

Το πρότυπο Unicode

Απόστολος Συρόπουλος

28ης Οκτωβρίου 366

671 00 Ξάνθη

E-mail: apostolo@obelix.ee.duth.gr

Δημήτριος Α. Φιλίππου

Κάτω Γατζέα

385 00 Βόλος

Ιωάννης Κ. Δημάκος

Παν/μιο Πατρών

E-mail: idiimakos@upatras.gr

1. Εισαγωγή

Ένας απλός τρόπος για να κρίνει κανείς κατά πόσο είναι σοβαρό ή μη ένα σύστημα στοιχειοθεσίας είναι να ελέγξει τις δυνατότητες στοιχειοθεσίας μαθηματικού και πολυγλωσσικού κειμένου με τη χρήση του συστήματος αυτού. Μπορούμε εύκολα να καταλάβουμε γιατί το μαθηματικό κείμενο είναι ιδιαίτερα απαιτητικό αν σκεφτούμε μόνο ότι οι διάφοροι μαθηματικοί τύποι (εξισώσεις, κ.λπ.) είναι στην πραγματικότητα διαδιάστατο κείμενο. Από την άλλη, η σύνταξη ενός πολύγλωσσου κειμένου είναι δύσκολη γιατί πρέπει να υποστηρίζεται όχι μόνο ένας μεγάλος αριθμός γραμματοσειρών, αλλά και οι τυπογραφικές ιδιαιτερότητες της κάθε γλώσσας.

Τόσο το μαθηματικό κείμενο όσο και η υλοποίηση των τυπογραφικών ιδιαιτεροτήτων είναι θέματα στα οποία το \TeX και οι επίγονοί του είναι πραγματικά πρωταθλητές. Όμως, το \TeX είναι «γνήσιο Αμερικανάκι» και, σαν τέτοιο, βασίζεται στο αμερικανικό πρότυπο ASCII. Σύμφωνα με το πρότυπο ASCII, κάθε χαρακτήρας αντιστοιχεί σε έναν θετικό ακέραιο αριθμό του οποίου το μέγεθος δεν μπορεί να ξεπεράσει τα οκτώ bit (ένα byte).¹ Συνεπώς, κατά το πρότυπο ASCII, δεν μπορούμε να γράψουμε κάτι στον υπολογιστή σε μια γλώσσα που περιέχει περισσότερα από $2^8 = 256$ γράμματα και σύμβολα.

¹ Στην αρχική του μορφή, το ASCII περιελάμβανε μόνο «χαρακτήρες μεγέθους 7 bit», δηλαδή σε $2^7 = 128$ χαρακτήρες. Κατόπιν διάφοροι κατασκευαστές Η/Υ και λειτουργικών συστημάτων επέκτειναν το ASCII στα 8 bit (1 byte) — ο καθένας σύμφωνα με τον δικό του τρόπο.

Για το T_EX, η λύση που επινοήθηκε για την αντιμετώπιση του προβλήματος της στοιχειοθεσίας κειμένων σε γλώσσες που έχουν περισσότερα από 256 γράμματα και σύμβολα ήταν η μεταγραφή, δηλαδή η γραφή ενός μη λατινικού αλφαβήτου χρησιμοποιώντας το λατινικό αλφάβητο στον κώδικα T_EX. Το μειονέκτημα αυτής της μεθόδου είναι ότι το αρχείο με τον κώδικα T_EX, γίνεται δυσανάγνωστο ακόμα και για τους μυημένους. Προβλήματα ακριβώς όπως αυτό, και κυρίως η ύπαρξη κωδικοποιήσεων (κωδικοσελίδων) για την ίδια γλώσσα διαφορετικών για κάθε λειτουργικό σύστημα, ήταν που οδήγησαν στη δημιουργία του προτύπου Unicode.

