# Approaches to Text Summarization: Questions and Answers

**Laura Alonso, Irene Castellón**

Dept. de Lingüística General
Universitat de Barcelona
Gran Via de les Corts Catalanes, 585
08007 Barcelona
{lalonso,castel}@lingua.fil.ub.es

**Salvador Climent**

Estudis d'Humanitats i Filologia
Universitat Oberta de Catalunya
Av. Tibidabo 39-43
08035 Barcelona
scliment@uoc.edu

**Maria Fuentes**
Dept. d'Informàtica i Matemàtica Aplicada
Universitat de Girona
Campus Montilivi
17071 Girona
maria.fuentes@udg.es

**Lluís Padró, Horacio Rodríguez**
TALP Research Center
Universitat Politècnica de Catalunya
Jordi Girona 1-3
08034 Barcelona
{padro,horacio}@lsi.upc.es

## Abstract

In this paper a comparative study of Automated Text Summarization (TS) Systems is presented. It describes the factors to be taken into account for evaluating those systems and outlines three alternative classifications. The paper provides extensive examples of working TS systems according to their characterizing features, performance, and obtained results, with a special emphasis on the multilingual aspect of summarization.

**Key Words**: Automated Text Summarization, Multilingual Systems

## 1 Introduction

The field of Text Summarization (TS) has experienced an exponential growth in the last years. That is why many comparative studies can be found in the literature, among the most comprehensive, Paice (1990) [123], Zechner (1997) [168], Sparck-Jones (1998) [147], Hovy and Marcu (1998) [70], Tucker (1999) [157], Radev (2000) [130], Mani (2001) [98] and Maybury and Mani (2001) [105]. Given that an upper bound of performance for TS systems is still far from being reached, task-based competitions are the main forum of discussion in the area. As follows, the SUMMAC (1998) [149] and especially DUC (2001, 2002, 2003) [45] contests provide a good overview of current working systems.

In this study, we provide an analysis of current work in TS, with special attention to the future developments of the field, like multilingual summarization. First, we present the factors affecting summarization in Section 2, and provide examples of how working systems handle each of these factors. In Section 3 three possible classifications of summarization systems are outlined, which are applied to concrete systems in Section 4, with a concrete example of multilingual summarization. To finish, we briefly discuss some burning issues in TS.

# 2 Some considerations on Summary Aspects

Summarization has traditionally been decomposed into three phases [147, 101, 58, 68, 98]:

- *analyzing* the input text to obtain text representation,
- *transforming* it into a summary representation,
- and *synthesizing* an appropriate output form to generate the summary text.

Effective summarizing requires an explicit and detailed analysis of context factors, as is apparent when we recognize that what summaries should be like is defined by what they are wanted for. The parameters to be taken into account in summarization systems have been widely discussed [101, 68, 98]. We will follow Sparck Jones (1998) [147], who distinguishes three main aspects that affect the process of TS: input, purpose and output, with a special focus on multilinguality.

## 2.1 Input Aspects

The features of the text to be summarized crucially determine the way a summary can be obtained. The following aspects of input are relevant to the task of TS:

***Document Structure.*** Besides textual content, heterogeneous documental information can be found in a source document, for example, labels that mark headers, chapters, sections, lists, tables, etc. If it is well systematized and exploited, this information can be of use to analyze the document. For example, Kan (2002) [74] exploits the organization of medical articles in sections to build a tree-like representation of the source. Teufel and Moens (2002) [156] systematize the structural properties of scientific articles to assess the contribution of each textual segment to the article, in order to build a summary from that enriched perspective.

However, it can also be the case that the information it provides is not the target of the analysis. In this case, document structure has to be removed in order to isolate the textual component of the document.

***Domain.*** Domain-sensitive systems are only capable of obtaining summaries of texts that belong to a pre-determined domain, with varying degrees of portability. The restriction to a certain domain is usually compensated by the fact that specialized systems can apply knowledge intensive techniques which are only feasible in controlled domains, as is the case of the multidocument summarizer SUMMONS [111], specialized in summaries in terrorism domain applying complex Information Extraction techniques. In contrast, general purpose systems are not dependant on information about domains, which usually results in a more shallow approach to the analysis of the input documents.

Nevertheless, some general purpose systems are prepared to exploit domain specific information. For example, the meta summarizer developed at Columbia University [19, 18, 61, 60, 108] applies different summarizers for different kinds of documents: MULTIGEN [19, 109] is specialized in simple events, DEMS [143] (with the bio configuration) deals with biographies, and for the rest of documents, DEMS has a default configuration that can be resorted to.

***Specialization level.*** A text may be broadly characterized as ordinary, specialized, or restricted, in relation to the presumed subject knowledge of the source text readers. This aspect can be considered the same as the *domain* aspect discussed above.

***Restriction on the language.*** The language of the input can be general language or restricted to a sublanguage within a domain, purpose or audience. It may be necessary to preserve the sublanguage in the summary.

***Scale.*** Different summarizing strategies have to be adopted to handle different text lengths. Indeed, the analysis of the input text can be performed at different granularities, for example, in determining meaning units. In the case of news articles, sentences or even clauses are usually considered the minimal meaning units, whereas for longer documents, like reports or books, paragraphs seem a more adequate unit of meaning. Also the techniques for segmenting the input text in these meaning units differ: for shorter texts, orthography and syntax, even discourse boundaries [103] indicate significant boundaries, for longer texts, topic segmentation [79, 63] is more usual.

***Media.*** Although the main focus of summarization is textual summarization, summaries of nontextual documents, like videos, meeting records, images or tables have also been undertaken in recent years. The complexity of multimedia

summarization has prevented the development of wide coverage systems, which means that most summarization systems that can handle multimedia information are limited to specific domains or textual genres [62, 104]. However, research efforts also consider the integration of information of different media [21], which allow a wider coverage of multimedia summarization systems by exploiting different kinds of documental information collaboratively, like metadata associated to video records [161].

*Genre.* Some systems exploit typical genre-determined characteristics of texts, such as the pyramidal organization of newspaper articles, or the argumentative development of a scientific article. Some summarizers are independent of the type of document to be summarized, while others are specialized on some type of documents: healthcare reports [48], medical articles [74], agency news [111], broadcast fragments [62], meeting recordings [169], e-mails [117, 3], web pages [132], etc.

*Unit.* The input to the summarization process can be a *single document* or *multiple documents*, either simple text or multimedia information such as imagery audio, or video [150].

*Language.* Systems can be language-independant, exploiting characteristics of documents that hold cross-linguistically [129, 125], or else their architecture can be determined by the features of a concrete language. This means that some adaptations must be carried out in the system to deal with different languages. As an additional improvement, some multi-document systems are able to deal simultaneously with documents in different languages [33, 34], which will be developed in Section 2.4.

## 2.2   Purpose Aspects

*Situation.* TS systems can perform general summarization or else they can be embedded in larger systems, as an intermediate step for another NLP task, like Machine Translation, Information Retrieval or Question Answering. As the field evolves, more and more efforts are devoted to task-driven summarization, in detriment of a more general approach to TS. This is due to the fact that underspecification of the information needs supposes a major problem for design and evaluation of the systems. As will be discussed in Section 5, evaluation is a major problem in TS.

Task-driven summarization presents the advantage that systems can be evaluated with respect to the improvement they introduce in the final task they are applied to.

*Audience.* In case a user profile is accessible, summaries can be adapted to the needs of specific users, for example, the user's prior knowledge on a determined subject. *Background* summaries assume that the reader's prior knowledge is poor, and so extensive information is supplied, while *just-the-news* are those kind of summaries conveying only the newest information on an already known subject. Briefings are a particular case of the latter, since they collect representative information from a set of related documents.

*Usage.* Summaries can be sensitive to determined uses: retrieving source text [75], previewing a text [88], refreshing the memory of an already read text, sorting...

## 2.3   Output Aspects

*Content.* A summary may try to represent all relevant features of a source text or it may focus on some specific ones, which can be determined by queries, subjects, etc. *Generic* summaries are text-driven, while *user-focused* (or query-driven) ones rely on a specification of the user's information need, like a question or key words.

Related to the kind of content that is to be extracted, different computational approaches are applied. The two basic approaches are top-down, using information extraction techniques, and bottom-up, more similar to information retrieval procedures. Top-down is used in query-driven summaries, when criteria of interest are encoded as a search specification, and this specification is used by the system to filter or analyze text portions. The strategies applied in this approach are similar to those of Question Answering. On the other hand, bottom-up is used in text-driven summaries, when generic importance metrics are encoded as strategies, which are then applied over a representation of the whole text.

*Format.* The output of a summarization system can be plain text, or else it can be formatted. Formatting can be targeted to many purposes: conforming to a pre-determined style (tags, organization in fields), improving readability (division in sections, highlighting), etc.

*Style.* A summary can be *informative*, if it covers the topics in the source text; *indicative*, if it

provides a brief survey of the topics addressed in the original; *aggregative*, if it supplies information non present in the source text that completes some of its information or elicits some hidden information [156]; or *critical*, if it provides an additional valoration of the summarized text.

***Production Process.*** The resulting summary text can be an *extract*, if it is composed by literal fragments of text, or an *abstract*, if it is generated. The type of summary output desired can be relatively polished, for example, if text is well-formed and connected, or else more fragmentary in nature (e.g., a list of key words).

There are intermediate options, mostly concerning the nature of the fragments that compose extracts, which can range from topic-like passages, paragraph or multiparagraph long, to clauses or even phrases. In addition, some approaches perform editing operations in the summary, overcoming the incoherence and redundancy often found in extracts, but at the same time avoiding the high cost of a NL generation system. Jing and McKeown (2000) [73] apply six re-writing strategies to improve the general quality of an extract-based summary by edition operations like deletion, completion or substitution of clausal constituents.

***Surrogation.*** Summaries can stand in place of the source as a surrogate, or they can be linked to the source [75, 88], or even be presented in the context of the source (e.g., by highlighting source text, [86]).

***Length.*** The targeted length of the summary crucially affects the informativeness of the final result. This length can be determined by a compression rate, that is to say, a ratio of the summary length with respect to the length of the original text. Traditionally, compression rates range from 1% to 30%, with 10% as a preferred rate for article summarization. In the case of multidocument summarization though, length cannot be determined as a ratio to the original text(s), so the summary always conforms to a pre-determined length. Summary length can also be determined by the physical context where the summary is to be displayed. For example, in the case of delivery of news of summaries to handhelds [23, 28, 39], the size of the screen imposes severe restrictions to the length of the summary. Headline generation is another application where the length of summaries is clearly determined [165, 41]. In very short summaries, coherence is usually sacrificed to informativeness, so lists of words are considered acceptable [80, 167].

## 2.4   Language coverage

As regards language coverage, systems can be classified as monolingual, multilingual, and crosslingual (a similar classification is commonly used in Information Retrieval systems). Monolingual summarization systems deal with only one language for both the input document and the summary. In the case of multilingual systems, input and output languages are also the same but in this case the system can cover several languages. Crosslingual systems are able to process input document in several languages, producing summaries in different languages.

Multilinguality does not imply additional difficulties. Most of the systems and techniques we will present below can be easily adapted to other languages, assuming, of course, the availability of the knowledge sources needed for the different methods. Roughly speaking, the more amount of linguistic knowledge is needed by a system, the more difficult is to transport it to another language.

A more complex challenge is crosslinguality. There are examples of single document crosslingual summarizers, implying a certain amount of translation, either on the input text or on the summary, but most crosslingual summarizers are multidocument. In this case a lot of problems specific of translinguality arise. Measures of similarity between documents and passages in different languages, for identifying relations or for clustering, have to be envisaged. Similarity between lexical units (words, NEs, multiword terms) belonging to different languages, have to be computed as well. Obviously, the more distant the involved languages are, the harder these problems turn to be, specially if the languages present different lexical units or character sets. Since this is a burning issue, it will be discussed at length in Section 5.

# 3   Approaches to Text Summarization

There are several ways in which one can characterize different approaches to text summarization. In this section, we present three possible classifications of text summarization systems, but many others can be found in the literature [70, 130, 105, 98]. The first classification, following Mani and Maybury (1999) [101], is based in

the level of processing that each system performs, the second, proposed in Alonso and Castellón (2001) [4], is based in the kind of information exploited, the third follows Tucker (1999) [157].

## 3.1 Classification 1: Level of Processing

One useful way to classify summarization systems is to examine the level of processing of the text. Based on this, summarization can be characterized as approaching the problem at the surface, entity, or discourse level [101].

### 3.1.1 Surface level

Surface-level approaches tend to represent information in terms of shallow features that are then selectively combined together to yield a salience function used to extract information, following the approach of Edmunson (1969) [47]. These features include:

***Term frequency*** statistics provide a thematic representation of text, assuming that important sentences are the ones that contain words that occur frequently. The score sentences increases for each frequent word. Early summarization systems directly exploit word distribution in the source [96].

***Location*** relies on the intuition that important sentences are located at positions that are usually genre-dependent, however, some general rules are the *lead method* and the *title-based method*. The lead method consists of just taking the first sentences. The title-based method assumes that words in titles and headings are positively relevant to summarization. A generalization of these methods is the OPP used by Hovy and Lin in their SUMMARIST system [91], where they exploit Machine Learning techniques to identify the positions where relevant information is placed within different textual genres. Many of the current systems, specially those applying machine learning techniques, take into account the location of meaning units in a document to assess their relevance.

***Bias.*** The relevance of meaning units is determined by the presence of terms from the title or headings, initial part of text, or user's query. For example, [37, 36, 144] use as features the position in the sentence, the number of tokens and the number of pseudo-query terms.

***Cue words*** and *phrases* are signals of relevance or irrelevance. They are typically meta-linguistic markers (e.g., cues: "in summary", "in conclusion", "our investigation", "the paper describes"; or emphasizers: "significantly", "important", "in particular", "hardly", "impossible"), as well as domain-specific bonus phrases and stigma terms. Although lists of these phrases are usually built manually [82, 154], they can also be detected automatically.

### 3.1.2 Entity-level

Entity-level approaches build an internal representation of the text by modeling text entities (simple words, compound nouns, named entities, etc.) and their relationships. These approaches tend to represent patterns of connectivity in the text (e.g., graph topology) to help determine saliency. Relations between entities include:

***Similarity.*** Similar words are those whose form is similar, for example, those sharing a common stem (e.g., "similar" and "similarity"). Similarity can be calculated with linguistic knowledge or by character string overlap. Myaeng and Jang (1999) [118] use two similarity measures for determining if a sentence belongs to the major content: a similarity between the sentence and the rest of the document and a similarity between the sentence and the title of the document. Also, in NTT [65, 66], CENTRIFUSER [75], several similarity measures are applied.

***Proximity.*** The distance between the text units where entities occur is a determining factor for establishing relations between entities.

***Cohesion.*** Cohesion can be defined in terms of *connectivity*. Connectivity accounts for the fact that important text units usually contain entities that are highly connected in some kind of semantic structure. Cohesion can be approached by:

- *Word co-occurrence*: words can be related if they occur in common contexts. Some applications are presented in Baldwin and Morton (1998), McKeown et al. (1999)[13, 109]. Salton et al. (1997), Mitra et al. (1997) [141, 113] apply IR methods at the document level, treating paragraphs in texts as documents are treated in a collection of documents. Using a traditional IR-based method, a word similarity measure is used to determine the set $S_i$ of paragraphs that each paragraph $P_i$ is related to. After

determining relatedness scores $S_i$ for each paragraph, paragraphs with the largest $S_i$ scores are extracted.

In SUMMAC [97], in the context of query-based summarization, Cornell's Smart-based approach expands the original query, compares expanded query against paragraphs, and selects top three paragraphs (max 25% of original) that are most similar to the original query.

