

Московский авиационный институт
(национальный исследовательский университет)

Институт информационных технологий и прикладной
математики

Кафедра вычислительной математики и программирования

Лабораторная работа №3 по курсу «Дискретный анализ»

Студент: И. С. Глушатов
Преподаватель: А. А. Кухтичев
Группа: М8О-207Б-19
Дата:
Оценка:
Подпись:

Москва, 2020

Лабораторная работа №3

Задача: Для реализации словаря из предыдущей лабораторной работы необходимо провести исследование скорости выполнения и потребления оперативной памяти. В случае выявления ошибок или явных недочётов, требуется их исправить. Результатом лабораторной работы является отчёт, состоящий из:

- Дневника выполнения работы, в котором отражено что и когда делалось, какие средства использовались и какие результаты были достигнуты на каждом шаге выполнения лабораторной работы.
- Выводов о найденных недочётах.
- Сравнение работы исправленной программы с предыдущей версии.
- Общих выводов о выполнении лабораторной работы, полученном опыте.

1 Описание

Инструменты, которыми я пользовался: valgrind, plot, perf

Схема директорий:

- main.cpp
- string.hpp
- avltree.hpp
- tests_creator.py
- makefile
- bash_script.sh
- gnuplot_script.txt
- perf.data
- plots (папка с графиками)
- tests (папка с тестами)
- scrins (папка с скриншотами)
- report (папка с отчетом)

Во-первых, для того, чтобы анализировать свою программу, время ее исполнения и утечки памяти, нужно было написать генератор тестов. Для данной задачи я выбрал язык Python.

```
1 import random
2 import copy
3 from random import choice
4 from string import ascii_uppercase
5
6 randstring = lambda x=1: ''.join(choice(ascii_uppercase) for i in range(x))
7
8 if __name__ == "__main__":
9
10     count_of_tests = 5
11     max_count_of_commands = 50
12     count_of_commands = 1
13
14     actions = ["+", "!", "?", "!"] # + - ? !
15     for test_number in range(1, count_of_tests+1):
16         tree_elems = {}
17         saved = 0
18         saved_file = {}
19         test_file = open("tests/" + str(test_number) + ".t", 'w')
20         file_with_answers = open("tests/" + str(test_number) + ".txt", "w")
21
22         for _ in range( count_of_commands ):
23             action = random.choice( actions )
24             if action == "+":
25                 key = randstring()
26                 value = random.randint( 1, 2**64-1 )
```

```

27         test_file.write("+ " + key + " " + str(value) + "\n")
28
29         key = key.lower()
30         answer = "Exist"
31         if key not in tree_elems:
32             answer = "OK"
33             tree_elems[key] = value
34         file_with_answers.write(answer + "\n")
35
36     elif action == "-":
37         key = randstring()
38         test_file.write("- " + key + "\n")
39
40         key = key.lower()
41         answer = "NoSuchWord"
42         if key in tree_elems:
43             del tree_elems[key]
44             answer = "OK"
45         file_with_answers.write(answer + "\n")
46
47     elif action == "?":
48         search_exist_element = random.choice( [ True, False ] )
49         key = random.choice([key for key in tree_elems.keys()]) if
            search_exist_element and len(tree_elems.keys()) > 0 else randstring
            ()
50         test_file.write(key + "\n")
51
52         key = key.lower()
53         if key in tree_elems:
54             answer = "OK: " + str(tree_elems[key])
55         else:
56             answer = "NoSuchWord"
57         file_with_answers.write(answer + "\n")
58
59     elif action == "!":
60         act_file = random.choice(["Load test", "Save test"])
61         if act_file == "Save test":
62             test_file.write(action + " " + act_file + "\n")
63             saved_file = tree_elems.copy()
64             file_with_answers.write("OK" + "\n")
65             saved = 1
66
67         elif saved == 1 and act_file == "Load test":
68             test_file.write(action + " " + act_file + "\n")
69             tree_elems = {}
70             tree_elems = saved_file.copy()
71             file_with_answers.write("OK" + "\n")
72
73     count_of_commands+=(max_count_of_commands-count_of_commands)/count_of_tests

