*Title of the Satellite Meeting:* Knowledge Management with Digital Humanities/Digital Scholarship
*Date: 22 August 2019*
*Location: Corfu, Greece*

# Artificial Intelligence: how knowledge is created, transferred and used

**Maria de Kleijn**
Elsevier, Amsterdam, the Netherlands.
m.dekleijn-lloyd@elsevier.com

**Mark Siebert**
Elsevier, Amsterdam, the Netherlands.

**Sarah Huggett**
Elsevier, Singapore, Singapore

**Abstract:**

*The growing importance and relevance of artificial intelligence (AI) to humanity is undisputed. However, AI does not seem to have a universally agreed definition, and different sectors of society use very different vocabulary to describe AI. Using AI to define AI, we were able to detect the relevant body of research, further structure it in sub-fields, and give a comprehensive overview of the research landscape.*

*There are strong regional differences in AI activity:*

- *China aspires to lead globally in AI and focuses on computer vision. It shows a rapid rise in scholarly output and citation impact. A net brain gain of AI researchers to China also suggests an attractive research environment.*

- *Europe is the largest producer of AI scholarly output, but appears to be losing academic AI talent. The broad spectrum of AI research in Europe reflects the diversity of European countries, each with their own agenda and specialties.*

- *AI research in the United States is robust, both in terms of scholarly output and talent retention. The US benefits from a strong corporate sector. The corpus shows less diversity in AI research than Europe but more than China.*

*A key area of further development in AI research worldwide is on ethical issues pertaining to AI. While a major topic in daily conversation, there is surprisingly little formal research published on AI ethics to date. We believe there is a need for more AI ethics research, which would bring many benefits to the field, its development, and its applications.*

## Put body of the paper here

### 1. Introduction

The growing importance and relevance of artificial intelligence (AI) to humanity is undisputed. AI impacts different parts of society from teaching to our daily lives.[i,ii]

As AI is increasingly seen as one of the real game-changers of today's economic world order,[iii,iv] many governments have formulated AI strategies aimed at creating a solid position that can withstand strong international competition.[v] The national policies we studied all agree that the different sectors in society – education, research, industry and public debate - need to work in unison, for society as a whole to reap the intended benefits. Research plays a driving role in AI development. A key question that pops up is "where do we stand compared to others", in other words, what is the AI research baseline?

To draw any conclusions on a global scale on research focus, intensity, or output in AI, the corpus definition is crucial. This means: which articles, conference proceedings, reviews, or other scholarly publications together make up the body of AI

research? A complicating factor to this from a technical point of view is that, unlike many established research fields, experts agree that AI lacks an agreed-upon definition, scope, and ontology[vi].

In addition, the close link between today's research and its application areas[vii] suggests that a significant volume of AI research may be done outside the traditional "home" of AI in Computer Science. AI research and AI-enabled research may be published in non-AI journals, and may not even use some of the typical computer science vocabulary. For these reasons, simple bibliometric approaches to define the AI research corpus – like looking only at AI journals, or selecting articles with the keyword "AI"- will not yield comprehensive enough results.

Instead, allowing insights across research fields and application areas, we used AI to delineate AI – bottom-up from the article level.

