

Методы сжатия данных

Багрин, Ратушин, Смирнов, Окин

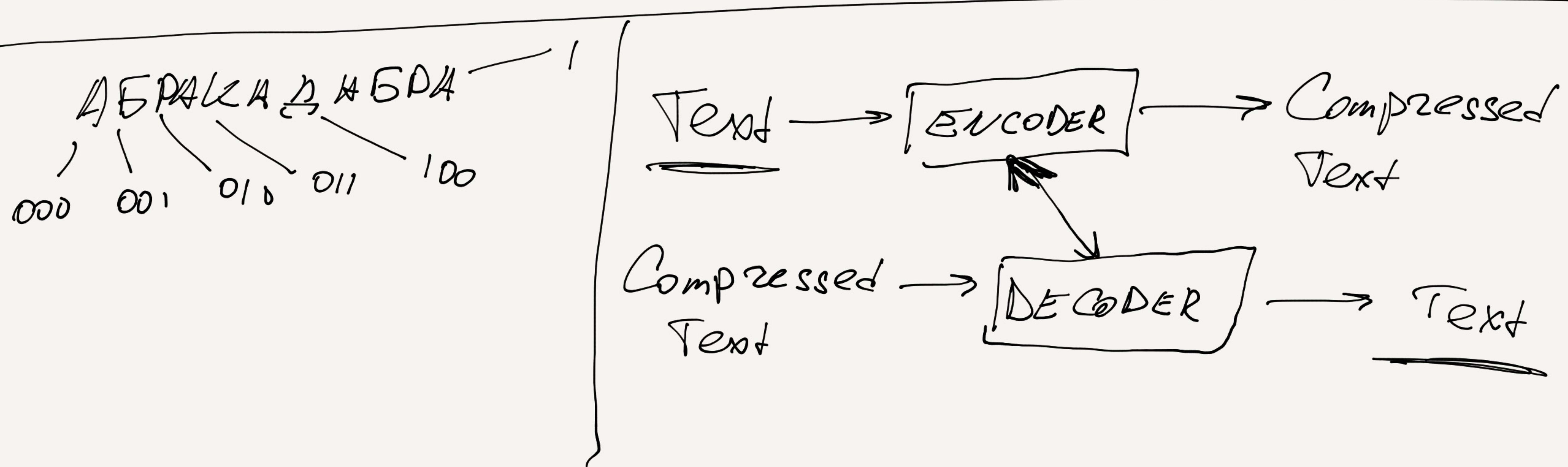
Managing Big Data

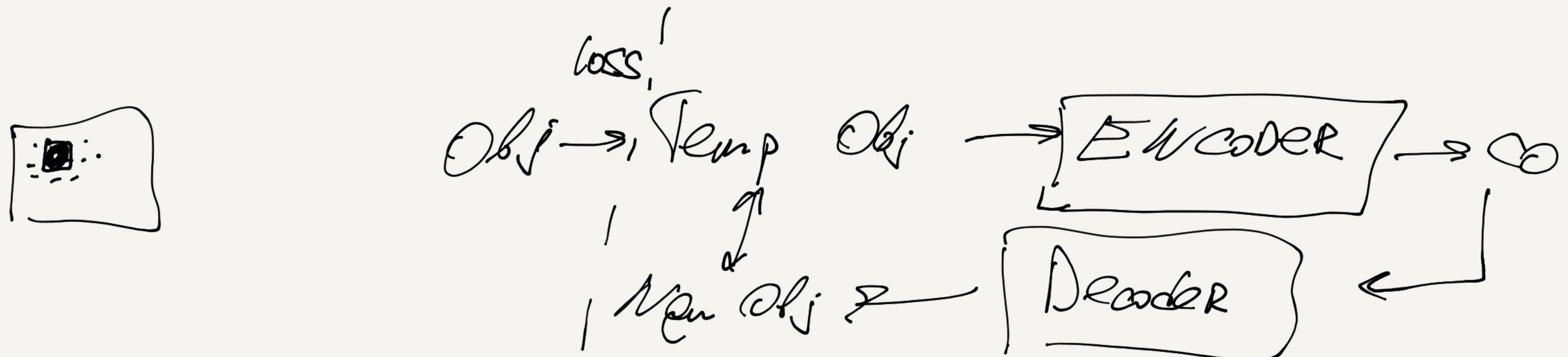
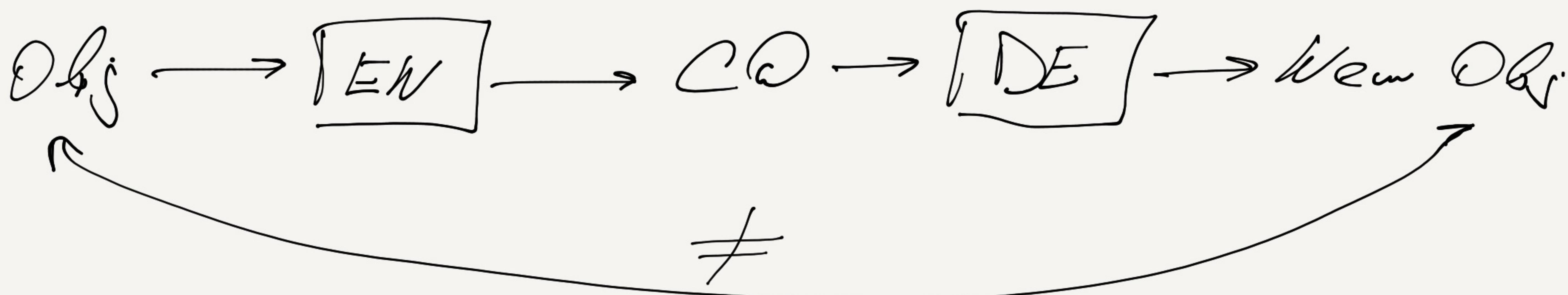
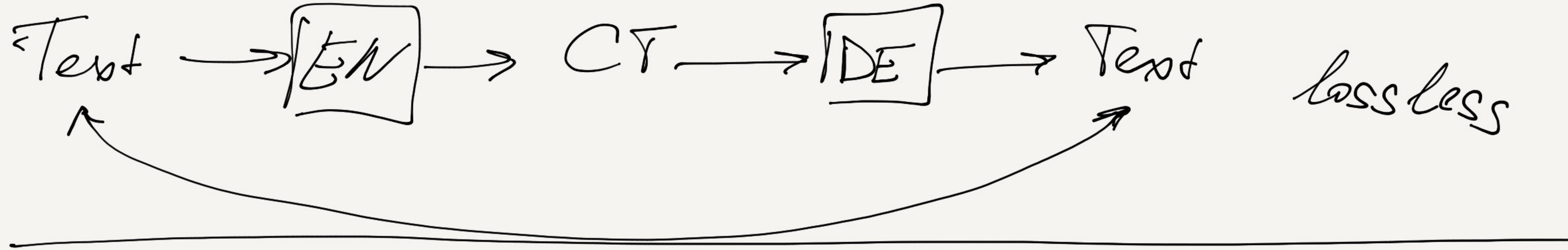
Compressing and Indexing Documents and Images // 1-2

Witten, Moffat, Bell

Algorithms and Theory of Computation Handbook // ch. 12

Atallah



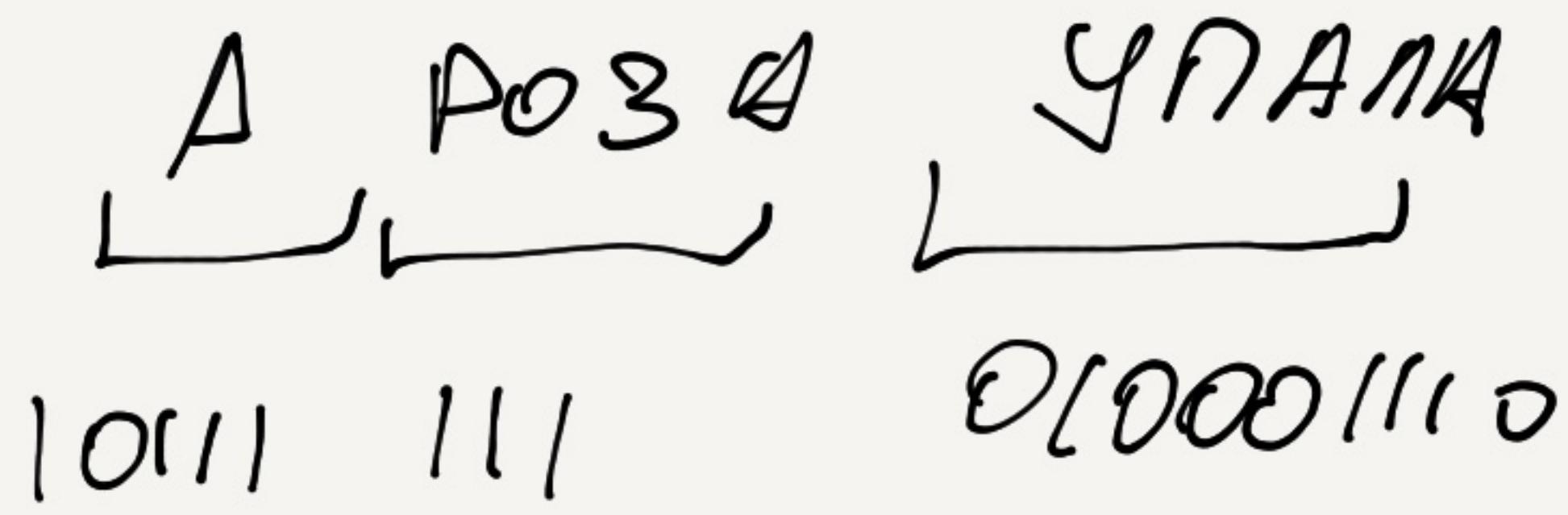


- Синтаксические

- Семантические

Характеристическое купирование

L2 - 77, L2 - 78.

