

# Automatic Speech Recognition

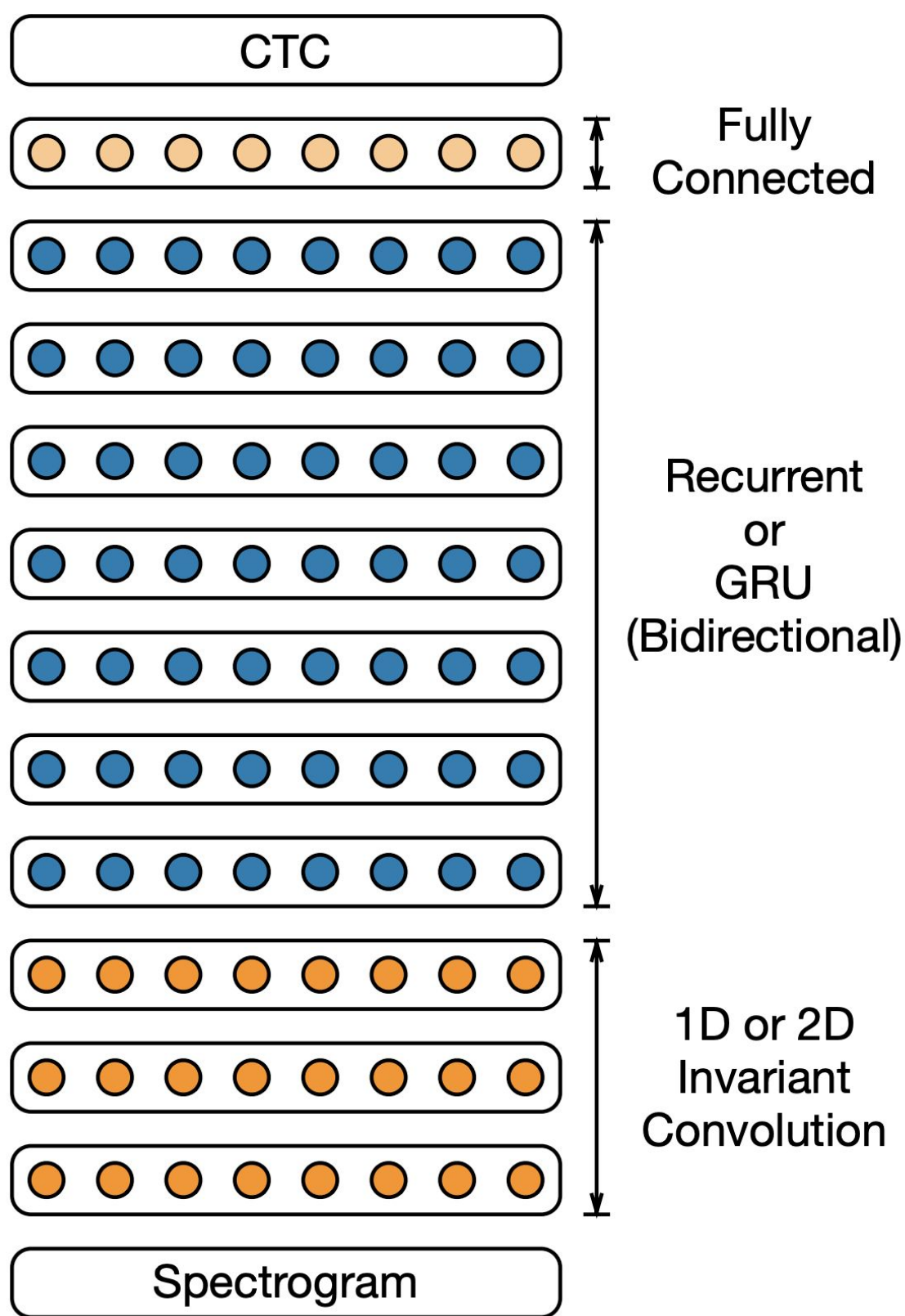
## лекция 2

Апарин Георгий

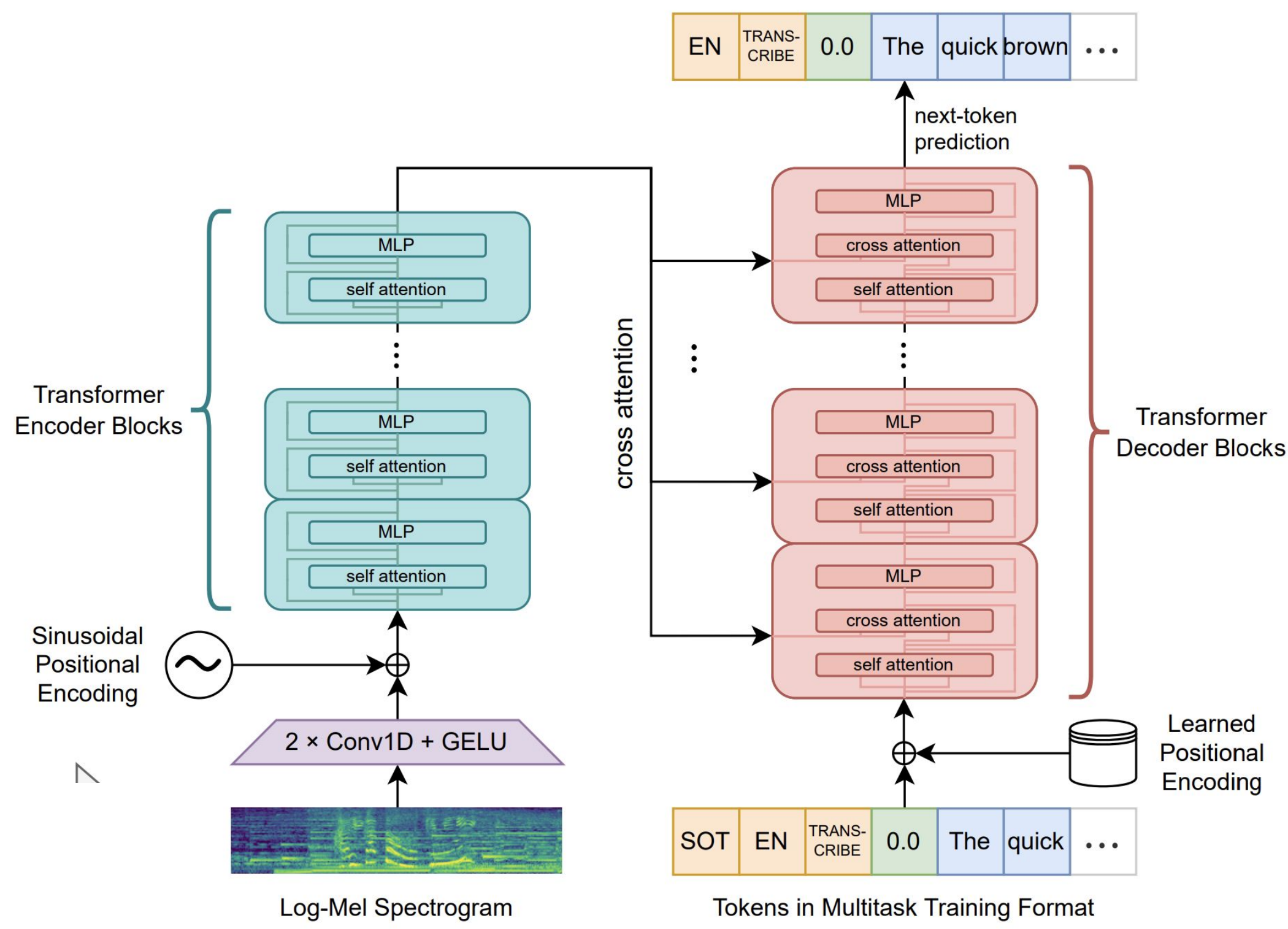
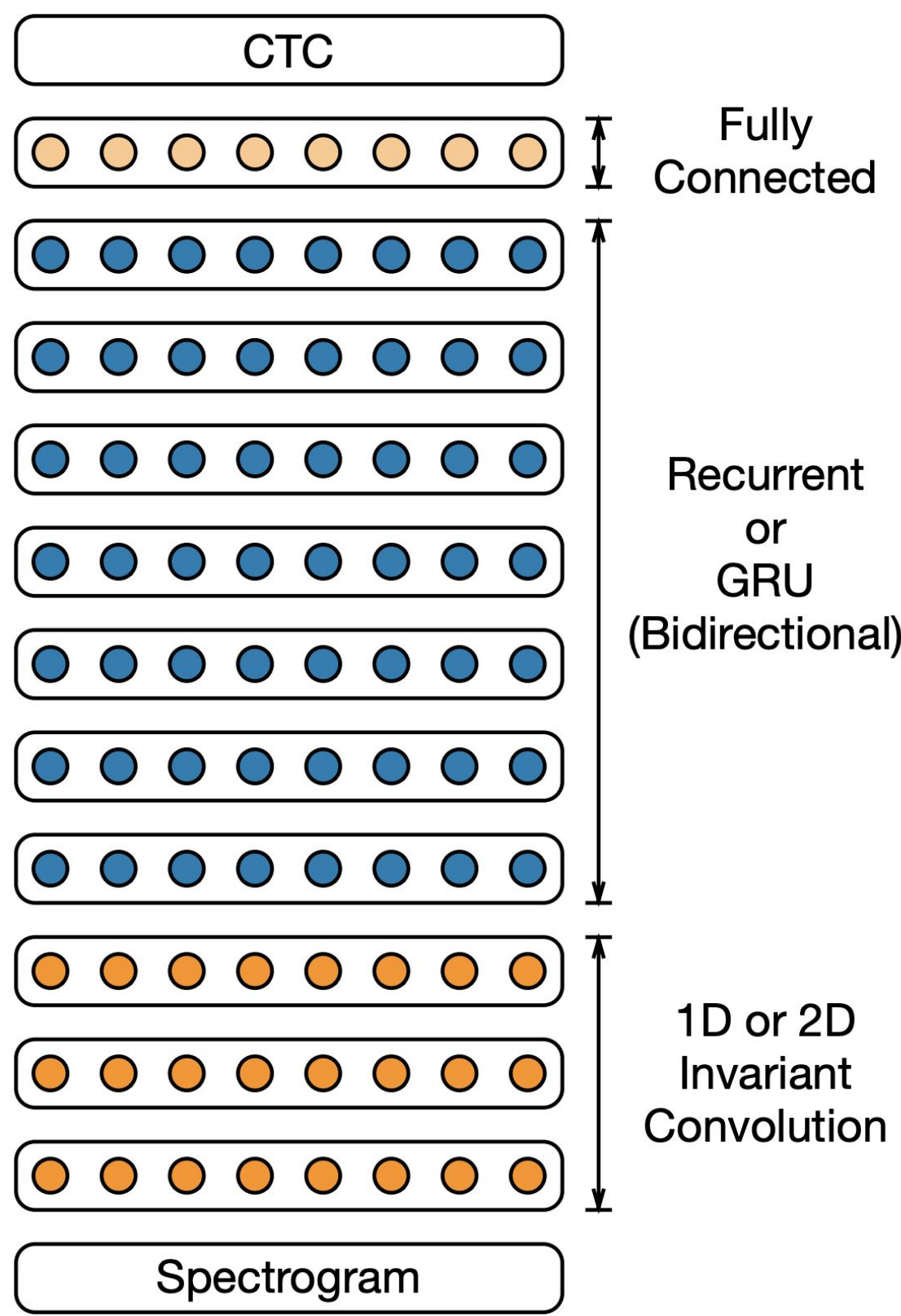
- Что такое доменные данные
- Какие бывают методы адаптации
- Shallow fusion
- Модели адаптации

# Доменные данные

# Вспомним архитектуры

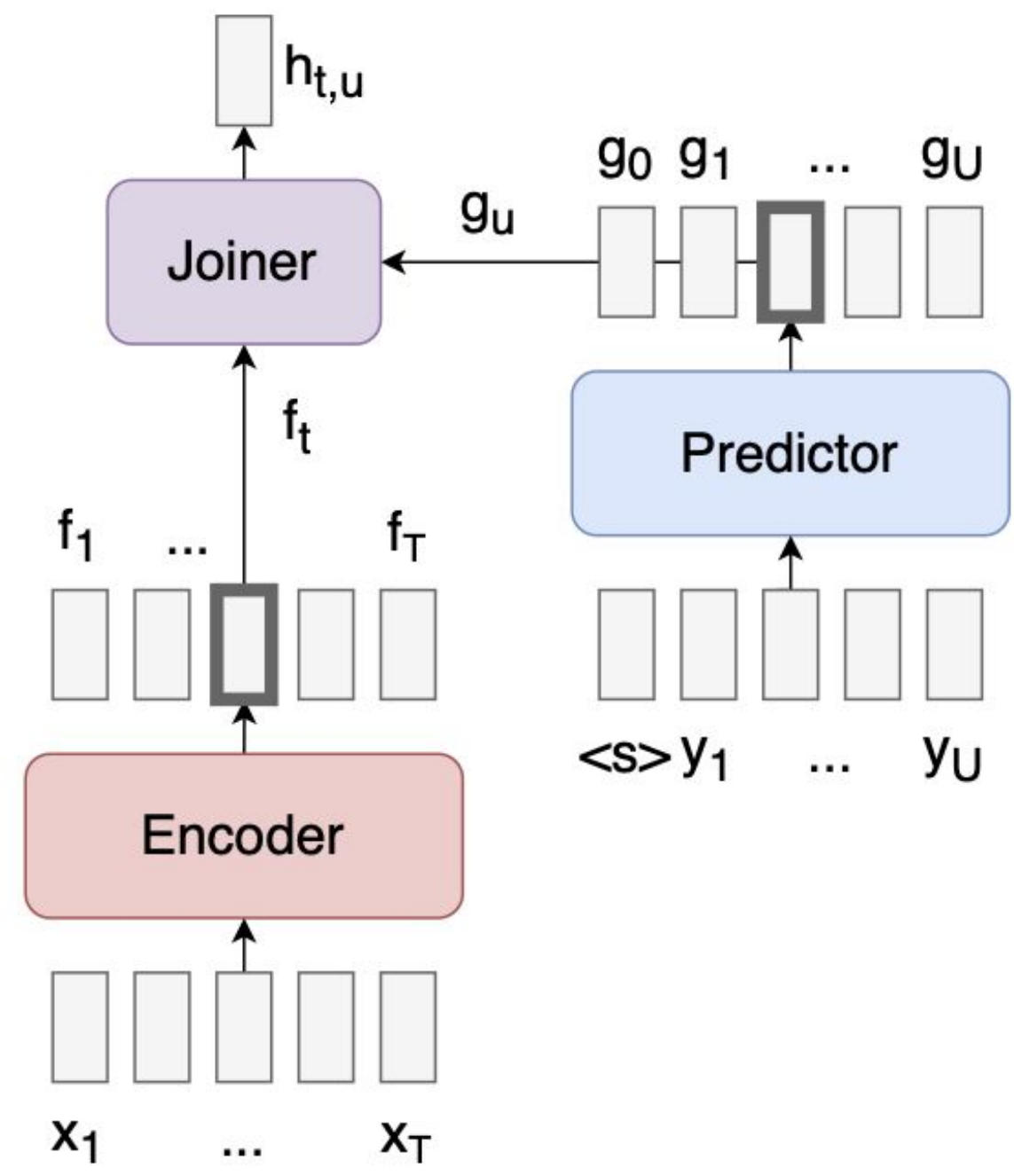
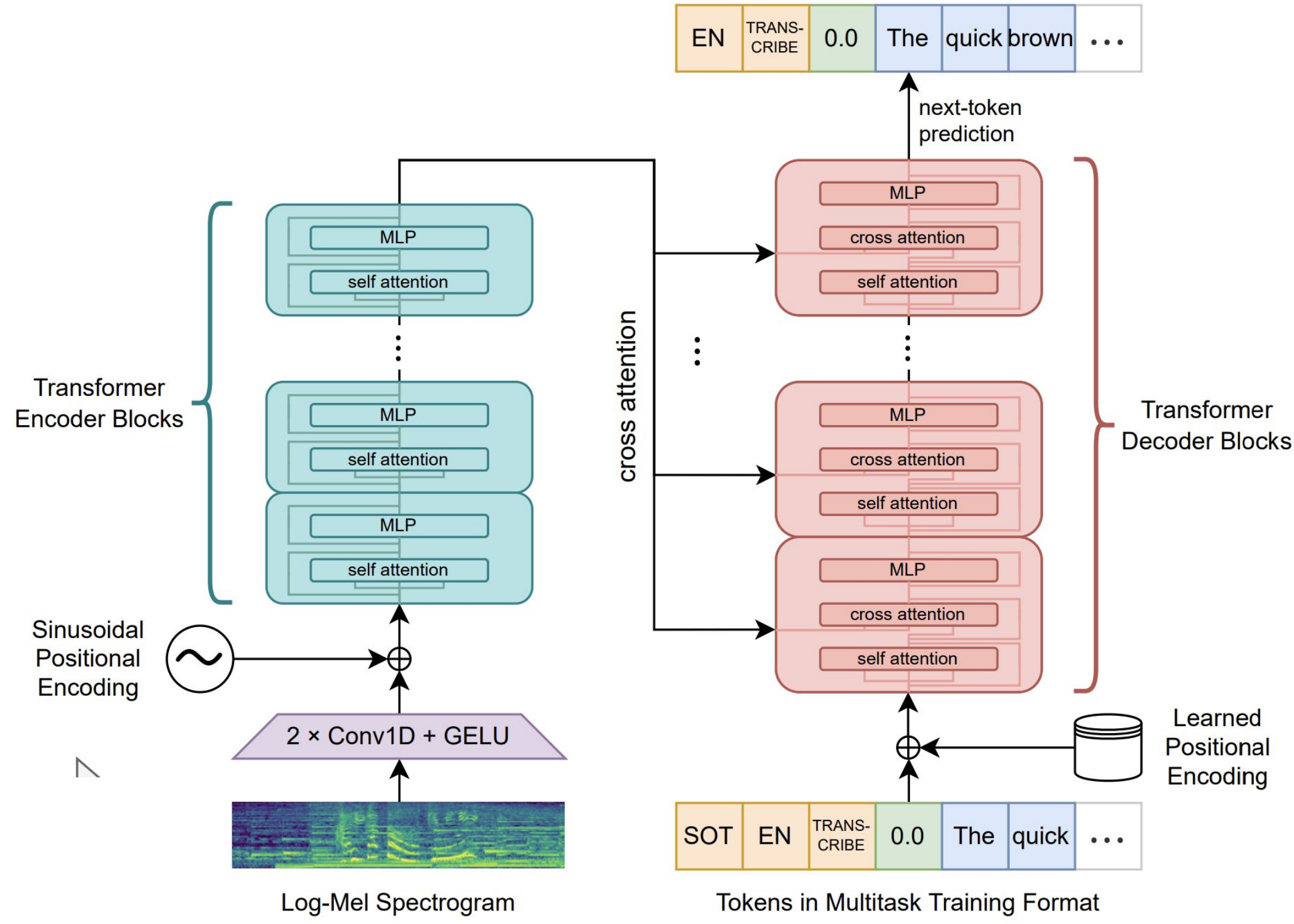
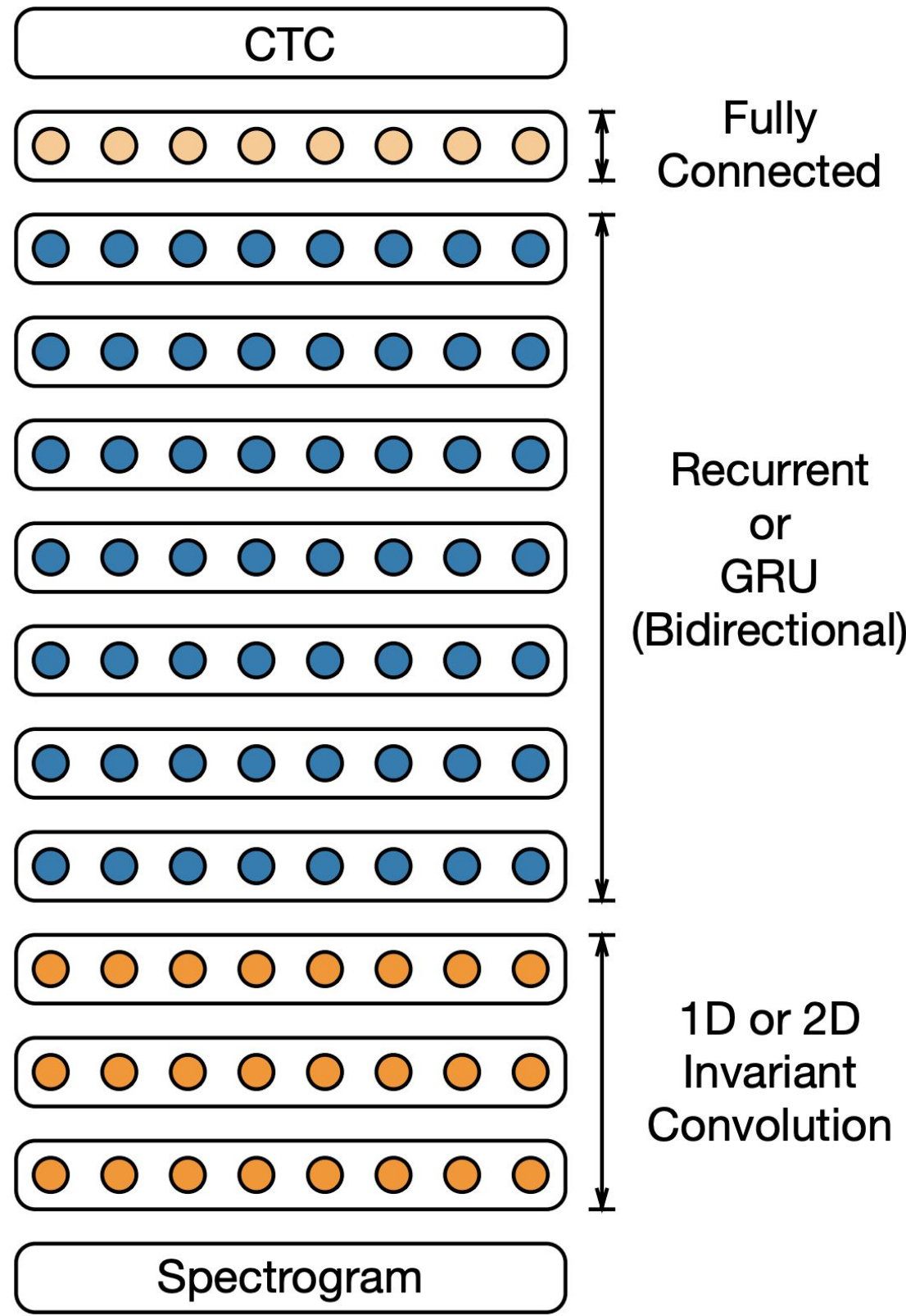


# Вспомним архитектуры



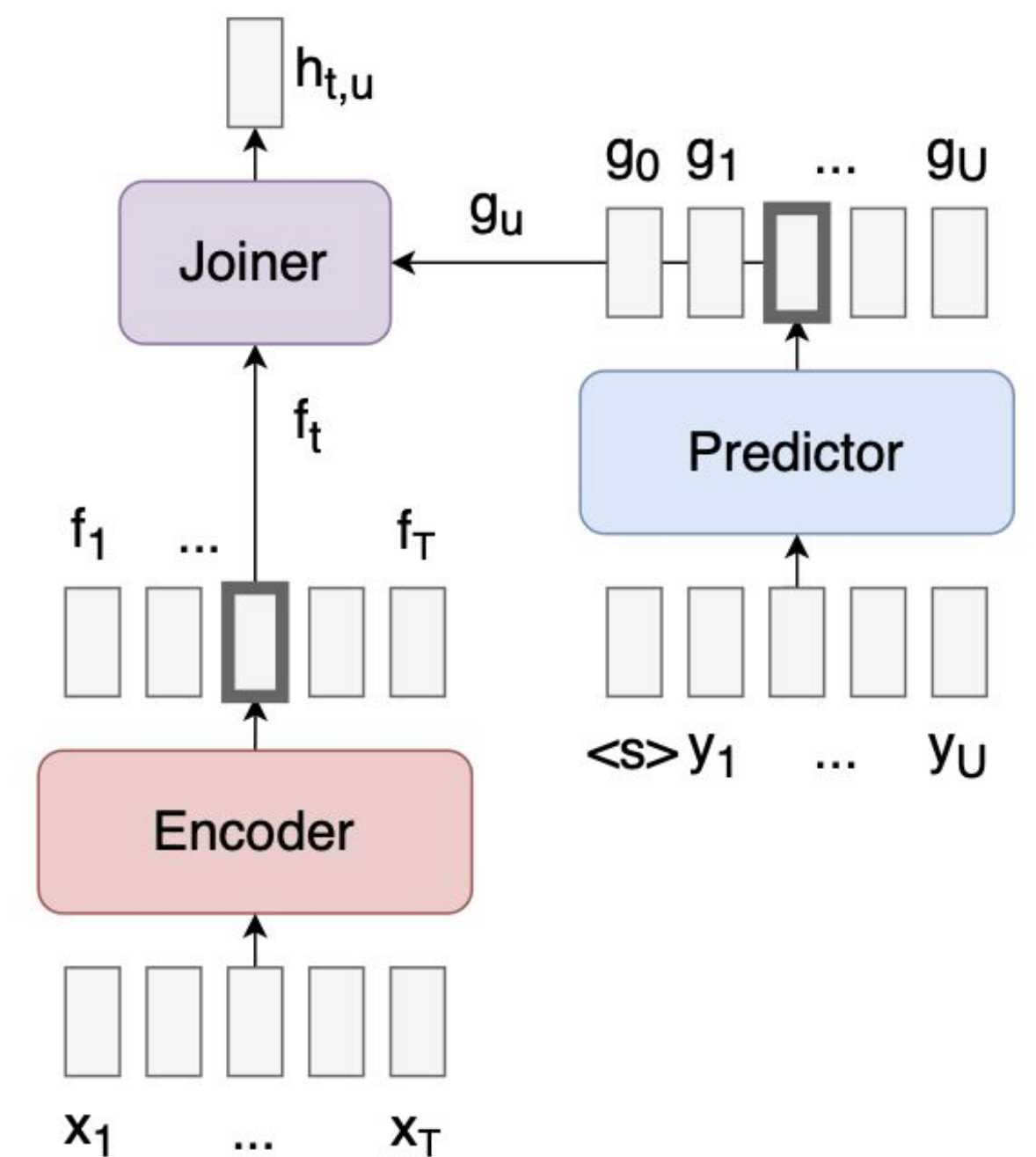
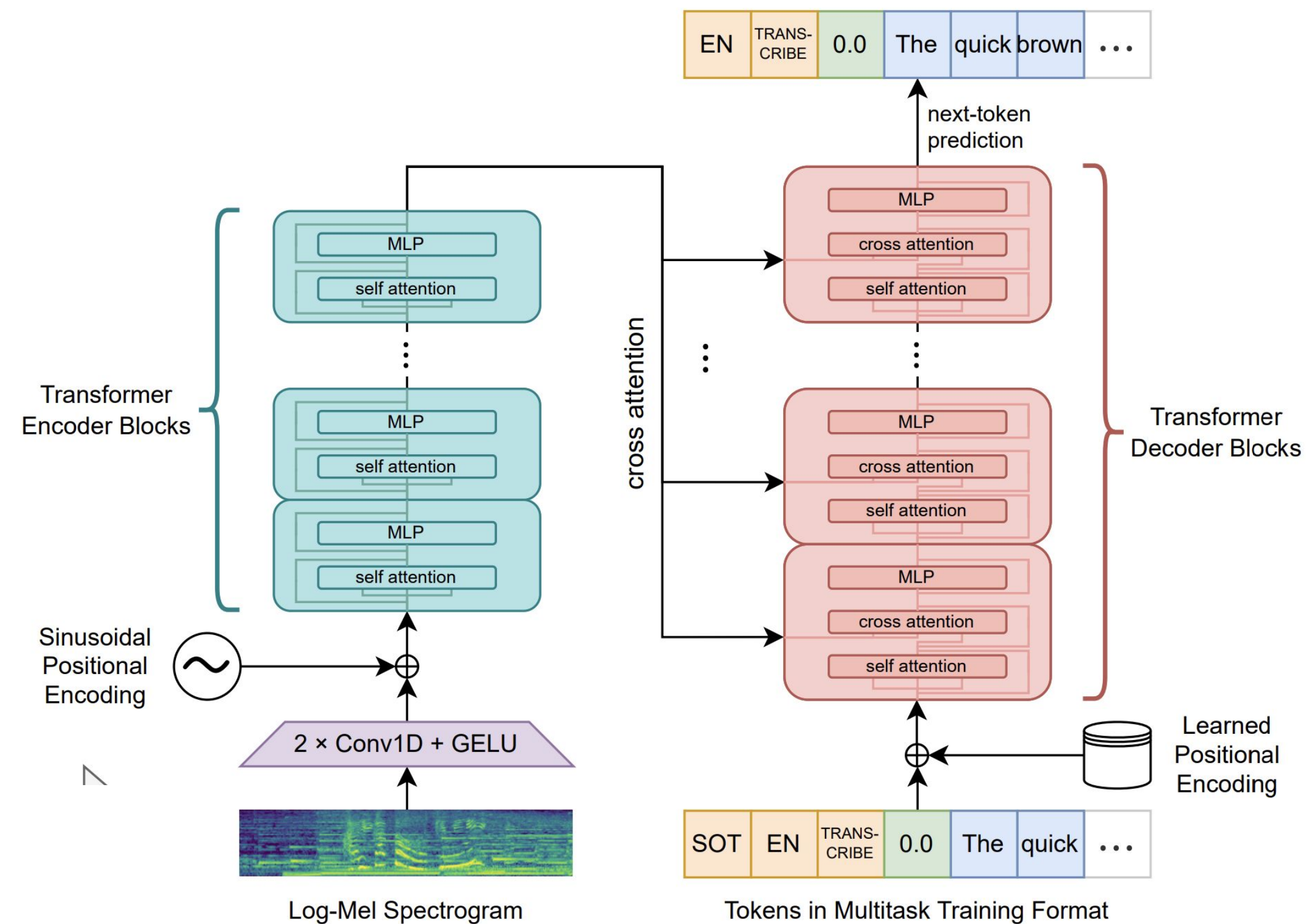
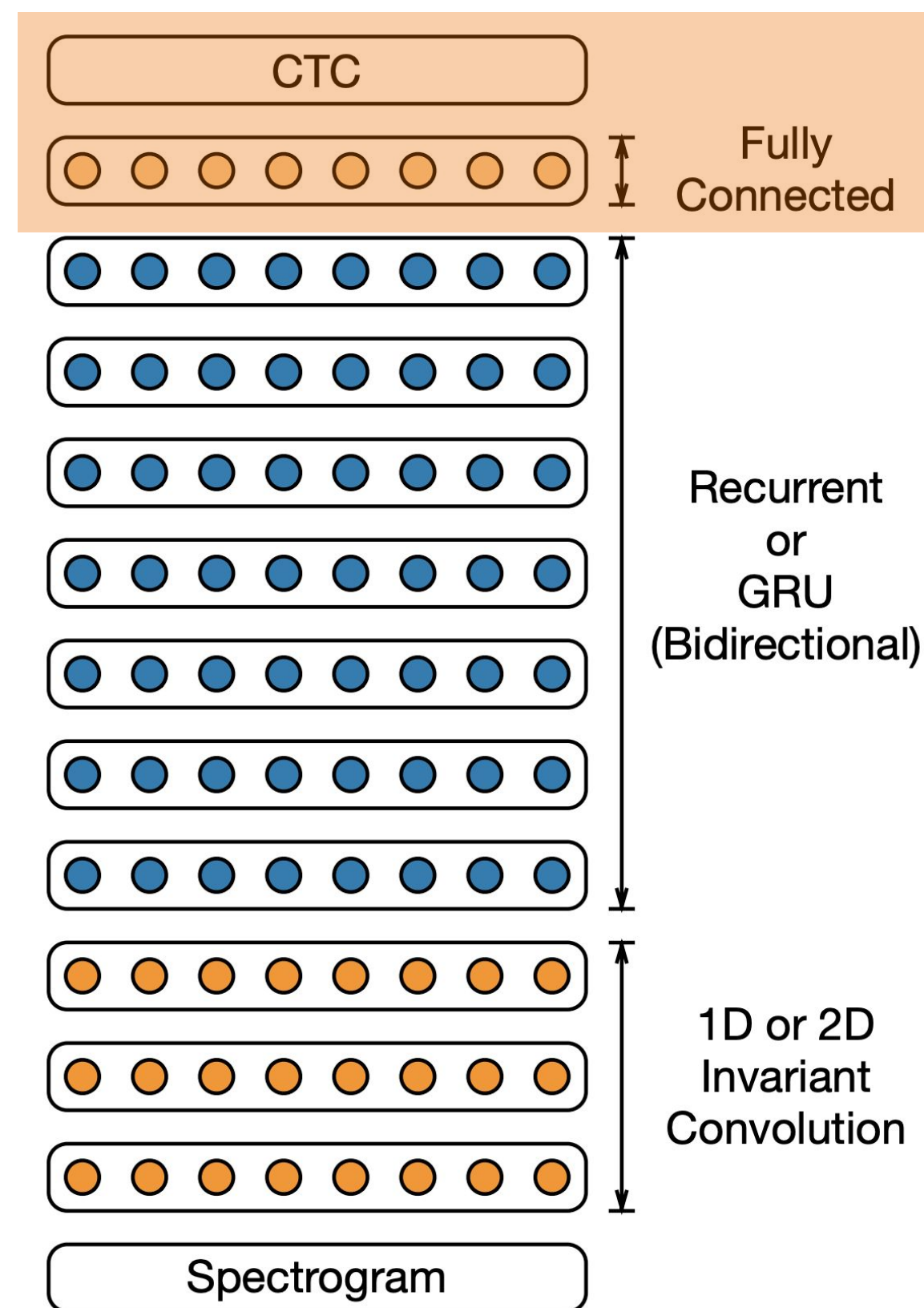


# Вспомним архитектуры



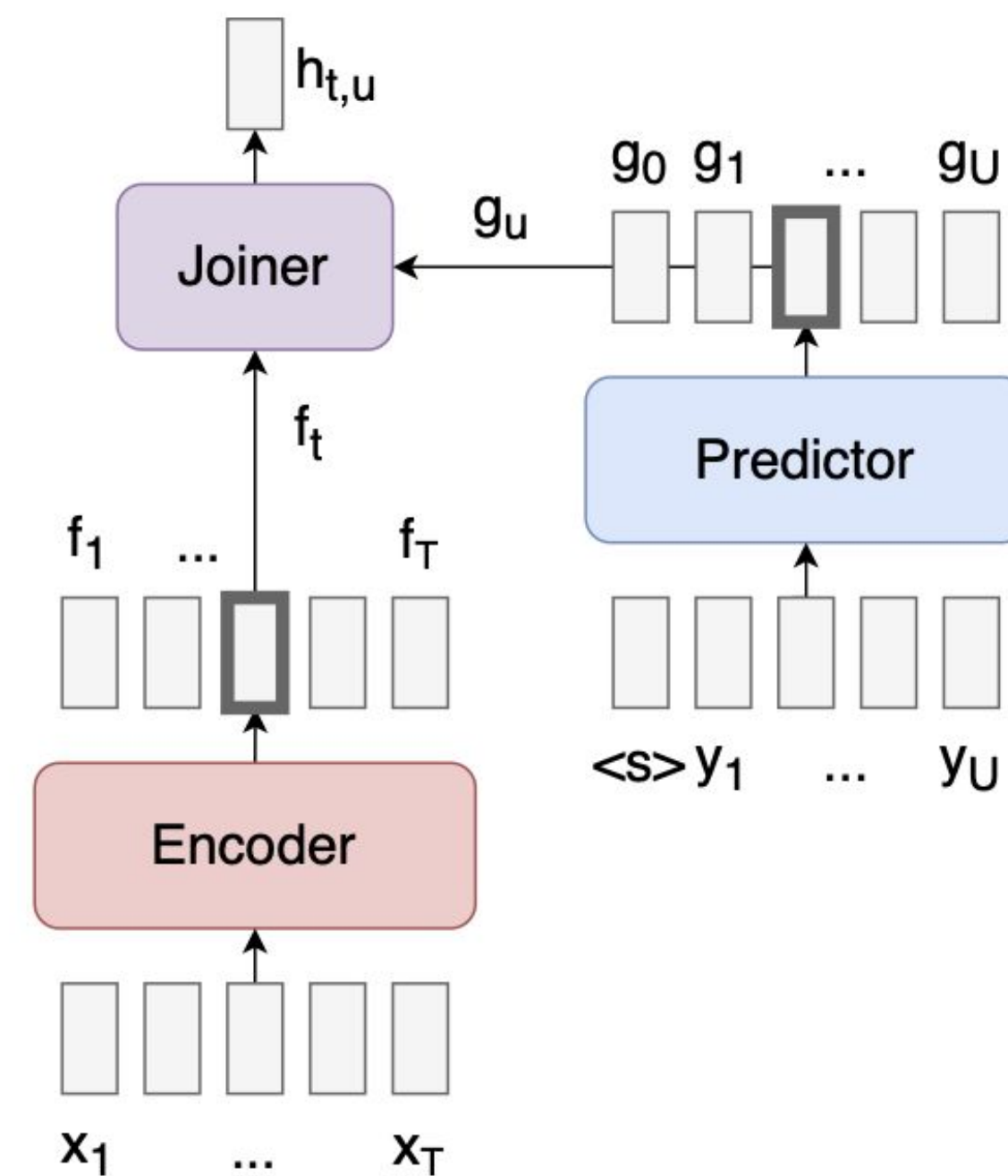
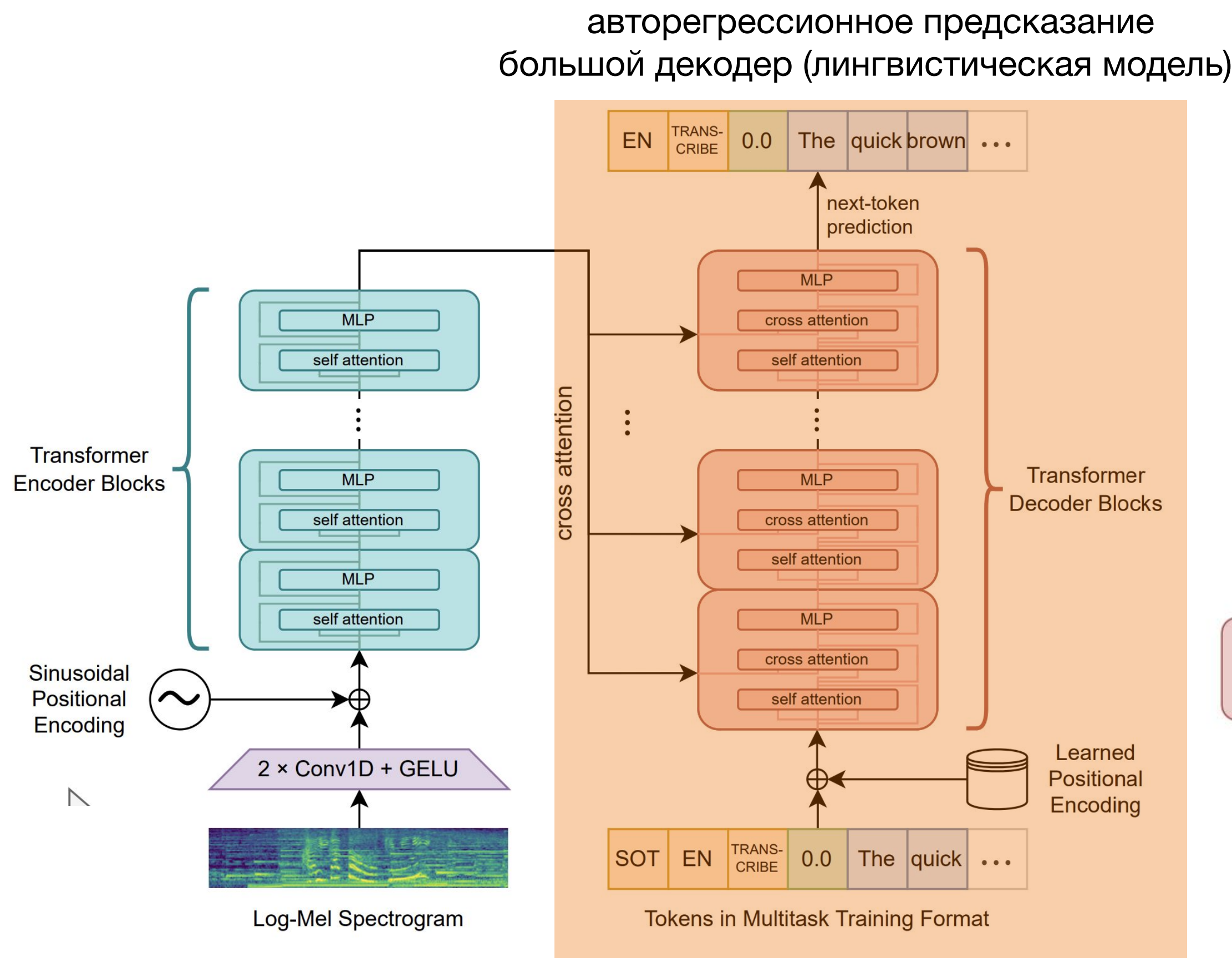
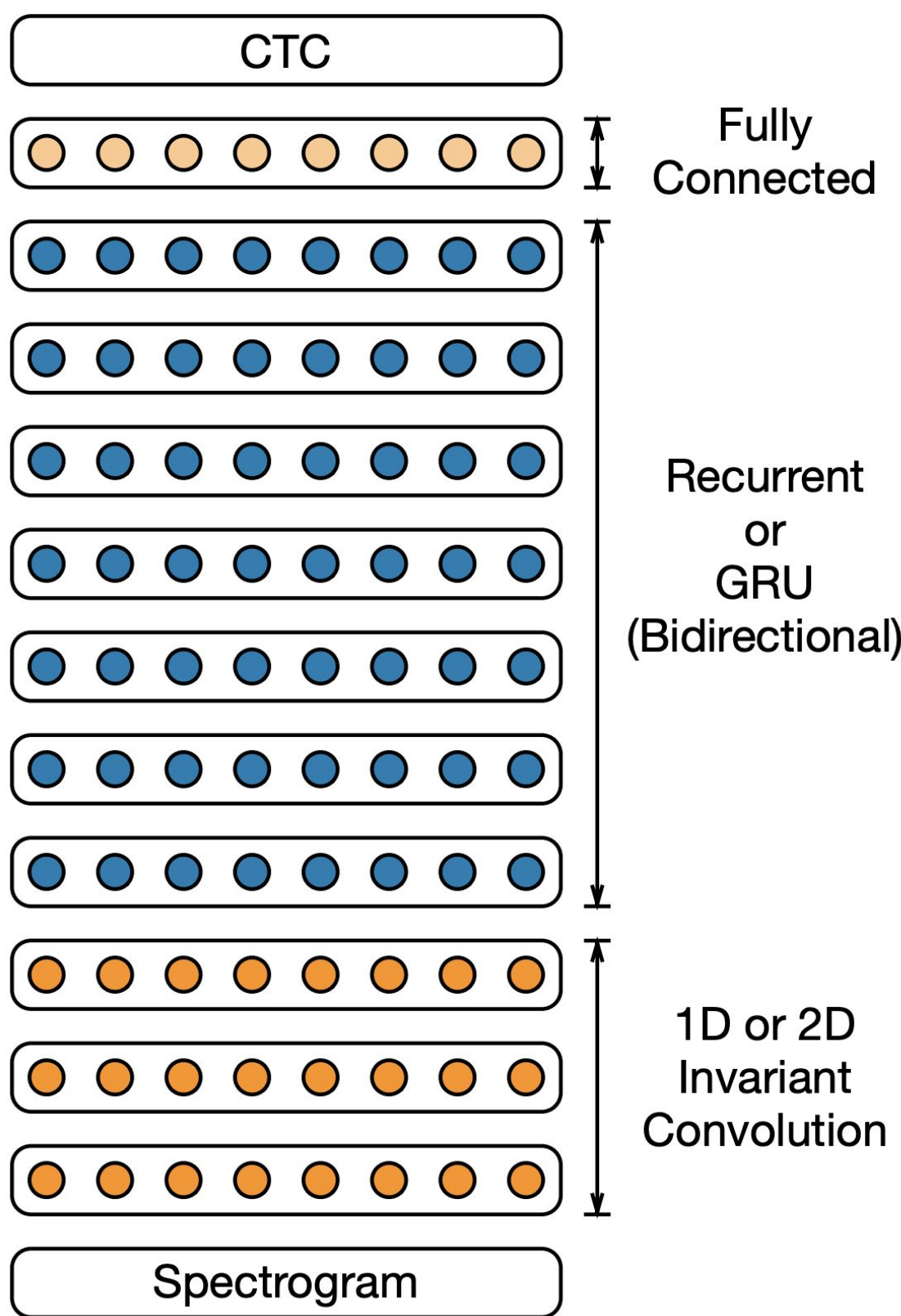
# Вспомним архитектуры

