

Uvod u organizaciju i arhitekturu računara 1

Stefan Mišković

5 Zapis realnih brojeva

5.1 Zapis u nepokretnom i pokretnom zarezu

Za zapis realnih brojeva u računaru, bez obzira koliko memorije je odvojeno za zapis broja, uvek može da se zapiše konačno mnogo brojeva. Zbog toga razlomljeni deo broja, ukoliko postoji, uvek mora biti zapisan sa konačnim brojem cifara. Svaki zapis određuje opseg brojeva koji se može zapisati, kao i broj dozvoljenih cifara u razlomljenom delu. Realni brojevi se mogu zapisati u nepokretnom (fiksnom) i pokretnom zarezu.

Kod zapisa u nepokretnom zarezu je unapred određen broj cifara za zapis broja i broj cifara za zapis razlomljenog dela. Nekad sam programer mora voditi računa o mestu zapisa decimalne tačke, zbog čega u praksi može doći do greške. Na primer, neka je dozvoljena dužina zapisa za razlomljeni deo 11 i neka je potrebno zapisati heksadekadni broj $(0.919)_{16}$. Za zapis razlomljenog dela ovog broja je potrebno 12 cifara, a na ovaj način greškom poslednji bit zapisa može biti odbačen. Tada se greškom može dobiti heksadekadni broj $(0.48C)_{16}$, što je ilustrovano sledećom šemom:

```
Ispravan zapis:           9       1       9
                        0|1 0 0 1|0 0 0 1|1 0 0 1
-----
Zapis u memoriji:       0 1 0 0 1 0 0 0 1 1 0 0 1
-----
                        0 1 0 0|1 0 0 0|1 1 0 0|1
Pogrešan zapis:        4       8       C
```

Pri izvođenju aritmetičkih operacija može doći do greške, ukoliko se ne vodi računa o broju cifara u razlomljenom delu. Na primer, pri sabiranju dekadnih brojeva $(0.63)_{10}$ i $(0.2)_{10}$, umesto da se brojevi ispravno potpišu pri sabiranju i dobije korektan rezultat $(0.83)_{10}$, može se zanemariti dodavanje nule kod drugog sabirka i time dobiti pogrešan zbir:

```
  0 6 3
+ 0 0 2
-----
  0 6 5
```

Kod brojeva u pokretnom zarezu, ove greške se ne mogu pojaviti, jer se uvek automatski vodi računa o mestu decimalne tačke. Zapis u pokretnom zarezu je znatno praktičniji od zapisa u nepokretnom, zbog čega se i danas najčešće koristi. Svaki realan broj se time može napisati u obliku $\pm f \cdot n^e$, gde \pm označava da se mogu zapisati i pozitivni

i negativni brojevi, f frakciju (značajni deo broja), n osnovu, a e eksponent. Zapis kod koga je $f = f_0.f_1f_2 \dots f_{p-1}$ i $f_0 \neq 0$ se naziva normalizovan. Broj cifara frakcije (ukupan broj cifara u celobrojnom i razlomljenom delu) se naziva preciznost, u oznaci p .

Na primer, ako je broj $(0.4)_{10}$ potrebno zapisati u osnovi $n = 10$ i sa preciznošću $p = 4$, zapis je $4.000 \cdot 10^{-1}$, pa je normalizovana frakcija $f = 4.000$, a eksponent $e = -1$. Ukoliko se ne zahteva da frakcija bude normalizovana, isti broj se može napisati na više načina, na primer, u obliku $0.040 \cdot 10^1$ ($f = 0.040$ i $e = 1$) ili $0.004 \cdot 10^2$ ($f = 0.004$ i $e = 2$). Ako je broj $(0.4)_{10}$ potrebno zapisati u osnovi $n = 2$ i sa preciznošću $p = 6$ u normalizovanom obliku, najpre je potrebno izvršiti prevođenje $(0.4)_{10} = (0.011001100110\dots)_2$. Da bi broj mogao biti zapisan sa 6 značajnih cifara, mora se izvršiti odsecanje ili zaokruživanje razlomljenog dela (ovde će biti izvršeno odsecanje). Tako je traženi zapis $(1.10011)_2 \cdot 2^{-2}$. Primetimo da je frakcija uvek zapisana u osnovi n , a eksponent u dekadnoj osnovi. Osnova n se uvek podrazumeva, tako da se pri kodiranju njen zapis izostavlja, o čemu će više biti reči kada se bude govorilo o konkretnim načinima zapisa.

