{i,T}^2$  is cited  $h-1$  times for each variable  $x_i$  of  $\phi$  and so that  $C_{j,z}^2$  in some pair  $\{C_{j,z}^1, C_{j,z}^2\}$  for a clause  $c_j$  is cited once, whereas  $C_{j,z}^1$  is cited  $h-1$  times. This concludes the construction of our H-INDEX MANIPULATION(fusionCite) instance. Note that, for each of the  $h = m+n$  clauses and variables, we created at least one article with  $h-1$  citations, which gives an instance with H-index  $h-1$ . We now prove the correctness of the presented reduction.

First, if we have an assignment for  $\phi$  that sets exactly one variable in each clause to true, then we merge the pair  $\{X_{i,T}^1, X_{i,T}^2\}$  for all true variables  $x_i$  and merge the pairs  $\{C_{j,i}^1, C_{j,i}^2\}$  for all clauses  $c_j$  containing  $x_i$ . We will thus get  $h = m+n$  articles with  $h$  citations. To show the converse, we first make two important observations:

1. For each variable  $x_i$  of  $\phi$ , at most one of the pairs  $p_1 := \{X_{i,T}^1, X_{i,T}^2\}$  and  $p_2 := \{X_{i,F}^1, X_{i,F}^2\}$  can be merged into an article  $P$  with  $\text{fusionCite}(P) \geq h$ : Observe that the sum of the citations of each pair is exactly  $h$ . However, if both  $p_1$  and  $p_2$  are merged, then the article resulting from merging  $p_2$  will get at most  $h-1$  citations: it gets a citation from each of the articles of  $p_1$ , which will be counted only as one citation after merging  $p_1$ .
2. For each clause  $c_j = (x_i \vee x_k \vee x_\ell)$ , at most one of the pairs  $p_1 := \{C_{j,i}^1, C_{j,i}^2\}, p_2 := \{C_{j,k}^1, C_{j,k}^2\}, p_3 := \{C_{j,\ell}^1, C_{j,\ell}^2\}$  can be merged into an article  $P$  with  $\text{fusionCite}(P) \geq h$ : Assume that  $p_x$  for some  $x \in \{1, 2, 3\}$  is merged into an article  $P$  with  $\text{fusionCite}(P) \geq h$ . Then, no pair  $p_y$  with  $y > x$  can be merged into such an article:  $p_y$  can get at most  $h-1$  citations when merged since  $p_y$  gets a citation from each of the two articles of  $p_x$ . By the same argument, if any  $p_y$  with  $y < x$  could be merged into one article  $P$  with  $\text{fusionCite}(P) \geq h$ , then this would contradict the assumption that merging  $p_x$  yields such an article.

Since we ask for merging articles in order to increase the H-index to  $h := m+n$ , which is precisely the number of variables and clauses in the input formula, we have to create at least one article with  $h$  citations for each variable gadget and for each clause gadget. That is, if we can achieve H-index  $h$ , then, for each variable gadget and each clause gadget, exactly one pair is merged into an article with at least  $h$  citations. Moreover, if, for some clause  $c_j = (x_i \vee x_k \vee x_\ell)$  the pair  $\{C_{j,z}^1, C_{j,z}^2\}$  is merged for some  $z \in \{i, k, \ell\}$ , then the pair  $\{X_{z,F}^1, X_{z,F}^2\}$  cannot be merged into an article with  $h$  citations since it gets one citation from each of  $C_{j,z}^1$  and  $C_{j,z}^2$ . It follows that  $\{X_{z,T}^1, X_{z,T}^2\}$  has to be merged. Moreover, for  $y \in \{i, k, \ell\} \setminus \{z\}$ , the pair  $\{X_{y,T}^1, X_{y,T}^2\}$  cannot be merged into an article with  $h$  citations, since it gets one citation from each of  $C_{j,z}^1$  and  $C_{j,z}^2$ . It follows that  $\{X_{y,F}^1, X_{y,F}^2\}$  has to be merged.

Thus, we obtain an assignment for  $\phi$  that sets exactly one variable of each clause to true by setting those variables  $x_i$  to true for which the pair  $\{X_{i,T}^1, X_{i,T}^2\}$  is merged into an article with at least  $h$  citations.  $\square$

### 3. Merging few articles or increasing the H-Index by one

In this section, we consider two variants of H-INDEX MANIPULATION: CAUTIOUS H-INDEX MANIPULATION, where we allow to merge at most  $k$  articles, and H-INDEX IMPROVEMENT, where we ask whether it is possible to increase the H-index at all.

CAUTIOUS H-INDEX MANIPULATION is motivated by the fact that insincere authors could try to conceal their manipulation by merging only few articles. Formally, the problem is defined as follows, where  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$  as before.

CAUTIOUS H-INDEX MANIPULATION( $\mu$ )

**Input:** A citation graph  $D = (V, A)$ , a compatibility graph  $G = (V, E)$ , a set  $W \subseteq V$  of articles, and non-negative integers  $h$  and  $k$ .

**Question:** Is there a partition  $\mathcal{P}$  of  $W$  that

- i) complies with  $G$ ,
- ii) has H-index at least  $h$  with respect to  $\mu$ , and
- iii) is such that the number  $\sum_{P \in \mathcal{P}} (|P| - 1)$  of merges is at most  $k$ ?



We show that CAUTIOUS H-INDEX MANIPULATION parameterized by  $k$  is fixed-parameter tractable only for the sumCite measure and when the compatibility graph is a clique. Allowing arbitrary compatibility graphs or using more complex measures leads to  $W[1]$ -hardness with respect to  $k$ .

The second problem considered in this section, H-INDEX IMPROVEMENT( $\mu$ ), was introduced by de Keijzer and Apt [9]; it is formally defined as follows.

H-INDEX IMPROVEMENT( $\mu$ )

**Input:** A citation graph  $D = (V, A)$ , a compatibility graph  $G = (V, E)$ , and a set  $W \subseteq V$  of articles.

**Question:** Is there a partition  $\mathcal{P}$  of  $W$  that complies with  $G$  and has a larger H-index with respect to  $\mu$  than the partition of  $W$  into singletons?

Theorem 2(iii) shows that H-INDEX IMPROVEMENT(fusionCite) is NP-hard even when the compatibility graph has connected components of size at most two. However, de Keijzer and Apt [9] gave a polynomial-time algorithm for H-INDEX IMPROVEMENT(sumCite) if the compatibility graph is a clique. We flesh out the boundary between hardness and tractability by proving that more general compatibility graphs lead to NP-hardness in the case of sumCite and that even clique compatibility graphs lead to NP-hardness in the case of unionCite and fusionCite.

First, we consider the case of general compatibility graphs. Here, we obtain hardness results for both problem variants by reductions from MULTICOLORED CLIQUE:

MULTICOLORED CLIQUE

**Input:** An  $\ell$ -partite undirected graph  $H$  along with the  $\ell$  partite sets.

**Question:** Is there a clique with  $\ell$  vertices contained in  $H$ ?

MULTICOLORED CLIQUE is known to be NP-hard and  $W[1]$ -hard with respect to  $\ell$  [17]. Both hardness results hold for sumCite and, thus, also for unionCite and fusionCite.

**Theorem 3.** *Parameterized by  $k$ , CAUTIOUS H-INDEX MANIPULATION(sumCite) is  $W[1]$ -hard. Moreover, H-INDEX IMPROVEMENT(sumCite) is NP-hard.*

The theorem directly follows from the following lemma, from which we will derive another hardness result in the next section.

**Lemma 1.** *There is a polynomial-time many-one reduction from MULTICOLORED CLIQUE to CAUTIOUS H-INDEX MANIPULATION(sumCite) with  $k = \ell - 1$  and  $h = \ell$  and to H-INDEX IMPROVEMENT.*

**Proof.** The two reductions from the MULTICOLORED CLIQUE problem differ only in specifying the H-index that we want to achieve and the upper bound on the number of merges for CAUTIOUS H-INDEX MANIPULATION. The reductions work as follows. We create a citation graph  $D$ , a compatibility graph  $G$ , and a set of articles  $W$ , such that the instance  $(D, G, W, h := \ell, k := \ell - 1)$  of CAUTIOUS H-INDEX MANIPULATION and the instance  $(D, G, W)$  of H-INDEX IMPROVEMENT are yes-instances if and only if  $(H, \ell)$  is a yes-instance for MULTICOLORED CLIQUE.

