



# Multi-Language Character Sets

**Bob Balaban, President  
Looseleaf Software, Inc.  
<http://www.looseleaf.net>**

**What They Are, How to  
Use Them**

Click here  
to add clip  
art



---

# Agenda

- Speaker introduction
- Background: terminology
- History of electronic character sets
- The need for multi-language character sets
- LMBCS - Lotus MultiByte Character Set
  - How it works
  - Why no one else uses it



---

# Agenda - 2

- Unicode
  - How it works
  - Why everyone uses it
- Software for multi-language character sets
  - LMBCS
  - Unicode
  - Java, C, LotusScript
- Web applications and MLCS



---

# Speaker Introduction

- Master's degree in Chinese History/Anthropology
- Professional software developer since 1978
- Engineer at Lotus Development Corp. 1987-97
- Reviewer of original LMBCS specification (1988)



---

## Speaker - 2

- Team leader for adding LMBCS/Kanji support to 123/G (os/2)
- Developer on Notes/Domino 1993-97
  - Author of LotusScript "back-end classes"
  - Author of Java APIs for back-end classes
- Founded Looseleaf Software, Inc., 1997
  - Custom development, training, architecture/design, consulting
  - Notes/Domino, J2EE, Groove



---

# Terminology

- Textual characters on a computer are really numbers
  - Like everything else
- The numbers, when representing text, are mapped to symbols on a screen
  - or on paper
- The number of symbols you can display depends on the number of bits assigned to the "character set", or "code space"



---

# Terminology - 2

- One "character" is a "code point"
  - Just a number
- The screen symbol is the "glyph"
- The "font" is the style in which the glyphs are displayed
  - serif, sans serif, etc.
- The mapping of numbers to code points defines a "character set", or "code page"



---

## Terminology - 3

- What is a "native character set"?
- The default character set used by the operating system
- Many OS's can handle multiple character sets (Ascii, CP850, etc.)
- Many OS versions depend on a "locale" specification
  - E.g., CP932 in Japan, CP850 in NAmerica





---

# Terminology - 4

- All mappings are essentially arbitrary
- But some have been agreed upon as standards

A = 1

B = 2

...

Z = 26

A possible mapping

---

A = 17

B = 18

....

K = 99

....

Z = 102

Another possible mapping



---

# Terminology - 5

- What are "control codes"?
- Generally, non-printing
  - No visible representation, no glyph
- Invented to manage early printer/terminal devices, and communication protocols
  - Ack/Nak
  - Carriage return, newline, tab, formfeed
  - Bell



---

# Background: History of electronic character sets

- Original character set was Binary Coded Decimal (BCD)
  - IBM
  - 64 characters, 6 bits
  - Known later as "SIXBIT ASCII"
- Developed for IBM punch cards
- Upper case letters, digits, punctuation



---

## Background - 2

- When IBM created the System 360, they extended BCD
  - 1965
  - EBCDIC
  - Extended Binary Coded Decimal Interchange Code
  - 256 code points
  - 64 control codes
  - Upper and lower case!
  - Some slots were empty, customers complained about wasted memory!

# IBM EBCDIC, ca. 1965

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

Done



---

# EBCDIC and Beyond

- The most popular character set through the 1970s
  - Until the PC was born
- IBM used 8 bits because the word size on S360 was 32 bits
- IBM machines (other than PCs) still use EBCDIC
  - There are a few variants, some code points may be different



---

## Beyond - 2

- Note the strange ordering, where the alphabetic sequences are interrupted
- Makes for lots of fun when writing string comparison/sorting code!
  - Or converting between upper and lower case
- Note, no accented characters
  - Or non-Latin glyphs



---

# The Rise of ASCII

- American Standard Code for Information Interchange
  - Or something like that
- Network bandwidth was very expensive, people wanted to save money and utilize 7-bit channels
- And wanted to conserve on memory, eliminate all that "wasted" space in EBCDIC





---

# ASCII - 2

- Adopted by PC OSs (DOS), and by Unix systems
- 32 control codes, punctuation, digits, upper and lower alphabets (0 is special)
  - This time letters were in a rational sequence (caps first)!
  - You could add/subtract 32 to go from upper to lower and back
- Still no accents or "foreign" glyphs

