

1

[gerasimov@mlab.cs.msu.su](mailto:gerasimov@mlab.cs.msu.su), [romaha@mlab.cs.msu.su](mailto:romaha@mlab.cs.msu.su), [mash@cs.msu.su](mailto:mash@cs.msu.su),  
[michael@cs.msu.su](mailto:michael@cs.msu.su), [tsarev@mlab.cs.msu.su](mailto:tsarev@mlab.cs.msu.su), [andy@mlab.cs.msu.su](mailto:andy@mlab.cs.msu.su)

« »

1

11-07-00616, 12-07-00585.

14.514.11.4016

1.

( - , , - )  
[[1]].

« » [[2]].  
Thomson Reuters ( Institute for Scientific Information, ISI): Science Citation Index, Social Sciences Citation Index Arts & Humanities Citation Index, Journal Citation Reports (JCR) [[3]].

CiteSeer, 1997 [[4]].

« » ( ) [[5], [6]] — « »  
[[7], [8]].

, 2008–2011 .  
(European 7th Framework Programme) European Educational Research Quality Indicators (EERQI) [[9]],

( , , ).  
« » , [[10]].

EERQI , « », « »  
[[11], [12]]

« », « » /

2.

( , )

,

-

-

: « » (

); « » (

); « » (

); « » (

); « » (

).

2

3

4

5

« » ,

( . . . ) [[2]].

(

« »

Thomson Scientific (Science Citation

Index, Social Sciences Citation Index Arts & Humanities Citation Index, Journal Citation Reports) [[3]], CiteSeer [[4]] Google Scholar [[13]].

(« » )

( , , , , , « » ).

2005 .

( )

[[5], [6]].

[[7], [8]]:

- « » ( , );

- « » ( , ),

( « » ),

2008–2011 . (European 7th Framework Programme) EERQI [[9]],

(

).

« » ,

:

- « » ( ) —

:

( . *rigour*) —

- ( . *originality*) —

( . *significance*) —

— « » ( , . ) —

Google Web Search ( Google Scholar, MetaGer):

- ( ),
- ( ),
- ,
- g- ( h- ),
- ( Google Web Search MetaGer).

( citeulike ([www.citeulike.org](http://www.citeulike.org)), LibraryThing ([www.librarything.com](http://www.librarything.com)), Connotea ([www.connotea.org](http://www.connotea.org)), Mendeley ([www.mendeley.com](http://www.mendeley.com))).

, EERQI

(« » ),

« »

[[11], [12]],

», « / ».

«

### 3.

—

— « » — ,

— « » — ;

— « » — ,

);

— « » —

— « » — ,

.

,

#### 3.1 « »

[[14], [15], [16], [17]].

( . *Latent semantic analysis, LSA*) [[14], [15], [16]].

“*bag-of-words*” [[14]].

*Value Decomposition, SVD*) [[14]]

( . *Non-negative Matrix Factorization, NMF*) [[15]].

- 1.
- 2.

(  
).  
[[18], [19]].

) [[20], [21], [22], [23]].

- 3.

1. : ( )  
,  
[19].  
( )  
[[20]],

10% 100% 10%,  
10

2. ( )

[[15], [24]].

### 3.2 « » « »

« » « » 1 ( )  
« »).  
« »  
« » [[6]].  
.1

	PageRank	
	PageRank	


PageRank

PageRank,

[[25]],

( . .

).

:

- ;

- , 0,25- );

- (0,5- );

- ( , 0,75- );

- .

### 3.3

*chunking*).

(

[[29]]),

n-

n-

[[27]]

*String Sequence Kernel (SSK)* [[28]],

$[k, n]$  (  $k$   $n$  — ),

### 3.4 « »

Web-  
( . *readability*) , ... « »  
[[30], [31], [32]].

« » ( ) ( ) [[32]].

[[30]] , [[30], [31]].

« » .

1. ( ) ;

2. ( ) ;

3. 7 ;

4. - ( , , ) ;

5. - ;

6. ;

7. («!», «?», «...»);

8.

### 3.5

1. —

, ...

[[33]].