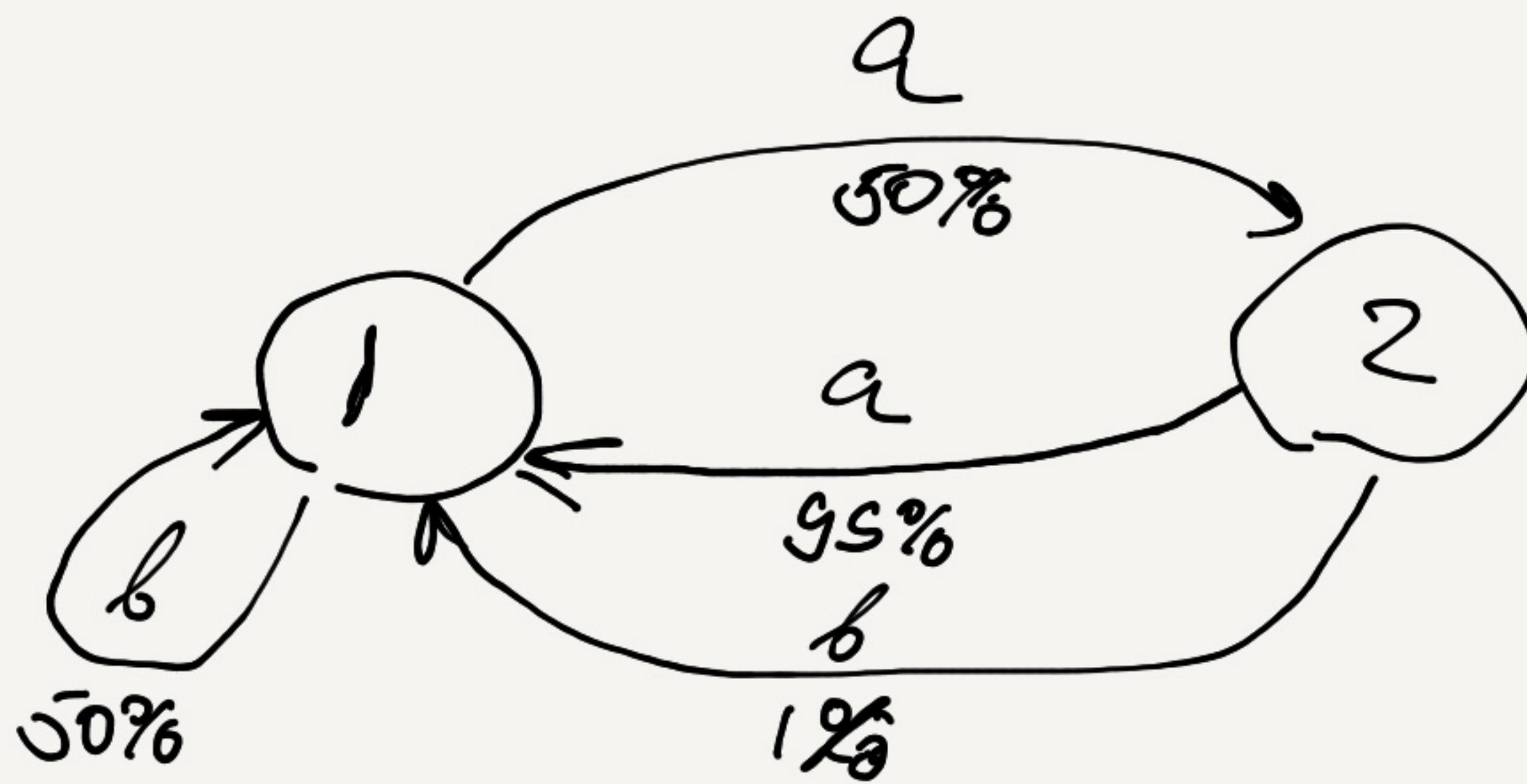
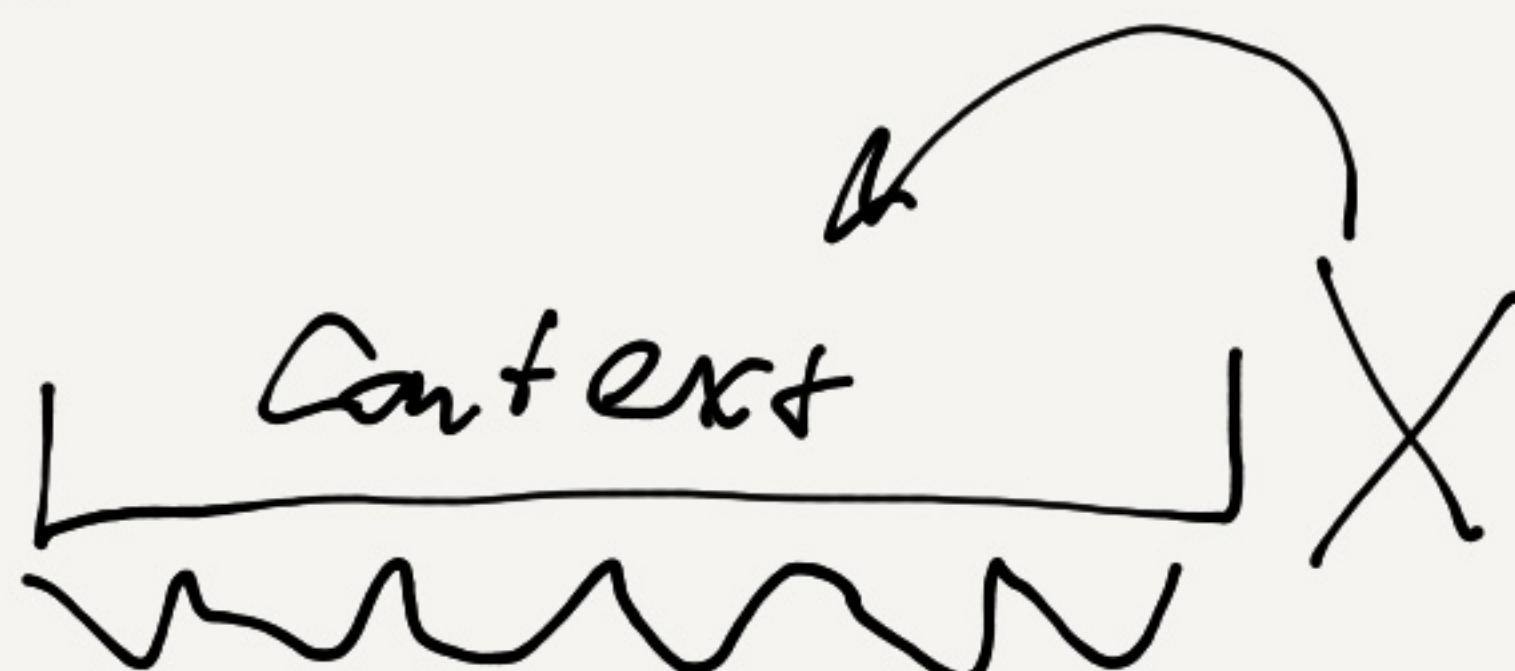


$u = 2\%$

7 бит.



1 бит.

