

جاءه دجاج مطحون قفيق حنطة حصاد

# Unicode going down the rabbit hole



# Unicode - going down the rabbit hole

A history of the written word

A history of the computerized word

Unicode

C++ and Unicode

Future perfect

2019 AD

# Unicode - going down the rabbit hole

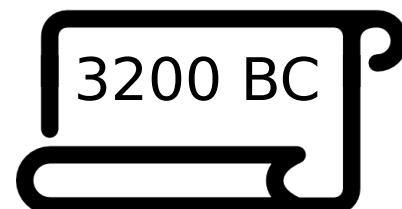
*A history of the written word*

A history of the computerized word

Unicode

C++ and Unicode

Future perfect



# الأَحْرُوفُ الْعَرَبِيَّةُ

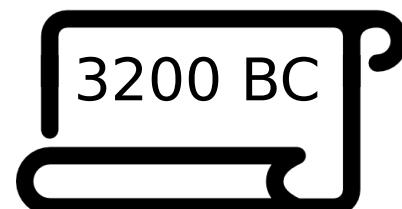
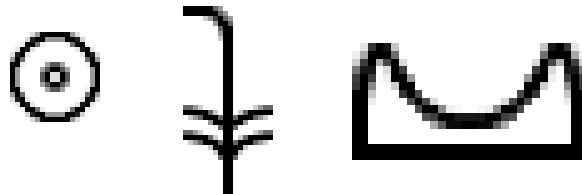


Source: George Hodan (Flickr)

3200 BC

# Egyptian Hieroglyphs

- Used from 32<sup>nd</sup> century BC until 394 AD
- Three kinds of ways to use ‘letters’:
  - Literal
  - Symbolic
  - As a ‘letter’
- Written on, or carved into, stone



أَنْجُونْفِي  
الْمُعَرِّبَةِ

- 这本书正文有 500 页。
- 我一得到任何消息，就立刻给你发短信。
- 我们能在周二前看到你演讲的全文吗？
- 《尤利西斯》为考试指定的必读书。

1200 BC?

# Trick question

- Chinese
- Japanese
- Korean
- Taiwanese

1200 BC?

# Chinese (Hanzi)

- 漢字 – Han character(s)
- Initially every ‘letter’ was a word
  - Monosyllabic
- Evolved to polysyllabic words
  - Most letters are one syllables

日 木

1200 BC?

# Hebrew

- Stopped being in common use in 400AD
- Reappeared in late 1800s
- Language is read and written right-to-left

סעיף א. כל בני אדם נולדו בני חורין ושווים בראכם  
ובזכויותיהם. כולם חוננו בתבונה ובמצפון, לפיכך חובה  
עליהם לנוהג איש ברעהו ברוח של אחווה

1000 BC

ÃfÆ'Ãcâ,-ÃiÃfâ€SA,Ã

# Greek

- Used from 8<sup>th</sup> century BC on
- Origin of 'alphabet' - **alpha, beta, gamma, ...**
- Adaptation of the Phoenician alphabet
- Evolved greatly during use

‘Όλοι οι ἀνθρώποι γεννιούνται ελεύθεροι

800 BC

# Latin

- The script we all know and use every day
- Multiple letters form a syllable
- Mapping from letters to syllables **very** region and language dependent

700 BC

# Japanese (Kanji)

- Same character set as Chinese Hanzi
  - Kanji ... also means "Han characters"

問話文動度県水安氏和政保



57 AD

子供の時分の事は、最も大抵忘れて了つたが、不思議なもので、覺えてゐる事だと、判然と昨日の事のやうに想はれる事もある。中市をして置く。其處で生れて其處で育つたのだ。

私は地方生れだ。戸籍を並べても仕方がないから、唯某縣の某十の時死別れた祖母の面だ。

今でも目を瞑ると、直ぐ顯然と目の前に浮ぶ。面長の老人だら、無論皺は寄つてゐたが、縮つた口元で、段鼻で、なか相だつたが眼が大きな眼で、女には強過る程權が有つて、古屋的眼だつたさうだ。成程然ういへば、何か氣に入らぬ事が有つて祖母が白眼で、ゾロリと睨むと、子供心にも何だか無氣味だつたやうな覺がまだ有る。

義とか云くて、併ても作者の経験した憑にも關かぬ事を聊かも技巧を加へず、有の儘に、だらくと、牛の涎のやうに書くのが流行るさうだ。好い事が流行る。私も矢張り其て行く。  
題は「平凡」、書方は牛の涎。  
まあ是からが本文だが、此處らて回を改めたが、好からうと思ふ。

### 三

Arabic script is a ligature-based language, written and read right-to-left.

# Arabic

Ligature-based language

- Written and read right-to-left

العربية الفصحى

العربية الفصحى

الأحرف العربية

~200 AD



# Ligatures

- Shuffle
- High-flying
- Gone fishing
- Sheffield
- Shuffle
- High-flying
- Gone fishing
- Sheffield

~200 AD

Ãf Ä'Äçâ,-Äi Äf Ä€SA, Ä

# Ligatures

•ffi ffl      •fi fi

•fl fl      • ffi ffi

~200 AD

ÃfÆ'Ãcâ,-ÃiÃfâ€SA,Ã

# Japanese (Hiragana)

- “Simplified” Kana (letters)
  - Only ~50 letters to learn
- Pretty close to an alphabet
- Old forms of Hiragana deprecated
  - Still used, name changed to Hentaigana

はっぴょうけっか ちょっかん ひっし

~500 AD

ÃfÆ'Ãcâ,-ÃiÃfâ€SA, A  
ÃfÆ'Ãcâ,-ÃiÃfâ€SA, A  
ÃfÆ'Ãcâ,-ÃiÃfâ€SA, A  
ÃfÆ'Ãcâ,-ÃiÃfâ€SA, A  
ÃfÆ'Ãcâ,-ÃiÃfâ€SA, A  
ÃfÆ'Ãcâ,-ÃiÃfâ€SA, A

# Japanese (Katakana)

- Used for transliterating foreign words

カタカナ - Ka-ta-ka-na

ペテルビヌデルス - Peter Bindels (approx.)

