

Acesta este capitolul 7 — *Codificări de interes practic* — al ediției electronice a cărții *Rețele de calculatoare*, publicată la Casa Cărții de Știință, în 2008, ISBN: 978-973-133-377-9.

Drepturile de autor aparțin subsemnatului, Radu-Lucian Lupșa.

Subsemnatul, Radu-Lucian Lupșa, acord oricui dorește dreptul de a copia conținutul acestei cărți, integral sau parțial, cu condiția atribuirii corecte autorului și a păstrării acestei notițe.

Cartea, integrală, poate fi descărcată gratuit de la adresa <http://www.cs.ubbcluj.ro/~rlupsa/works/retele.pdf>

Capitolul 7

Codificări de interes practic

7.1. Probleme privind reprezentarea numerelor întregi

Până aici am privit un mesaj transmis printr-o rețea ca un șir de simboluri dintr-un alfabet finit. Între simbolurile ce alcătuiesc șirul se stabilește o *ordine*, existând un prim element al șirului, un al doilea, etc. Transmiterea elementelor se face în ordinea în care apar ele în șir, primul simbol transmis fiind cel care ocupă prima poziție din șir. Până aici am considerat că transmiterea unui șir de la un dispozitiv la altul conservă ordinea între elemente.

Așa cum vom vedea însă în paragraful de față, din rațiuni legate de standardizarea reprezentării numerelor întregi, transmiterea unui șir nu conservă întotdeauna ordinea elementelor. În cele ce urmează, vom examina interferențele între reprezentarea numerelor în calculator și ordinea simbolurilor ce alcătuiesc un mesaj. Vom oferi cititorului o perspectivă, inspirată din [Cohen 1980] și mai puțin întâlnită în alte lucrări, asupra relațiilor dintre biți, octeți și reprezentarea numerelor.

7.1.1. Reprezentări pe biți

În paragraful de față vom face abstracție de anumite complicații constructive ale sistemelor de calcul reale. Vom considera reprezentări pe biți, ignorând deocamdată aspectele legate de gruparea biților în octeți și de faptul că, în memoria calculatorului, adresele identifică octeți și nu biți.

Ca urmare, rugăm cititorul să uite, pentru moment, noțiunea de *octet* (*byte*).

7.1.1.1. Bitul

Pentru reprezentarea diverselor date, alegerea unui alfabet cu două simboluri este avantajoasă din două motive. Pe de o parte, este cel mai mic alfabet posibil, ca urmare alegerea unui alfabet cu două elemente aduce o anumită simplitate și naturalețe construcției matematice. Pe de altă parte, din punct de vedere practic, al construcției echipamentelor fizice, dispozitive cu două stări stabile sunt mult mai ușor de construit decât dispozitive cu mai multe stări.

În scris, cele două simboluri sunt notate în mod obișnuit cu 0 și 1. Atragem atenția că:

- Alegerea celor două simboluri utilizate, precum și a corespondenței dintre starea dispozitivului fizic și simbolul asociat, poate fi făcută oricum. Odată însă făcută o alegere, aceasta trebuie respectată de toate entitățile implicate în comunicație.
- Numai în unele cazuri simbolurile au rol de cifră, adică au asociate valori numerice. Valoarea numerică a unui simbol este importantă doar dacă simbolul este interpretat ca număr sau intră în reprezentarea unui număr. În restul cazurilor, este important doar să existe simboluri distincte.

Un simbol dintr-un alfabet cu două elemente se numește *bit*, de la *binary digit* (rom. *cifră binară*).

7.1.1.2. Șiruri de biți

În cadrul unui șir de biți, avem nevoie să identificăm fiecare bit al șirului. Pentru aceasta, se stabilește o *ordine* a biților: avem un prim bit, un al doilea bit, etc. Trebuie remarcat însă că ordinea este o convenție: nu există o legătură directă între ordinea convențională a unui șir de biți și amplasamentul dispozitivelor fizice care memorează acei biți. Numărul de ordine al unui bit, în cadrul acestei ordini convenționale, se numește în mod obișnuit *poziția* (sau, eventual, *adresa* sau *deplasamentul*) bitului. Numerotarea pozițiilor se face de obicei începând de la 0 sau de la 1; în cele ce urmează vom utiliza numerotarea de la 0.

La transmiterea unui șir de biți, este natural ca primul bit transmis, considerând ordinea cronologică, să fie primul bit al șirului (poziția 0). La memorarea unui șir într-o memorie cu acces direct (memorie RAM sau fișier pe disc), este natural ca în celula cu adresa cea mai mică, dintre celulele alocate șirului, să fie plasat primul bit al șirului (bitul de pe poziția 0). În acest fel, *primul bit* al unui șir înseamnă, simultan, bitul de pe poziția (convențională) 0, bitul transmis primul (cronologic) și bitul memorat la adresa cea mai mică.

7.1.1.3. Reprezentarea pe biți a numerelor întregi

Reprezentarea numerelor naturale prin șiruri de biți se bazează pe ceea ce matematicienii numesc *reprezentare pozițională în baza 2*, pe care o presupunem cunoscută. În reprezentarea într-o bază de numerație, distingem cifra unităților, având ponderea $2^0 = 1$, cifra de pondere $2^1 = 2$ (în baza zece s-ar numi cifra zecilor; pentru baza 2 nu avem un nume), cifra de pondere $2^2 = 4$, etc.

Există două alegeri posibile cu privire la legătura dintre ponderile cifrelor în număr și pozițiile lor în șir:

1. *Primul bit al șirului este cifra de pondere maximă.* Aceasta alegere este identică celei utilizate în scrierea numerelor în limbile „obișnuite“ (indo-europene), cu scriere de la stânga spre dreapta. Se mai numește *big endian*.

În această schemă de reprezentare, un șir de biți $b_0b_1\dots b_{n-1}$ reprezintă numărul

$$b_0 \cdot 2^{n-1} + b_1 \cdot 2^{n-2} + \dots + b_{n-1} \cdot 2^0.$$

2. *Primul bit al șirului este cifra de pondere 1.* Această reprezentare este asemănătoare scrierii numerelor în limbile semite (araba și ebraica, cu scriere de la dreapta spre stânga și unde numerele sunt scrise tot cu cifra unităților în dreapta). Se mai numește *little endian*.

Această alegere are avantajul unei scrieri mai simple a relației dintre valoarea numărului și șirul de biți: valoarea unui număr reprezentat pe n biți este

$$b_0 \cdot 2^0 + b_1 \cdot 2^1 + \dots + b_{n-1} \cdot 2^{n-1}.$$

Fiind două scheme de reprezentare distincte, dacă un sistem transmite un număr în reprezentare *little endian*, iar celălalt sistem interpretează șirul de biți primit ca fiind într-o reprezentare *big endian*, receptorul „înțelege“ alt număr decât cel transmis de emițător. Ca urmare, orice protocol care specifică transmitere binară a numerelor trebuie să precizeze dacă se utilizează o reprezentare *little endian* sau una *big endian*.

EXEMPLUL 7.1: Fie șirul de biți 11001, în care am scris primul bit (în sensul din § 7.1.1.2) pe poziția cea mai din stânga.

Dacă acest șir este reprezentarea *big endian* a unui număr, numărul respectiv este 25. Dacă reprezentarea a fost făcută în format *little endian*, numărul este 19.

Este important de remarcat că distincția dintre schema de reprezentare *big endian* și schema *little endian* există numai acolo unde pe de o parte avem o ordine a biților dată de adresele lor în memorie sau de ordinea cronologică la transmiterea lor prin mediul fizic, iar pe de altă parte fiecare bit are o anumită pondere în reprezentarea unui număr întreg.

7.1.2. Reprezentări pe octeți

În paragraful precedent, am ignorat în mod deliberat noțiunea de *octet (byte)* și am presupus că, în memorie, fiecare bit ar avea o adresă individuală. Vom studia, în continuare, problemele legate de gruparea, în cadrul sistemelor de calcul reale, a biților în octeți și de faptul că, în general, ordinea biților în octeți nu este vizibilă programatorului.

7.1.2.1. Octeți

În memoria calculatoarelor, biții sunt grupați în grupuri de dimensiune fixă, în mod obișnuit câte 8 biți. Un astfel de grup se numește *octet* (denumire intrată pe filieră franceză) sau *bait* (adaptare a englezescului *byte*).

Un octet poate fi privit în două moduri distincte:

- ca un șir de 8 biți,
- ca un număr întreg cuprins între 0 și 255.

