

Analysis of bibliographic networks

Vladimir Batagelj

IMFM Ljubljana and IAM UP Koper

vladimir.batagelj@fmf.uni-lj.si

Daria Maltseva

NRU HSE Moscow

d_malceva@mail.ru

NetGloW workshop

St Petersburg, July 4, 2018

Workshop agenda

Part 1: Transforming bibliographic data into networks

Part 2: Analyzing bibliographic networks

Transforming bibliographic data into networks

1. Goals, research questions, and theory
2. Bibliographic data:
 - Structure of bibliographic data
 - Bibliographic databases and data collection
 - Networks from bibliographic data
3. Problems associated with bibliographic data collection
4. Tools for collection and maintenance of bibliographic data
5. Conversion to networks

1. Goals, research questions, and theory

Goals: study of social and cognitive structure of different scientific fields.

Research questions:

- How do scientists collaborate with each other? How different groups of scientists relate to each other? → *Co-authorship, co-citation network analysis.*
- How the certain fields in science develop trough time? → *Citation network analysis.*
- What is the topic structure of the scientific filed? → *Co-occurrence key words analysis.*

Different levels (authors, institutions, countries) and units (publications, journals) of analysis.

1. Goals, research questions, and theory

Theoretical background:

- The “philosophical” grounds of the field go back to the works of the sociologist R. Merton and the historians of science D. de Sola Price and G. Small.
- E. Garfield - the first scientific citation index - Science Citation Index (SCI) [Garfield, 1972]. Since its creation, the citation analysis has grown into an independent research field [Wilson, 1999, Bar-Ilan, 2008].
- D. Crane (Crane 1972) introduced the notion of “invisible college” - a core group of scientists who collaborated with each other and generated a disproportionate volume of new ideas - and showed that internal social structure of the scientific community influences the development of the ideas, and **study of informal social and communication structures can bring important results for understanding the modern development of scientific disciplines.**

2.1. Bibliographic data: some examples

- Survey papers

Granovetter, M. (1983). The strength of weak ties: A network theory revisited. *Sociological theory*, 201-233.

Albert, R., & Barabási, A. L. (2002). Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1), 47.

Batagelj V., Ferligoj A., Squazzoni F. The emergence of a field: a network analysis of research on peer review // *Scientometrics*. – 2017. – T. 113. – №. 1. – C. 503-532.

- Books

White, Harrison C. *Identity and Control: How Social Formations Emerge* (Second Edition). PRINCETON; OXFORD: Princeton University Press, 2008.

Burt R. S. *Structural holes: The social structure of competition*. – Harvard university press, 2009.

Batagelj, Vladimir, Andrej Mrvar, and Wouter de Nooy. "Exploratory social network analysis with Pajek." (2008).

2.1. Bibliographic data: some examples

- Survey bibliographies
 - [Web survey bibliography](#)
 - [Bibliography of Research Methods Texts](#)
 - [A survey and annotated bibliography of multiobjective combinatorial optimization](#)
 - [Community detection in graphs](#)
- Book bibliography
 - [Handbook of Product Graphs, Second Edition](#)
 - [Computational Geometry](#)
- Bibliography of scientific community
 - [Bibliography on Self-Organizing Map \(SOM\) method](#)
 - [Computational Geometry Bibliographies](#)

2.1. Bibliographic data: some examples

Imrich W, Klavžar S.
(1999) Graph products.

References

- [Abay-Asmerom, 1998] Abay-Asmerom, G. (1998). Imbeddings of the tensor product of graphs where the second factor is a complete graph. *Discrete Math.*, 182:13–19.
- [Aho et al., 1974] Aho, A. V., Hopcroft, J. E., and Ullman, J. D. (1974). *The Design and Analysis of Computer Algorithms*. Addison-Wesley, Reading, MA.
- [Aho et al., 1987] Aho, A. V., Hopcroft, J. E., and Ullman, J. D. (1987). *Data Structures and Algorithms*. Addison-Wesley, Reading, MA.
- [Albertson and Collins, 1985] Albertson, M. O. and Collins, K. L. (1985). Homomorphisms of 3-chromatic graphs. *Discrete Math.*, 54:127–132.
- [Alexe and Olaru, 1997] Alexe, G. and Olaru, E. (1997). The strongly perfectness of normal product of t -perfect graphs. *Graphs and Combinatorics*, 13:209–215.
- [Alles, 1985] Alles, P. (1985). The dimension of sums of graphs. *Discrete Math.*, 54:229–233.
- [Alon, 1986] Alon, N. (1986). Covering graphs by the minimum number of equivalence relations. *Combinatorica*, 6:201–206.
- [Alon et al., 1997] Alon, N., Yuster, R., and Zwick, U. (1997). Finding and counting given length cycles. *Algorithmica*, 17:209–223.
- [Alspach et al., 1990] Alspach, B., Bermond, J.-C., and Sotteau, D. (1990). Decompositions into cycles I: Hamilton decompositions. In Hahn, G., Sabidussi, G., and Woodrow, R. E., editors, *Cycles and Rays: Basic Structures in Finite and Infinite Graphs*, volume 301 of *NATO ASI Ser., Ser. C*, pages 9–18. Kluwer, Dordrecht.
- [Alspach and George, 1990] Alspach, B. and George, J. C. (1990). One-factorizations of tensor products of graphs. In Bodendiek, R. and Henn, R., editors, *Topics in Combinatorics and Graph Theory: Essays in Honour of Gerhard Ringel*, pages 41–46. Physica-Verlag, Heidelberg.

2.2. Bibliographic databases

Records from [BiBTeX](#) (reference management software for formatting lists of references, typically used together with the LaTeX document preparation system).

```
@Article{int:Mizuno1,
  author =      "S. Mizuno",
  title =       "An  $O(n^3L)$  algorithm using a sequence for
                linear complementarity problems",
  journal =     "Journal of the Operations Research Society of Japan",
  volume =     "33",
  year =       "1990",
  pages =      "66--75",
}

@InCollection{int:Vorst1,
  author =      "{J. G. G. van de} Vorst",
  title =       "An attempt to use parallel computing in large scale
                optimisation",
  booktitle =   "Logistics, Where Ends Have to Meet~: Proceedings of
                the Shell Conference on Logistics in Apeldoorn, The
                Netherlands, November 1988",
  editor =      "{C. F. H. van} Rijn",
  year =       "1989",
  pages =      "112--119",
  publisher =   "Pergamon Press",
  address =     "Oxford, United Kingdom",
}
```

BIBTEX

BiBTeX -> Pajek converter [Bib2Pajek.py](#)

2.2. Bibliographic databases

Records from [DBLP](#) (online database database of a computer science bibliography).

```
<article mdate="2004-01-15" key="journals/arscom/BeinekeGL97">  
<author>Lowell W. Beineke</author>  
<author>Wayne Goddard</author>  
<author>Marc J. Lipman</author>  
<title>Graphs with Maximum Edge-Integrity.</title>  
<year>1997</year>  
<volume>46</volume>  
<journal>Ars Comb.</journal>  
<url>db/journals/arscom/arscom46.html#BeinekeGL97</url>  
</article>
```

```
<inproceedings mdate="2004-12-09" key="conf/sigcse/BermanD96">  
<author>A. Michael Berman</author>  
<author>Robert C. Duvall</author>  
<title>Thinking about binary trees in an object-oriented world.</title>  
<pages>185-189</pages>  
<year>1996</year>  
<crossref>conf/sigcse/1996</crossref>  
<booktitle>SIGCSE</booktitle>  
<ee>http://doi.acm.org/10.1145/236536</ee>  
<url>db/conf/sigcse/sigcse1996.html#BermanD96</url>  
</inproceedings>
```



DBLP XML data to Pajek Convertor -> Pajek converter [DBLP2Pajek.py](#)

2.2. Bibliographic databases

Records from [Zentralblatt Math](#) (international reviewing service providing reviews and abstracts for articles in pure and applied mathematics).

```
an 00549739
ai gross.mark-d
is ISSN 0025-5874; ISSN 1432-1823
au Gross, Mark
py 1993
cc *14M15 14J15
ti Surfaces of bidegree  $(3,n)$  in  $\text{Gr}(1,\mathbb{P}^3)$ .
ut congruence; family of lines
so Math. Z. 212, No.1, 73-106 (1993).
an 01488230
ai tiras.yuecel; harmanci.abdullah; -
is ISSN 0092-7872; ISSN 1532-4125
au Tiras, Yücel; Harmancı, Abdullah; Smith, P.F.
py 2000
cc *13A15 13C05
ti Some remarks on dense submodules of multiplication modules.
ut multiplication module; dense submodule
so Commun. Algebra 28, No.5, 2291-2296 (2000).
se 00000057 Communications in Algebra Commun. Algebra 0092-7872; 1532-4125
```



