



Unicode Search of Dirty Data, Or: How I Learned to Stop Worrying and Love Unicode Technical Standard #18

By

Jon Stewart and Joel Uckelman

Presented At

The Digital Forensic Research Conference

DFRWS 2013 USA Monterey, CA (Aug 4th - 7th)

DFRWS is dedicated to the sharing of knowledge and ideas about digital forensics research. Ever since it organized the first open workshop devoted to digital forensics in 2001, DFRWS continues to bring academics and practitioners together in an informal environment. As a non-profit, volunteer organization, DFRWS sponsors technical working groups, annual conferences and challenges to help drive the direction of research and development.

<http://dfrws.org>

Unicode Search of Dirty Data, Or: How I Learned to Stop Worrying and Love Unicode Technical Standard #18

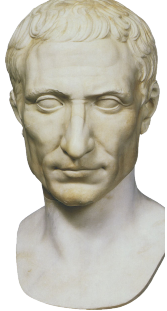
Jon Stewart Joel Uckelman*

Lightbox Technologies
Arlington, VA
`{jon,joel}@lightboxtechnologies.com`

DFRWS 2013

A Motivating Example

A Motivating Example



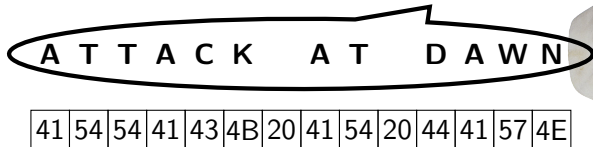
A Motivating Example



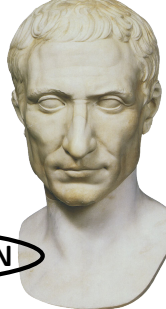
A Motivating Example



A Motivating Example



A Motivating Example

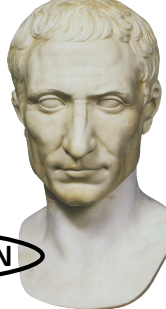


A T T A C K A T D A W N

latin1

41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

A Motivating Example



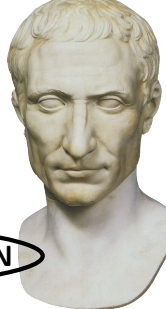
A T T A C K A T D A W N

latin1

41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

D3	63	63	D3	51	9D	4C	D3	63	4C	45	D3	4B	09
----	----	----	----	----	----	----	----	----	----	----	----	----	----

A Motivating Example



A T T A C K A T D A W N

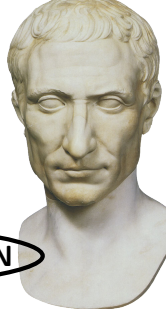
latin1

41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

OCE(latin1)

D3	63	63	D3	51	9D	4C	D3	63	4C	45	D3	4B	09
----	----	----	----	----	----	----	----	----	----	----	----	----	----

A Motivating Example



A T T A C K A T D A W N

latin1

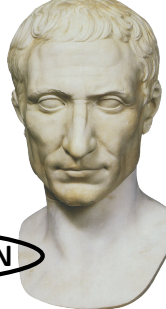
41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

OCE(latin1)

D3	63	63	D3	51	9D	4C	D3	63	4C	45	D3	4B	09
----	----	----	----	----	----	----	----	----	----	----	----	----	----

41	00	54	00	54	00	41	00	43	00	4B	00	20	00	41	00	54	00	20	00	44	00	41	00	57	00	4E	00
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

A Motivating Example



A T T A C K A T D A W N

latin1

41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

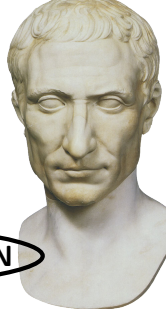
OCE(latin1)

D3	63	63	D3	51	9D	4C	D3	63	4C	45	D3	4B	09
----	----	----	----	----	----	----	----	----	----	----	----	----	----

41	00	54	00	54	00	41	00	43	00	4B	00	20	00	41	00	54	00	20	00	44	00	41	00	57	00	4E	00
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

UTF-16LE

A Motivating Example



A T T A C K A T D A W N

latin1

41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

OCE(latin1)

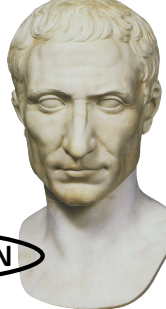
D3	63	63	D3	51	9D	4C	D3	63	4C	45	D3	4B	09
----	----	----	----	----	----	----	----	----	----	----	----	----	----

41	00	54	00	54	00	41	00	43	00	4B	00	20	00	41	00	54	00	20	00	44	00	41	00	57	00	4E	00
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

UTF-16LE

00	41	00	54	00	54	00	41	00	43	00	4B	00	20	00	41	00	54	00	20	00	44	00	41	00	57	00	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

A Motivating Example



A T T A C K A T D A W N

latin1

41	54	54	41	43	4B	20	41	54	20	44	41	57	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----

OCE(latin1)

