

Table des matières

Avant-propos.....	XVII
-------------------	------

Première partie – Introduction

Chapitre 1 – Concepts de base et terminologie	3
1.1 Pourquoi Unicode ?	3
1.2 Absence de codage universel	5
1.3 Langue et écriture.....	6
1.3.1 <i>Notation, écriture, transcription et translittération.....</i>	7
1.3.2 <i>Qu'est-ce qu'un caractère ?</i>	8
1.3.3 <i>Graphème, caractère et glyphe</i>	9
1.3.4 <i>Terminologie de l'apparence des caractères.....</i>	11
1.3.5 <i>Classification des caractères</i>	11
1.3.6 <i>Caractères romains, latins, italiques et gothiques</i>	12
1.3.7 <i>Écriture CJC</i>	13
1.4 Unicode, en quelques mots	14
1.4.1 <i>Ce qu'Unicode est... ..</i>	14
1.4.2 <i>Ce qu'Unicode n'est pas.....</i>	14
1.5 Apprivoiser les polices Unicode.....	16
1.5.1 <i>Afficher des caractères Unicode</i>	16
1.5.2 <i>S'assurer que sa police est une police Unicode</i>	17
1.5.3 <i>Où trouver des polices multi-écritures supplémentaires ?</i>	18
1.5.4 <i>Absence de glyphe pour un caractère</i>	20
1.5.5 <i>Incorporation des polices.....</i>	21
1.6 Saisir des caractères Unicode.....	22
1.6.1 <i>Claviers</i>	25

1.6.2	Méthodes d'entrée extrême-orientales	28
1.6.3	Tableau de caractères	29
1.7	Internationalisation et localisation	32
Chapitre 2 – Répertoires et jeux de caractères codés.		35
2.1	Terminologie.	35
2.1.1	Répertoire de caractères	36
2.1.2	Jeu de caractères codés, code et codage	36
2.2	L'ASCII.	37
2.2.1	ISO 646 et « variantes nationales de l'ASCII »	39
2.2.2	« ASCII 8 bits »	39
2.3	Les codes ISO/CEI 8859.	40
2.3.1	ISO/CEI 8859-1 ou ISO Latin-1.	41
2.3.2	ISO/CEI 8859-15 ou ISO Latin-9.	43
2.4	Windows Latin 1	44
2.5	EBCDIC	46
2.6	KOI8-R.	47
2.7	ISO/CEI 2022.	47
2.8	ISO/CEI 10646 et Unicode	49
2.9	GB 18030, l'« Unicode chinois »	51

Deuxième partie – L'essentiel d'Unicode

Chapitre 3 – Structure d'Unicode.		55
3.1	Principes directeurs.	55
3.2	Caractères normalisés.	56
3.2.1	Plan multilingue de base	57
3.2.2	Plans complémentaires	61
3.2.3	Nombre de caractères normalisés	61
3.2.4	Unicode 5.1.	62
3.3	Caractères	64
3.3.1	Caractère abstrait et caractère codé.	64
3.3.2	Caractères combinatoires et diacritiques	65
3.3.3	Suite de caractères de base et diacritiques.	66
3.3.4	Caractères combinatoires multiples	66
3.4	Principes de conception d'Unicode	67
3.4.1	Universalité	68

