Internationalization: A Guide for the Perplexed

by Peter Saint-André for RFC Editor Team, November 2013

Topic #1: What's the Problem?

in the beginning was ASCII

and ASCII was with the Internet

and ASCII was [a false] god :)

problem:ASCII is extremely limited

there are thousands of languages and scripts

can't we force everyone to use ASCII?

sorry, but that's incredibly naïve

we need to encode more than [A-Z][a-z]

i.e., we need internationalization

specifically, we need Unicode...

[unicode.org]

a set of every* character humans care about

[* sorry, no Klingon, Elvish, or Dwarvish]

technically, Unicode is a "coded character set" (RFC 6365)

i.e., each character has a alphanumeric code assigned to it

in Unicode, these codes are hexadecimal (i.e., base 16 instead of base 10)

we use the convention U+xxx [RFC 5137]

e.g., (ASCII) SPACE is the 32nd code point, i.e., U+0020

similarly, SOLIDUS "/" is the 47th code point, i.e., U+002F

P = U + 0050

P = U+0070

$\pi = U + 03C0$

平 = U+5E73

and so on up to ∞ [U+221e]

(well, up to 1,114,111;)

each character also has various properties

letter, number, symbol, punctuation, etc.

case: UPPER vs. lower vs. Title

status as a modifier (e.g., accent mark)

width: f u l l , half, narrow, even zero!

left-to-right vs. right-to-left

etc.

each character looks* and behaves differently

[* mostly – we'll talk about confusable characters...]

we handle characters based on properties

but there be dragons!

case folding

character equivalence

decomposition and recomposition

normalization

various encodings

string comparison

mappings (e.g., based on user locale)

right-to-left vs. left-to-right scripts and characters

confusable characters

enforcement in protocols (or documentation :)

registration policies

versioning

user interface issues

reliance on rendering engines

and more!

plus, many rules have exceptions!

lots of messy complexity

are you scared yet?

if not, you will be soon :)

Topic #2: Case

(note: symbols like 4 are only my personal convention)

[U+0049] ↓ [U+0069]

and so on, right?

not so fast!

only a few scripts have the notion of "case"

some characters don't map cleanly

e.g., German esszett, which is only lowercase

$B \Rightarrow SS \lor ss$

[U+00DF] ⇒ [U+0053] [U+0053] → [U+0073] [U+0073]

therefore: FUSSBALL → fussball, not fußball

e.g., Greek final sigma

even worse, locale & context matter...

e.g., "Turkish dotless i"

L L L & I & I → I U+0049] → [U+0131] & [U+0069] → [U+0130]

these differences can have consequences

thankfully, not *that* many exceptions

Topic #3: Character Equivalence

one character can be equivalent to another

$\overset{o}{A} \equiv \overset{o}{A}$ $[U+212B] \equiv [U+00C5]$ $(angstrom sign) \equiv (``a'' with ring above)$

one character can be equivalent to a sequence of characters

$\mathring{A} \equiv A + \circ$

[U+00C5] = [U+0041] + [U+030A]

$c = c + c_{s}$ [U+00E7] = [U+0063] + [U+0327]

$\check{R} \equiv R + \check{}$

[U+0158] = [U+0052] + [U+030C]

Ř is a "composite character"

R + is a combining sequence

composite characters are your friends :-)

two kinds of equivalence...

(a) Canonical Equivalence

κανών = rule, standard, measure

"this character is the standard for that one"

characters look and behave the same

$\overset{o}{A} \equiv \overset{o}{A}$ $[U+212B] \equiv [U+00C5]$ $(angstrom sign) \equiv (``a'' with ring above)$

$\overset{\circ}{A} \equiv A + \overset{\circ}{}$ $[U+00C5] \equiv [U+0041] + [U+030A]$

(2) Compatibility Equivalence

compati = "suffer with"

"this character suffers with that one"

a.k.a. we suffer with compatibility equivalence :-)

often for the sake of backward compatibility

characters might look and behave differently

 $[U+2163] \approx [U+0049] + [U+0056]$ (roman numeral four) \approx (uppercase "i") + (uppercase "v")

$\mathbf{fi} \approx \mathbf{f} + \mathbf{i}$

 $[U+FB01] \approx [U+0066] + [U+0069]$ (ligature "fi") \approx (lowercase "f") (lowercase "i")

$f \approx S$

 $[U+017F] \approx [U+0073]$ ("long s") \approx (lowercase "s")

canonical vs. compatible is a key to Unicode!

many forms of compatibility...

