

7 Bits Good, 8 Bits Bad or “The Eight-Bit Blight”

Malcolm Clark
Polytechnic of Central London
malcolmc@mole.pcl.ac.uk

Brian Hamilton Kelly
Royal Military College of Science
Shrivenham
tex@rmcs.cranfield.ac.uk

Niel Kempson
25 Whitethorn Drive
Cheltenham
tex@rmcs.cranfield.ac.uk

Abstract

Inter-networking and e-mail systems can usually be relied upon to permit faithful exchange of seven-bit ASCII data. Transfer of eight-bit binary data is not so reliable, especially when the data must traverse gateways or be exchanged between different system types.

Until now, it has been possible to exchange T_EX sources of papers by electronic mail without much difficulty. Now that T_EX and its relations support eight-bit input their source files will now suffer the same problems as binary data.

The proliferation of electronic archive services has highlighted the need to be able to exchange binary data between disparate systems, often connected via gateways. The authors introduce a new file-encoding standard that meets this need and far surpasses existing schemes.

Reliable and faithful exchange of binary files between computers over networks is a well-known problem, especially if the computers use different operating systems and are connected to different networks via a gateway. Unfortunately inter-networking and electronic mail are very much children of the '60s: they might have had to wait until the '70s for their naissance, but their progenitors were mentally locked-in to the concept of the 7-bit ASCII code for conveying textual information. The T_EX community has long been aware of this problem when trying to exchange “machine-independent” .dvi files and font-related data such as .tfm and .pk files. It has sometimes been possible to exchange this binary data by using encoding schemes that allow the data to be represented using a subset of the seven-bit ASCII character set.

Academics and authors in many fields have hitherto been able to pass .tex files back and forth

by electronic mail — apart from a few minor quirks and blemishes, such T_EX source files pass unharmed across the planet's networks. Problems are encountered when mail passes through certain gateway machines that introduce irreversible character corruptions. Particularly notorious is the Janet/Bitnet gateway, which has the unfortunate habit of converting ‘~’ to ‘^’ and ‘^’ to ‘%’. Since it leaves ‘%’ itself unaffected, this makes recovery of the original file a non-trivial exercise. It sometimes also changes the brace characters ‘{ }’ into odd characters above 128; this is particularly embarrassing, of course, for .tex files!

For some years, many T_EX users, particularly those working in languages other than English, and thus familiar with character set encodings containing other than the basic ASCII set, have been agitating for T_EX to be able to handle input in their

mother tongues, using their own languages’ character sets. In 1989, Knuth [1] announced T_EX v.3, and implementors world-wide beavered away to bring each implementation up to date. T_EX v.3 now supports eight-bit character sets and so .tex source files are now effectively ‘binary’ files and will therefore suffer from the same exchange problems experienced with .dvi files.

All those authors who had previously been able to cooperate, despite being separated by hundreds or thousands of miles, might once again be forced to entrust floppy disks to the vagaries of the world’s postal systems (although one shouldn’t underestimate the bandwidth of the Royal [or other] Mail system).

Unless or until the various e-mail protocols, networks and software are converted to support uncorrupted transmission of characters codes ‘040..’176 and ‘241..’376, it will have to become the norm for .tex sources to be encoded for transmission by e-mail.

The Aston Archive

All three authors are volunteer assistants to Peter Abbott in running the world’s principal repository of T_EX-related material at Aston University [2] in Birmingham. The archive holds several hundred megabytes of text and binary files including

- program sources for T_EX, METAFONT, DVI drivers and many other utilities;
- binary executables for a variety of popular operating systems (e.g., Atari, Macintosh, MS-DOS, UNIX, VAX/VMS and VM/CMS);
- METAFONT sources for Computer Modern and other fonts;
- binary font files (mainly .tfm and .pk) for a number of different output devices;
- text, macro and style files.

The archive provides access to these files via the following services:

- NIFTP¹ from Janet hosts. Typically 300 megabytes of data are transferred every month;

¹ Network Independent File Transfer Protocol — in the UK, one does not perform the pseudo-login that Internet users are accustomed to using with the FTP protocol. Instead, one issues a ‘transfer request’ for a file to be sent to or from the remote machine — the transfer itself takes place asynchronously. One nice consequence is that such transfers can be queued for overnight execution, leaving daytime bandwidth free for e-mail and true remote interactive logins.

this would probably be much greater if we were not limited by the bandwidth of our 9600-baud connection to Janet.

