

Introduction To Unicode

Making the digital world more inclusive

Alexandre Bergel

Computer Science Department - FCFM

University of Chile

<http://bergel.eu>

abergel@dcc.uchile.cl

[@AlexBergel](#)

◌ U+3002	🍒 U+26AA	ꣳ U+0AAC) U+207E	` U+02CB	♀ U+26A2	@ U+FF20	😘 U+1F618
› U+203A	୨ U+0E15	‘ U+2018	ଁ U+10E6	˘ U+FF9E	🎾 U+1F3BE	✕ U+00D7	😬 U+1F617
Ს U+12CE	↑ U+2191	<p>Everyone in the world should be able to use their own language on phones and computers.</p> <p>LEARN MORE ABOUT UNICODE</p>					
𑂔 U+0937	人 U+4EBA						
Δ U+0394	Ⅱ U+1111	🍟 U+1F35F	【 U+3010	ꣳ U+0CA5	♠ U+2660	ϑ U+04E9	☹ U+2639
₩ U+D1E1	☆ U+2606	👽 U+1F47D	◼ U+2B1B	˘ U+1D54	↩ U+21A9	🏏 U+0B37	😍 U+1F970
♂ U+26A3	➤ U+309D	◻ U+2752	⚡ U+1F923	🔒 U+1F512	⌘ U+2318	🎵 U+266B	# U+266F



<https://en.wikipedia.org/wiki/Teleprinter#/media/File:WACsOperateTeletype.jpg>



https://en.wikipedia.org/wiki/Teleprinter#/media/File:Siemens_t37h_without_cover.jpg

ASCII

- Supports meaningful exchange of text data
- Proposed in 1963
- Coded on 7-bits => 128 characters
 - **A** = 65
- Very limited, not even adequate for English
 - e.g., “résumé” is an English word
- Only letters, digits, and punctuation are considered as printable characters

ASCII Code Chart

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

ASCII Code Chart

Line feed

Carriage return

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

“Hello\nWorld” ➡ Hello
World

Line feed

Carriage return

ASCII Code Chart

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

“Hello\nWorld” ➡ Hello
World

Line feed

Carriage return

ASCII Code Chart

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3									FS	GS	RS	US
2		!	"	#									,	-	.	/
3	0	1	2	3									<	=	>	?
4	@	A	B	C									L	M	N	O
5	P	Q	R	S									\]	^	_
6	`	a	b	c									l	m	n	o
7	p	q	r	s										}	~	DEL



Many other standards

- ASCII has many limitations
- Many industrials proposed *their own improvement*
 - MacRoman from Apple
 - IBM's EBCDIC-based code pages
 - Microsoft, SAP, Oracle, ...

Unicode

- “Unicode is an information technology *standard* for the consistent encoding, representation, and handling of text expressed in most of the *world's writing systems*.”
- Designed to improve the mess inherited from *telegraph machine*
- *Enable world-wide interchange of data*
- Multilingual
- A single implementation
- Support legacy data

Writing direction

Gidi said, “	אם אין אני לי מי לי	”.
→	←	→

- ① אָ טיף געטע 1123 זייטע.
- ② נא ראה עמוד 1123.
- ③ راجع صفحة ١١٢٣ من فضلك.
- ④ Please see page 1123.
- ⑤ 1123ページをみてください。

- ⑥ 1123ページを
みてください。

Character composition

A + ö → Ä
0041 0308

2 + ◊ → ◊2
2621 20DF

☕ + ☹ → ☹☕
2615 20E0

Ɔ + ○ → ○Ɔ
062D 20DD

प/ूर्/र/र्/त/ि



Overview of Unicode

Unicode

- > 143,859 characters
- > 154 modern and historical scripts
- *Script*: collection of letters and other written signs used to represent textual information

Character model

- Four layers
- Level 1: Abstract character set => *What is a character?*
- Level 2: Coded character set => *How to name and enumerate abstract character?*
- Level 3: Character encoding forms => *How to represent coded characters in a computer?*
- Level 4: Character encoding schemes => *how to serialize characters into bytes?*

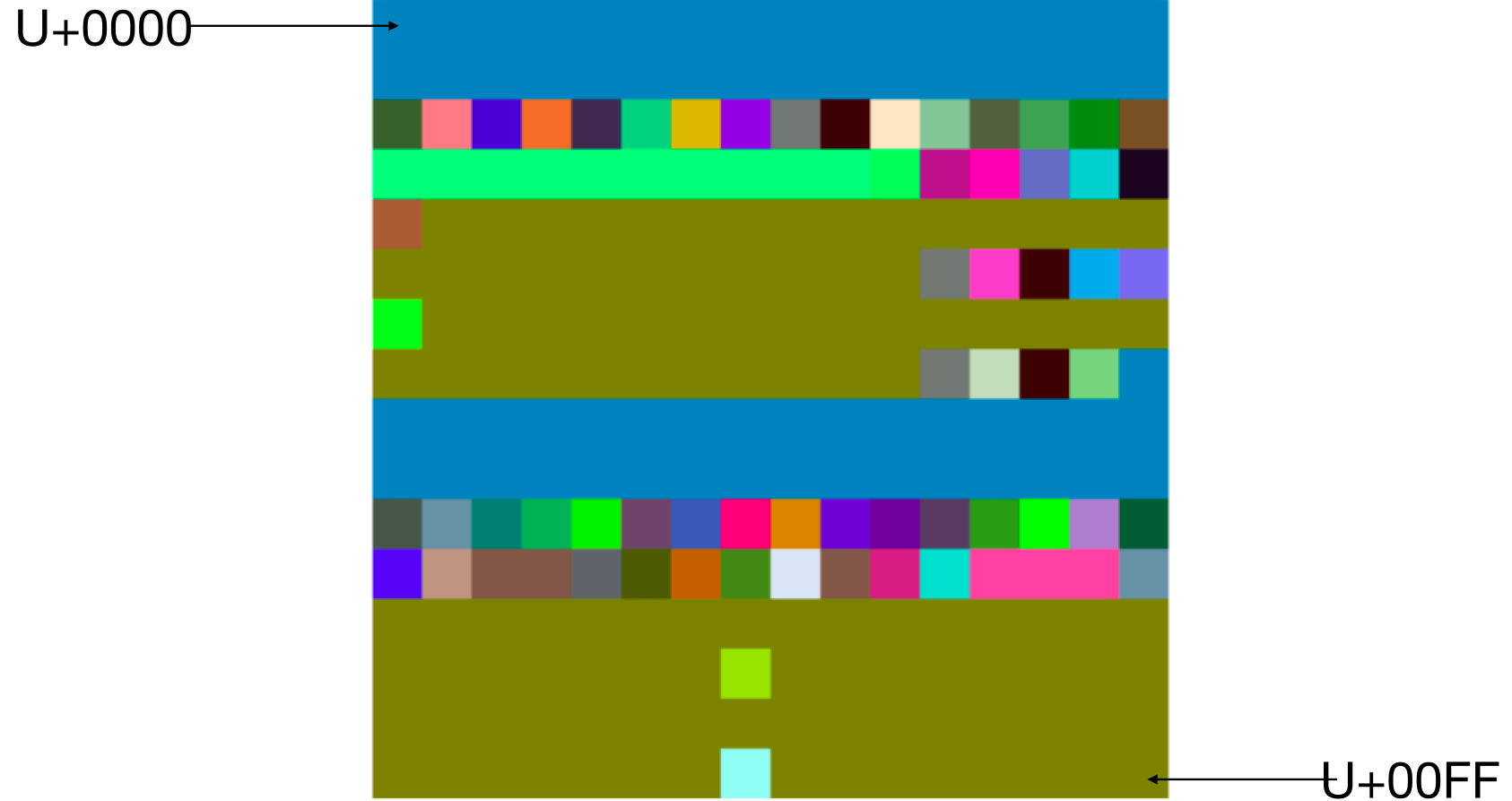
Abstract character set (level 1/4)

- Character: *The smallest component of written language that has semantic value*
- Wide variation across scripts:
 - *Alphabetic* => each character is a letter. Both *consonant* and *vowel* have equal status
 - *Syllabary* => each character is a syllab
 - E.g., Hiragana (あ, せ, め), Cherokee (Ꮝ, Ꮜ, Ꮚ), Vai (𞤎, 𞤍, 𞤌)
 - *Abjad* => each character is a consonant, vowel marking is absent
 - E.g., Hebrew (א, ב, ג), Punic (𐤀, 𐤁, 𐤂)
 - *Abugidas* => each character is a sequence of consonant - vowel, the vowel notation is secondary
 - *Logographic* => each character is a word
 - E.g., Egyptian hieroglyphs, emoji (although still debated)
- Abstract character: *a unit of information used for the organization, control, or representation of textual data*

Coded character set (level 2/4)

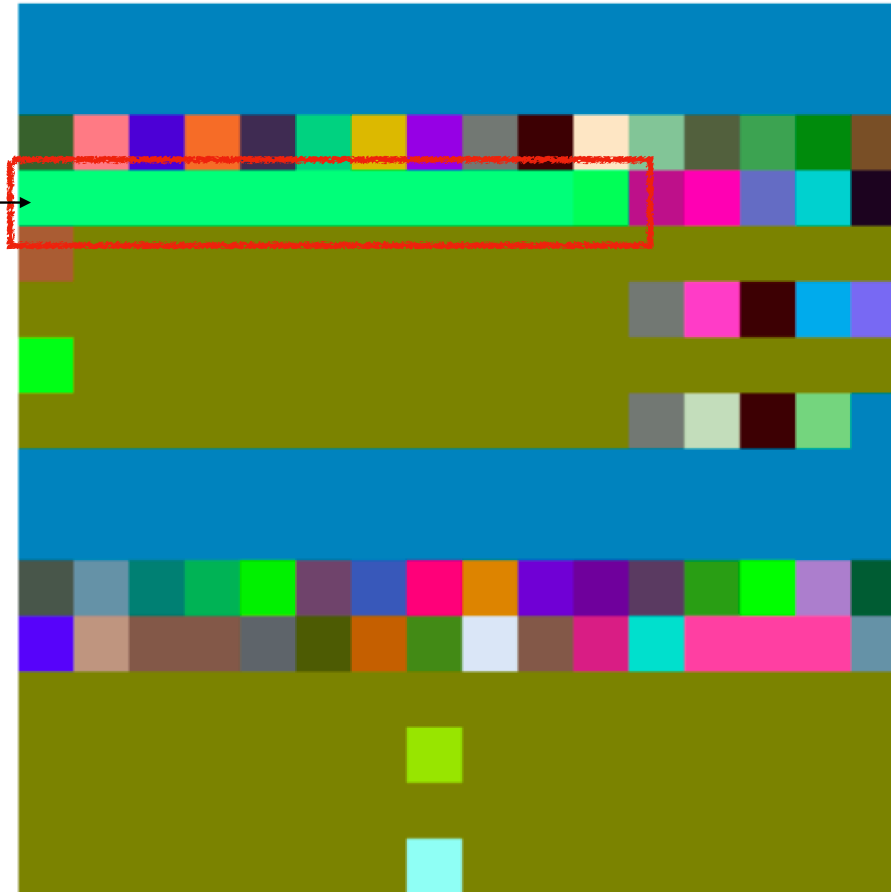
- Give a name and a code point to each abstract character
- Name: LATIN CAPITAL LETTER A
- Code point: pure number
 - Legal value: U+0000 - U+10FFFF
 - Space for >1M different characters
- Characters that are specific to a script *are mostly grouped*
- *No connection* to the computer