Our CAUTIOUS H-INDEX MANIPULATION and H-INDEX IMPROVEMENT instances have an article set  $W = W_{\geq} \uplus W_{<}$ , where  $W_{<} := V(H)$  and  $W_{\geq}$  consists of  $\ell - 1$  new articles. For each article  $w \in W_{\geq}$  we introduce a set of  $\ell$  articles that are not contained in  $W$  and that cite  $w$  and no other article. Similarly, for each article  $w \in W_{<}$  we introduce one article not in  $W$  that cites  $w$  and no other article. In this way, we have implicitly defined the citation graph  $D$ . Next, we construct the compatibility graph  $G$  from  $H$  by adding each article in  $W_{\geq}$  as an independent vertex. This concludes the construction. Clearly, we can carry it out in polynomial time. Note that the reduction is a parameterized reduction from MULTICOLORED CLIQUE parameterized by  $\ell$  to CAUTIOUS H-INDEX MANIPULATION parameterized by  $k$  since  $k = \ell - 1$ .

Now we prove the equivalence of the three instances. If the MULTICOLORED CLIQUE instance  $(H, \ell)$  is a yes-instance, then there is a clique  $S$  of size  $\ell$  in  $H$ . Merging the corresponding articles  $S \subseteq W_{<}$  complies with the compatibility graph and, hence, yields a merged article with  $\ell$  citations. Together with the  $\ell - 1$  articles in  $W_{\geq}$ , this results in  $\ell$  articles with  $\ell$  citations and, hence, H-index at least  $\ell = h$ . Furthermore, exactly  $\ell - 1$  merges are performed in this way, implying that the CAUTIOUS H-INDEX MANIPULATION instance is a yes-instance.

Note that the H-index of the singleton partition  $\mathcal{W}$  of  $W$  is  $\ell - 1$ . That is, the CAUTIOUS H-INDEX MANIPULATION instance asks to increase the H-index of  $\mathcal{W}$  by one. Thus, clearly, if the CAUTIOUS H-INDEX MANIPULATION instance is a yes-instance, then also the H-INDEX IMPROVEMENT instance is.

Finally, assume that the H-INDEX IMPROVEMENT instance is yes. Then there is a merged article  $S$  with  $\ell$  citations. Since only articles in  $W_{<}$  can be merged,  $S$  consists of at least  $\ell$  articles. Furthermore,  $G[S] = H[S]$  is a clique since the merging has to comply with  $G$ . Hence, the MULTICOLORED CLIQUE instance is yes, concluding the proof that all three instances are equivalent.  $\square$



Now we restrict the compatibility graph to be a clique, meaning that arbitrary pairs of articles can be merged. Recall that H-INDEX IMPROVEMENT(sumCite) is polynomial-time solvable in this case [9]. We also achieve a (fixed-parameter) tractability result for CAUTIOUS H-INDEX MANIPULATION(sumCite) parameterized by the number  $k$  of article merges.

**Theorem 4.** *If the compatibility graph  $G$  is a clique, then CAUTIOUS H-INDEX MANIPULATION(sumCite) is solvable in  $O(9^k k \cdot (k + n + m))$  time, where  $k$  is the number of allowed article merges.*

**Proof.** Assume that  $(D, G, W, h, k)$  is a yes-instance and let  $\mathcal{P}$  be a partition of  $W$  with H-index at least  $h$  and at most  $k$  merges. Let  $M := \{v \in W \mid v \in P, P \in \mathcal{P}, |P| \geq 2\}$  be the set of articles that have been merged with other articles, and let  $W' := \{v \in W \mid \deg_D^{\text{in}}(v) \leq h\}$  be the set of articles with at most  $h$  citations. Let  $v_1, v_2, \dots$  be the articles of  $W'$  ordered by non-increasing citation counts. We claim that we may assume that  $M = \{v_1, \dots, v_{|M|}\}$ . Otherwise, we are in one of the following cases:

- Case 1. There is an article  $v \in M$  with more than  $h$  citations. That is,  $v \in P \in \mathcal{P}$  and  $|P| \geq 2$ . In this case, we may simply split  $P$  into  $P \setminus \{v\}$  and  $\{v\}$  without dropping the H-index of  $\mathcal{P}$  below  $h$ .
- Case 2. There is an article  $v_i \in M$  with  $i > |M|$ . That is,  $v_i \in P \in \mathcal{P}$  with  $|P| \geq 2$ . Then, since the compatibility graph is a clique, we may replace  $v_i$  in  $P$  with an arbitrary article  $v_j \notin M$  and  $j \leq |M|$  (which clearly exists) without decreasing the H-index of  $\mathcal{P}$ .

Since at most  $k$  article merges are allowed, we have  $|M| \leq 2k$ . Hence, if there is a solution, then there is also one where all merged articles are within  $\{v_1, \dots, v_{2k}\}$ . Thus, we can remove all edges from the compatibility graph  $G$  that are incident with articles of at least  $h$  citations and discard all articles  $v_j$  with  $j > 2k$ . In this way, we obtain an instance with a compatibility graph that contains at most  $2k$  vertices. We now obtain the claimed fixed-parameter tractability result by adapting the dynamic programming algorithm behind Theorem 1.

Since the only nontrivial connected component  $C$  of the compatibility graph after the above preprocessing is a clique, we apply the algorithm only to  $C$ . Thus, the auxiliary table  $Q$ , used to store whether merging a set  $V'$  of articles creates an article with at least  $h$  citations, may ignore the compatibility graph. More formally, for all  $V' \subseteq V(C)$ , we let

$$Q[V'] := \begin{cases} 1 & \text{if } \mu(V') \geq h, \\ 0 & \text{otherwise.} \end{cases}$$

In order to ensure that we make at most  $k$  merges, we need an additional index in the main table  $T$ . More precisely, for a set  $V' \subseteq V(C)$  of vertices, let  $T[V', k]$  be the maximum number of parts  $P$  with  $\mu(P) \geq h$  in any partition of  $V'$  that can be obtained from the singleton partition by performing at most  $k$  merges. Then,

$$T[V', k] = \begin{cases} 0 & \text{if } V' = \emptyset, \\ 0 & \text{if } k \leq 0, \text{ and} \\ \max_{V'' \subsetneq V'} (T[V' \setminus V'', k - (|V''| - 1)] + Q[V'']) & \text{otherwise.} \end{cases}$$

Since the ground set  $V(C)$  of articles considered in the dynamic programming table has size at most  $2k$ , this algorithm has a running time of  $O(9^k k \cdot (k + n + m))$ .  $\square$

For the unionCite and fusionCite measure, we obtain hardness results for both CAUTIOUS H-INDEX MANIPULATION and H-INDEX IMPROVEMENT; the (parameterized) reductions are from the INDEPENDENT SET problem.

#### INDEPENDENT SET

**Input:** An undirected graph  $H$  and a non-negative integer  $\ell$ .

**Question:** Is there an independent set of size at least  $\ell$  in  $H$ , that is, a set of  $\ell$  pairwise nonadjacent vertices?

INDEPENDENT SET is NP-hard and W[1]-hard with respect to  $\ell$  [12].

**Theorem 5.** *For  $\mu \in \{\text{unionCite}, \text{fusionCite}\}$ , CAUTIOUS H-INDEX MANIPULATION( $\mu$ ) is W[1]-hard parameterized by  $k$  even if the compatibility graph is a clique.*

**Proof.** Let  $(H, \ell)$  be an instance of INDEPENDENT SET. We construct an instance  $(D, G, W, h, k := \ell - 1)$  of CAUTIOUS H-INDEX MANIPULATION that is a yes-instance if and only if  $(H, \ell)$  is a yes-instance for INDEPENDENT SET. Clearly, this is a parameterized reduction with respect to  $\ell$  and  $k$ .

Let  $n := |V(H)|$  and  $h := \ell n$ . Without loss of generality, we assume that  $n > \ell > 1$ . Our CAUTIOUS H-INDEX MANIPULATION instance has an article set  $W = W_{\geq} \uplus W_{<}$ , where  $W_{<} := V(H)$  and  $W_{\geq}$  consists of  $h - 1$  new articles. Next, for each article  $w \in W_{\geq}$ , we introduce  $h$  new articles not in  $W$  that cite  $w$  and no other article. The citations of the articles in  $W_{<}$

are defined as follows. For each pair of adjacent vertices  $u, v \in V(H)$ , we introduce a new article  $e_{\{u,v\}}$  not contained in  $W$  that cites the articles  $u, v \in W_{<}$  and no other articles. Furthermore, we increase the citation counts of each article in  $W_{<}$  to exactly  $n$ . That is, for each article  $w \in W_{<}$  we introduce new articles not contained in  $W$  that cite only  $w$  until  $w$  has  $n$  citations. The compatibility graph  $G$  is a clique. This concludes the construction.