- *Local salience*: important phrasal expressions are given by a combination of grammatical, syntactic, and contextual parameters [24].

- *Lexical similarity*: words can be related by thesaural relationships (synonymy, hypernymy, meronymy relations). Barzilay (1997) [16] details a system where Lexical Chains are used, based on Morris and Hirst (1991) [116]. This line has also been applied to Spanish, relying on EuroWordNet relations between words, by Fuentes and Rodríguez (2002) [53]. The assumption is that important sentences are those that are crossed by strong chains[1]. This approach provides a partial account of texts, since it focuses mostly on cohesive aspects. An integration of cohesion and coherence features of texts might contribute to overcome this, as Alonso and Fuentes (2002) [5] point out.

- *Co-reference*: referring expressions can be linked, and co-reference chains can be built with co-referring expressions. Both Lexical Chains and Co-reference Chains can be priorised if they contain words in a query (for query-based summaries) or in the title. So, the preference imposed on chain is: query > title > document. Baga and Baldwin (1998), Azzam et al. (1999) [11, 10] use coreference chains for summarization. Baldwin and Morton (1998) [13] exploit co-reference chains specifically for query-sensitive summarization.

  Connectedness method [100] represents map text with graphs. Words in the text are the nodes, and arcs represent adjacency, grammatical, co-reference, and lexical similarity-based relations.

**Logical relations** such as agreement, contradiction, entailment, and consistency.

**Meaning representation-based relations.** Establishing relations, such as predicate-argument, between entities in the text.

The system of Baldwin and Morton (1998) [13] uses argument detection in order to resolve co reference between the query and the text for performing summarization.

### 3.1.3 Discourse-level

Discourse-level approaches model the global structure of the text, and its relation to communicative goals. At this level, the following information can be exploited:

**Format** of the document (e.g., hypertext markup, document outlines).

**Threads of topics** can be revealed in the text. An example of this is SUMMARIST, which applies Topic identification [69, 95]. Topic identification implies previous acquisition of Topic Signatures (that can be automatically learned) and then the identification of a text span as belonging to a topic characterized by its signature. Topic identification, then, includes text segmentation and comparison of text spans with existing Topic Signatures. The topic identified are fused during the interpretation of the process. The fused topics are then expressed in new terms. Other systems are Boros et al. (2001) [25] and MEAD [133, 128, 121]. These systems assign a topic to the sentences in order to create clusters for selecting the sentences to appear in summary.

**Rhetorical structure** of the text, representing argumentation or narrative structure. The main idea is that the coherence structure of a text can be constructed, so that the 'centrality' of the textual units in this structure will reflect their importance. A tree-like representation of texts is proposed by the Rhetorical Structure Theory [102]. Ono et al. (1994) [120] and Marcu (1997) [103] attempt to use this kind of discourse representation in order to determine the most important textual units. They propose an approach to rhetorical parsing by discourse markers and semantic similarities in order to hypothesize rhetorical relations. These hypotheses are used to derive a valid discourse representation of the original text.

---

[1]Lexical chains have also been used in other NLP tasks, such as automatic extraction of interdocument links [56].

## 3.2 Classification 2: Kind of Information

Summarization systems can be classified by the kind of information they deal with [4]. According to this, we can distinguish between those exploiting lexical aspects of texts, those working with structural information and those trying to achieve deep understanding of texts.

### 3.2.1 Lexical

These approaches exploit the information associated to words in the texts. Some of them are very shallow, relying on the frequency of words, but some others apply lexical resources to obtain a deeper representation of texts. Beginning by the most shallow, the following main trends can be distinguished. A common assumption of these approaches is that repeated information could be a good indicator of importance:

***Word Frequency*** approaches assume that the most frequent words in text are the most representative of its content, and consequently fragments of text containing them are more relevant. Most systems apply some kind of filter to leave out of consideration those words that are very frequent but not indicative, for example, by the *tf\*idf* metric or by excluding the so-called *stop words*, words with grammatical but no meaning content.

***Domain Frequency*** tries to determine the relevance of words by first assigning the document to a particular domain. Domain specific words have a previous relevance score, which serves as a comparison ground to adequately evaluate their frequency in a given text.

***Concept Frequency*** abstracts from mere word-counting to concept-counting. By use of an electronic thesaurus or WordNet, each word in the text is associated to a more general concept, and frequency is computed on concepts instead of particular words.

***Cue words and phrases*** can be considered as indicators of relative relevance or non-relevance of fragments of text in respect to the others.

***Chains*** can be built from lexical items which are related by conceptual similarity according to a lexical resource (*lexical chains*) or by identity, if they co-refer to the same entity (*co-reference chains*). The fragments of text crossed by most chains or by most important chains or by most important parts of chains can be considered the most representative of the text.

### 3.2.2 Structural Information

A second direction in TS tries to exploit information from the texts as structured entities. Since texts are structured in different dimensions (documental, discursive, conceptual), different kinds of structural information can be exploited. Beginning by the most shallow:

***Documental Structure*** exploits the information that texts carry in their format, for example, headings, sections, etc.

***Textual Structure*** . Some positions in text systematically contain the most relevant information, for example, the beginning paragraph of news stories. These positions are usually genre- or domain-dependant.

***Conceptual structure*** . The chains mentioned in lexical approaches can be considered as a kind of conceptual structure.

***Discursive Structure*** can be divided in two main lines: linear or narrative and hierarchical or rhetoric. The first tries to account for *satisfaction-precedence*-like relations among pieces of text, the second explains texts as trees where fragments of text are related with each other by virtue of a set of rhetorical relations, mostly asymmetric.

### 3.2.3 Deep Understanding

Some approaches try to achieve understanding of the text in order to build a summary. Two main lines can be distinguished:

***Top-down*** approaches try to recognize predefined knowledge structures to texts, for example, templates or frames.

***Bottom-up*** approaches try to represent texts as highly conceptual constructs, such as scene. Others apply fragmentary knowledge-structures to clue parts of text, and then build a complete representation out of these small parts.

## 3.3 Classification 3: Richard Tucker 1999

This classification is taken from Tucker (1999) [157]. It considers four main directions in TS: summarizing from attentional networks, sentence by sentence, from informational content and from discourse structure.

The classes proposed here are even less disjunct than those in the two previous classifications, thus every system can be considered as an instance of more than one of the classes. This shows the inadequacy of a taxonomic perspective on summarization systems, due to the heterogeneous kinds of knowledge and techniques that systems tend to incorporate.

### 3.3.1 Attentional Networks

The approaches to summarization in this direction try to grasp what a text is 'about' by identifying concepts that are in some sense central to the text, on the basis of the occurrence of the same or related concepts in different parts of the source representation. *Aboutness* is represented as the links between these occurrences.

*Frequency-based* approaches exploit the frequency with which the concepts occur in the representation. In systems based in word frequency, attentional networks are only represented implicitly. Some systems account for frequency significance by applying IR techniques, such as the *tf\*idf* measure. Others apply corpus-based statistical natural language processing, such as collocation or proper noun identification. Sill others try to abstract from individual words to achieve concept frequency, by using lexicons or thesauri [69].

On the other hand, some systems identify and exploit the *cohesive links* holding between parts of the source text. These links can be represented as graph-like structures [145] as lexical chains.

### 3.3.2 Sentence by Sentence

Some summarizing systems decide for each sentence in the source text whether it is important for summarizing, rather independently of the text as a whole. To do that, they rely on relevance or irrelevance marks that can be found in sentences, for example, *cue words*.

However, it must be noted that most of the systems applying sentence-by-sentence relevance ranking do not rely entirely in this method, but use it in combination with other methods that tend to consider the text as a whole.

### 3.3.3 Informational Content

Some approaches to summarization have tried to understand the text, that is to say, to achieve a representation of some or all of its meaning whereupon reasoning can be applied. This approach requires deeper analysis of the source text but allows the production of sophisticated summaries, for example, by applying NL generation techniques. However, these methods tend to be highly domain-dependant, because of the huge amount of information they require.

### 3.3.4 Discourse Structure

Discourse structure is used by many systems in a limited way, for example, by trying to grasp a text's 'aboutness'. In contrast, some other methods apply discourse theories to the analysis of the source text in order to obtain a representation of their discourse structure. However, work in this area has been largely theoretical.

## 3.4 Combined Systems

The predominant tendency in current systems is to integrate some of the techniques mentioned so far. Integration is a complex matter, but it seems the appropriate way to deal with the complexity of textual objects. In this section, we are going to present some examples of combination of different techniques.

There are several systems where different methods are combined. Among the most interesting are: [82, 156, 69, 100] where title-based method is combined with cue-location, position, and word-frequency based methods.

As the field progresses, summarization systems tend to use more and deeper knowledge. For example, IE techniques are becoming widely used. Many systems do not rely any more in a single indicator of relevance or coherence, but take into account as many of them as possible. So, the tendency is that heterogeneous kinds of knowledge

are merged in increasingly enriched representations of the source text(s).

These enriched representations allow for adaptability of the final summary to new summarization challenges, such as multidocument, multilingual and even multimedia summarization. In addition, such a rich representation of text is a step forward generation or, at least, pseudo-generation by combining fragments of the original text. Good examples of this are [108, 93, 41, 84, 59], among others.

# 4    Summarization Systems

Table 1 shows how existing summarization systems would be classified according to each of the classifications presented in the previous section. However, it must be taken into account that most current summarization systems are very complex, resorting to very heterogeneous information and applying varied techniques, so a classification will never be clear cut. Moreover, systems tend to evolve with time, which makes their classification still more controversial.

Files with a more extense description of some of these systems (marked with an asterisk) can be found in the Annex (in electronic version only). Additionally, Table 2 lists on-line or downloadable systems.

Multilinguality of the systems is one of the features in each describing file. It is stated whether the system can summarize only a single language, a definite set of languages, or whether its architecture permits unrestricted multilinguality. In this latter case, it is also stated whether experiments with different languages are reported.

As a concrete example of an approach to multilingual summarization, we present the systems developed within project HERMES[2]. The target of project HERMES is to adapt and apply language technologies for Spanish, Catalan, Basque and English to improve access to textual information in digital libraries, Internet, documental Intranets, etc. Therefore, HERMES summarization system should integrate multiple languages in a common architecture. Since the resources available for every language are uneven, this architecture has to be flexible enough to adapt to knowledge-poor representations of text but also to exploit rich representations when available.

EuroWordNet [160] is a general resource available for these four languages, so a first approach to summarization exploited this resource. A Lexical Chain summarizer was developed for Spanish [53]. As can be seen in Figure 1, the architecture of the summarizer permits easy adaptation to other languages, provided there is at least a morphological analyzer and a version of EuroWordNet available for the language. If other NLP tools are available, like Named Entity Recognizers or co-reference solvers, they can be easily integrated within the system. Once the text has been analyzed and Lexical Chains have been obtained, a summary is built by extracting candidate textual units from the text. Candidate units are chosen applying a certain heuristic, weighting some aspects of Lexical Chains.

A second approach to the task of summarization, seen in Figure 2, [52] tries to overcome this dependency on lexic applying Machine Learning techniques. The system is trained with a corpus of sentences described with a set of features, like position in the text, length, and also being crossed by a Lexical Chain. For each of these sentences, it is previously determined whether it belongs to a summary of the text or not, so that it can be learned which combinations of features characterize summary sentences. In a text to be summarized, each sentence is described with the same set of features, and it is determined whether these describing features characterize the sentence as a summary sentence or not. The summary is composed with sentences qualifying as summary sentences.

This second system does not require any specific feature to produce a summary, not even Lexical Chains. However, the more information available, the more accurate the learning process will be, which will result in better summaries. This approach has been evaluated for English within DUC 2003 contest, but it can be used straightforwardly for any other language, as long as there is a training corpus available.

# 5    Burning Issues

The field has experienced an exponential growth since its beginnings, but some crucial questions are still open.

---

[2]http://terral.ieec.uned.es/hermes/

## 5.1 Coherence of Summary texts

Paice (1990) [123] pointed out that the main shortcomings of summarization systems up to the 1990s was their low representativity of the content in the source text and their lack of coherence.

Much of the work in this area has treated the problem of text summarization from a predominant information-theoretic perspective. Therefore, texts have been modeled as mathematical objects, where relevance and redundancy could be defined in purely statistical terms. This approach seems specially valuable to produce a satisfactory representation of the content of a text. However, it fails in producing coherent texts, acceptable for human users.

The shortcomings of purely statistical approaches to text summarization on handling textual coherence are addressed from two different perspectives:

- Applying *machine learning* techniques. They have been used mainly for two purposes: classifying a sentence from a source text into relevant or non-relevant [82, 8, 99, 90, 65] and transforming a source sentence considered relevant into a summary sentence [73, 78, 59]. Input for learning algorithms are usually texts with their corresponding abstracts. Therefore, the main shortcoming of this approach is to obtain large quantities of <text, abstract> tuples for a variety of textual genres.

- Resorting to *symbolic linguistic or world knowledge*. Understanding of texts, mainly through IE extraction techniques, seems a desirable way of producing quality summaries. Until recently, such techniques had only been applied for very restricted domains [111]. However, recent systems tend to incorporate IE extraction modules that perform a partial understanding of text, either by modeling the typical context of relevant pieces of information [84, 76], or by applying general templates to find, organize and use the typical content of a kind of text or event [59, 41]. This use of IE techniques has produced very good results, as is reflected in the high ranking of Harabagiu and Lacatusu (2002) [59] in DUC 2002. A combination of deeper knowledge with surface clues seems to yield good results, too [93].

## 5.2 Multidocument summarization

Multidocument summarization is one of the major challenges in current summarization systems. It consists of producing a single summary of a collection of documents dealing with the same topic. The work has been mostly determined by the corresponding DUC task. Therefore, it has mainly focused in collections of news articles with a given topic. Remarkable progresses have been achieved in avoiding redundancy, mainly based on the work in Carbonell and Goldstein (1998) [30].

When dealing with MDS new problems arise: lower compression factors implying a more aggressive condensation, anti-redundancy, temporal dimension, more challenging coreference task (inter-document), etc. Clustering of similar documents plays now a central role [30, 133, 60, 110]. Selecting the most relevant fragments from each cluster and assuring coherence of the summaries coming from different documents are other important problems, currently under development in MDS systems.

## 5.3 Multilingual summarization

As for multilingual summarization, not much work has been done yet, but the roadmap for the DUC contests [12] contemplates this challenge in the near future of the area.

The most well known Multilingual Summarization System is SUMMARIST [69]. The system extracts sentences in a variety of languages (English, Spanish, Japanese, etc.) and translates the resulting summaries. SUMMARIST proceeds in three steps: Topic identification, Interpretation and Summary generation. Topic identification implies previous acquisition of Topic Signatures and then the identification of a text span as belonging to a topic characterized by its signature. Topic Signatures are tuples of the form <Topic, Signature> where Signature is a list of weighted terms: $\{< t_1, w_1 >, < t_2, w_2 >, ..., < t_n, w_n >\}$. Topic signatures can be automatically learned [89, 95]. Topic identification, then, includes text segmentation (using Marti Hearst's TextTiling) and comparison of text spans with existing Topic Signatures. The identified topics are fused during interpretation, the second step of the process. The fused topics are then reformulated, that is to say, expressed in new terms. The last step is a conventional extractive task.

In order to face multilingual problems the involved knowledge sources have to be as much as possible language independent. In the case of SUMMARIST, sets of Topic Signatures have to be obtained for all the languages involved using the same procedures. Also the segmentation procedure is language independent. So, the accuracy of the resulting summaries depends heavily on the quality of the translators.

As has been said before, a more challenging issue is Crosslingual Multidocument Summarizers. Basically three main problems have to be addressed: 1) clustering of multilingual documents, 2) measuring the distance (or similarity) between multilingual units (documents, paragraphs, sentences, terms), and 3) automatic translation of documents or summaries. Most systems differ on the way they face these problems, the order of performance and the granularity of the units they deal with.

