# Check-in

Review: Predictive Parsing

Assume an LL(1) parser with…          this selector table:          this syntax stack:          and this **(** lookahead token:

|   | ( | ) | { | } |
|---|---|---|---|---|
| S | ( S ) | ) | { | } |

| S |
|---|
| ) |
| ) |
| **eof** |

**(**

Draw the configuration of the parser after it processes the tokens  **( )**

# Housekeeping
## Administrivia

**Projects**

- P2 (nominally) due Wednesday

- P3 out Friday

**Trials**

- Trial 1 due tonight

# Housekeeping
## Administrivia

**Labs**

- Based on the confusion about abstract classes, I've decided to shift the labs a bit

# EECS 665

# COMPILER CONSTRUCTION

# FIRST Sets

# Last Time
Review – Predictive Parsing

## Intro to Parsing

- Complexity

## A New Type of Language – LL(k)

- Intro

- LL(1) parsing

> ### You Should Know
>
> - What parsing is
> - What LL(1) languages are
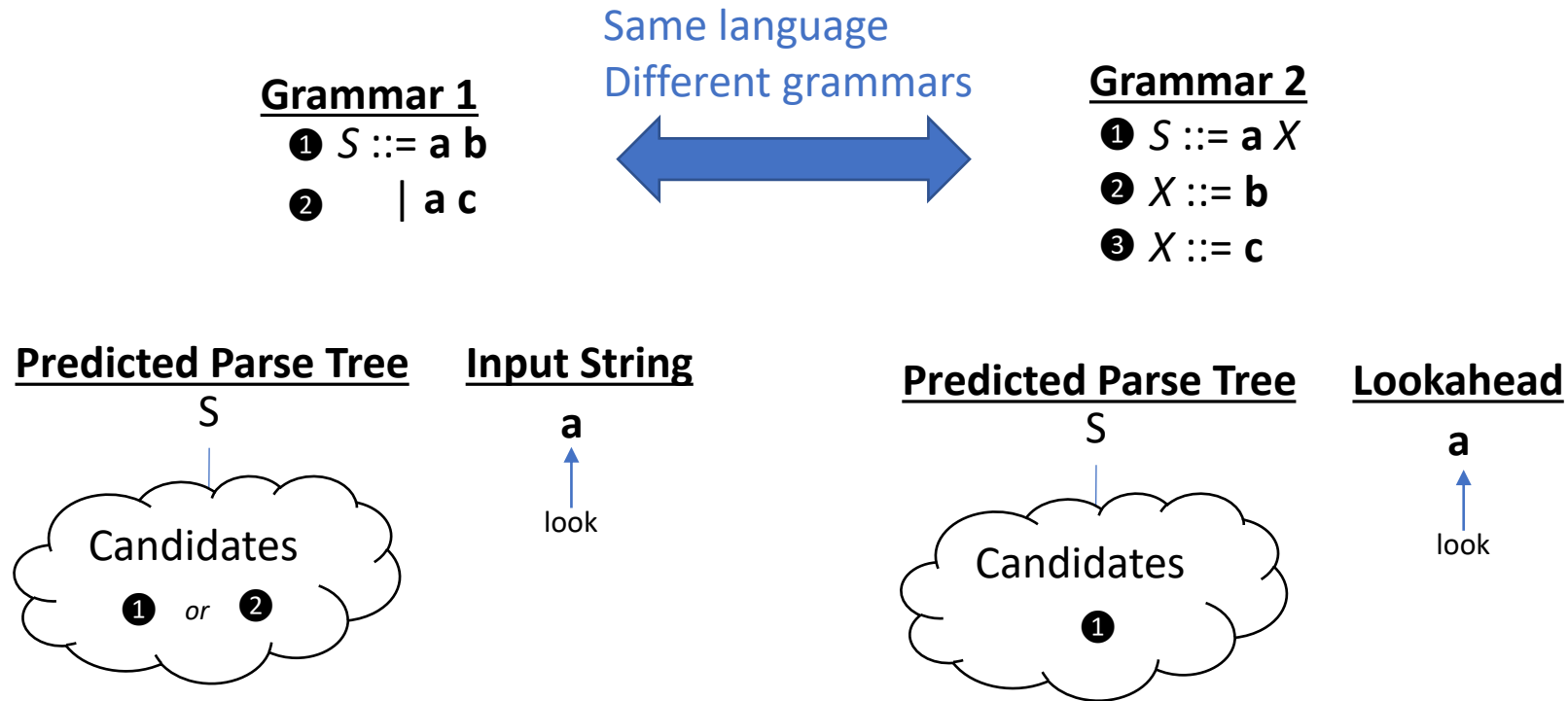> - How an LL(1) parser operates

**Parsing**

# Where we Left Off
Review – Predictive Parsing

## The language might be LL(1) … even when the grammar is not!

Same language
Different grammars

**Grammar 1**
❶ $S$ ::= **a b**
❷    | **a c**

**Grammar 2**
❶ $S$ ::= **a** $X$
❷ $X$ ::= **b**
❸ $X$ ::= **c**

**Predicted Parse Tree**
S

Candidates
❶ *or* ❷

**Input String**
**a**

look

**Predicted Parse Tree**
S

Candidates
❶

**Lookahead**
**a**

look

# Today's Outline
## Preview – FIRST Sets

**Transforming Grammars**

• Fixing LL(1) "near misses"

