

*A programozási nyelvek elemei*  
*(Lexikális elemek)*

- A programok fordítási egységei lexikális elemek sorozatai.
- A lexikális elemek karaktersorozatok, határoló jelekkel elválasztva.
- A nyelvekben általában megkülönböztetjük
  - a **grafikus** karaktereket (a Latin-1 vagy az angol abc betűi, számjegyek és speciális karakterek - pl. " # & ' ( ) , \* - + / : ; < = \_ [ ] { } | . ) --),
  - a **formátum vezérlő** karaktereket (az ISO 6429 szabványban: character tabulation (HT), line tabulation (VT), carriage return (CR), line feed (LF), and form feed (FF)), és
  - az **egyéb** (implementáció függő) **vezérlő** karaktereket.

## **Lexikális elem:**

- azonosító,**
- numerikus literál,**
- karakter literál,**
- string literál,**
- határoló,**
- megjegyzés**

# A jelkészlet

A **karakterkészlet** különböző karakterek egy halmaza.

- Nem teszünk fel semmit a számítógép belső karakterábrázolásáról. Még rendezést sem tételezünk fel a karakterek között,
- Megadjuk a karakterek neveit és a karakter megjelenítését látható formában.
- Tartalmazhat olyan karaktereket, amelyek bizonyos megjelenítéseknél ugyanúgy néznek ki, mégis logikailag különbözőnek tekintjük.

Például: a latin nagy „A”, a cirill nagy „A” és a görög nagybetűs Alfa „A”

Például:

EXCLAMATION !

QUESTION\_MARK ?

SEMICOLON ;

A **karakterkód** egy leképezés, amely kölcsönösen egyértelmű megfeleltetést ad a karakterkészlet karakterei és a nemnegatív egészek egy halmaza között.

Egy egyedi számkódot, egy kódpozíciót rendel a karakterkészlet minden karakteréhez. (Gyakran táblázat)

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32;	<b>Space</b>	64	40	100	&#64;	<b>@</b>	96	60	140	&#96;	<b>`</b>
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33;	<b>!</b>	65	41	101	&#65;	<b>A</b>	97	61	141	&#97;	<b>a</b>
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34;	<b>"</b>	66	42	102	&#66;	<b>B</b>	98	62	142	&#98;	<b>b</b>
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35;	<b>#</b>	67	43	103	&#67;	<b>C</b>	99	63	143	&#99;	<b>c</b>
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	<b>\$</b>	68	44	104	&#68;	<b>D</b>	100	64	144	&#100;	<b>d</b>
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	<b>%</b>	69	45	105	&#69;	<b>E</b>	101	65	145	&#101;	<b>e</b>
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	<b>&amp;</b>	70	46	106	&#70;	<b>F</b>	102	66	146	&#102;	<b>f</b>
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39;	<b>'</b>	71	47	107	&#71;	<b>G</b>	103	67	147	&#103;	<b>g</b>
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40;	<b>(</b>	72	48	110	&#72;	<b>H</b>	104	68	150	&#104;	<b>h</b>
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	<b>)</b>	73	49	111	&#73;	<b>I</b>	105	69	151	&#105;	<b>i</b>
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	<b>*</b>	74	4A	112	&#74;	<b>J</b>	106	6A	152	&#106;	<b>j</b>
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	<b>+</b>	75	4B	113	&#75;	<b>K</b>	107	6B	153	&#107;	<b>k</b>
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	<b>,</b>	76	4C	114	&#76;	<b>L</b>	108	6C	154	&#108;	<b>l</b>
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45;	<b>-</b>	77	4D	115	&#77;	<b>M</b>	109	6D	155	&#109;	<b>m</b>
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46;	<b>.</b>	78	4E	116	&#78;	<b>N</b>	110	6E	156	&#110;	<b>n</b>
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47;	<b>/</b>	79	4F	117	&#79;	<b>O</b>	111	6F	157	&#111;	<b>o</b>
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48;	<b>0</b>	80	50	120	&#80;	<b>P</b>	112	70	160	&#112;	<b>p</b>
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49;	<b>1</b>	81	51	121	&#81;	<b>Q</b>	113	71	161	&#113;	<b>q</b>
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50;	<b>2</b>	82	52	122	&#82;	<b>R</b>	114	72	162	&#114;	<b>r</b>
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51;	<b>3</b>	83	53	123	&#83;	<b>S</b>	115	73	163	&#115;	<b>s</b>
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52;	<b>4</b>	84	54	124	&#84;	<b>T</b>	116	74	164	&#116;	<b>t</b>
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	<b>5</b>	85	55	125	&#85;	<b>U</b>	117	75	165	&#117;	<b>u</b>
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	<b>6</b>	86	56	126	&#86;	<b>V</b>	118	76	166	&#118;	<b>v</b>
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	<b>7</b>	87	57	127	&#87;	<b>W</b>	119	77	167	&#119;	<b>w</b>
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56;	<b>8</b>	88	58	130	&#88;	<b>X</b>	120	78	170	&#120;	<b>x</b>
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57;	<b>9</b>	89	59	131	&#89;	<b>Y</b>	121	79	171	&#121;	<b>y</b>
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58;	<b>:</b>	90	5A	132	&#90;	<b>Z</b>	122	7A	172	&#122;	<b>z</b>
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59;	<b>;</b>	91	5B	133	&#91;	<b>[</b>	123	7B	173	&#123;	<b>{</b>
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60;	<b>&lt;</b>	92	5C	134	&#92;	<b>\</b>	124	7C	174	&#124;	<b> </b>
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61;	<b>=</b>	93	5D	135	&#93;	<b>]</b>	125	7D	175	&#125;	<b>}</b>
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62;	<b>&gt;</b>	94	5E	136	&#94;	<b>^</b>	126	7E	176	&#126;	<b>~</b>
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63;	<b>?</b>	95	5F	137	&#95;	<b>_</b>	127	7F	177	&#127;	<b>DEL</b>

