

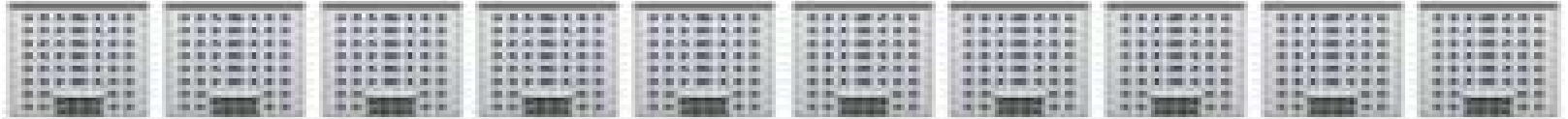
Unicode & Emoji

Mark Davis

President & Co-founder, Unicode Consortium

[@mark_e_davis](https://twitter.com/mark_e_davis)

<https://goo.gl/Lff11S>



What is the Unicode Consortium?

*Enable everybody,
speaking every language on the Earth,
to be able to use their language
on computers and smartphones*

Levels of Language Support

Core	Characters, Fonts, Keyboards
Content	Sorting, Numbers, Dates, Units, Currency, ...
UI	Translated UI, Help text, ...
NLP	Predictive typing, Voice input / output, OCR, Entity recognition, Handwriting input, ...

Full Members (Voting)



Institutional Members (Voting)



Supporting Members (Voting)



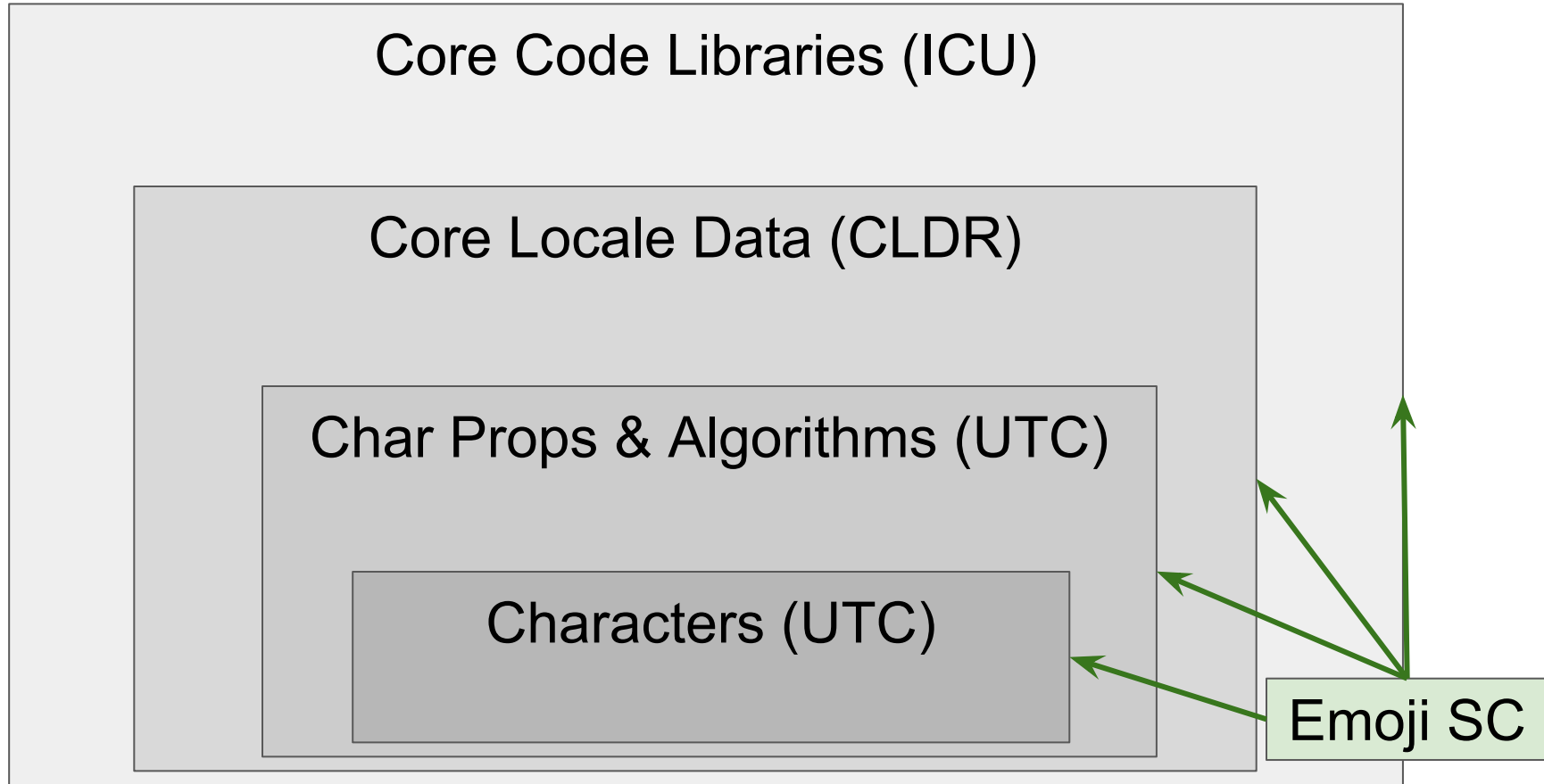
Associate Members

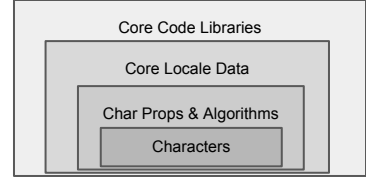


member-supported, non-profit



Unicode Consortium

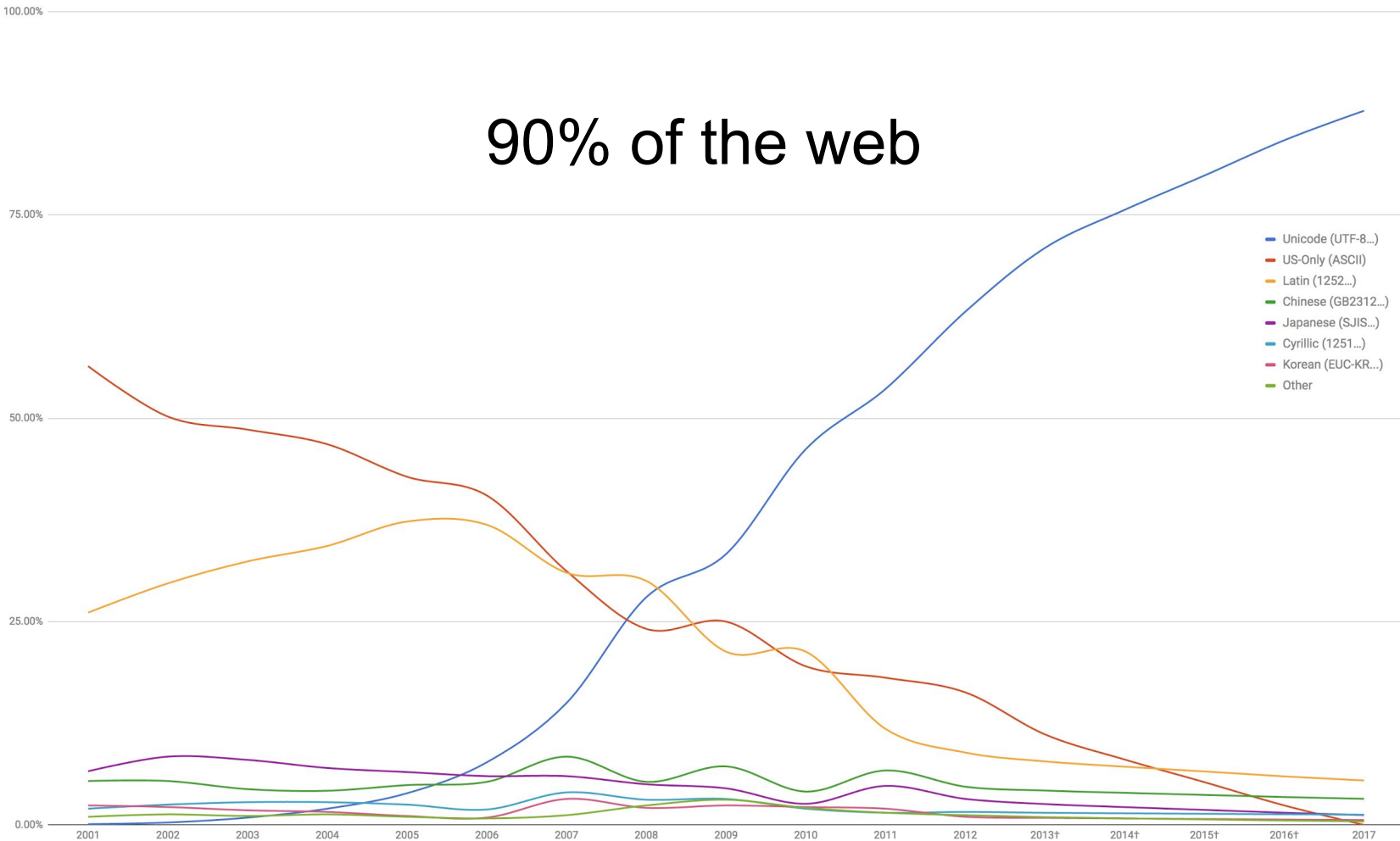




What is a Unicode Character?