Echivalența între aceste două moduri de-a privi un octet este o problemă ce necesită multă atenție. Dacă identificăm biții după ponderile lor, există o corespondență biunivocă între un astfel de grup de 8 biți și un număr întreg între 0 și 255. Dacă însă identificăm biții după poziția lor în șirul de 8 biți (bitul 0, bitul 1, ..., bitul 7), atunci corespondența biunivocă între număr și șir de biți există doar după ce am stabilit o corespondență între poziția unui bit în șir și ponderea sa (*big endian*, *little endian* sau eventual o corespondență mai complicată).

După modul în care se identifică biții unui octet în definiția operațiilor efectuate de diverse componente ale unui sistem de calcul, operațiile se pot împărți în trei categorii:

1. *Operații pentru care biții se identifică după pondere* sau, echivalent, octetul este privit ca număr. Aici se încadrează operațiile aritmetice și operațiile de deplasare pe biți. De remarcat că operațiile *deplasare la stânga* (engl. *shift left*), respectiv *deplasare la dreapta* (engl. *shift right*) pot fi descrise în termeni de operații aritmetice: deplasarea la stânga este o înmulțire cu 2, iar deplasarea la dreapta este o împărțire la 2. În acest context, „spre stânga” și „spre dreapta” înseamnă spre pozițiile cu pondere mai mare, respectiv mai mică, neavând nici o legătură cu

primele sau cu ultimele poziții. Aceste operații sunt efectuate de unitatea aritmetică din microprocesorul calculatorului.

2. *Operații pentru care biții sunt identificați după numărul lor de ordine* sau, echivalent, octetul este privit ca un șir arbitrar de biți, fără a avea asociată o valoare numerică. Aici intră transmiterea bit cu bit (transmitere serială) a unui octet. Această operație este efectuată de placa de rețea și de alte adaptoare seriale (de exemplu, adaptoarele USB). Tot aici s-ar încadra, dacă ar exista, o operație de obținere sau de modificare a unui bit (al octetului) identificat prin numărul său de ordine. O asemenea operație nu este oferită, în mod normal, într-un sistem de calcul — nu există o instrucțiune care să extragă, de exemplu, bitul numărul 5 dintr-un octet.
3. *Operații care pot fi definite fie identificând biții după ponderea lor, fie identificând biții după numărul lor de ordine.* În această categorie se încadrează transmiterea unui octet ca un tot unitar, verificarea egalității a doi octeți și operațiile logice pe bit (*și, sau, sau exclusiv* și negația).

Pentru oricare dintre aceste operații, dacă definim operația identificând biții după numărul lor de ordine, efectul ei asupra valorii numerice a octetului nu depinde de corespondența aleasă între pozițiile biților și ponderile lor.

În aceste condiții, în interiorul unui calculator, biții din cadrul unui octet sunt identificați după ponderea lor. În construcția calculatorului, proiectantul are grijă ca atunci când un bit având, într-un modul al calculatorului, o anumită pondere este transferat către alt modul al calculatorului, să ajungă acolo pe o poziție cu aceeași pondere.

La transmisia unui octet între două sisteme de calcul, mecanismele de transmisie sunt astfel construite încât să transmită valoarea octetului. Întrucât, prin mediul fizic al rețelei, biții sunt transmiși secvențial, biții unui octet sunt aici identificați prin numărul lor de ordine în cadrul transmisiei. Pentru a păstra valoarea octetului în timpul transmisiei prin mediul rețelei, corespondența dintre numărul de ordine al unui bit și ponderea sa (*little endian* sau *big endian*) trebuie să facă parte din specificațiile protocolului de nivel fizic al rețelei.

Numerotarea biților unui octet intervine, de asemenea, în descrierea unor scheme de reprezentare a datelor unde un număr este reprezentat pe un grup de biți ce nu formează un număr întreg de octeți. Este cazul schemelor de reprezentare pentru structuri de date ce conțin câmpuri de 1 bit, 2 biți, 12 biți, etc. Și aici este necesar să se specifice dacă numerotarea biților este *little*

endian sau *big endian*. Mai multe detalii despre astfel de reprezentări vor fi studiate în § 7.1.2.4.

7.1.2.2. Șiruri de octeți

Ca și în cazul biților (vezi § 7.1.1.2), și cu octeții putem construi șiruri. În cadrul unui șir de octeți, octeții sunt așezați într-o ordine, existând un prim octet (numerotat ca octetul 0), al doilea octet (poziția 1), etc. La transmisia printr-o legătură în rețea, primul octet al șirului este, cronologic, primul octet transmis. La memorare, primul octet este cel memorat la adresa cea mai mică.

În virtutea celor două moduri de-a privi un octet, un șir de n octeți poate fi, la rândul lui, privit ca:

- un șir de $8n$ biți,
- un șir de n numere, fiecare cuprins între 0 și 255.

Pentru a putea privi un șir de n octeți ca un șir de $8n$ biți, este necesar să avem o numerotare, bine definită, a biților în cadrul unui octet. Rezultă un șir de biți în care între poziția p_B a unui bit în șirul de $8n$ biți, poziția p_{BO} a bitului în cadrul octetului în care se găsește și poziția p_O a celui octet în șirul de octeți are loc relația:

$$p_B = 8 \cdot p_O + p_{BO}. \quad (7.1)$$

Relația de mai sus are această formă simplă dacă se utilizează numerotare de la 0; pentru numerotarea de la 1, forma relației e mai complicată.

Transmiterea unui șir de octeți printr-o conexiune, precum și memorarea șirului într-un fișier pe disc urmată de citirea lui înapoi în memorie, păstrează ordinea și valorile octeților din șir. Valorile octeților sunt păstrate dacă privim octeții ca numere între 0 și 255; dacă privim octeții ca șiruri de 8 biți, valorile octeților se păstrează numai dacă pe ambele sisteme utilizăm aceeași corespondență între numerele de ordine și ponderile biților.

7.1.2.3. Reprezentarea numerelor pe un număr întreg de octeți

Cel mai mare număr ce poate fi reprezentat pe un octet este 255, ceea ce este mult prea puțin pentru majoritatea aplicațiilor. Pentru a putea reprezenta numere din intervale mai largi, sunt necesare scheme de reprezentare pe mai mult de 8 biți. Schemele cele mai simple sunt cele care utilizează un număr întreg de octeți; acestea vor fi prezentate în continuare. Schemele de reprezentare ce utilizează șiruri de biți ce nu formează neapărat un număr întreg de octeți vor fi studiate în § 7.1.2.4.

Deoarece un octet are valoarea între 0 și 255, putem considera fiecare octet ca fiind o cifră în baza 256. Reprezentarea unui număr printr-un șir de octeți se face ca reprezentare pozițională în baza 256. Există două reprezentări posibile:

- *little endian*: primul octet are ponderea 1, al doilea octet are ponderea 256, al treilea octet are ponderea $256^2 = 65536$, etc.
- *big endian*: primul octet are ponderea 256^{n-1} (unde n este numărul de octeți ai reprezentării), al doilea octet are ponderea 256^{n-2} ș. a. m. d., penultimul octet are ponderea 256, iar ultimul octet are ponderea 1.

Reamintim că prin *primul octet* înțelegem octetul care este transmis primul, în ordine cronologică, de la un dispozitiv la altul și, totodată, octetul memorat la adresa cea mai mică.

EXEMPLUL 7.2: Descriem mai jos reprezentarea numărului 300 în schemele de reprezentare *little endian* și *big endian*, pe 2 și pe 4 octeți. Pentru fiecare dintre aceste patru scheme de codificare, este dat șirul de octeți ce reprezintă numărul 300.

poziție octet (nr. ordine)	Valorile octeților pentru diverse reprezentări			
	2 octeți big endian	2 octeți little endian	4 octeți big endian	4 octeți little endian
0	1	44	0	44
1	44	1	0	1
2	—	—	1	0
3	—	—	44	0

De exemplu, în cadrul reprezentării pe 2 octeți în format *big endian*, valoarea numărului reprezentat se regăsește ca valoarea octetului 0 înmulțită cu 256 plus valoarea octetului 1, anume: $1 \cdot 256 + 44 = 300$. În cadrul reprezentării pe 4 octeți în format *little endian*, octetul 0 are ponderea $256^0 = 1$, octetul 1 are ponderea $256^1 = 256$, octetul 2 are ponderea $256^2 = 65536$, iar octetul 3 are ponderea $256^3 = 16218368$. Valoarea numărului reprezentat este

$$44 \cdot 1 + 1 \cdot 256 + 0 \cdot 256^2 + 0 \cdot 256^3 = 300$$

Unitatea aritmetică a calculatorului poate efectua operații aritmetice asupra numerelor reprezentate pe 2 octeți sau, pentru unele calculatoare, pe 4 sau 8 octeți. Pe unele calculatoare, unitatea aritmetică lucrează cu numere reprezentate după sistemul *big endian*; pe alte calculatoare, unitatea

aritmetică cere reprezentare *little endian*. După acest criteriu, calculatoarele ale căror unități aritmetice lucrează cu numere reprezentate pe mai mult de un octet se împart în calculatoare *little endian* și calculatoare *big endian*.

