Pre-publication draft –

Huang, Y., Zhang, Y., Wu, M., Porter, A., Barrangou, R. 2021, Determination of factors driving the genome editing field in the CRISPR era using bibliometrics, *The CRISPR Journal*, to appear.

# Determination of factors driving the genome editing field in the CRISPR era using bibliometrics

Ying Huang 1,2,3#, Yi Zhang 4#, Mengjia Wu 4, Alan Porter 5,6, and Rodolphe Barrangou 7\*

- 10 ¹ Center for Studies of Information Resources, School of Information Management, Wuhan
   11 University, Wuhan 430072, China.
- <sup>2</sup> Research Center for Science, Technology & Education Management and Evaluation , Wuhan
   University, Wuhan 430072, China.
- <sup>3</sup> Department of MSI, Centre for R&D Monitoring (ECOOM), KU Leuven, Leuven 3000, Belgium.
- 15 ORCID 0000-0003-0115-4581.
- <sup>4</sup> Australian Artificial Intelligence Institute, Faculty of Engineering and Information Technology,
- 17 University of Technology Sydney, NSW 2007, Australia. ORCID 0000-0002-7731-0301 ORCID
- 18 0000-0003-3956-7808.
- 19 <sup>5</sup> Search Technology, Inc., Norcross, GA 30092, United States.
- Frogram in Science, Technology & Innovation Policy, Georgia Institute of Technology, Atlanta,
   GA 30332, United States. ORCID 0000-0002-4520-6518.
- Department of Food, Bioprocessing and Nutrition Sciences, North Carolina State University,
   Raleigh, NC 27695, United States. ORCID 0000-0002-0648-3504.

# Authors contributed equally to this work.

\* Correspondence: rbarran@ncsu.edu.

#### Abstract

 Over the past two decades, the discovery of CRISPR-Cas immune systems and the repurposing of their effector nucleases as biotechnological tools have revolutionized genome editing. The corresponding work has been captured by 90,000 authors representing 7,600 affiliations in 126 countries, who have published over 19,000 papers spanning medicine, agriculture and biotechnology. Here, we use tech mining and an integrated bibliometric and networks framework to investigate the CRISPR literature over three time periods. The analysis identified seminal papers, leading authors, influential journals and rising applications and topics interconnected through collaborative networks. A core set of foundational topics gave rise to diverging avenues of research and applications, reflecting a *bona fide* disruptive emerging technology. This analysis illustrates how bibliometrics can identify key factors, decipher rising trends and untangle emerging applications and technologies that dynamically shape a morphing field, and provides insights into the trajectory of the genome editing.

# Keywords

CRISPR, Cas, CRISPR-Cas, genome editing, bibliometrics, disruptive technology

## Introduction

While genome editing has been on the rise over the past two decades, the advent of CRISPRbased (Clustered Regularly Interspaced Short Palindromic Repeats) technologies has accelerated and democratized genome editing in the past 9 years. <sup>1,2</sup> Several Cas-based (CRISPR associated) molecular machines have been co-opted from the bacterial adaptive immune system 3 to generate CRISPR-based technologies, such as sgRNA:Cas9 4, that have enabled facile genome editing since 2013. <sup>5,6</sup> Recently, the leading developers of this genome editing technology were awarded the 2020 Nobel Prize in Chemistry, illustrating the tremendous potential and impact of this technology. Early work focused on deciphering the molecular processes that drive CRISPRbased adaptive immunity in bacteria 7, and the development of programmable Cas proteins, that laid a preparatory foundation for CRISPR-based technologies. 8 Subsequently, these Cas effectors were deployed to manipulate genomes, transcriptomes and epigenomes in a broad diversity of organisms across the tree of life, such as bacteria, plants, and humans. 9 More recently, these CRISPR-based technologies have been widely adopted to engineer model organisms and even develop gene therapies tested in clinical settings. <sup>10</sup> Besides Cas9, the CRISPR toolbox has been expanded to encompass various Cas effector proteins such as Cas9, Cas12, Cas13, and Cascade. <sup>9</sup> As tools continue to be optimized with regards to specificity, efficiency, and delivery modalities, the intellectual property landscape is being defined <sup>11-13</sup> to enable widespread exploitation in medicine (e.g. gene therapies and antimicrobials), agriculture (e.g. crop breeding and disease resistance in livestock), and biotechnology (e.g. enzyme engineering and biofuel genesis). The accessibility and dissemination of CRISPR tools via repositories such as Addgene have allowed broad access to the best tools by academics and non-profit organizations across the globe. <sup>2</sup>

Though the rise of genome editing and global spread of CRISPR tools is undeniable, relatively little is recognized about the geographical, topical, individual and collaborative patterns that drive this academic phenomenon and commercially disruptive technology. <sup>14</sup> Here, we implemented an integrated research framework, using a bibliometric approach <sup>15,16</sup>, augmented by text mining, analysis of abstract record compilations and a scientific evolutionary pathway analysis <sup>17,18</sup>, to investigate the underlying patterns that have driven the adoption and implementation of CRISPR technologies. Specifically, we analyzed publication trends and authorship patterns for the CRISPR and the genome editing literature over space and time, using queries in the *Web of Science*, to identify key contributors and influential papers, as well as the topics that have shaped and are currently driving the field.