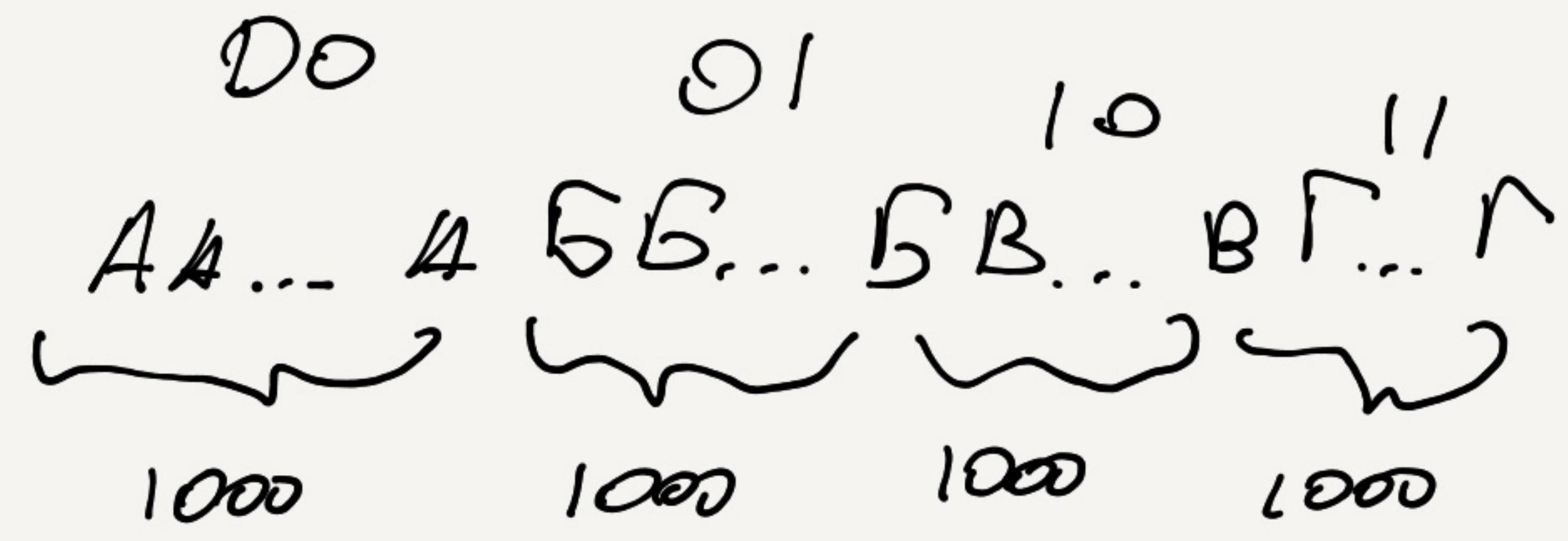


Context
[history + current]

Finite context model.

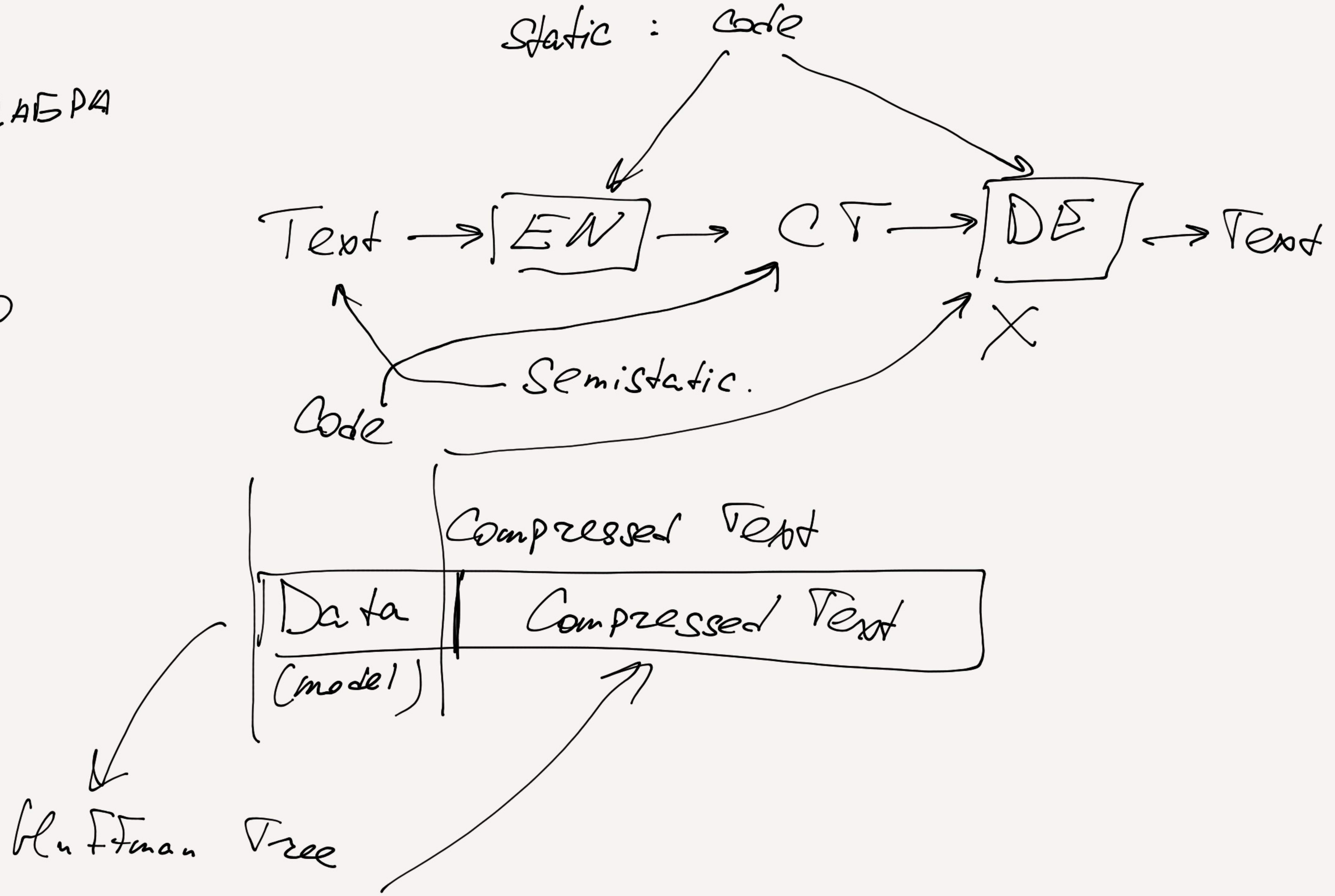
- PAGED
- DEMQ

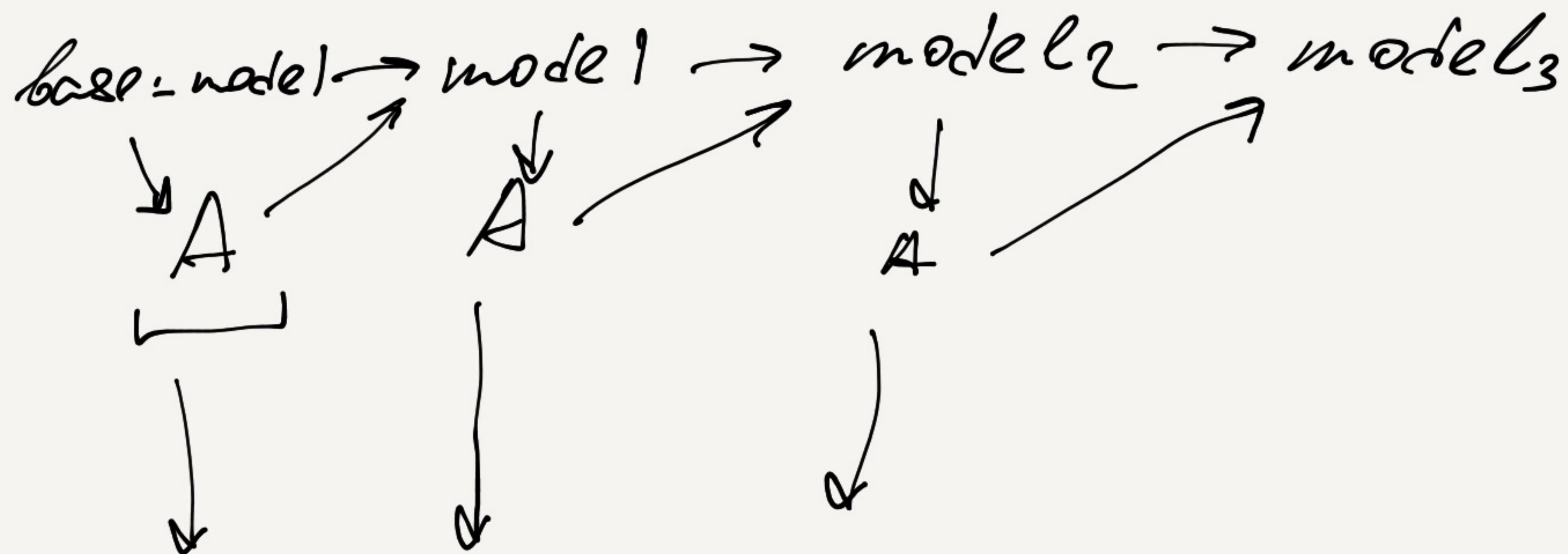
- СТАТИЧЕСКИЕ
- ПОЛУСТАТИЧЕСКИЕ / ДОЛЖНЫ АДАПТИРОВАТЬСЯ
- АДАПТИВНЫЕ.



