

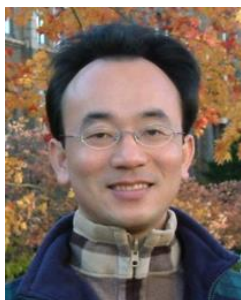
李御璽 教授

銘傳大學資訊工程學系

Big Data Analytics on Social Media (社群媒體大數據分析)

Time: 2015/12/25 (14:00-15:30)

Place: S402, Ming Chuan University



Min-Yuh Day

戴敏育

Assistant Professor

專任助理教授

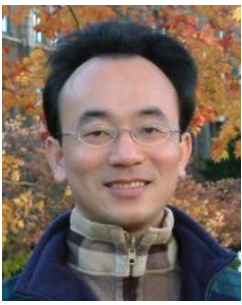
Dept. of Information Management, Tamkang University

淡江大學 資訊管理學系

<http://mail.tku.edu.tw/myday/>

2015-12-25





戴敏育 博士 (Min-Yuh Day, Ph.D.)

淡江大學資管系專任助理教授

中央研究院資訊科學研究所訪問學人

國立台灣大學資訊管理博士

Publications Co-Chairs, IEEE/ACM International Conference on
Advances in Social Networks Analysis and Mining (ASONAM 2013-)

Program Co-Chair, IEEE International Workshop on
Empirical Methods for Recognizing Inference in Text (IEEE EM-RITE 2012-)

Workshop Chair, The IEEE International Conference on
Information Reuse and Integration (IEEE IRI)



Outline

- Big Data Analytics on Social Media
- Analyzing the Social Web:
Social Network Analysis
- NTCIR 12 QALab-2 Task

Social Media



Social Media



Social Media



Facebook



Twitter



Twitter



LinkedIn



Google+



My Space



Tumblr



Bebo



Foursquare



Delicious



Digg



Stumbleupon



Reddit



Technorati



Slashdot



Share this



You Tube



Flickr



Instagram



Pinterest



Deviant Art



Soundcloud



Vimeo



Twylah



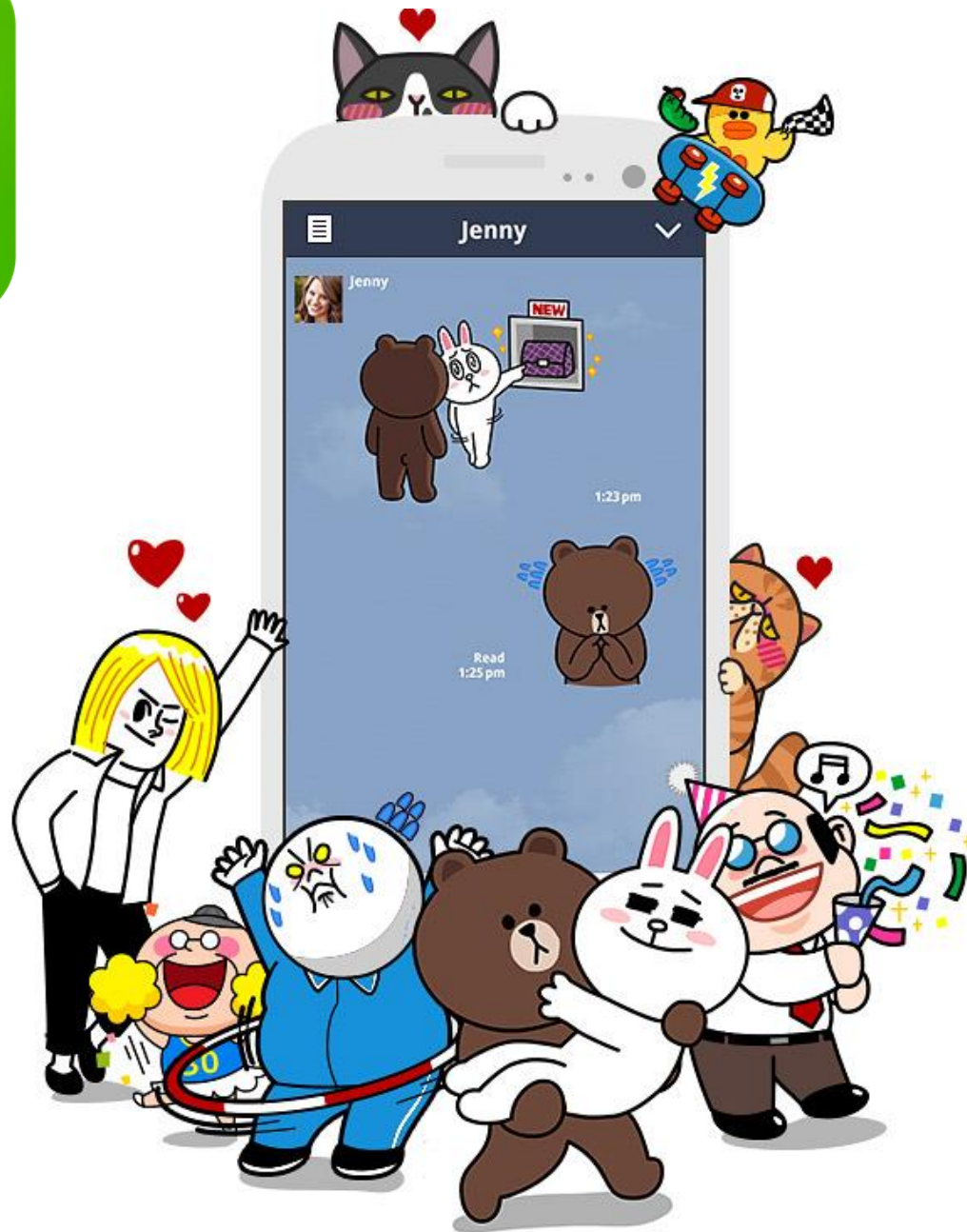
RSS



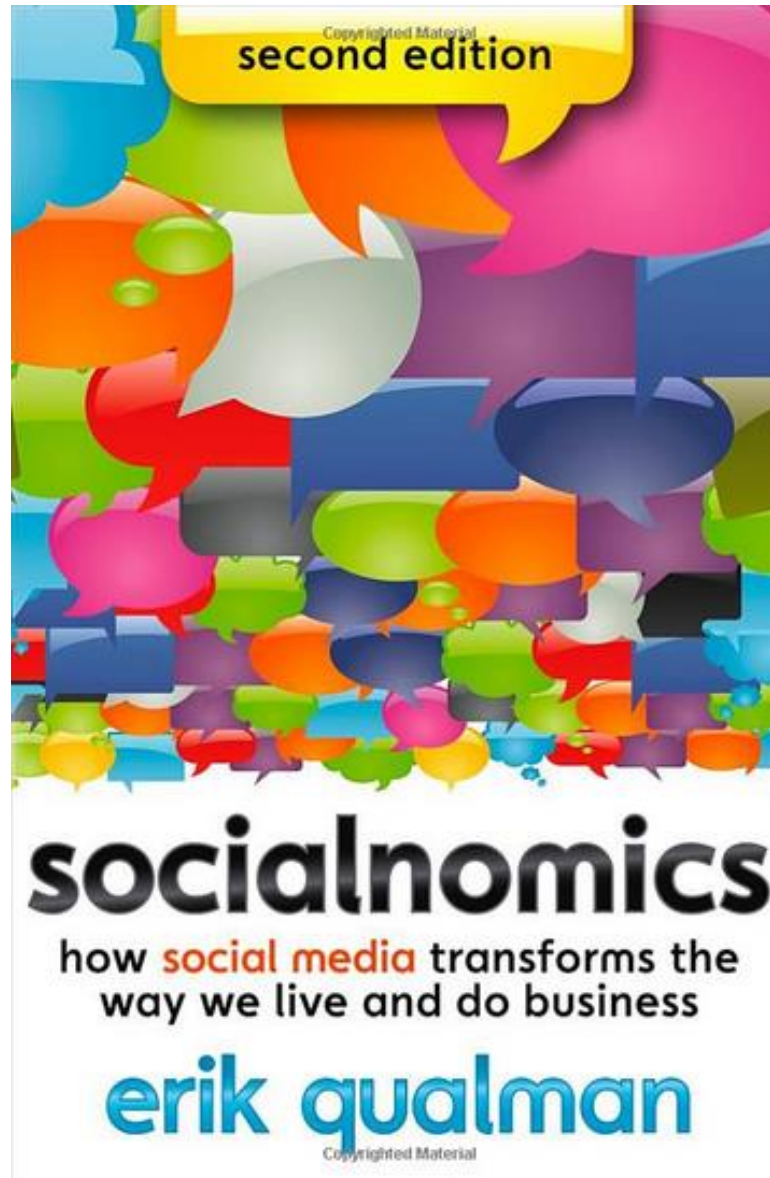
Skype



Line



Socialnomics



Emotions

Love

Anger

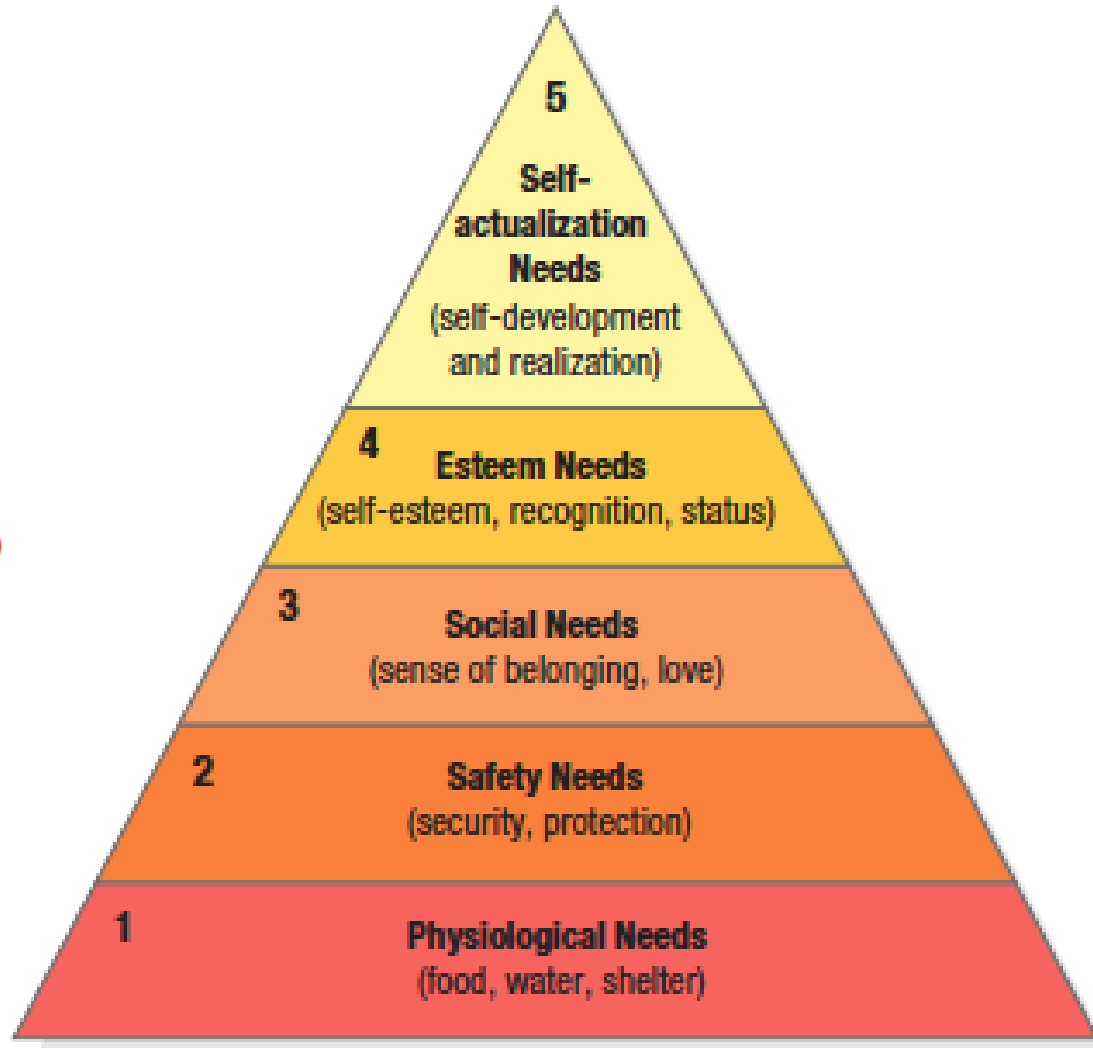
Joy

Sadness

Surprise

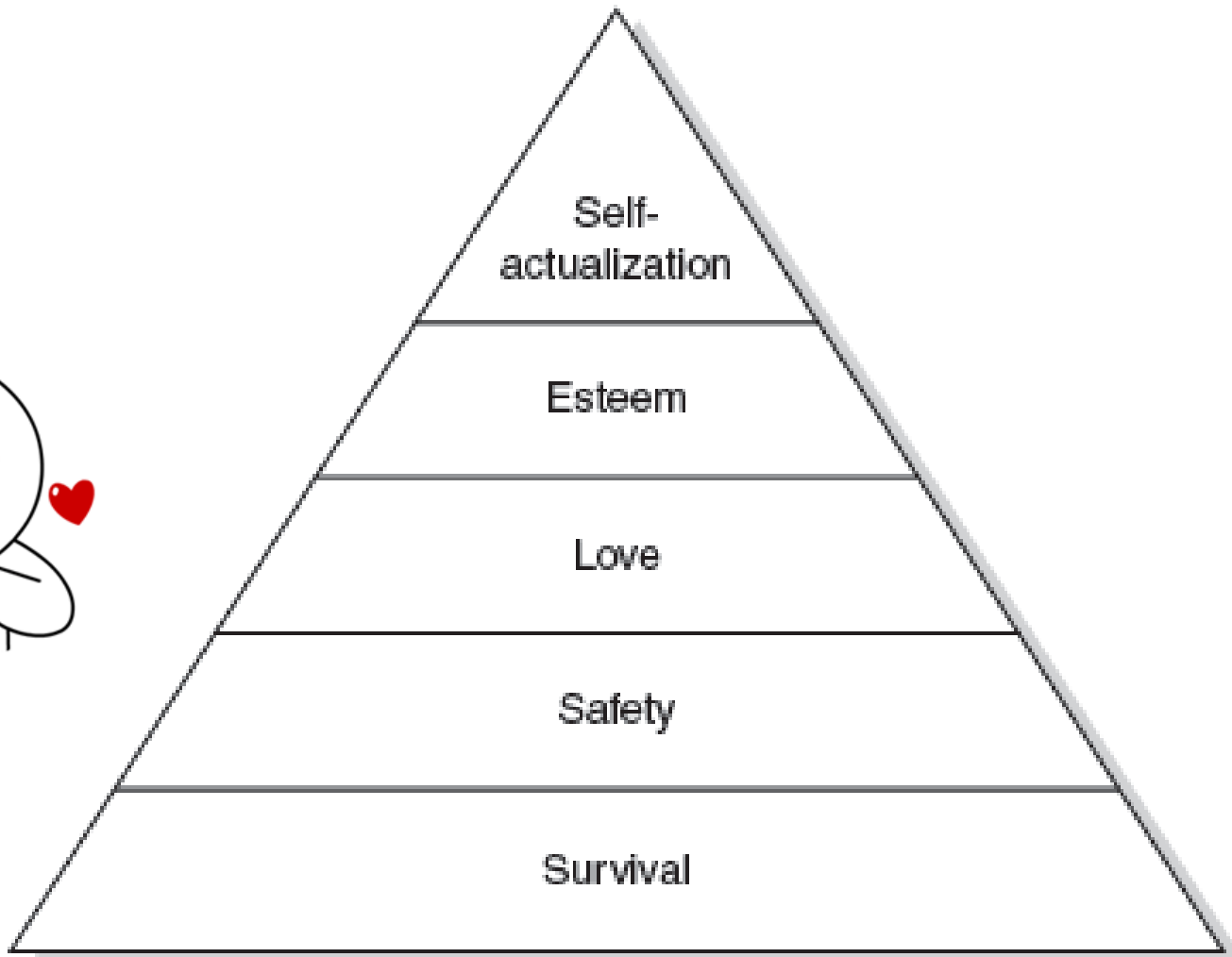
Fear

Maslow's Hierarchy of Needs

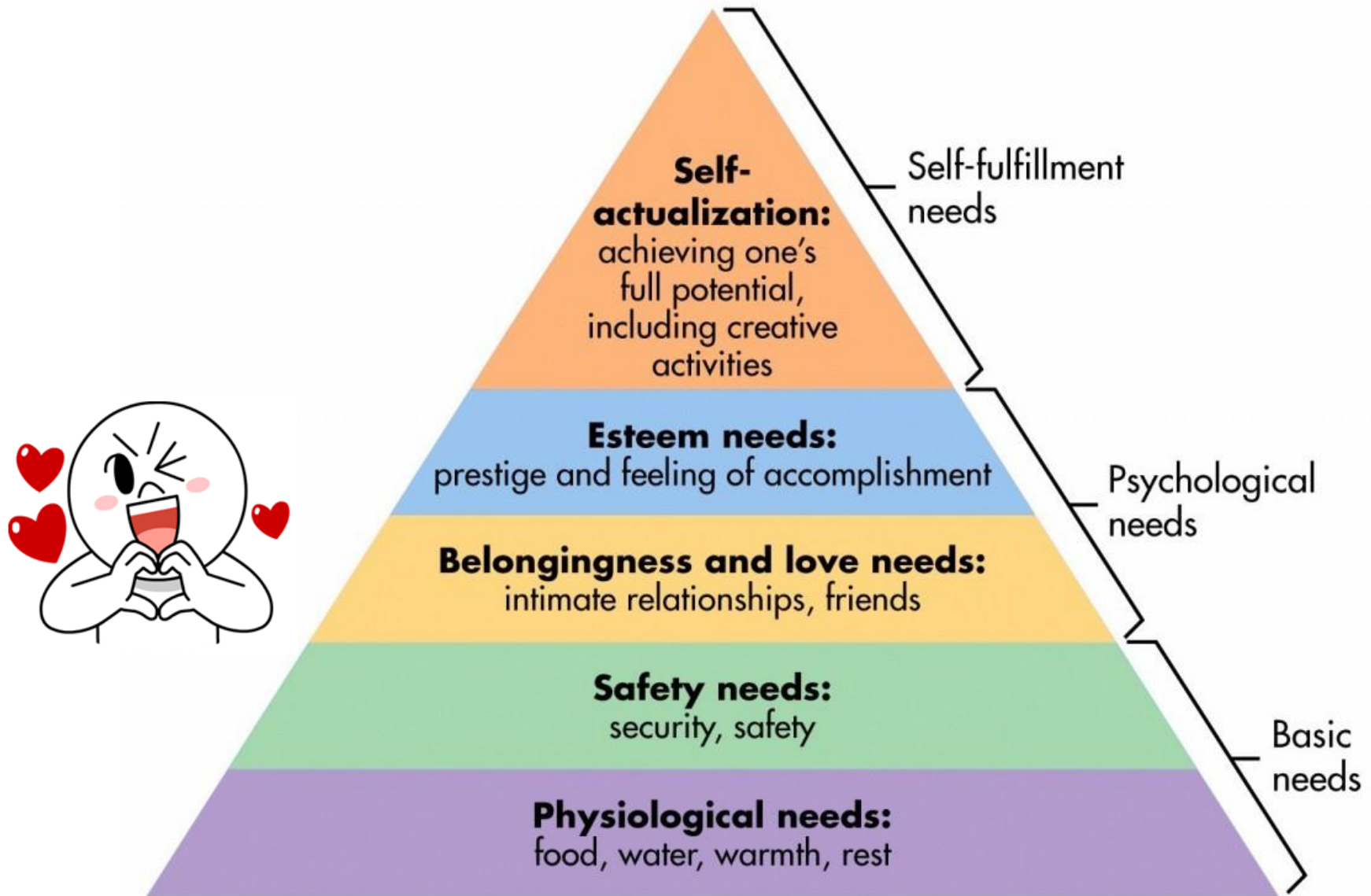


Maslow's hierarchy of human needs

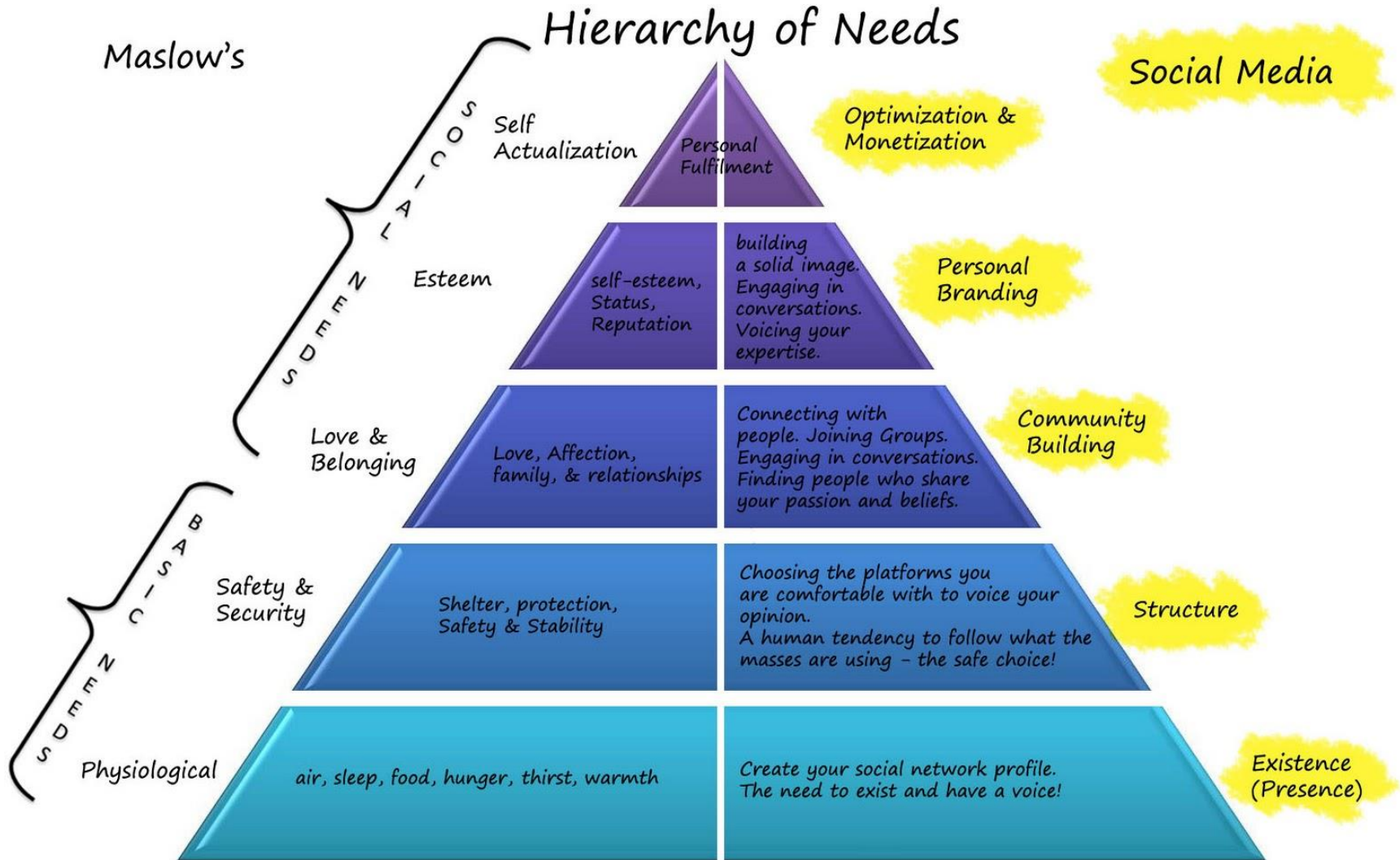
(Maslow, 1943)



