

Strings&
Strings&
Strings&
Strings&
Strings&
Strings&
Strings&
Strings&
Strings

Swift Strings Seven Ways
a talk presented by Mattt
at try! Swift NYC 2018



R D W H G 9 & P A 5 Y 6 B P
5 M 6 oV e) u W B A A o h G A
d y 4 F C 7 o U T Z o C R h
S V F U E O T O P & O J K N
A O G G Y L 6 S S G i O T h G
P P F H L o G L t o 8 O E

Cwy

/dʒaɪə'gɪ/



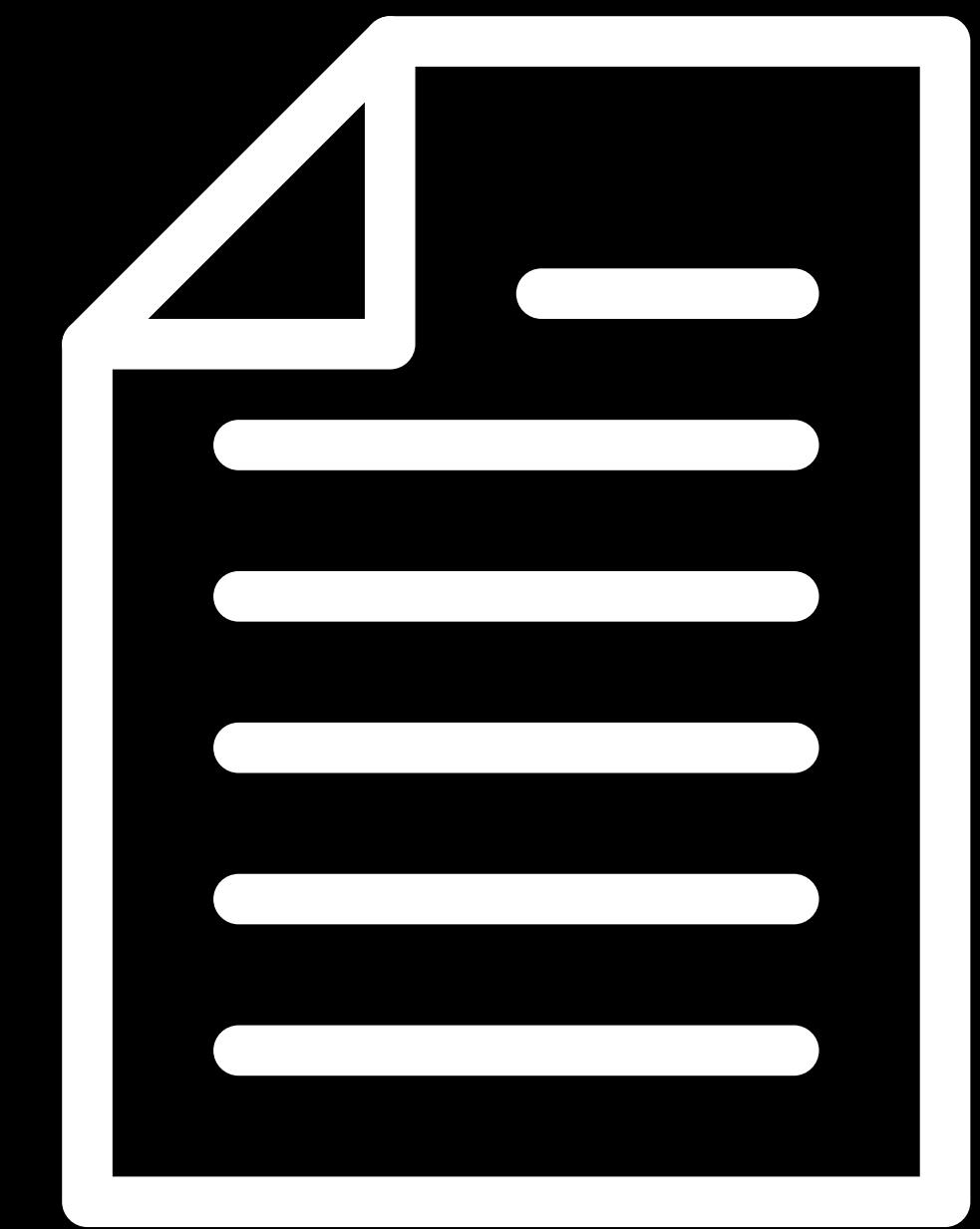
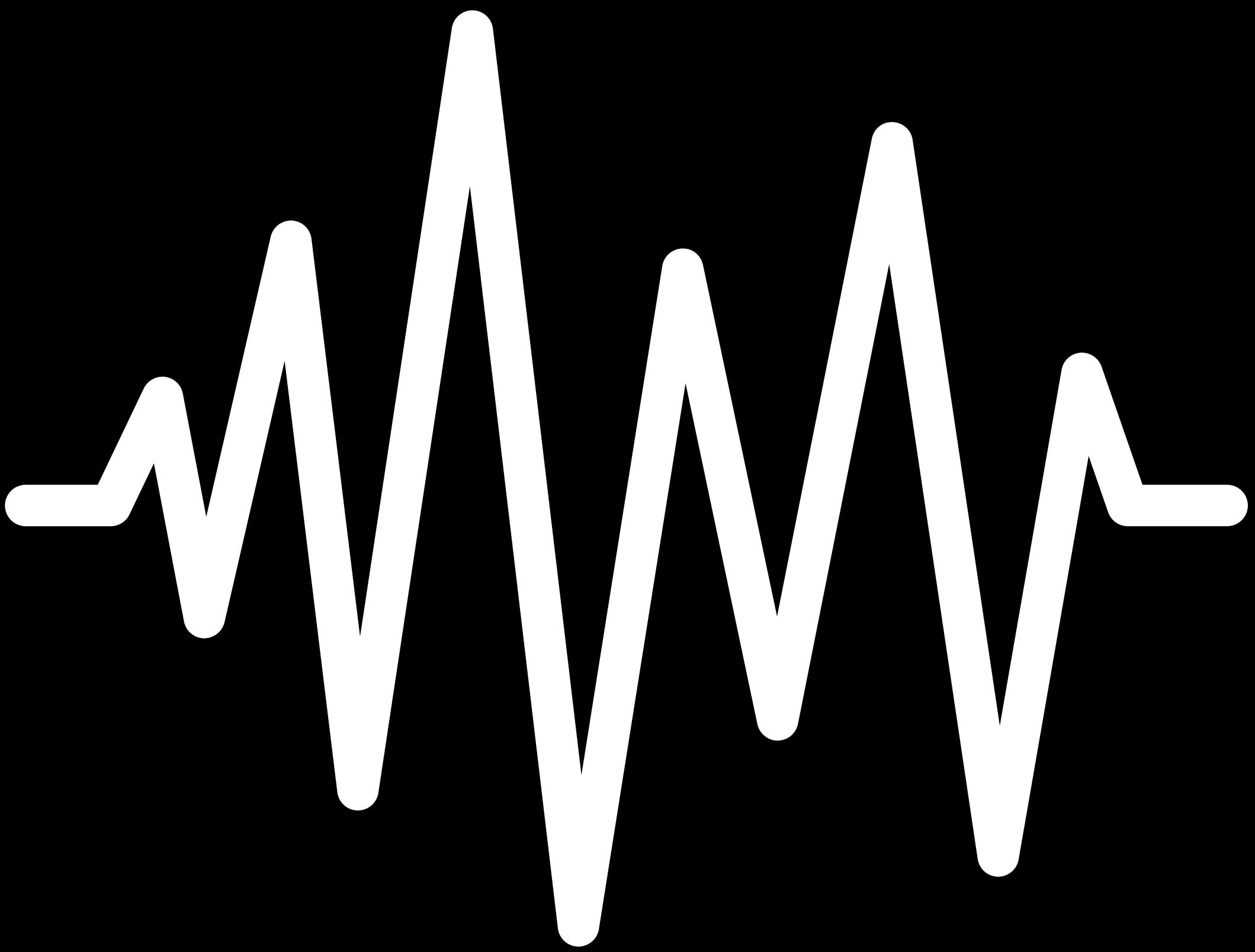
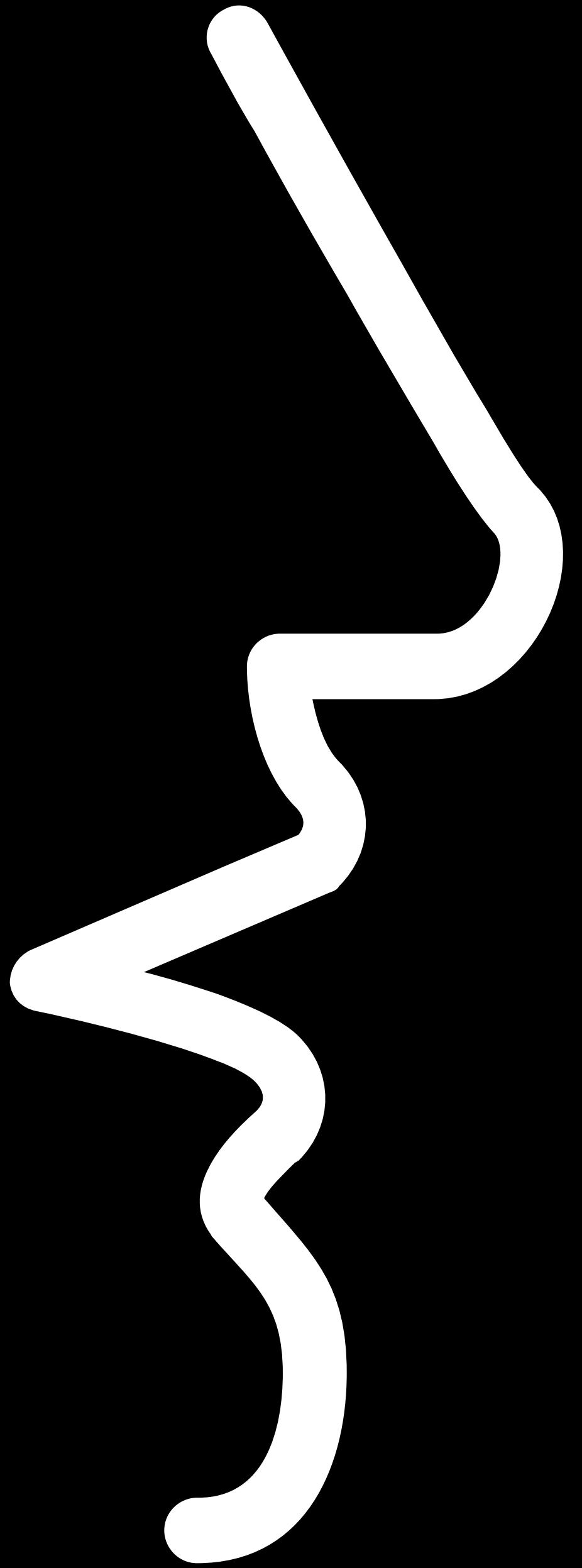
R D
A Z V B P W K
O O C W B M
H P A J S Q D
R H X Z Z Q F
T C S A P L E Q
Y B Q A J K P
G I D C N S R
H T C D B
Y L L



