

CSS125

notes and stuff

Petcham

Copyright © Petcham 2025

Table of contents

Programming Paradigms Study Guide	4
CSS 125 - Principles of Programming Languages	4
What is a Programming Paradigm?	4
The Four Main Programming Paradigms	4
Imperative vs. Declarative - A Key Distinction	10
Language Categories and Examples	11
Key Programming Concepts Explained	12
How Different Paradigms Solve the Same Problem	15
Choosing the Right Paradigm	16
Common Misconceptions	18
Study Tips	18
Conclusion	19
Lexical Analysis	20
1. Introduction to Lexical Analysis	20
2. The Role of the Lexical Analyzer	20
3. Why Separate Lexical Analysis from Parsing?	22
4. Tokens, Patterns, and Lexemes	23
5. Lexical Errors and Error Handling	27
6. Token Specifications and Regular Expressions	29
7. String Set Operations	34
8. Regular Expression Operator Precedence	35
9. Real-World Applications	38
Summary and Key Takeaways	42
Lexical Analysis Implementation	44
What is Lexical Analysis?	44
Recognition of Tokens	44
Transition Diagrams	46

Implementation Details	48
Identifiers and Symbol Table Integration	52
Finite State Automata (FSA)	53
Practical Examples and Applications	56
Key Concepts Summary	57
Syntax Analysis	60
What is Syntax?	60
Role of the Parser	61
Classes of Languages	63
Error Handling	65
Error Recovery Strategies	67
Practical Examples and Applications	74
Study Tips and Key Takeaways	77

Programming Paradigms Study Guide

CSS 125 - Principles of Programming Languages

What is a Programming Paradigm?

Think of a **programming paradigm** like different ways of giving directions to someone.

Simple Definition: A programming paradigm is a fundamental style or approach to writing computer programs - it's like choosing whether to write a recipe as step-by-step instructions, or as a collection of ingredients that work together, or as a set of rules that automatically create the final dish.

Real-World Analogy:

- **Procedural** = Like following a cooking recipe step-by-step: "First crack the eggs, then heat the pan, then pour the eggs..."
 - **Object-Oriented** = Like organizing a kitchen where different appliances (objects) have their own jobs: the oven bakes, the blender mixes, the refrigerator stores
 - **Functional** = Like using mathematical formulas: "If you put these ingredients into this formula, you'll get this result"
 - **Dynamic/Scripting** = Like having a flexible assistant that can adapt and change what it does based on the situation
-

The Four Main Programming Paradigms

1. Procedural Programming

What it is: Procedural programming is like writing a very detailed, step-by-step instruction manual.

Key Characteristics:

- **Imperative paradigm** - You tell the computer exactly what to do and how to do it
- **Step-by-step execution** - Programs run from top to bottom, line by line
- **Functions and modules** - You can group related instructions together
- **Any procedure can be called anytime** - Like being able to jump to any page in your instruction manual

Real-World Example: Imagine you're teaching someone to make a sandwich:

```
1. Get two slices of bread
2. Open the peanut butter jar
3. Get a knife
4. Spread peanut butter on one slice
5. Get jelly
6. Spread jelly on the other slice
7. Put the slices together
8. Cut in half
9. Serve
```

Why it's useful:

- Easy to understand and follow
- Natural for solving problems that have clear, sequential steps
- Good for beginners because it matches how we naturally think about tasks

Example Languages:

- **C** - Used to create operating systems like Windows and Linux
- **COBOL** - Still used by banks and government systems for processing millions of transactions
- **FORTRAN** - Used by scientists for complex calculations
- **BASIC** - Often the first language people learn

Real Applications:

- Banking software that processes your ATM transactions step-by-step
- Scientific calculations that must be done in a specific order
- Embedded systems in your microwave or car

2. Object-Oriented Programming (OOP)

What it is: Object-oriented programming is like organizing the world into different objects that each have their own characteristics and abilities.

Key Characteristics:

- **Imperative paradigm** - Still tells the computer what to do
- **Objects everywhere** - Everything is represented as objects with properties and behaviors
- **Objects communicate** - Objects send messages to each other through their interfaces
- **The Four Pillars of OOP:**

The Four Pillars of OOP (Explained Simply):

1. **Abstraction** - Hiding complex details

- **Like:** Using a TV remote without knowing how the TV's electronics work
- **In Programming:** You can use a "Car" object without knowing exactly how the engine works internally

2. **Encapsulation** - Keeping related things together and protecting them

- **Like:** A medicine bottle that keeps pills safe and has a child-proof cap
- **In Programming:** A "BankAccount" object keeps your balance private and only lets you access it through specific methods like "deposit" or "withdraw"

3. **Inheritance** - Children inherit traits from parents

- **Like:** How a puppy inherits traits from its parent dogs
- **In Programming:** A "SportsCar" inherits basic car features but adds its own special abilities like "turboBoost"

4. **Polymorphism** - Same action, different results depending on who does it

- **Like:** Both birds and airplanes can "fly," but they do it differently
- **In Programming:** Both cats and dogs can "makeSound," but cats say "meow" and dogs say "woof"

Real-World Example: Think of a video game:

- **Player object:** Has health, name, position, can move(), attack(), jump()
- **Enemy object:** Has health, strength, position, can move(), attack(), takeDamage()
- **Weapon object:** Has damage, durability, can fire(), reload()

Each object knows what it can do and can interact with other objects.

Example Languages:

- **Java** - Used to build Android apps and large business applications
- **C++** - Used for video games, web browsers, and system software
- **Python** - Used for web development, data science, and artificial intelligence
- **PHP** - Powers most websites including Facebook and WordPress

Real Applications:

- Video games where characters, weapons, and environments are all objects
- Banking software where accounts, transactions, and customers are objects
- Social media platforms where users, posts, and comments are objects

3. Functional Programming

What it is: Functional programming is like using mathematical formulas and rules to solve problems.

Key Characteristics:

- **Declarative paradigm** - You tell the computer what you want, not exactly how to do it
- **Functions are the building blocks** - Everything is done through functions
- **Function compositions** - You combine simple functions to create complex ones
- **No changing data** - Data stays the same; you create new data instead of modifying old data

Real-World Example: Think of mathematical functions:

- $f(x) = x + 2$
- $g(x) = x * 3$
- You can combine them: $h(x) = g(f(x)) = (x + 2) * 3$

If $x = 5$:

- $f(5) = 7$
- $g(7) = 21$
- So $h(5) = 21$

Another Analogy: Like a factory assembly line where each station does one specific job and passes the result to the next station. Each station is a "function."

Why it's powerful:

- Very reliable because functions always give the same output for the same input
- Easy to test because each function is independent
- Great for parallel processing (multiple tasks at once)
- Fewer bugs because data doesn't change unexpectedly

Example Languages:

- **Haskell** - Used in finance for complex mathematical calculations
- **Lisp** - One of the oldest languages, used in artificial intelligence research
- **Clojure** - Used by companies like Netflix for handling millions of users
- **Erlang** - Powers WhatsApp and other messaging systems that need to handle millions of messages

Real Applications:

- Financial systems that calculate interest, loans, and risk
- Data analysis and machine learning
- Cryptocurrency and blockchain systems
- Large-scale web services that need to be very reliable

4. Dynamic/Scripting Languages

What it is: Dynamic/scripting languages are like having a very flexible, adaptable assistant that can change what it does while it's working.

Key Characteristics:

- **Not exactly a paradigm** - More about how the language works at runtime
- **Dynamic runtime environment** - The program can change and adapt while it's running
- **Flexible and adaptable** - Can modify itself based on conditions
- **Usually interpreted** - Code is read and executed line by line

Real-World Example: Imagine a smart assistant that can:

- Learn new tasks while working
- Change its approach based on the situation
- Adapt to new information without stopping
- Handle unexpected situations gracefully

Like a chameleon that changes color based on its environment!**Why they're popular:**

- **Fast development** - You can write programs quickly
- **Easy to modify** - Changes can be made without recompiling
- **Interactive** - You can test ideas immediately
- **Flexible** - Can handle many different types of data and situations

Example Languages:

- **Python** - Used for web development, data science, AI, automation
- **JavaScript** - Powers all web browsers and many mobile apps
- **Ruby** - Used by GitHub, Shopify, and many web applications
- **R** - Used by statisticians and data scientists
- **MATLAB** - Used by engineers and scientists for calculations and simulations

Real Applications:

- Web development (making websites interactive)
 - Data analysis and visualization
 - Automation scripts (like automatically organizing your files)
 - Rapid prototyping (quickly testing new ideas)
 - Scientific research and analysis
-

Imperative vs. Declarative - A Key Distinction

Imperative (How to do it)

Like giving step-by-step directions: "Go straight for 2 blocks, turn right at the red building, go 3 more blocks, turn left at the gas station..."

