

# Application of bibliographic data-analysis and clustering for policy makers applied to hypersonic terms in Australia

Raoul Mazumdar

*Melbourne, Australia*

---

## Abstract

Data analysis and primarily, clustering, of bibliographic information can yield information for policymakers in technology canvassing, current research clusters, and provide a convenient interface for data interrogation. The use of clustering metrics can show researchers whom the most influential members or institutions are within a specific technology area, thereby assisting in the tendering process of any prospective technology and policy. Furthermore, the information can unveil the current nature of the research work in the domestic domain, and help policies focus on potential opportunities. The approach has been applied to general hypersonic related terms, including scramjets and ramjets in this work. The outcome shows that the University of Queensland and the University of New South Wales at Canberra are influential partners in this domain. This is not necessarily new information, as this is well depicted in the abundance of papers present in hypersonic literature. However, the clustering metrics highlight other influential authors within other governmental agencies and the University of Southern Queensland and the Royal Melbourne Institute of Technology. The clustering work also highlights several different private organisations who have contributed in some form to the Hi-Fire hypersonics research program, further emphasising the potential of clustering as a tool for policy development. The cluster's general structure and interaction and the applied nature of hypersonic research indicate that domestic collaboration is underutilised and could be enhanced if a policy is developed to address the unconnected research groups.

*Keywords:* Clustering, Bibliographic, Data-analysis, Hypersonics, Policy

---

**1 1. Introduction**

2 From a policy point of view, it is crucial to know the current research  
3 level and collaboration within the prospective technological area to address  
4 them better. If a field is shown to lack academic research and staff, the policy  
5 can be driven to develop the fundamental research. Whereas if a particular  
6 technology area shows a large degree of research and development, the policy  
7 can help technology diffusion towards the application phase. Understanding  
8 the interconnectivity of a research cluster, key staff members, and growth  
9 is valuable. In semiconductors, quantum computing and other niche areas,  
10 canvassing would often require a subject matter expert to disseminate the in-  
11 formation. On the other hand, clustering and bibliographic analysis can be  
12 re-purposed as a fast canvassing tool and visualisation aid, for the inspection  
13 of technologies, their clusters, and critical members for engagement.

14

15 Clustering of publications is not a new activity in itself, and have been ap-  
16 plied to several patent databases [1, 2]. In other instances, publications have  
17 been analysed concerning geography as per the works of Baptista [3]. The  
18 paper shows how proximity can be strongly correlated with enhanced tech-  
19 nology diffusion and innovation. Additional research by Liyanage [4] further  
20 reinforces a positive coupling between research collaboration and innovation  
21 clusters. The bodywork suggests that co-authorship of journal publications  
22 are linked to positive research outcomes. The research clustering's volume  
23 or size in different institutions also holds weight in terms of publication, pro-  
24 ductivity, and general output. Dundar and Lewis [5] discussed this, inves-  
25 tigating productivity in universities within the United States across several  
26 disciplines. The findings show again, correlation to the size of a depart-  
27 ment and group and their research output, including journal publications.  
28 The strength of the correlation happens to be specific to the scientific field.  
29 Fundamental research activities are shown to garner more international col-  
30 laboration. The paper also expresses an issue where domestic entities are less  
31 likely to collaborate due to competitive advantage [6]. Unfortunately, this  
32 lack of domestic collaboration is more pronounced with application-oriented  
33 research, in which mature hypersonic research would coincide within.

34

35 It is important to note that co-authorship does not necessarily indicate  
36 cross-collaboration, visiting PhD students, research exchanges, and one-time  
37 projects can produce those co-authorships without long term cooperation re-

38 search. However, the pattern of the co-authorship clusters can shed light as  
39 to the nature of the relationship. The works of Velden et al. [7] describe  
40 the reasoning behind the particular formation of clusters, indicating whether  
41 clusters are a result of cooperative agreements, staff migration or dominated  
42 by a single principal investigator. This leads to a theoretical study of clustering  
43 by Molontay and Nagy [8]. The work shows a strong correlation between  
44 individuals with high centrality in the cluster, and high citation counts and  
45 other scientometric values. Co-authorship is also shown to be weakly coupled  
46 to the citation, thereby higher co-authorship is likely to produce publications  
47 with more significant citations. Articles also receive higher citations when  
48 they expand upon a group of publications, thereby closely tied to multiple  
49 sources and hence the possibility of further citations[9]. Thus individuals  
50 with high centrality in bibliographic clusters are likely to contribute to large  
51 research outputs, regardless of the journal field.

