

Кодирование информации

1. Язык и кодирование
2. Двоичное кодирование
3. Кодирование чисел и символов
4. Кодирование рисунков
5. Кодирование звука и видео

Кодирование информации

Тема 1. Язык и кодирование

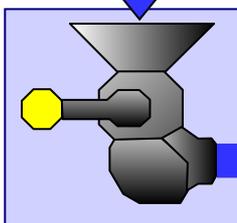
Что такое кодирование?

Кодирование – это запись информации с помощью некоторой знаковой системы (языка).



Зачем кодируют информацию?

кодирование



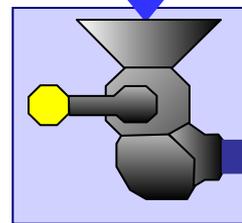
данные (код)

10101001010

передача

Информация передается, обрабатывается и хранится в виде кодов.

борьба с помехами
(специальные способы кодирования)



данные (код)

11111100010

передача

обработка



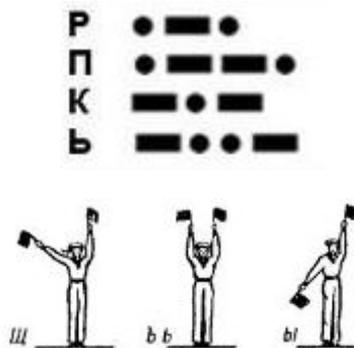
хранение

Языки

Язык – знаковая система, используемая для хранения и передачи информации.

- **естественные** (русский, английский, ...)
есть правила и исключения
- **формальные** (строгие правила)

$$E = mc^2$$



```
program qq;
begin
writeln("Привет!");
end.
```

$$16 = 10_{16} = 20_8 = 10000_2$$

Грамматика – правила по которым из символов алфавита строятся слова.

Синтаксис – правила, по которым из слов строятся предложения.

Азбука Морзе

Задача 1. Закодируйте свое имя с помощью азбуки Морзе.

А	● —	П	● — — — ●	Ь	— ● ● —
Б	— ● ● ●	Р	● — ●	Ы	— ● — —
В	● — —	С	● ● ●	Й	● — — —
Г	— — ●	Т	—		
Д	— ● ●	У	● ● —	1	● — — — —
Е	●	Ф	● ● — ●	2	● ● — — —
Ж	● ● ● —	Х	● ● ● ●	3	● ● ● — —
З	— — — ● ●	Ц	— ● — — ●	4	● ● ● ● —
И	● ●	Ч	— — — — ●	5	● ● ● ● ●
К	— ● — —	Ш	— — — — —	6	— ● ● ● ●
Л	● — — ● ●	Щ	— — — ● —	7	— — — ● ● ●
М	— — —	Э	● ● — — ● ●	8	— — — — ● ●
Н	— — ●	Ю	● ● — — —	9	— — — — — ●
О	— — — —	Я	● — — ● —	0	— — — — —

ВАСЯ



Код неравномерный, нужен разделитель!

Кодовые таблицы

Задача 2. Закодируйте свое имя с помощью кодовой таблицы (*Windows-1251*):

	0	1	2	3	4	5	6	7	8	9	А	В	С	Д	Е	Ф
С	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
Д	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я

В А С Я

ВАСЯ

С2 С0 Д1 ДF



Код равномерный, разделитель **HE** нужен!

Цели и способы кодирования

Текст:

- в России: **Привет, Вася!**
- Windows-1251: **CFF0E8E2E52C20C2E0F1FF21**
- передача за рубеж (транслит): **Privet, Vasya!**
- стенография: 
- шифрование: **Рсѷгжу-!Гбта”**

Числа:

- для вычислений: **25**
- прописью: **двадцать пять**
- римская система: **XXV**



Как зашифровано?



Информация (смысл сообщения) может быть закодирована разными способами!