не авторегрессионное предсказание



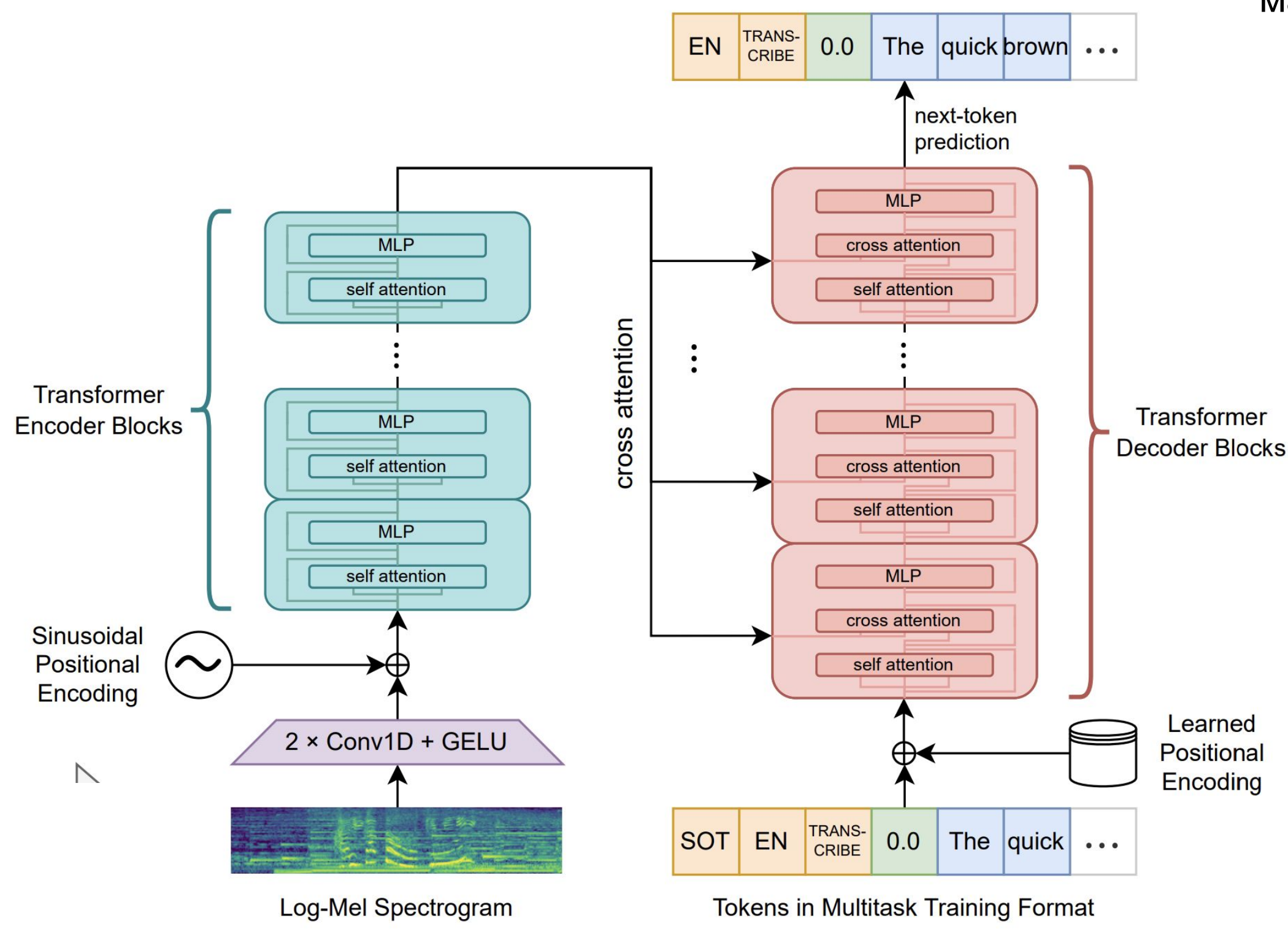
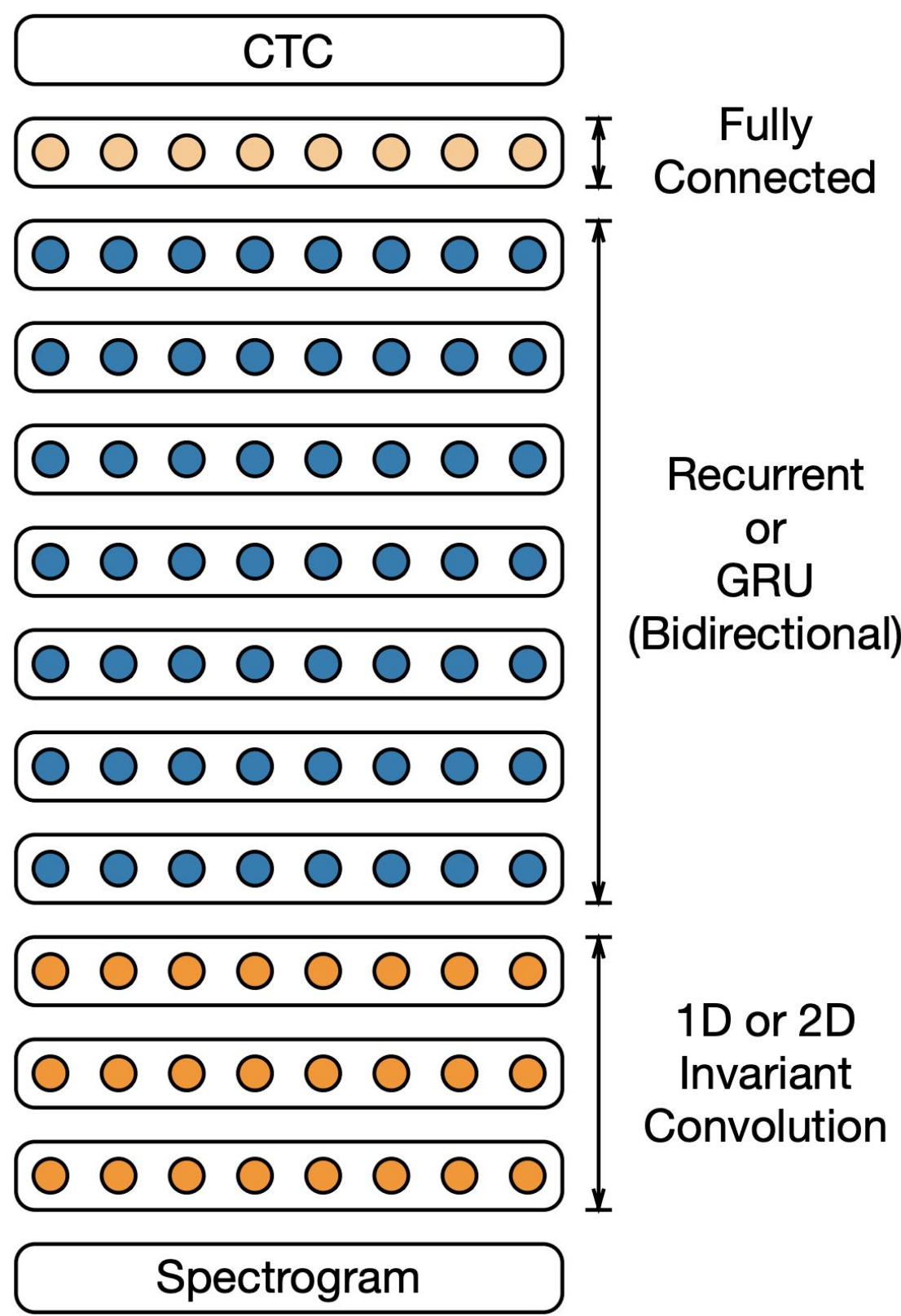


# Вспомним архитектуры

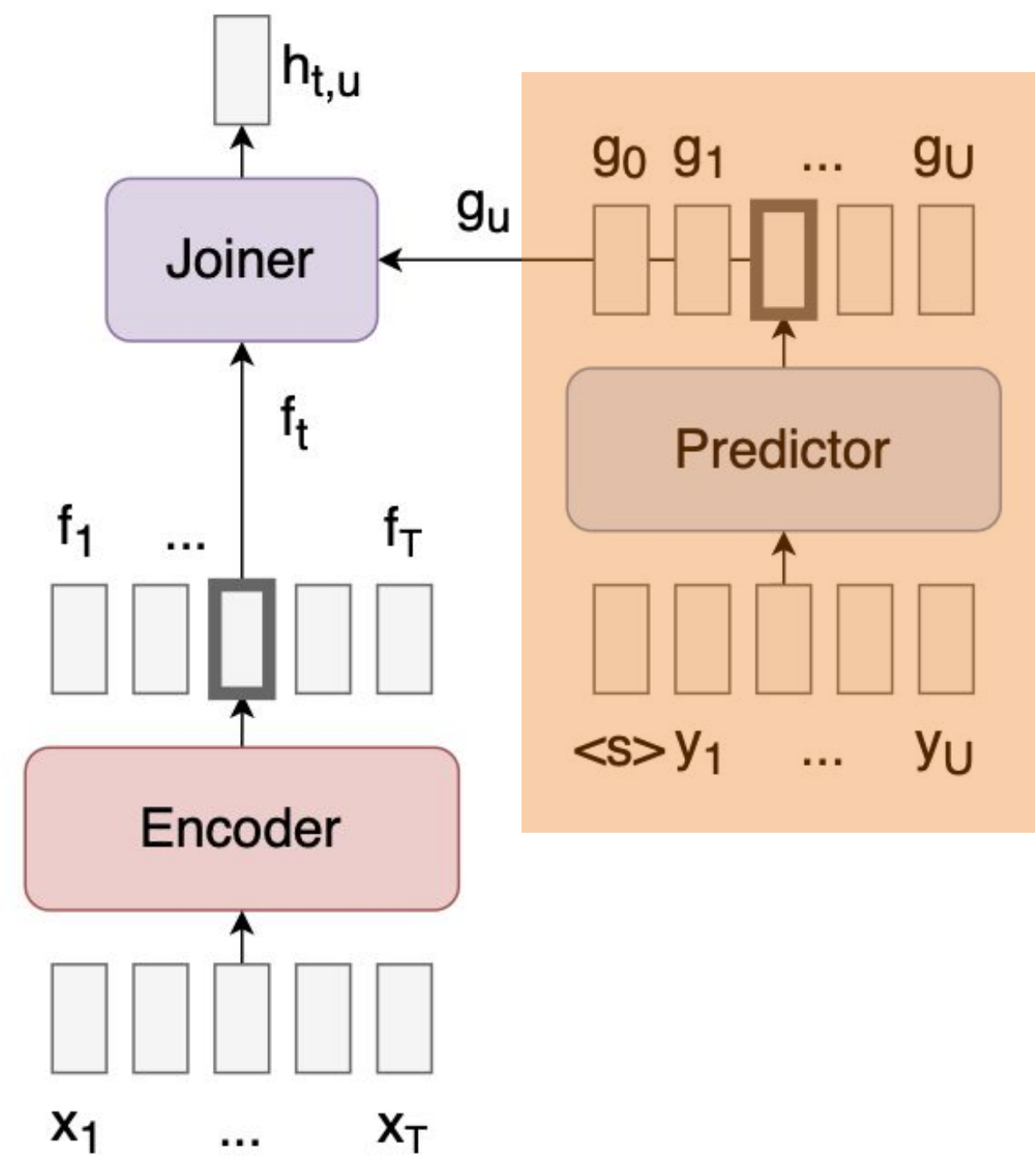




# Вспомним архитектуры

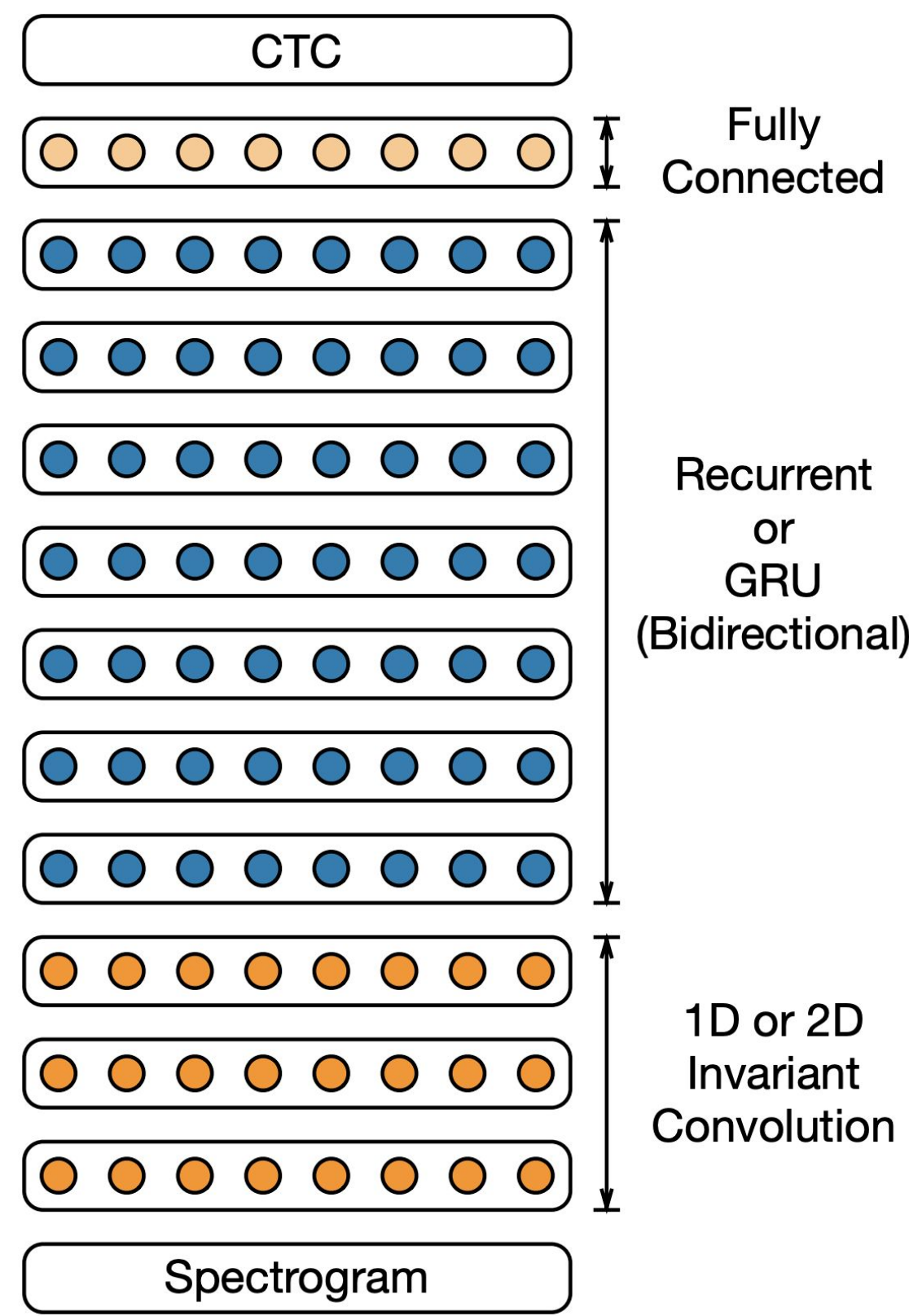


авторегрессионное предсказание  
маленький декодер (лингвистическая модель)

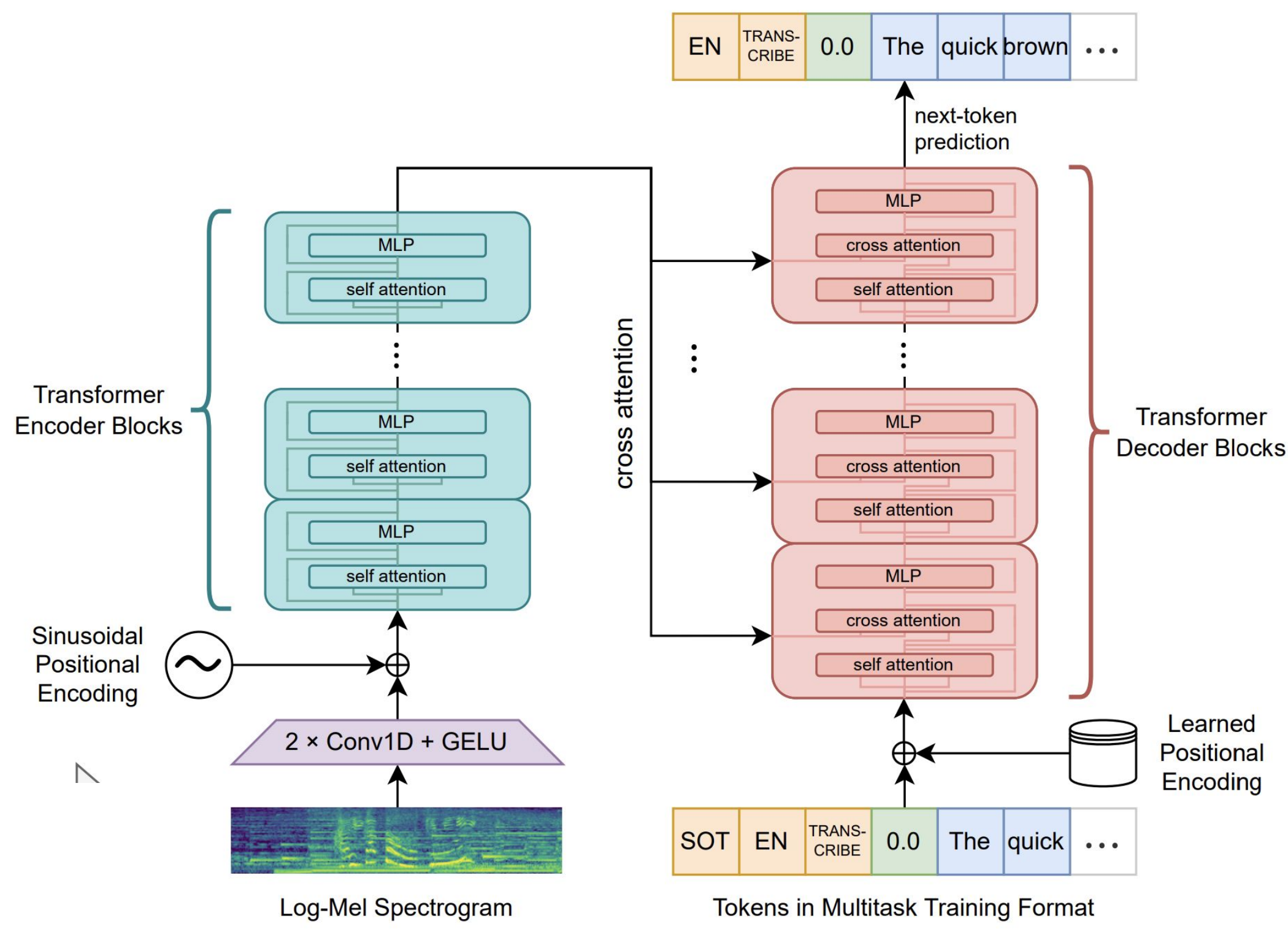


# Вспомним архитектуры

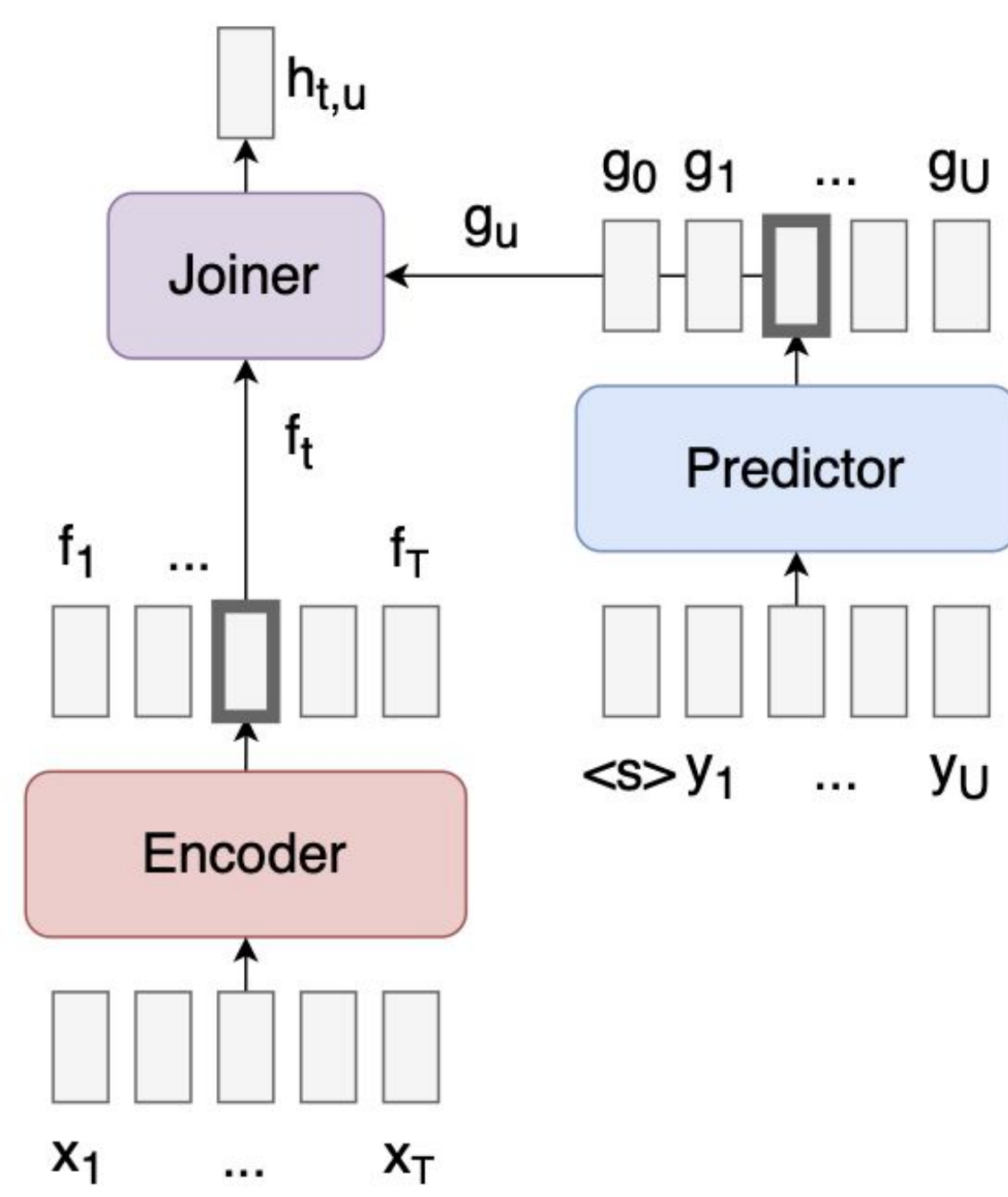
поиск лучшего выравнивания  
между фреймами и токенами



предсказание следующего токена

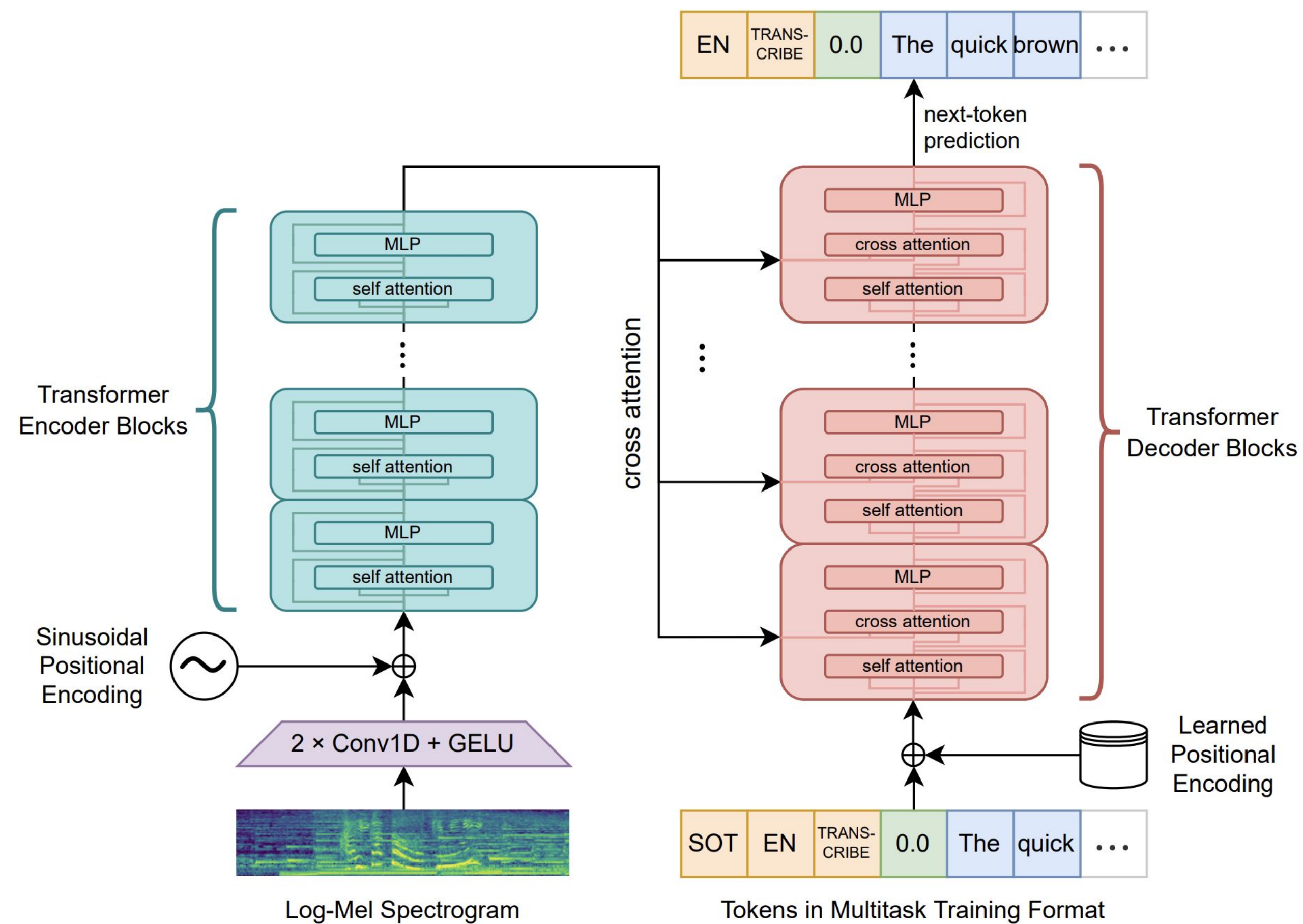


поиск лучшего выравнивания  
между фреймами и токенами





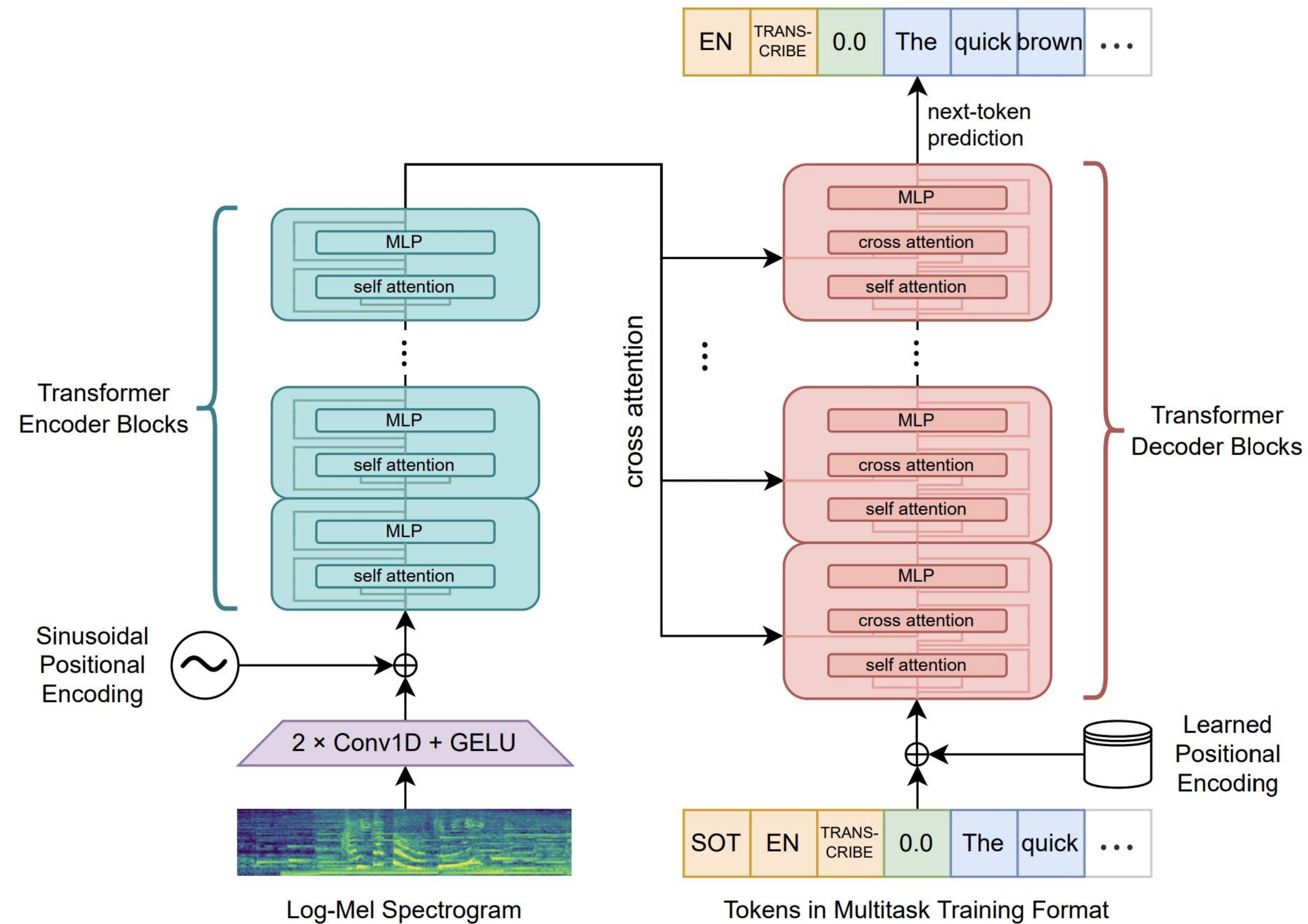
# Типы вне-доменных данных





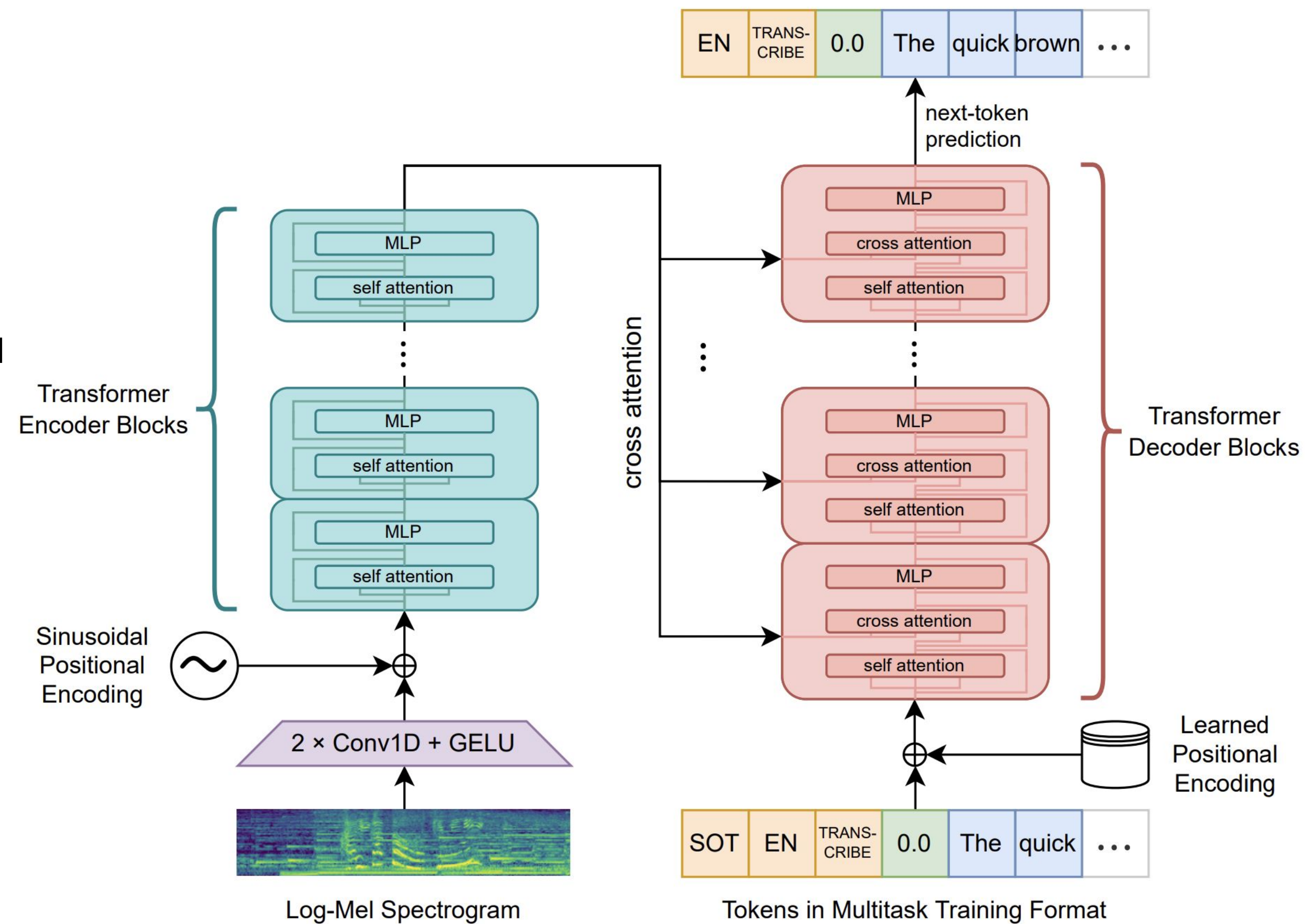
# Типы вне-доменных данных

- Различные шумы



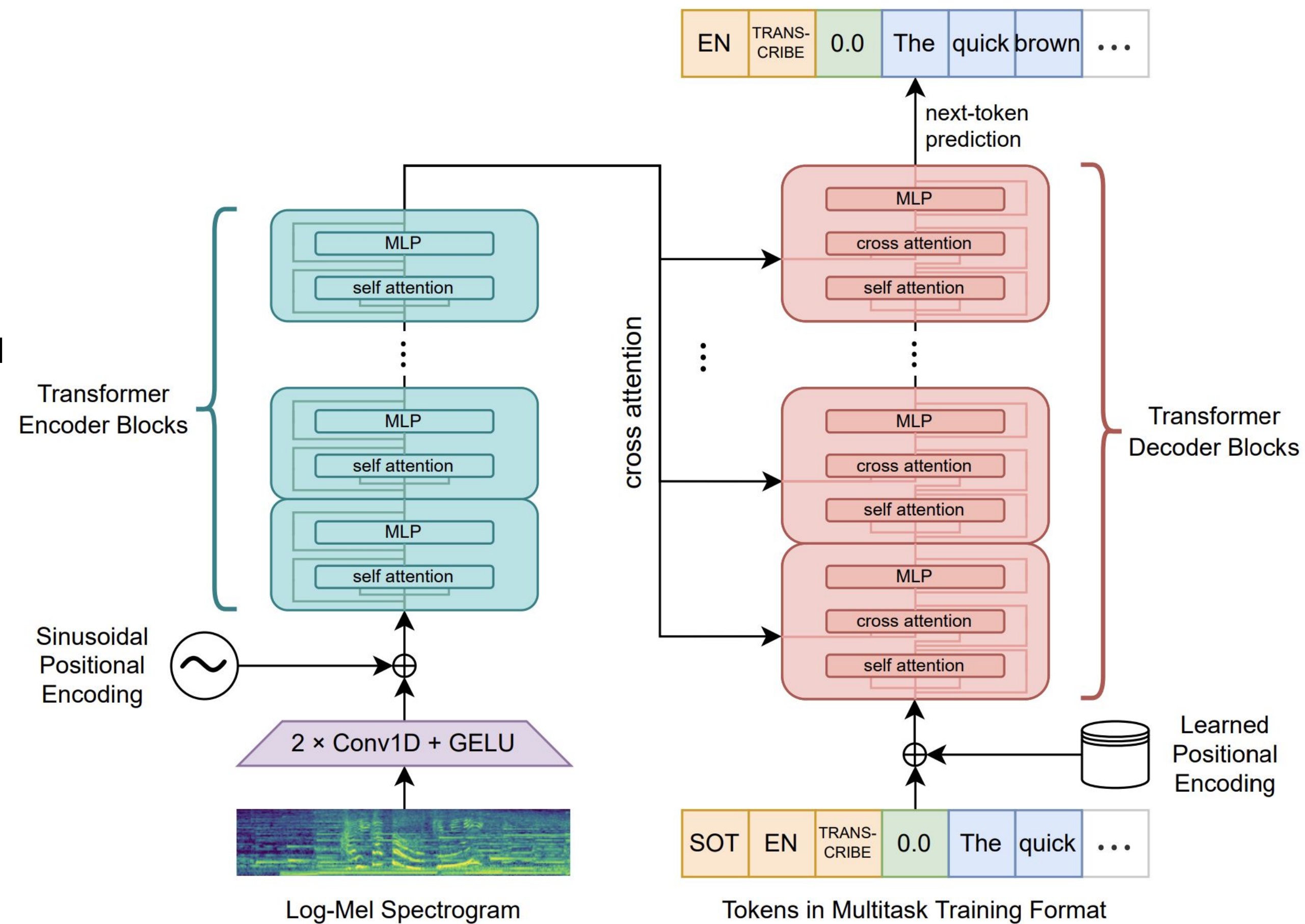
# Типы вне-доменных данных

- Различные шумы
- Новые звуковые эвенты



# Типы вне-доменных данных

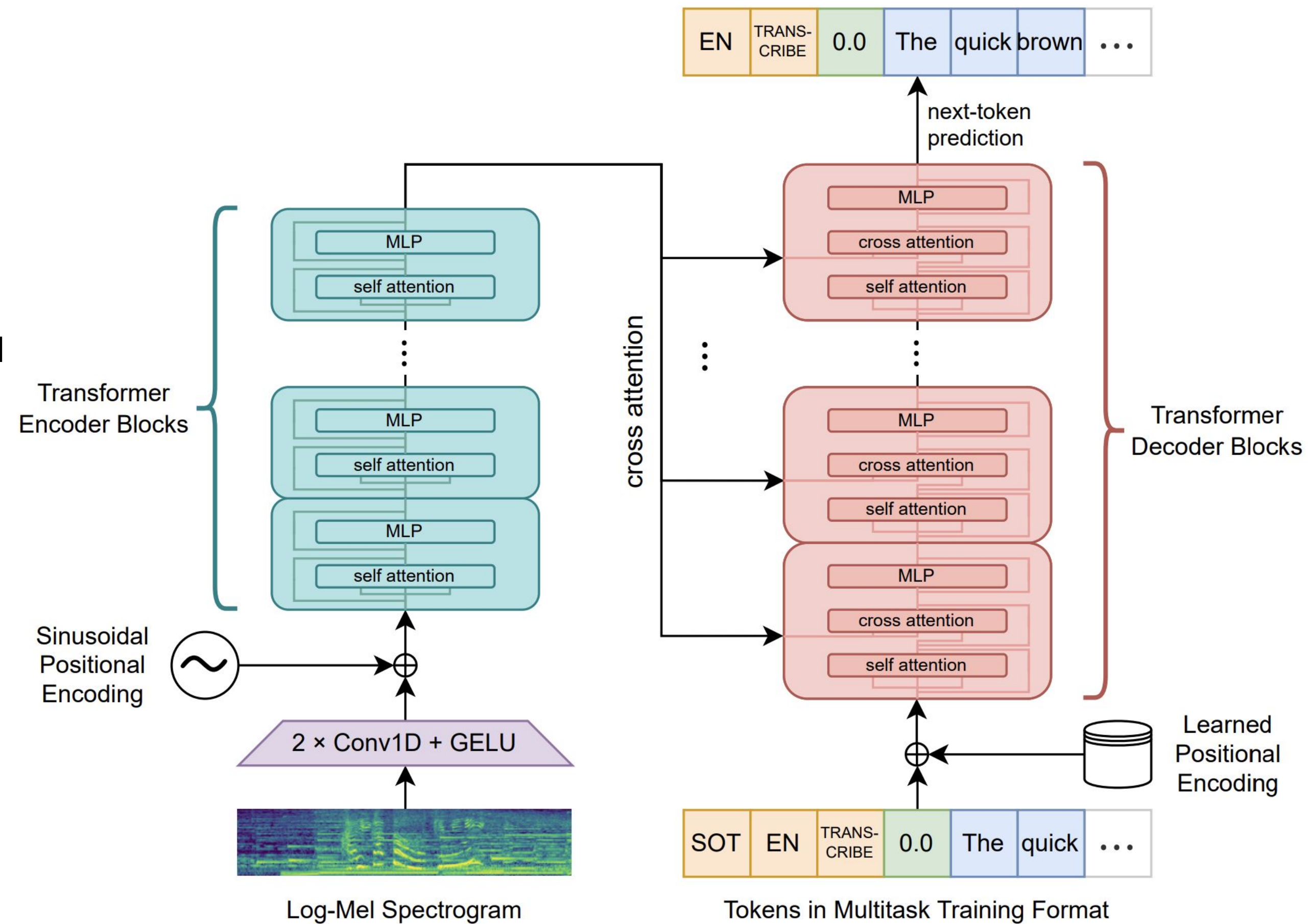
- Различные шумы
- Новые звуковые эвенты
- Эхо





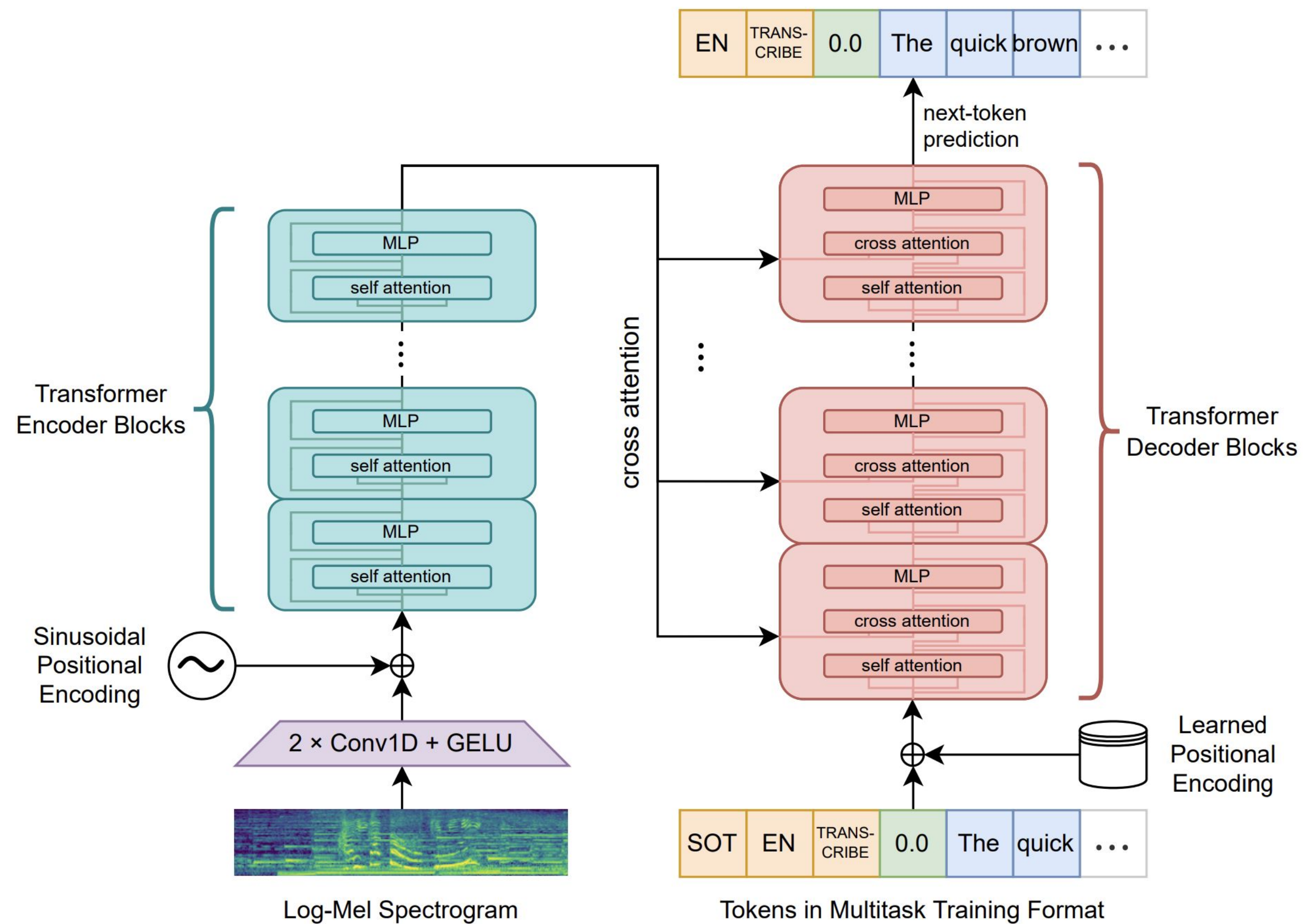
# Типы вне-доменных данных

- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент



# Типы вне-доменных данных

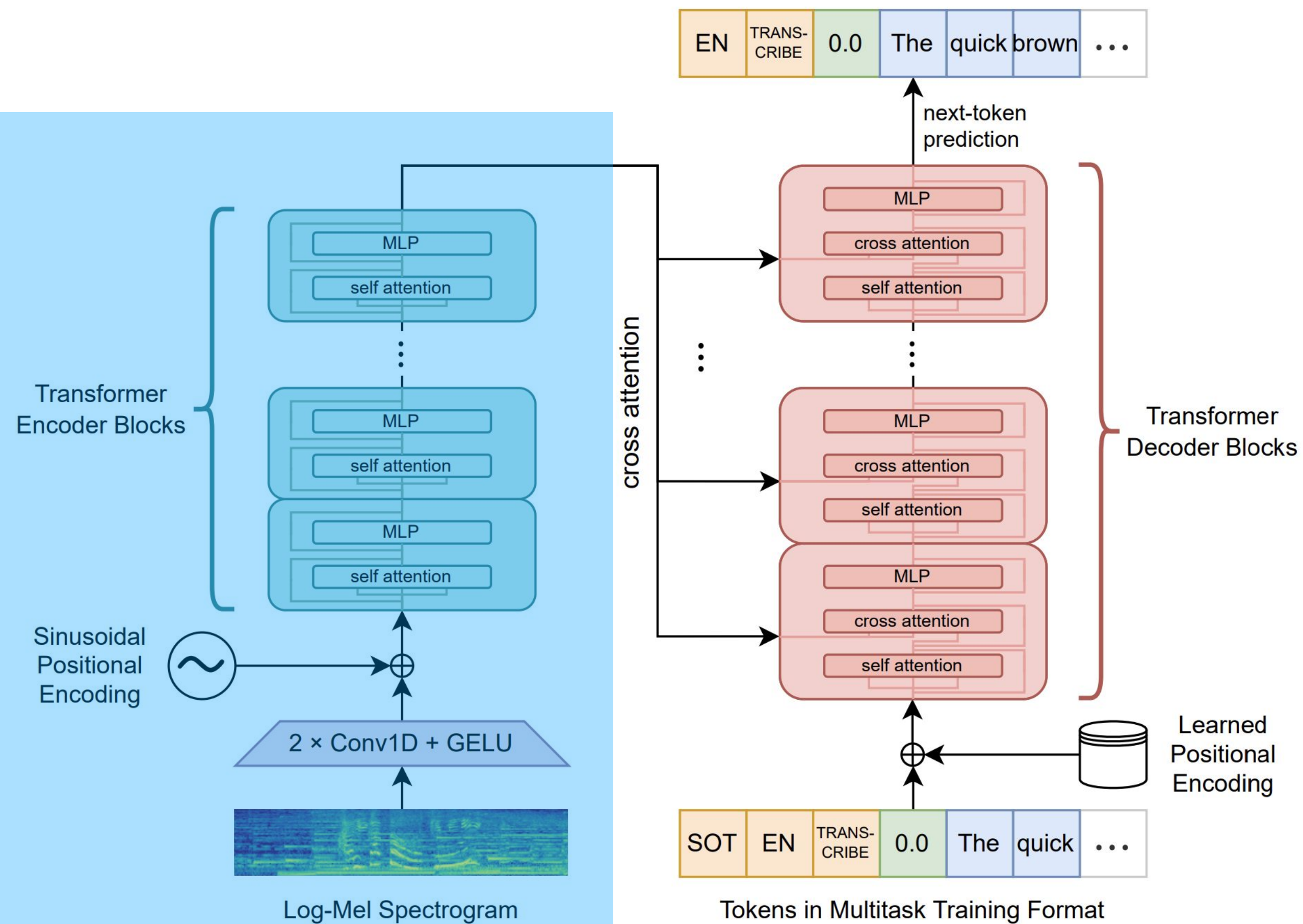
- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи





# Типы вне-доменных данных

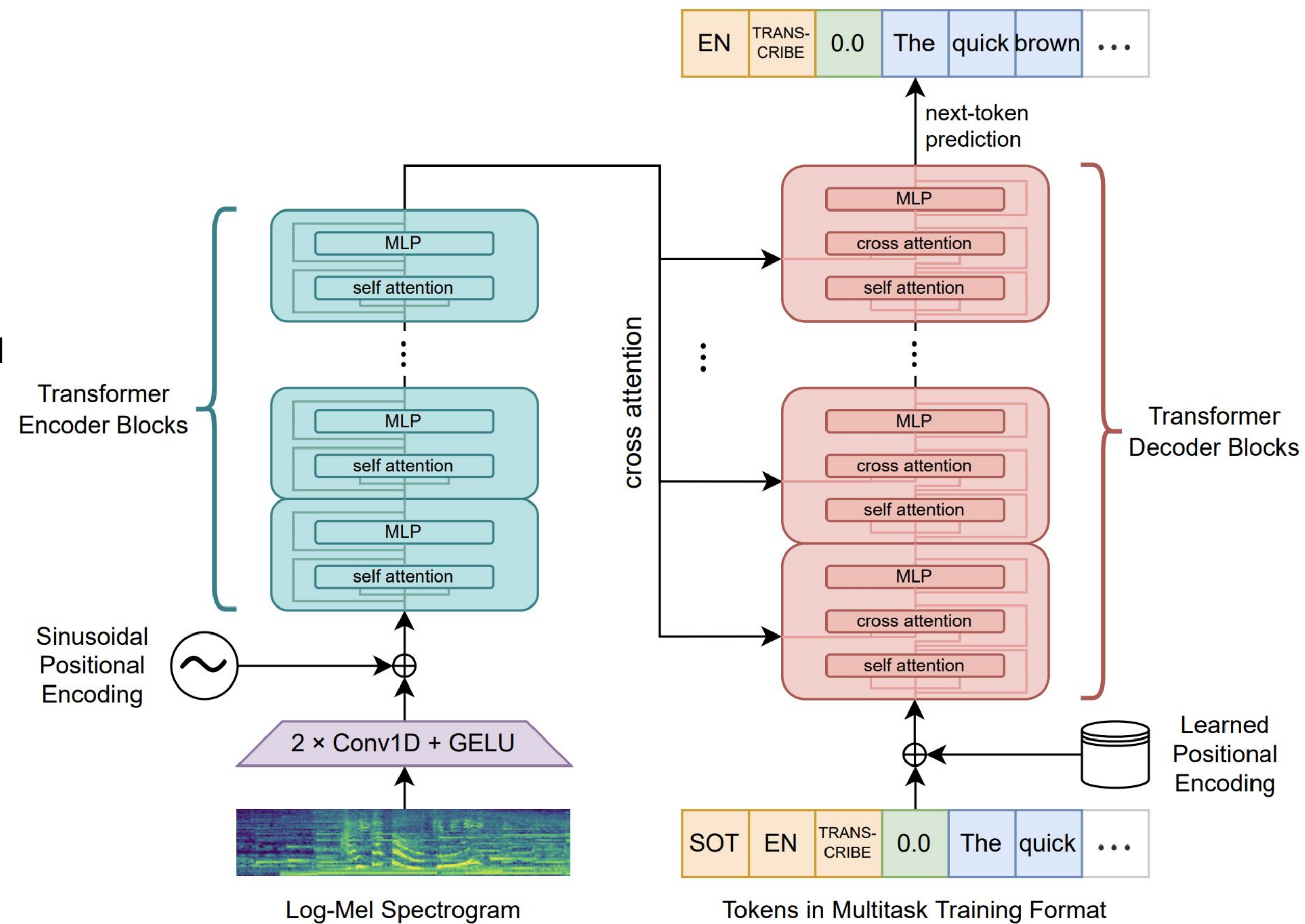
- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи





# Типы вне-доменных данных

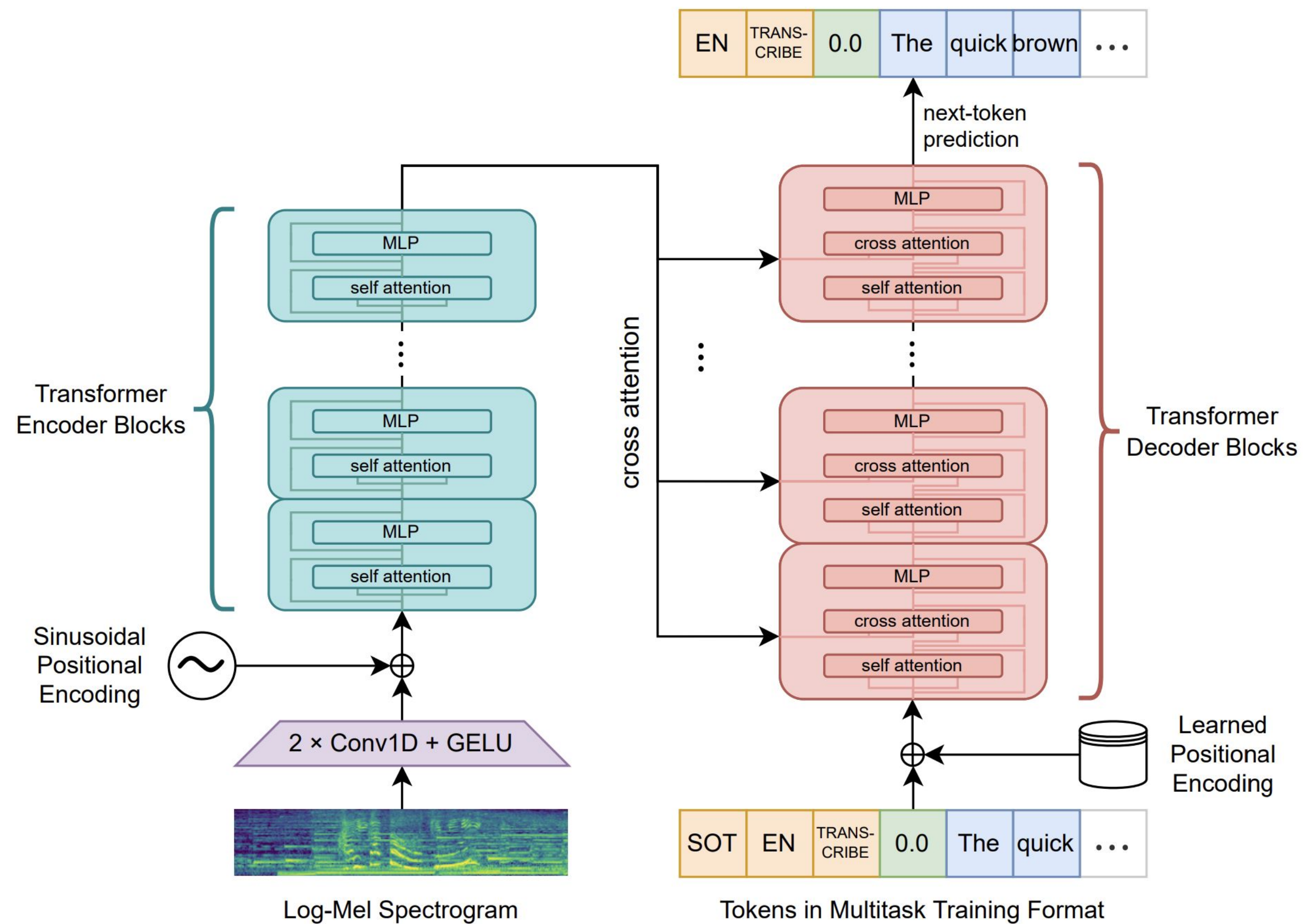
- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи



- Новые слова

# Типы вне-доменных данных

- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи

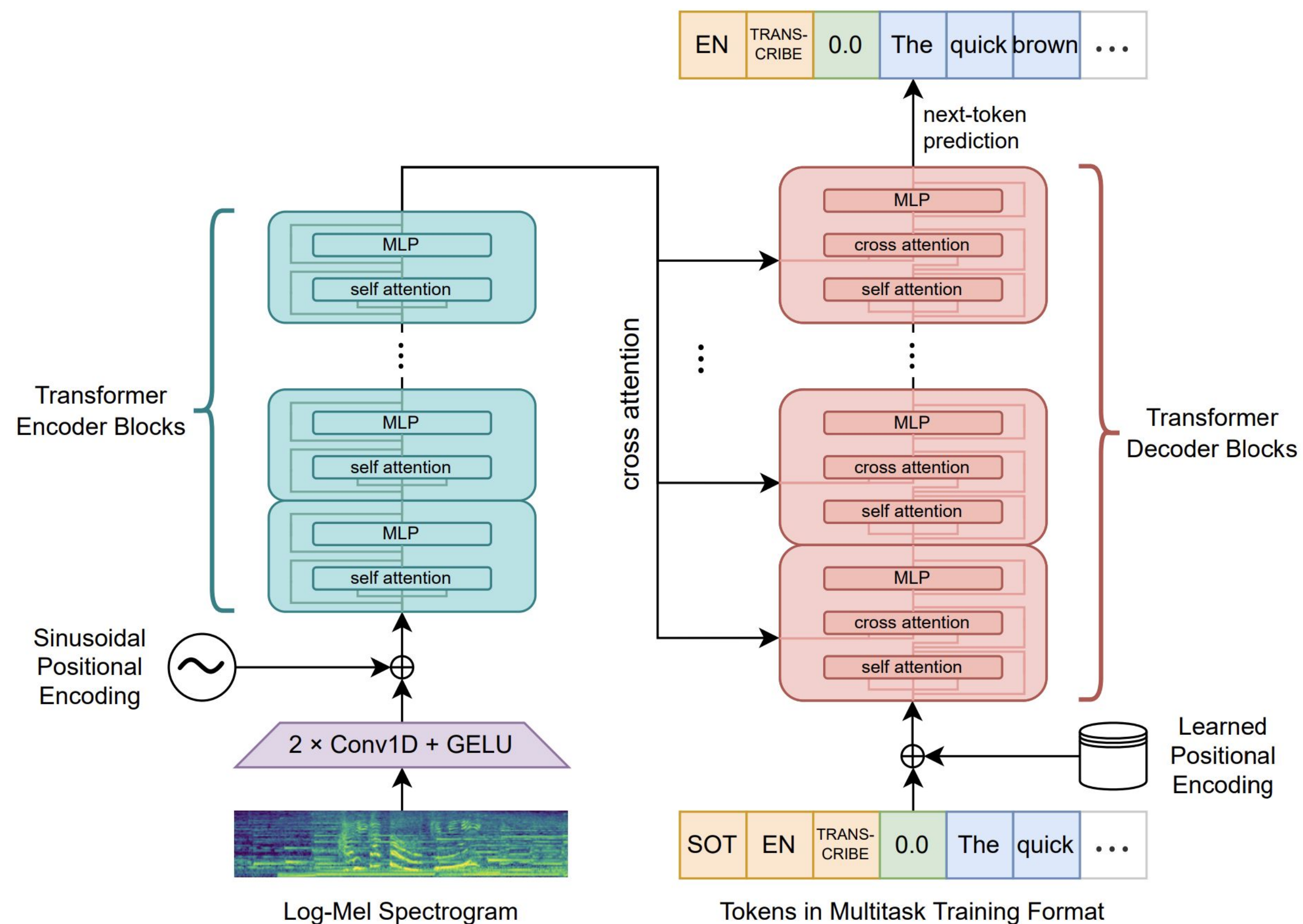


- Новые слова
- Проф. термины



# Типы вне-доменных данных

- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи

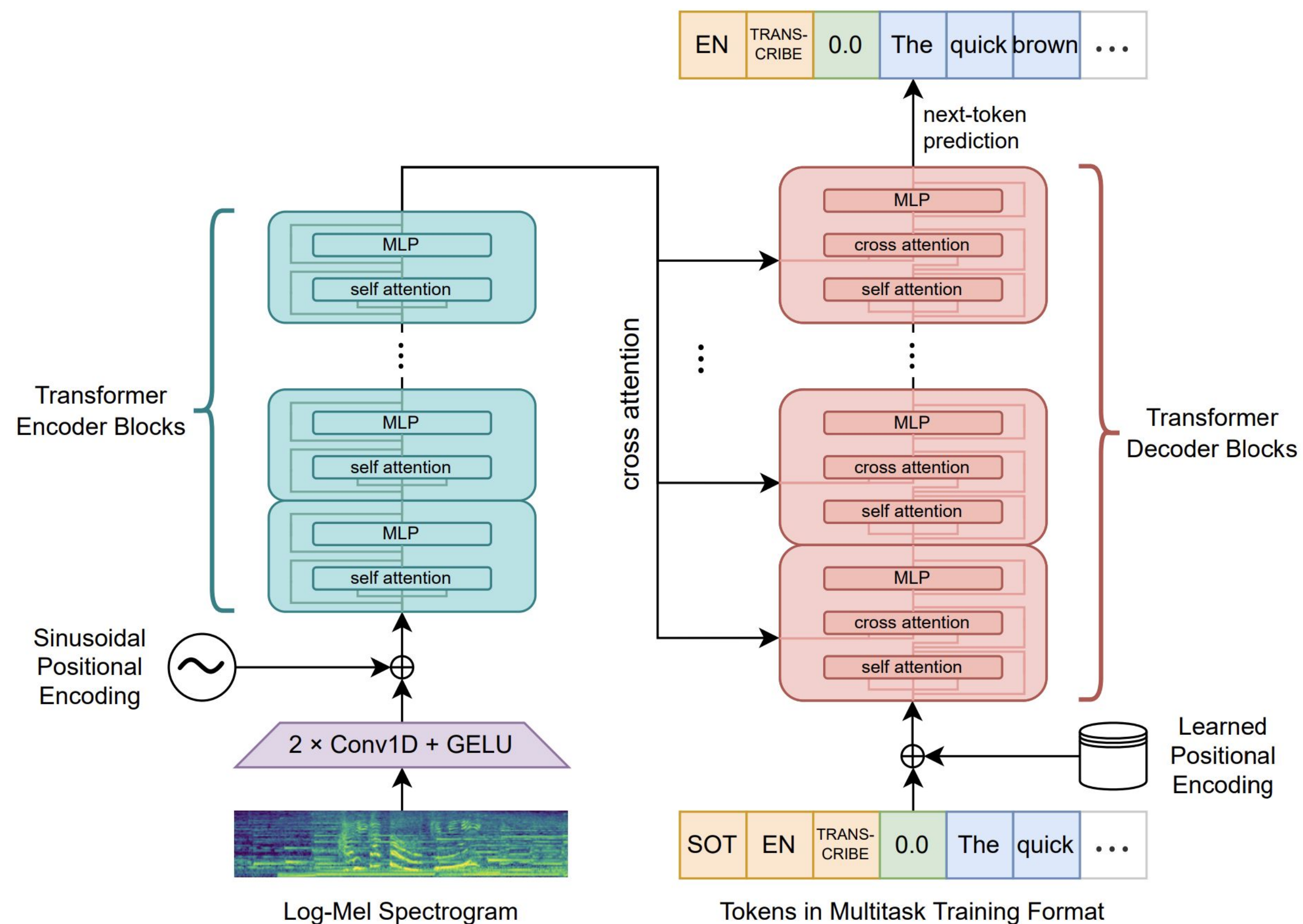


- Новые слова
- Проф. термины
- Сленг



# Типы вне-доменных данных

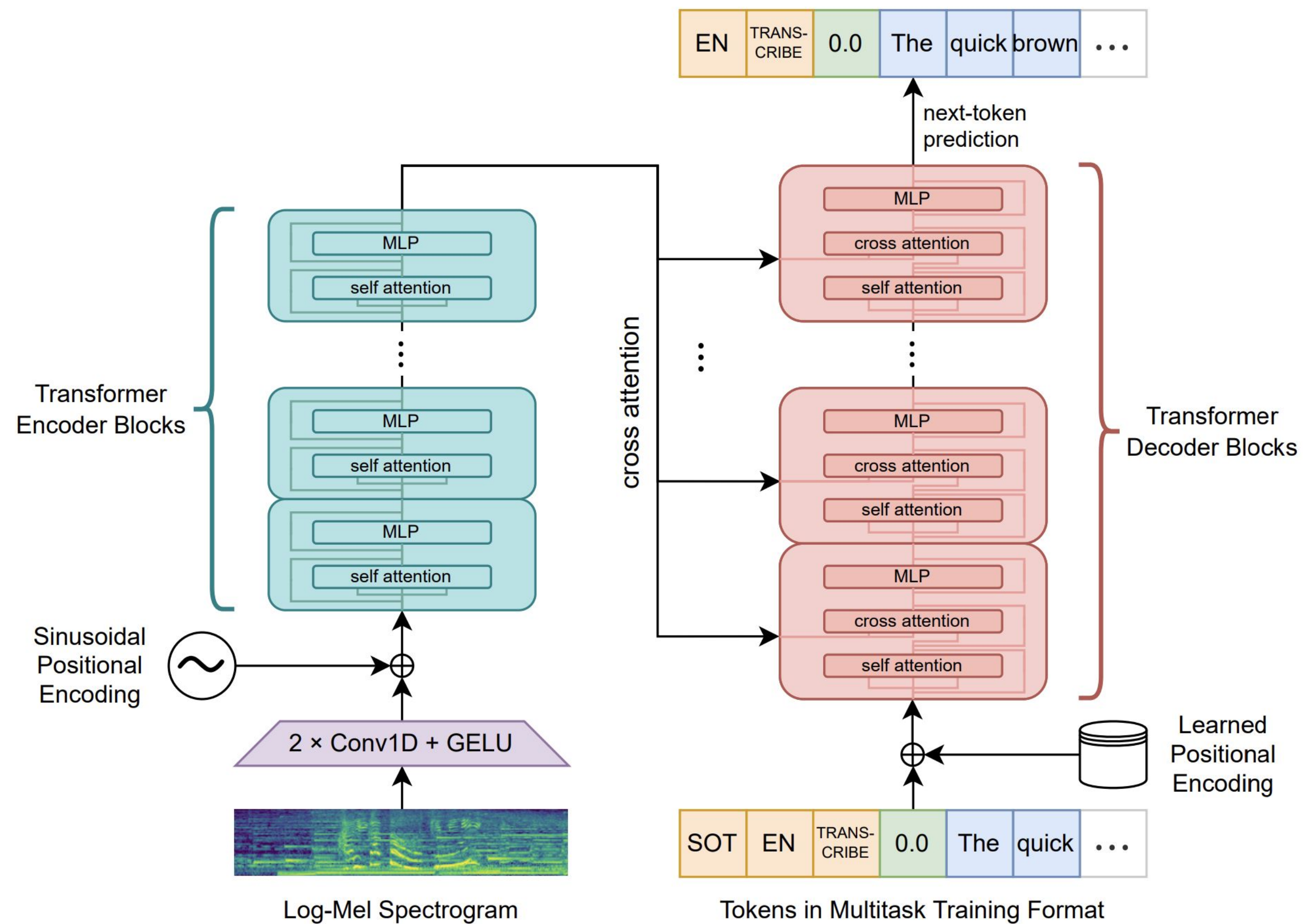
- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи



- Новые слова
- Проф. термины
- Сленг
- Акцент (Диалект)

# Типы вне-доменных данных

- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи

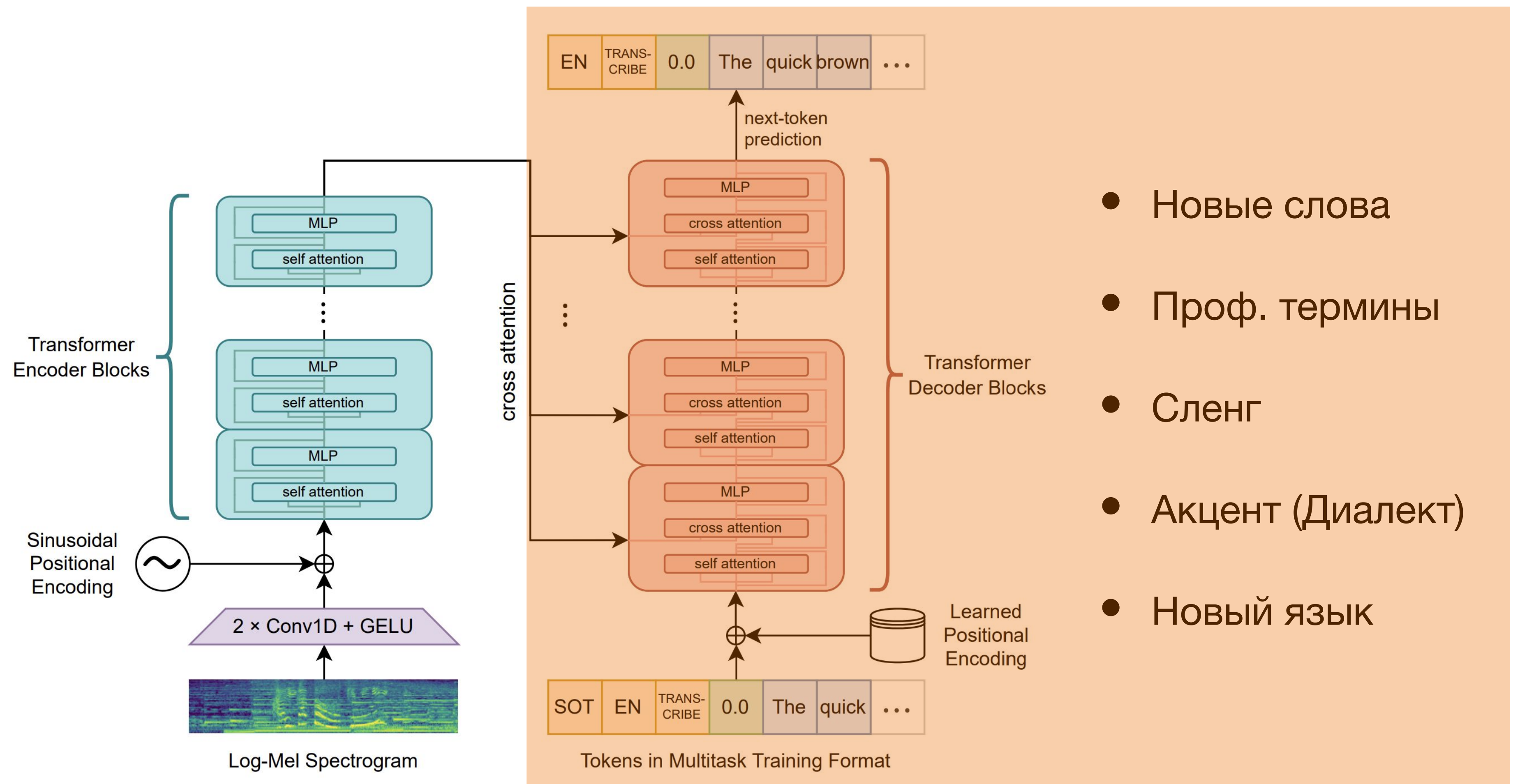


- Новые слова
- Проф. термины
- Сленг
- Акцент (Диалект)
- Новый язык



# Типы вне-доменных данных

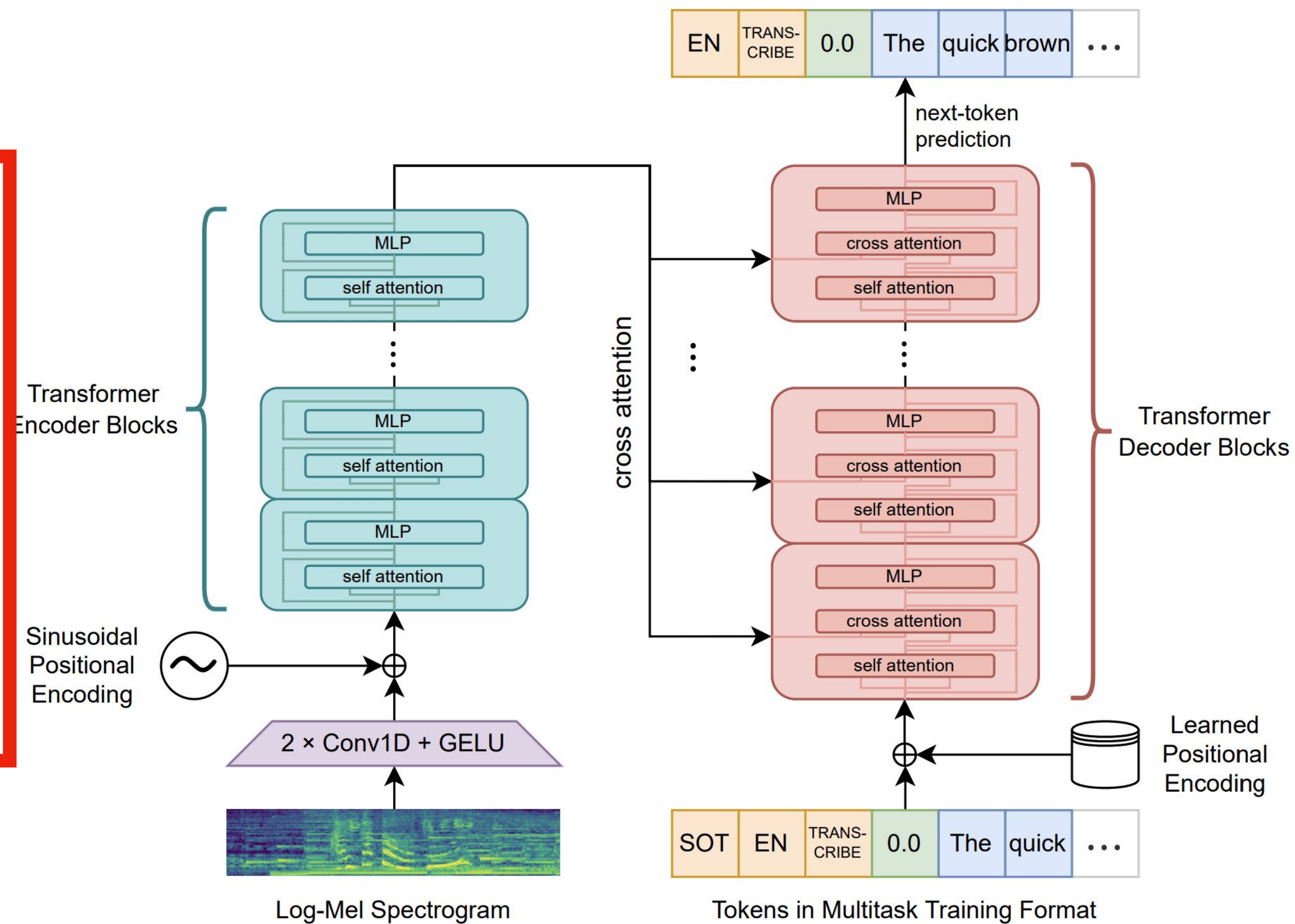
- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи





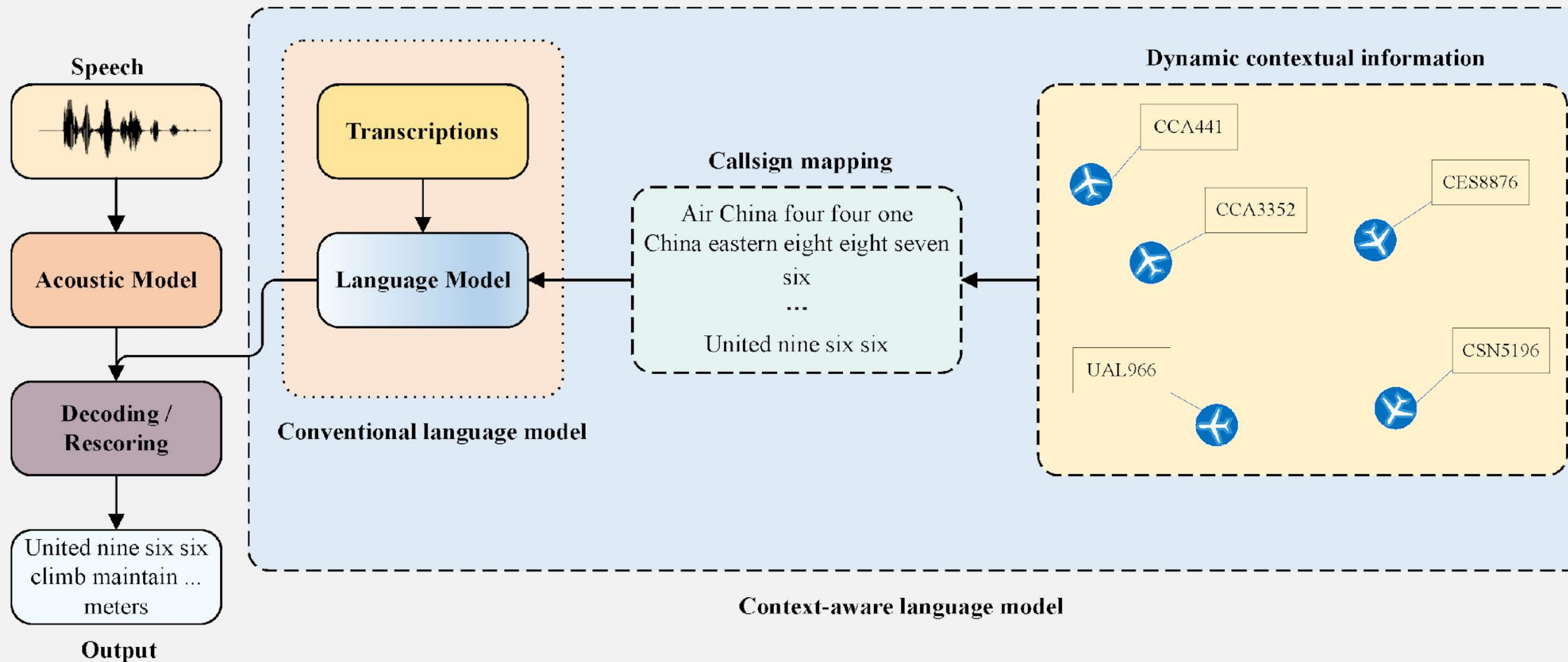
# Типы вне-доменных данных

- Различные шумы
- Новые звуковые эвенты
- Эхо
- Акцент
- Громкость речи



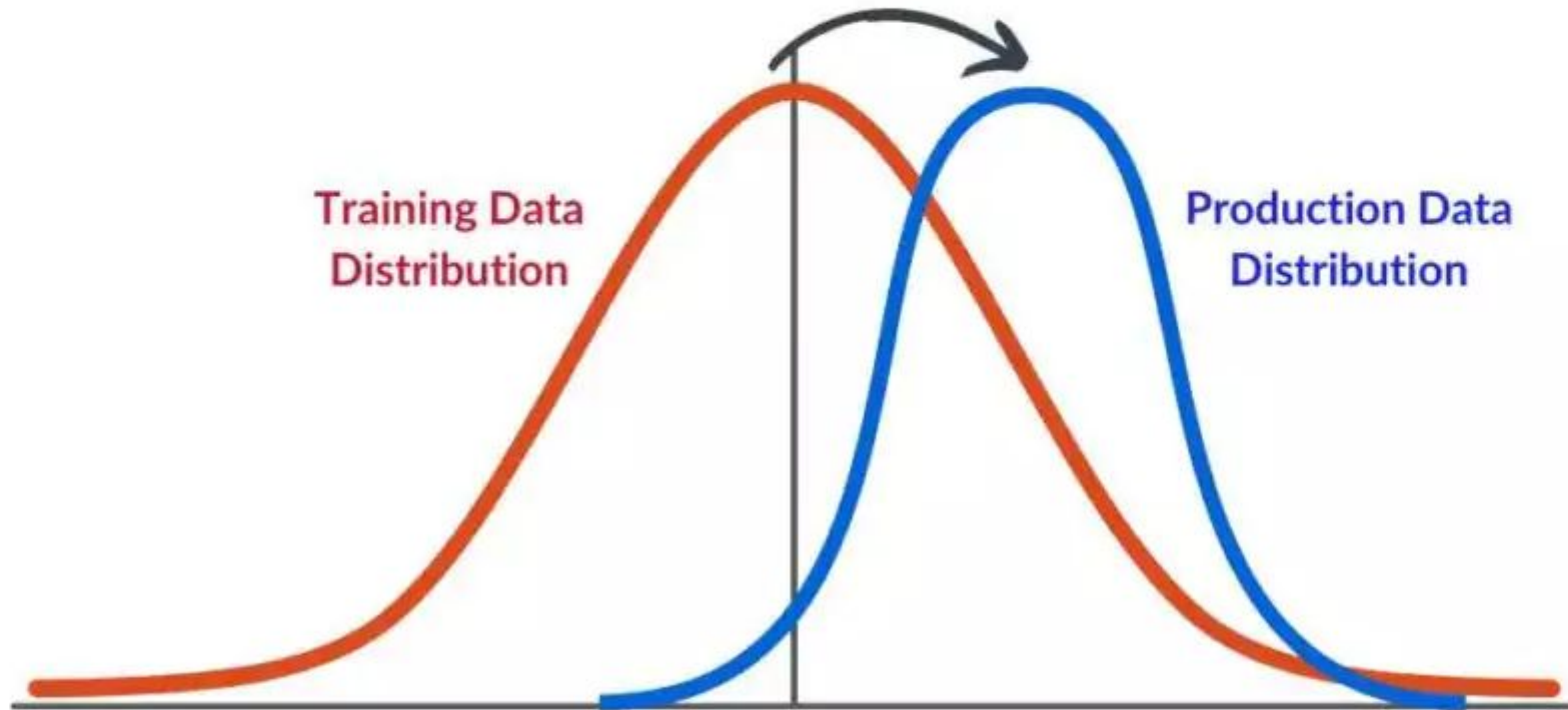
- Новые слова
- Проф. термины
- Сленг
- Акцент (Диалект)
- Новый язык

# Типы вне-доменных данных





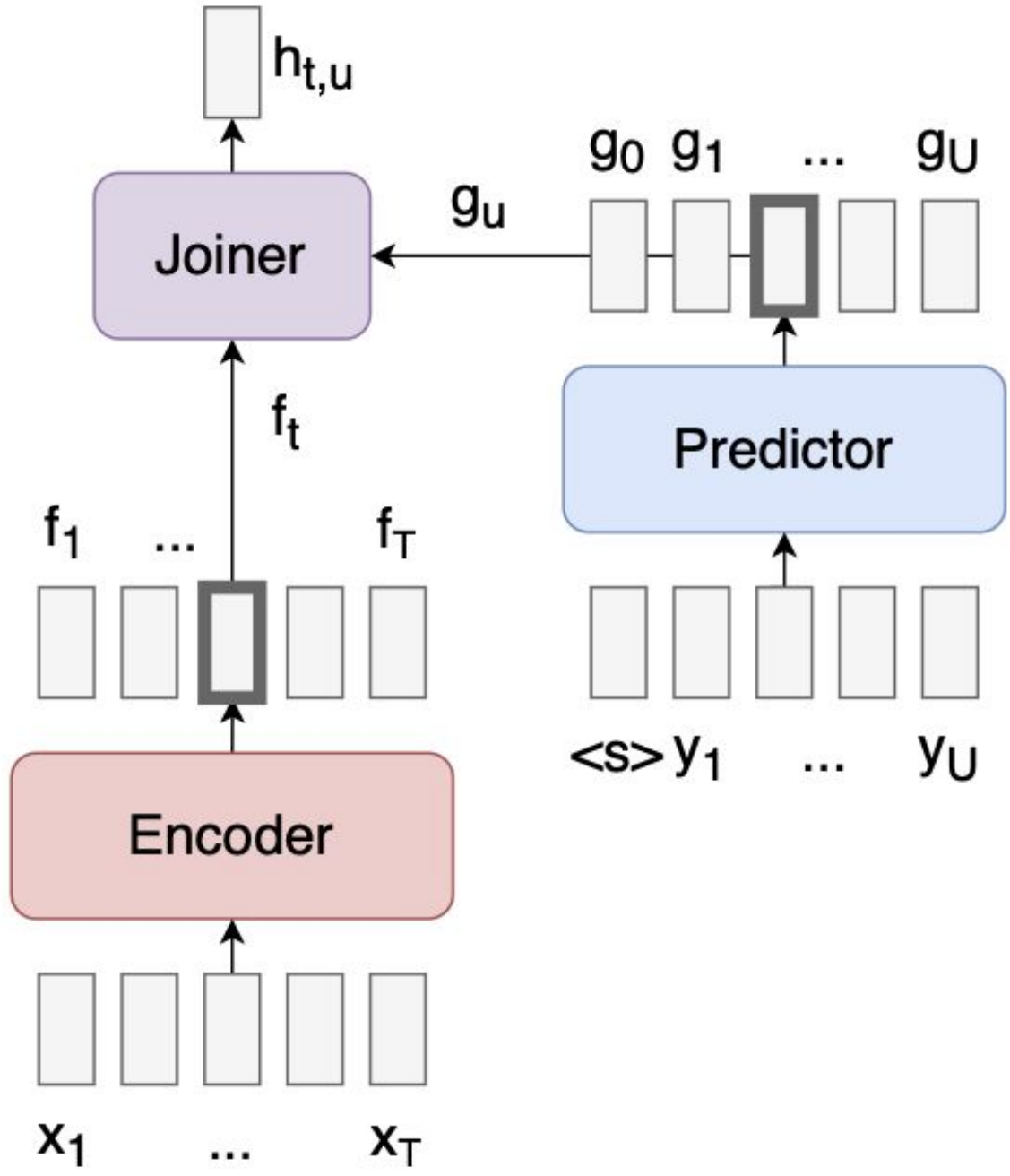
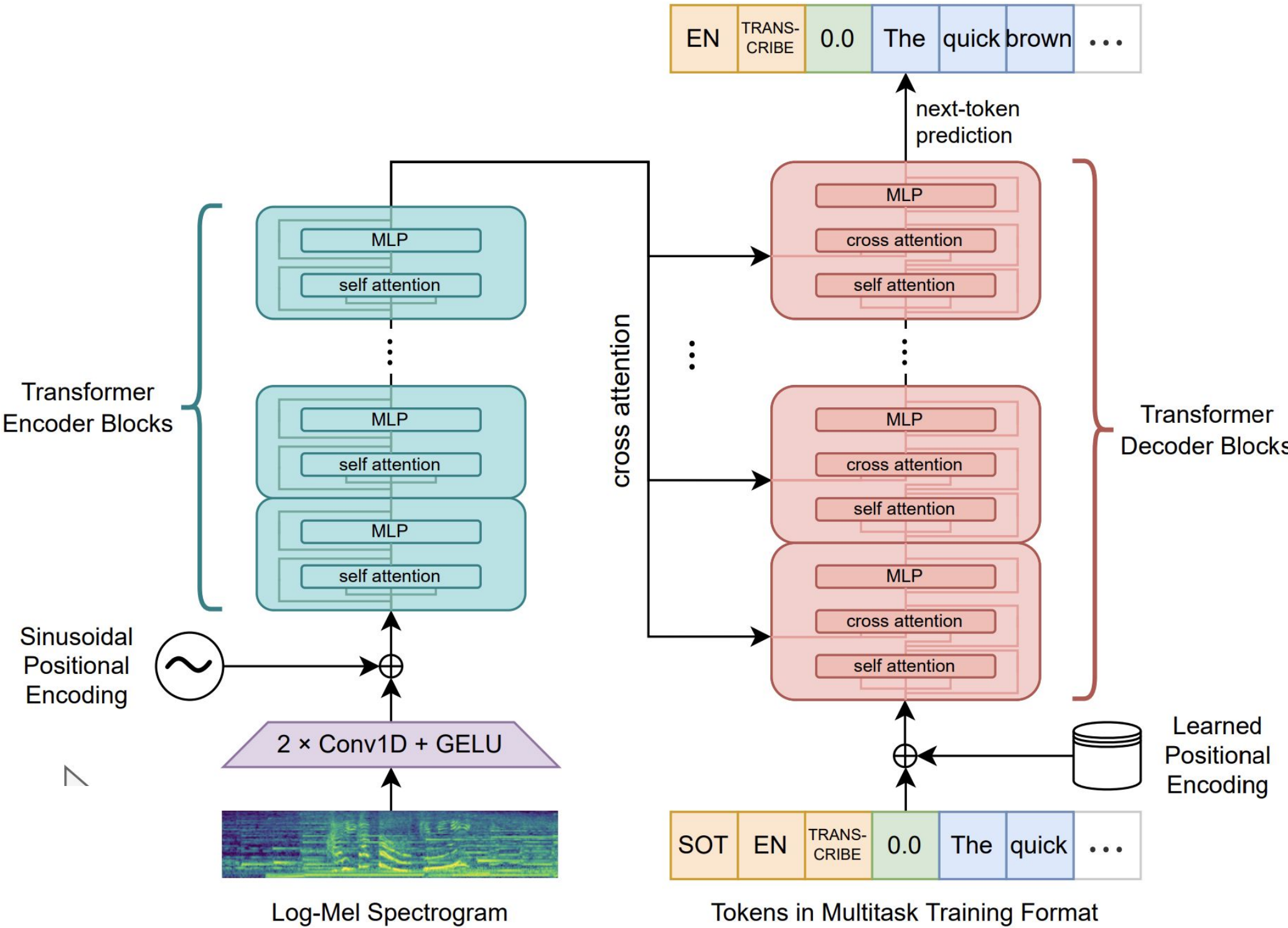
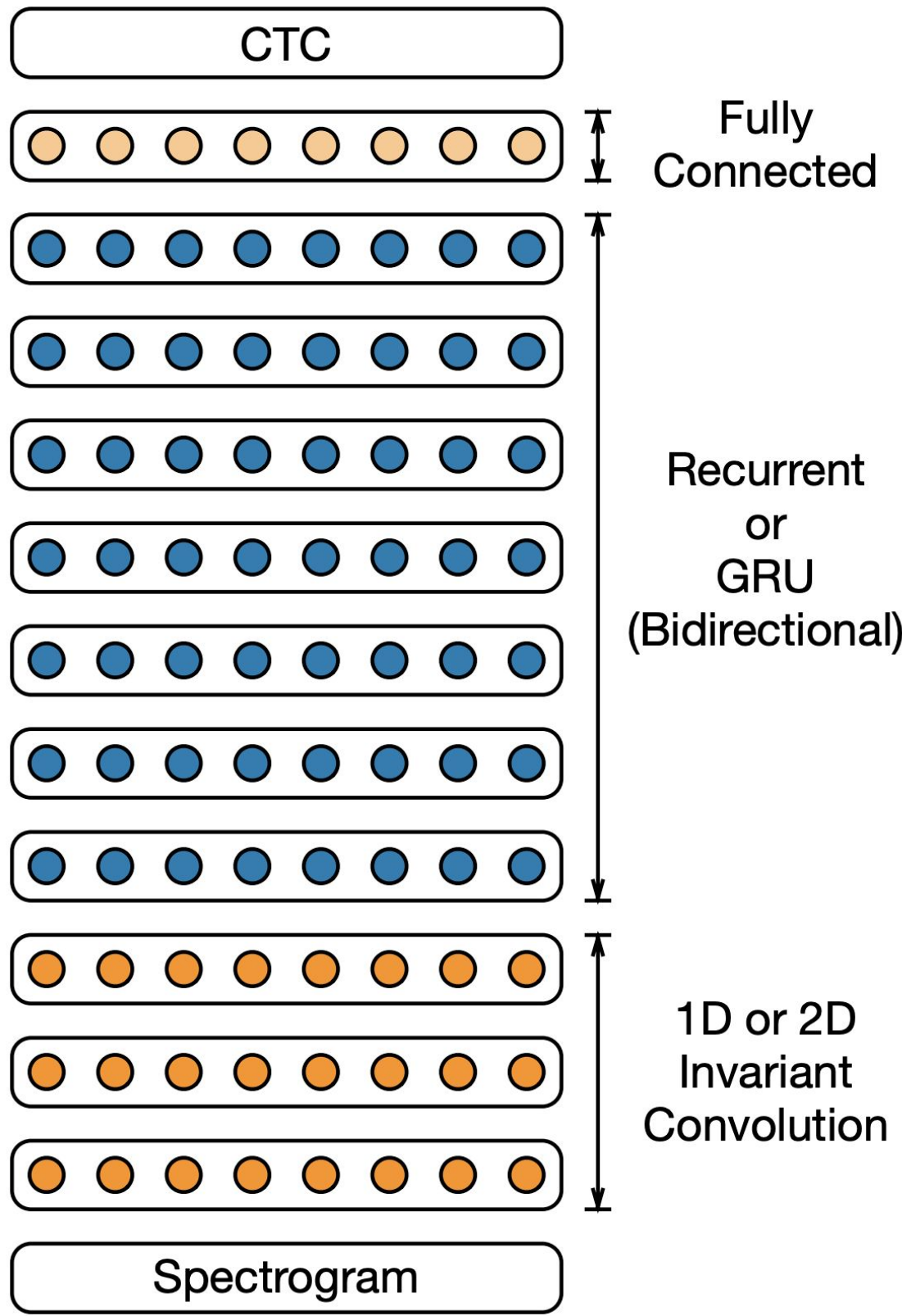
# Типы вне-доменных данных