*On a laptop, server, or mobile phone
every character you type,
every character you see
is Unicode*

90% of the web



googleblog.blogspot.ch/2012/02/unicode-over-60-percent-of-web.html

— graph updated to 2017

Not as simple as you think

Combined
glyphs



fish

Individual
glyphs



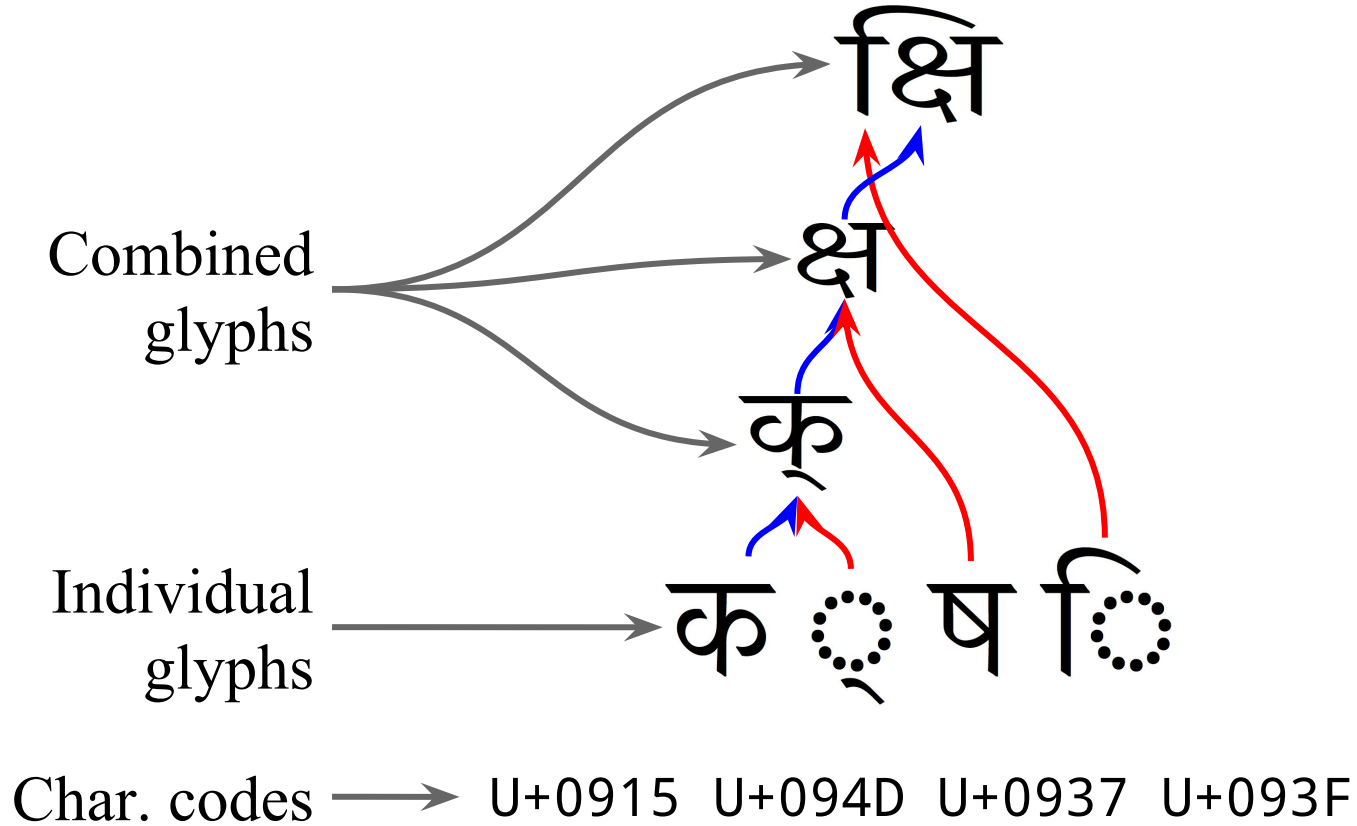
f i s h

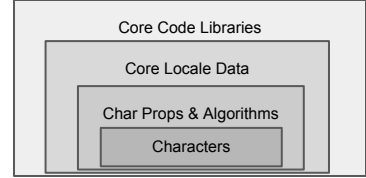
Char. codes



U+0066 U+0069 U+0073 U+0068

Not as simple as you think





Why Properties?

Full name

Zoë Straub|

Illegal Name

Before

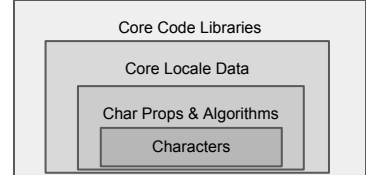
```
if  
    'a' ≤ X AND X ≤ 'z'  
then handle(X)
```

After

```
If  
    isLowercase(X)  
then handle(X)
```


Unicode Properties

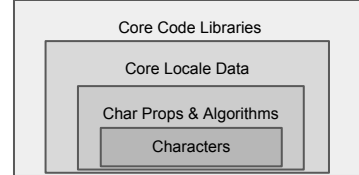
General	Case	Normalization	Shaping and Rendering
Name	Uppercase	Canonical_Combining_Class	Join_Control
Name_Alias	Lowercase	Decomposition_Mapping	Joining_Group
Block	Lowercase_Mapping	Composition_Exclusion	Joining_Type
Age	Titlecase_Mapping	Full_Composition_Exclusion	Line_Break
General_Category	Uppercase_Mapping	Decomposition_Type	Grapheme_Cluster_Break
Script	Case_Folding	NFC_Quick_Check	Sentence_Break
Script_Extensions	Simple_Lowercase_Mapping	NFKC_Quick_Check	Word_Break
White_Space	Simple_Titlecase_Mapping	NFD_Quick_Check	East_Asian_Width
Alphabetic	Simple_Uppercase_Mapping	NFKD_Quick_Check	Prepended_Concatenation_Mark
Hangul_Syllable_Type	Simple_Case_Folding	NFKC_Casefold	Bidirectional
Noncharacter_Code_Point	Soft_Dotted	Changes_When_NFKC_Casefolded	Bidi_Class
Default_Ignorable_Code_Point	Cased	Miscellaneous	Bidi_Control
Deprecated	Case_Ignorable	Math	Bidi_Mirrored
Logical_Order_Exception	Changes_When_Lowercased	Quotation_Mark	Bidi_Mirroring_Glyph
Variation_Selector	Changes_When_Uppercased	Dash	Bidi_Paired_Bracket
Identifiers	Changes_When_Titlecased	Sentence_Terminal	Bidi_Paired_Bracket_Type
ID_Continue	Changes_When_Casefolded	Terminal_Punctuation	CJK
ID_Start	Changes_When_Casemapped	Diacritic	Ideographic
XID_Continue	Numeric	Extender	Unified_Ideograph
XID_Start	Numeric_Value	Grapheme_Base	Radical
Pattern_Syntax	Numeric_Type	Grapheme_Extend	IDS_Binary_Operator
Pattern_White_Space	Hex_Digit	Indic_Positional_Category	IDS_Tertiary_Operator
	ASCII_Hex_Digit	Indic_Syllabic_Category	Unicode_Radical_Stroke



Unicode Algorithms

Unicode detection, conversion	Script detection
Normalization, equivalence	Case detection, conversion
Collation (sorting)	Segmentation
Regular expressions	Unicode domain names
Bidirectional text	Security mechanisms
Identifier parsing	Emoji validity
Shaping (Arabic,...)	Vertical orientation (CJK)
...	

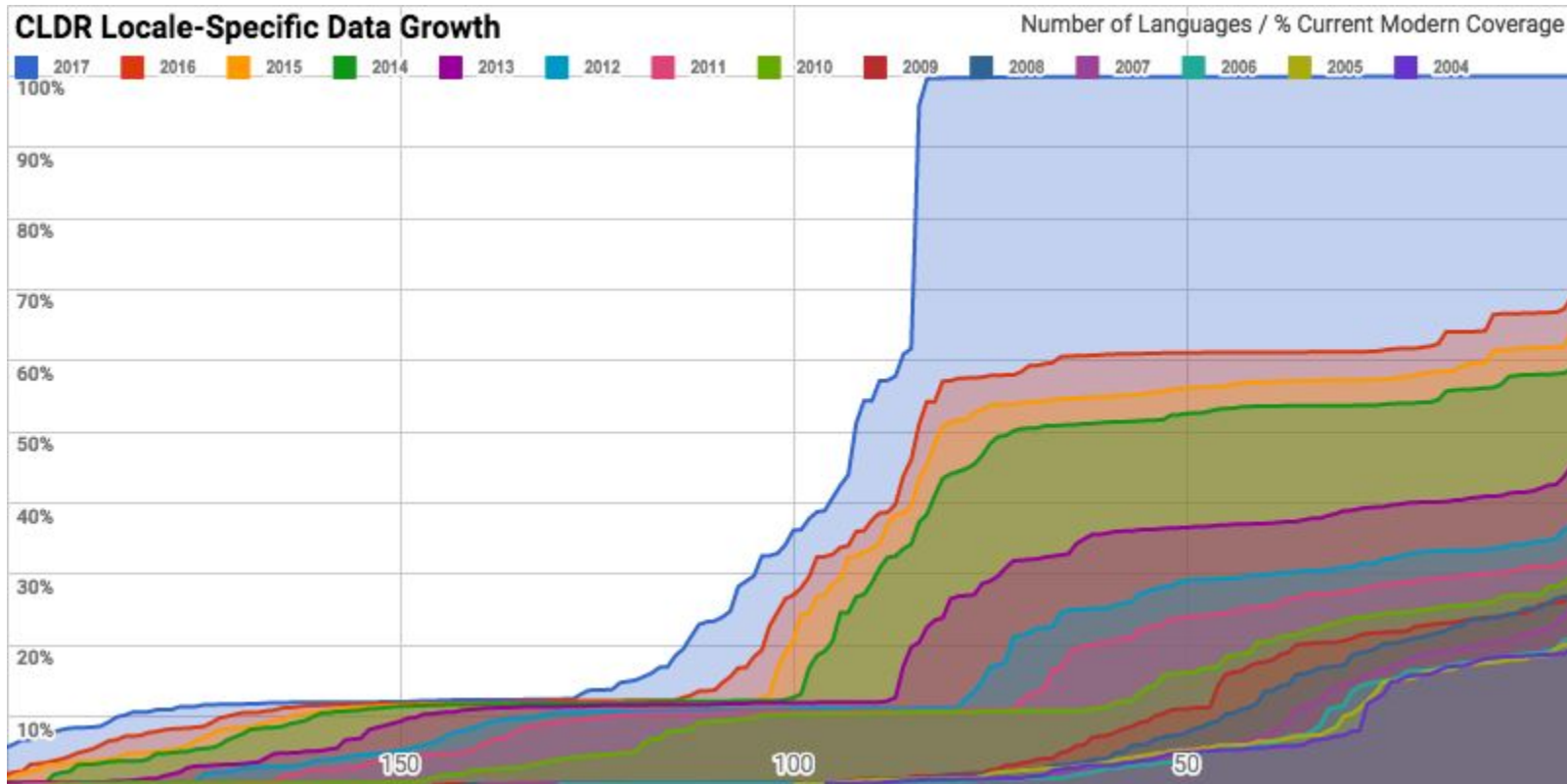
Language-Specific Data (CLDR)













Characters:	ä, ø,...
Sorting	ä < b vs ä > z
Dates & times	Dec 3–5; last Tuesday
Day periods, timezones	8:00 → morning1; 13:52 Italy Time
Numbers, units, currencies	1.234; \$3.5M; 123 Kanadische Dollar
Lists, ranges, compact forms	Dec 3–5; 3.5–4.2kg; 3.5B
Names of languages, scripts,...	FR → 法语; Cyril → Kyrillisch; 419 → Südamerika
Character labels, names,...	☞ → Arrows; ♉ → Taurus; bull ox zodiac
Transforms	Путин → Putin; Трамп → Trump
Spell out numbers	25 → twenty-five
Plurals (cardinals, ordinals)	1 book, 2 books; 1st book; 1–2 books
Keyboards	C03 → д; C03+caps → Д

Locale matching/normalization. Region info (currencies, containment, ...) ...

