International Journal on Document Analysis and Recognition manuscript No. (will be inserted by the editor)

Large Scale Address Recognition Systems

Truthing, Testing, Tools and other Evaluation Issues *

Srirangaraj Setlur¹, Alfred Lawson², Venugopal Govindaraju¹, Sargur Srihari¹

¹ Center of Excellence for Document Analysis and Recognition, 520 Lee Entrance Ste 202, Amherst NY 14228, USA;

 $e\text{-mail: } \{\texttt{setlur,govind,srihari}\} \\ \texttt{@cedar.buffalo.edu}$

² United States Postal Service, 8403 Lee Highway, Merrifield VA 22082, USA;

e-mail: alawson@email.usps.gov

Received: date / Revised version: date

Abstract. This paper describes the issues involved in the design of a system for evaluating improvements in the performance of a real-time address recognition system being used by the United States Postal Service for processing mail-piece images.

Evaluation of the performance of recognition systems is normally carried out by measuring the performance of the system on a representative sample of images. Designing a comprehensive and valid testing scenario is a complex task that requires careful attention.

Sampling live mail-stream to generate a deck of images representative of the general mail-stream for testing, truthing (generating reference data on a significant number of images), grading and evaluation, and designing tools to facilitate these functions are important topics that need to be addressed. This paper describes the

efforts of the United States Postal Service and CEDAR towards developing an infrastructure for sampling, truthing and testing of mail-stream images.

Key words: Truthing – Software tools – Evaluation systems – Recognition systems

1 Introduction

The United States Postal Service has invested significantly towards automated processing of mail-pieces to speed up sorting and reduce labor costs.

The letter mail automation program of the United States Postal Service (USPS) utilizes recognition software developed by a number of vendors. A brief introduction to the scale and nature of the Postal recognition systems program is essential to appreciate the issues in-

 $^{^\}star$ This work was supported by contracts from the USPS $Correspondence\ to:$ S. Setlur

volved in the development of a comprehensive system for evaluating the performance of these real-time large scale automation systems.

This paper describes the current status of postal automation of letter mail, provides an introduction to USPS addressing standards and databases and describes the working of a typical address recognition system. The paper then addresses issues involved in the evaluation of such address recognition systems, such as the testing process, generation of representative test decks, truthing of these test decks, tools developed to facilitate truthing and minimizing errors in truthing, and evaluation metrics.

2 Postal Automation - Background

Postal automation represents a fertile area for the application of image processing and pattern recognition techniques. Mail handling is a labor intensive process and the prohibitive cost of manual labor has made automatic sorting of mail an attractive economic proposition. The cost of processing per 1000 mail-pieces drops from \$47.78 for manual processing to \$27.46 for mechanized processing to \$5.30 for automated processing. The savings multiply rapidly given the volume of mail processed by the United States Postal Service (About 400 million pieces of letter mail per day). In addition to the cost factor, the knowledge level required for the sorting process is considerable and must be acquired.

The postal automation program focuses on three strategies: generating bar-coded mail, processing bar-coded mail in automated operations and adjusting the work force resulting in a reduction of work hours.

Bar-codes are generated by customers or by the Postal Service. In 1998, the United States Postal Service handled 41 % of the world's mail volume which translates to about 630 million pieces of mail every day. The next largest is Japan at 6 %. Letter mail accounted for nearly 60 % of this volume i.e., a staggering 135 billion letter mail-pieces in the year 1998.[1] Of these, customers generated nearly 70 billion bar-coded letters or 52 % of the letter mail volume. The pre-barcoded mail is easily sorted using fast bar-code sorting systems. The remaining 48 % of the letter mail is processed through various recognition systems.

The letter mail is first processed through the MLOCRs or the MultiLine Optical Character Readers which attempt to recognize addresses on mail-pieces (machine-printed) in real-time. The MLOCRs are augmented by a second recognition system called the Recognition Co-processor. Of the bar-codes generated by the Postal Service, the MLOCRs accounted for about 19 billion mail-pieces. The mail-pieces that cannot be read by the MLOCRs are captured and processed through the Remote Bar Coding System or RBCS.

The RBCS processing (Fig 1) is divided into two general categories: RCR (Remote Computer Reader) processing, which involves automatic reading and interpretation of the address information, and REC (Remote Encoding Center) keying, which requires human interpretation.

Handwritten as well as machine printed mail-pieces are processed through the RCR. CEDAR's Handwritten Address Interpretation Software (HWAI) is used for processing handwritten mail in the RCR.[2]

When all electronic means of resolving address information have been exhausted, the mail-piece image is sent to REC sites where operators use video display terminals and keyboards to manually enter the address information. This keyed information is fed back to the RBCS to allow bar-codes to be applied on to the mail-piece. The RBCS accounts for the bar-coding of an additional 35 billion letter mail-pieces a year. Providing partial RCR results to the REC sites with the image can allow the operator to process the address with minimal number of keystrokes.

The bar-code encapsulates the address that is present on the mail-piece. The objective of reading the address on the mail-piece is to generate a DPC or a delivery point code which can be applied as a bar-code on the mail-piece and bar code sorters can then use this bar-code to sort the mail according to the carrier walk sequence which is the order in which the mail-carriers do their rounds.

Automation is revolutionizing mail processing operations as the Postal Service continues to invest in technology infrastructure to reduce costs, improve service and increase operating efficiencies. Upgrades to the Remote Computer Readers (RCRs) have resulted in significant savings to the Postal Service, including a reduction of 12 million work hours annually and cumulative labor savings of more than \$208 million. [3]

3 USPS Databases & Addressing Standards

This section provides an insight into the types of postal data available, the wide variety of addresses possible and the kind of rules required to resolve the address on the mail-piece.

The Postal Service maintains different databases of postal addresses, two of which are the ZIP+4 file and the DPF (Delivery Point File). The ZIP+4 database has records representing a range of primary numbers (street numbers or PO box numbers as the case may be) whereas the DPF or the Delivery Point File has records representing every individual delivery point in the United States. Supplemental files provide information about ZIPtype, City-State-ZIP correspondence and ZIP translation.

The Postal Service and the bulk mailing industry have jointly developed postal addressing standards that are geared towards enhancing the processing and delivery of mail and reducing undeliverable mail thereby providing mutual cost reduction opportunities through improved efficiency. The standards are fairly well adhered to especially in the machine printed mail-stream. Handwritten mail finds greater deviations from the standards.

The last line or the last two lines of a mail-piece with a US domestic delivery address contain City, State and ZIP code information. The line above contains delivery information such as a street number, street name with a street suffix and secondary information such as apartments, suites, etc. The mail-piece could also have a firm or organization name or a personal name or PO box information.

