

Контекстное моделирование

- Й. Риссанен
- Г. Лэнгдон
- 1981г.

- Процесс сжатия сообщения реализуется в два этапа
 - моделирование источника
 - кодирование (в соответствии с полученной моделью)

- Практически всегда истинная статистика источника скрыта,
- поэтому для кодирования очередного символа сообщения необходимо строить вероятностное распределение или модель источника

- Чем точнее построенная модель (или предсказание очередного символа), тем меньше избыточность построенного на основе этой модели кода.

- Техника контекстного моделирования
Prediction by Partial Matching (PPM)

Клири, Уиттен 1984 г.

- Контекстное моделирование –
оценивание вероятности появления
текущего символа в зависимости от
контекста.

Контекст – последовательность
предыдущих символов сообщения

- Контекстная модель порядка m
КМ(m) – это набор оценок вероятностей символов, которые строятся на основании обычных счетчиков частот, связанных с текущим контекстом длины m .

КЛОКОКЛ

	К	О	Л
КМ(-1)	1	1	1
КМ(0)	3	2	2
КМ(1) контекст «К»	-	1	2
КМ(2) контекст «ОК»	-	1	1

- Оценка вероятности появления символа X в контексте S

$$P(X | S) = \frac{N(X | S)}{N(S)}$$

где $N(X|S)$ – количество появлений символа X после контекста S

$N(S)$ – общее количество подстрок S

КЛОКОКЛ

	К	О	Л
КМ(-1)	1/3	1/3	1/3
КМ(0)	3/7	2/7	2/7
КМ(1) контекст «К»	-	1/3	2/3
КМ(2) контекст «ОК»	-	1/2	1/2

- В случае если контекст еще не встречался в сообщении или очередной символ не появлялся в данном контексте, то нельзя оценить вероятность
- Для устранения такой ситуации в алфавит вводят дополнительный специальный символ.

- Символ ухода – виртуальный символ, который показывает необходимость перехода к контекстной модели более низкого порядка.
- Для кодирования символа ухода также необходимо оценить вероятность его использования.

- Вероятность ухода – это суммарная вероятность всех символов алфавита входного потока, еще ни разу не появившихся в контексте.
- Любая КМ должна давать отличную от нуля оценку вероятности ухода. Кроме случаев, когда все оценки вероятностей символов отличны от нуля.

- Оценка вероятности появления символа X в контексте S с учетом символа ухода

$$P(X | S) = \frac{N(X | S)}{N(S) + 1}$$

где $N(X|S)$ – количество появлений символа X после контекста S

$N(S)$ – общее количество подстрок S

Схема метода RRM

- 1. Если вероятность очередного символа X оценивается $KM(N)$ ненулевым числом, то код строится с использованием этой оценки.
- 2. Иначе выдается сигнал в виде символа ухода и происходит переход к шагу 1 с $KM(N-1)$ до тех пор пока символ X не будет оценен ненулевым числом.

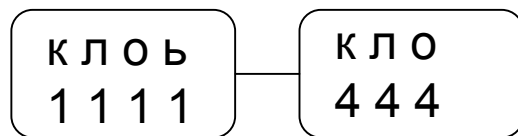
- КМ(-1) гарантирует, что вероятность очередного символа X будет оценена.
- Если в процессе оценки обнаруживается, что текущий рассматриваемый контекст встречается первый раз, то для него создается КМ.
- После кодирования символа X происходит обновление статистики всех КМ, которые использовались при оценке.

КЛОКОКЛОЛКОЛ?

