

CSCC09F

Programming on the Web



i18n and Unicode

extending the Web to all
written languages

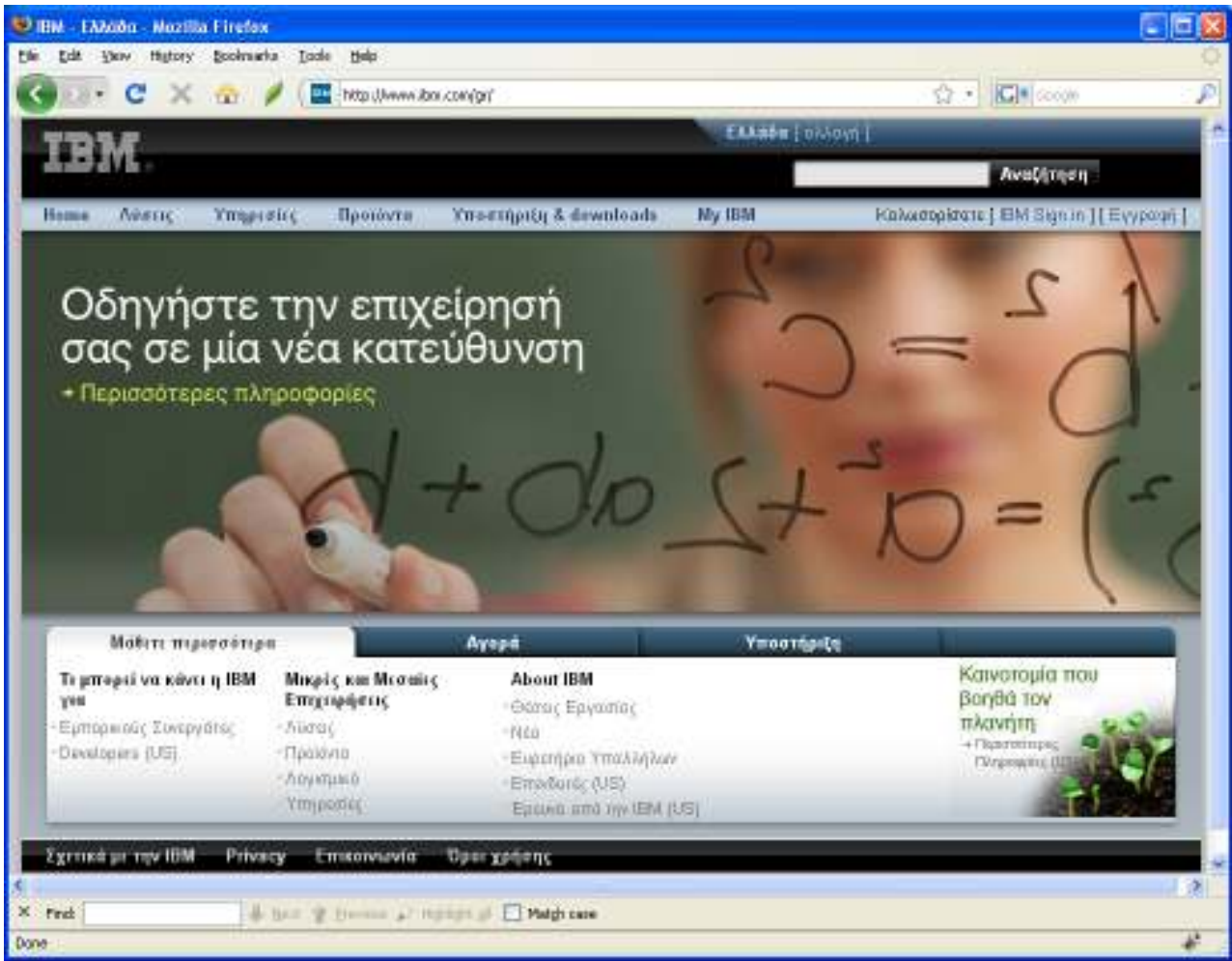
Internationalization

- ❑ 'i18n' for short (because we like acronyms!)
- ❑ deals with:
 - character sets (“code-points”)
 - text direction
 - local formats (dialects)
 - number formats (e.g. “.” vs “,”)
 - word-break rules (hyphenation conventions)
 - sorting order
 - calendars and dates
- ❑ the Web must be i18n, why?