Variabilele de tip întreg, în majoritatea limbajelor de programare, sunt reprezentate pe 2, 4 sau 8 octeți, în ordinea fixată de unitatea aritmetică.

Este posibilă utilizarea, pentru anumite variabile întregi, a unei reprezentări diferite de cea a unității aritmetice. De asemenea, se pot utiliza reprezentări pe mai mulți octeți decât permite unitatea aritmetică. La manipularea acestor variabile, programatorul trebuie să aibă în vedere că operațiile aritmetice „normale“ fie nu pot fi executate deloc, fie nu se efectuează corect asupra lor. Astfel de numere se manipulează, de obicei, prelucrând separat fiecare octet.

La memorarea pe disc sau la transmiterea printr-o conexiune, trebuie stabilit printr-un standard dacă se utilizează un format *little endian* sau *big endian*, precum și numărul de octeți pe care se reprezintă fiecare număr memorat sau, respectiv, transmis. Majoritatea protocoalelor pentru Internet prevăd formate *big endian* pentru numerele întregi transmise. Multe dintre formatele de fișiere prevăd formate *little endian*. Există și formate (de exemplu, formatul TIFF pentru imagini, formatele UTF-16 și UTF-32 pentru texte) care permit emițătorului să aleagă formatul dorit și prevăd un mecanism prin care emițătorul informează receptorului despre alegerea făcută.

Dacă formatul de pe disc sau de pe conexiune coincide cu formatul unității aritmetice locale, un program poate transfera direct șiruri de octeți între o variabilă întregă locală și fișierul de pe disc sau, respectiv, conexiunea spre celălalt calculator. Dacă formatul de pe disc sau de pe conexiune este invers față de cel local, un program care transferă date trebuie să inverseze ordinea octeților imediat înainte de scrierea pe disc sau de trimiterea pe conexiune, precum și imediat după citirea de pe disc sau recepționarea de pe conexiune.

7.1.2.4. Reprezentarea numerelor pe un șir arbitrar de biți

Ne vom ocupa în continuare de metode de reprezentare, pentru numere întregi, în care biții ce intră în reprezentarea unui număr nu formează neapărat un număr întreg de octeți. O astfel de schemă este o generalizare a schemei prezentate în paragraful precedent.

O astfel de metodă de reprezentare se bazează pe reprezentarea numerelor în baza 2 (vezi și § 7.1.1.3). Pentru ca o astfel de schemă să fie complet definită, este necesar să fie stabilită (standardizată) corespondența dintre poziția fiecărui bit din reprezentare și ponderea asociată. În descrierea

unei astfel de scheme de reprezentare, trebuie precizate trei lucruri:

- dacă reprezentarea numărului prin șirul de biți se face după metoda *big endian* sau *little endian*;
- dacă numerotarea biților în cadrul fiecărui octet se face începând de la bitul de pondere 1 (adică valoarea octetului este reprezentată după schema *little endian*) sau începând de la bitul de pondere 128 (adică valoarea octetului este reprezentată după schema *big endian*).
- dacă numerotarea biților în cadrul șirului de biți se face după relația (7.1) sau după o altă metodă.

Într-o schemă de reprezentare „rațională“, la primele două puncte se utilizează fie formatul *big endian* pentru amândouă, fie formatul *little endian* pentru amândouă, iar la punctul al treilea se utilizează relația (7.1).

Rezultă astfel două metode coerente:

- *big endian*. În această metodă, numerotarea biților începe cu cel mai semnificativ bit al primului octet, iar după cel mai puțin semnificativ bit al primului octet urmează cel mai semnificativ bit al celui de-al doilea octet. Orice număr, indiferent pe câți biți s-ar reprezenta, se reprezintă în format *big endian*.
- *little endian*. În această metodă, numerotarea biților începe cu cel mai puțin semnificativ bit al primului octet, iar după cel mai semnificativ bit al primului octet urmează cel mai puțin semnificativ bit al celui de-al doilea octet. Orice număr, indiferent pe câți biți s-ar reprezenta, se reprezintă în format *little endian*.

În cadrul acestor două metode, dacă reprezentăm un număr folosind un șir de biți ce formează un număr întreg de octeți, formatele de reprezentare rezultate coincid cu formatele de reprezentare pe octeți studiate în paragraful precedent.

EXEMPLUL 7.3: Considerăm o schemă de reprezentare pentru două numere întregi, a și b , în care a se reprezintă pe 4 biți și b se reprezintă pe 12 biți. În total avem $4 + 12 = 16$ biți, adică 2 octeți.

Dacă alegem metoda *big endian*, schema de reprezentare va utiliza cei mai semnificativi 4 biți ai primului octet pentru a-l reprezenta pe a , ceilalți 4 biți ai primului octet vor fi cei mai semnificativi 4 biți ai lui b , iar cel de-al doilea octet va conține cei mai puțin semnificativi 8 biți din b . Această schemă de reprezentare este ilustrată în figura 7.1, cu valori concrete $a = 11$ și $b = 300$.

Dacă alegem metoda *little endian*, schema de reprezentare va utiliza cei mai puțin semnificativi 4 biți ai primului octet pentru a-l reprezenta pe a , ceilalți 4 biți ai primului octet vor fi cei mai puțin semnificativi 4 biți ai lui b ,

iar cel de-al doilea octet va conține cei mai semnificativi 8 biți din b . Această schemă de reprezentare este ilustrată în figura 7.2, cu valori concrete $a = 11$ și $b = 300$.

b_0	b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9	b_{10}	b_{11}	b_{12}	b_{13}	b_{14}	b_{15}
1	0	1	1	0	0	0	1	0	0	1	0	1	1	0	0

(a) Reprezentarea privită ca șir de biți

Nr. octet	Valoare (binar)								Valoare (zecimal)
	c_0	c_1	c_2	c_3	c_4	c_5	c_6	c_7	
0	1	0	1	1	0	0	0	1	177
1	0	0	1	0	1	1	0	0	44

(b) Valorile octeților. La scrierea valorilor biților în octet (coloana din mijloc) s-a utilizat convenția obișnuită, de-a scrie cifrele mai semnificative în stânga.

Figura 7.1: Reprezentare *big endian* pentru numărul 11 pe 4 biți urmat de numărul 300 reprezentat pe 12 biți (exemplul 7.3).

b_0	b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9	b_{10}	b_{11}	b_{12}	b_{13}	b_{14}	b_{15}
1	1	0	1	0	0	1	1	0	1	0	0	1	0	0	0

(a) Reprezentarea privită ca șir de biți

Nr. octet	Valoare (binar)								Valoare (zecimal)
	c_7	c_6	c_5	c_4	c_3	c_2	c_1	c_0	
0	1	1	0	0	1	0	1	1	203
1	0	0	0	1	0	0	1	0	18

(b) Valorile octeților. La scrierea valorilor biților în octet (coloana din mijloc) s-a utilizat convenția obișnuită, de-a scrie cifrele mai semnificative în stânga.

Figura 7.2: Reprezentare *little endian* pentru numărul 11 pe 4 biți urmat de numărul 300 reprezentat pe 12 biți (exemplul 7.3).

Un alt exemplu în care avem de-a face cu numere reprezentate pe șiruri arbitrare de biți este legat de așa-zisa *codificare în baza 64*, descrisă în § 7.4.2.

7.1.3. Probleme privind reprezentarea lungimii șirurilor

Oridecâteori se transmite un șir de obiecte, este necesar ca receptorul să poată determina numărul de obiecte transmise. Acest lucru este valabil

indiferent de natura obiectelor: biți, octeți, cifre zecimale, caractere ale unui text, numere în cadrul unui șir de numere.

