



КОДИРОВАНИЕ ТЕКСТОВОЙ ИНФОРМАЦИИ

информатика

ТЕКСТОВУЮ ИНФОРМАЦИЮ КОДИРУЮТ ДВОИЧНЫМ КОДОМ ЧЕРЕЗ
ОБОЗНАЧЕНИЕ КАЖДОГО СИМВОЛА АЛФАВИТА ОПРЕДЕЛЕННЫМ ЦЕЛЫМ
ЧИСЛОМ. С ПОМОЩЬЮ ВОСЬМИ ДВОИЧНЫХ РАЗРЯДОВ ВОЗМОЖНО
ЗАКОДИРОВАТЬ 256 РАЗЛИЧНЫХ СИМВОЛОВ. ДАННОГО КОЛИЧЕСТВА
СИМВОЛОВ ДОСТАТОЧНО ДЛЯ ВЫРАЖЕНИЯ ВСЕХ СИМВОЛОВ
АНГЛИЙСКОГО И РУССКОГО АЛФАВИТОВ.

В ПЕРВЫЕ ГОДЫ РАЗВИТИЯ КОМПЬЮТЕРНОЙ ТЕХНИКИ ТРУДНОСТИ
КОДИРОВАНИЯ ТЕКСТОВОЙ ИНФОРМАЦИИ БЫЛИ ВЫЗВАНЫ
ОТСУТСТВИЕМ НЕОБХОДИМЫХ СТАНДАРТОВ КОДИРОВАНИЯ. В
НАСТОЯЩЕЕ ВРЕМЯ, НАПРОТИВ, СУЩЕСТВУЮЩИЕ ТРУДНОСТИ СВЯЗАНЫ
С МНОЖЕСТВОМ ОДНОВРЕМЕННО ДЕЙСТВУЮЩИХ И ЗАЧАСТУЮ
ПРОТИВОРЕЧИВЫХ СТАНДАРТОВ.

ДЛЯ АНГЛИЙСКОГО ЯЗЫКА, КОТОРЫЙ ЯВЛЯЕТСЯ НЕОФИЦИАЛЬНЫМ
МЕЖДУНАРОДНЫМ СРЕДСТВОМ ОБЩЕНИЯ, ЭТИ ТРУДНОСТИ БЫЛИ
РЕШЕНЫ. ИНСТИТУТ СТАНДАРТИЗАЦИИ США ВЫРАБОТАЛ И ВВЕЛ В
ОБРАЩЕНИЕ СИСТЕМУ КОДИРОВАНИЯ ASCII (AMERICAN STANDARD CODE FOR
INFORMATION INTERCHANGE – СТАНДАРТНЫЙ КОД ИНФОРМАЦИОННОГО
ОБМЕНА США).

Для кодировки русского алфавита были разработаны несколько вариантов кодировок:

1) Windows-1251 – введена компанией Microsoft; с учетом широкого распространения операционных систем (ОС) и других программных продуктов этой компании в Российской Федерации она нашла широкое распространение;

2) КОИ-8 (Код Обмена Информацией, восьмизначный) – другая популярная кодировка русского алфавита, распространенная в компьютерных сетях на территории Российской Федерации и в российском секторе Интернет;

3) ISO (International Standard Organization – Международный институт стандартизации) – международный стандарт кодирования символов русского языка. На практике эта кодировка используется редко.

Код — правило (алгоритм) сопоставления каждому конкретному сообщению строго определённой комбинации символов (знаков) (или сигналов)
Представляет собой систему условных знаков для представления информации.

Кодирование – перевод информации в удобную для передачи, обработки, хранения формы с помощью некоторого кода.