Clearly, the construction can be carried out in polynomial time. Moreover, the reduction is a parameterized reduction from INDEPENDENT SET parameterized by  $\ell$  to CAUTIOUS H-INDEX MANIPULATION parameterized by  $k$  since  $k = \ell - 1$ . Note that no article in  $W$  cites another article in  $W$  and, hence, for any part  $P$  in a partition of  $W$ , we have  $\text{unionCite}(P) = \text{fusionCite}(P)$ .

Let us prove the correctness of the reduction. Assume first that  $(H, \ell)$  is a yes-instance and let  $S$  be an independent set of size  $\ell$  in  $H$ . Then, merging all articles of  $S$  into one article in the CAUTIOUS H-INDEX MANIPULATION instance is valid since the compatibility graph  $G$  is a clique. Furthermore, it yields a merged article  $S$  with  $\text{unionCite}(S) \geq h$  citations: Since the vertices in  $S$  are independent in  $G$ , there is no article  $e_{\{u,v\}}$  citing both  $u, v \in S$  in the CAUTIOUS H-INDEX MANIPULATION instance. Thus, the citations of the articles in  $S$  are pairwise disjoint. Together with the  $h - 1$  atomic articles in  $W_{\geq}$  we have H-index  $h$ .

Conversely, assume that  $(D, G, W, h, \ell - 1)$  is a yes-instance. Since we are allowed to merge at most  $\ell - 1$  times in order to achieve an H-index of  $h = \ell n$  and since each article in  $W_{<}$  has exactly  $n$  citations, we need to merge  $\ell$  articles of  $W_{<}$  into one article. That is, there is a part  $S \subseteq W_{<}$  in any solution for CAUTIOUS H-INDEX MANIPULATION with  $\text{unionCite}(S) \geq h$  citations. This means that the articles contained in  $S$  have pairwise disjoint sets of citations because each of them has only  $n = h/\ell$  citations. Thus,  $S$  is an independent set in  $H$ .  $\square$

The reduction for Theorem 5 exploits the fact that at most  $k$  merges are allowed. Hence, to show NP-hardness for H-INDEX IMPROVEMENT, we need a different reduction. Note that the NP-hardness for the fusionCite measure and general compatibility graphs already follows from Theorem 2(iii). We complement this result by the following theorem.

**Theorem 6.** H-INDEX IMPROVEMENT( $\mu$ ) is NP-hard for  $\mu \in \{\text{unionCite}, \text{fusionCite}\}$  even if the compatibility graph is a clique.

**Proof.** We give a polynomial-time reduction from INDEPENDENT SET. Let  $(H, \ell)$  be an instance of INDEPENDENT SET and let  $q := |E(H)|$ . Without loss of generality, we assume that  $q \geq \ell > 2$ . We now construct an instance of H-INDEX IMPROVEMENT with citation graph  $D$ , a set  $V$  of articles, and a subset  $W \subseteq V$  of own articles. The compatibility graph  $G$  will be a clique on all articles. We introduce citations so that the H-index of the singleton partition of  $W$  will be  $q - 1$ , hence the goal in the constructed instance will be to achieve H-index at least  $q$ .

The article set  $W$  is partitioned into three parts  $W = W_{\geq} \uplus W_{-1} \uplus W_{<}$ . The first part,  $W_{\geq}$ , consists of  $q - \ell - 1$  articles, and for each article  $w \in W_{\geq}$  we introduce  $q$  articles not in  $W$  that cite  $w$  and no other article. The second part,  $W_{-1}$ , consists of  $\ell$  articles, and for each article  $w \in W_{-1}$  we introduce  $q - 1$  articles not in  $W$  that cite  $w$  and no other article. The last part,  $W_{<}$ , contains the vertices of the INDEPENDENT SET instance, that is,  $W_{<} := V(H)$ . Finally, for each edge  $\{u, v\} \in E(H)$  we introduce one article  $e_{\{u,v\}}$  not in  $W$  that cites both  $u$  and  $v$ . This concludes the construction of the citation graph  $D$ . Note that the singleton partition of  $W$  has H-index  $q - 1$ . Hence, we have created an instance  $(D, G, W)$  of H-INDEX IMPROVEMENT where we are looking to increase the H-index to at least  $q$ . Clearly, we can carry out this construction in polynomial time. Furthermore, since there are no self-citations, that is, no articles in  $W$  cite each other, for any subset  $P$  of  $W$  we have  $\text{unionCite}(P) = \text{fusionCite}(P)$ . Let us now prove the equivalence of the two instances.

Assume that  $(H, \ell)$  is a yes-instance. We claim that then also the H-INDEX IMPROVEMENT instance is a yes-instance. Choose an independent set  $S$  of size  $\ell$  in  $H$ . Take each of the corresponding articles in  $S$  and merge them with the articles in  $W_{-1}$ , pairing them one by one. This creates  $\ell$  merged articles with  $q$  citations each. Together with the articles in  $W_{\geq}$ , we now have  $q - 1$  articles with  $q$  citations, some of them merged. To create another article with  $q$  citations, simply merge all articles in  $W_{<} \setminus S$  into one article: Since  $S$  is an independent set, for each article  $e_{\{u,v\}}$  citing  $W_{<}$ , either  $u$  or  $v$  is not in  $S$ . Hence, the merged article  $W_{<} \setminus S$  has  $q$  citations. Thus,  $(D, G, W)$  is a yes-instance.

Now assume that  $(D, G, W)$  is a yes-instance and let us show that also  $(H, \ell)$  is. Take a partition  $\mathcal{P}$  of  $W$  with H-index at least  $q$ . Note that any subset  $R \subseteq W_{<}$  has  $\mu(R) \geq q$  only if  $R$  is a vertex cover of  $H$  (a vertex cover of a graph is a subset  $X$  of the vertices such that each edge is incident with some vertex in  $X$ ). Hence, as there are at most  $q - 1$  parts  $P \in \mathcal{P}$  with  $P \not\subseteq W_{<}$ , there is at least one part  $P \in \mathcal{P}$  such that  $P \cap W_{<}$  is a vertex cover of  $H$ . For the sake of contradiction, assume that there are two parts  $P_1, P_2$  such that  $P_1 \cap W_{<}$  and  $P_2 \cap W_{<}$  are vertex covers for  $H$ . Then  $P_1 \cup P_2 \supseteq V(H) = W_{<}$ . Furthermore, each remaining part of  $\mathcal{P}$  contains only articles in  $W_{\geq} \cup W_{-1}$ , that is, out of these parts, at most  $q - \ell - 1 + \lfloor \ell/2 \rfloor$  can have at least  $q$  citations. However, as  $\ell > 2$ , there are at most  $q - \lfloor \ell/2 \rfloor - 1 + 2 \leq q - 1$  parts with at least  $q$  citations in  $\mathcal{P}$ , a contradiction. Thus, there is exactly one part  $P \in \mathcal{P}$  such that  $R := P \cap W_{<}$  is a vertex cover of  $H$ .

Take  $S := V(H) \setminus R$ . Note that, since  $R$  is a vertex cover of  $H$ ,  $S$  is an independent set in  $H$ ; we claim that  $S$  has size at least  $\ell$ . Since there is exactly one part in  $\mathcal{P}$  that contains a vertex cover of  $H$ , each remaining part has at least  $q$  citations and there are at least  $q - 1$  of them. This means that no two articles in  $W_{\geq} \cup W_{-1}$  are merged. Hence, each article in  $W_{-1}$  is merged into an article in  $S$ , that is,  $S$  contains at least  $\ell$  articles.  $\square$

#### 4. Achieving a moderately large H-index

We now consider the H-index that we want to achieve as a parameter. This parameter is often not very large as researchers in the early stage of their career have an H-index below 20. Even for more experienced researchers the H-index seldom exceeds 50. Hence, in many cases, the value of a desired H-index is sufficiently low to serve as useful parameter in terms of gaining efficient fixed-parameter algorithms. However, note that [Lemma 1](#) immediately implies that  $\text{H-INDEX MANIPULATION}(\mu)$  is  $W[1]$ -hard with respect to the target H-index  $h$ . This hardness also transfers to the unionCite and fusionCite measures:

**Corollary 1.**  $\text{H-INDEX MANIPULATION}(\mu)$  is  $W[1]$ -hard with respect to the target H-index  $h$  for each  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$ .

In contrast, we now show that  $\text{H-INDEX MANIPULATION}(\mu)$  is fixed-parameter tractable for any citation measure  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$  if the compatibility graph is a clique. To this end, we describe a kernelization algorithm, that is, a polynomial-time data reduction algorithm that produces an equivalent instance whose size is upper-bounded by some function of the parameter  $h$ .