The first 256 characters



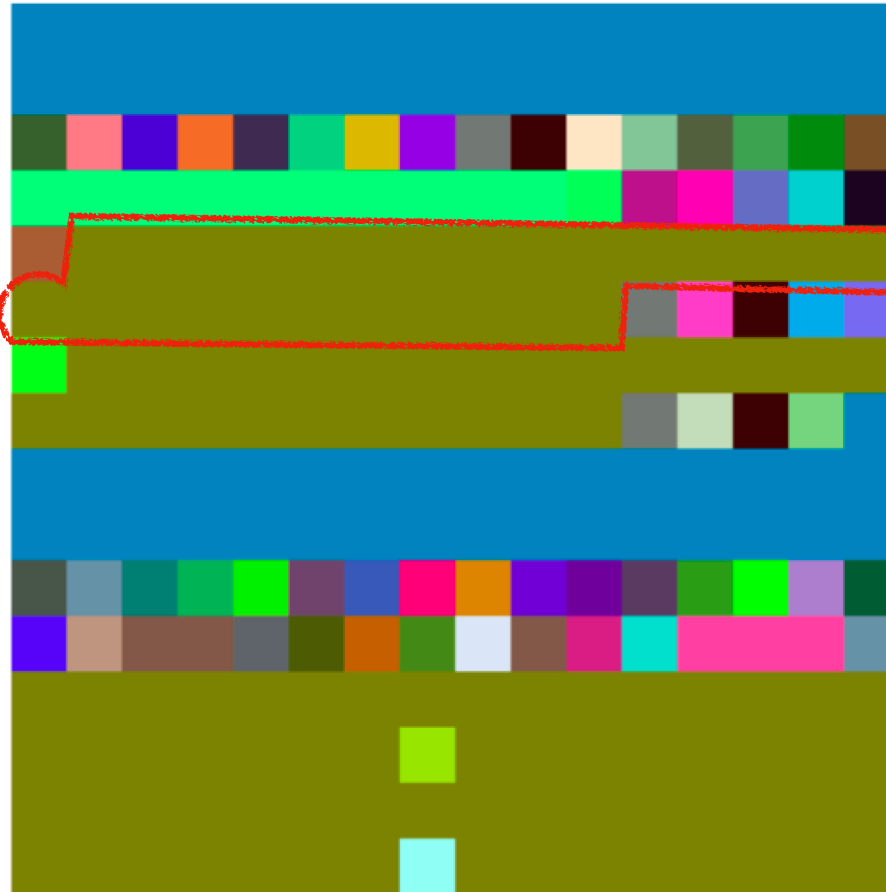
The first 256 characters

Digit 0 - 9

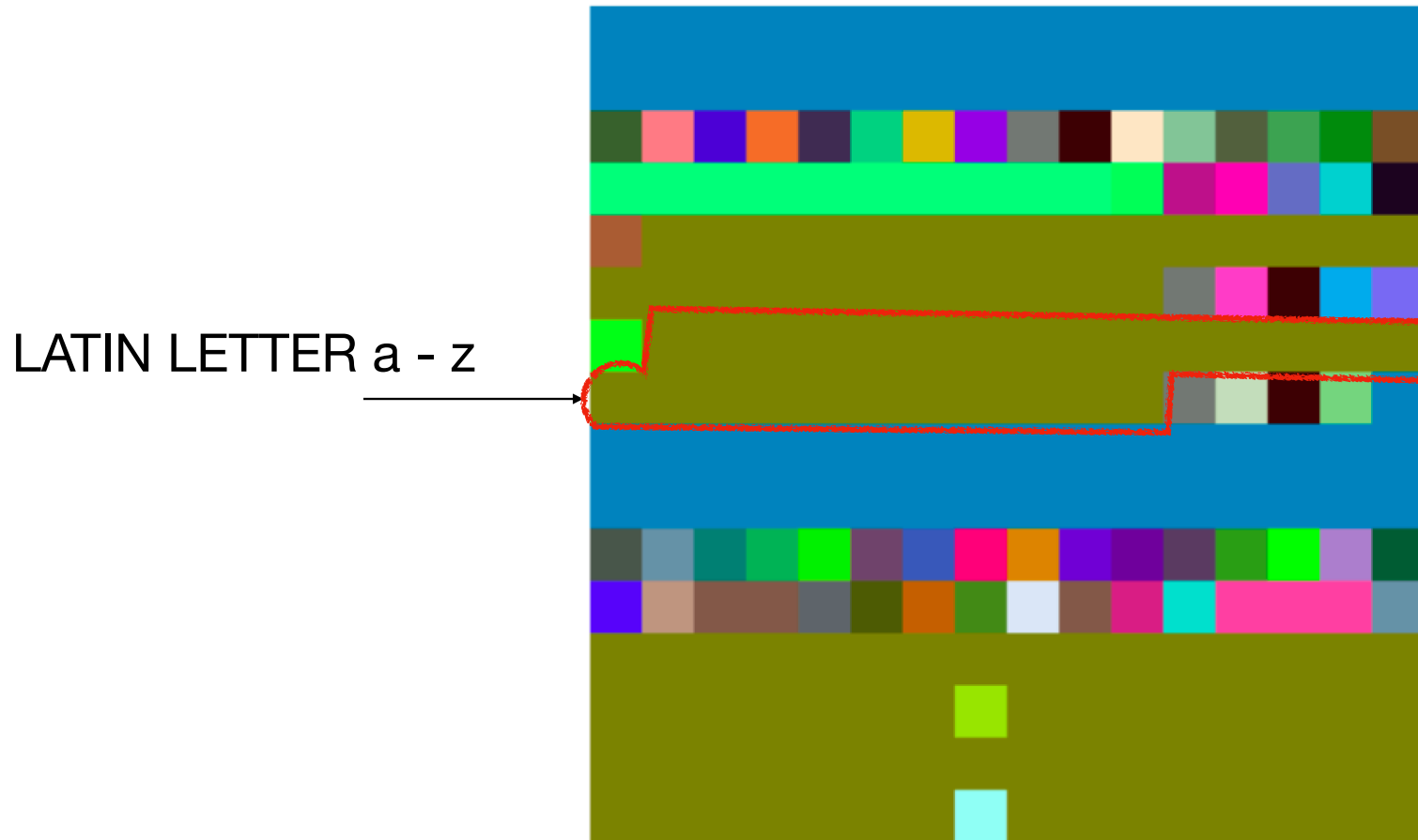


The first 256 characters

LATIN LETTER A - Z



The first 256 characters



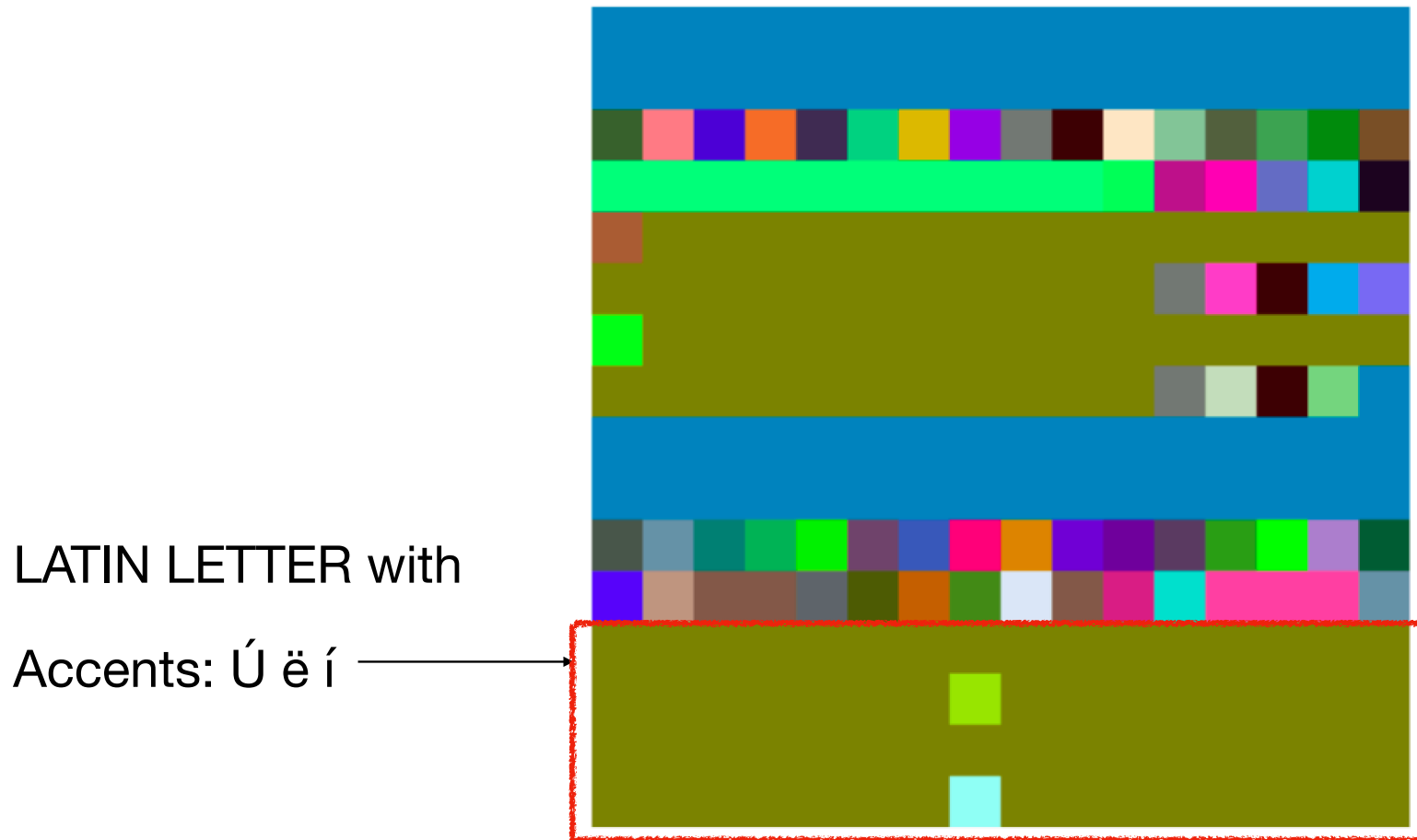
The first 256 characters

ASCII Code Chart																
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

