# Compiler Construction
# Lent Term 2021
# Lecture 5 : Theoretical foundations of Bottom-up (LR) parsing

1. **This lecture develops a general theory for <span style="color:red">non-deterministic</span> bottom-up parsing**
2. **Next lecture will present two techniques for imposing determinism --- SLR(1) parsing and LR(1) parsing.**

**Timothy G. Griffin**
**tgg22@cam.ac.uk**
**Computer Laboratory**
**University of Cambridge**

1

$$G_2 = (N_2, T_1, P_2, E')$$

$$N_2 = \{E', E, T, F\} \qquad T_1 = \{+, *, (,), \text{id}\}$$

$$P_2 : E' \rightarrow E$$

$$E \rightarrow E + T \mid T \qquad \text{(expressio ns)}$$

$$T \rightarrow T * F \mid F \qquad \text{(terms)}$$

$$F \rightarrow (E) \mid \text{id} \qquad \text{(factors)}$$

**Note: E' was added for convenience to ensure that there is a single starting production.**

2

# Rightmost derivations

$$w \in T^* \qquad \alpha, \beta \in (N \cup T)^*$$

Given : $\alpha A w$ and a production $A \to \beta$

a rightmost derivation step is written as

$$\alpha A w \Rightarrow_{rm} \alpha \beta w$$

3

# A rightmost derivation of *(x+y)*

$$E' \Rightarrow_{rm} E$$
$$\Rightarrow_{rm} T$$
$$\Rightarrow_{rm} F$$
$$\Rightarrow_{rm} (E)$$
$$\Rightarrow_{rm} (E + T)$$
$$\Rightarrow_{rm} (E + F)$$
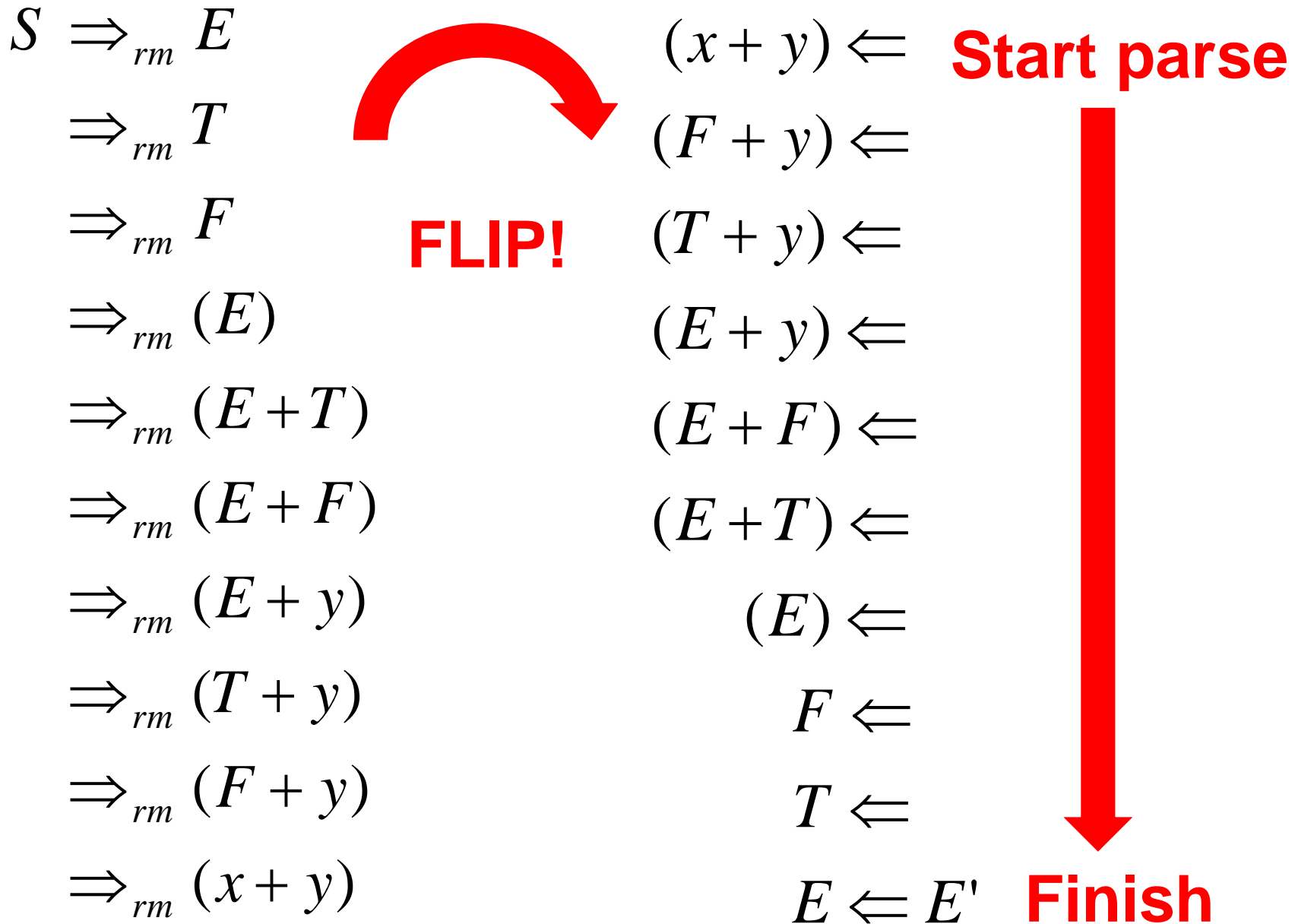$$\Rightarrow_{rm} (E + y)$$
$$\Rightarrow_{rm} (T + y)$$
$$\Rightarrow_{rm} (F + y)$$
$$\Rightarrow_{rm} (x + y)$$