**Building LL(1) Parsers**

• What the selector table needs

• FIRST Sets

**Parsing**

Transforming Grammars – Fixing LL(1) Near Misses

**Given a language, we can't always find an LL(1) grammar *even if one exists***

- Best we can do: simple transformations that remove "obvious" disqualifiers

# Checking if a Grammar is LL(1)

Transforming Grammars – Fixing LL(1) Near Misses

**If either of the following hold, the grammar is <u>not</u> LL(1):**

- The grammar **is** left-recursive

- The grammar **isn't** left-factored

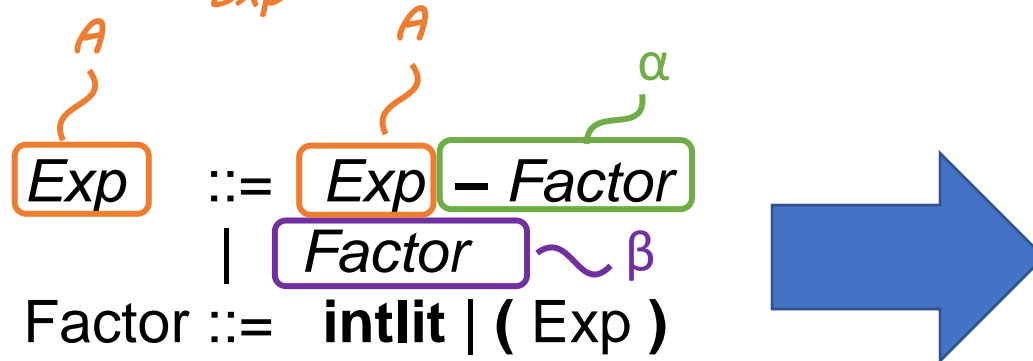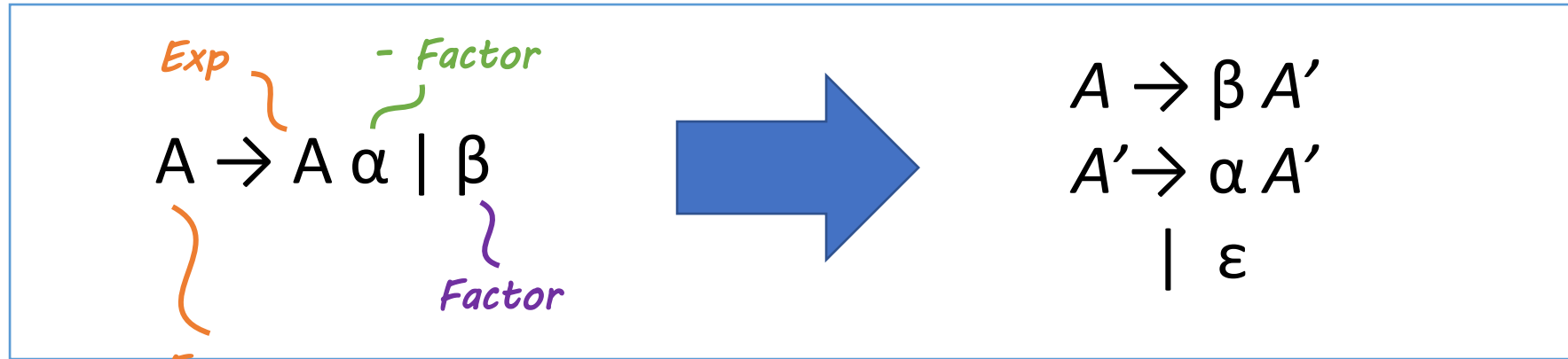**We can transform *some* grammars while preserving the recognized language**

# (Immediate) Left Recursion
## Transforming Grammars – Fixing LL(1) Near Misses

- Recall, a grammar such that $X \overset{+}{\Rightarrow} X\, \alpha$ is left recursive

- A grammar is immediately left recursive if this can happen in one step:
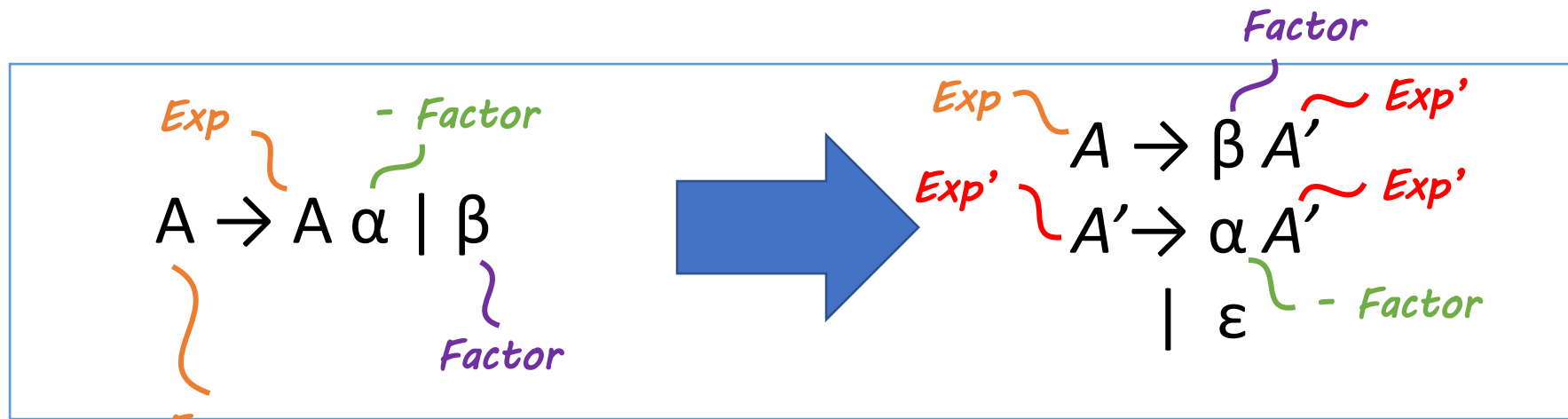
$$A \rightarrow A\ \alpha \mid \beta$$