К	Л	О	Ь
1	1	1	1

КМ(-1)

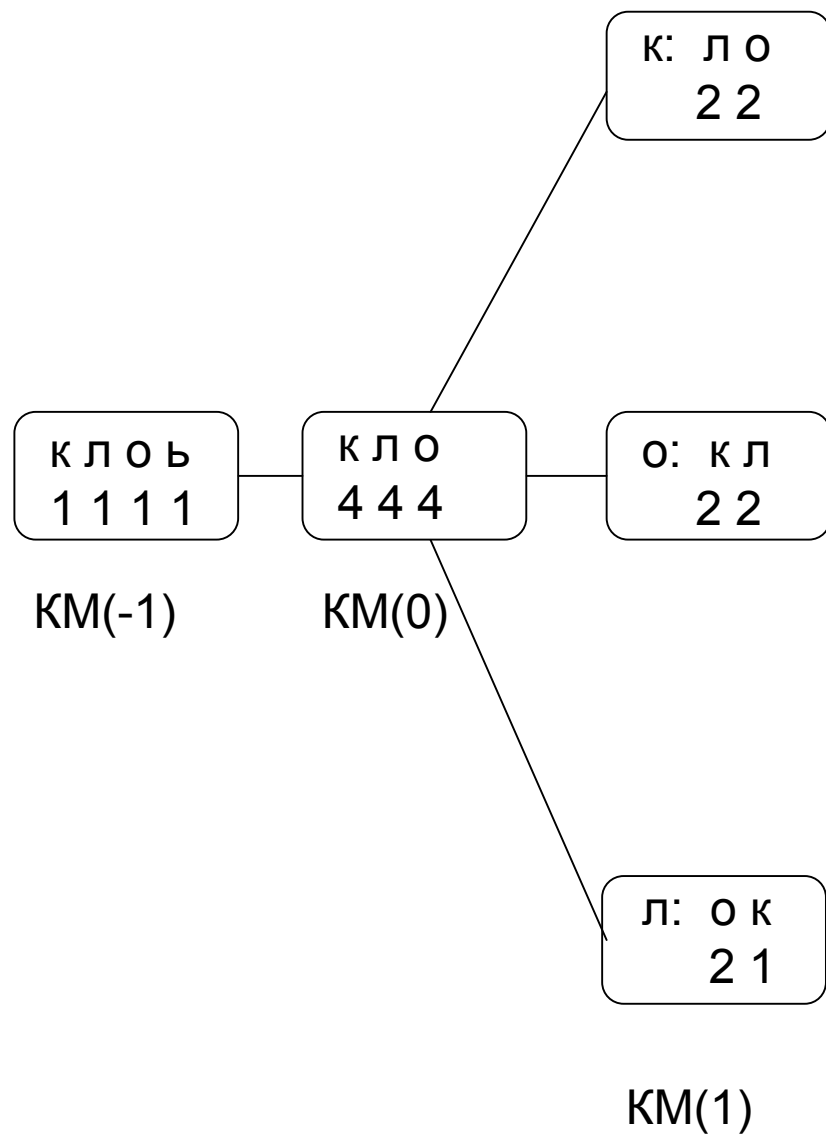
КЛОКОКЛОЛКОЛ?



