# Journal of Information Science

**Paradigms for abstracting systems**

María Pinto and Carmen Gálvez

The online version of this article can be found at:

Published by:

$SAGE Publications

http://www.sagepublications.com

On behalf of:

cilip

Chartered Institute of Library and Information Professionals

**Additional services and information for *Journal of Information Science* can be found at:**

**Email Alerts:** http://jis.sagepub.com/cgi/alerts

**Subscriptions:** http://jis.sagepub.com/subscriptions

**Reprints:** http://www.sagepub.com/journalsReprints.nav

**Permissions:** http://www.sagepub.com/journalsPermissions.nav

# Paradigms for abstracting systems

**María Pinto and Carmen Gálvez**

*University of Granada, Spain*

**Abstract.**

**The ever-growing amount of documents (electronic or other) has increased the value of abstracts, and abstracting systems and services, as instruments for concentrating and supplying relevant information. Ongoing research in different disciplines, with abstracting as the common subject matter, points to the usefulness of the paradigm concept. This perspective would provide the coherence necessary not only for the processes of abstracting but also for the products and services derived. We identify and describe four paradigms (communicational, physical, cognitive and systemic), comprising the most significant lines of investigation in this area, in order to offer a 'state of the art' analysis that may prove helpful for future research studies.**

## 1. Introduction

The abstract (produced by machine or human), as a concise statement of the central message of a document, has become an increasingly important tool for distinguishing truly relevant information from the bulk of information available. The importance of the abstract is increasing in all fields of study, because of:

(1) the increasing volume of machine-readable text: electronic documents entering the Internet every day, plus information retrieval (IR) systems that use more and more full-text documents;

(2) advances in natural language processing (NLP);

(3) the progress in automatic systems of abstracting, giving rise to new operational methods and models of the processes involved.

New approaches to abstracting are justified by the impossibility of simplifying operations within the systems of knowledge representation and IR in which they are carried out. These operations are not merely a set of mechanistic processes to be imparted without consideration for the greater, more complex functions through which knowledge is represented and transmitted. Consequently, there are many problems involved in abstracting. Most frequently, these difficulties have to do with comprehension of natural language, semantic representation, speech models or with knowledge of the world at a given time. However, there are other areas and matters involved as well: linguistics, logic, statistics, psychology and artificial intelligence (AI) all play some role in information processing (IP). In view of this complexity, the aim of our study is to provide a certain degree of organisation and coherence to the most representative accomplishments in this area to date, under the epistemological shelter offered by scientific paradigms.

To avoid confusion, it is necessary to clearly define the limits of the conceptual topic of departure: the abstracting system (AS) that we describe is the result of the general processes of abstracting within the context of a database or an information storage and retrieval system and, more specifically, it would constitute a specialised feature of knowledge-based information systems. The pragmatic dimension of this contextual factor is the origin of the current indexing and abstracting services. They feature specific applications of human and automatic abstracting methods and are conditioned by:

(1) input factors (type of documents or source forms);

(2) the processes involved (statistical, linguistic or cognitive);

(3) the function and purpose of the abstracts (i.e. type of users to which they are directed); and

---

*Correspondence to:* Professor M. Pinto, Department of Documentation, 18071 University of Granada, Granada, Spain. E-mail: mpinto@platon.ugr.es

(4) output factors (type of format and style of the abstracts).

The intentions of our study are threefold:

(1) to identify, from an appropriate perspective, the problems that arise in this complex system;

(2) to discover the properties of dynamic interaction that configure ASs as a whole; and

(3) to integrate the various lines of research into a unitary framework.

The scientific state of ASs can be perceived clearly only if we first identify the systematic research in the area of knowledge to be covered. This topic is tied in with transdisciplinary and pluridisciplinary structures as well as interdisciplinary ones; therefore, it is necessary to find criteria that will mark the boundaries of our subject matter. This paper not only analyses the paradigmatic level as a 'vertical' influence on ASs, it also describes the disciplinary convergence that functions syntagmatically as a 'horizontal' influence and harmonises within this subject as a single entity (as illustrated in Figs. 1–7). At the same time, the proposed selection of paradigms conforms to the criterion of relevancy in solving problems within our field of research and is not meant to be exclusive.

## 2. Basic paradigms in abstracting systems

As stated above, the subject content of ASs is included within the framework of information science (IS), an interdisciplinary domain that comprises all the areas of knowledge integrated in information studies, from the generation of information to its transmission to potential users through different channels. This particularity implies a complex relationship between IS and other disciplines that are much more consolidated. At first glance, the disciplinary overlapping seems to be a significant obstacle. However, it achieves unity under Kuhn's concept of paradigm [1], which refers to the whole of scientific performances providing research models, methods and goals to the scientific community. Paradigms are characterised by the homogeneity and coherence of their scientific postulates, enabling us to:

(1) formulate a *theory* that can serve as a basic frame of reference for a specific scientific community;

(2) respond to problems by means of an appropriate *methodology;* and

(3) set specific goals to be transformed into the *subject of research.*

From this point of view, there are four paradigms which comprise the epistemological and methodological postulates and regulate the criteria for ASs (see Fig. 1):
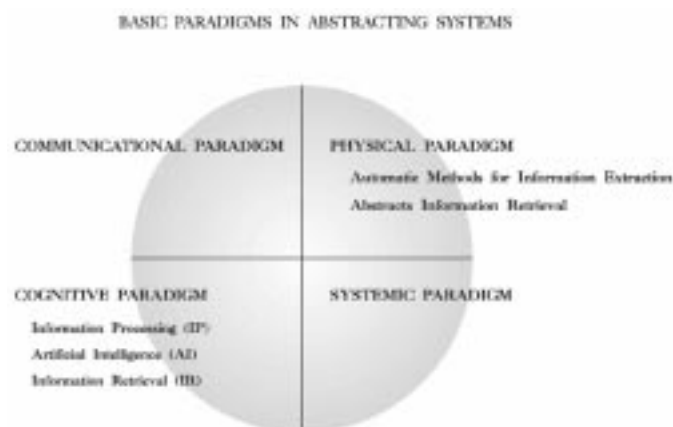


Fig. 1. Basic paradigms in abstracting systems.

- Communicational paradigm;
- Physical paradigm:
  – Automatic methods for information extraction and abstracting;
  – Information retrieval through abstracts;
- Cognitive paradigm:
  – Information processing;
  – Artificial intelligence;
  – Information retrieval;
- Systemic paradigm:
  – Quality management.

An eclectic attitude is needed to consider these as complementary paradigmatic contributions that are not incompatible. The present analysis focuses on the main ideas of each of the above-mentioned paradigms which, in our opinion, have a 'vertical' influence on ASs.

## 3. Communicational paradigm

The general theory concerning this paradigm, proposed by Shannon and Weaver [2], was initially called *The Mathematical Theory of Communication* and is now universally known as *information theory.* Its variables are used to measure and verify optimal conditions for data transmission and refer to the determination of:

(1) the amount of information that data transmission may contain;

(2) the communication channel or network by which more information may flow more quickly and to a greater number of users;

(3) the type of coding – the signal sequence organisation which allows discrimination of a larger variety of messages in a simpler way, without any ambiguity; and

(4) the effects on decoding due to disturbances (noise) occurring during transmission.

Under this paradigm, communication problems are tackled formally (the phenomena observed experimentally are described and explained through logical-mathematical laws) and general conditions for data transmission (seen as signal sequences) are determined, regardless of the message conveyed.

The methodology for the study of these aspects of communication is procured considering:

(1) a general model for representing the communicative flow of information;
(2) a general scale for calculating the information (number of messages) in each of the communicative transmission points; and
(3) a comparison between the amount of information expected and that corroborated in order to verify transmission fidelity.

Although information theory does not suffice to explain the complexity of communication carried out by people as compared with machines, it is indeed very effective when explaining aspects related to signal transmissions: protocols, interfaces or networks. Limiting information theory to a theory of message transmission, the abstracting system/service as a node in the data transmission network is sketched, including problems related with AS design and the variables that enhance the successful transfer of information (see Fig. 2). These services should be adapted to the changes introduced by transmission and telecommunications technology. For example, the paradigm will analyse problems arising from the eventual elimination of secondary services because of online full-text electronic dissemination [3].

The variables selected as relevant facts for information transfer are analysed within the content of ASs and synthesised as follows.

