

Until recently, it was not necessary for software to deal with supplementary code points, those from U+10000 to U+10FFFF. With the assignment of over 40,000 supplementary characters in Unicode 3.1 and the definition of new national codepage standards that map to these new characters, it is important to modify BMP-only software to handle the full range of Unicode code points.

Typically, only a small percentage of code needs to be changed. This affects mostly low-level handling of 16-bit code units and data structures containing per-character data.

This presentation discusses the changes required to handle all of Unicode vs. just the BMP subset, concentrating on 16-bit Unicode -- the most common processing form. It describes techniques for finding the small percentage of code that typically needs to be changed, and shows how to modify such code. Detailed examples for Java and C/C++ use the many helper functions from ICU to illustrate practical solutions.

Presentation Goals



The main goal of this paper is to discuss how to migrate UCS-2 code: e.g. code that uses 16-bit Unicode but that does not handle the surrogate code points that represent supplementary characters.

The early versions of the Unicode standard defined up to 65536 code points. Unicode 2.0 extended the this to a million code points, but no important characters were assigned. The first important supplementary code points were assigned with Unicode version 3.1. These new code points are required for interworking with GB 18030, JIS X 0213 and Big5-HKSCS.

The first problem is identifying the places that need to be changed in order to support supplementary code points. Not all the places require changes. Once the problematic spots are found, the code modification should take place.

Luckily, there are many techniques that allow transformation of existing code base to handle supplementary code points.



UTF-16 uses two surrogate code points. A key feature is that the lead surrogate, trail surrogate and singleton code units never overlap in values. This means that a lot of code doesn't care about surrogates, as we will see in the examples.

Note that some of the applications assume that UTF-8 encodes code points with up to 3 bytes. This automatically prevents correct handling of supplementary code points.

Supplementary vs Surrogate

- Supplementary code **point**
 - Values in 10000..10FFFF
 - Corresponds to character
 - Rare in frequency
- Surrogate code unit
 - Value in D800..DFFF
 - Does *not* correspond to character by itself
 - Used *in pairs* to represent supplementaries *in UTF-16*

21st International Unicode Conference



Used by itself, types for characters may need to be changed to handle supplementary code points, which means either making them 32-bit wide (like int in Java) or handling the surrogate pairs (if staying with 16-bits wide units).

However, pointer types, such as UChar* should be given more consideration, as they could be treated as strings (as it is done in ICU).

Some types can be compiler/OS dependant, like wchar_t. In these cases, they need to be changed only if it is not possible to store a 32-bit value in them.

Deciding When to Change

- Varies by situation
- Operations with strings alone are rarely affected
- Code using characters might have to be changed
 - Depends on the types of characters
 - Depends on the type of code
 - Key Feature: Surrogates don't overlap!
- Use libraries with support for supplementaries
- Detailed examples below

21st International Unicode Conference

Dublin, Ireland, May 2002 6

Supplementary character can be ignored if the application is not processing text.

Explicit search for BMP and ASCII characters not affected

Most modern scripts (Latin, Cyrillic, Greek, Arabic, Hindi, Thai) not affected

Chinese, Japanese, historic scripts and certain Math symbols encoded in the supplementary space. If these are used, the code has to be changed.



Always index by code unit for performance, so that doesn't change.

Supplementaries are handled in certain cases, as we will see below.



The International Components for Unicode(ICU) is a C and C++ library that provides robust and full-featured Unicode support on a wide variety of platforms.

ICU is a collaborative, open-source development project jointly managed by a group of companies and individual volunteers throughout the world, using the Internet and the Web to communicate, plan, and develop the software and documentation.

The ICU project is licensed under the <u>X License</u> (see also the <u>x.org original</u>), which is <u>compatible with GPL</u> but non-viral.

Using ICU for Supplementaries

- Wide variety of utilities for UTF-16
- All internationalization services handle supplementaries
 - Character Conversion, Compression
 - Collation, String search, Normalization, Transliteration
 - Date, time, number, message format & parse
 - Locales, Resource Bundles
 - Properties, Char/Word/Line Breaks, Strings (C)
 - Supplementary Character Utilities

21st International Unicode Conference



Sun licenses ICU code for all the JVMs starting from Java 1.0



Most of the code in a program does not need to be changed because of supplementaries. In this case, for example, no supplementary characters need to be detected, so the code does not need to be changed.