Maslow's Hierarchy of Needs

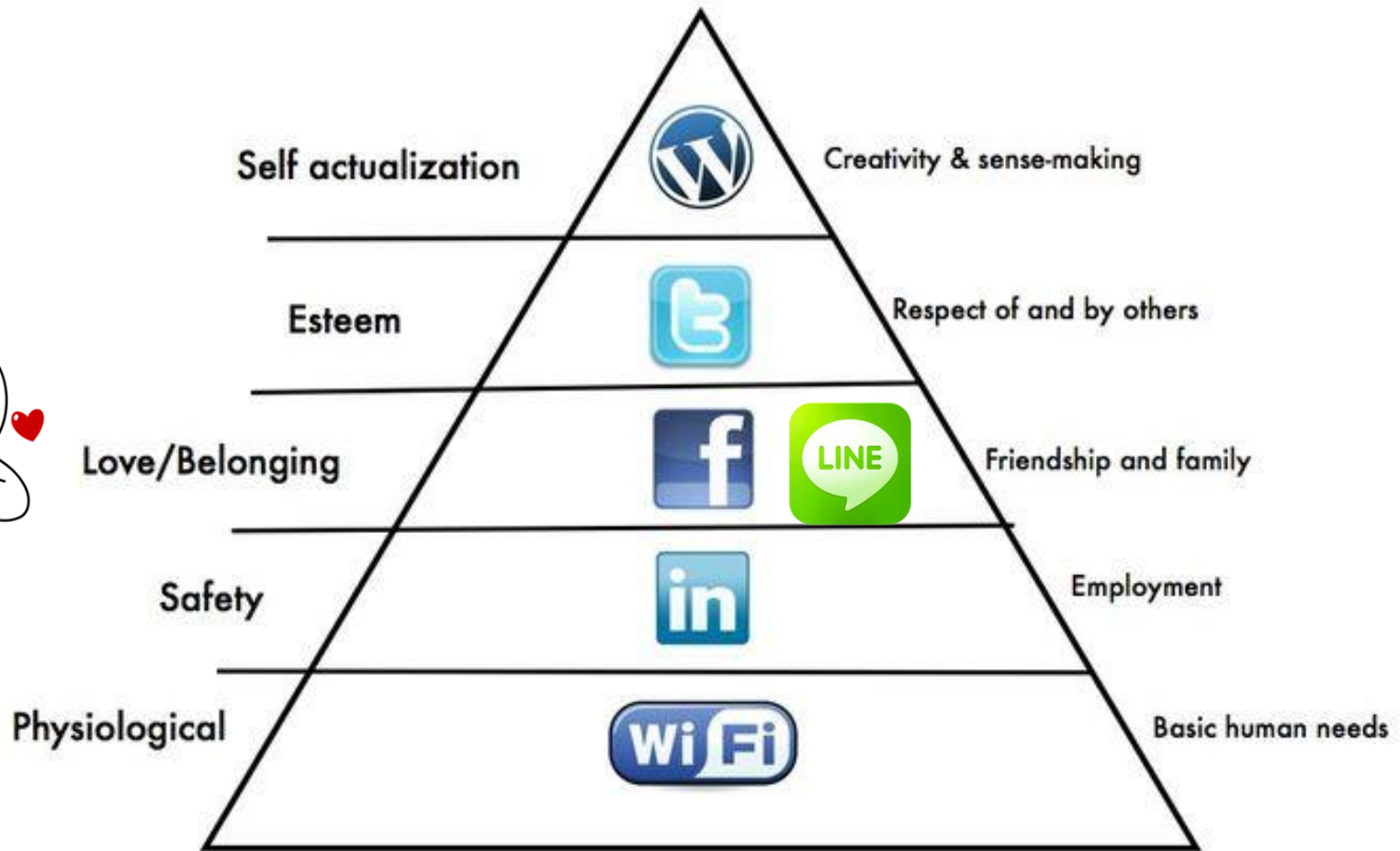


Social Media Hierarchy of Needs



Social Media Hierarchy of Needs - by John Antonios

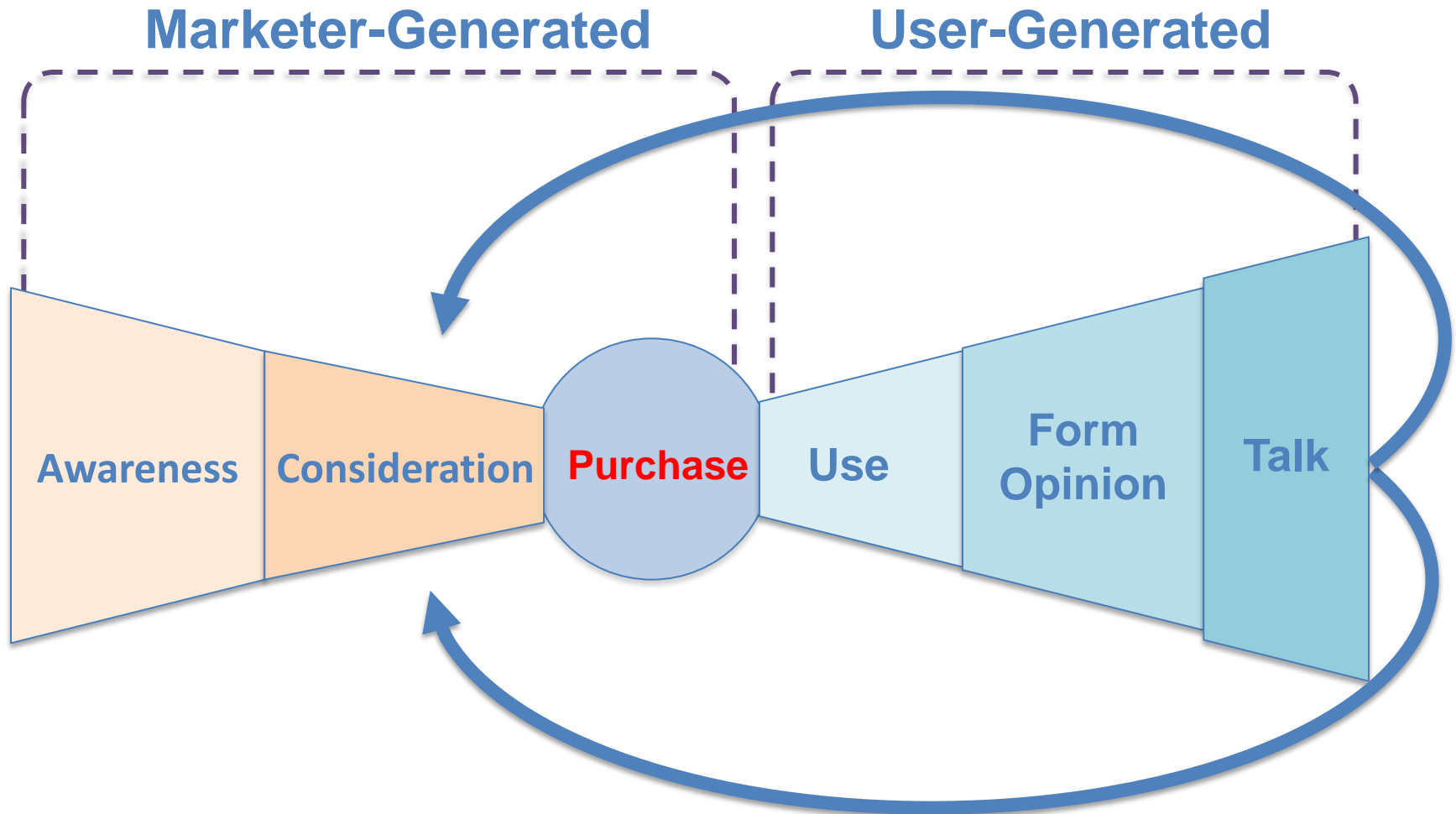
Social Media Hierarchy of Needs



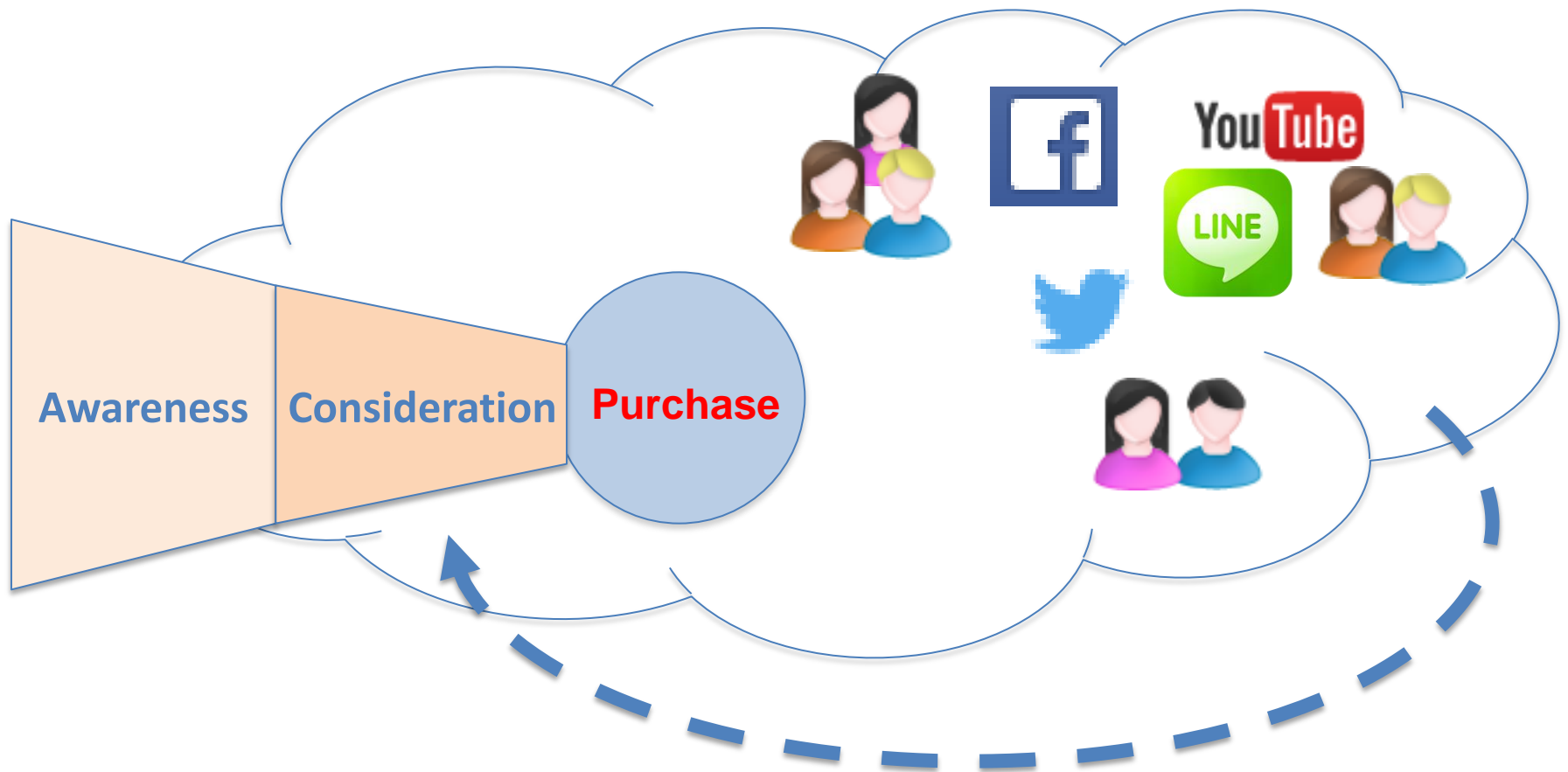
@daveduarte

The Social Feedback Cycle

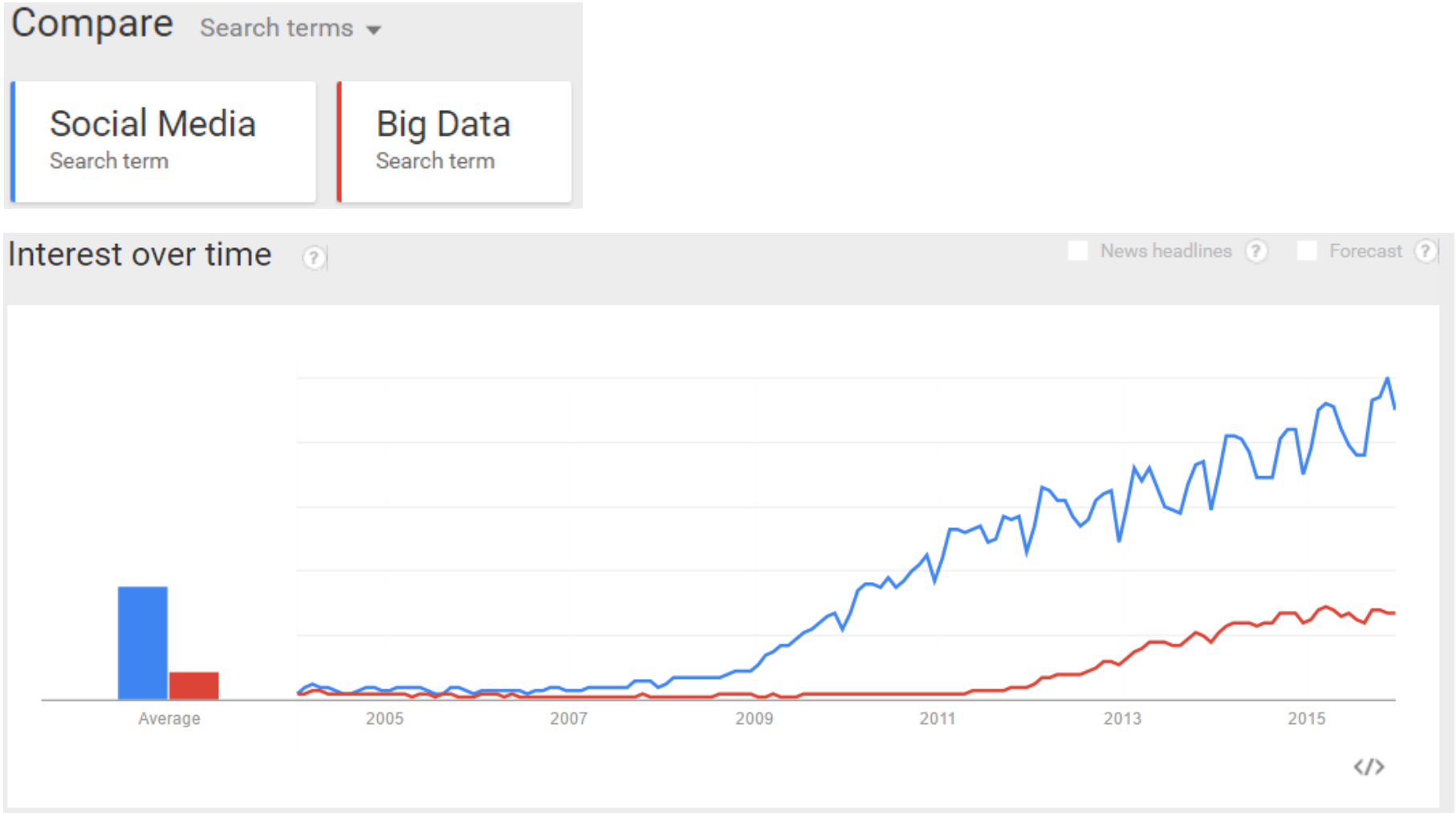
Consumer Behavior on Social Media



The New Customer Influence Path



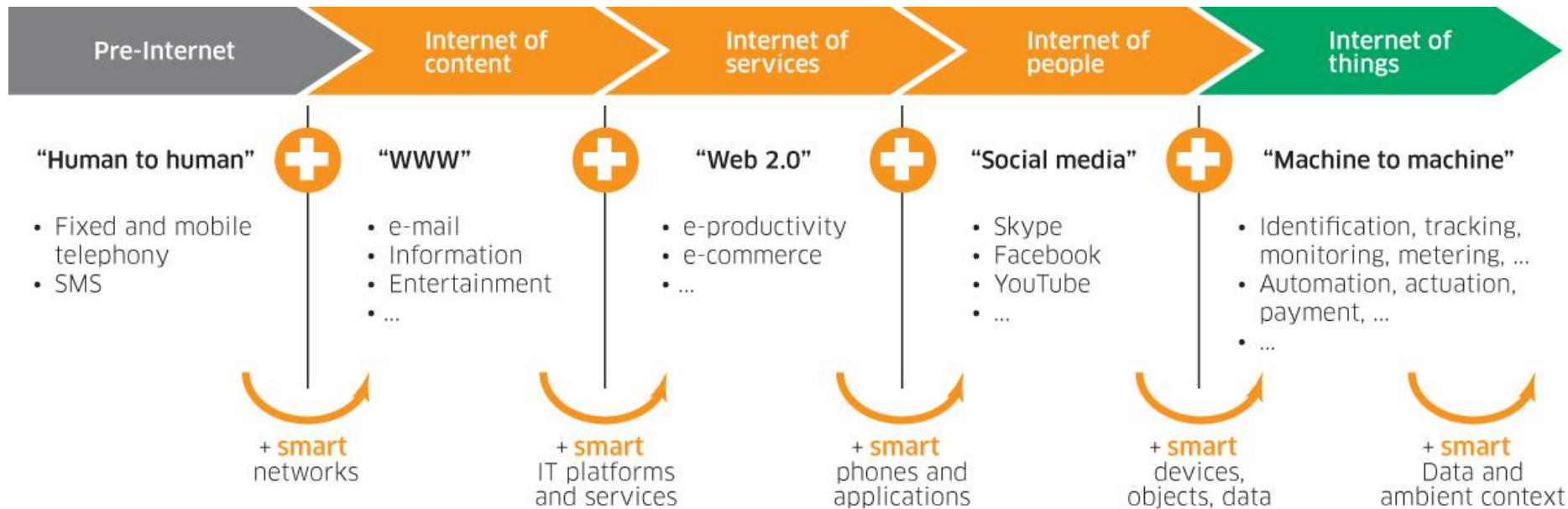
Google Trends on Social Media



Internet Evolution

Internet of People (IoP): Social Media

Internet of Things (IoT): Machine to Machine



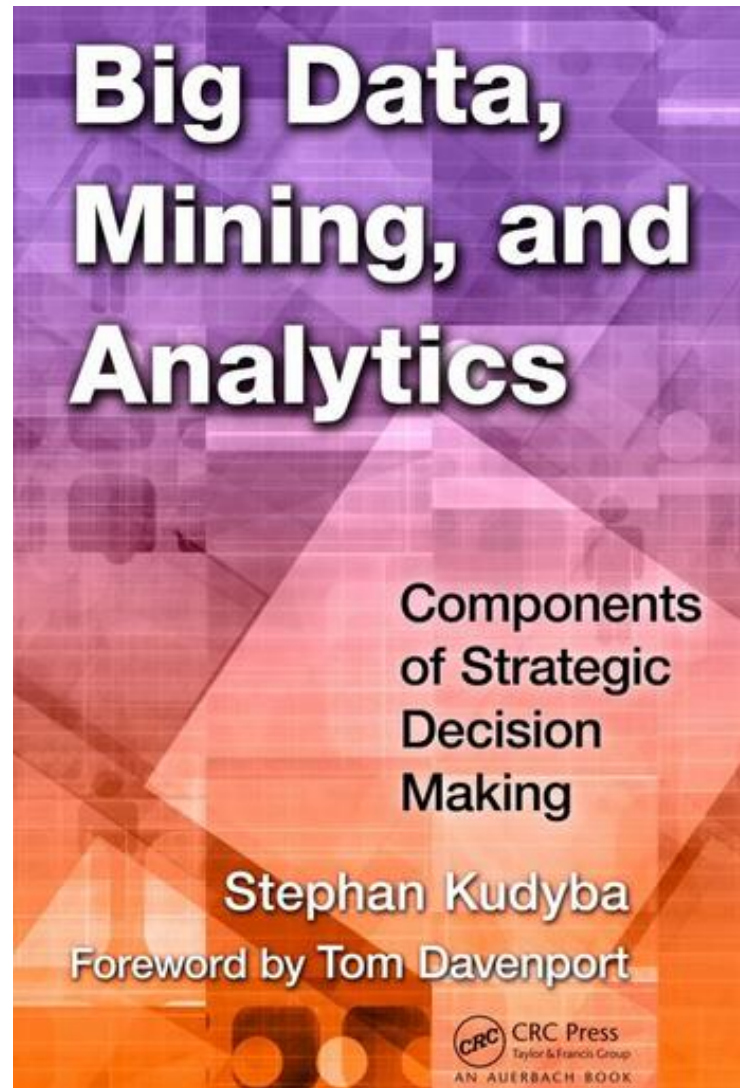
Source: Marc Jadoul (2015), The IoT: The next step in internet evolution, March 11, 2015

<http://www2.alcatel-lucent.com/techzine/iot-internet-of-things-next-step-evolution/>

Business Insights
with
Social Analytics

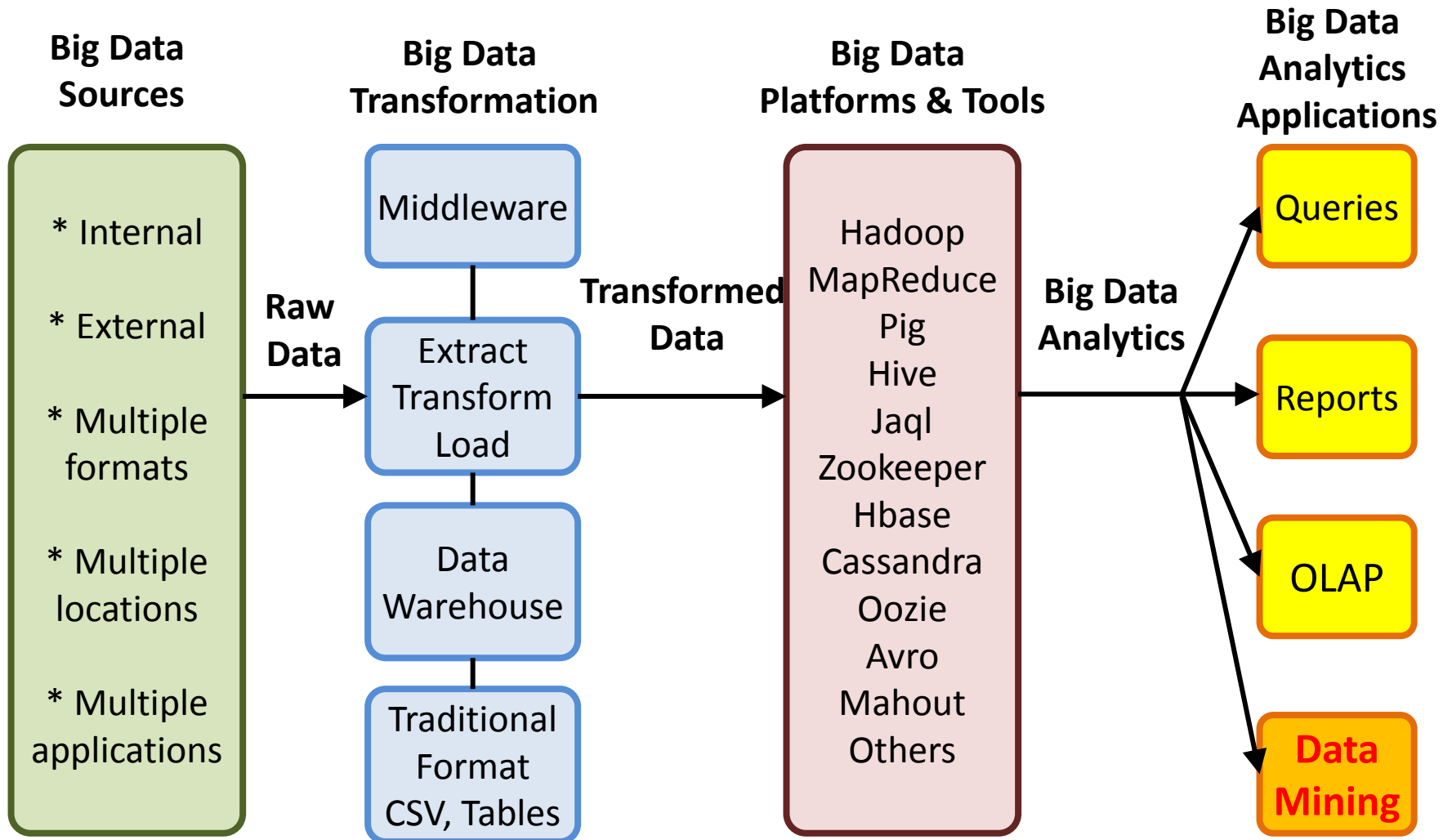
Big Data
Analytics
and
Data Mining

Stephan Kudyba (2014),
Big Data, Mining, and Analytics:
Components of Strategic Decision Making, Auerbach Publications



Source: <http://www.amazon.com/gp/product/1466568704>

Architecture of Big Data Analytics



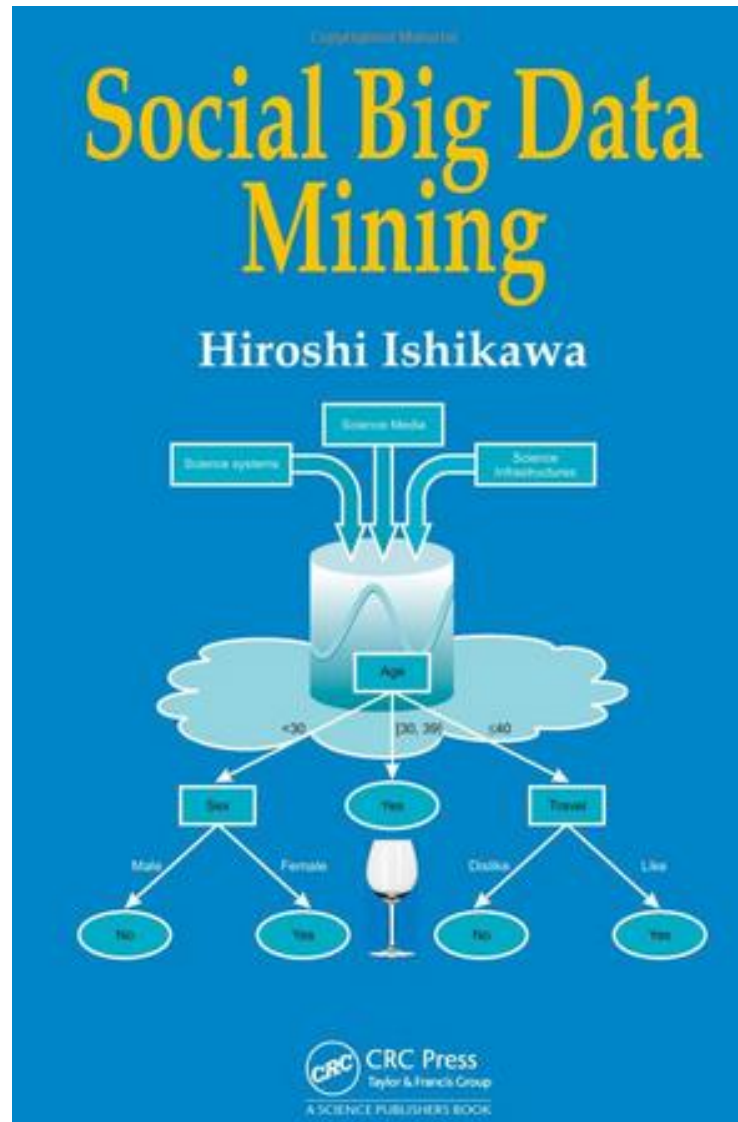
Architecture of Big Data Analytics



Source: Stephan Kudyba (2014), Big Data, Mining, and Analytics: Components of Strategic Decision Making, Auerbach Publications

Social Big Data Mining

(Hiroshi Ishikawa, 2015)

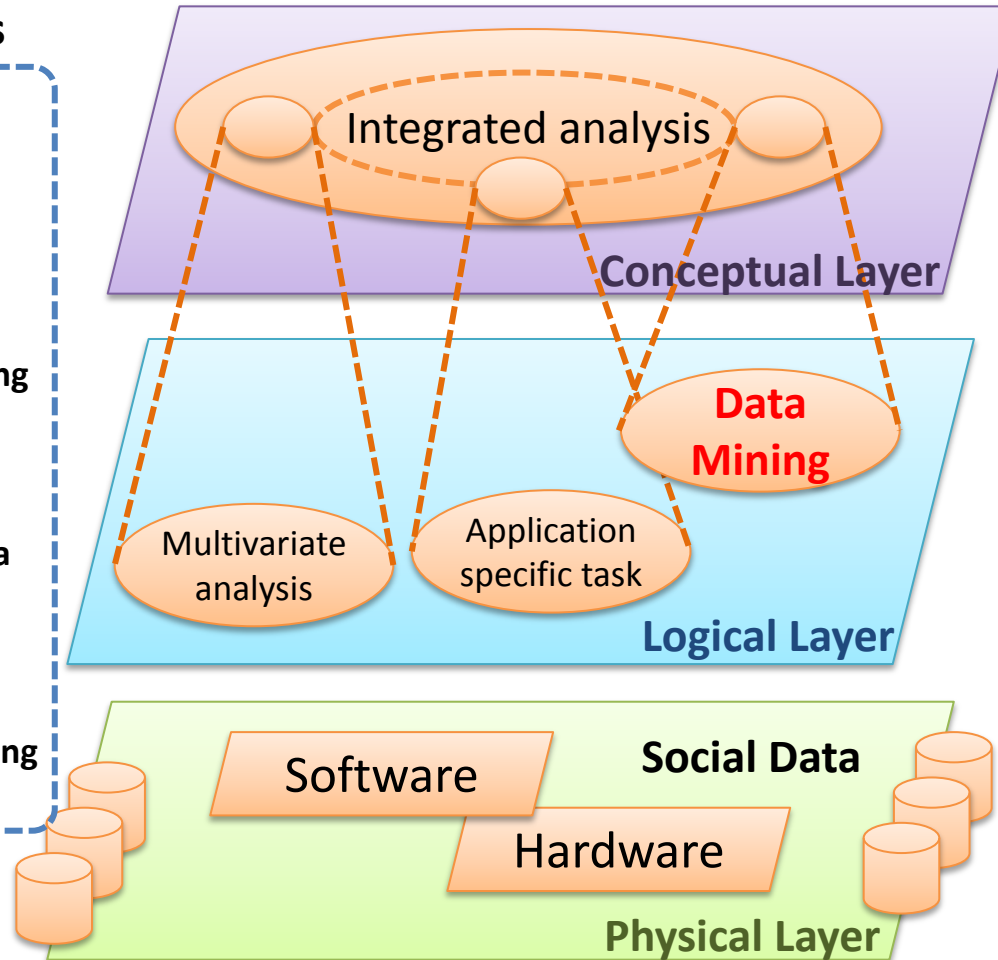


Architecture for Social Big Data Mining

(Hiroshi Ishikawa, 2015)

Enabling Technologies

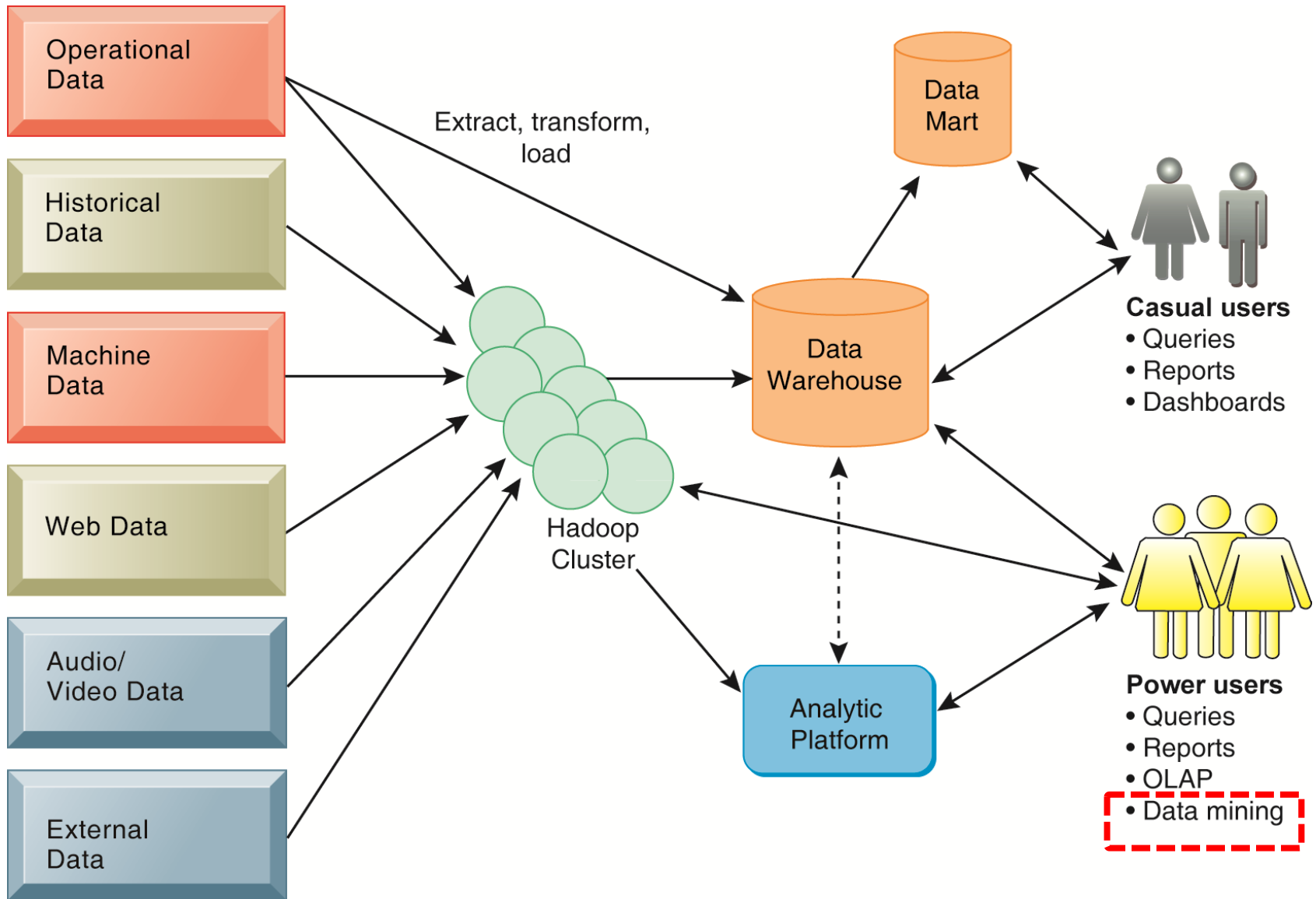
- Integrated analysis model
- Natural Language Processing
- Information Extraction
- Anomaly Detection
- Discovery of relationships among heterogeneous data
- Large-scale visualization
- Parallel distributed processing



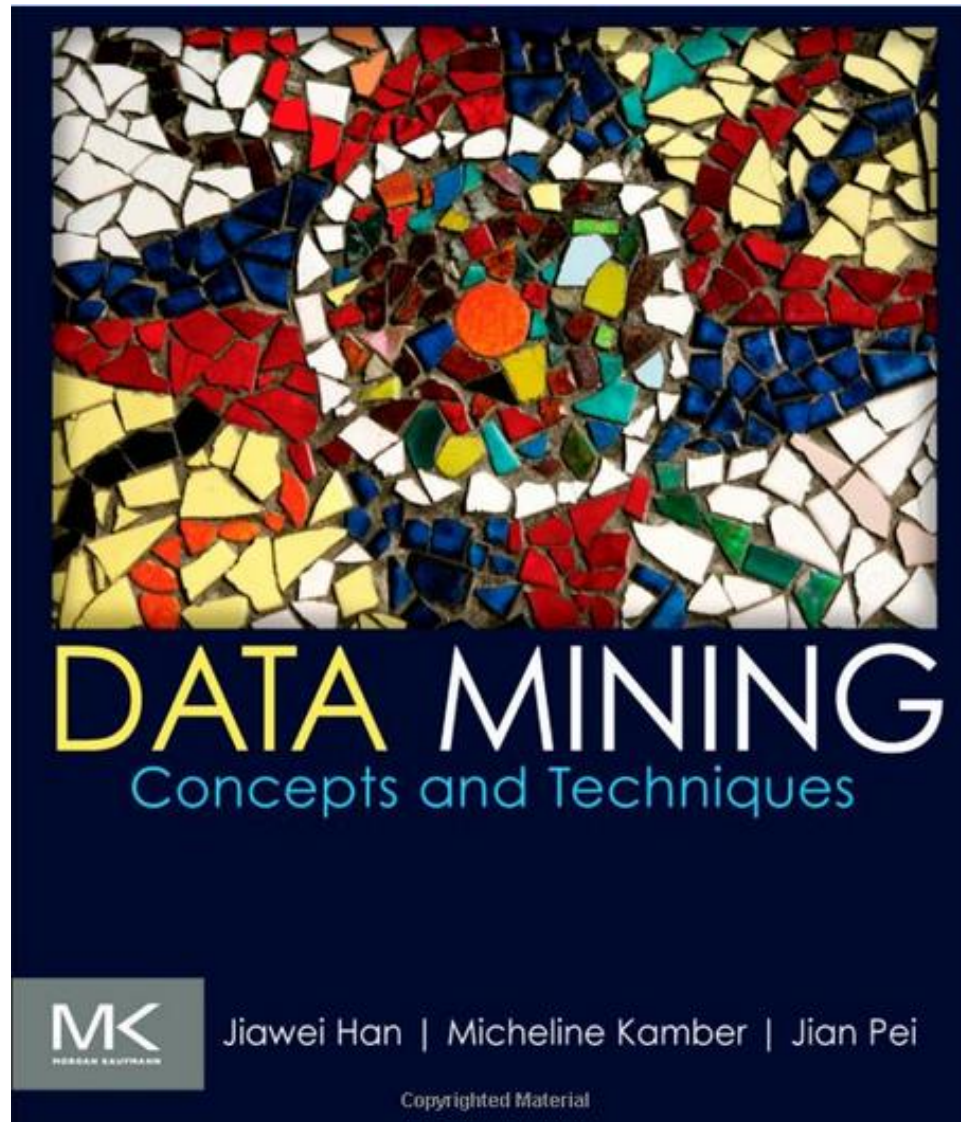
Analysts

- Model Construction
- Explanation by Model
- Construction and confirmation of individual hypothesis
- Description and execution of application-specific task

