

Compiler Construction

Prof. Dr.-Ing. Jeronimo Castrillon
TU Dresden – Cfaed
WS 2022

Chair for Compiler Construction
Helmholtzstrasse 18
jeronimo.castrillon@tu-dresden.de

cfaed.tu-dresden.de



TECHNISCHE
UNIVERSITÄT
DRESDEN

DRESDEN
concept



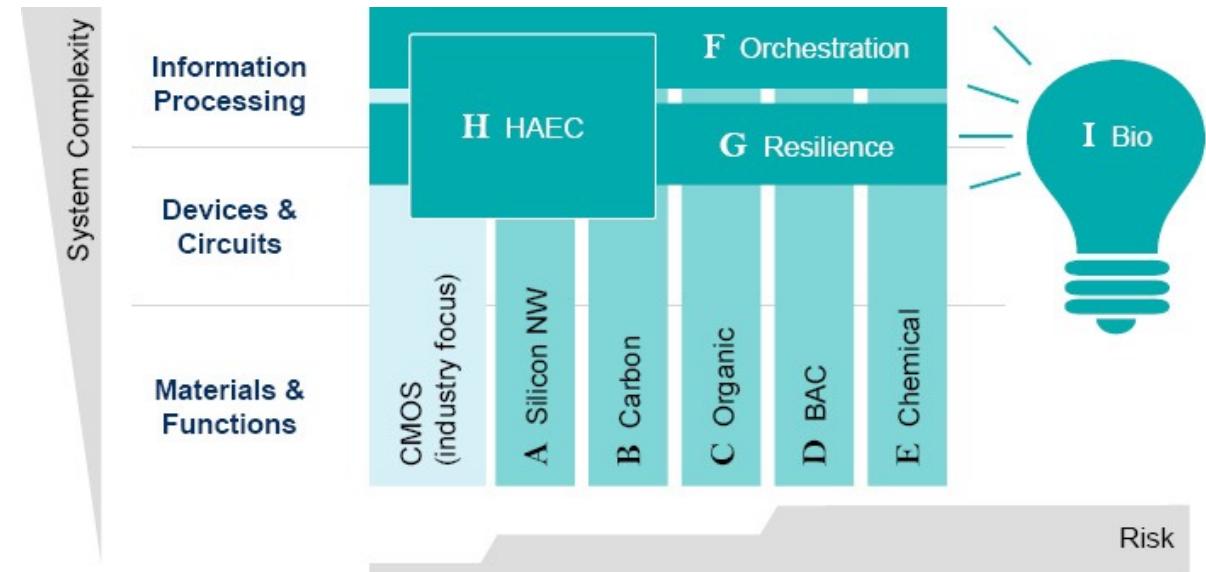
DFG

WR

WISSENSCHAFTSRAT

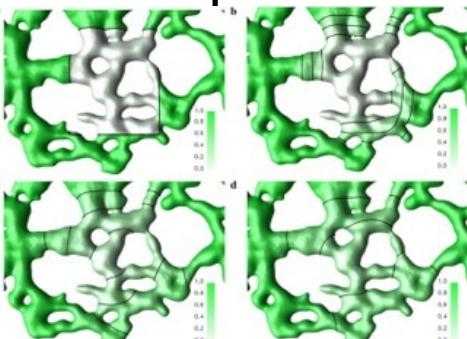
About the Chair for Compiler Construction (CCC)

- Cfaed: Center for Advancing Electronics Dresden
- CCC: Programming methods
 - Parallel heterogeneous systems
 - Post-CMOS systems
 - Domain-specific languages and optimization
 - Goals: Performance, energy efficiency, resilience

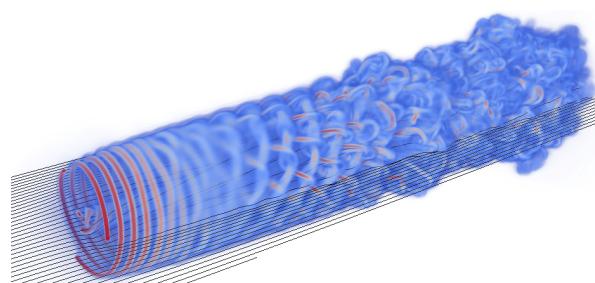
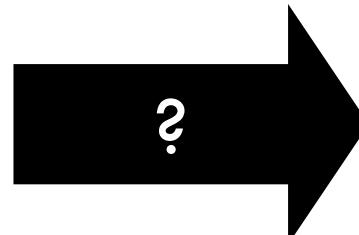


Domain-specific languages (DSLs)

- ❑ Raise the level of abstraction for particular application domains to enable more powerful optimizations for large parallel systems



Sbalzarini, Ivo F. et al. 2006. „Simulations of (An)isotropic Diffusion on Curved Biological Surfaces“.



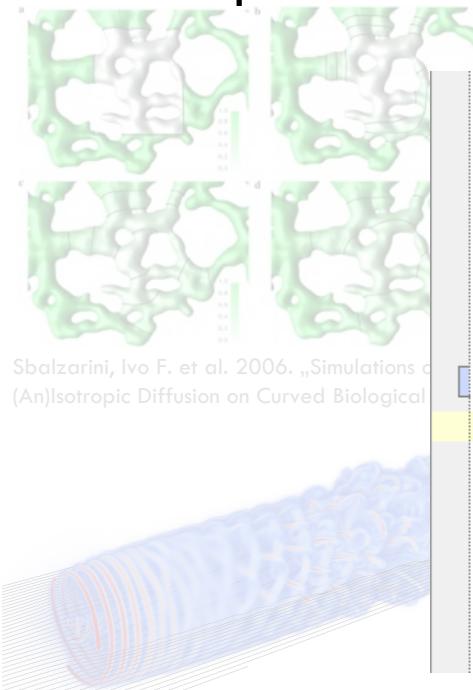
Chatelain et al. 2011. „Vortex particle-mesh methods with immersed lifting lines applied to the Large Eddy Simulation of wind turbine wakes“



Bull HPC-Cluster „Taurus“, ZIH,
TU Dresden

Domain-specific languages (2)

- ❑ Raise the level of abstraction for particular application domain to enable more powerful optimizations for large parallel systems



Sbalzarini, Ivo F. et al. 2006. „Simulations of (An)isotropic Diffusion on Curved Biological Tissues“

```
external real Dv = 1.0;      "diffusion constant of V"

phase initialize
  << ... >>
end initialize

phase solve
  field <real , 1> U;
  field <real , 1> V;
  integer t;

  
$$\frac{\partial U}{\partial t} = Du * \nabla^2 U - U * V^2 + F * (1 - U);$$

  
$$\frac{\partial V}{\partial t} = Dv * \nabla^2 V + U * V^2 - (F + kRate) * V;$$

end solve

phase finalize
  << ... >>
end finalize
```

Chatelain et al. 2011. „Vortex particle-mesh methods with immersed lifting lines applied to the Large Eddy Simulation of wind turbine wakes“



Contact: Nesrine Khouzami
nesrine.khouzami@tu-dresden.de

Bull HPC-Cluster „Taurus“, ZIH,
TU Dresden

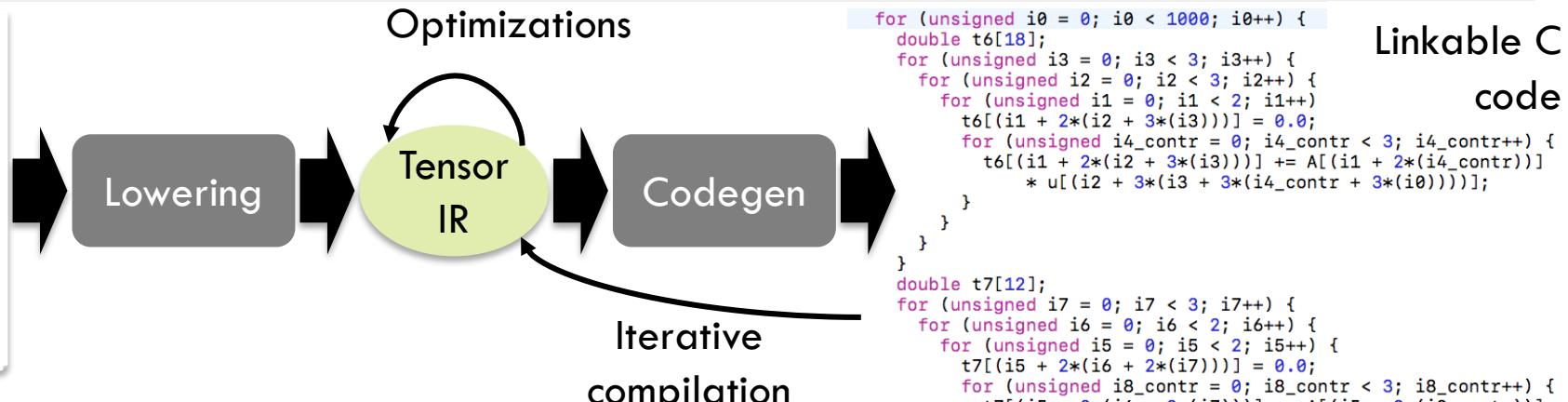
Tensor languages

```
source =  
type matrix : [mp np] &  
type tensorIN : [np np np ne] &  
type tensorOUT : [mp mp mp me] &  
  
var input A : matrix &  
var input u : tensorIN &  
var input output v : tensorOUT &  
var input alpha : [] &  
var input beta : [] &  
  
v = alpha * (A # A # A # u .  
[[5 8] [3 7] [1 6]]) + beta * v
```

Fortran embedding

```
auto A = Matrix("A", m, n);  
auto u = Tensor<3>("u", n, n, n);  
  
auto Interpolation = (A*A*A)(u);  
  
const CFDlang::Program *Prog =  
    Emitter.lowerExpr(Interpolation, "v");
```

C++ embedding (with templates)



- ❑ Language/IR design and implementation
- ❑ C++/Haskell
- ❑ Tensorflow
- ❑ Parallelization

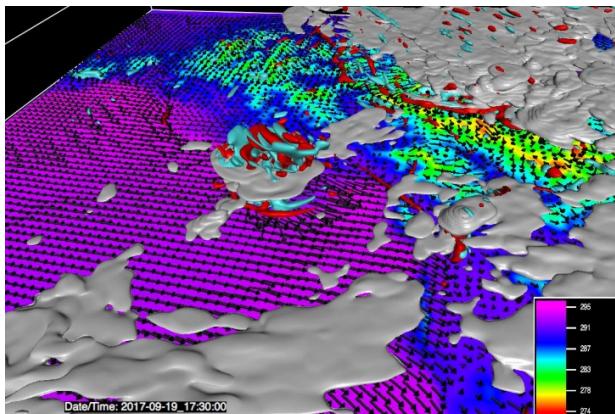
```
for (unsigned i0 = 0; i0 < 1000; i0++) {  
    double t6[18];  
    for (unsigned i3 = 0; i3 < 3; i3++) {  
        for (unsigned i2 = 0; i2 < 3; i2++) {  
            for (unsigned i1 = 0; i1 < 2; i1++)  
                t6[(i1 + 2*(i2 + 3*(i3)))] = 0.0;  
            for (unsigned i4_contr = 0; i4_contr < 3; i4_contr++) {  
                t6[(i1 + 2*(i2 + 3*(i3)))] += A[(i1 + 2*(i4_contr))]  
                    * u[(i2 + 3*(i3 + 3*(i4_contr + 3*(i0))))];  
            }  
        }  
    }  
}  
double t7[12];  
for (unsigned i7 = 0; i7 < 3; i7++) {  
    for (unsigned i6 = 0; i6 < 2; i6++) {  
        for (unsigned i5 = 0; i5 < 2; i5++) {  
            t7[(i5 + 2*(i6 + 2*(i7)))] = 0.0;  
            for (unsigned i8_contr = 0; i8_contr < 3; i8_contr++) {  
                t7[(i5 + 2*(i6 + 2*(i7)))] += A[(i5 + 2*(i8_contr))]  
                    * t6[(i6 + 2*(i7 + 3*(i8_contr)))];  
            }  
        }  
    }  
}  
double t8[1];  
double t9[1];  
for (unsigned i11 = 0; i11 < 2; i11++) {  
    for (unsigned i10 = 0; i10 < 2; i10++) {  
        for (unsigned i9 = 0; i9 < 2; i9++) {  
            t9[0] = 0.0;  
            for (unsigned i12_contr = 0; i12_contr < 3; i12_contr++) {  
                t9[0] += A[(i9 + 2*(i12_contr))] * t7[(i10 + 2*(i11 +  
                    2*(i12_contr)))];  
            }  
            t8[0] = alpha[0] * t9[0];  
            double t10[1];  
            t10[0] = beta[0] * v[(i9 + 2*(i10 + 2*(i11 + 2*(i10))))]  
                ;  
            v[(i9 + 2*(i10 + 2*(i11 + 2*(i10))))] = t8[0] + t10[0];  
        }  
    }  
}
```

Weather modeling: DSLs and synthesis flows

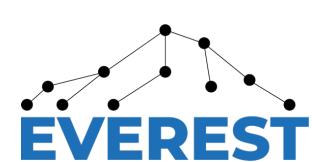
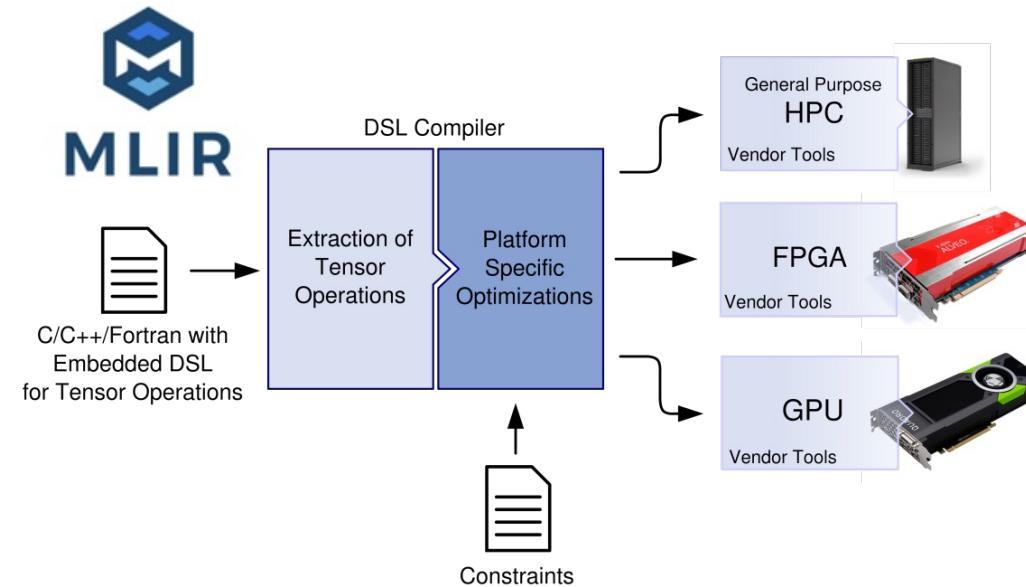
- Stencil and Tensor Operations in Weather Modelling
- Optimizations for computations in WRF (CPU, GPU, FPGA)
- Integration of stencils on top of MLIR dialects



Contact: Karl Friebel
karl.friebel@tu-dresden.de



Source: CIMA Foundation

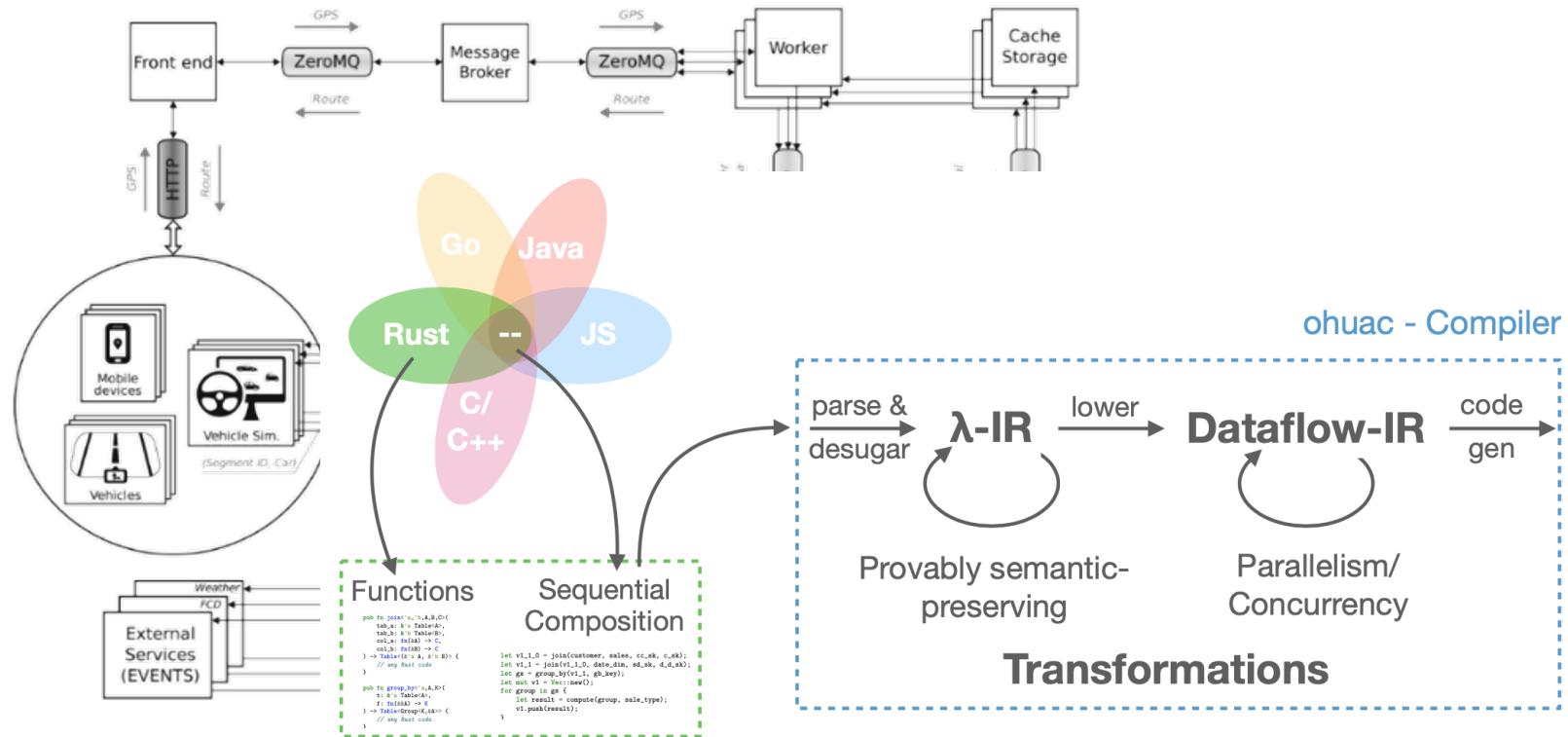


Traffic routing: Coordination + Determinism

- Dataflow functional abstraction
- Implicit parallelism: compiler controlled coordination



Contact: Felix Wittwer
felix.wittwer1@tu-dresden.de



A Toolflow for Programming Cyber-physical Systems (1)

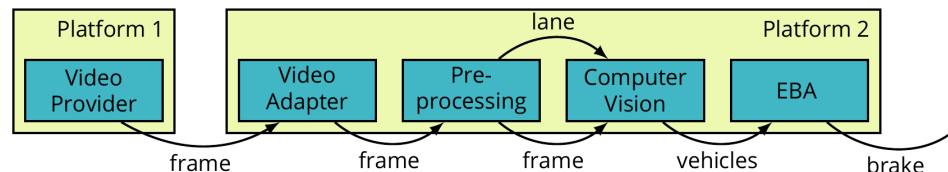
Context

- “A cyber-physical system is a computer without a keyboard that can kill you.”
- How can we build confidence in the safety of cyber-physical systems?

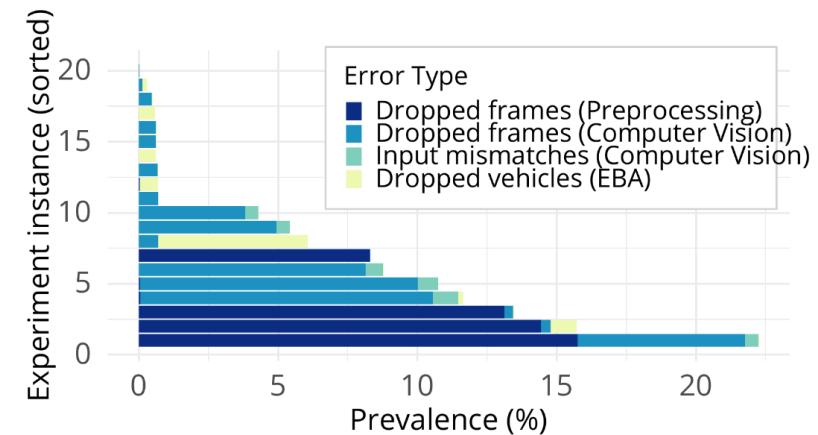


Example

- In today's distributed systems, nondeterminism is a big issue.
- Even the demonstrator application (brake assistant) of the upcoming automotive standard Adaptive AUTOSAR exposes nondeterminism:



Data may be lost on any of the communication channels. This could have fatal consequences!



Contact: Christian Menard
Christian.Menard@tu-dresden.de

Compiler Optimizations for Racetrack Memories

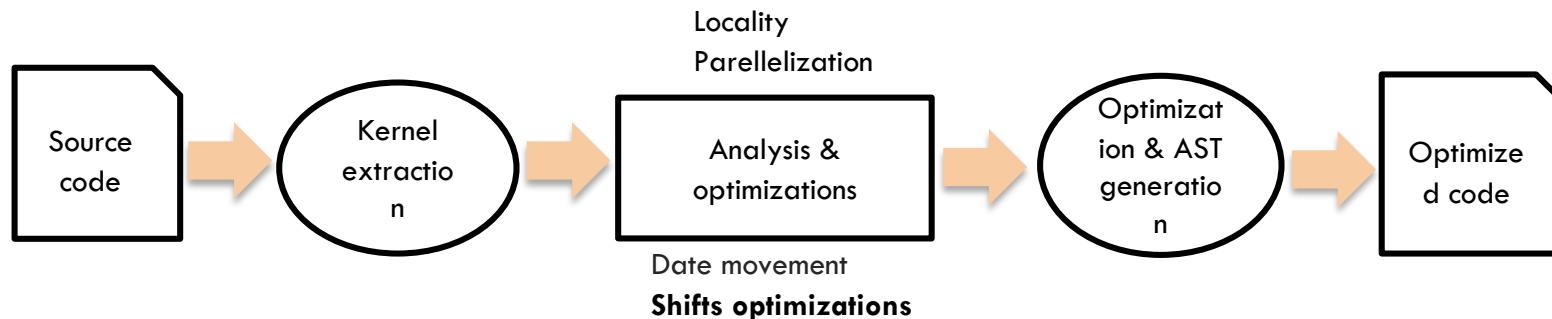
□ Context:

- Racetrack Memory (RTM) is a promising new class of the emerging non-volatile memories
- Shifts in RTMs are necessary but unwanted
- Goal: Generate efficient code that minimizes shifts in RTMs

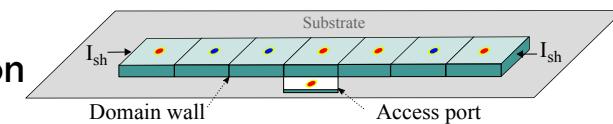


Contact: Asif Ali Khan

asif.ali.khan@tu-dresden.de

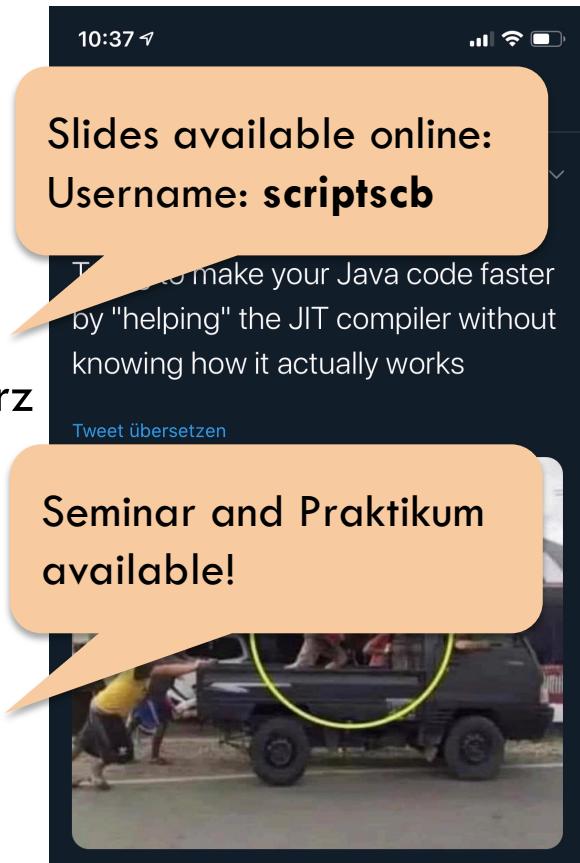


- Skills: Required C/C++; Preferable LLVM, MLIR, polyhedral compilation
- Suitable for DA or MA



About this course

- Classical compiler technology
 - Great “application” of computer science (more on this later)
 - Selected, representative optimization algorithms
- Content inspired by other lecturers
 - Prof. Leupers (RWTH), Prof. Franke (Edinburgh), Prof. Schwarz (Stanford) & others
 - Several books (more on this later)
 - Feedback is welcome!
- 2/2/0 – along **with exercise** (seminar in the website)
- **Examination:** Oral examination (normally)
- Hope to see some of you later for an *Abschlussarbeit!!*



About the exercise

- Main responsible persons: Lars Schütze
- Start: **next week**
- Structure
 - **Theoretical and practical sessions**
 - Topic-wise a week behind the lecture
 - Topic online a week before the exercise (prepare upfront)



About you

- Study programs: IST, EE, CS, CSE, DSE?
- Who knows what is a compiler, a preprocessor, a linker?
- Who has ever programmed in assembly?
- Who has ever written something like a compiler?
- Who knows what a deterministic finite state automaton is? And an LR parser?
- Who knows what is a cache? A pipeline? A register? An ALU?
- Who knows what RISC/VLIW/Superscalar are?

About compilers?

- Compilers are important, complex software systems
- Engineering: Optimization problems & approximation
- Great insight: Border between software and hardware (processor architecture)
- Learn how programming languages work & trade-offs in language design
- Ongoing: New architectures & goals → Require new compilers

- Good application of different areas of Computer Science
 - Graph theory: Coloring, topological sort, ...
 - Algorithms: Dynamic programming, greedy heuristics, ...
 - Machine Learning: Heuristic search
 - Theory: Automata, parsers, pattern matching, fixed-point algorithms, ...

Lecture overview



1. Introduction
2. Lexical analysis
3. Syntax analysis
4. Semantic analysis
5. Intermediate representation
6. Control & data-flow analysis
7. IR optimization
8. Target architectures
9. Code selection
10. Register allocation
11. Scheduling
12. Advanced topics

Compiler Construction: Frontend

Prof. Dr.-Ing. Jeronimo Castrillon
TU Dresden – Cfaed
SS 2022

Chair for Compiler Construction
Helmholtzstrasse 18
jeronimo.castrillon@tu-dresden.de

cfaed.tu-dresden.de



TECHNISCHE
UNIVERSITÄT
DRESDEN

DRESDEN
concept



DFG

WR

WISSENSCHAFTSRAT

Lecture overview – Frontend



1. Introduction
2. Lexical analysis
3. Syntax analysis
4. Semantic analysis
5. Intermediate representation
6. Control & data-flow analysis
7. IR optimization
8. Target architectures
9. Code selection
10. Register allocation
11. Scheduling
12. Advanced topics

1. Introduction

- History
- Terminology
- Compiler structure
- Trends
- Literature

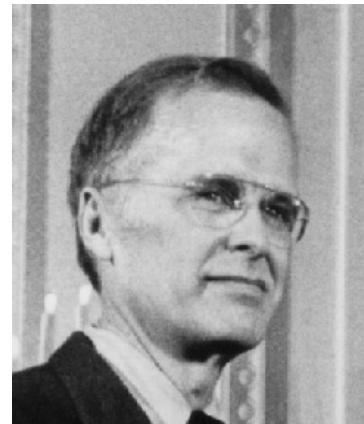


History of compilers

- < 1950 Programming in assembly
- 1950s First machine-independent languages (skepticism) – A-0 Language
To compile: “put things together” – Not really what it is today
- 1959 Complete compiler – J. Backus @ IBM for Fortran (in assembly)
Complexity: 15 man-year projects



Admiral
Grace Hopper
A-0 & COBOL
Source: wikipedia.org



John Backus
FORTRAN, IBM
Source:
columbia.edu/cU/computinghistory

History of compilers (2)

- 1960s Theoretical work for code analysis
- 1970 Bootstrapping becomes mainstream (C-compiler written in C)
First ever: LISP in 1962
- 1970 Tools to create parts of the compiler (lex & yacc)
- 1980s/90s Code generator generators
- > 2000 Optimizations (mostly backend)

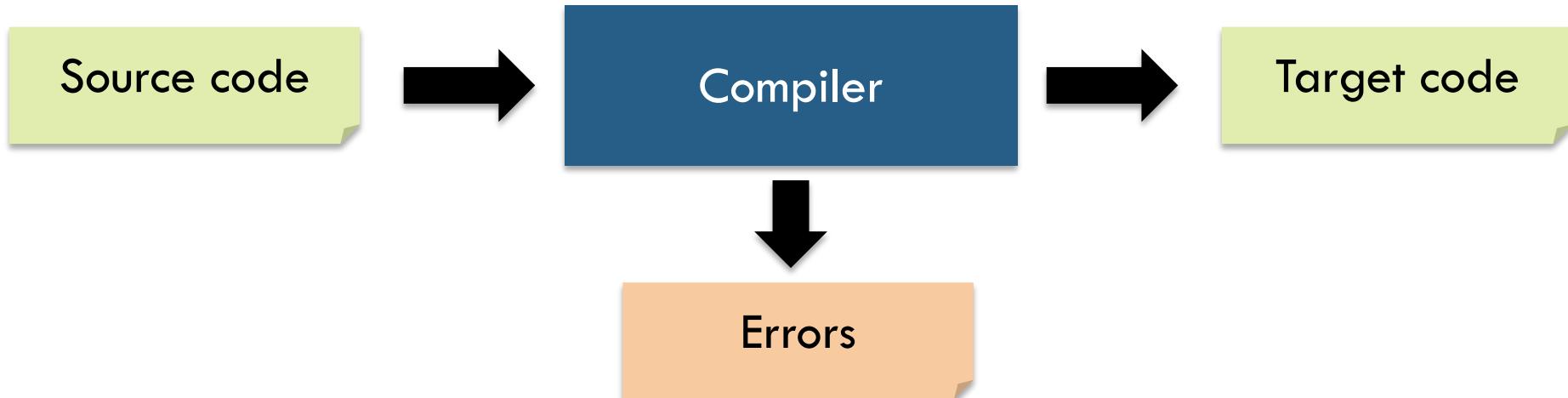
- Today SIMD, vectorization, continuous compilation, new goals,
parallelizing compilers, domain-specific languages, skeletons,
auto-tuning, ...

1. Introduction

- History
- Terminology
- Compiler structure
- Trends
- Literature



What is a compiler



- Compiler translates an input source code to a target code
- Typical: target code closer to machine code (e.g., C → assembly)
- Must recognize illegal code and generate correct code
- Must agree with lower layers (e.g., storage, linker and runtime)

Compiler vs. Interpreter

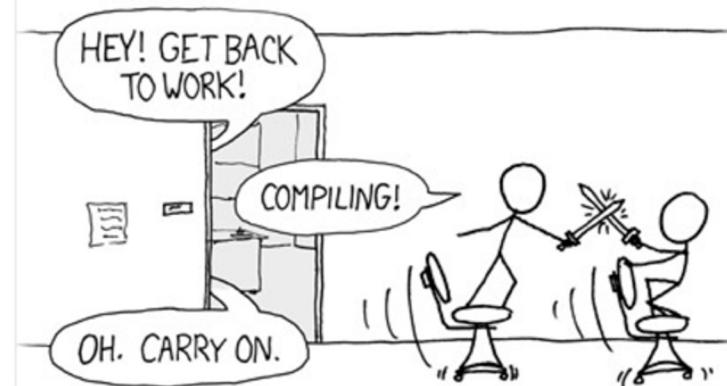


- Interpreter: Does not generate another program – stepwise translation
 - “Easier” to implement but less room for optimization
- Java byte-code is interpreted by the JVM (nowadays: just-in-time (JIT) compiled)

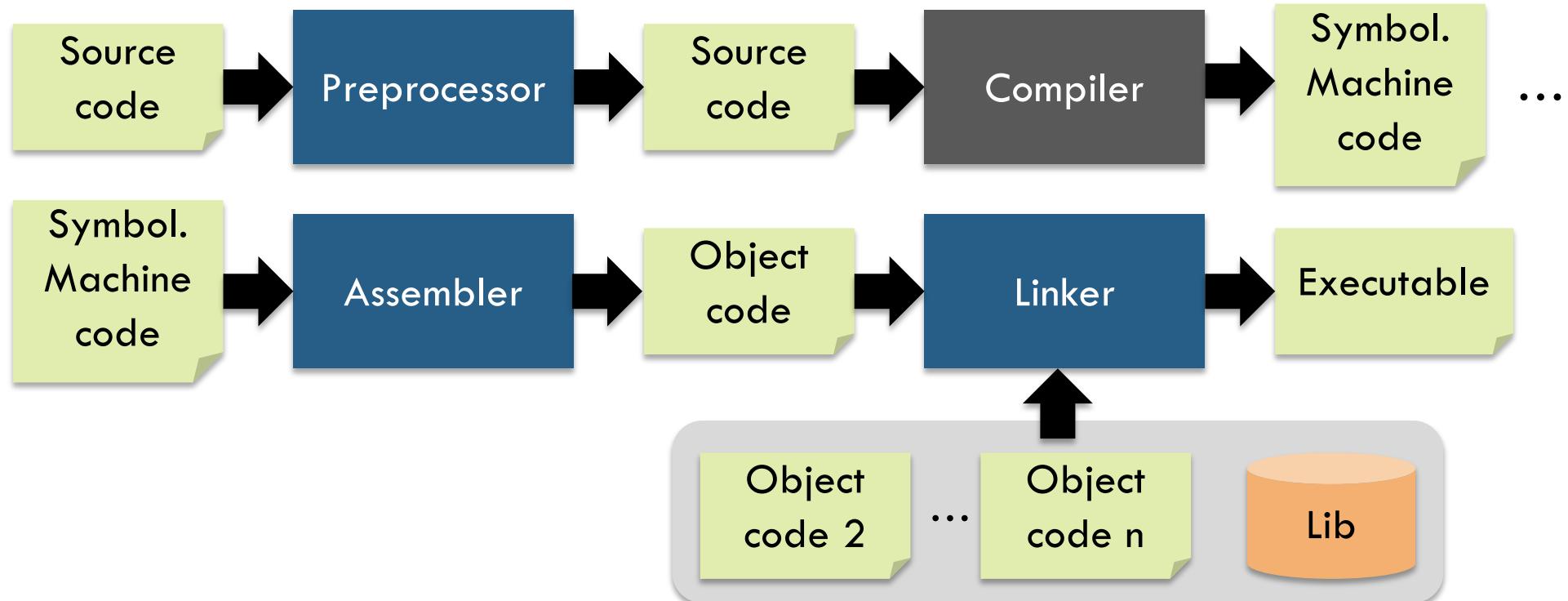
Compiler goals

- First and foremost: **Correctness**
 - Correct translation while preserving semantics
 - Incorrect code must be discarded
- Optimization goals
 - Typ.: Performance and code size (e.g., inlining)
 - New: Energy/power consumption, robustness
 - Optimality? – Undecidable, often NP complete (especially in the Backend)
- Off-line processing
 - Tolerable turn-around times – time complexity in $O(n) - O(n^2)$
 - Some domains are more patient (embedded)

THE #1 PROGRAMMER EXCUSE
FOR LEGITIMATELY SLACKING OFF:
"MY CODE'S COMPILING."



