

Кодировки символов

Кодировка символов – способ преобразования последовательности байтов в последовательность символов.

1. Однобайтовые кодировки

Один символ кодируется одним байтом

Максимально возможное количество символов в наборе символов – 256.

Во всех однобайтовых кодировках первая половина таблицы символов совпадает (00h-7Fh), но различается вторая половина таблицы символов (80h-FFh).

Примеры кодировок:

Windows-1251 (ANSI)

| | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 0A | 0B | 0C | 0D | 0E | 0F |
|----|---------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|-------------------|--------------------|--------------------|-------------------|-------------------|-------------------|--------------------|
| 00 | <u>NUL</u> 0000 | <u>STX</u> 0001 | <u>SOT</u> 0002 | <u>ETX</u> 0003 | <u>EOT</u> 0004 | <u>ENQ</u> 0005 | <u>ACK</u> 0006 | <u>BEL</u> 0007 | <u>BS</u> 0008 | <u>HT</u> 0009 | <u>LF</u> 000A | <u>VT</u> 000B | <u>FF</u> 000C | <u>CR</u> 000D | <u>SO</u> 000E | <u>SI</u> 000F |
| 10 | <u>DLE</u> 0010 | <u>DC1</u> 0011 | <u>DC2</u> 0012 | <u>DC3</u> 0013 | <u>DC4</u> 0014 | <u>NAK</u> 0015 | <u>SYN</u> 0016 | <u>ETB</u> 0017 | <u>CAN</u> 0018 | <u>EM</u> 0019 | <u>SUB</u> 001A | <u>ESC</u> 001B | <u>FS</u> 001C | <u>GS</u> 001D | <u>RS</u> 001E | <u>US</u> 001F |
| 20 | <u>SP</u> 0020 | ! | " | # | \$ | % | & | ' | (|) | * | + | , | - | . | / |
| 30 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | : | ; | < | = | > | ? |
| 40 | @ | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 50 | P | Q | R | S | T | U | V | W | X | Y | Z | [| \ |] | ^ | _ |
| 60 | ` | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o |
| 70 | p | q | r | s | t | u | v | w | x | y | z | { | | } | ~ | <u>DEL</u> 007F |
| 80 | Ђ | Ѓ | Ѕ | Ї | Љ | Њ | Ћ | Ќ | Ў | Ъ | Ы | Ь | Э | Ю | Я | а |
| 90 | Ђ | Ѓ | Ѕ | Ї | Љ | Њ | Ћ | Ќ | Ў | Ъ | Ы | Ь | Э | Ю | Я | а |
| A0 | <u>NBSP</u> 00A0 | Ў | Ў | Ј | Ќ | Љ | Њ | Ћ | Ќ | Ў | Ъ | Ы | Ь | Э | Ю | Я |
| B0 | ° | ± | І | і | ґ | µ | ¶ | · | ё | № | е | » | ј | ѕ | ѕ | ї |
| C0 | А | В | В | Г | Д | Е | Ж | З | И | Й | К | Л | М | Н | О | П |
| D0 | Р | С | Т | У | Ф | Х | Ц | Ч | Ш | Щ | Ъ | Ы | Ь | Э | Ю | Я |
| E0 | а | б | в | г | д | е | ж | з | и | й | к | л | м | н | о | п |
| F0 | р | с | т | у | ф | х | ц | ч | ш | щ | ъ | ы | ь | э | ю | я |

Windows-1251 используется в Windows по умолчанию. Кроме русского содержит символы большей части славянских языков.

CP866 (OEM)

