

DOI:10.1145/3376901

Identifying female CS scientists by combining a robust bibliographic database and name filtering tools.

BY SANDRA MATTAUCH, KATJA LOHMANN, FRANK HANNIG, DANIEL LOHMANN, AND JÜRGEN TEICH

A Bibliometric Approach for Detecting the Gender Gap in Computer Science

WOMEN ARE UNDERREPRESENTED in the fields of science, technology, engineering, and mathematics (STEM) in most countries, including Germany and the U.S.^{29,32} This was demonstrated in several surveys investigating the proportion of women in the STEM fields for specific populations. Some of these studies, for example, investigated the number of enrolled students^{10,30} or the percentage of female professors at universities. Other studies analyzed the disparities in research funding.²³ Nearly all these surveys selected a particular population of women in consideration of their university degree

or their nationality.^{11,34} Like many other studies investigating the gender gap and its reasons in science, these surveys are usually based on data records from several kinds of registrations or enrollments, for example, the enrollment as student or doctoral student, the registration of finished doctoral theses or the membership as professor in a certain country.^{1,14,16,28} However, researchers at the postdoctoral level or industrial researchers are often not registered and unfortunately drop out of the surveys.

Bibliometric approaches are widely used to detect the gender gap and to determine possible reasons for it,^{4,12,15,33} for example, the research performance or collaboration behavior^{1,2,4,18} or different cognitive or sociocultural determinants.^{9,13,16} In this study, we use a method to detect the gender gap in the group of scientifically active researchers regardless of the limitations mentioned and focused to a certain scientific field. The group of interest comprises scientists that are currently active in doing research and publishing their findings—regardless of their university degree, nationality, gender, age, or origin and irrespective of their employment level in university or industry. As a case study, we measured the gender gap in the scientific field of the Transregional Research Centre 89 Invasive Computing (CRC/Transregio 89),^a which investigates a novel paradigm for the design and programming of future parallel computing systems and covers research from diverse domains of computer science and

a <http://www.invasic.de>

» key insights

- The bibliometric approach allows to estimate the proportion of scientifically active women in CS, regardless of their degree, employment level, nationality, age, or origin.
- The percentage of women contributing to 19 representative conferences in CS within the last six years is, on average, below 10%.
- The percentage of women shows only small variations over individual years and conferences.

Table 1. Selected conferences.

Conference name and abbreviation

- International Conference on Applied Cryptography and Network Security (ACNS)
- International Conference on Architecture of Computing Systems (ARCS)
- International Conference on Application-specific Systems, Architectures and Processors (ASAP)
- Asia and South Pacific Design Automation Conference (ASP-DAC)
- International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)
- International Conference on Compilers, Architectures, and Synthesis for Embedded Systems (CASES)
- International Conference on Compiler Construction (CC)
- International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)
- Design Automation Conference (DAC)
- Design, Automation and Test in Europe (DATE)
- International European Conference on Parallel and Distributed Computing (Euro-Par)
- European Conference on Computer Systems (EuroSys)
- International Conference on Parallel Computing (ParCo)
- Symposium on Operating Systems Principles (SOSP)
- USENIX Annual Technical Conference (USENIX)
- International Conference on Virtual Execution Environments (VEE)
- Conference on Design and Architectures for Signal and Image Processing (DASIP)
- International Conference on Humanoid Robots (Humanoids)
- Network and Distributed System Security Symposium (NDSS)

Table 2. List of Excluded Classes.

Onomastic Classes

- Hong Kong
- China
- Taiwan
- Republic of Korea
- Viet Nam
- Democratic People's Republic of Korea

electrical engineering, such as computer engineering, operating systems, programming languages, security, robotics, and high-performance computing. To ensure only scientifically active scientists are taken into account, we decided to collect data of researchers that successfully published their results in proceedings of international conferences within the last six years.