52

53 Figure 1 a.) shows a co-authorship cluster centred around a primary  
54 node. This node generally has an abundance of co-authorships and is gen-  
55 erally the principal investigator institution or university. The second case  
56 in Fig. 1 b.) highlights the connection of a single node between two main  
57 clusters, the reasoning for this can stem from staff migration, visiting re-  
58 searchers, and or one-off funded projects. Generally speaking, they indicate  
59 temporary positions and academic staff’s movement, thus leading to the co-  
60 authoring pattern as presented. In Fig 1 c.) we can visualise what is likely  
61 to be an international cooperative research initiative or some domestic re-  
62 search hub. The multiple large nodes represent small interconnected research  
63 groups, and can be considered multiple principal investigators closely linked  
64 together. Multiple connections between separate nodes with the principal in-  
65 vestigator represent large cooperative agreements, broad research exchanges,  
66 joint-PhD students and or the application of complementary knowledge.

67

68 Following this, the work looks to investigate research clusters applied to  
69 an area of hypersonics, which includes the re-entry of space vehicles, and  
70 high-speed aircraft and rockets. The area has undergone intensive develop-  
71 ment since the 1960s but remains a resource-intensive field [10]. First, the  
72 flow-speeds and high enthalpies can only be recreated in shock tubes and  
73 shock-reflection tunnels, which generate milliseconds of applicable air-flow.  
74 Furthermore, the computational tools coupled with reacting chemistry is nu-  
75 matical burdensome. These factors, coupled with fundamental research and

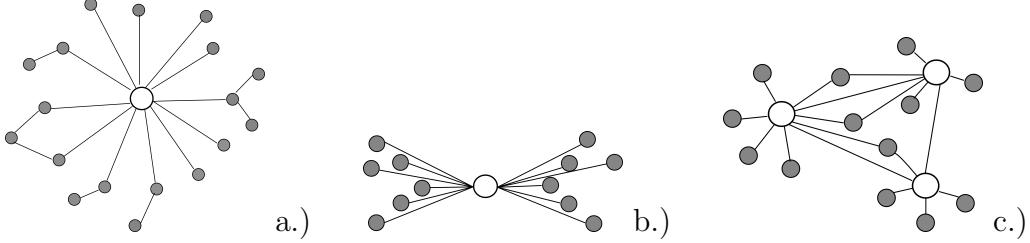


Figure 1: Typical clustering arrangements adapted from Velden et al. [7]

76 laser diagnostics' growth, have slowly allowed the maturity of the technological  
 77 field. That said, within Australia, the development of a mature hypersonic system is still pending. Thus, examining the research cluster within  
 78 79 Australia is an area of interest, and how academic and industry partners are  
 80 collaborating through journal publications.

81

82 In this body of work, bibliographic journal publications about hypersonics are analysed using general data analytics and a clustering-based method.  
 83 Due to the field's sensitivity, journal papers are only sourced from keyword  
 84 searches that include hypersonics, scramjets or ramjets. This happens to  
 85 exclude other fundamental research areas and complementary technological  
 86 areas. The research is conducted through the use of data-harvesting from  
 87 the Scopus journal database and followed by post-processing through python  
 88 based programs to sort names and university affiliation. The processed information  
 89 is then imported into Gephi1, open-source software for manipulating  
 90 and exploring networks as depicted by Fig. [11].  
 91

92

## 93 2. Method

### 94 2.1. Data Harvesting and Sorting

95 Various journal publications are sourced through Scopus using an application  
 96 programming interface written in Python for importing journal authors,  
 97 affiliations, abstracts, keywords and publication dates. However, as mentioned  
 98 in papers concerning bibliographic data-analysis, name identification  
 99 is a problem [7]. Different journal publications hold different name conventions  
 100 and express a single author in different formats, thereby alluding to  
 101 multiple authors instead of one. In order to account for this, each new name  
 102 in a data-set is checked for similarity for any other names within the data-set.

