

# Основы кодирования

Михаил Шихов

[m.m.shihov@gmail.com](mailto:m.m.shihov@gmail.com)

Лекция по дисциплине «информатика»  
(15 июня 2016 г.)

# Содержание

## 1 Оптимальное кодирование

- Энтропия
- Алгоритм Хаффмана для  $m = 2$
- Алгоритм Фано для  $m = 2$

## 2 Кодирование с целью сжатия информации

- Сжатие
- Алгоритм Лемпела-Зива

## 3 Кодирование с целью защиты свойств информации

- Защита целостности
- Защита конфиденциальности
- Защита принадлежности

# Определение информации

## Definition (Юридическое)

**Информация** — это сведения (сообщения, данные) независимо от формы их представления<sup>а</sup>.











<sup>а</sup> №149-ФЗ от 27 июля 2006 г «Об информации, информационных технологиях и о защите информации»

## Definition

**Информация** — это упорядоченная последовательность (цепочка) **кодовых символов**, принадлежащих конечному алфавиту. При этом каждый символ последовательности несёт определённую смысловую нагрузку.



# Уровни доступа к информации

Носителя			
Средств взаимодействия			
Представления (код, формат)	<b>СЛОВО</b> <b>word</b> كلمة 字	<b>NTFS FAT32</b> <b>ExFAT</b>  *.txt *.xls *.png *.htm *.doc *.svg	
Семантики (понимания)			

# Трудности оценки количества информации

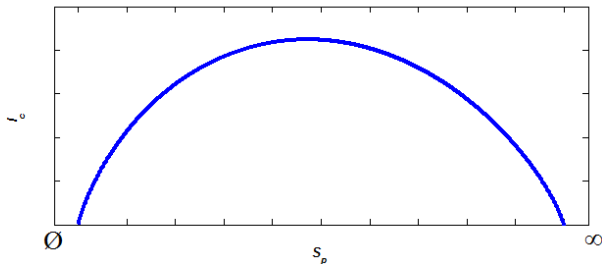
- Просто оценить количество явно заданной информации: достаточно посчитать количество кодовых символов последовательности.
- Сложно оценить необходимое количество информации для адекватного «отражения» исходного объекта.

# Семантический подход к количественной оценке

Зависимость  $I_c$  в сообщении от тезауруса  $S_p$  получателя

$$I = I_s + I_c + I_n.$$

где  $I_s$  — известна,  $I_c$  — неизвестна и понятна,  $I_n$  — шум.  $C = \frac{I_c}{I}$  называется коэффициентом содержательности информации.



# Прагматический подход к количественной оценке

Оценка Александра Александровича Харкевича

$$I = \log_m \frac{p_2}{p_1} = \log_m p_2 - \log_m p_1, \quad (1)$$

где  $m$  — основание логарифма, определяющее единицы измерения ( $m > 1$ ),  $p_1$  — вероятность достижения потребителем цели до получения информации,  $p_2$  — вероятность достижения потребителем цели после получения информации. Ценность информации в случае  $p_1 > p_2 > 0$  положительна, в случае  $0 < p_1 < p_2$  отрицательна, а в случае  $p_1 = p_2$  равна нулю.



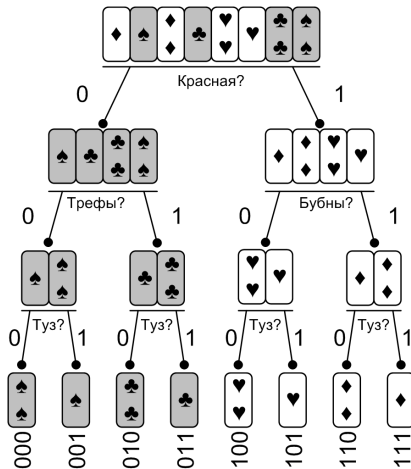
# Синтаксический подход

## Задача о картах (постановка)

### Example

Имеется колода из восьми карт. По две карты (туз и двойка) каждой масти. Некто вытягивает наугад карту и готов честно отвечать только да или нет на любые задаваемые вопросы. Требуется минимальным количеством вопросов угадать вытянутую карту.

# Задача о картах (решение)



# Количественная оценка Ральфа Хартли

$$m^H \geq N$$

$$H = \log_m N,$$

где  $m$  — количество **кодовых символов**;  $N$  — количество состояний **отражаемого объекта**.

## Example

В случае примера с картами: количество состояний  $N = 8$ , количество символов  $m = 2$ . Количество информации по Хартли:  $H = \log_2 8 = 3$  бита.

# Количественная оценка по Клоду Шеннону

## Постулаты

Отражаемый объект — источник **событий**.

- 1 Количество информации есть непрерывная функция от вероятности события.
- 2 Количество информации  $I_i$  одиночного  $i$ -го события  $s_i \in S$ ,  $1 \leq i \leq N$  происходящего с вероятностью  $p_i$ , имеет положительное значение.

$$I_i \geq 0; I_i = I(p_i); \sum_{i=1}^N p_i = 1.$$

- 3 Количество информации  $I_{ij}$  двух независимых событий  $s_i, s_j \in S$  с вероятностью  $p_{ij} = p_i \cdot p_j$ , равно сумме количеств информации событий в отдельности:  $I_{ij} = I_i + I_j$ ;  $I(p_i \cdot p_j) = I(p_i) + I(p_j)$ .

# Количественная оценка по Клоду Шеннону

## Зависимость информации от вероятности

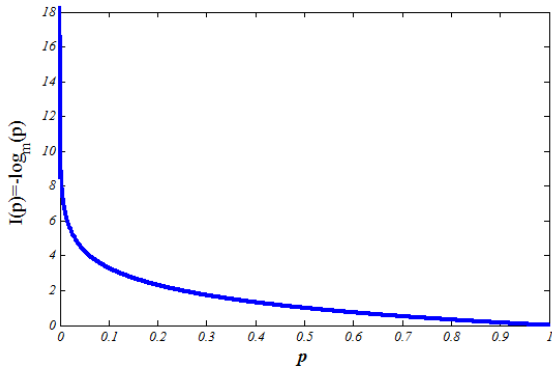
$$I(p) = -\log_m(p),$$

где  $I(p)$  — информация события, происходящего с вероятностью  $p$ ;  $m$  — количество **КОДОВЫХ СИМВОЛОВ**.

Example ( $m$  определяет единицы измерения информации)

- $m = 2$ . **бит**.
- $m = e$ . **нат**.
- $m = 3$ . **трит**.
- $m = 10$ . **дит**.

