

Compilation: Frontend

Copyright (C) 2018 - 2019
Flavien Breuvar
Licence Creative Commons Paternité
Partage à l'Identique 3.0 non transposé
(CC BY-SA 3.0)

8 avril 2020

1 Préliminaires

1.1 Récursion/induction et point fixe

Dans ce cours, on verra plusieurs algorithmes écrits en pseudo-code. Pour plus de clarté et d'uniformité, on utilisera deux structures particulières : l'induction et le point fixe.

Induction. L'induction est le pendant algorithmique de la récursion en mathématique où l'on "récurse" pas seulement sur des entiers, mais sur une structure de donnée finie (entier, mot, arbre). Pour ceux qui ont fait du Haskell ou du OCaml, il s'agit d'une fonction récursive commençant par un "pattern-matching" sur un argument. Une induction est donc une fonction (souvent auxiliaire) avec trois points :

- l'argument inductif : il s'agit d'un argument ou sous-argument (si l'argument est un tuple) qui va décroître à chaque appel,
- le ou les cas de base ($0, \epsilon, \text{feuille} \dots$),
- le ou les cas inductifs ($(n+1, t.\omega, \text{noeud} \dots)$) qui ont le droit de rappeler la fonction avec un argument inductif plus petit.

Le déroulement de l'algorithme est alors le suivant : on regarde l'argument inductif, si c'est l'un des cas de base, on exécute le code associé, et si c'est un cas inductif, on exécute le cas associé sauf que à chaque appel récursif on met en pause, on exécute l'algorithme sur ce cas plus petit, puis on revient à notre exécution et on utilise le résultat.

Exemple 1. *La taille d'un mot peut être calculée inductivement.*

Fonction $|\omega|$ calculée par induction :

- *argument inductif : mot ω ,*
- $|\epsilon| = 0$
- $|t\omega| = 1 + |\omega|$

Remarque 1. *En pseudo-code, on se passera de discussion d'optimisations telles que la récursion terminale.*

Point fixe. Le principe d'un algorithme par point fixe est le suivant : on se donne un ensemble fini E et on cherche le plus petit sous-ensemble $R \subseteq E$ qui contienne un certain $I \subseteq E$ et qui vérifie une condition croissante $R = R \cup f(R)$. C'est pourquoi il faut donner :

- un ensemble fini E ,
- une initialisation $I \subseteq E$, que l'on écrira $R := I$ avec dans l'idée que R est un mutable qui va être augmenté à chaque pas,
- un "pas" f , que l'on écrira $R := R \cup f(R)$ ou même $R := f(R)$ si f est croissant.

L'exécution associée est alors la suivante :

```
Set R := I;
Do {
  R' := union(R, f(R));
  fixe := isEmpty(R' - R);
  R := R'
} While (not fixe)
```

Exemple 2. Soit E un ensemble fini, $A \subseteq E$ et f une fonction de E dans les parties de E . La clôture de A par f (qu'on note $f^*(A)$) est l'ensemble des éléments de E accessibles depuis A par des applications successives de f . On la génère par point fixe :

- Ensemble : E
- Initialisation : $f^*(A) := A$,
- Pas : $f^*(A) := f^*(A) \cup \{x \in E \mid \exists a \in f^*(A), x \in f(a)\}$.

Remarque 2. La notion peut être généralisée à d'autres structures finies que les sous-ensembles d'un ensemble fini, cela marche pour tout ensemble fini ordonné. On peut même se passer de la condition "fini" en mettant des conditions plus faibles (depo, coinductifs), mais on doit alors procéder avec une exécution paresseuse et on n'a pas accès à certaines informations.

1.2 Tokens

On a vu, dans l'introduction et dans la description globale du compilateur, que les lettres de nos mots n'étaient pas des lettres, mais des "tokens", c'est à dire des lettres avec un contenu :

Définition 1.

Un alphabet de tokens Σ est la donnée d'un alphabet Σ et, pour chaque lettre $t \in \Sigma$, d'un ensemble $\text{tokens}(t)$.

Un token $\mathbf{t} \in \Sigma$ est un couple (t, T) pour $t \in \Sigma$ et $T \in \text{tokens}(t)$.

Notations 1.

On utilisera toujours version bleu et grasse \mathbf{t} pour désigner un token et la version noir t pour la lettre sous-jacente. De même, Σ est l'alphabet sous-jacent à l'alphabet de token Σ , et ω est le mot de lettres sous-jacent au mot de tokens ω .

Remarque 3. *La plupart du temps, on n'a pas besoin de distinguer le token de la lettre sous-jacente. En fait, le contenu du token ne sera utilisé que pour les actions des automates “shift-action”.*

Toutefois, il est impératif de ne jamais parler de token dans les règles de vos (variantes d') automates, car un alphabet de token peut contenir une infinité de token !

Dans ce cadre, on n'explicitera pas les token, même pour les études de complexité (ce devrait être nécessaire pour la complexité en espace, mais on ne pourrait pas dire grand chose sans spécifier le contenu des tokens).

Définition 2.

Un langage sur un alphabet Σ est simplement un ensemble de mots $L \subseteq \Sigma^$.*

Par contre on ne parlera jamais de langage de token : dans le frontend, on ne différencie jamais deux tokens par autre chose que le symbole associé. Par abus de langage, on utilisera des fois la notation $L \subseteq \Sigma^$, pour parler d'un langage $L \subseteq \Sigma^*$ et de l'ensemble $L := \{\omega \mid \omega \in L\}$.*

2 Grammaires

2.1 Définitions

Définition 3.

Une “grammaire hors contexte” est donnée par

- *Un ensemble \mathcal{N} de non-terminaux,*
- *Un ensemble (disjoint) \mathcal{T} de terminaux,*
- *Un non-terminal principal $S \in \mathcal{N}$,*
- *Pour chaque non-terminal $N \in \mathcal{N}$, un ensemble $\tau(N) \subseteq (\mathcal{T} \uplus \mathcal{N})^*$ de mots composés de terminaux et non-terminaux, appelés règles de N .*

Dans la suite, on appellera “grammaire” une grammaire hors contexte.

Notations 2 (Mathématique).

Pour poser la théorie et faire les preuves, on utilisera les notations suivantes :

- *les non-terminaux sont dénotés par les métavariabes majuscules S, N, \dots ,*
- *les terminaux sont dénotés par les métavariabes minuscules s, t, \dots qui peuvent correspondre à des tokens, dans les exemples et exercices, on utilisera aussi de symboles $(+, *, \lambda..)$ pour plus de lisibilité (attention le epsilon ϵ ne sont pas utilisés),*
- *les symboles qui peuvent être terminaux ou non-terminaux¹ sont dénotés par les métavariabes minuscules soulignées $\underline{a}, \underline{b}, \underline{c}, \dots$,*
- *le non-terminal principal est toujours S ,*
- *les règles sont notées $N \mapsto \omega$ pour $\omega \in \tau(N)$,*
- *le symbole ϵ est spécial : il désigne le mot vide comme résultat de la règle.*

1. Lors de la description des algorithmes on peut décrire des étapes qui se passent pareil sur des terminaux et non terminaux

Notations 3 (Informatique).

Pour décrire une grammaire formelle (réelle), on utilisera les notations suivantes :

- les non-terminaux sont dénotés par des noms entre crochets $\langle \text{expression} \rangle$, $\langle \text{toto} \rangle \dots$ commençant par des minuscules,
- les terminaux sont dénotés par des mots-clefs verbatim majuscule $\langle \text{VAR} \rangle$, $\langle \text{TOTO} \rangle \dots$, des caractères spéciaux verbatim $\mathbf{\text{£}}, (,), +, \dots$,
- lorsque $\tau(N) = \{\omega_1, \dots, \omega_n\}$, les règles sont notées $\langle \text{non-term} \rangle := \omega_1 \mid \dots \mid \omega_n$
- le non-terminal principal est celui dont les règles sont écrites en premières
- attention : le mot vide est ici représenté par une absence de mot et non pas par ϵ .

Notations 4 (Exercices). Pour faire des exercices, on se permettra de mélanger les deux types de notations tant qu'il n'y a pas d'ambiguïté.

Exemple 3. La grammaire du lambda calcul polymorphe (système F) est la suivante :

$$\begin{array}{ll}
 S \mapsto \mathbf{x} & T \mapsto \mathbf{t} \\
 S \mapsto \lambda \mathbf{x} : T. S & T \mapsto T \rightarrow T \\
 S \mapsto \Lambda \mathbf{t}. S & T \mapsto T \times T \\
 S \mapsto S S & T \mapsto (T) \\
 S \mapsto S T & \\
 S \mapsto (S) &
 \end{array}$$

où \mathbf{x} et \mathbf{t} représentent des terminaux encapsulant un token (les variables de termes pour le premier, et de type pour le second). Remarquez que seuls S et T sont des non terminaux, tous les symboles (λ , $.$, $:$, \rightarrow , times et même les parenthèses) sont des terminaux.

De façon plus pratique, on peut utiliser une formulation informatique :

$$\begin{array}{ll}
 \langle \text{terme} \rangle := \langle \text{VAR} \rangle & \langle \text{type} \rangle := \langle \text{VARTYPE} \rangle \\
 \mid \backslash \langle \text{VAR} \rangle : \langle \text{type} \rangle . \langle \text{terme} \rangle & \mid \langle \text{type} \rangle \langle \text{FLECHE} \rangle \langle \text{type} \rangle \\
 \mid \langle \text{LAMBDA} \rangle \langle \text{VARTYPE} \rangle . \langle \text{terme} \rangle & \mid \langle \text{type} \rangle * \langle \text{type} \rangle \\
 \mid \langle \text{terme} \rangle \langle \text{terme} \rangle & \mid (\langle \text{type} \rangle) \\
 \mid \langle \text{terme} \rangle \langle \text{type} \rangle & \\
 \mid (\langle \text{terme} \rangle) &
 \end{array}$$

Pour un exercice on se permettra de mélanger les notation pour aller au plus

lisible, par exemple :

$\langle \text{terme} \rangle := \mathbf{x}$ $\quad \lambda \mathbf{x} : \langle \text{type} \rangle . \langle \text{terme} \rangle$ $\quad \Lambda \alpha . \langle \text{terme} \rangle$ $\quad \langle \text{terme} \rangle \langle \text{terme} \rangle$ $\quad \langle \text{terme} \rangle \langle \text{type} \rangle$ $\quad (\langle \text{terme} \rangle)$	$\langle \text{type} \rangle := \alpha$ $\quad \langle \text{type} \rangle \rightarrow \langle \text{type} \rangle$ $\quad \langle \text{type} \rangle \times \langle \text{type} \rangle$ $\quad (\langle \text{type} \rangle)$
---	---

Dans tous les cas, la résolution est la même à renommage près, et vous pouvez utiliser une autre représentation sur votre copier si vous le souhaitez **à condition d'indiquer votre renommage en début d'exercice.**

Remarquez que pour chacun des cas, il faut préciser ce que sont les tokens \mathbf{x} (ou $\langle \text{VAR} \rangle$) et \mathbf{t} (ou $\langle \text{VARTYPE} \rangle$ ou α). Ici il s'agit de tokens contenant des chaînes de caractères.