Most string operations are safe, and String parameters can always handle supplementaries.

If two strings are both well formed, then their concatenation is.



Even substringing is ok, if the indices passed in are code point boundaries.

JAVA: API Problems

- You can't pass a supplementary character in function (1)
- You can't retrieve a supplementary from function (2)

1) void func1(char foo) {}

2) char func2() {}

21st International Unicode Conference



Ints are simpler for conversion, and can carry supplementaries. Changing to an int doesn't require call-site changes: if we call func('a'), it still works because Java widens.

However, often chars were originally a mistake, too narrow an interface. For example: having a currency symbol be a char is incorrect: you can't represent 'sFr' for Swiss Franc. Changing to a String is often a better approach, although String is much heavier weight than int, so it should be avoided in high-performance code. There are also pluses and minuses as far as your conversion goes.

Changing the API to have the parameter type be String will help reveal if any of the call-site code was not paying attention to surrogates when it should have. However, often this isn't needed. You may not have the freedom to change the API, either.

The alternative if you want String is to have an overload.



JAVA: Call Site Fixes



• Before 2. char x = myObject.func();

• After a) int x = myObject.func();

21st International Unicode Conference



A very common situation is where all the characters in a string are iterated. As a matter of fact, a majority of the code in ICU that required changes were in these situations, so it is worth taking a special look at them.

ICU4J: Looping Changes



Here is one style of change, that generally has the least impact on the body of the loop.

This change presumes that the function doSomething() has been changed (or overloaded) to accept supplementaries.

ICU4J: Tight Loops



For tight loops, sometimes other code is required.

Note: in this case the counter *i* is different in the body of the loop; it is in the middle of a supplementary character. Generally this is not important, but where it is, alternative styles need to be used.



We will go into more detail on these in the next slide.



- 1. Gets a 32-bit code point from an offset in string
- 2. Counts number of code units in a code point (could be 1 or 2)
- 3. Produces a string from a code point
- 4. Counts the number of code points in a string



Here is an example of converting indices.



The main functions for modifying a buffer of chars are here. There are parallel versions for plain char arrays.

Note: although it is not obvious, UTF16.setCharAt can change the length of the string. If a supplementary code point is replaced by a BMP code point the string will shrink. In opposite situation, it will grow.



The standard character properties are supplied. For ease of porting, these retain the same method names as in Java; the class name just has a U on the front.

What about Sun?

- Nothing in JDK 1.4
 - Except rendering; TextLayout does handle surrogates
- Expected support in next release
 - 2004?...
 - API?...
- In the meantime, ICU4J gives you the tools you need
- Code should co-exist even after Sun adds support

21st International Unicode Conference

ICU: C/C++

- Macros for UTF-16 encoding
- UnicodeString handles supplementaries
- UChar32 instead of UChar
- APIs enabled for supplementaries
- Very easy transition if the program is already using ICU4C

21st International Unicode Conference

Basic Data Types

- In C many types can hold a UTF-16 code unit
- Essentially 16-bit wide and unsigned
- ICU4C uses:
 - UTF-16 in UChar data type
 - UTF-32 in UChar32 data type

21st International Unicode Conference





Most of the code in a program does not need to be changed because of supplementaries. In this case, for example, no supplementary characters need to be detected, so the code does not need to be changed.

C++: Safe Code

No overlap with supplementaries



21st International Unicode Conference



Most string operations are safe, and String parameters can always handle supplementaries.

If two strings are both well formed, then their concatenation is.

The above example assumes that both s and t are NULL terminated that there is enough space in s to hold the concatenation result.



Most string operations are safe, and String parameters can always handle supplementaries.

If two strings are both well formed, then their concatenation is.



Supplementary characters cannot be passed as arguments to functions, nor can they be returned.



UChar32s are simpler for conversion, and can carry supplementaries. Changing to an UChar32 doesn't require call-site changes: if we call func1('a'), it still works because C/C++ widens.

