

# An Experimental Study on the Dynamic Evolution of Core Documents

Lin Zhang<sup>1</sup>, Wolfgang Glänzel<sup>2</sup>, Fred Y. Ye<sup>3</sup>

<sup>1</sup>zhanglin\_1117@126.com

<sup>1</sup>Dept. Management and Economics, North China University of Water Conservancy and Electric Power, Zhengzhou (China)

<sup>2</sup>Wolfgang.Glanzel@kuleuven.be

<sup>2</sup>Centre for R&D Monitoring (ECOOM) and Dept. MSI, KU Leuven, Leuven (Belgium)  
Dept. Science Policy & Scientometrics, Library of the Hungarian Academy of Sciences, Budapest (Hungary)

<sup>3</sup>yey@nju.edu.cn

<sup>3</sup>School of Information Management, Nanjing University, Nanjing 210023, (China)  
Jiangsu Key Laboratory of Data Engineering and Knowledge Service, Nanjing 210023 (China)

## Introduction

The concept of the core of documents had originally been introduced in connection of co-citation analysis (Small 1973). The term *core documents* has later been re-introduced in the context of bibliographic coupling (BC; see Glänzel & Czerwon, 1996) and hybrid BC and text based similarities (Glänzel & Thijs, 2011) in order to identify strongly interlinked papers that form important nodes in the network of scholarly communication. In order to study stability and dynamics of core-document sets we apply two different methods to h-index related literature in the period 2005–2013 for illustration.

## Data Sources and Processing

Data were retrieved from Thomson Reuters Web of Science Core Collection (WoS) following the strategy of Zhang et al. (2011), with extension of the period 2005–2013. We also added citing papers but removed duplicates and papers with less than 5 references to avoid biases in BC similarities. We obtained a final set of 3,270 documents. Figure 1 shows the annual increment of papers in this set.

## Research Questions, Methods and Results

In this study we apply two different methods to determine core documents, (Method I) the traditional one according to Glänzel & Czerwon (1996) with a fixed number of links ( $n = 15$ ) and Method II using the h-core of the network (Glänzel, 2012). In both cases we applied a *hybrid approach*. We used link strengths of 0.5 and 0.4 according to Salton's cosine measure. Using these parameters, we analysed the dynamics of core documents along the following questions.

- How is evolution of core documents reflected by the two methods?
- Do the two methods provide stable results?
- Do core documents adequately represent the evolution of the topic?

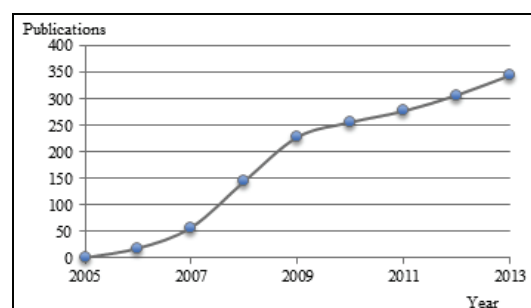


Figure 1. Distribution of h-related publications during 2005–2013.

Core document are by definition strongly interlinked with a large number of other documents in the set under study and thus represent the very core of the set. As expected, their number increases with expanding time spans, the average annual growth rate of the cumulative set amounted to 46% (Method I) and 25% (Method II), respectively. Not only the number of nodes in the network but also the number of their links is growing, however at a different pace. Indeed, we found that the complete h-related set increased at a large constant pace of 11% while the growth of the core sets was faster (see above), but its growth slowed down. This might in part be a consequence of the increasing age of references. In 2013 the core reached a representation of 2.0% and 2.4%, respectively. This characterizes the evolution of the core set with respect to the topic dynamics. The second question that arises from these figures is in how far do both methods mirror the same “core” of literature. In order to check the robustness of these methods, we compared the overlap of the sets of core documents obtained from the two methods. To this end we used BC with fixed number of links as reference standard. Concordance with Method I ranged between 83.8% and 95.2% with increasing trend from 2005–2007 to 2005–2013 and using Method II the shares ranged between 96.8% and 80.7%, however with decreasing trend.