A **karakterkódolás** egy **algoritmus** karakterek digitális formában történő megadására. Ez a karakterek kódszámainak sorozatait „oktetek” sorozatára képezi le.



Például az ISO 10646 karakterkódban az 'a' az 'ä' és az ‰ (ezrelék jel) karaktereknek megfelelő számkódok a 97, a 228 és a 8240.

A karakterkódolás adja meg, hogy a karakterkódokat hogyan képezzük le oktetek sorozatára.

Így például az ISO 10646 egy lehetséges kódolásában minden karakter kódolásához két oktetet használva a kódolási algoritmus ezeknek a (0, 97), a (0, 228) és a (32,48) oktetpárokat felelteti meg.

# Milyen karakterek használhatók a nyelvben?

- **Pascal (ISO 7185:1990):**

- csak minimum követelmények, pl.: '0'..'9' rendezett és folytonos kell legyen, 'A'..'Z' rendezett, de nem feltétlenül folytonos, 'a'..'z' nem kötelező, de ha van, akkor mint a nagybetűk.
- Nincs előírva a lehetséges karakterek halmaza, különböző implementációk különböző kódolást használhatnak(!) (pl. EBCDIC vagy ASCII) => különböző implementációk másképpen viselkedhetnek (pl. '}' ∈ 'A'..'Z' EBCDIC-ben és '}' ∉ 'A'..'Z' ASCII-ben, vagy 'A' < '1' EBCDIC-ben és '1' < 'A' ASCII-ben).
- => ez nem tesz jót a programok hordozhatóságának ☹

- **C++(IISO/IEC 14882:2011 ):**

- **Az alap karakterkészlet 96 elemű: a space, a horizontal tab, vertical tab, new-line, form feed vezérlő karakterei, és a következő 91 grafikus karakter: 'a'..'z' 'A'..'Z' '0'..'9'**

**" # & ' ( ) , \* - + / : ; < > = \_ [ ]  
{ } | . % ? ^ ~ ! \**

- **ISO 10646 szerint kódolva**

- **univerzális karakterek elvben megengedettek**

**\UNNNNNNNN a neve, annak, aminek a rövid neve az ISO/IEC 10646 NNNNNNNN, ill. \uNNNN annak, aminek 0000NNNN, ez a gyakorlatban fordítófüggő, hogy működik-e**

- **=> ez precízebb, mint az 1990-es, ahol nem volt a kód előírva.**

- **CLU, Eiffel:**
  - ASCII kód a megengedett.
- **ADA (ISO/IEC 8652:2012(E)):**
  - A teljes ISO/IEC 10646:2011
- **Java:**
  - "Java programs are written using the Unicode character set"
- **C#:**
  - „A source file is an ordered sequence of Unicode characters.”