# Immediate Left Recursion Removal
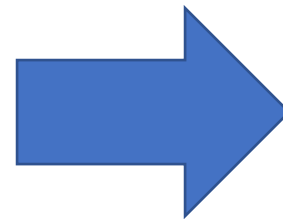## *(Predictive) Parsing - LL(1) Transformations*

(for a single immediately left-recursive rule)

$$A \rightarrow A\,\alpha \mid \beta$$

**Arbitrary Strings (nonterminal or terminal)**

$$A \rightarrow \beta\,A'$$
$$A' \rightarrow \alpha\,A'$$
$$\mid \varepsilon$$

# Immediate Left Recursion Removal
## (Predictive) Parsing - LL(1) Transformations

Exp

− Factor

$A \rightarrow A\ \alpha \mid \beta$

Exp

Factor

$A \rightarrow \beta\ A'$
$A' \rightarrow \alpha\ A'$
$\mid \varepsilon$

A

Exp

A

Exp

α

$\boxed{Exp}$  ::=  $\boxed{Exp}$ $\boxed{- \ Factor}$

$\mid$ $\boxed{Factor}$ $\beta$

Factor ::=  **intlit** $\mid$ **(** Exp **)**

# Immediate Left Recursion Removal
## (Predictive) Parsing - LL(1) Transformations

$$A \rightarrow A \, \alpha \mid \beta$$

$$A \rightarrow \beta \, A'$$
$$A' \rightarrow \alpha \, A'$$
$$\mid \varepsilon$$

$Exp$ ::= $Exp$ − Factor
| Factor
Factor ::= **intlit** | **(** Exp **)**

$Exp$ ::= Factor Exp'
$Exp'$ ::= - Factor Exp'
| ε
Factor ::= **intlit** | **(** Exp **)**

13

# Immediate Left Recursion Removal

## (Predictive) Parsing - LL(1) Transformations

(general rule)

Given Productions

$$A ::= \alpha_1$$
$$| \quad \alpha_2$$
$$| \quad \alpha_n$$
$$| \quad A\,\beta_1$$
$$| \quad A\,\beta_2$$
$$| \quad A\,\beta_m$$

Convert to

$$A ::= \alpha_1\,A'$$
$$| \quad \alpha_2\,A'$$
$$| \quad \alpha_n\,A'$$
$$A' ::= \beta_1\,A'$$
$$| \quad \beta_2\,A'$$
$$| \quad \beta_m\,A'$$
$$| \quad \varepsilon$$

# Left Factoring Grammar
## (Predictive) Parsing - LL(1) Transformations

- If a nonterminal has (at least) two productions whose RHS has a common prefix, the grammar is **not** left factored

  (and **not** an LL(1) grammar)

*Question: What makes this grammar not left-factored?*
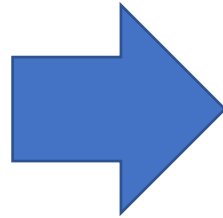
$$Exp ::= ( Exp )$$
$$| \; \{ Exp \}$$
$$| \; ()$$
$$| \; a \; b$$
$$| \; b \; b$$

# Left Factoring: Simple Rule
## *(Predictive) Parsing - LL(1) Transformations*

Given Productions

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$$

*Pull suffix into a new nonterminal*

Convert to

$$A \rightarrow \alpha\, A'$$
$$A' \rightarrow \beta_1 \mid \beta_2$$

*Add a new rule for suffixes*

$\alpha$  $\beta_1$

X ::= **a b** **c d**
X ::= **a b** **e f**

$\beta_2$

X ::= **a b** $X'$
$X'$ ::= **c d** | **e f**

*(Predictive) Parsing - LL(1) Transformations*

Remove immediate left-recursion

Left-factored

$Exp$ ::= **(** $Exp$ **)**
    | $Exp\ Exp$
    | **( )**

$Exp$ ::= **(** $Exp$ **)** $Exp'$
    | **( )** $Exp'$
$Exp'$ ::= $Exp\ Exp'$
    | ε

$A \rightarrow A\ \alpha\ |\ \beta$

becomes

$A \rightarrow \beta\ A'$
$A' \rightarrow \alpha\ A'$
    | ε

# Attempt LL(1) Conversion
## (Predictive) Parsing - LL(1) Transformations

Remove immediate left-recursion                    Left-factored

$Exp ::= ( Exp )$
$| Exp\ Exp$
$| ( )$

$Exp ::= ( Exp )\ Exp'$
$| ( )\ Exp'$
$Exp' ::= Exp\ Exp'$
$| \varepsilon$

$A \rightarrow \alpha\ \beta_1\ |\ \alpha\ \beta_2$ becomes

$A \rightarrow \alpha\ A'$
$A' \rightarrow \beta_1\ |\ \beta_2$

# Attempt LL(1) Conversion
## (Predictive) Parsing - LL(1) Transformations

Remove immediate left-recursion

Left-factored

$Exp ::= ( Exp )$
$\quad | \; Exp \; Exp$
$\quad | \; ( \, )$

$Exp \;\; ::= ( Exp ) \; Exp'$
$\qquad | \; ( \, ) \; Exp'$
$Exp' ::= Exp \; Exp'$
$\qquad\quad | \; \varepsilon$

