Langages de Programmation, Interprétation, Compilation

Christine Paulin

Christine.Paulin@universite-paris-saclay.fr

Département Informatique, Faculté des Sciences d'Orsay, Université Paris-Saclay

LDD3 Informatique, Mathématiques - Magistère Informatique

2025-26

Analyse syntaxique d'un langage impératif



Analyse syntaxique d'un langage impératif

- Affichage d'une expression
- Grammaires et dérivations
- Analyse récursive descendante
- Analyseur complet pour IMP
- Présentation alternative des dérivations
- Construction d'un analyseur

Affichage d'une expression avec notations infixes

- Transformer un arbre de syntaxe abstraite en une chaîne représentant la même expression
- La structure d'arbre supprime le besoin de parenthèses
- L'affichage de l'arbre doit remettre des parenthèses
 - partout? juste celles qui sont nécessaires?
 - notion de priorité entre les opérateurs et d'associativité
 - plusieurs fonctions d'affichage de type expr -> string en fonction des priorités

Grammaires

La grammaire d'un langage décrit les structures de phrases légales dans ce langage.

On la définit à l'aide de deux éléments en plus du lexique :

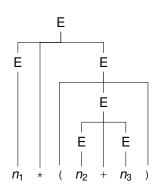
- les différentes catégories de (fragments de) phrases du langage,
- les manières dont différents éléments peuvent être regroupés pour former un fragment de l'une ou l'autre des catégories (données par des règles de combinaisons).

Les phrases bien formées sont alors celles dont les mots peuvent être regroupés d'une manière compatible avec les règles de combinaison.

Structure grammaticale d'une phrase

- les feuilles sont les mots de la phrase (dans l'ordre),
- chaque nœud interne dénote le regroupement de ses fils

$$n_1 * (n_2 + n_3)$$
.



Grammaire algébrique

Une grammaire algébrique est un quadruplet (T, N, S, R) où :

- T est un ensemble de symboles terminaux, désignant des mots,
- N est un ensemble de symboles <u>non terminaux</u>, désignant des groupes de mots,
- S est le symbole de départ, qui décrit une phrase complète,
- R est un ensemble de règles de production
 - paires $X \prec a_1 \ldots a_k$ avec X symbole non terminal et $a_1 \ldots a_k$ de symboles (terminaux ou non) formant un groupe de type X correct (qui peut être vide).

Grammaire simple pour les expressions arithmétiques

- les symboles terminaux +, *, (,) ainsi qu'un terminal n pour une constante entière,
- le symbole non terminal E,
- le symbole de départ E,
- les règles de production E ≺ n, E ≺ E+E, E ≺ E * E et E ≺ (E), encore écrites

Arbre de dérivation syntaxique

- nœud interne : symbole de N,
- feuille : symbole de T,
- si un nœud interne a l'étiquette $X \in N$, les étiquettes de ses fils prises dans l'ordre forment une séquence β telle que $X \prec \beta \in R$

Type Caml des arbres de dérivation (paramétré par le type des terminaux) :

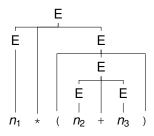
```
type 't deriv =
          Node of string * 't deriv list
          | Leaf of 't
```

Une phrase m est <u>dérivable</u> à partir d'un symbole non terminal X s'il existe un arbre de dérivation dont la racine est étiquetée par X et dont les feuilles, prises dans l'ordre, forment cette phrase.

Les phrases <u>dérivables</u> de la grammaire sont celles qui sont dérivables à partir du symbole de départ *S*.

Exemple

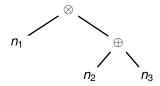
Arbre de dérivation pour l'expression arithmétique $n_1 \star (n_2 + n_3)$ avec la grammaire naïve



Arbre de syntaxe abstraite

Mul(Cst n1, Add(Cst n2, Cst n3))

Ne pas confondre arbre de dérivation avec l'arbre de syntaxe abstraite

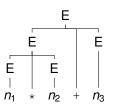


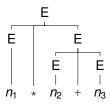
- L'analyseur syntaxique construit l'arbre de syntaxe abstraite.
- L'arbre de dérivation décrit les différentes étapes de l'analyse réalisée, mais n'est pas construit explicitement par l'analyseur.

