



Europäisches Patentamt
European Patent Office
Office européen des brevets



Publication number: **0 602 821 A2**

EUROPEAN PATENT APPLICATION

Application number: **93309575.4**

Int. Cl.⁵: **G06K 9/03**

Date of filing: **01.12.93**

Priority: **15.12.92 GB 9226137**

Date of publication of application:
22.06.94 Bulletin 94/25

Designated Contracting States:
DE FR GB

Applicant: **International Business Machines Corporation**
Old Orchard Road
Armonk, N.Y. 10504(US)

Inventor: **Chevion, Dan**
88 Shoshanat Harcarmel Street
Haifa 3423(IL)
Inventor: **Gilat, Ittai**
38 Henrietta Szold Street
Haifa(IL)
Inventor: **Heilper, Andre**
4/89 Nissenboim Street

Haifa(IL)
Inventor: **Kagan, Oren**
Rachel 2
Haifa(IL)
Inventor: **Kolsky, Amir**
39 Klebanov
Haifa(IL)
Inventor: **Medan, Yoav**
25 Hankin Street
Haifa(IL)
Inventor: **Walach, Eugene**
11 Netanyahu Street
Kiryat, Motzkin(IL)

Representative: **Lloyd, Richard Graham**
IBM (UK) Ltd,
UK Intellectual Property Department,
Hursley Park
Winchester, Hampshire SO21 2JN (GB)

Data entry system.

A data entry system generates an electronically stored coded representation of a character sequence from one or more electronically stored document images, comprising optical character recognition logic (90) for generating, from the document image or images, character data specifying one of a plurality of possible character values for corresponding segments of the document images; characterised by interactive display apparatus comprising: means (110) for generating and sequentially displaying, one or more types of composite image, each composite image comprising segments of the document image or images arranged according to the character data, and a correction mechanism responsive to a user input operation to enable the operator to correct the character data associated with displayed segments.

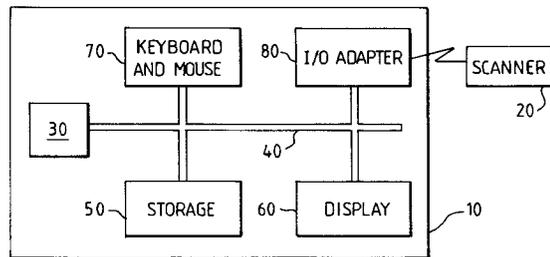


FIG. 2

EP 0 602 821 A2

The invention relates to a data entry system for the extraction of character information from scanned document images.

With the advance of computerised data processing it has become common practice to generate information in electronic form via the automated recognition of information from scanned documents. In other words, data acquisition is performed as follows. First documents are scanned, then predetermined fields are extracted and passed to recognition modules for translation of an image into a coded form, such as ASCII.

However, despite the fact that this recognition can now be performed for both printed and handwritten data, almost always the recognition process results in a certain number of errors and rejects.

Hence, there remains a need for manual checking and correction of recognised data. This is generally performed by the display to a human operator of the scanned images and the associated optical character recognition results. The human operators perform comparison and correction where necessary.

Such systems are disclosed in US 5,025,484, EP-A-O,107,083 and GB 2,149,171.

Conventionally, there are two approaches to this process. Either an operator is asked to review all the fields and key-in errors and rejects discovered in the OCR process, or they are asked to view and correct only those fields where at least one error or reject has been discovered.

In both cases, the process is as follows. A problematic field is displayed on the screen with the recognition results being displayed in the vicinity of the image. Problematic characters are generally emphasized in some way, for instance the cursor may stop under the first problematic character or reversed video mode may be used to emphasize the problematic characters. The operator checks and if necessary corrects the emphasized characters and the process is repeated until all the problematic characters are resolved.

A major disadvantage of the prior art is that the fields are viewed by the operator in the context of the original document image. Assuming the optical character recognition produces 5% of rejected characters and that each field has an average of 20 characters. In such a case about 35% of the fields will have at least one rejected character.

It follows that, in the prior art, in at least 35% of cases would the operator be required to display the field image, focus on it mentally, identify the problem and correct it. Thus prior art OCR assisted data-entry methods improve the productivity of data entry by at most a factor of 3 over a purely manual data entry process.

The present invention is directed to improving yet further the efficiency of OCR-assisted character

data entry.

