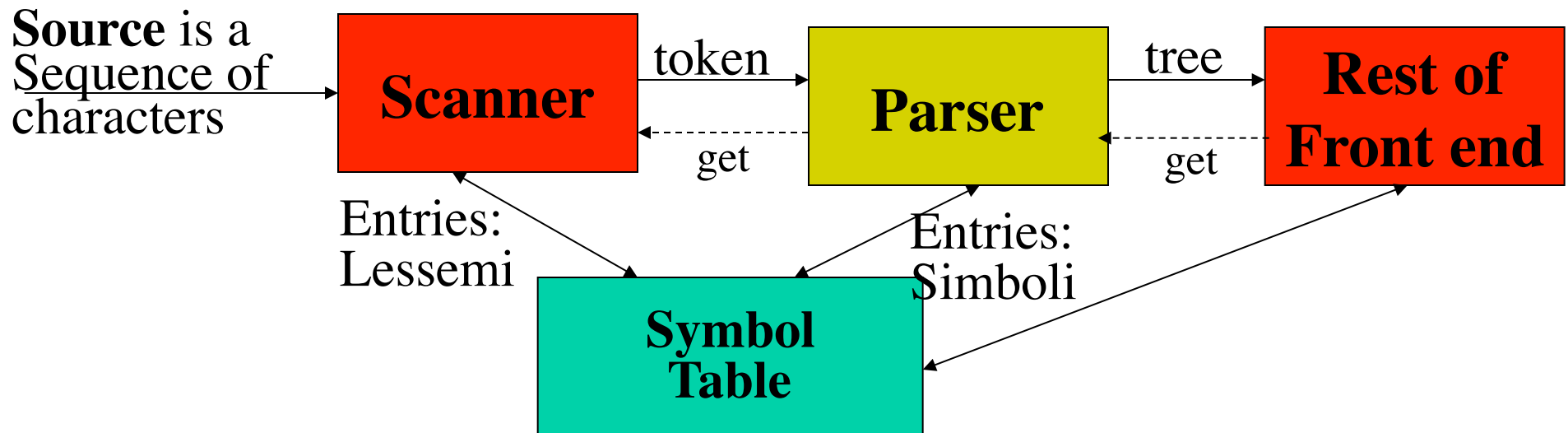# SYNTACTIC ANALYSIS

- **How to define Syntax;**
- **How to do Analysis;**
- **At what extent Analysis can be done in Linear Time;**
- **How to build (Linear) Parser Generators;**
- **Relationship between Analysis and Internal Representation (Tree Representation)**

# One Pass Structure (Two phase pipeline)

**Source** is a Sequence of characters → **Scanner**

Scanner — token → **Parser**

Scanner ← get -- Parser

Parser — tree → **Rest of Front end**

Parser ← get -- Rest of Front end

Entries: Lessemi

Entries: Simboli

**Symbol Table**

Syntactic Analysis (Parser) is driven from
Semantics Analysis which is asking for visiting a subtree not built yet

# How to define Syntax

- Syntactic Analysis and Syntactic Languages
- Syntactic Languages and Grammars
- Classification of Grammars
- Classification of Languages
- Foundations: Derivation, Sentential Form, Ambiguity, Tarski's Fixpoint Iteration

# Syntactic Analysis

- It scans sequences of tokens to check for *phrase structures* that belong to the Syntax of the Language