### 2. Method: using AI to define AI

We mined the text of representative textbooks, Massive Open Online Course (MOOC) syllabi, books, research outputs, patents, and news items to identify meaningful concepts using the Elsevier FingerPrint Engine (EFE).[viii] The choice to include sources outside of the scholarly arena is based on the realisation that different parts of society have a stake in AI, from teaching to research to innovation.
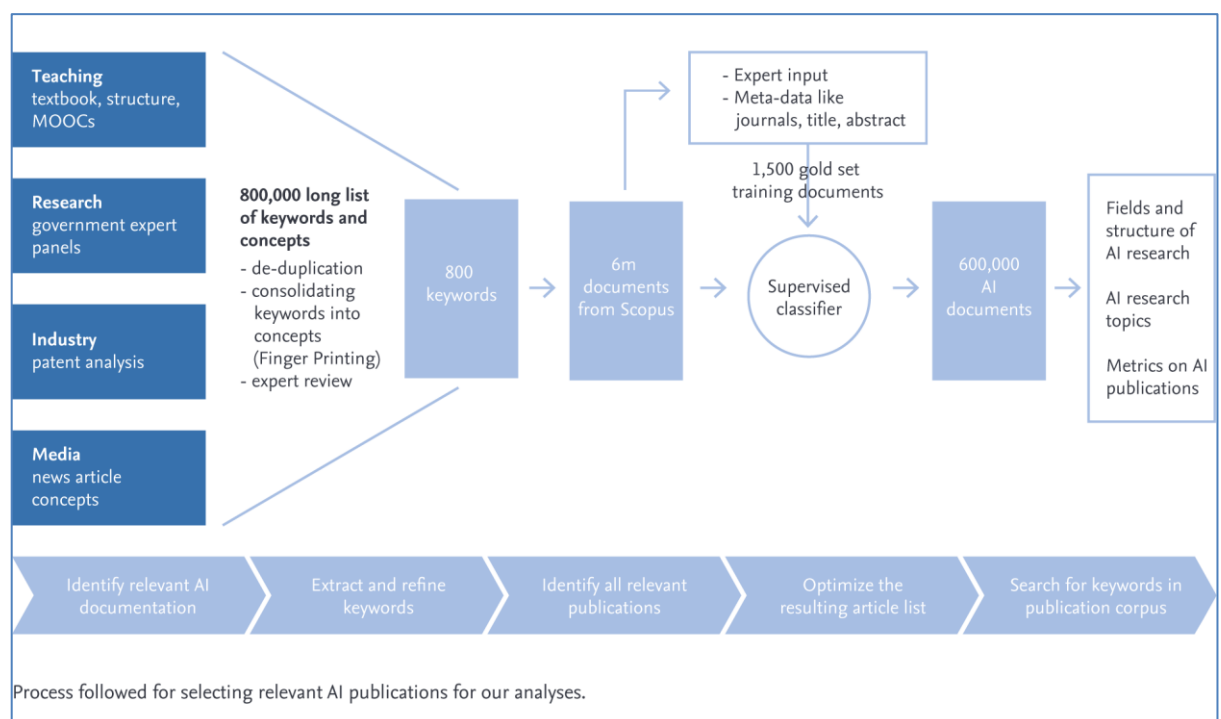
With manual expert feedback combined with automated steps, we reduced our initial mined longlist of 800,000 keywords down to 797 keywords and concepts.

We then searched for each keyword and concept in the titles, abstracts, and keywords of documents included in a Scopus May 2018 dataset, retrieving 5.7 million unique documents, with many false positives. These false positives were caused by

application terms (e.g. "finite elements"), broad terms (e.g. "ethical values"), or terms that also have relevance in other research fields (e.g. "neural networks" in biology).

To eliminate false positives from the corpus while retaining relevant AI papers we used supervised Machine Learning with further expert input on the training data set. The 797 concepts were ranked on a scale of High/Medium/Low with regards to relevance to the core field of AI and assigned a respective weight. In parallel, 1,500 documents were manually classified by internal experts as either "AI" or "not AI" to use as reference and training input for the algorithm to determine the classification. Once the model was trained and validated, it could identify AI papers with 85% precision. The complete set of 5.7 million documents was run through the model to generate predictions that were used to reduce the number of documents identified as AI Research to some 600,000 documents.

The full workflow is pictured below, in Figure 1.



*Figure 1. Process followed for selecting relevant AI publications.*

### 3. Method: scientometric analyses

Scopus® is Elsevier's abstract and citation database of peer reviewed literature, covering 71 million documents from more than 23,700 active journals, book series, and conference proceeding papers by 5,000 publishers. In the Scopus database each document has one or more authors and is linked to one or more institutions. In addition, the citation links between Scopus documents, and citing links to other databases like patent databases or news sources, are available. This makes a full range of scientometric analyses on the identified AI corpus possible. Key metrics are explained below, according to their published definitions[ix] and based on full counting.

**Output** (of an institution or country): the count of publications with at least one author from the author byline listing that institution/country as an affiliation.

**FWCI**: Field-Weighted Citation Impact (FWCI) is an indicator of the citation impact of a publication. It is calculated by comparing the number of citations actually received by a publication with the number of citations expected for a publication of the same document type, publication year, and subject. Field-Weighted Download Impact (**FWDI**) is a similar indicator of usage based on full text downloads.

For the **mobility analysis** included in this report, we use the affiliations that authors specify on their publications. We define active researchers as belonging to a specific mobility class as follows:

- **Sedentary**: active researchers whose Scopus author data for the period indicates that they have not published outside the region.

- **Transitory**: active researchers whose Scopus author data for the period indicates that they have been abroad or in the region for less than two years.

- **Inflow**: researchers whose publication history indicates that they first published outside of the region and then published inside of the region.

- **Outflow**: researchers whose publication history indicates that they first published inside the region and then published outside of the region.