- FTP from Internet hosts. At the time of writing, the Internet connection has been approved and should be available by the third quarter of 1991.
- Interactive browsing service via Janet PAD, including the facility to send files out using NIFTP (and later FTP).
- Interactive browsing service via dial-up modem lines, including the facility to download files using Kermit and similar protocols.
- An e-mail file server that typically sends 150 megabytes of data per month to sites all over the world (though predominantly to EARN/Bitnet sites).
- A magnetic-media distribution service via surface carriers. Copies of the entire archive have been sent to embryonic T_EX communities in Czechoslovakia, Hungary and Poland.

We have experienced many problems trying to support all of these file types, operating systems and access methods. The e-mail file server clearly needs a reliable method of encoding files if its many customers are not to be denied access to the non-text files in the archive.

Binary files such as .pk font files are stored in different ways to accommodate the requirements of the different operating systems supported. Currently we maintain multiple font directory trees for the Macintosh, MS-DOS, UNIX and VAX/VMS, with all the attendant problems of synchronization, disk space and archivists’ time. We need a single storage format that allows export to all of our supported operating systems.

Specification for a Coding Scheme

In mid-1990, the archivists came to the conclusion that a universal encoding scheme was required to accommodate the many different kinds of file and file organizations that needed to be supported by the archive.

Niel Kempson formulated the first draft of this specification in mid-1990; the requirements of the encoding scheme may be summarized as follows.

Preserving File Structure. It is insufficient, especially for an archive holding binary files for a variety of machine types, merely to encode data simply as a stream of bytes:

- Virtually all operating systems² make a distinction between binary and text files, so the coding system should recognize and maintain this distinction.
- UNIX and most PC-based operating systems treat files as streams of bytes with no further structure imposed. On the other hand, certain widely-used operating systems (e.g., VAX/VMS and VM/CMS) have record-oriented file systems where different types of file are stored in a format appropriate to the type of file.³

For these operating systems, we consider it essential that the encoding scheme identify, preserve and record the most commonly used file organizations. The decoding program should be able to use this information to create the output file using the organization appropriate to the operating system in use. If the information is of no consequence to the receiving system, the default file structure (if any) should be created. If the encoding system does not have structure in its files, the receiving system may provide suitable defaults automatically. In all cases, the programs should permit the user to override or supplement file structure information.

- Whenever possible, these details of structure should be determined automatically by the encoding program; at the very least, an indication of whether the file is text or binary shall be provided (even under an operating system such as UNIX that need make no such distinction for its own use), to allow decoding to an appropriate file organization on those systems that *do* make such a distinction.

Coding Scheme. Whatever method is used for ensuring that encoded data can be e-mailed:

- It should be possible to specify the coding table to be used to encode the data. The coding table used should be recorded with each part of the encoded data.
- If a recorded coding table is found while decoding, it should be used to construct an appropriate decoding table. Simple one-to-one character corruptions should be corrected as long as only one of the input characters is mapped to any one output character.
- The recommended encoding uses only the following characters:

² UNIX is a notable exception to this rule.

³ It is argued that the increase in efficiency more than offsets the increase in complexity.

```
+--0123456789  
abcdefghijklmnopqrstuvwxy  
z  
ABCDEFGHIJKLMNOPQRSTUVWXYZ
```

Such an encoding has been shown to pass successfully through all the gateways that are known to corrupt characters.

Integrity of Encoded Data. We want to ensure that the *whole* encoded file passes through the e-mail network.

- Encoded lines should be prefixed by an appropriate character string to distinguish them from unwanted lines, such as mail headers and trailers. Whilst not essential, this feature does assist the decoding program in ignoring these spurious data.
- Lines should not end with whitespace characters, as some mailers and operating systems strip off trailing whitespace.
- The encoding program should calculate input file parameters, such as the number of bytes and CRC (cyclic redundancy check), and record them at the end of the encoded data.

The decoding program should calculate the same parameters from the decoded data and compare the values obtained from those recorded at the end of the encoded data.

Making Files Mailable. A mechanism is needed to overcome some gateways' refusal to handle large files.

- The encoding program should be able to split the encoded output into parts, each no larger than a maximum specified size. Splitting the output into smaller parts is useful if the encoded data is to be transmitted using electronic mail or over unreliable network links that do not stay up long enough to transmit a large file. The recommended default maximum part size is 30kBytes.
- The decoding program should be able to decode a multi-part encoded file very flexibly. It should *not* be necessary to:
 1. strip out mail headers and trailers,
 2. combine all of the parts into one file in the correct order, and
 3. process each part of the encoded data as a separate file.