- Syntax, just like Lexics, is expressed by a Language: The Syntactic Language

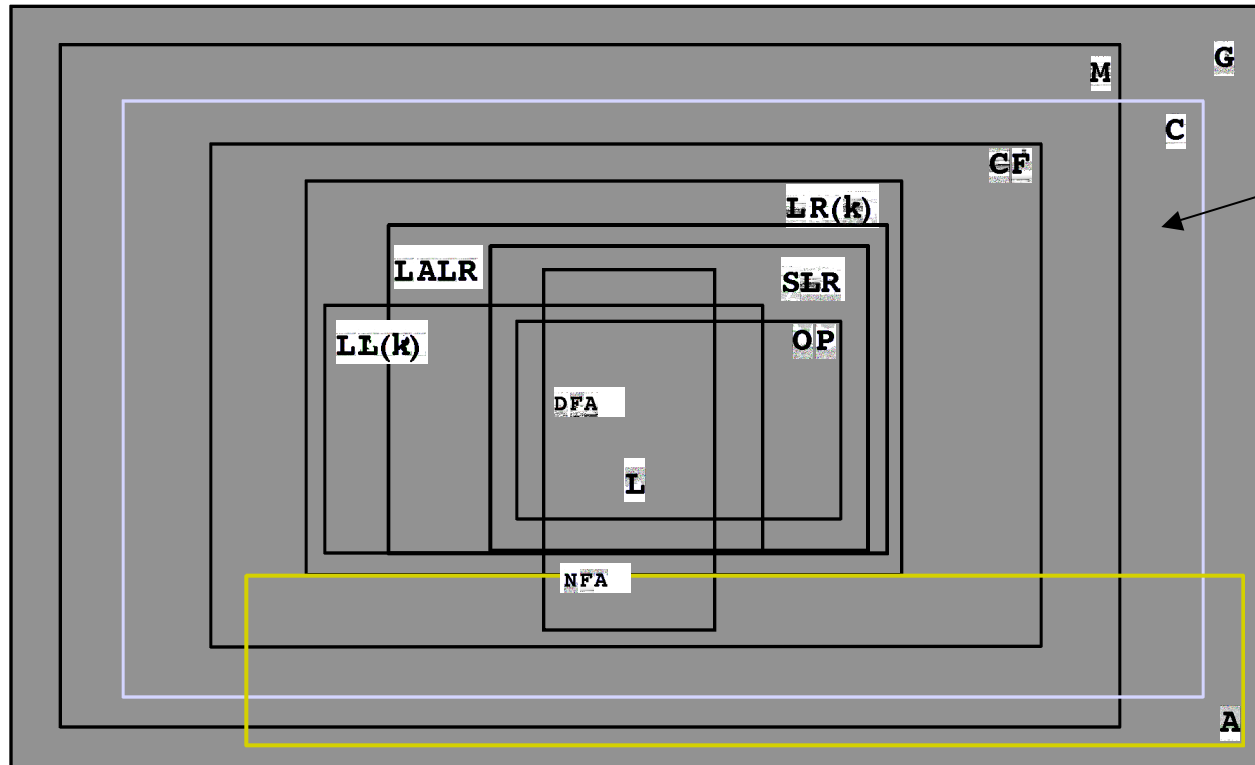- Syntactic Langs are much more complicated than Lexical ones

$\{u^n v^n \mid n \geq 0\}$   has been proved not be regular

**but**

num+(3-((id\*id)+num)/id) $\in \{(\alpha u)^n (\alpha v)^n \mid n \geq 0, \alpha \in L, u="(", v=")"\}$

# Grammar Classification (Chomsky)



**Grammars Inclusion**

Non-monotone **C**

**Kinds of Grammar [Defined Language Features]**

Diagram labels (nested boxes from outermost to innermost): G, M, C, CF, LR(k), LALR, SLR, LL(k), OP, DFA, L, NFA, A