# Зависимость количества информации от вероятности



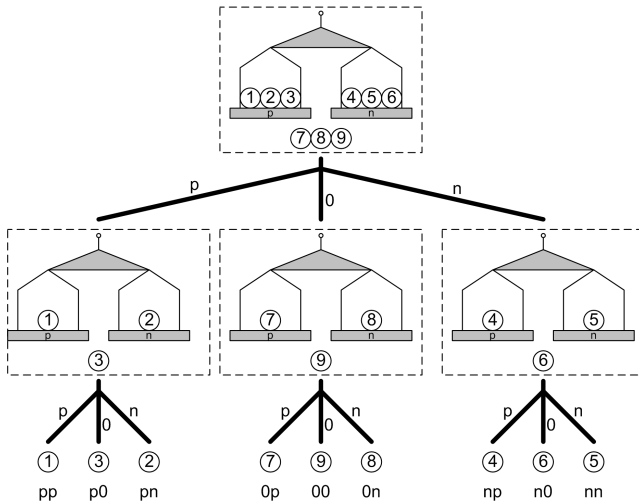
## Задача о бильярдных шарах (постановка)

### Example

Имеется восемь бильярдных шаров с номерами 1-8 соответственно. Все шары одинаковой массы, кроме одного, который тяжелее остальных. Имеются весы Фемиды (чашечные). Какое количество взвешиваний потребуется, чтобы определить номер тяжелого шара?

# Задача о бильярдных шарах (решение)

Решение.  $H = \log_3 9 = I(p) = -\log_3 \frac{1}{9} = 2$  трита





# Энтропия

Мера информативности источника событий (сколько выдаёт  $I$  в среднем за раз)

$$E = - \sum_{i=1}^N p_i \cdot \log_m p_i.$$

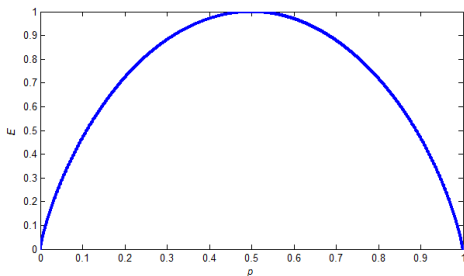


Рис.: Энтропия для источника с двумя состояниями

Война префиксов закончилась 19 марта 2005 года  
Принят стандарт IEEE 1541. 1000 байт — 1 kB (килобайт), 1024 байт — 1KiB (кибибайт)

Множитель	СИ/SI	Множитель	IEEE 1541
$10^3 = 1000^1$	<i>kilo</i> (k) кило	$2^{10} = 1024^1$	<i>kibi</i> (Ki) киби
$10^6 = 1000^2$	<i>mega</i> (M) мега	$2^{20} = 1024^2$	<i>mebi</i> (Mi) меби
$10^9 = 1000^3$	<i>giga</i> (G) гига	$2^{30} = 1024^3$	<i>gibi</i> (Gi) гиби
$10^{12} = 1000^4$	<i>tera</i> (T) тера	$2^{40} = 1024^4$	<i>tebi</i> (Ti) теби
$10^{15} = 1000^5$	<i>peta</i> (P) пета	$2^{50} = 1024^5$	<i>pebi</i> (Pi) пеби
$10^{18} = 1000^6$	<i>exa</i> (E) экса	$2^{60} = 1024^6$	<i>exbi</i> (Ei) эксби
$10^{21} = 1000^7$	<i>zetta</i> (Z) зетта	$2^{70} = 1024^7$	<i>zebi</i> (Zi) зеби
$10^{24} = 1000^8$	<i>yotta</i> (Y) йотта	$2^{80} = 1024^8$	<i>yobi</i> (Yi) йоби

# Кодирование

## Definition

**Кодирование** — процесс перехода от **источника событий** к **источнику информации**. Т.е. сопоставление **событиям** цепочек из **кодowych символов**.

Некоторые назначения кодирования:

- ❶ принципиальная возможность описания мира с помощью символов конечного алфавита;
- ❷ устранение избыточности, сжатие информации, экономия памяти и снижение нагрузки на каналы передачи информации;
- ❸ обеспечение устойчивости к помехам;
- ❹ защита важных свойств информации (конфиденциальность, целостность, принадлежность и т.д.).

# Формальное определение кодирования

## Definition

Дано:

- Алфавит **событий**:  $S = \{s_1, \dots, s_N\}$ ;
- Алфавит кодовых **символов**:  $T = \{t_1, \dots, t_m\}$ ;

Требуется задать отображение  $\delta : S \rightarrow T^+$  (таблицу кодов, **схему кодирования**):

$$\delta = \langle s_1 \mapsto \omega_1, \dots, s_N \mapsto \omega_N \rangle,$$

где  $\omega_i = t_{i_1} \cdots t_{i_{k_i}}$ , причем слово  $s_j = s_{j_1} \cdots s_{j_t}$  будет кодироваться символами кодового алфавита как  $s_j = \omega_{j_1} \cdots \omega_{j_t}$ . Множество кодовых слов  $\omega_i$ , соответствующих  $s_i$  называется множеством **элементарных кодов**:

$$\Omega = \{\omega_1, \dots, \omega_N\}.$$

# Примеры кодирования

## Example (Неоднозначное декодирование. $\delta$ — не биекция)

$S = \{A, B, C, D, E, F, G, H\}$ ;  $T = \{0, 1\}$ ;  $\delta = \langle A \rightarrow 0, B \rightarrow 1, C \rightarrow 10, D \rightarrow 11, E \rightarrow 100, F \rightarrow 101, G \rightarrow 110, H \rightarrow 111 \rangle$ .

- Кодирование однозначно:  $ABAB \mapsto 0101$ .
- Декодирование нет: 0101 разделяется на слова  $ABAB$ ,  $AF$  и  $ACB$ .

## Example (Однозначное декодирование. $\delta$ — биекция)

$S = \{A, B, C, D, E, F, G, H\}$ ;  $T = \{0, 1\}$ ;  $\delta = \langle A \mapsto 000, B \mapsto 001, C \mapsto 010, D \mapsto 011, E \mapsto 100, F \mapsto 101, G \mapsto 110, H \mapsto 111 \rangle$ .

- Кодирование:  $ABBA \mapsto 000111000111$ .
- Декодирование:  $000111000111 \mapsto ABBA$ .

# Схемы кодирования

## Definition

Схема кодирования  $\delta$  является **разделимой**, если любое слово  $\varsigma$ , составленное из элементарных кодов  $\omega_i$  единственным образом разлагается на элементарные коды.

Разделимая схема допускает декодирование. Важным частным случаем **разделимых** схем являются **префиксные** схемы.

## Definition

Схема называется **префиксной**, если ни один элементарный код  $\omega_i$  из множества  $\Omega$  не является префиксом<sup>a</sup> другого кода из того же множества.

---

<sup>a</sup>Префиксом слова  $\omega$  называется слово  $\omega_1$ , если  $\omega = \omega_1\omega_2$

## «Равномерное» кодирование

Наиболее простым вариантом кодирования является **равномерное** кодирование, когда все элементарные коды равной длины.

