# ALEF BIT

## Hebrew on RLIN

*Joan M. Aliprand*
Research Libraries Group, Inc.
Stanford, CA

The Research Libraries Group, Inc. (RLG)[1] has undertaken a series of projects to add major non-Roman scripts such as Chinese, Cyrillic, and Hebrew to its automated bibliographic system, the Research Libraries Information Network (RLIN).

This paper describes features of the RLIN system pertaining to non-Roman scripts, particularly Hebrew. It was written at the beginning of 1987 when the design for Hebrew on RLIN had almost been completed, but before programming had begun to make the design a reality. The present tense has been used throughout the paper for consistency of style, although the future tense should properly be used in describing the Hebrew enhancements to RLIN.

This paper does not describe how one would create a Hebrew record using RLIN, because that information belongs in training materials; but, in brief: the user will be able to key Hebrew characters, from right to left, and search with Hebrew words or phrases; the RLIN system will require the cataloger to include romanized equivalents for certain "core" fields when they are entered in the original Hebrew.

### Non-Roman Phase I:
### Chinese, Japanese and Korean

RLG's first non-Roman project was "CJK"— the ability to input and search with Chinese characters, Japanese *kana* and Korean *hangul*. The RLIN CJK capability was introduced in September 1983. Since then, over a quarter of a million CJK records have been added to the database by participating libraries (including the Library of Congress and the British Library).

Chinese, Japanese and Korean are related because Chinese characters are part of the orthography of Japanese and Korean. But the CJK scripts and alphabetic scripts, such as Cyrillic and Hebrew, have little similarity. Nevertheless, the implementation of any non-Roman script on RLIN is bound to a common, generalized design.

Some of the design issues faced in the CJK project were specific to the scripts; for example, how does one uniquely encode each of the thousands of Chinese logograms? (For the answer to this question, see *Smith-Yoshimura & Tucker, 1985*.)

But more of the design issues related to the incorporation of a script other than the Roman alphabet into the RLIN system. This class of problem would have had to be solved for *any* non-Roman script, even an alphabetic one. Typical problems of this kind are:

> How is a user to switch back and forth between scripts?
>
> How is a change of script to be indicated in the stored data?
>
> If a user's terminal lacks non-Roman capability, can the user search for and see records containing non-Roman data?
>
> How is the non-Roman data in a record in the RLIN database to be ouput on tape for bibliographic exchange? What changes to the USMARC formats need to be introduced?

RLG's solutions to these particular problems, which have been implemented for CJK and then for Cyrillic (and will apply to Hebrew) are discussed below.

### Non-Roman Phase II: Cyrillic

The second non-Roman script to be added to RLIN was Cyrillic; specifically, the Russian alphabet and additional letters for other Slavic languages. (RLIN's implementation of Cyrillic does not include letters unique to non-Slavic languages, nor the alphabet of Old Church Slavonic.) Cyrillic was released in May 1986.

Cyrillic was the first non-Roman script to be added to the RLG-written software which makes a personal computer function as an RLIN terminal. Hebrew will be added to this software. Currently, CJK requires special-purpose terminals manufactured for RLG by the Transtech International Corporation, but will, in the future, be available on specially equipped personal computers.

Cyrillic preceded Hebrew because its implementation was a simpler task. Languages written in Cyrillic characters run from left-to-right, therefore no major revisions of displays and input techniques were needed. International standards for the computer encoding of Cyrillic characters existed, so that RLG neither had to determine what characters were needed for Cyrillic cataloging nor assign codes to them.

The RLIN East Asian Character Code (REACC) contains all the graphic symbols used to write the Chinese, Japanese and Korean languages; it is a single non-Roman character set. The CJK project, therefore, did not address the problem of combining more than one non-Roman character set in a bibliographic record. This problem had to be addressed, and solved, for Cyrillic, which has two separate character sets—the Russian alphabet ("Basic Cyrillic") and supplementary letters for other Slavic languages ("Extended Cyrillic").

### Non-Roman Phase III: Hebrew

RLG's current non-Roman script project is to add a Hebrew capability to the RLIN system and the RLIN terminal emulation software. (Some initial work on Hebrew was done while Cyrillic was still under development.) The scope of the "Hebrew" project is not restricted to the Hebrew language, but covers Yiddish and other languages usually written with Hebrew characters.

Hebrew has introduced a number of additional problems:

> What coding standard should be used?
>
> How will text in opposite directions be input and displayed?
>
> Are there any unique features of the Hebrew language that affect indexing?

## Computer Encoding of Hebrew Characters

One significant problem associated with computer manipulation of Hebrew text is that there is no international standard for the encoding of Hebrew characters, although there is an Israeli standard. RLG was faced with a similar problem for CJK, and devised the RLIN East Asian Character Code (REACC) by amalgamating various standards.

The principles which underly RLG's work on character sets are:

*The character set must allow a cataloger to transcribe bibliographic data as fully and accurately as possible.*

*The character set must encompass the alphabets of all languages written in the script, not just the predominant language.*

*The character set must not be limited to the transcription of bibliographic data from modern printed books, but must also permit the bibliographic description of historical material (including manuscripts) in variant orthographies.*

The Israeli standard, SI 960 (1976), is limited to modern, unpointed Hebrew, and so does not meet RLG's needs. The character set to be used for Hebrew cataloging on RLIN must be adequate for all Hebraic languages, including those in which diacritics, vowel points, and digraphs are significant components of the orthography.

The International Organization for Standardization, through several subcommittees, has been working on standards for the computer encoding of Hebrew for over 10 years (ISO, 1979; ISO, 1985).

The current (1985) draft of the ISO proposed standard for Hebrew consists of two character sets: a "basic set" consisting of the alphabet of the Hebrew language, the vowel points and semi-vowel (*sheva*), and a number of other diacritical marks (e.g., the *rafeh*) and a "supplementary set" of ancient Tiberian, Babylonian, Palestinian and Samaritan vowel points.[2]

RLG commissioned Bella Weinberg (whom the National Endowment for the Humanities had proposed to RLG as a technical consultant for Hebrew on RLIN) to survey the graphic repertoire necessary for Hebrew cataloging and determine whether any characters were lacking from the ISO draft standard. (At the time this work was undertaken, only the 1st draft, dated 1979, was available.)

The character set which Dr. Weinberg compiled, with advice from language and bibliographic experts, included digraphs and "pointed" letters found in Hebrew and Yiddish, diacritical marks found in Judezmo (Ladino) and Judeo-Arabic, and the composite *ḥaṭaf* vowels. The initial RLG proposal for a Hebrew character set appears in Figure 1.

RLG also established an Advisory Group on Hebrew; the members were: Edith Degani (*Jewish Theological Seminary*), Leonard Gold (*The New York Public Library*), Paul Maher and Sally McCallum (*Library of Congress*), Jonathan Rogers (*Yale University*), Bella Weinberg (*St. John's University & YIVO*), and Herbert Zafren/Philip Miller (*Hebrew Union College*).