| | .0 | .1 | .2 | .3 | .4 | .5 | .6 | .7 | .8 | .9 | .A | .B | .C | .D | .E | .F |
|----|---|---|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 8. | А 410 | Б 411 | В 412 | Г 413 | Д 414 | Е 415 | Ж 416 | З 417 | И 418 | Й 419 | К 41A | Л 41B | М 41C | Н 41D | О 41E | П 41F |
| 9. | Р 420 | С 421 | Т 422 | У 423 | Ф 424 | Х 425 | Ц 426 | Ч 427 | Ш 428 | Щ 429 | Ъ 42A | Ы 42B | Ь 42C | Э 42D | Ю 42E | Я 42F |
| A. | а 430 | б 431 | в 432 | г 433 | д 434 | е 435 | ж 436 | з 437 | и 438 | й 439 | к 43A | л 43B | м 43C | н 43D | о 43E | п 43F |
| B. |  2591 |  2592 |  2593 | 2502 | └ 2524 | ┌ 2561 | ┐ 2562 | └ 2556 | ┌ 2555 | ┐ 2563 | ┐ 2551 | └ 2557 | ┐ 255D | ┐ 255C | └ 255B | └ 2510 |
| C. | └ 2514 | ┐ 2534 | └ 252C | └ 251C | — 2500 | └ 253C | └ 255E | └ 255F | └ 255A | └ 2554 | └ 2569 | └ 2566 | └ 2560 | = 2550 | └ 256C | └ 2567 |
| D. | └ 2568 | └ 2564 | └ 2565 | └ 2559 | └ 2558 | └ 2552 | └ 2553 | └ 256B | └ 256A | └ 2518 | └ 250C | ■ 2588 | ■ 2584 | ■ 258C | ■ 2590 | ■ 2580 |
| E. | р 440 | с 441 | т 442 | у 443 | ф 444 | х 445 | ц 446 | ч 447 | ш 448 | щ 449 | ъ 44A | ы 44B | ь 44C | э 44D | ю 44E | я 44F |
| F. | Ё 401 | ё 451 | Є 404 | є 454 | Ї 407 | ї 457 | Ў 40E | ў 45E | ° B0 | · 2219 | · B7 | √ 221A | № 2116 | х A4 | ■ 25A0 | A0 |

Кодировка, которая используется в консольном режиме (эмуляция DOS). В Интернет используется редко. Кодировка создавалась таким образом, чтобы псевдографические символы, которые используются для рисования таблиц в DOS, были совместимы со стандартной западноевропейской кодировкой (чтобы псевдографика была одинаковой во всех кодировках).

Позиции русских символов не совпадают с Windows-1251, требуется конвертация текстов из одной кодировки в другую.

KOI8-R

| | .0 | .1 | .2 | .3 | .4 | .5 | .6 | .7 | .8 | .9 | .A | .B | .C | .D | .E | .F |
|----|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 8. | — 2500 | 2502 | Г 250C | Г 2510 | Л 2514 | Л 2518 | Т 251C | Т 2524 | Т 252C | Т 2534 | Т 253C | ■ 2580 | ■ 2584 | ■ 2588 | ■ 258C | ■ 2590 |
| 9. | ▒ 2591 | ▒ 2592 | ▒ 2593 | ┌ 2320 | ■ 25A0 | · 2219 | √ 221A | ≈ 2248 | ≤ 2264 | ≥ 2265 | | ┘ 2321 | ○ B0 | ² B2 | · B7 | ÷ F7 |
| A. | = 2550 | 2551 | F 2552 | ё 451 | Г 2553 | Г 2554 | Г 2555 | Г 2556 | Г 2557 | Г 2558 | Г 2559 | Г 255A | Г 255B | Г 255C | Г 255D | Г 255E |
| B. | Г 255F | Г 2560 | Г 2561 | ё 401 | Г 2562 | Г 2563 | Г 2564 | Г 2565 | Г 2566 | Г 2567 | Г 2568 | Г 2569 | Г 256A | Г 256B | Г 256C | © A9 |
| C. | ю 44E | а 430 | б 431 | ц 446 | д 434 | е 435 | ф 444 | г 433 | х 445 | и 438 | й 439 | к 43A | л 43B | м 43C | н 43D | о 43E |
| D. | п 43F | я 44F | р 440 | с 441 | т 442 | у 443 | ж 436 | в 432 | ь 44C | ы 44B | з 437 | ш 448 | э 44D | щ 449 | ч 447 | ъ 44A |
| E. | Ю 42E | А 410 | Б 411 | Ц 426 | Д 414 | Е 415 | Ф 424 | Г 413 | Х 425 | И 418 | Й 419 | К 41A | Л 41B | М 41C | Н 41D | О 41E |
| F. | П 41F | Я 42F | Р 420 | С 421 | Т 422 | У 423 | Ж 416 | В 412 | Ь 42C | Ы 42B | З 417 | Ш 428 | Э 42D | Щ 429 | Ч 427 | Ъ 42A |

Кодировка создавалась для совместимости более старой 7-битовой кодировкой. Русские символы размещены таким образом, что если убрать старший бит байта, то русская буква превратится в латинскую, которая наиболее близка ей по звучанию (фонетически). Поэтому буквы следуют не в алфавитном порядке.