2. ( ,

[[34]] ) [[35]] (

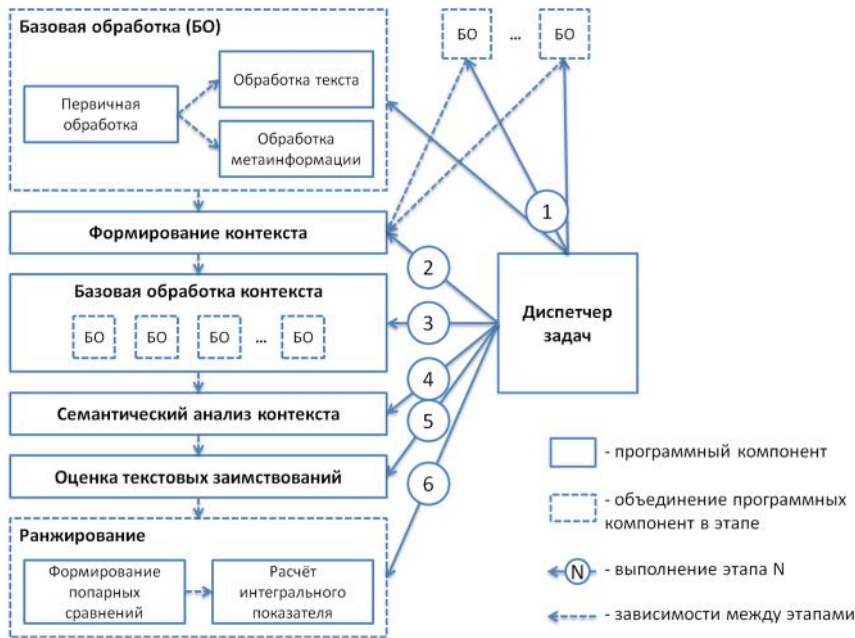
3. -

### 4.

( ) , -

1. «...» ( ) —  
:  
- :  
( )  
, ;  
- :  
- :  
;
2. «...»  
:  
, ,  
( )  
- );
3. «...»  
—
4. «...»  
:  
( ), —  
,
5. «...»  
—

6. «...»  
:  
-  
- ( )  
)  
, ( )  
«...»;  
( / )  
«...»  
«...»  
( )  
«...»  
«...»  
«...»  
( )  
«...»  
«...»  
( )  
( )
1. -
- 2.
3. ,
4. .
- 5.
- 6.



. 1.

Twisted Framework [[36]].

Twisted,

C++

Python

Twisted,  
C++

2.4 ) MS Windows.

Unix ( Linux

JavaScript

[[37]], Ajax-  
« » - gooxdoo  
Explorer 6+, Firefox 2+, Opera 9+, Safari 3.0+ Chrome 2+ : Internet

### 5.

« »  
2011 . [[38]]  
« »  
2012 . ;  
ICDM 2006 [[39]] ICMET 2011  
[[40]].  
10  
« » « » 5  
« »  
« »  
1 2»,  
5  
« » 5  
25 « » 5 « », ( ).

Ranking Loss [[41]],



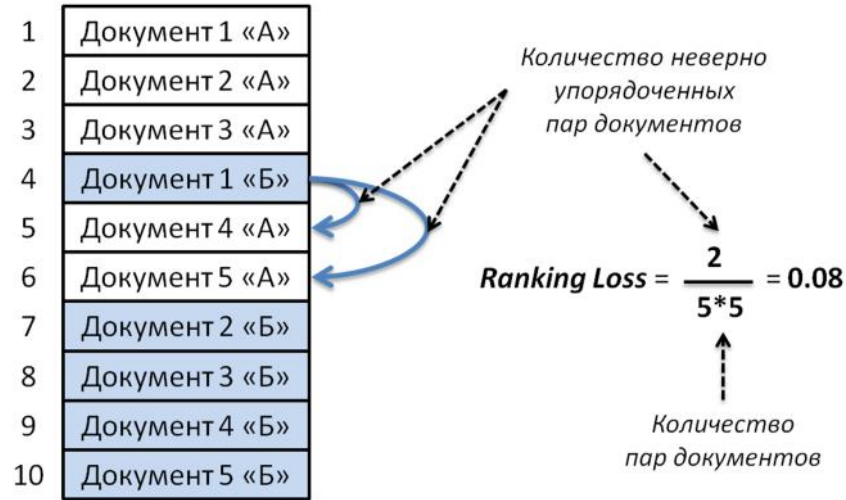
( .2).

