

A Constraint-based Approach to Postal Address Interpretation

Peter B. Cullen, Jonathan J. Hull, Sargur N. Srihari
Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo
226 Bell Hall
Buffalo, NY, 14260

Abstract

An approach, based on the constraint satisfaction paradigm, for postal address interpretation is described. A constraint satisfaction problem (CSP) is defined as a set of variables, each with its own domain of values and a set of constraints that specify legal combinations of variables and values. The solution to a CSP is a value assignment to each variable such that all the constraints are satisfied. Address interpretation can be formulated in this way. Segmented words from the address block are treated as variables and these word images are recognized as belonging to a small subset of a lexicon. These subsets or neighborhoods are then the domains of the variables. The constraints are provided by the ZIP+4 database which lists all valid addresses in the country. In our approach, an intelligent backtracking algorithm is applied to find the most compatible, highest confidence recognition choice for each word. By formulating the problem as a CSP, we can overcome some of the uncertainties associated with address interpretation.

A system that uses this constraint-based approach to address interpretation was recently tested on 1200 machine printed address blocks. This test set was skewed more towards difficult to read addresses. The system was able to assign ZIP+4 Codes to 82% of the addresses. A current postal OCR was only able to assign ZIP+4 Codes to 28% of the same addresses.

1 Introduction

Postal address interpretation is a problem of central importance in developing an address recognition system. This paper describes an approach, based on the constraint satisfaction paradigm, to interpreting postal addresses. A constraint satisfaction problem (CSP) can be thought of as a special kind of search problem. A (CSP) is defined as a set of n variables, each with its own domain of values and a set of constraints that specify the legal combinations of variables and values. A solution to a CSP is a value assignment to each variable such that all the constraints are satisfied. Solutions are usually found by employing some type of backtracking algorithm that searches the space of variable/value combinations. Network consistency algorithms can also be used to enforce local consistency among variables. These remove values from variable's domains that can not possibly participate in an overall solution to the problem.

Postal address interpretation can be formulated as a CSP. In such a formulation segmented words from the address are treated as variables. These words are then recognized to create their respective domains. The constraints are provided by the ZIP+4 database which lists all valid addresses in the country. This database can tell whether the values assigned to one or more variables are part of a legal address.

Postal address interpretation has been formulated as a general search problem; however, such a formulation limits the ability to overcome uncertainties that exist in the interpretation process. An address may be incomplete, in which case a partial match must be allowed. An address may contain incorrect information which must be compensated for in some way. Also, the automatic recognition results for each word may not provide the correct choice. By formulating address interpretation in the constraint satisfaction paradigm, each word is treated independently. Thus, if uncertainties do occur, the errors they may create are not propagated. In fact, with good recognition of the remaining information in the address, the constraints will point out the incorrect information. The incorrect piece of information can then be removed from consideration, allowing an interpretation to be found.

The rest of this paper describes the constraint-based approach to postal address interpretation used in our system to read machine-printed addresses. Section 2 describes the problem of postal address interpretation using a real example. Section 3 describes the constraint satisfaction problem. Section 4 then details how postal address interpretation can be formulated as a CSP and how uncertainties can be overcome in doing so. Finally, section 5 gives some experimental results and section 6 gives some conclusions of this work.

2 Postal Address Interpretation

The goal of postal image interpretation is, given an image of an address, identify the nine-digit ZIP Code that matches the address, even though the address does not contain a nine-digit ZIP Code. The finite list of possible interpretations for this problem is given in the form of a database of nine-digit ZIP Codes, called the ZIP+4

database. The national ZIP+4 database consists of 28 million records, each of which corresponds to a different nine-digit ZIP Code. The idea is to isolate a particular ZIP+4 record for a given address.

For example, the postal address shown in Figure 1 exactly matches the ZIP+4 record shown. The address image is first segmented into words. These word images are then recognized, usually character by character. The character recognition results are then combined and used to match into the ZIP+4 database. The result for this example would be the ZIP+4 Code of 76012-2499.