**G=General** [Recursive Enumerable but Non-Recursive -$\{u^n v^{akermann(n)}\}$]

**A=Ambiguous**

**M=Monotone** [Recursive Languages - $\{u^n v^{n!}\}$]

**C=Contextual** [$\{u^n v^n z^n\}$]

**LR(k)=Context-Free** [$\{u^n v^n\}$]

**LALR(k)=Context-Free** [$\{u^n v^n\}$]

**SLR(k)=Simple Left-to-right rightmost reversed** [Viable-Prefix; Bottom-Up/k symbols$\{u^n v^n\}$]

**LL(K)=Leftmost-Left Left-to-right** [Predictive; Top-Down/k symbols$\{u^n v^n\}$]

**OP=Operator-Precedence** [$\{u^n v^n\}$]
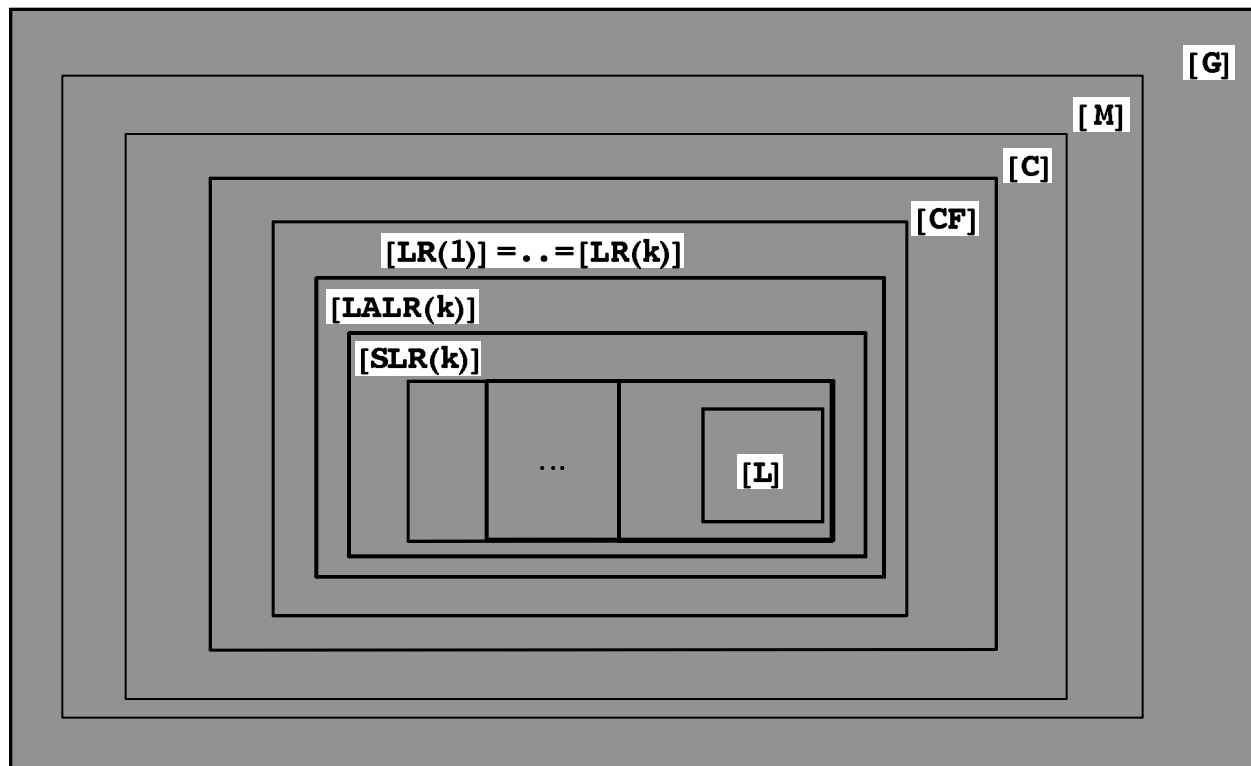
**L=Linear** [Recursive Grammars; Regular Languages]

**DFA= --** [Regular Grammars/Expressions; Regular Languages]

**NFS= --** [Regular Grammars/Expressions; Regular Languages]

5

# Language Classification



**Language Inclusion**

[G]
[M]
[C]
[CF]
[LR(1)] =..= [LR(k)]
[LALR(k)]
[SLR(k)]
...
[L]

[G] = Recursively Enumarable Languages
[M] = Recursive Languages
[C] = Contextual Languages: $\{u^n v^n z^n \mid n \geq 0\}$
[CF] = Context-Free Languages: $\{u^n v^m z^k \mid n,m,k \geq 0 \text{ and } (n=m \text{ or } m=k)\}$
[LR(k)] = LR/k symbols Languages: $\{u^m v^n \mid m > n \geq 0\}$
[LALR(k)] = LALR/k symbols Languages
[SLR(k)] = SLR/k symbols Languages
[LL(k)] = LL/k symbols Languages: $\{u^n v^n \mid n \geq 0\}$
[L] = Regular Languages: $\{u^n v^m \mid n \geq 0, m \geq 0\}$

# Definitions: Derivation, SF

Let $G = \langle V, \Sigma, s \in V, P \rangle$

**Derivation** is a binary relation $\Rightarrow_G$ su $(\Sigma \cup V)^* \times (\Sigma \cup V)^*$

$$\alpha A \beta \Rightarrow \alpha \gamma \beta \quad sse \quad A ::= \gamma \in P$$

Subscript, G, is omitted, in $\Rightarrow_G$, when the grammar G is clearly stated from the context

