

3 Caractere și citirea lor

În general, un program operează asupra unor date de intrare, efectuează prelucrări (calcule) și produce ca rezultat niște date de ieșire.

În exemplele de până acum, rezultatele au fost scrise la ieșire, folosind funcția `printf`, dar valorile de intrare prelucrate au fost specificate direct în program, de obicei ca argumente ale funcțiilor care efectuează prelucrarea. Aceasta necesită însă modificarea și recomplarea programului ori de câte ori dorim să calculăm o nouă valoare, de exemplu `fact(4)` în loc de `fact(6)`.

Discutăm în continuare cum se pot scrie programe cu care utilizatorul poate să interacționeze direct, introducând (la tastatură) datele pe care programul să le prelucreze. De exemplu, cum scriem un program care să calculeze și tipărească factorialul unui număr tastat de utilizator?

3.1 Reprezentarea caracterelor

Exprimându-ne precis, utilizatorul nu introduce de fapt un *număr* (care e o noțiune abstractă, matematică), ci *reprezentarea* sa ca o secvență de cifre (aici, în baza 10). Cifrele zecimale, împreună cu literele mari și mici, operatori, semne de punctuație, etc. formează un *alfabet* de *caractere* (simboluri) care pot fi introduse de la tastatură. Pentru ca programele să poată opera cu caracterele introduse, trebuie să definim o altă *reprezentare* a caracterelor, în calculator. Reprezentarea printr-o matrice de pixeli e potrivită pentru afișarea pe ecran, dar nu și pentru memorarea eficientă sau prelucrările de cuvinte (șiruri de caractere) care apar frecvent în programe. Calculatoarele fiind concepute să lucreze în principal cu *numere*, vom reprezenta caracterele prin numere întregi: dispunem toate simbolurile alfabetului într-o anumită ordine și reprezentăm un caracter prin numărul de ordine (poziția) în acest sir.

Cea mai utilizată reprezentare de acest tip este *codul ASCII* (American Standard Code for Information Interchange), care folosește pentru codificare întregii de la 0 la 127 (care se pot reprezenta în baza 2 pe 7 biți). Codul ASCII definește 95 de caractere tipăribile, de la codul 32 (pentru *spațiu*) până la 126; codul 127 (Delete) și cele de la 0 la 31 sunt folosite pentru *caractere de control*, netipăribile, care pot avea înțeles special la afișare sau transmiterea datelor. Alte tabele derivate, standardizate de ISO folosesc coduri până la 255 pentru a reprezenta caractere suplimentare (cu diacritice, alte simboluri, etc.)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	
0x0	\0										\a	\b	\t	\n	\v	\f	\r
0x10:																	
0x20:	!	"	#	\$	%	&	,	()	*	+	,	-	.	/		
0x30:	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?	
0x40:	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
0x50:	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	-	
0x60:	'	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	
0x70:	p	q	r	s	t	u	v	w	x	y	z	{		}	~		

În tabelul de mai sus, am dispus caracterele pe 8 linii și 16 coloane, pentru a obține ușor codul ASCII al fiecărui caracter ca valoare *în baza 16*.

Convențional, *cifrele hexazecimale* cu valori de la 10 la 15 se reprezintă prin literele de la A la F (fie mari, fie mici). În limbajul C, numerele hexazecimale se disting prin prefixul 0x sau 0X. Astfel, din tabel vedem că litera A are codul ASCII 0x41, adică 65. Litera z are codul 0x70 + 0xA, adică $7 \cdot 16 + 10 = 122$.

Observăm că în tabela ASCII, cifrele, litere mari și literele mici sunt reprezentate în trei grupuri contigute, dar separate.

Caracterele fiind reprezentate în calculator ca întregi, când scriem programe care prelucrează caractere, ar trebui să folosim codurile lor. Conform tabelului de mai sus, pentru a verifica dacă un caracter dintr-un sir este C ar trebui să-l comparăm cu numărul 0x43 (67). Scrierea numerică ar fi atât un inconvenient, cât și o posibilă sursă de erori și dependențe de implementare.