Для кодирования  $N$  событий требуется использовать цепочки длины

$$l(n) = \lceil \log_m(n) \rceil,$$

где  $m$  — количество кодовых символов;  $\lceil X \rceil$  — наименьшее целое, большее или равное  $X$ .

Эта же оценка на основе постулатов Шеннона:

$$l(n) = \left\lceil -\log_m \left( \frac{1}{n} \right) \right\rceil = \lceil \log_m(n) \rceil.$$

# Равномерное кодирование

## Example

В соревновании участвуют 33 спортсмена. Для регистрации пересечения финишной черты каждому спортсмену выдается радио-брелок. В момент пересечения финишной черты спортсменом, брелок передает двоичный код для идентификации спортсмена. Все брелки передают код одинаковой длины. Какое минимально необходимоме количество бит в общем случае должен передать брелок?



# Равномерное кодирование

## Example

В соревновании участвуют 33 спортсмена. Для регистрации пересечения финишной черты каждому спортсмену выдается радио-брелок. В момент пересечения финишной черты спортсменом, брелок передает двоичный код для идентификации спортсмена. Все брелки передают код одинаковой длины. Какое минимально необходимое количество бит в общем случае должен передать брелок?

## Решение.

$$\lceil \log_2(33) \rceil = 6.$$



# Сигнал

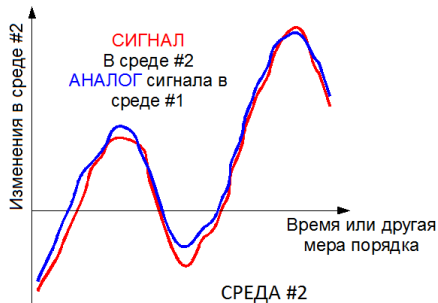
## Definition

**Сигнал** — это изменение (во времени или пространстве) физической величины, несущее информацию, т.е. способ, позволяющий фиксировать **символ** в материально-энергетическом носителе

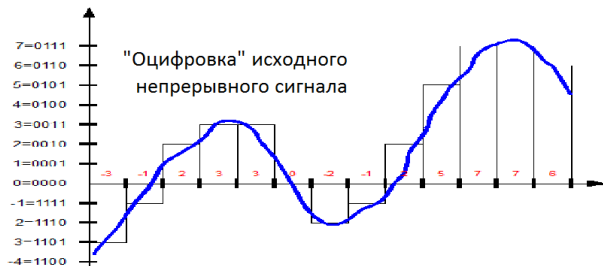
Выделяют два вида сигналов:

- 1 аналоговые;
- 2 цифровые.

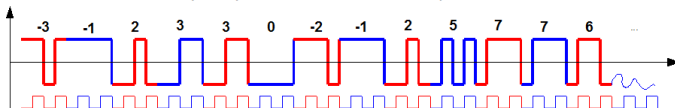
# Аналоговый сигнал



# Цифровой сигнал



Цифровой сигнал (цифры 0 и 1 задаются "высоким" и "низким" значениями некоторой физической величины)



# Информативность источника событий

Источнику событий после кодирования соответствует источник информации, выдающий коды событий. Оценку информативности **источника событий** дает величина, называемая **энтропией** источника событий:

$$E = - \sum_{i=1}^N p_i \cdot \log_m p_i, \quad (2)$$

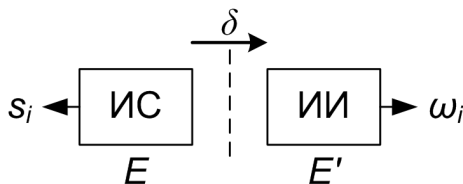
где  $p_i$  — вероятность  $i$ -го события  $s_i \in S$  на выходе источника событий,  $m$  — количество кодовых символов,  $N$  — количество событий  $N = |S|$ .

# Информативность источника информации

Так как вероятности появления кодов событий останутся прежними, то энтропия **источника информации**  $E'$  будет равна

$$E' = \sum_{i=1}^n p_i \cdot l_i, \quad (3)$$

где  $l_i$  — длина кода  $\omega_i$  для  $i$ -го события.



# Задача оптимального кодирования

## Аксиома

Энтропия источника информации всегда больше энтропии отражаемого источника событий.

Задача оптимального кодирования: максимально приблизить энтропию источника информации к энтропии источника событий.

# Оценка оптимальности кодирования

Пусть имеется источник событий  $s_i$ , о вероятности появления которых на его выходе известно следующее:

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1

Энтропия источника событий, формула (2), составляет:

$$\begin{aligned} E &= -(0.5 \cdot \log_2 0.5 + 0.3 \cdot \log_2 0.3 + 0.1 \cdot \log_2 0.1 + 0.1 \cdot \log_2 0.1) \approx \\ &\approx (0.5 + 0.521089678 + 0.332192809 + 0.332192809) \approx \\ &\approx 1.685475297 \text{ бит.} \end{aligned}$$



# Оценка оптимальности кодирования

Для равномерного кодирования битами может быть получен такой вариант:

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$	00	01	10	11

Энтропия данного источника информации составит, по формуле (3):

$$E' = (0.5 \cdot 2 + 0.3 \cdot 2 + 0.1 \cdot 2 + 0.1 \cdot 2) = 2 \text{ бита.}$$

## Оценка оптимальности кодирования

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$	0	10	110	111

Так же как и предыдущая, эта схема префиксная и разделимая, но неравномерная. Энтропия источника информации теперь составляет

$$E' = (0.5 \cdot 1 + 0.3 \cdot 2 + 0.1 \cdot 3 + 0.1 \cdot 3) = 1.7 \text{ бита.}$$

# Оценка оптимальности кодирования

Представленные источники **эквивалентны**. Если запустить источник информации на выдачу, например, 100 кодов событий, то первый вариант кодирования выдаст цепочку длины  $\approx 200$ , а второй —  $\approx 170$  бит.

# Алгоритм Хаффмана для $m = 2$

- ❶ События сортируются по убыванию вероятности.
- ❷ Два события с минимальными вероятностями объединяются в одно составное событие с суммарной вероятностью исходных. При этом одно из исходных событий помечается кодовым символом 0, а второе — символом 1. Исходные события исключаются из множества событий, вместо них остается одно составное.
- ❸ Шаги 1 и 2 последовательно повторяются до тех пор, пока все события не склеятся в единственное составное событие, вероятность которого равна 1. После этого кодовое слово  $\omega_i$  для исходного события  $s_i$  есть цепочка из кодовых символов, которыми помечены все составные события от корня до  $s_i$ .