Around the compiler



Around the compiler: Preprocessor



- Prepare code before actual compilation
- Examples: C-Preprocessor (try `gcc -E main.c`)

- Constant definitions

```
#define NUM_ELEMENTS 1000
```

- Macros

```
#define MIN(x,y) (x <= y ? x : y)
```

- Include directives

```
#include <stdio.h>
```

Around the compiler: Assembler



- ❑ Translates human-readable machine code into binary

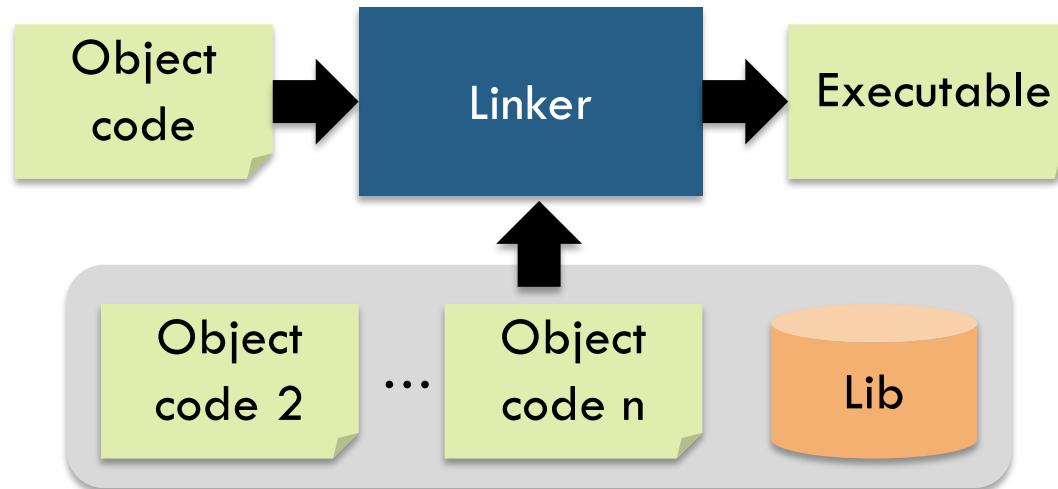
- ❑ Mnemonics (ADD) into opcodes (0011)
- ❑ Symbolic addresses into numeric values
- ❑ Generates re-locatable binary code

b = a + 2 →

MOV a,R1
ADD #2,R1
MOV R1,b

	opcode	reg.	Immediates, (*) Relocatable code
	0010	01	00000000 (*)
	0011	01	00000010
	0010	01	00000100 (*)

Around the compiler: Linker



- Link different object codes & common libraries
- Move code segments (e.g., global variables & stack)
- Fix external references (e.g., library calls)

1. Introduction

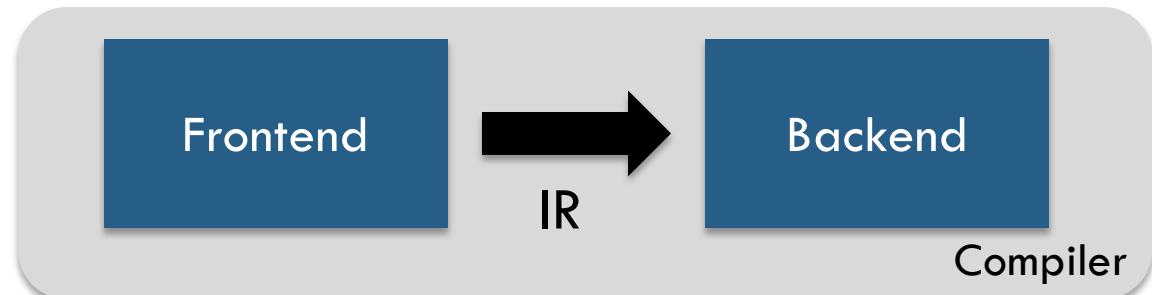
- History
- Terminology
- Compiler structure
- Trends
- Literature



Main Phases

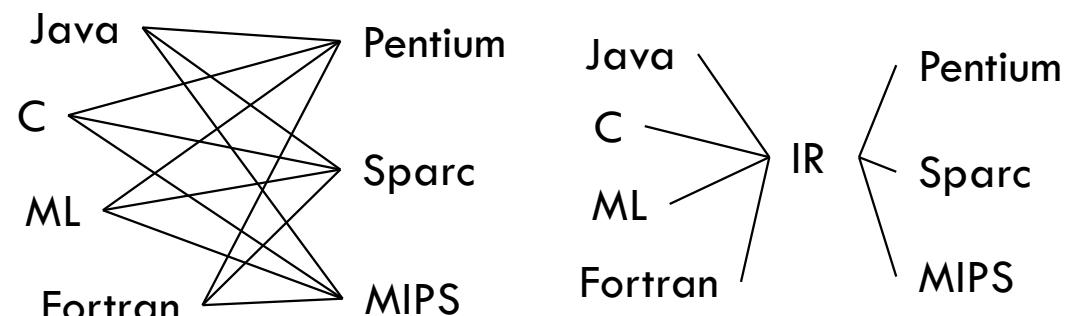
- Two phases

- IR: Intermediate representation
- Frontend: Legal code → IR
 - (Typ.) Target-independent
 - Complexity: $O(n)$ - $O(n \log n)$
- Backend: IR → Target code
 - Target-dependent
 - Complexity: (Typ.) NP complete



- Why IR? – Ideal case

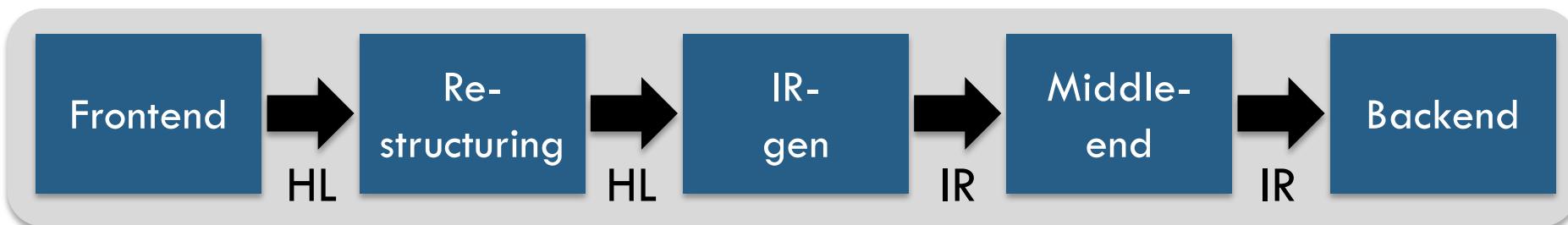
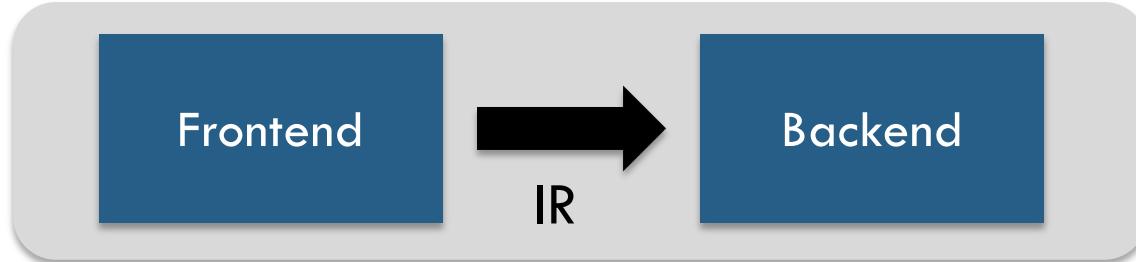
- Build $n \times m$ compilers with $n + m$ components



Adapted from: A. Appel: Modern Compiler Implementation in C.

Main Phases (2)

- Two phases
- Three phases
 - Middle-end: several optimization “passes”
- 4+ phases
 - Re-structure: almost source-level optimizations



Structure of a compiler: Frontend + Middle-end



- Lexical (scanner): Maps a character stream into words (tokens)
- Syntax (parser):
 - Recognizes “sentences of tokens” according to a grammar
 - Produces a representation of the application: **Syntax Tree (ST)**
- Semantic analysis: Adds information and checks – types, declarations, ...
- IR-generation: Abstract representation of the program, amenable for code generation (backend) – A (abstract) syntax tree is a form of IR
- IR-optimization: Simplification & improvements (e.g., remove redundancies)

Structure of a compiler: Backend



- Code selection: Decide which instructions should implement the IR
- Register allocation: Decide in which register to place variables
- Scheduling: Decide when to execute the instructions (e.g., ordering in assembly program) & ensure conformance with interfaces and constraints

Lexical Analysis – Example



```
while (y < z)
{
    int x = a + b;
    y += x;
}
```

T_While
T_LeftParen
T_Identifier y
T_Less
T_Identifier z
T_RightParen
T_OpenBrace
T_Int T_Identifier x

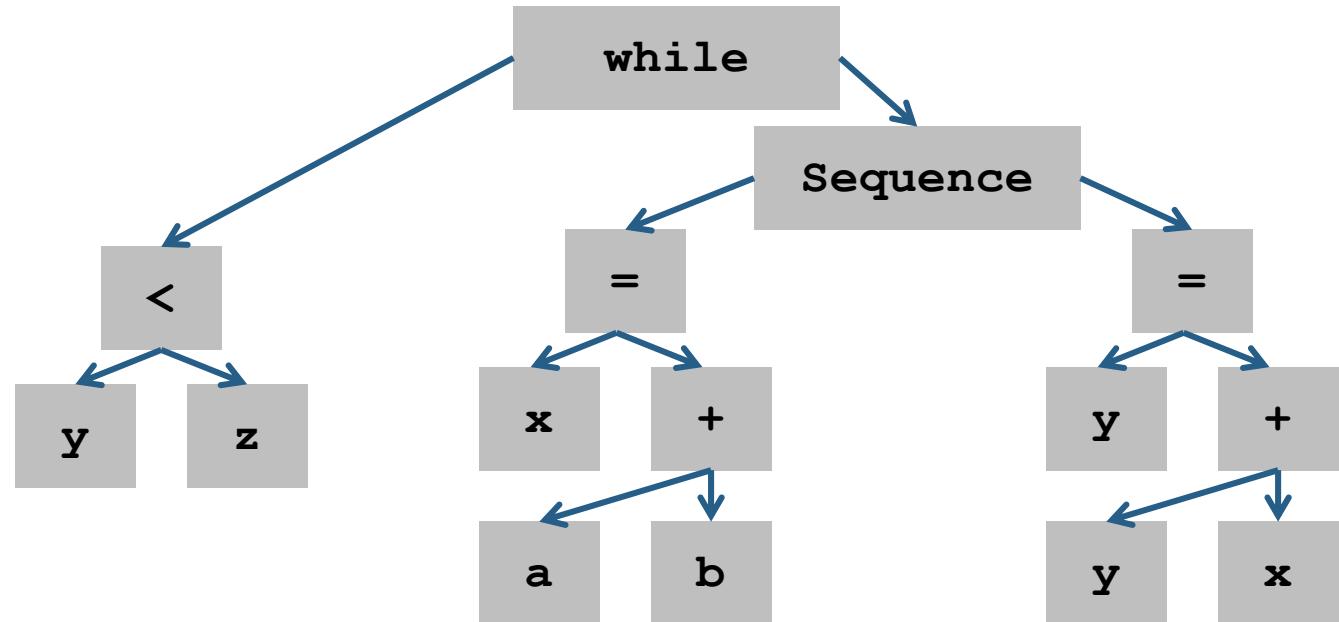
T_Assign T_Identifier a
T_Plus T_Identifier b
T_Semicolon
T_Identifier y
T_PlusAssign
T_Identifier x
T_Semicolon
T_CloseBrace

Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Syntax Analysis – Example

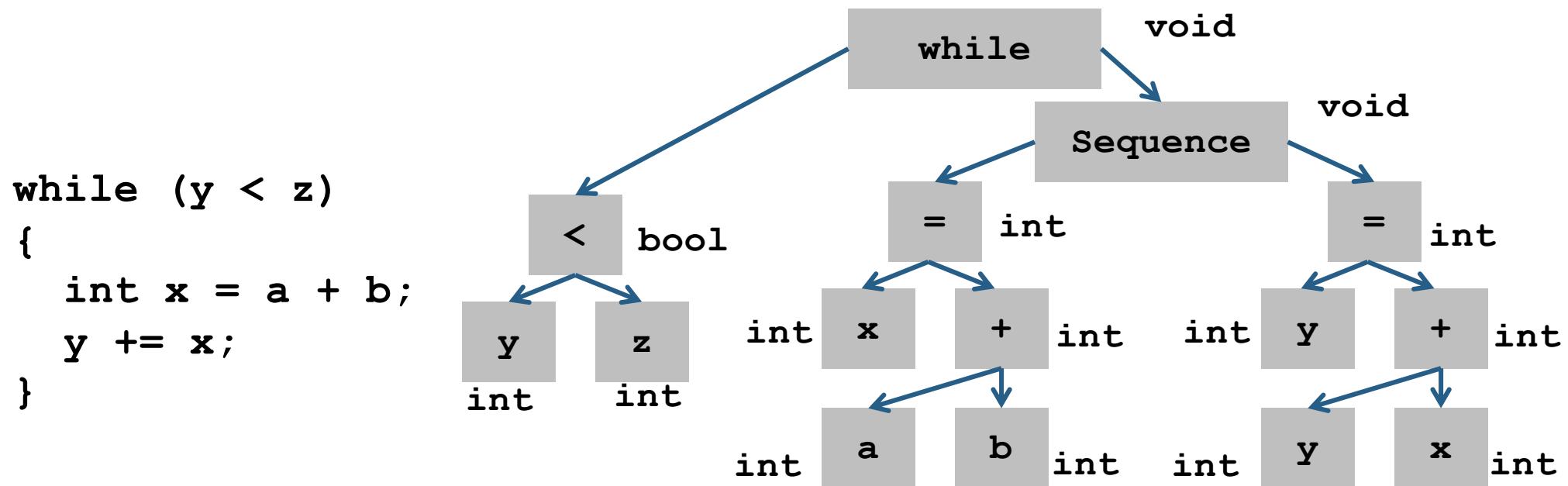


```
while (y < z)
{
    int x = a + b;
    y += x;
}
```



Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Semantic Analysis – Example



Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

IR generation – Example



```
while (y < z)
{
    int x = a + b;
    y += x;
}
```

```
<...> // First exit condition
Loop: x = a + b
      y = x + y
      _t1 = y < z
      if _t1 goto Loop
```

Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

IR Optimization – Example



```
while (y < z)
{
    int x = a + b;
    y += x;
}
```

```
Loop: x = a + b
      y = x + y
      _t1 = y < z
      if _t1 goto Loop
```

```
x = a + b
Loop: y = x + y
      _t1 = y < z
      if _t1 goto Loop
```

Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Code generation – Example



```
while (y < z)
{
    int x = a + b;
    y += x;
}
```

```
x = a + b
Loop: y = x + y
      _t1 = y < z
      if _t1 goto Loop
```

```
ADD R1, R2, R3
Loop: ADD R4, R1, R4
      SLT R6, R4, R5
      BEQ R6, loop
```

Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Why phases?

- The compiler is a complex software ($>10^6$ lines of code)
 - It must be very reliable for SW development!
 - Modularization is important
 - Good modularization allows to add more optimization passes later
- Mature theory and tools for automatic generation of phases (e.g., in frontend)
- Problem with phases: Phase-ordering problem limits optimization potential

1. Introduction

- History
- Terminology
- Compiler structure
- Trends
- Literature



Compiler trends

- Automation in compiler generation
 - Code generator generators
 - Retargetable compilers (from processor models)
- Code optimization: Whole program analysis
- New optimization goals: Real-time & energy constraints
- OOP Compilers for new architectures (VLIW & GPGPU)
- Parallelization for multi-core architectures
- Compilers for domain-specific languages
- Continuous compilation

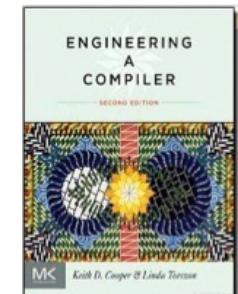
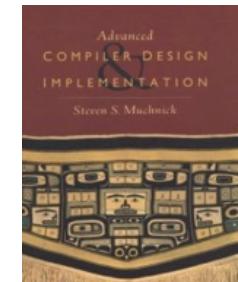
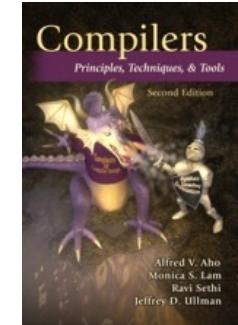
1. Introduction

- History
- Terminology
- Compiler structure
- Trends
- Literature



Recommended books

- A. Aho, M. Lam, R. Sethi, J. Ullman: **Compilers – Compilers: Principles, Techniques, and Tools**, 2nd Ed., Addison-Wesley, 2008, ISBN 978-3-8273-7097-6
- A. Appel: **Modern Compiler Implementation in C**, Cambridge University Press, 1998, ISBN 0-521-58390-X
- S. Muchnik: **Advanced Compiler Design and Implementation**, Morgan Kaufmann Publishers, 1997, ISBN 1-55860-320-4
- K. Cooper, L. Torczon. **Engineering a Compiler**, Second Edition. Morgan-Kaufmann, ISBN 1-55860-698-X
- R. Wilhelm, D. Maurer: **Übersetzerbau**, 2. Auflage, Springer, 1997, ISBN 3-540-61692-6



Where are we?

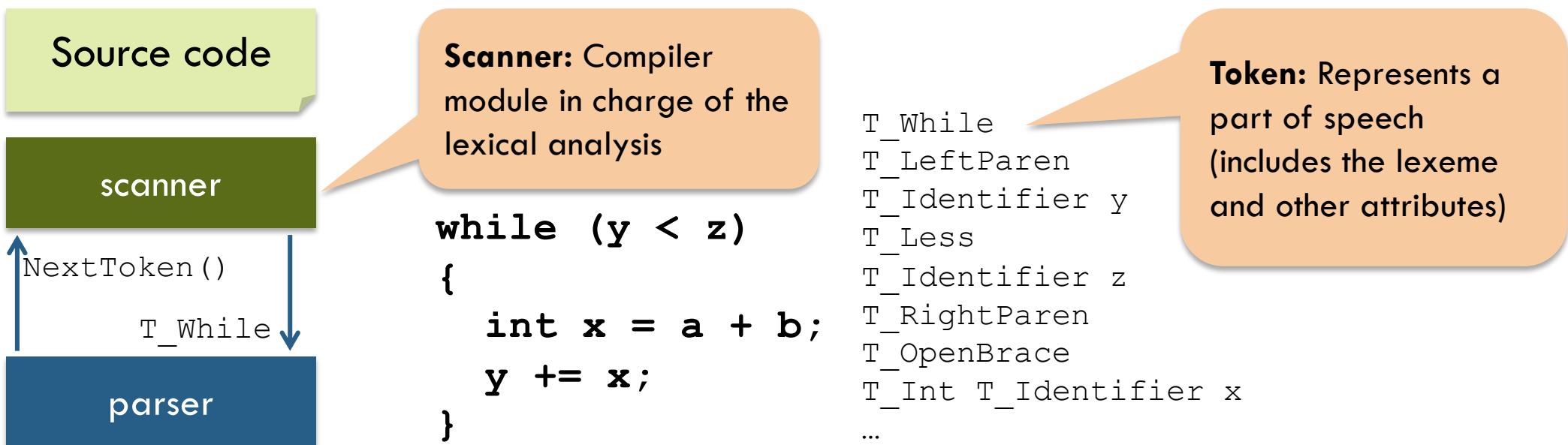
1. Introduction
2. Lexical analysis
3. Syntax analysis
4. Semantic analysis
5. Intermediate representation
6. Control & dataflow analysis
7. Compiler optimization
8. Target architectures
9. Code selection
10. Register allocation
11. Scheduling
12. Advanced topics

2. Lexical Analysis

- The scanner
- Regular expressions
- Finite automata
- Scanner generation with Lex



The scanner: Recall



The scanner: Implementation idea

- How to recognize a potentially infinite set of lexemes represented by a token?
 - Token: IDENTIFIER → a, Email_Adress, b12345, ~~1234abc~~, ...
- Need a formalism
 - Formal languages
 - Use regular expressions to identify the lexemes
 - Automata theory
 - Build recognizers automatically

2. Lexical Analysis

- The scanner
- Regular expressions
- Finite automata
- Scanner generation with Lex



Def.: An **alphabet** is a non-empty finite set of symbols

□ Examples

- Binary digits $\Sigma = \{0,1\}$
- Decimal digits $\Sigma = \{0,1,2,3,4,5,6,7,8,9\}$
- C keywords $\Sigma = \{\text{if, then, else, for, while, do, switch, ...}\}$

Def.: A **string** over an alphabet Σ is any finite sequence of symbols from Σ

□ Examples

- Binary strings $\Sigma = \{0,1\}$ 0, 01010, ...
- Decimal numbers $\Sigma = \{0,1,2,3,4,5,6,7,8,9\}$ 12, 1980, ...

□ Set of strings of any length over Σ is denoted Σ^* (Kleene closure of Σ)

$$\Sigma^* = \bigcup_{n \in N} \Sigma^n$$

ϵ : Empty word

- Example: $\{0,1\}^* = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots\}$

Def.: A **formal language** is a set of strings over an alphabet Σ
(i.e., a subset of Σ^*)

□ Alphabet and strings

- Binary strings: $\Sigma_1 = \{0,1\}$ 0, 01010, ...
- Decimal numbers: $\Sigma_2 = \{0,1,2,3,4,5,6,7,8,9\}$ 12, 1880, ...

□ Examples

- $L = \{\text{strings over } \Sigma_1 \text{ of length 3}\} = \{000,001,010,011,100,101,110,111\}$
- $L = \{\text{multiples of 3 over } \Sigma_2\} = \{3,6,12,\dots\}$

Regular Expressions – regexp

- A regexp is sequence of symbols (constant and operator symbols) that forms a search pattern
- regexps are a family of descriptions that can be used to describe so-called **regular languages** (generated by Type-3 grammars)
- Given a regexp R, we denote the language associated with R, as L(R)

Def.: A **regexp** over an alphabet Σ is a sequence of symbols from $(\Sigma \cup \epsilon)$ and operators $\{^*, ^\circ, | \}$, constructed according to the following rules:

1. Empty string: ϵ is a regexp, with $L(\epsilon) = \{\epsilon\}$
2. Literal: every symbol a in Σ is a regexp, with $L(a) = \{a\}$
3. ...

Regular Expressions – regexp (2)

Def.: A **regexp** over an alphabet Σ is a sequence of symbols from $(\Sigma \cup \varepsilon)$ and operators $\{^*, ^o, | \}$, constructed according to the following rules:

1. Empty string: ε is a regexp, with $L(\varepsilon) = \{\varepsilon\}$
2. Literal: Every symbol a in Σ is a regexp, with $L(a) = \{a\}$
3. Union: Let M and N be regexp, then $M | N$ is a regexp, with
 $L(M | N) = L(M) \cup L(N)$
4. Concatenation: Let M and N be regexp, then M^oN (MN) is a regexp, with
 $L(M^oN) = \{s_1s_2; s_1 \in L(M) \wedge s_2 \in L(N)\}$
5. Kleene star: Let M be a regexp, then M^* is a regexp, with
 $L(M^*) = (L(M))^*$

regexp: Precedence

- Precedence rules

(R)

R^*

$R_1 R_2$

$R_1 | R_2$

- Operators $(^, |)$ are left-associative

- Example: $01^*0 | 10 \rightarrow ((0(1^*))0) | (10)$

regexp: Examples

- $R = (0 \mid 1)^* 1$
- $R = (0 \mid 1) (0 \mid 1) (0 \mid 1)$
- $R = a a^* (b \mid \epsilon) a^*$
- ??
- $L(R)$: Odd binary numbers
- $L(R)$: Binary strings of length 3
- $L(R)$: Strings of ‘a’s, beginning with ‘a’ and with at most one ‘b’
- $L(R)$: binary strings with a ‘0’ in the middle and the same number of ‘1’s left and right (e.g., 101, 11011, 111101111,...)

Def.: Two regexps M and N are equivalent, $M=N$, if $L(M) = L(N)$

- Commutativity: $M|N = N|M$ ($MN \neq NM$)
- Associativity: $(M|N)|O = M|(N|O)$, holds also for \circ
- Distributivity: $M(N|O) = MN|MO$
- Examples
 - $a^*(a^*bb^*|a^*) = a^*a^*bb^*|a^*a^* = a^*bb^*|a^*$
 - $(a|b)^*(b|a)^* = (a|b)^*(a|b)^* = (a|b)^*$

2. Lexical Analysis

- The scanner
- Languages and regular expressions
- Finite automata
- Scanner generation with Lex



Match regexp: Scanner



```
while (y < z)
{
    int x = a + b;
    y += x;
}
```

`T_While T_LeftP T_ID T_Less T_ID T_RightP
T_OpenBrace
 T_Int T_ID T_eq T_ID T_plus T_ID
 T_ID T_plus T_eq T_ID
T_CloseBrace`

- Regular expressions can be implemented using finite automata (non-deterministic and deterministic)

Non-deterministic finite automata (NFA)

Def.: An **NFA** is a 5-tuple $(Q, \Sigma, \Delta, q_0, F)$, where Q is a finite set of states, Σ is an alphabet and

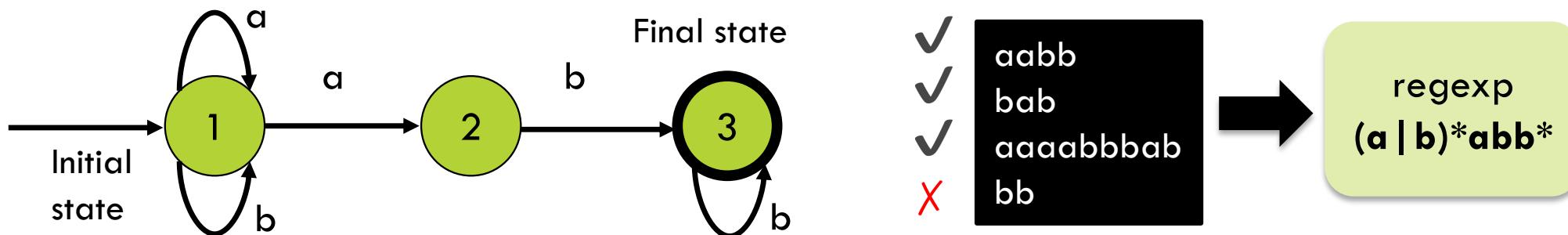
1. Δ : Transition **relation** $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$
2. q_0 : Is an initial state ($q_0 \in Q$)
3. F : Set of final states ($F \subseteq Q$)

- Given an NFA in current state q , the next state is decided by
 - Input: Move to q' by reading $a \in \Sigma$ if $(q, a, q') \in \Delta$
 - ϵ -moves: By moving to state q' , if $(q, \epsilon, q') \in \Delta$
- An NFA accepts a string if a state can be reached that is in F
- The set of all strings accepted by an NFA Z defines a language $L(Z)$

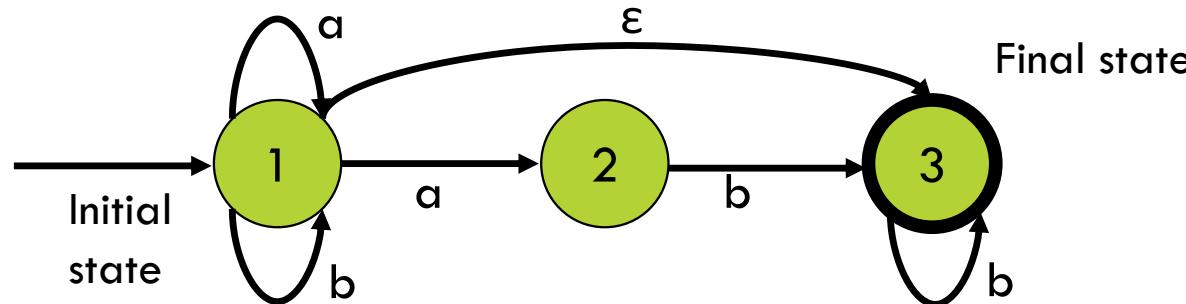
Graph representation

Def.: A **directed graph** (digraph) is a pair $G=(V,E)$ where V is a set of elements called vertices and $E \subseteq V \times V$ is a set of ordered pairs called edges (arcs)

- NFA digraph representation
 - States as graph nodes and transitions as graph edges
- Example with $Q=\{1,2,3\}$, $\Sigma = \{a,b\}$, $q_0=1$ and $F= \{3\}$



Non-determinism



- On state $q=1$, not exactly defined what to do on arrival of symbol ‘a’
- ϵ -moves introduce non-determinism as well
- Recall: An NFA **accepts** a string if there is at least a sequence of states that finishes with a state in F
- For **aaab**
 - 1-a-1-a-1-a-1- ϵ -3-b-3
 - 1-a-1-a-2-a-Error

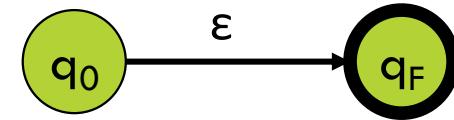
Relation regexp – NFA

For every regexp R there is a NFA Z , with $L(R) = L(Z)$

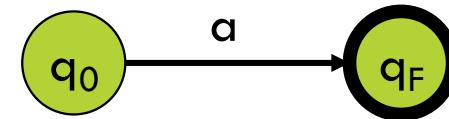
- Construction idea (Thompson's algorithm): Associate an NFA to each regexp construction rule, with
 - Exactly one accepting state q_F
 - No transitions out of q_F
 - No transitions into the starting state q_0

Thompson's algorithm: NFA construction rules

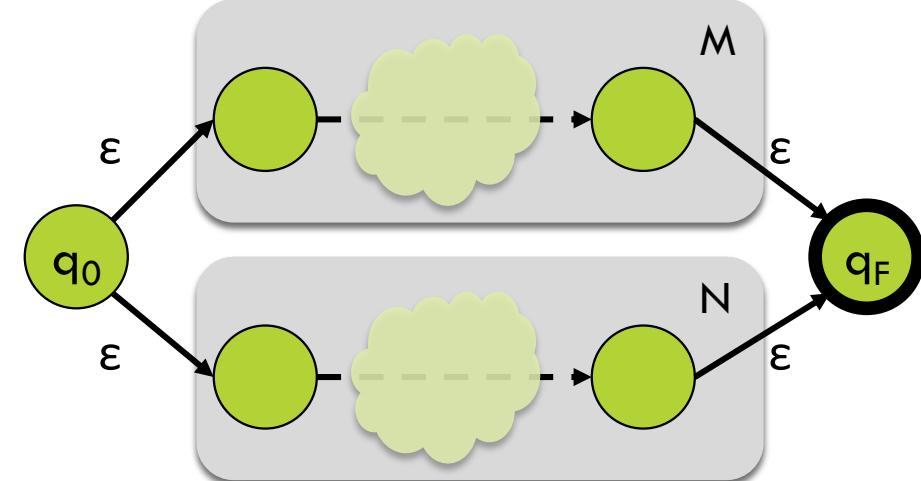
- Rule 1: Empty string $R = \epsilon$



- Rule 2: Alphabet symbol $R = a$

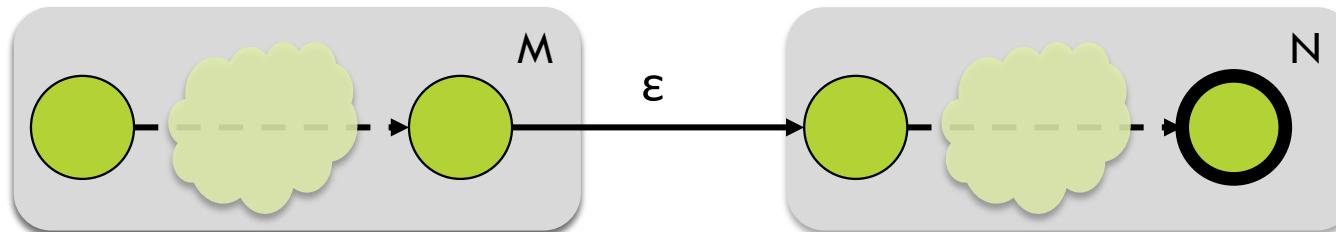


- Rule 3: Union $R = M | N$

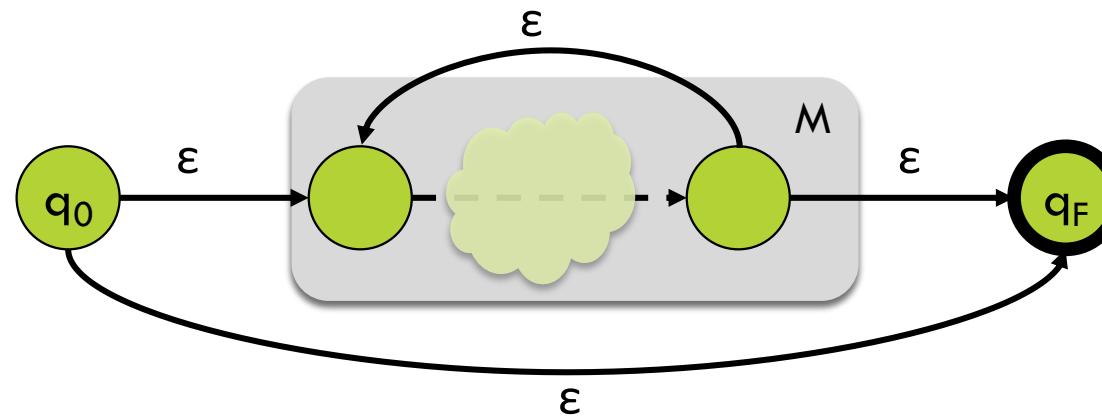


Thompson's algorithm: NFA construction rules (2)

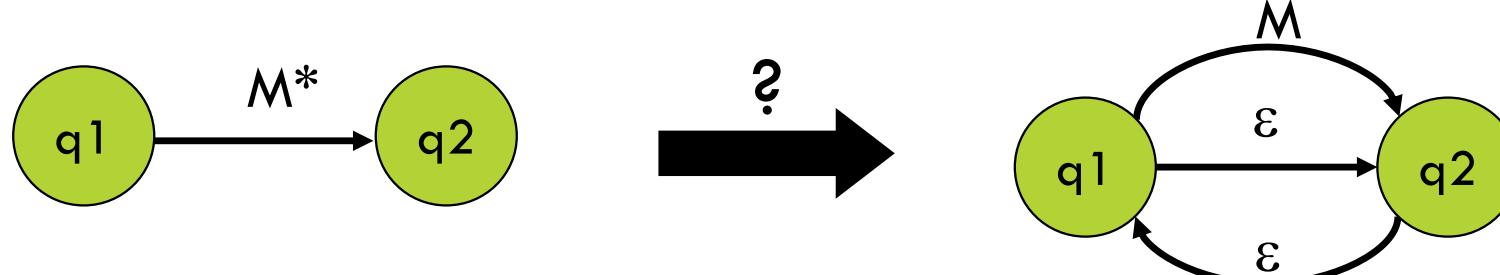
- Rule 4: Concatenation $R = MN$



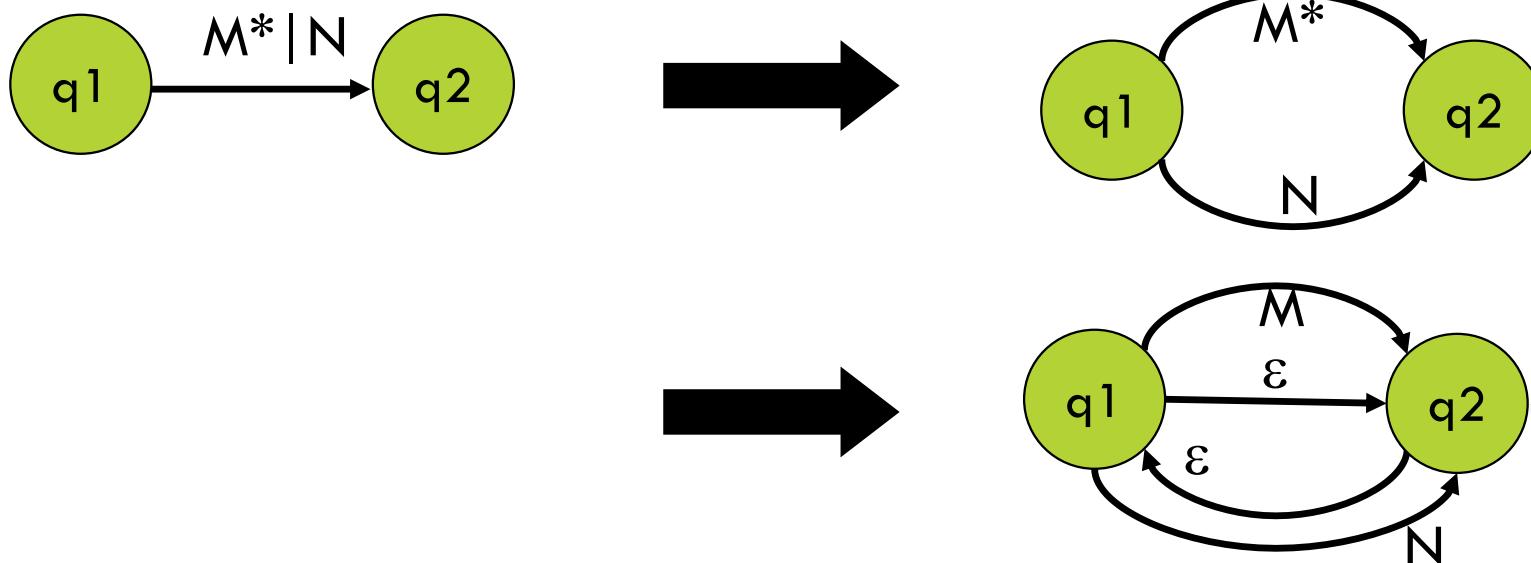
- Rule 5: Kleene star $R = M^*$



Kleene start: Easier way?