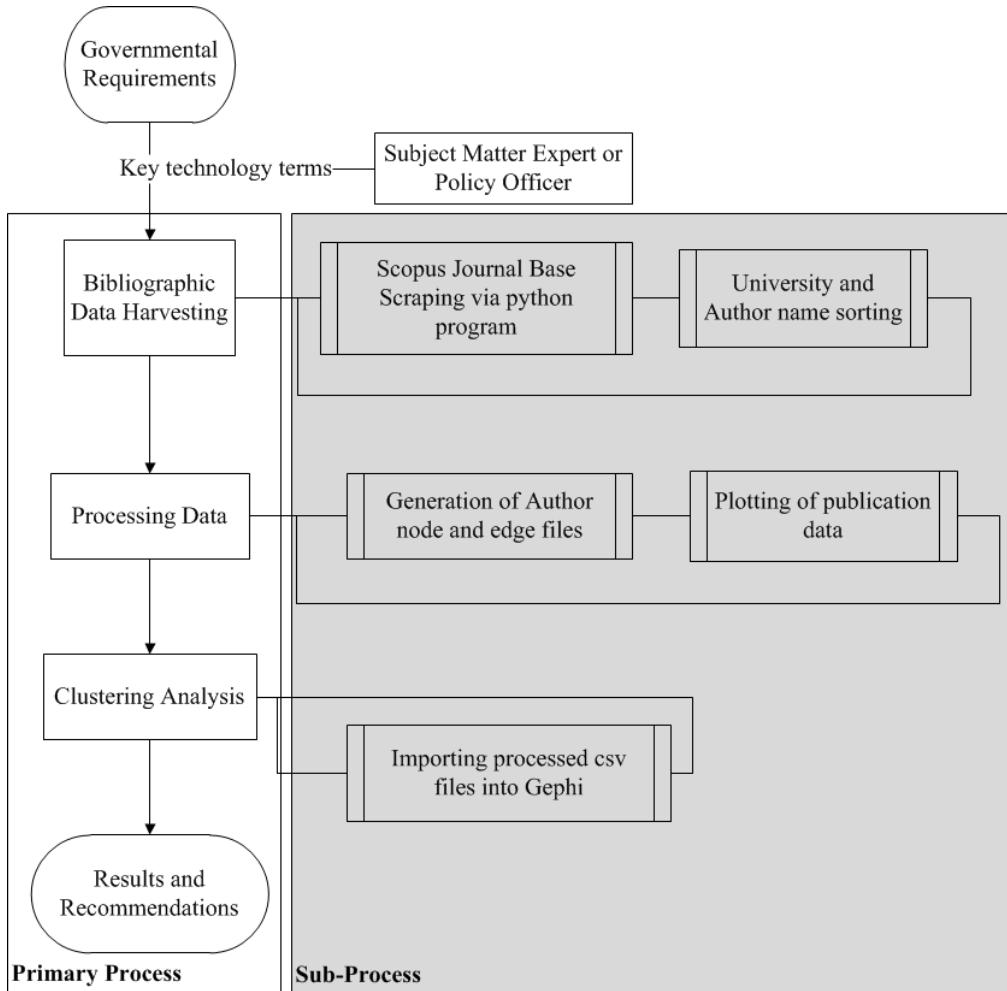


Figure 2: Data-analysis and method

<sup>103</sup> The process works by generating a list of possible permutations of a name  
<sup>104</sup> and checks the prospective strings using Levinstein distance [12].

<sup>105</sup> Similarly, domestic universities are also subject to different naming con-  
<sup>106</sup> ventions. For example, the University of New South Wales can also be rep-  
<sup>107</sup> resented as UNSW Sydney. A solution includes the checking names against  
<sup>108</sup> a common set of aliases within the python program using the Levinstein dis-  
<sup>109</sup> tance. Furthermore, prefixes and suffixes are added and checked for a high  
<sup>110</sup> degree of similarity from a target string. These words include common terms  
<sup>111</sup> that may be used, such as, school of engineering or engineering department.

112 While these approaches are effective for most data strings, there are errors  
113 due to the tolerance required for name similarities to be met.

114

### 115 *2.2. Clustering and Mapping*

116 Upon completing the previous task of sorting the data into two formats,  
117 nodes and edges. In which nodes represent an author affiliation, and edges  
118 represent an author's list of collaborators. The material is ready to be im-  
119 ported into an open-source software called Gephi. This software allows for  
120 the visualisation of nodes in clustered maps, accommodating various clus-  
121 tering modes. Modes such as force-atlas, where nodes act as point charges  
122 but edges act as attractive masses. Thereby leading to well-connected nodes  
123 forming at the heart of a group and disconnect authors pushed to the exte-  
124 rior. In this work, a Fruchterman Reingold approach is used, which provides  
125 better visualisation for independent groups. Although the clustering map is  
126 more comfortable to visualise, there can be thousands of nodes in some cases.  
127 While examining all the nodes might be fruitful, the time required is imprac-  
128 tical. Through Gephi, it is possible to conduct statistical analysis for nodes  
129 that are centrally position. The term, betweenness centrality, represents the  
130 shortest distance between other co-authors. They can represent principal  
131 investigators, mobile staff, or generally well connected or influential indi-  
132 viduals. Using this metric, we generate a list of well-positioned individuals  
133 which researcher can instantly engage on for policy or tendering processes.