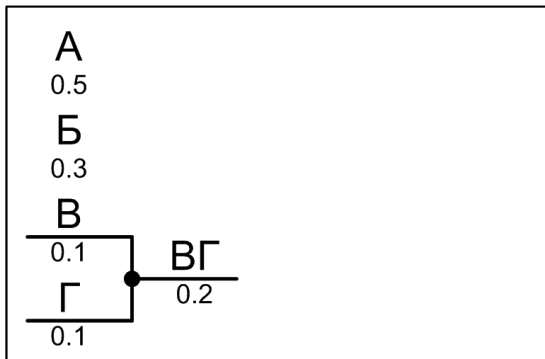
# Оптимальное кодирование по Хаффману

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				

А
0.5
Б
0.3
В
0.1
Г
0.1

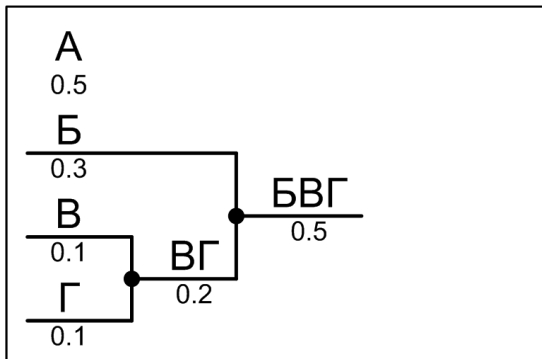
# Оптимальное кодирование по Хаффману

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				



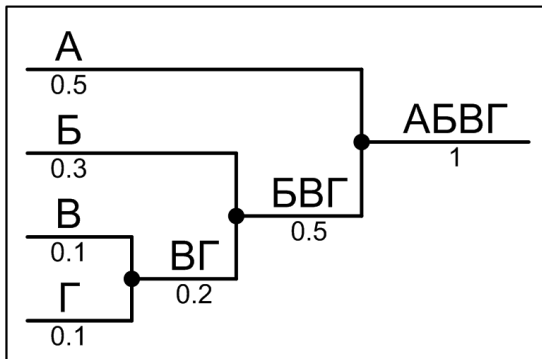
# Оптимальное кодирование по Хаффману

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				



# Оптимальное кодирование по Хаффману

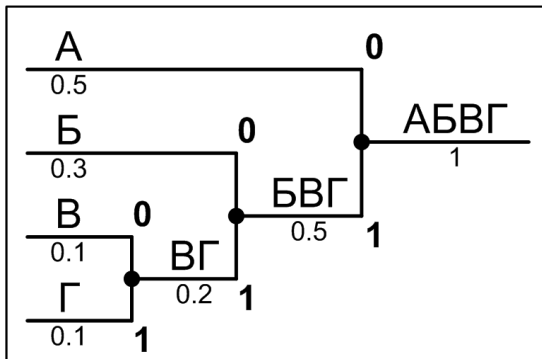
Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				





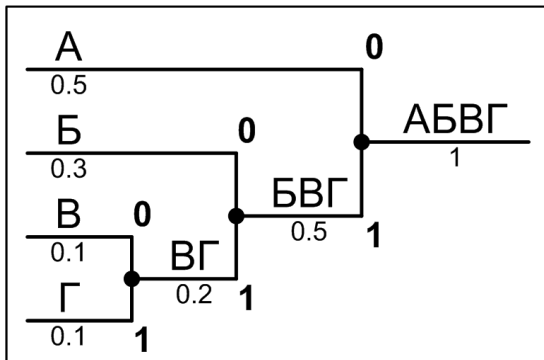
# Оптимальное кодирование по Хаффману

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				



# Оптимальное кодирование по Хаффману

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$	0	10	110	111

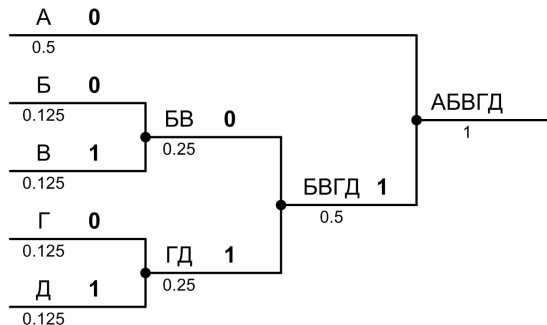


# Оптимальное кодирование по Хаффману (задача)

Событие $s_i$	А	Б	В	Г	Д
Вероятность $p_i$ события $s_i$	0.5	0.125	0.125	0.125	0.125
Код события $\omega_i$					

# Оптимальное кодирование по Хаффману (задача)

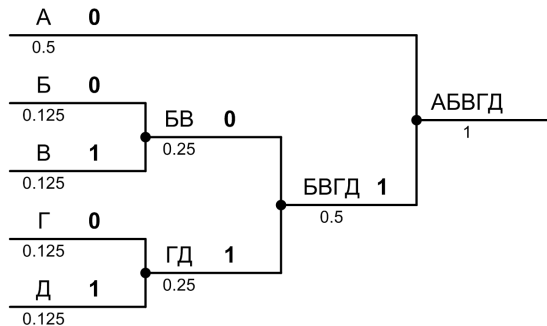
Событие $s_i$	А	Б	В	Г	Д
Вероятность $p_i$ события $s_i$	0.5	0.125	0.125	0.125	0.125
Код события $\omega_i$	0	100	101	110	111



100011011101111 →

# Оптимальное кодирование по Хаффману (задача)

Событие $s_i$	А	Б	В	Г	Д
Вероятность $p_i$ события $s_i$	0.5	0.125	0.125	0.125	0.125
Код события $\omega_i$	0	100	101	110	111



10001101110111  $\rightarrow$  БАГДАД

# Алгоритм Фано для $m = 2$

- 1 Исходный массив событий, сортируется в порядке убывания вероятностей.
- 2 Массив разбивается на две части, так, чтобы разница сумм вероятностей событий каждой части была минимальна. Первый кодовый символ элементарного кода  $\omega$ ; находится так: для всех событий левой части разбитого массива кодовый символ будет 0, а для всех событий правой части — 1.
- 3 Второй и последующие кодовые символы определяется так: каждая часть разбитого исходного массива, в которой более одного события, становится исходным массивом, и её разбиение выполняется так же, как исходного массива (шаг 2).