Za razliku od brojeva zapisanih u nepokretnom zarezu, ovde se mogu zapisati i veoma veliki brojevi, ukoliko nemaju veliki broj značajnih cifara. Na primer, broj $53 \cdot 10^{100}$, ako je $n = 10$ i $p = 5$, može se zapisati u obliku $5.3000 \cdot 10^{101}$. Realni brojevi u pokretnom zarezu se u računaru najčešće zapisuju u jednostrukoj (32 bita), dvostrukoj (64 bita) i četverostrukoj tačnosti (128 bitova).

5.2 Zapis sa binarnom osnovom na starijim računarima

Ilustracije radi, ovde će biti prikazano kako su se realni brojevi zapisivali na nekim od starijih verzija računara. Brojevi su bili zapisivani u binarnoj osnovi i podržavali su jednostruku i dvostruku tačnost. Frakcija proizvoljnog dekadnog broja se u ovom kontekstu piše u obliku $0.f_0f_1 \dots f_{p-1}$, pri čemu mora važiti $f_0 \neq 0$. Kako se ovde radi o binarnom sistemu, mora biti $f_0 = 1$, pa se ova cifra, budući da je podrazumevano uvek jednaka 1, izostavlja. Eksponent se piše na 8 mesta u zapisu višak 128. To znači da se eksponent, nakon dodavanja 128 u dekadnom sistemu, pretvara u binarni sistem i zapisuje na 8 mesta (sa eventualnim vodećim nulama, ukoliko je potrebno). Osnova 2 se podrazumeva i ona se ne kodira. Takođe, znak broja se kodira pomoću jednog bita. Ako je broj pozitivan, odgovarajući bit je 0, a ako je negativan, odgovarajući bit je 1.

U zapisu u jednostrukoj tačnosti, na raspolaganju su 32 bita. Sleva nadesno, oni su podeljeni u tri grupe. Prvi bit označava znak, narednih 8 eksponent, a preostala 23 bita frakciju. Kod dvostruke tačnosti se znak i eksponent zapisuju na isti način, s tim što su dodata još 32 bita za frakciju. Ovim se povećava preciznost. Specijalno, 0 se piše sa svim nulama u zapisu, bilo da se radi o jednostrukoj ili dvostrukoj tačnosti.

Primetimo da se na ovaj način može zapisati samo konačan interval brojeva. Ukoliko je broj veliki po apsolutnoj vrednosti tako da ne može da se zapiše, dolazi do prekoračenja. Nasuprot tome, broj može biti po apsolutnoj vrednosti veoma blizak nuli i nedovoljno veliki za zapis, kada dolazi do potkoračenja.

U nastavku će biti navedeno nekoliko primera kodiranja i dekodiranja brojeva u ovom zapisu:

- Broj 0 se piše kao sve nule u zapisu. Njegov zapis u jednostrukoj tačnosti je $0\ 00000000\ \underbrace{0\dots0}_{23\ \text{nule}}$.
- Za zapis dekadnog broja 15 u jednostrukoj tačnosti, najpre je potrebno izvršiti prevođenje $(15)_{10} = (1111)_2$. Pogodan oblik broja za zapis mora biti takav da

se frakcija sastoji samo od razlomljenih cifara, pri čemu je prva cifra jednaka 1. Pogodan zapis je $(0.1111)_2 \cdot 2^4$. Eksponent 4 je potrebno napisati sa uvećanjem 128. Kako je $4 + 128 = 132$, zapis eksponenta je $(132)_{10} = (10000100)_2$. Kao što je već navedeno, kod frakcije 0.1111 se deo 0.1 izostavlja, pa je njen zapis 111, sa odgovarajućim brojem nula u nastavku frakcije. Broj je pozitivan, pa je bit za znak jednak 0. Konačno, zapis je $0\ 10000100\ 111\ \underbrace{0\dots 0}_{20\ \text{nula}}$.

- Za zapis dekadnog broja -15 u dvostrukoj tačnosti se dobija $(-15)_{10} = (-1111)_2 = (-0.1111)_2 \cdot 2^4$. Kod eksponenta je $(10000100)_2$, kao u prethodnom primeru. Bit za znak je 1, jer se radi o negativnom broju, a frakcija se sastoji od 111, nakon čega treba nadovezati odgovarajući broj nula. Zapis broja je $1\ 10000100\ 111\ \underbrace{0\dots 0}_{52\ \text{nule}}$.
- Neka je dat zapis $0\ 01111011\ \underbrace{0\dots 0}_{23\ \text{nule}}$, koji je potrebno prevesti u dekadni sistem. Vidimo da je vrednost eksponenta $(1111011)_2 - 128 = 123 - 128 = -5$. Kako se frakcija sastoji samo od nula, njena binarna vrednost je 0.1. Broj je pozitivan, pa je njegova dekadna vrednost $(0.1)_2 \cdot 2^{-5} = 2^{-1} \cdot 2^{-5} = 2^{-6} = 1/64$.