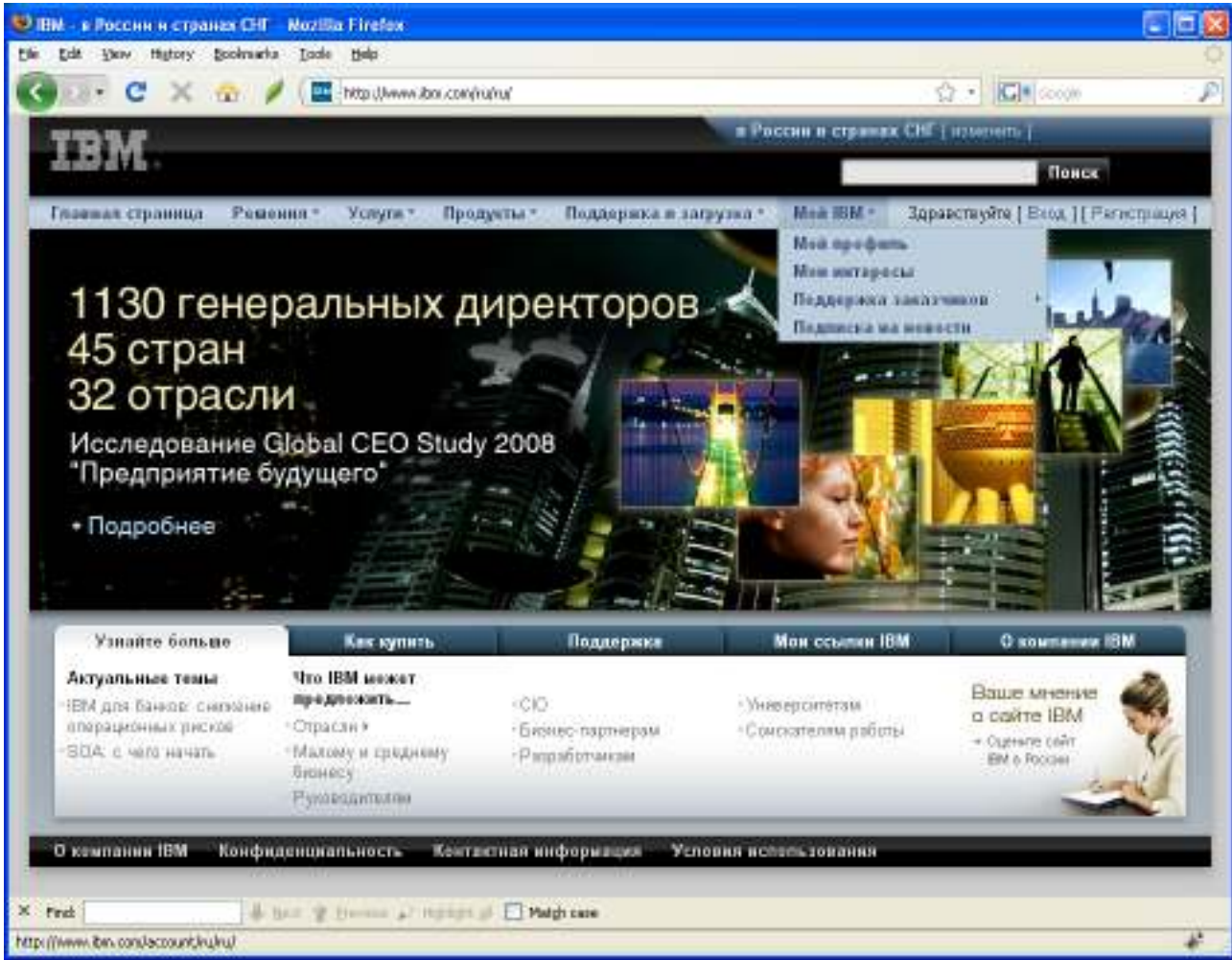
Internationalization (i18n)

