



Internationalization Guide

Version 6.1, December 2003

IONA, IONA Technologies, the IONA logo, Orbix, Orbix/E, Orbacus, Artix, Orchestrator, Mobile Orchestrator, Enterprise Integrator, Adaptive Runtime Technology, Transparent Enterprise Deployment, and Total Business Integration are trademarks or registered trademarks of IONA Technologies PLC and/or its subsidiaries.

Java and J2EE are trademarks or registered trademarks of Sun Microsystems, Inc. in the United States and other countries.

CORBA is a trademark or registered trademark of the Object Management Group, Inc. in the United States and other countries. All other trademarks that appear herein are the property of their respective owners.

While the information in this publication is believed to be accurate, IONA Technologies PLC makes no warranty of any kind to this material including, but not limited to, the implied warranties of merchantability and fitness for a particular purpose. IONA Technologies PLC shall not be liable for errors contained herein, or for incidental or consequential damages in connection with the furnishing, performance or use of this material.

COPYRIGHT NOTICE

No part of this publication may be reproduced, stored in a retrieval system or transmitted, in any form or by any means, photocopying, recording or otherwise, without prior written consent of IONA Technologies PLC. No third party intellectual property right liability is assumed with respect to the use of the information contained herein. IONA Technologies PLC assumes no responsibility for errors or omissions contained in this book. This publication and features described herein are subject to change without notice.

Copyright © 2003 IONA Technologies PLC. All rights reserved.

All products or services mentioned in this manual are covered by the trademarks, service marks, or product names as designated by the companies who market those products.

M 3 1 3 3

Updated: 03-Dec-2003

Contents

List of Tables	v
Preface	vii
Chapter 1 Orbix Internationalization	1
Code Sets	2
Locales	5
Orbix Internationalization	8
Chapter 2 CORBA Internationalization	11
Overview	12
Supported Code Sets	14
Code Set Negotiation	16
Configuring the Code Set Plugin	18
Java Internationalization	22
C/C++ Internationalization	24
Chapter 3 Web Services Internationalization	29
WSDL Test Client	30
Chapter 4 Restrictions	33
Translations	34
Property Values	35
COMet	36
Index	37

CONTENTS

List of Tables

Table 1: IANA Charset Names	3
Table 2: Popular code sets supported by Orbix	14
Table 3: Code sets supported by the light weight plugin	18
Table 4: IDL to C++ character mappings	24
Table 5: Mime Encodings for WSDL Test Client	31

LIST OF TABLES

Preface

Audience

This guide is intended for Orbix programmers who develop products that might be internationalized. This guide also provides information to administrators of Orbix deployments where internationalization is required.

Related documentation

The Orbix documentation set includes the following related documentation:

- *CORBA Programmer's Guide* (C++ and Java)
- *CORBA Programmer's Reference* (C++ and Java)
- Administrator's Guide

The latest versions of all Orbix documentation can be found online at <http://www.ionas.com/support/docs>.

Additional resources

The [IONA knowledge base](#) contains helpful articles, written by IONA experts, about Artix and other products. You can access the knowledge base at the following location:

The [IONA update center](#) contains the latest releases and patches for IONA products:

If you need help with this or any other IONA products, contact IONA at support@ionas.com. Comments on IONA documentation can be sent to docs-support@ionas.com.

Typographical conventions

This guide uses the following typographical conventions:

<i>Constant width</i>	<p>Constant width (courier font) in normal text represents portions of code and literal names of items such as classes, functions, variables, and data structures. For example, text might refer to the <code>CORBA::Object</code> class.</p> <p>Constant width paragraphs represent code examples or information a system displays on the screen. For example:</p> <pre>#include <stdio.h></pre>
<i>Italic</i>	<p>Italic words in normal text represent <i>emphasis</i> and <i>new terms</i>.</p> <p>Italic words or characters in code and commands represent variable values you must supply, such as arguments to commands or path names for your particular system. For example:</p> <pre>% cd /users/<i>your_name</i></pre> <p>Note: Some command examples may use angle brackets to represent variable values you must supply. This is an older convention that is replaced with <i>italic</i> words or characters.</p>

Keying conventions

This guide may use the following keying conventions:

No prompt	When a command's format is the same for multiple platforms, a prompt is not used.
%	A percent sign represents the UNIX command shell prompt for a command that does not require root privileges.
#	A number sign represents the UNIX command shell prompt for a command that requires root privileges.
>	The notation > represents the DOS, Windows NT, Windows 95, or Windows 98 command prompt.
...	Horizontal or vertical ellipses in format and syntax descriptions indicate that material has been eliminated to simplify a discussion.
.	
.	
.	

- [] Brackets enclose optional items in format and syntax descriptions.
- { } Braces enclose a list from which you must choose an item in format and syntax descriptions.
- | A vertical bar separates items in a list of choices enclosed in { } (braces) in format and syntax descriptions.

PREFACE

Orbix

Internationalization

Orbix includes features that enable development and deployment of applications that manipulate user data encoded in characters beyond the traditional ASCII and ISO 8859-1 (Latin 1) code sets.