Conferences and the appropriate conference proceedings are the common publication medium in computer science and have a much higher impact than journal papers. For this purpose, and for working with representative and high-quality data, we used the DBLP Computer Science Bibliography,⁸ which lists the major computer science journals and conference proceedings, as our database. Table 1 presents a summary and selection of the 19 most relevant conferences for different disciplines of our CRC/Transregio 89. Based on this selection, we developed a Perl script extracting the author names by the given

constraints (conference name and a period of six years). Based on the filtered results, we subsequently determined the country of origin and the gender of each author by NamSor Applied Onomastics.²⁰ We finally verified this approach by random sampling and manual classification of the sampled names. The extracted information was then used to detect the gender gap in the field of the CRC/Transregio 89 Invasive Computing.

Methods

Extraction of author names from the DBLP Computer Science Bibliography.

To gather the original population of all scientifically active researchers within the scientific field described, we extracted the names of authors contributing to most relevant conferences (Table 1) within the last six years from the DBLP Computer Science Bibliography.⁸

The DBLP Computer Science Bibliography provides bibliographic information on all major computer science journals and proceedings. This open-data service indexes more than 4 million articles, published by more than 2.1 million authors.⁸

To pull the author names from the DBLP database, we created a Perl script: This script, which is publicly available under MIT license^b extracts all author names—regardless of the

^b <https://github.com/luhsra/venueauthor>

order of authors—for all papers published at a certain conference. The conference is defined by the input variables *venue* and *year*. The *venues* are the acronyms of the conferences as listed in Table 1. For *years*, we chose the 2012 to 2017. The script displays a list with the authors' first and last name, and the conference name and year. The resulting population comprises of 18,116 authors. Some 242 authors used abbreviations instead of first names, so these names were excluded from the analyses, resulting in an original population of 17,874 names.

Data handling. The extracted author names from the DBLP database were subsequently classified by NamSor Applied Onomastics, a name recognition software provided by a private start-up company.²⁰ The specialized data mining software also recognizes the linguistic or cultural origin of each personal name in any alphabet/language and allocates an onomastic class and the gender to each author name. The innovative machine learning algorithm provides unmatched accuracy at a fine-grained level, with flexibility and integration capability, to filter through large databases and extract names. It recognizes which language or culture stands behind a given name.²⁰ It is already known that the cultural context and origin are important for the determination of gender by name. Therefore, some names cannot be clearly defined without the origin. The name *Andrea*, for example, is a male name in Italy, but a female name in Spain. Some more examples are Jean, Joan, Laurence, Sascha, and Maria. To ensure a high degree of accuracy in the classification of the author names and to take the cultural context and the origin into account, we decided to use NamSor Origin API first, followed by NamSor Gender API.

Determination of the likely country of origin of a name by NamSor Origin API. NamSor Origin API allows determining the likely country of origin of each author, based on the sociolinguistics of the name (language, culture). The anthroponomical classification can be summarized as follows: Judging from the name only and the publicly available list of all 150k Olympic athletes since 1896 (and other similar lists of names), for which national team would

the person most likely run? Here, the U.S., Australia, among others are typically considered as a melting pot of other cultural origins (Ireland, Germany, among others) and not as an onomastic class on its own.^{25,27}

Based on the NamSor Origin API algorithm, the basic population of 17,874 authors was classified into 71 onomastic classes. The 20 proportionally largest classes represent 82.3% of the basic population. 16 onomastic classes have less than 20 authors listed and represent together under 1% of the basic population. The classification of cultural and geographical provenience of the author names by the NamSor Origin API algorithm shows that our data set is reasonably diverse and shows an acceptable variability with respect to the origin.

Determination of the likely gender of a name by using NamSor Gender API. For this task, we used the NamSor Gender API. The software predicts the gender of a personal name on a -1 (male) to +1 (female) scale and covers the U.S., European, Indian, African, Chinese, Hebrew, Russian/Slavic/Cyrillic, and Arabic names. In this step, the software combines two algorithms to maximize accuracy. First, a unique global name sociolinguistics algorithm that recognizes the origin of the couple first name and last name and infers whether the name sounds male or female in that particular culture. Second, a query in a massive database (800,000 names), which contains statistical information about baby names in each country of the world.¹⁹ Nevertheless, NamSor recommends passing additional geography/local context to the names to improve the accuracy of classification.¹⁹ The reliability of this method was already investigated in several publications.^{6,25-27,31}

Figure 1 reveals that 67.7% of the author names are classified as male and only a small proportion of 9.9% are classified as female names. Some 22.4% of the names in the basic population are unclassified (scale 0). These not classified names mainly have two reasons: names like Kerry, Jean, or Maria that are not strongly correlated to gender, and the structure and usage of Asian names.

