

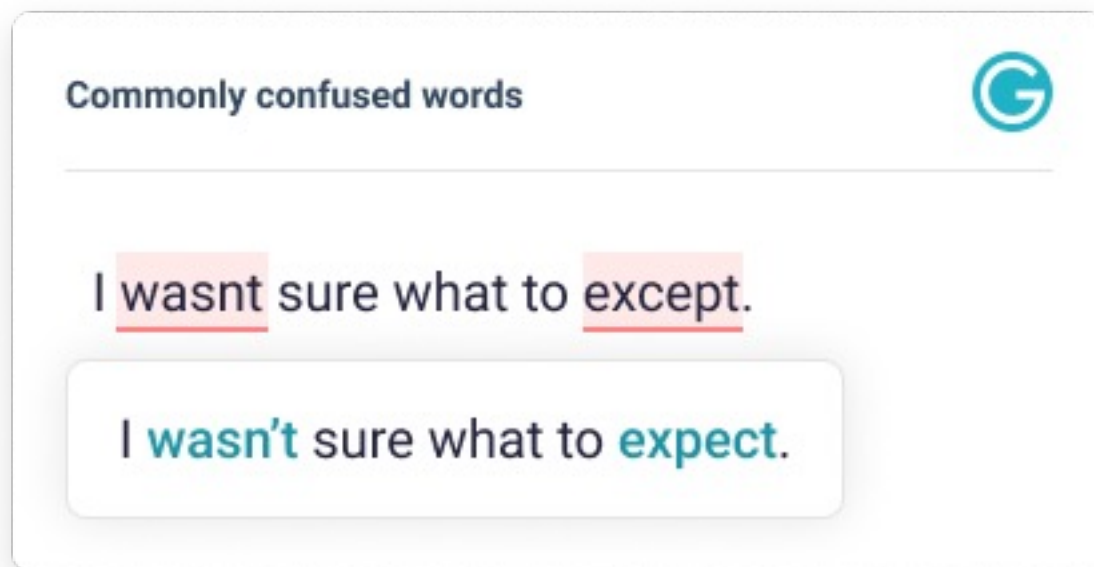
# **Задача исправления опечаток**

# План занятия

1. Постановка задачи
2. Модель Норвига
3. Модель на N-граммах
4. Трансформеры. Трюк с загрязнением данных
5. Библиотеки для решения задачи

## Постановка задачи

Задача исправления опечаток в машинном обучении заключается в разработке моделей, способных автоматически обнаруживать и исправлять ошибки в тексте, вызванные опечатками. Опечатки могут включать в себя ошибки при наборе, неверное распознавание слов, перестановку букв, альтернативное написание слов и другие виды ошибок, связанных с неверным вводом текста.



## Расстояние Левенштейна

Для поиска кандидатов для исправления слова часто используется расстояние Левенштейна.

Расстояние Левенштейна (или редакционное расстояние) — это метрика сходства между двумя строковыми последовательностями.

- Чем больше расстояние, тем более различны строки.
- Для двух одинаковых последовательностей расстояние равно нулю.

По сути это **минимальное число односимвольных преобразований (удаления, вставки или замены), необходимых, чтобы превратить одну последовательность в другую.**

Б	И	Б	А
Б	О	Б	А

$$\text{LEV}(\text{'БИБА'}, \text{'БОБА'}) = 1$$

## Расстояние Левенштейна

А	В	С	Т	Р			И	Я
А	В	С	Т	Р	А	Л	И	Я

$$\text{LEV}(\text{'АВСТРИЯ'}, \text{'АВСТРАЛИЯ'}) = 2$$

	К	О	Т	И	К	
С	К	О	Т	И	Н	А

$$\text{LEV}(\text{'КОТИК'}, \text{'СКОТИНА'}) = 3$$

## Расстояние Левенштейна

Чему равно расстояние Левенштейна между словами HONDA и HYUNDAI?

## Расстояние Левенштейна

Чему равно расстояние Левенштейна между словами HONDA и HYUNDAI?