ARLINGTON HOSPITAL
800 S RANDOM HILL RD
ACCOUNTS PAYABLE
ARLINGTON TEXAS 76012

76012-2499 800 S RANDOM HILL RD

Figure 1: Example Postal Address Image and Matching ZIP+4 Record

Provided the character recognition results are reliable, postal address interpretation becomes a problem of matching against the ZIP+4 database. However, there are usually several choices for each character and, therefore, several word-choices for each word image. Thus, in most cases, the interpretation process must search the ZIP+4 database using the word choices [1].

The need to search the ZIP+4 database is only one of the difficulties in postal address interpretation that must be overcome. The example shown in Figure 2 exhibits other types of uncertain information that make can make address interpretation difficult.

Jack Dannel
Southwest States Bankcard Assoc.
4555 Beltline Rd.
DALLAS TX 75001

Figure 2: Postal Address Containing Uncertainties

In this example the city name has been obscured, possibly during image acquisition. The author of this address has made an error in the ZIP Code. Also, due to poor image quality, the automated recognition system may fail to correctly recognize the suffix. In this case, not only does the interpretation process have to search the

ZIP+4 database, but it will need to overcome the problems mentioned above. One way to make this process easier is to constrain the word recognition lexicon by the recognition results of previous words. However, this can lead the entire search away from the correct answer. In Figure 2 the ZIP Code should be easy to recognize since it is written clearly. However, as it turns out, the given ZIP Code does not correspond to the street address. Constraining recognition results by the ZIP Code, in this case, may lead the search process astray. Therefore, we have used an alternative solution.

Our solution has been to recognize each word independently, then use the constraints provided by the ZIP+4 database to search for the correct ZIP+4 Code. In the case of the previous example, the correct city name can be inferred by using the available information to locate the correct ZIP+4 record. Additionally, our system would use the state and street information to overcome the recognition error as well as the ZIP Code error. This process will be examined in more detail in the following sections.

3 Constraint Satisfaction

A constraint satisfaction problem (CSP) consists of n variables, $x_1 \dots x_n$, each having an associated domain, $d_1 \dots d_n$ [2]. Each variable can take on any value in its domain but *constraints* are imposed on the legal combinations of values that two or more variables can take on. A solution to the constraint satisfaction problem consists of a set of values for the X_i such that all the constraints are simultaneously satisfied.

The discovery of a solution to a CSP is often performed by a search-like algorithm. Techniques such as backtracking [3] [4], network consistency [5] [6] [7] and associated heuristics [8] have been used before.

These techniques solve the CSP, but they do not solve the interpretation problem. Our solution to the interpretation problem that uses a form of backtracking is shown in Figure 3. After the words have been segmented and recognized, the top recognition choice is instantiated for each variable. The domains for each variable are then sorted in increasing order by the number of conflicts each value has with the instantiated values of the other variables. A conflict between two variables occurs when their values are not compatible according to the constraints. If two values in a domain have the same number of conflicts, the one with the higher recognition confidence is placed first. The constraint satisfaction problem is then solved using the informed-backtracking algorithm and the min-conflicts heuristic [9]. The min-conflicts heuristic dynamically orders the values in each variable's domain.

If the instantiated values of all the variables satisfy all the constraints, then a solution has been found. Otherwise, a backtracking routine chooses a variable for repair that participates in the most conflicts. This repair process consists of re-sorting that variable's domain in increasing number of conflicts with the instantiated values of the other variables. The informed-backtracking routine then iterates through the sorted domain, just like a regular backtracking algorithm, instantiating its top choice and checking for consistency among all instantiated values.