When the Advisory Group evaluated the proposed Hebrew character set, many of the composite characters were excluded, under the generally-accepted computing axiom that it should not be possible to encode a particular graphic in two separate ways. For example, the Yiddish vowel *pasekh alef* had been assigned its own code *5E* (in hexadecimal notation), but could also be encoded as the vowel *patah* (*50*) and unmarked *alef* (*60*), i.e., by the sequence *5060*.[3]

Most members of the Advisory Group felt that the digraphs *tsvey vovn* (double *vav*), *tsvey yudn* (double *yod*) and *vov yud* were necessary for correct Yiddish orthography. [Yiddish pronunciation is being used in this context.] As a very large number of publications within Hebraica are in the Yiddish language, it is important that a character set designed for cataloging purposes accommodate the special orthographic features of this language. The Yiddish digraphs are unique graphic symbols; the evidence that this is so comes from typewriters and children's primers, as well as from the authoritative works on Yiddish orthography. Furthermore, the graphic *pasekh tsvey yudn*, in which the vowel is centered under the double *yod*, cannot be properly encoded by repeating the single letter *yod* (*yud*) together with a *patah* (*pasekh*); the digraph *tsvey yudn must* therefore have its own unique code.

In sum, a Hebrew character set consisting only of the 22 consonants is inadequate for Yiddish, because its use distorts the standard orthography of the language. Firstly, the digraphs cannot be represented unambiguously. Secondly, "pointed" letters which are part of standard Yiddish orthography require a pointed Hebrew font; information crucial to pronunciation and romanization,

e.g., the distinction between *pey* and *fey*, is lost if unpointed Hebrew letters are substituted.[4]

The final version of the Hebrew character set which will be available on RLIN appears in Figure 2. The character set consists of Hebrew letters, Yiddish digraphs, vowel points and the semi-vowel, diacritical marks, numerals, punctuation marks and other symbols; the names of all the characters are given in Note 5. The ISO supplementary set was, in general, considered superfluous for cataloging purposes: those symbols that were required were incorporated into the single RLIN Hebrew character set.

RLG will work with the appropriate national and international standards-making organizations to have the RLIN Hebrew character set approved as a standard for the exchange of Hebrew bibliographic data. [*The paper by Kuperman which follows this one underscores the need for standard coding of Hebrew characters.—Ed. (B.H.W.)*]

## Computer Representation of Hebrew Text

The principles for intermixing scripts in computer-encoded data are set forth in the standard entitled *American National Standard Code Extension Techniques for Use with the 7-Bit Coded Character Set of American National Standard Code for Information Interchange (ANSI X3.41-1974)*. These principles are followed for all non-Roman scripts on RLIN.

Text in a particular character set is preceded by a sequence of computer codes which identifies the character set ("here beginneth Hebrew"). In RLIN, the proper sequence is inserted by the RLIN terminal software as the user switches from one script to another. The "escape sequence" (so called because it begins with the control character "escape"), which defines the character set of the following text, is concealed from the user; the user sees only the effect of its presence, as text in a non-Roman script.

Hebrew text is stored in the RLIN database in logical order. This is equivalent to the sequence of keystrokes (taking into account corrections made via deletion, insertion and overstriking). This is also the order of Hebrew text for records distributed in the USMARC format.

"When fields contain escape sequences to languages written from right to left, the field will still be given in its logical order. For example, the first letter of a Hebrew title would be the eighth character in a field (following the indicators, a delimiter, a subfield code, and a three-character

escape sequence). The first letter would *not* appear just before the end of field character and proceed backwards to the beginning of the field." (*JOLA, 1981, p. 215*).

## Input

The procedures for input and display of multi-script text were developed as part of the CJK project. On the CJK terminal, Chinese characters, Japanese *kana*, Korean *hangul* and the Roman character set have been assigned separate keys; a particular set of graphemes is invoked by pushing the appropriate key.

The RLIN terminal emulation program has been designed to support a number of non-Roman scripts. Each is invoked by pressing the "character set invocation" key (currently, the key marked "Alt"), followed by a letter to designate the script; for example, "h" for "Hebrew." The designated script is "locking," i.e., the keyboard operates in that script until another script is invoked.

In the current (Cyrillic) RLIN terminal software, depressing the subfield delimiter symbol ("double dagger") key automatically invokes the Roman character set. If non-Roman text follows the subfield code, it must be explicitly re-invoked by the user. It is planned that, in future releases of the software, the non-Roman script currently in effect will merely be suspended and will automatically return after the subfield code has been keyed. For example, the cataloger need not re-invoke Hebrew after inputting the title proper and marking the beginning of other title information or the author statement.

The Roman and Cyrillic character sets are supplemented by "alternate" character sets—respectively, the diacritical marks and special letters supplementing ASCII[6] (the so-called "ALA extension to Roman") and the Cyrillic letters used in Slavic languages other than Russian. An alternate character is obtained by pressing the non-locking "alternate character set" key (one of the function keys) followed by the appropriate keyboard key. After the alternate character has been keyed, the keyboard immediately reverts to the alphabetic keyboard for the currently designated script. As indicated above, there is no "alternate" character set for Hebrew on RLIN.

With Hebrew as the currently-designated script, the keyboard is arranged as shown in Figure 3. The layout of the Hebrew consonants on the keyboard follows the pattern generally found on Hebrew typewriters. The vowels, which are arranged in a-e-i-o-u sequence on the bottom row of keys, are obtained by shifting. The diacritical marks and

Figure 1. Original RLG proposal for a Hebrew character set, based on ISO draft proposal, prepared August 1985.

the Yiddish digraphs are positioned as the "upper case" equivalent of a related consonant. For example, the *varika* (used in Ladino) is the shifted equivalent for *gimel* (on which letter it frequently appears); shifted *yod* (*yud*) gives *tsvey yudn*. Vowel points and diacritical marks are typed *before* the letter to which the point or mark applies and appear as separate characters on the screen; there is no "dead key" capability on an RLIN terminal.

Prototype Hebrew terminal emulation software was demonstrated to a number of Judaica librarians in June 1986 at the New York Public Library. In the prototype software, Hebrew was input by being inserted while the cursor remained stationary, giving the impression that the Hebrew was being "pushed" rightwards from the cursor position.

Although the Judaica librarians did not reject this approach, certain RLG staff members were not satisfied. It was subsequently demonstrated that there was a flaw in this input methodology: it is impossible to enter a Hebrew phrase containing a Roman word in the correct logical sequence. When the text is entered so as to read correctly on the screen, the internal logical sequence of the characters is not right; if the internal order is correct, the text on the screen is lexically incorrect.