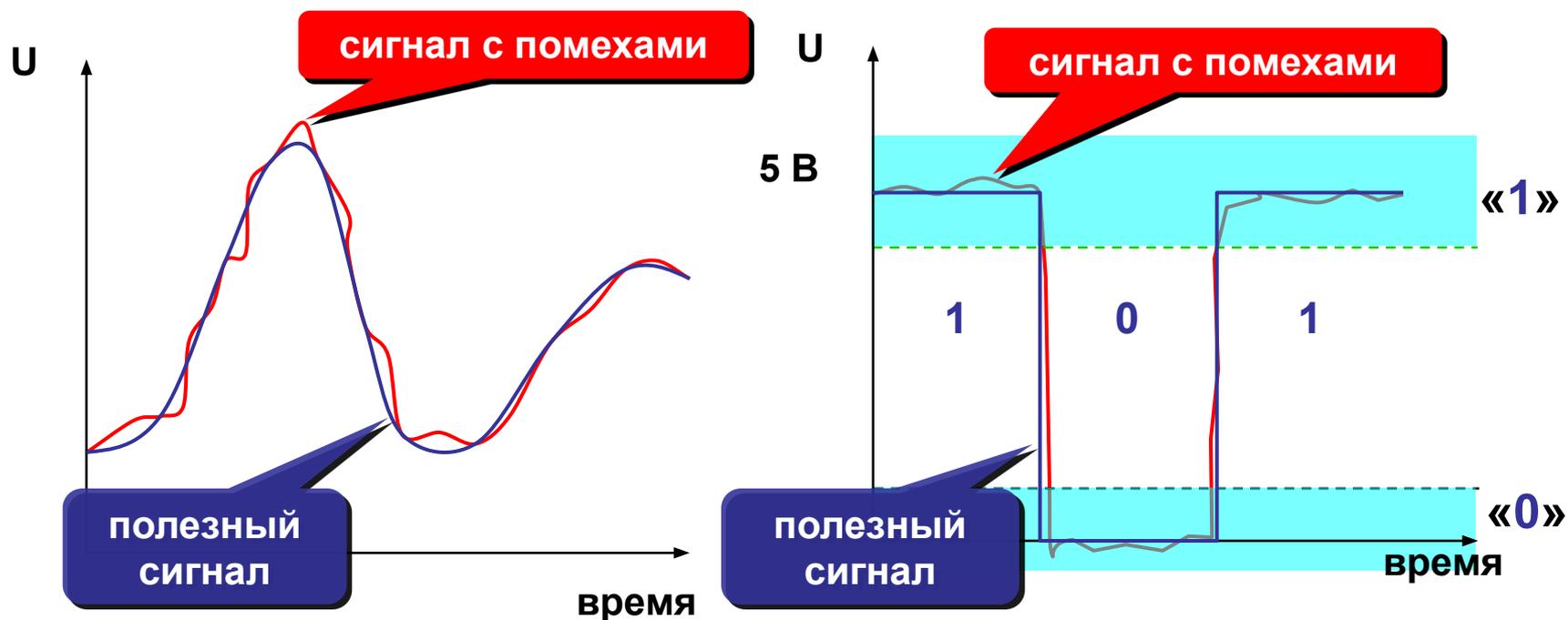
Кодирование информации

Тема 2. Двоичное кодирование

Двоичное кодирование

Двоичное кодирование – это кодирование всех видов информации с помощью двух знаков (обычно 0 и 1).

Передача электрических сигналов:



Двоичное кодирование



- в такой форме можно закодировать (почти) **все виды** информации
- нужны только устройства с **двумя состояниями**
- почти **нет ошибок** при передаче данных
- **компьютеру легче** обрабатывать данные



человеку сложно воспринимать двоичные коды



Можно ли использовать не «0» и «1», а другие символы, например, «А» и «Б»?

Декодирование

Декодирование – это восстановление сообщения из последовательности кодов.

М	А	Ы	Л	У	пробел
00	1	01	0	10	11

МАМА МЫЛА ЛАМУ → 00 1 00 1 11 00 01 0 1 11 0 1 00 10

Приняли сообщение:

0010011100010111010010 → ???

ЛЛАЛЛАААЛЛЛАЛАААЛЛАЛЛАЛ



Не все коды допускают однозначное декодирование!



Почему?

Равномерные коды

Равномерные коды – все кодовые слова (коды отдельных букв) имеют одинаковую длину.

М	А	Ы	Л	У	пробел
000	001	010	011	100	101

МАМА МЫЛА ЛАМУ:

000 001 000 001 101 000 010 011 001 101 011 001 000 100



Равномерные коды позволяют однозначно декодировать сообщения!

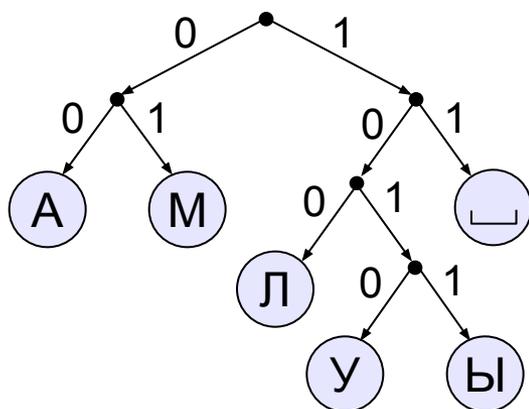


сообщения получаются длинными

Неравномерные коды

кодовые слова имеют разную длину

М	А	Ы	Л	У	пробел
01	00	1011	100	1010	11



0100010011011011100001110000011010
 М А М А ▭ М Ы Л А ▭ Л А М У

Префиксный код – ни одно кодовое слово не совпадает с началом другого кодового слова (условие Фано).



Любой префиксный код позволяет однозначно декодировать сообщения!

Постфиксные коды

Постфикс = окончание слова.

Постфиксный код – ни одно кодовое слово не совпадает с концом другого кодового слова («обратное» условие Фано).

М	А	Ы	Л	У	пробел
10	00	1101	001	0101	11



Любой постфиксный код позволяет однозначно декодировать сообщения (с конца)!



для декодирования нужно получить всё сообщение целиком

Задачи на построение кода

Для передачи по каналу связи сообщения, состоящего только из букв А, Б, В, Г, решили использовать неравномерный по длине код:

А	Б	В	Г
1	000	001	?

Как нужно закодировать букву Г, чтобы длина кода была минимальной и допускалось однозначное разбиение кодированного сообщения на буквы?