Evans and Klavans (2003) [49] present a platform for multilingual news summarization that extends the Columbia's Newsblaster system [106]. The system adds a new component, translation, to the original six major modules: crawling, extraction, clustering, summarization, classification and web page generation, that have been, in turn, modified for allowing multilinguality (language identification, different character encoding, language idiosyncrasy, etc.).

In this system multilingual documents are translated into English before clustering, so that clustering is performed only on English texts.

Translation is carried out at two levels. Because a low quality translation is usually enough for clustering purposes and assessing the relevance of the sentences, a simple and fast technique is applied for glossing the input documents prior to clustering. Higher (relatively) quality translation (using Altavista's Babelfish interface to Systran) is performed in a second step only over fragments selected to be part of the summary.

The system takes as well into account the possible degradation of the input texts as result of the translation process, since most of the sentences resulting from this process are simply not grammatically correct.

Chen et al. (2003) [34] consider three possibilities for scheduling the basic steps of document translation and clustering:

1. Translation before document clustering (as in Columbia's system), named one-phase strategy. This model clusters the multilingual multidocuments directly resulting in multilingual clusters.

2. Translation after document clustering, named two-phase strategy. This model clusters documents in each language separately and merges the clustering results.

3. Translation deferred to sentence clustering. First, monolingual clustering is performed at document level. All the documents in each cluster refer to the same event in a specific language. Then, for generating the extracted summary of an event all the clusters referring to this event are taken into account. Similar sentences of these multilingual clusters are clustered together, now at sentence level. Finally a representative sentence is chosen from each cluster and translated if needed.

The accuracy of this process depends basically on the form of computing the similarity between different multilingual units. Several forms of such functions are presented and empirially evaluated by the authors.

These measures are multilingual extensions of a baseline monolingual similarity measure. Sentences are represented as bag of words (only nouns and verbs are taken into account). The similarity measure is a function of the number of (approximate) matches between words and of the size of the bags. The matching function in the baseline reduces, except for NE, to the identity. In the multilingual variants of the formula, a bilingual dictionary is used as knowledge source for computing this matching.

Despite of its simplicity the position-free measure (the simplest one) seems to be the most accurate among the studied alternatives. In this approach the translations of all the words of the bag are collected and the similarity is computed as in the baseline. All the other alternatives constraint in some ways the possible mappings between words, using different greedy strategies. The results are, however, worse.

The two-phase strategy outperforms in the experiments the on-phase strategy. The third strategy, deferring the translation to sentence clustering, seems to be the most promising.

A system, covering English and Chinese, follow-

ing this approach is presented in Chen and Lin (2000) [35]. The main components of the system are a set of monolingual news clusterers, a unique multilingual news clusterer and a news summarizer. A central issue of the system is the definition and identification of meaningful units as base for comparison. For English these units can be reduced to sentences but for Chinese the identification of units and the associated segmentation of the text can be a difficult task. Another important issue of the system (general for systems covering distant languages or different encoding schemata) is the need of a robust transliteration of names (or words not occurring in the bilingual dictionary) for assuring an accurate matching.

## 5.4 Evaluation

Last but not least, evaluation of summaries is a major issue, because objective judgements are needed to assess the progress achieved by different approaches. Some contests have been carried out to evaluate summarization systems with common, public procedures: the SUMMAC contest and the series of DUC contests. Specially the last has provided sets of criteria to evaluate summary quality in many different dimensions: informational coverage (precision and recall), suitability to length requirements, grammatical and discursive coherence, etc.

An extensive investigation on the automatic evaluation of automatic summaries was carried out in a six-week workshop at Johns Hopkins University [134], where different evaluation metrics were proposed, including the *relative utility* method. Mani (2001) [98] provides a clear picture of the current state-of-the-art in evaluation, both with human judges and by automated metrics, with a special emphasis on content-based metrics. Hovy and Lin (2003) [94] show that the summaries produced by human judges are not reliable as a gold standard, because they strongly disagree with each other. A consensus summary obtained by applying content-based metrics, like unigram overlap, seems much more reliable as a golden standard against which summaries can be contrasted.

## Acknowledgements

## References

[1] Enrique Alfonseca and Pilar Rodríguez. Description of the UAM system for generating very short summaries at DUC-2003. In *HLT/NAACL Workshop on Text Summarization / DUC 2003*, 2003.

[2] D. Allport. The TICC: parsing interesting text. In *Proceedings of the Second Conference on Applied Natural Language Processing*, pages 211–218, 1988.

[3] Laura Alonso, Bernardino Casas, Irene Castellón, Salvador Climent, and Lluís Padró. CARPANTA eats words you don't need from e-mail. In *SEPLN, XIX Congreso Anual de la Sociedad Española para el Procesamiento del Lenguaje Natural*, 2003.

[4] Laura Alonso and Irene Castellón. Aproximació al resum automàtic per marcadors discursius. Technical report, CLiC, Universitat de Barcelona, Barcelona, 2001.

[5] Laura Alonso and Maria Fuentes. Collaborating discourse for text summarisation. In *Proceedings of the Seventh ESSLLI Student Session*, 2002.

[6] R. Angheluta, R. De Busser, and M-F. Moens. The use of topic segmentation for automatic summarization. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[7] Roxana Angheluta, Marie-Francine Moens, and Rik De Busser. K.u. leuven summarization system. In *DUC03*, Edmonton, Alberta, Canada, May 31 - June 1 2003. Association for Computational Linguistics.

[8] C. Aone, M. Okurowski, and J. Gorlinsky. Trainable scalable summarization using robust NLP and machine learning. In *COLING-ACL*, pages 62–66, 1998.

[9] Chinatsu Aone, Mary Ellen Okurowski, James Gorlinsky, and Bjornar Larsen. A scalable summarization system using robust NLP. In *Proceeding of the ACL'97/EACL'97 Workshop on Intelligent Scalable Text Summarization*, pages 66–73, 1997.

[10] Saliha Azzam, Kevin Humphrey, and Robert Gaizauskas. Using coreference chains for text summarisation. In Amit Bagga, Brek Baldwin, and Sara Shelton, editors, *Proceedings of the ACL'99 Workshop on Coreference and Its Applications*, pages 77 – 84, University of Maryland, College Park, Maryland, USA, June 1999. ACL.

[11] Amit Bagga and Breck Baldwin. Algorithms for scoring coreference chains. In *Proceedings of the Linguistic Coreference Workshop at The First International Conference on Language Resources and Evaluation (LREC'98)*, pages 536–566, Granada, 1998.

[12] Breck Baldwin, Robert Donaway, Eduard Hovy, Elizabeth Liddy, Inderjeet Mani, Daniel Marcu, Kathleen McKeown, Vibhu Mittal, Marc Moens, Dragomir Radev, Karen Sparck Jones, Beth Sundheim, Simone Teufel, Ralph Weischedel, and Michael White. An evaluation road map for summarization research. TIDES, TIDES 2000.

[13] Breck Baldwin and Thomas S. Morton. Dynamic coreference-based summarization. In *Proceedings of the Third Conference on Empirical Methods in Natural Language Processing*, Granada, Spain, June 1998.

[14] M. Banko, V. Mittal, and M. Witbrock. Headline generation based on statistical translation. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics, ACL*, 2000.

[15] Michele Banko, Vibhu Mittal, Mark Kantrowitz, and Jade Goldstein. Generating extraction-based summaries from handwritten summaries by aligning text spans. In *Proceedings of PACLING-9*, Waterloo, Ontario, July 1999.

[16] Regina Barzilay. Lexical chains for summarization. Master's thesis, Ben-Gurion University of the Negev, 1997.

[17] Regina Barzilay and Michel Elhadad. Using lexical chains for text summarization. In Inderjeet Mani and Mark Maybury, editors, *Intelligent Scalable Text Summarization Workshop (ISTS'97)*, pages 10–17, Madrid, 1997. ACL/EACL.

[18] Regina Barzilay, Noemie Elhadad, and Kathy McKeown. Sentence ordering in multidocument summarization. In *HLT'01*, 2001.

[19] Regina Barzilay, Kathy McKeown, and Michel Elhadad. Information fusion in the context of multi-document summarization. In *Proceedings of ACL 1999*, 1999.

[20] M. Benbrahim and K. Ahmad. Computer-aided lexical cohesion analysis and text abridgement. Technical Report Computing Sciences Report CS-94-11, University of Surrey, 1994.

[21] A. B. Benitez and S.-F. Chang. Multimedia knowledge integration, summarization and evaluation. In *Proceedings of the 2002 International Workshop On Multimedia Data Mining in conjuction with the International Conference on Knowledge Discovery and Data Mining (MDM/KDD-2002)*, Edmonton, Alberta, 2002.

[22] Adam Berger and Vibhu Mittal. Ocelot: A system for summarizing web pages. In *Proceedings of the 23rd Annual Conference on Research and Development in Information Retrieval (ACM SIGIR)*, Athens, 2001.

[23] Branimir Boguraev, Rachel Bellamy, and Calvin Swart. Summarisation miniaturisation: Delivery of news to hand-helds. In *NAACL'01*, 2001.

[24] Branimir Boguraev and Christopher Kennedy. Salience-based content characterisation of text documents. In *Proceedings of ACL'97 Workshop on Intelligent, Scalable Text Summarisation*, pages 2–9, Madrid, Spain, 1997.

[25] E. Boros, P.B. Kantor, and D.J. Neu. A clustering based approach to creating multi-document summaries. In *Workshop on Text Summarization in conjunction with the ACM SIGIR Conference 2001*, New Orleans, 2001.

[26] Ronald Brandow, Karl Mitze, and Lisa F. Rau. Automatic condensation of electronic publications by sentence selectio.

*Information Processing and Management,* 31(5):675–68, 1995.

[27] M. Brunn, Y. Chali, and B. Dufou. The University of Lethbridge text summarizer at DUC 2002. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization),* Philadelphia, July, 11-12 2002.

[28] Orkut Buyukkokten, Hector Garcia-Molina, and Andreas Paepcke. Text summarization of web pages on handheld devices. In *NAACL'01,* 2001.

[29] N. H. M. Caldwell. An investigation into shallow processing for summarisation. Technical Report Computer science tripos part II project, University of Cambridge Computer Laboratory, 1994.

[30] Jaime G. Carbonell and Jade Goldstein. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of SIGIR,* pages 335–336, 1998.

[31] J. Carroll, G. Minnen, Y. Canning, S. Devlin, and J. Tait. Practical simplification of english newspaper text to assist aphasic readers. In *AAAI-98 Workshop on Integrating Artificial Intelligence and Assistive Technology,* 1998.

[32] Y. Chali, M. Kolla, N. Singh, and Z. Zhang. The university of lethbridge text summarizer at DUC 2003. In *HLT/NAACL Workshop on Text Summarization / DUC 2003,* 2003.

[33] Hsin-Hsi Chen. Multilingual summarization and question answering. In *Workshop on Multilingual Summarization and Question Answering (COLING'2002),* 2002.

[34] Hsin-Hsi Chen, June-Jei Kuo, and Tsei-Chun Su. Clustering and visualization in a multi-lingual multi-document summarization system. In *Proceedings of the 25th European Conference on IR Research,* pages 266–280, 2003.

[35] Hsin-Hsi Chen and Chuan-Jie Lin. A multilingual news summarizer. In *Proceedings of 18th International Conference on Computational Linguistics, COLING 2000,* pages 159–165, 2000.

[36] John M. Conroy and Dianne P. O'Leary. Text summarization via Hidden Markov Models. In *SIGIR 2001,* 2001.

[37] John M. Conroy, Judith D. Schlesinger, Dianne P. O'Leary, and Mary Ellen Okurowski. Using HMM and Logistic Regression to generate extract summaries for DUC. In *Workshop on Text Summarization in conjunction with the ACM SIGIR Conference 2001,* New Orleans, Louisiana, 2001.

[38] T. Copeck, S. Szpakowicz, and N. Japkowic. Learning how best to summarize. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization),* Philadelphia, July, 11-12 2002.

[39] Simon H. Corston-Oliver. Text compaction for display on very small screens. In *NAACL'01,* 2001.

[40] R. E. Cullingford. SAM. In Schank and Riesbeck, editors, *Inside Computer Understanding.* Lawrence Erlbaum Assoc., Hillsdale, NJ, 1981.

[41] H. Daumé III, A. Echihabi, D. Marcu, D.S. Munteanu, and R. Soricut. GLEANS: A generator of logical extracts and abstracts for nice summaries. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization),* Philadelphia, July, 11-12 2002.

[42] Hal Daumé III and Daniel Marcu. A noisy-channel model for document compression. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics,* 2002.

[43] G. DeJong. An overview of the frump system. In W. G. Lehnert and M. H. Ringle, editors, *Strategies for natural language processing,* pages 149 – 176. Hillsdale, NJ: Lawrence Erlbaum, 1982.

[44] J. Dersy. Producing summary content indicators for retrieved texts. Master's thesis, University of Cambridge Department of Engineering, 1996.

[45] DUC. DUC–document understanding conference. http://duc.nist.gov/.

[46] Daniel M. Dunlavy, John M. Conroy, Judith D. Schlesinger, Sarah A. Goodman, Mary Ellen Okurowski, Dianne P. O'Leary, and Hans van Halteren. Performance of a three-stage system for multi-document summarization. In *DUC03*, Edmonton, Alberta, Canada, May 31 - June 1 2003. Association for Computational Linguistics.

[47] H. P. Edmunson. New methods in automatic extracting. *Journal of the Association for Computing Machinery*, 16(2):264 – 285, April 1969.

[48] Noemie Elhadad and Kathleen R. McKeown. Towards generating patient specific summaries of medical articles. In *NAACL'01 Automatic Summarization Workshop*, 2001.

[49] David Kirk Evans and Judith L. Klavans. A platform for multilingual news summarization. Technical Report CUCS-014-03, Computer Science, University of Columbia, 2003.

[50] A. Farzindar, G. Lapalme, and H. Saggion. Summaries with SumUM and its expansion for document understanding conference (DUC 2002). In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[51] Atefeh Farzindar and Guy Lapalme. Using background information for multi-document summarization and summaries in response to a question. In *DUC03*, Edmonton, Alberta, Canada, May 31 - June 1 2003. Association for Computational Linguistics.

[52] Maria Fuentes, Marc Massot, Horacio Rodríguez, and Laura Alonso. Mixed approach to headline extraction for DUC 2003. In *HLT/NAACL Workshop on Text Summarization / DUC 2003*, Edmonton, Canada, 2003.

[53] Maria Fuentes and Horacio Rodríguez. Using cohesive properties of text for automatic summarization. In *JOTRI'02*, 2002.

[54] P. Gladwin, S. Pulman, and K. Sparck-Jones. Shallow processing and automatic summarising: a first study. Technical Report 223, University of Cambridge Computer Laboratory, 1991.

[55] Jade Goldstein, Vibhu Mittal, Mark Kantrowitz, and Jaime Carbonell. Summarizing text documents: Sentence selection and evaluation metrics. In *SIGIR-99*, 1999.

[56] Stephen J. Green. *Automatically generating hypertext by computing semantic similarity*. PhD thesis, University of Toronto, 1997.

[57] U. Hahn. Topic parsing: Accounting for text macro structures in full-text analysis. *Information Processing and Management*, 26(1):135–170, 1990.

[58] Udo Hahn and Inderjeet Mani. The cahllenges of automatic summarization. *IEEE Computer*, 33(11):29–36, 2000.