$\Rightarrow^*$: Transitive and Reflexive Closure of $\Rightarrow$
- $\alpha \Rightarrow^* \alpha$
- if $\alpha_1 \Rightarrow ... \Rightarrow \alpha_n$
  then $\alpha_1 \Rightarrow^* \alpha_n$

**Sentential form** of G
$$SF = \{\gamma \mid s \Rightarrow^* \gamma\}$$

$\Rightarrow^+$: Transitive Closure of $\Rightarrow$
if $\alpha_1 \Rightarrow \alpha_2 \Rightarrow ... \Rightarrow \alpha_n$ and $\alpha_1 \neq \alpha_2 \neq ... \neq \alpha_n$
allora $\alpha_1 \Rightarrow^+ \alpha_n$

# L(G): Language Generated by a Grammar

Let **G = <V,Σ,s∈V,P>**

$$L(G)=\{w \in \Sigma^* \mid s =>^+ w\}$$
(where **=>** is **=>$_G$** )

**Example: Let G below**

*p1*: E::= E+E
*p2*: E::= E*E
*p3*: E::= id

**Then**

**id+id∈L(E)**

**A proof (Lefmost Derivation):**
E =>$_{(p1)}$ E+E =>$_{(p3)}$ id+E =>$_{(p3)}$ id+id

**A different proof (Rightmost Derivation):**
E =>$_{(p1)}$ E+E =>$_{(p3)}$ E+id =>$_{(p3)}$ id+id

# Ambiguous Grammars are
# Bad Definitions for Lang. Syntax

Example: Let G below

*p1*: E::= E+E
*p2*: E::= E*E
*p3*: E::= id      Then

**id+id*id∈L(E)**

**A proof (Lefmost Derivation):**
E =>$_{(p1)}$ E+E =>$_{(p3)}$ id+E =>$_{(p2)}$ id+E*E =>$_{(p3)}$ id+id*E =>$_{(p3)}$ id+id*id
**A different proof (another Leftmost Derivation):**
E =>$_{(p2)}$ E*E =>$_{(p1)}$ E+E*E =>$_{(p1)}$ id+E*E =>$_{(p1)}$ id+id*E =>$_{(p1)}$ id+id*id

**Different Leftmost (Rightmost) Derivations lead to different Parse Trees**

9

# Derivations (on the P-tree domain)

Let $G = \langle V, \Sigma, s \in V, P \rangle$
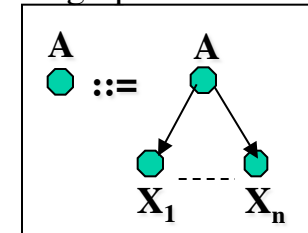
$(\Sigma \cup V)_T^*$

Smallest set such that, $\forall a \in \Sigma \cup V$:
- **leaf: $\langle [a,-] \rangle \in (\Sigma \cup V)_T^*$**
- $\forall \langle t1,\ldots,tn \rangle \in (\Sigma \cup V)_T^*$
    - **Tree: $\langle [a,\langle t1,\ldots,tn \rangle] \rangle \in (\Sigma \cup V)_T^*$**
    - **Forest: $\langle t1,\ldots,tn,u1,\ldots,um \rangle \in (\Sigma \cup V)_T^*$, $\forall \langle u1,\ldots,um \rangle \in (\Sigma \cup V)_T^*$**

A graphical view



## Productions on $(\Sigma \cup V)_T^*$

$\mathbf{A::=X_1\ldots X_n} \in P$   **sse**   $\langle [A,-] \rangle ::= \langle [A,\langle\langle [X_1,-] \rangle,\ldots,\langle [X_n,-] \rangle\rangle] \rangle \in P_T$
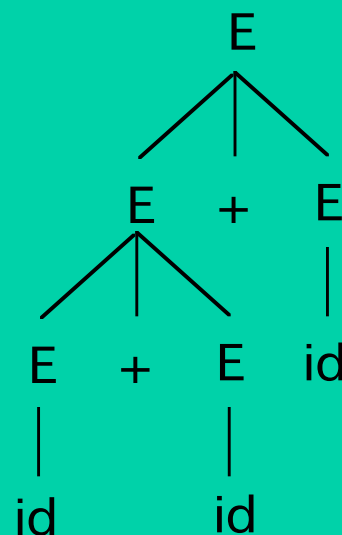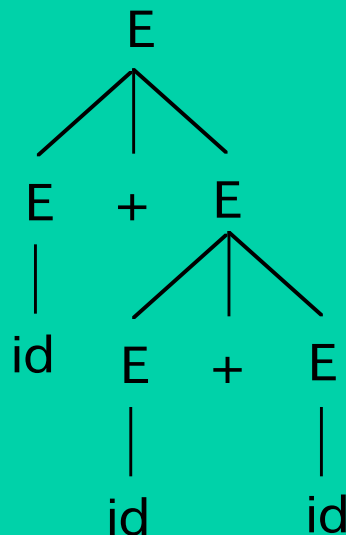
## Relation $\Rightarrow$ on $(\Sigma \cup V)_T^*$ x $(\Sigma \cup V)_T^*$