- 1) 00 2) 01 3) 11 4) 010

Решение:

- 1) для букв А-Б-В выполняются условие Фано
- 2) при Г=00 условие Фано нарушится (пары Г-Б, Г-В)
- 3) при Г=01 условие Фано выполняется**
- 4) при Г=11 условие Фано нарушится (пара А-Г)
- 5) при Г=010 условие Фано выполняется (но длиннее 01)

Кодирование информации

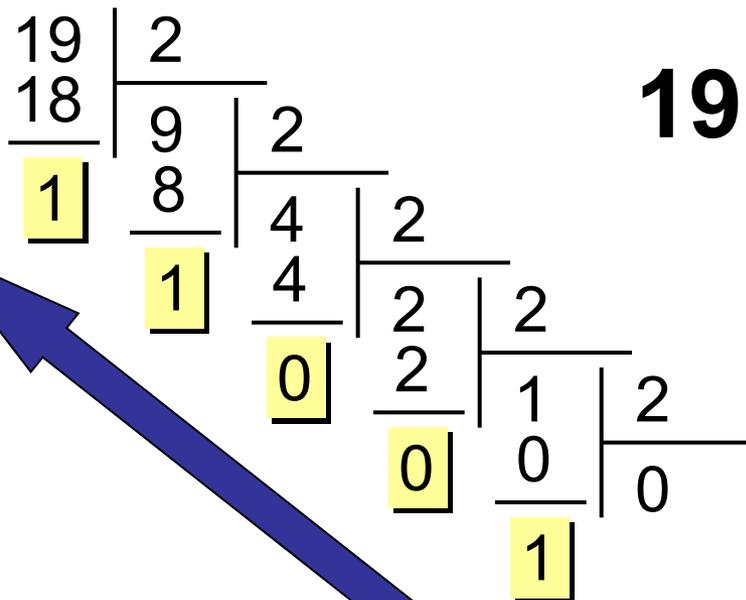
Тема 2. Кодирование чисел и СИМВОЛОВ

Кодирование чисел (двоичная система)

Алфавит: 0, 1

Основание (количество цифр): 2

10 → 2



$$19 = 10011_2$$

система
счисления

2 → 10

4 3 2 1 0 разряды

$$10011_2 = 1 \cdot 2^4 + \cancel{0 \cdot 2^3} + \cancel{0 \cdot 2^2} + 1 \cdot 2^1 + 1 \cdot 2^0$$

$$= 16 + 2 + 1 = 19$$

Кодирование символов

Текстовый файл

- на экране (символы)
- в памяти – двоичные коды



1000001 ₂	1000010 ₂	1000011 ₂	1000100 ₂
65	66	67	68



В файле хранятся не изображения символов, а их числовые коды в двоичной системе!

А где же хранятся изображения?

Кодирование символов

1. **Сколько символов** надо использовать одновременно? **256** или 65536 (UNICODE)
2. **Сколько места** надо выделить **на символ**:
 $256 = 2^8 \Rightarrow 8$ бит на символ
3. Выбрать **256 любых символов** (или 65536) - **алфавит**.
4. Каждому символу – **уникальный код 0..255** (или 0..65535). Таблица символов:

коды	65	66	67	68		
...	A	B	C	D	...	

5. Коды – в **двоичную систему**.

8-битные кодировки (1 байт на символ)

0	1		127	128		254	255
таблица ASCII (международная)				расширение (национальный алфавит)			

ASCII = *American Standard Code for Information Interchange*

0-31 управляющие символы:

7 – звонок, 10 – новая строка, 13 – возврат каретки, 27 – Esc.

32 пробел

знаки препинания: . , : ; ! ?

специальные знаки: + - * / () { } []

48-57 цифры 0..9

65-90 заглавные латинские буквы **A-Z**

97-122 строчные латинские буквы **a-z**

Кодовая страница (расширенная таблица ASCII)

для русского языка:

CP-866 для системы *MS DOS*

CP-1251 для системы *Windows* (Интернет)

КОИ8-Р для системы *UNIX* (Интернет)

8-битные кодировки (1 байт на символ)

-  1 байт на символ – файлы небольшого размера!
- просто обрабатывать в программах
-  нельзя использовать символы разных кодовых страниц одновременно (русские и французские буквы, и т.п.)
- неясно, в какой кодировке текст (перебор вариантов!)
- для каждой кодировки нужен свой шрифт (изображения символов)

Стандарт UNICODE



Идея: объединить все символы в одну таблицу!

- 110 182 символа (2012)
- каждому символу присвоен код
кириллица:
А – 0410_{16} , Б – 0411_{16} , ...
а – 0430_{16} , б – 0431_{16} , ...
- коды $0..10FFFF_{16}$, всего 1 114 112



UNICODE в Windows (UTF-16)

- общеупотребительные символы
 $0..65535 = 2^{16}-1$ ($0..FFFF_{16}$)
- ЭТИ СИМВОЛЫ МОЖНО закодировать с помощью 16 бит
- кодировка UTF-16 (почти все символы по 16 бит)

 можно одновременно использовать символы разных языков (Интернет)

 размер файла увеличивается

UNICODE в Linux (кодировка UTF-8)