```

Генератор тестов хорошо работает и не учитывает только случаи, когда нельзя будет открыть файл по каким-то причинам. Поэтому он не сможет проверить ошибки с ERROR. Так как я хочу испытывать произвольное число тестов и причем их большое количество, я написал bash файл, где будет запускаться моя программа, перенаправляя вывод в файл, а потом с помощью утилиты diff сравнивать файлы от питона и от моей программы.

```

1  #!/bin/bash
2  max=6
3
4  for ((i = 1; i < ${max}; i++))

```

```

5 do
6   echo "test: $i ";
7   ./main.out < tests/${i}.t > tests/${i}ans.txt;
8   echo "ans: $i ";
9   diff tests/${i}ans.txt tests/${i}.txt
10 done

```

Переменная `max` равняется количеству тестов, которые нужно прогнать.

В самой программе я установил таймер из библиотеки `<chrono>`, а так же счетчик команд. Результат времени работы программы и счетчик команд выводятся через пробел в файл `time_log.txt`, откуда `gnuplot` считывает данные и строит график зависимости время исполнения программы от количества входных данных. Для `gnuplot` так же написан скрипт. Мы устанавливаем название осей, расширение изображения, и директорию, куда будет выводиться картинка.

```

1 set xlabel "data count"
2 set ylabel "time of appending"
3 set terminal png
4 set output "plots/plot.png"
5 plot "time_log.txt" with lines
6 pause -1

```

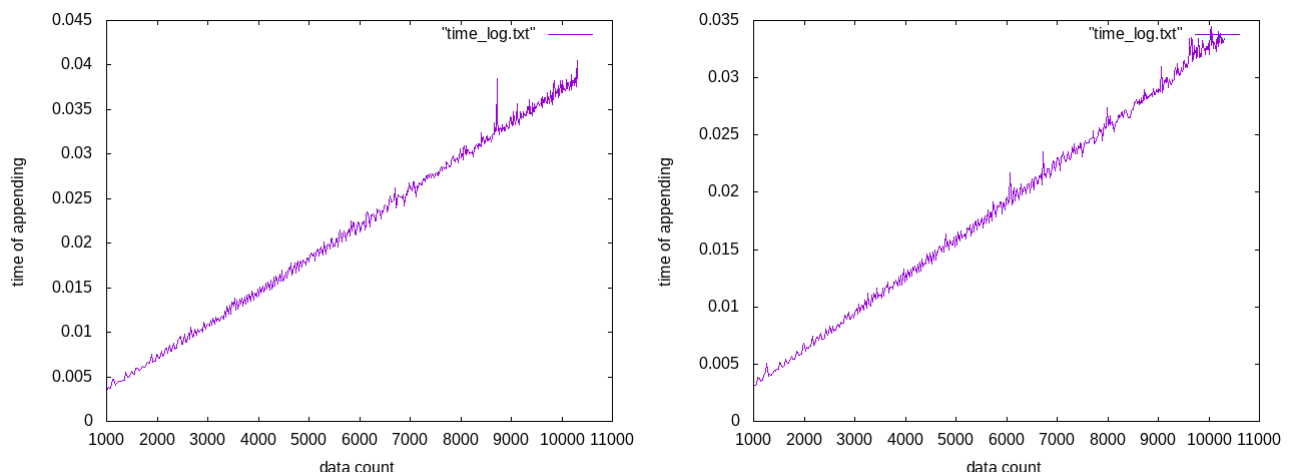
Makefile собирает все скрипты вместе и запускает их.

```

1 #!/bin/bash
2 CC = g++
3 FLAGS = -pedantic -Wall -std=c++11 -Werror -Wno-sign-compare -O2 -lm
4
5 all: main.out
6
7 main.out: main.cpp avltree.hpp string.hpp
8   ${CC} ${FLAGS} main.cpp -o main.out
9
10 run:
11   rm -rf time_log.txt
12   python tests_creator.py
13   ./bash_script.sh
14   gnuplot gnuplot_script.txt
15
16 testss:
17   python tests_creator.py

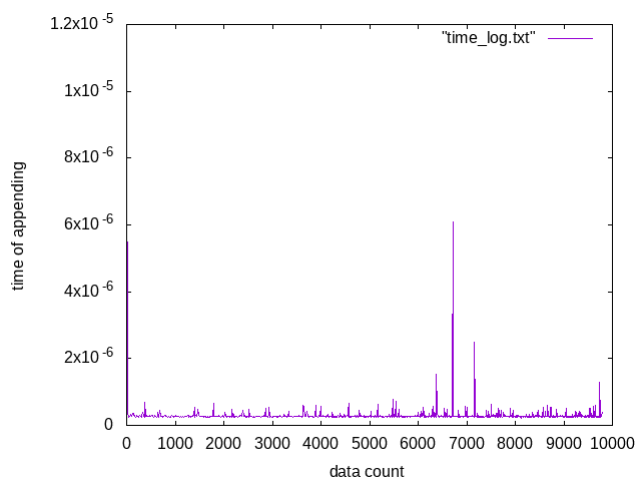
```

