## Лексический анализ

JICKON ICOKNIN GITGINIS

Лекция №3

#### Книга "Syntactic Structures", Avram Noam Chomsky

- Вышла в 1957 году
- Фундамент всех формальных языков

#### Книга "Syntactic Structures", Avram Noam Chomsky

- Вышла в 1957 году
- Фундамент всех формальных языков

Что есть в книге:

- 1. Порождающие грамматики (Generative Grammars)
- 2. Гипотеза о врождённых способностях к усвоению языка через порождающие грамматики

Виды разбора

в русском языке

Косил косой косой косой

#### Косил косой косой косой

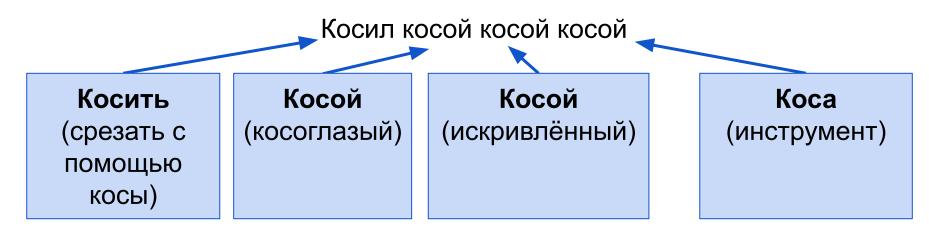
#### Лексический разбор:

- 1. Определяем смысл каждого слова (по словарю)
- 2. Для многозначных слов определяем значение

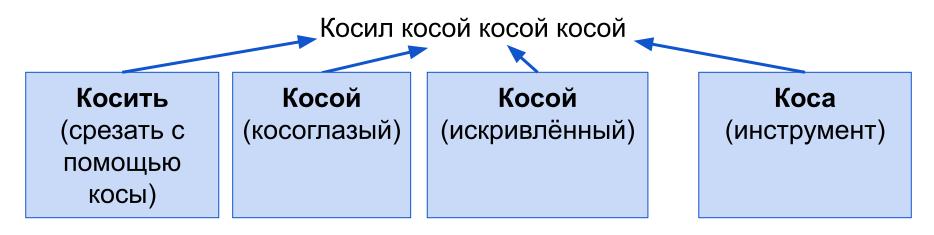
#### Косил косой косой косой

#### Лексический разбор:

- 1. Определяем смысл каждого слова (по словарю)
- 2. Для многозначных слов определяем значение
- Лемматизация замена всех слов на леммы
- Лемма словарная форма слова