If the informed-backtracking algorithm fails to find a solution to the CSP, the

Input: image, list of interpretations

Segment image into possible objects

Recognize objects (domain for each variable)

Sort domains

While (solution not found)

Informed-backtrack(variables)

If (no solution) **then**

Failure analysis using *max-conflicts* heuristic

endif

end while

Infer values for variables that have been removed or are missing

Figure 3: Image Interpretation Algorithm to Resolve Uncertainties

system uses the max-conflicts heuristic to determine the variable(s) that caused the failure. The failure inducing variables are termed *culprits*. The max-conflicts heuristic indicates the variable(s) that initially had the most conflicts among the top recognition choices, after the initial domain sortation. The results of the failure analysis are either to delete the culprits or modify their domains using contextual knowledge derived from the informed-backtracking routine. This process of backtracking and failure analysis continues until a solution is found or all variables have been modified or deleted. Once a consistent set of values is found, the interpretation is made based on those values. If, in solving the problem, some variables were removed or if some objects were missing, their values can be inferred from the determined solution.

The constraint satisfaction approach combined with failure analysis is able to resolve various uncertainties in the interpretation process. When an incomplete address is input to the system, it is able to use partial information to make an interpretation. Many addresses contain redundant information so the loss of some words will not prevent a correct interpretation from being found. If only a partial set of words are present, fully recognizing them provides enough specific information to make a complete interpretation of the image. Once a consistent set of values are identified, values for the missing or removed objects can be inferred from the available information and a complete interpretation can be made.

If patron errors have been introduced in the address, failure analysis can remove those variables from consideration. If this is the case, it is most likely that the recognition results for these words will not make a legal combination with the choices for other true words. The proposed CSP algorithm can identify patron errors, since they should conflict with many of the recognition choices for the other words. Once identified as being incorrect, their corresponding variables in the CSP can be assigned values based on the recognition results for the other variables.

Finally, recognition errors can be corrected during failure analysis by modifying a variable's domain or deleting it and inferring the correct choice later. Even though the informed-backtracking algorithm may fail to find a solution to the CSP, it can still provide useful information as to what a solution may look like. By removing it, the

other variables may form a consistent solution from which the value for the removed variable can be inferred. Even if this is not the case, this partial solution may limit the number of choices that the removed variable can assume. The word can be re-recognized with the set of possible choices limited to what the partial solution predicts. Thus, contextual knowledge from the address can be used to improve recognition results.

4 Postal Image Interpretation

Postal image interpretation can be stated to fit the previously outlined method of interpretation. The address block images consist of words which can be segmented using low-level image processing techniques. The segmented words from the image become the variables in the constraint satisfaction formalism. Once the words are segmented, our system uses a hierarchical recognition process that first identifies the class that the object belongs to and then recognizes the object within that class. We refer to the first level as *parsing* and to the second level as *word recognition*. The parser uses a statistical method to assign a type to each word such as ZIP Code or street name [10]. Figure 4 shows a typical machine-printed address block, along with the segmented words and their categories attached to the upper right of each word.

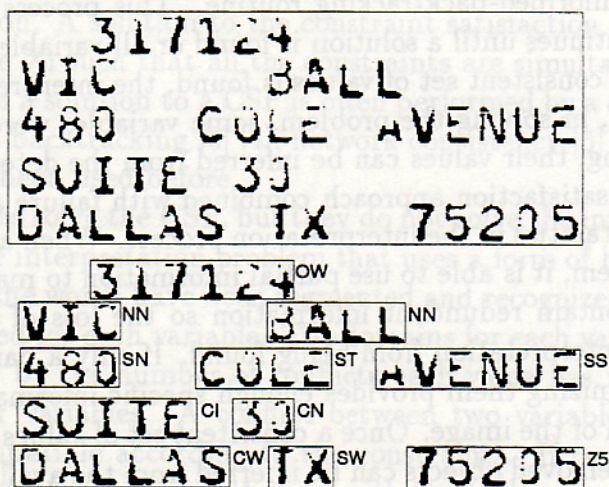


Figure 4: Word Segmentation and Parsing Results

The chart in Figure 5 shows a partial listing of the various word category abbreviations.