Définition 4 (Règles partielles).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire. On appelle “règles partielles de \mathcal{G} depuis un non-terminal $N \in \mathcal{N}$ ”, l'ensemble

$$\tau_{\bullet}(N) = \{(\omega_1, \omega_2) \mid \omega_1 \omega_2 \in \tau(N)\} \subseteq (\mathcal{T} \uplus \mathcal{N})^* \times (\mathcal{T} \uplus \mathcal{N})^*.$$

On dénote $N \mapsto \omega_1 \bullet \omega_2$ ces règles partielles $(\omega_1, \omega_2) \in \tau_{\bullet}(N)$.

On qualifie d’“initiale” une règle partielle de la forme $N \mapsto \bullet \omega$ et de “finale” une règle partielle de la forme $N \mapsto \omega \bullet$.

Définition 5 (Multistep).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

On appelle “réduction multistep” les couples $\omega \mapsto^* \omega'$ pour $\omega, \omega' \in (\mathcal{N} \uplus \mathcal{T})^*$ définis comme la clôture transitive réflexive de \mapsto :

- Ensemble : $(\mapsto^*) \subseteq ((\mathcal{N} \uplus \mathcal{T})^*)^2$,
- Initialisation : $\omega \mapsto^* \omega$ pour tout $\omega \in \mathcal{N} \uplus \mathcal{T}$, et $N \mapsto^* \omega$ lorsque $N \mapsto \omega$,
- Pas : $\omega_1 \omega_2 \mapsto^* \omega'_1 \omega'_2$ lorsque $\omega_1 \mapsto \omega'_1$ et $\omega_2 \mapsto \omega'_2$, de plus $\omega \mapsto^* \omega''$ lorsqu'il existe $\omega \mapsto^* \omega' \mapsto^* \omega''$.

Remarque 4. Ceci n'est pas un algorithme car l'ensemble $((\mathcal{N} \uplus \mathcal{T})^*)^2$ est infini. Ce point-fixe montre par contre que la relation (\mapsto^*) est récursivement énumérable, c'est à dire que si $\omega \mapsto^* \omega'$, on le trouvera en un temps fini.

Définition 6 (Langage d'une grammaire).

Le langage $L_{\mathcal{G}}$ d'une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ est l'ensemble des mots de terminaux accessibles depuis le non-terminal initial :

$$L_{\mathcal{G}} := \{\omega \in \mathcal{T}^* \mid S \mapsto^* \omega\}.$$

On dit qu'un mot $\omega \in \mathcal{T}^*$ est reconnu par \mathcal{G} si $\omega \in L_{\mathcal{G}}$.

2.2 Arbres syntactiques

Définition 7 (Arbre syntactique version mathématique).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

On définit les arbres syntactiques des non-terminaux de \mathcal{G} comme un ensemble \mathcal{N} de tokens sur \mathcal{N} défini par point fixe :

- Ensemble : $\mathcal{N} \subseteq \mathcal{N} \times (\mathcal{T} \uplus \mathcal{N})^*$,
- Initialisation : $\mathcal{N} = \emptyset$,
- Pas :

$$\mathcal{N} := \bigcup_{N \mapsto \underline{c_1} \dots \underline{c_n}} \{(N, \underline{c_1} \dots \underline{c_n}) \mid \forall i, \underline{c_i} \in \mathcal{T} \Rightarrow \underline{c_i} \in \mathcal{T} \text{ et } \underline{c_i} \in \mathcal{N} \Rightarrow \underline{c_i} \in \mathcal{N}\}$$

On appelle arbre syntactique de \mathcal{G} un arbre syntactique \mathbf{S} de S .

Notations 5. Un arbre syntactique $\mathbf{N} = (N, \underline{c_1} \dots \underline{c_n})$ peut se représenter sous la forme d'un arbre de la racine labellé par N et avec pour fils les arbres syntaxiques $\underline{c_1} \dots \underline{c_n}$.

Définition 8 (Arbre syntactique version informatique).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

On définit les arbres syntactiques `NonTerm` partant de `<non-term> ∈ \mathcal{N}`
`<non-term> := <a_1_1>...<a_1_k1> | ... | <a_n_1>...<a_n_kn>`
différemment selon le langage :

- comme un type-union (ex : en C) qui est une struct obtenue comme l'union des structs `NonTerm1...NonTermN` suivantes :

```
NonTermI := {a_i_1 : AI1, ... , a_i_ki : AIKI }
```

- comme le type algébrique (ex : OCaml, Haskell) :

```
NonTerm := Regle1 of (A11, ..., A1K1)
          | ...
          | RegleN of (AN1, ..., ANKN)
```

- comme une classe abstraite (ex : Java) hérité par n versions concrétisés :

```
public class NonTermI extends NonTerm {
    AI1 a_i_1;
    ...
    AIKI a_i_ki;
}
```

Remarque 5. De part la multiplicité des implémentations, on utilisera toujours la version mathématique. De plus, la version mathématique est, par définition, restreinte aux arbres finis, ce qui nous est nécessaire pour le théorème suivant :

Proposition 1 (Induction sur les arbres syntactique).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

Il y a un ordre bien fondé sur les tokens de terminaux et les non-terminaux tel que pour tout $\mathbf{N} = (N, \underline{c_1} \dots \underline{c_n}) \in \mathcal{N}$ et tout $i \leq n$, $\underline{c_i} < \mathbf{N}$.

Autrement dit, on peut faire une induction sur les arbres syntactiques.

Définition 9 (Mot reconnu par un arbre syntactique).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire et $N \in \mathcal{N}$ un arbre syntactique.

Le mot $\text{mot}(\mathbf{c}) \in \mathcal{T}^*$ reconnu par $\underline{\mathbf{c}} \in \mathcal{N} \uplus \mathcal{T}$ est défini par induction :

- argument inductif : $\underline{\mathbf{c}} \in \mathcal{N} \uplus \mathcal{T}$,
- cas $\mathbf{t} \in \mathcal{T}$: alors $\text{mot}(\mathbf{t}) = \mathbf{t} \in \mathcal{T}^*$,
- cas $N = (N, \mathbf{c}_1 \dots \mathbf{c}_n) \in \mathcal{N}$: alors $\text{mot}(N) = \text{mot}(\mathbf{c}_1) \dots \text{mot}(\mathbf{c}_n)$.

Proposition 2.

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire et $N \in \mathcal{N}$ un non-terminal.

L'ensemble des mots reconnus par des arbres syntaxiques issus de N est exactement l'ensemble des réduits terminaux de N :

$$N \mapsto^* \omega \quad \Leftrightarrow \quad \exists N, \omega = \text{mot}(N)$$

Définition 10 (Ambigüité).

Une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ est ambiguë si un non terminal a plusieurs arbres syntactiques reconnaissant un même mot, c'est à dire si :

$$\exists N \in \mathcal{N}, \exists (N, \mathbf{u}) \neq (N, \mathbf{v}) \in \mathcal{N}, \quad \text{mot}(N, \mathbf{u}) = \text{mot}(N, \mathbf{v}) .$$

On peut associer des règles d'associativité et de priorité sur une grammaire ambiguë pour la désambigüiser.

Définition 11 (Paire critique).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$.

Une paire critique est la donnée de deux règles $(N \mapsto u)$, $(N \mapsto v)$ sur un même non-terminal N , et de deux décompositions différentes $u = u_1 N u_2$ et $v = v_1 N v_2$ telles que $u_1 v_1 = v_1 u_1$ et $u_2 v_2 = v_2 u_2$.

On note $(N \mapsto \underline{N} u_2 \mid v_1 \underline{N} v_2)$ une telle paire critique.

Proposition 3. Une grammaire avec une paire critique $(N \mapsto u_1 \underline{N} u_2, N \mapsto v_1 \underline{N} v_2)$ dont les règles sont productives est nécessairement ambiguë.

Remarque 6. Ce n'est pas une équivalence, càd. qu'il existe des grammaires ambiguës sans paires critiques. Mais il s'agit de la majeure partie des cas.

Définition 12 (Priorité).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

Une priorité est un ordre $(N \mapsto u) > (N \mapsto v)$ entre deux règles différentes sur un même non-terminal.

Une telle priorité va casser toutes les paires critiques $(N \mapsto u_1 \underline{N} u_2 \mid v_1 \underline{N} v_2)$ entre ces deux règles (càd. $u = u_1 N u_2$ et $v = v_1 N v_2$) en interdisant les arbres syntaxiques contenant un sous-arbre de la forme $(N, \mathbf{u}_1 (N, \mathbf{v}_1 \underline{N} \mathbf{v}_2) \mathbf{u}_2)$ au profit de l'arbre équivalent de la forme $(N, \mathbf{v}'_1 (N, \mathbf{u}'_1 \underline{N} \mathbf{u}'_2) \mathbf{v}'_2)$ où $\mathbf{v}'_1 \mathbf{u}'_1 := \mathbf{u}_1 \mathbf{v}_1$ et $\mathbf{u}'_2 \mathbf{v}'_2 := \mathbf{v}_2 \mathbf{u}_2$.

Remarque 7. Étant donné deux règles $(N \mapsto u)$, $(N \mapsto v)$, il peut y avoir plusieurs décompositions $(u, v) = (u_1 N u_2, v_1 N v_2)$ avec $u_1 v_1 = v_1 u_1$ et $u_2 v_2 = v_2 u_2$. Dans ce cas, dire que la première règle est prioritaire sur la seconde impose la priorité pour toutes les décompositions.

Définition 13 (Associativité).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

Une associativité à gauche sur un opérateur $N \mapsto NuN$ va casser toutes les paires critiques ($N \mapsto u_1 \underline{NuN} \mid Nu \underline{N}$) de la règle avec elle même en interdisant les arbres syntaxiques contenant un sous-arbre de la forme $(N, \mathbf{Nu}(N, \mathbf{N}'u'N''))$ au profit de l'arbre équivalent $(N, (N, \mathbf{NuN}')u'N'')$.

Une associativité droite fera de même mais en privilégiant le premier sur le second.

Remarque 8. Cette définition peut se généraliser à toute règle $N \mapsto w$, mais dans la pratique ce n'est jamais utile.

En effet, une paire critique entre w et lui même est forcément de la forme $((Nu)^n \underline{N}(uN)^k \mid (Nu)^{n'} \underline{N}(uN)^{k'})$ pour $n + k = n' + k'$ mais $n \neq n'$. On interdit alors tous les arbres syntaxiques contenant un sous-arbre de la forme $(N, (\mathbf{Nu})^n(N, (\mathbf{Nu})^{n'} \mathbf{N}(\mathbf{Nu})^{k'})(uN)^k)$ au profit de l'arbre équivalent de la forme $(N, (\mathbf{Nu})^{n'}(N, (\mathbf{Nu})^n \mathbf{N}(\mathbf{Nu})^k)(uN)^{k'})$ lorsque $n' < n$.

