

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ
РОССИЙСКОЙ ФЕДЕРАЦИИ
федеральное государственное автономное образовательное учреждение
высшего образования «Санкт-Петербургский политехнический
университет Петра Великого»

Институт компьютерных наук и кибербезопасности

Высшая школа технологий искусственного интеллекта

Направление: 02.03.01 «Математика и компьютерные науки»

Монады в Haskell
Курсовая работа

Студент,
группы 5130201/20101

_____ Астафьев И. Е.

Преподаватель

_____ Моторин Д. Е.

«_____» _____ 2024г.

Санкт-Петербург, 2024

Содержание

| | |
|--|-----------|
| Введение | 3 |
| 1 Теоретические сведения | 4 |
| 1.1 Синтаксический анализатор | 4 |
| 1.2 Основные этапы работы парсера | 4 |
| 1.3 Функторы | 4 |
| 1.4 Аппликативные функторы | 5 |
| 1.5 Класс типов <code>Alternative</code> | 5 |
| 1.6 Монады | 5 |
| 1.7 N-граммы | 5 |
| 2 Часть 1. Парсер строк бинарных чисел и операций | 7 |
| 2.1 Реализация программы | 7 |
| 3 Часть 2. Генератор предложений с помощью N-грамм | 11 |
| 3.1 Реализация программы | 11 |
| 4 Результаты | 18 |
| 4.1 Часть 1. Парсер строк бинарных чисел и операций | 18 |
| 4.2 Часть 2. Генератор предложений с помощью N-грамм | 19 |
| Заключение | 21 |
| Приложение А. Парсер бинарных операций | 22 |
| Приложение Б. Генератор предложений | 26 |

Введение

В данном отчете описаны результаты выполнения курсовой работы: монады в Haskell.

Часть 1

Постановка задачи:

1. Написать синтаксический анализатор (парсер), разбирающий строки, прочитанного файла .txt. Файл должен содержать значения и бинарные операции.
 - Значения: строки битов
 - Бинарные операции: AND(&), OR(|), XOR(\oplus)
 - Пример строки в файле 01010 & 0011. Вычислить проанализированное выражение результат вычисления на экран. Пользователь вводит название файла.

Часть 2

Постановка задачи:

1. Прочитать текст из файла, указываемого пользователем. Синтаксически проанализировать текст согласно правилам: слова состоят только из букв; предложения состоят только из слов и разделены символами: !?;,. Разбить текст на предложения. Удалить все символы пунктуации и цифры из слов и предложений.
2. Составить модель N-грамм. Использовать модель биграмм и триграмм. По списку предложений составить словарь. Ключами являются: одно слово, либо пара слов. Значениями в словаре является список всех уникальных возможных продолжений триграммы (т.е. список пар слов или одиночных слов). Словарь сохранить в файл .txt.

Пример текста и словаря на его основе:

Текст: [a b, c d e! b c d? e b c # a d. a f; f.]

Словарь: ['a' : ['b', 'd', 'f', 'b c']; 'b' : ['c', 'c d', 'c a']; 'c' : ['d', 'a', 'd e', 'a d']; 'd' : ['e']; 'e' : ['b', 'b c']; f : []; 'a b' : ['c']; 'b c' : ['d', 'a']; 'c d' : ['e']; 'e b' : ['c']; 'c a' : ['d']; 'a d' : []; 'a f' : []; 'd e' : []].

3. Реализовать взаимодействие с пользователем. Пользователь вводит одно слово или пару слов. Программа возвращает стоку случайной длины в диапазоне от 2 до 15 слов, если задаваемого пользователем слова нет в ключах словаря, выдавать соответствующее сообщение. Фраза составляется путем добавления случайного слова (или пары) из списка значений текущего слова-ключа (или пары-ключа), до тех пор, пока либо не будет сформировано предложение нужной длины, либо не будет достигнут ключ, у которого нет значений.
4. Организовать диалог двух моделей N-грамм созданных на двух разных текстах. Тексты для второй модели выбрать самостоятельно. Пользователь задает начальное слово (или пару) и глубину M сообщений, которыми обмениваются модели. Ответ модели основывается на последнем слове из предложения оппонента (если последнее слово отсутствует в словаре, то предпоследнее и т.д. пока не будет найдено подходящее слово или не закончится предложение оппонента).

Автор текстов: Максим Горький.

1 Теоретические сведения

1.1 Синтаксический анализатор

Синтаксический анализатор, или парсер, является ключевым компонентом любой системы обработки текста, языков программирования или данных. Его задача заключается в разборе входной последовательности символов (или токенов, если предварительно был выполнен лексический анализ) и проверке её соответствия заданной грамматике, а также в построении структуры данных, описывающей синтаксис входной строки.

1.2 Основные этапы работы парсера

Работа синтаксического анализатора включает следующие этапы:

1. **Получение входных данных.** Парсер принимает последовательность символов или токенов на вход. Эти данные поступают либо напрямую из исходного текста, либо из лексического анализатора.
2. **Проверка соответствия грамматике.** Парсер проверяет, соответствует ли входная строка правилам грамматики, заданным в виде контекстно-свободной грамматики (КС-грамматики) или другого формализма.
3. **Построение выходной структуры.** Если строка соответствует грамматике, парсер формирует дерево разбора (дерево синтаксического анализа), абстрактное синтаксическое дерево (AST) или другую структуру данных, которая будет использоваться для последующей обработки.

Типы синтаксических анализаторов Существуют различные типы синтаксических анализаторов, которые различаются по методам работы и применяемым алгоритмам:

- **Восходящие парсеры (Bottom-up).** Эти парсеры строят дерево разбора, начиная с листьев (токенов) и продвигаясь вверх к корню. Примеры: алгоритмы LR, LALR и SLR.
- **Нисходящие парсеры (Top-down).** Такие парсеры начинают с корня дерева и рекурсивно обрабатывают правила грамматики, продвигаясь к листьям. Примером является рекурсивно-спускающийся парсер.
- **LL-парсеры.** Это подмножество нисходящих парсеров, которые обрабатывают грамматику слева направо (Left-to-right) и строят левую производную (Leftmost derivation).
- **Рекурсивно-спускающиеся парсеры.** Это парсеры, которые используют рекурсивные вызовы функций для обработки правил грамматики. Они просты в реализации и часто используются для разработки парсеров на языках программирования.

1.3 Функторы

Функторы — это абстракция, представляющая контейнеры, к элементам которых можно применять функции. Они принадлежат классу типов `Functor`, который определяет операцию `fmap`. Эта операция позволяет трансформировать содержимое контейнера, не изменяя его структуру. Например, с помощью `fmap` можно применить функцию ко всем элементам списка.

1.4 Аппликативные функторы

Аппликативные функторы расширяют возможности функторов, добавляя способ применения функций, заключённых в контексте, к значениям в том же контексте. Это реализуется через операции `<*>` и `pure` из класса типов `Applicative`. Операция `pure` помещает значение в контекст, а `<*>` применяет функцию из контекста к значениям, также находящимся в контексте. Аппликативные функторы полезны для работы с вычислениями, которые могут завершиться неудачей.

