

Rapport Projet ASCO:
Analyseur syntaxique et sémantique pour un
sous-langage de C

AZIZE Salma et RAJAT Chorouk

Semestre 3, 2025–2026

Table des matières

1	Introduction	3
2	Structure de l'AST	3
2.1	Organisation des types	3
2.2	Choix de conception	4
3	Lexer : analyse lexicale	4
3.1	Gestion de l'ordre des règles	4
3.2	Constantes flottantes	4
3.3	Chaînes de caractères et concaténation	5
3.4	Commentaires	5
4	Parser : analyse syntaxique	5
4.1	Organisation de la grammaire	6
4.2	Gestion des précédences	6
4.3	Déclarateurs et pointeurs	6
4.4	Résolution du conflit if/else	7
4.5	Le problème de l'entrelacement déclarations/fonctions	7
4.6	Suppression des conflits	8
5	Analyse sémantique	8
5.1	Table des symboles hiérarchique	8
5.2	Simplification des types	9
5.3	Règles de typage strictes	9
5.4	Vérification des fonctions	9
6	Débogage et tests	10
6.1	Outils de débogage	10
6.2	Visualisation de l'AST	10
6.3	Tests automatisés	11
6.4	Résultats des tests	11
7	Conclusion	11

1 Introduction

Ce projet consiste à développer un analyseur syntaxique et sémantique pour C⁻, un sous-langage simplifié du C. L'objectif est de construire un arbre de syntaxe abstraite (AST) à partir d'un fichier source, puis de vérifier la cohérence des portées et des types selon les règles définies dans la spécification.

On a choisi d'implémenter le projet en OCaml, en utilisant `ocamllex` et `ocamlyacc` comme imposé par le sujet. Le développement s'est fait sous Linux (via WSL) avec Visual Studio Code, ce qui nous a permis de compiler et tester facilement avec un Makefile.

2 Structure de l'AST

En premier lieu, on a défini les types nécessaires pour générer l'arbre de syntaxe abstraite dans le fichier `ast.ml`. Cette étape a été guidée directement par la grammaire fournie dans le sujet.

2.1 Organisation des types

On a structuré l'AST en plusieurs catégories de types :

Types de base et attributs. Les types de base (`void`, `char`, `int`, `float`, `double`) sont représentés par un variant simple, accompagnés d'une liste d'attributs optionnels (`signed`, `unsigned`, `short`, `long`). Les pointeurs sont gérés récursivement via un constructeur `Pointer`.

```
1 type base_type = TVoid | TChar | TInt | TFloat | TDouble
2 type type_attr = Signed | Unsigned | Short | Long
3 type typ =
4   | BaseType of base_type * type_attr list
5   | Pointer of typ
```

Constantes et expressions. Les constantes incluent les entiers, flottants, caractères et chaînes. Pour les expressions, on a séparé les opérateurs unaires et binaires dans deux types distincts, ce qui simplifie énormément la vérification de types par la suite.

```
1 type binop = Add | Sub | Mul | Div | Mod | Lt | Gt | ...
2 type unop = Deref | Addr | Neg | UnaryPlus | UnaryMinus
3 type expr =
4   | Identifier of string
5   | Constant of constant
6   | BinaryOp of binop * expr * expr
7   | UnaryOp of unop * expr
8   | ...
```

Les affectations composées (`+=`, `-=`, etc.) sont traitées comme des opérateurs binaires normaux, ce qui évite de dupliquer la logique de vérification.

Instructions et blocs. Les instructions couvrent toutes les structures de contrôle du langage. Un bloc contient à la fois des déclarations locales et des instructions, conformément à la structure du C.

Fonctions et fichiers. Une fonction est représentée par un record contenant le type de retour, le nom, les paramètres et le corps. Un fichier C⁻ complet est représenté par deux listes distinctes : déclarations globales et définitions de fonctions.

2.2 Choix de conception

Au départ, on avait une seule liste mixte pour les déclarations et fonctions, mais séparer les deux dès l'AST a grandement simplifié l'analyse sémantique. Ce choix s'est révélé particulièrement utile lors de la vérification de portée, où on doit d'abord enregistrer toutes les fonctions avant de vérifier leur corps.

