

# Projet Architecture des Ordinateurs

## Machine à Pile

Théophile Dano, Vincent Gouteux, Adrien Legros

31 janvier 2017

# Première partie

## Le compilateur assembleur

Liste des fichiers :

```
ast.c
ast.h
compiler.h
grammar.y
makefile
scanner.l
symtable.c
symtable.h
```

### 1 Approche générale

Dans la manière de concevoir ce petit compilateur, nous avons voulu utiliser les manières plus populaires pour le faire. C'est-à-dire créer un compilateur en plusieurs passes, qui toutes ont leur propres entrées/sorties et qui travaillent à la chaîne.

L'idée est que le travail du compilateur est divisé en 4 phases :

- L'analyse lexical (le "lexing")
- L'analyse grammaticale (le "parsing")
- L'optimisation
- La génération de code

Nous avons décidé d'ignorer la phase d'optimisation du code puisqu'elle n'est pas du tout requise pour ce projet, même si cela pourrait être intéressant.

#### 1.1 Analyse Lexicale

Le rôle du *lexer* est de générer depuis une chaîne de caractères une liste de *jetons* ("token" dans la suite) pour identifier chacune des suites de caractères qu'il reconnaît.

Le but du lexer est donc de découper ce qui lui est donné en une liste de tokens pour ensuite pouvoir en comprendre la structure.

Regardons sur un petit exemple. Disons que l'on veuille analyser une chaîne de caractères qui représente une phrase en français, et que notre lexer ne reconnaisse que les *VERBEs* et les *NOMs* et que tout autre mot soit identifié par un token de type *MOT*.

*Le chat mange la souris*

Si nous devons exécuter notre lexer sur cette phrase test, il nous reverrait la liste de tokens :

- MOT : "Le"
- NOM : "chat"
- VERBE : "mange"
- MOT : "la"
- NOM : "souris"

A ce stade le lexer a fini son travail. Il a atteint la fin de l'entrée et n'a rencontré aucune erreur. Il passe donc la main au parser.

## 1.2 Analyse grammaticale

**Grammaire** L'enjeu d'un parser est, étant donné une chaîne de tokens et leur contenu, vérifier qu'ils forment une structure conforme à une grammaire définie à l'avance. Ici il faut entendre le mot grammaire au sens classique du terme, c'est-à-dire une liste de règles définissant comment on peut organiser des mots d'un langage.

Si maintenant on introduit pour notre lexer le token ARTICLE alors les mots "le" et "la" seront reconnus comme tels. Cela va nous aider à définir une grammaire extrêmement simple avec une unique règle PHRASE-SIMPLE.

Le lexing de notre première phrase se transforme alors en :

- ARTICLE : "Le"
- NOM : "chat"
- VERBE : "mange"
- ARTICLE : "la"
- NOM : "souris"

En suivant la notation EBNF (*Extended Backus-Naur Form*) pour les grammaires, on peut définir PHRASE-SIMPLE comme :

```
PHRASE-SIMPLE
:
| ARTICLE NOM ARTICLE NOM
;
```

On remarque immédiatement que n'importe quelle phrase basée sur cette structure grammaticale sera reconnue et validée par le parser. (par exemple "Les taupes creusent le sol" est une entrée valide).

C'est aussi le rôle du parser de faire remonter les possibles erreurs de grammaires contenues dans le programme.

**Représentation interne du code** Il est évident que l'on ne veut pas juste vérifier que la structure du code soit correcte mais bien en faire quelque chose. Par exemple générer du code correspondant dans un autre langage, ou colorer les mots clés dans un éditeur etc.

Pour cela, il nous faut garder en mémoire chaque groupe de tokens qui a été reconnu et validé par le parser dans une structure de données appropriée. Les arbres sont pour cela bien adaptés puisqu'ils permettent d'organiser les informations de manière très pratique. On peut représenter une PHRASE-SIMPLE par un nœud ayant pour fils chacun des tokens repérés dans l'ordre.

De manière générale, on utilise un *A.S.T.* ("*Abstract Syntax Tree*" où Arbre de Syntaxe Abstraite).

### 1.3 Génération du code

Cette dernière passe d'un compilateur n'est généralement pas compliquée. En effet, une approche simple et efficace est de parcourir simplement l'arbre (AST) généré par le parser et pour chaque nœud de traverser ses fils récursivement.

