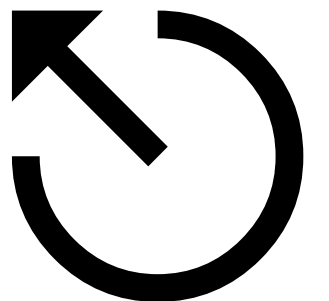


Do you want to print a ☺?

Mark Fowler, Albany Perl Mongers

Perl String Escapes Recap



é

```
# hex
```

```
$eacute = "\x{e9}";
```

```
# oct (<256 only)
```

```
$eacute = "\351";
```

```
# name
```

```
$eacute = "\N{LATIN SMALL LETTER E WITH ACUTE}";
```



```
# hex
```

```
$snowman = "\x{2603}";
```

```
# name
```

```
$snowman = "\N{SNOWMAN}";
```


Character Encoding Recap



Léon



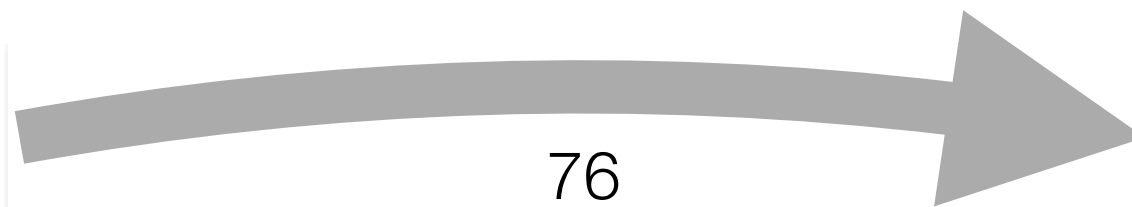
Léon





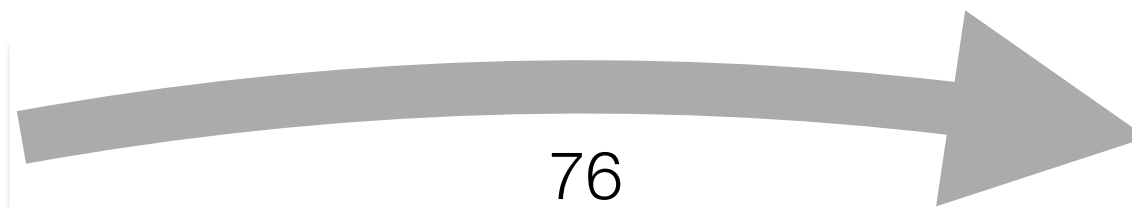
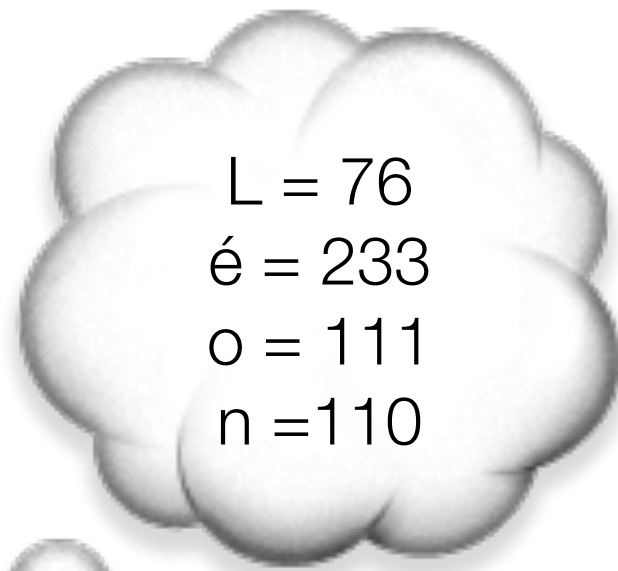
01001100
11101001
01101111
01101110



