Vincent Corruble
Masayuki Takeda
Einoshin Suzuki (Eds.)

# Discovery Science

**10th International Conference, DS 2007
Sendai, Japan, October 2007
Proceedings**

DS

Vincent Corruble   Masayuki Takeda
Einoshin Suzuki (Eds.)

# Discovery Science

10th International Conference, DS 2007
Sendai, Japan, October 1-4, 2007
Proceedings

 Springer

Volume Editors

Vincent Corruble
Université Pierre et Marie Curie (Paris 6)
Laboratoire d'Informatique de Paris 6
104 avenue du Président Kennedy, 75016 Paris, France
E-mail: Vincent.Corruble@lip6.fr

Masayuki Takeda
Einoshin Suzuki
Kyushu University
Department of Informatics
744 Motooka, Nishi, Fukuoka 819-0395, Japan
E-mail: {takeda, suzuki}@i.kyushu-u.ac.jp

# Lecture Notes in Artificial Intelligence (LNAI)

Vol. 4570: H.G. Okuno, M. Ali (Eds.), New Trends in Applied Artificial Intelligence. XXI, 1194 pages. 2007.

Vol. 4565: D.D. Schmorrow, L.M. Reeves (Eds.), Foundations of Augmented Cognition. XIX, 450 pages. 2007.

Vol. 4562: D. Harris (Ed.), Engineering Psychology and Cognitive Ergonomics. XXIII, 879 pages. 2007.

Vol. 4548: N. Olivetti (Ed.), Automated Reasoning with Analytic Tableaux and Related Methods. X, 245 pages. 2007.

Vol. 4539: N.H. Bshouty, C. Gentile (Eds.), Learning Theory. XII, 634 pages. 2007.

Vol. 4529: P. Melin, O. Castillo, L.T. Aguilar, J. Kacprzyk, W. Pedrycz (Eds.), Foundations of Fuzzy Logic and Soft Computing. XIX, 830 pages. 2007.

Vol. 4520: M.V. Butz, O. Sigaud, G. Pezzulo, G. Baldassarre (Eds.), Anticipatory Behavior in Adaptive Learning Systems. X, 379 pages. 2007.

Vol. 4511: C. Conati, K. McCoy, G. Paliouras (Eds.), User Modeling 2007. XVI, 487 pages. 2007.

Vol. 4509: Z. Kobti, D. Wu (Eds.), Advances in Artificial Intelligence. XII, 552 pages. 2007.

Vol. 4496: N.T. Nguyen, A. Grzech, R.J. Howlett, L.C. Jain (Eds.), Agent and Multi-Agent Systems: Technologies and Applications. XXI, 1046 pages. 2007.

Vol. 4483: C. Baral, G. Brewka, J. Schlipf (Eds.), Logic Programming and Nonmonotonic Reasoning. IX, 327 pages. 2007.

Vol. 4482: A. An, J. Stefanowski, S. Ramanna, C.J. Butz, W. Pedrycz, G. Wang (Eds.), Rough Sets, Fuzzy Sets, Data Mining and Granular Computing. XIV, 585 pages. 2007.

Vol. 4481: J. Yao, P. Lingras, W.-Z. Wu, M. Szczuka, N.J. Cercone, D. Ślęzak (Eds.), Rough Sets and Knowledge Technology. XIV, 576 pages. 2007.

Vol. 4476: V. Gorodetsky, C. Zhang, V.A. Skormin, L. Cao (Eds.), Autonomous Intelligent Systems: Multi-Agents and Data Mining. XIII, 323 pages. 2007.

Vol. 4456: Y. Wang, Y.-m. Cheung, H. Liu (Eds.), Computational Intelligence and Security. XXIII, 1118 pages. 2007.

Vol. 4455: S. Muggleton, R. Otero, A. Tamaddoni-Nezhad (Eds.), Inductive Logic Programming. XII, 456 pages. 2007.

Vol. 4452: M. Fasli, O. Shehory (Eds.), Agent-Mediated Electronic Commerce. VIII, 249 pages. 2007.

Vol. 4451: T.S. Huang, A. Nijholt, M. Pantic, A. Pentland (Eds.), Artifical Intelligence for Human Computing. XVI, 359 pages. 2007.

Vol. 4441: C. Müller (Ed.), Speaker Classification. X, 309 pages. 2007.