Il y a évidemment des problèmes qui se posent comme la gestion des identifiants (noms de variables etc.). Sont-ils utilisés dans un cadre correct ? Ont-ils bien été définis avant ? Si oui, sont-ils du type attendu et ainsi de suite.

Ces problèmes se résolvent souvent facilement par l'inspection d'une table de symboles ("*Symtable*" en anglais). Cette table (souvent une *Hash Map*) répertorie chaque noms de symboles définis dans un programme et leur associe des informations connexes (type, mutabilité, définition, etc.).

## 2 Application au projet

Dans cette partie nous allons voir dans quelle mesure il est possible d'appliquer ces techniques au projet en C à l'aide de bibliothèques et d'un *parser-generator* : le tandem *Flex/Bison*, version retravaillée de *Lex/Yacc*.

### 2.1 Lexer – Flex

Pour la première passe de notre compilateur, le lexing, nous avons choisi d'utiliser Flex un générateur de lexer simple à prendre en main, souvent utilisé dans les cours d'introductions à la compilation et à l'analyse de textes divers.

Nous n'allons pas détailler le fonctionnement interne de Flex mais simplement analyser quelques portions du fichier **scanner.l**.

Dans la première partie, nous avons détaillé comment on pouvait définir certains types de tokens. Il s'agit ici (avec Flex) de le faire avec du code C. Regardons le plus simple des tokens du projet : le mot clé **pop**.

```
| pop { return TOKEN_POP; }
```

La syntaxe utilisée n'est pas du tout pensée pour être du C pur. Seul les parties entre accolades doivent être du code C. Ici, le mot **pop** est utilisé pour désigner un *pattern* à reconnaître et l'action associée est tout simplement de renvoyer au parser l'information que l'on est tombé sur un token de type POP.

Utiliser Flex est d'une trivialité sans pareil. Cependant, pour identifier un type particulier d'entrée qui n'est pas défini clairement au préalable, il faut utiliser un outil obscur au premier regard, mais qui est très logique.

Les expressions régulières ("*regex*" de "*regular expressions*") sont très utilisées en informatique pour décrire des *patterns* et reconnaître des formations syntaxiques particulières. On appelle cela plus couramment le *pattern matching*, dans le sens où l'on essaie de reconnaître des patterns particuliers.

Prenons l'exemple des étiquettes (labels) d'instructions dans notre petit langage assembleur. On ne sait pas donner une liste de mots étant à reconnaître comme des étiquettes et donc on ne peut pas lister tous les labels possibles. Cependant, on connaît la manière dont sont définis les labels. Ce sont des mots

(combinaisons de caractères ASCII suivis de chiffres ou d'underscores) qui ne sont pas des mots-clés et qui sont suivis de deux points.

On peut alors définir ça comme suit à l'aide des expressions régulières et Flex :

```
[a-zA-Z_][a-zA-Z0-9_]*" : " {
    yytext[strlen(yytext) - 1] = 0; // pour ignorer les ":"
    yylval.ident = strdup(yytext);

    return TOKEN_LBL;
}
```

Sans rentrer dans les détails des expressions régulières, la règle utilisée ici est très simple. La règle `[a-zA-Z_]` précise que le premier caractère doit être une lettre peu importe sa casse, ou bien un underscore. Ensuite, `[a-zA-Z0-9_]*` permet de récupérer n'importe quel caractère alphanumérique ou underscore. L'étoile est utilisée pour indiquer que l'on accepte une répétition sans limite de cette règle. Enfin, `" : "` indique que le token doit se terminer par deux points.

Finalement, si la règle est "matchée", on exécute quelques actions et on retourne le token au parser qui saura quoi en faire.

## 2.2 Parser – Bison

Maintenant que notre fichier ou code source a été relu et découpé par le lexer, le parser peut commencer son travail. C'est-à-dire construire l'AST et vérifier que l'entrée soit conforme à la grammaire utilisée.

Ici, nous n'allons pas nous éterniser sur le pourquoi du comment des grammaires mais il est important de noter qu'il faut faire attention aux récursions des règles (une règle s'invoque elle-même, par la droite ou la gauche par exemple) car cela peut mener à des boucles infinies et planter le programme.

