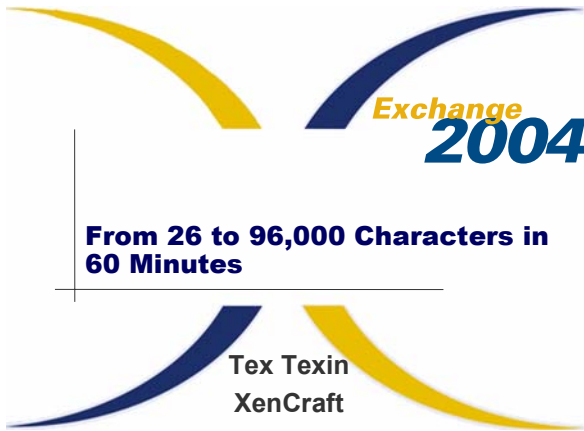


From 26 to 96,000 Characters in 60 Minutes



Objectives

- Modern software integrates with many technologies
- Unicode is the basis for text processing in many environments
- Objective is to inform on many aspects of working with Unicode

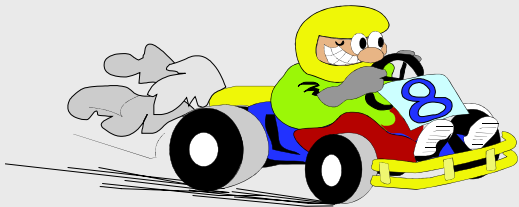
From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



On Your Mark! Get Set! GO! 26 Letters

- ABCDEFGHIJKLMNOPQRSTUVWXYZ



From 26 to 96,000 Characters in 60 Minutes

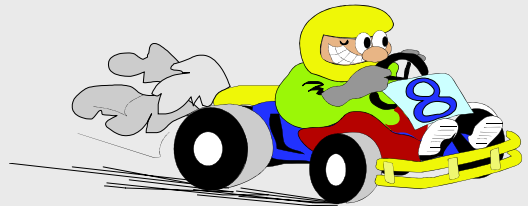
3

Copyright © 2004 Tex Texin. All rights reserved.



52 Letters, so I don't have to SHOUT!

- ABCDEFGHIJKLMNOPQRSTUVWXYZ
- abcdefghijklmnopqrstuvwxyz



From 26 to 96,000 Characters in 60 Minutes

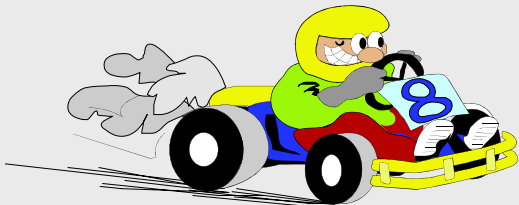
4

Copyright © 2004 Tex Texin. All rights reserved.



Characters Not Just Letters

- ABCDEFGHIJKLMNOPQRSTUVWXYZ
- abcdefghijklmnopqrstuvwxyz
- Punctuation : . , ; ? ! () - ...



From 26 to 96,000 Characters in 60 Minutes

5

Copyright © 2004 Tex Texin. All rights reserved.



Arithmetic

- ABCDEFGHIJKLMNOPQRSTUVWXYZ
- abcdefghijklmnopqrstuvwxyz
- Punctuation : . , ; ? ! () - ...
- 0 1 2 3 4 5 6 7 8 9 + - ± * × / ÷ < = > % ‰
- # ¼ ½ ¾ ⅓ ⅔ ⅛ ⅞ IV iii viii IX ix XII

From 26 to 96,000 Characters in 60 Minutes

6

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

Currency and Business Symbols

- ABCDEFGHIJKLMNOPQRSTUVWXYZ
- abcdefghijklmnopqrstuvwxyz
- Punctuation : . , ; ? ! () - ...
- 0 1 2 3 4 5 6 7 8 9 + - ± × ÷ < = > % ‰
- # ¼ ½ ¾ ⅓ ⅔ ⅛ ⅜ IV iii viii IX ix XII
- \$ ¢ £ ¥ € ₣ ₧ ₨ ₪ ₮ ₩ © ® ™ % ° °C °F ∞