« »  
Ranking Loss,

Ranking Loss = 0.

: 100\*(1 –

Ranking Loss).



.2.

Ranking Loss

20%  
98.8%.

4-5

16-

6.

7.

- [1]. Steve Lawrence, Kurt Bollacker, C . Lee Giles. Indexing and Retrieval of Scientific Literature // Eighth International Conference on Information and Knowledge Management, CIKM 99, Kansas City, Missouri, November 2–6, pp. 139–146, 1999.
- [2]. . . . , 2007, N1, . 128-140. //
- [3]. ISI Web of Knowledge ( Healthcare & Science business Thomson Reuters) // <http://www.webofknowledge.com>.
- [4]. CiteSeer // <http://citeseerx.ist.psu.edu>.
- [5]. // [http://elibrary.ru/project\\_risc.asp](http://elibrary.ru/project_risc.asp).
- [6]. . . . //

». — ., 2005.

- [7]. Meho L (Meho, Lokman); Yang K (Yang, Kiduk). Fusion approach to citation-based quality assessment // Proceedings Of Issi 2007: 11th International Conference Of The International Society For Scientometrics And Informetrics, Vols I And II : 568-581.

- [8]. Angela Vorndran, Alexander Botte. Analysis and evaluation of existing methods and indicators for quality assessment of scientific publications // [http://www.eerqi.eu/sites/default/files/Analysis\\_and\\_evaluation\\_of\\_existing\\_methods\\_and\\_indicators.pdf](http://www.eerqi.eu/sites/default/files/Analysis_and_evaluation_of_existing_methods_and_indicators.pdf) [PDF].
- [9]. EERQI – European Educational Research Quality Indicators // [www.eerqi.eu](http://www.eerqi.eu).
- [10]. EERQI Project Final Report (2011) // [http://eerqi.eu/sites/default/files/Final\\_Report.pdf](http://eerqi.eu/sites/default/files/Final_Report.pdf) [PDF].
- [11]. Moyses Szklo. Quality of scientific articles // Revista Saúde Pública vol.40 special issue São Paulo Aug. 2006.
- [12]. Dr Navneet Gupta BSc (Hons) PhD MCOptom FBCLA. How to Evaluate a Scientific Research Article // <http://www.optometry.co.uk/uploads/articles/ARTICLE%200309.pdf> [PDF].
- [13]. Google Scholar// <http://scholar.google.ru>.
- [14]. Berry M.W., Dumais S.T., O'Brien G.W. Using Linear Algebra for Intelligent Information Retrieval // University of Tennessee Knoxville. TN. USA, 1994.
- [15]. Lee D.D., Seung H.S. Learning the parts of objects by non-negative matrix factorization // Nature, 401, pp. 788-791, 1999.
- [16]. Rakesh P., Shivapratap G., Divya G., Soman KP. Evaluation of SVD and NMF Methods for Latent Semantic Analysis // International Journal of Recent Trends in Engineering, Vol. 1, No. 3, 2009.
- [17]. Griffiths T L, Steyvers M. Finding scientific topics // In: Proceedings of the National Academy of Sciences. USA, 2004, 101: 5228–5235.
- [18]. Steinberger J., Ježek K. Text Summarization and Singular Value Decomposition // In Lecture Notes for Computer Science vol. 2457, Springer-Verlag, 2004, pp. 245-254.
- [19]. Steinberger J. Text Summarization within the LSA Framework // PhD Thesis, University of West Bohemia in Pilsen, Czech Republic, January 2007.
- [20]. // 14, 2013. 91-102.
- [21]. Mashechkin I.V., Petrovskiy M.I., Popov D.S., Tsarev D.V. Automatic text summarization using latent semantic analysis // Programming and Computer Software, pp. 299-305, 2011.
- [22]. Tsarev D., Petrovskiy M., Mashechkin I. Using NMF-based text summarization to improve supervised and unsupervised classification // 11th International Conference on Hybrid Intelligent Systems (HIS), Malacca, MALAYSIA. P. 185-189, 2011.
- [23]. Dmitry Tsarev, Mikhail Petrovskiy and Igor Mashechkin, Supervised and Unsupervised Text Classification via Generic Summarization International Journal of Computer Information Systems and Industrial Management Applications. MIR Labs, Volume 5, 2013, pp. 509-515.
- [24]. Wei Xu, Xin Liu, Yihong Gong Document clustering based on non-negative matrix factorization // Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval, Toronto, Canada, 2003.
- [25]. Y. Ding. Applying weighted PageRank to author citation networks. In Proceedings of JASIST. 2011, pp. 236-245.
- [26]. M. Potthast, T. Gollub, M. Hagen, J. Graßegger, J. Kiesel, M. Michel, A. Oberländer, M. Tippmann, A. Barrón-Cedeño, P. Gupta, P. Rosso, B. Stein. Overview of the 4th International Competition on Plagiarism Detection. CLEF2012. 2012.
- [27]. S. Alzahrani, N. Salim. Fuzzy Semantic-Based String Similarity for Extrinsic Plagiarism Detection, Lab Report for PAN at CLEF2010, 2010.
- [28]. A. Martins. String kernels and similarity measures for information retrieval. 2006.
- [29]. Berry M.W., Browne M., Langville A.N., Pauca V.P., Plemmons R.J. Algorithms and applications for approximate nonnegative matrix factorization // Computational Statistics and Data Analysis, pp. 155-173, 2007.
- [30]. // 1996, c.768-820.
- [31]. Braslavski P. Document Style Recognition Using Shallow Statistical Analysis. In Proceedings of the ESSLLI 2004 Workshop on Combining Shallow and Deep Processing for NLP, Nancy, France, 2004, p. 1–9.
- [32]. DuBay, W.H. The Principles of Readability. Costa Mesa, CA: Impact Information. 2004.
- [33]. P.V. Rao and L.L. Kupper, “Ties in paired-comparison experiments: A generalization of the Bradley–Terry model”, Amer. Statist. Assoc, 62, 1967, pp. 194–204.
- [34]. Turner, H and Firth, D (2012). Bradley-Terry Models in R: The BradleyTerry2 Package. Journal of Statistical Software 48(9), 1–21.
- [35]. Hastie, Tibshirani and Friedman (2008). The Elements of Statistical Learning (2nd edition) Springer-Verlag. 763 pages.
- [36]. Twisted Framework // <http://twistedmatrix.com>.
- [37]. qooxdoo // <http://qooxdoo.org>.
- [38]. « » // <http://www.mmro.ru>.
- [39]. The IEEE International Conference on Data Mining (ICDM) // <http://www.cs.uvm.edu/~icdm>.
- [40]. International Conference on Mechanical and Electrical Technology (ICMET) // <http://www.icmet.ac.cn>.
- [41]. Zhang M.-L., Zhou Z.-H. A k-nearest neighbor based algorithm for multi-label classification // Proceedings of the 1st IEEE International Conference on Granular Computing (GrC'05). Beijing, China, 2005. pp. 718-721.