1. Unicode Characters
2. Identifiers
3. Paths
4. Natural Language Text
5. Structured Data
6. Encoded Data
7. Typographic Information

1

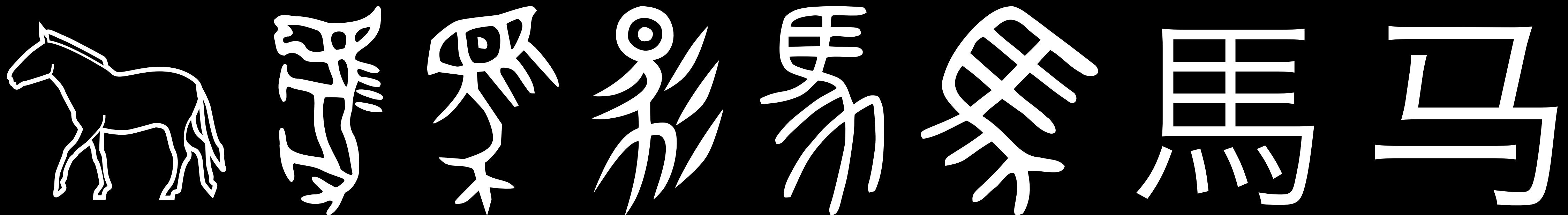
Collections of Unicode Characters



Western Writing Systems

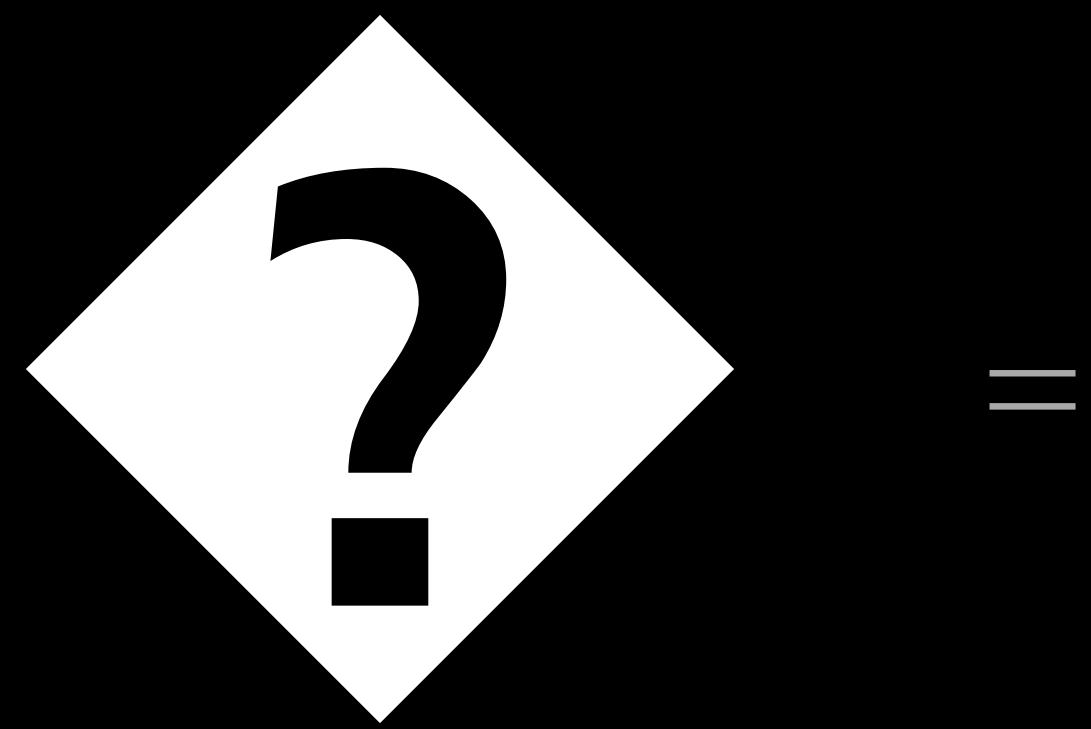
ox X A A

Eastern Writing Systems





=

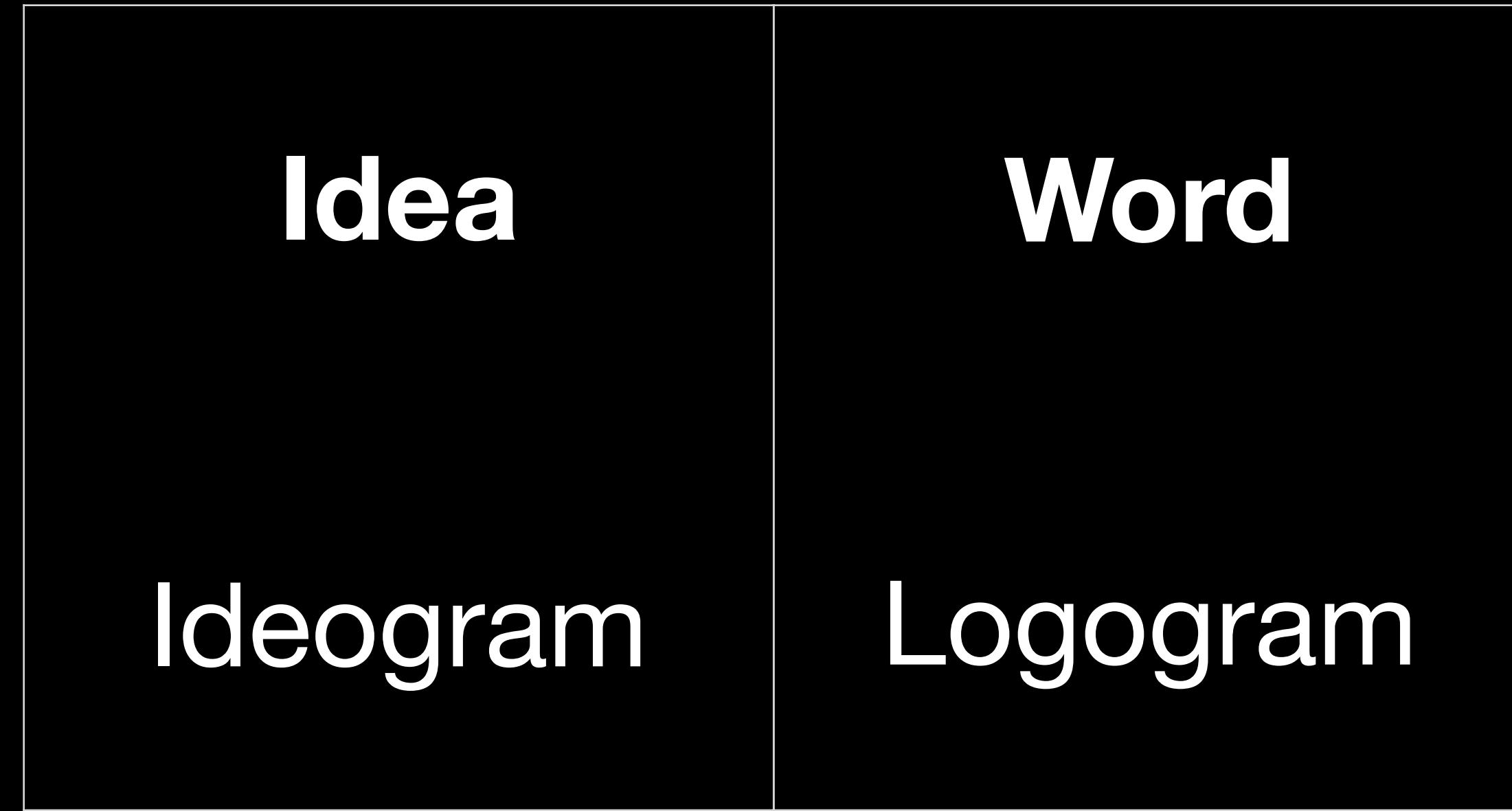


=

Idea
Ideogram



=





=

Idea	Word	Phoneme
Ideogram	Logogram	Alphabet



=

Idea	Word	Phoneme
Ideogram	Logogram	Alphabet
Consonant + Vowel		
Abugida		



=

Idea	Word	Phoneme
Ideogram	Logogram	Alphabet
Consonant + Vowel	Consonant	
Abugida	Abjad	



=

Idea	Word	Phoneme
Ideogram	Logogram	Alphabet
Consonant + Vowel	Consonant	Syllable / Mora
Abugida	Abjad	Syllabary

Alphabets

Latin / Кирилица / Ελληνικό / Հայութ / 한글

Abjads

العَرَبِيَّةُ / אַלְפָבֵיט עַבְרִי / اُردو

Abugidas

දෙවනාගරී / संस्कृतम् / ଶାହ / தெய்வு, தமிழ் / ՚ଓଡ଼ି

Syllaberies

ひらがな / カタカナ / Gwy

Logograms

漢字 / 汉字

Ideograms

👋😊 / * / ☢

Byte
char

ASCII (1960)

American Standard Code for
Information Interchange

	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	s _P	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
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

“Hello”

H	1000100
e	0101110
I	1100110
I	1100110
o	1111110
NUL	0000000

„Hvåd med mine økæler?”

— Dansk

「つまらない質問で恐縮ですが
¥はどこにありますか？」

– 日本語

ISO/IEC 646 (1964)

	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	s _P	!	"			%	&		()	*		+	-	.	/