Removal of Asian names. In most countries and cultures, the method of onomastics is very accurate, with a

precision in the range of 95%–99%—but we should pay attention to the structure of Asian names. The used Perl script generates a list of authors with first name and family name. In Asia, the family name comes first, followed by the first name. Although there are currently over 4,000 Chinese surnames, only 100 surnames still make up over 85% of China’s 1.3 billion citizens. In fact, just the top three Wang, Li, and Zhang cover more than 20% of the population.²² The situation is aggravated by the fact that a lot of Chinese names are not strongly correlated with gender. Moreover, if they were transliterated in Latin characters, even more information gets lost. The automatic determination of gender from Asian names with sufficient accuracy is not within the bounds of possibility of this work.³⁵ The analysis shows that 96.3% of the unclassified names come from these six onomastic classes. For these reasons, we decided to exclude all these Asian names from the onomastic classes listed in Table

2. Removal of these names reduces the population by 4,773 to 13,101 names. After the removal of Asian names, 149 unclassified names are remaining.

In Figure 2, the distribution of male, female, and unclassified authors after the removal of Asian names is shown. The percentage of female authors increases slightly to 11.3%, but the number of unclassified names has been reduced to 1.1%. The number of male authors has increased accordingly to 87.5%.

Validation of name sorting. After applying the procedure described earlier, we ended up with a population of 13,101 names (basic population): 1,486 names were classified as female names, 11,466 as male names. To test whether the names classified as female names really belong to women and—vice versa—those classified as male names really belong to men, we randomly selected samples from the basic population of men and women. The minimal sample sizes n of women and men is calculated using the following formula:

Figure 1. Distribution of female, male, and unclassified names as assorted by NamSor Gender API in the original population.

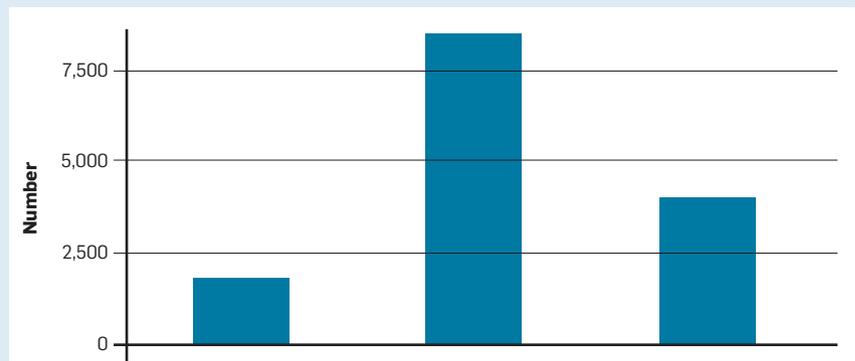
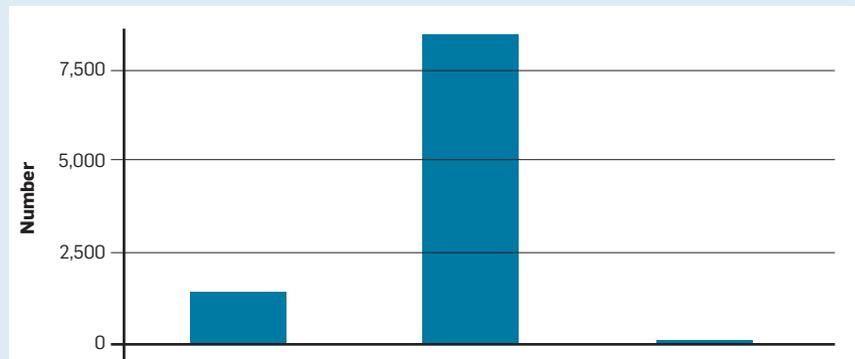