Με απλά λόγια, το πρότυπο Unicode είναι ένα μεγάλο σύνολο χαρακτήρων το οποίο υποστηρίζει τις περισσότερες από τις γραφές του κόσμου (νεκρές ή και χρησιμοποιούμενες) επιτρέποντας έτσι την εύκολη μετάδοση, επεξεργασία και προβολή πολυγλωσσικών κειμένων. Σε ό,τι ακολουθεί θα παρουσιάσουμε σε συντομία τις τεχνικές λεπτομέρειες που αφορούν το πρότυπο, τις ομάδες χαρακτήρων που υποστηρίζει αλλά και την ελληνική υποστήριξη.

2. Περιγραφή του προτύπου Unicode

Το πρότυπο Unicode δημιουργήθηκε με βάση μια σειρά από αρχές. Κύρια αρχή αποτέλεσε η απαίτηση ο κάθε χαρακτήρας να αναπαριστάται από 16 bit (δηλαδή δύο byte). Ο κάθε χαρακτήρας αντιστοιχεί σε ένα κωδικό σημείο το οποίο αναπαριστά τη θέση του χαρακτήρα στο πρότυπο και μπορεί να χρησιμοποιηθεί για συγκρίσεις. Το κωδικό σημείο κάθε χαρακτήρα σημειώνεται με ένα τετραψήφιο δεκαεξαδικό αριθμό, π.χ., AF01, 0308, κ.ο.κ. Επιπλέον, τα αρχεία Unicode θα έπρεπε να είναι εύκολα προσπελάσιμα, ενώ η επεξεργασία τους δεν θα έπρεπε να παρουσιάζει ιδιαίτερες δυσκολίες. Επειδή το Unicode είναι ουσιαστικά ένα (μεγάλο) σύνολο χαρακτήρων, ακριβώς όπως είναι και το ASCII, ο αναγνώστης θα πρέπει να έχει στο μυαλό του ότι το πρότυπο Unicode καθορίζει τον αριθμό στον οποίο αντιστοιχεί ένας χαρακτήρας η ένας γλύφος (δηλαδή ένας χαρακτήρας που προκύπτει από τον συνδυασμό δύο ή περισσότερων απλών χαρακτήρων), αλλά δεν καθορίζει διαφορετικές μορφές (σχεδιάσεις) χαρακτήρων. Για παράδειγμα, στο παρακάτω σχήμα δίνουμε τέσσερις διαφορετικές μορφές (σχεδιάσεις) του γλύφου LATIN SMALL LIGATURE FI, που αντιστοιχεί στον κωδικό αριθμό Unicode FB01:

fi fi fi fi

Επιπλέον, προς αποφυγή προβλημάτων ανάμειξης διαφορετικών γλωσσών, αρκετοί χαρακτήρες που μοιάζουν εμφανίζονται πολλές φορές. Μια τέτοια περίπτωση είναι ο χαρακτήρας A που απαντάται τουλάχιστον στις εξής μορφές: LATIN CAP-

ITAL LETTER A, GREEK CAPITAL LETTER ALPHA και CYRILLIC CAPITAL LETTER A.

Το πρότυπο Unicode έχει σχεδιασθεί έτσι ώστε ο κάθε χαρακτήρας να αντιστοιχεί σε μία σειρά 16 bit. Αυτή η σύμβαση οδηγεί στο συμπέρασμα ότι με το Unicode μπορούμε να αναπαραστήσουμε μέχρι $2^{16} = 65.536$ διαφορετικούς χαρακτήρες. Όμως στην πραγματικότητα μπορούμε άμεσα να απαραστήσουμε 63.486 χαρακτήρες, ενώ οι υπόλοιπες 2.048 δεκαεξάδες χρησιμοποιούνται για την αναπαράσταση άλλων 1.048.544 χαρακτήρων χρησιμοποιώντας ζεύγη δεκαεξάδων τα οποία ονομάζονται *αντικαταστάτες* (surrogates). Έτσι, το Unicode μάς επιτρέπει να αναπαραστήσουμε 1.112.030 διαφορετικούς χαρακτήρες.