Advanced Facer Canceler Image RCR Multi-Line OCR RCR Remote Encoding

Fig. 1. USPS Mail processing flow diagram

The Postal Service has developed an encoding scheme that allows for the sorting of mail to the block level, building level or delivery point level using different record types.

3.1 Types of ZIP Codes

There are four different types of ZIP codes that affect how the address is encoded.

1. Unique ZIP: Unique ZIPs are assigned to organizations or firms that receive large volumes of mail. The

distribution of mail to mail-stops within the organization is handled by the organization itself. These addresses are encoded to just 5 digits unless the patron writes the four digit addon on the mail-piece in which case the correct encode would be to return the entire nine-digit string.

2. Automated/Non-automated ZIP: The extent of automation available for sorting at the destination Post Office determines whether a mail-piece has to be en-

coded to just the 5 digit ZIP or whether the address has to be further interpreted to determine the addon and delivery point code. If the destination ZIP happens to be a non-automated ZIP, encoding can stop at the 5 digit ZIP. If the destination happens to be an automated ZIP, encoding should be to 9 or 11 digits depending on whether the automation equipment at the destination Post Office can sort the mail-piece to 9 or 11 digits. This information is also available in the postal databases.

- 3. Translated ZIP: As the number of delivery points within a ZIP increases, the Postal Service periodically undertakes ZIP re-organization efforts to distribute delivery points evenly amongst available ZIP codes. ZIPs could be either split or moved or new ZIPs created as required by the Postal Service. If the patron applied delivery ZIP on the mail-piece is a translated ZIP, a ZIP translation table has to be looked up to determine the list of ZIPs to which the ZIP in question is translated to and the mail-piece has to be encoded to the ZIP that contains the destination address on the mail-piece.
- APO/FPO ZIP: These are special ZIP codes assigned to military addresses. These ZIPs are encoded similar to mail-pieces with a Unique ZIP.

In the absence of a ZIP code on the mail-piece, City and State information on the mail-piece is to be used to locate the appropriate ZIP code. Once the ZIP code and ZIP type are determined, the delivery lines of the address are used to determine the appropriate 4 digit addon and the 2 digit delivery point code.

3.2 Types of addresses

The addressing scheme instituted by the Postal Service includes different types of addresses viz., (i) street, (ii) PO Box, (iii) rural route, and (iv) special address types such as General delivery and Postmaster.

- 1. Street Addresses: Data corresponding to street addresses are stored in the postal databases in a hierarchical scheme that consists of four different types of records. Street records are the coarsest followed by high-rise default records, high-rise specific records and firm records which are the finest. These types represent levels of sortation and while the addons corresponding to the coarser records are valid for a given mail-piece, the cost incurred by the Postal Service to sort it to the finest addon is the least and hence vendors of recognition systems are expected to return the finest level of sort possible for a given address.
- PO Box Addresses: The data corresponding to PO Box addresses are stored as PO Box type records.
- 3. Rural Route Addresses: The data corresponding to rural route addresses are stored as rural route records.
- 4. Other Address Types: General Delivery and Postmaster are two specialized types of addresses which are stored as General Delivery type and Postmaster type records.

3.3 Encoding a mail-piece

Using the databases and the supplemental files and a complex set of encoding rules the destination address is encoded.

We will use some sample addresses to illustrate the different addressing concepts and the process of deriving a correct encode.

In the address:

CEDAR

UB Commons

520 Lee Entrance Ste 202

Amherst NY 14228

the ZIP 14228 is an automated ZIP and is not UNIQUE and hence the desired level of encoding is 11 digits.

The records in the DPF postal database corresponding to the ZIP 14228 and street number 520 are shown in Table 1. It can be seen from the table that the addon corresponding to the address 520 Lee Entrance Ste 202 is 2583. A set of rules helps determine the DPC corresponding to this high-rise specific addon and for Ste 202 the corresponding DPC code is 52. The finest level of encode for the given address is 14228-2583-52 and this 11 digit code determines a unique delivery point in the United States. The coarser encodes, viz. the street encode of 14228-2500-20 and the high-rise default encode of 14228-2567-99 are also valid but do not represent the finest depth encode for the mail-piece. The finest depth encode allows sortation to the finest depth possible which results in the maximum savings to the Postal Service in terms of processing costs. The coarser depths do not allow sortation to the finest extent possible and hence require some manual sorting which adds to processing costs.

4 Address Recognition Systems

The task of reading addresses on mail-pieces involves solving challenging problems in the areas of image analysis, digit, word and phrase recognition, directory design and look-up and developing intelligent interpretation techniques.

After locating the destination address, the address recognition system splits the address block into its component lines and words, uses knowledge about addressing syntax and standards (see Section 3) to parse and label address elements such as ZIP code, state abbreviation, city name, street number and street name. Recognizers are then used to recognize these elements and access postal databases based on recognition results to determine the delivery point code. Figure 2 describes the flow of the CEDAR HWAI system.[2]

The automated address interpretation process requires an in-depth understanding of the structure of the postal data, the levels of sortation possible and the rules to be used to arrive at the best encoding for the address on the mail-piece. The fact that the patron (postal consumer) might have written conflicting data on the mail-piece increases the complexity of the rules required to determine the best encode for the mail-piece.

Table 1. Sample records from the DPF corresponding to the ZIP 14228 and Street number 520

| Primary No | Rec Type | St Name | St Suffix | Sec Des | Sec No | Sec/Firm Name | Addon | Carrier Route |
|------------|----------|------------------|-----------|---------|--------|---------------|-------|---------------|
| 520 | 10 | N ELLICOTT CREEK | RD | | | | 2323 | C059 |
| 520 | 10 | LEE ENTRANCE | | | | | 2500 | C050 |
| 520 | 10 | LEE ENTRANCE | | STE | 200 | UB COMMONS | 2500 | C050 |
| 520 | 10 | LEE ENTRANCE | | STE | 211 | UB COMMONS | 2500 | C050 |
| 520 | 20 | LEE ENTRANCE | | | | | 2567 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 101 | UB COMMONS | 2577 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 103 | UB COMMONS | 2577 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 109 | UB COMMONS | 2577 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 112 | UB COMMONS | 2577 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 204 | UB COMMONS | 2577 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 209 | UB COMMONS | 2578 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 215 | UB COMMONS | 2578 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 306 | UB COMMONS | 2578 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 215A | UB COMMONS | 2578 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 301 | UB COMMONS | 2579 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 302 | UB COMMONS | 2579 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 303 | UB COMMONS | 2579 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 105 | UB COMMONS | 2580 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 201 | UB COMMONS | 2580 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 203 | UB COMMONS | 2580 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 217 | UB COMMONS | 2580 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 304 | UB COMMONS | 2580 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 100 | UB COMMONS | 2583 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 110 | UB COMMONS | 2583 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 116 | UB COMMONS | 2583 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 202 | UB COMMONS | 2583 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 210 | UB COMMONS | 2583 | C050 |
| 520 | 20 | LEE ENTRANCE | | STE | 106 | UB COMMONS | 2584 | C050 |