The various words are then recognized within their respective types. That is, city words are recognized within a lexicon of city words, street words are recognized within a lexicon of street words and so on. The lexicons are only constrained by the particular application; in this case, city words and street words extracted from the ZIP+4 database. A distinction is made between words that contain digits (ZIP Codes, street numbers, etc.) and words that contain alphabets (cities, street names,

| Abbrev. | Meaning |
|---------|-------------------------------|
| CW | City Word |
| SW | State Word |
| Z5 | ZIP Code |
| Z9 | ZIP+4 Code |
| SN | Street Number |
| ST | Street Word |
| SS | Street Suffix |
| CI | Building Secondary Identifier |
| CN | Building Secondary Number |
| NN | Personal Name |
| OW | Other Word |
| FW | Firm Word |

Figure 5: Word Category Abbreviations and Meanings

etc.). The first kind are input to a digit recognition routine, while the second kind are recognized using a multiple classifier approach [11]. Both methods return a ranked list of choices for each word. This list of choices becomes the domain for the appropriate variable.

The domains are sorted by 1) the min-conflicts heuristic and 2) the recognition confidence. Each variable is then instantiated with its top choice. The informed backtracking procedure then tries to find a consistent set of values from the recognition choices for all the variables that satisfy the constraints imposed by the ZIP+4 database. Figure 6 shows the constraint graph provided by the ZIP+4 database. This fully connected graph shows how every word in the image constrains the choices for every other word. For example, given a particular ZIP Code (z5), only certain street suffixes (ss) occur within that ZIP Code (ie. those suffixes are compatible with that ZIP Code). These are examples of the binary constraints used in counting conflicts. In addition, this graph illustrates the queries [12] that can be made on a set of variables to test whether a solution satisfies all the constraints. If a solution cannot be found, any errors in the image caused by the aforementioned uncertainties are resolved by failure analysis.

As mentioned in the previous section our system is able to resolve the various uncertainties associated with address interpretation using failure analysis. Incomplete images can be the result of poor image acquisition, or the author of the mailpiece (the patron) not providing enough information. The patron can also make an error in writing the mailpiece, whether intentional or otherwise. They may write "St" instead of "Dr" for the suffix or accidentally transpose two digits in the ZIP Code, for example. Our recognition process may have also made errors. Failure analysis is able to identify each of these kinds of errors and make an attempt at correcting them.

Referring again to Figure 2, assume the following variables and recognition choices: X1=4555 for the street number, X2=Beltline for the street name, X3=Dr,St for the street suffix, X4=REJECT for the city name, X5=TX,TN for the state name, and

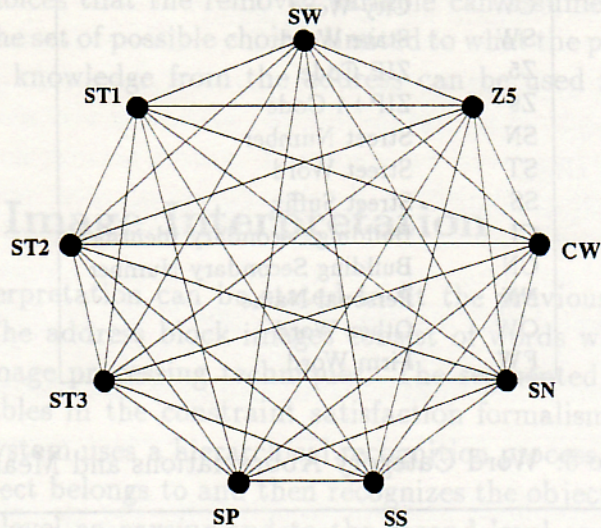


Figure 6: ZIP+4 Constraint Graph

X6=75001 for the ZIP Code. Each variable is then instantiated with its top choice. When informed-backtracking is called with these variables and values, it is found that they are not consistent so the number of conflicts each variable is participating in is counted. X4 is left out of this since it has no values. This process is done by querying the ZIP+4 database on the various combinations of variables, taken two at a time. For example, the street number(X1) whose instantiated value is "4555" agrees with the street name and the state name, but disagrees with the street suffix(X3) and the ZIP Code. What this means is that there is at least one record in the ZIP+4 database with a street number of "4555" and a street name of "Beltline", at least one record with a street number of "4555" and a state name of "TX", and that there is no record that has a suffix of Dr and a street number of 4555, nor is there one that has a ZIP Code of 75001 and a street number of 4555.