(1) *Original source*: printed documents; electronic documents (full text, electronic journals, websites, etc), accessibility; language or code for information storage, type (periodical printed publications, dissertations, periodical electronic publications, patents, etc).
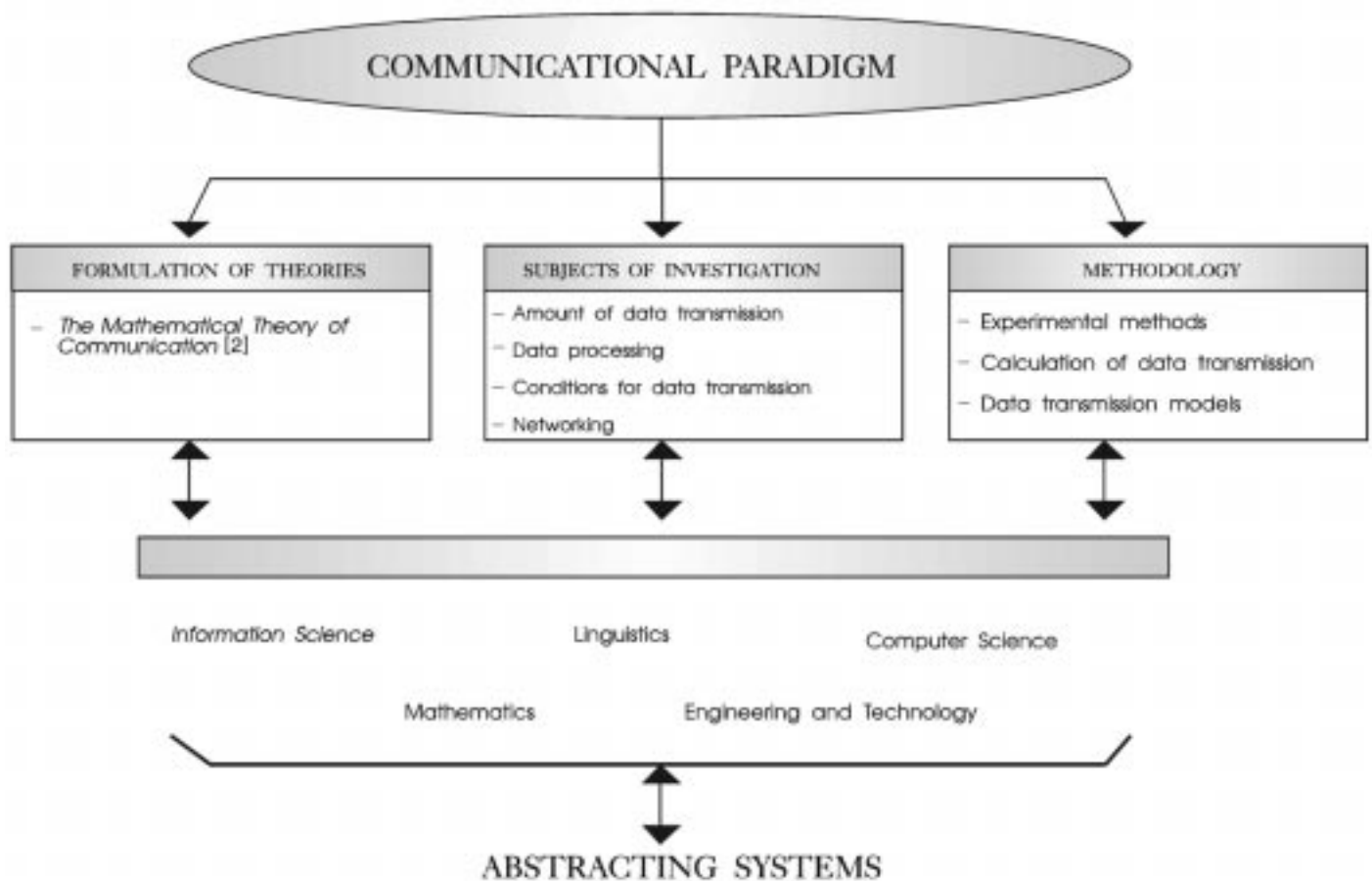


Fig. 2. Communicational paradigm.

(2) *Recipient*: type; need for information; location, etc.
(3) *Media*: source storage media; source transmission media; recipient's reception media; recipient's storage media.
(4) *Social context of information transfer*: need for information, information distribution channels, effects on the audience, national and foreign information policies.
(5) *Quality*: information selection quality; information presentation quality; information impact.

Most research under this paradigm is concerned with the transmission of abstracts; at one time, it was even limited to the evaluation of abstracting and indexing secondary services, as shown by Edwards [4], Gilchrist [5] and Lancaster [6]. Until now, secondary services have been disseminated through printed or online databases (LISA, ISI, ERIC). However, these conventional transmission media have now been surpassed by Internet subject trees (Excite, Lycos, Top 5%, Yahoo!, Infoseek) and gateways (OMNI, ADAM, EEVL). Wheatley and Armstrong have evaluated the new role of abstracts in a networked environment. According to them, an analogy exists between conventional and current services: 'Web search engines and catalogues exist in much the same positions as conventional abstracting and indexing services to conventional paper documents' [7, p. 207]. However, they also found some differences, since online database abstracts (ERIC, LISA) have 'a filtering out role in which users scan through the abstracts in a broadly appropriate answer set in order to reduce it to ideal records', while Internet subject trees and gateway abstracts 'are more active in the area of information discovery or filtering in'.

The key role of the specific social context in which relations among sources, recipients and channels take place must also be stressed. A contextual paradigm which would make it possible to analyse the different audiences' needs and their effects on abstracts could even be mentioned. Despite their significance, however, the ramifications of social context lie beyond the scope of the present review. Social context, including its political and economic effects, has vague boundaries and a complex structure that generally defy definition or general consensus. For the purposes of this paper, therefore, instead of suggesting a social sub-paradigm within the communicational paradigm, *context* will be considered as a continuous variable which plays either a main or a secondary role in each paradigm.

The methodology of the communicational paradigm is very useful for tackling specific problems involving coding, data transfer by means of electronic circuits and measurement of results through a binary signal code, the amount of information calculated in bits. Nonetheless, the strict application of this theory proves ineffective, because communication among agents of documentary information systems depends not only on signals but also on multiple 'chances' concerning linguistic and cognitive structures that are difficult to pin down, even outside this paradigmatic framework.

## 4. Physical paradigm

The main goal of ASs is to facilitate and enhance IR, allowing correct access to items of information [8]. Research on IR is commonly conducted from two different perspectives: the physical paradigm (or the Cranfield paradigm) and the cognitive paradigm [9].

The origins of IR research under the physical paradigm (see Fig. 3) can be traced to 1953, when separate groups of tests were carried out in Britain and the USA to evaluate the performance of the then controversial 'Uniterm' system of Mortimer Taube against more conventional approaches to subject indexing and retrieval (the Cranfield-Uniterm test of the UK and the Armed Services Technical Information Agency (ASTIA)-Uniterm test carried out in the USA). These tests were influential as archetypes for the later series done at, or in association with, the Cranfield Institute of Technology. The Cranfield tests, in turn, signalled the real beginning of IR research as an empirical discipline: they 'established the principle that arguments about the relative merits of different retrieval systems had to be empirically grounded, and, in this respect, they mark a historical change in consciousness from a philosophical and speculative approach to an experimental and empirical one' [10, p. 50].

Many problems concerning IR under the physical paradigm are related to information representation, storage and access (considering information as a physical entity that is measurable and quantifiable). The problems outlined in this paradigm are so wide that it is very difficult to synthesise them. They would encompass, among others, storage media and techniques, file layout, searching strategies, automatic techniques for information representation and extraction, advanced retrieval methods and retrieval evaluation systems.