Για την αναπαράσταση ενός χαρακτήρα του Unicode, χρησιμοποιούμε όλα τα 256 διαφορετικά byte. Αυτό όμως έχει το μειονέκτημα ότι αν έχουμε μια εφαρμογή η οποία αναμένει κάποια συγκεκριμένα byte και αυτή πρέπει να επεξεργαστεί ένα αρχείο Unicode, τότε είναι πολύ πιθανό να αποτύχει. Έτσι για λόγους συμβατότητας το πρότυπο Unicode ορίζει και μια μορφή κωδικοποίησης στην οποία χρησιμοποιούμε οκτάδες bit αντί για δεκαεξάδες για την αναπαράσταση των χαρακτήρων. Η μορφή αυτή είναι γνωστή ως UTF-8, ενώ η εξ ορισμού μορφή των 16 bit είναι γνωστή ως UTF-16.

Στην κωδικοποίηση UTF-8 οι λατινικοί χαρακτήρες που περιέχονται στο πρότυπο ASCII αναπαριστώνται όπως πριν, δηλαδή με ένα byte. Για την αναπαράσταση των ιαπωνικών, κινεζικών και κορεατικών χαρακτήρων χρησιμοποιούνται τρία byte, ενώ για τους χαρακτήρες που απαντώνται στα απλά φθογγικά αλφάβητα (π.χ., το ελληνικό, το κυριλλικό, κ.ά.) χρησιμοποιούνται δύο byte. Τέλος, για την αναπαράσταση των επιπλέον χαρακτήρων χρησιμοποιούνται τέσσερα byte.

Η κωδικοποίηση UTF-16 όταν δεν χρησιμοποιούνται οι αντικαταστάτες, είναι γνωστή ως κωδικοποίηση UCS-2. Ακόμη, επειδή υπάρχουν δύο τρόποι για την αποθήκευση των bit σ' ένα byte, θα πρέπει είτε ο πρώτος χαρακτήρας του αρχείου μας να είναι ο ZERO WIDTH NO-BREAK SPACE, είτε να καθορίζουμε εμείς τον τρόπο κωδικοποίησης των byte. Στην περίπτωση που τα λιγότερο σημαντικά bit βρίσκονται στο τέλος, λέμε ότι τα byte έχουν μορφή big endian, αλλιώς ότι έχουν μορφή little endian. Για να κατάλαβεται τη διαφορά, σας λέμε ότι αν τα byte έχουν μορφή big endian το γράμμα C αναπαριστάμε με τα ψηφία 00000000 01000011, ενώ στην άλλη περίπτωση με τα ψηφία 00000000 01100001. Έτσι έχουμε τις κωδικοποιήσεις: UTF-16BE, UTF-16LE, UCS-2BE και UCS-2LE.

Οι αντικαταστάτες σχηματίζονται από δύο ζεύγη byte τα οποία ονομάζονται *άνω αντικαταστάτης* (high surrogate) και *κάτω αντικαταστάτης* (low surrogates). Ο άνω αντικαταστάτης θα πρέπει να αντιστοιχεί σ' ένα κωδικό σημείο μεταξύ του D800 και του DBFF, ενώ ο κάτω αντικαταστάτης θα πρέπει να αντιστοιχεί σ' ένα κωδικό σημείο μεταξύ DC00 και DFFF.