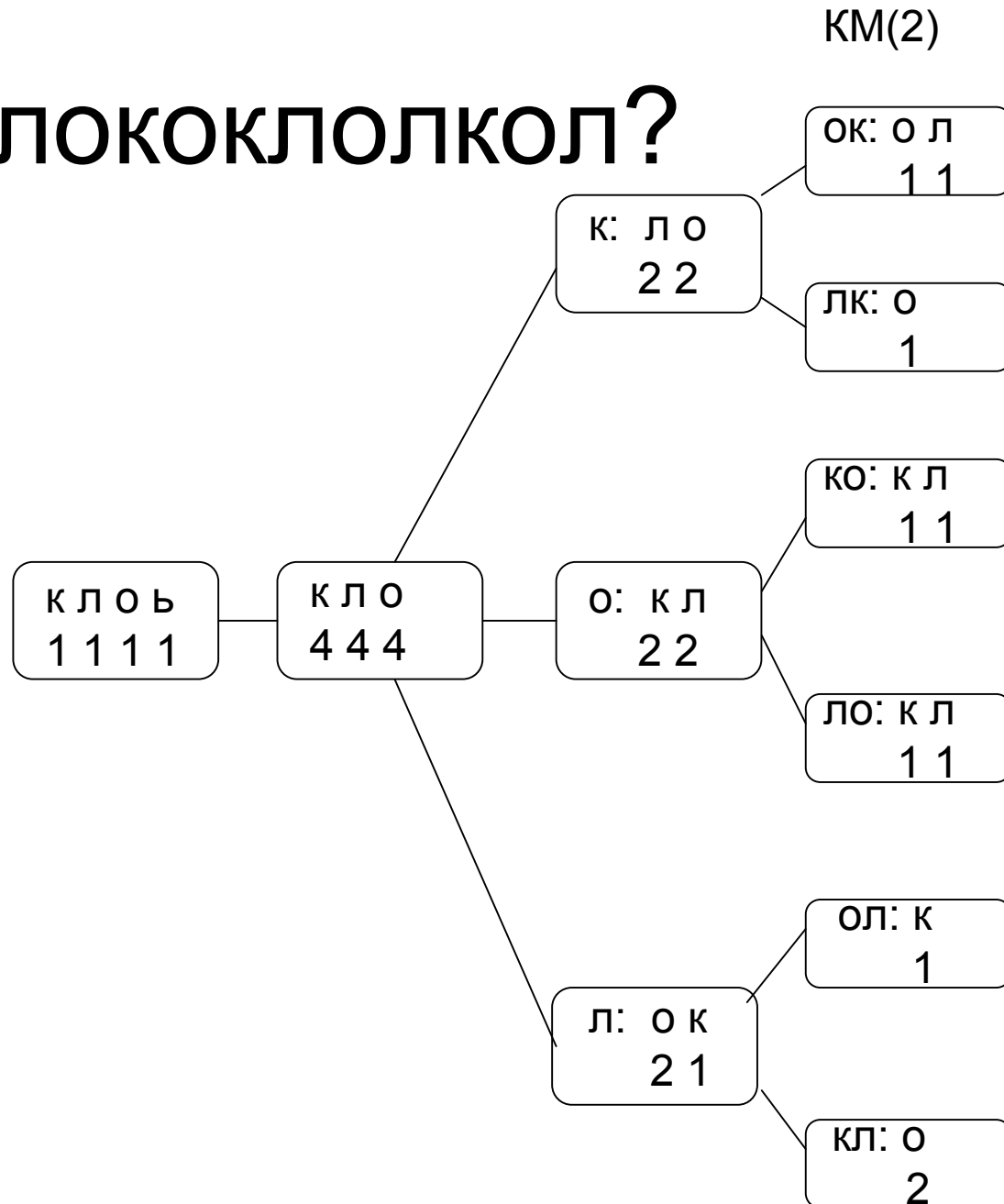
КМ(-1)

КМ(0)

КЛОКОКЛОЛКОЛ?



КЛОКОКЛОЛКОЛ?



КМ(3)

КЛОКОКЛОЛКОЛ?