The methodology applied in this approach towards IR systems is fundamentally empirical. Consequently, it is analogous to a physical system – mechanical and physical engineering systems – in which objects and facts extracted from reality are subject to specific experimental tests in an artificial environment. Once their effectiveness is demonstrated and evaluated, the same
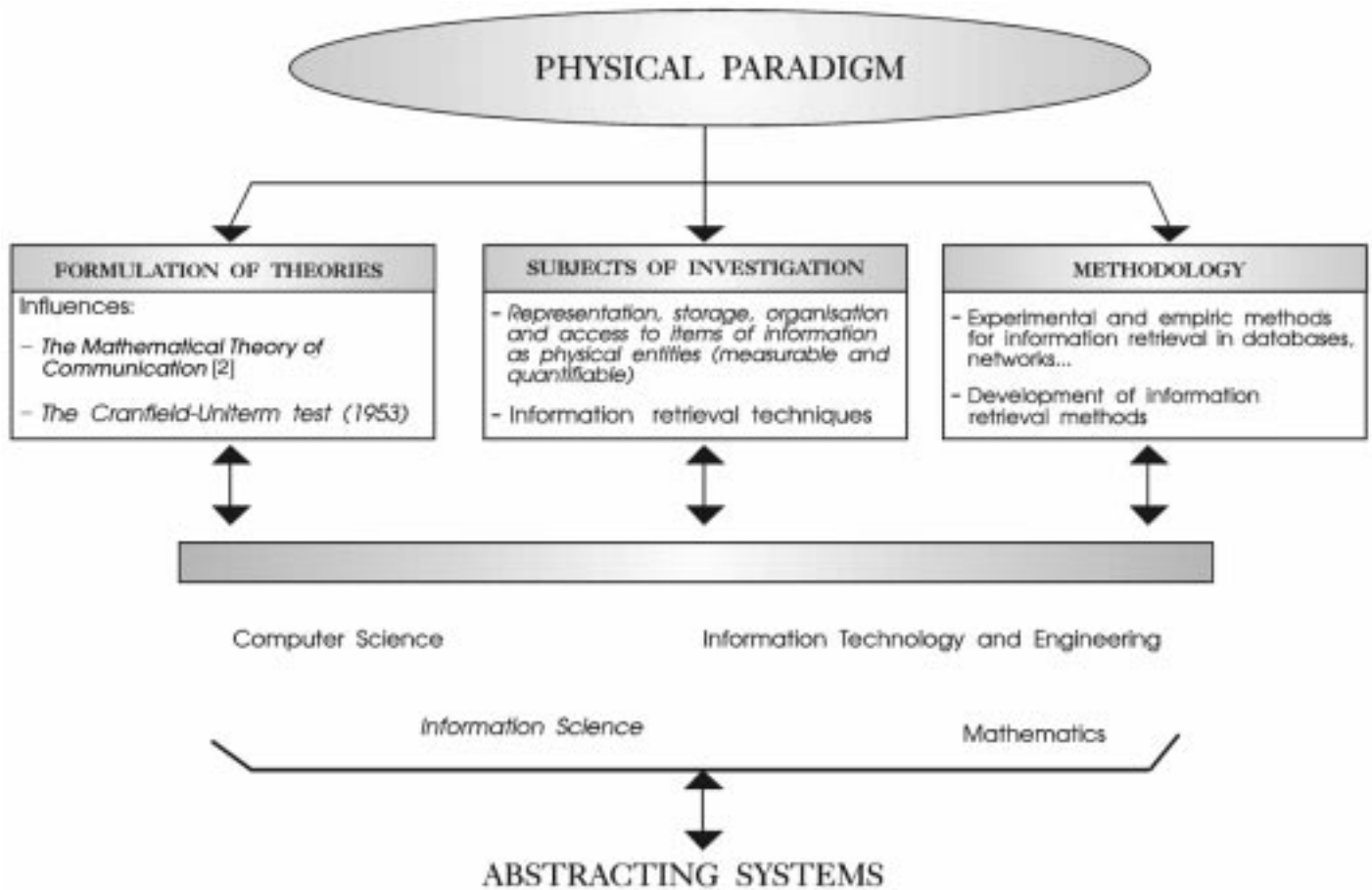
Fig. 3. Physical paradigm.

method is applied to real units. In this sense, an IR system and a physical system are similar not only in their nature but also in the experimental techniques appropriate for their study [9]. With respect to ASs, the utility of this methodology can be taken in two directions: (i) automatic methods for information extraction and abstracting and (ii) IR through abstracts.

### 4.1. Automatic methods for information extraction and abstracting

Most automatic methods for information extraction and abstracting consist of assigning a value or weight to each sentence from the source document according to the type of words and phrases included. Sentences having the highest values will be used to construct the abstract. These methods distinguish between 'words occurring in the text body, and words occurring in titles, captions and section headings' [10, p. 439]. Firstly, the system identifies the words and then calculates a term weight on the basis of word occurrence. The next step is the construction of phrases using the words that concur in the sentences, to which specific weight will also be assigned. Once the weight is calculated for terms and phrases, the sentence is constructed, taking these values into account. Finally, sentences are extracted for abstract composition. According to Paice [11], the systems using this method are those of Luhn [12], Oswald *et al.* [13] and Edmundson [14].

There are other methods which, besides calculating word frequency, account for additional factors such as sentence position within the text, word and phrase indicators of the more relevant sentences, and syntactic criteria that provide coherence to the text. Studies based on these systems are those of Baxendale [15], Rush *et al.* [16], Paice [17] and Black and Johnson [18].

Yet, these automatic methods for information extraction are plagued by a lack of coherence and balance, calling for the development of alternative systems that would solve or by-pass these problems, as well as deal with anaphoric relations and text structures. New research and systems came along to address these

difficulties [11]: ADAM Automatic Abstracting Program, proposed by Rush *et al.* [16]; studies of anaphora in scientific abstracts, carried out at Syracuse University, New York, by Liddy *et al.* [19]; GARP (Gareth's Anaphor Recognition Program), proposed by Gareth Husk [20] and developed at Lancaster University, UK.

All these methods are still deficient, however, because they do not provide for the coherence and balance needed to construct quality abstracts. Thus, improvement in this field became dependent on the existence of a satisfactory text structure theory, which led to the following contributions: McKeown's TEXT System [21], which classifies textual sentences and creates informative paragraphs according to a specific scheme; research on textual structures, such as that done by Meyer [22]; theories of text macrostructures, microstructures and superstructures [23] and studies of the function of abstract structures [24].

In addition, changes were brought about by the impact of networked electronic documents, leading to the reformulation of automatic methods for information extraction. Web search engines generate automatic text extraction, which would be equivalent to conventional abstracts. Abstracts from new sources of electronic information, consequently, will consist of a set of 'microtexts' extracted from source documents [7, p. 207], such as those developed in 'AltaVista's simplistic sampling of the first parts of web pages; Infoseek's extraction of text from the body of pages; or Yahoo's one-line or two-line characterisations of extensive web sites'.

These developments do not eliminate the usual problems caused by automatic methods for information extraction. On the contrary, they reproduce problematic structures more closely, with a greater lack of cohesion in paragraphs and more incoherent structures in extracted sentences. A more recent attempt to solve this old problem consists of generating coherent text segments by taking advantage of the convenient way in which automatic hypertext links are created at present. Salton *et al.* [25, p. 195] used IR techniques to generate automatic hypertext links for their use in automatic text summarisation. Thus, rather than creating inter-document links among various documents, they used 'automatic link generation techniques to generate intra-document links' among various paragraphs or sentences in a text. Afterwards, these paragraphs and inter-document links are located in a text relationship map in order to visualise the text structure and then to produce text summaries by passage extraction.

These research studies all emphasise a feature inherent to structured texts (abstracts and summaries included): elements originating in different levels are joined together in a complex way by an internal organisation which ensures cohesion and facilitates discourse progression by giving more information. On the other hand, text structures are not isolated but rather linked to different social contexts that lend meaning to them.

The type of discourse is determined by the pragmatic-sociological conditions of the different interactive situations in which texts are created and processed. Thus, the structure of a text is revealed by the exposition of an organised plan, which involves overall schemes and different representational patterns of social knowledge to ensure interconnection among all the parts.

### 4.2. Information retrieval through abstracts

Another important area of research is IR from databases, departing from the various fields that integrate a bibliographical record. Within this framework, free-text retrieval (or differently originating controlled terms retrieval) from the varied database fields has motivated a wide range of empirical studies, the most relevant of which are IR tests in the ERIC (Educational Resources Information Center) database, by Markey *et al.* [26], and IR evaluation, with the comparison of its effectiveness depending on the use of the full text, abstract or different controlled terms. Within this second group, the most important studies are those by Tenopir [27] and Ro [28], using the *Harvard Business Review* database, and the evaluation of retrieval effectiveness using the STAIRS system [29].

The role of abstracts in networked IR has begun to be evaluated on the basis of the alternatives proposed by metadata. First of all, 'within the body of a Web page, there is no provision for fielded data and so it is not possible to search, as one might in the conventional database, for title words, indexing terms, publication date, author or corporate source' [7, p. 206] and, consequently, there is also no provision for a field called 'abstract'. The HTML (hypertext mark-up language) format offers a solution to this problem by means of a <meta> tag for describing resources in Web documents. Metadata, i.e. data about data, located at the heading of a document could contain an attribute-value pair, with room for an abstract. This idea suggests the creation of a new abstracting model for IR in a networked environment from <meta> tags. Again according to Wheatley and Armstrong [7, p. 212], 'an ideal Internet abstract might include user guidance, assessment of authority, discussion of physical attributes (the design

of the site or the ease of navigation), judgements on quality, or pointers to alternative sources'.

## 5. Cognitive paradigm

This approach is influenced by *The Mathematical Theory of Communication* [2], especially in relation to the concept of information: an item that can be quantified as a set of units for its appropriate transmission. This idea was incorporated into cognitive processes and interpreted as an analogy of computer-assisted transmission systems. The basic theories of this paradigm may be synthesised as follows.

(1) Development of cybernetics by Wiener [30] as a general theory for computer systems by analogy to the human mind: the evolution of this theory implies the birth of the *connectionist paradigm*, related to the simulation of neuronal networks.

(2) Influence of the generative-transformational grammar by Chomsky [31], highlighting the form in which individuals structure and create language: this concept goes beyond the description of language studies and takes on more flexible linguistic structures, to be represented through a generative process with a set of finite rules as the starting point.

(3) Influence of the theory of algorithms, a prescribed set of well-defined rules for solving problems pertaining to the models that explain complex cognitive processes.