3. Οι χαρακτήρες που «υποστηρίζει» το Unicode

Σ' αυτή την ενότητα παρουσιάζουμε συνοπτικά τους διάφορους χαρακτήρες που υποστηρίζει προς το παρόν το Unicode. Κατ' αρχάς, οι γραφές που υποστηρίζονται είναι: αραβική, αρμενική, Bengali, Bopomofo, Buhid, καναδική συλλαβική γραφή, τσερόκι, κυριλλική, Deseret, Devanagari, αιθιοπική, γεωργιανή, γοτθική, ελληνική, Gujarati, Gurmukhi, κινεζική (ή Χαν, όπως είναι το σωστό), κορεατική, εβραϊκή, ιαπωνική (συλλαβική και ιδεογραφική), χμέρ, λατινικές γραφές, γραφή της γλώσσας του Λάος, Malayalam, μογγολικά, Myanmar, Ogham, ετρουσκική, Oriya, ρουνική, γι, θιβητιανή, ταϊλανδέζικη, κ.ά.²

Επιπλέον, υπάρχουν και αρκετά σύμβολα τα οποία δεν είναι γράμματα και περιλαμβάνονται στο πρότυπο Unicode. Τα σύμβολα αυτά είναι: αριθμοί, τονικά σύμβολα, σύμβολα στίξης, γενικά σύμβολα, μαθηματικά σύμβολα, μουσικά σύμβολα (δυτικά και βυζαντινά), τεχνικά σύμβολα, βέλη, πλαίσια και άλλα γεωμετρικά σύμβολα, τα σύμβολα του αλφαβήτου Braille για τους τυφλούς και τα Kangxi radicals που χρησιμοποιούνται στα κινεζικά λεξικά μιας και δεν είναι δυνατό να υπάρχουν χιλιάδες ενότητες.

Αν και το Unicode υποστηρίζει όλες τις γραφές που παρουσιάσαμε παραπάνω, είναι επιπλέον συμβατό και με τα διάφορα πρότυπα 8-bit, όπως για παράδειγμα το Latin-1, το ISO8859-7, κ.λπ. Ο Πίνακας 1 δείχνει την κωδικοποίηση της ελληνικής μονοτονικής γραφής: Οι χαρακτήρες που δεν φαίνονται είναι απλά εναλλακτικές μορφές διαφόρων γραμμάτων (π.χ., π και π) οι οποίες όμως δεν υποστηρίζονται από τη γραμματοσειρά που χρησιμοποιούμε.

Στον Πίνακα 2 στο τέλος του άρθρου, παρουσιάζουμε την κωδικοποίηση της ελληνικής πολυτονικής γραφής κατά το πρότυπο Unicode.

Αξίζει να αναφέρουμε ότι οι πίνακες ετούτου του άρθρου στοιχειοθετήθηκαν με το Λ και μια εν δυνάμει γραμματοσειρά τύπου Unicode η οποία στηρίζεται στις γραμματοσειρές που χρησιμοποιούνται στο πακέτο babel. Η γραμματοσειρά είναι τμήμα μιας σειράς εργαλείων τα οποία ελπίζουμε να διαθέσουμε προσεχώς στο ευρύ κοινό.

Περισσότερες πληροφορίες σχετικά με το πρότυπο Unicode μπορείτε να βρείτε στο δικτυακό τόπο: <http://www.unicode.org>.

Η Συντακτική Επιτροπή του περιοδικού θα ήθελε να προσκαλέσει όλους τους αναγνώστες του περιοδικού (αλλά και άλλους ενδιαφερόμενους) να μελετήσουν τους πίνακες του πρότυπου Unicode και να παρουσιάσουν τις δικές τους απόψεις. Οι συγγραφείς εύχονται το άρθρο τους και άλλα παρόμοια που έχουν δημοσιευτεί στο *Εϋτυπον* να δώσουν την αφορμή για μια γόνιμη και δημιουργική ανταλλαγή

² Οι γραφές που αναφέρουμε εδώ με λατινικά γράμματα είναι γραφές που δεν τις έχουμε συναντήσει ποτέ στα ελληνικά. Απαντώνται κυρίως στην Ινδική Χερσόνησο.