Copyright © 2004 Tex Texin. All rights reserved



Common or Specialty Symbols

- ABCDEFGHIJKLMNOPQRSTUVWXYZ
- abcdefghijklmnopqrstuvwxyz
- Punctuation : . , ; ? ! () - ...
- 0 1 2 3 4 5 6 7 8 9 + - ± * × / ÷ < = > % ‰
- # ¼ ½ ¾ ⅓ ⅔ ⅛ ⅞ IV iii viii IX ix XII
- \$ ¢ £ ¥ € ₣ ₧ ₨ ₪ ₮ ₩ © ® ™ ° °C °F ∞
- ¶ § ♀ ♂ ♠ ♥ ♦ ♣ ← ↑ ↓ ↔ ↵ ↶ ↷ ↸ ↹ ↺ ↻ ↼ ↽ ↾ ↿ ⇄ ⇅ ⇆ ⇇ ⇈ ⇉ ⇊ ⇋ ⇌ ⇍ ⇎ ⇏ ⇐ ⇑ ⇒ ⇓ ⇔ ⇕ ⇖ ⇗ ⇘ ⇙ ⇚ ⇛ ⇜ ⇝ ⇞ ⇟ ⇠ ⇡ ⇢ ⇣ ⇤ ⇥ ⇦ ⇧ ⇨ ⇩ ⇪ ⇫ ⇬ ⇭ ⇮ ⇯ ⇰ ⇱ ⇲ ⇳ ⇴ ⇵ ⇶ ⇷ ⇸ ⇹ ⇺ ⇻ ⇼ ⇽ ⇾ ⇿
- ▼ ▽ ▹ ▸ ◂ ▸ ◃ ◅ ◆ ◇ ◈ ◉ ◊ ○ ◌ ◍ ◎ ● ◐ ◑ ◒ ◓ ◔ ◕ ◖ ◗ ◘ ◙ ◚ ◛ ◜ ◝ ◞ ◟ ◠ ◡ ◢ ◣ ◤ ◥ ◦ ◧ ◨ ◩ ◪ ◫ ◬ ◭ ◮ ◯ ◰ ◱ ◲ ◳ ◴ ◵ ◶ ◷ ◸ ◹ ◺ ◻ ◼ ◽ ◾ ◿

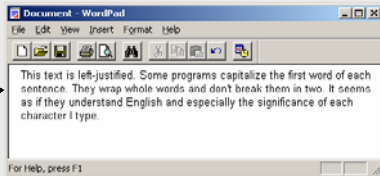
Copyright © 2004 Tex Texin. All rights reserved.



Writing and Computing

How do computers process text?

- Input? Display? Justification? Printing?



Copyright © 2004 Tex Texin. All rights reserved



Text is a sequence of numbers

- T e x t i s a s e q
- 54 65 78 74 20 69 73 20 61 20 73 65 71
- u e n c e o f n u
- 75 65 6E 63 65 20 6F 66 20 6E 75
- m b e r s .
- 6D 62 65 72 73 2E

Copyright © 2004 Tex Texin. All rights reserved.



“Properties” define behavior for each

- 54 Use the Glyph (image) “T”
- is-character TRUE
- is-digit FALSE
- is-punctuation FALSE
- is-upper TRUE
- tolower 74
- toupper 54
- sort-rank 94

Copyright © 2004 Tex Texin. All rights reserved



Why is T the number 54 (hex)?

- Arbitrary
- History (ASCII, ISO 646, ECMA 6)
- It isn't always (EBCDIC)

	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
10	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	ESC	FS	GS	RS	US	
20	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	
30	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
40	'	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
50	P	Q	R	S	T	U	V	W	X	Y	Z	[~]	^	_
60	@	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
70	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

IBM Extended Binary Coded Decimal Interchange Code (EBCDIC)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	Nul		DS		SP				a	j			A	J		0
1			SOS			/			b	k						1
2			FS						c	l	s		B	K	S	2
3		TM							d	m	t		C	L	T	3
4	PF	RES	BYP	PN	RS				e	n	u		D	M	U	4
5	HT	NL	LF	PN					f	o	v		E	N	V	5
6	LC	BS	EOB	UC					g	w	x		F	O	W	6
7	DEL	IL	PRE	EOT					h	p	y		G	P	X	7
8									i	q	z		H	Q	Y	8
9										r			I	R	Z	9
A			SM		g	!	:									
B					"	\$,	#								
C					<	.	%	@								
D					()	-	'								
E					+	>	=									
F						?	"									

Assigning Characters Numeric Values

- **Anyone can do it!**

- Computer, Printer, Phone, Monitor, and other hardware vendors
- Operating systems, font, and other software vendors
- Governments
- Standards Organizations (ISO, NCITS, ECMA, JIS, etc.)

But sometimes it is confusing

- Is ellipses ... one character or three?
 - How many fractions do I need? $\frac{1}{4}$ $\frac{1}{2}$ $\frac{3}{4}$ $\frac{1}{3}$
 - Is “viii” a digit?
 - How many characters is ™? %?
 - Are these writing symbols or drawing?
- ♀ ♂ ♠ ♣ ♥ ♦ ← ↑ ↓ ↔ ↵ ↶ ↷ ↸ ↹ ↻ ↺ ↻ ↻ ↻ ↻ ↻ ↻ ↻ ↻ ↻
- Are these: “ Π Σ ” mathematical symbols or Greek Letters? Does it matter?

Choosing numeric values

- Often it depends on your application(s)

- If you are writing a book about music or card games, suits and notes may be a natural and necessary part of text flow.
 - Math programs give “ $\prod \sum$ ” special semantics. They may have one value for math operators and another for letters.
 - Searching and Sorting are affected.
- Search for “.” or “/”. Return “...”, or “ $\frac{1}{2}$ ”?
- Search for “8” or “l”. Return “viii”?

Choosing numeric values

- Plain text requirements differ from rich text.
 - OpenEdge™
 - OpenEdgeTM
- Note- Characters from other languages haven't been discussed yet.