Within the cognitive paradigm, there is no general model valid for our documentary approach that satisfactorily explains how the documentalist/analyst processes information or how human knowledge is represented for the purpose of processing information. The lack of such a model does not allow one to identify a user's *cognitive state* with regard to his or her information requirements and needs. We will therefore use this paradigm to analyse the influence of three important sub-paradigms of our informative-documentary context: information processing, artificial intelligence and information retrieval.

### 5.1. Information processing

Previous influences have not only provided a new way of solving problems related to the description of human cognitive processes using computers as models; conversely, they have also given rise to the development of computer programs with human cognitive processes as models. This analogy gives rise to the IP sub-paradigm (see Fig. 4), which uses a research methodology based

on the simulation of computational formalisms and combines computer theory and technology in order to create processing models. Studies in this field include those by Minsky [32], Anderson [33], Pylyshyn [34], Fodor [35] and Rumelhart and McClelland [36]. In general terms, this methodology is based on the potential utility of the computer analogy in generating functional cognitive models.

The main interest of this paradigm for ASs lies in the construction of models to explain how information is comprehended, represented and synthesised; a methodological framework with considerable impact on the following variables:

(1) **comprehension models**: Winograd [37], Van Dijk and Kintsch [38];

(2) **reading/processing**:
   (a) *serial*: Forster [39];
   (b) *parallel*: Marslen-Wilson [40], Stemberger [41];
   (c) *ascendant*: Laberge and Samuels [42];
   (d) *descendant*: Goodman [43];
   (e) *interactive*: Rumelhart [44];
   (f) *modular*: Fodor [35], Anderson [45];

(3) **information storing**: Brown [46], Peterson [47], Atkinson and Shiffrin [48]; and

(4) **information production models**: Van Dijk and Kintsch [38], Flower and Hayes [49], Pinto [50], Endres-Niggemeyer [51], Pinto and Gálvez [52].

Some of the above models refer to how the social or situational context is joined to text processing and how it influences such a process. However, the great variety of dimensions and levels adopted have often rendered this variable meaningless or have caused the elimination of social context in an artificial and deliberate way. Obviously, social context cannot be processed through the application of algorithms or heuristic modelling. Even so, it is generally accepted that the situation interacts with all the agents involved in IP, the linguistic system and documents, although the way in which a given document is processed, observed and represented cannot be clearly stated. Most of these models, therefore, mention the application of scripts, plans and frames to communicative situations or the use of *situation models* regarding real contexts of IP.

### 5.2. Artificial intelligence

AI dates back to the mid-1950s, when Dartmouth held a constitutional conference (see Fig. 5) on new perspectives in computer research, with the aim of creating a multi-use machine able to perform intelligent actions. Primarily, AI could be described as an attempt to create computer programs that would accomplish tasks by
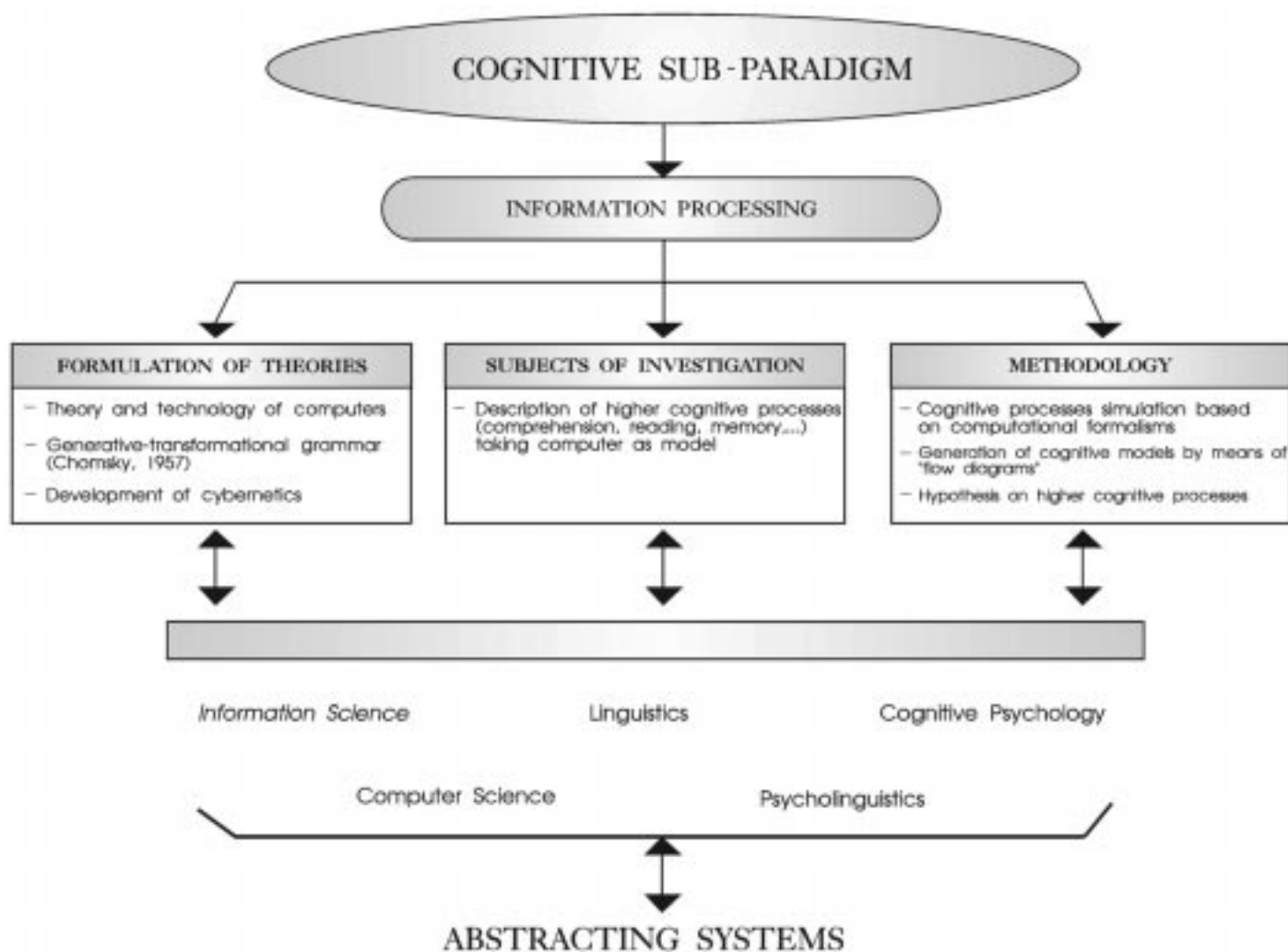
Fig. 4.  Information processing sub-paradigm.

means of processes similar to those carried out by the human mind. Although this notion has since been taken in numerous directions, the original idea remains: to understand the nature of intelligence on the one hand and, on the other, to create tools for simulating cognitive processes (perception, reasoning, comprehension, learning, etc).

The subjects of investigation handled by AI are very complex and especially concerned with the simulation of human reasoning by means of logical systems within computer programs. However, computers need to have specific information or 'base knowledge' in order to simulate reasoning. AI does not work with algorithms; instead, it deals with the problems associated with representations of knowledge. Along these lines, various scientific explanations of how knowledge is

represented have been developed: *schemes:* Rumelhart and Ortony [53]; *theory of scripts*: Schank and Abelson [54]; *semantic networks*: Quillian [55] and *frameworks*: Minsky [32]. To emulate reasoning, computers need not only a symbolic representation of knowledge, but also a mechanism that operates with such a representation, i.e. an inference strategy. These ideas led to the development of the so-called 'expert systems', which emulate the specific sort of knowledge managed by experts (human specialists) when solving problems in a particular area.

