

1 Абстрактные синтаксические деревья

1.1 Управление памятью на основе регионов

1.1.1 Мотивировка

Текущая реализация абстрактного синтаксического дерева имеет следующие недостатки:

1. Выделение памяти стандартным методом может значительно фрагментировать оперативную память, затрудняя доступ к ней.
2. Любое выделение и удаление памяти требует вмешательства системных вызовов, что может стать причиной дополнительных издержек во время работы программы.
3. Программист не имеет возможности ручного управления выделяемой им памятью.

Избавиться от этих недостатков можно используя различные оптимизации. В рамках этой работы воспользуемся управлением памятью на основе, так называемых, регионов (арен, зон) [1].

Под регионом далее будем понимать непрерывную область памяти, содержащую внутри себя объекты. При запуске программы выделим регион некоторого размера, при необходимости увеличивая его размер в некоторое постоянное число раз.

Этот подход имеет следующие преимущества:

1. Элементы располагаются последовательно, в связи с чем минимизируется фрагментация и упрощается доступ к объектам.
2. Выделение и освобождение памяти выполняется с минимальными издержками.
3. Программисту предоставляется большая свобода для управления выделенной памятью.

1.1.2 Построение

Формально определим требования к системе:

1. Регион должен представлять из себя некоторый непрерывный участок размера n байт (в начальный момент времени размер равен некоторой начальной величине n_0).
2. При обращении к региону он должен предоставить k байт памяти и вернуть некоторый идентификатор этого участка для последующего обращения.

3. При заполнении региона должна быть возможность увеличить объем доступной памяти в некоторое число раз, которое далее будем называть коэффициентом увеличения.
4. Должна быть доступна возможность эффективного освобождения всей выделенной регионом памяти.

Единственной сложной операцией над регионом является его увеличение. Так как выделение нового участка потенциально может сопровождаться изменением адресов объектов, то необходимо организовать доступ к ним независимо от первоначального адреса. Для этого для каждого объекта будем получать доступ к нему через некоторый индекс.

Кроме того, коэффициент увеличения должен быть выбран таким образом, чтобы был соблюден баланс между оптимальным объемом выделенной памяти и частотой системных вызовов.

1.1.3 Определение структуры

Определим нашу структуру следующим образом:

```
1 typedef struct arena {
2     // Указатель на начало региона
3     struct node* arena;
4     // Размер региона
5     unsigned int size;
6     // Объем выделенной регионом памяти
7     unsigned int allocated;
8 } arena;
```

1.1.4 Инициализация

Теперь определим функцию `arena_construct`, выполняющую начальную инициализацию состояния региона:

```
1 int arena_construct (arena* arena) {
2     // Начальный размер региона равен некоторой постоянной, равной
   ↪  DEFAULT_ARENA_SIZE
3     arena->size = DEFAULT_ARENA_SIZE;
4     arena->allocated = 0;
5     // Выделим необходимое число памяти
6     arena->arena = malloc(sizeof(node) * DEFAULT_ARENA_SIZE);
7     // Если выделение прошло неудачно - вернем в качестве кода ошибки отличное
   ↪  от 0 значение.
```

```

8     if (arena->arena == NULL) {
9         return (!0);
10    }
11    return 0;
12 }

```

1.1.5 Выделение памяти

После выделения некоторого объема памяти возможно обращение к ней. Определим это обращение с помощью функции `arena_allocate`:

```

1  int arena_allocate (arena* arena, unsigned int count) {
2      // Если места в регионе недостаточно
3      if (arena->allocated + count >= arena->size) {
4          // Определим новый размер региона
5          unsigned int newSize = MULTIPLY_FACTOR * arena->size;
6          // Выделим регион большего размера и освободим ранее занятую память
7          node* newArena = realloc(arena->arena,
8                                   newSize * sizeof(node));
9          if (NULL == newArena) {
10             return -1;
11         }
12         arena->arena = newArena;
13         arena->size = newSize;
14     }
15     // В качестве результата вернем индекс первого свободного участка региона
16     unsigned int result = arena->allocated;
17     // Сместим индекс на объем выделенной памяти
18     arena->allocated += count;
19     // Вернем результат
20     return result;
21 }

```

Отметим, что наиболее часто значением `MULTIPLY_FACTOR` оказывается числа 1.5 и 2. Это позволяет достичь амортизационно константного времени выполнения операции выделения памяти [2].

1.1.6 Освобождение выделенной памяти

Наконец, реализуем освобождение выделенной региону памяти с помощью функции `arena_free`

```

1 void arena_free (arena* arena) {
2     if (arena->arena != NULL)
3         free(arena->arena);
4     arena->arena = NULL;
5 }

```

1.1.7 Модификация абстрактного синтаксического дерева

Осталось изменить исходный код программы, чтобы обеспечить выделение памяти с помощью полученной нами структуры данных.