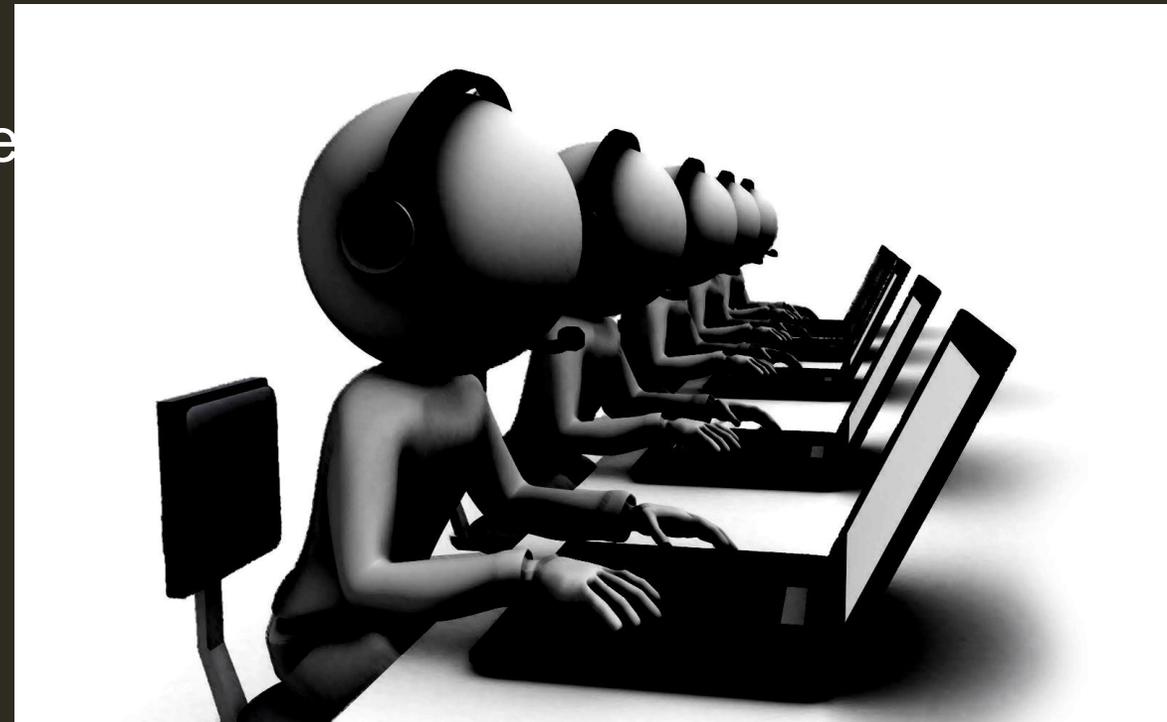
Декодирование – процесс восстановления содержания закодированной информации.

Алфавит – множество символов, с помощью которых записывается текст.

Мощность – число символов алфавите.

Проблемы кодирования:

- 1) Отсутствие информации о кодировке
- 2) Программа не знает кодировки
- 3) Избыток информации о кодировках



Ключевые свойства текстовых материалов:

- ценность
- новизна
- полезность
- адекватность
- истинность

Текстовая информация - последовательность символов, печатных знаков, которые принадлежат тому или иному набору символов. Может храниться в формативном или не нормативном виде.

ПОДРОБНЕЕ

Ценность информации – одно из важнейших свойств информации, оценка которого зависит от целей процессов её генерации и обработки.

Полезность информации – информация, которая имеет значение.

Адекватность информации – уровень соответствия образа, создаваемого с помощью информации, реальному объекту.



КЛАССИФИКАЦИЯ ИНФОРМАЦИИ:

- синтаксическая – отображает формально-структурные характеристики информации;
- прагматическая – отображает соответствие информации цели управления, реализуемой на её основе;
- семантическая – определяет степень соответствия образа объекта самому объекту.

На сегодняшний день большое количество пользователей при помощи компьютера обрабатывает текстовую информацию, которая состоит из: букв, цифр, знаков препинания и других элементов.

Обычно для кодирования одного символа, используется 1 байт памяти то есть 8 бит. По теории вероятностей с помощью простой формулы, которая связывает количество возможных событий (K) и количество информации (I), можно вычислить сколько не одинаковых символов можно закодировать: $K = 2^I = 2^8 = 256$.

Принцип данного кодирования заключается в том, что каждому символу (букве, знаку) соответствует свой двоичный код от 00000000 до 11111111, так-же текстовая информация может быть представлена в десятичном коде от 0 до 255.

Нужно запомнить, что на сегодняшний день для кодирования букв русского алфавита используют пять разных кодировочных таблиц (КОИ - 8, CP1251, CP866, Mac, ISO), запомните, что тексты закодированные с помощью одной таблицы не будут корректно отображаться в другой кодировке. Это можно увидеть в объединенной таблице кодировки символов.