Există trei metode de a face ca receptorul să poată determina numărul de obiecte ce-i sunt transmise:

- *Numărul de obiecte este fixat.* În acest caz, indiferent de valorile datelor ce se transmit, numărul de obiecte transmise este același. Metoda este utilizată frecvent la memorarea unui șir în memoria RAM sau pe disc, deoarece permite alocarea de la început a memoriei pentru reprezentarea lui și permite accesul direct la datele memorate după șirul în discuție. Dezavantajul principal al metodei este acela că dimensiunea fixată trebuie astfel aleasă încât să fie suficientă în orice caz ce poate să apară la execuție. De asemenea, trebuie să existe o valoare potrivită pentru pozițiile „neutilizate“ din șir; de exemplu, în reprezentarea numerelor, pozițiile cele mai semnificative se completează cu zerouri.

Deoarece la transmisia printr-o conexiune nu se poate pune problema accesului direct (adică altfel decât secvențial) la date, metoda este puțin utilizată în transmisia datelor prin rețea. Șiruri de lungime fixă se utilizează la reprezentarea binară a numerelor.

- *Numărul de obiecte este transmis separat, în fața șirului.* Această metodă ușurează munca receptorului, care știe exact ce să aștepte și poate alocă memorie pentru recepționarea datelor. În schimb, munca emițătorului este complicată prin faptul că acesta trebuie să cunoască de la început numărul de obiecte din șir. Acest fapt face metoda inaplicabilă în anumite situații.

Transmiterea de la început a numărului de obiecte este utilizată, de exemplu, de protocolul *HTTP* (§ 11.3.2) la transmiterea, de către server, a conținutului paginii cerute de client. Serverul transmite întâi numărul de octeți ai paginii și apoi șirul de octeți ce formează pagina.

- *După ultimul obiect din șirul propriu-zis, se transmite o valoare specială, cu rol de terminator.* Această metodă ușurează munca emițătorului, care poate să înceapă transmiterea șirului înainte de-a ști câte elemente are, în schimb îngreunează munca receptorului, care trebuie să citească elementele șirului unu câte unu și să verifice dacă nu a întâlnit terminatorul. De asemenea, trebuie fixată valoarea terminatorului, care trebuie să fie o valoare reprezentabilă în formatul pentru element, dar care nu apare niciodată ca valoare a unui element valid.

Metoda este utilizată frecvent în transmiterea unui șir de caractere. Rolul de terminator poate fi acordat caracterului *null* (caracterul cu codul ASCII zero), caracterului *newline* (sfârșit de rând), caracterului

spațiu, etc. Orice alegere s-ar face, caracterul sau caracterele astfel alese pentru a marca sfârșitul unui șir nu pot să apară în șirul propriu-zis. Ca urmare, transmiterea unui fișier cu conținut arbitrar (șir de octeți cu valori arbitrare) nu se poate face prin metoda cu terminator (decât dacă un octet al fișierului se codifică pe mai mult de 8 biți).

7.1.4. Alte metode de reprezentare a numerelor întregi

Schemele de reprezentare (formatele) pentru numere întregi, studiate în § 7.1.2.3 și § 7.1.2.4, se numesc formate binare. Pe lângă formatele binare, pentru reprezentarea numerelor întregi se mai utilizează următoarele tipuri de formate:

- Formatul *text*. În cadrul acestui format, reprezentarea numărului este un text format din caracterele corespunzătoare cifrelor reprezentării zecimale (obișnuite) a numărului.
- Formatul *binar-zecimal*, numit și *BCD* din engl. *binary coded decimal*. În cadrul acestui format, numărul este reprezentat mai întâi în baza 10, iar apoi fiecare cifră zecimală este reprezentată pe 4 biți conform metodelor din § 7.1.2.4.

Descriem puțin mai pe larg reprezentarea numerelor în format text, deoarece o astfel de reprezentare se utilizează frecvent în protocoalele în rețea. Motivul principal al utilizării formatului text este ușurința depanării aplicațiilor ce utilizează astfel de protocoale: comunicația poate fi înregistrată într-un fișier și examinată cu un program obișnuit pentru vizualizarea fișierelor text.

În format text, se utilizează convențiile de reprezentare a numerelor în textele scrise: se începe cu cifra cea mai semnificativă, numărul de cifre este variabil (depinde de valoarea numărului) și prima cifră (cea mai semnificativă) scrisă nu este zero, cu excepția cazului numărului 0 care se reprezintă ca o singură cifră zero.

Fiecare cifră se reprezintă ca un caracter, fiind necesară mai departe o schemă de reprezentare a textelor (vezi § 7.2). În cazul codificării ASCII, reprezentarea fiecărei cifre ocupă exact un octet. De remarcat însă că, în acest caz, cifra 0 nu se reprezintă ca un octet cu valoarea 0, ci ca un octet având ca valoare codul ASCII pentru caracterul „0”; acesta este 48 (sau, echivalent, 30₁₆).

Deoarece lungimea reprezentării este variabilă, este necesar să fie transmisă sub o formă sau alta informația privind lungimea reprezentării numărului (numărul de cifre). În acest scop, în reprezentările text, numerele sunt separate de obicei prin spații, caractere *tab*, caractere *newline* etc.

EXEMPLUL 7.4: Redăm mai jos reprezentările text ASCII terminat cu spațiu, BCD *big endian* pe 4 octeți și BCD *little endian* pe 4 octeți, pentru numărul 300. Valorile octeților sunt scrise în baza 2, bitul cel mai semnificativ fiind scris în stânga.

poziție octet (nr. ordine)	Valorile octeților		
	text ASCII	BCD <i>big endian</i>	BCD <i>little endian</i>
0	00110011	00000000	00000000
1	00110000	00000000	00000011
2	00110000	00000011	00000000
3	00100000	00000000	00000000

7.2. Codificarea textelor

Prin text înțelegem aici un text scris în limbaj natural sau într-un limbaj de programare, fără formatare avansată.

Un text este văzut în general ca o succesiune de *caractere*. Caracterele sunt în principal literele din alfabetul limbii în care este scris textul, semne de punctuație, cifre și diferite alte semne grafice. Nu se face distincție între diferitele variante de-a scrie o aceeași literă (litere „normale“, cursive („italice“), adline („bold“), etc).

Pe lângă caracterele grafice, descrise mai sus, sunt definite *caractere de control*, având rolul de a marca puncte (locuri) în cadrul textului sau fragmente din text. Utilizări ale caracterelor de control sunt, de exemplu, trecerea la rând nou sau interzicerea trecerii la rând nou (ruperea rândului) într-un anumit punct.

Un aspect discutabil legat de alegerea setului de caractere este dacă o literă cu un semn diacritic este cazul să fie caracter distinct față de litera simplă din care provine sau să fie format din caracterul corespunzător literei respective fără semne diacritice și un caracter de control care să marcheze semnul diacritic adăugat. În aceeași idee, s-ar putea face și distincția dintre literele mari (majuscule) și literele mici (minuscule) corespunzătoare tot pe baza unor caractere de control cu rol de modificador.

Operațiunile efectuate asupra textelor, care trebuie să fie permise de codificarea aleasă, sunt în principal următoarele:

- afișarea textului;
- concatenarea unor texte sau alte operații de sinteză;

- căutarea unui cuvânt, extragerea unor cuvinte sau unor fragmente de text și diverse alte operații de analiză a textului;
- sortarea alfabetică.

De notat că regulile de sortare alfabetică sunt complexe și depind de limbă. De exemplu, în limba română, literele cu diacritice sunt considerate imediat după literele fără diacritice. Următoarele cuvinte sunt sortate alfabetic: *sac*, *suc*, *șiret*; de notat că orice cuvânt ce începe cu *ș* este sortat după toate cuvintele ce încep cu *s*. În franceză, însă, literele cu diacritice sunt considerate, în prima fază, echivalente cu cele fără diacritice, intervenind în ordinea alfabetică doar pentru cuvinte care diferă doar prin diacritice. Exemplu: *été*, *être*, *étude*; de notat că apar amestecate cuvinte ce încep cu *é* și *ê*.

Majoritatea codificărilor sunt bazate pe reprezentarea una după alta a literelor (caracterelor) ce formează cuvintele textului.

Codificările caracterelor sunt de obicei descrise în două etape. În prima etapă, fiecărui caracter îi este asociat un număr întreg pozitiv, numit *codul caracterului*. În a doua etapă, fiecărui cod de caracter îi este asociată o codificare ca șir de biți sau ca șir de octeți.

Pentru o schemă de codificare trebuie așadar specificate trei elemente:

- setul de caractere;
- numerotarea (codificarea) caracterelor;
- reprezentarea pe biți sau pe octeți a codurilor caracterelor.