In the mobility analysis we refer to **relative productivity**. This is a measurement of the number of publications per year since the first appearance of each researcher as an author during the analytical period, relative to all regional researchers in the same period.
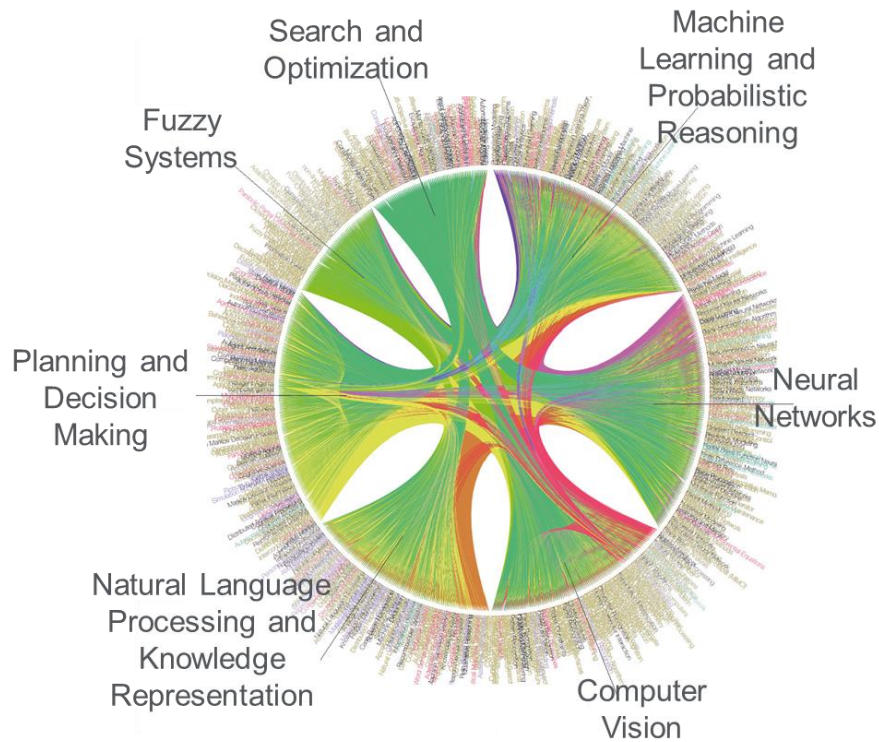
### 4. Semantic results

A first, striking observation was the limited overlap between the mined keywords from the four different perspectives: teaching, media, research, and industry. These different stakeholder groups have different ways to speak about AI, and partly focus on different aspects. Media for instance puts emphasis on the "personality" and the physical forms of AI, while teaching provides broad overviews of approaches, architectures, and tools. The overlap in keywords between each perspective is depicted in Figure 2.

Teaching
268

83

Media
82

3

52

153

6

1

17

10

444

Industry
641

Research
42

Keyword mapping (number of keywords) between AI perspectives

*Figure 2. Keyword mapping between four AI perspectives.*

Second, we applied unsupervised machine learning to cluster the 600,000 AI documents into subfields, to understand the AI research fields and focus of different geographies. To this end, we generated a matrix of co-occurring keywords in documents in our AI corpus. This was then transformed into a network structure and processed with a Louvain clustering algorithm, identifying closeness between keywords. A chord graph visualisation of the global result is depicted in Figure 3, with manually labelled keyword clusters.

*Figure 3. Clustering of AI, based on keyword co-occurrences of at least 500 co-occurring documents, 2017.*

A key learning from this analysis is that it is neither stable over time, nor identical across geographies. Rather, it shows how AI as a research field evolves. In the remainder of our analysis we have used clusterings such as the above to understand differences in research focus between geographies and over time.
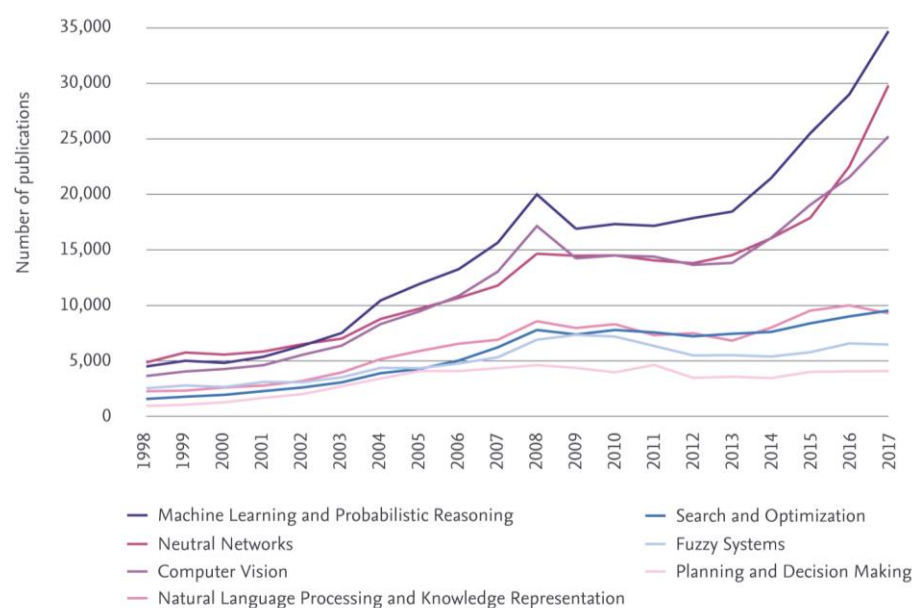
## 5. The AI research landscape

### 5.1 Global trends

Global AI research output has grown far quicker than the world average for all research; growth has accelerated in recent years to more than 12% per annum.