Compatibility

- "standard", denoted by <compat>
- <sub>, e.g., F₂ (U+2082)
- <super>, e.g., 2⁴ (U+2072)
- <circle>, e.g., (8) (U+2467)
- <fraction>, e.g., ³/₄ (U+00BE)
- and more!

Topic #4: Decomposition

two kinds of decomposition...

canonical decomposition

compatibility decomposition

decomposition can take more than one step...

$\mathring{A} \equiv \mathring{A} \equiv A + \circ$

[U+2I2B] = [U+00C5] = [U+004I] + [U+030A]

in decomposition, order matters!

$\tilde{\psi} \equiv \tilde{\psi} +$

[U+IFA7] = [U+IF67] + [U+0345]

$\tilde{\omega} = \dot{\omega} + \tilde{\omega}$

[U+IF67] = [U+IF61] + [U+0342]

$\dot{\omega} \equiv \omega + \dot{\omega}$

[U+IF6I] = [U+03C9] + [U+03I4]

$$\tilde{\psi} \equiv \omega + ' + ~ + '$$

[U+IFA7] = [U+03C9] + [U+03I4] + [U+0342] + [U+0345]

full decomposition can have both canonical and compatibility steps...

how "aggressive" do we want to be?

İ ≡ **I** + .

 $[U+IE9B] \equiv [U+0I7F] + [U+0307]$

(this is canonical equivalence)

$f \approx S$

$[\mathsf{U}\texttt{+}\mathsf{0}\mathsf{I}\mathsf{7}\mathsf{F}]\approx[\mathsf{U}\texttt{+}\mathsf{0}\mathsf{0}\mathsf{7}\mathsf{3}]$

(this is compatibility equivalence)

$$f \simeq S + \bullet$$

[U+2163] \simeq [U+0073] + [U+0307]

(full decomposition leads to a strange result)

some characters don't decompose as we might expect...

æ ≄ ae

[U+00E6] ≄ [U+0061] [U+0065]

(is this purely an æsthetic issue?;-)

Topic #5: (Re-)Composition

after a character is decomposed, we can put it back together

recomposition returns a composite character (well, usually)

output depends on which decomposition we used

i.e., canonical, compatibility, or both?

how "aggressive" were we in decomposition?

|V => |V|

[U+2163] => [U+0049] + [U+0056] (roman numeral four) => (uppercase "i") + (uppercase "v")

e.g., is Henry the Fourth the same as Henry Eye Vee?

fi => f + i

[U+FB01] => [U+0066] + [U+0069] (ligature "fi") => (lowercase "f") (lowercase "i")

[U+IE9B] => [U+IE6I]

İ => **İ**

³/₄ => 3/4

[U+00BE] => [U+0033] [U+2044] [U+0034]

(note: U+2044 is "fraction slash", not solidus!)

[U+2467] => [U+0038]

8 <= 8

Topic #6: Normalization Forms

process for determining equivalence

there are 4 forms of normalization

Normalization Forms D, C, KD, and KC

a.k.a. NFD, NFC, NFKD, NFKC

2 perform only decomposition (NFD and NFKD)

2 perform decomposition and recomposition (NFC and NFKC)

Normalization Forms

- NFD = canonical decomposition
- NFKD = canonical and compatibility decomposition ("K" is for compatibility!)
- NFC = canonical decomposition, then recomposition
- NFKC = canonical and compatibility decomposition, then recomposition

NFD

- Applies canonical equivalence rules only
- Performs decomposition only
- Does not return a composite character (usually)
- Can result in faster processing (no compatibility, no recomposition)
- The simplest of the normalization forms

NFKD

- Applies canonical equivalence rules and compatibility equivalence rules
- Performs decomposition only
- Does not return a composite character
- Can result in slower processing than NFD (compatibility, but still no recomposition)
- More Clever[™] than NFD

NFD vs. NFKD

input	NFD	NFKD
fi	fi	f + i
ŕ	ſ+'	s + .
IV	IV	I + V
3⁄4	3⁄4	3 + / + 4
8	8	8
2 ⁵	2 + 5	2 + 5
ယ့်	_ + ' + ~ +	+ `+ `+