ZBml files -> Pajek converter [ZBml.py](#)

2.2. Bibliographic databases

Records from [Web of Science](#) (online subscription-based scientific citation indexing service providing a comprehensive citation search).

```
PT J
AU Elmer, T
   Boda, Z
   Stadtfeld, C
TI The co-evolution of emotional well-being with weak and strong friendship
   ties
SO NETWORK SCIENCE
LA English
DT Article
DE social networks; ordered stochastic actor-oriented models [...]
ID ADOLESCENT DEPRESSIVE SYMPTOMS;[...]
AB Social ties are strongly related to well-being. But what characterizes this relationship
C1 [Elmer, Timon; Boda, Zsafia; Stadtfeld, Christoph] Swiss Fed Inst Technol,
   Chair Social Networks, Dept Humanities Social & Polit Sci, Zurich, Switzerland.
RP Elmer, T (reprint author), Swiss Fed Inst Technol, Chair Social Networks,
   Dept Humanities Social & Polit Sci, Zurich, Switzerland.
EM timon.elmer@gess.ethz.ch; [...]
CR Aharony N, 2011, PERVASIVE MOB COMPUT, V7, P643, DOI 10.1016/j.pmcj.2011.09.004
   Baerveldt C., 2004, CONNECTIONS, V26, P11
   Reis H. T., 2000, HDB RES METHODS SOCI
   Ripley Ruth M., 2015, MANUAL RSIENA
...
```



WEB OF SCIENCE™

2.2. Bibliographic databases

Records from [Web of Science](#) (online subscription-based scientific citation indexing service providing a comprehensive citation search).

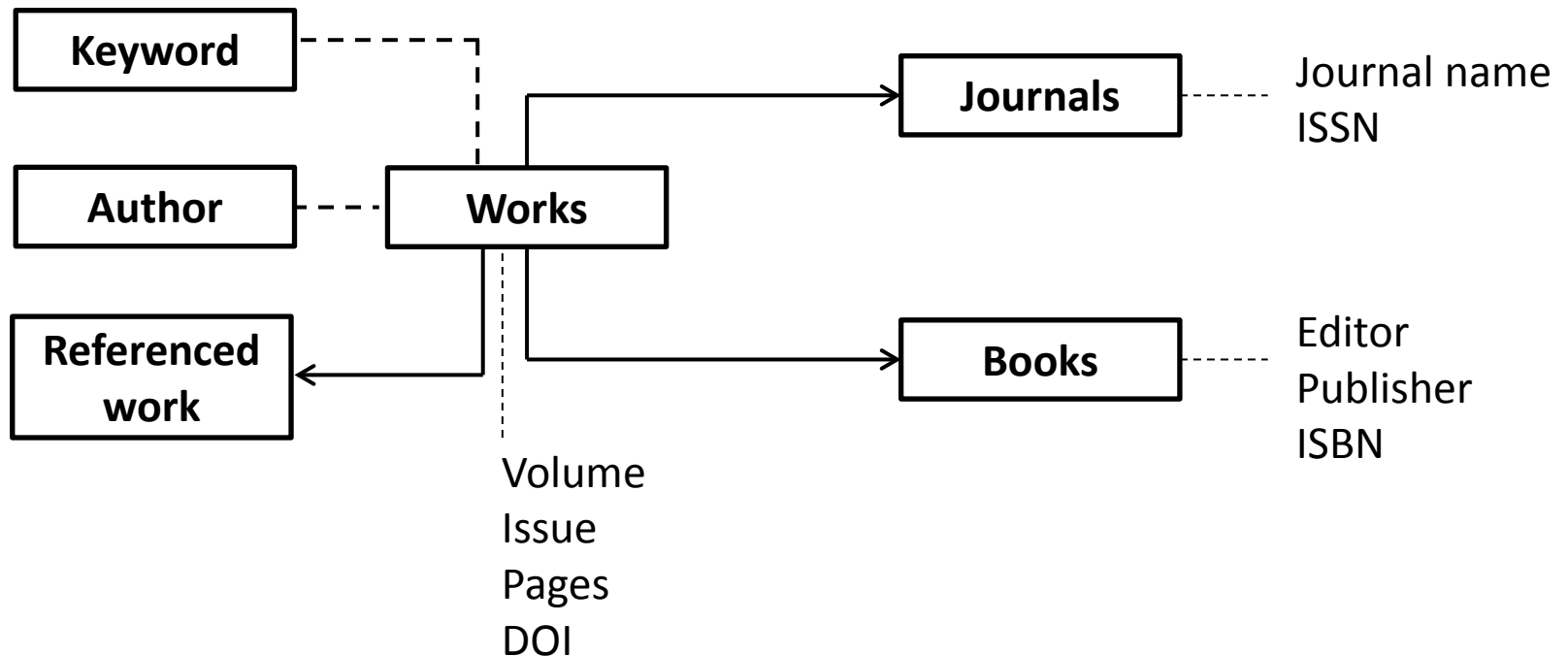
```
...
NR 83
TC 1
PU CAMBRIDGE UNIV PRESS
PI NEW YORK
PA 32 AVENUE OF THE AMERICAS, NEW YORK, NY 10013-2473 USA
SN 2050-1242
J9 NETW SCI
JI Netw. Sci.
PD SEP
PY 2017
VL 5
IS 3
BP 278
EP 307
DI 10.1017/nws.2017.20
PG 30
SC Social Sciences - Other Topics
GA FF0AM
UT WOS:000408564600003
ER
```



WEB OF SCIENCE™

Web of Science files -> Pajek converter [WoSPajek](#)

2.3. Structure of bibliographic data



2.4. Networks from bibliographic data

We can derive some two-mode networks on selected topics from bibliographies from:

- books and survey papers,
- special bibliographies ([BibTeX](#))
- bibliographic services
 - Web of Science
 - Scopus
 - SICRIS
 - CiteSeer
 - Zentralblatt MATH
 - Google Scholar
 - DBLP Bibliography
 - US patent office
 - IMDb

2.4. Networks from bibliographic data

Two-mode networks on selected topics:

- works \times authors (**WA**),
- works \times journals or book publishers (**WJ**);
- works \times keywords (**WK**);
- works \times classification (**WC**) - from some data;
- the one-mode citation network works \times works (**Ci**), where works include papers, reports, books, patents etc.;
- authors \times institutions (**AI**);
- authors \times countries (**AC**).

Besides this we get also at least the partition of works by the journal or publisher, the partition of works by the publication year, and the vector of number of pages.

2.4. Networks from bibliographic data

Creating your own bibliographic data base in Excel

How to describe a network $\mathbf{N} = (\mathbf{V}, \mathbf{L}, \mathbf{P}, \mathbf{W})$? In principle the answer is simple - we list its components: nodes \mathbf{V} , links \mathbf{L} , node properties \mathbf{P} , and link weights \mathbf{W} .

The simplest way is to describe a network \mathbf{N} by providing (\mathbf{V}, \mathbf{P}) and (\mathbf{L}, \mathbf{W}) in a form of two tables.

2.4. Networks from bibliographic data

Creating your own bibliographic data base in Excel

As an example, let us describe a part of network determined by the following works:

[Generalized blockmodeling](#), [Clustering with relational constraint](#), [Partitioning signed social networks](#), [The Strength of Weak Ties](#).

There are **nodes of different types** (modes): *persons, papers, books, series, journals, publishers*;

and **different relations** among them: *author_of, editor_of, contained_in, cites, published_by*.

For small bibliographies both tables can be maintained in Excel and exported as text in [CSV](#) (Comma Separated Values) format.

2.4. Networks from bibliographic data

Creating your own bibliographic data base in Excel

[bibNodes.csv](#)

```
name;mode;country;sex;year;vol;num;fPage;lPage;x;y
"Batagelj, Vladimir";person;SI;m;;;;;;809.1;653.7
"Doreian, Patrick";person;US;m;;;;;;358.5;679.1
"Ferligoj, Anuška";person;SI;f;;;;;;619.5;680.7
"Granovetter, Mark";person;US;m;;;;;;145.6;660.5
"Moustaki, Irini";person;UK;f;;;;;;783.0;228.0
"Mrvar, Andrej";person;SI;m;;;;;;478.0;630.1
"Clustering with relational constraint";paper;;;1982;47;;413;426;684.1;380.1
"The Strength of Weak Ties";paper;;;1973;78;6;1360;1380;111.3;329.4
"Partitioning signed social networks";paper;;;2009;31;1;1;11;408.0;337.8
"Generalized Blockmodeling";book;;;2005;24;;1;385;533.0;445.9
"Psychometrika";journal;;;;;;741.8;086.1
"Social Networks";journal;;;;;;321.4;236.5
"The American Journal of Sociology";journal;;;;;;111.3;168.9
"Structural Analysis in the Social Sciences";series;;;;;;310.4;082.8
"Cambridge University Press";publisher;UK;;;;;;534.3;238.2
"Springer";publisher;US;;;;;;884.6;174.0
```

In large networks, to avoid the empty cells, we split a network to some subnetworks
- a collection.