حُرُوفُ الْعَرَبِيَّةِ

~500 AD

ÃfÆ'Ãcâ,-ÃiÃfâ€SA,Ã

# Japanese (Katakana)

- Used for transliterating foreign words

カタカナ - Ka-ta-ka-na

ペテルビヌデルス - Peter Bindels (approx.)

حُرُوفُ الْعَرَبِيَّةِ

~500 AD

- Irish written language
- Usually carved into the edge of stone
- Only language with non-empty space symbol

# Ogham (ᚠᚢᚦᚢᚦ)



~500 AD

# Hangul (Korean)

- Korea used Hanja (Chinese Hanzi) before
  - Letters are not the units you think
  - ㄱ ㅏ ㄴ ㄹ ㅂ ㅓ ㄹ → 꿀벌
- 세계를 향한 대화 , 유니코드로 하십시오 .

1450 AD

# Japanese (Kanji)

- Writing of common characters simplified after WWII
  - “Shinjitai”, as opposed to old “kyūjitai”

1946 AD

# Simplified Chinese ( 简化字 ; jiǎn huà zì )

- People's Republic of China created
- Similar to Japan's simplifications, but different

- 棧 - stack (Hanzi / kyūjitai Kanji)
- 桧 - stack (Shinjitai Kanji)
- 栈 - stack (Simplified Chinese)

1960 AD

# Unicode - going down the rabbit hole

*A history of the written word*

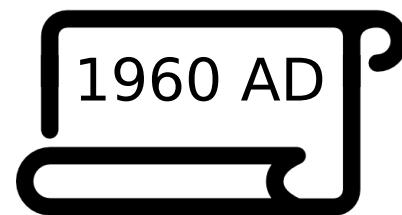
A history of the computerized word

Unicode

C++ and Unicode

Future perfect

1960 AD



# Unicode - going down the rabbit hole

A history of the written word

***A history of the computerized word***

Unicode

C++ and Unicode

Future perfect

1960 AD

# Basic classification

- Single-byte encodings
  - ASCII compatible
- Multi-byte encodings
- Variable-length encodings
  - Self-synchronizing
  - Non-self-synchronizing

1960 AD

# Simple Latin encodings

- ASCII
  - American Standard Code for Information Interchange
  - 7-bit, all characters in logical order
- EBCDIC
  - Extended from Binary-Coded Decimal (BCD) Interchange Code
  - 8-bit, compatible with punch cards

الأَحْرَفُ الْعَرَبِيَّةُ



1963 AD

# GB2312

- 8-bit, bottom half is ASCII
- Variable length, multibyte character set
- Guojia Biaozhun (国家标准 , national standard)
- Chinese character set

أَلْحَرُوفُ الْعَرَبِيَّةُ

1980 AD

# ANSI

- ANSI didn't have anything to do with it
- Also often called “Extended ASCII”
- Code Page 437 from DOS
- 8-bit, bottom half is ASCII
- Top half filled with line drawing and some math symbols

1981 AD

ÃfÆ'Ãcâ,-ÃiÃfâ€SA,Ã

# DOS CP437

- ANSI didn't have anything to do with it
- Also often called "Extended ASCII"
- Code Page 437 from DOS
- 8-bit, bottom half is ASCII
- Top half filled with line drawing and some math symbols

1981 AD

# DOS CP866

- 8-bit, bottom half is ASCII
- Top half filled with line drawing and cyrillic letters

الأَحْرُوفُ الْعَرَبِيَّةُ

1981 AD

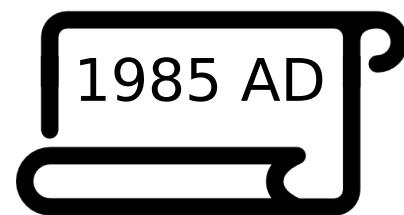
# DOS CP850

- 8-bit, bottom half is ASCII
- Top half filled with line drawing and accented letters

# Windows 1252

- 8-bit, bottom half is ASCII
- Top half filled with text processor symbols, accented letters and other things

and yes, LaTeX calls this “newansi”.



ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Windows 1255

- 8-bit, bottom half is ASCII
- Top half filled with Hebrew letters

الأُحْرَوْفُ الْعَرَبِيَّةُ

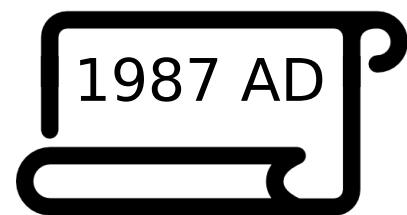
1985 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# ISO 8859

- Standardize 8-bit code pages
  - 15 defined
  - Cannot mix use

الأُحْرَوْفُ الْعَرَبِيَّةُ



# GBK

- Extension of GB2312
- Just adds way more Chinese characters

# Unicode - going down the rabbit hole

A history of the written word

*A history of the computerized word*

Unicode

C++ and Unicode

Future perfect

1993 AD

# Unicode - going down the rabbit hole

A history of the written word

A history of the computerized word

***Unicode***

C++ and Unicode

Future perfect

# Unicode

Unicode is intended to address the need for a workable, reliable world text encoding.