Address Interpretation System

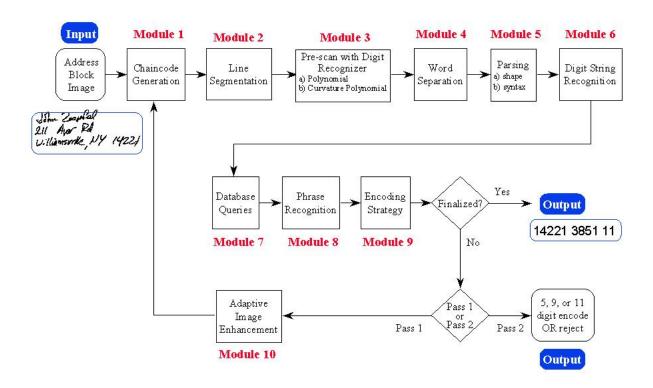


Fig. 2. Address Interpretation Process

5 Evaluation - Guidelines

The evaluation model for an address interpretation system has to also take into account some problems that are unique to the postal mail-stream scenario.

- 1. Postal addresses are not *static* for the following reasons:
 - New addresses are added to the database.
 - ZIP codes get translated into new ZIP codes and addresses move between ZIP codes.
 - PO Box numbers get cancelled or re-assigned.
 - New firms are added to the database.

The postal databases are a key component of the system encoding as well as the truthing processes. If the truthing is done using a database from a time different from the database that is used by the system being tested, the testing process would be flawed since the truth for the same mail-piece given one temporal instantiation of the database might not be the truth given a different temporal instantiation of the database. Hence, it is important to ensure that the database used by the system during the test is the same temporally as the database used for truthing.

- 2. Patrons sometimes write ambiguous, incomplete or erroneous addresses on mail-pieces. Comprehensive encoding rules are required to ensure that the resultant encode for a given address can be uniquely determined and the test deck has to be truthed in accordance with these rules so that the system can incorporate the same rules in its encoding scheme and the evaluation will not suffer as a result of the inherent ambiguities in patron addressing.
- Handwritten addresses present another source of confusion. Subjective interpretation of ambiguous digits and characters present problems in determining the truth unambiguously.

5.1 Letter Mail Scoring and Evaluation Rules

These guidelines are used for scoring and evaluating mail pieces processed by any letter mail system. Each rule is designed to achieve the best potential possible for a mail piece. The finest depth of code is determined by using the available mail piece information to get the best and most cost effective ZIP Code resolution. This is performed by using the Expert Zip+4 Finder (EZF) or other directory query programs.

The following rules are used to score the finest potential of a mail piece and to judge the resolutions achieved from a specific system.[4]

Unique ZIP Codes (consistent with the city and state)
 on mail pieces override all other address information.
 The system is expected to return the Unique ZIP
 Code on the mail piece. A 5-digit Unique ZIP Code

may have a +4 range assigned to it. Coding is required to the 9-digit level when indicated on the mail piece, regardless of whether or not the +4-digit code exists in the USPS database. Such a piece would be in error if not resolved to the 9-digit level.

EXAMPLE:

ENGINEERING

8403 LEE HWY

MERRIFIELD VA 22082-8133

The 5-digit ZIP Code of 22082 is a Unique ZIP Code representing the USPS Engineering building. However, the full potential of the mail piece is 22082-8133. APO (Army and Air Force Post Offices) and FPO (Navy Fleet Post Offices) ZIP Codes are similar to unique ZIP Codes, in the fact that mail is processed by the organization and is not delivered to the individual addressee by the USPS. Coding is required to the nine digit or eleven digit level when indicated on the mail piece or by the address.

BRM (Business Reply Mail) mail piece is considered direct. The 5- or 9- digit on the BRM mail piece must be duplicated. BRM should never be resolved to any line above the City/State line.

2. Addresses with strike through or the standard return to sender indicator shall be resolved as follows:

A mail piece with a strike through the destination address, but with no forwarding address or return to sender indicator shall be resolved to the destination address. A mail piece with or without a strike through the destination address and with a forwarding address but no return to sender indicator shall be resolved to the forwarding address.

A mail piece with a return to sender indicator shall be resolved to the return address.

A mail piece with a CFSII yellow label with or without Form 3547 printed on it shall be resolved to the address written on the label. However, when a copy of the Form is attached to the USPS letter going to the sender as an address correction service request, the address shall be resolved to the return address.

- Any correction to the customer ZIP should be validated by the city/state and street address.
- 4. In a non-automated zone (consistent with city/state),5-digit is the highest depth of potential.
- 5. Mail with both a street address and a PO Box number is delivered to the address immediately above the city/state line (or to the PO Box if both the street address and PO Box are on the same line). If a ZIP+4 or 5-digit ZIP is used, it must correspond to the address element immediately above the city/state line (or with the PO Box number if both the street address and PO Box are on the same line). These restrictions also apply to return addresses.
- 6. Addressing which indicates Caller Service (CS) or Caller Number (CN), Post Office Box, POB, PO, Drawer, Caller, Firm Caller, Bin, Lock Box, Lockbox, or Drawer CS is considered the same as a PO Box.

If the PO Box number does not exist within the ZIP Code of the mail piece address, the system should correct the customer error.

EXAMPLE:

JOHN DOE

PO BOX 990

AKRON OH 44311

The above should be corrected to 44309-0990-90, the only zone in the city where PO Box 990 is located.

A customer error in the ZIP Code cannot be corrected if the PO Box number exists within several zones of the addressed city. The potential for mail pieces that indicate more than one result is the customer ZIP, if the ZIP correlates to the City/State. If it does not correlate, the system should code to the city default code.

EXAMPLE:

April May June

PO BOX B

YOUNGSTOWN OH 44501

PO Box B is located in zones 44507 and 44515. Neither of these zones can be selected when there is more than one valid answer. Therefore, this mail piece should be resolved to the customer of ZIP 44501, which correlates to Youngstown, OH.

For mail pieces with the word BOX written in the address:

If the BOX number is a valid PO Box range for the ZIP Code on the mail piece in that city, code to the PO box, only if there is no street address.

If a dual address is on the mail piece, i.e., a BOX with a street address on the same line, code to the street address. If the zone is a PO Box zone only, code to the PO Box.

If the address indicates a rural route or other box type locations, do not code to the PO box. Code to the appropriate rural route / box ZIP+4.

7. The potential for a mail piece that contains rural route (RR) or highway contract route (HCR) with street addressing is Street DPS. A rural route default is considered a lower potential.