# Melyek az elhatároló jelek?

- A legtöbb programozási nyelvben a helyköz (space), a tab, a sorvége elhatároló jel.
- Van sok egy karakterből álló elhatároló jel, mint pl.:

Pascal: & ' ( ) + - \* / : . ; < = >

C++: & % ^ ' ( ) + - { } | ~ [ ] \ " \* / : . ; , < = > ! ?

Java: ( ) { } [ ] ; , . + - \* / : < = > ! ~ ? : & | ^ %

ADA: & ' ( ) \* + , - . / : ; < = > | stb.

# Melyek az elhatároló jelek?

- **Van számos több karakterből álló elhatároló jel is, mint pl. :**

Pascal: := >= <= <> << \*\* ..

C++: -> ++ -- .\* ->\* == <= >= && || << >> <<=  
>>= != .. += -= &= ^= |= :: \*= /=

Java: == <= >= << != && || ++ -- >> >>> += -=  
\*= /= &= ^= %= <<= >>= >>>=

ADA: => .. \*\* := /= >= <= << >> <> -- stb.

# Megkülönbözteti-e a nyelv a kis- és nagybetűket?

**Kutya = kutya = KUTYA = KuTYa = ...?**

- **Pascal, Fortran, CLU, ADA...**: ekvivalensek
- **C, C++, Modula, Perl, Java, C#,...**: minden betű különböző
- **Eiffel**: a kis- és a nagybetűk ekvivalensek, de használhatjuk ugyanazt a nevet egy objektumra és a típusára
  - pl. : a: A megengedett

# Azonosítók

- Milyen karakterek használhatóak az azonosítók leírására?
- Mi az azonosító megengedett szintaxisa?
- Van-e hosszúsági megkötés az azonosítókra?
- Vannak-e kötött szavak?  
Van-e különbség a kulcsszavak és az előre definiált szavak között?



# Milyen karakterek használhatóak azonosítókbán?

- **Pascal** (ISO 7185:1990):
  - betűk ('A'..'Z', 'a'..'z') és számjegyek ('0'..'9')
- **C++**(ISO/IEC 14882:2003):
  - betűk ('A'..'Z', 'a'..'z'), számjegyek ('0'..'9') és '\_'
- **ADA**
  - Az 1995-ös szabványban (ISO/IEC 8652:1995) még: az ISO 10646 BMP Row 00 bármely karaktere, amelynek a neve “Latin Capital Letter” és “Latin Small Letter” (ezek az ún. **identifier\_letters**), számjegyek ('0'..'9') és '\_'
  - új szabvány (Ada2012): minden, ami betű, számjegy! (Πλάτων, Чайковский stb. is)

# Milyen karakterek használhatóak azonosítókbán?

- **Java:**
  - A **Unicode** betűi és számjegyei, a '\$' és '\_'
- **CLU, Eiffel:**
  - betűk ('A'..'Z', 'a'..'z'), számjegyek ('0'..'9') és '\_'
- **Perl:**
  - az előzőeken túl \$, @ és % jellel is kezdődhetnek
    - \$ a skalárokat, @ számmal indexelt tömböt, % asszociatív tömböt jelöl.

# Milyen karakterek használhatóak azonosítókban?

- **C#:**

Az azonosítókat a Unicode szabvány ajánlásának megfelelően kell írni.

Érdekesség, hogy a kulcsszavak az @ bevezető jellel használhatóak azonosítóként is (például: @bool). Ez a jel azonosítót nem vezethet be.