Accordingly, the invention provides a data entry system for generating an electronically stored coded representation of a character sequence from one or more electronically stored document images, comprising optical character recognition logic for generating, from the document image or images, character data specifying one of a plurality of possible character values for corresponding segments of the document images; characterized by interactive display apparatus comprising: means for generating and sequentially displaying, one or more types of composite image, each composite image comprising segments of the document image or images arranged according to the character data, and a correction mechanism responsive to a user input operation to enable the operator to correct the character data associated with displayed segments.

The invention enables a substantial reduction of the time required to correct and complete the data entry of documents following an OCR process.

Preferably, the arrangement of the segments in one type of the composite images is into groups of segments for which the character data specifies the same character value.

The displaying of composite images to the operator made up entirely of characters which have been recognised by the OCR logic as being of the same type enables errors to be rapidly recognized and marked on an exception basis. Once recognised, these errors can then be corrected either immediately or during correction of characters rejected by the OCR logic.

Where the character data comprises a confidence value associated with each specified character value, which confidence value is indicative of the likelihood that the specified character value be correct, the composite images can comprise segments having corresponding confidence value indicating a low such likelihood.

For extra speed, the segments can be arranged in groups of n segments in such composite images, the operator being able to set the confidence value of each of the n segments in a group to indicate a high likelihood that the character code be correct through a single user input operation.

If the documents being processed are such that the character codes are arranged in fields according to predetermined criteria, such as location on the document, the correction mechanism can advantageously comprise logic for generating one or more composite images comprising groups of segments of the document image associated with fields of the same type and displaying the composite images for correction by the operator. This enables correction of those few errors which can only be corrected in context.

In one embodiment, the predetermined user input operation is actuation of the mouse at a time when a screen cursor is located over the displayed segment and the display of the composite images is via a video display unit, such as a cathode ray tube monitor.

Viewed from another aspect, the invention provides a method of operating data entry apparatus to generate an electronically stored coded representation of a character sequence from an electronically stored image of a document, the method comprising: generating character data specifying one of a plurality of possible character values for corresponding segments of the document images using optical character recognition logic; generating and sequentially displaying, one or more types of composite image from the character codes and the document image, each composite image comprising the segments of the document image arranged according to the character data; selecting segments of the document image displayed as part of the composite images; and correcting the character data corresponding to the selected segments to generate the coded representation.

In a particularly advantageous form of the invention, each of the document segments has a confidence value associated therewith indicative of the likelihood that the character code associated therewith be correct, and the character codes are arranged in fields each comprising one or more characters, the method comprises generating and sequentially displaying a plurality of composite images of a first type wherein the arrangement of the segments is into groups of segments for which the character data specifies the same character value; selecting segments displayed in the first type of composite image and setting the corresponding confidence value to indicate a low likelihood of the specified character value being correct; and generating and sequentially displaying a plurality of composite images of a second type comprising segments having a confidence value indicating a low likelihood associated therewith; selecting segments of the document image displayed as part of the composite images of the second type; and correcting the character data corresponding to the selected segments to generate the coded representation.

To correct any remaining errors the method can comprise generating and sequentially displaying a plurality of composite images of a third type comprising groups of segments associated with fields having the same type value associated therewith; selecting segments of the document image displayed as part of the composite images of the third type; and correcting the character data corresponding to the selected segments to generate the coded representation.

An embodiment of the invention will now be described, by way of example only, with reference to the accompanying drawings, wherein:

Fig. 1 shows the apparatus used in the embodiment;

Fig. 2 is a schematic view of the apparatus of fig. 1;

Fig. 3 is a schematic diagram of the software elements used;

Fig. 4 shows a basic key-in process module;

Fig. 5 is a flow chart showing the sequence of operation;

Fig. 6 shows an example of a composite image.

The data entry system of the embodiment of the invention is shown generally in Fig. 1. It comprises personal computer 10 which is connected to scanning device 20. Fig 2. is a schematic diagram showing the relevant parts of computer 10 and scanner 20. Processor 30 is connected via bus 40 to disk storage device 50, cathode ray tube display device 60, user input devices - keyboard and mouse - shown generally at 70. Processor 30 is connected via bus 40 and input/output adapter 80 to scanning device 20.

Scanning device 20 is a conventional device of a type well known in the art. It is used to scan document pages placed therein and pass digitised images thereof to computer 10 which are stored in storage device 50. Image digitisation is achieved in scanner 20 by measuring the brightness of light reflected from small rectangular areas of the pages, the areas being arranged in rows and columns to cover the whole page, and comparing the brightness levels corresponding to each area with a threshold parameter. In the digitised image produced by the scanner, pels above the threshold are set to 1 and pels below set to 0.