3 Lexer : analyse lexicale

La deuxième étape a été l'écriture du lexer avec `ocamllex`. Cette étape va de pair avec la définition des tokens du parser, car ce sont les valeurs de retour du lexer.

3.1 Gestion de l'ordre des règles

Le premier piège qu'on a rencontré concerne l'ordre des règles dans le lexer. Les opérateurs composés doivent impérativement être définis avant les opérateurs simples, sinon le lexer lit plusieurs tokens au lieu d'un seul.

```
1 | '+=', { PLUS_ASSIGN }
2 | '-=', { MINUS_ASSIGN }
3 | '==', { EQ }
4 | '<=', { LEQ }
5 | '>=', { GEQ }
6 | '+', { PLUS }
7 | '=', { ASSIGN }
8 | '<', { LT }
9 | '>', { GT }
```

Sans cette précaution, `==` était reconnu comme deux tokens `=` successifs, ce qui provoquait des erreurs difficiles à diagnostiquer dans le parser. Cette erreur nous a coûté plusieurs heures de débogage au début du projet.

3.2 Constantes flottantes

La partie la plus délicate du lexer a été la reconnaissance des constantes flottantes. La spécification de C⁻ autorise plusieurs formes (`1.5`, `.5`, `1.`, `1e10`, etc.) mais interdit certaines combinaisons comme `. seul ou .e10`.

Au début, on laissait passer ces cas invalides et on avait des erreurs bizarres plus loin dans le parser. Pour corriger ça, on a ajouté une règle qui capture explicitement les formes interdites *avant* les formes valides :

```
1 | invalid_float_no_parts {
2 |   raise (LexicalError 'Constante flottante invalide')
3 | }
4 | float_with_both_parts as f {
5 |   FLOAT_CONST(float_of_string f)
6 | }
7 | float_decimal_only as f {
8 |   FLOAT_CONST(float_of_string f)
9 | }
10 | float_exp_only as f {
11 |   FLOAT_CONST(float_of_string f)
12 | }
```

Cette approche permet de détecter les erreurs immédiatement avec un message clair, plutôt que d'avoir des comportements imprévisibles dans le parser.

3.3 Chaînes de caractères et concaténation

La gestion des chaînes a demandé plusieurs ajustements. C⁻ autorise la concaténation de deux chaînes successives ('hello' 'world' devient 'helloworld'), ce qui ne se fait pas automatiquement dans un lexer.

Pour gérer cela, on utilise un Buffer OCaml dans lequel on accumule les fragments successifs :

```

1 | '\' , {
2     Buffer.clear string_buffer;
3     string_literal lexbuf;
4     string_concat lexbuf
5 }
```

De plus, un retour à la ligne dans une chaîne doit être interdit. On a ajouté une règle spécifique qui détecte immédiatement ce cas :

```

1 and string_literal = parse
2 | '\n' {
3     raise (LexicalError
4         'Retour a la ligne interdit dans une chaine')
5 }
6 | '\\', '\' , {
7     Buffer.add_char string_buffer '\'';
8     string_literal lexbuf
9 }
10 | '\' , { STRING_CONST (Buffer.contents string_buffer) }
11 | _ as c {
12     Buffer.add_char string_buffer c;
13     string_literal lexbuf
14 }
```

3.4 Commentaires

Les commentaires /* ... */ et // ... sont pris en charge via des règles récursives. Le cas le plus délicat est celui du commentaire non fermé :

```

1 and comment = parse
2 | '*/' { token lexbuf }
3 | eof { raise (LexicalError 'Commentaire non ferme') }
4 | _ { comment lexbuf }
```

Cette vérification précoce évite que des commentaires mal formés atteignent le parser et génèrent des erreurs cryptiques.

4 Parser : analyse syntaxique

La troisième étape a été l'écriture du parser avec `ocamlyacc`. C'est l'étape qui a pris le plus de temps, notamment à cause de la gestion des précédences et des conflits shift/reduce.

4.1 Organisation de la grammaire

On a structuré la grammaire de manière hiérarchique, du plus simple au plus complexe :

1. Définition des types
2. Expressions (divisées en niveaux de précédence)
3. Déclarations
4. Instructions
5. Fonctions
6. Fichier complet

Ce découpage rend la grammaire plus lisible et limite les risques de conflits lors de la génération des tables LR.