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

ISO/IEC 646 (-INV)

US / IRV	#	\$	@	[\]	^	`	{		}	~
DE	#	\$	§	Ä	Ö	Ü	^	`	ä	ö	ü	ß
DK	#	¤	@	Æ	Ø	Å	Ü	`	æ	ø	å	ü
FR	£	\$	à	°	ç	§	^	μ	é	ù	è	..
JP	#	\$	@	[¥]	^	`	{		}	-
...												

C:¥Program Files¥



ISO/IEC 8859 (1987)

Part 1	Latin-1 Western European
Part 2	Latin-2 Central European
Part 3	Latin-3 South European
Part 4	Latin-4 North European
Part 5	Latin/Cyrillic
Part 6	Latin/Arabic
Part 7	Latin/Greek
Part 8	Latin/Hebrew
Part 9	Latin-5 Turkish
Part 10	Latin-6 Nordic
Part 11	Latin/Thai
Part 12	Latin/Devanagari
Part 13	Latin-7 Baltic Rim
Part 14	Latin-8 Celtic
Part 15	Latin-9
Part 16	Latin-10 South-Eastern European

8-bit $\Rightarrow 2^8 = 256$

Languages with Incomplete Coverage for ISO/IEC 8859

Latin 6

يولد جميع الناس أحراراً متساوين في الكرامة والحقوق، وقد وهبوا عقلاً
وضميراً وعليهم أن يعامل بعضهم بعضاً بروح الإخاء.

Latin 5

Bütün insanlar hür, haysiyet ve haklar bakımından eşit doğarlar.
Akıl ve vicdana sahiptirler ve birbirlerine karşı kardeşlik
zihniyeti ile hareket etmelidirler.

106,230

异体字字典

Entries in the Dictionary of Chinese Variant Form



Unicode (1991)

1,114,112

Unicode Code Points

UTF-8

UTF-16

UTF-32

A **String** is a collection of **Characters**

A **String** is a collection of **Characters**

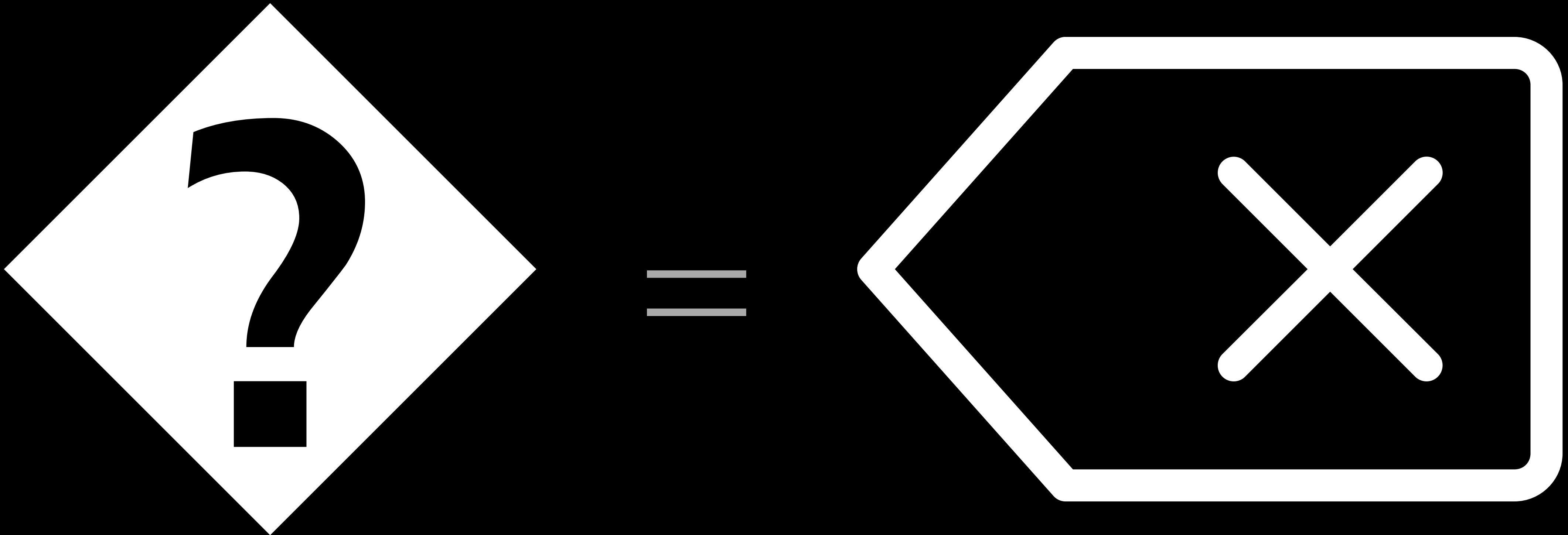
A **Character** consists of one or more **Unicode scalars**