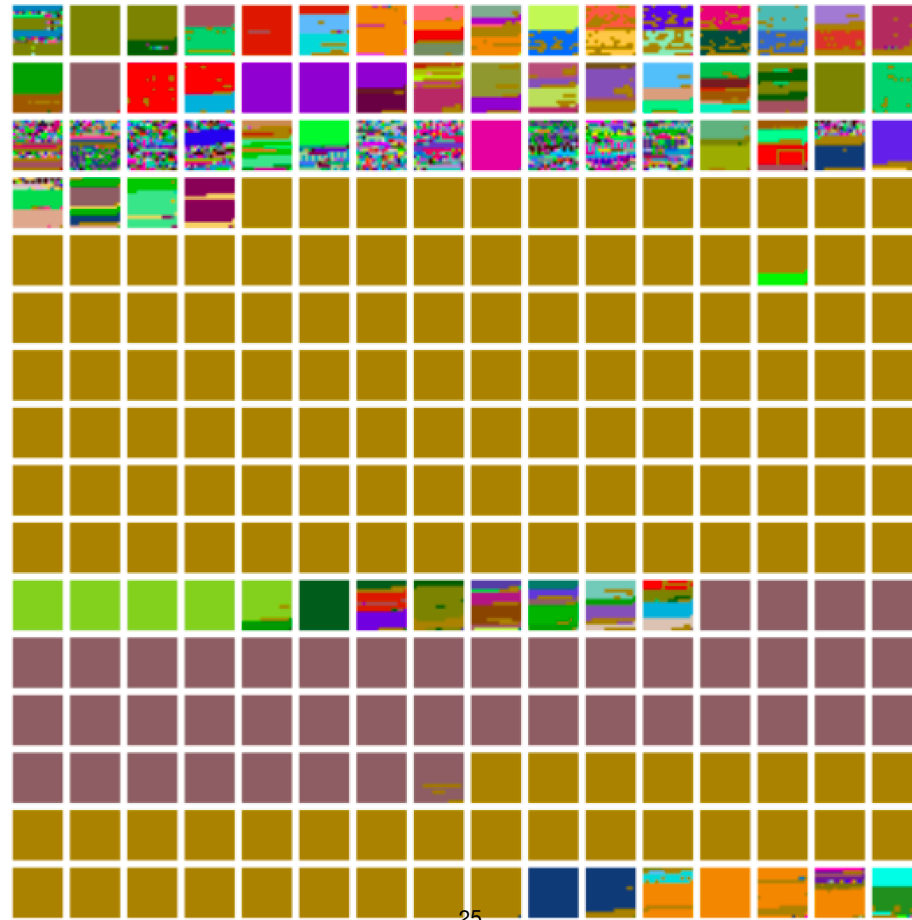
LATIN LETTER a - z



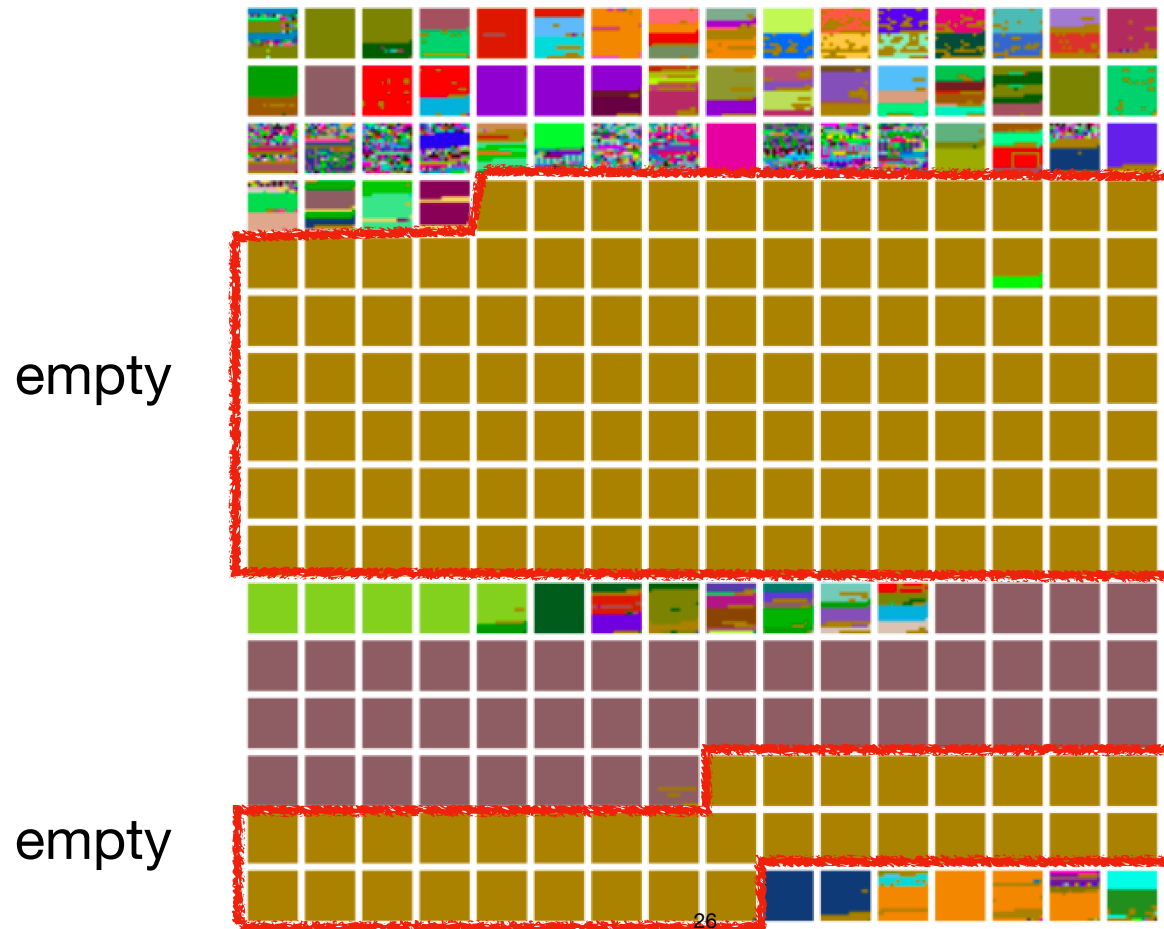
The first 256 characters



The first 65K characters



The first 65K characters

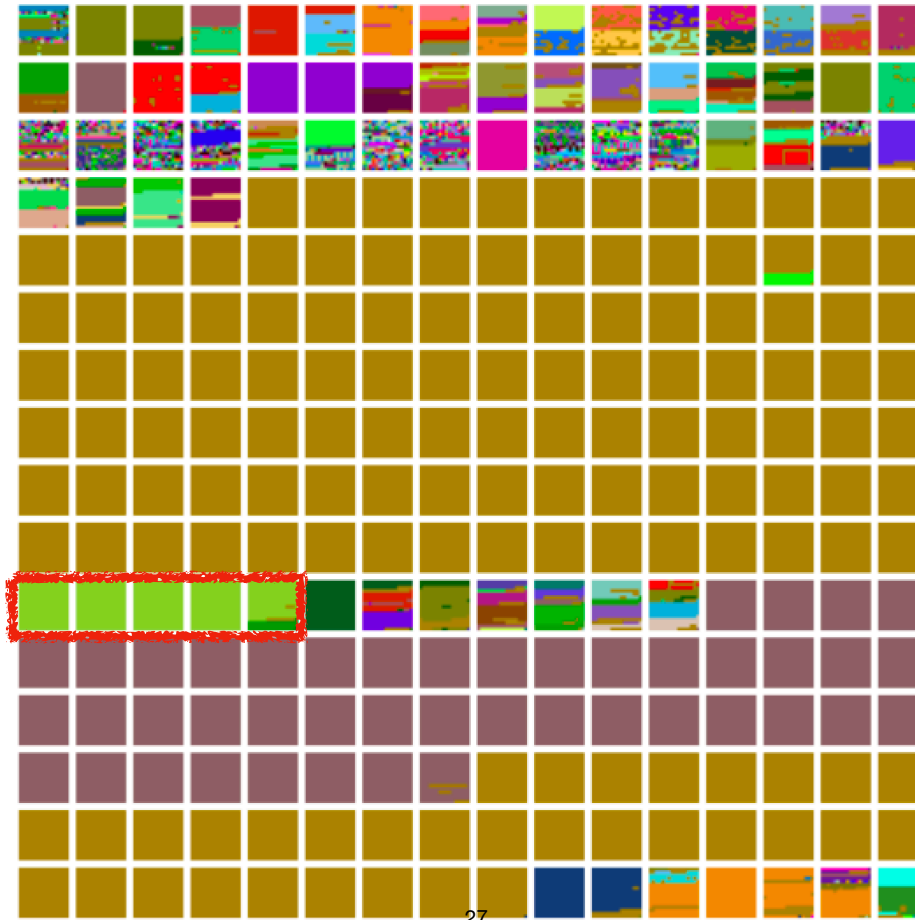


The first 65K characters



Yi

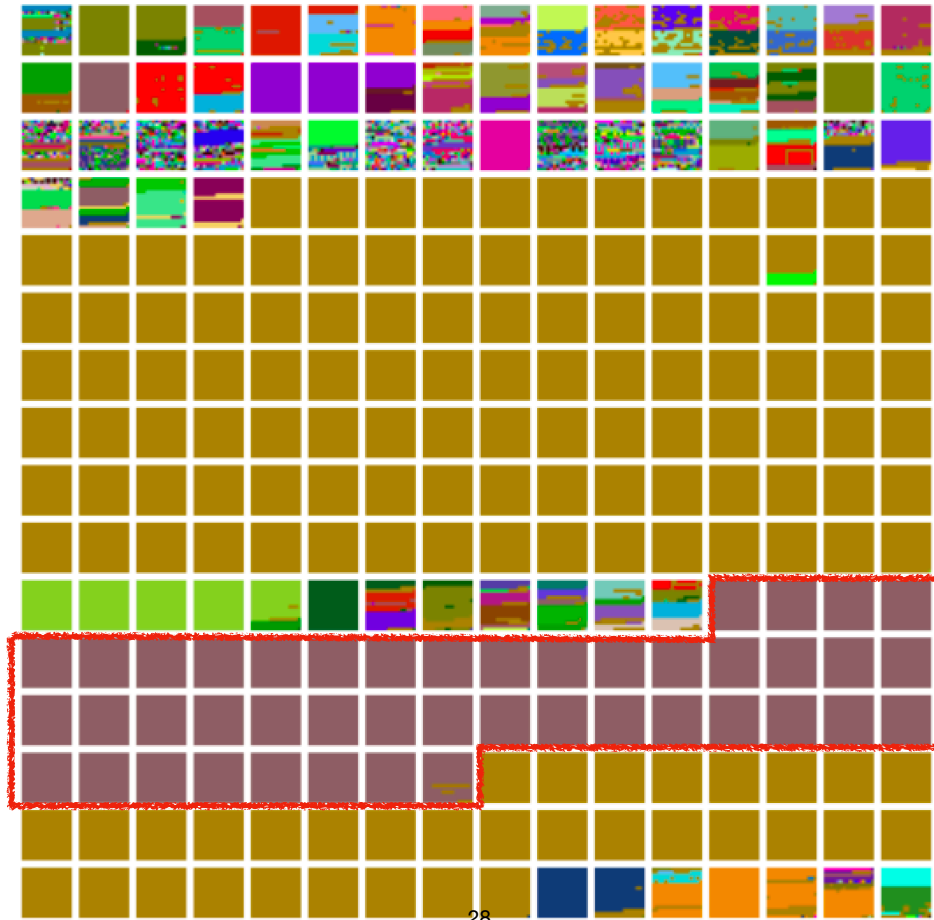
(ethnic group in China)

[illegible]

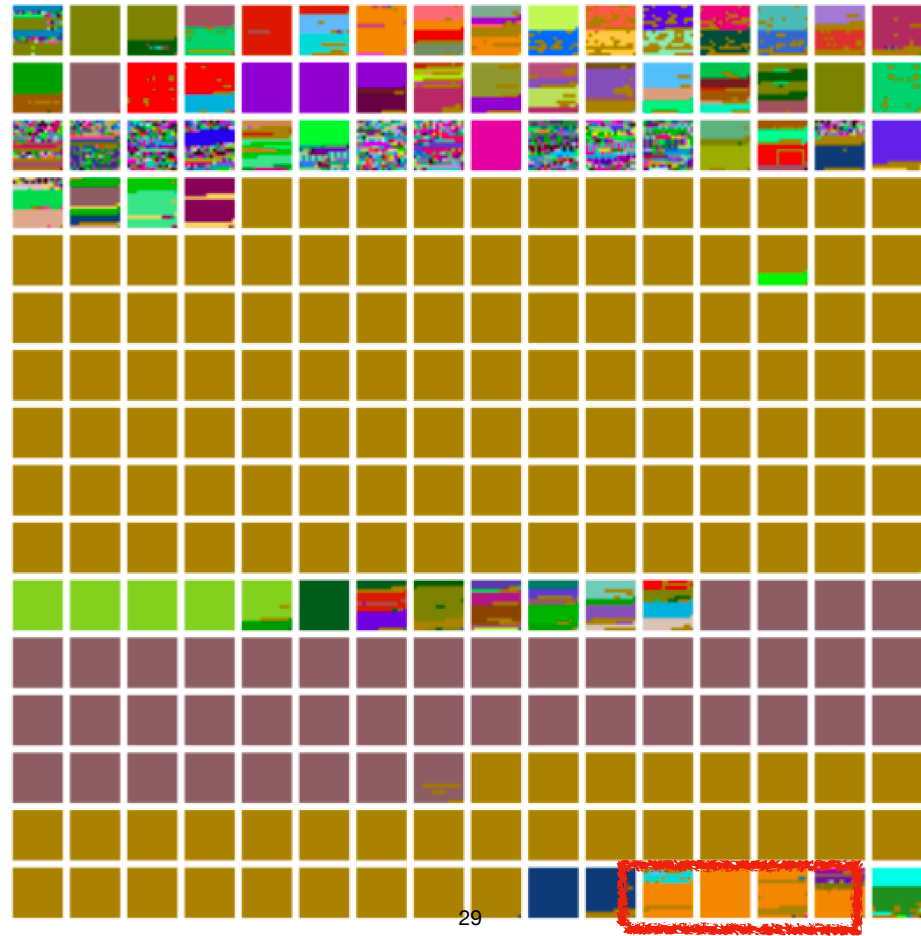
The first 65K characters

Hangul
(Korean alphabet)