In this chapter

This chapter discusses the following topics:

Code Sets	page 2
Locales	page 5
Orbix Internationalization	page 8

Code Sets

Overview

A *coded character set*, or *code set* for short, is a mapping between integer values and characters they represent. The best known code set is ASCII, which defines 94 graphic characters and 34 control characters using the 7-bit integer range.

European languages

The 94 characters defined by the ASCII code set are sufficient for English, but they are not sufficient for European languages, such as French, Spanish, and German.

To remedy the situation, an 8-bit code set, ISO 8859-1, also known as Latin-1, was invented. The lower 7-bit portion is identical to ASCII. The extra characters in the upper 8-bit range cover those languages used widely in the Western European region.

Many other code sets are defined under ISO 8859 framework. These cover languages in other regions of Europe as well as Russian, Arabic and Hebrew. The most recent addition is ISO 8859-15, which is a revision of ISO 8859-1 and it adds the Euro currency symbol and other letters while removing less used characters. For further information about ISO-8859-x encoding, refer to the Web site [“The ISO 8859 Alphabet Soup”](#).

Ideograms

Asian countries that use ideograms in their writing systems needed more characters than they fit in an 8-bit integer. Therefore, they invented a double-byte code set, where a character is represented by a bit pattern of 2 bytes.

These languages also needed to mix the double-byte code set with ASCII in a single text file. So, *character encoding schema*, or simply *encodings*, was invented as a way to mix characters of multiple code sets.

Some of the popular encodings used in Japan include:

- Shift JIS
- Japanese EUC
- Japanese ISO 2022

Unicode

Unicode is a new code set that is gaining popularity. It aims to assign a unique number, or code point, to every character that exists (and even once existed) in all languages. To accomplish this, Unicode, which began as a double-byte code set, has been expanded into a quadruple-byte code set.

Unicode, in pure form, can be difficult to use within existing computer architectures, because many APIs are byte-oriented and assume that the byte value 0 means the end of the string.

For this reason, Unicode Transformation Format for 8-bit channel, or UTF-8, is frequently used. When browsers list “Unicode” in its encoding selection menu, they usually mean UTF-8, rather than the pure form of Unicode.

Visit [Unicode Inc.](#) for more information about Unicode and its variants.

Charset names

To address the need for computer networks to connect different types of computers that use different encodings, the Internet Assigned Number Authority, or IANA, has a registry of encodings at <http://www.iana.org/assignments/character-sets>.

IANA names are used by many Internet standards including MIME, HTML, and XML.

[Table 1](#) lists IANA names for some popular charsets.

Table 1: *IANA Charset Names*

IANA Name	Description
US-ASCII	7-bit ASCII for US English
ISO-8859-1	Western European languages
UTF-8	Byte oriented transformation of Unicode
UTF-16	Double-byte oriented transformation of 4-byte Unicode
Shift_JIS	Japanese DOS & Windows
EUC-JP	Japanese adaptation of generic EUC scheme, used in Unix
ISO-2022-JP	Japanese adaptation of generic ISO 2022 encoding scheme

Note: IANA names are case insensitive. For example, US-ASCII can be spelled as us-ascii or US-ascii.

CORBA Names

In CORBA code sets are identified by numerical values registered with the Open Group's registry, OSF Code Set Registry:

ftp://ftp.opengroup.org/pub/code_set_registry/code_set_registry1.2g.txt.

Java Names

Java has its own names for charsets. For example, ISO-8859-1 is named `ISO8859_1`, Shift_JIS is named `SJIS`, and UTF-8 is named `UTF8`.

Java is transitioning to IANA charset names, to be aligned with MIME. JDK 1.3 and above recognizes both names.

Note: This guide uses IANA charset names even for CORBA code sets.

Locales

Concept of locale

Most of modern operating systems are multilingual. Users can choose their language, and the operating system behaves according to the linguistic convention of the chosen language. For example, the user the number 1234.56 can be displayed as 1,234.56 to the English users and 1.234,56 to the German users.

However, language alone is not enough to determine a behavior, especially for languages that are used in many countries. For example, French speakers in Canada expect to see 1,234.56 while European French speakers expect to see 1.234,56.

The concept of *locale* addresses these issues. A locale combines charsets and display behavior for specific regions.

ISO standards

The International Standard Organization, ISO, defines two standards to specify locale. In general, ISO 639 specifies the language code and ISO 3166 specifies the country code.

These standards can be down loaded from the ISO web site at www.iso.org.

Operating system locales

Windows

On Windows, the **Regional and Language Option** control panel is used to select a locale. Locales are listed by language. For languages that are spoken in multiple countries, such as English, you must choose the setting for your region. Regions are listed next to the language in parenthesis. For example a Spanish speaker in Mexico would select **Spanish (Mexico)**.

UNIX

On Solaris and other POSIX conformant platforms, a set of environment variables, `LANG` and `LC_XXXX`, are used to select a language. Usually, only `LANG` is used to set locale behavior to the specified locale.

Typical locale values on Solaris are `en` (English, generic), `en_GB` (British English), or `en_GB.ISO8859-15` (British English using ISO-8859-15 encoding). The two-letter code, `en`, means English, and is taken from the language code standard from ISO, and `GB` means Great Britain and is taken

from another ISO standard on the country code. As the last example indicates, an OS locale name can also include an encoding name (ISO8859-15).