Declarative (What you want)

Like giving a destination: "Take me to the mall" (and let the GPS figure out the best route)

Programming Examples:

- **Imperative:** "Loop through this list, check each item, if it's greater than 5, add it to the new list"
 - **Declarative:** "Give me all items greater than 5" (and let the system figure out how)
-

Language Categories and Examples

Procedural Languages

Language	Famous For	Real-World Use
C	System programming	Operating systems, embedded systems
COBOL	Business applications	Banking, government systems
FORTRAN	Scientific computing	Weather prediction, engineering
BASIC	Education	Learning programming, simple applications

Object-Oriented Languages

Language	Famous For	Real-World Use
Java	"Write once, run anywhere"	Android apps, enterprise software
C++	Performance + objects	Video games, browsers, system software
Python	Simplicity + power	AI, web development, data science
C#	Microsoft ecosystem	Windows applications, web services

Functional Languages

Language	Famous For	Real-World Use
Haskell	Pure functional programming	Financial modeling, research
Lisp	AI and symbolic computation	Research, specialized applications
Erlang	Fault-tolerant systems	Telecommunications, messaging systems
F#	.NET functional programming	Financial services, data analysis

Dynamic/Scripting Languages

Language	Famous For	Real-World Use
Python	Readability and versatility	Everything! Web, AI, automation
JavaScript	Web interactivity	All websites, mobile apps, servers
Ruby	Developer happiness	Web applications, automation
R	Statistical computing	Data analysis, research

Key Programming Concepts Explained

Data Types (Like Different Containers)

Think of data types as different types of containers for different kinds of information:

1. Integer (int) - Whole numbers

- **Like:** A jar that can only hold whole marbles
- **Examples:** 5, -2, 0, 1000
- **Use:** Counting things, positions, IDs

2. Float/Double - Decimal numbers

- **Like:** A container that can hold liquid (can be partial)
- **Examples:** 3.14, -2.5, 0.001
- **Use:** Measurements, calculations, percentages

3. String - Text

- **Like:** A box that holds letters and words
- **Examples:** "Hello", "John Smith", "CSS125"
- **Use:** Names, messages, descriptions

4. Boolean - True or False

- **Like:** A light switch (on or off)
- **Examples:** true, false
- **Use:** Yes/no questions, flags, conditions

5. Array/List - Multiple items

- **Like:** A egg carton that holds multiple eggs
- **Examples:** [1, 2, 3, 4], ["apple", "banana", "orange"]
- **Use:** Shopping lists, collections, sequences

Variable Handling

Variables are like labeled boxes that store information.

Explicit vs. Implicit Types

Explicit Types - You must tell the computer what type of box you want

```
int age = 25;           // I'm explicitly saying this box holds integers
string name = "John"; // I'm explicitly saying this box holds text
```

Like: Labeling moving boxes "Kitchen Items" or "Books"

Implicit Types - The computer figures out what type of box you need

```
age = 25           # Computer knows this is an integer
name = "John"      # Computer knows this is text
```

Like: Smart storage that automatically knows what type of items you're putting in

Static vs. Dynamic Typing

Static Typing - Once you label a box, it can only hold that type of item

```
int age = 25;
age = "twenty-five"; // ERROR! This box only holds numbers
```

Like: A labeled medicine bottle that should only hold specific pills

Dynamic Typing - Boxes can change what type of items they hold

```
age = 25  
age = "twenty-five" # OK! The box can change what it holds
```

Like: A flexible container that can hold different items at different times

Control Structures

Branching Structures (Making Decisions)

Like a fork in the road - you choose which path to take based on conditions.

If-Then-Else:

```
IF it's raining THEN  
    bring an umbrella  
ELSE  
    wear sunglasses
```

Switch/Case:

```
DEPENDING ON the day of the week:  
    CASE Monday: go to work  
    CASE Saturday: sleep in  
    CASE Sunday: go to church
```

Iterative Structures (Repeating Actions)

Like doing the same task multiple times.

For Loop - Do something a specific number of times

```
FOR each day of the week:  
    set alarm clock
```

While Loop - Keep doing something until a condition is met

```
WHILE there are dirty dishes:  
    wash a dish
```

Subprograms (Functions/Methods)

Like having specialized helpers that do specific jobs.

Think of it like this:

- You have a "Calculate Tip" helper
- You give it: bill amount and tip percentage
- It gives back: the tip amount
- You can use this helper anytime you need to calculate tips

Benefits:

- **Reusable** - Write once, use many times
 - **Organized** - Each helper has one job
 - **Testable** - You can test each helper separately
 - **Maintainable** - Fix problems in one place
-

How Different Paradigms Solve the Same Problem

Problem: Calculate the total price of items in a shopping cart with tax.

Procedural Approach:

1. Set total = 0
2. For each item in cart:
 - a. Add item price to total
3. Calculate tax = total * tax_rate
4. Add tax to total
5. Return total

Object-Oriented Approach:

ShoppingCart object:

- Has: list of items, tax rate
- Can: addItem(), removeItem(), calculateTotal()

Item object:

- Has: name, price, category
- Can: getPrice(), getCategory()

Customer object:

- Has: cart, payment method
- Can: checkout(), pay()

Functional Approach:

```
total = calculateTotal(
    addTax(
        sumPrices(
            getItems(cart)
        ),
        TAX_RATE
    )
)
```

Dynamic/Scripting Approach:

```
# Python example - very flexible and readable
def calculate_total(cart, tax_rate=0.08):
    subtotal = sum(item['price'] for item in cart)
    return subtotal * (1 + tax_rate)
```

Choosing the Right Paradigm

When to Use Procedural:

- **Simple, linear tasks** - Like data processing scripts

- **System programming** - Operating systems, embedded systems
- **Scientific computing** - Mathematical calculations
- **Learning programming** - Easy to understand

When to Use Object-Oriented:

- **Complex applications** - Video games, business software
- **Team projects** - Easy to divide work among programmers
- **Reusable code** - When you need to use similar objects many times
- **User interfaces** - Buttons, windows, menus are natural objects

When to Use Functional:

- **Data processing** - Analyzing large datasets
- **Mathematical computations** - Financial calculations, scientific research
- **Concurrent programming** - When you need to do many things at once
- **Reliability** - When correctness is critical (banking, aerospace)

When to Use Dynamic/Scripting:

- **Rapid development** - Quick prototypes and experiments
 - **Web development** - Interactive websites and web applications
 - **Automation** - Scripts to automate repetitive tasks
 - **Data analysis** - Exploring and visualizing data
-

Common Misconceptions

"One paradigm is better than others"

Reality: Each paradigm is a tool. Like having different tools in a toolbox - you use the right tool for the job.

"You must choose only one paradigm"

Reality: Many modern languages support multiple paradigms. Python can be procedural, object-oriented, and functional!

"Functional programming is too hard"

Reality: While it can be different from what you're used to, it's just another way of thinking about problems.

"Object-oriented is always the best for large projects"

Reality: It depends on the project. Sometimes a functional or procedural approach is better.

Study Tips

To Remember Paradigms:

1. **Procedural** = Recipe (step-by-step)
2. **Object-Oriented** = Organized kitchen (objects with roles)
3. **Functional** = Math formulas (input → function → output)
4. **Dynamic** = Flexible assistant (adapts while working)

Practice Exercises:

1. **Identify the paradigm** - Look at code examples and identify which paradigm they use
2. **Compare solutions** - Solve the same problem using different paradigms
3. **Real-world mapping** - Think of real-world processes and how each paradigm would handle them

Key Questions to Ask:

- How does this paradigm organize code?
 - What are the main building blocks?
 - How does data flow through the program?
 - When would this approach be most useful?
-

Conclusion

Programming paradigms are different ways of thinking about and solving problems with computers. Just like there are different ways to organize your room, cook a meal, or plan a trip, there are different ways to write programs.

The key is understanding that each paradigm has its strengths and is suited for different types of problems. As you learn programming, you'll develop intuition for when to use each approach.

Remember: **The goal isn't to memorize everything, but to understand the different ways of thinking about problems and solutions.**

Lexical Analysis

1. Introduction to Lexical Analysis

Lexical Analysis is the first phase of compilation that converts raw source code into a stream of tokens. Think of it as the "word recognition" phase - just like how you recognize individual words when reading a sentence, the lexical analyzer recognizes individual meaningful units in source code.

Real-World Analogy

Consider reading this sentence: "The quick brown fox jumps over the lazy dog."

Your brain automatically:

- Recognizes individual words
- Ignores spaces between words
- Identifies punctuation
- Groups characters into meaningful units

Similarly, a lexical analyzer takes source code like `int count = 42;` and recognizes:

- `int` as a keyword
 - `count` as an identifier
 - `=` as an assignment operator
 - `42` as a number literal
 - `;` as a semicolon
-

2. The Role of the Lexical Analyzer

The lexical analyzer (also called a "lexer" or "scanner") has several critical responsibilities:

Primary Functions

2.1 Converting Raw Text into Tokens

What it does: Transforms the continuous stream of characters in source code into discrete, meaningful units called tokens.

Example:

```
int x = 5 + 3;
```

Becomes the token stream:

```
[KEYWORD: int] [IDENTIFIER: x] [ASSIGN_OP: =] [NUMBER: 5] [PLUS_OP: +] [NUMBER: 3] [SEMICOLON: ;]
```

2.2 Stripping Out Whitespace and Comments

What it does: Removes unnecessary characters that don't affect program semantics.