$Exp \;\;\; ::= ( \; Exp''$
$Exp'' ::= Exp \, ) \; Exp'$
$\qquad\quad\;\; | \, ) \; Exp'$
$Exp' \;\; ::= Exp \; Exp'$
$\qquad\qquad | \;\; \varepsilon$

$$A \rightarrow \alpha \, \beta_1 \mid \alpha \, \beta_2 \quad \text{becomes}$$

$$A \rightarrow \alpha \, A'$$
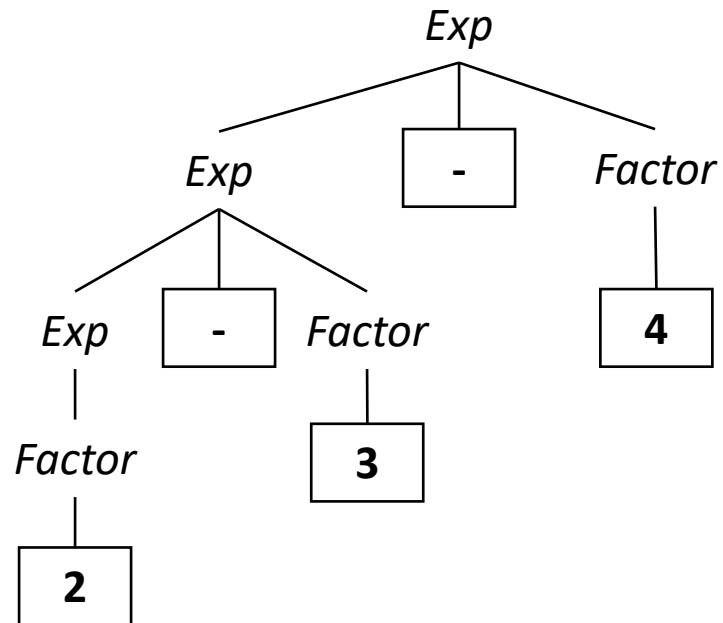$$A' \rightarrow \beta_1 \mid \beta_2$$
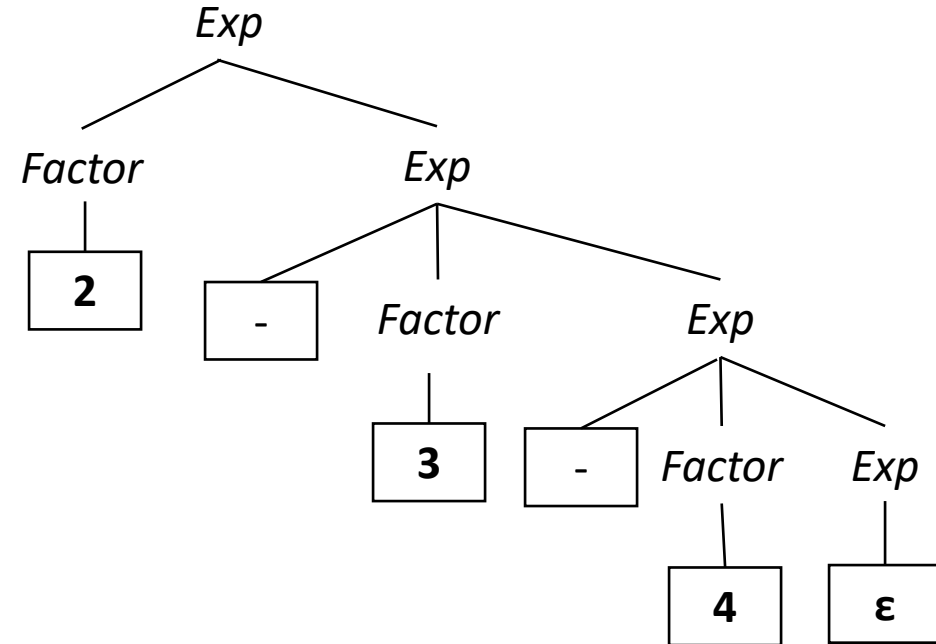
# Current Status
## *(Predictive) Parsing - LL(1) Transformations*

