

## MACHINE TRANSLATION: A BRIEF HISTORY

W.John Hutchins

The translation of natural languages by machine, first dreamt of in the seventeenth century, has become a reality in the late twentieth. Computer programs are producing translations - not perfect translations, for that is an ideal to which no human translator can aspire; nor translations of literary texts, for the subtleties and nuances of poetry are beyond computational analysis; but translations of technical manuals, scientific documents, commercial prospectuses, administrative memoranda, medical reports. Machine translation is not primarily an area of abstract intellectual inquiry but the application of computer and language sciences to the development of systems answering practical needs.

After an outline of basic features, the history of machine translation is traced from the pioneers and early systems of the 1950s and 1960s, the impact of the ALPAC report in the mid-1960s, the revival in the 1970s, the appearance of commercial and operational systems in the 1980s, research during the 1980s, new developments in research in the 1990s, and the growing use of systems in the past decade. This brief history can mention only the major and most significant systems and projects, and for more details readers are referred to the publications listed.

### 1. Basic features and terminology

The term 'machine translation' (MT) refers to computerized systems responsible for the production of translations with or without human assistance. It excludes computer-based translation tools which support translators by providing access to on-line dictionaries, remote terminology databanks, transmission and reception of texts, etc. The boundaries between machine-aided human translation (MAHT) and human-aided machine translation (HAMT) are often uncertain and the term computer-aided translation (CAT) can cover both, but the central core of MT itself is the automation of the full translation process.

Although the ideal goal of MT systems may be to produce high-quality translation, in practice the output is usually revised (post-edited). It should be noted that in this respect MT does not differ from the output of most human translators which is normally revised by a second translator before dissemination. However, the types of errors produced by MT systems do differ from those of human translators (incorrect prepositions, articles, pronouns, verb tenses, etc.). Post-editing is the norm, but in certain circumstances MT output may be unedited or only lightly revised, e.g. if it is intended only for specialists familiar with the text subject. Output might also serve as a rough draft for a human translator, i.e. as a 'pre-translation'.

The translation quality of MT systems may be improved either, most obviously, by developing more sophisticated methods or by imposing certain restrictions on the input. The system may be designed, for example, to deal with texts limited to the sublanguage (vocabulary and grammar) of a particular subject field (e.g. biochemistry) and/or document type (e.g. patents). Alternatively, input texts may be written in a controlled language, which restricts the range of vocabulary, and avoids homonymy and polysemy and complex sentence structures. A third option is to require input texts to be marked (pre-edited) with indicators of prefixes, suffixes, word divisions, phrase and clause boundaries, or of different grammatical categories (e.g. the noun *cónvict* and its homonymous verb *convíct*). Finally, the system itself may refer problems of ambiguity and selection to human operators (usually translators) for resolution during the processes of translation itself, in an interactive mode.

Systems are designed either for two particular languages (bilingual systems) or for more than a single pair of languages (multilingual systems). Bilingual systems may be designed to operate either in only one direction (unidirectional), e.g. from Japanese into English, or in both directions (bidirectional).

Multilingual systems are usually intended to be bidirectional; most bilingual systems are unidirectional.

In overall system design, there have been three basic types. The first (and historically oldest) type is generally referred to as the 'direct translation' approach: the MT system is designed in all details specifically for one particular pair of languages, e.g. Russian as the language of the original texts, the source language, and English as the language of the translated texts, the target language. Translation is direct from the source language (SL) text to the target language (TL) text; the basic assumption is that the vocabulary and syntax of SL texts need not be analyzed any more than strictly necessary for the resolution of ambiguities, the correct identification of TL expressions and the specification of TL word order; in other words, SL analysis is oriented specifically to one particular TL. Typically, systems consist of a large bilingual dictionary and a single monolithic program for analysing and generating texts; such 'direct translation' systems are necessarily bilingual and unidirectional.

The second basic design strategy is the interlingua approach, which assumes that it is possible to convert SL texts into representations common to more than one language. From such interlingual representations texts are generated into other languages. Translation is thus in two stages: from SL to the interlingua (IL) and from the IL to the TL. Procedures for SL analysis are intended to be SL-specific and not oriented to any particular TL; likewise programs for TL synthesis are TL-specific and not designed for input from particular SLs. A common argument for the interlingua approach is economy of effort in a multilingual environment. Translation from and into  $n$  languages requires  $n(n-1)$  bilingual 'direct translation' systems; but with translation via an interlingua just  $2n$  interlingual programs are needed. With more than three languages the interlingua approach is claimed to be more economic. On the other hand, the complexity of the interlingua itself is greatly increased. Interlinguas may be based on an artificial language, an auxiliary language such as Esperanto, a set of semantic primitives presumed common to many or all languages, or a 'universal' language-independent vocabulary.

The third basic strategy is the less ambitious transfer approach. Rather than operating in two stages through a single interlingual representation, there are three stages involving underlying (abstract) representations for both SL and TL texts. The first stage converts SL texts into abstract SL-oriented representations; the second stage converts these into equivalent TL-oriented representations; and the third generates the final TL texts. Whereas the interlingua approach necessarily requires complete resolution of all ambiguities in the SL text so that translation into any other language is possible, in the transfer approach only those ambiguities inherent in the language in question are tackled; problems of lexical differences between languages are dealt with in the second stage (transfer proper). Transfer systems consist typically of three types of dictionaries (SL dictionary/ies containing detailed morphological, grammatical and semantic information, similar TL dictionary/ies, and a bilingual dictionary relating base SL forms and base TL forms) and various grammars (for SL analysis, TL synthesis and for transformation of SL structures into TL forms).

Within the stages of analysis and synthesis (or generation), many MT systems exhibit clearly separated components involving different levels of linguistic description: morphology, syntax, semantics. Hence, analysis may be divided into morphological analysis (identification of word endings, word compounds), syntactic analysis (identification of phrase structures, dependency, subordination, etc.), semantic analysis (resolution of lexical and structural ambiguities); synthesis may likewise pass through semantic synthesis (selection of appropriate compatible lexical and structural forms), syntactic synthesis (generation of required phrase and sentence structures), and morphological synthesis (generation of correct word forms). In transfer systems, the transfer component may also have separate programs dealing with lexical transfer (selection of vocabulary equivalents) and with structural transfer (transformation into TL-appropriate structures). In some earlier forms of transfer systems analysis did not involve a semantic stage, transfer was restricted to the conversion of syntactic structures, i.e. syntactic transfer alone.

In many older systems, particularly those of the 'direct translation' type the components of analysis, transfer and synthesis were not always clearly separated. Some of them also mixed data (dictionary and grammar) and processing rules and routines. Later systems have exhibited various degrees of modularity, so that system components, data and programs can be adapted and changed without damage to overall system efficiency. A further stage in some recent systems is the reversibility of analysis and synthesis

components, i.e. the data and transformations used in the analysis of a particular language are applied in reverse when generating texts in that language.

The direct translation approach was typical of the "first generation" of MT systems. The indirect approach of interlingua and transfer based systems is often seen to characterise the "second generation" of MT system types. Both are based essentially on the specification of rules (for morphology, syntax, lexical selection, semantic analysis, and generation). Most recently, corpus-based methods have changed the traditional picture (see below). During the last five years, there is beginning to emerge a "third generation" of hybrid systems combining the rule-based approaches of the earlier types and the more recent corpus-based methods. The differences between direct and indirect, transfer and interlingua, rule-based, knowledge-based and corpus-based are becoming less useful for the categorization of systems. Transfer systems incorporate interlingual features (for certain areas of vocabulary and syntax); interlingua systems include transfer components; rule-based systems make increasing use of probabilistic data and stochastic methods; statistics- and example-based systems include traditional rule-based grammatical categories and features; and so forth. These recent developments underline what has always been true, namely that MT research and MT systems adopt a variety of methodologies in order to tackle the full range of language phenomena, complexities of terminology and structure, misspellings, 'ungrammatical' sentences, neologisms, etc. The development of an operational MT system is necessarily a long-term 'engineering' task applying techniques which are well known, reliable and well tested.