АБСОЛЮТНАЯ

A - 2 11
B - 3 010
P - 2
I - 4
K - 5





B

Code Code 1 Code 2

10 10 (())

$$1000 \cdot 2 = 2000$$

$$998 \cdot 1 + 2 \cdot 2 = 1002.$$

$a_1 \quad a_2 \dots \quad a_n$
 $p_1 \quad p_2 \dots \quad p_n$

$$H = - \sum_i p_i \log_2 p_i$$

$$\sum_i p_i = 1.$$

СЛОВАРНЫЕ

L2-77 / L2-78

Lempel, Ziv

C2-1 / C8-2

$z_L - 1 \rangle z_L - 2$



$\underline{<0, 0, a>}$

$<0, 0, b>$

$\underline{<2, 1, a>}$

$<3, 2, \underline{1}, \underline{b}>$

$<6, 4, \underline{b}, \underline{b}>$

$<2, 2, \underline{b}, \underline{b}>$

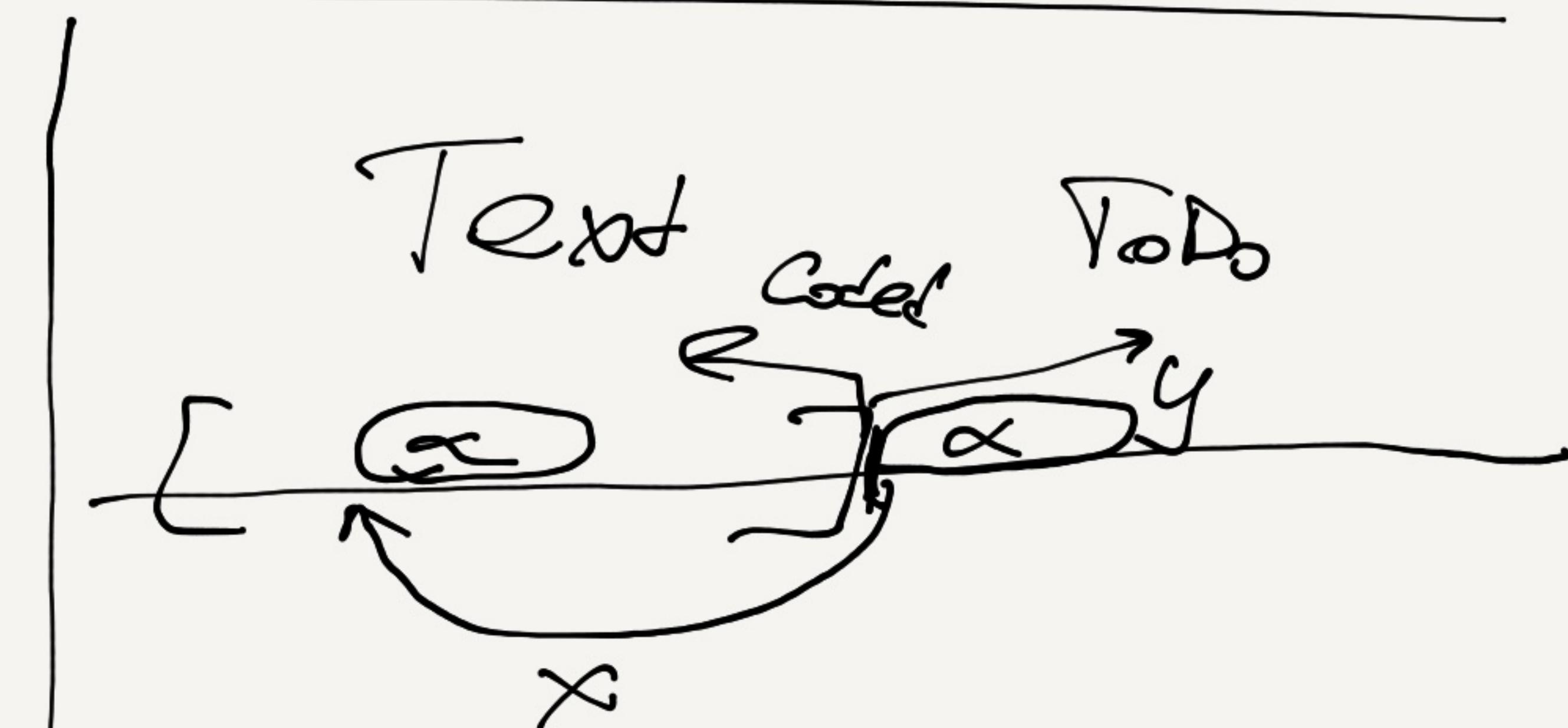
$<5, 5, a>$

a b a a b a b b a a b b b b b b b b b a
↑ ↑

<Shift, size, is_new>

$\leq 5, 2$, true, $\omega \geq$

`< √, 3, false >`



$\langle x, \alpha), y \rangle$



x ... a

a b a a b a b a b b b b b b b a

$\langle 0, 0, a \rangle$

$\langle 0, 0, b \rangle$

$\langle 2, 1, a \rangle$

$\langle 3, 2, b \rangle$

$\langle 2, 2, b \rangle$

$\langle 2, 2, b \rangle$

$\langle 1, 8, a \rangle$

b b . b b b
b b b b

$\langle \text{shift}, \text{size}, \text{symbol} \rangle$

$\text{Shift} \geq \text{size}$

b b b b b b

size