3.4.2	<i>Efficacité</i>	69
3.4.3	<i>Caractères et non glyphes</i>	70
3.4.4	<i>Sémantique</i>	70
3.4.5	<i>Texte brut</i>	71
3.4.6	<i>Ordre logique</i>	71
3.4.7	<i>Unification</i>	72
3.4.8	<i>Composition dynamique</i>	72
3.4.9	<i>Stabilité</i>	74
3.4.10	<i>Convertibilité</i>	75
Chapitre 4 – Modèle de codage, propriétés des caractères et tri		77
4.1	<i>Modèle de codage des caractères</i>	77
4.1.1	<i>Répertoire de caractères abstraits</i>	80
4.1.2	<i>Jeu de caractères codés</i>	80
4.1.3	<i>Mot, octet, seizet, codet</i>	80
4.1.4	<i>Forme en mémoire des caractères</i>	81
4.1.5	<i>Mécanisme de sérialisation de caractères</i>	84
4.2	<i>Propriétés des caractères</i>	86
4.2.1	<i>Voir les propriétés grâce à BabelMap</i>	86
4.2.2	<i>Catégorie générale</i>	88
4.2.3	<i>Casse</i>	91
4.2.4	<i>Classe combinatoire canonique</i>	94
4.2.5	<i>Directionnalité</i>	95
4.2.6	<i>Réflexion bidi</i>	96
4.2.7	<i>Classes de coupure de lignes</i>	96
4.2.8	<i>Coupure de lignes et segmentation de texte</i>	99
4.2.9	<i>Cohérence des propriétés</i>	100
4.3	<i>Formes normalisées</i>	101
4.3.1	<i>La concaténation n'est pas fermée</i>	103
4.3.2	<i>Formes de normalisation et base de données</i>	103
4.3.3	<i>Stabilité des formes normalisées</i>	103
4.3.4	<i>Invariance des caractères latins de base</i>	104
4.3.5	<i>Compositions exclues</i>	104
4.4	<i>Le tri et le repérage</i>	105
4.4.1	<i>Tri et comparaison binaires</i>	106
4.4.2	<i>Tri et comparaison lexicographiques</i>	106
4.4.3	<i>Ça ne peut quand même pas être si compliqué ?</i>	106
4.4.4	<i>La solution – un tri à quatre niveaux</i>	107

- 4.4.5 *Les éléments de tri et les clés de tri* 109
- 4.4.6 *La table DUCET* 110
- 4.4.7 *Expansions et contractions* 111
- 4.4.8 *Mise en œuvre du tri Unicode* 112
- 4.4.9 *Personnalisation* 113
- 4.5 *Conformité* 113
- 4.6 *Le standard Unicode : mode d'emploi* 115
 - 4.6.1 *Liste des noms de caractère* 115
 - 4.6.2 *Images dans les tableaux et dans les listes de caractères* 116
 - 4.6.3 *Renvois* 116
 - 4.6.4 *Renseignements sur les langues* 117
 - 4.6.5 *Décompositions* 117

Troisième partie – Caractères remarquables

- Chapitre 5 – Lettres et signes diacritiques** 123
 - 5.1 *Latin étendu et API* 123
 - 5.2 *Lettres modificatives* 124
 - 5.3 *Clones à chasse des diacritiques* 125
 - 5.4 *Écriture grecque* 125
 - 5.5 *Signes diacritiques* 127
 - 5.5.1 *Diacritiques généraux* 128
 - 5.5.2 *Supplément de diacritiques* 129
 - 5.5.3 *Diacritiques destinés aux symboles* 129
 - 5.5.4 *Demi-signes diacritiques* 130
- Chapitre 6 – Ponctuation** 131
 - 6.1 *Ponctuation Latin-1* 131
 - 6.1.1 *Guillemet anglais* 132
 - 6.1.2 *Croisillon* 132
 - 6.1.3 *Perluète* 132
 - 6.1.4 *Apostrophe* 132
 - 6.1.5 *Astérisque et obèle* 133
 - 6.1.6 *Trait d'union-signe moins* 134
 - 6.1.7 *Arrobe* 134
 - 6.1.8 *Clones de diacritique* 134
 - 6.1.9 *Paragraphe et pied-de-mouche* 135
 - 6.1.10 *Symbole degré et ordinal masculin* 136