1.5 Класс типов `Alternative`

Класс типов `Alternative` используется для работы с контейнерами, которые могут содержать либо значения, либо быть пустыми. Основные операции:

- `empty` — представляет пустой контейнер.
- `<|>` — предоставляет способ комбинирования двух контейнеров. Например, выбирает первый непустой контейнер.

`Alternative` часто используется в парсерах, где можно выбирать между различными вариантами анализа текста.

1.6 Монады

Монады — это обобщение аппликативных функторов, позволяющее связывать последовательные вычисления. Они определяются через операции `>=` (`bind`) и `return`. Операция `return` помещает значение в контекст, а `>=` позволяет передавать результат одного вычисления следующему.

Монады обеспечивают удобство работы с вычислениями, которые могут быть последовательными, контекстуальными или содержать побочные эффекты. Они находят широкое применение в обработке ошибок, работе с вводом-выводом, парсерах и других задачах.

1.7 N-граммы

N-граммы — это последовательности из N элементов, получаемые из данных. В контексте обработки естественного языка (NLP) элементы обычно представляют собой слова или символы. Основной задачей N-грамм является моделирование и анализ последовательностей в тексте, что может быть полезно для различных задач, таких как предсказание следующего слова, анализ вероятности перехода от одного состояния к другому или классификация текста.

Определение: N-грамма — это последовательность из N элементов, где каждый элемент может быть как словом, так и символом, в зависимости от задачи. Если $N = 1$, то это называется *унитарной граммой* или *униграммой*, если $N = 2$ — *биграммой*, $N = 3$ — *триграммой* и так далее. Например:

- Униграмма для текста "Я люблю читать": ["Я", "люблю", "читать"].
- Биграмма для того же текста: ["Я люблю", "люблю читать"].
- Триграмма: ["Я люблю читать"].

Применение N-грамм в обработке естественного языка:

- *Языковое моделирование:* N-граммные модели широко используются для оценки вероятности появления слова в определенном контексте. Например, биграммная модель будет оценивать вероятность появления слова на основе предыдущего слова.
- *Предсказание текста:* Используя вероятности N-грамм, можно предсказать следующее слово в предложении. Если модель обучена на больших корпусах текста, она может точно угадывать следующее слово, основываясь на предыдущих N-1 словах.
- *Классификация текста:* N-граммные признаки могут использоваться для представления текста в задачах классификации, таких как определение тональности текста или его тематической принадлежности.

В некоторых задачах моделирования языка или предсказания текста, для представления N-грамм удобно использовать словарь, где ключами являются N-граммы (представленные в виде строк), а значениями — последовательности элементов, которые могут следовать за данным контекстом.

Пример структуры такого словаря для биграмм может выглядеть следующим образом:

```
"a": ["b "d "b c"], "b c": ["f"], "f": []
```

Здесь:

- Ключ "a" ассоциируется с массивом значений ["b "d "b c"], что означает, что после слова "a" могут следовать слова "b "d" или "b c".
- Ключ "b c" ассоциируется с массивом ["f"], что означает, что после биграммы "b c" может следовать слово "f".
- Ключ "f" имеет пустой массив [], что означает, что после слова "f" не ожидаются другие слова или элементы.

Этот словарь может использоваться в различных задачах обработки текста, таких как генерация текста или анализ вероятности появления слов. Основные преимущества этого подхода заключаются в том, что такой словарь позволяет:

- **Хранить контексты** для каждой N-граммы, то есть, для каждого слова или последовательности слов хранить информацию о том, какие элементы могут следовать за ними.
- **Быстро искать возможные продолжения** для каждой N-граммы, что полезно при решении задач предсказания следующего слова или анализа текста.
- **Использовать для предсказания вероятности**, вычисляя вероятность появления следующего слова как частоту появления той или иной N-граммы в обучающем корпусе текста.

2 Часть 1. Парсер строк бинарных чисел и операций

2.1 Реализация программы

Полный исходный код представлен в [Приложение А. Парсер бинарных операций](#).

Описание Parser

В данном разделе представлен синтаксический анализатор (парсер), реализованный на языке Haskell. Его цель — поэтапный разбор текста с возможностью обработки ошибок и комбинирования простых парсеров в более сложные конструкции.

Парсер реализован как абстракция, представляющая собой функцию, которая принимает текст на вход и возвращает результат анализа в виде оставшегося текста и результата парсинга. В случае неудачного анализа результатом работы является специальное значение `Nothing`.

Основные возможности парсера обеспечиваются реализацией нескольких стандартных классов Haskell:

- Класс `Functor` позволяет применять функции к результату парсинга, не изменяя процесс самого анализа. Это достигается за счёт обёртки результата существующего парсера в новый парсер, который трансформирует результат с помощью переданной функции.
- Класс `Applicative` позволяет комбинировать несколько парсеров, применяя функции с несколькими аргументами к их результатам. Этот механизм используется для построения последовательных цепочек парсеров. Например, можно объединить парсеры для разбора чисел и операторов в одном выражении.
- Класс `Alternative` предоставляет средства для обработки альтернативных вариантов и восстановления после ошибок. С его помощью можно определять парсеры, которые пробуют несколько вариантов анализа и выбирают первый успешный результат.

Основной особенностью реализации является использование концепций функционального программирования, таких как иммутабельность данных и композиция функций, что делает код простым для комбинирования и тестирования. Такой подход позволяет строить мощные анализаторы для обработки текста с минимальными усилиями.

Разработанный парсер служит основой для построения высокоуровневых разборщиков, таких как парсер бинарных выражений, представленный в рамках данной работы.

Описание функций проверки и парсинга символов

Для реализации базовых операций парсинга в коде используются следующие функции:

- `isBinaryChar`

Данная функция проверяет, является ли символ допустимым двоичным символом (`'0'` или `'1'`). Она используется в парсерах, которые разбирают двоичные числа.

- **isOperation**

Эта функция проверяет, относится ли символ к допустимым логическим операциям: '&' (AND), '|' (OR) или '⊕' (XOR). Она помогает распознать оператор в выражении.

- **isSpace**

Проверяет, является ли символ пробелом (' '). Функция используется для обработки пробелов между элементами выражения.

- **satisfy**

Универсальная функция-парсер, которая принимает предикат (функцию `Char -> Bool`) и создаёт парсер, возвращающий первый символ текста, удовлетворяющий предикату. Если текст пустой или первый символ не подходит под предикат, парсер возвращает `Nothing`.

Алгоритм работы:

1. Сначала извлекается первый символ строки с помощью функции `T.uncons`.
2. Если строка пуста (`Nothing`), парсер возвращает `Nothing`.
3. Если символ присутствует, проверяется выполнение предиката.
 - При успешной проверке возвращается пара: оставшийся текст и символ.
 - При неудаче возвращается `Nothing`

Простые парсеры

В данном разделе описаны базовые парсеры, которые используются для распознавания и обработки текстовых данных, таких как двоичные числа, логические операции и пробелы. Эти парсеры являются ключевыми элементами для построения более сложных синтаксических анализаторов.

- **Парсер двоичных чисел.**

Этот парсер считывает последовательность двоичных символов (0 или 1) из входного текста. Если первый символ строки является двоичным, он добавляется к результату. Парсер рекурсивно продолжает обработку оставшейся строки.