Unicode could be roughly described as "wide-body ASCII" that has been stretched to 16 bits to encompass the characters of all the world's living languages. In a properly engineered design, 16 bits per character are more than sufficient for this purpose.

- Joe Becker

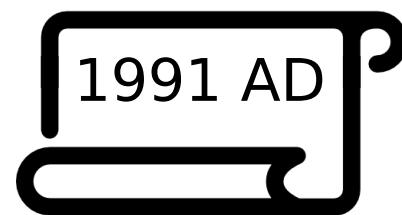
1991 AD

Ãf Ä'Äçâ,-Äi Äf â€SA, A  
Ãf Ä'Äçâ,-Äi Äf â€SA, A

# Unicode

- encodes graphemes rather than glyphs
- "wide-body ASCII" stretched to 16 bits
- incorporated on 3 January 1991
- first volume published in October 1991

أَلْحُرُوفُ الْعَرَبِيَّةُ



# UCS2

- 16-bit, bottom bit is ASCII
- Intended to contain all the world's languages
- No more conversions
- Byte order dependent
  - Byte Order Mark (BOM) recommended at start

1991 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# UTF-7

- Not ASCII compatible per se
- Intended to transport Unicode across (7-bit) ASCII-only communication channels

# UTF-8

- ASCII compatible
- Intended as a temporary encoding
  - C-style functions with 8-bit encoding support
  - More efficient than UCS2 for ASCII
  - Less efficient than UCS2 for all Eastern languages
  - Variable-length

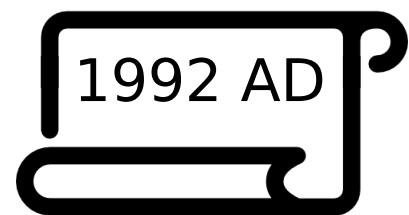
1991 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Han unification

- 16-bit character set
- GBK has 21,886
- Hanja, Hanzi, Kanji and Simplified Chinese...

الأَحْرُوفُ الْعَرَبِيَّةُ



# Han unification

- Not all characters unified
- Most characters no issue
- Some have different glyphs, but the same grapheme now...

1992 AD

Code point	Chinese (simplified) ( zh-Hans )	Chinese (traditional) ( zh-Hant )	Chinese (traditional, Hong Kong) ( zh-Hant-HK )	Japanese ( ja )	Korean ( ko )	Vietnamese ( vi-Hani )	English
U+4ECA	今	今	今	今	今	今	now
U+4EE4	令	令	令	令	令	令	cause/command
U+514D	免	免	免	免	免	免	exempt/spare
U+5165	入	入	入	入	入	入	enter
U+5168	全	全	全	全	全	全	all/total
U+5177	具	具	具	具	具	具	tool
U+5203	刃	刃	刃	刃	刃	刃	knife edge
U+5316	化	化	化	化	化	化	transform/change
U+5916	外	外	外	外	外	外	outside
U+60C5	情	情	情	情	情	情	feeling
U+624D	才	才	才	才	才	才	talent
U+62B5	抵	抵	抵	抵	抵	抵	arrive/resist
U+6B21	次	次	次	次	次	次	secondary/follow

1992 AD

# What is a character

- a or a ?
- a or a ?
- 丟 vs 丢 ?

1992 AD

# What is a character

- a or a ? → Same thing, different glyph
- a or a ? → Different thing, same glyph
- 丟 vs 丢 ? → Different thing



# Wingdings

1990 AD

الأَحْرُوفُ الْعَرَبِيَّةُ



# Wingdings

1990 AD

# Emoji

- e ( 絵 , "picture") + moji ( 文字 , "character")
- Popular because of instant-messaging and texting