In terms of constraints, "4555" and "Beltline" satisfy the constraint between the street number and the street name, "4555" and "TX" satisfy the constraint between the street number and the state name, but "4555" and "Dr" do not satisfy the constraint between the street number and the street suffix, nor does "4555" and "75001" satisfy the constraint between the ZIP Code and the street number. Thus the street number is participating in two conflicts. The same process is repeated for the other variables. The chart in Figure 7 shows the current values of each variable along with the conflicts that exist, designated by asterisks. The street name(X2) is in conflict with the street suffix, and the ZIP code, so it has two conflicts. The street suffix(X3) is in conflict with the street name and number, as well as the ZIP Code, so it has three conflicts. The state word(X5) agrees with all the other variables and the ZIP Code(X6) disagrees with the street name, number and suffix so it has three conflicts.

| var=value | X1=4555 | X2=Beltline | X3=Dr | X5=TX | X6=75001 | total conflicts |
|-------------|---------|-------------|-------|-------|----------|-----------------|
| X1=4555 | | | * | | * | 2 |
| X2=Beltline | | | * | | * | 2 |
| X3=Dr | * | * | | | * | 3 |
| X5=TX | | | | | | 0 |
| X6=75001 | * | * | * | | | 3 |

Figure 7: Instantiated values and their Conflicts

Therefore the suffix or the ZIP Code would be repaired first, since they are both tied with the highest number of conflicts. In this case the ZIP Code is chosen for repair. Since there is only one choice, 75001, it is again instantiated. The process of testing a solution and repairing a variable continues until a solution is found or the informed-backtracking returns without finding one. The latter happens in this example due to the errors on the address block.

Since no solution was found, failure analysis looks at the ZIP Code and the suffix because they had the most initial conflicts. Both of these are removed and the informed-backtracking is repeated, this time with only the street number, street name and state name. In this second trial, the top choices for these three variables combine to form a consistent solution. From this we are able to resolve the various uncertainties in this image. The city word(X4) can be inferred from the partial solution as "Addison" since there is only one city in Texas with a street name of Beltline and a street number of 4555. Thus, we are able to interpret the image even though some information was missing. The ZIP Code can also be inferred from the partial solution as "75244", so we have corrected the error made by the patron. Finally, we can constrain the choices for the suffix using the partial solution. In this case, there is only one choice, "Rd", so we have corrected our own recognition error. The final interpretation for this address is: 4555 Beltline Rd., Addison, Texas 75244-2416. The key point here is the additional four digits. These along with the ZIP Code constitute an interpretation for this address.

5 Experimental Results

This section compares the performance of our system to an optical character reader (OCR) of the type currently used in over 100 Post Offices in the United States. The input for this test consisted of 1200 address block images. These images were biased towards poor quality and difficulty in interpretation by a human. The poor image quality led to uncertainties of the first kind, incomplete images. In addition, many of the images contained what were called patron errors earlier. These are incomplete addresses or addresses that contain conflicting information, such as a ZIP

| Version | Encode | | Error | | Accept | | Error | | Reject | |
|---------|--------|-------|-------|-------|--------|-------|-------|-------|--------|-------|
| | N | % | N | % | N | % | N | % | N | % |
| OCR | 332 | 27.67 | 51 | 15.36 | 808 | 67.33 | 98 | 12.13 | 392 | 32.67 |
| CSP | 983 | 81.92 | 167 | 16.99 | 1171 | 97.58 | 198 | 16.91 | 29 | 2.42 |

Figure 8: Summary of Experimental Results

Code that does not match the city name. Also, while our word recognition results are encouraging, errors were made there as well.