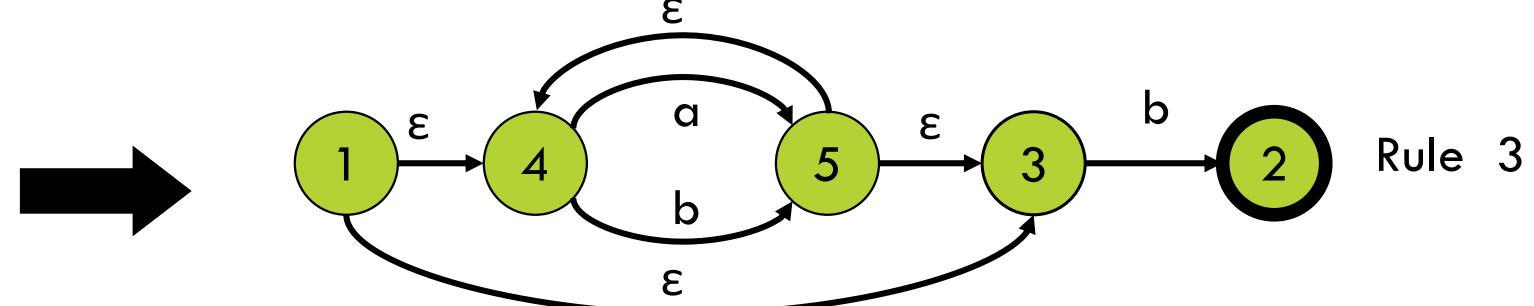
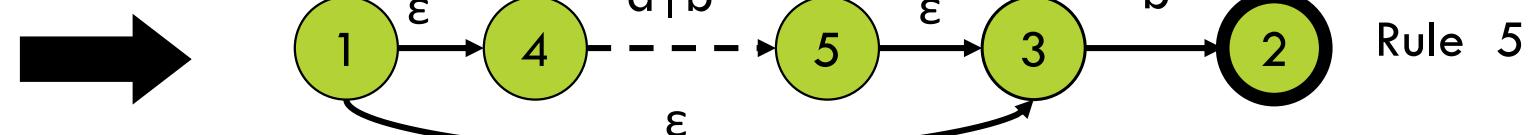
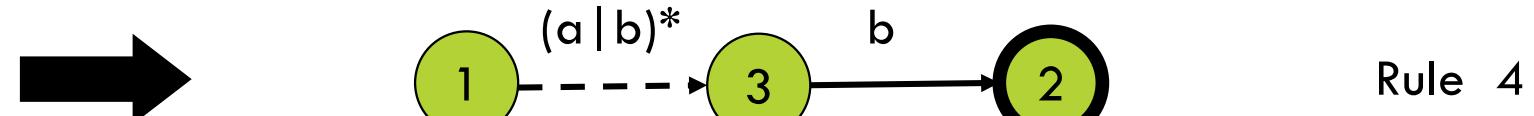
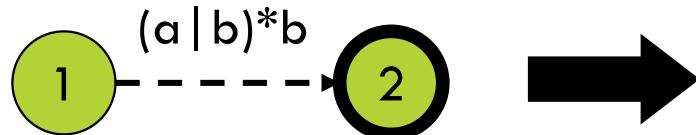


Applying rules for $M^* | N$



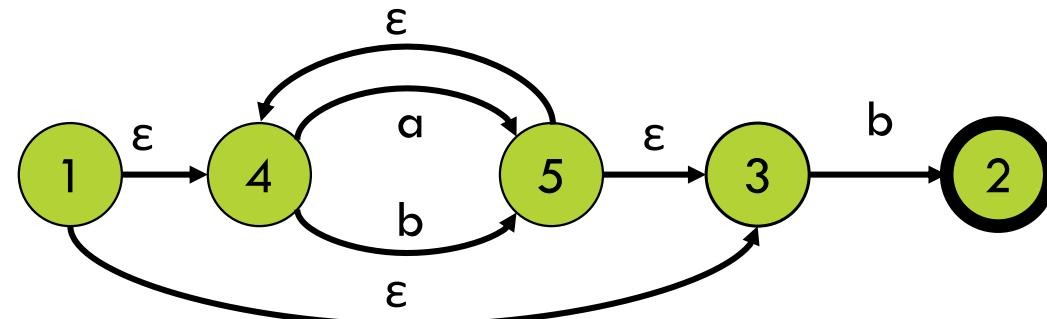
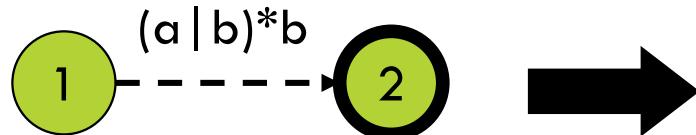
NFA construction: Example

- Construction for $(a \mid b)^*b$

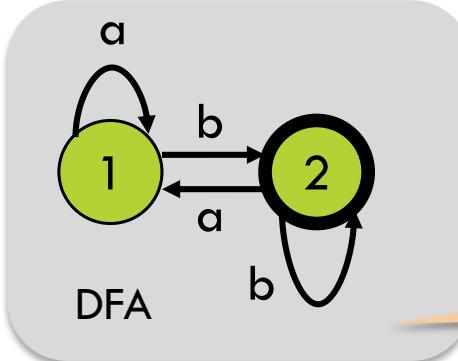
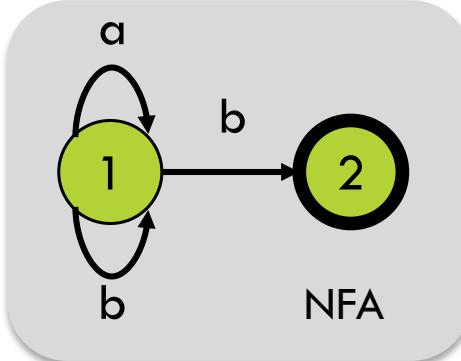


NFA construction: Example (2)

- Construction for $(a \mid b)^*b$



- Equivalents



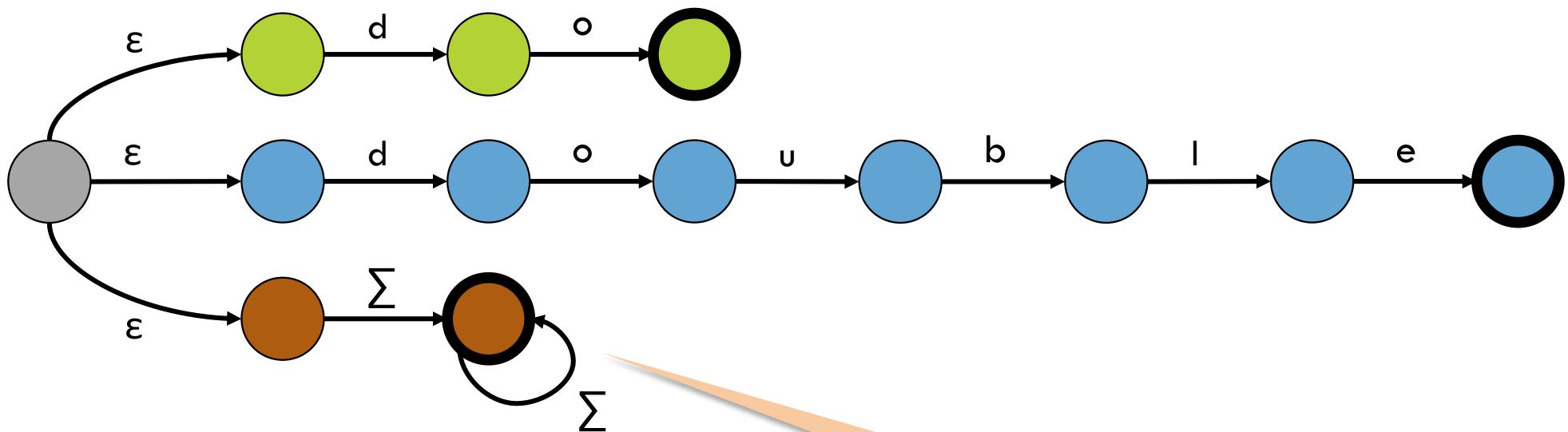
Different options:
Which one to implement?
What is the runtime complexity?
Important: **Automation**

DFA: **Deterministic** finite state
automaton (upcoming)

Scanner operation with NFAs

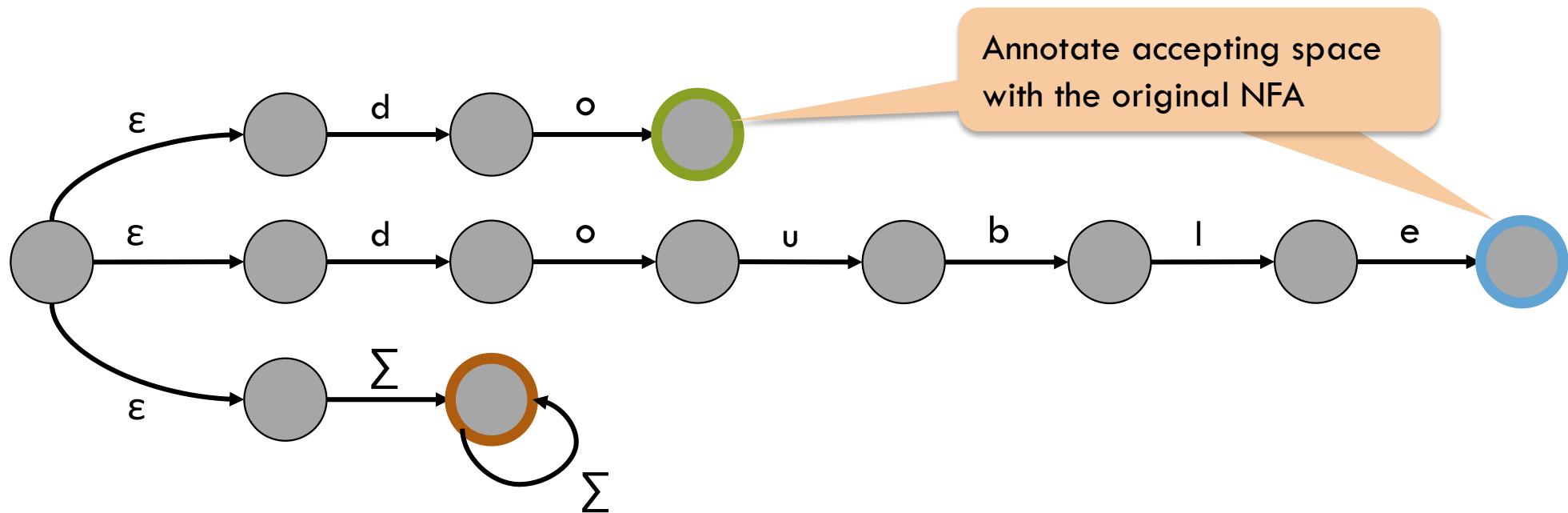
- A scanner = Multiple NFAs running in parallel
- Ambiguities
 - Tie-breaking rule: Maximal munch
 - More ties? – lexeme matches two tokens: Use priorities (first rule wins)
- If nothing matches?
 - Add a catch-up rule that reports an error

Scanner operation with NFAs: Example



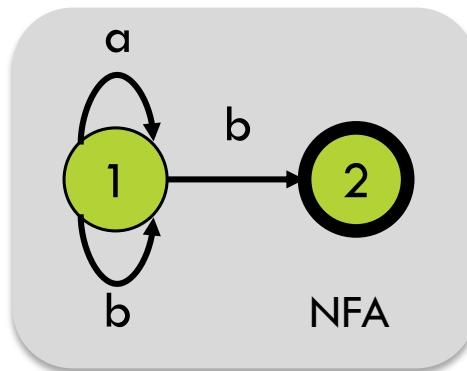
Error-catch rule

Scanner operation with NFAs: Example

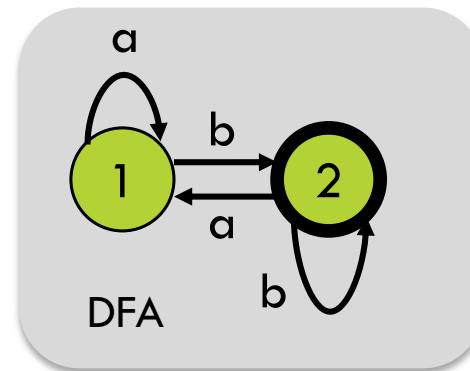


Scanner operation with NFAs: Complexity

- What to do if there are multiple ways to scan the input?
 - Simulate the NFA with n states on string of length m: Complexity $O(mn^2)$
 - Option: Build a **deterministic** finite automaton (DFA)



vs.



Deterministic Finite Automata (DFA)

Def.: A **DFA** is a 5-tuple $(Q, \Sigma, \Delta, q_0, F)$, where Q is a finite set of states, Σ is an alphabet and

1. Δ : Transition **function** $\Delta : Q \times \Sigma \rightarrow Q$
2. q_0 : Initial state ($q_0 \in Q$)
3. F : Set of final states ($F \subseteq Q$)

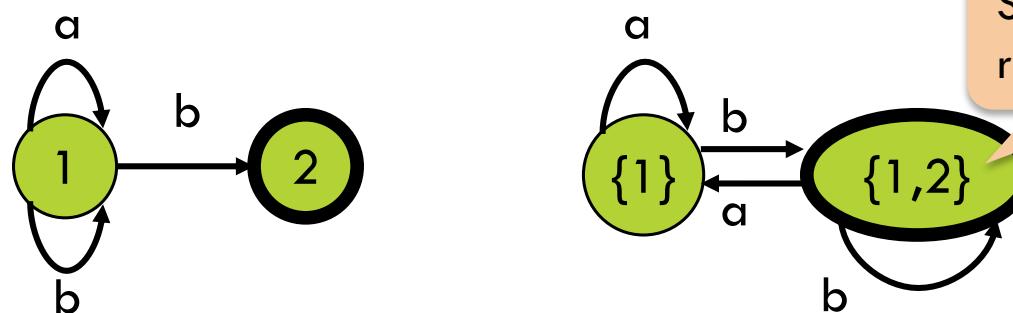
- DFA has no ϵ -moves and Δ is a function, i.e., $\Delta(q, a) = q'$ (no other options)
 - Simulating a DFA has complexity **O(m)**, with m the length of the input string

NFA → DFA: Subset construction

For every NFA Z there is a DFA Y, with $L(Z) = L(Y)$

- Construction idea: Make the DFA simulate the NFA
 - DFA states represent subsets of NFA states
 - Transitions in the DFA represent transitions between subsets in the NFA
 - Problem: For n NFA states potentially 2^n DFA states (power set)
- Example

- $(\{1\}, a) \rightarrow \{1\}$
- $(\{1\}, b) \rightarrow \{1, 2\}$
- $(\{1, 2\}, a) \rightarrow \{1\}$
- $(\{1, 2\}, b) \rightarrow \{1, 2\}$



DFA Construction: ϵ -closure

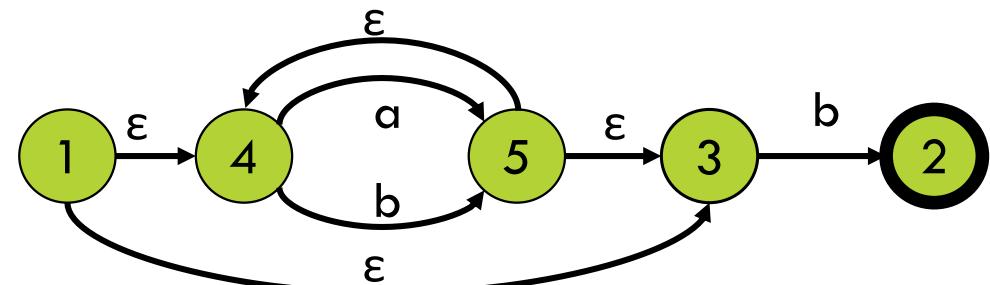
Def.: Given an NFA $Z=(Q,\Sigma,\Delta,q_0,F)$, the ϵ -closure($q \in Q$) is the set of states reachable from q by ϵ -moves, including q itself

Def.: Given an NFA $Z=(Q,\Sigma,\Delta,q_0,F)$ and a subset $S \subseteq Q$, the **ϵ -closure**(S) is:

$$\bigcup_{q \in S} \epsilon\text{-closure}(q)$$

□ Example

- ϵ -closure(1)
 $= \{1,3,4\}$
- ϵ -closure({1,5})
 $= \{1,3,4,5\}$



DFA Construction

Def.: Given an NFA $Z=(Q,\Sigma,\Delta,q_0,F)$, the $\text{move}(q \in Q, a \in \Sigma)$ is the set of states reachable from q when reading symbol ' a ', i.e.,

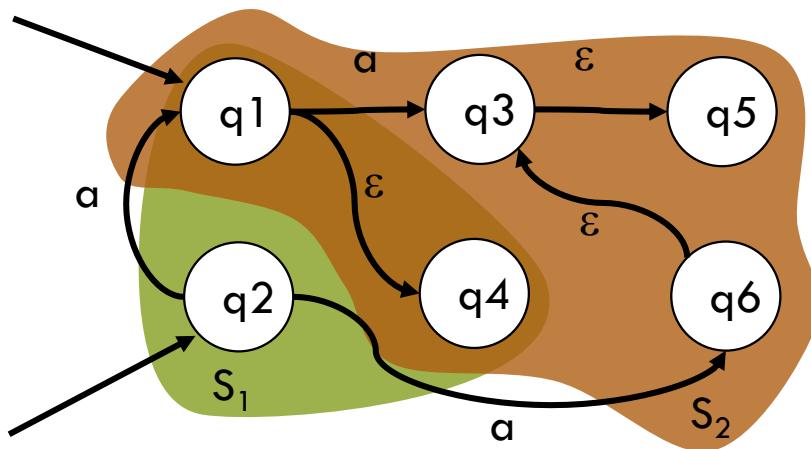
$$\text{move}(q,a) = \{p \mid (q,a,p) \in \Delta\}$$

Def.: Similarly, for a subset $S \subseteq Q$, $\text{move}(S,a)$ is: $\bigcup_{q \in S} \text{move}(q, a)$

- Given an NFA $Z=(Q,\Sigma,\Delta,q_0,F)$ we construct a DFA $Y=(Q',\Sigma, \Delta', q'_0, F')$, where:
 - $Q' \subseteq P(Q)$ – With $P()$ the power set
 - $q'_0 = \epsilon\text{-closure}(q_0)$
 - $F' = \{S \in Q'; S \cap F \neq \emptyset\}$ – The DFA accepts when the NFA could have accepted
 - Δ' $\Delta'(S, a) = \epsilon\text{-closure}(\text{move}(S, a))$

DFA construction: Example

- Consider following portion of an NFA and input ‘a’



$$\text{Move}(S_1, a) = \{q1, q3, q6\}$$

$$\begin{aligned}\varepsilon\text{-closure}(\text{move}(S_1, a)) &= \varepsilon\text{-closure}(\{q1, q3, q6\}) \\ &= \{q1, q3, q4, q5, q6\} = S_2\end{aligned}$$

$$\Delta'(S_1, a) = S_2$$

DFA Construction: Algorithm

Input NFA $Z = (\Sigma, Q, \Delta, q_0, F)$

Output DFA $Y = (\Sigma, Q', \Delta', q'_0, F')$

$q'_0 = \epsilon\text{-closure}(q_0); Q' = \{q'_0\}; \Delta' = \emptyset;$

while (Q' still changing)

foreach $S \in Q'$

foreach $a \in \Sigma$

$S_? = \epsilon\text{-closure}(\text{move}(S, a))$

if $S_? \notin Q'$ **then**

$Q' = Q' \cup \{S_?\}$

$\Delta' = \Delta' \cup \{(S, a) \rightarrow S_?\}$

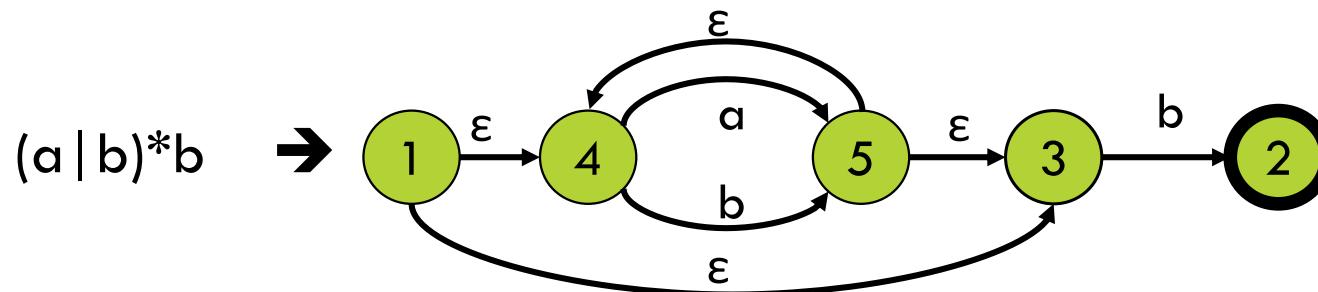
$F' = \{S \in Q'; S \cap F \neq \emptyset\}$

Fixed-point computations:

- Q' contains no duplicates
- Monotone while loop
- $2^{|Q|}$ is finite

Similar algorithms in parsing
and dataflow analysis

DFA construction: Example



Final states?

$$q'0 = \{1,3,4\}$$

$$\epsilon\text{-closure}(\text{Move}(q'0,a)) = \{3,4,5\} \rightarrow q'1 \text{ (new)}$$

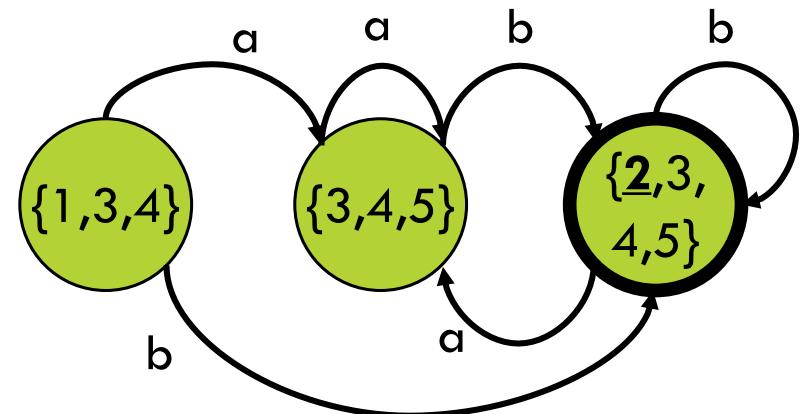
$$\epsilon\text{-closure}(\text{Move}(q'0,b)) = \{2,3,4,5\} \rightarrow q'2 \text{ (new)}$$

$$\epsilon\text{-closure}(\text{Move}(q'1,a)) = \{3,4,5\} \rightarrow q'1$$

$$\epsilon\text{-closure}(\text{Move}(q'1,b)) = \{2,3,4,5\} \rightarrow q'2$$

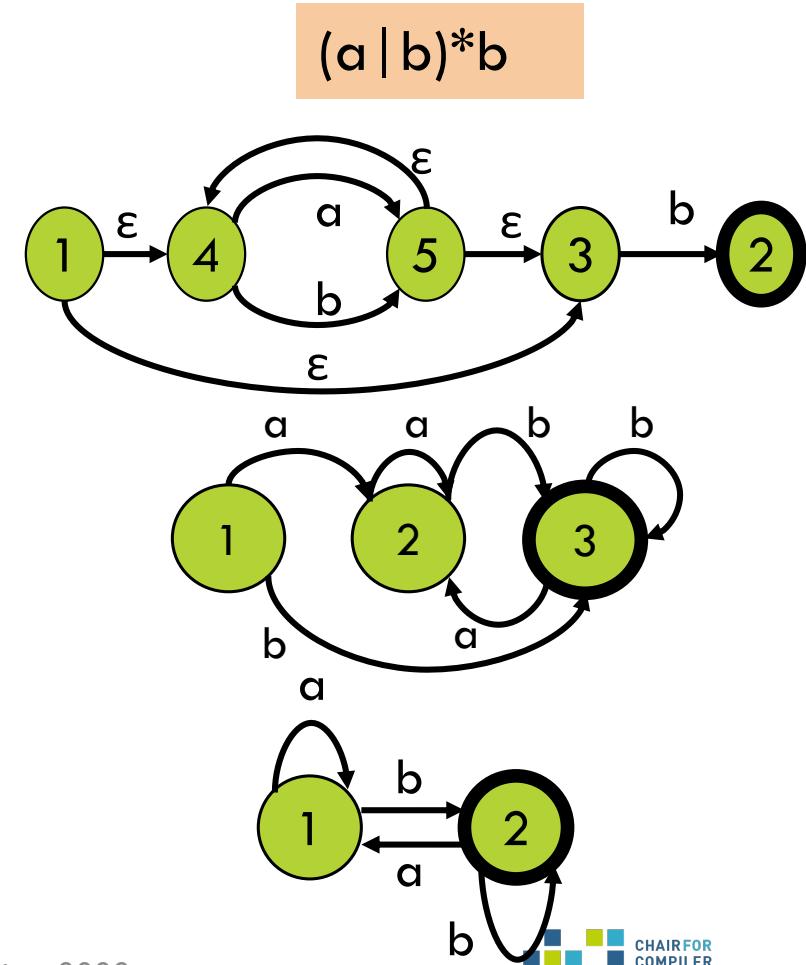
$$\epsilon\text{-closure}(\text{Move}(q'2,a)) = \{3,4,5\} \rightarrow q'1$$

$$\epsilon\text{-closure}(\text{Move}(q'2,b)) = \{2,3,4,5\} \rightarrow q'2$$



DFA minimization

- So far: regexp → NFA (Thompson) → DFA (Powerset construction)
 - Problem: State explosion – Time complexity? – Memory requirements?
- DFA minimization: Hopcroft's algorithm
 - Identify & merge equivalent states
 - Paths leading to states are equivalent
 - Transition on a given symbol 'a' lead to equivalent states
 - Transition on different symbols lead to distinct sets



2. Lexical Analysis

- The scanner
- Languages and regular expressions
- Finite automata
- Scanner generation with Lex



- Lex is the standard scanner generator in unix systems (1975)
 - Input: lex specification
 - Output: Scanner's C-code (produces tokens on demand for the parser)
- Input specification

Definition section

%%

Rules section

%%

C code section

flex: Example

Regexp
 $(0|1\dots|9)(0|1\dots|9)^*$

```
/** Definition section **/  
  
%{  
/* C code to be copied verbatim */  
#include <stdio.h>  
%}  
  
/* This tells flex to read only one input file */  
%option noyywrap  
  
%%  
     /** Rules section **/  
  
     /* [0-9]+ matches a string of one or more digits */  
[0-9]+ {  
    /* yytext is a string containing the matched text. */  
    printf("Saw an integer: %s\n", yytext);  
}  
  
.|\n    { /* Ignore all other characters. */ }  
  
%%  
/** C Code section **/  
  
int main(void)  
{  
    /* Call the lexer, then quit. */  
    yylex();  
    return 0;  
}
```

Action: Executed in case of detection

Source: wikipedia.org

Where are we?

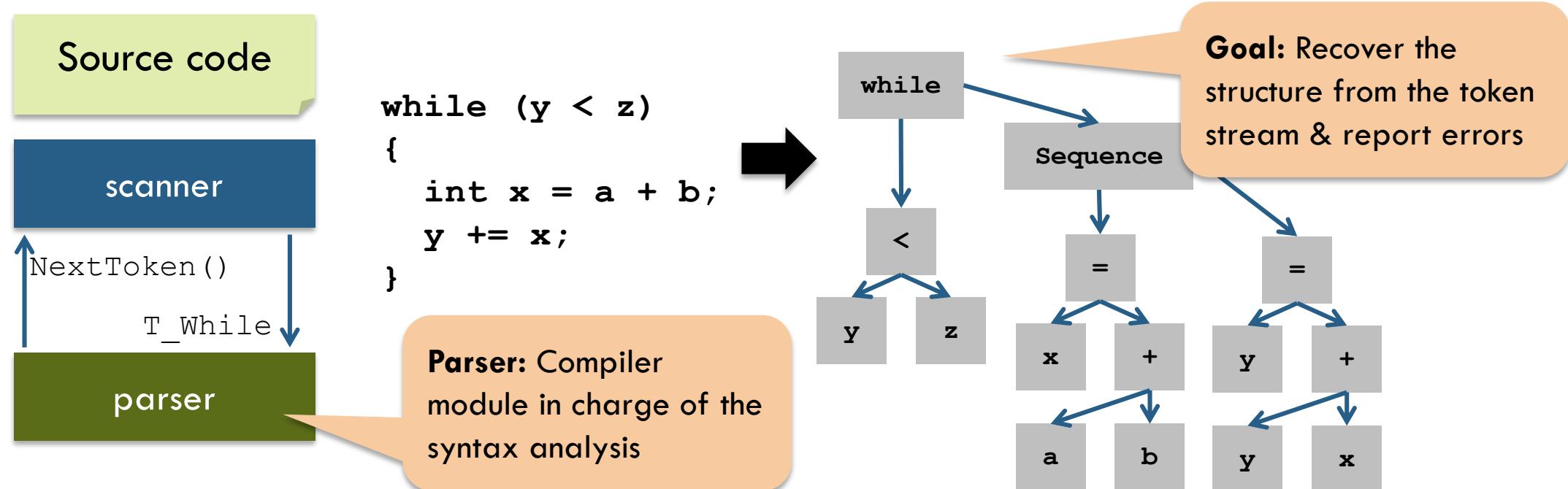
1. Introduction
2. Lexical analysis
3. Syntax analysis
4. Semantic analysis
5. Intermediate representation
6. Control & data-flow analysis
7. IR optimization
8. Target architectures
9. Code selection
10. Register allocation
11. Scheduling
12. Advanced topics

3. Syntax Analysis

- Introduction
- Context free grammars
- Top-down parser
- Bottom-up parser
- Error handling
- Parser generation with yacc



The parser: Recall



Why not use regexp and DFAs?

- Usually programming languages aren't regular languages
 - Balanced parenthesis & nested blocks in functions
- Example: $L = \{a^n b^n : n > 0\}$, regular?
 - If it is regular, then there is a regexp → and a DFA, with k states (finite)
 - For an input with $n > k$, at least one repeated state while reading a^n
 - Let repeated state be s , after reading a^p and a^q
 - The DFA should accept $a^p b^p$ and $a^q b^q$
 - Consequently, the DFA would also accept $a^q b^p$ and $a^p b^q$
- (Alternative proof: By contradiction to the **pumping lemma**)

→ Need more powerful formalism: **Context-free grammars**

Terminology

	Scanner	Parser
Alphabet	Characters, e.g., ASCII	Tokens, e.g., ID, numbers, keywords
String	Character string, e.g., keyword	Entire program
Language	Subset of strings, e.g., floating numbers	Programming language
Formalism	Regular expressions	Context-free grammars
Recognition	DFA	Parsers, e.g., LL(k) & LR(k)

3. Syntax Analysis

- Introduction
- Context free grammars
- Top-down parser
- Bottom-up parser
- Error handling
- Parser generation with yacc



Context-free grammar

Def.: A **context-free grammar (CFG)** is a 4-tuple (N, T, P, S) , where

1. N : **Alphabet** of non-terminal symbols (help-symbols)
2. T : **Alphabet** of terminal symbols (of the language defined by the CFG), $N \cap T = \emptyset$
3. $P \subseteq N \times (N \cup T)^*$: Set **rules** or **productions**
4. $S \in N$: Is the start symbol

□ Notation

- A, B, C, \dots : Non-terminals
- a, b, c, \dots : Terminals
- $\alpha, \beta, \gamma, \delta, \dots$: Strings of $(N \cup T)^*$
- Productions: A relation element $(A, \alpha) \in P$, is denoted $A \rightarrow \alpha$

Derivation

Def.: For any strings $\alpha, \beta \in (N \cup T)^*$, β can be **directly derived** from α , denoted $\alpha \rightarrow \beta$, if there are $\gamma_1, \gamma_2 \in (N \cup T)^*$ and $(A, \gamma_3) \in P$ such that $\alpha = \gamma_1 A \gamma_2$ and $\beta = \gamma_1 \gamma_3 \gamma_2$

Def.: A sequence of derivations $\alpha_1 \rightarrow \alpha_2 \rightarrow \alpha_3 \rightarrow \dots \rightarrow \alpha_n$ of zero or more steps is denoted $\alpha_1 \rightarrow^* \alpha_n$ (transitive closure of \rightarrow)

Def.: For a CFG G , a string $\alpha \in (N \cup T)^*$ is called a **sentential form** iff $S \rightarrow^* \alpha$

Deutsch: **Satzform**

Example: Chemistry

Form → Cmp | Cmp Ion

Cmp → Term | Term Num | Cmp Cmp

Term → Elem | (Cmp)

Elem → H | He | Li | Be | B | C | ...

Ion → + | - | IonNum + | IonNum -

IonNum → 2 | 3 | 4 | ...

Num → 1 | IonNum

P

$G = (N, T, P, S)$

$N = \{Form, Cmp, Ion, Term, Num, IonNum\}$

$T = \{H, He, \dots, +, -, 1, 2, 3, \dots\}$

$S = Form$

“Form → Cmp | Cmp Ion” is a shortcut for:

Form → Cmp

Form → Cmp Ion

Adapted from:

<http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Example: Chemistry

Form → Cmp | Cmp Ion

Cmp → Term | Term Num | Cmp Cmp

Term → Elem | (Cmp)

Elem → H | He | Li | Be | B | C | ...

Ion → + | - | IonNum + | IonNum -

IonNum → 2 | 3 | 4 | ...

Num → 1 | IonNum

P

Sample derivation

Form

⇒ Cmp Ion

⇒ Cmp Cmp Ion

⇒ Cmp Term Num Ion

⇒ Term Term Num Ion

⇒ Elem Term Num Ion

⇒ Mn Term Num Ion

⇒ Mn Elel Num Ion

⇒ MnO Num Ion

⇒ MnO IonNum Ion

⇒ MnO₄ Ion

⇒ MnO₄⁻

Adapted from:

<http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Example: Programming languages

BLOCK → **STMT**
| { **STMTS** }

STMTS → ϵ
| **STMT STMTS**

STMT → **EXPR ;**
| **if (EXPR) BLOCK**
| **while (EXPR) BLOCK**
| **do BLOCK while (EXPR) ;**
| **BLOCK**
| ...

EXPR → **identifier**
| **constant**
| **EXPR + EXPR**
| **EXPR - EXPR**
| **EXPR * EXPR**
| ...