## Methods

 Publication records were retrieved using text queries mining the *Web of Science* records as of March 25<sup>th</sup>,2021, spanning manuscripts published between 2000 and 2020. Records were retrieved and cross-indexed using entries providing information with regards to manuscript

authors, affiliated institutions, publication journal, year, title and abstracts. For scientific evolutionary pathways (SEP) analysis, we used the method pioneered by Zhang *et al.* <sup>19</sup>, to trace the evolution of scientific topics into different subtopics by identifying a predecessor-descendant relationship from this bibliometric data. We then used this SEP approach to track the convergence and divergence of research topics on genome editing research and discover potential connections between these topics within a knowledge flow.

Generally, we ascribed six definitions as follows:

Definition 1: An article is represented by a vector (article vector): its feature space consists of terms of the entire dataset and its cell represents the frequency of a given term appearing in this article.

Definition 2: A topic is a collection of articles sharing similar semantic content, and is geometrically represented as a circle, with a centroid measured by the mean of all involved article vectors, and a boundary measured by the largest Euclidean distance between the centroid and all other article vectors.

Definition 3: Articles published in the same year are organized in one time slice. The entire dataset is analyzed as a bibliometric stream, that is, the SEP algorithm is to sequentially analyze each time slice according to the order of publication year, and for each time slice the algorithm is to sequentially analyze each article according to the order of unified publication ID.

Definition 4: Initial topics are topics consisting of articles in the first time slice and are starting points of the evolutionary pathways. Initial topics usually represent the root (e.g., original ideas and concepts) of the case (i.e. CRISPR in this paper).

Definition 5: A topic has two status categories, either 'live' or 'dead', as defined by 'sleeping beauties' <sup>20</sup>, for which a topic could 'die' if it does not receive new articles in certain sequential time slices, and a 'dead' topic could be revived and 'alive' again if a newly born topic shares the highest similarity with it.

Definition 6: A community is a group of proximate topics in a network – usually a branch in a SEP map-, which represents a subfield of the case.

Based on the above definitions, we implemented a stepwise algorithm to create the SEP as follows:

Step 1: All articles in the first time slice are grouped as one initial topic, which is set as the starting point of the evolutionary pathways. The algorithm moves to the second time slice and analyzes its involved articles one by one.

Step 2: We measure the cosine similarity between a current article and the centroids of all 'live' topics.

Step 3: We assign the article to its most similar topic. If the Euclidean distance between the article and the centroid of the assigned topic is smaller than its boundary, this article will be directly involved in the topic, or else, it will be labeled as 'drift.' Then, we return to Step 2 and analyze the next article until the end of this time slice.

Step 4: After analyzing all articles in one time slice, we check the status of each topic – i.e., set topics as 'dead' if they meet with the constraint in Definition 4 (a parameter is used here to decide the length of sequential time slices). For each 'live' topic, an unsupervised K-means approach is introduced to group its assigned 'drift' articles into certain sub-topics (an interval for seeking the local-optimal number of topics is required).

Step 5: We measure the cosine similarity between each sub-topic and two sets of topics - its assigned 'live' topic and all 'dead' topics. If the most similar topic of the sub-topic is its assigned one, their relationship is defined as 'predecessor-descendent,' or else, the most similar 'dead' topic will be revived and set as 'live,' and, then, becomes the predecessor of the sub-topic.

Step 6: We label a new topic (i.e., a sub-topic in Step 5) via the term with the highest similarity with all other terms in the topic - if the term has already been used before, choose the term with the second highest similarity, et cetera.

Step 7, We update the centroid and boundary of all 'live' topics, and the algorithm moves to the next time slice, and we return to Step 2.