[59] S.M. Harabagiu and F. Lacatusu. Generating single and multi-document summaries with GISTEXTER. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[60] V. Hatzivassiloglou, J. Klavans, M. Holcombe, R. Barzilay, M.Y. Kan, and K.R. McKeown. Simfinder: A flexible clustering tool for summarization. In *NAACL'01 Automatic Summarization Workshop*, 2001.

[61] Vassileios Hatzivassiloglou, Judith Klavans, and Eleazar Eskin. Detecting text similarity over short passages: Exploring linguistic feature combinations via machine learning. In *EMNLP/VLC'99*, Maryland, 1999.

[62] A. G. Hauptmann and M. J. Witbrock. Informedia: News-on-demand multimedia information acquisition and retrieval. In M. Maybury, editor, *Intelligent Multimedia Information Retrieval*, pages 215–239. AAAI/MIT Press, 1997.

[63] Marti Hearst. Multi-paragraph segmentation of expository text. In *32nd Annual Meeting of Association for Computational Linguistics*, 1994.

[64] Ulf Hermjakob. *Learning Parse and Translation Decisions From Examples With Rich Context*. PhD thesis, University of Texas at Austin, 1997.

[65] T. Hirao, Y. Sasaki, H. Isozaki, and E. Maeda. NTT's Text Summarization system for DUC-2002. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[66] T. Hirao, J. Suzuki, H. Isozake, and E. Maeda. NTT's multiple document summarization system for DUC2003. In *HLT/NAACL Workshop on Text Summarization / DUC 2003*, 2003.

[67] Michael Hoey. *Patterns of Lexis in Text. Describing English Language.* Oxford University Press, 1991.

[68] Eduard Hovy. *Handbook of Computational Linguistics*, chapter 28: Text Summarization. Oxford University Press, 2001.

[69] Eduard Hovy and Chin-Yew Lin. Automated Text Summarization in SUMMARIST. In Mani and Maybury, editors, *Advances in Automatic Text Summarization*. 1999.

[70] Eduard Hovy and Daniel Marcu. Automated Text Summarization. COLING-ACL, 1998. tutorial.

[71] Hongyan Jing. Sentence simplification in automatic text summarization. In *ANLP-2000*, 2000.

[72] Hongyan Jing. *Cut-and-Paste Text Summarization.* PhD thesis, Graduate School of Arts and Sciences, Columbia University, 2001.

[73] Hongyan Jing and Kathleen McKeown. Cut and paste based text summarization. In *1st Conference of the North American Chapter of the Association for Computational Linguistics*, 2000.

[74] Min-Yen Kan. *Automatic text summarization as applied to information retrieval: Using indicative and informative summaries.* PhD thesis, Columbia University, 2003.

[75] Min-Yen Kan, Judith L. Klavans, and Kathleen R. McKeown. Domain-specific informative and indicative summarization for information retrieval. In *Workshop on Text Summarization in conjunction with the ACM SIGIR Conference 2001*, New Orleans, 2001.

[76] Min-Yen Kan and Kathleen McKeown. Information extraction and summarization: Domain independence through focus types. Technical report, Computer Science Department, Columbia University, New York, 1999.

[77] M. Karamuftuoglu. An approach to summarization based on lexical bonds. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[78] Kevin Knight and Daniel Marcu. Statistics-based summarization - step one: Sentence compression. In *The 17th National Conference of the American Association for Artificial Intelligence AAAI'2000*, Austin, Texas, 2000.

[79] Hideki Kozima. Text segmentation based on similarity between words. In *Proceedings of the 31th Annual Meeting of the Association for Computational Linguistics*, pages 286–288, 1993.

[80] W. Kraaij, M. Spitters, and A. Hulth. Headline extraction based on a combination of uni- and multidocument summarization techniques. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[81] W. Kraaij, M. Spitters, and M. van der Heijden. Combining a mixture language model and naive bayes for multi-document summarisation. In *Workshop on Text Summarization in conjunction with the ACM SIGIR Conference 2001*, New Orleans, Louisiana, 2001.

[82] Julian Kupiec, Jan Pedersen, and Francine Chen. A trainable document summarizer. In *Proceedings of ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 68–73. ACM Press, 1995.

[83] Finley Lacatusu, Paul Parker, and Sanda Harabagiu. Lite-GISTexter: Generating short summaries with minimal resources. In *DUC03*, Edmonton, Alberta, Canada, May 31 - June 1 2003. Association for Computational Linguistics.

[84] P. Lal and S. Rueger. Extract-based summarization with simplification. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[85] Abderrafih Lehmam. Text structuration leading to an automatic summary system: Rafi. *Information Processing and Management*, 35(2):181–191, 1999.

[86] Abderrafih Lehmam and Philippe Bouvet. Évaluation, rectification et pertinence du résumé automatique de texte pour une utilisation en réseau. In S. Chaudiron and C. Fluhr, editors, *III Colloque d'ISKO-France: Filtrage et résumé automatique de l'information sur les réseaux*, 2001.

[87] W. G. Lehnert. Plot units: a narrative summarization strategy. In W. G. Lehnert and M. H. Ringle, editors, *Strategies for natural language processing*, pages 375 – 412. Hillsdale, NJ: Lawrence Erlbaum, 1982.

[88] Anton Leuski, Chin-Yew Lin, and Eduard Hovy. iNeATS: Interactive multi-document summarization. In *ACL'03*, 2003.

[89] C-Y. Lin. *Robust Automated Topic Identification*. PhD thesis, University of Southern California, 1997.

[90] Chin-Yew Lin. Training a selection function for extraction. In *ACM-CIKM*, pages 55–62, 1999.

[91] Chin-Yew Lin and Eduard Hovy. Identifying topics by position. In *Proceedings of the Applied Natural Language Processing Conference (ANLP-97)*, pages 283–290, Washington, DC, 1997.

[92] Chin-Yew Lin and Eduard Hovy. NeATS: A multidocument summarizer. In *Workshop on Text Summarization in conjunction with the ACM SIGIR Conference 2001*, New Orleans, 2001.

[93] Chin-Yew Lin and Eduard Hovy. NeATS in DUC 2002. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[94] Chin-Yew Lin and Eduard Hovy. Automatic evaluation of summaries using n-gram co-occurrence statistics. In Marti Hearst and Mari Ostendorf, editors, *HLT-NAACL 2003: Main Proceedings*, pages 150–157, Edmonton, Alberta, Canada, May 27 - June 1 2003. Association for Computational Linguistics.

[95] Chin-Yew Lin and Eduard H. Hovy. The automated acquisition of topic signatures for Text Summarization. In *COLING-00*, Saarbrücken, 2000.

[96] H. P. Luhn. The automatic creation of literature abstracts. *IBM Journal of research and development*, 2(2):159 – 165, 1958.

[97] I. Mani, D. House, G. Klein, L. Hirschman, L. Obrst, T. Firmin, M. Chrzanowski, and B. Sundheim. The tipster SUMMAC text summarization evaluation: Final report. Technical report, DARPA, 1998.

[98] Inderjeet Mani. *Automatic Summarization*. Nautral Language Processing. John Benjamins Publishing Company, 2001.

[99] Inderjeet Mani and Eric Bloedorn. Machine learning of generic and user-focused summarization. In *AAAI*, pages 821–826, 1998.

[100] Inderjeet Mani and Eric Bloedorn. Summarizing similarities and differences among related documents. *Information Retrieval*, 1(1-2):35–67, 1999.

[101] Inderjeet Mani and Mark T. Maybury, editors. *Advances in automatic text summarisation*. MIT Press, 1999.

[102] William C. Mann and Sandra A. Thompson. Rhetorical structure theory: Toward a functional theory of text organisation. *Text*, 3(8):234–281, 1988.

[103] Daniel Marcu. From discourse structures to text summaries. In Mani and Maybury, editors, *Advances in Automatic Text Summarization*, pages 82 – 88, 1997.

[104] M. Maybury and A. Merlino. Multimedia summaries of broadcast news. In *International Conference on Intelligent Information Systems*, 1997.

[105] Mark T. Maybury and Inderjeet Mani. Automatic summarization. ACL/EACL'01, 2001. tutorial.

[106] K. McKeown, R. Barzilay, D. Evans, V. Hatzivassiloglou, J. Klavans, C. Sable, B. Schiffman, and S. Sigelman. Tracking and summarizing news on a daily basis with Columbia's Newsblaster. In *Proceedings of the Human Language Technology Conference*, 2002.

[107] K. McKeown, S.-F. Chang, J. Cimino, S. Feiner, C. Friedman, L. Gravano, V. Hatzivassiloglou, S. Johnson, D. Jordan, J. Klavans, A. Kushniruk, V. Patel, and S. Teufel. Persival, a system for personalized search and summarization over multimedia healthcare information. In *ACM+IEEE Joint Conference on Digital Libraries (JCDL 2001)*, 2001.

[108] K. McKeown, D. Evans, A. Nenkova, R. Barzilay, V. Hatzivassiloglou, B. Schiffman, S. Blair-Goldensohn, J. Klavans, and S. Sigelman. The columbia multi-document summarizer for DUC 2002. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[109] Kathleen McKeown, Judith Klavans, Vassileios Hatzivassiloglou, Regina Barzilay, and Eleazar Eskin. Towards multidocument summarization by reformulation: Progress and prospects. In *AAAI 99*, 1999.

[110] Kathleen R. McKeown, Regina Barzilay, David Evans, Vasileios Hatzivassiloglou, Min-Yen Kan, Barry Schiffman, and Simone Teufel. Columbia multi-document summarization: Approach and evaluation. In *Proceedings of the Workshop on Text Summarization, ACM SIGIR Conference*, 2001.

[111] Kathleen R. McKeown and Dragomir R. Radev. Generating summaries of multiple news articles. In *ACM Conference on Research and Development in Information Retrieval SIGIR'95*, Seattle, WA, 1995.

[112] Jean-Luc Minel, Jean-Pierre Desclés, Emmanuel Cartier, Gustavo Crispino, Slim Ben Hazez, and Agata Jackiewicz. Résumé automatique par filtrage sémantique d'informations dans des textes. présentation de la plate-forme filtext. *Revue Technique et Science Informatique*, 2001.

[113] M. Mitra, A. Singhal, and C. Buckley. Automatic Text Summarization by paragraph extraction. In Inderjeet Mani and Mark Maybury, editors, *Intelligent Scalable Text Summarization Workshop (ISTS'97)*, pages 39 – 46, Madrid, 1997. ACL/EACL.

[114] V. Mittal, M. Kantrowitz, J. Goldstein, and J. Carbonell. Selecting text spans for document summaries: Heuristics and metrics. In *AAAI 1999*, 1999.

[115] Vibhu Mittal and Adam Berger. Query-relevant summarization using faqs. In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL 2000)*, Hong Kong, 2000.

[116] Jane Morris and Graeme Hirst. Lexical cohesion, the thesaurus, and the structure of text. *Computational linguistics*, 17(1):21–48, 1991.

[117] S. Muresan, E. Tzoukermann, and J. Klavans. Combining linguistic and machine learning techniques for email summarization. In *ACL-EACL'01 CoNLL Workshop*, 2001.

[118] Sung Hyon Myaeng and Myung-Gil Jang. Integrating digital libraries with cross-language ir. In *Proceedings of the 2nd Conference on Digital Libraries*, 1999.

[119] Ani Nenkova, Barry Schiffman, Andrew Schlaiker, Sasha Blair-Goldensohn, Regina Barzilay, Sergey Sigelman, Vasileios Hatzivassiloglou, and Kathleen McKeown. Columbia at the duc 2003. In *DUC03*, Edmonton, Alberta, Canada, May 31 - June 1 2003. Association for Computational Linguistics.

[120] K. Ono, K. Sumita, and S. Miike. Abstract generation based on rhetorical structure extraction. In *Proceedings of the 15th International Conference on Computational Linguistics (COLING-94)*, pages 344 – 348, Kyoto, Japan, 1994.

[121] J.C. Otterbacher, A.J. Winkel, and D.R. Radev. The michigan single and multi-document summarizer for DUC 2002. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[122] Chris D. Paice. The automatic generation of literature abstracts: an approach based on the identification of self-indicating phrases. In R. N. Oddy, C. J. Rijsbergen, and P. W. Williams, editors, *Information Retrieval Research*, pages 172 – 191. London: Butterworths, 1981.

[123] Chris D. Paice. Constructing literature abstracts by computer. *Information Processing & Management*, 26(1):171 – 186, 1990.

[124] T.A.S. Pardo and L.H.M. Rino. DMSumm: Review and assessment. In E. Ranchhod and N. J. Mamede, editors, *Advances in Natural Language Processing*, pages 263–273. Springer-Verlag, 2002.

[125] T.A.S. Pardo, L.H.M. Rino, and M.G.V. Nunes. GistSumm: A summarization tool based on a new extractive method. In N.J. Mamede, J. Baptista, I. Trancoso, and M.G.V. Nunes, editors, *6th Workshop on Computational Processing of the Portuguese Language - Written and Spoken*, number 2721 in Lecture Notes in Artificial Intelligence, pages 210–218. Springer-Verlag, 2003.

[126] J. J. Pollock and A. Zamora. Automatic abstracting research at chemical abstracts service. *Journal of Information and Computer Sciences*, 15(4):226–23, 1975.

[127] K. Preston and S. Williams. Managing the information overload. physics in business. Institute of Physics, 1994.

[128] Dragomir Radev, Sasha Blair-Goldensohn, and Zhu Zhang. Experiments in single and multi-document summarization using MEAD. In *First Document Understanding Conference*, New Orleans, LA, September 2001.

[129] Dragomir Radev, Jahna Otterbacher, Hong Qi, and Daniel Tam. MEAD ReDUCs: Michigan at DUC 2003. In *DUC03*, Edmonton, Alberta, Canada, May 31 - June 1 2003. Association for Computational Linguistics.

[130] Dragomir R. Radev. Text Summarization. ACM SIGIR, 2000. tutorial.

[131] Dragomir R. Radev, Sasha Blair-Goldensohn, Zhu Zhang, and Revathi Sundara Raghavan. Interactive, domain-independent identification and summarization of topically related news articles. In *5th European Conference on Research and Advanced Technology for Digital Libraries*, Darmstadt, 2001.

[132] Dragomir R. Radev, Weiguo Fan, and Zhu Zhang. Webinessence: A personalized web-based multi-document summarization and recommendation system. In *NAACL Workshop on Automatic Summarization*, Pittsburgh, 2001.

[133] Dragomir R. Radev, Hongyan Jing, and Malgorzata Budzikowska. Centroid-based summarization of multiple documents: sentence extraction, utility-based evaluation, and user studies. In *ANLP/NAACL Workshop on Summarization*, Seattle, Washington, 2000.

[134] Dragomir R. Radev, Simone Teufel, Horacio Saggion, Wai Lam, John Blitzer, Arda Çelebi, Hong Qi, Elliott Drabek, and Danyu Liu. Evaluation of Text Summarization in a Cross-lingual Information Retrieval Framework. Technical report, Center for Language and Speech Processing, Johns Hopkins University, Baltimore, MD, June 2002.

[135] Lisa F. Rau, Paul S. Jacobs, and Uri Zernik. Information extraction and text summarisation using linguistic knowledge acquisition. *Information Processing & Management*, 25(4):419 – 428, 1989.

[136] RIPTIDES. RIPTIDES: Rapidly Portable Translingual Information Extraction and Interactive Multidocument Summarization. http://www.cs.cornell.edu/Info/People/cardie/tides/, 2002.

[137] J. E. Rush and et al. Automatic abstracting and indexing. ii. production of abstracts by application of contextual inference and syntactic coherence criteria. *Journal of the American Society for Information Science*, 22(4):260 – 274, 1971.

[138] Horacio Saggion and Guy Lapalme. Generating Indicative-Informative Summaries with SumUM. *Computational Linguistics*, 28(4), 2002.