The domain of AI includes methodology based on heuristic approaches for knowledge representation, inference methods operating on a knowledge base, design of user interfaces, *ad hoc* abstracting methods for texts, and the development of communication
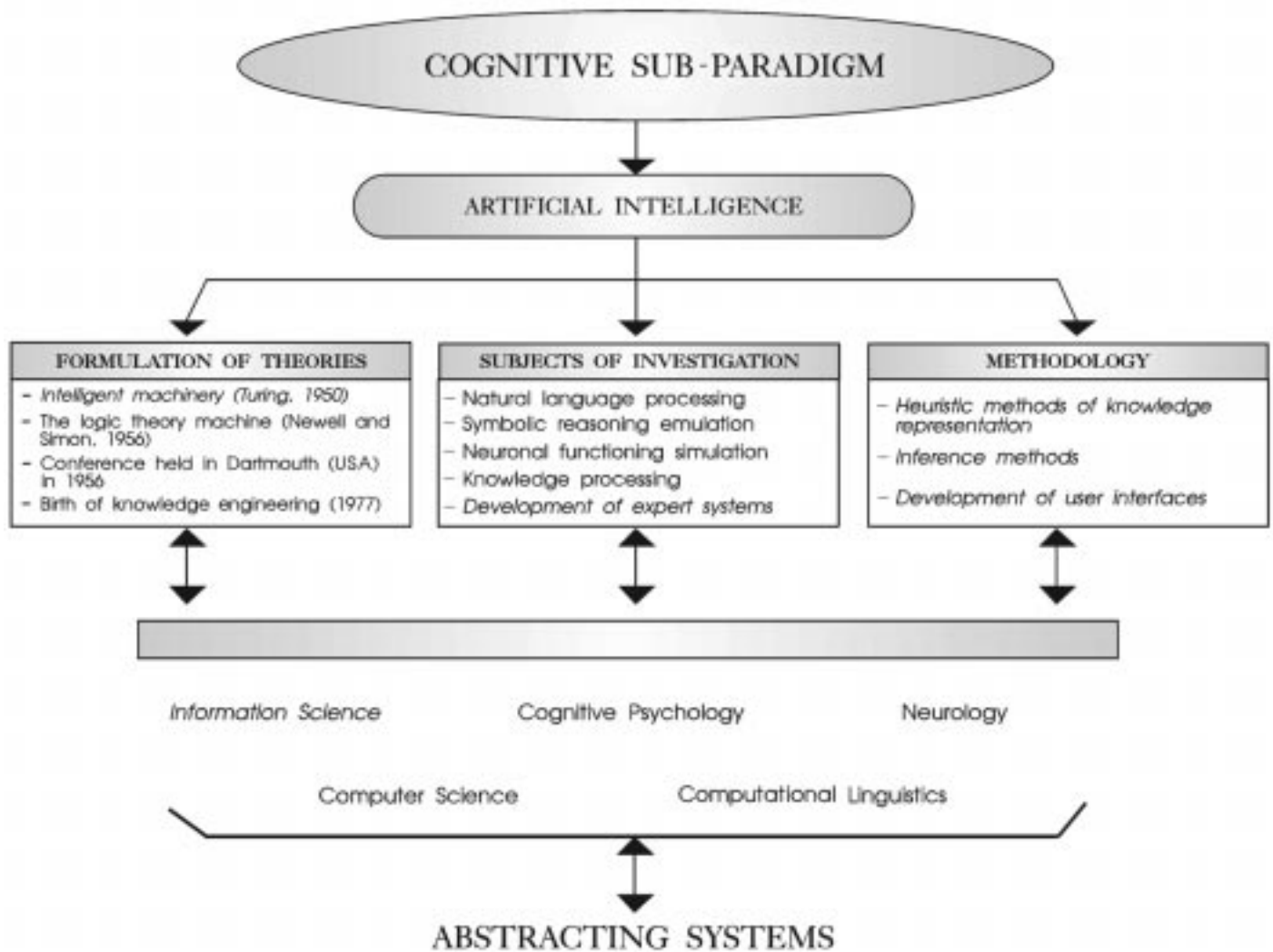
Fig. 5. Artificial intelligence sub-paradigm.

models that incorporate knowledge schemes related to the communicative situation. The impact of these methods on automatic abstracting within the framework of AI has not lived up to expectations. Moreover, studies on this subject are restricted to specific fields that deal with 'detailed semantic information' of well-known characteristics [11, p. 172]. This is the approach adopted by the following systems:

(1) DeJong's FRUMP system [56], which analyses news articles by means of slots in one of a set of pre-defined frames. When the analysis is complete, a script is used to generate a summary of the information held in the relevant frame;

(2) Rau's SCISOR (System for Conceptual Information Summarisation Organisation and Retrieval) [57]: the detailed linguistic analysis of a text (or, indeed,

of several interrelated texts) results in the construction of a semantic graph, which is convenient for intermediate storage. A natural language generator may then produce summaries from the stored material;

(3) Hahn and Reimer's TOPIC (Text-Oriented Procedures for Information Management and Condensation of Expository Texts) system [58]: a project designed to summarise texts about microprocessor systems.

Another important area of research includes the problems observed in NLP: text comprehension, automatic translation, parsing (syntactic analysis), the development of interfaces for improving the relation between users and IR systems, and the creation of documents in natural language. All these processes rely on

sets of grammatical rules, especially semantic ones, in order to analyse expressions according to logical principles or semantic networks. They also require an internal representation of the world or knowledge base (formalised as a set of rules for a specific field). This approach is the basis of important studies, some of which have been pointed out by Lancaster [3]:

(1) Fum *et al.* model [59], which integrates a system based on weighting and parsing procedures for automatic abstracting;

(2) Hahn and Reimer's model [60], based on the application of knowledge and parsing structures for the creation of what they call 'textual condensation';

(3) Grishman *et al.* PROTEUS (PROtotype TExt Understanding System) [61], which analyses real texts and builds a structured abstract from the information contained therein (in the form of a database file). PROTEUS makes a full syntactic analysis of each sentence, providing a regularised structure that serves as the starting point for semantic analysis and reduces the message to a thematic structure. Finally, the pattern generator transfers the interpretation to the different database fields [62].

Meanwhile, AI continues in its attempts to emulate symbolic reasoning as performed by neuronal networks, i.e. by means of new programming forms based on neuronal functioning simulation. Despite the great impact of this new branch of research, known as the *connectionist paradigm*, it is not yet consolidated with respect to information extracting and abstracting methods.

### 5.3 Information retrieval

From an epistemological viewpoint, this sub-paradigm affirms that 'any processing of information, whether perceptual or symbolic, is mediated by a system of categories of concepts which, for the information-processing device, are a model of its world' [63, p. 48]. This system consists of knowledge or cognitive structures as determined by individuals and their environment. For IS, 'taking the cognitive view has typically meant considering its scope as being concerned with some sort of human communication system, in which texts play a key role, and of individuals within that system in their interactions with texts (or information), and with one another in relation to such texts' [64, p. 12].

In this context, it is suggested that knowledge structures are determined by individuals (system users). The main consequence deriving from this approach is an IR system design simulating or constructing an idea of the potential user's needs. As a result, the problems presented by IR under the cognitive paradigm are very

closely related to the user (see Fig. 6). Basic areas of research focus on the user's information needs, the variations in these needs as a result of interaction between the user and the retrieval system, the design of interfaces facilitating the user-IR system relationship and, in general, the development of search tools to ensure successful IR.

Methodologies adopting the cognitive viewpoint in IR are synthesised by Daniels [65] in three groups, which comprise the representation of:

(1) *users and their problems*, which stems from the hypothesis proposed by Belkin [64] on the 'anomalous states of knowledge' (ASK), according to which the user searches for information;

(2) *search strategies*, which compile the different ways search strategies and processes are carried out, depending on the variables involved – user, intermediary, IR system. (Ingwersen's research [66, 67] deserves special mention here); and

(3) *documents and information*, which is considered a major goal of current IR research, since it embraces the whole corpus of studies about user models intended to eliminate the intermediary's role in retrieval systems. The aim of this approach is to allow users direct access to the system by means of the representation of documents and intelligent interfaces.

The main research objective in abstract IR is the development of techniques for modelling the cognitive structures of authors, systems designers, abstractors, indexers or users as an interactive part of IR, in order to meet specific needs for information. According to Fidel [68], abstracts of the most widely available bibliographic databases can be searched in the free-text mode, which allows users to search online for the occurrence of any terms they think appropriate. Among the abstractors using cognitive models and specialising in empirical abstracts, Liddy [69] discovered the problem with searching in the free-text mode in most current IR systems: terms and/or phrases are tracked down as isolated fragments, with minimal facilities provided for searching for concepts that occur in particular semantic relations to each other or fulfil special semantic roles in the respective text. This rather unrefined approach for selecting documents in response to a query results in the retrieval of many irrelevant documents. The search terms do, of course, occur in all the retrieved documents, but the roles or functions played by the concepts represented by these terms may not be what was desired by the user**.**

