

# Building Unicode Scanners with GPLEX

(Version 1.1.0 March 2009)

John Gough QUT

March 3, 2009

## *Documentation Map*

This paper documents the unicode-specific features of the *gplex* scanner generator. The complete documentation for *gplex* consists of the following files —

- \* *Gplex.pdf* – the main documentation file
- \* *Gplex-Input.pdf* – description of the input language of *gplex*
- \* *Gplex-Unicode.pdf* – documentation of the unicode-specific features of *gplex*, (this file)
- \* *Gplex-Changlog.pdf* – change log for *gplex*

## 1 Overview

Gardens Point *LEX* (*gplex*) is a scanner generator which accepts a “*LEX*-like” specification, and produces a *C#* output file. The scanners produced by *gplex* can operate in two modes —

- \* *Byte Mode*, in which patterns of seven or eight-bit bytes are specified, and the input source is read byte-by-byte. This mode corresponds to the traditional semantics of *LEX*-like scanner generators.
- \* *Unicode Mode*. In this mode the patterns are specified as regular expressions over the unicode alphabet. The generated scanner matches sequences of code-points. Traditional *LEX* has no equivalent semantics.

The choice between byte-mode and unicode-mode is made at scanner generation time, either by a command-line option to *gplex*, or an option marker in the specification file.

For unicode mode scanners, the input to the generated scanner must be decoded according to some known encoding scheme. This choice is made at scanner-runtime. Unicode text files with a valid unicode prefix (sometimes called a *Byte-Order-Mark*, “*BOM*”) are decoded according to the scheme specified by the prefix. Files without a prefix are interpreted according to a “*fallback code page*” option. This option may be specified at scanner generation time. The scanner infrastructure also provides methods

to allow scanner applications to override the default at scanner runtime, or even to defer choice until after scanning the entire file.

## 1.1 Gplex Options for Unicode Scanners

The following options of *gplex* are relevant to the unicode features of the tool.

### ***/codePage:Number***

In the event that an input file does not have a unicode prefix, the scanner will map the bytes of the input file according to the code page with the specified number. If there is no such code page an exception is thrown and processing terminates.

### ***/codePage:Name***

In the event that an input file does not have a unicode prefix, the scanner will map the bytes of the input file according to the code page with the specified name. If there is no such code page an exception is thrown and processing terminates.

### ***/codePage:default***

In the event that an input file does not have a unicode prefix, the scanner will map the bytes of the input file according to the default code page of the host machine. This option is the default for unicode scanners.

### ***/codePage:guess***

In the event that an input file does not have a unicode prefix, the scanner will rapidly scan the file to see if it contains any byte sequences that suggest that the file is either *utf-8* or that it uses some kind of single-byte code page. On the basis of this scan result the scanner will use either the default code page on the host machine, or interpret the input as a *utf-8* file. See Section 2.5 for more detail.

### ***/codePage:raw***

In the event that an input file does not have a unicode prefix, the scanner will use the uninterpreted bytes of the input file. In effect, only codepoints from 0 to u+00ff will be delivered to the scanner.

### ***/unicode***

By default *gplex* generates byte-mode scanners that use 8-bit characters, and read input files byte-by-byte. This option allows for unicode-capable scanners to be created. Using this option implicitly uses character classes.

### ***/noUnicode***

This negated form of the */unicode* option is the default for *gplex*.

**/utf8default**

This option is deprecated. It will continue to be supported in version 1.0. However, the same effect can be obtained by using `"/codePage:utf-8"`.

**/noUtf8default**

This option is deprecated. It will continue to be supported in version 1.0. However, the same effect can be obtained by using `"/codePage:raw"`.

## 1.2 Unicode Options for Byte-Mode Scanners

Most of the unicode options for *gplex* have no effect when a byte-mode scanner is being generated. However, the code page options have a special rôle in the special case of character set predicates.

The available character set predicates in *gplex* are those supplied by the *.NET* base class libraries. These predicates are specified over the unicode character set. On a machine with that uses a single-byte code page *gplex* must know what that code page is, in order to correctly construct character sets such as `[":ISpunctuation:"]`.

The available options are —

**/codePage:Number**

If a character set predicate is used, the set will include all the byte values which correspond in the code page mapping to unicode characters for which the predicate is true. The nominated code page must have the single-byte property.