CLDR Growth

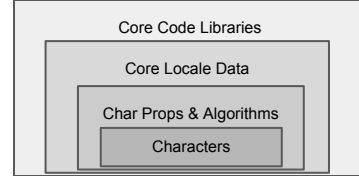


Survey Tool

 -name	panda face	✓	Pandagesicht
 -name	paw prints	✓	Tatzenabdrücke
Animal-Bird			
 -name	turkey	✓	Truthahn
 -name	chicken	✓	Huhn
 -name	rooster	✓	Hahn
 -name	hatching chick	✓	schlüpfendes Küken
 -name	baby chick	✓	Küken
 -name	front-facing baby chick	✓	Küken in Frontansicht
 -name	bird	✓	Vogel
 -name	penguin	✓	Pinguin

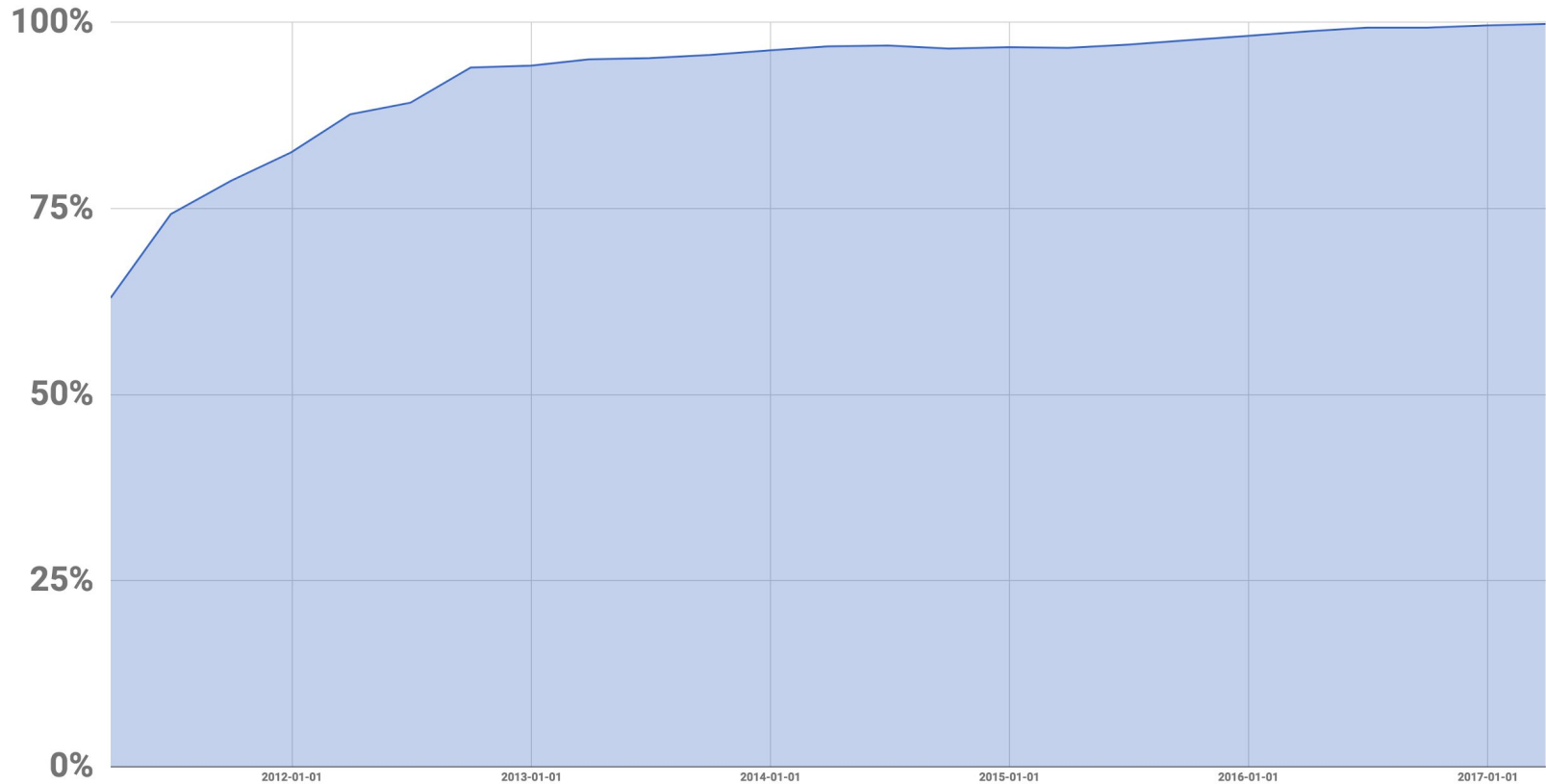
BTW, Unicode looking for contractors: FE & Performance

Core Code Libraries (ICU)



Unicode text handling	Formatting
Charset	Date & time
Conversions (200+)	Durations, intervals
Detection	Messages
Collation & Searching	Numbers, currencies
Resource Bundles	Measurement units
Calendar & Time Zones	Plurals
Unicode Regular Expressions	Transforms
Boundaries: word, line, ...	Normalization
	Casing
	Transliterations
	...

ICU on Smartphones, Market Share



IDC: Worldwide smartphone shipments — 2011Q1..2017Q1

In a nutshell

- Mature, widely used set of C/C++ and Java libraries
- Identical results on all platforms/programming languages
 - C/C++ (ICU4C): 30+ platforms/compilers
 - Java (ICU4J): Oracle, IBM JRE, Android
 - PyICU & other wrappers
- Customizable & Modular

Oxford Dictionaries

**WORD
OF THE YEAR**



1999: Japanese Emoji



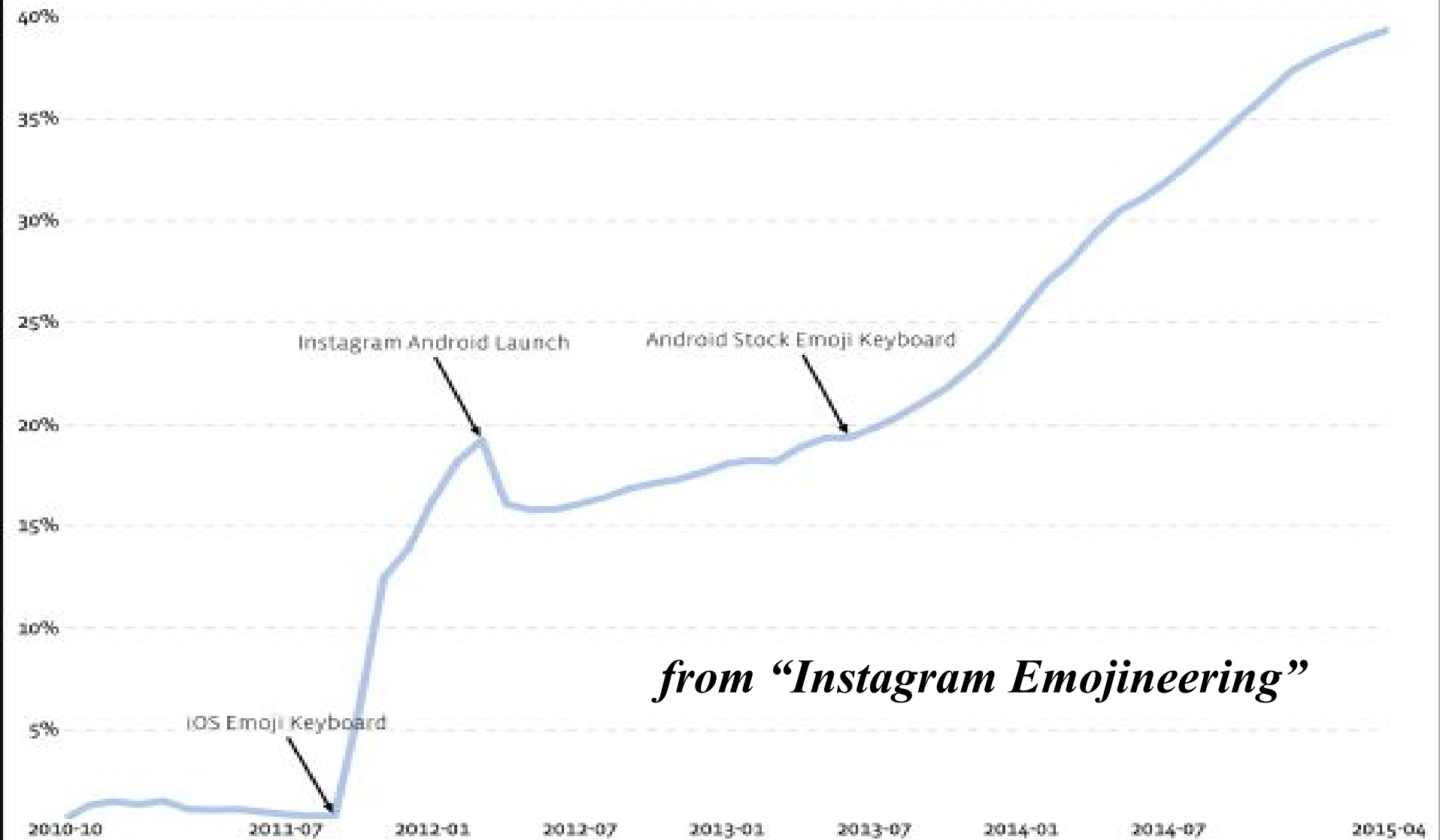
2008: Gmail / iPhone / ...