Example:

```
int    x    =    5;    // This assigns 5 to x
/* Multi-line comment
   explaining the code */
int y = 10;
```

Becomes:

```
[KEYWORD: int] [IDENTIFIER: x] [ASSIGN_OP: =] [NUMBER: 5] [SEMICOLON: ;]
[KEYWORD: int] [IDENTIFIER: y] [ASSIGN_OP: =] [NUMBER: 10] [SEMICOLON: ;]
```

2.3 Error Message Correlation

What it does: Tracks line numbers and column positions to help correlate compiler error messages with the original source code.

Example: If there's a syntax error, instead of just saying "syntax error," the compiler can say:

```
Error on line 15, column 8: Expected ';' after expression
```

Interaction with Other Compiler Phases

The lexical analyzer works as a "token supplier" for the parser:

```
Source Code → [Lexical Analyzer] → Token Stream → [Parser] → Parse Tree
                ↑
                ↓
            Symbol Table
```

Real-World Example: Think of an assembly line where:

- Raw materials (source code) come in
- First station (lexical analyzer) sorts materials into standard parts (tokens)
- Second station (parser) assembles parts into subcomponents (expressions, statements)
- Later stations build the final product (executable code)

3. Why Separate Lexical Analysis from Parsing?

Separating lexical analysis from parsing might seem like extra work, but it provides crucial benefits:

3.1 Simplicity in Design

Benefit: Each phase handles one specific task well.

Real-World Analogy: Consider building a house:

- You don't simultaneously dig the foundation, frame walls, and install plumbing
- Each phase has specialists who focus on their expertise
- Similarly, lexical analysis focuses on character-level patterns, while parsing focuses on structural relationships

Technical Example: Without separation, a parser would need to handle:

```
while ( x < 10 ) {
```

As individual characters: w, h, i, l, e, , (, , x, ...

With separation, the parser receives clean tokens: [WHILE] [LPAREN] [ID: x] [LESS] [NUM: 10] [RPAREN] [LBRACE]

3.2 Improvement in Efficiency

Benefit: Specialized algorithms for each task perform better than general-purpose solutions.

Character Processing: Lexical analysis uses efficient character-scanning techniques (like finite automata) optimized for linear text processing.

Structure Processing: Parsing uses context-free grammar techniques optimized for handling nested structures.

Performance Example:

- Lexical analysis: $O(n)$ where n = number of characters
- If combined with parsing: $O(n^2)$ or worse due to backtracking

3.3 Enhanced Compiler Portability

Benefit: Only the lexical analyzer needs to change when porting to different character encodings or platforms.

Practical Example:

- Moving from ASCII to Unicode: Only lexical analyzer changes
 - Different comment styles (`//` vs `#`): Only lexical analyzer changes
 - Different platforms (Windows vs Unix line endings): Only lexical analyzer changes
-

4. Tokens, Patterns, and Lexemes

These three concepts are fundamental to understanding lexical analysis:

4.1 Definitions and Relationships

Token

Definition: A category or class of lexical units. The "type" of a piece of text.

Think of it as: The grammatical category (like "noun," "verb," "adjective" in English)

Pattern

Definition: The rule or regular expression that describes what text can belong to a token.

Think of it as: The spelling rules that define the category

Lexeme

Definition: The actual sequence of characters in the source code that matches a pattern.

Think of it as: The specific word instance

4.2 Detailed Examples

Example 1: Identifiers

```
Token: IDENTIFIER  
Pattern: [A-Za-z_][A-Za-z0-9_]*  
Lexemes: "count", "userName", "_temp", "MAX_SIZE"
```

Real-World Context:

```
int userAge = 25;           // "userAge" is a lexeme of token IDENTIFIER  
char firstName[50];        // "firstName" is a lexeme of token IDENTIFIER  
int _private = 0;          // "_private" is a lexeme of token IDENTIFIER
```


Example 2: Number Literals

```
Token: INTEGER_LITERAL
Pattern: [0-9]+
Lexemes: "42", "0", "999", "2024"
```

```
Token: FLOAT_LITERAL
Pattern: [0-9]+\.[0-9]+
Lexemes: "3.14", "0.5", "99.99"
```

Real-World Context:

```
int year = 2024;           // "2024" is a lexeme of INTEGER_LITERAL
double pi = 3.14159;       // "3.14159" is a lexeme of FLOAT_LITERAL
float price = 19.99;       // "19.99" is a lexeme of FLOAT_LITERAL
```

Example 3: Keywords vs Identifiers

This demonstrates why patterns must be carefully designed:

```
int if_count = 5;          // "if_count" is IDENTIFIER (not keyword "if")
if (x > 0) return;         // "if" is KEYWORD
```

Pattern Priority: Keywords typically have higher priority than general identifier patterns.

4.3 Handling Multiple Pattern Matches with Attributes

When a lexeme could match multiple patterns, we use **attributes** to provide additional information:

Token-Attribute Pairs

Instead of just tokens, the lexical analyzer produces token-attribute pairs:

```
<TOKEN_TYPE, attribute_value>
```

Example from the PDF: For the Fortran statement: `E = M * C ** 2`

The token stream becomes:

```
<id, pointer to symbol-table entry for E>
<assign_op, >
<id, pointer to symbol-table entry for M>
<mult_op, >
<id, pointer to symbol-table entry for C>
<exp_op, >
<num, integer value 2>
```

Real-World Application: Symbol Table Integration

```
int count = 42;
int total = count + 10;
```

Token-Attribute Stream:

```
<KEYWORD, "int">
<IDENTIFIER, ptr_to_symbol_table["count"]>
<ASSIGN, >
<INTEGER, 42>
<SEMICOLON, >
<KEYWORD, "int">
<IDENTIFIER, ptr_to_symbol_table["total"]>
<ASSIGN, >
<IDENTIFIER, ptr_to_symbol_table["count"]> // Same pointer as before
<PLUS, >
<INTEGER, 10>
<SEMICOLON, >
```

4.4 Reserved Strings (Keywords)

Definition: Lexemes that have specific meaning in the programming language and cannot be used as identifiers.

Examples in C:

- `if, else, while, for, int, char, return, void`

Examples in Python:

- `def, class, import, from, if, elif, else`

Implementation Strategy:

1. First check if a word matches a keyword pattern
2. If not a keyword, then check if it matches identifier pattern
3. Keywords have higher precedence than identifiers

5. Lexical Errors and Error Handling

5.1 Nature of Lexical Errors

Key Insight: Few errors are actually detectable at the lexical level alone.

Why? The lexical analyzer only sees character patterns, not program logic.

5.2 Examples of Lexical vs Non-Lexical Errors

Lexical Errors (Detected by Lexer)

```
int x = 3.14.159;           // Invalid number format
char c = 'unterminated string
int count = 999999999999999999; // Number too large
```

Non-Lexical Errors (Not Detected by Lexer)

```
fi(a == f(x))      // "fi" is valid identifier, just misspelled "if"
int x = y + z;      // All valid tokens, but 'y' and 'z' might not be declared
return x + ;        // Valid tokens, but invalid syntax
```

5.3 Error Recovery Strategies

When the lexical analyzer encounters invalid character sequences, it needs to recover and continue processing:

5.3.1 Panic Mode Recovery

Strategy: Delete characters until a well-formed token can be created.

Example:

```
int x@ = 5;
```

Recovery: Delete '@' and continue, treating it as `int x = 5;`

5.3.2 Character-Level Corrections

The lexer can attempt these corrections:

Deletion: Remove an invalid character

```
in@t x = 5; → int x = 5;    // Delete '@'
```

Insertion: Add a missing character

```
int x == 5; → int x = 5;    // Assume one '=' was intended
```

Replacement: Replace a wrong character

```
int x -= 5; → int x = 5;    // Replace '-' with nothing (or report as compound operator)
```

Transposition: Swap adjacent characters

```
int x = 5;; → int x = 5;    // Remove duplicate semicolon
```

5.4 Real-World Error Handling Example

Consider this malformed code:

```
int main() {
    int x = 3.14.159;    // Lexical error: invalid number
    int y = x +;         // Syntax error: incomplete expression
    retur x;             // Lexical: "retur" is valid identifier, not "return"
}
```

Lexer Response:

1. Line 2: Error - invalid floating point number, suggest `3.14159`
 2. Line 3: Produces tokens `int, y, =, x, +, ;` - parser will catch syntax error
 3. Line 4: Produces `retur` as identifier - semantic analyzer will catch undefined identifier
-

6. Token Specifications and Regular Expressions

6.1 Formal Foundations

Alphabet (Σ)

Definition: A finite set of symbols used to build strings.

Examples:

- **ASCII:** 128 characters including letters, digits, punctuation
- **Unicode:** Millions of characters from world languages
- **Binary:** $\{0, 1\}$
- **Simple:** $\{a, b, c\}$

String

Definition: A finite sequence of symbols from an alphabet.

Examples over alphabet $\{a, b\}$:

- ϵ (empty string)
- `a`
- `ab`
- `abba`
- `aaabbbaaaa`

Length of String $|S|$

Definition: Number of symbols in string S .

Examples:

- $|\epsilon| = 0$
- $|\text{"hello"}| = 5$
- $|\text{"C++"}| = 3$

Language

Definition: A set of strings over an alphabet.