4.2 Gestion des précédences

Une grande partie de la complexité provient des expressions, dont la précédence doit respecter exactement les règles du C. Pour éviter les conflits shift/reduce, on a défini un ensemble complet de directives de précédence :

```
1 %right ASSIGN PLUS_ASSIGN MINUS_ASSIGN ...
2 %left OR
3 %left AND
4 %left EQ NEQ
5 %left LT GT LEQ GEQ
6 %left PLUS MINUS
7 %left STAR SLASH PERCENT
8 %right NOT UMINUS UPLUS USTAR UAMPERSAND CAST
9 %nonassoc LBRACKET
```

Ces règles nous ont permis d'obtenir une grammaire compacte, sans avoir à multiplier les catégories intermédiaires.

4.3 Déclarateurs et pointeurs

Gérer les déclarateurs avec pointeurs (`int **p`) a posé un défi particulier. Pour éviter une explosion du nombre de règles, on a choisi de compter le nombre d'étoiles avant l'identifiant :

```
1 simple_declarator:
2   | IDENTIFIER { (\$1, 0) }
3   | STAR simple_declarator {
4     let (name, ptr_count) = \$2 in
5     (name, ptr_count + 1)
6 }
```

Puis on reconstruit le type final avec une fonction auxiliaire :

```
1 let rec add_pointers t n =
2   if n = 0 then t
3   else add_pointers (Pointer t) (n - 1)
```

Cette approche évite de dupliquer la logique pour chaque niveau de pointeur et reste fidèle à l'intention du langage.

4.4 Résolution du conflit if/else

Le problème classique du `if/else` ambigu se pose également en C⁻. Un `else` doit se rattacher au `if` le plus proche qui n'a pas déjà de `else`. On gère ça avec `%prec` :

```
1 %nonassoc THEN
2 %nonassoc ELSE
3
4 stmt:
5   | IF LPAREN expr RPAREN stmt %prec THEN
6     { If(\$3, \$5, None) }
7   | IF LPAREN expr RPAREN stmt ELSE stmt
8     { If(\$3, \$5, Some(\$7)) }
```

Sans ces directives, ocamllyacc génère un conflit shift/reduce que le compilateur ne peut pas résoudre automatiquement.

4.5 Le problème de l'entrelacement déclarations/fonctions

Le gros problème qu'on a rencontré lors des tests était lié à la structure du fichier. Notre grammaire initiale imposait d'abord toutes les déclarations, puis toutes les fonctions. Mais le fichier `alloc.c` des tests a une fonction, puis des déclarations, puis d'autres fonctions.

Notre première tentative :

```
1 file:
2   | declarations functions EOF { File(\$1, \$2) }
3   | functions EOF { File([], \$1) }
4   | declarations EOF { File(\$1, []) }
```

Cette grammaire ne permettait pas d'entrelacer les éléments. Après plusieurs essais infructueux, on a complètement réécrit cette partie :

```
1 top_level_item:
2   | decl      { 'Decls \$1' }
3   | func_def { 'Func \$1' }
4
5 top_level_list:
6   | /* vide */ { [] }
7   | top_level_item top_level_list { \$1 :: \$2 }
8
9 file:
10  | top_level_list EOF {
11    let partition items =
12      List.fold_right (fun item (ds, fs) ->
13        match item with
14        | 'Decls d -> (d @ ds, fs)
15        | 'Func f -> (ds, f :: fs)
16      ) items ([] , [])
17    in
18    let (decls, funcs) = partition \$1 in
19    File(decls, funcs)
20  }
```

On accumule tout dans une liste hétérogène pendant le parsing, puis on sépare les déclarations des fonctions à la fin. Cette solution respecte parfaitement la spécification qui dit explicitement : “Un fichier C⁻ contient une suite de déclarations et de définitions de fonctions, dans un ordre quelconque”.

