

Béatrice Creusillet <bcreusillet@quarkslab.com>
Philippe Viroulet <pviroulet@quarkslab.com>

Compilation pour le reverseur

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Ma première passe

Pré-requis

- ▶ `docker pull viroulep/tp-obfu:latest`
- ▶ TDs: `git clone git@github.com:bac4tout/compil.git`
- ▶ `docker run -it -v /path/to/sources:/home/tp/sources -i viroulep/tp-obfu /bin/bash`

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

Quelques idées folles ... ou pas

Ma première passe

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

- Compilation pour le reverseur, vraiment ?

- Compilation vs. interprétation

- Compilation statique : du code source au binaire

- Options de compilation

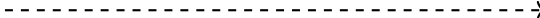
- Différents compilateurs, différents binaires

- Un peu de mécanique : les internes d'un compilateur

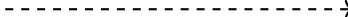
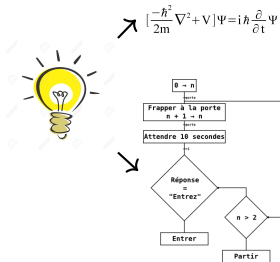
- Quelques idées folles ... ou pas

Ma première passe

Comment peut-on être binaire ?

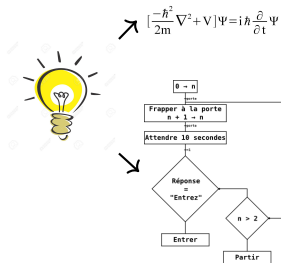


Comment peut-on être binaire ?



```
0101010110000111100000
0101011010100010101000
011000011100001000001111
101101010101110001011001
0001111010111101010111
010101110001101000101111
0001011101001010100000
110110001011010101111111
110110010110101010111111
110001000001000101111111
101001100001010110001111
1011110111111111000000
100101100111001100010111
100101100111001100010111
110101011010111101010000
000100011011010101100000
010111111000101000110001
000101110011111001101000
100111000101111010011111
```

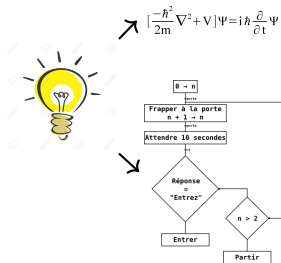
Comment peut-on être binaire ?



Code

```
0101010110000111100000  
01010110101000101101000  
11100101110001000001111  
101101010101110001011001  
000111101011110101011  
01010111000110100010111  
0001011101001010110100  
11010001011010010111011  
11010001000100101011111  
1010010000101011000101110  
1011110111111111000100  
1001011001110111000100110  
110101010101111010100000  
00010011011010101100000  
0101111100010100110100  
000101110011111001101000  
10011000101111010001110
```


Comment peut-on être binaire ?



→ Code

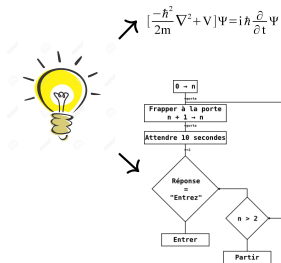
→ Traduction

→

```
01101011011000011111000000
10110111011010001011101000
11110010111000010000011111
10110110110111000101100111
00011111011011111101011011
01011011110001110100010111
10001011110100101101101100
11011000101110100101111011
11110010000010101101111111
1010011000010101110001011110
11011111011111111110001000
10010111001111011100010111
110110110110111110101100000
00010011111101010111100000
01011111100101100111110111
10010111101111110011101000
10011100010111110100011111
```

Comment peut-on être binaire ?

Spécialisation



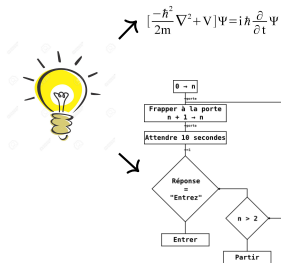
Code

Traduction



Comment peut-on être binaire ?