Note: The name of the encoding does not follow IANA charset registry.

When no encoding name appears as a component of the locale name, a default encoding is implied. For example, the locale `ja` on Solaris implies use of EUC-JP encoding, `eucJP` in the Solaris naming convention. Therefore, `ja` and `ja_JP.eucJP` are synonymous on Solaris.

All Unix-derived operating systems share similar locale semantics although the naming conventions vary widely.

Java locales

Java has its own locale mechanism. Its notational convention is similar to Solaris, except that there is no encoding specifier as Java's internal encoding is always UTF-16. Typical locales include:

- `en_US`
- `fr_FR`
- `de_DE`
- `zh_CN`
- `zh_TW`
- `ja_JP`
- `ko_KR`

Java's default locale is automatically inherited from the JVM's current operating system locale. [Example 1](#) shows how to use the Java Locale class to determine the locale and file encoding settings for your system.

Example 1: *Printing the Java locale and file encoding setting*

```
// Java
import java.util.Locale;

public class printLocale
{
    public static void main(String [])
    {
        Locale default_locale = Locale.getDefault();
        System.out.println("Default locale: " +
            default_locale.toString());
    }
}
```

Example 1: *Printing the Java locale and file encoding setting*

```
String file_encoding = System.getProperty("file.encoding");
System.out.println("File encoding: " + file_encoding);
}
}
```

For more information, see the JDK API document for the `Locale` class.

Language on the Internet

The Hyper Text Transfer Protocol, HTTP, has two headers, `Accept-Language` and `Content-Language`, that relate to locale. These two headers take a language value (or list of values) which has the form *language[-subtag]*.

According to the HTTP specification, *language* and *subtag* can be any string value. In practice, the ISO two-letter language code and country code are used. Therefore, the values in these two fields are almost the same as the Java locales.

Orbix Internationalization

Feature list

The following internationalization features are available in Orbix:

- CORBA Internationalization Features
This includes IDL `wchar` and `wstring` datatypes, code set negotiation, and extended code set support.
- Java 2 Enterprise Edition (J2EE) 1.3
This includes the servlet `response.setContentType()` and `request.setCharacterEncoding()` methods and JSP `pageEncoding` and `contentType` attributes in the `page` directive.
- Locale-to-encoding mapping enhancement to J2EE
An encoding can be associated with a locale via configuration. When a JSP or Servlet specifies a locale attribute for a response using the `response.setLocale()` method, the encoding associated with the locale is used in the `charset` component of the `Content-Type` header of the response.
- Configurable fall-back encoding of servlet for J2EE
A fall-back encoding for requests on a servlet or JSP can be configured using the URL mapping for the servlet. If the servlet does not associate an encoding with a request programmatically (for example, using `request.setContentType()`) then the fall-back encoding is used, if configured.
- IANA-charset-name to Java-converter-name mapping for J2EE
In a small number of cases it is necessary to decode HTTP request body data using a Java converter with a different name from the name of the IANA charset used in the request header. ASP provides a mechanism to map IANA charset names to Java converter names via configuration.
- Character encoding support in WSDL Test Client
WSDL Test Client has a menu to specify the encoding
- Large number of code sets/character encoding supported

163 built-in code sets are supported for CORBA. All character encodings (code sets) supported by the underlying JDK can be used by J2EE programs.

Enabling tools

With Orbix, you can enable applications from different locales or using different code sets to interoperate. However, Orbix does not provide tools to help you applications capable of working with multiple locales or code sets.

CORBA

Internationalization

Orbix lets you run applications in numerous locales.

In this chapter

This chapter discusses the following topics:

Overview	page 12
Supported Code Sets	page 14
Code Set Negotiation	page 16
Configuring the Code Set Plugin	page 18
Java Internationalization	page 22
C/C++ Internationalization	page 24

Overview

Wide characters

CORBA 2.1 introduced the datatypes `wchar` (*wide character*) and `wstring` (*wide string*). Wide characters allow a character to be stored in a fixed length datatype whether it is a one byte character, a double-byte character, or a quad-byte character. This makes programming for multiple languages easier.

Note: The actual size of a wide character varies by operating system. Typical sizes are two bytes or four bytes.

[Example 2](#) shows an IDL definition that uses wide datatypes.

Example 2: Sample IDL using wide datatypes

```
// IDL
interface WideEcho
{
    wstring echo(in wstring ws);
    wstring echoSingleWChar(in wchar wc);
};
```

Mixing wide and narrow data

The traditional `string` datatype, sometimes called a *narrow string*, can also represent a multibyte character string. A character in this form has a varying byte length, usually 1 to 3 bytes. The number of bytes depends on the code set in use.

The `char` datatype, however, cannot store a multibyte character because its size is limited to one byte. [Example 3](#) shows an IDL definition of an interface that mixes narrow and wide strings.

Example 3: Mixing narrow and wide strings

```
// IDL
interface WideEcho
{
    wstring echo(in wstring ws);
    wstring echoSingleWChar(in wchar wc);
    wstring echoNarrowString(in string ns);
};
```

Code set negotiation

Because CORBA is designed to work in a heterogeneous networking environment, the server's native code set might differ from the client's native code set. CORBA defines a mechanism for ensuring that both client and server can exchange meaningful data efficiently. This process is called *code set negotiation*.