D3	63	63	D3	51	9D	4C	D3	63	4C	45	D3	4B	09
----	----	----	----	----	----	----	----	----	----	----	----	----	----

41	00	54	00	54	00	41	00	43	00	4B	00	20	00	41	00	54	00	20	00	44	00	41	00	57	00	4E	00
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

UTF-16LE

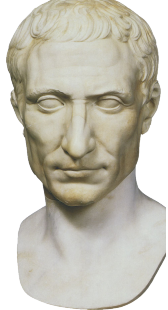
00	41	00	54	00	54	00	41	00	43	00	4B	00	20	00	41	00	54	00	20	00	44	00	41	00	57	00	4E
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

UTF-16BE

A Motivating Example

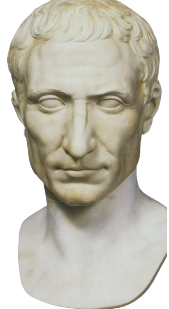


A Motivating Example



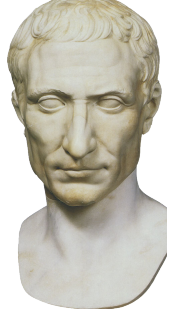
- Many possible encodings for text, all different

A Motivating Example



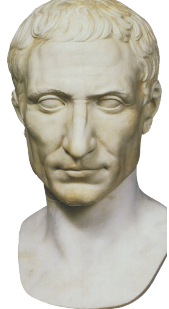
- Many possible encodings for text, all different
- We don't want to search serially for each encoding

A Motivating Example



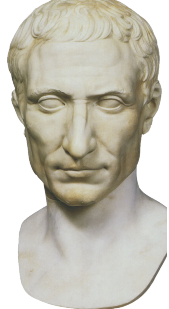
- Many possible encodings for text, all different
- We don't want to search serially for each encoding
- We'd like to have encoding-independent patterns

A Motivating Example



- Many possible encodings for text, all different
- We don't want to search serially for each encoding
- We'd like to have encoding-independent patterns
- We have a multipattern search tool (Lightgrep)

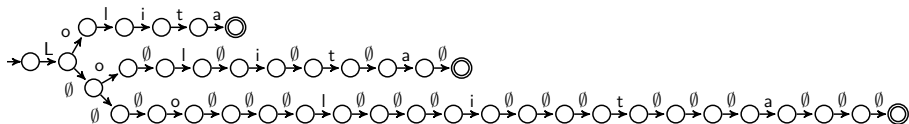
A Motivating Example



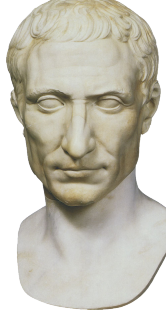
- Many possible encodings for text, all different
- We don't want to search serially for each encoding
- We'd like to have encoding-independent patterns
- We have a multipattern search tool (Lightgrep)

What to do?

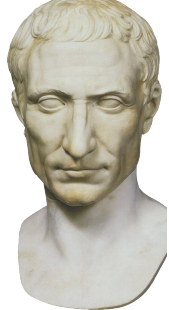
Multiencoding search is multipattern search.



A Motivating Example



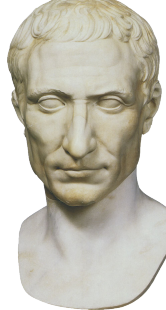
A Motivating Example



A T T A C K A T V I I I

A Motivating Example

A	T	T	A	C	K	A	T	I	V
A	T	T	A	C	K	A	T	V	
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	I	X
A	T	T	A	C	K	A	T	X	
A	T	T	A	C	K	A	T	X	I



A Motivating Example

A	T	T	A	C	K	A	T	I	V
A	T	T	A	C	K	A	T	V	
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	I	X
A	T	T	A	C	K	A	T	X	
A	T	T	A	C	K	A	T	X	I

ATTACK AT [IVX]+



A Motivating Example

A	T	T	A	C	K	A	T	I	V
A	T	T	A	C	K	A	T	V	
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	I	X
A	T	T	A	C	K	A	T	X	
A	T	T	A	C	K	A	T	X	I

ATTACK AT [IVX]+



A Motivating Example

A	T	T	A	C	K	A	T	I	V
A	T	T	A	C	K	A	T	V	
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	I	X
A	T	T	A	C	K	A	T	X	
A	T	T	A	C	K	A	T	X	I

ATTACK AT [IVX]+



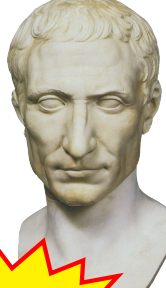
VIII ROMAN NUMERAL EIGHT (U+2167)



A Motivating Example

A	T	T	A	C	K	A	T	I	V
A	T	T	A	C	K	A	T	V	
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	V	I
A	T	T	A	C	K	A	T	I	X
A	T	T	A	C	K	A	T	X	
A	T	T	A	C	K	A	T	X	I

ATTACK AT [IVX]4



VIII ROMAN NUMERAL EIGHT (U+2167)



A Motivating Example

A T T A C K A T I V
A T T A C K A T V
A T T A C K A T V I
A T T A C K A T V I I
A T T A C K A T V I I I
A T T A C K A T X
A T T A C K A T X
A T T A C K A T X I

ATTACK AT [IVX]+

VIII ROMAN NUMERAL EIGHT (U+2167)



What is Unicode?