Spécialisation



Code

Traduction



Reversing
Généralisation

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

Quelques idées folles ... ou pas

Ma première passe

Un gradient d'approches

- ▶ Compilation : traduction statique en binaire
 - ▶ réutilisable
 - ▶ pour des inputs variés
 - ▶ adapté à un environnement particulier (windows, linux, osx,...)
 - ▶ ex: C, C++, Fortran

Un gradient d'approches

- ▶ Compilation : traduction statique en binaire
 - ▶ réutilisable
 - ▶ pour des inputs variés
 - ▶ adapté à un environnement particulier (windows, linux, osx,...)
 - ▶ ex: C, C++, Fortran

- ▶ Interprétation : traduction dynamique = à la volée
 - ▶ non-réutilisable
 - ▶ optimisée pour des inputs particuliers
 - ▶ ex: Basic, R, bash

Un gradient d'approches

- ▶ Compilation : traduction statique en binaire
 - ▶ réutilisable
 - ▶ pour des inputs variés
 - ▶ adapté à un environnement particulier (windows, linux, osx,...)
 - ▶ ex: C, C++, Fortran
- ▶ Pré-compilation : traduction statique vers un 'byte-code'
 - ▶ byte-code générique
 - ▶ adaptation à l'environnement au run-time
 - ▶ ex: Java
- ▶ Interprétation : traduction dynamique = à la volée
 - ▶ non-réutilisable
 - ▶ optimisée pour des inputs particuliers
 - ▶ ex: Basic, R, bash

Un gradient d'approches

- ▶ Compilation : traduction statique en binaire
 - ▶ réutilisable
 - ▶ pour des inputs variés
 - ▶ adapté à un environnement particulier (windows, linux, osx,...)
 - ▶ ex: C, C++, Fortran
- ▶ Pré-compilation : traduction statique vers un 'byte-code'
 - ▶ byte-code générique
 - ▶ adaptation à l'environnement au run-time
 - ▶ ex: Java
- ▶ JIT : Just-in-Time compilation
 - ▶ traduction dynamique partiellement 'réutilisable'
 - ▶ parfois limitée aux hot spots
 - ▶ ex: Matlab, Julia, Java
- ▶ Interprétation : traduction dynamique = à la volée
 - ▶ non-réutilisable
 - ▶ optimisée pour des inputs particuliers
 - ▶ ex: Basic, R, bash

Compilateurs commerciaux

- ▶ Intel C/C++ compiler (icc) (Windows, linux, MacOS, Intel)
- ▶ Visual Studio C++ (Windows, Microsoft)
- ▶ IAR C/C++ compiler (Windows, Linux, others, IAR Systems)
- ▶ Edison Design Group (Windows, Linux, others, EDG)
- ▶ VisualAge C++, XL C/C++ (Windows, Linux, Aix, OS/2, OS/400,... IBM)
- ▶ ...

OpenSource ou Freeware

- ▶ GNU C/C++ (**gcc**, MinGW) (Windows, linux, MacOS, GNU Project)
- ▶ Clang/LLVM (**clang**) (Windows, linux, MacOS, LLVM Project)
- ▶ ...

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

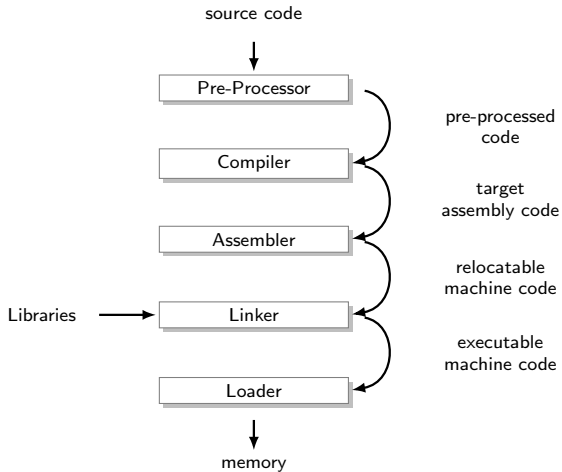
Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