2.3 AST

Un AST, pour *Abstract Syntactic Tree*, est ce qui va être rendu en sortie de votre parseur pour être traité par la suite de votre compilateur.

C'est une abstraction de l'arbre syntaxique, cela signifie :

- que l'on peut y supprimer les informations inutiles pour la suite de la compilation, par exemple, les parenthèses ne sont utiles que pour l'analyse syntaxique (parseur) et peuvent donc être enlevées de l'AST,
- que l'on peut modifier la structure de l'arbre pour refléter une grammaire plus aisée, à utiliser, par exemple, les différentes opérations peuvent être remplacées par des noeuds spécifiques,
- que l'on peut ajouter des informations utile pour la suite (ex : dans le projet, on va avoir besoin d'y ajouter la taille attendu du code généré, on peut aussi y ajouter les types).

2.4 Firsts et Follows

Définition 14 (union-epsilon).

Soit $e, f \subseteq \mathcal{T} \uplus \{\epsilon\}$ et , on appelle union-epsilon l'opérateur :

$$e \uplus f := \begin{cases} (e - \{\epsilon\}) \cup f & \text{si } \epsilon \in e \\ e & \text{sinon} \end{cases}$$

correspondant aux élément non- ϵ de e et aux élément de f uniquement ajoutés si $\epsilon \in e$.

Algorithme 1 (firsts).

Donnée : Une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$,

Résultat : pour chaque mot $\omega \in (\mathcal{T} \uplus \mathcal{N})^*$ un ensemble $\mathbf{first}(\omega) \subseteq \mathcal{T} \uplus \{\epsilon\}$,

Spécification 1 : $\epsilon \in \mathbf{first}(\omega)$ ssi $\omega \mapsto^* \epsilon$,

Spécification 2 : $t \in \mathbf{first}(\omega)$ ssi il existe ω' tel que $\omega \mapsto^* t\omega'$

Point fixe :

- *Ensemble* : $\mathbf{First} \subseteq \mathcal{N} \times (\mathcal{T} \uplus \{\epsilon\})$
- *Initialisation* : $\mathbf{First} := \emptyset$
- *Pas* : $\mathbf{First} := \{(N, t) \mid \exists(N \mapsto \omega), t \in \mathbf{first}(\omega)\}$

Macros :

$\mathbf{first}(\epsilon) := \emptyset$
 $\mathbf{first}(tu) := \{t\}$
 $\mathbf{first}(Nu) := \{t \mid (N, t) \in \mathbf{First}\} \uplus \mathbf{first}(u)$

Algorithme 2 (follows).

Donnée : Une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$,

Résultat : pour chaque non-terminal $N \in \mathcal{N}$ un ensemble $\mathbf{follow}(N) \subseteq \mathcal{T} \uplus \{\$\}$,

Spécification 1 : $\$ \in \mathbf{follow}(N)$ ssi il existe une réduction $S \mapsto^* \omega N$,

Spécification 2 : $t \in \mathbf{follow}(N)$ ssi il existe une réduction $M \mapsto^* \omega_1 N t \omega_2$

Point fixe :

- *Ensemble* : $\mathbf{Follow} \subseteq \mathcal{N} \times (\mathcal{T} \uplus \{\$\})$
- *Initisation* : $\mathbf{Follow} := \{(S, \$)\}$
- *Pas* : $\mathbf{Follow} := \{(N, t) \mid \exists(M \mapsto \omega_1 N \omega_2), t \in \mathbf{first}(\omega_2) \uplus \mathbf{follow}(M)\}$

Macros :

$\mathbf{follow}(N) := \{t \mid (N, t) \in \mathbf{Follow}\}$

3 Automates : rappels et préliminaires

3.1 Automates Finis

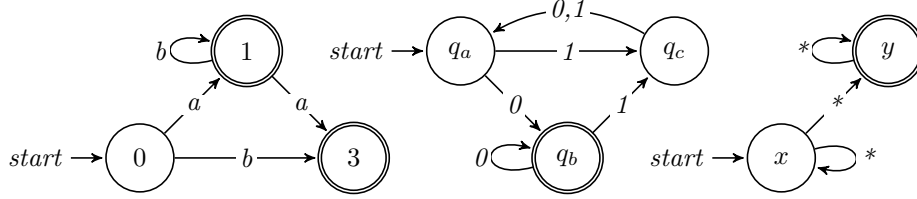
Définition 15 (Automate fini).

Un automate fini \mathcal{A} est la donnée :

- d'un alphabet (de tokens) Σ ,
- d'un ensemble d'états Q ,
- d'un état $q_0 \in Q$ initial,
- d'un ensemble d'états $F \subseteq Q$ finaux,
- d'une transition $\mathbf{trans}(q, t) \in Q$ pour chaque état $q \in Q$ et caractère $t \in \Sigma$.

Notations 6. On représente un automate fini $\mathcal{A} = (\Sigma, Q, q_0, F, \mathbf{trans})$ comme un graphe dirigé étiqueté dont les noeuds sont les états $q \in Q$ et les arêtes sont les $p \xrightarrow{t} q$ tels que $q = \mathbf{trans}(p, t)$, l'état initial est indiqué par une flèche entrante sans origine, et les états finaux par un double cercle (ou double carré lorsque l'on écrit des états carrés) autour de l'état. Attention, \mathbf{trans} étant une fonction, il doit y avoir exactement une arête $p \xrightarrow{t} q$ pour chaque état p et chaque lettre t . Pour la lisibilité, on peut fusionner deux arêtes en une seule en plaçant les deux labels côte à côte, séparés par une virgule.

Exercice 1. Des graphes ci-dessous, lesquels représentent des automates finis ? Donner leurs définitions formelles



Définition 16 (reconnaissance sur un automate).

Un mot $\omega \in \Sigma^*$ sur l'alphabet Σ est reconnu par un automate fini $\mathcal{A} = (\Sigma, Q, q_0, F, \text{trans})$ si

- ou bien $\omega = \epsilon$ est le mot vide et $q_0 \in F$,
- ou bien $\omega = t.\omega'$ et ω' est reconnu par $(\Sigma, Q, \text{trans}(q_0, t), F, \text{trans})$.

L'ensemble des mots reconnus par \mathcal{A} est son langage, désigné par $L_{\mathcal{A}}$

Proposition 4 (reconnaissance).

La complexité en temps de cet algorithme de reconnaissance est

$$O(|\omega| * \log(|Q| * |\Sigma|)) \stackrel{\Sigma}{\simeq} O(|\omega| * \log(|Q|)).$$

Sa complexité en espace est (en considérant que les tokens sont vides) :

$$O(|\omega| * \log(|\Sigma|) + (|Q| * |\Sigma|) * \log(Q)) \stackrel{\Sigma}{\simeq} O(|\omega| + |Q| \log(|Q|))$$

Démonstration. En temps :

À chaque lettre du mot, on teste si le mot est vide (constant), on récupère la première lettre (constant), on appelle la fonction $\text{trans}(q_0, t)$ qui est en $\log(|Q| * |\Sigma|)$ dans le cas général. Le fait que la fonction $\text{trans}(q_0, t)$ qui est en $\log(|Q| * |\Sigma|)$ dans le cas général vient du cas où trans est un dictionnaire (associant les clés dans $Q * \Sigma$ à des résultats dans Q) codé avec un arbre équilibré (AVL, ARN, A23, R-tree, etc...), la complexité peut être prouvée optimale. En espace :

Il faut stocker le mot (en $O(|\omega| \cdot \log(|\Sigma|))$ si les tokens sont vides, car une lettre prend une taille $\log(|\Sigma|)$), un pointeur sur le mot (en $O(\log(|\omega|))$) et l'automate lui même. Pour l'automate, on doit stocker (dans le pire cas) le nom d'un état (en $O(\log(|Q|))$) pour chaque entrée état/lettre (en $O(|Q| * |\Sigma|)$). \square

Remarque 9. On arrive ici aux limites entre la théorie et la pratique de la complexité : bien que les facteurs logarithmiques soit toujours là, on les ignore généralement en pratique, et à raison ! En effet, en utilisant une structure de tableau, on perd en flexibilité, mais on a un accès qui semble être en temps constant. Je dis bien "semble" car il n'est en temps constant que tant qu'il y a de la place sur votre RAM... ce qui en pratique est toujours vérifié (à moins d'avoir fait exploser votre complexité en espace...).

Exercice 1 (Automates partiels²). Un automate fini partiel est défini de la même façon, sauf que q_0 , ainsi que $\text{trans}(q, t)$ ne sont pas toujours définis (renvoient \perp), lors de la reconnaissance, on renvoi faux lorsque l'on essaye de calculer un état non défini.

Montrez que l'on peut réduire tout automate partiel \mathcal{A} en un automate fini \mathcal{A} de même langage. L'augmentation du nombre d'états doit être constante.

Notations 7. *Puisqu'ils sont si proche, on dessinera généralement un automate fini partiel pour parler de l'automate fini. En effet, l'état puits/poubelle a plein de transition qui lui arrive, ce qui complexifie le dessin pour rien.*

Les automates finis ont 5 propriétés remarquables qui justifient leur étude et qui expliquent pourquoi on essaye de les étendre :

- la complexité linéaire (en pratique) de la reconnaissance que l'on vient de voir,
- l'équivalence avec les langages rationnels, que l'on va voir tout de suite,
- la minimisation, que l'on va voir la sous section suivante,
- la détermination, que l'on va voir la sous section encore après et qui sera le centre de la partie,
- les propriétés de vérifications (langage vide, etc...) que l'on ne verra pas.

3.2 Expressions régulières

3.2.1 Théorie

Définition 17. *Les expressions régulières $\text{Reg}(\Sigma)$ sur l'alphabet Σ sont données par la grammaire suivante :*

- $\mathcal{T} = \Sigma \cup \{(\cdot), *, +, \epsilon\}$,
- $\mathcal{N} = \{S\}$,
- les règles sont :

$$\begin{array}{lll} S \mapsto t & S \mapsto \epsilon & S \mapsto S + S \\ S \mapsto SS & S \mapsto S^* & S \mapsto (S) \end{array}$$

Ou de façon équivalente (ici \mathfrak{t} représente tous les symboles de $t \in \Sigma$) :

$$\begin{array}{lll} \langle \text{regexpr} \rangle := \langle \text{CHAR} \rangle & | \langle \text{EPS} \rangle & | \langle \text{regexpr} \rangle + \langle \text{regexpr} \rangle \\ & | \langle \text{regexpr} \rangle^* & | \langle \text{regexpr} \rangle \langle \text{regexpr} \rangle & | (\langle \text{regexpr} \rangle) \end{array}$$

De plus, la somme $S + S$ et la composition SS sont associatif à droite, avec l'étoile S^* prioritaire sur la composition elle même prioritaire sur la somme.