Examples:

```
L1 = { $\epsilon$ , a, aa, aaa, ...}           // All strings of a's (including empty)
L2 = {ab, aabb, aaabbb, ...}       // Strings with equal a's followed by equal b's
L3 = {"if", "else", "while", ...}  // C keywords
```

6.2 String Set Operations

Concatenation (L_1L_2)

Definition: Set of all strings xy where $x \in L_1$ and $y \in L_2$

Example:

```
L1 = {a, ab}
L2 = {c, cd}
L1L2 = {ac, acd, abc, abcd}
```

Programming Example:

```
Keywords = {"int", "char"}
Identifiers = {"x", "count"}
// Concatenation creates sequences like "intx", "charcount"
// (Not useful for programming languages, but demonstrates the concept)
```

Union ($L_1 \cup L_2$)

Definition: Set of all strings x where $x \in L_1$ OR $x \in L_2$

Example:

```
Keywords = {"int", "if", "else"}
Operators = {"+", "-", "*"}
Tokens = Keywords  $\cup$  Operators = {"int", "if", "else", "+", "-", "*"}
```

String Power (L^n)

Definition: String concatenated with itself n times.

Examples:

```
"ab"1 = "ab"
"ab"2 = "abab"
"ab"3 = "ababab"
```

Kleene Closure (Σ^*)

Definition: Set of all strings comprised of zero to many symbols from Σ .

Example with $\Sigma = \{a, b\}$:

```
 $\Sigma^* = \{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, abb, baa, bab, bba, bbb, \dots\}$ 
```

Programming Application:

```
Digits = {0, 1, 2, 3, 4, 5, 6, 7, 8, 9}
Digits* = All possible number strings:  $\{\epsilon, 0, 1, \dots, 42, 123, 999, \dots\}$ 
```

Positive Closure (Σ^+)

Definition: Set of all strings comprised of one to many symbols from Σ .

Difference from Kleene Closure: Excludes the empty string ϵ .

```
Digits+ = {0, 1, 2, ..., 42, 123, 999, ...} // No empty string
```

6.3 Regular Expression Definitions

Regular expressions provide a formal way to define token patterns:

Basic Regular Expressions

ϵ (Empty String)

Regular Expression: ϵ
Language: $\{\epsilon\}$
Example Use: Optional elements

Single Symbol

Regular Expression: a (where $a \in \Sigma$)
Language: $\{a\}$
Example: The letter 'a' matches only the string "a"

Concatenation (rs)

Regular Expression: rs
Language: All strings from concatenating strings matching r with strings matching s
Example: ab matches only "ab"

Union ($r|s$)

Regular Expression: $r|s$
Language: All strings matching either r or s
Example: $a|b$ matches "a" or "b"

Kleene Star (r^*)

Regular Expression: r^*
Language: Zero or more concatenations of strings matching r
Example: a^* matches "", "a", "aa", "aaa", ...

6.4 Practical Token Patterns

Identifiers

Pattern: `[A-Za-z_][A-Za-z0-9_]*`

Meaning: Letter or underscore, followed by zero or more letters, digits, or underscores

Examples: `userName`, `_temp`, `MAX_SIZE`, `x`, `y1`

Integer Literals

Pattern: `[0-9]+`

Meaning: One or more digits

Examples: `0`, `42`, `999`, `2024`

Floating Point Literals

Pattern: `[0-9]+\.[0-9]+`

Meaning: Digits, decimal point, more digits

Examples: `3.14`, `0.5`, `99.99`

String Literals

Pattern: `\("[^"]*"`

Meaning: Quote, zero or more non-quote characters, quote

Examples: `"hello"`, `"C++ Programming"`, `""`

Comments (C-style)

Single-line: `//[^\\n]*`

Multi-line: `/*.*?*/`

Whitespace

Pattern: `[\\t\\n\\r]+`

Meaning: One or more spaces, tabs, newlines, or carriage returns

7. String Set Operations

7.1 Concatenation in Detail

Formal Definition: For languages L_1 and L_2 :

$$L_1 L_2 = \{xy \mid x \in L_1 \wedge y \in L_2\}$$

Programming Language Example:

```
PrefixOperators = {"++", "--"}
Variables = {"x", "y", "count"}
PreIncrements = PrefixOperators · Variables
Result = {"++x", "++y", "++count", "--x", "--y", "--count"}
```

Real Code:

```
++count;      // Matches "++count" from concatenation
--x;          // Matches "--x" from concatenation
```

7.2 Union in Language Design

Practical Application: Defining all valid tokens

```
Numbers = {[0-9]+}
Identifiers = {[A-Za-z][A-Za-z0-9]*}
Keywords = {"if", "else", "while", "for"}
Operators = {"+", "-", "x", "/", "="}

AllTokens = Numbers U Identifiers U Keywords U Operators
```

7.3 Kleene and Positive Closure Applications

Kleene Closure (*) - Zero or More

Use Case: Optional repeated elements

```
int x;           // Zero dimensions: int x
int array[10];   // One dimension: int array[10]
int matrix[10][20]; // Two dimensions: int matrix[10][20]
```

Pattern: `[number]` can appear zero or more times **Regular Expression:** `int identifier (\[.*?\])*`

Positive Closure (+) - One or More

Use Case: Required repeated elements

```
// Valid: one or more digits
123, 42, 7, 999999

// Invalid: empty string not allowed for numbers
// Pattern: [0-9]+ ensures at least one digit
```

8. Regular Expression Operator Precedence

Understanding operator precedence prevents ambiguity in regular expressions:

8.1 Precedence Rules (Highest to Lowest)

1. Kleene Star (*) - Highest Precedence, Left-Associative

```
Expression: ab*
Interpretation: a(b*) // NOT (ab)*
Matches: "a", "ab", "abb", "abbb", ...
Does NOT match: "", "ab", "abab", "ababab", ...
```

2. Concatenation - Middle Precedence, Left-Associative

Expression: `abc`
 Interpretation: `(ab)c = a(bc)` // Left-associative doesn't matter for concatenation
 Matches: Only `"abc"`

3. Union (|) - Lowest Precedence, Left-Associative

Expression: `a|bc`
 Interpretation: `a|(bc)` // NOT `(a|b)c`
 Matches: `"a"` or `"bc"`
 Does NOT match: `"ac"` or `"bb"`

8.2 Complex Examples with Precedence

Example 1: Identifier or Keyword

Expression: `if|[A-Za-z][A-Za-z0-9]*`
 Interpretation: `(if)|([A-Za-z][A-Za-z0-9]*)`
 Matches: `"if"` or any identifier starting with letter

Example 2: Optional Exponent in Numbers

Expression: `[0-9]+\.[0-9]+E[+-]?[0-9]+|[0-9]+`
 Interpretation: `(([0-9]+\.[0-9]+E[+-]?[0-9]+)|([0-9]+))`
 Matches: `"3.14E+2"`, `"1.5E-3"`, `"42"`

Example 3: C-style Comments

Expression: `//[^\n]*|/*.*?*/`
 Interpretation: `(//[^\n]*)|(/*.*?*/)`
 Matches: Single-line OR multi-line comments

8.3 Using Parentheses to Override Precedence

Without Parentheses (Following Precedence)

```
ab|cd*   →   a(b)|(c(d*))   →   "ab" or "c", "cd", "cdd", ...
```

With Parentheses (Explicit Grouping)

```
(ab|cd)* → (ab|cd)* → "", "ab", "cd", "abab", "cdcd", "abcd", ...
```

8.4 Real-World Pattern Examples

Email Address Pattern (Simplified)

```
[A-Za-z0-9]+@[A-Za-z0-9]+\.[A-Za-z]+
```

Precedence Analysis:

1. * applied to character classes: [A-Za-z0-9]+, [A-Za-z0-9]+, [A-Za-z]+
2. Concatenation: combines all parts in sequence
3. Result: user@domain.com format

URL Pattern (Simplified)

```
https?://[A-Za-z0-9.-]+/[A-Za-z0-9/*]*
```

Precedence Analysis:

1. ? applied to 's': http or https
2. + applied to domain characters
3. * applied to path characters
4. Concatenation: combines protocol + :// + domain + / + path

9. Real-World Applications

9.1 Programming Language Compilers

Every major programming language uses lexical analysis:

C/C++ Compiler

```
#include <stdio.h>
int main() {
    int count = 42;
    printf("Count: %d\n", count);
    return 0;
}
```

Lexical Analysis Output:

```
[PREPROCESSOR: #include] [ANGLE_BRACKET: <] [IDENTIFIER: stdio.h] [ANGLE_BRACKET: >]
[KEYWORD: int] [IDENTIFIER: main] [LPAREN: (] [RPAREN: )] [LBRACE: {]
[KEYWORD: int] [IDENTIFIER: count] [ASSIGN: =] [INTEGER: 42] [SEMICOLON: ;]
...
```

Python Interpreter

```
def calculate_area(radius):
    pi = 3.14159
    return pi * radius ** 2
```

Token Stream:

```
[KEYWORD: def] [IDENTIFIER: calculate_area] [LPAREN: (] [IDENTIFIER: radius] [RPAREN: :]
[IDENTIFIER: pi] [ASSIGN: =] [FLOAT: 3.14159]
[KEYWORD: return] [IDENTIFIER: pi] [MULTIPLY: *] [IDENTIFIER: radius] [POWER: **] [INTEGER: 2]
```