shift



a b a a b a b a b b b b b b a

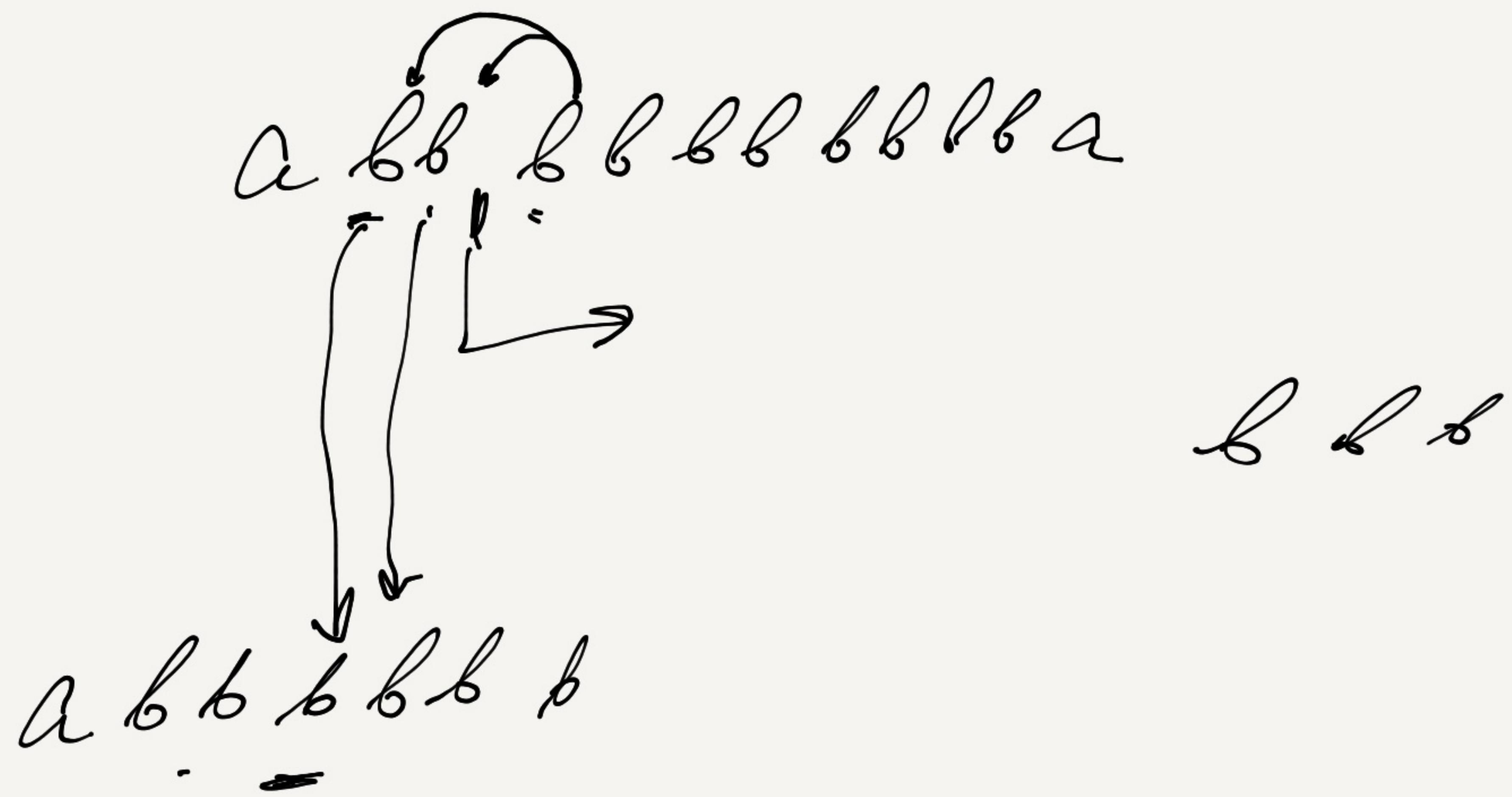
$\langle 1, 8, a \rangle$

a b a a b a b b b b b b a

a b a a b a b b b b b b a

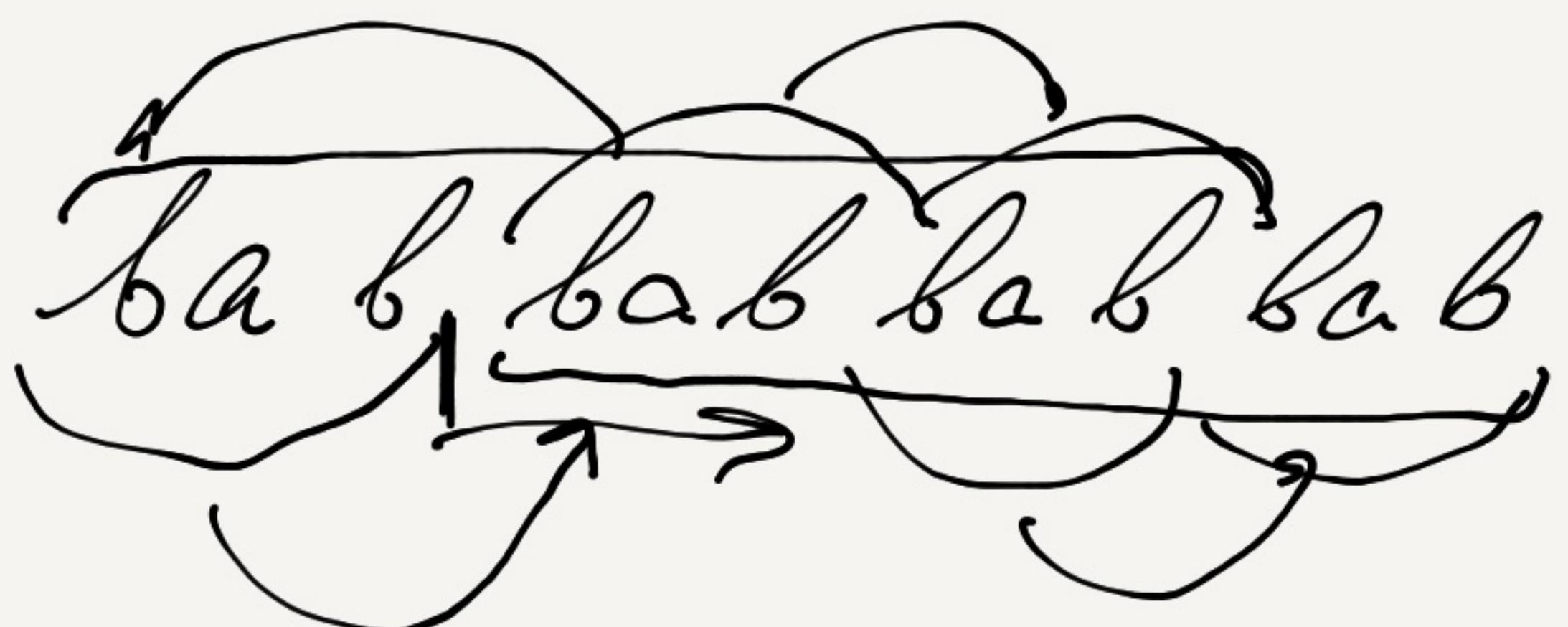
b a b b b b b a b b a b

$\langle 3, 9, \text{EOF} \rangle$
EOF.



shift, size, symbol

a b b b b b b b b a

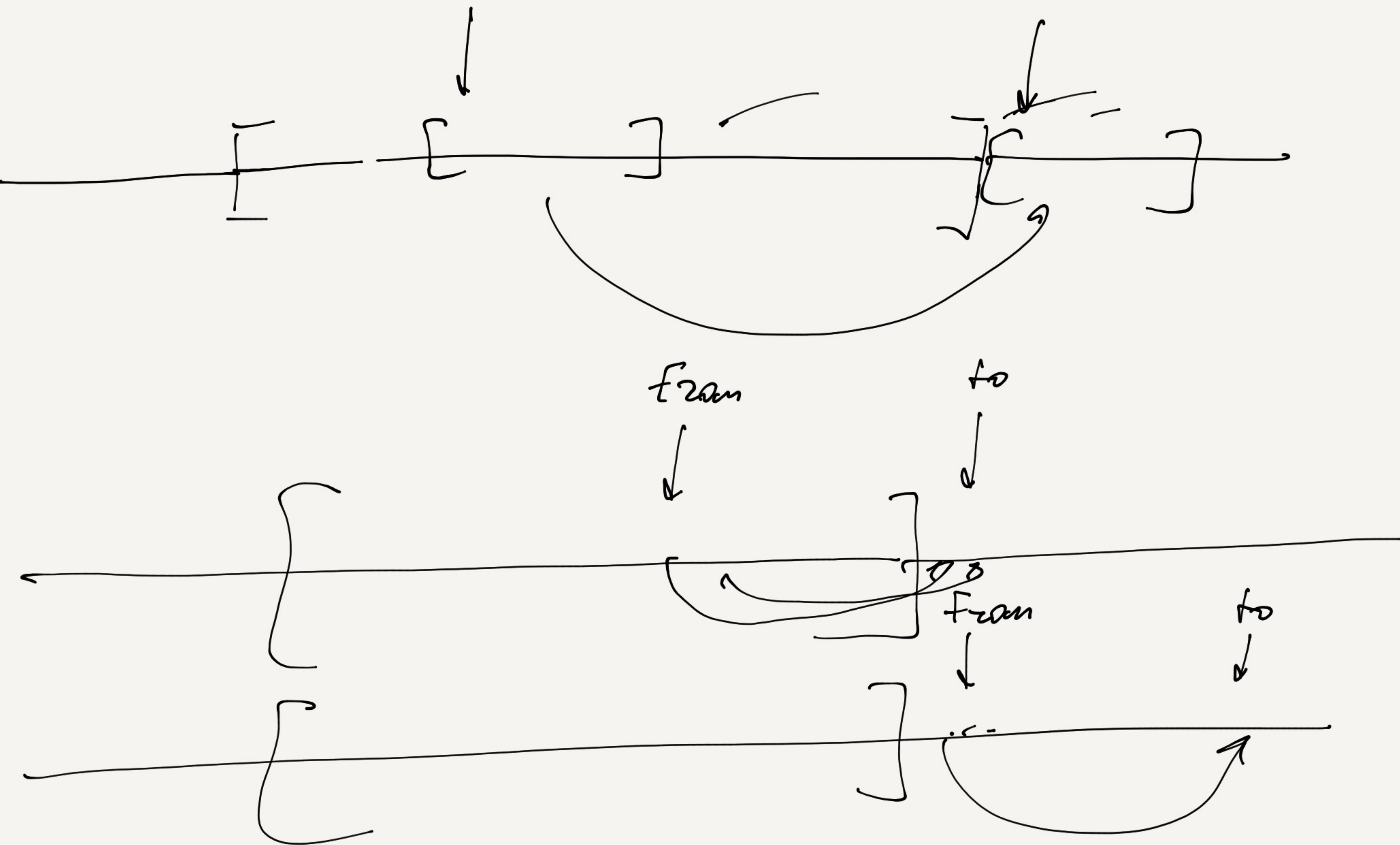


Shift

3, 9, EOM

$\dots \overbrace{bab, bab}^{\text{bab}}, \overbrace{bab}^{\text{bab}}, \overbrace{bab}^{\text{bab}}, \dots$

$\langle 3, g, \text{EOM} \rangle$



LZ-78

ACAGAATAGAGA
└─└─└─└─└─└─└─┘

$\begin{cases} <0, A> \quad A \\ <0, C> \quad C \\ <1, G> \quad AG \\ <1, A> \quad AA \\ <0, T> \quad T \\ <3, A> \quad AGA \\ <0, G> \quad G \\ <1, EOM> \end{cases}$

A-1
C-2
AG-3.
AA-4
T-5
AGA-6
G-7
 $\begin{cases} <0, A> \end{cases}$

$\begin{cases} <0, A> \\ <0, C> \\ <1, G> \\ <1, A> \\ <0, T> \\ <3, A> \\ <0, G> \\ <1, EOM> \end{cases}$

A-1
C-2
AG-3
AA-4
T-5
AGA-6
G-7

ACAGAATAGAGA

Flush

$\angle ZW$

ACAGA ATAGAGA

1, 2, 1, 3, 1, 1, 4, 7, 12.



word Y \notin Dict



\rightarrow word Y \rightarrow Dict:

A - 1

C - 2

G - 3

T - 4.