5.3 IEEE 754 standard

U početku nije postojao jedinstven standard za zapis realnih brojeva, pa je bilo moguće da različite arhitekture implementiraju zapis na različite načine, zbog čega može doći do grešaka i nepodudaranja pri prenosu podataka. Zbog potrebe sa unifikacijom načina zapisa realnih brojeva, nastao je IEEE 754 standard. Njime je precizno propisan način zapisa brojeva, kao i algoritmi za različite operacije nad njima, ponašanje u graničnim situacijama, i slično. Kako danas sve mašine podržavaju ovaj standard, smanjena je verovatnoća greške pri prenosu podataka. IEEE 754 standard podržava zapis brojeva sa dekadnom i binarnom osnovom, a mogući su zapisi sa jednostrukom, dvostrukom i četverostrukom tačnošću. Oni redom zauzimaju 32, 64 i 128 bitova.

Zaokruživanje. U računaru uvek može da se zapiše konačno mnogo realnih brojeva, zbog čega se neki brojevi ne mogu precizno zapisati, već se pamti njihova približna vrednost, pri čemu se vrši zaokruživanje. Postoji nekoliko načina na koji se može izvršiti zaokruživanje, među kojima su:

- Zaokruživanje na parnu cifru – broj se zaokružuje na najbližu pretpostavljenu vrednost, a ako je na sredini intervala, onda se zaokruživanje vrši na parnu cifru. Tako se, na primer, brojevi 1.212, 1.218, 1.225 i 1.235 zaokružuju na dve decimale redom na vrednosti 1.21, 1.22, 1.22 i 1.24.
- Zaokruživanje ka $+\infty$ – broj se zaokružuje tako da mu je zaokružena vrednost uvek veća ili jednaka od stvarne. Tako se, na primer, pozitivni brojevi 1.21 i 1.29 na jednu decimalu zaokružuju redom na 1.3 i 1.3, a negativni brojevi -1.21 i -1.29 redom na -1.2 i -1.2 .
- Zaokruživanje ka $-\infty$ – broj se zaokružuje tako da mu je zaokružena vrednost uvek manja ili jednaka od stvarne. Tako se, na primer, pozitivni brojevi 1.21 i 1.29 na jednu decimalu zaokružuju redom na 1.2 i 1.2, a negativni brojevi -1.21 i -1.29 redom na -1.3 i -1.3 .

- Zaokruživanje ka nuli – broj se zaokružuje tako da mu je zaokružena vrednost uvek po apsolutnoj vrednosti manja ili jednaka od stvarne. Tako se, na primer, pozitivni brojevi 1.21 i 1.29 na jednu decimalu zaokružuju redom na 1.2 i 1.2, a negativni brojevi -1.21 i -1.29 redom na -1.2 i -1.2 .

Neka 0.0574367 predstavlja matematički tačnu vrednost, a neka približan zapis broja nakon zaokruživanja za preciznost $p = 3$ iznosi $5.74 \cdot 10^{-2}$. Možemo primetiti da je razlika stvarne i zaokružene vrednosti 0.0000367 , a ako to posmatramo u jedinicama koje su u odnosu na poslednje mesto zaokružene zapisa, greška iznosi 0.367 . Veličina koja predstavlja broj jedinica na poslednjem mestu se naziva ULP. Drugi način merenja greške je relativna greška. Ona predstavlja apsolutnu vrednost razlike stvarne i zaokružene vrednosti broja, podeljenu sa apsolutnom vrednošću broja. U našem primeru ona iznosi $|0.0574367 - 0.0574|/|0.0574367| \approx 0.0006$. Osnovna razlika između ULP-a i relativne greške se ogleda u tome šta oni mere. Dok se pomoću ULP meri stvarna greška, pomoću relativne greške se meri odgovarajući odnos, zbog čega se ona najčešće koristi kod analize greške koja se javlja pri različitim izračunavanjima. Time, ako se broj poveća nekoliko puta, toliko puta će se povećati i ULP, dok relativna greška ostaje ista.