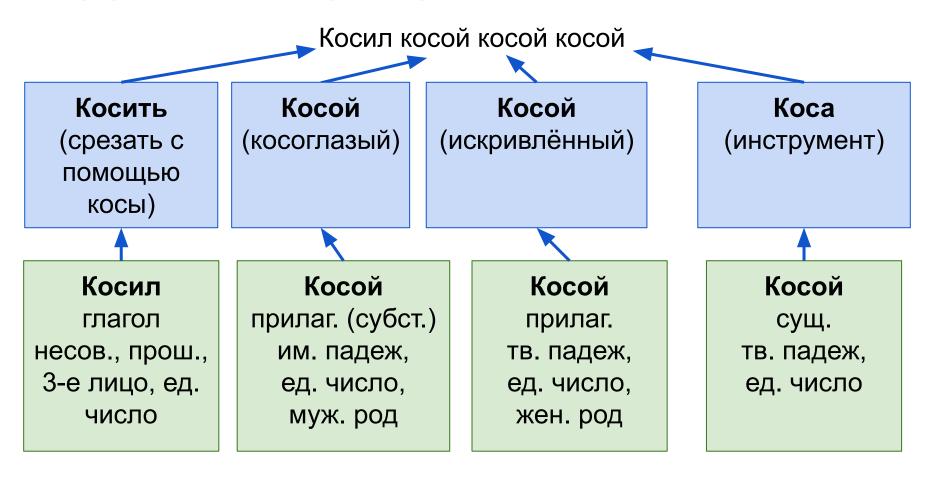
#### Морфологический разбор слов



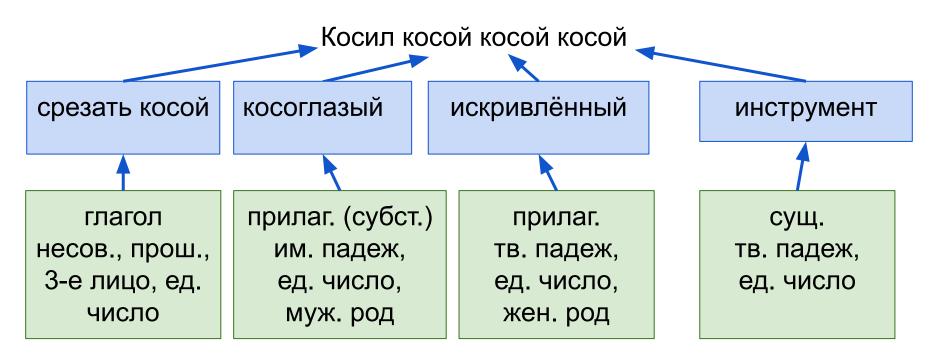
#### Морфологический разбор:

- 1. Определяем часть речи слова
- 2. Определяем признаки
  - Глаголы имеют вид, время, лицо, число
  - Прилагательные имеют род, падеж, число

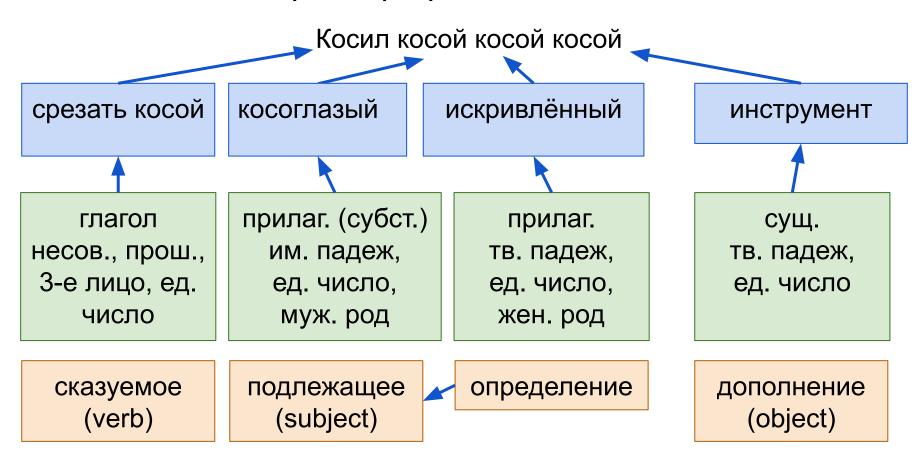
### Морфологический разбор слов



#### Синтаксический разбор предложения



#### Синтаксический разбор предложения



#### Виды анализа для определения синтаксиса

#### Естественные языки

Формальные языки

???

- 1. Лексический анализ
- 2. Морфологический анализ
- 3. Синтаксический анализ

#### Виды анализа для определения синтаксиса

#### Естественные языки

- 1. Лексический анализ
- 2. Морфологический анализ
- 3. Синтаксический анализ

#### Формальные языки

- 1. Лексический анализ
- 2. Синтаксический анализ

Морфологический анализ не нужен, потому что лексемы однозначны.

Лексическая структура SQL

SELECT first\_name, last\_name FROM student;

#### Лексемы:

- 1. Ключевые слова SELECT, FROM
- 2. Идентификаторы first\_name, last\_name
- 3. Разделители ";", ","

SELECT 2 \* 3.14159 \* radius FROM circle;

Лексемы:

???

SELECT 2 \* 3.14159 \* radius FROM circle;

#### Лексемы:

- 1. Ключевые слова SELECT, FROM
- 2. Идентификатор radius
- 3. Литералы чисел 2, 3.14159
- 4. Разделитель "**;**"
- 5. Оператор "\*"

-- Outputs squares for all circles
SELECT 2 \* 3.14159 \* radius FROM circle;

Лексемы: ???