4.6 Suppression des conflits

Au début, ocamllyacc générait plusieurs conflits shift/reduce. Pour les identifier, on a utilisé l'option -v qui génère le fichier `parser.mly` contenant tous les états de l'automate LR :

```
ocamllyacc -v parser.mly
cat parser.output
```

Chaque conflit était accompagné d'une description précise de l'état et des tokens concernés. Par exemple :

```
43: shift/reduce conflict (shift 73, reduce 47) on ELSE
state 43
stmt : IF LPAREN expr RPAREN stmt . (47)
stmt : IF LPAREN expr RPAREN stmt . ELSE stmt (48)
```

Ce fichier nous a permis d'identifier rapidement les problèmes et d'ajouter les directives `%prec` appropriées.

5 Analyse sémantique

L'analyse sémantique constitue la dernière étape du front-end. Elle vérifie deux aspects essentiels : la bonne utilisation des identifiants via une gestion correcte des portées, et le respect des règles de typage.

5.1 Table des symboles hiérarchique

On a implémenté une table des symboles avec un pointeur vers la table parente pour gérer les blocs imbriqués :

```
1 type symbol_info =
| Variable of typ
| Function of typ list * typ
2
3
4 type symbol_table = {
5   symbols: (string * symbol_info) list;
6   parent: symbol_table option;
7 }
8 }
```

Chaque entrée dans un nouveau bloc crée une nouvelle portée via `create_scope`, et les recherches d'identifiants remontent la chaîne des parents jusqu'à trouver le symbole ou atteindre la portée globale.

Cette structure nous a permis de détecter efficacement :

- Les redéclarations dans le même bloc
- Les variables masquant des fonctions
- Les identifiants non déclarés

Par exemple, la détection d'une redéclaration :

```
1 match find_symbol_current_scope tbl name with
2 | Some _ ->
3   raise (SemanticError
4     ('Variable `' ^ name ^ '\'` deja declaree'))
5 | None -> add_variable tbl name typ
```

5.2 Simplification des types

La spécification impose de ne considérer que trois types de base après l'analyse syntaxique : entier, flottant et void. On a donc défini un type intermédiaire simplifié :

```
1 type simple_type =
2   | SInt | SFloat | SVoid | SPointer of simple_type
3
4 let rec simplify_type = function
5   | BaseType (TInt, _) -> SInt
6   | BaseType (TChar, _) -> SInt
7   | BaseType (TFloat, _) -> SFloat
8   | BaseType (TDouble, _) -> SFloat
9   | Pointer t -> SPointer (simplify_type t)
```

Cette simplification rend les règles de compatibilité très explicites et faciles à contrôler.

5.3 Règles de typage strictes

La spécification est très claire : il n'y a **aucune conversion implicite** entre types. Les deux côtés d'une opération arithmétique doivent avoir exactement le même type :

```
1 match op with
2   | Add | Sub | Mul | Div ->
3     if t1 = t2 && (t1 = SInt || t1 = SFloat)
4       then t1
5     else raise (SemanticError
6                  'Types incompatibles pour operation arithmetique')
7   | Mod ->
8     if t1 = SInt && t2 = SInt then SInt
9     else raise (SemanticError
10        'Modulo necessite des types entiers')
```

Les comparaisons retournent toujours un entier, et les opérateurs logiques `&&` et `||` n'acceptent que des entiers.

Pour l'indexation de tableau, on vérifie qu'on a un pointeur à gauche et un entier à droite :

```
1 | ArrayAccess (arr, index) ->
2   match check_expr table arr with
3   | SPointer t ->
4     if check_expr table index <> SInt then
5       raise (SemanticError
6              'Index de tableau doit etre entier');
7     t
8   | _ -> raise (SemanticError
9      'Acces a un non-tableau')
```

5.4 Vérification des fonctions

Pour chaque fonction, on crée une nouvelle portée, on ajoute les paramètres, puis on vérifie le corps. On vérifie aussi que le type du `return` correspond au type déclaré :

```
1 let return_expr_type = check_expr table expr in
2 match find_symbol table 'current_function' with
3 | Some (Function (_, func_return_type)) ->
```

```

4     if not (types_compatible
5         return_expr_type func_return_type) then
6     raise (SemanticError
7         'Type de retour incompatible')

```

Cette vérification nous a permis de détecter plusieurs tests KO fournis dans l'archive du projet.

