# *COS320: Compiling Techniques*

Zak Kincaid

February 23, 2022

- Reminder: HW2 due today
- HW3 on course webpage. Due March 21. **Start early!**
  - You will implement a compiler for a simple imperative programming language (Oat), targeting LLVMlite.
  - You may work individually or in pairs
- Midterm next Thursday
  - Covers material in lectures up to February 23rd (this Wednesday)
    - Interpreters, program transformation, X86, IRs, lexing, parsing
  - How to prepare:
    - Sample exams on Canvas later today
    - Start on HW3
    - Review slides
    - Review example code from lectures (try re-implementing!)
  - Review next Monday: come prepared with questions

*Parsing II: LL parsing*
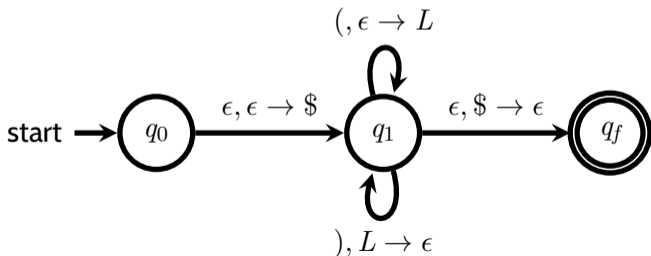
# Recall: Context-free grammars

- A *context-free grammar* $G = (N, \Sigma, R, S)$ consists of:
  - $N$: a finite set of *non-terminal symbols*
  - $\Sigma$: a finite alphabet (or set of *terminal symbols*)
  - $R \subseteq N \times (N \cup \Sigma)^*$ a finite set of *rules* or *productions*
  - $S \in N$: the starting non-terminal.

# Recall: Context-free grammars

- A *context-free grammar* $G = (N, \Sigma, R, S)$ consists of:
  - $N$: a finite set of *non-terminal symbols*
  - $\Sigma$: a finite alphabet (or set of *terminal symbols*)
  - $R \subseteq N \times (N \cup \Sigma)^*$ a finite set of *rules* or *productions*
  - $S \in N$: the starting non-terminal.
- A word $w$ is accepted by $G$ if is derivable in zero or more steps from the starting non-terminal
  - Write $\gamma \Rightarrow \gamma'$ if $\gamma'$ is obtained from $\gamma$ by replacing a non-terminal symbol in $\gamma$ with the right-hand-side of one of its rules
  - Write $\gamma \Rightarrow^* \gamma'$ if $\gamma'$ can be obtained from $\gamma$ using 0 or more derivation steps
  - A word $w \in \Sigma^*$ is accepted by $G$ if $S \Rightarrow^* w$

# Parsing

- Context-free grammars are *generative*: easy to find strings that belongs to $\mathcal{L}(G)$, not so easy determine whether a *given* string belongs to $\mathcal{L}(G)$
- *Pushdown automata* (PDA) are a kind of automata that recognize context-free languages
- Pushdown automaton recognizing <S> ::= <S><S> | (<S>) | $\epsilon$:
  - *Stack alphabet*: $ marks bottom of the stack, $L$ marks unbalanced left paren

# Recall: pushdown automata

- A *push-down automaton* $A = (Q, \Sigma, \Gamma, \Delta, s, F)$ consists of
  - $Q$: a finite set of states
  - $\Sigma$: an (input) alphabet
  - $\Gamma$: a (stack) alphabet
  - $\Delta \subseteq \underbrace{Q}_{\text{source}} \times \underbrace{(\Sigma \cup \{\epsilon\})}_{\text{read input}} \times \underbrace{\Gamma^*}_{\text{read stack}} \times \underbrace{Q}_{\text{dest}} \times \underbrace{\Gamma^*}_{\text{write stack}}$ , the transition relation
  - $s \in Q$: start state
  - $F \subseteq Q$: set of final (accepting) states