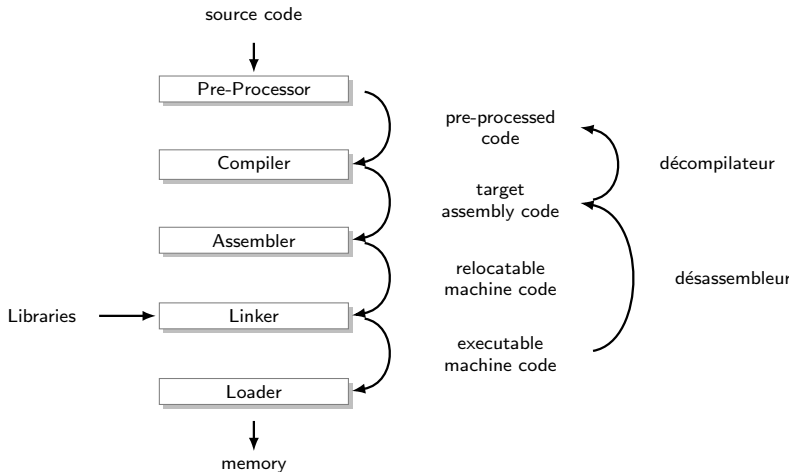
Quelques idées folles ... ou pas

Ma première passe

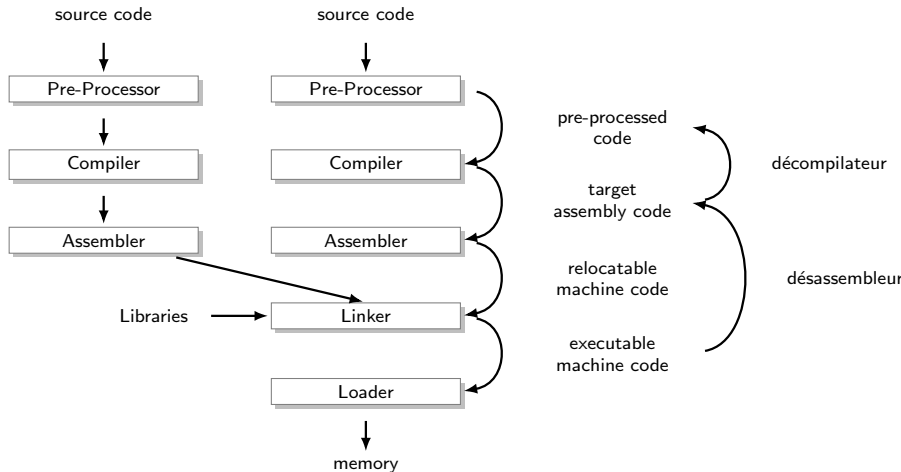
Une traduction en plusieurs étapes



Une traduction en plusieurs étapes



Une traduction en plusieurs étapes



TD 1 : les étapes de la génération de binaire

welcome.c :

```
1  #include <stdio.h>
2  void welcome(char const* nom) {
3      printf("welcome %s!\n", nom);
4  }
```

main_welcome.c :

```
1  #ifndef DEFAULT
2  #define DEFAULT "everybody"
3  #endif
4
5  extern void welcome(char const*);
6
7  int main(int argc, char const* argv[]) {
8      if      (argc == 1) { welcome(DEFAULT); }
9      else if (argc == 2) { welcome(argv[1]); }
10     else { return 1; }
11     return 0;
12 }
```

td1 : exécution

Version par défaut :

```
1 clang welcome.c main_welcome.c
2 ./a.out
3 ./a.out alice
```

td1 : exécution

Version par défaut :

```
1  clang welcome.c main_welcome.c
2  ./a.out
3  ./a.out alice
```

```
1  welcome everybody!
2  welcome alice!
```


td1 : exécution

Version par défaut :

```
1 clang welcome.c main_welcome.c
2 ./a.out
3 ./a.out alice
```

```
1 welcome everybody!
2 welcome alice!
```