$\alpha A \beta \Rightarrow \alpha \gamma \beta$   **sse**   $A::=\gamma \in P_T$

# Ambiguous Grammars
# A Graphical View

**A different proof of ambiguity that uses: The *last trees* of two different Tree-derivations**



*p1*: E::= E+E
*p2*: E::= E*E
*p3*: E::= id

Exercise: Show the same but using leftmost

(rightmost) derivations

# From Grammars to Languages
## A methodology for finding L(G), given G

$$A_0 ::= e_0$$
$$A_1 ::= e_1$$
$$\ldots$$
$$A_n ::= e_n$$

1) Partial ordering $\geq^G$ on non-Terminals

$$\forall i \geq 0, \quad e_i \equiv f(A_{i_1}, \ldots, A_{i_{n_i}}) \quad \text{then:} \quad A_{i_1}, \ldots, A_{i_{n_i}} \geq^G A_i$$

Removal of Mutual Recursion, when possible

$$A_j ::= g(A_{j_1}, \ldots, A_i, \ldots, A_{j_{n_j}}) \quad \text{con} \quad A_{j_1}, \ldots, A_{j_{n_j}} \geq^G A_i \geq^G A_j$$
$$A_i ::= f(A_{i_1}, \ldots, A_j, \ldots, A_{i_{n_i}}) \quad \text{con} \quad A_{i_1}, \ldots, A_{i_{n_i}} \geq^G A_i$$

$$A_j ::= g(A_{j_1}, \ldots, A_i, \ldots, A_{j_{n_j}})$$
$$A_i ::= f(A_{i_1}, \ldots, g(A_{j_1}, \ldots, A_i, \ldots, A_{j_{n_j}}), \ldots, A_{i_{n_i}})$$

# Example

S::= u A B v | B u
A::= v u v | B
B::= u S v | u v u

$B \geq^G A \geq^G S$

S::= u A B v | B u
B::= u S v | u v u

↓

S::= u A B v | B u
B::= u (u A B v | B u) v | u v u

S::= u A B v | B u
A::= v u v | B
B::= uu A B vv | u B uv | u v u

A::= v u v | B
B::= uu A B vv | u B uv | u v u

↓

A::= v u v | B
B::= uu (v u v | B) B vv | u B uv | u v u

S::= u A B v | B u
A::= v u v | B
B::= uu vuv B vv | uu B B vv | u B uv | u v u

13

$\forall i \geq 0, \quad L(A_i) = L(e_i)$

- $L(e)$ is an expression on $2^{\Sigma^*}$ containing only:

    X (finite products)

    $\cup$ (possibly, denumerable unions)

- $L(e)$ is continuos on $2^{\Sigma^*}$

Whenever $L(A_i) = L(e_i)$ is recursive: $L(e_i) \equiv E(L(A_i))$

Recursive equations $X=E(X)$, have to be solved in the variable $X \equiv L(A_i)$ on $2^{\Sigma^*}$ using the **Tarski's (Fixpoint) Iteration** below:

$$X = \bigcup_{i \in \aleph} E(\bot)^i$$
$$E(\bot)^0 = E(X \leftarrow \bot)$$
$$E(\bot)^{k+1} = E(X \leftarrow E(\bot)^k)$$