-- Outputs squares for all circles

SELECT 2 \* 3.14159 \* radius FROM circle;

Комментарий — не лексема, потому что не важен для синтаксиса

-- Outputs squares for all circles

SELECT 2 \* 3.14159 \* radius FROM circle;

Комментарий — не лексема, потому что не важен для синтаксиса

#### Лексемы:

- 1. Ключевые слова SELECT, FROM
- 2. Идентификатор radius
- 3. Литералы чисел 2, 3.14159
- 4. Разделитель "**;**"
- 5. Оператор "\*"

## BNF u EBNF

#### Идентификатор

Идентификатор в языке SQL:

- 1. Начинается с буквы или "\_"
- 2. Далее могут идти буквы, цифры или "\_"

## BNF — форма Бэкуса-Haypa (Backus-Naur) identifier ::= head | head tail

identifier ::= head | head tail
head ::= letter | "\_"
tail ::= non\_head | non\_head tail

non\_head ::= letter | digit | "\_"

## BNF — форма Бэкуса-Наура (Backus-Naur)

```
identifier ::= head | head tail
head ::= letter | "_"
tail ::= non_head | non_head tail
non_head ::= letter | digit | "_"
```

В BNF правило содержит:

- альтернативы, соединённые через "|"
- символы в двойных кавычках (терминалы)
- список других правил (нетерминалы)

## EBNF — расширенная форма Бэкуса-Наура

identifier = head, { non\_head } ;
head = letter | "\_" ;
non\_head = letter | digit | "\_";

# EBNF — расширенная форма Бэкуса-Наура identifier = head, { non\_head } ; head = letter | "\_" ;

Что добавляет EBNF:

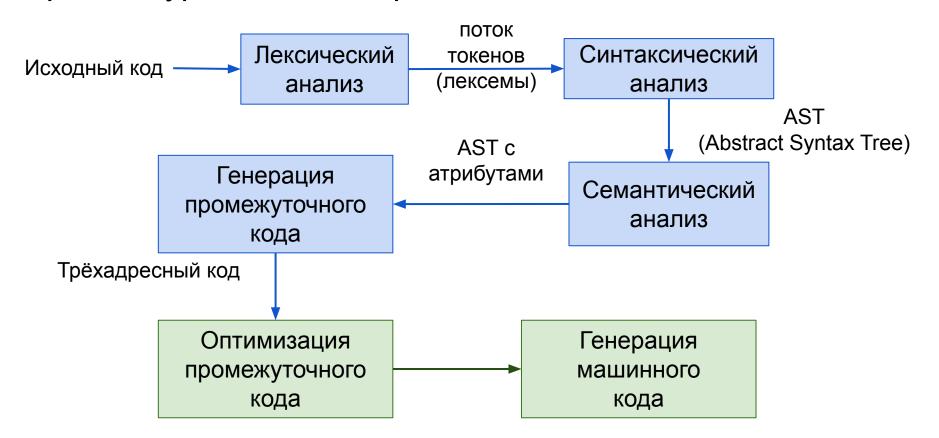
- { symbol\_1 symbol\_2 } для повтора 0, 1 или много раз
- [symbol\_1 symbol\_2] для опциональных частей
- (symbol\_1 symbol\_2) для группировки

non\_head = letter | digit | "\_";

Лексический анализатор

в разработке компилятора

## Архитектура компилятора



## Фронтенд компилятора — слои и функции

	Выражения	Блоки кода	Ветвления	Циклы	Функции
Лексика	X	X	X	×	X
Синтаксис	X	X	X	×	X
Семантика	X	X	X	×	X
Генерация IR	×	X	×	X	×

#### Фронтенд компилятора — разработка по слоям

	Выражения	Блоки кода	Ветвления	Циклы	Функции
Лексика	V	V	V	V	V
Синтаксис	X	X	X	X	×
Семантика	X	X	X	X	×
Генерация IR	×	×	×	X	×