6.1.11	<i>Point médian</i>	137
6.1.12	<i>Trait d'union conditionnel</i>	138
6.1.13	<i>Ponctuation appariée</i>	140
6.2	<i>Guillemets</i>	140
6.2.1	<i>Usages européens</i>	140
6.2.2	<i>Usage extrême-oriental</i>	142
6.3	<i>Espaces</i>	142
6.3.1	<i>Espace mot et espace insécable</i>	143
6.3.2	<i>Les différents caractères d'espacement</i>	144
6.3.3	<i>Ajustement de l'espacement</i>	146
6.3.4	<i>Espaces fines en français</i>	146
6.3.5	<i>Espace sans chasse</i>	149
6.3.6	<i>Disposition des espaces</i>	149
6.4	<i>Autres signes typographiques</i>	151
6.4.1	<i>Points de suspension et points de conduite</i>	151
6.4.2	<i>Traits d'union et tirets</i>	152
6.4.3	<i>Puces, barre de fraction, ponctuation doublée</i>	154
6.4.4	<i>Ponctuation archaïque</i>	154
6.5	<i>Caractères de coupure de lignes</i>	155
Chapitre 7 – Symboles et notations		157
7.1	<i>Symboles de type lettre</i>	157
7.2	<i>Symboles monétaires</i>	159
7.3	<i>Mathématique</i>	159
7.4	<i>Musique</i>	161
Chapitre 8 – Caractères techniques spéciaux		163
8.1	<i>Caractères de commande</i>	163
8.1.1	<i>Commandes C0 et suppression</i>	163
8.1.2	<i>Commandes C1</i>	164
8.2	<i>Gluon et diacritique invisible bloquant</i>	165
8.2.1	<i>Gluon de mots (U+2060)</i>	165
8.2.2	<i>Diacritique invisible bloquant (U+034F)</i>	165
8.3	<i>Caractères spéciaux</i>	166
8.3.1	<i>Délimiteurs d'annotation interlinéaire</i>	167
8.3.2	<i>Non-caractères</i>	168
8.4	<i>Positions non attribuées</i>	168

8.5	Caractères déconseillés et désuets	169
8.6	Zones à usage privé	169
8.7	Indicateur d'ordre des octets	170
8.8	Étiquettes linguistiques	171

Quatrième partie – Applications et techniques liées à Unicode

Chapitre 9 – Préciser la langue, l'écriture et le pays	177	
9.1	ISO 639 – indicatifs de langue	178
9.2	ISO 3166 – indicatifs de pays	180
9.3	M.49 – Indicatifs de pays et de régions	181
9.4	ISO 15924 – indicatifs d'écriture	181
9.5	RFC 4646 – Étiquettes de langue	182
9.7	BCP 47	184
Chapitre 10 – Unicode et les protocoles Internet	187	
10.1	De l'utilité des métadonnées	187
10.2	Les premiers protocoles Internet	188
10.3	Type de médias Internet/type MIME	189
10.3.1	Visualiser les entêtes	189
10.3.2	Les types de média ou types de contenu	191
10.3.3	L'information sur le codage de caractères (« charset »)	192
10.3.4	Les entêtes relatifs au surcodage de transfert	193
10.3.5	Le surcodage des entêtes	197
10.3.6	Recettes de dépannage	198
10.4	Codage de caractères sur le Web	200
10.4.1	Entêtes HTTP	200
10.4.2	Préciser le codage de caractères dans HTTP	201
10.4.3	Vérifier les entêtes HTTP	202
10.4.4	Quel codage utiliser pour mes pages web ?	204
10.4.5	Balise meta	204
10.4.6	Préciser le codage dans XHTML et XML	205
10.4.7	Conflits des définitions de codage	206
10.4.8	Configuration de serveurs web	207
10.5	HTTP internationalisé	210
10.5.1	La négociation de langue	210
10.5.2	Les entêtes reliés aux caractères	219