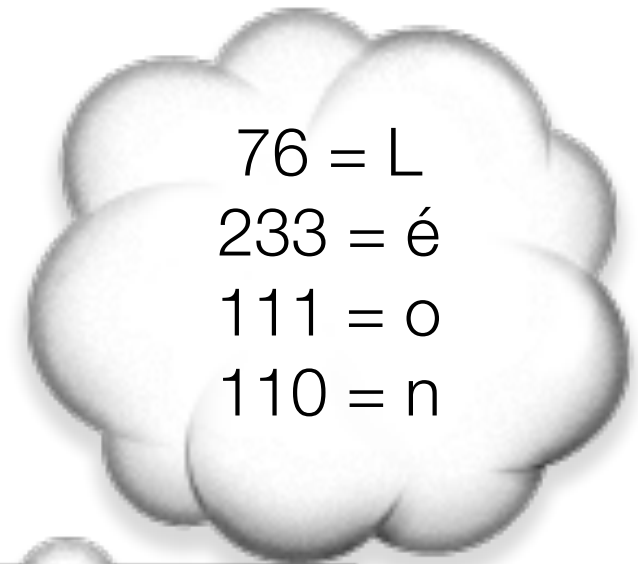


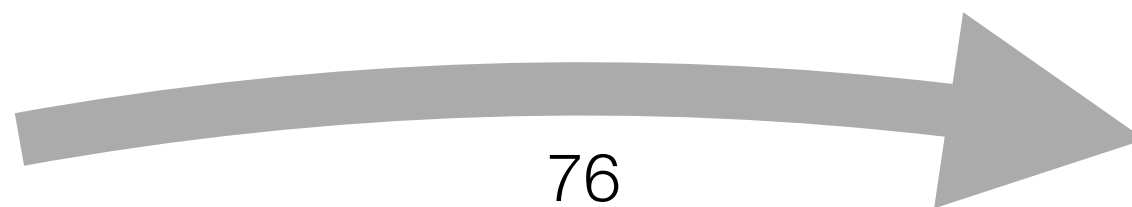
76
233
111
110





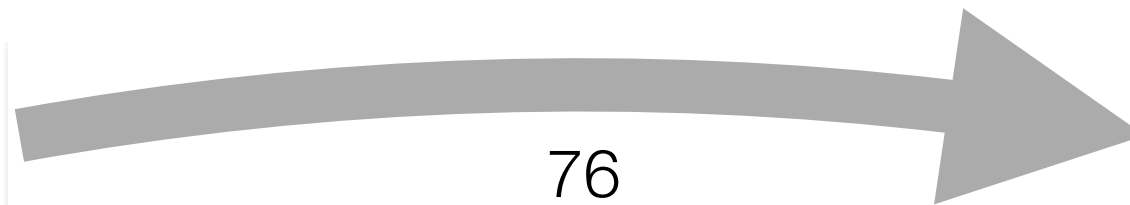
76
233
111
110





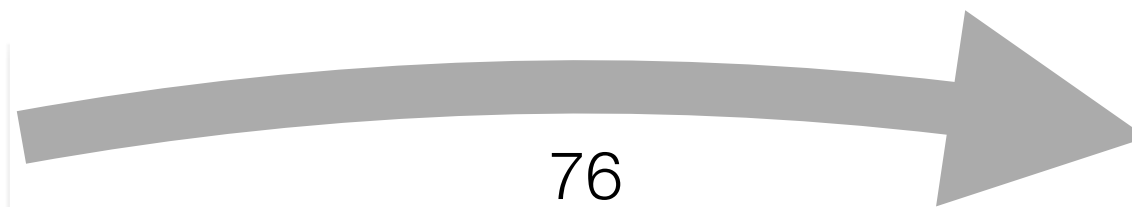
76
233
111
110





76
233
111
110

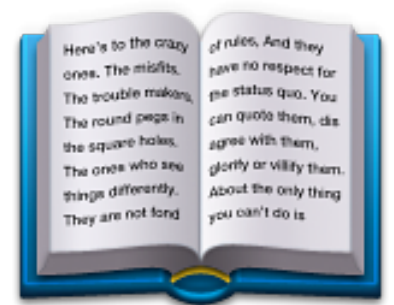




76
233
111
110



Terminology

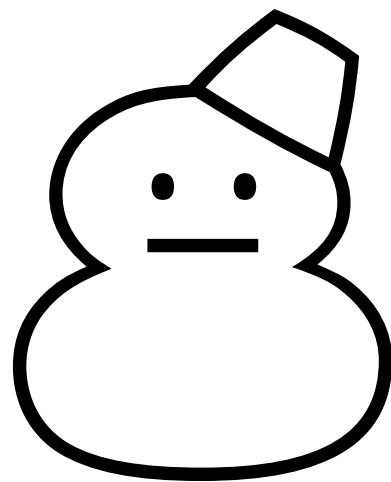


Character

A single letter or symbol

a

é



Character String

One letter or symbol followed by another,
followed by another, and so on...an array of
characters

anna

elsa

Characters are
abstract

Byte

An eight bit number, between 0-255 inclusive

72

144

253

Byte String

A chunk of memory, a series of bytes, that may somehow have characters encoded in it...an array of bytes

76 233

111 110

Bytes are what
your code reads

Bytes are what
your code **prints**

Encoding

Going from a Character String to a Byte String

Decoding

Going from a Byte String to a Character String

Character Set Encoding

What bytes are used to represent a character
in a byte string

Common Character Sets



ASCII

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	
0	0 <small>\000 0x00</small>	1 <small>\001 0x01</small>	2 <small>\002 0x02</small>	3 <small>\003 0x03</small>	4 <small>\004 0x04</small>	5 <small>\005 0x05</small>	6 <small>\006 0x06</small>	7 <small>\007 0x07</small>	8 <small>\010 0x08</small>	9 <small>\011 0x09</small>	10 <small>\012 0x0A</small>	11 <small>\013 0x0B</small>	12 <small>\014 0x0C</small>	13 <small>\015 0x0D</small>	14 <small>\016 0x0E</small>	15 <small>\017 0x0F</small>	0
	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	
	☐	┐	└	┘	↯	✚	✓	⤵	↶	➤	≡	∇	⇓	⇐	⊗	⊙	
1	16 <small>\020 0x10</small>	17 <small>\021 0x11</small>	18 <small>\022 0x12</small>	19 <small>\023 0x13</small>	20 <small>\024 0x14</small>	21 <small>\025 0x15</small>	22 <small>\026 0x16</small>	23 <small>\027 0x17</small>	24 <small>\030 0x18</small>	25 <small>\031 0x19</small>	26 <small>\032 0x1A</small>	27 <small>\033 0x1B</small>	28 <small>\034 0x1C</small>	29 <small>\035 0x1D</small>	30 <small>\036 0x1E</small>	31 <small>\037 0x1F</small>	1
	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US	
	☐	⌚	⌚	⌚	⌚	↯	⌚	┐	⌚	⌚	?	⊖	☐	☐	☐	☐	
2	32 <small>\040 0x20</small>	33 <small>\041 0x21</small>	34 <small>\042 0x22</small>	35 <small>\043 0x23</small>	36 <small>\044 0x24</small>	37 <small>\045 0x25</small>	38 <small>\046 0x26</small>	39 <small>\047 0x27</small>	40 <small>\050 0x28</small>	41 <small>\051 0x29</small>	42 <small>\052 0x2A</small>	43 <small>\053 0x2B</small>	44 <small>\054 0x2C</small>	45 <small>\055 0x2D</small>	46 <small>\056 0x2E</small>	47 <small>\057 0x2F</small>	2
	SP										*	+	,	-	.	/	
	△	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	
3	48 <small>\060 0x30</small>	49 <small>\061 0x31</small>	50 <small>\062 0x32</small>	51 <small>\063 0x33</small>	52 <small>\064 0x34</small>	53 <small>\065 0x35</small>	54 <small>\066 0x36</small>	55 <small>\067 0x37</small>	56 <small>\070 0x38</small>	57 <small>\071 0x39</small>	58 <small>\072 0x3A</small>	59 <small>\073 0x3B</small>	60 <small>\074 0x3C</small>	61 <small>\075 0x3D</small>	62 <small>\076 0x3E</small>	63 <small>\077 0x3F</small>	3
	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?	
	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?	
4	64 <small>\060 0x40</small>	65 <small>\061 0x41</small>	66 <small>\062 0x42</small>	67 <small>\063 0x43</small>	68 <small>\064 0x44</small>	69 <small>\065 0x45</small>	70 <small>\066 0x46</small>	71 <small>\067 0x47</small>	72 <small>\068 0x48</small>	73 <small>\069 0x49</small>	74 <small>\070 0x4A</small>	75 <small>\071 0x4B</small>	76 <small>\072 0x4C</small>	77 <small>\073 0x4D</small>	78 <small>\074 0x4E</small>	79 <small>\075 0x4F</small>	4
	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
5	80 <small>\080 0x50</small>	81 <small>\081 0x51</small>	82 <small>\082 0x52</small>	83 <small>\083 0x53</small>	84 <small>\084 0x54</small>	85 <small>\085 0x55</small>	86 <small>\086 0x56</small>	87 <small>\087 0x57</small>	88 <small>\088 0x58</small>	89 <small>\089 0x59</small>	90 <small>\090 0x5A</small>	91 <small>\091 0x5B</small>	92 <small>\092 0x5C</small>	93 <small>\093 0x5D</small>	94 <small>\094 0x5E</small>	95 <small>\095 0x5F</small>	5
	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_	
	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_	
6	96 <small>\060 0x60</small>	97 <small>\061 0x61</small>	98 <small>\062 0x62</small>	99 <small>\063 0x63</small>	100 <small>\064 0x64</small>	101 <small>\065 0x65</small>	102 <small>\066 0x66</small>	103 <small>\067 0x67</small>	104 <small>\068 0x68</small>	105 <small>\069 0x69</small>	106 <small>\070 0x6A</small>	107 <small>\071 0x6B</small>	108 <small>\072 0x6C</small>	109 <small>\073 0x6D</small>	110 <small>\074 0x6E</small>	111 <small>\075 0x6F</small>	6
	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	
	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	
7	112 <small>\080 0x70</small>	113 <small>\081 0x71</small>	114 <small>\082 0x72</small>	115 <small>\083 0x73</small>	116 <small>\084 0x74</small>	117 <small>\085 0x75</small>	118 <small>\086 0x76</small>	119 <small>\087 0x77</small>	120 <small>\088 0x78</small>	121 <small>\089 0x79</small>	122 <small>\090 0x7A</small>	123 <small>\091 0x7B</small>	124 <small>\092 0x7C</small>	125 <small>\093 0x7D</small>	126 <small>\094 0x7E</small>	127 <small>\095 0x7F</small>	7
	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL	
	p	q	r	s	t	u	v	w	x	y	z	{		}	~	///	
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	