Figure 2. Distribution of female, male, and unclassified names as assorted by NamSor Gender API in the population when disregarding Asian names.



$$n \geq \frac{N}{1 + \frac{(N_1) \cdot e^2}{z^2 \cdot P \cdot (1 - P)}}$$

In Equation (1), N is the number of elements in the stock population, e the margin of error (5%), z is the z-score (1.96 for a confidence level of 95%), and P the prior judgment of the correct distribution (0.5, no prior judgment).

This gives us a sample size of 306 for the group of female names and 372 for the group of male names. The gender of scientists from these sample groups was manually verified by searching them on the Internet—assuming scientifically active persons to have an Internet presence. We determined the gender of the scientists by photos and the usage of gender-specific keywords (he, she, him, her, among others) on the personal homepages, on platforms like LinkedIn¹⁷ or ResearchGate²⁴ or pages referring to the scientist, for instance, as authors.

The results are shown in Figure 3. The estimation of the likely gender of a person by “NamSor Gender API” works quite well for male scientist but noticeably not as good for the group of female scientists: In the group of men, 84% were correctly verified to be male, only 0.3% were female, and 15.6% could not be verified due to no Internet presence. In the group of women, only 70% were correctly verified as female, yet 17.6% were male and 12% could not be found on the Internet.

In addition to the determination of the likely gender on the basis of the country of origin, we evaluated the gender classification accuracy when alternatively using the affiliation country extracted from the Scopus database. Scopus, Elsevier’s abstract and citation database generates precise citation search results and automatically updates researcher and institution profiles, unlike the DBLP database. The Python script we developed to extract the affiliation country of authors is publicly available under MIT license.^c To compare the classification accuracy of both approaches, the hand-verified set of 314 male and 269 female names serves as input. The percentage of true classifications for the first method, for example, is obtained as the number of correctly classified samples (528) in re-

lation to the total of 583 samples. We could show that there is a difference of less than 1% in classification accuracy when either using the country of origin or the affiliation country as input for the NamSor Gender API (see Figure 3).

Based on these random experiments, we decided to correct the automatically determined number of female and male authors accordingly using the following term:

$$F_{corr} = F \cdot corr_{ff} + M \cdot corr_{fm} \quad (2)$$

$$M_{corr} = M \cdot corr_{mm} + F \cdot corr_{mf} \quad (3)$$

In Eqs. (2) and (3), F_{corr} and M_{corr} denote the corrected numbers of women and men, F and M are the original values obtained from the name-sorting procedure, and $corr_x$ are the correction factors estimated from the results of the verification of name sorting:

$$corr_{ff} = \text{females in female group} = 0.70$$

$$corr_{mm} = \text{males in male group} = 0.84$$

$$corr_{fm} = \text{females in male group} = 0.003$$

$$corr_{mf} = \text{males in female group} = 0.17$$

The results shown next present corrected percentages of female and male researchers using Eqs. (2) and (3).

Case Study

For the 19 representative computer science conferences selected for our analysis as shown in Table 1, we extracted from the DBLP Computer Science Bibliography a total of 18,116 names of authors contributing to these conferences within the last six years and removed 242 authors that used initials instead of the full first names (original population). The names were then classified by origin and gender using the NamSor Applied Onomastics. From the original population, 4,773 author names assigned to Hong Kong, China, Taiwan, Republic of Korea, Viet Nam, and the Democratic People’s Republic of Korea were removed due to the infeasibility of automatic classification. A small number of 149 names (0.8%) were left unclassified for unknown reasons.

After applying the presented stochastic sampling of this population and subsequently applying the correction according to Eqs. (2) and (3) on the resulting basic population of 13,101 names, we could finally estimate that the percentage of women contributing

to the 19 conferences within the last six years is, on average, below 10% (as illustrated in Figure 4). On a per year basis, the percentage of female authors shows only small variations between 8.68% in 2012 and 10.1% in 2016.