A **String** is a collection of **Characters**

A **Character** consists of one or more **Unicode scalars**

Characters are **encoded** into a **sequence of bytes**

```
let string = "Café ☕"  
  
for character in string {  
    print(character)  
}
```



C

a

f

é

s_P



e + ' = é

e + ' = é

é = é

```
let character = string[2]
```

```
let character = string[2]  'subscript' is unavailable: cannot subscript
```

```
for character in string {
    print(character)

    for byte in "\u{character}".utf8 {
        print(String(byte, radix: 16, uppercase: true))
    }

    for byte in "\u{character}".utf16 {
        print(String(byte, radix: 16, uppercase: true))
    }

    for scalar in character.unicodeScalars {
        print(scalar.value)
    }
}
```

Character	C	a	f	é	s_p					
Codepoint	U+0043	U+0061	U+0066	U+00E9	U+0020		U+1F956			
Scalar Value	67	97	102	233	32		129366			
UTF-8	0x43	0x61	0x66	0xC3	0xA9	0x20	0xF0	0x9F	0xA5	0x96
UTF-16	0x0043	0x0061	0x0066	0x00E9	0x0020	0xD83E		0xDD56		

2

Identifiers

```
struct Product {  
    let identifier: String  
}
```

```
let coffee = Product(identifier: "☕")
```

```
let coffee = Product(identifier: "☕")
```

```
coffee.identifier.
```

M	Void append(c: Character)
M	Void append(contentsOf: Sequence)
M	Void append(contentsOf: String)
M	Void append(contentsOf: Substring)
M	Void append(other: String)
V	String.CharacterView characters
V	Int count
V	Mirror customMirror

Appends the given character to the string.

```
struct Product {  
    struct Identifier: Equatable, Hashable {  
        private let value: String  
    }  
  
    let identifier: Identifier  
}
```

```
let tea = Product(identifier:  
    Product.Identifier(value: "🍵"))
```

```
let string = "Hello!"
```

```
let id: Product.Identifier = "🍵"
```

```
extension Product.Identifier:  
    ExpressibleByStringLiteral  
{  
    init(stringLiteral value: String) {  
        self.init(value: value)  
    }  
}
```

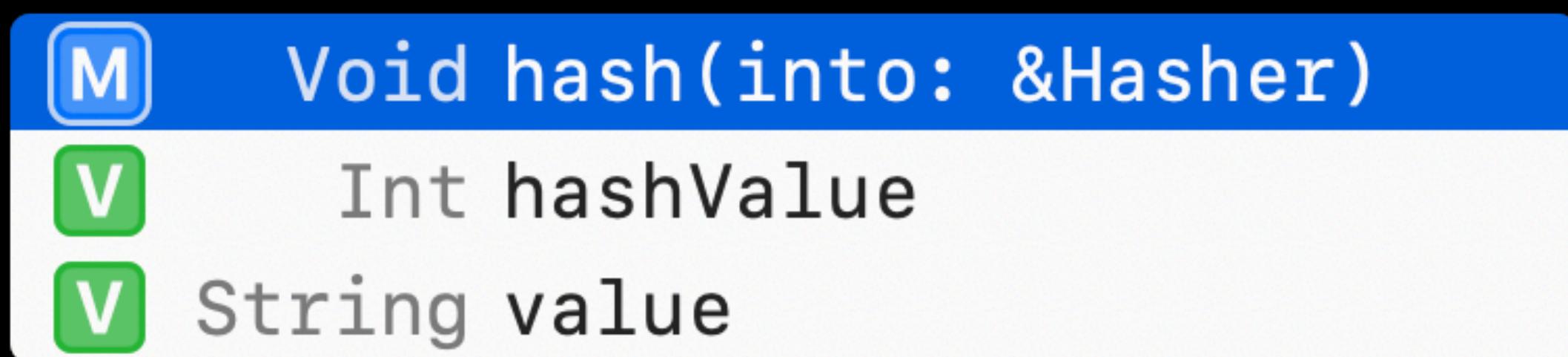
```
let id: Product.Identifier = "🍵"  
let tea = Product(identifier: id)
```

```
let tea = Product(identifier: "🍵")
```

```
let tea = Product(identifier: "🍵")
```

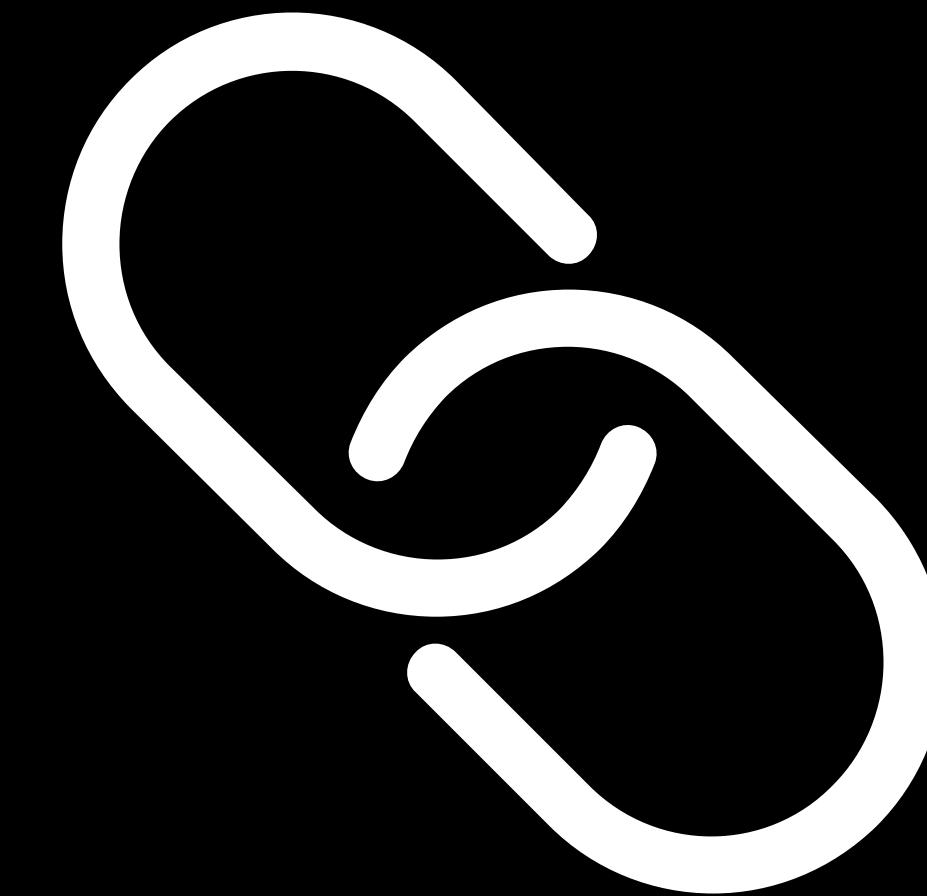
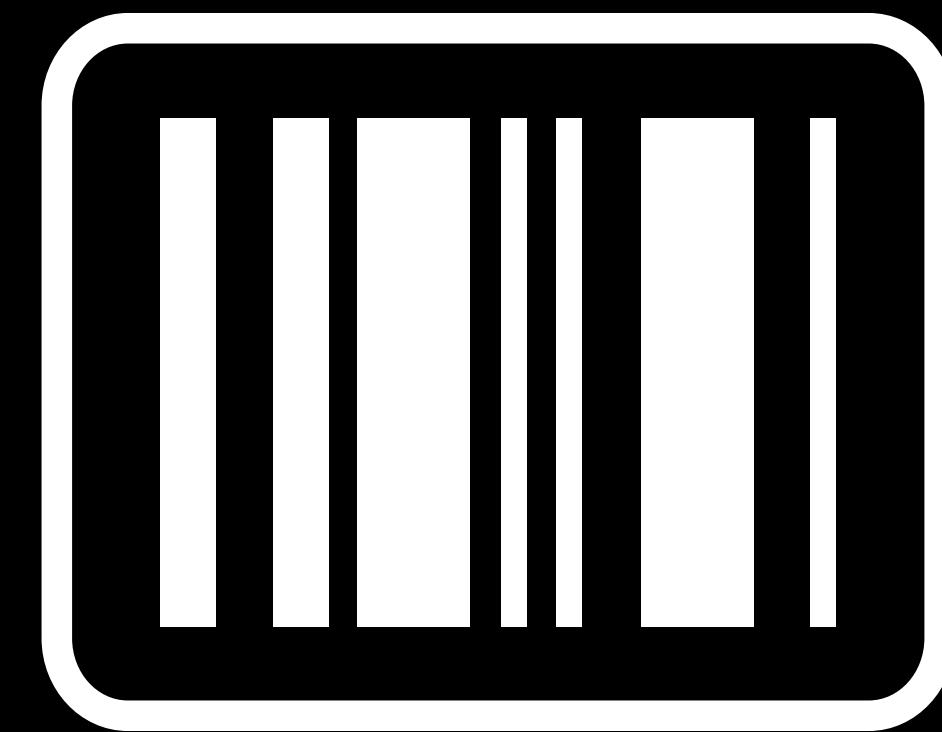
```
let tea = Product(identifier: "🍵")
```

```
tea.identifier.
```



```
struct Identifier<Value>: Equatable, Hashable
    where Value : Equatable & Hashable
{
    private let value: Value
}
```

```
struct Product {
    let identifier: Identifier<String>
}
```