## Фронтенд компилятора — разработка по функциям

	Выражения	Блоки кода	Ветвления	Циклы	Функции
Лексика		X	×	×	X
Синтаксис		×	×	×	X
Семантика		×	×	×	X
Генерация IR		×	×	X	×

#### Фронтенд компилятора — разработка по функциям

	Выражения	Блоки кода	Ветвления	Циклы	Функции
Лексика		X	X	×	X
Синтаксис	V	X	X	×	X
Семантика		X	X	×	X
Генерация IR		X	×	X	×
		1			

«Представьте, что история – это тонкий вертикальный срез, проходящий сквозь горизонтальные слои проекта.»

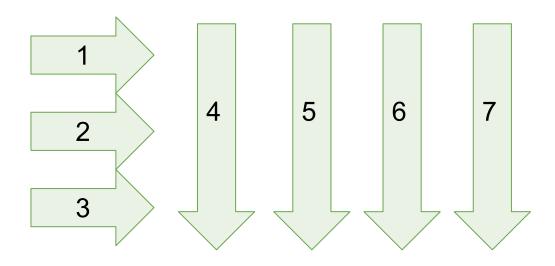
— Роберт Мартин, «Чистый Agile»

## Фронтенд компилятора — смешанный подход

	Выражения	Блоки кода	Ветвления	Циклы	Функции
Лексика	V	V		×	X
Синтаксис		X	X	×	X
Семантика		X	×	×	X
Генерация IR		×	×	X	×

#### Когда предметная область незнакома

- 1. Лучше разрабатывать по функциям, а не слоям
- 2. ... но в начале проекта строим каркас по слоям



Каноничный TDD

## Что такое модульный тест?

- Модульный тест — проверяет модуль через внутренние интерфейсы

#### Что такое модуль?

- Модульный тест проверяет модуль через внутренние интерфейсы
- Модуль класс, функция или файл
  - имеет интерфейс и реализацию

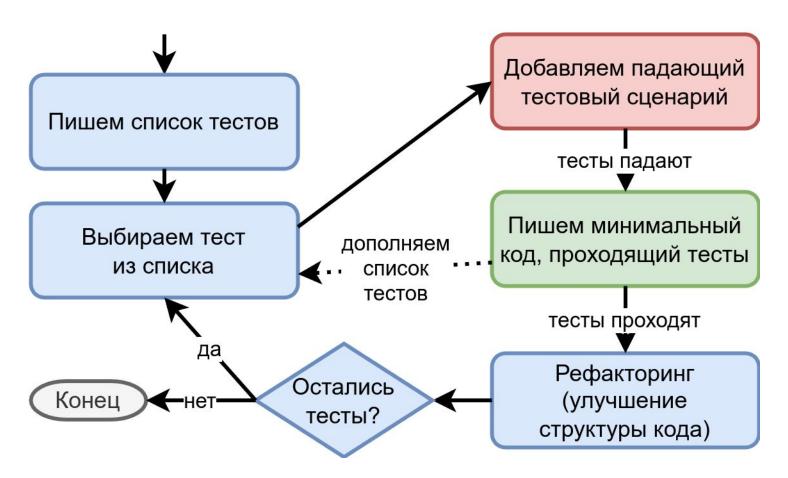
#### Что такое TDD?

- Модульный тест проверяет модуль через внутренние интерфейсы
- Модуль класс, функция или файл
  - имеет интерфейс и реализацию
- TDD (Test Driven Development) тесты ведут кодирование
  - тесты определяют дизайн кода
  - тесты задают спецификацию поведения
  - тесты создают качество

#### Что такое TDD?

- Модульный тест проверяет модуль через внутренние интерфейсы
- Модуль класс, функция или файл
  - имеет интерфейс и реализацию
- TDD (Test Driven Development) тесты ведут кодирование
  - тесты определяют дизайн кода
  - тесты задают спецификацию поведения
  - тесты создают качество
- Можно иметь модульные тесты и не иметь TDD