Adapted from:
<http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Context-free languages

Def.: Given a CFG $G = (N, T, P, S)$ the associated context-free **language** is defined by $L(G) = \{w \in T^*; S \xrightarrow{*} w\}$

Def.: A language L is called **context-free** if there exists a CFG G , so that $L = L(G)$

- All regular languages are also context-free languages, i.e., for every regexp R there is a CFG G so that $L(R) = L(G)$
- Simple rules
 - $(a \mid b) \rightarrow A \rightarrow a; A \rightarrow b$
 - $a^* \rightarrow A \rightarrow aA; A \rightarrow \epsilon$
- Example: $(a \mid b)^*b$
 - $S \rightarrow AB$
 - $A \rightarrow CA \mid \epsilon$
 - $C \rightarrow a$
 - $C \rightarrow b$
 - $B \rightarrow b$

Derivations: Right and leftmost

Def. A **leftmost derivation** is a derivation in which each step expands the **leftmost** non-terminal (analogous for a rightmost derivation)

□ Example: Expressions

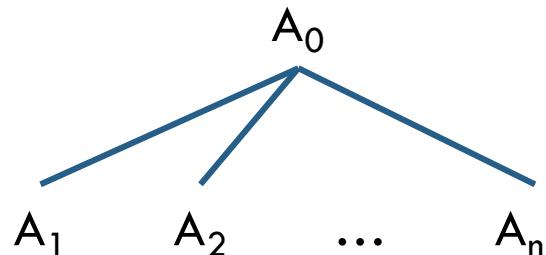
- (1) $E \rightarrow E + E$
- (2) $E \rightarrow E * E$
- (3) $E \rightarrow (E)$
- (4) $E \rightarrow -E$
- (5) $E \rightarrow id$

- Leftmost
 - $E \rightarrow -E \rightarrow -(E) \rightarrow -(E+E) \rightarrow -(id + E) \rightarrow -(id + id)$
- Rightmost
 - $E \rightarrow -E \rightarrow -(E) \rightarrow -(E+E) \rightarrow -(E + id) \rightarrow -(id + id)$

Parse tree

Def.: A **parse tree (concrete syntax tree)** encodes the steps of a derivation. The root node represents the start-symbol, internal nodes represent non-terminal symbols and leaves represent terminal symbols

- In-order left-to-right traversal of the leaves represent the program
- Parse trees are invariant to leftmost and rightmost derivation
- When applying a production $A_0 \rightarrow A_1A_2..A_n$, with $A_0 \in N$, $A_{i:1,...,n} \in (N \cup T)$, the children $A_{i:1,...,n}$ of node A_0 are placed from left to right



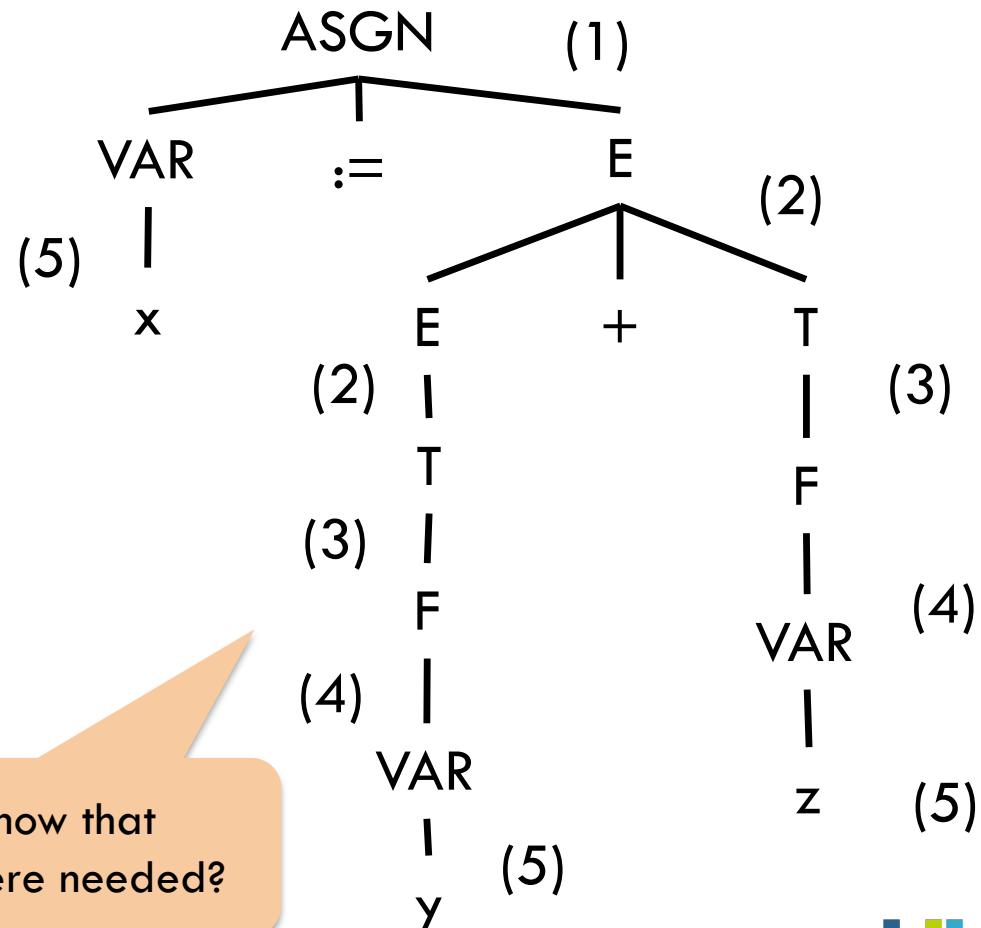
Parse tree: Example

Given the rules

- (1) ASGN \rightarrow VAR := E
- (2) E \rightarrow E + T | T
- (3) T \rightarrow T x F | F
- (4) F \rightarrow (E) | VAR
- (5) VAR \rightarrow x | y | z

Parse tree for: x := y + z

How did we know that
these steps were needed?



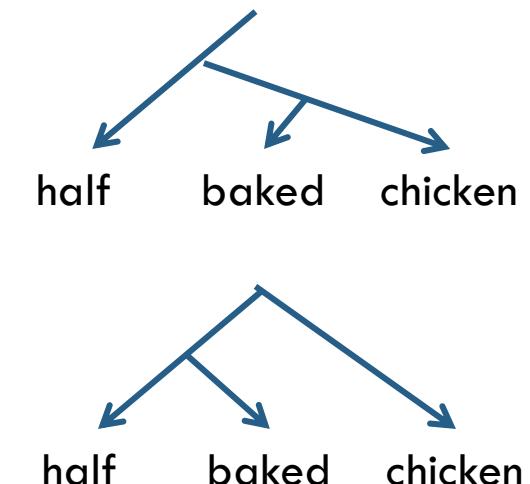
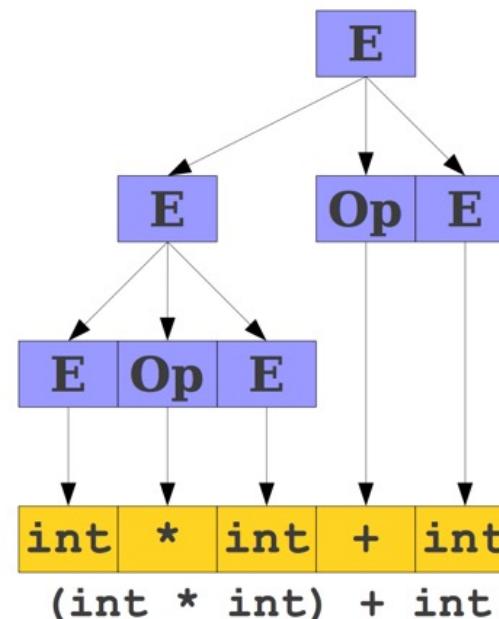
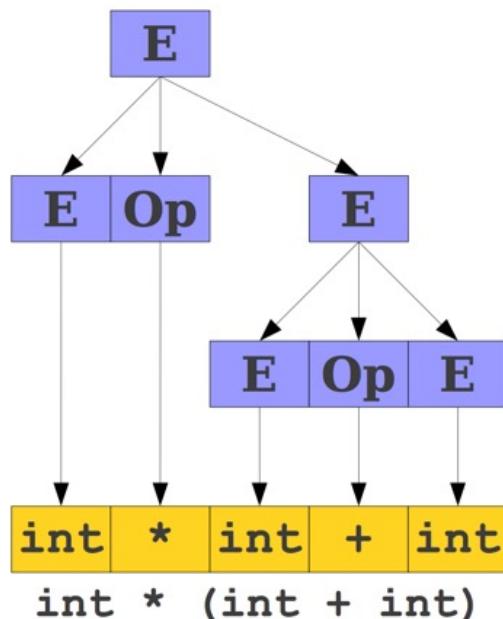
Parse challenge: Ambiguity

Def.: A CFG G is called ambiguous if there is at least one word w in $L(G)$, i.e.,
 $S \xrightarrow{*} w$, with two or more parse trees

- Ambiguity is a property of the grammar
 - A language is **inherently ambiguous** if there is no unambiguous grammar
- Detecting ambiguity: Undecidable

Parse challenge: Ambiguity (2)

- Why is this a problem? – different semantics → How to translate then?
 - Examples



Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Resolving ambiguity

- Layering: Break derivation into different layers in the tree
 - Enforce precedence: $\text{id} + \text{id} \times \text{id}$

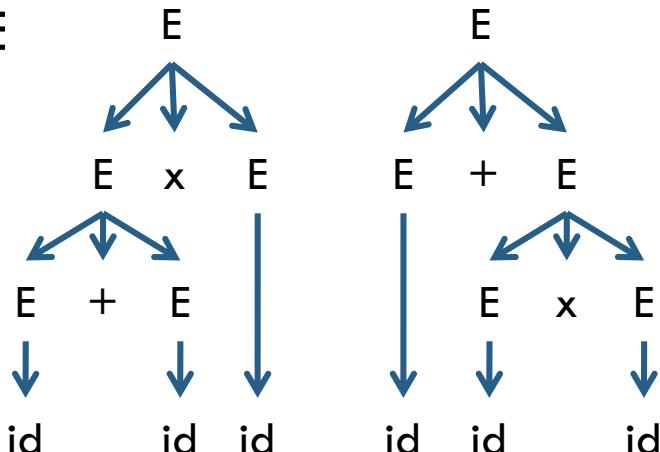
$$(1) E \rightarrow E + E$$

$$(2) E \rightarrow E \times E$$

$$(3) E \rightarrow (\text{id})$$

$$(4) E \rightarrow -\text{id}$$

$$(5) E \rightarrow \text{id}$$



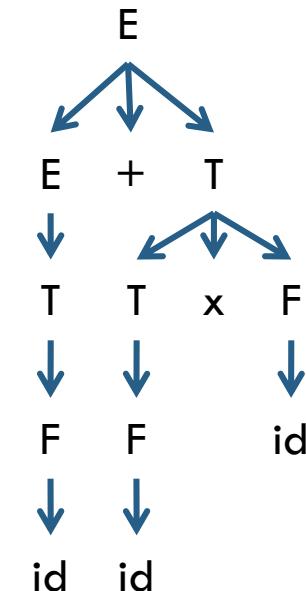
$$(1) E \rightarrow E + T$$

$$(2) E \rightarrow T$$

$$(3) T \rightarrow T \times F$$

$$(4) T \rightarrow F$$

$$(5) F \rightarrow \text{id}$$



3. Syntax Analysis

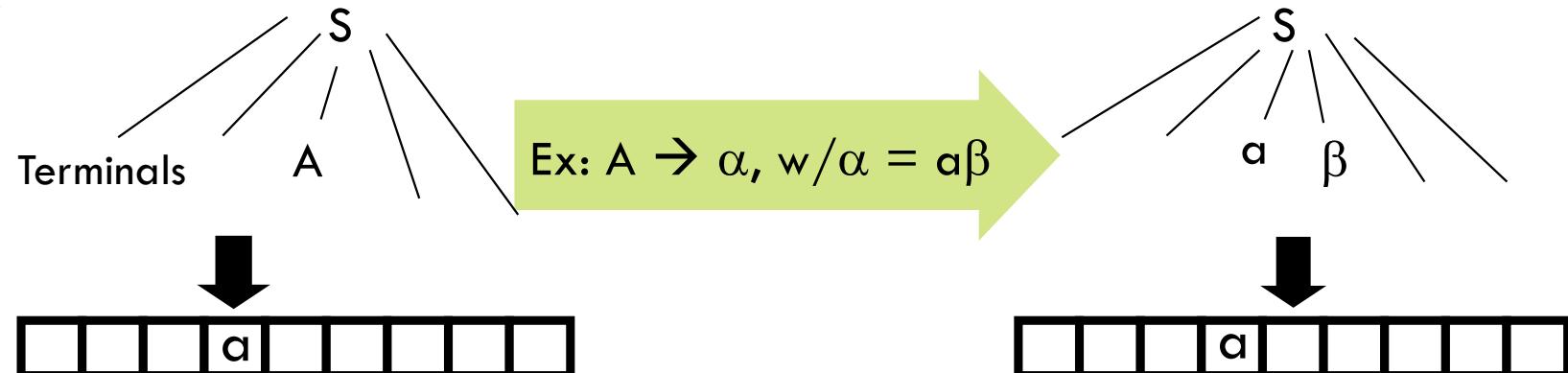
- Introduction
- Context free grammars
- Top-down parser
- Bottom-up parser
- Error handling
- Parser generation with yacc



Top-down parser

- Recall – Parser goal: Given a program find a parse tree (or report error)
- Overview:
 - Start with the start symbol (parse tree root)
 - Grow towards the leaves using leftmost derivation, picking rules to match the input
 - Bad selection: backtrack (expensive)
 - Some grammars are backtrack-free (later: predictive parsing)

Graphically



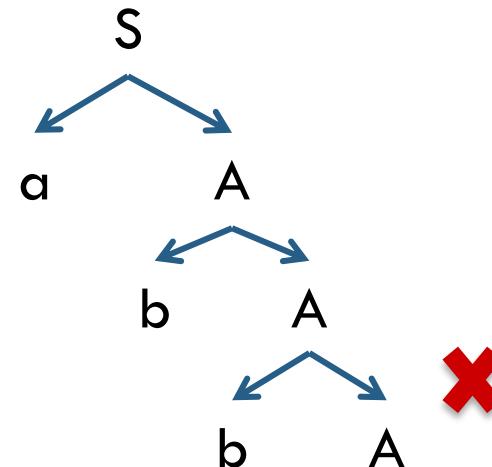
Top-down parser: Example

□ Rules

- $S \rightarrow aA$
- $A \rightarrow bA$
- $A \rightarrow bc$

□ Input: **abbc**

- Need a: $S \rightarrow aA$
- Need b: $A \rightarrow bA$
- Need b: $A \rightarrow bA$
- Need c: Impossible → backtrack



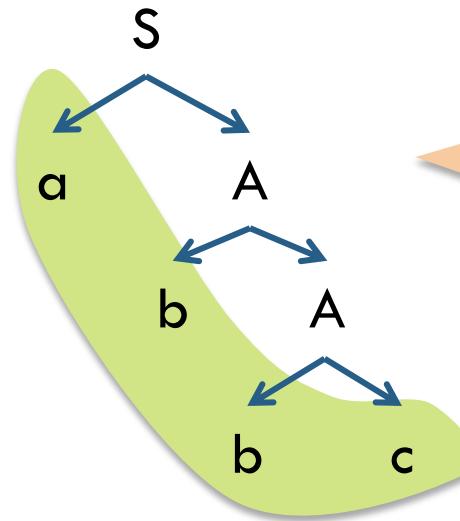
Top-down parser: Example

□ Rules

- $S \rightarrow aA$
- $A \rightarrow bA$
- $A \rightarrow bc$

□ Input: **abbc**

- Need a: $S \rightarrow aA$
- Need b: $A \rightarrow bA$
- Need b: $A \rightarrow bA$
- Need c: Impossible → backtrack
- Take $A \rightarrow bc$ (done!)



What kinds of grammars would be problematic?

Problem 1: Left recursion

Def.: A CFG G is left recursive if there is a non-terminal symbol A and a sentential $\alpha \in (N \cup T)^*$ so that $A \xrightarrow{*} A\alpha$ (analogous for right recursion)

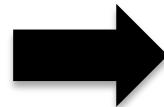
- Top-down parsers cannot handle left recursion → leads to non-termination!
- Solution: Grammar transformation into right recursive

$$\begin{array}{ll} (1) & E \rightarrow E\alpha \\ (2) & E \rightarrow \beta \end{array} \quad \longrightarrow \quad \begin{array}{l} (1) \quad E \rightarrow \beta T \\ (2) \quad T \rightarrow \alpha T \\ (3) \quad T \rightarrow \epsilon \end{array}$$

Problem 2: Multiple options

- Having multiple rules for a non-terminal in the partial parse tree may lead to backtracking
- Solution: Factorization

$$\begin{array}{l} (1) \quad E \rightarrow AB \\ (2) \quad E \rightarrow ABC \end{array}$$



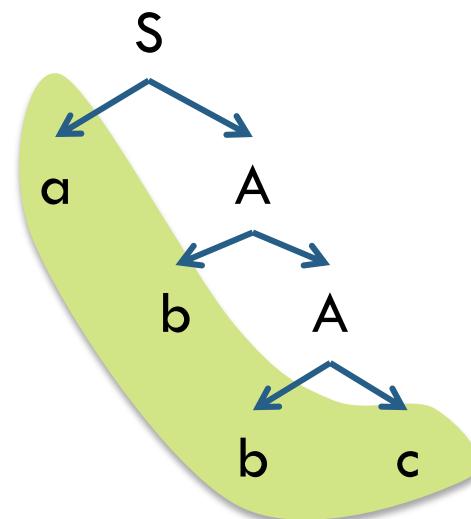
$$\begin{array}{l} (1) \quad E \rightarrow ABX \\ (2) \quad X \rightarrow C \\ (3) \quad X \rightarrow \epsilon \end{array}$$

Predictive parsing (recursive decent parser)

- Idea: Use the input to better predict what rule to use and prevent backtracking
 - Use **look-ahead tokens** from the input stream
 - The more tokens, the better
- Tradeoff: Easy parser construction, but low expressiveness
- Example with **two** look-ahead tokens

Input **abbc**

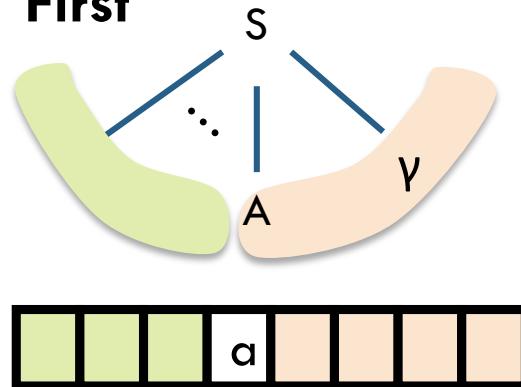
- $S \rightarrow aA$ → Input: ab: $S \rightarrow aA$
- $A \rightarrow bA$ → Input: bb: $A \rightarrow bA$
- $A \rightarrow bc$ → Input: bc: $A \rightarrow bc$



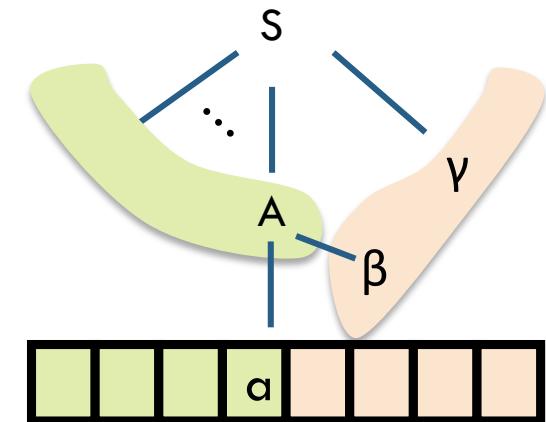
First and Follow sets: Intuition

- When to apply a rule $A \rightarrow \alpha$ when looking at a token 'a'?

1. First

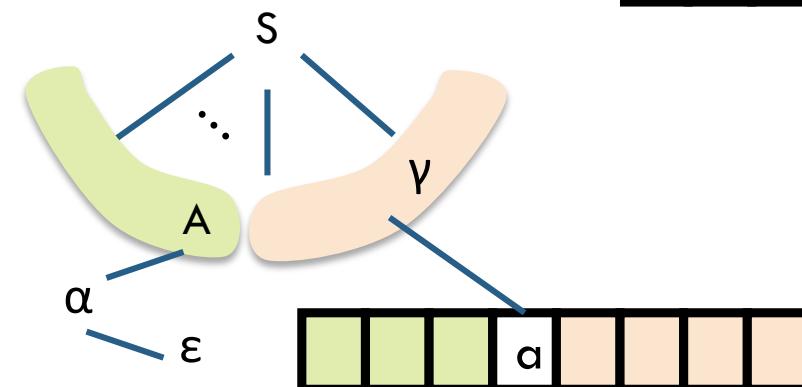


If α can **start** with 'a', e.g.,
 $A \rightarrow a\beta$, then use rule and
derive the rest from $\beta\gamma$



2. Follow

If $\alpha \rightarrow^* \epsilon$ and $S \rightarrow^* \beta A \alpha \gamma$, then
use rule and derive the rest
from $\alpha \gamma$



First and Follow sets: Definition

Def.: Given a CFG $G = (N, T, P, S)$ and a string $\alpha \in (N \cup T)^*$ the **first set** is defined as $\text{FIRST}(\alpha) = \{\alpha \in T; \alpha \xrightarrow{*} \alpha\beta, \beta \in (N \cup T)^*\} \cup \{\varepsilon \mid \alpha \xrightarrow{*} \varepsilon\}$

Def.: Given a CFG $G = (N, T, P, S)$ and non-terminal $A \in N$ the **follow set** is $\text{FOLLOW}(A) = \{\alpha \in T; S \xrightarrow{*} \beta A \alpha \gamma, \beta, \gamma \in (N \cup T)^*\}$

First and Follow sets: Example

$$\text{FIRST}(\alpha) = \{a \in T; \alpha \xrightarrow{*} a\beta, \beta \in (N \cup T)^*\} \cup \{\epsilon \mid \alpha \xrightarrow{*} \epsilon\}$$
$$\text{FOLLOW}(A) = \{a \in T; S \xrightarrow{*} \beta A a \gamma, \beta, \gamma \in (N \cup T)^*\}$$

$$S \rightarrow uBDz$$

$$\text{FIRST}(uBDz) = \{ u \}$$

$$\text{FOLLOW}(D) = \{ z \}$$

$$B \rightarrow Bv$$

$$\text{FIRST}(Bv) = \{ w \}$$

$$\text{FOLLOW}(B) = ??$$

$$B \rightarrow w$$

$$\text{FIRST}(w) = \{ w \}$$

$$D \rightarrow EF$$

$$\text{FIRST}(EF) = \{ \epsilon, x, y \}$$

Easy in "Bv", but in "uBDz"
 $\text{FOLLOW}(B) \supset \text{FIRST}(Dz)-\{\epsilon\}$

$$E \rightarrow y$$

$$\text{FIRST}(y) = \{ y \}$$

$$E \rightarrow \epsilon$$

$$\text{FIRST}(x) = \{ x \}$$

$$F \rightarrow x$$

FIRST(α) for all right-hand sides

$$F \rightarrow \epsilon$$

First and Follow sets: Example

$$\text{FIRST}(\alpha) = \{\alpha \in T; \alpha \xrightarrow{*} \alpha\beta, \beta \in (N \cup T)^*\} \cup \{\epsilon \mid \alpha \xrightarrow{*} \epsilon\}$$
$$\text{FOLLOW}(A) = \{\alpha \in T; S \xrightarrow{*} \beta A \alpha \gamma, \beta, \gamma \in (N \cup T)^*\}$$

$$S \rightarrow uBDz$$

$$\text{FIRST}(uBDz) = \{ u \}$$

$$\text{FOLLOW}(D) = \{ z \}$$

$$B \rightarrow Bv$$

$$\text{FIRST}(Bv) = \{ w \}$$

$$\text{FOLLOW}(B) = ??$$

$$B \rightarrow w$$

$$\text{FIRST}(EF) = \{ \epsilon, x, y \}$$

$$D \rightarrow EF$$

$$\text{FIRST}(\alpha \in T) = \{ \alpha \}$$

Easy in "Bv", but in "uBDz"
 $\text{FOLLOW}(B) \supset \text{FIRST}(Dz)-\{\epsilon\}$

$$E \rightarrow y$$

$$\text{FIRST}(D) = \{ \epsilon, x, y \}$$

$$E \rightarrow \epsilon$$

$$\text{FIRST}(B) = \{ w \}$$

$$F \rightarrow x$$

$$\text{FIRST}(E) = \{ \epsilon, y \}$$

$$F \rightarrow \epsilon$$

$$\text{FIRST}(F) = \{ \epsilon, x \}$$

$$\text{FIRST}(D) = \text{FIRST}(EF)$$

First and Follow sets: Example

$$\text{FIRST}(\alpha) = \{\alpha \in T; \alpha \xrightarrow{*} \alpha\beta, \beta \in (N \cup T)^*\} \cup \{\epsilon \mid \alpha \xrightarrow{*} \epsilon\}$$
$$\text{FOLLOW}(A) = \{\alpha \in T; S \xrightarrow{*} \beta A \alpha \gamma, \beta, \gamma \in (N \cup T)^*\}$$

$$S \rightarrow uBDz$$

$$\text{FIRST}(uBDz) = \{ u \}$$

$$\text{FOLLOW}(D) = \{ z \}$$

$$B \rightarrow Bv$$

$$\text{FIRST}(Bv) = \{ w \}$$

$$\text{FOLLOW}(B) = \{ v, x, y, z \}$$

$$B \rightarrow w$$

$$\text{FIRST}(EF) = \{ \epsilon, x, y \}$$

$$D \rightarrow EF$$

$$\text{FIRST}(\alpha \in T) = \{ \alpha \}$$

Easy in "Bv", but in "uBDz"
 $\text{FOLLOW}(B) \supset \text{FIRST}(Dz)-\{\epsilon\}$

$$E \rightarrow y$$

$$\text{FIRST}(D) = \{ \epsilon, x, y \}$$

$$E \rightarrow \epsilon$$

$$\text{FIRST}(B) = \{ w \}$$

$$F \rightarrow x$$

$$\text{FIRST}(E) = \{ \epsilon, y \}$$

$$F \rightarrow \epsilon$$

$$\text{FIRST}(F) = \{ \epsilon, x \}$$

$\text{FIRST}(D) = \text{FIRST}(EF)$

First and Follow sets: Example

$$\text{FIRST}(\alpha) = \{\alpha \in T; \alpha \xrightarrow{*} \alpha\beta, \beta \in (N \cup T)^*\} \cup \{\epsilon \mid \alpha \xrightarrow{*} \epsilon\}$$
$$\text{FOLLOW}(A) = \{\alpha \in T; S \xrightarrow{*} \beta A \alpha \gamma, \beta, \gamma \in (N \cup T)^*\}$$

$S \rightarrow uBDz$	$\text{FIRST}(uBDz)$	$= \{ u \}$	$\text{FOLLOW}(D) = \{ z \}$
$B \rightarrow Bv$	$\text{FIRST}(Bv)$	$= \{ w \}$	$\text{FOLLOW}(B) = \{ v, x, y, z \}$
$B \rightarrow w$	$\text{FIRST}(EF)$	$= \{ \epsilon, x, y \}$	$\text{FOLLOW}(E) =$
$D \rightarrow EF$			$\text{FIRST}(F)-\{\epsilon\} \cup \text{FOLLOW}(D)$
$E \rightarrow y$	$\text{FIRST}(\alpha \in T)$	$= \{ \alpha \}$	Since $F \xrightarrow{*} \epsilon$
$E \rightarrow \epsilon$	$\text{FIRST}(D)$	$= \{ \epsilon, x, y \}$	
$F \rightarrow x$	$\text{FIRST}(B)$	$= \{ w \}$	
$F \rightarrow \epsilon$	$\text{FIRST}(E)$	$= \{ \epsilon, y \}$	
	$\text{FIRST}(F)$	$= \{ \epsilon, x \}$	

First and Follow sets: Example

$$\text{FIRST}(\alpha) = \{\alpha \in T; \alpha \xrightarrow{*} \alpha\beta, \beta \in (N \cup T)^*\} \cup \{\epsilon \mid \alpha \xrightarrow{*} \epsilon\}$$
$$\text{FOLLOW}(A) = \{\alpha \in T; S \xrightarrow{*} \beta A \alpha \gamma, \beta, \gamma \in (N \cup T)^*\}$$

$S \rightarrow uBDz$	$\text{FIRST}(uBDz)$	$= \{ u \}$	$\text{FOLLOW}(D) = \{ z \}$
$B \rightarrow Bv$	$\text{FIRST}(Bv)$	$= \{ w \}$	$\text{FOLLOW}(B) = \{ v, x, y, z \}$
$B \rightarrow w$	$\text{FIRST}(EF)$	$= \{ \epsilon, x, y \}$	$\text{FOLLOW}(E) = \{ x, z \}$
$D \rightarrow EF$			$\text{FOLLOW}(F) = \{ z \}$
$E \rightarrow y$	$\text{FIRST}(\alpha \in T)$	$= \{ \alpha \}$	$\text{FOLLOW}(S) = \{ \$ \}$
$E \rightarrow \epsilon$	$\text{FIRST}(D)$	$= \{ \epsilon, x, y \}$	
$F \rightarrow x$	$\text{FIRST}(B)$	$= \{ w \}$	
$F \rightarrow \epsilon$	$\text{FIRST}(E)$	$= \{ \epsilon, y \}$	
	$\text{FIRST}(F)$	$= \{ \epsilon, x \}$	

Help symbol: End of string (EOF)

First and Follow sets: Systematic

- **FIRST(X):** for $X \in T$, $\text{FIRST}(X) = \{X\}$, for $X \in N$ and $X \rightarrow \alpha$:
 - If $\alpha \xrightarrow{*} \varepsilon$, then $\varepsilon \in \text{FIRST}(X)$
 - If $\alpha = Y_1Y_2\dots Y_n$ and there is a $k \in \{1, \dots, n-1\}$ for which $Y_i \xrightarrow{*} \varepsilon$ for all $i \in \{1, \dots, k\}$,
then $(\text{FIRST}(Y_j) - \{\varepsilon\}) \subset \text{FIRST}(X)$, for $j \in \{1, \dots, k+1\}$
- **FOLLOW(X):** for $X \in N$,
 - Look at every rule $A \rightarrow \alpha X \beta$ where A is reachable from S:
 - $(\text{FIRST}(\beta) - \{\varepsilon\}) \subset \text{FOLLOW}(X)$
 - If $\beta \xrightarrow{*} \varepsilon$, then $\text{FOLLOW}(A) \subset \text{FOLLOW}(X)$
 - $\$ \in \text{FOLLOW}(S)$

Table-based predictive parser

- Given a CFG $G = (N, T, P, S)$, build a table $M[A \in N, a \in T] = r \in P$
 - Given the leftmost non-terminal A and look-ahead ‘ a ’ use rule ‘ r ’
- Construction: Entry for $M[A, a]$
 - Entry $r: A \rightarrow \alpha$ if $a \in \text{FIRST}(\alpha)$
 - Entry $r: A \rightarrow \alpha$ with $\alpha \xrightarrow{*} \varepsilon$ if $a \in \text{FOLLOW}(A)$
 - Entry is error otherwise

Def.: A so constructed parser is called **LL(1) parser** – it parses the input from Left to right, performing a Leftmost derivation with 1 look-ahead token

Table-based predictive parser: Example

$S \rightarrow uBDz$
 $B \rightarrow Bv \mid w$
 $D \rightarrow EF$
 $E \rightarrow y \mid \epsilon$
 $F \rightarrow x \mid \epsilon$

$\text{FIRST}(uBDz)$	$= \{ u \}$
$\text{FIRST}(Bv)$	$= \{ w \}$
$\text{FIRST}(EF)/\text{FIRST}(D)$	$= \{ \epsilon, x, y \}$
$\text{FIRST}(\alpha \in T)$	$= \{ \alpha \}$
$\text{FOLLOW}(D)$	$= \{ z \}$
$\text{FOLLOW}(E)$	$= \{ x, z \}$
$\text{FOLLOW}(F)$	$= \{ z \}$
$\text{FOLLOW}(B)$	$= \{ x, y, v, z \}$
$\text{FOLLOW}(S)$	$= \{ \$ \}$

M	u	v	w	x	y	z
D				$D \rightarrow EF$	$D \rightarrow EF$	$D \rightarrow EF$
E					$E \rightarrow y$	$E \rightarrow \epsilon$
F				$F \rightarrow x$		$F \rightarrow \epsilon$
B			$B \rightarrow Bv$ $B \rightarrow w$			
S	$S \rightarrow uBDz$					

Table-based predictive parser: Example

$S \rightarrow uBDz$
 $B \rightarrow Bv \mid w$
 $D \rightarrow EF$
 $E \rightarrow y \mid \epsilon$
 $F \rightarrow x \mid \epsilon$

LL(1) parser not enough
 for this grammar!
 Reason: Left-recursion
 “Homework”: Rewrite
 grammar and build
 table again

M	u	v	w	x	y	z
D				$D \rightarrow EF$	$D \rightarrow EF$	$D \rightarrow EF$
E				$E \rightarrow \epsilon$	$E \rightarrow y$	$E \rightarrow \epsilon$
F				$F \rightarrow x$		$F \rightarrow \epsilon$
B			$B \rightarrow Bv$ $B \rightarrow w$			
S	$S \rightarrow uBDz$					

Def.: A CFG G is in the class LL(1), denoted $G \in \text{LL}(1)$ if the parse table has at most one entry per cell

- The concept of LL(1) can be extended to k look-ahead tokens: LL(k)
 - The table grows exponentially with k ($|T|^k$)
 - The more tokens, the more expressive the grammar, i.e.,

$\text{LL}(0) \subset \text{LL}(1) \subset \dots \subset \text{LL}(k) \dots \subset \text{unambiguous grammars} \subset \text{all grammars}$

What kind of language would this be?

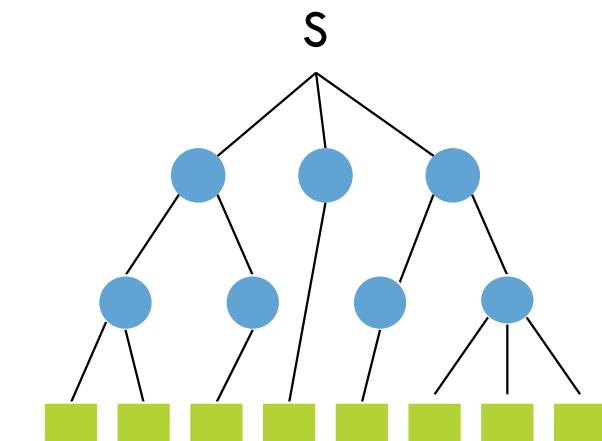
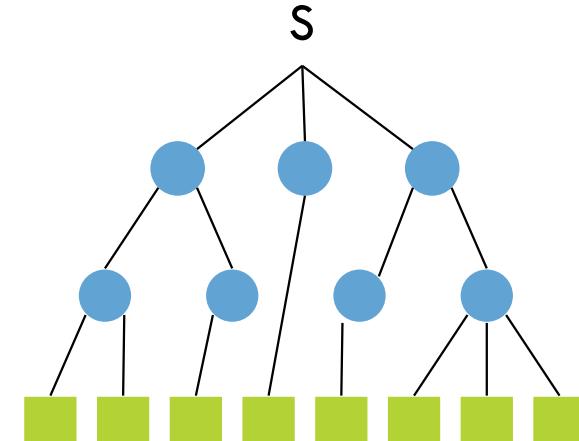
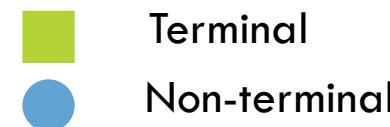
3. Syntax Analysis

- Introduction
- Context free grammars
- Top-down parser
- Bottom-up parser
- Error handling
- Parser generation with yacc



Top-down vs. Bottom-up

- Top-down
 - Start from the root and grow toward the leaves
 - Pick a rule and try to match input (backtrack on mismatch)
- Bottom-up
 - Start from the leaves and grow toward the root
 - Try to match right-hand side of rules
 - Finish once the start symbol is reached



LR: Inverse Left-to-right Rightmost derivation

- Why the name? – Example

CFG:

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

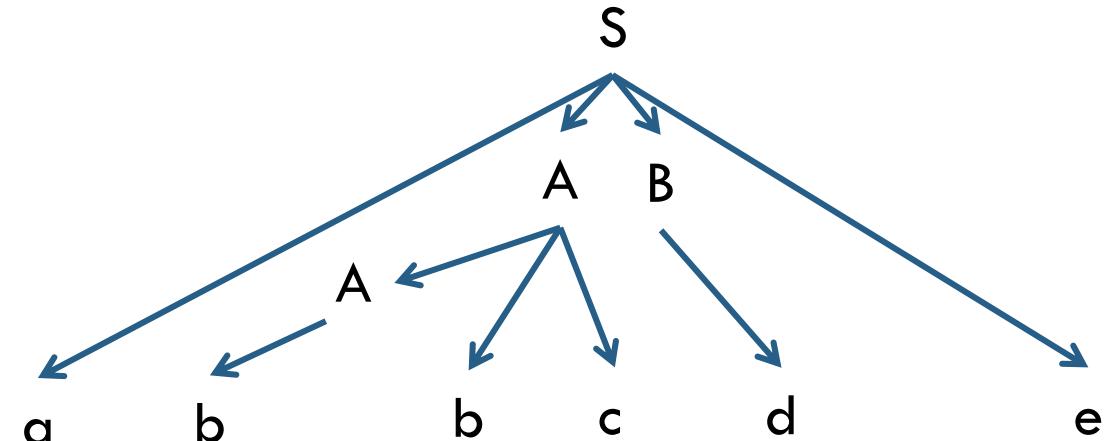
$B \rightarrow d$

$A \rightarrow b$

$A \rightarrow Abc$

$B \rightarrow d$

$S \rightarrow aABe$

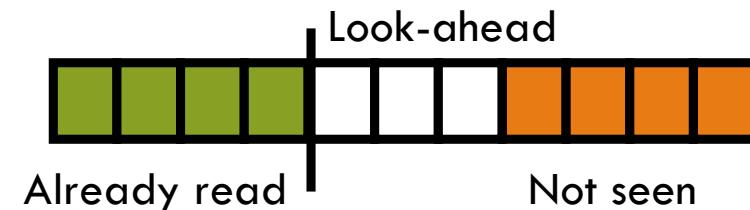


Input:



Construction idea

- The string is split into three parts



- “Already read”: Encode as state
 - Partial derivation (recall from example `aAbcde`)
 - Options to continue the current parsing (e.g., tokens expected to complete a rule)
- Look-ahead + “Not seen”: String made only of terminals
- Parser (main) actions: **Shift/Reduce**
 - **Shift:** No rule completed → Move terminals from right to the left
 - **Reduce:** Rule $A \rightarrow \alpha$ completed → Replace partial sentential form on the left with A

Construction idea: Encoding the state

❑ Idea

- ❑ Do not store the entire partial sentential form
- ❑ Use a stack of partial derivations, i.e. using **locations** in the rules

❑ Example

$S \rightarrow .ABC$

$S \rightarrow A.BC$

Dot: Marks the **location** at which the parser is

Location at parse begin:
We have not read anything

CFG rules

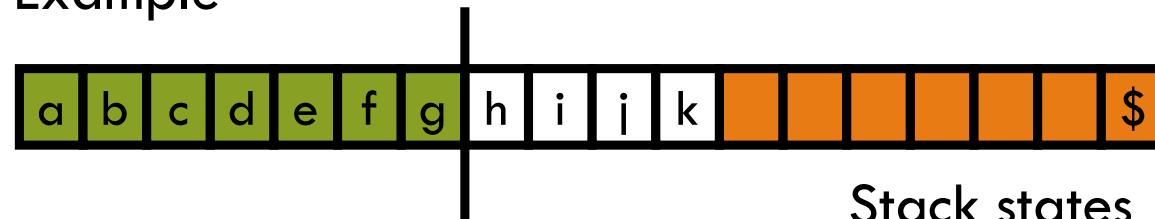
1. $S \rightarrow ABC$
2. $A \rightarrow abcd$
3. $B \rightarrow efgX$
4. $B \rightarrow efgD$
5. ...