2010: Unicode emoji



Emoji Usage over Time



2017



[U+1F631](#)

Details: unicode.org/emoji/charts

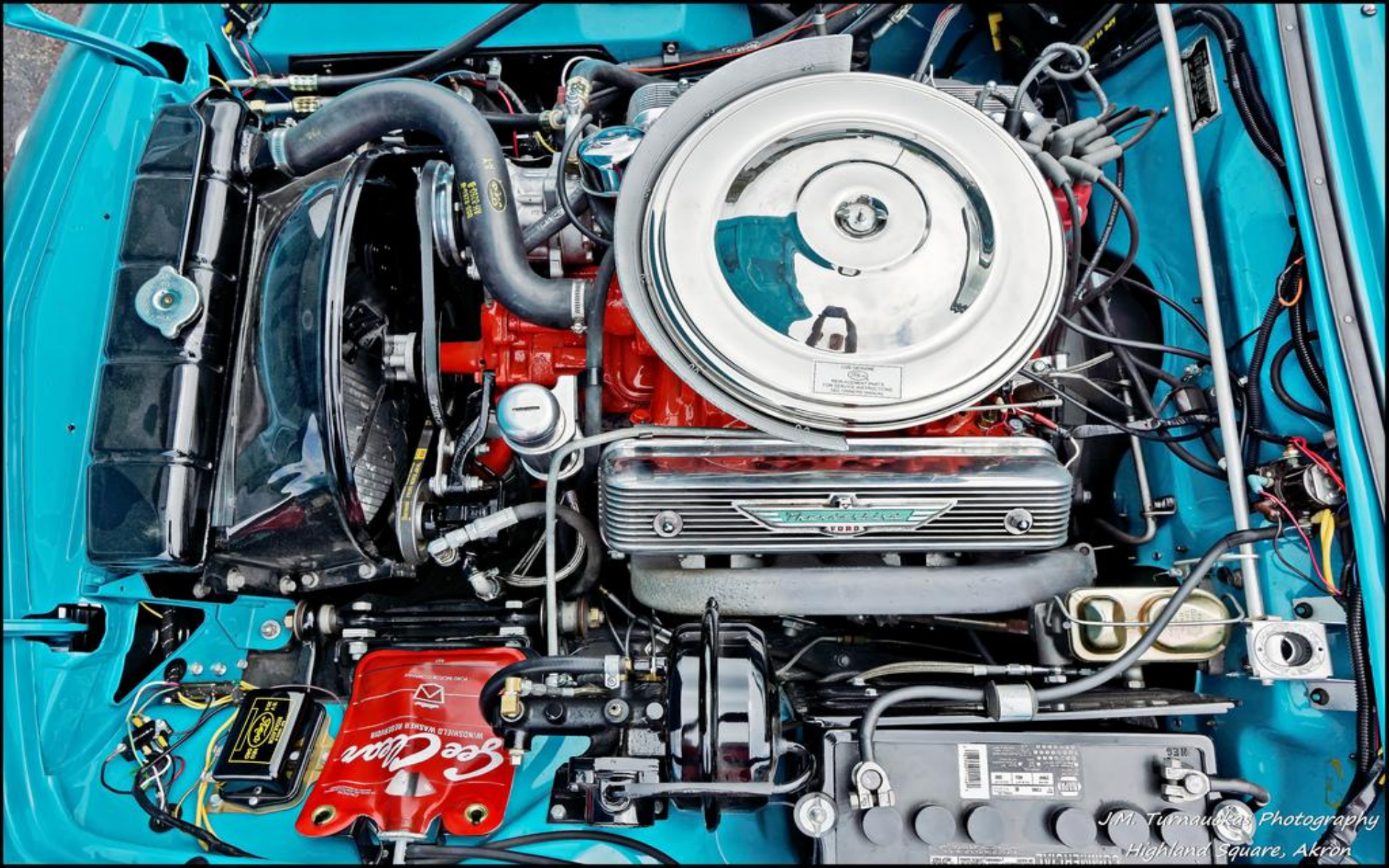
Emoji By-Product

MySQL 8.0 RC1 Now Available

Sep 27, 2017

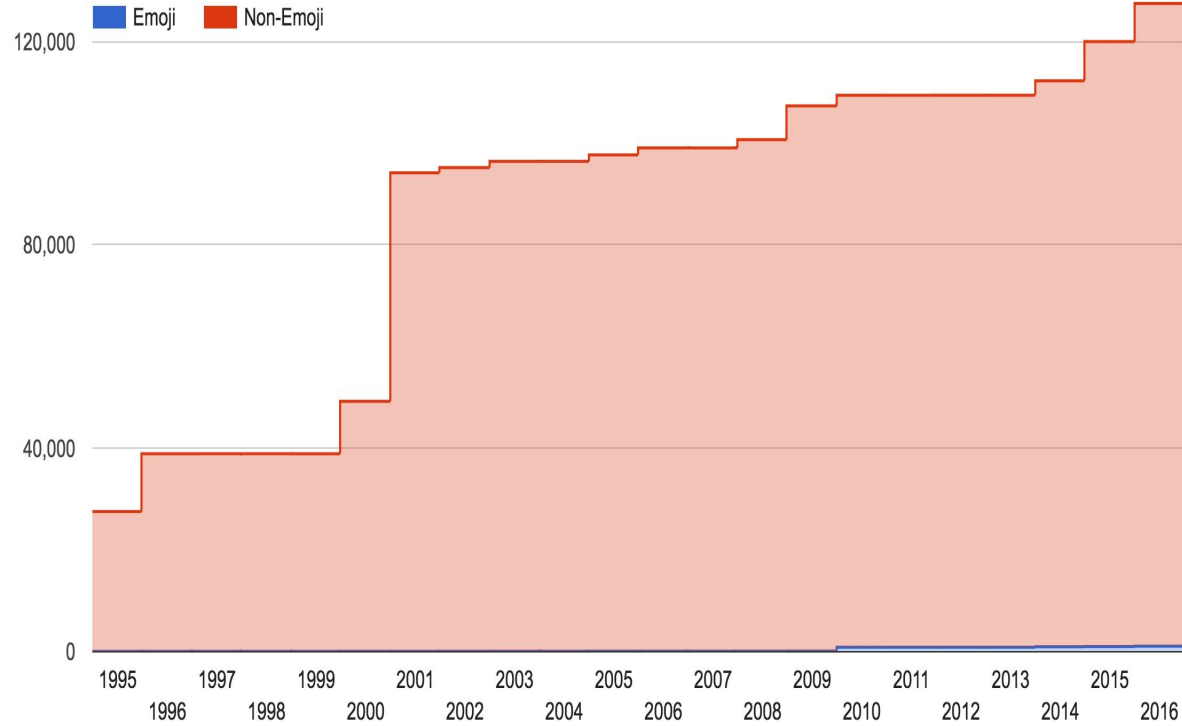
...

In addition, he adds, "Unicode (or more specifically UTF-8 encoding) has become universal even in English speaking markets. A key driver is mobile applications, where **emojis** are frequently used as character input. To support modern applications, it is important to have first-class support for UTF-8 out of the box."



*J.M. Turnauckas Photography
Highland Square, Akron*

Growth of Unicode Characters



How many? (Emoji 5.0)

	Smileys & People	Animals & Nature	Food & Drink	Travel & Places	Activities	Objects	Symbols	Flags	Other	Total
character	297	113	102	207	60	162	193	5	0	1,139
keycap seq	0	0	0	0	0	0	12	0	0	12
flag seq	0	0	0	0	0	0	0	258	0	258
tag seq	0	0	0	0	0	0	0	3	0	3
mod seq	510	0	0	0	0	0	0	0	0	510
zwj seq + gender	43	0	0	0	0	0	0	0	0	43
zwj seq + modifier	160	0	0	0	0	0	0	0	0	160
zwj seq + gender, modifier	195	0	0	0	0	0	0	0	0	195
zwj seq other	61	0	0	0	0	0	0	1	0	62
Subtotal	1,266	113	102	207	60	162	205	267	0	2,382
component	5	0	0	0	0	0	0	0	38	43
typical dup	241	0	0	0	0	0	0	0	0	241
Total	1,512	113	102	207	60	162	205	267	38	2,666

Emoji Properties

Property	Description
Emoji	Emoji characters, also in sequences
Emoji_Presentation	Presented as emoji by default
Emoji_Modifier	Skin tones
Emoji_Modifier_Base	Bases for skin tones
Emoji_Component	For sequences, not typically on keyboard
Extended_Pictographic	CLDR: for proper (and future-proofed) segmentation.

Details: unicode.org/reports/tr51

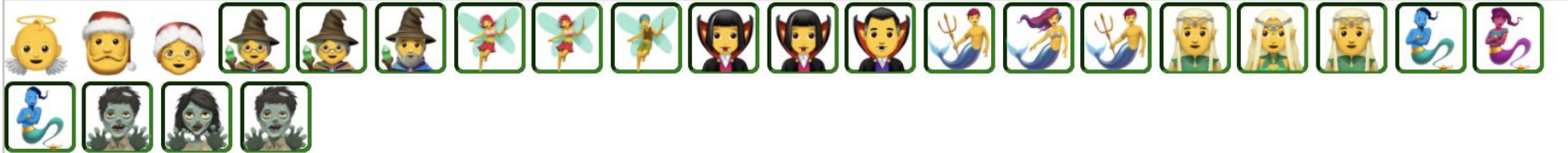
Sort Order (CLDR)