The new design of the terminal software introduces the concept of a right-to-left directional "flag" for each field. The normal (default) direction for a field is left-to-right. However, the cataloger has the option of flagging a field as "right-to-left" by depressing the "field direction" key before beginning to type the contents of the field. Pressing the directional key causes the cursor to appear at the right-hand side of the screen. As Hebrew letters are keyed in normal order, the cursor advances to the left. Hebrew typing in a right-to-left field is the mirror image of Roman typing in a left-to-right field.

The cursor also advances when left-to-right text (e.g., Roman) is keyed in a right-to-left field, but, as letters are typed, the text grows at the end of the left-to-right text instead of above the cursor. (The text appears in normal—not backwards—order.)

Insertion and deletion are on a strictly logical basis. This sometimes causes the cursor to reposition itself at the end of a line, when insertion or deletion occurs at a directional boundary (as illustrated in Figure 4). This may be disconcerting at first, but is entirely logical; the logical approach is the only way to achieve consistent behavior in bi-directional character manipulation.



**Figure 2. Hebrew character set for RLIN, November 1986. (The names of the characters are given in Note 5.)**

## Parallel Fields

Traditionally, the body of the entry on a cata-log card has been in the original script; Rule 1.0E of *AACR2* (p. 15) states:

> "In the following areas, give information transcribed from the item itself in the lan-guage and script (wherever practicable) in which it appears there:
> Title and statement of responsibility
> Edition
> Publication, distribution, etc.
> Series"

This rule has been waived (*LCRI, 1984*) for cataloging in machine-readable form. All non-Roman data is represented by its romanized equivalent in computer records. Since the introduction of CJK on RLIN, how-ever, libraries have been able to revert to the original intent of Rule 1.0E, at least for their East Asian cataloging. The inclusion of non-Roman data in an RLIN record is a policy decision made by the cataloging li-brary. RLG cataloging standards do not mandate the use of RLIN's non-Roman capabilities.

A question which faced the designers of CJK was: what should be shown to a user on a Roman-only terminal when a CJK rec-ord is retrieved by a search? Many RLIN terminals that are used for searching can only display standard ASCII—the English alphabet, Western-style arabic numbers, English punctuation, and miscellaneous symbols such as "%" and "$". Transcribing the non-Roman data in non-Roman script only would mean that CJK (and other, fu-ture non-Roman) records would be essen-tially lost to the majority of users who would be unable to view the body of the entry.

RLG therefore developed the concept of "core fields" for non-Roman cataloging. If any of the following "core" fields—245 [Title statement], 250 [Edition statement], 260 [Publication, distribution, etc. (Imprint)], 4XX [Series statement]—exists in a record in ver-nacular (non-Roman) form, the RLIN sys-tem requires the user to precede that field with a parallel romanized equivalent. (Com-pare the "core fields" with the fields listed in Rule 1.0E of AACR2 cited above.) This re-quirement for parallel romanized fields en-ables a user to see an essentially complete— albeit romanized—record on any ASCII ter-minal (although the terminal may not be able to display diacritical marks, which are often part of a romanization scheme).

Non-Roman data may be included in any of the fields numbered 100 through 899 (i.e., main entry through tracings). The fields in an RLIN non-Roman record are defined as "core" or "non-core." Core fields are either
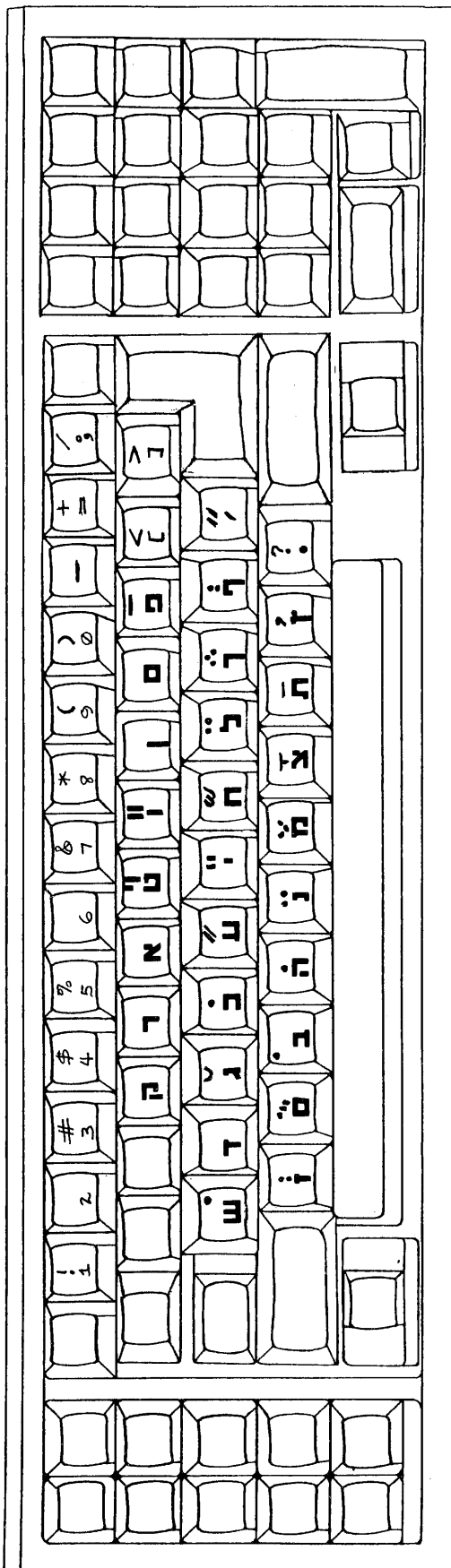


Figure 3. Hebrew keyboard for RLIN.

Roman/non-Roman pairs or unpaired Roman fields (e.g., English-language imprint or series title). Non-core fields may be any of: Roman/non-Roman pairs, unpaired Roman fields, or unpaired non-Roman fields. (Note that an unpaired non-Roman field cannot be seen on an ASCII terminal, because it lacks a romanized equivalent).

RLG defines "parallel field" conservatively: the romanized equivalent of a non-Roman field must be established by systematic romanization; it cannot be a translation or non-standard romanization from the work itself. For example, if a book has an edition statement in Hebrew and English, and (under the provisions of AACR2 Rule 1.2B5) the Hebrew edition statement is transcribed, the parallel romanized 250 field must contain the systematic romanization of the Hebrew edition statement, not the English edition statement that appears in the work.

This conservative interpretation of "parallel field" has particular implications for name access. Headings established according to AACR2 are in Roman script. If the AACR2 heading is a systematic romanization of the author's name, the non-Roman form may be included as a parallel field (at the library's option). If the AACR2 heading is not established by systematic romanization using the ALA/LC scheme, a parallel non-Roman form is not permitted. The Library of Congress uses the alternative to AACR2 Rule 22.3C2 for names written in Hebraic script, so that systematic romanization is used for a heading only when a romanized form of the author's name does not appear in the work nor in specified reference sources.