Cyrillic Alphabet

Аа	Бб	Вв	Гг	Дд	Ее	Ёё	Жж	Зз
Ии	Йй	Кк	Лл	Мм	Нн	Оо	Пп	Рр
Сс	Тт	Уу	Фф	Хх	Цц	Чч	Шш	Щщ
Ъъ	Ыы	Ьь	Ээ	Юю	Яя			

	А	Б	В	Г	Д	Е	З	И	Й	К	Л	М	Н	О	П	Р	С	Т	У	Ф	Х	Ц	Ч	Щ	Ъ	Ы	Э	Ю
0-7F	ASCII																											
80-9F	а	б	в	г	д	е	з	и	й	к	л	м	н	о	п	р	с	т	у	ф	х	ц	ч	щ	ъ	ы	э	ю
A0-BF	а	б	в	г	д	е	з	и	й	к	л	м	н	о	п	р	с	т	у	ф	х	ц	ч	щ	ъ	ы	э	ю
C0-DF	А	Б	В	Г	Д	Е	З	И	Й	К	Л	М	Н	О	П	Р	С	Т	У	Ф	Х	Ц	Ч	Щ	Ъ	Ы	Э	Ю
E0-FF	а	б	в	г	д	е	з	и	й	к	л	м	н	о	п	р	с	т	у	ф	х	ц	ч	щ	ъ	ы	э	ю

From 26 to 96,000 Characters in 60 Minutes

Greek Alphabet

- ΑΒΓΔΕΖΗΘΙΚΑΜΝΞΟΠΡΣΤΥΦΧΨΩ
- αβγδεζηθικλμνξοπρςστυφχψω
- Should A, B, E, etc. be same as Latin A, B, E?

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
8	€		,	f	„	...	†	‡	^	0/00	§	‘	œ			ÿ
9		‘	’	“	”	•	—	™	~							
A	“	”	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í
B	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	Ø
C	ƒ	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î
D	Π	Ρ	Σ	Τ	Υ	Φ	Χ	Ψ	Ω	Ι	Κ	Λ	Μ	Ν	Ξ	Ο
E	ϐ	α	β	γ	δ	ε	ζ	η	θ	ι	κ	λ	μ	ν	ξ	ο
F	π	ρ	σ	τ	υ	φ	χ	ψ	ω	ϊ	ϋ	ϗ	ϙ	ϛ	ϝ	ϟ

From 26 to 96,000 Characters in 60 Minutes

19

Copyright © 2004 Tex Texin. All rights reserved.



Small character sets

- Character Set: Collection of “characters” needed in your applications.
- Small sets (<256) require less memory and simpler property implementations.
 - Small number of glyphs
 - Property arrays indexed by numeric value
- May be incomplete and not extensible
 - Language is always evolving
 - ISO 8859-1 is an example that is incomplete

From 26 to 96,000 Characters in 60 Minutes

20

Copyright © 2004 Tex Texin. All rights reserved.



ISO 8859-1 is full, but Euro Needed

- ISO 8859-15 deleted characters
 - acute accent “ ´ ” cedilla “ ¸ ”
 - broken bar “ ¯ ” diaeresis “ ¨ ”
 - international currency symbol “ ₧ ”
 - Fractions ¼ ½ ¾
- To make room for new ones
 - € Ÿ Š š Œ œ Ž ž
- Similarity to ISO 8859-1 leads to errors

From 26 to 96,000 Characters in 60 Minutes

21

Copyright © 2004 Tex Texin. All rights reserved.



ISO 8859-15

80																
90																
A0	NSP	ı	ç	€	€	Ÿ	Š	š	Œ	œ	Ž	ž	ı	ı	ı	ı
B0	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
C0	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
D0	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
E0	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
F0	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı

From 26 to 96,000 Characters in 60 Minutes

22

Copyright © 2004 Tex Texin. All rights reserved.



If a character is not in the current code page, how can it be expressed?

ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı

West European Code Page ISO 8859-1

ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı

Russian Code Page 1251

From 26 to 96,000 Characters in 60 Minutes

23

Copyright © 2004 Tex Texin. All rights reserved.



If a file contains a 255, what character does it represent?

ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı

West European Code Page ISO 8859-1

ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı
ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı	ı

Russian Code Page 1251

From 26 to 96,000 Characters in 60 Minutes

24

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

Characters are not in sorted order

0 – 7F: ASCII

80 □ □ , f , n ... t i % § < E □ □ □ □ , m , - " m j > o □ □ Ÿ

i j k l m n o p q r s t u v w x y z

A A A A A A E C E E E E I I I I N O O O O X U U U U Y Z

a á â ã ä å æ ç è é ê ë ì í î ï ð ñ ò ó ô õ ö ø ÷ ù ú û ü ý þ

FF

West European Code Page ISO 8859-1

FOR EACH CUSTOMER WHERE NAME <= CHR(255):

This does not return all records.

For example, it misses names beginning with “Z” or “z”.

From 26 to 96,000 Characters in 60 Minutes

25

Copyright © 2004 Tex Texin. All rights reserved.



Does language determine encoding?

- In Spanish, “ch” is a letter.
- It sorts after “c” and before “d”
 - Color, Charlar, Dar
- “ch” is not (usually) encoded separately.
- These ligatures are:
 - IJ ij Œ œ Æ æ ff fi fl ffi ffl

From 26 to 96,000 Characters in 60 Minutes

26

Copyright © 2004 Tex Texin. All rights reserved.



Bidirectional Scripts

- **Hebrew, Arabic, Urdu and other languages are written right to left.**
- **Often they are mixed with left to right text.**
- **Encodings can be visual or logical**
 - Visual – text stored as displayed
 - Logical- text stored as logically advancing, but may display differently.