На выходе:

- Если парсинг завершился успешно, возвращается `Just (remainingText, parsedText)`, где `remainingText` — остаток строки после обработки, а `parsedText` — строка, состоящая из всех двоичных символов, считанных подряд.
- Если первый символ не является двоичным или строка пустая, возвращается `Nothing`.

- **Парсер операций.**

Этот парсер использует функцию `satisfy`, чтобы найти первый символ, который соответствует одному из предопределённых операторов (&, | или ?). Функция `isOperation` проверяет, является ли символ одним из допустимых операторов. Если символ подходит, он возвращается как результат парсера.

На выходе:

- Если первый символ является одним из допустимых операторов, парсер возвращает `Just (remainingText, op)`, где `remainingText` — остаток строки после обработки, а `op` — сам символ оператора.
- Если первый символ не является допустимым оператором, парсер возвращает `Nothing`.

- **Парсер одиночного пробела.**

Этот парсер использует функцию `satisfy`, чтобы найти первый символ, который является пробелом. Функция `isSpace` проверяет, является ли символ пробелом. Если символ подходит, он возвращается как результат парсера.

На выходе:

- Если первый символ — пробел, парсер возвращает `Just (remainingText, ' ')`, где `remainingText` — остаток строки после обработки, а `' '` — символ пробела.
- Если первый символ не является пробелом, парсер возвращает `Nothing`.

- **Парсер последовательности пробелов.**

Этот парсер использует комбинацию парсеров для обработки пробелов. Он рекурсивно использует `oneSpace` для поиска первого пробела, затем применяет себя (парсер `spaces`) для обработки оставшихся пробелов. В случае, если пробелы больше не найдены, используется `pure T.empty`, чтобы вернуть пустую строку.

На выходе:

- Если строка начинается с пробела, парсер возвращает строку из пробелов. Например, `Just (remainingText,)` в случае одного пробела или `Just (remainingText,)` в случае двух пробелов.
- Если строка не содержит пробелов или они закончились, парсер возвращает `Just (remainingText,)`, то есть пустую строку.
- Если первый символ не является пробелом, парсер возвращает `Nothing`.

Обработка логических операций и форматирование выражений

Данный раздел описывает функции, которые выполняют вычисление логических операций над двоичными числами и форматируют выражения для удобного представления результата.

- **Функция для применения операций.** Функция `applyOperation` принимает символ, обозначающий логическую операцию (`&`, `|`, `^`), и возвращает соответствующую функцию, которая применяется к двум целым числам.

- Оператор `&` интерпретируется как побитовое *AND*.
- Оператор `|` интерпретируется как побитовое *OR*.
- Оператор `^` интерпретируется как побитовое *XOR*.

Этот подход обеспечивает компактный и удобный способ выбора операции в зависимости от символа.

- **Парсер для выражений с двоичными числами.**

Этот парсер представляет собой комбинацию нескольких простых парсеров, который парсит строку, содержащую два бинарных числа и операцию между ними, игнорируя пробелы. Он использует парсер `spaces` для обработки возможных пробелов в начале, середине и конце выражения.

Сначала парсится первое бинарное число с помощью парсера `binary`, который распознаёт последовательности символов '0' и '1'. Затем, между ними парсится операция (например, `&`, `|` или `^`) с помощью парсера `operation`. После этого снова пропускаются пробелы с помощью `spaces`. В завершение, парсится второе бинарное число с использованием того же парсера `binary`.

Вся эта структура композиционно собрана с помощью оператора `<*>` из типа `Applicative`, что позволяет комбинировать несколько парсеров, обеспечивая удобную работу с каждым компонентом выражения. Полученные результаты передаются в функцию `formatExpression`, которая форматирует их в строку, включая результат вычисления операции.

- **Форматирование и вычисление выражения.** Функция `formatExpression` преобразует разобранные элементы выражения (два числа и операцию) в итоговый текст, включающий результат вычисления. Алгоритм:
 - Преобразует оба двоичных числа в целые числа.
 - Применяет указанную операцию к этим числам с помощью функции `applyOperation`.
 - Преобразует результат обратно в двоичный формат.
 - Формирует строку вида `<число1> <операция> <число2> = <результат>`.

Пример: строка `"101 & 010"` будет преобразована в `"101 & 010 = 0"`.

Эти функции обеспечивают корректную обработку двоичных выражений, их вычисление и представление результата в удобочитаемой форме, что делает парсер удобным и эффективным инструментом для анализа логических выражений.

3 Часть 2. Генератор предложений с помощью N-грамм

3.1 Реализация программы

Полный исходный код представлен в [Приложение Б. Генератор предложений](#).

Описание Parser

Описание Parser идентично описанному в разделе [Описание Parser](#) в Части 1.

Описание функций проверки и парсинга символов

Для реализации базовых операций обработки текста в коде используются следующие функции:

- **isPunctuation**

Эта функция проверяет, является ли символ знаком препинания. Она возвращает **True**, если символ входит в список `['.', '!', '?', ';', ':', '(', ')']`, и **False** в противном случае. Эта функция используется для определения наличия знаков препинания в строках, что может быть полезно при анализе текста или парсинге.

- **joinWords**

Функция принимает список строк `[Text]` и объединяет их в одну строку, вставляя между словами пробелы. Для этого используется функция `T.intercalate`, которая вставляет заданный разделитель (в данном случае пробел) между элементами списка и возвращает единую строку. Эта функция полезна для объединения слов или фраз в одно предложение или текст.

- **isWordChar**

Эта функция проверяет, является ли символ частью слова. Она возвращает **True**, если символ является буквой (проверяется с помощью функции `isLetter`) или если символ — это апостроф (`'`). Таким образом, эта функция может быть полезна для распознавания частей слов, например, для обработки таких случаев, как апострофы в английских словах (например, `"don't"` или `"it's"`).

- **satisfy**

Эта функция принимает предикат `pr` (функцию, проверяющую условие для символа) и возвращает парсер, который пытается применить этот предикат к первому символу строки. Если первый символ удовлетворяет предикату, парсер возвращает оставшуюся строку и сам символ, иначе возвращает **Nothing**. Функция полезна для парсинга отдельных символов, удовлетворяющих заданному условию, и может быть использована в различных парсерах для обработки символов, таких как цифры, пробелы или буквы.

Простые парсеры

В данном разделе описаны парсеры, которые используются для распознавания слов, пробелов, пунктуации и игнорирования ненужных символов. Эти парсеры являются важными для обработки текста и составляют основу для более сложных анализаторов.

- **Парсер слов.**

Этот парсер использует `some` и `satisfy`, чтобы считать последовательность символов, удовлетворяющих предикату `isWordChar`. Он обрабатывает символы, которые могут быть буквами или апострофами. После того как слово собрано, оно очищается от апострофов с обеих сторон с помощью `T.dropAround`. Если слово состоит только из апострофов, парсер возвращает `Nothing`.

На выходе:

- Если слово успешно разобрано, парсер возвращает `Just (remainingText, word)`, где `remainingText` — остаток строки после обработки, а `word` — разобранное слово.
- Если слово состоит только из апострофов, парсер возвращает `Nothing`.