Validation

```
init(value: String) {  
    precondition(!value.isEmpty, "cannot be blank")  
    precondition(value.count < 100, "too long")  
    for scalar in value.unicodeScalars {  
        guard (0x1F344...0x1F37F).contains(scalar.value) else {  
            preconditionFailure("must be food emoji")  
        }  
    }  
  
    self.value = value  
}
```

Normalization

```
init(value: String) {  
    self.value = value.lowercased()  
}
```

3

Paths

scheme://user:password@
hostname:port/path#fragment?query

URL (RFC 3986)

```
import Foundation  
  
let url = URL(string: "https://tryswift.co/")!
```

```
import Foundation

var components =
    URLComponents(string: "https://tryswift.co/")!

components.path = "/events/2018/nyc/"

components.url
// https://tryswift.co/events/2018/nyc/
```

Percent-Encoding / URI Encoding

```
import Foundation

var components =
    URLComponents(string: "https://tryswift.co/")!
components.path = "/events/2018/ /"

components.url
// https://tryswift.co/events/2018/%F0%9F%97%BD/
```

%F0%9F%97%BD



STATUE OF LIBERTY

(U+1F5FD)

UTF-8

F0

9F

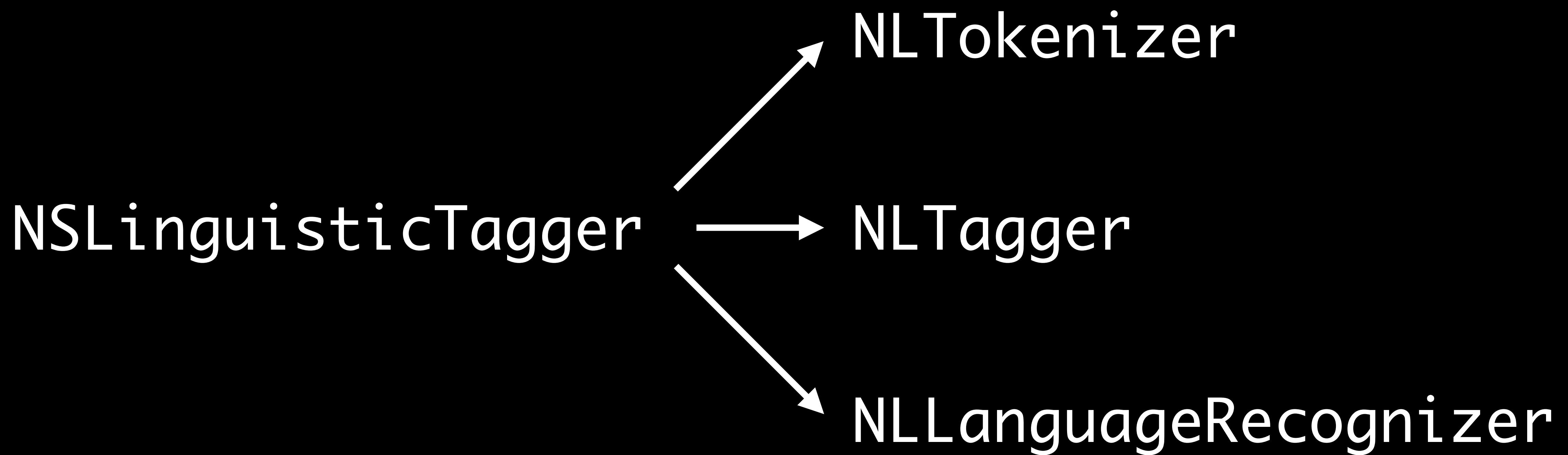
97

BD



**Natural
Language
Text**

Natural Language Framework



Tokenization

All human beings are born free and equal in dignity and rights.
They are endowed with reason and conscience and should act
towards one another in a spirit of brotherhood.

Tokenization

All human beings are born free and equal in dignity and rights.
They are endowed with reason and conscience and should act
towards one another in a spirit of brotherhood.

Tokenization

All human beings are born free and equal in dignity and rights.

They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.

Tokenization

All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.

Tokenization

All **human** beings are born free and equal in dignity and rights.
They are endowed with reason and conscience and should act
towards one another in a spirit of brotherhood.

Tokenization

All human **beings** are born free and equal in dignity and rights.
They are endowed with reason and conscience and should act
towards one another in a spirit of brotherhood.

```
let string = """  
    All human beings are born free and equal in dignity and rights.  
"""
```

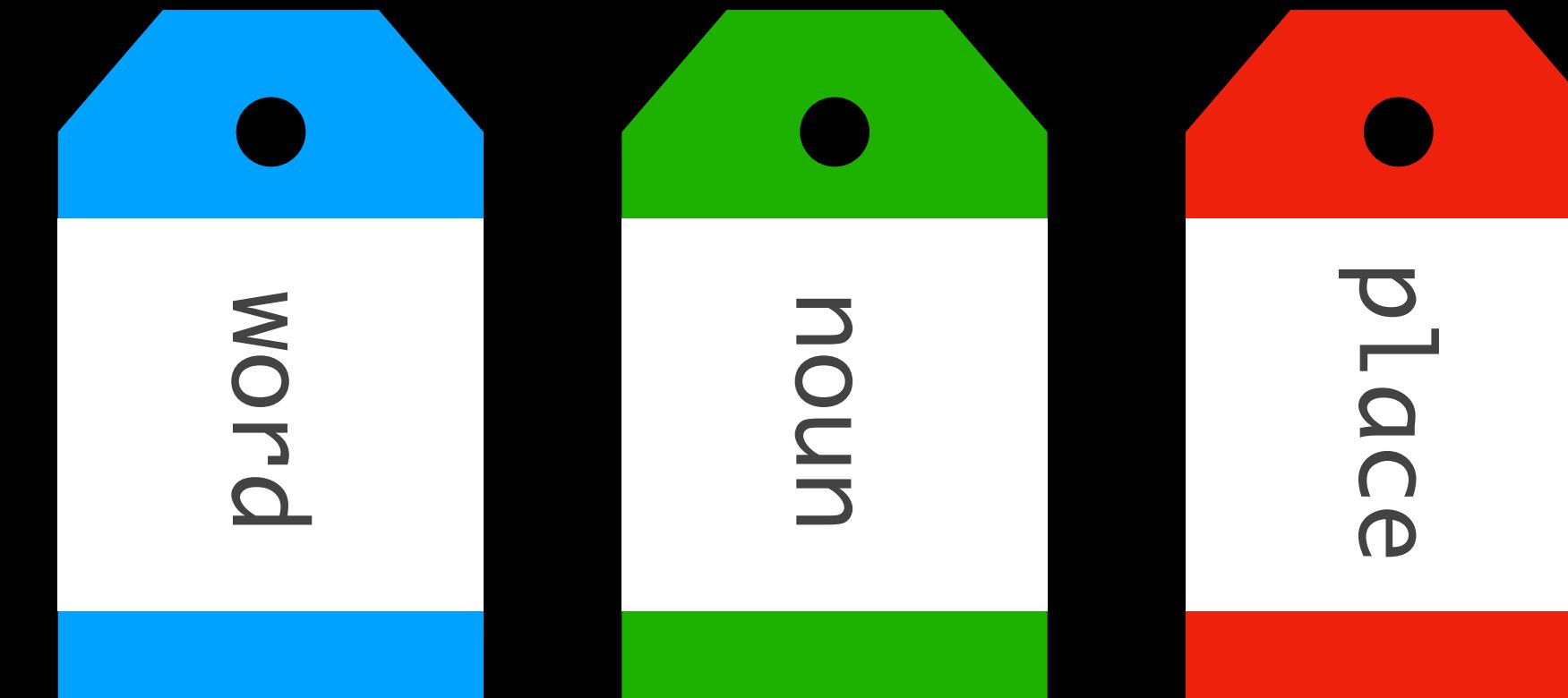
```
let paragraphs = string.split(separator: "\n")  
let sentences = paragraphs.flatMap{ $0.split(separator: ".") }  
let words = sentences.flatMap{ $0.split(separator: " ") }
```

```
let string = """  
    人人生而自由，在尊严和权利上一律平等。  
    ""
```

```
let tokenizer = NLTokenizer(unit: .word)
tokenizer.string = string
```

```
let stringRange = string.startIndex..
```