Cifre čuvari. Cifre čuvari predstavljaju dodavanje još jedne cifre na zapis broja pri primeni računskih operacija, gde se povećava preciznost za 1 i time smanjuje greška. Neka je potrebno oduzeti brojeve 10.1 i 9.93 , ako su zapisani u obliku $1.01 \cdot 10^1$ i $0.99 \cdot 10^1$. Njihovim oduzimanjem se dobija rezultat $0.02 \cdot 10^1$, a stvarna vrednost bi trebalo da bude 0.17 . Preciznija razlika se može dobiti uvođenjem dodatne cifre čuvara. Tada se oduzimanjem zapisa $1.010 \cdot 10^1$ i $0.993 \cdot 10^1$ dobija vrednost $0.017 \cdot 10^1$.

Tačno zaokružene operacije. Kada se koriste cifre čuvari, u računaru se preciznost uvećava za 1, pa dobijeni rezultat, iako precizniji od slučaja kada se one ne koriste, nije naročito precizniji. U idealnom slučaju, ukoliko to resursi dozvoljavaju, operacije se mogu izvršiti nad stvarnim vrednostima brojeva, a kasnije zaokružiti rezultat. Za takve operacije kažemo da su tačno zaokružene.

5.4 IEEE 754 standard sa dekadnom osnovom

Kod IEEE 754 standarda sa dekadnom osnovom, osnova je uvek 10, a frakcija može biti dekadna ili binarna. Ako je frakcija dekadna, radi se o DPD kodiranju, a ako je binarna, o BID kodiranju. Za obe vrste kodiranja će biti opisan zapis u jednostrukoj tačnosti, kada se brojevi kodiraju sa 32 bita.

DPD kodiranje. Ukoliko su i frakcija i osnova dekadne, IEEE 754 standard koristi DPD kodiranje. Prethodno smo se upoznali kako se mogu iskoristiti tabele za DPD kodiranje i dekodiranje, pri čemu se parovi od 3 dekadne cifre mogu napisati sa 10 binarnih cifara. Podsetimo se kako one izgledaju (iznad je navedena tabela za kodiranje, a ispod tabela za dekodiranje):

aei	pqr	stu	v	wxy
000	bcd	fgh	0	jkl
001	bcd	fgh	1	00l
010	bcd	jkh	1	01l
100	jdk	fgh	1	10l
110	jdk	00h	1	11l
101	fgd	01h	1	11l
011	bcd	10h	1	11l
111	00d	11h	1	11l

vwxst	abcd	efgh	ijkl
0...	Opqr	Ostu	Owxy
100..	Opqr	Ostu	100y
101..	Opqr	100u	Osty
110..	100r	Ostu	Opqy
11100	100r	100u	Opqy
11101	100r	Opqu	100y
11110	Opqr	100u	100y
11111	100r	100u	100y

U nastavku ćemo kroz primere opisati kako se vrši kodiranje iz dekadne vrednosti u zapis i dekodiranje iz zapisa u dekadnu vrednost. Neka je, na primer, potrebno kodirati broj 123.4 po IEEE 754 standardu u dekadnoj osnovi u jednostrukoj tačnosti koristeći DPD kodiranje. Broj najpre treba napisati kao proizvod celobrojne frakcije i stepenovane osnove. Ovde je $123.4 = 1234 \cdot 10^{-1}$. Celobrojnu frakciju je uvek potrebno zapisati na 7 mesta. Ako ima manje od 7 značajnih cifara, dopunjuje se nulama, a ako ima više, vrši se zaokruživanje. Ovde je prilagođena frakcija oblika 0001234. Ona se sada deli na tri dela, gde prvi deo sadrži jednu cifru, a ostala dva dela po 3 cifre. U našem primeru frakcija se deli na 0, 001 i 234. Poslednja dva dela frakcija se kodiraju sa po 10 bitova koristeći tabelu za kodiranje. Kodirajmo najpre trojku 001. Zapisi 8421 dekadnih cifara 0, 0 i 1 su redom 0000, 0000 i 0001. Ako označimo da je $abcd = 0000$, $efgh = 0000$ i $ijkl = 0001$, budući da je $aei = 000$, iz tabele za kodiranje sledi da je $pqr = bcd$, $stu = fgh$, $v = 0$ i $wxy = jkl$, pa je traženi kod $pqrstuvwxy = 0000000001$. Na sličan način se dobija da se dekadna trojka 234 kodira sa 0100110100. Preostala je još prva cifra frakcije. Ukoliko je ona između 0 i 7, kodira se sa tri binarne cifre, tako što se cifra pretvori u binarni zapis na tri mesta, uz eventualno dodavanje vodećih nula (na primer, 0 se kodira kao 000, 5 kao 101 i 7 kao 111). Ukoliko je cifra frakcije 8 ili 9, ona se kodira sa jednom binarnom cifrom. Osmica se kodira sa 0, a devetka sa 1. U našem primeru je prva cifra frakcije 0, pa se ona kodira kao 000. Ostalo je još da se kodira eksponent. Eksponent -1 je potrebno zapisati u binarnom sistemu na 8 mesta sa uvećanjem 101. Dakle, zapis eksponenta je oblika $-1 + 101 = 100 = (01100100)_2$. Na osnovu koda prve cifre frakcije (3 bita) i eksponenta (8 bitova), formira se kombinacija koja se sastoji od 11 bitova. (Ako je prva cifra 8 ili 9, ona se kodira samo sa jednom binarnom cifrom, pa se kombinacija od 11 bitova dobija na nešto drugačiji način, što će biti objašnjeno u narednom primeru.) Kombinacija se dobija tako što se eksponent razdvoji na dva dela, levi koji sadrži prva dva i desni koji sadrži poslednjih šest bitova, a između njih se ubaci trobitni kod početka frakcije. U našem slučaju, kombinacija je 01000100100. Broj je pozitivan, pa je bit za znak 0. Konačno, prvi bit zapisa je za znak, narednih 11 bitova čini kombinacija, a preostalih 20 ostatak