Miscellaneous. Further considerations include:

- Support for character sets other than ASCII is essential if the encoding scheme is to be useful to IBM hosts. The encoding program should label the character set used by the encoded data,

and both encoder and decoder should enable the conversion between the local character set and another character set. For example, a user on an EBCDIC host should be able to encode text files for transmission to another EBCDIC host, or to convert them to ASCII before encoding and transmission to an ASCII host. Similarly, that user should be able to decode text files from ASCII and EBCDIC machines, creating EBCDIC output files.

- Where possible, the original file’s timestamp should be encoded and used by the decoding program when recreating the file; this will permit archives to retain the originator’s time of creation for files, and thus permit the users (not to mention the archivists) to identify more clearly when a new version of a file has been made available.
- The encoding and decoding schemes should be able to read and write files compatible with one or more of the well-established coding schemes.
- The source code for the programs should be freely available. It should also be portable and usable with as many computers, operating systems and compilers as possible.

The Search Commences

Naturally, the first step was to examine the existing coding schemes in comparison with the above ideal specification. Such schemes fell into two broad classes: *portable schemes*, which were intended to permit the encoding of files on any computer architecture into a form that could be transmitted electronically, and decoded on the same or a different architecture; and *platform-specific schemes*, which provided rather better support for transferring files between two computers using the same architecture and operating system.

Portable Coding Schemes. The most commonly used coding schemes supported by a variety of platforms are:

- boo
- UUcode
- XXcode

Most implementations of these schemes known to the authors are designed for use with stream file systems. These programs have no means of recording, let alone preserving, record structure and are thus unsuitable for our purposes. This is not surprising since UUcode and its mutation, XXcode, were developed specifically for exchanging files between UNIX systems. In fairness to these schemes, they are well

suitable to the transmission of text files and certain unstructured binary files.

Standard UUcode encodes files using characters ‘.’ ‘_’ of ASCII. This can result in one or more spaces appearing at the ends of lines; some mailers decide that this is information not worth transmitting, with consequent inability to reconstruct the original file.

Files containing characters such as ‘:’ are often irreversibly corrupted by mail gateways; this problem led to the development of XXcode, which uses a rather more robust character set, namely:

```
+-01234567890
abcdefghijklmnopqrstuvwxy
ABCDEFGHIJKLMNQRSTUWXYZ
```

The encoding table used is recorded with the encoded data to allow the detection of character corruptions, and the correction of reversible character transpositions. Whilst superficially a step forward, XXcode offered little more than most existing versions of UUcode, which already supported coding tables. Its major contribution was in formalizing the encoding table, and in particular its default table was proof against all the known gateway-induced corruptions.

Platform-Specific Coding Schemes. Encoding schemes have been developed to support transfer of files possessing some structure that therefore cannot be reconstructed correctly when encoded by the portable schemes. When the encoding and decoding programs of such a platform-specific scheme are each used on the same computer and operating system type, files may be encoded and transmitted with a great deal of confidence that the decoded file will reproduce the original’s structure and attributes in their entirety.

Examples of such programs are TELCODE and MFTU for VAX/VMS, NETDATA for IBM mainframes, and Stuffit and MacBinary for the Macintosh. But these programs have the major disadvantage that they have each been implemented *only* on the single architecture for which they were designed; thus the only two of these schemes that could be used on the VAX/VMS-based Aston Archive would be of minimal interest elsewhere!

The Archive’s content is in some respects artificially inflated by the presence of .hqx files for Macintoshes, .boo for MS-DOS, etc., which have to be held in pre-encoded form for transfer by those requiring them.

VVcode is Born

Realizing that none of the existing portable schemes were close enough to our ideal, an early version of our specification was circulated on various mailing lists by Niel Kempson towards the end of 1990. When the anticipated 'nil return' was all that resulted, Brian Hamilton Kelly went ahead and created a rudimentary `VVencode` by modifying an existing VAX-PASCAL implementation of `uuencode`. After generating the companion `VVdecode`, he then re-implemented the programs in Turbo C under the MS-DOS operating system on the IBM-PC, and thereby was able to prove that the new scheme was both viable and sufficient.

A Production VVcode. Following the minor feasibility study, Niel Kempson re-engineered the pair of programs from scratch (adding certain features of the evolving specification), paying particular attention to making the code⁴ portable across a wide variety of operating systems. Particular care was taken to avoid the use of supposedly standard C functions that experience had shown behaved differently under individual manufacturer's implementations, or were even non-existent in some. Therefore, the code may sometimes appear to be performing certain operations in a very long-winded way; it's very easy to look at it and say, "Why didn't the author use the ... function, which does this much more efficiently?" But this function may not even exist under another implementation of C, or it may behave in a subtly different manner.