- We've removed 2 disqualifiers from LL(1)
  - Left-recursive grammar
  - **Not** Left-Factored grammar
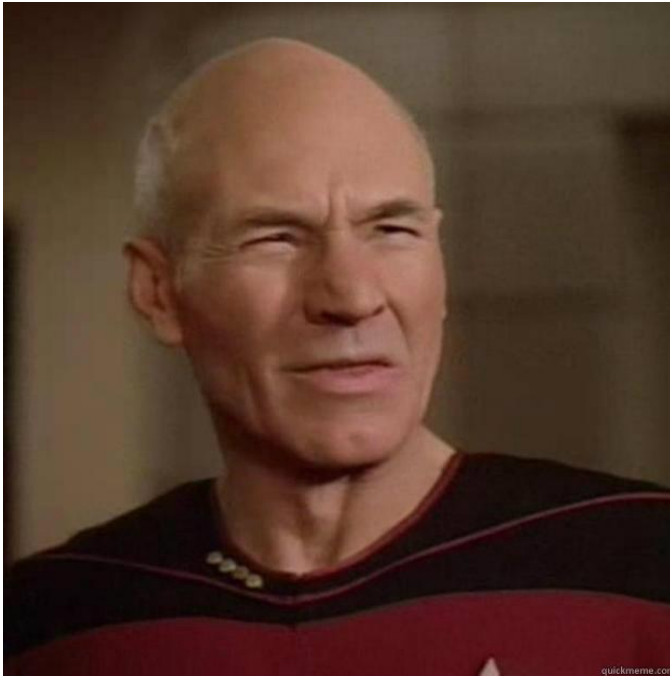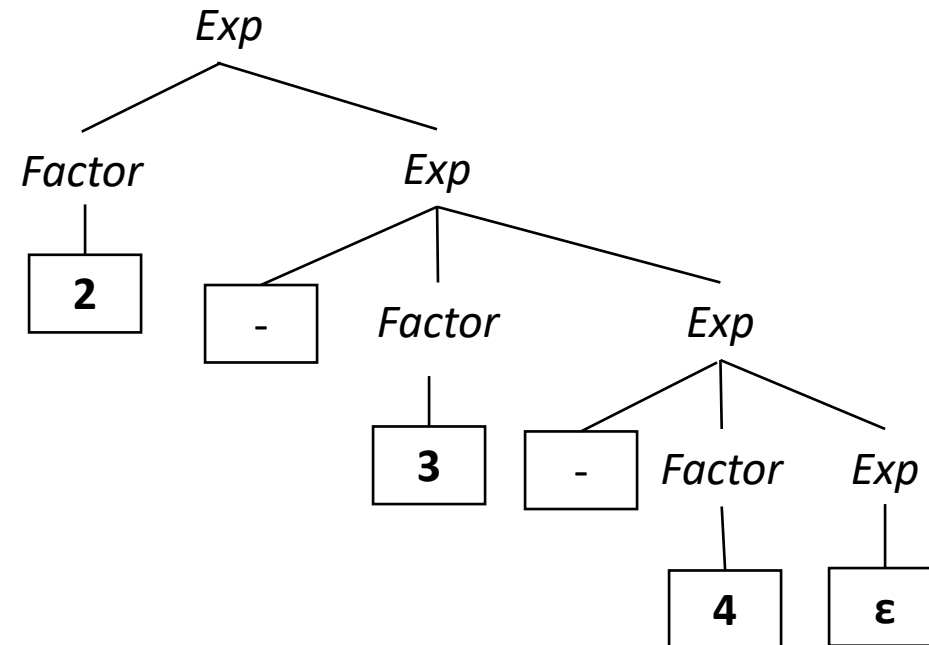
# Let's Check on the Parse Tree
## LL(1) Grammar Transformations

Exp $\quad\rightarrow\quad$ Exp – Factor
$\quad\quad$ | $\quad$ Factor
Factor $\rightarrow\quad$ **intlit** | **(** Exp **)**

Exp $\quad\rightarrow\quad$ Factor Exp'
Exp' $\quad\rightarrow\quad$ **-** Factor Exp'
$\quad\quad\quad$ | $\quad$ ε
Factor $\rightarrow\quad$ **intlit** | **(** Exp **)**

# Let's Check on the Parse Tree
### *LL(1) Grammar Transformations*



$$Exp \rightarrow Factor \; Exp'$$
$$Exp' \rightarrow \textbf{-} \; Factor \; Exp'$$
$$| \quad \varepsilon$$
$$Factor \rightarrow \textbf{intlit} \mid \textbf{(} \; Exp \; \textbf{)}$$

# Nevermind, We'll Fix Parse Trees Later
## *LL(1) Grammar Transformations*

¯\\_(ツ)_/¯

# Today's Outline
### Lecture 9 – FIRST sets

## Transforming Grammars

- Fixing LL(1) "near misses"

## Building LL(1) Parsers

- Understanding LL(1) Selector Tables
- FIRST Sets

**Parsing**

# Recall the LL(1) Parser's Operation

### Building LL(1)Selector Table

**LL(1)**

- Processes **L**eft-to-right
- **L**eftmost derivation
- **1** token of lookahead

**Predictive Parser: "guess & check"**

- Starts at the root, *guesses* how to unfold a nonterminal (derivation step)
- *Checks* that terminals match prediction