Strict vorbind, standardul C99 nu prevede folosirea codificării ASCII, ci doar cerința de a putea reprezenta cifrele, literele mari și mici, majoritatea celorlalte semne grafice uzuale, spațiul și sfârșitul de linie, ca reprezentarea fiecărui caracter să încapă pe un octet, iar cifrele să aibă valori succesive în reprezentare.

În limbajul C, ne putem referi la *valoarea* unui caracter (codul său numeric) inclusiv caracterul între apostroafe, de exemplu: 'c', '7', ':'. Aceste valori se numesc *constante caracter*, al treilea tip de constantă întâlnit până acum, pe lângă cele întregi și reale. Între apostroafe putem încadra orice caracter tipăribil (inclusiv spațiul), dar apostroful însuși trebuie prefixat cu caracterul \ (backslash): '\' pentru a-l deosebi de cele două apostroafe delimitator.

Caracterele de control (netipăribile) uzuale sunt reprezentate prin secvențe speciale formate din \ urmat de o literă care sugerează numele caracterului:

'\a'	alert	'\n'	newline	'\r'	carriage return
'\b'	backspace	'\v'	vertical tab	'\'	apostrof
'\t'	tab	'\f'	form feed	'\\'	backslash

Caracterul \ trebuie dublat pentru a fi interpretat ca atare și nu ca început de secvență specială, o convenție ușuală, folosită și pentru % în printf.

Ca aplicație, scriem o funcție simplă care returnează valoarea (codul) cifrei hexazecimale care reprezintă un număr între 0 și 15. Funcția va returna deci în cod de caracter fie între '0' și '9', fie între 'A' și 'F'. Cifrele zecimale fiind reprezentate consecutiv, codul pentru cifra de valoare n este '0' + n (cu n poziții mai mare decât codul lui 0). Pentru litere, codul va fi cu n - 10 poziții mai mare decât cel al lui A (de exemplu, 'A' + (14 - 10) == 'E').

```
// returneaza cifra hexazecimala pentru 0 <= n < 16, altfel '?'
int xdigit(unsigned n)
{
    return n < 10 ? '0' + n
                  : n < 16 ? 'A' + n - 10
                  : '?';
}
```

Am tratat cazul de eroare (apelul cu un parametru mai mare decât 15) prin returnarea valorii '?', care ar fi evidentă la tipărire. O altă variantă ar fi returnarea valorii -1 care nu reprezintă un caracter. Decizia e la latitudinea proiectantului funcției, dar *cazurile de eroare trebuie tratate* și documentate, împreună cu funcționalitatea, în comentarii.

3.2 Scrierea de caractere

Caracterele fiind reprezentate intern ca întregi, se pune întrebarea cum se poate produce *reprezentarea fizică* a unui caracter pornind de la codul său, adică cum se poate scrie (tipări). Acest lucru e realizat în C de funcția standard

```
int putchar(int c);           // declarata in stdio.h
```

Care scrie caracterul cu valoarea (codificarea numerică) dată la ieșirea standard (în mod implicit, ecranul). Funcția putchar returnează chiar valoarea parametrului (caracterul scris), convenție ușuală la unele funcții standard care, deși au scopul principal de a produce un *efect vizibil* (sau *efect lateral*, noțiune discutată ulterior), returnează și un rezultat utilizabil la nevoie în program.

Apelul putchar('0') va tipări deci cifra 0 (fără apostroafe, folosite doar la reprezentarea constanțelor caracter în program). Nu vom scrie putchar(48), deoarece e mai greu de înțeles și depinde de tabelul de codificare folosit. Apelul putchar(xdigit(12)) (cu funcția xdigit de mai sus) va scrie litera (cifra hexazecimală) C.

Reamintim că o constantă caracter trebuie încadrată între apostroafe; astfel putchar(9) nu scrie cifra 9, ci caracterul cu valoare 9 (în codul ASCII: tab); 'a' + 5 este 'f', dar în expresia a + 5 se presupune că a este un *identificator* (declarat anterior), a cărui valoare se adună cu 5, etc.

3.3 Citirea de caractere

Într-un program C, citirea unui caracter se poate face cu funcția standard

```
int getchar(void);           // declarata in stdio.h
```

Funcția nu are nevoie de parametri (citirea se face de la *intrarea standard*, implicit tastatura). Ea returnează valoarea caracterului citit (codul acestuia).