Vol. 4438: L. Maicher, A. Sigel, L.M. Garshol (Eds.), Leveraging the Semantics of Topic Maps. X, 257 pages. 2007.

Vol. 4434: G. Lakemeyer, E. Sklar, D.G. Sorrenti, T. Takahashi (Eds.), RoboCup 2006: Robot Soccer World Cup X. XIII, 566 pages. 2007.

Vol. 4429: R. Lu, J.H. Siekmann, C. Ullrich (Eds.), Cognitive Systems. X, 161 pages. 2007.

Vol. 4428: S. Edelkamp, A. Lomuscio (Eds.), Model Checking and Artificial Intelligence. IX, 185 pages. 2007.

Vol. 4426: Z.-H. Zhou, H. Li, Q. Yang (Eds.), Advances in Knowledge Discovery and Data Mining. XXV, 1161 pages. 2007.

Vol. 4411: R.H. Bordini, M. Dastani, J. Dix, A.E.F. Seghrouchni (Eds.), Programming Multi-Agent Systems. XIV, 249 pages. 2007.

Vol. 4410: A. Branco (Ed.), Anaphora: Analysis, Algorithms and Applications. X, 191 pages. 2007.

Vol. 4399: T. Kovacs, X. Llorà, K. Takadama, P.L. Lanzi, W. Stolzmann, S.W. Wilson (Eds.), Learning Classifier Systems. XII, 345 pages. 2007.

Vol. 4390: S.O. Kuznetsov, S. Schmidt (Eds.), Formal Concept Analysis. X, 329 pages. 2007.

Vol. 4389: D. Weyns, H. Van Dyke Parunak, F. Michel (Eds.), Environments for Multi-Agent Systems III. X, 273 pages. 2007.

Vol. 4386: P. Noriega, J. Vázquez-Salceda, G. Boella, O. Boissier, V. Dignum, N. Fornara, E. Matson (Eds.), Coordination, Organizations, Institutions, and Norms in Agent Systems II. XI, 373 pages. 2007.

Vol. 4384: T. Washio, K. Satoh, H. Takeda, A. Inokuchi (Eds.), New Frontiers in Artificial Intelligence. IX, 401 pages. 2007.

Vol. 4371: K. Inoue, K. Satoh, F. Toni (Eds.), Computational Logic in Multi-Agent Systems. X, 315 pages. 2007.

Vol. 4369: M. Umeda, A. Wolf, O. Bartenstein, U. Geske, D. Seipel, O. Takata (Eds.), Declarative Programming for Knowledge Management. X, 229 pages. 2006.

Vol. 4363: B.D. ten Cate, H.W. Zeevat (Eds.), Logic, Language, and Computation. XII, 281 pages. 2007.

Vol. 4343: C. Müller (Ed.), Speaker Classification I. X, 355 pages. 2007.

Vol. 4342: H. de Swart, E. Orłowska, G. Schmidt, M. Roubens (Eds.), Theory and Applications of Relational Structures as Knowledge Instruments II. X, 373 pages. 2006.

Vol. 4335: S.A. Brueckner, S. Hassas, M. Jelasity, D. Yamins (Eds.), Engineering Self-Organising Systems. XII, 212 pages. 2007.

Vol. 4334: B. Beckert, R. Hähnle, P.H. Schmitt (Eds.), Verification of Object-Oriented Software. XXIX, 658 pages. 2007.

Vol. 4333: U. Reimer, D. Karagiannis (Eds.), Practical Aspects of Knowledge Management. XII, 338 pages. 2006.

Vol. 4327: M. Baldoni, U. Endriss (Eds.), Declarative Agent Languages and Technologies IV. VIII, 257 pages. 2006.

Vol. 4314: C. Freksa, M. Kohlhase, K. Schill (Eds.), KI 2006: Advances in Artificial Intelligence. XII, 458 pages. 2007.

Vol. 4304: A. Sattar, B.-h. Kang (Eds.), AI 2006: Advances in Artificial Intelligence. XXVII, 1303 pages. 2006.

¥516.⁰⁰元

# Preface

This volume contains the papers presented at DS-2007: The Tenth International Conference on Discovery Science held in Sendai, Japan, October 1–4, 2007.