# ASCII

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32;	Space	64	40	100	&#64;	@	96	60	140	&#96;	`
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33;	!	65	41	101	&#65;	A	97	61	141	&#97;	a
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34;	"	66	42	102	&#66;	B	98	62	142	&#98;	b
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35;	#	67	43	103	&#67;	C	99	63	143	&#99;	c
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	\$	68	44	104	&#68;	D	100	64	144	&#100;	d
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	%	69	45	105	&#69;	E	101	65	145	&#101;	e
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	&	70	46	106	&#70;	F	102	66	146	&#102;	f
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39;	'	71	47	107	&#71;	G	103	67	147	&#103;	g
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40;	(	72	48	110	&#72;	H	104	68	150	&#104;	h
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	)	73	49	111	&#73;	I	105	69	151	&#105;	i
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	*	74	4A	112	&#74;	J	106	6A	152	&#106;	j
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	+	75	4B	113	&#75;	K	107	6B	153	&#107;	k
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	,	76	4C	114	&#76;	L	108	6C	154	&#108;	l
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45;	-	77	4D	115	&#77;	M	109	6D	155	&#109;	m
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46;	.	78	4E	116	&#78;	N	110	6E	156	&#110;	n
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47;	/	79	4F	117	&#79;	O	111	6F	157	&#111;	o
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48;	0	80	50	120	&#80;	P	112	70	160	&#112;	p
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49;	1	81	51	121	&#81;	Q	113	71	161	&#113;	q
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50;	2	82	52	122	&#82;	R	114	72	162	&#114;	r
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51;	3	83	53	123	&#83;	S	115	73	163	&#115;	s
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52;	4	84	54	124	&#84;	T	116	74	164	&#116;	t
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	5	85	55	125	&#85;	U	117	75	165	&#117;	u
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	6	86	56	126	&#86;	V	118	76	166	&#118;	v
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	7	87	57	127	&#87;	W	119	77	167	&#119;	w
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56;	8	88	58	130	&#88;	X	120	78	170	&#120;	x
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57;	9	89	59	131	&#89;	Y	121	79	171	&#121;	y
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58;	:	90	5A	132	&#90;	Z	122	7A	172	&#122;	z
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59;	;	91	5B	133	&#91;	[	123	7B	173	&#123;	{
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60;	<	92	5C	134	&#92;	\	124	7C	174	&#124;	
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61;	=	93	5D	135	&#93;	]	125	7D	175	&#125;	}
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62;	>	94	5E	136	&#94;	^	126	7E	176	&#126;	~
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63;	?	95	5F	137	&#95;	_	127	7F	177	&#127;	DEL



---

# The need for multi-language character sets

- Eventually, manufacturers realized that they could maybe sell more software and hardware outside NAmerica if they supported languages other than English
- But how to deal with those "foreign" symbols?
  - While still allowing for interoperability with Latin text?



---

## Multi-language - 2

- Easy! Go back to 8-bit "codepages"!
  - Use the "upper half" of the 8-bit space (another 128 characters beyond ASCII)
- Each computer sold in a "foreign" country would carry a different code point mapping (a "Codepage")
- Again, this system was invented by IBM
  - And done very systematically
  - Each codepage got a numeric designation
  - And a glyph mapping chart in a book









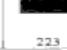
---

# Codepages - continued

- ASCII was always the "lower half" of the codepage
  - So everyone in the world could benefit from using English
- Windows ASCII (again, 8 bits) is cp437
- IBM "international English" is cp850
- Japanese is cp932 (more on this later)
- PRC Simplified Chinese: cp936
- Korean: 949
- Cyrillic: cp1251
- Latin1: cp1252
- etc.

# Codepage 850

850 MS-DOS LATIN 1

	20	30	40	50	60	70	80	90	A0	B0	C0	D0	E0	F0
0		0	@	P	`	p	Ç	É	á		Ł	đ	Ó	Š
	32	48	64	80	96	112	128	144	160	176	192	208	224	240
1	!	1	A	Q	a	q	ü	æ	í		Ł	Đ	ß	±
	33	49	65	81	97	113	129	145	161	177	193	209	225	241
2	"	2	B	R	b	r	é	Æ	ó		τ	Ê	Ô	=
	34	50	66	82	98	114	130	146	162	178	194	210	226	242
3	#	3	C	S	c	s	â	ô	ú		†	Ë	Ò	¾
	35	51	67	83	99	115	131	147	163	179	195	211	227	243
4	\$	4	D	T	d	t	ä	ö	ñ	‡	—	È	õ	¶
	36	52	68	84	100	116	132	148	164	180	196	212	228	244
5	%	5	E	U	e	u	à	ò	Ñ	Á	+	ı	Õ	§
	37	53	69	85	101	117	133	149	165	181	197	213	229	245
6	&	6	F	V	f	v	â	û	ª	Â	ã	Í	µ	÷
	38	54	70	86	102	118	134	150	166	182	198	214	230	246
7	'	7	G	W	g	w	ç	ù	º	À	Ã	Î	þ	ˆ
	39	55	71	87	103	119	135	151	167	183	199	215	231	247
8	(	8	H	X	h	x	ê	ÿ	¿	©	℥	Ï	Þ	°
	40	56	72	88	104	120	136	152	168	184	200	216	232	248
9	)	9	I	Y	i	y	ë	Ö	®	¶	℥	Ɔ	Ú	˚
	41	57	73	89	105	121	137	153	169	185	201	217	233	249
A	*	:	J	Z	j	z	è	Û	—		℥	Ɔ	Û	•
	42	58	74	90	106	122	138	154	170	186	202	218	234	250
B	+	;	K	[	k	{	ï	ø	½	¶	π		Ü	¹
	43	59	75	91	107	123	139	155	171	187	203	219	235	251
C	,	<	L	\	l		î	£	¼	¶	¶		ý	³
	44	60	76	92	108	124	140	156	172	188	204	220	236	252
D	-	=	M	]	m	}	ì	Ø	;	¢	=	!	Ý	²
	45	61	77	93	109	125	141	157	173	189	205	221	237	253
E	.	>	N	^	n	~	Ä	×	«	¥	¶	Î	—	
	46	62	78	94	110	126	142	158	174	190	206	222	238	254
F	/	?	O	_	o	□	Å	f	»	¿	α		'	NBSP
	47	63	79	95	111	127	143	159	175	191	207	223	239	255



---

## But.... Japanese?

- What if you can't fit a culture's "character set" into 256 slots?
- Asian languages, for example:
  - Thai: 40+ alphabetic/phonetic symbols
  - Chinese: 5000 common ideographs, 60,000+ total, TWO different systems!
  - Japanese: same (almost) ideographs as Chinese (trad.), PLUS 1 alphabetic series (Katakana), PLUS 1 phonetic (Hiragana), PLUS double-wide latin letters (so they line up with Kanji)!