H	O		N	D	A	
H	Y	U	N	D	A	I

H	Y	U	N	D	A	I
H		O	N	D	A	

# Расстояние Левенштейна

Для вычисления расстояния Левенштейна используется алгоритм Вагнера-Фишера, хорошо рассказанный [здесь](#).

Но как правило мы будем пользоваться готовыми его реализациями – это следующие питоновские библиотеки:

- **strsimpy**
- **python-Levenshtein**
- **NLTK**

[https://colab.research.google.com/drive/1AzpUnreWsUxL\\_w64yEChDlvnQ9HrB2H6?usp=sharing](https://colab.research.google.com/drive/1AzpUnreWsUxL_w64yEChDlvnQ9HrB2H6?usp=sharing)



## Модель Норвига “How to Write a Spelling Corrector” (2007)

В статье представлен простой метод исправления опечаток на основе статистики слов. В методе используется частота встречаемости слов и расстояние Левенштейна для предложения наилучшего варианта исправления опечаток.

## Модель Норвига “How to Write a Spelling Corrector” (2007)

- 1. Корпус текста:** Используется некоторый корпус текста, который предоставляет информацию о частоте встречаемости слов в естественном языке.
- 2. Генерация кандидатов:** Для данного слова с опечаткой генерируются кандидаты исправлений (на основе расстояния Левенштейна). Кандидаты получаются из слов-оригиналов операциями вставки, удаления, замены и транспозиции букв.
- 3. Оценка вероятности:** Для каждого кандидата вычисляется вероятность (посчитанная по большому корпусу текстов), основанная на частоте встречаемости слов в корпусе. Более вероятные слова имеют более высокую вероятность быть правильными исправлениями.
- 4. Выбор наилучшего кандидата:** Выбирается кандидат с наивысшей вероятностью в качестве наилучшего исправления.
- 5. Возврат результата:** Возвращается исправленное слово.

## Модель на N-граммах

Эта модель – развитие модели Норвига. В ней вероятности считаются с учетом контекста.

Идея: если мы встречаем слово не из словаря, то среди слов из словаря, находящихся от него на расстоянии 1 или 2, выбираем слово с наибольшей вероятностью находиться на этой позиции, зная предыдущие несколько слов.

## Модель на N-граммах

### **1.Создание модели N-грамм:**

1. Собираем обучающий корпус текста, на котором можно обучить модель N-грамм.
2. Делим текст на N-граммы, где N - это число элементов в последовательности.  
Например, для биграмм  $N=2$ , для триграмм  $N=3$ .
3. Подсчитываем частоту встречаемости каждой N-граммы в корпусе.

### **2.Генерация кандидатов:**

1. Для данного слова с опечаткой генерируем кандидатов на основе операций вставки, удаления, замены и транспозиции букв, так же, как и в методе Норвига.