9.2 Text Editors and IDEs

Syntax Highlighting

Modern text editors use lexical analysis for syntax highlighting:

```
// Different colors for different token types:
int count = 42; // Blue for keywords, black for identifiers, red for numbers
```

Implementation:

1. Lexical analyzer identifies token types
2. Editor applies color schemes based on token types
3. Real-time analysis as you type

Code Completion

IDEs use lexical analysis to provide intelligent suggestions:

```
str|           // Typing 'str' triggers completion
  ↓
string         // Suggests 'string' keyword
strlen         // Suggests 'strlen' function
strcpy         // Suggests 'strcpy' function
```

9.3 Search Engines and Text Processing

Search Query Processing

```
Query: "machine learning" AND python -java
```

Lexical Analysis:

```
[QUOTED_PHRASE: "machine learning"] [AND_OPERATOR: AND] [TERM: python] [EXCLUDE_OPERATOR: -]
[TERM: java]
```

Log File Analysis

```
2024-01-15 10:30:42 ERROR: Connection failed to database server 192.168.1.100
```

Token Extraction:

```
[DATE: 2024-01-15] [TIME: 10:30:42] [LEVEL: ERROR] [MESSAGE: Connection failed...] [IP: 192.168.1.100]
```

9.4 Data Validation and Processing

Email Validation

```
Pattern: [A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Za-z]{2,}
Valid: john.doe@example.com, user123@test.org
Invalid: @example.com, user@, user@.com
```

Phone Number Parsing

```
Input: "(555) 123-4567", "555-123-4567", "5551234567"
Pattern: \(?([0-9]{3})\)?[-. ]?([0-9]{3})[-. ]?([0-9]{4})
Output: Groups (555, 123, 4567) for all formats
```

9.5 Configuration File Processing

JSON Parser

```
{
  "name": "John Doe",
  "age": 30,
  "active": true
}
```


Lexical Analysis:

```
[LBRACE: {] [STRING: "name"] [COLON: :] [STRING: "John Doe"] [COMMA: ,]
[STRING: "age"] [COLON: :] [NUMBER: 30] [COMMA: ,]
[STRING: "active"] [COLON: :] [BOOLEAN: true] [RBRACE: }]
```

Configuration Files

```
[database]
host=localhost
port=5432
username=admin
```

Token Stream:

```
[SECTION: database] [KEY: host] [ASSIGN: =] [VALUE: localhost]
[KEY: port] [ASSIGN: =] [VALUE: 5432] [KEY: username] [ASSIGN: =] [VALUE: admin]
```

9.6 Web Development

CSS Parsing

```
.button {
  background-color: #ff0000;
  padding: 10px;
}
```

Lexical Analysis:

```
[SELECTOR: .button] [LBRACE: {]
[PROPERTY: background-color] [COLON: :] [COLOR: #ff0000] [SEMICOLON: ;]
[PROPERTY: padding] [COLON: :] [SIZE: 10px] [SEMICOLON: ;]
[RBRACE: }]
```

HTML Template Processing

```
<div class="{{cssClass}}" data-value="{{value}}">{{content}}</div>
```

Template Token Extraction:

```
[HTML_TAG: <div] [ATTRIBUTE: class] [TEMPLATE_VAR: {{cssClass}}]
[ATTRIBUTE: data-value] [TEMPLATE_VAR: {{value}}] [HTML_TAG: >]
[TEMPLATE_VAR: {{content}}] [HTML_TAG: </div>]
```

9.7 Database Query Processing

SQL Parser

```
SELECT name, age FROM users WHERE age > 21 ORDER BY name;
```

Lexical Analysis:

```
[KEYWORD: SELECT] [IDENTIFIER: name] [COMMA: ,] [IDENTIFIER: age]
[KEYWORD: FROM] [IDENTIFIER: users] [KEYWORD: WHERE] [IDENTIFIER: age]
[OPERATOR: >] [NUMBER: 21] [KEYWORD: ORDER] [KEYWORD: BY] [IDENTIFIER: name] [SEMICOLON: ;]
```

Summary and Key Takeaways

Essential Concepts to Remember

1. **Lexical Analysis is Pattern Recognition:** It converts character sequences into meaningful tokens using pattern matching.
2. **Separation of Concerns:** Lexical analysis handles character-level patterns while parsing handles structural patterns.
3. **Token-Pattern-Lexeme Relationship:** Tokens are categories, patterns are rules, lexemes are actual text instances.
4. **Regular Expressions are Fundamental:** They provide formal notation for specifying token patterns.
5. **Error Recovery is Important:** Lexical analyzers must handle invalid input gracefully and continue processing.
6. **Real-World Applications are Everywhere:** From compilers to search engines to text editors, lexical analysis is foundational to modern computing.

Study Tips

1. **Practice Pattern Writing:** Create regular expressions for common programming constructs.
2. **Trace Through Examples:** Manually tokenize code snippets to understand the process.
3. **Understand Precedence:** Master operator precedence in regular expressions to avoid ambiguity.
4. **Connect to Real Tools:** Use tools like grep, sed, or regex testers to experiment with patterns.
5. **Think About Edge Cases:** Consider how lexical analyzers handle errors, unusual input, and boundary conditions.

The lexical analysis phase, while conceptually straightforward, requires careful attention to detail and thorough understanding of pattern matching principles. Master these concepts, and you'll have a solid foundation for understanding how programming languages and text processing tools work under the hood.

Lexical Analysis Implementation

What is Lexical Analysis?

Lexical analysis is the first phase of compilation where the source code is broken down into meaningful units called **tokens**. Think of it like reading a sentence and identifying individual words, punctuation, and their meanings.

Real-World Analogy

Imagine you're reading this sentence: `int x = 42;`

- A human recognizes: "int" (keyword), "x" (identifier), "=" (operator), "42" (number), ";" (delimiter)
 - A lexical analyzer does the same thing for programming languages
-

Recognition of Tokens

Basic Assumptions

The lexical analyzer operates under two fundamental assumptions:

1. Reserved words cannot be used as identifiers

- Keywords like `int`, `if`, `while` have special meanings
- You can't name a variable `int` in most languages

2. Whitespace is ignored

- Spaces, tabs, and newlines are typically discarded
- They serve only to separate tokens, not as meaningful tokens themselves

Token Recognition Process

The lexical analyzer follows a systematic 5-step process:

1. Order all token patterns in terms of precedence

- Keywords come before identifiers (so `if` isn't mistaken for an identifier)
- Longer operators before shorter ones (so `<=` isn't read as `<` and `=`)

2. Maintain lexeme-beginning pointer and forward pointer

- **Lexeme-beginning pointer:** Marks the start of the current token being analyzed
- **Forward pointer:** Moves through the input to build the token

3. Forward pointer moves through a transition diagram

- Each character is processed according to predefined patterns
- The pointer advances as it matches expected patterns

4. If valid lexeme is detected, return token

- When a complete, valid token is recognized, it's returned to the parser
- Example: Reading `123` completely forms a NUMBER token

5. If invalid, return forward pointer to initial pointer and try next pattern

- If the current pattern fails, backtrack and try the next pattern
- This prevents getting stuck on partial matches

Detailed Example of Token Recognition

Consider the input: `count = 42;`

```
Input: c o u n t   =   4 2 ;
      ↑               (lexeme-beginning pointer)
      ↑               (forward pointer)
```

Step 1: Try keyword patterns first

- "c" → doesn't match any keyword, continue
- "co" → doesn't match any keyword, continue
- "cou" → doesn't match any keyword, continue
- "coun" → doesn't match any keyword, continue
- "count" → doesn't match any keyword, try identifier pattern

Step 2: Try identifier pattern

- "count" → matches identifier pattern (letter followed by letters/digits)
- Return IDENTIFIER token with value "count"

Step 3: Skip whitespace, move to "="

- "=" → matches assignment operator pattern
- Return ASSIGN_OP token

Step 4: Skip whitespace, move to "42"

- "4" → start number pattern
- "42" → complete number
- Return NUMBER token with value 42

Step 5: Process ";"

- ";" → matches semicolon pattern
- Return SEMICOLON token

Transition Diagrams

What are Transition Diagrams?

Transition diagrams are visual representations of how a lexical analyzer moves between different states while processing input characters. They're like flowcharts for token recognition.