1997 AD



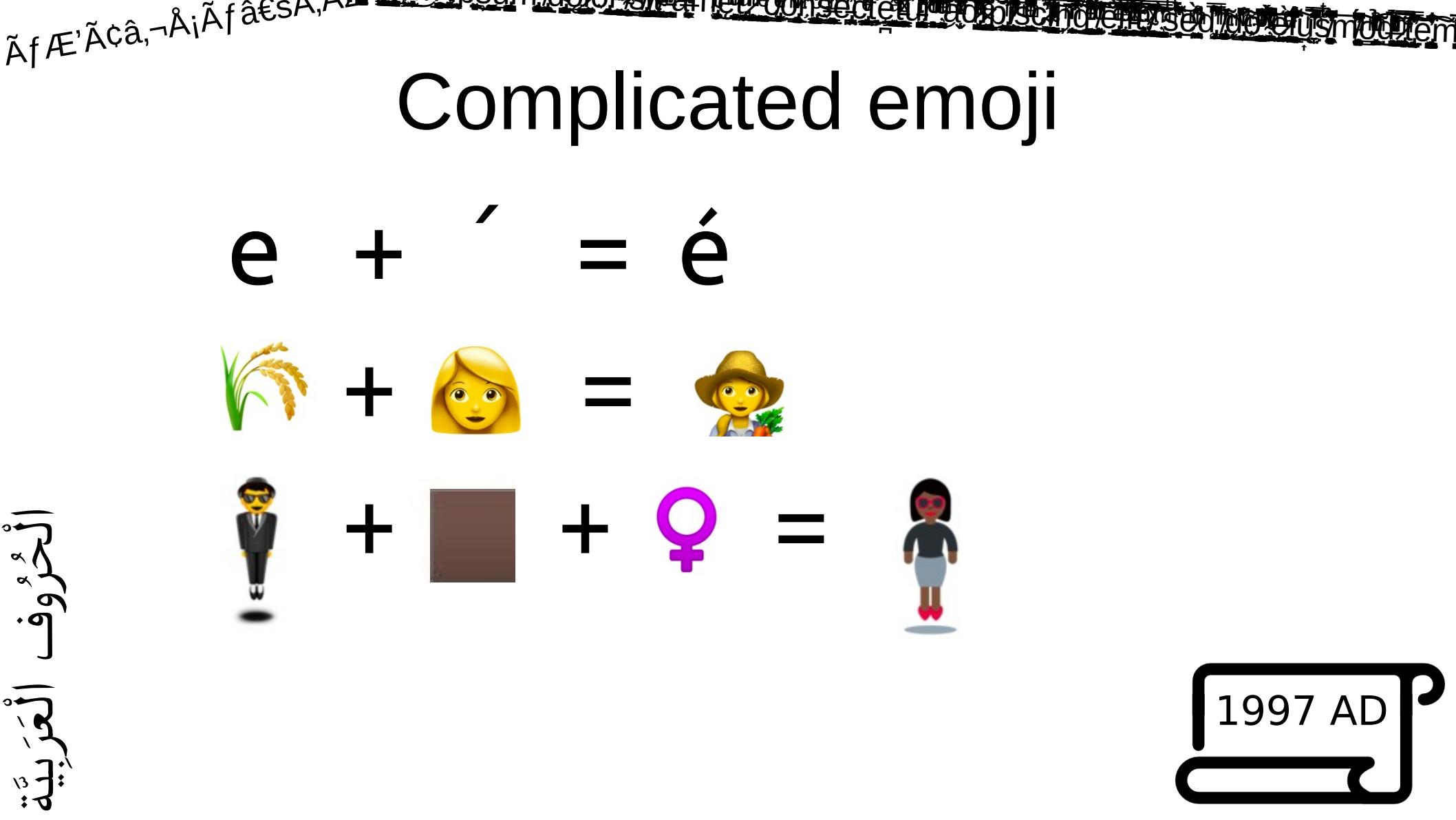
2017 AD

Ãf Ä'Äçâ,-Äi Äf â€SA,A  
Ãf Ä'Äçâ,-Äi Äf â€SA,A

# Controversies around emoji

-  is mostly replaced with a water gun
- ,  and  lend themselves really well to crime
-  and  used to represent anatomy instead of a fruit and a vegetable
- Lawsuits threatened for allowing use of 

1997 AD



# Complicated emoji

e + ' = é

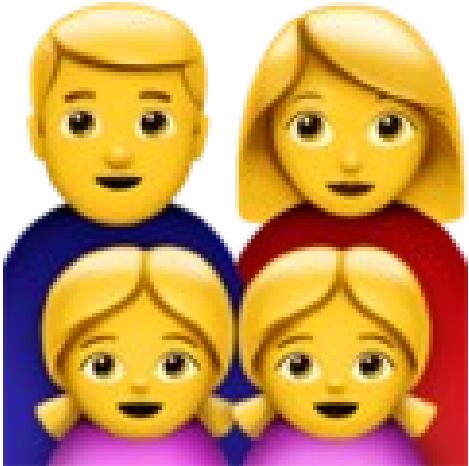
🌾 + 🧑 = 🌾🧑🌾

👨 + 🤵 + ♀ = 👩

1997 AD



# Emoji



1997 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# UTF-16

- 16-bit, bottom bit is ASCII
- Surrogate pairs create characters not in UCS2
- Byte order dependent
  - Byte Order Mark (BOM) recommended at start
- Variable length
- Cannot encode code points beyond 1'114'112

أَلْحَرُوفُ  
الْعَرَبِيَّةُ

1997 AD

# UTF-32 (UCS-4)

- 32-bit, bottom bit is ASCII
- Surrogate pair characters invalid
- Byte order dependent
  - Byte Order Mark (BOM) recommended at start
- Fixed length
- Very inefficient for any text
- May not contain code points inaccessible in UTF-16

1997 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# UTF-8

- ASCII compatible
- Intended as a temporary encoding
- No byte order dependency
- Pretty efficient, only Asian languages are less efficient than UTF-16
- May not contain code points inaccessible in UTF-16

ـ حـرـوفـ الـعـرـبـيـةـ

1997 AD

# Encoding unicode

- Storing UTF-16 as UTF-8
  - CESU-8
- Storing UTF-16 with null bytes as UTF-8
  - Modified UTF-8

2011 AD

# Encoding unicode

- But what if my system uses EBCDIC?
  - UTF-EBCDIC
- We can't just switch from GBK
  - GB18030