EXAMPLES:

Terry Bull

43901 GIRDLAND COURT

ASHBURN VA 22011

Terry Bull

43901 GIRDLAND CT RURAL RT 5

ASHBURN VA 22011

Terry Bull

43901 Girland CT

R.R. 5

Ashburn VA 22011

All three of these addresses can be resolved to an 11-digit ZIP of 22011-3922-01, which indicates a street address. A rural route default of 22011-9805-99 is a lower potential.

8. Addressing indicating Postmaster or General Delivery should be resolved to the 11-digit ZIP repre-

senting the General Delivery or Postmaster delivery point. The following example should achieve a 11-digit ZIP of 22180-9998-99 to indicate Postmaster:

POSTMASTER

Vienna VA 22180

Absence of a suffix may limit the possible depth of code:

EXAMPLE:

Mr. & Mrs. XYZ

7313 MASONVILLE

ANNANDALE VA 22003

7313 Masonville Court in Annandale has a ZIP Code of 22003-1629-13. 7313 Masonville Drive has a ZIP Code of 22003-1637-13.

Neither should be selected based upon the address information as it appears in the example. The potential is 5 digit.

This example, with a missing suffix, should result in a default code for Annandale VA, 22003, unless the system can identify name or unique information on the mail piece to break the tie. If the $\rm ZIP + 4$ was present, the directional/suffix could be assumed.

10. Prefixes and directional may also be important elements in mail piece resolution:

EXAMPLE:

Tim Burr

322 East St.

Vienna VA 22180

This address does not indicate direction. However, 322 East St. SE would achieve a street with a DPS

code of 22180-4821-22, and 322 East St. NE would achieve a street with DPS code of 22180-3618-22. Since more than one answer is valid with the given address information, only the 5 digit can be selected unless additional information is used to break the tie. If ZIP + 4 was present, the directional could be assumed.

11. One line or two line addressing:

EXAMPLE:

OHIO EDISON 45 S MAIN ST AKRON OH

A unique 5-digit ZIP, 44399, should be achieved based upon the firm name in this example.

EXAMPLE:

Chanda Lear

45 S MAIN ST AKRON OH

This address would achieve a Street with DPS ZIP of 44308-1803-45 since the firm name was not included.

12. Default add-ons:

EXAMPLE:

Pete Moss

RURAL ROUTE 3 180 NEWTON MILTON RD
NEWTON FALLS OH 44444

The rural route default add-on code of 44444-9803-99, indicating the entire rural route rather than the particular sector/segment for the rural route box range, is considered a lower potential. The correct ZIP, 44444-9303-80, indicating Street with DPS should be achieved.

EXAMPLE:

Davey Crockett

505 ROOSEVELT BLVD B-721

FALLS CHURCH VA 22044

The correct ZIP of 22044-3122-XX, indicating a range of apartments (the secondary address unit) within the building, should be achieved. The ZIP indicating the entire building, 22044-3141-99 (a Building with No Secondary Code) or the Street with DPS, 22044-3114-05, are considered to be lower potential.

A mail piece with building secondary potential (building non-default) that is only coded to the Street DPS or to the building default is a lower potential.

- 14. At any multi-coded address, the add-on on the mail piece will be the finest depth of code if there is correlation between the add-on and the street primary information on the mail piece; as long as there is no correlation between any secondary information on the mail piece and any other directory entries at that address.
- 15. There must be a street correlation to match a firm.

 The only exception is an exact firm name match on a mail piece that contains a 5-digit ZIP Code that correlates with City/State information, contains no street information, and is the only firm listing for that firm in the zone. The use of Postal Abbreviations is acceptable. NOTE: If there is a ZIP+4 on the mail piece, an exact firm name match is not necessary.

13. Building/High-rise add-ons:

16. A bar code will be treated in the same manner as the customer-applied ZIP Code in all rules dealing with performance evaluation.

17. A partial result will be considered correct if the fivedigit ZIP code returned is consistent with the customer applied city/state and street information.

EXAMPLE:

7301 Westwind

Bowie MD 20715

Result should be 20715 as there is a Westwind Ct & Westwind Dr

7301 Westwind

Bowie MD

Result should be 20715 which is the default for Bowie MD

7301 Westwind

Bowie MD 20718

Result should be 20718 which is the customer ZIP, can not change because there is no confident result

7301 Westwind Dr

Bowie MD 20718

Result should be 20715-1735-01 because Westwind Dr is in 20715. The partial result shall be 20715.

18. Delivery Point Sequence (DPS) Rules

(a) General Rule:

Address: 234 Main Street, or PO Box 234, or RR $\,$

1 Box 234, or HC 1 Box 234

DPS: 34

Print the ZIP code characters in the Delivery Point Sequence representing the last two digits of the primary street number, PO Box number, or rural box number in the address.

(b) Single-Digit Numbers:

Address: 8 Main Street, or PO Box 8, or RR 1 Box 8, or HC 1 Box 8

DPS: 08

Add a leading zero before the single digit street number, print the code characters in the Delivery Point Sequence. In this example, it would be 08.

(c) Fractional Numbers:

Address: 234 Main Street, or PO Box 234, or RR

1 Box 234

DPS: 34

Ignore fractions in the address; print the code characters in the DPS to represent the two digits preceding the fraction. In this example, it would be 34.

(d) Trailing Alphas:

Address: 235A Main Street, or PO Box 235A, or

RR 1 Box 235A

DPS: 35

Ignore trailing alpha characters in the address; print the code characters in the DPS to represent the two digits preceding the alpha(s). If a single digit is to the left of the alpha(s), add a leading zero. In this example it would be 35.

(e) Space and Alphas:

Address: 12 A Main Street, or 12 AA Main Street,

or PO Box 12 A

14

DPS: 12

Ignore space if followed by a trailing alpha character(s) in the address; ignore the alpha character(s) as well (see Rule 4.4 above); print the code characters in the DPS representing the two digits preceding the space(s) and alpha(s). If a single digit is to the left of the space(s) and alpha(s), add a leading zero.

(f) Embedded Alphas

Address: 123S4 Main Street

DPS: 04

Print the code characters in the DPS representing the last two digits to the right of the alpha(s). Replace an embedded alpha character with 0 if the alpha is to the immediate left of the last digit in the primary street number of the address; print the code characters representing zero (0) and the last number of the primary street number in the DPS. In this example it would be 04.

(g) Other Embedded Symbols:

Address: 1.23 Main Street, or 1-23 Main Street, or 1*23 Main, or PO Box 1-23

DPS: 23

Print the code characters in the DPS representing the last two digits to the right of all other symbols other than slashes, such as periods (.) and hyphens (-), appearing in the primary street

number. If there is a single digit to the right, add

a leading zero (0).