Cela étant dit, regardons comment on peut mettre en place une grammaire et son parser avec Bison. Nous allons regarder qu'un seul exemple pour pouvoir ensuite parler de la génération de l'AST. Dans le fichier **grammar.y** on trouve une liste de règles de grammaire parmi lesquelles celle qui définit comment un programme assembleur doit être organisé.

```
program
: %empty
| stat NL program
| TOKEN_LBL {
    tok = $1;
    if (yydolog) fprintf(yylog, "%3d: _label_ -> %s\n", pc, $1);
    add_sym_entry(labels, $1, pc);
} stat NL program
| NL program
;
```

Comme pour Flex, les commandes en langage C sont définies entre accolades. Ici, NL correspond à un retour à la ligne, stat est n'importe quelle instruction assembleur valide, et TOKEN\_LBL le token associé aux labels (vu plus haut).

La règle **%empty** signifie qu'un programme valide peut ne rien être (chaîne vide). Sinon (la barre verticale signifie une alternative), un label suivi d'une instruction et d'un retour à la ligne est accepté, une instruction simple et un retour à la ligne de même.

On peut observer que la règle **program** est récursive (par la droite). C'est-à-dire qu'une suite de **program** est un **program** (eux mêmes constitués d'instruction et d'étiquettes).

Regardons maintenant l'action faite lorsque l'on rencontre un label. On récupère le token dans la variable **tok** grâce au symbole **\$1** qui sera remplacé par Bison par le pointeur adapté. Ensuite on ajoute le label à la Symtable (table des symboles) en indiquant comme information connexe son adresse dans le code.

En jetant un rapide coup d'œil à la règle **stat**,

```
stat
: TOKEN_POP INT {
    tok = $1;
    if (yydolog) fprintf(yylog, "%3d: _pop_%d\n", pc, $2);
    ast_tail->next = new_instr(POP, $2, pc++);
    ast_tail = ast_tail->next;
| ...
;
}
```

On voit comment est gérée la génération de l'AST par le parser. Lorsqu'une instruction correcte est rencontrée, on génère un nœud dans l'arbre qui lui correspond. Ici, l'instruction **pop** est valide si le token POP est suivi du token INT (pour un nombre entier). Dans ce cas, on génère un simple nœud instruction avec son adresse dans le code et son argument (le nombre entier).

## 2.3 AST et génération de code

Dans cette dernière partie liée au compilateur, nous allons nous intéresser à l'AST et la génération de code pour la machine virtuelle (à pile) qui sera discutée dans la prochaine partie.

Pour cela, intéressons-nous à la structure utilisée pour les nœuds de l'AST.

```
struct node {
    int lno;
    int sym;
    union {
        int arg;
        char *name;
    };
    struct node *next;
};
```

On remarque immédiatement qu'il s'agit en réalité d'une liste chaînée et non d'un arbre (même si une liste chaînée est un arbre ...). En effet, comme le langage assembleur utilisé a une syntaxe très linéaire, il n'y a aucune structure de code qui soit assez complexe pour justifier l'utilisation de plusieurs fils par nœuds (si l'on avait pensé le langage en acceptant de simples expressions, on aurait dû modifier cette structure).

Le membre **int** *lno* correspond à la ligne où l'instruction apparaît. **int** *sym* dénote du type de l'instruction (pop, push, call, etc.). L'**union** est utilisée pour donner l'information de soit l'argument numérique de l'instruction, ou l'identifiant (par exemple dans **jmp** *label*).

**Pour construire l'AST** On utilise ensuite les trois fonctions déclarées dans le fichier **ast.h** :

```
struct node *new_ctrl(int sym, char *to, int lno);
struct node *new_instr(int sym, int arg, int lno);
struct node *new_ph(int sym, char *name, int lno);
```

qui servent respectivement à :

- Ajouter un nœud de contrôle (**jmp**, **call** ...)
- Ajouter un nœud pour une instruction classique (**pop** 10, ...)
- Ajouter un "*placeholder*" pour signaler la position d'un label (ignoré à l'étape de la génération de code, mais utilisé pour repérer un label dans le code).

Deuxième partie

## La machine a pile



Troisième partie

## Exemples d'exécution