2011 AD

# Mojibake

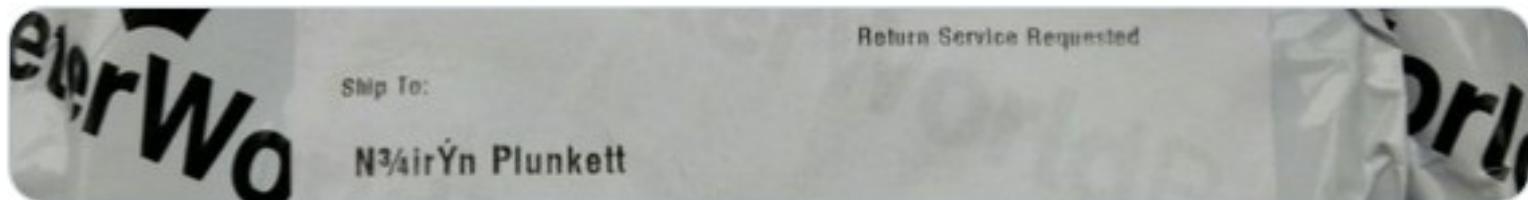
- The result of decoding text in one encoding, with a different encoding
  - Most common: UTF8 read as some 1-byte encoding
  - Billy O'Neill

2011 AD



Noirín Plunkett  
@noirinp

Didn't realise that in 2015, there were still ways to misencode my name that I hadn't seen before.



9:49 PM · May 19, 2015 · Twitter for Android

1 Retweet 3 Likes

2015 AD

Preview Untitled-1 - Visual Studio Code

Untitled-1 ●

```
1 = Вывод уравнения упругой линии балки для
2 :nofooter:
3 :lq: pass:[&laquo;]
4 :rq: pass:[&raquo;]
5 :stem: latexmath
6
7 Hello
8 Привет
9
10 {asciidocor-version}
```

Preview Untitled-1

Hello

2.0.1

2015 AD

# الأْحَرُوفُ الْعَرَبِيَّةُ



"This letter was sent to a Russian student by her French friend, who manually wrote the address that she received by e-mail." Mojibake diacritics translated to Cyrillic by the postal employees via [The New Aesthetic](#)

2015 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Displaying unicode

- Displaying most encodings is simple:
  - Map each character to a glyph
  - Calculate positions for each glyph
  - Render glyphs at positions

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Displaying unicode

- Unicode is different
  - Combining diacritics
  - RTL languages – sometimes mixed!
  - Korean, Devanagari, ... complex characters
  - Ligatures

Ãf ÄE'Ãcâ,-ÃiÃf â€SA,Ã

# Displaying unicode

- 1)Decode to unicode code points
- 2)Normalize code points
- 3)Extract grapheme clusters
- 4)Shape to determine glyphs & positions
- 5)Render glyphs at appropriate positions

أَلْحَرْفُ  
الْعَرَبِيَّةُ

2011 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Unicode abuse

- Zalgo text
  - Take regular text (Lorem ipsum)
  - Add some combining diacritics (Łöřëm íþşüṁ)
  - Add tons more ( Łöřëmíþşüṁ)
  - And even more ( Łöřëmíþşüṁ)

أَلْحُرُوفُ الْعَرَبِيَّةُ

2010 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Unicode abuse

- ¯\\_(ツ)\_/¯
  - Starts with macron (U+00AF), backslash, underscore, open parentheses, Hiragana character, closing parentheses, underscore, slash, macron

Ãf Ä'Äçâ,-Äi Äf â€SA, Ä

# Unicode abuse

Just take letters from other languages that look like Latin upside-down.

languages that look like latin upside-down.  
Just take letters from other

الأَحْرُوفُ الْعَرَبِيَّةُ

2010 AD

# Unicode abuse

- Using imitation letters to resemble something else
- Often abused in domain names
  - <https://www.apple.com>
  - <https://www.apple.com>

2017 AD

# Unicode abuse

- Using imitation letters to resemble something else
- Often abused in domain names
  - <https://www.apple.com> → Actually Apple
  - <https://www.apple.com> → Cyrillic letters

2017 AD

# Unicode - going down the rabbit hole

A history of the written word

A history of the computerized word

***Unicode***

C++ and Unicode

Future perfect

# Unicode - going down the rabbit hole

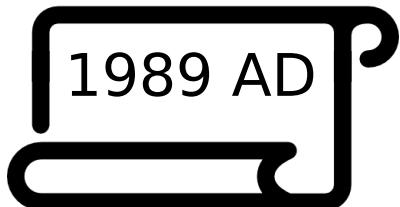
A history of the written word

A history of the computerized word

Unicode

***C++ and Unicode***

Future perfect



1989 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# C++ in 1989

- No support
  - Remember, unicode doesn't get created until 1991

# C++ in 1998

- std::string
  - Current locale with 8-bit characters
- std::wstring
  - Current wide locale, with sizeof(wchar\_t) characters
    - 16 bit on Windows, 32 bit on Unix
  - wchar\_t definition: big enough to hold any character in the encoding

1998 AD

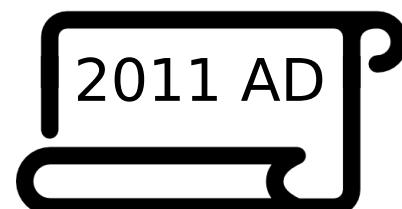
# C++ in 1998

- Not much to handle unicode
  - C95 gave us what we have
- wcstombs
- mbstowcs

1998 AD

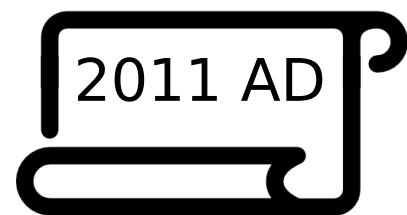
# C++ in 2011

- std::u16string, char16\_t
  - UTF16 string
- std::u32wstring, char32\_t
  - UTF32 string
- std::regex
  - No Unicode considerations
- No UTF8 string or character type



