

Truthing, Testing and Evaluation Issues in Complex Systems*

Srirangaraj Setlur Venu Govindaraju

Sargur Srihari

Center of Excellence for Document Analysis and Recognition, Amherst NY , USA

email:{setlur, govind, srihari}@cedar.buffalo.edu

Alfred Lawson

United States Postal Service, Merrifield VA, USA

email:alawson@email.usps.gov

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 per-

formance 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 pa-

*This work was supported by contracts from the USPS

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

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

2 Postal Automation

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. 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. The mail-

pieces that cannot be read by the MLOCs 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.[1]

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. 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.

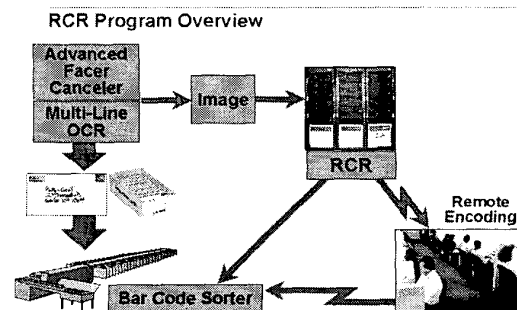


Figure 1: USPS Mail processing flow diagram

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, in-

cluding a reduction of 12 million work hours annually and cumulative labor savings of more than \$208 million. [2]

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 stan-

dards 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.

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 Encoding a mail-piece

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

is encoded.

e.g. 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 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

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	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	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

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.

3. Handwritten addresses present another source of confusion. Subjective interpretation of ambiguous digits and characters present problems in determining the truth unambiguously.

A complex set of scoring and evaluation rules has been designed by the USPS 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 (EZ4) or other directory query

programs.

5 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 2 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 various 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.

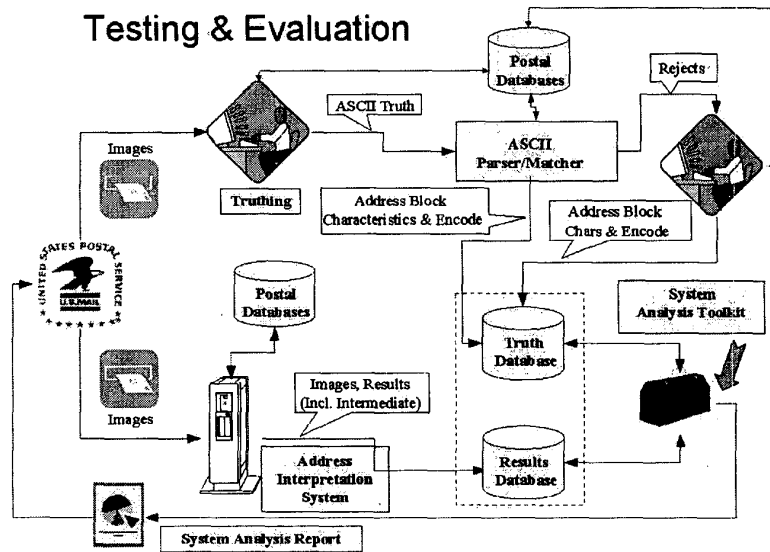


Figure 2: Testing and Evaluation System

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.

6 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 against the truth and a detailed performance 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 contemporary 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.

7 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 mail-piece 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*.

Due to the complex nature of the rules that govern the encoding of an address 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. A good training program and an efficient monitoring system have helped maintain truthing quality.

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 Summary

A 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.

References

- [1] V. Govindaraju, A. Shekhawat and S.N. Srihari, "Interpretation of Handwritten Addresses in US Mail Stream", *Proceedings of the Third IWFHR*, pp 197-206, 1993.
- [2] *USPS Annual Report* - United States Postal Service, 1999.
- [3] *USPS RCR2000 Scoring Rules* - United States Postal Service, 2000.