조선글
한글



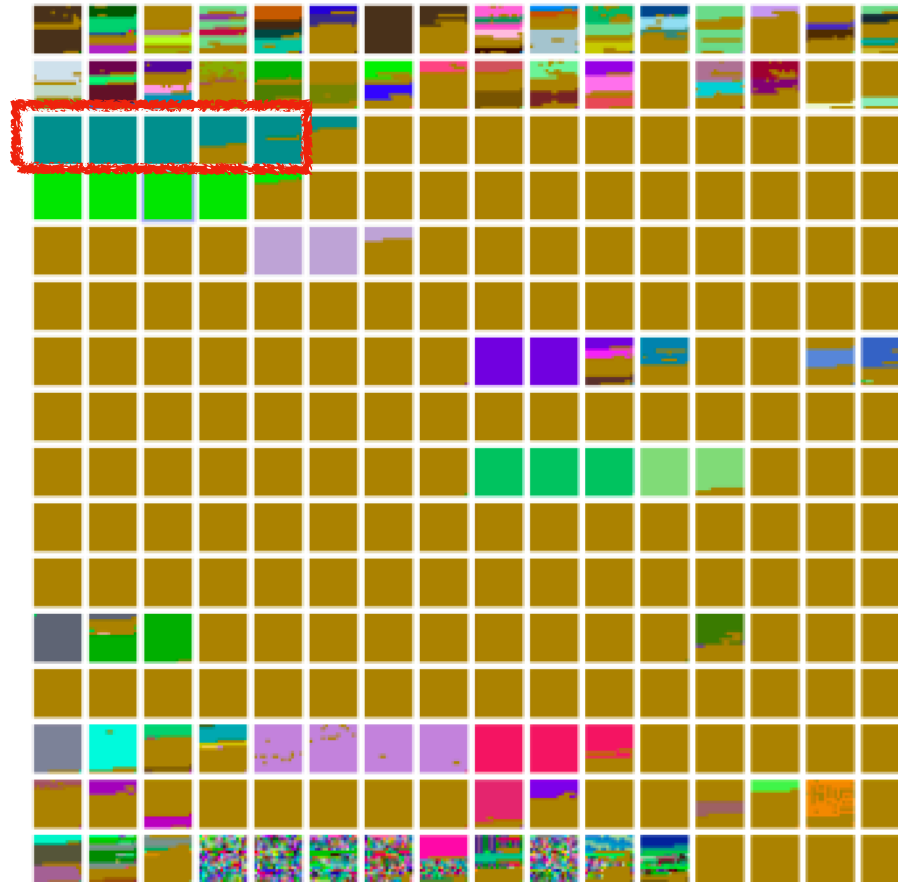
The first 65K characters



Arabic

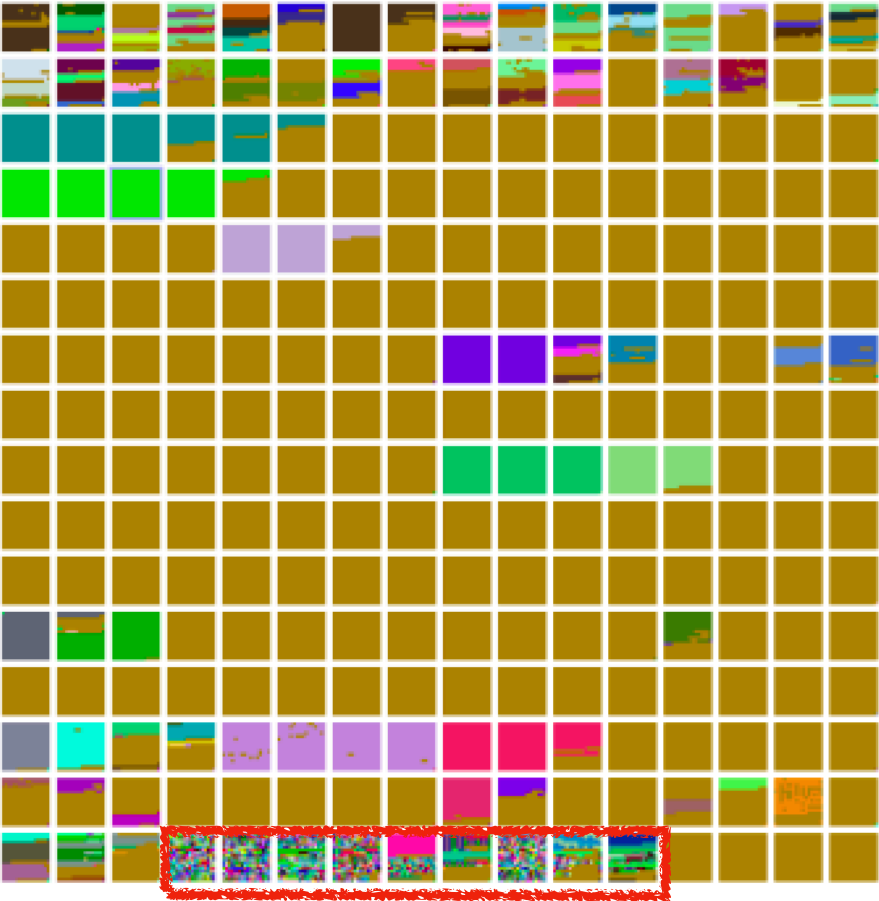
The second 65K characters

Cuneiform, invented
by Sumerians in
ancient
Mesopotamia



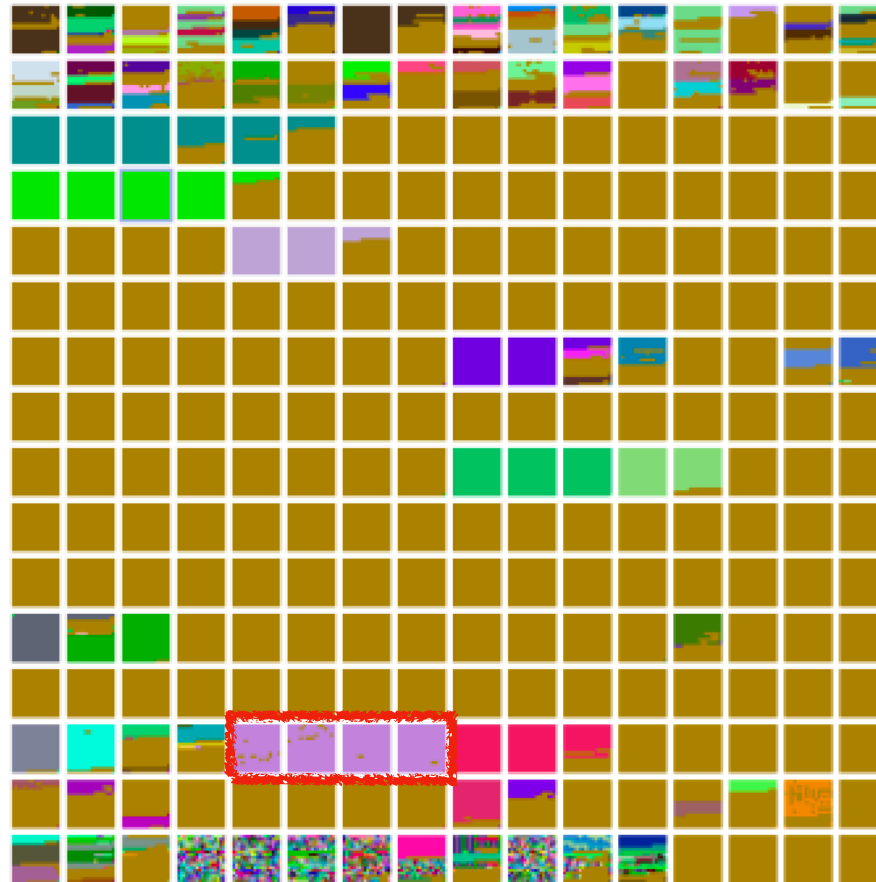
The second 65K characters

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
U+1F6Dx	👊	🍅	🛒			🔥	📦	📦								
U+1F6Ex	🔧	🛡️	🍷	🏠	✉️	🚢				✈️		🛩️	🛩️			
U+1F6Fx	🦋			🚢	🛴	🛵	🛶	🛶	🛶	🛶	🛶	📦	📦			
U+1F7Ex	🟠	🟡	🟢	🟣	🟤	🟥	🟦	🟧	🟨	🟩	🟪	🟫				
U+1F90x													📦	💜	💛	👉
U+1F91x	😄	😄	😄	😄	😄	😄	😄	👉	👉	👉	👉	👉	👉	👉	👉	👉
U+1F92x	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄	😄
U+1F93x	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉
U+1F94x	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉
U+1F95x	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉
U+1F96x	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉	👉