10.6	Adresses internationalisées.	220
10.6.1	Noms de domaine et DNS	220
10.6.2	Internationaliser les URI.	222
10.6.3	Noms de domaine internationalisés (NDI)	224
10.6.4	Menaces informatiques : hameçonnage et parodie	226
10.6.5	Caractères non ASCII dans les chemins des IRI	229
10.7	La locale POSIX.	232
Chapitre 11	– Unicode et (X)HTML, XML, CSS.	239
11.1	Préciser le codage en (X)HTML	239
11.1.1	HTML	240
11.1.2	XHTML – le prologue XML	240
11.1.3	Préciser le codage en CSS.	244
11.2	Préciser la langue	245
11.2.1	HTML et XML.	245
11.2.2	Passages dans une autre langue.	245
11.2.3	Documents bilingues.	246
11.2.4	La langue dans l'entête HTTP ou l'attribut lang ?	247
11.3	Préciser la directionnalité.	248
11.4	Stylage sensible à la langue.	249
11.4.1	Les sélecteurs CSS	249
11.4.2	Utilisation des sélecteurs de langue	250
11.5	Schémas XML internationalisés.	251
11.5.1	Texte dans les attributs XML	251
11.5.2	Les éléments qui ne contiennent qu'une chaîne	252
11.5.3	Prévoir l'attribut xml:lang.	252
11.5.4	Prévoir un élément de type span	252
11.5.5	Ne pas créer d'éléments de présentation.	252
11.5.6	Prévoir xml:id sur tous les éléments traduisibles	253
11.6	Notation des caractères	253
11.6.1	Appels d'entités HTML : souvent de peu d'utilité	255
11.6.2	Appels de caractère : à n'utiliser que rarement.	256
11.6.3	Quand les appels de caractère et d'entité sont utiles	256
11.6.4	Entités en XHTML	258
11.7	Caractère ou balisage ?	258
11.7.1	Unicode contient trop de caractères.	258
11.7.2	Caractères de commande en HTML et XHTML	259

11.7.3	<i>Autres caractères permis et interdits en XML</i>	260
11.7.4	<i>Caractères incompatibles avec le balisage</i>	261
11.7.5	<i>Caractères de compatibilité Unicode</i>	263
11.8	<i>Réglage de l'algorithme bidi</i>	263
11.8.1	<i>Bref rappel de l'algorithme bidi</i>	264
11.8.2	<i>Désactiver l'algorithme</i>	268
11.8.3	<i>Cas problématiques – les neutres mal placés</i>	269
11.8.4	<i>Éditer du texte bidi balisé</i>	270
11.8.5	<i>Emboîtement des passages bidi</i>	271
11.8.6	<i>Caractères de commande bidi et balisage</i>	273
11.8.7	<i>Les feuilles de styles CSS et le bidi</i>	273
11.9	<i>Formulaire « universel »</i>	275
11.9.1	<i>Créer un formulaire</i>	275
11.9.2	<i>Accept-charset sur la balise form</i>	278
11.9.3	<i>Inclusion de caractères étrangers au charset</i>	279
11.9.4	<i>Solution : n'envoyer et n'accepter que de l'UTF-8</i>	279
Chapitre 12	– Internationalisation des logiciels	283
12.1	<i>Internationaliser ?</i>	283
12.1.1	<i>PNB par langue</i>	283
12.1.2	<i>L'adaptation de logiciels à la pièce</i>	284
12.1.3	<i>Internationalisation</i>	284
12.1.4	<i>Quelques a priori culturels</i>	285
12.1.5	<i>Localisation ou adaptation culturelle</i>	286
12.2	<i>Les langages de programmation</i>	287
12.2.1	<i>Les caractères dans C et C++</i>	287
12.2.2	<i>Les caractères dans Java</i>	289
12.2.3	<i>Les caractères dans C# et la plateforme .NET</i>	291
12.2.4	<i>Bibliothèques d'internationalisation, le cas ICU</i>	292
12.2.5	<i>Propriétés de caractères en Java/ICU</i>	293
12.2.6	<i>Unicode dans les expressions régulières</i>	295
12.2.7	<i>UnicodeSet en ICU</i>	296
12.2.8	<i>Normalisation</i>	297
12.2.9	<i>Comment supprimer les accents d'un texte ?</i>	298
12.2.10	<i>Créer des noms de domaine internationalisés</i>	299
12.2.11	<i>Les transformations ICU</i>	299
12.3	<i>Profil culturel ou locale</i>	302
12.3.1	<i>Concept et nécessité</i>	302