Latin-1

	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
00	NUL 0000	STX 0001	SOT 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
10	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
20	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	(0028) 0029	* 002A	+ 002B	, 002C	- 002D	. 002E	/ 002F
30	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	: 003A	; 003B	< 003C	= 003D	> 003E	? 003F
40	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
50	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[005B	\ 005C] 005D	^ 005E	_ 005F
60	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
70	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F
80	€ 20AC		/ 201A		// 201E	… 2026	† 2020	‡ 2021		‰ 2030	Š 0160	< 2039	Š 015A	Ť 0164	Ž 017D	Ž 0179
90		\ 2018	/ 2019	“ 201C	” 201D	• 2022	– 2013	— 2014		™ 2122	š 0161	> 203A	ś 015B	ţ 0165	ž 017E	ž 017A
A0	NBSP 00A0	ˆ 02C7	˘ 02D8	Ł 0141	ł 00A4	À 0104	Á 00A6	Â 00A7	Ã 00A8	Ä 00A9	Å 015E	« 00AB	¬ 00AC	– 00AD	® 00AE	Ž 017B
B0	° 00B0	± 00B1	² 02DB	‡ 0142	´ 00B4	µ 00B5	¶ 00B6	· 00B7	¸ 00B8	¹ 0105	º 015F	» 00BB	ƒ 013D	˜ 02DD	ı 013E	ž 017C
C0	Ř 0154	Á 00C1	Â 00C2	Ă 0102	Ä 00C4	Í 0139	Ć 0106	Ç 00C7	Č 010C	É 00C9	Ê 0118	Ë 00CB	Ě 011A	Í 00CD	Î 00CE	Ď 010E
D0	Đ 0110	Ň 0143	Ň 0147	Ó 00D3	Ô 00D4	Õ 0150	Ö 00D6	× 00D7	Ř 0158	Ů 016E	Ú 00DA	Û 0170	Ü 00DC	Ý 00DD	Ť 0162	ß 00DF
E0	ř 0155	á 00E1	â 00E2	ă 0103	ä 00E4	í 013A	ć 0107	ç 00E7	č 010D	é 00E9	ê 0119	ë 00EB	ě 011B	í 00ED	î 00EE	ď 010F
F0	č 0111	ń 0144	ň 0148	ó 00F3	ô 00F4	õ 0151	ö 00F6	÷ 00F7	ř 0159	ů 016F	ú 00FA	û 0171	ü 00FC	ý 00FD	ţ 0163	· 02D9

Latin-1

	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
00	NUL 0000	STX 0001	SOT 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
10	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
20	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	(0028) 0029	* 002A	+ 002B	, 002C	- 002D	. 002E	/ 002F
30	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	: 003A	; 003B	< 003C	= 003D	> 003E	? 003F
40	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
50	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[005B	\ 005C] 005D	^ 005E	_ 005F
60	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
70	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F
80	€ 20AC		/ 201A		// 201E	… 2026	† 2020	‡ 2021		‰ 2030	Š 0160	< 2039	Š 015A	Ť 0164	Ž 017D	Ž 0179
90		\ 2018	/ 2019	“ 201C	” 201D	• 2022	— 2013	— 2014		™ 2122	š 0161	> 203A	š 015B	ť 0165	ž 017E	ž 017A
A0	NBSP 00A0	˘ 02C7	˘ 02D8	Ł 0141	ł 00A4	À 0104	¡ 00A6	Â 00A7	Ã 00A8	Ä 00A9	Å 015E	« 00AB	¬ 00AC	— 00AD	® 00AE	Ž 017B
B0	° 00B0	± 00B1	² 02DB	‡ 0142	´ 00B4	µ 00B5	¶ 00B6	· 00B7	¸ 00B8	¸ 0105	Š 015F	» 00BB	ƒ 013D	˜ 02DD	ı 013E	ž 017C
C0	Ř 0154	Á 00C1	Â 00C2	Ă 0102	Ä 00C4	Í 0139	Ć 0106	Ç 00C7	Č 010C	É 00C9	Ê 0118	Ë 00CB	Ě 011A	Í 00CD	Î 00CE	Ď 010E
D0	Đ 0110	Ň 0143	Ň 0147	Ó 00D3	Ô 00D4	Õ 0150	Ö 00D6	× 00D7	Ř 0158	Ů 016E	Ú 00DA	Û 0170	Ü 00DC	Ý 00DD	Ť 0162	ß 00DF
E0	ř 0155	á 00E1	â 00E2	ă 0103	ä 00E4	í 013A	ć 0107	ç 00E7	č 010D	é 00E9	ê 0119	ë 00EB	ě 011B	ï 00ED	î 00EE	ď 010F
F0	č 0111	ń 0144	ň 0148	ó 00F3	ô 00F4	õ 0151	ö 00F6	÷ 00F7	ř 0159	ů 016F	ú 00FA	û 0171	ü 00FC	ý 00FD	ţ 0163	· 02D9