#### Каноничный TDD

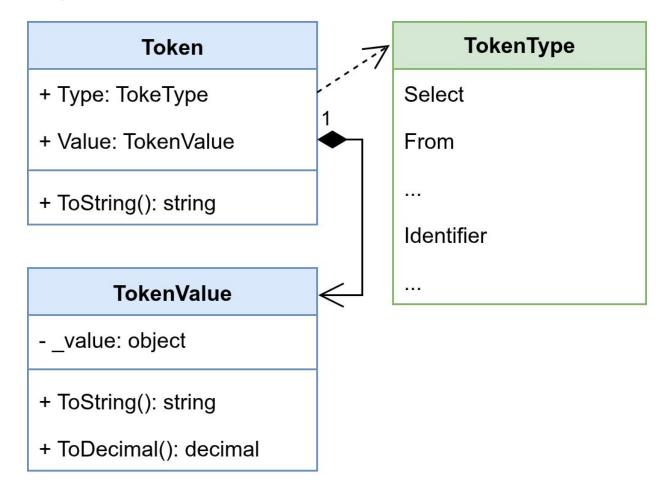


Проектирование Lexer

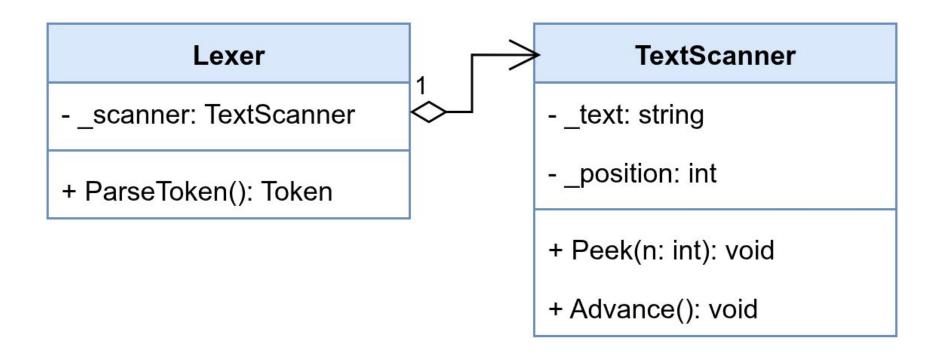
#### Ключевые решения

- 1. Целиком читаем исходный код в память
- 2. Разделяем на лексемы вручную
- 3. Выделяем абстракции

#### Token — представляет данные лексемы



# TextScanner — абстракция над строкой



Элементы реализации

#### Совет: пишите правило в комментарии

```
/// <summary>
/// Распознаёт идентификатор по правилам:
/// identifier = [letter | '_' ] { letter | digit | '_' }
/// letter = "a" | "b" | ... | "z" | unicode_letter
/// digit = "0" | "1" | ... | "9"
/// unicode_letter — любая буква Unicode.
/// </summary>
private Token ParseIdentifierOrKeyword()
```

# Ключевое слово — это зарезервированное имя

```
if (Keywords.TryGetValue(
  value.ToUpper(CultureInfo.InvariantCulture),
  out TokenType type)
   return new Token(type); // Ключевое слово
return new Token(
  TokenType.Identifier, new TokenValue(value)
```

## Разбор целых чисел

```
decimal value = GetDigitValue(_scanner.Peek());
_scanner.Advance();
for (char c = _scanner.Peek();
    char.IsAsciiDigit(c);
    c = _scanner.Peek())
  value = value * 10 + GetDigitValue(c);
  _scanner.Advance();
```

### Пропуск комментариев

```
private void SkipWhiteSpacesAndComments() {
  do {
       SkipWhiteSpaces();
   } while (
       TryParseMultilineComment()
       || TryParseSingleLineComment()
```

#### Исходный код примера

https://sourcecraft.dev/sshambir-public/memsql

- docs/ описание лексической структуры
- src/ исходный код
- tests/ тесты

# Конец Вопросы?