Для этого воспользуемся директивой `%param` и заявим в качестве параметра переменную типа `arena*`. В функциях `eval`, `newnum`, `newast` внесем изменения, чтобы обеспечить выделение памяти с помощью написанных ранее функций.

С полным кодом программы можно ознакомиться в приложении А.

1.1.8 Сборка проекта

Теперь проект можно собрать, незначительно изменив `Makefile`:

```

1 calc.out: calc.l calc.y arena_ast.h
2     bison -d calc.y
3     flex calc.l
4     cc -o $@ calc.tab.c lex.yy.c arena_ast.c arena.c

```

и запустить. Результат работы программы представлен на рис. 1



Рисунок 1 – Демонстрация работы программы

2 Сравнение полученных реализаций

Проведем анализ производительности полученных версий анализатора. В качестве данных для тестирования возьмем выражения вида $\underbrace{2 + 2 + 2 \dots + 2}_n$ для $n = 1 \dots 100$ с шагом 1. Для вычисления времени выполнения воспользуемся библиотекой `time` Python 3.9.5. Автоматизацию обеспечим с помощью библиотеки `subprocess`. Получим следующий код:

```
1 import subprocess as sb
2 import time
3 import sys
4 def run(args):
5     return sb.run(args,
6                   capture_output=True, ).stdout.decode().strip()
7
8 def main():
9     if (len(sys.argv) < 6):
10         print("Invalid arguments")
11         return
12     exe_path = sys.argv[1]
13     out_path = sys.argv[2]
14     right_bound = int(sys.argv[3])
15     step = int(sys.argv[4])
16     iter = sys.argv[5]
17
18     f = open(out_path, "w")
19     f.write(f"{exe_path} \n ")
20     f.close()
21     for expr_len in range(1, right_bound, step):
22         test_string = "+".join(['2'] * expr_len)
23         args = [exe_path, test_string, iter]
24         t = time.monotonic()
25         run(args)
26         end_t = time.monotonic()
27         f = open(out_path, "a")
28         f.write(f"{expr_len} {(end_t - t) / (int(iter))} \n ")
29         f.close()
30         print(f" Step {expr_len} finished")
31
32 if __name__ == "__main__":
33     main()
```

Кроме того, отметим, что в ранее написанные программы были внесены некоторые изменения для проведения эксперимента. Ознакомиться с ними можно в приложении А.

Ознакомиться с полным исходным кодом программы, осуществляющей исследование производительности можно в приложении Б.

Для большей наглядности графики интерполированы полиномом с помощью функции `polyfit` библиотеки `numpy`.

Ознакомиться с полным исходным кодом программы, осуществляющей анализ полученных результатов можно в приложении В.

Результаты исследования изображены на рис. 2:

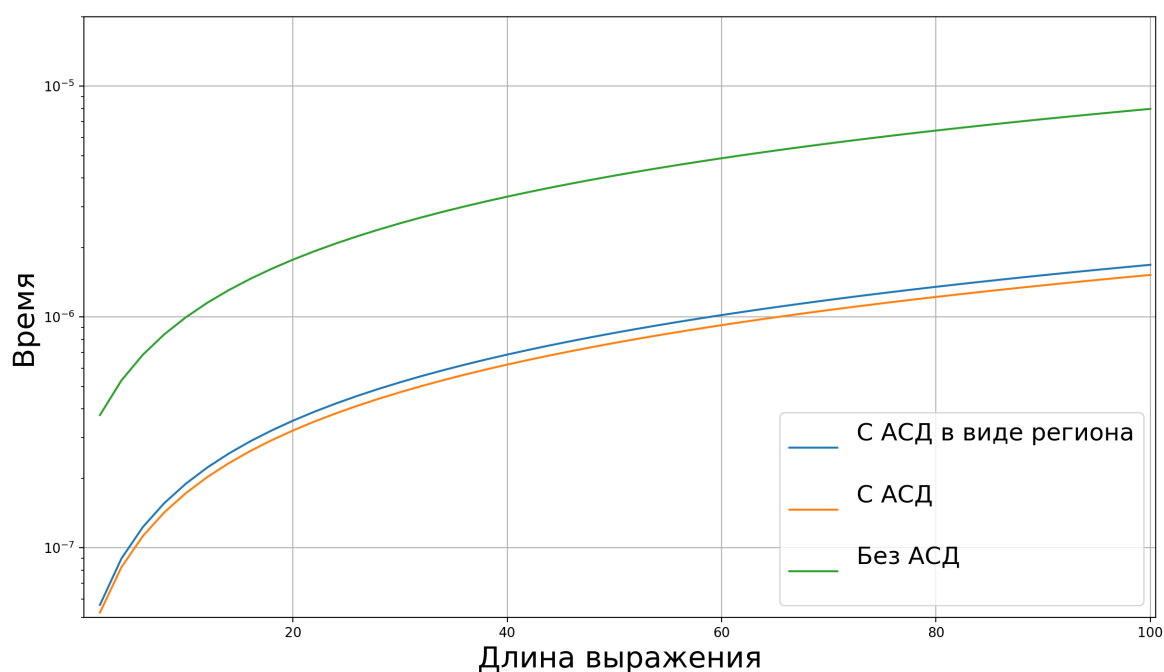


Рисунок 2 – Сравнение полученных результатов

Исследование показало, что использование абстрактных синтаксических деревьев позволяет уменьшить время работы программы более чем в 5 раз, что существенно заметно для выражений любой длины.