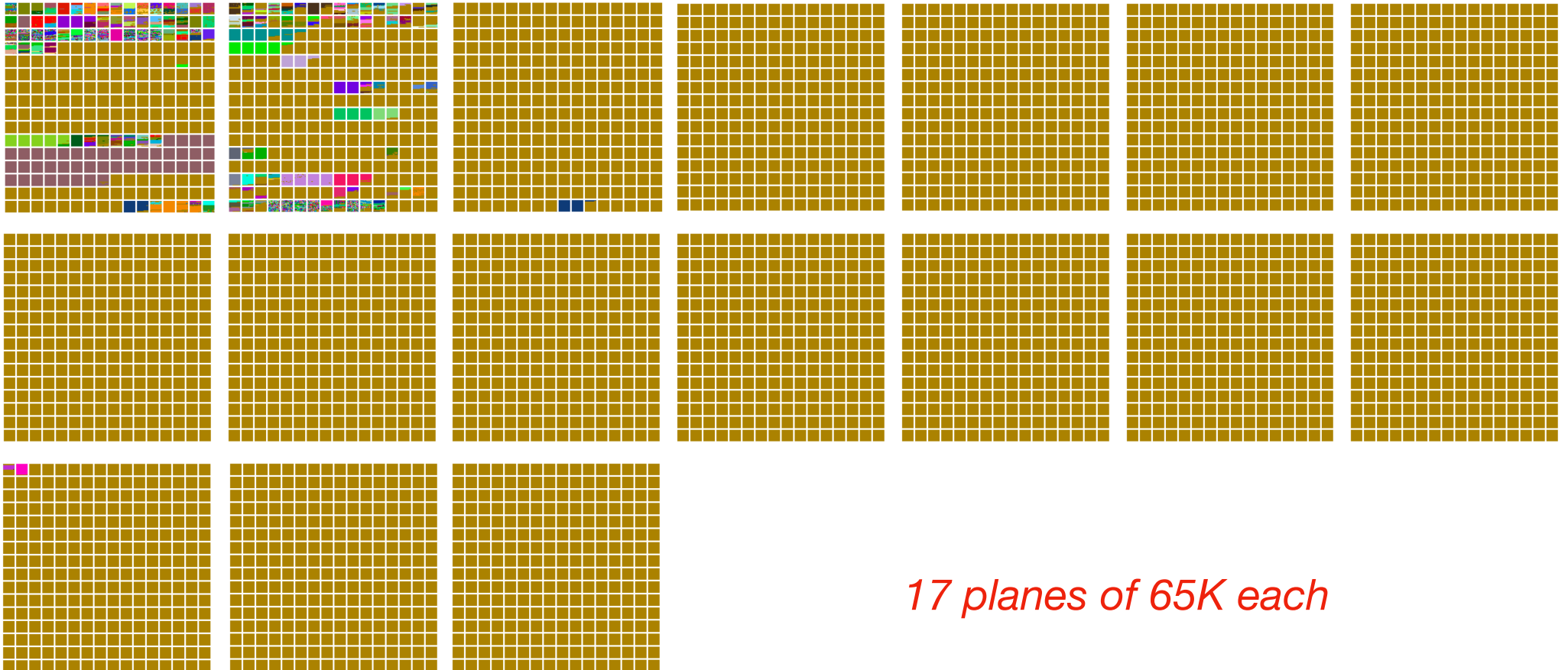
31emoji

The second 65K characters



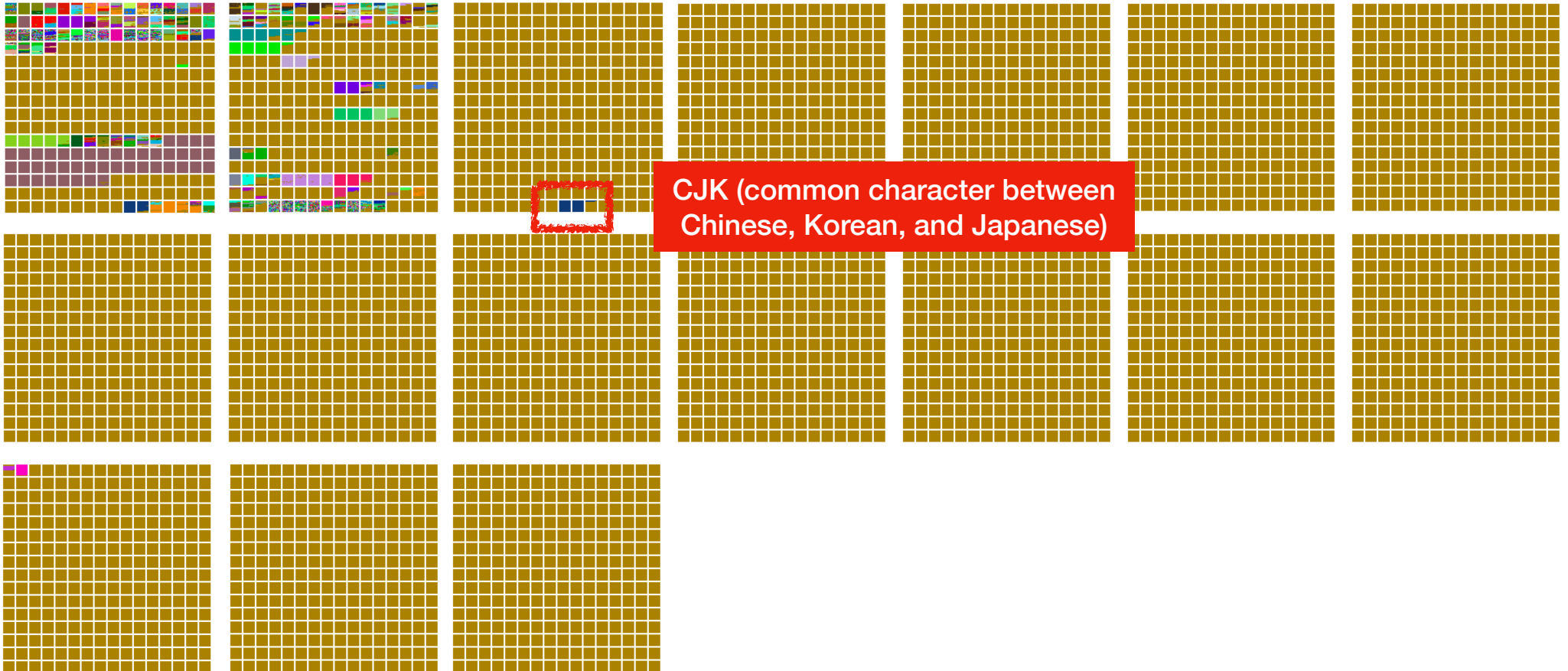
Mathematics

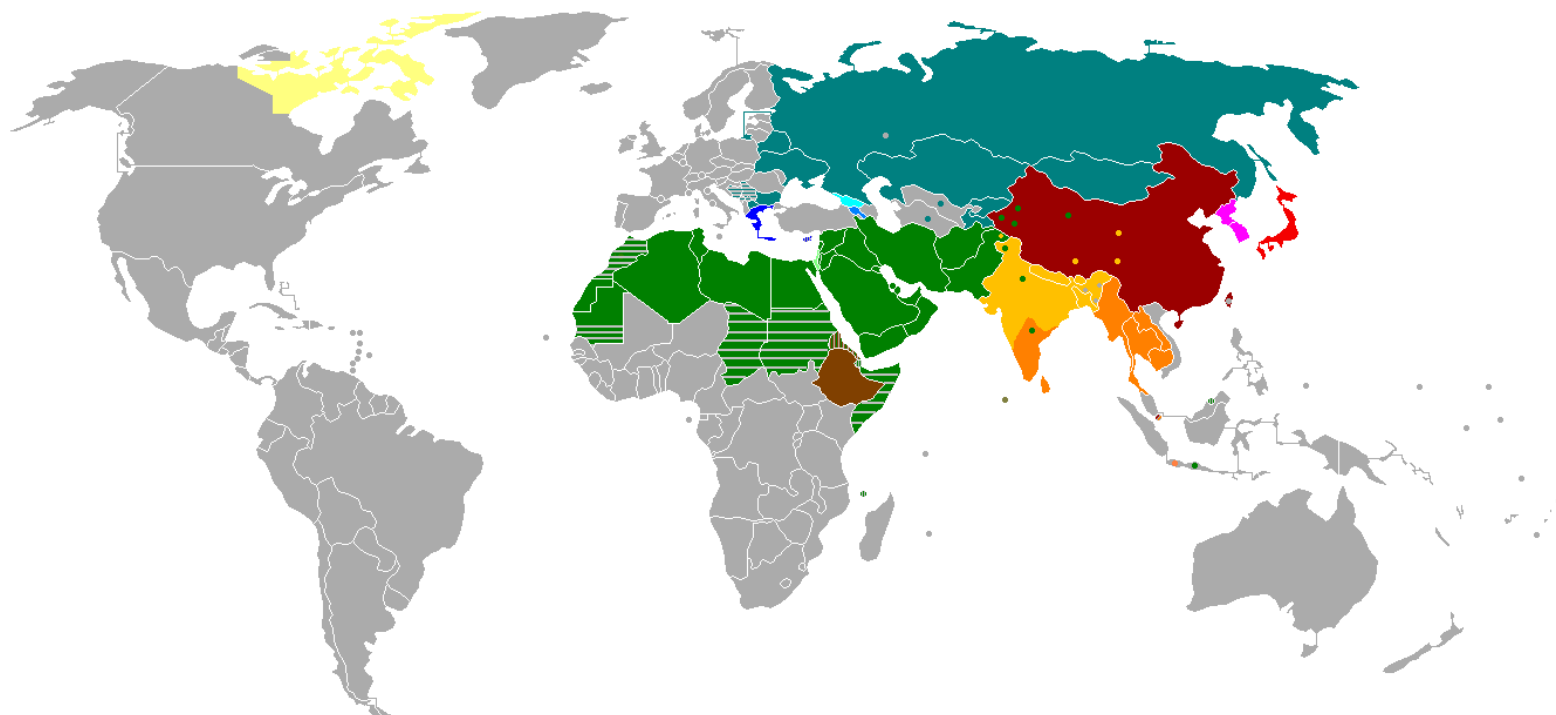
The 1 114 111 Unicode characters



17 planes of 65K each

The 1 114 111 Unicode characters





Index of predominant national and selected regional or minority scripts			
Alphabetic	[L]ogographic and [S]yllabic	Abjad	Abugida
Latin	Hanzi [L]	Arabic	North Indic
Cyrillic	Kana [S] / Kanji [L]	Hebrew	South Indic
Greek	Hanja ^b [L]		Ethiopic
Armenian			Thaana
Georgian			Canadian syllabic
Hangul ^a			

Character encoding forms: UTFs (level 3/4)

- Representation of a scalar value in a computer
- No escape: a simple juxtaposition is a concatenation

Character latin **A**

- abstract character:
 - the letter **A** of the Latin script
- coded character:
 - name: LATIN CAPITAL LETTER A
 - code point: U+0041
- encoding forms:
 - UTF-8: 41
 - UTF-16: 0041

Character Hiragana MA

- abstract character:
 - the letter ま of the Hiragana script (Japanese, each made of 3 strokes)
- coded character:
 - name: HIRAGANA LETTER MA
 - code point: U+307E
- encoding forms:
 - UTF-8: E3 81 BE
 - UTF-16: 307E

Unicode encodes characters, not glyphs

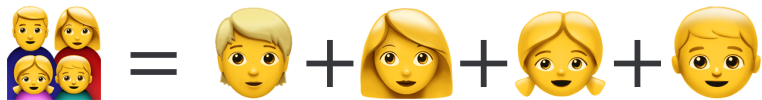
- The character U+0041 can equally well be displayed as **A**, **A**, **A**, **A**, **Ä**, ...
- Sometimes different glyphs may be required
 - Egg in French is written **œuf**
- going from characters to glyphs: *shaping*

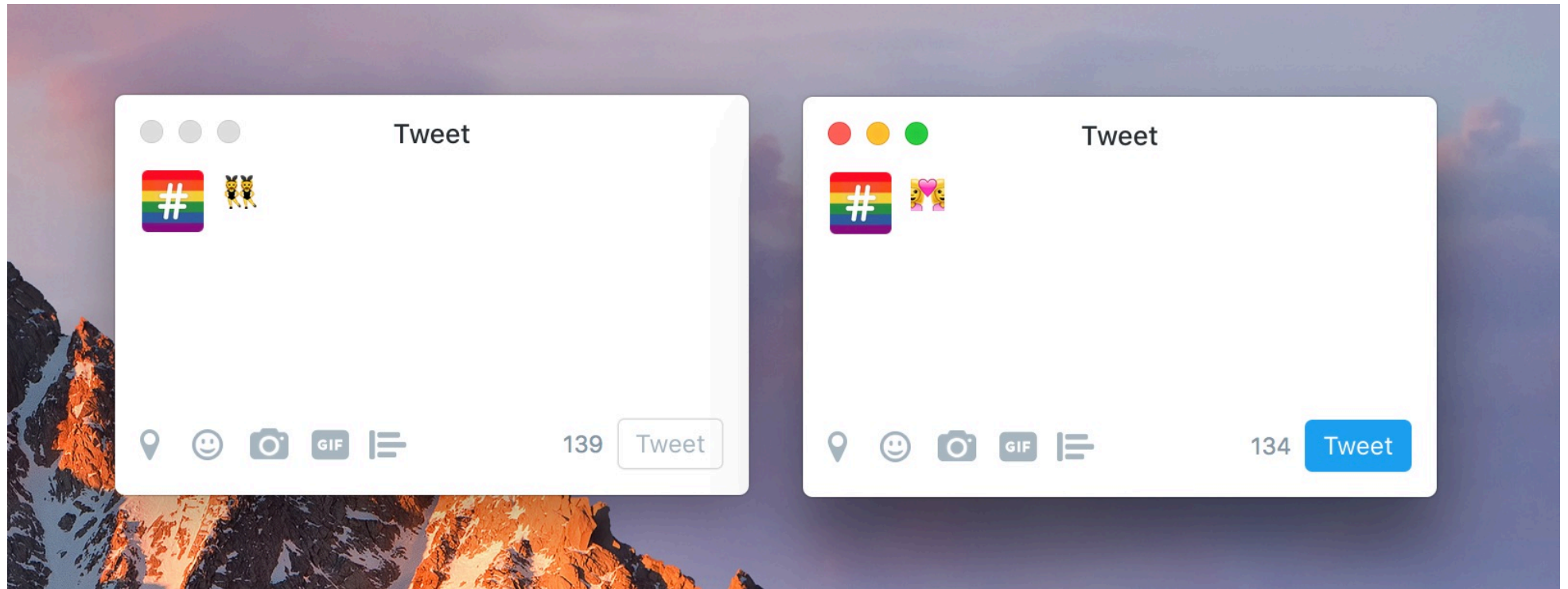
Zero Width Joiner (ZWJ)



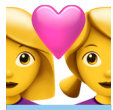


The invisible glue character is called a Zero Width Joiner (ZWJ)[1] and a sequence of emojis joined together with a ZWJ character is known as an Emoji ZWJ Sequence.

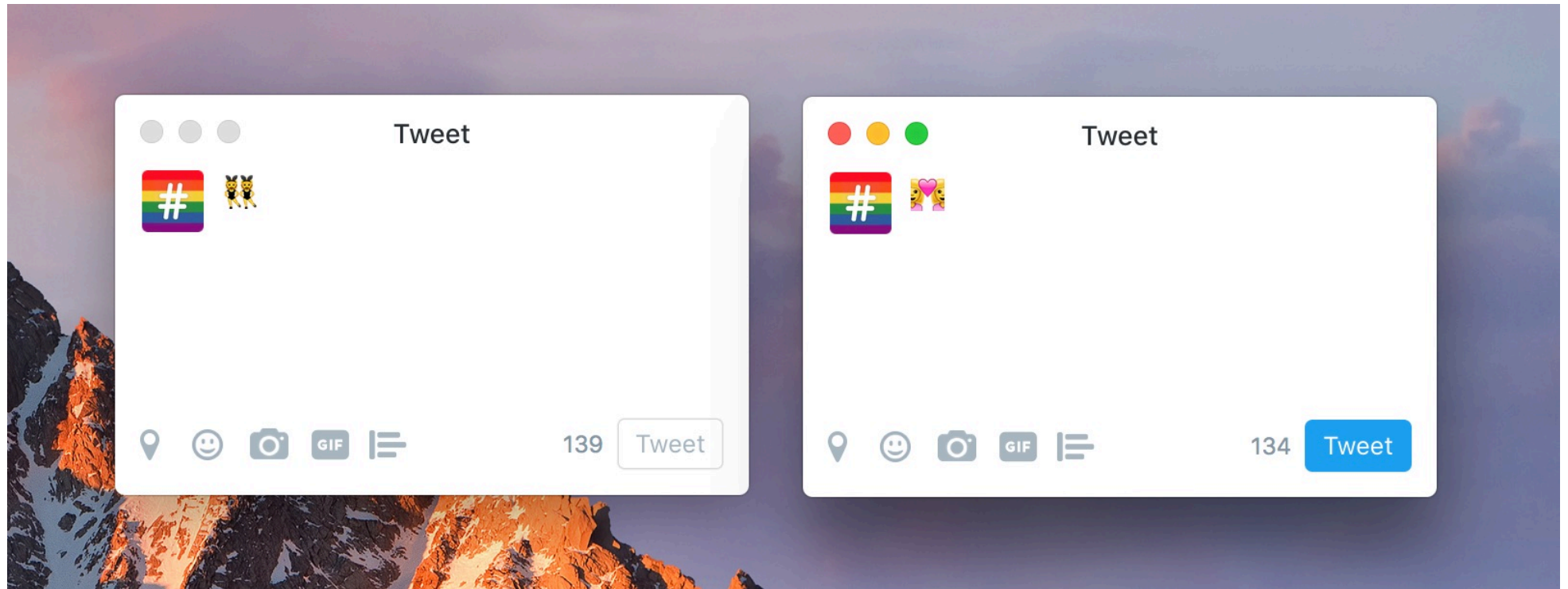




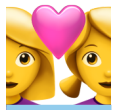
takes one character (i.e., it has a single code point)



includes three emojis but a total of six code points



takes one character (i.e., it has a single code point)