Version avec définition de variable de pré-processing :

```
1 clang welcome.c main_welcome.c -DDEFAULT=' "bob" '
2 ./a.out
3 ./a.out alice
```

td1 : exécution

Version par défaut :

```
1 clang welcome.c main_welcome.c
2 ./a.out
3 ./a.out alice
```

```
1 welcome everybody!
2 welcome alice!
```

Version avec définition de variable de pré-processing :

```
1 clang welcome.c main_welcome.c -DDEFAULT=' "bob" '
2 ./a.out
3 ./a.out alice
```

```
1 welcome bob!
2 welcome alice!
```

Le pré-processeur : traitement des macros

```
1 clang -E main_welcome.c
```

Le pré-processeur : traitement des macros

```
1 clang -E main_welcome.c
```

```
1 int main(int argc, char const* argv[]) {  
2     if (argc == 1) { welcome("everybody"); }  
3     else if (argc == 2) { welcome(argv[1]); }  
4     else { return 1; }  
5     return 0;  
6 }
```

Le pré-processeur : traitement des macros

```
1 clang -E main_welcome.c
```

```
1 int main(int argc, char const* argv[]) {  
2     if (argc == 1) { welcome("everybody"); }  
3     else if (argc == 2) { welcome(argv[1]); }  
4     else { return 1; }  
5     return 0;  
6 }
```

```
1 clang -E main_welcome.c -DDEFAULT="'bob'"
```

Le pré-processeur : traitement des macros

```
1 clang -E main_welcome.c
```

```
1 int main(int argc, char const* argv[]) {  
2     if (argc == 1) { welcome("everybody"); }  
3     else if (argc == 2) { welcome(argv[1]); }  
4     else { return 1; }  
5     return 0;  
6 }
```

```
1 clang -E main_welcome.c -DDEFAULT='"bob"'
```

```
1 int main(int argc, char const* argv[]) {  
2     if (argc == 1) { welcome("bob"); }  
3     else if (argc == 2) { welcome(argv[1]); }  
4     else { return 1; }  
5     return 0;  
6 }
```

Le pré-processeur : inclusion de fichiers *header*

```
1 clang -E welcome.c | wc -l
```

```
1 797
```

Le pré-processeur : inclusion de fichiers *header*

```
1 clang -E welcome.c | wc -l
```

```
1 797
```

```
1 clang -E welcome.c
```


Le pré-processeur : inclusion de fichiers *header*

```
1 clang -E welcome.c | wc -l
```

```
1 797
```

```
1 clang -E welcome.c
```

- ▶ les fichiers sont intégralement inclus, récursivement
- ▶ attention aux include redondants

```
1 #ifndef MY_HEADER_H_  
2 #define MY_HEADER_H_  
3 .... // header content  
4 #endif // MY_HEADER_H_
```

- ▶ attention à l'expansion de macros !

Le compilateur : génération de code assembleur

```
1 clang -S main_welcome.c
2 more main_welcome.s
```

L'assemblage: génération de code objet

```
1  as main_welcome.s -o main_welcome.o
2  file main_welcome.o
3  nm main_welcome.o
```

L'assemblage: génération de code objet

```
1  as main_welcome.s -o main_welcome.o
2  file main_welcome.o
3  nm main_welcome.o
```

```
1  main_welcome.o: ELF 64-bit LSB relocatable, x86-64, version 1 (
    SYSV), not stripped
2
3          U _GLOBAL_OFFSET_TABLE_
4  0000000000000000 T main
5          U welcome
```

L'assemblage: génération de code objet

```
1  as main_welcome.s -o main_welcome.o
2  file main_welcome.o
3  nm main_welcome.o
```

```
1  main_welcome.o: ELF 64-bit LSB relocatable, x86-64, version 1 (
    SYSV), not stripped
2
3              U _GLOBAL_OFFSET_TABLE_
4  0000000000000000 T main
5              U welcome
```

```
1  clang -c welcome.c -o welcome.o
2  file welcome.o
3  nm welcome.o
```