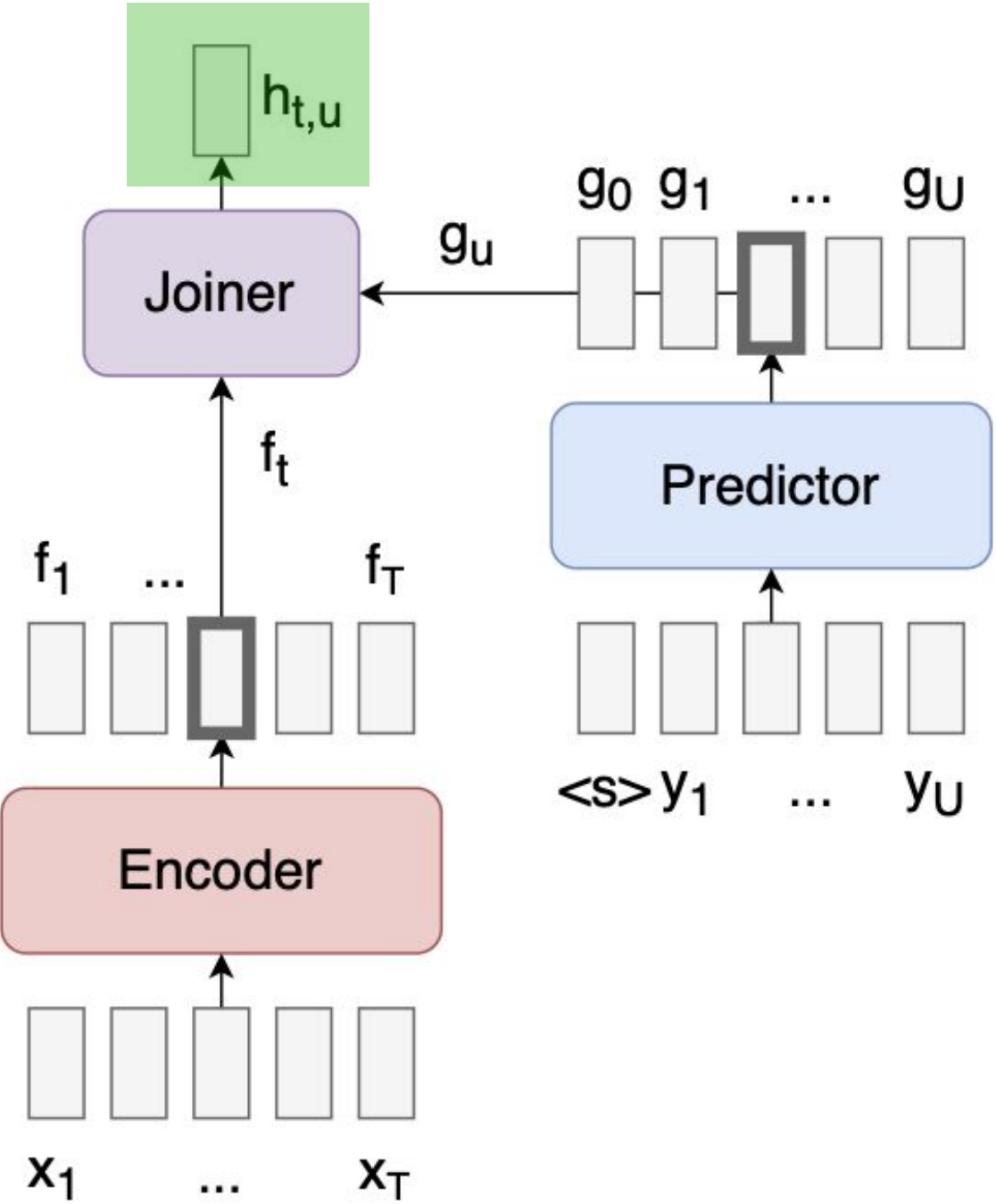
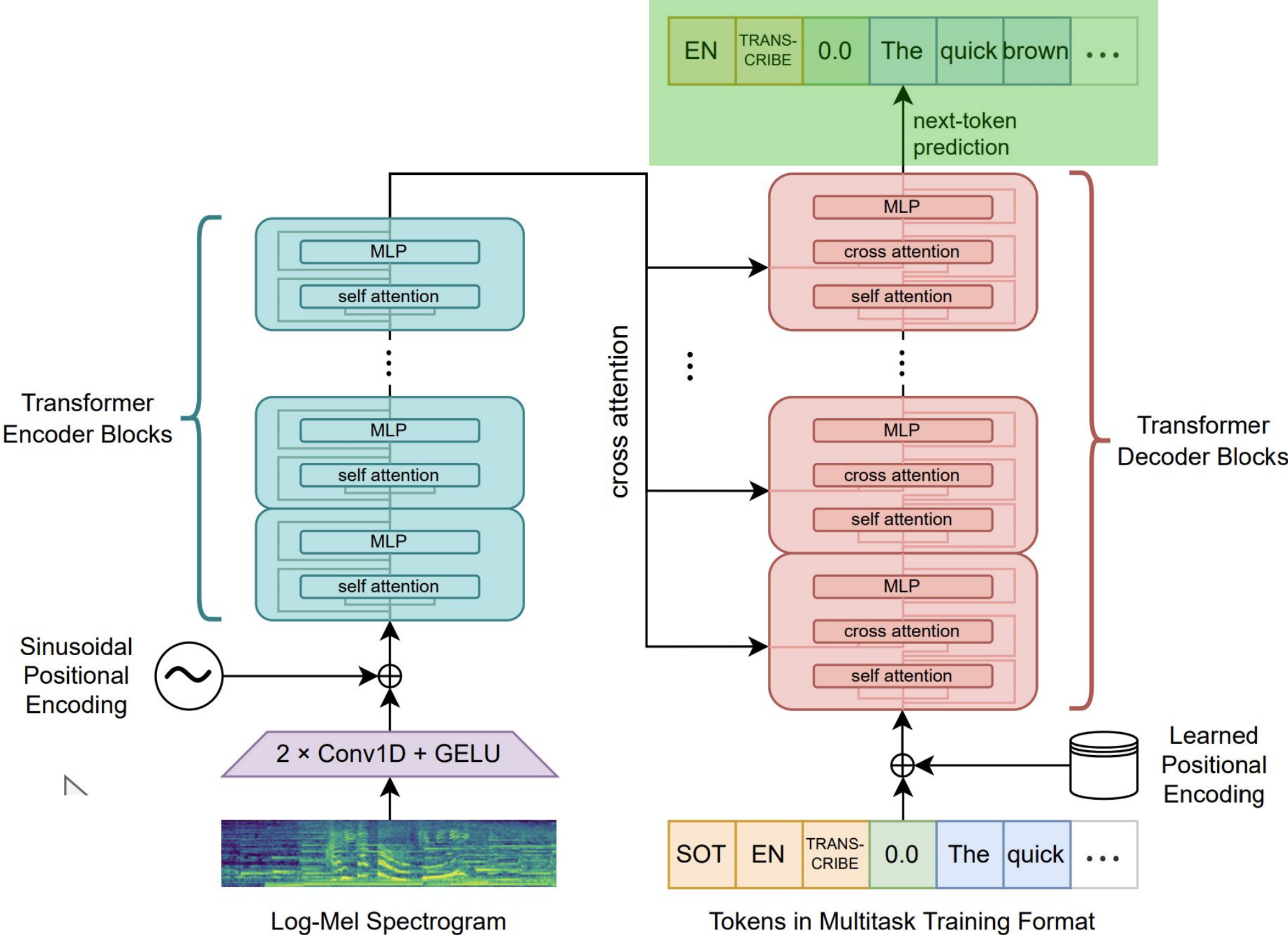
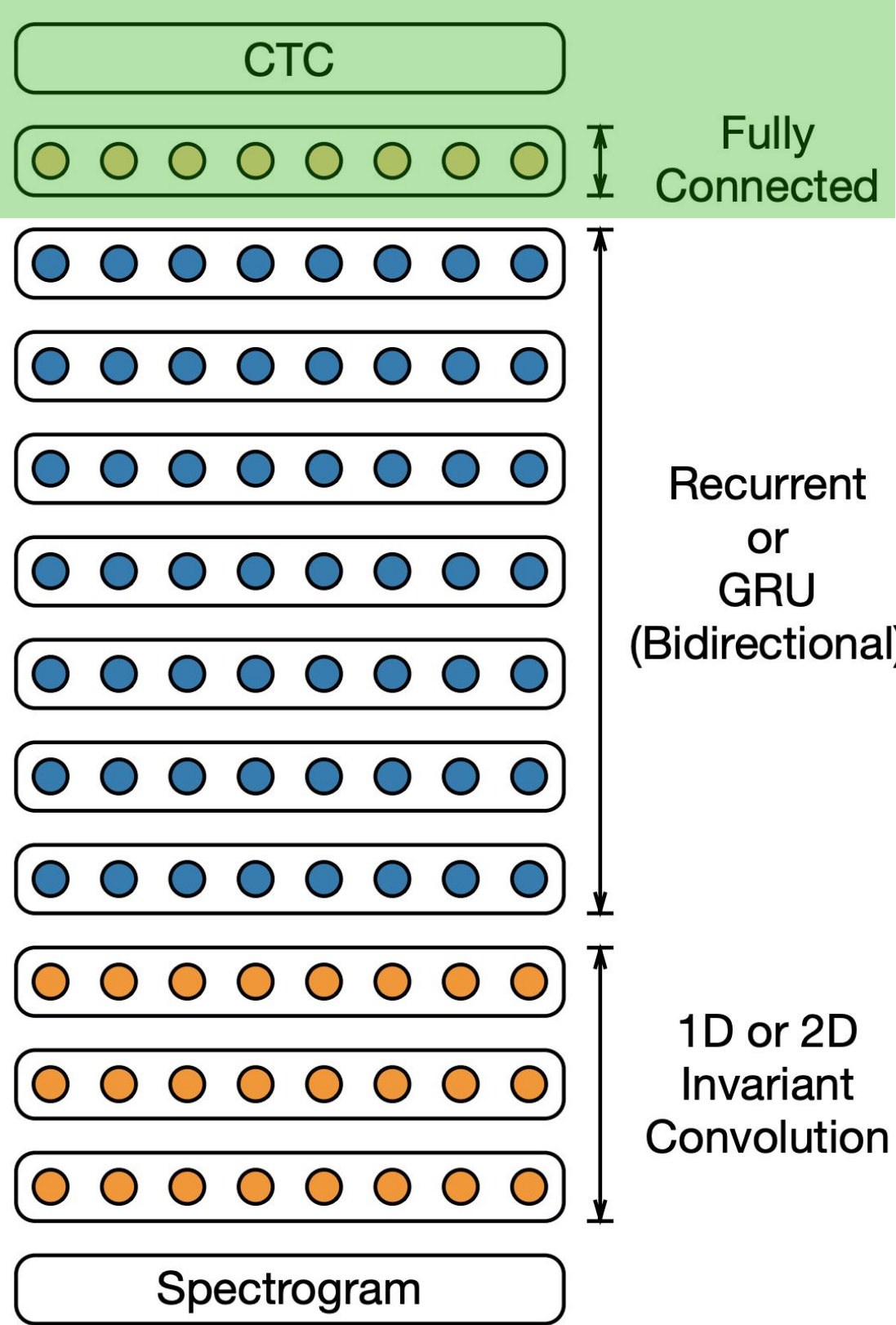


# Методы адаптации

# Методы адаптации



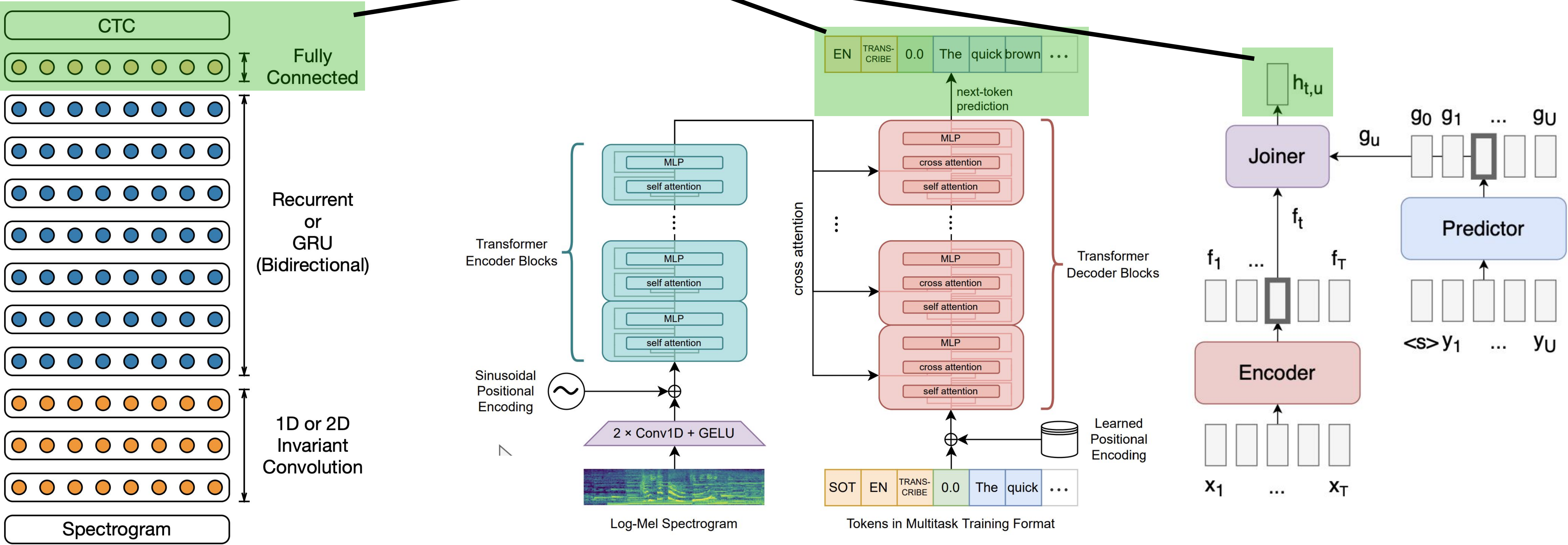
# Методы адаптации



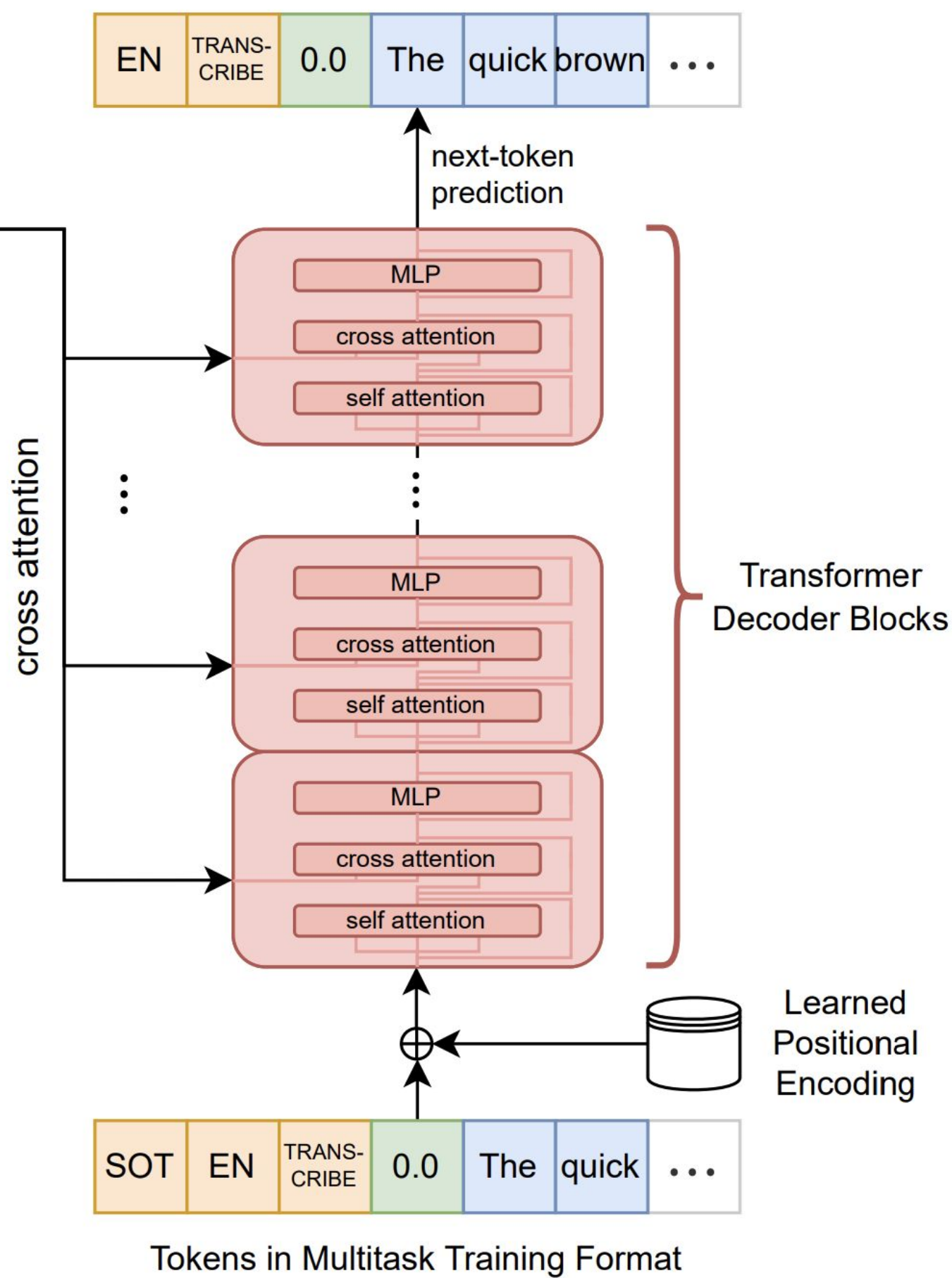


# Методы адаптации

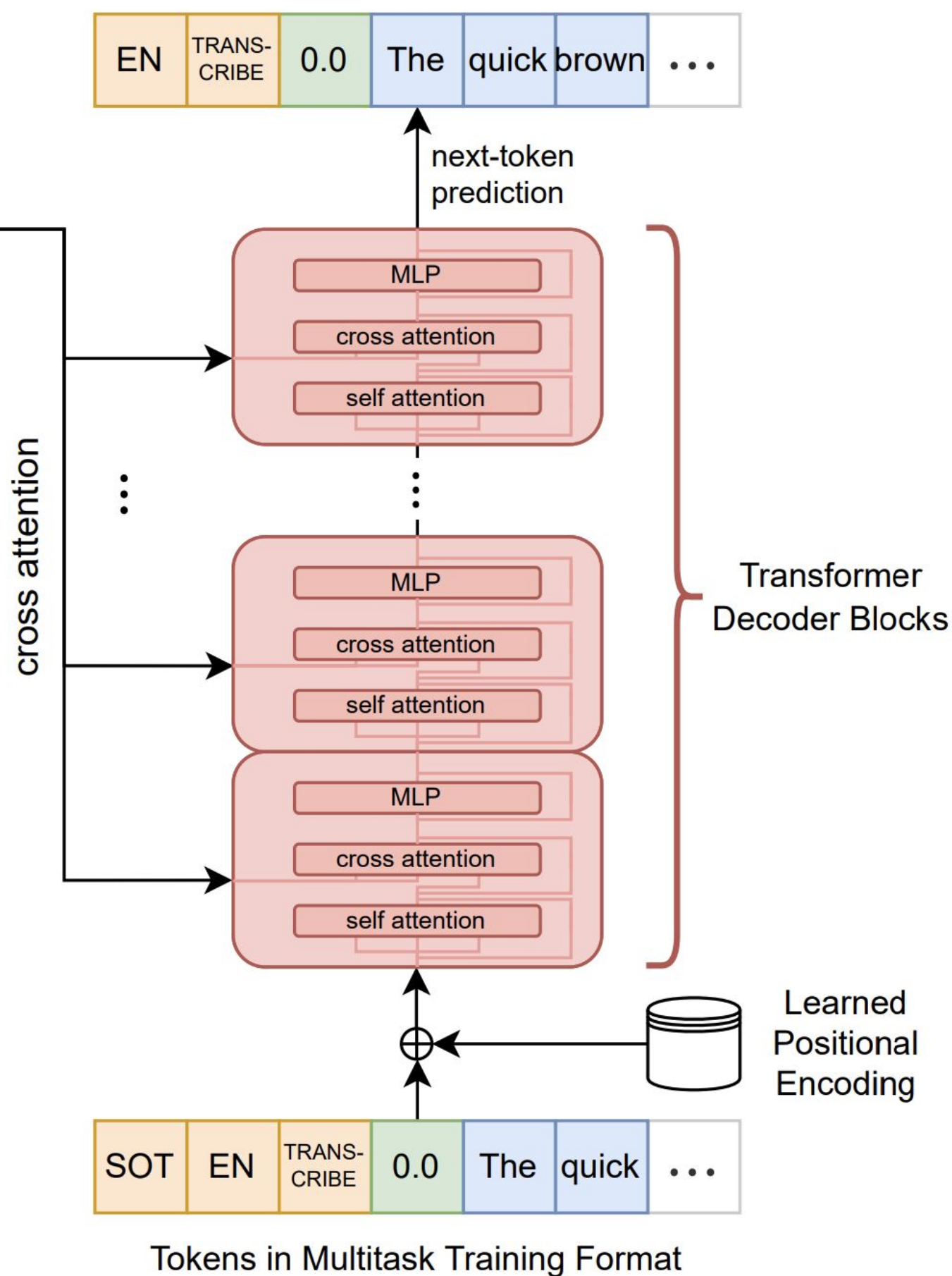
вероятность следующего токена



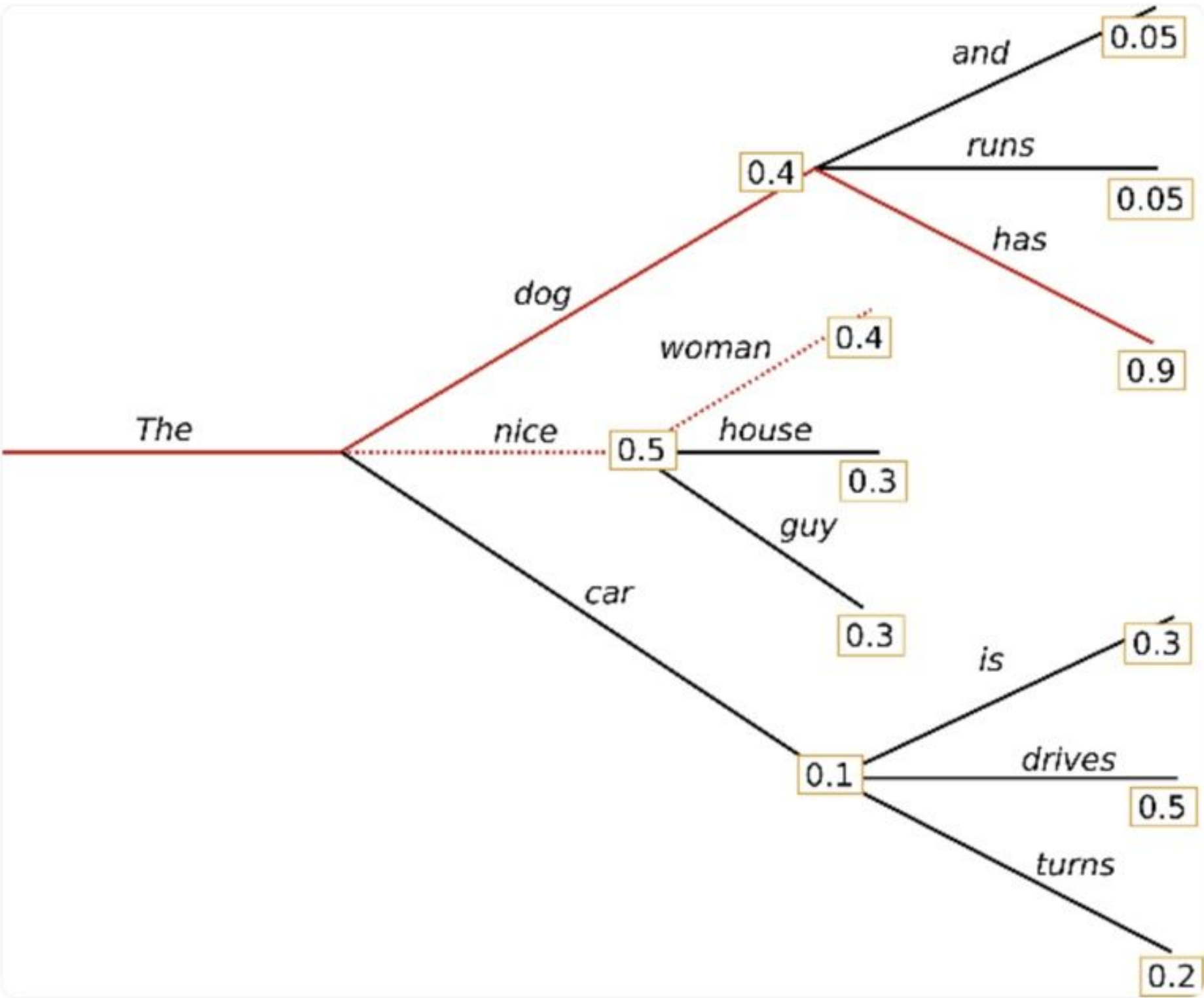
# Методы адаптации



# Методы адаптации

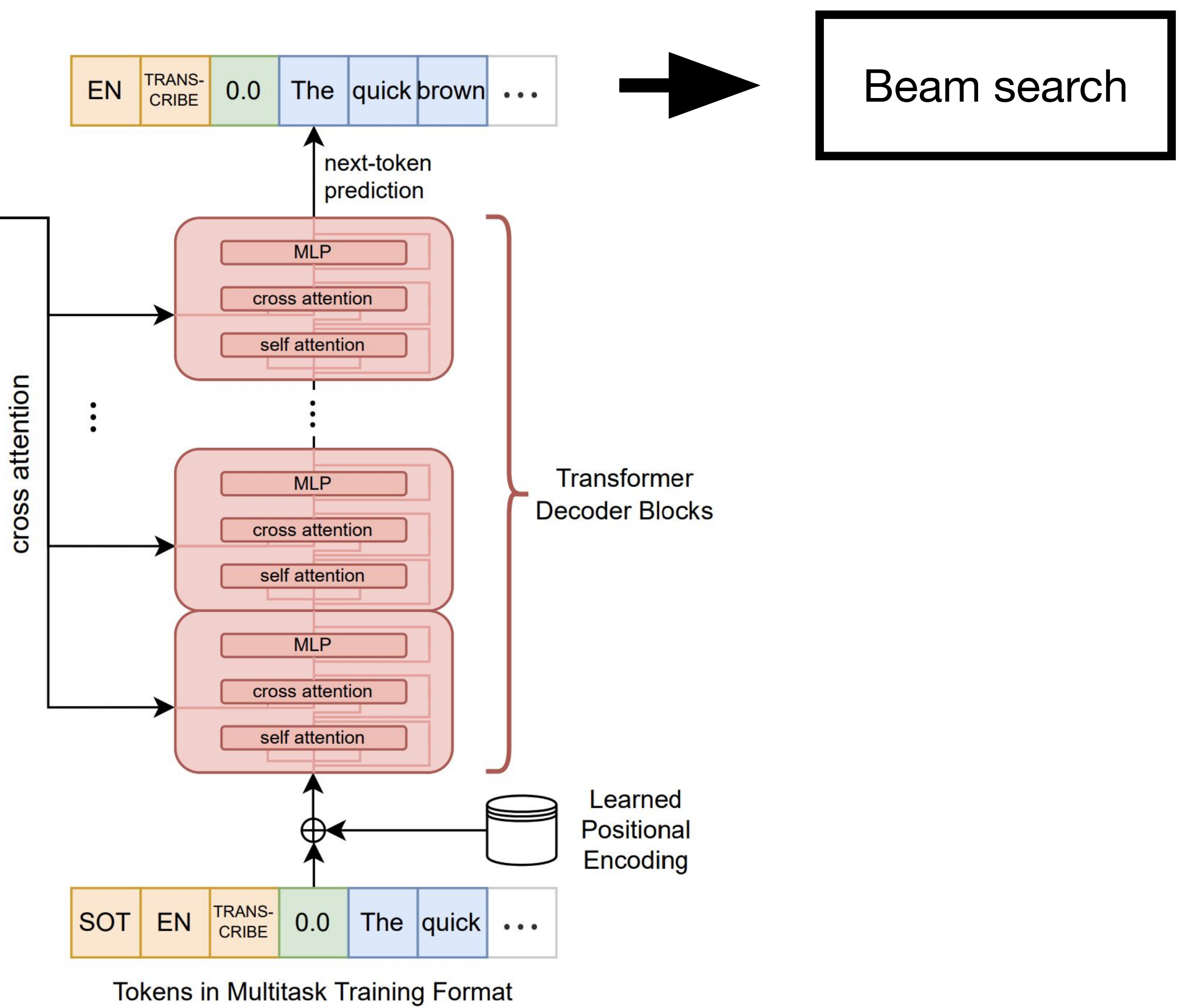


## Beam

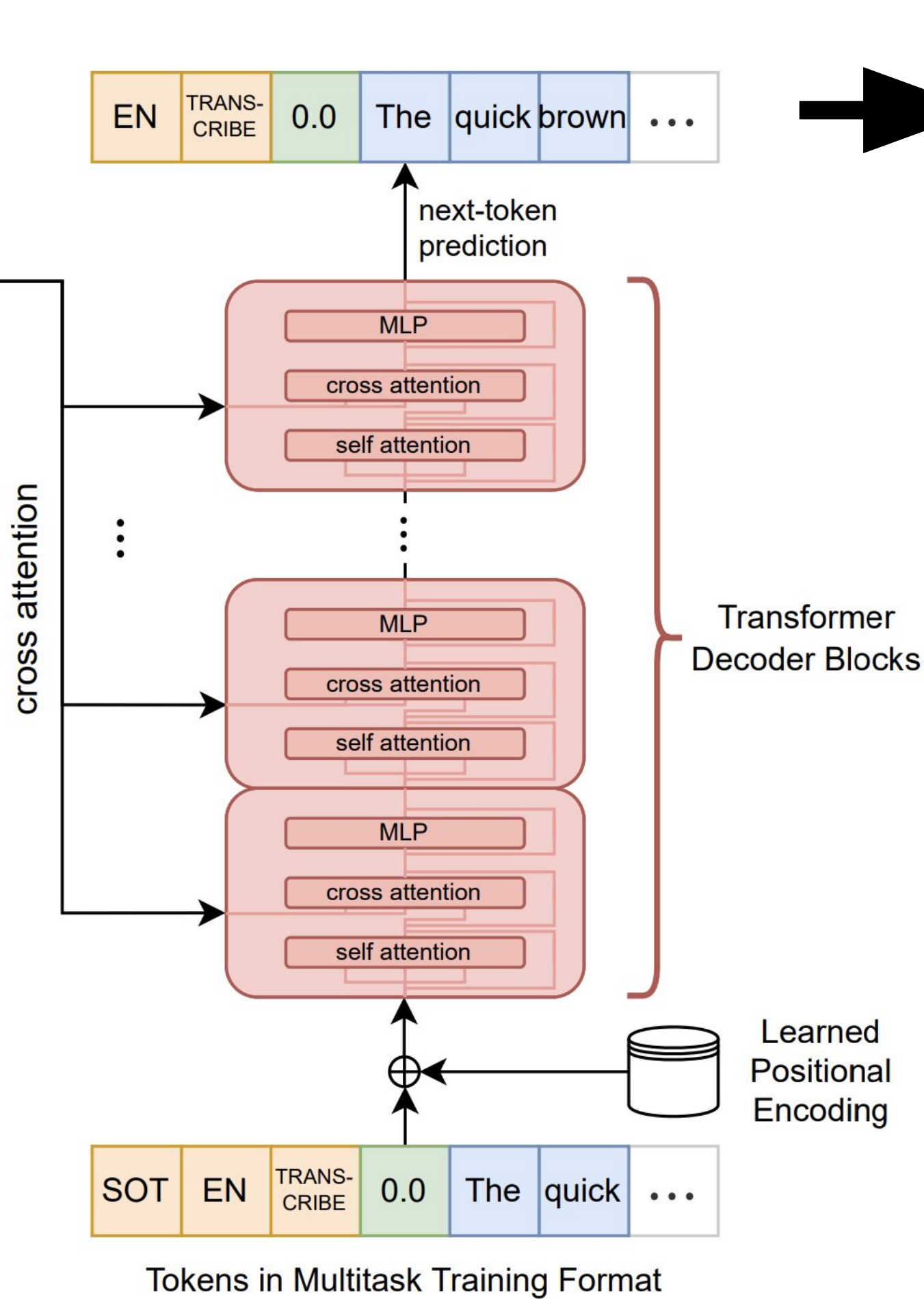




# Методы адаптации



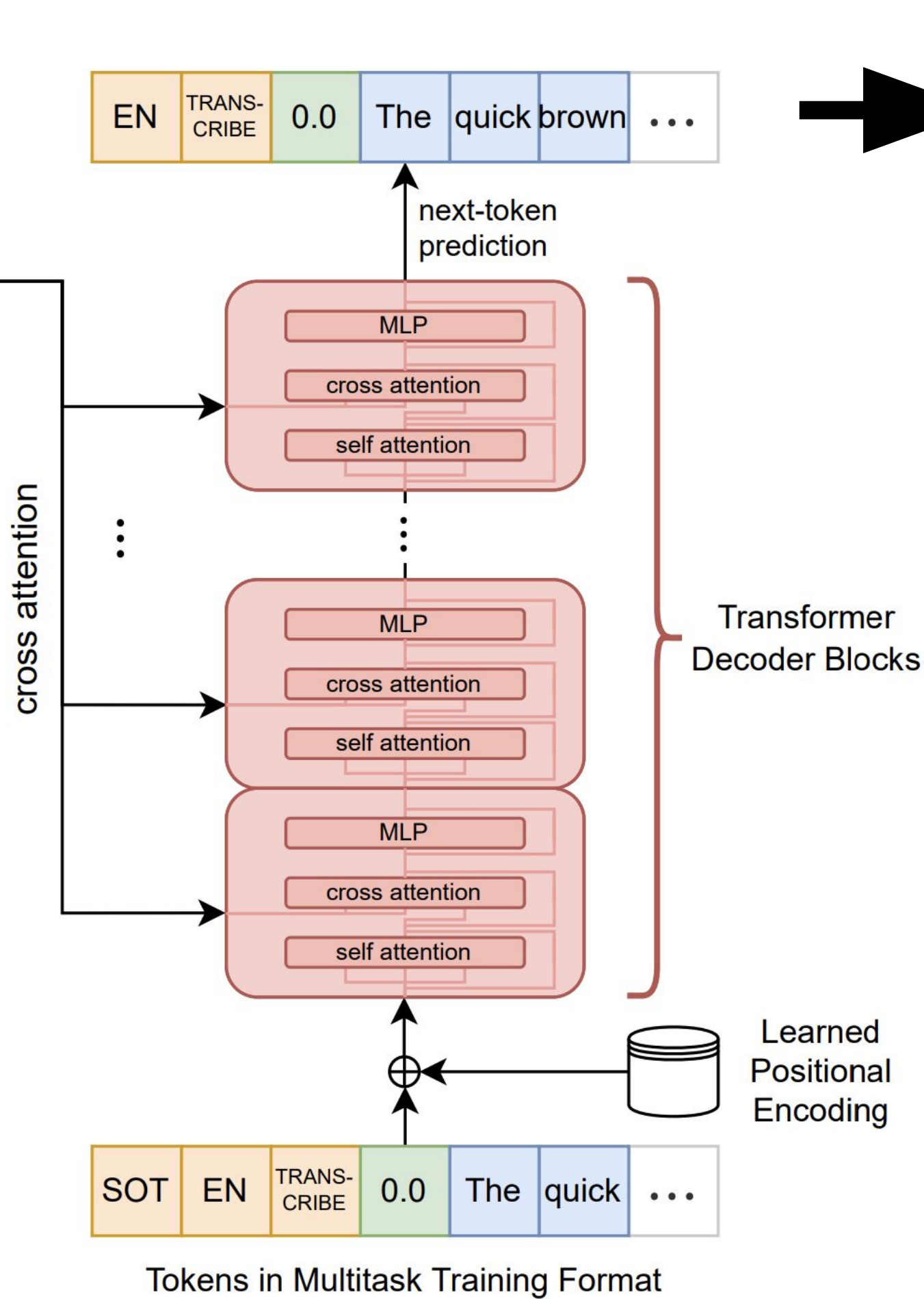
# Методы адаптации



Beam search

Адаптация может происходить на уровне пересчёта скоров в лучах и на уровне реранкинга итоговых последовательной

# Методы адаптации

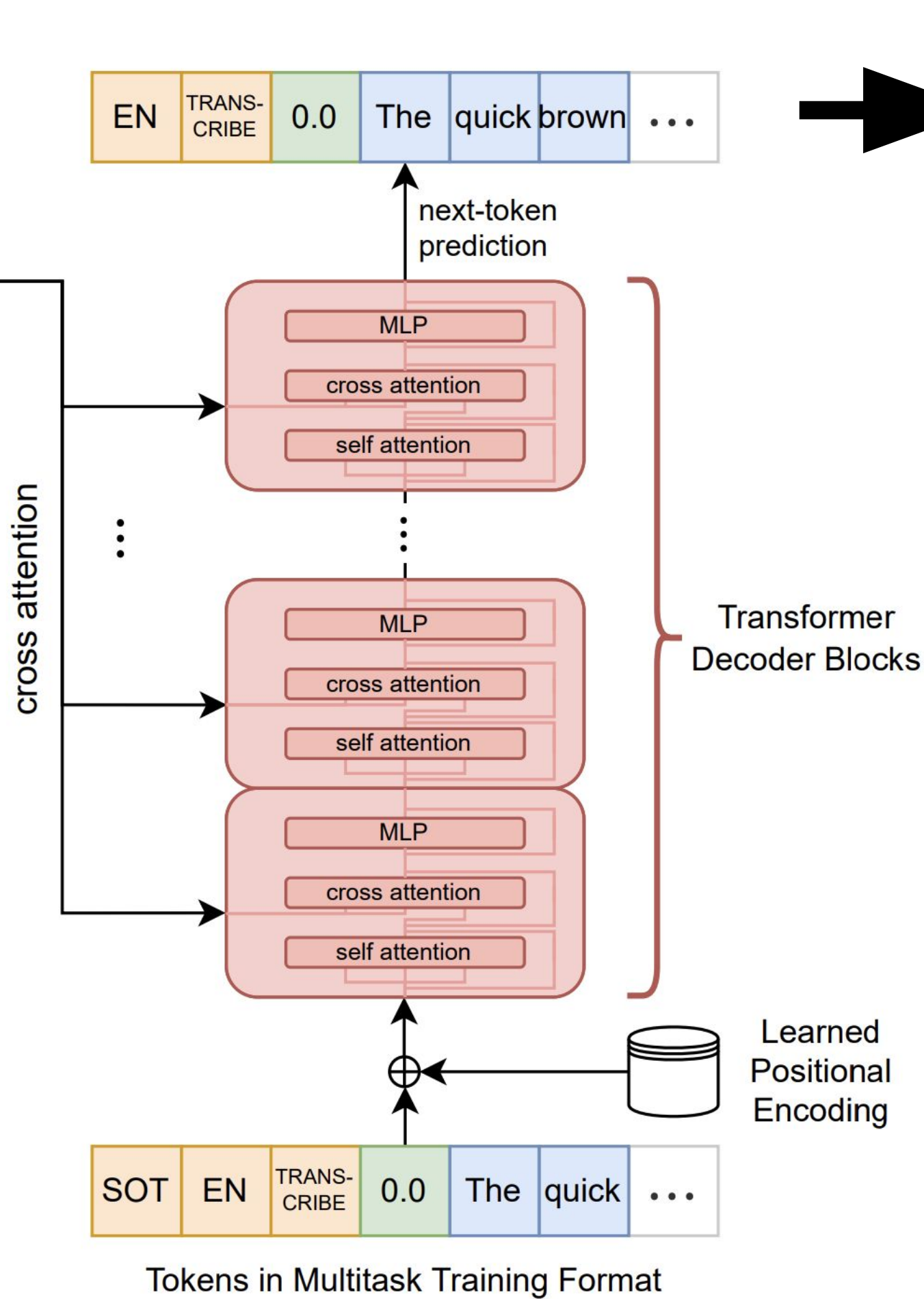


Beam search

Адаптация может происходить на уровне пересчёта скоров в лучах и на уровне **переранкинг итоговой последовательной**



# Методы адаптации

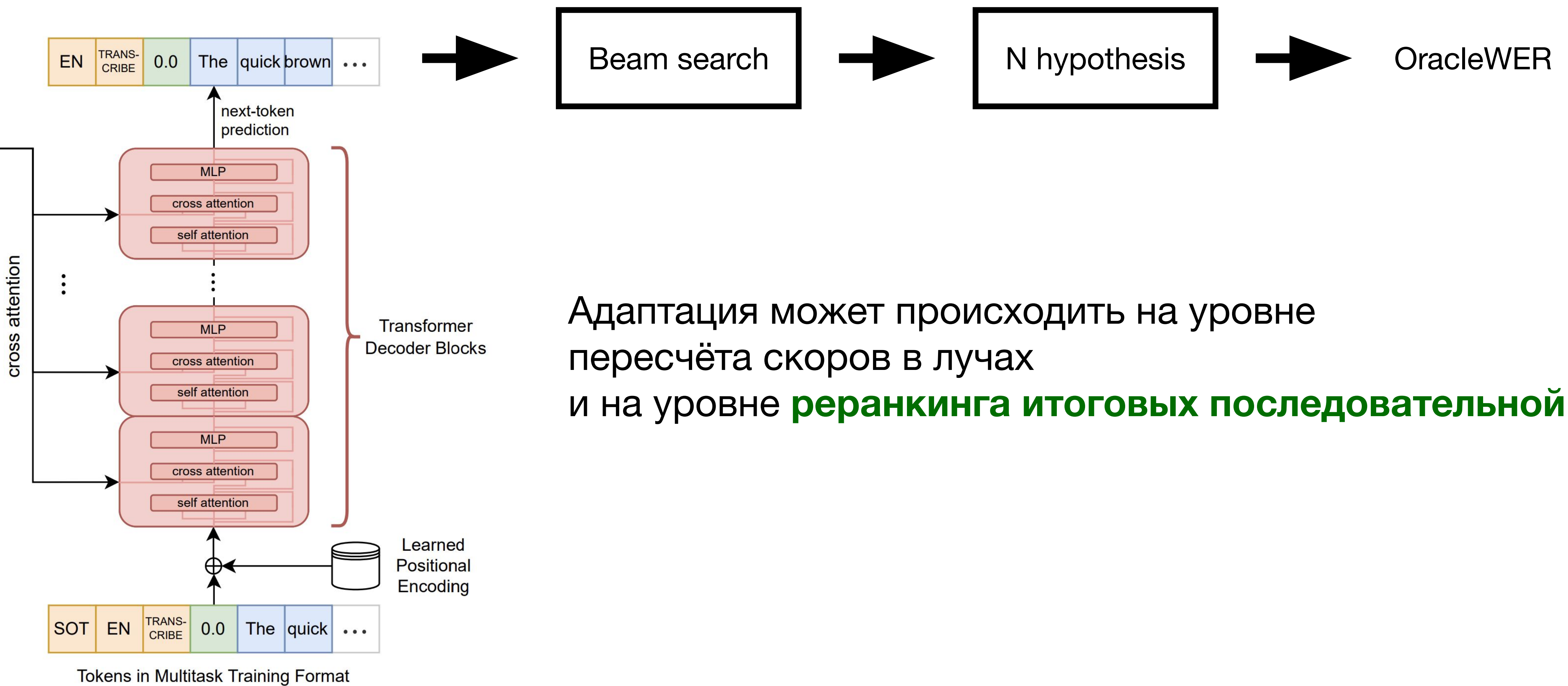


Beam search

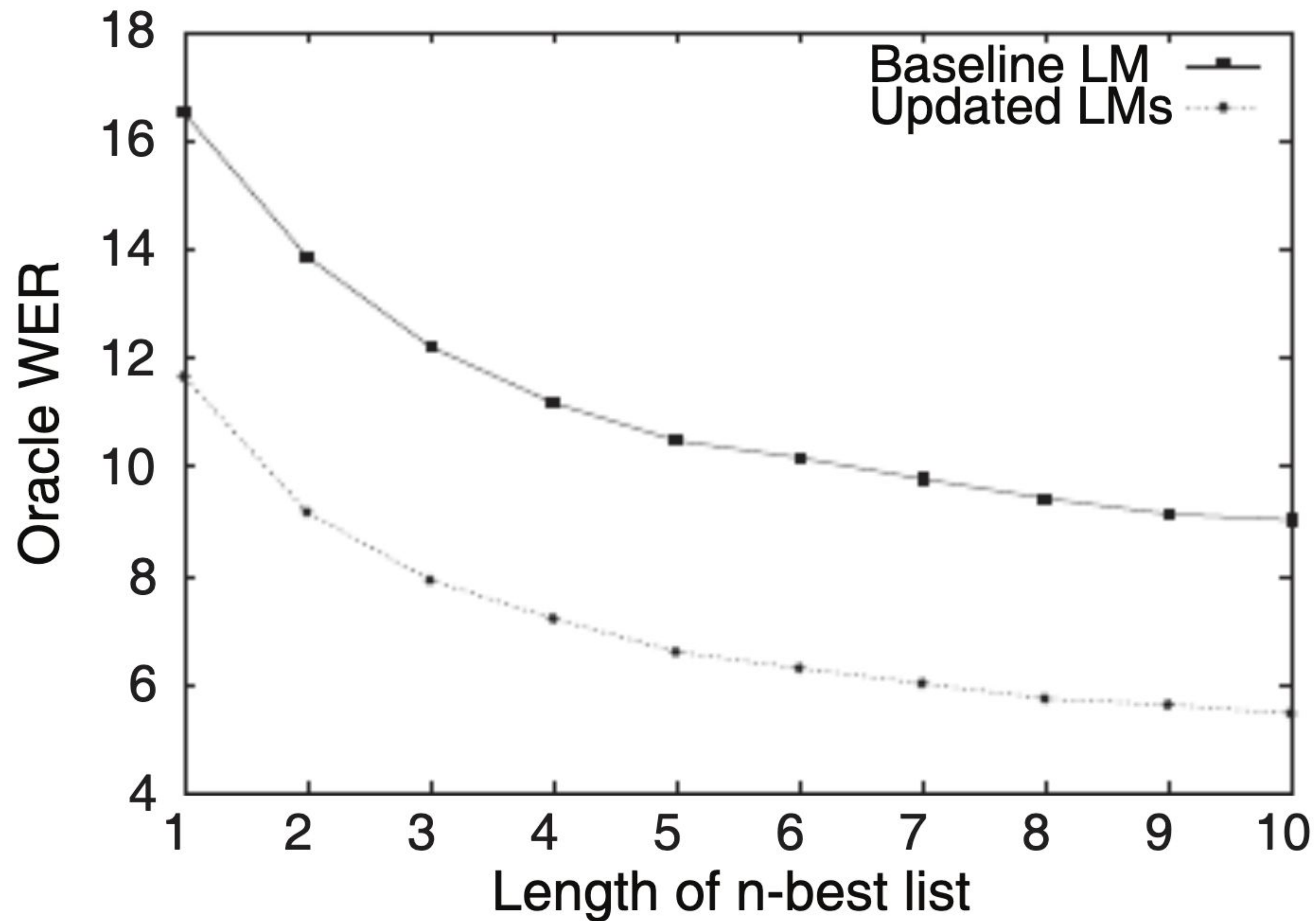
N hypothesis

Адаптация может происходить на уровне пересчёта скоров в лучах и на уровне **перанкинга итоговой последовательной**