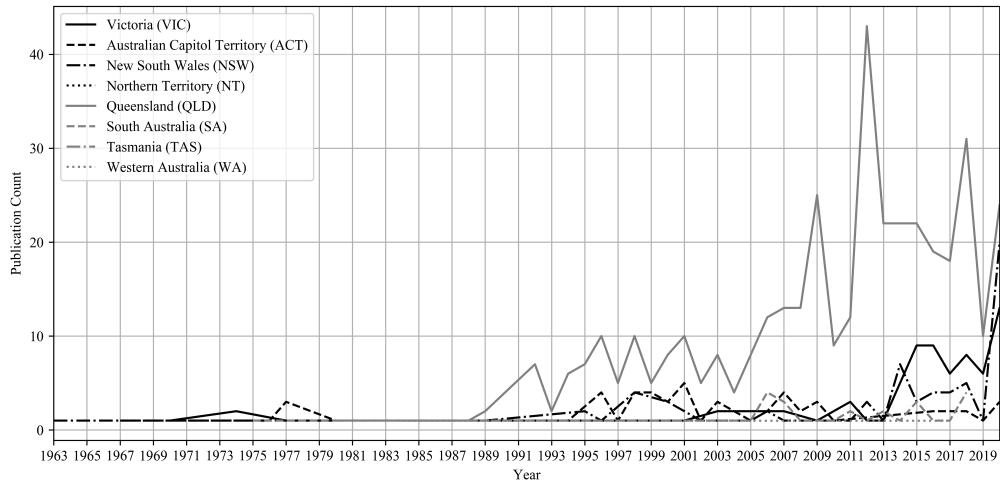
### 134 *2.3. Result*

135 The bibliographic information's initial findings show a large publication  
136 volume from the state of Queensland, followed in part by New South Wales,  
137 and Victoria as per 2.3 a.). The Australian Capital Territory is highly in-  
138 fluential in having an early publication history dating back to the 70s, but  
139 is under-represented in the publication count. The Australian National Uni-  
140 versity and the University of New South Wales at Canberra have been active  
141 in hypersonic related research. However, the University of New South Wales  
142 at Canberra is a satellite campus of the University New South Wales at Syd-  
143 ney due to the data-analysis and arrangement. Hence, publications from  
144 that satellite campus are merged with the primary university, thus overes-  
145 timating the contribution of New South Wales while under-representing the  
146 Australian Capital territory. That aside, we can also see the country affil-  
147 iation of co-authors in the data-set in Fig. 2.3 b.). The entire volume of