From 26 to 96,000 Characters in 60 Minutes

27

Copyright © 2004 Tex Texin. All rights reserved.



Visual vs. Logical

Be logical, not visual

- ▶ Logical & visual contrasted

פעילות הבינאום, W3C

W3C, מואניבה תוליעפ

פעילות הבינאום, W3C

From 26 to 96,000 Characters in 60 Minutes

28

Copyright © 2004 Tex Texin. All rights reserved.



Direction

- **Direction is sometimes ambiguous.**
 - The number 764 in Hebrew is: תשטז.
 - (764, תשטז)
- **Characters have a directional property**
 - Strong, Weak, Neutral
 - Neutral chars take direction of the “run”.
 - Special override characters exist.
- **Some Hebrew characters have 2 forms:**

כך צץ פה נזמם

From 26 to 96,000 Characters in 60 Minutes

29

Copyright © 2004 Tex Texin. All rights reserved.



Arabic has more complex rendering

- Right to left like Hebrew
- Both have diacritic marks
- Characters can have 4 shapes:
 - Isolate , Initial, Medial, Final
 - Beh:
- Complex rendering
 - lam ﻝ alef ا lam-alef ﻻ
 - kashida (elongation of connectors for justification)

From 26 to 96,000 Characters in 60 Minutes

36

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

Arabic example page from Quran



From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



Japanese, Chinese, Korean

Languages with >255 characters

- Japanese ~ 50,000 characters
- Korean 2172 Hangul
- Simplified Chinese ~ 6000+
- Traditional Chinese ~ 17000 encoded
- GB 18030 ~ Unicode ~ 70,000 Ideographs

Writing direction

- Vertically, columns going right to left.
- Sometimes columns left to right.
- Or written horizontally left to right.

A b c 日 本 語 d e

From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



Japanese, Chinese, Korean

- Spaces are not used for word separation.
- No upper or lower case.
- Sorting is done a number of ways
 - binary, stroke-radical, radical-stroke, phonetic.
 - dictionaries disagree on stroke counts and orderings
 - 日 + 月 = 明
- Input methods needed for large character set
- Intricate characters require taller fonts
- Trie and other techniques for efficiency with large or sparse data are used.

From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



Japanese Encodings

Several approaches to large character sets

- JIS X208-1990 Character Set
 - 6879 character grid 95x95 (33-126x33-126)
- Encoded as ISO 2022-jp, SJIS, EUC-JIS
- ISO 2022-jp shift-sequenced (aka escape sequenced) 7-bit
- EUC-JIS 8-bit conversion of JIS x208 values (add 128 to each byte)
- SJIS Shift-JIS complex reassignments plus halfwidth katakana

From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



Escape-sequenced character sets

- Escape sequence changes character set
 - Shifter hex values character set
 - Esc (B ASCII
 - Esc (J JIS Roman (JIS X 0201-1976)
 - Esc \$ B JIS X 0208-1983
- ISO 2022 defines sequences for other non-Japanese character sets as well.
 - Esc(B English Text Esc\$B Nihongo...

From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



Double-byte Character Sets (DBCS) Japanese, Chinese, Korean

Mixed size characters: 1 or 2 bytes

ASCII compatible

	A	b	c	日	本	語	d	e
Text:	S	S	S	L	T	L	T	S
# Chars:	1	2	3	4	4	5	5	6
# Bytes:	1	2	3	4	5	6	7	8

How long is a DBCS string?

ABC123456789 A=61
A B C 1 2 3 中區ΦΨ A=A2CF

From 26 to 96,000 Characters in 60 Minutes

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

Double-byte Character Sets

- There are many double-byte character sets
 - Shift-JIS for Japanese,
 - EUCJIS for Japanese
 - GB2312, CP936 for Simplified Chinese,
 - KSC 5601, CP 949, CP 1361 for Korean,
 - Big-5, CP950 for Traditional Chinese

From 26 to 96,000 Characters in 60 Minutes

37

Copyright © 2004 Tex Texin. All rights reserved.



Unicode

- Created and maintained by consortium www.unicode.org
- Unicode covers all major living scripts
- Version 4.0 has 96,000+ characters
- Capacity for 1 million+ characters
- Unicode Character Set = ISO 10646
 - Unicode adds character properties and algorithms
 - ISO and Unicode work together to synchronize
 - ISO support enhances international acceptance

From 26 to 96,000 Characters in 60 Minutes

38

Copyright © 2004 Tex Texin. All rights reserved.



Unicode Characteristics

- 16 bit design originally
- Now has 3 equivalent forms
 - UTF-8: 8-bit variable width, multi-byte (max. 4)
 - UTF-16: 16-bit, variable width, surrogates (max 2)
 - UTF-32: 32-bit, fixed width (max 1)
- Designed to avoid multi-byte performance problems
- Precise algorithm specifications provide interoperability
- Allows one binary program image to be used worldwide
- Developers do not need to be linguists to implement

From 26 to 96,000 Characters in 60 Minutes

39

Copyright © 2004 Tex Texin. All rights reserved.