Results of the SEP approach include a list of topics and their predecessor-descendant relationships. These topics are then visualized in a network via Gephi <sup>20</sup>. In the network, each topic is represented by a node, and the size of a node represents its importance, as measured by the value of term frequency inverse document frequency (tf-idf) analysis. A directed edge represents the predecessor-descendant relationship between its connected nodes, and the weight of an edge reveals the strength of the relationship (e.g., semantic similarity). The color of nodes reflects their communities identified by an approach of community detection integrated in Gephi as "modularity" <sup>21</sup>. Similarity measurements were carried out for the 119 topics identified across the three distinct time periods (9 topics pre-2013, 64 topics between 2013 and 2018, 46 topics since 2019), using semantic similarity coefficients. Details are available at: https://github.com/IntelligentBibliometrics/Gene-editing.

#### Results

CRISPR technology fueled the rise of the genome editing literature

To provide quantitative and qualitative insights into the drivers of the CRISPR craze <sup>22</sup>, we first defined the genome editing lexicon of interest and quantified relevant publications over the past twenty years, focusing on articles, reviews and letters comprising 26,484 records (Supplemental Table S1). Results show that the CRISPR literature (over 19,000 papers published since 2000 by

90,000 authors from around 7,600 institutions located in 126 countries; Supplemental Table S2) is rapidly growing, and that CRISPR-based tools impressively overtook incumbent technologies such as ZFNs, TALENs, and Meganucleases in 2013 (Figure 1A), within months of publication of the first proof of concept for CRISPR-based genome editing in human cells. <sup>5,6</sup> Currently, CRISPR-related publications account for the near totality of the genome editing field, and are over ten times more numerous than ZFN, TALEN, and Meganuclease papers combined (Figure 1A). Indeed, publications related to these first-generation genome editing technologies have been in decline since the advent of CRISPR-based genome editing technologies in 2012 (Figure 1A).

Amazingly, despite this rapid early adoption pattern, especially in the US and China, the CRISPR literature continues to expand at an impressive rate (Figure 1A), perhaps suggesting that genome editing is yet to hit maturity as a field, which is consistent with the continued dissemination of CRISPR tools across the planet. <sup>1,2</sup> Importantly, this shows how CRISPR as a field evolved from a relatively small "niche" microbiology topic into the major driver of genome editing in 2013, establishing a "before CRISPR" era <sup>23</sup>, and perhaps an "after displacement" of incumbent technologies period thereafter. This rise was fueled by the advent of the guide RNA technology in 2012, which quickly enabled genome editing (Figure 1B) and prompted an explosion in genome editing studies and citations (Figure 1C), as recognized by the 2020 Chemistry Nobel selection committee. Critical advances achieved in the past two years are also notable, with development of novel base editing tools and polished technologies such as prime editing <sup>24,25</sup>, as well as the transition of the technology from research laboratories into clinical settings with *bona fide* CRISPR-based therapeutics. <sup>10</sup> These tipping points triggered by specific publications and technology development define distinct time-periods that provide useful to assess the dynamic evolution of the field. <sup>23,26</sup>

## An interwoven network of collaborative authors

Next, we carried out a co-authorship network analysis to delve into the collaborative efforts driving contributions by the 48 most prolific and impactful authors, over time (Figure 2, Table 1). On a global basis, investigating publication patterns across these authors (as defined by number of publications, citations and h-index within the field), we note extensive and inter-connected collaborative networks with most authors engaged in several collaborative efforts. Actually, it appears the most influential authors collaborate with other key contributing authors in interconnected and overlapping authorship networks (Figure 2). Interestingly, many "early" authors who were active in the field prior to 2013 originally focused on CRISPR biology and mechanisms of action continue to do so (Figure 2), whereas distinct collaborative networks that fueled the rise of CRISPR-based genome editing technologies in parallel (Figure 2A) now directly overlap in topics of interest (Figure 2B). Noteworthy, the early community-wide focus on Cas9-based genome editing was comprising both overlapping and competitive interests, which created an intellectual property challenge regarding licensing and freedom to operate for the technology, 11-<sup>13</sup> which presumably prompted searches for novel Cas effectors. Interestingly, while some believe that the CRISPR IP challenges are a scientific hurdle that may have stifled innovation, the data suggests that it may rather have pushed the community towards actively mining for alternatives, while not precluding its broad adoption by diverse academic groups across the globe. Those initially established Cas12 as an alternative technology and recently unearthed new CRISPR-Cas types based on Cas13, Cas14 and others <sup>9,27</sup>, suggesting a need-based innovative push rather than a limiting competitive constraint.