**/codePage:Name**

If a character set predicate is used, the set will include all the byte values which correspond in the code page mapping to unicode characters for which the predicate is true. The nominated code page must have the single-byte property.

**/codePage:default**

If a character set predicate is used, the set will include all the byte values which correspond to unicode characters for which the predicate is true. In this case the mapping from byte values to unicode characters is performed according to the default code page of the *gplex* host machine. The default code page must have the single-byte property.

**/codePage:raw**

If a character set predicate is used, the set will include all the byte values which numerically correspond to unicode codepoints for which the predicate is true.

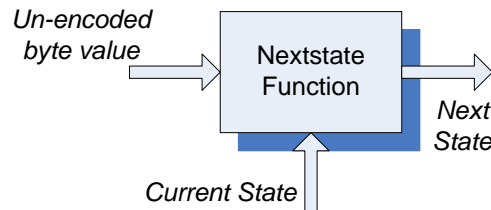
**Caution**

Character set predicates should be used with caution in byte-mode scanners. The potential issue is that the byte-mode character sets are computed at scanner generation time. Thus, unlike the case of unicode scanners, the code page of the scanner host machine must be known at scanner generation time rather than at scanner runtime (see also section 2.2).

## 2 Specifying Scanners

The scanning engine that *gplex* produces is a finite state automaton (FSA)<sup>1</sup> This FSA deals with codepoints from either the *ASCII* or *Unicode* alphabet. Byte-mode scanners have the conceptual form shown in Figure 1. The un-encoded byte values of the input

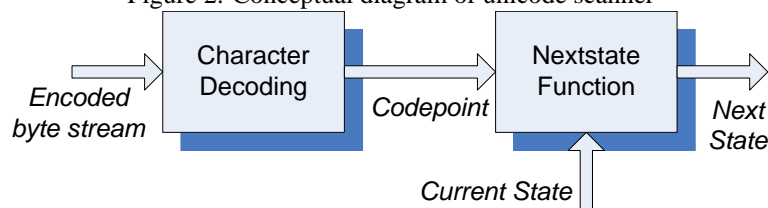
Figure 1: Conceptual diagram of byte-mode scanner



are used by the “next-state” function to compute the next state of the automaton.

In the unicode case the sequence of input values may come from a string of *System.Char* values, or from a file. Unicode codepoints need 21-bits to encode, so some interpretation of the input is required for either input form. The conceptual form of the scanner is shown in Figure 2 for file input. The corresponding diagram for *string* input

Figure 2: Conceptual diagram of unicode scanner



differs only in that the input is a sequence of *System.Char* values rather than bytes.

For *gplex* version 1.0 the scanner that *gplex* uses to read its own input (the “\*.lex” file) operates in byte-mode. Nevertheless, the input byte-mode text may specify either a byte-mode scanner as *gplex-output*, or a unicode-mode scanner as output.

Because of the choice of byte-mode for *gplex* input, literal characters in specifications denote precisely the codepoint that represents that character in the input file. Characters that cannot denote themselves in character literals must be specified by character escapes of various kinds.

In this section we consider the way in which byte-mode scanners and unicode scanners respectively are specified while complying with this constraint. Issues of portability of specifications and generated scanners across globalization boundaries are also discussed.

<sup>1</sup>(Note for the picky reader) Well, the scanner is *usually* an FSA. However, the use of the “/stack” option allows state information to be stacked so that in practice such *gplex*-generated recognizers can have the power of a push-down automaton.

## 2.1 Byte Mode Scanners

In byte-mode scanners, the only valid codepoints are in the range from ‘\0’ to ‘\xff’. When *gplex* input specifies a byte-mode scanner, character literals in regular expression patterns may be: literals such as ‘a’, one of the traditional control code escapes such as ‘\0’ or ‘\n’, or any other of the allowed numeric escapes.

The allowed numeric escapes are octal escapes ‘\ddd’, where the *d* are octal digits; hexadecimal escapes ‘\xhh’, where the *h* are hexadecimal digits; unicode escapes ‘\uhhhh’ and ‘\Uhhhhhhhh’, where the *h* are hexadecimal digits. If the specification is for a byte-mode scanner the numerical value of any character literal must be less than 256, or an error occurs.

It is important to see that even for byte-mode scanners, these choices lead to certain kinds of portability issues across cultures. Let us examine an example.