Remarque : Ci-dessus, on utilise le symbole $\epsilon \neq \epsilon$ pour bien indiquer que la règle ne crée pas le mot vide ϵ , mais crée un symbole ϵ qui est celui de l'expression régulière vide.³

2. Cet exercice est pour ceux qui ne maîtrisent pas les automates, ce ne sera pas demander au partiel.

3. Si vous pensez que l'on mixe les niveaux d'abstraction, attendez la suite.

Proposition 5. *La grammaire des expression régulière est non ambiguë modulo les règles de priorités et d'associativité, $\text{Reg}(\Sigma)$ est donc bien défini.*

Notations 8. *On appelle expression régulière à la fois les arbres syntactiques de la grammaire et les mots dérivés. En effet, la grammaire n'étant pas ambiguë il y a une correspondance exacte.*

Définition 18. *Soit r une expression régulière sur l'alphabet Σ .*

Le langage L_r reconnu par r est défini par induction sur r :

- $L_t = \{t\}$,
- $L_\varepsilon = \Sigma$,
- $L_{r_1+r_2} = L_{r_1} \cup L_{r_2}$,
- $L_{r_1r_2} = \{\omega_1\omega_2 \mid \omega_1 \in L_{r_1}, \omega_2 \in L_{r_2}\}$,
- $L_{s^*} = \{\omega_1\dots\omega_n \mid n \geq 0, \forall i \leq n, \omega_i \in L_s\}$.

Remarque : La somme et la composition sont commutatives au sens où choisir une associativité à gauche n'aurait pas changé le langage reconnu.

Définition 19. *Soit r une expression régulière sur l'alphabet Σ .*

L'automate \mathcal{A}_r associée à r est défini comme la déterminisation et minimisation de l'automate $\mathcal{A}_r^0 = (\Sigma, Q_r, q_{r,0}, F_r, \text{eps}_r, \text{trans}_r)$ défini par induction sur r :

– cas t :

$$Q_t = \{q_{t,0}, 1\}, \quad F_t = \{1\}, \quad \text{eps}_t(q, s) = \emptyset,$$

$$\text{trans}_t(q, s) \begin{cases} \{1\} & \text{si } q = q_{t,0} \text{ et } s = t \\ \emptyset & \text{sinon} \end{cases}$$

– cas ε :

$$Q_\varepsilon = \{q_{\varepsilon,0}, 1\}, \quad F_\varepsilon = \{1\}, \quad \text{epst.}(q, s) = \emptyset,$$

$$\text{trans}_\varepsilon(q, s) \begin{cases} \{1\} & \text{si } q = q_{\varepsilon,0} \\ \emptyset & \text{sinon} \end{cases}$$

– cas $r_1 + r_2$:

$$Q_{r_1+r_2} = Q_{r_1} \uplus Q_{r_2} \uplus \{q_{r_1+r_2,0}\}, \quad F_{r_1+r_2} = F_{r_1} \uplus F_{r_2}$$

$$\text{eps}_{r_1+r_2}(q) \begin{cases} \{q_{r_1}, q_{r_2}\} & \text{si } q = q_{r_1+r_2,0} \\ \text{eps}_{r_1}(q) & \text{si } q \in Q_{r_1} \\ \text{eps}_{r_2}(q) & \text{si } q \in Q_{r_2} \end{cases} \quad \text{trans}_{r_1+r_2}(q) \begin{cases} \emptyset & \text{si } q = q_{r_1+r_2,0} \\ \text{trans}_{r_1}(q) & \text{si } q \in Q_{r_1} \\ \text{trans}_{r_2}(q) & \text{si } q \in Q_{r_2} \end{cases}$$

– cas r_1r_2 :

$$Q_{r_1r_2} = Q_{r_1} \uplus Q_{r_2}, \quad q_{r_1r_2,0} = q_{r_1,0}, \quad F_{r_1r_2} = F_{r_2}$$

$$\text{eps}_{r_1r_2}(q) \begin{cases} \text{eps}_{r_1}(q) & \text{si } q \in Q_{r_1} - F_{r_1} \\ \text{eps}_{r_1}(q) \cup \{q_{r_2,0}\} & \text{si } q \in F_{r_1} \\ \text{eps}_{r_2}(q) & \text{si } q \in Q_{r_2} \end{cases} \quad \text{trans}_{r_1r_2}(q) \begin{cases} \text{trans}_{r_1}(q) & \text{si } q \in Q_{r_1} \\ \text{trans}_{r_2}(q) & \text{si } q \in Q_{r_2} \end{cases}$$

— cas r^* :

$$Q_{r^*} = Q_r \uplus \{q_{r^*,0}\}, \quad F_{r^*} = F_r$$

$$\text{eps}_{r^*}(q) \begin{cases} \{q_{r,0}\} & \text{si } q = q_{r^*,0} \\ \text{eps}_r(q) & \text{si } q \in Q_{r_1} - F_{r_1} \\ \text{eps}_r(q) \cup \{q_{r,0}\} & \text{si } q \in F_r \end{cases} \quad \text{trans}_{r^*}(q) \begin{cases} \emptyset & \text{si } q = q_{r^*,0} \\ \text{trans}_r(q) & \text{sinon} \end{cases}$$

Proposition 6. *Le langage reconnu par l'automate associé à une expression régulière est exactement le langage reconnu par celle-ci :*

$$L_{\mathcal{A}_{r^*}} = L_r.$$

Remarquez que l'on peut aussi aller dans l'autre sens :

Proposition 7. *Pour tout automate, il existe une expression régulière qui reconnaît le même langage*

Cela veut dire que le pouvoir expressif des automates et celui des expressions régulières sont les mêmes, on parle de langages réguliers. Mais ce second sens n'a pas vraiment d'intérêt pratique.⁴

3.2.2 Pratique

Dans la pratique, grammaires et expressions régulières sont intimement liées et définies l'une à partir de l'autre... Nous allons essayer de démêler tout ça.

D'abord, notons qu'en pratique, on utilise principalement les expressions régulières sur des caractères, et que l'on ne veut pas s'amuser à écrire 'a'+ 'b'+ 'c'+ ... + 'z' lorsque l'on veut représenter toutes les lettres minuscule, on utilise plutôt 'a' - 'z' qui désigne tous les caractères entre 'a' à 'z'.

Définition 20 (Ensembles de caractères). *Le token <CHAR> sera utilisé pour désigner un caractère (typiquement ASCII, mais on pourrait imaginer d'autres encodages) entre deux quotes ' et avec échappements.*

On utilise alors la grammaire suivante pour décrire les ensembles de caractères :

$$\langle \text{chars} \rangle := \langle \text{CHAR} \rangle \mid \langle \text{CHAR} \rangle - \langle \text{CHAR} \rangle \mid \langle \text{CHAR} \rangle, \langle \text{chars} \rangle \mid \langle \text{CHAR} \rangle - \langle \text{CHAR} \rangle, \langle \text{chars} \rangle$$

Remarques : D'abord, remarquez que l'on ne décrit pas bien ce qu'est un "caractère", normalement, il faudrait l'expression régulière définissant les caractères, qui elle-même nécessite de parler d'ensemble de caractères... si on ne veut pas boucler, il faudrait être plus formel, mais ce serait trop...

Aussi utilisé pour les automates : Les ensembles de caractères sont aussi utilisés pour les transitions des automates sur des caractères.

⁴ En tout cas pas en compilation...

Une grammaire plus riche : En pratique, les outils utilisant des expressions régulières utilisent chacun leur grammaire d'expressions régulières, la très grande majorité de ces notations se traduisent avec les opérateurs que nous avons vus, mais certaines sont des “expressions régulières étendues” et ne peuvent s'exprimer que en utilisant une grammaire.

3.3 Minimisation

Définition 21 (Accessibilité et co-accessibilité).

Soit $\mathcal{A} = (\Sigma, Q, q_0, F, \text{trans})$ automate fini (partiel).

Un état $q \in Q$ est dit accessible s'il existe un mot $t_1 \dots t_n \in \Sigma^n$ une séquence $q_0, q_1, \dots, q_n = q$ tels que $q_{i+1} = \text{trans}(q_i, t_{i+1})$, c'est à dire que q peut être accédé depuis l'état initial q_0 .

Un état $q \in Q$ est dit co-accessible s'il existe un mot $t_1 \dots t_n \in \Sigma^n$ une séquence $q = p_0, p_1, \dots, p_n \in F$ tels que $p_{i+1} = \text{trans}(p_i, t_{i+1})$, c'est à dire qu'un état final $p_n \in F$ peut être accédé depuis q .

Algorithme 3 (Co-accessibilité).

Donnée : Un automate fini $\mathcal{A} = (\Sigma, Q, q_0, F, \text{trans})$,

Résultat : un automate fini partiel $\text{cac}(\mathcal{A}) = (\Sigma, \text{cac}(Q), q'_0, F, \text{trans})$,

Spécification : $L_{\mathcal{A}} = L_{\text{cac}(\mathcal{A})}$ et $\text{cac}(Q) \subseteq Q$ et tous les états de $\text{cac}(\mathcal{A})$ co-accessibles.

Point fixe :

— Ensemble : $\text{cac}(Q) \subseteq Q$

— Initialisation : $\text{cac}(Q) := F$,

— Pas : $\text{cac}(Q) := \text{cac}(Q) \cup \{q \mid \exists t \in \Sigma, \text{trans}(q, t) \in \text{cac}(Q)\}$.

À la fin, on doit vérifier que $q_0 \in \text{cac}(Q)$, si c'est le cas $q'_0 = q_0$ sinon $q'_0 = \perp$ est non définit.

Algorithme 4 (Accessibilité).

Donnée : Un automate fini partiel $\mathcal{A} = (\Sigma, Q, q_0, F, \text{trans})$,

Résultat : un automate fini partiel $\text{acc}(\mathcal{A}) = (\Sigma, \text{acc}(Q), q_0, F', \text{trans})$

Spécification : $L_{\mathcal{A}} = L_{\text{acc}(\mathcal{A})}$ et $\text{acc}(Q) \subseteq Q$ et tous les états de $\text{acc}(\mathcal{A})$ accessibles.

Point fixe :

— Ensemble : $\text{acc}(Q) \subseteq Q$

— Initialisation : $\text{acc}(Q) := \{q_0\}$ si $q_0 \neq \perp$ et \emptyset sinon,

— Pas : $\text{acc}(Q) := \text{acc}(Q) \cup \{\text{trans}(q, t) \mid q \in \text{acc}(Q), t \in \Sigma\}$.

À la fin on pose $F' = F \cap \text{acc}(Q)$.

Algorithme 5 (Minimisation).

Donnée : Un automate fini $\mathcal{A} = (\Sigma, Q, q_0, F, \text{trans})$,

Résultat : un automate fini $\text{min}(\mathcal{A}) = (\Sigma, \text{min}(Q), q_0^m, F^m, \text{trans}^m)$ ou une erreur,

Spécification : $L_{\mathcal{A}} = L_{\text{min}(\mathcal{A})}$ et $\text{min}(\mathcal{A})$ est minimal parmi les automates ayant cette propriété.