# C++ in 2011

- uchar.h, cuchar
  - Multibyte to u16string, u32string and reverse

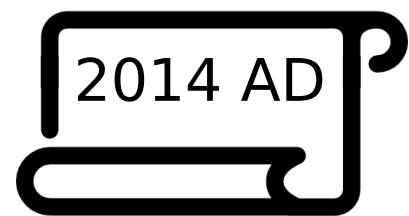


الأُحْرُوفُ الْعَرَبِيَّةُ

Ãf Ä'Ãçâ,-ÃiÃf â€SA,Ã

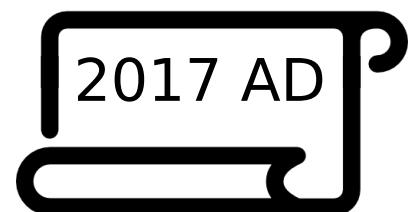
# C++ in 2014

- ...



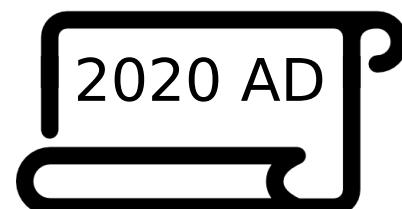
# C++ in 2017

- u8"", u"" and U"" string literals
  - resulting in char[], char16\_t[] and char32\_t[]
- SG16 formed to handle proper Unicode in C++
  - ISO C++ committee study group



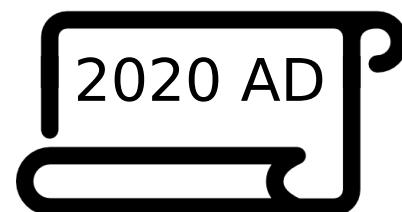
# C++ in 2020

- std::u8string, char8\_t
  - UTF8 string and character type
- u8"" changed to char8\_t[]
  - Breaking change
- uchar.h, cuchar
  - Multibyte to u8string and reverse



# C++ in 2020

- std::string<T>
  - No easy construction from mbs or wcs inputs
  - No logical output conversion for printf, std::cout, ...



```
#include <string>
#include <iostream>

int main() {
    std::u16string s = u"Hello ";
    std::u32string s2 = U"CppCon ";
    std::u8string s3 = u8"World ";
    std::cout << s;
    std::cout << s2;
    std::cout << s3;
}
```

# C++ in 2020

2020 AD

```
#include <string>
#include <iostream>

int main() {
    std::u16string s = u"Hello ";
    std::u32string s2 = U"CppCon ";
    std::u8string s3 = u8"World ";
    std::wcout << s;
    std::wcout << s2;
    std::wcout << s3;
```

]

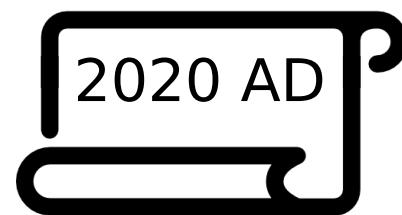
2020 AD

ÃfÆ'Ãçâ,-ÃiÃfâ€SA,Ã

# Displaying unicode

- 1)Decode – <mostly not in C++>
- 2)Normalize – <not in C++>
- 3)Extract graphemes – <not in C++>
- 4)Shape – Harfbuzz, or similar
- 5)Render – Freetype, or similar