Libraries are, however, at liberty to include local access points in their records. When a non-Roman parallel to a particular heading is not permitted, vernacular name access may be achieved by entering as local headings the Roman/non-Roman pair for the systematically romanized form of the heading plus the Hebrew heading, or an unpaired Hebrew heading. The addition of vernacular access points for author and translator is shown on the cover.

The authority control for such local access points is the responsibility of the library. (Data in the RLIN Name Authorities File (NAF)—which includes name and series authority records created by the Library of Congress or by NACO participants—is currently limited to the Roman character set, and so cannot serve as a resource file for headings in Hebrew. RLG is well aware of the need to accommodate data in non-Roman scripts in the Authorities subsystem.)

## Indexing

The most notable feature of the RLIN bibliographic system is its powerful indexes. Searches may be done on names, words, or phrases, as well as on identifying numbers (such as ISBN). In addition, one may search on combinations of elements through logical (Boolean) operators (*and, or, and not*).

Truncation may be used to search for just the initial part of a word or phrase. For example, the truncated title word search Y'SRAL# (*the "pound sign" is the symbol for truncation*) will retrieve titles containing words such as Y'SRALYT, Y'SRALY, and Y'SRALDYQ&R, as well as those containing the word Y'SRAL. [This is the result of an actual RLIN search; the romanized Hebrew conforms to the modification of the ANSI standard for reversible romanization used by the New York Public Library.]

Non-Roman fields in RLIN records are indexed as fully as their romanized equivalents. The terms in RLIN searches may be in any of the character sets available on an RLIN terminal, or even a combination of scripts, for example, a Roman name plus a Hebrew title phrase.

At a minimum, Hebrew records will be searchable by Hebrew phrases or words from the title (title proper and subtitle) and any series notes. Since the title and the series note(s) are core fields, the RLIN system requires that they be in parallel romanized and non-Roman forms. Searching on Hebrew personal names, corporate names, and subjects is possible only if these optional access points have been included in records.

The RLIN indexing methodology was evaluated to assess its effect on Hebrew data, and to determine whether any special indexing procedures were needed for Hebrew. It was found that the following linguistic features would affect the indexing of Hebrew data: positional (medial vs. final) forms of letters; vocalization and diacritical marks; and prefixed particles.

In indexing, the medial and final forms of letters (e.g., *nun* and *nun sofit*) are not differentiated, although they *are* distinct in input and display. The process of substituting one character for another is called "normalization." Because of normalization, a search on the Yiddish word *oyf* will retrieve both titles containing the standard Yiddish form and titles containing the Soviet form using medial *fey* in final position. (In the search argument, the word may be spelled with either a medial or a final *fey*.) Normalization is also applied to the Yiddish digraphs, since there is no way for the searcher to know whether the cataloger keyed them as digraphs or as the component letters (LC intends to do the latter).

Normalization is akin to the procedures applied when filing in a Hebrew script catalog. The final form of a letter is filed before the medial form because it is usually followed by a space (or a hyphen regarded as a space in filing) and, according to the general rule of word-by-word filing, "nothing [the space] files before something [a letter]." A medial form in final position interfiles with the final form. The Yiddish digraphs are not treated as alphabetically discrete letters, but are filed as the component Hebrew letters. All these substitutions are done cerebrally by a filer.

Vowel points and diacritical marks are dropped in indexing. (The librarians advising RLG on Hebrew rejected another set of Hebrew indexes in which vowel points and diacritical marks would have been significant.) Unvocalized and vocalized Hebrew will both be retrieved, whether the search uses vocalized or unvocalized spelling. RLIN will not, however, automatically merge full and defective Hebrew spelling, just as it does not unify British and American spelling. The dropping of diacritics will cause "false drops" in Yiddish, since different letters will be normalized to the same unmarked letter in indexing. [This is what is generally done in library filing, although most Yiddish reference works separate *pey* and *fey*, etc.]

The RLIN title word index can be used for topical searching when standard subject headings are inadequate; therefore, it seemed worthwhile to eliminate grammatical particles for Hebrew word indexing, so that basic word forms, primarily nouns, could be used for searching. RLIN already has this on a minor scale: the romanizations of the Hebrew definite article *ha-* and *he-*, as well as the romanized Arabic prefixes *al-* and *el-*, are ignored in title word indexing. (Truncated searching—described above—can be used to deal with another grammatical problem: inflections or suffixes at the end of words.)

Several alternative strategies for ignoring the particles were considered. A crude solution would be to index words more than once: with and without the letters which could possibly be particles (e.g., *he*). It was recognized that doing this would result in "false drops" for words in which the letter is significant, e.g., *histadrut*.

A strategy which was explored in considerable detail before being rejected was the

use of the parallel romanized field as a guide to the presence of particles. ALA/LC romanization requires that the particle(s) be separated from the word proper by a hyphen (e.g., *shebe-, uve-*). It was proposed that when a romanized word begins with a hyphenated prefix, the initial letter(s) of the corresponding Hebrew word be skipped in indexing. Particles embedded in hyphenated Hebrew words could also be detected and eliminated, and the two individual words, minus particle(s), would be indexed separately in the appropriate word index. The convention of using the hyphen to demarcate the (romanized) particle(s) would have to be extended to reversible romanization.

The use of the romanized title or series note as a template for indexing the parallel Hebrew field for the title word index was found to fail in a number of situations; for example, when the number of words in the romanized and vernacular fields is unequal. This happens when a Hebrew abbreviation using *gershayim* is romanized to two separate abbreviations, e.g., *sh. u-t.* (responsa). An example using a proper name is *Ar. ha-B.*, from the Hebrew ARH"B, the abbreviation for *Artsot ha-Berit*, the Hebrew term for "United States." Another case where ALA/LC romanization changes the word count is when a number is pronounced as several words and is romanized as pronounced. The exceptions, in themselves, did not cause the strategy to be rejected, since they were clearly defined cases which could be documented in the searching manual.

Technically, the comparison process was feasible, but its programming would be a major task (which had not been included in time-estimates for Hebrew). The size of the task, and the unreliability of the algorithm under certain circumstances, caused this strategy to be abandoned.

The chosen strategy is to have a special "separator" character. The cataloger will position the "separator" between the particle(s) and the word proper. In word indexing, the Hebrew letter(s) before the separator (and the separator itself) will be disregarded. A composite (hyphenated) Hebrew word, with a particle in the second half, must be split into its component words for word indexing. Because the second half is being treated as an individual word, the particle is detected. For example, *bi-yeme — ha-benayim* is separated into *bi-yeme* and *ha-benayim*, which are indexed as *yeme* and *benayim* (with a normalized medial *mem*).

**Search Results and "Clustering"**

The RLIN database does not have a "master record" structure, but retains cataloging contributed by libraries as individual, discrete records. An incoming LC MARC record does not "bump" an existing record for the same work, but is treated as simply another record being contributed to the database.

