CS345H: Programming Languages

Lecture 4: Implementation of Lexical Analysis

Thomas Dillig

Announcements



Announcements

- WA1 and PA0 are due Today
- WA2 and PA1 out today :-)

Announcements

- WA1 and PA0 are due Today
- WA2 and PA1 out today :-)
- If you are not very, very busy right now, get started now

Last time: Specifying lexical structure using regular expressions

- Last time: Specifying lexical structure using regular expressions
- Today: How to recognize strings matching regular expressions using finite automata.

- Last time: Specifying lexical structure using regular expressions
- Today: How to recognize strings matching regular expressions using finite automata.
- We will see determinist finite automata (DFAs) and non-deterministic finite automata (NFAs)

- Last time: Specifying lexical structure using regular expressions
- Today: How to recognize strings matching regular expressions using finite automata.
- We will see determinist finite automata (DFAs) and non-deterministic finite automata (NFAs)
- High-level story: RegEx -> NFA -> DFA -> Tables

• Last lecture: How to specify the predicate $s \in L(R)$

- Last lecture: How to specify the predicate $s \in L(R)$
- But yes/no answer is not enough!

- Last lecture: How to specify the predicate $s \in L(R)$
- But yes/no answer is not enough!
- We really want to partition input into tokens

- Last lecture: How to specify the predicate $s \in L(R)$
- But yes/no answer is not enough!
- We really want to partition input into tokens
- We adapt regular expressions for this goal

Step 1: Write a regular expression for the lexemes of each token

- Step 1: Write a regular expression for the lexemes of each token
 - Integer constant: digit⁺

- Step 1: Write a regular expression for the lexemes of each token
 - Integer constant: digit⁺
 - Identifier: letter (letter + digit)*

- Step 1: Write a regular expression for the lexemes of each token
 - Integer constant: digit⁺
 - Identifier: letter (letter + digit)*
 - Lambda: 'lambda'

- Step 1: Write a regular expression for the lexemes of each token
 - Integer constant: digit⁺
 - Identifier: letter (letter + digit)*
 - Lambda: 'lambda'

▶ ...

Step 2: Construct R, matching lexemes for all tokens

- Step 2: Construct R, matching lexemes for all tokens
- ► R = Integer constant + Identifier + Lambda + ...

▶ Let the input be characters *x*₁...*x*_n

- ▶ Let the input be characters *x*₁...*x*_n
- **Step 3**: For each $1 \le i \le n$ check $x_1 \dots x_j \in L(R)$ for some j

- ▶ Let the input be characters *x*₁...*x*_n
- **Step 3**: For each $1 \le i \le n$ check $x_1 \dots x_j \in L(R)$ for some j
- ▶ Then, remove *x*₁...*x_j* from input and repeat

> There are ambiguities in this algorithm. Where?

- There are ambiguities in this algorithm. Where?
- ▶ How much input is used? What if $x_1...x_i \in L(R)$ and $x_1...x_j \in L(R)$?

- There are ambiguities in this algorithm. Where?
- ▶ How much input is used? What if $x_1...x_i \in L(R)$ and $x_1...x_j \in L(R)$?
- ► Example: identifier = letter (letter + digit)*, if = 'i' 'f'

- There are ambiguities in this algorithm. Where?
- ▶ How much input is used? What if $x_1...x_i \in L(R)$ and $x_1...x_j \in L(R)$?
- ► Example: identifier = letter (letter + digit)*, if = 'i' 'f'
- Rule: Pick longest possible string in L(R)

- There are ambiguities in this algorithm. Where?
- ▶ How much input is used? What if $x_1...x_i \in L(R)$ and $x_1...x_j \in L(R)$?
- ► Example: identifier = letter (letter + digit)*, if = 'i' 'f'
- Rule: Pick longest possible string in L(R)
- This is known as "maximal munch"

What if two rules match with the same number of characters?

- What if two rules match with the same number of characters?
- $x_1...x_i \in L(R_1)$ and $x_1...x_i \in L(R_2)$?