NFC

- Applies canonical equivalence rules only (first decomposition, then recomposition)
- Compared to NFKC:
 - Produces more matches during comparison operations
 - Requires less time and processing
 - Less Clever[™] (but smarter than NFD)

NFKC

- Applies canonical equivalence rules and then compatibility equivalence rules (first decomposition, then recomposition)
- Compared to NFC:
 - Produces more false negatives
 - Requires more time and processing
 - It's Really Clever™

NFC vs. NFKC

input	NFC	NFKC
fi	fi	fi
ŕ	Ì	Ś
IV	IV	IV
3⁄4	3⁄4	3 / 4
8	8	8
2 ⁵	2 ⁵	2 5
ယို	ယို	ယ်ို

which normalization form is best?

it depends on what you want to accomplish :)

NFC is generally recommended

[RFC 5198]

think long and hard about using something other than NFC

Topic #7: Encoding

Unicode is not an Internet technology

Unicode is not even a computing technology

nothing we've talked about yet has anything to do with computers

a code point just identifies a character

a character could be written, spoken, etc.

computers need characters to be encoded as bits and bytes

e.g., ASCII was originally a 7-bit system (2⁷ gives us 128 code points)

8-bit ASCII gives us $2^8 = 256$ code points

Unicode has many more code points, we need fancier encodings

UTF-8, UTF-16, UTF-32

UTF-8 (RFC 3629) is the IETF-preferred encoding (RFC 2277 / RFC 5198)

each character is encoded using I-4 8-bit "octets"

for the ASCII range, UTF-8 preserves the old 7-bit assignments

e.g., P = ASCII decimal code 80 (i.e., UTF-8 hex code 50)

for characters above decimal code 128, we need 2+ 8-bit "octets"

most modern characters are encoded with two or three octets (up to U+FFFF)

a.k.a. the "Basic Multilingual Plane" (BMP)

higher planes are available (a.k.a. the "astral planes")

e.g., $\frac{1}{2} = U + 10133$

[i.e., AEGEAN NUMBER NINETY THOUSAND]

however, these are unlikely to be used on the Internet

although UTF-8 is very common, there are exceptions...

especially the internal data representation in Java, JavaScript, and Windows

[note: some systems insert a "byte order mark" (BOM) at the start of UTF-8 data]

think long and hard before using something other than UTF-8

Topic #8: String Comparison

some strings are special

e.g., addresses and other identifiers

many reasons to compare strings...

authentication

authorization

registration

data storage

and many other operations

first attempt: stringprep (~2002)

designed for domain names ("IDNA")

applied to many other identifier types

addresses, usernames, passwords, file paths, nicknames, etc....

each has different uses, needs, and structure

Stringprep Basics (I)

- Choose a Unicode version (oops, 3.2 only!)
- Choose a normalization form (NFKC?!)
- Specify how to handle whitespace
- Specify whether to use case folding
- Specify bidirectional handling
- Specify prohibited characters

Stringprep Basics (2)

- Specify handling via comprehensive tables that capture:
 - Mappings (e.g., whitespace, case folding)
 - Prohibited characters (e.g., controls, spaces, symbols)
 - Bidirectionality

Why Not Stringprep?

- Version agility is important (latest = 6.3)
- NFKC can lead to unintuitive results, as we've seen (e.g., $f \approx s$)
- Accepting registration of all characters and scripts can cause problems (e.g., phishing)
- Big tables are hard to maintain and update
- See RFC 4690 for details

IDNA2008 (I)

- No more stringprep
- Decisions based on properties of Unicode characters
- Algorithms, not huge tables
- Version agility

IDNA2008 (2)

- Four "buckets" based on properties:
 - PROTOCOL-VALID
 - CONTEXT RULE REQUIRED
 - DISALLOWED
 - UNASSIGNED

IDNA2008 (3)

- Basically, PVALID = "letter-digit-hyphen"
- The "inclusion approach" of IDNA2008 works because domain names have always traditionally been "letter-digit-hyphen", not just any random symbols, punctuation, etc.
- Domain names are mnemonics, not random strings of characters

Challenges

- The dividing line between user interface and protocol has moved substantially
- Applications need to take more responsibility
- Can't just hand things off to stringprep and expect good things to happen
- Mappings are out of scope for IDNA2008

Stringprep Customers

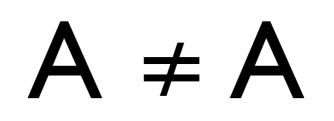
- IDNA was the main stringprep "customer"
- Other customers: LDAP, SASL, iSCSI, XMPP, etc.
- A new approach developed in the PRECIS WG ("Preparation and Comparison of Internationalized Strings")

PRECIS

- Follow "inclusion approach" like IDNA
- Define two "string classes" (IdentifierClass, FreeformClass)
- Enable "profiling" for particular protocols (case mapping, normalization, etc.)
- draft-ietf-precis-framework etc.