The correct interpretation for each image was chosen by a human operator. These were used in grading the performance of the two systems. Out of the 1200 addresses, the human could find ZIP+4 Codes for 92%. The results of this test are summarized in Figure 8.

Of the 1200 images, the OCR was able to encode 28% to nine digits with an 15% error rate. Our system was able to encode 82% to nine digits with a 17% error rate. In postal address interpretation, the nine-digit encoding is the best level of interpretation, but not the only correct one. It is also correct to encode only to the five-digit ZIP Code. In such a case the specific ZIP+4 record cannot be determined, but there exists enough contextual knowledge on the mailpiece to develop an interpretation for the five-digit ZIP Code. The OCR was able to encode 67% of the images to at least five digits with a 12% error rate. Our system encoded 98% to at least five digits with a 17% error rate.

6 Conclusions

This paper presented a new approach to postal address interpretation. This approach is able to resolve various uncertainties found in address interpretation. Incomplete addresses, patron errors and recognition errors make address interpretation a challenging problem.

We proposed a new method for dealing with these uncertainties which combines the constraint satisfaction paradigm and failure analysis on the results of informed-backtracking. By formulating the problem in this way, we are able to generate partial solutions and infer the missing information, leave out words from the address and infer values for them later, and, also, overcome recognition errors.

The experimental results of this approach are encouraging and demonstrate an improvement over a previous system used in this domain. The technique is general enough to be applicable to other areas of image interpretation.

Acknowledgments

This research was supported by the United States Postal Service Office of Advanced Technology under task order 104230-88D-2576. The authors would like to acknowledge the support of the Office of Advanced Technology of the United States Postal Service. Carl O'Connor of USPS as well as Gilles Houle and Gerardo Garcia of Arthur D. Little, Inc. have provided guidance and encouragement.

References

- [1] P B Cullen, J J Hull, and S N Srihari. A constraint satisfaction approach to the resolution of uncertainty in image interpretation. In *Proc. of the Eighth Conference on Artificial Intelligence for Applications*, pages 127-133, March, 1992.
- [2] A K Mackworth. Constraint satisfaction. In *Shapiro S. C. (ed.), Encyclopedia of Artificial Intelligence*, volume 1, 2nd Edition, pages 205-211, 1990. John Wiley & Sons, New York.
- [3] J Bitner and E M Reingold. Backtrack programming techniques. *Communications of the ACM*, 18(11):651-656, 1975.
- [4] R M Stallman and G J Sussman. Forward reasoning and dependency-directed backtracking in a system for computer-aided circuit analysis. *Artificial Intelligence*, 9(2):135-196, 1977.
- [5] U Montanari. Networks of constraints: fundamental properties and applications to picture processing. *Information Sciences*, 7:95-132, 1974.
- [6] A K Mackworth. Consistency in networks of relations. *Artificial Intelligence*, 8:98-118, 1977.
- [7] R Mohr and T C Henderson. Arc and path consistency revisited. *Artificial Intelligence*, 28:225-233, 1986.
- [8] R Dechter and J Pearl. Network-based heuristics for constraint-satisfaction problems. *Artificial Intelligence*, 34:1-38, 1988.
- [9] S Minton, M D Johnston, A B Philips, and P Laird. Solving large-scale constraint satisfaction and scheduling problems using a heuristic repair method. In *Proceedings of the Eighth National Conference on Artificial Intelligence*, pages 17-24, July 1990. Boston, Massachusetts.
- [10] M Prussak and J J Hull. A multi-level pattern matching method for text image parsing. In *Proceedings of the Seventh IEEE Conference on Artificial Intelligence Applications*, February 1991. Miami, Florida.
- [11] T Ho, J J Hull, and S N Srihari. Combination of structural classifiers. In *Proceedings of IAPR Workshop on Syntactic and Structural Pattern Recognition*, pages 123-136, June 1990. Murray Hill, New Jersey.
- [12] P Filipski, J J Hull, and S N Srihari. Compression for fast read-only access of a large memory-resident structured database. In *Proc. of the Fifth Advanced Technology Conference*, November 1992.