- ❑ want worldwide reach for products and services
- ❑ increase user/customer satisfaction, which in turn can increase business
- ❑ improve quality of customer support communication, in turn reducing cost
- ❑ open new markets for goods and services, where English does not meet local needs
- ❑ decreased expenses: develop once, deploy in multiple markets

Internationalization (i18n)



Internationalization (i18n)

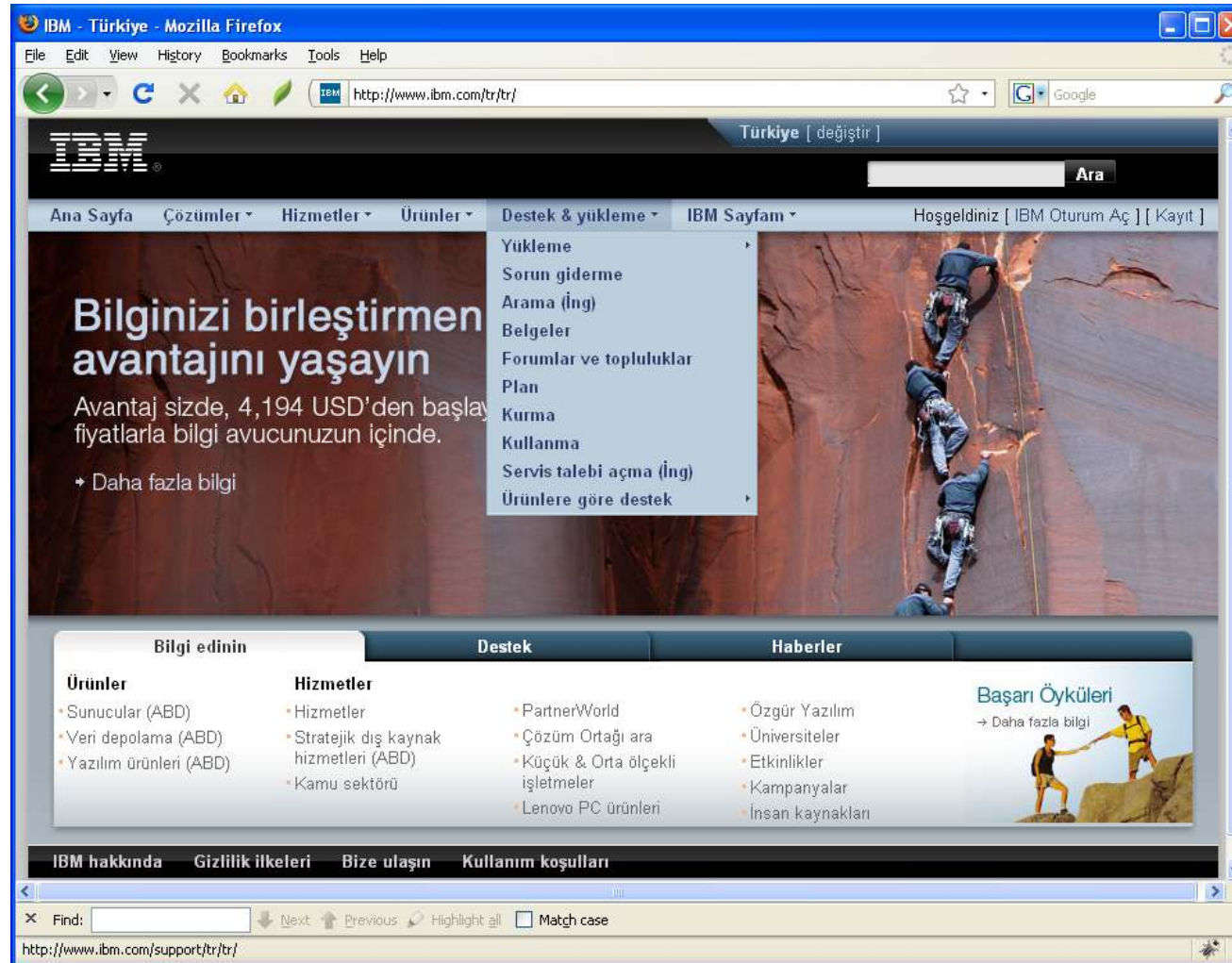


18 - Unicode

CSCC09

5

Internationalization (i18n)



18 - Unicode

CSCC09

6

Internationalization (i18n)



Language Transliteration

- ❑ ஒலீன் (Tamil)
- ❑ অ্যালেন (Bengali)
- ❑ Алэн (Russian)
- ❑ 阿蘭 (Chinese)
- ❑ Алан (Mongolian)
- ❑ آلان (Arabic)
- ❑ אֵלָן (Hebrew)
- ❑ एलन (Nepali)
- ❑ Αλαν (Greek)
- ❑ アラン (Japanese)

Character Encodings

❑ Character Set:

- a collection of characters grouped together
 - ❑ code point \leftrightarrow abstract character
 - e.g. the lower-case letter 'n' as used in English

❑ Character Encoding:

- a way of encoding code points into a byte stream
 - ❑ byte stream \leftrightarrow code points

❑ Font

- a collection of character "pictures"
 - ❑ code point \leftrightarrow glyph

Single-Byte Character Encoding

❑ ASCII

- American Standard Code for Information Interchange
 - ❑ 7-bit code (0-127), why 7 bits?
 - ❑ because it was designed in the 1960's to work on systems (h/w) with 7-bit bytes, and for 8-bit bytes with one bit as a parity (error checking) bit