2.4. Networks from bibliographic data

Creating your own bibliographic data base in Excel

[bibLinks.csv](#)

```
from;relation;to
"Batagelj, Vladimir";authorOf;"Generalized Blockmodeling"
"Doreian, Patrick";authorOf;"Generalized Blockmodeling"
"Ferligoj, Anuška";authorOf;"Generalized Blockmodeling"
"Batagelj, Vladimir";authorOf;"Clustering with relational constraint"
"Ferligoj, Anuška";authorOf;"Clustering with relational constraint"
"Granovetter, Mark";authorOf;"The Strength of Weak Ties"
"Granovetter, Mark";editorOf;"Structural Analysis in the Social Sciences"
"Doreian, Patrick";authorOf;"Partitioning signed social networks"
"Mrvar, Andrej";authorOf;"Partitioning signed social networks"
"Moustaki, Irini";editorOf;"Psychometrika"
"Doreian, Patrick";editorOf;"Social Networks"
"Generalized Blockmodeling";containedIn;"Structural Analysis in the Social Sciences"
"Clustering with relational constraint";containedIn;"Psychometrika"
"The Strength of Weak Ties";containedIn;"The American Journal of Sociology"
"Partitioning signed social networks";containedIn;"Social Networks"
"Partitioning signed social networks";cites;"Generalized Blockmodeling"
"Generalized Blockmodeling";cites;"Clustering with relational constraint"
"Structural Analysis in the Social Sciences";publishedBy;"Cambridge University Press"
"Psychometrika";publishedBy;"Springer"
```

2.4. Networks from bibliographic data

Factorization and description of large networks

To save space and improve the computing efficiency we often replace values of categorical variables with integers. In R this encoding is called a **factorization**.

We enumerate all possible values of a given categorical variable (coding table) and afterwards replace each its value by the corresponding index in the coding table.

This approach is used in most programs dealing with large networks. Unfortunately the coding table is often a kind of meta-data.

2.4. Networks from bibliographic data

Factorization and description of large networks

[CSV2Pajek.R](#)

```
# transforming CSV file to Pajek files
# by Vladimir Batagelj, June 2016
setwd("C:/Users/batagelj/work/Python/graph/SVG/EUSN")
colC <- c(rep("character",4),rep("integer",7)); nas <- c("", "NA", "NaN")
nodes <- read.csv2("bibNodes.csv",encoding='UTF-8',colClasses=colC,na.strings=nas)
n <- nrow(nodes); M <- factor(nodes$mode); S <- factor(nodes$sex)
mod <- levels(M); sx <- levels(S); S <- as.numeric(S); S[is.na(S)] <- 0
links <- read.csv2("bibLinks.csv",encoding='UTF-8',colClasses="character")
F <- factor(links$from,levels=nodes$name,ordered=TRUE)
T <- factor(links$to,levels=nodes$name,ordered=TRUE)
R <- factor(links$relation); rel <- levels(R)
net <- file("bib.net","w"); cat('*vertices ',n,'\n',file=net)
clu <- file("bibMode.clu","w"); sex <- file("bibSex.clu","w")
cat('%',file=clu); cat('%',file=sex)
for(i in 1:length(mod)) cat(' ',i,mod[i],file=clu)
cat('\n*vertices ',n,'\n',file=clu)
for(i in 1:length(sx)) cat(' ',i,sx[i],file=sex)
cat('\n*vertices ',n,'\n',file=sex)
for(v in 1:n) {
  cat(v,' ',nodes$name[v],""'\n',sep='',file=net);
  cat(M[v],'\n',file=clu); cat(S[v],'\n',file=sex)
}
for(r in 1:length(rel)) cat('*arcs : ',r,' ',rel[r],""'\n',sep='',file=net)
cat('*arcs\n',file=net)
for(a in 1:nrow(links))
  cat(R[a],': ',F[a],', ',T[a],', 1 1 ',rel[R[a]],""'\n',sep='',file=net)
close(net); close(clu); close(sex)
```

2.4. Networks from bibliographic data

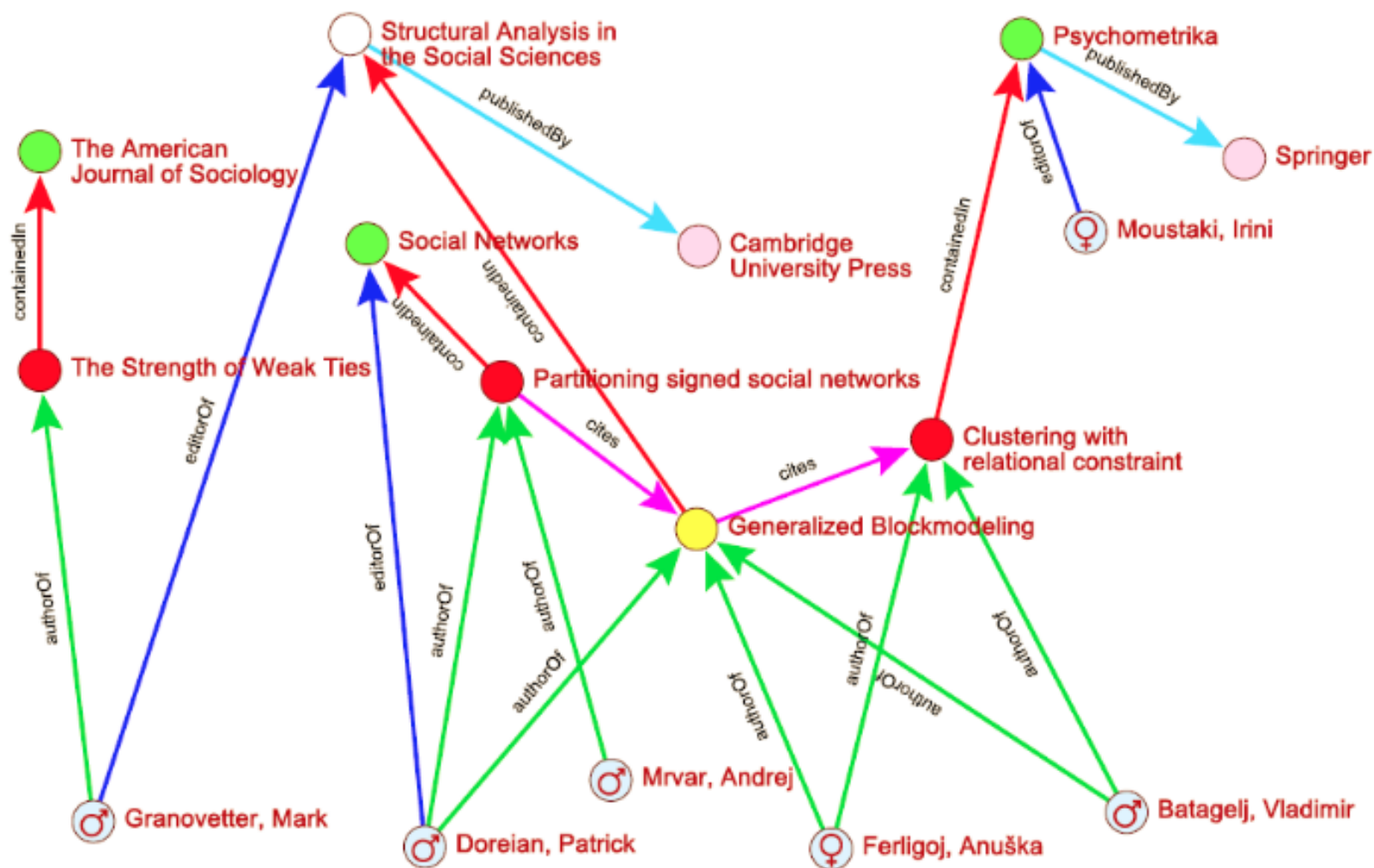
[bib.net](#)

```
*vertices 16
1 "Batagelj, Vladimir"
2 "Doreian, Patrick"
3 "Ferligoj, Anuška"
4 "Granovetter, Mark"
5 "Moustaki, Irini"
6 "Mrvar, Andrej"
7 "Clustering with relational constraint"
8 "The Strength of Weak Ties"
9 "Partitioning signed social networks"
10 "Generalized Blockmodeling"
11 "Psychometrika"
12 "Social Networks"
13 "The American Journal of Sociology"
14 "Structural Analysis in the Social Sciences"
15 "Cambridge University Press"
16 "Springer"
*arcs :1 "authorOf"
*arcs :2 "cites"
*arcs :3 "containedIn"
*arcs :4 "editorOf"
*arcs :5 "publishedBy"

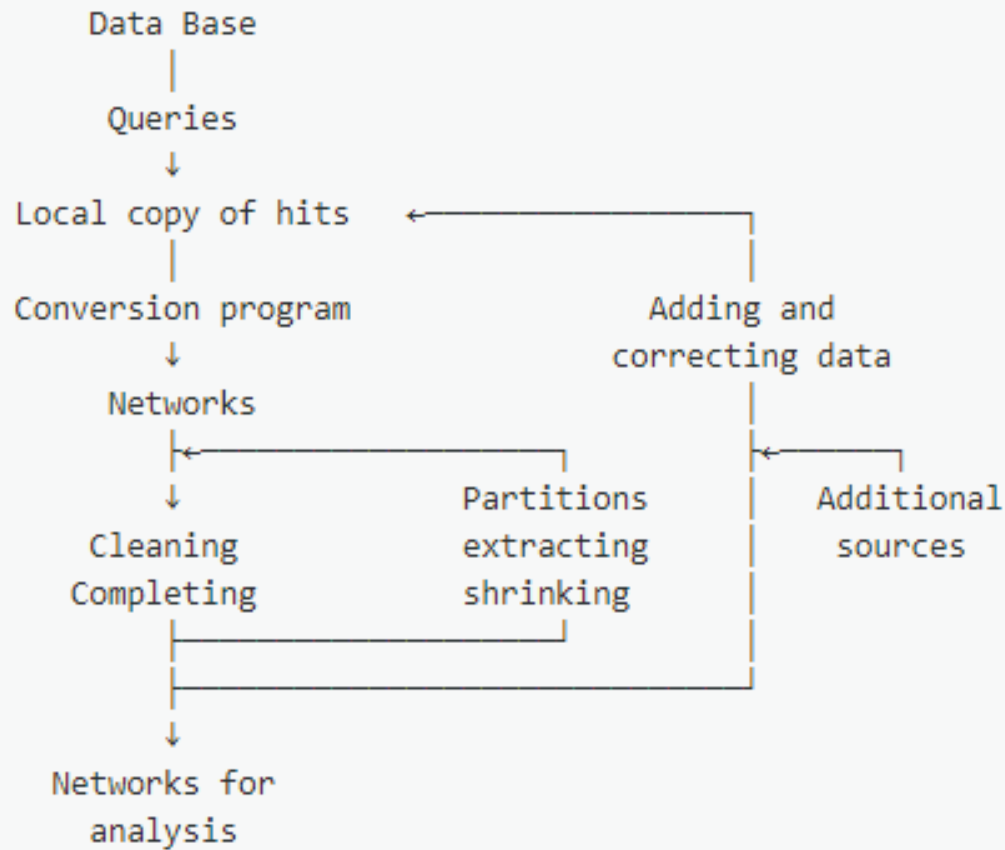
*arcs
1: 1 10 1 1 "authorOf"
1: 2 10 1 1 "authorOf"
1: 3 10 1 1 "authorOf"
1: 1 7 1 1 "authorOf"
1: 3 7 1 1 "authorOf"
1: 4 8 1 1 "authorOf"
4: 4 14 1 1 "editorOf"
1: 2 9 1 1 "authorOf"
1: 6 9 1 1 "authorOf"
4: 5 11 1 1 "editorOf"
4: 2 12 1 1 "editorOf"
3: 10 14 1 1 "containedIn"
3: 7 11 1 1 "containedIn"
3: 8 13 1 1 "containedIn"
3: 9 12 1 1 "containedIn"
2: 9 10 1 1 "cites"
2: 10 7 1 1 "cites"
5: 14 15 1 1 "publishedBy"
5: 11 16 1 1 "publishedBy"
```