Before Unicode:

- Hundreds of character encodings

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

Unicode is an attempt to represent all characters used by humans.

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

Unicode is an attempt to represent all characters used by humans.

After Unicode:

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

Unicode is an attempt to represent all characters used by humans.

After Unicode:

- Only a handful of actively-used encodings

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

Unicode is an attempt to represent all characters used by humans.

After Unicode:

- Only a handful of actively-used encodings
- Character set \neq character encoding

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

Unicode is an attempt to represent all characters used by humans.

After Unicode:

- Only a handful of actively-used encodings
- Character set \neq character encoding
- Unicode is a character set.
UTF-8, UTF-16LE, etc. are character encodings for Unicode.

What is Unicode?



Before Unicode:

- Hundreds of character encodings
- Character set = character encoding
- Small, often single-byte, character sets (≤ 256 characters)
- Interoperability was messy
- Multilanguage documents were difficult

Unicode is an attempt to represent all characters used by humans.

After Unicode:

- Only a handful of actively-used encodings
- Character set \neq character encoding
- Unicode is a character set.
UTF-8, UTF-16LE, etc. are character encodings for Unicode.
- Huge character set—space for 1.1m *code points*

Code point A (possibly named) character, identified by a hex number

Code point A (possibly named) character, identified by a hex number

E.g.: U+0058, LATIN CAPITAL LETTER X

Code point A (possibly named) character, identified by a hex number

E.g.: U+0058, LATIN CAPITAL LETTER X

Code points range from U+0000 to U+10FFFF

- Unicode contains many weird and wonderful things.

ð ୯ Д € 獮 Φ ﷺ Ɔ :: fff 🍷 📞

- Unicode contains many weird and wonderful things.

ø ॢ Д € 獮 Φ ﷺ Ɔ :: fff 🍷 📞

- Regular expressions were designed for ASCII text.

- Unicode contains many weird and wonderful things.

ð ୯ Д € 獮 Φ ﷺ Ɔ :: fff 🍷 📞

- Regular expressions were designed for ASCII text.
- Unicode Technical Standard #18 gives guidelines for supporting Unicode in regular expressions:

- Unicode contains many weird and wonderful things.

ð ୯ Д € 獮 Φ ﷺ Ɔ :: fff 🍷 📞

- Regular expressions were designed for ASCII text.
- Unicode Technical Standard #18 gives guidelines for supporting Unicode in regular expressions:
 - ▶ Level 1: Basic Unicode Support

- Unicode contains many weird and wonderful things.

ø ლ Д € 獮 Φ ﷺ ூ :: fff 🍷 📞

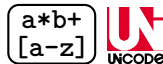
- Regular expressions were designed for ASCII text.
- Unicode Technical Standard #18 gives guidelines for supporting Unicode in regular expressions:
 - ▶ Level 1: Basic Unicode Support
 - ▶ Level 2: Extended Unicode Support

- Unicode contains many weird and wonderful things.

ø ॢ Д € 獮 Φ ﷺ Ɔ :: fff 🍷 📞

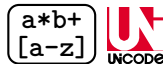
- Regular expressions were designed for ASCII text.
- Unicode Technical Standard #18 gives guidelines for supporting Unicode in regular expressions:
 - ▶ Level 1: Basic Unicode Support
 - ▶ Level 2: Extended Unicode Support
 - ▶ Level 3: Tailored Support

What is that? How do I type it?



- Some characters are indistinguishable from one another.

What is that? How do I type it?

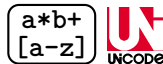


- Some characters are indistinguishable from one another.

K LATIN CAPITAL LETTER K (U+004B)

K KELVIN SIGN (U+212A)

What is that? How do I type it?



- Some characters are indistinguishable from one another.

K LATIN CAPITAL LETTER K (U+004B)

K KELVIN SIGN (U+212A)

tab (U+0009)

SPACE (U+0020)

MONGOLIAN VOWEL SEPARATOR (U+180E) FOUR-PER-EM SPACE (U+2005)

EN QUAD (U+2000)

SIX-PER-EM SPACE (U+2006)

EM QUAD (U+2001)

FIGURE SPACE (U+2007)

EN SPACE (U+2002)

PUNCTUATION SPACE (U+2008)

EM SPACE (U+2003)

THIN SPACE (U+2009)

THREE-PER-EM SPACE (U+2004)

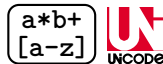
HAIR SPACE (U+200A)

NARROW NO-BREAK SPACE (U+202F)

MEDIUM MATHEMATICAL SPACE (U+205F)

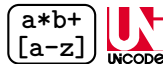
IDEOGRAPHIC SPACE (U+3000)

What is that? How do I type it?



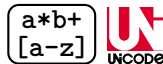
- Some characters are indistinguishable from one another.
- Some characters will be unfamiliar to you.

What is that? How do I type it?



- Some characters are indistinguishable from one another.
- Some characters will be unfamiliar to you.
- Some characters you won't know how to type.

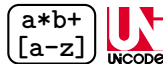
What is that? How do I type it?