- most computers today use 8-bit bytes, opening up an additional 128 values for use as character encodings

Single-Byte Character Encoding

❑ ISO-8859 Series

- A set of 8-bit code-points backward compatible with ASCII
 - ❑ 96 extra character slots (other 32 slots for control codes)
- ISO-8859-0 through ISO-8859-15 (new default)
 - ❑ e.g., ISO-8859-1 / Latin 1
 - Afrikaans, Albanian, Basque, Catalan, Danish, Dutch, English, Faroese, Finish, French, German, Icelandic, Irish, Italian, Norwegian, Portuguese, Rhaeto-Romanic, Scottish, Spanish, Swahili, Swedish
 - Default code set for HTML
- ISO-8859-6 (Arabic) and ISO-8859-7 (Greek) do not fully handle all their respective languages

ISO-8859 Series

- ❑ For all members of the series, code-points 0-127 identical to US-ASCII
- ❑ 128-159 for less-used control characters
 - ISO-6429
- ❑ Upper portion of code-points, 160-255, differ for each variant (Latin, Cyrillic, Arabic, etc.)
- ❑ Good information at
<http://czyborra.com/>

ISO-8859-1 (Latin1)

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	í	φ	£	¤	¥	¦	§	¨	©	ª	«	¬	­	®	¯
B0	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾
C0	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î
D0	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ
E0	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î
F0	ö	ñ	õ	ö	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	ÿ

Latin1 is also represented as the first 256 code points of ISO-10646 (Unicode)

ISO-8859-5 (Cyrillic)