Suppose that a specification file is being prepared with an editing system that uses the Western European (Windows) code page 1252. In this case the user can enter a literal character ‘ß’, the *sharp s* character. This character will be represented by a byte 0xdf in the specification file. The byte-mode scanner which is generated will treat any 0xdf byte as corresponding to this character. To be perfectly clear: when the specification is viewed in an editor it may display a *sharp s* but *gplex* neither knows nor cares about how characters are displayed on the screen. When *gplex* reads its input it will find a 0xdf byte, and will interpret it as meaning “a byte with value 0xdf”.

Suppose now that the same specification is viewed on a machine which uses the Greek (Windows) code page 1253. In this case the same character literal will be displayed as the character *í*, *small letter iota with tonos*. Nevertheless, the scanner that *gplex* generates on the second machine will be identical to the scanner generated on the first machine.

Thus the choice of a byte-mode scanner for *gplex*-input achieves portability in the sense that any specification that does not use character predicates will generate a precisely identical scanner on every host machine. However, it is unclear whether, in general, the *meaning* of the patterns will be preserved across such boundaries.

In summary, byte-mode scanners handle the full 8-bit character set, but different code pages may ascribe different meanings to character literals for the upper 128 characters. Byte-mode scanners are inherently non-portable across cultures.

## 2.2 Character Predicates in Byte-Mode Scanners

Scanner specifications may use character set literals in the familiar form, the archetypical example of which is “[a-zA-Z]”. In *gplex* character set definitions may also use character predicates, such as “[[:IsLetter:]]”. In traditional *LEX*, the names of the character predicates are those available in “libc”. In *gplex* the available predicates are from the *.NET* base class library, and apply to unicode codepoints.

Consider the following example: a byte-mode specification declares a character set

```
PunctuationChars [[:IsPunctuation:]]
```

Now, the base class library function allows us to easily generate a set of *unicode* codepoints *p* such that the static predicate

```
Char.IsPunctuation(p);
```

returns true. Sadly, this is not quite what we need for a byte-mode scanner. Recall that byte-mode scanners operate on uninterpreted byte-values, as shown in figure 1. What we need is a set of byte-values *v* such that

```
Char.IsPunctuation(Map(v));
```

returns true, for the mapping *Map* defined by some code page.

For example, in the Western European (Windows) character set the ellipsis character ‘...’ is byte 0x85. The ellipsis is a perfectly good punctuation character, however

```
Char.IsPunctuation((char)0x85);
```

is false! The problem is that the ellipsis character is unicode codepoint u+2026, while unicode codepoint u+0085 is the “newline” control character *NEL*. All of the characters of the iso-8859 encodings that occupy the byte-values from 0x80 to 0x9f correspond to unicode characters from elsewhere in the space.

The character set “[ :IsLetter: ]” provides another example. For a byte-mode scanner using the Western European code page 1252, this set will contain 126 members. The same set has only 123 members in code page 1253. In the uninterpreted, raw case the set has only 121 members.

Nevertheless, it is permissible to generate character sets using character predicates in the byte-mode case. When this is done, the user may specify the code page that maps between the byte-values that the generated scanner reads, and the unicode codepoints to which they correspond.

If no code page is specified, the mapping is taken from the default code page of the *machine on which gplex is running*. This poses no problem if the machine on which the generated scanner will run has the same culture settings as the generating machine, or if the code page of the scanner host is known with certainty at scanner generation time. Other cases may lack portability.

### 2.3 Unicode Mode Scanners

The unicode standard ascribes unique 21-bit *codepoints* for every defined character<sup>2</sup>. Thus, if we want to recognize *both* the ‘ß’ character *and* the ‘ı’ character then we must use a unicode scanner. In unicode ß has codepoint u+00df, while ı has codepoint u+03af.

In unicode-mode scanners, the valid codepoints are in the range from u+0000 to u+10ffff. As was the case for byte-mode, character literals in the specification file may be literals such as ‘a’, one of the traditional control code escapes such as ‘\0’, or ‘\n’, or any other of the allowed numeric escapes.