Unicode is Generative

- Character composition can create “new” characters
- Base character + non-spacing character(s)

$A + \circ = \mathring{A}$
U+0041 + U+030A = U+00C5

From 26 to 96,000 Characters in 60 Minutes

40

Copyright © 2004 Tex Texin. All rights reserved.



Unicode Characteristics

- Multilingual
 - All scripts, one encoding
- Semantic properties
 - Case, digit, alpha/letter/ideogram, directional class, mirroring, combining class, etc. provided by Unicode
- Logical order
 - Text is stored in logical order, not visual
- Convertibility
 - Accurate round-trip for legacy character sets
- Byte Order Mark (BOM) for endianness, identifier

From 26 to 96,000 Characters in 60 Minutes

41

Copyright © 2004 Tex Texin. All rights reserved.



Unicode — Han Unification

- 70,000 Han characters
 - More than 21 national standards (around 120,000 original characters)
 - Same Chinese root for Kanji, Hanzu and Hanja ideograms
- Drawbacks:
 - No national sort order
 - Disagreement among some scholars

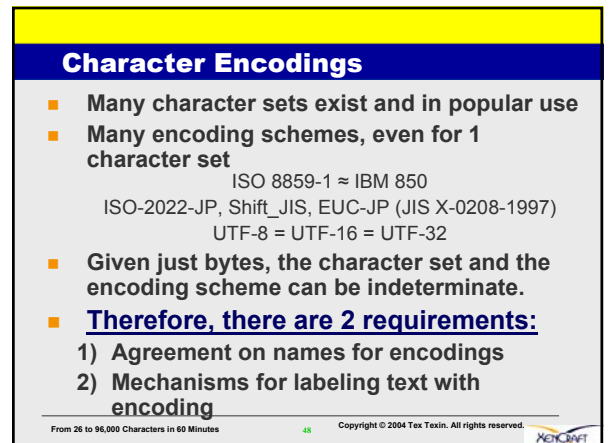
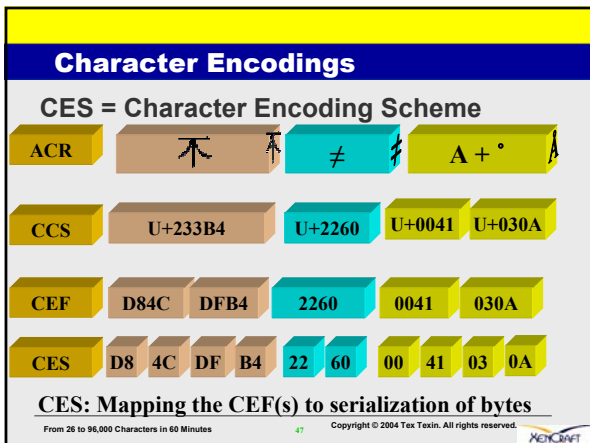
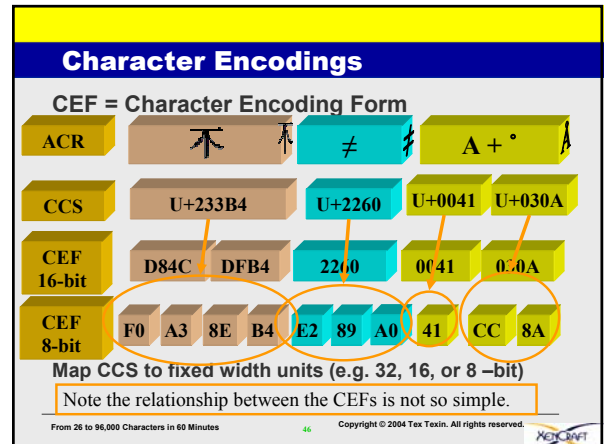
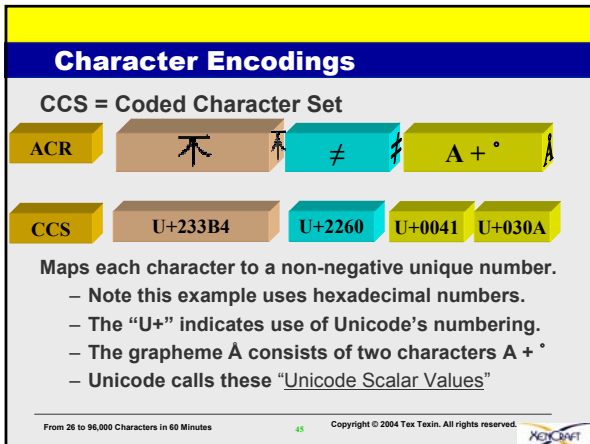
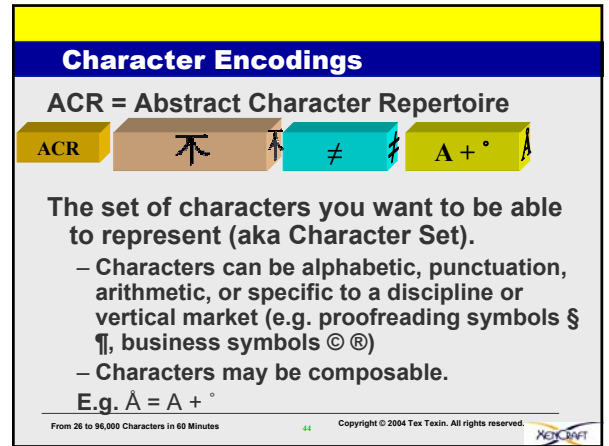
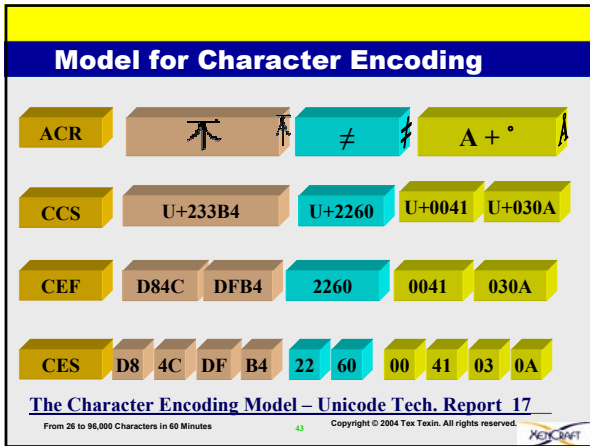
From 26 to 96,000 Characters in 60 Minutes