0374	´	0375	˘	037A	˙	037E	;	0386	Ά
0388	Έ	0389	Ή	038A	Ί	038C	Ό	038E	Ύ
038F	Ω	0390	ι	0391	Α	0392	Β	0393	Γ
0394	Δ	0395	Ε	0396	Ζ	0397	Η	0398	Θ
0399	Ι	039A	Κ	039B	Λ	039C	Μ	039D	Ν
039E	Ξ	039F	Ο	03A0	Π	03A1	Ρ	03A3	Σ
03A4	Τ	03A5	Υ	03A6	Φ	03A7	Χ	03A8	Ψ
03A9	Ω	03AA	Ϊ	03AB	Ϋ	03AC	ά	03AD	έ
03AE	ή	03AF	ί	03B0	ύ	03B1	α	03B2	β
03B3	γ	03B4	δ	03B5	ε	03B6	ζ	03B7	η
03B8	θ	03B9	ι	03BA	κ	03BB	λ	03BC	μ
03BD	ν	03BE	ξ	03BF	ο	03C0	π	03C1	ρ
03C2	ς	03C3	σ	03C4	τ	03C5	υ	03C6	φ
03C7	χ	03C8	ψ	03C9	ω	03CA	ϊ	03CB	ϋ
03CC	ό	03CD	ύ	03CE	ώ	03D0		03D1	
03D2		03D3		03D4		03D5		03D6	
03D7	κῶ	03DA	ξ	03DB	Ϝ	03DC	Ϝ	03DD	ϝ
03DE	ϝ	03DF	Ϟ	03E0	ϟ	03E1	ϟ		

Πίνακας 1: Η κωδικοποίηση των χαρακτήρων της ελληνικής μονοτονικής γραφής κατά το πρότυπο Unicode.

1FB8	Ά	1FB9	Ἀ	1FBA	Α	1FBB	Α	1FBC	Α _τ
1FBD		1FC2	ἥ	1FC3	η	1FC4	ἥ	1FC6	ἦ
1FC7	ἦ	1FC8	Έ	1FC9	Ε	1FCA	Ή	1FCB	Ή
1FCC	Η _τ	1FD0	ι	1FD1	ι	1FD2	ι	1FD3	ι
1FD6	ι	1FD7	ϊ	1FD8	Ϊ	1FD9	Ϊ	1FDA	Υ
1FDB	Ί	1FE0	ύ	1FE1	υ	1FE2	υ	1FE3	υ
1FE4	ϋ	1FE5	ϋ	1FE6	ϋ	1FE7	ϋ	1FE8	Ϋ
1FE9	Υ	1FEA	Ύ	1FEB	Ύ	1FEC	Ρ	1FF2	Ϙ
1FF3	Ϙ	1FF4	ϙ	1FF6	ϙ	1FF7	ϙ	1FF8	Ό
1FF9	Ό	1FFA	Ω	1FFB	Ω	1FFC	Ω _τ		

Πίνακας 2: Η κωδικοποίηση των χαρακτήρων της ελληνικής πολυτονικής γραφής κατά το πρότυπο Unicode.

απόψεων και προτάσεων για την περαιτέρω ορθότερη ανάπτυξη του προτύπου αυτού όσον αφορά την ελληνική γλώσσα.