# Recall: pushdown automata

- A *push-down automaton* $A = (Q, \Sigma, \Gamma, \Delta, s, F)$ consists of
  - $Q$: a finite set of states
  - $\Sigma$: an (input) alphabet
  - $\Gamma$: a (stack) alphabet
  - $\Delta \subseteq \underbrace{Q}_{\text{source}} \times \underbrace{(\Sigma \cup \{\epsilon\})}_{\text{read input}} \times \underbrace{\Gamma^*}_{\text{read stack}} \times \underbrace{Q}_{\text{dest}} \times \underbrace{\Gamma^*}_{\text{write stack}}$ , the transition relation
  - $s \in Q$: start state
  - $F \subseteq Q$: set of final (accepting) states
- A word $w$ is accepted by $A$ if there is a $w$-labeled accepting path in $A$
  - A *configuration* of $A$ is a pair $(q, v)$ consisting of a state $q \in Q$ and a stack $v \in \Gamma^*$
  - Write $(q, v) \xrightarrow{w} (q', v')$ if there is some $t \in \Gamma^*$ such that $v = at$, $v' = bt$, and $(q, w, a, q', b) \in \Delta$
  - Write $(q, v) \xrightarrow{w}^* (q', v')$ if there is some $w_1, \dots, w_n$ and $(q_1, v_1), \dots, (q_{n-1}, v_{n-1})$ such that $w = w_1 \cdots w_n$ and

$$(q, v) \xrightarrow{w_1} (q_1, v_1) \xrightarrow{w_2} (q_2, v_2) \xrightarrow{w_3} \dots \xrightarrow{w_{n-1}} (q_{n-1}, v_{n-1}) \xrightarrow{w_n} (q', v')$$

  - A word $w$ is accepted iff $(s, \epsilon) \xrightarrow{w}^* (q, a)$ for some $q \in F$.
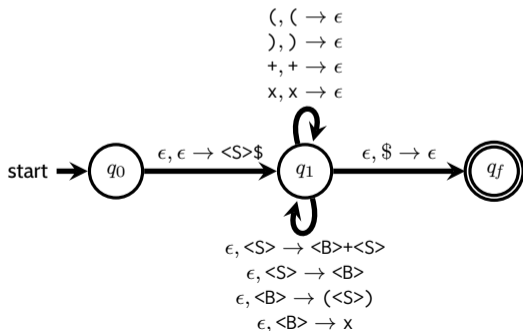
# Context free languages

- Claim: a language is recognized by a context-free grammar if and only if it is recognized by a pushdown automaton
  - Say that a language is *context free* if it is recognized by a context-free grammar (equiv. pushdown automaton).
- Consequence: can "compile" context-free grammars to pushdown automata in order to implement parsers

# Context free languages

- Claim: a language is recognized by a context-free grammar if and only if it is recognized by a pushdown automaton
  - Say that a language is *context free* if it is recognized by a context-free grammar (equiv. pushdown automaton).
- Consequence: can "compile" context-free grammars to pushdown automata in order to implement parsers
- Two strategies, which correspond to different ways to implement parsers:
  - Top-down (LL parsing)
  - Bottom-up (LR parsing)

# Top-down parsing

- Stack represents intermediate state of a derivation, minus the consumed part of the input string.
- Start with $S$ on the stack
- Any time top of the stack is a non-terminal $A$, non-deterministically choose a rule $A ::= \gamma \in R$. Pop $A$ off the stack, and push $\gamma$
- If the top of the stack is a terminal $a$, consume $a$ from the input string and pop $a$ off the stack
- Accept when stack is empty

```
<S> ::= <B>+<S> | <B>
<B> ::= (<S>) | x
```

$$(, ( \rightarrow \epsilon$$
$$), ) \rightarrow \epsilon$$
$$+, + \rightarrow \epsilon$$
$$x, x \rightarrow \epsilon$$

start $\rightarrow q_0 \quad \xrightarrow{\epsilon, \epsilon \rightarrow <S>\$} \quad q_1 \quad \xrightarrow{\epsilon, \$ \rightarrow \epsilon} \quad q_f$

$$\epsilon, <S> \rightarrow <B>+<S>$$
$$\epsilon, <S> \rightarrow <B>$$
$$\epsilon, <B> \rightarrow (<S>)$$
$$\epsilon, <B> \rightarrow x$$