Первые результаты бенчмарка я опробовал на тестах с только вставкой и на тестах с вставкой, удалением и поиском.



Для меня было удивлением увидеть такие результаты, ведь я ожидал график с ас-

симптоматикой $n \log_2 n$, а увидел прямопропорциональную зависимость времени исполнения программы от количества элементов. Поэтому я решил провести тест, где строился график зависимости времени вставки в дерево от количества уже вставленных элементов в него, т.е. время вставки последнего элемента в дерево.



И по данному графику я понял, почему два до этого давали такой результат. На таком маленьком числе элементов (порядка 10000) время вставки было практически константным. Скорее всего, если бы я мог задействовать одно ядро компьютера полностью под свою программу и взять тесты побольше, то я бы увидел плавный рост, однако не обладая такими возможностями, я удовлетворился этими результатами.

Несмотря на такие хорошие показатели, моя программа не сразу прошла на чекере и причина тому теперь как никогда ясна - команды Save и Load сильно тормозили процесс. Поэтому пришлось время провести тест, добавив эти две команды. Сначала я использовал версию программы, где загрузка дерева проходила с помощью вставки в обычное бинарное дерево. Потом я сделал еще одну версию, где сохранение не сильно отличалось от старой, а загрузка производилось с помощью флагов, которые означали есть ли ребенок слева и справа.

```

1 void TAVLTree::Save(const TAVLTree *head, FILE *s) {
2     if (head) {
3         fprintf(s, "%s %llu ", head->key.str, head->value);
4         Save(head->left, s);
5         Save(head->right, s);
6     }
7 }
8
9 TAVLTree* TAVLTree::BinInsert(TAVLTree *head, const NMystd::TString& key, const
    unsigned long long value) {
10     if (head) {
11         if (key > head->key) {
12             head->right = BinInsert(head->right, key, value);
13         } else {
14             head->left = BinInsert(head->left, key, value);
15         }
16         return head;
17     }
18     return new TAVLTree(key, value);
19 }
20
21 TAVLTree* TAVLTree::Load(std::ifstream &in) {
22     TAVLTree *res = 0;
23     char str[256];
24     unsigned long long val;
25     while (fscanf(s, "%s %llu ", str, &val)>0) {