Because IR, under the cognitive paradigm, takes the user into account in a high-priority way, the role of
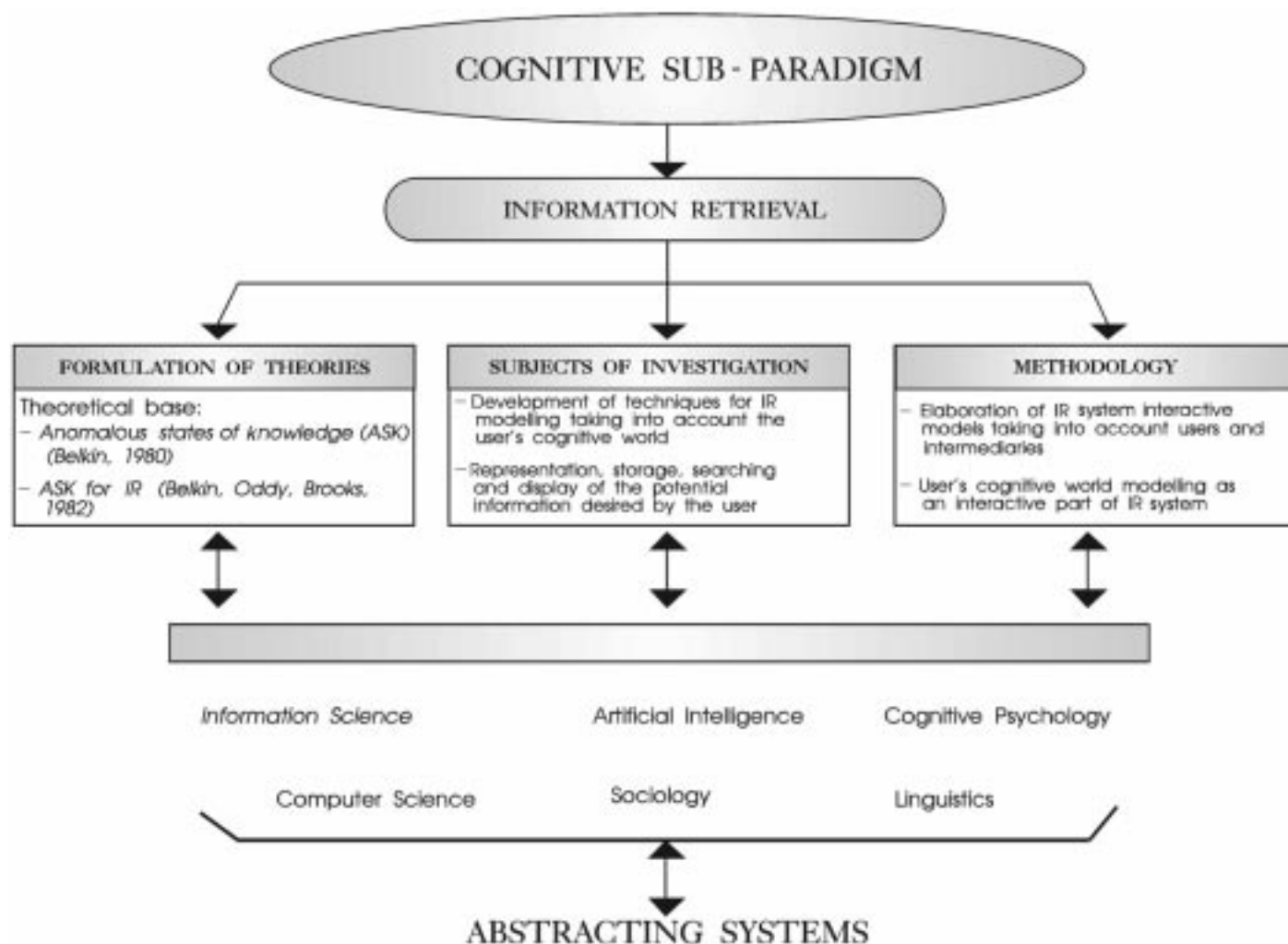
Fig. 6.  Information retrieval sub-paradigm.

interface efficiency is likewise emphasised. Although research studies on the human-computer interface (HCI) and abstracting are scarce, notable contributions to efficiency have been made using *concept-maps* and *structured abstracts* in the visualisation of information. A concept-map presents a graphic representation of the concepts and relationships that make up the textual unit. It is a concrete framework that facilitates not only the codification of information, but also its retrieval. The key advantage of the concept-map is its visual impact: it clearly and concisely shows the relationships between the principal ideas of a text, appealing to the human capacity for obtaining information through visual impressions [52].

The evaluation of structured abstracts, with a modified format that includes sub-headings within the

abstract (such as 'Background', 'Aims', 'Method', 'Results' and 'Conclusions'), has been the focal point of a number of projects, such as the one carried out by Hartley, Sydes and Blurton [70].

## 6. Systemic paradigm

Information was first identified as the result of a process of measurable and quantifiable data transmission (communicational and physical paradigms) and, later, of a process of knowledge transmission (cognitive paradigm). Returning to communicational and physical paradigms, and considering the influence of *The Mathematical Theory of Information* [2], the latter can be understood as something restricted to a system – a

set of organised components (people, procedures and equipment) that work together to achieve the goal of transforming input elements into output elements. In this sense, information is a basic unit of measurement of communicational and physical systems.

The extension of the concept of information beyond the domain of physics not only occurs under the theoretical and methodological proposals of the cognitive paradigm, but also within *General Systems Theory* [71], the foundation of the systemic paradigm. In other paradigms, the problem of communication was the accurate reproduction of a message (a selected message from another document) which, in the case of the communicational and physical paradigms, shares the probabilistic and empirical axiom of what is physically transmitted; or the representation and cognitive processing of information (cognitive paradigm).

General Systems Theory considers the selected message and the reproduced message as the two end-points of a process of interaction. Thus, communication is no longer a question of message selection and reproduction, but rather of the process taking place between a source and a target, i.e. what *interacts* between them.

The general systems theory constitutes the conceptual framework for the development of a large typology of systems: automatic services systems, information storage and retrieval systems, information management systems, computer networks and systems, and others. Our concern here is focused on the *knowledge representation systems*, which take in documents and information needs and, after an *interaction/transformation* process, inform people.

Under this paradigm, any interaction system (organisations, institutions, etc) may be analysed and evaluated informatively. With that same criterion, different participant variables may be planned and designed, in order to achieve the effective management of documents and to fulfil information needs and produce documentary representations.

The methodology needed for the creation of an AS follows the general models elaborated for information storage and retrieval systems: Willitts [72], Yourdon [73], Checkland [74], Bunge [75], Soergel [76] and Meadow [77]. In addition, it takes into account all the processes, variables and components that interact together to obtain an objective (see Fig. 7). An AS requires the installation of a set of applicable development methods to the phases that should be covered, the activities to be accomplished, the products or abstracts to be obtained and the human and automatic techniques to be carried out in each one of the activities designed to produce such products. However, these

integral approximations constitute the weak point of the different scientific contributions.

A well-designed device for information processing and transfer implies economic considerations. To develop such a device, one may resort to modern management methods and to the emergent *paradigm of quality management* (QM), which is based on a single goal (continuous improvement), three principles (user-oriented approach, continuous improvement and full involvement) and six support elements (leadership, education/training, support structure, communications, recognition/reward and measurement). QM key features are (among others) quality improvement as a way of life, collaboration with providers and users, recognition of internal users, identification of key performance indicators, participation of employees, priorities of work groups, elimination of internal barriers and simplification-normalisation of processes and procedures. QM is a systematic and holistic approach to the management problems of an organisation: the whole is greater than the sum of the parts. Though not a goal in itself, QM is a commitment that should be established throughout the organisation for the long term.