[bibMode.clu](#), [bibSex.clu](#); [bib.paj](#), [bib.ini](#)

2.4. Networks from bibliographic data

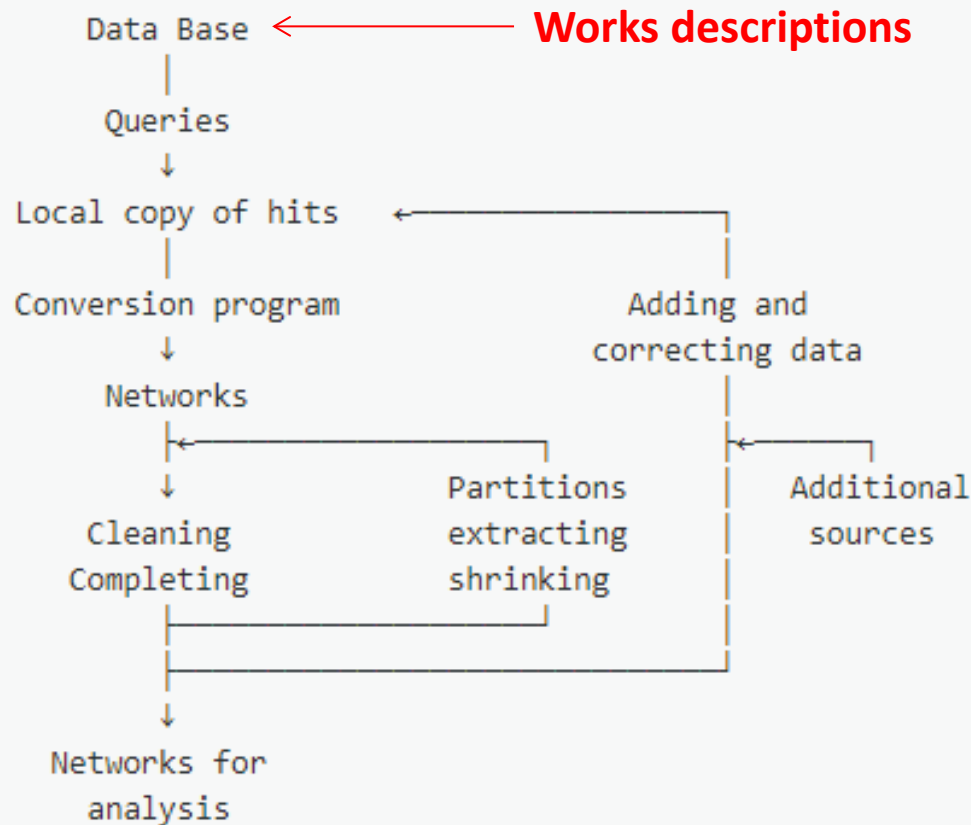


2.4. Networks from bibliographic data: the general procedure of transformation



3. Problems with bibliographic data.

Problem 1: Descriptions



3. Problems with bibliographic data.

Problem 1: Descriptions

Most of the source bibliographic data are semi-structured – they are available in the form of records from some data base.

Selected fields in the record represent different units: names of people, names of journals, keywords, IDs of works, countries, institutions, etc. Unfortunately the names of these units are usually not stored in a standardized way.

- **Detail of description** (list of attributes)
- **Completeness of description** (all entities are included - authors)

3. Problems with bibliographic data.

Problem 1: Descriptions

- **Citation formats:**

Different academic styles, guided by different associations: [MLA](#) (Modern Language Association of America), [APA](#) (American Psychological Association), [Chicago](#) (University of Chicago Press), [AMA](#) (American Medical Association), [SCE](#) (Council of Science Editors), [GOST](#) (Russian state standard), etc.

APA

White, H. (2008). Identity and Control: How Social Formations Emerge (Second Edition). PRINCETON; OXFORD: Princeton University Press.

MLA

White, Harrison C. Identity and Control: How Social Formations Emerge (Second Edition). Princeton University Press, 2008.

Chicago

White, Harrison C. Identity and Control: How Social Formations Emerge (Second Edition). PRINCETON; OXFORD: Princeton University Press, 2008.

GOST

White H. C. Identity and control. - Princeton University Press, 2008.

- **Bibliographic formats:** [BibTex](#) (LaTeX style), [EndNote](#) (Clarivate Analytics), [RIS](#) (Research Information Systems style), etc.

```
@book{white2012identity,  
  title={Identity and control},  
  author={White, Harrison C},  
  year={2012},  
  publisher={Princeton University Press}  
}
```

```
%0 Book  
%T Identity and control  
%A White, Harrison C  
%D 2012  
%I Princeton University Press
```

```
TY - BOOK  
TI - Identity and Control  
AU - White, Harrison C.  
AB - <p>In this completely revised edition ...</p>  
PB - Princeton University Press  
PY - 2008  
SN - 9780691137155  
T1 - How Social Formations Emerge (Second Edition)  
UR - http://www.jstor.org/stable/j.ctt1r2fg1  
ER -
```

3. Problems with bibliographic data.

Problem 1: Descriptions

A typical description in bibliographies from books and (survey) papers contains **the following elements**:

- Names of authors; sometimes not complete (et al.)
- Title
- Publication year (date)

WASSERMAN S, 1994, SOCILA NETWORK ANAL

Wasserman, S., & Faust, K. (1994). Social network analysis: Methods and applications (Vol. 8). Cambridge university press.