- Some characters are indistinguishable from one another.
- Some characters will be unfamiliar to you.
- Some characters you won't know how to type.

UTS #18 requires:

What is that? How do I type it?



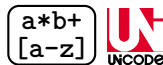
- Some characters are indistinguishable from one another.
- Some characters will be unfamiliar to you.
- Some characters you won't know how to type.

UTS #18 requires:

- Code points can be specified by hex notation (RL1.2)

$\backslash x\{10DA\} = \text{𐍚}$

What is that? How do I type it?



- Some characters are indistinguishable from one another.
- Some characters will be unfamiliar to you.
- Some characters you won't know how to type.

UTS #18 requires:

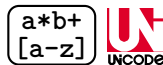
- Code points can be specified by hex notation (RL1.2)

`\x{10DA}` = ☞

- Code points can be specified by name (RL2.5)

`\N{latin small letter sharp s}` = ß

What is that? How do I type it?



- Some characters are indistinguishable from one another.
- Some characters will be unfamiliar to you.
- Some characters you won't know how to type.

UTS #18 requires:

- Code points can be specified by hex notation (RL1.2)

`\x{10DA}` = ☞

- Code points can be specified by name (RL2.5)

`\N{latin small letter sharp s}` = ß

- All code points are supported (RL1.7)

- How many characters match `[a-z]` case insensitively? 52

- How many characters match `[a-z]` case insensitively?

~~52~~

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!
 - ▶ `a-z` and `A-Z`, as expected but also

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!
 - ▶ `a-z` and `A-Z`, as expected but also
 - ▶ `K` KELVIN SIGN (U+212A), and

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!
 - ▶ a–z and A–Z, as expected but also
 - ▶ K KELVIN SIGN (U+212A), and
 - ▶ ſ LATIN SMALL LETTER LONG S (U+017F)

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!
 - ▶ `a-z` and `A-Z`, as expected but also
 - ▶ `K` KELVIN SIGN (U+212A), and
 - ▶ `ſ` LATIN SMALL LETTER LONG S (U+017F)
- UTS #18 requires proper handling of case-insensitivity (RL1.5)

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!
 - ▶ `a-z` and `A-Z`, as expected but also
 - ▶ `℄` KELVIN SIGN (U+212A), and
 - ▶ `ſ` LATIN SMALL LETTER LONG S (U+017F)
- UTS #18 requires proper handling of case-insensitivity (RL1.5)
- E.g., `happiness` case-insensitively matches the last word in “the pursuit of Happiness”

- How many characters match `[a-z]` case insensitively? ~~52~~ 54 !!!
- What?!
 - ▶ `a-z` and `A-Z`, as expected but also
 - ▶ `℄` KELVIN SIGN (U+212A), and
 - ▶ `ſ` LATIN SMALL LETTER LONG S (U+017F)
- UTS #18 requires proper handling of case-insensitivity (RL1.5)
- E.g., `happiness` case-insensitively matches the last word in “the pursuit of Happiness”
- No missed hits when case-insensitively searching unfamiliar scripts

- *Properties* are named sets of code points. E.g.:

Alphabetic	Letter	Mark	Separator
Uppercase	Number	Symbol	Assigned
Lowercase	Punctuation	Control	...

- *Properties* are named sets of code points. E.g.:

Alphabetic	Letter	Mark	Separator
Uppercase	Number	Symbol	Assigned
Lowercase	Punctuation	Control	...

- *Scripts* are properties:

Latin	Arabic	Hangul	Hebrew
Greek	Han	Katakana	Devanagari
Cyrillic	Thai	Hiragana	...

- *Properties* are named sets of code points. E.g.:

Alphabetic	Letter	Mark	Separator
Uppercase	Number	Symbol	Assigned
Lowercase	Punctuation	Control	...

- *Scripts* are properties:

Latin	Arabic	Hangul	Hebrew
Greek	Han	Katakana	Devanagari
Cyrillic	Thai	Hiragana	...

- UTS #18: properties accessible as character classes (RL1.2, RL2.7)

- *Properties* are named sets of code points. E.g.:

Alphabetic	Letter	Mark	Separator
Uppercase	Number	Symbol	Assigned
Lowercase	Punctuation	Control	...

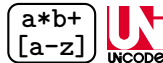
- *Scripts* are properties:

Latin	Arabic	Hangul	Hebrew
Greek	Han	Katakana	Devanagari
Cyrillic	Thai	Hiragana	...

- UTS #18: properties accessible as character classes (RL1.2, RL2.7)

```
\p{Cyrillic}{4,}  
[\p{Arabic}\p{Whitespace}\p{Punctuation}]+
```


Set Operations on Character Classes



- How do I match *lowercase* Greek?

- How do I match *lowercase* Greek?

`[α-ω]` misses `ï`, `ü`, `ó`, `ú`, `ώ`

- How do I match *lowercase* Greek?

`[α-ω]` misses `ï`, `ü`, `ó`, `ú`, `ώ`

- UTS #18 requires union, intersection, subtraction for CCs (RL1.3)

- How do I match *lowercase* Greek?