# Оптимальное кодирование по алгоритму Фано

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				

$s_i$	$p_i$
А	0.5
Б	0.3
В	0.1
Г	0.1

# Оптимальное кодирование по алгоритму Фано

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				

$s_i$	$p_i$	$\omega_i$
А	0.5	0
Б	0.3	1
В	0.1	1
Г	0.1	1



# Оптимальное кодирование по алгоритму Фано

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				

$s_i$	$p_i$	$\omega_i$	
А	0.5	0	
Б	0.3	1	0
В	0.1	1	1
Г	0.1	1	1

# Оптимальное кодирование по алгоритму Фано

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$				

$s_i$	$p_i$	$\omega_i$		
А	0.5	0		
Б	0.3	1	0	
В	0.1	1	1	0
Г	0.1	1	1	1

# Оптимальное кодирование по алгоритму Фано

Событие $s_i$	А	Б	В	Г
Вероятность $p_i$ события $s_i$	0.5	0.3	0.1	0.1
Код события $\omega_i$	0	10	110	111

$s_i$	$p_i$	$\omega_i$		
А	0.5	0		
Б	0.3	1	0	
В	0.1	1	1	0
Г	0.1	1	1	1

# Оптимальное кодирование по алгоритму Фано (задача)

$s_i$	А	Б	В	Г	Д	Е	Ж
$p_i$	0.135	0.24	0.25	0.125	0.0635	0.124	0.0625
$\omega_i$							

## Оптимальное кодирование по алгоритму Фано (задача)

$s_i$	А	Б	В	Г	Д	Е	Ж
$p_i$	0.135	0.24	0.25	0.125	0.0635	0.124	0.0625
$\omega_i$							

$s_i$	$p_i$	$\omega_i$			
В	0.25	0	0		
Б	0.24	0	1		
А	0.135	1	0	0	
Г	0.125	1	0	1	
Е	0.124	1	1	0	
Д	0.0635	1	1	1	0
Ж	0.0625	1	1	1	1

# Оптимальное кодирование по алгоритму Фано (задача)

$s_i$	А	Б	В	Г	Д	Е	Ж
$p_i$	0.135	0.24	0.25	0.125	0.0635	0.124	0.0625
$\omega_i$	100	01	00	101	1110	110	1111

$s_i$	$p_i$	$\omega_i$			
В	0.25	0	0		
Б	0.24	0	1		
А	0.135	1	0	0	
Г	0.125	1	0	1	
Е	0.124	1	1	0	
Д	0.0635	1	1	1	0
Ж	0.0625	1	1	1	1

# Сжатие

Кодирование с целью сжатия (или просто **сжатие**) ставит себе в задачу уменьшить количество информации, не теряя (или оставаясь в допустимых рамках) при этом свойство адекватности отражаемому объекту. В случае сжатия, события  $s_i$  уже представляют собой слова в алфавите  $T$  (коды). При сжатии информация **перекодируется** в том же алфавите  $T$ .

Классы алгоритмов сжатия:

- сжатие с потерями (адекватности);
- сжатие без потерь.

# Алгоритм Лемпела-Зива (LZ)

## Кодирование

- 1 В словарь нулевым элементом помещается пустая цепочка  $\varepsilon$ .  
Пустое слово  $\varepsilon$  не содержит букв, и для любого слова  $\omega$  справедливо  $\omega = \varepsilon\omega = \omega\varepsilon$ .
- 2 От исходной цепочки  $t$  отделяется слово  $\omega a$ , где  $\omega$  — максимально длинное слово из словаря,  $a$  — расширяющая буква. Т.е.  $t = \omega a t'$ .
- 3 В конец словаря добавляется новое слово  $\omega a$ . К коду  $c$  добавляется пара  $\langle i_\omega, a \rangle$ , где  $i_\omega$  — индекс слова  $\omega$  в словаре. От исходного текста отделяется слово  $\omega a$ :  $t = t'$ .
- 4 Пункты 2-3 последовательно повторяются до тех пор, пока в тексте  $t$  остается хоть одна буква.

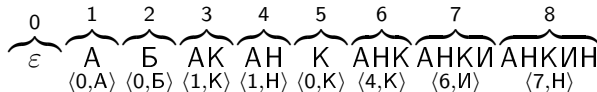
В результате получается код  $c = \langle i_1, a_1 \rangle \cdots \langle i_n, a_n \rangle$ .



## Алгоритм Лемпела-Зива

Пример сжатия текста: «АБАКАНКАНКАНКИАНКИН»

$i$	$t$	$\omega a$	$c = \langle i_\omega, a \rangle$
		$0 \rightarrow \varepsilon$	
1	$\varepsilon$ АБАКАНКАНКАНКИАНКИН	$1 \rightarrow A$	$\langle 0, A \rangle$
2	$\varepsilon$ БАКАНКАНКАНКИАНКИН	$2 \rightarrow Б$	$\langle 0, Б \rangle$
3	АКАНКАНКАНКИАНКИН	$3 \rightarrow АК$	$\langle 1, К \rangle$
4	АНКАНКАНКИАНКИН	$4 \rightarrow АН$	$\langle 1, Н \rangle$
5	$\varepsilon$ КАНКАНКИАНКИН	$5 \rightarrow К$	$\langle 0, К \rangle$
6	АНКАНКИАНКИН	$6 \rightarrow АНК$	$\langle 4, К \rangle$
7	АНКИАНКИН	$7 \rightarrow АНКИ$	$\langle 6, И \rangle$
8	АНКИН	$8 \rightarrow АНКИН$	$\langle 7, Н \rangle$



Дать оценку длин кода и текста

# Алгоритм Лемпела-Зива

## Задание

Сжать текст:

«тартарарамитамтамывтартарарах»

# Алгоритм Лемпела-Зива

## Задание

Сжать текст:

«тартарарамитамтамывтартарарах»

0	1	2	3	4	5	6	7	8	9	10	11	12	13
$\varepsilon$	т	а	р	та	ра	рам	и	там	тамы	в	тар	тара	рах
	$\langle 0, \text{т} \rangle$	$\langle 0, \text{а} \rangle$	$\langle 0, \text{р} \rangle$	$\langle 1, \text{а} \rangle$	$\langle 3, \text{а} \rangle$	$\langle 5, \text{м} \rangle$	$\langle 0, \text{и} \rangle$	$\langle 4, \text{м} \rangle$	$\langle 8, \text{ы} \rangle$	$\langle 0, \text{в} \rangle$	$\langle 4, \text{р} \rangle$	$\langle 11, \text{а} \rangle$	$\langle 5, \text{х} \rangle$

Код:

$\langle 0, \text{т} \rangle, \langle 0, \text{а} \rangle, \langle 0, \text{р} \rangle, \langle 1, \text{а} \rangle, \langle 3, \text{а} \rangle, \langle 5, \text{м} \rangle, \langle 0, \text{и} \rangle,$   
 $\langle 4, \text{м} \rangle, \langle 8, \text{ы} \rangle, \langle 0, \text{в} \rangle, \langle 4, \text{р} \rangle, \langle 11, \text{а} \rangle, \langle 5, \text{х} \rangle$

# Алгоритм Лемпела-Зива (LZ)

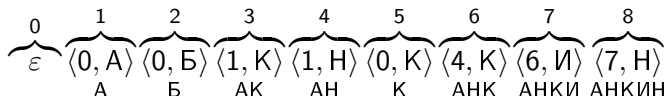
## Декодирование

- 1 В словарь нулевым элементом помещается пустая цепочка  $\varepsilon$ .  
Текст  $t$  не содержит букв:  $t = \varepsilon$ .
- 2 От исходного кода  $s$  отделяется пара  $\langle i, a \rangle$ , в словарь добавляется слово  $w_i a$ , где  $w_i$  —  $i$ -е слово из словаря. Восстанавливается текст  $t = t w_i a$ .
- 3 Пункт 2 последовательно повторяется до тех пор, пока в коде  $s$  остается хоть одна пара.