7.2.1. Codificarea ASCII

Codificarea (codul) ASCII este codificarea cea mai des întâlnită pentru texte.

Setul de caractere cuprinde 128 caractere dintre care:

- 33 de caractere de control;
- caracterul *spațiu* (considerat de unii ca fiind caracter imprimabil și de alții ca fiind caracter de control);
- 94 de caractere imprimabile, cuprinzând: 52 de litere (cele 26 litere ale alfabetului latin, cu cele două forme, majuscule și minuscule) cele 10 cifre zecimale și un număr de 32 de semne de punctuație și alte simboluri.

Codurile asociate caracterelor ASCII sunt cuprinse între 0 și 127, caracterele de control primind codurile 0–31 și 127, spațiul are codul 32, iar

(celelalte) caractere imprimabile au codurile cuprinse între 33 și 126. Pentru a ușura sortarea alfabetică, codurile sunt grupate astfel:

- literele mari de la 65 (41 hexa) pentru A la 90 (5A hexa) pentru Z;
- literele mici de la 97 (61 hexa) pentru a la 122 (7A hexa) pentru z;
- cifrele de la 48 (30 hexa) pentru 0 la 57 (39 hexa) pentru 9.

De remarcat și că diferența dintre codul oricărei litere mici și codul literei mari corespunzătoare este 32 (20 hexa).

Pentru reprezentarea unui caracter ASCII sunt necesari doar 7 biți, însă cel mai adesea un caracter ASCII se reprezintă pe un octet, al cărui cel mai semnificativ bit este întotdeauna 0.

Datorită faptului că pe de o parte caracterele ASCII se reprezintă pe un octet, iar pe de altă parte că dintre caracterele de control multe nu sunt utilizate deloc în majoritatea aplicațiilor, rămân multe coduri reprezentabile (cca. 140) care nu sunt utilizate. Se poate extinde setul de caractere, asociind noilor caractere coduri între 128 și 255 sau coduri între 0 și 31 a căror caractere corespunzătoare nu sunt folosite efectiv. Toate aceste codificări rezultate se numesc generic *seturi ASCII extinse*.

7.2.2. Codificările ISO-8859

ISO-8859 este o familie de coduri, construite toate ca extensii (alternative) ale codificării ASCII.

Fiecare cod din familie cuprinde câte 256 caractere: cele 128 caractere ASCII, plus 128 de caractere alese pentru a acoperi alfabetul câte unui grup de limbi. Limbile acoperite de câteva dintre codificările ISO-8859 sunt:

- ISO-8859-1, alfabetul latin pentru limbile din vestul Europei;
- ISO-8859-2, alfabetul latin pentru limbile din estul Europei;
- ISO-8859-5, alfabetul chirilic;
- ISO-8859-6, alfabetul arab;
- ISO-8859-7, alfabetul grecesc;
- ISO-8859-8, alfabetul ebraic.

Codurile asociate caracterelor sunt codurile din codificarea ASCII pentru cele 128 de caractere din setul ASCII și numere de la 128 la 255 pentru caracterele suplimentare.

Reprezentarea pe octeți pentru un text ISO-8859-*n* se face cu câte un octet pentru fiecare caracter, octetul conținând codul caracterului.

Fiecare cod din familie este extensie a codului ASCII în sensul că mulțimea caracterelor din fiecare astfel de cod include mulțimea caracterelor

Caracter	Cod (hexa)	Caracter	Cod (hexa)
Ă	C3	ă	E3
Â	C2	â	E2
Î	CE	î	EE
Ș	AA	ș	BA
Ț	DE	ț	FE

Tabelul 7.1: Caracterele cu diacritice din alfabetul limbii române și codificările ISO-8859-2 corespunzătoare

ASCII și codurile asociate caracterelor comune cu setul ASCII coincid cu codurile ASCII. Ca urmare, un text ASCII este întotdeauna interpretat corect ca text ISO-8859- n . Pe de altă parte, un text scris în ISO-8859- n și interpretat ca ISO-8859- m va fi evident interpretat greșit.

Ordinea numerică a codurilor din oricare dintre codificările ISO-8859 este diferită de ordinea alfabetică. În general, în ordinea alfabetică, literele cu diacritice își au locul imediat lângă literele similare fără diacritice; în codificările ISO-8859-1 sau ISO-8859-2, de exemplu, literele cu diacritice au codurile mai mari de 128 în vreme ce literele fără diacritice au coduri între 65 și 123.

7.2.3. Codificările Unicode

Unicode este un set de caractere ce se dorește să cuprindă litere din toate scrierile de pe Pământ. Numărul de caractere din unicode este limitat, datorită modurilor de codificare definite, la aproximativ un milion (mai exact, la $110000_{16} = 1114112$). Nu toate aceste coduri sunt definite în prezent, codurile încă nedefinite putând fi definite în versiuni următoare ale standardului.

Codurile unicode sunt numere de la 0 la $2^{20} + 2^{16} - 1$. Codurile de la 0 la 127 corespund aceluiași caractere ca și în codificarea ASCII.

Reprezentarea codurilor unicode ca șiruri de octeți poate fi făcută în mai multe moduri. Cele mai răspândite codificări sunt:

- *UTF-8* este o codificare de lungime variabilă, între 1 și 4 octeți pentru un caracter;
- *UTF-16*, *UTF-16LE*, *UTF-16BE* sunt codificări de lungime variabilă, 2 sau 4 octeți pentru un caracter;
- *UTF-32*, *UTF-32LE*, *UTF-32BE* sunt codificări de lungime fixă, de 4 octeți pentru fiecare caracter.

Carac- ter	Cod <i>unicode</i> (hexa)	Cod <i>unicode</i> (zecimal)	UTF-8 (hexa)
Ă	102	258	C4 82
ă	103	259	C4 83
Â	C2	194	C3 82
â	E2	226	C3 A2
Î	CE	206	C3 8E
î	EE	238	C3 AE
Ș	218	536	C8 98
ș	219	537	C8 99
Ț	21A	538	C8 9A
ț	21B	539	C8 9B
Ş	15E	350	C5 9E
ş	15F	351	C5 9F
Ț	162	354	C5 A2
ț	163	355	C5 A3

Tabelul 7.2: Caracterele cu diacritice din alfabetul limbii române și codificările *unicode* corespunzătoare. Notă: caracterele Ș, ș, Ț și ț au câte două forme utilizate: una cu virgulă dedesupt, cealaltă cu sedilă. Conform normelor stabilite de Academia Română, forma corectă este cea cu virgulă. Codificarea formei cu virgulă a fost standardizată mai recent, motiv pentru care multe documente utilizează încă forma cu sedilă.

7.2.3.1. Codificarea UTF-8

Correspondența de la codul caracterului la șirul de octeți este dată în tabelul 7.3.

Valorile lui c (în baza 16)	reprezentarea UTF-8 (în baza 2)
0–7F	$0c_7c_6c_5c_4c_3c_2c_1c_0$
80–7FF	$110c_{10}c_9c_8c_7c_6$ $10c_5c_4c_3c_2c_1c_0$
800–FFFF	$1110c_{15}c_{14}c_{13}c_{12}$ $10c_{11}c_{10}c_9c_8c_7c_6$ $10c_5c_4c_3c_2c_1c_0$
10000–1FFFFF	$11110c_{20}c_{19}c_{18}$ $10c_{17}c_{16}c_{15}c_{14}c_{13}c_{12}$ $10c_{11}c_{10}c_9c_8c_7c_6$ $10c_5c_4c_3c_2c_1c_0$

Tabelul 7.3: Codificarea UTF-8. c reprezintă codul *unicode* al caracterului; $c_{20} \dots c_0$ reprezintă cifrele reprezentării binare a lui c , cu c_{20} reprezentând cifra cea mai semnificativă și c_0 cea mai puțin semnificativă. Codificarea există doar pentru $0 \leq c < 2^{21}$.

De remarcat că schema pentru coduri mari (de exemplu, schema pentru c între 80_{16} și $7FF_{16}$) poate fi principial aplicată și la coduri mai mici (de exemplu, pentru $c = 41_{16}$, rezultând doi octeți, $C1_{16}$ urmat de 81_{16}). O astfel de codificare este însă interzisă de standard pentru a asigura unicitatea codificării UTF-8.