---

## Japanese - 2

- But a codepage only has an 8-bit namespace!!?
- So, reserve a sequence (or 2, or 3) in the upper half of the table as "pointers" to another 8-bit space
- As in CP932 (Japanese)



# Codepage 932

## Extract for Lead Byte E9-EA

E9	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
4	顛 9871	顛 9874	顛 9873	風 98AA	風 98AF	颯 98B1	颯 98B6	颯 98C4	颯 98C3	颯 98C6	飢 98E9	飢 98EB	餃 9903	餡 9909	餡 9912	餡 9914
5	餘 9918	餡 9921	饒 991D	饒 991E	饒 9924	餅 9920	餅 992C	饗 992E	饒 993D	饒 993E	饒 9942	饒 9949	饒 9945	饒 9950	饒 994B	饒 9951
6	饒 9952	饒 994C	饒 9955	馗 9997	馗 9998	馗 99A5	馗 99AD	馗 99AE	馗 99BC	馗 99CF	駝 99DB	駝 99DD	駝 99DE	駝 99D1	駝 99ED	駝 99EE
7	駝 99F1	駝 99F2	駝 99FB	駝 99F8	駝 9A01	駝 9A05	駝 99E2	駝 9A19	駝 9A2B	駝 9A37	駝 9A45	駝 9A42	駝 9A40	駝 9A43		
8	駝 9A3E	駝 9A55	駝 9A4D	駝 9A5B	駝 9A57	駝 9A5F	駝 9A62	駝 9A65	駝 9A64	駝 9A69	駝 9A6B	駝 9A6A	駝 9AAD	駝 9AB0	駝 9ABC	駝 9AC0
9	體 9ACF	體 9AD1	體 9AD3	體 9AD4	體 9ADE	體 9ADF	體 9AE2	體 9AE3	體 9AE6	體 9AEF	體 9AEB	體 9AEE	體 9AF4	體 9AF1	體 9AF7	體 9AFB
A	體 9B06	體 9B18	體 9B1A	體 9B1F	體 9B22	體 9B23	體 9B25	體 9B27	體 9B28	體 9B29	體 9B2A	體 9B2E	體 9B2F	體 9B32	體 9B44	體 9B43
B	魏 9B4F	魏 9B4D	魏 9B4E	魏 9B51	魏 9B58	魏 9B74	魏 9B93	魏 9B93	魏 9B91	魏 9B96	魏 9B97	魏 9B9F	魏 9BA0	魏 9BA8	魏 9BB4	魏 9BC0
C	鯊 9BCA	鯊 9BB9	鯊 9BC6	鯊 9BCF	鯊 9BD1	鯊 9BD2	鯊 9BE3	鯊 9BE2	鯊 9BE4	鯊 9BD4	鯊 9BE1	鯊 9C3A	鯊 9BF2	鯊 9BF1	鯊 9BF0	鯊 9C15
D	鯊 9C14	鯊 9C09	鯊 9C13	鯊 9C0C	鯊 9C06	鯊 9C08	鯊 9C12	鯊 9C0A	鯊 9C04	鯊 9C2E	鯊 9C1B	鯊 9C25	鯊 9C24	鯊 9C21	鯊 9C30	鯊 9C47
E	鯊 9C32	鯊 9C46	鯊 9C3E	鯊 9C5A	鯊 9C60	鯊 9C67	鯊 9C76	鯊 9C78	鯊 9CE7	鯊 9CEC	鯊 9CF0	鯊 9D09	鯊 9D08	鯊 9CEB	鯊 9D03	鯊 9D06
F	鵠 9D2A	鵠 9D26	鵠 9DAF	鵠 9D23	鵠 9D1F	鵠 9D44	鵠 9D15	鵠 9D12	鵠 9D41	鵠 9D3F	鵠 9D3E	鵠 9D46	鵠 9D48			

EA	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
4	鵠 9D5D	鵠 9D5E	鵠 9D64	鵠 9D51	鵠 9D50	鵠 9D59	鵠 9D72	鵠 9D89	鵠 9D87	鵠 9DAB	鵠 9D6F	鵠 9D7A	鵠 9D9A	鵠 9DA4	鵠 9DA9	鵠 9DB2
5	鵠 9DC4	鵠 9DC1	鵠 9DBB	鵠 9DB8	鵠 9DBA	鵠 9DC6	鵠 9DCF	鵠 9DC2	鵠 9DD9	鵠 9DD3	鵠 9DF8	鵠 9DE6	鵠 9DED	鵠 9DE7	鵠 9DFD	鵠 9E1A
6	鵠 9E1B	鵠 9E1E	鵠 9E75	鵠 9E79	鵠 9E7D	鵠 9E81	鵠 9E88	鵠 9E8B	鵠 9E8C	鵠 9E92	鵠 9E95	鵠 9E91	鵠 9E9D	鵠 9EA5	鵠 9EA9	鵠 9EB8
7	鵠 9EAA	鵠 9EAD	鵠 9761	鵠 9ECC	鵠 9ECE	鵠 9ECF	鵠 9ED0	鵠 9ED4	鵠 9EDC	鵠 9EDE	鵠 9EDD	鵠 9EE0	鵠 9EE5	鵠 9EE8	鵠 9EEF	
8	鵠 9EF4	鵠 9EF6	鵠 9EF7	鵠 9EF9	鵠 9EFB	鵠 9EFC	鵠 9EFD	鵠 9F07	鵠 9F08	鵠 76B7	鵠 9F15	鵠 9F21	鵠 9F2C	鵠 9F3E	鵠 9F4A	鵠 9F52
9	鵠 9F54	鵠 9F63	鵠 9F5F	鵠 9F60	鵠 9F61	鵠 9F66	鵠 9F67	鵠 9F6C	鵠 9F6A	鵠 9F77	鵠 9F72	鵠 9F76	鵠 9F95	鵠 9F9C	鵠 9FA0	鵠 582F
A	鵠 69C7	鵠 9059	鵠 7464	鵠 51DC	鵠 7199											
B																
C																