Tagging



Part of Speech Tagging

```
let string = "The quick brown fox jumps over the lazy dog."
```

```
let tagger = NLTagger(tagSchemes: [.lexicalClass])
tagger.string = string
```

```
let options: NLTagger.Options = [.omitWhitespace, .omitPunctuation]
tagger.enumerateTags(in: string.startIndex..
```

Part of Speech Tagging

The	quick	brown	fox	jumped	over	the	lazy	dog
Determiner	Adjective	Adjective	Noun	Verb	Preposition	Determiner	Adjective	Noun

Named Entity Recognition

```
let string = "Tim Cook is the CEO of Apple Inc, based in Cupertino, California."  
  
let tagger = NLTagger(tagSchemes: [.nameType])  
tagger.string = string  
  
let options: NLTagger.Options = [.omitPunctuation, .omitWhitespace, .joinNames]  
let tags: Set<NLTag> = [.personalName, .placeName, .organizationName]  
  
// ...
```

Named Entity Recognition

Tim Cook	PersonalName
Apple Inc	OrganizationName
Cupertino	PlaceName
California	PlaceName

Custom Tagging Model

```
import CreateML

let trainingData: [(tokens: [MLWordTagger.Token], labels: [String])] = [
    (
        ["Check", "out", "my", "new", "iPhone!"],
        [ "", "", "", "", "apple" ]
    ),
    (
        ["How", "do", "you", "like", "your", "iPad?"],
        [ "", "", "", "", "", "apple" ]
    )
]

let wordTagger = try MLWordTagger(trainingData: trainingData)
```

Custom Tagging Model

```
import NaturalLanguage

let customModel = try NLModel(mlModel: wordTagger.model)

let appleProductScheme = NLTagScheme(rawValue: "appleProduct")

let string = "Try out the app on your iPhone or iPad."

let tagger = NLTagger(tagSchemes: [appleProductScheme])
tagger.setModels([customModel], forTagScheme: appleProductScheme)
tagger.string = string
```

iPhone



iPad



Tråñslítératiôñ

Transliteration


```
import Foundation
```

```
let string = "모든 인간은 태어날 때부터 자유로우며 그 존엄과 권리에 있어 동등하다."
```

```
string.applyingTransform(.toLatin, reverse: false)  
/* "modeun ingan-eun taeeonal ttaebuteo jayuloumyeo  
geu jon-eomgwa gwonlie iss-eo dongdeunghada." */
```



Search

5

Structured Data



```
import Foundation

let string = """
    try! Swift NYC starts at 9 AM on September 5th.
    The venue is at 340 W 50th St.
    For more information, see https://tryswift.co/
"""

let detector =
    try NSDataDetector(types: NSTextCheckingAllTypes)
```

```
detector.enumerateMatches(in: string,
                           options: [],
                           range: range) { (match, flags, _) in
    guard let match = match else {
        return
    }

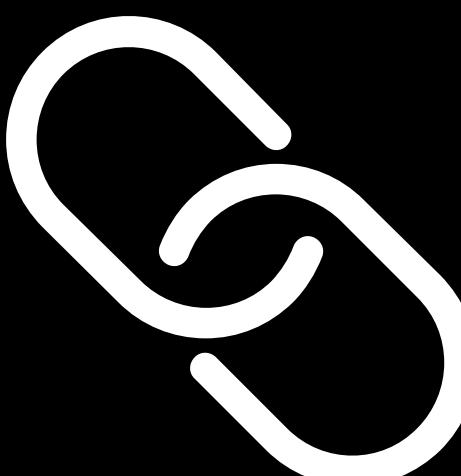
    switch match.resultType {
        case .date:
            let date = match.date
            print(date)
        case .address:
            let street = match.components? [.street]
            print(street)
        case .link:
            let url = match.url
            print(url)
        default:
            return
    }
}
```