### **3.Оценка вероятности с использованием N-грамм:**

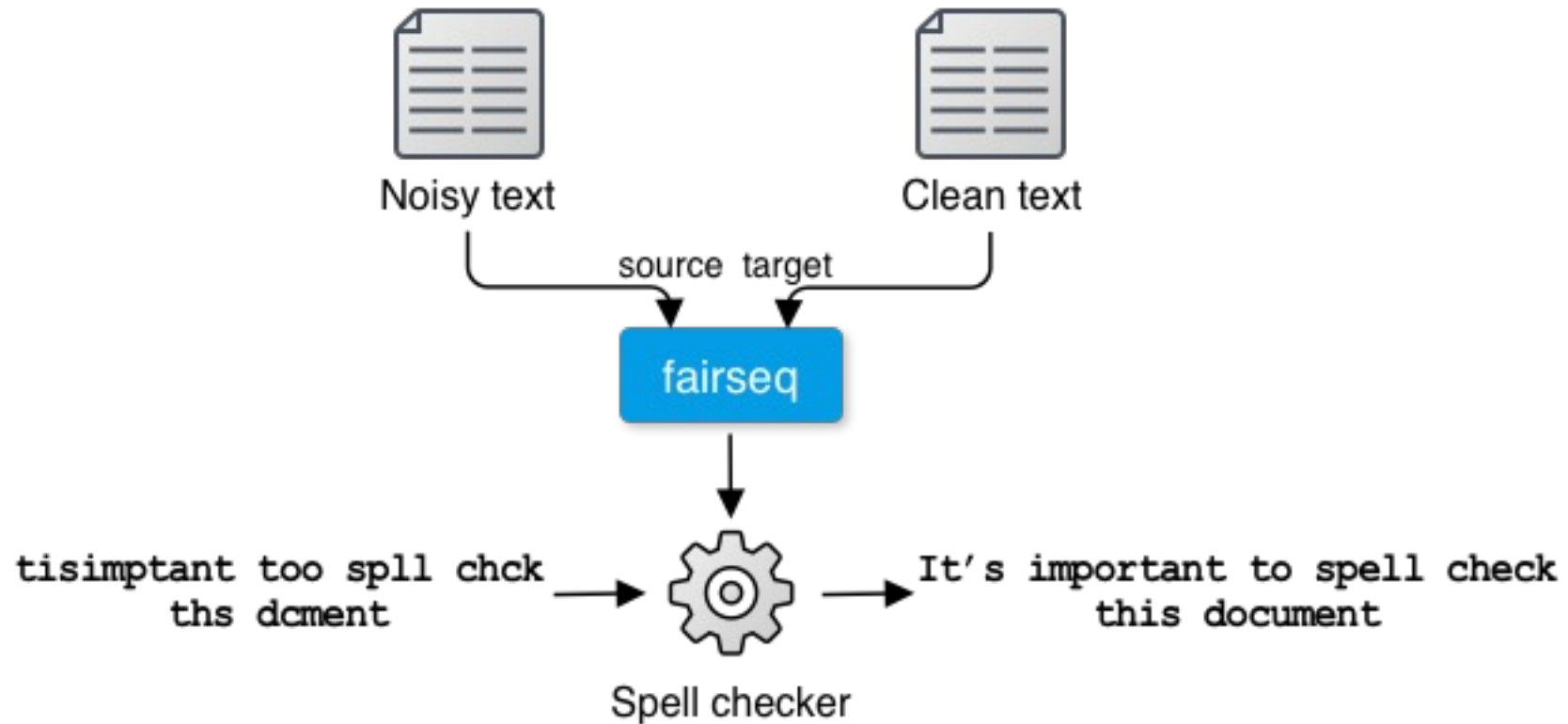
1. Для каждого кандидата оцениваем вероятность его появления в тексте с использованием модели N-грамм: это можно сделать при помощи оценки вероятности всей фразы, содержащей слово-кандидата, основанной на частоте встречаемости соответствующих N-грамм (по формуле условной вероятности).

### **4.Выбор наилучшего кандидата:**

1. Выбираем кандидата с наибольшей вероятностью в качестве наилучшего исправления.