---

## CP 932, continued


- This gives us an effective 16-bit space
  - Almost, sort of
  - It's really a partially-chained space
- But requires that we allow for 2-byte (and later, even 3-byte) sequences to describe one "character"
- You'd better know which CodePage you're dealing with for every string



---

## CP 932, continued

- BUT: still no way to combine multi-CP values in one string
- Unless you start inventing special "opcodes" to switch codepages in the middle of a byte stream
- No one thought that was a good idea
  - So that particular problem was ignored for a while



---

# LMBCS - Lotus MultiByte Character Set

- Lotus invented a new scheme as part of the work done for 123 Release 3
  - The first Lotus spreadsheet coded in C
  - Ran on both DOS and OS/2 (NOT Presentation Manager, which didn't exist yet)
  - Shipped in March, 1989



---

# LMBCS, continued

- Essentially also a "lead-byte" system, taking advantage of the fact that the ASCII codepoints 01 through 1F are non-printing characters
  - 00 is still special, particularly for C strings
- So, let's assign every glyph to a "code group", every code group has a unique identifier, a value between 1 and 1F
- Each code group corresponds (approximately) to a CodePage



---

# LMBCS, continued

- Code group 1 is Latin1
  - Very close to IBM CP850
- Each code group is either single byte, or multi byte
- Each code group's "lower half" is ASCII
- Code group 16 (Japanese) is multi byte, with designated "lead bytes", very similar to CP932
- And so on...



---

# LMBCS, continued

- The big LMBCS innovation:
  - Every character in "canonical" LMBCS is preceded by its code group identifier

"ABC" in hex format:  
01-41-01-42-01-43-0



---

## LMBCS, continued

- Characters in multi-byte code groups could therefore be 2, 3, or 4 bytes long
- 0 still used (in C) as a string terminator
- Guaranteed to be no embedded 0s in a string

Example of 2-byte Japanese character:  
10-<lead byte>-<2nd byte>-0





---

## LMBCS, continued

- But this is obviously flawed!
- All strings double in size!
- So we apply an "optimization"
- Each application file comes with a "default code group" designation
- Any character where the code group value is omitted is assumed to be in the default code group
- We detect code group identifiers because they are non-printing characters in a specified range
- Any character in a non-default code group must have its code group value specified



---

# LMBCS, continued

- This is known as "compressed", or "optimized" LMBCS

"ABC" in default code group 01:

41-42-43-0

"ABC" in default code group 16:

01-41-01-42-01-43-0

Hiragana small o in code group 16:

82-A7-0

Hiragana small o in code group 1:

10-82-A7-0



---

# Special Topics

- What about diacriticals?
- What about BIDI (bi-directional) character sets?
- What about "ligatures"?
- How to handle the huge Chinese/Japanese/Korean character namespaces?
- These are issues that all multi-language systems must handle properly
- (There are even more issues, but these are the common ones)



---

# Diacriticals

- We have a choice for dealing with accent marks
  - Create a code point and glyph for every combination

a á à

e é è



---

## Diacriticals - 2

- Or,
  - Invent the "non-spacing character"
- The "base" character has its own code point, each accent has its own
- Everyone knows not to move the cursor for the non-spacing character

a + ` = à

e + ' = é



---

## Diacriticals - 3

- You might want it either way
- Using non-spacing characters might make sorting easier (accent insensitive sorting)
- Or, might want each code point to be separate (accent sensitive sorting)
- LMBCS does not use non-spacing characters



---

# "BIDI"

- Some character sets are written right-to-left
  - Hebrew
  - Arabic
- But numbers are still left-to-right
- Some character sets are traditionally top-to-bottom, right-to-left
  - Japanese
  - Chinese



---

## BIDI - 2

- Mostly it's an issue for the rendering software
  - And for the input software
- But, how are the code points represented in memory?
  - Which order?
- The convention is that strings are represented in memory in "logical" order
  - I.E., in the order that you input them





---

## BIDI - 3

- So, if you have a RTL sequence of 3 characters
  - Let's represent as: "ZYX"
- Followed by 3 digits ("123")
- You would see those code points in memory as:
  - XYZ123
- The rendering software has to know to reverse the letters, not the digits



---

# Ligatures

- Historically, physical typefaces have combined certain individual letter combinations in to a single glyph
  - ff, ffl, fl, fi, ffi
  - Because the tip of the 'f' overhangs the next character
- Most computer fonts ignore this
- Some don't, they want to render as historically accurate
  - You need to know the character following the 'f' before you can render



---

## Ligatures - 2

- High-end publishing systems want to do this
- So room was made in some code pages for special glyphs representing the ligatures
- Makes sorting very complicated!