The main objective of the Discovery Science (DS) conference series is to provide an open forum for intensive discussions and the exchange of new ideas and information among researchers working in the area of automating scientific discovery or working on tools for supporting the human process of discovery in science. It has been a successful arrangement in the past to co-locate the DS conference with the International Conference on Algorithmic Learning Theory (ALT). This combination of ALT and DS allows for a comprehensive treatment of the whole range, from theoretical investigations to practical applications. Continuing this tradition, DS 2007 was co-located with the 18th ALT conference (ALT 2007). The proceedings of ALT 2007 were published as a twin volume 4754 of the LNCS series.

The International Steering Committee of the Discovery Science conference series provided important advice on a number of issues during the planning of Discovery Science 2007. The members of the Steering Committee are Einoshin Suzuki (Kyushu University, Chair), Achim G. Hoffmann (University of New South Wales, Vice Chair), Setsuo Arikawa (Kyushu University), Hiroshi Motoda (Osaka University), Masahiko Sato (Kyoto University), Satoru Miyano (University of Tokyo), Thomas Zeugmann (Hokkaido University), Ayumi Shinohara (Tohoku University), Alberto Apostolico (Geogia Institute of Technology and University of Padova), Massimo Melucci (University of Padova), Tobias Scheffer (Max Planck Institute for Computer Science), Ken Satoh (National Institute of Informatics), Nada Lavrac (Jozef Stefan Institute), Ljupco Todorovski (University of Ljubljana), and Hiroki Arimura (Hokkaido University).

In response to the call for papers 55 manuscripts were submitted. The Program Committee selected for publication 17 submissions as long papers and 10 submissions as regular papers. Each submission was reviewed by at least two members of the Program Committee, which consisted of international experts in the field. The selection was made after careful evaluation of each paper based on clarity, significance, technical quality, and originality, as well as relevance to the field of discovery science. This volume consists of three parts. The first part contains the papers/abstracts of the invited talks, the second part contains the accepted long papers, and the third part contains the accepted regular papers.

We are deeply indebted to the Program Committee members as well as their subreferees who played the critically important role of reviewing the submitted papers and contributing to the intense discussions which resulted in the selection of the papers published in this volume. Without their enormous effort, ensuring the high quality of the work presented at Discovery Science 2007 would not have been possible. Furthermore, we would like to thank all individuals and

institutions who contributed to the success of the conference: the authors for submitting papers, the invited speakers for their acceptance of the invitation and their stimulating contributions to the conference, the Steering Committee, and the sponsors for their support. In particular, we acknowledge the generous financial support from the Air Force Office of Scientific Research (AFOSR), Asian Office of Aerospace Research and Development (AOARD)[1]; the Graduate School of Information Sciences (GSIS), Tohoku University for providing secretarial assistance and equipment; the Research Institute of Electrical Communication (RIEC), Tohoku University; New Horizons in Computing, MEXT Grant-in-Aid for Scientific Research on Priority Areas; and the Semi-Structured Data Mining Project, MEXT Grant-in-Aid for Specially Promoted Research.

July 2007                                                      Vincent Corruble
                                                              Masayuki Takeda
                                                              Einoshin Suzuki

---

# Conference Organization

## Conference Chair

Ayumi Shinohara      Tohoku University, Japan

## Program Committee

| | |
|---|---|
| Vincent Corruble (Co-chair) | Université Pierre et Marie Curie, Paris, France |
| Masayuki Takeda (Co-chair) | Kyushu University, Japan |
| Jean-Francois Boulicaut | INSA Lyon, France |
| Will Bridewell | CSLI, Stanford, USA |
| Simon Colton | Imperial College London, UK |
| Antoine Cornuejols | Université Paris-Sud, France |
| Andreas Dress | Shanghai Institutes for Biological Sciences, China |
| Saso Dzeroski | Jozef Stefan Institute, Slovenia |
| Tapio Elomaa | Tampere University of Technology, Finland |
| Johannes Fuernkranz | Technical University of Darmstadt, Germany |
| Dragan Gamberger | Rudjer Boskovic Institute, Hungary |
| Ricard Gavalda | Technical University of Catalonia, Spain |
| Gunter Grieser | Technical University of Darmstadt, Germany |
| Fabrice Guillet | Ecole Polytechnique of the University of Nantes, France |
| Mohand-Said Hacid | Université Lyon 1, France |
| Udo Hahn | Jena University, Germany |
| Makoto Haraguchi | Hokkaido University, Japan |
| Tomoyuki Higuchi | The Institute of Statistical Mathematics, Japan |
| Kouichi Hirata | Kyushu Institute of Technology, Japan |
| Tu Bao Ho | Japan Advanced Institute of Science and Technology, Japan |
| Achim Hoffmann | University of New South Wales, Australia |
| Tamás Horváth | Fraunhofer Institute for Intelligent Analysis and Information Systems, Germany |
| Daisuke Ikeda | Kyushu University, Japan |
| Kentaro Inui | Nara Institute of Science and Technology, Japan |
| Szymon Jaroszewicz | National Institute of Telecommunications, Poland |
| Hisashi Kashima | IBM Research, Tokyo Research Laboratory, Japan |
| Kristian Kersting | Universitaet Freiburg, Germany |
| Ross King | University of Wales, UK |
| Andras Kocsor | University of Szeged, Hungary |
| Kevin Korb | Monash University, Australia |
| Stefan Kramer | Technical University of Munich, Germany |
| Nicolas Lachiche | Université de Strasbourg, France |