Our approach now allows us to have a closer look at the proportion of scientifically active women in different individual conferences, and thus areas of computer science and not only to calculate the overall proportion of women in computer science as a whole. To illustrate the percentage of female authors in individual conferences, we picked out three of them: The International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS), the Design, Automation and Test in Europe (DATE) and the International Conference on Compiler Construction (CC). The percentage of female authors varies here between 6.2% for the CC and 11.7% for the CODES+ISSS conference. For the DATE conference, the percentage of female authors amounts to an average value of 9.6%.

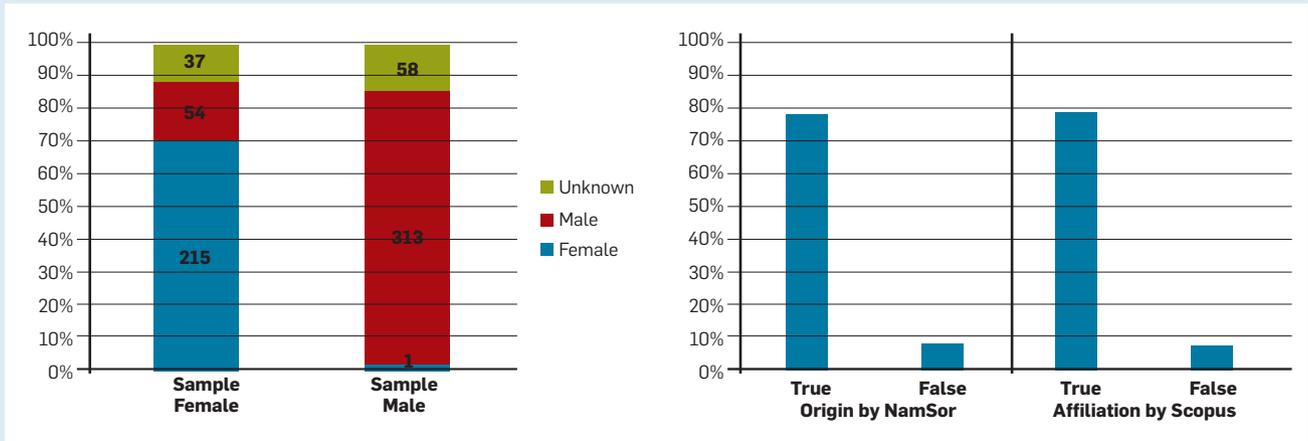
A closer look at the participation of women in all 19 conferences finally reveals a nearly symmetrical distribution. Five of the investigated conferences have a percentage of female authors above 10%, and five conferences have a proportion of female authors below 8.0% (see Table 3).

Table 3. Percentage of female authors in the examined conferences.

Conference	Percentage of female authors
CODES+ISSS	11.73
ACNS	11.51
Humanoids	10.98
DAC	10.41
CASES	10.26
ASP-DAC	9.94
Euro-Par	9.80
VEE	9.72
DATE	9.63
NDSS	9.61
DASIP	9.39
PARCO	8.68
EuroSys	8.48
USENIX	8.32
ASPLOS	7.96
SOSP	7.00
ARCS	6.72
ASAP	6.70
CC	6.15

c <https://github.com/luhsra/venueauthor>

Figure 3. Results of manual verification of gender classification using samples of names classified as female, respectively male (left) and comparison of accuracy when either using country of origin determined by NamSor or affiliation country extracted from Scopus as an alternative (right).