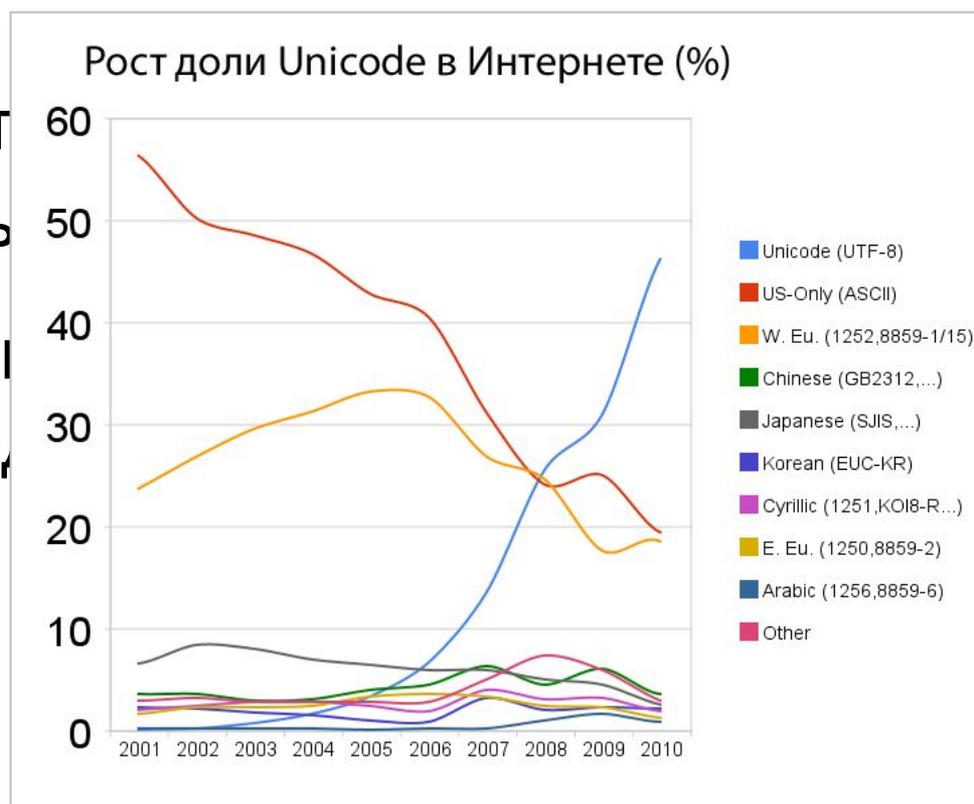
- символы ASCII – 1 байт на символ
- остальные символы от 2 до 4 байт
- более 50% сайтов используют UTF-8



• текст
(кодировка)



• перекодировка
• замедление



• кодировка ASCII
• в размере
• символ

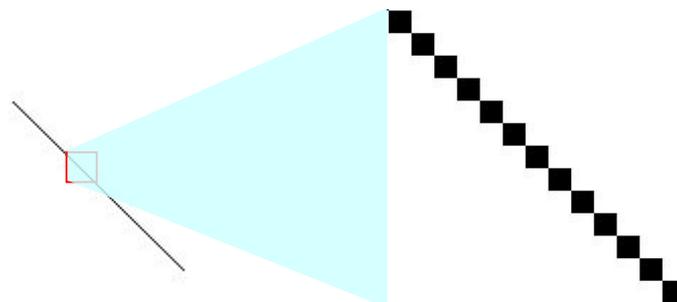
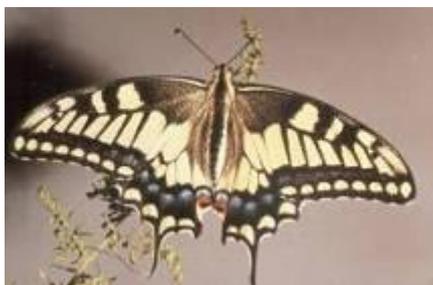
Кодирование информации

Тема 4. Кодирование рисунков

Два типа кодирования рисунков

• растровое кодирование

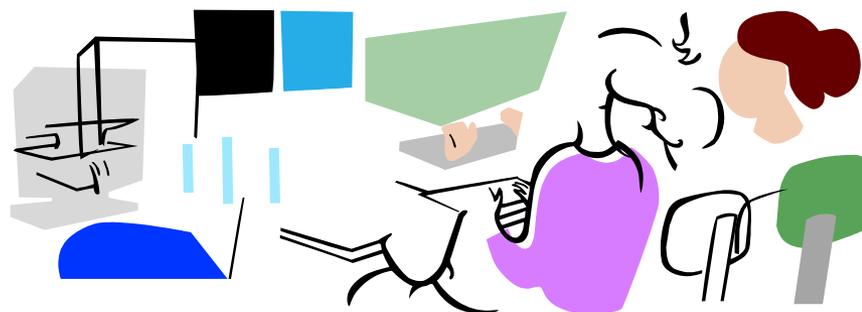
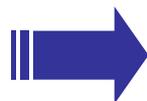
точечный рисунок, состоит из **пикселей**