Supported Code Sets

Popular code sets

Table 2 shows some of the code sets that Orbix supports.

Table 2: Popular code sets supported by Orbix

OSF code set name	OSF code set id	IANA charset	Java encoding
ISO 8859-1:1987	0x00010001	ISO-8859-1	ISO8859_1
UCS-2, Level 1	0x00010100	UCS-2	UTF-16
UCS-4, Level 1	0x00010104	UCS-4	UCS-4
X/Open UTF-8	0x05010001	UTF-8	UTF8
JIS eucJP	0x00030010	EUC-JP	EUC_JP
OSF Japanese SJIS-1	0x05000011	Shift_JIS	SJIS

Java CORBA

For Java CORBA, Orbix supports 163 code sets. It uses the Java native encoding converters listed for JDK 1.4 (<http://java.sun.com/j2se/1.4/docs/guide/intl/encoding.doc.html>) and for JDK 1.3 (<http://java.sun.com/j2se/1.3/docs/guide/intl/encoding.doc.html>).

Because Java does not use OSF code set IDs to name the encodings, an OSF code set ID must first be mapped to a Java encoding name.

Because the mapping between OSF code set IDs and IANA and Java naming schema is not one-to-one, Orbix maps the most popular code set among the code sets that are almost identical. For example, Orbix supports code set id 0x05000011 (OSF Japanese SJIS-1) but not 0x05020002 (JVC_SJIS). Please make sure to use the code set that Orbix supports.

WARNING: The mapping is subject to change without notice in future releases.

One of the most popular code sets used in Japan, `ISO-2022-JP`, is missing from OSF registry therefore Orbix does not support it. `ISO-2022-JP` is mainly used in e-mail, and it is rarely used in inter-process communication or in storage.

C/C++ CORBA

For C and C++ CORBA, Orbix supports 112 code sets. It uses the code set converters from ICU (<http://oss.software.ibm.com/icu>), an open-source project supported by IBM.

Custom code set plugins

The Java encoders and the ICU encoders built in to Orbix do not support all code sets in use today. For situations where conversion is needed for an unsupported code set, Orbix has a pluggable code set converter architecture which makes it possible to write and add a custom code set converter plugin.

IONA Professional Service can write a custom code set plugin for you. Contact your local IONA sales office for details.

Code Set Negotiation

Overview

Code set negotiation is the process by which two CORBA processes which use different *native code sets* determine which code set to use as a *transmission code set*. Occasionally, the process requires the selection of a *conversion code set* to transmit data between the two processes. The algorithm is defined in section 13.10.2.6 of the CORBA specification (<http://cgi.omg.org/docs/formal/02-12-06.pdf>).

Note: For CORBA programming in Java, you can specify codeset other than the true native codeset.

Native code set

A native code set (NCS) is a code set that a CORBA program speaks natively. For Java, this is UTF-8 (0x05010001) for `char` and `String`, and UTF-16(0x00010109) for `wchar` and `wstring`. For C and C++, this is the encoding that is set by `setlocale()`, which in turn depends on the `LANG` and `LC_XXXX` environment variables.

Conversion code set

A conversion code set (CCS) is an alternative code set that the application registers with the ORB. More than one CCS can be registered for each of the narrow and wide interfaces. CCS should be chosen so that the expected input data can be converted to and from the native code set without data loss. For example, Windows code page 1252 (0x100204e4) can be a conversion code set for ISO-8859-1 (0x00010001), assuming only the common characters between the two code sets are used in the data.

Each application has its own native code set and a set of conversion code sets for `char` and `string`. Each application also has a separate native code set and conversion code sets for `wchar` and `wstring`. The CCS for `wchar` and `wstring` can be same as or different from those for `char` and `string`.

Transmission code set

A transmission code set (TCS) is the code set agreed upon after the code set negotiation. The data on the wire uses this code set. It will be either the native code set, one of the conversion code sets, or UTF-8 for the narrow interface and UTF-16 for the wide interface.

Negotiation algorithm

Code set negotiation uses the following algorithm to determine which code set to use in transferring data between client and server:

1. If the client and server are using the same native code set, no translation is required.
2. If the client has a converter to the server's code set, the server's native code set is used as the transmission code set.
3. If the client does not have an appropriate converter and the server does have a converter to the client's code set, the client's native code set is used as the transmission code set.
4. If neither the client nor the server has an appropriate converter, the server ORB tries to find a conversion code set that both server and client can convert to and from without loss of data. The selected conversion code set is used as the transmission code set.
5. If no conversion code set can be found, the server ORB determines if using UTF-8 (narrow characters) or UTF-16 (wide characters) will allow communication between the client and server without loss of data. If UTF-8 or UTF-16 is acceptable, it is used as the transmission code set. If not, a `CODESET_INCOMPATIBLE` exception is raised.

Code set compatibility

The last steps involves a compatibility test, but the CORBA specification does not define when a code set is compatible with another. The compatibility test algorithm employed in Orbix is outlined below:

1. ISO 8859 Latin-*n* code sets are compatible.
2. UCS-2 (double-byte Unicode), UCS-4 (four-byte Unicode), and UTF-*x* are compatible.
3. All other code sets are not compatible with any other code sets.