Elements of a Transition Diagram

1. **States** (represented as circles)

- Each state represents a stage in token recognition
- **Start state**: Where processing begins (marked with an arrow)
- **Accept states**: Final states where tokens are recognized (double circles)
- **Intermediate states**: Processing states between start and accept

2. **Edges or Transitions** (arrows between states)

- Show how to move from one state to another
- Labeled with conditions for the transition

3. **Symbols** (labels on edges)

- Specific characters or character classes that trigger transitions
- Examples: 'a', digit, letter, other

4. **Determinism**

- Each state has exactly one transition for each input symbol
- No ambiguity about which path to take

Example Transition Diagram Analysis

From the PDF's diagram with states A, B, C:



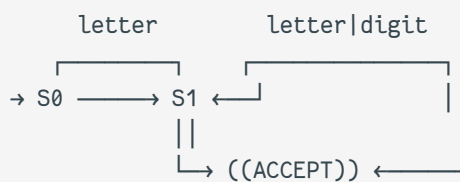
What this diagram recognizes:

- Strings that end with 'a' (state C is accepting)
- 'b' characters keep you in intermediate states
- Only 'a' can lead to acceptance

Trace Examples:

- Input "ba": $A \rightarrow A \rightarrow B$ (not accepted, no final 'a')
- Input "bba": $A \rightarrow A \rightarrow A \rightarrow B$ (accepted, ends in 'a')
- Input "abba": $A \rightarrow B \rightarrow A \rightarrow A \rightarrow B$ (accepted)

Real-World Application: Identifier Recognition

**How it works:**

1. Start in S0
2. First character must be a letter → go to S1
3. Subsequent characters can be letters or digits → stay in S1
4. When no more valid characters, accept as IDENTIFIER

Examples:

- `variable` → IDENTIFIER
- `count123` → IDENTIFIER
- `123abc` → fails (starts with digit)

Implementation Details

Code Structure

Lexical analyzers are typically implemented with these key components:

1. State-Based Processing

- Each state has a code segment
 - Different code executes depending on current state
 - State-specific logic for character processing

```
switch(currentState) {  
    case START:  
        // Handle initial character  
        break;  
    case IN_NUMBER:  
        // Process digits  
        break;  
    case IN_IDENTIFIER:  
        // Process letters/digits  
        break;  
}
```

2. State Transitions

- Each edge may change the state
 - Transitions based on input characters
 - Updates to current processing state

3. Character Processing Function

- Use `nextchar()` function to:
 - Read next character from input buffer
 - Advance the forward pointer
 - Return the character for processing

```
function nextchar() {  
    char c = buffer[forwardPointer];  
    forwardPointer++;  
    return c;  
}
```

Pointer Management

Two Critical Pointers:

1. **Lexeme-beginning pointer:** Marks start of current token
2. **Forward pointer:** Scans ahead to build tokens

Error Handling Rules:

Case 1: Invalid token encountered

- **Action:** Return forward pointer to lexeme-beginning pointer
- **Next step:** Try next transition diagram/pattern
- **Purpose:** Backtrack when current pattern fails

Case 2: Valid token completed

- **Action:** Retract forward pointer by one character
- **Next step:** Save current lexeme as token
- **Purpose:** Don't consume character that belongs to next token

Detailed Implementation Example

```
function getNextToken() {
    skipWhitespace();
    lexemeStart = forwardPointer;

    // Try each token pattern in order of precedence
    for each pattern in tokenPatterns {
        resetPointer(forwardPointer, lexemeStart);

        if (tryPattern(pattern)) {
            token = createToken(pattern, getLexeme());
            return token;
        }
    }

    // No pattern matched
    error("Invalid token at position " + lexemeStart);
}

function tryPattern(pattern) {
    state = pattern.startState;

    while (hasMoreInput()) {
        char c = nextchar();

        if (hasTransition(state, c)) {
            state = getNextState(state, c);
        } else {
            if (isAcceptState(state)) {
                retractPointer(); // Don't consume last character
                return true;      // Valid token found
            } else {
                return false;     // Invalid token
            }
        }
    }

    return isAcceptState(state);
}
```

Identifiers and Symbol Table Integration

Identifier Processing

Key Requirements:

1. Identifiers must be installed into the symbol table

- Symbol table stores all identifiers and their attributes
- Used later by semantic analyzer and code generator

2. Keyword optimization strategy

- **Store all keywords in a special collection** (hash table/array)
- **Benefits:** Reduces states needed in transition diagram
- **Process:** First recognize as identifier, then check if it's a keyword

Identifier vs Keyword Resolution

```
function processIdentifier(lexeme) {  
    if (isKeyword(lexeme)) {  
        return createKeywordToken(lexeme);  
    } else {  
        symbolTableEntry = installIdentifier(lexeme);  
        return createIdentifierToken(symbolTableEntry);  
    }  
}  
  
function isKeyword(lexeme) {  
    return keywordSet.contains(lexeme);  
}
```

Symbol Table Operations

```
function installIdentifier(name) {  
    if (symbolTable.contains(name)) {  
        return symbolTable.lookup(name);  
    } else {  
        entry = new SymbolTableEntry(name);  
        symbolTable.insert(name, entry);  
        return entry;  
    }  
}
```

Finite State Automata (FSA)

Definition and Components

A **Finite State Automaton** is a mathematical model used for input processing, consisting of:

1. **Finite set of states** (Q)
 - Limited number of distinct processing states
2. **Input alphabet** (Σ)
 - Set of valid input symbols
3. **Transition function** (δ)
 - Rules for moving between states: $\delta(\text{state}, \text{input}) \rightarrow \text{next_state}$
4. **Start state** (q_0)
 - Initial state where processing begins
5. **Set of accepting states** (F)
 - States where input is accepted as valid

Mathematical Notation

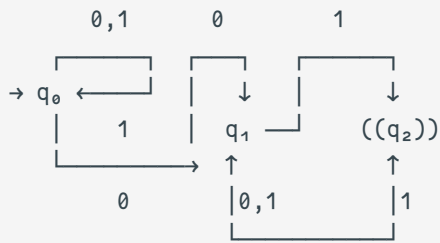
FSA Definition: $M = (Q, \Sigma, \delta, q_0, F)$

Where:

- $Q = \{q_0, q_1, q_2, \dots\}$ (states)
- $\Sigma = \{'a', 'b', '0', '1', \dots\}$ (alphabet)
- $\delta: Q \times \Sigma \rightarrow Q$ (transition function)
- $q_0 \in Q$ (start state)
- $F \subseteq Q$ (accepting states)

Practical FSA Examples

FSA 1: Strings ending with "01"



States:

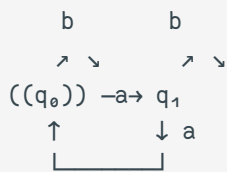
- q_0 : Start state (no specific pattern matched)
- q_1 : Just read '0' (potential start of "01")
- q_2 : Just read "01" (accepting state)

Transitions:

- q_0 on '0' $\rightarrow q_1$ (start of potential "01")
- q_0 on '1' $\rightarrow q_0$ (restart)
- q_1 on '0' $\rightarrow q_1$ (still potential start of "01")
- q_1 on '1' $\rightarrow q_2$ (complete "01" pattern)
- q_2 on '0' $\rightarrow q_1$ (new potential "01")
- q_2 on '1' $\rightarrow q_0$ (restart completely)

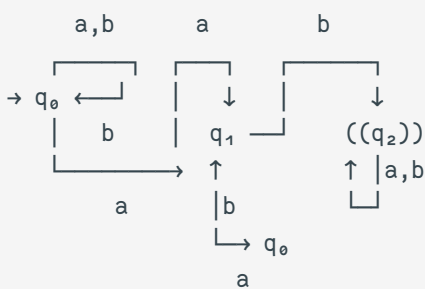
Example Traces:

- "101": $q_0 \rightarrow q_0 \rightarrow q_1 \rightarrow q_2$ ✓ (accepted)
- "1001": $q_0 \rightarrow q_0 \rightarrow q_1 \rightarrow q_1 \rightarrow q_2$ ✓ (accepted)
- "010": $q_0 \rightarrow q_1 \rightarrow q_2 \rightarrow q_1$ ✗ (not accepted)

FSA 2: Even number of 'a's**States:**

- q_0 : Even number of 'a's seen (accepting)
- q_1 : Odd number of 'a's seen (not accepting)

Key insight: Each 'a' toggles between even/odd states

FSA 3: Contains "ab" as substring**States:**

- q_0 : No relevant progress toward "ab"
- q_1 : Just saw 'a' (potential start of "ab")
- q_2 : Saw complete "ab" (accepting, stays accepting)

Practical Examples and Applications

Real Compiler Application

Consider analyzing this C code snippet:

```
int main() {
    int count = 0;
    return count + 1;
}
```

Lexical Analysis Output:

Token Type	Lexeme	Line	Column
KEYWORD	int	1	1
IDENTIFIER	main	1	5
L_PAREN	(1	9
R_PAREN)	1	10
L_BRACE	{	1	12
KEYWORD	int	2	5
IDENTIFIER	count	2	9
ASSIGN_OP	=	2	15
NUMBER	0	2	17
SEMICOLON	;	2	18
KEYWORD	return	3	5
IDENTIFIER	count	3	12
PLUS_OP	+	3	18
NUMBER	1	3	20
SEMICOLON	;	3	21
R_BRACE	}	4	1

Modern Applications

1. Integrated Development Environments (IDEs)

- **Syntax highlighting:** Lexical analysis identifies keywords, strings, comments
- **Auto-completion:** Recognizes partial identifiers
- **Error detection:** Highlights invalid tokens in real-time

2. Code Formatters and Linters

- **Prettier/ESLint:** Use lexical analysis to understand code structure
- **Style enforcement:** Identify and fix spacing, naming conventions

3. Programming Language Tools

- **Interpreters:** Python, JavaScript interpreters use lexical analysis
- **Compilers:** GCC, Clang start with lexical analysis
- **Transpilers:** TypeScript→JavaScript, CoffeeScript→JavaScript

Performance Considerations

Why FSAs are Efficient:

- **Linear time complexity:** $O(n)$ where n is input length
- **Constant memory:** Fixed number of states regardless of input size
- **Predictable behavior:** Deterministic transitions

Real-world Performance:

- Modern compilers can tokenize millions of lines per second
 - Lexical analysis typically represents <5% of total compilation time
-

Key Concepts Summary

Essential Terms

Token: Meaningful unit of source code (keyword, identifier, operator, etc.)