For papers:

- Journal
- Volume
- Issue
- Pages

Granovetter, M. (1983). The strength of weak ties: A network theory revisited. Sociological theory, 201-233.

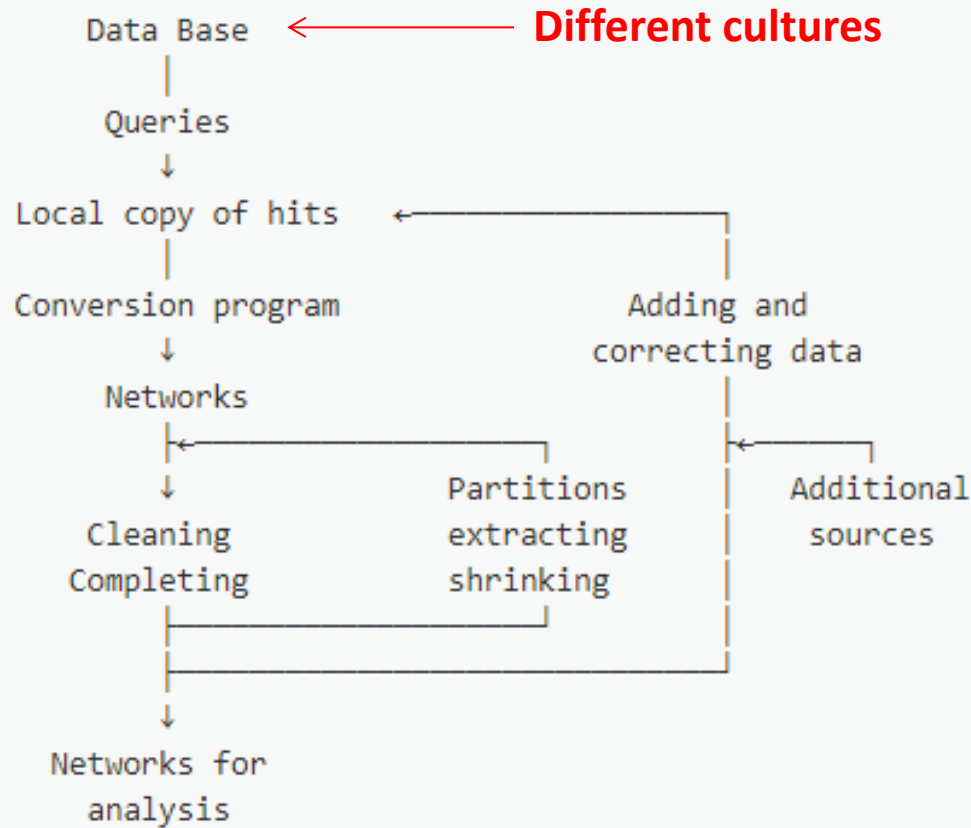
For books:

- Publisher (Company, Place)

White, H. (2008). Identity and Control: How Social Formations Emerge (Second Edition). PRINCETON; OXFORD: Princeton University Press.

3. Problems with bibliographic data.


Problem 2: Different cultures



3. Problems with bibliographic data.

Problem 2: Different cultures

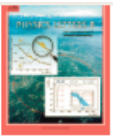
- Different number of co-authors;
- Russia – PhD-candidates supposed to publish as the only authors.



Physics Letters B

Volume 716, Issue 1, 17 September 2012, Pages 1-29

open access



Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC ☆

This paper is dedicated to the memory of our ATLAS colleagues who did not live to see the full impact and significance of their contributions to the experiment.

This paper is dedicated to the memory of our ATLAS colleagues who did not live to see the full impact and significance of their contributions to the experiment.

ATLAS Collaboration *

G. Aad⁴⁸, T. Abajyan²¹, B. Abbott¹¹¹, J. Abdallah¹², S. Abdel Khalek¹¹⁵, A.A. Abdelalim⁴⁹, O. Abdinov¹¹, R. Aben¹⁰⁵, B. Abi¹¹², M. Abolins⁸⁸, O.S. AbouZeid¹⁵⁸, H. Abramowicz¹⁵³, H. Abreu¹³⁸, B.S. Acharya^{164a, 164b}, L. Adamczyk³⁸, D.L. Adams²⁵, T.N. Addy⁵⁶, J. Adelman¹⁷⁶, S. Adomeit⁹⁸, P. Adragna⁷⁵, T. Adye¹²⁹, S. Aefsky²³, J.A. Aguilar-Saavedra^{124b, a}, M. Agustoni¹⁷, M. Aharrouche⁸¹, S.P. Ahlen²², F. Ahles⁴⁸, A. Ahmad¹⁴⁸, M. Ahsan⁴¹, G. Aielli^{133a, 133b}, T. Akdogan^{19a}, T.P.A. Åkesson⁷⁹, G. Akimoto¹⁵⁵, A.V. Akimov⁹⁴, M.S. Alam², M.A. Alam⁷⁶, J. Albert¹⁶⁹, S. Albrand⁵⁵, M. Aleksa³⁰, I.N. Aleksandrov⁶⁴, F. Alessandria^{89a}, C. Alexa^{26a}, G. Alexander¹⁵³, G. Alexandre⁴⁹, T. Alexopoulos¹²⁰, B.M.M. Allbrooke¹⁸, P.P. Allport⁷³, S.E. Allwood¹²⁰, A. Alonso⁷⁹, F. Alonso⁷⁰, A. Altheimer³⁵, B. Alvarez Gonzalez⁸⁸, M.G. Alviggi^{102a, 102b}, K. Amako⁶⁵, C. Amelung²³, V.V. Ammosov^{128, *}, S.P. Amor Dos Santos^{124a}, A. Amorim^{124a, b}, N. Amram¹⁵³, C. Anastopoulos³⁰, L.S. Ancu¹⁷, N. Andari¹¹⁵, T. Andeen³⁵, C.F. Anders^{58b}, G. Anders^{58a}, K.J. Anderson³¹, A. Andreazza^{89a, 89b}, V. Andrei^{58a}, M.-L. Andrieux⁵⁵, X.S. Anduaga⁷⁰, S. Angelidakis⁹, P. Anger⁴⁴, A. Angerami³⁵, F. Anghinolfi³⁰, A. Anisenkov¹⁰⁷, N. Anjos^{124a}, A. Annovi⁴⁷, A. Antonaki⁹, M. Antonelli⁴⁷, A. Antonov⁹⁶, J. Antos^{144b}, F. Anulli^{132a}, M. Aoki¹⁰¹, S. Aoun⁸³, L. Aperio Bella⁵, R. Apolle^{118, c}, G. Arabidze⁸⁸, I. Aracena¹⁴³, Y. Arai⁶⁵, A.T.H. Arce⁴⁵, S. Arfaoui¹⁴⁸, J.-F. Arguin⁹³, E. Arik^{19a, *}, M. Arik^{19a}, A.J. Armbruster⁸⁷, O. Arnaez⁸¹, V. Arnal⁸⁰, C. Arnault¹¹⁵, A. Artamonov⁹⁵, G. Artoni

More than 3 000 co-authors

3. Problems with bibliographic data.

Problem 2: Different cultures

- Writing of names (initial/full name, first name first/last)
- The order of first and last name (French, Spanish, Arabian names etc., names with prefixes).
- Some journals have special rules about abbreviations of journal names

```
Bon G., 1896, CROWD STUDY POPULAR
Le Bon G, 1897, CROWD STUDY POPULAR
LeBon G., 1960, CROWD STUDY POPULAR
Lebon G., 2011, PSIHOLOGIJA NARODOV
Le Bon Gustave, 1930, CROWD STUDY POPULAR
Gustave Le Bon, 1982, PSYCHOL MASSEN
```

```
GRANOVET.MS, 1973, AM J SOCIOL, V78, P1360
GRANOVETTER M, 1983, SOCIOLOGICAL THEORY, V1, P203
```

```
Newman, M. E. (2001). Scientific collaboration networks.
II. Shortest paths, weighted networks, and centrality. Physical review E, 64(1), 016132.

M.E.J. Newman, preceding paper, Phys. Rev. E 64, 016131 (2001).
```


3. Problems with bibliographic data.

Problem 2: Different cultures

Examples of diverse citation practices

- Vol, Issue, Pages
- Paper number
- Citation without paper title

Volume
Issue
First page

Freeman, L. C., & White, D. R. (1993). Using Galois lattices to represent network data. Sociological methodology, 127-146.