Stored document images are processed to generate a sequence of coded character data, which contains the character information extracted from the document images but which is in a form, such as ASCII, which is easier to store and manipulate in the computer.

Fig 3 is a schematic diagram showing the software elements used in this embodiment. OCR logic 90 takes the document images as input and generates a sequence of character codes and associated confidence values which are stored along with the segmented character images in database 100. Database 100 contains and manages the following information for each document:

1. The document image.
2. an image of the all the segmented characters that were processed by the OCR logic 90.
3. a record of the recognition results, ie the character codes and confidence values.
4. a form definition file which contains the information from which the location of each char-

acter and OCR result can be computed. This is actually a list of field names with their geometric location on the form, their type (numeral, alpha, etc.) and the number of characters in the field.

Basic Key-in Process (BKP) modules 110, 120 and 130 extract and modify the database data as will be described below.

The data entry process generally comprises five stages:

1. The document images comprise a number of fields each corresponding to a particular pre-determined area of the document image. Each field to be recognised is segmented into images of individual characters to be recognised.
2. The individual characters are fed into OCR logic 90 which either recognises the characters, possibly with errors, or rejects them. Many methods of optical character recognition are well known in the art and the details of the operation of the character recognition logic not be described further herein. The OCR logic generates for each character image a character code and a confidence value. Those characters generated with a low confidence value are considered rejected.
3. Images of individual characters are pasted together to form composite mosaic images suitable for optimal operator productivity. Thus characters are merged from different fields and even different documents into the composite images.
4. The character codes associated with these composite images are corrected by an operator in a key-in session.
5. Data entry results are recast into their proper sequence within the original document for final processing.

The success of the process depends on being able to achieve higher operator productivity for the correction of the composite images than for the display of the individual characters in the context of the document itself.

In the embodiment three different types of composite image are used in different stages of the correction process. These types of composite image have been designated by the inventors Exemption Data Entry, Memory Span Data Entry and Field Type Data Entry.

1. Exemption Data Entry.

The characters are extracted, sorted and displayed in composite images according to the recognition results. All the characters recognised with a particular value of the character code are grouped together in the composite images, ie the images of all the characters recognised as '1's are displayed together as a composite image, likewise

the images of all the characters recognised as '2's and so on.

These results are verified for correctness by the operator. Errors are marked using a pointing device, ie by moving the cursor onto the incorrect image and clicking the mouse, and tagged as rejects. Using this method all OCR errors can be tagged.

2. Memory Span Data Entry

A composite image is formed composed of n rejected character images displayed on the screen simultaneously. The rejected character images can either be rejected by the OCR engine, ie recognised with a low confidence value, or they can have been tagged as errors by the exemption data entry procedure.

The operator can then manually correct these images. The parameter n is chosen so that the display screen format is easy for an operator to grasp to maximise data entry speed. It has been found by the inventors that a value of $n = 3$ is optimal. However, this number is adjustable by the operator according to their preferences.

It is advantageous to perform this memory span data entry of data sorted according to the OCR results so that the fields can be corrected by the exemption data entry and it is easy to identify those sets of rejected characters where no correction is required. If 80% of the characters recognised by the OCR with a low confidence level are actually correct, then for $n = 3$, in 50% of the sets the OCR results will be correct and a single key-stroke will be required to accept the entire set.

Field Type Data Entry.

A few errors can only be corrected by an operator in their context. This type of data entry is designed to handle these errors. Characters are grouped into fields of the same type and the operator can use the information they have about the field-type to assist in correcting the error. For example in a given batch of documents, all the dates are grouped together, all the names are grouped together and so on. As a result it is easier for an operator to concentrate on a certain constant mode of operation. Moreover there is an increased chance of similarity between fields, for example if many documents were generated on the same date it is easy to fill in this date on all the fields of this type and then only exceptions need to be corrected.

In the method of operation used in the embodiment, these three types of data entry are combined in an advantageous way as follows.

1. OCR operates on all characters. Some characters may be rejected by this process.
2. Recognized characters are reviewed using the exemption data entry technique. In this way all OCR errors are detected and tagged.
3. Rejected characters and those tagged as errors are corrected using the memory span data entry sorted into fields of the same character type.
4. The few remaining characters which could not be corrected by step 3 are corrected in context by the field-type data entry.