Top-down (LL) parsing is based on
<u>left-most</u> derivations.

Bottom-up (LR) parsing is based on
<u>right-most</u> derivations.

4

# But Bottom-up parsers perform the derivation in reverse!

$$S \Rightarrow_{rm} E$$
$$\Rightarrow_{rm} T$$
$$\Rightarrow_{rm} F$$
$$\Rightarrow_{rm} (E)$$
$$\Rightarrow_{rm} (E+T)$$
$$\Rightarrow_{rm} (E+F)$$
$$\Rightarrow_{rm} (E+y)$$
$$\Rightarrow_{rm} (T+y)$$
$$\Rightarrow_{rm} (F+y)$$
$$\Rightarrow_{rm} (x+y)$$

**FLIP!**

$$(x+y) \Leftarrow$$
$$(F+y) \Leftarrow$$
$$(T+y) \Leftarrow$$
$$(E+y) \Leftarrow$$
$$(E+F) \Leftarrow$$
$$(E+T) \Leftarrow$$
$$(E) \Leftarrow$$
$$F \Leftarrow$$
$$T \Leftarrow$$
$$E \Leftarrow E'$$

**Start parse**

**Finish**

5

**Can we transform a backwards derivation into an execution of a stack machine?**

| stack | input |
|---|---|
| $ | $(x+y)$$ |
| $(F | $+y)$$ |
| $(T | $+y)$$ |
| $(E | $+y)$$ |
| $(E+F | $)$$ |
| $(E+T | $)$$ |
| $(E) | $$ |
| $F | $$ |
| $T | $$ |
| $E | $$ |
| $E' | $$ |

$(x+y) \Leftarrow$
$(F+y) \Leftarrow$
$(T+y) \Leftarrow$
$(E+y) \Leftarrow$
$(E+F) \Leftarrow$
$(E+T) \Leftarrow$
$(E) \Leftarrow$
$F \Leftarrow$
$T \Leftarrow$
$E \Leftarrow E'$

View the reversed derivation as a stack machine (use $ as stack bottom and end-of-input).

Can we make this work?

# Let's try to formalize such a parser

An <u>LR parser configuration</u> has the form

$$\$\alpha, x\$$$

$$(\alpha \text{ is the stack, } x \text{ the remaining input})$$

The configuration is <u>valid</u> when there exists a right-most derivation of the form

$$S \underset{rm}{\overset{*}{\Longrightarrow}} \alpha x$$

7

Suppose

$$\alpha A x \Rightarrow_{rm} \alpha \beta B z x$$

Our "backwards" parser MIGHT move

from one configuration to another like so :

$$\$\alpha \beta B z, x\$ \xrightarrow{\quad reduce \quad} \$\alpha A, x\$$$

This action is called a reduction

using production $A \rightarrow \beta B z$

8

Suppose we have the derivation

$$\alpha A x \Rightarrow_{rm} \alpha \beta B z x \Rightarrow_{rm} \alpha \beta \gamma z x$$

using $A \rightarrow \beta B z$ and then $B \rightarrow \gamma$.