One particular difficulty in designing data reduction rules for the fusionCite measure is that citations in  $D[W]$  are somewhat fragile as they may be “destroyed”, for example, if two adjacent vertices in  $W$  are merged. Thus, we take the following route to obtain the problem kernel. First, in  $O(n + m)$  time, we use a greedy strategy to compute a maximal matching in the undirected graph underlying the citation graph  $D$ . If this matching has size at least  $h^2$ , then we show that there is a partition achieving H-index  $h$ . Otherwise, we use the fact that the articles that do not participate in the matching do not cite each other to design further data reduction rules.

**Reduction Rule 1.** Let  $(D, G, W, h)$  be an instance of  $\text{H-INDEX MANIPULATION}(\mu)$  for  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$  such that  $G$  is a clique. Compute a maximal matching in  $D$  by iteratively putting an arc into the matching as long as possible. If the resulting matching has size at least  $h^2$ , then accept.

**Lemma 2.** *Reduction Rule 1 is correct and an application can be performed in  $O(n + m)$  time.*

**Proof.** Let  $M$  be a matching of size  $h^2$  in  $D$ . Let  $W'$  denote the set of  $h^2$  vertices that are the heads of the arcs in  $M$ . The articles in  $W'$  are cited by the tails of the respective arcs in  $M$ . Thus, we may assume  $W' \subseteq W$  since only citations to vertices in  $W$  are counted. Consider a partition  $\mathcal{P}$  of  $W$  that is obtained by partitioning  $W'$  into exactly  $h$  sets, each of size  $h$ , and not merging any other articles in  $W$ . Since  $M$  is a matching, there are, for each merged article  $P \in \mathcal{P}$ , at least  $h$  independent arcs from an article in  $V \setminus W'$  to an article in  $P$ . Thus,  $\mathcal{P}$  has H-index  $h$  with respect to sumCite and unionCite. Moreover, since the articles in  $V \setminus W'$  are not merged, there are thus  $h$  distinct unmerged articles that cite an article of  $P$ . Hence,  $\mathcal{P}$  has H-index  $h$  with respect to fusionCite.

To see the claim about the running time, observe that it suffices to iterate once over all edges, maintaining a label for each vertex that indicates whether it has an incident edge in the matching.  $\square$

Now assume that we have computed a maximal matching of size at most  $h^2$  in  $D$  in  $O(n + m)$  time. Then, the vertices incident to the matching arcs form a vertex cover  $C$  of size at most  $2h^2$  for  $D$ . It remains to upper bound the number of articles in the independent set  $V \setminus C$ . To this end, we first give a data reduction rule that ensures that each article in  $C$  cites only few articles in  $W \setminus C$ . To do this, we need the following lemma, which enables us to assume that a solution merges only few articles in the independent set.

**Lemma 3.** *Let  $(D, G, W, h)$  be a yes-instance of  $\text{H-INDEX MANIPULATION}(\mu)$ , where  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$ , such that  $G$  is a clique and let  $X \subseteq W$  such that no article in  $X$  cites any other article in  $X$ . Then, there is a partition  $\mathcal{P}$  of  $W$  that has H-index at least  $h$  with respect to  $\mu$  and such that at most  $h^2$  atomic articles from  $X$  are not singletons in  $\mathcal{P}$ .*

**Proof.** Consider a partition  $\mathcal{P}^*$  that has H-index at least  $h$  with respect to  $\mu$  and such that the number of atomic articles from  $X$  that are not singletons in  $\mathcal{P}^*$  is minimum. If  $\mathcal{P}^*$  has more than  $h$  merged articles or a merged article  $P$  with  $\mu(P) < h$ , then one of the merged articles can be split into its atomic articles without decreasing the H-index below  $h$ . Thus, assume that  $\mathcal{P}^*$  contains at most  $h$  merged articles  $P$ , each with  $\mu(P) \geq h$ . To prove the lemma, it is enough to show that there is no merged article  $P^* \in \mathcal{P}^*$  containing more than  $h$  atomic articles  $a_1, \dots, a_\ell \in X$ , where  $\ell > h$ . We will lead the existence of such an article  $P^*$  to a contradiction.

If  $\mu = \text{sumCite}$ , then it is clear that  $P^*$  can be replaced by two articles  $P^* \setminus \{a_i\}, \{a_i\}$  in  $\mathcal{P}^*$  for some  $i \in \{1, \dots, \ell\}$  so that still  $\text{sumCite}(\mathcal{P}^* \setminus \{a_i\}) \geq h$ , contradicting our choice of  $\mathcal{P}^*$ .

For  $\mu = \text{unionCite}$ , consider the articles  $P_j := (P^* \setminus X) \cup \{a_1, \dots, a_j\}$  for  $j \in \{1, \dots, \ell\}$ . If  $\text{unionCite}(P_j) < \text{unionCite}(P_{j+1})$  for each  $j \in \{1, \dots, \ell - 1\}$ , then  $\text{unionCite}(P_h) \geq h$  and replacing  $P^*$  by  $P_h$  in  $\mathcal{P}^*$  and adding the articles of  $P^* \setminus P_h$  as singletons yields a partition with more singletons from  $X$  and H-index at least  $h$ , thus contradicting our choice of  $\mathcal{P}^*$ .

Otherwise, there is a  $j \in \{1, \dots, \ell - 1\}$  such that  $\text{unionCite}(P_j) \geq \text{unionCite}(P_{j+1})$ , meaning that  $a_{j+1}$  is cited by a subset of the atomic articles that cite  $P_j$ . Then, replacing  $P^*$  by the two articles  $P^* \setminus \{a_{j+1}\}$  and  $\{a_{j+1}\}$  in  $\mathcal{P}^*$  yields a partition with H-index at least  $h$  and more singletons from  $X$ , contradicting our choice of  $\mathcal{P}^*$ .

Similarly, for  $\mu = \text{fusionCite}$ , consider the series of partitions  $\mathcal{P}_j$  for  $j \in \{1, \dots, \ell - 1\}$  arising from replacing  $P^*$  by the  $\ell - j + 1$  (merged) articles  $P_j$  and  $\{a_{j+1}\}, \dots, \{a_\ell\}$  in  $\mathcal{P}^*$ . If  $\text{fusionCite}_{\mathcal{P}_j}(P_j) < \text{fusionCite}_{\mathcal{P}_{j+1}}(P_{j+1})$  for each  $j \in \{1, \dots, \ell - 1\}$ , then  $\text{fusionCite}_{\mathcal{P}_h}(P_h) \geq h$  and each remaining merged article  $P' \in \mathcal{P}_h \cap \mathcal{P}^*$  has  $\text{fusionCite}_{\mathcal{P}_h}(P') \geq \text{fusionCite}_{\mathcal{P}}(P')$ . Yet  $\mathcal{P}_h$  has less atomic articles from  $X$  that are not singletons than  $\mathcal{P}^*$ , which contradicts our choice of  $\mathcal{P}^*$ . If, otherwise, there is a  $j \in \{1, \dots, \ell - 1\}$  such that  $\text{fusionCite}_{\mathcal{P}_j}(P_j) \geq \text{fusionCite}_{\mathcal{P}_{j+1}}(P_{j+1})$ , then we consider two cases.

In the first case,  $a_{j+1}$  does not cite any article in  $P^* \setminus X$ . Since articles in  $X$  do not cite each other and  $a_{j+1} \in X$ , it follows that  $a_{j+1}$  does not cite any article in  $P^*$ . Therefore,  $\text{fusionCite}_{\mathcal{P}_j}(P_j) \geq \text{fusionCite}_{\mathcal{P}_{j+1}}(P_{j+1})$  implies that the set of (merged) articles that cite  $a_{j+1}$  is a subset of the (merged) articles that cite  $P_j$ . This implies that the set of (merged) articles that cite  $a_{j+1}$  is a subset of the (merged) articles that cite  $P_\ell$ . Thus, replacing the merged article  $P^* = P_\ell$  by the two articles  $P^* \setminus \{a_j\}$  and  $\{a_j\}$  in  $\mathcal{P}$ , we obtain a partition with H-index at least  $h$  that has one less atomic article from  $X$  that is not a singleton. This contradicts our choice of  $\mathcal{P}^*$ .