# Table of Contents

## Regular Papers

# Challenge for Info-plosion

Masaru Kitsuregawa

Institute of Industrial Science, The University of Tokyo
4-6-1 Komaba, Meguro-ku, Tokyo 153-8505, Japan
kitsure@tkl.iis.u-tokyo.ac.jp

**Abstract.** Information created by people has increased rapidly since the
year 2000, and now we are in a time which we could call the "information-
explosion era." The project "Cyber Infrastructure for the Information-
explosion Era" is a six-year project from 2005 to 2010 supported by
Grant-in-Aid for Scientific Research on Priority Areas from the Min-
istry of Education, Culture, Sports, Science and Technology (MEXT)
of Japan. The project aims to establish the following fundamental tech-
nologies in this information-explosion era: novel technologies for efficient
and trustable information retrieval from explosively growing and hetero-
geneous information resources; stable, secure, and scalable information
systems for managing rapid information growth; and information uti-
lization by harmonized human-system interaction. It also aims to design
a social system that cooperates with these technologies. Moreover, it
maintains the synergy of cutting-edge technologies in informatics.

## 1 New IT Infrastructure for the Information Explosion Era

The volume of information generated by mankind has increased exponentially,
i.e., "exploded" since 2000. The purpose of our research project, "Cyber Infras-
tructure for the Info-plosion Era" in the Ministry of Education, Culture, Sports,
Science and Technology (MEXT) Grand-in-Aid for Scientific Research on Prior-
ity Area, is to build advanced IT infrastructure technologies for this information
explosion era [1]. According to the research by the University of California at
Berkeley, the volume of information created by human is explosively increas-
ing [2, 3]. Huge volume of data is also created by sensors and machines. We have
considered that the most important theme for researchers in the field of computer
science is the research on new IT infrastructure for the Information-explosion
Era.

The project has three major research components (Research Groups) to
achieve this goal: build technologies to search for needed information efficiently,
without bias and without being at risk from the rapidly growing volume of in-
formation (A01); build new and sustainable technologies that can operate large-
scale information systems managing enormous amounts of information safely and
securely (A02), and build human-friendly technologies to enable flexible dialogue
between men and machines and enable everyone to utilize information (A03).

Underlying these three components is the research into new social systems that can facilitate the use of advanced, information-based IT services (B01). In addition, Large-scale Info-plosion Platform (LIP) is implemented as the shared platforms used by all Research Groups. The project is interdisciplinary in its structure, bringing together advanced research methods in information-related areas. (Project leader: Masaru Kitsuregawa)

## 2  Infrastructure for Information Management, Fusion and Utilization in the Information Explosion Era (A01)

This Group focuses on the shortfalls of present internet searches and looks into new search methods, including better ranking systems (where minority opinions are not overlooked), interactive searches, reliability assessments and time-space searches. Currently, for knowledge workers, about 30% of their time on intellectual activities is spent just for retrieval information [4]. The Group will attempt to create a system of search platforms to search a massive volume of web page contents.