фотографии, размытые изображения

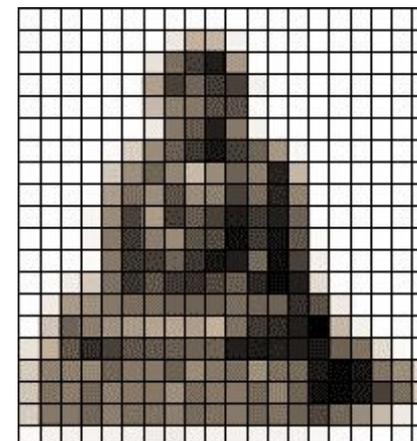
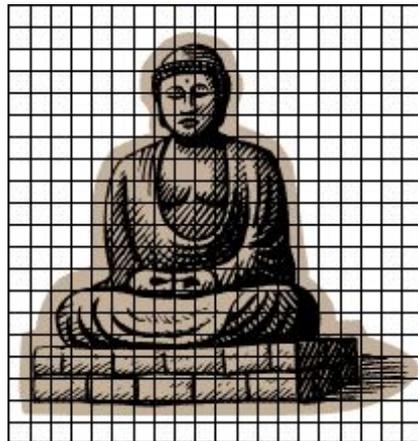
• векторное кодирование

рисунок, состоит из **отдельных геометрических фигур**



чертежи, схемы, карты

Растровое кодирование



Шаг 1. Дискретизация:
разбивка на *пиксели*.

Пиксель – это наименьший элемент рисунка, для которого можно независимо установить цвет.

Шаг 2. Для каждого пикселя определяется **единый цвет**.



Есть потеря информации!

- почему?
- как ее уменьшить?

Разрешение: число пикселей на дюйм, *pixels per inch (ppi)*
экран **96** ppi, печать **300-600** ppi, типография **1200** ppi

Растровое кодирование (*True Color*)

Шаг 3. От цвета – к числам: модель RGB

цвет = **R** + **G** + **B**

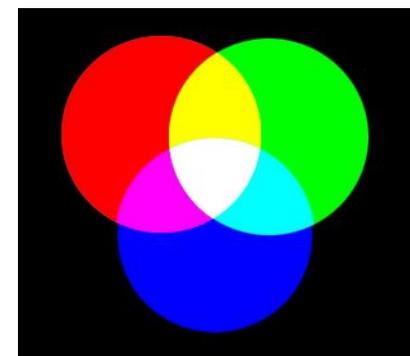
<i>red</i>	<i>green</i>	<i>blue</i>
красный	зеленый	синий
0..255	0..255	0..255



R = 218
G = 164
B = 32



R = 135
G = 206
B = 250



Шаг 4. Числа – в двоичную систему.



Сколько разных цветов можно кодировать?

$256 \cdot 256 \cdot 256 = 16\,777\,216$ (*True Color*)

Глубина
цвета



Сколько памяти нужно для хранения цвета 1 пикселя?

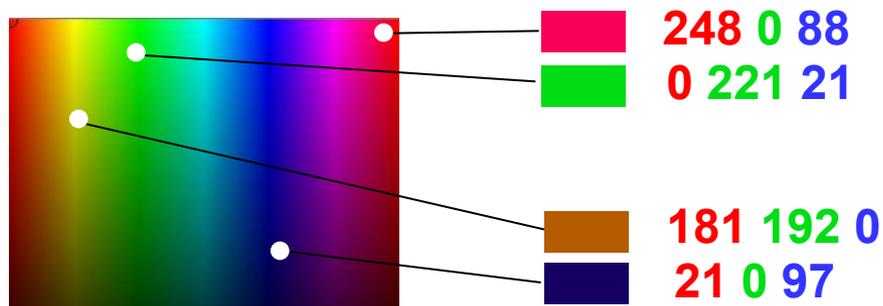
R: $256=2^8$ вариантов, нужно 8 бит = 1 байт

R G B: всего 3 байта

Растровое кодирование с палитрой

Шаг 1. Выбрать количество цветов: 2, 4, ... 256.

Шаг 2. Выбрать 256 цветов из палитры:



Шаг 3. Составить палитру (каждому цвету – номер 0..255)
палитра хранится в начале файла

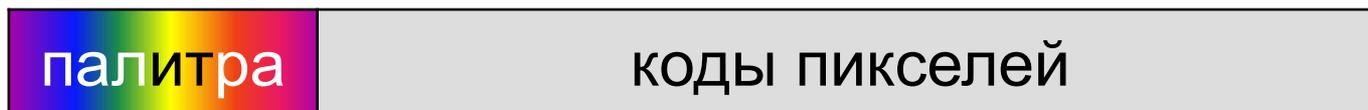
0	1	...	45	...	65	...	254	255
248 0 88	0 221 21	...	181 192 0	...	21 0 97	...	161 12 20	19 23 90

Шаг 4. Код пикселя = номеру его цвета в палитре

65	1	45	14	...	12	23
----	---	----	----	-----	----	----

Растровое кодирование с палитрой

Файл с палитрой:



Сколько занимает палитра и основная часть?