When we split the global output by subfield as in Figure 4, we see that the rapid growth in recent years is largely due to 3 clusters: machine learning and probabilistic reasoning, neural networks, and computer vision.
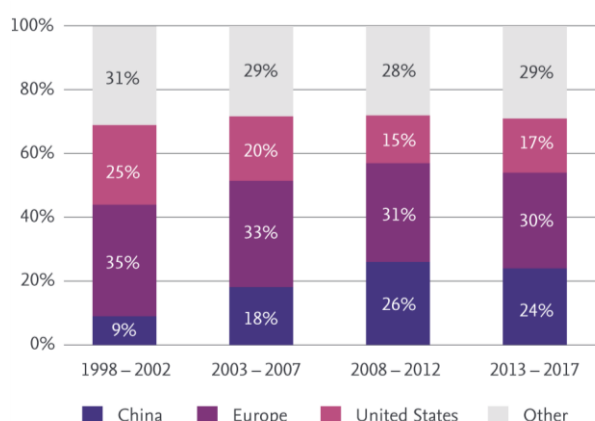
*Figure 4. Annual number of AI publications by keyword co-occurrence cluster (all document types), 1998-2017; sources: Scopus and Elsevier clustering.*
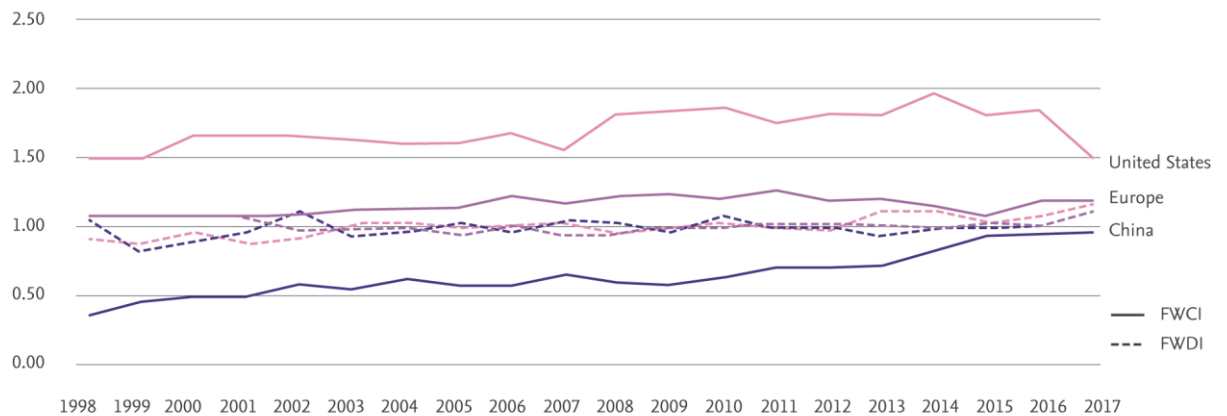
When comparing the size of the research enterprise across the world, by share of world publications as in Figure 5, it becomes apparent that Europe – defined as EU44, the countries fully eligible under the Horizon programme[x] and its

predecessors – is the largest contributor to the world's output, but also gradually losing share. China shows significant growth, although the 2008-2012 figure is somewhat inflated due to an unusually high number of conferences.



*Figure 5. Share of global publication output in AI (all document types) for periods 1998-2002, 2003-2007, 2008-2012, and 2013-2017, per region; source: Scopus*

Further comparing China, the US, and Europe by field-weighted citation impact (rebased to 1.0 within the AI corpus), as in Figure 6, we see that while citation impact of the US and Europe is largely stable, China has managed a strong increase in FWCI to reach the world average, but remains below the US and Europe in this indicator.



*Figure 6. Rebased AI Field-Weighted Citation Impact and Field-Weighted Download Impact per region, 1998-2017; source: Scopus*

There are strong regional differences in AI activity. The next section zooms in on some key findings per region.