Another important issue for information retrieval is the dangers associated with information rankings. Search engines are extensively used in the web world. When a general word is given, it may hit millions of pages, and only about 10
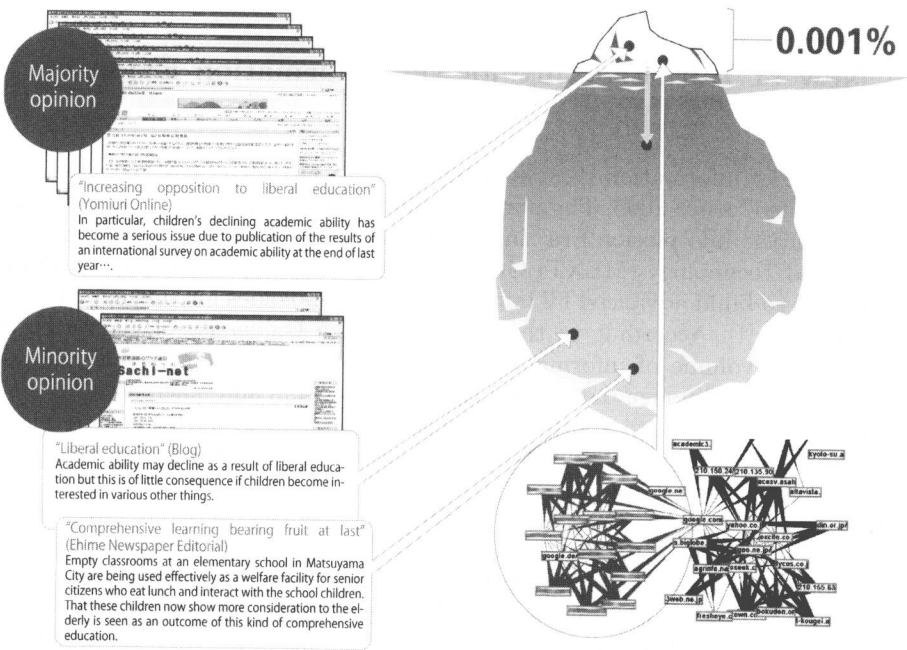


**Fig. 1.** Dangers Associated with Information Rankings

probable candidates are listed as top rankers on the first page of the search result. Now, is this ranking really reasonable? Who guarantee the correctness of it? Actually, this is controlled by just a private company. There should be a possibility that the ranking is sold and bought. In addition, it is possible to control the ranking deliberately. We have found, as shown in the Fig.1, that the red island of links suddenly appear on February 2004, while only the black island of links is found in 2003. This is an example of a trick to raise their own rank by linking them to a well-known site.



**Fig. 2.** Next generation searches

Also, in the current ranking system, majority opinions are highlighted while minority opinions are buried. For example, when "Yutori Kyouiku" (liberal education) is searched, majority opinions are found easily like "increasing opposition to liberal education". Although there are minority opinions such as "comprehensive learning bearing fruit at last", they could be completely neglected.

Various technologies, including information searching, natural language processing, machine learning, artificial intelligence (AI) and database technology, are integrated into this system to enable quantitative analysis to be performed (Fig.2). We hope to enable a remarkable level of interdisciplinary synergy among different fields. It is better for human to understand by presenting relational information also than information itself. Several innovative researches are expected in this area, for example, presenting comparative information [5] and information

on time sequence [6]. Research on the methods of management, integration and processing of "exploding" real-world information, including cyber information and information obtained from remote sensors, will also be conducted. (Group leader: Masaru Kitsuregawa, University of Tokyo)

## 3   Infrastructure for Information-Explosion-Proof IT Systems (A02)

The exponential increase in the amount of information requires far larger IT systems to handle the volume. According to [7], Google has now over 450,000 servers over 25 locations around the world. The ratio of computers used for search engines among shipped computers is 5% in total, according to MSRA (Microsoft Research Asia) Summit in 2006. Data on the Internet is dispersed over millions of nodes, making the overall system unstable and vulnerable to information overload.

In order to keep stable operation of such huge systems, real-time monitoring of behavior of software is inevitable. For such a purpose, explosive volume of information extracted by software sensors should be analyzed so as to point out anomaly behavior of the system and stabilize it. Researches on mining huge volume of data yielded from monitoring very-large-scale systems are particularly important for the age in which a nation-wide cyber attack becomes reality like Estonian case.