(h) Embedded Spaces:

Address: 1 23 Main Street, or PO Box 1 23

DPS: 23

Treat spaces that are embedded in the primary street number like Other Symbols (see Rule 4.7 above); print the code characters in the DPS representing the last two digits to the right of the space. If a single numeric digit is to the right, add a leading zero (0).

(i) Numeric Street Names:

Address: 12 8th Ave

DPS: 12

Ignore numeric street names (in this example "8th").

Print the code characters in the DPS representing the last two digits of the actual primary street number (see Rule 4.1 above).

(j) Non-Anomaly:

123-3C Main Street

This is not a DPS address format anomaly. Apartment 3C at 123 Main Street could receive a ZIP+4 building record bar code (if the address is assigned a building record), or the DPS code, in which case 23 would be coded to produce the delivery point bar code.

- (k) For Postmaster, General Delivery, and Rural Route Default the delivery point code will be 99. For Firm records the delivery point code will be 99 except where the firm has a suite number, the delivery point shall be to the suite using High-rise Delivery Point Code algorithms.
- Deriving Delivery Point Codes For Highrise Secondary
 Addresses

The USPS has identified the potential for significant savings and service improvements through the implementation of Delivery Point Sequencing of high-rise deliveries.

The following rules were found to be effective for creating unique deliver sequences form the current coding of high-rise secondary addresses in the AMS database.

The rules apply only to input addresses matching to high-rise records in the ZIP+4 or DSF file (Record Type H).

(a) Numeric Simple Rule:

This rule applies when the secondary address is numeric and the numeric value of the combination of hundreds and thousands digits equals zero. Simply stated, the last two digits of the secondary number become the delivery point code.

| APT/STE | DPC | APT/STE | DPC |
|---------|-----|---------|-----|
| | | | |
| 1 | 01 | 10001 | 01 |
| 2 | 02 | 10002 | 02 |
| 3 | 03 | 10003 | 03 |
| 98 | 98 | 10098 | 98 |
| 99 | 99 | 10099 | 99 |
| | | | |

(b) Alphabetic Rule:

This rule is used when the secondary address contains only alphabetic characters, and the rightmost character is within the range form A through Z. The rule simply defines each character of the alphabet as a delivery point code from 73 through 98 (i.e., A=73, B=74, C=75, , Z=98) using the rightmost alphabetic character only.

| APT/STE | DPC | APT/STE | DPC | APT/STE | DPC |
|---------|-----|---------|-----|---------|-----|
| | | | | | |
| A | 73 | RF | 78 | ABA | 73 |
| В | 74 | RR | 90 | A-C | 75 |
| С | 75 | SE | 77 | A-Q | 89 |
| AA | 73 | SW | 95 | AAAB | 74 |
| AB | 74 | W | 95 | | |
| | | | | | |

(c) Alphanumeric Rule - Trailing Alpha:

This rule applies to alphanumeric secondary addresses with trailing alphabetic character(s). It forms the delivery point code from the secondary address by converting the trailing alphabetic character to a numeric value using the following conversion of the 26 alphabetic characters, and then adding 10 times the least significant numeric character.

A = 1 F = 6 K = 21 P = 26 U = 41
B = 2 G = 7 L = 22 Q = 27 V = 42
C = 3 H = 8 M = 23 R = 28 W = 43
D = 4 I = 9 N = 24 S = 29 X = 44
E = 5 J = 10 O = 25 T = 30 Y = 45
Z = 46

Formula When: Trailing alpha conversion = X, rightmost numeric = Y

Delivery Point Code: DPC = X + (10 * Y)

NOTE: Derive the delivery point code by determining the modulus (remainder) of

$$((10 * Y) + X) / 100$$

Formula When: Alpha conversion value = X, trailing numeric value = Y (X*10) + Y For unit umber N-101 = (24*10) + 1 = 241 MOD (241/100) = 41

| | | | APT/ST E | (10 * Y) + X MOD OF 100 DPC | | |
|----------|---------------------------|-----|----------|-----------------------------|----|--|
| APT/ST E | (10 * Y) + X MOD OF 100 I |)PC | A1 | (1*10)+1 = 11 | 11 | |
| 1 A | (10*1)+1 = 11 | 11 | A-33 | (1*10)+3 = 13 | 13 | |
| 4-4M | (10*4)+23 = 63 | 63 | В3 | (2*10)+3 = 23 | 23 | |
| 12B | (10*2)+2 = 22 | 22 | 3A-174 | (1*10)+4 = 14 | 14 | |
| A78Z | (10*8)+46 = 126 | 26 | 4G5 | (7*10)+5 = 75 | 75 | |
| A4A | (10*4)+1 = 41 | 41 | 3V-175 | (42*10)+5 = 425 | 25 | |
| 44M3A | (10*3)+1 = 31 | 31 | Q37 | (27*10)+7 = 277 | 77 | |
| 99Q | (10*9)+27 = 117 | 17 | 440222 | (3*10)+2 = 32 | 32 | |
| 78-AAB | (10*8)+2 = 82 | 82 | | | | |
| | | | | | | |

(d) Alphanumeric Rule - Trailing Numeric:

This rule applies to alphanumeric secondary addresses with trailing numeric character(s). It forms the delivery point code from the secondary address by converting the leading alphabetic character using the alphabetic conversion table above, multiplying by 10 and adding the trailing numeric digit.

NOTE: Derive the delivery point code by determining the modulus (remainder) of $((X\ *\ 10)\ +\ Y)\ /\ 100$

(e) Numeric Computed Rule:

This rule applies to numeric secondary addresses when the numeric value of the combination of hundreds and thousands digits is greater than zero. It computes the delivery point code from the secondary address by right justifying, zero filling, computing the modulus-4 remainder of the combination of thousands and hundreds digits, multiplying the remainder by 25, and adding the modulus-25 remainder of the combination of the tens and units digits.

The algorithm is as follows:

- * Determine the modulus (remainder) when dividing the thousands-hundreds (x) value by 4.
- * Multiply the remainder by 25.
- * Determine the modulus when dividing the tensunits value by 25.
- * Add the two results.