Since each record is individually indexed, a search may retrieve more than one record with identical cataloging. To eliminate this duplication, records have been "clustered," which means that records for the same bibliographic entity have been grouped together. The elements of comparison used to establish matching records include identifying numbers (such as ISBN and LCCN), title, and imprint, but not headings.

The record with the highest level of cataloging is used as the representative record for all the records which match ("cluster together"). This representative record is termed the "primary cluster member"; all records which match it are "secondary cluster members." The primary cluster member is used in displays unless a particular secondary cluster member record is specifically requested.

The romanized forms of fields in non-Roman records are used in the matching process. Because of this, completely romanized records and those with non-Roman data can belong to the same cluster. In clustering, no preference is given to records containing non-Roman data. Thus the primary cluster member may be a completely romanized record, even though there are records containing non-Roman data in the cluster. (In the Primary display, the secondary cluster members which contain non-Roman data are identified; these records may be displayed individually.)

In a clustered file, a search term which matches an access point in *any* of the records in a cluster retrieves the entire cluster. For a Hebrew search, this means that, if only one of the records in a cluster contains Hebrew data, all the other, completely romanized records in the cluster will also be retrieved.

The clustering of the records for Hebrew language works in the RLIN database is an oddity, since two different romanization schemes have been used. The university libraries romanize according to the ALA/LC scheme. The New York Public Library uses (modified) ANSI reversible romanization for

```
SUCCESSIVE DELETION OF ROMAN CHARACTERS:

abcdefg ----> abcefg  ----> abcfg


SUCCESSIVE DELETION OF HEBREW CHARACTERS:

ZWHDGB@ ---->  ZWHGB@ ---->   ZWGB@


SUCCESSIVE DELETION BEGINNING AT HEBREW/ROMAN BOUNDARY:

Right to left field

abcGB@ ---->  abcB@ ---->   bcB@

   This looks extraordinary, but the letters are being deleted
   in the order in which they were keyed; i.e., in "logical"
   order, thus:

123456        12456        1256
@BGabc ---->  @Babc ---->  @Bbc
```

**Figure 4. Deletion at a Script Boundary**

the body of the entry.[7] (For the Library's printed *Dictionary Catalog* (NYPL, 1972–81), the Hebrew was re-created from the reversibly romanized data.) Since the romanized fields used for clustering (e.g., the title) do not match, reversibly romanized and ALA/LC romanized records cannot cluster together. So, for every Hebrew language work, there may be two clusters in the database. (This ignores the problem of multiple clusters resulting from variant application of the ALA/LC romanization scheme.[8])

The New York Public Library establishes name headings according to AACR2 and ALA/LC romanization. A name search using ALA/LC romanization will retrieve both clusters. A title word or phrase search using romanized Hebrew will retrieve only one cluster, unless both romanizations are strung together in the search with a logical "*or.*" When both clusters include records containing vernacular Hebrew, a Hebrew title word or phrase search will retrieve both clusters.

### Transcription of Bibliographic Data

Although there is great interest in the keying of Hebrew on RLIN (and dismay at the prospect of having to key both a Hebrew field and its romanized counterpart), much of the cataloging on RLIN is done with a DERive or CREate * command after a successful search for cataloging copy.

The ability to transcribe another library's cataloging by a single command is a powerful and productive feature of RLIN. Once a library catalogs a title to the agreed-upon RLG standard, no other library has to repeat the task. The CJK experience is that, even with "double keying" of romanized as well as vernacular fields, the productivity of catalogers was increased through the use of RLIN (*Reardon-Anderson, 1985*).

It might be thought that the efficiency of shared cataloging cannot be realized until there is an accumulation of Hebrew records in the RLIN database, and that, initially, Hebrew records must be keyed in their entirety. The RLIN database already contains an undetermined but significant number of completely romanized records for Hebrew works (including LC MARC records). Many are in ALA/LC romanization; those entered by the New York Public Library are in reversible romanization. Brandeis University has recently begun a retrospective conversion project to add its Hebrew cataloging to RLIN. A completely romanized record may be transcribed by command, to serve as the basis for a Hebrew record; all that must be

added are the parallel Hebrew fields (when the romanization and headings in the source record are acceptable).

In the transcription of bibliographic data, it should be noted that Hebrew vocalization is transcribed only as mandated by *AACR2*; that is, vocalization that is not present on the piece should not be added. (ALA/LC romanization, on the other hand, does require the cataloger to supply vocalic values.)

### "Double Keying"

The requirement for parallel core fields has caused some libraries to be reluctant to use RLIN's non-Roman capabilities because of "double keying;" that is, the input of the romanized equivalent of non-Roman text in a parallel field.

The reasons for romanized core fields paralleling the non-Roman fields of the body of the entry have been described above: to allow a non-Roman record to be seen on a Roman-only terminal, and to allow completely romanized records to cluster with non-Roman records for the same bibliographical item.

The choice of romanization scheme to be used for the core fields is up to each library. RLG's BibTech Committee is evaluating the consequences of allowing ANSI reversible romanization as an alternative to ALA/LC romanization in "standard" Hebrew records. This discussion of "double keying" does not address the pros and cons of the two romanization schemes, but merely the additional work imposed by the RLIN requirement that each non-Roman core field be preceded by a romanized equivalent.

Double keying is not necessary when a Hebrew record for the title being cataloged already exists in the RLIN database. Keying of both romanized and Hebrew vernacular fields must be done *only* for titles not in the database or when there is only a record completely romanized according to an unacceptable scheme.

The RLIN system requires "double keying" only for the "core fields." The inclusion of all other paired romanized and non-Roman fields is at the cataloger's option. Additional "double keying" to provide vernacular access points in the record is certainly worthwhile, even if the headings are "uncontrolled," but vernacular headings are not mandatory in the RLIN system.

Every record includes a title (245 field), and almost all records include details of publication or production (260 field). The two other core fields—edition statement and series note(s)—do NOT occur in every record. Edition statements occur in 15% of RLIN

records, and series notes in, at most, 30% of records (figures based on *Crawford, 1986, p. 307*). Thus, just over half the time, only the title and imprint will have to be "double keyed," and only in records for which suitable copy cannot be found in the database. The "title" is not just the title proper, but the complete title statement including subtitle, parallel title(s) and statement(s) of responsibility.

### Output

The data in an RLIN record can be presented in various ways and on different media. There are a number of predefined RLIN displays which may be used to view a record on a terminal. The design of the RLIN displays for bidirectional text follows the provisions of ISBD(G) (*UBC, s.d.*) In essence, this means breaking to a new line when the direction of the text changes.

The card-like "Long" display is illustrated on the cover. The "Partial" display shows a library's record up to and including any series notes, plus information about the library's holdings of the title (and also its acquisition, if it was ordered through RLIN). This display is useful in technical processing and in reference work, since it shows location and call number for each copy. A later addition to RLIN CJK was a display to simulate a Library of Congress card. This was added at the request of the Library of Congress, to support its Asian card printing operation. Although the Library of Congress has requested a Hebrew variant of this display, it will not be present in the initial implementation of Hebrew.