```

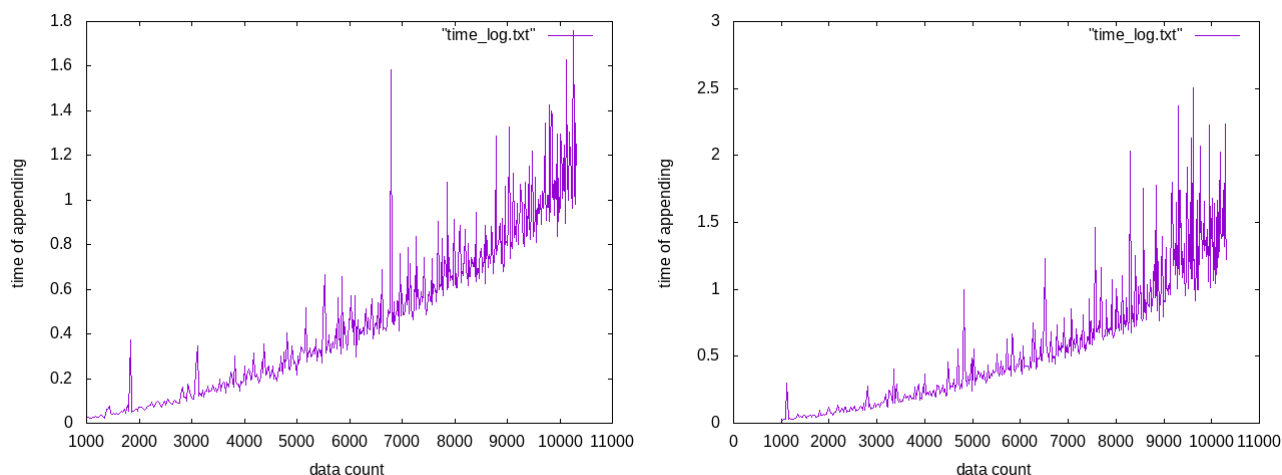
```

26     res = res->BinInsert(res, NMystd::TString(str), val);
27 }
28 return res;
29 }

1 void TAVLTree::Save(const TAVLTree *head, FILE *s) {
2     if (!head) {
3         return;
4     }
5
6     fprintf(s, "%s %llu %u ", head->key.str, head->value, head->h);
7     if (head->left) {
8         fprintf(s, "1 ");
9         Save(head->left, s);
10    } else {
11        fprintf(s, "0 ");
12    }
13    if (head->right) {
14        fprintf(s, "1 ");
15        Save(head->right, s);
16    } else {
17        fprintf(s, "0 ");
18    }
19 }

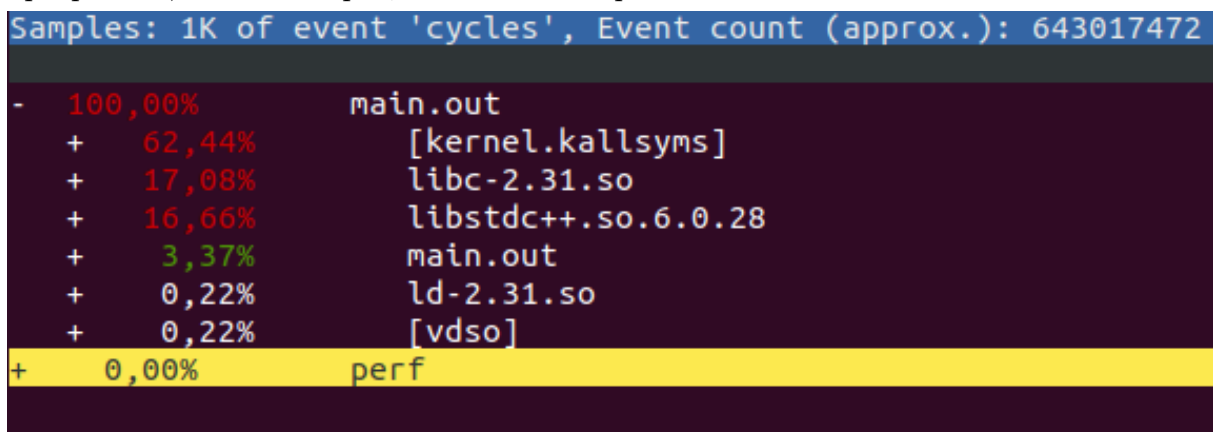
20
21 TAVLTree* TAVLTree::Load(std::ifstream &in) {
22     TAVLTree *res;
23     int flag;
24     char str[256];
25     unsigned long long val;
26     unsigned int h;
27
28
29     if (in >> str >> val >> h) {
30         res = new TAVLTree(NMystd::TString(str), val);
31         res->h = h;
32         in >> flag;
33         if (flag==1) {
34             res->left = Load(in);
35         } else {
36             res->left = NULL;
37         }
38         in >> flag;
39         if (flag==1) {
40             res->right = Load(in);
41         } else {
42             res->right = NULL;
43         }
44     } else {
45         res = NULL;
46     }
47     return res;
48 }