- Étape 1 : calculer $\text{cac}(\mathcal{A})$ avec l'Algorithme 3,
 Étape 2 : calculer $\text{acc}(\text{cac}(\mathcal{A}))$ avec l'Algorithme 4,
 Étape 3 : rajouter un état puit pour retrouver un automate total $\mathcal{A}' = (\Sigma, Q', q'_0, F', \text{trans})$
 Étape 4 : on cherche une relation d'équivalence (\equiv) sur Q' appelée bisimulation,
 pour ça on procède par point fixe :
- Ensemble : $(\equiv) \subseteq Q'^2$,
 - Initialisation : $(\equiv) := (=)$ c'est à dire que deux états sont équivalents ssi ils sont égaux,
 - Pas : $p \equiv q$ lorsque :
 - $(p \in F \Leftrightarrow q \in F)$,
 - et pour tout $t \in \Sigma$, $\text{trans}(p, t) \equiv \text{trans}(q, t)$.
- Étape 5 : On pose alors
- $\min(Q) := Q_{/\equiv}$ est l'ensemble des classes d'équivalences de \equiv ,
 - $q_0^m := [q_0]$ est la classe d'équivalence de q_0 ,
 - $F^m := F_{/\equiv}$ est l'ensemble des classes d'équivalences de \equiv sur F ,
 - $\text{trans}^m([q], t) := [\text{trans}(q, t)]$

3.4 Automates non déterministes et déterminisation

Définition 22 (Automate fini non déterministe).

Un automate fini non déterministe \mathcal{A} est la donnée :

- d'un alphabet (de tokens) Σ ,
- d'un ensemble d'états Q ,
- d'un état $q_0 \in Q$ initial,
- d'un ensemble d'état $F \subseteq Q$ finaux,
- d'un ensemble d' ϵ -transitions $\text{eps}(q) \subseteq Q$ pour chaque état $q \in Q$,
- d'un ensemble de transitions $\text{trans}(q, t) \subseteq Q$ pour chaque état $q \in Q$ et caractère $t \in \Sigma$.

Notations 9. On représente un automate fini non déterministe $\mathcal{A} = (\Sigma, Q, q_0, F, \text{eps}, \text{trans})$ comme un graphe dirigé étiqueté dont les noeuds sont les états $q \in Q$ et les arêtes sont les $p \xrightarrow{t} q$ tels que $q \in \text{trans}(p, t)$ et les $p \xrightarrow{\epsilon} q$ tels que $q \in \text{eps}(p)$.

Définition 23 (reconnaissance sur un automate non déterministe).

Un mot $\omega \in \Sigma^*$ sur l'alphabet Σ est reconnu par un automate fini non déterministe $\mathcal{A} = (\Sigma, Q, q_0, F, \text{eps}, \text{trans})$ si

- ou bien $\omega = \epsilon$ est le mot vide et $q_0 \in F$,
- ou bien ω est reconnu par $(\Sigma, Q, q', F, \text{trans})$ pour un état $q' \in \text{eps}(q_0)$,
- ou bien $\omega = t.\omega'$ et ω' est reconnu par $(\Sigma, Q, q', F, \text{trans})$ pour un état $q' \in \text{trans}(q_0, t)$.

L'ensemble des mots reconnus par \mathcal{A} est son langage, désigné par $L_{\mathcal{A}}$

Proposition 8 (reconnaissance non déterministe).

Cet algorithme naïf de reconnaissance est en

$$|Q|^{O(|\omega|*|Q|)} * \log(|\Sigma|)^{O(|\omega|)} \stackrel{\Sigma}{\simeq} |Q|^{O(|\omega|*|Q|)}.$$

à retenir : plus que exponentiel en $|Q|$ et en $|\omega|$

Démonstration. Pour éviter de boucler, on remarque que l'on ne peut pas prendre une ϵ -transition plus de $|Q|$ fois de suite.

Soit $f_k(\omega)$ le temps de reconnaissance du mot ω depuis un état quelconque alors que l'on vient de prendre k ϵ -transitions de suite.

On a alors

$$\begin{aligned} f_{|Q|}(\epsilon) &= O(\log(|Q|)) \\ f_k(\epsilon) &= O(\log(|Q|) + \log(|Q|) * |Q| * f_{k+1}(\epsilon)) \\ f_{|Q|}(t\omega) &= O(|Q| * \log(Q * |\Sigma|) * f_0(\omega)) \\ f_k(t\omega) &= O(|Q| * \log(Q * |\Sigma|) * f_0(\omega) + \log(|Q|) * |Q| * f_{k+1}(\epsilon)) \end{aligned}$$

On a donc :

$$\begin{aligned} f_k(\epsilon) &= O\left(\log(|Q|) * \frac{(|Q| * \log(|Q|))^{|Q|-k+1} - 1}{|Q| * \log(|Q|) - 1}\right) \\ &= O\left(|Q|^{|Q|+1-k}\right) \\ f_k(t\omega) &= O\left(\log(|Q| * |\Sigma|) * |Q| * \frac{(|Q| * \log(|Q|))^{|Q|-k+1} - 1}{|Q| * \log(|Q|) - 1} * f_0(\omega)\right) \\ &= O\left(|Q|^{|Q|+2-k} \log(|\Sigma|) * f_0(\omega)\right) \end{aligned}$$

soit :

$$f_0(\omega) = O\left(|Q|^{(|Q|+2)*(|\omega|+1)-1} \log(|\Sigma|)^{|\omega|}\right) = |Q|^{O(|\omega|*|Q|)} * \log(|\Sigma|)^{O(|\omega|)}$$

□

Theorem 1 (déterminisation).

Tout automate fini non déterministe se réduit à un automate fini de même langage. Le nombre d'états de l'automate créé est au plus $2^{|Q|}$ où $|Q|$ est le nombre d'états de celui de départ.

Démonstration. On commence éliminer les ϵ -transitions, càd. par réduire notre automate fini non déterministe en un automate fini non déterministe sans ϵ -transitions avec le même nombre d'état.

On nomme cet automate non déterministe sans ϵ -transitions $\mathcal{A} = (\Sigma, Q, I, F, \text{trans})$.

On utilise alors l'automate fini $\text{det}(\mathcal{A}) := (\Sigma, Q^*, I, F^*, \text{trans}^*)$ où :

- Q^* est l'ensemble $\mathcal{P}(Q)$ des sous-ensembles de Q ,
- $I \in \mathcal{P}(Q)$ est bien un état de $\text{det}(\mathcal{A})$,
- $F^* := \{A \subseteq Q \mid A \cap F \neq \emptyset\}$ est l'ensemble des ensembles d'états contenant un terminal,
- $\text{trans}^*(A, t) := \bigcup_{q \in A} \text{trans}(q, t)$ est l'ensemble des états accessibles depuis l'un des états de A en prenant une t -transition.

On prouve par induction sur $|\omega|$ que tout $\omega \in \Sigma^*$ est reconnu par \mathcal{A} si et seulement s'il est reconnu par $\det(\mathcal{A})$:

- ϵ est reconnu par \mathcal{A} ssi $I \cap F \neq \emptyset$, c'est à dire ssi $I \in F^*$,
- $t.\omega'$ est reconnu par \mathcal{A} ssi ω' est reconnu par $(\Sigma, Q, \bigcup_{q \in I} \text{trans}(q, t), F, \text{trans}) = (\Sigma, Q, \text{trans}^*(I, t), F, \text{trans})$; par hypothèse d'induction sur $|\omega'| < |\omega|$, cela est équivalent à ω' reconnu par $(\Sigma, Q^*, \text{trans}^*(I, t), F^*, \text{trans}^*)$, or cela signifie exactement que ω est reconnu par $\det(\mathcal{A})$.

□

Algorithme 6 (déterminisation).

Donnée : Un automate fini non déterministe $\mathcal{A} = (\Sigma, Q, q_0, F, \text{eps}, \text{trans})$,

Résultat : un automate fini déterministe $\det(\mathcal{A}) = (\Sigma, Q_*, q_*, F_*, \text{trans}^*)$

Spécification : $L_{\mathcal{A}} = L_{\det(\mathcal{A})}$ et $|Q_*| \leq 2^{|Q|}$.

Point fixe :

- *Ensemble* : $Q_* \subseteq Q^* := \mathcal{P}(Q)$,
- *Initialisation* : $Q_* := \{\text{eps}^*({q_0})\}$,
- *Pas* : $Q_* := Q_* \cup \{\text{eps}^*(\text{trans}^*(A, t)) \mid A \in Q_*, t \in \Sigma\}$.

Sous fonctions :

- $\text{eps}^*(A) \subseteq Q$ est la clôture de $A \subseteq Q$ par eps comme définit Exemple 2.
- $\text{trans}^*(A, t) := \bigcup_{q \in A} \text{trans}(q, t)$ pour $A \subseteq Q$ et $t \in \Sigma$.

Proposition 9. Une fois déterminisé, la reconnaissance d'un mot ω par un automate fini non déterministe $\mathcal{A} = (\Sigma, Q, q_0, F, \text{eps}, \text{trans})$ est en :

$$O(|\omega| * (|Q| + \log|\Sigma|)) \stackrel{\Sigma \text{ cst}}{\simeq} O(|\omega| * |Q|) .$$

Remarque 10. L'ensemble Q_* d'états trouvés est inférieur à la borne Q^* donnée dans la preuve du théorème :

$$Q_* \subseteq Q^* .$$

Bien qu'en théorie, elle peut atteindre cette borne, en pratique elle est TRÈS inférieure : souvent polynomial en $|Q|$ et non exponentielle ! Cela signifie, qu'en pratique, la reconnaissance après déterminisation est souvent en

$$O(|\omega| * \log(|Q| * |\Sigma|)) \stackrel{\Sigma \text{ cst}}{\simeq} O(|\omega| * \log(|Q|)) .$$

4 Automates shift-action

Intuition : Un automate shift-action ressemble a un automate habituel, sauf qu'il n'a pas d'états finaux, mais des transitions spécifiques, appelées actions, qui vont lancer une sous-routine. Celles-ci peuvent avoir lieux au milieu du mot (si on lit un caractère que l'on ne comprends pas par exemple), ou à la fin (comme l'acceptation sur l'automate habituel). On utilisera le symbole \$ pour désigner la fin de mot (correspond à EOF ou `end_of_file` dans les outils commerciaux).

Définition 24 (Automate shift-action).

Un automate shift-action \mathcal{A} est la donnée :

- d'un alphabet (de tokens) Σ ,
- d'un ensemble d'actions possibles Δ avec le code associé dont une spécifique **Err** qui échoue,
- d'un ensemble d'états Q ,
- d'un état $q_0 \in Q$ initial,
- d'un choix d'action $C(q) \subseteq \Sigma$ pour chaque état $q \in Q$,
- d'une transition (ou shift) $\text{shift}(q, t) \in Q$ pour chaque état $q \in Q$ et caractère $t \in \Sigma - C(q)$,
- d'une action $\text{action}(q, t) \in \Delta$ pour chaque état $q \in Q$ et caractère $t \in C(q) \cup \{\$\}$.