## Алгоритм Лемпела-Зива

Пример декодирования «0A,0B,1K,1H,0K,4K,6I,7H»

$i$	$c = \langle i_w, a \rangle$	$\omega a$	$t$
		$0 \rightarrow \varepsilon$	
1	$\langle 0, A \rangle$	$1 \rightarrow A$	$\varepsilon A$
2	$\langle 0, B \rangle$	$2 \rightarrow B$	$A \varepsilon B$
3	$\langle 1, K \rangle$	$3 \rightarrow AK$	$ABAK$
4	$\langle 1, H \rangle$	$4 \rightarrow AH$	$ABAKAH$
5	$\langle 0, K \rangle$	$5 \rightarrow K$	$ABAKAH \varepsilon K$
6	$\langle 4, K \rangle$	$6 \rightarrow ANK$	$ABAKANKAN$
7	$\langle 6, I \rangle$	$7 \rightarrow ANKI$	$ABAKANKANANKI$
8	$\langle 7, H \rangle$	$8 \rightarrow ANKIH$	$ABAKANKANANKIH$



# Алгоритм Лемпела-Зива

## Задание

Восстановить текст из кода:

$\langle 0, \text{в} \rangle, \langle 0, \text{о} \rangle, \langle 0, \text{т} \rangle, \langle 1, \text{а} \rangle, \langle 0, \text{м} \rangle, \langle 4, \text{р} \rangle, \langle 6, \text{ы} \rangle, \langle 2, \text{т} \rangle, \langle 6, \text{ы} \rangle$

# Алгоритм Лемпела-Зива

## Задание

Восстановить текст из кода:

$\langle 0, \text{в} \rangle, \langle 0, \text{о} \rangle, \langle 0, \text{т} \rangle, \langle 1, \text{а} \rangle, \langle 0, \text{м} \rangle, \langle 4, \text{р} \rangle, \langle 6, \text{ы} \rangle, \langle 2, \text{т} \rangle, \langle 6, \text{ы} \rangle$

0	1	2	3	4	5	6	7	8	9
$\underbrace{\quad}_{\varepsilon}$	$\underbrace{\langle 0, \text{в} \rangle}$	$\underbrace{\langle 0, \text{о} \rangle}$	$\underbrace{\langle 0, \text{т} \rangle}$	$\underbrace{\langle 1, \text{а} \rangle}$	$\underbrace{\langle 0, \text{м} \rangle}$	$\underbrace{\langle 4, \text{р} \rangle}$	$\underbrace{\langle 6, \text{ы} \rangle}$	$\underbrace{\langle 2, \text{т} \rangle}$	$\underbrace{\langle 6, \text{ы} \rangle}$
	в	о	т	ва	м	вар	вары	от	вары

Текст:

«вотвамварварыотвары»

# Свойства информации с точки зрения её защиты

- целостность:
- конфиденциальность:
- принадлежность:
- доступность:



# Свойства информации с точки зрения её защиты

- целостность: контрольные суммы; корректирующие коды;
- конфиденциальность: шифрование; скрытая передача;
- принадлежность: цифровая подпись;
- доступность: надежность информационных систем.

# Классификация ошибок

Ошибки, возникающие в цифровом (двоичном) канале могут быть следующими:

- замещения кодового символа;
- вставка кодового символа;
- выпадение кодового символа.

Далее рассматриваются только ошибки замещения. Существуют две стратегии защиты от ошибок замещения:

- с обнаружением и запросом на повторную передачу (ARQ — Automatic Repeat Request);
- с обнаружением и непосредственным исправлением на стороне получателя (FEC — Forward Error Correction).

# ARQ

Примером стратегии ARQ может считаться контроль по четности (нечетности).

$$p_{\text{чётн}} = d_{n-1} \oplus \dots \oplus d_1 \oplus d_0,$$

$$p_{\text{нечётн}} = d_{n-1} \oplus \dots \oplus d_1 \oplus d_0 \oplus 1.$$

# FEC

Можно кодировать каждый бит исходной последовательности по схеме

$$\delta = \{0 \mapsto 000, 1 \mapsto 111\},$$

а декодировать по схеме

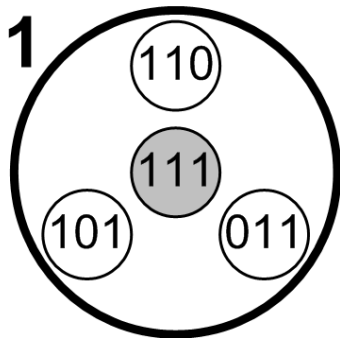
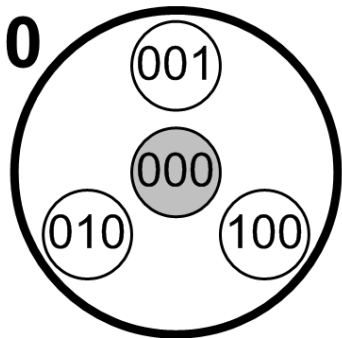
$$\delta' = \{000 \mapsto 0, 001 \mapsto 0, 010 \mapsto 0, 100 \mapsto 0, \\ 111 \mapsto 1, 110 \mapsto 1, 101 \mapsto 1, 011 \mapsto 1\},$$

## Example

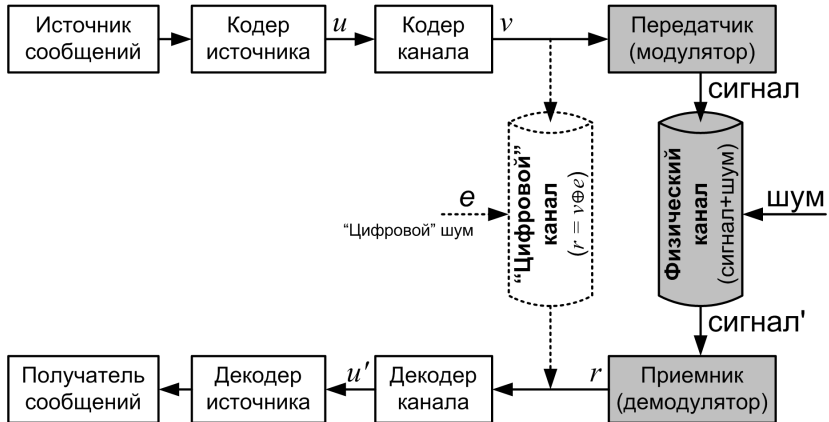
Пусть передается слово 101. Кодировается: 111000111. Поступает в канал. Возникает одиночная ошибка: 11 $\boxed{0}$ 000111. Декодируется: 101. При этом декодер обнаруживает и исправляет одиночную ошибку.  $\square$