6 Débogage et tests

La phase de débogage a joué un rôle important dans l'aboutissement du projet. On a progressivement mis en place des outils pour isoler chaque type de problème.

6.1 Outils de débogage

Débogage du lexer.

Pour isoler les problèmes lexicaux, on peut tester le lexer indépendamment en affichant simplement les tokens :

```
ocamllex lexer.mll
ocamlc -c lexer.ml
```

Débogage du parser.

L'option `-v` d'ocamlyacc génère le fichier `parser.output` qui contient les états de l'automate LR et la liste des conflits. Cela a été essentiel pour corriger tous nos conflits shift/reduce.

On a également utilisé le mode debug d'OCaml :

```
export OCAMLRUNPARAM=p
```

Le parser affiche alors toutes ses actions (shift, reduce, goto), ce qui permet de comprendre exactement pourquoi un fichier est rejeté.

Débogage de l'analyse sémantique.

On a ajouté des messages d'erreur explicites dans toutes les situations : déréférencement d'un non-pointeur, type de retour incorrect, variable masquée, pointeurs incompatibles, etc.

Ces erreurs sont levées via l'exception `SemanticError`, ce qui facilite leur repérage lors de l'exécution.

6.2 Visualisation de l'AST

Pour mieux comprendre ce que produit notre parser, on a ajouté deux fonctions d'impression de l'AST :

- `string_of_file` : reconstruit le code source à partir de l'AST
- `print_ast_tree` : affiche la structure d'arbre de l'AST

La première fonction permet de vérifier que le parser a bien compris le programme. La seconde montre la vraie structure de l'AST, avec l'arborescence complète des nœuds.

Par exemple, pour un simple `if` :

```
If
Condition:
BinaryOp:
Left: Identifier: x
```

```

Right: IntConst: 3
Then:
ExprStmt
BinaryOp: =
Left: Identifier: x
Right: IntConst: 2

```

Cette visualisation nous a beaucoup aidé à comprendre comment l'AST encode réellement la structure du programme.

6.3 Tests automatisés

On a mis en place une exécution automatisée des tests via le Makefile. Les tests sont organisés en deux catégories : OK (doivent être acceptés) et KO (doivent être rejetés).

```
make test
```

Cette commande exécute tous les tests et génère des fichiers de sortie dans `../out_tests/` contenant soit l'AST complet (pour les OK), soit le message d'erreur (pour les KO).

Le Makefile utilise des boucles pour tester automatiquement tous les fichiers :

```

1 test: test-ok test-ko
2     @echo ''
3     @echo 'Tous les tests sont terminés.'
4     @echo 'Les résultats sont dans \$(OUT_TESTS)/'
5     @echo ''

```

Les tests KO ont été particulièrement utiles pour repérer des failles dans l'analyse sémantique, notamment sur la comparaison de types et la gestion des arguments de fonctions.

6.4 Résultats des tests

Au final, tous les tests OK fournis sont correctement acceptés et la majorité des tests KO sont rejettés conformément à la spécification. Les quelques tests KO qui passaient au début nous ont permis d'identifier et corriger des bugs subtils, notamment :

- L'incompatibilité entre pointeurs de types différents
- La vérification du nombre d'arguments dans les appels de fonctions
- La détection des variables utilisées avant d'être déclarées
- Le typage strict du modulo (seulement pour les entiers)

7 Conclusion

Ce projet nous a permis de comprendre en profondeur le fonctionnement d'un compilateur. Les parties les plus difficiles ont été :

- La gestion des flottants dans le lexer, avec toutes leurs formes autorisées et interdites
- L'entrelacement des déclarations et fonctions dans le parser
- Les règles strictes de typage sans conversion implicite dans l'analyse sémantique
- La résolution de tous les conflits shift/reduce du parser

Au final, on a obtenu un analyseur complet et robuste qui respecte toute la spécification de C⁻. Le code est bien structuré, commenté, et accompagné d'outils de test et de visualisation qui facilitent le débogage.

Les choix de conception qu'on a faits (séparation des opérateurs, table des symboles hiérarchique, simplification des types) se sont révélés judicieux et nous ont permis d'implémenter les fonctionnalités demandées de manière claire et maintenable.