This Group aims to establish a new "resilient grid" infrastructure which can automatically allocate computer resources, handle large-scale system faults over
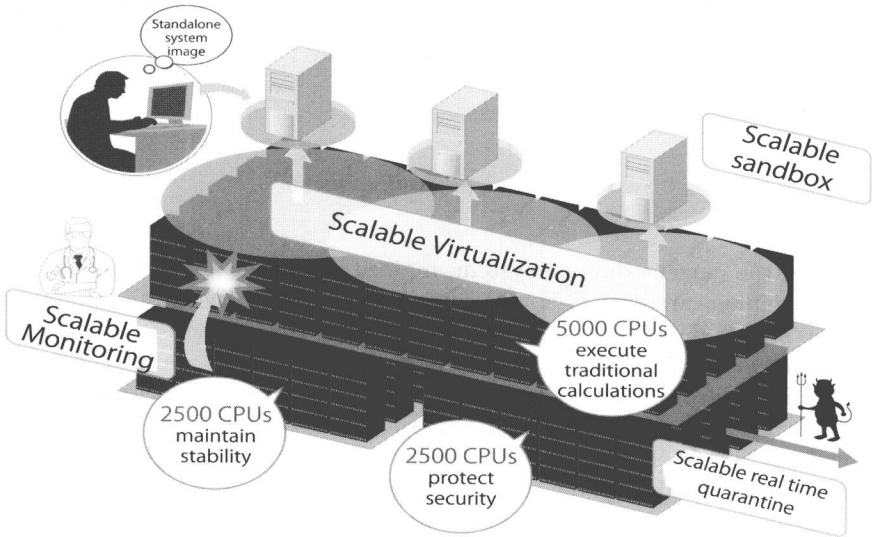


**Fig. 3.** Infrastructure for large scale system

the network without human intervention and without the modern-day concerns of security breaches and intrusion (Fig.3). This is partially based on Autonomic Computing [8] and interesting research works performed by UC Berkeley and Stanford University [9]. The resilient grid will allow a high-performance virtual computing environment to be configured autonomously, on which applications can be deployed safely and securely. (Group leader: Satoshi Matsuoka, Tokyo Institute of Technology)

## 4  Infrastructure for Human Communication in the Information Explosion Era (A03)

The information explosion has two aspects: qualitative (volume) and quantitative (complexity). This Group proposes studies of the advancement of human communication to address the issues related to complexity. The underlying concept is a mutually adaptable multi-modal interaction that can fill the communication gaps between people and information systems. This is the key to overcoming the complexity resulting from highly functional and multi-functional information systems, and establishing a secure and user-friendly interactive environment (Fig.4).

Searching information from explosively huge size of information space still requires advanced skills, since existing tools for such a purpose is not necessarily easy to use for naive users. Human-friendly interfaces as well as communications
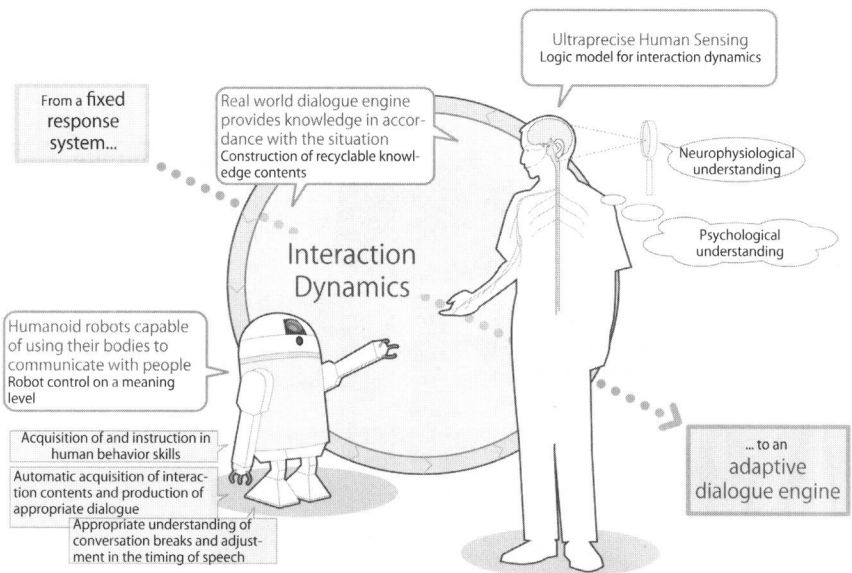


**Fig. 4.** Infrastructure for human