An RLIN terminal display can be reproduced on paper if there is a printer attached to the terminal. The full range of non-Roman characters in the RLIN terminal emulation software is supported on the printers recommended by RLG; Hebrew characters are shown in Figure 5.

Hebrew catalog cards are not part of RLG's Hebrew project. While RLG is well aware of the demand for non-Roman catalog cards (from East Asian libraries, as well as from Judaic libraries), the inclusion of card formatting and printing would expand each non-Roman project considerably. In particular, the provision of CJK cards was judged to be a formidable and expensive task, and was therefore excluded from the CJK project.

The trend in general academic libraries is towards closing the card catalog and substituting online devices for patron access. Although the RLIN system was not intended to serve as an online catalog, its indexes provide all the access points found in a traditional card catalog, plus those which would
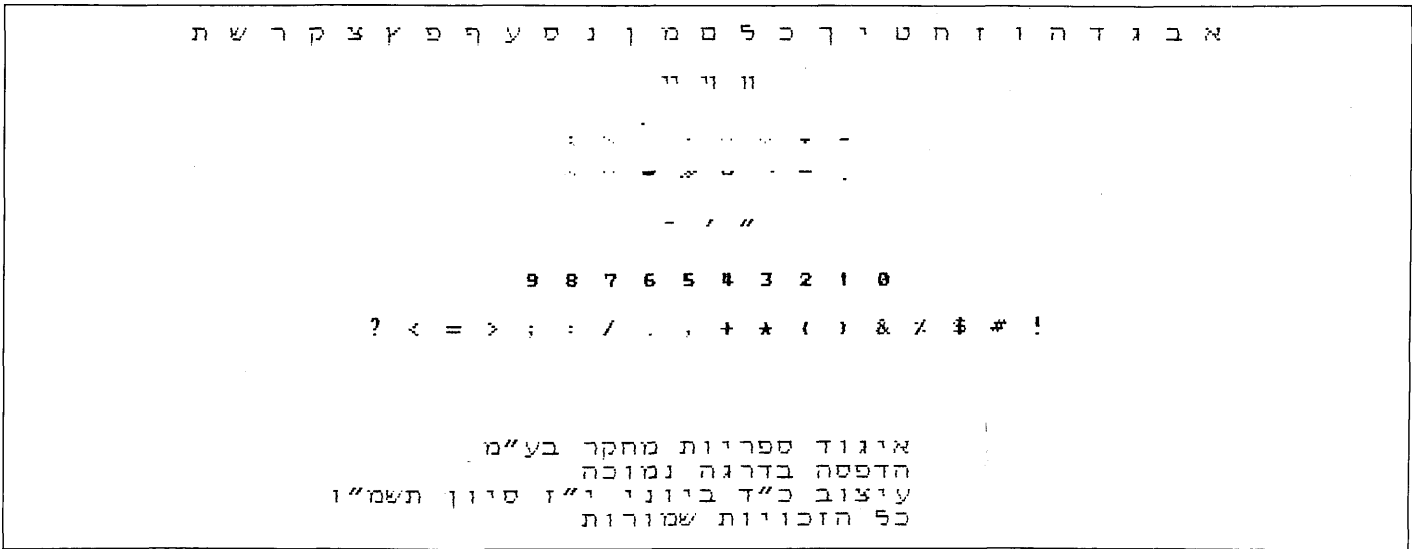
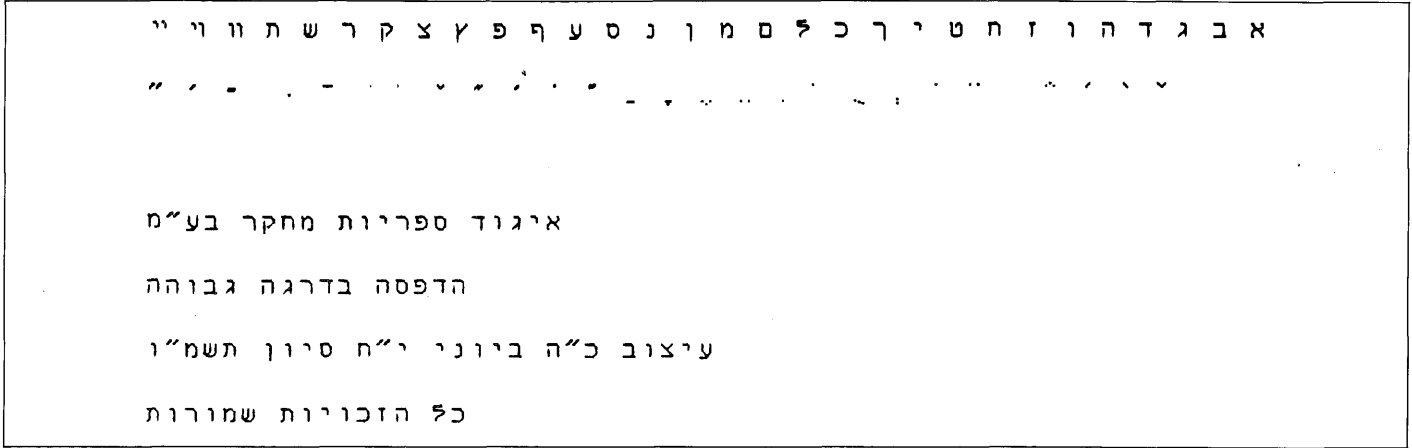**Figure 5a. Hebrew character set, designed June 1986, printed on IBM 80 CPS Matrix/Graphics Printer.**

**Figure 5b. Hebrew character set, designed June 1986, printed on C. Itoh CI-3500 Series High Speed Printer, Model 20.**

be impossible in a card catalog (e.g., access by each word in the title).

Is it practical to use RLIN as a substitute for the card catalog? Each library must answer this question for itself. How does the cost of the card catalog—including the labor for filing new cards and "maintaining" existing ones—compare with RLIN costs? What capabilities does the card catalog have that RLIN does not, and vice versa? Can added expense be justified by unique capabilities?

Hebrew data in RLIN can also be supplied in machine-readable form, on tapes in the USMARC format. The CJK project necessitated extensions to USMARC, to accommodate non-Roman data. The 066 field—

Character Sets Present—identifies every non-Roman character set that occurs in a record, by the escape sequence which introduces each. The non-Roman data itself (which occurs in parallel fields online) is in successive occurrences of the 880 field—Alternate Graphic Representation—in a USMARC record. The contents of an 880 field include the tag (e.g., 245) of the romanized field which this non-Roman data parallels. Online and tape versions of the variable field data in an RLIN Hebrew record are illustrated in Figure 6.