Codificarea UTF-8 permite recuperarea sincronismului (dacă receptorul pierde câțiva octeți poate regăsi unde începe un caracter nou), deoarece fiecare caracter nou începe cu un octet cuprins între 0 și 127 sau între 192 și 255, iar ceilalți octeți din codificarea unui caracter sunt cuprinși între 128 și 191. O altă proprietate este că lungimea codificării UTF-8 a unui caracter poate fi determinată după citirea primului octet.

7.2.3.2. Codificările UTF-16

Codificarea UTF-16 este descrisă în două etape: într-o primă etapă, codul *unicode* este transformat într-unul sau două numere de câte 16 biți, iar în a doua etapă fiecare astfel de număr este scris ca 2 octeți consecutivi.

Caracterele cu codul *unicode* între 0 și $D7FF_{16}$ sau între $E000_{16}$ și $FFFF_{16}$ se scriu ca un singur întreg pe 16 biți.

Caracterele cu codul *unicode* între 10000_{16} și $10FFFF_{16}$ se scriu ca doi întregi de câte 16 biți astfel: Mai întâi, din codul *unicode* se scade 10000_{16} , rezultând o valoare între 0 și $FFFF_{16}$ (20 biți). Primul întreg de 16 biți se formează punând cifrele 110110 urmate de primii 10 din cei 20 de biți. Al doilea întreg se formează punând cifrele 110111 urmate de ultimii 10 din cei 20 de biți. De exemplu, codul *unicode* 10302_{16} se scrie ca doi întregi astfel: $D83C_{16}$ $DF02_{16}$.

Într-o a doua etapă este definită scrierea fiecărui întreg de 16 biți ca un șir de doi octeți. Există două modalități de a reprezenta fiecare astfel de întreg, începând de la octetul mai semnificativ (de rang mai mare) sau începând de la octetul mai puțin semnificativ. Pentru a reflecta aceste variante diferite de alegere, există trei codificări distincte numite generic *UTF-16*:

- *UTF-16LE*: Primul octet este cel mai puțin semnificativ (*little endian*);
- *UTF-16BE*: Primul octet este cel mai semnificativ (*big endian*);
- *UTF-16*: Ordinea octeților poate fi fie *big endian*, fie *little endian*, la alegerea emițătorului. Primul caracter codificat trebuie să fie caracterul cu codul $FEFF_{16}$ (definit inițial ca fiind un caracter de control ce interzice ruperea în rânduri în acel punct, dar este utilizat în prezent doar ca marcaj pentru identificarea ordinii octeților). Ordinea octeților este dedusă de receptor prin examinarea primilor doi octeți: dacă aceștia sunt FE_{16} urmat de FF_{16} , înseamnă că ordinea octeților este *big endian*; dacă este FF_{16} urmat de FE_{16} , înseamnă că ordinea octeților este *little endian*.

7.2.3.3. Codificările UTF-32

Codificarea UTF-32 constă în codificarea fiecărui caracter ca un întreg pe 32 de biți, reprezentat la rândul lui ca un șir de 4 octeți. Ca și în cazul codificărilor *UTF-16*, există trei codificări *UTF-32*:

- *UTF-32LE*: Primul octet este cel mai puțin semnificativ (*little endian*);
- *UTF-32BE*: Primul octet este cel mai semnificativ (*big endian*);
- *UTF-32*: Ordinea octeților poate fi fie *big endian*, fie *little endian*, la alegerea emițătorului. Primul caracter codificat trebuie să fie caracterul cu codul $FEFF_{16}$. Ordinea octeților este dedusă de receptor prin examinarea primilor patru octeți: dacă aceștia sunt 0, 0, FE_{16} , FF_{16} , înseamnă că ordinea octeților este *big endian*; dacă este FF_{16} , FE_{16} , 0, 0, înseamnă că ordinea octeților este *little endian*.

7.3. Reprezentarea datei și orei

Determinarea datei și orei producerii unui eveniment, precum și memorarea sau transmiterea acestora, sunt operații frecvente într-o rețea de calculatoare.

Problema reprezentării datei și orei este mult mai dificilă decât pare la prima vedere. Din acest motiv, vom începe prin a studia ce se poate înțelege

prin „ora curentă“, iar apoi vom vedea ce scheme de reprezentare ale datei și orei există și ce avantaje și dezavantaje aduce fiecare dintre ele.

7.3.1. Măsurarea timpului

Există două metode utilizate pentru indicarea timpului curent (datei și orei curente):

- pe baza unor fenomene astronomice, anume alternanța zi-noapte (în termeni astronomici, *ziua solară mijlocie*), alternanța anotimpurilor (*anul tropic*) și, eventual, fazele lunii (*luna sinodică*);
- pe baza unui fenomen fizic repetabil, de exemplu oscilația unui pendul sau vibrația unui cristal de cuarț.

Prima variantă este de interes practic imediat pentru sincronizarea activităților umane. Are însă complicații inerente legate de următoarele fapte:

- alternanța zi-noapte nu este simultană pe tot Pământul ci este decalată pe longitudine;
- anul, luna și ziua sunt incommensurabile (rapoartele duratelor lor sunt numere iraționale);
- anul, luna și ziua nu au durate constante (fenomenele corespunzătoare nu sunt perfect periodice) și nici măcar previzibile exact (în special rotația Pământului are neuniformități imprevizibile datorate redistribuirii masei în interiorul Pământului).

A doua variantă măsoară direct timpul ca mărime fizică și oferă avantaje atunci când avem de determinat ordinea cronologică a unor evenimente sau de calculat duratele de timp dintre ele. Timpul (fizic) a ajuns să poată fi definit independentă de mișcarea Pământului doar după dezvoltarea, începând cu anii 1950, a ceasurilor atomice, mai precise decât mișcările Pământului. Măsurarea timpului se face pe baza *secunde* definite în Sistemul Internațional de unități (SI) ca *9192631770 de perioade ale oscilației corespunzătoare tranziției între cele două nivele hiperfine ale stării fundamentale a atomului de cesiu 133*.

Ca urmare a acestor complicații, există mai multe standarde de măsurare și reprezentare a timpului:

Timpul atomic internațional (TAI) este dat de numărul de secunde SI scurse de la un anumit moment ales ca reper. Secundele TAI se grupează în minute, ore, zile, etc.

Timpul universal UT1 este de fapt măsura unui anumit unghi, legat de rotația Pământului, exprimată în unități de timp (24 h în loc de 360°). (Unghiul respectiv este unghiul orar, pentru un observator aflat

pe meridianul 0° , al soarelui mijlociu.) Curge neuniform datorită neuniformității mișcării de rotație a Pământului; după media ultimilor câțiva ani, 24 h UT1 este aproximativ 86400,002 s SI.

Timpul universal coordonat (UTC) este bazat pe secunda SI, dar gruparea secundelor în zile este modificată pentru a menține diferența dintre UT1 și UTC la sub o secundă.

Astfel, o zi UTC normală are 24 ore a 60 minute a 60 secunde SI fiecare, adică 86400 s. Dacă UT1–UTC se apropie de -1 s, se adaugă o secundă de corecție (engl. *leap second*) la o zi, astfel încât acea zi UTC are 86401 s, proces aproape echivalent cu a muta UTC cu o secundă înapoi. În acest scop, ultimul minut al zilei are 61 de secunde în loc de 60, după ora 23:59:59 urmează, la o secundă, 23:59:60 și abia după încă o secundă ora 0:00:00 a zilei următoare. Dacă UT1–UTC se apropie de 1 s, se elimină o secundă din ultimul minut al unei zile, astfel că la o secundă după 23:59:58 urmează ora 0:00:00 a zilei următoare. Din anul 1972 (de la introducerea UTC în forma actuală) până în 2008 au fost adăugate 23 de secunde de corecție și nu a fost eliminată nici una. A 24-a secundă de corecție se va adăuga la sfârșitul anului 2008, astfel încât ziua de 31 decembrie 2008 va avea 86401 secunde. Datorită unei diferențe inițiale de 10 s între TAI și UTC, diferența TAI–UTC este în prezent de 33 s.

Timpul legal în fiecare țară este definit fie pe baza UT1, fie pe baza UTC (diferența este neglijabilă pentru uzul practic), ca fiind UTC (sau UT1) plus sau minus un anumit număr de ore și uneori și fracțiuni de oră (exemplu, India are ora legală UTC+5h30min).

În țările în care există oră de vară, la trecerea de la ora de iarnă la cea de vară și invers, diferența dintre ora legală și UTC crește, respectiv scade, cu o oră (de notat că UTC nu are oră de vară). De exemplu, ora legală în România este UTC+2 h în timpul iernii (din ultima duminică din octombrie până în ultima duminică din martie) și UTC+3 h în timpul verii.