$\frac{AC}{AC}$ - 5

$\frac{CA}{CA}$ - 6

$\frac{AG}{AG}$ - 7

$\frac{GA}{GA}$ - 8

$\frac{AA}{AA}$ - 9

$\frac{AT}{AT}$ - 10

$\frac{TA}{TA}$ - 11

$\frac{AGA}{AGA}$ - 12

AA - 9

1, 2, 1, 3, 1, 1, 4, 7, 12

ACAGA ATAGAGA

.

A - 1

AT - 10

C - 2

TA - 11

G - 3

AGA - 12

T - 4

AC - 5

CA - 6

AG - 7

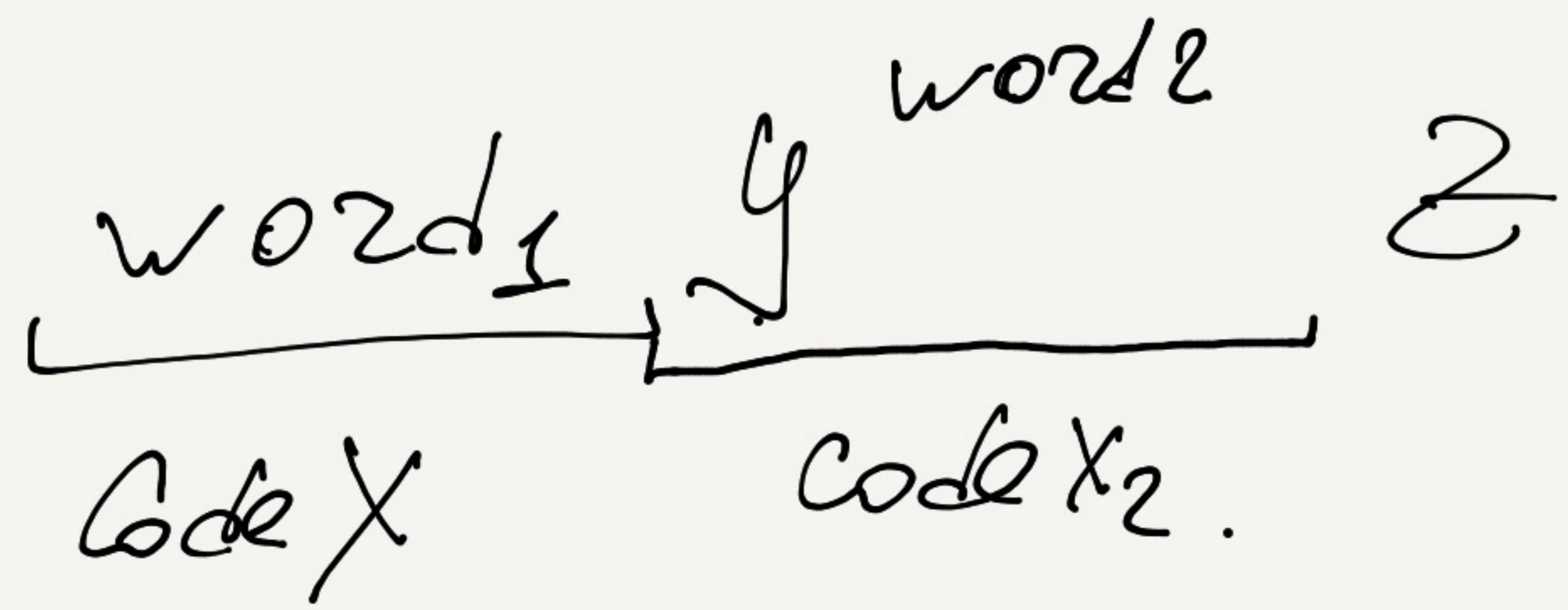
GA - 8

AA - 9

word X $\xrightarrow{\times}$ word X'

word X - Code

word \rightarrow AGA
 $\frac{AGA}{AGA}$ \rightarrow AGA



word₁ - code X

word₂ - code X₂.

word₃ ∈ Dict

word₃ → code Y

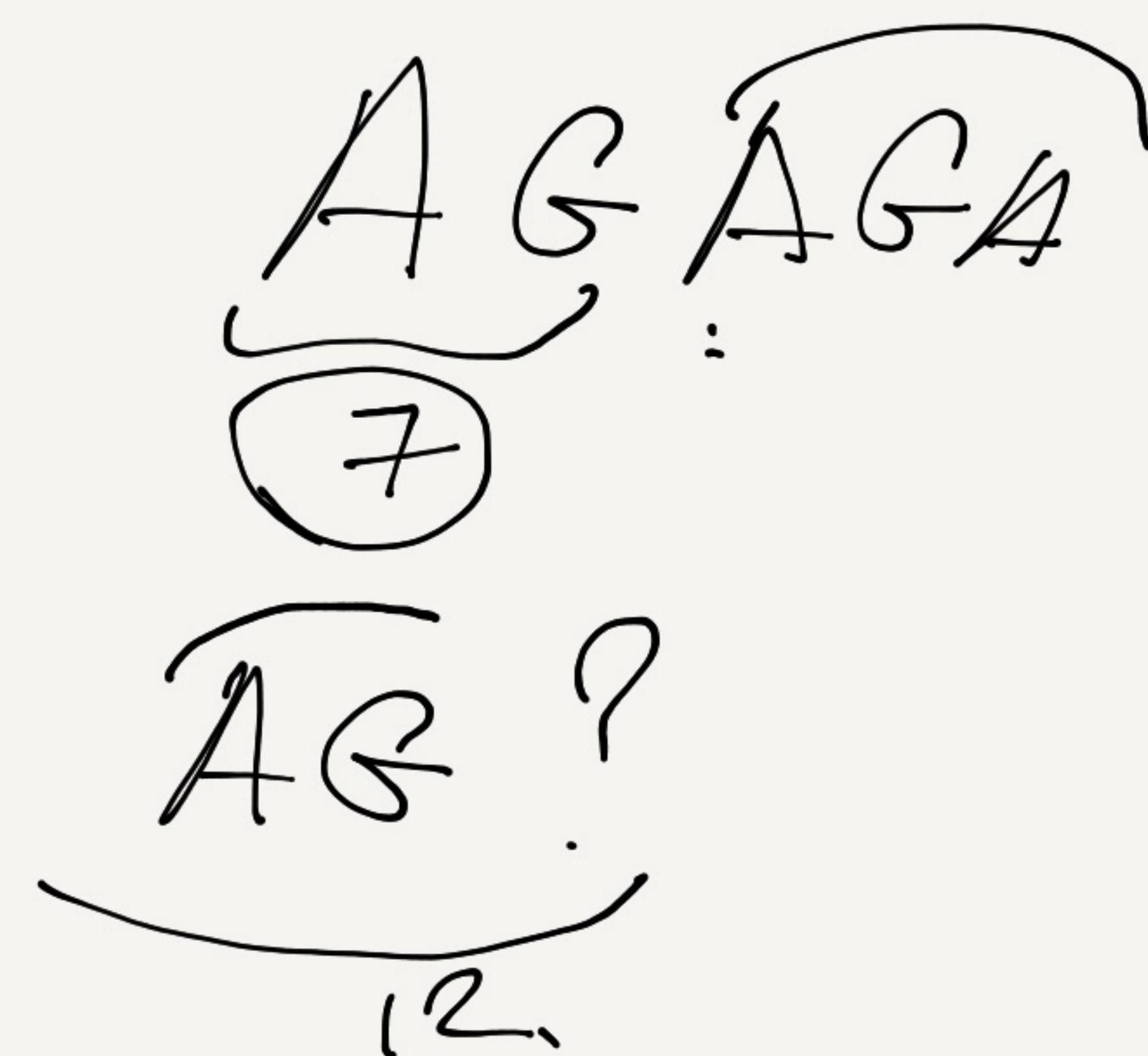
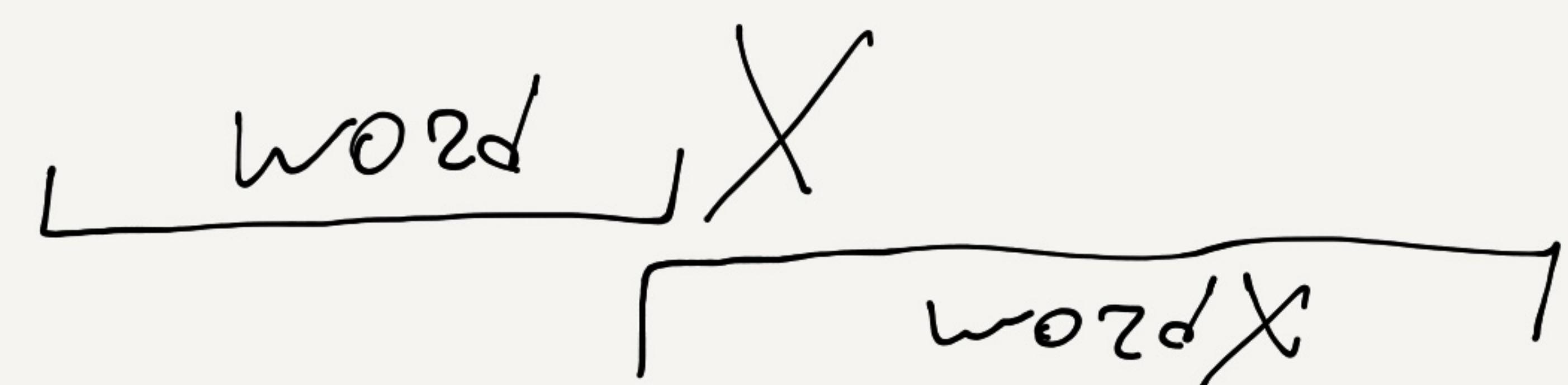
word₄ ∉ Dict