The structure of each BKP module is shown in Fig. 4. They comprise the following three modules:

1. Pre-processing module 140. This module takes the set of document images and rearranges the characters in all these documents according to some input rules. It is a data base query processing module which interrogates the data base 100 of document images and OCR results, and the result of a query is a set of composite images consisting of characters of the database which answer the query. The queries accepted by this module are:

- a. "Take all characters recognized as 'x' with high confidence. and arrange them in a mosaic" ('x' is some character). In this case, the result of this query would be all the characters in the set which were recognized as 'x' arranged in a set of images, each of which is a matrix of character images. The characters in the mosaic would eventually be all mixed up, each coming from one form or another, but all were recognized as 'x'. This query refers to the "key-in by exemption" method described above.
- b. "Take all characters recognized with low confidence, and arrange them in a mosaic". In this case, each composite image consists of n characters, where n is a parameter which can be set by the user. This query refers to the "key-in by memory span" method described above.
- c. "Take all fields labelled 11,12,13,..., and arrange them in a mosaic", where 1i is a symbolic name of a specific field in the form. In this case, the composite image consists of a full field. "Take all fields representing a year" is an example of such a query, which results in mosaics consisting of fields representing years (year of birth, year of current date, etc). This query refers to the "key-in according to the field type" method described above.

The output of preprocessing module 140 comprises the following elements:

- a. the composite images.

b. the corresponding OCR results, arranged according to the mosaics.

c. pointers back to the data base for enabling its update, once the correct recognition value for each character in the mosaics has been obtained. These pointers are generated in a file containing this information,i.e., for each character in the mosaic, from which document, and which character in that document it was taken. For example, if a key-in by exemption mosaic is generated, the character in location i,j in the mosaic is the k'th character in field m in document n.

2. Key-in user interface 150. This module receives the pre-processing module's output as input and, after a key-in session, generates a set of key-in results for the mosaics. User interface 150 is a Graphical User Interface which displays the composite image to be keyed-in and which provides data entry boxes for entering the key-in data, a data entry box for each element in the mosaic. The key-in results have exactly the same format as the OCR results.

3. Post-processor 160 which takes as input the key-in results and the pointer file and updates the recognition results of the initial forms data base according to the key-in results.

These three modules are connected to each other according to Figure 4. The pre-processing module 140 generates the mosaics, and a set of pointers back to the data base. From these, the key-in user interface 150 generates a set of keyed-in results. These are passed to the post processing module 160 and the OCR records are updated using the pointer file from the pre-processing module.

Three different BKP's 110, 120, 130 are constructed using the pre-processing modes described in 1a, 1b, and 1c and connected as described in the flow chart of Figure 5. After processing 170 by the OCR logic, BKP with key-in by exemption is performed 180 so that in the data base 100 all characters having high OCR confidence which were actually found to have a wrong OCR result are set to low confidence ie tagged as rejects.

Fig 6. shows an example of a composite image produced during key-in by exemption. The composite image comprises all the characters which have been recognised as the letter 'O' by the OCR. Errors have been identified and marked in this example, as shaded segments.

Next, BKP with key-in by memory span is performed 190 so that in data base 100 all characters having low confidence are corrected as far as possible. Finally, BKP with key-in by field type is performed 200 to correct those few errors which need to be corrected in context to generate a correct data base.

The embodiment has been described in terms of a suitably programmed general purpose computer used in conjunction with a conventional scanning device but, in practice, it will be understood that the invention could be implemented in hardware, either as part of the scanner or as a special purpose adapter for use with the computer or a stand alone hardware device or be implemented as any combination of hardware and software.

It will be appreciated that many variations of the apparatus and methods describe here are possible within the scope of the attached claims. The methods described here could work with a wide variety of supplementary modules for field extraction, OCR of either printed or handwritten data, character segmentation, image display and key-in.

The interactive display apparatus could take many forms. For example the composite images could be displayed by being printed using a printing device. The operator could then mark errors on the printed sheets and feed the marked sheets into a scanning apparatus which would detect the location of the marked errors and update the character data accordingly. Alternatively, the location of the errors on the printed copy could be keyed-in by hand.

The inventors have used and tested their apparatus for the extraction of data from hand printed filled-in forms and for the recognition of addresses from digitised address labels in the context of automated mail sorting. However, it will be understood there may be wide variation in the types of data to be keyed-in.