Some of the most impactful contributions made by these influential authors can be captured by analyzing the most cited papers in the field (Table 2), over the three aforementioned eras, and the journals in which they have been published (Supplementary Table S3). The early contributions primarily consist of seminal studies establishing CRISPR-Cas as the adaptive immune system in bacteria <sup>7,26</sup>, providing DNA-encoded, RNA-mediated, nucleic acid targeting, culminating in 2012 with the development of the sgRNA:Cas9 programmable CRISPR effector. <sup>4</sup> This technology was used in 2013 for genome editing <sup>5,6</sup>, and shortly thereafter for transcriptional control and high-throughput screens. In the past two years, base editing technologies have been on the rise, primarily fueled by the rapid ascent of engineered Cas effectors from the David Liu lab (Table 1, Table 2). <sup>9,24,25</sup> Inevitably, the most cited manuscripts have been research papers published in high-profile journals contributed by prolific authors, together with a few noteworthy reviews and resource-focused papers (Table 2).

Predictably, citation patterns for most highly cited papers in the space reflect the rise of genome editing, notably the rapid explosion in 2013-2014 (Figure 1); these papers were published in the most influential journals in the world (Supplementary Table S3). Impressively, the most cited early CRISPR studies were also published in these journals, and they have been and continue to be the most influential journals in this field (Figure 1, Table S3), despite fundamental shifts in topics of interest and the vast expansion of the contributing authors pool, as well as a diversified and more global readership (Figure 2). To date, these papers reflect early work, mostly on development of the sgRNA:Cas9 technology, and its use and rapid adoption for genome editing in human cells, with the majority of the most cited papers published within the first 2 years of the CRISPR craze (Figure 1B).

In order to delve more into the key organisms, topics and genes subjected to the most attention in genome editing, we mined the published data and show that human cells are the primary organism of interest for the bulk of genome editing studies, predictably, followed by mouse, as the canonical proxy animal model for human studies (Supplementary Figure S1). Noteworthy, studies focused on humans and mice represent 10 times more than all other organisms of interest in CRISPR research, reflecting the heavy focus on human disease and medical applications, notwithstanding interest in and potential for other areas such as agriculture. Actually, this suggests that there is perhaps perplexing under-exploitation, or an adoption lag in other areas of interest, such as microbiology, which is ironically where these systems broadly occur and were originally characterized and repurposed. Next, we focused on key diseases of interest in these studies and determined that cancer-related research accounts for the majority of the studies, followed by genetic disease, and infectious disease, including viral infections (Supplementary Figure S1). This is further corroborated by the top 10 list of genes most associated with genome editing research (Supplementary Figure S1), notably the most studied trio: TP53 (the most popular tumor suppressor), AKT (protein kinase B), and MYC (protooncogene transcription factor).

# Emergence of networks of divergent genome editing topics

265266

267268

269

270

271

272

273

274

275

276277

278

279280

281

282

283

284

285286

287

288 289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305 306

307

308

To gain bibliometric insights into how the field evolved and morphed over time, we used a scientific evolutionary pathway (SEP) analysis (see Methods) to trace the evolution of topics of scientific interest in these published studies by identifying clusters of predecessor-descendant topical relationships. <sup>19</sup> This allowed tracking of convergence and divergence of research topics on genome editing and connections among these topics over time (Figure 3, Figure S1). This analysis revealed the existence of 9 topic communities that have evolved over the three time periods discussed previously. First, the field started with seminal bacterial work that occurred prior to 2012, which focused on adaptive immunity. This community topic is at the core of the network, and initially encompassed foundational topics such as "Cas nuclease", "acquired immunity", and "E. coli" (see the pink cluster at the center of Figure 3 and Supplementary Figure 1). This core gave rise to the sgRNA:Cas9 genome editing technology, a tipping point for the field, which emerged as a new topic in 2013, centered on "guide RNA", and links to incumbent genome editing technologies such as ZFNs and TALENs (see the green cluster, Figure 3). Over time, the core also gave rise to a community focused on screens (genetic screens, high-throughput screens, center right purple cluster). Likewise, the core cluster also gave rise to a community topic focused on transcriptional control, relatively early on with the rise in 2014 of a transcription-focused cluster encompassing gene expression, gene regulation, transcription factors and transcriptional regulators (center left, blue). Later on, as the technology evolved and matured, applicationfocused clusters arose, focusing on gene therapies, viral diseases, and neurodegenerative diseases.

Analysis of similarity measurements (Supplementary Table S4) between these topic communities reveals how disruptive CRISPR technology is, given the diversity of distinct clusters that arose from the original core cluster, and the relatively low level of similarity observed between and across these 119 topics. This is further supported by the low level of similarity observed between topics across time periods (Supplementary Table S4). The recent increase in topics in the past two years (46 new topics in two years, compared with 64 topics spanning the explosive 2013-2018 period) likely indicates continued disruptive innovation and expansion of this technology into new areas of research, as well as novel and diversified applications. This is consistent with the development of novel technologies (e.g. base editing), the continued dissemination of CRISPR technologies across the globe (e.g. Addgene distributions) and the transition to applications, especially in therapeutic settings with CRISPR-based diagnostics, antivirals and gene therapies all with clinical ambition in the short term. Critically, it is important to note the cross-referencing of the various visualization modalities and tabular lists of entries throughout our tables and figures, that consistently identify the same key factors fueling the genome editing revolution, and robustly establish the seminal studies and technological developments that have shaped this morphing subject over time.