L'assemblage: génération de code objet

```
1  as main_welcome.s -o main_welcome.o
2  file main_welcome.o
3  nm main_welcome.o
```

```
1  main_welcome.o: ELF 64-bit LSB relocatable, x86-64, version 1 (
    SYSV), not stripped
2
3          U _GLOBAL_OFFSET_TABLE_
4  0000000000000000 T main
5          U welcome
```

```
1  clang -c welcome.c -o welcome.o
2  file welcome.o
3  nm welcome.o
```

```
1  welcome.o: ELF 64-bit LSB relocatable, x86-64, version 1 (SYSV)
    , not stripped
2
3          U _GLOBAL_OFFSET_TABLE_
4          U printf
5  0000000000000000 T welcome
```

L'édition de lien: génération de binaire

```
1  ld welcome.o main_welcome.o
```

L'édition de lien: génération de binaire

```
1  ld welcome.o main_welcome.o
```

```
1  ld : avertissement: le symbole d'entree _start est introuvable  
    ; utilise par défaut 00000000004000b0  
2  welcome.o: Dans la fonction « welcome » :  
3  welcome.c:(.text+0x20) : reference indefinie vers « printf »
```


L'édition de lien: génération de binaire

```
1  ld welcome.o main_welcome.o
```

```
1  ld : avertissement: le symbole d'entree _start est introuvable  
    ; utilise par défaut 00000000004000b0  
2  welcome.o: Dans la fonction « welcome » :  
3  welcome.c:(.text+0x20) : reference indefinie vers « printf »
```

```
1  ld welcome.o main_welcome.o -lc  
2  file a.out
```

L'édition de lien: génération de binaire

```
1 ld welcome.o main_welcome.o
```

```
1 ld : avertissement: le symbole d'entree _start est introuvable  
    ; utilise par défaut 00000000004000b0  
2 welcome.o: Dans la fonction « welcome » :  
3 welcome.c:(.text+0x20) : reference indefinie vers « printf »
```

```
1 ld welcome.o main_welcome.o -lc  
2 file a.out
```

```
1  
2 ld: avertissement: le symbole d'entree _start est introuvable;  
    utilise par défaut 00000000004002c0  
3  
4 a.out: ELF 64-bit LSB executable, x86-64, version 1 (SYSV),  
    dynamically linked, interpreter /lib/ld64.so.1, not  
    stripped
```

L'édition de lien: génération de binaire

```
1 clang -v welcome.o main_welcome.o 2>&1 | grep ld
```

L'édition de lien: génération de binaire

```
1 clang -v welcome.o main_welcome.o 2>&1 | grep ld
```

```
1  "/usr/bin/ld" -z relro --hash-style=gnu --eh-frame-hdr -m
    elf_x86_64 -dynamic-linker /lib64/ld-linux-x86-64.so.2 -o
    a.out /usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/../../x86_64-linux-gnu/crt1.o /usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/../../x86_64-linux-gnu/crti.o /usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/crtbegin.o -L/usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0 -L/usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/../../x86_64-linux-gnu -L/lib/../../lib64 -L/usr/lib/x86_64-linux-gnu -L/usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/../../x86_64-linux-gnu -L/usr/lib/llvm-6.0/bin/../../lib -L/lib -L/usr/lib welcome.o
    main_welcome.o -lgcc --as-needed -lgcc_s --no-as-needed -lc -lgcc --as-needed -lgcc_s --no-as-needed /usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/crtend.o /usr/bin/../../lib/gcc/x86_64-linux-gnu/7.3.0/../../x86_64-linux-gnu/crtn.o
```

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

Quelques idées folles ... ou pas

Ma première passe

Piloter son compilateur : les options

Contrôle des inputs/outputs

- ▶ fichiers sources : paramètres positionnels
- ▶ bibliothèques : `-l<libname>`
- ▶ chemins d'accès
 - ▶ fichiers headers : `-I<path>`
 - ▶ bibliothèques : `-L<path>`
- ▶ standard du langage : `-std=c++11` par ex.
- ▶ warnings : `-W<warning>` (ex: `-Wall`, `-Werror`)
- ▶ infos : `--help`, `--version`, `-print-search-dirs`, ...