Один цвет в палитре: **3 байта (RGB)**

256 = 2⁸ цветов:

палитра	$256 \cdot 3 = 768$ байт
рисунок	8 бит на пиксель

16 цветов:

палитра	$16 \cdot 3 = 48$ байт
рисунок	4 бита на пиксель

2 цвета:

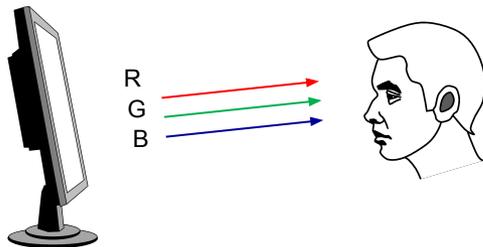
палитра	$2 \cdot 3 = 6$ байт
рисунок	1 бит на пиксель

Глубина
цвета

Форматы файлов (растровые рисунки)

Формат	True Color	Палитра	Прозрачность
BMP			
JPG			
GIF			
PNG			

Кодирование цвета при печати



Белый – красный

= голубой

C = Cyan

Белый – зелёный

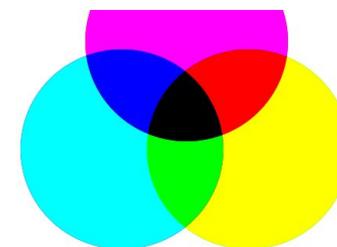
= пурпурный

M = Magenta

Белый – синий

= желтый

Y = Yellow



Модель CMY



Модель CMYK: + **Key color**

Меньший расход краски и лучшее качество для чёрного и серого цветов.

Растровые рисунки

-  • лучший способ для хранения **фотографий** и изображений без четких границ
 - **спецэффекты** (тени, ореолы, и т.д.)
-  • есть **потеря информации** (почему?)
 - при изменении размеров рисунка он **искажается**
 - **размер файла** не зависит от сложности рисунка (а от чего зависит?)



Какие свойства цифрового рисунка определяют его качество?

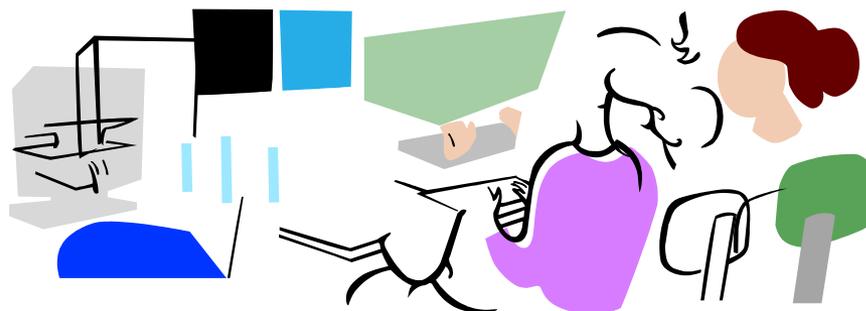
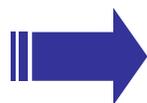
Векторные рисунки

Строятся из геометрических фигур:

- отрезки, ломаные, прямоугольники
- окружности, эллипсы, дуги
- сглаженные линии (кривые Безье)

Для каждой фигуры в памяти хранятся:

- размеры и координаты на рисунке
- цвет и стиль границы
- цвет и стиль заливки (для замкнутых фигур)



Форматы файлов:

- **WMF** (*Windows Metafile*)
- **AI** (*Adobe Illustrator*)
- **CDR** (*CorelDraw*)
- **SVG** (*Inkscape*)

для Web

Векторные рисунки

```
<svg>
```

прямоугольник

размеры

```
<rect width="135" height="30"  
  x="0" y="10"  
  stroke-width="1" stroke="rgb(0,0,0)"  
  fill="rgb(255,255,255)" />
```

координаты

контур

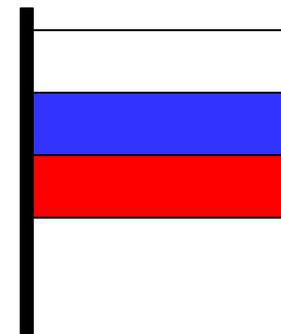
заливка

```
<rect width="135" height="30" x="0" y="40"  
  stroke-width="1" stroke="rgb(0,0,0)"  
  fill="rgb(0,0,255)" />
```

```
<rect width="135" height="30" x="0" y="70"  
  stroke-width="1" stroke="rgb(0,0,0)"  
  fill="rgb(255,0,0)" />
```

```
<line x1="0" y1="0"  
  x2="0" y2="150"  
  stroke-width="15" stroke="rgb(0,0,0)" />
```

```
</svg>
```



Векторные рисунки



- лучший способ для хранения **чертежей, схем, карт**;
- при кодировании **нет потери информации**;
- при изменении размера **нет искажений**;
- меньше **размер файла**, зависит от сложности рисунка;

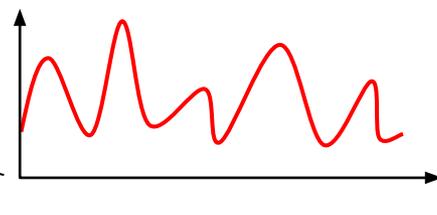
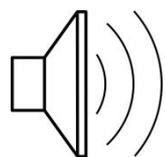