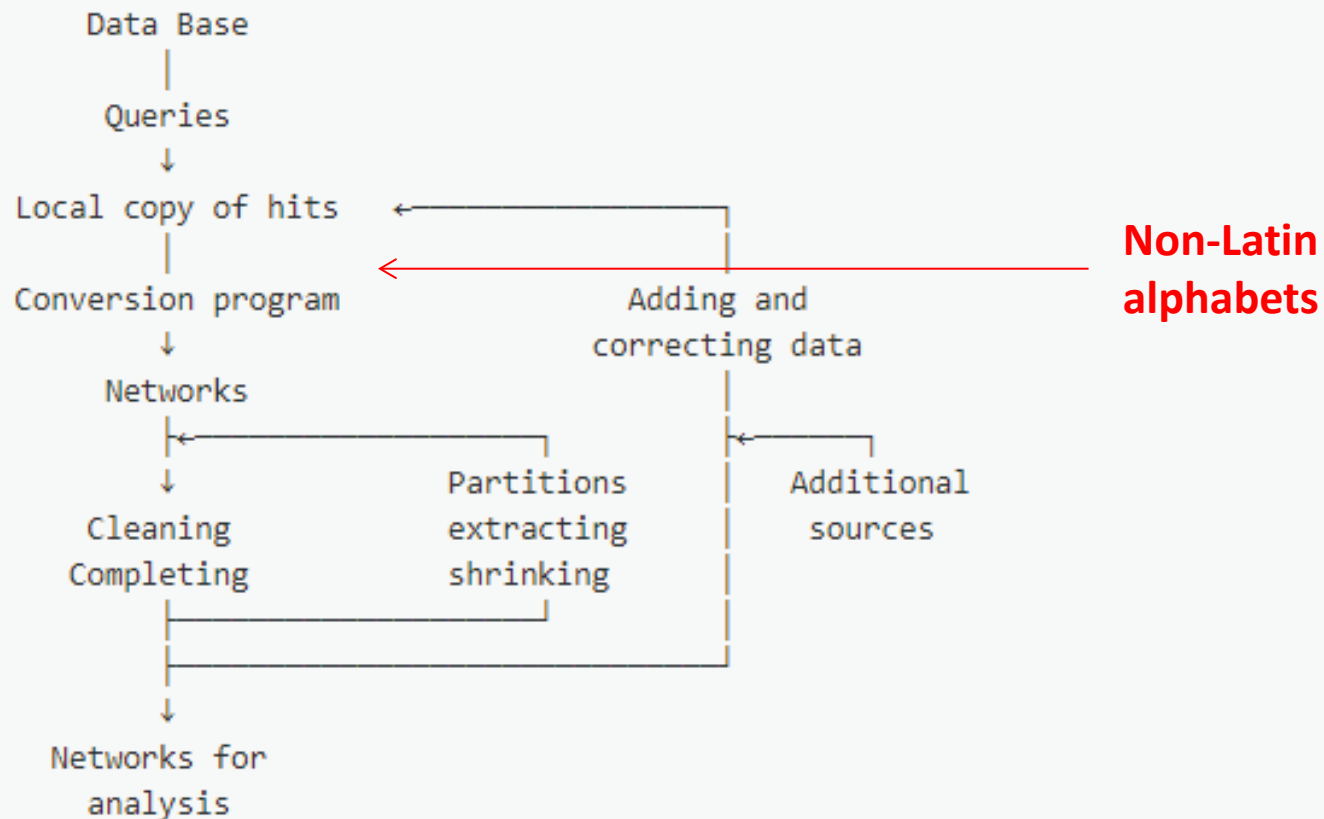
Newman, M. E. (2001). Scientific collaboration networks.
II. Shortest paths, weighted networks, and centrality. Physical review E, 64(1), 016132.

P. Erdos and A. Renyi, Publ. Math. Inst. Hung. Acad. Sci. 5, 17, 1960.

Volume
First page
Year

3. Problems with bibliographic data.

Problem 3: Non-Latin alphabets



3. Problems with bibliographic data.

Problem 3: Non-Latin alphabets

Some names can be written in several languages – the procedure of author disambiguation is needed.

Cyrillic to Latin (Unicode, automatic transcription)

R, stringi library

```
> tail(N)
[1] "ГОМЗИН А"      "НЕДУМОВ Я"      "IVANOV I"      "АСТРАХАНЦЕВ Н"
[5] "ТРИПУТИНА В"   "МАКАГОНОВА Н"
> tail(R)
[1] "GOMZIN A"      "NEDUMOV A"      "IVANOV I"      "ASTRAHANCEV N"
[5] "TRIPUTINA V"   "MAKAGONOVA N"
```

Problems with
character "Ъ"

```
> N[44]
[1] "ЗОРЪКИНА К"
> R[44]
[1] "ZOR'KINA K"
> utf8ToInt(R[44])
[1] 90 79 82 697 75 73 78 65 32 75
> T <- sapply(R,function(w)gsub(intToUtf8(697),"",w),USE.NAMES=FALSE)
> T[44]
[1] "ZOR'KINA K"
> utf8ToInt(T[44])
[1] 90 79 82 39 75 73 78 65 32 75
```

3. Problems with bibliographic data.

Problem 3: Non-Latin alphabets

Transliteration ([different approaches](#))

Cyrillic		Latin		Unicode	
А	а	A	a		
Ä	ä	Ä	ä	00C4	00E4
Ǟ	ǟ	Ǟ	ǟ	00C4+0323	00E4+0323
Ǻ	ǻ	Ǻ	ǻ	0102	0103
Ā	ā	Ā	ā	0100	0101
Æ	æ	Æ	æ	00C6	00E6
Á	á	Á	á	00C1	00E1
À	à	À	à	00C5	00E5
Б	б	B	b		
В	в	V	v		
Г	г	G	g		