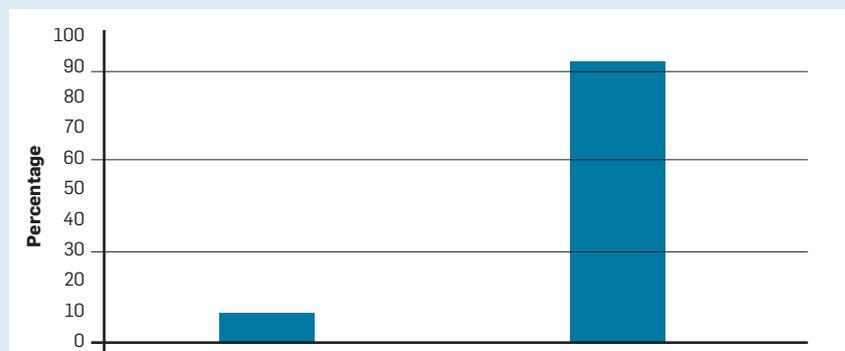
Discussion

In this work, we used a bibliometric approach to estimate the proportion of scientifically active women in the specific scientific field of computer science. In contrast to previous studies in the STEM fields that refer to limited data records, our method provides a more general approach with reduced limitations:

- We make sure to take all authors with publication activity in the last six years into account—independent of their university degree. Along with professors and postdoctoral, and industrial researchers, the examined group includes senior lecturers, doctoral students, career changers, and even employees without an academic degree like technicians or qualified IT specialists. Our approach allows us to exclude researchers that are not scientifically active anymore, for example, due to a change in their scientific field or job. Also, researchers active in administration or management are omitted, as well as students at the beginning of their studies. We cannot exclude that our results are partly influenced by an imbalance in research activity between female and male researchers, as has been shown for other scientific areas.^{5,15,21,28} However, since we consider the number of authors and not the number of publications, we assume this influence is relatively small.

- We generate our population independent of the origin of the authors. On the selected international conferences, one can find conference delegates from all over the world. As expected, we found author names from 71 different on-

Figure 4. Final distribution of female and male names for 19 conferences in computer science and electrical engineering after removal of Asian and unclassified names, and correction using stochastic samples and applying Eqs. (2) and (3).



mastic classes on our list, reflecting the likely country of origin of the authors. Our approach also provides the possibility to generate a population of authors only for national conferences or for individual conferences.

Compared to many previous studies searching for female scientists in computer science, our approach makes it possible to focus the analysis to a single conference further, a set of conferences representative for a specific scientific field, or to limit the data to a certain period of time. For the case study presented here, we examined representative conferences suggested by the researchers of the CRC/Transregio 89 Invasive Computing, which covers computer engineering, operating systems, programming languages, security, and the field of application including robotics and high-performance computing. By the selection of conferences, it would be pos-

sible to investigate other scientific fields or to limit further the scientific area (for example, to operating systems or computer security).

Despite these advantages of the method, we are not able to directly extract the gender or origin of the authors from the DBLP Computer Science Bibliography, one reason being that DBLP does not list these properties. By applying NamSor Applied Onomastics, we were able to determine the gender of the authors automatically. Yet, after testing the accuracy of this fully automatic classification on random samples from the group of men and women, we found out that although only one man was wrongly classified, 17.6% of those classified as women were in fact men. A more thorough inspection indicated that 24.1% of wrongly classified women were from India. These differences in accuracy between men and women through verification by random

sampling is not explained by NamSor Gender API. Indeed, they do not provide any information about the classification of names from India. To take the wrong classifications into account, we determined corrective factors.

The most significant disadvantage and a potential source of error of our approach is the removal of names classified as Asian names. The excluded group comprises a total of 4,773 names, which amounts to 26.7% of all names in the original population obtained from the DBLP Computer Science Bibliography. The removal of these names may distort the results. However, there is no evidence so far that the proportion of women in the group of removed Asian names is significantly higher than in the investigated group. In fact, several studies on women in the STEM disciplines in Asia indicate that the proportion of female students is even lower than in other parts of the world.^{3,30} For the approach introduced in this study, there was no possibility to determine the gender on the basis of an Asian name, as explained in detail previously. The use of the *Chinese Name Gender Guesser*⁷ or other software platforms was not taken into consideration because these take the traditional Chinese characters of the name to classify the gender.

For our analysis, we also removed 391 additional names of unknown gender due to missing information. For example, 242 authors submitted only a single character as the first name. There is obviously no way to determine the gender by one letter. However, there is no evidence that there is a disproportionate percentage of women in this group. These names reflect 2.2% of the entire population and were therefore neglected.