Ambiguïtés

Une grammaire est <u>ambiguë</u> s'il existe une phrase qui peut être dérivée par deux arbres différents.

Pour la phrase $n_1 * n_2 + n_3$:



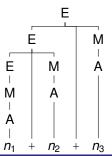


Grammaire non-ambigüe

Reconnait le même langage

$$E ::= E + M$$
 $| M |$
 $M ::= M * A$
 $| A ::= n$
 $| (E)$

Arbre de dérivation pour la phrase $n_1+n_2+n_3$.



Partie 1: Introduction

- Analyse syntaxique d'un langage impératif
 - Affichage d'une expression
 - Grammaires et dérivations
 - Analyse récursive descendante
 Analyseur complet pour IMP
 - Présentation alternative des dérivation
 - Construction d'un analyseur

Analyseur syntaxique

L'analyseur syntaxique

- prend en entrée une séquence de lexèmes,
- doit vérifier que cette séquence est cohérente avec les règles de grammaire du langage et renvoie
 - en cas de succès, un arbre de syntaxe abstraite donnant la structure de la séquence lue en entrée,
 - en cas d'échec, la localisation et l'explication de ce qui est grammaticalement incohérent dans l'entrée.

Analyse descendante

- Construire l'arbre de la racine vers les feuilles pour arriver à la suite de terminaux passée en entrée. On parle d'analyse descendante, ou top-down
- On définit pour chaque symbole non terminal X, une fonction f_X :
 - prend en entrée une séquence de lexèmes I,
 - cherche dans un préfixe de I une phrase m pour une règle $X \prec m$,
 - renvoie les lexèmes qui n'ont pas été utilisés
 - échoue si elle ne parvient pas à reconnaître de fragment m pour aucune règle

Construction manuelle d'un analyseur

Grammaire non ambiguë des expressions arithmétiques

$$E ::= E + M \\ | M |$$

$$M ::= M * A \\ | A |$$

$$A ::= n \\ | (E)$$

Type de données pour les lexèmes

Reconnaissance d'un atome

Reconnaître un atome en préfixe d'une suite de léxèmes :

```
let validate_atom: token list -> token list = function
```

- Le premier léxème est soit une constante soit une parenthèse
- Si constante n, on réussit et on renvoie le reste de la liste
- Si parenthèse ouvrante (, alors la règle qui peut s'appliquer est $A \prec (E)$. Reconnaître une expression puis vérifier la parenthèse fermante

Changement de grammaire

Pouvoir faire le choix de la règle en regardant le premier léxème de l'entrée

Partie 1: Introduction

- Analyse syntaxique d'un langage impératif
 - Affichage d'une expression
 - Grammaires et dérivations
 - Analyse récursive descendante
 - Analyseur complet pour IMP
 - Construction d'un analyseur

Analyseur complet pour IMP: symboles

Symboles terminaux:

- un symbole *n* pour chaque constante entière,
- un symbole x pour chaque nom de variable,
- les symboles +, *, <, &&, (,), {, }, :=, ;,
- les mots-clés print, while, if, else.

Symboles non-terminaux

- Symbole de départ P : une séquence d'instructions.
- Symbole S pour une suite d'instructions (peut être vide)
- Symbole I pour une instruction
- Symbole E pour une expression (opération binaire ou expression atomique)
- Symbole A pour une expression atomique (constante, variable ou expression parenthésée)

Analyseur complet pour IMP : Règles

```
I ::= x := E;
     | while (E) {S}
| if (E) {S} else {S}
        \mid print (E);
egin{array}{ccccc} E & | & A & & & \\ & | & E o E & & & avec \ o \in \{\,+,\,\star,\,<,\,\&\,\&\,\} \end{array}
```