person-role



person-fantasy



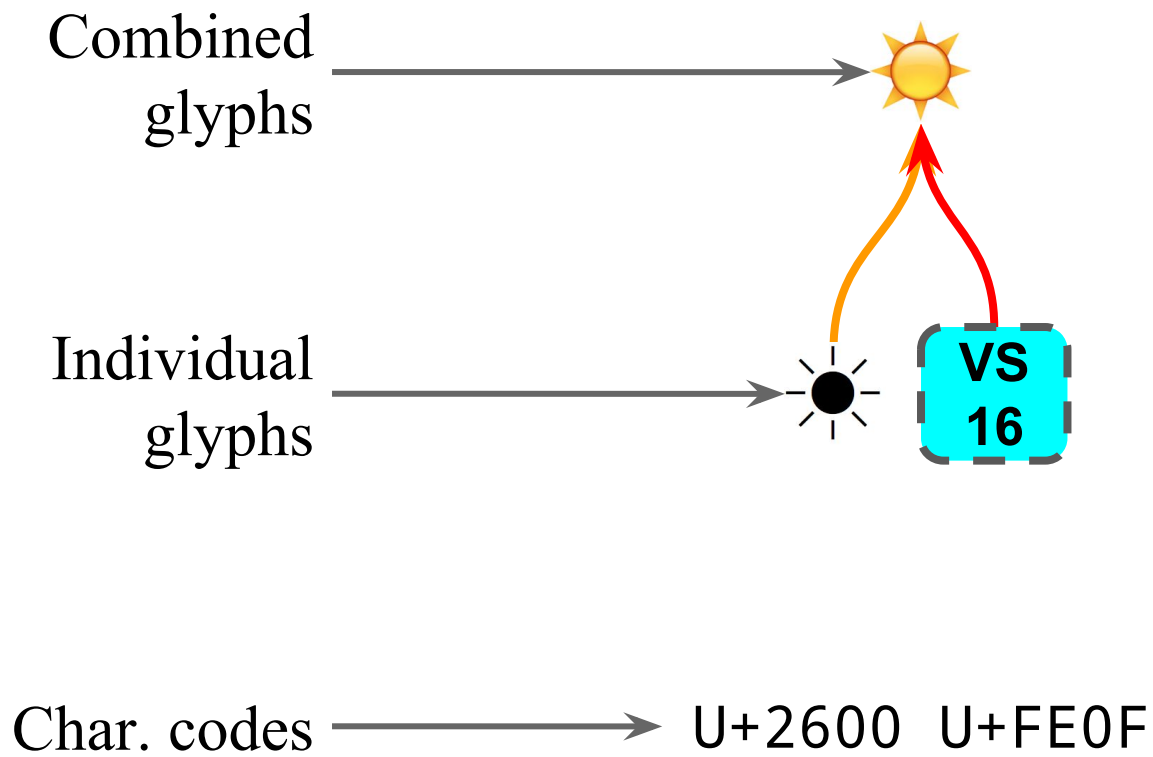
Names, Keywords (CLDR)



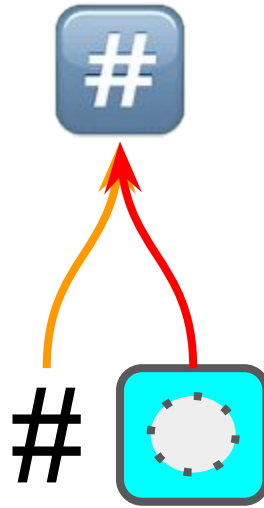
nerd face	nørdansigt	nerderig gezicht	Nerd- Smiley	nördaandlit
face geek nerd	ansigt nørd	geek gezicht nerd	Gesicht Nerd	andlit lúði nørd



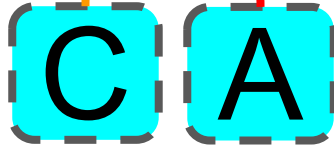
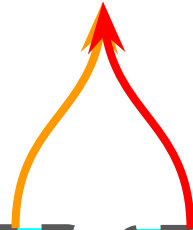
Variation



Keycaps



Flags



III Formed

C



Invalid (but well-formed)

C B



Valid (but not RGI)