$<S> ::= <B>+<S> \mid <B>$

$<B> ::= (<S>) \mid x$



| State | Stack | Input |
|---|---|---|
| $q_0$ | $\epsilon$ | (x+x)+x |
| $q_1$ | <S>\$ | (x+x)+x |
| $q_1$ | <B>+<S>\$ | (x+x)+x |
| $q_1$ | (<S>)+<S>\$ | (x+x)+x |
| $q_1$ | <S>)+<S>\$ | x+x)+x |
| $q_1$ | <B>+<S>)+<S>\$ | x+x)+x |
| $q_1$ | x+<S>)+<S>\$ | x+x)+x |
| $q_1$ | +<S>)+<S>\$ | +x)+x |
| $q_1$ | <S>)+<S>\$ | x)+x |
| $q_1$ | <B>)+<S>\$ | x)+x |
| $q_1$ | x)+<S>\$ | x)+x |
| $q_1$ | )+<S>\$ | )+x |
| $q_1$ | +<S>\$ | +x |
| $q_1$ | <S>\$ | x |
| $q_1$ | <B>\$ | x |
| $q_1$ | x\$ | x |
| $q_1$ | \$ | $\epsilon$ |
| $q_f$ | $\epsilon$ | $\epsilon$ |

# Bottom-up parsing

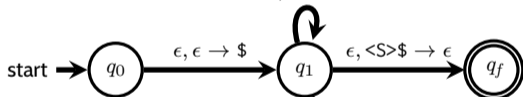- Stack holds a word in $(N \cup \Sigma)^*$ such that it is possible to derive the part of the input string that has been consumed from its reverse.
- At any time, may read a letter from input string and push it on top of the stack
- At any time, may non-deterministically choose a rule $A ::= \gamma_1 \ldots \gamma_n$ and apply it in reverse: pop $\gamma_n \ldots \gamma_1$ off the top of the stack, and push $A$.
- Accept when stack just contains start non-terminal

$$<S> ::= <B>+<S> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

<S> ::= <B>+<S> | <B>

<B> ::= (<S>) | x

$$(, \epsilon \to ($$
$$), \epsilon \to )$$
$$+, \epsilon \to +$$
$$x, \epsilon \to x$$

start $\to$ $q_0$ $\xrightarrow{\epsilon, \epsilon \to \$}$ $q_1$ $\xrightarrow{\epsilon, <S>\$ \to \epsilon}$ $q_f$

$$\epsilon, <S>+<B> \to <S>$$
$$\epsilon, <B> \to <S>$$
$$\epsilon, )<S>( \to <B>$$
$$\epsilon, x \to <B>$$

| State | Stack | Input |
|---|---|---|
| $q_0$ | $\epsilon$ | (x+x)+x |
| $q_1$ | $ | (x+x)+x |
| $q_1$ | ($ | x+x)+x |
| $q_1$ | x($ | +x)+x |
| $q_1$ | <B>($ | +x)+x |
| $q_1$ | +<B>($ | x)+x |
| $q_1$ | x+<B>($ | )+x |
| $q_1$ | <B>+<B>($ | )+x |
| $q_1$ | <S>+<B>($ | )+x |
| $q_1$ | <S>($ | )+x |
| $q_1$ | )<S>($ | +x |
| $q_1$ | <B>$ | +x |
| $q_1$ | +<B>$ | x |
| $q_1$ | x+<B>$ | $\epsilon$ |
| $q_1$ | <B>+<B>$ | $\epsilon$ |
| $q_1$ | <S>+<B>$ | $\epsilon$ |
| $q_1$ | <S>$ | $\epsilon$ |
| $q_f$ | $\epsilon$ | $\epsilon$ |

# Parsing overview

- Basic problem with both top-down and bottom-up construction: *non-determinism*
  - Non-deterministic search is inefficient
    - E.g., consider <S> ::= <S>a | <S>b | $\epsilon$. Top-down parser must "guess" the entire input string at the beginning (breadth-first backtracking search takes exponential time in length of input string, depth-first does not terminate).
  - Algorithms for parsing any context free grammar in cubic[1] time, based on dynamic programming (Earley, and Cocke-Younger-Kasami).