## **2. Precursors and pioneers, 1933-1956**

The use of mechanical dictionaries to overcome the barriers of language was first suggested in the 17th century. However, it was not until the 20th century that the first concrete proposals were made, in patents issued independently in 1933 by George Artsrouni, a French-Armenian, and by a Russian, Petr Smirnov-Troyanskii. Artsrouni designed a storage device on paper tape which could be used to find the equivalent of any word in another language; a prototype was apparently demonstrated in 1937. The proposals by Troyanskii were in retrospect more significant. He envisioned three stages of mechanical translation: first, an editor knowing only the source language was to undertake the 'logical' analysis of words into their base forms and syntactic functions; secondly, the machine was to transform sequences of base forms and functions into equivalent sequences in the target language; finally, another editor knowing only the target language was to convert this output into the normal forms of his own language. Troyanskii envisioned both bilingual and multilingual translation. Although his patent referred only to the machine which would undertake the second stage, Troyanskii believed that "the process of logical analysis could itself be mechanized".

Troyanskii was ahead of his time and was unknown outside Russia when, within a few years of their invention, the possibility of using computers for translation was first discussed by Warren Weaver of the Rockefeller Foundation and Andrew D. Booth, a British crystallographer. On his return to Birkbeck College (London) Booth explored the mechanization of a bilingual dictionary and began collaboration with Richard H. Richens (Cambridge), who had independently been using punched cards to produce crude word-for-word translations of scientific abstracts. However, it was a memorandum from Weaver in July 1949 which brought the idea of MT to general notice (Weaver 1949). He outlined the prospects and suggested various methods: the use of war-time cryptography techniques, statistical methods, Shannon's information theory, and the exploration of the underlying logic and universal features of language, "the common base of human communication".

Within a few years research had begun at the University of Washington (Seattle), at the University of California at Los Angeles and at the Massachusetts Institute of Technology. It was at MIT in 1951 that the first full-time researcher in MT was appointed, Yehoshua Bar-Hillel. A year later he convened the first MT conference, where the outlines of future research were already becoming clear. There were proposals for dealing with syntax by Victor Oswald and by Bar-Hillel (his categorial grammar), suggestions that texts should be written in MT-oriented restricted languages, and arguments for the construction of sublanguage systems. It was obvious that fully automatic translation would not be achieved without long-term basic research, and (in the interim) human assistance was essential, either to prepare texts or to revise the output (known already as pre- and post-editing.) A number of participants

considered that the first requirement was to demonstrate the feasibility of MT. Accordingly, at Georgetown University Leon Dostert collaborated with IBM on a project which resulted in the first public demonstration of an MT system in January 1954. A carefully selected sample of 49 Russian sentences was translated into English, using a very restricted vocabulary of 250 words and just 6 grammar rules. Although it had little scientific value, it was sufficiently impressive to stimulate the large-scale funding of MT research in the USA and to inspire the initiation of MT projects elsewhere in the world, notably in the USSR.

### **3. The decade of high expectation and disillusion, 1956-1966**

In the 1950s and 1960s research tended to polarize between empirical trial-and-error approaches, often adopting statistical methods with immediate working systems as the goal, and theoretical approaches involving fundamental linguistic research and aiming for long-term solutions. The contrastive methods were usually described at the time as 'brute-force' and 'perfectionist' respectively. Any evaluation of the period must remember that computer facilities were frequently inadequate; much effort was devoted to improving basic hardware (paper tapes, magnetic media, access speeds, etc.) and to devising programming tools suitable for language processing. Some groups were inevitably forced to concentrate on theoretical issues, particularly in Europe and the Soviet Union. For political and military reasons, most US research was for Russian-English translation, and most Soviet research on English-Russian systems.

The research under Erwin Reifler at the University of Washington (Seattle) epitomized the word-for-word approach; it involved the construction of large bilingual dictionaries where lexicographic information was used not only for selecting lexical equivalents but also for solving grammatical problems without the use of syntactic analysis. Entries gave English translations with rules for local reordering of output. The huge lexicon made extensive use was made of English 'cover terms' for Russian polysemes, the inclusion of phrases and clauses and the classification of vocabulary into sublanguages. After initial work on German and English, the group was engaged on the foundations of a Russian-English system for the 'photoscopic store', a large memory device rather similar to the laser disk. From 1958 practical development was directed by Gilbert King at the IBM Corporation (Yorktown Heights, New York). A system was installed for the US Air Force which produced translations for many years, until replaced in 1970 by Systran (see 7 below). By any standards the output was crude and sometimes barely intelligible, but, unlike some MT researchers at the time, King never made excessive claims for his system. With all its deficiencies, it was able to satisfy basic information needs of scientists.

The empirical attitude was exemplified at the RAND Corporation (1950-1960), which distrusted current linguistic theory and emphasized statistical analyses. From a large corpus (Russian physics texts) were prepared bilingual glossaries with grammatical information and simple grammar rules; a computer program was written for a rough translation; the result was studied by post-editors who indicated errors; the revised text was analyzed; the glossaries and the rules were revised; the corpus was translated again; and so the process continued in cycles of translation and post-editing. The main method of analysis was statistical distribution, but it was at RAND that David Hays developed the first parser based on dependency grammar.

The research under Leon Dostert at Georgetown University had a more eclectic approach, undertaking empirical analyses of texts only when traditional grammatical information was inadequate. The Georgetown group was the largest in the USA and there were considerable differences of viewpoint among its members. Four groups were set up, each encouraged to submit their methods for testing in open competition on a Russian chemistry text. Ariadne Lukjanow's 'code-matching' method produced excellent results but apparently with ad hoc corpus-specific rules; the 'syntactic analysis' method by Paul Garvin (who had prepared the linguistic basis for the 1954 demonstration system) was not ready in time for the test; and the 'sentence-by-sentence' method by Anthony Brown was an example of the empirical cyclical method. The 'general analysis' method by a group under Michael Zarechnak was the method adopted and named Georgetown Automatic Translation (GAT). This had three levels of analysis: morphological (including identification of idioms), syntagmatic (agreement of nouns and adjectives, government of verbs, modification of adjectives, etc.), and syntactic (subjects and predicates, clause relationships, etc.) GAT was initially implemented on the SERNA system, largely the work of Petr Toma (later designer of

Systran), and then with the programming method developed by Brown for his own separate French-English system. In this form it was successfully demonstrated in 1961 and 1962, and as a result Russian-English systems were installed at Euratom in Ispra (Italy) in 1963 and at the Oak Ridge National Laboratory of the US Atomic Energy Commission in 1964.

Further development of his parsing method was continued by Paul Garvin at the Ramo-Wooldridge Corporation from 1960 to 1967. Garvin sought a middle way between the empiricists and the perfectionists. His fulcrum method was essentially a dependency parser, a linguistic pattern recognition algorithm which identified the fulcrum of a structure and the relationships of dependent elements to the fulcrum. At a later stage, Garvin introduced heuristic methods, employing statistical information when appropriate. The method was also applied at Wayne State University (1958-1972) in the development of a Russian-English system.