frakcije. Zapis traženog broja je 0 01000100100 0000000001 0100110100.

Pretpostavimo da je potrebno kodirati na isti način broj $-982.5294 \cdot 10^{42}$. Pogodan zapis broja za kodiranje je $-9825294 \cdot 10^{38}$. Kod cifre 9 je 1, a desetobitni kodovi trojki 825 i 294 su redom 1000101101 i 0101011010. Ovi kodovi su dobijeni koristeći tabelu za kodiranje, slično kao u prethodnom primeru. Zapis eksponenta je $38 + 101 = 139 = (10001011)_2$. Kombinacija od 11 bitova se i ovde dobija od eksponenta i koda prve cifre frakcije, ali je ovde problem što zapis početka frakcije ima samo 1 bit. U slučaju kada je prva cifra frakcije 8 ili 9 (kada se ona ustvari kodira jednim bitom), na početku kombinacije se dodaje 11, a ostatak se dobija analogno, razdvajanjem eksponenta i ubacivanjem cifre frakcije. Dakle, u ovom slučaju kombinacija se sastoji od 11, početka eksponenta 10, početka frakcije 1 i nastavka eksponenta 001011, pa je ona oblika 11101001011. Uzimajući u obzir i da je broj negativan, njegov zapis je 1 11011001011 1000101101 0101011010.

Objasnimo sada obratni proces, kada se vrši prevođenje iz zapisa u dekadni broj na primeru 32-bitnog zapisa 0 01111011110 0000000000 0000110101. Najpre treba posmatrati kombinaciju, jer tu mogu postojati dva slučaja – kada je prva cifra frakcije između 0 i 7, odnosno kada je 8 ili 9. Ukoliko ona počinje sa 11, radi se o početnoj cifri 8 ili 9, a inače, početna cifra je između 0 i 7. U našem primeru prve dve cifre i poslednjih 6 cifara kombinacije su cifre eksponenta, a preostala trojka 111 kodira prvu cifru frakcije čija je vrednost 7. Eksponent iznosi $(01011110)_2 - 101 = 94 - 101 = -7$. Na osnovu tabele za dekodiranje (videti lekciju gde se prvi put pominje DPD), dobija se da kodovi 0000000000 i 000011010 redom predstavljaju trojke dekadnih cifara 000 i 035. Kako je broj pozitivan, konačno se dobija da je tražena vrednost $7000035 \cdot 10^{-7} = 0.7000035$.

Pored zapisa konačnih brojeva, IEEE 754 standard sa dekadnom osnovom (i kod DPD i kod BID kodiranja) podržava zapis specijalnih vrednosti. Kada se vrši dekodiranje, potrebno je prvo razmotriti da li se radi o nekoj specijalnoj vrednosti, pa ako se ustanovi da je u pitanju konačan broj, treba pribeći samom algoritmu. Specijalne vrednosti koje imaju poseban zapis u ovom standardu su:

- Nula se može pojaviti kao konačna vrednost ili kao neka veoma mala vrednost koja se ne može zapisati, kada dolazi do potkoračenja, gde se ta vrednost zapisuje da je približno jednaka nuli. Može postojati pozitivna ili negativna nula, pa joj bit za znak može biti 0 ili 1. Na svim bitovima kojim je označena frakcija je potrebno da stoje nule, a na bitovima gde je eksponent može biti proizvoljna vrednost. To znači da su treći, četvrti i peti bit kombinacije, kao i poslednjih 20 bitova zapisa jednaki 0.
- Beskonačno se može pojaviti kada se neki veoma veliki broj ne može zapisati ili kao rezultat operacija poput $5 - \infty = -\infty$, $\infty + \infty = \infty$, $-5/0 = -\infty$, i slično. U zavisnosti od toga da li se radi o pozitivnoj ili negativnoj beskonačnosti, bit za znak može biti 0 ili 1. Za beskonačno važi da je početak kombinacije oblika 11110.
- NaN vrednosti se javljaju u izuzetnim situacijama. Postoji signalni NaN (sNaN) i tihi NaN (qNaN). Signalni NaN se javlja pri greškama u inicijalizaciji ili konverziji, a tihi NaN pri greškama u aritmetičkim operacijama. Primeri za pojavu qNaN vrednosti su nedefinisane vrednosti poput $\infty - \infty$, $0 \cdot \infty$, $0/0$, itd. Ukoliko je argument aritmetičke operacije sNaN ili qNaN, rezultat je qNaN. Kod sNaN vrednosti

kombinacija počinje sa 111111, a kod qNaN vrednosti sa 111110. Ostale cifre zapisa mogu biti proizvoljne.

BID kodiranje. Ukoliko je osnova dekadna, a frakcija binarna, IEEE 754 standard koristi BID kodiranje. U nastavku će biti navedeni primeri kodiranja i dekodiranja u jednostrukoj tačnosti. Specijalne vrednosti se zapisuju na isti način kao u slučaju DPD kodiranja.

Neka je, na primer, potrebno odrediti zapis broja -14.37 . Pogodan zapis broja za kodiranje je

$$-14.37 = -1437 \cdot 10^{-2} = (-10110011101)_2 \cdot 10^{-2}.$$

Dakle, broj treba zapisati tako da mu je frakcija celobrojna, a zatim tu frakciju prevesti u binarni sistem. Osnova ostaje sve vreme dekadna. Ukoliko frakcija sadrži manje od 23 bita, potrebno je dopisati vodeće nule, kako bi bila dopunjena do tog broja bitova. U ovom slučaju, zapis frakcije bi bio 0000000000010110011101. Eksponent se i ovde piše na 8 mesta sa uvećanjem 101 i iznosi $-2 + 101 = 99 = (01100011)_2$. Zapis broja se sastoji iz tri dela, jednobitnog znaka, kombinacije od 11 bitova koja se formira nadovezivanjem cifara eksponenta i prve tri cifre frakcije i nastavka frakcije od 20 bitova. Kako je u pitanju negativan broj, zapis je 1 01100011000 00000000010110011101.

Neka je potrebno izvršiti prevođenje broja 9000001. Pogodan zapis za kodiranje je

$$9000001 \cdot 10^0 = (100010010101010001000001)_2 \cdot 10^0.$$

Kada frakcija ima 24 bita, što je ovde slučaj, njena prva tri bita se ingorišu (podrazumevano su uvek jednaki 100). Kombinacija se dobija dodavanjem para bitova 11 na početak, zatim 8 cifara eksponenta i četvrtog bita frakcije. Ostalih 20 bitova se dodaju u nastavku frakcije. Eksponent je $0 + 101 = 101 = (01100101)_2$, pa je kombinacija 11011001010. Broj je pozitivan, pa je njegov zapis 0 11011001010 10010101010001000001.

Izvršimo sada dekodiranje na primeru zapisa 1 01100110000 00000000000000001011. Nakon što se utvrdi da se ne radi ni o kakvoj specijalnoj vrednosti, treba proveriti da li kombinacija počinje sa 11. Ako počinje, razmatra se slučaj dekodiranja gde frakcija ima 24 cifre, a inače slučaj kada frakcija ima 23 cifre. Ovde prvih 8 bitova kombinacije predstavlja eksponent, a poslednja tri početak frakcije. Eksponent je $(01100110)_2 - 101 = 102 - 101 = 1$, a frakcija $(-1011)_2 = -11$. U pitanju je broj $-11 \cdot 10^1 = -110$.