Another assumption taken in this study is the Internet presence of the authors for the estimation of the correction factors in Eqs. (2)-(3). This assumption, however, turned out not to be critical since the percentages of authors not found on the Internet are in the same range for female and male authors.

In conclusion, we are presenting a bibliometric method to capture and classify female scientists that are currently active in research and in each a specific field of computer science. The group of female authors we captured with our method includes those female scientists successfully publishing their

research findings in peer-reviewed publications and, thus, having an impact on their scientific community. The data was collected regardless of the university degree and irrespective of whether the scientist is employed at a university or industry. The data provided by the presented method is closing the gap of postdoctoral researchers in industry and university existing in many other surveys of women in science. The method allows estimating the number of female candidates suitable for recruiting them as high-potential postdocs or professors and could also be used to address other questions of interest in the area of gender research as well as in a more general context of university research.

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - Project number 146371743 - TRR 89: Invasive Computing. 

References

1. Abramo, G., D'Angelo, C.A., and Caprasecca, A. The contribution of star scientists to sex differences in research productivity. *Scientometrics* 81, 1 (2009), 137–156; doi: 10.1007/s11192-008-2131-7.
2. Abramo, G., D'Angelo, C.A., and Murgia, G. Gender differences in research collaboration. *J. Informetrics* 7, 4 (2013), 811–822; doi: 10.1016/j.joi.2013.07.002.
3. Association for Academics and Societies in Sciences in Asia. *Women in Science and Technology in Asia*. Panmun Education Co., Ltd, Aug. 2015; http://bit.ly/38ZfIU.
4. Araújo, T. and Fontainha, E. The specific shapes of gender imbalance in scientific authorships: A network approach. *J. Informetrics* (Feb. 2017), 88–102; doi: 10.1007/s11192-011-0369-y.
5. Arensbergen, P., van der Weijden, I. and van den Besselaar, P. Gender differences in scientific productivity: A persisting phenomenon? *Scientometrics* 93, 3 (2012), 857–868; doi: 10.1007/s11192-012-0712-y.
6. Carsenat, E. What's the gender gap in the European Union Whoiswho? Sept. 2014; http://blog.namsor.com/2014/09/09/whats-the-gender-gap-in-the-european-union-whoiswho/.
7. Chinese Name Gender Guesser. July 2018; http://www.chinesetools.com/tools/gender-guesser.html.
8. *dblp computer science bibliography*. July 2018; http://dblp.uni-trier.de/db/.
9. De Paola, M. and Scoppa, V. Gender discrimination and evaluators' gender: Evidence from Italian academia. *Economica* (Dec. 2013), doi: 10.1111/ecca.12107.
10. Denisco, A. *The state of women in computer science: An investigative report*. Sept. 2017; https://tek.io/2QcvKCm.
11. European Commission. *She Figures 2015—Gender in Research and Innovation*; doi: 10.2777/744106.
12. Holman, L., Stuart-Fox, D., and Hauser, C.E. The gender gap in science: How long until women are equally represented? *PLoS Biology*, 2018.
13. Hyde, J.S., Fennema, E., and Lamon, S.J. Gender differences in mathematics performance: A meta-analysis. *Psychological Bulletin* 107 (Apr. 1990), 139–55; doi: 10.1037/0033-2909.107.2.139.
14. Jonathan, C. and Zuckerman, H. The productivity puzzle. *Advances in Motivation and Achievement* (Jan. 1984), 217–258.
15. Larièvre, V., Ni, C., Gingras, Y., Cronin, B., and Sugimoto, C. Bibliometrics: Global gender disparities in science. *Nature* 504 (Dec. 2013), 211–213; doi: 10.1038/504211a.
16. Larièvre, V., Vignola-Gagné, E., Villeneuve, C., Gélinas, P., and Gingras, Y. Sex differences in research funding, productivity and impact: An analysis of Québec university professors. *Scientometrics* 87, 3 (June 2011), 483–498; doi: 10.1007/s11192-011-0369-y.
17. LinkedIn. July 2018; url: https://www.linkedin.com.