Anthony Oettinger at Harvard University believed in a gradualist approach. From 1954 to 1960 the group concentrated on the compilation of a massive Russian-English dictionary, to serve as an aid for translators (a forerunner of the now common computer-based dictionary aids), to produce crude word-for-word translations for scientists familiar with the subject, and as the basis for more advanced experimental work. From 1959 research focused on the 'predictive syntactic analyzer' - originally developed at the National Bureau of Standards under Ida Rhodes - a system for the identification of permissible sequences of grammatical categories (nouns, verbs, adjectives, etc.) and the probabilistic prediction of following categories. The results were often unsatisfactory, caused primarily by the enforced selection at every stage of the 'most probable' prediction. The system was revised to examine all possible predictions. (This was the Multiple-path Predictive Analyzer, from which was later developed William Woods' familiar Augmented Transition Network parser.) However, the results were equally unsatisfactory: multiple parsings were produced of even apparently quite unambiguous sentences; some kind of semantic analysis to 'filter' out undesirable parses was clearly needed. By 1965 the group had effectively ceased MT research.

Research at MIT, started by Bar-Hillel in 1951, was directed by Victor Yngve from 1953 until its end in 1965. Whereas other groups saw syntax as an adjunct to lexicographic transfer, as a means of resolving ambiguities and rearranging TL output, Yngve placed syntax at the centre: translation was a three-stage process, a SL grammar analyzed input sentences as phrase structure representations, a 'structure transfer routine' converted them into equivalent TL phrase structures, and the TL grammar rules produced output text. There is some evidence of the influence of transformational grammar (Chomsky was associated with the project for two years), but in many respects the practicalities of MT led MIT researchers away from Chomskyan theory. Much was achieved both in basic syntactic research and on developing the first string-handling programming language (COMIT). But eventually the limitations of the 'syntactic transfer' approach became obvious. By the mid-1960's Yngve acknowledged that MT research had come up against the 'semantic barrier... and that we will only have adequate mechanical translations when the machine can "understand" what it is translating' (Yngve 1964).

The Linguistic Research Center (LRC) at the University of Texas was founded by Winfried Lehmann in 1958 and, like MIT, concentrated on basic syntactic research of English and German. Efforts were made to devise reversible grammars to achieve bidirectional translation within an essentially 'syntactic transfer' approach. The foundations were laid for the later successful development of the METAL system (sections 5 and 6 below.)

At the University of California, Berkeley, the project under the direction of Sydney Lamb stressed the importance of developing maximally efficient dictionary routines and a linguistic theory appropriate for MT. Lamb developed his stratificational grammar with networks, nodes and relations paralleling the architecture of computers. Translation was seen as a series of decoding and encoding processes, from the graphemic stratum of SL text through its morphemic and lexemic strata to a sememic stratum, from which TL text could be generated by passing through a similar series of strata. Translation was characterized as word-by-word, each word examined within the broadest possible environment and not limited by sentence boundaries or immediate contexts.

There were no American groups taking the interlingua approach. This was the focus of projects elsewhere. At the Cambridge Language Research Unit, Margaret Masterman and her colleagues adopted two basic lines of research: the development of a prototype interlingua producing crude 'pidgin' (essentially word-for-word) translations, and the development of tools for improving and refining MT output, primarily by means of the rich semantic networks of a thesaurus (conceived as lattices of interlocking meanings.) At Milan, Silvio Ceccato concentrated on the development of an interlingua based on conceptual analysis of words (species, genus, activity type, physical properties, etc.) and their possible correlations with other words in texts.

In the Soviet Union research was as vigorous as in the United States and showed a similar mix of empirical and basic theoretical approaches. At the Institute of Precision Mechanics research began in 1955 shortly after a visit by D.Y. Panov to see a demonstration of the IBM-Georgetown system. Experiments on English-Russian translation were on similar lines to the approach at Georgetown, but with less practical success. More basic research was undertaken at the Steklov Mathematical Institute under Lyapunov, Kulagina and others, mainly towards French-Russian translation. As at MIT attention was paid to the development of programming tools for linguistic processes. The research by Igor Mel'chuk and others at the Institute of Linguistics in Moscow was more theoretical in nature, including work on interlinguas and leading to the stratificational 'meaning-text' model as a basis for MT (see sect. 5 below). However, the main centre for interlingua investigations was Leningrad University, where a team under Nikolai Andreev conceived of an interlingua not as an abstract intermediary representation (which was Mel'chuk's approach) but as an artificial language complete in itself with its own morphology and syntax, and having only those features statistically most common to a large number of languages.

By the mid-1960s MT research groups had been established in many countries throughout the world, including most European countries (Hungary, Czechoslovakia, Bulgaria, Belgium, Germany, France, etc.), China, Mexico, and Japan. Many of these were short-lived and with no subsequent influence. But some groups created at the time became important later, in particular the project which began in 1960 at Grenoble University.

#### **4. The ALPAC report and its consequences**

In the 1950s optimism was high; developments in computing and in formal linguistics, particularly in the area of syntax, seemed to promise great improvement in quality. There were many predictions of imminent breakthroughs and of fully automatic systems operating within a few years. However, disillusion grew as the complexity of the linguistic problems became more and more apparent. In a review of MT progress, Bar-Hillel (1960) criticized the prevailing assumption that the goal of MT research should be the creation of fully automatic high quality translation (FAHQT) systems producing results indistinguishable from those of human translators. He argued that it was not merely unrealistic, given the current state of linguistic knowledge and computer systems, but impossible in principle. He demonstrated his argument with the word *pen*. It can have at least two meanings (a container for animals or children, and a writing implement). In the sentence *The box was in the pen* we know that only the first meaning is plausible; the second meaning is excluded by our knowledge of the normal sizes of (writing) pens and boxes. Bar-Hillel contended that no computer program could conceivably deal with such 'real world' knowledge without recourse to a vast encyclopedic store. His argument carried much weight at the time. Many researchers were already encountering similar 'semantic barriers' for which they saw no straightforward solutions. Bar-Hillel recommended that MT should adopt less ambitious goals, it should build systems which made cost-effective use of man-machine interaction.

In 1964 the government sponsors of MT in the United States formed the Automatic Language Processing Advisory Committee (ALPAC) to examine the prospects. In its famous 1966 report it concluded that MT was slower, less accurate and twice as expensive as human translation and that 'there is no immediate or predictable prospect of useful machine translation.' It saw no need for further investment in MT research; instead it recommended the development of machine aids for translators, such as automatic dictionaries, and continued support of basic research in computational linguistics. The ALPAC report was widely condemned as narrow, biased and shortsighted. It is true that it failed to recognize, for example, that revision of manually produced translations is essential for high quality, and it was unfair to criticize MT

for needing to post-edit output. It may also have misjudged the economics of computer-based translation, but large-scale support of current approaches could not continue. The influence of the ALPAC report was profound. It brought a virtual end to MT research in the USA for over a decade and MT was for many years perceived as a complete failure.

## 5. The quiet decade, 1967-1976.

In the United States the main activity had concentrated on English translations of Russian scientific and technical materials. In Canada and Europe the needs were quite different. The Canadian government's bicultural policy created a demand for English-French (and to a less extent French-English) translation beyond the capacity of the market. The problems of translation were equally acute in Europe, in particular within the European Communities with growing demands for translations of scientific, technical, administrative and legal documentation from and into all the Community languages. The focus of MT activity switched from the United States to Canada and to Europe.

At Montreal, research began in 1970 on a syntactic transfer system for English-French translation. The TAUM project (Traduction Automatique de l'Université de Montréal) had two major achievements: firstly, the Q-system formalism, a computational metalanguage for manipulating linguistic strings and trees and the foundation of the Prolog programming language widely used in natural language processing; and secondly, the Météo system for translating weather forecasts. Designed specifically for the restricted vocabulary and limited syntax of meteorological reports, Météo has been successfully operating since 1976 (albeit since 1984 in a new version running under a different programming language GramR developed subsequently by the original designer of Météo, John Chandiooux). An attempt by TAUM to repeat its success with another sublanguage, that of aviation manuals, failed to overcome the problems of complex noun compounds and phrases (e.g. *hydraulic ground test stand pressure and return line filters*), problems which would defeat human translators without the relevant subject knowledge. TAUM came to an end in 1981.