# Recall the LL(1) Parser's Operation
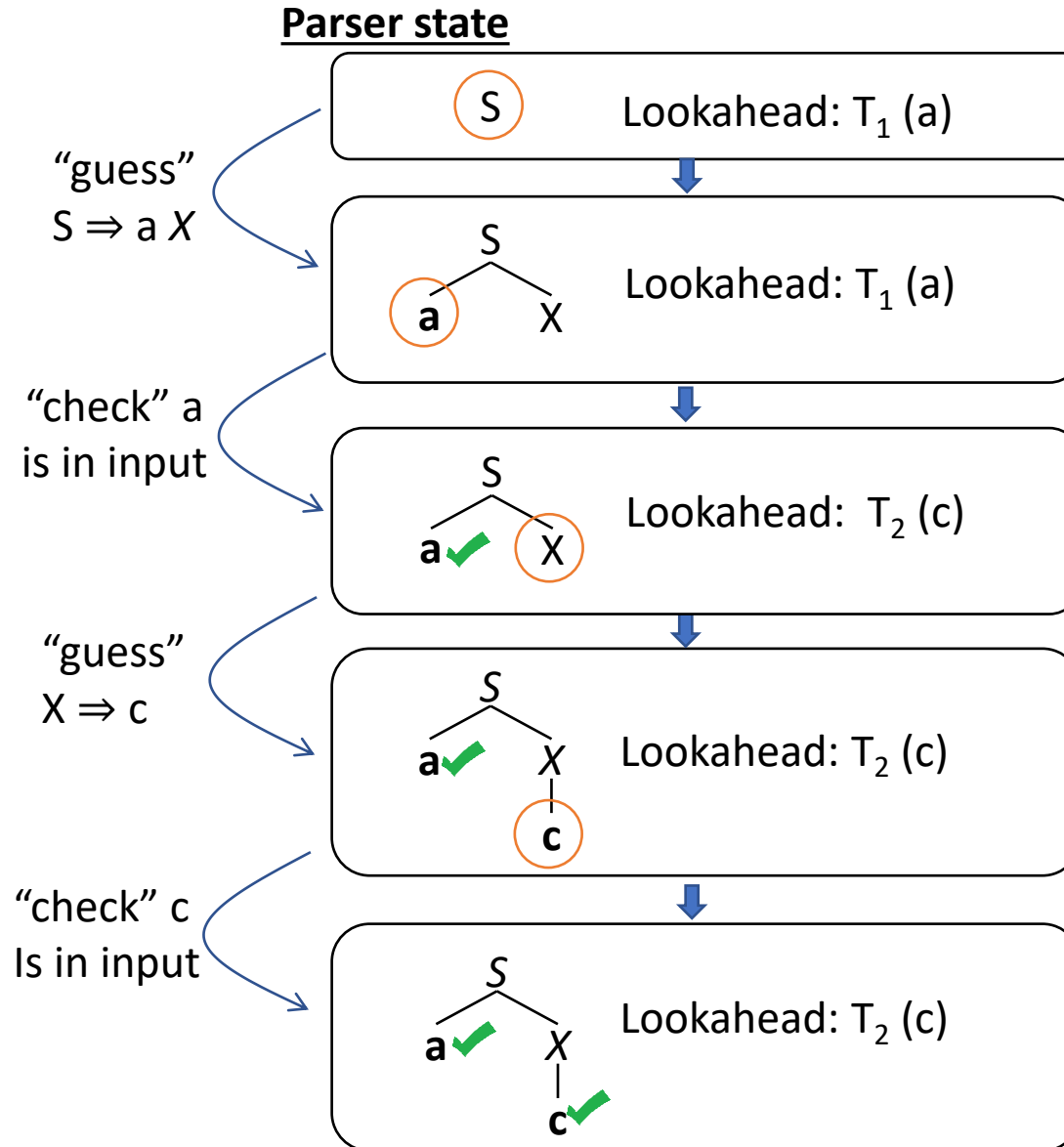## *Building LL(1)Selector Table*

**Parser state**

**Example LL(1) Grammar:**

$S ::= \mathbf{a}\ X$

$X ::= \mathbf{b}\ \mathbf{a} \mid \mathbf{c}$

**Example Input:**

**a   c**

↑   ↑

$T_1$   $T_2$

In practice, table-driven parser uses a stack to match this tree

S          Lookahead: $T_1$ (a)

"guess"
$S \Rightarrow a\ X$

S
a        X          Lookahead: $T_1$ (a)

"check" a
is in input

S
a ✔   X          Lookahead: $T_2$ (c)

"guess"
$X \Rightarrow c$

S
a ✔   X          Lookahead: $T_2$ (c)
c

"check" c
Is in input

S
a ✔   X          Lookahead: $T_2$ (c)
c ✔

*Building Parser Tables*

**The intuition is a bit tricky**

• We need to get into the mindset of the parser



*Pretend your consciousness has been transported inside an LL(1) parser*
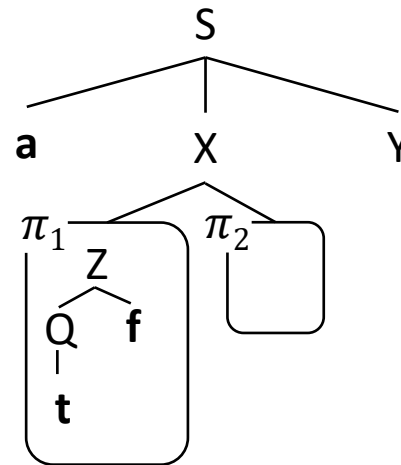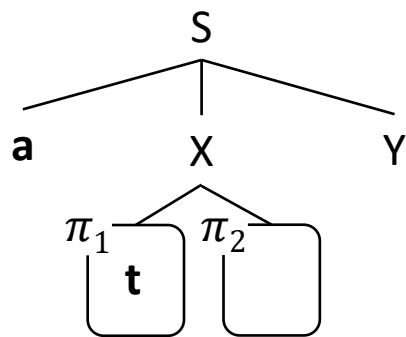
# Become the Parser
## *Building Parser Tables*

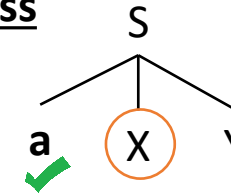You need to unfold a nonterminal *X* with lookahead token **t**

Assume there's an *X* production $X ::= \pi_1\ \pi_2$ (where $\pi_1$ and $\pi_2$ are some kind of symbol)

How do we know to guess this production?

Case 1: $\pi_1$ subtree may start with **t**

**Parse in Progress**

Lookahead:  $T_2$ (**t**)

**Grammar Fragment**

...

$X ::= \pi_1\ \pi_2$
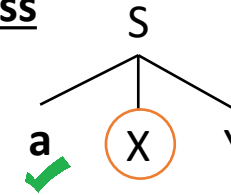
...