In the described embodiment the OCR logic associates a low confidence value with a rejected character, some OCR methods provide a second best character code, which was recognized with the next highest confidence value. This information could be used in the correction mechanism in order to speed up the data entry process by presenting the second choice to the user and enabling them to accept or reject this with a single key-stroke.

Various other logical tests may be applied to fields such as dates, checksums and other code sums. Also, values that exist in databases or dictionaries, such as ID numbers, can be verified for validity. These tests can improve both the speed and quality of the data entry process and eliminate manual intervention. Accordingly, they may be used before and after each of the above data entry stages as an automatic error correction/detection mechanism.

For applications where accuracy is of paramount importance, the method can be used as part of a double pass process. Results of data entry using this present method would be compared with results obtained via manual data entry, with any

differences being pointed out to and corrected by an operator.

The following is an approximate analysis, provided for the purposes of illustration only, of the performance improvement which can potentially be achieved using the present approach.

Assume that a given recognition package rejects 5% of characters and accepts the remainder, with 1% substitution error and that 10,000 documents having 20 characters in a single field have to be keyed-in and that the operator's productivity is 10000 key-strokes per hour, yielding an average of 1s for correction of each field.

Using manual data entry 200,000 character have to be keyed-in using a double pass. Each pass requires 20 hours of operator time, yielding a total of 40 hours operator time.

Using the prior art OCR assisted data entry all data would need to be reviewed manually to pick up the 1% substitution errors from the OCR. Assuming a first pass of OCR assisted data entry with correction of fields containing rejected characters, followed by a second manual pass of all the data,

$65\% = 100\% \times (1 - 0.95^{20})$
of fields must be viewed in the first pass. Since we have assumed 1s for the correction of each field, this would take

$10,000 \text{ docs} \times 0.65 \times 1\text{s} = 1.8 \text{ hrs.}$

The second manual pass would take 20 hours as above and thus the total time would be 21.8 hours

Using the present technique starting with exemption data entry of all recognised characters ($200,000 \times 0.95 = 190,000$ characters) if each composite image contains 80 characters and is processed in 1s $190,000/80 = 2,375$ screens are needed taking 0.66 hours. As a result of this 1% (1900) of characters are tagged as rejected. These characters are added to the 10,000 recognition rejects to yield a total of 11,900 characters which need correction using memory span data entry. This would require $11,900/10,000 = 1.19$ hours. Therefore the total operator time is $0.66 + 1.19 = 1.85$ hours.

It can thus be seen that the invention provides a dramatic improvement over both the double pass fully manual data entry and the double pass OCR assisted data entry, both of which may be expected to produce an equivalent precision level.

Claims

1. Data entry system for generating an electronically stored coded representation of a character sequence from one or more electronically stored document images, comprising optical character recognition logic (90) for generating, from the document image or images, character data specifying one of a plurality of

- possible character values for corresponding segments of the document images; characterised by interactive display apparatus comprising:
 means (110) for generating and sequentially displaying, one or more types of composite image, each composite image comprising segments of the document image or images arranged according to the character data, and a correction mechanism responsive to a user input operation to enable the operator to correct the character data associated with displayed segments.
2. Data entry apparatus as claimed in claim 1 wherein the arrangement of the segments in one type of the composite images is into groups of segments for which the character data specifies the same character value.
 3. Data entry system as claimed in claim 1 or claim 2 wherein the character data comprises a confidence value associated with each specified character value, which confidence value is indicative of the likelihood that the specified character value be correct, wherein the composite images comprise segments having corresponding confidence value indicating a low such likelihood.
 4. Data entry system as claimed in claim 3 wherein in the composite images the segments are arranged in groups of n segments wherein the correction mechanism enables the operator to set the confidence value of each of the n segments in a group to indicate a high likelihood that the character code be correct through a single user input operation.
 5. Data entry system as claimed in any preceding claim wherein the character data are arranged in fields each comprising one or more characters, the character data comprising a type value associated with each field, and wherein one type of the composite images comprise groups of segments associated with fields having the same type value associated therewith.
 6. Data entry system as claimed in any preceding claim wherein the predetermined user input operation is actuation of the mouse at a time when a screen cursor is located over the displayed segment.
 7. Data entry system as claimed in any preceding claim wherein the display of the composite images is via a video display unit.
 8. Method of operating data entry apparatus to generate an electronically stored coded representation of a character sequence from an electronically stored image of a document, the method comprising:
 generating character data specifying one of a plurality of possible character values for corresponding segments of the document images using optical character recognition logic;
 generating and sequentially displaying, one or more types of composite image from the character codes and the document image, each composite image comprising the segments of the document image arranged according to the character data;
 selecting segments of the document image displayed as part of the composite images; and
 correcting the character data corresponding to the selected segments to generate the coded representation.
 9. Method as claimed in claim 8 for operating data entry apparatus in which the character data comprises a confidence value for each specified character value, which confidence value is indicative of the likelihood that the character code associated therewith be correct,
 the method comprising generating and sequentially displaying a plurality of composite images of a first type wherein the arrangement of the segments is into groups of segments for which the character data specifies the same character value;
 selecting segments displayed in the first type of composite image and setting the corresponding confidence value to indicate a low likelihood of the specified character value being correct; and
 generating and sequentially displaying a plurality of composite images of a second type comprising segments having a confidence value indicating a low likelihood associated therewith;
 selecting segments of the document image displayed as part of the composite images of the second type; and
 correcting the character data corresponding to the selected segments to generate the coded representation.
 10. Method as claimed in claim 9 of operating a data entry system in which the character data is arranged in fields each comprising a one or more characters and comprises a type value associated with each field, which type value can take one of a plurality of values, compris-