Orice funcție care interacționează cu mediul extern (în acest caz, intrarea furnizată de utilizator) este supusă la posibilele *situări de eroare* pe care nu le poate controla, și trebuie să poată raporta aceste erori. La citire, e posibil ca intrarea să nu poată furniza un caracter: de exemplu, dacă utilizatorul a încheiat introducerea datelor, sau dacă programul este rulat cu intrarea standard redirectată, citind datele dintr-un fișier în locul tastaturii, și s-a ajuns cu citirea la sfârșitul fișierului.

În caz de eroare, funcția `getchar` returnează o valoare specială identificată prin numele `EOF` ([end of file](#)). Valoarea `EOF` e specificată de standard ca un întreg negativ (pentru a fi diferită de valorile caracterelor, care sunt nenegative), și e definită în fișierul antet `stdio.h` care conține și declarația funcției `getchar`, și a celorlalte funcții standard de intrare/iesire (de ex. `printf`).

În multe implementări, valoarea `EOF` este `-1`, dar nu trebuie să ne bazăm pe această presupunere. Ne referim la această valoare specială doar prin nume.

3.4 Funcții pentru clasificarea caracterelor

În multe probleme dorim să prelucrăm doar caractere de un anumit fel. De exemplu, dacă citim de la intrare un cuvânt, caracter cu caracter, trebuie să ne oprim când caracterul citit nu mai este o literă.

Limbajul C oferă o serie de funcții care testează din ce categorie face parte un caracter. Toate aceste funcții sunt declarate în fișierul antet `ctype.h`. Ele iau ca parametru un întreg (codul caracterului de testat), și returnează un întreg care e: nenul dacă caracterul se încadrează în categoria dorită (corespunzătoare cu numele funcției), sau zero în caz contrar. De exemplu, funcția

```
int isdigit(int c);           // declarata in ctype.h
```

returnează o valoare nenulă dacă parametru este un caracter de la `'0'` la `'9'` și zero pentru orice alt caracter.

Dăm ca exemplu o funcție care returnează valoarea (de la 0 la 9) a unui caracter cifră zecimală, sau `-1` (ca și cod de eroare) pentru orice alt caracter:

```
int digitvalue(int c)
{
    return isdigit(c) ? c - '0' : -1;
}
```

Expresia `c - '0'` numără câte cifre sunt între `c` și `0` în tabela de caractere, ceea ce e tocmai valoarea cifrei `c` (de exemplu, `'7' - '0' == 7`). Este tocmai calculul invers celui făcut în funcția `xdigit`.

Acvest exemplu necesită o clarificare privind valorile logice în limbajul C. Ca și condiție pentru operatorul `? :` am folosit pâna acum doar expresii (cum ar fi comparațiile) care pot da rezultatele logice *adevărat* sau *fals*. În limbajul C, nu există tip special pentru valori logice. Operatorii care furnizează valori logice (cum ar fi relaționali și de comparație, `<`, `!=`, etc.) produc de fapt o valoare întreagă: 1 pentru rezultat *adevărat* și 0 pentru *fals*. Acolo unde se *solicită* o valoare logică (de exemplu în expresia condițională), se poate folosi mai general *orice* expresie numerică: orice valoare *nenulă* (nu doar 1) semnifică *adevărat*, iar valoarea 0 înseamnă *fals*.

Lista funcțiilor declarate în `ctype.h` cuprinde:

```
int isalnum(int c);   // este litera sau cifra?
int isalpha(int c);  // este litera? (mare sau mica)
int isdigit(int c);  // este cifra?
int islower(int c);  // este litera mica?
int isspace(int c);  // spatiu alb? (spatiu + \t\n\v\f\r)
int isupper(int c);  // este litera mare?
int isxdigit(int c); // este cifra hexa? (0-9, A-F, a-f)

int iscntrl(int c);  // este caracter de control?
int isgraph(int c);  // caracter tiparibil? (excl. spatiu)
int isprint(int c);  // caracter tiparibil? (incl. spatiu)
int ispunct(int c);  // este semn de punctuatie?
```

Toate aceste funcții acceptă ca parametru orice întreg din domeniul de valori care reprezintă caracter, sau valoarea EOF (e nedefinit însă comportamentul dacă parametrul depășește acest domeniu de valori). În consecință, rezultatul returnat de `getchar` poate fi testat direct cu aceste funcții fără a trebui verificată separat condiția de eroare/sfârșit de fișier. Evident, EOF, nefind caracter nu intră în nicicare din clasele de mai sus.