word₄ → code Y₂

?

12

AGA - 12



ACAGA ATAGA GAA

Encode

1, 2, 1, 3, 1, 1, 4, 7, 12

1, 2, 1, 3, 1, 1, 4, 7, 12

decode

ACAGA ATAGA GAA

A - 1

C - 2

G - 3

T - 4

AC - 5

CA - 6

AG - 7

GA - 8

AA - 9

AT - 10

TA - 11

AGA - 12

AGA

A-1
C-2
G-3
T-4
AC-5
CA-6
AG-7
GA-8
AA-9
AT-10
TA-11

AGA-12

AG AGA
C

7, 12.

AG - 7

TA - 11

AGA - 12.

, 7, 12

AG - 7.

TA - 11

12?

A G.

AG A C A

AG - 7

AGA - 12.

7.

ACAGA ATAGACA

1, 2, 1, 3, 1, 1, 4, 7, 5, 1

A-1

C-2

G-3

T-4

AC-5

CA-6

AG-7.

GA-8

AA-9

AT-10

TA-11

AGA-12

ACA-13

1, 2, 1, 3, 1, 1, 4, 7, 5, 1

ACA GA ATAGACA

A-1

C-2

G-3

T-4

AC-5

CA-6

AG-7

GA-8

AA-9

AT-10

TA-11

AGA-12

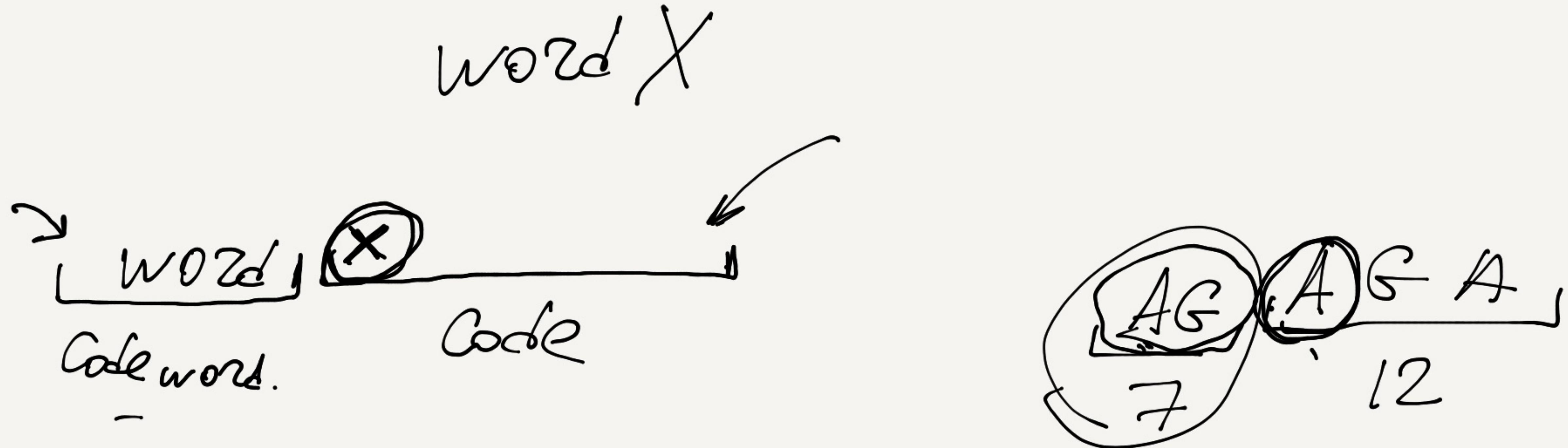
ACA-13

AGAGC

7, 7

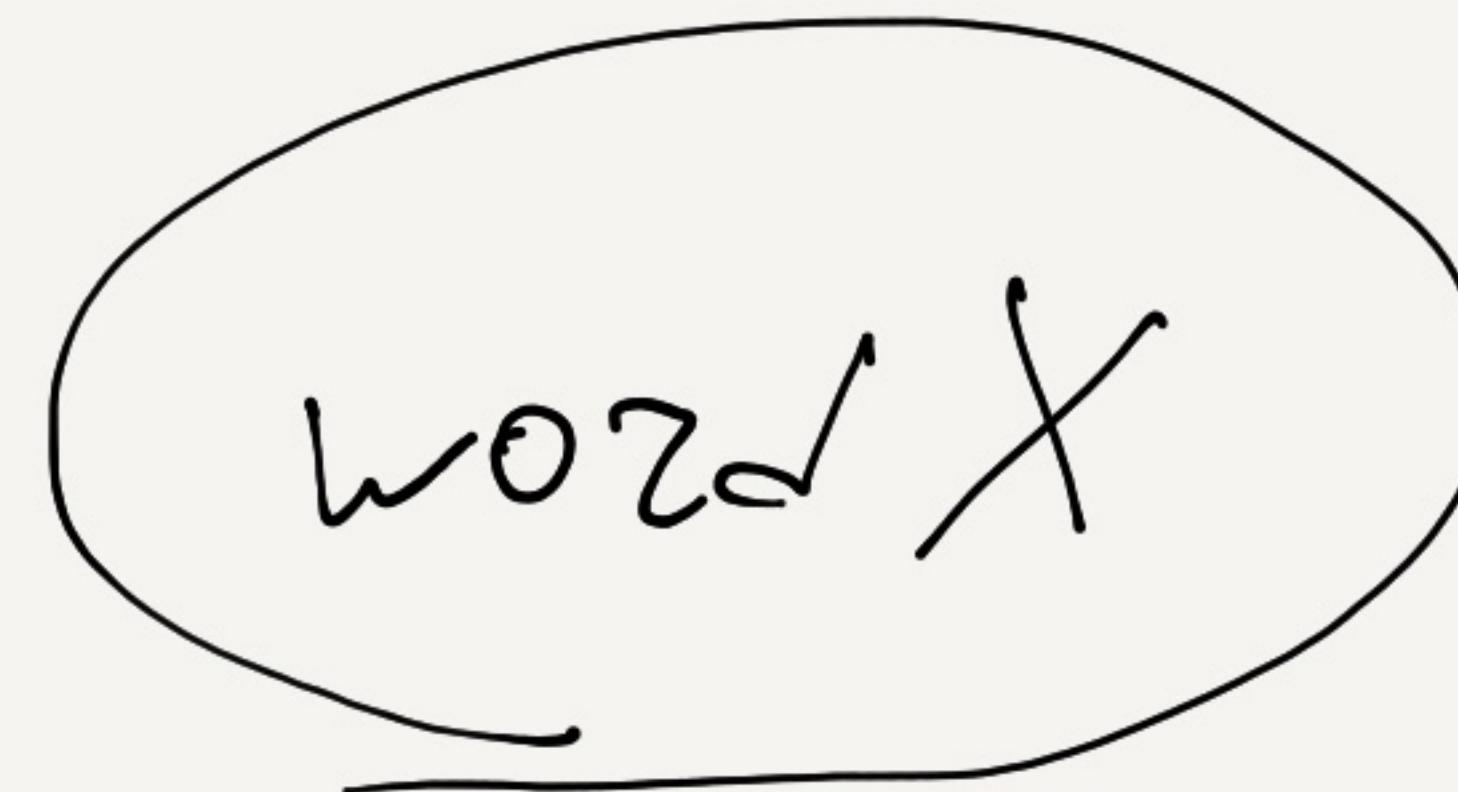
AGA - 12.