Construction idea: Encoding the state

❑ Idea

- ❑ Do not store the entire partial sentential form
- ❑ Use a stack of partial derivations, i.e. using **locations** in the rules

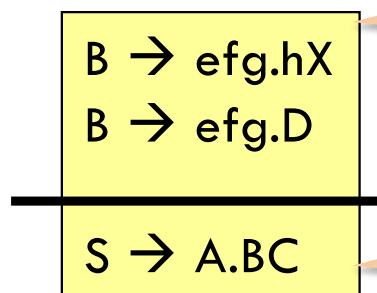
❑ Example



CFG rules

1. $S \rightarrow ABC$
2. $A \rightarrow abcd$
3. $B \rightarrow efgX$
4. $B \rightarrow efgD$
5. ...

Stack states



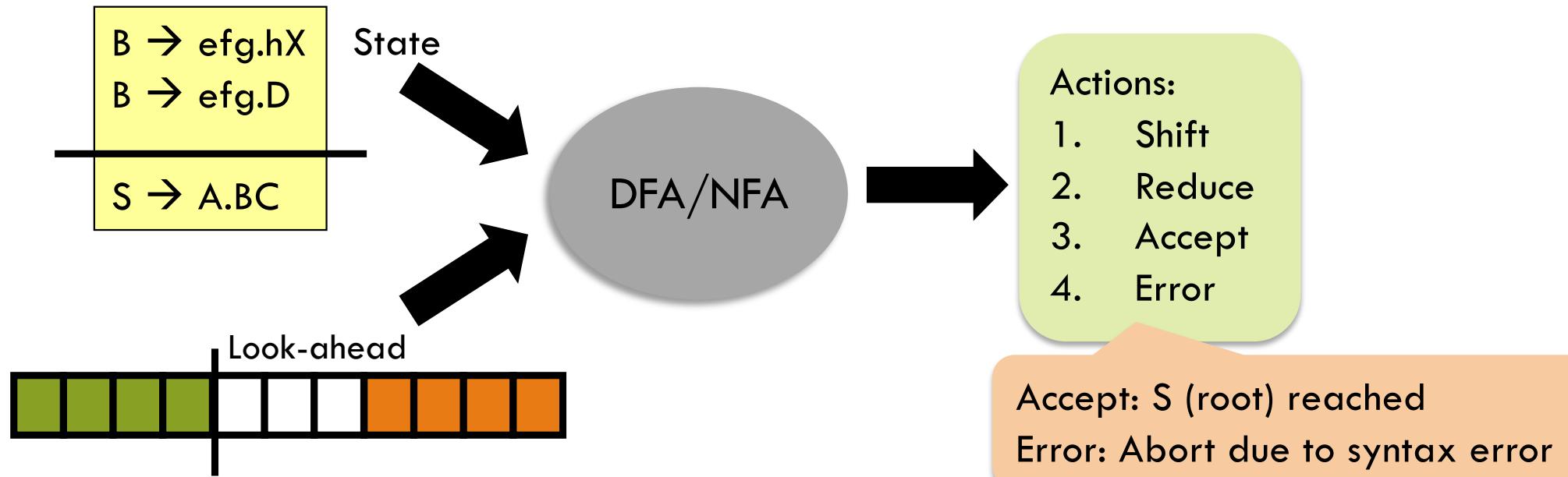
“Trying” to match right-hand side for B: Read efg, the rest from hX or D

“abcd” can be reduced to A, rest should come from BC

Construction idea: Deciding actions (shift/reduce)

□ Idea

- There are finite productions, and finite locations within them
- Use a DFA (or NFA) to track state changes and select **actions**



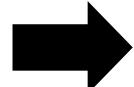
States: LR(0)-Items

Def.: Given a CFG $G = (N, T, P, S)$, an **LR(0)-item** is a rule from P with a special dot added somewhere in the right-hand side

□ Example

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$



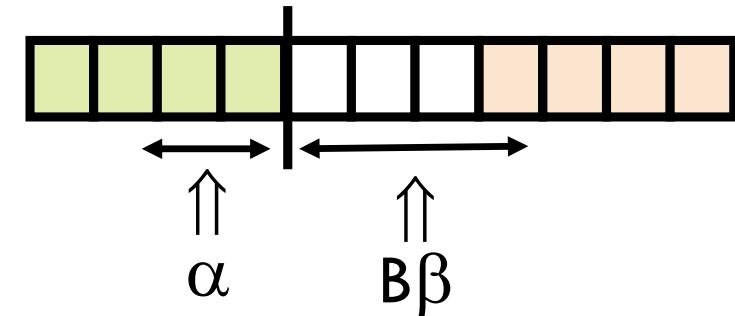
For rule $E \rightarrow T + E$

1. $[E \rightarrow .T + E]$
2. $[E \rightarrow T. + E]$
3. $[E \rightarrow T + .E]$
4. $[E \rightarrow T + E.]$

Content to the left represents
input that has already been read

States: Sets of LR(0)-Items

- Given an LR(0)-item $[A \rightarrow \alpha.B\beta]$
 - Part of the stack was derived from α
 - The rest should have as prefix something derivable from $B\beta$
 - Since $B \in N \rightarrow$ a prefix of the remaining input could also be derived from r_i : $B \rightarrow \beta_i$ for all $r_i \in P$
 - This corresponds to a closure operator
- Example



CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T + E$
4. $T \rightarrow x$

$[S \rightarrow .E\$]$
 $[E \rightarrow .T + E]$
 $[E \rightarrow .T]$
 $[T \rightarrow .x]$

Closure for LR(0)-items

Def.: Given a CFG $G=(N,T,P,S)$ and a set I of LR(0)-items, the closure of the set is defined by $cl(I) = I \cup \{[B \rightarrow .\gamma] \mid [A \rightarrow \alpha.B\beta] \in cl(I) \text{ and } B \rightarrow \gamma \in P\}$

□ Algorithm

```
Input    I
Output   cl(I) = I
While   cl(I) still changing
        foreach [A → α.Bβ] ∈ cl(I)
            foreach B → γ ∈ P
                cl(I) = cl(I) ∪ {[B → .γ]}
return  cl(I)
```

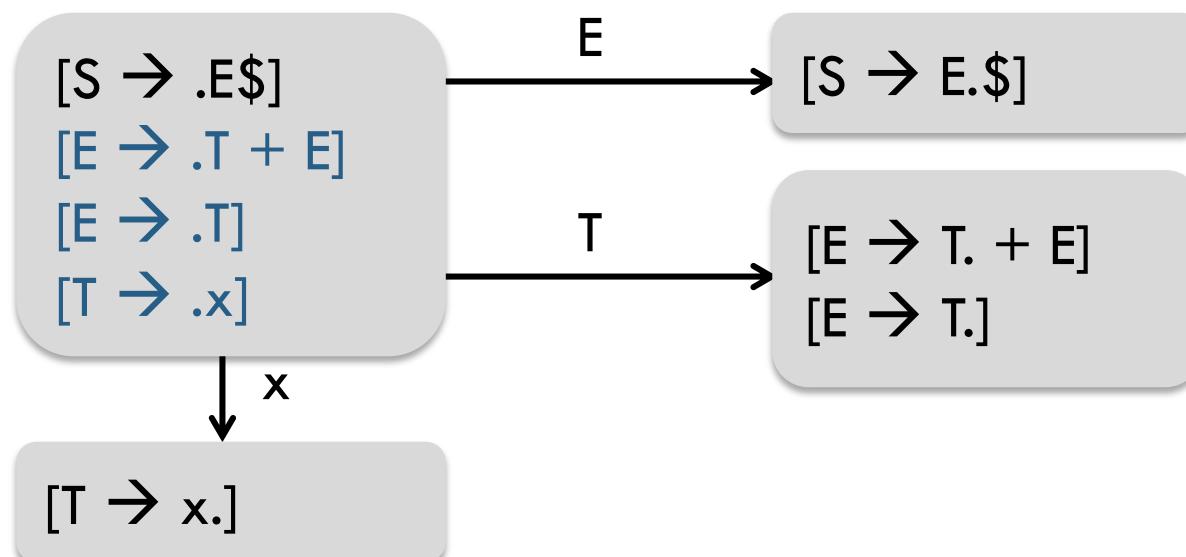
Transitions

Def.: Given a set I of LR(0)-items and a symbol $X \in (N \cup T)$, we define the transition function $GOTO(I, X) = cl(\{ [A \rightarrow \alpha X \beta] \mid [A \rightarrow \alpha X \beta] \in I \})$

- Intuition: Just move the cursor “.” to the right and calculate the closure
- Example

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T + E$
4. $T \rightarrow x$



DFA for LR Parser: Example

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

Entire Right-hand side observed →
Reduce (& roll-back)

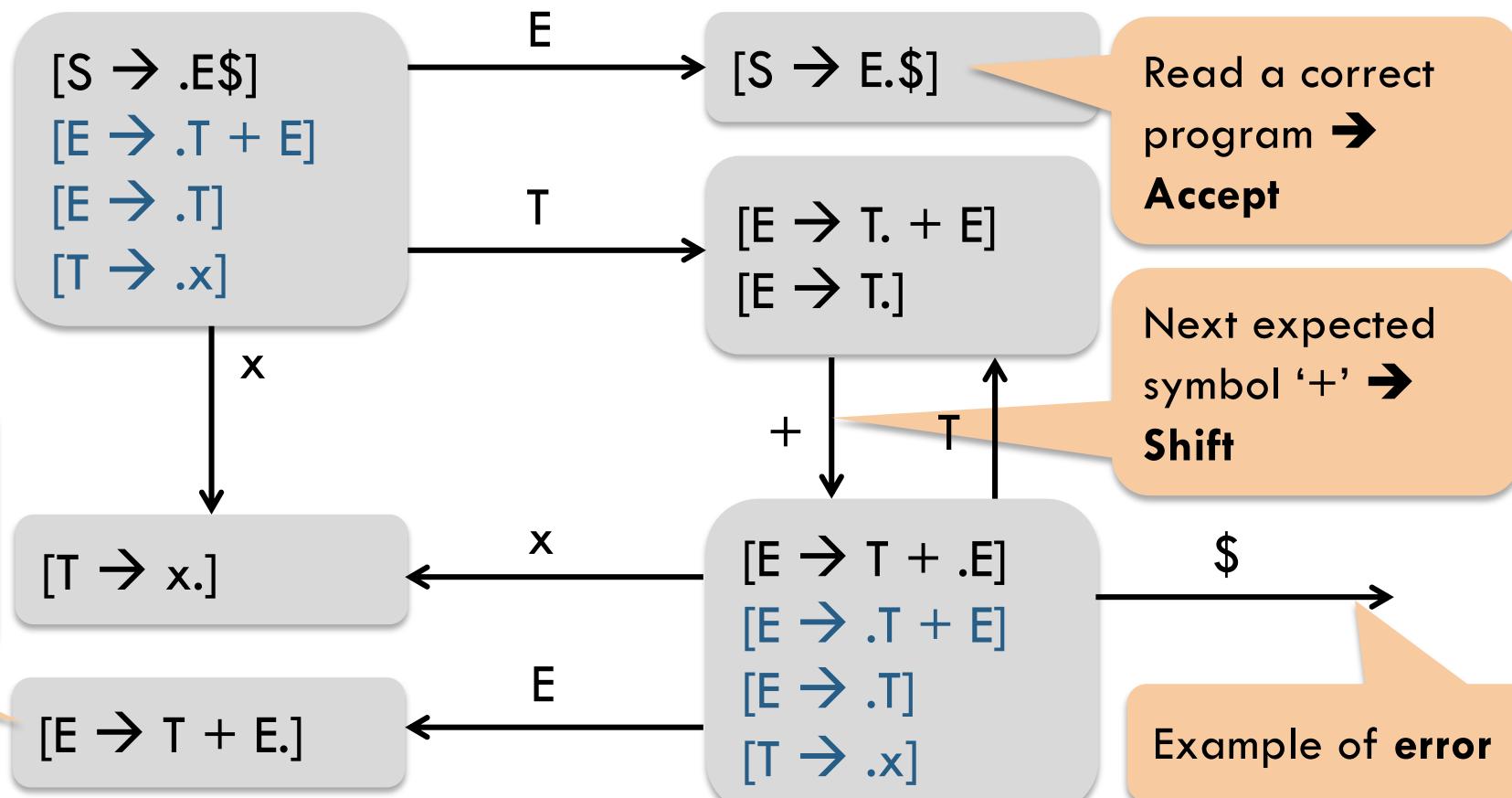


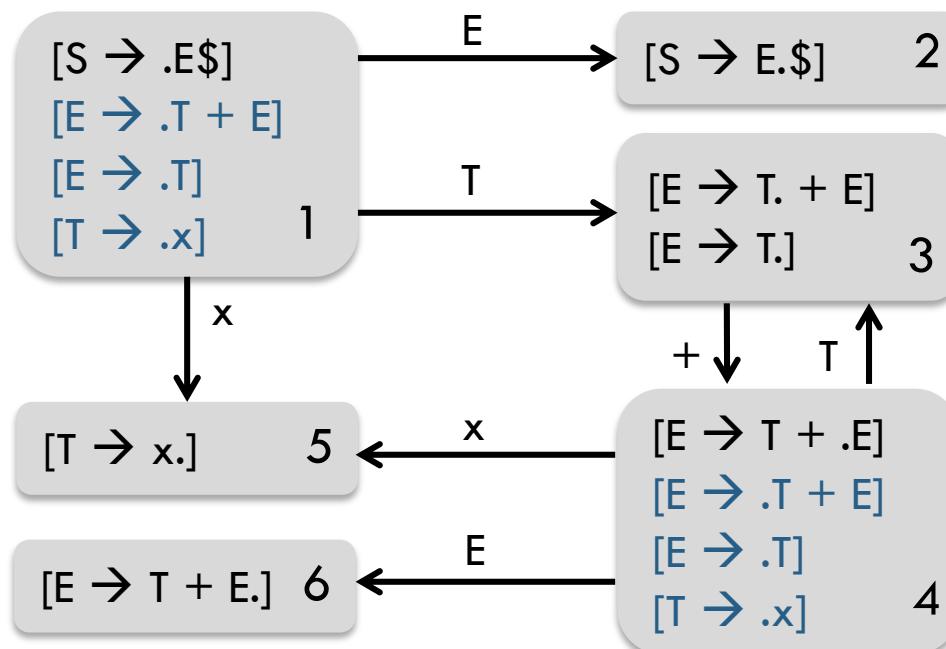
Table-based LR-parser

- Parser implementation through table (similar to LL(k) parsers)
- Given a CFG $G = (N, T, P, S)$, and a transition (I, X, J) of the DFA – I, J : Sets of LR(0)-items & $X \in (N \cup T)$, build table $M[I, X]$ as follows:
 - If $X \in T$ $M[I, X] = \text{shift state } J$
 - If $X \in N$ $M[I, X] = \text{goto state } J$
 - If $[S' \rightarrow S.\$] \in I$ $M[I, \$] = \text{accept}$
 - If $[A \rightarrow \alpha.] \in I$ $M[I, X] = \text{reduce rule } n$ (with $r_n: A \rightarrow \alpha \in P$)
 - Otherwise $M[I, X] = \text{error}$

LR table: Example

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

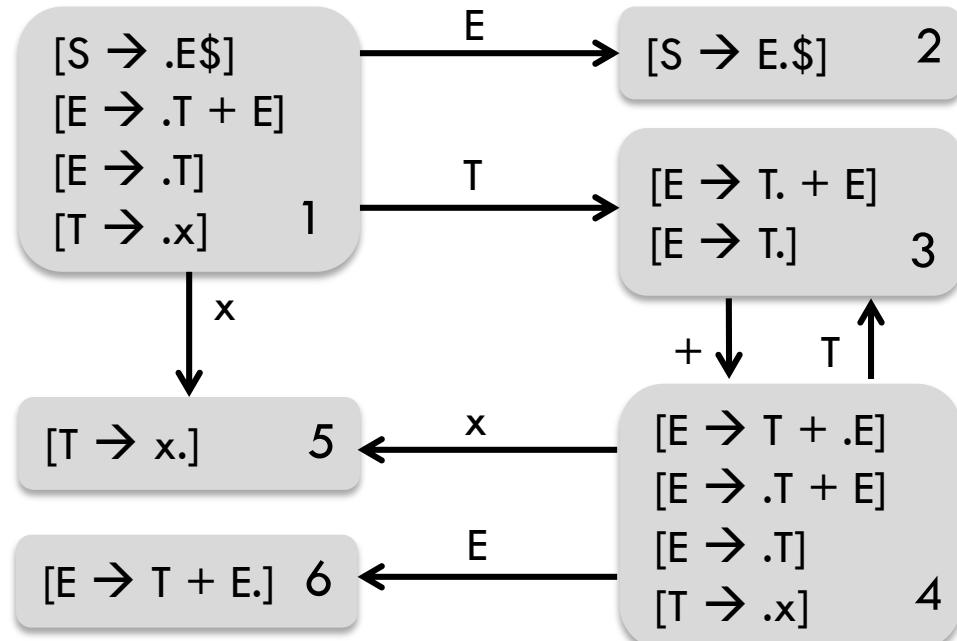


M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

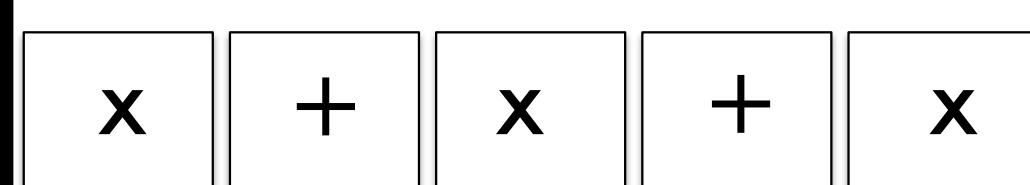
s/g n: shift/goto state n
r n: reduce rule n

Prefer shift
over reduce

LR Parser operation: Intuition



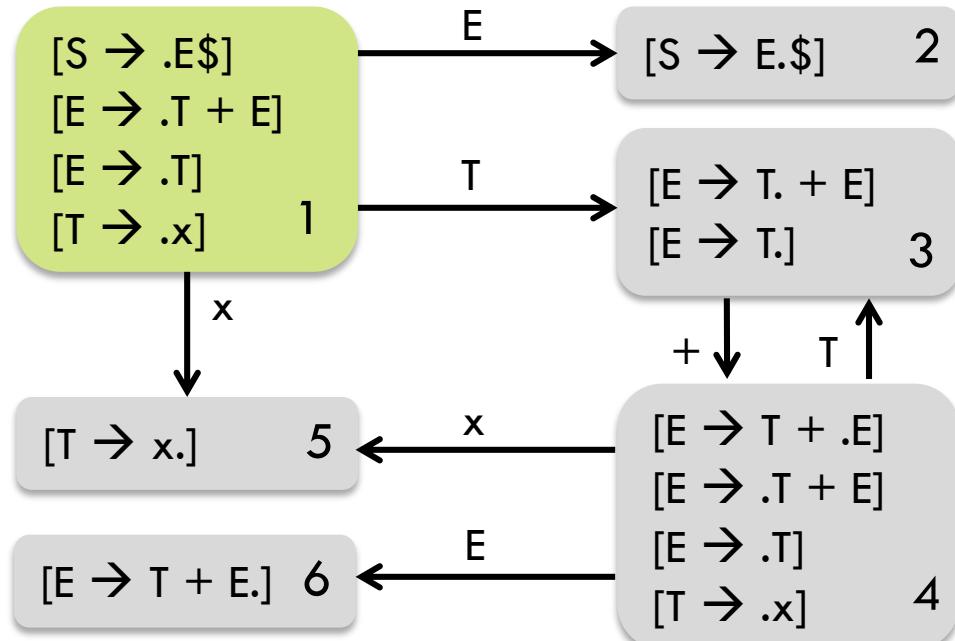
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



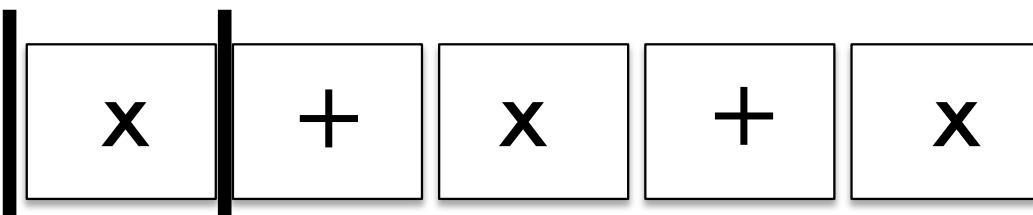
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



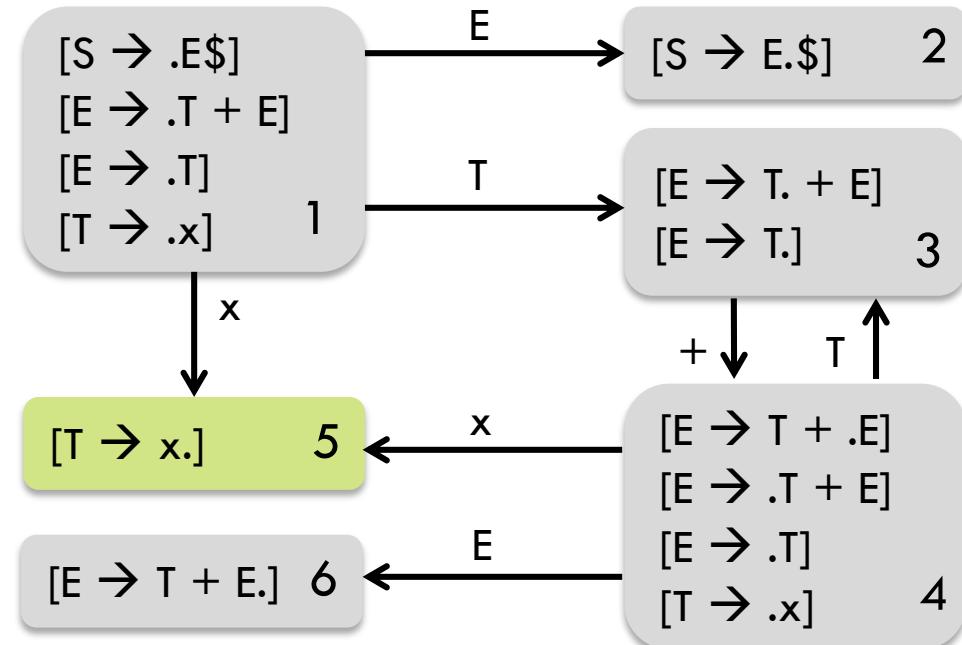
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



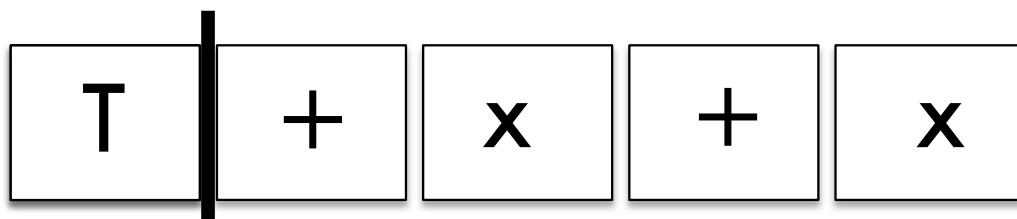
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



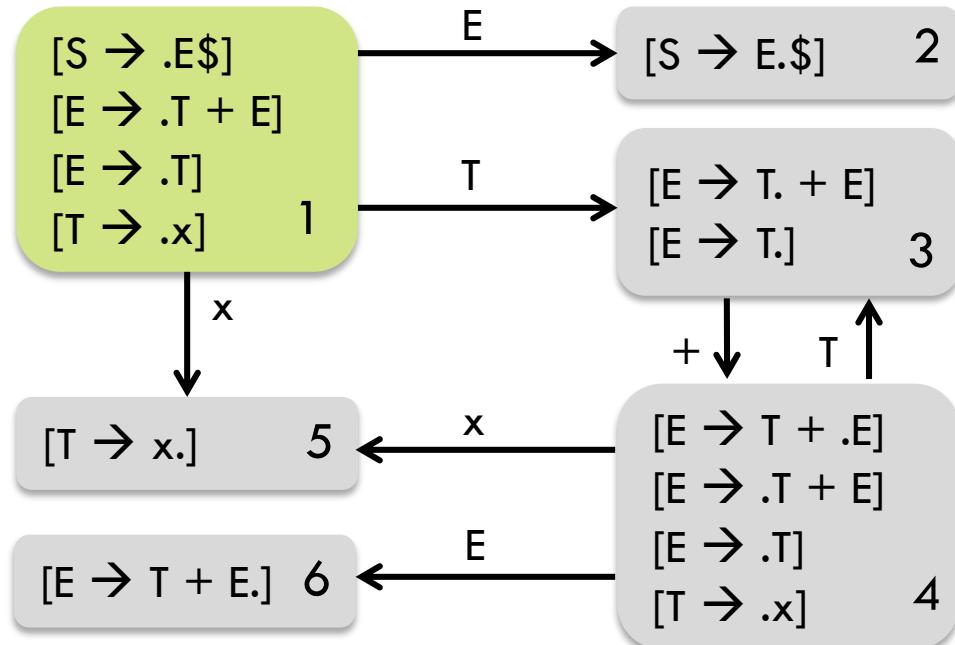
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



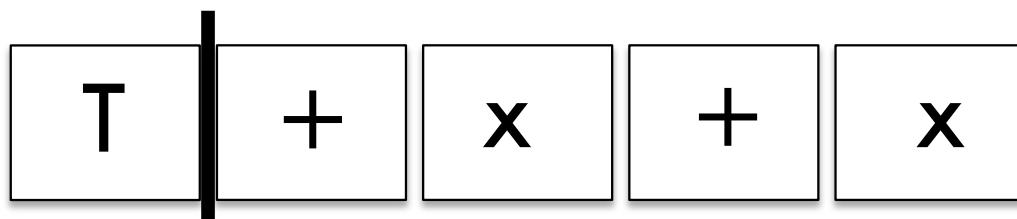
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



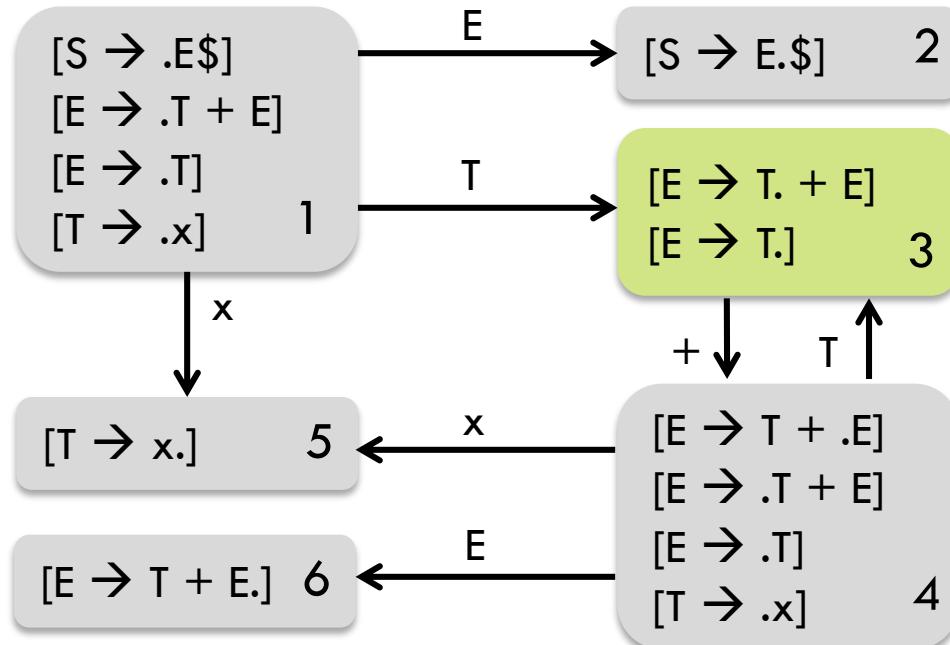
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



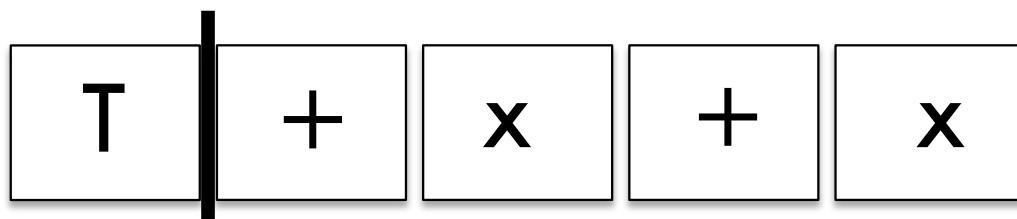
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



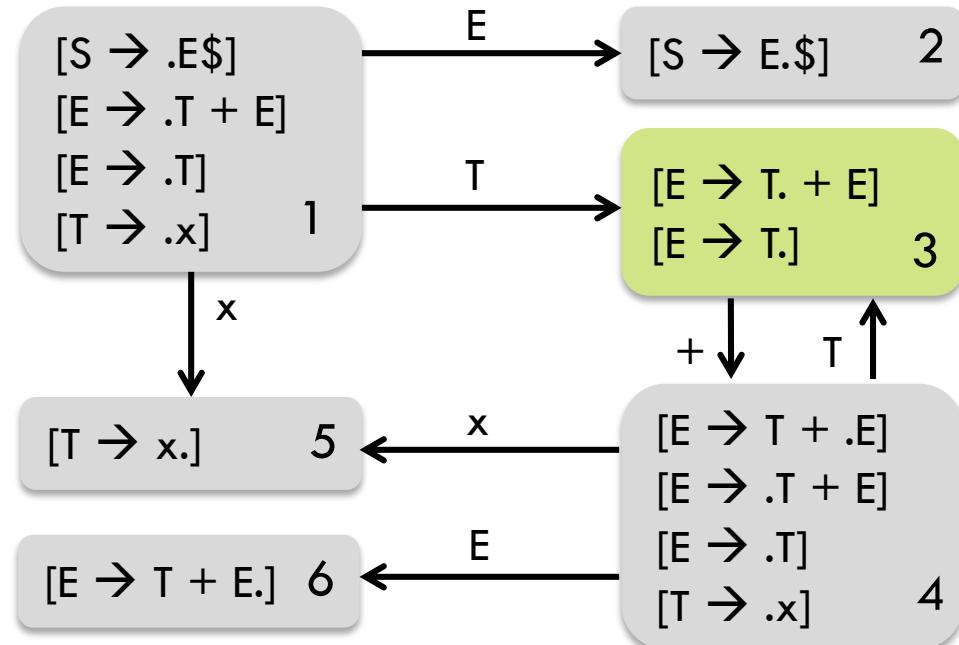
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



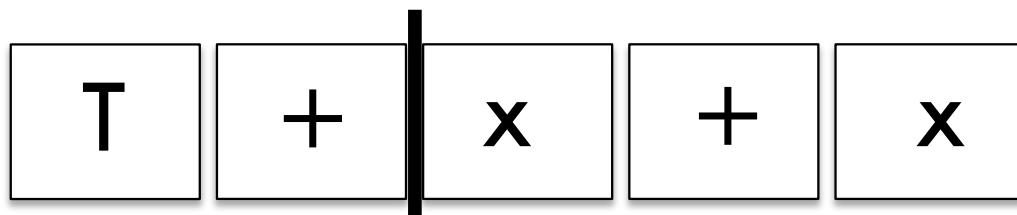
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



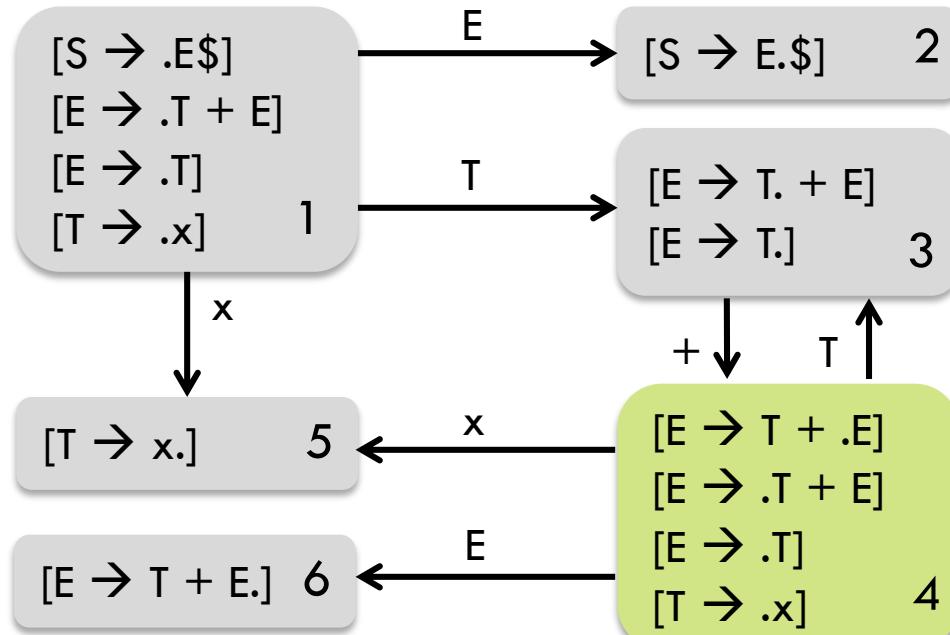
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



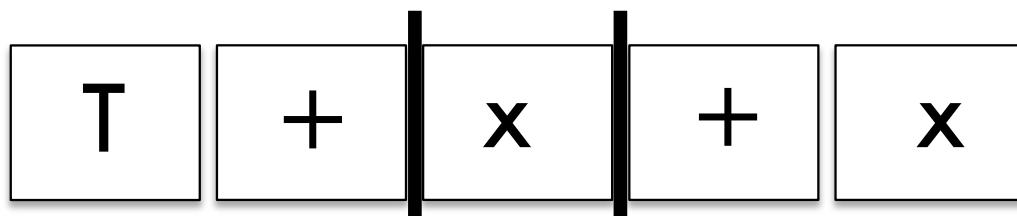
CFG rules

- | | |
|------------------------|------------------------|
| 1. $S \rightarrow E\$$ | 3. $E \rightarrow T+E$ |
| 2. $E \rightarrow T$ | 4. $T \rightarrow x$ |

LR Parser operation: Intuition



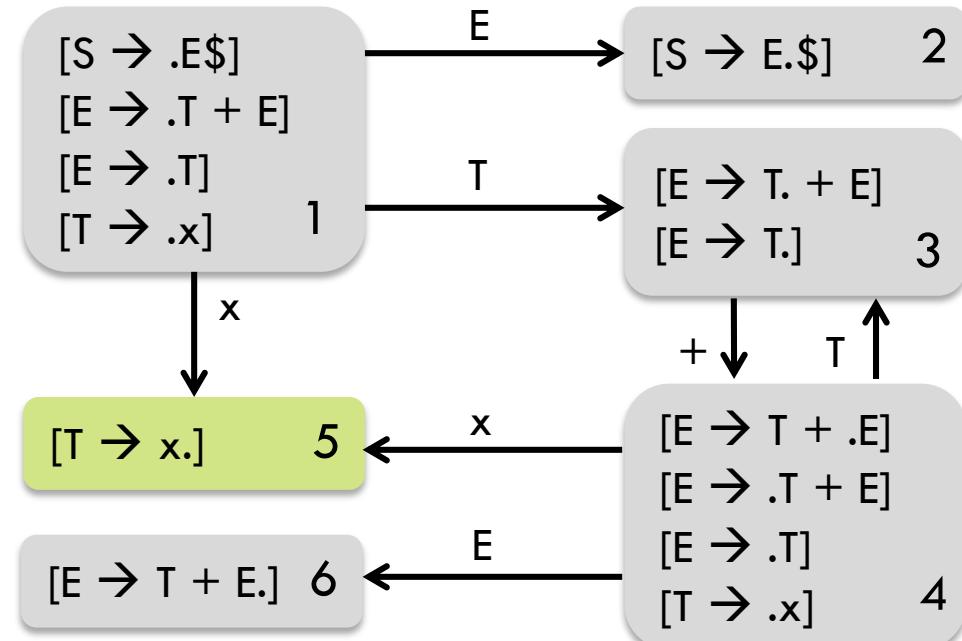
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