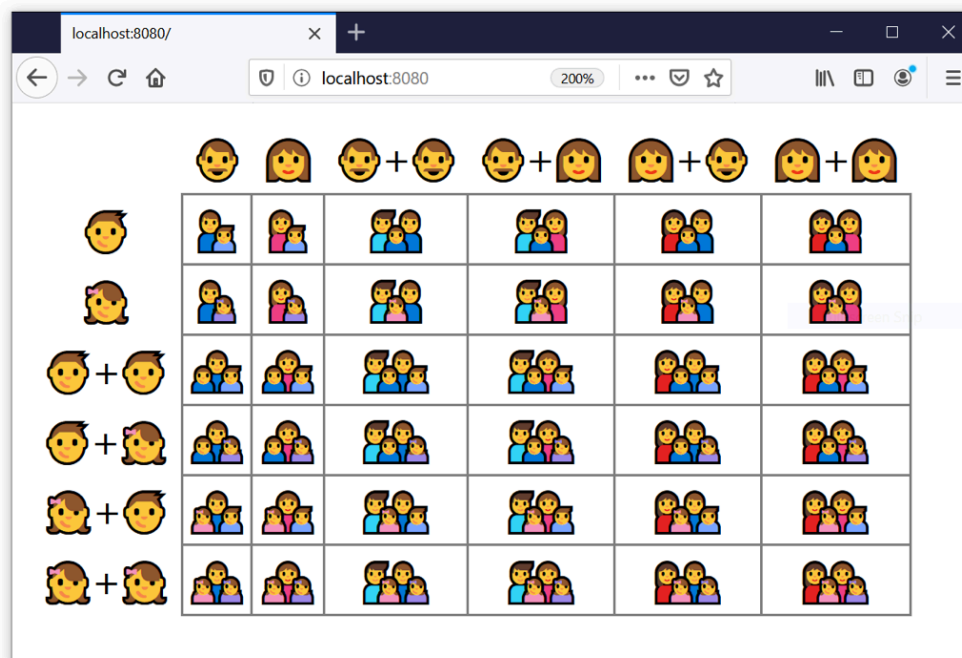


includes three emojis but a total of six code points

<https://blog.emojipedia.org/emoji-zwj-sequences-three-letters-many-possibilities/>

Building Multi-Character Emojis

Unicode has always supported building accented characters by combining letters and diacritics. This idea has been extended to meet the growing demand for emojis.



<https://www.fluentpython.com/extra/multi-character-emojis/>

Common practical issues



En alianza con **PayPal**

¡Estás listo! Ya puedes comenzar a utilizar nuestro servicio de retiros

Hola Alexandre Bergel,

Comienza a transferir tus fondos en dólares desde PayPal a tu cuenta bancaria nacional en pesos chilenos.

Solicita un retiro de tus fondos PayPal:

1. Inicia sesión en nuestro sitio web con tu RUT y Clave de 4 dígitos [clic aquí](#).
2. Ingresas el monto de tus fondos PayPal que quieres retirar.
3. Revisa el monto en dólares que vas a retirar, la tarifa y el monto en pesos que será depositado en tu cuenta bancaria.
4. Confirma la solicitud del retiro.

Recibirás tu dinero en tu cuenta bancaria en Chile en 5 días hábiles.

Te saluda atentamente,

Multicaja en asociación con PayPal

Para asegurarte de recibir correctamente nuestros emails, agrega contacto@multicaja.cl a tu lista de direcciones.

Por favor no respondas este email. Si necesitas contactarte con nosotros llámanos al 600 363 20 20



La información contenida en este mensaje y/o sus archivos adjuntos es de propiedad del emisor y de carácter confidencial o privilegiada y está destinada al uso exclusivo de éste. Si usted no es el destinatario, le informamos que cualquier almacenamiento, divulgación, distribución o copia de esta información está prohibida y sancionada por la ley. Si usted recibe este mensaje por error, le rogamos comunicarlo al remitente y destruirlo. El emisor no puede asegurar que este mensaje se encuentre libre de error, virus o interferencia.

[REDACTED]
Re: CIEL 2014 : Demande d'Ã©valuation du papier N°6

To: Alexandre Bergel, Cc: [REDACTED]

Bonjour,

Oui, si cela vous est plus commode, pas de soucis.

- Christelle et Marie-Agnès

----- Mail original -----

[See More](#) from Alexandre Bergel

```

public class Test {

    public static void main(String[] args) throws Exception {
        String french = "Les élèves ont des œufs";
        String spanish = "Las niñas y los niños";
        String japanese = "日本語";

        System.out.println(french);
        System.out.println(spanish);
        System.out.println(japanese);

        System.out.println("UTF-8 French: " + new String(french.getBytes("UTF-8")));
        System.out.println("UTF-8 Spanish: " + new String(spanish.getBytes("UTF-8")));
        System.out.println("UTF-8 Japanese: " + new String(japanese.getBytes("UTF-8")));

        System.out.println("ISO-8859-1 French: " + new String(french.getBytes("ISO-8859-1")));
        System.out.println("ISO-8859-1 Spanish: " + new String(spanish.getBytes("ISO-8859-1")));
        System.out.println("ISO-8859-1 japanese: " + new String(japanese.getBytes("ISO-8859-1")));
    }
}

```

```

public class Test {

    public static void main(String[] args) throws Exception {
        String french = "Les élèves ont des œufs";
        String spanish = "Las niñas y los niños";
        String japanese = "日本語";

        System.out.println(french);
        System.out.println(spanish);
        System.out.println(japanese);

        System.out.println("UTF-8 French: " + new String(french.getBytes("UTF-8")));
        System.out.println("UTF-8 Spanish: " + new String(spanish.getBytes("UTF-8")));
        System.out.println("UTF-8 Japanese: " + new String(japanese.getBytes("UTF-8")));

        System.out.println("ISO-8859-1 French: " + new String(french.getBytes("ISO-8859-1")));
        System.out.println("ISO-8859-1 Spanish: " + new String(spanish.getBytes("ISO-8859-1")));
        System.out.println("ISO-8859-1 Japanese: " + new String(japanese.getBytes("ISO-8859-1")));
    }
}

```

Les élèves ont des œufs

Las niñas y los niños

日本語

UTF-8 French: Les élèves ont des œufs

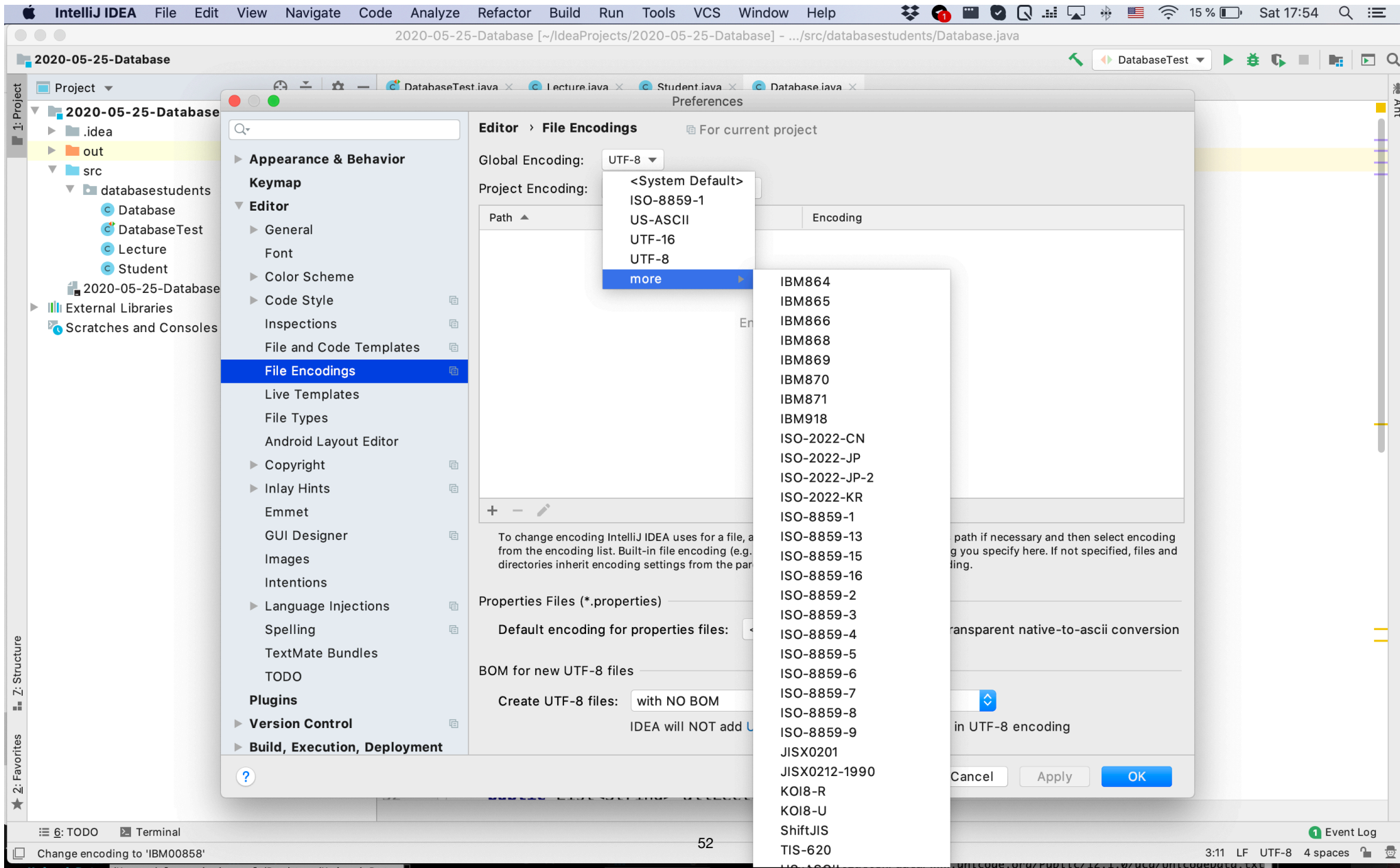
UTF-8 Spanish: Las niñas y los niños

UTF-8 Japanese: 日本語

ISO-8859-1 French: Les ?l?ves ont des ?ufs

ISO-8859-1 Spanish: Las ?i?as y los ?i?os

ISO-8859-1 Japanese: ???