Business Intelligence (BI) Infrastructure

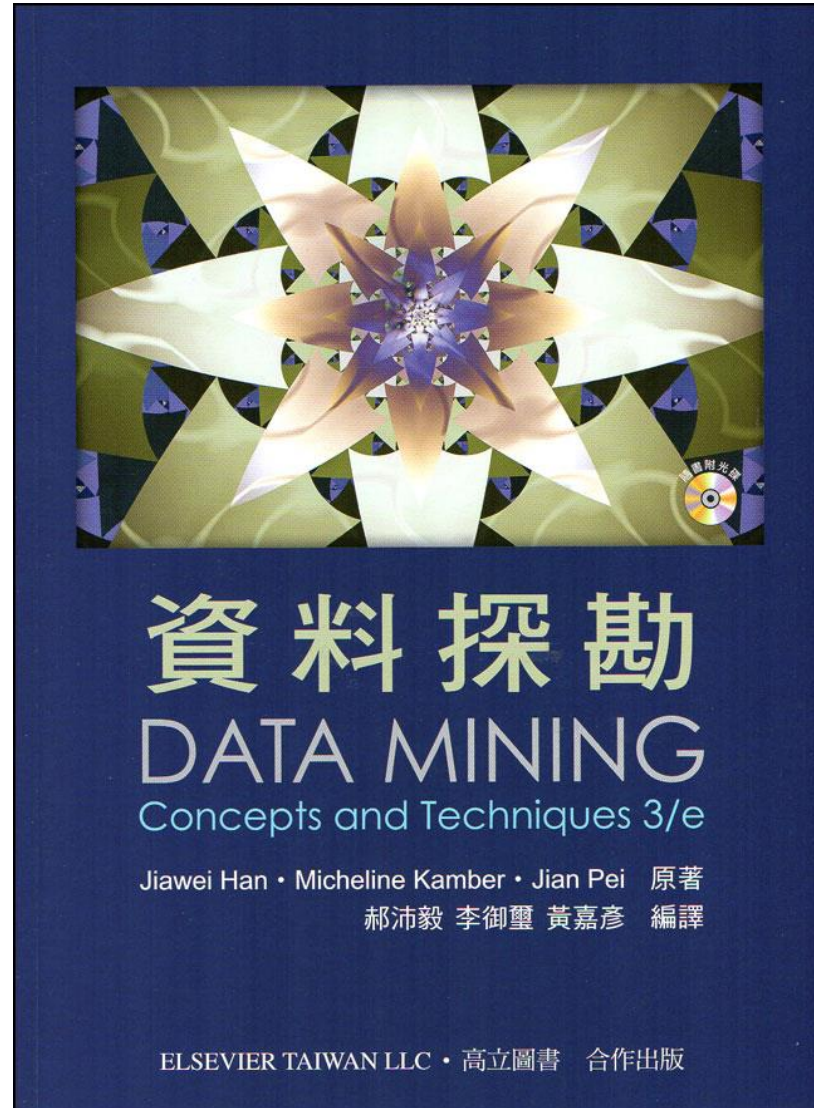


Data Mining



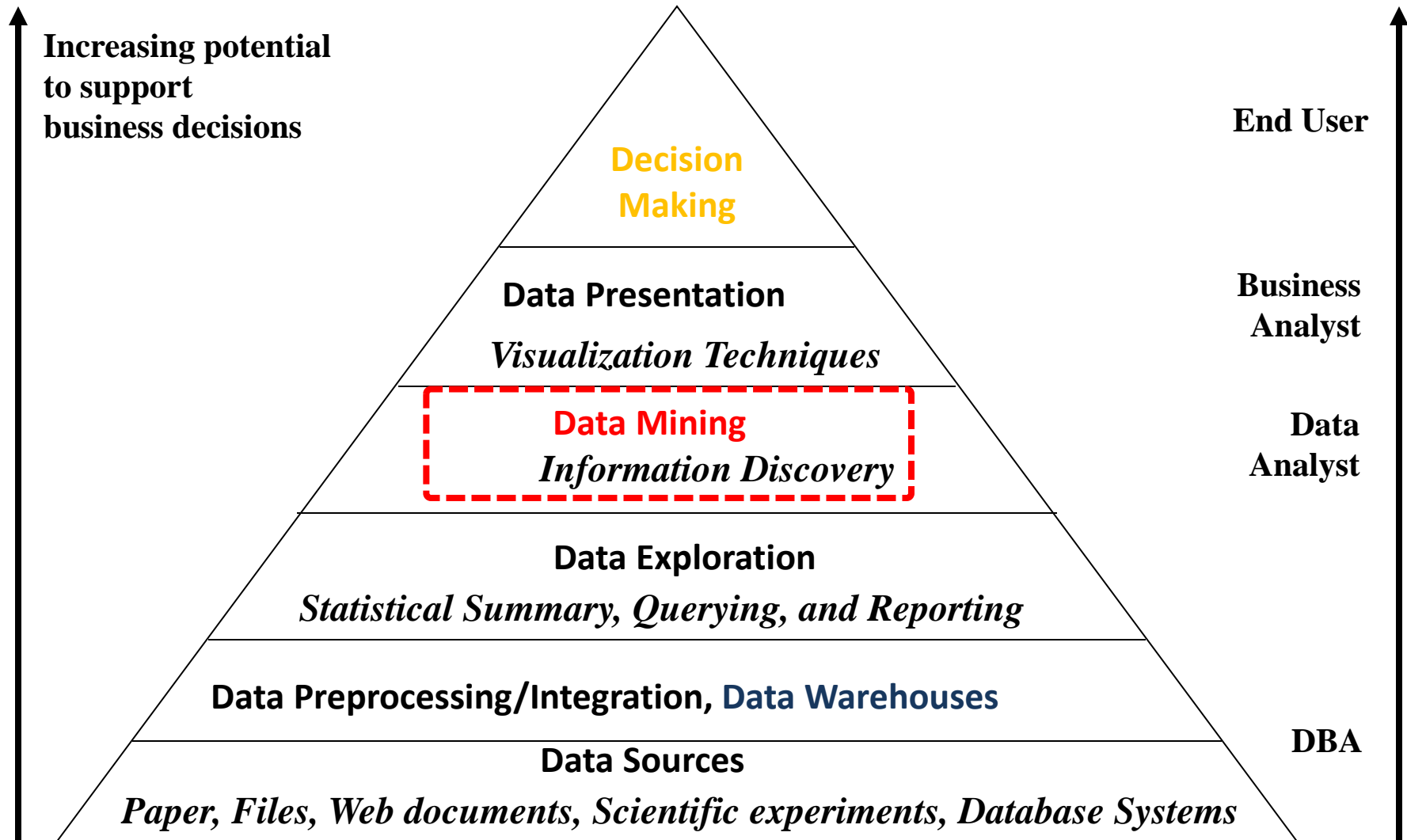
郝沛毅, 李御璽, 黃嘉彥 編譯, 資料探勘

(Jiawei Han, Micheline Kamber, Jian Pei, Data Mining - Concepts and Techniques 3/e),
高立圖書, 2014

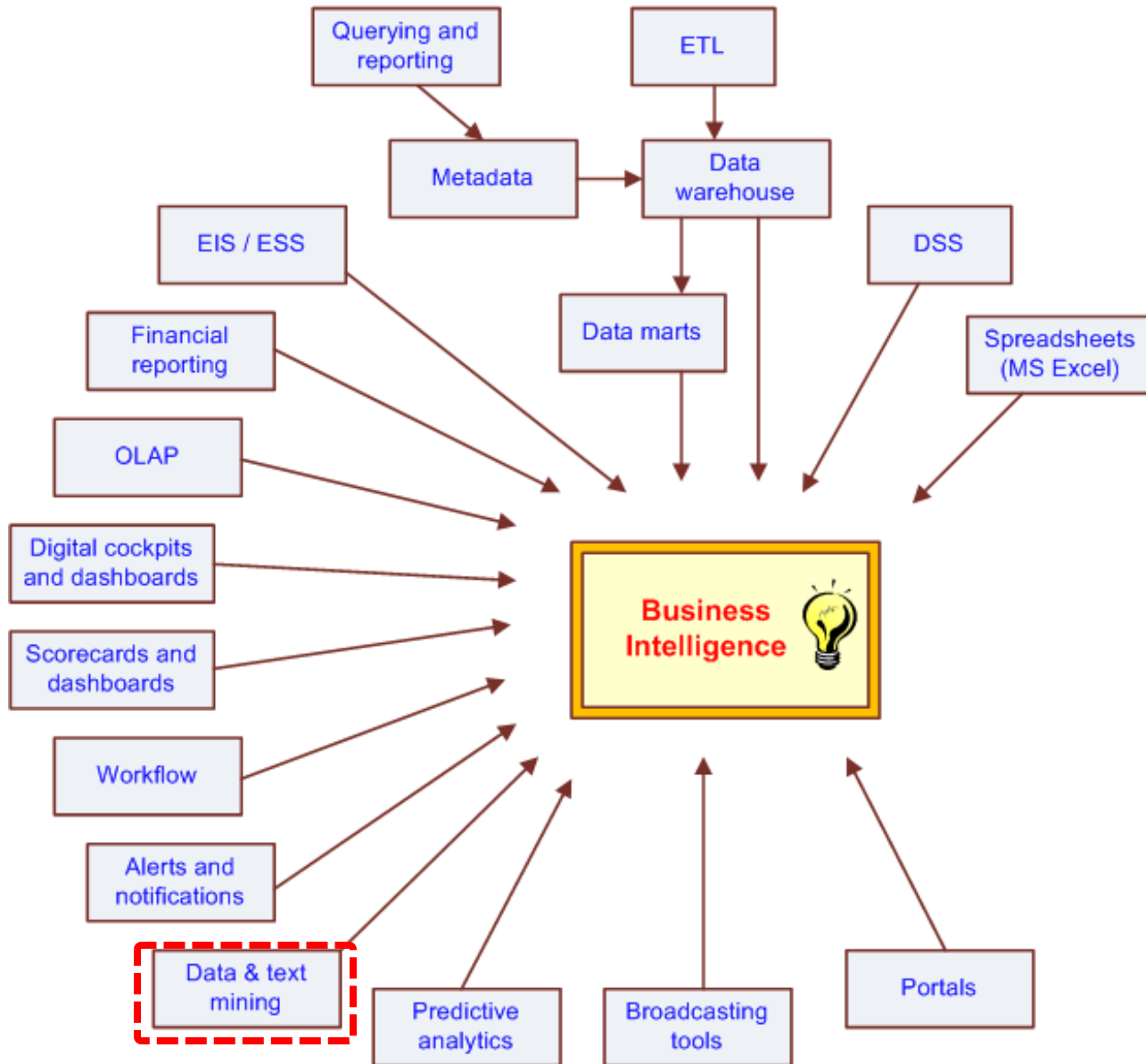


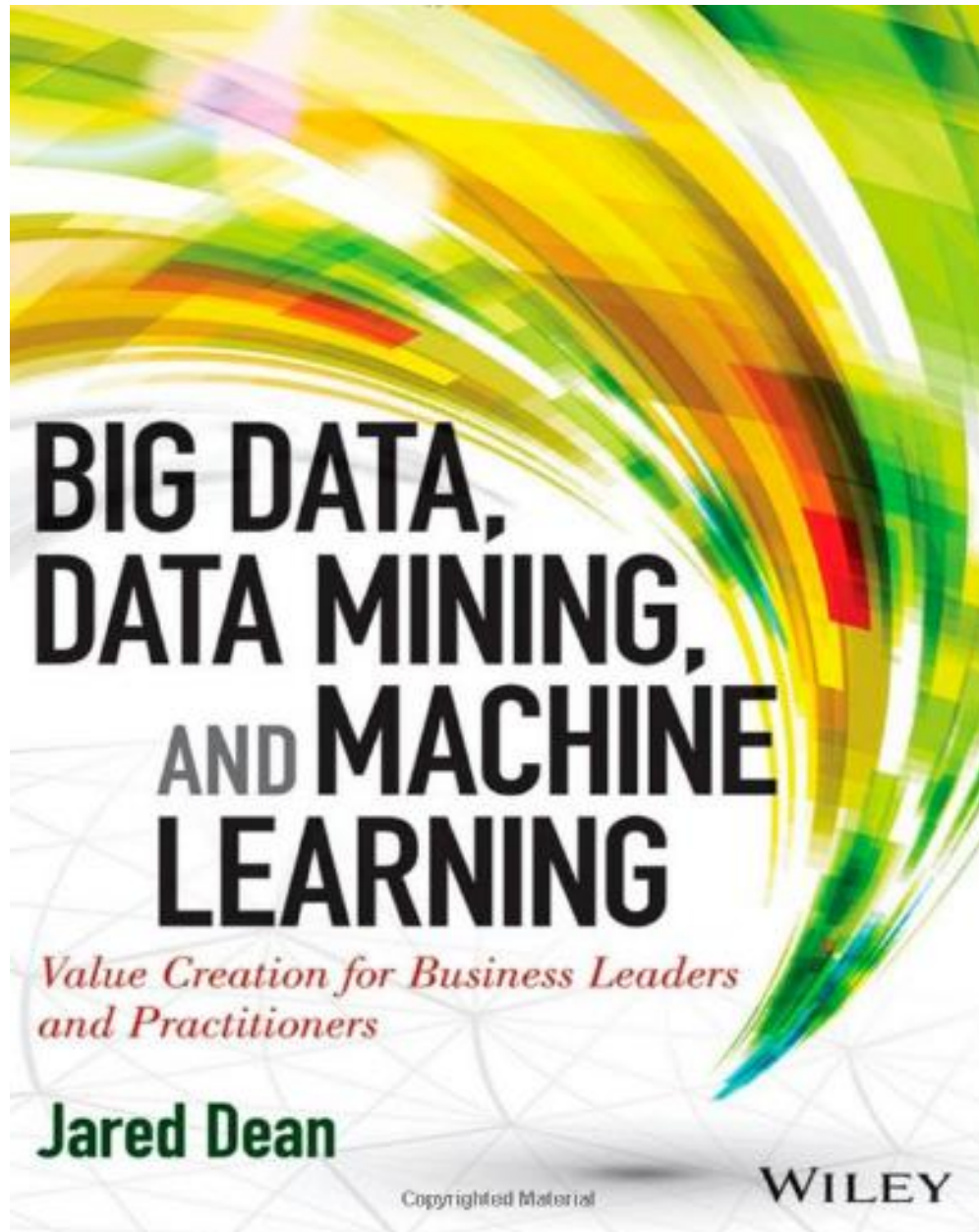
Data Warehouse

Data Mining and Business Intelligence



The Evolution of BI Capabilities

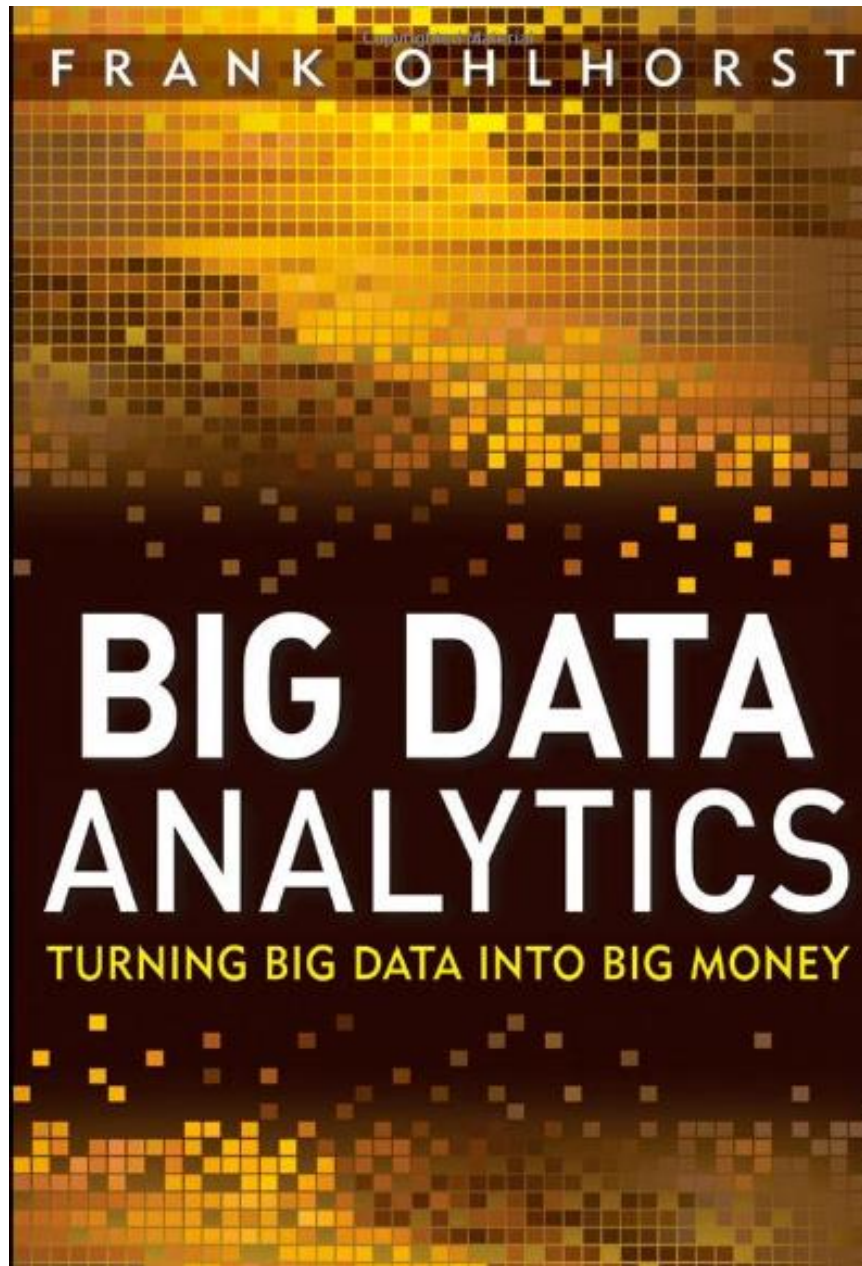




Deep Learning

Intelligence from Big Data







National Security

Cyber security

Maritime security

Smarter Transport

...

VISUAL ANALYTICS

DYNAMIC & INTERACTIVE

Dashboard Graph
Map

ENHANCE

Understanding Investigation
User Experience



BIG ANALYTICS

QUERY & FILTER

Complex queries
R²I²

DETECT

Anomalies
Communities
Typologies

PREDICT

Trending
Real-time
Prediction

DECIDE

Simulation
Optimization



BIG DATA – Batch



BIG DATA – Real Time



Complex by nature



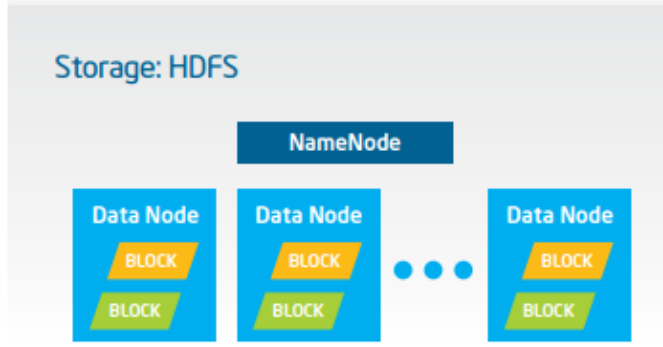
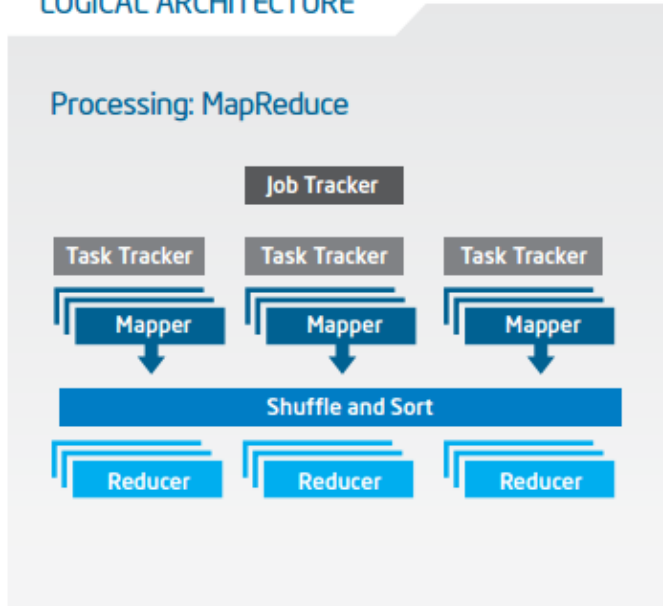
DATA