Les actions $\delta \in \Delta$ sont des codes code arbitraires qui dépendent de deux paramètres : le buffer $u \in \Sigma^*$ et le mot qui restera à lire $\omega \in \Sigma^*$.

Exemple 4.

1. L'automate shift-action $\mathcal{A}_\perp = (\Sigma, \text{Err}, \{\perp\}, \perp, C, \text{shift}, \text{action})$ où :
 - C définit par $C(\perp) = \Sigma$
 - on n'a pas de shift car $\Sigma - C(\perp) := \emptyset$,
 - on renvoi toujours une erreur avec $\text{action}(\perp, t) := \text{Err}$ pour tout t .
2. Un automate \mathcal{A} est un automate shift-action si on considère $(\Sigma, \{\text{Acc}, \text{Err}\}, \{\perp\}, \perp, C, \text{shift}, \text{action})$ où :
 - **Acc** se contente d'accepter le mot,
 - C définit par $C(a) = \emptyset$,
 - $\text{shift}(q, a) = \delta_A(q, a)$,
 - $\text{action}(q, \$) := \text{Acc}$ si $q \in F_A$ et $\text{action}(q, \$) := \text{Err}$ sinon.
3. Un automate \mathcal{A} avec puits positifs et négatifs est un automate shift-action $(\Sigma, \{\text{Acc}, \text{Err}\}, \{\perp\}, \perp, C, \text{shift}, \text{action})$ où :
 - **Acc** se contente d'accepter le mot,
 - C définit par $C(a)$ est l'ensemble des transitions vers un puit (positif ou négatif),
 - $\text{shift}(q, a) = \delta_A(q, a)$ sont les transitions n'allant pas vers des puits,
 - $\text{action}(q, \$) := \text{Acc}$ si $q \in F_A$ et $\text{action}(q, \$) := \text{Err}$ sinon,
 - $\text{action}(q, a) := \text{Acc}$ si la transition $\delta(q, a)$ est un puits positif et $\text{action}(q, a) := \text{Err}$ si c'est un puits négatif.

Notations 10. Pour des raisons de lisibilité, on ne dessinera généralement pas les actions **Err**. De plus, il y a généralement peu d'états avec des actions (autres que **Err**), on notera ces états avec un double cercle (ou double carré lorsque l'on écrit des états carrés) et les action qui en partent avec des flèches étiquetées sortantes ; on peut se permettre de ne pas marquer les actions sortantes lorsqu'elles sont décrites par ailleurs.

Définition 25 (Exécution sur un automate shift-action).

On considère un automate shift-action $\mathcal{A} = (\Sigma, \Delta, Q, q_0, C, \text{shift}, \text{action})$, un mot ω et un buffer $u \in \Sigma^*$ (vide au départ). L'exécution de ω sur \mathcal{A} se déroule par induction :

- argument inductif : ω ,
- Cas ϵ : on exécute l'action $\text{action}(q_0, \$)(\mathbf{u}, \epsilon)$,
- Cas $t\omega$ lorsque $t \in C(q)$: on exécute l'action $\text{action}(q_0, t)(\mathbf{u}, \omega)$,
- Cas $t\omega$ lorsque $t \notin C(q)$: on place t dans le buffer \mathbf{u} , et on exécute l'automate shift-action $(\Sigma, Q, \Delta, \text{shift}(q_0, t), C, \text{shift}, \text{action})$ sur ω .

On note $L_{\mathcal{A}}$ l'ensemble des quadruplets $(\omega, \delta, \omega', \mathbf{u}) \in \Sigma^* \times \Delta \times \Sigma^* \times \Sigma^*$ telle que \mathcal{A} exécuté sur ω va lancer δ sur le buffer \mathbf{u} et le résidu ω' .

Remarque 11. Si $(\omega, \delta, \omega', \mathbf{u}) \in L_{\mathcal{A}}$, alors $\omega = \mathbf{u}.\omega'$.

Remarque 12. Ce formalisme est très (trop ?) lax, au sens où une action est capable de faire n'importe quoi, y compris relancer l'automate ou utiliser le buffer \mathbf{u} (c'est d'ailleurs pour ça qu'il est là). Ce qui est important est que la phase avec les shifts se comporte comme un automate fini, on peut donc la minimiser, ainsi que déterminer sa version non déterministe !

Algorithme 7 (Co-accessibilité).

Donnée : Un automate shift-action $\mathcal{A} = (\Sigma, \Delta, Q, q_0, C, \text{shift}, \text{action})$,

Résultat : un automate shift-action $\text{cac}(\mathcal{A}) = (\Sigma, \Delta, \text{cac}(Q), q'_0, C', \text{shift}, \text{action}')$,

Spécification : $L_{\mathcal{A}} = L_{\text{cac}(\mathcal{A})}$ et $\text{cac}(Q) \subseteq Q$ et tous les états de $\text{cac}(\mathcal{A})$ co-accessibles.

Point fixe :

— Ensemble : $\text{cac}(Q) \subseteq Q$

— Initialisation : $\text{cac}(Q) := \{q \mid \exists t \in C(q) \cup \{\$\}, \text{action}(q, t) \neq \text{Err}\}$,

— Pas : $\text{cac}(Q) := \text{cac}(Q) \cup \{q \mid \exists t \in \Sigma - C(q), \text{trans}(q, t) \in \text{cac}(Q)\}$.

À la fin, on doit vérifier que $q_0 \in \text{cac}(Q)$, si c'est le cas $q'_0 = q_0$ sinon on rend l'automate shift-action qui rend toujours une erreur $\text{cac}(\mathcal{A}) := \mathcal{A}_{\perp}$. De plus, on complète les transitions libres avec des erreurs :

$F'(q) := \{t \mid t \in C(q) \text{ où } \text{shift}(q, t) \notin Q_*\}$,

$\text{action}'(q, t) := \begin{cases} \text{action}(q, t) & \text{si } t \in C(q) \\ \text{Err} & \text{si } \text{shift}(q, t) \notin Q_* \end{cases}$

Algorithme 8 (Accessibilité).

Donnée : Un automate shift-action $\mathcal{A} = (\Sigma, \Delta, Q, q_0, F, \text{shift}, \text{action})$,

Résultat : un automate shift-action partiel $\text{acc}(\mathcal{A}) = (\Sigma, \Delta, \text{acc}(Q), q_0, F, \text{shift}, \text{action})$

Spécification : $L_{\mathcal{A}} = L_{\text{acc}(\mathcal{A})}$ et $\text{acc}(Q) \subseteq Q$ et tous les états de $\text{acc}(\mathcal{A})$ accessibles.

Point fixe :

— Ensemble : $\text{acc}(Q) \subseteq Q$

— Initialisation : $\text{acc}(Q) := \{q_0\}$ si $q_0 \neq \perp$ et \emptyset sinon,

— Pas : $\text{acc}(Q) := \text{acc}(Q) \cup \{\text{shift}(q, t) \mid q \in \text{acc}(Q), t \in \Sigma - C(q)\}$.

Algorithme 9 (Minimisation).

Donnée : Un automate shift-action $\mathcal{A} = (\Sigma, \Delta, Q, q_0, F, \text{shift}, \text{action})$,

Résultat : un automate shift-action $\text{min}(\mathcal{A}) = (\Sigma, \Delta, \text{min}(Q), q_0^m, F^m, \text{shift}^m, \text{action}^m)$

Spécification : $L_{\mathcal{A}} = L_{\text{min}(\mathcal{A})}$ et $\text{min}(\mathcal{A})$ est minimal parmi les automates ayant cette propriété.

Étape 1 : calculer $\text{cac}(\mathcal{A})$ avec l'Algorithme 7,

Étape 2 : calculer $\text{acc}(\text{cac}(\mathcal{A}))$ avec l'Algorithme 8,

Étape 3 : rajouter un état puits et une erreur pour retrouver un automate total $\mathcal{A}' = (\Sigma, Q', q'_0, F', \text{trans})$

Étape 4 : on cherche une relation d'équivalence (\equiv) sur Q' appelée bisimulation, pour ça on procède par point fixe :

- Ensemble : $(\equiv) \subseteq Q'^2$,
- Initialisation : $(\equiv) := (=)$ c'est à dire que deux états sont équivalents ssi ils sont égaux,
- Pas : $p \equiv q$ lorsque :
 - $C(p) = C(q)$,
 - pour tout $t \in \Sigma - C(p)$, $\text{shift}(p, t) \equiv \text{shift}(q, t)$,
 - et pour tout $t \in C(p) \cup \{\$\}$, $\text{action}(p, t) = \text{action}(q, t)$.

Étape 5 : On pose alors

- $\min(Q) := Q_{/\equiv}$ est l'ensemble des classes d'équivalences de \equiv ,
- $q_0^m := [q_0]$ est la classe d'équivalence de q_0 ,
- $F^m([p]) := [C(p)]$,
- $\text{shift}^m([q], t) := [\text{shift}(q, t)]$
- $\text{action}^m([q], t) := \text{action}(q, t)$

Définition 26 (Automate shift-action non déterministe).

Un automate shift-action non déterministe \mathcal{A} est la donnée :

- d'un alphabet (de tokens) Σ ,
- d'un ensemble d'états Q ,
- d'un ensemble d'actions possibles Δ avec le code associé,
- d'un état $q_0 \in Q$ initial,
- d'un ensemble d' ϵ -transitions $\text{eps}(q) \subseteq Q$ pour chaque état $q \in Q$,
- d'un ensemble de transitions $\text{shift}(q, t) \subseteq Q$ pour chaque état $q \in Q$ et caractère $t \in \Sigma$
- d'un ensemble d'actions $\text{action}(q, t) \subseteq \Delta$ pour chaque état $q \in Q$ et caractère $t \in \Sigma \cup \$$.

Définition 27 (Automate shift-action psedo-déterministe).

Un automate shift-action psedo-déterministe \mathcal{A} est la donnée :

- d'un alphabet (de tokens) Σ ,
- d'un ensemble d'actions possibles Δ ,
- d'un ensemble d'états Q ,
- d'un état $q_0 \in Q$ initial,
- d'un choix $C(q) \subseteq \Sigma$ pour chaque état $q \in Q$,
- d'une transition $\text{trans}(q, t) \in Q$ pour chaque état $q \in Q$ et caractère $t \in \Sigma - C(q)$
- d'un ensemble d'actions $\text{action}(q, t) \subseteq \Delta$ pour chaque état $q \in Q$ et caractère $t \in \Sigma \cup \$$.

Algorithme 10 (psedo-déterminisation).

Donnée : Un automate shift-action non déterministe $\mathcal{A} = (\Sigma, Q, \Delta, q_0, \text{eps}, \text{shift}, \text{action})$,

Résultat : un automate shift-action psedo-déterministe $\text{det}(\mathcal{A}) = (\Sigma, \Delta \cup \{\text{Err}\}, Q_*, q_0^*, C, \text{shift}^*, \text{action}^*)$

Spécification : $L_{\mathcal{A}} = L_{\det(\mathcal{A})}$ et $|Q'| \leq 2^{|Q|}$.

