

Лекция 4

Взаимодействие с окружением программы на Си/Си++

Реализации Си

- Freestanding реализация – поддерживается ограниченный набор заголовочных файлов и стандартных функций (например, `memset`)
 - Для ядер операционных систем
 - Для встроенных систем (embedded) без управления ОС
- Hosted реализация – полный набор (возможно, кроме опциональных) заголовочных файлов и библиотечных функций
 - Программирование на уровне пользовательских программ ОС

Стандартная библиотека Си (hosted в Linux)

- В Unix системах традиционно называется libc, является частью ОС
- Заголовочные файлы размещаются в /usr/include
- Бинарный динамически загружаемый файл: /lib/libc.so.6 (Linux)
- Помимо функций библиотеки Си содержит и функции POSIX и расширения
- Библиотека математических функций отдельно – libm – требуется опция -lm при компиляции

Взаимодействие программы на Си с окружением

- Стандартные потоки ввода и вывода `stdin`, `stdout`, `stderr`
- Аргументы командной строки
- Переменные окружения
- Код завершения программы

Обработка ошибок

- Библиотечные функции и системные вызовы в случае ошибки возвращают специальное значение (например, `open` возвращает `NULL`, часто возвращается `-1`)
- В этом случае переменная `errno` содержит код ошибки, например, `EPERM`, `EAGAIN`
- Переменная `errno` и коды ошибок определены в `<errno.h>`
- `strerror` из `<string.h>` возвращает строку, соответствующую ошибке
- Сообщения об ошибках должны выводиться на `stderr`

Взаимодействие со средой

- Процесс завершается системным вызовом `_exit(exitcode)` или `exit` или `_Exit`
- Или возвращаемое значение `return` из `main`
- Значение в диапазоне `[0;255]` — код завершения процесса, он доступен процессу-родителю
- Код 0 — успешное завершение (`/bin/true`)
- Ненулевой код — ошибка (`/bin/false`)

Аргументы командной строки

- Функция `main` получает аргументы командной строки:

```
int main(int argc, char *argv[])
```

- `argv` — массив указателей на строки Си

```
./prog foo 1 bar
```

```
argv[0] → "./prog";    путь к программе
```

```
argv[1] → "foo";
```

```
argv[2] → "1";
```

```
argv[3] → "bar";
```

```
argv[4] → NULL;
```

- Передаются на стеке процесса

argv[0]

- Обычно argv[0] – путь, использованный для запуска программы
- Некоторые программы анализируют argv[0] и модифицируют свое поведение (например, busybox)

Переменные окружения

- Именованные значения доступные процессу
- По умолчанию передаются неизменными порождаемым процессам

```
char *getenv(const char *name);
```

- В процесс передаются на стеке
- Глобальная переменная `environ` содержит указатель на массив переменных

Строки в Си

- Null-terminated strings – в конце строки находится байт 0 (или '\0') – признак конца строки
- Строковые литералы “abcd” содержат “невидимый” \0 в конце
 - `char s[] = “abcd”; // sizeof(s) == 5`
- Если под строковый литерал память явно не выделяется, он размещаются в read-only памяти
 - `char *s = “abcd”; // sizeof(s) = sizeof(void*)`
`s[2] = 'd'; // undefined behavior`

Pros & contras

- (+) для работы со строкой достаточно одного указателя
- (+) сдвигая указатель по строке вперед все равно получаем строку
- (-) получение длины строки (strlen) выполняется за линейное время
 - **НИКОГДА!**
for (int i = 0; i < strlen(s); ++i) {...}
- (-) нельзя использовать \0 в строке

Альтернативы

- Хранить пару <указатель, длина> (std::string)
 - (+) нет проблемы байта \0
 - (-) размер такой структуры в два раза больше (а размер самой строки на один байт меньше)
- Хранить длину в начале строки (pascal style)
 - Либо ограниченный размер (если длина – 1 байт), либо неэффективное использование памяти (4 байта длины для коротких строк - много)

Управление памятью

- В Си практически все управление памятью возложено на программиста
- При работе с указателями важно понимать, как и где выделена память, на которую он указывает:
 - Глобальная/статическая память
 - Thread-local storage
 - Автоматическая память
 - Динамическая память (куча)