Ora suplimentară introdusă la trecerea de la ora de vară la cea de iarnă nu are notație distinctă, de tipul secundelor de corecție din UTC. Ca urmare, la trecerea de la ora de vară la ora de iarnă, există perechi de momente de timp care sunt notate la fel și ca urmare ora legală este ambiguă. Exemplu: dacă prin trecerea de la ora de vară la cea de iarnă ora 4:00:00 devine 3:00:00, atunci notația 3:30:00 poate corespunde la două momente de timp, ora de vară 3:30:00 (la 30 min înaintea schimbării orei) și ora de iarnă 3:30:00 (la 30 min după schimbarea orei).

Pentru gruparea zilelor în unități mai mari, în special în ani, sunt utilizate mai multe sisteme (calendare):

Calendarul gregorian, introdus în anul 1582 și în vigoare în România din anul 1924, este calendarul actual. Anii bisecți (de 366 de zile) sunt anii cu numărul anului divizibil cu 4, cu excepția celor divizibili cu 100 fără a fi divizibili cu 400. Ani bisecți sunt 1600, 2000, 2400 etc; ani nebisecți divizibili cu 100 sunt 1700, 1800, 1900, 2100, 2200 etc. Durata medie a anului gregorian este 365,2425 zile, ceva mai lung decât anul tropic de aproximativ 265,2422 zile.

Calendarul iulian, predecesorul calendarului gregorian, introdus în anul 45 î.e.n. și având regula mai simplă cum că sunt bisecți toți anii cu număr divizibil cu 4. Este utilizat adesea de istorici pentru a data și evenimente dinainte de anul 45 î.e.n., caz în care el este numit *calendar iulian proleptic*. Cu o durată medie a anului de 365,25 zile, calendarul iulian rămâne în urmă cu 1 zi la aproximativ 128 de ani.

Ziua iuliană este un simplu număr ce arată numărul de zile scurse de la o dată de referință. Acest sistem este utilizat frecvent în astronomie deoarece permite ușor calculul duratelor dintre două date; din același motiv reprezintă o schemă potrivită pentru reprezentarea datei în calculator. Există în două variante. Prima, JD (*julian day*), are ca referința data de 1 ianuarie 4713 î.e.n. conform calendarului iulian proleptic, la amiaza UT1. Momentul respectiv este JD 0,0, miezul nopții următoare este JD 0,5, etc. Cealaltă, MJD (*modified julian day*), are ca referință 17 noiembrie 1858 ora 0, adică este $JD - 2400000,5$.

7.3.2. Obiectivele în alegerea reprezentării timpului în calculator

De obicei, operațiile ce trebuiesc efectuate asupra reprezentării timpului sunt:

1. Citirea sau scrierea timpului ca oră legală conform calendarului gregorian în formatul obișnuit, precum și efectuarea de operații aritmetice de genul adunării sau scăderii unei durate formale (exemplu: mâine la aceeași oră; aceasta înseamnă în mod obișnuit peste 24 de ore, dar poate însemna peste 23 sau 25 de ore dacă intervine trecerea de la ora de iarnă la cea de vară sau invers).
2. determinarea ordinii cronologice a două momente de timp;
3. determinarea exactă, ca timp fizic, a duratei între două momente de timp,

4. pentru aplicații speciale, citirea sau afișarea timpului în alte formate (timpul legal al altui fus orar, UTC, TAI, JD, etc).

Punctul 1 este cerut de toate sistemele. Punctul 2 este important pentru foarte multe aplicații și rezolvarea lui corectă interzice mutarea ceasului înapoi. Punctul 3 este important în aplicațiile în timp real; de asemenea, funcționarea ceasului sistem presupune, în mod repetat, adunarea unei durate de timp la un moment de timp.

Reprezentarea directă a orei legale, sub forma an, lună, zi, oră, minut, secundă, fracțiuni de secundă, rezolvă simplu punctul 1. Ea ridică însă probleme la punctul 2 dacă sunt implicate calculatoare aflate pe fusuri orare distincte sau dacă se efectuează operații în intervalul de o oră în jurul trecerii de la ora de vară la cea de iarnă; pentru tratarea corectă a acestor cazuri este necesar să se știe, despre fiecare oră, pe ce fus orar este considerată și care sunt regulile privind ora de vară. Punctul 3 ridică, pe lângă problemele comune cu cele legate de punctul 2, complicații legate de saltul cu o oră înainte la trecerea de la ora de iarnă la ora de vară și calculele legate de calendar; de asemenea, pentru calcule exacte ale duratelor, sunt necesare informații cu privire la secunde de corecție.

Reprezentarea orei UTC permite determinarea ordinii cronologice și a duratelor fără a necesita date despre fusurile orare sau regulile privind ora de vară, în schimb aceste date sunt necesare la conversia între reprezentarea UTC și timpul legal.

Reprezentarea TAI ca număr de unități de timp scurse de la un anumit moment fixat rezolvă extrem de simplu punctele 2 și 3 în schimb mută dificultățile la rezolvarea punctului 1.

7.3.3. Formate utilizate în practică

Deoarece într-o rețea pot fi prezente calculatoare situate pe fusuri orare distincte, aproape orice format util în rețea fie transmite direct ora UTC sau TAI, fie transmite suficientă informație pentru ca receptorul să poată calcula ușor ora UTC.

7.3.3.1. Formatul utilizat de poșta electronică

Pentru poșta electronică (§ 11.1), reprezentarea datei se face ca text și conține, în ordine:

- opțional ziua din săptămână, ca prescurtare de 3 litere din limba engleză),
- ziua, ca număr între 1 și 31,
- luna, ca șir de trei litere, prescurtare din engleză,

- anul, ca șir de 4 cifre,
- ora, totdeauna ca 2 cifre, între 00 și 23,
- minutul, ca două cifre, între 00 și 59,
- opțional, secunda, ca două cifre între 00 și 60,
- diferența dintre ora legală conform căreia a fost scrisă data și ora UTC; aceasta este dată ca 4 cifre, 2 pentru numărul de ore și 2 pentru numărul de minute, cele patru cifre fiind precedate de semnul + sau -. Componentele datei sunt separate printr-un amestec de virgule, spații și caractere *două puncte*.

De exemplu, data:

Thu, 25 Oct 2007, 17:22:19 +0300

înseamnă că la momentul scrierii mesajului ora locală a expeditorului era joi, 25 octombrie 2007, ora 17:22:19 și că ora respectivă este cu 3 ore în avans față de UTC. Ca urmare, ora UTC la acel moment era 14:22:19.

Data considerată în acest exemplu este plauzibilă conform orei legale a României, în 25 octombrie 2007 fiind încă în vigoare ora de vară care este cu 3 ore în avans față de UTC.

Orele astfel specificate sunt ușor de comparat și nu există ambiguități legate de trecerea de la ora de vară la cea de iarnă. De exemplu, un mesaj trimis înainte de trecerea la ora de iarnă ar fi datat

Sun, 28 Oct 2007, 03:40 +0300

urmat la jumătate de oră, după trecerea la ora de iarnă, de

Sun, 28 Oct 2007, 03:10 +0200

7.3.3.2. ISO-8601 și RFC-3339

ISO-8601 este o standardizare a modului de scriere ca text a datei și orei. Standardul fiind foarte complex, a apărut RFC-3339 care cuprinde cazurile mai utile și mai frecvent folosite din ISO-8601.

RFC-3339 prevede reprezentarea datelor astfel:

- anul, ca patru cifre (nu sunt admise prescurtări de genul 07 pentru 2007),
- luna, ca 2 cifre (01 pentru ianuarie, 12 pentru decembrie),
- ziua, ca 2 cifre (de la 01 până la 31 sau mai puțin, în funcție de lună).