TestUnicode [~/IdeaProjects/TestUnicode] - .../src/Test.java

TestUnicode

- src
 - Test

Test.java

```
1 public class Test {
2
3     public static void main(String[] args) throws Exception {
4         String french = "Les Ã©lÃ¨ves ont des œufs";
5         String spanish = "Las niñas y los niños";
6         String japanese = "日本語";
7
8         System.out.println(french);
9         System.out.println(spanish);
10        System.out.println(japanese);
11
12        System.out.println("UTF-8 Spanish: " + new String(spanish.getBytes( charsetName: "UTF-8")));
13        System.out.println("UTF-8 Japanese: " + new String(japanese.getBytes( charsetName: "UTF-8")));
14
15        System.out.println("UTF-8 French: " + new String(french.getBytes( charsetName: "UTF-8")));
16        System.out.println("UTF-8 Spanish: " + new String(spanish.getBytes( charsetName: "UTF-8")));
17        System.out.println("UTF-8 Japanese: " + new String(japanese.getBytes( charsetName: "UTF-8")));
18
19        System.out.println("ISO-8859-1 French: " + new String(french.getBytes( charsetName: "ISO-88
```

Run: Test

Les Ã©lÃ¨ves ont des œufs
Las niñas y los niños
æ ¥æ -èª
UTF-8 Spanish: Las Ã Ä±Ã Ä±as y los Ã Ä±Ã Ä±os
UTF-8 Japanese: Ã¡Ã Ä¥Ã¡Ã Ä-Ã-ÃªÃ
UTF-8 French: Les Ã Ä©Ã Ä¨ves ont des Ã Ä ufs
UTF-8 Spanish: Las Ã Ä±Ã Ä±as y los Ã Ä±Ã Ä±os
UTF-8 Japanese: Ã¡Ã Ä¥Ã¡Ã Ä-Ã-ÃªÃ
ISO-8859-1 French: Les Ã©lÃ¨ves ont des œufs

Terminal Messages Run TODO

Build completed successfully in 3 s 266 ms (7 minutes ago)

53

14:1 LF UTF-8 4 spaces

Useful tips

- Never, ever use Windows' Notepad to write code
 - It uses ISO 8859-1, which leads to numerous problems
- *Use UTF-8* all the way down:
 - For the text files used in your program
 - For the source code used in your program
- Make sure that your programming environment is set to UTF-8
- Always specify UTF-8 when *opening for reading or writing a textual file*

Useful tips

Text files:

```
Reader reader = new InputStreamReader(new FileInputStream("/tmp/foo.txt", "UTF-8"));  
Writer writer = new OutputStreamWriter(new FileOutputStream("/tmp/foo.txt", "UTF-8"));
```

String files:

```
byte[] bytesInDefaultEncoding = someString.getBytes(); // May generate corrupt bytes.  
byte[] bytesInUTF8 = someString.getBytes("UTF-8"); // Correct.  
String stringUsingDefaultEncoding = new String(bytesInUTF8); // Unknown bytes becomes "?".  
String stringUsingUTF8 = new String(bytesInUTF8, "UTF-8"); // Correct.
```

Fun

```
→ T cat Unicode.java
\u0070\u0075\u0062\u006C\u0069\u0063\u0020\u0020\u0020\u0063\u006C\u0061\u0073\u0073\u0020\u0020
\u0055\u006E\u0069\u0063\u006F\u0064\u0065\u0020\u007B\u0020\u0070\u0075\u0062\u006C\u0069\u0063
\u0020\u0020\u0073\u0074\u0061\u0074\u0069\u0063\u0020\u0020\u0076\u006F\u0069\u0064\u0020\u0020
\u006D\u0061\u0069\u006E\u0020\u0028\u0020\u0053\u0074\u0072\u0069\u006E\u0067\u0020\u005B\u005D
\u0061\u0072\u0067\u0073\u0020\u0029\u0020\u007B\u0020\u0053\u0079\u0073\u0074\u0065\u006D\u002E
\u006F\u0075\u0074\u002E\u0070\u0072\u0069\u006E\u0074\u006C\u006E\u0028\u0022\u0049\u0022\u002B
\u0022\u0020\u0020\u002665\u0020\u0055\u006E\u0069\u0063\u006F\u0064\u0065\u0022\u0029\u003B\u007D\u007D
→ T javac Unicode.java
→ T java Unicode
I ♥ Unicode
```

```
\u0070\u0075\u0062\u006C\u0069\u0063\u0020\u0020\u0020\u0063\u006C\u0061\u0073\u0073\u0020\u0020
\u0055\u006E\u0069\u0063\u006F\u0064\u0065\u0020\u007B\u0020\u0070\u0075\u0062\u006C\u0069\u0063
\u0020\u0020\u0073\u0074\u0061\u0074\u0069\u0063\u0020\u0020\u0076\u006F\u0069\u0064\u0020\u0020
\u006D\u0061\u0069\u006E\u0020\u0028\u0020\u0053\u0074\u0072\u0069\u006E\u0067\u0020\u005B\u005D
\u0061\u0072\u0067\u0073\u0020\u0029\u0020\u007B\u0020\u0053\u0079\u0073\u0074\u0065\u006D\u002E
\u006F\u0075\u0074\u002E\u0070\u0072\u0069\u006E\u0074\u006C\u006E\u0028\u0022\u0049\u0022\u002B
\u0022\u0020\u0020\u002665\u0020\u0055\u006E\u0069\u0063\u006F\u0064\u0065\u0022\u0029\u003B\u007D\u007D
```

The Unicode Consortium

Mission

- “making the digital world more inclusive”
- “Everyone in the world should be able to use their own language on phones and computers.”


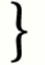




























Membership

- The effort of the Unicode Consortium is sponsored by membership fees
- Anyone can be a membership (student = 35 USD, corporation up to 21K USD)
- A member
 - has access to the technical discussions
 - can vote on many aspects of the consortium
- Actual members are large IT companies and some governments

Adopt a character


- Do you want your company to be associated with the 🍔 emoji?
- Do you want to declare your love with a 💍 or 💘 dedication for your partner?

54 Gold Sponsors

 Elastic	 Elastic	 Adobe Inc.
 Catalyst IT Limited	 Ann Lewnes and Greg Welch	 Suzi Slavik
 Yaya Wang	 Enjoy the present (致独一无二的莎莎)	 Cal Henderson
 Jason Jenkins	 Digital Love Media	 discourse.org
 Vinton G. Cerf	 In loving memory of Ryan Neale Cross (1984-2017)	 R. Martin Chavez & Adam Norbury
 White Unicorn Agency	 Oakland Athletics	 Michael D'Errico
 Norton Pedee Family	 Indian Type Foundry	 Candice DeStefano
 Oakland Athletics	 USA Pears	 S&A Group
 Zespri Kiwifruit	 Avocados From Mexico	 Scott Spears
 Richard Shupak	 Buffalo Wild Wings	 Google



Google Design  @GoogleDesign · Feb 14, 2018

No matter where the cheese goes, now it's official → [#GoogleFonts](#) adopted the [@unicode](#) [#EmojiBurger](#) on behalf of [@Google](#)!  



 8

 35

 132



Final words

Unicode

- *Awesome environment* to learn about a central aspect to many other cultures
- UTF-8, UTF-8, UTF-8, UTF-8, UTF-8, UTF-8, UTF-8, UTF-8, ...
- Make sure you use
 - a proper professional programming environment and
 - text editing tool



NEWS, ANNOUNCEMENTS, RELEASE INFO, AND CALENDAR UPDATES FROM THE UNICODE CONSORTIUM


TUESDAY, SEPTEMBER 14, 2021

Announcing The Unicode® Standard, Version 14.0

Version 14.0 of the Unicode Standard is now available, including the core specification, annexes, and data files. This version adds 838 characters, for a total of 144,697 characters. These additions include five new scripts, for a total of 159 scripts, as well as 37 new emoji characters.

The new scripts and characters in Version 14.0 add support for modern language groups in Bosnia, India, Indonesia, Iran, Java, Malaysia, Mongolia, Myanmar, Pakistan, and the Philippines, plus other languages in Africa and North America, including:

- Arabic script additions that include honorifics and additions for Quranic use, and characters used to write languages across Africa, the Balkans, and South and Southeast Asia

 10570	 10580	 10590
 10571	 10581	 10591
 10572	 10582	 10592

LINKS OF INTEREST

[What is Unicode?](#)

[The Unicode Consortium](#)

[Archived Announcements](#)

BLOG ARCHIVE

▼ [2021](#) (27)

► [November](#) (2)

► [October](#) (3)

▼ [September](#) (3)

[Announcing The Unicode® Standard, Version 14.0](#)

[Unicode CLDR v40 Alpha available for testing](#)

[Unicode Consortium Announces Version 14.0 Cover De...](#)

► [August](#) (1)

► [July](#) (2)

► [June](#) (1)

<http://blog.unicode.org/2021/09/announcing-unicode-standard-version-140.html>

References

- <https://www.unicode.org>
- <http://reedbeta.com/blog/programmers-intro-to-unicode/>
- <https://balusc.omnifaces.org/2009/05/unicode-how-to-get-characters-right.html>
- <https://www.joelonsoftware.com/2003/10/08/the-absolute-minimum-every-software-developer-absolutely-positively-must-know-about-unicode-and-character-sets-no-excuses/>
- <https://www.joelonsoftware.com/2003/10/08/the-absolute-minimum-every-software-developer-absolutely-positively-must-know-about-unicode-and-character-sets-no-excuses/>
- <https://en.wikipedia.org/wiki/Mojibake>
- <https://www.unicode.org/notes/tn23/Muller-Slides+Narr.pdf>
- <http://www.unicode.org/versions/Unicode13.0.0/ch02.pdf#G14527>

Introduction To Unicode

Making the digital world more inclusive

Alexandre Bergel

Computer Science Department - FCFM

University of Chile

<http://bergel.eu>

abergel@dcc.uchile.cl

[@AlexBergel](#)

Overview

- 29 publicly available version
- Version 1.1.5 established in 1995
- Version 14.0.0 defined in March 2020

Evolution of the Unicode standard