However, often UChars were originally a mistake, too narrow an interface. For example: having a currency symbol be a char is incorrect: you can't represent 'sFr' for Swiss Franc. Changing to a UnicodeString is often a better approach, although UnicodeString is much heavier weight than UChar32, so it should be avoided in high-performance code. There are also pluses and minuses as far as your conversion goes.

Changing the API to have the parameter type be UnicodeString will help reveal if any of the call-site code was not paying attention to surrogates when it should have. However, often this isn't needed. You may not have the freedom to change the API, either.

The alternative if you want UnicodeString is to have an overload.



UChar32s are simpler for conversion, and can carry supplementaries. Changing to an UChar32 doesn't require call-site changes: if we call func1('a'), it still works because C/C++ widens.

However, often UChars were originally a mistake, too narrow an interface. For example: having a currency symbol be a char is incorrect: you can't represent 'sFr' for Swiss Franc. Changing to a UnicodeString is often a better approach, although UnicodeString is much heavier weight than UChar32, so it should be avoided in high-performance code. There are also pluses and minuses as far as your conversion goes.

Changing the API to have the parameter type be UnicodeString will help reveal if any of the call-site code was not paying attention to surrogates when it should have. However, often this isn't needed. You may not have the freedom to change the API, either.

The alternative if you want UnicodeString is to have an overload.

C/C++: Return Value Fixes

- Return values are trickier.
 - a) If you can change the API, then you can return a different value (String/int).
 - b) Otherwise, you have to have a variant name.
- Either way, you have to change the call sites.

21st International Unicode Conference





C/C++: Use Compiler

- Changes needed to address argument and return value problems easy to make, but error prone
- Compiler should be used to verify that all the changes are correct
- Investigate all the warnings!

21st International Unicode Conference



Function u_isalpha() expects a UChar32. UnicodeString::charAt function returns char. If a supplementary code point is in the string, it won't be picked up correctly.



This change presumes that the function doSomething() has been changed (or overloaded) to accept supplementaries.

In this loop *i* holds the offset of the code unit to be processed.



This change presumes that the function doSomething() has been changed (or overloaded) to accept supplementaries.

After UTF_NEXT_CHAR, i holds the offset to the next code unit to be processed, unlike the C++ version.



21st International Unicode Conference

ICU4C : Basic String Utilities

 Methods of UnicodeString class and macros defined in utf*.h.

1. cp = s.char32At(offset);

- 2. UTF_GET_CHAR(p, start, offset, length, cp)
- 3. cpLen = s.countChar32();
- 4. count = UTF_CHAR_LENGTH(cp);
- 5. s = cp;
- 6. UTF_APPEND_CHAR(p, offset, length, cp)
- 7. offset = s.indexOf(cp);
- 8. offset = s.indexOf(uchar);

21st International Unicode Conference



All the C++ methods have a C counterpart that works on an array of Unicode characters.



C macros are defined in utf.h, utf8.h, utf16.h, utf32.h. They allow for easy iterating over arrays containing one of these forms, as well as for converting between representation forms



ICU4C : String Modification



Character Iterator

- Convenience class, allows for elegant looping over strings
- Subclasses can be instantiated from:
 - UChar array
 - UnicodeString class
- Performance worse than previous examples
- Provides APIs parallel to UTF_* macros

21st International Unicode Conference



Instead of StringCharacterIterator, we could have used UCharCharacterIterator it(UCharArray, UCharArrayLen).

Very useful when a function takes a CharacterIterator reference as an argument.



C++ APIs exist, but are deprecated, as they are 1-1 wrappers around C APIs

Summary

- Because of the design of UTF-16, most code remains the same.
- Conversion is fairly straightforward...

With the right tools!

21st International Unicode Conference

Q & A

21st International Unicode Conference



To iterate over UTF-8 strings, one can use one of the macros that support different encoding forms. Do note, however, that the _UNSAFE functions are unsafe (both in regard to potential bounds breakage and malformation of the strings). These are to be used if and only if one is sure that the strings that are to be processed are well formed. Otherwise, go with _SAFE variants.



Converter is very fast and gives additional security. If dealing with external strings that are UTF-8 encoded, use a converter



When processing well formed data – provided by other APIs or trusted sources, you can use a faster converter – $u_strFromUTF8$, which avoids the overhead imposed by initializing and using converters.