Simulating this in reverse, our parser gets stuck :

$$\$\alpha\beta\gamma, zx\$$$

$$\xrightarrow{reduce} \$\alpha\beta B, zx\$$$

$$\xrightarrow{???} ???$$

We want $\beta B z$ on top of the stack!  9

# We need an action that <u>shifts</u> a terminal onto the stack!

$$\alpha A x \Rightarrow_{rm} \alpha \beta B z x \Rightarrow_{rm} \alpha \beta \gamma z x$$

$\$\alpha\beta\gamma, zx\$$

$\xrightarrow{\quad reduce \quad} \$\alpha\beta B, zx\$$

$\xrightarrow{\quad shift(s) \quad} \$\alpha\beta B z, x\$$

$\xrightarrow{\quad reduce \quad} \$\alpha A, x\$$

How do we know when to stop shifting? Here we don't want to gobble up $x$!

Let's make sure that this can work when $B$ does
not appear in the right-hand side of $A$'s production,

$$\alpha BxAz \Rightarrow_{rm} \alpha Bxyz \Rightarrow_{rm} \alpha\gamma xyz$$

using production $A \rightarrow y$, then $B \rightarrow \gamma$.

Our parser's possible actions :

$\$\alpha\gamma, xyz\$$

$\xrightarrow{reduce} \$\alpha B, xyz\$$

$\xrightarrow{shift(s)} \$\alpha Bxy, z\$$

$\xrightarrow{reduce} \$\alpha BxA, z\$$

All good! But again, how do we know when to reduce and when to stop shifting?

11

The previous two slides demonstrat e that if
we have a derivation

$$S \Rightarrow_{rm}^{*} w$$

Then we can always "replay it" in reverse using
shift/redu ce actions

$$\$, w\$ \rightarrow^{*} \$S, \$$$

This tells us that shift and reduce are sufficient .
However, when we are parsing a $w$ we won' t
have access to a derivation to replay! So our
parser wil l be non - determinis tic and GUESS what
the future holds!

12

| stack | input | action[X, a] |
|-------|-------|--------------|
| $ | $(x + y)\$$ | shift |
| $($ | $x + y)\$$ | shift |
| $(x$ | $+ y)\$$ | reduce $F \rightarrow id$ |
| $(F$ | $+ y)\$$ | reduce $T \rightarrow F$ |
| $(T$ | $+ y)\$$ | reduce $E \rightarrow T$ |
| $(E$ | $+ y)\$$ | shift |
| $(E +$ | $y)\$$ | shift |

13

| stack | input | action[X, a] |
| --- | --- | --- |
| $\$(E + y$ | $)\$$ | reduce $F \rightarrow id$ |
| $\$(E + F$ | $)\$$ | reduce $T \rightarrow F$ |
| $\$(E + T$ | $)\$$ | reduce $E \rightarrow E + T$ |
| $\$(E$ | $)\$$ | shift |
| $\$(E)$ | $\$$ | reduce $F \rightarrow (E)$ |
| $\$F$ | $\$$ | reduce $T \rightarrow F$ |
| $\$T$ | $\$$ | reduce $F \rightarrow E$ |
| $\$E$ | $\$$ | reduce $S \rightarrow E$ |
| $\$E'$ | $\$$ | accept! |

14

# How do we decide when to shift and when to reduce?

Suppose $A \rightarrow \beta\gamma$ is a production . When

our parser is in the configurat ion

$$\$\alpha\beta\gamma, x\$$$

we MIGHT want to reduce with $A \rightarrow \beta\gamma$.

However, if we have

$$\$\alpha\beta, x\$$$

we MIGHT want to continue parsing with the hope

of eventually getting $\beta\gamma$ on top of the stack so

that we can then reduce to $A$.

For every grammar production

$$A \to \beta\gamma \qquad (\beta, \gamma \in (N \cup T)^*)$$

produce the LR(0) item

$$A \to \beta \bullet \gamma$$

Interpretation of $A \to \beta \bullet \gamma$ : we have already parsed some input $x$ derivable from $\beta$ ($\beta \Rightarrow_{rm}^* x$) and we MIGHT next see some input derivable from $\gamma$.

