



유니 코드

국민대학교 소프트웨어학부

강 승 식

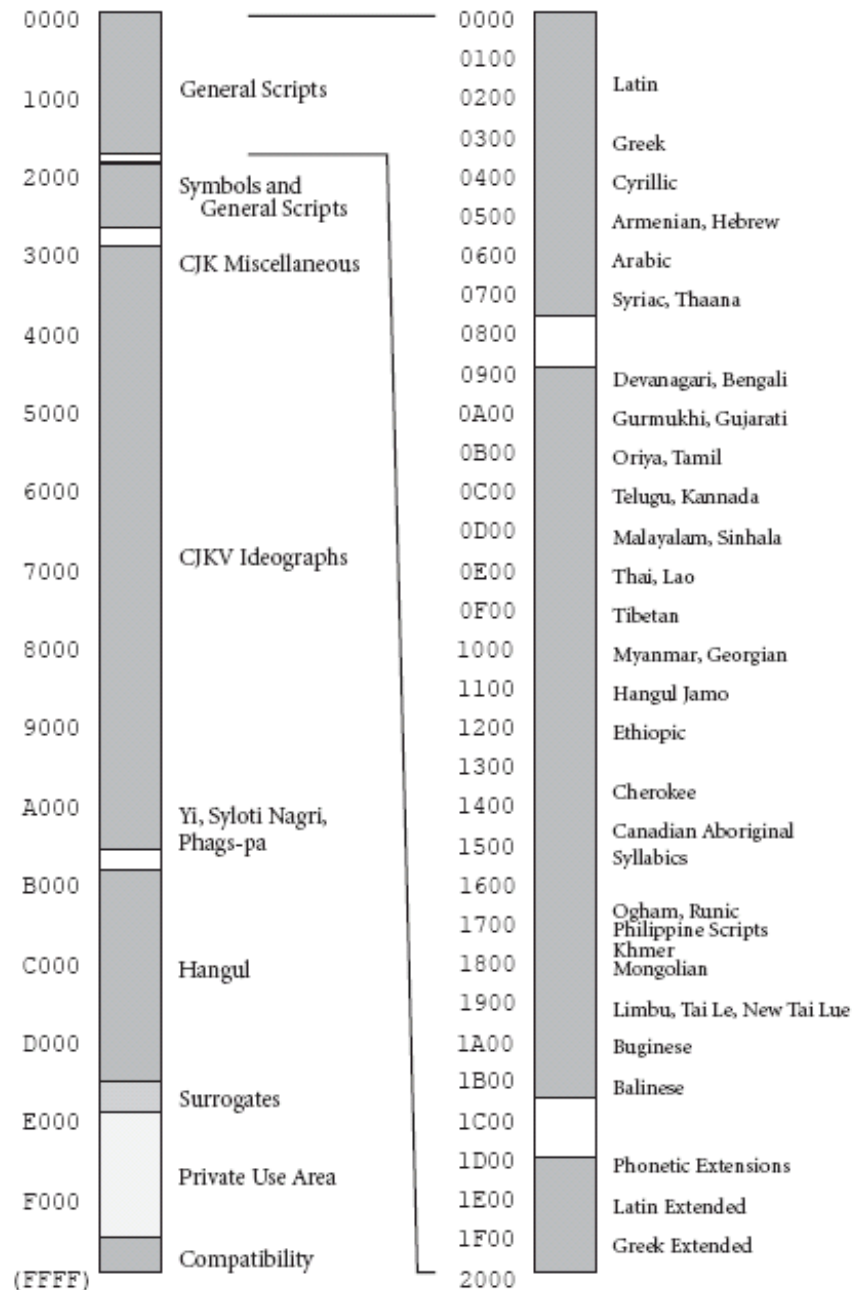
- What is unicode?