# What about a system of equations

A system of Recursive equations:

$$\{X_1 = E_1(X_1, \ldots X_n), \ldots, X_n = E_n(X_1, \ldots X_n)\}$$

$X_i \equiv L(A_i)$ on $2^{\Sigma^*}$ using the **Tarski's (Fixpoint) Iteration** below:

$$X_j = \bigcup_{i \in \aleph} E_j(\bot, \ldots, \bot)^i$$

$$E_j(\bot, \ldots, \bot)^0 = E_j(X_1 \leftarrow \bot, \ldots, X_n \leftarrow \bot)$$

$$E_j(\bot, \ldots, \bot)^{k+1} = E_j(X_1 \leftarrow E_1(\bot)^k, \ldots, X_n \leftarrow E_n(\bot)^k)$$

# Example1: Tarski's Fixpoint Iteration

S::=u S | ε $\longrightarrow$ X= $\{u\}\times X\cup\{\lambda\}$

$E(X)$

$E(X\leftarrow \bot)^0 = \{u\}\times\bot\cup\{\lambda\} = \bot\cup\{\lambda\} = \{\lambda\}$

$E(X\leftarrow \bot)^1 = \{u\}\times\{\lambda\}\cup\{\lambda\} = \{u, \lambda\}$

$E(X\leftarrow \bot)^2 = \{u\}\times\{u,\lambda\}\cup\{\lambda\} = \{uu, u, \lambda\}$

$E(X\leftarrow \bot)^3 = \{u\}\times \{uu, u, \lambda\} \cup\{\lambda\} = \{u^3,u^2,u, \lambda\}$

$E(X\leftarrow \bot)^n = \{u^n,u^{n-1},\ldots,u, \lambda\}$

$L(S)=\{u^n \mid n\in\aleph\}$ $= u^*$

# Example2: Tarski's Fixpoint Iteration

$S ::= u \, S \, v \mid z$ $\longrightarrow$ $X = \{u\} \times X \times \{v\} \cup \{z\}$

$E(X)$

$E(X \leftarrow \bot)^0 = \{u\} \times \bot \times \{v\} \cup \{z\} = \bot \cup \{z\} = \{z\}$

$E(X \leftarrow \bot)^1 = \{u\} \times \{z\} \times \{v\} \cup \{z\} = \{uzv, z\}$

$E(X \leftarrow \bot)^2 = \{u\} \times \{uzv, z\} \times \{v\} \cup \{z\} = \{u^2zv^2, uzv, z\}$

$E(X \leftarrow \bot)^n = \{u^n z v^n, u^{n-1} z v^{n-1}, \dots, z\}$

$L(S) = \{u^n z v^n \mid n \geq 0\}$

# Example3: Tarski's Fixpoint Iteration

A::= A+A
A::= A*A
A::= id

$\Longrightarrow$  X= {x+y | x,y∈X}∪{x*y | x,y∈X}∪{id}

E(X)

E(X← ⊥)$^0$ = {x+y | x,y∈⊥}∪{x*y | x,y∈⊥}∪{id}
  = ⊥∪⊥∪{id} = {id}

E(X← ⊥)$^1$ = {x+y | x,y∈ {id}}∪{x*y | x,y∈ {id}}∪{id}
  = {id+id}∪{id*id}∪{id} = {id, id+id, id*id}

E(X← ⊥)$^2$ = {x+y | x,y∈E(X← ⊥)$^1$}∪{x*y | x,y∈E(X← ⊥)$^1$}∪{id}
  = {id, id+id, id*id,id+id+id, id+id*id, id*id+id, id*id*id,…, id*id*id*id}

E(X← ⊥)$^n$ = {id, id t$^k$ | t∈{+id, *id}, k ∈[1..2$^n$-1]}

L(A)={id t$^n$ | t∈{+id, *id}, n∈ℵ }

18

# How to do Syntactic Analysis

- Top-Down and leftmost derivation
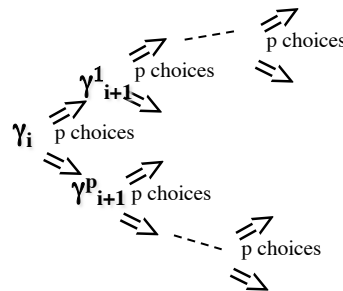- Bottom-Up and reversed roghtmost dervation

# TOP-DOWN and BOTTOM-UP  Parsers

Let **G=<V,Σ,s∈V,P>** be a (context free) grammar. Let **w** be a sequence of words in **Σ.**

- Analysis has to answer to the following question:
  is **w ∈L(G) or not** ?
- or, equivalently:
  is **s =>*w or not** ?
- **Membership:** is this Decision Problem, computable?
    **-- Yes.** It is decidable for all classes of **Monotone Grammars**.