Lexeme: Actual string of characters that forms a token

Pattern: Rule describing how tokens are formed

Symbol Table: Data structure storing identifiers and their attributes

Transition Diagram: Visual representation of token recognition logic

FSA: Mathematical model for pattern recognition with states and transitions

Critical Implementation Points

1. **Precedence matters:** Keywords before identifiers, longer operators before shorter
2. **Backtracking is essential:** Must be able to undo failed pattern attempts
3. **Pointer management is crucial:** Careful tracking of lexeme boundaries
4. **Error recovery:** Graceful handling of invalid input

Common Pitfalls and Solutions

Problem: Keyword vs Identifier confusion **Solution:** Check keyword table after identifier pattern matches

Problem: Operator precedence issues (`<=` vs `< + =`)

Solution: Order patterns by length/specificity

Problem: Whitespace handling **Solution:** Explicit whitespace-skipping functions

Problem: End-of-file handling **Solution:** Special EOF handling in all states

Study Tips

1. **Practice drawing transition diagrams** for different patterns
2. **Trace through examples** step by step with pointers
3. **Understand the relationship** between FSAs and implementation code
4. **Focus on error cases** - what happens when patterns don't match?
5. **Connect to real tools** - examine output from actual lexical analyzers

Review Questions

1. How does the lexical analyzer handle the sequence `<=` vs `< =` ?
2. What happens when an identifier matches a keyword?

3. Why is backtracking necessary in lexical analysis?
4. How do accepting states work in transition diagrams?
5. What's the difference between a lexeme and a token?

This comprehensive understanding of lexical analysis forms the foundation for all subsequent compiler phases and is essential for anyone working with programming language tools or building parsers and interpreters.

Syntax Analysis

What is Syntax?

Definition and Core Concepts

Syntax refers to the rules that define the structure of valid statements in a programming language. It's essentially the "grammar" of a programming language that determines how symbols, keywords, and operators can be combined to form meaningful programs.

Grammar - The Foundation of Syntax

Grammar is a precise, yet easy-to-understand specification that defines the syntactic rules of a language. Think of grammar as a formal rulebook that:

- **Precisely** defines what constitutes a valid program
- **Clearly** communicates these rules to both humans and machines
- **Systematically** organizes the structure for efficient processing

Why Grammar Matters

A properly designed grammar provides structure that is useful for translating source programs into executable logic. Without clear grammatical rules, compilers wouldn't know how to interpret your code.

Real-World Examples

Example 1: English vs Programming Language Syntax

- English: "The cat sits on the mat" (Subject-Verb-Object structure)
- C++: `int x = 5;` (Type-Identifier-Assignment-Value-Semicolon structure)

Example 2: Syntax Rules in Different Languages

- **Python:** `if condition:` (colon required, indentation-based blocks)
- **Java:** `if (condition) { }` (parentheses and braces required)
- **JavaScript:** `if (condition) { }` or `if (condition) statement;` (flexible bracing)

Why Syntax is Important in Language Design

1. **Readability:** Well-designed syntax makes code easier to read and understand
 2. **Writability:** Good syntax reduces the cognitive load on programmers
 3. **Reliability:** Clear syntactic rules help prevent ambiguity and errors
 4. **Tool Support:** Consistent syntax enables better IDE support, syntax highlighting, and automated tools
 5. **Maintenance:** Clean syntax makes code easier to maintain and debug
-

Role of the Parser

What is a Parser?

A **parser** (also called a **syntax analyzer**) is a crucial component of a compiler that takes a stream of tokens from the lexical analyzer and determines if they form a valid program according to the language's grammar rules.

Position in the Compiler Pipeline

```
Source Code → [Lexical Analyzer] → [Parser] → [Rest of Frontend] → Intermediate Representation
                ↓ tokens                ↓ parse tree                ↓
                                Symbol Table
```

Detailed Flow Explanation

1. **Input:** The parser receives tokens from the lexical analyzer

2. **Processing:** It applies grammar rules to check if the token sequence is valid
3. **Output:** If valid, it produces a parse tree; if invalid, it reports syntax errors
4. **Interaction:** It communicates with the symbol table to store identifiers and their properties

Parser Functions in Detail

Primary Responsibilities

- **Syntax Validation:** Ensures the input conforms to the language grammar
- **Structure Recognition:** Identifies the hierarchical structure of the program
- **Parse Tree Generation:** Creates a tree representation showing how the program is structured
- **Error Detection:** Identifies and reports syntax errors with location information

Real-World Example: Parsing an Assignment Statement

Input Code: `int result = x + y * 2;`

Lexical Analysis Output (tokens):

```
[INT] [IDENTIFIER:result] [ASSIGN] [IDENTIFIER:x] [PLUS] [IDENTIFIER:y] [MULTIPLY] [NUMBER:2]
[SEMICOLON]
```

Parser's Work:

1. Recognizes `int result` as a variable declaration
2. Identifies `= x + y * 2` as an assignment expression
3. Builds a parse tree showing operator precedence (multiplication before addition)
4. Validates that the semicolon properly terminates the statement

Parse Tree Structure (simplified):

```

Declaration
├─ Type: int
├─ Identifier: result
└─ Assignment
    └─ Expression
        └─ Identifier: x
            └─ Addition
                └─ Identifier: y
                    └─ Multiplication
                        └─ Number: 2

```

Classes of Languages

Programming languages can be categorized based on how their parsers are typically implemented. The two main classes are:

LL Languages - Manual Implementation

LL stands for "Left-to-right scan, Leftmost derivation"

Characteristics:

- **Top-down parsing:** Starts with the start symbol and works down to terminals
- **Predictive parsing:** Can determine which production rule to use by looking ahead
- **Manual implementation:** Often hand-written by compiler developers
- **Left recursion issues:** Cannot handle left-recursive grammars directly

Examples of LL Languages:

- **Pascal:** Simple, predictable syntax
- **Early versions of C:** Relatively straightforward parsing
- **Some domain-specific languages:** Designed with LL parsing in mind

Real-World Example - LL Parsing:

```
// Pascal-style if statement (LL-friendly)
if condition then
    statement1
else
    statement2;
```

The parser can easily predict what comes next:

- Sees `if` → expects condition, then `then`
- Sees `else` → expects statement
- No ambiguity about what production rule to apply

LR Languages - Automated Implementation

LR stands for "Left-to-right scan, Rightmost derivation in reverse"

Characteristics:

- **Bottom-up parsing:** Starts with terminals and builds up to the start symbol
- **More powerful:** Can handle a broader class of grammars
- **Automated tools:** Typically generated using tools like YACC, Bison, or ANTLR
- **Handles ambiguity better:** Can resolve many ambiguous constructs

Examples of LR Languages:

- **C/C++:** Complex syntax with operator precedence
- **Java:** Object-oriented features require sophisticated parsing
- **Python:** Complex expression handling and dynamic features
- **JavaScript:** Flexible syntax with many edge cases

Real-World Example - LR Parsing:

```
// C-style expression (requires LR parsing)
result = a + b * c++;
```

The parser must:

1. Handle operator precedence (`*` before `+`)
2. Process postfix increment (`++`)
3. Resolve potential ambiguities
4. Build the correct parse tree structure

Why Use Automated Tools for LR?

Manual Implementation Challenges:

- LR parse tables can have thousands of states
- Complex conflict resolution
- Difficult to maintain and debug

Tool Benefits:

- **YACC/Bison**: Generates efficient parsers from grammar specifications
- **ANTLR**: Modern tool with excellent error handling and tree building
- **PLY**: Python-based parser generator

Error Handling

Syntax errors are among the most common compilation errors. A well-designed parser must handle these errors gracefully.

The Reality of Syntax Errors

Statistical Fact: A majority of compilation errors occur during the syntax analysis phase. This makes robust error handling crucial for user experience.

Error Handler Goals

1. Clear and Accurate Reporting

What this means: Error messages should tell the programmer:

- **Type of error:** What went wrong?
- **Location:** Where did it happen?
- **Context:** What was expected?

Good Error Message Example:

```
Error at line 15, column 23: Expected ';' after expression
    result = x + y
                ^
```

Bad Error Message Example:

```
Syntax error
```

2. Error Recovery and Continuation

Purpose: Don't stop at the first error; find as many errors as possible in one compilation pass.

Why important:

- Saves developer time (fix multiple errors at once)
- Provides better overview of code problems
- More efficient development workflow

3. Performance Preservation

Requirement: Error handling shouldn't significantly slow down compilation of correct programs.

Implementation consideration: Error handling code paths should be optimized to be fast when no errors occur.