Topic #9: Confusable Characters

many characters look alike ("confusables")...



[U+0041] ≠ [U+0410]

4 ≠ **4**

[U+I3CE] ≠ [U+0034]

199

[U+I3DA] [U+I3A2] [U+I3B5] [U+I3CB] [U+I3A2] [U+I3CB] [U+I3D2] ≠ [U+0053] [U+0054] [U+0050] [U+0045] [U+0054] [U+0045] [U+0052]

$STPETER \neq STPETER$

these are *not* equivalents

humans usually can't distinguish

4 vs. 4 looks like a font difference

no foolproof solutions for confusables

prohibiting mixed scripts can help

see RFC 4690/5890 and draft-ietf-precis-framework

Topic #10: Rules and Responsibilities

Who's Your Registrar?

- In IDNA, domain registrars have policies
- E.g., Hungarian registrar likely won't accept characters from Korean code block
- Do providers of (say) email and IM services also need to define such policies?

Enforcement

- Who enforces the rules?
 - Server?
 - Client?
 - Any network endpoint?
 - Needs to be clear for each protocol!

Topic #11: Versioning

Version Changes

- New Unicode versions can add new characters, deprecate old characters, etc.
- Character properties can change between Unicode versions (e.g., from number to letter), but this should be rare
- A character could change from PVALID to DISALLOWED (etc.)

Version Mismatches

- Possibility of problems with authentication, message delivery, etc.
- In practice, not a concern because most of the modern characters we need are mapped in a stable way

Topic #12: User Interface

Good UI is Hard

- Account for application type, string types, locale, scripts, culture, input methods, output methods, graphical capabilities, etc.
- Probably not much that protocol geeks can say in the matter :)
- We need input from UI experts

Rendering

- UTF-8 encoded Unicode characters are rendered in a UI by a rendering engine
- These have improved dramatically over time! (Fewer "renderings" via □)
- In general, support for Unicode is improving all the time

Topic #13: Directionality

LTR and RTL

- Most scripts are rendered left-to-right
- Some scripts are rendered right-to-left (e.g., Arabic and Hebrew)
- Each Unicode code point is LTR or RTL
- What if they're mixed? Hard problem!

BiDi Policies

- If a string contains any LTR character, the entire string is left-to-right
- BiDi rule from RFC 5893
- Other rules are possible (but there be major dragons here!)

THE END

STPETER @ STPETER.iM

References: Unicode

- Unicode 6.3 spec @ unicode.org
- UAX 15: Unicode Normalization Forms
- UTR 17: Unicode Character Encoding Model
- UTR 36: Unicode Security Considerations

References: General

- RFC 6365: Internationalization terminology
- RFC 2277: IETF policy on characters sets and languages
- RFC 3629: UTF-8
- RFC 5137:ASCII escaping for Unicode
- RFC 5198: Unicode format for networks

References: Stringprep & IDNA2003

- RFC 3454: Stringprep
- RFC 3490: IDNA2003
- RFC 3491: Nameprep
- RFC 3492: Punycode
- RFC 4690: IDN review by IAB

References: IDNA2008

- RFC 5890: Definitions
- RFC 5891: Protocol
- RFC 5892: Unicode and IDNA
- RFC 5893: Right-to-Left Scripts
- RFC 5894: Background

References: PRECIS

- http://datatracker.ietf.org/wg/precis/
- RFC 6885
- draft-ietf-precis-framework
- draft-ietf-precis-mappings
- PRECIS-related I-Ds on usernames, passwords, nicknames, JabberIDs, etc.

Useful Tools and Websites

- Unicode Checker (Mac OS X)
- unicode-table.com
- Wikipedia pages about Unicode, UTF-8, and related topics

Acknowledgements

- Martin Dürst
- Joe Hildebrand
- John Klensin
- PRECISWG
- XMPP community