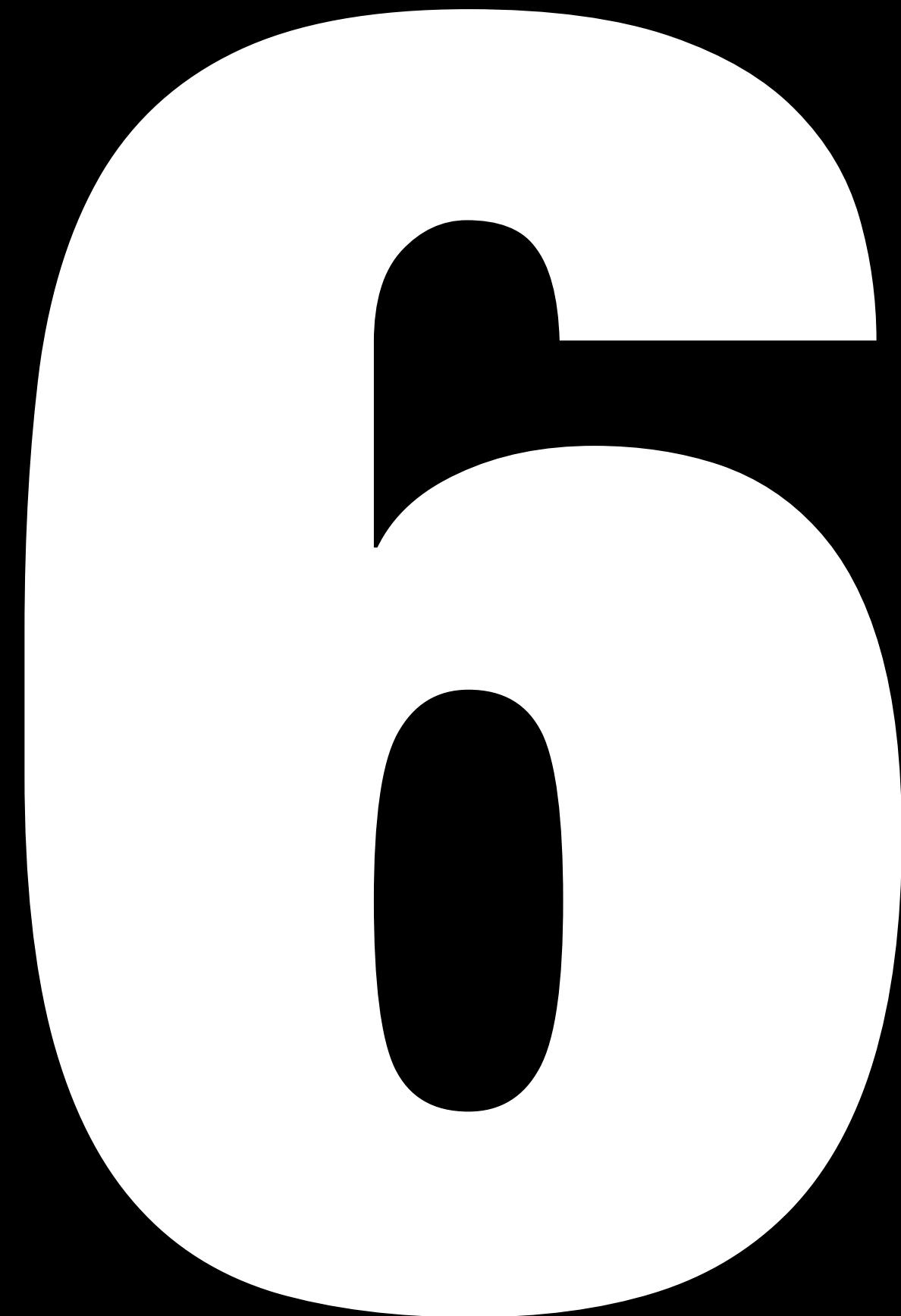
2018-09-05 13:00:00 +0000



340 W 50th St.



<https://tryswift.co/>



Encoded
Data

	0	1	2	3	4	5	6	7
0	A	B	C	D	E	F	G	H
1	I	J	K	L	M	N	O	P
2	Q	R	S	T	U	V	W	X
3	Y	Z	a	b	c	d	e	f
4	g	h	i	j	k	l	m	n
5	o	p	q	r	s	t	u	v
6	w	x	y	z	0	1	2	3
7	4	5	6	7	8	9	+	/

Base64

**All human beings are born free
and equal in dignity and rights.**

QWxsIGH1bWFuIGJlaW5ncy
BhcmUgYm9ybiBmcmVlIGFu
ZCBlcXVhbCBpbIBkaWduaX
R5IGFuZCByaWdodHMu



To: to@swift.org
From: from@swift.org
Subject: MIME Email Example
MIME-Version: 1.0
Content-Type: multipart/alternative;boundary = "+++"
Message-ID: <20180901000000.123456789@iMac.local>
Date: Wed, 5 Sep 2018 16:10:00 -0500 (EST)

Content-Type: text/txt; charset=utf-8
Content-Transfer-Encoding: base64

QWxsIGh1bWFuIGJlaW5ncy+
BhcmUgYm9ybiBmcmVlIGFu+
ZCBlcXVhbCBpbkBkaWduaX+
R5IGFuZCByaWdodHMu



```
@font-face {  
    font-family: 'CustomTypeface';  
    src: url(data:font/woff2;  
              charset=utf-8;  
              base64,...==)  
         format('woff2'),  
    font-weight: normal;  
    font-style: normal;  
}
```



```
import Foundation
import UIKit

let base64EncodedPNG = """
iVBORw0KGgoAAAANSUhEUgAAAAEAAAAB
CAAAAAAA6fptVAAAACKlEQVR4nGP6DwAB
BQECz6AuzQAAAABJRU5ErkJggg==
"""

let data = Data(base64Encoded: base64EncodedPNG) !
let image = UIImage(data: data)
```



Typographic Information

Characters

≠

Glyphs

A LATIN CAPITAL LETTER A
(U+0041)

A

LATIN CAPITAL LETTER A
(U+0041)

A LATIN CAPITAL LETTER A
(U+0041)

A

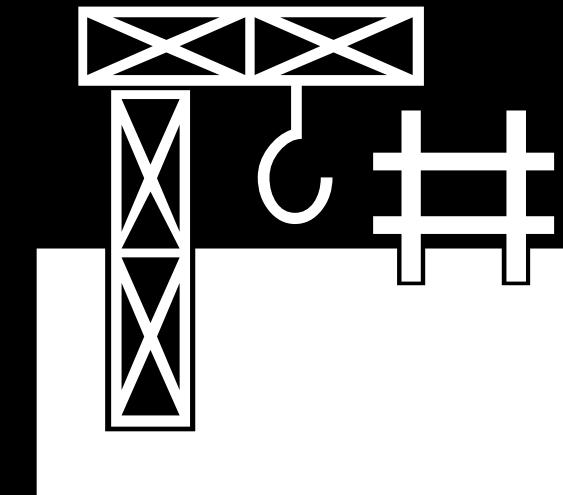
LATIN CAPITAL LETTER A

(U+0041)

A

A

A



Georgia

Marker Felt

Zapfino

Webdings

ع

Initial

ع

Medial

ع

Final

ع

Isolated

ffí

ffí

U

+

S

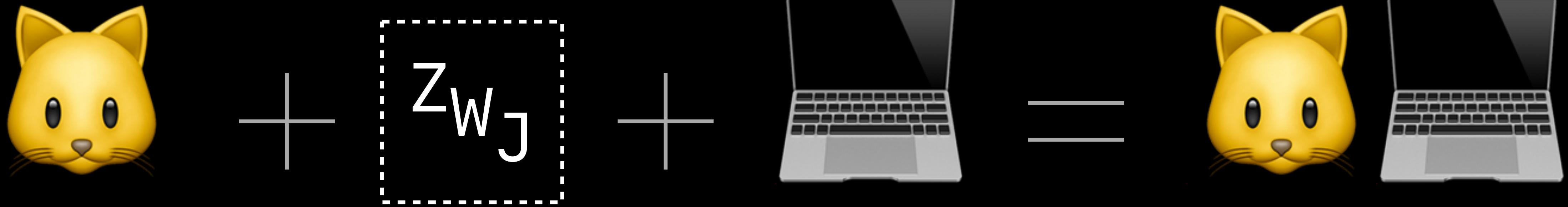
=



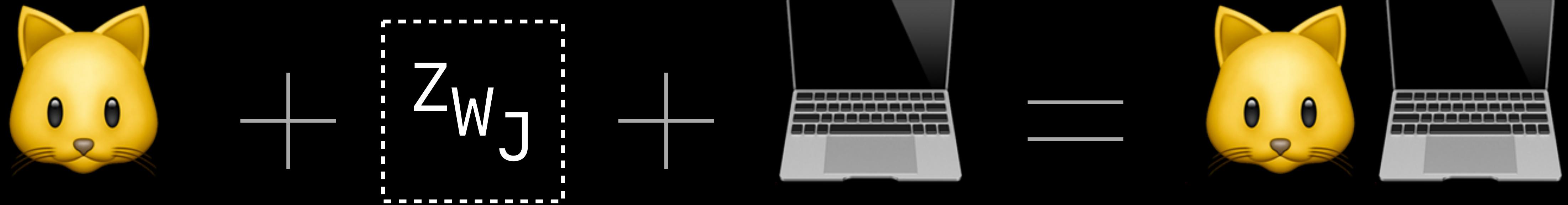
$$\begin{array}{c} \text{White flag} \\ + \\ \boxed{z_{W_J}} \\ + \\ \text{Rainbow} \\ = \\ \text{Rainbow flag} \end{array}$$

$$\begin{array}{c} \text{Black flag} \\ + \\ \boxed{z_{W_J}} \\ + \\ \text{Skull} \\ = \\ \text{Skull flag} \end{array}$$