Analyse lexicale

```
type token =
   NUM of int
   | PRINT | WHILE | IF | ELSE
    EOI
let tokenizer (s: string): unit -> token =
  let i = ref 0 in
  let rec next token () = match s.[!i] with
   | '+' -> incr i; PLUS
   | '*' -> incr i; STAR
   ':' \rightarrow if s.[!i+1] = '=' then (i := !i+2; SET)
           else failwith "lexical error"
   | ``&` -> if s.[!i+1] = ``&` then (i := !i+2; AND)
           else failwith "lexical error"
```

Analyse lexicale mots clés

```
| 'a'...'z' \rightarrow let k = scan word 1 in
                 let t = String.sub s !i k in
                 i := !i + k;
                  (match t with
                   | "print" -> PRINT
                   | "while" -> WHILE
                   | "if" -> IF
                   | "else" -> ELSE
                    -> VAR t)
  | ' ' | ' \mid n' \rightarrow incr i; next token ()
   c -> failwith (Printf.sprintf "unknown character %c" c)
    exception e -> EOI
and scan word k = match s.[!i+k] with
    a'...'z' \rightarrow scan word (k+1)
    -> k
    exception e -> k
```

Syntaxe abstraite

```
type prog = instr list
and instr =
      Print of expr
     Set of string * expr
    | While of expr * instr list
     If of expr * instr list * instr list
and expr =
     Cst of int
     Var of string
      Bop of bop * expr * expr
and bop = Add | Mul | Lt | And
On associe à chaque opérateur une priorité
let prio = function
     AND -> 1 | LT -> 2 | PLUS -> 3 | STAR -> 4
      -> failwith "prio: expected operator"
let root prio = 0
```

Analyse syntaxique

- Une fonction interne tokenizer qui lit le prochain léxème
- une référence current_token pour le léxème courant
- next : unit -> unit : mise à jour par par le prochain léxème
- take token : unit -> token : renvoie le léxème courant et lit le suivant
- expect (s : token) : unit : vérifie que le léxème courant est s et lit le suivant

Analyse descendante : suite d'instructions

Suite d'instructions, peut se terminer par la fin de l'entrées EOI ou par l'accolade fermante

Analyse descendante: instruction

Pour une instruction, le premier léxème détermine la règle à appliquer

```
and parse instr () = match take token() with
   | VAR x -> expect SET:
              let e = parse expr root prio in
              expect SEMI:
              Set(x, e)
     PRINT -> let e = parse cond () in
                 expect SEMI:
                  Print(e)
     WHILE -> let e = parse cond () in
                  let is = parse block () in
                 While (e, is)
    IF \rightarrow let e = parse cond () in
                  let is1 = parse block () in
                 expect ELSE;
                  let is2 = parse_block () in
                  If (e, is1, is2)
                  -> failwith "syntax_error"
```

Analyse descendante: instruction (suite)

```
and parse_cond () =
    expect LPAR;
    let e = parse_expr root_prio in
    expect RPAR;
    e
and parse_block () =
    expect BEGIN;
    let is = parse_instrs () in
    expect END;
    is
```

Analyse descendante : expressions

Utilisation de fonctions paramétrées par le niveau de priorité

- parse_expr p : reconnaît une expression dont l'opérateur principal est d'une priorité > p
- expand_op p : cherche à compléter l'expression par des opérateurs de priorité > p
- initialisé par root_prio (plus petit que tous les opérateurs)

(on parle d'analyse par priorit és ascendantes). À l'initialisation, on utilise une priorité root_prio qui est le niveau le plus bas.