- **Парсер одиночного пробела.**

Этот парсер использует функцию `satisfy`, чтобы найти первый символ, который является пробелом. Он использует предикат `isSpace` для проверки, является ли символ пробелом. Если символ является пробелом, он возвращается как результат парсера.

На выходе:

- Если первый символ — пробел, парсер возвращает `Just (remainingText, ' ')`, где `remainingText` — остаток строки после обработки, а `' '` — символ пробела.
- Если первый символ не является пробелом, парсер возвращает `Nothing`.

- **Парсер последовательности пробелов.**

Этот парсер использует комбинацию парсеров для обработки одного или нескольких пробелов. Он рекурсивно вызывает `oneSpace`, чтобы найти первый пробел, а затем повторно использует себя (`spaces`), чтобы обработать оставшиеся пробелы. Если пробелы закончились, используется `pure T.empty`, чтобы вернуть пустую строку.

На выходе:

- Если строка содержит пробелы, парсер возвращает строку из пробелов, например, `Just (remainingText,)` или `Just (remainingText,)` в случае нескольких пробелов.
- Если строка не содержит пробелов, парсер возвращает `Just (remainingText,)`, то есть пустую строку.
- Если первый символ не является пробелом, парсер возвращает `Nothing`.

- **Парсер пунктуации.**

Этот парсер использует `satisfy`, чтобы найти первый символ, который является одним из символов пунктуации. Для этого используется предикат `isPunctuation`, который проверяет, является ли символ одним из символов пунктуации, таких как `'.'`, `'!'`, `'?'`, `','`, `':'`, `'('`, `')'`.

На выходе:

- Если первый символ является символом пунктуации, парсер возвращает `Just (remainingText, punctuation)`, где `remainingText` — остаток строки после обработки, а `punctuation` — сам символ пунктуации.

- Если первый символ не является символом пунктуации, парсер возвращает `Nothing`.

- **Парсер для пропуска ненужных символов.**

Этот парсер использует `T.dropWhile` для пропуска всех символов, которые не являются буквами или символами пунктуации. Он продолжает обработку строки до тех пор, пока не встретит допустимый символ. После этого возвращает остаток строки. Этот парсер полезен для игнорирования символов, которые не имеют значения для анализа (например, пробелы или другие разделители).

На выходе:

- Если строка содержит ненужные символы, которые можно пропустить, парсер возвращает остаток строки, начиная с первого допустимого символа. Результат парсинга — `()`.
- Если строка не содержит допустимых символов, парсер возвращает `Nothing`.

Разбиение текста на предложения

Данный раздел описывает парсеры, которые анализируют текст, извлекая из него предложения, а затем преобразуют их в различные форматы для дальнейшей обработки.

- **Парсер для предложения.**

Этот парсер разбивает текст на слова и пунктуацию, создавая список слов, которые составляют предложение. Он использует комбинацию парсеров `word` для получения слов и `punctuation` для обработки знаков препинания. Для обработки ненужных символов применяется парсер `skipJunk`.

Парсер работает следующим образом:

- Сначала он парсит одно или несколько слов, используя парсер `word`, каждый из которых может быть отделён ненужными символами, которые пропускаются с помощью `skipJunk`.
- Затем он ожидает один или несколько символов пунктуации (например, точку, восклицательный знак, вопросительный знак).
- После этого снова пропускаются лишние символы с помощью `skipJunk`.
- Если после этих шагов не осталось текста, возвращается пустой результат (используется `empty`).

На выходе:

- Если предложение успешно разобрано, парсер возвращает `Just (remainingText, words)`, где `remainingText` — остаток строки, а `words` — список слов.
- Если не удалось распарсить предложение, возвращается `Nothing`.

- **Парсер для всех предложений.**

Этот парсер использует `some`, чтобы рекурсивно разобрать несколько предложений с помощью парсера `sentence`. Он находит все предложения в тексте, разделённые символами, и возвращает их как список предложений.

На выходе:

- Если предложения успешно разобраны, возвращается список всех предложений: `Just (remainingText, sentences)`.
 - Если не удалось распарсить предложения, возвращается `Nothing`.
- **Парсер для предложения, представленного как текст.**
 Этот парсер схож с парсером для предложения, но вместо списка слов он возвращает строку, состоящую из слов, соединённых пробелами. Он использует парсер `joinWords` для объединения списка слов в одну строку. Этот парсер может быть полезен, если нужно собрать предложение в текстовом формате, а не в виде списка отдельных слов.
 На выходе:
 - Если предложение успешно разобрано, парсер возвращает строку, представляющую предложение в текстовом виде.
 - Если не удалось распарсить предложение, возвращается `Nothing`.

- **Парсер для всех предложений в текстовом виде.**

Этот парсер использует `some`, чтобы разобрать все предложения с помощью парсера `sentenceAsText`. Он извлекает все предложения, преобразуя каждое в строку, и возвращает их в виде списка текстовых строк.

На выходе:

- Если предложения успешно разобраны, возвращается список строк, представляющих предложения.
- Если не удалось распарсить предложения, возвращается `Nothing`.

Создание N-грамм

В данном разделе представлен набор функций для генерации N-грамм из списка слов.

- **Тип `NGramMap`.** Определён тип `NGramMap` как список пар `((Text, [Text]))`, где каждый элемент представляет собой пару: слово (или несколько слов) и список слов, которые следуют за ним в тексте. Это основная структура данных для представления N-грамм.
- **Функция `toBiGrams`.** Эта функция генерирует биграммы (двуграммы) из списка слов. Биграммы — это пары, состоящие из текущего слова и следующего за ним. Функция использует стандартную функцию `zip` для формирования пар слов: первое слово из списка и следующее за ним.
- **Функция `toBiGramsJoined`.** В этой функции генерируются биграммы, где второе слово в паре состоит из двух слов, соединённых пробелом. Для этого используется вспомогательная функция `triple`, которая разбивает список на тройки подряд идущих слов. После чего слова во второй части тройки объединяются в одну строку с пробелом между ними.
- **Функция `triple`.** Эта функция используется для разбиения списка на тройки подряд идущих элементов.

- **Функция `toTriGrams`.** Генерирует триграммы из списка слов. В отличие от биграмм, триграммы состоят из двух первых слов, соединённых пробелом, и третьего слова, которое следует за ними. Триграммы представляют собой пары: первые два слова объединяются в строку, а третье слово идёт отдельно.
- **Функция `groupPairs`.** Эта функция принимает список пар и группирует их по первому элементу. Например, если пары содержат одинаковое первое слово, они будут собраны в одну группу, а вторые элементы (слова) будут собраны в список. Это позволяет собрать N-граммы для каждого слова, указывающего на следующие слова или фразы в тексте.
- **Функция `makeNGrams`.** Функция `makeNGrams` комбинирует все три предыдущие функции, генерируя полный набор N-грамм из списка слов. Она создаёт:
 - Биграммы, где каждое слово сопоставляется с последующим.
 - Биграммы, где каждое слово сопоставляется с двумя следующими словами, соединёнными пробелом.
 - Триграммы, где комбинация из двух слов сопоставляется с третьим словом.

Все эти данные собираются в общий список N-грамм.