The principal experimental efforts of the decade focused on interlingua approaches, with more attention to syntactic aspects than previous projects at Cambridge, Milan and Leningrad. Between 1960 and 1971 the group established by Bernard Vauquois at Grenoble University developed an interlingua system for translating Russian mathematics and physics texts into French. The 'pivot language' of CETA (Centre d'Etudes pour la Traduction Automatique) was a formalism for representing the logical properties of syntactic relationships. It was not a pure interlingua as it did not provide interlingual expressions for lexical items; these were translated by a bilingual transfer mechanism. Syntactic analysis produced first a phrase-structure (context-free) representation, then added dependency relations, and finally a 'pivot language' representation in terms of predicates and arguments. After substitution of TL lexemes (French), the 'pivot language' tree was converted first into a dependency representation and then into a phrase structure for generating French sentences. A similar model was adopted by the LRC at Texas during the 1970s in its METAL system: sentences were analyzed into 'normal forms', semantic propositional dependency structures with no interlingual lexical elements.

The research by Mel'chuk in the Soviet Union towards an interlingua system was more ambitious. His influential 'meaning-text' model combined a stratificational dependency approach (six strata: phonetic, phonemic, morphemic, surface syntactic, deep syntactic, semantic) with a strong emphasis on lexicographic aspects of an interlingua. Fifty universal 'lexical functions' were identified at the deep syntactic stratum covering paradigmatic relations (synonyms, antonyms, conversives (*fear: frighten*), verbs and corresponding agentive nouns (*write: writer, prevent: obstacle*), etc.) and a great variety of syntagmatic relations (inceptive verbs associated with given nouns, *conference: open, war: break out*; idiomatic causatives, *compile: dictionary, lay: foundations*, etc.) Although Mel'chuk emigrated to Canada in 1976 his ideas continue to inspire Russian MT research to the present, as well as receiving wider attention elsewhere.

By the mid-1970s, however, the future of the interlingua approach seemed to be in doubt. Both LRC and CETA had problems which were attributed to the rigidity of the levels of analysis (failure at any stage meant failure to produce any output), the inefficiency of parsers (too many partial analyses which had to

be 'filtered' out), and in particular the loss of information about the surface forms of SL input which could have been used to guide the selection of TL forms and the construction of acceptable TL sentence structures. As a consequence, it became widely accepted that the less ambitious transfer approach offered better prospects.

## **6. Operational and commercial systems, 1976-1989**

During the 1980s MT advanced rapidly on many fronts. Many new operational systems appeared, the commercial market for MT systems of all kinds expanded, and MT research diversified in many directions.

The revival was laid in the decade after ALPAC. Systems were coming into operational use and attracting public attention. The Georgetown systems had been operating since the mid-1960s. As well as Météo, two other sublanguage systems had appeared: in 1970 the Institut Textile de France introduced TITUS, a multilingual system for translating abstracts written in a controlled language, and in 1972 came CULT (Chinese University of Hong Kong) for translating mathematics texts from Chinese into English, a 'direct translation' system requiring extensive pre- and post-editing. More significant, however, were the installations of Systran in 1970 by the US Air Force for Russian-English translation, and in 1976 by the European Communities for English-French translation.

Systran has been the most successful operational system so far. Developed by Petr Toma, who had previously worked for the Georgetown project, initially as a 'direct translation' system, its oldest version is the Russian-English system at the USAF Foreign Technology Division (Dayton, Ohio) which translates over 100,000 pages a year. At the Commission of the European Communities (CEC) the English-French version was followed shortly by systems for French-English, English-Italian and subsequently for combinations of most other languages of the European Communities (now European Union). The original design has been greatly modified, with increased modularity and greater compatibility of the analysis and synthesis components of different versions, permitting cost reductions when developing new language pairs. Outside the CEC, Systran has been installed at a number of intergovernmental institutions, e.g. NATO and the International Atomic Energy Authority, and at a number of major companies, e.g. General Motors of Canada, Dornier, and Aérospatiale. The application at the Xerox Corporation is particularly noteworthy: post-editing has been virtually eliminated by controlling the input language of technical manuals for translation from English into a large number of languages (French, German, Italian, Spanish, Portuguese, and Scandinavian languages).

Another long-established commercial system is that of the Logos Corporation, directed by Bernard E.Scott. This company's first effort in MT was an English-Vietnamese system for translating aircraft manuals during the 1970s. Experience gained in this ultimately short-term project was applied in the development of a German-English system which appeared on the market in 1982. Initially, Logos systems were based on a direct translation approach, but later systems are closer to a transfer design and incorporate sophisticated means for recording and applying semantic features.

Systems such as Systran and Logos are in principle designed for general application, although in practice their dictionaries are adapted for particular subject domains. Systems specifically designed for one particular environment were also developed during the 1970s and 1980s. The Pan American Health Organization in Washington built two mainframe systems, one for Spanish into English (SPANAM, basically a direct system) and the other for English into Spanish (ENGSPAN, a transfer system). Both were developed essentially by just two researchers, Muriel Vasconcellos and Marjorie León - showing what could be achieved at the time with limited resources using well-tested and reliable techniques.

Large tailor-made systems have been the speciality of the Smart Corporation (New York) since the early 1980s. Customers have included Citicorp, Ford, and largest of all, the Canadian Department of Employment and Immigration. The principal feature of Smart systems is (as at Xerox) strict control of SL vocabulary and syntax. Texts are written in restricted English (interactively at terminals); manuals are clear and unambiguous, and translation is regarded as almost a by-product.



During the 1980s, the greatest commercial activity was in Japan, where most of the computer companies developed software for computer-aided translation, mainly for the Japanese-English and English-Japanese market, although they did not ignore the needs for translation to and from Korean, Chinese and other languages. Many of these systems are low-level direct or transfer systems with analysis limited to morphological and syntactic information and with little or no attempt to resolve lexical ambiguities. Often restricted to specific subject fields (computer science and information technology are popular choices), they rely on substantial human assistance at both the preparatory (pre-editing) and the revision (post-editing) stages. Examples are systems from Oki (PENSEE), Mitsubishi (MELTRAN), Sanyo, Toshiba (AS-TRANSAC), Hitachi (HICATS) and Fujitsu (ATLAS). Japanese input demands considerable pre-editing, but it is acceptable to operators of Japanese word processors who have to interpret Japanese script, with two vernacular alphabets (hiragana and katakana), Chinese characters, no capitals and no indicators of word boundaries. As consequence, however, good knowledge of Japanese is essential for usable results from Japanese-English systems.

The most sophisticated commercially available system during the 1980s was, however, the METAL German-English system, which had originated from research at the University of Texas University. After its interlingua experiments in the mid 1970s this group adopted an essentially transfer approach, with research funded since 1978 by the Siemens company in Munich (Germany). The METAL system, intended for translation of documents in the fields of data processing and telecommunications, appeared in 1988 and has been followed in the early 1990s by systems involving Dutch, French and Spanish as well as English and German.

Some of the Japanese systems were designed for microcomputers. But they were not the first. The earliest came from the USA at the beginning of the 1980s. Often linguistically crude, but sometimes capable of providing economically viable results, microcomputer systems have secured for MT a higher public profile than the mainframe systems had done. Most have been rightly marketed not as full MT systems but as computer aids for translators.

Earliest were the American Weidner and ALPS systems, which became commercially available in 1981 and 1983 respectively. The ALPS system offered three levels of assistance: multilingual word-processing, automatic dictionary and terminology consultation, and interactive translation. In the latter case, translators could work with MT-produced rough drafts. However, the ALPS products were not profitable, and from the mid 1980s onwards the company diverted into providing a translation service rather than selling computer aids for translators. The Weidner (later World Communications Center) systems proved to be more successful commercially, offering translation packages for a large number of language pairs, with its Japanese-English systems being particularly popular. One version MicroCAT was intended for small personal computers, the other MacroCAT for larger machines. In the late 1980s Weidner was acquired by Bravice and shortly afterwards the MacroCAT version was sold to the Intergraph Corporation (see below). By this time, other systems for personal computers had come onto the market (PC-Translator from Linguistic Products, GTS from Globalink and the Language Assistant series from MicroTac), which were to have a major impact in the next decade.