`[α-ω]` misses `ï`, `ü`, `ó`, `ú`, `ώ`

- UTS #18 requires union, intersection, subtraction for CCs (RL1.3)

`[\p{Greek}&&\p{Lowercase}]`

- How do I match *lowercase* Greek?

`[α-ω]` misses `ï, ü, ó, ú, ó`

- UTS #18 requires union, intersection, subtraction for CCs (RL1.3)

`[\p{Greek}&&\p{Lowercase}]`

`[\p{ASCII}--\p{Whitespace}]`

	lightgrep	ICU 50	Perl 5.16	Java 7	Python regex
Level 1: Basic Unicode Support	○	●	○	●	●
RL1.1 Hex Notation	●	●	●	●	●
RL1.2 Properties	●	●	●	●	●
RL1.3 Subtraction and Intersection	●	●	○	●	●
RL1.4 Simple Word Boundaries		●	●	●	●
RL1.5 Simple Loose Matches	●	●	●	●	●
RL1.6 Line Boundaries		●	○	●	●
RL1.7 Supplementary Code Points	●	●	●	●	●
Level 2: Extended Unicode Support	○	○	○	○	○
RL2.1 Canonical Equivalents				●	
RL2.2 Extended Grapheme Clusters		○	○		●
RL2.3 Default Word Boundaries		●			
RL2.4 Default Case Conversion					
RL2.5 Name Properties	●	●	●		●
RL2.6 Wildcards in Property Values					
RL2.7 Full Properties	●		●		●
Level 3: Tailored Support	○		○		
RL3.1 Tailored Punctuation					
RL3.2 Tailored Grapheme Clusters					
RL3.6 Context Matching			○		
RL3.7 Incremental Matches	●				
RL3.9 Possible Match Sets	●				
RL3.11 Submatchers					

● = full support, ○ = partial support

	lightgrep	ICU 50	Perl 5.16	Java 7	Python regex
Level 1: Basic Unicode Support	○	●	○	●	●
RL1.1 Hex Notation	●	●	●	●	●
RL1.2 Properties	●	●	●	●	●
RL1.3 Subtraction and Intersection	●	●	○	●	●
RL1.4 Simple Word Boundaries	●	●	●	●	●
RL1.5 Simple Loose Matches	●	●	●	●	●
RL1.6 Line Boundaries	●	●	○	●	●
RL1.7 Supplementary Code Points	●	●	●	●	●
Level 2: Extended Unicode Support	○	○	○	○	○
RL2.1 Canonical Equivalents				●	
RL2.2 Extended Grapheme Clusters		○	○		●
RL2.3 Default Word Boundaries		●			
RL2.4 Default Case Conversion					
RL2.5 Name Properties	●	●	●		●
RL2.6 Wildcards in Property Values					
RL2.7 Full Properties	●		●		●
Level 3: Tailored Support	○		○		
RL3.1 Tailored Punctuation					
RL3.2 Tailored Grapheme Clusters					
RL3.6 Context Matching			○		
RL3.7 Incremental Matches	●				
RL3.9 Possible Match Sets	●				
RL3.11 Submatchers					

● = full support, ○ = partial support

Encoding Chains

- Sometimes data is transformed, not just (character) encoded:

Encoding Chains

- Sometimes data is transformed, not just (character) encoded:

ROT-13

OCE

any substitution cipher

URL-encoding

XOR(n)

...

Encoding Chains

- Sometimes data is transformed, not just (character) encoded:

ROT-13

OCE

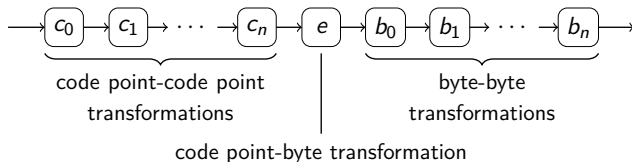
any substitution cipher

URL-encoding

XOR(n)

...

- We can generalize:



Encoding Chains

- Sometimes data is transformed, not just (character) encoded:

ROT-13

OCE

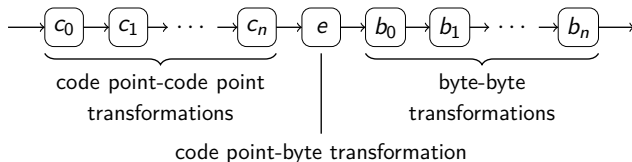
any substitution cipher

URL-encoding

XOR(n)

...

- We can generalize:



- We support arbitrarily-long encoding chains:

Encoding Chains

- Sometimes data is transformed, not just (character) encoded:

ROT-13

OCE

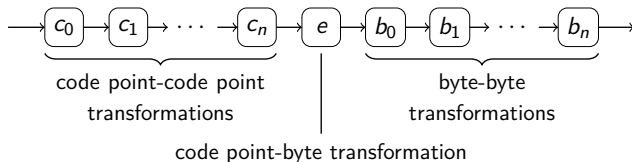
any substitution cipher

URL-encoding

XOR(n)

...