Same
As ASCII

Latin-1

	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
00	NUL 0000	STX 0001	SOT 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
10	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
20	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	(0028) 0029	* 002A	+ 002B	, 002C	- 002D	. 002E	/ 002F
30	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	: 003A	; 003B	< 003C	= 003D	> 003E	? 003F
40	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
50	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[005B	\ 005C] 005D	^ 005E	_ 005F
60	` 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
70	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F
80	€ 20AC		/ 201A		// 201E	… 2026	† 2020	‡ 2021		‰ 2030	Š 0160	< 2033	Š 015A	Ť 0164	Ž 017D	Ž 0179
90		\ 2018	/ 2019	“ 201C	” 201D	• 2022	— 2013	— 2014		™ 2122	Š 0161	> 203A	Š 015B	ť 0165	ž 017E	ž 017A
A0	NBSP 00A0	˘ 02C7	˘ 02D8	Ł 0141	ł 00A4	À 0104	¡ 00A6	Â 00A7	Ã 00A8	Ä 00A9	Å 015E	« 00AB	¬ 00AC	— 00AD	® 00AE	Ž 017B
B0	° 00B0	± 00B1	µ 02DB	ł 0142	´ 00B4	µ 00B5	¶ 00B6	· 00B7	¸ 00B8	¸ 0105	Š 015F	» 00BB	Ł 013D	˘ 02DD	ı 013E	ž 017C
C0	Ř 0154	Á 00C1	Â 00C2	Ă 0102	Ä 00C4	Í 0139	Ć 0106	Ç 00C7	Č 010C	É 00C9	Ê 0118	Ë 00CB	Ě 011A	Í 00CD	Î 00CE	Ď 010E
D0	Đ 0110	Ň 0143	Ň 0147	Ó 00D3	Ô 00D4	Õ 0150	Ö 00D6	× 00D7	Ř 0158	Ů 016E	Ú 00DA	Û 0170	Ü 00DC	Ý 00DD	Ť 0162	ß 00DF
E0	ř 0155	á 00E1	â 00E2	ă 0103	ä 00E4	í 013A	ć 0107	ç 00E7	č 010D	é 00E9	ê 0119	ë 00EB	ě 011B	ï 00ED	î 00EE	ď 010F
F0	č 0111	ñ 0144	ň 0148	ó 00F3	ô 00F4	õ 0151	ö 00F6	÷ 00F7	ř 0159	ů 016F	ú 00FA	û 0171	ü 00FC	ý 00FD	ţ 0163	· 02D9

Just in
Latin-1

Unicode

1,111,998 Different Characters

109,384 assigned in version 6.0
First 256 are the same as Latin-1

Unicode

1,111,998 Different Characters

109,384 assigned in version 6.0
First 256 are the same as Latin-1

byte \neq codepoint

Since a single byte can only
map to 256 code points

Unicode

1,111,998 Different Characters

109,384 assigned in version 6.0
First 256 are the same as Latin-1

byte != codepoint

Since a single byte can only
map to 256 code points

multi-byte encodings

UTF-8

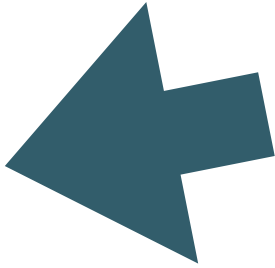
a → \x{61}

F → \x{70}

é → \x{c3}\x{a9}

☺ → \x{e2}\x{98}\x{83}

UTF-8

a	→	\x{61}	
F	→	\x{70}	
é	→	\x{c3}\x{a9}	
☺	→	\x{e2}\x{98}\x{83}	

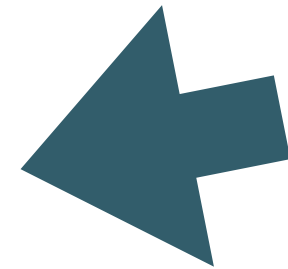
UTF-8

a → \x{61}

F → \x{70}

é → \x{c3}\x{a9}

☹ → \x{e2}\x{98}\x{83}



Perl Unicode Basics



Meanwhile In Java

```
byte[] bytes = getBytes();  
String str = new String(bytes, "UTF-8");
```

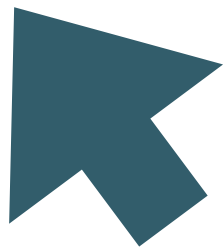
Meanwhile In Java



```
byte[] bytes = getBytes();  
String str = new String(bytes, "UTF-8");
```

Meanwhile In Java

```
byte[] bytes = getBytes();  
String str = new String(bytes, "UTF-8");
```



Meanwhile In Java

```
byte[] bytes = getBytes();  
String str = new String(bytes, "UTF-8");
```



Back in Perl: Just Scalars

```
use Encode qw(decode);
```

```
my $bytes = getBytes();  
$str = decode("utf-8", $bytes);
```

Back in Perl: Just Scalars

```
use Encode qw(decode);
```

```
my $bytes = getBytes();
```

```
$str = decode("utf-8", $bytes);
```

It's all up to you

- You must remember which scalars are:
 - Character Strings
 - Byte Strings
- You must remember which characters are:
 - Really Characters
 - Characters containing bytes