Atenție! Funcțiile de mai sus au ca rezultat 0 sau o valoare nenulă, nu 0 sau 1 ca și rezultatul comparațiilor. Ele se folosesc *direct* în condiții: `isalnum(...)` ? ... după cum sugerează și numele (este caracter alfanumeric?) și nu comparând cu 1: `isalnum(...) == 1` ? ... deoarece rezultatul *adevărat* ar putea să fie reprezentat prin altă valoare nenulă.

Se recomandă folosirea acestor funcții standard în locul testelor bazate pe compararea valorii caracterelor (cum ar fi verificarea că un caracter e literă mare prin încadrarea între 'A' și 'Z'). Un motiv e eficiența implementării (discutată ulterior), dar mai importantă e corectitudinea. Funcțiile din `ctype.h` sunt concepute să țină cont de definirea diferită a alfabetului de litere în funcție de mediul internaționalizat în care e rulat programul, pe când comparația limitată la literele latine ar putea da erori prin omisiune în acest caz.

Antetul `ctype.h` declară și *funcții de conversie* între litere mici și mari:

```
int tolower(int c); // returnează litera mica corespondentă
int toupper(int c); // returnează litera mare corespondentă
```

Dacă parametrul este literă mare (resp. mică) și are în alfabetul mediului de rulare un corespondent literă mică (resp. mare), se returnează acesta, altfel se returnează chiar valoarea neschimbată a argumentului.

3.5 Citirea unui număr

Scriem o funcție care citește de la intrare un număr natural reprezentat în baza 10. Calculăm valoarea numărului pas cu pas, odată cu citirea fiecărei cifre. Dacă numărul are n cifre c_1, c_2, \dots, c_n (de la stânga la dreapta) putem forma succesiv numerele: $r_0 = 0$, $r_k = r_{k-1} \cdot 10 + c_k$, pentru $k = 1, 2, \dots, n$. De exemplu, pentru numărul 247 cu $n = 3$ cifre, avem: $r_0 = 0$, $r_1 = 0 \cdot 10 + 2 = 2$, $r_2 = 2 \cdot 10 + 4 = 24$, $r_3 = 24 \cdot 10 + 7 = 247$ care e valoarea dorită.

Transcriem relația recursivă de mai sus într-o funcție care ia doi parametri: rezultatul parțial deja calculat (r_{k-1}) și ultimul caracter citit (care dă cifra c_k). Ramura de oprire corespunde unui caracter care nu e cifră, caz în care se returnează rezultatul parțial deja calculat. Altfel se returnează rezultatul apelului recursiv cu rezultatul parțial actualizat cu cifra citită, și un caracter nou.

Dispunem deja de funcția `getchar` pentru citirea unui caracter, și am văzut că valoarea numerică reprezentată de caracterul cifră zecimală `c` este `c - '0'` (deplasamentul față de cifra 0 în tabela de caractere). Cum nu știm dinainte câte cifre va avea numărul introdus de utilizator, folosim funcția `isdigit` pentru a testa când caracterul citit nu mai este cifră și a ne opri.

```
unsigned readnat_rc(unsigned r, int c)
{
    return isdigit(c) ? readnat_rc(r*10 + (c-'0'), getchar()) : r;
}
```

Pentru utilizare, această funcție trebuie apelată cu valori inițiale corespunzătoare pentru parametri. Rezultatul parțial inițial este $r_0 = 0$, iar primul caracter e citit cu `getchar`, deci apelul se scrie `readnat_rc(0, getchar())`.