## 5.2 Regional findings - China

China aspires to lead globally in AI and is supported by ambitious national policies.[xi] China's AI focuses on computer vision, as is shown in Figure 7 below. It shows robust growth of its research ecosystem, with a rapid rise in scholarly output. China's AI research has a rapidly increasing yet still comparatively low citation impact (Figure 6 above), which could be a symptom of regional, rather than global, reach.
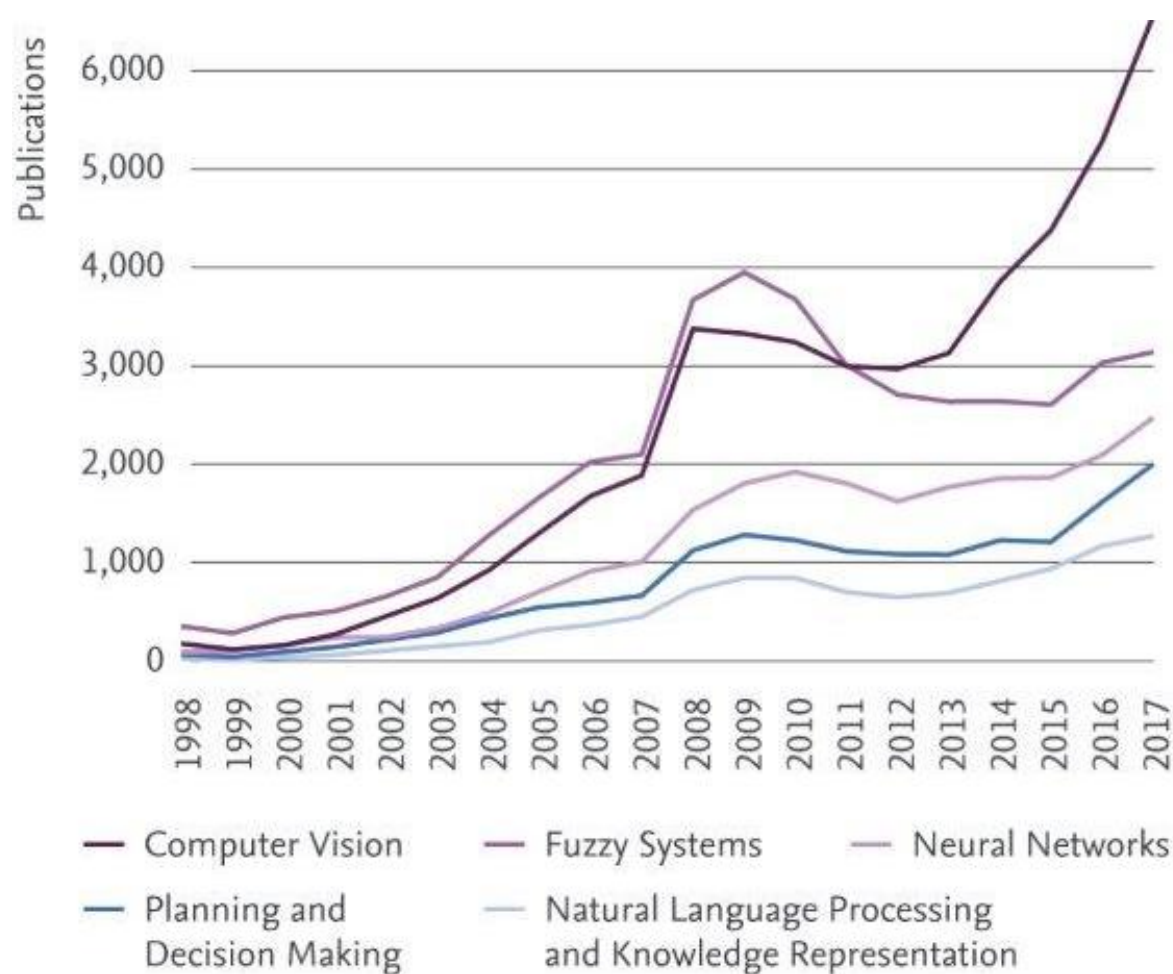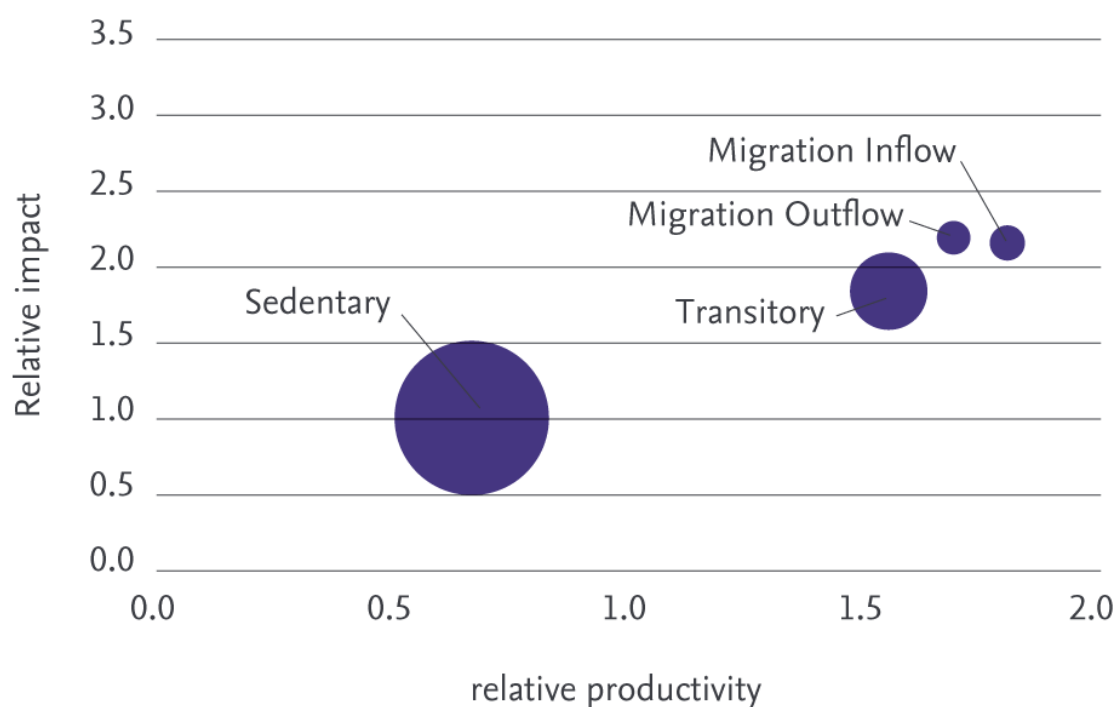