This compatibility algorithm is subject to change without notice in future releases. Therefore, it is best to configure the code set variables as explicitly as possible in order to reduce dependency on the compatibility algorithm.

Configuring the Code Set Plugin

Configuration variables

In order for an ORB to transmit character data in a code set other than ISO 8859-1, the ORB must be configured properly. Four configuration variables control the code set plugin:

plugins.codeset.char.ncs specifies the native code sets used to represent narrow characters.

plugins.codeset.char.ccs specifies the list of conversion code sets supported for narrow characters.

plugins.codeset.wchar.ncs specifies the native code sets used to represent wide characters.

plugins.codeset.wchar.ccs specifies the list of conversion code sets supported for wide characters.

For more information on these variables, see the *Orbix Configuration Reference*.

Note: For CORBA programming in Java, you can specify a codeset other than the true native codeset. Please see [“Native code set” on page 23](#).

Light weight code set plugin

The standard code set plugin for C/C++ requires 8MB of memory. If memory is limited or if you are planning to use only the code sets listed in [Table 3](#), a light weight replacement, `it_basic_codeset`, is available.

Table 3: Code sets supported by the light weight plugin

UCS-2
UTF-4
UTF-16
ISO-8859-1
EBCDIC

To use this plugin instead of the standard one, modify the configuration with this setting

```
initial_references:IT_CodeSet_Registry:plugin= "basic_codeset" ;
```

Note: Changing the code set plugin in this way has no effect on Java.

Choice of conversion code sets

When choosing code sets to use as conversion code sets, three points should be considered:

- **Data compatibility**
Conversion code sets should be chosen to minimize the of loss or corruption of data. At minimum, the code sets chosen must cover all the characters that need to be handled by the application. Shift_JIS and EUC-JP are based on the same base character sets, except that EUC-JP includes an extra code set that is in rare use. If the application does not expect to handle the extra code set, these code sets can be treated as compatible. Usually, Shift_JIS and ISO-8859-1 are not considered compatible because many of the letters with diacritics in ISO-8859-1 do not exist in Shift_JIS. However, if the application is not expected to handle data that includes these incompatible characters, you can consider them compatible. Ultimately, it is up to the application designer to decide whether a code set is compatible with another.
- **Performance**
The choice of the transmission code set greatly affects the performance. In general, fixed-length code sets such as ISO-8859-x, EBCDIC, UCS-2 and UCS-4 can achieve better performance than the variable-length code sets such as UTF-8, EUC-JP and Shift_JIS.

Note: Although UTF-16 is a variable length code set, Orbix implements it as fixed-width code set because the characters beyond Basic Multilingual Plain (BMP) are not supported.

- **Compatibility with legacy CORBA services**
It is often the case that an application uses multiple CORBA servers. Some servers might not support code sets other than ISO-8859-1,

which is the only code set that is mandated to be supported by CORBA. Some servers might also not support `wchar` or `wstring`. In fact, most of the IONA services such as naming, locator, etc. belong to this category. For this reason, ISO-8859-1 should be included in `plugins:codeset:char:ccs`.

Example configurations

[Example 4](#) shows a basic configuration that works in a mixed Java/C++ environment. This configuration works in Latin-1 based locales.

Example 4: *Basic code set configuration*

```
plugins:codeset:char:ncs = "0x05010001"; # UTF-8
plugins:codeset:char:ccs = ["0x00010001"]; # ISO-8859-1;
plubins:codeset:wchar:ncs = "0x00010109"; # UTF-16
plugins:codeset:wchar:ccs = [];
```

[Example 5](#) shows a configuration for a heterogeneous system environment where some servers are written in Java and not configured to accept ISO-8859-1, and some servers are hosted on a mainframe. Because EBCDIC (IBM code page 037) does not include characters needed for European languages, this configuration can only be used for English.

Example 5: *Mixed environment configuration*

```
plugins:codeset:char:ncs = "0x00010001"; # ISO-8859-1
plugins:codeset:char:ccs = ["0x05010001", "0x10020025"]; # UTF-8, EBCDIC(IBM-037)
plubins:codeset:wchar:ncs = "0x00010001"; # ISO-8859-1
plugins:codeset:wchar:ccs = ["0x00010109"]; # UTF-16
```

For the language such as Japanese, where multiple code sets are used in a heterogeneous system environment, all of the used code sets should be put in the conversion code set list. [Example 6](#) shows a configuration for a Japanese system with servers running on both Windows and Solaris systems.

Example 6: *Japanese mixed environment configuration*

```
plugins:codeset:char:ncs = "0x05010001"; # UTF-8
plugins:codeset:char:ccs = ["0x00030010", "0x05000011", "0x00010001"]; # JIS eucJP, OSF SJIS,
    ISO-8859-1
plugins:codeset:wchar:ncs = "0x00010109"; # UTF-16
plugins:codeset:wchar:ccs = ["0x00030010", "0x05000011", "0x00010104", "0x00010100"]; # JIS
    eucJP, OSF SJIS, UCS-4 Level 1, UCS-2 Level 1
```

Logging

The code set plugin outputs informational event messages using the event subsystem `IT_CODESET`. To view these events make sure the configuration variable `event_log:filters` includes the entry `"IT_CODESET=*" or "IT_CODESET=INFO"`.