In the second case,  $a_{j+1}$  cites at least one article in  $P^* \setminus X$ . Let  $c$  be the number of (merged) articles outside of  $P_j$  that cite  $a_{j+1}$  and none of the articles in  $P_j$  and let  $d$  be the number of (merged) articles outside of  $P_j$  that cite both  $a_{j+1}$  and an article in  $P_j$ . We have  $c \leq d + 1$  since, otherwise,  $\text{fusionCite}_{\mathcal{P}_j}(P_j) < \text{fusionCite}_{\mathcal{P}_{j+1}}(P_{j+1})$ . Since no article in  $X$  cites  $a_{j+1}$  and  $P_\ell \setminus P_j \subseteq X$ , we also have  $c' \leq d' + 1$ , where  $c'$  is the number of (merged) articles outside of  $P_\ell$  that cite  $a_{j+1}$  and none of the articles in  $P_\ell \setminus \{a_{j+1}\}$  and where  $d'$  is the number of (merged) articles outside of  $P_\ell$  that cite both  $a_{j+1}$  and an article in  $P_\ell \setminus \{a_{j+1}\}$ . Replacing the merged article  $P_\ell$  by the two articles  $P_\ell \setminus \{a_j\}$  and  $\{a_j\}$  in  $\mathcal{P}^*$ , we obtain a partition with H-index at least  $h$  that has one less atomic article from  $X$  that is not a singleton. This is a contradiction to the choice of  $\mathcal{P}^*$ .  $\square$

The idea for the following data reduction rule is that, if an article in a vertex cover  $C$  cites many articles outside of  $C$ , then only few of these are in merged articles and only few of the remaining articles are needed to maintain the citations of the merged articles. Hence, superfluous citations can be removed.

**Reduction Rule 2.** Let  $(D, G, W, h)$  be an instance of H-INDEX MANIPULATION( $\mu$ ), where  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$ , such that  $G$  is a clique and let  $C$  be a vertex cover in  $D$ . If there is an article  $v \in C$  that cites more than  $2h^2 + 2h$  articles in  $W \setminus C$ , then remove an arbitrary citation  $(v, w)$  for some  $w \in W \setminus C$ .

**Lemma 4.** Reduction Rule 2 is correct and can be exhaustively applied in  $O(n + m)$  time.

**Proof.** We first prove the correctness. Clearly, if the instance resulting from an application of Reduction Rule 2 is a yes-instance, then also the original instance is a yes-instance. For the converse, consider a partition  $\mathcal{P}$  with H-index  $h$  for  $(D, G, W, h)$  that does not have H-index  $h$  after removing  $(v, w)$  from  $D$ . Without loss of generality, we can assume that there are no merged articles  $P \in \mathcal{P}$  with  $\mu(P) < h$ , because unmerging such articles can only increase the number of citations of other articles (in the case of fusionCite) and, hence, cannot decrease the H-index of  $\mathcal{P}$ .

Observe that, after the deletion of citation  $(v, w)$  from  $D$ , there is at most one merged article  $P \in \mathcal{P}$  such that  $\mu(P) < h$ . We claim that, among the articles cited by  $v$ , there is an atomic article with less than  $h$  citations that we can add to  $P$  so that  $P$  has  $h$  citations again, thus yielding a partition with H-index at least  $h$ . Let  $U$  denote the set of more than  $2h^2 + 2h$  articles in  $W \setminus C$  that are cited by  $v$ . By Lemma 3, we may assume that at most  $h^2$  of the articles in  $U$  are in merged articles. Thus, since  $\mathcal{P}$  does not have H-index  $h$ , there is an article  $u \in U$  that is a singleton in  $\mathcal{P}$  and satisfies  $\mu(\{u\}) < h$ . If  $\mu = \text{sumCite}$ , then adding  $u$  to  $P$  yields a partition with H-index  $h$  because  $u$  has at least one citation (from  $v$ ). If  $\mu = \text{unionCite}$ , then observe that  $v$  does not cite  $P$ . Hence, adding  $u$  to  $P$  yields a partition with H-index  $h$  since  $v$  cites  $P \cup \{u\}$ . It remains to handle the case  $\mu = \text{fusionCite}$ .

Note that having  $\text{fusionCite}_{\mathcal{P}}(P) < h$  after deleting  $(v, w)$  from  $D$  means  $v \notin P$ . Recall that, by Lemma 3, we may assume that at most  $h^2$  of the articles in  $U$  are in merged articles. Furthermore, there are at most  $h - 1$  articles in  $U$  that cite  $P$  and at most  $h - 1$  articles  $u \in U$  with  $\text{fusionCite}_{\mathcal{P}}(u) \geq h$ . Denote the remaining articles of  $U$  by  $u_1, \dots, u_\ell$ . That is, each  $u_i$  is cited by  $v$ , is a singleton in  $\mathcal{P}$ , does not cite  $P$ , and has  $\text{fusionCite}_{\mathcal{P}}(u_i) < h$ . Observe that, if one of these articles, say  $u_i$ , does not cite any merged article in  $\mathcal{P}$ , then adding  $u_i$  to  $P$  yields a partition with H-index  $h$ . Hence, assume that each article  $u_i$  for  $i \in \{1, \dots, \ell\}$  cites at least one merged article. Observe furthermore that, if there is some  $u_i$  such that each merged article  $P' \neq P$  that is cited by  $u_i$  receives  $h$  citations from  $u_1, \dots, u_{i-1}$ , then adding  $u_i$  to  $P$  yields a partition with H-index  $h$  (recall that  $u_i$  does not cite  $P$ ). Call such an article  $u_i$  *good*. We claim that there is at least one good article. Assign to each  $u_i$  the integer  $c_i := \sum_{P' \in \mathcal{P} \setminus \{P\}} \min\{h, \text{cites}(i, P')\}$ , where  $\text{cites}(i, P')$  is the number of citations of  $P'$  from  $u_1, \dots, u_{i-1}$ . Observe that each  $u_i$  either cites at least one merged article  $P' \neq P$  that receives less than  $h + 1$  citations from  $u_1, \dots, u_{i-1}$  or it is good. Hence, either  $c_i > c_{i-1}$  or  $u_i$  is good. Furthermore, if  $c_i \geq (h - 1)h$ , then  $u_i$  is good. Thus, if  $\ell > (h - 1)h$ , then there is a good article. Because  $\ell \geq |U| - h^2 - (h - 1) - (h - 1)$  and  $|U| \geq 2h^2 + 2h$ , there is a good article indeed.

Regarding the running time, for each article  $v \in V$ , it can be checked in  $O(\deg_D^{\text{out}}(v))$  time whether the rule is applicable. In the same time, all citations exceeding the number  $2h^2 + 2h$  can be deleted. Since application of the reduction rule to one article cannot make it applicable to other articles, it follows that it can exhaustively be applied in  $O(\sum_{v \in V} \deg_D^{\text{out}}(v)) = O(n + m)$  time.  $\square$

Finally, we need the following cleanup rule.

**Reduction Rule 3.** Let  $(D, G, W, h)$  be an instance of H-INDEX MANIPULATION( $\mu$ ), where  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$ , such that  $G$  is a clique.

- If there are  $h$  articles in  $W$  with  $h$  citations each, then accept.
- If there is an article in  $W$  that is not cited, then remove this article.
- If there is an article in  $V \setminus W$  that cites no other article, then remove this article.
- If there is an article in  $W$  that is cited by more than  $h$  articles in  $V \setminus W$ , then remove an arbitrary one of these citations.

**Lemma 5.** *Reduction Rule 3 is correct and can be exhaustively applied in  $O(n + m)$  time.*

**Proof.** It is clear that Reduction Rule 3 can be exhaustively applied in  $O(n + m)$  time. For the correctness, the only non-obvious part is the last one. To see that it is correct, let  $v$  be an article to which it has been applied, and observe that every merged article that  $v$  can be contained in has at least  $h$  citations before applying the rule, as well as after applying the rule.  $\square$

Combining all data reduction rules above, we can give the promised polynomial-size problem kernel.

**Theorem 7.** *If the compatibility graph is a clique, then a  $(4h^4 + 6h^3 + 5h^2)$ -article problem kernel for H-INDEX MANIPULATION( $\mu$ ) is computable in  $O(n + m)$  time, where  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$ .*

**Proof.** To compute the problem kernel, apply exhaustively Reduction Rules 1 to 3. By the corresponding lemmas, the resulting instance is a yes-instance if and only if the input instance is a yes-instance, and the rules can be carried out in  $O(n + m)$  time.

To see the upper bound on the size, let  $C$  be the vertex cover of  $D$  computed from the matching of Reduction Rule 1. Note that  $|C| \leq 2h^2$ . We upper bound the size of  $W$  and, due to reducedness with respect to Reduction Rule 3, it then suffices to upper bound the number of articles that cite or are cited by articles in  $C$ . We divide the articles into four groups:

- The set  $W_{\geq} \subseteq W$  of articles with at least  $h$  citations from articles in  $V \setminus W$ ,
- the set  $W_{<} \subseteq W$  of articles with less than  $h$  citations from articles in  $V \setminus W$ ,
- the set  $V_{<} \subseteq V \setminus (W \cup C)$  of articles that cite articles in  $W_{<}$ , and
- the set  $V_{\geq} \subseteq V \setminus (W \cup C \cup V_{<})$  of articles that cite articles in  $W_{\geq}$  but no article in  $W_{<}$ .