---

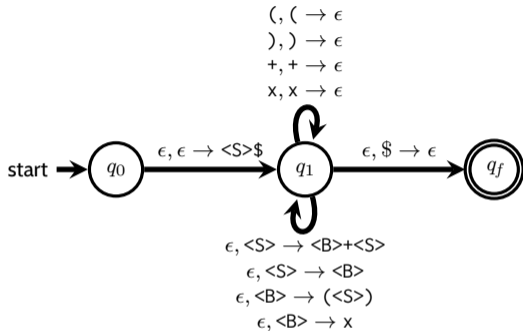[1]Also sub-cubic galactic algorithms: Valiant 1975

# Parsing overview

- Basic problem with both top-down and bottom-up construction: *non-determinism*
  - Non-deterministic search is inefficient
    - E.g., consider <S> ::= <S>a | <S>b | $\epsilon$. Top-down parser must "guess" the entire input string at the beginning (breadth-first backtracking search takes exponential time in length of input string, depth-first does not terminate).
  - Algorithms for parsing any context free grammar in cubic[1] time, based on dynamic programming (Earley, and Cocke-Younger-Kasami).
- Parser generators use these same ideas, but restricted to cases where we can eliminate non-determinism.
- Possible for both top-down and bottom-up style
  - **Today:** *LL* (*L*eft-to-right, *L*eftmost derivation) parsers: top-down
    - Easy to understand & write by hand
  - **Next time:** *LR* (*L*eft-to-right, *R*ightmost derivation) parsers: bottom-up
    - More general, (variations) implemented in parser generators

---

[1]Also sub-cubic galactic algorithms: Valiant 1975

# LL parsing

<S> ::= <B>+<S> | <B>
<B> ::= (<S>) | x



- "Any time top of the stack is a non-terminal $A$, non-deterministically choose a production $A ::= \gamma \in R$. Pop $A$ off the stack, and push $\gamma$"
    - Key problem: need to deterministically choose which production to use
    - Solution: Look at the next input symbol, but don't consume it (*lookahead*)
        - This is $LL(1)$ parsing. $LL(k)$ allows $k$ lookahead tokens

- We say that a grammar is $LL(k)$ if when we look ahead $k$ symbols in a top-down parser, we know which rule we should apply.
  - Let $G = (N, \Sigma, R, S)$ be a grammar. $G$ is $LL(k)$ iff: for any $S \Rightarrow^* \alpha A \beta$, for any word $w \in \Sigma^k$, if there is some $A ::= \gamma \in R$ such that $\gamma \beta \Rightarrow^* w \beta'$ (for some $\beta'$), then $\gamma$ is unique.
- Not every context-free language has an $LL(k)$ grammar.
  - $\{a^i b^j : i = j \lor 2i = j\}$ is not $LL(k)$ for any $k$
- Which of the following are $LL(1)$ grammars?
  - <S> ::= a<S> | b<S> | $\epsilon$

  - <S> ::= <S>a | <S>b | $\epsilon$
  - <S> ::= <B>+<S> | <B>
    <B> ::= (<S>) | x

- We say that a grammar is $LL(k)$ if when we look ahead $k$ symbols in a top-down parser, we know which rule we should apply.
  - Let $G = (N, \Sigma, R, S)$ be a grammar. $G$ is $LL(k)$ iff: for any $S \Rightarrow^* \alpha A \beta$, for any word $w \in \Sigma^k$, if there is some $A ::= \gamma \in R$ such that $\gamma\beta \Rightarrow^* w\beta'$ (for some $\beta'$), then $\gamma$ is unique.
- Not every context-free language has an $LL(k)$ grammar.
  - $\{a^i b^j : i = j \lor 2i = j\}$ is not $LL(k)$ for any $k$
- Which of the following are $LL(1)$ grammars?
  - <S> ::= a<S> | b<S> | $\epsilon$
    More generally, any grammar that results from our DFA→CFG conversion
  - <S> ::= <S>a | <S>b | $\epsilon$
  - <S> ::= <B>+<S> | <B>
    <B> ::= (<S>) | x