```



Несмотря на то, что изначальный вариант программы имеет лучшие результаты, видно, что ассимптотика у них одинаковая, а из-за шума посторонних процессов компьютера сложно наверняка определить, какой из вариантов лучше. В первом случае сложность загрузки $n \log_2 n$, так как мы вытаскиваем из файла по одному элементу ($O(n)$), и вставляем в дерево за $\log_2 n$. Во втором случае мы встречаемся с элементом только один раз, вставляем и больше к нему не возвращаемся, из чего следует сложность ($O(n)$). Однако как мы убедились на прошлых графиках, вставка на таком количестве элементов производится почти за константное время, поэтому графики выходят с практически одинаковой ассимптотикой. Почему не линейный вид графика? Потому что считается время выполнения всей программы, а значит чем больше команд подается, тем больше сохранений и загрузок будут появляться в разных местах теста, когда в дереве то 10 элементов, то 10000. Точный расчет сложности от таких случайных тестов должен включать теорию вероятностей, поэтому на данный момент остается для меня нерешаемой задачей.

Утилита Perf предоставляет инструмент анализа производительности в Linux посредством счетчиков. К сожалению, документация этой утилиты крайне сложная и размытая, однако я попытался чуть-чуть разобраться. Задача, которую я хочу достичь с помощью Perf - узнать распределение нагрузки частей программы на процессор. Что работает медленнее (а значит дольше), а что работает хорошо. Первая сложность появилась на этапе запуска этой утилиты. Выдавалась ошибка об отказе доступа, которая решалась запуском с `sudo`. Разобравшись со множеством непонятного синтаксиса команд, я смог запустить в итоге трассировку именно моей пользовательской программы, а не всех процессов компьютера.



| | | |
|---|---------|-----------------------------|
| - | 100,00% | main.out |
| - | 62,44% | [kernel.kallsyms] |
| + | 2,50% | [k] do_syscall_64 |
| + | 1,88% | [k] crc32c_pcl_intel_update |
| + | 1,45% | [k] psi_task_change |
| | 1,05% | [k] entry_SYSCALL_64 |
| + | 1,00% | [k] ext4_do_update_inode |
| | 0,96% | [k] memset_erms |

Из первых двух скриншотов мы видим, что большая часть нагрузок пала на системные вызовы ядра, что не удивительно. В основном все они внутри распределены равномерно и каждому уходит не более 2%, кроме `do_syscall_64`.

| | | |
|---|--------|--|
| + | 17,00% | libc-2.31.so |
| - | 16,66% | libstdc++.so.6.0.28 |
| - | 5,25% | [.] std::num_get<char, std::istreambuf_iterator<char, std::char_traits<char> > >::_M_extract_int<unsigned long long> |
| | 5,09% | 0x48fd8948550010bd |
| | | std::_numpunct_cache<char>::~~numpunct_cache |
| | | - 0x1 |
| | 5,02% | std::num_get<char, std::istreambuf_iterator<char, std::char_traits<char> > >::_M_extract_int<unsigned long long> |
| + | 2,05% | [.] std::istream::sentry::sentry |

Из библиотек `std` больше всего уходит нагрузка на вызов функции `num_get`, связанная с вводом и выводом, записью и чтением чисел, так как изначально они представляются в виде символов. То есть это своего рода парсер, а так как ввода и вывода, записи и чтения при командах загрузки и сохранения дерева много, то такой расход ресурсов вполне объясняется.

| | | |
|---|-------|--------------------------|
| + | 1,33% | [.] TAVLTree::Load |
| + | 0,84% | [.] TAVLTree::Save |
| | 0,48% | [.] main |
| | 0,31% | [.] TAVLTree::RemoveTree |
| | 0,14% | [.] 0x000000000000002224 |
| | 0,13% | [.] TAVLTree::Insert |
| | 0,07% | [.] 0x000000000000002384 |
| | 0,06% | [.] TAVLTree::Balancing |

Как и предполагалось, функция `Load` использует больше всего ресурсов, нежели функции вставки, удаления, перебалансировки и даже сохранения в файл.