16

$E' \to \bullet E$         $E' \to E \bullet$

$E \to \bullet E + T$      $T \to \bullet T * T$      $F \to \bullet (E)$

$E \to E \bullet + T$      $T \to T \bullet * F$      $F \to (\bullet E)$

$E \to E + \bullet T$      $T \to T * \bullet F$      $F \to (E \bullet)$

$E \to E + T \bullet$      $T \to T * F \bullet$      $F \to (E) \bullet$

$E \to \bullet T$          $T \to \bullet F$          $F \to \bullet id$

$E \to T \bullet$          $T \to F \bullet$          $F \to id \bullet$

Definition . Item $A \rightarrow \beta \bullet \gamma$ is valid for $\phi\beta$
if there exists a derivation
$$S \Rightarrow^{*}_{rm} \phi A x \Rightarrow_{rm} \phi \beta \gamma x$$

If item $A \rightarrow \beta \bullet \gamma$ is valid for $\phi\beta$

then our parser could use the item

as a guide when in configurat ion

$\$\phi\beta, z\$$.

Suppose $A \to \beta B \gamma$ and $B \to \alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_k$. Consider the ways in which items for these productions might be used as parsing guides.

| Derivation | Parse | Possible guides |
|---|---|---|
| $S$ | $\$S, \$^*$ | |
| $\Rightarrow_{rm}^* \phi A x$ | $^* \leftarrow \$\phi A, x\$$ | |
| $\Rightarrow_{rm} \phi \beta B \gamma x$ | $\leftarrow \$\phi \beta B \gamma, x\$$ | $A \to \beta B \gamma \bullet$ |
| $\Rightarrow_{rm}^* \phi \beta B z x$ | $^* \leftarrow \$\phi \beta B, z x\$$ | $A \to \beta B \bullet \gamma$ |
| $\Rightarrow_{rm} \phi \beta \alpha_i z x$ | $\leftarrow \$\phi \beta \alpha_i, z x\$$ | $B \to \alpha_i \bullet$ |
| $\Rightarrow_{rm}^* \phi \beta u z x$ | $^* \leftarrow \$\phi \beta, u z x\$$ | $A \to \beta \bullet B \gamma, B \to \bullet \alpha_i$ |

# Using items as parsing guides

Suppose our parser is in the config

$$\$\phi\beta, cz\$$$

and $A \rightarrow \beta \bullet c\gamma$ is valid for $\phi\beta$.

Then we MIGHT shift c onto the stack :

$$\$\phi\beta, cz\$ \xrightarrow{shift} \$\phi\beta c, z\$$$

---

Suppose our parser is in the config

$$\$\phi\beta, z\$$$

and $A \rightarrow \beta \bullet$ is valid for $\phi\beta$.

Then we MIGHT perform a reduction :

$$\$\phi\beta, z\$ \xrightarrow{reduce} \$\phi A, z\$$$

20

Suppose our parser is in the config

$$\$\phi\beta, z\$$$

which we will assume is valid, so $S \Rightarrow^*_{rm} \phi\beta z$.

Suppose $A \rightarrow \beta \bullet \gamma$ is valid for $\phi\beta$.

Then $\gamma$ MIGHT capture the future of our parse (the past of that derivation ). That is, it MIGHT be that

$$S \Rightarrow^*_{rm} \phi A x \Rightarrow_{rm} \phi\beta\gamma x \Rightarrow^*_{rm} \phi\beta y x = \phi\beta z$$

If so, our parser MIGHT proceed like so :