Генерация предложений

Этот код реализует функционал для обработки текста, создания N-грамм и генерации предложений на основе этих N-грамм.

- **Функция `processText`.** Эта функция принимает текст и генерирует набор N-грамм для этого текста. В процессе работы:
 - Текст парсится с использованием парсера `allSentences`, который извлекает все предложения.
 - Все слова из предложений извлекаются и дублируются (с помощью `pub`, чтобы исключить повторения).
 - N-граммы генерируются с помощью функции `makeNGrams`, а дубликаты удаляются.
 - Добавляются слова, которые не имеют следующих слов в N-граммах. Эти слова добавляются в виде пар с пустыми списками.
 - Результат сортируется по алфавиту по ключам слов.

Выходом функции является отсортированный список N-грамм, который затем используется для генерации предложений.

- **Функция `generatePhrase`.** Эта функция генерирует фразу на основе заданного начального слова и N-грамм, используя генератор случайных чисел. Процесс генерации фразы включает следующие этапы:
 - Если первое слово не найдено в словаре N-грамм, возвращается ошибка с сообщением, что слово не найдено.
 - Если первое слово найдено в словаре, генерируется случайная длина фразы (между 2 и 15 словами).

- С помощью рекурсивной функции `generatePhraseHelper` фраза строится поэтапно, начиная с первого слова и добавляя к фразе подходящие следующие слова.

На каждом шаге выбирается следующее слово случайным образом из списка возможных слов (на основе текущего слова). Функция рекурсивно добавляет слова в итоговую фразу, пока не будет достигнута заданная длина.

- **Функция `generatePhraseHelper`.** Эта вспомогательная функция реализует рекурсивную генерацию фразы. Она продолжает добавлять слова в фразу, пока не будет достигнута максимальная длина. Если для текущего слова нет подходящих вариантов продолжения, рекурсия прекращается. Для каждого шага выбирается случайное слово из возможных, добавляется к фразе, и процесс повторяется для нового слова.
- **Функция `findLastValidWord`.** Эта функция ищет первое подходящее слово в фразе, которое встречается в словаре N-грамм. Процесс начинается с последнего слова в фразе и движется к первому. Если слово найдено в словаре, оно возвращается, в противном случае поиск продолжается для предыдущего слова. Если подходящее слово не найдено, возвращается `Nothing`.
- **Типы `DialogueResponse` и `DialogueTurn`.**
 - `DialogueResponse` — это тип, представляющий возможный ответ в диалоге: либо ошибку (`Left`), либо сгенерированную фразу (`Right`).
 - `DialogueTurn` — это кортеж, состоящий из номера модели и её ответа в виде `DialogueResponse`.

Генерация диалога моделей

Этот код реализует алгоритм для генерации диалога между двумя моделями, которые чередуют свои ответы, используя два набора N-грамм. Алгоритм генерирует последовательность фраз, опираясь на случайные выборы из N-грамм, при этом ограничивает глубину диалога и переключает модели после каждого ответа.

- **Тип `DialogueResponse`.** Этот тип представляет возможный ответ в диалоге. Он может быть:
 - `Left` — ошибка, если не удалось сгенерировать фразу.
 - `Right` — успешный ответ в виде списка слов (фразы).
- **Тип `DialogueTurn`.** Этот тип представляет один оборот в диалоге, состоящий из:
 - Номера модели (1 или 2).
 - Ответа модели, представленного типом `DialogueResponse`.
- **Функция `generateDialogue`.** Основная функция для генерации диалога. Она принимает случайный генератор (`gen`), начальное слово (`firstWord`), два набора N-грамм (`dict1` и `dict2`), и максимальную глубину диалога (`depth`).
 - Сначала генерируется первый ответ модели 1 с использованием функции `generatePhrase`.
 - Затем начинается рекурсивная генерация диалога с помощью вспомогательной функции `generateDialogueHelper`.

В результате функция возвращает список диалогов, каждый из которых состоит из номера модели и её ответа.

- **Функция `generateDialogueHelper`.** Эта функция рекурсивно генерирует следующие фразы диалога.
 - Каждый оборот в диалоге начинается с выбора следующей модели (переключение между моделью 1 и моделью 2).
 - Для каждой модели выбирается соответствующий набор N-грамм.
 - Для каждого ответа извлекается последнее слово, которое подходит для продолжения диалога (с помощью функции `findLastValidWord`).
 - Если подходящее слово найдено, генерируется новый ответ с использованием `generatePhrase`.
 - Если слово не найдено, возвращается ошибка с сообщением о невозможности продолжения диалога.
 - После каждого ответа переключается модель, и процесс повторяется до достижения максимальной глубины диалога (`depth`).

Рекурсия завершается, когда глубина диалога достигает нуля.

- **Функция `findLastValidWord`.** Эта функция ищет последнее подходящее слово из предыдущего ответа, которое существует в словаре N-грамм для следующей модели. Она перебирает слова с конца списка ответа и проверяет их наличие в словаре. Если подходящее слово найдено, оно возвращается, в противном случае возвращается `Nothing`.

4 Результаты

4.1 Часть 1. Парсер строк бинарных чисел и операций

Для проверки работы парсера используется следующий текстовый файл:

```
101010 & 00010
1101 | 1011
1010 ~ 1001
1111 & 1100
10101 | 11011
110 ~ 100
10102 & 00010
1101 $ 1011
abcdeftrhg
asd 10101 || 11011
10101 || 11011
1111 & 1100 extra
10101 |
| 11011
110 ~ ~ 100
10101 | 11011
101010 & 00010
101010 & 00010
1101 | 1011
101010&00010
```

Результаты прохождения тестов, описанных выше, представлены на Рис. 1.

```

Enter file name:
test.txt
101010 & 00010 = 10
1101 | 1011 = 1111
1010 ^ 1001 = 11
1111 & 1100 = 1100
Wrong string
Wrong string
10101 | 11011 = 11111
110 ^ 100 = 10
Wrong string
Wrong string
Wrong string
Wrong string
Wrong string
1111 & 1100 = 1100
Wrong string
Wrong string
Wrong string
10101 | 11011 = 11111
101010 & 00010 = 10
101010 & 00010 = 10
1101 | 1011 = 1111
101010 & 00010 = 10

```

Рис. 1. Результаты прохождения тестов

При запуске программы, пользователь вводит путь до файла, после чего для каждой строки из файла выводится результат вычисления выражения или сообщение о неверном формате строки.

4.2 Часть 2. Генератор предложений с помощью N-грамм

Для проверки работы парсера используются текстовые 2 файла содержащих: главы 1 и 2 произведения «Старуха Изергиль» на английском языке Максима Горького и стенограмму эпизода 304 «Твик против Крейга» мультсериала «Южный парк».

Сначала пользователь вводит начальное слово, с которого будет начинаться предложение первой модели (Рис. 2).

```

Enter the first word for dialogue:
just

```

Рис. 2. Ввод начального слова

Затем пользователь вводит глубину сообщений в диалоге (Рис. ??).

```
Enter the first word for dialogue:
just
Enter the number of exchanges:
5
```

Рис. 3. Ввод глубины сообщений диалога

После этого программа выводит сгенерированный диалог состоящий из реплик двух моделей (Рис. ??).

```
Generated dialogue:
Model 1: just goes Larra meaning the
Model 2: the daydream
Model 1: the hunt was over
Model 2: over there and Cartman I wonder why can't you open coffin Richard
Model 1: open eyes and blood stained mouth and expected me
```

Рис. 4. Результат работы программы

По выводу программы видно, что каждая реплика начинается с последнего слова предыдущей реплики или, при отсутствии такого слова в словаре модели, используется предпоследнее слово и так далее.