- BMP(Basic Multilingual Plane), SMP
- UCS-2, UCS-4
- UTF-8, UTF-16, UTF-32
- BOM(Byte Order Mark)
- Big Endian, Little Endian

- 유니코드 문자셋 정의 영역: 언어별로

- 각 언어별 문자 인식 방법
 - 한글/한자/아랍어/러시아어 등 어떤 영역의 문자인지 확인 방법
- 다국어 문자가 혼합되어 있을 때 언어별로 **tokenization**

Figure 2-14. Allocation on the BMP



http://www.unicode.org/

The Unicode Consortium: Home Page - Beta - Windows Internet Explorer

http://www.unicode.org/

Unicode Home Page Beta

Older Version | Feedback | Site Map | Search

New to Unicode

- General Information
- The Consortium
- The Unicode Standard
- Key Specifications
- Technical Publications
- Work in Progress
- For Members Only

- What is Unicode?
- How to Use this Site
- FAQ
- Glossary of Unicode Terms

Unicode enables people around the world to use [our members](#) develop the [Unicode Standard](#), and other standards. These [specifications](#) form the foundation for Unicode internationalization in all major operating systems, search engines, applications, and the Web.

"Unicode... is a solid foundation for e-business in a global economy." — Steve Mills ([more](#))

Some of our Members

- NetApp
- THE CHURCH OF JESUS CHRIST OF LATTER-DAY SAINTS

News & Announcements

Internationalization & Unicode Conference 32
September 8–10, 2008 • San Jose, CA USA

Recent releases:

- Translated Terminology Section, Guide to Abbreviations in Standards and FAQ on Standards Developing Organizations (2008.07.29)
- Unicode enabled products and useful resources pages reorganized for easier access (2008.07.29)
- Unicode CLDR1.6.1 release for Olson Timezone Changes (2008.07.25)
- Updated version of UTR #36: Unicode Security Considerations now available (2008.07.17)
- Unicode Releases Common Locale Data Repository, Version 1.6 (2008.07.02)

Other news:

- Program for IUC 32 Available Online (2008.06.03)
- NetApp joins Unicode at the Full level (2008.05.08)

(Archived announcements...)

Public review issues:

- Proposed Update UAX #44: Unicode Character Database
- Proposed Update UTR #17: Unicode Character Encoding Model
- Proposed Update UTR #33: Unicode Conformance Model
- Proposed Update UTR #23: Unicode Character Property Model
- Bengali Currency Numerator Values

Technical Reports - Windows Internet Explorer

http://www.unicode.org/reports/index.html

Technical Reports

Unicode Technical Reports

For more information see [About Unicode Technical Reports](#) and the [Specifications](#)

Standard Annexes

- UAX 9 [Unicode Bidirectional Algorithm](#)
- UAX 11 [East Asian Width](#)
- UAX 14 [Unicode Line Breaking Algorithm](#)
- UAX 15 [Unicode Normalization Forms](#)
- UAX 24 [Unicode Script Property](#)
- UAX 29 [Unicode Text Segmentation](#)
- UAX 31 [Unicode Identifier and Pattern Syntax](#)
- UAX 34 [Unicode Named Character Sequences](#)
- UAX 38 [Unicode Han Database \(UniHan\)](#)
- UAX 41 [Common References for Unicode Standard Annexes](#)
- UAX 42 [Unicode Character Database in XML](#)
- UAX 44 [Unicode Character Database](#)
See also [Proposed Update](#)

Technical Standards

- UTS 6 [A Standard Compression Scheme for Unicode](#)
- UTS 10 [Unicode Collation Algorithm](#)
- UTS 18 [Unicode Regular Expressions](#)
See also [Proposed Update](#)
- UTS 22 [Character Mapping Markup Language](#)
- UTS 35 [Unicode Locale Data Markup Language \(LDML\)](#)
- UTS 37 [Ideographic Variation Database](#)
- UTS 39 [Unicode Security Mechanisms](#)

Technical Reports

- UTR 16 [UTF-EBCDIC](#)
- UTR 17 [Character Encoding Model](#)
See also [Proposed Update](#)
- UTR 20 [Unicode in XML and other Markup Languages](#)
- UTR 23 [The Unicode Character Property Model](#)
See also [Proposed Update](#)
- UTR 25 [Unicode Support for Mathematics](#)
See also [Proposed Update](#)
- UTR 26 [Compatibility Encoding Scheme for UTF-16: 8-Bit \(CESU-8\)](#)
- UTR 30 [Unicode Character Foldings Status: Draft](#)