AGC - 13

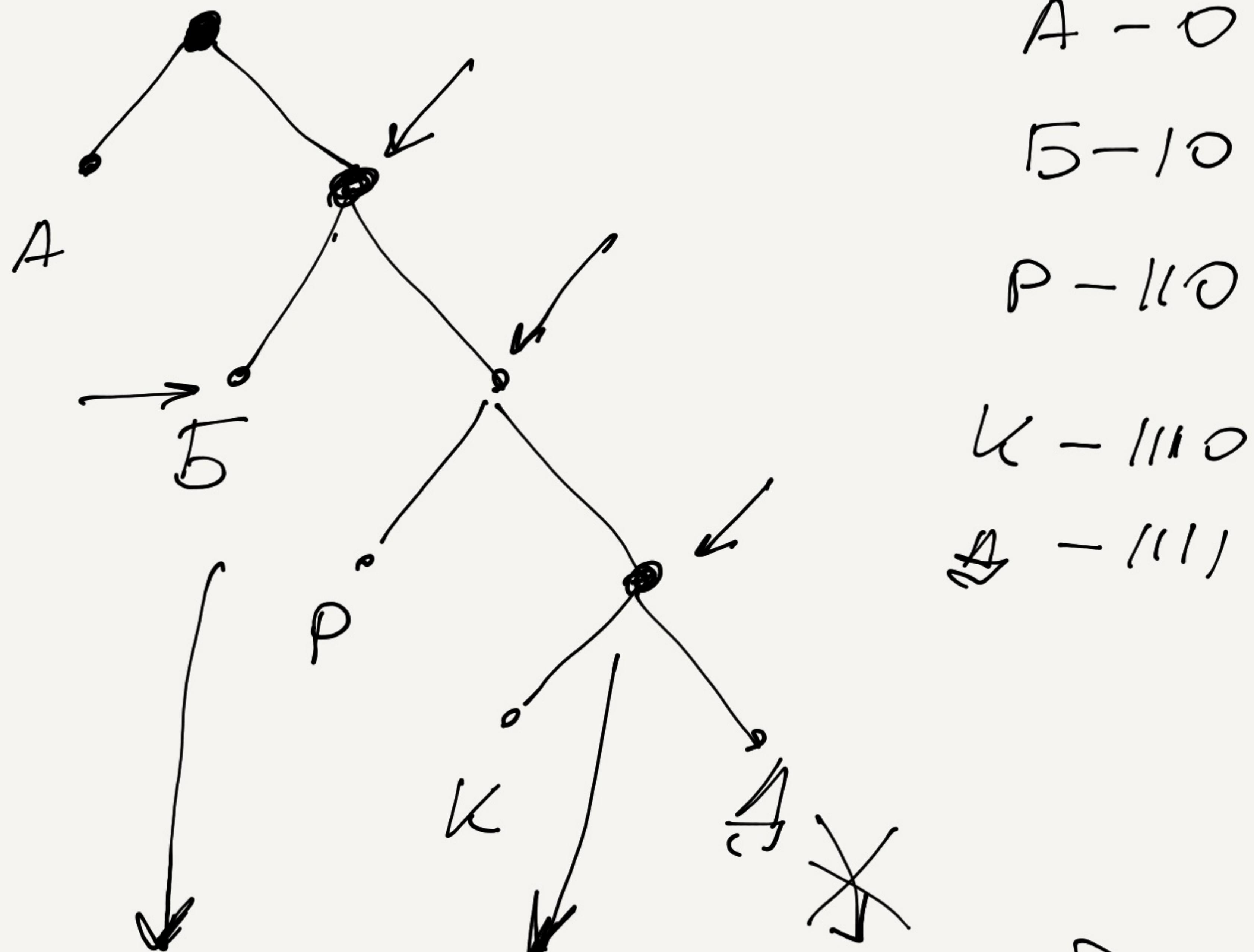


AGA - 12

Code & Dict.



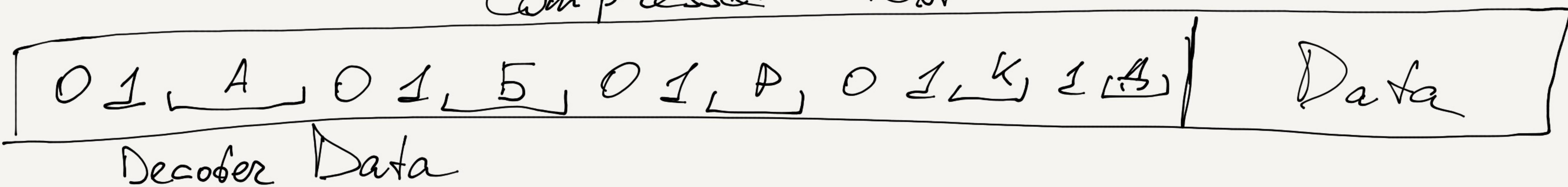
Codeword , Code



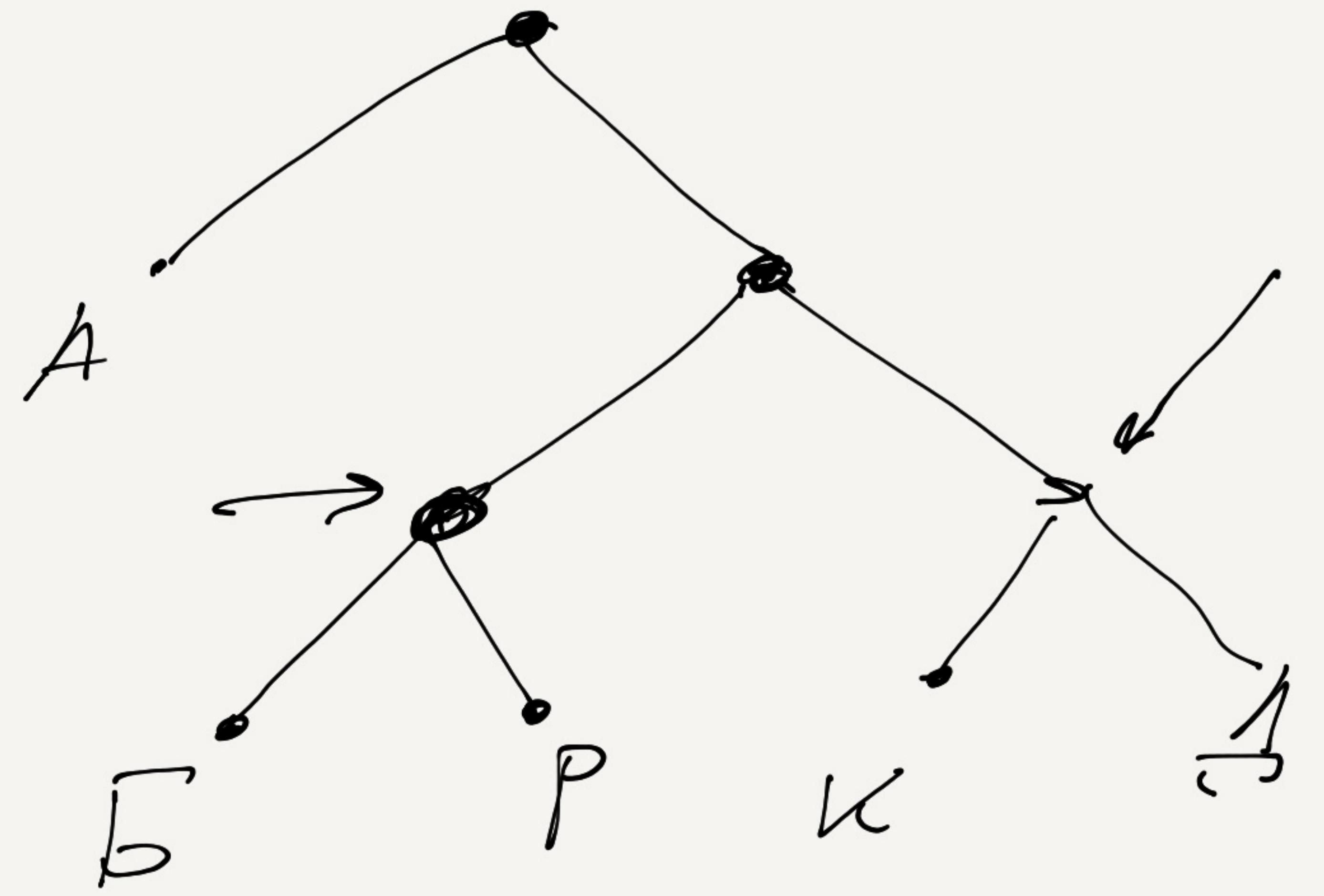
$A - 0$
 $B - 10$
 $P - 110$
 $K - 1110$
 $L - 1111$

Text \rightarrow EN \rightarrow CT \rightarrow DE \rightarrow D

Compressed Text



01_A_01_5_01_P_01_K_1(A) [Data]



O 1, A O O 1, B, 1, P O 1, K, 1 ~~A~~