Clearly,  $V = C \cup W_{\geq} \cup W_{<} \cup V_{<} \cup V_{\geq}$ . To upper bound the size of  $V$ , first note that  $|W_{\geq}| \leq h - 1$  by Reduction Rule 3. For  $W_{<}$ , note that each of these articles is either contained in  $C$  or cited by at least one article in  $C$ . By reducedness with respect to Reduction Rule 2, there are hence at most  $2h^2 + (2h^2 + 2h)|C| \leq 4h^4 + 4h^3 + 2h^2$  articles in  $W_{<}$ . Since each article in  $V_{<}$  cites at least one article in  $W_{<} \cap C$ , and since these articles receive at most  $h - 1$  such citations each, there are at most  $2h^3$  articles in  $V_{<}$ . Finally, each article in  $V_{\geq}$  cites at least one article in  $W_{\geq}$ , and these articles receive at most  $h$  citations from articles in  $V \setminus W$  each by Reduction Rule 3. Thus, there are at most  $h^2$  articles in  $V_{\geq}$  and, overall, there are at most

$$|C| + |W_{\geq}| + |W_{<}| + |V_{<}| + |V_{\geq}| \leq 2h^2 + (4h^4 + 4h^3 + 2h^2) + 2h^3 + h^2 = 4h^4 + 6h^3 + 5h^2$$

articles in a reduced instance.  $\square$

Applying an arbitrary algorithm that decides H-INDEX MANIPULATION to the instances resulting from the problem kernel in Theorem 7 now yields the following classification result.

**Corollary 2.** *If the compatibility graph is a clique, then H-INDEX MANIPULATION( $\mu$ ) with  $\mu \in \{\text{sumCite}, \text{unionCite}, \text{fusionCite}\}$  is linear-time solvable for constant  $h$ .*

## 5. Experiments

In this section, we examine by how much authors can increase their H-indices when allowing only merges of articles with similar titles or when fixing the allowed number of merges. To this end, we gathered article and citation data of AI



**Table 2**

Properties of the three data sets. Here,  $p$  is the number of profiles for each data set,  $|W|$  is the average number of atomic articles,  $\bar{c}$  is the average number of citations,  $\bar{h}$  is the average H-index in the data set, and  $h/a$  is the average (unmanipulated) H-Index increase per year; the ‘max’ subscript denotes the maximum of these values.

	$p$	$ W $	$ W _{\max}$	$\bar{c}$	$c_{\max}$	$\bar{h}$	$h_{\max}$	$h/a$
AI’s 10 To Watch 2011	5	170.2	234	1614.2	3725	34.8	46	2.53
AI’s 10 To Watch 2013	8	58.25	144	542.0	1646	14.0	26	2.77
IJCAI 2013	22	45.91	98	251.5	547	10.36	16	1.24

researchers, computed compatibility graphs based on similarity of article titles, and implemented heuristics and exact algorithms for maximizing the H-index. Herein, we focused on the islands of tractability that we determined in our theoretical analysis, that is, the cases of small number of merges and small connected components in the compatibility graph for the measures sumCite and unionCite. These cases are also practically relevant, as unionCite is the measure used by Google Scholar. The implemented algorithms are mainly based on [Theorems 1 and 4](#).

**Data acquisition** We crawled Google Scholar data of 22 selected authors of IJCAI’13. Our (biased) selection was based on capturing authors in their early career, for whom H-index manipulation would seem most attractive. Specifically, we selected authors who have a Google Scholar profile, an H-index between 8 and 20, between 100 and 1000 citations, who are active between 5 and 10 years, and do not have a professor position.

In addition, we crawled Google Scholar data of *AI’s 10 to Watch*, a list of young accomplished researchers in AI that is compiled every two years by *IEEE Intelligent Systems*. The dataset contains five profiles from the 2011 and eight profiles from the 2013 edition of the list [1,33]. Some profiles were omitted due to difficulties in the crawling process, for example, because of articles that could not be attributed unambiguously to the respective author due to non-unique author names. Compared to the IJCAI 2013 author set, AI’s 10 To Watch 2011 contains researchers who are more experienced and AI’s 10 To Watch 2013 falls in between these two data sets in this regard. [Table 2](#) gives an overview of the properties of the data sets.

For each author, we computed upper and lower bounds for the H-index increase when allowing at most  $k = 1, \dots, 12$  merges of arbitrary articles and the maximum possible H-index increase when merging only articles whose titles have a similarity above a certain compatibility threshold  $t = 0.1, 0.2, \dots, 0.9$ . The compatibility thresholding is described in more detail below.

**Generating compatibility graphs** Compatibility graphs are constructed using the following simplified ‘bag of words model’: Compute for each article  $u$  the set of words  $T(u)$  in its title. Draw an edge between articles  $u$  and  $v$  if  $|T(u) \cap T(v)| \geq t \cdot |T(u) \cup T(v)|$ , where  $t \in [0, 1]$  is the *compatibility threshold*. For  $t = 0$ , the compatibility graph is a clique. For  $t = 1$ , only articles with the same words in the title are adjacent. Inspection showed that, for  $t \leq 0.3$ , already very dissimilar articles are considered compatible.

**Implemented algorithms** We implemented our algorithms for the parameter “maximum connected component size  $c$  of the compatibility graph” ([Theorem 1](#)) and for the parameter  $k$  of allowed merges ([Theorem 4](#)). We ran both algorithms using both the sumCite and unionCite measures. Note that, when applied with the unionCite measure, the algorithm for [Theorem 4](#) does not necessarily compute the maximum possible H-index increase (cf. [Theorem 5](#)), but we note that it yields a lower bound. Moreover, running it with sumCite yields an upper bound for the maximum achievable H-Index with unionCite and thus, we obtain both a lower and upper bound on the achievable H-index with respect to unionCite using  $k$  merges.

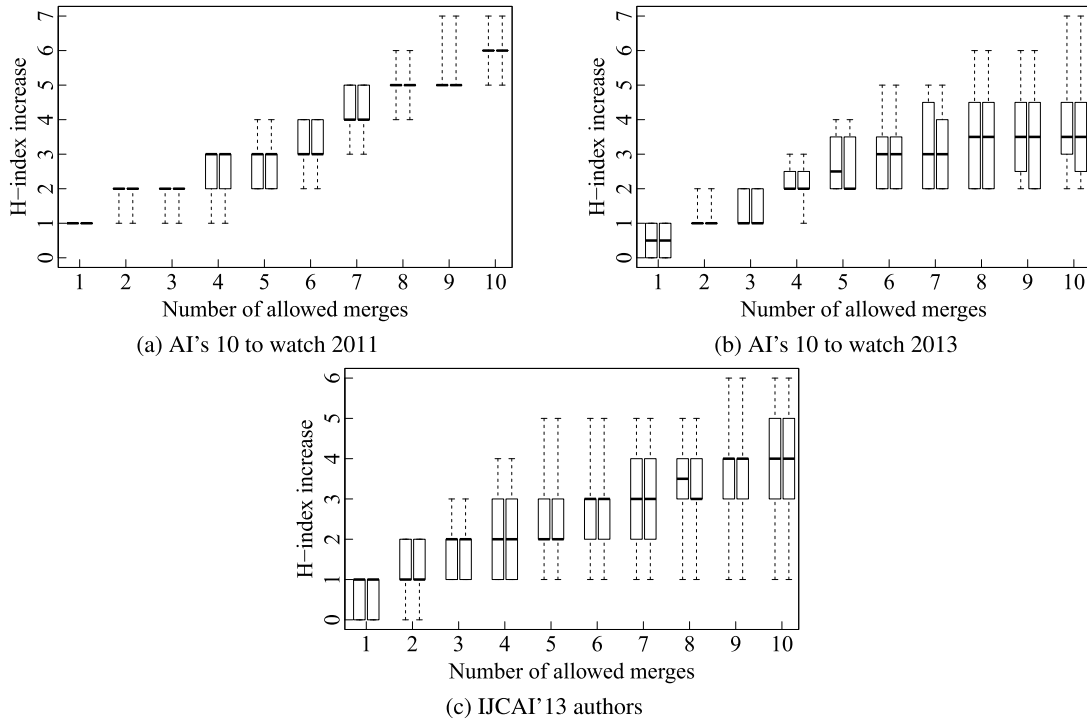
The fixed-parameter algorithm for parameter  $c$ , the size of the connected components of the compatibility graph, is not able to solve all instances. In particular, it fails for  $t = 0.2$ , where it runs out of memory in most cases. We thus implemented an alternative solution strategy that is based on the enumeration of cliques in the compatibility graph, exploiting the fact that any merged article is a clique in the compatibility graph  $G$ . Thus, a partition of the article set  $W$  that complies with  $G$  directly corresponds to a set of vertex-disjoint cliques in  $G$ .