Analyse descendante : expressions (code)

```
and parse expr p = parse a expr () |> expand op p
and expand op p e1 = match ! current token with
    (PLUS | STAR | LT | AND as op) when (prio op > p) -
     next();
     let e2 = parse expr (prio op) in
     let e = match op with
        PLUS -> Bop(Add, e1, e2)
      | STAR -> Bop(Mul, e1, e2)
       LT \rightarrow Bop(Lt, e1, e2)
        AND \rightarrow Bop(And, e1, e2)
        _ -> assert false
     in
     expand op p e
  | -> e1
```

Partie 1: Introduction

- Analyse syntaxique d'un langage impératif
 - Affichage d'une expression
 - Grammaires et dérivations
 - Analyse récursive descendante
 - Analyseur complet pour IMPPrésentation alternative des dérivations
 - Construction d'un analyseur

Dérivations

- grammaire (T, N, S, R)
- suite de phrases, du symbole de départ à une suite de terminaux
- Définition : $m \rightarrow m'$ si
 - $m = m_1 X m_2, m' = m_1 a_1 \dots a_n m_2$
 - $X \prec a_1 \dots a_k \in R$
- Fermeture transitive : $m \rightarrow^+ m'$ (au moins une étape)
- Fermeture réflexive-transitive : $m \rightarrow^* m'$ (0 ou plusieurs étapes)
- Une séquence m est dérivable lorsqu'il existe une dérivation $S \rightarrow^* m$.
- On parle de <u>dérivation gauche</u> lorsque on expanse à chaque étape le symbole non terminal le plus à gauche.
- Il n'y a plus de dérivation possible à partir d'une suite de terminaux
- Plusieurs dérivation peuvent donner le même arbre (mais une seule dérivation gauche)

Exemples

E*(E) se dérive en E*(E+E) en une étape

$$m_1 = E*($$
 $m_2 =)$
 $X = E$
 $\beta = E+E$

Dérivation de $n_1 * (n_2 + n_3)$:

Partie 1: Introduction

- Analyse syntaxique d'un langage impératif
 - Affichage d'une expression
 - Grammaires et dérivations
 - Analyse récursive descendante
 - Analyseur complet pour IMP
 - Présentation alternative des dérivations
 - Construction d'un analyseur

Construction systématique d'un analyseur descendant

- Systématiser la construction de l'analyseur à partir de la grammaire
- Pour un non terminal X choisir la règle X

 β à expanser de manière déterministe en fonction de l'entrée
 - possible seulement si la grammaire est non ambigüe
- Faire ce choix en fonction uniquement du premier caractère de l'entrée!

(Couvert dans le cours de PIL de L2)

Premiers et annulables

- grammaire (T, N, S, R),
- séquence $\alpha \in (T \cup N)^*$ de symboles terminaux ou non terminaux.
- On note Premiers(α) l'ensemble des symboles terminaux qui peuvent se trouver en tête d'une séquence dérivée à partir de α .

Premiers(
$$\alpha$$
) = { $a \in T \mid \exists \beta, \alpha \rightarrow^* a\beta$ }

 On note Annulables l'ensemble des symboles non terminaux à partir desquels il est possible de dériver la séquence vide :

$$\mathsf{Annulables} = \{\, \textit{X} \in \textit{N} \mid \textit{X} \rightarrow^* \varepsilon \,\}$$

Caractérisation récursive

Premiers

```
\begin{array}{lll} \operatorname{Premiers}(\varepsilon) &=& \emptyset \\ \operatorname{Premiers}(a\beta) &=& \{a\} & \operatorname{si} \ a \in T \\ \operatorname{Premiers}(X) &=& \bigcup_{X \prec \beta \in R} \operatorname{Premiers}(\beta) \\ \operatorname{Premiers}(X\beta) &=& \left\{ \begin{array}{ll} \operatorname{Premiers}(X) & \operatorname{si} \ X \not \in \operatorname{Annulables} \\ \operatorname{Premiers}(X) \cup \operatorname{Premiers}(\beta) & \operatorname{si} \ X \in \operatorname{Annulables} \end{array} \right. \end{array}
```

Annulables Le symbole X est annulable :

- s'il existe une règle X ≺ ε, ou
- s'il existe une règle $X \prec X_1 \dots X_k$ où tous les X_i sont des symboles non terminaux qui sont également annulables.