# Left-factoring

- The grammar

  $$\text{<S> ::= <B>+<S> | <B>}$$
  $$\text{<B> ::= (<S>) | x}$$

  is not LL(1): ( lookahead can't distinguish the two <S> rules
- However, there is an LL(1) grammar for the language

# Left-factoring

- The grammar

$$\langle S\rangle ::= \langle B\rangle+\langle S\rangle \mid \langle B\rangle$$
$$\langle B\rangle ::= (\langle S\rangle) \mid x$$

  is not LL(1): ( lookahead can't distinguish the two <S> rules
- However, there is an LL(1) grammar for the language

$$\langle S\rangle ::= \langle B\rangle\langle R\rangle$$
$$\langle R\rangle ::= +\langle S\rangle \mid \epsilon$$
$$\langle B\rangle ::= (\langle S\rangle) \mid x$$

- General strategy: factor out rules with common prefixes ("left factoring")

# Eliminating left recursion

- A grammar is **left-recursive** if there is a non-terminal $A$ such that $A \Rightarrow^+ A\gamma$ (for some $\gamma$)
- Left-recursive grammars are not $LL(k)$ for any $k$
- Consider:

$$<S> ::= <S>+<B> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

# Eliminating left recursion

- A grammar is **left-recursive** if there is a non-terminal $A$ such that $A \Rightarrow^+ A\gamma$ (for some $\gamma$)
- Left-recursive grammars are not $LL(k)$ for any $k$
- Consider:

$$<S> ::= <S>+<B> \mid <B>$$
$$<B> ::= (<S>) \mid x$$

Can remove left recursion as follows:

$$<S> ::= <B><S'>$$
$$<S'> ::= +<B><S'> \mid \epsilon$$
$$<B> ::= (<S>) \mid x$$

(Recognizes the same language, but parse trees are different!)

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'.S \Rightarrow \gamma A a \gamma'\}$

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'. S \Rightarrow \gamma A a \gamma'\}$
- Transition table $\delta$ for $G$ can be computed using **first**, **follow**, and **nullable**:
    1. For each non-terminal $A$ and letter $a$, initialize $\delta(A, a)$ to $\emptyset$
    2. For each rule $A ::= \gamma$
        - Add $\gamma$ to $\delta(A, a)$ for each $a \in$ **first**$(\gamma)$
        - If $\gamma$ is nullable, add $\gamma$ to $\delta(A, a)$ for each $a \in$ **follow**$(A)$

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'.S \Rightarrow \gamma A a \gamma'\}$
- Transition table $\delta$ for $G$ can be computed using **first**, **follow**, and **nullable**:
  1. For each non-terminal $A$ and letter $a$, initialize $\delta(A, a)$ to $\emptyset$
  2. For each rule $A ::= \gamma$
     - Add $\gamma$ to $\delta(A, a)$ for each $a \in$ **first**$(\gamma)$
     - If $\gamma$ is nullable, add $\gamma$ to $\delta(A, a)$ for each $a \in$ **follow**$(A)$
- $G$ is $LL(1)$ iff $\delta(A, a)$ is empty or singleton for all $A$ and $a$

# Mechanical construction of LL(1) parsers

- Fix a grammar $G = (N, \Sigma, R, S)$
- For any word $\gamma \in (N \cup \Sigma)^*$, define **first**$(\gamma) = \{a \in \Sigma : \gamma \Rightarrow^* aw\}$
- For any word $\gamma \in (N \cup \Sigma)^*$, say that $\gamma$ is **nullable** if $\gamma \Rightarrow^* \epsilon$
- For any non-terminal $A$, define **follow**$(A) = \{a \in \Sigma : \exists \gamma, \gamma'. S \Rightarrow \gamma A a \gamma'\}$
- Transition table $\delta$ for $G$ can be computed using **first**, **follow**, and **nullable**:
  1. For each non-terminal $A$ and letter $a$, initialize $\delta(A, a)$ to $\emptyset$
  2. For each rule $A ::= \gamma$
     - Add $\gamma$ to $\delta(A, a)$ for each $a \in$ **first**$(\gamma)$
     - If $\gamma$ is nullable, add $\gamma$ to $\delta(A, a)$ for each $a \in$ **follow**$(A)$
- $G$ is $LL(1)$ iff $\delta(A, a)$ is empty or singleton for all $A$ and $a$
- Operation of the parser on a word $w$:
  - Start with stack <S>
  - While $w$ not empty
    - If top of the stack is a terminal $a$ and $w = aw'$, pop and set $w = w'$
    - If top of the stack is a non-terminal $A$ and $w = aw'$, pop and push (singleton) $\delta(A, w)$ (or reject if $\delta(A, w)$ is empty)
  - Accept if stack is empty; reject otherwise.