Despite the observed congruence, the SEP algorithm relies on natural language processing techniques that are impacted by writing style and biases, as well as inconsistent use of terminology by different groups of authors, which can lead to synonyms being redundant and

separately accounted for. For example, there are entries related to transcription that encompass: "transcriptional control", "gene expression", "gene regulation" and "transcriptional regulation". There are also several connections between seemingly un-related topics due to language biases and topic-related complexity inherent to the same technology being used in unrelated organisms. There are also multiple examples of confounding coverage of topics that are often discussed together, but are not systematically linked, such as "human embryos" and "clinical trials" being discussed together without being co-dependent. Thus, the complexity of a broadly applicable tool must be deciphered and interpreted by the expert reader to account for otherwise unrelated topics and verbiage. Human interpretation is also important to fully assess the impact and influential contributions of individual authors and select manuscripts, to account for quantitative shortcomings and biases inherent to citation numbers, indices and impact factors. Indeed, qualitative insights should be used by the reader to complement quantitative metrics in the spirit of the *Leiden Manifesto*. <sup>28</sup> This manifesto highlights the need to rely on expert assessment to overcome bias tendencies and untangle conceptual ambiguity and uncertainty.

In several instances, there are connections that seems counter-intuitive and reflect high semantic similarity, but not technical dependence nor scientific derivation. Indeed, sets of authors can share similar language biases, such as clinically-relevant settings for patient sampling in medical applications for the epidemiological study of Mycobacterium tuberculosis and the implementation of genome editing for human gene therapies, linking two seemingly unrelated clusters because the authors share linguistic biases and keywords. Likewise, the link between Cas nucleases and DNA fingerprinting reflects the early use of CRISPR spacer hypervariability for genotyping and not the use of Cas proteins for molecular fingerprinting. This high semantic similarity need not reflect bona fide technical overlap or dependency, and can reveal linguistic biases, or indicate subsequent uses and applications of derived tools and technologies, including their eventual use in diverse model organisms. The latter explains the unexpected appearance of Saccharomyces cerevisiae, Caenorhabditis elegans, zebrafish, Chinese hamster ovary cells, and others throughout topic clusters. Some of the topical lineages shown reflect topical descendance within the CRISPR literature that evolved from a technical basis (using various Cas effectors as tools) to applications of these technologies in model organisms and cells. To a similar extent, select topics of interest to specific groups of authors and readers can be linked through SEP analyses such as "human embryos" and "clinical trials", though they need not be co-dependent (current clinical trials are not based on CRISPR-edited human embryos), so both applications and implications can entangle topic connections. In some cases, the appearance of a newly coined term reveals tipping points that created new sets of topics, notably the development of the "guide RNA" technology and the nomenclature update that reclassified Cas5/Csn1 as Cas9.

While some literature topics have arisen faster than CRISPR, such as the recent COVID19-related literature<sup>29</sup>, the speed of the adoption of the CRISPR technology, as much as the rise of the CRISPR-related literature, is noteworthy. The speed of the work in this field has been invoked as a distinguishing feature, but perhaps the most striking aspect is the adoption and democratization of the technology itself, which is captured by the rise in the number of citations and publications, as well as Addgene shipments. <sup>1-2</sup>