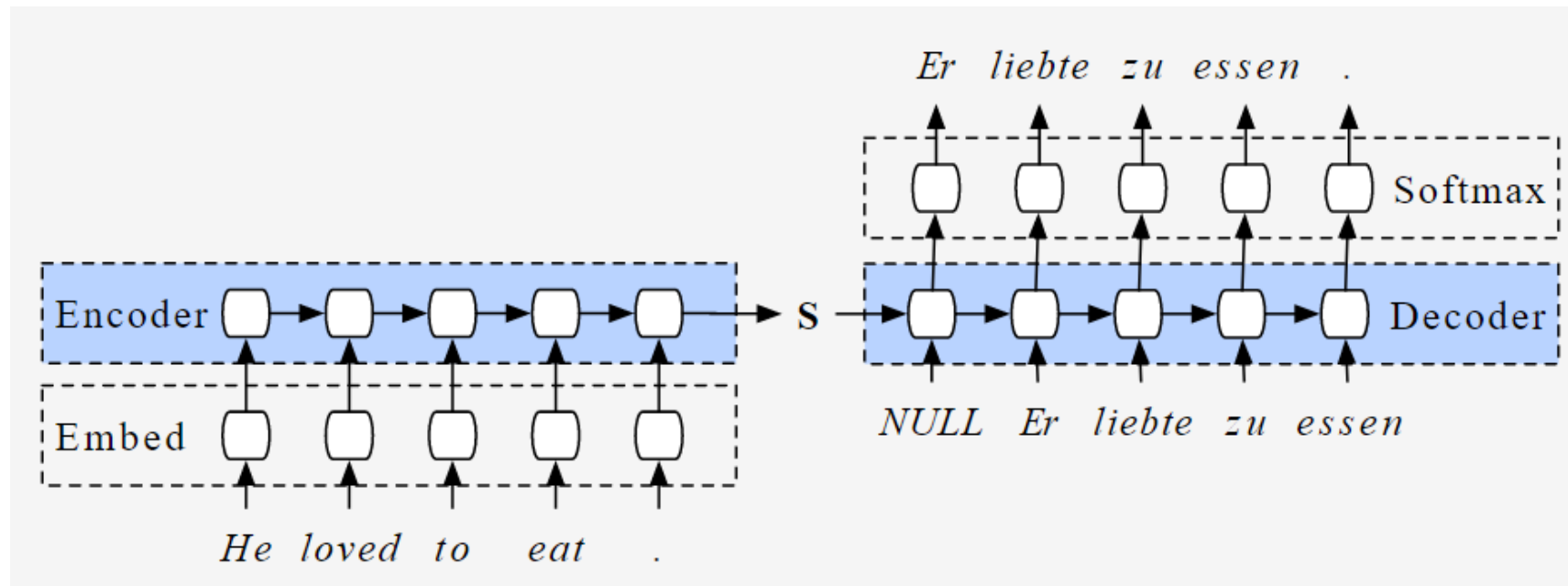
## Трансформеры для задачи исправления опечаток

Общая схема обучения при помощи фреймворка fairseq:



# Spell Correction as Machine Translation

Архитектура модели:



Мы решаем задачу как задачу **sequence-to-sequence** путем генерации слов без опечаток.

## Трансформеры для задачи исправления опечаток

Для обучения модели, которую мы будем обсуждать, использовался большой размеченный корпус:

<https://github.com/mhagiwara/github-typo-corpus>

# Трансформеры для задачи исправления опечаток

Модель обучается на буквах (то есть один токен — это один символ):

```

668 If_you_are_using_a_low-end_CPU_or_your_GP
669-#segment#nuked#bytes##Size_in_bytes_of_s
670-tenant#_#public#
671-ensure_#_#present,
672-echo_#_#h3#_H_Y_M_L_Intermediates_#_#h3#_#_#_i
673-updated_#_#true
674-#_#_For_#P#PMC_members#_Committer_#_#P#PMC_
675--_buildslave29
676 for_s_in_shapes#
677-#_inform_listeners_about_an_update_#stat
678-#_#_#_#_#MySQL#_#_#_#_#
679-#_#param_out_the_underlaying_writer
680-#_Locate_a_non-protected_child_node_def_d
681-The_model_API_in_mxnet_as_not_really_an_A
682-dateStr#_'Feb_22th',_2019',
683-#_#_Less_boilerplate_when_bootstrapping_#l
684-#div_class##rowd##
685-For_a_tutorial_on_using_the_Pulsar_Go_cli
686-#_#_anonymous_is_only_granted_READ_premiss
687-for_all_other_states_the_definition_must_be
688-We_have_following_receivers,_and_#default
689-LocalSpan_represents_a_normal_Java_method
690-into_oak-run,_or_be_specified_separately_
691-#_The_regular_expression_pattern_used_to_
692-pulsarFunctionsCluster#pulsar-cluster-1
693-Oak_module_providing_exercises_for_develo
694-The_Jackrabbit_main_project_is_located_in
695-scheduler_#_3.0.0-SNAPSHOT_#_#_standard-3
696 #Flip_Update_bit_#Queue_Updates#_for_serv