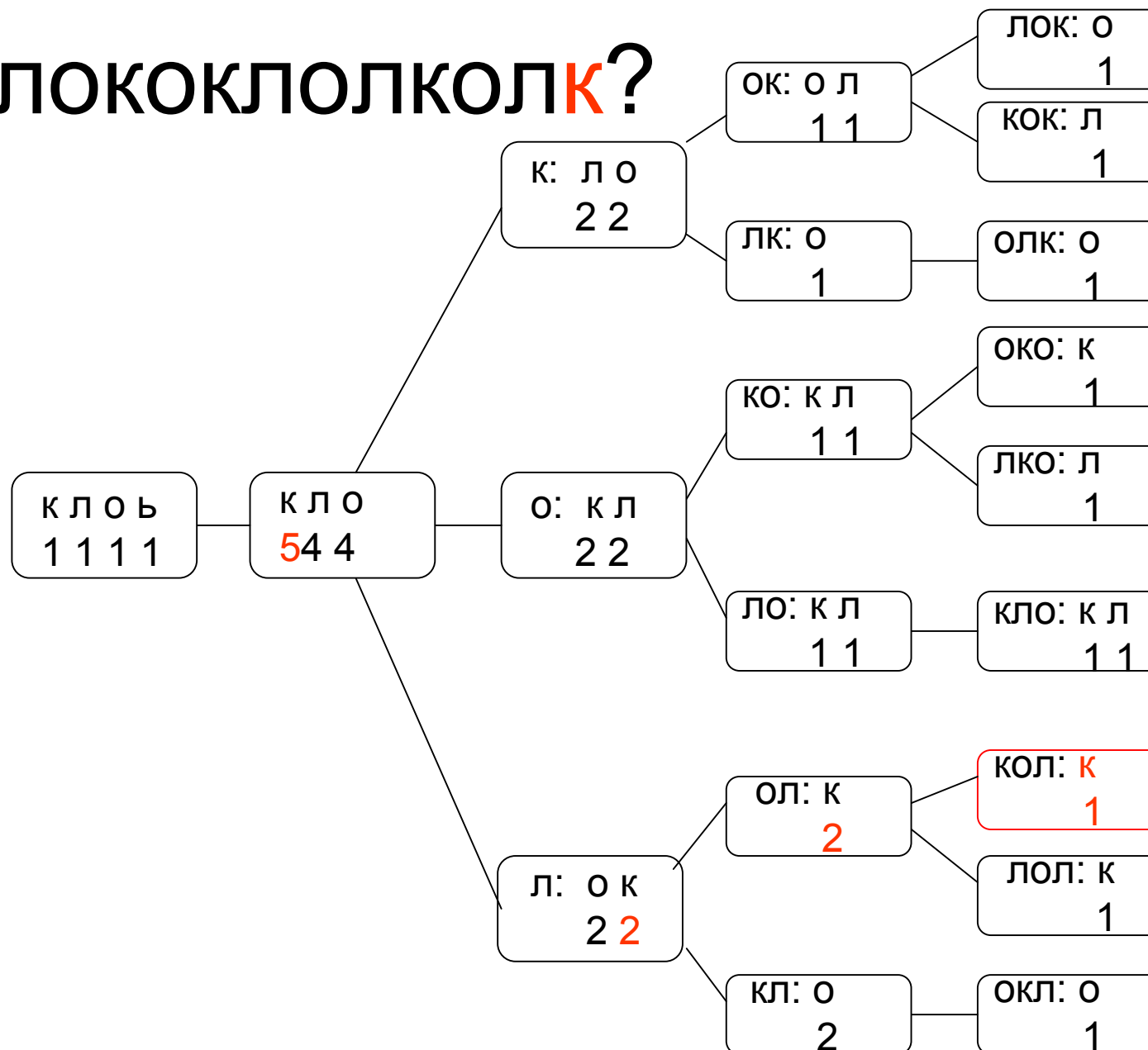


Оценки вероятностей

	3	2	1	0	-1	оценка
	«КОЛ»	«ОЛ»	«Л»	«»		
К	-	$1/(1+1)$	-	-	-	$1/2$
Л	-	$1/(1+1)$	$1/(2+1)$	$4/(4+1)$	-	$2/15$
О	-	$1/(1+1)$	$2/(2+1)$	-	-	$1/3$
Ь	-	$1/(1+1)$	$1/(2+1)$	$1/(4+1)$	1	$1/30$

КМ(3)