[139] Horacio Saggion and Guy Lapalme. Generating informative and indicative summaries with SumUM. *Computational Linguistics*, 28(4), 2002. Special Issue on Automatic Summarization.

[140] Gerard Salton, James Allan, and Chris Buckley. Automatic structuring and retrival of large text files. *CACM*, 37(2):97–108, 1994.

[141] Gerard Salton, Amit Singhal, M. Mitra, and C. Buckley. Automatic text structuring and summarization. *Information Processing and Management*, 33(3):193 – 207, 1997.

[142] R. Schank and R. Abelson. *Scripts, Plans, Goals, and Understanding*. Lawrence Erlbaum, Hillsdale, NJ, 1977.

[143] Barry Schiffman, Inderjeet Mani, and Kristian J. Concepcion. Producing biographical summaries: Combining linguistic knowledge with corpus statistics. In *EACL'01*, 2001.

[144] J.D. Schlesinger, J.M. Conroy, M.E. Okurowski, H.T. Wilson, D.P. O'Leary, A. Taylor, and J. Hobbs. Understanding machine performance in the context of human performance for multi-document summarization. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[145] E. F. Skorokhod'ko. Adaptive method of automatic abstracting and indexing. *Information processing*, 71, 1971.

[146] K. Sparck Jones, S. Walker, and S. Robertson. A probabilistic model of information retrieval: Development and status. Technical Report N 446, University of Cambridge Computer Laboratory, 1998.

[147] Karen Sparck-Jones. Automatic summarising: factors and directions. In Inderjeet Mani and Mark Maybury, editors, *Advances in Automatic Text Summarization*. MIT Press, 1999.

[148] Tomek Strzalkowski, Jin Wang, and Bowden Wise. A robust practical text summarization. In Eduard Hovy and Dragomir Radev, editors, *AAAI Spring Symposium on Intelligent Text Summarisation*, pages 26 – 33, Stanford, California, March 23-25 1998. American Association for Artificial Intelligence, AAAI Press.

[149] SUMMAC. SUMMAC, the final report. http://www.itl.nist.gov/iaui/894.02/related_projects/tipster_summac/, 1998.

[150] H. Sundaram. *Segmentation, Structure Detection and Summarization of Multimedia Sequences*. PhD thesis, Graduate School of Arts and Sciences, Columbia University, 2002.

[151] SweSum. http://www.nada.kth.se/ xmartin/swesum/index-eng.html, 2002.

[152] J. L. Tait. Automatic summarizing of english texts. Technical Report 47, University of Cambridge Computer Laboratory, 1983.

[153] S. L. Taylor. *Automatic abstracting by applying graphical techniques to semantic networks*. PhD thesis, Northwestern University, 1975.

[154] Simone Teufel and Marc Moens. Sentence extraction as a classification task. In Inderjeet Mani and Mark Maybury, editors, *Intelligent Scalable Text Summarization Workshop (ISTS'97)*, pages 58 – 59, Madrid, 1997. ACL/EACL.

[155] Simone Teufel and Marc Moens. Sentence extraction and rhetorical classification for flexible abstracts. In *AAAI Spring Symposium on Intelligent Text Summarisation*, pages 16 – 25, 1998.

[156] Simone Teufel and Marc Moens. Summarizing scientific articles – experiments with relevance and rhetorical status. *Computational Linguistics*, 28(4), 2002. Special Issue on Automatic Summarization.

[157] Richard Tucker. *Automatic Summarising and the CLASP system*. PhD thesis, University of Cambridge, 1999.

[158] E. Tzoukermann, S. Muresan, and J. Klavans. Gist-it: Summarizing email using linguistic knowledge and machine learning. In *ACL-EACL'01 HLT/KM Workshop*, 2001.

[159] H. van Halteren. Writing style recognition and sentence extraction. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[160] Piek Vossen, editor. *Euro WordNet: a multilingual database with lexical semantic networks*. Kluwer Academic Publishers, 1998.

[161] H. Wactlar. Multi-document summarization and visualization in the informedia digital video library, 2001.

[162] M. White, D. McCullough, C. Cardie, V. Ng, and K. Wagstaff. Detecting discrepancies and improving intelligibility: Two preliminary evaluations of riptides. In *Workshop on Text Summarization in conjunction with the ACM SIGIR Conference 2001*, New Orleans, 2001.

[163] Michael White and Claire Cardie. Selecting sentences for multidocument summaries using randomized local search. In *ACL Workshop on Automatic Summarization*, 2002.

[164] Michael White, Tanya Korelsky, Claire Cardie, Vincent Ng, David Pierce, and Kiri Wagstaff. Multi-document summarization via information extraction. In *Proceedings of the First International Conference on Human Language Technology Research*, 2001.

[165] M. Witbrock and V. Mittal. Ultra-summarization: A statistical approach to generating highly condensed nonextractive summaries. In *Proceedings of the 22nd International Conference on Research and Development in Information Retrieval (SIGIR-99)*, 1999.

[166] S. R. Young and P. J. Hayes. Automatic classification and summarisation of banking telexes. In *Second Conference on Artificial Intelligence Applications*, pages 402–408, New York, 1985.

[167] D. Zajic, B. Door, and R. Schwartz. Automatic headline generation for newspaper stories. In *Workshop on Text Summarization (In Conjunction with the ACL 2002 and including the DARPA/NIST sponsored DUC 2002 Meeting on Text Summarization)*, Philadelphia, July, 11-12 2002.

[168] Klaus Zechner. A literature survey on information extraction and Text Summarization. term paper, Carnegie Mellon University, 1997.

[169] Klaus Zechner. *Automatic Summarisation of Spoken Dialogues in Unrestricted Domains*. PhD thesis, Carnegie Mellon University, 2001.

| System | Processing Level | Information Kind | Tucker 1999 |
|---|---|---|---|
| Adam [137, 126] | surface | structural | sentencewise |
| Alfonseca and Rodríguez [1] | surface | structural | sentencewise |
| * Anes [26] | surface | lexical | att. networks |
| Barzilay and Elhadad 1997 [17] | entity | lexical | att. networks |
| Boguraev and Kennedy 1997 [24] | entity | lexical | att. networks |
| Caldwell 1994 [29] | entity | lexical | att. networks |
| * CENTRIFUSER [48] | discourse | understanding | info. content |
| * Chen and Lin (2000) [35] | surface | lexical | info. content |
| * Columbia MDS [108, 38, 119] | entity/discourse | understanding/structural | info. content |
| Copeck et al. 2002 [38] | surface | lexical | att. networks |
| * Cut-and-Paste [72] | surface | structural | info. content |
| Darsy 1996 [44] | entity | lexical | att. networks |
| * DiaSumm [169] | surface | lexical | discourse structure |
| DimSum [9] | surface | lexical | att. networks |
| * DMSumm [124] | discourse | structural | disc. structure |
| Edmunson 1969 [47] | surface | structural | sentencewise |
| FilText [112] | surface | structural | info. content |
| * FociSum [76] | entity | understanding | att. networks |
| Frump [43] | entity | understanding | info. content |
| GISTEXTER [59, 83] | discourse/entity | understanding | info. content |
| GISTSumm [125] | surface | lexical | att. networks |
| Gladwin et al. 1991 [54] | entity | lexical | att. networks |
| * GLEANS [41] | entity/discourse | understanding | info. content |
| * NTT [65, 66] | surface | structural/lexical | att. networks |
| * Karamuftuoglu 2002 [77] | surface | structural | att. networks |
| * Kraaij et al. 2002 [80] | surface | lexical | att. networks |
| K. U. Leuven [6, 7] | entity | lexical | att. networks |
| * Lal and Rueger 2002 [84] | entity/discourse | understanding | info. content |
| Lehnert 1982 [87] | entity | understanding | info. content |
| * Univ. of Lethbridge [27, 32] | entity | structural/lexical | att. networks |
| Luhn 1958 [96] | surface | lexical | att. networks |
| Marcu 1997 [103] | discourse | structural | disc. structure |
| * MEAD [128, 129] | surface | lexical | att. networks |
| * MultiGen [109, 19] | entity | structural | info. content |
| * NeATS [92, 93, 88] | entity | structural | info. content |
| * Newsblaster [106] | entity/discourse | structural/understanding | info. content |
| NewsInEssence [131] | surface | lexical | att. networks |
| Ono et al. 1994 [120] | discourse | structural | disc. structure |
| NetSumm [127] | surface | lexical | att. networks |
| Paice 1981 [122] | surface | structural | sentencewise |
| * PERSIVAL [107] | | understanding | info. content |
| Rafi [85] | surface | structural | att. networks |
| * RIPTIDES [136, 163] | entity/discourse | understanding | info. content |
| SAM [142, 40] | entity | understanding | info. content |
| Dunlavy et al. 2003 [144, 46] | surface | lexical | att. networks |
| Scisor [135] | entity | understanding | info. content |
| Scrabble [152] | entity | understanding | info. content |
| Skorochod'ko 1971 [145] | entity | lexical | att. networks |
| Smart [140, 113] | entity | lexical | att. networks |
| * SUMMARIST [69] | surface | lexical | att. networks |
| SUMMONS [111] | entity | understanding | info. content |
| SumUM [50, 138, 51] | discourse | structural | discourse structure |
| * SweSum [151] | surface | lexical | att. networks |
| Taylor 1975 [153] | entity | understanding | info. content |
| Tele-Pattan [20] | entity | lexical | att. networks |
| Tess [166] | entity | understanding | info. content |
| Teufel and Moens [155, 156] | discourse | structural | disc. structure |
| TICC [2] | entity | understanding | info. content |
| TOPIC [57] | discourse | structural | disc. structure |
| van Halteren 2002 [159] | surface | lexical | att. networks |
| WebInEssence [132, 167] | surface | lexical | att. networks |

Table 1: Classification of summarization systems

| On-line or Downloadable Demos | |
|---|---|
| Centrifuser<br><br>on-line demo | English<br>multi-document (specific-topic: medical documents)<br>http://centrifuser.cs.columbia.edu/centrifuser.cgi |
| Copernic<br><br>downloadable demo | English, French, German<br>single document (many formats)<br>http://www.copernic.com/desktop/products/summarizer/download.html |
| DMSumm<br><br>downoadable demo | English, Brazilian Portuguese<br>single document<br>http://www.nilc.icmc.usp.br/ thiago/DMSumm.zip |
| Extractor<br><br>downloadable demo | English, French, Spanish, German, Japanese, Korean<br>single document (many formats)<br>http://www.dbi-tech.com/dbi_extractor.asp |
| GISTexter<br><br>no straightforward access | English<br>Single and Multi-Document<br>form at: http://www.languagecomputer.com/demos/summarization/index.html |
| GistSumm<br><br>downloadable demo | multilingual<br>single document<br>http://www.nilc.icmc.usp.br/ thiago/Install_GistSum.zip |
| Newsblaster<br><br>on-line demo | Multilingual<br>multi-document<br>http://www1.cs.columbia.edu/nlp/newsblaster/ |
| Island InText<br><br>no straightforward downloading | English<br>single document<br>form at: http://www.islandsoft.com/orderform.html |
| Inxight Summarizer /<br>LinguistX / Xerox PARC<br><br>no straightforward downloading | Chinese, Danish, Dutch, English, Finnish, French, German, Italian,<br>Japanese, Korean, Norwegian, Portuguese, Spanish and Swedish<br>single document<br>form at: http://www.inxight.com/products/oem/summarizer/contact_sales.php |
| Kmaritime<br>on-line demo | Korean<br>http://nlplab.kmaritime.ac.kr/demo//f_ats.html |
| Lal and Rüger (2002)<br><br>on-line demo | English<br>single document<br>http://rowan.doc.ic.ac.uk:8180/summarizer/demo.html |
| MEAD / NewsInEssence / CLAIR<br><br>on-line and dowloadable demo | English and Chinese<br>multi-document, multi-lingual<br>http://www.clsp.jhu.edu/ws2001/groups/asmd/<br>multiple news summ. demo at: http://www.newsinessence.com/nie.cgi |
| MS-Word Autosummarize | supposedly any language<br>single document<br>included in MS-Word |
| Pertinence Summarizer<br><br><br>on-line demo | English, French, Spanish, German, Italian, Portuguese, Japanese,<br>Chinese, Korean, Arabic, Greek, Dutch, Norwegian and Russian<br>single document<br>http://www.pertinence.net |
| Sinope Summarizer Personal Edition<br><br>30-day trial downloadable | English, Dutch and German<br>single document<br>http://www.sinope.nl/en/sinope/index.html |
| Summ-It<br><br>on-line demo | probably English only<br>pasted text<br>http://www.mcs.surrey.ac.uk/SystemQ/summary/ |
| Surfboard<br>30-day trial downloadable demo | probably English only<br>single web pages (Mac OS X.1 only)<br>http://www.glu.com/binaries/surfboard/surfboard.dmg.gz |
| SweSum<br><br>on-line demo | Danish, English, French, German, Spanish, Swedish<br>single document (Web pages or pasted text)<br>http://www.nada.kth.se/ xmartin/swesum/index-eng.html |
| TextWise<br>Content Repurposing Suite<br>no straightforward access | probably English only<br>single document or e-mail<br>form at: http://www.textwise.com/technology/crs/demo.html |