User satisfaction should be viewed as the essence of quality, since it is the main goal of any abstracting service that aims to provide quality. This satisfaction is better understood as an emotional reaction to the documentary product/service, where expectations play a crucial role – despite the abundance of empirical investigations surrounding non-emotional measures. Consequently, it seems that quality goals should be objectified and, above all, measurable. Appropriate indicators should be established for the purpose of comparing the quality level of a specific characteristic to predefined standards. Requirements for an information service have to be clearly defined in terms of observable attributes that are subject to the user's evaluation. Qualitative appreciations should be transformed into quantitative material, as improvement can only by achieved by means of numerical values. The greatest difficulty for an AS arises from its cognitive nature, which has a unique means of expression in language. Both language and cognition play a main role in the process; however, they seem to defy quantification. How might we measure the abstractor's level of comprehension of a given text? How can we evaluate whether the interpretation is objective enough? How do we know whether the abstract meets the user's needs? It would take endless questions to address the great variety of components (material, functional, productive, human, automatic, commercial, etc) involved in abstracting.
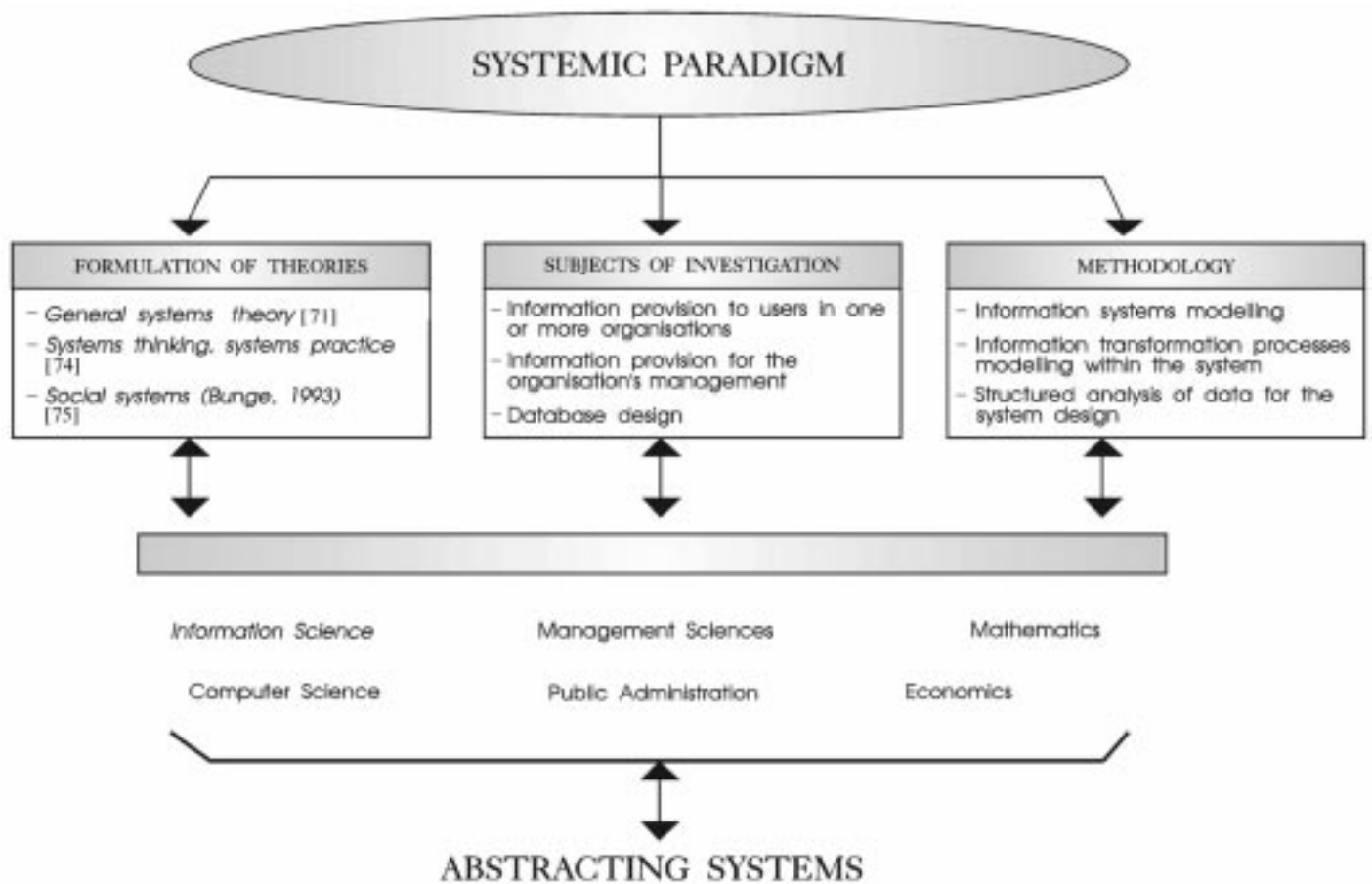
Fig. 7. Systemic paradigm.

ASs need an organisational transformation, which may be undertaken by evaluating activities, analysing deficiencies and establishing *quality programs*. These quality programs, based on a widely accepted rule, include modifications of managerial structures, image changes, the development of new services and, above all, a human commitment to satisfying users' needs. Clearly, the adoption of this paradigm implies transformations at a personal and organisational level. As quality is an issue closely related to orientation, leadership, worker participation and training, its improvement constitutes an endless process which must be carried out step by step and which should not be expected to provide immediate solutions.

## 7. Conclusion

Of all the features that characterise this updated review of ASs, two deserve special consideration: *multi-*

*dimensionality*, meaning structural, functional and procedural complexity, and *interdependence*, or the interaction between the different components, functions and processes. Consequently, the scientific basis of abstracting lies in a theoretical and methodological pluralism, without which a thorough analysis of this complex process is not possible. Research in this field is mainly characterised by two apparently contradictory aspects: epistemological unity, and diversity. In order to unite these opposing terms, the systematic context that justifies them must be accounted for.

The four basic paradigms – communicational, physical, cognitive and systemic – under which abstracting is analysed as a knowledge representation system must be considered as a totality in order to be of use to researchers; even so, this approach is insufficient. Consequently, we must bear in mind that the creation of explanatory models is not an isolated scientific process, but rather a research method at the disposal of other disciplines and thus subject to their developments. For this

reason, apart from considering the multi-paradigmatic integration of different dimensions and interrelations, it is important to establish their epistemological basis; in this way, the process can be influenced, reinforced and improved when a scientific development involving ASs comes along.

This primarily descriptive attempt to systematically reorganise projects in relation to the topic abstract-abstracting, by grouping them in paradigms, may serve as the starting point for further paradigmatic research. A continually deepened and updated analysis of the state of the AS art will, no doubt, contribute to enlightened activity in the design, production and diffusion of that modern-day treasure, the abstract.

With the consolidation of the epistemological basis and the integration of paradigms as our more immediate goals, this paper constitutes an eclectic overview of all the processes involved in ASs. Nevertheless, coordination of all the dimensions analysed implies the reformulation of many of the objectives and methods involved in this process from beginning to end. With this finality in mind, we may conclude that the many processes involved in the systematic creation and dissemination of abstracts possess specific characteristics that should be located within a communicative, linguistic, cognitive and documentary context for their description.

## Acknowledgements

## References

[1] T.S. Kuhn, *The Structure of Scientific Revolutions* (University of Chicago Press, Chicago, IL, 1970).

[2] C.E. Shannon and W. Weaver, *The Mathematical Theory of Communication* (University of Illinois Press, Urbana, IL, 1949) .

[3] F.W. Lancaster, *Indexing and Abstracting in Theory and Practice* (University of Illinois, Graduate School of Library and Information Science, Urbana-Champaign, IL, 1998).

[4] T. Edwards, *A Comparative Analysis of the Major Abstracting and Indexing Services for Library and Information Science* (Unesco, Paris, 1975).

[5] A. Gilchrist, Documentation of documentation: a survey of leading abstracts services in documentation and an identification of key journals, *Aslib Proceedings* 18 (1966) 62–80.

[6] F.W. Lancaster, Some considerations relating to the cost effectiveness of online services in libraries, *Aslib Proceedings* 33 (1981) 10–14.

[7] A. Wheatley and C.J. Armstrong, Metadata, recall, and abstracts: can abstracts ever be reliable indicators of document value? *Aslib Proceedings* 49(8) (1997) 206–213.

[8] G. Salton and M.J. McGill, *Introduction to Modern Information Retrieval* (McGraw-Hill, New York, 1983).

[9] D. Ellis, The physical and cognitive paradigms in information retrieval research, *Journal of Documentation* 48(1) (1992) 45–64.

[10] G. Salton, *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer* (Addison-Wesley, Reading, MA, 1989).

[11] C.D. Paice, Constructing literature abstracts by computer: techniques and prospects, *Information Processing and Management* 26(1) (1990) 171–186.

[12] H.P. Luhn, The automatic creation of literature abstracts, *IBM Journal of Research and Development* 2 (1958) 156–165.

[13] V.A. Oswald *et al.*, *Automatic Indexing and Abstracting of Contents of Documents* (Planning Research Corporation, Los Angeles, CA, 1959).

[14] H.P. Edmundson, New methods in automatic extracting, *Journal of the Association for Computing Machinery* 16 (1969) 264–289.

[15] P.B. Baxendale, Machine-made index for technical literature: an experiment, *IBM Journal of Research and Development* 2 (1958) 354–361.

[16] J.E. Rush, R. Salvador and A. Zamora, Automatic abstracting and indexing. II. Production of indicative abstracts by application of contextual inference and syntactic coherence criteria, *Journal of the American Society for Information Science* 22(4) (1971) 260–274.

[17] C.D. Paice, The automatic generation of literature abstracts: an approach based on the identification of self-indicating phrases. In: R.N. Oddy, S.E. Roberson, C.J. Van Rijsbergen and P.W. Williams (eds), *Information Retrieval Research* (Butterworths, London, 1981), pp. 172–191.

[18] W.J. Black and F.C. Johnson, A practical evaluation of two rule-based automatic abstracting techniques, *Expert Systems for Information Management* 1(3) (1988) 159–177.

[19] E.D. Liddy *et al.*, A study of discourse anaphora in scientific abstracts, *Journal of the American Society for Information Science* 34(4) (1987) 255–261.

[20] G.D. Husk, *Techniques for Automatic Abstraction of Technical Documents Using Reference Resolution and Self-Indication Phrases* (Lancaster University, Lancaster, UK, 1988).

[21] K. McKeown, Discourse strategies for generating natural language text, *Artificial Intelligence* (1985) 1–41.

[22] B.J.F. Meyer, Prose analysis: purpose, procedures and problems. In: B.K. Britton and J.B. Black (eds),

*Understanding Expository Text: A Theoretical and Practical Handbook for Analysing Explanatory Texts* (Lawrence Erlbaum, Hillsdale, NJ, 1985), pp. 11–64.