ing

generating and sequentially displaying a plurality of composite' images of a third type comprising groups of segments associated with fields having the same type value associated therewith; 5

selecting segments of the document image displayed as part of the composite images of the third type; and

correcting the character data corresponding to the selected segments to generate the coded representation. 10

15

20

25

30

35

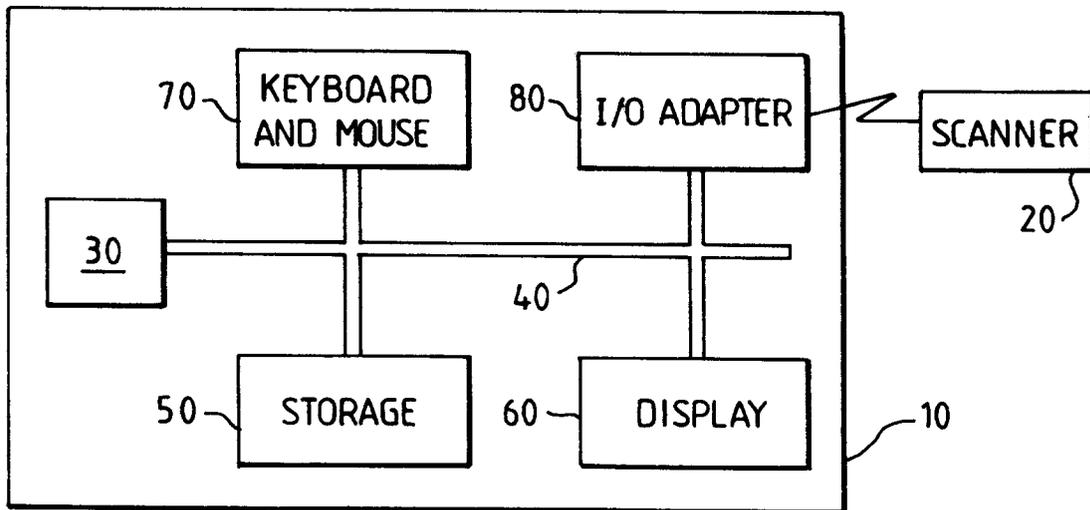
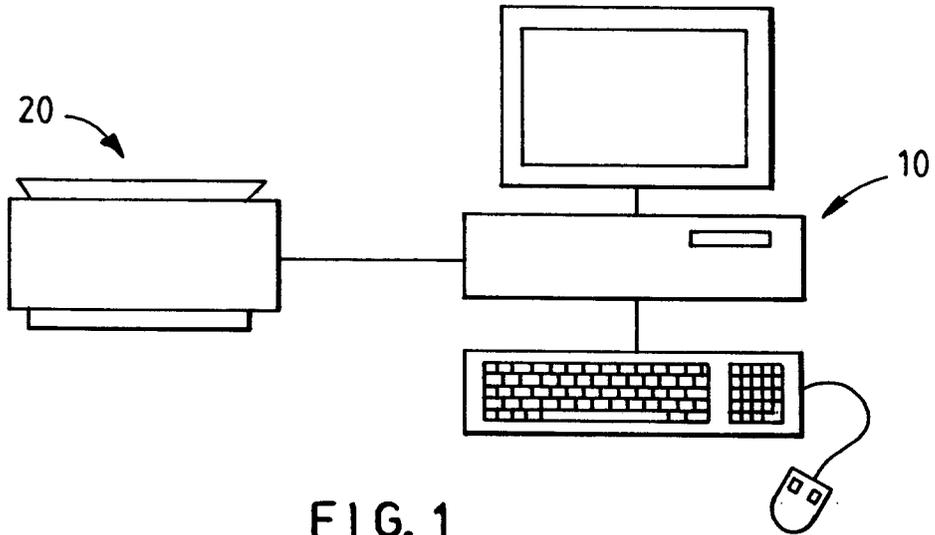
40