```

```

668 If_you_are_using_a_low-end_CPU_or_your_GP
669+#segment#nuked#bytes##Size_in_bytes_of_s
670+tenant#_#public#
671+ensure_#_#present#
672+echo_#_#h3#_H_T_M_L_Intermediates_#_#h3#_#_#_i
673+updated_#_#True
674+#_#_For_#P#PMC_members#_Committer_#_#P#PMC_
675+-_buildslave
676 for_s_in_shapes#
677+#_inform_listeners_about_an_update_#statu
678+#_#_#_#_#MySQL#_#_#_#_#
679+#_#param_out_the_underlying_writer
680+#_Locate_a_non-protected_child_node_def_d
681+The_model_API_in_mxnet_is_not_really_an_A
682+dateStr#_'Feb_22nd',_2019',
683+#_#_Less_boilerplate_when_bootstrapping_#l
684+#div_class##row##
685+For_a_tutorial_on_using_the_Pulsar_Go_cli
686+#_#_anonymous_is_only_granted_READ_permiss
687+for_all_other_states_the_definition_must
688+We_have_following_receivers,_and_#default
689+LocalSpan_represents_a_normal_Java_method
690+into_oak-run,_or_be_specified_separately_
691+#_The_regular_expression_pattern_used_to_
692+pulsarFunctionsCluster#_pulsar-cluster-1
693+Oak_module_providing_exercises_for_develo
694+The_Jackrabbit_main_project_is_located_in
695+scheduler_#_3.0.0_#_#_standard-3.0.0_#_Pr
696 #Flip_Update_bit_#Queue_Updates#_for_serv

```

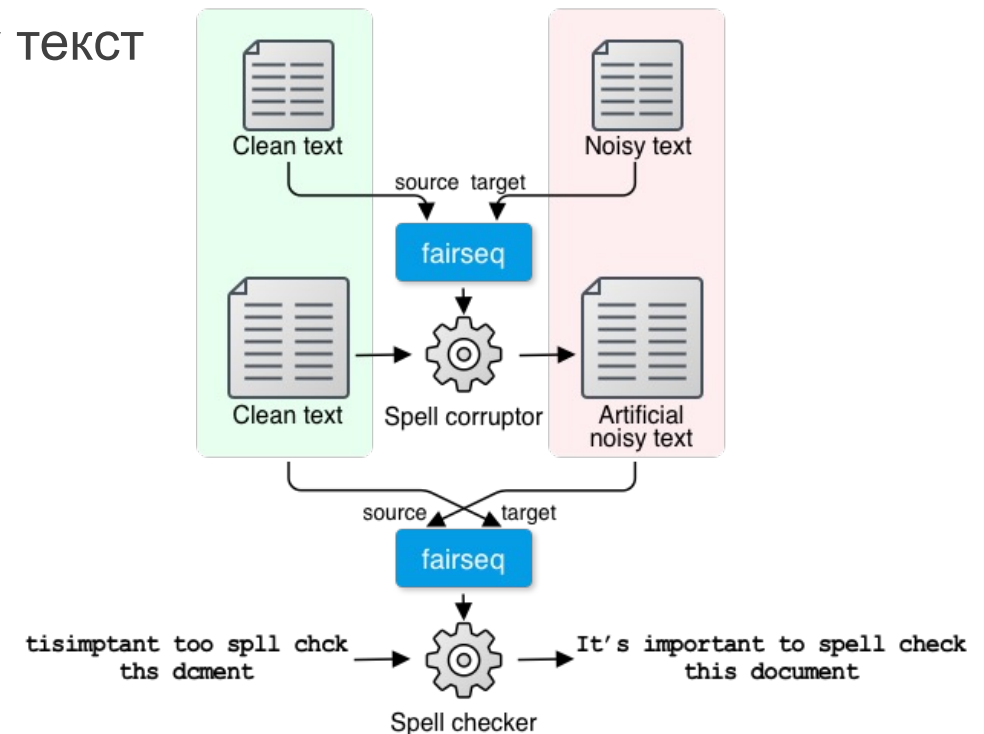


# Искусственное расширение датасета

Мы можем искусственно расширить датасет, добавив опечатки. Но сложность состоит в том, чтобы опечатки были не случайными, а такими, какие обычно делают люди.