What if two rules match with the same number of characters?

•
$$x_1...x_i \in L(R_1)$$
 and $x_1...x_i \in L(R_2)$?

Example: "if"

What if two rules match with the same number of characters?

•
$$x_1...x_i \in L(R_1)$$
 and $x_1...x_i \in L(R_2)$?

- Example: "if"
- Rule: Use rule listed first

- What if two rules match with the same number of characters?
- $x_1...x_i \in L(R_1)$ and $x_1...x_i \in L(R_2)$?
- Example: "if"
- Rule: Use rule listed first
- This is how "if" is matched as a keyword, not identifier

What if no rule matches a prefix of the input?

- What if no rule matches a prefix of the input?
- Solution 1: Get stuck

- What if no rule matches a prefix of the input?
- Solution 1: Get stuck \Rightarrow Unacceptable

- What if no rule matches a prefix of the input?
- Solution 1: Get stuck \Rightarrow Unacceptable
- Better Solution: Write a rule matching all "bad" strings
Error Handling

- What if no rule matches a prefix of the input?
- Solution 1: Get stuck \Rightarrow Unacceptable
- Better Solution: Write a rule matching all "bad" strings
- Question: What kind of rule and where to place it?

We now know how we can partition input string into tokens assuming we can decide if a string is in the language described by a regular expression.

- We now know how we can partition input string into tokens assuming we can decide if a string is in the language described by a regular expression.
- Next: How to decide if $s \in L(R)$

► Regular Expressions ⇔ Specification

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation
- A finite automata formally consists of:

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation
- A finite automata formally consists of:
 - \blacktriangleright An input alphabet Σ

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation
- A finite automata formally consists of:
 - \blacktriangleright An input alphabet Σ
 - \blacktriangleright A set of states S

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation
- A finite automata formally consists of:
 - An input alphabet Σ
 - \blacktriangleright A set of states S
 - ► A start state n

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation
- A finite automata formally consists of:
 - \blacktriangleright An input alphabet Σ
 - \blacktriangleright A set of states S
 - ► A start state n
 - A set of accepting states $F \subseteq S$

- ► Regular Expressions ⇔ Specification
- ► Finite Automata ⇔ Implementation
- A finite automata formally consists of:
 - An input alphabet Σ
 - \blacktriangleright A set of states S
 - ► A start state n
 - A set of accepting states $F \subseteq S$
 - A set of transitions state \rightarrow^{input} state

• Transition $S_1 \rightarrow^{\alpha} S_2$

- Transition $S_1 \rightarrow^{\alpha} S_2$
- This means: In state S_1 and input character α , go to state S_2

- Transition $S_1 \rightarrow^{\alpha} S_2$
- This means: In state S_1 and input character α , go to state S_2
- If end of input and in accepting state \Rightarrow accept

- Transition $S_1 \rightarrow^{\alpha} S_2$
- This means: In state S_1 and input character α , go to state S_2
- If end of input and in accepting state \Rightarrow accept
- Otherwise \Rightarrow reject

It is much easier to imagine finite automata visually:

It is much easier to imagine finite automata visually:

A state:

It is much easier to imagine finite automata visually:

A state:

The start state:



It is much easier to imagine finite automata visually:

A state:

The start state:

An accepting state:



It is much easier to imagine finite automata visually:

A state:

The start state:

An accepting state:

A transition:

Here is an automaton that only accepts the string "1":

▶ Here is an automaton that only accepts the string "1":



Another Simple Example

 A finite automaton accepting any number of 1's followed by a single 0

Another Simple Example

- A finite automaton accepting any number of 1's followed by a single 0
- ► Alphabet: {0,1}

Another Simple Example

- A finite automaton accepting any number of 1's followed by a single 0
- ► Alphabet: {0,1}



And Another Example

• Alphabet: $\{0,1\}$

And Another Example

- ► Alphabet: {0,1}
- What language does this automata recognize?