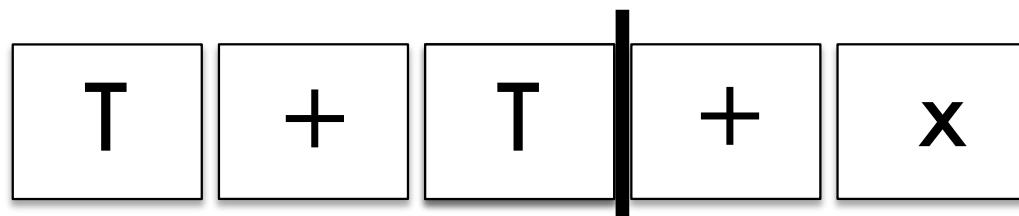
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



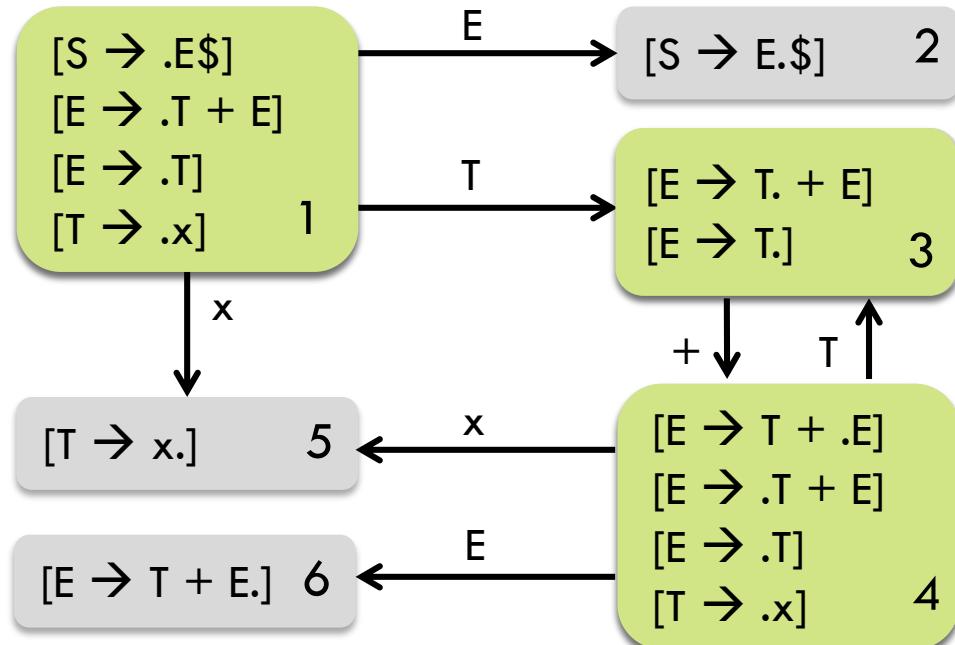
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



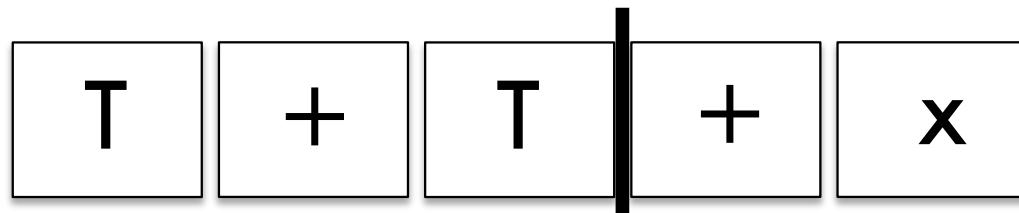
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR Parser operation: Intuition



M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

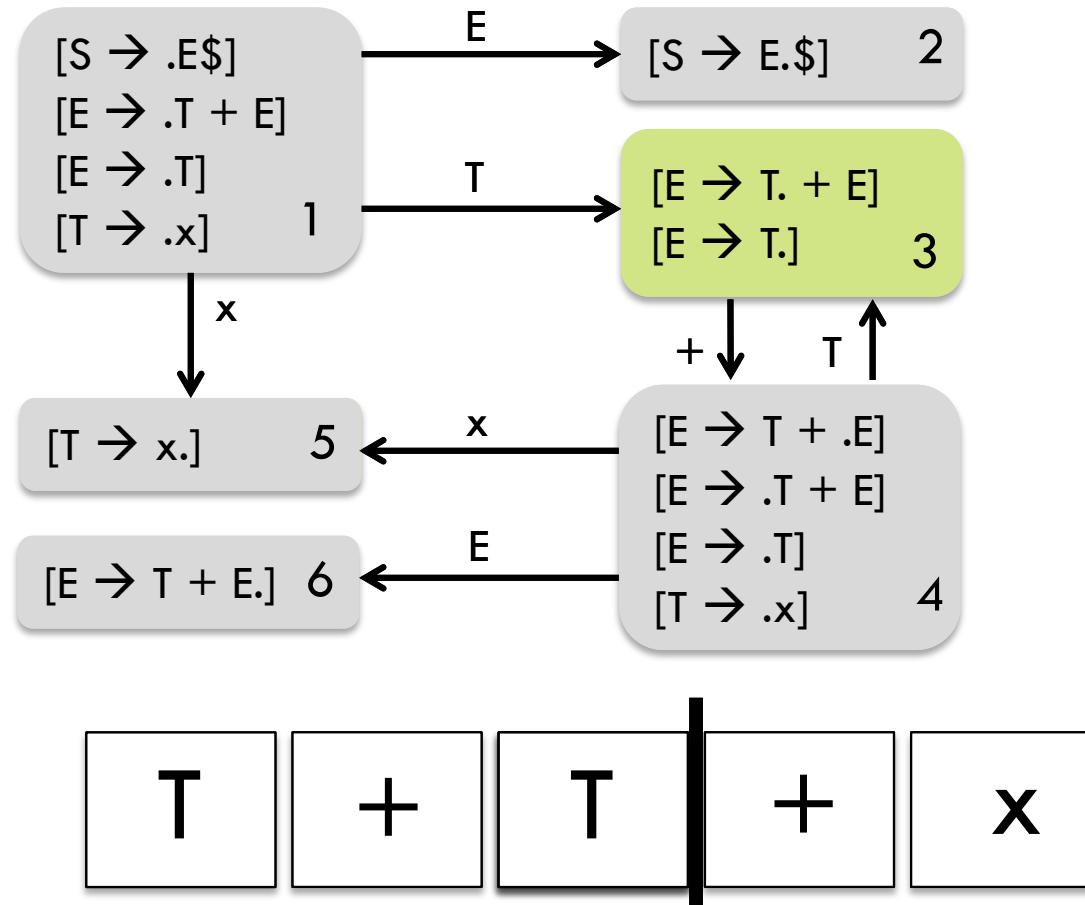


CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

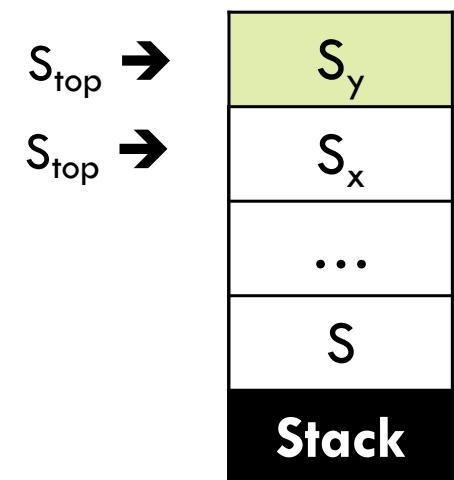
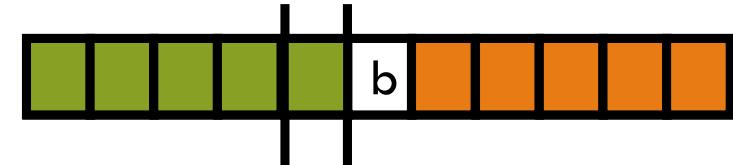
Optimization

- Idea: Do not restart the automaton after every reduce
 - Keep track of the states of the DFA with a stack
 - On reduce $A \rightarrow \alpha$: Go back to the previous state, i.e., roll-back as many transitions as symbols are in α
- Example:
 - Reduce $[T \rightarrow x]$
 - Roll back 1 state (back to 4)
 - Apply GOTO(4, T) = 3



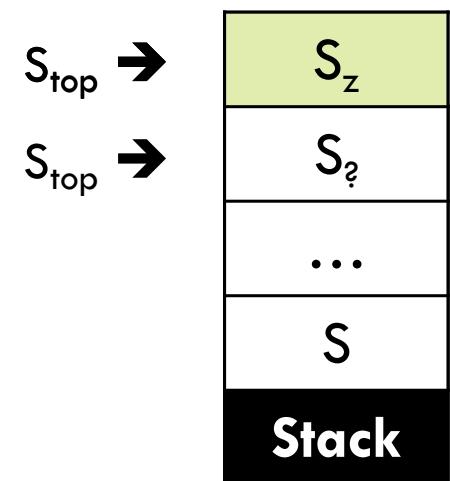
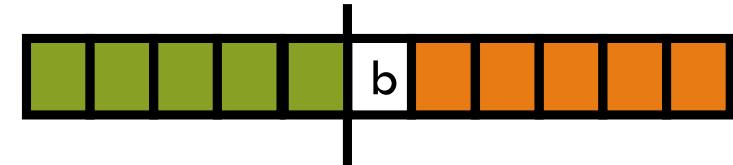
LR operation

- Init: Push start symbol to the stack
- Repeat until “accept” or “error”
 - Let S_{top} be the top of the stack and ‘a’ the next input
 - Compute action from table $act = M[S_{top}, a]$
 - If $act = \text{Shift } S_y$
 - Push S_y to stack and read next symbol

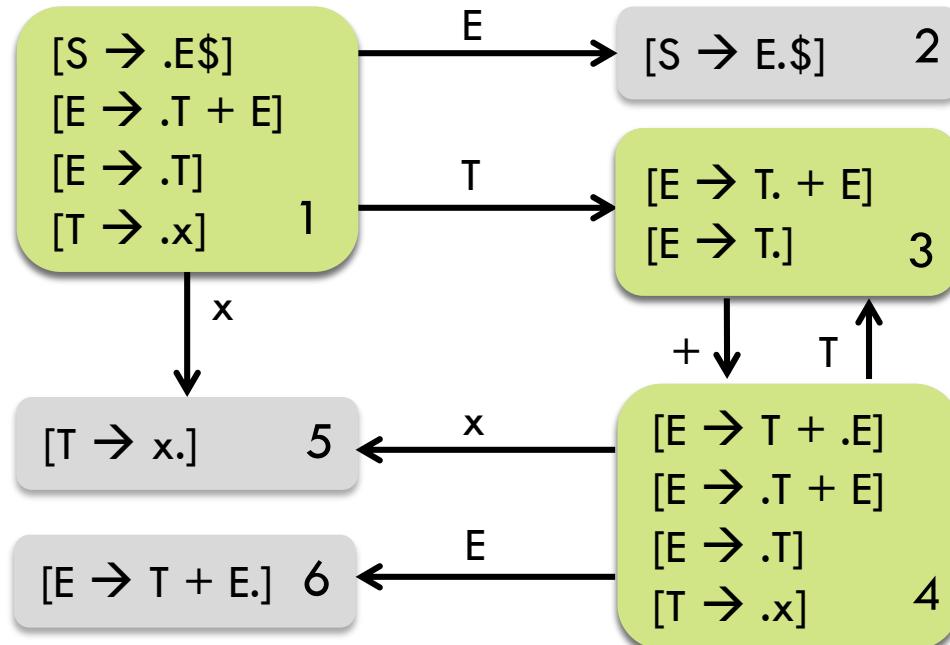


LR operation

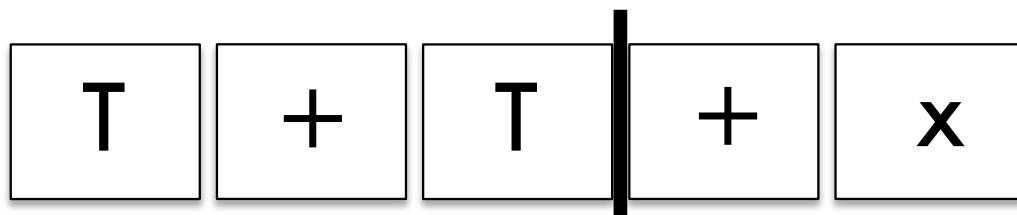
- Init: Push start symbol to the stack
- Repeat until “accept” or “error”
 - Let S_{top} be the top of the stack and ‘a’ the next input
 - Compute action from table $act = M[S_{top}, a]$
 - If $act = \text{Shift } S_y$
 - Push S_y to stack and read next symbol
 - If $act = \text{Reduce } n$ (with $r_n: A \rightarrow \alpha$)
 - Remove $|\alpha|$ states from stack (i.e., roll back)
 - **Do not advance** input cursor
 - Compute new state $M[S_{top}, A] = \text{goto } S_z$, and push to stack
 - If $act = \text{error/accept}$: Report and stop



LR operation: Example



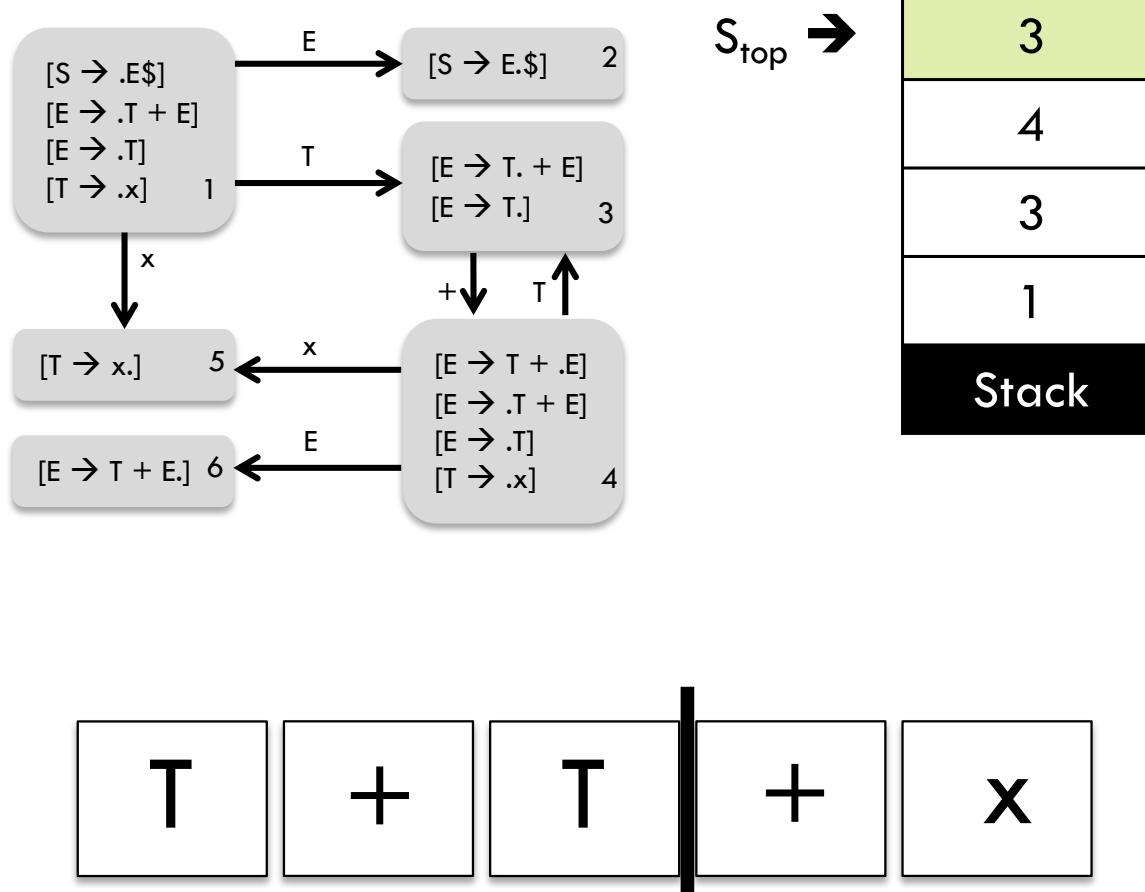
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4	r2	r2	
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR operation: Example

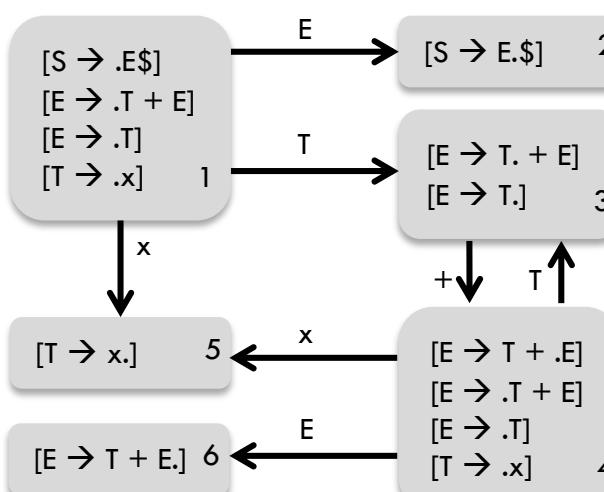


M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

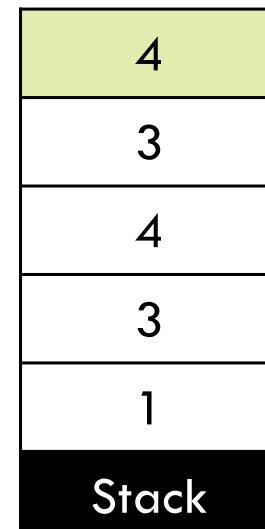
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T + E$
4. $T \rightarrow x$

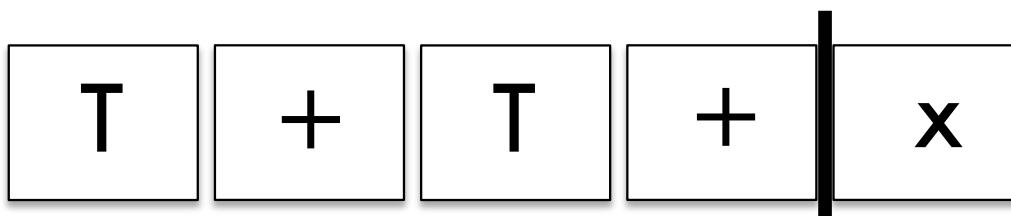
LR operation: Example



$S_{top} \rightarrow$



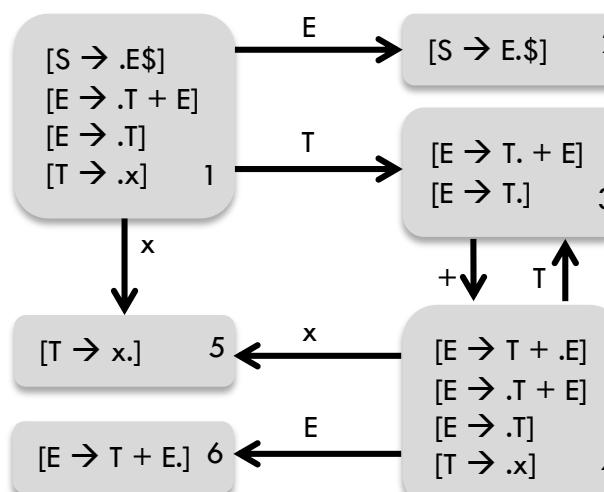
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



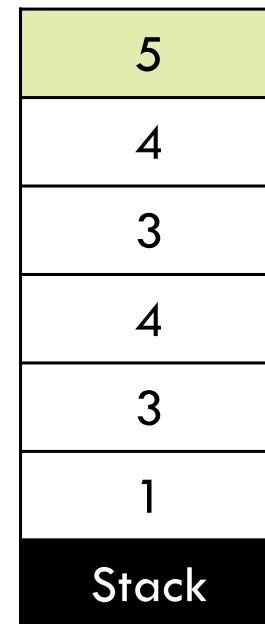
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

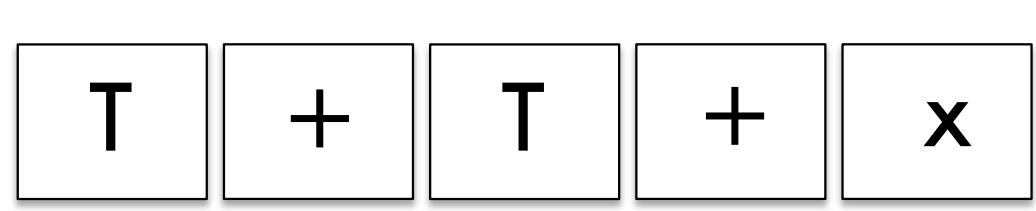
LR operation: Example



$S_{top} \rightarrow$



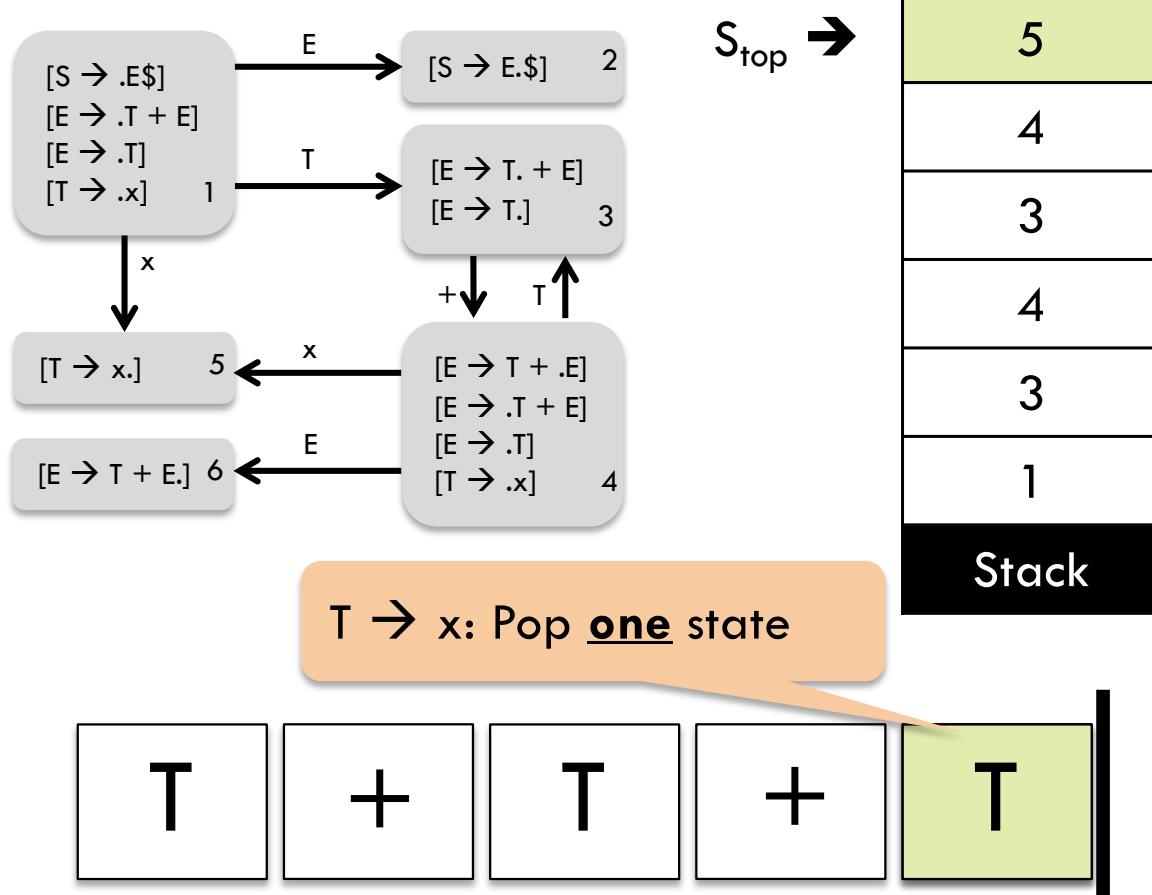
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR operation: Example

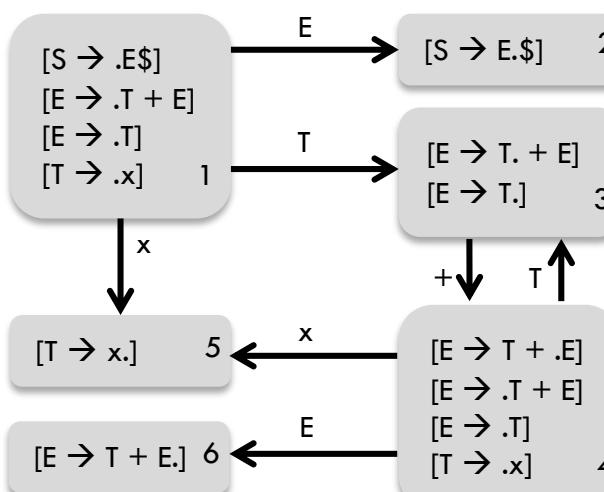


M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

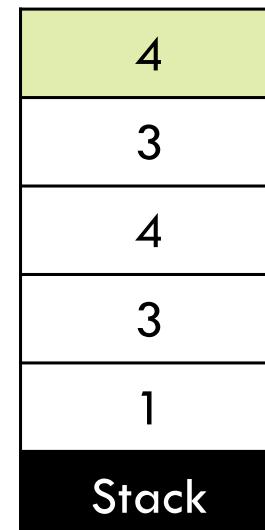
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

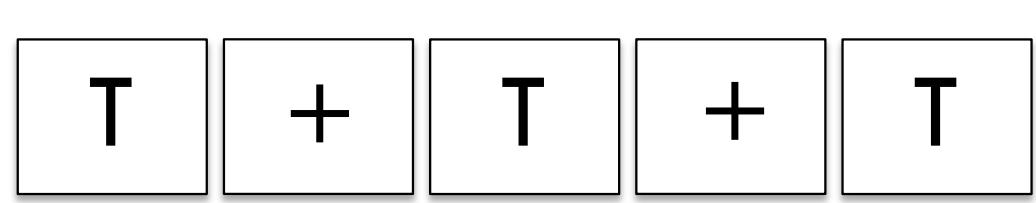
LR operation: Example



$S_{top} \rightarrow$



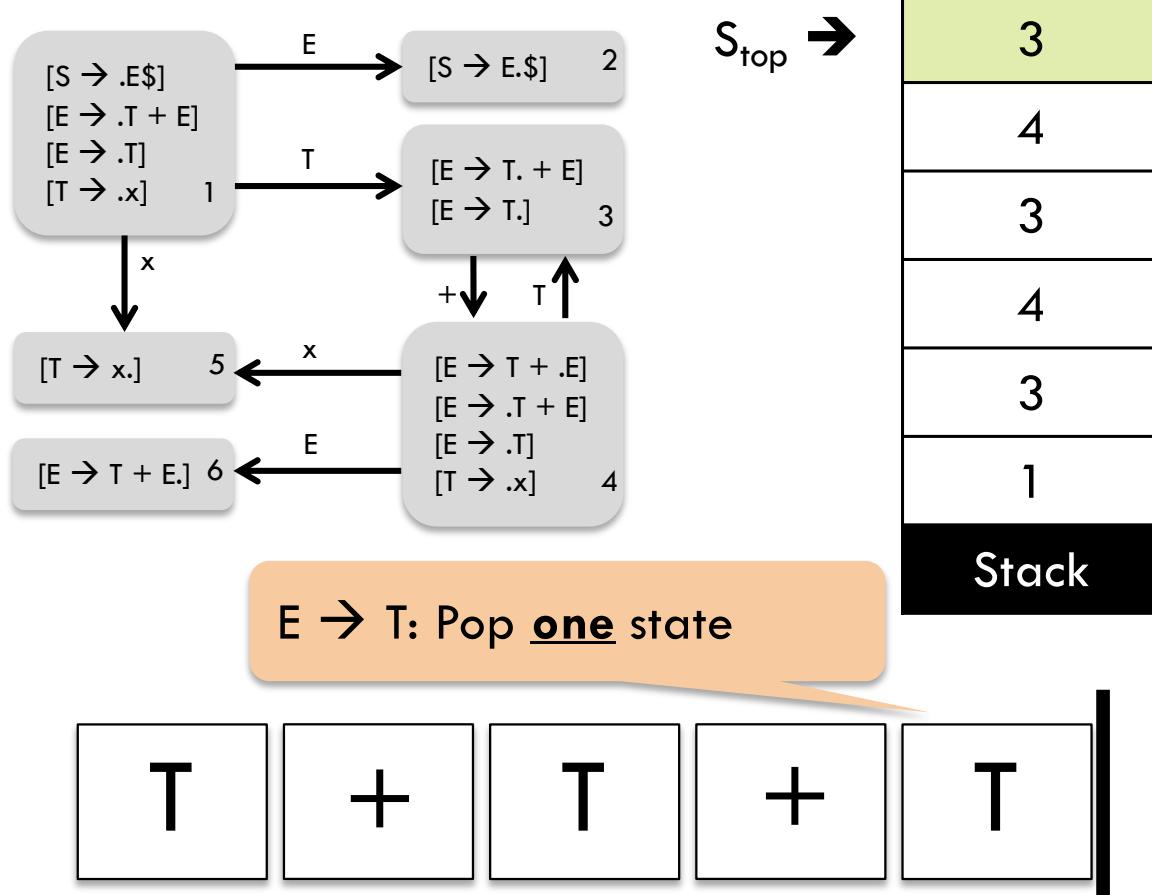
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR operation: Example

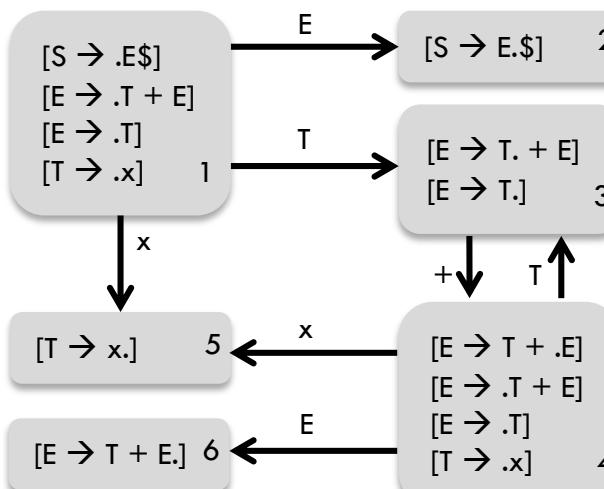


M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

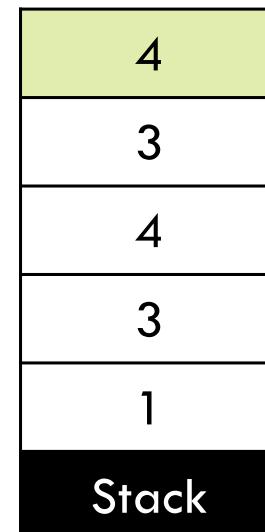
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T + E$
4. $T \rightarrow x$

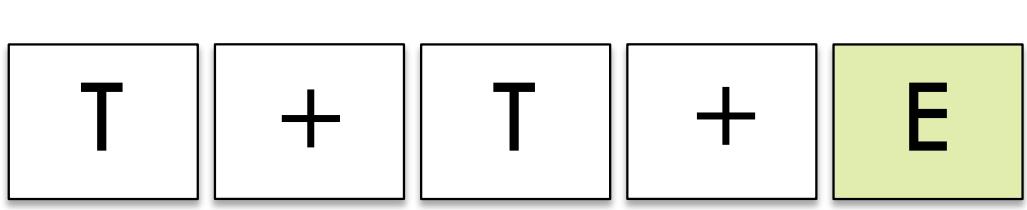
LR operation: Example



$S_{top} \rightarrow$



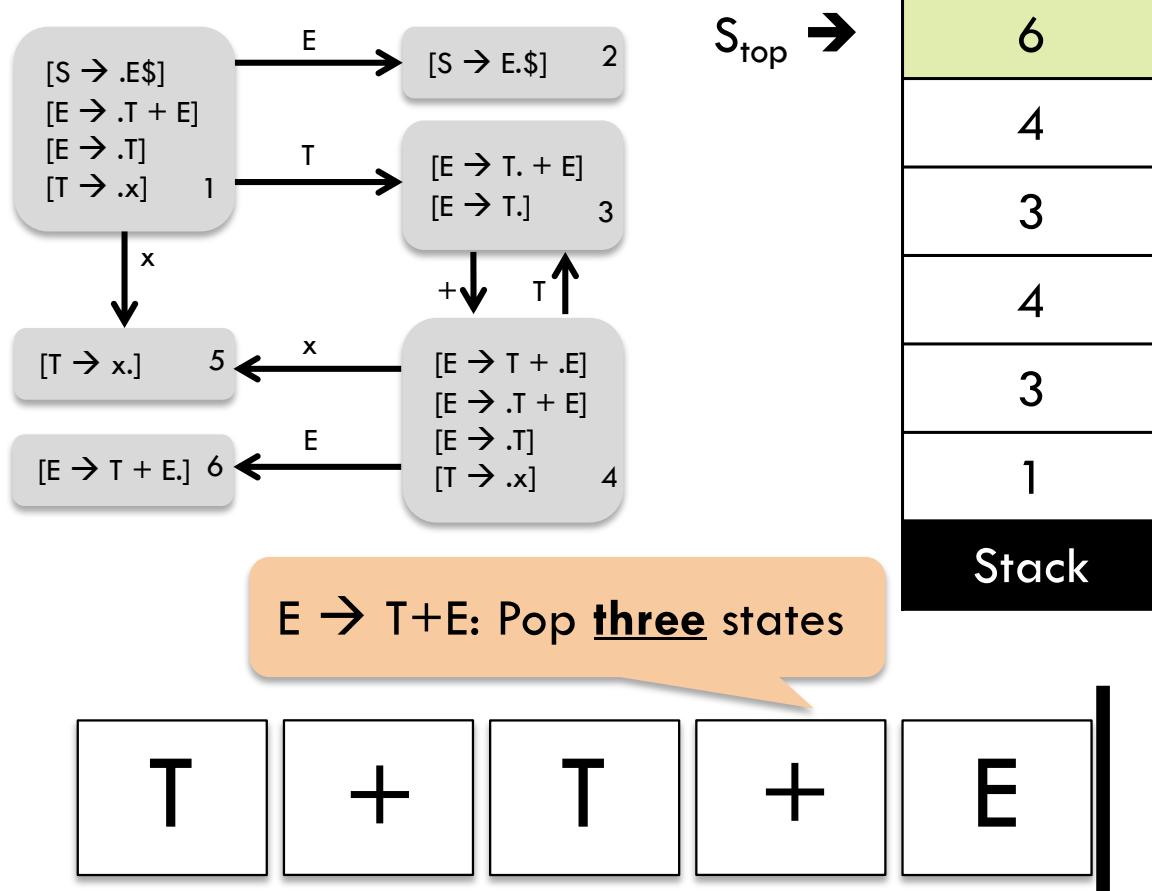
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR operation: Example

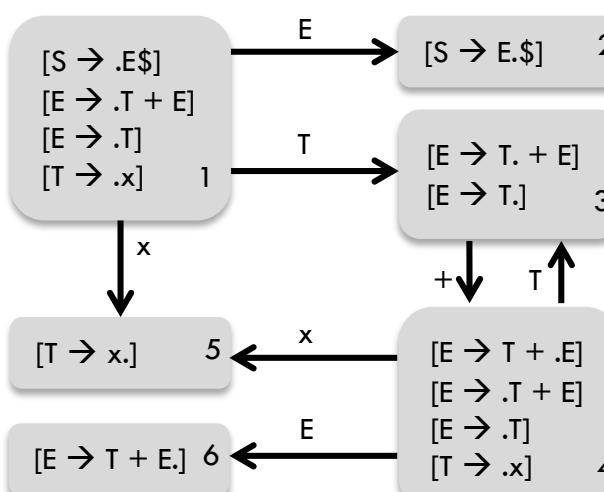


M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

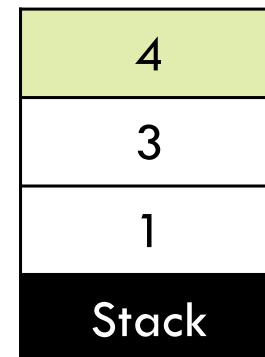
CFG rules

- $S \rightarrow E\$$
- $E \rightarrow T$
- $E \rightarrow T+E$
- $T \rightarrow x$