*Figure 7. Annual publications per cluster for China (all document types), 1998-2017; sources: Scopus and Elsevier clustering.*

A net brain gain of AI researchers in China (Figure 8) also suggests an attractive research environment. While the share of transitory and inflow researchers is modest, they do have comparatively high citation impact and productivity in terms of number of publications per author.



*Figure 8. Relative productivity and relative citation impact per mobility class for China, 1998-2017; bubble size represents the percentage of researchers in each mobility class. Source: Scopus*

### 5.3 Regional findings – Europe

Europe is the largest region in AI scholarly output, with high and rising levels of international collaborations outside of Europe. The broad spectrum of AI research in Europe reflects the diversity of European countries, each with their own agenda and specialties. Focus areas of European AI research include genetic programming for pattern recognition, fuzzy systems, and speech and face recognition.
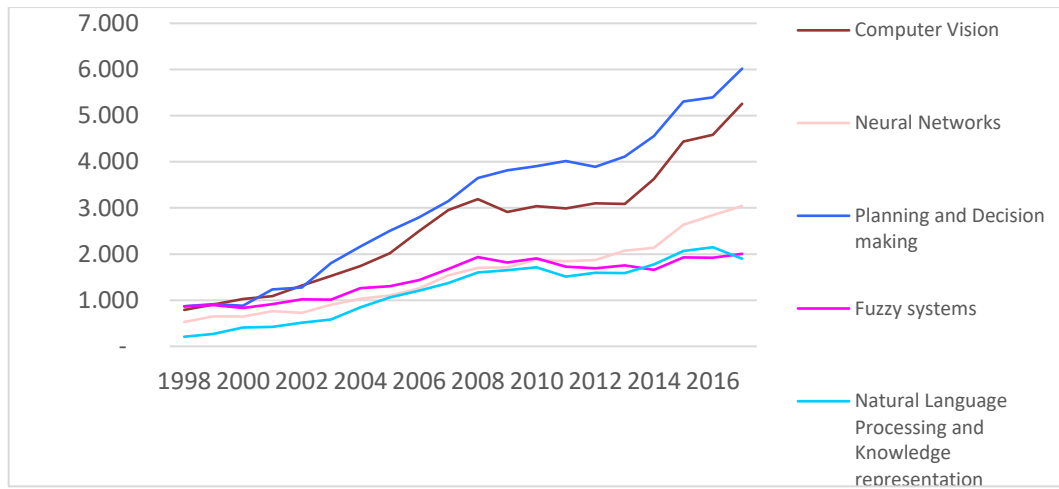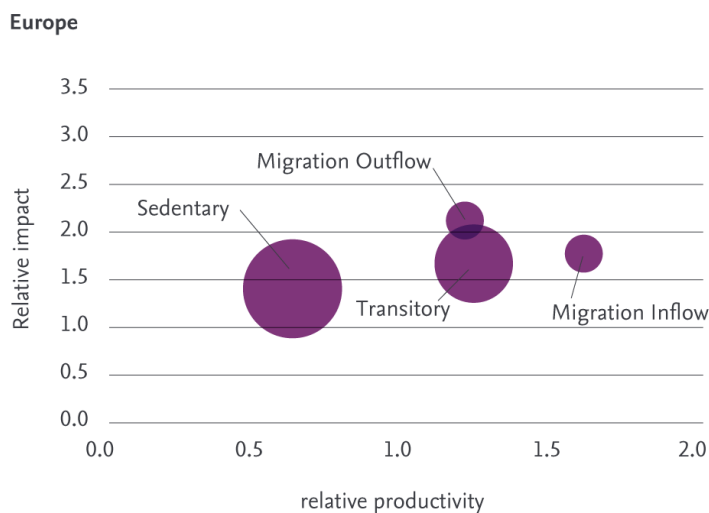
*Figure 9. Annual publications per cluster for Europe (all document types), 1998-2017; sources: Scopus and Elsevier clustering.*
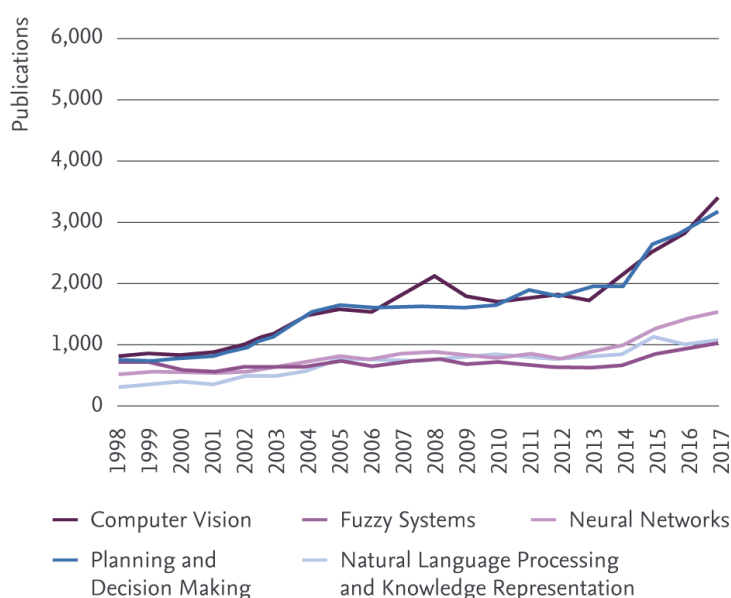
A key concern in Europe is its ability to retain top researchers.[xii] Our analysis, depicted in Figure 10, suggests a modest but steady outflow of researchers, often with high citation impact. However, Europe is clearly also able to attract top researchers from across the world for shorter or longer stays.