C S

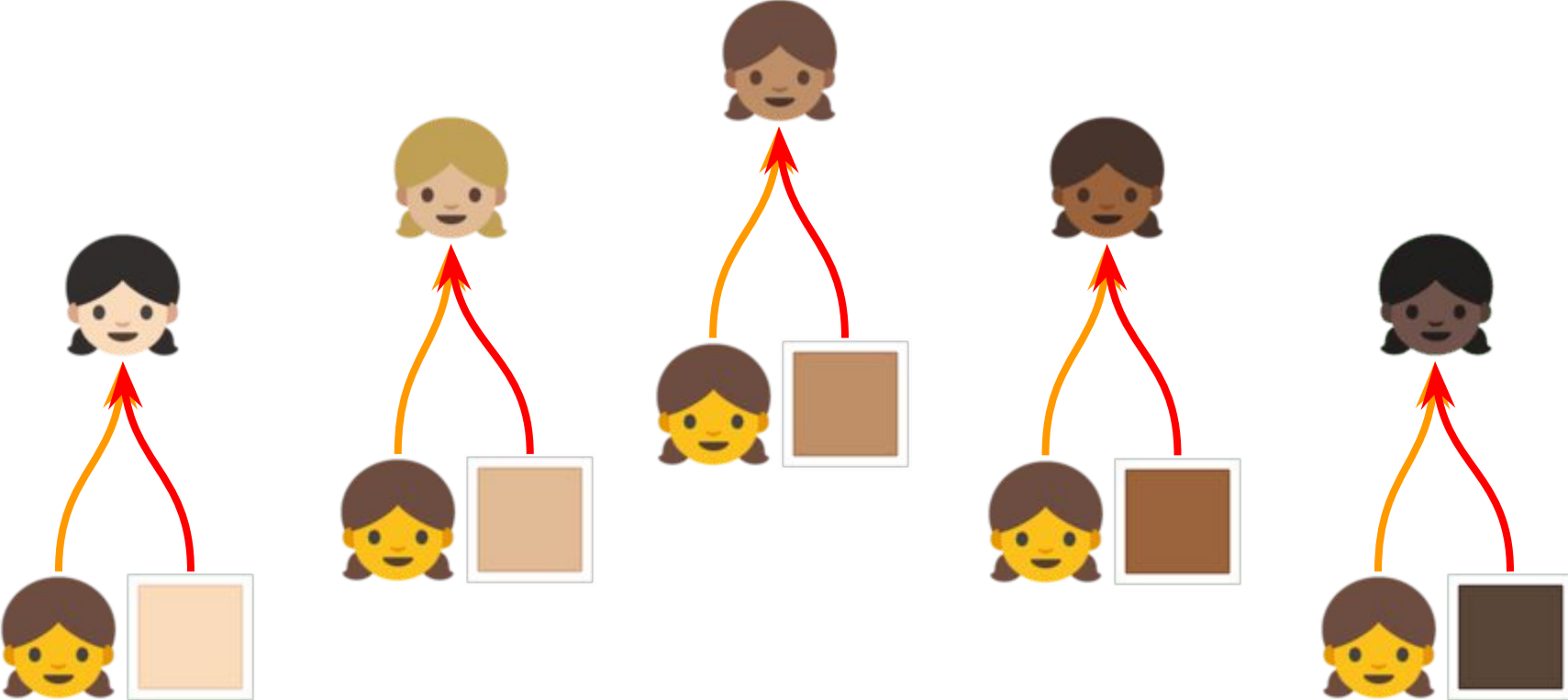


RGI

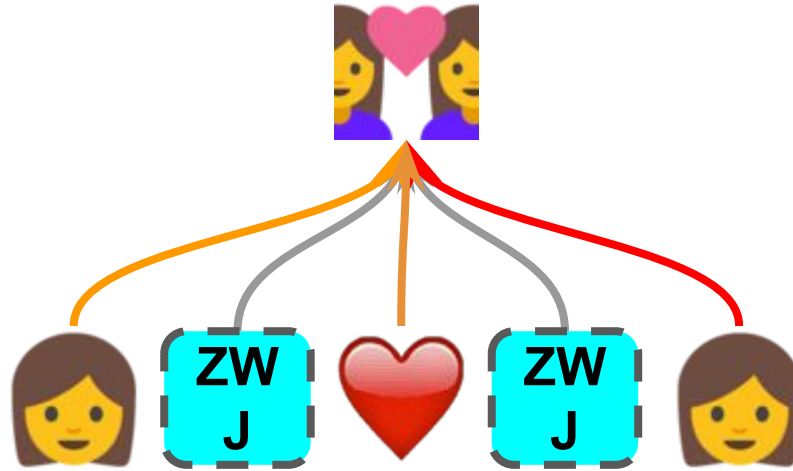
C A



Skin Tones

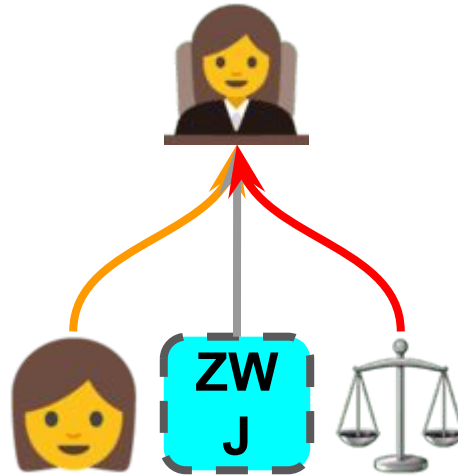


Emoji ZWJ Sequence

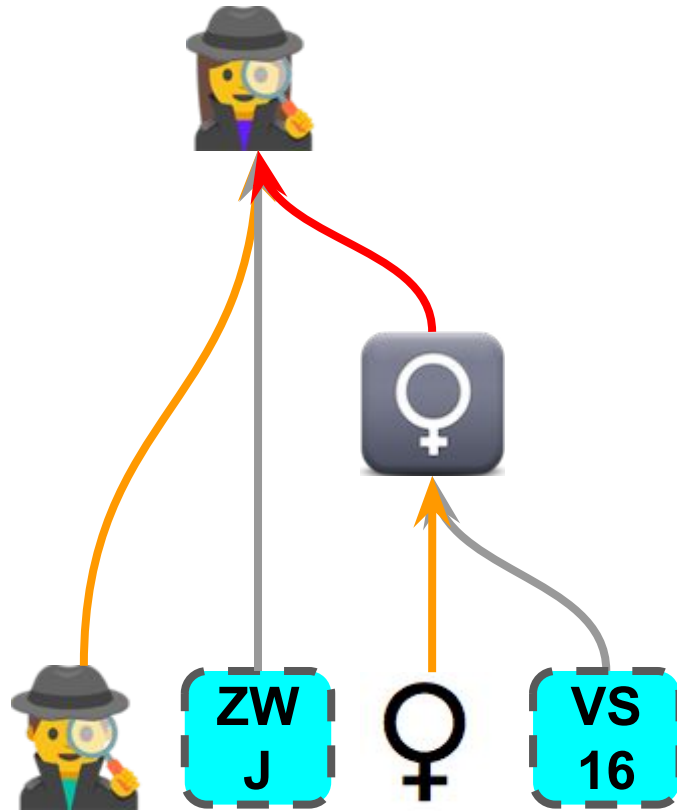




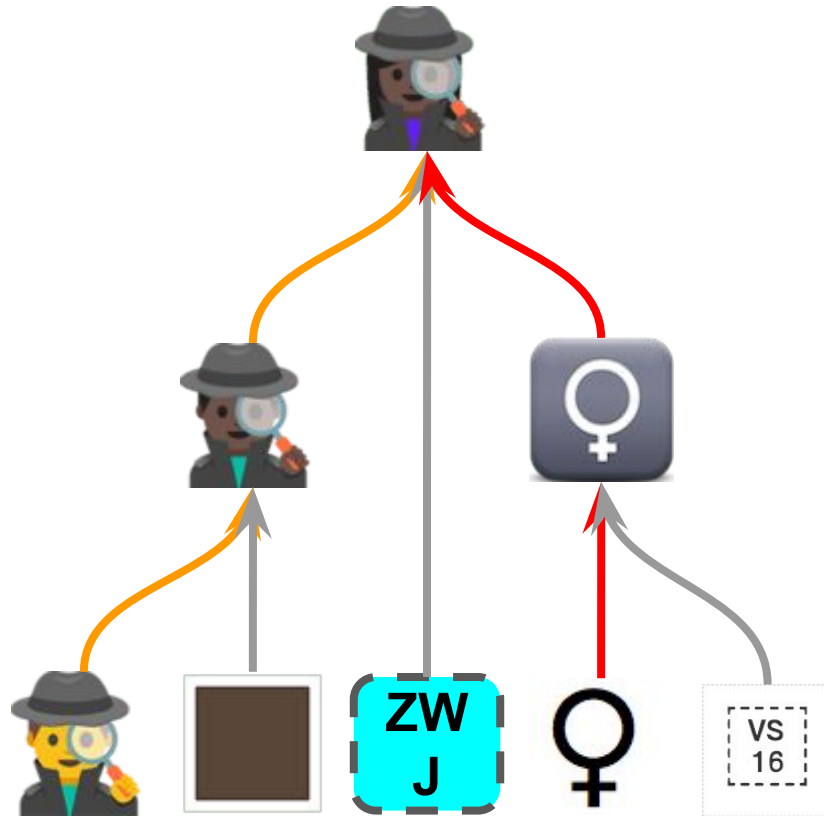
Gendered w/ Object



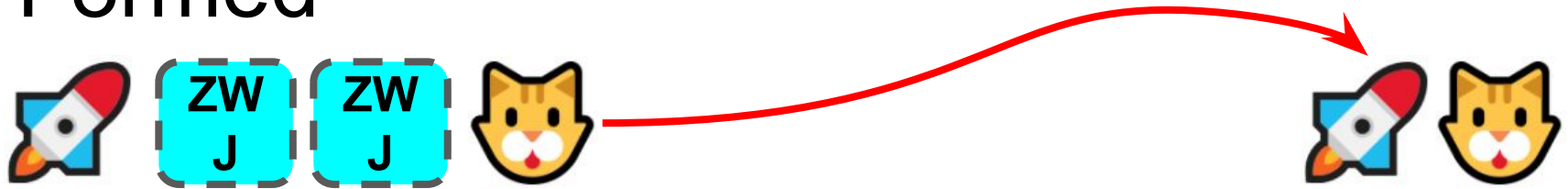
Gendered w/ Sign



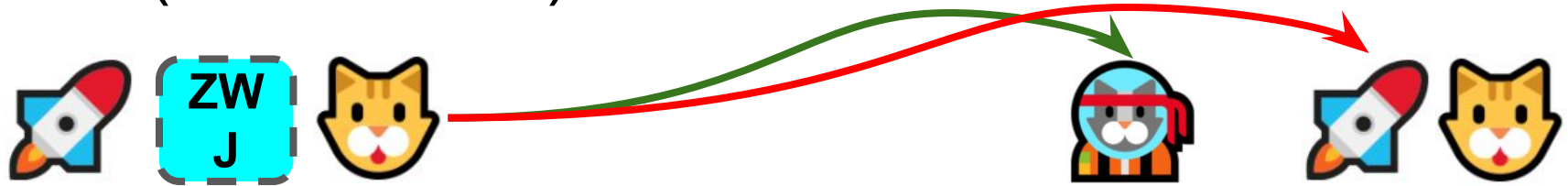
Gendered w/ Skintone & Sign



III Formed



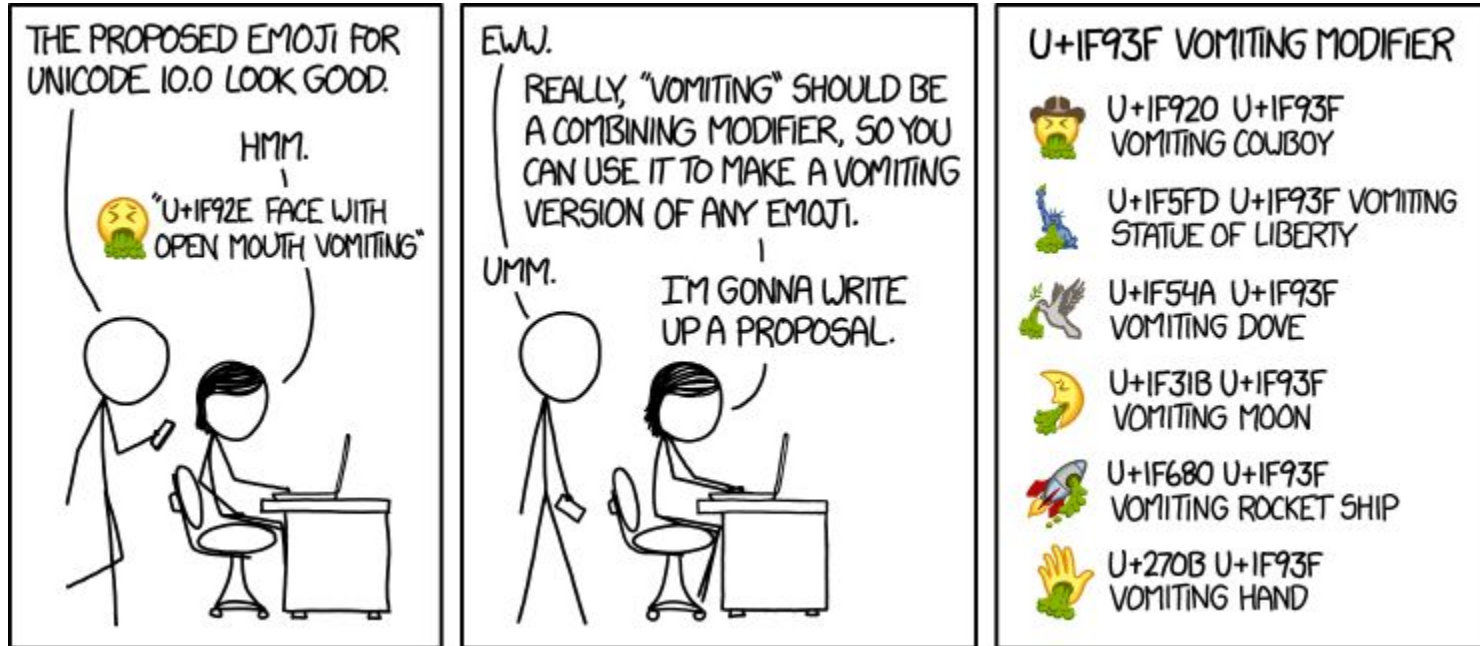
Valid (but not RGI)



RGI



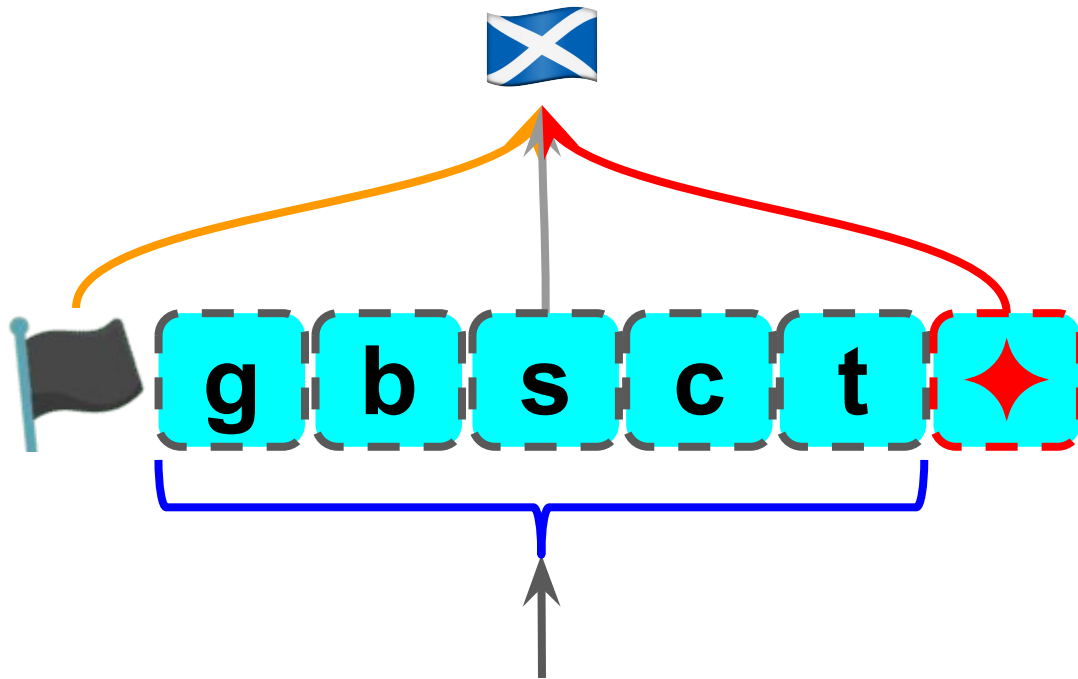
“Valid but not RGI?”



See <https://xkcd.com/1813/>

“Modifier” not needed: **Cowboy + ZWJ + Vomiting Face**

Tag Sequences

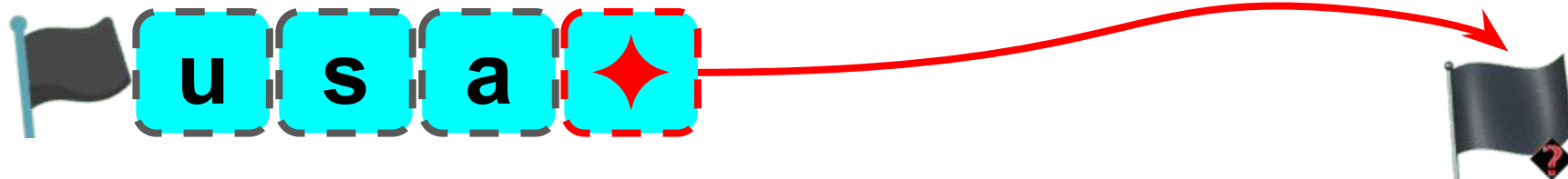


Subdivision = Province, State, Canton,...