Complex by structure

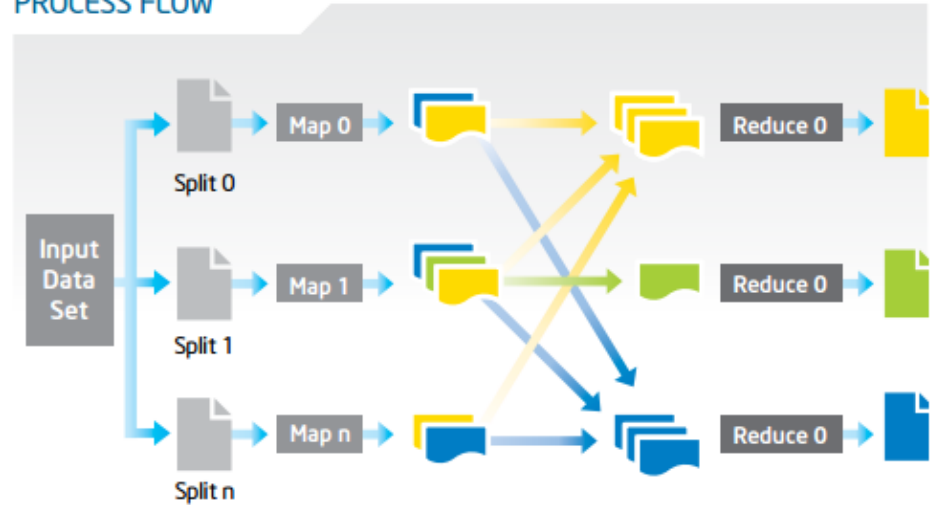


Big Data with Hadoop Architecture

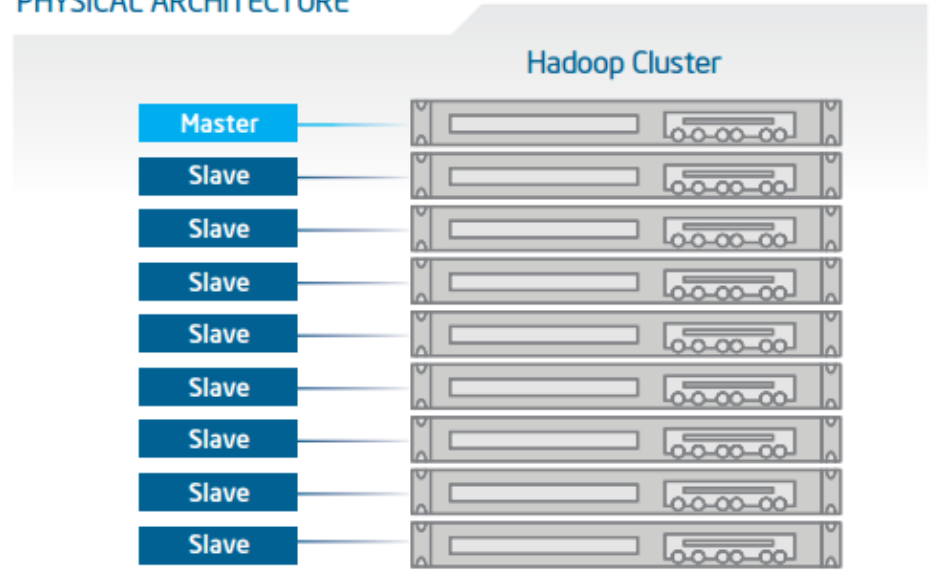
LOGICAL ARCHITECTURE



PROCESS FLOW



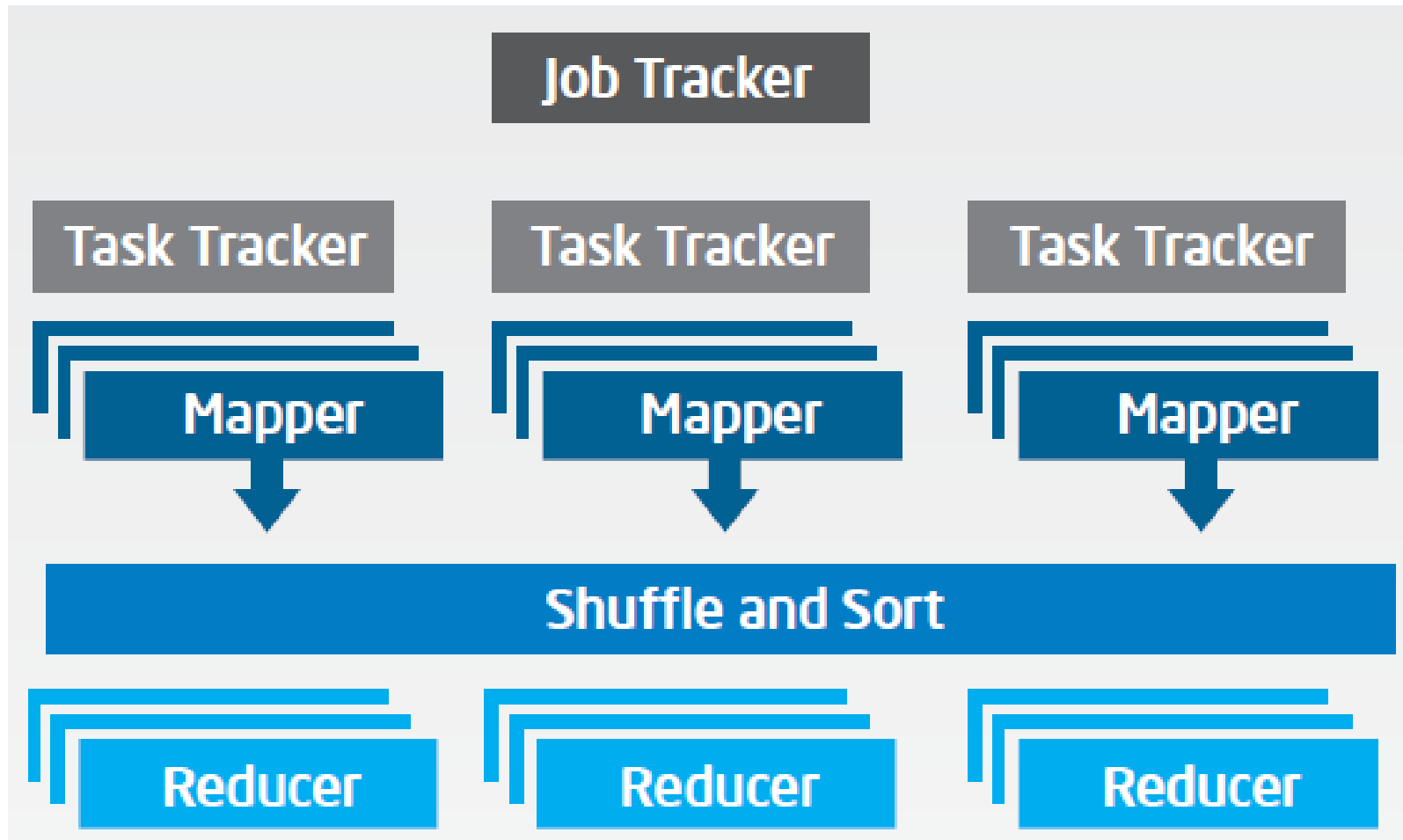
PHYSICAL ARCHITECTURE



Big Data with Hadoop Architecture

Logical Architecture

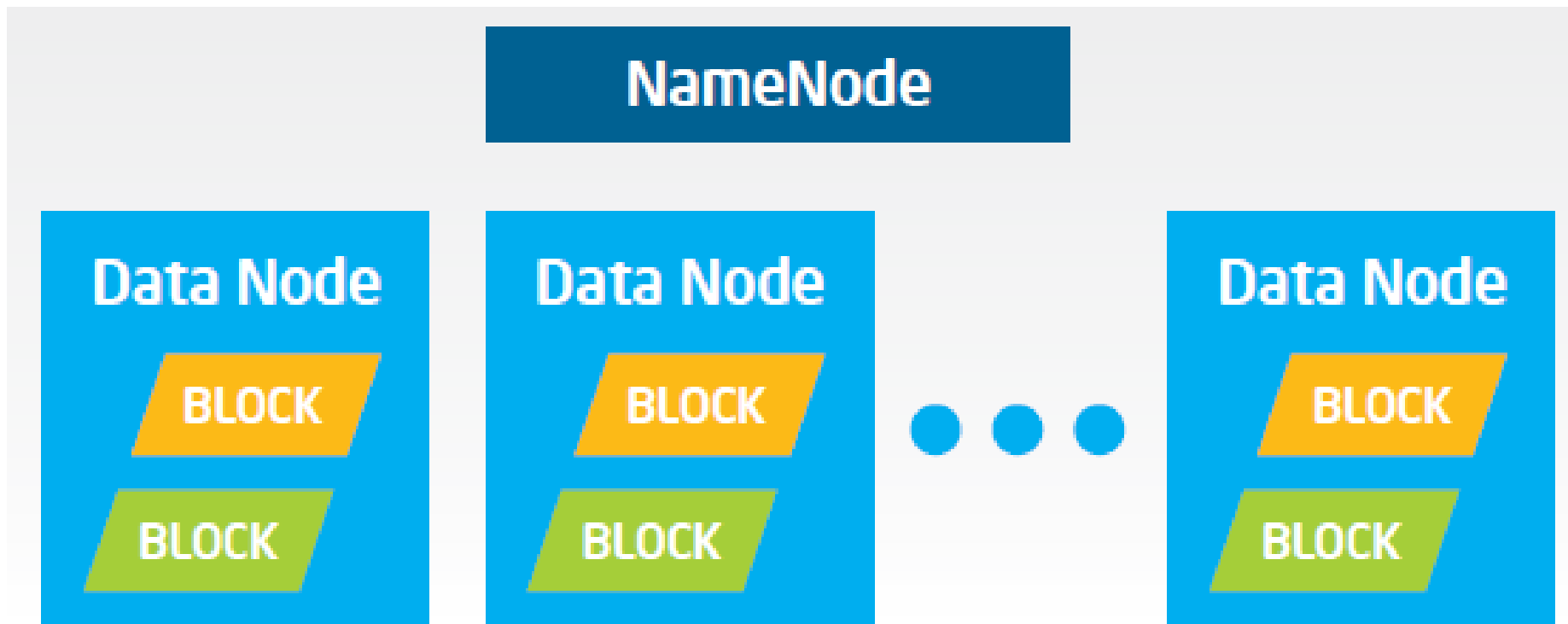
Processing: MapReduce



Big Data with Hadoop Architecture

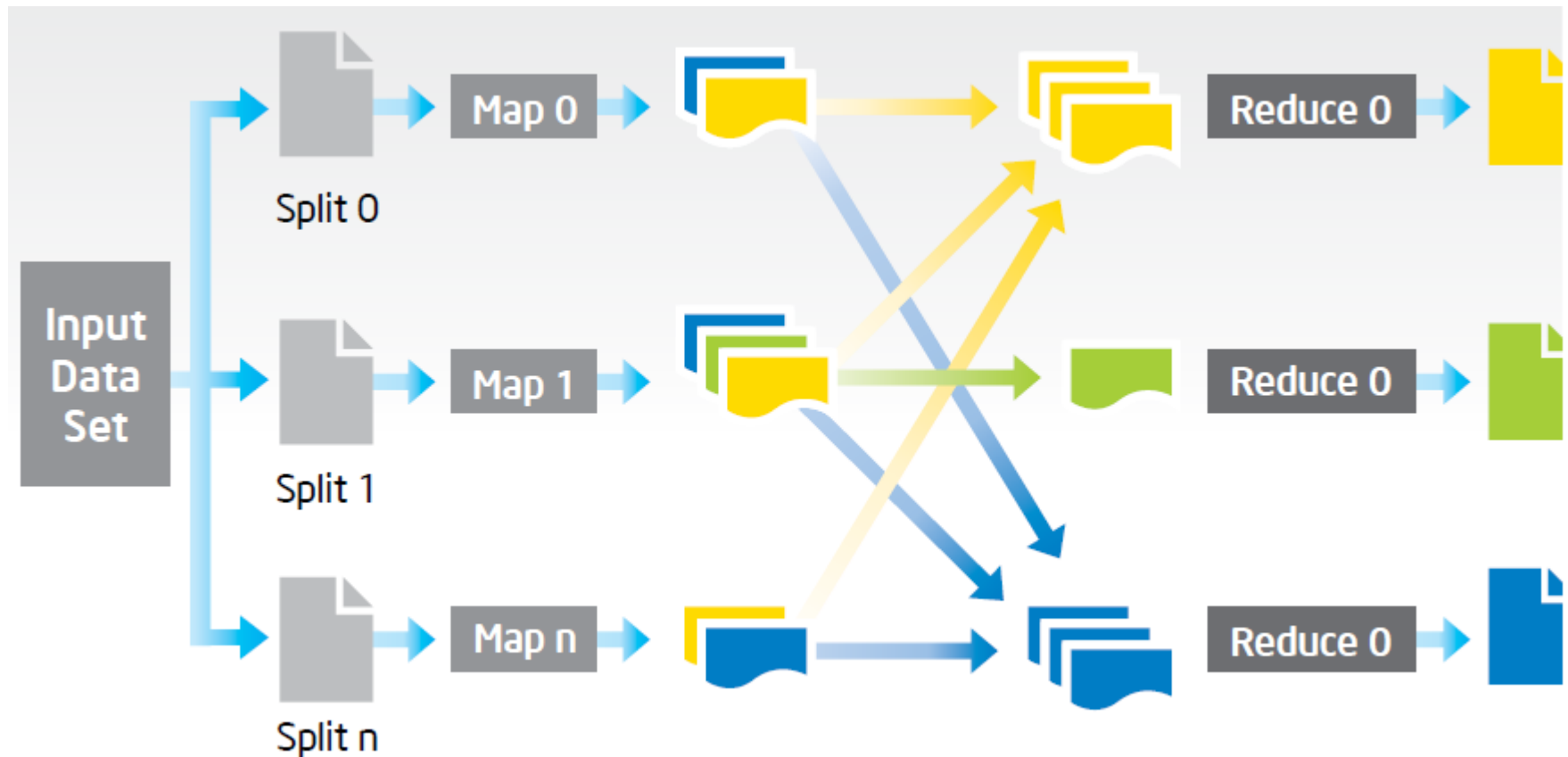
Logical Architecture

Storage: HDFS



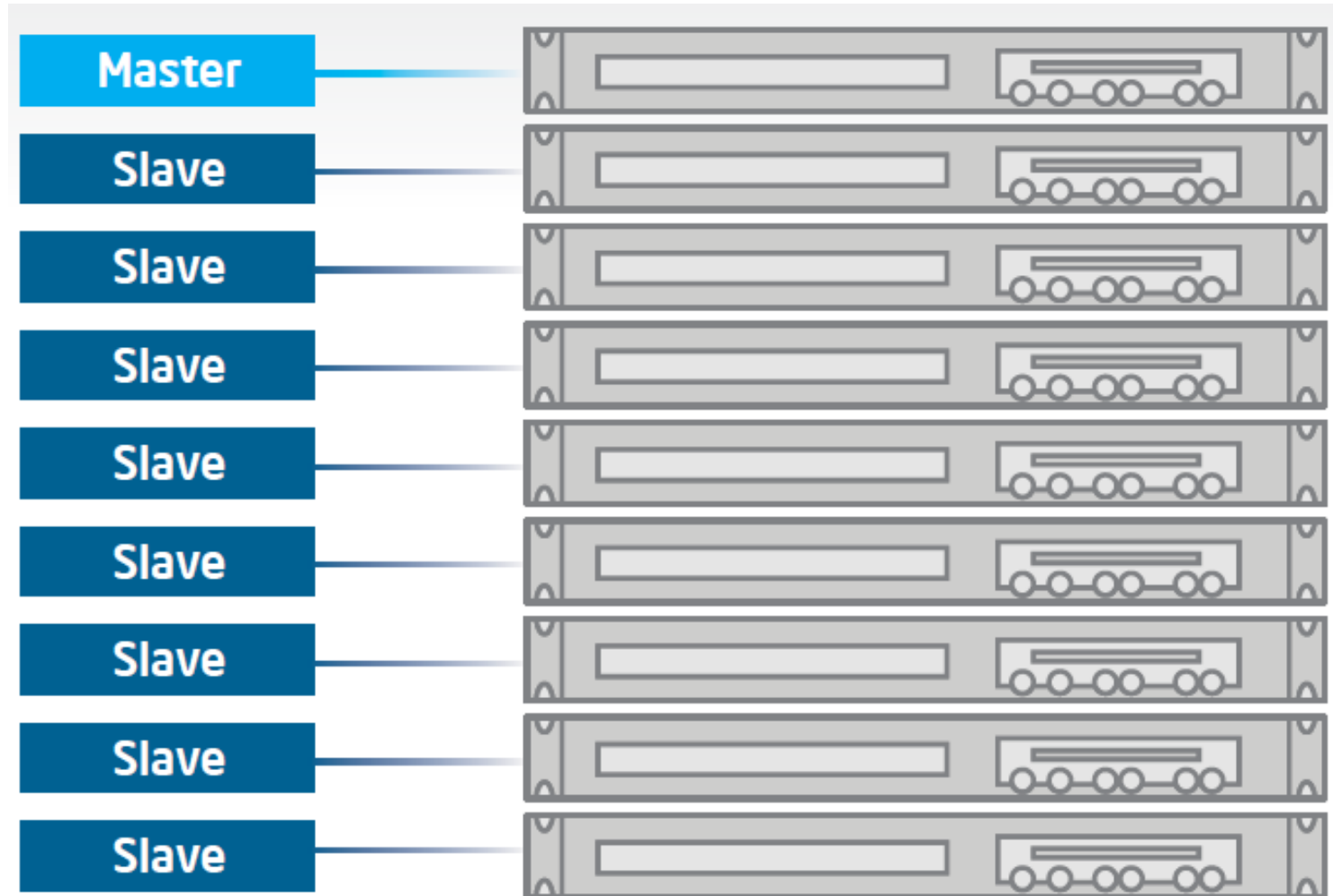
Big Data with Hadoop Architecture

Process Flow

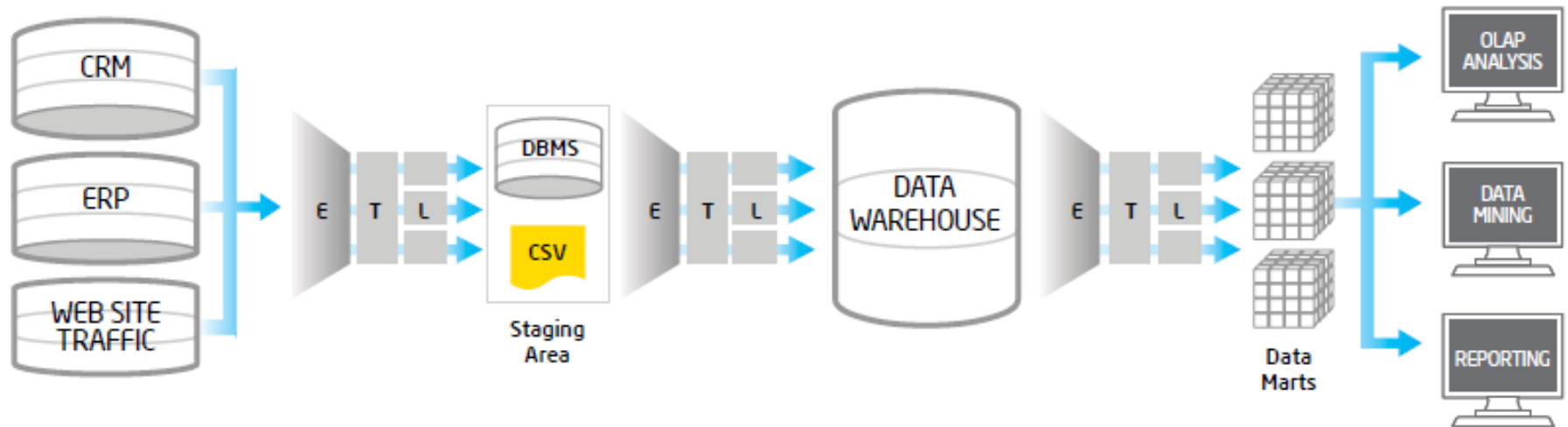


Big Data with Hadoop Architecture

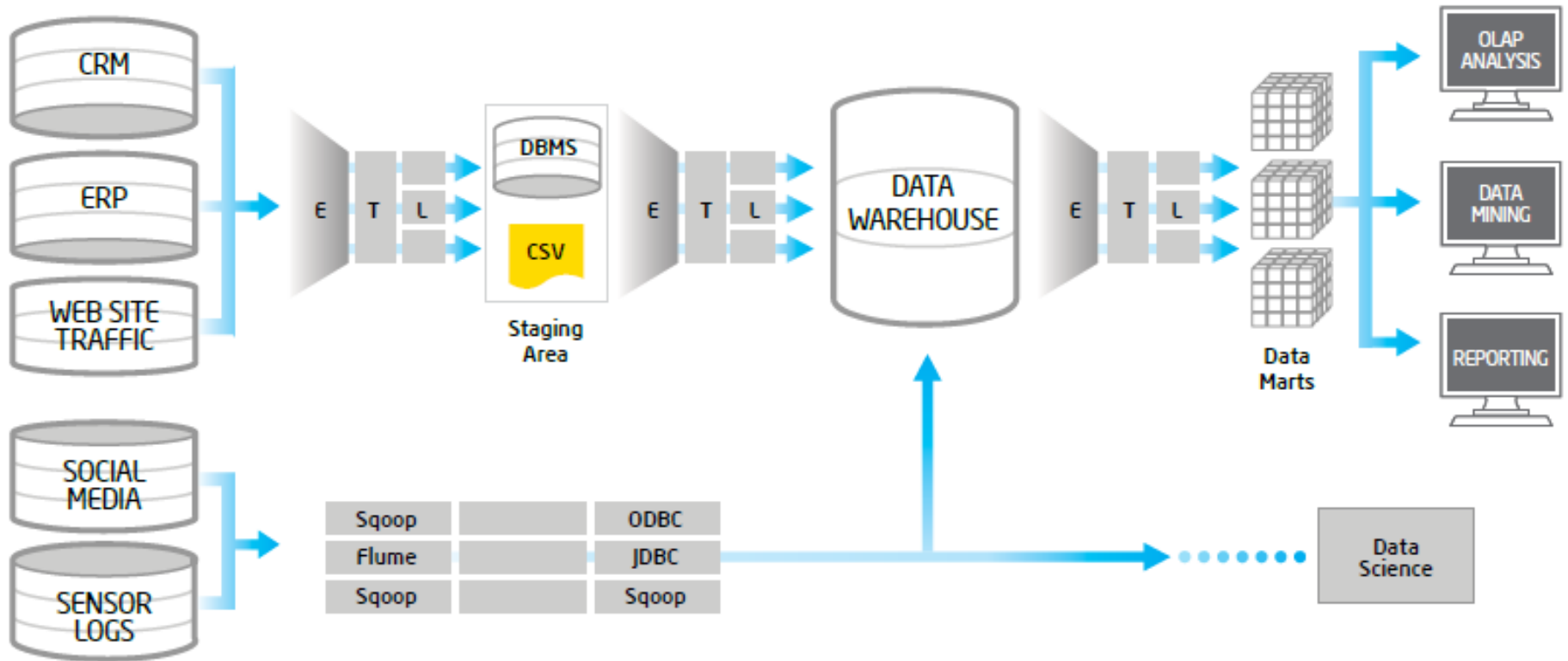
Hadoop Cluster



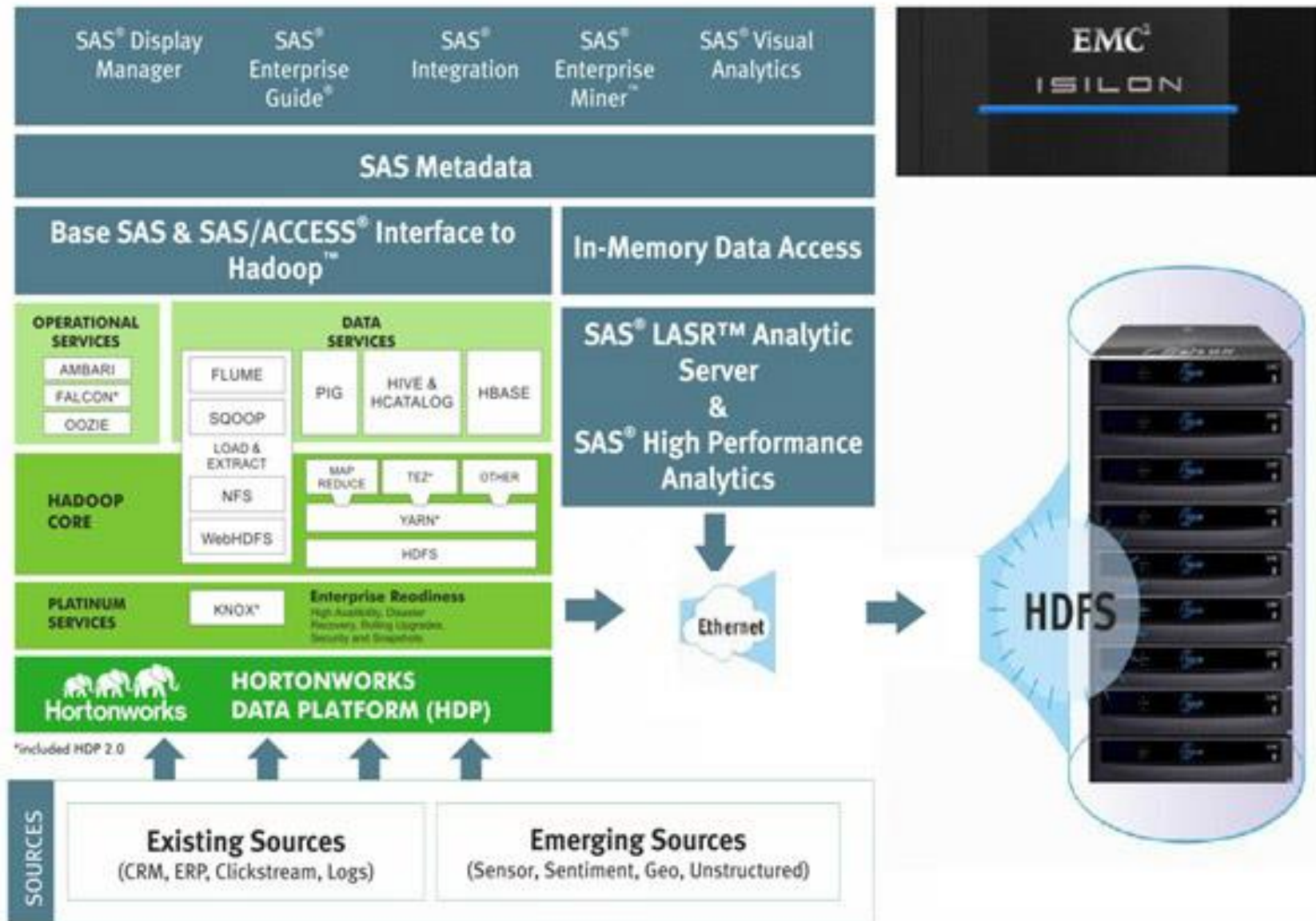
Traditional ETL Architecture



Offload ETL with Hadoop (Big Data Architecture)

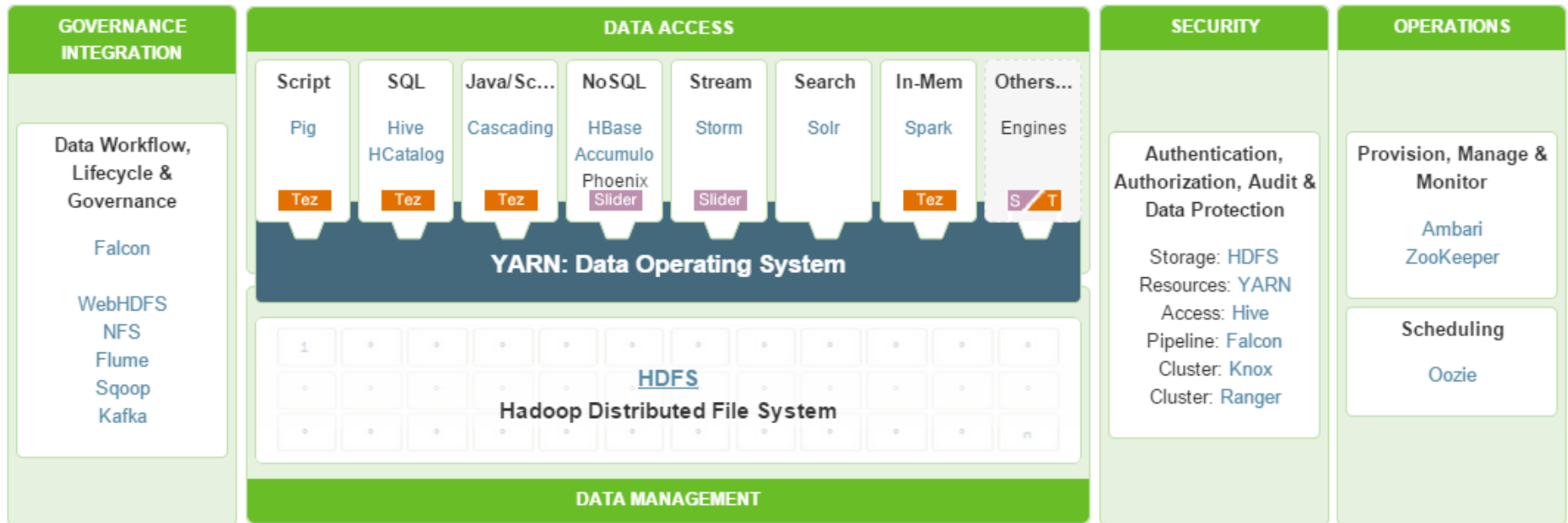


Big Data Solution













HDP

A Complete Enterprise Hadoop Data Platform

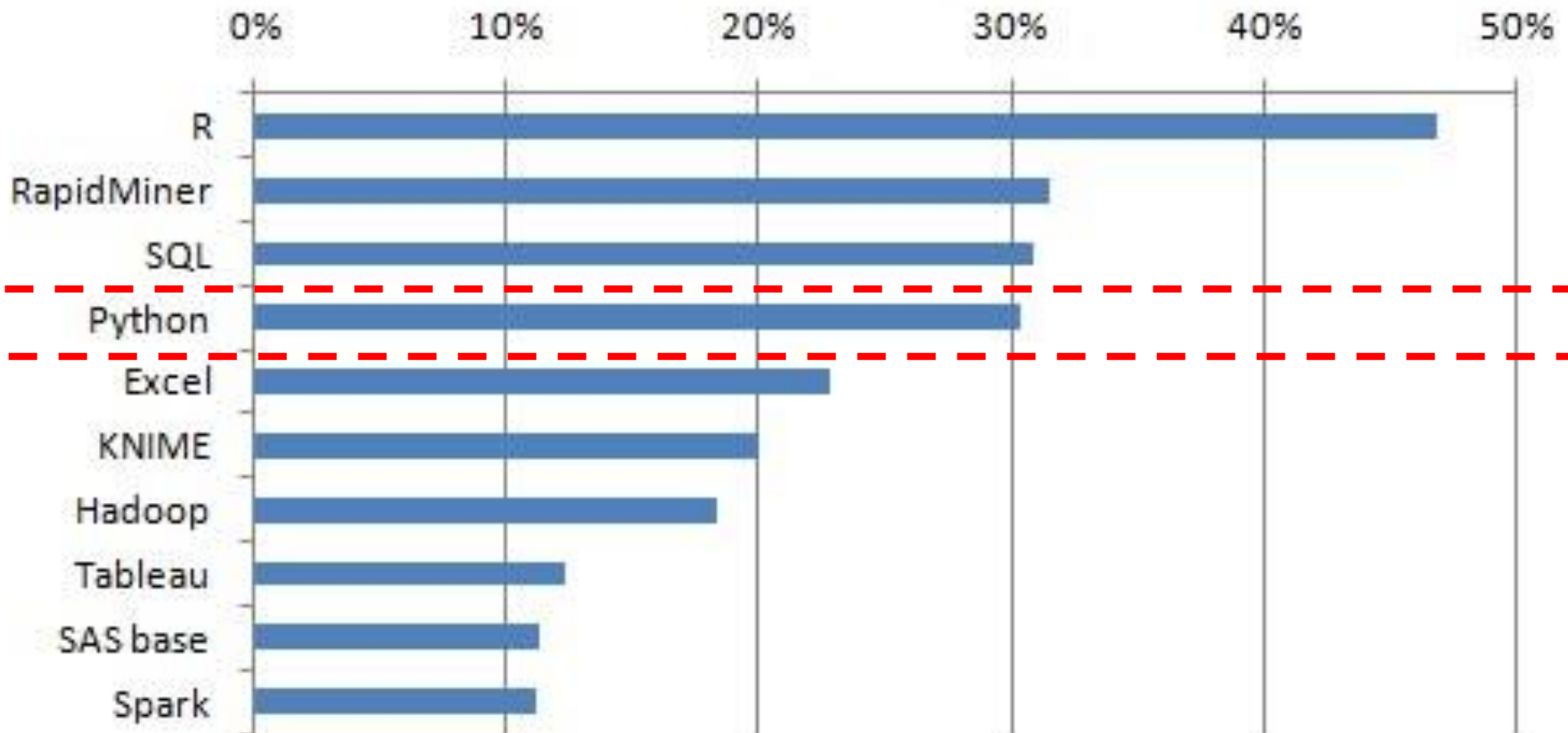


Python for Big Data Analytics

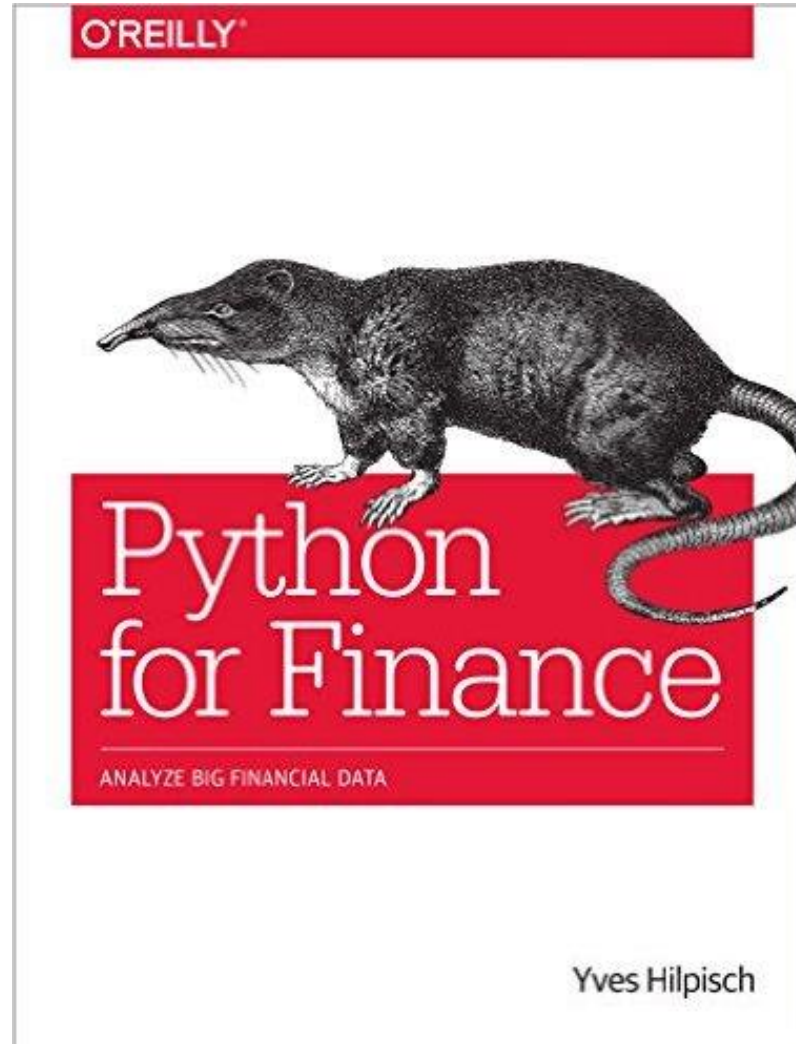
(The column on the left is the 2015 ranking; the column on the right is the 2014 ranking for comparison)

Language Rank	Types	2015 Spectrum Ranking	2014 Spectrum Ranking
1. Java		100.0	100.0
2. C		99.9	99.3
3. C++		99.4	95.5
4. Python		96.5	93.5
5. C#		91.3	92.4
6. R		84.8	84.8
7. PHP		84.5	84.5
8. JavaScript		83.0	78.9
9. Ruby		76.2	74.3
10. Matlab		72.4	72.8

Top Analytics, Data Mining, Data Science software used, 2015

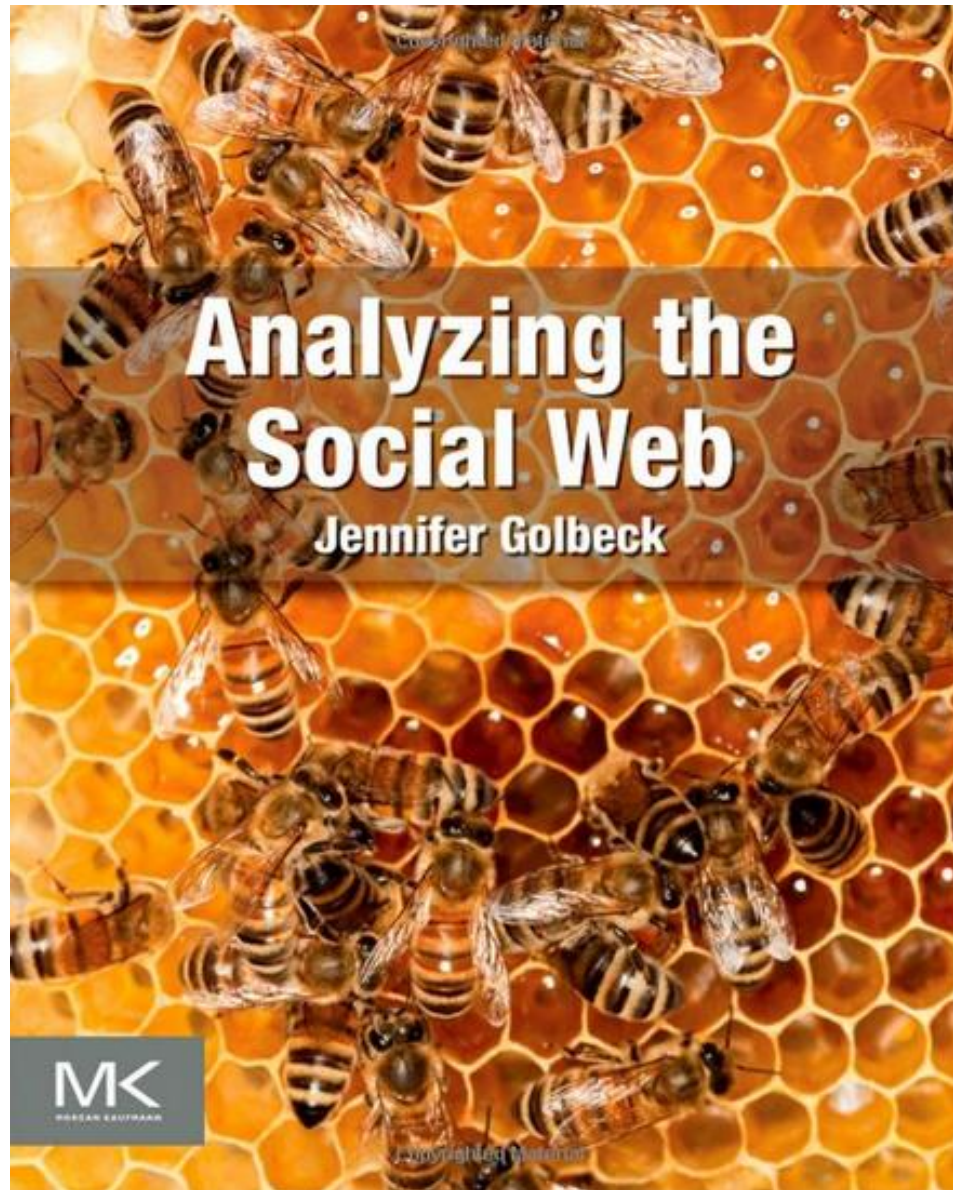


Yves Hilpisch, Python for Finance: Analyze Big Financial Data, O'Reilly, 2014



Analyzing the Social Web: Social Network Analysis

Jennifer Golbeck (2013), *Analyzing the Social Web*, Morgan Kaufmann



Source: <http://www.amazon.com/Analyzing-Social-Web-Jennifer-Golbeck/dp/0124055311>

Social Network Analysis (SNA)

Facebook TouchGraph

TouchGraph Photos x

box.touchgraph.com/facebook/TGFacebookBrowser.php?&signed_request=Gi-L3_6HrZ0S3SjxAXGdHR0rhMzqBjUnvFJ9vE4W6vg.eyJhbGdvcm0aG0iOiJITUFDI☆

Profiles Networks

Show Top 100 Friends Show All Friends Upload Advanced Restart

Zoom: Spacing:

Min-Yuh Day
 Networks: None
 Mutual Friends: 681

Facebook Profile

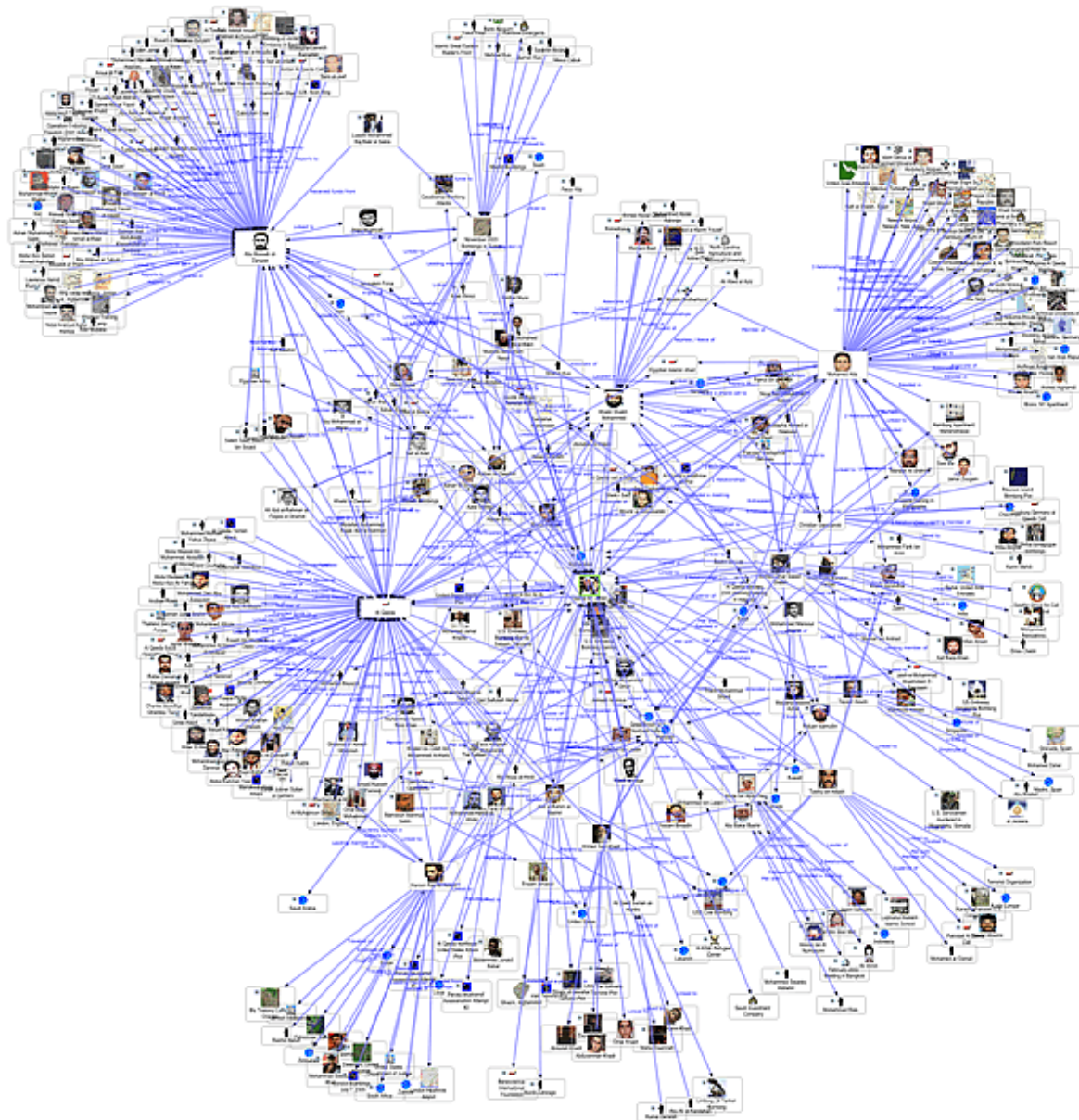
Network All All List Photo

Name	Rank #	Friend #
Min-Yuh Day	1	681
Gladys Hsieh	2	85
黃西田	3	74
施盛賓	4	67
John Lee	5	104
Kevin Tu	6	61
Yung Yu Shih	7	45
Wei Chen	8	107
Chichang Jou	9	50
Allen Green	10	81
黃煒勳	11	65
梁德昭	12	44
Eric Chen	13	51
吳錦波	14	39
Jessica Tien	15	49
蔡名宜	16	112
Enrico Lu	17	59
YaHan Hsieh	18	64
王慧雯	19	56
薛聖譚	20	80
蝦米	21	73