- неэффективно использовать для **фотографий** и размытых изображений

Кодирование информации

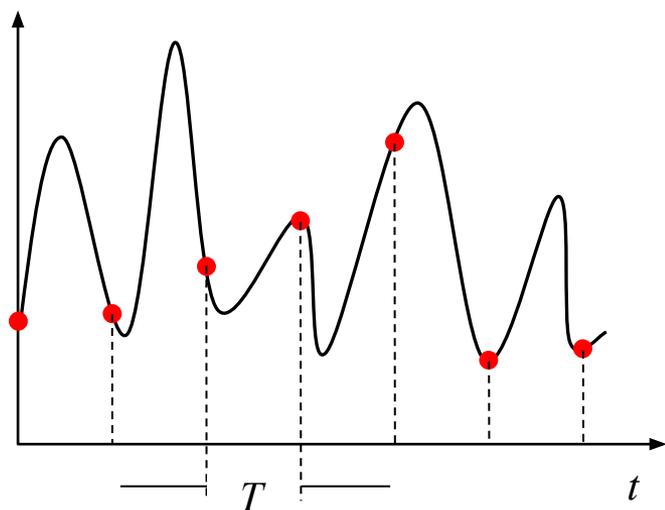
Тема 5. Кодирование звука и видео

Оцифровка звука



аналоговый
сигнал

Оцифровка – это преобразование аналогового сигнала в цифровой код (дискретизация).



T – интервал дискретизации (с)
 $f = \frac{1}{T}$ – частота дискретизации
 (Гц, кГц)

8 кГц – минимальная частота для
распознавания речи

11 кГц, 22 кГц,

44,1 кГц – качество CD-дисков

48 кГц – фильмы на DVD

96 кГц, 192 кГц

Человек слышит

16 Гц ... 20 кГц

Оцифровка звука: квантование

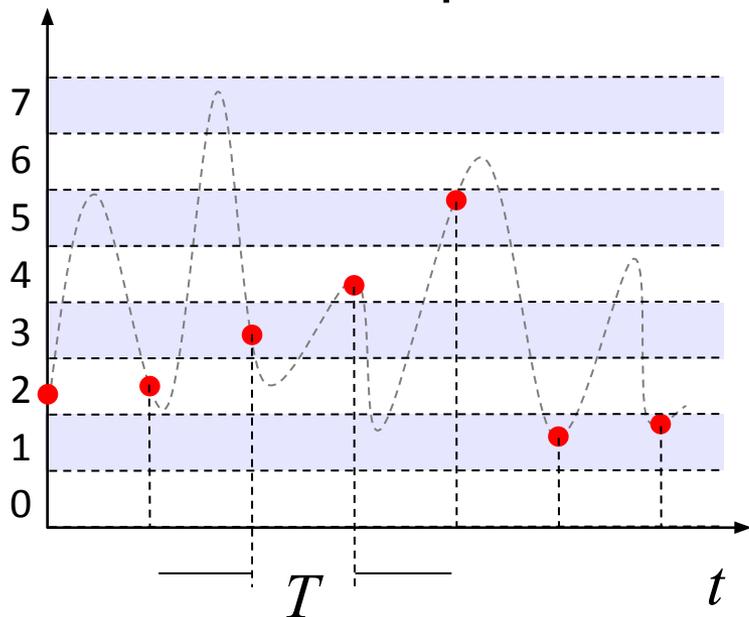


Сколько битов нужно, чтобы записать число 0,6?

Квантование (дискретизация по уровню) – это представление числа в виде цифрового кода конечной длины.

АЦП = Аналого-Цифровой Преобразователь

3-битное кодирование:



8 битов = 256 уровней

16 битов = 65536 уровней

24 бита = 2^{24} уровней

Разрядность кодирования — это число битов, используемое для хранения одного отсчёта.

Оцифровка звука

Задача. Определите информационный объем данных, полученных при оцифровке звука длительностью 1 минута с частотой 44 кГц с помощью 16-битной звуковой карты. Запись выполнена в режиме «стерео».

За 1 сек *каждый* канал записывает 44000 значений,
каждое занимает 16 битов = 2 байта
всего $44000 \cdot 2 \text{ байта} = 88000 \text{ байтов}$