18. Mihaljevic-Brandt, H., Santamaria, L. and Tullney, M. The effect of gender in the publication patterns in mathematics. *PLoS ONE* 11 (Oct. 2016), e0165367; doi: 10.1371/journal.pone.0165367.
19. NamSor Gender API. Sept. 2015; http://blog.namsor.com/api/.
20. NamSor Origin API. Sept. 2015; http://blog.namsor.com/name-recognition-software/.
21. O'Brien, K.R. and Hapgood, K.P. The academic jungle: Ecosystem modelling reveals why women are driven out of research. *Oikos* 121, 7 (2012), 999–1004; doi: 10.1111/j.1600-0706.2012.20601.x.
22. *People's Daily*. Chinese surname shortage sparks rethink. May 2007; http://en.people.cn/200706/19/eng20070619_385861.html.
23. Ranga, M., Gupta, N., and Etkowitz, H. *Gender Effects in Research Funding*, Mar. 2012.
24. Researchgate. *LinkedIn*. July 2018. url: https://www.researchgate.net.
25. Santamaria, L. and Mihaljević, H. Comparison and benchmark of name-to-gender inference services. In: *PeerJ Computer Science* (2018). doi: 10.7717/peerj-cs.156.
26. Science-Matrix, Inc. *Analytical support for bibliometrics Indicators to measure women's contribution to scientific publications*. Jan. 2018; http://bit.ly/35MNIao.
27. Shokhenmayer, E. and Carsenat, E. *Onomastics to Measure Cultural Bias in Medical Research Sing Scientists' Personal Name*. (Aug. 2014); http://bit.ly/2tGdSbr.
28. Stack, S. Gender, children and research productivity. *Research in Higher Education* 45, 8 (Dec. 2004), 891–920; doi: 10.1007/s11162-004-5953-z.
29. Stephen, C.J. and Williams, W.M. Understanding current causes of women's underrepresentation in science. In *Proceedings of the National Academy of Sciences* 108, 8 (2011), 3157–3162; doi: 10.1073/pnas.10148711108.
30. *Studentinnenanteile in Mathematik, Naturwissenschaften und Informatik sowie Ingenieurwissenschaften im internationalen Vergleich*. Center of Excellence Woman and Science, 2016.
31. Vichnevskaja, T. *Applying Onomastics to Scientometrics*. Jan. 2015; https://inserm.academia.edu/taniavichnevskaja.
32. Wang, M.T. and Degol, J.L. Gender gap in science, technology, engineering, and mathematics (STEM): Current knowledge, implications for practice, policy, and future directions. *Educational Psychology Review* 29, 1 (Mar. 2017), 119–140; doi: 10.1007/s10648-015-9355-x.
33. West, J., Jacquet, J., King, M., Correll, S.J. and Bergstrom, C.T. The role of gender in scholarly authorship. *PLoS ONE* 8 (July 2013), e66212; doi: 10.1371/journal.pone.0066212.
34. *Women, Minorities, and Persons with Disabilities in Science and Engineering*. National Science Foundation, National Center for Science and Engineering Statistics, 2018; http://www.nsf.gov/statistics/wmpdp/.
35. Zhao, H. and Kamareddine, F. Recursion identify algorithm for gender prediction with Chinese names. In *Proceedings of the Intern. Conf. Data Science* (Las Vegas, NV, USA, July 3–Aug. 2, 2018), 137–142.

Sandra Mattauch is a postdoctoral researcher at Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany.

Katja Lohmann is a postdoctoral researcher at Leibniz Universität Hannover, Germany.

Frank Hannig is a lecturer and senior researcher at Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany.

Daniel Lohmann is a professor at Leibniz Universität Hannover, Germany.

Jürgen Teich is a professor at Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany.

Copyright held by ACM.



Watch the authors discuss this work in the exclusive *Communications* video. <https://cacm.acm.org/videos/gender-gap>