45

50

55

8



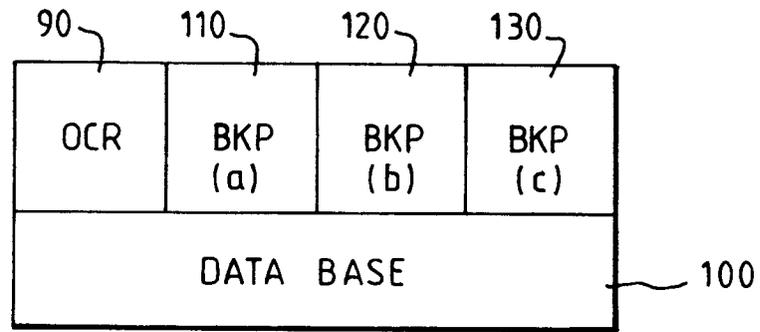


FIG. 3

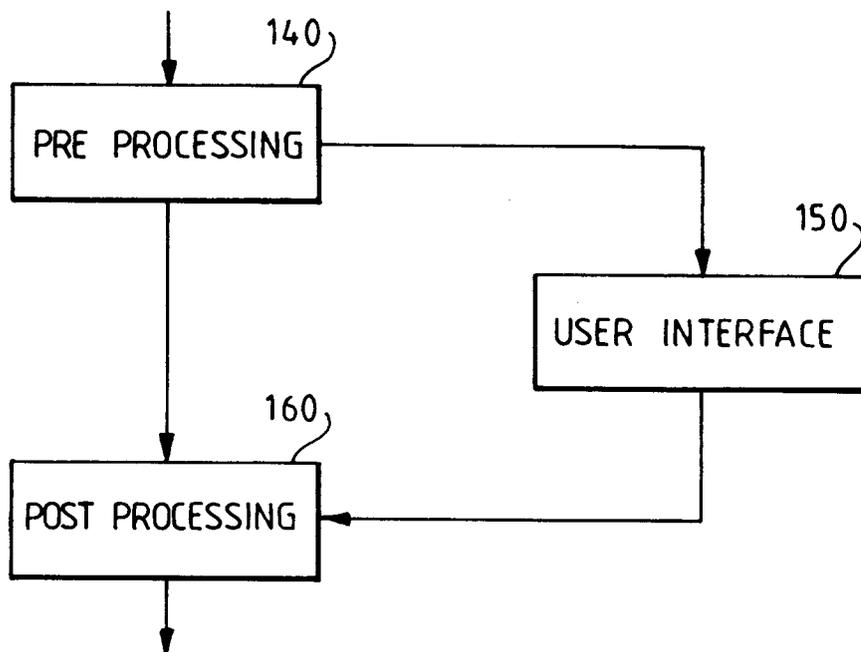