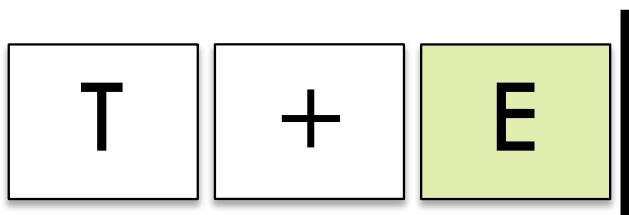
LR operation: Example



$S_{top} \rightarrow$



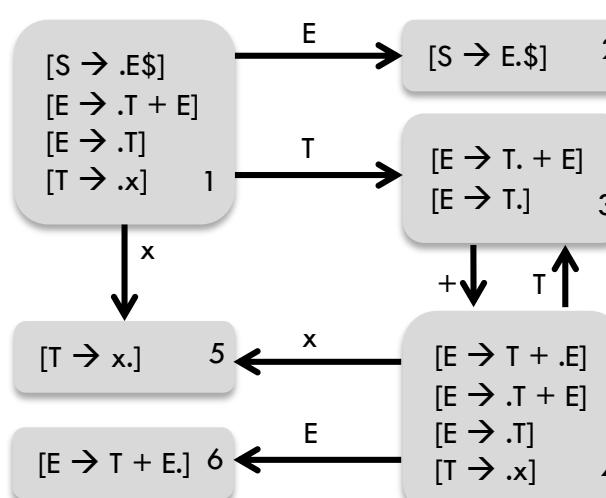
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



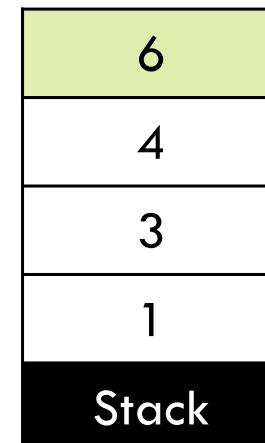
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

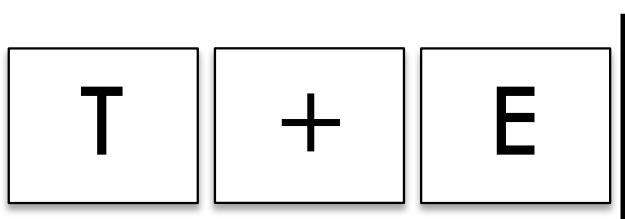
LR operation: Example



$S_{\text{top}} \rightarrow$



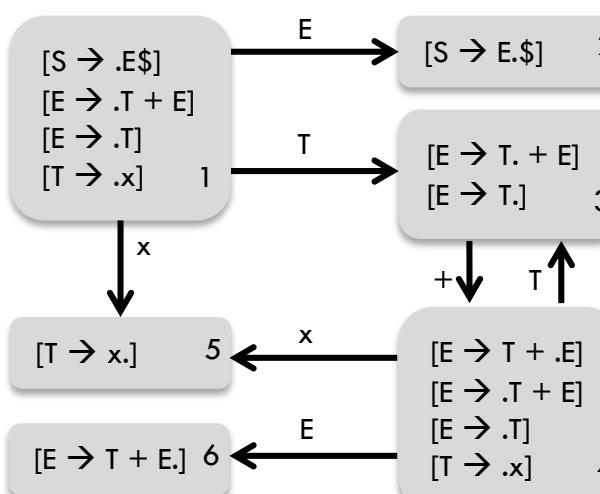
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		



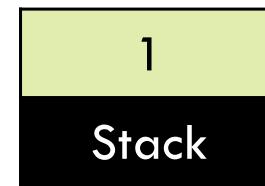
CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$

LR operation: Example



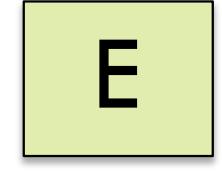
$S_{top} \rightarrow$



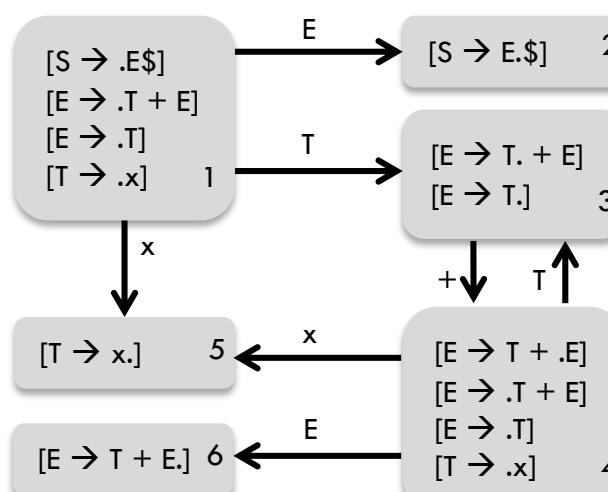
M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

CFG rules

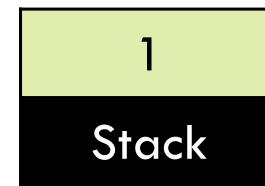
1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$



LR operation: Example



$S_{top} \rightarrow$



M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4 r2	r2		
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$



LR(k) grammars

Def.: A CFG G is in the class $LR(0)$, denoted $G \in LR(0)$ if the parse table has at most one entry per cell (no conflicts shift/reduce or reduce/reduce)

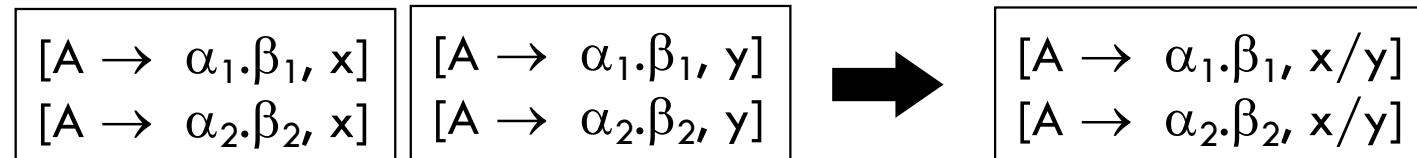
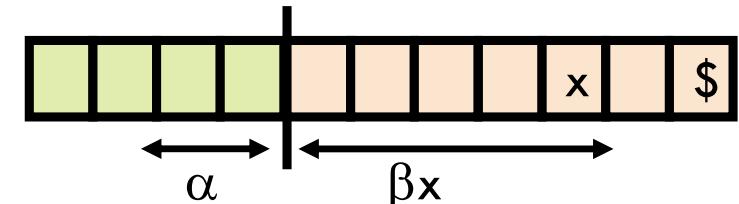
- From the example...
- Shift-reduce conflicts typically resolved by choosing shift (yacc)
- Reduce-reduce are more complex
- Solution: As with $LL(k)$, use look-ahead tokens $LR(k)$
 - Table grows exponentially with k
- Why is it call $LR(0)$?

M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4, r2	r2		
4	s5			g6	g3
5	r4		r4		
6					

Would it help to use semi-colon?

LR(k) and other grammars

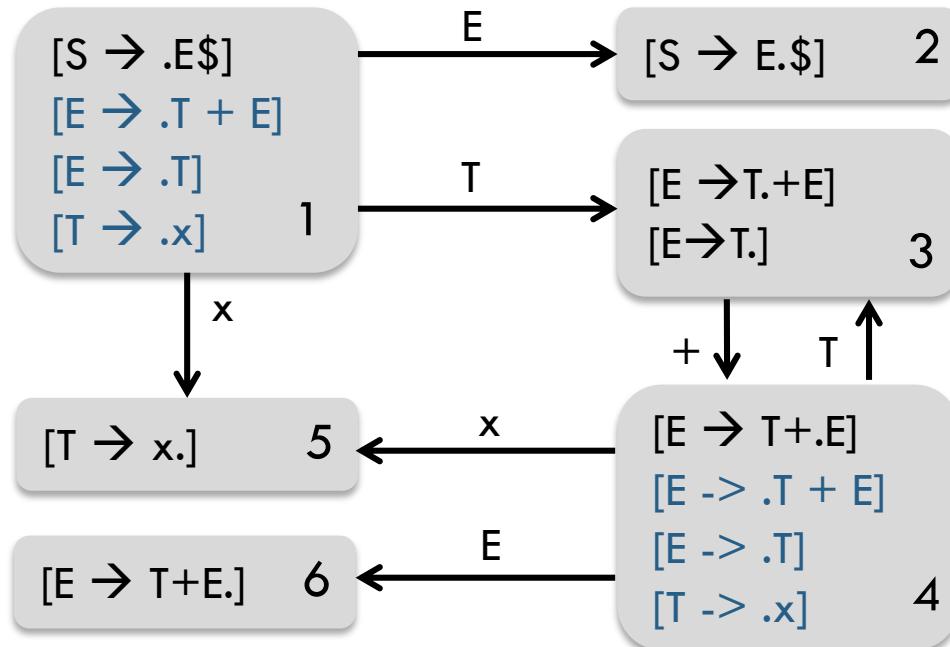
- For practical languages LR(1) is usually enough
 - New definitions: LR(1)-items & GOTO function
 - LR(1)-item are of the form $[A \rightarrow \alpha.\beta, x]$
 - Intuition: Use the look-ahead to disambiguate in case of conflicts
- Complexity reduction
 - SLR: Simple LR
 - Extension to LR(0): Reduce $A \rightarrow \alpha$ only for $\alpha \in \text{FOLLOW}(A)$
 - LALR(1): Look-Ahead(1) LR(0) – (The one used in yacc)
 - Start from LR(1) and merge states with the same **core**



Example: SLR for the example

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$



[$E \rightarrow T+E$]
[$E \rightarrow T.$] 3

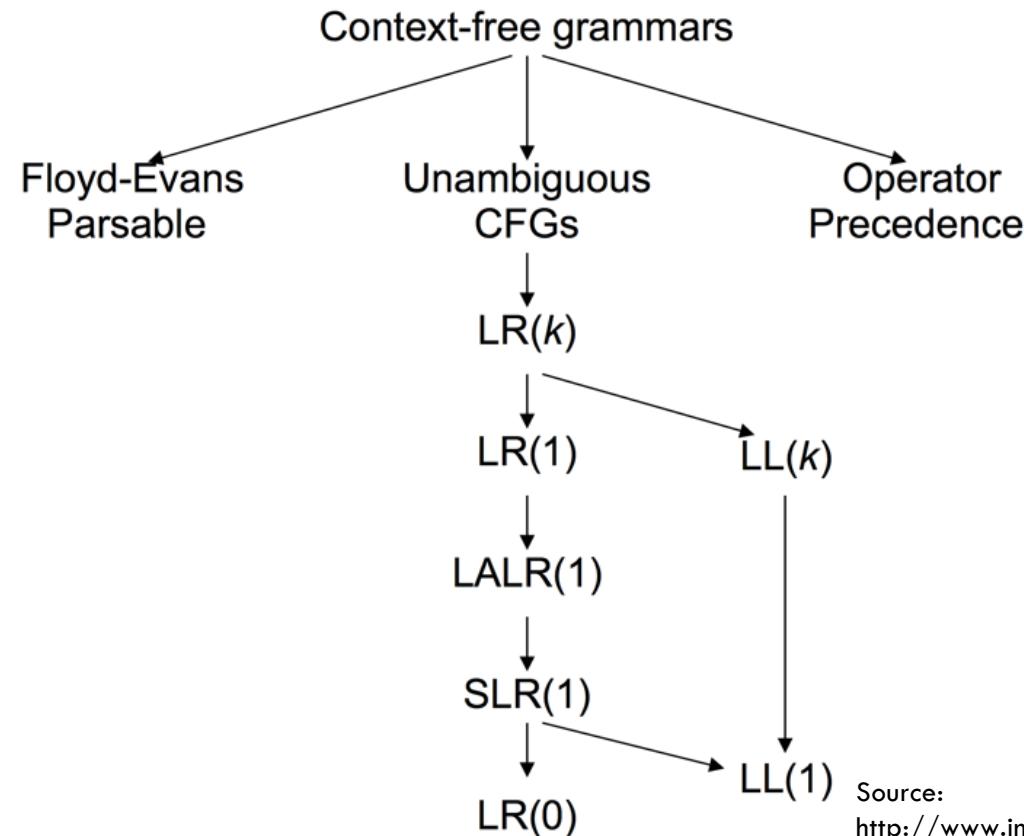


$\text{FOLLOW}(E) = \{ \$ \}$
 $\text{FOLLOW}(T) = \{ +, \$ \}$

M	X	+	\$	E	T
1	s5			g2	g3
2			accept		
3	r2	s4	r2	r2	
4	s5			g6	g3
5	r4	r4	r4		
6	r3	r3	r3		

Grammars so far and some others

- Inclusion hierarchy of context free grammars



Source:
<http://www.inf.ed.ac.uk/teaching/courses/ct/slides/Lecture07.pdf>

3. Syntax Analysis

- Introduction
- Context free grammars
- Top-down parser
- Bottom-up parser
- Error handling
- Parser generation with yacc



Error handling so far

M	x	+	\$	E	T
1	s5	E	E	g2	g3
2	E	E	accept	E	E
3	E	s4	r2	E	E
4	s5	E	E	g6	g3
5	E	r4	r4	E	E
6	E	E	r3	E	E

Error handling so far

- Desired
 - More information about the error
 - Resume after finding one error
- Ideas
 - Add common error rules to the grammar
 - Add error token to grammar and synchronization tokens to restart

M	x	+	\$	E	T
1	s5	E	E	g2	g3
2	E	E	accept	E	E
3	E	s4	r2	E	E
4	s5	E	E	g6	g3
5	E	r4	r4	E	E
6	E	E	r3	E	E

Error rules

- Useful to catch user programming mistakes
- Extending the grammar
 - Introduce sources of ambiguity
 - Make language fall outside of class, e.g., LALR(1)

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow T+E$
4. $T \rightarrow x$
5. $E \rightarrow EE$
6. $E \rightarrow E+$

Synchronization tokens

- Synchronization tokens: Those after error token
- Technique used by yacc
- Behavior on error
 - Set error flag
 - Pop from stack until $M[s,error] = \text{shift}$
 - Jump in the input until the synchronization token is found
 - Report error
 - Continue parsing

CFG rules

1. $S \rightarrow E\$$
2. $E \rightarrow T$
3. $E \rightarrow (E)$
4. $E \rightarrow T+E$
5. $T \rightarrow x$
6. $E \rightarrow (\text{error})$
7. $T \rightarrow \text{error } x$

3. Syntax Analysis

- Introduction
- Context free grammars
- Top-down parser
- Bottom-up parser
- Error handling
- Parser generation with yacc



yacc (bison)

- yacc: “yet another compiler compiler” developed in 1970 by S. Johnson, AT&T
 - Input: yacc specification of the CFG with semantic actions
 - Output: LALR(1) parser in C code – call it with **yyparse()**
 - Output: Conflict reports
 - Requisite: yacc expects the lexer to be available **yylex()**
- Input format

```
%{  
C-declarations  
%}  
Tokens  
%%  
CFG Rules  
%%  
C code section
```

yacc example

□ Sample rule for expressions

```
expr
  : expr[L] TOKEN_PLUS expr[R]
    {
      $$ = createOperation( ePLUS, $L, $R );
    }
  | expr[L] TOKEN_MULTIPLY expr[R]
    {
      $$ = createOperation( eMULTIPLY, $L, $R );
    }
  | TOKEN_LPAREN expr[E] TOKEN_RPAREN { $$ = $E; }
  | TOKEN_NUMBER { $$ = createNumber($1); }
;
```

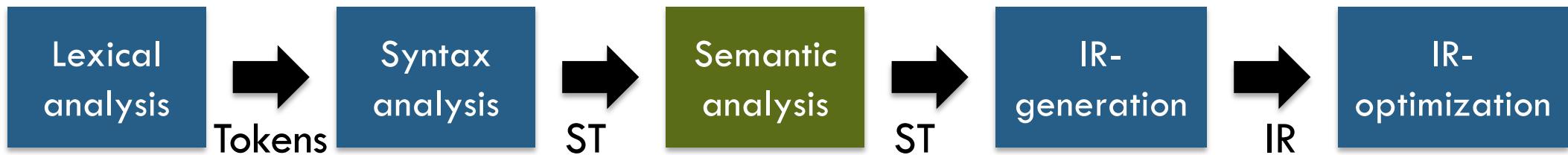
Where are we?

1. Introduction
2. Lexical analysis
3. Syntax analysis
4. Semantic analysis
5. Intermediate representation
6. Control & data-flow analysis
7. IR optimization
8. Target architectures
9. Code selection
10. Register allocation
11. Scheduling
12. Advanced topics

4. Semantic Analysis

- Introduction
- Attribute grammars
- Applications

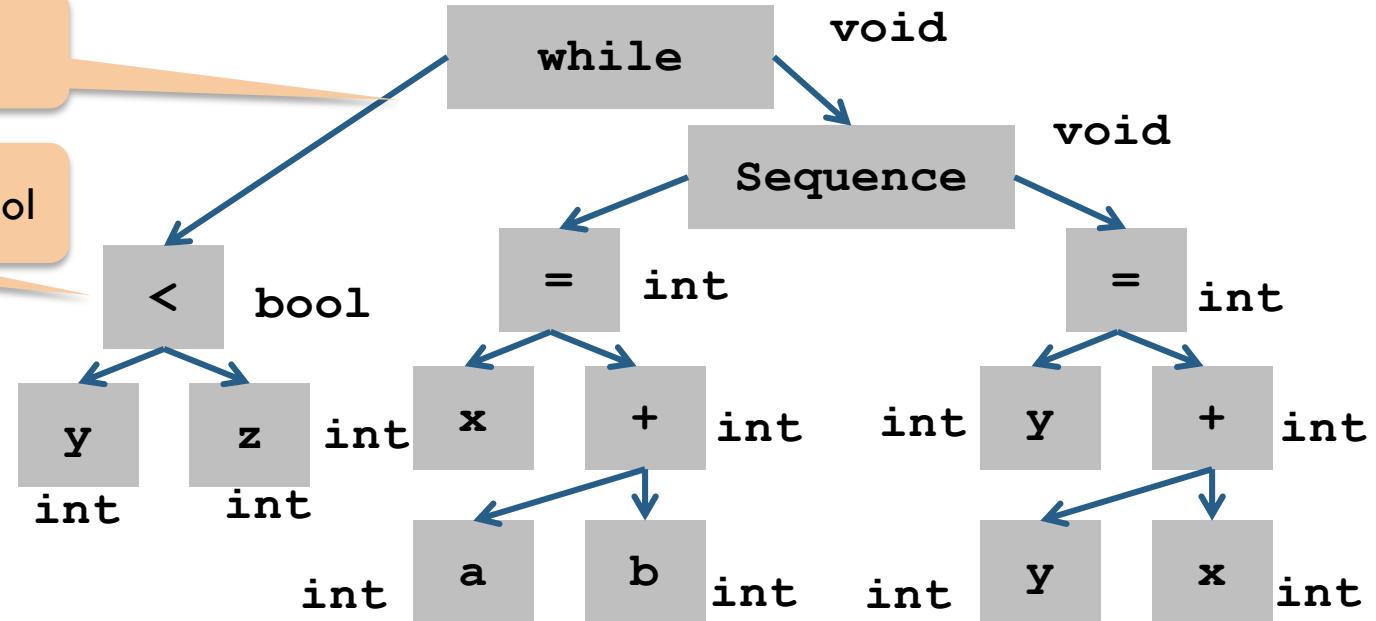
Semantic analysis recall



ST: Syntax Tree (= parse tree)

Semantic: int < int → produces a bool

```
while (y < z)
{
    int x = a + b;
    y += x;
}
```



Why semantic analysis?

```
foo(int a,int b,int c,int d)
{ ... }
bar()
{
    int f[3],g[0], h, i, j, k;
    char *p;
    foo(h,i,"ab",j, k);
    k = f * i + j;
    h = g[17];
    printf("<%s,%s>.\n", p,q);
    p = 10;
}
```

To generate code, the compiler must understand meaning

Wrong number of arguments and type

Wrong dimension for f

Accessing g[17], but declared g[0]

q unknown

p is "char*"

Adapted from: Keith D. Cooper, Ken Kennedy & Linda Torczon

Can't we use CFGs for all this?

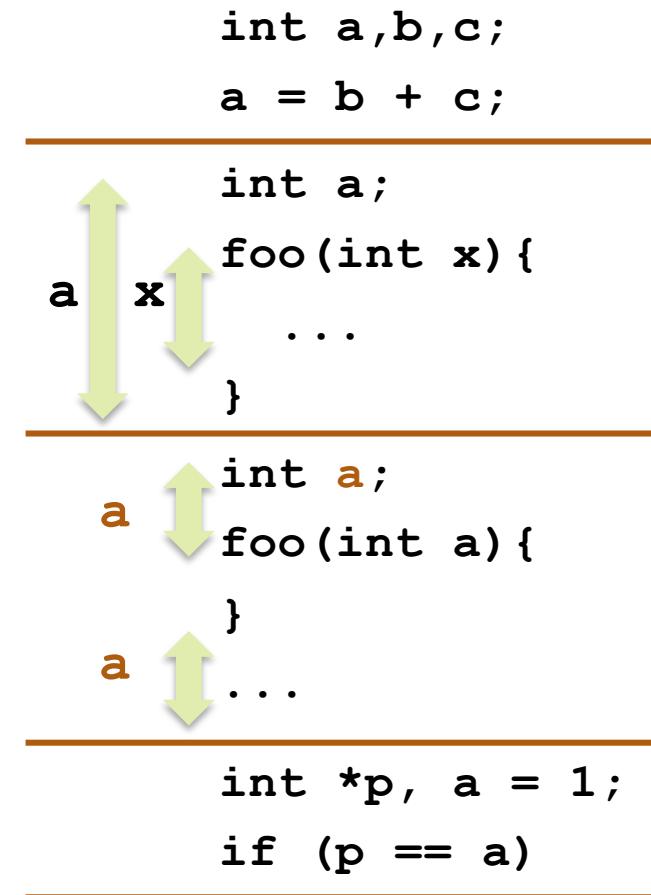
- Declaration, code and uses

```
int a,b,c;  
...  
a = b + c;
```

- Language abstraction: $L = \{wcw; w \in L((a|b)^*)\}$
 - Example: aaaabbbbabababbcaaabbababab
 - Context sensitive: Whether the second 'w' is valid depends on the first occurrence
 - It is more difficult to parse context sensitive grammars
- There are similar abstractions for checking function parameters

Semantic analysis: Sample tasks

- **Declaration:** Check that all identifiers were already declared
- **Validity:** Check whether in a portion of a program it is valid to use an identifier
- **Visibility:** Decide which valid “object” is visible when accessed
- **Type-consistency:** Check that the operands obey the type rules of operators



Approaches to semantic analysis

- Formal approach (apart from context-sensitive grammars): Attribute grammars
 - Automatic evaluation of attributes (declarative)
 - But: No globals & too many copies
 - Automatic generation of analyzer: ox
- Ad-hoc approach:
 - Associate code snippet with each production (yacc/bison)
 - Simple, functional but more rigid evaluation (actions given by parsing)
- Modern approach: Tree walk (used in LLVM)
 - Build the *abstract syntax tree*
 - Use virtual functions to walk through it (e.g., visitor pattern)
 - Need to associate rules to every tree-node type: more complicated

We will talk about this: Helps to understand the principles

4. Semantic Analysis

- Introduction
- Attribute grammars
- Applications

Attribute grammars: Intuition

□ Idea

- Syntax-controlled evaluation: semantic actions are associated to grammar rules
- Symbols ($N \cup T$) are given **attributes** (variables)
- According to the semantic actions concrete values are computed for the attributes

□ Example

Attributes	type	type	type
Semantic action	$\text{type} = \max(\text{type}, \text{type})$		
After parsing	$\text{float} = \max(\text{int}, \text{float})$		

Assume: Type
total order

...



float



int



char

Attribute grammars: Definition

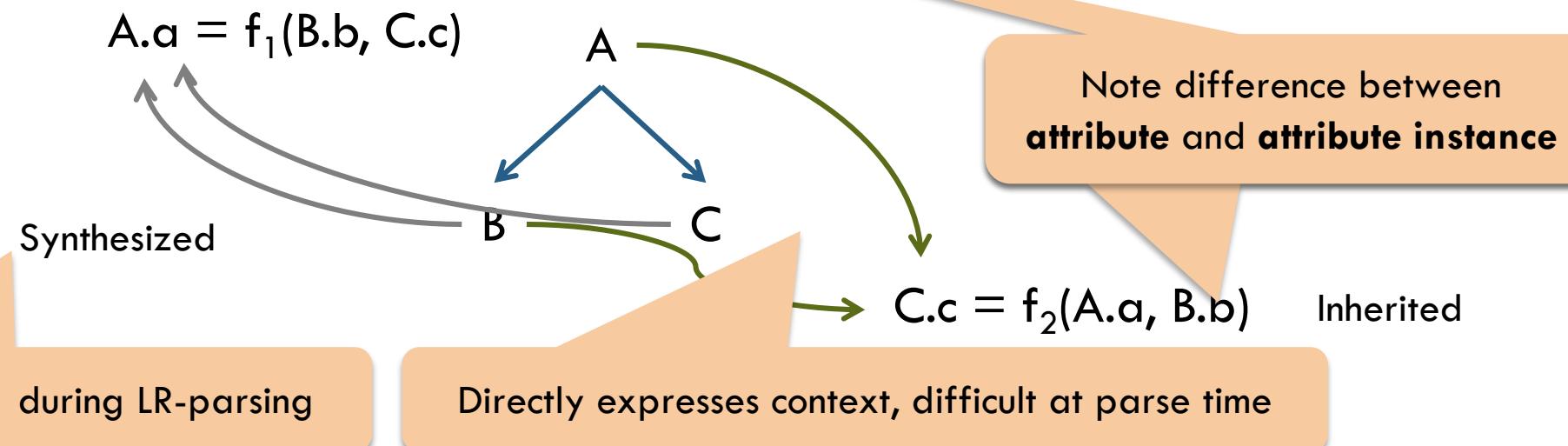
Def.: An **attribute grammar** is composed of a grammar CFG $G = (N, T, P, S)$ and

- Two disjoint sets of attributes SYN (synthesized) and INH (inherited)
 - SYN: Attributes that depend on children of the syntax tree
 - INH: Attributes that depend on parents and siblings of the syntax tree
- Notation: For $X \in (N \cup T)$, $\text{SYN}(X) \subseteq \text{SYN}$ and $\text{INH}(X) \subseteq \text{INH}$ denote the attributes associated with symbol X . $\text{ATTR}(X) = \text{SYN}(X) \cup \text{INH}(X)$
- A domain $D(a)$ for each attribute $a \in (\text{SYN} \cup \text{INH})$
- Semantic rules for every *computable* attribute in every production $p \in P$, i.e., given $X_0 \rightarrow X_1 X_2 \dots X_n$, with $X_0 \in N$, $X_{i=1,\dots,n} \in (N \cup T)$, a rule $a = f(c_1, \dots, c_m)$ with $a \in \text{SYN}(X_0) \cup \{\bigcup_{i=1}^n \text{INH}(X_i)\}$, and $c_{k=1,\dots,m} \in \bigcup_{i=0}^n \text{ATTR}(X_i)$

Attribute grammars: SYN & INH

□ Example:

- Grammar symbols A, B, C
- Attributes $a \in \text{ATTR}(A)$, $b \in \text{ATTR}(B)$, $c \in \text{ATTR}(C)$
- Sample production: $A \rightarrow BC$



Computing attributes: Example

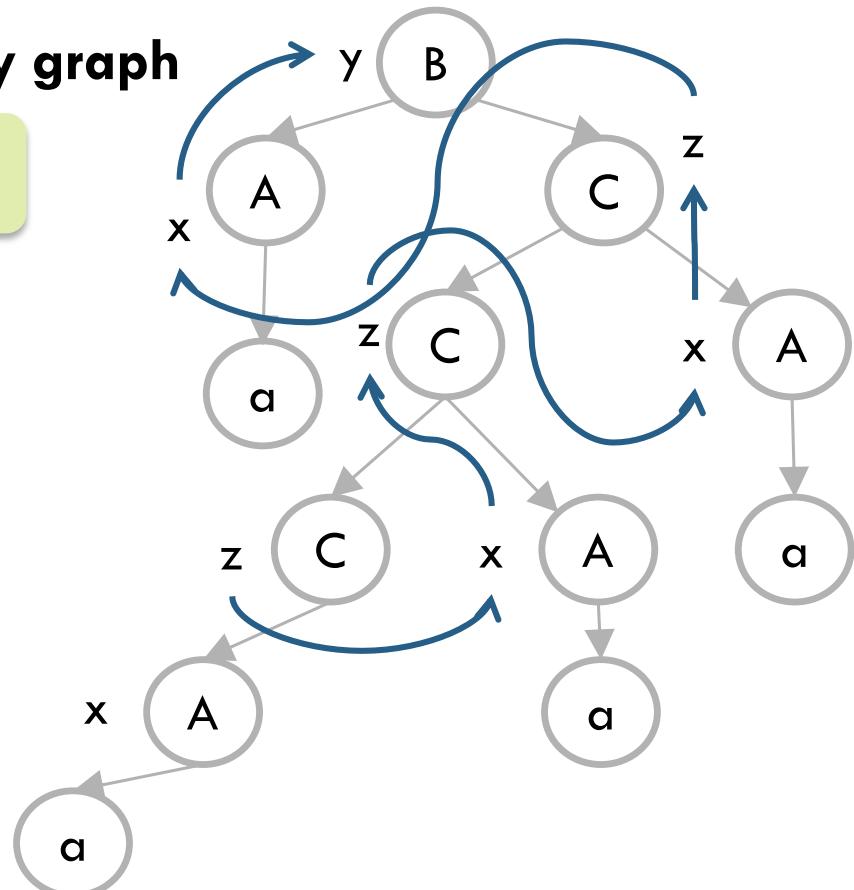
- ❑ Semantic rules and parse tree → dependency graph
 - ❑ Example

Production	Semantic rules
1. $B \rightarrow AC$	$A.x = C.z$ $B.y = A.x + 1$
2. $B \rightarrow A$	$A.x = 0$ $B.y = 1$
3. $C \rightarrow CA$	$A.x = C1.z$ $C0.z = A.x + 1$
4. $C \rightarrow A$	$A.x = 0$ $C.z = 1$
5. $A \rightarrow a$	

Input: aaaa

A.x: INH and
B.y & C.z: SYN

Numbers to distinguish instance

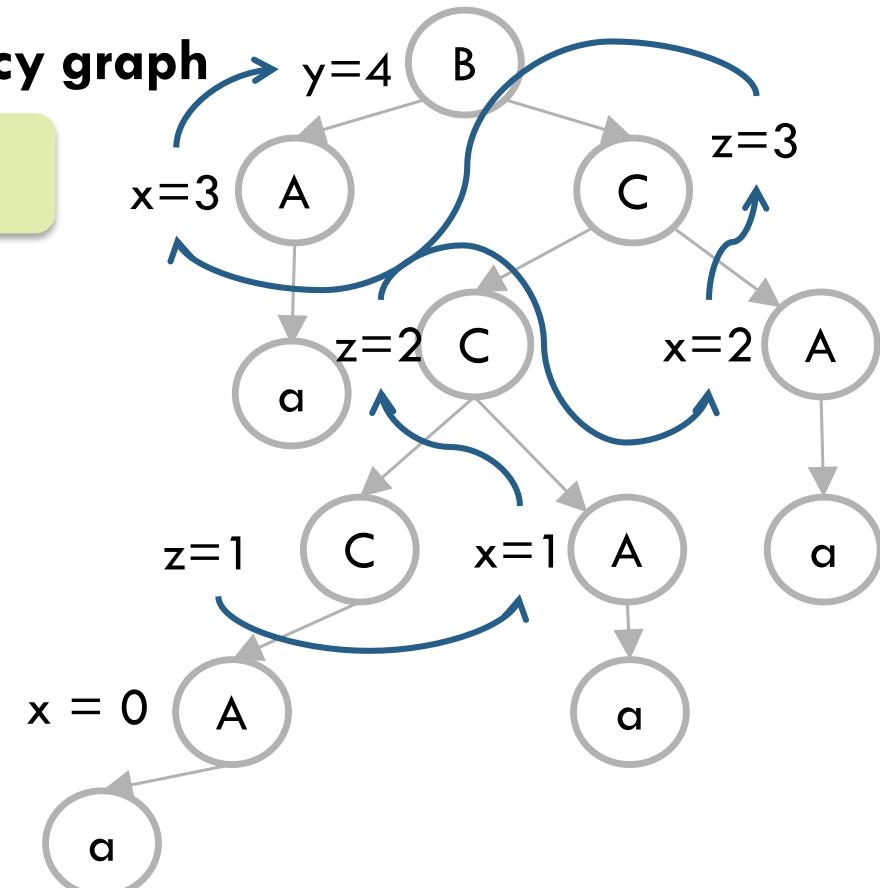


Computing attributes: Example

- Semantic rules and parse tree → dependency graph
- Example

Production	Semantic rules
1. $B \rightarrow AC$	$A.x = C.z$ $B.y = A.x + 1$
2. $B \rightarrow A$	$A.x = 0$ $B.y = 1$
3. $C \rightarrow CA$	$A.x = C.z$ $C.z = A.x + 1$
4. $C \rightarrow A$	$A.x = 0$ $C.z = 1$
5. $A \rightarrow a$	

Input: aaaa



Attribute dependency graph

Def.: An **attribute dependency graph** is a directed graph $G = (V, E)$ where nodes represent attribute instances and edges dependencies. There is an edge (c_i, a) for every attribute and every semantic rule $a = f(c_1, \dots, c_m)$

- Approach: Compute the attributes in the order given by a topological sort
- Cannot handle cyclic graphs
 - Largest class: Strongly non-circular grammars (SNC)
 - SNC-test: Polynomial-time (failing is not conclusive)

Topological sort

□ DFS - Algorithm

Input acyclic graph $G = (V, E)$

Output Node list L

$L = \{\}$

While $|L| \neq |V|$ **do**

 select $n \in (V - L)$

 append(n, L)

return L

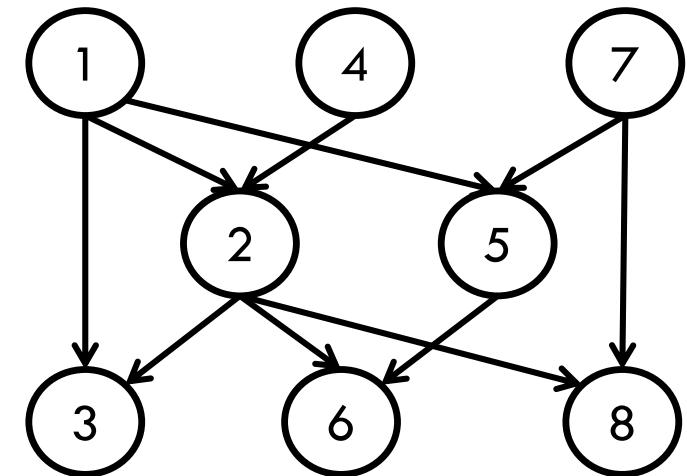
Function append(n, L)

foreach $(m, n) \in E, m \notin L$

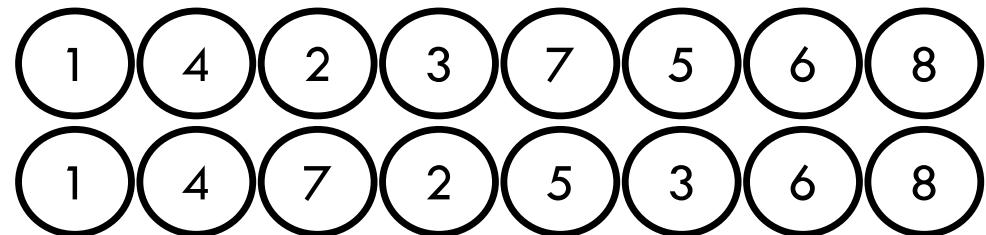
 append(m, L)

$L = \{L, n\}$

Example:



Possible orders:



4. Semantic Analysis

- Introduction
- Attribute grammars
- Applications



Application (examples) of attribute grammars



- In the coming slides
 - Synthetic example: Binary numbers
 - Abstract syntax trees
 - Label/goto correctness
 - Type checking

