

Programiranje 1
Beleške sa vežbi
Školska 2007/2008 godina

Matematički fakultet, Beograd

Jelena Tomašević

October 23, 2007

Sadržaj

| | | |
|----------|---|----------|
| 1 | | 5 |
| 1.1 | Reprezentacija znakovnih podataka | 5 |
| 1.1.1 | Tekst je niz karaktera | 5 |
| 1.1.2 | Zapis karaktera u računaru | 5 |
| 1.1.3 | Skupovi karaktera | 5 |
| 1.1.4 | Kodiranje ostalih jezika | 6 |
| 1.1.5 | Osobine YUSCII koda | 6 |
| 1.1.6 | Kodne strane | 7 |
| 1.1.7 | Kodne strane kod nas | 7 |
| 1.1.8 | Višebajtni karakterski kodovi | 7 |
| 1.1.9 | Karakter, Glifovi, Fontovi | 8 |

1

1

1.1 Reprezentacija znakovnih podataka

1.1.1 Tekst je niz karaktera

- Iako obično tekst zamišljamo kao dvodimenzioni objekat, u računarima se tekst predstavlja kao jednodimenzioni (linearni) niz karaktera.
- Potrebno je, dakle, uvesti specijalne karaktere koji označavaju prelazak u novi red, tabulator, kraj teksta i slično

1.1.2 Zapis karaktera u računar

- Računari su zasnovani na binarnoj aritmetici
- Cele brojeve je moguće predstaviti u binarnom sistemu
- Osnovna ideja je svakom karakteru pridružiti određeni ceo broj na unapred dogovoreni način
- Ove brojeve zovemo kodovima karaktera (character codes)

1.1.3 Skupovi karaktera

- Koliko karaktera želimo da predstavimo u računarima? Tokom razvoja računarstva broj karaktera je postajao sve veći
- Pošto je u početku razvoja englesko govorno područje bilo dominantno osnovno je bilo predstaviti sledeće karaktere :
 - Velika slova engleskog alfabeta : A,B,...,Z
 - Mala slova engleskog alfabeta : a,b,...,z
 - Cifre : 0,1,...,9
 - Interpunkcijske znake : ., ; i slično
 - Kontrolne znake : kraj reda, tabulator i slično
- Standardni karakterski kodovi: Sedamdesetih godina su se pojavile tabele standardnih karakterskih kodova dovoljne za zapis pomenutih karaktera
Najpoznatiji je ASCII Standard iz koga se razvila većina današnjih standarda

¹Zasnovano na materijalu "Zapis tekstova u računar" Filipa Marića

- ASCII (American Standard Code for Information Interchange)
ASCII je sedmobitan (broj karaktera koji je njime predstavljen je $128 = 2^7$)
- Npr. kod za A je $(41)_{16}$ tj. $0x41$ što je $(65)_{10}$ tj. $(1000001)_2$,
kod za a je $(61)_{16}$ tj. $0x61$ što je $(95)_{10}$ tj. $(1100001)_2$.
- PRIMER: transformisanje malih slova u velika
- Razmak SP se zapisuje kao $(20)_{16}$ što je $(32)_{10}$ tj. $(0100000)_2$.
- **Osobine ASCII koda:** Prvih 32 karaktera (kodovi $0x00-0x1F$) i poslednji karakter (kod $0x7F$) su kontrolni karakteri. To su karakteri bez grafije, kao CR (kod $0x0D$), LF (kod $0x0A$).
- Prvi karakter sa grafijom je blanko (kod $0x20$). Njegova grafija je belina.
- Skup velikih slova A-Z (kodovi $0x41-0x5A$), kao i skup malih slova a-z (kodovi $0x61-0x7A$), je u alfabetskom redosledu unutar kolacione sekvencije
($0x41 < 0x42$, prema tome $A < B$ to odgovara alfabetskom redosledu).
- Skup cifara 0 – 9 (kodovi $0x31 – 0x39$) je u rastućem brojčanom redosledu unutar kolacione sekvencije
($0x31 < 0x32$, prema tome $1 < 2$ što odgovara brojčanom redosledu).
- Skup velikih slova A-Z (kodovi $0x41 – 0x5A$), skup malih slova a-z (kodovi $0x61 – 0x7A$) i skup cifara 0-9 (kodovi $0x31 – 0x39$) su kontingentni unutar kolacione sekvencije (između slova A i slova Z nema drugih karaktera osim onih koji odgovaraju velim slovima engleske abecede).
- Sve cifre prethode svim velikim slovima, sva velika slova prethode svim malim slovima u kolacionoj sekvenciji. Specijalni i interpunkcijski znaci su izmešani između njih.
- Kod svakog velikog slova je za 32 (ili $0x20$) manji od koda odgovarajućeg malog slova. Na primer, za slovo E važi da je $0x45 + 0x20 = 0x65$, odnosno, $01000101 + 00100000 = 01100101$. Prema tome, binarni kodovi velih i malih slova razlikuju se samo u jednoj cifri, onoj koja odgovara petom stepenu osnove 2.