Буфер строки

- Буфер – область памяти, отведенная для хранения строки
- Буфер имеет ограниченный размер, но размер может изменяться
- При обработке строки “на чтение” достаточно только указателя на строку
- При формировании строки в памяти важен и адрес буфера, и размер буфера

Переполнение буфера

- Если не контролируется размер данных, записываемых в буфер, возможно **переполнение буфера**
- Может иметь катастрофические последствия для безопасности системы (arbitrary code execution)

Переполнение буфера

- Если не контролируется размер данных, записываемых в буфер, возможно **переполнение буфера**
- Может иметь катастрофические последствия для безопасности системы (arbitrary code execution)
- CVE-2015-2712 (Firefox)
- CVE-2010-1117 (IE)
- CVE-2016-5157 (Chrome)

Good vs evil

- “Плохие” функции: записывают строку, но не принимают параметр размера буфера: `gets`, `scanf(“%s”, ...)`, `strcpy`, `sprintf`
 - `gets`, `scanf` – запрещены; `strcpy`, `sprintf` – крайне осторожно
- “Хорошие” функции: записывают строку и принимают размер буфера строки: `fgets`, `snprintf`, `scanf(“%100s”, ...)`

Кодировки текста

- Исторически: ASCII – символы с кодами 0-127; достаточно для английского языка
- Недостаточно для других языков
 - Использование “верхней” части байта, коды 128-255: (iso8859-X, cpYYY, koī8-r, ...)
 - Специальные символы-переключатели состояния (JIS, EUC-JP – японский язык)
- Обмен текстовыми документами сложен

Unicode

- Определяет 1,114,112 кодовых позиций (Code Points) (обозначаются U+0 ... U+10FFFF)
- U+D800 – U+DFFF – недопустимы в корректном Unicode (UCS4, UTF8)
- Кодовые позиции содержат глифы всех известных письменностей, диакритические знаки
- U+0 – U+7F совпадает с ASCII
- Возможны разные кодировки (битовое представление для Code Points)

Кодировки Unicode

- UCS-4 (один CodePoint – 32-битный int)
 - (+) фиксированный размер – удобно обрабатывать
 - (-) 4 байта на все codepoints
 - (-) много байтов \0 в тексте – несовместим с ASCII
- UCS-2 (один CodePoint – uint16_t) – только для U+0 – U+FFFF
- UTF-16 (один CodePoint – один или два uint16_t)

UTF-8

- Байтовый поток
- Один CodePoint кодируется от 1 до 4 байт
- U+0 – U+7f кодируются 1 байтом (совместимость с ASCII)
- Байт \0 всегда обозначает U+0 и может использоваться как терминатор строки – совместимость с Си-строками
- По любому месту в потоке можно найти начало кодировки соответствующего CodePoint

UTF-8

- Кодирование Code points в UTF-8
- Overlong encoding (длина последовательности больше минимальной, например 0xC0 0xAF → '/') запрещен

UTF-8 (2003)

Number of bytes	Bits for code point	First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
1	7	U+0000	U+007F	0xxxxxxx			
2	11	U+0080	U+07FF	110xxxxx	10xxxxxx		
3	16	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
4	21	U+10000	U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Поддержка в C/C++

- `wchar_t` тип данных для хранения unicode codepoints во внутреннем представлении (unsigned short – Windows, int или long – Unix)
- В Unix внутреннее представление – UCS4
- Wide-char literals: `L'a'`
- Wide-char string literals: `L"Привет"`
- Функции: `getwc`, `fgetws`, `wscanf`, `wprintf`, `wcslen`, ...

Локаль (Locale)

- Определяет кодировку в системе, региональные особенности, язык взаимодействия с пользователем
- Переменные окружения LANG и LC_*
- LANG=en_US.utf8 – язык/регион – американский английский, кодировка UTF8
- LANG=ru_RU.UTF-8 – русский/Россия, UTF8
- LANG=C – по умолчанию ASCII

Setlocale

- `setlocale` позволяет установить локаль для выполняющейся программы
- По умолчанию – C, не позволяет обрабатывать символы вне ASCII
- Для установки системной локали:
`setlocale(LC_ALL, "");`