Пётр Ильич Чайковский

English: Pyotr Ilyich Tchaikovsky

German: Pjotr Iljitsch Tschaikowski

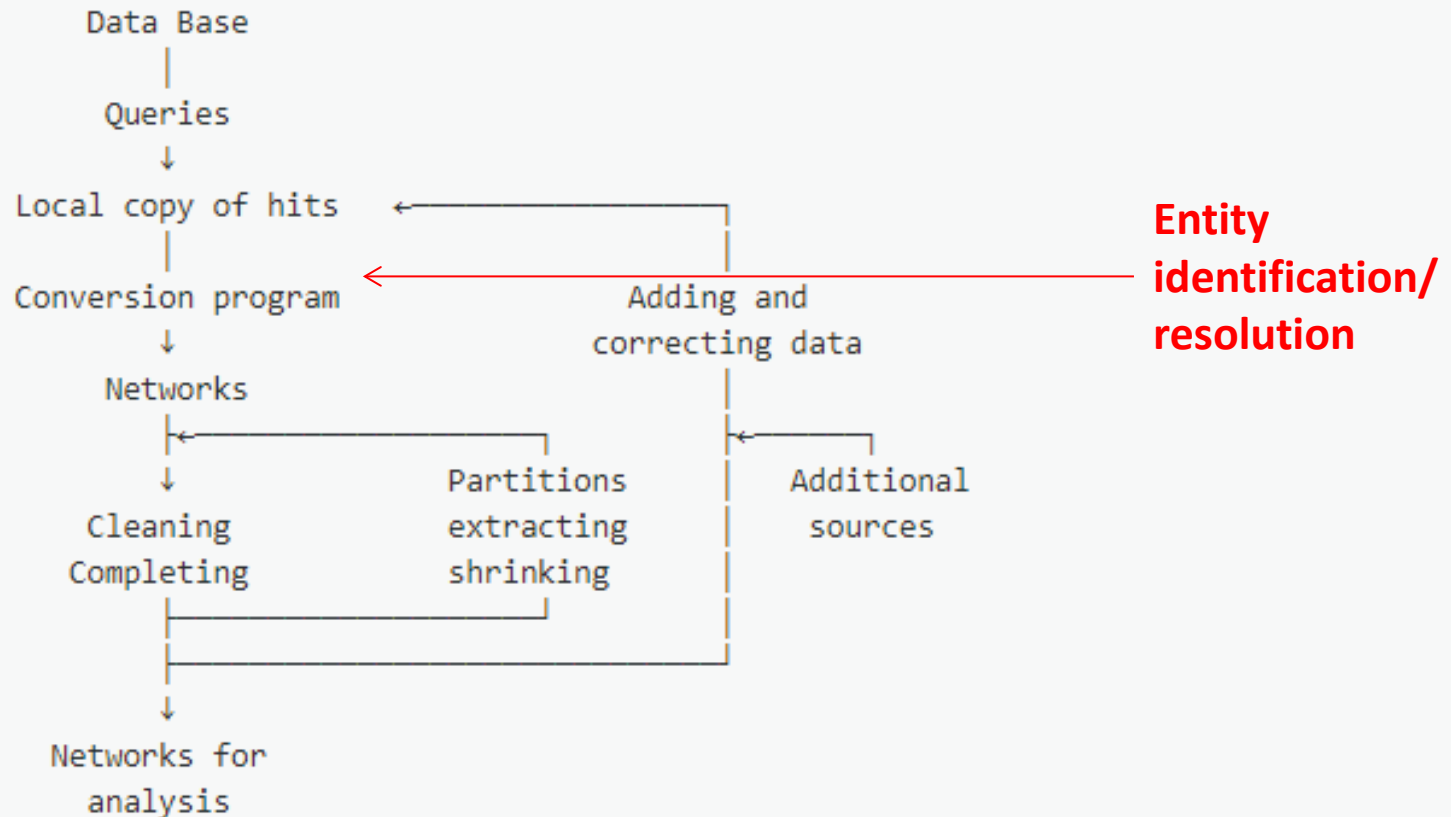
French: Piotr Ilitch Tchaïkovski

Spanish: Piotr Ilich Chaikovsky

Italian: Pëtr Il'ič Čajkovskij

3. Problems with bibliographic data.

Problem 4: Entity identification/resolution



3. Problems with bibliographic data

Problem 4: Entity identification/resolution

Author names

Lorenzo Bartolini from [Letters to Juliet](#)

- Many ways to write the name
"Krivoshe\u{i}n, Leonid Evgen\cprime evich" (using the TeX codes) = 20 distinct name variations for this author.
- Chinese, 100 names – ["three Zhang, four Li"](#) (there are at least 623 different mathematicians with the name Zhang, Li in the MathSciNet Database)



Hanwen Zhang

Universidad Santo Tomás

Подтвержден адрес электронной почты в домене usantotomas.edu.co

[estadística](#) [series de tiempo](#) [estadística bayesiana](#)

Цитируется: 204958



Yiguo Zhang

Professor of Cell Biochemistry and Gene Regulation

[cell biology](#) [gene regulation](#) [transcription factors](#) [live-cell imaging](#)

Цитируется: 203158



Yun Zhang

Professor of Geomatics, University of New Brunswick

Подтвержден адрес электронной почты в домене unb.ca

[Remote Sensing](#) [Image Processing](#) [Computer Vision](#) [Photogrammetry](#)

Цитируется: 121506



Zibin Zhang

Zhejiang University

Цитируется: 117994

3. Problems with bibliographic data

Problem 4: Entity identification/resolution

Author names

- Some data bases are trying to standardize the names (DBLP, ZB, ResearcherId).

[MathSciNet](#); [Orcid](#) - Enter author name in Search field

[Scopus](#); [eLibrary](#) - Click on author's name

and take the number after authorid

```
https://orcid.org/0000-0002-0240-9446
```

```
https://elibrary.ru/author_items.asp?authorid=155240
```

- Variations in the first names: Sort (last name, first name).
- Multi-alphabet (names written in different languages) - convert names to selected alphabet or use "dictionary".

[AMS approach](#) – look for details.

3. Problems with bibliographic data

Problem 4: Entity identification/resolution

Journal names and Books

- Form a key from initials of journal name and sort (key, journal name)
- International Standard Serial Number [ISSN](#); Digital Object Identifier [DOI](#); International Standard Book Number [ISBN](#).

Keywords

Provided in data or extracted from the text (title, abstract). Key phrases.

- Errors (typos) in the data base -- correct them in your copy of the data base data.

3. Problems with bibliographic data

Problem 4: Entity identification/resolution

Synonymy: Different names for the same units (people):

- **Otfried Cheong** (formerly **Otfried Schwarzkopf**): German computational geometer working in South Korea at KAIST
- **Michel Marie Deza** (formerly **Mikhail Efimovich Tylkin**): a Soviet and French mathematician, specializing in combinatorics, discrete geometry and graph theory.
- Borštnik, N. S. Mankoč; Mankoč Borštnik, N.; Mankoč-Borštnik, Norma; Mankoč Borštnik, Norma Susana; Mankoc-Borstnik, N.S.; and Mankoč Borštnik, N.S. = same author.
- NUCLEIC ACIDS RES, NUCL ACIDS RES, NUCLEIC ACIDS RES S, NUCLEIC ACIDS RES S2, NUCL ACID RES, NUCL ACIDS RES S2, NUCL ACIDS S SER, NUCL ACIDS RES S, NUCL AC RES, NUCLEIC ACIDS RES S1, Nucleic Acids Res, NUCL ACIDS RES S1 = same journal

Homonymy (ambiguity): Same names for different units (people).

- Smith, John W. - publications of the author(s) with this name spanned from 1868 to 2007.

3. Problems with bibliographic data

Problem 4: Entity identification/resolution

Description of equivalences - a condensed dictionary

id	canon	alter1	alter2	alter3	alter4	alter5	alter6
ORCID:0000-0002-0240-9446	Batagelj, Vladimir	Batagelj, Vlado	Batagelj, V	MR:32440	Scopus:56037441100		
Scopus:35615877200	Batagelj, Valentin	Batagelj, V					
ORCID:0000-0003-4467-7075	Cheong, Otfried	Schwarzkopf, Otfried	Cheong, O	Schwarzkopf, O	Scopus:57191986875		
MR:57370	Deza, Michel-Marie	Deza, MM	Deza, M	Deza, Mikhail	Scopus:7003745115		
MR:57370	Deza, Michel-Marie	Tylkin, Mikhail Efimovich	Tylkin, ME	Тылкин, Михаил Ефимович	Тылкин, МЕ	Деза, Мишел	Деза, М
ORCID:0000-0002-4294-9017	Zweig, Katharina Anna	Zweig, KA	Zweig, K	Lehmann, Katharina Anna	Lehmann, K	Scopus:25928162000	
eLib:696348	Maltseva, Daria	Мальцева, Дарья Васильевна	Мальцева, ДВ	Мальцева, Д	Maltseva, DV		
eLib:155240	Maltseva, Diana	Мальцева, Диана Васильевна	Мальцева, ДВ	Мальцева, Д	Maltseva, DV		

3. Problems with bibliographic data

Problem 4: Entity identification/resolution

How to deal with synonymy and homonymy:

- **normalization** - at data entry
- **standardization** - use standards whenever possible ([ORCID](#), [DOI](#), [ISBN](#), [ISSN](#) , standard abbreviations [JAS](#), [LTWA](#), [WoS](#), [Caltech](#))

- **"dictionaries"**

When the unit names are extracted from the text the so called *stopwords* are omitted.

The equivalence is automatically determined using stemming or [lemmatization](#) – replacement by the canonical forms of words.

- [lemmatization lists](#) (dictionaries)
- [keywords](#) - keyword recommendations
- **for synonymy**: sort labels of units, manually/visually identify equivalent units, create partition, (shrink) equivalent units.
- **for homonymy**: correct the data in your copy of the data base.

3. Problems with bibliographic data

Problem 4: Entity identification/resolution

ISI names: The usual ISI name of a work (field CR in WoS)

```
LEFKOVITCH LP, 1985, THEOR APPL GENET, V70, P585
```

has the following structure

```
AU + ', ' + PY + ', ' + SO[:20] + ', V' + VL + ', P' + BP
```

In WoS the same work can have different ISI names. To improve the precision the program WoS2Pajek supports also short names. They have the format:

```
LastNm[:8] + ' ' + FirstNm[0] + '(' + PY + ')' + VL + ':' + BP
```