Equations récursives dont on cherche le plus petit point fixe.

Théorème de point fixe de Tarski (fini)

Si on considère un ensemble ordonné A <u>fini</u> et doté d'un plus petit élément \bot , et une fonction $F: A \to A$ croissante, alors :

- la fonction F admet au moins un point fixe, c'est-à-dire qu'il existe un élément x ∈ A tel que F(x) = x, et
- en itérant F à partir de ⊥, c'est-à-dire en calculant F(⊥), puis F(F(⊥)), puis F(F(F(⊥))), etc, on finit nécessairement par trouver un point fixe de F (et même plus précisément : le plus petit des points fixes).

Ce théorème est vrai dans un cadre plus large (existence d'une borne sup pour toute chaîne croissante)

Calcul des annulables

- On cherche un ensemble de symboles non terminaux.
- A est l'ensemble des parties de N, ordonné par l'ordre d'inclusion.
- Le plus petit élément est ∅.
- La fonction croissante est $F : \mathcal{P}(N) \to \mathcal{P}(N)$ définie par

$$F(E) = \{ X \mid X \to \varepsilon \} \cup \{ X \mid X \to X_1 \dots X_k, \ X_i \in E \}$$

- Exemple des expressions arithmétiques : F(∅) = { E', M' } et
 F({ E', M' }) = { E', M' }.
 Le point fixe est { E', M' }, seuls symboles non terminaux annulables.
- Présenté dans un tableau de booléens donnant les annulables trouvés à chaque itération.

On s'arrête lorsque deux lignes consécutives sont identiques.

	E	E'	М	M'	<i>A</i>
0.	F	F	F	F	F
1.	F	٧	F	V	F
2.	F	٧	F	V	F

Calcul des premiers

- Calculer pour chaque symbole non terminal, un ensemble de symboles terminaux.
- On considère l'ensemble $\mathcal{P}(T) \times \ldots \times \mathcal{P}(T)$ (avec une composante par symbole non terminal), ordonné par le produit de l'ordre inclusion.
- Son plus petit élément est $(\emptyset, \dots, \emptyset)$,
- la fonction croissante utilisée est celle qui recalcule les informations pour chaque symbole non terminal en fonction des informations déjà connues.
- Présentation dans un tableau, chaque colonne correspond à un symbole non terminal, et chaque ligne à une itération.
- Le calcul s'arrête lorsque deux lignes deviennent identiques.

	E	E'	М	M'	A
0.	Ø	Ø	Ø	Ø	Ø
1.	Ø	{+}	Ø	{ * }	{ n , (}
2.	Ø	{+}	{ n , (}	{ * }	$\mid \{ n, (\} \mid$
3.	{ n , (}	{+}	{ n , (}	{ * }	{ n , (}
4.	{ n , (}	{+}	{ n , (}	{ * }	{ n , (}

Suivants

- une règle X ≺ β peut s'appliquer si le premier terminal de l'entrée appartient à Premiers(β)
- que se passe-t-il si $\beta \to^* \varepsilon$ (en particulier si $\beta = \varepsilon$)?
- le premier symbole de l'entrée est alors un symbole qui peut <u>suivre</u> X dans une dérivation
- En analysant la grammaire, on peut caractériser l'ensemble des <u>suivants</u> d'un symbole non terminal $X \in N$: peut apparaître immédiatement après X dans une dérivation.

Suivants(
$$X$$
) = { $a \in T \mid \exists \alpha, \beta, S \rightarrow^* \alpha X a \beta$ }

- Exemple des expressions arithmétiques :
 - $S \rightarrow M E' \rightarrow M + E$ et

$$S \rightarrow M E' \rightarrow M + E \rightarrow A M' + E \rightarrow A \varepsilon + E = A + E$$

 $+ \in Suivants(M)$ et $+ \in Suivants(A)$

Equations pour les suivants

Pour calculer l'ensemble des suivants de X, on regarde toutes les apparitions de X dans un membre droit de règle $Y \prec \alpha X \beta$

- on inclut dans Suivants(X):
 - les symboles terminaux qui peuvent apparaître en tête d'une séquence dérivée à partir de β , c'est-à-dire Premiers(β),
 - dans le cas où $\beta \to^* \varepsilon$ tous les suivants de Y peuvent également suivre X

On a à nouveau des équations récursives, à résoudre par la même technique que précédemment.