A0	A1 .. Ё	A2 ѐ	A3 ѓ	A4 ё	A5 ѕ	A6 і	A7 ї	A8 ј	A9 љ	AA њ	AB ћ	AC ќ	AD –	AE џ	AF џ
B0 А	B1 Б	B2 В	B3 Г	B4 Д	B5 Е	B6 Ж	B7 З	B8 И	B9 Й	BA К	BB Л	BC М	BD Н	BE О	BF П
C0 Р	C1 С	C2 Т	C3 У	C4 Ф	C5 Х	C6 Ц	C7 Ч	C8 Ш	C9 Щ	CA Ъ	CB Ы	CC Ь	CD Э	CE Ю	CF Я
D0 а	D1 б	D2 в	D3 г	D4 д	D5 е	D6 ж	D7 з	D8 и	D9 й	DA к	DB л	DC м	DD н	DE о	DF п
E0 р	E1 с	E2 т	E3 у	E4 ф	E5 х	E6 ц	E7 ч	E8 ш	E9 щ	EA ъ	EB ы	EC ь	ED э	EE ю	EF я
F0 ѐ	F1 ё	F2 ђ	F3 ѓ	F4 ё	F5 ѕ	F6 і	F7 ї	F8 ј	F9 љ	FA њ	FB ћ	FC ќ	FD –	FE џ	FF џ

ISO-8859-6 (Arabic)

A0				A4	ﻩ							AC	ﺍ	AD	—		
											BB	ﺓ				BF	؟
	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF		
	ﺔ	ﻑ	ﻏ	ﻮ	ﺍ	ﻛﻲ	ﺍ	ﺏ	ﻮ	ﻧ	ﻧ	ﺝ	ﺝ	ﺡ	ﺩ		
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA							
ﺩ	ﺭ	ﺯ	ﺳ	ﺸﺮ	ﺺ	ﺾ	ﻁ	ظ	ﻋ	ﻐ							
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF		
—	ﻑ	ﻘﻲ	ﻛ	ﻝ	ﻫﻢ	ﻥ	ﻮ	ﻮ	ﻲ	ﻲ	ﺶ	ﻲ	ﺶ	ﺶ	ﺶ		
F0	F1	F2															
ﺶ	ﻮ	ﻮ															

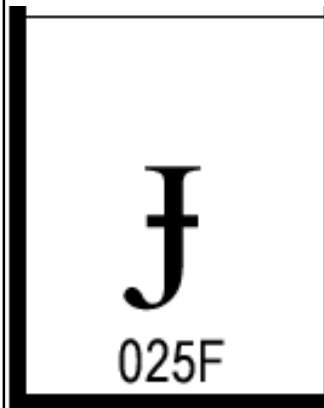
Multi-byte Code Sets

- ❑ Not every language can be shoehorned into 256 code-points (8-bit rep)
 - e.g. languages such as Chinese, Japanese, and Korean (CJK) require tens of thousands of code-points
 - also desirable to represent multi-language documents – thus need to have enough code points to segregate languages.
- ❑ 32-bits could handle up to 4 billion characters
 - more than sufficient, but ... don't want to impose this space inefficiency except where needed – not a reasonable generic standard
- ❑ Is there a happy medium between these?

The Unicode Standard

- ❑ Defines a code-point set with a fixed-width, 16-bit character-encoding scheme
 - characters from all the world's major scripts are uniformly supported
 - ❑ can combine Arabic, French, Japanese, and Russian characters all in the same string/document
 - ❑ even musical symbols, mathematics, historical languages, aboriginal languages, ecclesiastical scripts, ...
 - 65,536 unique slots (16-bits)
 - note, total estimated # characters introduced since the dawn of writing \approx 500,000 (more on this later)
 - for ease of interoperability, first 128 Unicode code point values are assigned to match ASCII, first 256 Unicode code point values match ISO-8859-1 (Latin-1)