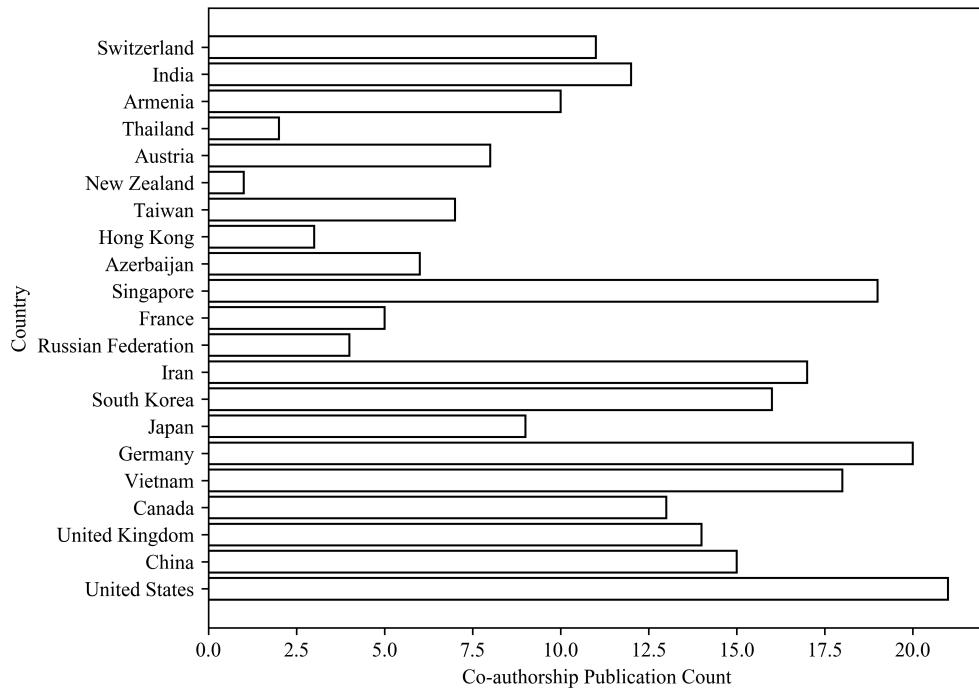
148 overseas authors is relatively small in the data-set, compared to the volume  
149 of papers published. Interestingly the largest volume of foreign co-authorship  
150 is from the United States, Germany and Singapore in that order. Examining  
151 the geography information can provide us with insight into the general  
152 research trends, but this excludes various other information. If we wish to  
153 understand the key players in the data-set, a clustering approach is further  
154 required.

155 A more visually informative approach is shown by clustering the data  
156 using the open-source software Gephi and the imported processed datasets.  
157 Figure 2.3 shows the full co-authorship cluster, which holds 972 authors rep-  
158 resented as nodes, including 4990 edges representing co-authorship of papers.  
159 Around 42% of the nodes represent either an international entity or a do-  
160 mestic industry as a co-author. The rest of the nodes are derived from some  
161 author within an Australian university. The second-largest number of nodes  
162 belongs to the University of Queensland, at 22% of the total. The University  
163 of New South Wales follows them at 11.83%, Monash University at 4.8%,  
164 the University of Sydney at 4.53%, RMIT university at 2.67%, and the Uni-  
165 versity of Melbourne at 2.26%. The rest of the universities make up the  
166 remainder of the total. In the clustering map, the University of Queensland  
167 and the University of New South Wales are depicted in the cluster's centre  
168 with interconnected nodes and generally are well represented.

169  
170 The defence science institute is detected in clustering on the lower right  
171 end of the map. The authors are shown to interface predominately with the  
172 University of New South Wales at Canberra, and the University of Queens-  
173 land. Other universities are generally disconnected from the co-authorships  
174 in hypersonic material from the institute. Within the cluster are a number  
175 of authors from the Boeing corporation, co-authoring papers regarding the  
176 Hi-Fire program. In the upper corner of the map is a dense, evenly spaced  
177 cluster, this happens to represent a paper concerning the luminosity of a stel-  
178 lar body. Although this is outside of hypersonic systems' general scope, the  
179 paper includes mention of hypersonic shockwaves and thus was included in  
180 the bibliographic data harvest. In a similar instance, another group, discon-  
181 nected from the majority of papers in the central map, discusses hypersonic  
182 photons, and again is outside of the search's primary aims. Thus, using this  
183 particular clustering algorithm, papers in the centre of the map represent  
184 similar research fields. In this case, those fields include hypersonics for the  
185 means of propulsion, whereas separate groups on the fringes may investigate



a.) Publication frequency with respect to Australian states and territories



b.) Publications by volume from foreign countries

Figure 3: Geographic breakdown of hypersonic related publications

186 hypersonics outside of the intended purposes.