Также из графиков видно, что в рамках данной работы не удалось добиться большей производительности при управлении памятью на основе регионов. Тем не менее, она все еще может считаться более предпочтительной ввиду перечисленных ранее преимуществ.

ЗАКЛЮЧЕНИЕ

В ходе данной работы:

1. Были изучены теоретические основы построения лексических и синтаксических анализаторов.
2. Проанализированы особенности реализации лексических и синтаксических анализаторов.
3. Были изучены принципы работы генераторов лексического и синтаксического анализа на примере Flex и GNU Bison.
4. Были созданы лексический и синтаксический анализаторы для анализа математического выражения.
5. Было изучено понятие абстрактного синтаксического дерева.
6. Проведен анализ производительности полученных реализаций.

Таким образом, все поставленные в рамках работы задачи выполнены.

Результаты исследования показали, что абстрактные синтаксические деревья позволяют добиться увеличения производительности в 5–6 раз.

А это, в свою очередь, позволяет утверждать о том, что концепция абстрактных синтаксических деревьев является крайне важной в информатике и ее приложениях, в частности, при создании синтаксических анализаторов.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 Wang, D. C.-A. Managing memory with types / D. C.-A. Wang. — Princeton University, 2002. — Pp. 29–55.
- 2 Facebook open-source library documentation [Электронный ресурс]. — URL: <https://github.com/facebook/folly/blob/main/folly/docs/FBVector.md> (Дата обращения 25.05.2021). Загл. с экр. Яз. англ.

ПРИЛОЖЕНИЕ А

Flash-носитель с исходным кодом программ, использующихся в работе

Папка `src` содержит оригинальный исходный код программы:

Папка `naive` — реализация без АСД

Папка `naiveast` — реализация с АСД

Папка `arena` — реализация с АСД на основе региона

Папка `extsrc` содержит измененный исходный код, необходимый для исследования производительности:

Папка `naive` — реализация без АСД

Папка `naiveast` — реализация с АСД

Папка `arena` — реализация с АСД на основе региона

ПРИЛОЖЕНИЕ Б

Исходный код программы на Python, осуществляющей исследование производительности полученных реализаций

```
1 import subprocess as sb
2 import time
3 import sys
4 def run(args):
5     return sb.run(args,
6         capture_output=True, ).stdout.decode().strip()
7
8 def main():
9     if (len(sys.argv) < 6):
10         print("Invalid arguments")
11         return
12     exe_path = sys.argv[1]
13     out_path = sys.argv[2]
14     right_bound = int(sys.argv[3])
15     step = int(sys.argv[4])
16     iter = sys.argv[5]
17
18     f = open(out_path, "w")
19     f.write(f "{exe_path} \n ")
20     f.close()
21     for expr_len in range(1, right_bound, step):
22         test_string = "".join(['2'] * expr_len)
23         args = [exe_path, test_string, iter]
24         t = time.monotonic()
25         run(args)
26         end_t = time.monotonic()
27         f = open(out_path, "a")
28         f.write(f "{expr_len} {(end_t - t) / (int(iter))} \n ")
29         f.close()
30         print(f " Step {expr_len} finished")
31
32 if __name__ == "__main__":
33     main()
```

ПРИЛОЖЕНИЕ В

Исходный код программы на Python, осуществляющей анализ полученных результатов

```
1 import subprocess as sb
2 from time import time
3 import matplotlib.pyplot as plt
4 import sys
5 import numpy as np
6
7 legend = []
8 for index in range(1, len(sys.argv)):
9     file_name = sys.argv[index]
10    f = open(file_name, "r")
11    try:
12        parser_type = f.readline()
13    except StopIteration:
14        parser_type = "Undefined parser"
15    legend.append(parser_type)
16
17    x_axis = []
18    y_axis = []
19    for line in f:
20        try:
21            w, h = [float(x) for x in next(f).split()]
22        except StopIteration:
23            break
24        x_axis.append(w)
25        y_axis.append(h)
26    plt.xlabel("Длина выражения", size = 22)
27    plt.ylabel("Время", size = 22)
28    p = np.polyfit(x_axis, y_axis, 1)
29    for i in range(len(x_axis)):
30        y_axis[i] = p[0] * x_axis[i] + p[1]
31
32    plt.plot(x_axis, y_axis)
33    plt.xlim(0.5, 100.5)
34    plt.ylim(5e-8, 2e-5)
35    plt.rcParams.update({'font.size': 18})
36    #plt.yscale("log")
37    plt.grid(True)
38    plt.legend(legend, loc="lower right")
```

39 `plt.show()`