2. Кодировка Unicode (www.unicode.org)

Разрабатывается консорциумом «Юникод».

Для кодировки символа используется 4 байта (от 1 до 4 байтов).

Символы в кодировке Unicode соответствуют универсальному набору символов (Universal Character Set - UCS), он определен в ISO10646 (файл ISO_2003_charts_0000_33FF_BasicLatin_CJKCompat.pdf).

Латинские буквы и цифры – номера с 0000 по 007E

Русские буквы – номера с 0400 по 052F

(0400H = 0000 01XX XXXX XXXX)

В юникод существует возможность комбинации символов из нескольких, например, к символу какой-либо буквы можно добавить символ ударения.

В ссылках на символы (å) используются четырехбайтовые символы, определенные в ISO10646.

(Каждый символ занимает 4 байта и может не быть совместимости с однобайтовыми кодировками.)

Для решения проблемы в Unicode существует 3 формы кодирования: UTF-32, UTF-16, UTF-8

UTF-32

Для кодирования символа всегда используется 4 байта. Прямое соответствие номеру символа.

Наиболее простой вариант кодировки.

UTF-16

Если код символа содержит 16 бит (2 байта) то он кодируется напрямую, если 21 бит то с помощью четырех байтов. Коды символов в Unicode занимают не более 21 бита.

(файл ch03.pdf, стр.41)

Table 3-5. UTF-16 Bit Distribution

| Scalar Value | UTF-16 |
|-----------------------------|------------------------------------|
| xxxxxxxxxxxxxxxxxx | xxxxxxxxxxxxxxxxxx |
| 000uuuuuuxxxxxxxxxxxxxxxxxx | 110110wwwxxxxxxxx 110111xxxxxxxxxx |

Note: www = uuuu - 1

Поскольку большая часть символов занимает не более 16 бит, то по сравнению с UTF-32 получается выигрыш в размере данных (примерно в 2 раза).

UTF-8

Преимущество – кодировка первых 127 символов совпадает с набором символов ASCII (однобайтовой кодировкой). Символы латиницы не изменяются. Наиболее частый вариант кодировки.

Table 3-6. UTF-8 Bit Distribution

| Scalar Value | First Byte | Second Byte | Third Byte | Fourth Byte |
|----------------------------|------------|-------------|------------|-------------|
| 00000000 0xxxxxxx | 0xxxxxxx | | | |
| 00000yyy yxxxxxxx | 110yyyyy | 10xxxxxx | | |
| zzzyyyyy yxxxxxxx | 1110zzzz | 10yyyyyy | 10xxxxxx | |
| 000uuuuu zzzzyyyy yxxxxxxx | 11110uuu | 10uuzzzz | 10yyyyyy | 10xxxxxx |

Текстовые данные должны передаваться с использованием принятого в Интернет порядка байт ("big-endian", байт высшего порядка следует первым).

Также существует порядок "little-endian", байт низшего порядка следует первым.

В начале файла обычно располагаются несколько служебных байтов BOM – byte order mark, которые определяют big-endian или little-endian.

В UTF-8 BOM не имеет значения, так как все передается в виде байтов и их порядок передачи не важен. В 16 и 32 используются многобайтовые значения, поэтому важен порядок байтов и используется BOM.

Указание кодировки в HTML

Для задания кодировки в HTML используется элемент META (в заголовочной секции).

<META http-equiv="Content-Type" content="text/html; charset=КОДИРОВКА">

Этот элемент не производит перекодировку. Он только указывает браузеру в какой кодировке сохранен документ. Если кодировка документа не соответствует META (например документ в кодировке Windows-1251, а в META указана KOI8-R), то документ будет отображаться неправильно.