Calcul des suivants

On ajoute le symbole terminal spécial # de fin d'entrée (correspond à EOI), il appartient aux suivants du symbole de départ S.

Dans notre exemple des expressions arithmétiques :

```
\begin{array}{lll} \text{Suivants}(E) &=& \{\#,\,) \;\} \cup \text{Suivants}(E') \\ \text{Suivants}(E') &=& \text{Suivants}(E) \\ \text{Suivants}(M) &=& \{+\} \cup \text{Suivants}(E) \cup \text{Suivants}(M') \\ \text{Suivants}(A') &=& \{*\} \cup \text{Suivants}(M) \end{array}
```

	E	E'	М	M′	<i>A</i>
0.	{ # }	Ø	Ø	Ø	Ø
1.	{ #,) }	{ # }	{ #, + }	Ø	{ * }
2.	{ #,) }	{ #,) }	{ #, +,) }	{ #, + }	{ #, +, * }
3.	{ #,) }	{ #,) }	{ #, +,) }	{ #, +,) }	{ #, +, *,) }
					{ #,+,*,) }

Analyse LL

Avec ces éléments, on peut construire une <u>table d'analyse LL(1)</u>. Pour chaque non terminal X (en cours d'analyse) et symbole terminal a (en tête de l'entrée), on cherche la(les) règle(s) à appliquer Le principe de construction est le suivant :

- Pour chaque règle $X \prec \beta$ et chaque symbole terminal $a \in \text{Premiers}(\beta)$ on indique la règle $X \prec \beta$ dans la case de la ligne X et de la colonne a.
- Pour chaque règle X ≺ β avec β →* ε et chaque symbole terminal a ∈ Suivants(X) on indique la règle X ≺ β dans la case de la ligne X et de la colonne a.

Si chaque case contient au plus une règle, on obtient une table d'analyse déterministe.

La grammaire utilisée est dite <u>LL(1)</u>

Exemple

	n	+	*	()	#
Ε	M E'			M E'		
E'		+ <i>E</i>			ε	ε
М	AM'			A M'		
M'		ε	* M		ε	ε
Α	n			(E)		

Conflits

- Si une case de la table contient plusieurs règles, on dit que la table contient un conflit, et que la grammaire n'est pas LL(1).
- Pour obtenir un analyseur dans ce cas, il faut revoir la grammaire ou utiliser une autre technique.
- être LL(1) est une propriété de la grammaire, et non du langage décrit par la grammaire.
- une grammaire LL(1) est non-ambigüe

Bilan sur l'analyse récursive descendante

- La technique d'analyse descendante est facile à programmer
- Contraintes fortes sur la grammaire pour choisir la règle à appliquer
 - fonctionne bien pour les instructions qui débutent par un mot clé
 - moins naturel pour les notations infixes (comme les opérateurs dans les expressions)
- D'autres techniques existent, avec des caractéristiques différentes.

Synthèse

A retenir

- Définition d'une grammaire, d'un arbre de dérivation syntaxique
- Notion d'ambiguïté
- Principales grammaires pour les expressions arithmétiques (grammaires naturelles ambigüe, non-ambigüe, grammaire LL(1)), grammaires pour des séquences d'instructions (séparées ou non par un caractère)
- Différence entre arbre de dérivation syntaxique et arbre de syntaxe abstraite
- Schéma général de l'analyse descendante (recherche du prochain lexème, fonctions mutuellement récursives de reconnaissance de non-terminaux)
- Analyseur descendant pour des expressions arithmétiques avec priorité

Savoir faire

Adapter l'analyseur descendant de IMP à de nouvelles constructions