The formula when the thousands-hundreds value = X and the tens-units value = Y is:

$$DPC = 25*(X MOD 4) + (Y MOD 25)$$

For apartment 1105, the DPC = 25 (11 MOD 4) + (05 MOD 25) = 25 (3) + 5 = 80

For apartment 1205, the DPC = 25 (12 MOD 4) + (05 MOD 25) = 25 (0) + 5 = 05

EXAMPLES:

| PRIMARY ADDRESS | SEC \# | DPC |
|-----------------|--------|-----|
| | | |
| 123 MAIN ST | 901 | 26 |
| 123 MAIN ST | 919 | 44 |
| 123 MAIN ST | 1001 | 51 |
| 123 MAIN ST | 1019 | 69 |
| 123 MAIN ST | 1201 | 01 |
| 123 MAIN ST | 1219 | 19 |
| | | |

This provides unique delivery point bar codes for all secondary addresses within a 4 X 25 range of secondary addresses. Stated simply, if an H (highrise) ZIP+4 Code covers 400 ranges (four floors) and the low-to-high range of secondary addresses

on each floor is not greater than 25, then no reassignment of ZIP+4 Codes will be necessary. For example, if the high-rise ZIP+4 Code includes apartments 101 through 125, 201 through 225, 301 through 325, and 401 through 425, all apartments will have a unique delivery point code. If the ZIP+4 Code includes more than 400 ranges, or if the low-to-high within a hundred range is greater than 25, additional ZIP+4 Codes may have to be assigned to obtain unique delivery point codes.

(f) Address Matching a Record With Blank Secondary Ranges

If the software matches to a record with a secondary descriptor but no secondary ranges, the software must return the delivery point code 99.

Some of the descriptors that may exist without a secondary value are:

BSMT, LBBY, REAR, LOWR, FRNT, OFC, SIDE, PH, UPPR

- (g) Address Matching a High-rise Default Record

 If software was unable to match the secondary
 address information to a high-rise record at the
 primary address, but matches instead to the H
 default record at the primary address, the delivery point code must be 99.
- (h) Miscellaneous:

For apartment numbers with fractions, the fractional component of the secondary address must be ignored when computing the delivery point bar code. If the input secondary address contains embedded special characters such as spaces, dashes, or periods these characters must be ignored when calculating the delivery point code.

20. Definitions:

- (a) Potential: Potential is the finest depth of sort to which an address can be resolved during truthing using the information available on the mail piece and in the postal databases. The following describes the categories of potential for any given mail piece:
 - Unique (U) Unique Zip Codes that are assigned to large volume mailers for their use in receiving mail. Included in this category are Business Reply Mail (BRM) and APO/FPO ZIP Codes.
 - Firm (F) An address with a firm name match
 - Post Office Box (P) A locked box, located in the post office lobby or other authorized place, that customers rent for delivery of their mail.
 - Delivery Point Sequence (DP) A street address record or a rural route with city addressing in an automated zone, which is in the directory. Street DPS codes are ZIP+4 codes with the addition of two (2) digits representing the last two (2) digits of the primary street number.
 - Building with Secondary (HN) A building address with secondary information that can

- be coded to the floor, suite, or apartment range.
- Building Default (H) A building within a sector segment in an automated zone that has been assigned a separate +4 code indicating the structure.
- 9-Digit (S) A mail piece coded to ZIP + 4 except for a mail piece falling under U category as defined above.
- 5-Digit (FD) A non-automated zone. In an automated zone addressing with insufficient information on the image for determining either the add-on or a higher potential.
- City/ /State Default (DF) A zone could not be determined due to insufficient or conflicting address information. The officially assigned default code for the City/State on the image is required.
- Sectional Center Facility SCF (DT) A postal facility that serves as the distribution and processing center for post offices in a designated geographic area, which is defined by the first three digits of the ZIP Code of those offices. This facility may serve more than one 3-digit ZIP Code range.
- Foreign 9-digit Plus (F9P) A foreign mail piece resolved beyond the destination country level, to a destination province, state, city or post-code.

- Foreign 9-digit (F9) A foreign mail piece resolved to the destination country name or alias country name.
- Foreign 5 (F5) A foreign mail piece destinating to a country having 5-digit potential
- Not Resolvable (N) Images with addresses,
 which cannot be determined.
- (b) Finalized Image: A finalized image is one that is not transmitted for further resolution. The following images are finalized:
 - * Images resolved to a 9-digit ZIP Code or better level.
 - * Images resolved to a 5-digit ZIP Code in a nonautomated zone and foreign mail piece having 5digit potential.
 - * Images falling under Unique ZIP Code category.
- (c) Finest Sort Rate: The finest sort rate is the number of images correctly resolved to their potential divided by the number of images in the sample.
- (d) Finalization Rate: The finalization rate is the number of images finalized, including errors, divided by the number of images in the sample.
- (e) Finalized Error Rate: Finalized errors made in the respective error category divided by the number of images finalized.
- (f) Partials: A partial is a resolved image but is not finalized.
- (g) Partial Error: A partial error is defined as the error in a partial that is not consistent with the

customer-applied city/state and street information.

21. Classification of errors:

The errors for the letter mail recognition system are classified as under:

(a) Error type 1:

An error in the first three digits of the ZIP Code. Addresses with strike through or the standard return to sender indicator shall be resolved as stated earlier. Errors made in this category shall be considered as an error type 1.

Addresses resolved to the return address when there is no strike through or the standard return to sender indicator.

An error in any of the first five digits of a foreign mail piece.

- (b) Error type 2: is an error in the fourth and fifth digits of the ZIP Code except in case of a foreign mail piece.
- (c) Error type 3: is an error in the four-digit add-on to the ZIP Code (digits six through nine).
- (d) Error type 4: is an error in the last two digits (digits ten and eleven) delivery point code.
- (e) If there are multiple errors in the ZIP Code, ZIP+4, or DPS resolution, the error occurring first (from left to right) is charged.

6 Testing Scenario

Address recognition systems have been installed at over 250 mail processing centers across the country. Turnaround

mail or mail destined to the sender's geographical area forms a substantial percentage of mail at each of these processing centers. Especially with handwritten mail, addressing characteristics could differ widely based on geographical regions and ethnic make-up of the populace.

Given the wide variety of handwritten mail, it would be difficult to model the variations and create synthetic data that could be used to effectively evaluate these recognition systems. Therefore, there is no alternative but to generate a representative test deck using real data.

It is vital for the test deck to reflect the national mail-stream in order to ensure that the performance of the address interpretation systems on the test deck can be expected to be reflected over most of the post offices where the software is deployed.

Figure 3 describes the testing and evaluation system. Images from the test deck are transferred to the truthing stations. The truthers create an ASCII representation of the address on the mail-piece and also record some characteristics of the mail-piece and the address block. The ASCII truth is then processed through an ASCII address matching engine that uses the rules of encoding determined by the Postal Service to interpret the ASCII address and generates the appropriate encodes at the various levels of sortation.

Patron errors on mail-pieces, incomplete, incorrect, and ambiguous addresses can cause the address matching engine to fail. If the address matching engine is unable to encode the mail-piece beyond 5 digits (for non-Unique ZIPs), the image is forwarded to a second truthing station where the truther uses additional queries to vari-

ous databases to try to resolve the mail-piece. The truth data is stored in a database for use during the evaluation phase.

The images from the test deck are also processed through the address interpretation engines provided by the vendors. The results from these runs are stored in a results database.