ICCU

powered by TouchGraph

Social Network Analysis



Social Network Analysis

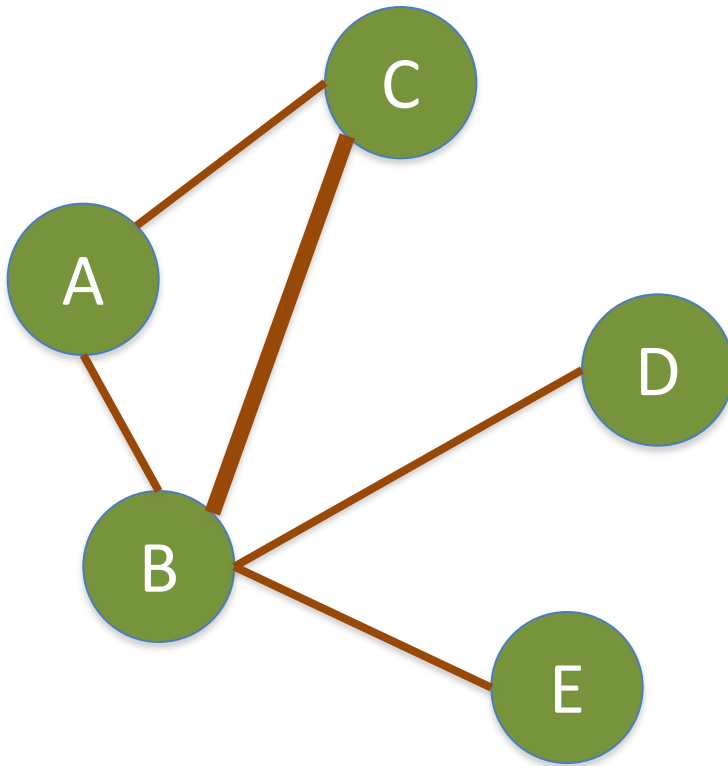
- A **social network** is a social structure of people, related (directly or indirectly) to each other through a common relation or interest
- **Social network analysis (SNA)** is the study of social networks to understand their structure and behavior

Social Network Analysis (SNA)

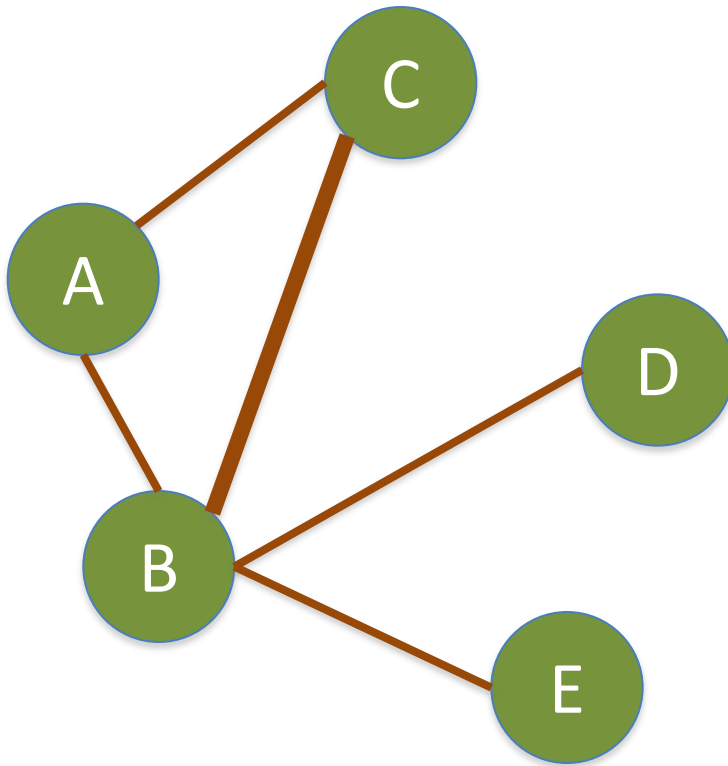
Centrality

Prestige

Degree



Degree



A: 2

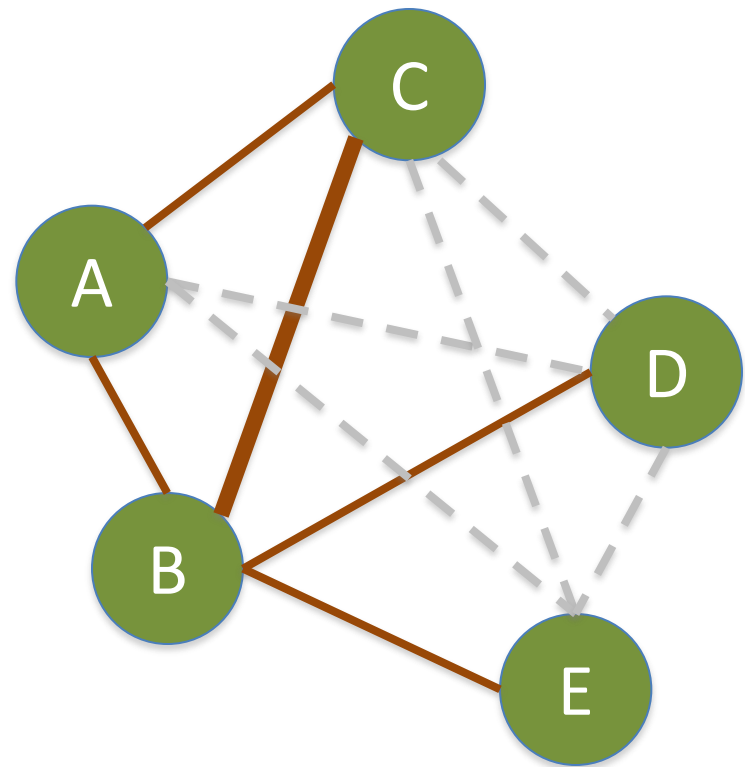
B: 4

C: 2

D: 1

E: 1

Density

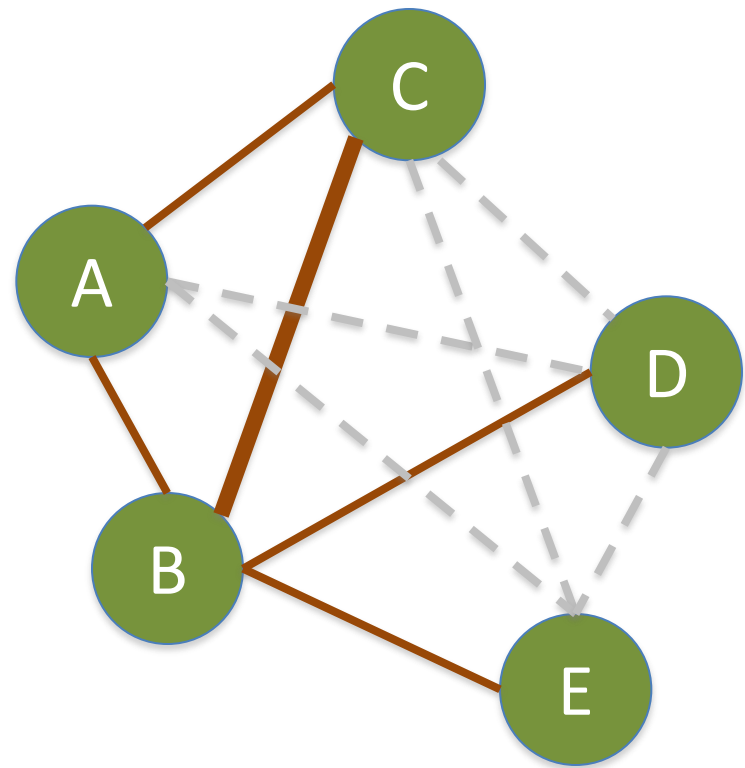


Density

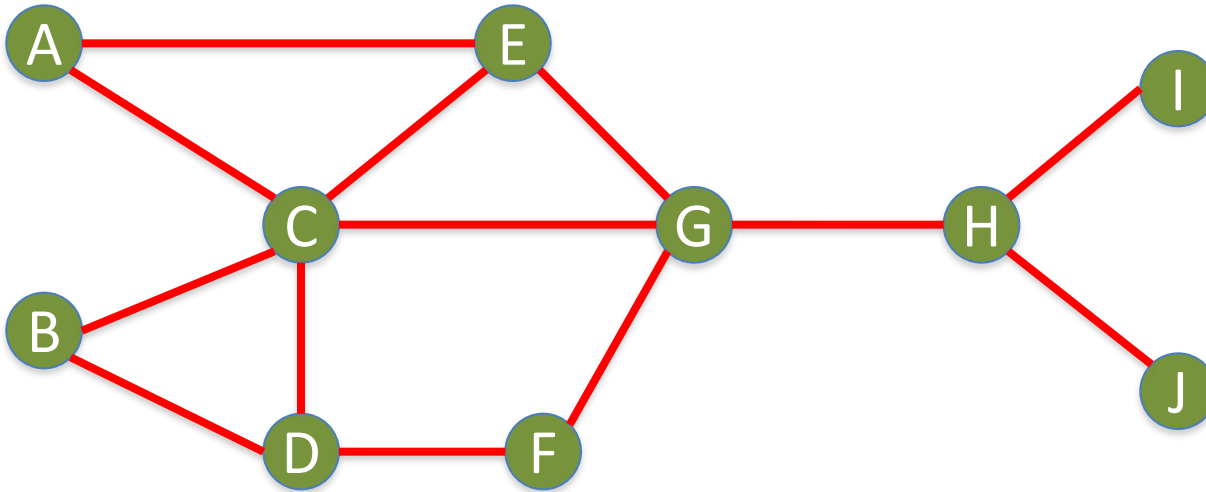
Edges (Links): 5

Total Possible Edges: 10

Density: $5/10 = 0.5$



Density



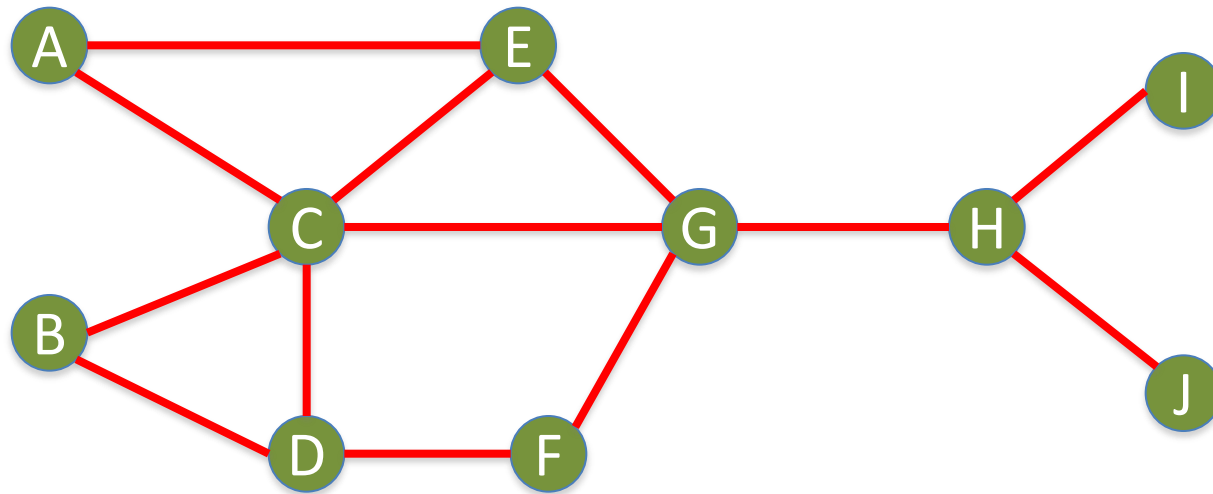
Nodes (n): 10

Edges (Links): 13

Total Possible Edges: $(n * (n-1)) / 2 = (10 * 9) / 2 = 45$

Density: $13/45 = 0.29$

Which Node is Most **Important**?



Centrality

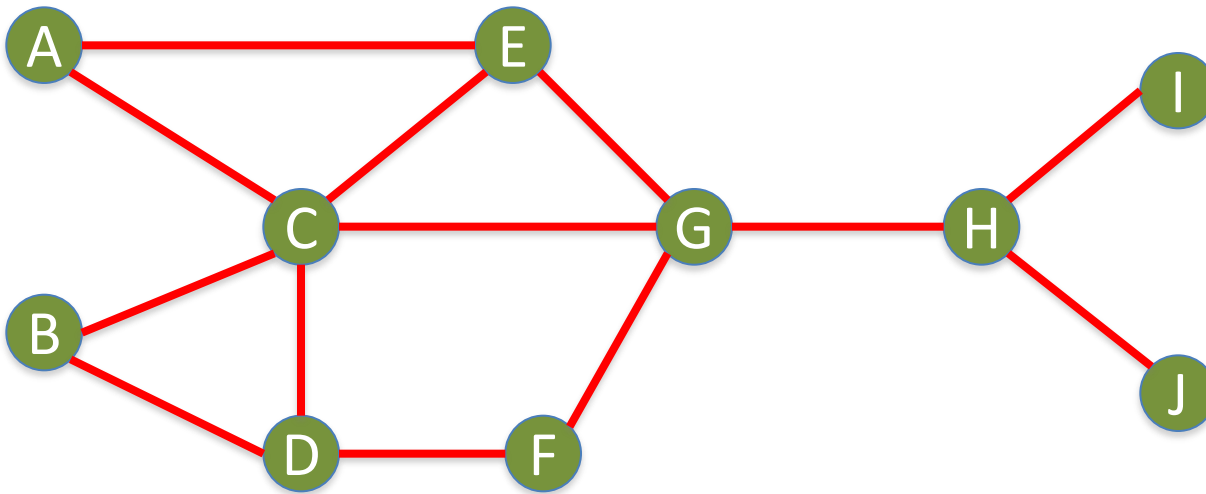
- **Important or prominent actors** are those that are linked or involved with other actors extensively.
- A person with extensive contacts (links) or communications with many other people in the organization is considered more important than a person with relatively fewer contacts.
- The links can also be called **ties**.
A **central actor** is one involved in many ties.

Social Network Analysis (SNA)

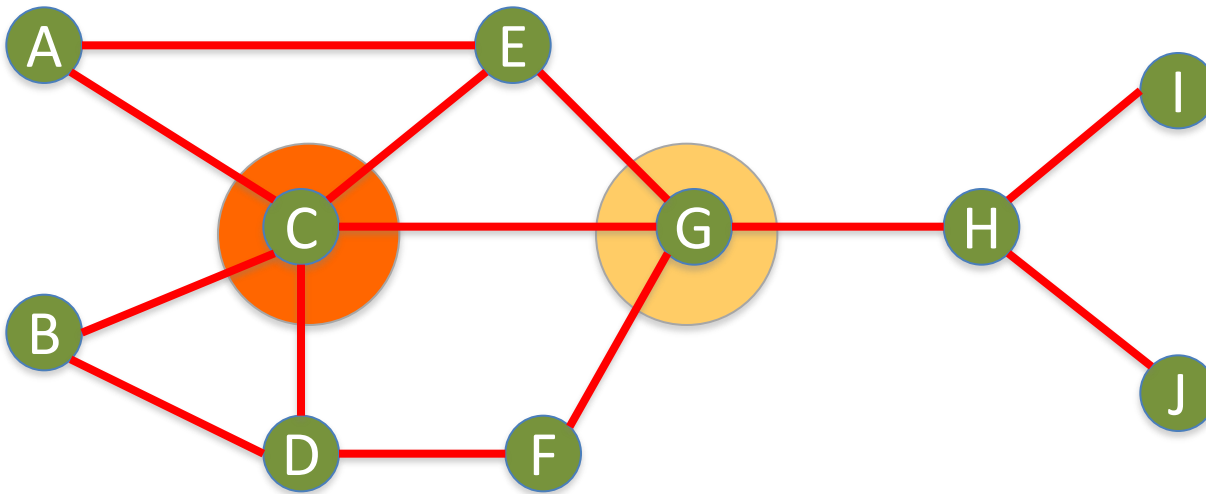
- Degree Centrality
- Betweenness Centrality
- Closeness Centrality

Degree Centrality

Social Network Analysis: Degree Centrality



Social Network Analysis: Degree Centrality



Node	Score	Standardized Score
A	2	$2/10 = 0.2$
B	2	$2/10 = 0.2$
C	5	$5/10 = 0.5$
D	3	$3/10 = 0.3$
E	3	$3/10 = 0.3$
F	2	$2/10 = 0.2$
G	4	$4/10 = 0.4$
H	3	$3/10 = 0.3$
I	1	$1/10 = 0.1$
J	1	$1/10 = 0.1$

Betweenness Centrality

Betweenness centrality:

Connectivity

Number of shortest paths
going through the actor

Betweenness Centrality

$$C_B(i) = \sum_{j < k} g_{ik}(i) / g_{jk}$$

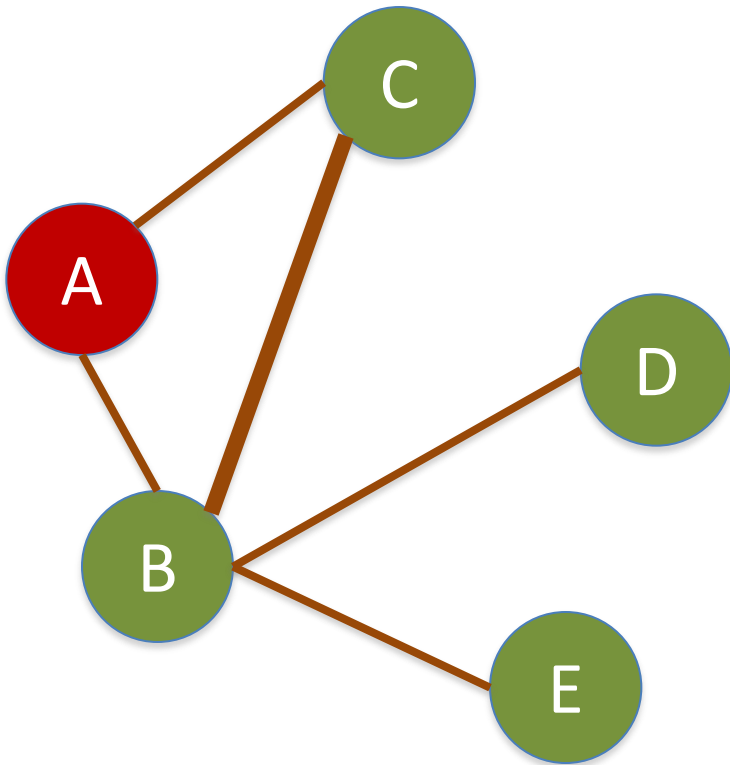
Where g_{jk} = the number of shortest paths connecting jk
 $g_{jk}(i)$ = the number that actor i is on.

Normalized Betweenness Centrality

$$C'_B(i) = C_B(i) / [(n-1)(n-2) / 2]$$

**Number of pairs of vertices
excluding the vertex itself**

Betweenness Centrality



A:

$$B \rightarrow C: 0/1 = 0$$

$$B \rightarrow D: 0/1 = 0$$

$$B \rightarrow E: 0/1 = 0$$

$$C \rightarrow D: 0/1 = 0$$

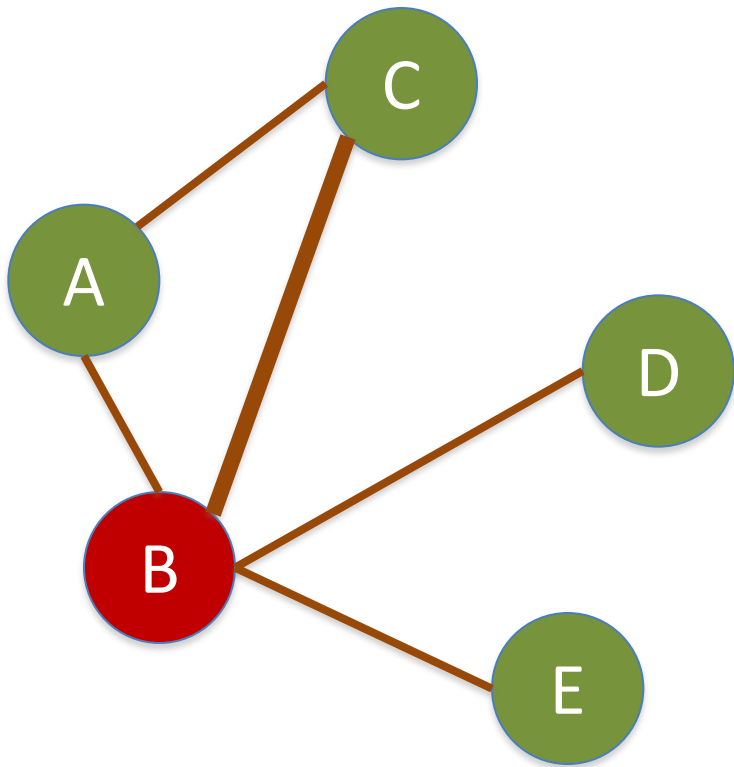
$$C \rightarrow E: 0/1 = 0$$

$$D \rightarrow E: 0/1 = 0$$

Total: 0

A: Betweenness Centrality = 0

Betweenness Centrality



B:

$$A \rightarrow C: 0/1 = 0$$

$$A \rightarrow D: 1/1 = 1$$

$$A \rightarrow E: 1/1 = 1$$

$$C \rightarrow D: 1/1 = 1$$

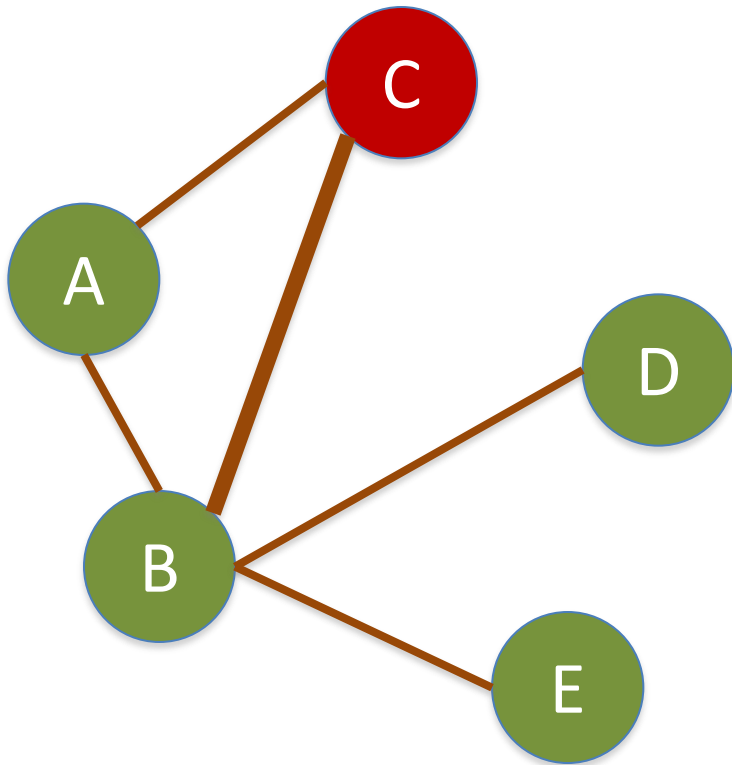
$$C \rightarrow E: 1/1 = 1$$

$$D \rightarrow E: 1/1 = 1$$

Total: 5

B: Betweenness Centrality = 5

Betweenness Centrality



C:

$$A \rightarrow B: 0/1 = 0$$

$$A \rightarrow D: 0/1 = 0$$

$$A \rightarrow E: 0/1 = 0$$

$$B \rightarrow D: 0/1 = 0$$

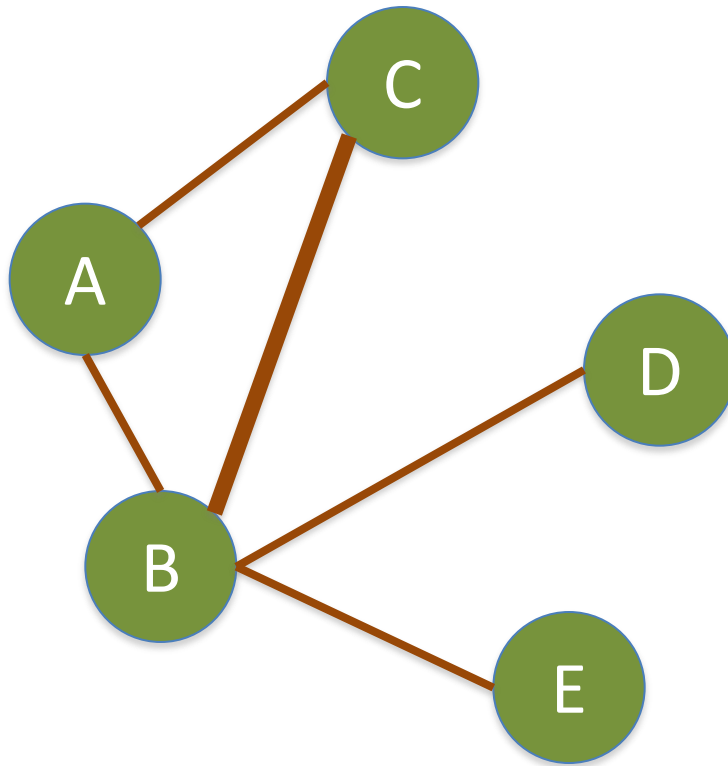
$$B \rightarrow E: 0/1 = 0$$

$$D \rightarrow E: 0/1 = 0$$

$$\text{Total: } \quad \quad \quad \underline{\quad 0 \quad}$$

C: Betweenness Centrality = 0

Betweenness Centrality



A: 0

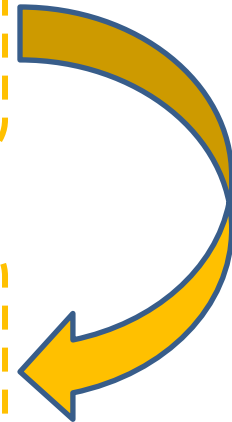
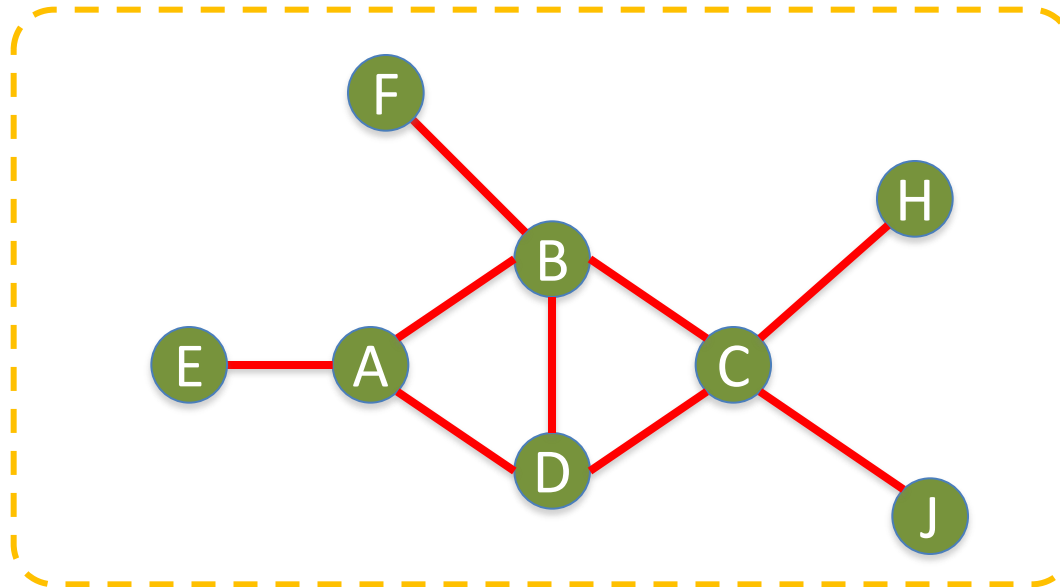
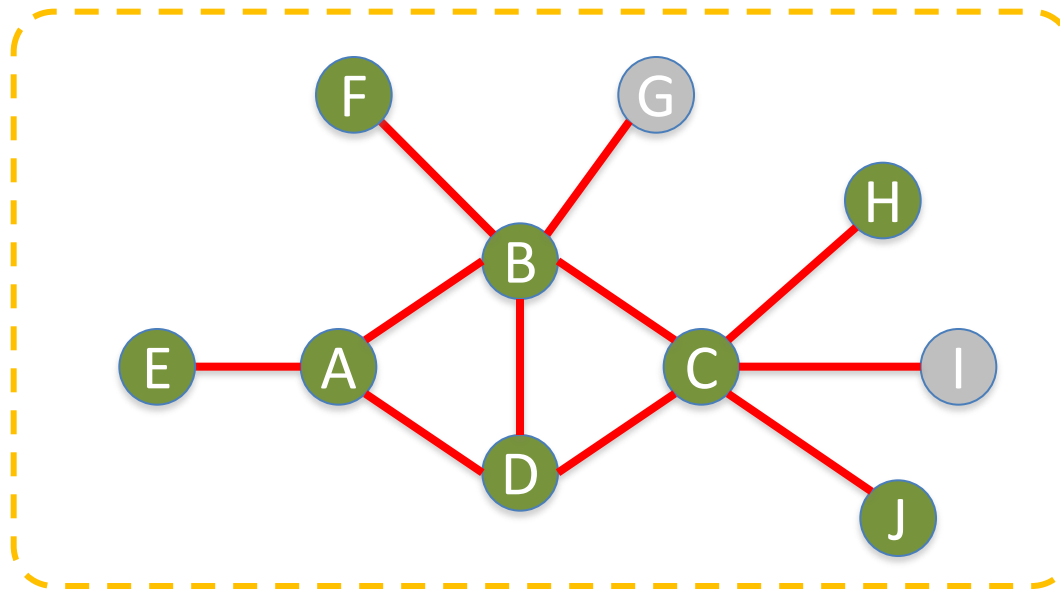
B: 5

C: 0

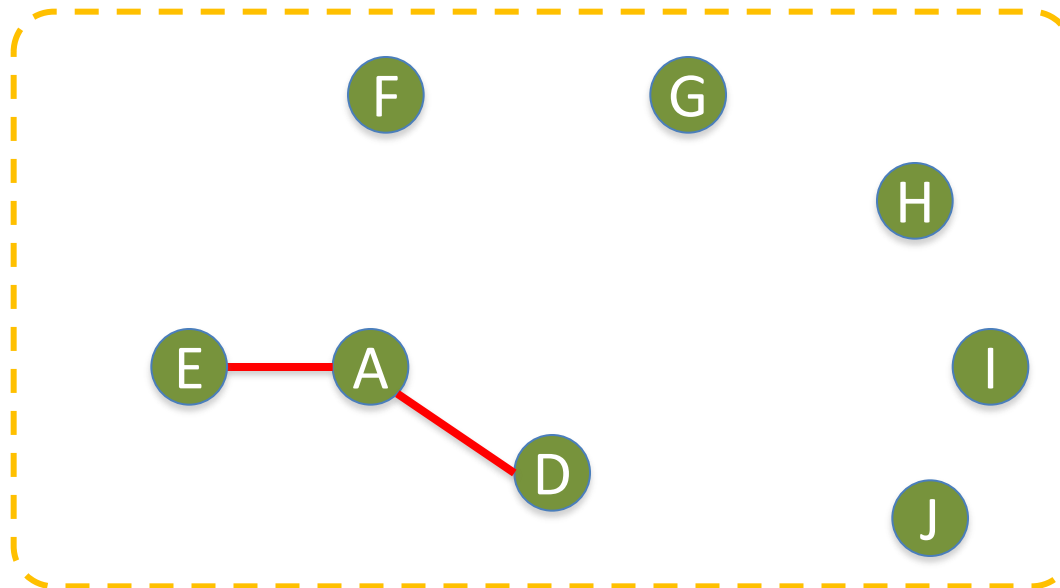
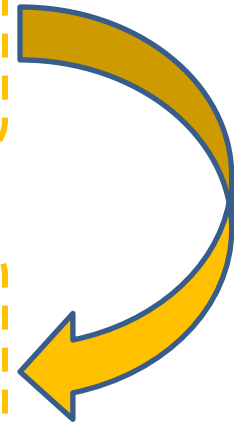
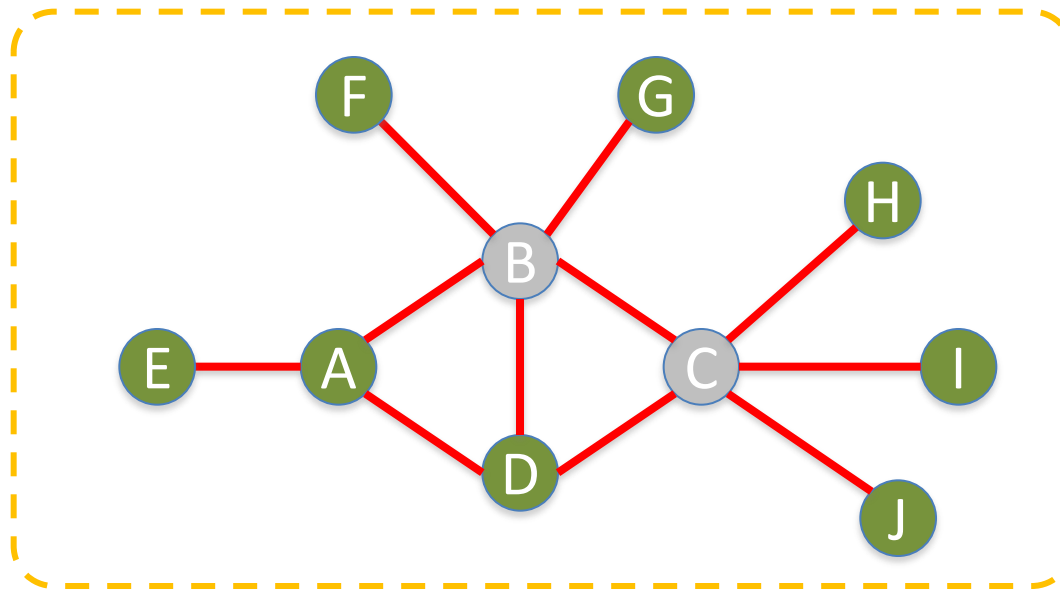
D: 0

E: 0

Which Node is Most **Important**?