*Figure 10. Relative productivity and relative citation impact per mobility class for Europe, 1998-2017; bubble size represents the percentage of researchers in each mobility class. Source: Scopus*
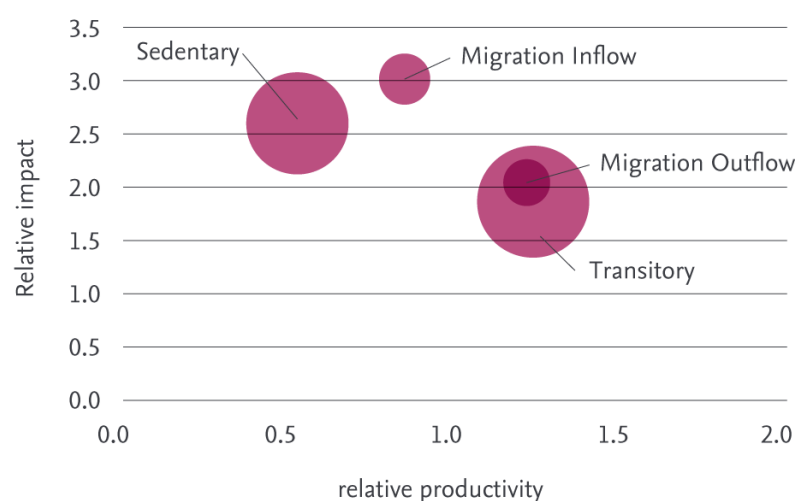
## 5.4 Regional findings – US



Computer Vision, Fuzzy Systems, Neural Networks, Planning and Decision Making, Natural Language Processing and Knowledge Representation

AI in the United States has a strong focus on specific algorithms and separates speech and image recognition into distinct clusters. The corpus shows less diversity in AI research than Europe but more diversity than China, as shown in Figure 11.

*Figure 11. Annual publications per cluster for the US (all document types), 1998-2017; sources: Scopus and Elsevier clustering.*

The mobility profile in the US (Figure 12) is interesting: sedentary researchers, while less productive than their mobile peers in terms of average number of



publications, show very high citation impact. This suggests that the US academic sector itself is both highly competitive, and able to retain top researchers.

*Figure 12. Relative productivity and relative citation impact per mobility class for the US, 1998-2017; bubble size represents the percentage of researchers in each mobility class. Source: Scopus*

**6. Suggested further research**

This case study on AI shows how to scope, and make sense of, a strongly evolving research field. However, the method can be further refined – for instance, sensitivity to specific keywords, the size of the training set, and so on. In addition, stronger links between research and business or society applications and implications are not yet widely understood.