---

# Chinese, Japanese, Korean

- Present a special problem, because to fully represent the writing system would consume over 150,000 code points
- Even fully 16-bit systems only give you 65,000 code points
- Compromise: Consolidate
  - A large percentage of the Traditional Chinese, Japanese Kanji and Traditional Korean character sets are held in common
  - So we can represent the set of common glyphs only once
- Referred to as CJK Consolidation



---

# Limitations of LMBCS

- It's a very robust technique for representing multi-lingual text
  - Including within a single string!
  - In a reasonably compressed format (no embedded 0s)
- But...
  - The rest of the world adopted Unicode
- LMBCS remains a "proprietary" Lotus technology



---

# Where is LMBCS?

- LMBCS is used uniformly throughout Lotus products
  - SmartSuite
  - Notes
  - Domino
- And nowhere else
  - Major investment in supporting software tools
  - More on this later



---

# What Does Everyone Use Now?

- Unicode
  - Lotus also supports Unicode, as we'll see
- Unicode is an open standard
  - See <http://www.unicode.org>



---

# What is Unicode?

- Earlier character encoding schemes are often referred to as SBCS or DBCS or MBCS
  - Single Byte Character Set (e.g., Ascii)
  - Double Byte Character Set (e.g., CP932)
  - Multi-Byte Character Set
- Unicode is none of these





---

# What is Unicode?

- Unicode is not byte oriented, so calling it even MBCS is misleading
  - Though most of the time this is a true statement
  - There are times when you can treat it as an SBCS too
- My working definition:
  - A character encoding scheme that assigns a unique numeric value to every character
  - All written characters for all human languages (and many special characters as well) are accommodated



---

# What is Unicode?

- Note: this definition says nothing about how big the values might be
  - Or how many "bytes" might be needed to represent any given character
  - Because that's not the important issue
- The key point is that the "name space" for characters in Unicode is flat
  - 0 is still a special value
  - Upper value is essentially unbounded
  - And therefore very extensible



---

# What is Unicode?

- The current standard defines an upper limit of 0x10FFFF (using 3 8-bit bytes)
  - More than one million code points
  - Of which about 5% is allocated
  - 13% is reserved for private use
  - 2% reserved
  - 5% planning underway
  - Plenty of room for growth
- Of course, the devil is in the details



---

# Unicode Mappings

- The very large Unicode character space is laid out in regions
- Designed for interoperability with Ascii and ISO Latin-1
  - For ease of translation
- Ascii values (up to 0xff) are maintained unchanged
- Standards committee determines the rest of the values



---

# Unicode Encodings

- Question on computer representations:
  - 0x10ffff is 3 bytes, 24 bits
  - Most computers do not align on 3-byte boundaries
  - So we would "round up" to 4 bytes, 32 bits
- So, do we need to allocate 32 bits per character in all our strings?
  - And, if we do, won't there be embedded 0s?



---

## Encodings - 2

- Yes, there would
- But we don't necessarily need all 32 bits for each character.
- There are alternate Unicode "encodings":
  - UTF-8: each char is 1, 2, 3, or 4 units
  - UTF-16: each char is 1 or 2 units
  - UTF-32: each char is 1 unit
- UTF-16 is the "default" encoding



---

## Encodings - 3

- A char's real value does not change with the encoding
  - Only the number of bytes used to represent that value

"A" in UTF-8: 0x41 (byte)  
in UTF-16: 0x0041 (short)  
in UTF-32: 0x00000041 (DWORD)



---

# Encodings - 4

- Implication:
  - In UTF-8 and (to a lesser extent) in UTF-16
  - You may STILL need multiple units to represent a given code point
  - Therefore, Unicode can still be "multi-byte", in a sense
  - Meaning, "multiple positions"
- Zero still used by convention to terminate strings
  - Though this can vary by programming language
  - And the number of bytes must conform to the encoding (UTF-8, 16, 32)





---

# A Point on Diacriticals

- Unicode also gives you a choice on representation
  - One glyph for a "pre-composed" character
    - U-umlaut is 0x00fc
  - Base char + non-spacing char ("composed")
    - U + umlaut is 0x0075 + 0x0308
- So even in UTF-32 you may need 2 positions



---

# Referencing Characters in Unicode

- All characters have a unique name
  - Based on ISO 10646
  - E.g., "Bengali Digit 5"
- The convention for representing a character's code point value in Unicode
  - U+xxxx
  - Always assumes Hex representation



---

# Software for multi-language character sets

- The "standard" C library string functions are no longer usable
  - They all assume 1 byte == 1 char
  - strlen() tells you the length of a string in bytes, but not how many chars are in it
  - strcmp() is only valid for *some* SBCS, would never work with (e.g.) EBCDIC
  - Cannot add 0x20 to go from uc to lc
- You *must* use a software library that supports the character set in use



---

## Software for multi-language character sets - 2

- What kinds of things do you need to do?
  - Length of string (chars, bytes)
  - Search for char in string
  - Search for substring in string
  - Move cursor within string (by char):
    - First, Last
    - Next, Prev
  - Truncate at a certain byte position



---

# Software for multi-language character sets - 3

- Tasks, continued:
  - Test type of char (digit, alpha, upper, lower, accented, not accented...)
  - Modify casing (upper, lower)
  - Modify accenting (with, without)
  - Sort (compare 2 strings lexically)
  - Translate to/from native character set



---

# A Word on Sorting

- All lexical sorting requires is the ability to compare any 2 characters and say one is < the other
- BUT: it's a complex topic
- Easy sorting simply compares raw code point values
  - 'A' < 'a'
  - 'a' < 'b'
  - '1' < '9'
- But where are numbers relative to letters?
  - In Ascii, all digits are < all letters
  - But not universally true!