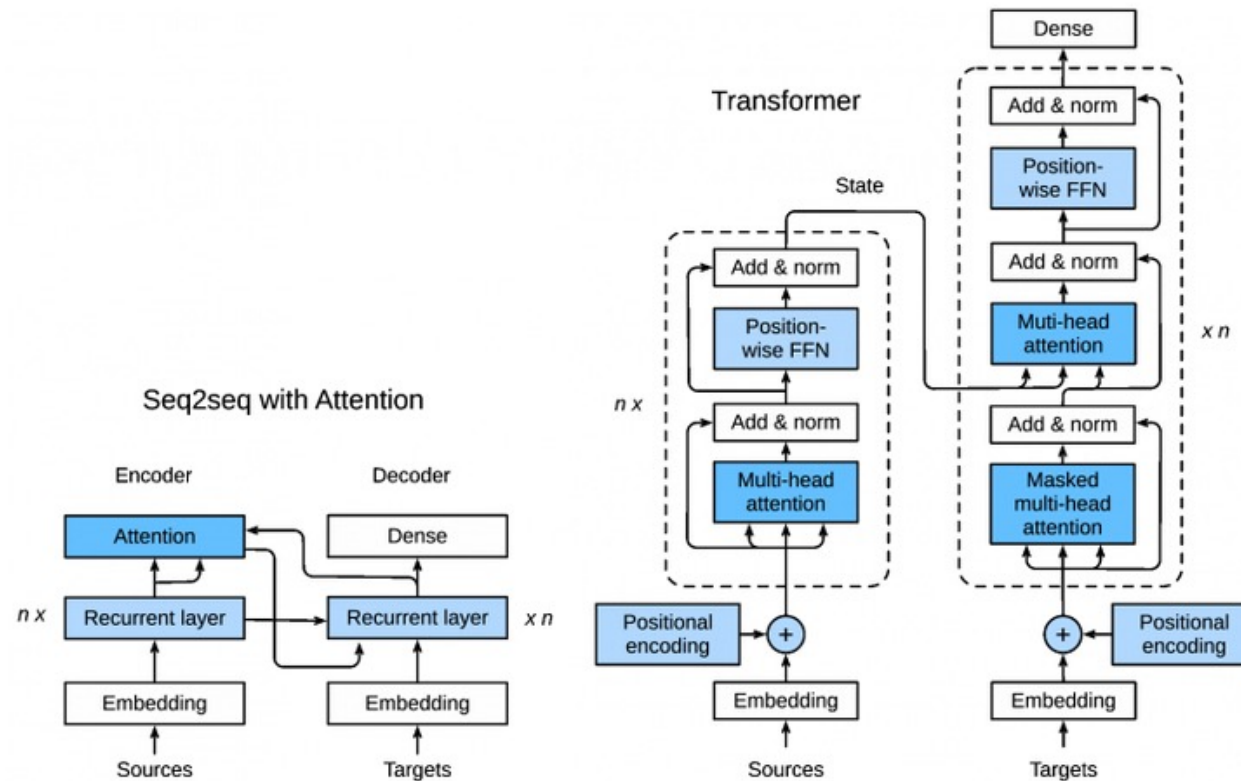
Идея:

- Обучим Spell Corruptor, поменяв местами Clean и Noisy текст из наших размеченных данных.
- Затем мы можем “испортить” любой чистый текст, и получим дополнительные данные для обучения.

