



**Exercises for *Foundations in Data Engineering*, WiSe 23/24**

Alexander Beischl, Maximilian Reif (i3fde@in.tum.de)

<http://db.in.tum.de/teaching/ws2324/foundationsde>

**Sheet Nr. 03**

**Exercise 1**

Show off how well you can handle bash, pipes and gnu tools by now!

Build a persistent key-value store using only these tools. It should support two actions: **set** and **get**.

**Set** takes a key and a value and stores it.

**Get** retrieves the value stored for a given key.

You may assume that keys do not contain a , (comma). Also, it is not necessary to achieve  $O(1)$  runtime for set and get as this task is about finding a very short solution. The functions can each be implemented in one line.

Once the implementation is done, do simple performance measurements:

- How long does it take to insert 1000 keys?
- How long does it take to retrieve 1000 keys afterwards?

**Solution:**

```
#!/usr/bin/env bash

set () {
    echo "$1,$2" >> database
}

get () {
    grep "^$1," database | tail -n 1 | sed -e "s/^$1,/"
}

$1 "$2" "$3"
```

Based On: Martin Kleppmann. "Designing Data-Intensive Applications."

Time measurements:

```
time (for i in $(seq 1 1000); do ./<script_name.sh> set $((
    $RANDOM)) $(( $RANDOM )); done;)
time (for i in $(seq 1 1000); do ./<script_name.sh> get $((
    $RANDOM)); done;) | wc -l
```

**Exercise 2**

On Linux, a plethora of useful information can be found in the system *manual pages* (*man pages*). They are especially relevant to a systems programmer, who often needs to interact with the kernel through system calls.

While man pages can of course also be found on the internet, we prefer to use the **man** executable to view man pages offline. That way, we can be sure that we are viewing the version of a man page that matches the programs installed on our system.

1. `man` does not display the text of the man page itself but it uses a helper program. Explain how `man` decides which helper program to use (Hint: Read the man page of `man`) in one sentence. How is the program called that is used on your system? Briefly describe your steps for finding the answer (1–2 sentences). List all commands you use.

**Solution:**

- `man` tries to use the `less` program by default.
- `man` uses the `-P/-pager` option or the `$MANPAGER` or
- `$PAGER` variables to determine which program to use.

Most systems usually use the program `less` to display man pages (although your system might be different). As man pages can be rather large, it is essential to be able to quickly search through man pages for relevant information.

2. Check the man page of `less`. How can you quickly search for all occurrences of a keyword (name the command)? How can you jump to specific lines in a file (name the command)? Briefly describe your steps for finding the answer (1–2 sentences).

**Solution:**

- The command `/<pattern>` can be used to search for a pattern. The `n` and `N` commands can be used to jump between occurrences.
- The `g` command followed by a line number can be used to jump to a specific line.

### Exercise 3

The `proc` filesystem (also called `procfs`) provides an interface to kernel data structures which can be queried for information about your system.

1. Use your knowledge about man pages and find a way to obtain information about your CPU with the `procfs`. Answer the following questions: Which file in the `procfs` contains the information about the available CPUs? Which helper program can be used to display the CPU infos in a nicer format?

**Solution:**

- The `/proc/cpuinfo` file contains information about the CPU.
- The `lscpu` helper program displays this information in a nicer format.

In particular, the kernel exposes CPU flags which may (e.g., on ARM) influence the behavior of the `g++` option `-march=native` which was briefly introduced in the lecture. This option potentially sets a large number of other `-m<option>` options specific to the current CPU.

2. Check the `g++` man page for information about the various `--help` options. Use this information to answer the following question: Which `-m<option>` options does `g++` set on your system when using `-march=native`? Briefly describe your steps for finding the answer (2–3 sentences, also list all commands that you used).

**Solution:**

- Check man page of `g++` to find that we want to use the `--help=target` option.
  - Check man page to find that the `g++ -Q --help=target -march=native` command lists the options that are set by the `-march=native` flag.
  - List the options that are set on the system.
3. Select 5 options that `g++` sets on your system when using `-march=native` from the last question. Briefly summarize their effect on the program that is produced by `g++` (one sentence per option).

**Solution: (Example, can vary on each machine)**

- `mavx2`: Use `avx2` instructions for vectorization.
  - `mcrc32`: This option enables built-in functions `__builtin_ia32_crc32qi`, `__builtin_ia32_crc32hi`, `__builtin_ia32_crc32si` and `__builtin_ia32_crc32di` to generate the `crc32` machine instruction.
  - More explanations can be found for example at: <https://gcc.gnu.org/onlinedocs/gcc/x86-Options.html>
4. If portability is an issue, it is usually not desirable to blindly apply the `-march=native` option during compilation. Explain why this is the case, and briefly outline a more flexible approach that leverages a build system (3–5 sentences).

**Solution:**

- Mention that `-march=native` sets options based on the current CPU, which could lead to portability problems when using a very recent CPU.
- A build system could be used to configure which CPU should be targeted.

**Exercise 4**

In the lecture, we introduced the C++ reference documentation at:

<https://en.cppreference.com/w/cpp>

This documentation should be the first place to go for any questions regarding the C++ language or standard library. Thus, being able to navigate and understand the reference documentation is an essential skill for a C++ programmer.

1. In the lecture, we introduced the value categories *lvalues* and *rvalues* but noted that this classification is inaccurate. Consult the reference documentation about value categories, and list the five value categories that C++ actually knows (no further explanation required).

**Solution:** `glvalues`, `prvalues`, `xvalues`, `lvalues`, `rvalues`

2. Consider the following code fragment.

```
for (unsigned i = 0, j = 0; i < 10; ++i, j *= 2) {
    /* do something */
}
```

Consult the reference documentation on `for` loops, and name the *syntactic* classification (e.g. statement, expression, ...) of the code fragments `unsigned i = 0, j = 0;` and `++i, j *= 2` in one sentence.

In one of these code fragments, the comma (,) is an operator. Explain in which fragment this is the case and briefly outline the semantics of the comma operator (3–5 sentences).

**Solution:** `unsigned i = 0, j = 0;` is a (declaration) statement, `++i, j *= 2` is an expression. In the second fragment the comma is an operator. The comma operator first evaluates the left operand and then discards its value. Then it evaluates the second operand. The result of the comma operator is the result of the second operand. The first expression is sequenced before the second.

3. Check the reference documentation on *implicit conversions* and briefly explain the key difference between numeric promotion and numeric conversion (3–5 sentences). Give two examples for each of these two conversions.

**Solution:**

- In numeric promotion integral types are implicitly converted to other integral types that can hold the initial value. This means that integral promotion usually leads to types that are larger.

Example

```
short i = 1;
int j = 2;
i + j; // i is promoted to int
```

- Numeric conversion also changes values from one integral to another integral type but may lose information by converting to a type that cannot hold the original value.

Example

```
int i = 1;
short j = i; // i is converted to short
```

4. Check the reference documentation on *I/O manipulators* and answer the following questions with a short example (1 sentence each):

- How can we print hexadecimal numbers?
- How can we print numbers with a fixed width and leading zeros?
- How can we right align numbers with a fixed width?

**Solution:**

- Print hexadecimal numbers: `std::hex`
- Fixed width: `std::setw`
- Leading zeros: `std::setfill`
- Right align: `std::right`