The core functions of `VVcode` are implemented as a collection of routines written in as portable a fashion as possible, with a separate module of a few routines that are operating-system specific.⁵ Porting `VVcode` to a new platform should require only that this latter module be re-implemented, in most cases by adapting an existing one.

`VVcode` implements all of the features listed in the specification, apart from the ability to generate `UUcode`- and `XXcode`-compatible files. However, the decoding program is backwards compatible and can decode files generated by `UUcode` and `XXcode`.

Arguments against VVcode. When the advent of the `VVcode` system was first aired in the various electronic digests, some heated debate followed, along the lines that a new encoding scheme was unnecessary, since `UUcode`/`XXcode` sufficed for them.

⁴ That written by BHK was, in Niel's words, "PASCAL written as C"!

⁵ Such as file I/O, timestamping, command-line or other interface, etc.

However, all these correspondents were UNIX users who had interpreted the 'VV' as meaning 'VAX-to-VAX' (by analogy with 'uu'⁶), and thus felt that such a scheme should be private to VAXen. The authors' response is that the encoding scheme was intended to support the needs of archives like Astons, and as such, must provide:

1. an automated tool (it would be somewhat difficult to expect our users to be able to tell the encoder what sort of file structure it is handling, when this concept is entirely alien to many of them);
2. facilities to encode binaries for many operating systems;
3. mail server features, such as splitting of large files; and
4. operation across the widest possible combination of platforms.

The overhead of using the `VVcode` system is at most a couple of hundred bytes over using `UUcode`, and the extra functionality and *universality* with respect to `UUcode` or `XXcode` thereby comes almost for free.

Availability of VVcode

At present, the `VVcode` system is only available in C, but it has been shown to run successfully on the following combinations of hardware, operating system, and compiler:

Unix

- DEC Mips; Ultrix (BSD 4.2); native C
- HP9000; HPUNIX 6.5; native C and GNU C
- IBM RS-6000 (BSD 4.3); native C
- ICL DRS6000 (SPARC); System V (Rel 4); AT&T C
- Masscomp 5600; native C
- MIPS M/2000 (MIPS R3000); RiscOS 4.51; native C
- Sun; SUNOS 3.x and 4.0.3; native C and GNU C
- Sun Sparcstation 1; SUNOS 4.0.3; native C and GNU C

VAX/VMS

- All VAXen; VMS 5.2-5.4-1; VAX/C v3.0-v3.1-51 and GNU C

MS-DOS

- IBM PS/2, PC (and clones); MS-DOS 3.3, 4.01; Borland Turbo C 1.5, 2.0 and Turbo C++ 1.0

⁶ 'V' was chosen simply because it followed 'U'; at one time, we had seriously considered calling it YAFES—Yet Another File Encoding Scheme!

- IBM PS/2, PC (and clones); MS-DOS 3.3, 4.01; Microsoft C 5.1 and 6.0

VM/CMS

- VM/CMS; Whitesmith C compiler v1.0 (This implementation was ported by Rainer Schöpf; basing it upon the UNIX implementation, this took him about one day.)

Macintosh

- At the time of writing (May 1991), John Rawnsley of the University of Warwick had commenced development of a Macintosh port. This will encode the resource and data forks in a manner that will permit the former to be ignored by non-Macintosh systems.

Who’s Going to Use VVcode?

Obviously, since the whole concept was invented by the archivists at Aston, the Aston Archive will use VVcode when honouring e-mail requests, and the programs will also be available to browsers calling from sites without a binary NIFTP capability.

Rainer Schöpf has indicated that he will support VVcode on the Heidelberg server, as has George Greenwade at Sam Houston State University in Texas. Nelson Beebe intends to provide it as part of the TUGlib archive at Utah.

Naturally, all of these archives will also provide the sources of the programs, and will, wherever possible, provide complete distribution kits for transfer by (NI)FTP; these kits will include “load-and-go” executables for at least MS-DOS, UNIX, VAX/VMS and VM/CMS. The MS-DOS kit will be included on all physical distributions of T_EX for the PC from Aston.

References

- [1] Knuth, Donald E. “The New Versions of T_EX and METAFONT.” *TUGboat* 10#3, pages 325 – 328, 1989.
- [2] Abbott, Peter. “The UKT_EX Archive at the University of Aston.” *TUGboat* 10#4, pages 675 – 680, 1989.
- [3] Abbott, Peter. “A UK-Based T_EX Mail Archive Server.” *TUGboat* 9#3, pages 263 – 264, 1988.