| | | |
|---|--------|--|
| + | 17,08% | libc-2.31.so |
| + | 16,66% | libstdc++.so.6.0.28 |
| - | 3,37% | main.out |
| - | 1,33% | [.] TAVLTree::Load |
| | | 0x48fb894853000c92 |
| | | std::basic_ifstream<char, std::char_traits<char> >::~~basic_ifstream |
| - | 0 | |
| | 1,28% | TAVLTree::Load |
| + | 0,84% | [.] TAVLTree::Save |
| | 0,48% | [.] main |
| | 0,31% | [.] TAVLTree::RemoveTree |

| | |
|-------|--|
| 51,19 | <pre> xor %eax,%eax shr \$0x3,%ecx rep stos %rax,%es:(%rdi) mov %r14d,0x110(%r12) movslq %r14d,%rcx </pre> |
|-------|--|

И если мы пойдем глубже, то увидим, что у нас слишком часто используется ассемблерная команда REP повторения строковой операции, а именно REP STOS - заполнение блока по адресу содержимым. Станным оказалось то, что я не увидел нагрузку на сравнение строк, хотя казалось бы, что это занимает достаточное количества времени, и по-хорошему было бы их хэшировать, чтобы добиться оптимизации кода.

Привести интересные проверки с valgrind'ом я не смогу, так как все утечки памяти были устранены ранее. Класс NMystd::string, написанный для второй лабораторной работы по условию задачи имеет ограничение на 256 символов, поэтому я не делал динамическое выделение памяти для этого типа и никаких утечек там быть не могло, а AVL-дерево достаточно простое и утечки памяти, к примеру, при удалении узла было легко отследить и устранить.

2 Выводы

В ходе второй лабораторной работы я изучил инструменты `gnuplot` и `perf`. `Gnuplot` оказался не таким уж и сложным, а скорее даже приятным для построения графиков и анализа сложностей алгоритмов. Надеюсь, я и дальше не поленюсь и буду пользоваться данным инструментом. `Perf` оказался сложной утилитой. Для того, чтобы сделать тривиальные вещи пришлось потратить какое-то время, однако эта сложность определяет мощь всего инструмента. Утилита показывает все узкие горлышки программы и анализирует вплоть до кода на ассемблере. Благо из-за простоты поиска в интернете, можно изучить, за что именно отвечает кусок кода и примерно понять, как оптимизировать код, если на то будут причины. `Valgrind` старый инструмент, которым я пользуюсь, а в данной лабораторной работе даже не оказалось ошибок утечек памяти, однако при написании второй лабораторной работы эта утилита сильно экономит время и помогает предотвратить ошибки, которые могут выявиться только спустя какое-то время. В моем опыте были ситуации, что программы, написанные на ОС Linux Ubuntu не происходило никаких ошибок, однако при портировании кода на Windows 10 они случались. Оказалось, что Ubuntu снисходительнее относится к использованию неинициализированной памяти и поэтому на ней сложнее было находить утечки, однако `Valgrind` позволяет мониторить данные ситуации и вовремя замечать ошибки. Так же в данной лабораторной работе пришлось вспомнить `bash`, который тоже бывает полезен при разработке сложных программ и тестов для них. Результаты эффективности программы также оказались неожиданностью. Я совсем не ожидал увидеть такую скорость вставки в AVL-дерево, однако это можно было предположить, так как на чекере моя программа на самом длительном 13 тесте прошла за 40 секунд, в то время как максимально на выполнение дается 130 секунд.

Список литературы

- [1] *Поисковик - Google.*
URL: <https://www.google.com/>
- [2] *Сайт с подробной документацией библиотек C++*
URL: <https://en.cppreference.com/>
- [3] *Сайт с постами про программирование (Perf)*
URL: <https://habr.com/ru/company/first/blog/442738/>
- [4] *Сайт с постами про программирование (Gnuplot)*
URL: <https://habr.com/ru/company/ruvds/blog/517450/>
- [5] *Еще немного про Perf*
URL: <https://selectel.ru/blog/perf-i-flamegraphs/>
- [6] *Намного больше про Perf*
URL: <https://www.protocols.ru/WP/perf/>
- [7] *Сайт с постами про программирование (Bash)*
URL: <https://habr.com/ru/company/ruvds/blog/325928/>