# Методы адаптации

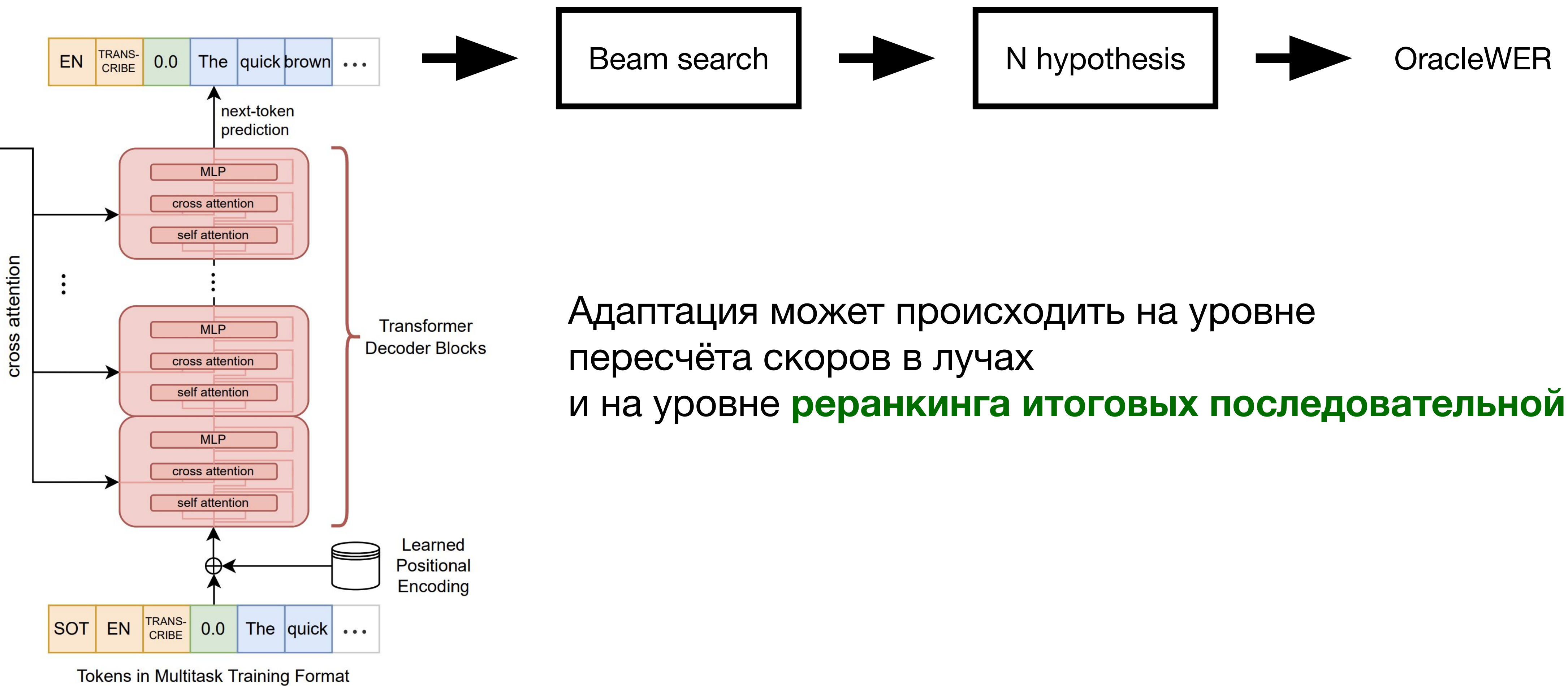


# Методы адаптации

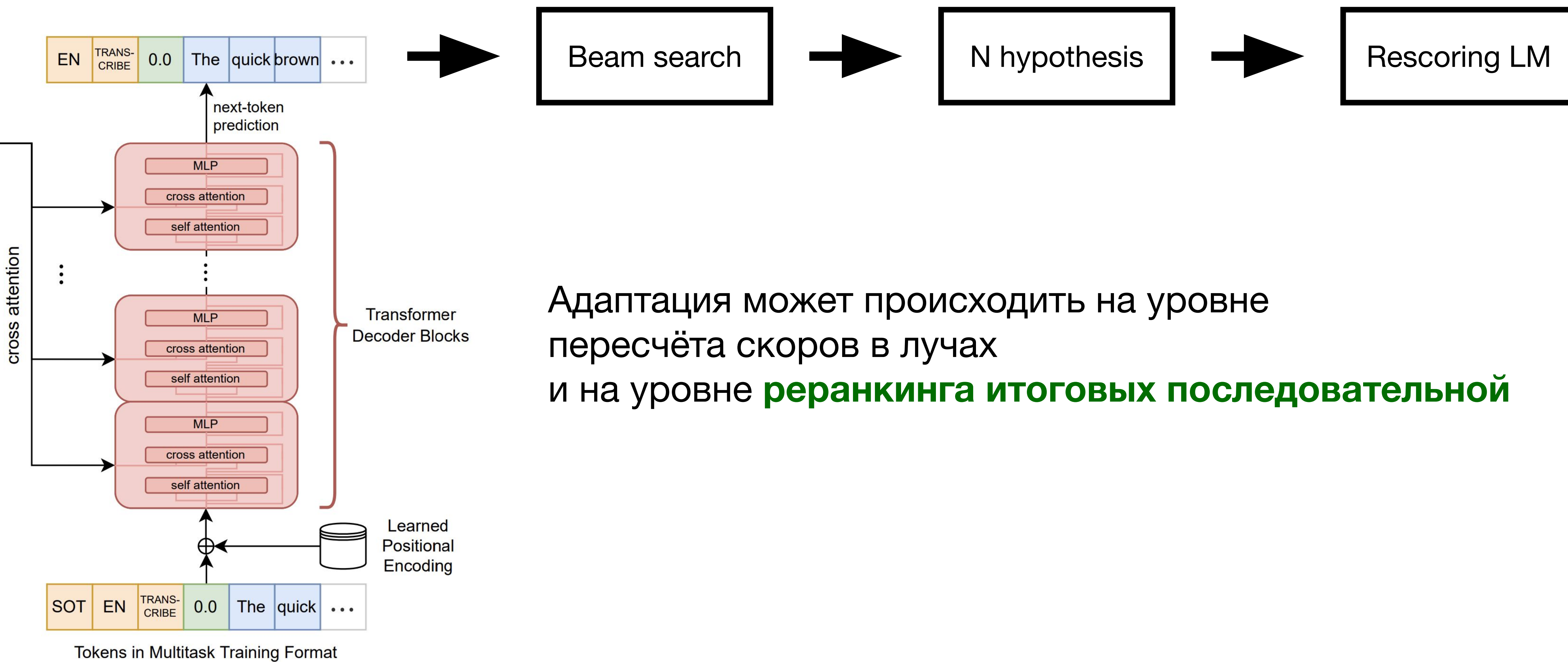




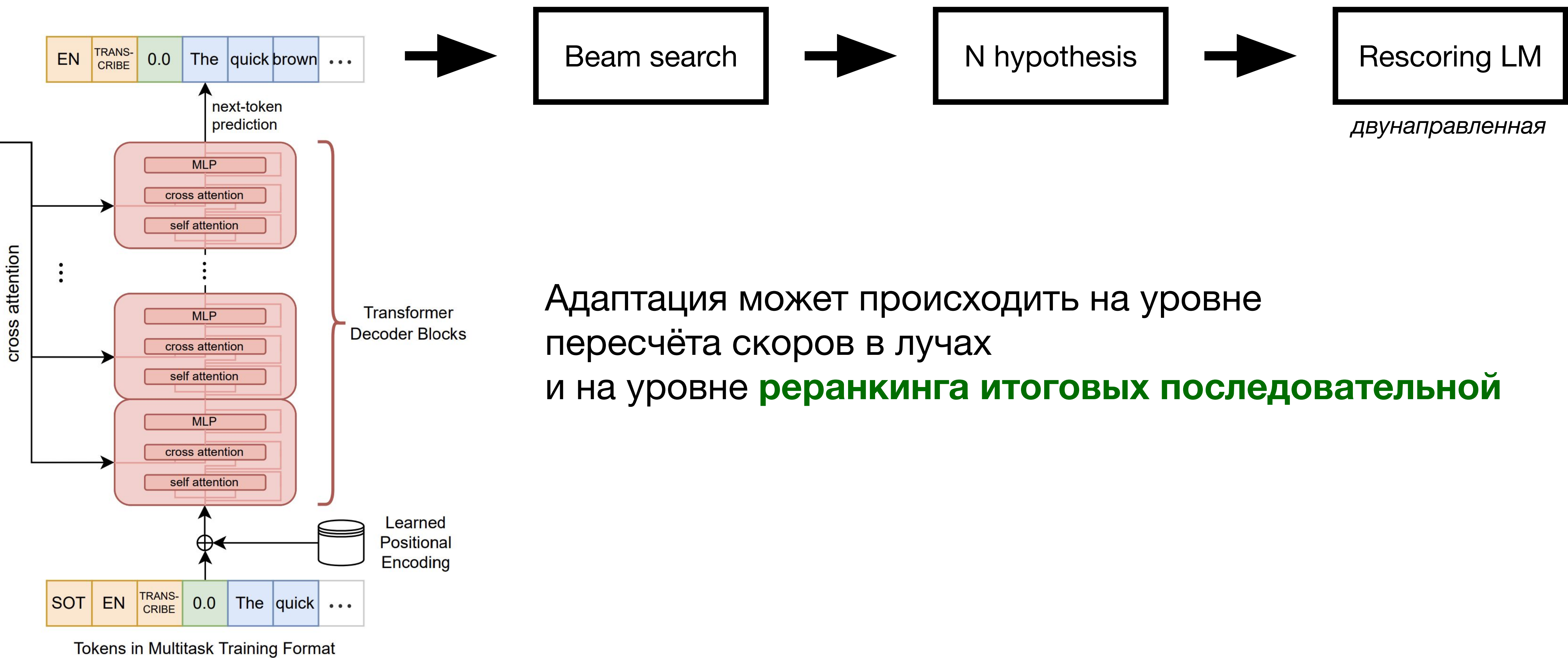
# Методы адаптации



# Методы адаптации

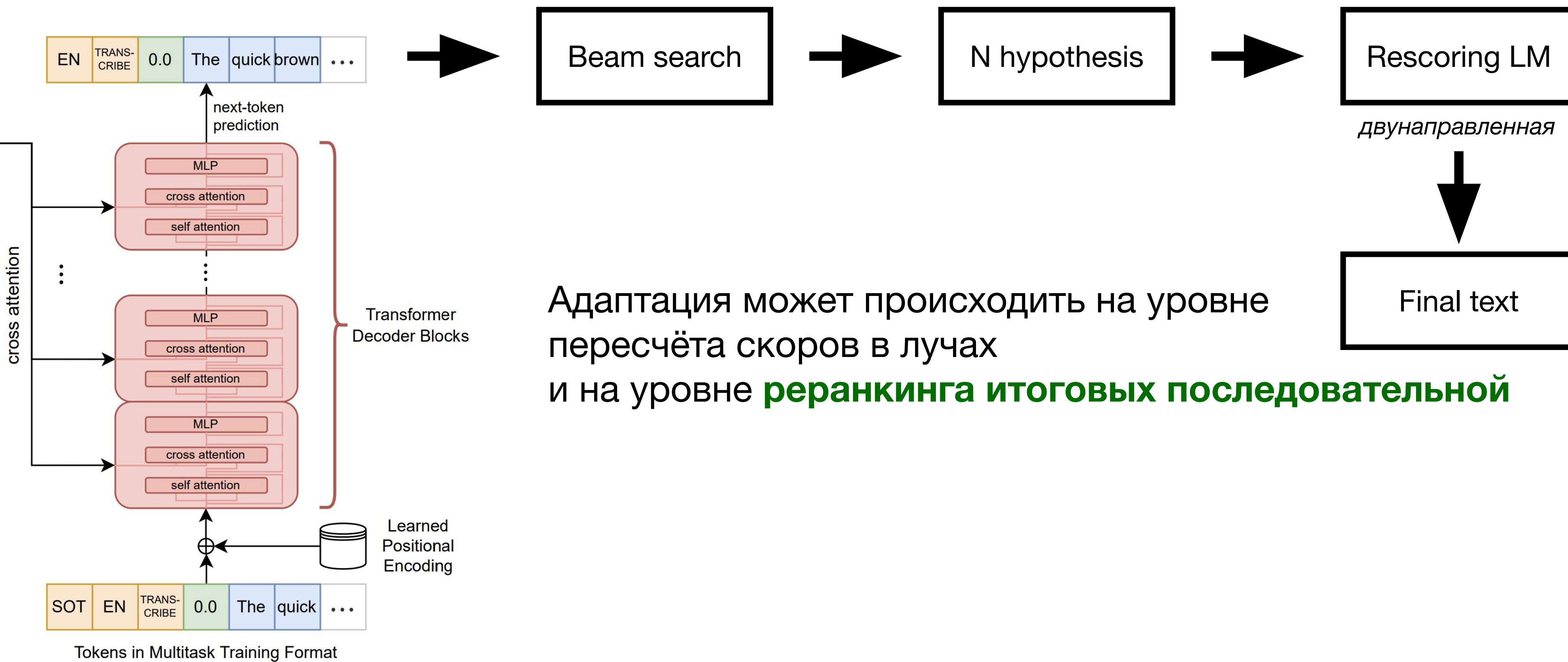


# Методы адаптации

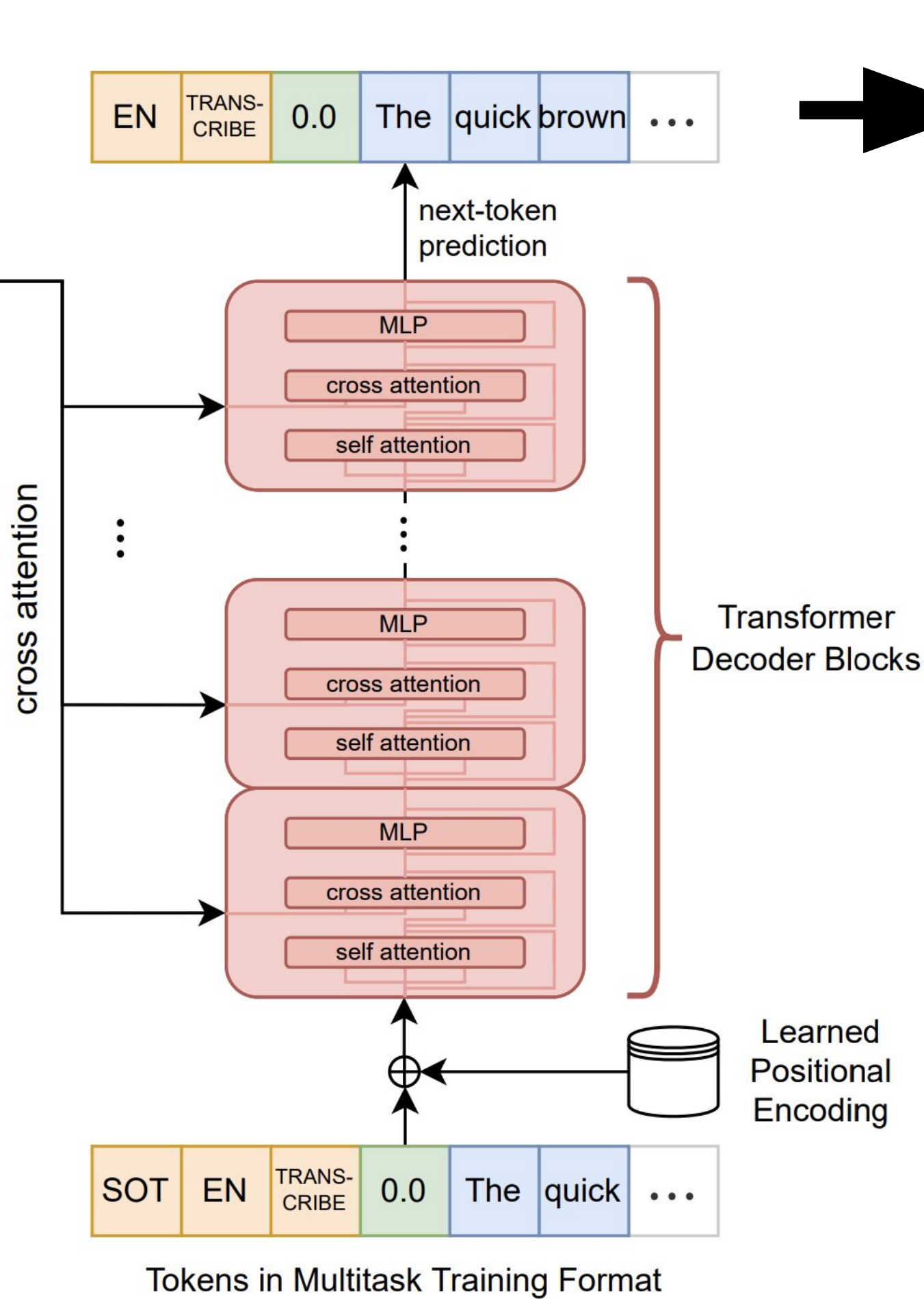




# Методы адаптации



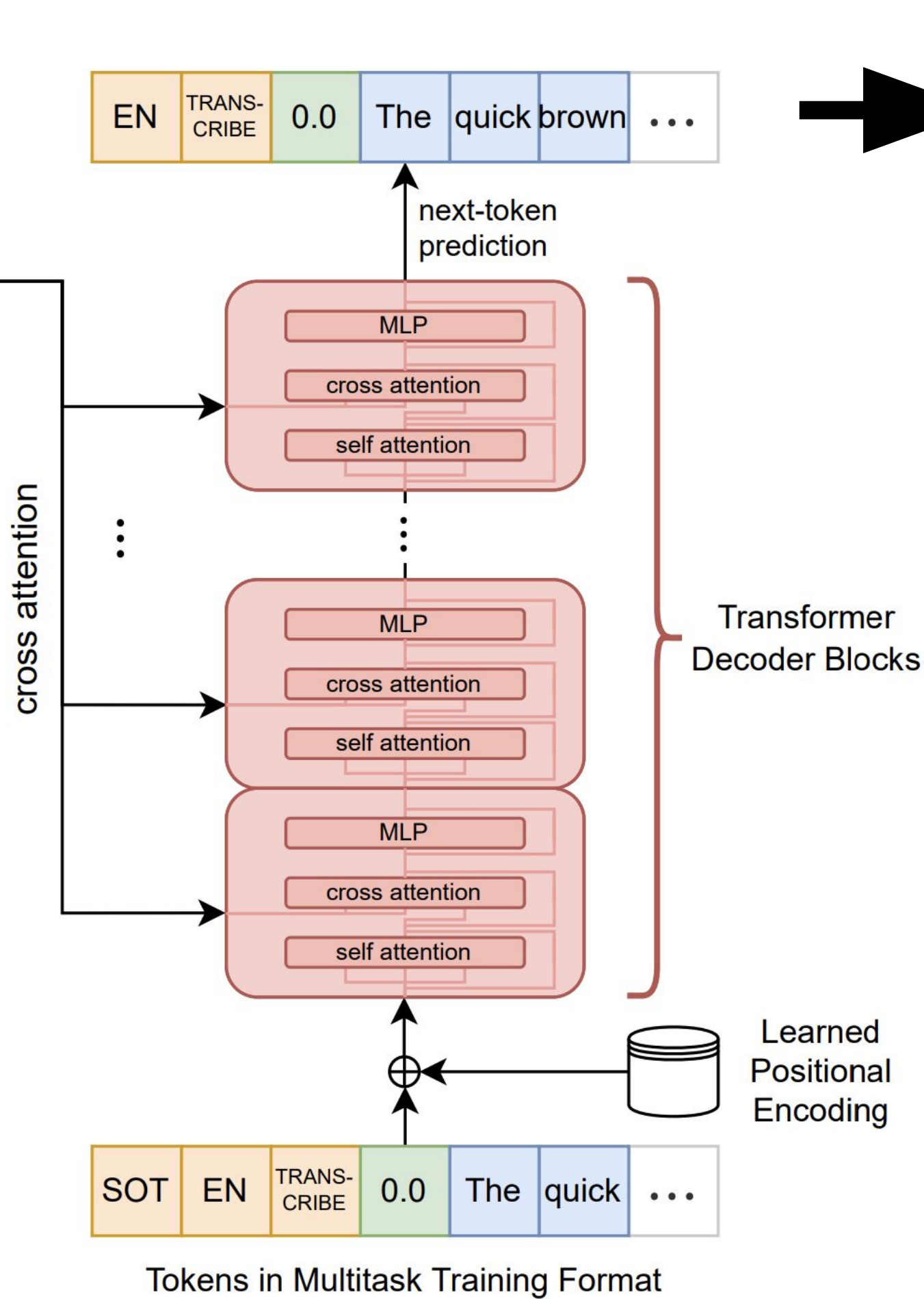
# Методы адаптации



Beam search

Адаптация может происходить на уровне пересчёта скоров в лучах и на уровне реранкинга итоговых последовательной

# Методы адаптации

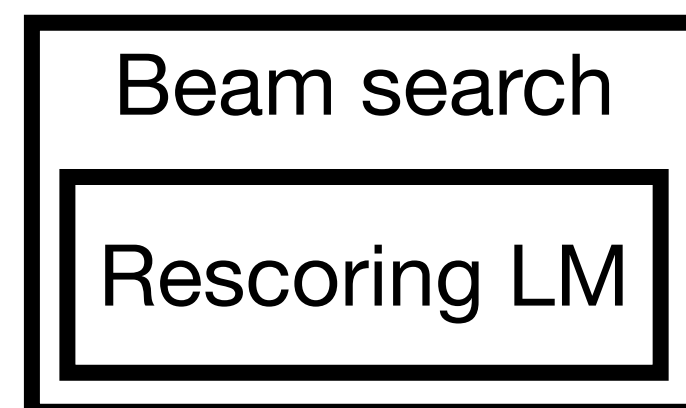
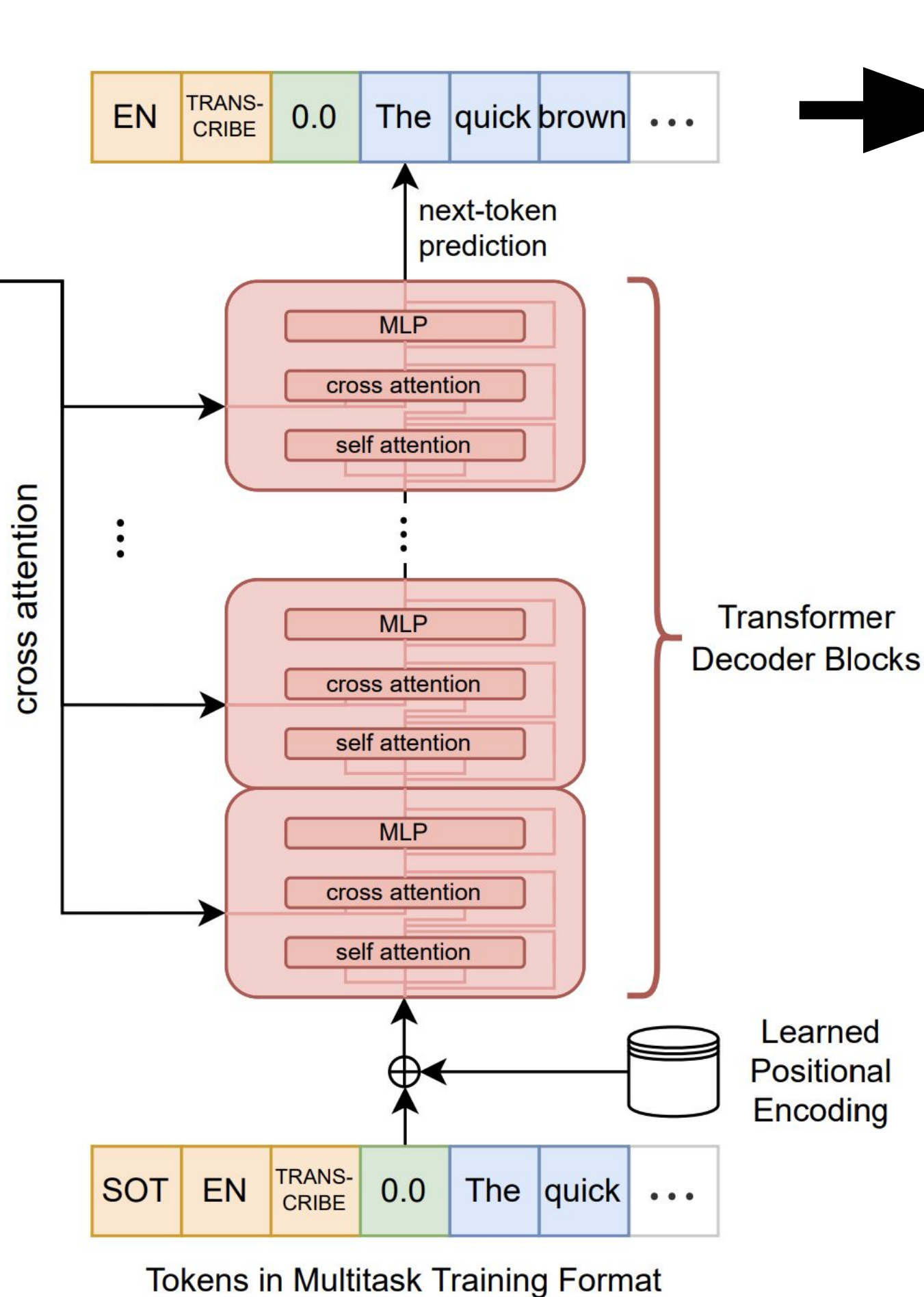


Beam search

Адаптация может происходить на уровне  
**пересчёта скоров в лучах**  
и на уровне реранкинга итоговых последовательной

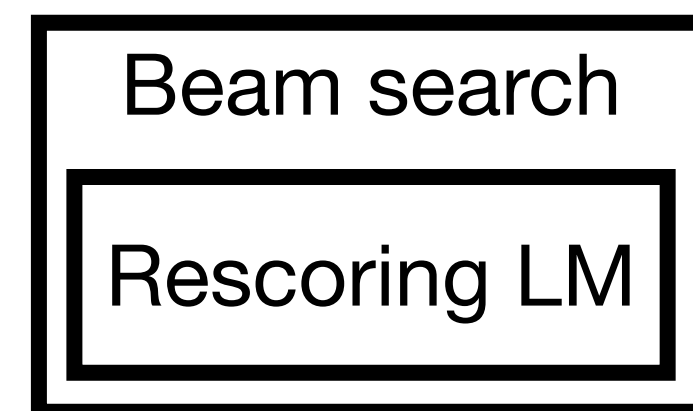
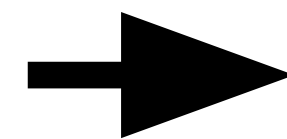
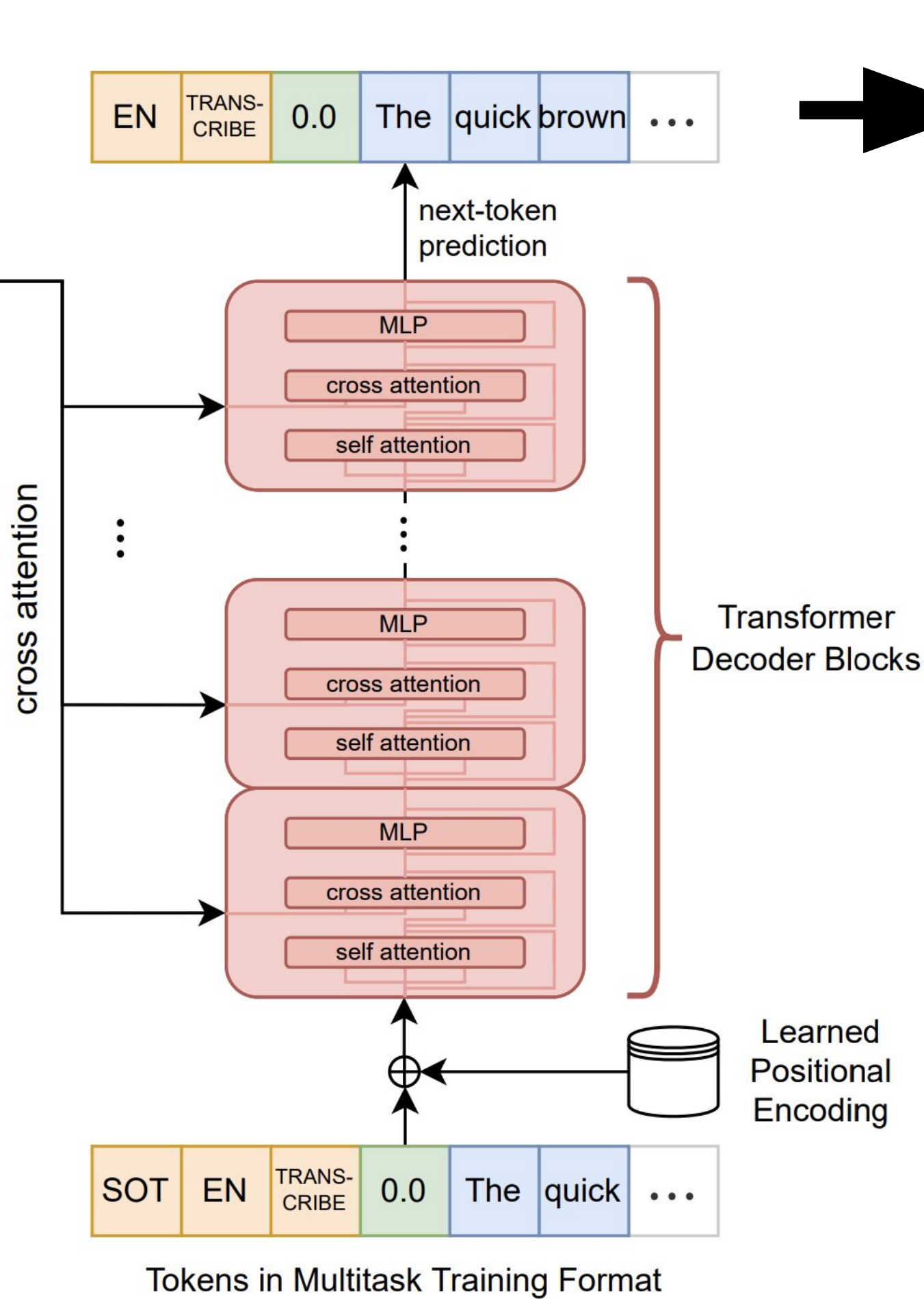


# Методы адаптации



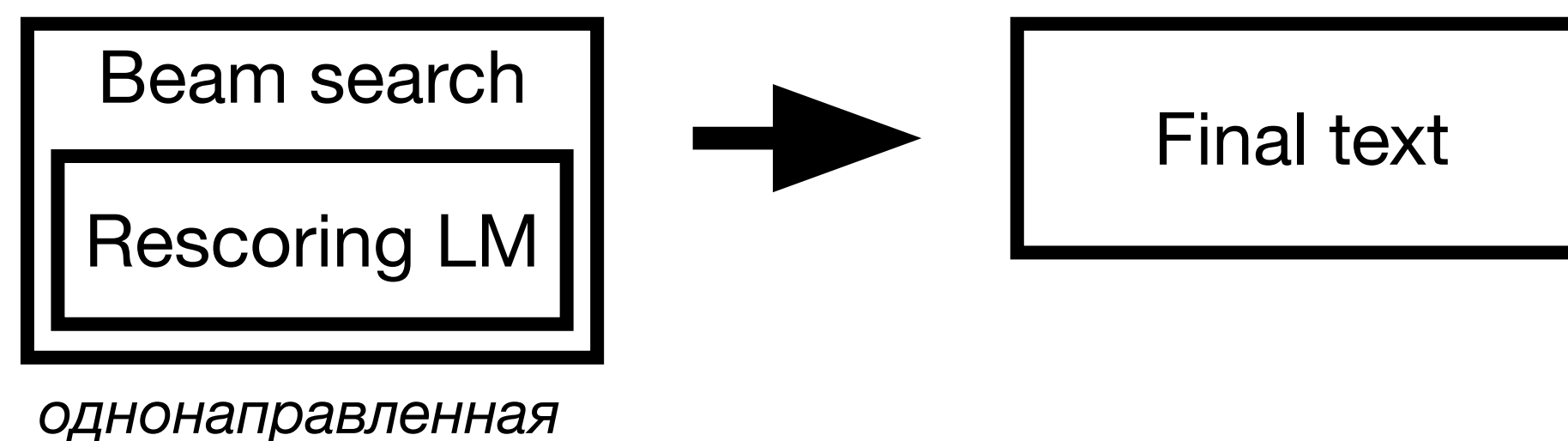
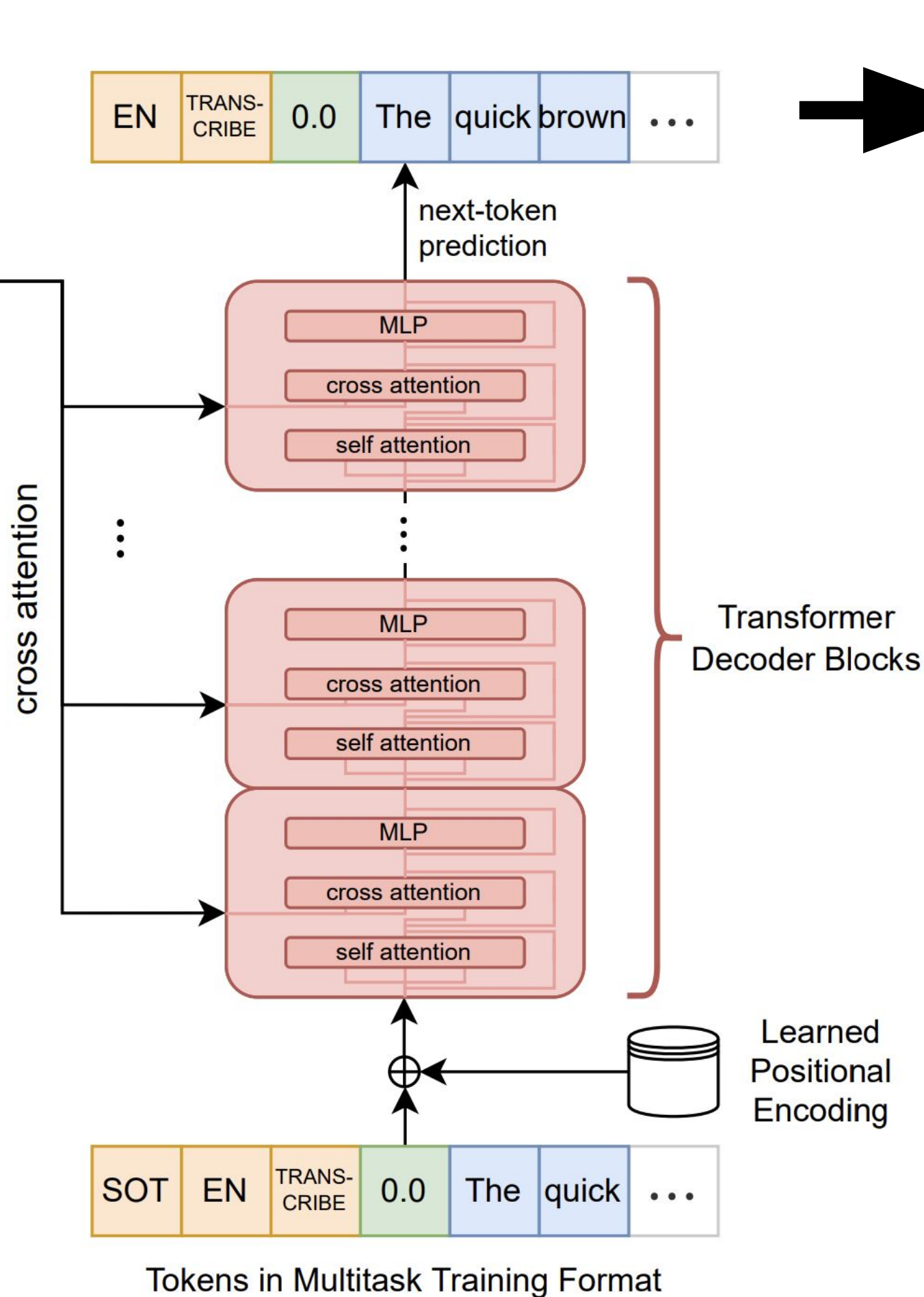
Адаптация может происходить на уровне  
**пересчёта скоров в лучах**  
и на уровне реранкинга итоговых последовательной

# Методы адаптации



Адаптация может происходить на уровне  
**пересчёта скоров в лучах**  
и на уровне реранкинга итоговых последовательной

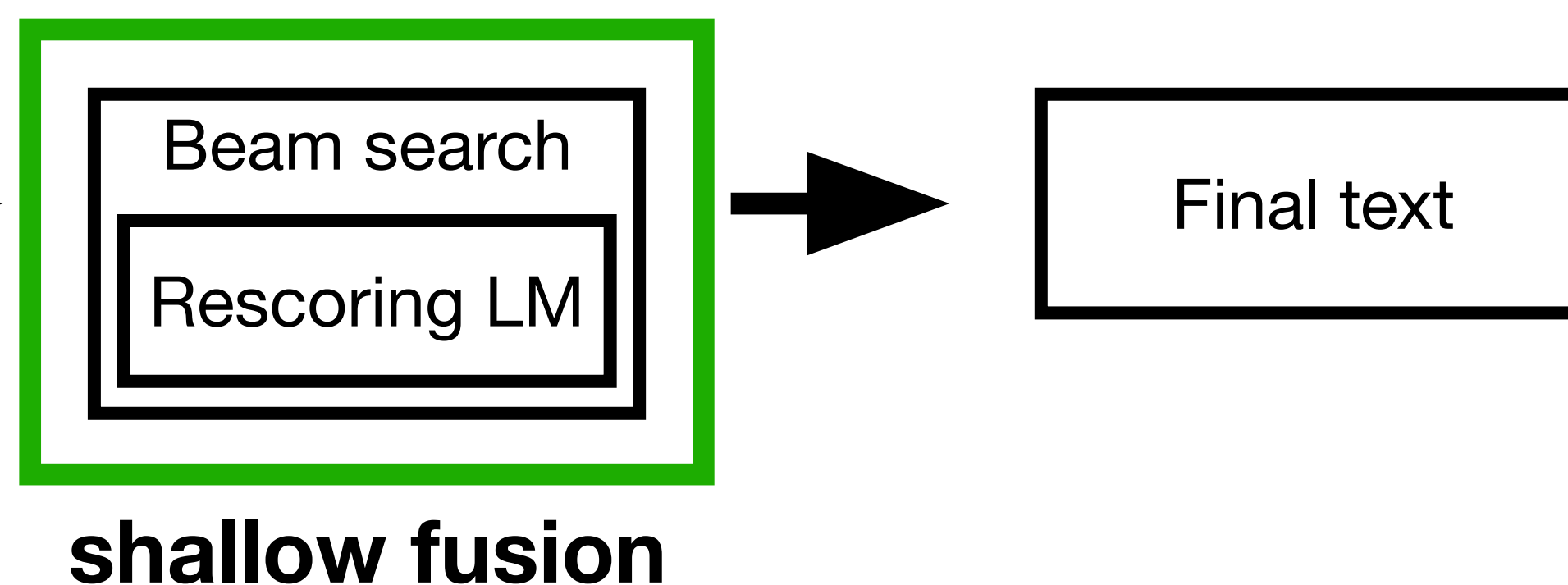
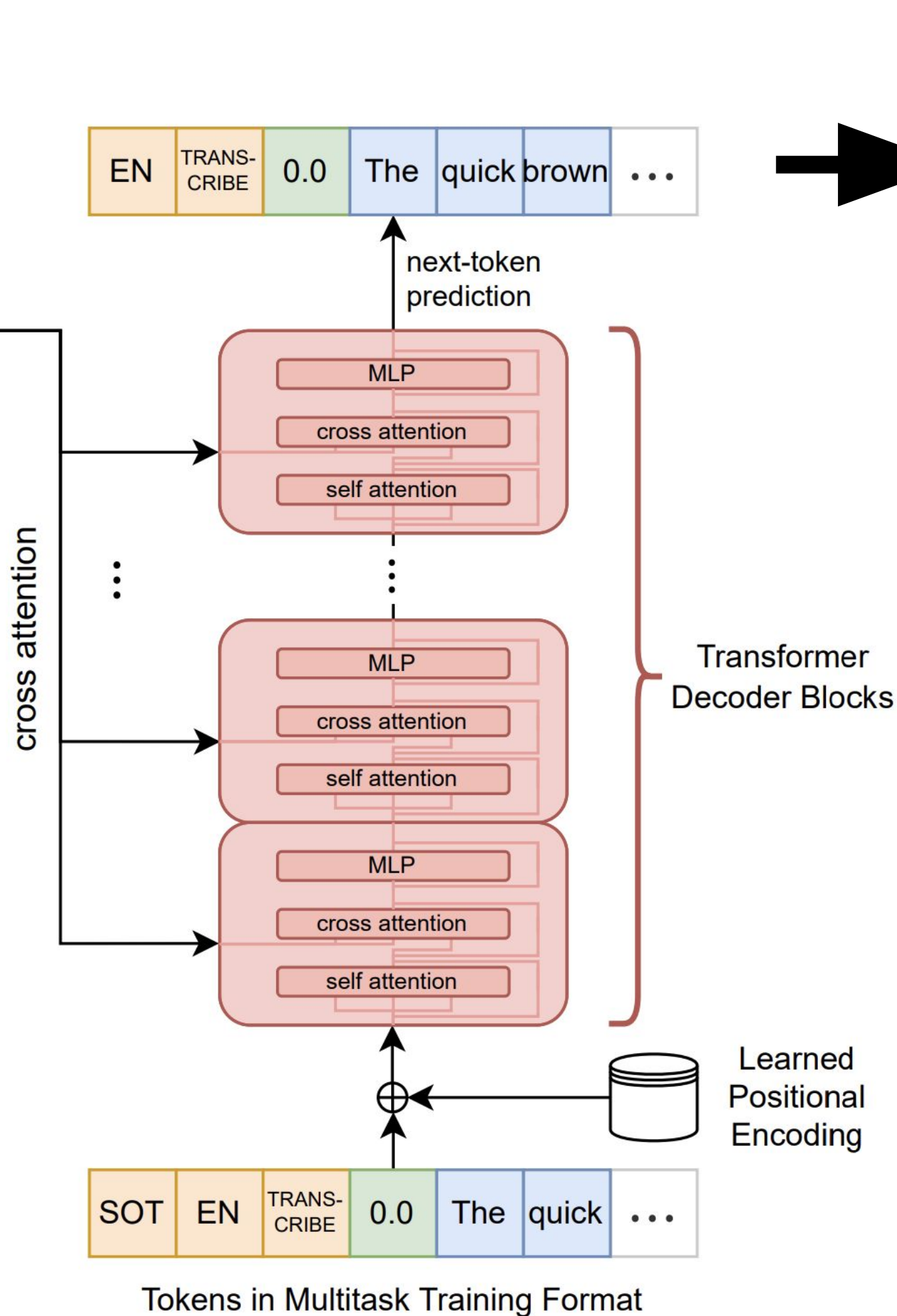
# Методы адаптации



Адаптация может происходить на уровне **пересчёта скоров в лучах** и на уровне реранкинга итоговых последовательной



# Методы адаптации



Адаптация может происходить на уровне  
**пересчёта скоров в лучах**  
и на уровне реранкинга итоговых последовательной

# Shallow fusion

# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$



# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$

# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$

# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$



# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$

# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$

# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$



# Функция смешения

Listen, Attend, Spell: 2015

$$s(\mathbf{y}|\mathbf{x}) = \frac{\log P(\mathbf{y}|\mathbf{x})}{|\mathbf{y}|_c} + \lambda \log P_{\text{LM}}(\mathbf{y})$$

# Функция смешения

Google: 2017

$$\mathbf{y}^* = \arg \max_y \log p(y|x) + \lambda \log p_{LM}(y) + \gamma c(x, y)$$

# Функция смешения

Google: 2017

$$\mathbf{y}^* = \boxed{\arg \max_y \log p(y|x)} + \lambda \log p_{LM}(y) + \gamma c(x, y)$$



# Функция смешения

Google: 2017

$$\mathbf{y}^* = \arg \max_y \log p(y|x) + \lambda \log p_{LM}(y) + \gamma c(x, y)$$

# Функция смешения

Google: 2017

$$\mathbf{y}^* = \arg \max_y \log p(y|x) + \lambda \log p_{LM}(y) + \gamma c(x, y)$$

# Функция смешения

Google: 2017

$$c(x, y) = \sum_j \log(\min(\sum_i a_{i,j}, 0.5))$$



# Функция смешения


Google: 2017

$$c(x, y) = \sum_j \log(\min(\sum_i a_{i,j}, 0.5))$$

↑  
frames

# Функция смешения

Google: 2017

$$c(x, y) = \sum_j \log(\min(\sum_i a_{i,j}, 0.5))$$


The diagram illustrates the variables in the equation. An upward-pointing arrow connects the word "tokens" to the index  $j$  in the summation  $\sum_j$ . Another upward-pointing arrow connects the word "frames" to the index  $i$  in the summation  $\sum_i$ .