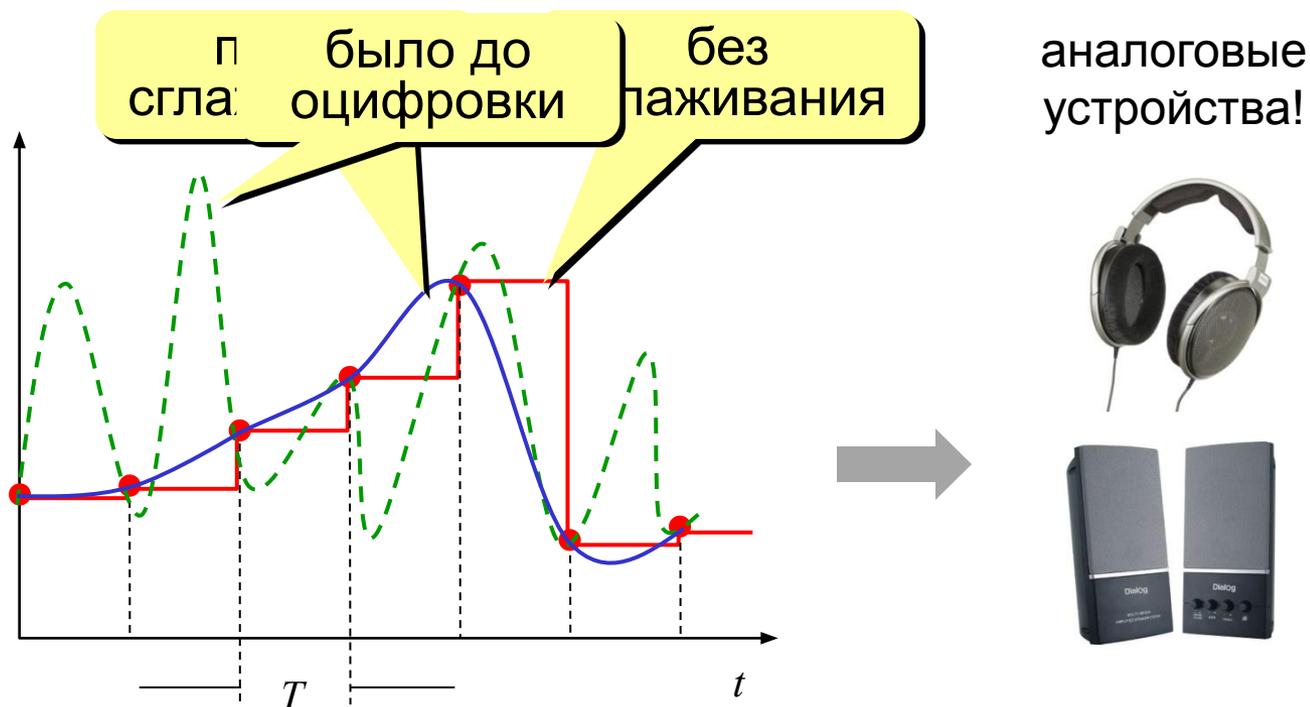
С учётом «стерео»
всего $88000 \cdot 2 = 176000 \text{ байтов}$

За 1 минуту
 $176000 \cdot 60 = 1056000 \text{ байтов}$
 $\approx 10313 \text{ Кбайт} \approx 10 \text{ Мбайт}$

Оцифровка звука

Как восстановить сигнал?

ЦАП = Цифро-Аналоговый Преобразователь



? Какой улучшить качество?

уменьшать T

? Что при этом ухудшится?

↑ размер файла

Оцифровка – итог

 можно закодировать **любой звук** (в т.ч. ГОЛОС, СВИСТ, шорох, ...)

 • есть **потеря информации**
• большой **объем файлов**

 Какие свойства оцифрованного звука определяют качество звучания?

Форматы файлов:

WAV (*Waveform audio format*), часто без сжатия (размер!)

MP3 (*MPEG-1 Audio Layer 3*, сжатие с учётом восприятия человеком)

AAC (*Advanced Audio Coding*, 48 каналов, сжатие)

WMA (*Windows Media Audio*, потоковый звук, сжатие)

OGG (*Ogg Vorbis*, открытый формат, сжатие)

Инструментальное кодирование

MIDI (*Musical Instrument Digital Interface* — цифровой интерфейс музыкальных инструментов).

в файле **.mid**:

- нота (высота, длительность)
- музыкальный инструмент
- параметры звука (громкость, тембр)
- до 1024 каналов

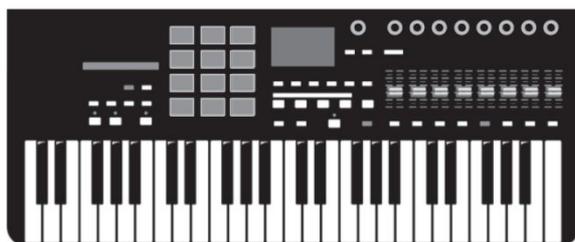
128 мелодических
и 47 ударных

программа для
звуковой карты!

в памяти звуковой карты:

- образцы звуков (волновые таблицы)

MIDI-клавиатура:



- **нет потери информации** при кодировании инструментальной музыки
- **небольшой размер файлов**



невозможно закодировать нестандартный звук, голос

Трекерная музыка

В файле (модуле):

- образцы звуков (*сэмплы*)
- нотная запись, трек (*track*) – дорожка
- музыкальный инструмент
- до 32 каналов

Форматы файлов:

MOD разработан для компьютеров *Amiga*

S3M оцифрованные каналы + синтезированный звук, 99 инструментов

XM, STM, ...

Использование: демосцены (важен размер файла)

Кодирование видео



Видео = изображения + звук Синхронность!

изображен

- ≥ 25 кадр
- **PAL**: 768
за 1
за 1
- **HDTV**: 1
- **ИСХОДНЫ**
- **сжатие** (
- **DivX**, **Xv**

звук:

- 48 кГц, 1
- **сжатие** (
- **MP3**, **AA**



Форматы видеофайлов

- AVI** – *Audio Video Interleave* – чередующиеся звук и видео; контейнер – могут использоваться разные кодеки
- MPEG** – *Motion Picture Expert Group*
- WMV** – *Windows Media Video*, формат фирмы *Microsoft*
- MP4** – *MPEG-4*, сжатое видео и звук
- MOV** – *Quick Time Movie*, формат фирмы *Apple*
- WebM** – открытый формат, поддерживается браузерами

Конец фильма

ПОЛЯКОВ Константин Юрьевич

д.т.н., учитель информатики высшей категории,
ГОО СОШ № 163, г. Санкт-Петербург

kpolyakov@mail.ru