Passer des options aux différentes étapes

- ▶ pré-processeur: `-Wp,<option>`
- ▶ assembleur: `-Wa,<option>`
- ▶ linker: `-Wl,<option>`

Liste des optimisations :

► gcc/g++:

```
1      g++ --help=optimizers
```

► clang/llvm:

```
1      clang -help  
2      opt -help
```

Optimisations globales : -O<level>

- -O0 : pas d'optimisations, pratique pour le debug
- -O3 : plus fort niveau d'optimisations
- -Os : -O2, mais minimise la taille du code
- -Ofast: -O3 + -ffastmath

TD 2 : Piloter l'optimiseur

```
1  #ifdef BIG_N
2  #define n 100
3  #else
4  #define n 5
5  #endif
6
7  int main() {
8      int a[n];
9      for (int i = 0; i < n; ++i) { a[i] = i+1; }
10
11     int s = 0;
12     for (int i = 0; i < n; ++i) { s +=a [i]; }
13
14     if ( s== (n*(n+1))/2 )      { return s; }
15     return -1;
16 }
```

- ▶ avec *#define n 5*
Comparer les assembleurs générés avec -O0 et -O2
- ▶ avec *#define n 100*
Comparer les binaires générés avec -O0 et -O3

TD 3 : les options pour le debug (1)

```
1  #include <assert.h>
2  #define MAX_SIZE 512
3
4  int count_char(const char* str) {
5      int c = 0;
6      while ( (c < MAX_SIZE) && (str[c] != '\0')) { ++c; }
7      return c;
8  }
9
10 int main(int argc, char** argv) {
11     assert(argc == 2 && "compliant number of arguments");
12     return count_char(argv[1]);
13 }
```

```
1  clang debug.c -o debug
2  ./debug toto
3  echo $?
4  ./debug
```

TD 3 : les options pour le debug (2)

L'impact de `-DNDEBUG`

```
1 clang debug.c -o debug-n -DNDEBUG
2 ./debug-n toto
3 echo $?
4 ./debug-n
```

Comparer les binaires générés avec et sans `-DNDEBUG`

L'impact de `-g`

```
1 clang debug.c -o debug-ng -DNDEBUG -g
```

Comparer les binaires générés avec et sans `-g`

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

Quelques idées folles ... ou pas

Ma première passe

TD 4 : clang vs. gcc

Reprendre le programme `array.c` du TD 2

Comparer les assembleurs générés par gcc et clang pour `n = 100` et différents niveaux d'option `-O<level>` :

`td4/gcc_vs_clang.sh` :

```
1  LEVEL=$1
2
3  # assembleur
4  clang array.c -DBIG_N -O$LEVEL -S -o array_clang_O$LEVEL.s
5  gcc array.c -DBIG_N -O$LEVEL -S -o array_gcc_O$LEVEL.s
6
7  # binaire
8  clang array.c -DBIG_N -O$LEVEL -o array_clang_O$LEVEL
9  gcc array.c -DBIG_N -O$LEVEL -o array_gcc_O$LEVEL
```

```
1  bash gcc_vs_clang.sh 1
2  bash gcc_vs_clang.sh 2
3  bash gcc_vs_clang.sh 3
```

Comparer aussi les versions générées par clang avec `-O2` et gcc avec `-O3`.