- We can generalize:



- We support arbitrarily-long encoding chains:

UTF-16LE | OCE

UTF-8 | XOR(1F)

ROT-13 | ASCII | OCE | XOR(D7)

Encoding Chains

- Sometimes data is transformed, not just (character) encoded:

ROT-13

OCE

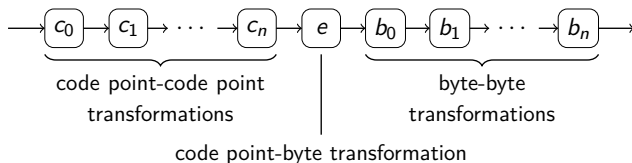
any substitution cipher

URL-encoding

XOR(n)

...

- We can generalize:



- We support arbitrarily-long encoding chains:

UTF-16LE | OCE

UTF-8 | XOR(1F)

ROT-13 | ASCII | OCE | XOR(D7)

- All the work happens at pattern-compile time.

Hit Context

- Seeing *decoded* context around search hits is handy

Hit Context

- Seeing *decoded* context around search hits is handy

A T T A C K A T D A W N

Hit Context

- Seeing *decoded* context around search hits is handy

A T T A C K A T D A W N

- But context could be *mojibake*!

Hit Context

- Seeing *decoded* context around search hits is handy

A T T A C K A T D A W N

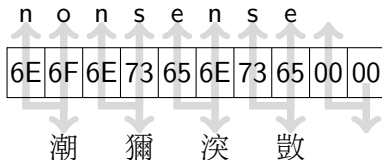
- But context could be *mojibake*!
- Is it nonsense or is it 潮獮湊數?

Hit Context

- Seeing *decoded* context around search hits is handy

A T T A C K **A T** D A W N

- But context could be *mojibake*!
- Is it nonsense or is it 潮獮湊數?

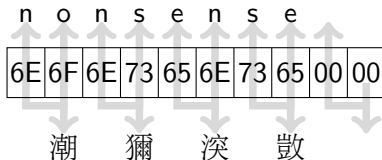


Hit Context

- Seeing *decoded* context around search hits is handy

A T T A C K **A T** D A W N

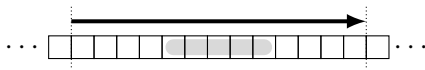
- But context could be *mojibake*!
- Is it nonsense or is it 潮獮湊數?



- User discretion is advised.