Table 2: Some on-line demos of summarization systems, both commercial and academic
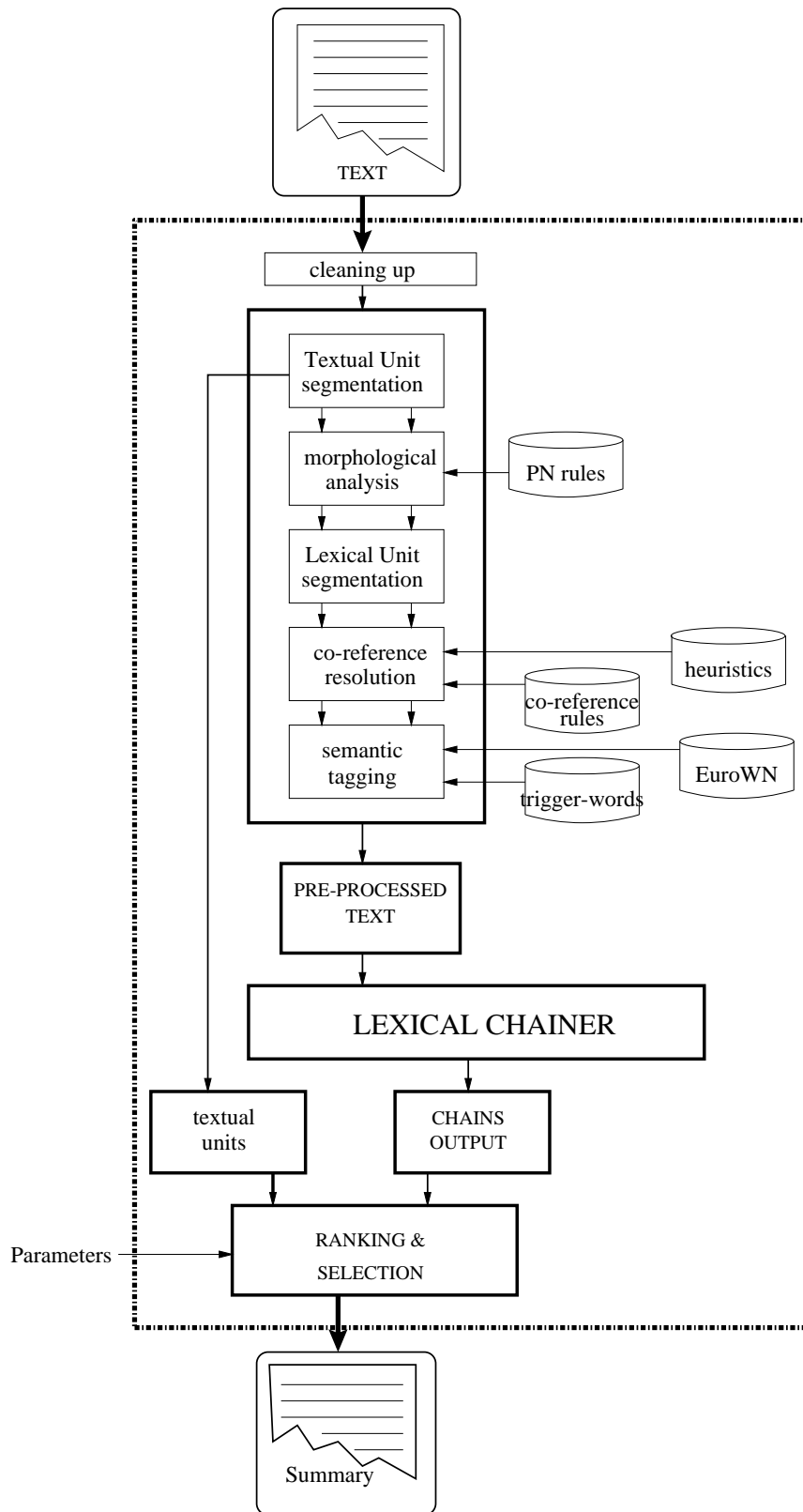
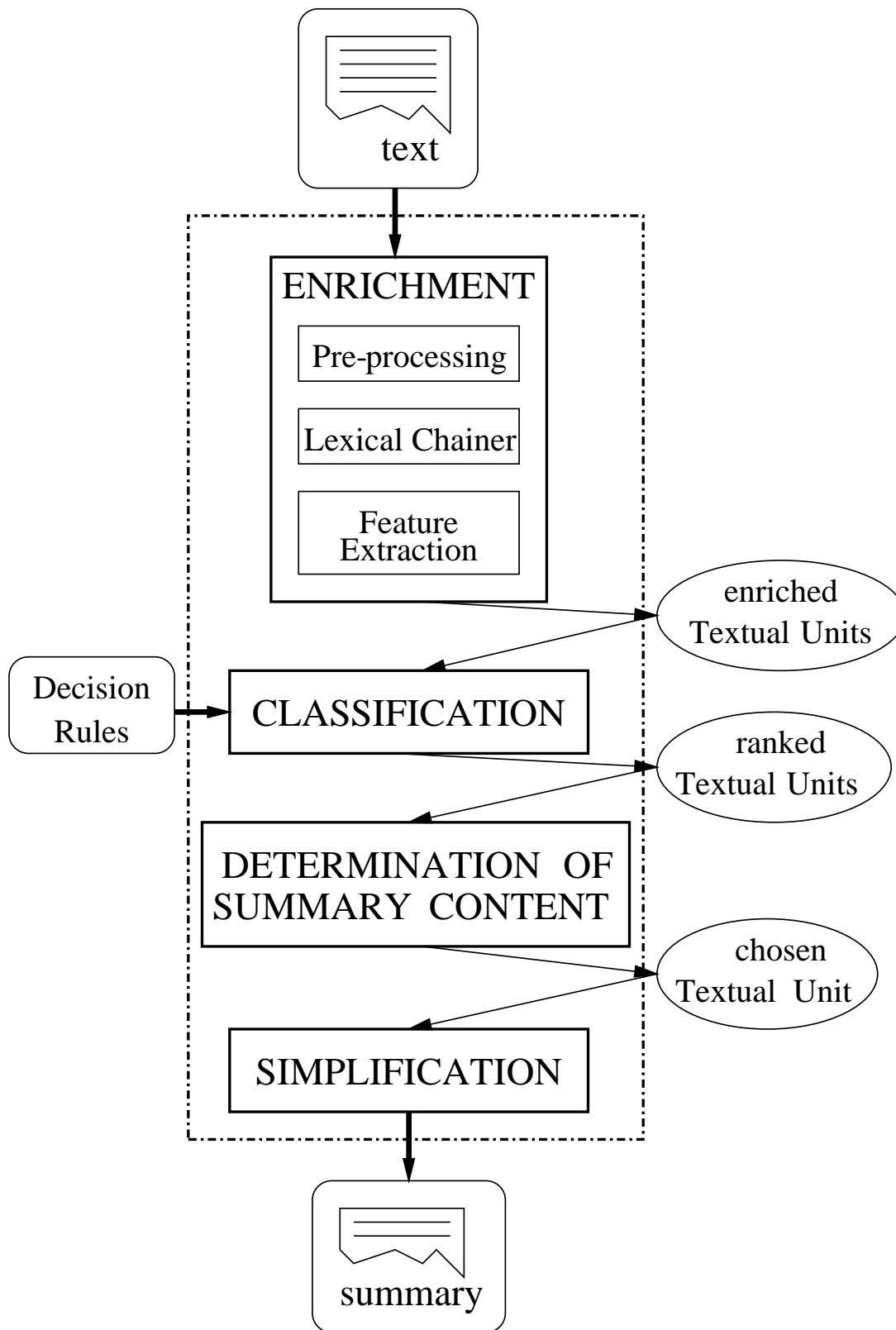Figure 1: Architecture of HERMES Lexical Chain Summarizer.

Figure 2: Architecture of HERMES Machine Learning Summarizer.

# Annex: Some Summarization Systems

## Alfonseca and Rodríguez 2003

- **Name:**

- **Reference:** [1]

- **Short description:** produces very short summaries (headline-like) of single documents applying genetic algorithms

- **System Features**

  - **Input:**

  - **Architecture:** The processing has two steps: first, the most relevant sentences of a document are extracted, applying a genetic algorithm that selects sentences according to their values for a series of features indicative of their relevance: sentence length, position in the document, order of the sentences, representativity, syntactical structure, redundancy. The algorithm is trained on the data from past DUC contests.

    Once the sentences are extracted, a headline is created by concatenation of portions of these sentences. To determine which portions should be extracted and which can be left aside, sentences are parsed, and hand-crafted rules are applied to guarantee well-formedness (extracting the main verb and its arguments) and informativity (extracting highly connected lexical items).

  - **Output facilities and constraints:**

  - **Language coverage:** English, potentially multilingual

- **Evaluation:** obtained average results (ranked in the middle of all systems) in DUC 2003

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): structural

  - within classification 3 (Tucker, 1999): sentence by sentence

- **Comments:**

# Baldwin and Morton 1998

- **Name**

- **Reference**: [13]

- **Short description**: Uses co reference between the query and the text for performing indicative, user-focused (query-sensitive) summarization

- **System Features**

    - **Input**:
    - **Architecture**: The system is based on a rich linguistics processing that includes the following tasks:
        * NER
        * Tokenization
        * Sentence segmentation
        * POS tagging
        * Morphological analysis
        * Parsing
        * Argument detection
        * Co-reference resolution: Identity and Part-Whole, including nominal and verbal phrases, acronyms, events
    - **Language coverage**: English
    - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

    - within classification 1 (level of processing): entity
    - within classification 2 (kind of information): lexical
    - within classification 3 (Tucker, 1999): sentence by sentence

- **Comments**:

## Banko et al. 1999, Mittal et al. 1999

- **Name**:

- **Reference**: [15], [114]

- **Short description**: Extraction-based summarization from hand-written summaries, i.e. going from abstracts to extracts, of single documents, by aligning text spans.

- **System Features**

  - **Input**:

  - **Architecture**: A *tl\*tf* (term length * term frequency) measure is used for weighting the relevance of terms and NE. [114] focuses on the selection of spans for document summaries. Sentences from the original document are ranked according to their salience using two parameters for tuning the process: i) granularity, e.g. paragraph, sentence, etc. and ii) metric for ranking. Features at discourse level include:
    * length of the span
    * density of NEs
    * complexity of NPs
    * punctuation
    * thematic phrases
    * anaphora density

    There are also features at subdocument level (sentence, phrase and word). These include:
    * word length
    * communicative actions
    * thematic phrases
    * use of honorifics, auxiliary verbs, negation, prepositions, etc.
    * type of sentence (interrogative, evaluative, etc.)

  - **Language coverage**: English

  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): discourse
  - within classification 2 (kind of information): understanding
  - within classification 3 (Tucker, 1999): informational content

- **Comments**: Related work includes Headline production [14, 22] and Ultrasummarization [165].

## Boros et al. 2001

- **Name:**

- **Reference:** [25]

- **Short description:** Multi-document summarization system

- **System Features**

  - **Input:**

  - **Architecture:** The system proceeds through the following steps i) From a document set a finite number of topics are extracted, ii) topics are ordered by importance, iii) a unique sentence is extracted from the collection for covering each topic; salience of sentences is computed using *tf*\**idf*, iv) sentences are clustered (several clustering techniques both hierarchical and non-hierarchical are experimented) and, finally, v) the summary is produced.

  - **Language coverage:** English, potentially multilingual

  - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): lexical

  - within classification 3 (Tucker, 1999): informational content

- **Comments:**

# Carbonell and Goldstein 1998, Goldstein et al. 1999

- **Name:**

- **Reference:** [30], [55]

- **Short description:** CMU approach to both SDS and MDS combines criteria of query relevance and novelty.

- **System Features**

  - **Input:**

  - **Architecture:** The base of the system is the MMR (Maximal Marginal Relevance) metric. Important issues are the diversity-based re-ranking for reordering documents (in MDS), the relevant passage extraction, the anti-redundancy measures, the way of combining criteria of relevance and novelty (relevant novelty vs. declining relevance to users"s query). In the case of SDS the system ranks sentences from the original document according to their salience or their likelihood of being part of the summary. For doing so, a weighted score of both linguistic and statistical features is used. The weights are optimised according to application genres. Among linguistic features we can find: name, place, honorifics, quotations, thematic phrases, etc. Statistical features include cosine, *tf\*idf*, pseudo-relevance feedback, query expansion, user interest profiles, etc. In the case of MDS different types of summaries can be produced using:
    * Common sections of documents
    * Common sections + unique sections of documents
    * Centroid
    * Centroid + outliers
    * Common sections + unique sections + time weighting factor

  - **Language coverage:** English, potentially multilingual

  - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): informational content

- **Comments:**

# CENTRIFUSER

- **Name**: CENTRIFUSER

- **Reference**: [75]

- **Short description**: Multi-document Summarizer. CENTRIFUSER meets the needs of browsers and searchers in highly structured domains.

- **System Features**

  - **Input**:

  - **Architecture**: The system uses SIMFINDER [61, 60], a flexible clustering tool for summarization (used also in MULTIGEN). This tool detects text similarity over short passages exploring linguistic features combinations via Machine Learning techniques. Among the primitive linguistic features we can find word co-occurrence, shared proper nouns, linked noun phrases, WN synonyms and semantically similar verbs. Composite features consist of pairs of simple features. An automatic feature detection system is applied and then the well-known ILP system, RIPPER, is performed. After clustering, the system uses key-terms for selecting one sentence or paragraph from each cluster (using the centroid method of [133]). The selected sentences are finally reordered by reformulation (in a similar way as in MULTIGEN).

  - **Language coverage**: English, parts of it potentially multilingual

  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): discourse

  - within classification 2 (kind of information): understanding

  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

# Chen and Lin (2000), Chen et al. 2003

- **Name:**

- **Reference:** [35, 34]

- **Short description:**

- **System Features**

  - **Input:** multidocument
  - **Architecture:** The main components of the system are a set of monolingual news clusterers, a unique multilingual news clusterer and a news summarizer. A central issue of the system is the definition and identification of meaningful units as base for comparison. For English these units can be reduced to sentences but for Chinese the identification of units and the associated segmentation of the text can be a difficult task. Another important issue of the system (general for systems covering distant languages or different encoding schemata) is the need of a robust transliteration of names (or words not occurring in the bilingual dictionary) for assuring an accurate matching.
  - **Output facilities and constraints:**
  - **Language coverage:** crosslingual: English and Chinese, potentially any language

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): informational content

- **Comments:**

# Columbia MDS

- **Name**: Columbia MDS

- **Reference**: [19, 18, 110, 61, 60, 108, 119]

- **Short description**: Enhanced version of MULTIGEN. Complex system that can be applied to different sources. It can be considered a sort of meta-summarizer.

- **System Features**

  - **Input**: Four different types of input that are identified in a way that the most appropriate summarizer is applied in each case. The system can deal with simple event, biography, multi-event and others.

  - **Architecture**: There is a pre-processing phase followed by a router that depending on the kind of input triggers the appropriate summarizer. For simple events the summarizer used is the conventional MULTIGEN, for biographies, DEMS [143] with the bio configuration, for multi-event and others, DEMS with the default configuration.

  - **Language coverage**: English

  - **Output facilities and constraints**:

- **Evaluation**: DUC 2002, consistently among the top three systems (second or third). For extracts, it ranked second precisionwise and third recallwise. For abstracts, it ranked second coveragewise and third precisionwise. Also participated in DUC 2003, and obtained good results for coverage and quality questions in some of the tasks.

- **Classification**

  - within classification 1 (level of processing): entity/discourse

  - within classification 2 (kind of information): structural/understanding

  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

# Conroy et al. 2001

- **Name:**

- **Reference:** [37, 36, 144, 46]

- **Short description:** Statistical approaches to summarisation

- **System Features**

  - **Input:**
  - **Architecture:** Two different techniques are used in [37]:
    * HMM, using as features the position in the sentence, the number of tokens and the number of pseudo-query terms.
    * Logistic Regression (LRM), using as features the number of query terms occurring in the sentence, the number of tokens (sentence length), the distance to the query terms and the position of the sentence.

    [36] use pivoted GR matrix decomposition. A token-sentence matrix is built and from it the columns giving good coverage of the tokens are selected. Two different approaches are used for this process: a greedy election and a pivoted QR factorisation. [144] merged the LRM and HMM by including all the features of the LRM in the HMM. An additional feature was the conditional probability that a sentence is a summary sentence given that the previous sentence is. A post-process is run on extracted sentences to remove sentence starting discourse markers and boilerplate, to improve cohesiveness. An extensive investigation was carried out to account for human performance in multi-document summarization. Conclusions were that single-document summaries could be used as a base for multi-document, but had to be enriched, possibly wiht discourse structure. Sentence pruning techniques were also found useful.

  - **Language coverage:** English, potentially multilingual
  - **Output facilities and constraints:**

- **Evaluation:** participated in DUC'01, DUC'02 and DUC'03. In DUC'02, it was ranked among the first systems, but did not beat the baselines. In DUC'03, it performed among the top systems.

- **Classification**

  - within classification 1 (level of processing): surface/entity
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): informational content

- **Comments:**

## Cut-and-Paste

- **Name**: Cut-and-Paste

- **Reference**: [71], [73]

- **Short description**: Sentence Reduction for automatic text summarization. The system relates the phrases occurring in a summary written by a professional summarizer and the phrases occurring in the original document.

- **System Features**

  - **Input**:

  - **Architecture**: 6 editing operations (learned from the performance of human summarizers) are used for sentence reduction:
    * removing extraneous phrases
    * combining a reduced sentence with other reduced sentences
    * syntactic transformations
    * substitution with paraphrases
    * substitution with more general or more specific descriptors
    * reordering

  - **Language coverage**: English

  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): structural

  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## DiaSumm

- **Name**: DiaSumm

- **Reference**: [169]

- **Short description**: Automatic Summarization of Spoken Dialogues in Unrestricted Domains

- **System Features**: Dealing with non textual documents implies that additional problems have to be faced. If the input comes from ASR (with or without confidence scores), speech disfluencies have to be detected and removed. Besides, sentence boundaries have to be detected and inserted. Topic segmentation plays a more important role in this situation. In addition, in the case of multi-party dialogs, relations between moves have to be identified (e.g. linking of question/answering pairs).