Заключение

Часть 1

В данной работе был разработан синтаксический анализатор для разбора строк, содержащих битовые выражения с операциями AND, OR и XOR. Анализатор принимает на вход файл формата .txt, содержащий такие выражения, и выводит результаты их вычислений на экран. Реализованный парсер позволяет корректно обрабатывать строки с указанными операциями, обеспечивая удобство использования для конечного пользователя.

Часть 2

В ходе выполнения данного задания была разработана программа, которая решает задачу синтаксического анализа текста, построения модели N-грамм и генерации фраз на основе этих моделей. В рамках первой части задачи был реализован алгоритм чтения текста из файла, разбивки его на предложения и удаления символов пунктуации и цифр. Во второй части было создано два типа моделей N-грамм: биграммы и триграммы, а также сформирован словарь, который сохранен в файл формата .txt.

Третья часть задачи включала реализацию взаимодействия с пользователем, где пользователь мог вводить одно слово или пару слов, а программа генерировала фразу случайной длины от 2 до 15 слов. Если заданного пользователем слова не находилось среди ключей словаря, выводилось соответствующее сообщение.

Четвертая часть заключалась в организации диалога между двумя моделями N-грамм, созданными на основе различных текстов Максима Горького. Пользователю предоставлялась возможность задать начальное слово или пару слов и указать количество сообщений, которыми будут обмениваться модели. Модели строят свои ответы на основе последнего слова предыдущего сообщения оппонента.

Результаты работы программы демонстрируют корректность алгоритмов обработки текста и построения моделей N-грамм. Полученная система может быть использована для создания простых чат-ботов или других приложений, основанных на анализе естественного языка.

Приложение А. Парсер бинарных операций

Lib.hs

```
1 module Lib
2   ( parseAndPrint
3     , readFileLines
4   ) where
5
6 import Data.Text (Text)
7 import qualified Data.Text as T
8 import qualified Data.ByteString as B
9 import qualified Data.Text.Encoding as TE
10 import qualified Data.Text.IO as TIO
11 import Control.Applicative
12 import Data.Char (digitToInt, intToDigit)
13 import Data.Bits ((.&.), (.|.), xor)
14 import Numeric (showIntAtBase)
15
16 newtype Parser a = Parser { runParser :: Text -> Maybe (Text, a) }
17
18 instance Functor Parser where
19   -- хотим применить функ над результатом парсера p
20   fmap func (Parser p) = Parser f where
21     -- парсер f возвращает:
22     f origText = case p origText of
23       Nothing -> Nothing -- Nothing, если парсер p возвращает Nothing
24       Just (remainingP, resP) -> Just (remainingP, func resP) -- (остаток, resP
25         ↳ обработанный функцией func), если p возвращает (остаток, resP)
26
27 instance Applicative Parser where
28   -- возвращаем всегда (изначальная строка, передаваемое значение)
29   pure text = Parser (\orig -> Just(orig, text))
30
31   -- хотим чтобы был парсер, который применяет к остатку 1 парсера 2 парсер,
32   -- а затем применяет 1 парсер ко 2
33   (Parser u) <*> (Parser v) = Parser f where
34     f origText = case u origText of
35       Nothing -> Nothing
36       -- remainingU - остаток 1 парсера
37       Just (remainingU, resU) -> case v remainingU of
38         Nothing -> Nothing
39         -- remainingV - итоговый остаток, resU применяем над resV
40         Just (remainingV, resV) -> Just (remainingV, resU resV)
41
42 instance Alternative Parser where
43   -- парсер всегда возвращающий Nothing
44   empty = Parser $ \_ -> Nothing
45
46   Parser u <|> Parser v = Parser f where
```

```

46     -- пытаемся применить парсер u
47     f origText = case u origText of
48         -- если он вернул Nothing, то применяю парсер v
49         Nothing -> v origText
50         -- если вернул какой то результат, то оставляем результат
51         res -> res
52
53 isBinaryChar :: Char -> Bool
54 isBinaryChar c =
55     c == '0' || c == '1'
56
57 isOperation :: Char -> Bool
58 isOperation c =
59     c == '&' || c == '|' || c == '?'
60
61 isSpace :: Char -> Bool
62 isSpace c = c == ' '
63
64 satisfy :: (Char -> Bool) -> Parser Char
65 satisfy pr = Parser f where
66     -- берем первый символ
67     f cs = case T.uncons cs of
68         Nothing -> Nothing
69         -- если первый элемент соответствует предикату pr
70         Just (fstChar, remainingText)
71             | pr fstChar -> Just (remainingText, fstChar) -- то возвращаем (остаток,
72                 ↳ подходящий первый элемент)
73             | otherwise -> Nothing
74     f _ = Nothing
75
76 binary :: Parser Text
77 binary = Parser $ \text ->
78     -- если первый элемент
79     case runParser (satisfy isBinaryChar) text of
80         Nothing -> Nothing -- не подошел, то строка очевидно сразу не подходит
81         -- (рекурсивно) если парсер на остатке
82         Just (remaining, c) -> case runParser binary remaining of
83             Nothing -> Just (remaining, T.singleton c) -- первый элемент не подошел, то
84                 ↳ берем (старый остаток, единственный подошедший элемент)
85             Just (remaining', rest) -> Just (remaining', T.cons c rest) -- сработал
86                 ↳ штатно, то (новый остаток, добавляем подошедший элемент к остальным)
87
88 binToInt :: Text -> Int
89 binToInt text = T.foldl1' (\acc c -> acc * 2 + digitToInt c) 0 text
90
91 intToBin :: Int -> Text
92 intToBin n = T.pack (showIntAtBase 2 intToDigit n "")
93
94 binaryInt :: Parser Int
95 binaryInt = binToInt <$> binary

```

```

93
94 operation :: Parser Char
95 operation = satisfy isOperation
96
97 oneSpace :: Parser Char
98 oneSpace = satisfy isSpace
99
100 spaces :: Parser Text
101 spaces = (T.cons) <$> oneSpace <*> spaces <|> pure T.empty
102
103 applyOperation :: Char -> Int -> Int -> Int
104 applyOperation op
105   | op == '&' = (.&.)
106   | op == '|' = (.|. )
107   | op == '?' = xor
108
109 binaryExpression :: Parser Int
110 binaryExpression =
111   (\_ b1 _ op _ b2 -> applyOperation op b1 b2)
112   <$> spaces
113   <*> binaryInt
114   <*> spaces
115   <*> operation
116   <*> spaces
117   <*> binaryInt
118
119 binaryExpressionFormatted :: Parser Text
120 binaryExpressionFormatted =
121   (\_ b1 _ op _ b2 -> formatExpression b1 op b2)
122   <$> spaces
123   <*> binary
124   <*> spaces
125   <*> operation
126   <*> spaces
127   <*> binary
128
129 formatExpression :: Text -> Char -> Text -> Text
130 formatExpression b1 op b2 =
131   let int1 = binToInt b1
132       int2 = binToInt b2
133       result = applyOperation op int1 int2
134   in b1 <> (T.pack " ") <> T.singleton op <> (T.pack " ") <> b2 <> (T.pack " = ") <>
135     ↪ intToBin result
136
137 readFileLines :: FilePath -> IO [T.Text]
138 readFileLines filePath = fmap (T.lines . TE.decodeUtf8) (B.readFile filePath)
139
140 parseAndPrint :: T.Text -> IO ()
141 parseAndPrint input =
142   case runParser binaryExpressionFormatted input of