Binary numbers: CFG & attributes

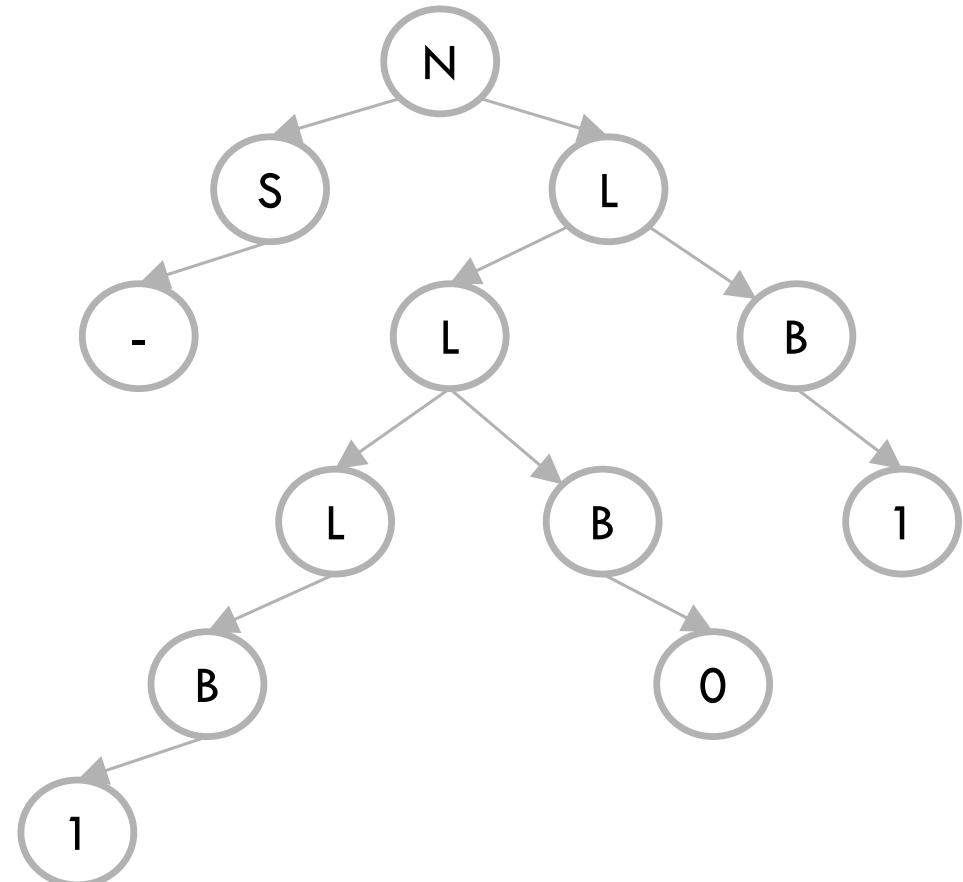
□ Binary representations of integers

CFG rules

1. $N \rightarrow S\ L$
2. $L \rightarrow L\ B$
3. $L \rightarrow B$
4. $B \rightarrow 0$
5. $B \rightarrow 1$
6. $S \rightarrow +$
7. $S \rightarrow -$

(N: number, S: sign,
L: list, B: bit)

Input: -101



Binary numbers: CFG & attributes

□ Binary representations of integers

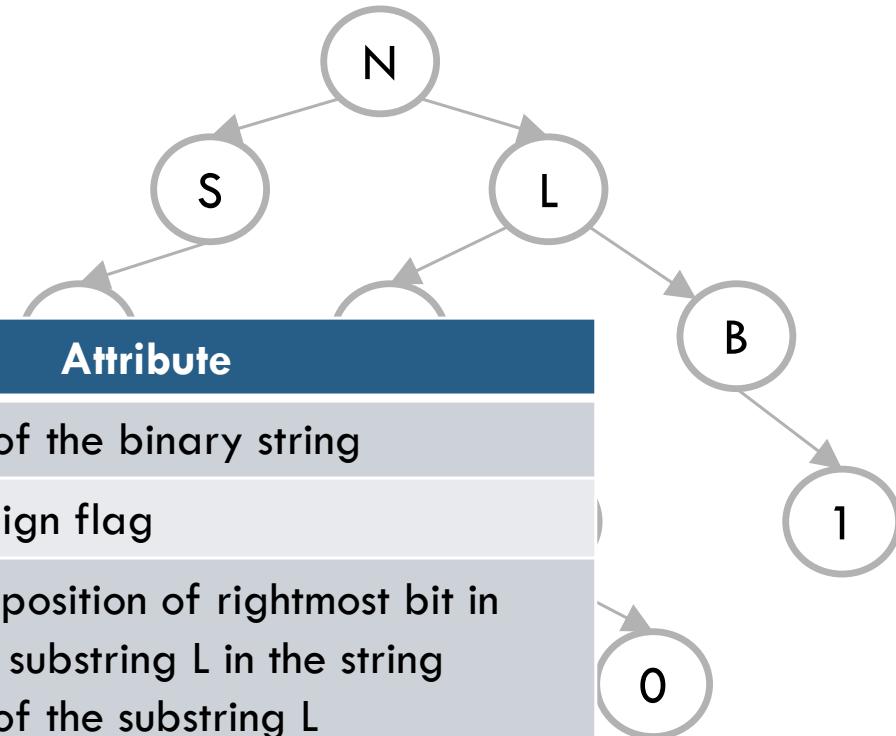
CFG rules

1. $N \rightarrow S\ L$
2. $L \rightarrow L\ B$
3. $L \rightarrow B$
4. $B \rightarrow 0$
5. $B \rightarrow 1$
6. $S \rightarrow +$
7. $S \rightarrow -$

(N: number, S: sign,
L: list, B: bit)

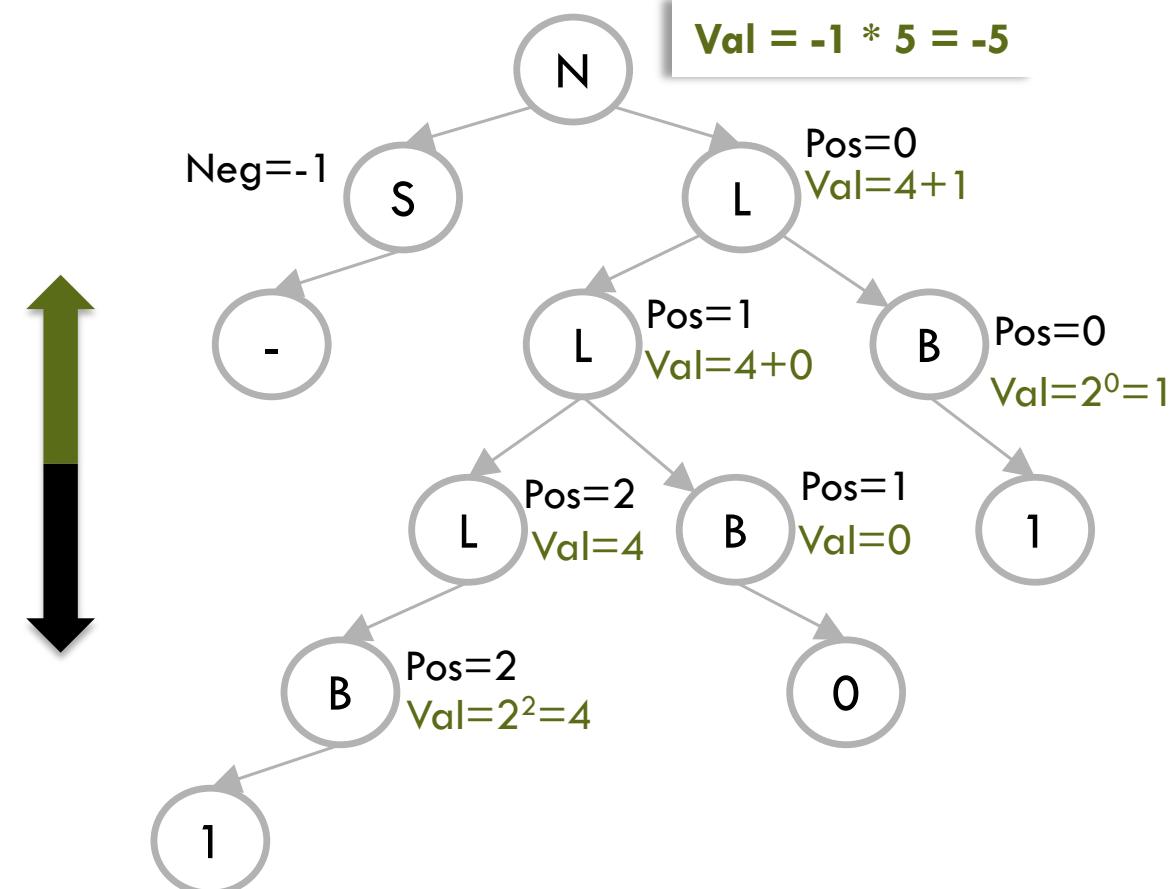
Input: -101

Symbol	Attribute
N	Val (Z): Value of the binary string
S	Neg ({-1,1}): Sign flag
L	Pos ($Z^+ \cup \{0\}$): position of rightmost bit in substring L in the string Val (Z): Value of the substring L
B	Pos ($Z^+ \cup \{0\}$): bit position in string Val (Z): Value of the bit (2^{Pos})



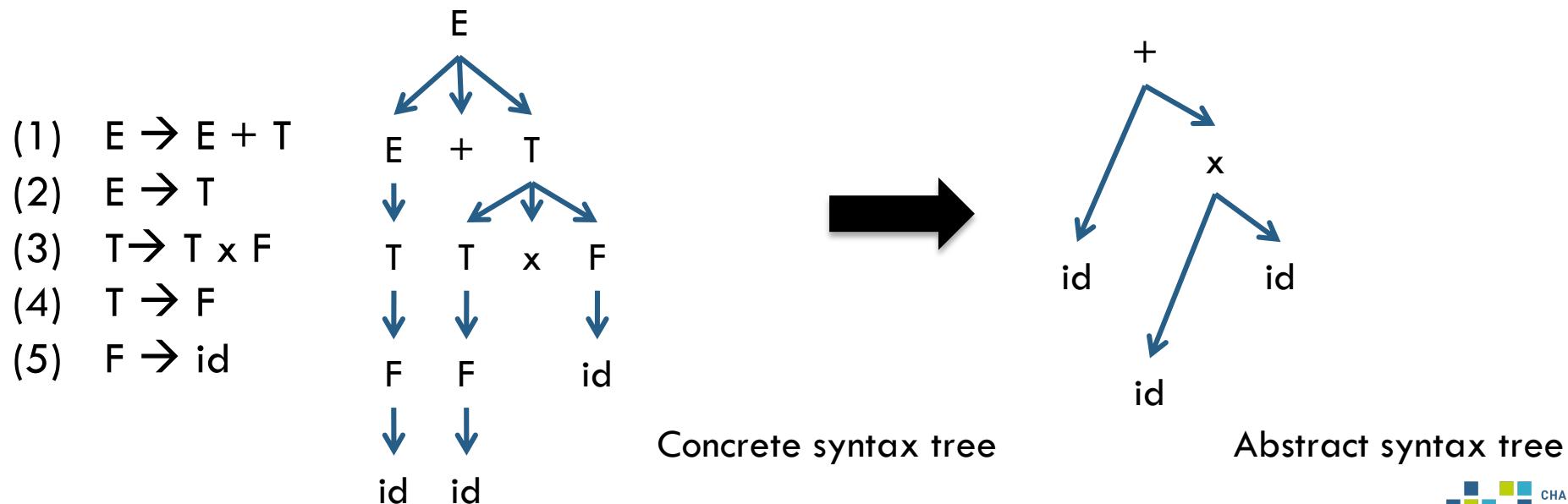
Binary numbers: Semantic rules

CFG rule	Semantic rule
$N \rightarrow S \ L$	$N.Val = S.Neg * L.Val$ $L.Pos = 0$
$L_0 \rightarrow L_1 \ B$	$L_1.Pos = L_0.Pos + 1$ $B.Pos = L_0.Pos$ $L_0.Val = L_1.Val + B.Val$
$L \rightarrow B$	$L.Val = B.Val$ $B.Pos = L.Pos$
$B \rightarrow 0$	$B.Val = 0$
$B \rightarrow 1$	$B.Val = 2^{B.Pos}$
$S \rightarrow +$	$S.Neg = +1$
$S \rightarrow -$	$S.Neg = -1$



Abstract syntax tree

- Parse tree = Concrete syntax tree
 - Contains symbols from the grammar
 - Symbols are irrelevant after parsing
- Abstract syntax tree (AST): Drop symbols and retain only what matters



AST: Attribute CFG

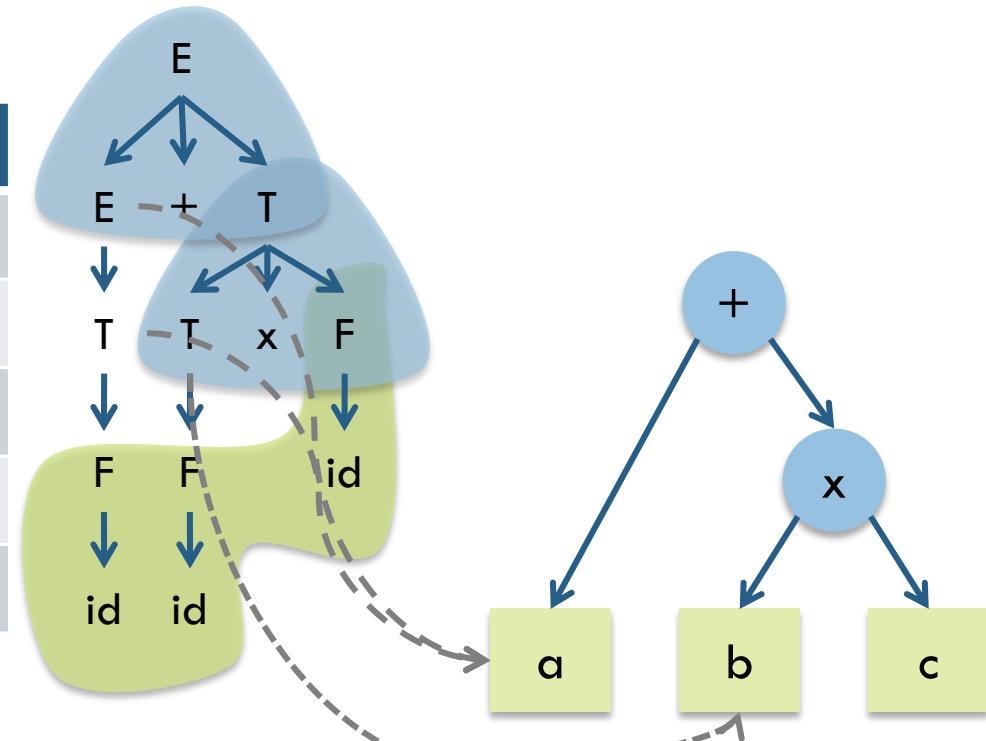
Symbol	Attribute	CFG rule	Semantic rule
E		$E_0 \rightarrow E_1 + T$	$E_0.ptr = \text{mknode}('+', E_1.ptr, T.ptr)$
T	ptr: pointer to AST-node	$E \rightarrow T$	$E.ptr = T.ptr$
F		$T_0 \rightarrow T_1 \times F$	$T_0.ptr = \text{mknode}('x', T_1.ptr, F.ptr)$
id	name: actual lexeme	$T \rightarrow F$	$T.ptr = F.ptr$
		$F \rightarrow id$	$F.ptr = \text{mkleaf}(id.name)$

- **mknode(op,left,right):** Help function to create internal AST nodes
- **mkleaf(name):** Help function to create leaf nodes for identifiers

AST: Example

Input: $a+b\times c$

CFG rule	Semantic rule
$E_0 \rightarrow E_1 + T$	$E_0.\text{ptr} = \text{mknode}('+', E_1.\text{ptr}, T.\text{ptr})$
$T \rightarrow T \times F$	$T_0.\text{ptr} = \text{mknode}(' \times ', T_1.\text{ptr}, F.\text{ptr})$
$E \rightarrow T$	$E.\text{ptr} = T.\text{ptr}$
$T \rightarrow F$	$T.\text{ptr} = F.\text{ptr}$
$F \rightarrow \text{id}$	$F.\text{ptr} = \text{mkleaf}(\text{id.name})$



Label/goto correctness

□ Grammar

- (1) $S' \rightarrow S$
- (2) $S \rightarrow L\ S$
- (3) $S \rightarrow G\ S$
- (4) $S \rightarrow \epsilon$
- (5) $L \rightarrow id\ :$
- (6) $G \rightarrow goto\ id;$

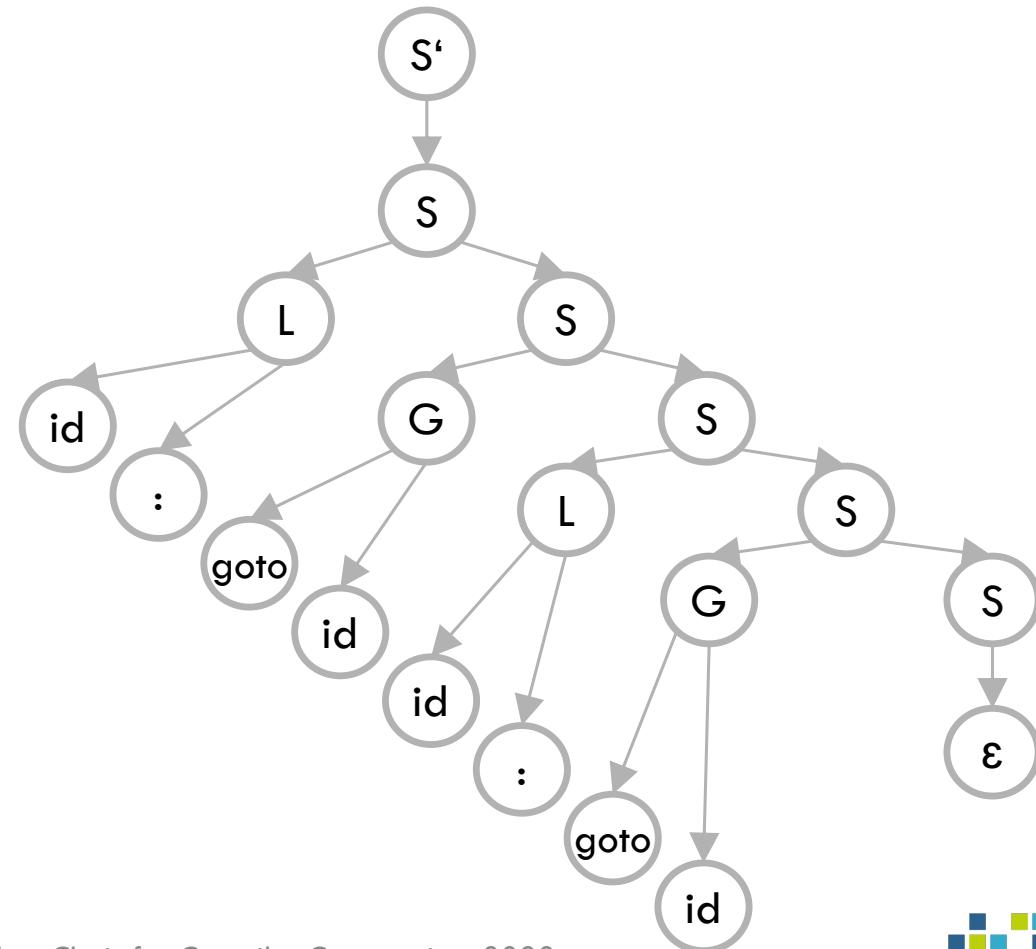
Input:

L1:

goto L2;

L2:

goto L1;



Label/goto: Attribute grammar

Symbol	Attribute
S'	S_{lab} : Set of all declared labels OK: flag \rightarrow program correct
S	S_{lab} : Set of all declared labels S'_{lab} : Set of declared labels below OK: flag \rightarrow statement correct
G	S_{lab} : Set of all declared labels OK: flag \rightarrow statement correct
L	name
id	name

CFG rule	Semantic rule
$S' \rightarrow S$	$S'.S_{lab} = S. S'_{lab}$ $S.S_{lab} = S'.S_{lab}$ $S.OK = S.OK$
$S_0 \rightarrow L S_1$	$S_0.S'_{lab} = S_1. S'_{lab} \cup \{L.name\}$ $S_1. S_{lab} = S_0.S_{lab}$ $S_0.OK = S_1.OK$
$S_0 \rightarrow G S_1$	$S_0.S'_{lab} = S_1. S'_{lab}$ $S_1. S_{lab} = S_0.S_{lab}$ $G.S_{lab} = S_0.S_{lab}$ $S_0.OK = S_1.OK \text{ AND } G.OK$
$S \rightarrow \epsilon$	$S.S'_{lab} = \emptyset, S.OK = \text{true}$
$L \rightarrow id :$	$L.name = id.name$
$G \rightarrow \text{goto id};$	$G.OK = id.name \in G.S_{lab}$

Label/goto correctness: Example

Input:

```
L1:  
  goto L2;  
L2:  
  goto L1;
```

$$S' \rightarrow S$$

$$S_0 \rightarrow L S_1$$

$$S_0 \rightarrow G S_1$$

$$S_0.S'_{lab} = S_1.S'_{lab}$$

$$S_1.S_{lab} = S_0.S_{lab}$$

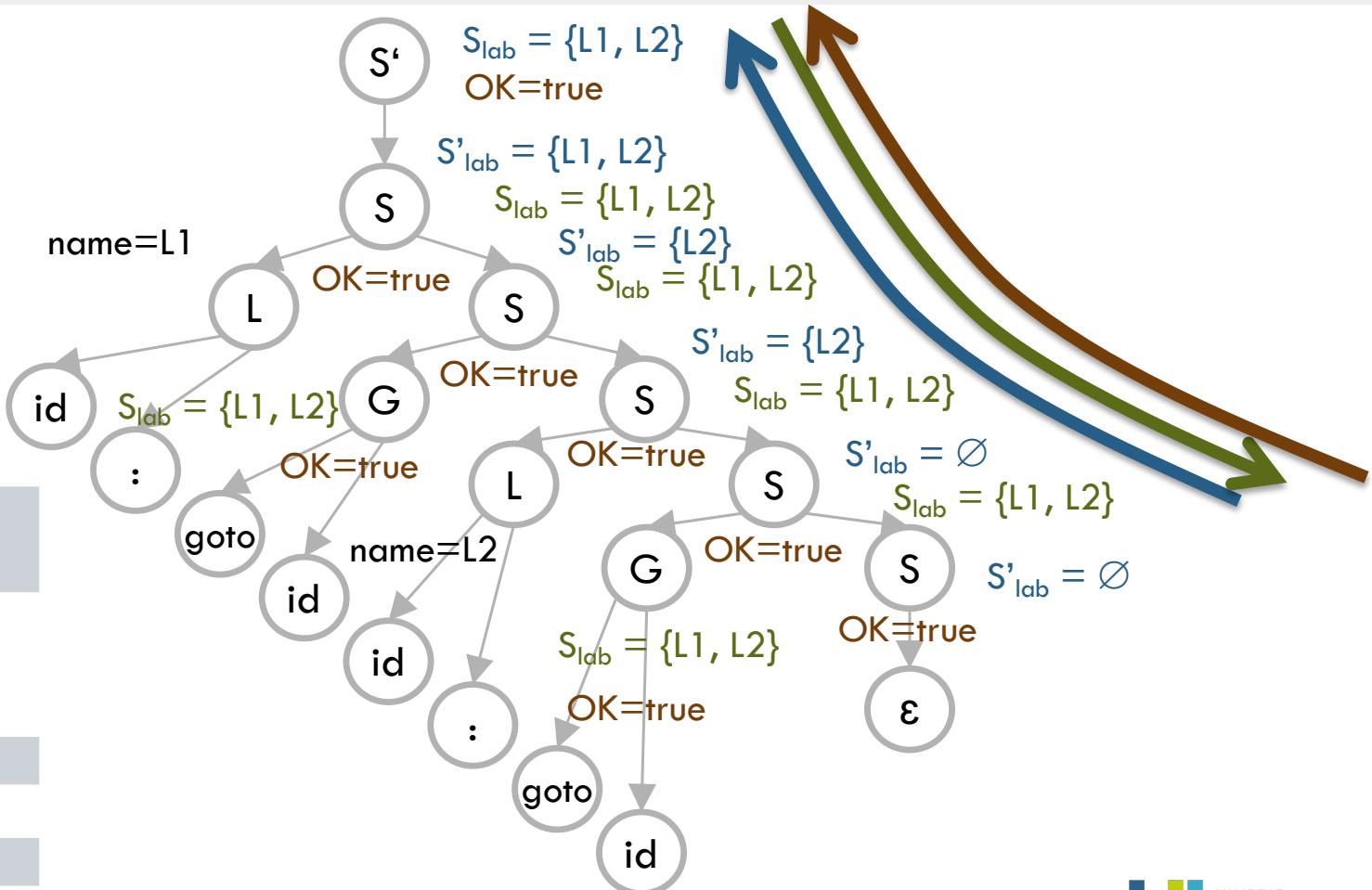
$$G.S_{lab} = S_0.S_{lab}$$

$$S_0.OK = S_1.OK \wedge$$

$S \rightarrow \epsilon$ $S.S'_{|gb} = \emptyset, S.OK = \text{true}$

$L \rightarrow id : L.name = id.name$

$$G \rightarrow \text{goto } id; \quad G.\text{OK} = id.\text{name} \in G.S_{\text{lab}}$$



Approaches to semantic analysis: Recall

- Formal approach (apart from context-sensitive grammars): Attribute grammars
 - Automatic evaluation of attributes (declarative)
 - **But: No globals & too many copies**
 - Automatic generation of analyzer: ox
- Ad-hoc approach:
 - Associate code snippet with each production (yacc/bison)
 - Simple, functional but more rigid evaluation (actions given by parsing)
- Modern approach: Tree walk (used in LLVM)
 - Build the *abstract syntax tree*
 - Use virtual functions to walk through it (e.g., visitor pattern)
 - Need to associate rules to every tree-node type: more complicated

Type-checking: Introduction

□ Recall motivation

```
int x,y;  
float foo(float y) {  
    return x + y; }
```

$$E \rightarrow E + E$$

Attributes	type	type	type
Semantic action	type = max (type, type)		
After parsing	float = max (int, float)		

□ What the compiler needs to know

- Scope information: Which ‘x’ is meant in “x + y”
- Type information (& other attributes) of the variables
- Type checking rules for ‘+’

Use so-called **symbol tables**

Use so-called **type systems**

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x, y, z);
4:     {
5:         int x, z;
6:         z = y;
7:         x = z;
8:         {
9:             int y = x;
10:            printf("%d,%d,%d\n", x, y, z);
14:        }
15:        printf("%d,%d,%d\n", x, y, z);
16:    }
17:}
```

Symbol table

x	0
z	1

Symbol

Line number (& attributes in reality)

Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x, y, z);
4:     {
5:         int x, z;
6:         z = y;
7:         x = z;
8:         {
9:             int y = x;
10:            printf("%d,%d,%d\n", x, y, z);
14:        }
15:        printf("%d,%d,%d\n", x, y, z);
16:    }
17:}
```

Symbol table		
x	0	
z	1	
x	2	
y	2	

New scope

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x@2, y@2, z@1);
4:     {
5:         int x, z;
6:         z = y;
7:         x = z;
8:         {
9:             int y = x;
10:            printf("%d,%d,%d\n", x, y, z);
14:        }
15:        printf("%d,%d,%d\n", x, y, z);
16:    }
17:}
```

Symbol table		
x		0
z		1
x		2
y		2

Look-up ↑

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x@2, y@2, z@1);
4:     {
5:         int x, z;
6:         z@5 = y@2;
7:         x@5 = z@5;
8:         {
9:             int y = x;
10:            printf("%d,%d,%d\n", x, y, z);
14:        }
15:        printf("%d,%d,%d\n", x, y, z);
16:    }
17:}
```

Symbol table		
x		0
z		1
x		2
y		2
x		5
z		5

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x@2, y@2, z@1);
4:     {
5:         int x, z;
6:         z@5 = y@2;
7:         x@5 = z@5;
8:         {
9:             int y = x@5;
10:            printf("%d,%d,%d\n", x@5, y@9, z@5);
14:        }
15:        printf("%d,%d,%d\n", x, y, z);
16:    }
17:}
```

Symbol table		
x		0
z		1
x		2
y		2
x		5
z		5
y		9

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x@2, y@2, z@1);
4:     {
5:         int x, z;
6:         z@5 = y@2;
7:         x@5 = z@5;
8:         {
9:             int y = x@5;
10:            printf("%d,%d,%d\n", x@5, y@9, z@5);
11:        }
12:        printf("%d,%d,%d\n", x@5, y@2, z@5);
13:    }
14: }
```

Symbol table		
x		0
z		1
x		2
y		2
x		5
z		5

Symbol tables: Intuition

```
0: int x = 137;
1: int z = 42;
2: int MyFunction(int x, int y) {
3:     printf("%d,%d,%d\n", x@2, y@2, z@1);
4:     {
5:         int x, z;
6:         z@5 = y@2;
7:         x@5 = z@5;
8:         {
9:             int y = x@5;
10:            printf("%d,%d,%d\n", x@5, y@9, z@5);
11:        }
12:        printf("%d,%d,%d\n", x@5, y@2, z@5);
13:    }
14: }
```

Symbol table

x	0
z	1

Symbol tables

- **Symbol table:** Mapping from a name to the object the name refers to
 - Filled during semantic analysis (with info from lexical and syntax analysis)
 - *Stack of maps:* Each map represents a scope
- **Object attributes:** Name, type, size, alignment, ... (see next slide)
- **Operations**
 - Push scope: Enter a new scope
 - Pop scope: Leave a scope and delete its declarations
 - Insert symbol
 - Look-up symbol

Symbol tables: Attributes

- Attributes depend on the language, e.g., for C:

Attribute	Representation	
Name	String	Compiler needs to insert <i>memcpy</i> when assigning structs <code>struct { ... } S, T; S = T;</code>
Type	Type-expression, e.g., int, float, array(100,int), struct(float,...), (float, short) → pointer(double), ...	Compiler has to know what are correct addresses in HW (0x0, 0x4, 0x8, ...)
Size	Integer	
Alignment	Integer	
Qualifier	{const, volatile}	
Storage class	{extern, static, auto, register}	Recognize illegal assignment, e.g., <code>const int x = 1; x = 2;</code>

Symbol tables: Implementation

- Need for speed: n insertions and m look-ups ($m > n$)
 - On every insert: Check for existence
- Linked lists
 - Insert: Existence check $O(n)$, actual insertion $O(1)$
 - Look-up: $O(n)$
 - Time complexity $O(n^2 + mn)$
- Hash tables (assumption: open hashing, balanced)
 - Insert: Existence check $O(1)$, actual insertion $O(1)$
 - Look-up: $O(1)$
 - Time complexity $O(n + m)$

Type-checking: Reloaded

□ Recall motivation

```
int x,y;  
float foo(float y) {  
    return x + y; }
```

$$E \rightarrow E + E$$

Attributes	type	type	type
Semantic action	type = max (type, type)		
After parsing	float	= max (int, float)	

□ What the compiler needs to know

- Scope information: Which ‘x’ is meant in “x + y”
- Type information (& other attributes) of the variables

- Type checking rules for ‘+’

Type-checking: Intuition

- Include checking rules in the semantic actions of the attribute grammar
- Example

- (1) $S \rightarrow D ; E$
- (2) $D \rightarrow D ; D \mid T \text{id}$
- (3) $T \rightarrow \text{int} \mid \text{float} \mid T[\text{num}] \mid T^*$
- (4) $E \rightarrow \text{num} \mid \text{num.num} \mid \text{id} \mid E + E \mid E[E] \mid *E$

```
int a;
int[10] b;
int* c;
A + b[1] + *c + 42;
```

Type-checking: Sample attribute grammar

CFG rule	Semantic rule
$D \rightarrow T \text{ id}$	
$T \rightarrow \text{int} \mid \text{float}$	
$T_0 \rightarrow T_1 [\text{num}]$	
$T_0 \rightarrow T_1^*$	
$E \rightarrow \text{num} \mid \text{num.num}$	
$E \rightarrow \text{id}$	

Type-checking: Sample attribute grammar

CFG rule	Semantic rule
$D \rightarrow T \text{id}$	$\text{symtab}[\text{id.name}] = T.\text{type}$
$T \rightarrow \text{int} \mid \text{float}$	$T.\text{type} = \text{int} \mid \text{float}$
$T_0 \rightarrow T_1 [\text{num}]$	$T_0.\text{type} = \text{array}(\text{num.val}, T_1.\text{type})$
$T_0 \rightarrow T_1^*$	$T_0.\text{type} = \text{ptr}(T_1.\text{type})$
$E \rightarrow \text{num} \mid \text{num.num}$	$E.\text{type} = \text{int} \mid \text{float}$
$E \rightarrow \text{id}$	$E.\text{type} = \text{symtab}[\text{id.name}]$

Type-checking: Sample attribute grammar (2)

CFG rule	Semantic rule
$E_0 \rightarrow E_1 + E_2$...
$E_0 \rightarrow E_1[E_2]$	
$E_0 \rightarrow *E_1$	

Type-checking: Sample attribute grammar (2)

CFG rule	Semantic rule
$E_0 \rightarrow E_1 + E_2$	$E_0.type = \text{if}$ $(E_1.type == \text{int} \&\& E_2.type == \text{int}) \text{ then int else if}$ $(E_1.type == \text{float} \&\& E_2.type == \text{float}) \text{ then float else errortype}$
$E_0 \rightarrow E_1[E_2]$	$E_0.type = \text{if}$ $(E_1.type == \text{array}(s,t) \&\& E_2.type == \text{int}) \text{ then t else errortype}$
$E_0 \rightarrow *E_1$	$E_0.type = \text{if}$ $E_1.type == \text{pointer}(t) \text{ then t else errortype}$

Type-checking as proofs: Inference rules

- Define inference rules to determine the type
- Notation

$$\frac{\begin{array}{c} S \vdash e_1: \text{int} \\ S \vdash e_2: \text{int} \end{array}}{S \vdash e_1 + e_2: \text{int}}$$

Read: If we can show that e_1 and e_2 are of **type int** in **scope S**, then we can show that $e_1 + e_2$ has **type int** as well

- Other examples

$$\frac{\begin{array}{c} S \vdash e_1: T[] \\ S \vdash e_2: \text{int} \end{array}}{S \vdash e_1[e2]: T}$$

f is an identifier for a function in scope S

f has type $(T_1, \dots, T_n) \rightarrow U$

$S \vdash e_i: T_i$ for $1 \leq i \leq n$

$$\frac{}{S \vdash f(e_1, \dots, e_n): U}$$

Partial order for types

- Introduce partial order of types (for example for class hierarchies)
- Given two types A and B, we say $A \leq B$ if A is convertible to B
 - $A \leq A$
 - If $A \leq B$ and $B \leq C \rightarrow A \leq C$
 - If $A \leq B$ and $B \leq A \rightarrow A = B$
- Example: Rule for comparison

$$\frac{\begin{array}{c} S \vdash e_1 : T \\ S \vdash e_2 : T \\ T \text{ is primitive} \end{array}}{S \vdash e_1 == e_2 : \text{bool}} \quad \longrightarrow \quad \frac{\begin{array}{c} S \vdash e_1 : T_1 \\ S \vdash e_2 : T_2 \\ T_1 \leq T_2 \text{ or } T_2 \leq T_1 \end{array}}{S \vdash e_1 == e_2 : \text{bool}}$$

Type-checking: Example



```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

Adapted from: <http://web.stanford.edu/class/archive/cs/cs143/cs143.1128/>

Type-checking: Example

```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

a is identifier

a is a variable in scope S of type T

S ⊢ a: T

Facts

S ⊢ x: int

S ⊢ y: int

S ⊢ z: int

Type-checking: Example

```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

$S \vdash e_1 : T_1$

$S \vdash e_2 : T_2$

$T_1 \leq T_2 \text{ or } T_2 \leq T_1$

$S \vdash e_1 == e_2 : \text{bool}$

Facts

$S \vdash x : \text{int}$

$S \vdash y : \text{int}$

$S \vdash z : \text{int}$

$S \vdash x == y : \text{bool}$

Type-checking: Example

```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

$S \vdash e_1 : \text{int}$

$S \vdash e_2 : \text{int}$

$S \vdash e_1 + e_2 : \text{int}$

Facts

$S \vdash x : \text{int}$

$S \vdash y : \text{int}$

$S \vdash z : \text{int}$

$S \vdash x == y : \text{bool}$

$S \vdash x + y : \text{int}$

Type-checking: Example

```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

$S \vdash e_1 : T_1$

$S \vdash e_2 : T_2$

$T_1 \leq T_2 \text{ or } T_2 \leq T_1$

$S \vdash e_1 < e_2 : \text{bool}$

Facts

$S \vdash x : \text{int}$

$S \vdash y : \text{int}$

$S \vdash z : \text{int}$

$S \vdash x == y : \text{bool}$

$S \vdash x + y : \text{int}$

$S \vdash x + y < z : \text{bool}$

Type-checking: Example

```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

$S \vdash e_1 : \text{int}$

$S \vdash e_2 : \text{int}$

$S \vdash e_1 > e_2 : \text{bool}$

> Error: Cannot compare int and bool

Facts

$S \vdash x : \text{int}$

$S \vdash y : \text{int}$

$S \vdash z : \text{int}$

$S \vdash x == y : \text{bool}$

$S \vdash x + y : \text{int}$

$S \vdash x + y < z : \text{bool}$

Type-checking: Example

```
int x, y, z;  
if (((x == y) > 5 && x + y < z)) { /* ... */ }
```

$S \vdash e_1 : \text{bool}$

$S \vdash e_2 : \text{bool}$

$S \vdash e_1 \&& e_2 : \text{bool}$

```
> Error: Cannot compare int and bool  
> Error: Unrecognized ??? && bool
```

Facts

$S \vdash x : \text{int}$

$S \vdash y : \text{int}$

$S \vdash z : \text{int}$

$S \vdash x == y : \text{bool}$

$S \vdash x + y : \text{int}$

$S \vdash x + y < z : \text{bool}$

Error type (bottom type) to prevent errors from propagating (analog to error token in parsing)

Where are we?

1. Introduction
2. Lexical analysis
3. Syntax analysis
4. Semantic analysis
5. Intermediate representation
6. Control & data-flow analysis
7. IR optimization
8. Target architectures
9. Code selection
10. Register allocation
11. Scheduling
12. Advanced topics