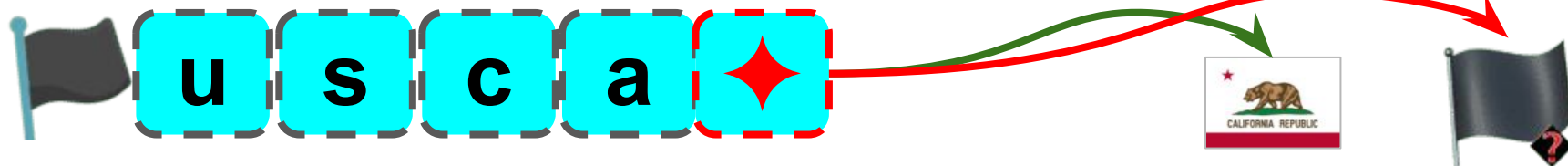
Ill Formed



Invalid (but well formed)



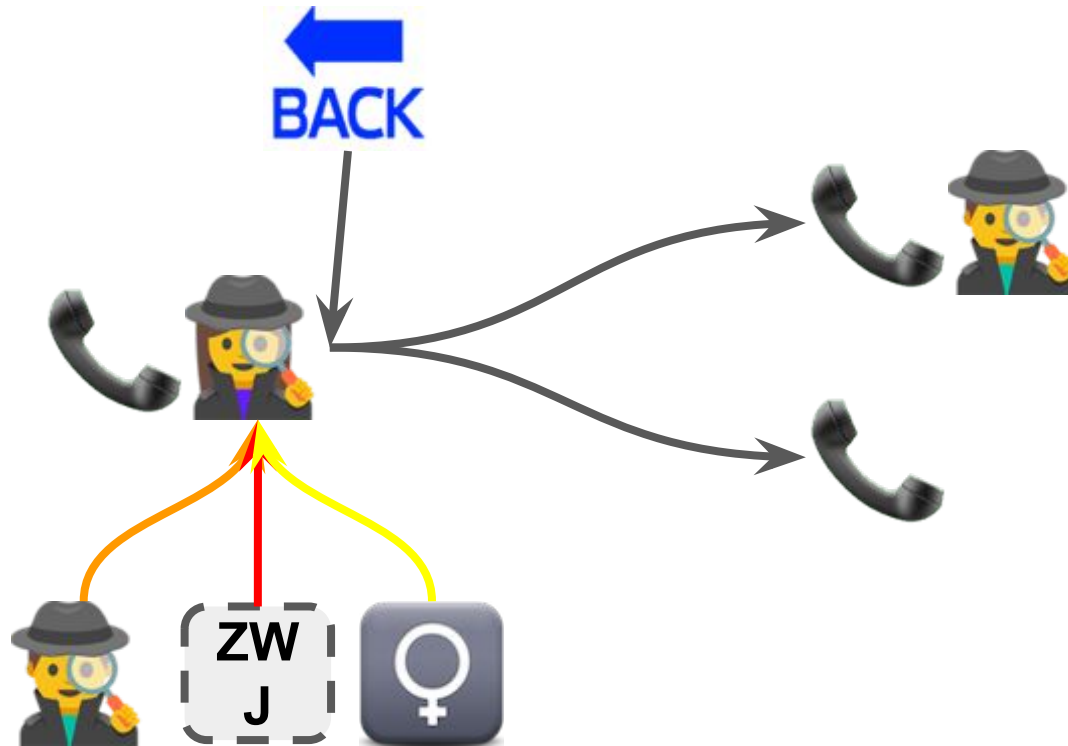
Valid (but not RGI)



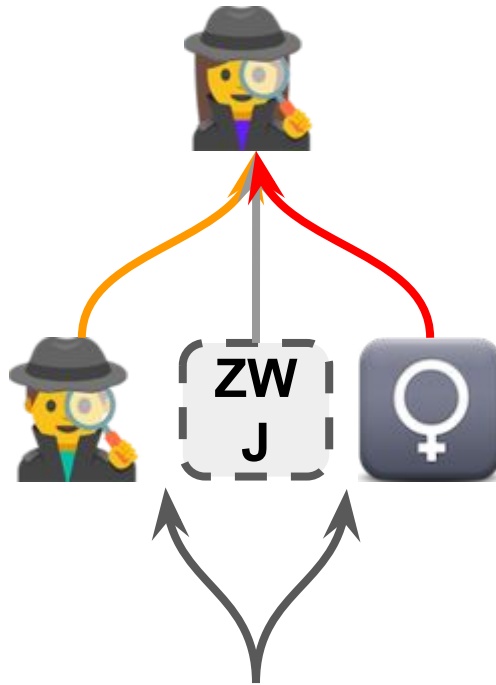
RGI



UI Actions: Backspace Example



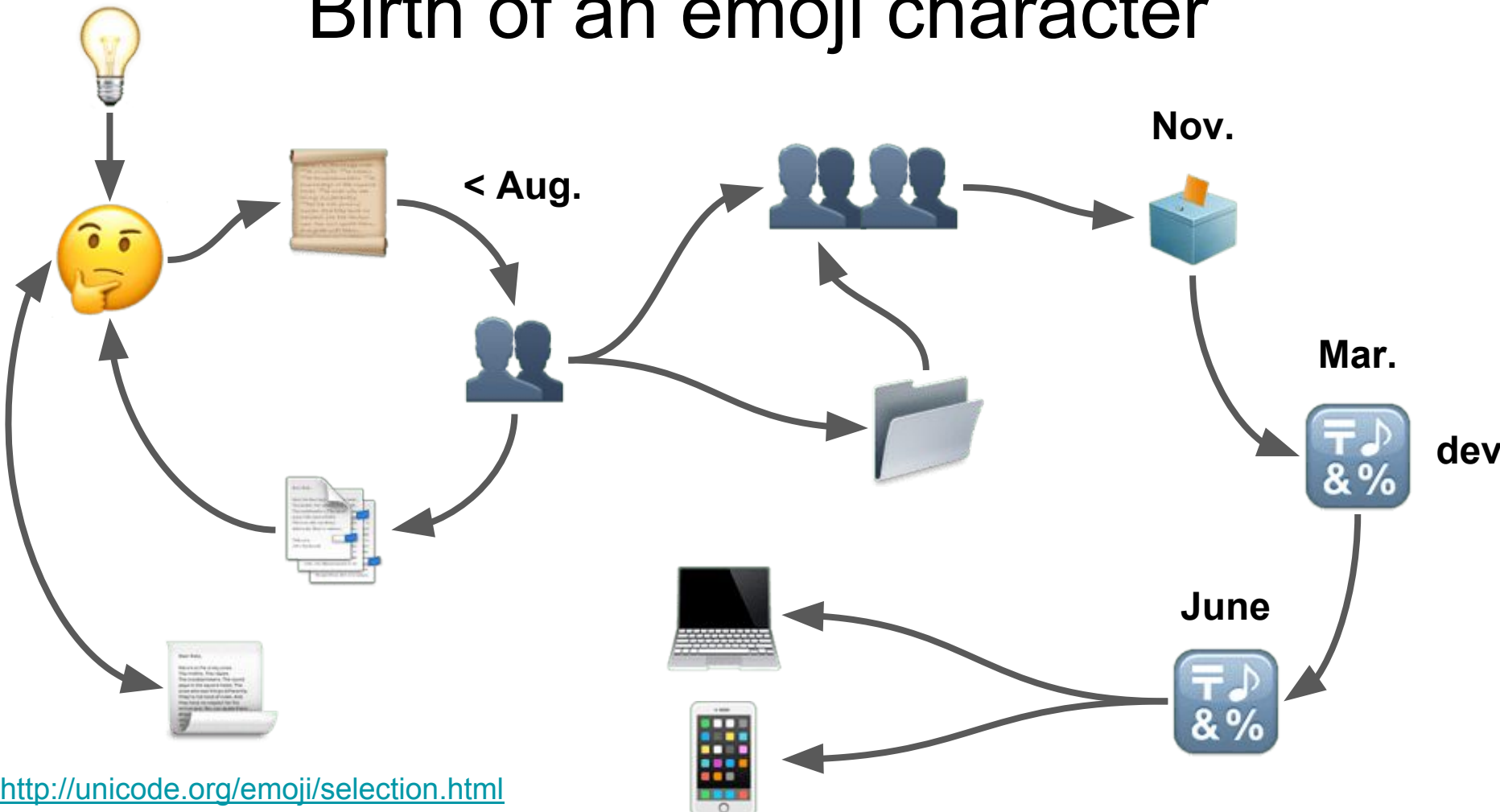
Segmentation



No linebreak here!



Birth of an emoji character



Selection Factors

Factors for Inclusion +

- A. Compatibility
- B. Expected usage level
- C. Image distinctiveness
- D. Completeness
- E. Frequently requested

Factors for Exclusion -

- F. Overly specific
- G. Open-ended
- H. Already Representable
 - I. Logos, brands, ... specific people, deities
- J. Transient
- K. Like compatibility emoji

Proposal for a New Emoji: Cricket

Susan Burtner

<susan.burtner@gmail.com>

September 15, 2016

Abstract

The intention of this proposal is to advocate for the inclusion of emoji, the cricket, as a Unicode emoji character. The symbolic significance of a cricket is multicultural. To some it can be viewed as a symbol of illness, while to others it can be seen as a symbol of prosperity. In forms of American media, the cricket, particularly the *sound* of a chirping, implies an awkward silence. None of the current emoji replace the varied symbolism implied by the cricket, or express the idea of “awkwardness” as clearly as the cricket does. The cricket emoji current design captures all of these meanings while still being visible enough to stand for a normal cricket.

2 Factors for Inclusion

2.1 Compatibility

Unfortunately, no emoji in the Unicode or any other system has created a cricket emoji like this. While this emoji cannot speak to the completeness and compatibility of other systems, it represents a unique opportunity to begin emoji representation for the ideas of “prosperity”, “luck”, and/or “awkwardness”.