Another area that warrants further study is AI ethics. Expert interviews and national AI policies confirm ethics being a key area of concern for AI worldwide. While a major topic in daily conversation, in media, public policy, [xiii, xiv] and between researchers, there is surprisingly little formal research published on the topic to date.[xv] Issues around algorithmic bias, governance, accountability, impersonation, privacy, for instance, leave a world of future research opportunities.

## References

[i] Schwab, K. (2016) *The Fourth Industrial Revolution*. Currency. Available at: https://www.weforum.org/about/the-fourth-industrial-revolution-by-klaus-schwab.

[ii] Adams, R. L. (2017) *10 Powerful Examples of Artificial Intelligence in Use Today*, *Forbes*. Available at: https://www.forbes.com/sites/robertadams/2017/01/10/10-powerful-examplesof-%0Aartificial-intelligence-in-use-today/#5590a7c9420d (Accessed: 19 November 2018).

[iii] Lee, K.-F. (2018) *AI Superpowers: China, Silicon Valley, and the New World Order*. New York: Houghton Mifflin Harcourt. Available at: https://aisuperpowers.com/about/about-the-book.

[iv] Manyika, J. *et al.* (2017) *Jobs Lost, Jobs Gained: Workforce Transitions in a Time of Automation_Mckinsey Global Institute*, *McKinsey Global Institute*. Available at: https://www.mckinsey.com/~/media/McKinsey/Featured Insights/Future of Organizations/What the future of work will mean for jobs skills and wages/MGI-Jobs-Lost-Jobs-Gained-Report-December-6-2017.ashx.

[v] Dutton, T. (2018) *An Overview of National AI Strategies*, *Medium*. Available at: https://medium.com/politics-ai/an-overview-of-national-ai-strategies-2a70ec6edfd (Accessed: 19 November 2018).

[vi] Marr, B. (2018) *The Key Definitions Of Artificial Intelligence (AI) That Explain Its Importance*. Available at: https://www.forbes.com/sites/bernardmarr/2018/02/14/the-key-definitions-of-artificial-intelligence-ai-that-explain-its-importance/#3bb229374f5d.

[vii] Cockburn, I. M., Henderson, R. and Stern, S. (2017) *The Impact of Artificial Intelligence on Innovation: An Exploratory Analysis*. Available at: http://www.nber.org/chapters/c14006.

[viii] The Elsevier Fingerprint Engine™ identifies concepts and their importance in any given text by using a wide range of thesauri and data-driven controlled vocabularies covering all scientific disciplines, and by applying a variety of Natural Language Processing (NLP) techniques. The resulting terms are of high quality and more representative than standard sets of keywords, which often contain duplicates, synonyms, and inclusion of irrelevant terms.

[ix] Elsevier (2018), Research metrics guidebook. Available at https://www.elsevier.com/research-intelligence/resource-library/research-metrics-guidebook

[x] The Framework Programmes for Research and Technological Development, also called Framework Programmes or abbreviated FP1 to FP7 with "FP8" being named "Horizon 2020" and "FP9" to "Horizon Europe", are funding programmes created by the European Union/European Commission to support and foster research in the European Research Area (ERA). See https://ec.europa.eu/info/designing-next-research-and-innovation-framework-programme/what-shapes-next-framework-programme_en

[xi] China Copyright and Media (2017) *A Next Generation Artificial Intelligence Development Plan. Retrieved from The law and policy of media in China*. Available at:

https://chinacopyrightandmedia.wordpress.com/2017/07/20/a-next-generation-artificial-intelligence-development-plan/.

xii Sample, I. (2017) 'Big tech firms' AI hiring frenzy leads to brain drain at UK universities', *The Guardian*, 2 November. Available at: https://www.theguardian.com/science/2017/nov/02/big-tech-firms-google-ai-hiring-frenzy-brain-drain-uk-universities.

xiii European Group on Ethics in Science and New Technologies (2018) *Statement on Artificial Intelligence, Robotics and 'Autonomous' Systems*. Luxembourg. Available at: https://ec.europa.eu/research/ege/pdf/ege_ai_statement_2018.pdf.

xiv Jirotka, M. *et al.* (2017) 'Responsible Research and Innovation in the Digital Age', *Communications of the ACM*, 60(5), pp. 62–68. doi: 10.1145/3064940.

xv Stahl, B. C., Timmermans, J. and Mittelstadt, B. D. (2016) 'The Ethics of Computing: A Survey of the Computing-Oriented Literature', *ACM Computing Surveys (CSUR)*, 48(4), p. 38

xvi Elsevier (2018) 'Artificial Intelligence: How knowledge is created, transferred, and used'. Available at https://www.elsevier.com/connect/resource-center/artificial-intelligence