Grading utilities which are part of the System Analysis Toolkit are then used to match the vendor results against the truth and a detailed performance analysis report is generated.

7 Test Deck Generation

A fundamental goal of the test deck generation process is to ensure that the images collected for the test deck are representative of the image population observed in actual field operation.

The evaluation process measures the performance of an address interpretation system on this test deck of images.

Two key metrics used in the performance evaluation of an address interpretation system are the encode rate of the system and its error rate.

To evaluate the error rate, it is necessary to have a deck which is "truthed" so that the accuracy of the encode can be checked. However, truthing is an extremely expensive process. So, a further down-sampling of the test deck is done to generate a smaller deck that is truthed.

Keeping in mind the twin objectives of measuring the encode rate and the error rate, a test deck of a million handwritten images was generated, consisting of contem-

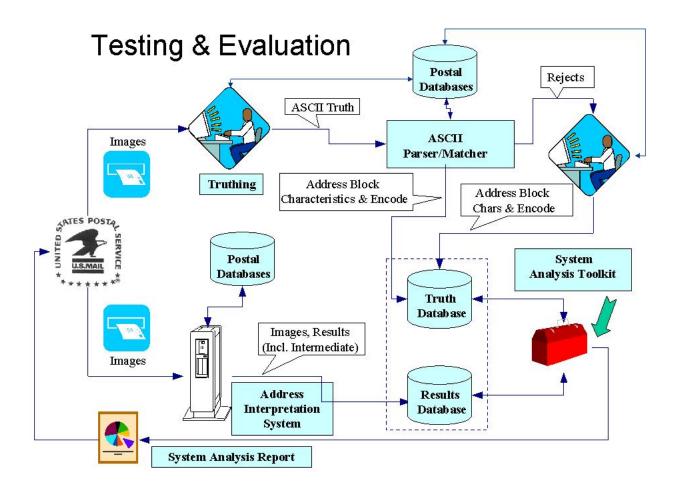


Fig. 3. Testing and Evaluation System

porary images from a minimum of 120 of the 251 Remote Computer Reader sites distributed all over the United States. Since the mail-stream consists of a mix of handwritten and machine-print images and one of the goals of the system is to separate the handwritten images from the machine-print images and process them appropriately, a machine print test deck of at least 250,000 pieces from the same 120 sites was also dispersed through the test deck prior to testing.

The handwritten image test deck consists of a collection of two sets of 500,000 images each. The first set of 500,000 images is made up of 25,000 samples collected

from each of 20 selected sites. 50,000 of these images selected randomly were truthed. The second image set consists of 500,000 images from 100 different field locations. 50,000 of these images selected randomly were also truthed.

A single machine print test deck with at least 250,000 images collected from the same 120 sites as the handwritten test deck including 100,000 truthed images dispersed throughout the deck will be used to evaluate machine print encoding performance.

8 Truthing

The error rate of an address interpretation system is a very crucial metric in the performance evaluation of the system. Incorrect encodes would not only increase the cost of the processing necessary to deliver the mail-piece but would also tarnish the reputation of the postal service. Hence, it is essential that the error rate of the system be measured accurately and a very low threshold of error rate imposed on the system.

The most effective way of measuring the error rate is to compare the encodes returned by the system against the encode determined by a human looking at the mailpiece with access to the relevant postal databases. This process of recording the correct encode corresponding to a mail-piece image by humans is termed *truthing*.

Truthing is by nature a labor-intensive and expensive process. Hence, it is important to come up with as small a sub-set of images as possible that is representative of the population of the entire mail-stream. The address interpretation system has to process an average of a 1,000,000 pieces of mail per day at over 250 mail processing centers around the country. Hence, the smallest number of images that must be sampled for truthing is still quite large.

8.1 CEDAR Truthing Laboratory

The CEDAR Truthing Laboratory is equipped with the following:

A SPARC-Center1000 server with 512 MB RAM and
 a 200 GB RAID array.

- 15 SUN workstations as clients inter-connected with a switched ethernet network.
- The SPARC-Center1000 is also set up as a ftp server and the images sent over the T1 line from USPS are down-loaded directly onto the disks attached to this server.
- The laboratory is set up as a sub-net separate from the rest of the CEDAR network and the images can be accessed only by physically logging into one of the machines in the truthing laboratory. This is to ensure that the access to the images is secure and that the images cannot be copied elsewhere or used outside of the sub-net.
- Truthers are assigned quotas and their output is monitored to ensure maximum truthing throughput.
- Access to the truthing laboratory is controlled with swipe card access which is logged and all possible efforts have been made to ensure the security of the images.
- A truthing manual addressing most of the problems that truthers encounter is available in the truthing lab. This includes a description of the available tools and their usage, truthing guidelines, and examples of cases that are not straight-forward and require special guidelines in resolving to their valid encodes..
 Efforts are made to ensure that the manual is as comprehensive as possible.

8.2 Truthing Tools

Truthing is a monotonous task and hence prone to errors caused by fatigue and inattention. Design of proper truthing tools can minimize these errors as well as increase the throughput.

It is quite clear from the complexity of the rules listed in section 5 that it would be difficult for the truthers to memorize all the rules for encoding. Therefore, it is important that sufficient care be taken in the design of the tools to eliminate the need for memorization of these rules as far as possible and also to provide visual cues and design ergonomic interfaces to minimize the possibility of error.

CEDAR has developed an integrated Image Evaluation System that allows entire sets of images to be truthed in a distributed environment. The system is written in Java for enhanced portability between platforms and uses a JDBC driver to interact with a Java database back-end for storing truth and result data.

The Image Evaluation System is based on a client server model. A multi-purpose client-server infrastructure handles the communication of data. This is used for the following tasks:

- Authentication of truthing clients
- Accessing image data sets
- Database queries
- Retrieval and modification of truth data

The Image Evaluation System provides the following utilities to assist the truthing process.

1. Address Block Truthing Interface

Fig 4 shows the interface used to record the truth of the address present on the mail-piece. When the truthers type in the ZIP code and the primary number (PO Box or street number) into the appropriate text fields, the city field gets highlighted and the pull down menu offers a list of valid city names and their variations based on a database lookup into the National City State file. The street name field also gets highlighted with a list of valid street names in that ZIP that have the user input street number. This database assisted tool allows the truthers to use the database to allow them to decipher illegible or ambiguous writing as well as reduce keystrokes by eliminating the need to type in the street and city names. Once the city and street names are selected, the secondary field and the firm name field get highlighted if and only if the building specified by the ZIP, street number and street name triple happens to contain these elements as per the database. These visual cues incorporated in the tool prevent the possibility of the truther overlooking possible information on the mailpiece. The tool also allows the recording of the ground truth as it appears on the mail-piece. ZIPtype information is also displayed prominently to attract the attention of the truther.