Default configuration

If any code set configuration variables are missing from the configuration, the default values shown in the *Orbix Configuration Reference* are used.

Java Internationalization

IDL-to-Java mapping

As specified by the OMG's IDL-to-Java mapping specification, the `string` and `wstring` IDL datatypes are mapped to the Java `String` class, and the `char` and `wchar` IDL datatypes are mapped to the Java `char` datatype.

Note: Although a Java `char` is two-bytes wide, any attempt to transmit a Java `char` that is bigger than `0xff` through the IDL `char` interface will throw the exception of `org.omg.CORBA.DATA_CONVERSION`.

The latest version of the CORBA *IDL-to-Java Language Mapping Specification* can be obtained at <http://www.omg.org/cgi-bin/doc?formal/02-08-05.pdf>.

Coding

Because Java treats both narrow and wide IDL datatypes alike, there is little that a CORBA developer using Java needs to consider when thinking about internationalization.

To use `WideEcho.idl`, shown in [Example 3 on page 12](#), apply `idlgen` using the following command:

```
idlgen java_poa_genie.tcl -jP WideEchoDemo -all WideEcho.idl
```

To implement the server, `WideEchoImpl.java` needs to have its operation defined. Each operation in the `WideEcho` class simply echoes the string that is passed into it. [Example 7](#) shows the implementation of the `WideEcho` operations.

Example 7: *WideEcho* Java server

```
// Java
public java.lang.String echo(java.lang.String ws)
throws org.omg.CORBA.SystemException
{
    return ws;
}
```

Example 7: *WideEcho Java server*

```

public java.lang.String echoSingleWChar(char wc)
throws org.omg.CORBA.SystemException
{
    // Returns a String consisting of just one char.
    char [] x = new char[1];
    x[0] = wc;
    return new String(x);
}
public java.lang.String echoNarrowString(java.lang.String ns)
throws org.omg.CORBA.SystemException
{
    return ns;
}

```

The Java implementation of the IDL operation `wstring echo(wstring)` is identical to the Java implementation of the IDL operation `wstring echoNarrowString(string)`. This is because Java makes no distinction between an IDL `string` and IDL `wstring`.

The Java implementation of the IDL operation `wstring echoSingleWChar(char)` is more complicated than the others only because it needs to convert a Java `char` to a Java `String`.

On the client side, the code is similar to a traditional Java CORBA program. IDL `wchar` and IDL `char` are both represented as Java `char`. IDL `wstring` and IDL `string` are represented as Java `String`.

Native code set

For Java-written CORBA programs, the real native code set is always UTF-16, as mandated by the Java language specification. However, you can declare any code set as the native code set for the purpose of CORBA, as long as the code set supports the language you need to support.

One limitation applies: only byte-oriented code sets. code sets that do not include null, can be set to `plugins:codeset:char:ncs` and `plugins:codeset:char:ccs`. For `plugins:codeset:wchar:ncs` and `plugins:codeset:wchar:ccs`, any code set can be used, whether it is byte-oriented or not.

C/C++ Internationalization

IDL-to-C++ mapping

Although both C and C++ have different IDL mapping specifications, most CORBA platforms using C map `wchar` and `wstring` with the C++ mappings. Table 4 shows IDL to C++ mappings for both narrow and wide character data.

Table 4: IDL to C++ character mappings

IDL	C++
<code>char</code>	<code>char</code>
<code>string</code>	<code>char *</code>
<code>wchar</code>	<code>wchar_t (CORBA::WChar)</code>
<code>wstring</code>	<code>wchar_t * (CORBA::WChar *)</code>

The latest version of the CORBA C *Language Mapping Specification* can be obtained at <http://www.omg.org/cgi-bin/doc?formal/99-07-35.pdf>.

The latest version of the CORBA C++ *Language Mapping Specification* can be obtained at <http://www.omg.org/cgi-bin/doc?formal/99-07-41.pdf>.

Coding

Because C++ maps wide character data to special datatypes, coding internationalizable applications requires a few modifications:

- C and C++ programs must include the file `locale.h`. This file contains the definitions and headers used to support code sets and locale functionality.
- `setlocale(LC_ALL, "")` must be called in the application's `main()` routine. This operation determines the system's locale settings and initializes the appropriate code sets. The encoding of the `char[]` and `wchar_t` is determined by `setlocale()`.
- The stream output operator `<<` does not take `wchar_t` (IDL `wchar`) or `wchar_t *` (IDL `wstring`) properly when applied to `cout`. In order to print these datatypes, use the wide-oriented stream object `wcout`.

- `idlgen` generates a class named `IT_GeniePrint` which has two methods that appear to print wide data, `print_wstring()` and `print_wchar()`. However, these methods print all characters in the hex form `\xNNNN` and can only be used when a hex dump is needed during debug. All wide data must be converted before being printed.

Example

To implement the `WideEcho.idl` example, [Example 3 on page 12](#), run `idlgen` as follows:

```
idlgen cpp_poa_genie.tcl -all WideEcho.idl
```

This generates a makefile for the Visual C++ `nmake` utility. To run the IDL compiler on Unix systems, run `nmake`.

