

Session 11: EQUIPMENT

MODERN TRENDS IN CHARACTER RECOGNITION MACHINES

Lt. Col. Dimitri A. Kellogg, USA

Army Research Office

In talking about automatic character recognition, popularly called "machine reading", I shall discuss the following:

- A. What it is
- B. Why we need it
- C. What methods are being tried
- D. The problems
- E. Our specific requirements
- F. The current technology

I shall try to stay away from proprietary matters, since this field is becoming highly competitive.

A. What is it?--A device into which one can feed pages of a foreign journal, say Russian, and from which comes a magnetic tape ready to feed into a translation program, transliterated if desired, and with some provision for non-word occurrences (display formulas and equations, superscripts and subscripts, figures) by marking and bypassing, photography, or conversion to digital code. The reading process consists of: feed, positioning (page, line, character), scan, comparison and identification (which involves decision-making), output.

B. Why do we need it?--Because keypunch input is too slow. Six-hundred words an hour are pitifully slow when machine translation rates run 10, 000 to 20, 000 words an hour. Keypunch input is also relatively expensive.

C. What methods are being tried?--All the present methods, except those using magnetic tape input, are optical methods using photo-cell scan. Some of these methods are:

1. Whole-character mask-matching
2. Mask with circular holes in critical areas
3. Matrix
4. Vertical strip-scan with pulse analysis
5. Crossings or radial scans
6. Identification of pieces and assembly thereof

## Session 11: EQUIPMENT

7. Curve tracing and area integration
8. Curve tracing and angle measurement.

Variations which may be applied to the above methods are:

1. Weighing of critical points of a character
2. Photography, to increase contrast and solve page feed
3. Closest-fit criteria, to permit multiple-font reading
4. Extra fonts
5. Size reduction of upper case characters
6. Edge trimming
7. Clustering (smoothing of outlines)
8. Logic, to deduce identity of ambiguous letters
9. Learning techniques involving feedback.

I want to mention in passing that direct voice sensing with printed or tape output is being considered, using pulse analysis, but this is still somewhat further in the future.

D. The problems. --Some of the major problems encountered in making an automatic reading device are:

1. Quality of input (paper, printing, type)
2. Multiple fonts
3. Non-word occurrences - especially pictures
4. Spacing between lines and characters
5. Page feeding and positioning of page, line, and character
6. Errors (in my hand translation of Russian scientific articles, I have found on the average one serious printing error per one and one-half pages, such as incorrect mathematical symbols or omitted numerals)
7. Ambiguities (a reading machine will have to use a logic-program to differentiate between the Russian soft sign ь and the Russian letter ѣ by looking ahead for the ѣ whenever it sees a ь ).

E. Our specific requirements. --Now I should like to present a set of specifications to show what we need in a page reader.

Session 11: EQUIPMENT

PAGE READER SPECIFICATIONS

	Minimum	Optimum
Process Rate (characters per second)	100	1000
Errors (per 5000 characters)	2	1
Fonts (journals handled)	1	all
Non-word occurrences	omitted but marked	handled
Operating cost per word	1/2¢	1/10¢
Machine cost (including development)	\$500,000	\$250,000

The "minimum" column is the worst we can settle for in any category; if a machine touches the left column in any item, it had better be at the so-called "optimum" in most others. A machine just satisfying each minimum criterion is not good enough.

F. The current technology. --No one machine exists today that satisfies even the so-called minimum specifications; however, the minimum is within the state-of-the-art, and so are parts of the optimum requirements. Here are some specific examples of what does exist:

1. One machine in operation has a speed of 200 characters per second in page-reading single-spaced typing of a single special alphanumeric font with an error rate of 1 per 4,000 characters. It is priced in the \$100,000-plus range.
2. An experimental model reads several fonts and uses size reduction. Although the cost is somewhat higher, the reading speed is not reduced since parallel circuits are used.
3. There is another machine which has a speed over the so-called optimum, but which reads only a special 14-character numeric-symbol font on a single line. It is priced in the \$20,000 range.
4. Devices exist which can put a picture on magnetic tape in digital code.

Thus a complete machine satisfying our needs does not exist, but the various required parts either exist or else appear to be feasible. Time estimates for development run from 9 to 12 months, though of course there is a difference between a verbal estimate and a firm contract.

## Session 11: EQUIPMENT

It appears reasonable, therefore, to plan to have automatic reader input for mechanical translation for both research purposes and eventual production.

Since the sticky word "production", has been mentioned, here is my position. I favor useful production of machine translations as soon as possible. I do not favor starting production until we have both assured automatic reading and better translations.

One of my objections would vanish with the introduction of automatic reading in place of keypunching. I firmly believe in the worth and attainability of automatic reading.

As regards machine translations, I do not believe that the best ones available today are acceptable for their intended uses without uneconomical postediting. A test of the unedited translations by non-linguists knowledgeable in the science being translated is obviously in order. If such a test should prove us wrong, it would be both our duty and desire to push for all-out production of machine translations.