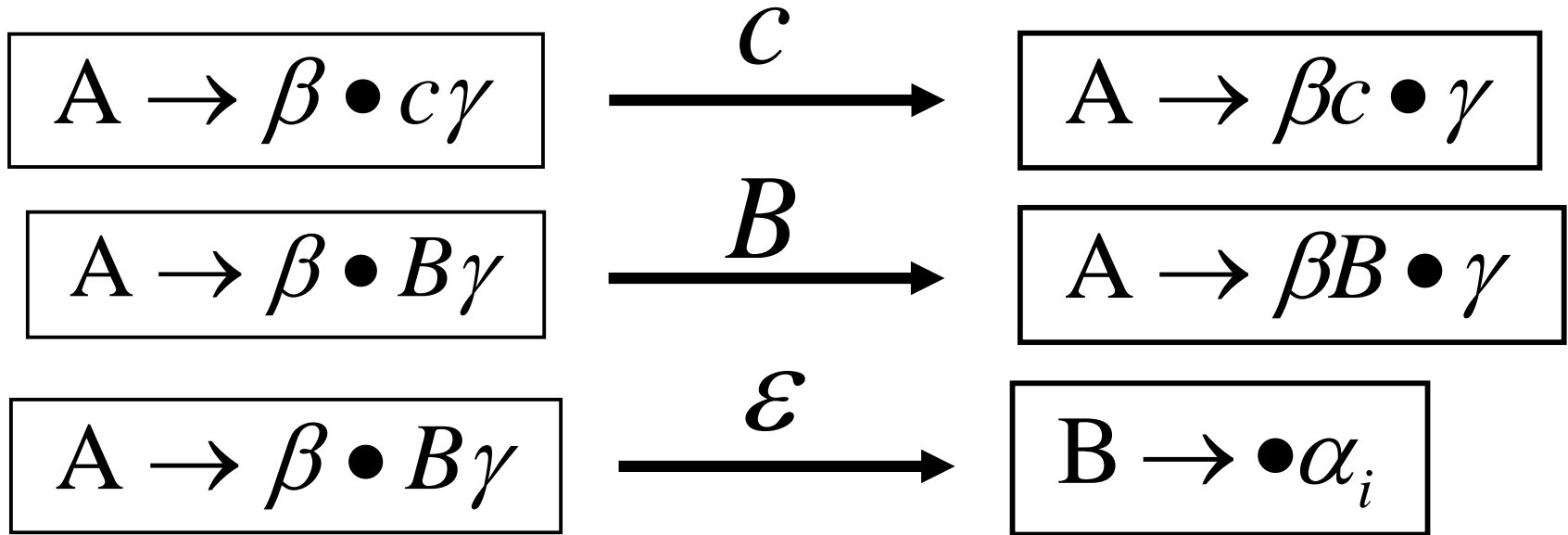$$\$\phi\beta, z\$ = \$\phi\beta, yx\$ \rightarrow^* \$\phi\beta\gamma, x\$ \xrightarrow{reduce} \$\phi A, x\$.$$

That is, our parser could guess that $\gamma$ will derive a prefix of the remaining input z.

21

Augment our shift/reduce parser in such
a way that in every configuration it can
derive the set of all items valid for the
contents of the current stack.

Then at each step the parser can
(non - deterministically) select an item
from this set to use as a guide.

22

$$A \to \beta \bullet c\gamma \quad \xrightarrow{c} \quad A \to \beta c \bullet \gamma$$

$$A \to \beta \bullet B\gamma \quad \xrightarrow{B} \quad A \to \beta B \bullet \gamma$$

$$A \to \beta \bullet B\gamma \quad \xrightarrow{\varepsilon} \quad B \to \bullet \alpha_i$$

The initial  state $q_0$ is this  item constructe d

from the unique  starting  production

$$E' \to \bullet E \qquad \text{(for example)}$$

and every item (state) is a final  state.