Jouer avec les compilateurs

<https://godbolt.org/>

Outil permettant de visionner les assembleurs générés par différents compilateurs pour différentes plateformes (x86_64, armxx,...)
(développé par Matt Godbolt)

<https://cppinsights.io/>

Outil basé sur Clang, qui permet de voir comment le front-end transforme le code
(développé par Andreas Fertig)

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

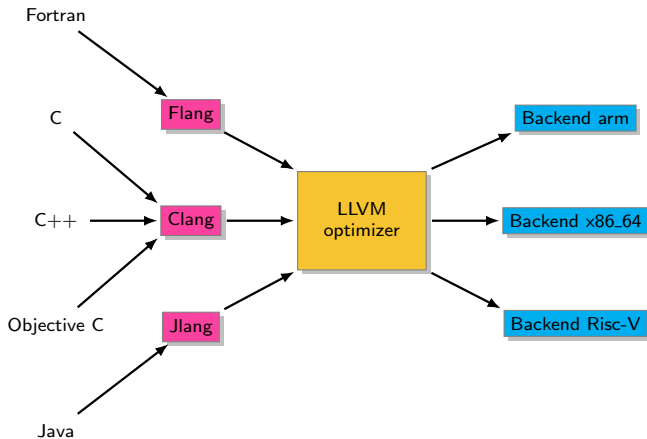
Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

Quelques idées folles ... ou pas

Ma première passe

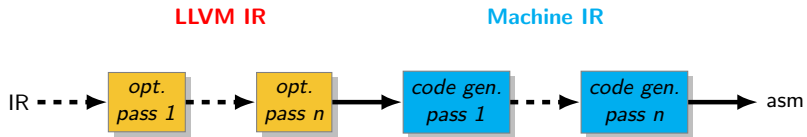
Clang/LLVM



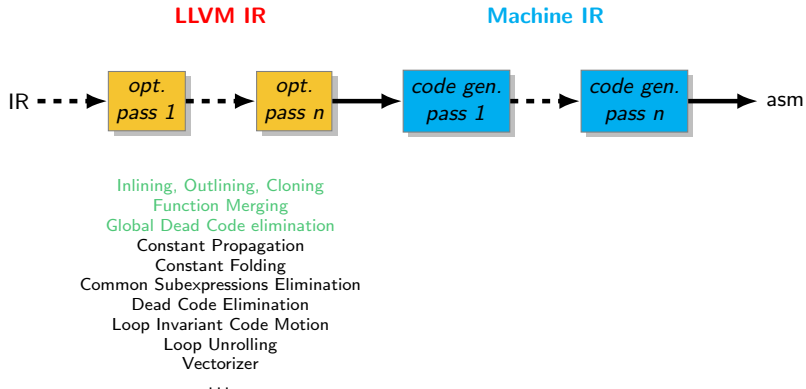
LLVM : organisation en passes



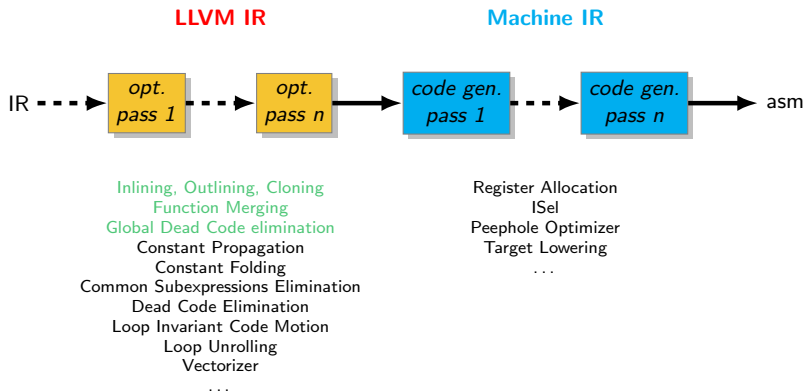
LLVM : organisation en passes



LLVM : organisation en passes



LLVM : organisation en passes



<https://llvm.org/docs/Passes.html>

<https://github.com/llvm-mirror/llvm/blob/master/lib/Transforms/IPO/PassManagerBuilder.cpp>

Transformations interprocédurales

- ▶ Inlining, outlining, merge de fonctions, cloning, global dce, . . .
- ▶ sur le Call Graph
- ▶ sur les variables globales

Transformations intraprocédurales

- ▶ Constant propagation, Common subexpressions elimination, loop transformations, ...
- ▶ nécessite de pouvoir déplacer les instructions les unes par rapport aux autres
- ▶ en respectant leurs dépendances :
 - ▶ **spatiales**
(est-ce que je touche la même zone mémoire ou le même registre?)
 - ▶ **temporelles**
(est-ce que ce store est exécuté avant de load?)