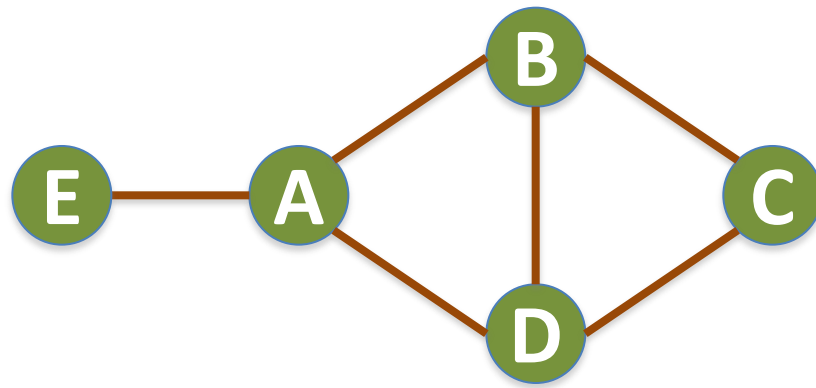


Which Node is Most **Important**?

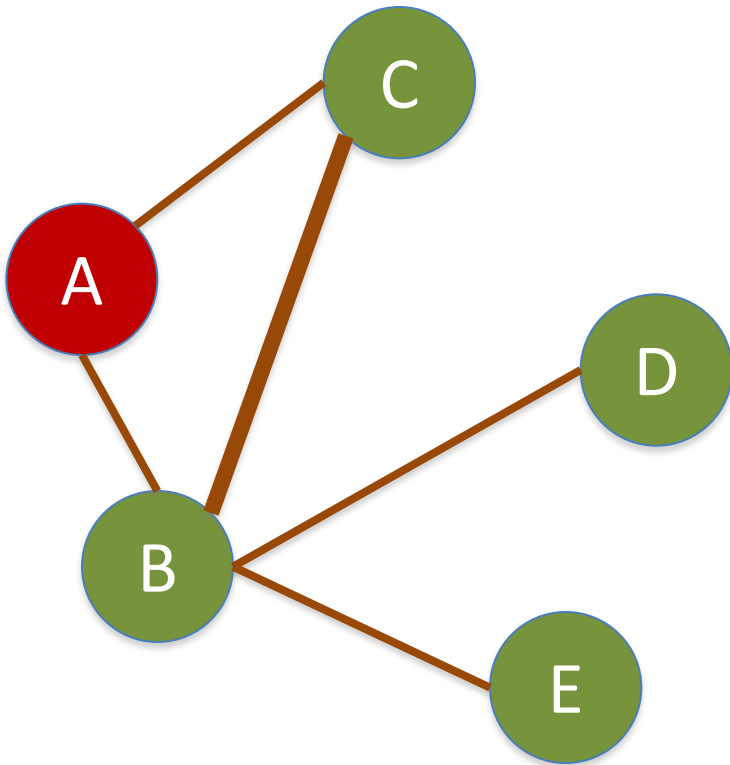


Betweenness Centrality

$$C_B(i) = \sum_{j < k} g_{ik}(i) / g_{jk}$$



Betweenness Centrality



A:

$$B \rightarrow C: 0/1 = 0$$

$$B \rightarrow D: 0/1 = 0$$

$$B \rightarrow E: 0/1 = 0$$

$$C \rightarrow D: 0/1 = 0$$

$$C \rightarrow E: 0/1 = 0$$

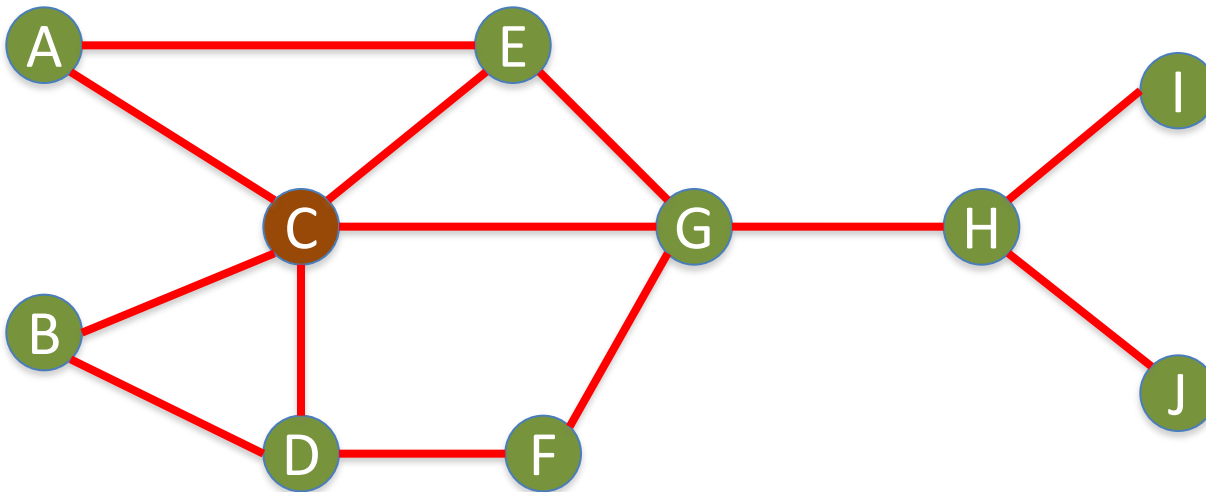
$$D \rightarrow E: 0/1 = 0$$

$$\text{Total: } \quad \quad \quad \underline{\quad 0 \quad}$$

A: Betweenness Centrality = 0

Closeness
Centrality

Social Network Analysis: Closeness Centrality

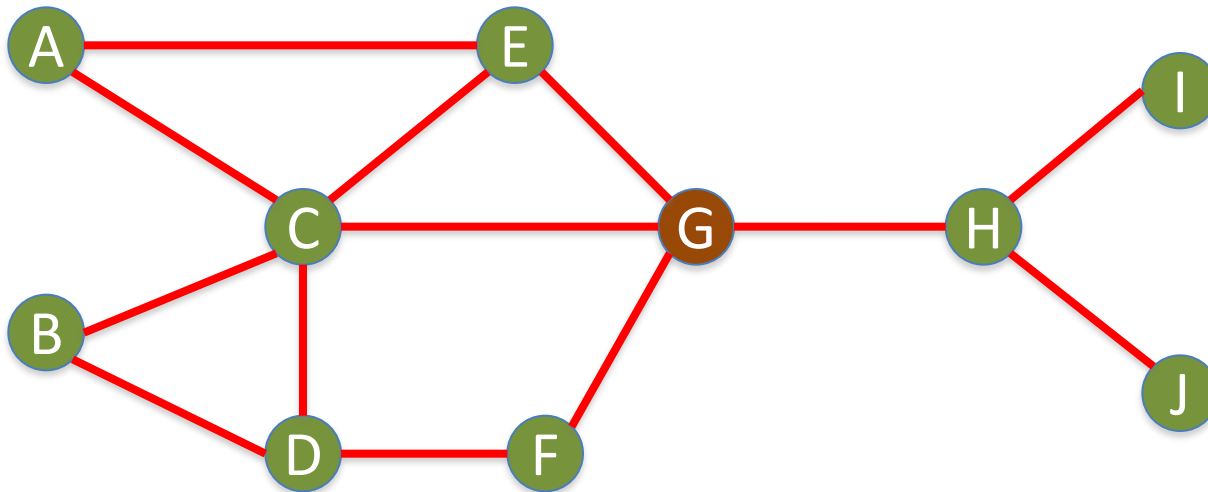


C→A: 1
C→B: 1
C→D: 1
C→E: 1
C→F: 2
C→G: 1
C→H: 2
C→I: 3
C→J: 3

Total=15

C: Closeness Centrality = $15/9 = 1.67$

Social Network Analysis: Closeness Centrality

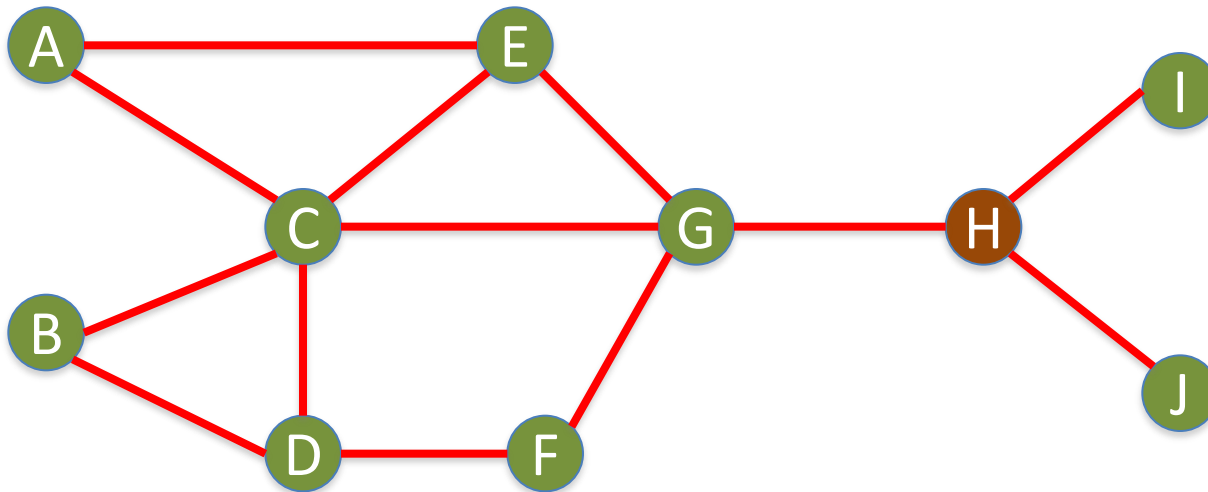


G→A: 2
G→B: 2
G→C: 1
G→D: 2
G→E: 1
G→F: 1
G→H: 1
G→I: 2
G→J: 2

Total=14

G: Closeness Centrality = $14/9 = 1.56$

Social Network Analysis: Closeness Centrality

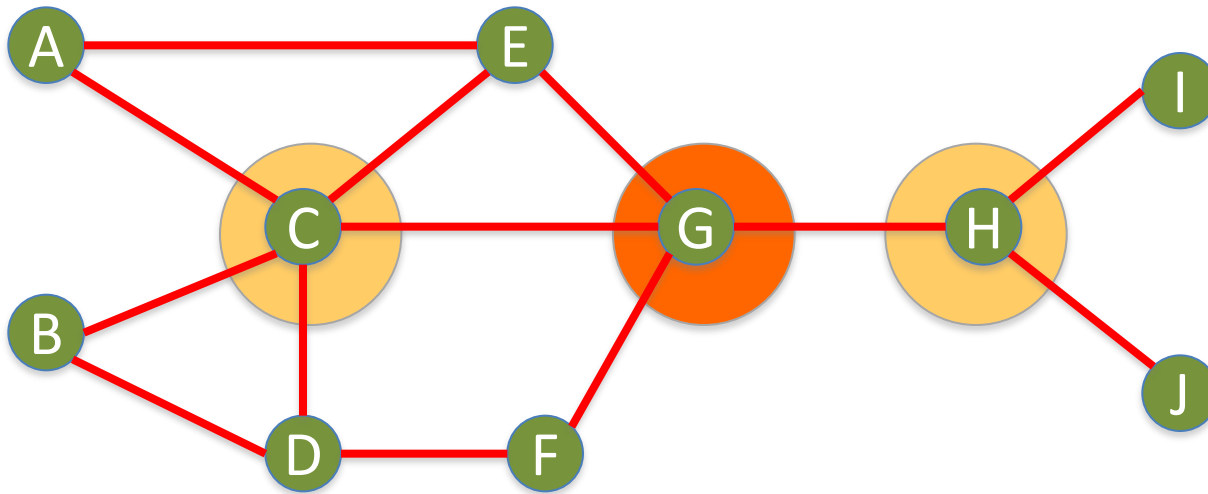


H→A: 3
H→B: 3
H→C: 2
H→D: 2
H→E: 2
H→F: 2
H→G: 1
H→I: 1
H→J: 1

Total=17

H: Closeness Centrality = $17/9 = 1.89$

Social Network Analysis: Closeness Centrality

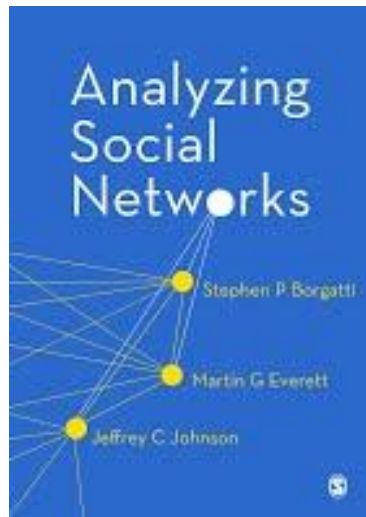


G: Closeness Centrality = $14/9 = 1.56$ ①

C: Closeness Centrality = $15/9 = 1.67$ ②

H: Closeness Centrality = $17/9 = 1.89$ ③

Social Network Analysis (SNA) Tools



- **UCINet**
- **Pajek**



Application of SNA

Social Network Analysis of Research Collaboration in Information Reuse and Integration

Example of SNA Data Source


















dblp
computer science bibliography













home | browse | search | about



IRI 2010: Las Vegas, NV, USA

-    **Proceedings of the IEEE International Conference on Information Reuse and Integration, IRI 2010, 4-6 August 2010, Las Vegas, Nevada, USA.** IEEE Systems, Man, and Cybernetics Society 2010
-    Reda Alhajj, James B. D. Joshi, Mei-Ling Shyu: **Message from Program Co-Chairs.** 1
-    Stuart Harvey Rubin, Shu-Ching Chen: **Forward.** 1
-    Lotfi A. Zadeh: **Precision of meaning - toward computation with natural language.** 1-4
-    Reda Alhajj, Shu-Ching Chen, Gongzhu Hu, James B. D. Joshi, Gordon K. Lee, Stuart Harvey Rubin, Mei-Ling Shyu, Lotfi A. Zadeh: **Panel title: Critical need for funding of basic and applied research in large-scale computing.** 1

Automation, Integration and Reuse across Various Apps

-    László István Etesi, André Csillaghy, Lin-Ching Chang: **A message-based interoperability framework with application to astrophysics.** 1-6
-    Awny Alnusair, Tian Zhao, Eric Bodden: **Effective API navigation and reuse.** 7-12
-    Manabu Ohta, Ryohei Inoue, Atsuhiko Takasu: **Empirical evaluation of active sampling for CRF-based analysis of pages.** 13-18
-    Qunzhi Zhou, Viktor K. Prasanna: **Workflow management of simulation based computation processes in transportation domain.** 19-24

Source: <http://www.informatik.uni-trier.de/~ley/db/conf/iri/iri2010.html>

Research Question

- RQ1: What are the scientific **collaboration patterns** in the IRI research community?
- RQ2: Who are the **prominent researchers** in the IRI community?

Methodology

- Developed a simple **web focused crawler** program to download literature information about all IRI papers published between **2003 and 2010** from **IEEE Xplore** and **DBLP**.
 - **767** paper
 - **1599** distinct author
- Developed a program to convert the list of coauthors into the **format of a network file** which can be readable by social network analysis software.
- **UCInet** and **Pajek** were used in this study for the social network analysis.

Top10 prolific authors (IRI 2003-2010)

1. Stuart Harvey Rubin
2. Taghi M. Khoshgoftaar
3. Shu-Ching Chen
4. Mei-Ling Shyu
5. Mohamed E. Fayad
6. Reda Alhajj
7. Du Zhang
8. Wen-Lian Hsu
9. Jason Van Hulse
10. Min-Yuh Day

Data Analysis and Discussion

- **Closeness Centrality**
 - Collaborated widely
- **Betweenness Centrality**
 - Collaborated diversely
- **Degree Centrality**
 - Collaborated frequently
- **Visualization of Social Network Analysis**
 - Insight into the structural characteristics of research collaboration networks

Top 20 authors with the highest **closeness** scores

Rank	ID	Closeness	Author
1	3	0.024675	Shu-Ching Chen
2	1	0.022830	Stuart Harvey Rubin
3	4	0.022207	Mei-Ling Shyu
4	6	0.020013	Reda Alhajj
5	61	0.019700	Na Zhao
6	260	0.018936	Min Chen
7	151	0.018230	Gordon K. Lee
8	19	0.017962	Chengcui Zhang
9	1043	0.017962	Isai Michel Lombera
10	1027	0.017962	Michael Armella
11	443	0.017448	James B. Law
12	157	0.017082	Keqi Zhang
13	253	0.016731	Shahid Hamid
14	1038	0.016618	Walter Z. Tang
15	959	0.016285	Chengjun Zhan
16	957	0.016285	Lin Luo
17	956	0.016285	Guo Chen
18	955	0.016285	Xin Huang
19	943	0.016285	Sneh Gulati
20	960	0.016071	Sheng-Tun Li

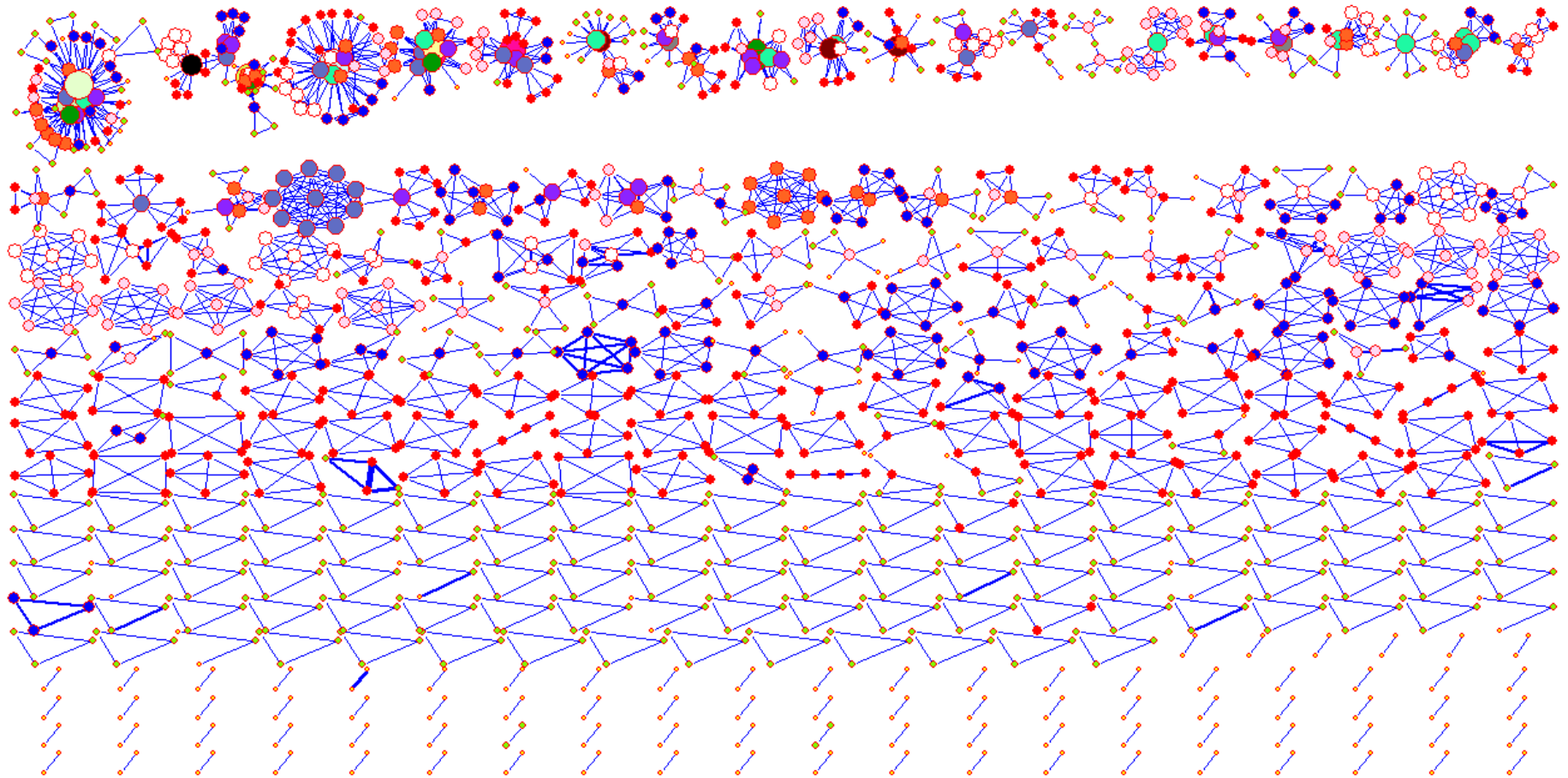
Top 20 authors with the highest **betweenness** scores

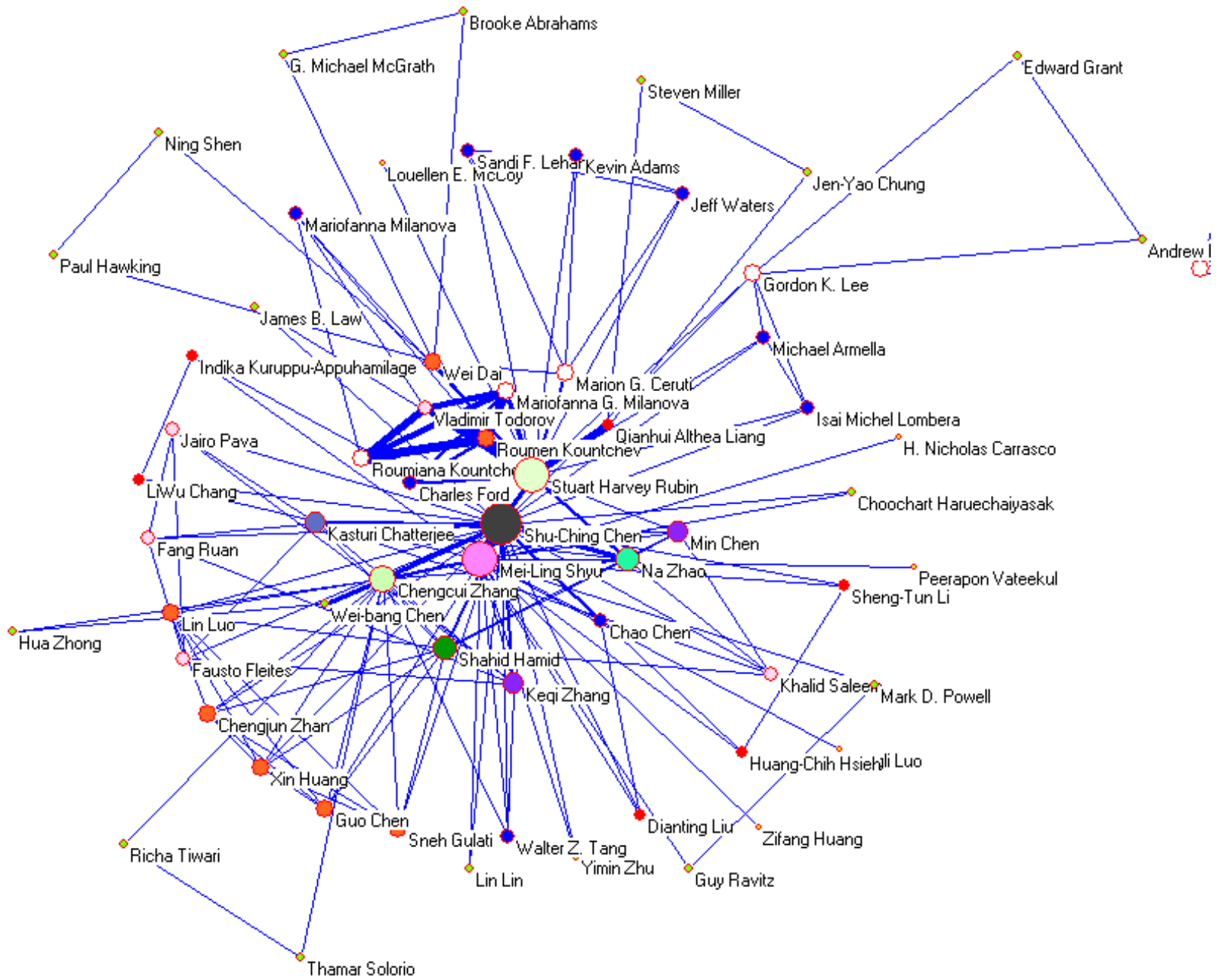
Rank	ID	Betweenness	Author
1	1	0.000752	Stuart Harvey Rubin
2	3	0.000741	Shu-Ching Chen
3	2	0.000406	Taghi M. Khoshgoftaar
4	66	0.000385	Xingquan Zhu
5	4	0.000376	Mei-Ling Shyu
6	6	0.000296	Reda Alhajj
7	65	0.000256	Xindong Wu
8	19	0.000194	Chengcui Zhang
9	39	0.000185	Wei Dai
10	15	0.000107	Narayan C. Debnath
11	31	0.000094	Qianhui Althea Liang
12	151	0.000094	Gordon K. Lee
13	7	0.000085	Du Zhang
14	30	0.000072	Baowen Xu
15	41	0.000067	Hongji Yang
16	270	0.000060	Zhiwei Xu
17	5	0.000043	Mohamed E. Fayad
18	110	0.000042	Abhijit S. Pandya
19	106	0.000042	Sam Hsu
20	8	0.000042	Wen-Lian Hsu

Top 20 authors with the highest degree scores

Rank	ID	Degree	Author
1	3	0.035044	Shu-Ching Chen
2	1	0.034418	Stuart Harvey Rubin
3	2	0.030663	Taghi M. Khoshgoftaar
4	6	0.028786	Reda Alhajj
5	8	0.028786	Wen-Lian Hsu
6	10	0.024406	Min-Yuh Day
7	4	0.022528	Mei-Ling Shyu
8	17	0.021277	Richard Tzong-Han Tsai
9	14	0.017522	Eduardo Santana de Almeida
10	16	0.017522	Roumen Kountchev
11	40	0.016896	Hong-Jie Dai
12	15	0.015645	Narayan C. Debnath
13	9	0.015019	Jason Van Hulse
14	25	0.013767	Roumiana Kountcheva
15	28	0.013141	Silvio Romero de Lemos Meira
16	24	0.013141	Vladimir Todorov
17	23	0.013141	Mariofanna G. Milanova
18	5	0.013141	Mohamed E. Fayad
19	19	0.012516	Chengcui Zhang
20	18	0.011890	Waleed W. Smari

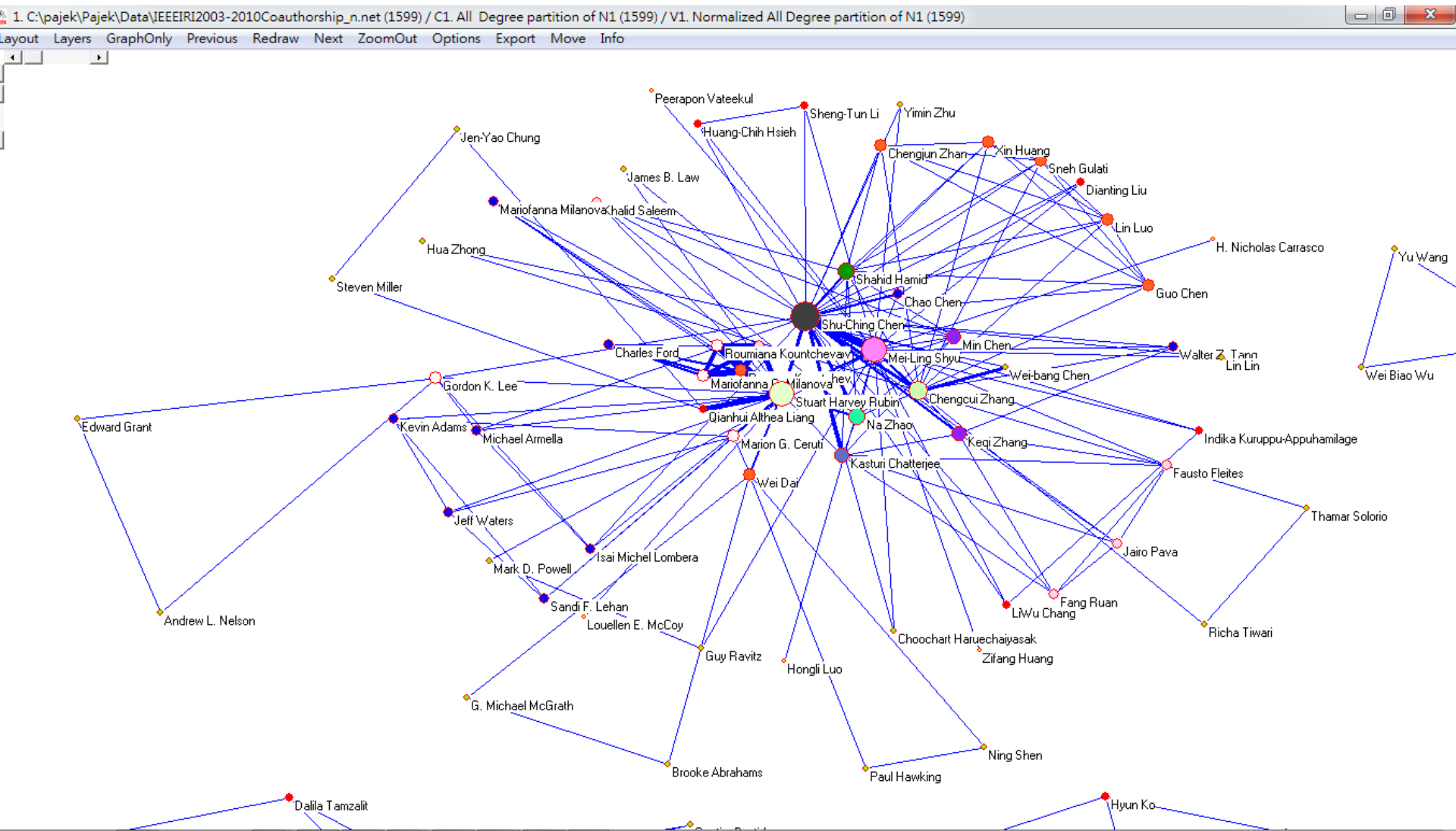
Visualization of IRI (IEEE IRI 2003-2010) co-authorship network (global view)