- The solution consists in defining a procedure (**The Parser Core**) able to construct a derivation $s=>\gamma_1=>\dots=>\gamma_k\equiv w$, if one exists.

# Construction of a Derivation

- The solution consists in defining a procedure (**The Parser Core**) able to construct a derivation $s => \gamma_1 => \ldots => \gamma_k \equiv w$, if one exists.

- The construction of a derivation could be done in a non-efficient way, and even worse, at a non-linear, up to exponential, complexity time (/space) cost.

Trying $p$ optional productions at each $\gamma_i$ leads to:



construction of (exponential) $(O(p^n))$ derivations to find the one right or to answer "no-accept".
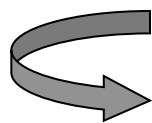
# Top-Down
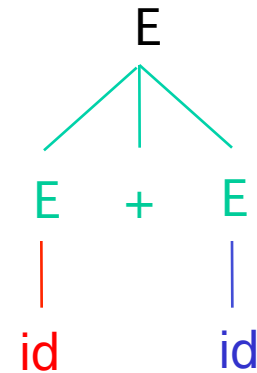## Simple for Handmade Constructions,
## Few Grammars

$E \Rightarrow_1 E+E \Rightarrow_1 id +E \Rightarrow_1 id+id$

Leftmost non-terminal of Left-Sentential-Form

*First* Applicable Production

*Failure*: Backward to the last alternative

```
        E
       /|\
      E + E
      |   |
      id  id
```

**Step 1**
**Step 2**
**Step 3**

Top-Down = Leftmost

*p1*: E::= E+E
*p2*: E::= E*E
*p3*: E::= id

# LSF forms a Complete Base for Context-Free Grammars

$G = <V, \Sigma, s \in V, P>$

Left Sentential Form (of G): $LSF_G$

$$\alpha\beta\gamma \in LSF_G \qquad iff \qquad s \mid=>^+ \alpha\beta\gamma$$

$$\alpha A\beta \mid=> \alpha\gamma\beta \qquad iff \qquad A::=\gamma \in P \ \& \ \alpha \in \Sigma^*$$
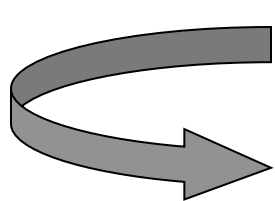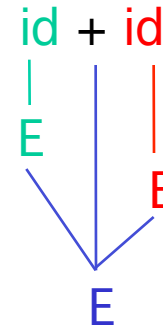
Only $LSF_G$

$$L(G) = \{w \in \Sigma^* \mid s =>^+ w\}$$
$$= \{w \in \Sigma^* \mid s \mid=>^+ w\}$$

# Bottom-Up
## More Complicated Techniques
## Many More Grammars - Many More Languages

E => E+E => E+id => id+id

id + id

Looking for Handle

reduction

*failure*: backward for "true" Handle

**Step 1**
**Step 2**
**Step 3**

Bottom-Up = Rightmost Reversed

*p1*: E::= E+E
*p2*: E::= E*E
*p3*: E::= id

# RSF forms a Complete Base for Context-Free Grammars

$G = <V,\Sigma,s{\in}V,P>$

Right Sentential Form (of G): $RSF_G$

$$\alpha\beta\gamma \in RSF_G \quad \textit{iff} \quad s \; {}_r{=}{>}^+ \; \alpha\beta\gamma$$

$$\mathbf{\alpha A\beta} \; {}_r{=}{>} \; \mathbf{\alpha\gamma\beta} \quad \textit{iff} \quad A{::=}\gamma \in P \; \& \; \beta{\in}\Sigma^*$$

Only RSF
$$L(G) = \{w \in \Sigma^* \mid s \; {=}{>}^+ \; w\}$$
$$= \{w \in \Sigma^* \mid s \; {}_r{=}{>}^+ \; w\}$$

$\mathbf{B{::=}\beta} \in P$ is **Handle** of $\mathbf{\alpha\beta\gamma} \in RSF_G$

*if and only if* $\quad \alpha B\gamma \in RSF_G$