[23] T.A. van Dijk, *Macro-Structures* (Lawrence Erlbaum, Hillsdale, NJ, 1980).

[24] E.D. Liddy, Discourse-level structure in abstracts. In: *Proceedings of the 50th Annual Meeting of ASIS* (Learned Information, Medford, NJ, 1987), pp. 138–147.

[25] G. Salton *et al.*, Automatic text structuring and summarization, *Information Processing and Management* 33(2) (1997) 193–207.

[26] K. Markey *et al.*, An analysis of controlled vocabulary and free-text search statements in online searches, *Online Review* 4 (1980) 225–236.

[27] C. Tenopir, Full text database retrieval performance, *Online Review* 9 (1985) 149–164.

[28] J.S. Ro, An evaluation of the applicability of ranking algorithms to improve the effectiveness of full-text retrieval, *Journal of the American Society for Information Science* 39 (1988) 73–78.

[29] D.C. Blair and M.E. Maron, An evaluation of retrieval effectiveness for a full-text document-retrieval system, *Communications of the ACM* 28 (1985) 289–299.

[30] N. Wiener, *Cybernetics* (John Wiley, New York, 1948).

[31] N. Chomsky, *Syntactic Structures* (Mouton, The Hague, 1957).

[32] M.A. Minsky, Framework for representing knowledge. In: P.H. Winston (ed.), *The Psychology of Computer Vision* (Academic Press, New York, 1975).

[33] J.R. Anderson, Acquisition of cognitive skills, *Psychological Review* 89(3) (1982) 369–406.

[34] Z.W. Pylyshyn, *Computation and Cognition: Toward a Foundation for Cognitive Science* (MIT Press, Cambridge, MA, 1984).

[35] J.D. Fodor, On modularity in syntactic processing, *Journal of Psycholinguistic Research* 17(2) (1988) 125–168.

[36] D.E. Rumelhart and J.L. McClelland, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition* (MIT Press, Cambridge, MA, 1986).

[37] T. Winograd, A framework for understanding discourse. In: M.A. Just and P.A. Carpenter (eds), *Cognitive Processes in Comprehension* (Wiley, New York, 1977).

[38] T.A. van Dijk and W. Kintsch, *Strategies of Discourse Comprehension* (Academic Press, New York, 1983).

[39] K.I. Forster, Levels of processing and the structure of language processors. In: W.E. Cooper and E.C. Walker (eds), *Sentence Processing* (Lawrence Erlbaum, Hillsdale, NJ, 1979).

[40] W.D. Marslen-Wilson, Functional parallelism in spoken word-recognition, *Cognition* 25 (1987) 71–102.

[41] J.P. Stemberger, An interactive activation model of language production. In: A.W. Ellis (ed.), *Progress in the Psychology of Language, Vol. 1* (Lawrence Erlbaum, London, 1985).

[42] D. Laberge and S.J. Samuels, Toward a theory of automatic information processing in reading, *Cognitive Psychology* 6 (1974) 293–323.

[43] K. Goodman, Psycholinguistic universals in the reading process. In: P. Primsleurs and T. Quinn (eds), *The Psychology of Second Language Learning* (Cambridge University Press, Cambridge, 1971).

[44] D.E. Rumelhart, *Human Information Processing* (Wiley, New York, 1977).

[45] J.R. Anderson, *The Architecture of Cognition* (Harvard University Press, Cambridge, MA, 1983).

[46] J. Brown, Some tests of the decay theory of immediate memory, *Quarterly Journal of Experimental Psychology* 10 (1958) 12–21.

[47] L.R. Peterson, Short-term retention of individual items, *Journal of Experimental Psychology* 58 (1959) 12–21.

[48] R.C. Atkinson and R.M. Shiffrin, Human memory: a proposed system and its control processes. In: K.W. Spence and J.T. Spence (eds), *The Psychology of Learning and Motivation: Advances in Research Theory* (Academic Press, New York, 1968).

[49] L. Flower and J.R. Hayes, A cognitive process theory of writing, *College Composition and Communication* 32 (1988) 365–387.

[50] M. Pinto, Interdisciplinary approaches to the concept and practice of written text documentary content analysis (WTDCA), *Journal of Documentation* 50(2) (1994) 113–133.

[51] B. Endres-Niggemeyer, Professional summarising: no cognitive simulation without observation. In: *Fourth International Colloquium on Cognitive Science* (1995).

[52] M. Pinto and C. Gálvez, *Análisis Documental de Contenido* [Documentary content analysis] (Síntesis, Madrid, 1996).

[53] D.E. Rumelhart and A. Ortony, The representation of knowledge in memory. In: A. Anderson, R.J. Spiro and N.E. Montague (eds), *Schooling and the Acquisition of Knowledge* (Lawrence Erlbaum, Hillsdale, NJ, 1977), pp. 99–135.

[54] R.C. Schank and R.P. Abelson, *Scripts, Plans, Goals and Understanding* (Lawrence Erlbaum, Hillsdale, NJ, 1977).

[55] M.R. Quillian, The teachable language comprehender: a simulation program and theory of language, *Communications of the ACM* 12(8) (1969) 459–476.

[56] G. DeJong, An overview of the FRUMP systems. In: W.G. Rehnert and M.H. Ringle (eds), *Strategies for Natural Language Processing* (Lawrence Erlbaum, London, 1982), pp. 149–172.

[57] L.F. Rau, Knowledge organization and access in a conceptual information system, *Information Processing and Management* 23(4) (1987) 269–283.

[58] U. Hahn and U. Reimer, *The TOPIC Project: Text-Oriented Procedures for Information Management and Condensation of Expository Tests* (University of Constance, Constance, Germany, 1985).

[59] D. Fum *et al.*, Forward and backward reasoning in automatic abstracting. In: *Proceedings of the Ninth International Conference on Computational Linguistics* (North-Holland Publishing, Amsterdam, 1982), pp. 83–88.

[60] U. Hahn and U. Reimer, Heuristic text parsing in TOPIC: methodological issues in a knowledge-based text condensation system. In: H.J. Dietschmann (ed.), *Representation and Exchange of Knowledge as a Basis of Information Processes* (North-Holland Publishing, Amsterdam, 1984), pp. 143–163.

[61] R. Grishman, J. Sterling and C. Macleod, Description of the PROTEUS systems as used for MUC-3. In: *Proceedings of the Message Understanding Conference–3* (Morgan Kaufmann, San Mateo, CA, 1991), pp. 183–190.

[62] A. Moreno Sandoval *et al.*, PROTEUS: un sistema multilingüe de extracción de información [PROTEUS: a multilingual system for information extraction]. In: *VIII Congreso de la Sociedad Española para el Procesamiento del Lenguaje Natural (SEPLN)* (SEPLN, Granada, 1993), pp. 47–56.

[63] M. De Mey, The relevance of the cognitive paradigm for information science. In: O. Harbo and L. Kajberg (eds), *Theory and Application of Information Research. Proceedings of the 2nd International Research Forum in Information Science* (Mansell, London, 1980), pp. 49–61.

[64] N.J. Belkin, The cognitive viewpoint in information science, *Journal of Information Science* 16 (1990) 11–15.

[65] P.J. Daniels, Cognitive models in information retrieval: an evaluation review, *Journal of Documentation* 42(4) (1986) 272–304.

[66] P. Ingwersen, Search procedures in the library analysed from the cognitive point of view, *Journal of Documentation* 38 (1982) 165–191.

[67] P. Ingwersen, *Information Retrieval Interaction* (Taylor Graham, London, 1992).

[68] R. Fidel, Writing abstracts for free-text searching, *Journal of Documentation* 42(1) (1986) 11–21.

[69] E.D. Liddy, The discourse-level structure of empirical abstracts: an exploratory study, *Information Processing and Management* 27(1) (1991) 55–81.

[70] J. Hartley, M. Sydes and A. Blurton, Obtaining information accurately and quickly: are structured abstracts more efficient? *Journal of Information Science* 22(5) (1996) 349–356.

[71] L. von Bertalanffy, *General Systems Theory* (G. Braziller, New York, 1968).

[72] J. Willitts, *Database Design and Construction: An Open Learning Course for Students and Information Managers* (Library Association Publishing, London, 1992).

[73] E. Yourdon, *Modern Structured Analysis* (Prentice-Hall, Englewood Cliffs, NJ, 1988).

[74] P.B. Checkland, *Systems Thinking, Systems Practice* (Wiley, Chichester, UK, 1981).

[75] M. Bunge, Social systems. In: R. Rodríguez Delgado and B.H. Banathy (eds), *International Systems Science Handbook: An Introduction to Systems Science for Everbody* (Systemic Publications, Madrid, 1993), pp. 211–221.

[76] D. Soergel, *Organizing Information: Principles of Data and Retrieval Systems* (Academic Press, London, 1985).

[77] C.T. Meadow, *Text Information Retrieval Systems* (Academic Press, London, 1992).