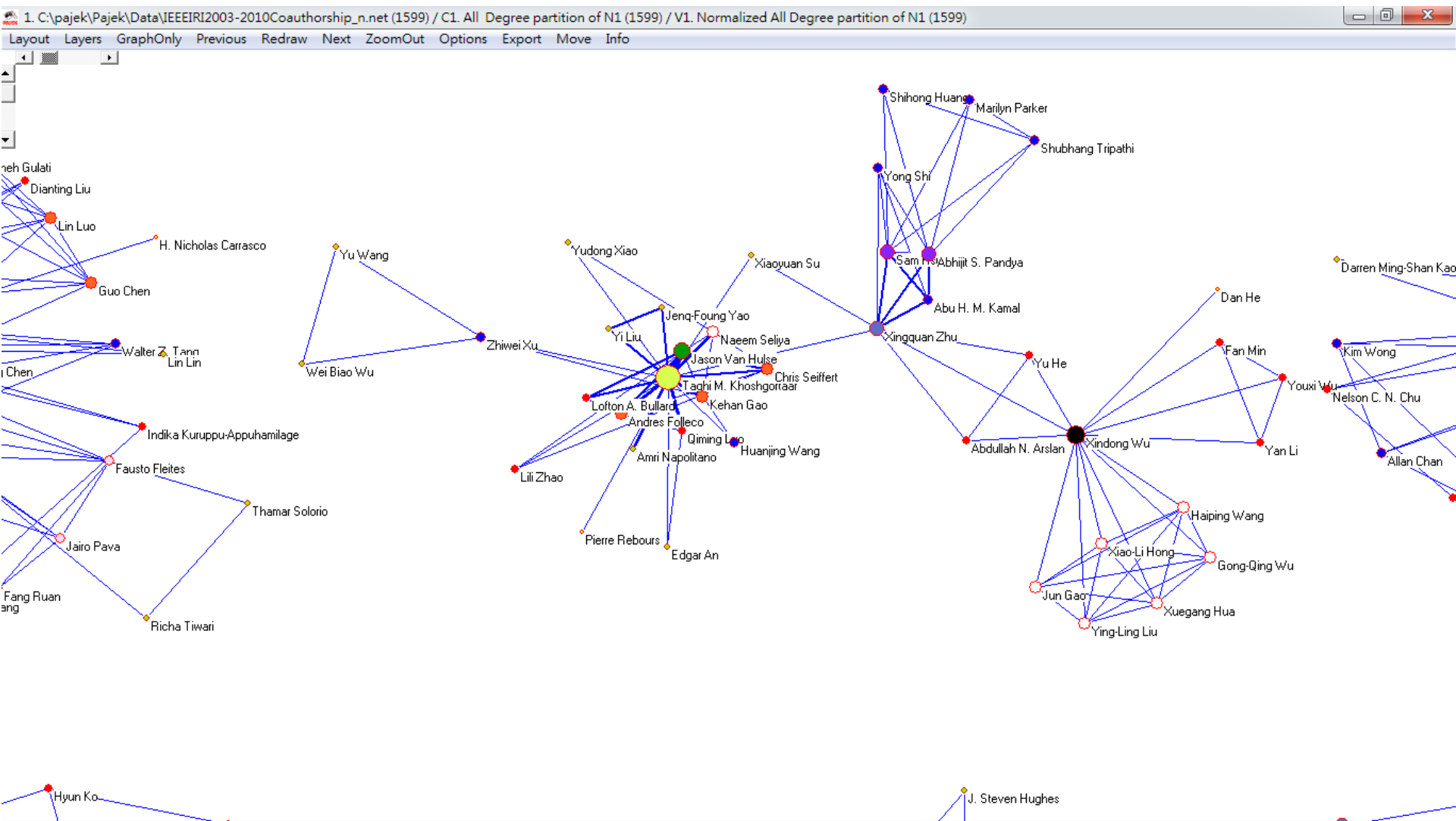
Source: Min-Yuh Day, Sheng-Pao Shih, Weide Chang (2011),
 "Social Network Analysis of Research Collaboration in Information Reuse and Integration"

Visualization of Social Network Analysis

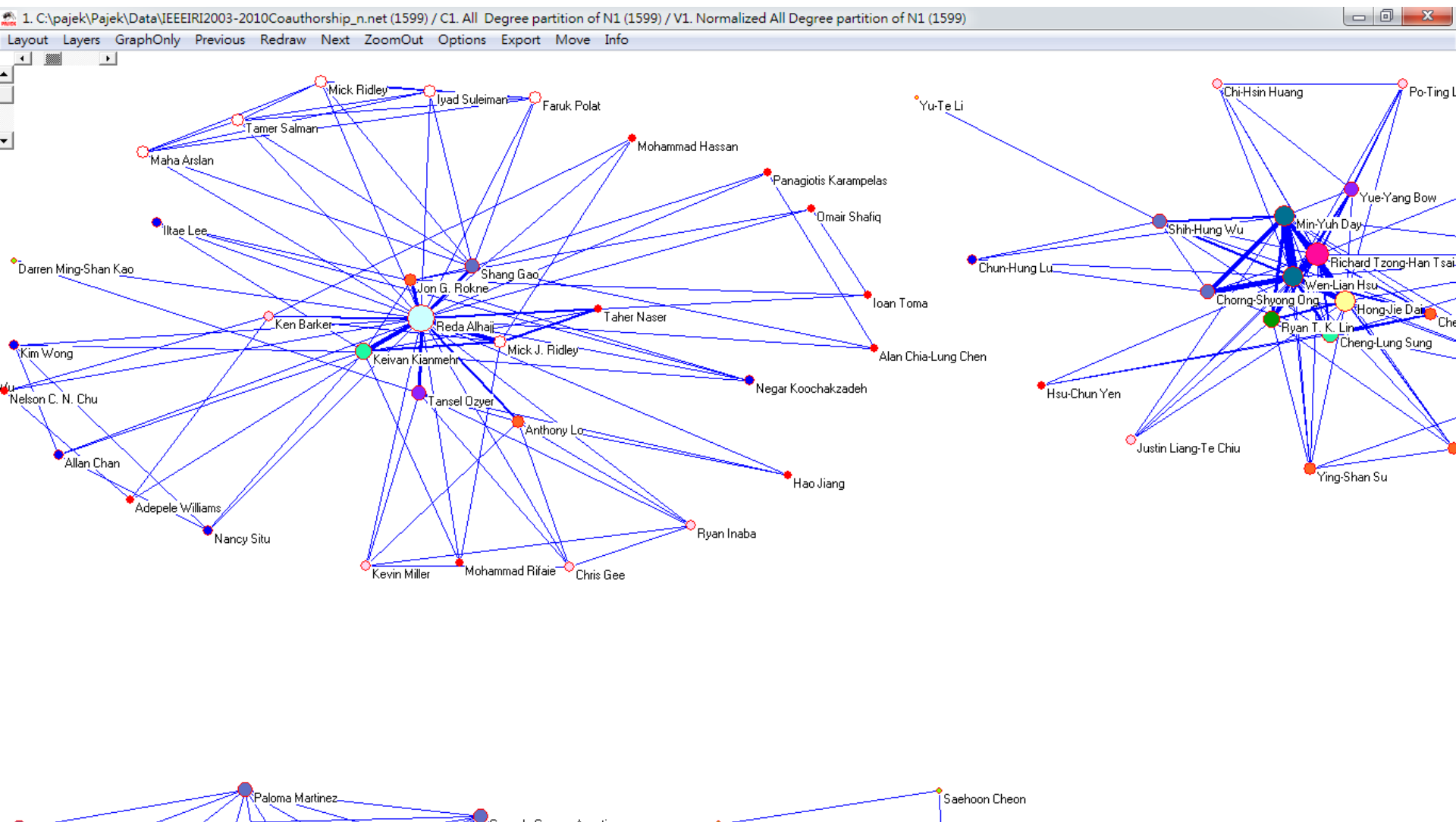


Source: Min-Yuh Day, Sheng-Pao Shih, Weide Chang (2011), "Social Network Analysis of Research Collaboration in Information Reuse and Integration"

Visualization of Social Network Analysis



Visualization of Social Network Analysis



Source: Min-Yuh Day, Sheng-Pao Shih, Weide Chang (2011),
"Social Network Analysis of Research Collaboration in Information Reuse and Integration"

NTCIR 12

QALab-2 Task



**Overview of
NTCIR
Evaluation Activities**



NTCIR

**NII Testbeds and Community for
Information access Research**

NII: National Institute of Informatics



NTCIR

Research Infrastructure for Evaluating Information Access

- A series of evaluation workshops designed to enhance research in information-access technologies by providing an infrastructure for large-scale evaluations.
- Data sets, evaluation methodologies, forum

NTCIR

- **Project started in late 1997**
 - 18 months Cycle

The 12th NTCIR (2015 - 2016)

Evaluation of Information Access Technologies

January 2015 - June 2016

Conference: June 7-10, 2016, NII, Tokyo, Japan



About NTCIR



FAQ

Search

Publications/
Online Proceedings

Data/Tools

NTCIR CMS Site

Related URL's

Contact us

🏠 [NTCIR Home](#) > [NTCIR-12](#)

NTCIR 12

NTCIR-12 Conference

NEWS

NTCIR-12 Aims

Call for Task Proposals

Call for Task
Participation

How to Participate

Task Participation

Task Overview

Data

User Agreement Forms

Organization

NTCIR-12

The 12th NTCIR (2015 - 2016)

Evaluation of Information Access Technologies

January 2015 - June 2016

Conference: June 7-10, 2016, NII, Tokyo, Japan

What's New

July 31, 2015 **NTCIR-12 Task Registration is still possible in each tasks:**

[IMine](#) : until October 31, 2015 (CLOSED)

[MedNLPDoc](#) : until October 31, 2015 (CLOSED)

[MobileClick](#) : until February 1, 2016

[SpokenQuery&Doc](#) : Registration is still possible. Due date please confirm.

[Temporalia](#) : until October 31, 2015 (CLOSED)

[MathIR](#) : until October 31, 2015 (Registration is still possible. Please confirm)

[Lifelog](#) : until the release of the formal run data.

[QALab](#) : until a week before each of the three Formal Runs. [\[detailed schedule\]](#).

[STC](#) : until October 31, 2015

NTCIR 12 (2015-2016) Tasks

- IMine
- MedNLPDoc
- MobileClick
- SpokenQuery&Doc
- Temporalia
- MathIRNEW
- Lifelog
- QA Lab (QA Lab for Entrance Exam; QALab-2)
- STC

NTCIR 12 (2015-2016) Schedule

- 31/July302015: Task Registration Due (extended deadline. Registration is still possible in each task. Please see here.)
- 01/July/2015: Document Set Release *
- July-Dec./2015: Dry Run *
- **Sep./2015-Feb./2016:Formal Run ***
- 01/Feb./2016: Evaluation Results Return
- 01/Feb./2016: Early draft Task Overview Release
- 01/Mar./2016: Draft participant paper submission Due
- 01/May/2016: All camera-ready paper for the Proceedings Due
- **07-10/June/2016:NTCIR-12 Conference & EVIA 2016 in NII, Tokyo, Japan**

QA Lab for Entrance Exam (QALab-2)(2015-2016)

- The goal is investigate the real-world complex Question Answering (QA) technologies using Japanese university entrance exams and their English translation on the subject of "World History (世界史)".
- The questions were selected from two different stages - The National Center Test for University Admissions (センター試験, multiple choice-type questions) and from secondary exams at multiple universities (二次試験, complex questions including essays).

RITE

(**R**ecognizing **I**nference in **T**ext)

NTCIR-9 RITE (2010-2011)

NTCIR-10 RITE-2 (2012-2013)

NTCIR-11 RITE-VAL (2013-2014)

RITE-2

Recognizing
Inference in
Text@NTCIR10

Overview of the Recognizing Inference in Text (RITE-2) at NTCIR-10



**Yotaro
Watanabe**
Tohoku
University



**Yusuke
Miyao**
NII



**Junta
Mizuno**
Tohoku
University



**Tomohide
Shibata**
Kyoto
University



**Hiroshi
Kanayama**
IBM
Research



**Cheng-
Wei Lee**
Academia
Sinica



**Chuan-
Jie Lin**
National Taiwan
Ocean University



**Shuming
Shi**
MSRA



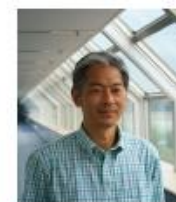
**Teruko
Mitamura**
CMU



**Noriko
Kando**
NII



**Hideki
Shima**
CMU



**Kohichi
Takeda**
IBM
Research

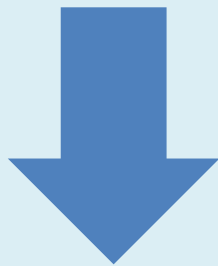
Source: Yotaro Watanabe, Yusuke Miyao, Junta Mizuno, Tomohide Shibata, Hiroshi Kanayama, Cheng-Wei Lee, Chuan-Jie Lin, Shuming Shi, Teruko Mitamura, Noriko Kando, Hideki Shima and Kohichi Takeda, Overview of the Recognizing Inference in Text (RITE-2) at NTCIR-10, Proceedings of NTCIR-10, 2013,

<http://research.nii.ac.jp/ntcir/workshop/OnlineProceedings10/pdf/NTCIR/RITE/01-NTCIR10-RITE2-overview-slides.pdf>

Overview of RITE-2

- RITE-2 is a generic benchmark task that addresses a common semantic inference required in various NLP/IA applications

t_1 : **Yasunari Kawabata** won the Nobel Prize in Literature for his novel “**Snow Country**.”



Can t_2 be inferred from t_1 ?
(entailment?)

t_2 : **Yasunari Kawabata** is the writer of “**Snow Country**.”

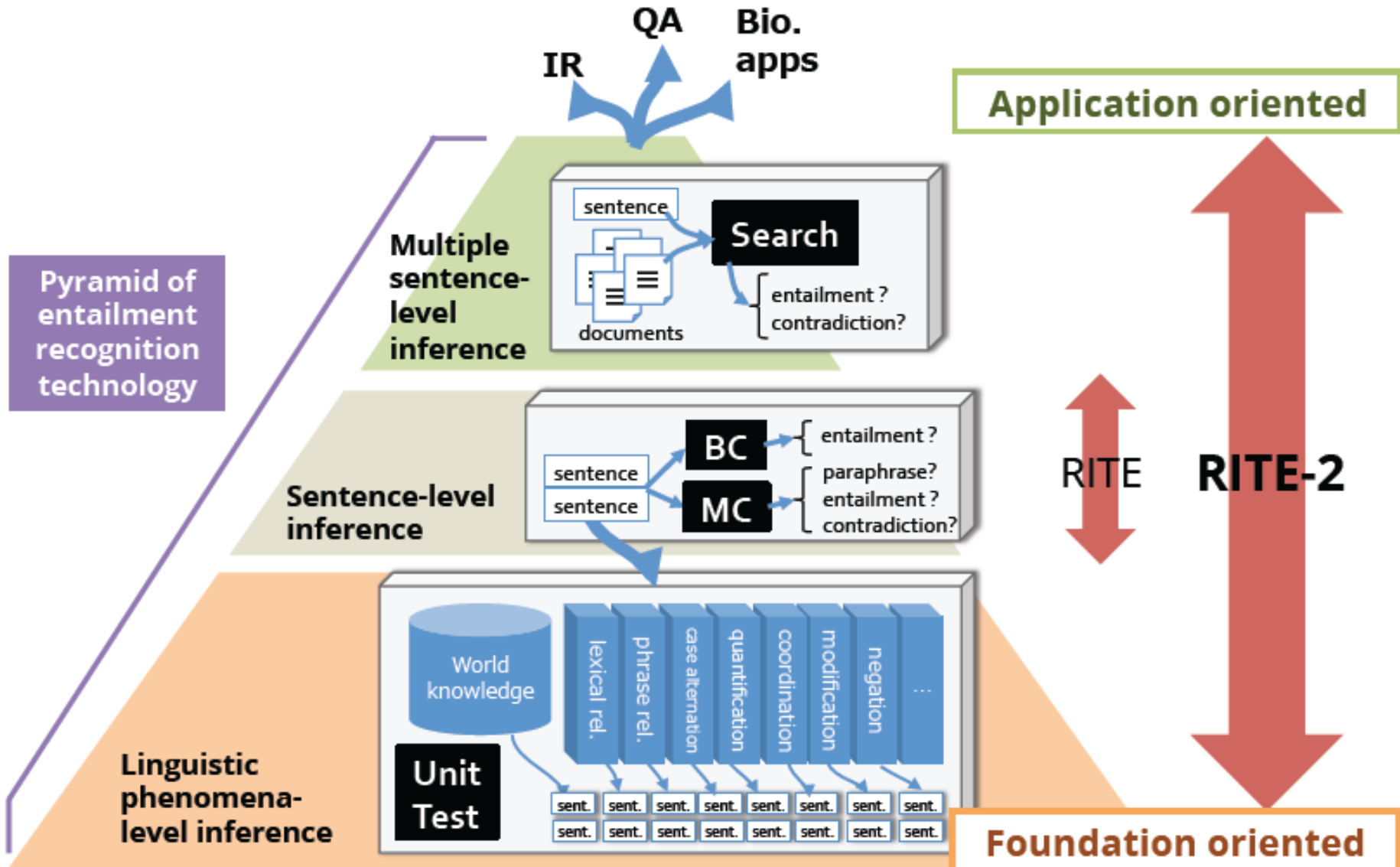


Yasunari Kawabata

Writer

Yasunari Kawabata was a Japanese short story writer and novelist whose spare, lyrical, subtly-shaded prose works won him the Nobel Prize for Literature in 1968, the first Japanese author to receive the award.

RITE vs. RITE-2




Motivation of RITE-2


- Natural Language Processing (NLP) / Information Access (IA) applications
 - Question Answering, Information Retrieval, Information Extraction, Text Summarization, Automatic evaluation for Machine Translation, Complex Question Answering
- The current entailment recognition systems have not been mature enough
 - The highest accuracy on Japanese BC subtask in NTCIR-9 RITE was only 58%
 - There is still enough room to address the task to advance entailment recognition technologies

BC and MC subtasks in RITE-2

t_1 : **Yasunari Kawabata** won the Nobel Prize in Literature for his novel “**Snow Country**.”

t_2 : **Yasunari Kawabata** is the writer of “**Snow Country**.”

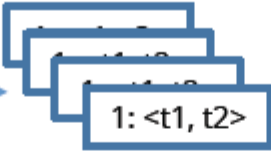
- **BC subtask** 
 - Entailment (t_1 entails t_2) or Non-Entailment (otherwise)

- **MC subtask** 
 - Bi-directional Entailment (t_1 entails t_2 & t_2 entails t_1)
 - Forward Entailment (t_1 entails t_2 & t_2 does not entail t_1)
 - Contradiction (t_1 contradicts t_2 or cannot be true at the same time)
 - Independence (otherwise)

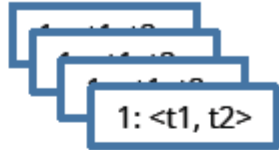
Development of BC and MC data



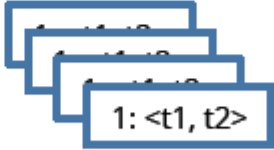
retrieve pairs
of sentences



edit pairs
if needed



for each example,
5 annotators
assigned its
semantic label



accept an example if
4 or more annotators
assigned the same label
to the example



Entrance Exam subtasks (Japanese only)

Entrance exam problem

National Center Test for University Admission (Daigaku Nyushi Center Shiken)

第1問 モニュメントや歴史的建造物について述べた次の文章A～Cを読み、下の問い(問1～11)に答えよ。(配点 33)

A 現在、アテネの中心部の丘にその偉容を誇る①**パルテノン神殿**は、古代ギリシアを象徴する歴史的建造物である。この神殿は、②**オスマン帝国**の支配下でモスクとして利用されたこともあったが、18世紀には廃墟となっていた。1799年にイギリスの大使としてイスタンブールに赴任したエルギン卿は、③**ギリシア**を訪れ、パルテノン神殿の遺跡から彫刻額を収集し、本国に送った。今日、大英博物館で「エルギン・マーブル」として展示されているものがそれである。1987年、パルテノン神殿は、世界文化遺産として登録された。

問3 下線部②の国について述べた文として最も適当なものを、次の①～④のうちから一つ選べ。

- ① スレイマン1世の時代が最盛期であった。
- ② 国教はシーア派のイスラーム教であった。
- ③ パルカン半島に誕生した後、小アジアへ進出した。
- ④ ベルリン会議により、ボスニア＝ヘルツェゴヴィナの統治権を得た。



WIKIPEDIA
The Free Encyclopedia

スレイマン1世

スルタン・スレイマン1世 (Kanuni Sultan Süleyman, **オスマン語** سلطان Sülaymān, **トルコ語** Süleyman, 1494年11月6日 - 1566年9月5日) は、オスマン帝国の第10代皇帝(在位: 1520年 - 1566年)。

46年の長期にわたる在位の中で13回もの対外遠征を行い、数多くの軍事的成功を収めてオスマン帝国を最盛期に導いた。英語では、「**壮麗帝**(the Magnificent)のあだ名で呼ばれ、日本ではしばしば「**スレイマン大帝**」と称される。トルコでは法典を編纂し帝国の制度を整備したことから「**立法帝**(カーヌニー **القانونی** al-Qanuni) / Kanuni)」のあだ名で知られている。

t_1 : スレイマン1世は数多くの軍事的成功を収めてオスマン帝国を最盛期に導いた。 (Suleiman I contributed in a lot of military successes and led the Ottoman Empire to its peak.)

t_2 : オスマン帝国ではスレイマン1世の時代が最盛期であった。 (The Ottoman Empire's peak was during the reign of Suleiman I).

Entrance Exam subtask: BC and Search

- **Entrance Exam BC**
 - Binary-classification problem (Entailment or Nonentailment)
 - t1 and t2 are given
- **Entrance Exam Search**
 - Binary-classification problem (Entailment or Nonentailment)
 - t2 and a set of documents are given
 - Systems are required to search sentences in Wikipedia and textbooks to decide semantic labels

UnitTest (Japanese only)

- **Motivation**
 - Evaluate how systems can handle linguistic
 - phenomena that affects entailment relations
- **Task definition**
 - Binary classification problem (same as BC subtask)

RITE4QA (Chinese only)

- **Motivation**

- Can an entailment recognition system rank a set of unordered answer candidates in QA?

- **Dataset**

- Developed from NTCIR-7 and NTCIR-8 CLQA data
 - t1: answer-candidate-bearing sentence
 - t2: a question in an affirmative form

- **Requirements**

- Generate confidence scores for ranking process

Evaluation Metrics

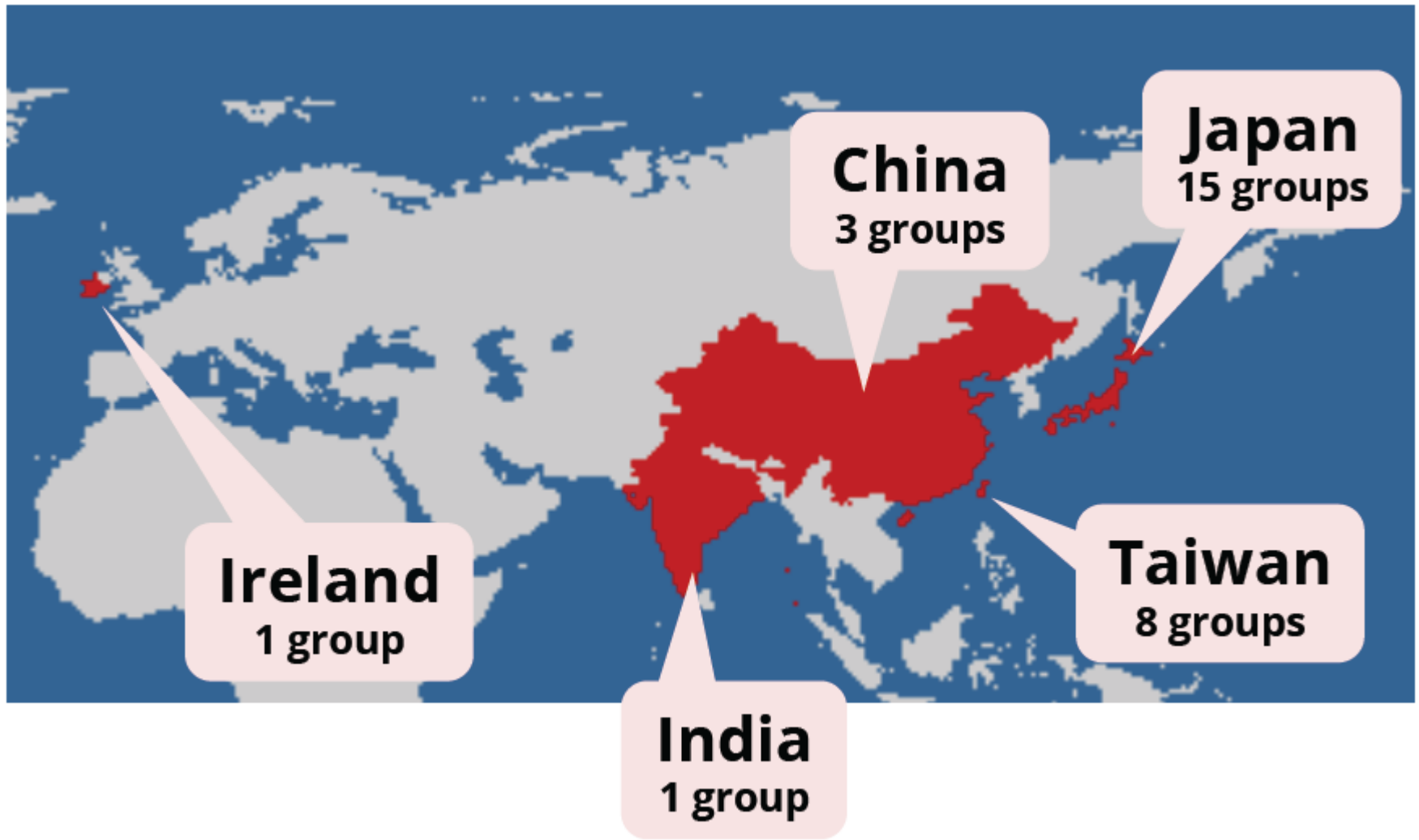
- **Macro F1 and Accuracy**
(BC, MC, ExamBC, ExamSearch and UnitTest)

$$MacroF1 = \frac{1}{|C|} \sum_{c \in C} F1_c \quad Accuracy = 100 \times \frac{N_{correct}}{N_{examples}}$$

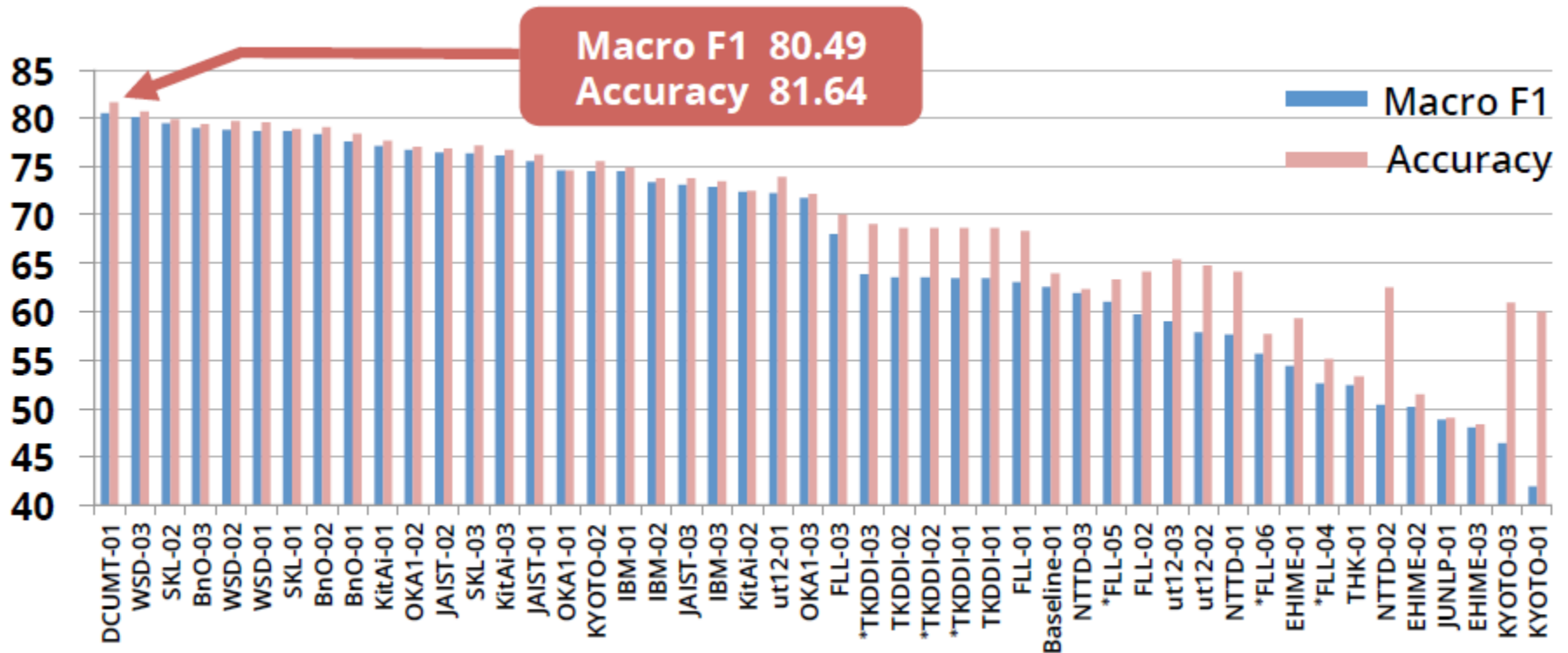
- **Correct Answer Ratio (Entrance Exam)**
 - Y/N labels are mapped into selections of answers and calculate accuracy of the answers
- **Top1 and MRR (RITE4QA)**

$$Top1 = \frac{1}{|Q|} \sum_{i=1}^{|Q|} [\text{top answer is correct}] \quad MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$

Countries/Regions of Participants

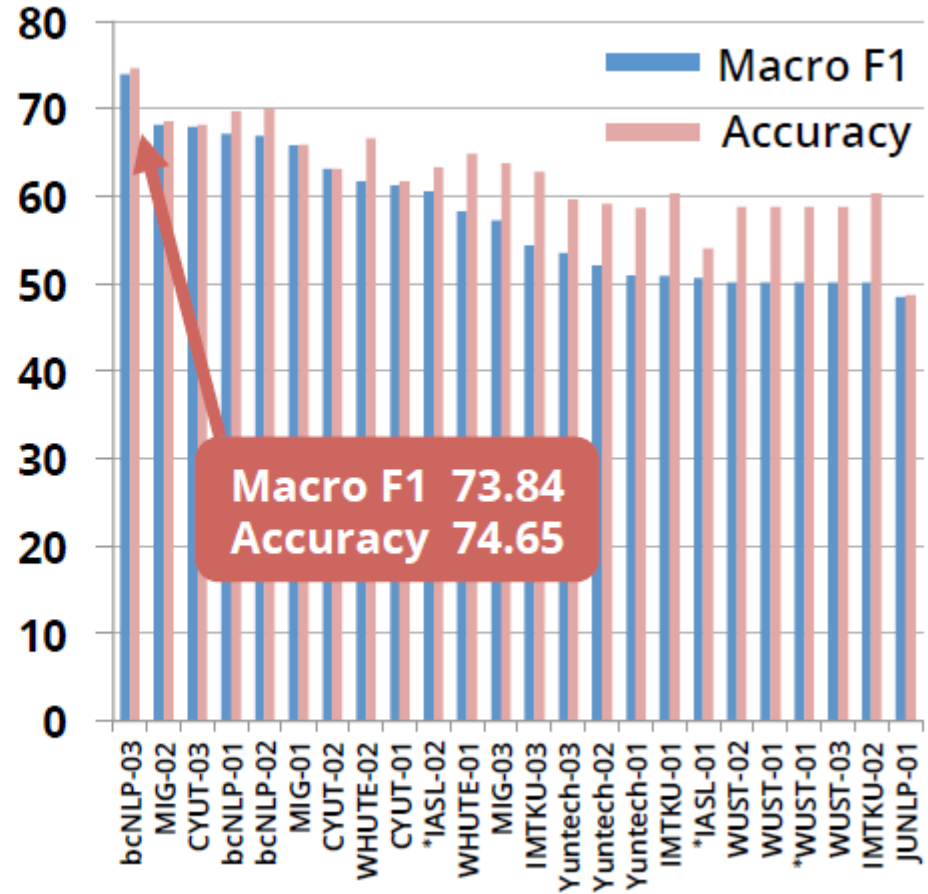
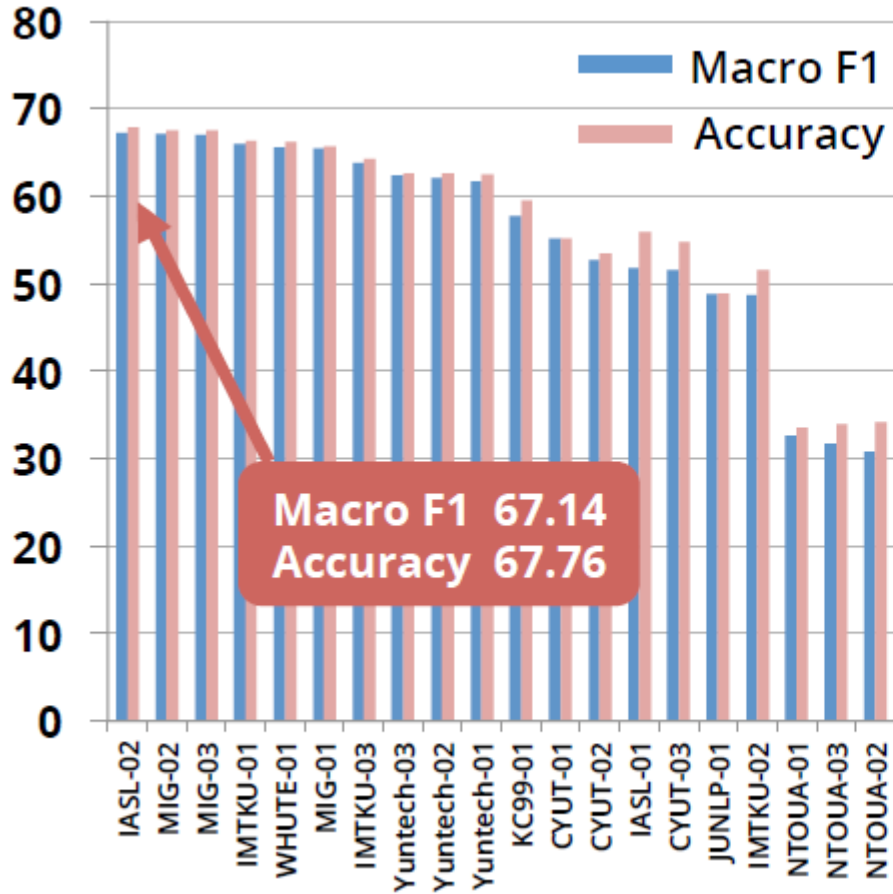