## **7. Research from 1976 to 1989**

Research after the mid-1970s had three main strands: first, the development of advanced transfer systems building upon experience with earlier interlingua systems; secondly, the development of new kinds of interlingua systems; and thirdly, the investigation of techniques and approaches from Artificial Intelligence.

After the failure of its interlingua system, the Grenoble group (GETA, Groupe d'Etudes pour la Traduction Automatique) began development of its influential Ariane system. Regarded as the paradigm of the "second generation" linguistics-based transfer systems, Ariane influenced projects throughout the world in the 1980s. Of particular note was its flexibility and modularity, its powerful tree-transducers, and its conception of static and dynamic grammars. Different levels and types of representation (dependency, phrase structure, logical) could be incorporated on single labelled tree structures and thus provide considerable flexibility in multilevel transfer representations. GETA was particularly prominent in the

encouragement of international collaboration and in the training of MT researchers, particularly in South-East Asia. Ariane as such did not become an operational system (despite hopes in the mid 1980s from its involvement in the French national project Calliope) and active research on the system ceased in the late 1980s. It was tested thoroughly by B'VITAL and became part of the EuroLang project in the 1990s, while the Grenoble team has continued MT research on other lines (see sect. 10 below).

Similar in conception to the GETA-Ariane design was the Mu system developed at the University of Kyoto under Makoto Nagao. Prominent features of Mu were the use of case grammar analysis and dependency tree representations, and the development of a programming environment for grammar writing (GRADE). The Kyoto research has been very influential in many subsequent Japanese MT research projects and in many of the Japanese commercial systems of the 1980s. Since 1986, the research prototype has been converted and elaborated into an operational system for use by the Japanese Information Center for Science and Technology for the translation of abstracts.

Experimental research at Saarbrücken (Germany) began in 1967. From the mid 1970s until the late 1980s the group developed SUSY (Saarbrücker Übersetzungssystem), a highly-modular multilingual transfer system displaying an impressive heterogeneity of linguistic techniques: phrase structure rules, transformational rules, case grammar and valency frames, dependency grammar, and variety of operation types: rule-driven, lexicon-driven, table-driven, the use of statistical data, preference rules, etc. Its main focus was the in-depth treatment of inflected languages such as Russian and German, but many other languages were also investigated, particularly English. Other projects included a French-German system (ASCOF) using semantic networks; and the development of a German generator (SEMSYN) to convert output from the Fujitsu ATLAS system in order to translate titles of Japanese scientific articles into German.

One of the best known projects of the 1980s was the Eurotra project of the European Communities. Its aim was the construction of an advanced multilingual transfer system for translation among all the Community languages. Its general design owed much to GETA-Ariane and SUSY systems. Like them, it was a linguistics-based modular transfer system intended for multilingual translation producing good quality but not perfect output. The design combined lexical, logico-syntactic and semantic information in multilevel interfaces at a high degree of abstractness. No direct use of extra-linguistic knowledge bases or of inference mechanisms was made, and no facilities for human assistance or intervention during translation processes were to be incorporated. It assumed batch processing and human post-editing.

During the period Eurotra stimulated much other innovative theoretical linguistic and computational-linguistic research, particularly in the Netherlands, Belgium, Denmark, Germany and Great Britain. Eurotra researchers advanced substantially the theoretical foundations of MT and made important contributions to syntactic theory, formal parsing theory, and discourse analysis. One of the aims of the Eurotra project was to stimulate such research, and in this it succeeded. However, it did not produce a working prototype, and attempts towards the end of the project to involve industrial partnerships were largely unfruitful. Like Ariane, a major defect, readily conceded by those involved, was the failure to tackle problems of the lexicon, both theoretically and practically. Whereas at the end of the 1970s, Eurotra was seen as representing the best 'linguistics-based' design, at the end of the 1980s it was seen by many as basically obsolete in conception, and by 1992 Eurotra had effectively come to an end.

During the latter half of the 1980s there was a general revival of interest in interlingua systems. Some were small and short-lived, e.g. the ATAMIRI system from Bolivia (primarily for English and Spanish) based on a South American language Aymara. Much more important were projects in the Netherlands, Japan and the United States, some beginning to use knowledge-based methods from research on artificial intelligence.

The DLT (Distributed Language Translation) system at the BSO software company in Utrecht (The Netherlands) was a six-year project from 1985 under the general direction of Toon Witkam. It was intended as a multilingual interactive system operating over computer networks, where each terminal was to be a translating machine from and into one language only; texts were to be transmitted between terminals in an intermediary language. As its interlingua, DLT chose a modified form of Esperanto.

Analysis was restricted primarily to morphological and syntactic features (formalised in a dependency grammar); there was no semantic analysis of the input. Disambiguation took place in the central interlingua component, where semantico-lexical knowledge was represented in an Esperanto database. From a combination of linguistic and extra-linguistic information the system computed probability scores for pairs of dependency-linked interlingual words. The project made a significant effort in the construction of large lexical databases, and in its final years proposed the building of a Bilingual Knowledge Bank from a corpus of (human) translated texts (Sadler 1989). In a number of respects DLT was a precursor of developments which became more prominent in the 1990s.

A second interlingua project in the Netherlands was innovative in another respect. This was the Rosetta project at Philips (Eindhoven) directed by Jan Landsbergen. The designers of this experimental system, involving three languages (English, Dutch and Spanish), opted to explore the use of Montague grammar in interlingual representations. A fundamental feature was the derivation of semantic representations from the syntactic structure of expressions, following the principle of compositionality; for each syntactic derivation tree there was to be a corresponding semantic derivation tree, and these semantic derivation trees were the interlingual representations. A second feature was the exploration of the reversibility of grammars, a feature of many subsequent MT projects.

MT research in Japan, initially greatly influenced by the Mu project at Kyoto University, showed a wide variety of approaches. While transfer systems predominated there were also a number of interlingua systems (e.g. the PIVOT system from NEC, now available commercially) and knowledge-based experiments (e.g. the LUTE project at NTT, and the Lamb system of Canon). Japan also launched its own multilingual multinational project in the mid 1980s, with participants from China, Indonesia, Malaysia and Thailand and the involvement of major Japanese research institutes, including the governmental Electrotechnical Laboratory (ETL) in Tokyo. An ambitious interlingua approach has been adopted with knowledge-based contextual analysis for disambiguation. The project continues to the present day.

Outside North America, Western Europe, and Japan, MT research was also becoming vigorous during the latter half of the 1980s in Korea (sometimes in collaborative projects with Japanese and American groups), in Taiwan (e.g. the ArchTran system), in mainland China at a number of institutions, and in Southeast Asia, particularly in Malaysia. Some of this research has been stimulated by the multilingual multinational project MMT led by the Japanese CICC (Center of the International Cooperation for Computerization), but also, more generally, it has been driven by commercial expansion particularly in the electronics and computing industries.

Until the end of the decade there was also an increase in activity in the Soviet Union. The ALPAC report had a negative impact during the 1970s, and a number of advanced MT projects lost support. From 1976 most research was concentrated at the All-Union Centre for Translation in Moscow. Systems for English-Russian (AMPAR) and German-Russian translation (NERPA) were developed based on the direct approach, and although work under the direction of Yurii Apresyan on a more advanced transfer system for French-Russian (FRAP) did continue, most activity in the Soviet Union was focused on the production of relatively low-level operational systems, often involving the use of statistical analyses.