The Unicode Specification



□ <http://www.unicode.org/>

□ e.g. 025F

- LATIN SMALL LETTER DOTLESS J WITH STROKE

- voiced palatal stop

- typographically a flipped f, but better thought of as a form of j

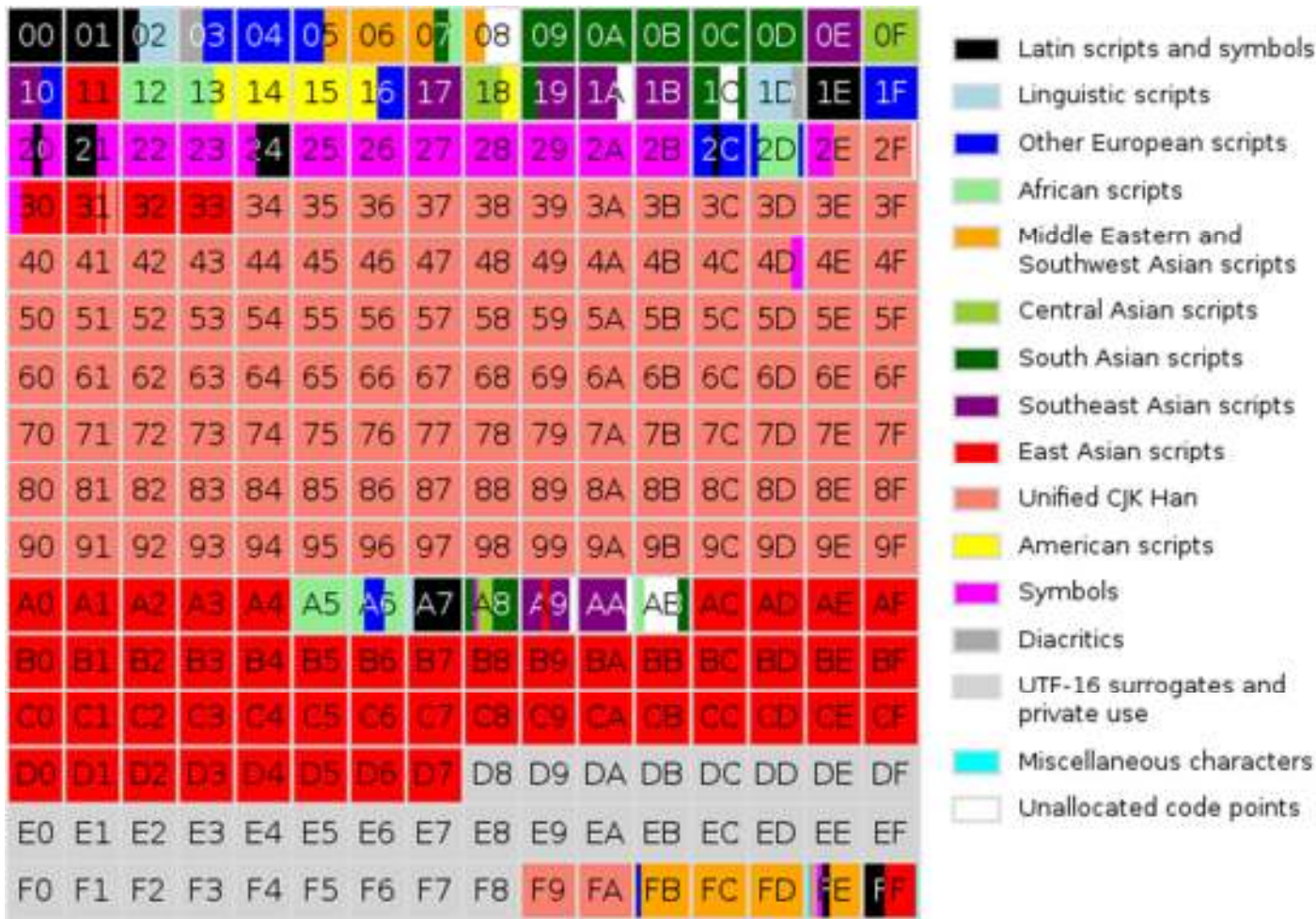
- "gy" in Hungarian orthography

- also archaic phonetic for palatoalveolar affricate 02A4



□ See www.unicode.org/charts for code charts including all encoded languages

Unicode Code-point Distribution



Unicode Design Principles

- ❑ 16-bit characters (mostly)
- ❑ non-modal – no notion of being in certain state or language mode
- ❑ characters, not glyphs (a glyph is the visual presentation of one or more characters)
- ❑ Semantic properties associated with character
 - map lower to upper
 - maps mirrored character (), { }, etc
 - indicates numeric values
 - indicates directionality
 - sample glyph, ...