FIG. 4

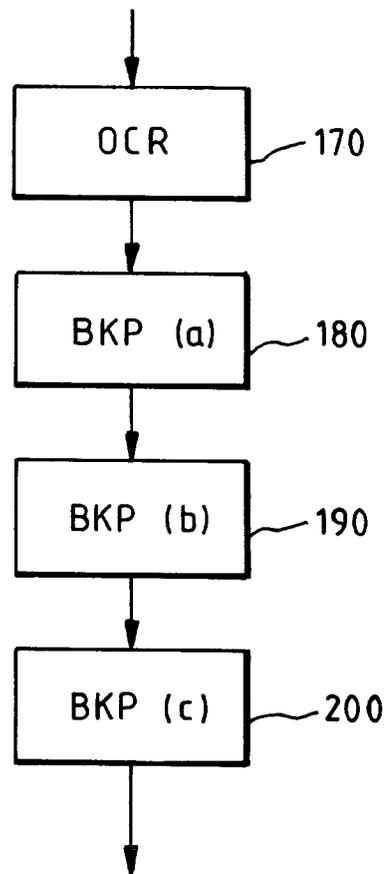


FIG. 5

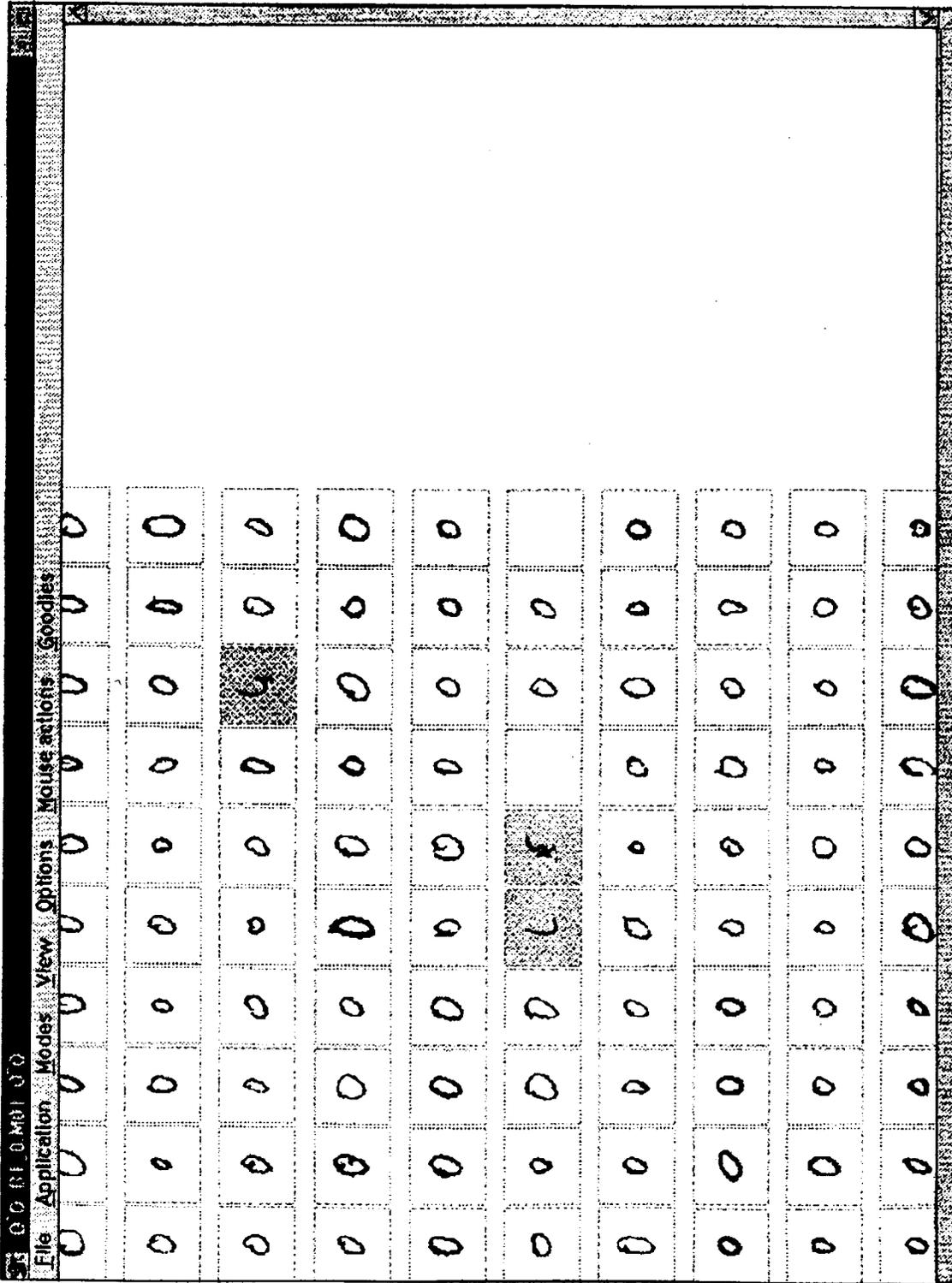


FIG. 6