#### Discussion

353

354 355

356

357

358

359

360

361362

363

364

365

366

367

368

369

370

371372

373

374

375

376

377

378

379

380

381

382

383 384

385

386 387

388 389

390

391392

393 394

395

396

Altogether, these results provide insights into the key factors driving the evolution of CRISPR, and illustrate how a diverse community of collaborative scientists is globally adopting this disruptive technology, and implementing it in various organisms of interest across applications. This analysis illustrates how bibliometrics can identify key individuals, topics and papers that dynamically shape a morphing research field, and decipher rising trends impacting the historical trajectory of a field and untangle emerging applications. The data presented here provide strong support that this is a bona fide emerging technology as defined by key attributes. 30 Indeed, all five defining elements of an emerging technology are met, with: (1) radical novelty: near-instant replacement of incumbent editing technologies, with aggressive pursuit of IP and topic diversification; (2) fast growth, as documented by publications, citations, and Addgene distribution patterns; (3) coherence, supported by overlapping collaborative authorship networks, as well as interconnected topics derived from a common core; (4) prominent impact, with enthusiastic commercialization in several industries spanning medicine, agriculture and biotechnology, as well as global adoption in academia and industry and the momentous 2020 Nobel Prize in Chemistry for two selected CRISPR pioneers; and (5) uncertainty and ambiguity, as documented by intellectual property issues, discussions related to regulatory frameworks for, and societal implications of, the various applications of genome editing. <sup>30</sup> Importantly, the evolution of the topic map over the three aforementioned time periods further endorses the emerging technology attributes of genome editing. Indeed, predecessor-topics created during the first time period established a scientific foundation for the field (coherence), with evolution over the next two time periods radically spearheading into various directions (radical novelty), with rapidly increasing number of descendant topics (fast growth), giving rise to diverse research foci. The eclectic community diversity is noteworthy, in terms of institutional affiliations, geographical location and scientific topics of interest, which collaborations transcend, as illustrated by coauthorship patterns. Yet, the overall primary focus is mostly on human therapeutic applications, reflecting the tremendous potential of genome editing implementation in the clinic, and the need to deploy CRISPR therapies for patients afflicted by genetic diseases. With FDA-enabled trials actively underway, confidence in regulatory agencies and progressing public engagement dialogues encompassing ethical, legal and societal implications <sup>31, 32</sup>, we anticipate the literature will continue to expand and hopefully document larger and broad clinical success in the near future, as well as fuel applications in agriculture and sustainability.

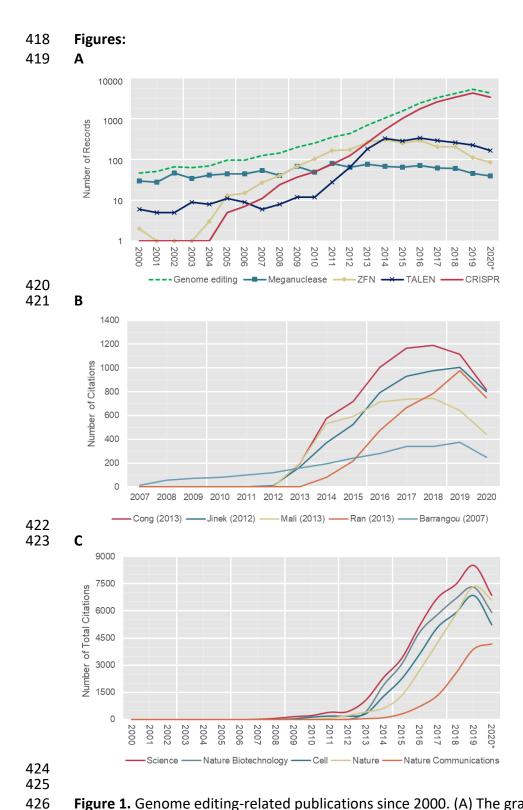
## **Acknowledgements**

The authors would like to acknowledge their lab members, collaborators and colleagues throughout the community for fruitful discussions and insightful opinions.

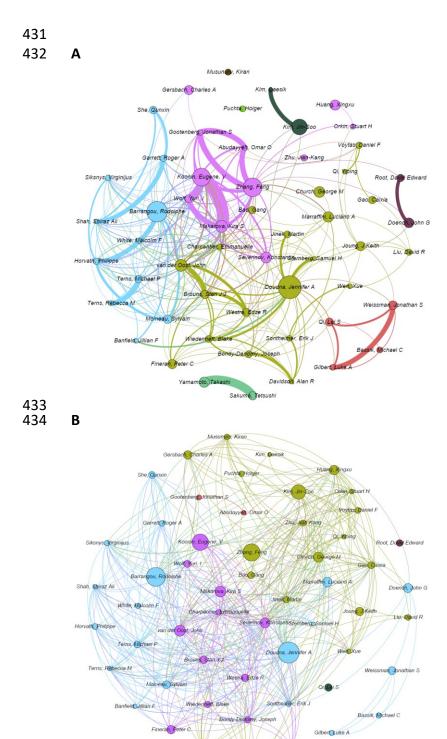
## **Author Contributions**

All authors contributed to the conceptualization of the study, interpretation of the results and drafting of the manuscript; YH, YZ and MW carried out analyses and generated figures.

