# From Source Code to Executable

## Dr. Axel Kohlmeyer

Research Professor, Department of Chemistry
Associate Director, ICMS
College of Science and Technology
Temple University, Philadelphia
**axel.kohlmeyer@temple.edu**

# Pre-process / Compile / Link

- Creating an executable includes multiple steps
- The "compiler" (gcc) is a wrapper for <u>several</u> commands that are executed in succession
- The "compiler flags" similarly fall into categories and are handed down to the respective tools
- The "wrapper" selects the compiler language from source file name, but links "its" runtime
- We will look into a C example first, since this is the language the OS is (mostly) written in

# A simple C Example

- Consider the minimal C program 'hello.c':
  ```c
  #include <stdio.h>
  int main(int argc, char **argv)
  {
          printf("hello world\n");
          return 0;
  }
  ```

- i.e.: what happens, if we do:
  ```
  > gcc -o hello hello.c
  ```
  (try: `gcc -v -o hello hello.c`)

**3**

# Step 1: Pre-processing

- Pre-processing is <u>mandatory</u> in C (and C++)
- Pre-processing will handle '#' directives
    - File inclusion with support for nested inclusion
    - Conditional compilation and Macro expansion
- In this case: **/usr/include/stdio.h**
  - and all files are included by it - are inserted and the contained macros expanded
- Use -E flag to stop after pre-processing:
  > **cc -E -o hello.pp.c hello.c**

**4**

# Step 2: Compilation

- Compiler converts a high-level language into the specific instruction set of the target CPU

- Individual steps:

  - Parse text (lexical + syntactical analysis)

  - Do language specific transformations

  - Translate to internal representation units (IRs)

  - Optimization (reorder, merge, eliminate)

  - Replace IRs with pieces of assembler language

- Try:> `gcc -S hello.c` (produces `hello.s`)

# Compilation cont'd

```
        .file   "hello.c"
        .section    .rodata
.LC0:
        .string "hello, world!"
        .text
.globl main
        .type   main, @function
main:
        pushl   %ebp
        movl    %esp, %ebp
        andl    $-16, %esp
        subl    $16, %esp
        movl    $.LC0, (%esp)
        call    puts
        movl    $0, %eax
        leave
        ret
        .size   main, .-main
        .ident  "GCC: (GNU) 4.5.1 20100924 (Red Hat 4.5.1-4)"
        .section        .note.GNU-stack,"",@progbits
```

gcc replaced `printf` with `puts`

try: gcc -fno-builtin -S hello.c

```c
#include <stdio.h>
int main(int argc,
         char **argv)
{
 printf("hello world\n");
 return 0;
}
```

**6**

# Step 3: Assembler / Step 4: Linker

- Assembler (as) translates assembly to binary
  - Creates so-called object files (in ELF format)

```
Try: > gcc -c hello.c
Try: > nm hello.o
00000000 T main
         U puts
```

- Linker (ld) puts binary together with startup code and required libraries

- Final step, result is executable.
```
Try: > gcc -o hello hello.o
```

# Adding Libraries

- Example 2: exp.c

```
#include <math.h>
#include <stdio.h>
int main(int argc, char **argv)
{    double a=2.0;
     printf("exp(2.0)=%f\n", exp(a));
     return 0;
}
```

- > gcc -o exp exp.c
  Fails with "undefined reference to 'exp'". Add: -lm

- > gcc -O3 -o exp exp.c
  Works due to inlining at high optimization level.