1F00	ǎ	1F01	ǎ̇	1F02	ǎ̈	1F03	ǎ̄	1F04	ǎ̆
1F05	ǎ̇	1F06	ǎ̈	1F07	ǎ̄	1F08	ǎ̆	1F09	ǎ̈
1F0A	ǎ̄	1F0B	ǎ̆	1F0C	ǎ̈	1F0D	ǎ̄	1F0E	ǎ̆
1F0F	ǎ̆	1F10	ǎ̈	1F11	ǎ̄	1F12	ǎ̆	1F13	ǎ̈
1F14	ǎ̄	1F15	ǎ̆	1F18	ǎ̈	1F19	ǎ̄	1F1A	ǎ̆
1F1B	ǎ̈	1F1C	ǎ̄	1F1D	ǎ̆	1F20	ǎ̈	1F21	ǎ̄
1F22	ǎ̆	1F23	ǎ̈	1F24	ǎ̄	1F25	ǎ̆	1F26	ǎ̈
1F27	ǎ̄	1F28	ǎ̆	1F29	ǎ̈	1F2A	ǎ̄	1F2B	ǎ̆
1F2C	ǎ̈	1F2D	ǎ̄	1F2E	ǎ̆	1F2F	ǎ̈	1F30	ǎ̄
1F31	ǎ̆	1F32	ǎ̈	1F33	ǎ̄	1F34	ǎ̆	1F35	ǎ̈
1F36	ǎ̄	1F37	ǎ̆	1F38	ǎ̈	1F39	ǎ̄	1F3A	ǎ̆
1F3B	ǎ̈	1F3C	ǎ̄	1F3D	ǎ̆	1F3E	ǎ̈	1F3F	ǎ̄
1F40	ǎ̆	1F41	ǎ̈	1F42	ǎ̄	1F43	ǎ̆	1F44	ǎ̈
1F45	ǎ̄	1F48	ǎ̆	1F49	ǎ̈	1F4A	ǎ̄	1F4B	ǎ̆
1F4C	ǎ̈	1F4D	ǎ̄	1F50	ǎ̆	1F51	ǎ̈	1F52	ǎ̄
1F53	ǎ̆	1F54	ǎ̈	1F55	ǎ̄	1F56	ǎ̆	1F57	ǎ̈
1F59	ǎ̄	1F5B	ǎ̆	1F5D	ǎ̈	1F5F	ǎ̄	1F60	ǎ̆
1F61	ǎ̈	1F62	ǎ̄	1F63	ǎ̆	1F64	ǎ̈	1F65	ǎ̄
1F66	ǎ̆	1F67	ǎ̈	1F68	ǎ̄	1F69	ǎ̆	1F6A	ǎ̈
1F6B	ǎ̄	1F6C	ǎ̆	1F6D	ǎ̈	1F6E	ǎ̄	1F6F	ǎ̆
1F70	ǎ̈	1F71	ǎ̄	1F72	ǎ̆	1F73	ǎ̈	1F74	ǎ̄
1F75	ǎ̆	1F76	ǎ̈	1F77	ǎ̄	1F78	ǎ̆	1F79	ǎ̈
1F7A	ǎ̄	1F7B	ǎ̆	1F7C	ǎ̈	1F7D	ǎ̄	1F80	ǎ̆
1F81	ǎ̈	1F82	ǎ̄	1F83	ǎ̆	1F84	ǎ̈	1F85	ǎ̄
1F86	ǎ̆	1F87	ǎ̈	1F88	ǎ̄	1F89	ǎ̆	1F8A	ǎ̈
1F8B	ǎ̄	1F8C	ǎ̆	1F8D	ǎ̈	1F8E	ǎ̄	1F8F	ǎ̆
1F90	ǎ̈	1F91	ǎ̄	1F92	ǎ̆	1F93	ǎ̈	1F94	ǎ̄
1F95	ǎ̆	1F96	ǎ̈	1F97	ǎ̄	1F98	ǎ̆	1F99	ǎ̈
1F9A	ǎ̄	1F9B	ǎ̆	1F9C	ǎ̈	1F9D	ǎ̄	1F9E	ǎ̆
1F9F	ǎ̈	1FA0	ǎ̄	1FA1	ǎ̆	1FA2	ǎ̈	1FA3	ǎ̄
1FA4	ǎ̆	1FA5	ǎ̈	1FA6	ǎ̄	1FA7	ǎ̆	1FA8	ǎ̈
1FA9	ǎ̄	1FAA	ǎ̆	1FAB	ǎ̈	1FAC	ǎ̄	1FAD	ǎ̆
1FAE	ǎ̈	1FAF	ǎ̄	1FB0	ǎ̆	1FB1	ǎ̈	1FB2	ǎ̄
1FB3	ǎ̆	1FB4	ǎ̈	1FB6	ǎ̄	1FB7	ǎ̆		

Πίνακας 2: συνέχεια.