---

## Sorting - 2

- Need to handle comparison of characters using more than one byte
- Need to handle composed characters, too
- Real-world products need to offer multiple compare/sort options:
  - Case sensitive/insensitive
  - Accent sensitive/insensitive
  - Width sensitive/insensitive (for Asian double-wide characters)
  - Number first/numbers last



---

# Software for multi-language character sets - 4

- Topics:
  - Unicode and Java
  - MLCS and Fonts: Rendering text
  - LMBCS and Unicode in Notes





---

# Unicode and Java

- Java is all Unicode internally
  - Generally UTF-16 encoding
  - "char" is a 16-bit quantity in most cases



---

# Editing Java Code

- Be careful using non-Ascii characters as literals in Java code
  - 'a' is fine
  - Special characters ('\t', '\n') are ok
  - 'KK' (some Kanji character) is dangerous
- This is editor-dependent
  - If your text editor supports Unicode, then it's ok
  - Otherwise, your code may not be compiled on all systems



---

## Editing Java Code - 2

- Not an issue for Java code entered directly in Domino Designer
  - It's a LMBCS editor
  - Source code is stored in UTF-16
- This is not an issue for LotusScript or @Function code
  - Because you are always using a LMBCS editor



---

## Unicode and Java - 2

- Java lets you specify the code point value for a single char:
  - `'\uXXXX'`
  - Must supply 4 digits, implying UTF-16
- byte is an 8-bit quantity in Java
  - Be careful! Smaller than a char!
- int is a 32-bit quantity in Java
  - So cannot always freely convert
  - `c = (char) i;` works *most of the time*
  - `i = (int) c;` always works



---

# Unicode and Java - 3

- Character manipulation is built-in to String/StringBuffer classes
  - `char c = string.charAt(i)`
  - `toUpperCase/LowerCase()`
  - `toCharArray()`
  - `search, replace, etc., etc.`
  - `append()`



---

# Java Type Conversion

- Java will convert between String and char[] transparently
- Java will convert between String and byte[] with a specified encoding

```
byte[] arrOfBytes = null;  
String st = "Some string";  
arrOfBytes = st.getBytes("UTF-8");
```

```
String st2 = new String(arrOfBytes, "UTF-8");
```



# Java Type Conversion - 2

- Java will convert between String and char[] transparently
- Java will convert between String and byte[] with a specified encoding

```
byte[] arrOfBytes = null;  
String st = "Some string";  
arrOfBytes = st.getBytes("UTF-8");
```

```
String st2 = new String(arrOfBytes, "UTF-8");
```

NOTE: length of  
st2 (chars) may  
not = length of  
arr (bytes)!



---

## Java Type Conversion - 3

- Useful for converting byte-oriented strings (e.g., from C programs) to "real" ML strings
- Example: Notes converts all strings to UTF-8 when passing them to Java
  - But what if it's Kanji?

```
String utf8 = document.getItemValueString("abc");  
byte[] bytes = utf8.getBytes("UTF-16");  
String kanji = new String(bytes, "UTF-16");
```





---

# Java Supports Localization

- `java.util.Locale`
- Represents:
  - A "language code", e.g., "en", "de"
    - ISO 639
  - A "country code", e.g., "US", "CA"
    - ISO 3166
  - (Optional) A "variant", e.g., "posix", "MAC"
    - vendor specific



---

## Localization - 2

- A Locale instance implies something about the usage/formatting of:
  - dates (9/10/01 vs. 10/09/01 vs. 10.09.01)
  - currency symbols
  - accented characters
- So, these are all different:
  - en/US, en/UK, en/IR
  - fr/CA, fr/FR, fr/BE



---

## Localization - 3

- But: Locale is independent of display (font)
- Can query Java for the current, "default" locale
  - Dependent on system configuration
- Locale instances can be used to auto-format strings for display



---

# Java, LotusScript for File I/O

- If the file system is Unicode (e.g., Windows NT), then there are no issues
- Otherwise, you have 2 choices:
  - Write binary data in Unicode format
    - Not readable by other system tools
  - Convert going in and out
- LotusScript does Choice 2 automatically



---

# Java and File I/O

- With Java, you can control how you do it
- Package `java.io.*`
  - Base layer consists of `InputStream` and `OutputStream`
  - Byte oriented, you figure it out
  - Other layers support richer data types



---

# Java and File I/O

- What's wrong with this code?
  - Works, *sometimes...*

```
FileInputStream fis = new FileInputStream("d:\\temp\\abc.txt");  
char c; int i;  
while ((i = fis.read()) >= 0)  
{  
    c = (char) i;  
    // etc....  
}  
fis.close();
```



---

# Java and File I/O

- It works when:
  - File contents are Ascii
  - File contents are Unicode
- Otherwise, you can't assume that raw bytes are cast-able to Unicode
  - Correctly, that is
  - You are likely to get garbage



---

# Java and File I/O

- For "native" file i/o, use Reader/Writer classes
  - char oriented, not byte oriented

```
FileReader fr = new FileReader("d:\\temp\\abc.txt");
char c; int i;
while ((i = fr.read()) >= 0)
{
    c = (char) i;
    // etc....
}
fr.close();
```





---

# MLCS and Fonts: Rendering Text

- Rendering means:
  - Displaying a stream of code point values on a screen, or on paper
- Requires mapping a code point value to a particular glyph
  - Or a sequence of code point values to a single glyph
- Requires a displayable glyph for the code point value
  - With attributes too (bold, italic, superscript, etc.)