# Tools for Quality Assessment of Scientific and Technical Documents

S.V. Gerasimov, R.V. Kuryinin, I.V. Mashechkin, M.I. Petrovskiy, Tsarev D.V.,  
A.A. Shestimerov  
Moscow State University, Moscow, Russia  
gerasimov@mlab.cs.msu.su, romaha@mlab.cs.msu.su, mash@cs.msu.su,  
michael@cs.msu.su, tsarev@mlab.cs.msu.su, andy@mlab.cs.msu.su

**Abstract:** In the paper the complex approach to scientific and technical document quality assessment is proposed based on various automatically calculated document quality characteristics as widely used bibliometric and scientometric (based on citation indices), and the new types of characteristics based on the text semantic analysis, heuristics, and also on plagiarism detection methods. The integrated indicator of scientific and technical document quality assessment is formed on the basis of the received basic characteristics with use of machine learning methods similar to the problem of ranking in information retrieval. The developed prototype system based on offered approach is presented, and also the experimental investigations of the developed system directed on check of scientific and technical document quality assessment accuracy are carried out. The analysis of the state of art researches of scientific and technical document quality assessment showed the offered approach based on enhanced list of basic characteristic groups was considered by nobody in so broad statement and as a whole is innovative. The main part of the paper has the following structure. The second section contains an analytical overview of existing approaches to assess quality of scientific and technical documents. The third section provides detail of a proposed approach to assess quality of scientific and technical documents. The fourth section describes a prototype system based on the proposed approach. The fifth section discusses results of experiments.