  - **Input**: Spoken dialogues
  - **Architecture**: DiaSumm is organised in the following modules:
    1. speech disfluency detection and removal
    2. identification and insertion of sentence boundaries
    3. identification and linking of Question-Answer regions
    4. topical segmentation
    5. information condensation (using MMR)
  - **Language coverage**: English, German
  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): discourse structure

- **Comments**:

## DMSumm

- **Name**: DMSumm (Discourse Modeling SUMMarizer)

- **Reference**: [124]

- **Short description**: a three-layered discourse-based summarizer

- **System Features**

  - **Input**: single document
  - **Architecture**: DMSumm is a deep approach to the summarization problem. It has three steps: content selection, text planning and linguistic realization. The content selection process select the information to be communicated in the summaries; the text planning makes a mapping of semantic and intentional relations onto rhetorical relations, building rhetorical text plans; the linguistic realization expresses the plans in the written summaries. It is based on a discourse model composed of three different knowledge sources, i.e., the semantic, intentional and rhetorical levels. Some basic generation restrictions are supposed to be verified: the communicative goal satisfaction and the central proposition preservation.
  - **Output facilities and constraints**:
  - **Language coverage**: English and Brazilian Portuguese

- **Classification**

  - within classification 1 (level of processing): discourse
  - within classification 2 (kind of information): structural
  - within classification 3 (Tucker, 1999): discourse structure

- **Evaluation**:

- **Comments**:

# eSseNSe, NewsInESSence, WebInESSence

- **Name**: eSseNSe, NewsInESSence, WebInESSence

- **Reference**: [131], [132]

- **Short description**: eSseNSe is basically a system for clustering documents after/before retrieval, summarization single/multi-document, personalization and recommendation of documents. From it two systems applied respectively to news (NewsInESSence) and Web pages (WebInESSence) have been derived.

- **System Features**

    - **Input**:
    - **Architecture**: These systems are based on the CST (Cross-Document Structure Theory). CST (that is related to RST for single documents) proposes a taxonomy of the informational relationships between documents in clusters of related documents. In NewsInESSence the aim is finding, visualizing and summarizing a topic-based cluster of news stories. A user selects a single news story from a news Web site. The system searches for other live sources of news for other stories related to this one and presents summaries
    - **Language coverage**: English, potentially multilingual
    - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

    - within classification 1 (level of processing): surface
    - within classification 2 (kind of information): lexical
    - within classification 3 (Tucker, 1999): attentional networks

- **Comments**:

## FociSum

- **Name**: FociSum

- **Reference**: [76], [75]

- **Short description**: Summarizing long documents. Domain specific informative and indicative summarization for Information Retrieval. Closely related to CENTRIFUSER.

- **System Features**

    - **Input**:

    - **Architecture**: Summarization of long documents presents interesting characteristics that do not occurs in conventional summarization systems (usually applied to summarize news, articles, Web pages and so). In long documents summarization sentences to be extracted occurs in distant locations. So, coherence properties are of less importance here. Focisum is an hybrid system that merges: i) Information Extraction techniques (template-based), ii) Sentence extraction (including both sentence-based and lead-based strategies) and iii) based on the dynamically determined foci of the text (in this context focus is the topic). Foci are built from NE and multiword terms.

    - **Language coverage**: English

    - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

    - within classification 1 (level of processing): entity

    - within classification 2 (kind of information): understanding

    - within classification 3 (Tucker, 1999): attentional networks

- **Comments**:

# GISTexter

- **Name**: GISTexter

- **Reference**: [59, 83]

- **Short description**: produces multidocument extracts and abstracts by template-driven IE. Templates are chosen by their adequacy to the topic of the document or collection of documents. Single document summaries by sentence extraction and compression.

- **System Features**

    - **Input**: collections of documents dealing with the same topic.
    - **Architecture**: for single documents, the most relevant sentences are extracted and compressed by rules that are learned from a corpus of human-written abstracts and their source texts (no further detail of these processes is given). For multi-document summarization, the system:
        * the IE system CICERO extracts relevant information by applying templates that are determined by the topic of the collection. Each template keeps a record of the text snippets where the information has been extracted from. If one of these snippets contains an anaphoric element, its co-reference chain is also recorded. If no template is provided for a given topic, a template is generated ad-hoc, based on the topical relations of the words in WordNet.
        * the *dominant event* of the collection is determined, and templates are classed depending on how central the dominant event is in the template and in the document where the template is extracted from.
        * within each class, templates are ordered by their representativeness. Highly representative templates are those that have the same slot fillers in the same slots as the majority of templates. Also those templates related to text snippets crossed by co-reference chains are more representative.
        * the summary is made from the text snippets recorded by the most representative template in the class of templates most closely related to the dominant event in the collection, in their order of appearance in the text. If they contain an anaphoric element, sentences containing the antecedent are also included. If the summary is too long, the linguistic form of dates and locations is shortened, unimportant coordinated phrases are dropped or, finally, the last sentence is dropped until the targeted length is achieved. If the summary is too short, the same process is applied to the most representative templates to the other classes of templates, in order of closeness to the dominant event.
    - **Language coverage**: English
    - **Output facilities and constraints**:

- **Evaluation**: participated in DUC 2002 and was ranked among the first. The best coverage rates for single and multi-document summarization, only surpassed by one system as to precision in multi-document summarization. In DUC 2003 they participated with Lite-Gistexter, which uses minimal lexico-semantic resources, obtaining good results for one of the four tasks.

- **Classification**

    - within classification 1 (level of processing): entity/discourse
    - within classification 2 (kind of information): understanding
    - within classification 3 (Tucker, 1999): informational content

- **Comments**: the mentioned reference does not provide much detail on some of the modules of the system.

# GISTSumm

- **Name:**

- **Reference:** [125]

- **Short description:** an automaitc text summarizer that tries to identify the text main idea, i.e., the gist, for generating the corresponding summary.

- **System Features**

  - **Input:**

  - **Architecture:** It is based in the assumptions that it is possible to:

    * find a sentence that represents the main idea of a text, the gist.
    * find the gist by statistical methods.
    * produce coherent abstracts relating the gist with other sentences of the original text

    It has two methods to summarize: via key words or via a metric to find the most representative words of a text (*tf\*isf*, term frequency - inverse sentence frequency).

  - **Output facilities and constraints:**

  - **Language coverage:** multilingual

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): lexical

  - within classification 3 (Tucker, 1999): attentional networks

- **Evaluation:**

- **Evaluation:**

- **Comments:**

## GLEANS

- **Name**: GLEANS

- **Reference**: [41]

- **Short description**: IE-based multi-document summarizer, makes explicit the main entities and relations in a document collection. It produces headlines, extracts and a reduced form of abstract.

- **System Features**

  - **Input**:
  - **Architecture**: summarization in four steps:
    * documents are parsed [64], the main constituents of each sentence are identified, some anaphoric expressions are resolved, and finally mapped into a canonical representation that explicits their main entities and relations
    * each collection of documents is classified by its content into *person, single event, multiple event* or *natural disaster*
    * given the collection type and the canonical representation of the documents, the core entities and relations are extracted, by choosing the most salient words in the collection.
    * a headline is created, based on the type of collection and teh core entities and relations. For *multiple event* collections, a short abstract can also be generated with the mechanisms to generate headlines.
    * an abstract is generated by applying a library of canonical schemas obtained from manual analysis of abstracts in a training corpus. These schemas determine which sentences of a source text fulfill the requirements of a canonical summary, and extract them. Chronological coherence, redundancy and dangling discourse references are treated.
    * in a post-process, dangling discourse markers are removed, decisions are made on which anaphoric expressions to use for each entity and temporal expressions are represented in a canonical form.
  - **Language coverage**: English
  - **Output facilities and constraints**:

- **Evaluation**: performance in DUC 2002 not high: low coverage, but improved when document collections were correctly classified. Specially bad on headline generation.

- **Classification**

  - within classification 1 (level of processing): entity/discourse
  - within classification 2 (kind of information): understanding
  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

# Knight and Marcu 2000

- **Name:**

- **Reference:** [78]

- **Short description:** This system is not a full summarizer but a sentence compressor. Sentence compressing is presented as a fundamental component of any high-quality non extractive summarizer

- **System Features**

  - **Input:**

  - **Architecture:** The system follows a statistical approach. Sentence compression is considered as a process of translation from a source language (full text) into a target language (summary). The process is accomplished following two different approaches: a conventional noise channel model and decision trees (using C4.5). The probabilistic models are trained on a corpus of ¡full text, summary¿ pairs.

  - **Language coverage:** English, potentially multilingual

  - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): sentence by sentence

- **Comments:** an enhancement of this approach was carried out later on, applying the same technique to rhetorical parse trees, with a scope beyond the sentence [42].

# Kraaij et al. 2001

- **Name**:

- **Reference**: [81]

- **Short description**: Probabilistic single document extractive summarizer.

- **System Features**

  - **Input**:
  - **Architecture**: The system follows a probabilistic approach. Two different statistical models are applied and their results are combined for selecting the sentences that have to be included in the summary. The former is a content-based language model (unigrams + smoothing) and the latter is based on non-content features (being or not the first sentence, containing cue phrases, sentence length, etc.)
  - **Language coverage**: English, potentially multilingual
  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## Lal and Rüger 2002

- **Name**:

- **Reference**: [84]

- **Short description**: single-document, extract-based summarizer, applies anaphora resolution and text simplification.

- **System Features**

  - **Input**:

  - **Architecture**: following the approach of [82], it works as a Bayesian pattern classifier over sentences trained from an annotated corpus. The features that are taken into account are: length of the sentence, position of the sentence within the paragraph and the paragraph within the document, mean *tf\*idf* of named entities, co-reference with named entities in headline, inclusion of highly co-refered named entities. Some dangling anaphors are replaced by their referent. Lexical simplification is performed with tools from the PSET project [31]. Background knowledge on people and places, taken from sources on the web, can also be included.

  - **Output facilities and constraints**: English

- **Evaluation**: DUC 2002, performed well except for grammaticality and coherence.

- **Classification**

  - within classification 1 (level of processing): entity/discourse

  - within classification 2 (kind of information): lexical/structural

  - within classification 3 (Tucker, 1999): sentence by sentence

- **Comments**: A demonstration can be found at http://km.doc.ic.ac.uk/pr-p.lal-2002/, and the system can be downloaded as a CREOLE Repository for GATE users.

# Lethbridge, University of

- **Name**: University of Lethbridge

- **Reference**: [27, 32]

- **Short description**: single- and multidocument lexical chain summarizer by extraction. It filters out chain candidates in subordinate clauses.

- **System Features**

    - **Input**:
    - **Architecture**: for multidocument summaries, the procedure is the same as for single document (below), but all segments in the collection are pooled together, assigning a time stamp to each.
        * topic segmentation of the text
        * removing unimportant nouns from text (nouns in subordinate clauses).
        * lexical chaining
        * sentence extraction
        * surface repairs: add previous sentence to a sentence containing a dangling anaphora, remove short sentences or sentences with question or quotation marks.
        *
    - **Language coverage**: English, potentially multilingual
    - **Output facilities and constraints**:

- **Evaluation**: DUC 2002, but no results reported in reference. They also participated in DUC 2003, obtaining "reasonable results" but admitting that "some improvements are still required when considering multi-document summarization".

- **Classification**

    - within classification 1 (level of processing): entity
    - within classification 2 (kind of information): lexical
    - within classification 3 (Tucker, 1999): attentional networks

- **Comments**:

# K. U. Leuven

- **Name**:

- **Reference**: [6, 7]

- **Short description**: adapts a hierarchical topic segmentation algorithm to text summarization

- **System Features**

  - **Input**:
  - **Architecture**: For multi-document summarization, a combination of topic segmentation and clustering techniques is applied, while for single-document headline generation, topic segmentation is combined with sentence scoring and compression.

    Thematic structures in texts are detected using generic text structure cues:

    * lexical chains are built following [17] but using only WordNet synonymy relations.
    * the topic of each sentence is determined, by general topicality mechanisms of English (initial position, persistency).
    * topics are distinguished from subtopics, because the first spread throughout the whole text, while the second have local scope.
    * for single document summarization, the number of levels of the topic hierarchy is restricted by the targeted summary length, so that only sentences in higher levels are included.
    * for multiple document summarization, headline-kind summaries are produced by listing non-redundant topic terms. For longer summaries, open-class words of every sentence in the collection are clustered.

    Key terms are associated to each topic, and a tree-like table of content is produced.

  - **Language coverage**: English, potentially multilingual
  - **Output facilities and constraints**: oriented to tables of contents, lacks cohesion for texts.

- **Evaluation**: DUC 2002, average scores, bad for short abstracts. In DUC 2003, the strategy for very short abstracts (headlines) was significantly improved, combining the informativeness of topic terms with hand-crafted grammatical rules for sentence compression, which resulted in very good results for the task of headline generation. In the other tasks, results were average.

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): lexic
  - within classification 3 (Tucker, 1999): attentional networks

- **Comments**:

# Lexical Bonds

- **Name**: Lexical Bonds

- **Reference**: [77]

- **Short description**: extractive single-document system based on analysis of lexical bonds between sentences in a text and a classification of sentences into important and unimportant using SVM.

- **System Features**

  - **Input**: single documents
  - **Architecture**: the original design includes a transformation phase that should compact the text extracted in the first phase and resolve anaphoric references, but it is not yet developed. The current architecture is:
    * sentences are splitted and stopwords are removed
    * record of features for every sentence: sentence position, number of words, number of backward, forward and total lexical bonds and lexical links, and information content
      · a lexical link between two sentences is found when a word stem occurs in both of them, a lexical bond is found when there are two or more lexical links between a pair of sentences [67].
      · the information content of a sentence is the IR function BM25 [146], which indicates the importance of the sentence with respect to the document.
    * SVM are used to select sentences according to these features (trained on DUC'02 manually selected extracts)
    * summaries are generated by following lexical bonds from a given sentence. Some constraints are: only sentences in the upper half of the document and selected by SVM are considered.

    The system produces cohesive summaries, but they are very redundant.
  - **Language coverage**: English, potentially multilingual
  - **Output facilities and constraints**: compactation process is under development.

- **Evaluation**: participated in DUC 2002, with good results in quality.

- **Classification**

  - within classification 1 (level of processing): surface/entity
  - within classification 2 (kind of information): discourse
  - within classification 3 (Tucker, 1999): attentional networks

- **Comments**:

## MEAD

- **Name**: MEAD

- **Reference**: [133, 128, 121, 129]

- **Short description**: Centroid-based multi-document summarization

- **System Features**

  - **Input**:
  - **Architecture**: MEAD begins identifying all the articles related to an emerging event (using the CIDR Topic Detection and Tracking system). CIDR produces a set of clusters. From each cluster a centroid is built. Then the sentences closest to the each of the centroids are selected to be included in the summary. CBSU (Centroid-based sentence utility) scores the degree of relevance of a particular sentence to the general topic of the entire cluster. CSIS (Cross-sentence informational subsumption) measures the overlap between the informational content of the sentences. CSIS is a similar measure than MMR. The difference is that CSIS is multi-document and query-independent while MMR is single-document and query-based. More recent versions of MEAD use a linear combination of three features: the centroid score and it assigns higher scores to sentences closer to the beginning of the document and to longer sentences.
  - **Language coverage**: multilingual: English, Chinese, potentially any language
  - **Output facilities and constraints**:

- **Evaluation**: DUC 2001, 2002 and 2003. In DUC 2002 they had format problems (SGML tags). In DUC 2003 they had the best score for question-focused multi-document summaries, and performed among the top 3 systems for all multi-document summarization tasks.