الْحُرُوفُ  
الْعَرَبِيَّةُ



# Unicode - going down the rabbit hole

A history of the written word

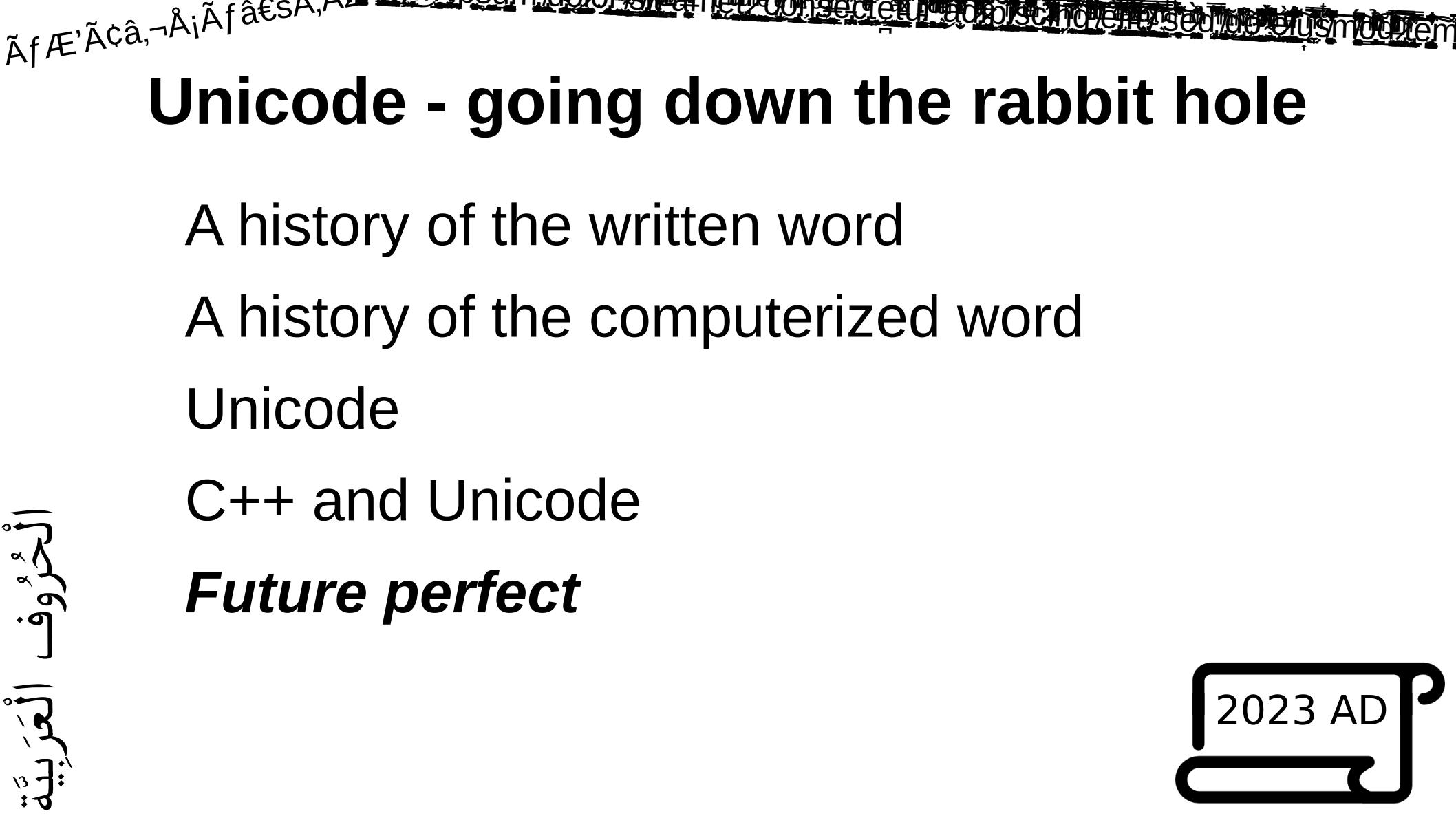
A history of the computerized word

Unicode

**C++ and Unicode**

Future perfect

2020 AD



# Unicode - going down the rabbit hole

A history of the written word

A history of the computerized word

Unicode

C++ and Unicode

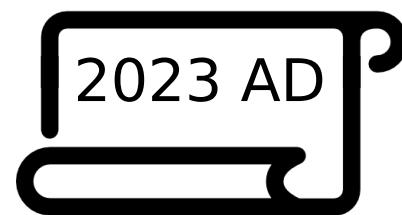
*Future perfect*

الأَحْرُوفُ  
الْعَرَبِيَّةُ

2023 AD

# SG16 plans

- C++20: Lay the groundwork
  - `char8_t` needs to be in
  - `u8""` needs to produce `char8_t[]`
  - Lots of detail fixes
  - Get references to Unicode standards (ISO10646) in the C++ standard



# SG16 plans

- C++23: Introduce text capable types
  - Text: a string that knows code points and graphemes
  - Enable unicode support for other parts of the standard
    - std::regex is unfixable
    - Hana Dusíková's CTRE library slated for C++23 with unicode support
- <https://wg21.link/p1238> (Unicode Direction paper)

2023 AD

# SG16 plans

- When to convert?
  - Current proposals: all internal logic is UTF8
  - Check your legacy encodings at the border
    - (convert whenever you receive data or send it out)
  - Preferably, normalize at the border too
    - But need normalizing during string edits too

# UTF-8 EVERYWHERE

## MANIFESTO

### 1 Purpose of this document

Our goal is to promote usage and support of the UTF-8 encoding and to convince that it should be the default choice of encoding for storing text strings in memory or on disk, for communication and all other uses. We believe that our approach improves performance, reduces complexity of software and helps prevent many Unicode-related bugs. We suggest that other encodings of Unicode (or text, in general) belong to rare edge-cases of optimization and should be avoided by mainstream users.