187

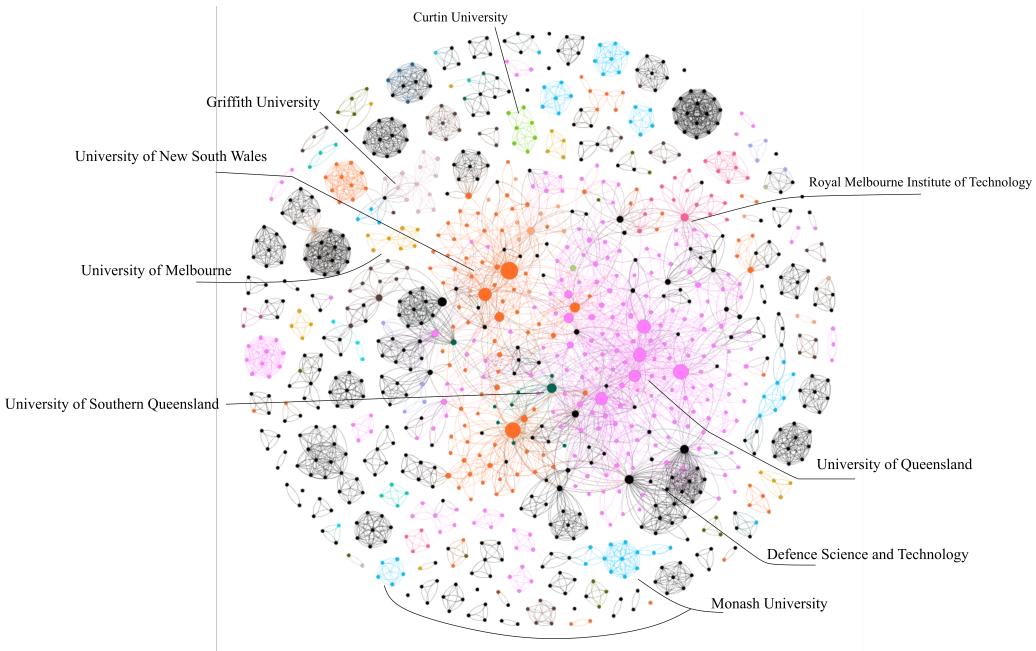


Figure 4: Co-authorship clustering

188 The full clustering map still includes over 800 authors, and analysing  
189 each author might be too time-consuming. But it is possible to identify key  
190 authors in the group, not through citations, or the number of publications,  
191 but due to their centrality within the co-authors. For example, Table 1 lists  
192 the key authors in the grouping with respect to university affiliation, ranked  
193 based on the normalised betweenness centrality. This depicts the top 17  
194 central authors in the cluster, prominent members in a group and the first  
195 contact for engaging hypersonics research in this instance. Authors are desig-  
196 nated A through to Q, with the highest betweenness centrality being sourced  
197 from an author at the University of New South Wales at the Australian  
198 Defence Force Academy, now known as the University of New South Wales  
199 at Canberra. The institution has significant influence in the field by being  
200 represented numerous times by authors A, B, E, and I. Followed in equal  
201 abundance by the University of Queensland as authors C, F, G, H, J, and Q.  
202 Notable additions include authors from the European Space Agency, Defence  
203 Science and Technology Group, University of Southern Queensland, RMIT

204 University, and the University of New South Wales at the Sydney based  
 205 campus. Ultimately, these authors represent a select few who are critical or  
 206 influential co-authors in the mapping network, which can be due to various  
 207 reasons. From the standpoint of policy-driven action, these individuals rep-  
 208 resent the first point in canvassing technology initiatives for the keywords  
 209 searched. The ease of constructing clustering maps through bibliographic  
 210 data helps users identify groupings, dominant players in the area, and their  
 211 associated connections.

212

Table 1: Key authors with respect to clustering centrality

	<b>Institution</b>	<b>Papers</b>	<b>Citations</b>	<b>Co-author</b>	<b>Rank</b>
A	University of New South Wales at Australian Defence Force Academy	4	71	31	1
B	University of New South Wales at Australian Defence Force Academy	2	1	36	0.86
C	University of Queensland	7	20	29	0.83
D	University of Queensland	5	35	23	0.71
E	University of New South Wales at Australian Defence Force Academy	7	9	29	0.70
F	University of Queensland	2	3	40	0.68
G	University of Queensland	7	49	33	0.63
H	University of Queensland	2	0	35	0.56
I	University of New South Wales at Australian Defence Force Academy	5	12	17	0.46
J	University of Queensland	2	3	6	0.44
K	University of Southern Queensland	3	2	22	0.42
L	UNSW Sydney	4	18	16	0.42
M	European Space Agency	3	35	13	0.40
N	Defence Science and Technology	3	11	25	0.39
O	Defence Science and Technology	6	119	29	0.36
P	RMIT University	5	36	21	0.29
Q	University of Queensland	2	1	13	0.27

213 In comparison to metrics about geography, the clustering shares specific  
 214 observations. There is an abundance of researchers from New South Wales,  
 215 the Australian Capitol Territory, and Queensland institutions. The clus-  
 216 tering using the centrality metrics also shows that the United States and  
 217 Europe are close contacts, as touched upon in papers by overseas authors.