42

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes



From 26 to 96,000 Characters in 60 Minutes

Character Encodings

Names of encodings

IANA (Internet Assigned Numbers Authority) maintains registry of official names for character sets (actually encodings) used on the internet
<http://www.iana.org/assignments/character-sets>

From 26 to 96,000 Characters in 60 Minutes

49

Copyright © 2004 Tex Texin. All rights reserved.



Character Encodings

- Names consist of ASCII, printable characters
 - case-insensitive, 40 characters max.
 - Aliases are also registered
- Unregistered Names
 - Begin with the name convention "x-"
 - Example: "x-Tex-Yves-encoding"
 - Useful for private encodings or very new encodings
 - Not useful on the web, except for private exchange
- MIME also uses IANA charset names

From 26 to 96,000 Characters in 60 Minutes

50

Copyright © 2004 Tex Texin. All rights reserved.



Character Encodings

- IANA Name Examples (Aliases)
 - ISO_8859-1:1987 (ISO_8859-1, ISO-8859-1, latin1, L1, IBM819, CP819, csISOLatin1)
 - Windows-1252, GB2312, BIG5, BIG5-HKSCS
 - SHIFT_JIS, HP-Legal
 - Extended_UNIX_Code_Packed_Format_for_Japanese
 - Adobe-standard-encoding
 - UTF-8, UTF-16, UTF-16BE, UTF-16LE, UTF-32
- Note- Registry contains many useless names
- Note- Preferred names indicated. Use them.

From 26 to 96,000 Characters in 60 Minutes

51

Copyright © 2004 Tex Texin. All rights reserved.



Storage and Serialization Formats

- UTF-32
 - 32 bits per character
 - Uses Unicode scalar value as-is
 - Big-endian and Little-endian
 - Unicode only goes to 10FFFF (21 bits)
- UTF-16
 - 16 bits per character
 - To support more than 65K characters, uses two "surrogate" values

From 26 to 96,000 Characters in 60 Minutes

52

Copyright © 2004 Tex Texin. All rights reserved.



Properties of UTF-16

Unicode Scalar	UTF-16	
0-FFFF	0-FFFF (Unchanged)	
	High	Low
10000	D800	DC00
10FFFF	DBFF	DFFF

21 Scalar value bits: **u uuuu xxxx xx**yy yyyy yyyy
www = uuuuu - 1
Format: High Surrogate / Low Surrogate
1101 10ww wwxx xxxx / 1101 11yy yyyy yyyy
– Note: Surrogates don't conflict with other values or each other

From 26 to 96,000 Characters in 60 Minutes

53

Copyright © 2004 Tex Texin. All rights reserved.



Properties of UTF-8

- Transforms Unicode to sequences of octets:
 - 7-bit ASCII is in 8-bits with 8th bit = 0
 - All other byte values have 8th bit = 1
 - Non-ASCII characters are either 2, 3 or 4 bytes
- Result:
 - Algorithms searching for ASCII characters (e.g., / \ < > ? + - a b c d etc.) work correctly
 - String length is not greatly increased
 - All of Unicode supported

From 26 to 96,000 Characters in 60 Minutes

54

Copyright © 2004 Tex Texin. All rights reserved.



Properties of UTF-8

- **First byte**
 - 0XXX XXXX = one byte character
 - 110X XXXX = two byte character
 - 1110 XXXX = three byte character
 - 1111 0XXX = four byte character
- **Bytes after first**
 - 10XX XXXX

From 26 to 96,000 Characters in 60 Minutes 55 Copyright © 2004 Tex Texin. All rights reserved.



Properties of UTF-8

Unicode Scalar: X XXXX XXXX XXXX XXXX XXXX

1-7F: 0xxxxxxx

80-7FF: 110xxxxx 10xxxxxx

800-FFFF: 1110xxxx 10xxxxxx 10xxxxxx

10000-10FFFF: 11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

From 26 to 96,000 Characters in 60 Minutes 56 Copyright © 2004 Tex Texin. All rights reserved.