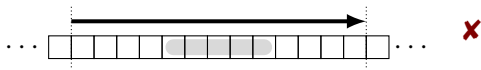
How We Decode Hit Context

Left to right



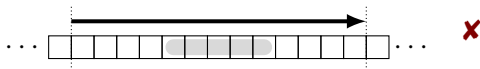
How We Decode Hit Context

Left to right

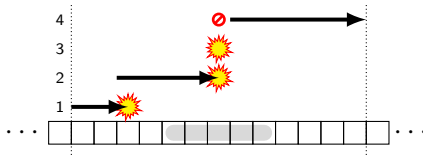


How We Decode Hit Context

Left to right

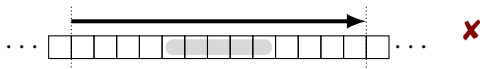


Left to right, restarting on error

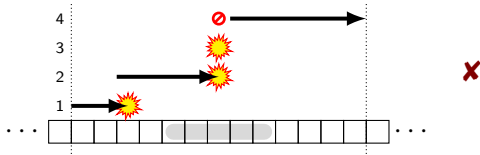


How We Decode Hit Context

Left to right

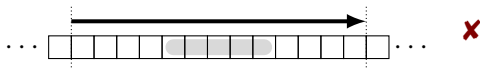


Left to right, restarting on error

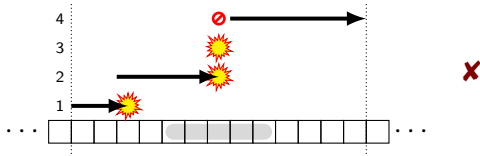


How We Decode Hit Context

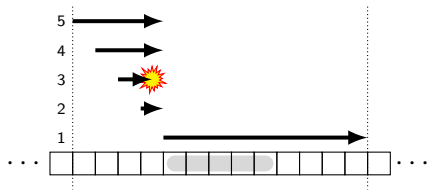
Left to right



Left to right, restarting on error

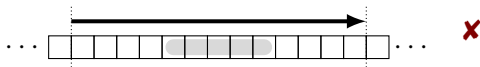


Leading context right to left, hit and trailing context left to right

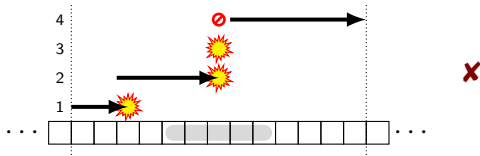


How We Decode Hit Context

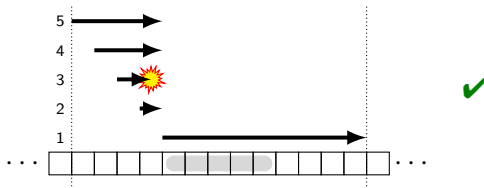
Left to right



Left to right, restarting on error

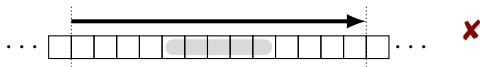


Leading context right to left, hit
and trailing context left to right

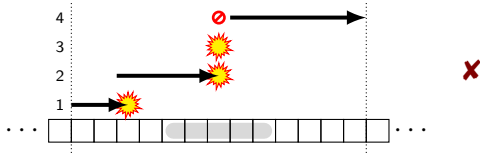


How We Decode Hit Context

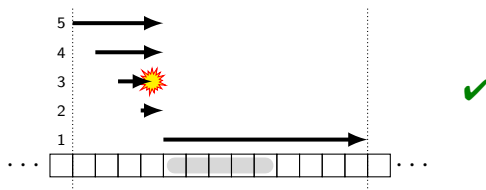
Left to right



Left to right, restarting on error



Leading context right to left, hit and trailing context left to right



$O(n^2)$ but can happen after (or in parallel with) search

How fast is search?

Image	Size	Hits	Bandwidth
Charlie-2009-11-12.E01	10GB	543,719	131.7 MB/s
Jo-2009-11-12.E01	14GB	536,906	123.8 MB/s
Pat-2009-11-12.E01	12GB	522,220	124.7 MB/s

Search performance using lightgrep via bulk_extractor
Images from the NPS 2009 M57 Patents-Redacted scenario

How fast is search?

Image	Size	Hits	Bandwidth
Charlie-2009-11-12.E01	10GB	543,719	131.7 MB/s
Jo-2009-11-12.E01	14GB	536,906	123.8 MB/s
Pat-2009-11-12.E01	12GB	522,220	124.7 MB/s

Search performance using lightgrep via bulk_extractor
Images from the NPS 2009 M57 Patents-Redacted scenario

Multiencoding/Unicode search can be fast.

Conclusions



- Unicode is complex because natural language is complex

Conclusions



- Unicode is complex because natural language is complex
- Encodings are messy for historical reasons

Conclusions



- Unicode is complex because natural language is complex
- Encodings are messy for historical reasons
- **GORY DETAILS** \rightsquigarrow Don't implement this stuff yourself

Conclusions



- Unicode is complex because natural language is complex
- Encodings are messy for historical reasons
- **GORY DETAILS** \rightsquigarrow Don't implement this stuff yourself

BUT:

Conclusions



- Unicode is complex because natural language is complex
- Encodings are messy for historical reasons
- **GORY DETAILS** \rightsquigarrow Don't implement this stuff yourself

BUT:

- You can *use*
 - multiencoding support
 - multilanguage support
- without knowing much of anything about the details**

Conclusions



- Unicode is complex because natural language is complex
- Encodings are messy for historical reasons
- **GORY DETAILS** \rightsquigarrow Don't implement this stuff yourself

BUT:

- You can *use*
 - multiencoding support
 - multilanguage supportwithout knowing much of anything about the details
- Encoding chains make more data directly searchable



- Unicode is complex because natural language is complex
- Encodings are messy for historical reasons
- **GORY DETAILS** \rightsquigarrow Don't implement this stuff yourself

BUT:

- **You can *use***
 - **multiencoding support**
 - **multilanguage support****without knowing much of anything about the details**
- Encoding chains make more data directly searchable
- Hit context is handy, not entirely trivial to produce

Conclusions



A problem to think about

- Latin alphabet is used by dozens of languages, but with accent marks

e è é ê ë ē ė ě ě ê ě

A problem to think about

- Latin alphabet is used by dozens of languages, but with accent marks

e è é ê ë ē ě ě ě ě ě

- Unicode has *confusable* characters

A problem to think about

- Latin alphabet is used by dozens of languages, but with accent marks

e è é ê ë ē ě ě ě ě ě ě

- Unicode has *confusable* characters

UNDUTCHABLES
THE RECRUITMENT AGENCY FOR INTERNATIONALS

A problem to think about

- Latin alphabet is used by dozens of languages, but with accent marks

e è é ê ë ē ě ě ě ě ě ě

- Unicode has *confusable* characters

UNDUTCHABLES
THE RECRUITMENT AGENCY FOR INTERNATIONALS

A problem to think about

- Latin alphabet is used by dozens of languages, but with accent marks

e è é ê ë ē ě è ę ě ê ò

- Unicode has *confusable* characters

UNDUTCHABLES
THE RECRUITMENT AGENCY FOR INTERNATIONALS

- Unicode has characters which *decompose* to character sequences

VIII \Rightarrow V·I·I·I IJ \Rightarrow I·J ﷺ \Rightarrow و س ل م ع ل ي ه ا ل ل ه ص ل ي

A problem to think about

- Latin alphabet is used by dozens of languages, but with accent marks

e è é ê ë ē ė ě ě ê ò

- Unicode has *confusable* characters



- Unicode has characters which *decompose* to character sequences

VIII \Rightarrow V·I·I·I IJ \Rightarrow I·J ﷺ \Rightarrow وسلم عليّه الله صلى

- Wouldn't it be neat to have a “looks-like” operator or mode?