# Функция смешения

Google: 2017

$$c(x, y) = \sum_j \log(\min(\sum_i a_{i,j}, 0.5))$$

Diagram illustrating the function  $c(x, y)$  with annotations:

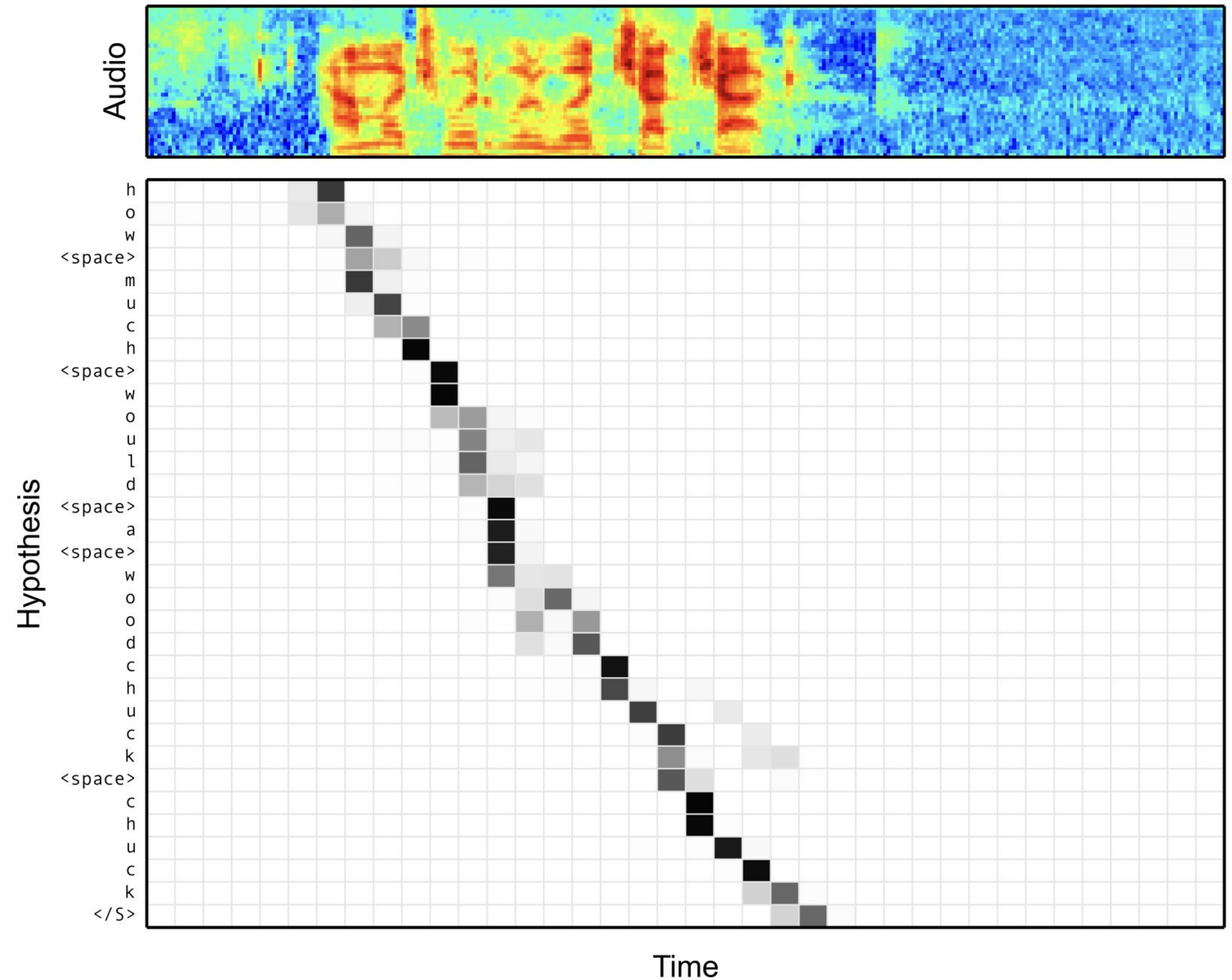
- The variable  $j$  is labeled "tokens" with an upward arrow.
- The variable  $i$  is labeled "frames" with an upward arrow.
- The term  $a_{i,j}$  is labeled "attention prob" with a downward arrow.



# Функция смешения

Alignment between the Characters and Audio

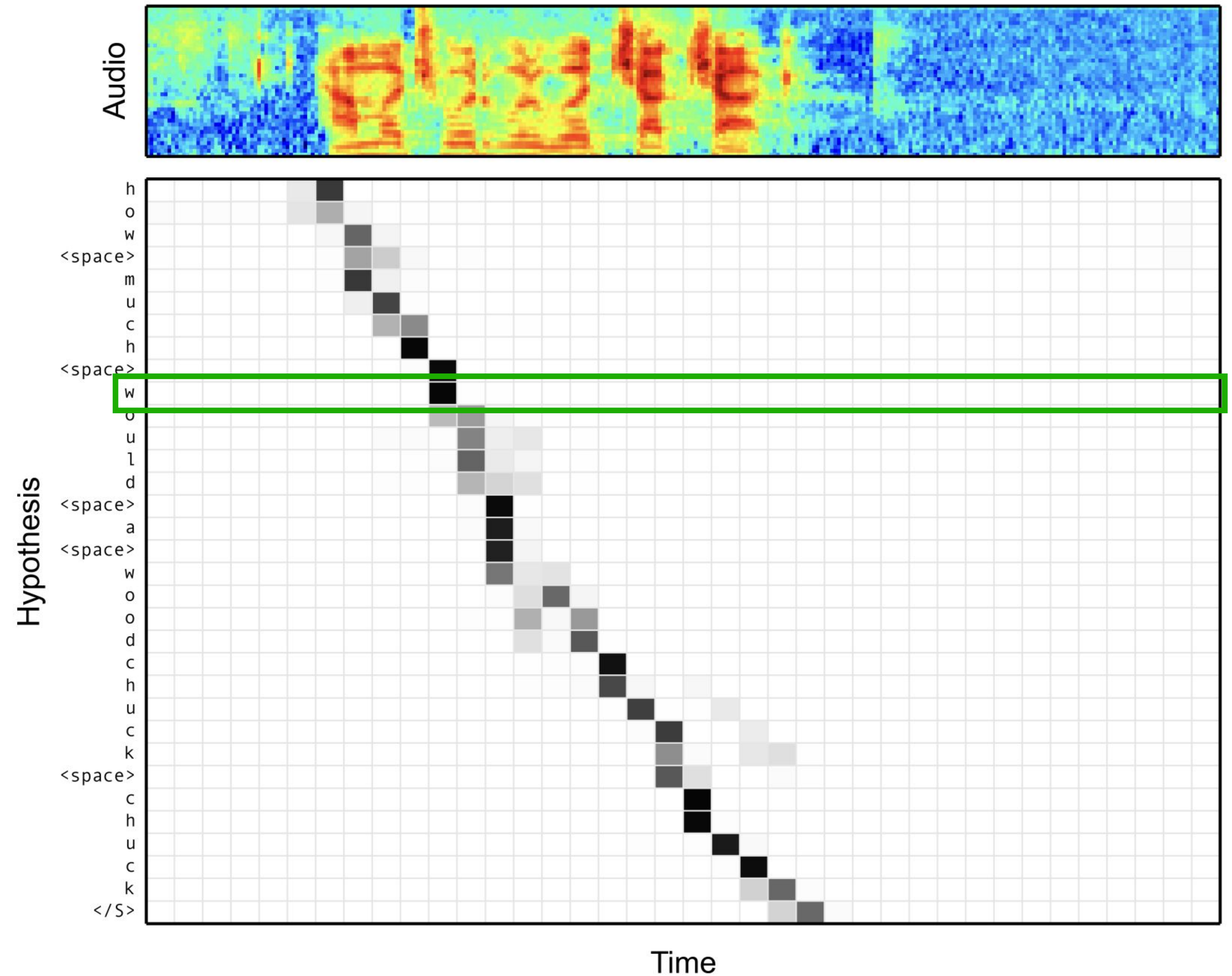
Google: 2017



# Функция смещения

Alignment between the Characters and Audio

Google: 2017

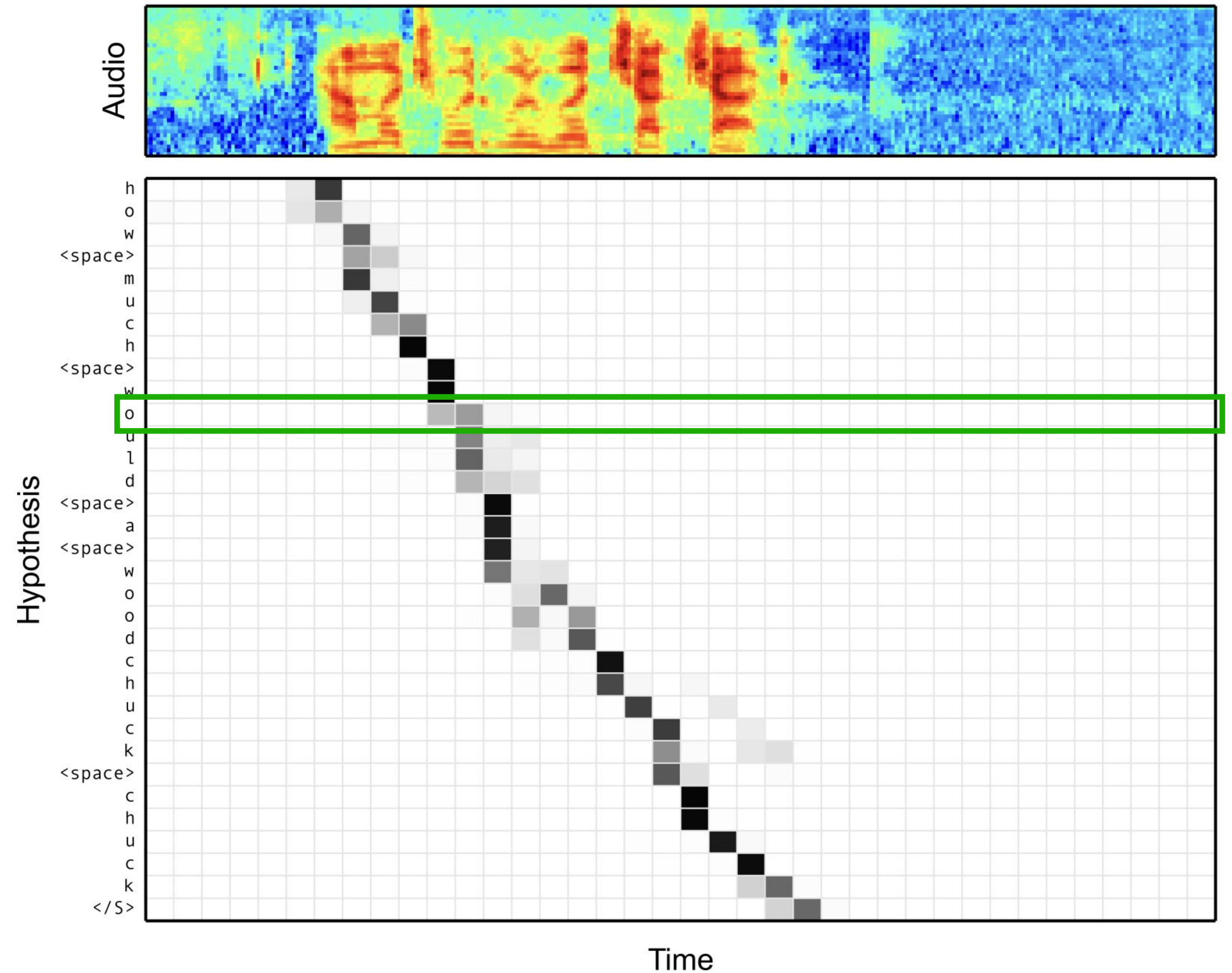




# Функция смещения

Alignment between the Characters and Audio

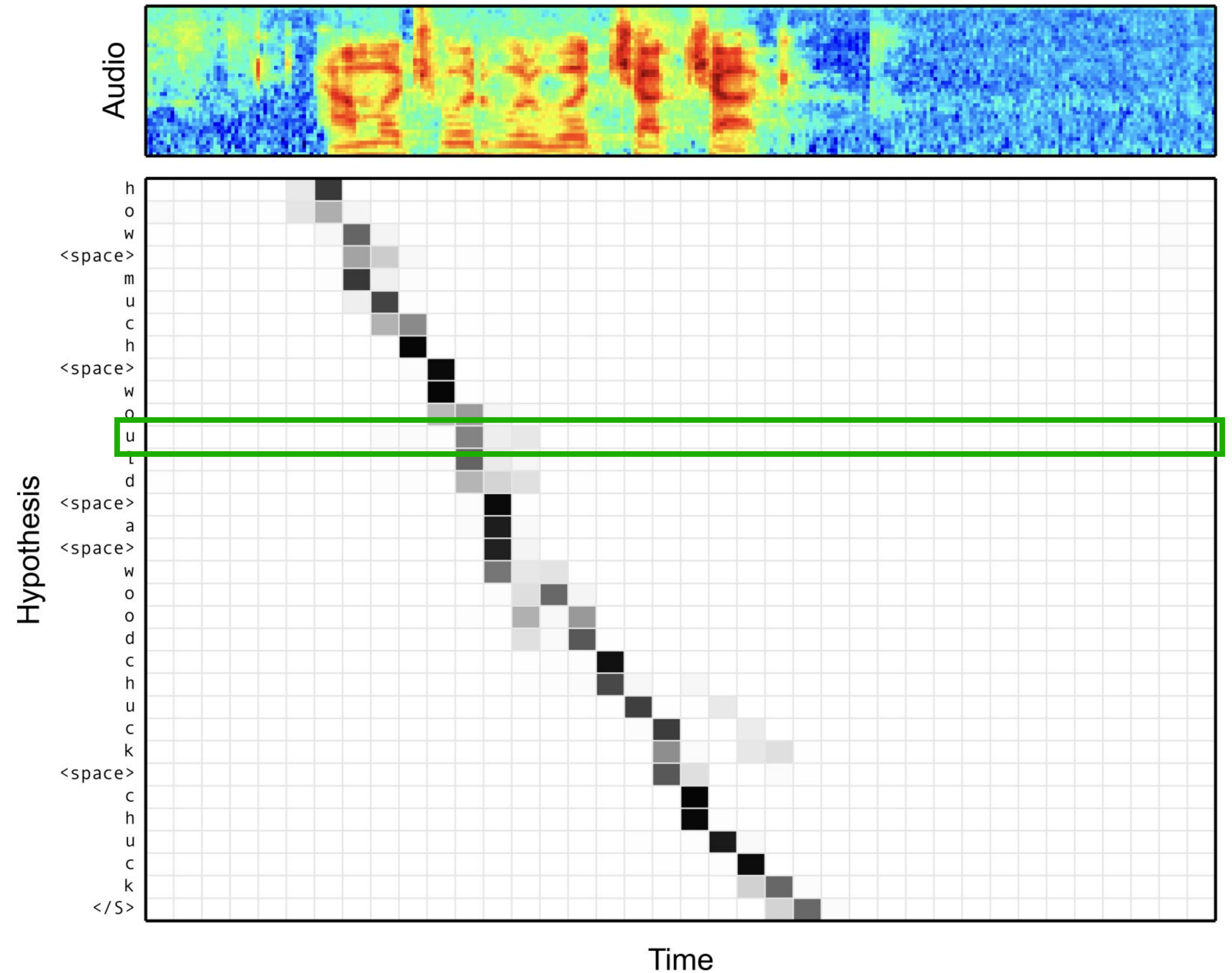
Google: 2017



# Функция смещения

Alignment between the Characters and Audio

Google: 2017

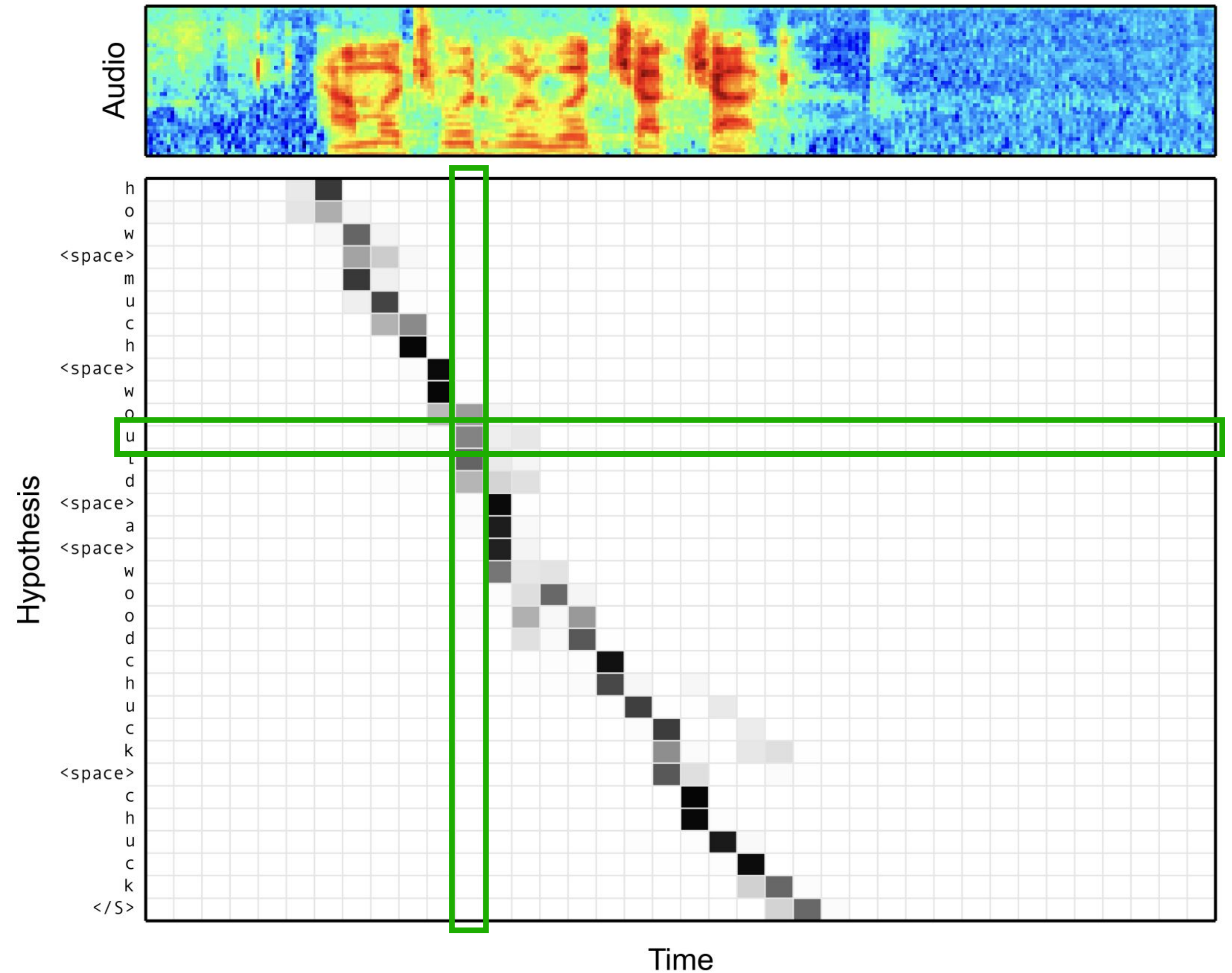




# Функция смещения

Alignment between the Characters and Audio

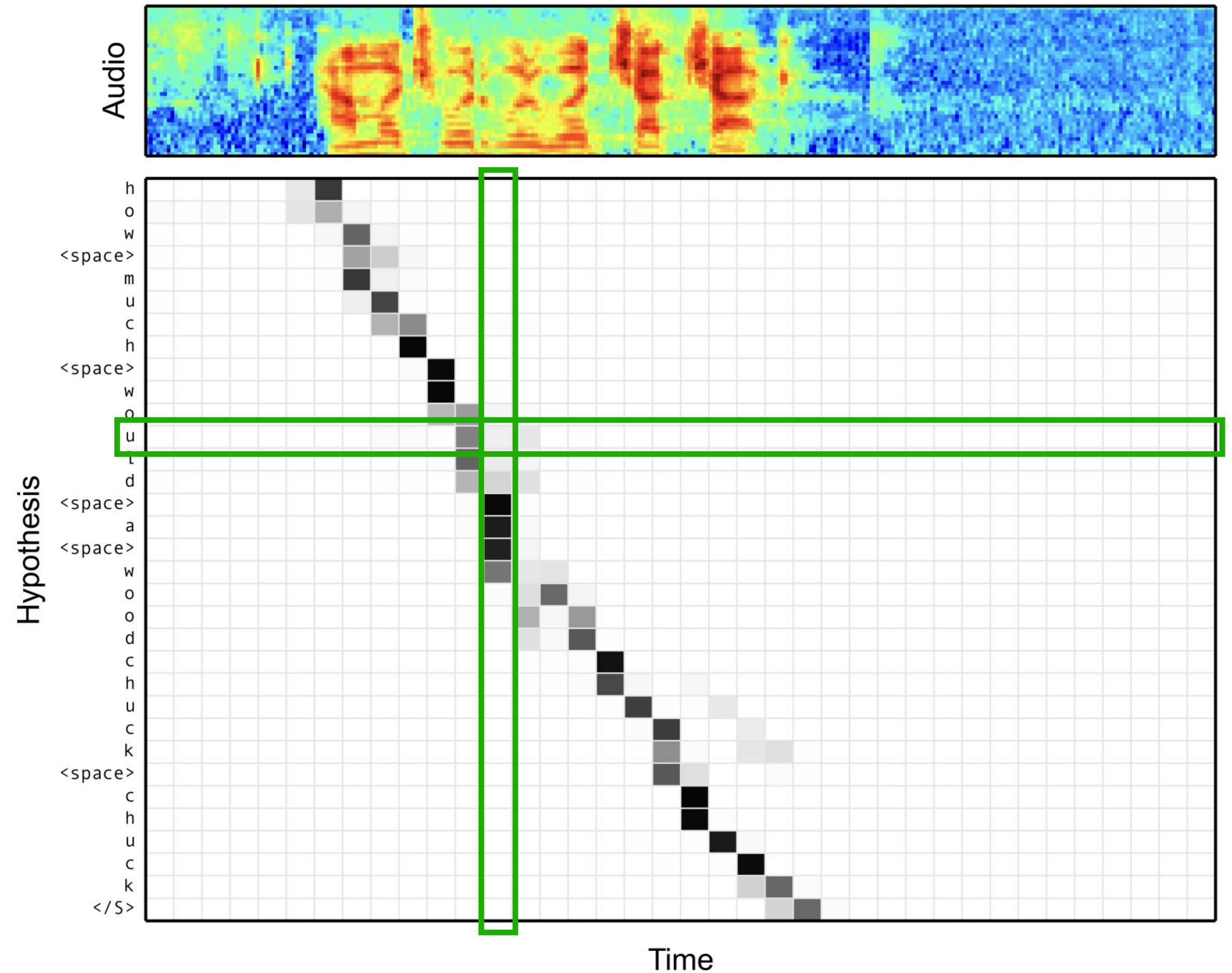
Google: 2017



# Функция смещения

Alignment between the Characters and Audio

Google: 2017

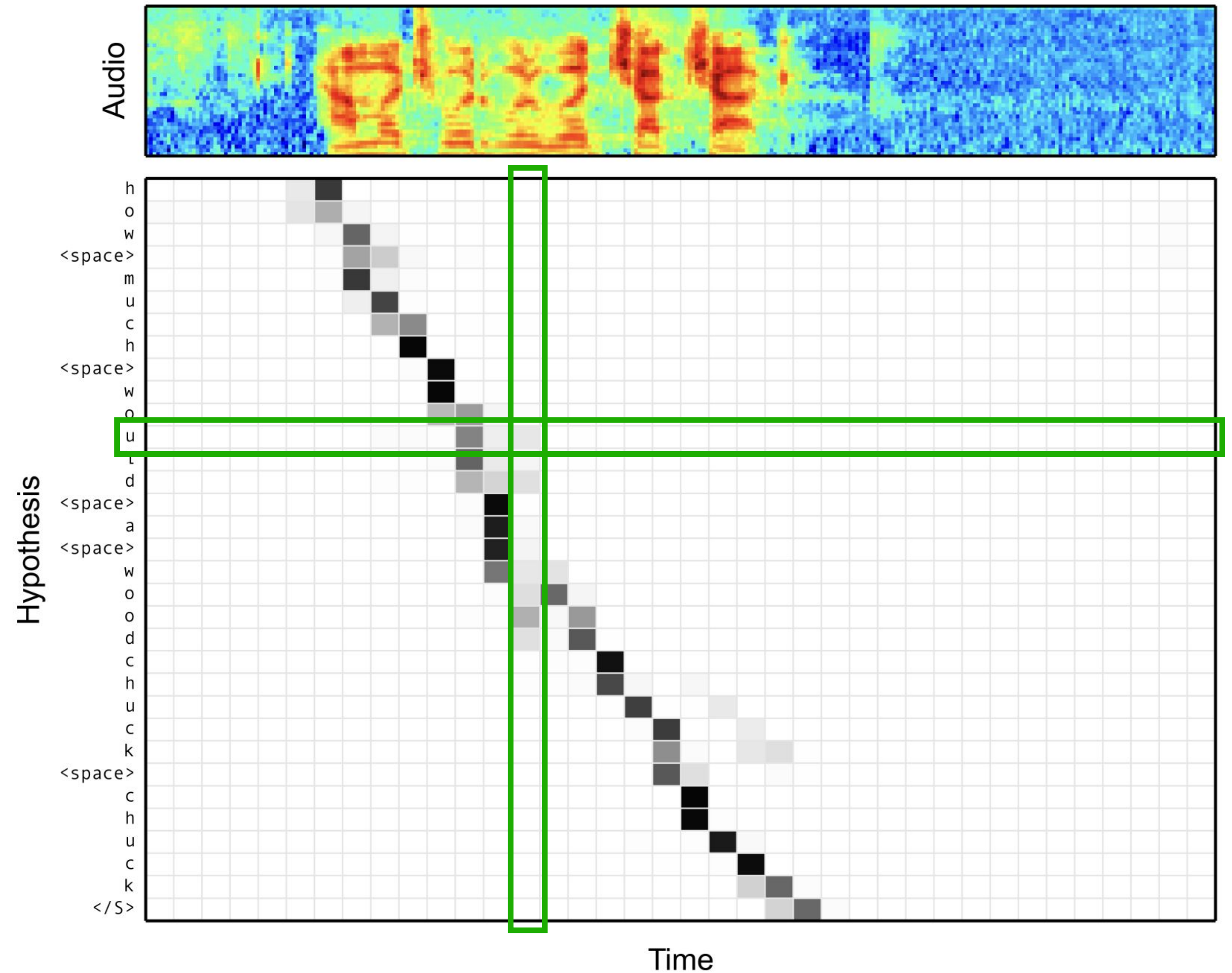




# Функция смешения

Alignment between the Characters and Audio

Google: 2017





# Функция смешения

Google: 2017

$$c(x, y) = \sum_j \log(\min(\sum_i a_{i,j}, 0.5))$$

tokens

frames

attention prob

# Функция смешения

Google: 2017

$$c(x, y) = \sum_j \log(\min(\sum_i a_{i,j}, 0.5))$$

The diagram illustrates the components of the mixing function formula. An arrow labeled "tokens" points to the index  $j$  in the outer summation. An arrow labeled "frames" points to the index  $i$  in the inner summation. A green rectangular box highlights the expression  $\min(\sum_i a_{i,j}, 0.5)$ , with an arrow labeled "attention prob" pointing to it from above.

# Функция смешения

Google: 2017

$$\mathbf{y}^* = \arg \max_y \log p(y|x) + \lambda \log p_{LM}(y) + \gamma c(x, y)$$



# Функция смешения

Google: 2017

$$\mathbf{y}^* = \arg \max_y \log p(y|x) + \lambda \log p_{LM}(y) + \gamma c(x, y)$$

# Функция смешения

Google: 2018

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \log P(\mathbf{y}|\mathbf{x}) + \lambda \log P_{LM}(\mathbf{y}) + \gamma \text{len}(\mathbf{y})$$

# Функция смешения

Google: 2018

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \log P(\mathbf{y}|\mathbf{x}) + \lambda \log P_{LM}(\mathbf{y}) + \gamma \text{len}(\mathbf{y})$$



# Функция смешения

SPMI: 2022

$$\hat{Y} = \arg \max_Y [\log P_{\text{RNN-T}}(Y|X) + \lambda_0 \log P_{\text{ILM}}(Y) + \lambda_1 \log P_{\text{ELM}}(Y) + \beta|Y|]$$

# Функция смешения

SPMI: 2022

Internal LM



$$\hat{Y} = \arg \max_Y [\log P_{\text{RNN-T}}(Y|X) + \lambda_0 \log P_{\text{ILM}}(Y) + \lambda_1 \log P_{\text{ELM}}(Y) + \beta|Y|]$$

# Функция смешения

SPMI: 2022

$$\hat{Y} = \arg \max_Y [\log P_{\text{RNN-T}}(Y|X) + \lambda_0 \log P_{\text{ILM}}(Y) + \lambda_1 \log P_{\text{ELM}}(Y) + \beta|Y|]$$

Diagram illustrating the components of the mixture function:

- Internal LM** (Internal Language Model) points to  $P_{\text{ILM}}(Y)$ .
- external LM** (External Language Model) points to  $P_{\text{ELM}}(Y)$ .



# Функция смешения

SPMI: 2022

$$\hat{Y} = \arg \max_Y [\log P_{\text{RNN-T}}(Y|X) + \lambda_0 \log P_{\text{ILM}}(Y) + \lambda_1 \log P_{\text{ELM}}(Y) + \beta |Y|]$$

Diagram illustrating the components of the mixture function:

- Internal LM (points to  $P_{\text{ILM}}(Y)$ )
- external LM (points to  $P_{\text{ELM}}(Y)$ )
- регуляризация на длину (points to  $\beta |Y|$ )

# Модели адаптации

# Лингвистическая модель (LM)



# Лингвистическая модель (LM)

- Оценивает лог-вероятность луча на каждом шаге

# Лингвистическая модель (LM)

- Оценивает лог-вероятность луча на каждом шаге
- Токенизация не обязательно должна быть как в основной модели

# Лингвистическая модель (LM)

- Оценивает лог-вероятность луча на каждом шаге
- Токенизация не обязательно должна быть как в основной модели
- Модель однонаправленная, учитывается только предыдущий контекст

# Лингвистическая модель (LM)

- Оценивает лог-вероятность луча на каждом шаге
- Токенизация не обязательно должна быть как в основной модели
- Модель однонаправленная, учитывается только предыдущий контекст
- Для быстрого инференса обучают преимущественно маленькие модели



# Типы лингвистических моделей

Аспект	N-Gram	RNN	Transformer

# Типы лингвистических моделей

<i>Аспект</i>	N-Gram	RNN	Transformer
<i>Обучение</i>	Быстро	Медленно	Медленно

# Типы лингвистических моделей

Аспект	N-Gram	RNN	Transformer
Обучение	Быстро	Медленно	Медленно
Инференс	Мгновенно	Медленно	Медленно

# Типы лингвистических моделей

Аспект	N-Gram	RNN	Transformer
Обучение	Быстро	Медленно	Медленно
Инференс	Мгновенно	Медленно	Медленно
Видеопамять	Не требуется	Фиксирована (веса + hidden)	Растёт от длины контекста



# Типы лингвистических моделей

Аспект	N-Gram	RNN	Transformer
Обучение	Быстро	Медленно	Медленно
Инференс	Мгновенно	Медленно	Медленно
Видеопамять	Не требуется	Фиксирована (веса + hidden)	Растёт от длины контекста
Качество	Среднее	Высокое	Высокое

# Типы лингвистических моделей

Аспект	N-Gram	RNN	Transformer
Обучение	Быстро	Медленно	Медленно
Инференс	Мгновенно	Медленно	Медленно
Видеопамять	Не требуется	Фиксирована (веса + hidden)	Растёт от длины контекста
Качество	Среднее	Высокое	Высокое

# Типы лингвистических моделей

Аспект	N-Gram	RNN	Transformer
Обучение	Быстро	Медленно	Медленно
Инференс	Мгновенно	Медленно	Медленно
Видеопамять	Не требуется	Фиксирована (веса + hidden)	Растёт от длины контекста
Качество	Среднее	Высокое	Высокое

# N-Gram LM



# N-Gram LM

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

# N-Gram LM

<s> Yes, I am watching the TV right now </s>

<s> The laptop I gave you is very good </s>

<s> I am very happy to be here </s>

$$P(w_i | w_{i-1}) = \frac{\textit{count}(w_{i-1}, w_i)}{\textit{count}(w_{i-1})}$$

# N-Gram LM

<s> Yes, I am watching the TV right now </s>

<s> The laptop I gave you is very good </s>

<s> I am very happy to be here </s>

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

$$P(\text{am} | \text{I}) = \text{count}(\text{"I am"}) / \text{count}(\text{"I"}) = 2 / 3 = 0.667$$



# N-Gram LM

<s> Yes, I am watching the TV right now </s>

<s> The laptop I gave you is very good </s>

<s> I am very happy to be here </s>

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

$$P(\text{am} | \text{I}) = \text{count}(\text{"I am"}) / \text{count}(\text{"I"}) = 2 / 3 = 0.667$$

$$P(\text{I} | \text{<s>}) = \text{count}(\text{"<s> I"}) / \text{count}(\text{"<s>"}) = 1 / 3 = 0.333$$



# N-Gram LM

<s> Yes, I am watching the TV right now </s>

<s> The laptop I gave you is very good </s>

<s> I am very happy to be here </s>

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

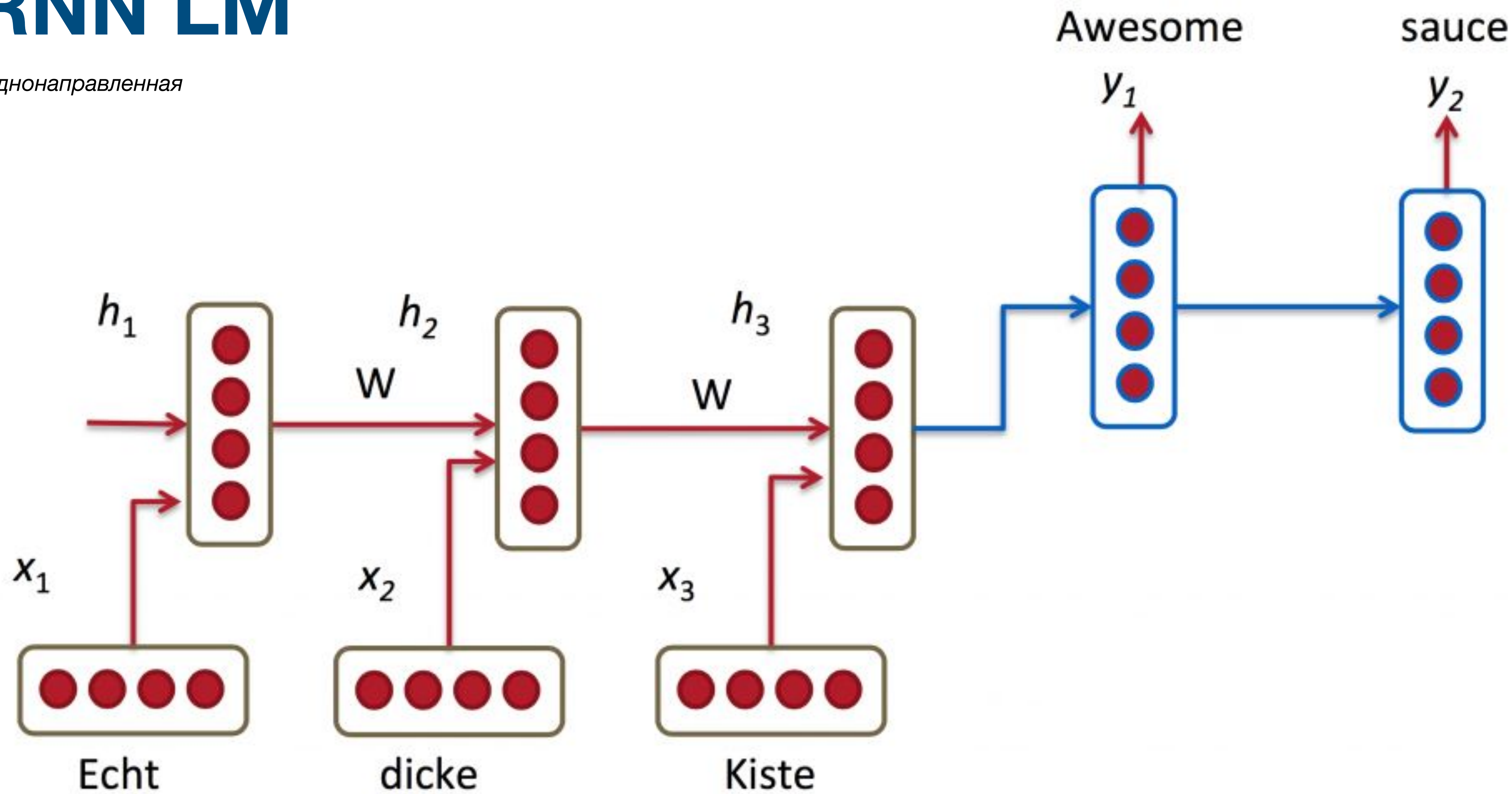
$$P(\text{am} | \text{I}) = \text{count}(\text{"I am"}) / \text{count}(\text{"I"}) = 2 / 3 = 0.667$$

$$P(\text{I} | \text{<s>}) = \text{count}(\text{"<s> I"}) / \text{count}(\text{"<s>"}) = 1 / 3 = 0.333$$

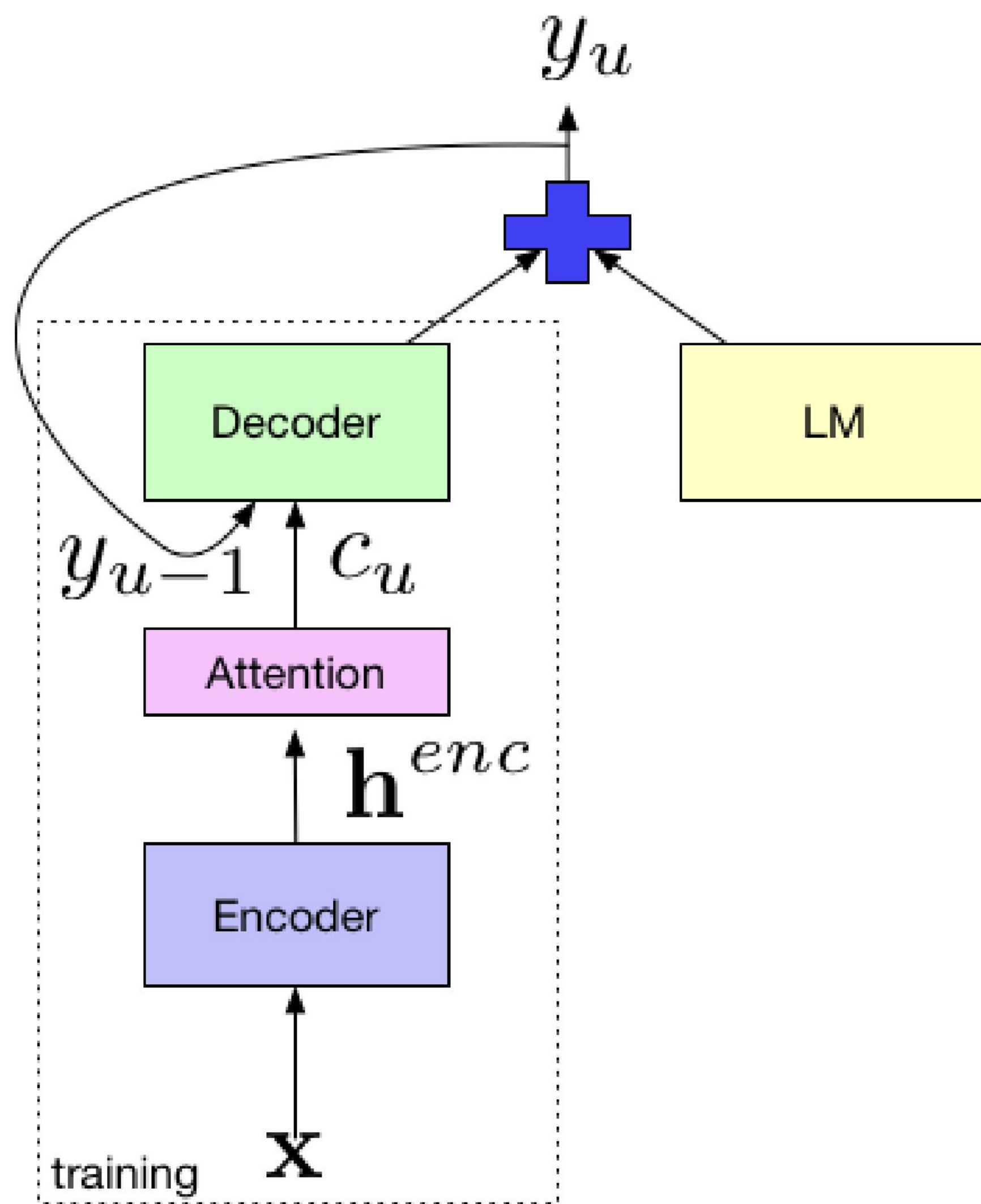
$$P(\text{now} | \text{right}) = \text{count}(\text{"right now"}) / \text{count}(\text{"right"}) = 1 / 1 = 1$$

# RNN LM

Однонаправленная

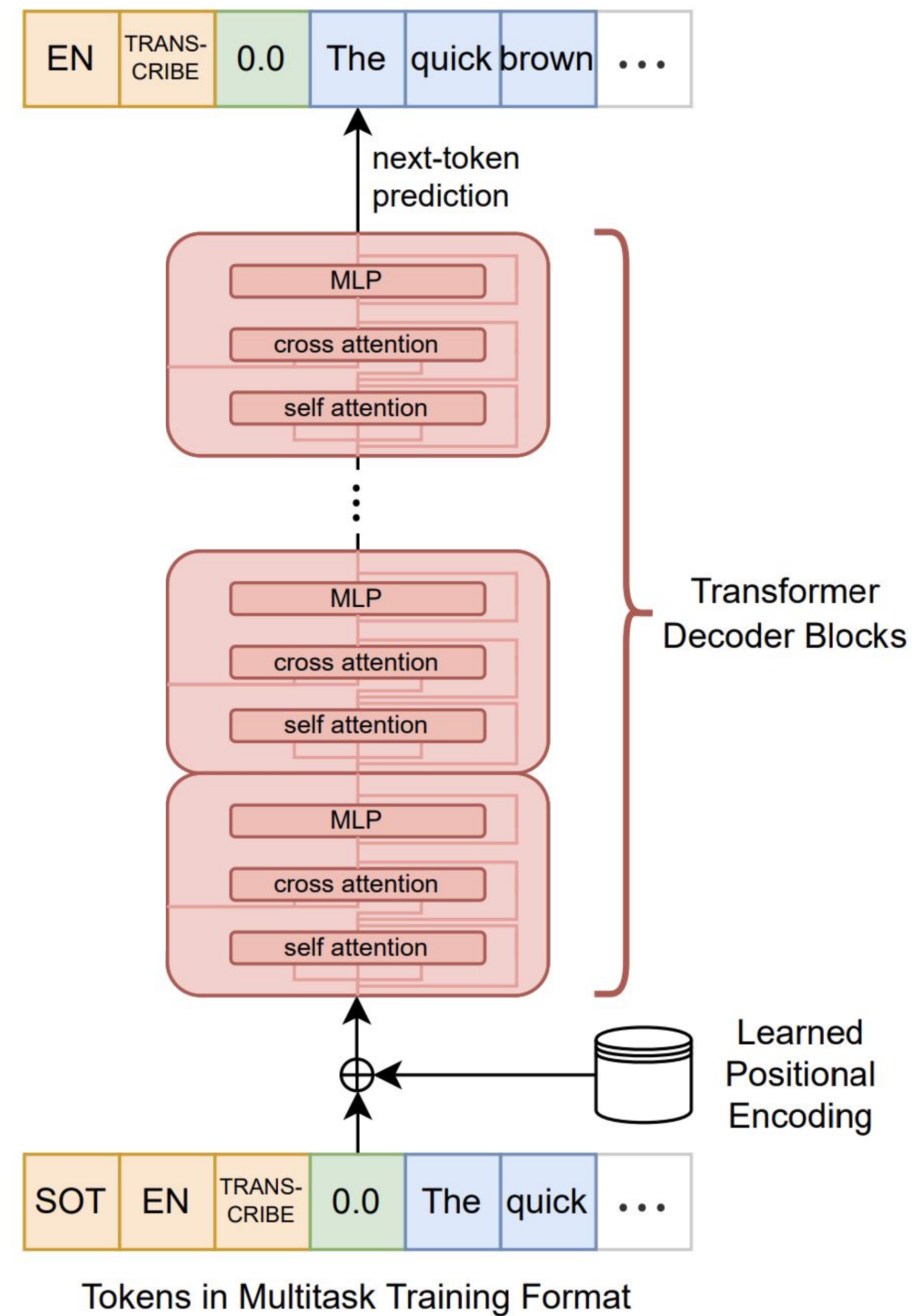
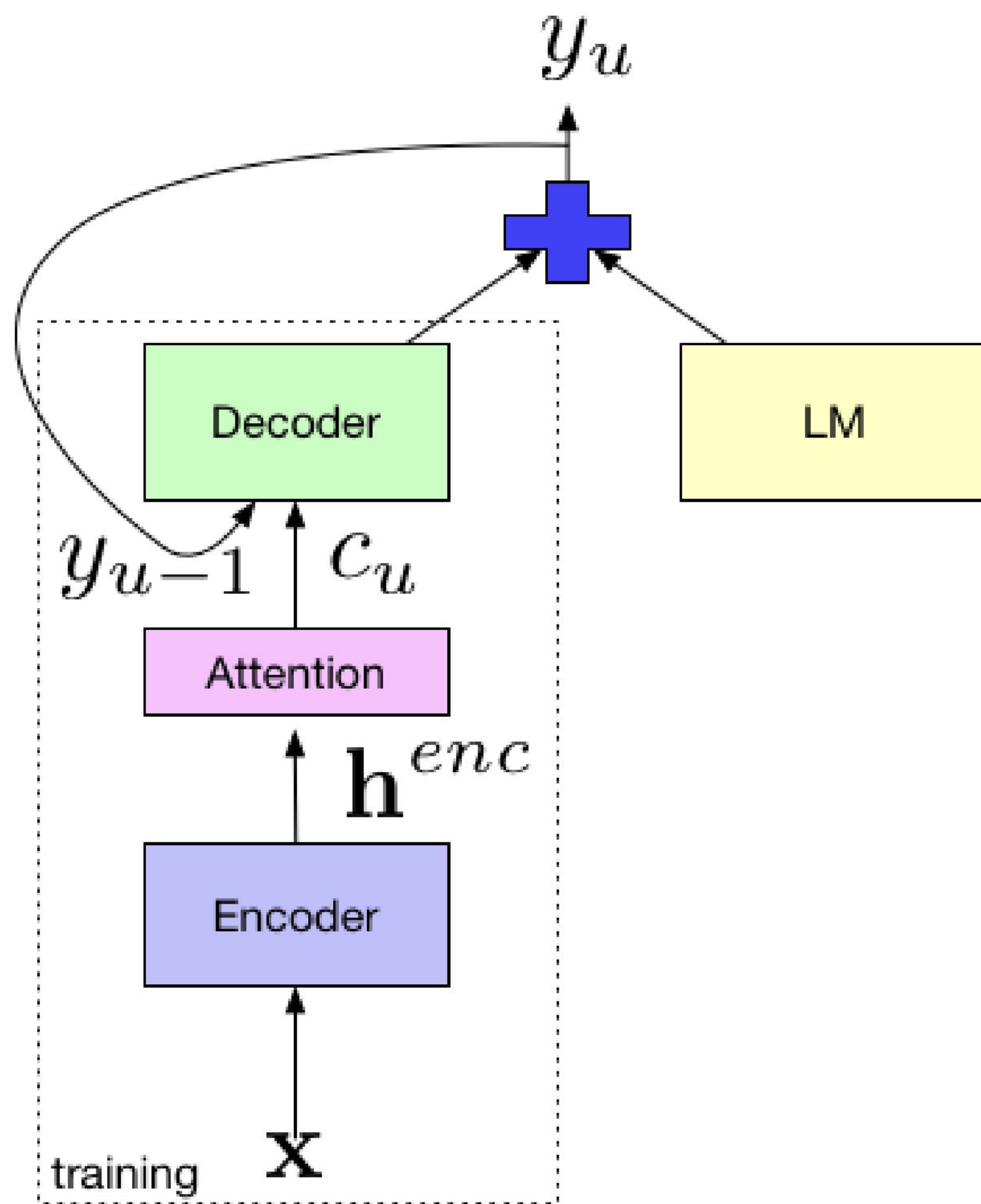