Starting with  $h = 1$ , we do the following.

1. Enumerate all minimal sets  $P$  such that  $P$  is a clique in the compatibility graph and  $\mu(P) > h$ . Each set  $P$  is a *potential merged article* in a merged profile that achieves H-index  $h$ ; clearly, we can restrict attention to minimal sets.
2. Find a maximum-cardinality set  $\mathcal{P}'$  of potential merged articles such that  $P \cap P' = \emptyset$  for each pair  $P, P' \in \mathcal{P}'$ .
3. If  $|\mathcal{P}'| > h$ , then an H-index of at least  $h$  can be achieved via merging. Continue with  $h \leftarrow h + 1$ . Otherwise, an H-index of  $h$  cannot be achieved; return  $h - 1$  as the maximum H-index that can be achieved via merging.

In the implementation of Step 1, we first enumerate all maximal cliques of the compatibility graph and then check for each subset of each maximal clique whether it is a minimal set such that  $\mu(P) > h$ . In Step 2, the size of  $\mathcal{P}'$  is computed by constructing an auxiliary graph whose vertices are the potential merged articles and where edges are added between





**Fig. 3.** Box plots for our three data sets of the achievable H-index increases when the number of merges is restricted but arbitrary pairs of articles can be merged (that is, the compatibility graph is a clique). For each number  $k$  of allowed merges, the left box shows the H-index increase for sumCite, the right box shows lower bounds on the possible H-index increase for unionCite. The lower edge of a box is the 25th percentile and the upper edge is the 75th percentile, a thick bar is the median. The whiskers above and below each box extend to the maximum and minimum observed values.

potential merged articles that have nonempty intersection. In this graph,  $\mathcal{P}'$  is a maximum-cardinality independent set. We compute  $\mathcal{P}'$  by computing a minimum-cardinality vertex cover via a simple fixed-parameter algorithm for the parameter vertex cover size.

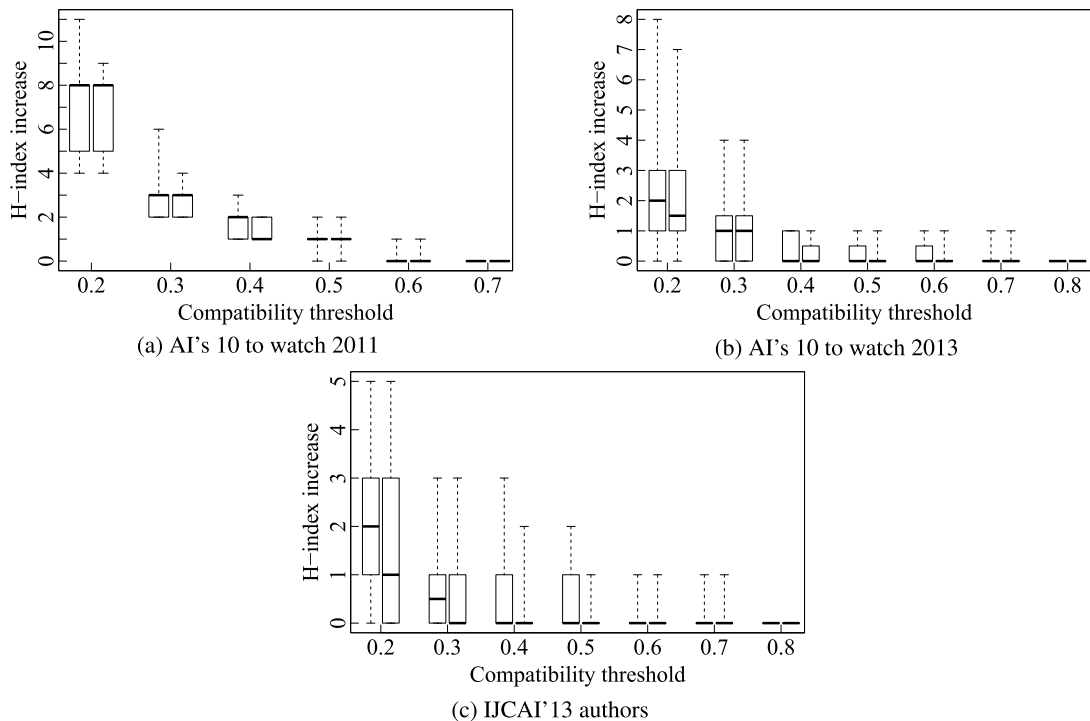
This algorithm has a higher worst-case running time than the fixed-parameter algorithm for parameter  $c$ : the overall number of potential merged articles  $p$  may be exponential in  $c$  and we solve INDEPENDENT SET on a graph of order  $p$ . Nevertheless, it works for the three data sets as the number of potential merged articles is much lower than in a worst-case instance.

Source code and data are freely available at <http://ftp.akt.tu-berlin.de/hindex>.

**Experimental results** We ran our algorithms under a time limit of one hour on a 3.6 GHz Intel Xeon E5-1620 processor and a memory limit of 64GB. Under these limits, the fixed-parameter algorithm for parameter  $k$ , the number of allowed merges, failed to solve instances with  $k \geq 11$ . Thus, Fig. 3 shows results for  $k \leq 10$  only. The fixed-parameter algorithm for parameter  $c$ , the size of the connected components of the compatibility graph, failed to solve instances with a compatibility threshold  $t \leq 0.2$ . Instances with  $k \leq 10$  and  $t \geq 0.3$  were usually solved within few seconds and using at most 100 MB of memory. The algorithm based on clique enumeration solved each instance with  $t \geq 0.2$  in several minutes. None of our algorithms was able to solve all instances with  $t = 0.1$ . Thus, Fig. 4 shows results for  $t \geq 0.2$  only.

Fig. 3 shows the H-index increase over all authors for each number  $k = 1, \dots, 10$  of allowed article merges when the compatibility graph is a clique. Remarkably, three merges are sufficient for all of our sample authors to increase their H-index by at least one. Let us put this number into perspective: as shown in Table 2, we measured that, without manipulation, on average the H-index in each group of our sample authors grows between 1.24 and 2.77 per year (which is higher than the one-per-year increase observed by Hirsch [22] in physics). Thus, from Fig. 3, one can conclude that two merges could save about nine months of work for half of our AI's 10 To Watch 2011 group, about four months of work for half of our AI's 10 To Watch 2013 group, and 19 months of work for half of our IJCAI'13 group.

Fig. 4 shows the H-index increase over all authors for unionCite and each compatibility threshold  $t = 0.2, 0.3, \dots, 0.9$ . Remarkably, when using a compatibility threshold  $t \geq 0.6$ , 75% of our sample authors cannot increase their H-index by merging compatible articles. We conclude that increasing the H-index substantially by article merges should be easy to discover since it is necessary to merge articles with highly dissimilar titles for such a manipulation.



**Fig. 4.** Box plots for our three data sets of the achievable H-index increases when compatibility of articles is restricted but an arbitrary number of articles can be merged. For each compatibility threshold  $t$ , the left box shows the H-index increase for sumCite, the right box for unionCite. The lower edge of a box is the 25th percentile and the upper edge is the 75th percentile, a thick bar is the median. The whiskers above and below each box extend to the maximum and minimum observed values.

## 6. Outlook

Clearly, it is interesting to consider merging articles in order to increase other measures than the H-index, like the  $g$ -index [14,31], the  $w$ -index [30], or the  $i10$ -index of a certain author. The  $i10$ -index, the number of articles with at least ten citations, is also currently used by Google Scholar. Elkind and Pavlou [27] performed a study in this direction and, among other results, showed that the  $g$ -index and the  $i10$ -index seem somewhat easier to manipulate than the H-index. In addition, they also studied a scenario where the manipulator wants to take into account the impact of the manipulation actions on other researchers (distinguishing between friends and competitors).

Moreover, merging articles in order to increase one index might decrease other indices, like the overall number of citations. Hence, it is also interesting to study the problem of increasing the H-index by merging without decreasing the overall number of citations or the  $i10$ -index below a predefined threshold. A systematic study of computing Pareto optimal solutions could also be interesting.

The computational problems related to optimal merging of articles in the different measures are quite natural as evidenced for example by their relation to BIN COVERING and MACHINE COVERING. Thus, improvements over the presented algorithms would be desirable as well as a study of further parameterizations in a broad multivariate complexity analysis [18,23,26].