$\backslash L\{e\}$ $\backslash L\{VIII\}$ $\backslash L\{UNDUTCHABLES\}$

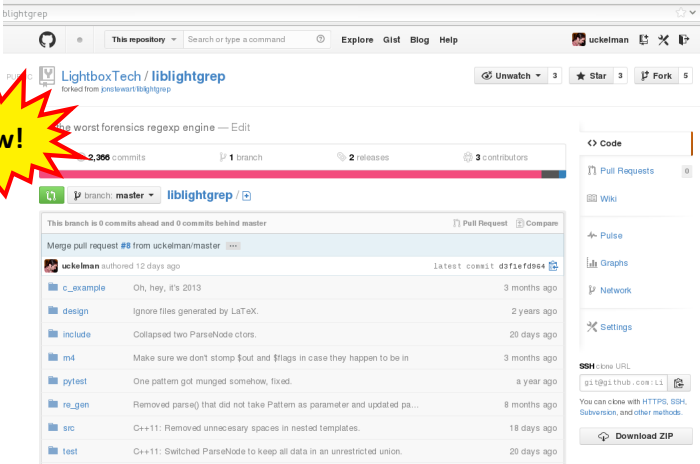
Where can I find this amazing software?

`https://github.com/LightboxTech/liblightgrep`

The screenshot shows the GitHub repository page for `liblightgrep` by `LightboxTech`. The repository is public and has 3 stars, 3 forks, and 5 watchers. It is forked from `jonstewart/liblightgrep`. The repository description is "not the worst forensics regexp engine — Edit". The repository has 2,366 commits, 1 branch, 2 releases, and 3 contributors. The current branch is `master`. The repository is 0 commits ahead and 0 commits behind master. A pull request #8 from `uckelman/master` is open. The latest commit is `d3f1efd964` by `uckelman` 12 days ago. The repository contains the following files and folders: `c_example` (Oh, hey, it's 2013, 3 months ago), `design` (Ignore files generated by LaTeX, 2 years ago), `include` (Collapsed two ParseNode ctors, 20 days ago), `m4` (Make sure we don't stomp \$out and \$flags in case they happen to be in, 3 months ago), `pytest` (One pattern got munged somehow, fixed, a year ago), `re_gen` (Removed parse() that did not take Pattern as parameter and updated pa..., 8 months ago), `src` (C++11: Removed unnecessary spaces in nested templates, 18 days ago), and `test` (C++11: Switched ParseNode to keep all data in an unrestricted union, 20 days ago). The right sidebar shows the "Code" tab selected, with options for Pull Requests, Wiki, Pulse, Graphs, Network, and Settings. The SSH clone URL is `git@github.com:LightboxTech/liblightgrep`. The repository can be cloned with HTTPS, SSH, Subversion, and other methods. A "Download ZIP" button is also present.

Where can I find this amazing software?

`https://github.com/LightboxTech/liblightgrep`



New!

LightboxTech / liblightgrep
forked from jstewart/liblightgrep

the worst forensics regexp engine — Edit

2,366 commits 1 branch 2 releases 3 contributors

branch: master liblightgrep

This branch is 0 commits ahead and 0 commits behind master

Merge pull request #8 from uckelman/master

uckelman authored 12 days ago latest commit d3f1efd964

c_example	Oh, hey, it's 2013	3 months ago
design	Ignore files generated by LaTeX.	2 years ago
include	Collapsed two ParseNode ctors.	20 days ago
m4	Make sure we don't stomp \$out and \$flags in case they happen to be in	3 months ago
pytest	One pattern got munged somehow, fixed.	a year ago
re_gen	Removed parse() that did not take Pattern as parameter and updated pa...	8 months ago
src	C++11: Removed unnecessary spaces in nested templates.	18 days ago
test	C++11: Switched ParseNode to keep all data in an unrestricted union.	20 days ago

SSH clone URL
git@github.com:LightboxTech/liblightgrep

You can clone with HTTPS, SSH, Subversion, and other methods.

Download ZIP

Where can I find this amazing software?

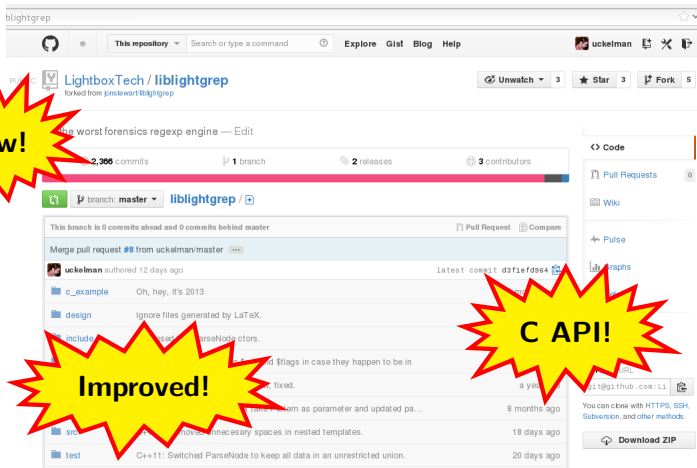
<https://github.com/LightboxTech/liblightgrep>

New!

Improved!

Where can I find this amazing software?

<https://github.com/LightboxTech/liblightgrep>



Where can I find this amazing software?

<https://github.com/LightboxTech/liblightgrep>