5.5 IEEE 754 standard sa binarnom osnovom

Kod IEEE 754 standarda sa binarnom osnovom, frakcija je binarna, a osnova 2. Zapis se sastoji iz tri dela – zapisa znaka, eksponenta i frakcije. Ovde će biti razmotrena jednostruka i dvostruka tačnost. Kod jednostruke tačnosti broj se kodira sa 32 bita, a kod dvostruke sa 64 bita. IEEE 754 standard sa binarnom osnovom propisuje zapis normalnih i subnormalnih brojeva, kao i specijalnih vrednosti. I u jednostrukoj i u dvostrukoj tačnosti, podržane specijalne vrednosti su 0, beskonačno i NaN vrednost. Nula se piše sa svim nulama u eksponentu i frakciji, a bit za znak može biti 0 ili 1, u zavisnosti od znaka nule. Beskonačno se piše sa svim jedinicama u eksponentu i svim nulama u frakciji. Znak za beskonačno takođe može biti 0 ili 1. Obe NaN vrednosti se pišu sa svim jedinicama u eksponentu. Kod sNaN vrednosti, prvi bit frakcije je 0, a kod ostatka frakcije nisu svi

I u dvostrukoj tačnosti je podržan zapis subnormalnih brojeva. Za dvostruku tačnost, eksponent iznosi -1022 . Na primer, na osnovu zapisa

0 0000000000 0001000,

zaključujemo da se radi o denormalizovanom broju. Na osnovu početka zapisa frakcije 0001 (ostatak su nule, pa nisu od značaja), zaključujemo da je frakcija oblika 0.0001. Broj je pozitivan, pa je vrednost datog denormalizovanog broja $(0.0001)_2 \cdot 2^{-1022} = 2^{-4} \cdot 2^{-1022} = 2^{-1026}$.

5.6 Aritmetičke operacije

U nastavku će biti opisane operacije sabiranja, oduzimanja, množenja i deljenja u IEEE 754 standardu na primeru zapisa sa binarnom osnovom u jednostrukoj tačnosti. Pre izvođenja bilo koje operacije, treba ustanoviti da li je rezultat eventualno neka specijalna vrednost. Ukoliko nije, operacija se može izvršiti odgovarajućim algoritmom. Na kraju izvršene operacije, eventualno može doći do prekoračenja ili potkoračenja, kada se rezultat redom proglašava za beskonačno ili nulu.

Sabiranje. Neka su data dva broja u IEEE 754 standardu u binarnoj osnovi u obliku $f_1 \cdot 2^{e_1}$ i $f_2 \cdot 2^{e_2}$, gde su f_1 i f_2 frakcije, a e_1 i e_2 eksponenti. Za sabiranje je brojeve potrebno dovesti na isti eksponent, a u ovom algoritmu se uzima onaj koji je veći. Pretpostavimo da je $e_1 \geq e_2$ i da je nakon svođenja na eksponent e_1 drugi sabirak oblika $f'_2 \cdot 2^{e_1}$. Sada se sabiranje svodi na

$$f_1 \cdot 2^{e_1} + f_2 \cdot 2^{e_2} = f_1 \cdot 2^{e_1} + f'_2 \cdot 2^{e_1} = (f_1 + f'_2) \cdot 2^{e_1}.$$

Na primer, neka je potrebno izvršiti sabiranje 0 1000010 100110000000000000000000 i 0 1000001 101100000000000000000000. Bit za znak rezultata je 0, budući da se vrši sabiranje dva pozitivna broja. Eksponent rezultata je jednak većem od eksponenata sabiraka i iznosi 1000010. Sada je potrebno sabrati odgovarajuće frakcije. Kako je eksponent rezultata jednak eksponentu prvog sabirka, frakcija prvog sabirka ostaje nepromenjena i iznosi 1.10011. Eksponent drugog sabirka treba povećati za 1, pa je potrebno izvršiti modifikaciju njegove frakcije 1.1011 pomeranjem decimalne tačke za jedno mesto ulevo. Dakle, njegova frakcija je oblika 0.11011, pa je frakcija zbira jednaka $1.10011 + 0.11011 = 10.0111$. Kada se ne dobije normalizovana frakcija, potrebno je izvršiti njenu normalizaciju i u skladu sa tim promeniti vrednost eksponenta. Važi da je $(10.0111)_2 = (1.00111)_2 \cdot 2^1$, pa eksponent rezultata treba uvećati za 1. On sada iznosi 1000011. Zapis zbira je 0 1000011 001110000000000000000000.