# Become the Parser
## *Building Parser Tables*

You need to unfold a nonterminal *X* with lookahead token **t**

Assume there's an *X* production $X ::= \pi_1 \; \pi_2$ (where $\pi_1$ and $\pi_2$ are some kind of symbol)

How do we know to guess this production?

**Parse in Progress**

Lookahead: $T_2$ (**t**)

**Grammar Fragment**

...

$X ::= \pi_1 \; \pi_2$

...

Case 1: $\pi_1$ subtree may start with **t**

Case 2: $\pi_1$ subtree may be empty and $\pi_2$ starts with **t**

# Become the Parser
## *Building Parser Tables*

You need to unfold a nonterminal *X* with lookahead token **t**

Assume there's an *X* production $X ::= \pi_1\ \pi_2$ (where $\pi_1$ and $\pi_2$ are some kind of symbol)

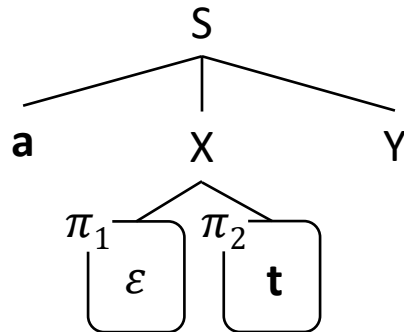How do we know to guess this production?
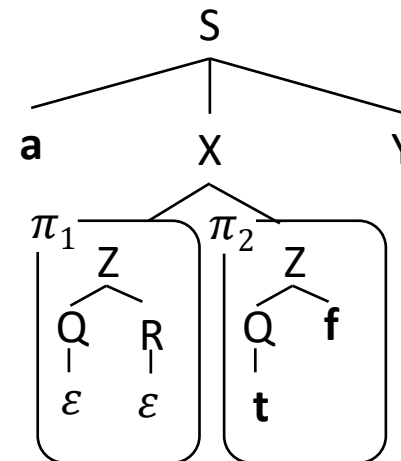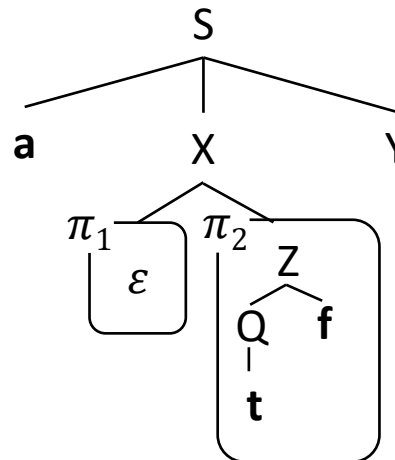
**Parse in Progress**

Lookahead:  $T_2$ (**t**)

**Grammar Fragment**

...

$X ::= \pi_1\ \pi_2$

...

Case 1: $\pi_1$ subtree may start with **t**

Case 2: $\pi_1$ subtree may be empty and $\pi_2$ starts with **t**

Case 3: both $\pi_1$ and $\pi_2$ may be empty and the sibling may start with **t**

# Become the Parser
## *Building Parser Tables*

**Parse in Progress**

You need to unfold a nonterminal *X* with lookahead token **t**



Lookahead: $T_2$ (**t**)

Assume there's an *X* production $X ::= \pi_1\ \pi_2$ (where $\pi_1$ and $\pi_2$ are some kind of symbol)

**Grammar Fragment**

How do we know to guess this production?

...

$X ::= \pi_1\ \pi_2$

...

Case 1: $\pi_1$ subtree may start with **t**

Case 2: $\pi_1$ subtree may be empty and $\pi_2$ starts with **t**

Case 3: both $\pi_1$ and $\pi_2$ may be empty and the sibling may start with **t**

*How can we account for these cases when building the parser?*

# Become the Parser
## Building Parser Tables

**Parse in Progress**

S

a ✓  (X)  Y

Lookahead:  $T_2$ (**t**)

You need to unfold a nonterminal *X* with lookahead token **t**

Assume there's an *X* production $X ::= \pi_1\ \pi_2$ (where $\pi_1$ and $\pi_2$ are some kind of symbol)
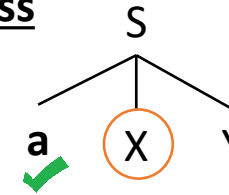
**Grammar Fragment**

...
$X ::= \pi_1\ \pi_2$
...

How do we know to guess this production?

Case 1: $\pi_1$ subtree may start with **t**

Case 2: $\pi_1$ subtree may be empty and $\pi_2$ starts with **t**

Case 3: both $\pi_1$ and $\pi_2$ may be empty and the sibling may start with **t**

**FIRST Sets**

**FOLLOW Sets**

Two sets are sufficient to capture these cases and to build the selector table

# Today's Outline

**Transforming Grammars**

- Fixing LL(1) "near misses"

**Building LL(1) Parsers**

- Reverse-Engineering Selector Tables

- FIRST Sets



**Parsing**

# An Informal Definition

FIRST($\alpha$) = The set of terminals that begin strings derivable
from $\alpha$, and also, if $\alpha$ can derive $\varepsilon$, then $\varepsilon$ is in FIRST(X).

*Building LL(1) Selector Table: FIRST sets, single symbol*

FIRST(α) = The set of terminals that begin strings derivable
from α, and also, if α can derive ε, then ε is in FIRST(X).