Din punct de vedere al utilizatorului, citirea unui număr natural de la intrarea standard nu ar avea nevoie de informații suplimentare: funcția ar trebui să fie fără parametri, la fel ca și `getchar` pentru citirea unui caracter. Scriem de aceea o funcție care încapsulează apelul cu valorile inițiale necesare stabilite și este astfel mai natural și ușor de folosit:

```
unsigned readnat(void) { return readnat_rc(0, getchar()); }
```

Continuăm implementând citirea unui întreg precedat optional de semn.

```
int readint_c(int c) // tine cont de semn; c: primul caracter
{
    return c == '-' ? - readnat() :
        c == '+' ? readnat() : readnat_rc(0, c);
}
```

Dacă primul caracter citit este semn, funcția apelează `readnat`, fără parametri, care va începe citirea cu un caracter nou. Altfel, se apelează `readnat_rc` cu valoarea inițială 0, și caracterul tocmai citit. Își pentru această funcție scriem o versiune fără parametri, care citește singură primul caracter necesar:

```
int readint(void) { return readint_c(getchar()); }
```

3.6 Refacerea unui caracter din intrare

Funcțiile scrise mai sus se opresc din citirea numărului la primul caracter necorespunzător. Acesta este însă consumat din intrare, și se pierde, nemaifiind disponibil după revenirea din funcție. Acest aspect nedorit trebuie evitat; funcția ar trebui să citească doar cât are nevoie, fără a afecta citirile ulterioare.

Limbajul C pune la dispoziție o funcție standard pentru a “pune înapoi” un caracter care a fost citit dintr-un flux de intrare. Declarația ei este:

```
int ungetc(int c, FILE *stream); // in stdio.h
```

Vom prezenta ulterior tipul `FILE *` al parametrului 2, reprezentând un fișier, în acest caz de intrare. Pentru moment dorim să punem înapoi un caracter în *intrarea standard*, de unde am efectuat citirea. Pentru aceasta, argumentul al doilea va fi `stdin`, un identificator declarat în `stdio.h` care reprezintă fișierul standard de intrare. Deci, pentru a pune înapoi în intrarea standard un caracter, reprezentat prin valoarea sa întreagă `int c`, vom folosi apelul

```
ungetc(c, stdin); // pune înapoi caracterul c
```

Funcția returnează valoarea caracterului dat în caz de succes, sau `EOF` pentru a semnala eroare. Caracterul pus înapoi devine următorul care va fi citit din fluxul de intrare. Limbajul C garantează că un caracter poate fi pus înapoi în acest fel înainte de a fi consumat printr-un apel de citire, cum ar fi `getchar()`.

3.7 Secvențierea pentru expresii

Reexaminând funcția de citire a unui număr natural, caracterul trebuie pus înapoi în intrare atunci când nu e cifră (`isdigit(c)` indică “fals”), deci pe a doua ramură a expresiei condiționale. În scrierea curentă, pe această ramură expresia ia valoarea rezultatului parțial `r` primit ca parametru.

```
isdigit(c) ? readnat_rc(r*10 + (c-'0'), getchar()) : r;
```

Cu elementele de limbaj prezentate pâna acum, în locul unei expresii nu se pot efectua *două* acțiuni: întâi apelul funcției `ungetc` și apoi să-i dăm expresiei valoarea `r`. Limbajul C permite însă execuția succesivă a două *instrucțiuni*. E normal să existe noțiunea de *secvențiere* și pentru expresii, mai ales că instrucțiunea cea mai simplă este chiar o evaluare de expresie.

În C, evaluarea secvențială a două expresii se face cu operatorul virgulă ,
expresie1 , expresie2

Se evaluatează întâi `expresie1`, iar rezultatul e ignorat: construcția este deci utilă dacă expresia are un *efect lateral* (un apel de intrare-ieșire sau o atribuire). Rezultatul și tipul întregii expresii se obțin evaluând `expresie2`.

Operatorul virgulă are cea mai scăzută precedență dintre operatori; pentru a-l folosi în ultimul operand al expresiei condiționale grupăm deci expresia între paranteze. Rescriem astfel citirea unui număr natural obținând o funcție care nu consumă caractere suplimentare din intrare.

```
unsigned readnat_rc(unsigned r, int c)
{
    return isdigit(c) ? readnat_rc(r*10 + (c-'0'), getchar())
                      : (ungetc(c, stdin), r);
}
```