Ezt a lehetőséget a programok hordozhatósága, illetve "átjárhatósága" miatt vezették be. Az @ prefixű szimbólumok neve valójában @ nélkül kerül fordításra. Így attól függetlenül, hogy egy másik nyelv nem tartalmazza például a 'sealed' kulcsszót, és egy programban azonosítóként használják, C#-ban lehetőség van az azonosító közvetlen használatára.

# Mi az azonosítók megengedett szintaxisa?

- A leggyakrabban:

**letter { letter | digit }**

néha az '\_' beszúrása is megengedett.

- **Pascal:** letter{letter|digit}
- **ADA:** letter{[\_]letter|[\_]digit},
  - itt a betűk halmaza nagyobb, minden, aminek a nevében szerepel az, hogy „letter”
- **Java:** Java\_letter{[\_]Java\_letter|[\_]digit},
  - αβγδε vagy El\_Niño is OK.

# Van-e hosszúsági megkötés az azonosítókra?

- **Fortran66, Fortran77:** max. 6 jel
- **Fortran90:** max. 31 jel
- **ADA:** be kell férjen egy sorba...
- **Pascal, Java, ...:** tetszőleges hosszúság

# Vannak-e fenntartott szavak?

- **Pascal, C, C++, Java, Perl, CLU, ADA** etc.:
  - a fenntartott szavak nem használhatók azonosítóként
- **ADA:**
  - különbséget tesz a kulcsszavak és az előredefiniált szavak között, mint pl. : Integer, True, etc..
  - Átdefiniálhatóak:  
type Integer is range -999\_999..+999\_999;  
és így a program implementáció-független lesz  
de: True: Integer ... -- nincs értelme...

## **Fortran, PL/1:**

megengedett a kulcsszavakat azonosítóként használni! ☹️

# Literálok

- Milyen numerikus literálok vannak?
- Milyen más alapok megengedettek a 10-esen kívül?
- Mi az egész, ill. valós literálok szintaxisa?
- Többsoros sztring literálok megengedettek-e?
- Vezérlőkértékek sztring literálokban megengedettek-e, és hogyan írjuk le őket?

## Numerikus literálok:

- A gyakorlatban csak a '0'..'9' és az 'A'..'F'
- néha '\_'
- ' ' (space) soha :-)

## Mi a számok szerkezete?

- **Egészek**: decimális: [-]digit{digit} -123, 456789  
pl. **ADA, Perl, Eiffel**: '\_' is: 456\_789  
**ADA-ban**: exponenciális alak is - pl. 1E6 -, ahol a kitevő nem negatív  
**C, C++, Java, C#**: *int* és *long* különbsége: 'L' a literál végén, pl. 2 vagy 2L. 'U' az unsigned



- **Valóságok:**

`[-]digit{digit}.digit{digit}[exponent]`

`-123.456E+3`

Majdnem minden nyelv ad ehhez valami specialitást:

**Pascal:** `[-]digit{digit}[exponent]`      `12E+3`

**ADA:** ‘\_’ is lehet:      `-1_234.0`

**Eiffel:** `[-]{digit}.{digit}[exponent]`      `-1.`

**Java:** “float\_suffix” is lehet a végén



# Karakterek és karaktersorozatok

- Karakterek: 'A' vagy „escape szekvenciák”: '\n', '\\', '\\ooo' stb.
- Sztringek "A"
- tárolásuk különböző lehet.

# Megjegyzések

- Szerkezete befolyásolja a program megbízhatóságát
- Szokásos lehetőségek:
  - egy speciális oszloptól (jeltől) a sor végéig
    - pl. Fortran 'C' az első oszlopba
  - speciális jel(ek) elején - végén
  - speciális jel(ek) elején - vége a sor vége

# Megjegyzések

- **Pascal:** (\* és \*) vagy { és }
- **CLU:** %-tól sor végéig
- **C:** /\* és \*/
- **C++:** // -tól sor végéig  
vagy /\* és \*/
- **Java:** mint C++-ban, vagy:  
/\*\* documentation \*/
- **C#:** mint C++-ban, vagy:  
/// egysoros dokumentációs megjegyzés
- **ADA, Eiffel:** -- -tól sor végéig