---

## Rendering - 2

- A "Font" (in computer terms) is a set of glyphs
  - Generally a table of displayable glyphs
  - Upper and lower case letters, digits, etc.
  - Kanji characters
- Can have multiple Fonts for any language character set
  - Just a different way of rendering each glyph
  - Times New Roman vs. Helvetica



---

## Rendering - 3

- Not all fonts have glyphs for the entire Unicode space
  - Most don't
- "Missing" glyphs default to a "fallback" character
  - The character data is still all there!
  - It's only the rendering software that can't handle the "missing" glyph



---

# Rendering - 4

- The character set that a font maps to is vendor-dependent
- Some are Ascii only
- Some are Asian only
  - Though most Asian fonts also include Ascii
- Others handle (some subset of) Unicode
- Rendering software is responsible for directionality
  - BIDI
  - Top-down



---

## Rendering - 5

- The amount of space a given string uses on the screen (or on paper) depends on:
  - Font characteristics (monospaced, proportional)
  - Attributes (bold, plain)
  - Point size (10, 12, 48...)
  - Pixel resolution of the display (DPI)
  - Actual chars used (esp. if proportional)
- Can be very complex to compute!



---

# Java & LMBCS in Notes


- All system display software requires "native" character set
- All Lotus product code bases use LMBCS as the character set
  - Internal manipulation (sorting, searching, etc.)
  - Persistent storage (on disk)
- With one exception:
  - LotusScript stores all strings as Unicode (UTF-16)
  - Must convert between:
    - LMBCS (for Notes)
    - Native (for display)



---

# Calling External Code From Java & LotusScript

- LotusScript calls externally using "Declare" statements
  - Describes library where code resides, name of entry point, and arguments
  - For String types, can declare the arg as:
    - "ByVal stringarg As *LMBCS* String"
    - Does conversion automatically
    - MUST be "ByVal"
- Otherwise, String is Unicode



---

# Calling External Code From Java & LotusScript - 2

- In Java, you can't invoke C entry points directly
  - Must use Java Native Interface (JNI)
  - C entry point receives a pointer to a Java services vector, plus arguments
    - Args are scalars or Java object handles
- You use the JVM services vector to manipulate the arguments
  - E.g., convert String objects to UTF-8 string buffers





---

# UI Support for Composed Characters

- In the Notes Client you can use Alt-F1 to compose chars
  - Alt-F1 + a + ` = à
  - Alt-F1 + c + , = ç
  - and so on



---

## Notes - 2

- Most conversions are transparent to the LS developer
  - Java Strings converted to Notes LMBCS
  - Notes LMBCS strings converted to Unicode (UTF-8)
  - LS strings converted to Notes LMBCS
  - Notes LMBCS converted to Unicode (UTF-16)



---

## Notes - 3

- When using Notes C API from LotusScript, need to worry about string type
  - Declare .... (....., st byval As LMBCS string)
  - Converts automatically to LMBCS
- When using custom C DLLs called from LotusScript
  - Either do same Declare trick
  - Or receive strings as UTF-16 Unicode



---

## Notes - 4

- For Java, when calling C:
  - String objects received as object handles
  - Use Java Native Interface (JNI) services to extract string value the way you want it
    - UTF-8, UTF-16
- Notes has many entry points in the C API for dealing with multi-language strings



---

# Multi-language APIs

- Categories:
  - Translation (OSTranslate, NLS\_Translate)
  - Parsing, searching (NLS\_XXX)
- All calls are documented in the Notes C API toolkit
- NOTE: All "char \*" decls in the C API mean "unsigned char"
  - LMBCS



---

# OSTranslate

- The key is the "TranslateMode" options
  - Note: result-length may vary from input-length!

```
WORD LNPUBLIC OSTranslate(  
    WORD TranslateMode,  
    char far *In,  
    WORD InLength,  
    char far *Out,  
    WORD OutLength);
```



---

## OSTranslate - 2

- Assumes that either source or destination string (or both) is LMBCS Group 1
  - Notes assumes Group 1 strings are optimized
    - Leading group indicator-bytes compressed out
  - All others non-optimized



---

# OSTranslate - 3

- OS\_TRANSLATE\_NATIVE\_TO\_LMBCS - Convert input string from machine's native character set to LMBCS
- OS\_TRANSLATE\_LMBCS\_TO\_NATIVE - Convert input string from LMBCS to machine's native character set.
- OS\_TRANSLATE\_LOWER\_TO\_UPPER - Current international case table.
- OS\_TRANSLATE\_UPPER\_TO\_LOWER - Current international case table.
- OS\_TRANSLATE\_UNACCENT - International unaccenting table. NO REVERSAL!!
- OS\_TRANSLATE\_OSNATIVE\_TO\_LMBCS - Same as NATIVE (not documented!)
- OS\_TRANSLATE\_LMBCS\_TO\_OSNATIVE - Same as TO\_NATIVE (not documented!)
- OS\_TRANSLATE\_LMBCS\_TO\_ASCII - Convert the input string from LMBCS to character text.
- OS\_TRANSLATE\_LMBCS\_TO\_UNICODE - Convert the input string from LMBCS to UNICODE.
- OS\_TRANSLATE\_LMBCS\_TO\_UTF8 - Convert the input string from LMBCS to UTF8.
- OS\_TRANSLATE\_UNICODE\_TO\_LMBCS - Convert the input string from UNICODE to LMBCS.
- OS\_TRANSLATE\_UTF8\_TO\_LMBCS - Convert the input string from UTF8 to LMBCS.