The allowed numeric escapes are just as for the byte-mode case: octal escapes ‘\ddd’, where the *d* are octal digits; hexadecimal escapes ‘\xhh’, where the *h* are hexadecimal digits; unicode escapes ‘\uhhhh’ and ‘\Uhhhhhhhh’, where the *h* are hexadecimal digits. However, in this case the unicode escapes may evaluate to a codepoint up to the limit of 0x10ffff.

Since unicode scanners deal with unicode codepoints, it is best practise to always use unicode escapes to denote characters beyond the (7-bit) *ASCII* boundary. Thus our two example characters should be denoted ‘\u00df’ and ‘\u03af’ respectively.

#### Reading Scanner Input

The automata of unicode scanners deal only with unicode codepoints. Thus the scanners that *gplex* produces must generate the functionality inside the left-hand box in

<sup>2</sup>This is not the same as saying that every character has an unambiguous meaning. For example, in the *CJK compatibility* region of unicode ideograms with different meanings in Chinese, Japanese and Korean may share the same codepoint provided they share the same graphical representation.

figure 2. This *Character Decoding* function maps the bytes of the input file (or the characters of a string) into the codepoints that the scanner automaton consumes.

In the best of all worlds, the problem is simple. If the scanner’s input file is encoded using “little-endian” utf-16 our two example characters will each take two bytes. The  $\beta$  character will be denoted by two bytes  $\{0xdf, 0x00\}$ , while the  $\acute{i}$  character will be denoted by the two bytes  $\{0xaf, 0x03\}$ .

If the scanner’s input file is encoded using utf-8 our two example characters will again take two bytes each. The  $\beta$  character will be denoted by two bytes  $\{0xc3, 0x9f\}$ , while the  $\acute{i}$  character will be denoted by the two bytes  $\{0xce, 0x9f\}$ .

In both of these cases, the files should begin with a prefix which unambiguously indicates the format of the file. If a file is opened which does not start with a prefix then there is a problem.

Consider the case of a byte file prepared using either code page 1252 or code page 1253. Of course, such a file cannot contain both  $\beta$  and  $\acute{i}$  characters, since both of these are denoted by the same byte value 0xdf. The question is — if such a file is being scanned and a 0xdf byte is found — what codepoint should be delivered to the automaton<sup>3</sup>? Note that unlike the “utf-with-prefix” cases there is no certain way to know what code page a file was encoded with, and hence no certain way to know what decoding to use.

At the time that *gplex* generates a scanner, either a command line “/codePage:” option or a “%option” declaration in the specification may specify the fall-back code page that should be used if an input file has no unicode prefix. A common choice is “/codePage:default”, which treats files without prefix as 8-bit byte files encode according to the default code page on the host machine. This is a logical choice when the input files are prepared in the same culture as the scanner host machine. In fact, this is the fallback that *gplex* uses in the event that no code page option is specified.

The other common choice is “/codePage:utf-8”, which treats files without prefix as utf-8 files anyway.

If it is known for certain that input files will have been generated using a code page that is different to the host machine, then that known code page may be explicitly specified as the fallback. Note however, that this fallback will be applied to *every* file that the scanner encounters that does not have a prefix. In such cases it is more useful to allow the fallback to be specified to the scanner application on a file-by-file basis. How to do this is the subject of the next section.

What may we conclude from this discussion?

- \* Use unicode scanners for global portability whenever possible.
- \* Input files to unicode scanners should always be in one of the utf formats, whenever that is possible. Always place a prefix on such files.
- \* Consider using the default fallback to the host-machine code page unless it is known at scanner generation time that input files will originate from another culture.
- \* Applications that use *gplex* scanners should allow users to override the code page fallback when it is known that a particular input file originates from another culture.

---

<sup>3</sup>We have discussed only two possibilities here. Other code pages will give many additional meanings to the same 0xdf byte value.

## 2.4 Overriding the Codepage Fallback at Application Runtime

The fallback code page that is specified at scanner generation time is hardwired into the code of the generated scanner. However, an application that uses a *gplex* scanner may need to have its fallback code page changed for a particular input file when the encoding of that file is known.

Scanners generated by *gplex* implement a static method with the following signature —

```
public static int GetCodePage(string command);
```

This method takes a string argument, which is a code page-setting command from the calling application. If the command begins with the string “code page:” this prefix is removed, and the remaining string is converted to a code page index. The command may specify either a code page name or a number, or the special values “raw”, “default” or “guess”. Raw denotes no interpretation of the raw byte values, while “default” decodes according to the default code page of the host machine. Finally, “guess” attempts to determine the code page from the byte-patterns in the file. These semantics are the same as the */codePage:* option of *gplex*, which indeed invokes this same method.