In the latter part of the 1980s developments in syntactic theory, in particular unification grammar, Lexical Functional Grammar and Government Binding theory, began to attract researchers, although their principal impact was to come in the 1990s. At the time, many observers believed that the most likely source of techniques for improving MT quality lay in research on natural language processing within the context of artificial intelligence (AI). Investigations of AI methods in MT research began in the mid-1970s with Yorick Wilks' work on 'preference semantics' and 'semantic templates'. Further inspiration came from the research of Roger Schank at Yale University, and particularly from the development of expert systems and knowledge-based approaches to text 'understanding'.

A number of projects applied knowledge-based approaches, particularly the use of knowledge banks - some in Japan (e.g. the ETL research for the Japanese multilingual project), others in Europe (e.g. at Saarbrücken and Stuttgart), and many in North America. The most important research has been undertaken at Carnegie-Mellon University in Pittsburgh which under the leadership of Jaime Carbonell

and Sergei Nirenburg has experimented with a number of knowledge-based MT systems since the mid 1980s to the present time. The methodology is described as "meaning-oriented MT in an interlingua paradigm". A working prototype for English and Japanese in both directions was designed for translation of personal computer manuals. The basic components included a small concept lexicon for the domain, analysis and generation lexicons for the two languages, a syntactic parser with semantic constraints, a semantic mapper (for semantic interpretation), an interactive 'augmentor', a semantic generator producing TL syntactic structures with lexical selection, and a syntactic generator for producing target sentences. The concept lexicon and the semantic information in the analysis and generation lexicons (i.e. defining semantic constraints) are language-independent but specific to the domain. The core of the system is the interlingual representation of texts, in the form of networks of propositions. They are derived from the processes of semantic analysis and of interactive disambiguation performed by the 'augmentor' by reference to the domain knowledge of the 'concept lexicon'. By the end of the 1980s, the Carnegie-Mellon team had fully elaborated its KANT prototype system and was ready to develop an operational knowledge-based system (see sect.10 below.)

## **9. Corpus-based MT research since 1989 to the present**

The dominant framework of MT research until the end of the 1980s was based on essentially linguistic rules of various kinds: rules for syntactic analysis, lexical rules, rules for lexical transfer, rules for syntactic generation, rules for morphology, etc. The rule-based approach was most obvious in the dominant transfer systems (Ariane, Metal, SUSY, Mu and Eurotra), but it was at the basis of all the various interlingua systems - both those which were essentially linguistics-oriented (DLT and Rosetta), and those which were knowledge-based (KANT).

Since 1989, however, the dominance of the rule-based approach has been broken by the emergence of new methods and strategies which are now loosely called 'corpus-based' methods. Firstly, a group from IBM published in 1988 the results of experiments on a system based purely on statistical methods. The effectiveness of the method was a considerable surprise to many researchers and has inspired others to experiment with statistical methods of various kinds in subsequent years. Secondly, at the very same time certain Japanese groups began to publish preliminary results using methods based on corpora of translation examples, i.e. using the approach now generally called 'example-based' translation. For both approaches the principal feature is that no syntactic or semantic rules are used in the analysis of texts or in the selection of lexical equivalents.

The most dramatic development has been the revival of the statistics-based approach to MT in the Candide project at IBM. Statistical methods were common in the earliest period of MT research, in the 1960s (see sect.3 above), but the results had been generally disappointing. With the success of newer stochastic techniques in speech recognition, the IBM team at Yorktown Heights began to look again at their application to MT. The distinctive feature of Candide is that statistical methods are used as virtually the sole means of analysis and generation; no linguistic rules are applied. The IBM research is based on the vast corpus of French and English texts contained in the reports of Canadian parliamentary debates (the Canadian Hansard). The essence of the method is first to align phrases, word groups and individual words of the parallel texts, and then to calculate the probabilities that any one word in a sentence of one language corresponds to a word or words in the translated sentence with which it is aligned in the other language.

What surprised most researchers (particularly those involved in rule-based approaches) was that the results were so acceptable: almost half the phrases translated either matched exactly the translations in the corpus, or expressed the same sense in slightly different words, or offered other equally legitimate translations. Obviously, the researchers have sought to improve these results, and the IBM group proposes to introduce more sophisticated statistical methods, but they also intend to make use of some minimal linguistic information, e.g. the treatment of all morphological variants of a verb as a single word, and the use of syntactic transformations to bring source structures closer to those of the target language.

The second major 'corpus-based' approach - benefiting likewise from improved rapid access to large databanks of text corpora - is what is known as the 'example-based' (or 'memory-based') approach.

Although first proposed in 1984 by Makoto Nagao, it was only towards the end of the 1980s that experiments began, initially in some Japanese groups and during the DLT project (as already mentioned). The underlying hypothesis is that translation often involves the finding or recalling of analogous examples, i.e. how a particular expression or some similar phrase has been translated before. The example-based approach is founded on processes of extracting and selecting equivalent phrases or word groups from a databank of parallel bilingual texts, which have been aligned either by statistical methods (similar perhaps to those used by the IBM group) or by more traditional rule-based morphological and syntactic methods of analysis. For calculating matches, some MT groups use semantic methods, e.g. a semantic network or a hierarchy (thesaurus) of domain terms. Other groups use statistical information about lexical frequencies in the target language. The main advantage of the approach is that since the texts have been extracted from databanks of actual translations produced by professional translators there is an assurance that the results will be accurate and idiomatic.

The availability of large corpora has encouraged experimentation in methods deriving from the computational modelling of cognition and perception, in particular research on parallel computation, neural networks or connectionism. In natural language processing connectionist models are 'trained' to recognise the strongest links between grammatical categories (in syntactic patterns) and between lexical items (in semantic networks). The potential relevance to MT is clear enough for both analysis and transfer operations, given the difficulties of formulating accurate grammatical and semantic rules in traditional approaches. As yet, however, within MT only a few groups have done some small-scale research in this framework.

Connectionism offers the prospect of systems 'learning' from past successes and failures. So far the nearest approach has been for systems to suggest changes on the basis of statistical data about corrections made by users, e.g. during post-editing. This model is seen in the commercial Tovna system and in the experimental PECOF 'feedback' mechanism in the Japanese MAPTRAN system. A similar mechanism has been incorporated in the NEC PIVOT system.

## **10. Rule-based systems since 1990.**

Although the main innovation since 1990 has been the growth of corpus-based approaches, rule-based research continues in both transfer and interlingua systems. For example, a number of researchers involved in Eurotra have continued to work on the theoretical approach developed, e.g. the CAT2 system at Saarbrücken, and one of the fruits of Eurotra research has been the PaTrans transfer-based system developed in Denmark for Danish/English translation of patents.

In a rather different form, both theoretical foundations of Ariane and Eurotra survive in the EuroLang project which is based at SITE, a French company which had purchased B'VITAL, the Grenoble company founded to develop Ariane in the French National MT project Calliope. The EuroLang project involves also the German company Siemens-Nixdorf which is to contribute with its transfer-based METAL system. EuroLang is developing ten language pairs, English into and from French, German, Italian and Spanish, and French into and from German. However, the first product of EuroLang has not been an MT system as such but a translator's workstation, the Optimizer, which incorporates an MT module (not just for METAL but also for the Logos systems.)

Otherwise, the most important work on the linguistics-based transfer approach is the research continuing on the LMT project which began under Michael McCord in the mid-1980s, and which is based at a number of IBM research centres in Germany, Spain, Israel and the USA. In LMT, translation is via four steps implemented in Prolog: lexical analysis, producing descriptions of input words and their transfers; syntactic analysis of source texts, producing representations of both surface and deep (logical) relations; transfer, involving both isomorphic structural transfer and restructuring transformations; and morphological generation of target texts. LMT is characterised as combining a lexicalist approach to grammar and logic programming - LMT stands for 'Logic programming MT'. The language pairs under investigation include English-German, German-English, and English-Spanish. In 1994 the LMT programs were marketed as modules for the IBM computer-based workstation TranslationManager/2.