```



```
142     Just (_, result) -> TIO.putStrLn result
143     Nothing -> putStrLn "Wrong string"
```

Main.hs

```
1  module Main (main) where
2
3  import Lib
4  import System.IO
5  import qualified Data.Text.Encoding as TE
6  import qualified Data.ByteString as B
7
8  main :: IO ()
9  main = do
10     hSetEncoding stdout utf8
11     putStrLn "Enter file name:"
12     fileName <- getLine
13     lines <- readFileLines fileName
14     -- mapM_ (TE.decodeUtf8 . B.putStrLn) lines
15     mapM_ parseAndPrint lines
```

Приложение Б. Генератор предложений

Lib.hs

```
1 module Lib
2   ( allSentences
3   , runParser
4   , generatePhrase
5   , processText
6   , NGramMap
7   , generateDialogue
8   ) where
9
10 import Control.Applicative
11 import Data.Text (Text)
12 import qualified Data.Text as T
13 import Data.Char (isLetter, isSpace)
14 import Data.List (sortBy, groupBy, nub)
15 import System.Random (RandomGen, randomR, next, split)
16
17 -- Parser
18 -- instances
19
20 newtype Parser a = Parser { runParser :: Text -> Maybe (Text, a) }
21
22 instance Functor Parser where
23   -- хотим применить функ над результатом парсера p
24   fmap func (Parser p) = Parser f where
25     -- парсер f возвращает:
26     f origText = case p origText of
27       Nothing -> Nothing -- Nothing, если парсер p возвращает Nothing
28       Just (remainingP, resP) -> Just (remainingP, func resP) -- (остаток, resP
29         ↳ обработанный функцией func), если p возвращает (остаток, resP)
30
31 instance Applicative Parser where
32   -- возвращаем всегда (изначальная строка, передаваемое значение)
33   pure text = Parser (\orig -> Just(orig, text))
34
35   -- хотим чтобы был парсер, который применяет к остатку 1 парсера 2 парсер,
36   -- а затем применяет 1 парсер ко 2
37   (Parser u) <*> (Parser v) = Parser f where
38     f origText = case u origText of
39       Nothing -> Nothing
40       -- remainingU - остаток 1 парсера
41       Just (remainingU, resU) -> case v remainingU of
42         Nothing -> Nothing
43         -- remainingV - итоговый остаток, resU применяем над resV
44         Just (remainingV, resV) -> Just (remainingV, resU resV)
45
46 instance Alternative Parser where
```

```

46  -- парсер всегда возвращающий Nothing
47  empty = Parser $ \_ -> Nothing
48
49  Parser u <|> Parser v = Parser f where
50      -- пытаемся применить парсер u
51      f origText = case u origText of
52          -- если он вернул Nothing, то применяем парсер v
53          Nothing -> v origText
54          -- если вернул какой то результат, то оставляем результат
55          res -> res
56
57  -- functions
58
59  isPunctuation :: Char -> Bool
60  isPunctuation c = c `elem` ['.', '!', '?', ';', ':', '(', ')']
61
62  joinWords :: [Text] -> Text
63  joinWords = T.intercalate (T.pack " ")
64
65  -- parsers
66
67  satisfy :: (Char -> Bool) -> Parser Char
68  satisfy pr = Parser f where
69      -- берем первый символ
70      f cs = case T.uncons cs of
71          Nothing -> Nothing
72          -- если первый элемент соответствует предикату pr
73          Just (fstChar, remainingText)
74              | pr fstChar -> Just (remainingText, fstChar) -- то возвращаем (остаток,
75                  ↳ подходящий первый элемент)
76              | otherwise -> Nothing
77      f _ = Nothing
78
79  isWordChar :: Char -> Bool
80  isWordChar c = isLetter c || c == '\''
81
82  word :: Parser Text
83  word = Parser f where
84      f input = case runParser (some (satisfy isWordChar)) input of
85          Nothing -> Nothing
86          Just (remaining, chars) ->
87              -- проверяем что слово начинается или заканчивается апострофом
88              let word = T.pack chars
89                  cleanWord = T.dropAround (== '\') word
90                  in if T.null cleanWord
91                      then Nothing -- если слово состоит только из апострофов, то возвращаем
92                          ↳ Nothing
93                      else Just (remaining, word)
94
95  oneSpace :: Parser Char

```

```

94 oneSpace = satisfy isSpace
95
96 spaces :: Parser Text
97 spaces = (T.cons) <$> oneSpace <*> spaces <|> pure T.empty
98
99 punctuation :: Parser Char
100 punctuation = satisfy isPunctuation
101
102 skipJunk :: Parser ()
103 skipJunk = Parser f where
104     f input = Just (T.dropWhile (\c -> not (isLetter c || isPunctuation c)) input, ())
105
106 sentence :: Parser [Text]
107 sentence = (\words _ -> words)
108     <$> some (word <*> skipJunk)
109     <*> some punctuation
110     <*> skipJunk
111     <|> empty
112
113 allSentences :: Parser [[Text]]
114 allSentences = some sentence
115
116 sentenceAsText :: Parser Text
117 sentenceAsText = (\words _ -> joinWords words)
118     <$> some (word <*> skipJunk)
119     <*> some punctuation
120     <*> skipJunk
121     <|> empty
122
123 allSentencesAsText :: Parser [Text]
124 allSentencesAsText = some sentenceAsText
125
126 -- N-gram
127
128 type NGramMap = [(Text, [Text])]
129
130 toBiGrams :: [Text] -> [(Text, Text)]
131 toBiGrams words = zip words (tail words)
132
133 toBiGramsJoined :: [Text] -> [(Text, Text)]
134 toBiGramsJoined ws =
135     [(w1, T.concat [w2, T.pack " ", w3]) | (w1,w2,w3) <- triple ws]
136
137 triple :: [a] -> [(a, a, a)]
138 triple (x:y:z:rest) = (x,y,z) : triple (y:z:rest)
139 triple _ = []
140
141 toTriGrams :: [Text] -> [(Text, Text)]
142 toTriGrams ws =
143     [(T.concat [w1, T.pack " ", w2], w3) | (w1,w2,w3) <- triple ws]