Viable-Prefix Property

Definition: The parser should maintain the ability to recognize valid prefixes of the language even when errors occur.

Practical meaning: After encountering an error, the parser should be able to synchronize and continue parsing the rest of the program meaningfully.

Offloading to Semantic Analysis

Some errors are better caught during semantic analysis rather than syntax analysis:

Symbol Table References

```
int main() {  
    undeclared_variable = 5; // Semantic error, not syntax error  
}
```

- **Syntactically correct:** Follows grammar rules
- **Semantically incorrect:** Variable not declared

Type Checking

```
int x = "hello"; // Type mismatch - semantic error
```

- **Syntactically valid:** Assignment statement structure is correct
 - **Semantically invalid:** String assigned to integer variable
-

Error Recovery Strategies

When syntax errors occur, parsers need strategies to recover and continue parsing. Here are the four main approaches:

1. Panic Mode Recovery

How it Works

1. **Discard tokens** until a "synchronizing token" is found
2. **Synchronizing tokens** are usually delimiters like `;`, `}`, `end`, etc.
3. **Resume parsing** from the synchronizing token

Language-Specific Synchronization

- **C/C++:** `;`, `}`, `{`
- **Pascal:** `;`, `end`, `begin`
- **Python:** Newline, dedent tokens
- **Java:** `;`, `}`, `class`, `public`

Real-World Example

```
int main() {  
    int x = 5 + + +; // Error: malformed expression  
    int y = 10;      // Parser recovers here after finding ';'   
    return 0;  
}
```

Recovery process:

1. Error detected at malformed expression
2. Parser discards tokens until it finds `;`
3. Parsing resumes with `int y = 10;`

Advantages and Disadvantages

Pros:

- Simple to implement
- Fast recovery
- Works well for many common errors

Cons:

- May skip over valid code
- Can miss subsequent errors in the discarded section
- Recovery quality depends on synchronizing token choice

2. Phrase Level Recovery

How it Works

Performs local correction by making small changes to the input stream around the error location.

Types of Local Corrections

1. **Insertion:** Add missing tokens
2. **Deletion:** Remove extra tokens
3. **Replacement:** Change incorrect tokens

Real-World Examples

Missing Semicolon:

```
// Input (error)
int x = 5
int y = 10;

// Phrase-level correction: Insert ';'
int x = 5;
int y = 10;
```

Extra Comma:

```
// Input (error)
function call(a, b, , c)

// Phrase-level correction: Remove extra comma
function call(a, b, c)
```

Wrong Operator:

```
// Input (error)
if (x = = 5)

// Phrase-level correction: Fix spacing
if (x == 5)
```

Implementation Considerations

- Requires heuristics to guess programmer intent
- Must be careful not to change semantics
- Works best for simple, obvious errors

3. Error Productions

How it Works

Add special grammar rules that recognize common error patterns and provide meaningful error messages.

Implementation Strategy

Extend the grammar with productions that match common mistakes:

```
// Normal grammar rule
statement: IF '(' expression ')' statement
         | IF '(' expression ')' statement ELSE statement
         ;

// Error production for missing parentheses
statement: IF expression statement
         {
             yyerror("Missing parentheses around if condition");
             // Generate corrected code
         }
         ;

// Error production for assignment instead of comparison
expression: IDENTIFIER '=' expression
          {
              yyerror("Did you mean '==' instead of '='?");
              // Treat as comparison for continued parsing
          }
          ;
```

Real-World Examples

Missing Braces:

```
// Common error
if (condition)
    statement1;
    statement2; // This looks like it's in the if, but it's not!

// Error production can catch and warn about this
```

Assignment in Condition:

```
if (x = 5) { // Common mistake: = instead of ==
    // Error production provides helpful message
}
```

Advantages

- Provides specific, helpful error messages
- Can suggest corrections
- Maintains good parsing state

Disadvantages

- Increases grammar complexity
- Can create ambiguities
- Requires anticipating common errors

4. Global Correction

How it Works

Find the minimal set of edits (insertions, deletions, substitutions) needed to transform the erroneous input into a syntactically correct program.

The Algorithm Concept

1. **Generate alternatives:** Consider all possible single edits
2. **Cost calculation:** Assign costs to different types of edits
3. **Optimization:** Find the minimum-cost sequence of edits
4. **Validation:** Ensure result is syntactically correct

Theoretical Example

```
// Input with errors
int main( {
    int x = 5
    return x
}

// Possible corrections (with edit costs):
// 1. Insert ')' after 'main(' - cost: 1
// 2. Insert ';' after 'x = 5' - cost: 1
// 3. Insert ';' after 'return x' - cost: 1
// Total minimum cost: 3 edits

// Corrected result:
int main() {
    int x = 5;
    return x;
}
```

Why Rarely Used in Practice

- **Computationally expensive:** Exponential complexity
- **Ambiguous results:** Multiple corrections may have the same cost
- **Over-correction:** May change programmer intent
- **Implementation complexity:** Very difficult to implement correctly

Modern Alternative: IDE Integration

Instead of global correction in compilers, modern development environments provide:

- Real-time syntax highlighting
 - Immediate error underlining
 - Auto-completion suggestions
 - Quick-fix suggestions
-

Practical Examples and Applications

Example 1: Building a Simple Expression Parser

Let's trace through parsing a mathematical expression:

Input: 3 + 4 * 2

Lexical Analysis Output:

```
[NUMBER:3] [PLUS] [NUMBER:4] [MULTIPLY] [NUMBER:2]
```

Grammar Rules:

```
Expression → Term (('+' | '-') Term)*  
Term       → Factor (('*' | '/') Factor)*  
Factor     → NUMBER | '(' Expression ')'
```

Parsing Steps:

1. **Start with Expression**
2. **Parse first Term:** Recognizes 3
3. **See '+' operator:** Continue with Expression rule
4. **Parse second Term:** 4 * 2
 - Factor: 4
 - Operator: *
 - Factor: 2
 - Result: (4 * 2) due to precedence
5. **Final result:** 3 + (4 * 2)

Example 2: Real Compiler Error Messages

GCC C++ Compiler:

```
int main() {
    int x = 5
    return 0;
}
```

Error Message:

```
error: expected ';' before 'return'
```

Rust Compiler:

```
fn main() {
    let x = 5
    println!("{}", x);
}
```

Error Message:

```
error: expected `;`, found `println`
--> main.rs:3:5
  |
2 |     let x = 5
  |           ^ help: add `;` here
3 |     println!("{}", x);
  |     ^^^^^^ expected `;`
```

Notice how Rust provides:

- Exact location with line numbers
- Visual indicator with `^`
- Helpful suggestion with "help: add `;` here"

Example 3: Parser Generator Tools

Using ANTLR (ANother Tool for Language Recognition)

Grammar file (Calculator.g4):

```
grammar Calculator;

// Parser rules
expr:  expr ('*'|'/') expr    # MulDiv
      |  expr ('+'|'-') expr    # AddSub
      |  INT                  # int
      |  '(' expr ')'          # parens
      ;

// Lexer rules
INT:   [0-9]+ ;
WS:    [ \t\r\n]+ -> skip ;
```

Generated parser automatically handles:

- Operator precedence
- Parse tree construction
- Error reporting
- Recovery strategies

Example 4: Modern IDE Error Handling

Visual Studio Code with TypeScript:

```
function calculateTotal(price: number, tax: number) {
    return price + tax;
    console.log("Calculation done"); // Unreachable code warning
}

let result = calculateTotal("100", 0.08); // Type error
```

IDE provides:

- Red squiggly lines under errors
- Hover tooltips with error descriptions
- Quick-fix suggestions
- Real-time error checking as you type

Example 5: Domain-Specific Language (DSL) Parsing

SQL Query Parsing:

```
SELECT name, age
FROM users
WHERE age > 18
ORDER BY name;
```

Parser must handle:

- Keywords (SELECT , FROM , WHERE , ORDER BY)
 - Identifiers (column and table names)
 - Operators (> , =)
 - Expressions and conditions
 - Syntax variations across SQL dialects
-

Study Tips and Key Takeaways

Core Concepts to Remember

1. **Syntax is the grammar of programming languages** - it defines valid program structure
2. **Parsers bridge the gap** between human-readable code and machine-processable structures
3. **LL and LR represent different parsing strategies** with different trade-offs
4. **Error handling is crucial** for developer productivity and experience
5. **Recovery strategies balance simplicity and effectiveness**

Practice Questions

1. **Design Exercise:** Create grammar rules for a simple if-statement in your favorite language
2. **Error Analysis:** Given a syntax error, determine which recovery strategy would work best
3. **Tool Comparison:** Research and compare different parser generators (ANTLR, Yacc, etc.)
4. **Real-world Application:** Find examples of syntax errors in actual code and analyze how different compilers handle them

Further Exploration

- **Implement a simple recursive descent parser** for arithmetic expressions
- **Study the grammar specifications** of your favorite programming language
- **Experiment with parser generator tools** like ANTLR or PLY
- **Analyze error messages** from different compilers and IDEs for the same syntax errors

This study guide covers the fundamental concepts of syntax analysis as presented in the course material, with detailed explanations and practical examples to reinforce understanding.