And Another Example

- ► Alphabet: {0,1}
- What language does this automata recognize?



Epsilon Moves

• Another kind of transition: ε -moves

Epsilon Moves

• Another kind of transition: ε -moves



Epsilon Moves

• Another kind of transition: ε -moves



► Machine can move from state A to B without reading any input

Deterministic Finite Automata (DFA)

At most one transition per input on any state

- At most one transition per input on any state
- ▶ No ε moves

- At most one transition per input on any state
- ▶ No ε moves
- Nondeterministic Finite Automate (NFA)

- At most one transition per input on any state
- ▶ No ε moves
- Nondeterministic Finite Automate (NFA)
 - Can have multiple transitions for one input in a given state
Deterministic and Nondeterministic Automata

Deterministic Finite Automata (DFA)

- At most one transition per input on any state
- No ε moves
- Nondeterministic Finite Automate (NFA)
 - Can have multiple transitions for one input in a given state
 - ► Can have *ε*-moves

 A DFA can only take one path through the state graph that is completely determined by the input

- A DFA can only take one path through the state graph that is completely determined by the input
- NFAs can choose:

- A DFA can only take one path through the state graph that is completely determined by the input
- NFAs can choose:
 - Whether to make ε moves

- A DFA can only take one path through the state graph that is completely determined by the input
- NFAs can choose:
 - Whether to make ε moves
 - Which one of multiple transitions for a single input to take

This means: A NFA can get into multiple states at the same time

- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.

- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.



- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.
- Consider input 1 0 0



- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.
- Consider input 1 0 0



- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.
- Consider input 1 0 0



- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.
- Consider input 1 0 0



- This means: A NFA can get into multiple states at the same time
- ► Consider again the alphabet ∑ = {0,1} and the language of all strings ending in at least two 0s.





 Fundamental Result: NFAs and DFAs recognize the same set of languages (regular languages)

NFAs vs. DFAs

- Fundamental Result: NFAs and DFAs recognize the same set of languages (regular languages)
- DFAs are faster to execute, since there are no choices to consider

NFAs vs. DFAs

- Fundamental Result: NFAs and DFAs recognize the same set of languages (regular languages)
- DFAs are faster to execute, since there are no choices to consider
- But NFAs can be much simpler for the same language

NFAs vs. DFAs

- Fundamental Result: NFAs and DFAs recognize the same set of languages (regular languages)
- DFAs are faster to execute, since there are no choices to consider
- But NFAs can be much simpler for the same language
- Result: DFAs can be exponentially larger than NFA recognizing same language

High-Level Sketch:

- High-Level Sketch:
 - Lexical Specification

- High-Level Sketch:
 - Lexical Specification
 - Regular Expressions

- High-Level Sketch:
 - Lexical Specification
 - Regular Expressions
 - NFA

- High-Level Sketch:
 - Lexical Specification
 - Regular Expressions
 - NFA
 - DFA

- High-Level Sketch:
 - Lexical Specification
 - Regular Expressions
 - NFA
 - DFA
 - Implementation of DFA

- High-Level Sketch:
 - Lexical Specification
 - Regular Expressions
 - NFA
 - DFA
 - Implementation of DFA



 For each kind of regular expression, define an NFA and combine

- For each kind of regular expression, define an NFA and combine
- ▶ Will use the following notation: NFA for regular expression M:
 → M ○

- For each kind of regular expression, define an NFA and combine
- ▶ Will use the following notation: NFA for regular expression M:
 → M ○
- Base cases:

- For each kind of regular expression, define an NFA and combine
- ▶ Will use the following notation: NFA for regular expression M:
 → M ○
- Base cases:



- For each kind of regular expression, define an NFA and combine
- ▶ Will use the following notation: NFA for regular expression M:
 → M ○

Base cases:









• Consider the regular expression $(1+0)^*1$

• Consider the regular expression $(1+0)^*1$



• Consider the regular expression $(1 + 0)^*1$



• Consider the regular expression $(1+0)^*1$


Example of Regular Expression to NFA conversion

• Consider the regular expression $(1+0)^*1$



Example of Regular Expression to NFA conversion

• Consider the regular expression $(1+0)^*1$



Insight: Simulate the NFA

- Insight: Simulate the NFA
- At any given time, the NFA is in a set of states

- Insight: Simulate the NFA
- At any given time, the NFA is in a set of states
- \blacktriangleright States in the DFA \Rightarrow all (reachable) subsets of states in the NFA

- Insight: Simulate the NFA
- At any given time, the NFA is in a set of states
- \blacktriangleright States in the DFA \Rightarrow all (reachable) subsets of states in the NFA
- Start State:

- Insight: Simulate the NFA
- At any given time, the NFA is in a set of states
- \blacktriangleright States in the DFA \Rightarrow all (reachable) subsets of states in the NFA
- \blacktriangleright Start State: the set of states reachable through ε moves from the NFA start state

- Insight: Simulate the NFA
- At any given time, the NFA is in a set of states
- \blacktriangleright States in the DFA \Rightarrow all (reachable) subsets of states in the NFA
- \blacktriangleright Start State: the set of states reachable through ε moves from the NFA start state
- Add transition $A \rightarrow^a B$ to DFA iff:

- Insight: Simulate the NFA
- At any given time, the NFA is in a set of states
- \blacktriangleright States in the DFA \Rightarrow all (reachable) subsets of states in the NFA
- \blacktriangleright Start State: the set of states reachable through ε moves from the NFA start state
- Add transition $A \rightarrow^a B$ to DFA iff:
 - B is in the set of states reachable from any state in A after seeing input a, considering ε moves as well



























































NFA to DFA: How many states?

We need a state in the DFA for each set of states the NFA can be in
- We need a state in the DFA for each set of states the NFA can be in
- How many different states?

- We need a state in the DFA for each set of states the NFA can be in
- How many different states?
- ► If there are N states, the NFA must be in some subset of those N states

- We need a state in the DFA for each set of states the NFA can be in
- How many different states?
- ► If there are N states, the NFA must be in some subset of those N states
- How many subsets of N states?

- We need a state in the DFA for each set of states the NFA can be in
- How many different states?
- ► If there are N states, the NFA must be in some subset of those N states
- How many subsets of N states? 2^N

- A DFA can be implemented by a 2D table T
 - One dimension is "states"

- A DFA can be implemented by a 2D table T
 - One dimension is "states"
 - Other dimension is "input symbols"

- One dimension is "states"
- Other dimension is "input symbols"
- ▶ For every transition $A \rightarrow^{c} B$, define T[A,c]=B

- One dimension is "states"
- Other dimension is "input symbols"
- ▶ For every transition $A \rightarrow^{c} B$, define T[A,c]=B
- DFA "execution": If in state A and input c, read T[A,c] = B and skip to state B

- One dimension is "states"
- Other dimension is "input symbols"
- ▶ For every transition $A \rightarrow^{c} B$, define T[A,c]=B
- DFA "execution": If in state A and input c, read T[A,c] = B and skip to state B
- Very efficient

Table Implementation of a DFA



Table Implementation of a DFA



	0	1
S	Т	U
Т	Т	U
U	Т	U

 Writing regular expressions as NFAs and converting them to DFAs is exactly what flex does

- Writing regular expressions as NFAs and converting them to DFAs is exactly what flex does
- In fact, if you open the auto-generated flex file lex.yy.c, you will see these tables emitted

- Writing regular expressions as NFAs and converting them to DFAs is exactly what flex does
- In fact, if you open the auto-generated flex file lex.yy.c, you will see these tables emitted
- But, these DFAs can be huge

- Writing regular expressions as NFAs and converting them to DFAs is exactly what flex does
- In fact, if you open the auto-generated flex file lex.yy.c, you will see these tables emitted
- But, these DFAs can be huge
- In practice, flex-like tools trade off speed for space in the choice of NFA and DFA representations