218 However, the clustering approach with the centrality metrics does not neces-  
219 sarily highlight Singapore or Germany, in the top number of centrally placed  
220 contacts. Furthermore, the map itself shows that the Defence Science and  
221 Technology group tend to interact with the University of New South Wales  
222 and the University of Queensland. This could be an opportunity to interact  
223 with other similar research groups that are presently not directly connected  
224 such as those at the University of Southern Queensland, Griffith Univer-  
225 sity, Monash University, and the Royal Melbourne Institute of Technology.  
226 Since hypersonics is an applied research area, domestic collaboration is less  
227 likely to occur unless framework is set in place to actively support collabo-  
228 ration due to competitive advance. Thus taking an approach that has been  
229 demonstrated in the United States with the clustering of universities under  
230 particular hypersonic research activities and centres of excellence.

231

232

233

### 234 **3. Conclusion**

235 The body of work examines the use of bibliographic clustering to aide in  
236 the canvassing of future technologies for policymakers. Using data-harvesting  
237 of the Scopus journal database and open-source software Gephi, it is readily  
238 possible for policy makers to examine research clusters in any field. This  
239 approach has been demonstrated using a few critical terms associated with  
240 hypersonic propulsion and has shown the dominant authors and institutions  
241 operating in this field. The clustering shows strong co-authorships between  
242 the University of New South Wales and the University of Queensland, indicat-  
243 ing cooperative agreements or strong cross-collaborations. That said, various  
244 other researchers are also picked up through the betweenness centrality mea-  
245 sures, including authors at the Defence Science Technology, European Space  
246 Agency, the University of Southern Queensland, and the Royal Melbourne  
247 Institute of Technology.

248

249 The canvassing demonstrated here also allows more detailed exploration  
250 of the data, which is visually easier to investigate. Thereby increasing pro-  
251 ductivity, and allowing those within the policy to examine the relationship  
252 between authors and minor entities. The work can potentially be applied  
253 to canvassing relevant staff and institutions for tendering, regarding relevant

254 technologies based on initiatives set out through federal and state govern-  
255 ment.

256

257

258 **References**

- 259 [1] H.-j. Lee, S. Lee, B. Yoon, Technology clustering based on evolutionary  
260 patterns: The case of information and communications technologies,  
261 Technological Forecasting and Social Change 78 (2011) 953–967.
- 262 [2] S. Valverde, R. V. Solé, M. A. Bedau, N. Packard, Topology and evo-  
263 lution of technology innovation networks, Physical Review E 76 (2007)  
264 056118.
- 265 [3] R. Baptista, Geographical clusters and innovation diffusion, Technolog-  
266 ical Forecasting and Social Change 66 (2001) 31–46.
- 267 [4] S. Liyanage, Breeding innovation clusters through collaborative research  
268 networks, Technovation 15 (1995) 553–567.
- 269 [5] H. Dundar, D. R. Lewis, Determinants of research productivity in higher  
270 education, Research in higher education 39 (1998) 607–631.
- 271 [6] G. Abramo, C. A. D’Angelo, F. Di Costa, Research collaboration and  
272 productivity: is there correlation?, Higher education 57 (2009) 155–171.
- 273 [7] T. Velden, A.-u. Haque, C. Lagoze, A new approach to analyzing pat-  
274 terns of collaboration in co-authorship networks: mesoscopic analysis  
275 and interpretation, Scientometrics 85 (2010) 219–242.
- 276 [8] R. Molontay, M. Nagy, Twenty years of network science: A bibliographic  
277 and co-authorship network analysis, arXiv preprint arXiv:2001.09006  
278 (2020).
- 279 [9] C. Biscaro, C. Giupponi, Co-authorship and bibliographic coupling net-  
280 work effects on citations, PloS one 9 (2014) e99502.
- 281 [10] S. Murthy, E. Curran, Scramjet propulsion, American Institute of Aero-  
282 nautics and Astronautics, 2001.

- 283 [11] M. Bastian, S. Heymann, M. Jacomy, Gephi: an open source software  
284 for exploring and manipulating networks, in: Proceedings of the Inter-  
285 national AAAI Conference on Web and Social Media, volume 3, 2009.
- 286 [12] F. P. Miller, A. F. Vandome, J. McBrewster, Levenshtein distance:  
287 Information theory, computer science, string (computer science), string  
288 metric, damerau? levenshtein distance, spell checker, hamming distance  
289 (2009).