- ¹“The Folklore and Mythology Surrounding Crickets”: <http://www.crystalwisdom.com/magical-legends-fables-and-lore/folklore/the-folklore-and-mythology-surrounding-crickets/>
- ²“Chinese Cricket Culture”: http://www.insects.org/ced3/chinese_crucul.html
- ³“Native American Cricket Mythology”: http://www.universalsky.com/Art/cricket_intuition_power_belief.htm

2.3 Image distinctiveness

People often confuse the cricket and the grasshopper. Some people think the cricket is a grasshopper or the other makes noises, and they both do, in fact. The cricket chirps with its wings, and the grasshopper makes a similar chirping sound with its legs.⁴ Fortunately for this proposal, neither the grasshopper nor the cricket exists as an emoji. The other emoji that come close, the ant, honeybee, or lady beetle, could still not be used as suitable alternatives for the cricket. A comparable emoji is the snail, since it is also a slow-moving creature with a more distinct shape and bodily characteristics. When the cricket is compared to a snail in a Google Search it can be seen that the two have different levels of popularity. A cricket search in relation to a snail search has a 3:1 ratio worldwide over the past five years (about 56.52%) but the cricket is on a somewhat of an upward trend. When compared to a spider (excluding the tarantula) within the same time frame this ratio is 2/40 or about 5%, and compared to a lady beetle is 45/1 or about 4500%.

- ⁴“Grasshoppers and Crickets (Order: Orthoptera)”: <https://www.amentsoc.org/files/orders/orthoptera.html>

2.2 Expected usage level

2.2.1 Frequency

It is the author’s expectation that the cricket emoji will enjoy widespread use. As mentioned above, the cricket can symbolize a variety of things, from prosperity as well as happiness from its singing presence. The cricket has long been a sign for things to come, and in conjunction with other emoji can symbolize how one perceives a future outcome. In the same vein, a cricket could be used to wish a friend luck or good fortune for a future event.

2.2.2 Multiple usages

One of the greatest strengths of the cricket emoji is that it can stand for a wide variety of things. In combination with mythological significance, an awkward situation, and a cricket that you might see in your home. A particularly popular cricket symbol in American media is its expression of awkwardness.

The cricket emoji could be used in situations in which no words or phrases in combination or otherwise, suggest that very common and distinct situations of awkwardness. A search of “cricket” among the Unicode Code Chart yields one result: the cricket bat and ball. However, there are multiple references to the popular association of crickets and awkwardness online. Several GIFs use the text “*crickets*” or “*cricket* *cricket*” to demonstrate awkwardness in dialogue in the image. A Google search for “crickets meme” yields results mostly pertaining to this idea of awkwardness. When searching for “awkward silence sound effect” in YouTube, 13 of the 20 results on the first page are related to crickets.

stand for a wide variety of things.

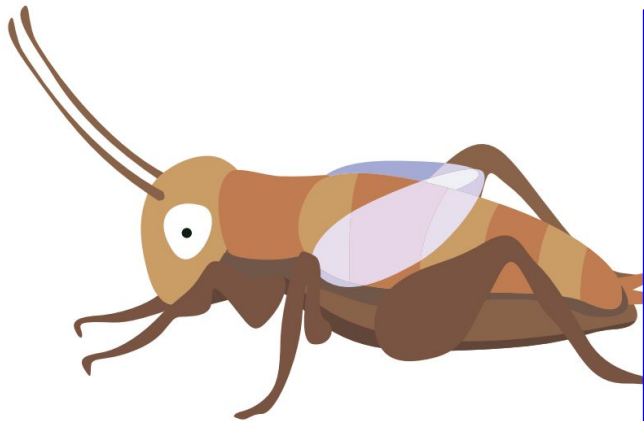
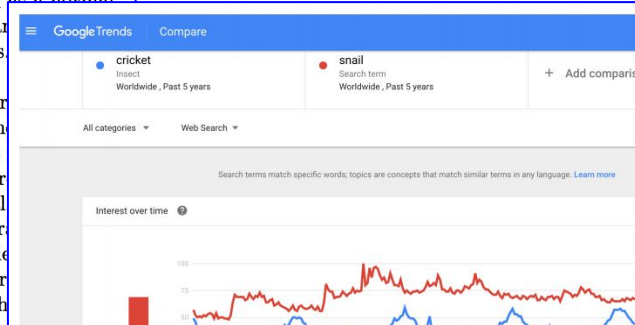


Figure 1: A cricket created by Susan Burtner and Yohji Rosen, free for use and distribution by the Unicode Consortium

1 Introduction

In many cultures around the world, the sound and sight of a cricket has a variety of good and bad prospects. In Brazil, a cricket can be seen

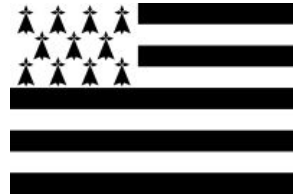
[Emoji Candidates \(click here\)](#)

Emoji 6.0 - Sample Requests

- Emojification



- RGI Tag Sequence



- RGI ZWJ Sequence

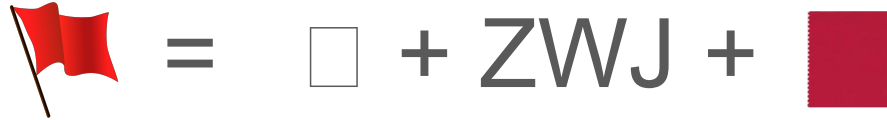


1F3F4 + ZWJ + 2620











Issues

Gender Strategy — Complete Gender-Neutral forms

Color Swatches



Directions

Appears to User	Internal Representation	Fallback Appearance
	 ZW J 	 
	 ZW J 	 

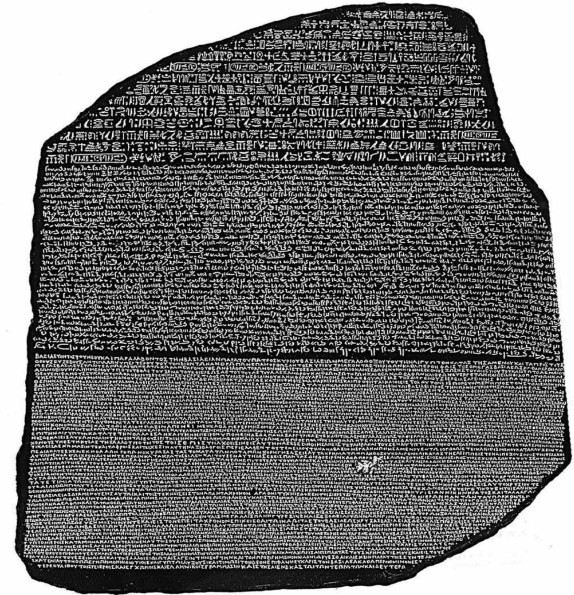
[List of Proposals \(click here\)](#)

A Greater Mission

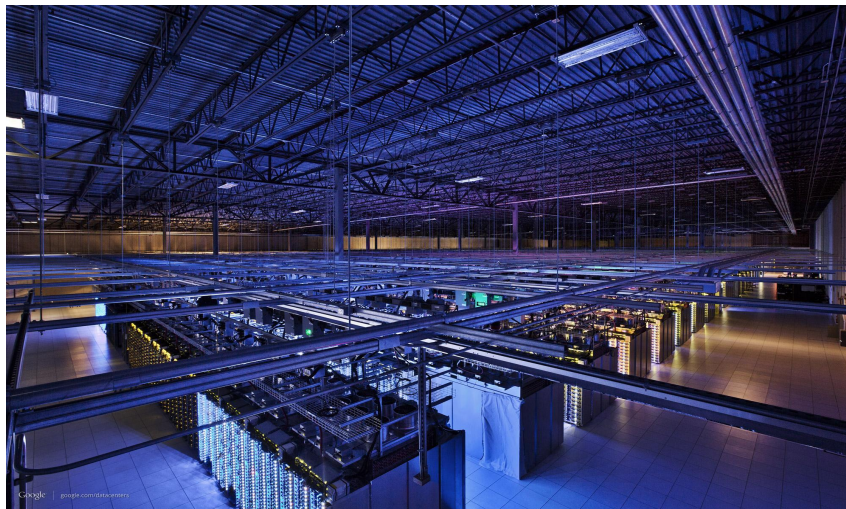
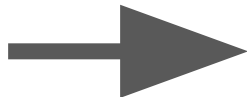
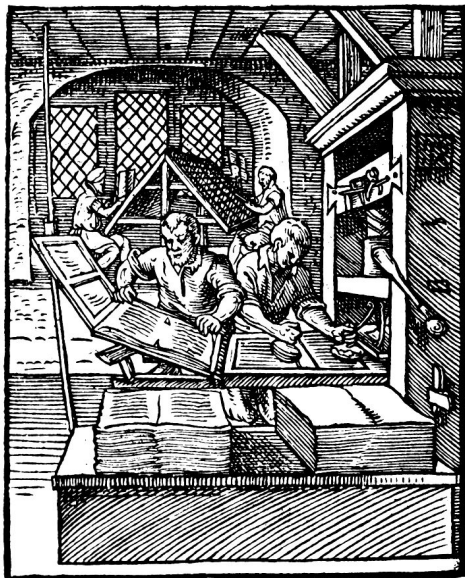
- Already vital to language support on computers:

Every keystroke, every font, every language on every computer uses Unicode

- Vision: “No language left behind”
 - Existing membership focused on ‘short tail’ (more commercially relevant)
 - Funding for ‘long tail’ languages/characters (culturally significant)



Preserving Cultural Heritage

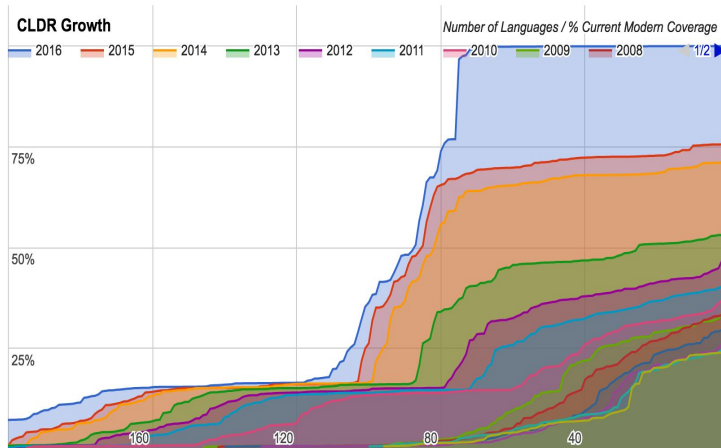


Reduce Digital Divide



Benefiting *Digitally-Disadvantaged Languages*

Emoji pushed products to improve their Unicode handling, but we have a long way to go in supporting more languages.



So we created a fundraising program for DDLs...



S&A Group

Questions?
(aka [AMA](#))



Unicode Consortium

Core Code Libraries (ICU)

Core Locale Data (CLDR)

Char Props & Algorithms (UTC)

Characters (UTC)