Formally, FIRST(α) =

$$\left\{ \hat{\alpha} \,\middle|\, \left( \hat{\alpha} \in \Sigma \wedge \alpha \overset{*}{\Rightarrow} \hat{\alpha}\beta \right) \vee \left( \hat{\alpha} = \varepsilon \wedge \alpha \overset{*}{\Rightarrow} \varepsilon \right) \right\}$$

*Building LL(1) Selector Table: FIRST sets, single symbol*

FIRST($\alpha$) = The set of terminals that begin strings derivable from $\alpha$, and also, if $\alpha$ can derive $\varepsilon$, then $\varepsilon$ is in FIRST(X).

What does the parse tree say about FIRST($A$)?



FIRST(A) includes { **b** }

Again, FIRST(A) includes { **b** }

FIRST(A) Includes { **f** }

FIRST(A) includes { **g** }

FIRST(A) includes { $\varepsilon$ }

If these were the only possible parse trees, then FIRST($A$) = { **b**, f, **g**, $\varepsilon$ }

*Building LL(1) Selector Table: FIRST sets, single symbol*

FIRST($\alpha$) = The set of terminals that begin strings derivable
        from $\alpha$, and also, if $\alpha$ can derive $\varepsilon$, then $\varepsilon$ is in FIRST(X).

**This isn't how you build FIRST sets**

- Looking at parse trees is illustrative for concepts only
- We need to derive FIRST sets directly from the grammar

# Building FIRST Sets: Methodology
## *Building Parser Tables*

**First sets exist for any arbitrary string of symbols** α

- Defined in terms of FIRST sets for a single symbol
  - FIRST of an alphabet terminal
  - FIRST for ε
  - FIRST for a nonterminal
- Use single-symbol FIRST to construct symbol-string FIRSTS

# Rules for Single Symbols
## *Building Parser Tables*

FIRST(X) = The set of terminals that begin strings derivable
from X, and also, if X can derive ε, then ε is in FIRST(X).

**Building FIRST for terminals**

FIRST($t$) = { $t$ } for $t$ in $\Sigma$

FIRST($\varepsilon$) = { $\varepsilon$ }

**Building FIRST($X$) for nonterminal $X$**

For each X ::= $\alpha_1 \alpha_2 \dots \alpha_n$

$C_1$: add FIRST($\alpha_1$) - $\varepsilon$

$C_2$: If $\varepsilon$ could "prefix" FIRST($\alpha_k$), add FIRST($\alpha_k$)- $\varepsilon$

$C_3$: If $\varepsilon$ is in every FIRST set $\alpha_1 \dots \alpha_n$, add $\varepsilon$

# Rules for Single Symbols
### Building LL(1) Parsers

**Building FIRST($X$) for nonterminal $X$**

For each X ::= $\alpha_1$ $\alpha_2$ ... $\alpha_n$

$C_1$: add FIRST($\alpha_1$) - $\varepsilon$

$C_2$: If $\varepsilon$ could "prefix" FIRST($\alpha_k$), add FIRST($\alpha_k$)- $\varepsilon$

$C_3$: If $\varepsilon$ is in every FIRST set $\alpha_1$ ... $\alpha_n$, add $\varepsilon$

# Rules for Single Symbols
## Building LL(1) Parsers

**Building FIRST($X$) for nonterminal $X$**

For each $X ::= \alpha_1\ \alpha_2\ ...\ \alpha_n$

    $C_1$: add FIRST($\alpha_1$) - $\varepsilon$

    $C_2$: If $\varepsilon$ could "prefix" FIRST($\alpha_k$), add FIRST($\alpha_k$)- $\varepsilon$

    $C_3$: If $\varepsilon$ is in every FIRST set $\alpha_1\ ...\ \alpha_n$, add $\varepsilon$

Say there's a production

    $X ::= Y\ Z\ R\ T$

and we know

    FIRST($Y$) = { $\varepsilon$, **a** }

    FIRST($Z$) = { $\varepsilon$, **b, m** }

    FIRST($R$) = { **c** }

    FIRST($T$) = { **d** }

By $C_2$ clause FIRST($X$) includes **b, m** and **c**

**b**,**m**  because FIRST of every symbol before the 2nd includes $\varepsilon$)

*Z in this case* ↗

**c**  because FIRST of every symbol before the 3rd includes $\varepsilon$)

*R in this case* ↗

FIRST($X$) does not add **d** in this clause because not every FIRST set before the T includes $\varepsilon$

# Building FIRST Sets for Symbol Strings

Building LL(1) Parsers

**Building FIRST($\alpha$)**

Let $\alpha$ be composed of symbols $\alpha_1$ $\alpha_2$ ... $\alpha_n$

$C_1$: add FIRST($\alpha_1$) - $\varepsilon$

$C_2$: If $\alpha_1$ ... $\alpha_{k-1}$ is nullable, add FIRST($\alpha_k$) - $\varepsilon$

$C_3$: If $\alpha_1$ ... $\alpha_n$ is nullable, add $\varepsilon$

**Base Cases:**

$\alpha_i$ is is a terminal **t**. Add **t**

$\alpha_i$ is is a nonterminal *X*. Add every leaf symbol that could begin an *X* subtree

**(this gets a bit complicated due to dependencies)**

# Summary: Explored the LL(1) Mindset
## FIRST Sets

**LL(1) "Parseability" Qualification**

- Knowing the leftmost terminal of a parse (sub)tree is enough to pick the next derivation step

**Elusive Conditions**

- Two different rules could start with the same terminal (not left factored)

- The same rule(s) could be applied repeatedly (left recursive)

**Began choosing matching productions to input**

- What terminal could the production be the start of (FIRST)?