인터넷

FAQ - UTF-8, UTF-16, UTF-32 & BOM - Windows Internet Explorer

http://unicode.org/faq/utf_bom.html#28

FAQ - UTF-8, UTF-16, UTF-32 & BOM

UTF-8, UTF-16, UTF-32 & BOM

General questions, relating to UTF or Encoding Forms:

- [Can Unicode text be represented in more than one way?](#)
- [What is a UTF?](#)
- [Where can I get more information on encoding forms?](#)
- [How do I write a UTF converter?](#)
- [Which of the UTFs do I need to support?](#)
- [What are some of the differences between the UTFs?](#)
- [Why do some UTFs have a BE or LE in their label, as in UTF-16LE?](#)
- [Are there any byte sequences that are not generated by a UTF? How should I interpret them?](#)
- [Is there a standard method to package a Unicode character so it fits an 8-Bit ASCII stream?](#)
- [Which of these approaches is the best?](#)
- [Which of these formats is the most standard?](#)

UTF-8 FAQ:

- [What is the definition of UTF-8?](#)
- [Is the UTF-8 encoding scheme the same irrespective of whether the underlying processor is little endian or big endian?](#)
- [Is the UTF-8 encoding scheme the same irrespective of whether the underlying system uses ASCII or EBCDIC encoding?](#)
- [How do I convert a UTF-16 surrogate pair such as <D800 DC00> to UTF-8? A one four byte sequence or as two separate 3-byte sequences?](#)
- [How do I convert an unpaired UTF-16 surrogate to UTF-8?](#)

UTF-16 FAQ:

- [What is UTF-16?](#)
- [What are surrogates?](#)
- [What is the algorithm to convert from UTF-16 to character codes?](#)
- [Why are some people opposed to UTF-16?](#)
- [Will UTF-16 ever be extended to more than a million characters?](#)
- [Are there any 16-bit values that are invalid?](#)
- [Are there any paired surrogates that are invalid?](#)
- [Since the surrogate pairs will be rare, does that mean I can dispense with them?](#)
- [When will most implementations of Unicode support surrogates?](#)

인터넷

http://www.unicode.org/ - Windows Internet Explorer

http://www.unicode.org/

Code Charts

Look up by character code: Go

Home | Site Map | Search

The Unicode Character Code Charts By Script

updated for Unicode 5.1

SYMBOLS AND PUNCTUATION | NAME INDEX | HELP AND LINKS

European Alphabets	African Scripts	Indic Scripts	East Asian Scripts	Central Asian Scripts
(see also Comb. Marks)	Ethiopic	Bengali	Han Ideographs	Kharoshthi
Armenian	Ethiopic	Devanagari	Unified CJK Ideographs (5MB)	Mongolian
Armenian	Ethiopic Supplement	Gujarati	CJK Ideographs Ext. A (2MB)	Phags-Pa
Armenian Ligatures	Ethiopic Extended	Gurmukhi	CJK Ideographs Ext. B (13MB)	Tibetan
Coptic	Other African scripts	Kannada		Ancient Scripts
Coptic	N'Ko	Lepcha	Compatibility Ideographs (.5MB)	Ancient Greek
Coptic in Greek block	Osmanya	Limbu	... Supplement (.5MB)	Ancient Greek Numbers
Cyrillic	Tifinagh	Malayalam	Kanbun	Ancient Greek Musical
Cyrillic	Vai	Oj Chiki	(see also Unihan Database)	Cuneiform
Cyrillic Supplement	Middle Eastern Scripts	Oriya	Radicals and Strokes	Cuneiform
Cyrillic Extended A	Arabic	Saurashtra	CJK Radicals	Cuneiform Numbers
Cyrillic Extended B	Arabic	Sinhala	KangXi Radicals	Old Persian
Georgian	Arabic Supplement	Syloti Nagri	CJK Strokes	Ugaritic
Georgian	Arabic Present. Forms A	Tamil	Ideographic Description	Linear B
Georgian Supplement	Arabic Present. Forms B	Telugu	Chinese-specific	Linear B Syllabary
Greek	Hebrew	South East Asian	Bopomofo	Linear B Ideograms
Greek	Hebrew	Balinese	Bopomofo Extended	Other Ancient Scripts
Greek Extended	Hebrew Present. Forms	Buginese	Japanese-specific	Aegean Numbers
(see also Ancient Greek)	Syriac	Cham	Hiragana	Ancient Symbols
Latin	Syriac	Kayah Li	Katakana	Carian
Basic Latin	Thaana	Khmer	Katakana Phonetic Ext.	Counting Rod Numerals
Latin-1	Thaana	Khmer Symbols	Halfwidth Katakana	Cypriot Syllabary
Latin Extended A	American scripts	Lao	Korean-specific	Glagolitic
Latin Extended B	Canadian Syllabics	Myanmar	Hangul Syllables (4MB)	Gothic
Latin Extended C	Cherokee	New Tai Lue	Hangul Jamo	Lycian
Latin Extended D	Deseret	Rejang	Hangul Compatibility Jamo	Lydian
Latin Extended Additional	Philippine Scripts	Sundanese	Halfwidth Jamo	Ogham
Latin Ligatures	Buhid	Tai Le	Yi	Old Italic
Fullwidth Latin Letters	Hanunoo	Thai	Yi (.6MB)	Phaistos Disc
Small Forms	Tagalog	Other Scripts	Yi Radicals	Phoenician
(see also Phonetic Symbols)	Tagbanwa	Shavian		Runic