1.1.4 Kodiranje ostalih jezika

- Razvojem računarstva se javlja potreba kodiranja tekstova i na drugim jezicima
- Kroz istoriju su postojala mnoga rešenja, od kojih su se neka zadržala, a neka su nestala

1.1.5 Osobine YUSCII koda

- ASCII kod je jedna verzija međunarodnog standarda ISO 646 IRV koji predstavlja međunarodnu referentnu verziju za 7-bitni kod. Ovaj standard propisuje da se pozicijama $0x40$, $0x5B – 0x5E$, $0x60$ i $0x7B – 0x7E$ ne pridružuje obavezna grafija već se da se one u nacionalnim verzijama standarda i u određenim aplikacijama mogu slobodno koristiti.
- Standard JUS.B1.002 koristi ovih 10 pozicija za kodiranje slova specifičnih za srpsku latinicu Ž, Š, Đ, Ć, Č.
- Yu-ASCII skup zadržava sve navedene osobine ASCII koda osim jedne, a ta je da ni velika ni mala slova nisu u alfabetskom redosledu unutar kolacione sekvencije. Naime, $0x40 < 0x41$, dakle Ž < A što ne odgovara alfabetskom redosledu unutar kolacione sekvencije. Nakon slova Ž slede velika slova engleske abecede, zatim slova Š($0x5B$), Đ($0x5C$), Ć($0x5D$), Č($0x5E$), ž($0x60$), zatim slede mala slova abecede, i na kraju slova š($0x7B$), đ($0x7C$), ć($0x7D$), č($0x7E$).

Zadatak 1 Poredati slova vašeg prezimena koristeći YUSCII kodnu šemu.

1.1.6 Kodne strane

- Pod *kodnom* stranom (Code page) tj. *skupom karaktera* (Character set, charset) podrazumevamo uređenu listu karaktera predstavljenih svojim karakterskim kodovima
- Podaci se u računarima obično zapisuju bajt po bajt
- ASCII je sedmobitni standard
- ASCII karakteri se zapisuju tao što se u svakom bajtu bit najveće težine postavi na 0
- To ostavlja prostor za novih 128 karaktera čiji binarni zapis počinje sa 1
- Ovaj prostor se može popuniti na razne načine
- Rešenje nije univerzalno, jer svakako na svetu postoji više od 256 različitih karaktera
- Postavljeni su razni standardi dopunjavanja ovih 128 karaktera
- Svim ovim kodnim stranama je zajedničko prvih 128 karaktera i oni se poklapaju sa ASCII
- Ovako napravljene kodne strane obično omogućuju kodiranje tekstova na više srodnih jezika (obično i geografski bliskih)
- Nama su uglavnom važne kodne strane napravljene za centralno-evropske (Central European) latinice, kao i ćirilicne kodne strane

1.1.7 Kodne strane kod nas

- Najčešće korišćene kodne strane kod nas i koje su delo međunarodne organizacije za standardizaciju (International Standard organization)
 - ISO 8859-2 (Latin2)
 - ISO 8859-5 (Ćirilicna)
- **Latin 1:** Poželjno je poznavati i osnovnu kodnu stranu ISO 8859-1 (Latin1) jer je veoma često postavljena kao podrazumevana kodna strana. Ona se koristi za zapis tekstova na zapadno evropskim jezicima (Western European)

1.1.8 Višebajtni karakterski kodovi

- Azijskim jezicima nije dovoljno 256 mesta za zapis svih karaktera.
- Zbog toga se uvode višebajtni karakterski kodovi
- **ISO 10646** je zamišljen kao 4 bajtni standard. Pri tome se prvih 65536 karaktera koriste kao osnovni višejezični skup karaktera dok je ostali prostor ostavljen kao proširenje za drevne jezike, celokupnu naucnu notaciju i slično.
- **UNICODE:** svakom karakteru dodeljuje dvobajtni kod
- Prvih 128 karaktera se poklapaju sa ASCII standardom, dok su sledećih 128 napravljene tako da se poklapaju sa Latin1 standardom

- UTF: A Unicode transformation format (UTF) algoritam koji svakom UNICODE karakteru dodeljuje određeni niz bajtova čija dužina varira od 1 do najviše 6.
- Ovo značajno otežava tumačenje podataka
- UTF je ASCII kompatibilan, što znaci da se ASCII karakteri zapisuju pomoću jednog bajta, na standardni način.
- Najčešće korišćena varijanta ovog algoritma je UTF-8 koja je dovoljna za zapis svih dvobajtnih UNICODE karaktera
- Pored ovoga ISO uvodi i UTF-16, UTF-32, kao i standard UCS-4

1.1.9 Karakteri, Glifovi, Fontovi

- Vrlo često se ne pravi jasna razlika između karaktera i njihove graficke reprezentacije
- Grafičku reprezentaciju karaktera nazivamo glifovima (*glyph*) Skupove glifova nazivamo fontovima (*font*)
- Korespondencija između karaktera i glifova ne mora biti jednoznačna
- Jedan glif može da predstavi više karaktera (ligature)
- Isti karakter može da se predstavlja različitim glifovima u zavisnosti od svoje pozicije u reči
- Za razliku od tradicionalnih fontova koji u sebi sadže glifove za karaktere jedne kodne strane, TrueType fontovi koji podržavaju WGL4 standard sadrže glifove za sve evropske karaktere

Zadatak 2 Zapisati cifru 3 u ASCII kodu.

Rešenje:

Broj 3 se zapisuje kao $(33)_{16}$ tj. $0x33$ što je $(51)_{10}$ tj. $(1010001)_2$

Zadatak 3 Zapisati reč Fakultet u ASCII kodu.