Let $\delta_G$ be the transitio n function  of this  NFA.    23
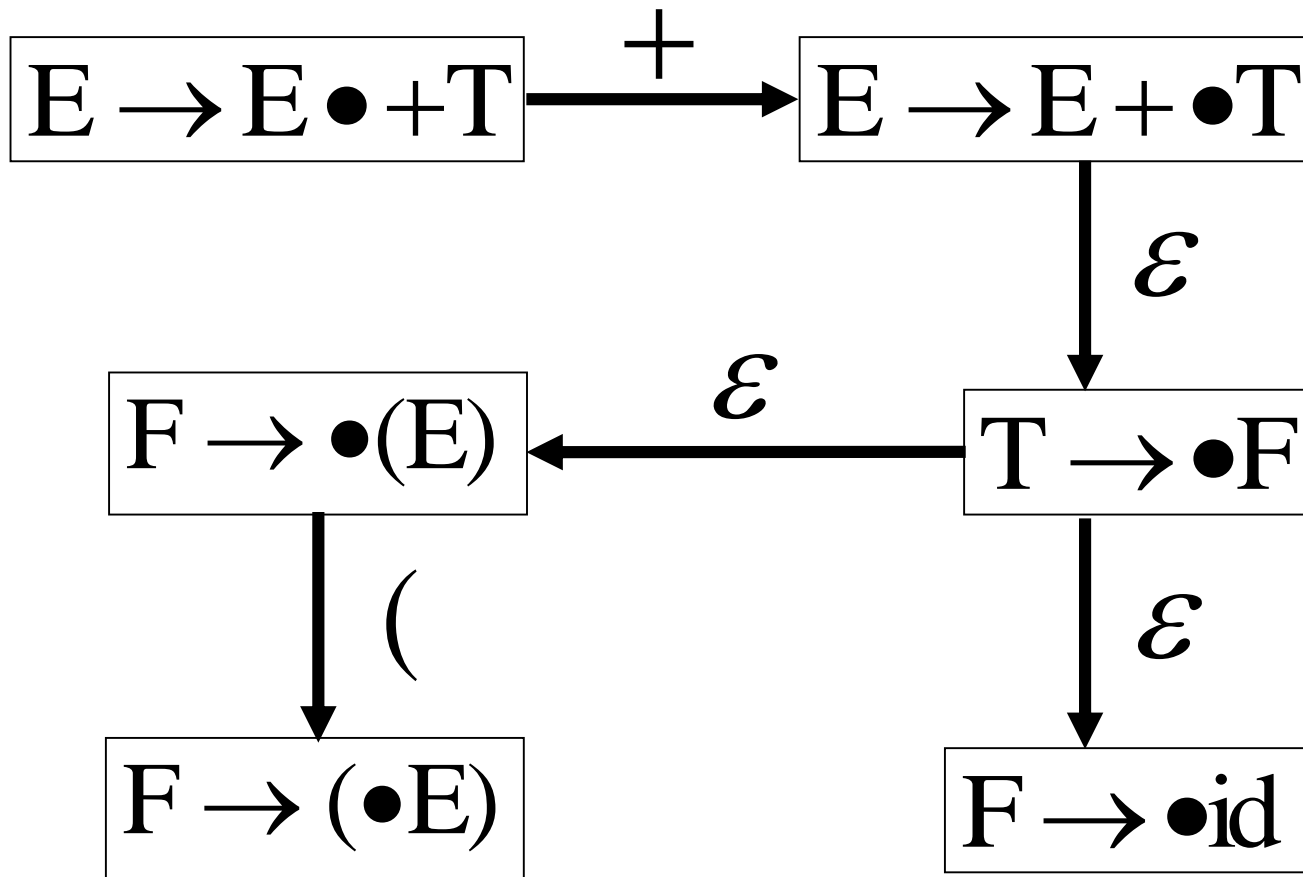
# Main LR parsing theorem

Theorem. $A \rightarrow \beta \bullet \gamma \in \delta_G(q_0, \phi\beta)$ if and only if $A \rightarrow \beta \bullet \gamma$ is valid for $\phi\beta$.

Amazing fact : the

language of the stack

is regular!

See proof (not examinable) in Introduction to Automata Theory, Languages, and Computation. Hopcroft and Ullman.

24

$c :=$ first symbol of input w$

while(true)

$\quad \alpha :=$ the stack

$\quad$ if $A \rightarrow \beta \bullet c\gamma \in \delta_G(q_0, \alpha)$

$\quad$ then shift $c$ onto the stack

$\quad\quad c :=$ next input token;

$\quad$ if $A \rightarrow \beta \bullet \in \delta_G(q_0, \alpha)$

$\quad$ then reduce : pop $\beta$ off the stack

$\quad\quad$ and then push $A$ onto the stack;

$\quad$ if $S \rightarrow \beta \bullet \in \delta_G(q_0, \alpha)$

$\quad$ then accept and exit if no more input;

$\quad$ if none of the above then ERROR

This is non-deterministic since multiple conditions can be true and multiple items can match any condition.

26

# How can we make the algorithm deterministic?

1. **The easy part: convert the NFA to a DFA**
2. **When there are shift/reduce or reduce/reduce conflicts, find some way of making a deterministic choice.**
3. **For (2), peek into the input buffer.**
4. **For (3), use FIRST and/or FOLLOW!**

Note : no matter how we do this there will be non-ambiguous grammars for which our deterministic parser will fail.

Next lecture : we will look at two popular approaches, SLR(1) and LR(1).