Oduzimanje. Neka su data dva broja u obliku $f_1 \cdot 2^{e_1}$ i $f_2 \cdot 2^{e_2}$. Kod oduzimanja je potrebno eksponent umanjioca dovesti na eksponent umanjenika. Pretpostavimo da je nakon svođenja na eksponent e_1 umanjioc oblika $f'_2 \cdot 2^{e_1}$. Sada se oduzimanje svodi na

$$f_1 \cdot 2^{e_1} - f_2 \cdot 2^{e_2} = f_1 \cdot 2^{e_1} - f'_2 \cdot 2^{e_1} = (f_1 - f'_2) \cdot 2^{e_1}.$$

Na primer, neka je potrebno izvršiti oduzimanje 0 1000001 010110000000000000000000 i 0 0111110 001000000000000000000000. Bit za znak razlike je 0, budući da se vrši oduzimanje dva pozitivna broja, gde se oduzima manji broj od većeg. Eksponent razlike je

jednak eksponentu umanjnika. Frakciju umanjioaca treba pomeriti za odgovarajući broj mesta ulevo kako bi se eksponent izjednačio sa eksponentom umanjnika. U ovom slučaju, frakcija umanjioaca 1.001 postaje 0.001001. Decimalnu tačku je potrebno pomeriti tri mesta ulevo zbog toga što se njegov eksponent povećava za 3, budući da je $(10000001)_2 - (01111110)_2 = 3$. Oduzimanjem frakcija se dobija $1.010110 - 0.001001 = 1.001101$. Frakcija razlike je normalizovana, pa je rezultat $0\ 10000001\ 001101000000000000000000$.

Množenje. Neka su data dva broja u obliku $f_1 \cdot 2^{e_1}$ i $f_2 \cdot 2^{e_2}$. Njihovim množenjem se dobija

$$(f_1 \cdot 2^{e_1}) \cdot (f_2 \cdot 2^{e_2}) = (f_1 f_2) \cdot 2^{e_1 + e_2}.$$

Iz poslednjeg sledi da je frakcija proizvoda jednaka proizvodu frakcija, a eksponent proizvoda zbiru eksponenata. Pri sabiranju eksponenata treba voditi računa da su oni u višku 127, koji nakon sabiranja treba anulirati.

Na primer, neka je potrebno izvršiti množenje $0\ 10000011\ 001000000000000000000000$ i $0\ 10000010\ 001100000000000000000000$. Kako se množe dva pozitivna broja, rezultat je takođe pozitivan. Kako su eksponenti oba činioaca zapisani u višku 127, a eksponent proizvoda je potrebno takođe napisati u višku 127, pri sabiranju eksponenata činioaca, da bi se dobio zapis eksponenta rezultata, potrebno je od dobijenog zbira oduzeti 127. Dakle, dobija se:

$$\begin{array}{r} 10000011 \\ + 10000010 \\ \hline 100000101 \\ - 01111111 \\ \hline 10000110 \end{array}$$

Zatim se množenjem frakcija brojeva (postupak množenja je sličan množenju dva dekadna broja) dobija $1.001 \cdot 1.0011 = 1.0101011$. Dobijenu frakciju nije potrebno normalizovati, jer je već svedena na oblik $1.f$. Proizvod je $0\ 10000110\ 010101100000000000000000$.

Deljenje. Neka su data dva broja u obliku $f_1 \cdot 2^{e_1}$ i $f_2 \cdot 2^{e_2}$. Njihovim deljenjem se dobija

$$(f_1 \cdot 2^{e_1}) / (f_2 \cdot 2^{e_2}) = (f_1 / f_2) \cdot 2^{e_1 - e_2}.$$

Iz poslednjeg sledi da je frakcija količnika jednaka količniku frakcija, a eksponent količnika razlici eksponenata. Pri oduzimanju eksponenata treba voditi računa da su oni u višku 127, koji nakon oduzimanja treba dodati.

Na primer, neka je potrebno izvršiti deljenje $0\ 10000100\ 111011100000000000000000$ i $0\ 10000001\ 101000000000000000000000$. Kako se dele dva pozitivna broja, rezultat je takođe pozitivan broj. Kako su eksponenti i deljenika i delioca zapisani u višku 127, a eksponent količnika je potrebno takođe napisati u višku 127, pri oduzimanju eksponenata deljenika i delioca, da bi se dobio zapis eksponenta rezultata, potrebno je dobijenoj razlici dodati 127. Dakle, dobija se:

$$\begin{array}{r}
10000100 \\
- 10000001 \\
\hline
00000011 \\
+ 01111111 \\
\hline
10000010
\end{array}$$

Zatim se deljenjem frakcija brojeva (postupak deljenja je analogan deljenju dva dekadna broja) dobija $1.1110111/1.101 = 1.0011$. Dobijenu frakciju nije potrebno normalizovati, jer je već svedena na oblik $1.f$. Dobijeni količnik je

$$0\ 10000010\ 0011000000000000000000.$$