For example: `LEFKOVIT L(1985)70:585`

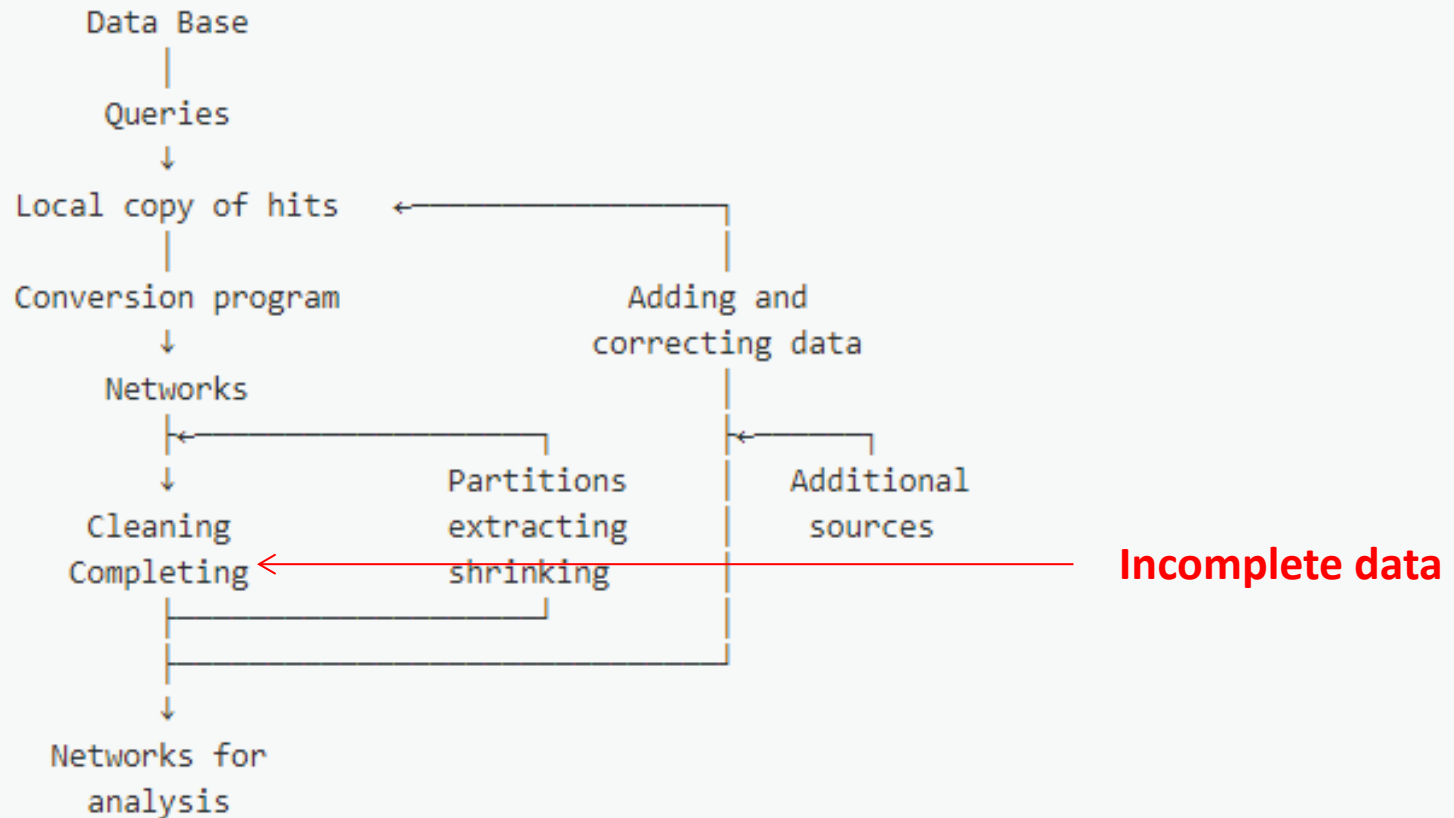
From the last names with prefixes VAN, DE, . . . the space is deleted.

```
CANTANZARO M, 2005, PHYS REV E, V71, UNSP 027103  
CANTAZARO M, 2005, PHYS REV E, V71, UNSP 056104  
CATANZARO M, 2005, PHYS REV E 2, V71, ARTN 056104
```

The best/final solution is to enter data in bibliographic data base in standardized way resolving homonyms.

3. Problems with bibliographic data.

Problem 5: Incomplete data



3. Problems with bibliographic data.

Problem 5: Incomplete data / Boundary problem

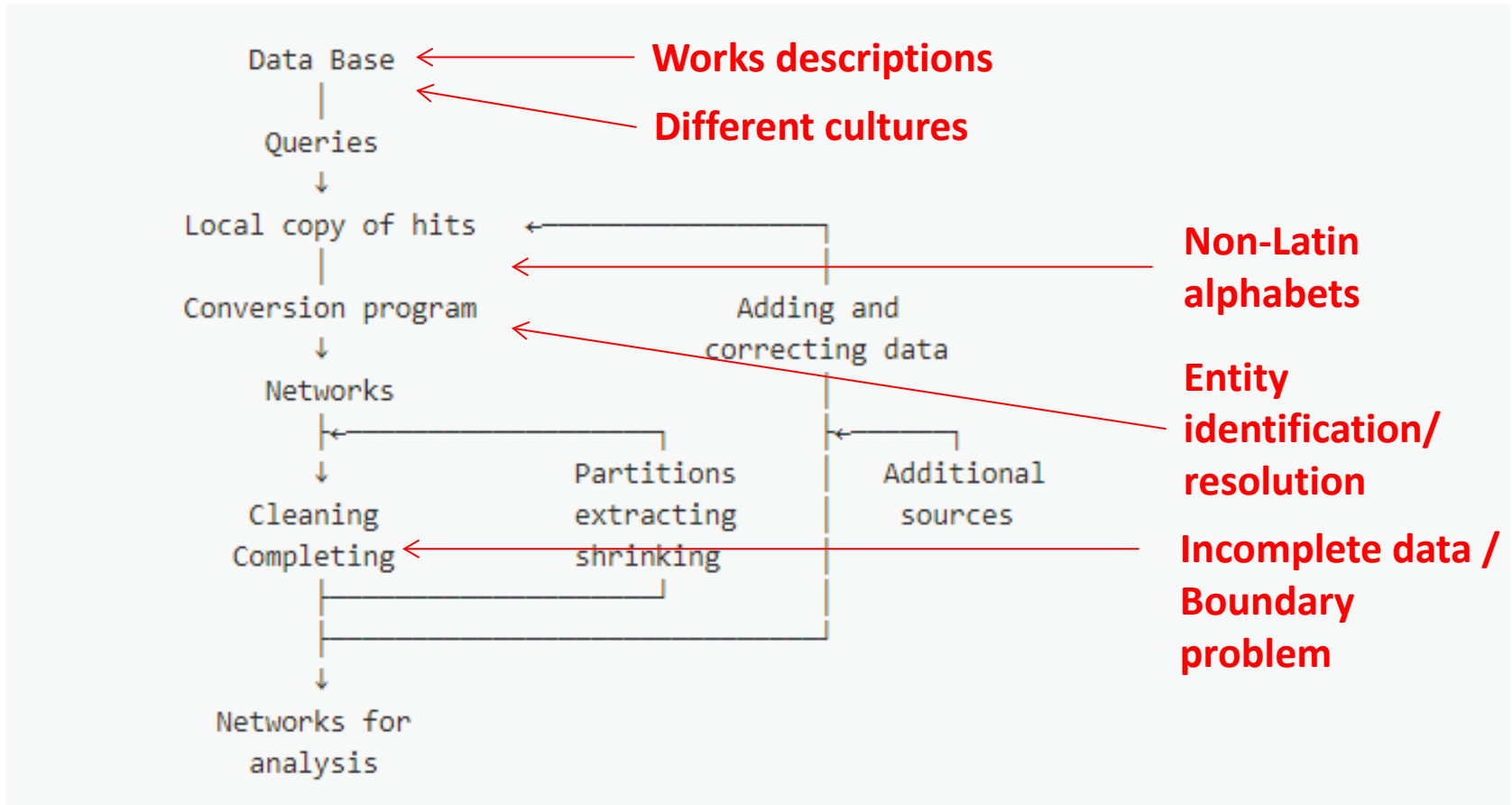
Incomplete data

- missing data: some parts of the units (works) descriptions which are not presented (additional authors, title, year, volume, issue, pages, publishers, etc.).
 - To add important missing parts of data manually
- Missing data: units (works) important for the studied topic which are not presented in the data set at all.
 - To search for them - look at the most cited works from the references and include them into the analysis

For small bibliographies where we can inspect, accept and "correct" each entry. The "Excel table" approach is sufficient.

- when extracting subset of data
- preliminary citation network analysis; manually completing the important data

3. Problems with bibliographic data



Different iterations are needed before the data set is complete!

4. Tools for collection and maintenance of bibliographic data

Bibliographic tools:

- [JabRef - open source bibliography reference manager](#)
- [Bibliographic management tools](#)
- [Bibliographic Conversion Tools](#)
- [Computer Science and Engineering: Bibliographic Tools](#)
- [12 Best Free Online Bibliography And Citation Tools](#)
- [Bibliographic Tools](#)
- [Bibliographic Software Overview](#)
- [Compare Some of the Popular Bibliographic Software Tools](#)

5. Conversion to networks

For conversion of bibliographic data special programs in languages such as Python and R are written.

Example:

1. export data from WoS
2. combine files into WoS file
3. run WoS2Pajek
4. compute indegrees in citation network