Formal Run Results: BC (Japanese)



- The best system achieved over 80% of accuracy (The highest score in BC subtask at RITE was 58%)
- The difference is caused by
 - Advancement of entailment recognition technologies
 - Strict data filtering in the data development

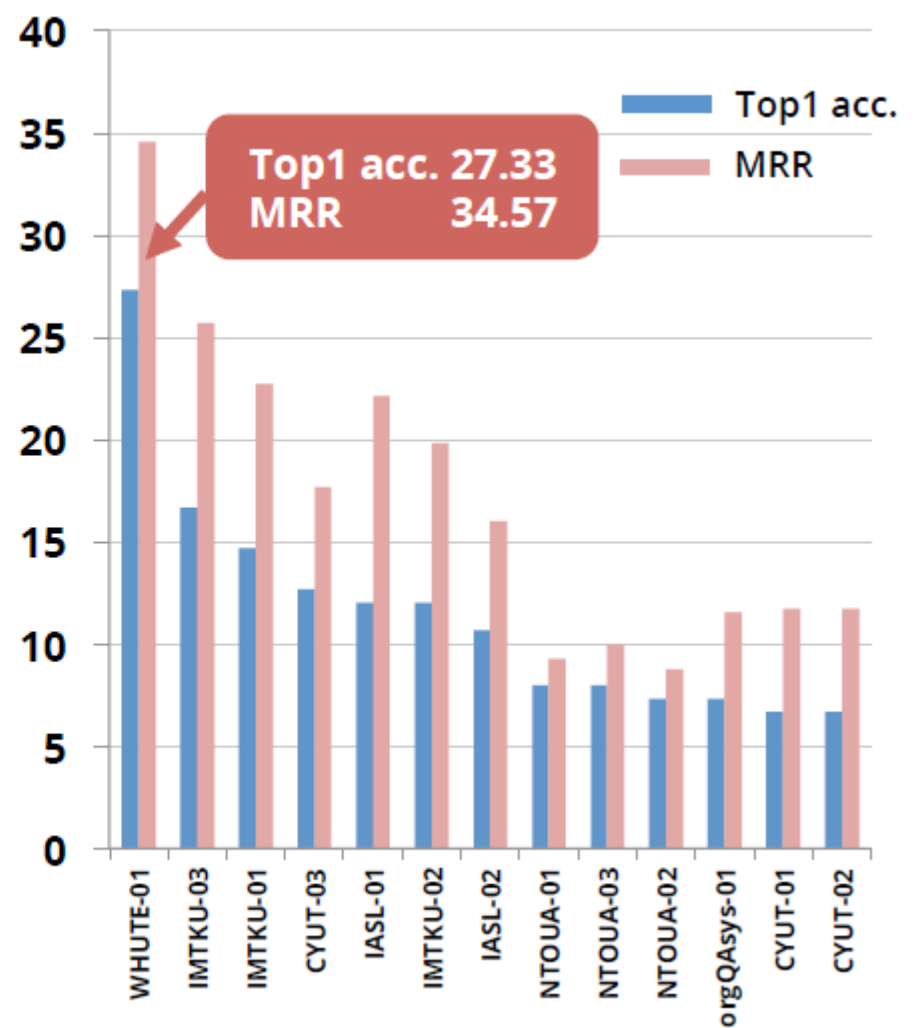
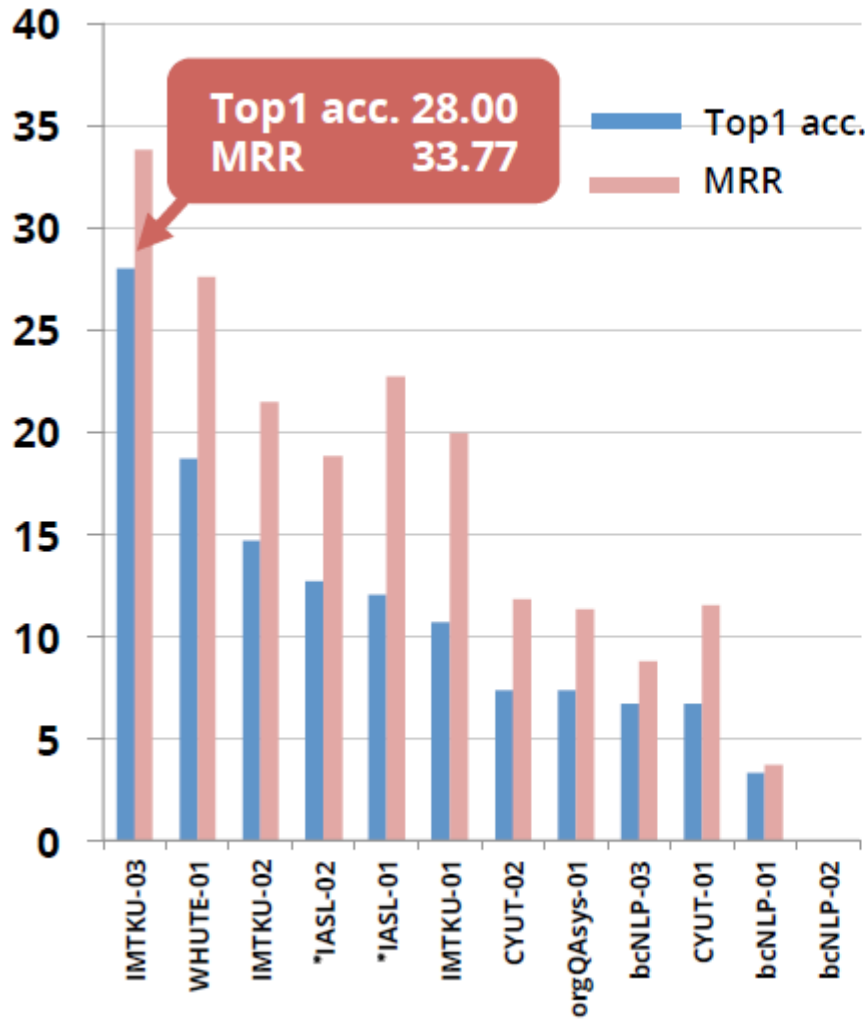
BC (Traditional/Simplified Chinese)



The top scores are almost the same as those in NTCIR-9 RITE

RITE4QA

(Traditional/Simplified Chinese)



Participant's approaches in RITE-2

- **Category**
 - Statistical (50%)
 - Hybrid (27%)
 - Rule-based (23%)
- **Fundamental approach**
 - Overlap-based (77%)
 - Alignment-based (63%)
 - Transformation-based (23%)

Summary of types of information explored in RITE-2

- Character/word overlap (85%)
- Syntactic information (67%)
- Temporal/numerical information (63%)
- Named entity information (56%)
- Predicate-argument structure (44%)
- Entailment relations (30%)
- Polarity information (7%)
- Modality information (4%)

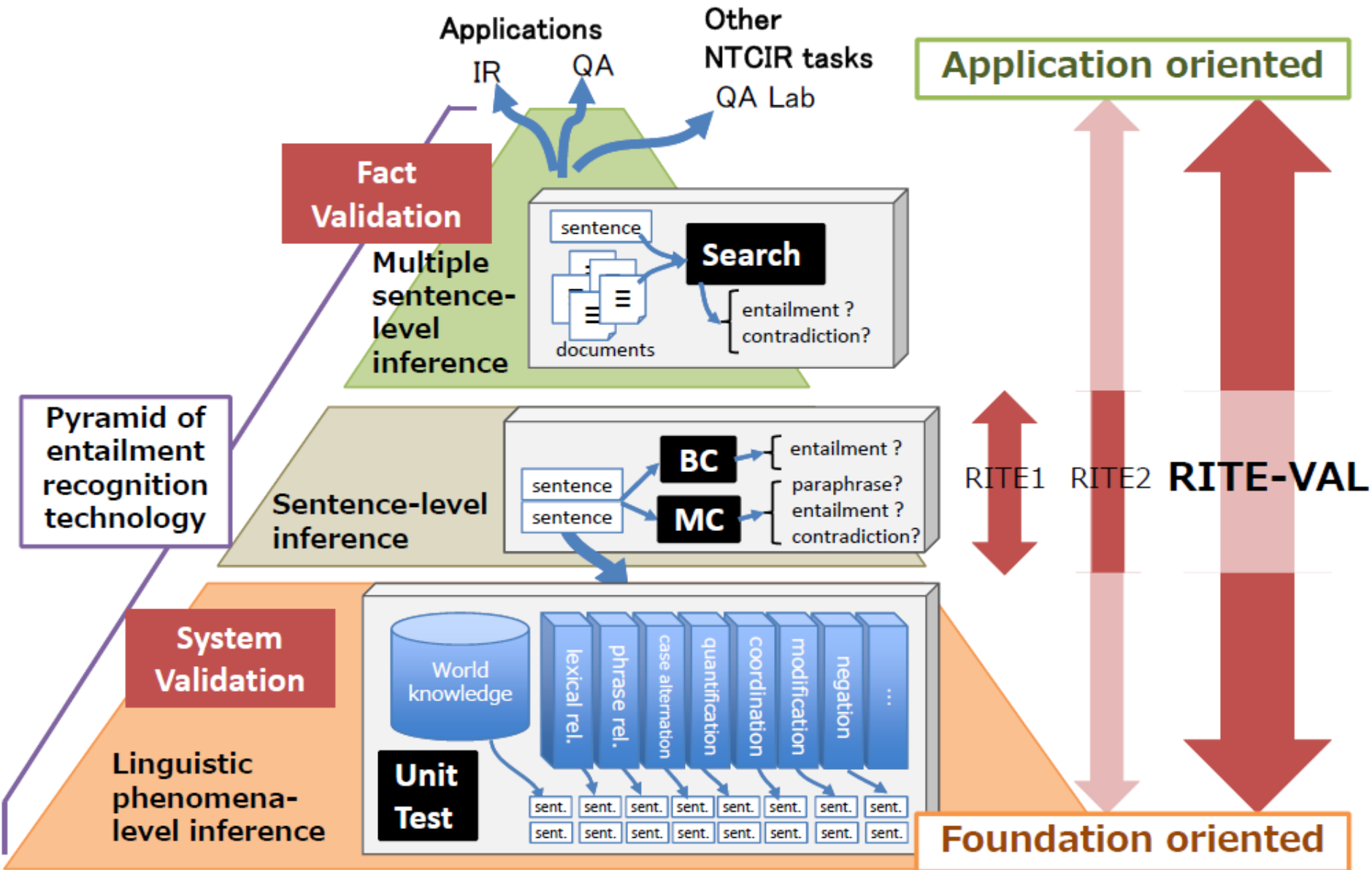
Summary of Resources Explored in RITE-2

- **Japanese**
 - Wikipedia (10)
 - Japanese WordNet (9)
 - ALAGIN Entailment DB (5)
 - Nihongo Goi-Taikei (2)
 - Bunruigoihyo (2)
 - Iwanami Dictionary (2)
- **Chinese**
 - Chinese WordNet (3)
 - TongYiCi CiLin (3)
 - HowNet (2)

Advanced approaches in RITE-2

- **Logical approaches**
 - Dependency-based Compositional Semantics (DCS) [BnO], Markov Logic [EHIME], Natural Logic [THK]
- **Alignment**
 - GIZA [CYUT], ILP [FLL], Labeled Alignment [bcNLP, THK]
- **Search Engine**
 - Google and Yahoo [DCUMT]
- **Deep Learning**
 - RNN language models [DCUMT]
- **Probabilistic Models**
 - N-gram HMM [DCUMT], LDA [FLL]
- **Machine Translation**
 - [JUNLP, JAIST, KC99]

RITE-VAL



Main two tasks of RITE-VAL

Fact Validation



t_2 : *The Kamakura Shogunate began in Japan in the 12th century.*



Search for evidence or counter-evidence for t_2 .

Docs entail t_2 .

Docs contradict t_2 .

System Validation

Category: modification

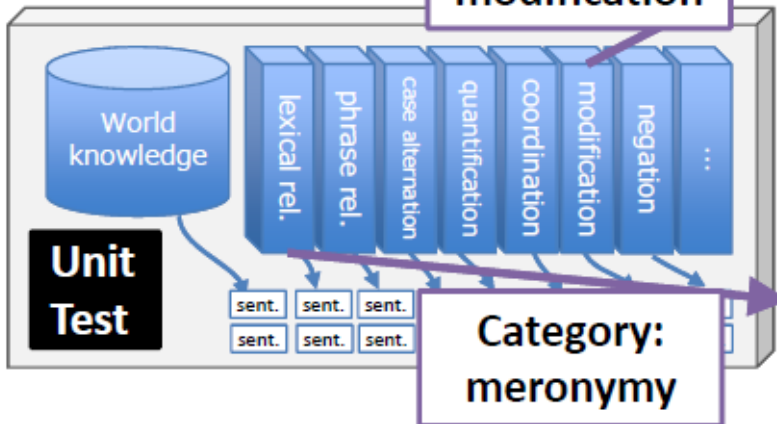
t_1 : *In the Meiji Constitution, legal clear distinction between the Imperial Family and Japan had been allowed.*

t_2 : *In the Meiji Constitution, distinction between the Imperial Family and Japan had been allowed.*

Category: meronymy

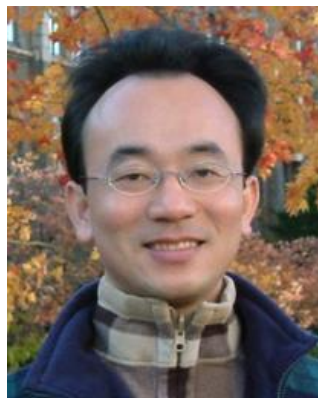
t_1 : *In the Meiji Constitution, distinction between the Imperial Family and Japan had been allowed.*

t_2 : *In the Meiji Constitution, distinction between the Emperor and Japan had been allowed.*



IMTKU Textual Entailment System for Recognizing Inference in Text at NTCIR-9 RITE

Department of Information Management
Tamkang University, Taiwan



Min-Yuh Day

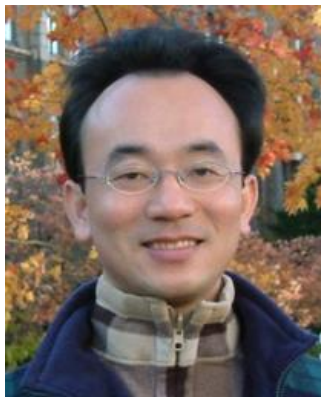


Chun Tu

myday@mail.tku.edu.tw

IMTKU Textual Entailment System for Recognizing Inference in Text at NTCIR-10 RITE-2

Department of Information Management
Tamkang University, Taiwan



Min-Yuh Day



Chun Tu



Hou-Cheng Vong



Shih-Wei Wu



Shih-Jhen Huang

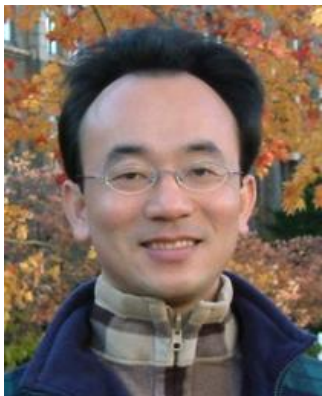
myday@mail.tku.edu.tw

IMTKU Textual Entailment System for Recognizing Inference in Text at NTCIR-11 RITE-VAL

Tamkang University

淡江大學

2014



Min-Yuh Day



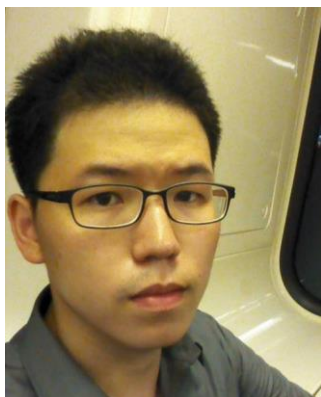
Ya-Jung Wang



Che-Wei Hsu



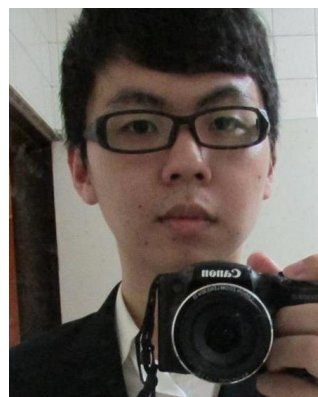
En-Chun Tu



Huai-Wen Hsu



Yu-An Lin



Shang-Yu Wu



Yu-Hsuan Tai



Cheng-Chia Tsai

IMTKU Question Answering System for Entrance Exam at NTCIR-12 QALab-2

Tamkang University

淡江大學

2016



Min-Yuh Day



Cheng-Chia Tsai



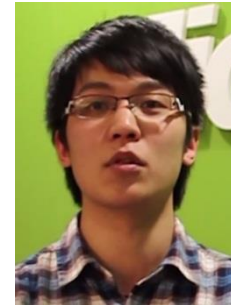
Wei-Chun Chung



Hsiu-Yuan Chang



Lin-Jin Kun



Yuan-Jie Tsai



Tzu-Jui Sun



Yu-Ming Guo



Yue-Da Lin



Wei-Ming Chen



Yun-Da Tsai



Cheng-Jih Han



Yi-Jing Lin

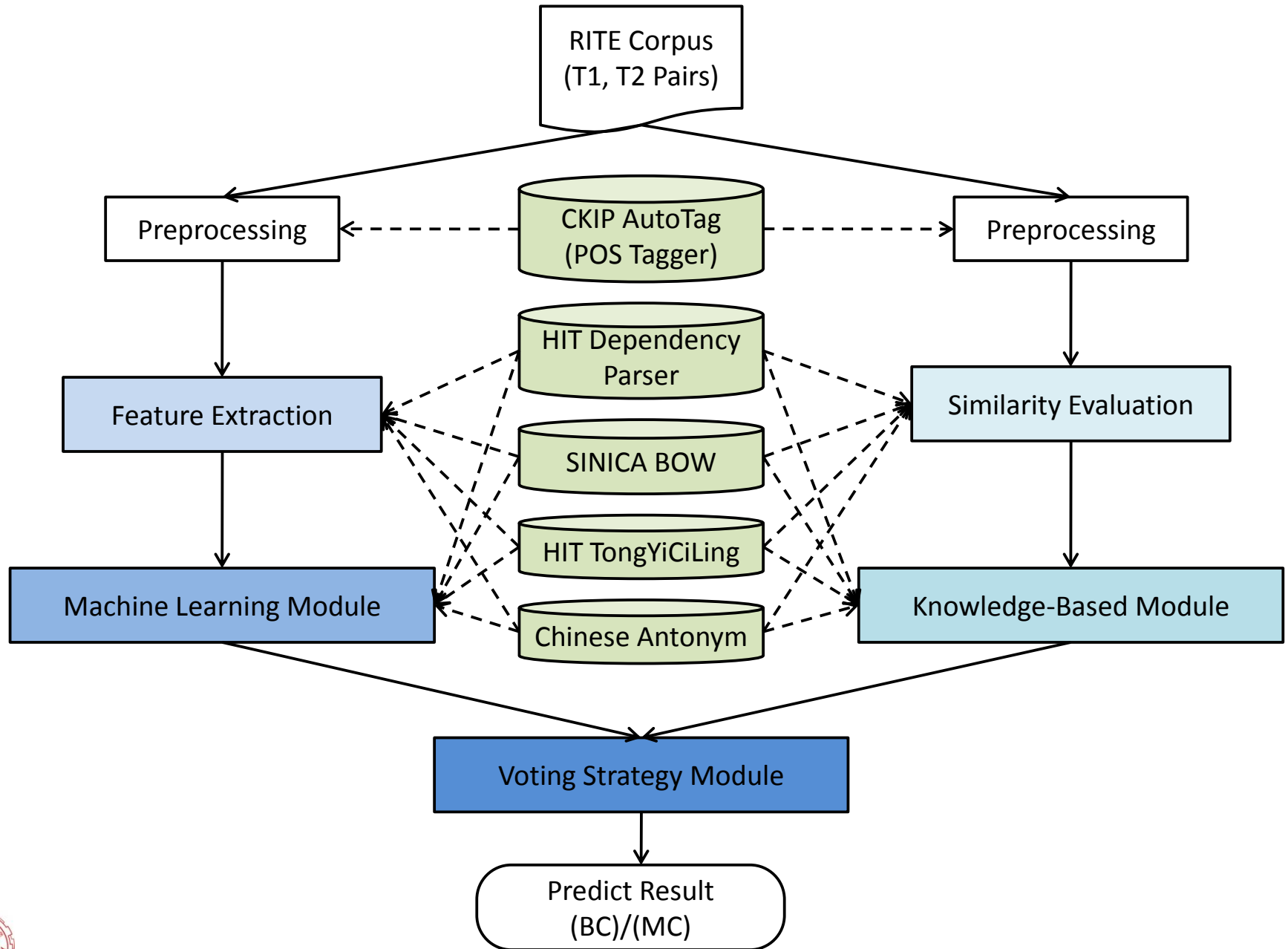


Yi-Heng Chiang

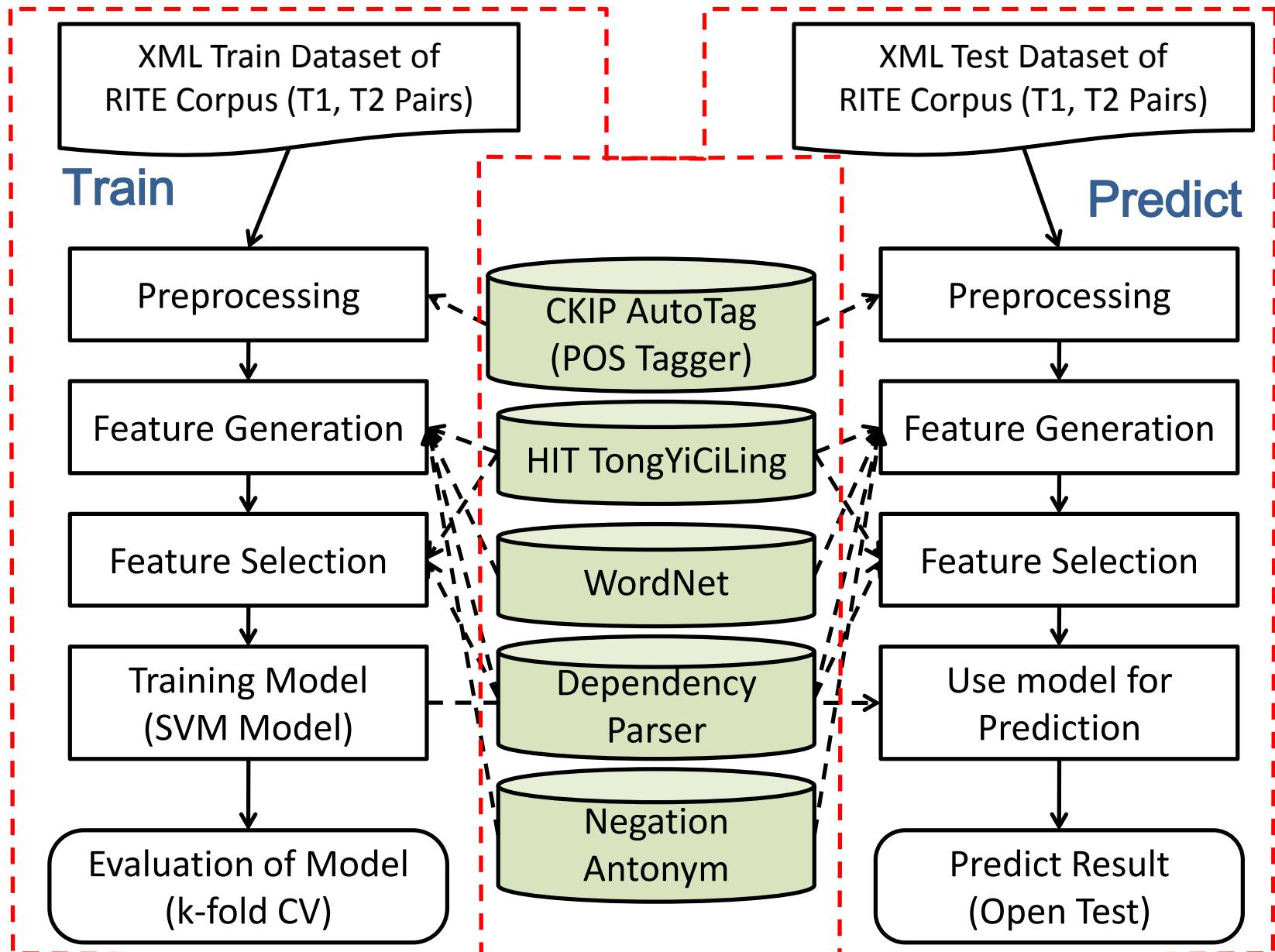


Ching-Yuan Chien

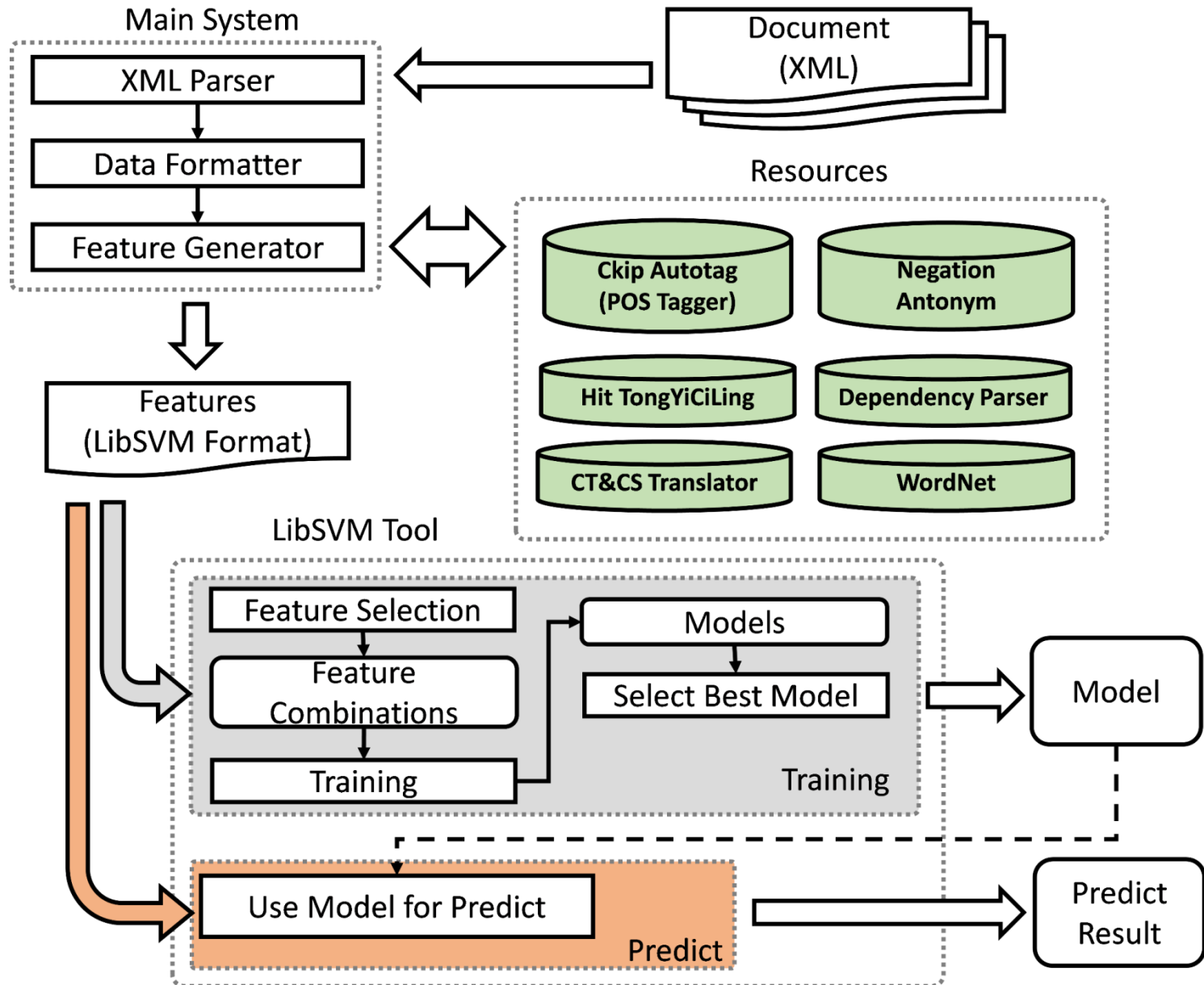
IMTKU System Architecture for NTCIR-9 RITE



IMTKU System Architecture for NTCIR-10 RITE-2



IMTKU System Framework for NTCIR -11 RITE-VAL



IMTKU at NTCIR

- The first place in the CS-RITE4QA subtask of the NTCIR-10 Recognizing Inference in TExt (RITE) task. (2013)
- The second place in the CT-RITE4QA subtask of the NTCIR-10 Recognizing Inference in TExt (RITE) task. (2013)
- The first place in the CT-RITE4QA subtask of the NTCIR-9 Recognizing Inference in TExt (RITE) task. (2011)
- The first place in the CS-RITE4QA subtask of the NTCIR-9 Recognizing Inference in TExt (RITE) task. (2011)
- The second place in the CT-MC subtask of the NTCIR-9 Recognizing Inference in TExt (RITE) task. (2011)

Summary

- Big Data Analytics on Social Media
- Analyzing the Social Web:
Social Network Analysis
- NTCIR 12 QALab-2 Task

References

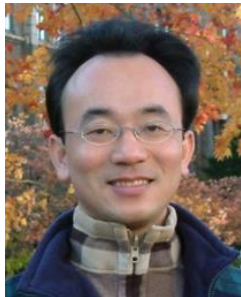
- Jiawei Han and Micheline Kamber (2011),
Data Mining: Concepts and Techniques, Third Edition, Elsevier
- Jennifer Golbeck (2013),
Analyzing the Social Web, Morgan Kaufmann
- Stephan Kudyba (2014),
Big Data, Mining, and Analytics: Components of Strategic
Decision Making, Auerbach Publications
- Hiroshi Ishikawa (2015),
Social Big Data Mining, CRC Press

Q & A

Big Data Analytics on Social Media (社群媒體大數據分析)

Time: 2015/12/25 (14:00-15:30)

Place: S402, Ming Chuan University



Min-Yuh Day

戴敏育

Assistant Professor

專任助理教授

Dept. of Information Management, Tamkang University

淡江大學 資訊管理學系

<http://mail.tku.edu.tw/myday/>

2015-12-25