**Keywords:** scientific and technical document quality assessment; bibliometrics; scientometrics; latent semantic analysis; non-negative matrix factorization; topic model; machine learning

## References

- [1]. Steve Lawrence, Kurt Bollacker, C. Lee Giles. Indexing and Retrieval of Scientific Literature. Eighth International Conference on Information and Knowledge Management, CIKM 99, Kansas City, Missouri, November 2–6, pp. 139–146, 1999.
- [2]. V.V. Pisljakov. Metody ocenki nauchnogo znanija po pokazateljam citirovanija [Methods of assessment of scientific knowledge in terms of citation]. M.: Sociologicheskij zhurnal [Journal of Sociology], 2007, N1, str. 128-140 (in Russian).
- [3]. ISI Web of Knowledge. <http://www.webofknowledge.com>.
- [4]. CiteSeer. <http://citeseerx.ist.psu.edu>.
- [5]. Rossijskij Indeks Nauchnogo Citirovanija [Russian Science Citation Index]. [http://elibrary.ru/project\\_risc.asp](http://elibrary.ru/project_risc.asp) (in Russian).
- [6]. Pisljakov V. V. Naukometricheskie metody i praktiki, rekomenduemye k primeneniju v rabote s rossijskim indeksom nauchnogo citirovanija [Scientometric methods and

practices that are recommended for use in working with the Russian Science Citation Index]. Otchjot o nauchno-issledovatel'skoj rabote (promezhutochnyj) po teme «Razrabotka sistemy statisticheskogo analiza rossijskoj nauki na osnove dannyh rossijskogo indeksa citirovanija». — M., 2005 (in Russian).

- [7]. Meho L (Meho, Lokman); Yang K (Yang, Kiduk). Fusion approach to citation-based quality assessment. Proceedings Of Issi 2007: 11th International Conference Of The International Society For Scientometrics And Informetrics, Vols I And II : 568-581.
- [8]. Angela Vorndran, Alexander Botte. Analysis and evaluation of existing methods and indicators for quality assessment of scientific publications. [http://www.eerqi.eu/sites/default/files/Analysis\\_and\\_evaluation\\_of\\_existing\\_methods\\_and\\_indicators.pdf](http://www.eerqi.eu/sites/default/files/Analysis_and_evaluation_of_existing_methods_and_indicators.pdf) [PDF].
- [9]. EERQI – European Educational Research Quality Indicators. [www.eerqi.eu](http://www.eerqi.eu).
- [10]. EERQI Project Final Report (2011). [http://eerqi.eu/sites/default/files/Final\\_Report.pdf](http://eerqi.eu/sites/default/files/Final_Report.pdf) [PDF].
- [11]. Moyses Szklo. Quality of scientific articles. Revista Saúde Pública vol.40 special issue São Paulo Aug. 2006.
- [12]. Dr Navneet Gupta BSc (Hons) PhD MCOptom FBCLA. How to Evaluate a Scientific Research Article. <http://www.optometry.co.uk/uploads/articles/ARTICLE%200309.pdf> [PDF].
- [13]. Google Scholar. <http://scholar.google.ru>.
- [14]. Berry M.W., Dumais S.T., O'Brien G.W. Using Linear Algebra for Intelligent Information Retrieval. University of Tennessee Knoxville. TN. USA, 1994.
- [15]. Lee D.D., Seung H.S. Learning the parts of objects by non-negative matrix factorization. Nature, 401, pp. 788-791, 1999.
- [16]. Rakesh P., Shivapratap G., Divya G., Soman KP. Evaluation of SVD and NMF Methods for Latent Semantic Analysis. International Journal of Recent Trends in Engineering, Vol. 1, No. 3, 2009.
- [17]. Griffiths T L, Steyvers M. Finding scientific topics. In: Proceedings of the National Academy of Sciences. USA, 2004, 101: 5228–5235.
- [18]. Steinberger J., Ježek K. Text Summarization and Singular Value Decomposition. In Lecture Notes for Computer Science vol. 2457, Springer-Verlag, 2004, pp. 245-254.
- [19]. Steinberger J. Text Summarization within the LSA Framework. PhD Thesis, University of West Bohemia in Pilsen, Czech Republic, January 2007.
- [20]. Mashechkin I.V., Petrovskij M.I., Carjov D.V. Metody vychislenija relevantnosti fragmentov teksta na osnove tematiceskikh modelej v zadache avtomaticheskogo annotirovanija [Methods for calculating the relevance of text fragments on the basis of thematic patterns in the problem of automatic annotation]. Vychislitel'nye metody i programirovanie. Tom 14, 2013. 91-102 [in Russian].
- [21]. Mashechkin I.V., Petrovskiy M.I., Popov D.S., Tsarev D.V. Automatic text summarization using latent semantic analysis. Programming and Computer Software, pp. 299-305, 2011.
- [22]. Tsarev D., Petrovskiy M., Mashechkin I. Using NMF-based text summarization to improve supervised and unsupervised classification. 11th International Conference on Hybrid Intelligent Systems (HIS), Malacca, MALAYSIA. P. 185-189, 2011.
- [23]. Dmitry Tsarev, Mikhail Petrovskiy and Igor Mashechkin, Supervised and Unsupervised Text Classification via Generic Summarization International Journal of Computer Information Systems and Industrial Management Applications. MIR Labs, Volume 5, 2013, pp. 509-515.