Cele trei componente sunt separate prin liniuțe (ISO-8601 permite și alipirea lor):

2007-10-28

Ora se reprezintă prin 2 cifre pentru oră (00–23), 2 cifre pentru minut (00–59), două cifre pentru secundă (00–60, în funcție și de prezența unei secunde de corecție), eventual fracțiunile de secundă și eventual specificarea fusului orar. Ora, minutul și secunda sunt separate prin *două puncte*, fracțiunile de secundă se separă de câmpul pentru secunde prin *punct*, iar specificarea fusului orar se face printr-un semn plus sau minus urmat de două cifre pentru numărul de ore diferență urmat de caracterul *două puncte* și încă două cifre pentru numărul de minute diferență. Exemplu:

21:12:58.342+02:00

reprezintă același moment de timp, dar pe alt fus orar, cu

14:12:58.342-05:00

Data și ora se specifică împreună punând litera T între ele:

2007-10-28T14:12:58.342-05:00

7.3.3.3. Timpul POSIX

În sistemele de tip UNIX (conforme standardului POSIX), reprezentarea timpului este făcută printr-un număr întreg considerat de obicei că reprezintă numărul de secunde scurse de la 1 ianuarie 1970 ora 0:00 UTC. Ora UTC în formatul obișnuit se obține grupând secunde în minute, ore, zile, luni și ani conform regulilor obișnuite. De fapt, numărul dat ca dată nu este exact numărul de secunde scurse de la 1 ianuarie 1970, ci diferă de acesta prin numărul de secunde de corecție adăugate pentru menținerea în sincronism a UTC cu rotația Pământului. De aceea, „timpul unix“ are salturi înapoi de câte o secundă la fiecare introducere a unei astfel de secunde de corecție. Timpul POSIX este reprezentat în mod obișnuit ca întreg cu semn pe 32 de biți, și ca urmare valoarea cea mai mare ce poate fi reprezentată corespunde datei de 19 ianuarie 2038, ora 3:14:07 UTC.

7.3.3.4. TAI 64

TAI 64 este un standard ce presupune reprezentarea timpului ca număr de secunde, incluzând secunde de corecție. Numărul de secunde este reprezentat pe 63 de biți (plus un bit rezervat), cu valoarea 2^{62} corespunzând datei de 1 ianuarie 1970 ora 0 TAI. Intervalul de timp reprezentabil este imens, de ordinul a 10^{11} ani.

Pentru aplicații ce au nevoie de rezoluție mai bună de o secundă, standardul prevede încă două câmpuri de câte 32 de biți (total 128 de biți), reprezentând respectiv numărul de nanosecunde și de attosecunde (valori între 0 și $10^9 - 1$).

7.4. Recodificări

Este necesar uneori să codificăm un șir mai mult sau mai puțin arbitrar de octeți sub forma unui șir de caractere supus unor restricții. Astfel de situații apar:

- La trimiterea fișierelor atașate la mesajele de poștă electronică, întregul mesaj, inclusiv partea ce cuprinde fișierele atașate, trebuie să îndeplinească anumite restricții, între altele, să nu conțină caractere ASCII de control cu excepția perechilor *carriage return–line feed* de la finalul fiecărui rând, să nu aibă rânduri prea lungi, etc. Pe de altă parte, fișierele atașate pot fi fișiere binare cu conținut arbitrar.
- La stocarea în fișiere text a unor informații reprezentate natural ca șir arbitrar de octeți, de exemplu la stocarea în fișiere text a unor chei de criptare, semnături electronice, etc.
- În limbaje de programare, la scrierea în șirurile de caractere a unor caractere cu rol special, ca de exemplu a ghilimelelor (care în mod normal sunt interpretate ca terminatorul șirului de caractere).

În astfel de situații, este necesar să se codifice un șir arbitrar de octeți (fiecare octet poate lua orice valori între 0 și 255) pe un alfabet constând în litere, cifre și câteva simboluri speciale. Deoarece numărul de simboluri disponibile este mai mic decât numărul de valori posibile ale octeților, un octet nu se va putea codifica numai pe un singur caracter.

7.4.1. Codificarea hexazecimală

Codificarea hexazecimală prevede ca fiecare octet să se reprezinte pe două caractere, fiecare caracter putând fi din mulțimea cuprinzând cifrele zecimale (0–9) și literele A–F.

Exemplu: șirul de octeți 120, 0, 23, 45, 20 se scrie: 7800172D14.

Uneori, în șirul de cifre hexa se inserează spații sau caractere *newline* pentru lizibilitate sau pentru evitarea rândurilor foarte lungi.

Deoarece un caracter se reprezintă de obicei pe un octet, prin această recodificare rezultă o dublare a lungimii șirului.

7.4.2. Codificarea în baza 64

Codificarea *în baza 64* codifică un șir de 3 octeți cu valori arbitrare ca un șir de 4 caractere din mulțimea cuprinzând literele mari, literele mici, cifrele și caracterele +, / și =.

Codificarea se face în modul următor:

1. Șirul inițial de octeți se completează la un multiplu de 3 octeți prin adăugarea a 0, 1 sau 2 octeți cu valoarea 0.
2. Se formează un șir de 24 de biți punând unul după altul cei câte 8 biți din fiecare octet; din fiecare octet se începe cu bitul cel mai semnificativ.
3. Șirul de 24 de biți se împarte în 4 grupuri de câte 6 biți.
4. Fiecare șir de 6 biți se consideră ca fiind un număr cuprins între 0 și 63, considerând primul bit din șir ca fiind cel mai semnificativ.
5. Fiecare număr obținut la pasul anterior se reprezintă printr-un caracter. Cele 64 de valori posibile, de la 0 la 63, se reprezintă ca: litere mari (0→A, 25→Z), litere mici (26→a, 51→z), cifre (52→0, 61→9) și caracterele + și / (62→+, 63→/). Dacă o valoare 0 provine din biți 0 proveniți integral dintr-un octet completat la pasul 1, în loc de A se scrie =.

De exemplu, șirul de octeți 120, 0, 23, 45, 20 devine șirul de biți

```
01111000 00000000 00010111 00101101 00010100 00000000
```

ultimul octet fiind adăugat la pasul 1. Șirul de biți se regroupează

```
011110 000000 000000 010111 001011 010001 010000 000000
```

rezultând șirul de numere „în baza 64”: 30, 0, 0, 23, 11, 17, 16, 0, care se codifică eAAXLRQ=

7.4.3. Codificări bazate pe secvențe de evitare

Recodificările prin secvențe de evitare sunt utilizate în situația în care majoritatea octeților sau caracterelor din textul de recodificat sunt codurile unor caractere ce se regăsesc în alfabetul în care se face recodificarea. În această situație, este favorabil ca, pe cât posibil, octeții sau caracterele din textul de recodificat să fie codificate prin ele însele, în special pentru ca textul să poată fi înțeles (parțial, cel puțin) de către un utilizator uman direct în forma recodificată.

Recodificarea se face astfel. Se distinge un caracter în alfabetul destinație, caracter ce este denumit *caracter de evitare* (enlg. *escape character*).

- Orice caracter din alfabetul sursă care se regăsește în alfabetul destinație și este diferit de caracterul de evitare se recodifică ca el însuși.

- Caracterul de evitare (dacă face parte din alfabetul sursă), precum și caracterele din alfabetul sursă ce nu se găsesc în alfabetul destinație, se codifică ca secvențe de caractere ce încep cu un caracter de evitare.

Exemple:

- Pentru poșta electronică, un caracter ce nu pot fi transmise direct este codificat ca o secvență de trei caractere, un caracter *egal* (=) și două cifre hexa reprezentând codul caracterului de transmis. Această recodificare se numește *quoted printables*. Caracterele ASCII imprimabile, cu excepția caracterului *egal*, se transmit direct (fără recodificare). Ca urmare, un text în care caracterele ce trebuie recodificate sunt rare poate fi înțeles de către un utilizator uman fără prea mari dificultăți. Ca exemplu, fraza precedentă se scrie (presupunând o codificare ISO-8859-2 pentru caractere și apoi o recodificare *quoted printables*):

Ca urmare, un text =EEn care caracterele ce trebuie
recodificate sunt rare poate fi =EEn=FEeles de c=E3tre
un utilizator uman f=E3r=E3 prea mari difficult=E3=FEi.

- În *URL*-uri, caracterele ce au rol sintactic (*spațiu*, */*, etc) se codifică prin caracterul *procent* (%) urmat de două cifre hexa reprezentând valoarea caracterului respectiv. De notat că aceste coduri sunt în cadrul codificării UTF-8; ca urmare, o pagină cu numele *șir* ar avea un *URL* de forma

`http://example.com/%C8%98ir`

- În limbajul C și în alte limbaje derivate din el, ghilimelele (") servesc la delimitarea șirurilor de caractere. Dacă se dorește introducerea unor astfel de caractere într-un șir, sau a caracterelor ASCII de control, se introduc secvențe *escape* cum ar fi: *"* pentru ghilimele ("), ** pentru *backslash* (\), *\n* pentru *newline* (caracterul ASCII cu codul 10).
- În HTML, caracterele *unicode* sunt scrise prin secvențe *&#cod;*. De exemplu, litera ț se scrie *ț*.