Back in Perl: Just Scalars

```
use Encode qw(decode);
```

```
my $bytes = getBytes();  
$str = decode("utf-8", $bytes);
```

Back in Perl: Just Scalars

```
use Encode qw(encode decode);
```

```
my $bytes = getBytes();
```

```
$str = decode("utf-8", $bytes);
```

```
my $bytes2 = encode("utf-8", $str);
```

Automatic With Filehandles

```
open my $fh1, $filename_to_read;  
my $string = decode(scalar(<$fh1>), "utf-8");
```

```
open my $fh2, $filename_to_write;  
print $fh2, encode("utf-8",$string);
```

Automatic With Filehandles

```
open my $fh1, "<:utf-8", $filename_to_read;  
my $string = <$fh1>;
```

```
open my $fh2, ">:utf-8", $filename_to_write;  
print $fh2, $string
```

Automatic With Filehandles

```
open my $fh1, "<:utf-8", $filename_to_read;  
my $string = <$fh1>;
```

```
open my $fh2, ">:utf-8", $filename_to_write;  
print $fh2, $string
```

Automatic With Filehandles

```
binmode $fh1, "utf-8";  
binmode $fh2, "utf-8";
```

“Th-th-th-that's all, folks!”

– Porky Pig

Behind The Curtain



Behind the Curtain

76 233 111 110

→ Léon

Behind the Curtain

76 233 111 110



→ Léon

76 195 169 111 110



→ Léon

226 152 131



→ 

Debugging Perl Unicode



“There’s always another encoding bug, you just haven’t found it yet”

– Fowler’s Rule on Encoding

bytey

bytey

```
Mark@travis:~$ cat /tmp/example.txt
Sven the m  se and Olaf the   
Mark@travis:~$ xxd /tmp/example.txt
Sven the m  se and Olaf the   
Mark@travis:~$ cat /tmp/example.txt | xxd
Sven the m  se and Olaf the   
Mark@travis:~$
```

<http://tinyurl.com/getbytey>

bytey

```
#!/usr/bin/env perl -s

use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```

bytey

```
#!/usr/bin/env perl -s

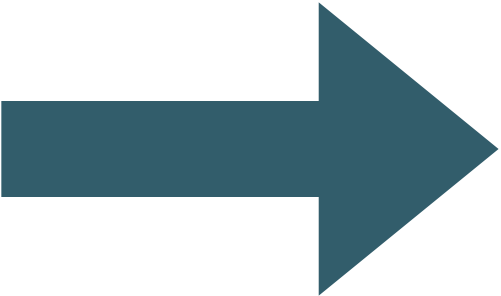
use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```


bytey



```
#!/usr/bin/env perl -s

use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```

bytey

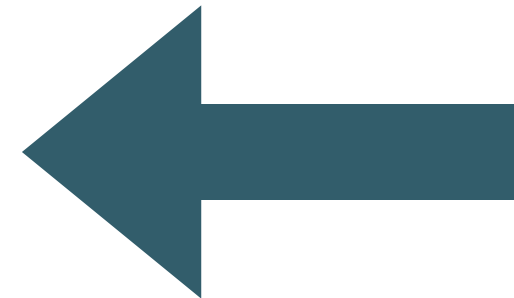
```
#!/usr/bin/env perl -s

use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```



bytey

```
#!/usr/bin/env perl -s

use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```



bytey

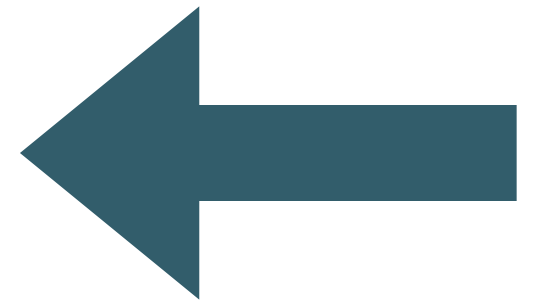
```
#!/usr/bin/env perl -s

use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```



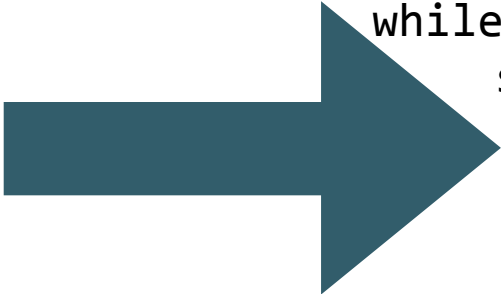
bytey

```
#!/usr/bin/env perl -s
```

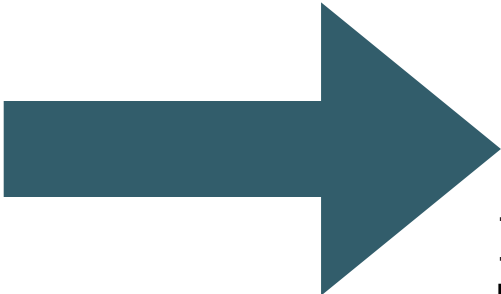
```
use strict;  
use warnings;
```

```
use Term::ANSIColor qw(colored);
```

```
$::r = $::r ? "" : '(![\x{1b}][\d+m)';
```



```
while (<>) {  
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{  
        colored(  
            join(  
                "",  
                map {  
                    '\x{' . sprintf('%02x',ord($_)) . '}'  
                } split(//, $1)  
            },  
            'yellow',  
        )  
    }eg;  
    print  
}
```



bytey

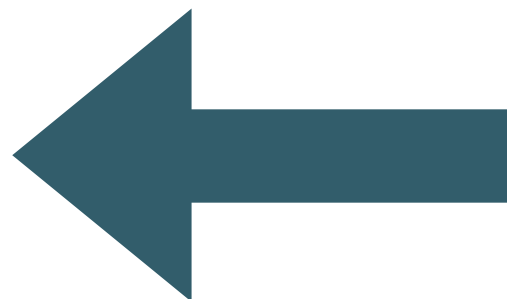
```
#!/usr/bin/env perl -s

use strict;
use warnings;

use Term::ANSIColor qw(colored);

$::r = $::r ? "" : '(?![\\x{1b}]\\[\\d+m)';

while (<>) {
    s{$::r([^\x{0a}\x{20}\x{21}-\x{7e}]+)}{
        colored(
            join(
                "",
                map {
                    '\\x{' . sprintf('%02x',ord($_)) . '}'
                } split(//, $1)
            ),
            'yellow',
        )
    }eg;
    print
}
```