The interlingua approach thrives with greater vigour, even with the end of the DLT and Rosetta projects. The leading group at Carnegie Mellon University continues with its knowledge-based approach (Nirenburg et al. 1992), developing several models over the years. In 1992, it announced the beginning of a collaborative project with the Caterpillar company with the aim of creating a large-scale high-quality system CATALYST for multilingual translation of technical manuals in the specific domain of heavy earth-moving equipment. This system combines the knowledge-based approach with controlled input.

Other interlingua-based systems are, e.g. the ULTRA system at the New Mexico State University, and the UNITRAN system based on the linguistic theory of Principles and Parameters (Dorr 1993). This experiment has shown that the application of abstract theory, albeit in a small-scale project, can raise significant issues for MT research in the area of lexicon construction. There is also the Pangloss project, an interlingual system restricted to the vocabulary of mergers and acquisitions, a collaborative project involving experts from the universities of Southern California, New Mexico State and Carnegie Mellon. Pangloss is itself one of three MT projects supported by ARPA, the others being the IBM statistics-based project mentioned above, and a system being developed by Dragon Systems, a company which has been particularly successful in speech research but with no previous experience in MT. The restitution of US government support for MT research signals perhaps the end of the damaging impact of the ALPAC report (sect.4 above).

Since the mid 1980s there has been a general trend towards the adoption of 'unification' and 'constraint-based' formalisms (e.g. Lexical-Functional Grammar, Head-Driven Phrase Structure Grammar, Categorical Grammar, etc.) In "second generation" rule-based systems (such as Eurotra and Ariane) there were series of complex multi-level representations and large sets of rules for the transformation, mapping and testing of labelled tree representations. Many rules applied only in very specific circumstances and to specific representations, i.e. grammars and transduction rules defined the 'constraints' determining transfer from one level to another and hence from SL text to TL text. The introduction of the unification and constraint-based approaches has led to the simplification of the rules (and hence the computational processes) of analysis, transformation and generation. Now, there are monostratal representations and a restricted set of abstract rules, with conditions and constraints incorporated into specific lexical entries. At the same time, the components of these grammars are in principle reversible, so that it is no longer necessary to construct for the same language different grammars of analysis and generation.

Much of this research has not focused directly on MT as such but on the development of general-purpose systems for natural language processing based on unification and constraint-based grammars. Some of these have been applied to translation tasks, e.g. the Core Language Engine system has been used for translation from Swedish into English and vice versa; the PLNLP (Programming Language for Natural Language Processing) system has provided the foundation for translation systems involving English, Portuguese, Chinese, Korean and Japanese; and the ELU engine (Environnement Linguistique d'Unification) developed at Geneva in Switzerland has formed the basis for a bidirectional system for translating avalanche bulletins between French and German.

The syntactic orientation which characterised transfer systems in the past has thus been replaced by 'lexicalist' approaches, with a consequential increase in the range of information attached to lexical units the lexicon: not just morphological and grammatical data and translation equivalents, but also information on syntactic and semantic constraints and non-linguistic and conceptual information. The expansion of lexical data is seen most clearly in the lexicons of interlingua-based systems, which include large amounts of non-linguistic information. Many groups are investigating and collaborating on methods of extracting lexical information from readily available lexicographic sources, such as bilingual dictionaries intended for language learners, general monolingual dictionaries, specialised technical dictionaries, and the terminological databanks used by professional translators. A notable effort in this area is the Electronic Dictionary Research project, which began in the late 1980s with support from several Japanese computer manufacturing companies.

## **11. New areas of research in the 1990s.**

One consequence of developments in example-based methods has been that much greater attention is now

paid to questions of generating good quality texts in target languages than in previous periods of MT activity when it was commonly assumed that the most difficult problems concerned analysis, disambiguation and the identification of the antecedents of pronouns. In part, the impetus for this research has come from the need to provide natural language output from databases, i.e. translation from the artificial and constrained language used to represent database contents into the natural language of database users. Some MT teams have researched multilingual generation (e.g. in Montreal, for marine forecasts and summaries of statistical data on the labour force.)

More directly, it has been the recognition of a demand for types of translations which have not previously been studied. Since 1990, various groups (at UMIST (Manchester), the University of Brussels, Grenoble University and the Science University of Malaysia) have experimented with 'dialogue-based MT' systems where the text to be translated is composed or written in a collaborative process between man and machine. In this way it is possible to construct a text which the system is known to be capable of translating without further reference to an author who does not know the target language, who cannot revise the output and therefore needs assurance of good quality output. The most obvious application is the fairly standardised messages of business communication.

Probably the most significant development of the last five years has been the growing interest in spoken language translation, with the challenge of combining speech recognition and linguistic interpretation of conversation and dialogue. A small-scale experiment was conducted at British Telecom in the late 1980s, using a pattern-matching approach to translate a small set of standard business phrases from English into French and vice versa. More significant has been the Japanese project at ATR Interpreting Telecommunications Research Laboratories (based at Nara, near Osaka), which began in 1986 and is funded until the end of the century. The team is developing a system for telephone registrations at international conferences and for telephone booking hotel accommodation.

Carnegie Mellon had also experimented with a spoken language translation in its JANUS project during the late 1980s, and since the early 1990s the University of Karlsruhe has been involved in an expansion of JANUS. In 1992, these groups joined ATR in a consortium C-STAR (Consortium for Speech Translation Advanced Research), and in January 1993 gave a successful public demonstration of telephone translation from English, German and Japanese into each of the three languages, within the limited domain of conference registrations. Recently a consortium has added further groups in France (LMSI), the United Kingdom (SRI) and other groups in North America and Asia.

More recently still, in May 1993, began the German funded Verbmobil project, an 8-10 year project aiming to develop a transportable aid for face to face English-language commercial negotiations by Germans and Japanese who do not know English fluently.

A distinctive feature of the last decade has been the globalisation of MT research. Within the last five years, research activity has grown rapidly in China, Taiwan, Korea, India and South East Asia. By contrast, MT research in Eastern Europe has been affected profoundly by the political changes since 1989. In many cases, researchers have successfully joined collaborative projects with Western European groups, e.g. researchers from the Czech Republic and from Bulgaria; others have been able to continue only at a much reduced level, particularly researchers from the former Soviet Union, although the important research under Yurii Apresyan continues on the Russian/English ETAP system.

## **12. Operational systems since 1990.**

The use of MT accelerated in the 1990s. The increase has been most marked in commercial agencies, government services and multinational companies, where translations are produced on a large scale, primarily of technical documentation. This is the major market for the mainframe systems: Systran, Logos, METAL, and ATLAS. All have installations where translations are being produced in large volumes. Indeed, it has been estimated that in 1993 over 300 million words a year were translated by such services: for example, one Logos operation alone (at Lexi-Tech, Canada) was translating annually more than 25 million words of technical manuals.

In these companies, the operators are often not the traditional professional translators, for whom perhaps earlier MT systems were intended. Before the 1980s it was often assumed that the aim of MT research was the (partial) replacement of human translation. Now the aim is focused on special domain-restricted mass-volume systems and on systems for non-translators - areas where professional translators have not been active. It is quite clear from recent developments that what the professional translators need are tools to assist them: provide access to dictionaries and terminological databanks, multilingual word processing, management of glossaries and terminology resources, input and output communication (e.g. OCR scanners, electronic transmission, high-class printing). For these reasons, the most appropriate and successful developments of the last few years have been the translator workstations (e.g. IBM's TranslationManager, the TRADOS TWB, the STAR Transit systems, and the Eurolang Optimizer, and PTT from the Canadian Translation Services.) MT research has contributed to their development most significantly through the statistical research on bilingual text alignment, e.g. at AT&T Bell Laboratories (Murray Hill, NJ) and Centre d'Innovations en Technologies de l'Information (Montreal), which has enabled translators to store and access previous translations for their later (partial) reuse or revision or as sources for examples of translations - this facility is known generally as 'translation memory'.