On the server side, the operations `echo()`, `echoSingleWchar()` and `echoNarrowString()` must be implemented in `WideEchoImpl.cxx`.

[Example 8](#) shows an implementation of the operations.

Example 8: C++ server implementation of `WideEcho.idl`

```
// C++
#include <locale.h>
CORBA::WChar* WideEchoImpl::echo(const CORBA::WChar* ws)
IT_THROW_DECL((CORBA::SystemException))
{
    return CORBA::wstring_dup(ws);
}
CORBA::WChar* WideEchoImpl::echoSingleWChar(CORBA::WChar wc)
IT_THROW_DECL((CORBA::SystemException))
{
    wchar_t x[2];
    x[0] = wc;
    x[1] = (wchar_t) 0;
    return CORBA::wstring_dup(x);
}
```

Example 8: C++ server implementation of *WideEcho.idl*

```

CORBA::WChar* WideEchoImpl::echoNarrowString(const char* ns)
IT_THROW_DECL((CORBA::SystemException))
{
    CORBA::WChar* _result;

    int xlen = strlen(ns)+1; // Max len of buf needed.
    wchar_t *x = new wchar_t[xlen]; // Temp buffer
    mbstowcs(x, ns, xlen);
    _result = CORBA::wstring_dup(x);
    delete [] x; // Clean up temp buffer
    return _result;
}

```

Note: CORBA::WChar is equivalent to wchar_t. idlgen generates code using CORBA::WChar.

Because C++ maps string and wstring to different datatypes, the implementation of echoNarrowString() must explicitly convert ns from char * into wchar_t * using the ANSI C standard function mbstowcs().

On the client side, in client.cxx, the calling convention is no different from the traditional CORBA convention. [Example 9](#) shows a client implementation.

Example 9: C++ client implementation for *WideEcho.idl*

```

// C++
#include <locale.h>
main(int argv, char[] argc)
{
    CORBA::Object_var obj;

    setlocale(LC_ALL, "") // set the locale
    ...
    obj = read_reference("WideEcho.ref");
    WideEcho_var WideEcho1 = WideEcho::_narrow(obj);
    ...
    // Replace Hello with your language equivalent
    wcout << WideEcho1->echo(L"Hello in wstring") << endl
    wcout << WideEcho1->echoSingleWChar(L"H") << endl
    wcout << WideEcho1->echoNarrowString("Hello in string") <<
        endl;
}

```

Native code set

The native code set of C and C++ applications is determined by the platform's locale setting. You must set the Orbix native code set to `plugins:codeset:char:ncs` and `plugins:codeset:wchar:ncs`.

On Windows in Western European locales, the native code set for the narrow char/string (NCS-C) is Windows Code Page 1252 which is ISO 8859-15. Since the OSF registry lacks Window Code Page 1252 or ISO 8859-15, you can use ISO 8859-1 as the best approximation and set the configuration variable as follows:

```
plugins:codeset:char:ncs = "0x00010001"; # ISO-8859-1
```

On Windows in a Japanese locale, the NCS-C is Window Code Page 932, an extension of Shift_JIS. Because the OSF registry also lacks Window Code Page 932, you can use OSF SJIS as the closest approximation and set the configuration variable as follows:

```
plugins:codeset:char:ncs = "0x05000011"; # OSF SJIS
```

On Windows, the native code set for the `wchar/wstring` (NCS-W) is always UCS-2 regardless of locale. Set the configuration variable as follows:

```
plugins:codeset:wchar:ncs = "0x00010100"; # UCS-2 Level 1
```

On Solaris, both NCS-C and NCS-W are determined by the current locale. For the ISO 8859-1 based locales such as `C`, `en`, `fr`, `de`, `es`, `it` and `pt`, both NCS-C and NCS-W should be set to `ISO-8859-1`. So, the configuration variables are set as follows:

```
plugins:codeset:char:ncs = "0x00010001"; # ISO-8859-1
plugins:codeset:wchar:ncs = "0x00010001"; # ISO-8859-1
```

For Solaris in the Japanese `ja` locale, both NCS-C and NCS-W should be set to the OSF code set equivalent of `EUC-JP` as follows:

```
plugins:codeset:char:ncs = "0x00030010"; # JIS eucJP
plugins:codeset:wchar:ncs = "0x00030010"; # JIS eucJP
```

For Solaris in UTF-8 based locales such as `en_US.UTF-8`, `ja_JP.UTF-8`, `ko_KR.UTF-8`, `zh_CN.UTF-8` and `zh_TW.UTF-8`, NCS-C is UTF-8 and NCS-W is UCS-4. Set the configuration variables as follows:

```
plugins:codeset:char:ncs = "0x05010001"; # UTF-8
plugins:codeset:wchar:ncs = "0x00010104"; # UCS-4 Level 1
```

Web Services Internationalization

The Orbix Web services container supports internationalized applications.

In this chapter

This chapter discusses the following topics:

WSDL Test Client

page 30

WSDL Test Client

Overview

Character encodings supported in WSDL Test Client depend on the JDK version installed on your machine. XMLBUS takes a subset of the character encodings supported by the JDK as its own character encodings.

Generally the default character encoding for WSDL test client is UTF-8.