Unicode Encoding Methods

- ❑ UTF-16 (UCS Transformation Format), UCS-2
 - default encoding method; straightforward
 - suited for storing UNICODE in memory (easy access)
- ❑ UTF-7
 - modal encoding method that uses only first 7-bits of a byte – introduced to cope with 7-bit apps like e-mail
- ❑ UTF-8
 - 8-bit, variable-length encoding
 - ❑ 1, 2, or 3 bytes could represent a UNICODE code element
 - ❑ ASCII compatible, mandated for use by Internet protocols
 - ❑ space efficient for Euro-text but costs 50% more for Asian
- ❑ UTF-32, UCS-4
 - Simplest, easy to access on 32-bit machines: 32-bit integer per Unicode character, inefficient for Latin/ASCII

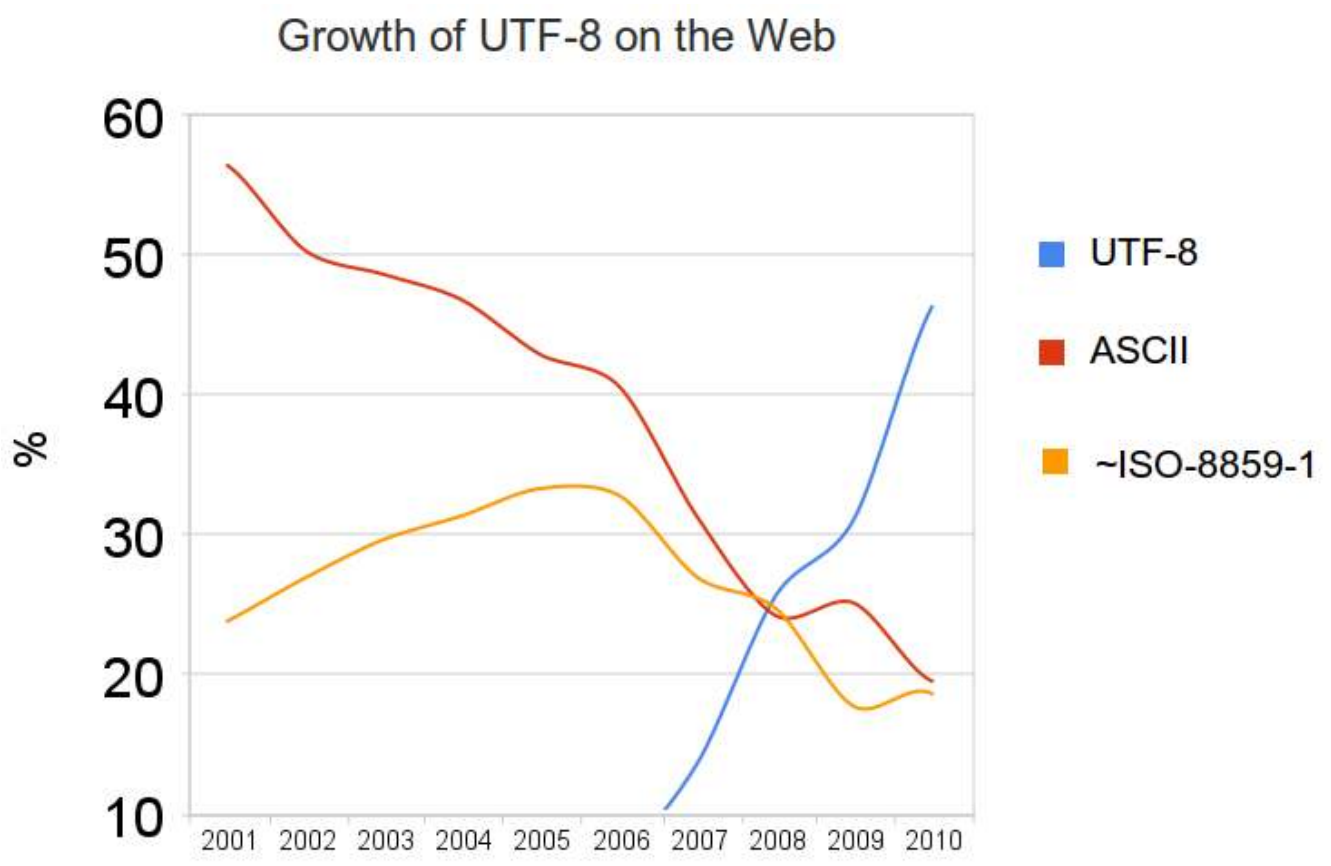
UCS-2 in UTF-8

UCS-2 Encoding Range	UTF-8 Bit Patterns
\u0000 - \u007F	0xxxxxxx
\u0080 - \u07FF	110xxxxx 10xxxxxx
\u0800 - \uFFFF	1110xxxx 110xxxxx 10xxxxxx

Choosing an Encoding

- ❑ Commonly used character encodings on the Web:
 - ASCII
 - ISO-8859-1 (Latin-1)
 - ISO-8859-n (e.g. n=5 supports Cyrillic)
 - UTF-8
 - UTF-16
 - SHIFT_JIS (a Japanese encoding)
 - EUC_JP (another Japanese encoding)
- ❑ Conforming UA's must correctly map supported encodings to UNICODE code-points

Character Encodings on the Web



18 - Unicode

CSCC09

29

Specifying Document Character Encoding

- ❑ In the HTTP header

Content-Type: text/html; charset=EUC-JP

- ❑ server must somehow determine encoding of document

- ❑ To address server or config limitations, HTML documents may include explicit info in meta

**<meta http-equiv="Content-Type"
content="text/html; charset=EUC-JP">**