A notable impetus has been the requirements of large computer software companies which sell in international markets. Both the software itself and its accompanying documentation must be translated quickly and accurately into the local language if companies are to maintain competitive advantage. Organisations set up to provide the 'localisation' and 'globalisation' of computer products are major users of MT systems and translation workbenches adapted and designed for their particular needs.

At the same time as the large MT systems are being used more widely, there has been an even more rapid growth in the use and availability of systems for personal computers. Several Japanese companies now market MT systems for microcomputers (e.g. Toshiba, Nova, Sharp, Hitachi, Mitsubishi, NEC, NTT, Oki, Brother, Catena), many for both IBM-compatibles and Macintosh machines. In the United States systems for personal computers are becoming very popular, in particular the systems from Globalink and MicroTac (recently merged in a single company) and PC-Translator - all first appearing in the late 1980s - and each now offering a wide variety of language pairs. Others have joined them, e.g. Finalsoft, Toltran, LogoVista, Winger - the latter a system originally developed in Denmark for a specific client. The growth is not confined to the USA and Japan; there are PC systems from Korea, Taiwan, China, Finland (Kielikone) and Russia (PARS and STYLUS). In face of this competition, producers of the older systems designed for mainframes are now beginning to downsize their products for personal computers. For example, Fujitsu have marketed a version of ATLAS, Systran has launched a PC-version of its software, and this trend will undoubtedly accelerate.

A further feature of the current decade is the increasing availability of MT on telecommunications networks. First in this field was Systran when it gave access in the mid 1980s to its systems over the French Minitel network. It followed by providing network access to subscribers in the United States, and similar access is available in Europe to the versions of Systran developed for the Commission of the European Communities. More recently, the CompuServe network has launched a translation service for its members founded upon the DP/Translator system from Intergraph; and Globalink has announced a similar service on Internet.

At the same time, the development of systems for specific subject domains and users has also expanded rapidly in recent years - often with controlled languages and based on specific sublanguages. Some of these systems have been developed for software companies for clients. For example, Volmac Lingware Services has produced MT systems for a textile company, an insurance company, and for translating aircraft maintenance manuals; Cap Gemini Innovation developed TRADEX to translate military telex messages for the French Army; in Japan, CSK developed its own ARGO system for translation in the area of finance and economics, and now offers it also to outside clients; and in Denmark, a similar history lies behind the commercial development of the Winger system. User-designed systems are a sign that the computational methods of MT are now becoming familiar outside the limited circles of researchers. Though rarely innovative from a theoretical viewpoint, they are often computationally advanced.

As the 1990s progress it has become clear that different types of MT systems are required to meet widely



differing translation needs. Those identified so far include the traditional MT systems for large organisations, usually within a restricted domain; the translation tools and workstations (with MT modules as options) designed for professional translators; the cheap PC systems for occasional translations; the use of systems to obtain rough gists for the purposes of surveillance or information gathering; the use of MT for translating electronic messages; systems for monolinguals to translate standard messages into unknown languages; systems for speech translation in restricted domains. It is equally clear that as MT systems of many varieties become more widely known and used the range of possible translation needs and possible types of MT systems will also become wider and stimulate further research and development, quite probably in directions not yet envisioned.

## References

The general history of MT is covered by Hutchins (1986), updated by Hutchins (1988, 1993, 1994), where full references for the systems and projects mentioned will be found. Basic sources for the early period are Locke & Booth (1955), Booth (1967) and Bruderer (1982). For the period after ALPAC (1966) there are good descriptions of the major MT systems in Slocum (1988), Nirenburg (1987) and King (1987), while details of later systems will be found in the proceedings of the latest biennial MT Summit conferences (Washington 1991, Kobe 1993) and of the conferences on Theoretical and Methodological Issues in Machine Translation (TMI) at Montreal 1992, and Kyoto 1993. Vasconcellos (1988), Newton (1992) and the Aslib conferences (since 1979) provide the wider perspectives of commercial developments and translators' experiences. For general introductions see the books by Lehrberger & Bourbeau (1988), Nagao (1989), Hutchins & Somers (1992), and Arnold et al. (1994).

ALPAC 1966 *Language and machines: computers in translation and linguistics*. A report by the Automatic Language Processing Advisory Committee. National Academy of Sciences, Washington, DC.

Arnold D et al 1994 *Machine translation: an introductory guide*. NCC/Blackwell, Manchester/Oxford.

Aslib 1979 (to date) *Translating and the computer*. [various editors]. Aslib, London.

Bar-Hillel Y 1960 The present status of automatic translation of languages. *Advances in Computers* 1: 91-163.

Booth A D (ed) 1967 *Machine translation*. North-Holland, Amsterdam.

Bruderer H E (ed) 1982 *Automatische Sprachübersetzung*. Wissenschaftliche Buchgesellschaft, Darmstadt.

Dorr B J 1993 *Machine translation: a view from the lexicon*. MIT Press, Cambridge, Mass.

Hutchins W J 1986 *Machine translation: past, present, future*. Ellis Horwood, Chichester, UK. (Halstead Press, New York)

Hutchins W J 1988 Recent developments in machine translation: a review of the last five years. In: Maxwell D et al.(eds) *New directions in machine translation*. Foris, Dordrecht, 7-62

Hutchins W J, Somers H L (1992) *An introduction to machine translation*. Academic Press, London.

Hutchins W J 1993 Latest developments in machine translation technology. In: MT Summit 4 (1993), 11-34.

Hutchins W J 1994 Research methods and system designs in machine translation: a ten-year review, 1984-1994. In: Machine Translation, Ten Years On, 12-14 November 1994, Cranfield University. 16pp.

King M (ed) 1987 *Machine translation today: the state of the art*. Edinburgh University Press, Edinburgh.

- Lehrberger J, Bourbeau L 1988 *Machine translation: linguistic characteristics of MT systems and general methodology of evaluation*. Benjamins, Amsterdam.
- Locke W N, Booth A D (eds) 1955 *Machine translation of languages*. MIT Press, Cambridge, Mass.
- MT Summit 3: *MT Summit III, July 1-4 1991*, Washington D.C., USA.
- MT Summit 4: *MT Summit IV: International Cooperation for Global Communication, July 20-22, 1993*, Kobe, Japan.
- Nagao M 1989 *Machine translation: how far can it go?* Oxford University Press, Oxford.
- Newton J (ed) 1992 *Computers in translation: a practical appraisal*. Routledge, London.
- Nirenburg S (ed) 1987 *Machine translation: theoretical and methodological issues*. Cambridge University Press, Cambridge.
- Nirenburg, S. et al. (1992): *Machine translation: a knowledge-based approach*. Morgan Kaufmann, San Mateo, Ca.
- Sadler, V. (1989): *Working with analogical semantics: disambiguation techniques in DLT*. Dordrecht: Foris.
- Slocum J (ed) 1988 *Machine translation systems*. Cambridge University Press, Cambridge.
- TMI-92: *Fourth International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages. Empiricist vs Rational Methods in MT*, June 25-27, 1992, Montréal, Canada.
- TMI-93: *Fifth International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages. MT in the Next Generation*, July 14-16, 1993, Kyoto, Japan.
- Vasconcellos M (ed) 1988 *Technology as translation strategy*. State University of New York, Binghamton, NY.
- Weaver W 1949 Translation. In: Locke & Booth (1955): 15-23.
- Yngve V H 1964 Implications of mechanical translation research. *Proceedings of the American Philosophical Society* **108**: 275-281. Repr. with addendum in: Bruderer H (1982), 33-49