# Инференс LM



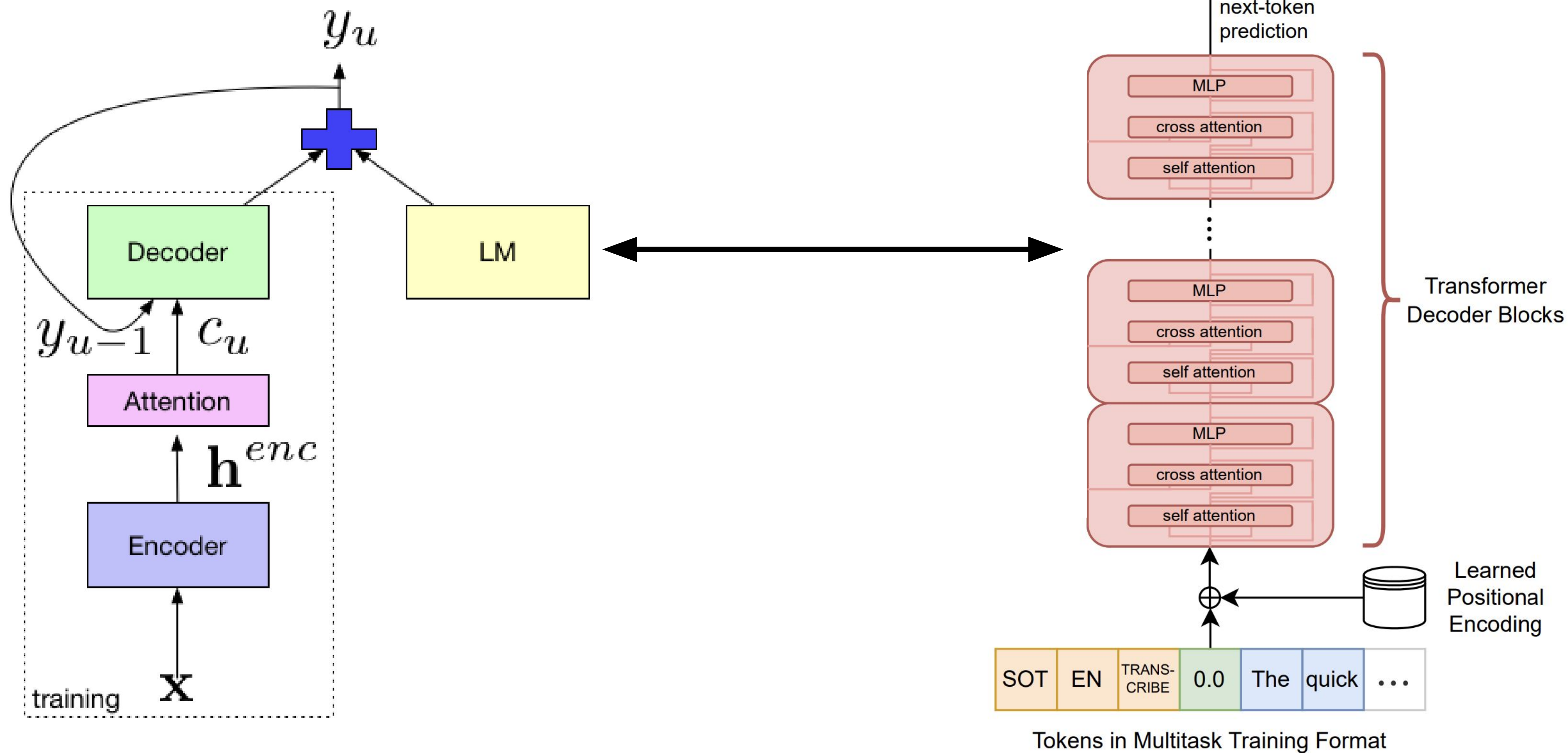


# Инференс LM

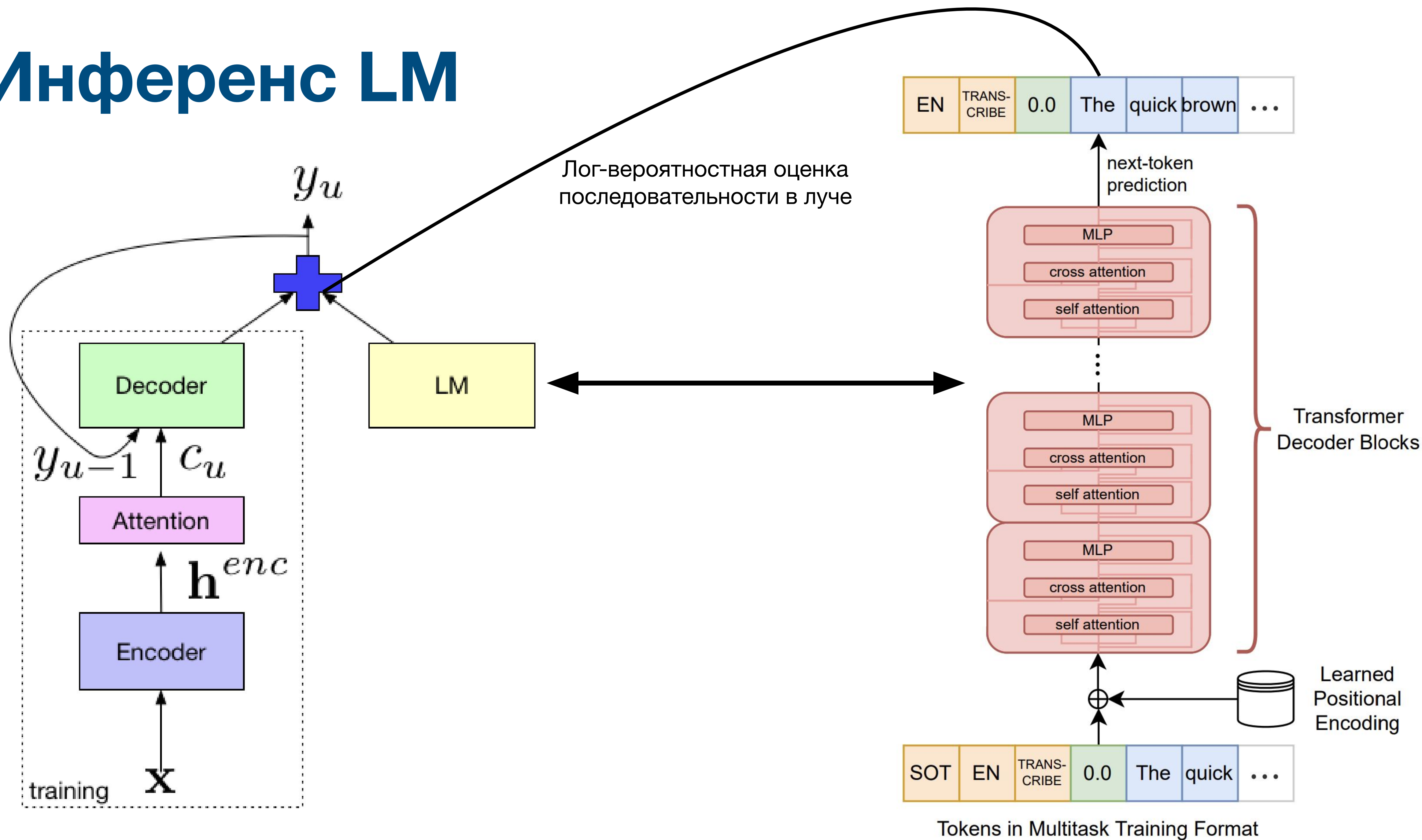




# Инференс LM

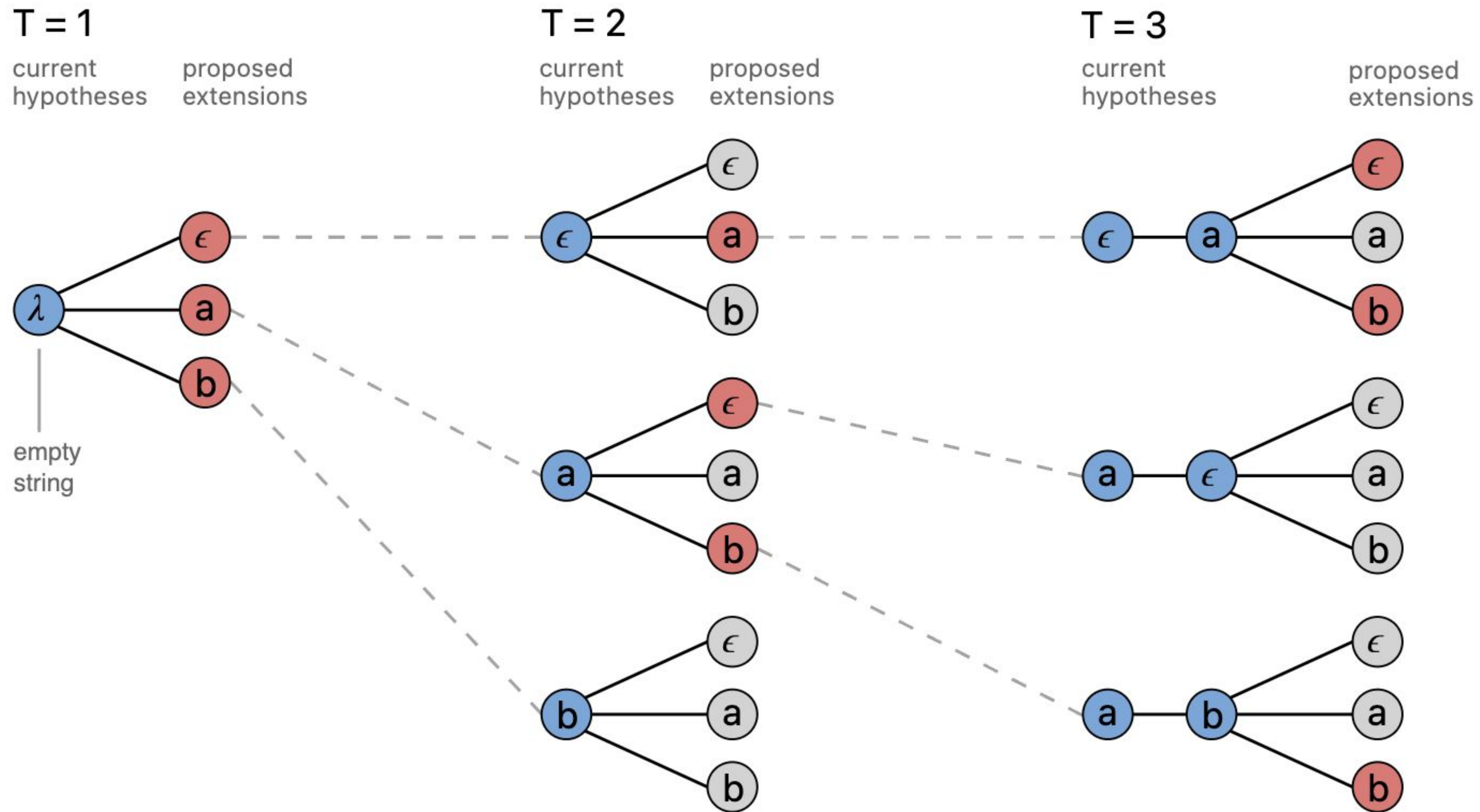


# Инференс LM





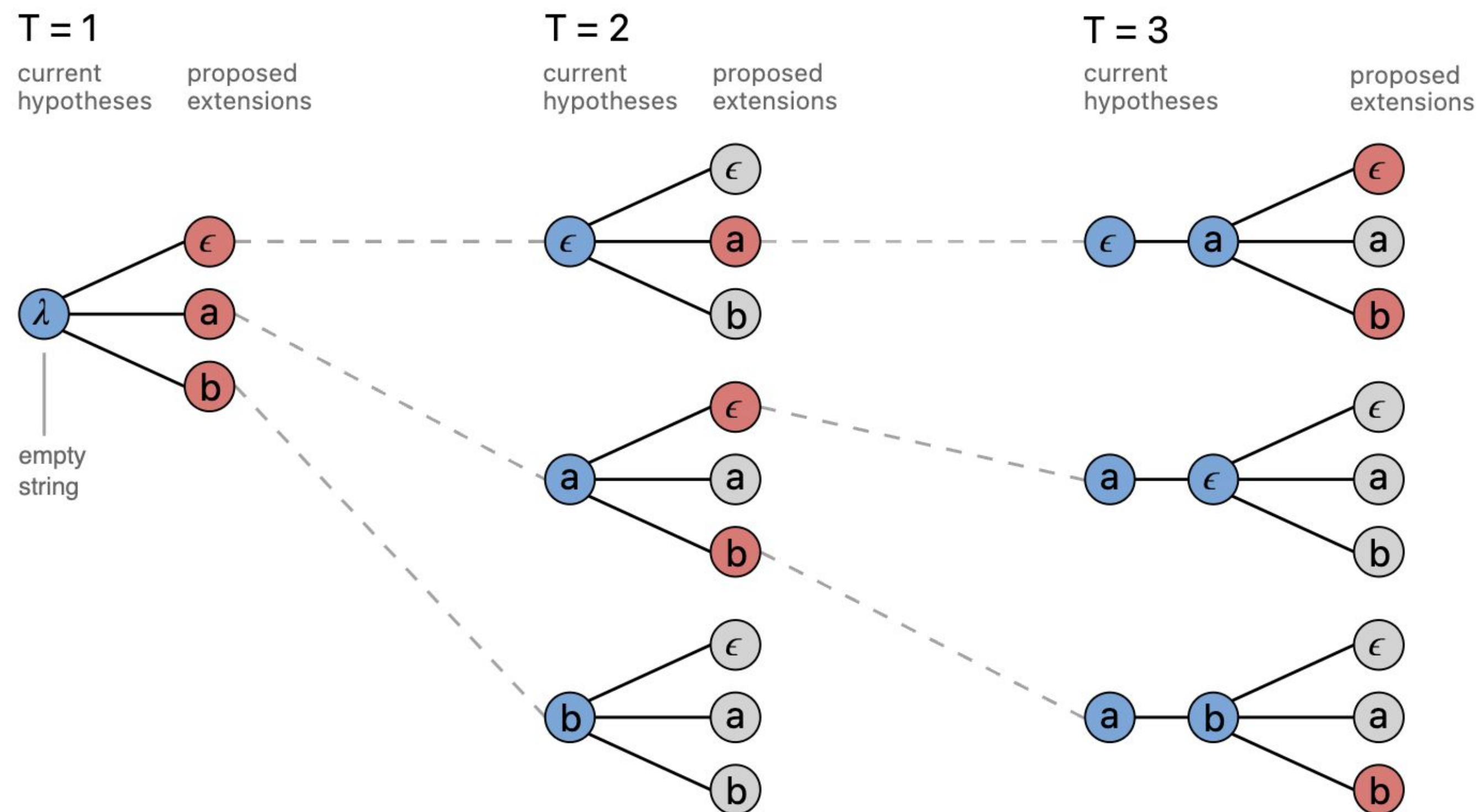
# Инференс LM



A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

# Инференс LM

## N-Gram



A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.



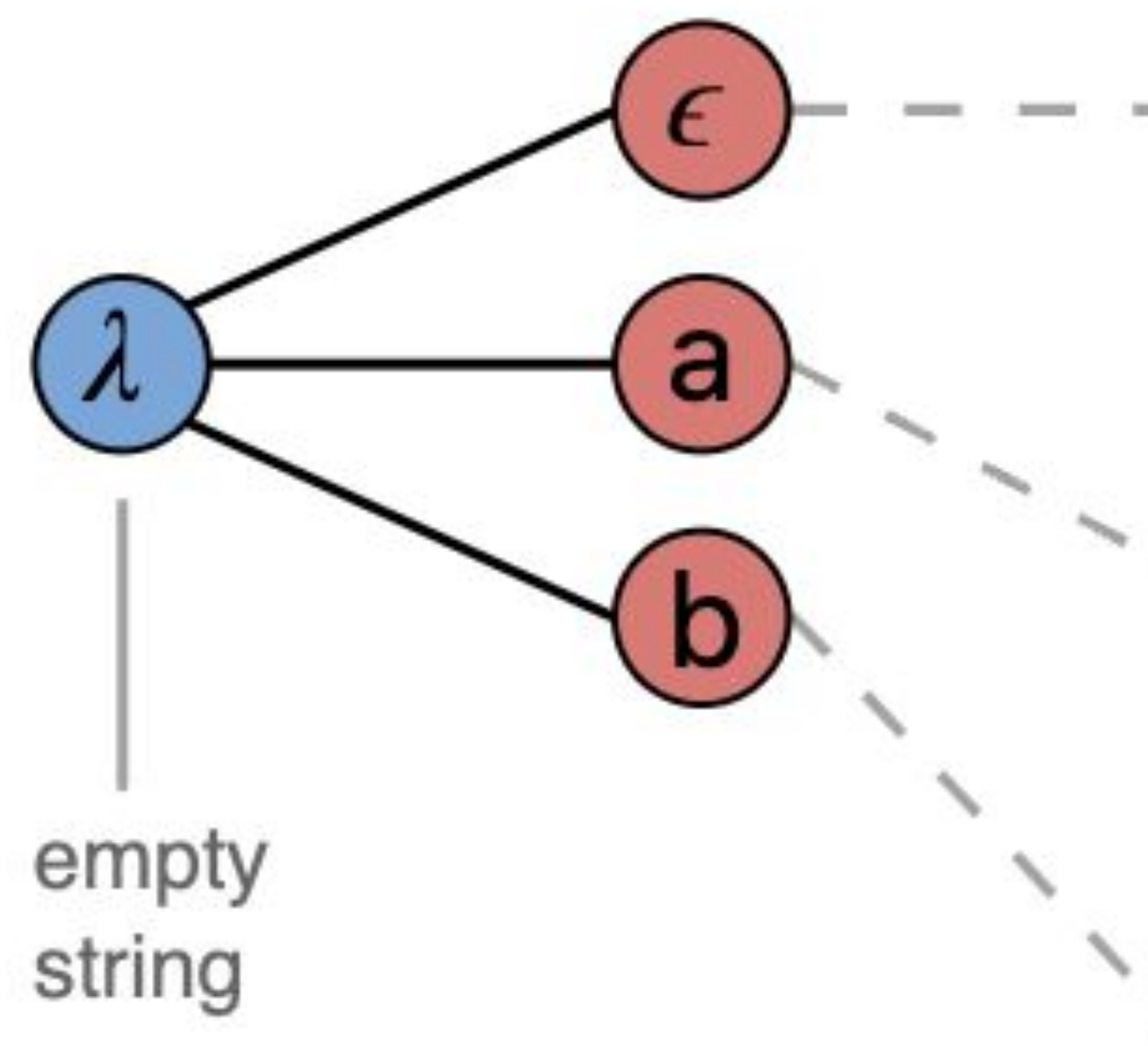
# Инференс LM

## N-Gram

$T = 1$

current  
hypotheses

proposed  
extensions



# Инференс LM

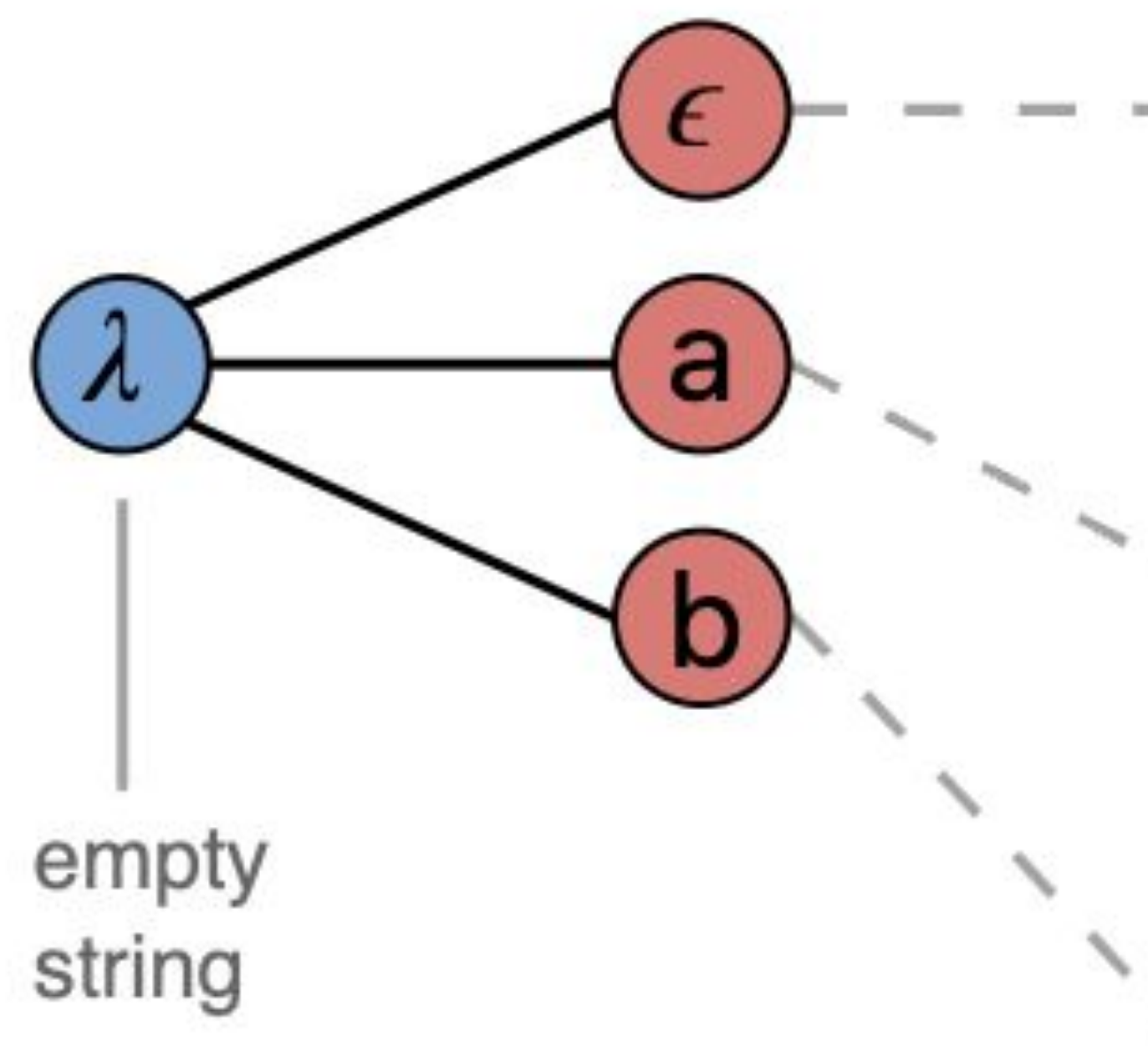
## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

$T = 1$

current  
hypotheses

proposed  
extensions



# Инференс LM

## N-Gram

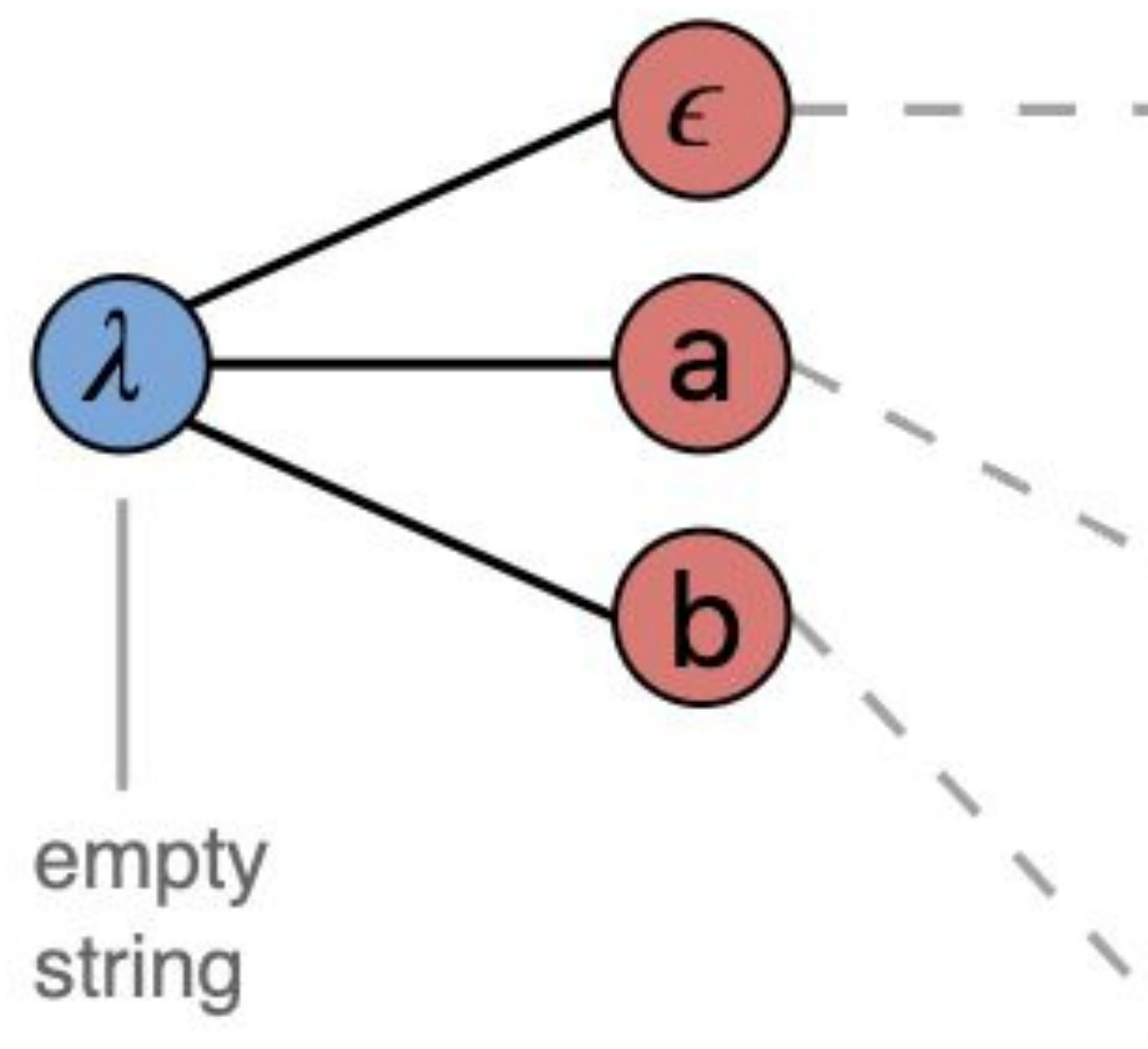
A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

$N=3$

$T = 1$

current  
hypotheses

proposed  
extensions



# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

$N=3$

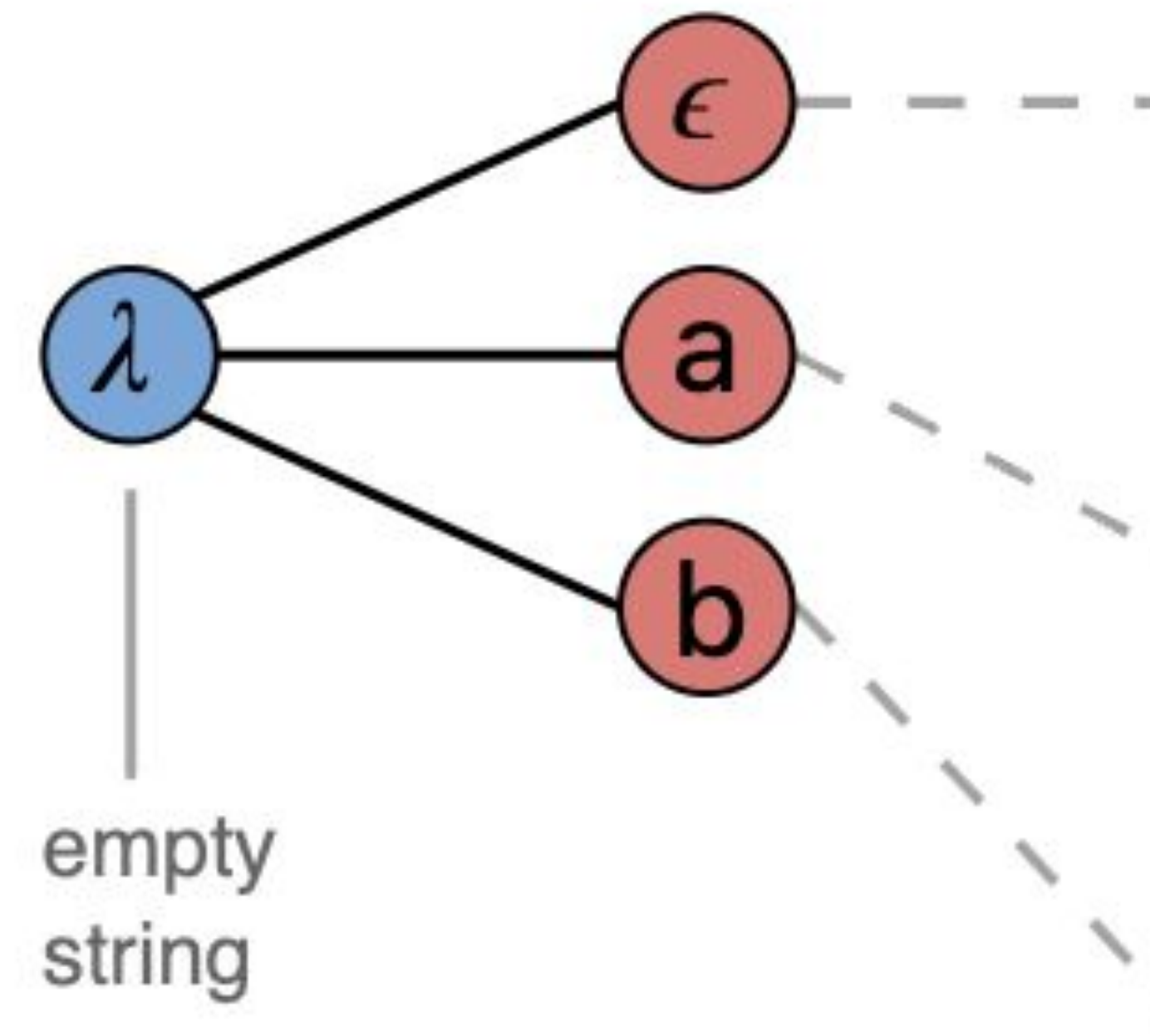
Шаг 1: [ $\langle \text{SOS} \rangle$ ]

$p(\text{blank}) = p(\text{blank}) * p(\text{blank} \mid \langle \text{SOS} \rangle)$

$T = 1$

current  
hypotheses

proposed  
extensions





# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

$N=3$

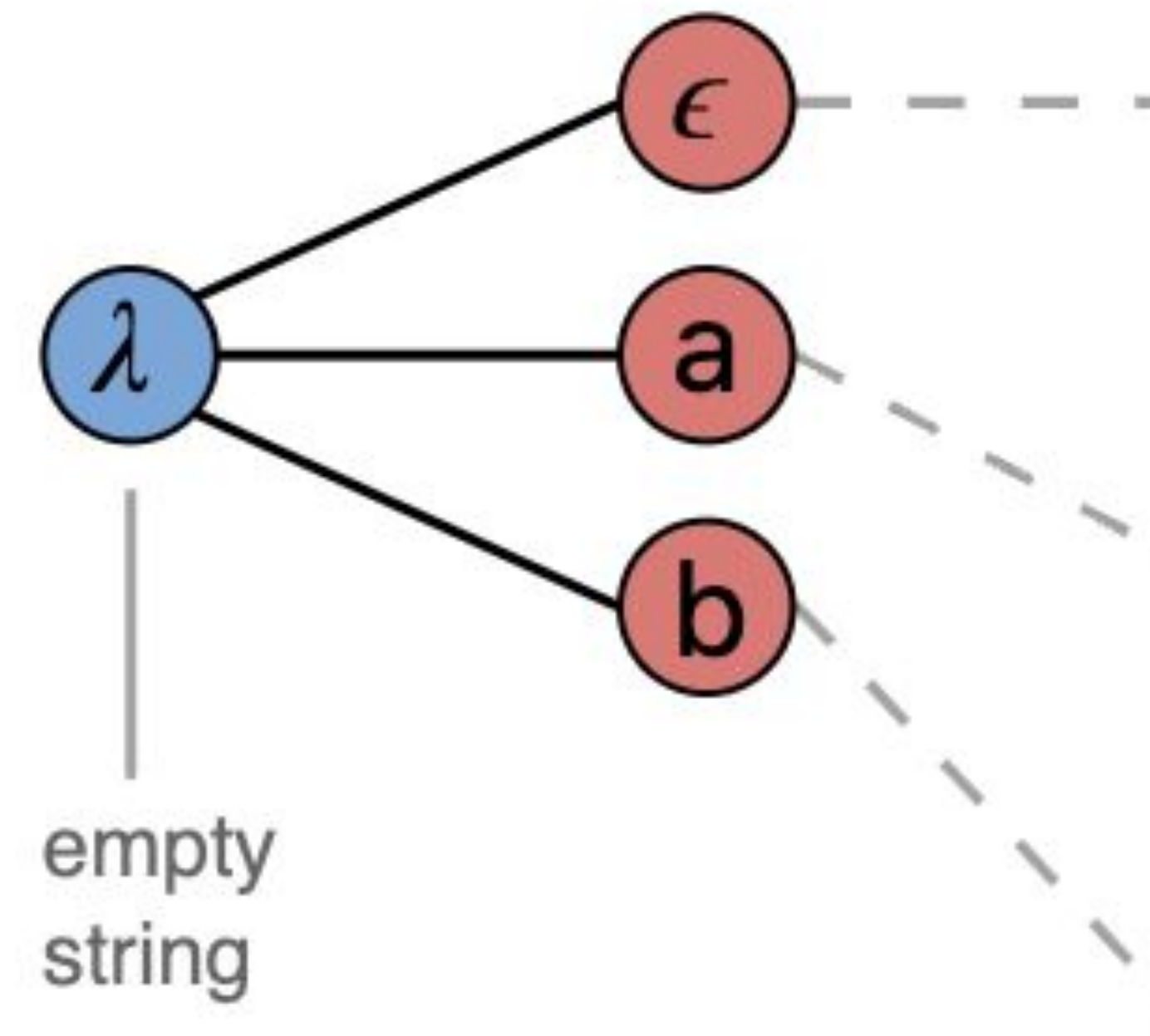
Шаг 1: [ $\langle \text{SOS} \rangle$ ]

$p(\text{blank}) = p(\text{blank}) * p(\text{blank} \mid \langle \text{SOS} \rangle)$

$T = 1$

current  
hypotheses

proposed  
extensions



# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

N=3

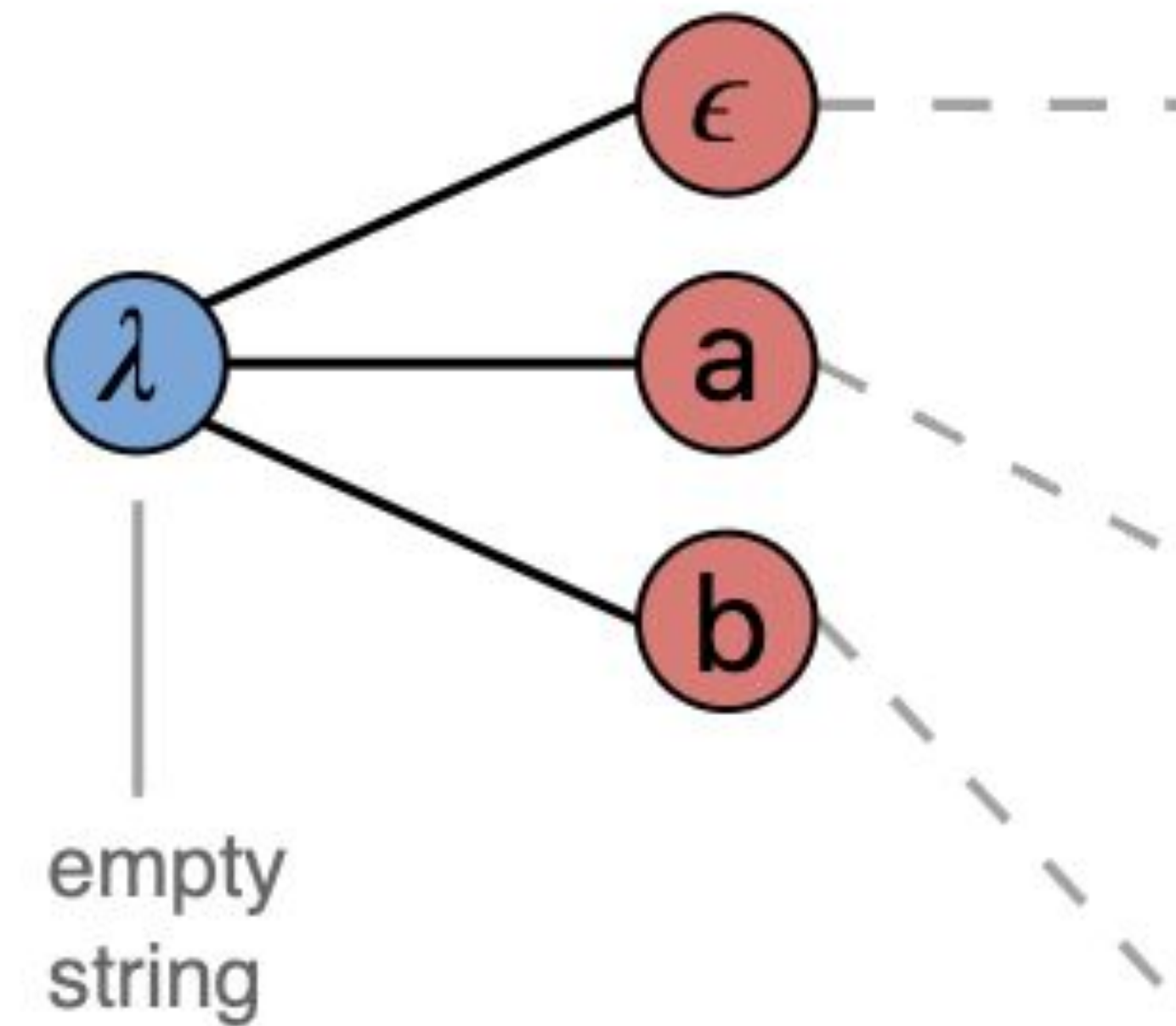
Шаг 1: [ $\epsilon$ ]  
 $p(\text{blank}) = p(\text{blank}) * p(\text{blank} | \epsilon)$

1-gram      2-gram

T = 1

current  
hypotheses

proposed  
extensions



# Инференс LM

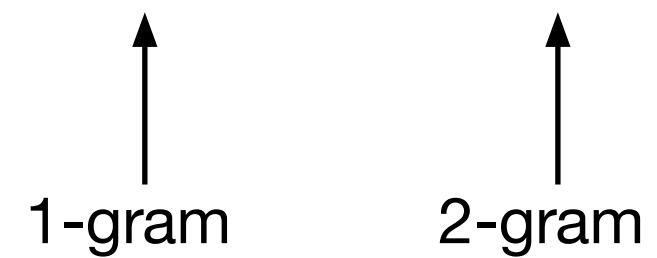
## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

N=3

Шаг 1: [ $\langle \text{SOS} \rangle$ ]

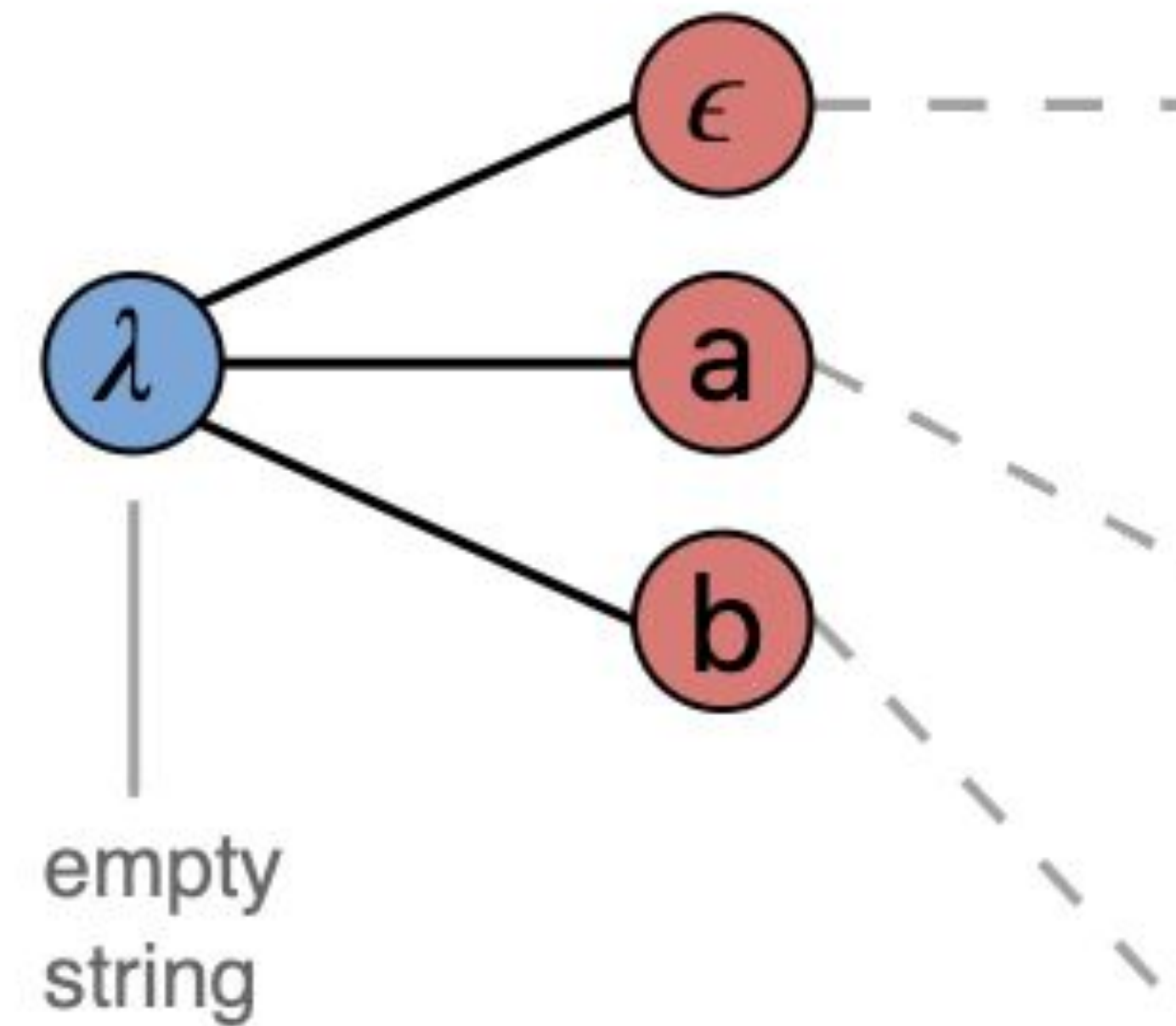
$$\begin{aligned} p(\text{blank}) &= p(\text{blank}) * p(\text{blank} \mid \langle \text{SOS} \rangle) \\ p(a) &= p(a) * p(a \mid \langle \text{SOS} \rangle) \end{aligned}$$



$T = 1$

current  
hypotheses

proposed  
extensions



# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

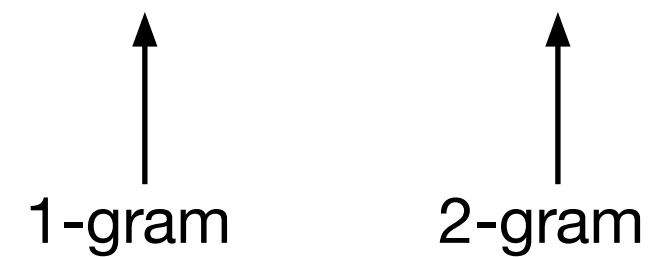
N=3

Шаг 1: [ $\langle \text{SOS} \rangle$ ]

$p(\text{blank}) = p(\text{blank}) * p(\text{blank} \mid \langle \text{SOS} \rangle)$