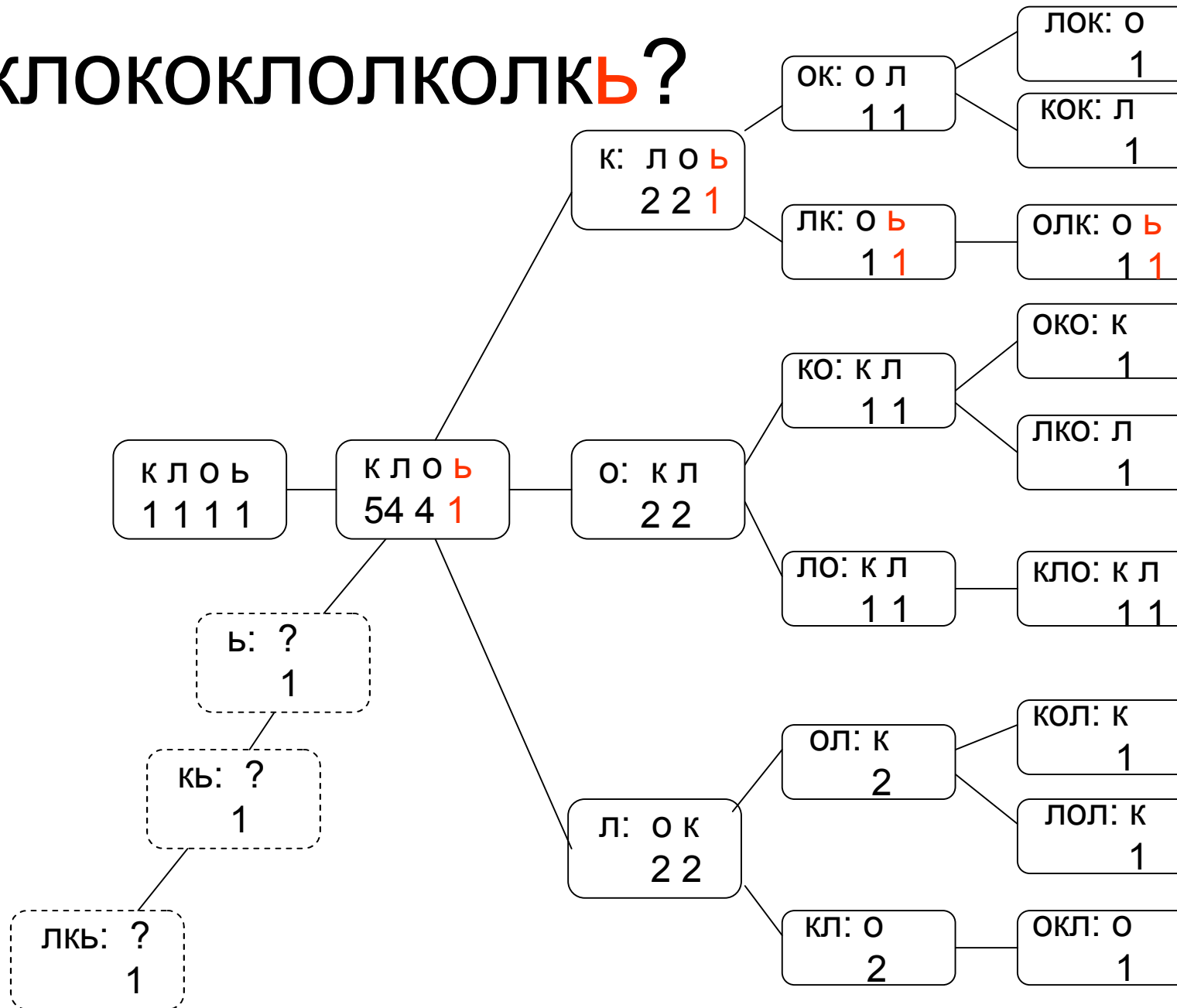
КЛОКОКЛОЛКОЛК?



Оценки вероятностей

	3	2	1	0	-1	оценка
	«ОЛК»	«ЛК»	«К»	«»		
К	$1/(1+1)$	1	$1/(2+1)$	$5/(5+1)$	-	$5/36$
Л	$1/(1+1)$	1	$2/(2+1)$	-	-	$1/3$
О	$1/(1+1)$	-	-	-	-	$1/2$
Ь	$1/(1+1)$	1	$1/(2+1)$	$1/(5+1)$	1	$1/36$

КЛОКОКЛОЛКОЛКЬ?



Оценки вероятностей

	3	2	1	0	-1	оценка
	«ЛКЪ»	«КЪ»	«Ъ»	«>»		
К	-	-	-	5/14	-	5/14
Л	-	-	-	4/14	-	4/14
О	-	-	-	4/14	-	4/14
Ъ	-	-	-	1/14	-	1/14