LLVM apporte ces informations sous différentes formes:

- ▶ sa représentation interne (la LLVM IR)
 - ▶ Basic blocks et CFG
 - ▶ forme SSA
transforme les dépendances de registres en chaînes de valeurs (use-defs chains)
attention : ne prend pas en compte les dépendances de mémoire
- ▶ des analyses que les passes de transformation peuvent requérir
 - ▶ AliasAnalysis (AA)
 - ▶ MemorySSA
 - ▶ ...

Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Compilation pour le reverseur, vraiment ?

Compilation vs. interprétation

Compilation statique : du code source au binaire

Options de compilation

Différents compilateurs, différents binaires

Un peu de mécanique : les internes d'un compilateur

Quelques idées folles ... ou pas

Ma première passe

Deux idées folles ? 1. Compiler pour reverser

Désassembleur/Décompileur = Compilateur ?

- ▶ Concepts communs (transformer du code, RI)
- ▶ Outils communs
 - ▶ CFG, CG
 - ▶ Analyses de flot de valeurs, forme SSA

Utiliser une RI commune (par exemple la LLVM IR) ?

- ▶ RetDec, McSema/ReMill
- ▶ Pour capitaliser l'existant
- ▶ Pour bénéficier d'outils d'analyse / transformation externes
 - ▶ ex: Klee, Souper
 - ▶ cf. Saturn (SPro'19)
- ▶ Pas forcément toujours souhaitable / possible
 - ▶ avoir une IR multi-level (= MLIR ;p)

Deux idées folles ? 2. Compiler pour obfusquer

Deux idées folles ? 2. Compiler pour obfusquer

Obfusquer = Transformer = Compiler

Obfusquer = Complexifier mais Compiler = Optimiser



Deux idées folles ? 2. Compiler pour obfusquer

Obfusquer = Transformer = Compiler

Obfusquer = Complexifier mais Compiler = Optimiser

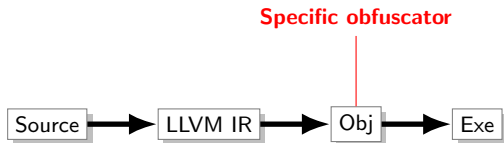
Specific obfuscator



Deux idées folles ? 2. Compiler pour obfusquer

Obfusquer = Transformer = Compiler

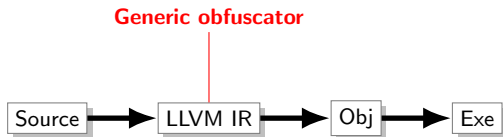
Obfusquer = Complexifier mais Compiler = Optimiser



Deux idées folles ? 2. Compiler pour obfusquer

Obfusquer = Transformer = Compiler

Obfusquer = Complexifier mais Compiler = Optimiser



Deux idées folles ? 2. Compiler pour obfusquer

Obfusquer = Transformer = Compiler

Obfusquer = Complexifier mais Compiler = Optimiser

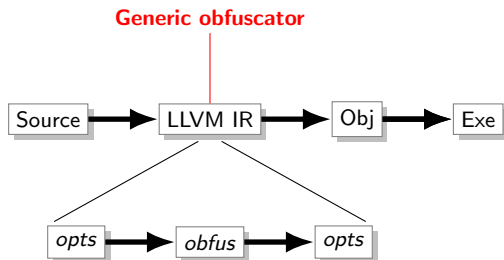


Table of Contents

Avant de rentrer dans le vif du sujet

Introduction générale

Ma première passe