Choosing a UTF

- **UTF-8**
 - Good choice for migrating legacy software and file formats (ASCII compatibility, multi-byte encoding)
 - Best storage form for European languages
- **UTF-16**
 - More efficient for sorting, processing
 - Best storage for Asian languages
 - Requires wide character datatypes
 - Good choice for new implementations
- **UTF-32**
 - efficient processing, wastes memory

From 26 to 96,000 Characters in 60 Minutes 57 Copyright © 2004 Tex Texin. All rights reserved.



FAQ

- **Does Text Storage Double?**
 - Although text expands, not by two
- **Performance**
 - UTF-8 algorithms are higher performance than DBCS.
 - UTF-16 surrogate characters are rarely used.
 - (But surrogates must be supported.)
- **Tip for merging regional databases**
 - Remember to realign sequences

From 26 to 96,000 Characters in 60 Minutes 58 Copyright © 2004 Tex Texin. All rights reserved.



UNICODE Byte Order Mark

- Indicates byte ordering or endian-ness for UTF-16 and UTF-32
- Identifies data in file as UNICODE
- Specifies UNICODE encoding
 - ASCII 74
 - UTF-16
 - Big Endian FE FF 00 74
 - Little Endian FF FE 74 00
 - UTF-8 EF BB BF 74

From 26 to 96,000 Characters in 60 Minutes 59 Copyright © 2004 Tex Texin. All rights reserved.



When Do You Drop the BOM?

- BOM will be first character in the file
- Be sure to remove before concatenation
- Not needed for UTF-8. But must be accepted.
- Must not be output if encoding is UTF-16BE, UTF-16LE, UTF-32BE, UTF-32LE

From 26 to 96,000 Characters in 60 Minutes 60 Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

String Indexing Or When 2+2=5?

- Which units should be used for counting?

Graphemes 3	天	天	≠	天	天	天
Characters 4	U+233B4	U+2260	U+0041	U+030A		
Code units 5	D84C	DFB4	2260	0041	030A	
Bytes 10	D8	4C	DF	B4	22	60

From 26 to 96,000 Characters in 60 Minutes

61

Copyright © 2004 Tex Texin. All rights reserved.



String Indexing

Character Model indexing recommendations

- Character counting is recommended for most programming interfaces (e.g. XML Path)
- Code unit counting may be used for internal efficiency (e.g. DOM)
- Graphemes may be useful for user interaction, once a suitable definition exists
- Avoid creating API with single unit arguments
e.g. "SS" = Uppercase("ß")

From 26 to 96,000 Characters in 60 Minutes

62

Copyright © 2004 Tex Texin. All rights reserved.



What is Normal?

Choosing a Canonical Form

- Representing data in more than 1 way leads to errors
- E.g. The Mars Climate Orbiter mission was disastrous. Information expected to be metric, was sent in English units
- Solution- Adopt 1 standard representation-
Normalize



From 26 to 96,000 Characters in 60 Minutes

63

Copyright © 2004 Tex Texin. All rights reserved.



Normalization

Unicode characters can have more than 1 representation

- Canonical equivalence
 - Indistinguishable, fundamental equivalence
 - E.g. combining sequences, singletons
 - "Å" U+00C5 (A-ring pre-composed)
 - "A+" U+0041 + U+030A (A + combining ring above)
 - "Å" U+212B (Angstrom)
- Compatibility equivalence
 - E.g. Formatting differences, ligatures
 - "力" U+FF76 "力" U+30AB (KA half and full width)
 - "fi" U+FB01 (ligature fi)

From 26 to 96,000 Characters in 60 Minutes

64

Copyright © 2004 Tex Texin. All rights reserved.



Early Uniform Normalization

- Unicode Consortium has defined canonical and compatibility decomposition formats and 4 different sets of rules for normalization: "Unicode Normalization Forms"
<http://www.unicode.org/unicode/reports/tr15/>
 - NFD - Canonical Decomposition
 - NFKD - Compatibility Decomposition
 - NFC - Canonical Decomposition, followed by Canonical Composition
 - NFKC - Compatibility Decomposition, followed by Canonical Composition
- Each form gives different results

From 26 to 96,000 Characters in 60 Minutes

65

Copyright © 2004 Tex Texin. All rights reserved.



Normalization

- The W3C Character Model recommends Normalization Form C (NFC) for Web
 - Brings canonical equivalences to composed form
 - Leaves compatibility forms as distinct
 - Most legacy text is composed, and is unchanged
- International Domain Names use NFKC
 - so similar looking characters won't be distinct
 - 'Å' U+00C5 vs. 'Å' U+0212B vs. 'A' + '°' U+0041+U+030A
 - 'u' + '°' vs. 'ü'
 - 'P' vs. 'P' (fullwidth vs halfwidth)
 - 'œ' vs. 'oe'

From 26 to 96,000 Characters in 60 Minutes

66

Copyright © 2004 Tex Texin. All rights reserved.



Web Normalization

Basic principles

- Without agreement on text representation, binary matching and string indexing fail
- Consequences are significant
 - E.g. Comparison of contracts, security
- Most existing text is composed
- Encrypted strings require normalization first
- Producers have information about the strings they create, simplifying normalization.
- Character escapes taken into account

From 26 to 96,000 Characters in 60 Minutes

67

Copyright © 2004 Tex Texin. All rights reserved.