- [24]. Wei Xu, Xin Liu, Yihong Gong Document clustering based on non-negative matrix factorization. Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval, Toronto, Canada, 2003.
- [25]. Y. Ding. Applying weighted PageRank to author citation networks. In Proceedings of JASIST. 2011, pp. 236-245.
- [26]. M. Potthast, T. Gollub, M. Hagen, J. Graßegger, J. Kiesel, M. Michel, A. Oberländer, M. Tippmann, A. Barrón-Cedeño, P. Gupta, P. Rosso, B. Stein. Overview of the 4th International Competition on Plagiarism Detection. CLEF2012. 2012.
- [27]. S. Alzahrani, N. Salim. Fuzzy Semantic-Based String Similarity for Extrinsic Plagiarism Detection, Lab Report for PAN at CLEF2010, 2010.
- [28]. A. Martins. String kernels and similarity measures for information retrieval. 2006.
- [29]. Berry M.W., Browne M., Langville A.N., Pauca V.P., Plemmons R.J. Algorithms and applications for approximate nonnegative matrix factorization. Computational Statistics and Data Analysis, pp. 155-173, 2007.
- [30]. Fomenko V.P., Fomenko T.G. Avtorskij invariant russkih literaturnyh tekstov [Author invariant Russian literary texts]. Predislovie A.T. Fomenko. Fomenko A.T. Novaja hronologija Grecii: Antichnost' v srednevekov'e. T. 2. M.: Izd-vo MGU, 1996, c.768-820 (in Russian).
- [31]. Braslavski P. Document Style Recognition Using Shallow Statistical Analysis. In Proceedings of the ESSLLI 2004 Workshop on Combining Shallow and Deep Processing for NLP, Nancy, France, 2004, p. 1–9.
- [32]. DuBay, W.H. The Principles of Readability. Costa Mesa, CA: Impact Information. 2004.
- [33]. P.V. Rao and L.L. Kupper, “Ties in paired-comparison experiments: A generalization of the Bradley–Terry model”, Amer. Statist. Assoc, 62, 1967, pp. 194–204.
- [34]. Turner, H and Firth, D (2012). Bradley-Terry Models in R: The BradleyTerry2 Package. Journal of Statistical Software 48(9), 1–21.
- [35]. Hastie, Tibshirani and Friedman (2008). The Elements of Statistical Learning (2nd edition) Springer-Verlag. 763 pages.
- [36]. Twisted Framework. <http://twistedmatrix.com>.
- [37]. qooxdoo. <http://qooxdoo.org>.
- [38]. Konferencija «Matematicheskie metody raspoznavanija obrazov» [The Conference «Mathematical Methods of Pattern Recognition»]. <http://www.mmro.ru> (In Russian).
- [39]. The IEEE International Conference on Data Mining (ICDM). <http://www.cs.uvm.edu/~icdm>.
- [40]. International Conference on Mechanical and Electrical Technology (ICMET). <http://www.icmet.ac.cn>.
- [41]. Zhang M.-L., Zhou Z.-H. A k-nearest neighbor based algorithm for multi-label classification. Proceedings of the 1st IEEE International Conference on Granular Computing (GrC'05). Beijing, China, 2005. pp. 718-721.