Figure 1 displays three network graphs illustrating the evolution of research on the relationship between the immune system and aging, categorized by time periods: 2005-2007, 2008-2010, and 2011-2013. The nodes represent authors, and the edges represent research connections.

**2005-2007:** This graph shows a dense network of authors. Key authors include Rau\_2007, Imperial\_et\_al\_2007, Saad\_2006, Bommann\_et\_al\_2007, Glanzel\_2006, Liang\_2006, Van\_Raan\_2006, Totto-Alves\_et\_al\_2007, Burrell\_2007, Schreiber\_2007a, Schreiber\_2007b, Jiri\_et\_al\_2007, Ferrand\_2007, Eggho\_2006, Eggho\_2007a, Eggho\_2007b, Frangopol\_2005, Chapron\_et\_al\_2005, Braun\_et\_al\_2006, and Eggho\_et\_al\_2006.

**2008-2010:** This graph shows a more dispersed network. Key authors include Woeginger\_2008a, Kulasgarah\_et\_al\_2010, Demeke\_et\_al\_2009, van\_Leeuwen\_2008, Dorta-Gonzalez\_et\_al\_2010, Woeginger\_2008b, Eggho\_et\_al\_2008b, Molinari\_et\_al\_2008, van\_Eck\_2008, Burrell\_2009, Gogolewski\_et\_al\_2008, Eggho\_2008a, Eggho\_2008b, Liu\_YX\_et\_al\_2009, Schreiber\_2008, Eggho\_et\_al\_2008a, Eggho\_2009a, Arenobita-Jorge\_et\_al\_2008, Zhang\_2010, Tol\_2009, Rousseau\_2008, and Schreiber\_2009.

**2011-2013:** This graph shows a highly interconnected network. Key authors include Badran\_et\_al\_2013, Slidder\_et\_al\_2013, Rezek\_et\_al\_2011, Prathap\_2011b, Prathap\_2011c, Prathap\_2011d, Prathap\_2011a, Cargiol\_et\_al\_2012, Quesada\_2011a, Quesada\_2011b, Zhang\_2013, Kuan\_et\_al\_2011b, Kuan\_et\_al\_2011a, Kuan\_et\_al\_2012, Liu\_et\_al\_2013, Eggho\_2011b, Eggho\_2011a, Eggho\_2011c, Eggho\_2012, Dorta-Gonzalez\_et\_al\_2011, and Eggho\_2013.

Core nodes in Figure 2 are based on BC but hybrid similarities are used to measure the links between the nodes. This can be done because of the strong concordance between the sets obtained from the two methods. The links between core nodes in Figure 2 are denser and stronger than in the BC approach, which is due to the inclusion of textural information. The interpretation of Figure 2 is not straightforward, but the structural changes of the networks during different periods presented here are quite clear and noteworthy. The network in the first sub-period (2005–2007) comprises above all theoretical publications. The network of 2008–2010 already reflects a different picture. While most theoretical papers are still located in the centre of the network, also ‘applied studies’ started to appear in the core-documents set. These are distributed at the periphery of the network, which indicates that the topic starts to expand from pure theory to more

Batagelj, V., & Mrvar, A. (2003). *Pajek-Analysis and visualization of large networks*. In M. Jünger & P. Mutzel (Eds.), *Graph drawing software* (pp. 77-103). Berlin: Springer.

Glänzel, W., & Czerwon, H. J. (1996). A new methodological approach to bibliographic coupling and its application to the national, regional and institutional level. *Scientometrics*, 37(2), 195–221.

Glänzel, W. & Bart Thijs, B. (2011). Using ‘core documents’ for the representation of clusters and topics. *Scientometrics*, 88(1), 297-309.

Glänzel, W. (2012). The role of core documents in bibliometric network analysis and their relation with h-type indices. *Scientometrics*, 93(1), 113–123.

Small, H. (1973). Cocitation in scientific literature – new measure of relationship between 2 documents. *JASIS*, 24(4), 265-269.

Zhang, L., Bart Thijs, B. & Glänzel, W. (2011). The diffusion of H-related literature. *Journal of Informetrics*, 5(3), 583-593.