- ❑ may only be used where the character-encoding is ASCII compatible (at least until the META element is parsed)

- ❑ For certainty, may include a 'charset' attribute on an element that designates an external resource (e.g. <a> and <link> elements)

Character References

- ❑ Numeric (decimal or hex)
 - `&#D;` refers to UNICODE char decimal D
 - `&#xH;` refers to UNICODE char hex H
 - ❑ `å = å = å = å`
- ❑ Character Entity References (defined in XML DTD)
 - symbolic names rather than code positions
 - small subset of char's predefined as “entities”, e.g.,
 - ❑ `<` (useful to escape a `<`)
 - ❑ ` ` (a non-breaking space)
 - ❑ `©` (copyright symbol ©)
 - can define own, e.g. `<!ENTITY euro "&x20AC;">`

Example

- ❑ Example file iso8859.html inserts a symbol in 3 ways:
 - `¾`;
 - `¾`
 - the raw byte value 190 inserted into the file
- ❑ The character encoding should only effect the way the last character is displayed.
- ❑ iso8859.html
 - change encoding to Cyrillic (ISO-8859-5) to get a big O
 - change encoding to Greek (ISO-8859-7) to get a big Y

Unicode Surrogate Pairs

- ❑ Unicode goal is to represent all language chars, but have only 16-bit code space, how does this work?
- ❑ 2048 unused slots in the Unicode hex range D800-DFFF are reserved to accommodate the overflow characters
- ❑ provide room for 1,048,576 (less commonly used) extra characters via "surrogate pairs"
 - Unicode 3.1 assigns more than 40,000 characters to the surrogate pair coding
- ❑ D800-DBFF followed by DC00-DFFF
 - no meaning in isolation (unpaired)