Devel::Peek

Devel::Peek

```
perl -e 'use Devel::Peek; Dump "L\x{e9}on"'  
SV = PV(0x7fd812004280) at 0x7fd812011368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK)  
  PV = 0x7fd811c0dde0 "L\351on"\0  
  CUR = 4  
  LEN = 16
```


Devel::Peek



```
perl -e 'use Devel::Peek; Dump "L\x{e9}on"'  
SV = PV(0x7fd812004280) at 0x7fd812011368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK)  
  PV = 0x7fd811c0dde0 "L\351on"\0  
  CUR = 4  
  LEN = 16
```

Devel::Peek

```
perl -e 'use Devel::Peek; Dump "L\x{e9}on"'  
SV = PV(0x7fd812004280) at 0x7fd812011368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK)  
  PV = 0x7fd811c0dde0 "L\351on"\0  
  CUR = 4  
  LEN = 16
```



Devel::Peek

```
perl -e 'use Devel::Peek; Dump "\N{SNOWMAN}"'  
SV = PV(0x7f91399b49c0) at 0x7f9139811368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK,UTF8)  
  PV = 0x7f913940dde0 "\342\230\203"\0  
    [UTF8 "\x{2603}"]  
  CUR = 3  
  LEN = 16
```

Devel::Peek

```
perl -e 'use Devel::Peek; Dump "\N{SNOWMAN}"'  
SV = PV(0x7f91399b49c0) at 0x7f9139811368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK,UTF8)  
  PV = 0x7f913940dde0 "\342\230\203"\0  
    [UTF8 "\x{2603}"]  
  CUR = 3  
  LEN = 16
```




Devel::Peek

```
perl -e 'use Devel::Peek; Dump "\N{SNOWMAN}"'  
SV = PV(0x7f91399b49c0) at 0x7f9139811368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK,UTF8)  
  PV = 0x7f913940dde0 "\342\230\203"\0  
    [UTF8 "\x{2603}"]  
  CUR = 3  
  LEN = 16
```



Devel::Peek

```
perl -e 'use Devel::Peek; Dump "\N{SNOWMAN}"'  
SV = PV(0x7f91399b49c0) at 0x7f91399b49c0:1368  
  REFCNT = 1  
  FLAGS = (POK,READONLY,pPOK,UTF8)  
  PV = 0x7f913940dde0 "\342\230\203"\0  
    [UTF8 "\x{2603}"]  
  CUR = 3  
  LEN = 16
```



Test::utf8

Test::utf8

check the string is good

```
is_valid_string($string);    # check the string is valid
is_sane_utf8($string);       # check not double encoded
```

check the string has certain attributes

```
is_flagged_utf8($string1);   # has utf8 flag set
is_within_ascii($string2);   # only has ascii chars in it
isnt_within_ascii($string3); # has chars outside the ascii range
is_within_latin_1($string4); # only has latin-1 chars in it
isnt_within_ascii($string5); # has chars outside the latin-1 range
```


Test::utf8

check the string is good

`is_valid_string($string);` # check the string is valid
`is_sane_utf8($string);` # check not double encoded

check the string has certain attributes

`is_flagged_utf8($string1);` # has utf8 flag set
`is_within_ascii($string2);` # only has ascii chars in it
`isnt_within_ascii($string3);` # has chars outside the ascii range
`is_within_latin_1($string4);` # only has latin-1 chars in it
`isnt_within_ascii($string5);` # has chars outside the latin-1 range

Test::utf8

check the string is good

```
is_valid_string($string);    # check the string is valid
is_sane_utf8($string);       # check not double encoded
```

check the string has certain attributes

```
is_flagged_utf8($string1);   # has utf8 flag set
is_within_ascii($string2);   # only has ascii chars in it
isnt_within_ascii($string3); # has chars outside the ascii range
is_within_latin_1($string4); # only has latin-1 chars in it
isnt_within_latin_1($string5); # has chars outside the latin-1 range
```

Test::utf8

check the string is good

is_valid_string(\$string); # check the string is valid
is_sane_utf8(\$string); # check not double encoded

check the string has certain attributes

is_flagged_utf8(\$string1); # has utf8 flag set
is_within_ascii(\$string2); # only has ascii chars in it
isnt_within_ascii(\$string3); # has chars outside the ascii range
is_within_latin_1(\$string4); # only has latin-1 chars in it
isnt_within_ascii(\$string5); # has chars outside the latin-1 range

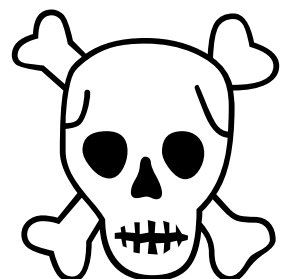
Test::utf8

```
# This will pass as it's a normal latin-1 string
is_sane_utf8("Hello L\x{e9}on");
```

```
# this will fail because the \x{c3}\x{a9} looks like the
# utf8 byte sequence for e-acute
my $string = "Hello L\x{c3}\x{a9}on";
is_sane_utf8($string);
```

```
# this will pass because the utf8 is correctly interpreted as utf8
Encode::_utf8_on($string)
is_sane_utf8($string);
```