An additional feature of the tool is that it allows the truthers to record additional information about the mail-piece such as its orientation, type of script and any other relevant information. Fig. 6 shows a screen shot of this interface. The list of tags displayed can be configured easily.



Fig. 4. Address Block Truthing

The encodes corresponding to the address information input by the truther can also be viewed using this interface. (Fig. 5)

2. Database Query Interface

Fig 7 shows the database query interface. This interface provides access to the following useful database queries:

- (a) Query for all records matching a given ZIP code and primary number
- (b) Query for all records matching a given ZIP code and addon

- (c) Query for all records matching the finance number of a given ZIP and primary number This query is useful in resolving cases with patron ZIP code errors.
- (d) Query for ZIPtype
- (e) Query for list of Firms with their associated primary numbers given a ZIP code
- (f) Query for list of street names given a ZIP code
- (g) Query for foreign country

The query interface provides additional information to the truthers by color coding records of different record types.



Fig. 5. Encodes: Levels of Sortation

3. ASCII Address Matching Interface The tool also provides access to an ASCII address matching engine which takes as input a free-form ASCII address, parses it and applies the appropriate encoding rules and returns the valid levels of sort. This engine also has the ability to correct some patron errors.

Some of the functions supported in the CEDAR truthing environment are as follows:

 Image Scoring - The addresses on the mail-piece are transcribed using a GUI editor and the address blocks are passed through an ASCII address matching engine that encodes the address based on encoding rules provided by the USPS. The ASCII address blocks which cannot be encoded by the address matching engine beyond 5 digits are passed to a second truthing stage where human truthers use advanced database lookup tools to encode the address on the mail-piece.

- Image Resolution Images are "truthed" by two different truthers and differences if any are resolved by a third truther.
- 3. Vendor Review and Resolution A result file (in a specified format) from a vendor system run is imported and the vendor can review and resolve the vendor results against the official truth.

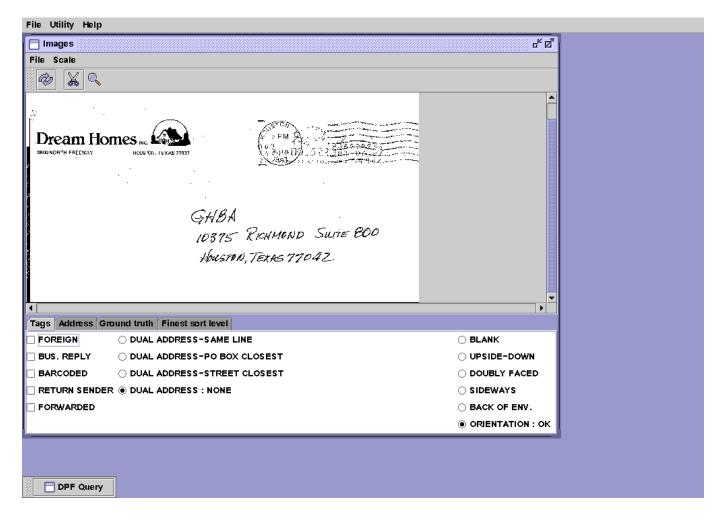


Fig. 6. Image Classification

- 4. Final Resolution The results produced by multiple vendors can be imported and the results from each vendor can be reviewed against the official truth simultaneously.
- Vendor Verification Grading of the results generated by a vendor system against the official truth to generate performance statistics.
- 6. Image Database System(IDS) A database system with functions to generate test decks from a given set of images based on user specified parameters.
- 7. Database Look-up Engine A customized USPS database is used for database look-ups. Since postal data changes

- over time and postal databases are constantly updated, the look-up engine should allow for the configuration of databases so that the database being used by vendor systems matches the database being used for truthing.
- 8. System Aided Truthing When it is necessary to reduce the turnaround time for truthing, the Address Interpretation System itself is used to generate results and then the system generated encode is verified and corrected if necessary. This provides a significant reduction in the time required to truth mail-pieces.

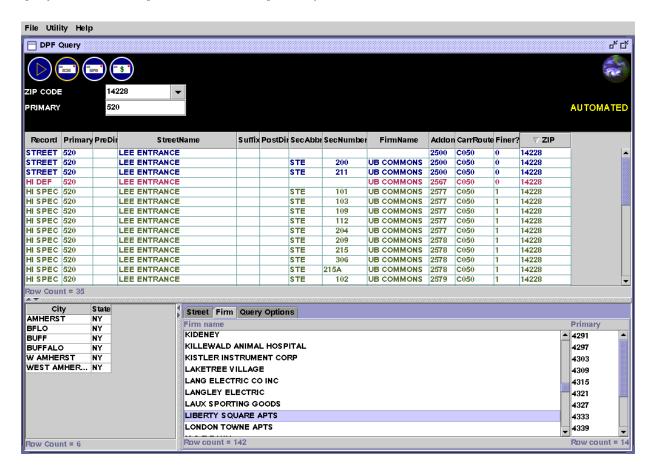


Fig. 7. Directory Look-up Interface

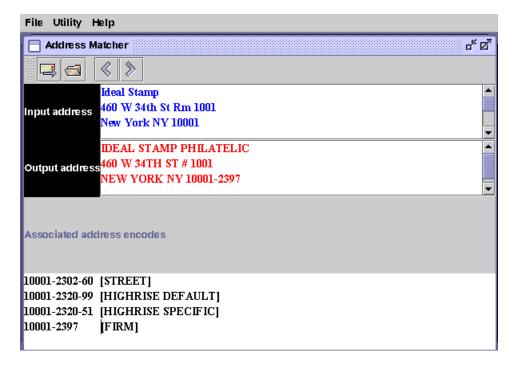


Fig. 8. ASCII Address Matching Utility

9 Summary

A comprehensive truthing and evaluation scenario for a large scale image processing system has been presented. It is clear that the nature of data being processed and how well the test set can mirror the characteristics of the real data population determine whether one should use a synthetic data set or representative samples of real data. Manual truthing of real data is by nature a monotonous and error prone activity. Design of ergonomic truthing tools with visual cues and redundant data integrity checks should help reduce truthing errors. A good training program and an efficient monitoring system can help maintain truthing quality.

References

- USPS Pub. 201 Postal Facts United States Postal Service, 1999.
- V. Govindaraju, A. Shekhawat and S.N. Srihari, "Interpretation of Handwritten Addresses in US Mail Stream", Proceedings of the Third International Workshop on Frontiers in Handwriting Recognition, pp 197-206, 1993.
- 3. USPS Annual Report United States Postal Service, 1999.
- 4. $USPS\ RCR2000\ Scoring\ Rules$ United States Postal Service, 2000.