$p(a) = p(a) * p(a \mid \langle \text{SOS} \rangle)$

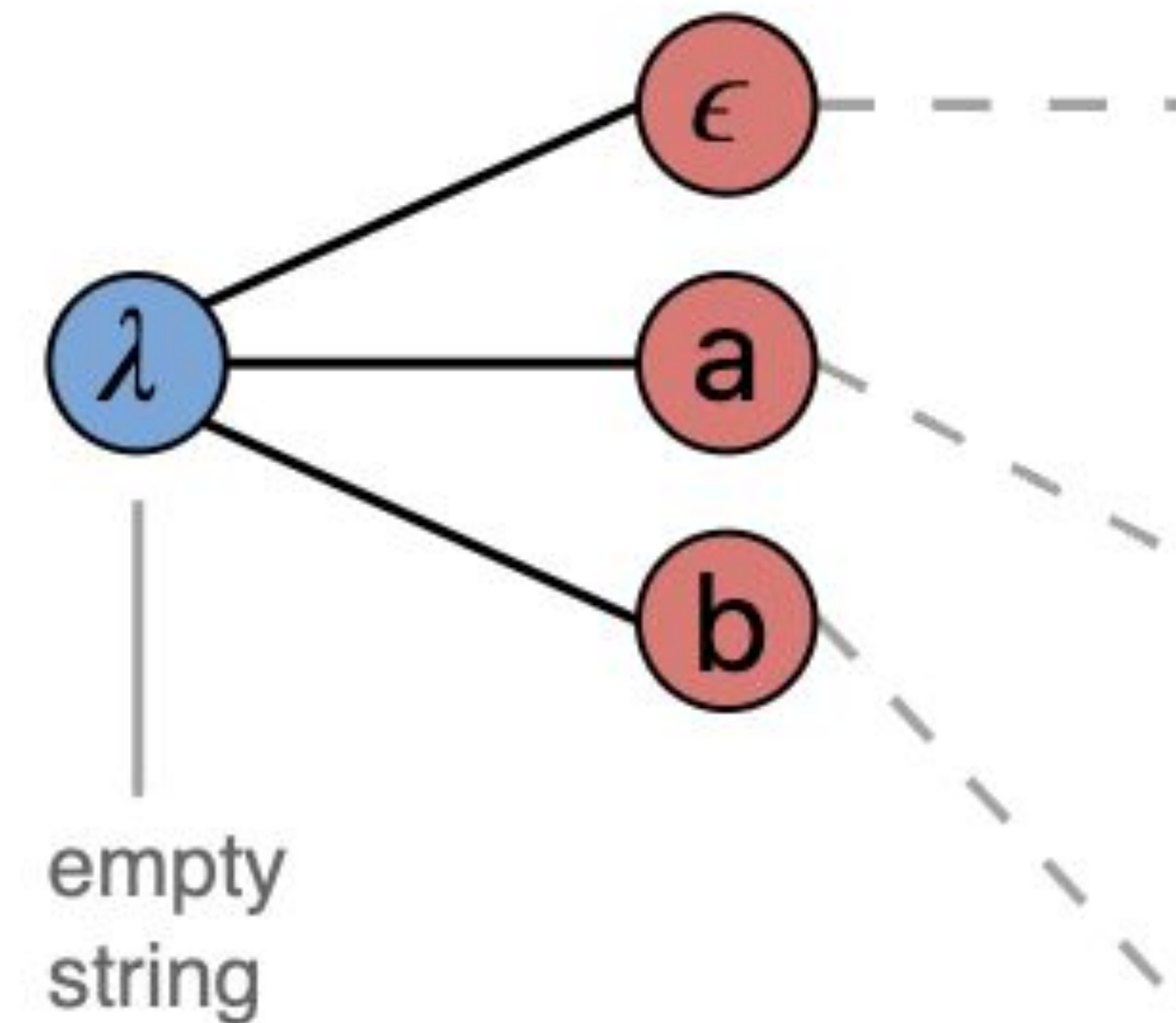
$p(b) = p(b) * p(b \mid \langle \text{SOS} \rangle)$



$T = 1$

current  
hypotheses

proposed  
extensions





# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

N=3

Шаг 1: [ $\epsilon$ ]

$$p(\text{blank}) = p(\text{blank}) * p(\text{blank} \mid \epsilon)$$

$$p(a) = p(a) * p(a \mid \epsilon)$$

$$p(b) = p(b) * p(b \mid \epsilon)$$

↑  
1-gram

↑  
2-gram

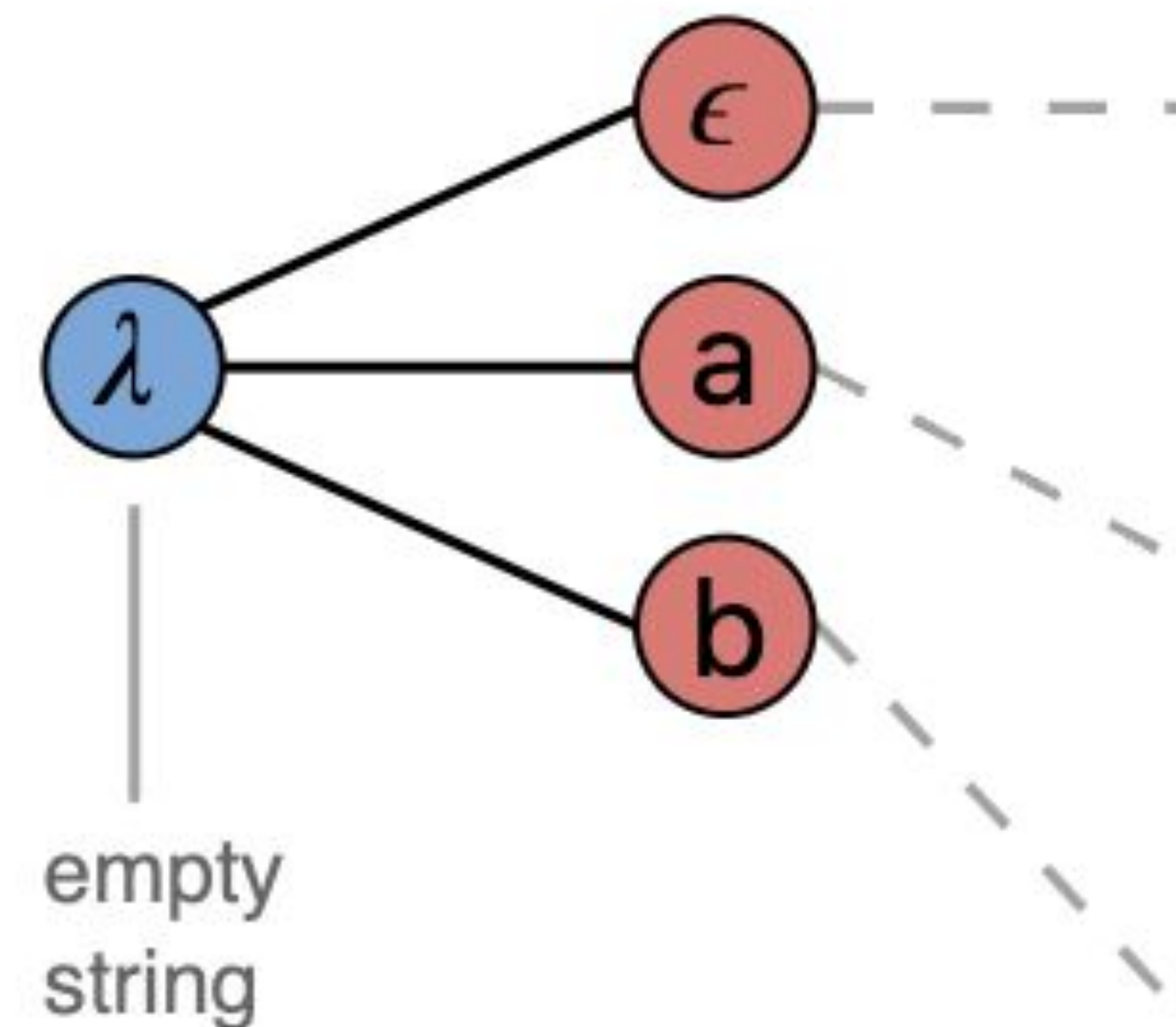
Шаг 1: [ $\epsilon$ ]

$$\log(p(\text{blank})) = \log(p(\text{blank})) + \log(p(\text{blank} \mid \epsilon))$$

T = 1

current  
hypotheses

proposed  
extensions



# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

N=3

Шаг 1: [ $\epsilon$ ]

$$p(\text{blank}) = p(\text{blank}) * p(\text{blank} \mid \epsilon)$$

$$p(a) = p(a) * p(a \mid \epsilon)$$

$$p(b) = p(b) * p(b \mid \epsilon)$$

↑  
1-gram

↑  
2-gram

Шаг 1: [ $\epsilon$ ]

$$\log(p(\text{blank})) = \log(p(\text{blank})) + \log(p(\text{blank} \mid \epsilon))$$

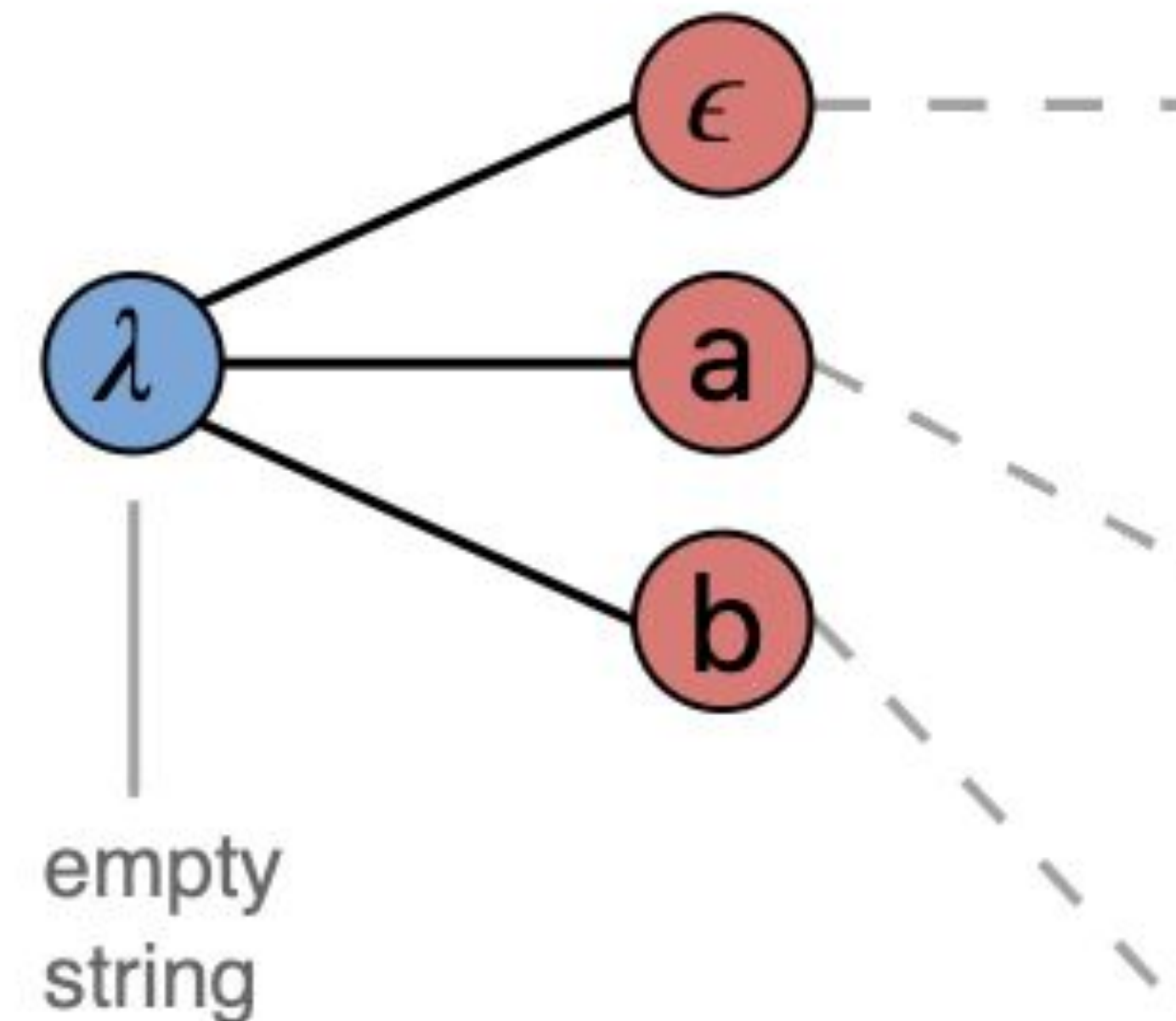
Общий вид:

$$\log(p(t)) = \text{Sum}(\log(p(t \mid t-1, t-2, \dots, t-N+1)))$$

T = 1

current  
hypotheses

proposed  
extensions



# Инференс LM

## N-Gram

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

N=3

Шаг 2: [ $\epsilon$ , blank]

$p(\text{blank}) = p(\text{blank}) * p(\text{blank} \mid \text{blank}) * p(\text{blank} \mid \text{blank}, \epsilon)$

$p(a) = p(a) * p(a \mid \text{blank}) * p(a \mid \text{blank}, \epsilon)$

$p(b) = p(b) * p(b \mid \text{blank}) * p(b \mid \text{blank}, \epsilon)$

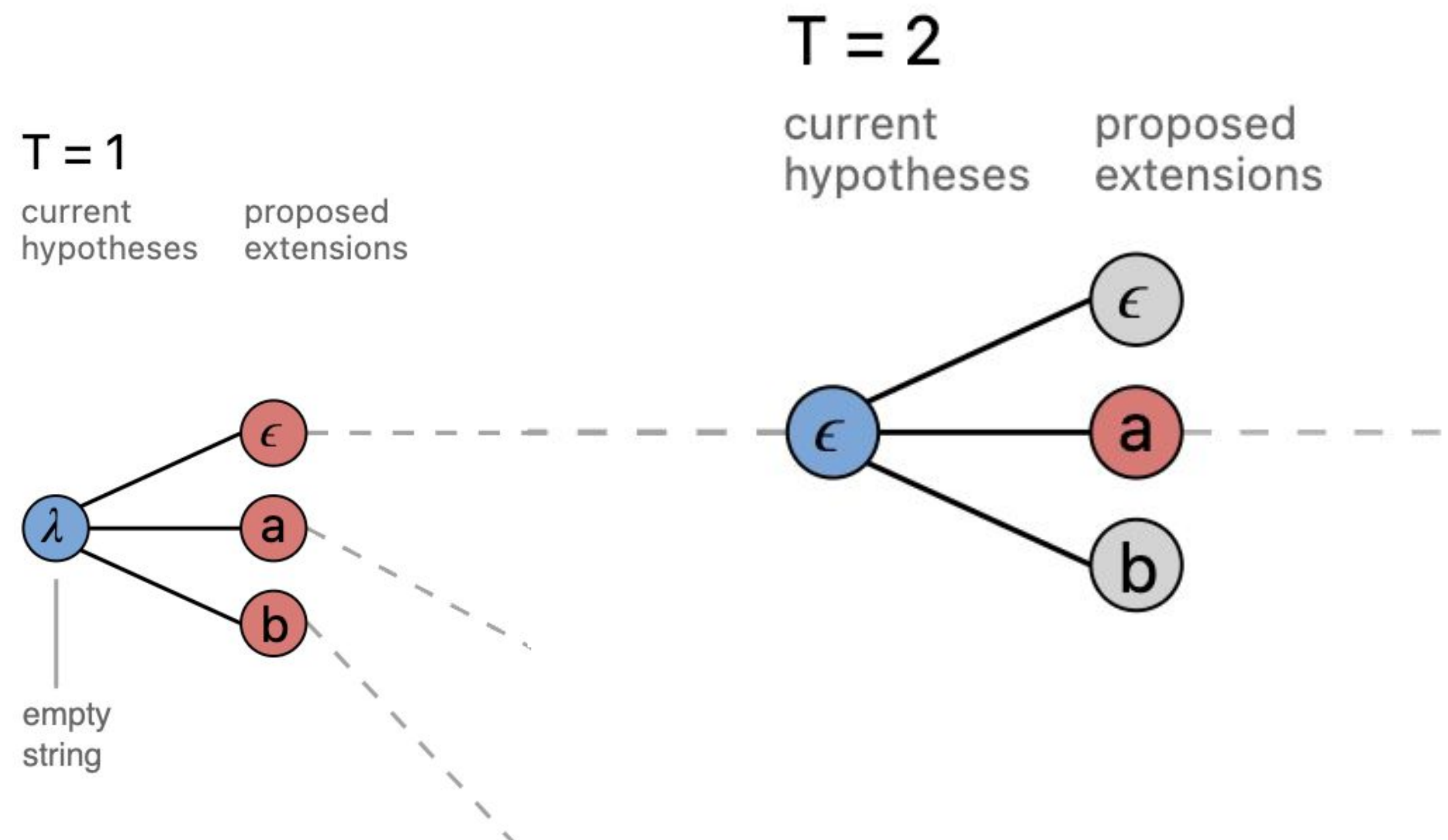
1-gram      2-gram      3-gram

Шаг 2: [ $\epsilon$ , blank]

$\log(p(a)) = \log(p(a)) + \log(p(a \mid \text{blank})) + \log(p(a \mid \text{blank}, \epsilon))$

Общий вид:

$\log(p(t)) = \text{Sum}(\log(p(t \mid t-1, t-2, \dots, t-N+1)))$

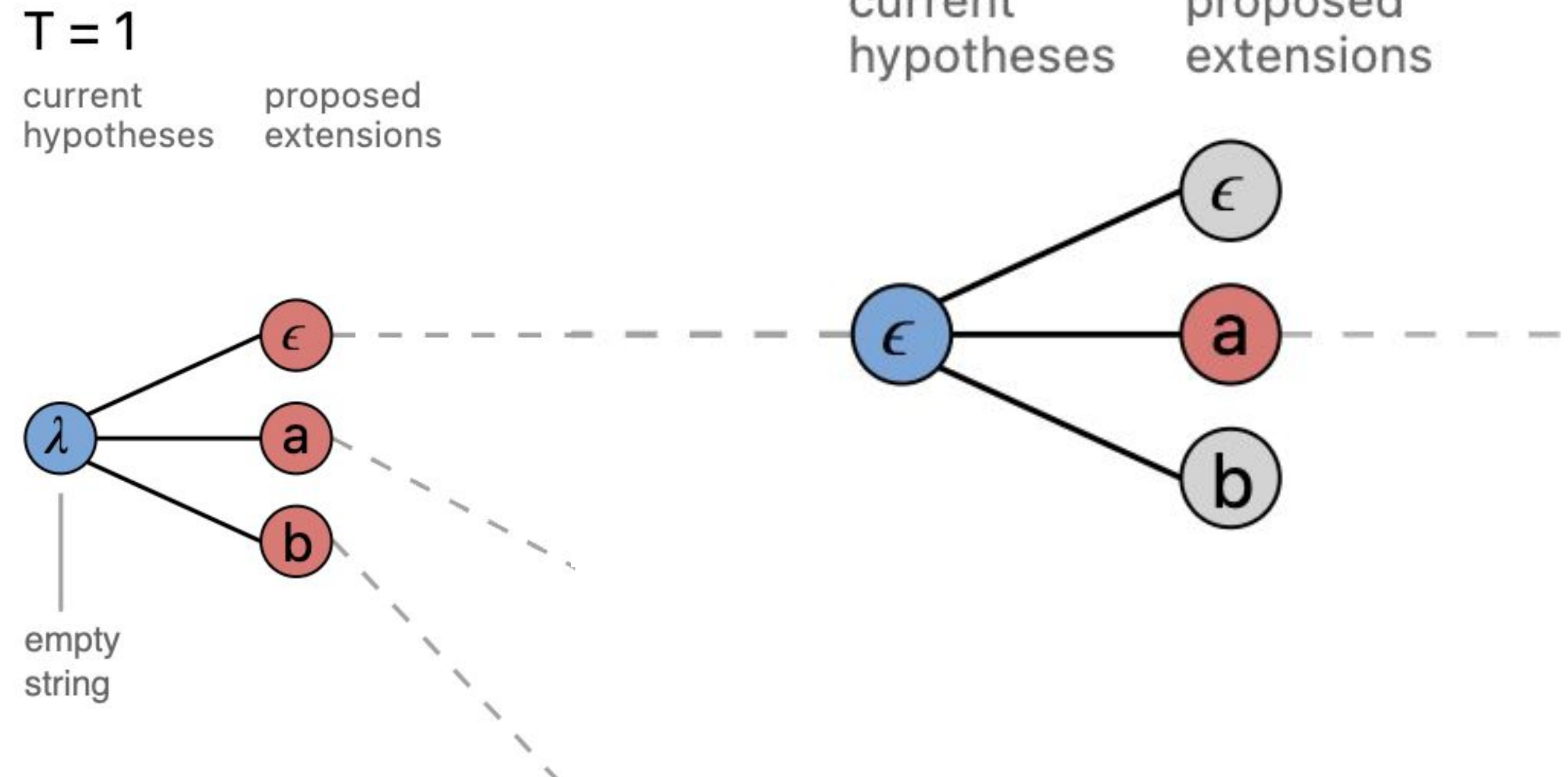


# Инференс LM

## RNN

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

Step 1: [ $\langle \text{SOS} \rangle$ ]  
hidden =  $h_0$  (zeros)  
log-score-1,  $h_1 = \text{log\_softmax}(\text{RNN}(\langle \text{SOS} \rangle, h_0))$





# Инференс LM

## RNN

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

Step 2: [ $\langle \text{SOS} \rangle$ , blank]

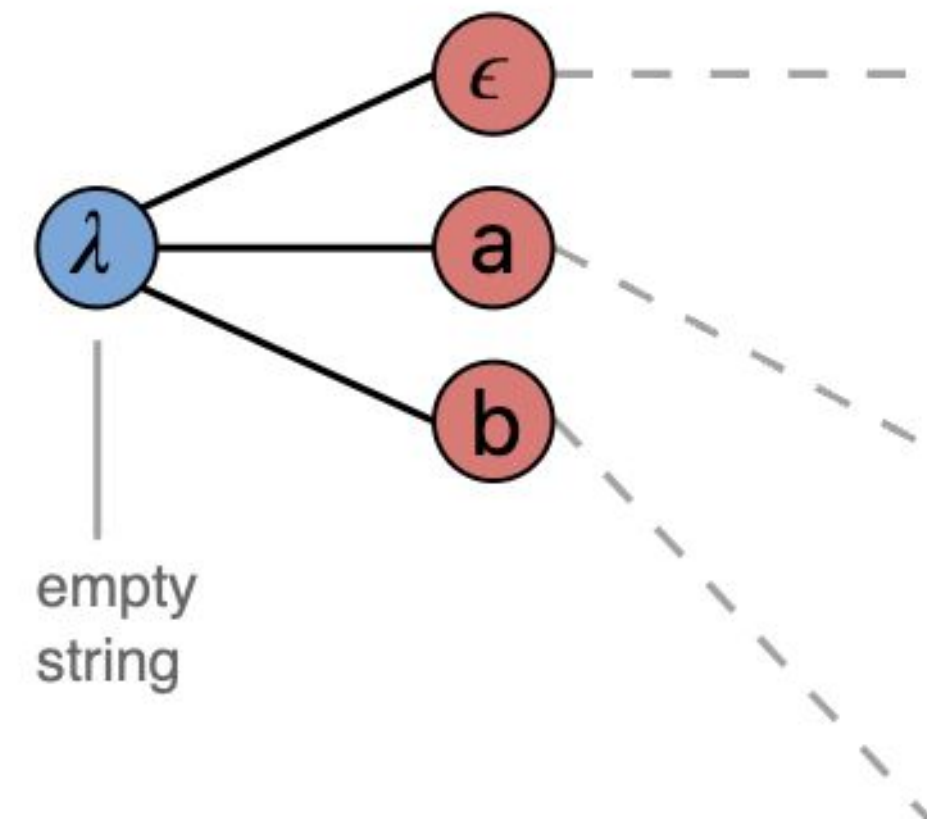
hidden =  $h_1$

log-score-2,  $h_2 = \text{log\_softmax}(\text{RNN}(\text{blank}, h_1))$

$T = 1$

current hypotheses

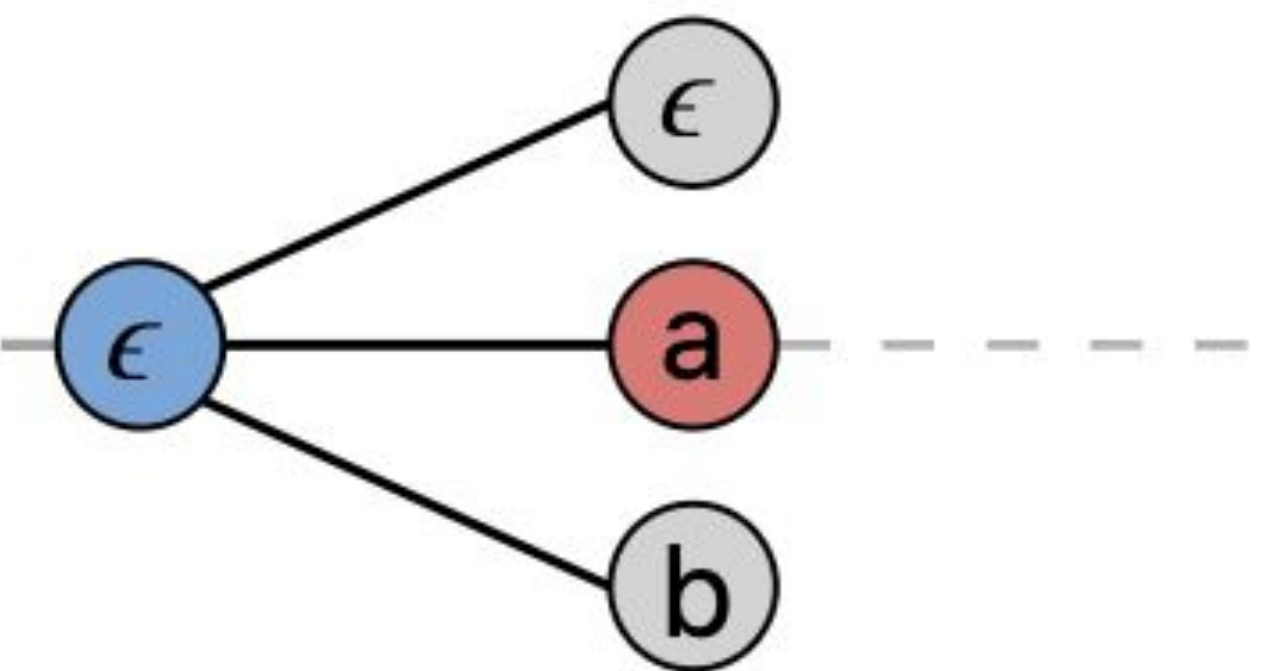
proposed extensions



$T = 2$

current hypotheses

proposed extensions



# Инференс LM

## RNN

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

Шаг 2: [ $\epsilon$ , blank]

hidden =  $h_1$

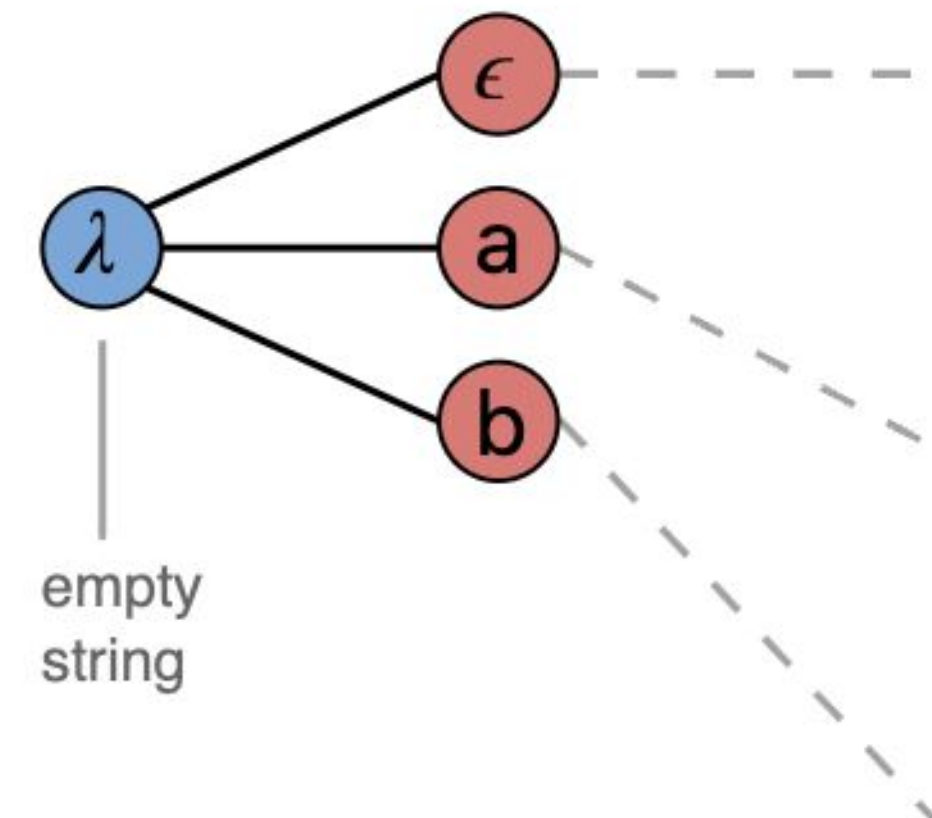
log-score-2,  $h_2 = \log\_softmax(RNN(blank, h_1))$

log-score-2 - вектор log-вероятностей размера словаря

$T = 1$

current hypotheses

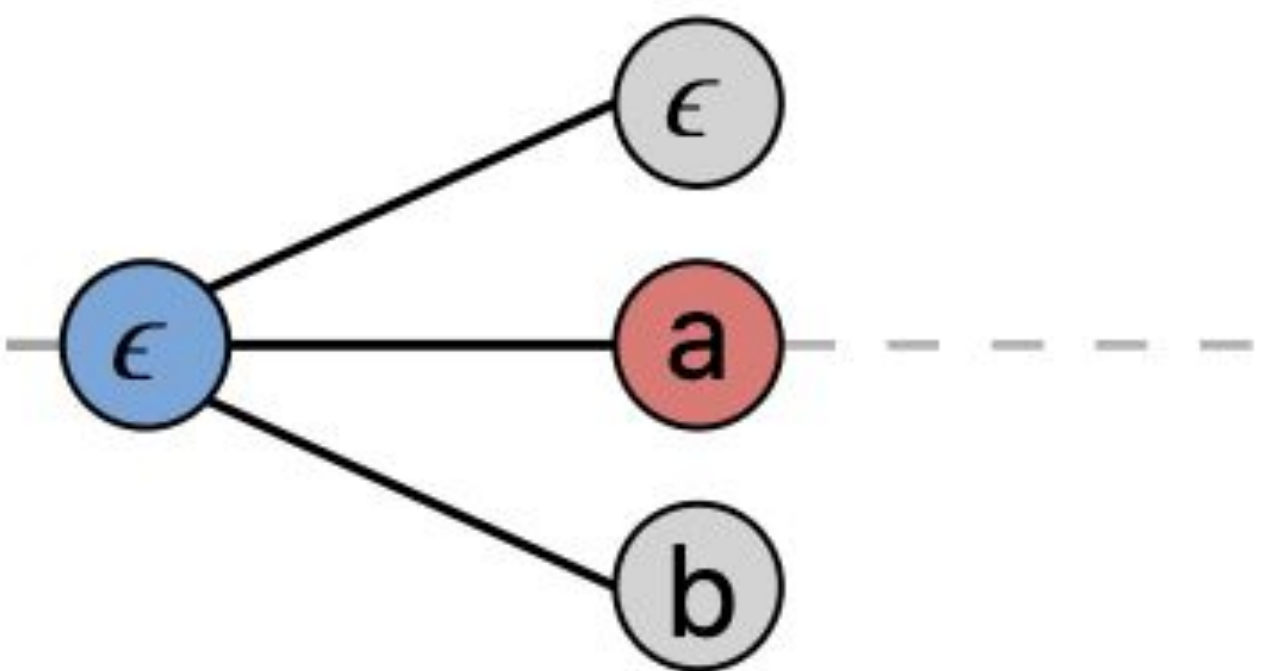
proposed extensions



$T = 2$

current hypotheses

proposed extensions



# Инференс LM

## RNN

A standard beam search algorithm with an alphabet of  $\{\epsilon, a, b\}$  and a beam size of three.

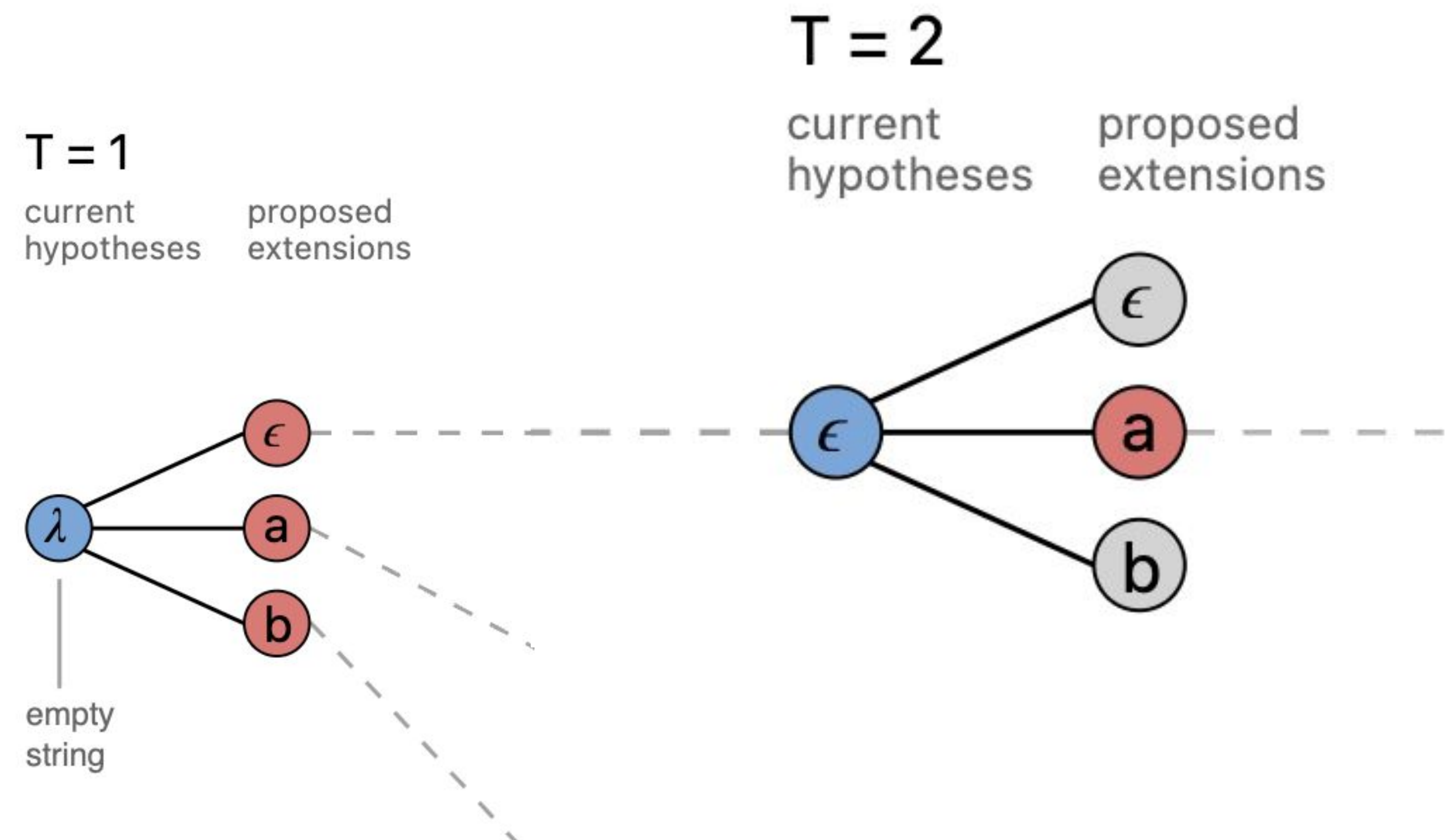
Шаг 2: [ $\epsilon$ , blank]

hidden =  $h_1$

log-score-2,  $h_2 = \text{log\_softmax}(\text{RNN}(\text{blank}, h_1))$

log-score-2 - вектор log-вероятностей размера словаря

для каждого луча хранится свой hidden



# Keywords / Hotwords

hotword - DLS



# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1:

score 1:

гипотеза 2:

score 2:

гипотеза 3:

score 3:

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: a

score 1: s1\_1

гипотеза 2: b

score 2: s2\_1

гипотеза 3: c

score 3: s3\_1



# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: ab

score 1: s1\_2

гипотеза 2: bc

score 2: s2\_2

гипотеза 3: ca

score 3: s3\_2

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: ab_	score 1: s1_3
-----------------	---------------

гипотеза 2: bc_	score 2: s2_3
-----------------	---------------

гипотеза 3: ca_	score 3: s3_3
-----------------	---------------

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: ab\_d

score 1:  $s1\_4 + \text{boost\_score} * 1$

гипотеза 2: bc\_d

score 2:  $s2\_4 + \text{boost\_score} * 1$

гипотеза 3: ca\_d

score 3:  $s3\_4 + \text{boost\_score} * 1$

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: ab\_dl

score 1:  $s1\_5 + \text{boost\_score} * 2$

гипотеза 2: bc\_dl

score 2:  $s2\_5 + \text{boost\_score} * 2$

гипотеза 3: ca\_dl

score 3:  $s3\_5 + \text{boost\_score} * 2$



# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: ab\_dls

score 1:  $s1\_6 + \text{boost\_score} * 3$

гипотеза 2: bc\_dls

score 2:  $s2\_6 + \text{boost\_score} * 3$

гипотеза 3: ca\_dla

score 3:  $s3\_6 + 0$

# Keywords / Hotwords

**hotword - DLS**

словарь = {a, b, c, d, l, s}

N-beams = 3

гипотеза 1: ab\_dls\_

score 1:  $s1\_7 + \text{boost\_score} * 3$

гипотеза 2: bc\_dlss

score 2:  $s2\_7 + \text{boost\_score} * 3$

гипотеза 3: ca\_dlas

score 3:  $s3\_7$

# Финальная схема

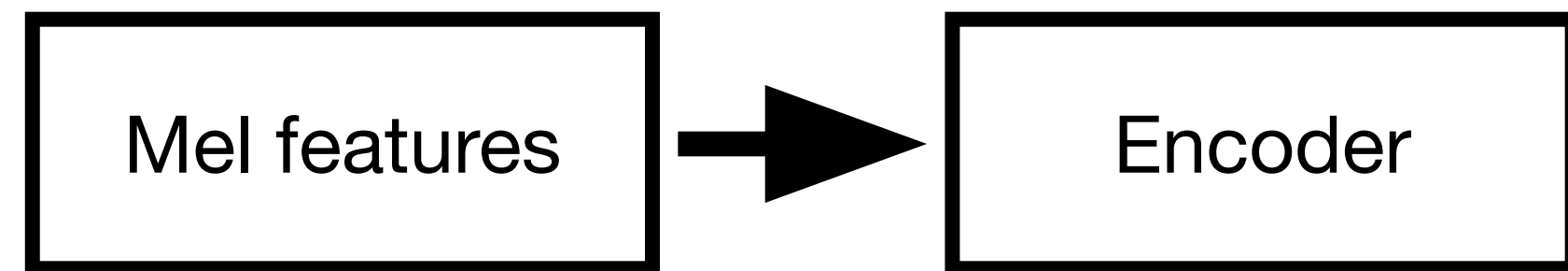
# Полный пайплайн



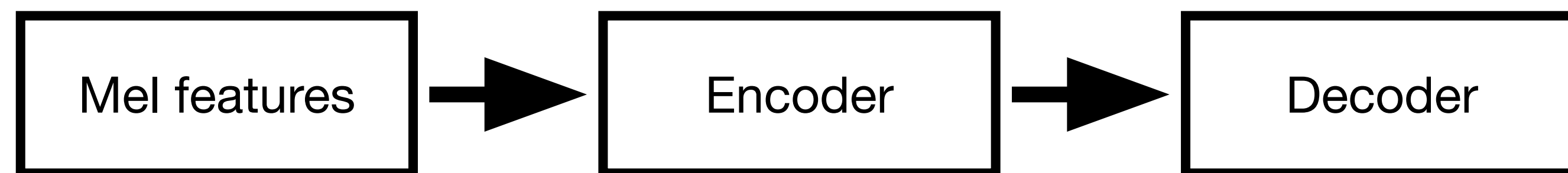
# Полный пайплайн

Mel features

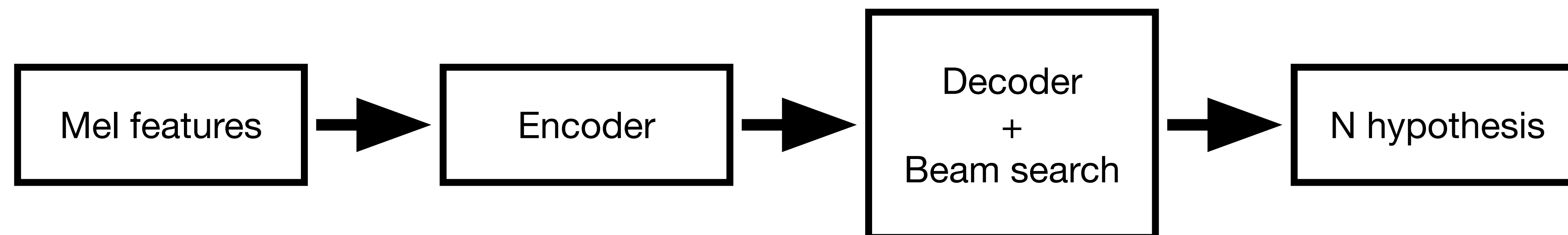
# Полный пайплайн



# Полный пайплайн

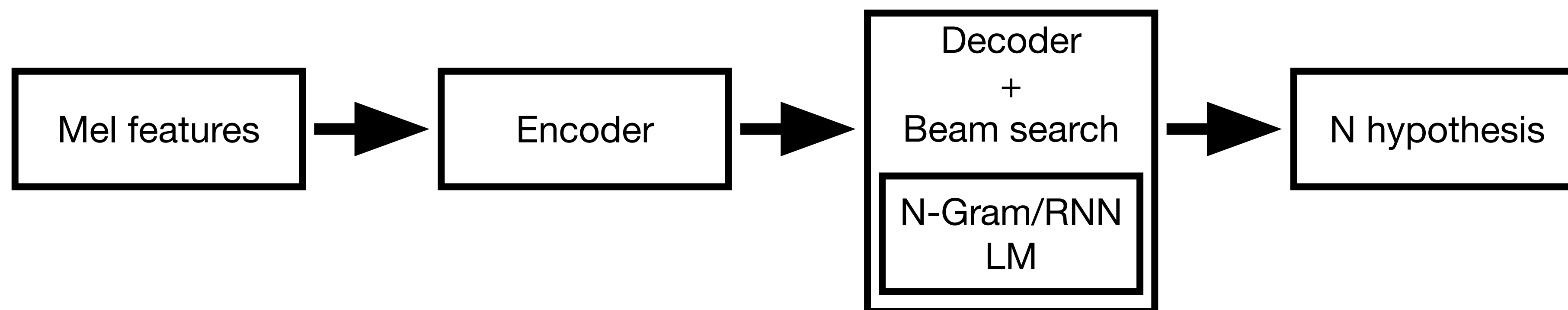


# Полный пайплайн

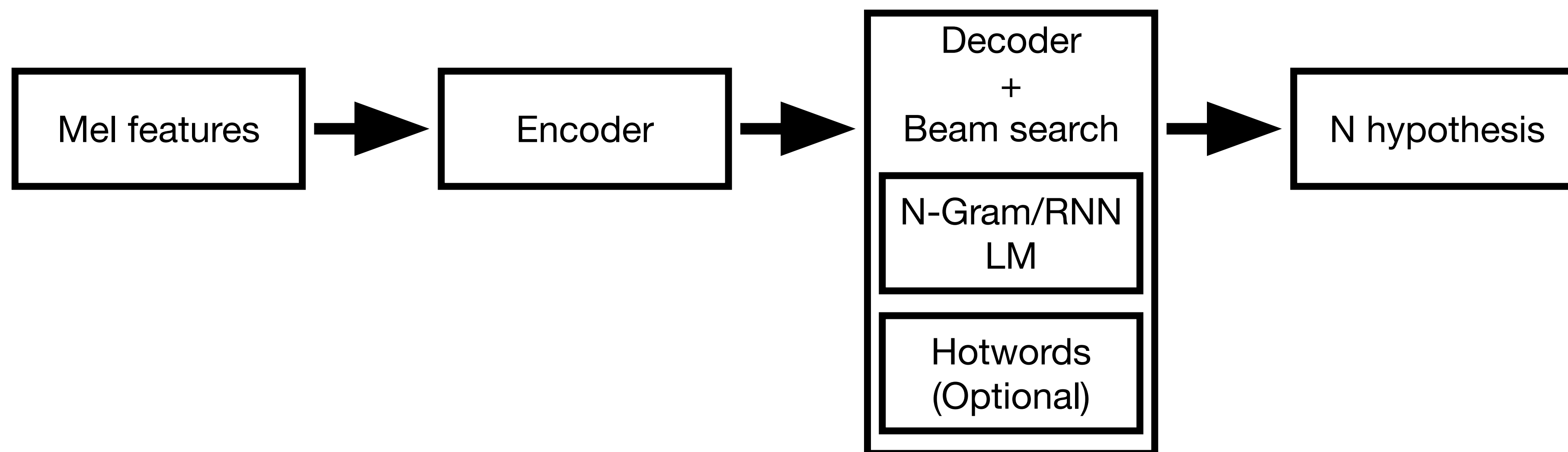




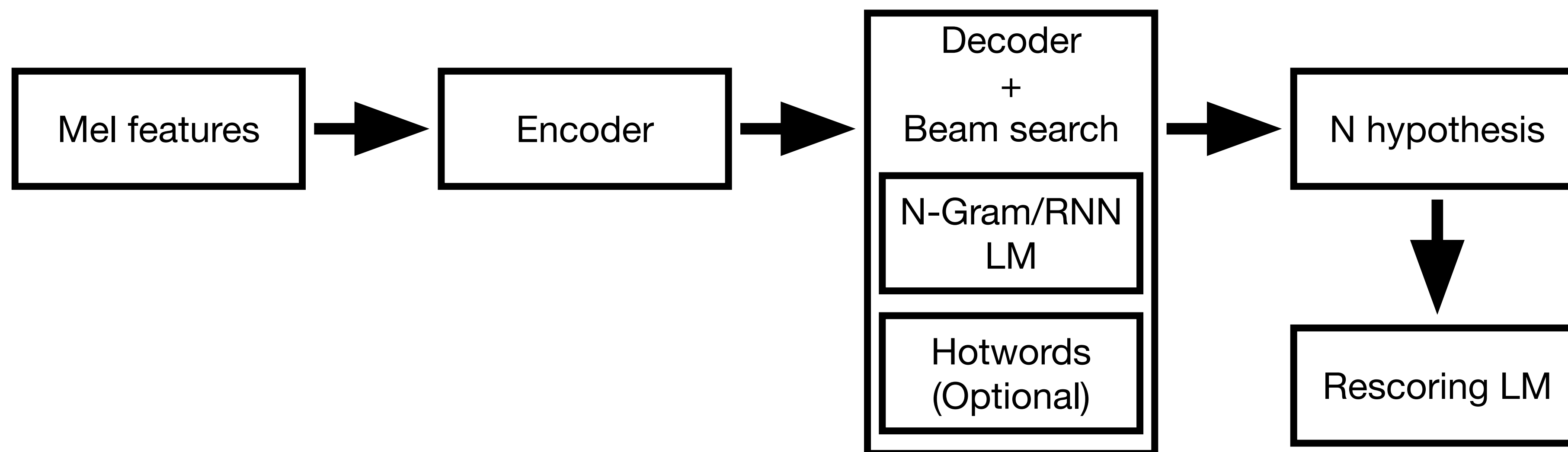
# Полный пайплайн



# Полный пайплайн



# Полный пайплайн



# Полный пайплайн