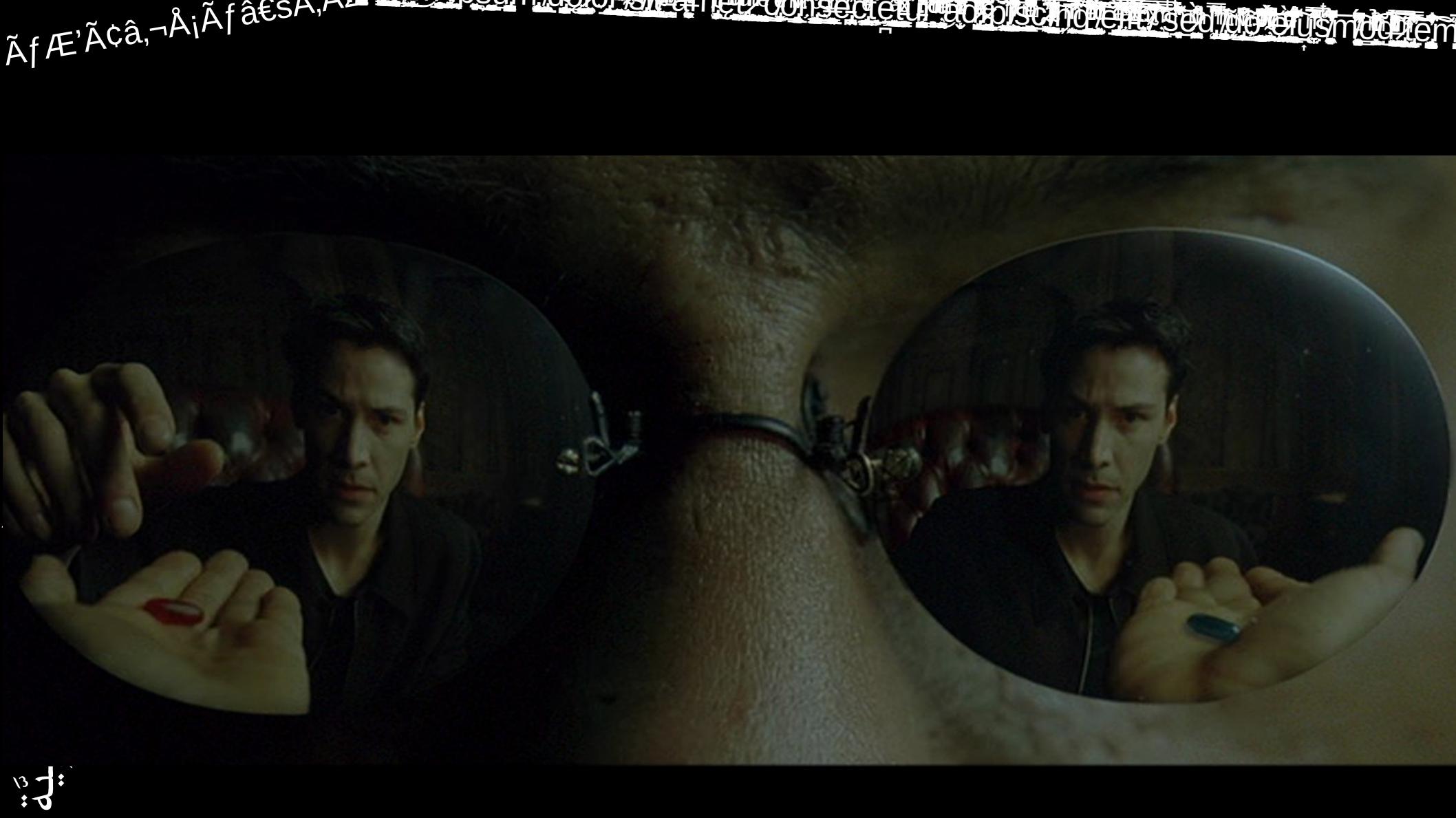
This document contains special characters. Without proper rendering support, you may see question marks, boxes, or other symbols.

In particular, we believe that the very popular UTF-16 encoding (often mistakenly referred to as ‘widechar’ or simply ‘Unicode’ in the Windows world) has no place in library APIs except for specialized text processing libraries, e.g.

2023 AD

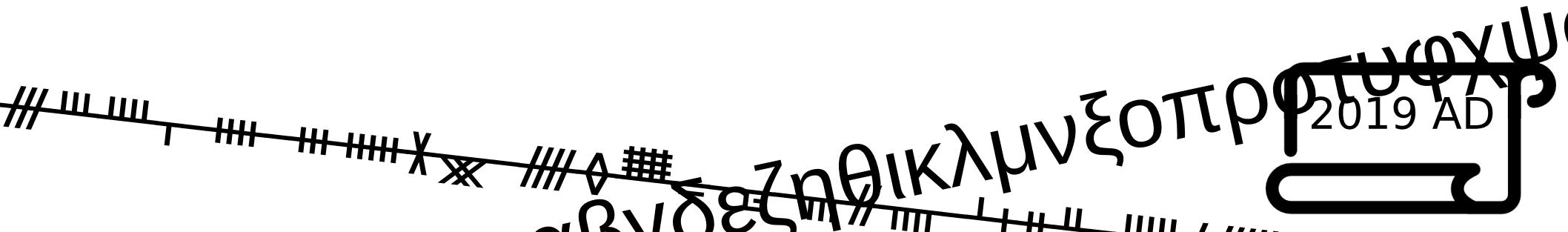
# SG16

- For any questions about SG16's current trajectory, ask on our mailing list
  - [unicode@isocpp.open-std.org](mailto:unicode@isocpp.open-std.org)
- Organisational stuff
  - Ask Tom Honermann, chair of SG16
  - <[tom@honermann.net](mailto:tom@honermann.net)>

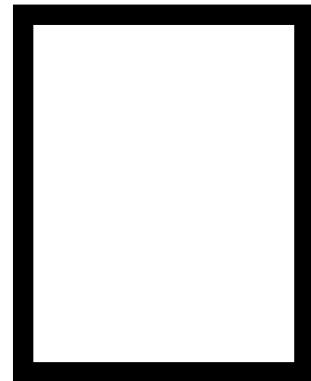


جاءه دجاج مطربي حزوه قمني فـ GHIJKLMNOPQRSTUVWXYZ

This was  
**Unicode**  
going down the rabbit hole



جاءه دعوه وحزنه يطحي قسم قشر حنثه حصد



2019 AD



ΤΟΥΡΚΙΚΗ ΜΝΗΜΟΝΙΑ  
2019 AD

جاءه دجاج مطحون قفيق قشر حنطة ABCDEFGHIJKLMNOPQRSTUVWXYZ

# Unicode

going down the rabbit hole

Presented at CppCon 2019

Peter Bindels (@dascandy42)

© TomTom