---

## OSTranslate - 4

- Note that multiple translations are "lossless"
  - If you go one way, then come back, all data is preserved
- EXCEPT:
  - Accented to unaccented
  - LMBCS to Ascii (maybe)
- Complicates display input and editing



---

# National Language System (NLS) APIs

- Supports low-level character based parsing, searching and translation
- Any character set, not just LMBCS/1 or current native



---

# NLS\_PINFO

- Use NLS\_load\_charset to read a charset descriptor table from disk
- Returns a NLS\_PINFO \*
  - Opaque data structure, but you will need it
- List of character set IDs found in NLS.H
- Used in NLS\_translate



---

# NLS-translate

- More functionality than in OSTranslate
  - Arbitrary translation to/from any character set
  - NLS\_translate

```
NLS_STATUS WINAPI NLS_translate(  
    BYTE far *pString,  
    WORD Len,  
    BYTE far *pStringTarget,  
    WORD far *pSize,  
    WORD ControlFlags,  
    NLS_PINFO pInfo);
```



---

# NLS Translation Options

- NLS\_NONULLTERMINATE - Does not add a NULL to the end of the translated result.
- NLS\_NULLTERMINATE - Adds a NULL to the end of the translated result.
- NLS\_STRIPUNKNOWN - Strips unknown characters.
- NLS\_TARGETISLMBCS - Converts target string to LMBCS.
- NLS\_SOURCEISLMBCS - Converts source string from LMBCS.
- NLS\_TARGETISUNICODE - Converts target string to UNICODE.
- NLS\_SOURCEISUNICODE - Converts source string from UNICODE.
- NLS\_TARGETISPLATFORM - Converts target string to the encoding used by the host OS.
- NLS\_SOURCEISPLATFORM - Converts source string from the encoding used by the host OS.



---

## Other NLS API Calls

- NLS\_string\_chars
  - Number of chars in a string
- NLS\_string\_bytes
  - Number of bytes in a string (same as strlen)
- Test a character for type:
  - isdigit, isupper, islower, isspace, ispunct, iscntrl, isalnum, isalpha, isarith,



---

## Other NLS API Calls - 2

- isleadbyte
  - If True, then the next 1 (or 2, or 3) bytes are part of the same character
- Enumerate through a string
  - NLS\_get
- Parse strings
  - goto\_next\_break
  - goto\_next/prev\_word\_end/start



---

## Other NLS API Calls - 3

- Search
  - find, find\_substr





---

# String Truncation


- String truncation: don't ignore this!
  - Failure to do this properly can cause crashes
- Situation:
  - You have an N-byte long MLCS string
  - You must truncate to M bytes
    - $M < N$
    - For display, to fit in a fixed buffer, whatever



---

# String Truncation - 2

- Can't just chop it off at M bytes!
  - Might be the middle of a multi-byte sequence
  - Will result (at least) in garbage on the display
- Correct technique:
  - Get pointer to max slot in string
  - Call NLS\_goto\_prev
    - Pass in \*\* for current position, \* for start of string
  - Moves your pointer to the start of the previous char
  - Can truncate there, or at start of next char



---

# MLCS APIs - Summary

- All the old C techniques are now useless
  - Worse: wrong!
- You will need to accustom yourself to a new set of assumptions and functions
- But once you get comfortable with it, it's not difficult
  - Make it habitual!
- NLS calls somewhat limited:
  - Only a few of the string-compare and search equivalents for MLCS are exposed!
  - Testers for full/half pitch not exposed!



---

# Web applications and MLCS

- What happens Notes data served to the Web?
- HTTP allows specification of the page charset
  - Client can specify desired charsets on request
  - Server specifies charset on result page



---

## Web Applications - 2

- Notes defaults to ISO Latin-1 for output
  - Automatically converts from LMBCS
- But, how do you know what the Client's real charset is?
  - Look at requested charsets
  - Look at country/language code
- Can you tell Domino what charset to use on output?



---

## Web Applications - 3

- Yes!
- Domino examines Web agent output for HTTP headers
  - Must occur BEFORE any data
  - Domino will notice them and pass them through
  - If you specify a Content-Type header, Domino will try to translate



---

## Web Applications - 4

- Just print out the header (in the correct format) as the first thing:

```
print "Content-Type: text/html; charset=ISO-8859-4"
print ""      ' null to force a newline: REQUIRED
print "Yo, world!"
```



---

## Web Applications - 5

- Unfortunately, only works from agents
- No way (that I know of) to build this into a form or page





---

# Web applications and MLCS

- Another trick for using in browsers
- HTTP specifies that Unicode values can be substituted for any char
  - 'Space' == U+0020
  - Can encode for HTTP as
    - %0020
    - Or %20
- So you can encode arbitrary Unicode chars in your Notes form or page!
  - The browser will interpret it correctly
  - Very useful when encoding pass-thru URLs



---

# Summary

- If everyone spoke (and wrote!) only one language, it would be easy
- Remember: character set encodings are independent of rendering
  - Though typographical legacies complicate encodings (e.g., ligatures)
- Unicode is an important technological advance
  - Everything defined in a single "space"



---

## Summary - 2

- But the tools are not yet sophisticated enough
  - Developers still need some knowledge of how character sets work
  - And Unicode implementations are not always uniform
    - UTF-8? 16? 32?
    - Java UTF-16 still has multi-char code points



---

## Summary - 3

- It will keep getting better
  - A good thing
- It will keep changing
  - Not always so good



---

**Click here to type page title**

**Q & A**