Point fixe :

- *Ensemble* : $Q_* \subseteq \mathcal{P}(Q) - \{\emptyset\}$ les ensembles non vides d'états,
- *Initialisation* : $Q_* := \{\text{eps}^*(\{q_0\})\}$,
- *Pas* : $Q_* := Q_* \cup \{\text{eps}^*(\text{shift}^*(A, t)) \mid A \in Q_*, t \in \Sigma\}$.

Sous fonctions :

- $\text{eps}^*(A) \subseteq Q$ est la clôture de $A \subseteq Q$ par **eps** comme définit Exemple 2,
- $C(A) := \{t \mid \text{shift}^*(A, t) = \emptyset\}$,
- $\text{shift}^*(A, t) := \bigcup_{q \in A} \text{shift}(q, t)$ pour $A \subseteq Q$ et $t \notin C(A)$,
- $\text{action}^*(A, t) := \bigcup_{q \in A} \text{action}(q, t)$ pour $A \subseteq Q$ et $t \notin C(A)$,
- $\text{action}^*(A, t) := \bigcup_{q \in A} \text{action}(q, t)$ pour $A \subseteq Q$ et $t \in C(A)$ si $\bigcup_{q \in A} \text{action}(q, t)$ est non vide,
- $\text{action}^*(A, t) := \{\text{Err}\}$ pour $A \subseteq Q$ et $t \in C(A)$ si $\bigcup_{q \in A} \text{action}(q, t)$ est vide,

Proposition 10 (déterminisation des automates shift-action).

Soit $\det(\mathcal{A}) = (\Sigma, Q, \Delta, q_0, \text{shift}, \text{action})$ un automate psedo-déterministe.

Cet automate est déterministe si pour tout $q \in Q$ et tout $t \in \Sigma$:

1. si $t \in C(q)$, $\text{action}(q, t)$ est un singleton,
2. si $t \notin C(q)$, $\text{action}(q, t)$ est un vide.

Définition 28 (conflits shift/action et action/action).

Soit $\det(\mathcal{A}) = (\Sigma, Q, \Delta, q_0, \perp, \text{shift}, \text{action})$ un automate psedo-déterministe.

On dit qu'il y a

1. un conflit action/action s'il existe $q \in Q$ et $t \in \Sigma$ tel que $|\text{action}(q, t)| > 1$,
2. un conflit shift/action s'il existe $q \in Q$ et $t \in \Sigma$ tel que $|\text{action}(q, t)| \geq 1$ et $\text{shift}(q, t) \neq \perp$.

Définition 29 (Automate shift-action avec priorité).

Un automate shift-action avec priorité est un automate shift-action non déterministe $\mathcal{A} = (\Sigma, Q, \Delta, q_0, \text{eps}, \text{shift}, \text{action})$ avec une notion de priorité sur les actions. Celle-ci prend la forme d'un ordre stricte ($<$) sur Δ .

Proposition 11 (déterminisation avec priorité).

Donnée : Un automate shift-action avec priorité $\mathcal{A} = (\Sigma, Q, (\Delta, <), q_0, \text{eps}, \text{shift}, \text{action})$,

Résultat : un automate shift-action $\det(\mathcal{A}) = (\Sigma, Q_*, \Delta, q_0^*, C, \text{shift}^*, \text{action}^*)$

Spécification : $L_{\mathcal{A}} = L_{\det(\mathcal{A})}$ et $|Q'| \leq 2^{|Q|}$.

Point fixe :

- *Ensemble* : $Q_* \subseteq Q^* := \mathcal{P}(Q) - \{\emptyset\}$,
- *Initialisation* : $Q_* := \{\text{eps}^*(\{q_0\})\}$,
- *Pas* : $Q_* := Q_* \cup \{\text{eps}^*(\text{shift}^*(A, t)) \mid A \in Q_*, t \in \Sigma\}$.

Sous fonctions :

- $\text{eps}^*(A) \subseteq Q$ est la clôture de $A \subseteq Q$ par **eps** comme définit Exemple 2.
- $C(A) := \{t \mid \exists q \in A, \text{shift}(q, t) \neq \emptyset\}$,

- $\text{shift}^*(A, t) := \bigcup_{q \in A} \text{shift}(q, t) \in \mathcal{P}(Q) - \{\emptyset\}$ pour $A \subseteq Q$ et $t \in \Sigma - C(A)$,
- $\text{action}^*(A, t) := \bigwedge \text{action}(q, t)$ pour $A \subseteq Q$ et $t \in C(A) \cup \{\$\}$, est l'action la plus basse dans l'ordre de priorité, avec $\bigwedge \emptyset = \text{Err}$.

5 Analyse lexicale

5.1 Lexeur

Un fonction de lexing est un programme qui prend en entrée une chaîne de caractères et rend en sortie un mot de token. Ici (pour plus d'uniformité) on la généralise comme un programme qui prend en entrée un mot de token et rend en sortie un autre mot de token.

Définition 30 (Schéma de lexing).

Un schéma de lexing est la donnée :

- d'un alphabet (de tokens) Σ d'entrée,
- d'un alphabet (de tokens) Γ de sortie,
- d'une séquence $r_1, \dots, r_n \in \text{Reg}(\Sigma)$ d'expressions régulières sur Σ ,
- pour tout $i \leq n$, d'une fonction $\delta_i(\omega) \in \Gamma$ qui à chaque mot (de token) $\omega \in L_{r_i}$ reconnu par le l'expression régulière r_i , associe un token de Γ

Définition 31 (Fonction de lexing).

Une fonction de lexing $\mathcal{L} : \Sigma^* \rightarrow \Gamma^* \cup \{\text{Err}\}$ est correct pour un schéma de lexing $(\Sigma, \Gamma, (r_i)_{i \leq n}, (\delta_i)_{i \leq n})$ si pour tout $\omega \in \Sigma^*$:

- ou bien $\mathcal{L}(\omega) = \text{Err}$,
- ou bien $\omega = \omega_1 \dots \omega_k$ et il exists $i_1, \dots, i_k \in [1, n]$ tels que pour tout $j \leq k$, $\omega_j \in L_{r_{i_j}}$ est reconnu par r_{i_j} et $\mathcal{L}(\omega) = \delta_{i_1}(\omega_1) \dots \delta_{i_k}(\omega_k)$.

Définition 32 (Lexeur non déterministe).

Le lexeur non déterministe associé au schéma de lexing $(\Sigma, \Gamma, (r_i)_{i \leq n}, (\delta_i)_{i \leq n})$ est l'automate shift-action avec priorité dont :

- l'alphabet de token est Σ ,
- les états sont :
 - l'état initial q_0 ,
 - les couples (i, q) tels que $i \leq n$ et $q \in \mathcal{Q}_{A_{r_i}} - \{\perp\}$ est un état de l'automate reconnaissant r_i (sauf \perp que l'on enlève),
- les actions possibles $\Delta := \{\delta'_i \mid i \leq n\}$, ordonnés par i , sont les actions δ'_i qui sur le buffer \mathbf{u} et le résiduel ω vont renvoyer $\delta_i(\mathbf{u})$ sur un buffer de sortie, puis soit terminer si $\omega = \epsilon$, soit relancer l'exécution du lexeur sur ω ,
- l'état initial est q_0 ,
- les ϵ -transitions vont l'état initial vers chacun des états initiaux des automates des expressions régulières :

$$\text{eps}(q_0) = \{(i, q_0) \mid i \leq n\}, \quad \text{eps}((i, q)) = \emptyset,$$

— les transitions sont celles des différents automates finis :

$$\mathbf{shift}(q_0, t) = \emptyset, \quad \mathbf{shift}((i, q), t) = \{(i, \mathbf{trans}(q, t)) \mid \mathbf{trans}(q, t) \neq \perp\},$$

— les actions sont activées lorsque l'on arrive plus à lire un caractère :

$$\mathbf{action}(q_0, t) = \emptyset, \quad \mathbf{action}((i, q), t) = \{\delta_i \mid \mathbf{trans}(q, t) = \perp\}.$$

Définition 33 (Lexeur).

Le lexeur associé au schéma de lexing $(\Sigma, \Gamma, (r_i)_{i \leq n}, (\delta_i)_{i \leq n})$ est la détermination et minimisation du lexeur non-déterministe.

Theorem 2. La fonction de lexing calculé par l'exécution du lexeur associé à un schéma est correcte pour le même schéma.

6 Analyse Syntaxique

6.1 Parseur LR

6.1.1 Parseur pour une grammaire

Définition 34 (Parseur LR).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire. Un parseur (resp. non-déterministe) LR pour \mathcal{G} est un automate shift-action (resp. non-déterministe) dont :

- l'alphabet $\mathcal{T} \uplus \mathcal{N}$ est formé par les terminaux et non terminaux,
- les actions sont de deux types $\Delta = \{\mathbf{err}, \mathbf{accept}\} \cup \{\mathbf{reduce}_{N \mapsto \omega} \mid N \in \mathcal{N}\omega \in \tau(N)\}$:
 - \mathbf{err} renvoie toujours une erreur,⁵
 - \mathbf{accept} renvoie le buffer u si $u = S$ et une erreur sinon,
 - $\mathbf{reduce}_{N \mapsto \omega}$ renvoie une erreur si $u \neq v\omega$, c.à.d. si le buffer u ne termine par ω , si on a bien $u = v\omega$ alors on forme un arbre syntaxique $\mathbf{N} := N(\omega)$ et on relance l'automate shift-action sur $v\mathbf{N}\omega'$ où ω' était le résiduel.
- l'action \mathbf{accept} ne peut être renvoyé que en voyant la fin du mot : $\mathbf{accept} \in \mathbf{action}(q, t) \Rightarrow t = \$$.

Définition 35. Un parseur (resp. non-déterministe) LR est correcte si :

- les accepts et reduces ne déclencheront jamais d'erreurs (on peut avoir des erreur \mathbf{err}),
- un mot est accepté par le parseur si et seulement si il est accepté par sa grammaire,
- un mot accepté renvoie un de ses arbres de dérivation dans la grammaire.

6.1.2 Parseur canonique d'une grammaire

Définition 36 (Parseur (non déterministe) canonique).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire. Le parseur canonique de \mathcal{G} est le parseur non-déterministe LR sur \mathcal{G} donné par :

5. on peut aussi mettre plusieurs erreurs plus expressives.