397 **Author Disclosure Statement** 398 399 RB is a shareholder of Caribou Biosciences, Intellia Therapeutics, Locus Biosciences, Inari Ag, 400 TreeCo, Ancilia Biosciences and CRISPR Biotechnologies; AP is a shareholder of Search 401 Technology Inc. 402 403 **Funding Statement** 404 YH acknowledges support from the National Natural Science Foundation of China (Grant No. 405 72004169); AP acknowledges support from the US National Science Foundation (Award 406 #1759960) to Search Technology, Inc., and Georgia Tech; YZ and MW acknowledges the 407 Discovery Early Career Researcher Award granted by the Australian Research Council (Grant No. DE190100994). 408 409 410 3 Figures and 2 Tables 411 412 **Supplementary Material** 413 Supplementary Table S1 414 Supplementary Table S2 415 Supplementary Table S3 416 Supplementary Table S4 Supplementary File S1 417

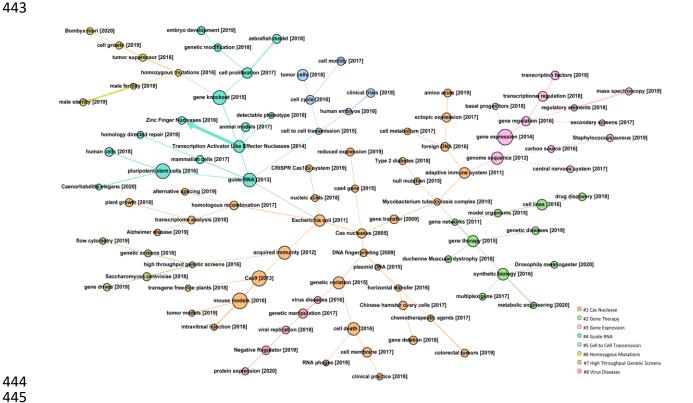


**Figure 1.** Genome editing-related publications since 2000. (A) The graph shows the number of publications related to genome editing and their various effectors, including Meganucleases, ZFNs, TALENs and CRISPR. The number of publications is showcased in a log10 scale. (B) Citations over time for the 5 most cited CRISPR papers; (C) citations for CRISPR papers published in selected journals, over time.



**Figure 2.** Collaborative authorship networks between the 48 most impactfully-prolific CRISPR researchers whose H-index within this topic is more than 20 since 2000. (A) co-authorship network, where node size reflects the number of records published by authors, lines reflect co-authorships, and the cluster colors reflect community detection algorithm-based groups; (B)

cosine similarity network, with cluster colors reflecting topic similarities (mesh terms allocated to the publications); only lines with similarities higher than 0.3 are shown.



**Figure 3.** Scientific evolutionary pathway (SEP) analysis of CRISPR and genome editing topics over time. Nine topic communities are represented using distinct colors, connected over time. Topics are linked using predecessor-descendant relationships defined by the literature patterns.