The method is found in the buffer code of the generated scanner. If the */noEmbed-Buffers* option is in force the method will be in the class *QUT.GplexBuffers.CodePage-Handling*. For the default, embedded buffer case, the class *CodePage-Handling* is directly nested in the same namespace as the *Scanner* class.

There are two constructors for the scanner objects in each unicode scanner that *gplex* generates. One takes a stream object as its sole argument, while the other takes a stream object and a command string denoting the fallback code page. The second constructor passes the string argument to *GetCodePage*, and then sends the resulting integer to the appropriate call of *SetSource*<sup>4</sup>. Alternatively, the application may directly call *SetSource* itself, as shown below.

An application program that wishes to set the fallback code page of its scanner on a file-by-file basis should follow the example of the schema in Figure 3. If the application passes multiple input files to the same scanner instance, then an appropriate value for the fallback code page should be passed to each subsequent call of *SetSource* in the same way as shown in the figure.

## 2.5 Adaptively Setting the Codepage

There are occasions in which it is not possible to predict the code page of input files that do not have a unicode prefix. This is the case, for example, with programming language scanners that deal with input that has been generated by a variety of different text editing systems.

In such cases, if an input file has no prefix, a last resort is to scan the input file to see if it contains some byte value sequences that unambiguously indicate the code page. In principle the problem has no exact solution, so we may only hope to make a correct choice in the majority of cases.

Version 1.0.0 of *gplex* contains code to automate this decision process. In this first release the decision is only made between the *utf-8* code page and the default code page of the host machine. The option is activated either by using the command line option “*/codePage:guess*”, or by arranging for the host application to pass this command to the *GetCodePage* method.

<sup>4</sup>In the case of byte-mode scanners there is no fallback code page, so only the first constructor is generated.

Figure 3: Using the *GetCodePage* method

```

string codePageArg = null;
...
// Process the code page argument
if (arg.StartsWith("codepage:"))
    codePageArg = arg;
...
// Instantiate a scanner
FileStream file = new FileStream(...);
Scanner scnr = new Scanner();
if (codePageArg != null) {
    int cp = CodePageHandling.GetCodePage(codePageArg);
    scnr.SetSource(file, cp);
}
else // Use machine default code page, arg1 = 0
    scnr.SetSource(file, 0);
...

```

The code that implements the decision procedure scans the whole file. The “*guesser*” is a very lean example of a *gplex*-generated byte-mode *FSA*. This *FSA* searches for byte sequences that correspond to well-formed two, three and four-byte utf-8 codepoints. The automaton forms a weighted sum of such occurrences. The automaton also counts bytes with values greater than 128 (“*high-bytes*”) which do not form part of any legal utf-8 codepoint.

If a file has an encoding with the single-byte property there should be many more high-bytes than legal utf-8 sequences, since the probability of random high-bytes forming legal utf-8 sequences is very low. In this event the host machine code page is chosen.

Conversely if a file is encoded in utf-8 then there should be many multi-byte utf-8 patterns, and a zero high-byte count. In this event a utf-8 decoder is chosen for the scanner.

Note that it is possible to deliberately construct an input that tricks the *guesser* into a wrong call. Nevertheless, the statistical likelihood of this occurring without deliberation is very small.

There is also a processing cost involved in scanning the input file twice. However, the auxiliary scanner is very simple, so the extra processing time will generally be significantly less than the runtime of the final scanner.

### 3 Input Buffers

Whenever a scanner object is created, an input buffer holds the current input text. There are three concrete implementations of the abstract *ScanBuff* class. Two are used for string input, and the last for any kind of file input.

The *ScanBuff* class in Figure 4 is the abstract base class of the stream and string buffers of the generated scanners. The important public features of this class are the property that allows setting and querying of the buffer position, and the creation of

Figure 4: Features of the *ScanBuff* Class

```

public abstract class ScanBuff {
    ...
    public abstract int Pos { get; set; }
    public abstract int Read();
    public abstract string GetString(int begin, int end);
}

```

strings corresponding to all the text between given buffer positions. The *Pos* property returns the character index in the underlying input stream.