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## MULTIGEN

- **Name**: MULTIGEN

- **Reference**: [19], [109]

- **Short description**: Multi-document Summarization using Information Fusion and Reformulation

- **System Features**

  - **Input**: News articles presenting different descriptions of the same event.

  - **Architecture**:
    * identify similarities and differences across documents by statistical techniques [111]
    * extract sets of similar sentences: THEMES
    * shallow syntactic analysis
    * order sets of similar sentences (Reformulation). Two different forms of implementing ordering are included: majority ordering and chronological ordering.
    * generation: Sentence generation begins with phrases, with paraphrases rules derived from corpus analysis. MULTIGEN takes profit of the experience of Columbia's group in NL Generation for building high quality summaries (not extracts but abstracts).

  - **Language coverage**: English

  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): entity

  - within classification 2 (kind of information): structural

  - within classification 3 (Tucker, 1999): informational content

- **Comments**: MULTIGEN has been extended in several directions. See Columbia MDS [18, 110] PERSIVAL [107] and CENTRIFUSER [75] among others.

## Muresan et al. 2001, Tzoukermann et al. 2001

- **Name:**

- **Reference:** [117], [158]

- **Short description:** e-mail summarization combining Machine Learning and linguistic information.

- **System Features**

  - **Input:**
  - **Architecture:** The basic process consists on learning the salient NPs occurring in the text. The following features are used for the learning task:
    * for the head of the NP:
      · head-tf*idf (relevance)
      · head-focc (position of the first occurrence of head)
    * for the whole NP
      · np-tf*idf
      · np-focc
      · np-length-words
      · np-length-chars
      · sentence-position
      · paragraph-position
      · all constituents in the NP equally weighted

    Different ML methods have been applied including decision trees (C4.5) and rule induction (Ripper). The linguistic process include:
    * inflectional morphology processing
    * removing unimportant modifiers
    * removing common words
    * removing empty words
  - **Language coverage:** English, potentially multilingual
  - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): understanding
  - within classification 3 (Tucker, 1999): attentional networks

- **Comments:**

# Myaeng and Jang 1999

- **Name:**

- **Reference:** [118]

- **Short description:** Single document summarizer based on statistical techniques

- **System Features**

    - **Input:**

    - **Architecture:** The system uses two similarity measures for determining if a sentence belongs to the major content: a similarity between the sentence and the rest of the document and a similarity between the sentence and the title of the document. Two statistical techniques are applied, a Bayesian model based on 14 features (signature terms and positional information) and the Dempter-Shafer combination rule.

    - **Language coverage:** English, potentially multilingual

    - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

    - within classification 1 (level of processing): surface

    - within classification 2 (kind of information): lexical

    - within classification 3 (Tucker, 1999): informational content

- **Comments:**

# NeATS, iNeATS

- **Name**: NeATS

- **Reference**: [92, 93, 88]

- **Short description**: Multi-document summarizer presented in DUC'01, DUC'02

- **System Features**

  - **Input**:
  - **Architecture**: NeATS proceeds in the following steps:
    1. extracting and ranking passages
       * Identification of key concepts for each topic group
       * Computing of unigram, bigram, trigram Topic Signatures
       * Removing words or phrases occurring in less than the half of texts
       * Saving signatures in a tree
       * Webclopedia query formation
       * Sentence-level IR giving to a ranked list of sentences
    2. Filtering for content: remove all sentences that are not within the first 10 sentences of a document, decrease ranking score of sentences containing stigma words.
    3. Enforcing cohesion and coherence by pairing each sentence with the lead sentence of the document
    4. Filtering for length: include sentences (paired with the corresponding lead sentence) that are most different from the included ones, until targeted length is satisfied.
    5. Ensuring chronological coherence

    As an additional enhancement, Leuski et al. (2003) [88] provide a graphical interface to improve the navigation and modification of the summaries produced by NeATS.
  - **Language coverage**: English, potentially multilingual
  - **Output facilities and constraints**:

- **Evaluation**: in DUC 2002, it was the system with highest precision and F1 measure, although it performed low in recall.

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): structural
  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## Newsblaster

- **Name**: Multilingual Columbia's Newsblaster

- **Reference**: [49]

- **Short description**:

- **System Features**

  - **Input**: multidocument

  - **Architecture**: A platform for multilingual news summarization that extends the Columbia's Newsblaster system [106]. The system adds a new component, translation, to the original six major modules: crawling, extraction, clustering, summarization, classification and web page generation, that have been, in turn, modified for allowing multilinguality (language identification, different character encoding, language idiosyncrasy, etc.).

    In this system multilingual documents are translated into English before clustering and, so, clustering is performed only on English texts.

    Translation is carried out at two levels. As a low quality translation is usually enough for clustering purposes and assessing the relevance of the sentences, a simple and fast technique is applied for glossing the input documents prior to clustering. Higher (relatively) quality translation (using Altavista's Babelfish interface to Systran) is performed in a second step only over fragments selected to be part of the summary.

    The system takes as well into account the possible degradation of the input texts as result of the translation process (most of the sentences resulting from this process are simply not grammatically correct).

  - **Output facilities and constraints**:

  - **Language coverage**: crosslingual

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): entity

  - within classification 2 (kind of information): structural

  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## NTT

- **Name**: NTT

- **Reference**: [65, 66]

- **Short description**: extractive summarizer based on classification of sentences by Support Vector Machines (SVM) and Maximal Marginal Relevance (MMR).

- **System Features**

  - **Input**:

  - **Architecture**: each sentence in a document is described with the following features: position, length, weight (*tf\*idf* score of the words in the sentence), similarity with the headline and presence of certain prepositions or verbs.

  - **Language coverage**: English, potentially multilingual

  - **Output facilities and constraints**:

- **Evaluation**: participated in DUC'02, with good results in coverage but low quality. For DUC 2003, NTT achieved the highest metrics for readability in the two multidocument summarization tasks it took part in, and got average positions for coverage.

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): lexical

  - within classification 3 (Tucker, 1999): sentence by sentence

- **Comments**:

## OCELOT

- **Name**: OCELOT

- **Reference**: [115]

- **Short description**: Summarizing of Web pages. Gist of Web document based on probabilistic models.

- **System Features**

  - **Input**:

  - **Architecture**: OCELOT is one of the applications of a general probabilistic approach that models summarisation as a translation process between two languages, the language of full text and the language of summaries. Berger in his thesis applies conventional stochastic translation methods for summarizing. Three different examples of application are provided and OCELOT is one of them.

  - **Language coverage**: English, potentially multilingual

  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): lexical

  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## PERSIVAL

- **Name**: PERSIVAL

- **Reference**: [107]

- **Short description**: PERSIVAL (Personalized Retrieval and Summarization of Image, Video and Language). The system builds patient specific (tailored access for both patients and physicians) summaries of medical articles contained in a distributed multimedia patient care digital library. It is a Digital Library project.

- **System Features**

  - **Input**: Multimedia collections in the medical domain
  - **Architecture**: Multimedia search triggered by a concept from patient's data. The system includes the annotation and organization of large collections of video data. Video documents are segmented and a storyboard summary is produced. Video are indexed at syntactic and semantic levels. A set of content-based video search tools has been developed. The system includes the use of DEFINDER tool (for looking for definitions).
  - **Language coverage**: English
  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): understanding
  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## RIPTIDES

- **Name**: RIPTIDES

- **Reference**: [164], [162]

- **Short description**: user directed document summarizer combining the application of techniques of Information Extraction, Extraction-based Summarization and Natural Language Generation. The former reference refers to single-document summarization while the latter to multi-document summarization.

- **System Features**

  - **Input:**
  - **Architecture**: The system proceeds in the following steps:
    1. User information needs are acquired from the system
    2. Scenario templates are filled by an IE system
    3. IE output templates are merged into an event-oriented structure where comparable facts are grouped. For doing so SimFinder is used.
    4. Importance scores are assigned to slot/sentences based on a combination of document position, document recency and group/cluster membership.
    5. Content selection
    6. Summary generation
  - **Language coverage**: English
  - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): understanding
  - within classification 3 (Tucker, 1999): informational content

- **Comments:**

## Schiffman et al. 2001

- **Name**:

- **Reference**: [143]

- **Short description**: Multi-document summarizer producing Biographical Summaries combining linguistic knowledge with corpus statistics.

- **System Features**

  - **Input**:

  - **Architecture**: A number of modules co-operate for producing the summaries:
    * Sentence tokenizer
    * Alembic POS tagger
    * Nametag NER
    * Cass parser
    * Cross-document co-reference
    * Appositives
    * Relative clause weighting
    * Sentential description, following [Sagion, Lapalme, 2000]

  - **Language coverage**: English

  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): entity

  - within classification 2 (kind of information): understanding

  - within classification 3 (Tucker, 1999): informational content

- **Comments**:

## SUMMARIST

- **Name**: SUMMARIST

- **Reference**: [69, 95]

- **Short description**: Extractive single document summarisation system

- **System Features**

  - **Input**:
  - **Architecture**: The system proceeds in three steps: Topic identification, Interpretation and Summary generation.
    * Topic identification implies previous acquisition of Topic Signatures and then the identification of a text span as belonging to a topic characterised by its signature. Topic Signatures are tuples of the form ¡Topic, Signature¿ where Signature is a list of weighted terms: ¡t1,w1¿, ¡t2,w2¿, ..., ¡tn,wn¿. Topic signatures can be automatically learned ([Lin, 1997], [Lin, Hovy, 2000]). Topic identification, then, includes text segmentation (using TextTiling) and comparison of text spans with existing Topic Signatures.
    * The topic identified are fused during the interpretation (2nd step) of the process. The fused topics are then reformulated (expressed in new terms).
    * The last step is a conventional extractive task.
  - **Language coverage**: multilingual: English, Japanese, Spanish, Arabic, Indonesian, Korean, potentially any language
  - **Output facilities and constraints**:

- **Evaluation**:

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexical
  - within classification 3 (Tucker, 1999): attentional networks

- **Comments**:

# SumUM

- **Name**: SumUM

- **Reference**: [50, 139, 51]

- **Short description**: generates single-document abstracts of scientific papers, based on shallow syntactic and semantic analysis oriented to conceptual identification and hand-made templates for text-regeneration. It interacts with the user. For DUC, an adaptation has been made to obtain biased multi-document summaries.

- **System Features**

  - **Input**: single-document, scientific or technical articles with the following structure: title, author and affiliation, introduction, main section, references. There is also an adaptation for multi-document summarization.
  - **Architecture**:
    * transducers identify concepts in text: domain transducers identify author, references, etc., and linguistic transducers identify noun groups and verb groups.
    * concepts are tagged semantically, marking discourse domain relations
    * sentences of indicative and informative type are identified
    * an indicative abstract is composed, by re-generation of text using pre-defined summary templates
    * based on the first, indicative abstract, an informative abstract can be composed, elaborating a specific query of the user
  - **Language coverage**: English
  - **Output facilities and constraints**: an interactive system: the user is presented with a short indicative abstract and a list of topics available for expansion, and an informative abstract can be produced, focusing on the topics chosen by the user.

- **Evaluation**: it was formally adapted to participate in DUC 2002, but with no adaptation to the news domain. It was ranked among the three first in quality, and the second in length-adjusted coverage, most probably due to the efficiency of templates. In DUC 2003, SumUM was adapted for biased multi-document summarization, achieving good scores for coverage but with a decrease on the quality of the resulting summaries.

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): understanding
  - within classification 3 (Tucker, 1999): informative content

- **Comments**:

## Strzalkowski 1998

- **Name:**

- **Reference:** [148]

- **Short description:** Query-based single document non-extractive summarizer

- **System Features**

  - **Input:**
  - **Architecture:** The system proceeds in two steps, Analysis and Generation. Analysis phase consists of three tasks: Feature extraction, feature synthesis and rule induction. As result a set of themes is identified. The system uses both simple and composite features. Simple features include word co-occurrence, noun phrases (detected with linkIT), WN synonyms and common semantic classes for verbs (following Levin's, see [Klavans, Kan, 1998]). Generation phase includes the performance of a content planner (based on the intersection of themes obtained in the previous phase and on a sentence planner) and a sentence generator.
  - **Language coverage:** English
  - **Output facilities and constraints:**

- **Evaluation:**

- **Classification**

  - within classification 1 (level of processing): entity
  - within classification 2 (kind of information): structural
  - within classification 3 (Tucker, 1999): informational content

- **Comments:**

## Teufel and Moens

- **Name**:

- **Reference**: [155, 156]

- **Short description**: analyzes the rhetorical structure of scientific articles and produces extractive summaries with the main contributions.

- **System Features**

  - **Input**: scientific articles (specialized in computational linguistics domain)

  - **Architecture**: Each sentence in an article is described with a number of features, like its length (in words) or its position in the document. But the main emphasis is put in describing the contribution of each sentence to the rhetorical structure of the document. To do that, a number of linguistic knowledge sources are exploited, among others: document layout, section titles, lexico-syntactical structures, citation procedures and cue phrases typical of the genre of scientific articles.

    Then, a machine learning algorithm is applied to classify each sentence as one of a number of rhetorical categories that account for the rhetorical status of the sentence with respect to the whole text. A parallel classification is carried out to determine the relevance of each sentence.

  - **Output facilities and constraints**:

  - **Language coverage**: English

- **Evaluation**: an evaluation by comparison with a human-made golden standard is presented in [156], with good results.

- **Classification**

  - within classification 1 (level of processing): discourse

  - within classification 2 (kind of information): structural

  - within classification 3 (Tucker, 1999): discourse structure / sentence by sentence

- **Comments**:

## TNO-TPD summarizer

- **Name**: TNO-TPD summarizer

- **Reference**: [81], [80]

- **Short description**: extractive multi-document summarizer. Sentences are selected according to a statistical language model and applying a bayesian classifier.

- **System Features**

  - **Input**:
  - **Architecture**:
    * an unigram language model of a cluster of documents determines content-based salience of each sentence
    * each sentence is assigned values for some surface features: sentence position, length, presence of positive or negative cue phrases, and the mentioned content score.
    * sentences are classified by a Naive Bayes classifier into summary and non-summary sentences.
    * redundancy is reduced by applying MMR [30]
    * to generate headlines, the most frequent word in the highest ranked sentence for every document and the titles is considered a *trigger word*. Then, the sentences in the whole cluster are ranked according to their importance. The highest ranked noun phrase that contains the trigger word is chosen as the headline.
  - **Language coverage**: English, potentially multilingual
  - **Output facilities and constraints**:

- **Evaluation**: participated in DUC 2002 in the multi-document extract and abstract tracks, with "disappointing performance". In addition, a self-evaluation applying relative utility [133], which reports better results. An investigation on the individual contribution of each feature was also performed, revealing that *position in the sentence* is highly indicative, while *negative cue phrase* was not well-defined.

- **Classification**

  - within classification 1 (level of processing): surface
  - within classification 2 (kind of information): lexic
  - within classification 3 (Tucker, 1999): attentional networks / sentence by sentence

- **Comments**:

## van Halteren 2002

- **Name**:

- **Reference**: [159]

- **Short description**: multi-document, extractive summarizer. Sentences are classified by feature sets used for writing style recognition.

- **System Features**

  - **Input**:

  - **Architecture**: each sentence is described by a set of features: distance between occurrences of the same word, distribution of words, relative position of words, sentence length, sentence position and context of POS tags. A classifier trained for a writing style recognition task exploits these features for sentence scoring and extraction.

  - **Language coverage**: English, potentially multilingual

  - **Output facilities and constraints**:

- **Evaluation**: participated in DUC 2002, but obtained not so good results.

- **Classification**

  - within classification 1 (level of processing): surface

  - within classification 2 (kind of information): lexical

  - within classification 3 (Tucker, 1999): sentence by sentence

- **Comments**: the system was trained on materials not oriented to the summarization task