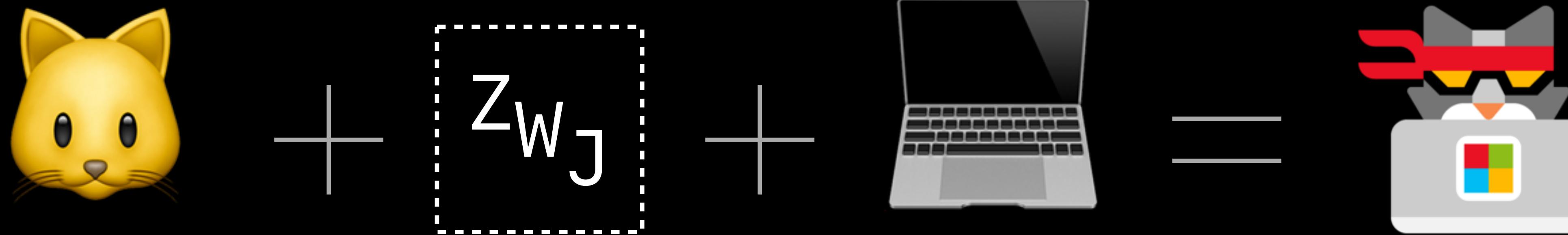
Apple Platforms



Apple Platforms



Microsoft Platforms





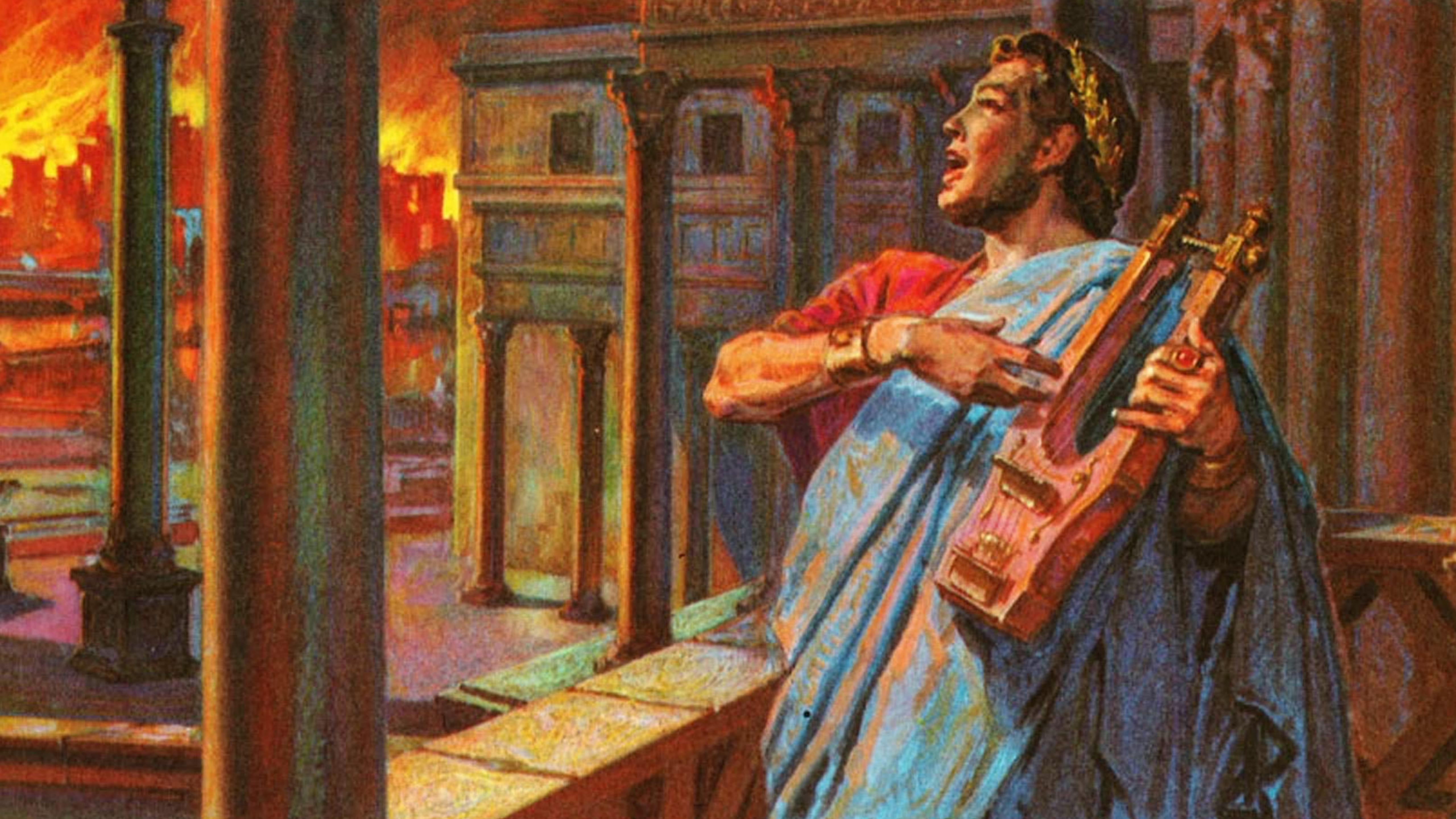


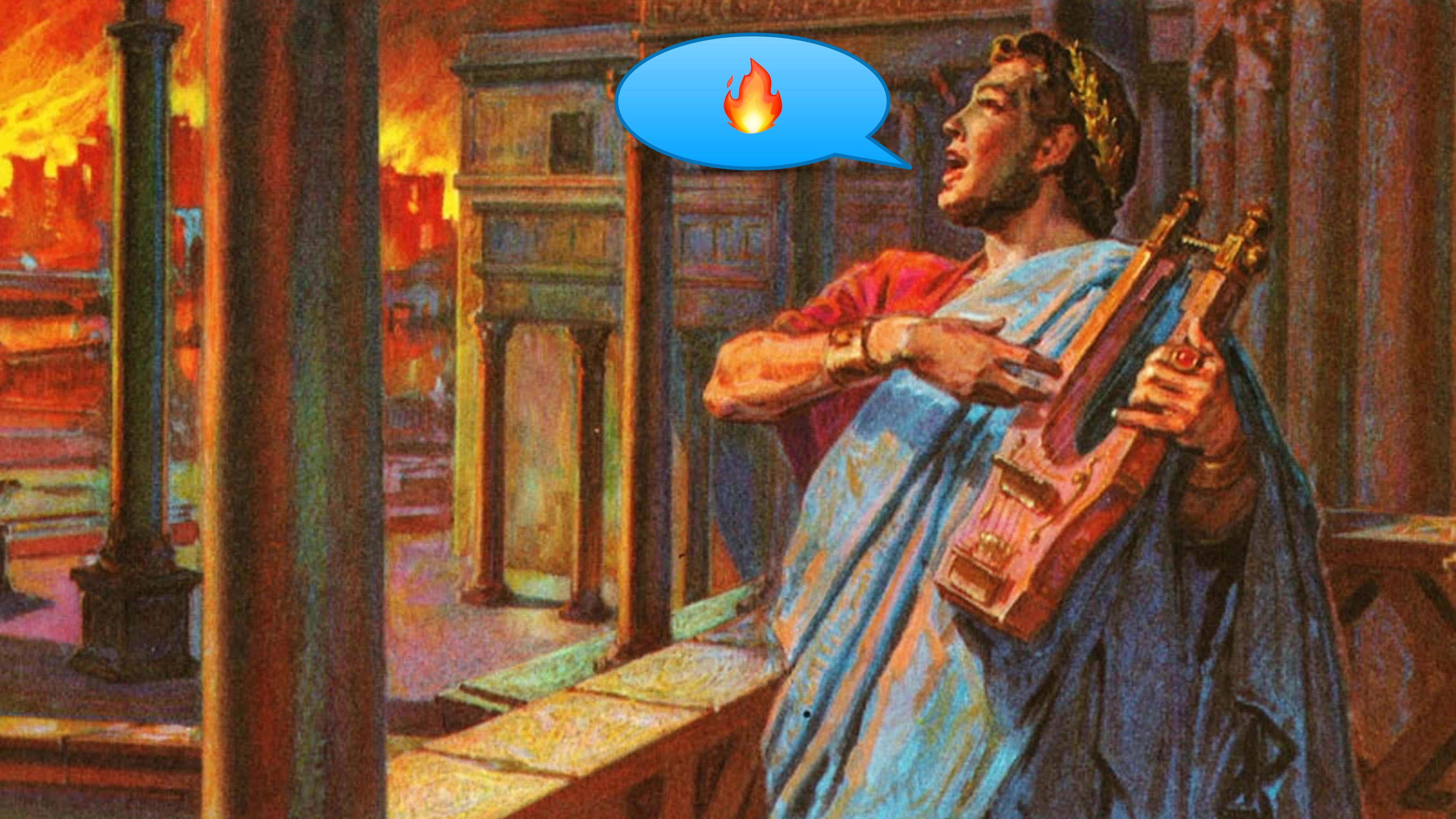


*

A A







Thanks&
Thanks&
Thanks&
Thanks&
Thanks&
Thanks&
Thanks&
Thanks&
Thanks