# Symbols in Object Files & Visibility

- Compiled object files have multiple sections and a symbol table describing their entries:
    - "Text": this is executable code
    - "Data": pre-allocated variables storage
    - "Constants": read-only data
    - "Undefined": symbols that are used but not defined
    - "Debug": debugger information (e.g. line numbers)
- Entries in the object files can be inspected with either the "nm" tool or the "readelf" command

# Example File: visbility.c

```c
static const int val1 = -5;
const int val2 = 10;
static int val3 = -20;
int val4 = -15;
extern int errno;

static int add_abs(const int v1, const int v2) {
    return abs(v1)+abs(v2);
}

int main(int argc, char **argv) {
    int val5 = 20;
    printf("%d / %d / %d\n",
            add_abs(val1,val2),
            add_abs(val3,val4),
            add_abs(val1,val5));
    return 0;
}
```

```
nm visibility.o:
00000000 t add_abs
         U errno
00000024 T main
         U printf
00000000 r val1
00000004 R val2
00000000 d val3
00000004 D val4
```

# What Happens During Linking?

- Historically, the linker combines a "startup object" (crt1.o) with all compiled or listed object files, the C library (libc) and a "finish object" (crtn.o) into an executable (a.out)

- With current compilers it is more complicated

- The linker then "builds" the executable by matching undefined references with available entries in the symbol tables of the objects

- crt1.o has an undefined reference to "main" thus C programs start at the main() function

# Static Libraries

- Static libraries built with the "ar" command are collections of objects with a global symbol table

- When linking to a static library, object code is <u>copied</u> into the resulting executable and all direct addresses recomputed (e.g. for "jumps")

- Symbols are resolved "from left to right", so circular dependencies require to list libraries multiple times or use a special linker flag

- When linking only the <u>name</u> of the symbol is checked, not whether its argument list matches

# Shared Libraries

- Shared libraries are more like executables that are missing the main() function

- When linking to a shared library, a marker is added to load the library by its "generic" name (soname) and the list of undefined symbols

- When resolving a symbol (function) from shared library all addresses have to be recomputed (relocated) on the fly.

- The shared linker program is executed first and then loads the executable and its dependencies

# Differences When Linking

- Static libraries are fully resolved "left to right"; circular dependencies are only resolved between explicit objects or inside a library -> need to specify libraries multiple times or use: **-Wl,--start-group (...) -Wl,--end-group**

- Shared library symbols are **<u>not</u>** fully resolved at link time, only checked for symbols required by the object files. **<u>Full check</u>** only at runtime.

- Shared libraries may depend on other shared libraries whose symbols will be globally visible

# Dynamic Linker Properties

- Linux defaults to dynamic libraries:
  **> ldd hello**
  **linux-gate.so.1 =>  (0x0049d000)**
  **libc.so.6 => /lib/libc.so.6**
  **(0x005a0000)**
  **/lib/ld-linux.so.2 (0x0057b000)**
- **/etc/ld.so.conf, LD_LIBRARY_PATH**
  define where to search for shared libraries
- **gcc -Wl,-rpath,/some/dir** will encode
  **/some/dir** into the binary for searching

# What is Different in Fortran?

- Basic compilation principles are the same
  => preprocess, compile, assemble, link

- In Fortran, symbols are <u>case insensitive</u>
  => most compilers translate them to lower case

- In Fortran symbol names may be modified to make them different from C symbols
  (e.g. append one or more underscores)

- Fortran entry point is not "main" (no arguments)
  PROGRAM => MAIN__ (in gfortran)

- C-like main() provided as startup (to store args)

# Pre-processing in C and Fortran

- Pre-processing is <u>mandatory</u> in C/C++
- Pre-processing is <u>optional</u> in Fortran
- Fortran pre-processing enabled implicitly via file name: name.F, name.F90, name.FOR
- Legacy Fortran packages often use /lib/cpp: /lib/cpp -C -P <span style="color:blue">-traditional</span> -o name.f name.F
  - -C : keep comments (may be legal Fortran code)
  - -P : no '#line' markers (not legal Fortran syntax)
  - -traditional : don't collapse whitespace (incompatible with fixed format sources)

# Common Compiler Flags

- Optimization: -O0, -O1, -O2, -O3, -O4, ...
    - Compiler will try to rearrange generated code so it executes faster
    - Aggressive compiler optimization may not always execute faster or may miscompile code
    - High optimization level (> 2) may alter semantics
- Preprocessor flags: -I/some/dir -DSOM_SYS
- Linker flags: -L/some/other/dir -lm
  -> search for libm.so/libm.a also in /some/dir