Unicode Vs. Markup

- 96,000+ characters as of Unicode 4.0
 - Should we use them all?
 - Are there any we shouldn't use?
 - Does Unicode's capabilities, needed for plain text, interfere with markup?
- Markup can do some things better than character codes. Not all Unicode characters are needed.

From 26 to 96,000 Characters in 60 Minutes

68

Copyright © 2004 Tex Texin. All rights reserved.



Unicode Vs. Markup

Potential problem areas

- Redundancies impact searching
 - “Å” A-ring “A+” A+ring “Å” Angstrom
- Formatting characters vs. Markup
 - E.g. Bidi controls, interlinear annotation characters
- Characters with style vs. Markup
 - E.g. Superscript, subscript
- Object Replacement Character vs. Markup
 - Better to use markup to include an image

From 26 to 96,000 Characters in 60 Minutes

69

Copyright © 2004 Tex Texin. All rights reserved.



Unicode Vs. Markup

Solution types

- Restrict characters so they cannot be used
- Replace redundancies (normalization)
- Replace with Markup
 - Extensible
 - presentation can be separate from content

Joint W3C and Unicode recommendations in:
“Unicode in XML and other Markup Languages”

<http://www.w3.org/TR/unicode-xml/>

<http://www.unicode.org/unicode/reports/tr20/>

From 26 to 96,000 Characters in 60 Minutes

70

Copyright © 2004 Tex Texin. All rights reserved.



Resource Identifiers: URIs, IRIs

- Currently: URIs encode bytes, not characters
- Most ASCII bytes expressed as ASCII
- Other bytes are %HH
- Non-ASCII characters are ambiguous
 - Character encoding is not taken into consideration
- IRI-Internationalized Resource Identifiers
 - Transcode to UTF-8, then encode as URI
 - Adopters: IE, XLink, XPointer, URN, XML, XML Schema
 - <http://www.w3.org/International/O-URL-and-ident.html>
 - <http://www.w3.org/International/iri-edit/draft-duerst-iri-07.txt>

From 26 to 96,000 Characters in 60 Minutes

71

Copyright © 2004 Tex Texin. All rights reserved.



Code page conversions

- Many standards are not standard
 - E.g. Variants of SHIFT-JIS
 - Even ASCII has variants
 - Currency symbols often = 5C (backslash)
- Round trip mappings are sometimes not possible
 - eg. encoding has A-ring and Angstrom, the other has one

From 26 to 96,000 Characters in 60 Minutes

72

Copyright © 2004 Tex Texin. All rights reserved.



From 26 to 96,000 Characters in 60 Minutes

SHIFTJIS chars that vary

- 0x5C(YEN SIGN),
- 0x7E(OVERLINE),
- 0x815C(FULLWIDTH EM DASH),
- 0x815F(REVERSE SOLIDUS),
- 0x8160(WAVE DASH),
- 0x8161(DOUBLEVERTICAL LINE),
- 0x817C(MINUS SIGN),
- 0x8191(CENT SIGN),
- 0x8192(POUND SIGN), and
- 0x81CA(NOT SIGN).

From 26 to 96,000 Characters in 60 Minutes

73

Copyright © 2004 Tex Texin. All rights reserved.



Code page conversions

Solutions

- Accurately identify the source encoding variation
- Base conversion on application or semantics
 - e.g. Yen in currency, Backslash in Filename
- Adopt internal standard
- Unicode Normalization
 - Identifies preferred encoding values
- Maintain the original with the conversion
 - To enable round tripping

From 26 to 96,000 Characters in 60 Minutes

74

Copyright © 2004 Tex Texin. All rights reserved.



Unicode does not equal internationalization

- Unicode simplifies development
 - Single source code
 - Enables multilingual processing
 - Properties reduce research for each language
- Unicode does not fix all internationalization
 - E.g. Date, time, number and other formats
 - Linguistic processing can require additional algorithms, data (e.g. word breaking)
 - Continue identify, support cultural requirements
 - Conversion to native encodings for interface to legacy software, systems can impose limitations

From 26 to 96,000 Characters in 60 Minutes

75

Copyright © 2004 Tex Texin. All rights reserved.



Unicode Support

- Employed by all modern technologies
- Standard for web technologies
 - Web consortium's Reference Processing Model
 - Java, HTML, XML, DOM, etc.
 - Supported by browsers IE, Mozilla, Opera
- Support on most platforms and databases
 - Internal code set for Windows NT, 2000, XP
 - Supported by Sun, IBM, HP, Oracle, SQL Server
- Supported by development tools and libraries
 - IBM ICU (International Components for Unicode)

From 26 to 96,000 Characters in 60 Minutes

76

Copyright © 2004 Tex Texin. All rights reserved.



Nice Places To Visit

- Հայաստան Armenia
- Nunavut
- 中国 China
- مصر Egypt
- Ethiopia
- Ελλάδα Greece
- 대한민국 S. Korea
- Россия Russia

From 26 to 96,000 Characters in 60 Minutes

78

Copyright © 2004 Tex Texin. All rights reserved.



Conclusions

- Unicode is well-supported, ubiquitous, and often required in integrated environments.
- Unicode simplifies working with different languages
- But the large character set requires some additional considerations.

From 26 to 96,000 Characters in 60 Minutes

79

Copyright © 2004 Tex Texin. All rights reserved.