# FEC

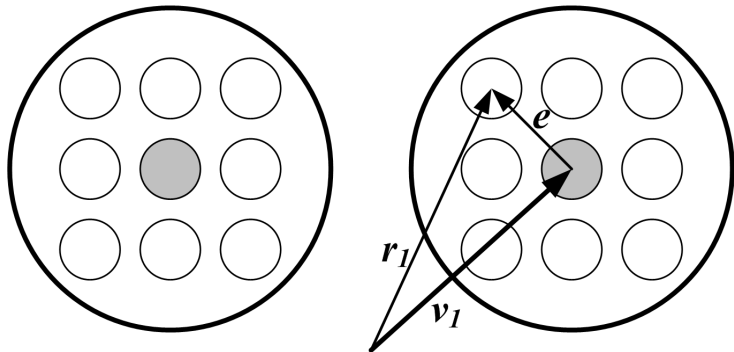
## Пример кодирования «утроением»



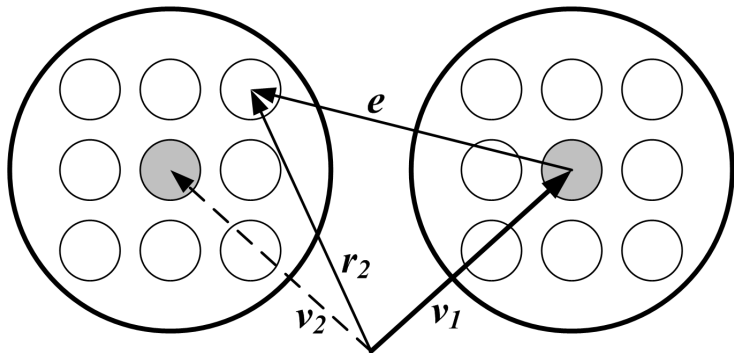
# Схема канала передачи данных



# FEC — Ошибка обнаружена и верно исправлена

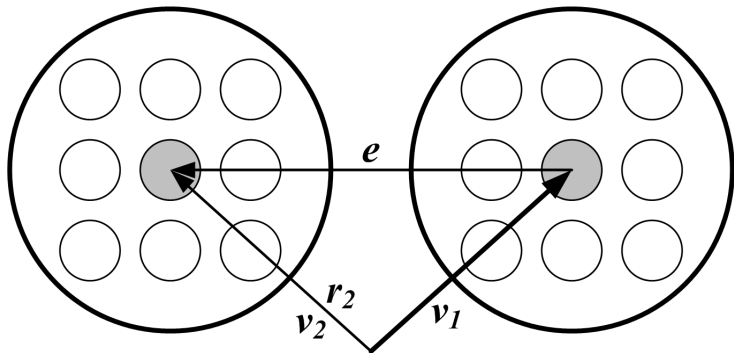


# FEC — ошибка обнаружена, но исправлена неверно





# FEC — необнаружимая ошибка



# FEC: Код Хемминга

- Код Хемминга формирует номер ошибочного разряда.
- Признаком отсутствия ошибок является нулевой номер.
- Поэтому вводится «фиктивный» нулевой разряд.
- Если исходное слово имеет длину  $n$  бит, тогда к нему нужно добавить  $m$  дополнительных бит, исходя из неравенства

$$2^m \geq n + m + 1, \quad (4)$$

где левая часть неравенства — количество  $m$ -разрядных двоичных чисел, правая — общая длина кода с учетом «фиктивного» разряда. Выбирается минимальное  $m$  из возможных.

# Кодирование и декодирование по Хеммингу

## Алгоритм кодирования

- 1 В двоичном числе длиной  $m + n$  бит (без фиктивного разряда) контрольные  $m$  бит размещаются в разрядах с номерами, равными степени двойки ( $2^i, 0 \leq i < m$ ). А  $n$  бит исходного слова размещаются в оставшихся разрядах. Контрольные биты при этом инициализируются нулевыми значениями.
- 2 Каждый контрольный бит  $c_{2^i}$  в разряде  $2^i$  пересчитывается как сумма по XOR бит кода, находящихся в разрядах с номерами, двоичное представление которых содержит единицу в  $i$  разряде (включая и сам контрольный разряд).

При декодировании контрольные разряды пересчитываются в соответствии с пунктом 2 алгоритма кодирования. В результате в контрольных разрядах будет получено двоичное представление номера разряда ошибочного бита.

# Построить код Хемминга для слова $u = 0011$

$n = 4$ . Исходя из формулы (4) выбирается  $m = 3$ .

$$u = \begin{bmatrix} 0 & 0 & 1 & 1 \end{bmatrix}$$

$$v = \begin{bmatrix} 0 & 0 & 1 & 0 & 1 & 0 & 0 & \end{bmatrix}$$

7	6	5	4	3	2	1	0	
1	0	1	0	1	0	1	0	$v_1$
1	1	0	0	1	1	0	0	$v_2$
1	1	1	1	0	0	0	0	$v_4$

$$v_1 = v_1 \oplus v_3 \oplus v_5 \oplus v_7 = 0 \oplus 1 \oplus 1 \oplus 0 = 0$$

$$v_2 = v_2 \oplus v_3 \oplus v_6 \oplus v_7 = 0 \oplus 1 \oplus 0 \oplus 0 = 1$$

$$v_4 = v_4 \oplus v_5 \oplus v_6 \oplus v_7 = 0 \oplus 1 \oplus 0 \oplus 0 = 1$$

$$\begin{bmatrix} 1 & 1 & 0 \end{bmatrix}$$

$$v = \begin{bmatrix} 0 & 0 & 1 & 1 & 1 & 1 & 0 & \end{bmatrix}$$

# Код Хемминга

## Обнаружение и исправление одиночных ошибок

$v =$ 

0	0	1	1	1	1	0	
---	---	---	---	---	---	---	--

$r =$ 

0	0	0*	1	1	1	0	
---	---	----	---	---	---	---	--

7	6	5	4	3	2	1	0	
1	0	1	0	1	0	1	0	$r_1$
1	1	0	0	1	1	0	0	$r_2$
1	1	1	1	0	0	0	0	$r_4$

$$r_1 = r_1 \oplus r_3 \oplus r_5 \oplus r_7 = 0 \oplus 1 \oplus 0 \oplus 0 = 1$$

$$r_2 = r_2 \oplus r_3 \oplus r_6 \oplus r_7 = 1 \oplus 1 \oplus 0 \oplus 0 = 0$$

$$r_4 = r_4 \oplus r_5 \oplus r_6 \oplus r_7 = 1 \oplus 0 \oplus 0 \oplus 0 = 1$$

1	0	1
---	---	---

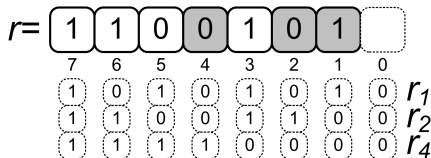
$u' =$ 

0	0	1	1
---	---	---	---

# Код Хемминга

## Задание

Получена последовательность бит  $r$ . Перед передачей из исходного 4-х битного слова был получен код Хемминга. Выяснить, были ли ошибки в процессе передачи. Если были, то выполнить коррекцию, предполагая, что ошибка одиночная. Выделить исходное 4-х битное слово.