12.3.2	<i>Définir une Locale</i>	303
12.3.3	<i>Définir une Locale avec ICU</i>	303
12.3.4	<i>Locale sur .NET</i>	304
12.3.5	<i>Locale implicite.</i>	305
12.3.6	<i>Changer la casse d'une chaîne de caractères.</i>	305
12.3.7	<i>Correspondance de casse non localisée.</i>	306
12.4	<i>Isoler les données culturelles</i>	307
12.4.1	<i>Pourquoi extraire ?</i>	307
12.4.2	<i>Les « ResourceBundle »</i>	307
12.4.3	<i>Définition de ResourceBundle</i>	308
12.4.4	<i>Accès aux ResourceBundle</i>	309
12.4.5	<i>Découverte et recherche des ResourceBundle</i>	309
12.4.6	<i>Où sont stockées les ressources ?</i>	310
12.5	<i>Formater les messages</i>	312
12.5.1	<i>Formats prédéfinis de date et heure</i>	312
12.5.2	<i>Formater la date et l'heure à l'aide de motifs</i>	313
12.5.3	<i>Formater et analyser chiffres et montants</i>	316
12.5.4	<i>Variabilité de l'ordre des mots dans les langues.</i>	318
12.5.5	<i>Messages variables et internationalisés.</i>	319
12.5.6	<i>Gestion de l'accord en nombre</i>	322
12.6	<i>Comparaison et tri</i>	325
12.6.1	<i>Tri en Java ou ICU pour Java.</i>	325
12.6.2	<i>Comparer des chaînes en ignorant les accents</i>	327
12.6.3	<i>Personnaliser le tri</i>	327
12.6.4	<i>Ignorer la ponctuation.</i>	330
12.6.5	<i>Tri des tableaux de données à plusieurs champs</i>	331
12.6.6	<i>Améliorer la performance</i>	333
12.7	<i>Frontières de texte</i>	334
12.7.1	<i>La classe BreakIterator</i>	334
12.7.2	<i>Exemples : détecter les frontières de phrase et de mot</i>	336
12.7.3	<i>Personnaliser la détection de frontières.</i>	338
12.8	<i>CLDR</i>	339
12.9	<i>Les exceptions et l'i18n</i>	340
12.10	<i>Conversion de données</i>	342
12.10.1	<i>Entrées/sorties</i>	342
12.10.2	<i>Écrire de l'Unicode avec des OutputStream.</i>	344

12.10.3 <i>String.getBytes()</i> et l'analyseur XML	345
12.10.4 Maîtrise de la conversion grâce à <i>java.nio</i>	346
12.11 L'interface utilisateur	347
12.11.1 Foisonnement du texte	347
12.11.2 Neutralité culturelle	348
12.11.3 Internationalisation des images et des icônes	349
12.12 À ne pas internationaliser	350
Chapitre 13 – Unicode et les polices.	355
13.1 Caractères et variantes de glyphes	355
13.2 Sélecteurs de variante.	356
13.3 Impact sur le choix de police	356
13.3.1 <i>Police de repli</i>	357
13.3.2 <i>Police de substitution</i>	357
13.3.3 <i>Police liée</i>	358
13.3.4 CSS et ses « <i>polices liées</i> »	358
13.3.5 <i>Utiliser un équivalent canonique</i>	359
13.4 Ligatures	359
13.4.1 <i>Liant et antiliant</i>	360
13.4.2 <i>Liaison cursive</i>	362
13.4.3 <i>Liant, antiliant et les écritures brahmies</i>	363
13.4.4 <i>Filtrage des liants et antiliants</i>	364
13.4.5 <i>Liant et antiliant dans les polices</i>	364
13.5 Pas d'expédients ASCII, de l'Unicode !.	365
13.6 Passage des caractères aux glyphes	366
13.6.1 <i>Des caractères et non des glyphes</i>	366
13.6.2 <i>Fonctionnement d'un moteur de rendu</i>	367
13.6.3 <i>Les polices de glyphes</i>	371
13.7 Processus de rendu	378
13.8 Un moteur de rendu : Uniscribe	380
13.9 Adobe et Uniscribe	382
13.10 Fonctionnalités et règles OpenType.	382
13.11 Intégrer le tout	385
Bibliographie	389
Index	391