**Cover: Simulation of "Long" display of a Hebrew record on RLIN. "Long", a card-like display, is only one of several different ways to view an RLIN record online.**

**The romanized Hebrew in this example was copied from the Library of Congress MARC record; the Yiddish appears on the LC card in romanized form and was reconstructed.**

```
            As seen online                          As communicated
           ─────────────                        in the USMARC format
                                                ──────────────────────

                                          066    ‡c(2
100 10 ‡aAACR2 main entry                 100 10 ‡6880-01‡aAACR2 main entry
100 10 ‡aHebrew main entry
245 14 ‡aRomanized title                  245 14 ‡6880-02‡aRomanized title
245 10 ‡aHebrew title
250    ‡aRomanized edition statement      250    ‡6880-03‡aRomanized edition statement
250    ‡aHebrew edition statement
260    ‡aRomanized imprint                260    ‡6880-04‡aRomanized imprint
260    ‡aHebrew imprint
300    ‡aCollation                        300    ‡aCollation
500    ‡aNote in Roman text               500    ‡aNote in Roman text
505    ‡aRomanized contents note          505    ‡6880-05‡aRomanized contents note
505    ‡aHebrew contents note
700 10 ‡aRomanized added entry            700 10 ‡6880-06‡aRomanized added entry
700 10 ‡aPaired Hebrew added entry
700 10 ‡aSolo Roman added entry           700 10 ‡aSolo Roman added entry
700 10 ‡aSolo Hebrew added entry

                                          880 10 ‡6100-01/(2‡aHebrew main entry
                                          880 10 ‡6245-02/(2‡aHebrew title
                                          880    ‡6250-03/(2‡aHebrew edition statement
                                          880    ‡6260-04/(2‡aHebrew imprint
                                          880    ‡6505-05/(2‡aHebrew contents note
                                          880 10 ‡6700-06/(2‡aPaired Hebrew added entry
                                          880 10 ‡6700-00/(2‡aSolo Hebrew added entry
```

Figure 6. Variable fields of a Hebrew record: on RLIN and in USMARC format.

## Summation

The next few years will be eventful ones in Judaica librarianship. It is exciting to see Hebrew on a terminal, but RLIN is far more than a substitute for a Hebrew typewriter. Using RLIN to acquire and catalog Hebraica means that libraries will be able to share information for collection development, cataloging, and inter-library loan.

The active field of Hebraic/Judaic cataloging will become even more so when librarians have rapid access to other catalogers' work. Questions about the application of AACR2 to non-Roman cataloging, and to unique problems found in Hebraic material, have already arisen. Arbitration on aspects of ALA/LC romanization will undoubtedly be necessary. The number of LC Rule Interpretations will be increased as questions relating to Hebrew cataloging are answered.

Using Hebrew for searching is an exciting prospect, but it will be of limited utility if libraries only provide the required "core" fields. RLIN allows (but does not require) headings in the vernacular for persons, corporate bodies, places, and topical subjects which cannot be adequately named in English.

The long-term impact of RLIN will be in collection development, preservation and retrospective conversion. Use of RLIN will not only facilitate data capture, but will make the data available to all RLIN users. National (and even international) cooperation in the field of Hebraica will be stimulated through libraries' use of RLIN and through their participation in activities sponsored by RLG.

## Notes

¹The Research Libraries Group, Inc., is a non-profit corporation owned and operated by its members—the libraries of major universities and research institutions in the United States. In addition, RLG has many "programmatic members"—institutions that participate in one or more of the corporation's programs.

RLG's integrated set of cooperative programs aids members in the areas of collection management and development, shared resources, preservation, general bibliographic access and control, and access to and management of specific forms of research information. RLG's automated information system, RLIN, combines data bases and computer systems to support these programs. RLIN, a nationwide network, serves both RLG members and non-member institutions, including public, academic, and special libraries.

²I crave the indulgence of the reader with regard to the names and romanizations used to identify characters. In discussing Yiddish usage, I have used Yiddish terms; otherwise, I have followed ALA/LC.

³A particular graphic image is referenced by its column and row coordinates. Thus, the medial *nun* in the upper right-hand corner is designated by the code "70" (column 7, row 0). Rows 10 through 15 can also be designated as "A" through "F." This designation conforms to hexadecimal (base 16) notation.

⁴The Library of Congress is, however, of a different opinion, and has stated:

> 'Also after discussion we feel it advisable not to use the double yod, double vav, and vav yod, but to key these characters as separates. Our people feel that generally the typographies of items do not clearly indicate the double characters. They would thus be constantly problematic to key "correctly."' (*Letter to the author from Henriette Avram, 12/11/85*).

This is at variance with the ALA/LC romanization scheme for Hebrew and Yiddish (*Cataloging Service, Bulletin*, 118 (Summer 1976), p. 63), where double *yod*, double *vav*, and *vav yod* have distinctive romanizations, and are *not* romanized as the component single characters.

⁵The characters of the RLIN Hebrew character set are:

| Position | Name | Notes |
|---|---|---|
| 2/0 | SPACE | |
| 2/1 | exclamation mark | |
| 2/2 | gershayim | Quotes in ASCII |
| 2/3 | number sign | |
| 2/4 | dollar sign | |
| 2/5 | percent sign | |
| 2/6 | ampersand | |
| 2/7 | geresh | Apostrophe in ASCII |
| 2/8 | opening parenthesis | Consistent with ASMO Arabic (ISO Regn. 89) |
| 2/9 | closing parenthesis | Consistent with ASMO Arabic (ISO Regn. 89) |
| 2/10 | asterisk | |
| 2/11 | plus sign | |
| 2/12 | Hebrew comma | Comma in ASCII |
| 2/13 | makef | Hyphen in ASCII |
| 2/14 | period | |
| 2/15 | slash | |
| 3/0 | zero | 3/0-3/9: numerals zero through nine in ASCII |
| 3/1 | one | |
| 3/2 | two | |
| 3/3 | three | |
| 3/4 | four | |
| 3/5 | five | |
| 3/6 | six | |
| 3/7 | seven | |
| 3/8 | eight | |
| 3/9 | nine | |
| 3/10 | colon | |
| 3/11 | semi-colon | |
| 3/12 | Hebrew less-than sign | |
| 3/13 | equals sign | |
| 3/14 | Hebrew greater-than sign | |
| 3/15 | question mark | |
| 4/0 | patah | |
| 4/1 | kamats | |
| 4/2 | segol | |
| 4/3 | tsereh | |
| 4/4 | hirik | |
| 4/5 | holom | |
| 4/6 | kubuts | |
| 4/7 | sheva | |
| 4/8 | dagesh/mapik | |
| 4/9 | rafeh | |
| 4/10 | right shin dot | Use holom for left sin dot |
| 4/11 | varika | |
| 4/12 | double acute | |
| 4/13 | shaddah | ISO Hebrew Set 2: |
| 4/14 | super-script tsereh (ta'-marbutah) | Babylonian zere, Palestinian qibbuz |
| 4/15 | inverted segol | 5/0-5/10, 5/12, 5/14 and 5/15 have no assigned character |
| 5/11 | opening bracket | |
| 5/13 | closing bracket | |
| 6/0 | alef | |
| 6/1 | unmarked bet | |
| 6/2 | gimel | |
| 6/3 | dalet | |
| 6/4 | he | |
| 6/5 | vav | |
| 6/6 | zayin | |
| 6/7 | het | |