Common Problems



Raw Printing a non-ASCII Unicode Character

Raw printing Unicode

Léon & 🧊

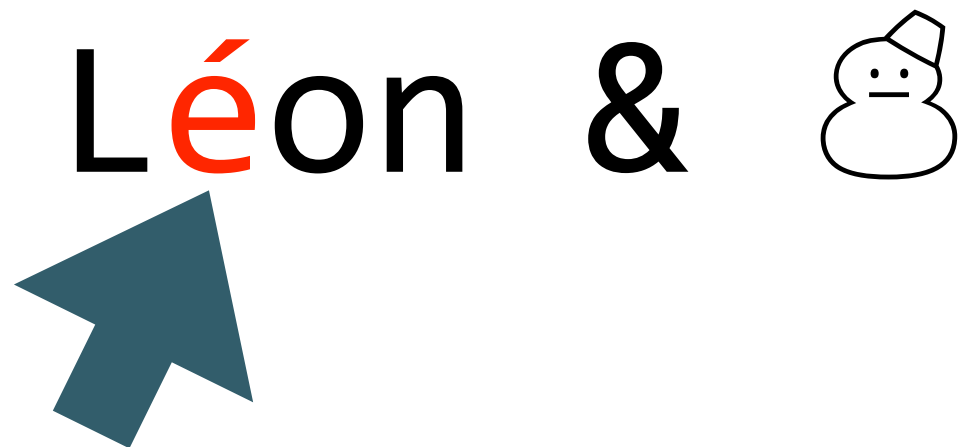
Raw printing Unicode

“send byte `\x{4c}`”



Léon & 🤖

Raw printing Unicode



“send byte `\x{e9}`”

Raw printing Unicode

Léon & 



“send byte ...?”

Raw printing Unicode

Léon & 



“Ah screw it send utf-8 byte
sequence `\x{e2}\x{98}\x{83}`”

Raw printing Unicode

Wide character in print at -e line 1.

Léon & 



“Ah screw it send utf-8 byte
sequence `\x{e2}\x{98}\x{83}`”

Accidental Unicode Character in Byte String

UTF-8 Upgrade (good)

```
perl -e 'use Devel::Peek; Dump "L\x{e9}on" . "\N{SNOWMAN}"'
SV = PV(0x7fa428804290) at 0x7fa428811320
  REFCNT = 1
  FLAGS = (POK,READONLY,pPOK,UTF8)
  PV = 0x7fa428505b60 "L\303\251on\342\230\203"\0
    [UTF8 "L\x{e9}on\x{2603}"]
  CUR = 8
  LEN = 16
```

```
Mark@travis:~$ perl -E 'binmode STDOUT, 'utf8'; say "L\x{e9}on" . "\N{SNOWMAN}"'
Léon☸
Mark@travis:~$ perl -E 'binmode STDOUT, 'utf8'; say "L\x{e9}on" . "\N{SNOWMAN}"' | byte
L\x{c3}\x{a9}on\x{e2}\x{98}\x{83}
```

UTF-8 Upgrade (bad)

```
perl -e 'use Devel::Peek; Dump "L\x{c3}\x{a9}on" . "\N{SNOWMAN}"'
SV = PV(0x7fc072004290) at 0x7fc072011320
  REFCNT = 1
  FLAGS = (POK,READONLY,pPOK,UTF8)
  PV = 0x7fc071c0dc10 "L\x{303}\x{203}\x{302}\x{251}on\x{342}\x{230}\x{203}"\0
    [UTF8 "L\x{c3}\x{a9}on\x{2603}"]
  CUR = 10
  LEN = 16
```

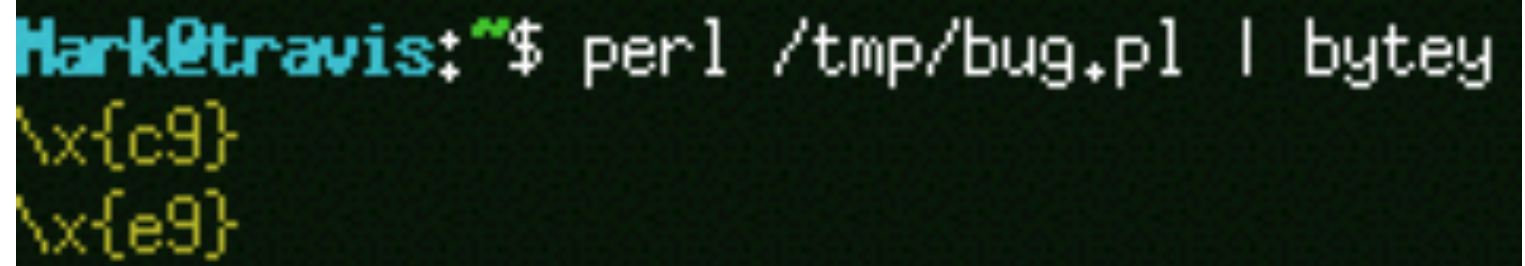
```
Mark@travis:~$ perl -E 'binmode STDOUT, 'utf8'; say "L\x{c3}\x{a9}on" . "\N{SNOWMAN}"'
LÃon❄
Mark@travis:~$ perl -E 'binmode STDOUT, 'utf8'; say "L\x{c3}\x{a9}on" . "\N{SNOWMAN}"' | byte
L\x{c3}\x{83}\x{c2}\x{a9}on\x{e2}\x{98}\x{83}
```

“The Unicode Bug”

The Unicode Bug

```
my $string = "\x{c9}";  
print lc $string, "\n";
```

```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```

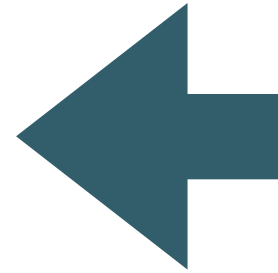


A terminal window with a black background. The prompt is 'Mark@travis:~\$'. The command entered is 'perl /tmp/bug.pl | byte'. The output consists of two lines: '\x{c9}' and '\x{e9}'. The prompt and command are in white, and the output is in yellow.