# References

- LaManna CM & Barrangou R. Enabling the Rise of a CRISPR World. *CRISPR J* 2018; **1:**205-208. doi:10.1089/crispr.2018.0022.
- LaManna CM, Pyhtila B, Barrangou R. Sharing the CRISPR Toolbox with an Expanding Community. *CRISPR J* 2020; 3:248-252. doi:10.1089/crispr.2020.0075.
- Barrangou R. *et al.* CRISPR provides acquired resistance against viruses in prokaryotes. 5 Science 2007; 315:1709-1712. doi:10.1126/science.1138140.
- 458 4 Jinek M *et al.* A programmable dual-RNA-guided DNA endonuclease in adaptive bacterial immunity. *Science* 2012; 337: 816-821. doi:10.1126/science.1225829.
- Cong L et al. Multiplex genome engineering using CRISPR/Cas systems. Science 2013;
   339:819-823. doi:10.1126/science.1231143.
- 462 6 Mali P *et al.* RNA-guided human genome engineering via Cas9. *Science* 2013; 339:823-463 826. doi:10.1126/science.1232033.
- Hille F et al. The Biology of CRISPR-Cas: Backward and Forward. Cell 2018; 172: 1239 1259. doi:10.1016/j.cell.2017.11.032.
- Hsu PD, Lander ES, Zhang F. Development and applications of CRISPR-Cas9 for genome engineering. *Cell* 2014; 157:1262-1278. doi:10.1016/j.cell.2014.05.010.
- 468 9 Anzalone AV, Koblan L W, Liu DR. Genome editing with CRISPR-Cas nucleases, base editors, transposases and prime editors. *Nat Biotechnol* 2020; 38:824-844. 470 doi:10.1038/s41587-020-0561-9.
- 471 10 Frangoul H *et al.* CRISPR-Cas9 Gene Editing for Sickle Cell Disease and beta-Thalassemia. 472 *N Engl J Med* 2020; doi:10.1056/NEJMoa2031054 (2020).
- 473 11 Egelie K J, Graff GD, Strand SP, Johansen B. The emerging patent landscape of CRISPR-474 Cas gene editing technology. *Nat Biotechnol* 2016; 34: 1025-1031. 475 doi:10.1038/nbt.3692.
- Sherkow JS. The CRISPR Patent Landscape: Past, Present, and Future. *CRISPR J* 2018; 1:5-477 9. doi:10.1089/crispr.2017.0013.
- 478 13 Martin-Laffon ., Kuntz M, Ricroch AE. Worldwide CRISPR patent landscape shows strong geographical biases. *Nat Biotechnol* 2019; 37:613-620. doi:10.1038/s41587-019-0138-7.
- Huang Y, Porter A, Zhang Y, Barrangou R. Collaborative networks in gene editing. *Nat Biotechnol* 2019; 37: 1107-1109. doi:10.1038/s41587-019-0275-z.
- 482 15 Porter AL, Kongthon A, Lu JC. Research profiling: improving the literature review.
  483 *Scientometrics* 2002; 5:351-370. doi: 10.1023/A:1014873029258
- 484 16 DeBellis N. Bibliometrics and citation analysis. *The Scarecrow Press* Lanham, MD 2009.
- Porter AL, Cunningham SW. Tech mining: exploiting new technologies for competitive advantage. *Wiley, NY* 2005. doi:10.1016/j.ipm.2005.01.005
- Porter AL, Youtie J. Where does nanotechnology belong in the map of science? *Nat Nanotechnol* 2009; 4:534-536, doi:10.1038/nnano.2009.207.
- Zhang Y, Zhang G, Zhu D, Lu, J. Scientific evolutionary pathways: identifying and
   visualizing relationships for scientific topics. *J. Ass. Info. Sci. Tech.* 2017; 68:1925-1939.
   doi:10.1002/asi.23814.
- 492 20 Bastian M, Heymann S, Jacomy M. Proc. Third Int. ICWSM Conf. 361-362.

493	21	Newman ME. Modularity and community structure in networks. <i>Proc Natl Acad Sci U S A</i>
494		2006; 103:8577-8582. doi:10.1073/pnas.0601602103.
495	22	Pennisi E. The CRISPR craze. <i>Science</i> 2013; 341:833-836,
496		doi:10.1126/science.341.6148.833.
497	23	Urnov FD. Genome Editing B.C. (Before CRISPR): Lasting Lessons from the "Old
498		Testament". CRISPR J 2018; 1:34-46. doi:10.1089/crispr.2018.29007.fyu.
499	24	Anzalone AV et al. Search-and-replace genome editing without double-strand breaks or
500		donor DNA. <i>Nature</i> 2019; 576: 149-157. doi:10.1038/s41586-019-1711-4.
501	25	Gaudelli NM et al. Directed evolution of adenine base editors with increased activity and
502		therapeutic application. Nat Biotechnol 2020; 38:892-900. doi:10.1038/s41587-020-
503		0491-6.
504	26	Barrangou R, Horvath P. A decade of discovery: CRISPR functions and applications. Nat
505		Microbiol 2017; 2:17092. doi:10.1038/nmicrobiol.2017.92.
506	27	Makarova KS et al. Evolutionary classification of CRISPR-Cas systems: a burst of class 2
507		and derived variants. Nat Rev Microbiol 2020; 18:67-83. doi:10.1038/s41579-019-0299-
508		X.
509	28	Hicks D, Wouters P, Waltman L, deRijke S, Rafols I. Bibliometrics: the Leiden Manifesto
510		for research metrics. <i>Nature</i> 2015; 520:429-32. doi:10.1038/520429a
511	29	Porter AL, Zhang Y, Huang Y, Wu M. Tracking and mining the COVID-19 research
512		literature. Front. Res. Metr. Anal. 2020; 5:594060.
513		https://doi.org/10.3389/frma.2020.594060
514	30	Rotolo D, Hicks D, Martin BR. What is an emerging technology? <i>Res. Policy</i> 2015;
515		441827-1843. doi: 10.1016/j.respol.2015.06.006
516	31	Sherkow JS. Controlling CRISPR Through Law: Legal Regimes as Precautionary Principles.
517		CRISPR J 2019; 2:299-303. doi:10.1089/crispr.2019.0029 (2019).
518	32	Howell EL, Yang S, Beets B, Brosard D, Scheufele DA, Xenos MA. What do we (not) know
519		about global views of human genome editing? Insights and blind spots in the CRISPR era.
520		CRISPR J 2020; 3:148-155. doi: 10.1089/crispr.2020.0004
521		