Selecting an encoding

You can select a character encoding for WSDL Test Client from the encoding drop-down list. Using the Interop Test demo as an example, the following steps walk you through selecting an encoding for WSDL Test Client on the Windows platform.

1. Create a configuration domain, sample-domain, by selecting all available services.
2. Start domain services.
3. Set domain environment by running
`IT_PRODUCT_DIR\etc\sample-domain_env.bat.`
4. Go to `IT_PRODUCT_DIR\asp\6.1\bin`, and run `itws_builder.bat`.
5. Open the web page <http://localhost:53205/xmlbus>.
6. Select `Interop Test` from the list and click on **List Services**.
7. Select `Interop TestService` from the services list and click on **List Endpoints**.
8. Click on `Interop TestPort` from **Endpoint Information**.
9. Select `echoString(string)` from listed operations and click **GET TEST FORM**.
10. Click on the **Mime Encoding** drop-down list to see the supported character encodings in the WSDL Test Client.
11. For example, select the encoding `US_ASCII` from the list and enter `Hello World` as the value for the string.
12. Click **INVOKE OPERATION**. The string should be returned correctly.
13. Select another character encoding and input a string value.
14. Invoke the operation to see whether the string displays correctly.

Most names in the **Mime Encoding** drop-down list are IANA names. However, there are some exceptions. Most notably, `EBCDIC-CP-xx` is used for the EBCDIC family of encodings, as listed in [Table 5](#).

Table 5: *Mime Encodings for WSDL Test Client*

Mime Encoding List	Java Encoding	Meaning
EBCDIC-CP-YU	Cp870	IBM Multilingual Latin-2
EBCDIC-CP-US	Cp037	USA, Canada, Netherlands, Portugal, Brazil, Australia
EBCDIC-CP-SE	Cp278	IBM Finland, Sweden
EBCDIC-CP-ROECE	Cp870	IBM Multilingual Latin-2
EBCDIC-CP-NO	Cp277	IBM Denmark, Norway
EBCDIC-CP-NL	Cp037	USA, Canada, Netherlands, Portugal, Brazil, Australia
EBCDIC-CP-IT	Cp280	IBM Italy
EBCDIC-CP-IS	Cp871	IBM Iceland
EBCDIC-CP-HE	Cp424	IBM Hebrew
EBCDIC-CP-GB	Cp285	IBM United Kingdom, Ireland
EBCDIC-CP-FR	Cp297	IBM France
EBCDIC-CP-FI	Cp278	IBM Finland, Sweden
EBCDIC-CP-ES	Cp284	IBM Catalan/Spain, Spanish Latin America
EBCDIC-CP-DK	Cp277	IBM Denmark, Norway
EBCDIC-CP-CH	Cp500	EBCDIC 500V1
EBCDIC-CP-CA	Cp037	USA, Canada, Netherlands, Portugal, Brazil, Australia
EBCDIC-CP-AR1	Cp420	IBM Arabic
EBCDIC-CP-AR2	Cp918	IBM Pakistan (Urdu)

Restrictions

Orbix has some limitations in its internationalization support.

In this chapter

This chapter discusses the following topics:

Translations	page 34
Property Values	page 35
COMet	page 36

Translations

Orbix is internationalized, but it is not localized. All the GUI applications and messages remain in English.

There are some exceptions. Some GUI elements and messages that originate from the underlying operating system or Java run-time environment are localized automatically. For example, the **OK** button is translated in some dialog boxes.

Similarly, some messages from the operating systems are in the language of the locale.

Property Values

Generally speaking, various properties that Orbix uses are restricted to the traditional ASCII range.

For example, the following properties must be in ASCII in order for IONA to guarantee their proper behavior:

- File path
- User ID
- Password
- URL
- Repository name
- Channel name
- Configuration domain name
- Configuration variable name
- Configuration variable value
- Scope name
- Cluster name
- (J2EE) Application name

The only Orbix property that can have non-ASCII characters is the role names if Orbix is configured to use an LDAP server that supports non-ASCII.

COMet

COMet supports non-ASCII characters in the same manner as regular ORBs using the IDL `string` and `char` datatypes. It does not support `wstring` or `wchar`.

Be careful that a multibyte character is not passed to a `char` parameter because a `char` is a single byte storage. Attempts to do so result in an abnormal program termination on the COMet side.

Index

A

ASCII 2

C

character encoding schema 2
code set 2
CODESET_INCOMPATIBLE 17
code set negotiation 13, 17
Conversion code set 16
CORBA::WChar 24

I

International Standard Organization 5
Internet Assigned Number Authority 3
ISO 8859 2
it_basic_codeset 18

L

Light weight codeset plugin 18
locale 5
locale.h 24

N

native code set 16

O

org.omg.CORBA.DATA_CONVERSION 22
OSF Code Set Registry 4

P

plugins.codeset.char.ccs 18
plugins.codeset.char.ncs 18, 27
plugins.codeset.wchar.ccs 18
plugins.codeset.wchar.ncs 18, 27
print_wchar() 25
print_wstring() 25

S

setlocale() 24
String 22

T

transmission code set 16

U

Unicode 3

W

wchar 22
wchar_t 24
wide character 12
wide string 12
wstring 22