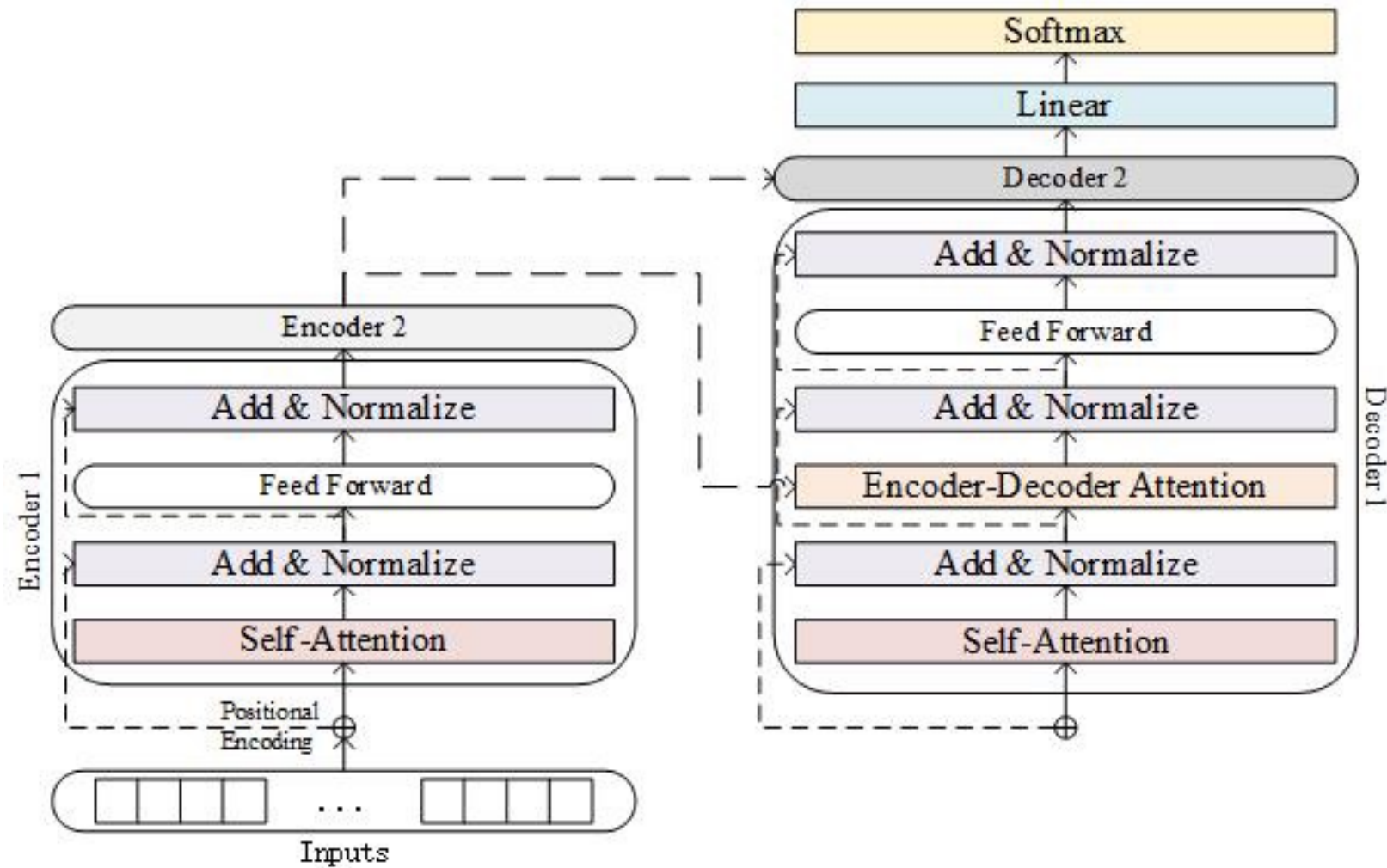


# M2M100-модель

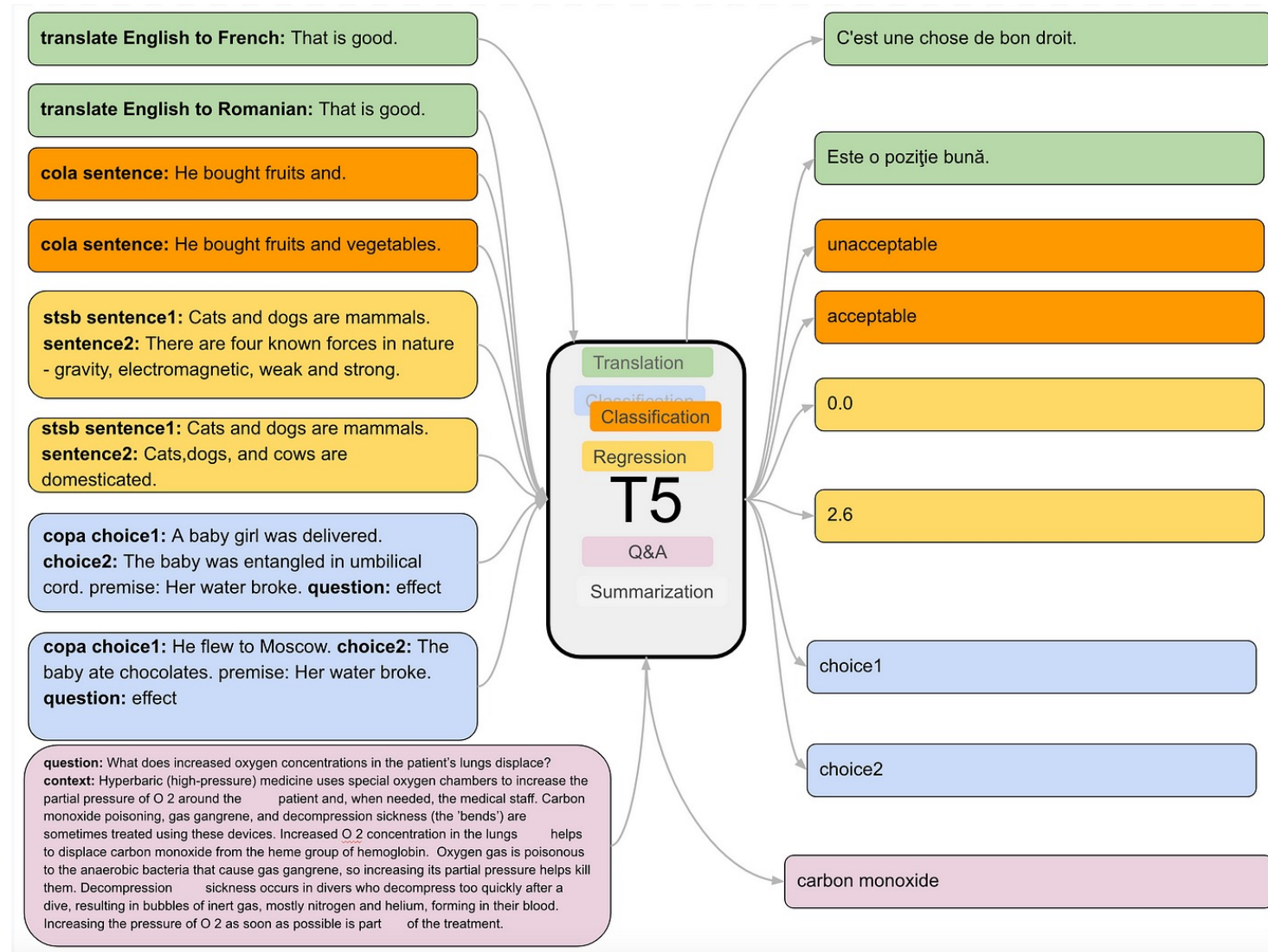
M2M100 is a multilingual encoder-decoder (seq-to-seq) model primarily intended for translation tasks. As the model is multilingual it expects the sequences in a certain format: A special language id token is used as prefix in both the source and target text. The source text format is `[lang_code] X [eos]`, where `lang_code` is source language id for source text and target language id for target text, with `X` being the source or target text.



# T5-модель



# T5-модель



# Библиотеки для Spell Correction

Статистические подходы к исправлению опечаток:

- TextBlob
- Pyspellchecker (Norwig model)
- Spylls (N-gram model)

Языковые модели для исправления опечаток:

- JamSpell
- SpaCy
- Модели из HuggingFace

## Практика! Библиотеки для Spell Correction

- <https://colab.research.google.com/drive/1N3UK004WVR1-vhF4GP35ms06L4blhHaM?usp=sharing>
- [https://colab.research.google.com/drive/1dd-IcVXXY1OrTtB5xKPF\\_cZgP0NI1Q7I?usp=sharing](https://colab.research.google.com/drive/1dd-IcVXXY1OrTtB5xKPF_cZgP0NI1Q7I?usp=sharing)