**Note 5** *(continued)*

| Position | Name |
|---|---|
| 6/8 | ṭet |
| 6/9 | yod |
| 6/10 | final khaf |
| 6/11 | unmarked kaf |
| 6/12 | lamed |
| 6/13 | final mem |
| 6/14 | mem |
| 6/15 | final nun |
| 7/0 | nun |
| 7/1 | samekh |
| 7/2 | 'ayin |
| 7/3 | final fe |
| 7/4 | unmarked pe |
| 7/5 | final tsadi |
| 7/6 | tsadi |
| 7/7 | ḳof |
| 7/8 | resh |
| 7/9 | unmarked shin |
| 7/10 | unmarked taṿ |
| 7/11 | tsvey vovn |
| 7/12 | vov yud |
| 7/13 | tsvey yudn |

7/14 has no assigned character

Substitutions:
diacritical prime: use geresh
thousands dot: use geresh
acute: use Roman acute
grave: use Roman grave
hacek: use Roman hacek
supra-linear dot: use Roman dot above

[6]The acronym ASCII stands for "American Standard Code for Information Interchange" (*ANSI X3.4-1977*) and includes the following graphic characters: the letters of the English alphabet, Western-style arabic numbers, English punctuation, and miscellaneous symbols such as "%" and "$".

[7]RLG's Library Technical Systems and Bibliographic Control Program Committee (BibTech) is currently considering Hebrew romanization as it relates to RLG's standards for cataloging. If a record containing Hebrew script is to be "standard," must the ALA/LC scheme be used for the romanization of the required "core fields" and any optional notes, *or* will reversible romanization be permitted for these fields? If BibTech decrees that only the ALA/LC scheme is "standard," a library that prefers to use reversible romanization for non-heading fields will be penalized, i.e., must pay the rate for non-standard cataloging.

Rates for RLIN Cataloging are levied according to the contribution that the cataloging makes to the RLIN database, primarily, its usefulness to other libraries as a source of cataloging copy. There is no charge for a record that represents a title new to the database and that is fully content designated in accordance with RLG's standards for cataloging (which are based on AACR2,

the USMARC formats, and authorized subject heading lists such as LCSH). The highest charge for cataloging is imposed for pure copy cataloging ("derivative—not upgraded") or for non-standard original cataloging. Thus, there is a financial incentive to contribute "standard" original cataloging to the RLIN database, or to upgrade existing cataloging to standard.

[8]In an unpublished experiment conducted at an AJL Cataloging Workshop with a group of experienced Judaica catalogers, 22 variant romanizations were produced for one title (illustrated on the cover). These would have formed 21 separate clusters in the RLIN database. The ability to view other catalogers' work on RLIN might lead to less variation.

---

## References

**[AACR2]** *Anglo-American Cataloging Rules.* 2d ed. Chicago: American Library Association, 1978.

**[ANSI X3.4-1977]** American National Standards Institute. *American National Standard Code for Information Interchange.* New York, 1977.

**[ANSI X3.41-1974]** American National Standards Institute. *American National Standard Code Extension Techniques for Use with the 7-Bit Coded Character Set of American National Standard Code for Information Interchange.* New York, 1974.

The international form of this standard is:

International Organization for Standardization. *Information Processing—ISO 7-Bit and 8-Bit Coded Character Sets—Code Extension Techniques.* 2d ed. [s.l.], 1982. (International Standard, ISO 2022-1982).

**Crawford, Walt.** *Bibliographic Displays in the Online Catalog.* White Plains: Knowledge Industry Publications, 1986.

**[ISO, 1979]** International Organization for Standardization. Technical Committee 46. Sub-Committee 4. Working Group 1. *Hebrew Alphabet Character Sets for Bibliographic Use.* 2d draft. Ghent, 1979. (ISO/TC46/SC4/WG1 N98).

**[ISO, 1985]** International Organization for Standardization. Technical Committee 46. Sub-Committee 4. Working Group 1. *Hebrew Alphabet Character Sets for Bibliographic Information Interchange.* [Draft proposal]. [s.l.], 1985. (ISO/TC46/SC4/WG1 DP).

**[JOLA]** "Inclusion of Nonroman Character Sets." *Journal of Library Automation,* vol. 14, no. 3 (Sept. 1981), pp. 210-15.

**[LCRI]** "Library of Congress Rule Interpretations (LCRI)—1.0E. Language and script of the description." [Rev.] In: *Cataloging Service Bulletin,* No. 27 (Winter 1985), p. 9-17. Originally promulgated in Nov. 1984 by the Library of Congress.

**[NYPL]** New York Public Library. Research Libraries. *Dictionary Catalog of the Research Libraries: a cumulative list of authors, titles, and subjects representing books and book-like materials added to the collections since January 1, 1971.* [New York]: New York Public Library, [1972-81].

**Reardon-Anderson, J.** "RLIN/CJK: the First Year." New York: C.V. Starr East Asian Library, Columbia University, March 1985. (*Unpublished paper, available from Publications Coordinator, RLG*).

**[SI 960]** The Standards Institution of Israel. *7 Bit Coded Character Set for Information Interchange (Translation).* Tel Aviv, 1976. (Israel Standard SI 960). [*English translation issued by the Standards Institution of Israel. "Only the original standard in Hebrew is authentic."* ]

**Smith-Yoshimura, K.; Tucker, A.** "RLIN East Asian Character Code and the RLIN CJK Thesaurus." *In:* Asian-Pacific Conference on Library Science (2nd : 1985 : Seoul, Korea). *Proceedings of the Second Asian-Pacific Conference on Library Science: 20–24 May, 1985, Seoul.* Seoul, Korea: Cultural and Social Centre for the Asian and Pacific Region, c1985, pp. 105–37.

**[UBC]** IFLA International Office for UBC. *Changes to ISBD(G) as a Consequence of Decisions Taken at the ISBD Review Committee Meeting, January 1983.* London, [198-]. (IFLA/UBC/C7.5.6).

---

*Joan Aliprand is employed by the Research Libraries Group, Inc. as a Library Systems Analyst; she is currently writing the external design specifications for adding Hebrew to RLIN. Ms. Aliprand is a graduate of the School of Librarianship, University of New South Wales and also studied at the Graduate Library School, University of Chicago. She has held professional librarian positions at the University of California, Berkeley, the University of Chicago, and Macquarie University (in Sydney, Australia).*