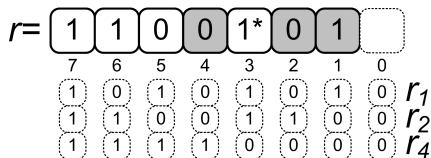
$u' =$

--	--	--	--

# Код Хемминга

## Задание

Получена последовательность бит  $r$ . Перед передачей из исходного 4-х битного слова был получен код Хемминга. Выяснить, были ли ошибки в процессе передачи. Если были, то выполнить коррекцию, предполагая, что ошибка одиночная. Выделить исходное 4-х битное слово.



$$r_1 = r_1 \oplus r_3 \oplus r_5 \oplus r_7 = 1 \oplus 1 \oplus 0 \oplus 1 = 1$$

$$r_2 = r_2 \oplus r_3 \oplus r_6 \oplus r_7 = 0 \oplus 1 \oplus 1 \oplus 1 = 1$$

$$r_4 = r_4 \oplus r_5 \oplus r_6 \oplus r_7 = 0 \oplus 0 \oplus 1 \oplus 1 = 0$$

0	1	1
---	---	---

$u' =$

1	1	0	0
---	---	---	---

# Базовая схема передачи информации



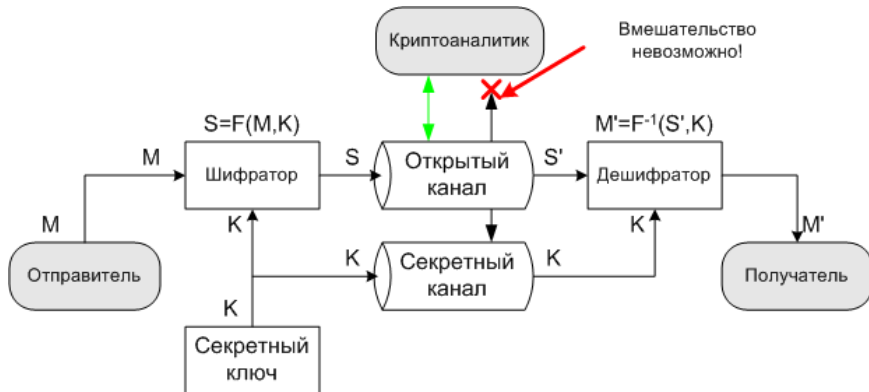
Можно выделить следующие виды **каналов связи**:

- **Секретный** гарантирует конфиденциальность, целостность и принадлежность  $M$ ;
- **Аутентичный** гарантирует только целостность и принадлежность  $M$ ;
- **Открытый** не гарантирует ничего в отношении  $M$ .



# Схемы шифрования

## Симметричная схема

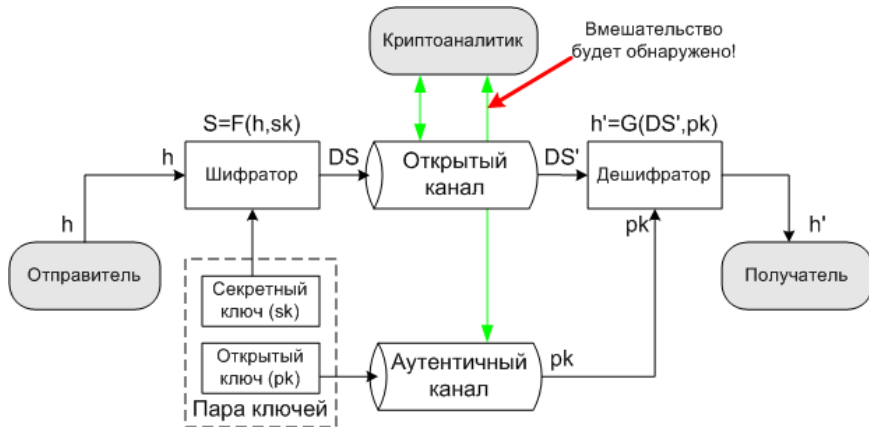


# Схемы шифрования

## Ассиметричная схема



# Цифровая подпись



## В заключение

Изложение математических основ кодирования можно найти, например, в [1, 2]. По основам теории информации можно рекомендовать книгу [3]. Основы кодирования подробно изложены в [4]. Заинтересовавшимся алгоритмами сжатия можно рекомендовать книгу [5].

Как введение по вопросам безопасности информации можно рекомендовать [6]. Обзор задач и протоколов информационной безопасности прекрасно описан в [7]. Математические основы шифрования для сильно интересующихся можно найти в [7, 8, 9].

# Библиография I



*Ф.А.Новиков.* Дискретная математика для программистов /  
Ф.А.Новиков. —  
СПб.: Питер, 2000. —  
304 с.



*С.В.Яблонский.* Введение в дискретную математику: учебное  
пособие для вузов / С.В.Яблонский; Под ред.  
В.А.Садовниченко. —  
М.: Высшая школа, 2001. —  
384 с.

## Библиография II



*В.В.Панин. Основы теории информации: учебное пособие для вузов / В.В.Панин. — 3 изд. — М.: БИНОМ, 2009. — 438 с.*



*М.Вернер. Основы кодирования: учебник для ВУЗов / М.Вернер. — М.: Техносфера, 2004. — 288 с.*



*Д.Сэлмон. Сжатие данных, изображений и звука / Д.Сэлмон. — М.: Техносфера, 2004. — 368 с.*

## Библиография III



*Э.Танненбаум.* Современные операционные системы /  
*Э.Танненбаум.* —  
3 изд. —  
СПб.: Питер, 2010. —  
1120 с.




*Б.Шнайер.* Прикладная криптография. Протоколы, алгоритмы,  
исходные тексты на языке Си / *Б.Шнайер.* —  
М.: Триумф, 2002. —  
816 с.



*Н.Смарт.* Криптография / *Н.Смарт.* —  
М.: Техносфера, 2006. —  
528 с.

## Библиография IV

-  *В.Мао. Современная криптография: теория и практика / В.Мао. — М.: Издательский дом «Вильямс», 2005. — 786 с.*