To get a list of code charts for a character, enter its code in the search box at the top. To access a chart for a given block, click on its entry in the table. The charts are [PDF](#) files, and some of them may be very large. For frequent access to the same chart, right-click and save the file to your disk. For an alphabetical index of character and block names, use the [Unicode Character Names Index](#). For terms of use, conventions used in this table, access to additional charts and ways to access the code charts, see [Character Code Chart Help and Links](#).

유니코드 2.0

- 한글 11,172자를 순서대로 코드 부여
 - 코드 범위 : ‘가’(0xAC00) ~ ‘힉’(0xD7A3)
- 자모 코드(0x31xx 영역)
 - 자음 : 0x3131 ~ 0x314E
 - 모음 : 0x314F ~ 0x3163
 - 채움 코드 : 0x3164
 - 옛글 자모 : 0x3165 ~ 0x318E
- 초성-중성-종성 코드(0x11xx 영역)
- CJKV 한자 코드
 - 한-중-일 공통 한자는 1개의 코드 부여
 - 음가 2 이상인 한자에 하나의 코드만 부여

유니코드에서 초성/중성/종성 인식

- 초성

$$\left(\left(\text{코드값} - 0xAC00 \right) / 28 \right) / 21 \right) \% 19$$

- 중성

$$\left(\left(\text{코드값} - 0xAC00 \right) / 28 \right) \% 21$$

- 종성

$$\left(\text{코드값} - 0xAC00 \right) \% 28$$

유니코드 <-> KS완성형 변환

- 유니코드 <-> KS완성형 변환
 - 한글 11,172자
 - 자모
 - 문장부호: 가운데점 등
 - 한자
- 코드변환 테이블 작성
 - Windows의 notepad(메모장) 이용
- Java, Perl에서 유니코드-완성형 변환 방법
- “iconv” 라이브러리 사용

Topics

1. Unicode란 무엇인가?

- Unicode의 특징, EUC-KR과 다른 점 등
- Unicode의 한글영역: 음절, 자음과 모음, 특수문자(가운데점 등)

2. Unicode의 국가별 문자셋 영역(한글 제외)

- 아시아 언어 영역 -- 일본어 영역, 중국어(한자) 영역 등
- 유럽어 영역 -- 독일어의 움라우트, 러시아어 등
- 아랍어 영역
- 기타 언어 영역

3. Tokenization program 작성

한/중/일/아랍어/러시아어/문장부호 등이 혼합된 입력 문자열에 대해 각 언어 영역별로 문자열을 분리하는 프로그램

4. Unicode Encoding Forms

UTF-8, UTF-16, UTF-32

5. Unicode 관련 용어 설명

<http://www.unicode.org/>의 "Glossary of Unicode Terms" 참조

- Big Endian, Little Endian
- BMP Code Point, Supplementary Code Point, Surrogate Code Point, High-Surrogate Code Point
- BOM(Byte Order Mark)
- CJK
- Code Page, Code Point, Codespace
- DBCS, SBCS
- Consonant, Vowel, Diphthong, Jamo, Leading Consonant
- Syllable block, Hangul Syllable Block