Altogether, our experiments show that the merging option leaves some room for manipulation but that *substantial* manipulation requires merging visibly unrelated articles. Hiring committees that use the H-index in their evaluation thus should either examine the article merges more closely or rely on databases that do not allow article merges.

## Acknowledgements

We thank the anonymous referees for their helpful comments, in particular for pointing out Theorem 2(iii).

René van Bevern was supported by the German Research Foundation (DFG), project DAPA (NI 369/12), at TU Berlin, and by the Russian Foundation for Basic Research (RFBR), project 16-31-60007 mol\_a-dk, at Novosibirsk State University. Christian Komusiewicz was supported by the DFG project MAGZ (KO 3669/4-1) and Manuel Sorge was supported by the DFG project DAPA (NI 369/12). Toby Walsh was supported by the Alexander von Humboldt Foundation, Bonn, Germany, while at TU Berlin. The main work was done while Toby Walsh was affiliated with University of New South Wales and Data61, Sydney, Australia.

## References

- [1] AI's 10 to watch, *IEEE Intell. Syst.* 26 (1) (2011) 5–15, <http://dx.doi.org/10.1109/MIS.2011.7>.
- [2] S. Assmann, D. Johnson, D. Kleitman, J.-T. Leung, On a dual version of the one-dimensional bin packing problem, *J. Algorithms* 5 (4) (1984) 502–525, [http://dx.doi.org/10.1016/0196-6774\(84\)90004-X](http://dx.doi.org/10.1016/0196-6774(84)90004-X).
- [3] C. Bartneck, S. Kokkermans, Detecting *h*-index manipulation through self-citation analysis, *Scientometrics* 87 (1) (2011) 85–98, <http://dx.doi.org/10.1007/s11192-010-0306-5>.
- [4] R. van Bevern, C. Komusiewicz, R. Niedermeier, M. Sorge, T. Walsh, *H*-index manipulation by merging articles: models, theory, and experiments, in: *Proceedings of the 24th Conference on Artificial Intelligence, IJCAI'15*, AAAI Press, 2015, pp. 808–814.
- [5] R. van Bevern, C. Komusiewicz, H. Molter, R. Niedermeier, M. Sorge, T. Walsh, *H*-index manipulation by undoing merges, in: *Proceedings of the 22nd European Conference on Artificial Intelligence, ECAI'16*, in: *Frontiers in Artificial Intelligence and Applications*, IOS Press, 2016, in press.
- [6] H.L. Bodlaender, M. van Kreveld, Google scholar makes it hard—the complexity of organizing one's publications, *Inf. Process. Lett.* 115 (12) (2015) 965–968, <http://dx.doi.org/10.1016/j.ipl.2015.07.003>.
- [7] E.G. Coffman Jr., J. Csirik, G. Galambos, S. Martello, D. Vigo, Bin packing approximation algorithms: survey and classification, in: *Handbook of Combinatorial Optimization*, Springer, New York, 2013, pp. 455–531.
- [8] M. Cygan, F.V. Fomin, L. Kowalik, D. Lokshtanov, D. Marx, M. Pilipczuk, M. Pilipczuk, S. Saurabh, *Parameterized Algorithms*, Springer, 2015.
- [9] B. de Keijzer, K.R. Apt, The *H*-index can be easily manipulated, *Bull. EATCS* 110 (2013) 79–85.
- [10] E. Delgado López-Cózar, N. Robinson-García, D. Torres-Salinas, The Google Scholar experiment: how to index false papers and manipulate bibliometric indicators, *J. Assoc. Inf. Sci. Technol.* 65 (3) (2014) 446–454, <http://dx.doi.org/10.1002/asi.23056>.
- [11] R. Denman, S. Foster, Using clausal graphs to determine the computational complexity of *k*-bounded positive one-in-three SAT, *Discrete Appl. Math.* 157 (7) (2009) 1655–1659, <http://dx.doi.org/10.1016/j.dam.2008.09.011>.
- [12] R.G. Downey, M.R. Fellows, *Fundamentals of Parameterized Complexity*, Springer, 2013.
- [13] B. Dutta, M.O. Jackson, M.L. Breton, Strategic candidacy and voting procedures, *Econometrica* 69 (4) (2001) 1013–1037, <http://dx.doi.org/10.1111/1468-0262.00228>.
- [14] L. Egghe, Theory and practise of the *g*-index, *Scientometrics* 69 (1) (2006) 131–152, <http://dx.doi.org/10.1007/s11192-006-0144-7>.
- [15] P. Faliszewski, A.D. Procaccia, AI's war on manipulation: are we winning?, *AI Mag.* 31 (4) (2010) 53–64.
- [16] P. Faliszewski, E. Hemaspaandra, L.A. Hemaspaandra, Using complexity to protect elections, *Commun. ACM* 53 (11) (2010) 74–82, <http://dx.doi.org/10.1145/1839676.1839696>.
- [17] M.R. Fellows, D. Hermelin, F. Rosamond, S. Viallette, On the parameterized complexity of multiple-interval graph problems, *Theor. Comput. Sci.* 410 (1) (2009) 53–61, <http://dx.doi.org/10.1016/j.tcs.2008.09.065>.
- [18] M.R. Fellows, B.M.P. Jansen, F.A. Rosamond, Towards fully multivariate algorithmics: parameter ecology and the deconstruction of computational complexity, *Eur. J. Comb.* 34 (3) (2013) 541–566, <http://dx.doi.org/10.1016/j.ejc.2012.04.008>.
- [19] J. Flum, M. Grohe, *Parameterized Complexity Theory*, Springer, 2006.
- [20] D.K. Friesen, B.L. Deuermeier, Analysis of greedy solutions for a replacement part sequencing problem, *Math. Oper. Res.* 6 (1) (1981) 74–87, <http://dx.doi.org/10.1287/moor.6.1.74>.
- [21] J. Guo, R. Niedermeier, Invitation to data reduction and problem kernelization, *ACM SIGACT News* 38 (1) (2007) 31–45, <http://dx.doi.org/10.1145/1233481.1233493>.
- [22] J.E. Hirsch, An index to quantify an individual's scientific research output, *Proc. Natl. Acad. Sci. USA* 102 (46) (2005) 16569–16572, <http://dx.doi.org/10.1073/pnas.0507655102>.
- [23] C. Komusiewicz, R. Niedermeier, New races in parameterized algorithmics, in: *Proceedings of the 37th International Symposium on Mathematical Foundations of Computer Science, MFCS '12*, in: *Lecture Notes in Computer Science*, vol. 7464, Springer, 2012, pp. 19–30.
- [24] S. Kratsch, Recent developments in kernelization: a survey, *Bull. EATCS* 113 (2014) 58–97.
- [25] R. Niedermeier, *Invitation to Fixed-Parameter Algorithms*, Oxford Lecture Series in Mathematics and Its Applications, vol. 31, Oxford University Press, 2006.
- [26] R. Niedermeier, Reflections on multivariate algorithmics and problem parameterization, in: *Proceedings of the 27th International Symposium on Theoretical Aspects of Computer Science, STACS '10*, in: *LIPIcs*, vol. 5, Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2010, pp. 17–32.
- [27] C. Pavlou, E. Elkind, Manipulating citation indices in a social context, in: *Proceedings of the 2016 International Conference on Autonomous Agents and Multiagent Systems, AAMAS'16*, International Foundation for Autonomous Agents and Multiagent Systems, 2016, pp. 32–40.
- [28] P. Vinkler, Would it be possible to increase the Hirsch-index,  $\pi$ -index or CDS-index by increasing the number of publications or citations only by unity?, *J. Informetr.* 7 (1) (2013) 72–83, <http://dx.doi.org/10.1016/j.joi.2012.08.001>.
- [29] R. Walter, M. Wirth, A. Lawrinenko, Improved approaches to the exact solution of the machine covering problem, *J. Sched.* (2016), <http://dx.doi.org/10.1007/s10951-016-0477-x>, in press.
- [30] G.J. Woeginger, An axiomatic characterization of the Hirsch-index, *Math. Soc. Sci.* 56 (2) (2008) 224–232, <http://dx.doi.org/10.1016/j.mathsocsci.2008.03.001>.
- [31] G.J. Woeginger, An axiomatic analysis of Egghe's *g*-index, *J. Informetr.* 2 (4) (2008) 364–368, <http://dx.doi.org/10.1016/j.joi.2008.05.002>.
- [32] A. Yong, Critique of Hirsch's citation index: a combinatorial Fermi problem, *Not. AMS* 61 (9) (2014) 1040–1050, <http://dx.doi.org/10.1090/noti1164>.
- [33] D. Zeng, AI's 10 to watch, *IEEE Intell. Syst.* 28 (3) (2013) 86–96, <http://dx.doi.org/10.1109/MIS.2013.57>.