The method *Read* returns an integer corresponding to the ordinal value of the next character, and advances the input position by one or more input elements. *Read* returns  $-1$  for end of file.

New buffers are created by calling one of the *SetSource* methods of the scanner class. The signatures of these methods are shown in Figure 5.

Figure 5: Signatures of *SetSource* methods

```

// Create a string buffer and attach to the scanner. Start reading from offset ofst
public void SetSource(string source, int ofst);

// Create a line buffer from a list of strings, and attach to the scanner
public void SetSource(ICollection<string> source);

// Create a stream buffer for a byte-file, and attach to the scanner
public void SetSource(Stream src);

// Create a text buffer for an encoded file, with the specified encoding fallback
public void SetSource(Stream src, int fallbackCodepage);

```

### 3.1 String Input Buffers

There are two classes for string input: *StringBuff* which holds a single string of input, and *LineBuff* which holds a list of lines.

Scanners that accept string input should always be generated with the */unicode* option. This is because non-unicode scanners will throw an exception if they are passed a codepoint greater than 255. Unless it is possible to guarantee that no input string will contain such a character, the scanner will be unsafe.

#### The *StringBuff* Class

If the scanner is to receive its input as a single string, the user code passes the input to the first of the *SetSource* methods, together with a starting offset value —

```

public void SetSource(string s, int ofst);

```

This method will create a buffer object of the *StringBuff* type. Colorizing scanners for *Visual Studio* always use this method.

Buffers of this class consume either one or two characters for each call of *Read*, unless the end of string has been found, in which case the *EOF* value  $-1$  is returned. Two characters are consumed if they form a surrogate pair, and the caller receives a single codepoint which in this case will be greater than  $u+ffff$ . Calls directly or indirectly to *GetString* that contain surrogate pairs will leave the pair as two characters.

### The *LineBuff* Class

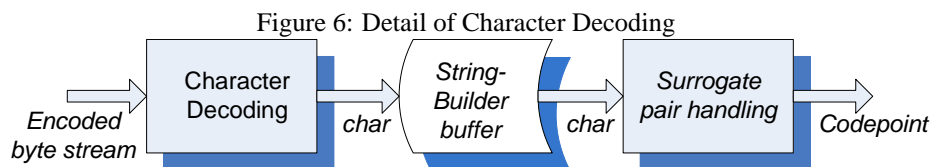
An alternative string interface uses a data structure that implements the *IList<string>* interface —

```
public void SetSource(IList<string> list);
```

This method will create a buffer object of the *LineBuff* type. It is assumed that each string in the list has been extracted by a method like *ReadLine* that will remove the end-of-line marker. When the end of each string is reached the buffer *Read* method will report a ‘\n’ character, for consistency with the other buffer classes. In the case that tokens extend over multiple strings in the list *buffer.GetString* will return a string with embedded end of line characters.

## 3.2 File Input Buffers

All file input to *gplex* is held in a buffer of the *BuildBuffer* class. In every case the sequence of bytes in the file is transformed into a sequence of code points supplied to the scanner by the scheme shown in Figure 6. The various generation-time options and



scanner-runtime code page settings simply modify the processing in the rectangular boxes of the figure.

The various possibilities are —

- \* *GetBuffer* is called with a single, *Stream* argument. This is the only possibility in the case of a byte-mode scanner. In this case the character decoding is trivial, with the bytes of the stream added unmodified to the buffer. In this case surrogate pairs cannot arise, so the right-hand box in the figure is empty also.
- \* *GetBuffer* is called with the *Stream* and a fallback code page argument. In all such cases the scanner checks if the stream begins with a valid *utf*-prefix. If a prefix is detected an appropriate *StreamReader* object is created, and transforms the bytes of the stream to the characters in the buffer. The buffer is filled with block-read operations rather than character by character.
- \* If the two-argument version of *GetBuffer* is called but no prefix is found then there are three special cases, and a general case. The general case is to create a *StreamReader* using whatever encoding is specified by the fallback code page.

The three special cases are the distinguished fallback values “raw”, “default” and “guess”. The raw value reverts to byte-mode decoding. The default value uses the default code page of the runtime host machine. Finally, the guess value causes the entire file to be scanned before a choice is made between *utf-8* and the default code page of the host machine. See also the discussion in section 2.5.