```

```

144
145
146 groupPairs :: [(Text, Text)] -> NGramMap
147 groupPairs pairs = map (\group -> (fst $ head group, map snd group))
148                     $ groupBy (\x y -> fst x == fst y)
149                     $ sortBy (\x y -> compare (fst x) (fst y)) pairs
150
151 makeNGrams :: [Text] -> [(Text, Text)]
152 makeNGrams words =
153     -- создаем биграммы вида (word -> next word)
154     toBiGrams words ++
155     -- создаем биграммы вида (word -> two next words joined)
156     toBiGramsJoined words ++
157     -- создаем триграммы вида (two words -> next word)
158     toTriGrams words
159
160 processText :: Text -> NGramMap
161 processText text = case runParser allSentences text of
162     Nothing -> []
163     Just (_, sentences) -> let
164         -- берем все слова из всех предложений
165         allWords = nub $ concat sentences
166         -- создаем n-граммы из всех предложений и удаляем дубликаты
167         allNGrams = groupPairs
168                     $ nub -- удаляем дубликаты
169                     $ concatMap makeNGrams sentences
170         -- добавляем слова, которые не имеют следующих слов
171         singleWords = map (\w -> (w, []))
172                         $ filter (\w -> not $ any (\(prefix, _) -> prefix == w) allNGrams)
173                         $ allWords
174         in sortBy (\x y -> compare (fst x) (fst y))
175             $ allNGrams ++ singleWords
176
177 generatePhrase :: RandomGen g => g -> Text -> NGramMap -> Either Text [Text]
178 generatePhrase gen firstWord nGrams =
179     -- если первого слова нет в словаре, то возвращаем ошибку (Left ошибка, Right [слова
180     -- ↪ для фразы])
181     case lookup firstWord nGrams of
182         Nothing -> Left $ T.concat [T.pack "Word '", firstWord, T.pack "' not found in
183         ↪ the dictionary"]
184         -- если первое слово есть в словаре, то начинаем генерировать фразу
185         Just nextWords ->
186             let (targetLength, newGen) = randomR (2, 15) gen
187             in Right $ generatePhraseHelper newGen [firstWord] firstWord nGrams
188                 ↪ targetLength
189
190 generatePhraseHelper :: RandomGen g => g -> [Text] -> Text -> NGramMap -> Int -> [Text]
191 generatePhraseHelper gen acc lastKey nGrams targetLength
192     | length acc >= targetLength = take targetLength acc -- достигли максимальной
193     ↪ длины

```

```

190 | null possibleNextWords = acc -- нет больше слов
191 | otherwise =
192     let (idx, newGen) = randomR (0, length possibleNextWords - 1) gen
193         nextWord = possibleNextWords !! idx
194         -- если два слова, то разбиваем на два
195         nextWords = T.words nextWord
196         -- и добавляем в список слов фразы
197         newAcc = acc ++ nextWords
198         -- используем следующее слово (или два слова) как ключ для следующего шага
199         newKey = nextWord
200     -- рекурсивно вызываем функцию для следующего шага
201     in generatePhraseHelper newGen newAcc newKey nGrams targetLength
202 where
203     -- находим список слов по данному ключу (Just [Text]), если такого ключа в словаре
204     -- ↪ нет, то Nothing
205     possibleNextWords = maybe [] id $ lookup lastKey nGrams
206
207 -- Найми первое подходящее слово от конца фразы, которое есть в словаре
208 findLastValidWord :: [Text] -> NGramMap -> Maybe Text
209 findLastValidWord [] _ = Nothing
210 findLastValidWord (word:rest) nGrams =
211     case lookup word nGrams of
212         Just _ -> Just word
213         Nothing -> findLastValidWord rest nGrams
214
215 type DialogueResponse = Either Text [Text] -- Left - ошибка, Right - фраза
216 type DialogueTurn = (Int, DialogueResponse) -- (номер модели, ответ)
217
218 generateDialogue :: RandomGen g =>
219     g ->
220     Text ->
221     NGramMap ->
222     NGramMap ->
223     Int ->
224     [DialogueTurn] -- возвращает список (номер модели, ответ)
225 generateDialogue gen firstWord dict1 dict2 depth =
226     let firstResponse = generatePhrase gen firstWord dict1
227         (_, newGen) = next gen
228     in (1, firstResponse) : generateDialogueHelper newGen firstResponse 1 dict1 dict2
229     ↪ depth []
230
231 generateDialogueHelper :: RandomGen g =>
232     g ->
233     DialogueResponse -> -- последний ответ
234     Int -> -- номер модели
235     NGramMap -> -- первый словарь
236     NGramMap -> -- второй словарь
237     Int -> -- остаток глубины
238     [DialogueTurn] -> -- накопленный диалог
239     [DialogueTurn] -- финальный диалог

```

```

238 generateDialogueHelper _ _ _ _ 0 acc = reverse acc
239 generateDialogueHelper gen lastResponse speaker dict1 dict2 depth acc =
240     let
241         -- следующая модель
242         nextSpeaker = if speaker == 1 then 2 else 1
243
244         -- словарь для следующей модели
245         currentDict = if nextSpeaker == 1 then dict1 else dict2
246
247         -- последнее подходящее слово из предыдущего ответа
248         lastWords = case lastResponse of
249             Right phrase -> reverse phrase -- если слово есть
250             Left _ -> [] -- если нет, то пустой список
251
252         -- последнее подходящее слово из предыдущего ответа
253         nextWord = findLastValidWord lastWords currentDict
254
255         -- следующий ответ (даже если не нашли подходящего слова)
256         nextResponse = case nextWord of
257             Nothing -> Left (T.pack "No valid word found to continue dialogue")
258             Just word ->
259                 let (newGen1, _) = split gen
260                 in generatePhrase newGen1 word currentDict
261
262         (_, newGen2) = split gen
263     in generateDialogueHelper
264         newGen2
265         nextResponse
266         nextSpeaker
267         dict1
268         dict2
269         (depth - 1)
270         ((nextSpeaker, nextResponse) : acc)

```

Main.hs

```

1  module Main (main) where
2
3  import Lib
4  import qualified Data.Text as T
5  import qualified Data.Text.IO as TIO
6  import System.Random (newStdGen)
7
8  main :: IO ()
9  main = do
10
11      input1 <- TIO.readFile "input.txt"
12      input2 <- TIO.readFile "input2.txt"
13

```

```

14 let nGrams1 = processText input1
15 let nGrams2 = processText input2
16
17 -- print nGrams1
18 -- print "\n"
19 -- print nGrams2
20 -- print "\n"
21
22 TIO.putStrLn (T.pack "Enter the first word for dialogue:")
23 firstWord <- T.strip <$> TIO.getLine
24
25 TIO.putStrLn (T.pack "Enter the number of exchanges:")
26 depthStr <- getLine
27 let depth = (read depthStr :: Int) - 1
28
29 gen <- newStdGen
30 let dialogue = generateDialogue gen firstWord nGrams1 nGrams2 depth
31 TIO.putStrLn (T.pack "\nGenerated dialogue:")
32 mapM_ (printDialogueTurn . formatDialogueTurn) dialogue
33 where
34   formatDialogueTurn (speaker, response) =
35     (T.pack $ "Model " ++ show speaker ++ ": ",
36      case response of
37        Left err -> err
38        Right phrase -> T.unwords phrase)
39
40   printDialogueTurn (prefix, message) =
41     TIO.putStr prefix >> TIO.putStrLn message

```