```
Mark@travis:~$ perl /tmp/bug.pl | byte  
\x{c9}  
\x{e9}
```

The Unicode Bug

```
my $string = "\x{c9}";  
print lc $string, "\n";
```



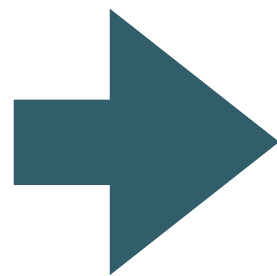
```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```

```
Mark@travis:~$ perl /tmp/bug.pl | byte  
\x{c9}  
\x{e9}
```

The Unicode Bug

```
my $string = "\x{c9}";  
print lc $string, "\n";
```

```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```

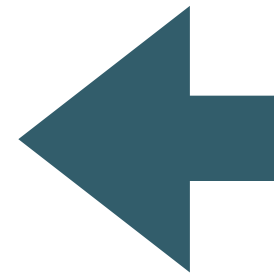


```
Mark@travis:~$ perl /tmp/bug.pl | byte  
\x{c9}  
\x{e9}
```

The Unicode Bug

```
my $string = "\x{c9}";  
print lc $string, "\n";
```

```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```

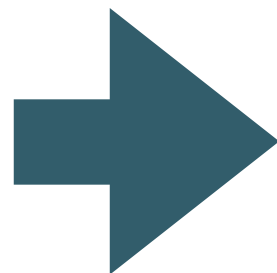


```
Hark@travis:~$ perl /tmp/bug.pl | byte  
\x{c9}  
\x{e9}
```

The Unicode Bug

```
my $string = "\x{c9}";  
print lc $string, "\n";
```

```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```



```
Mark@travis:~$ perl /tmp/bug.pl | byte  
\x{c9}  
\x{e9}
```

The Unicode Bug

128-255 = ???

The Unicode Bug

```
use feature 'unicode_strings';
```

```
my $string = "\xc9";  
print lc $string, "\n";
```

```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```

The Unicode Bug

```
use 5.12.0;
```

```
my $string = "\x{c9}";  
print lc $string, "\n";
```

```
$string .= "\N{SNOWMAN}";  
chop $string;  
print lc $string, "\n";
```


XS routines loses the UTF-8 flag

XS and UTF-8

- Bad XS code copies SV's PV but not UTF-8 flag
 - Flip flag with `Encode::_utf8_on`
 - or Re-do the decode again
 - or use `Data::Structure::Util` in anger
 - or PATCH THE DARN XS CODE

Re-encode a byte string a second time

Re-decode a character string a second time

“Doctor, it hurts when I do this.”

“Don't do that.”

– Dr. Kronkheit and His Only Living Patient, Smith & Dale, 1920s

Tricks and Tips



Unicode in Perl Source Code

Unicode In Source Files

```
#!/usr/bin/perl

use 5.18.0;
use warnings;

say length "Léon";

{
    use utf8;
    say length "Léon";
}
```

Unicode In Source Files

```
Hark@travis:~$ bytey /tmp/unicode.txt
#!/usr/bin/perl

use 5.18.0;
use warnings;

say length "L\x{c3}\x{a9}on";

{
    use utf8;
    say length "L\x{c3}\x{a9}on";
}
Hark@travis:~$
```


Unicode In Source Files

```
#!/usr/bin/perl
```

```
use 5.18.0;  
use warnings;
```

```
say length "Léon";
```

```
# prints "5"
```

```
{  
    use utf8;  
    say length "Léon";  
}
```

```
# prints "4"
```

Unicode Regex

Matching Digits With a Regex

`/^\d+$/`

Matching Digits With a Regex

```
perl -E 'say "100" =~ /^\\d+$/ ? "yes" : "no"';  
yes
```

```
perl -E 'say "abc" =~ /^\\d+$/ ? "yes" : "no"';  
no
```

Matching Digits With a Regex

```
perl -E 'say "100" =~ /\d+$/ ? "yes" : "no";'  
yes
```

```
perl -E 'say "abc" =~ /\d+$/ ? "yes" : "no";'  
no
```

Matching Digits With a Regex

```
perl -E 'say "100" =~ /\d+$/ ? "yes" : "no"';  
yes
```

```
perl -E 'say "abc" =~ /\d+$/ ? "yes" : "no"';  
no
```

Matching Digits With a Regex

```
perl -E 'say "\x{09EA}" =~ /\d+$/ ? "yes" : "no"';
```

Matching Digits With a Regex

```
perl -E 'say "\x{09EA}" =~ /\d+$/ ? "yes" : "no";'  
yes
```

09EA BENGALI DIGIT FOUR

Matching Digits With a Regex



```
perl -E 'say "\x{09EA}" =~ /\d+$/a ? "yes" : "no";  
no
```

Matching Digits With a Regex

/a

Make `\d`, `\s`, `\w` match only in ASCII range

Matching A-Z With a Regex

`/^[a-z]+$ /i`

Matching A-Z With a Regex

```
perl -E 'say "\x{212A}" =~ /^[a-z]+$/i ? "yes":"no"';
```

Matching A-Z With a Regex

```
perl -E 'say "\x{212A}" =~ /^[a-z]+$/i ? "yes":"no"';  
yes
```

212A KELVIN SIGN

Matching A-Z With a Regex



```
perl -E 'say "\x{212A}" =~ /^[a-z]+$/i'aa ? "yes":"no";  
no
```

212A KELVIN SIGN

Matching A-Z With a Regex

`/aa`

Make `\d`, `\s`, `\w` match only in ASCII range

Make ASCII/non-ASCII never match

Data::Structure::Utils

Handy Dandy Utility Functions

`has_utf8($var)`

Returns `$var` if the utf8 flag is enabled for `$var` or any scalar that a data structure passed in `$var` contains.

`_utf8_off($var)`

Recursively disables the utf8 flag on all scalars within `$var`.

`_utf8_on($var)`

Recursively enables the utf8 flag on all scalars within `$var`.

`utf8_on($var)`

This routine performs a "sv_utf8_upgrade" on each scalar string in the passed data structure that does not have the utf8 flag turned on.

`utf8_off($var)`

This routine performs a "sv_utf8_downgrade" on each scalar string in the passed data structure that has the utf8 flag turned on.

Converting Anything To ASCII

Converting Anything To ASCII

```
use Text::Unidecode;  
say unidecode("北京");
```

Bei Jing

Normalized Form

md5 bytes only

```
my $code = md5sum( $passwd_str );
```

md5 bytes only

```
my $passwd_bytes = encode("utf-8", $passwd_str);  
my $code = md5sum( $passwd_bytes );
```

When é isn't é

é → \x{e9}

When é isn't é

é → \x{65} \x{0301}

When é isn't é

“LATIN SMALL LETTER E”



é → \x{65} \x{0301}



“COMBINING ACUTE ACCENT”

md5 bytes only

```
use Unicode::Normalize;
$passwd_str = NFC($string);
my $passwd_bytes = encode("utf-8", $passwd_str);
my $code = md5sum( $passwd_bytes );
```

ETX