- les états $Q = \mathcal{N} \uplus \bigcup_{N \in \mathcal{L}} \tau_\bullet(N) \uplus \{(* \mapsto \bullet S), (* \mapsto S \bullet)\}$ sont les non-terminaux, les règles partielles, et deux états particuliers se comportant comme des règles partielles spéciales : l'état initial $q_0 = (* \mapsto \bullet S) \in L$, l'état final $\text{accept} = (* \mapsto S \bullet)$,
- pas de shift ni d'action depuis un non-terminal :

$$\text{shift}(N, \underline{t}) = \text{action}(N, \underline{t}) = \emptyset,$$

- par contre il y a des ϵ -transitions :

$$\text{epsilon}(N) = \{N \mapsto \bullet \omega \mid \omega \in \tau(N)\},$$

- pas d'action depuis une règle partielle non-final :

$$\text{action}(N \mapsto \omega_1 \bullet \underline{c} \omega_2, \underline{t}) = \emptyset,$$

- les transitions de shift depuis une règle partielle correspondent à avancer dans la règle :

$$\text{shift}((N \mapsto \omega_1 \bullet \underline{c} \omega_2), \underline{c}') = \begin{cases} \{N \mapsto \omega_1 \underline{c} \bullet \omega_2\} & \text{si } \underline{c} = \underline{c}' \\ \emptyset & \text{sinon} \end{cases},$$

- l'ensemble des ϵ -transitions de depuis une règles partielle est le non-terminal qu'il faut lire si besoin :

$$\text{epsilon}(N \mapsto \omega_1 \bullet \underline{c} \omega_2) = \begin{cases} \{\underline{c}\} & \text{si } \underline{c} \text{ est un non-terminal} \\ \emptyset & \text{si } \underline{c} \text{ est un terminal} \end{cases},$$

- pas de shift ni de epsilon depuis une règle partielle finale :

$$\text{shift}(N \mapsto \omega \bullet, \underline{c}) = \text{epsilon}(N \mapsto \omega \bullet) = \emptyset,$$

- les reduce des règles partielle finales sont les réécriture indiquées par les états, la seul transition final est depuis la règle $* \mapsto S \bullet$:

$$\text{shift}((N \mapsto \omega \bullet), \underline{c}) = \begin{cases} \{\text{reduce}_{N \mapsto \omega}\} & \text{si } N \neq * \\ \{\text{accept}\} & \text{si } N = * \text{ et } \underline{c} = \$ \\ \emptyset & \text{sinon} \end{cases},$$

Définition 37 (Parseur spécial canonique).

Il s'agit du parseur canonique d'une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ sauf que l'on restreint les reduces en utilisant les follows :

$$\text{shift}((N \mapsto \omega \bullet), \underline{c}) = \begin{cases} \{\text{reduce}_{N \mapsto \omega}\} & \text{si } N \neq * \text{ et } \underline{c} \in \text{follow}(N) \\ \{\text{accept}\} & \text{si } N = * \text{ et } \underline{c} = \$ \\ \emptyset & \text{sinon} \end{cases},$$

Theorem 3. Le parseur canonique de \mathcal{G} est correcte (pour \mathcal{G}).

6.1.3 (Pseudo-)parseur LR_0

L'algorithme LR_0 est simple : il s'agit de pseudo-déterminiser le parseur canonique de notre grammaire, on échoue si le résultat n'est pas déterministe.

Algorithme 11 (Construction du pseudo-parseur LR_0).

Donnée : une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$

Résultat : un parseur pseudo-déterministe $LR_0(\mathcal{G})$ sur \mathcal{G} ,

Spécification : $LR_0(\mathcal{G})$ est correcte pour \mathcal{G} .

Étape 1 : construire le parseur canonique de \mathcal{G} ,

Étape 2 : le pseudo-déterminiser,

Étape 3 : vérifier qu'il n'y a pas de conflit,

Étape 4 (optionnelle) : le minimiser.

Définition 38 (Parseur et grammaire LR_0).

Une grammaire est LR_0 si $LR_0(\mathcal{G})$, son pseudo-parseur LR_0 , se trouve être déterministe au sens de la Proposition 10.

Le parseur correspondant est nommé le parseur LR_0 de la grammaire.

6.1.4 (Pseudo-)parseur SLR

Pour l'algorithme SLR, on prend en compte les follows pour éviter quelques conflits :

Algorithme 12 (Construction du pseudo-parseur SLR).

Donnée : une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$

Résultat : un parseur pseudo-déterministe $SLR(\mathcal{G})$ sur \mathcal{G} ,

Spécification : $SLR(\mathcal{G})$ est correcte pour \mathcal{G} .

Étape 1 : construire le parseur spécial canonique de \mathcal{G} ,

Étape 2 : le pseudo-déterminiser,

Étape 3 : vérifier qu'il n'y a pas de conflit,

Étape 4 (optionnelle) : le minimiser.

Définition 39 (Parseur et grammaire SLR).

Une grammaire est SLR si $SLR(\mathcal{G})$, son pseudo-parseur SLR, se trouve être déterministe au sens de la Proposition 10.

Le parseur correspondant est nommé le parseur SLR de la grammaire.

6.1.5 L'algorithme LR_1^\dagger

Le principe est simple : on modifie la grammaire en dupliquant les non-terminals de sorte à ce que chaque non-terminal n'ait qu'un seul follow.

Définition 40 (Grammaire LR_1 -étendue).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire.

La LR_1 -extension de \mathcal{G} est la grammaire $\text{exten}_{LR_1}(\mathcal{T}, \mathcal{N}_{LR_1}, \tau_{LR_1})$ définie par :

— les non-terminaux sont dupliqués pour chaque follow :

$$\mathcal{N}_{LR_1} := \{(N, t) \in \mathcal{N} \times \mathcal{T} \mid t \in \text{follow}(N)\},$$

on écrira N_t pour (N, t) ,

— règles sont dupliquées et modifiées en conséquent :

$$\tau_{LR_1}(N_t) := \bigcup_{\omega \in \tau(N)} \omega_t$$

où $\omega_t \subseteq (\mathcal{T}_{LR_1} \uplus \mathcal{N})_{LR_1}$ est défini par induction :

$$(\omega s)_t := \{\omega' s \mid \omega' \in \omega_u\} \quad (\omega M)_t := \left\{ \omega' M_t \mid \omega' \in \bigcup_{t' \in \text{first}(Mt)} v_{t'} \right\}.$$

Proposition 12. La LR_1 -extension $\text{exten}_{LR_1}(G)$ d'une grammaire G à des singletons comme follow :

$$\forall N_t \in \mathcal{N}_{LR_1}, \text{follow}(N_t) = \{t\}.$$

Algorithme 13 (Construction du psedo-parseur LR_1).

Donnée : une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$

Résultat : un parseur psedo-déterministe LR $LR_1(\mathcal{G})$ sur \mathcal{G} ,

Spécification : $LR_1(\mathcal{G})$ est correcte pour \mathcal{G} .

Étape 1 : calculer les firsts et follows,

Étape 2 : calculer la LR_1 -extension $\text{exten}_{LR_1}(\mathcal{G})$ de \mathcal{G} ,

Étape 3 : construction du psedo-parseur SLR de $\text{exten}_{LR_1}(\mathcal{G})$,

Étape 4 : remplacer les reduces par des G -reduces en enlevant les annotations de follow,

Étape 5 : vérifier qu'il n'y a pas de conflit.

Définition 41 (Parseur et grammaire LR_1).

Une grammaire est LR_1 si $LR_1(\mathcal{G})$, son psedo-parseur LR_1 , se trouve être déterministe au sens de la Proposition 10.

Le parseur correspondant est nommé le parseur LR_1 de la grammaire.

6.1.6 L'algorithme $LALR^\dagger$

Cette fois, on modifie la grammaire en dupliquant les non-terminaux de sorte à ce qu'à chaque non-terminal N corresponde un seul shift dans le parseur résultant :

$$\forall N, \exists! q, \text{shift}(q, N) \neq \emptyset$$

Remarquez que l'on a l'équivalence :

Proposition 13. Dans une psedo-extention SLR , on a la propriété suivante :

$$\text{shift}(q, N) \neq \emptyset \Leftrightarrow \exists \omega, (N \mapsto \bullet \omega) \in q.$$

Une autre façon de décrire l'extension grammaticale est de dire qu'après la détermination du parseur LR spécial canonique, chaque non-terminal N ne se retrouve que dans un seul état du pseudo-parseur SLR.

Définition 42 (Grammaire LALR-étendue).

Soit $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$ une grammaire et $\text{SLR}(\mathcal{G}) = (\mathcal{T}, \mathcal{N}, q_0, \perp, \text{shift}, \text{action})$ le pseudo-parseur SLR de G .

La LALR-extension de G est la grammaire $\text{exten}_{\text{LALR}}(\mathcal{T}, \mathcal{N}_{\text{LALR}}, \tau_{\text{LALR}})$ définie par :

les non-terminaux sont dupliqués pour chaque utilisation dans \mathcal{A} :

$$\mathcal{N}_{\text{LALR}} := \{(N, q) \in \mathcal{N} \times \mathcal{T} \mid \text{shift}(q, N) \neq \emptyset\},$$

on écrira N_q pour (N, q) ,

règles sont dupliquées et modifiées en conséquent :

$$\tau_{\text{LR}_1}(N_q) := \bigcup_{\omega \in \tau(N)} \text{LA}(q, N \mapsto \omega \bullet \omega)$$

où l'ensemble $\text{LA}(q, N \mapsto \omega' \bullet \omega)$ est défini par induction pour tout $(N \mapsto \omega' \bullet \omega) \in q$:

- $\text{LA}(q, N \mapsto \omega' \bullet) = \epsilon$,
- $\text{LA}(q, N \mapsto \omega' \bullet t\omega) = t.\text{LA}(q', N \mapsto \omega't \bullet \omega)$ avec $\{q'\} = \text{shift}(q, t)$ si t est un terminal,
- $\text{LA}(q, N \mapsto \omega' \bullet M\omega) = M_q.\text{LA}(q', N \mapsto \omega'M \bullet \omega)$ avec $\{q'\} = \text{shift}(q, M)$ si M est un non-terminal.

Proposition 14. $\text{SLR}(\text{exten}_{\text{LALR}}(\mathcal{G}))$ est au pire aussi gros (en nombre d'états) que $\text{SLR}(\mathcal{G})$.

Algorithme 14 (Construction du pseudo-parseur LR_1).

Donnée : une grammaire $\mathcal{G} = (\mathcal{T}, \mathcal{N}, S, \tau)$

Résultat : un parseur pseudo-déterministe LR $\text{LALR}(\mathcal{G})$ sur \mathcal{G} ,

Spécification : $\text{LALR}(\mathcal{G})$ est correcte pour \mathcal{G} .

Étape 1 : calculer $\text{SLR}(\mathcal{G})$, le pseudo-parseur SLR de \mathcal{G} ,

Étape 2 : calculer $\text{exten}_{\text{LALR}}(\mathcal{G})$ la LALR-extension de \mathcal{G} ,

Étape 3 : construction du pseudo-parseur SLR de $\text{exten}_{\text{LR}_1}(\mathcal{G})$,

Étape 4 : remplacer les réduces par des G -réduces en enlevant les annotations,

Étape 5 : vérifier qu'il n'y a pas de conflit.

Définition 43 (Parseur et grammaire LALR).

Une grammaire est LALR si $\text{LALR}(\mathcal{G})$, son pseudo-parseur LALR, se trouve être déterministe au sens de la Proposition 10.

Le parseur correspondant est nommé le parseur LR_1 de la grammaire.