# Computing nullable

- **nullable** is the *smallest set* of non-terminals such that if there is some $A ::= \gamma_1 \ldots \gamma_n \in R$ with $\gamma_1, \ldots, \gamma_n \in$ **nullable** implies $A \in$ **nullable**
    - Fixpoint computation:
        - $\textbf{nullable}_0 = \emptyset$
        - $\textbf{nullable}_{i+1} = \{A : \exists \gamma_1, \ldots, \gamma_n \in \textbf{nullable}_i . A ::= \gamma_1 \ldots \gamma_n \in R\}$
        - $\textbf{nullable} = \bigcup\limits_{i=0}^{\infty} \textbf{nullable}_i$

    nullable $\leftarrow \emptyset$;
    changed $\leftarrow$ true;
    **while** *changed* **do**
        changed $\leftarrow$ false;
        **for** $A := \gamma_1 \ldots \gamma_n \in R$ **do**
            **if** $A \notin$ *nullable* $\land \gamma_1, \ldots, \gamma_n \in$ *nullable* **then**
                *nullable* $\leftarrow$ *nullable* $\cup \{A\}$;
                changed $\leftarrow$ true;

# Computing nullable

- **nullable** is the *smallest set* of non-terminals such that if there is some $A ::= \gamma_1 \ldots \gamma_n \in R$ with $\gamma_1, \ldots, \gamma_n \in$ **nullable** implies $A \in$ **nullable**
  - Fixpoint computation:
    - $\text{nullable}_0 = \emptyset$
    - $\text{nullable}_{i+1} = \{A : \exists \gamma_1, \ldots, \gamma_n \in \text{nullable}_i . A ::= \gamma_1 \ldots \gamma_n \in R\}$
    - $\text{nullable} = \bigcup\limits_{i=0}^{\infty} \text{nullable}_i$

  ```
  nullable ← ∅;
  changed ← true;
  while changed do
      changed ← false;
      for A := γ₁ … γₙ ∈ R do
          if A ∉ nullable ∧ γ₁, …, γₙ ∈ nullable then
              nullable ← nullable ∪ {A};
              changed ← true;
  ```

- Fixpoint computations appear everywhere!
  - Later we will see how they are used in dataflow analysis

# Computing first and follow

- **first** is the *smallest function*[2] such that
  - For each $a \in \Sigma$, $\textbf{first}(a) = \{a\}$
  - For each $A ::= \gamma_1 \ldots \gamma_i \ldots \gamma_n \in R$, with $\gamma_1, \ldots, \gamma_{i-1}$ nullable, $\textbf{first}(A) \supseteq \textbf{first}(\gamma_i)$
- **follow** is the *smallest function* such that
  - For each $A ::= \gamma_1 \ldots \gamma_i \ldots \gamma_n \in R$, with $\gamma_{i+1}, \ldots, \gamma_n$ nullable, $\textbf{follow}(\gamma_i) \supseteq \textbf{follow}(A)$
  - For each $A ::= \gamma_1 \ldots \gamma_i \ldots \gamma_j \ldots \gamma_n \in R$, with $\gamma_{i+1}, \ldots, \gamma_{j-1}$ nullable, $\textbf{follow}(\gamma_i) \supseteq \textbf{first}(\gamma_j)$
- Both can be computed using a fixpoint algorithm, like **nullable**

---

[2]Pointwise order: $f \leq g$ if for all $x$, $f(x) \leq g(x)$