

C Programming II

2024 Spring

Homework 02

Instructor: Po-Wen Chi

Due: 2024.xx.xx PM 11:59

Policies:

- **Zero tolerance** for late submission.
- **Plagiarism is not allowed.** Both source and copycat will be **zero**.
- You need to prepare a README file about how to make and run your program. Moreover, you need to provide your name and your student ID in the README file.
 - Your Name and Your ID.
 - The functional description for each code.
 - Anything special.
- Please pack all your submissions in one zip file.
- For convenience, your executable programs must be named following the rule hw**XXYY**, where the red part is the homework number and the blue part is the problem number. For example, **hw0102** is the executable program for homework #1 problem 2.
- I only accept **PDF**. MS Word is not allowed.
- **Do not forget your Makefile. For convenience, each assignment needs only one Makefile.**

1 Subtitle Player (20 pts)

When you watch an English video, do you need its subtitle? This time, I want you to develop a subtitle player. The subtitle format we use in this problem is **Advanced SubStation Alpha (ASS)** . Please see the following link for its detail.

[illegible]

I will also provide you a subtitle example on the website. The program usage is as follows.

```

1 $ ./hw0201
2 Please enter the file name: test.ass
3 Time Shift ( -10 ~ 10 ): 0           // <-- integer
4 Speed (0.25,0.5,0.75,1,1.25,1.5,1.75,2): 1 // <-- only support these
   candidate

```

The subtitle is the example on the website. After the above input and 25.06 seconds, your should output the following string.

```

1 ... Three minutes.

```

After 3.8 seconds, please **clear the screen**. Then, after 0.008 second, you should print the following string.

```

1 Oui, Chef. Turn on the stove. Come on. Hurry up. Oui, Chef.

```

Your program should repeat this operation till the end of file. When printing strings, do not forget to print them with **colors**¹.

2 Premier League (20 pts)

Given a Premier League match data for some season, please print the final table of that season. You can download data from the following site.

<https://github.com/datasets/football-datasets/tree/master/datasets/premier-league>

```

1 ./hw0202
2 Please enter the data file name: season-1819.csv
3      Team                W    D    L    GF    GA    GD    Pts
4 01) Manchester City FC   32    2    4   95   23   +72   98
5 02) Liverpool FC        30    7    1   89   22   +67   97
6 03) Chelsea FC          21    9    8   63   39   +24   72
7 ...
8 19) Fulham FC            7     5   26   34   81   -47   26
9 20) Huddersfield Town FC  3     7   28   22   76   -54   16

```

Note that your output must be aligned.

3 Wordle Solver (20 pts)

Do you know what **Wordle** is? If not, please access the following link and play the Game.

<https://www.nytimes.com/games/wordle/index.html>

As you know, my English is poor. Therefore I want you to develop a tool for me to win this game. I will show you what I want.

1. First, you need an English dictionary. You can install **hunspell-en-us** in Ubuntu and get the English dictionary **en_US.dic**. Or you can directly download it from the following link.

<https://cgkit.freedesktop.org/libreoffice/dictionaries/tree/en>

2. Find all possible 5-letter words from the dictionary.

¹According to the PrimaryColour defined in Styles.

3. Propose the most possible candidate word to the user.
 - You must use the letter frequency table from the following link to calculate the summation of 5-letter frequency. You should use the text one. The 5-letter word in the candidate with the highest frequency is the most possible one.
https://en.wikipedia.org/wiki/Letter_frequency
4. Let the user to input the feedback from the Wordle site. The feedback will be "XXXXX" where X can be:
 - G: in the word and in the correct spot.
 - Y: in the word but in the wrong spot.
 - B: not in the word in any spot.

For your convenience,

- You do not need to provide the dictionary file when submitting your homework. Our TAs will put a dictionary file in your program directory and you can simply open the dictionary file directly in your program.
- Undoubtedly, there is no complete dictionary in this world. In this problem, you can just treat the dictionary as the complete one.
- You do not need to care about any variations of the word.
- The input should be case insensitive.

The program interface is as follows. Note that this is from my playing record and I do not do frequency analysis when playing.

```
1 ./hw0203
2 Advice:  AUDIO    // <-- Advice from your program
3 Feedback: YBYBB   // <-- User input
4 Advice:  SHADE    // <-- Advice from your program
5 Feedback: GGGGG   // <-- User input
6 Congratulations!!
```

4 BMP (20 pts)

I want you to develop a program to distort a BMP file with a given angle.

```
1 $ ./hw0204
2 Please input a BMP file: doraemon.bmp
3 Please input the output BMP file name: doraemon_out.bmp
4 Angle (0-90): 45
```

The distortion is as in figure 1.

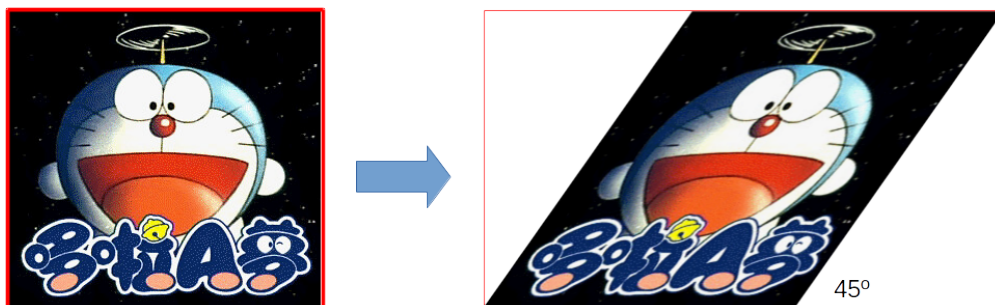


Figure 1: Distortion Example. Note that the boundary red line should not be presented. The additional part should be colored with white.

5 Large Language Model (20 pts)

In 2022, ChatGPT garnered widespread attention; in 2023, following the release of the foundational model **LLaMa**, thousands of Large Language Models (LLMs) emerged, signaling the rise of generative AI. This trend has continued into 2024.

While the global community races to develop increasingly powerful models reliant on massive GPU clusters with thousands of GPU, a segment of researchers and the open-source community is taking a different approach. Their goal is to make large language models accessible on a wide range of devices, including smartphones, with or without specialized acceleration hardware.

A standout project in this effort is [LLaMa.cpp](#). This project utilizes CPU acceleration technologies (e.g., NEON, AVX) and GPU acceleration frameworks (e.g., Metal, CUDA) to facilitate model inference on devices like laptops and smartphones. It employs the GGML compute library at its core, which is a pure C library designed for executing machine learning computations efficiently on CPU, GPU, and other devices.

LLaMa.cpp uses a special format to store the metadata and the model weights, called GGUF. Fig. 2 is a beautiful visualization of GGUF v3, provided in the official documentation.

Don't worry, I am not asking you to build an inference engine in pure C. (Some people really do so!) Instead, you are going to implement a GGUF file reader which is able to read the **model.gguf** in the current directory and output the metadata, the tensors, and the total parameters (summed up by all tensors) of the model.

The structure of GGUF is as simple as what you see in the above image:

```

1 struct gguf_file_t {
2     // The header of the file.
3     gguf_header_t header;
4
5     // Tensor infos, which can be used to locate the tensor data.
6     gguf_tensor_info_t tensor_infos[header.tensor_count];
7
8     // Other data that you don't need to worry about this time.
9     uint8_t other_binary_things[];
10 };

```

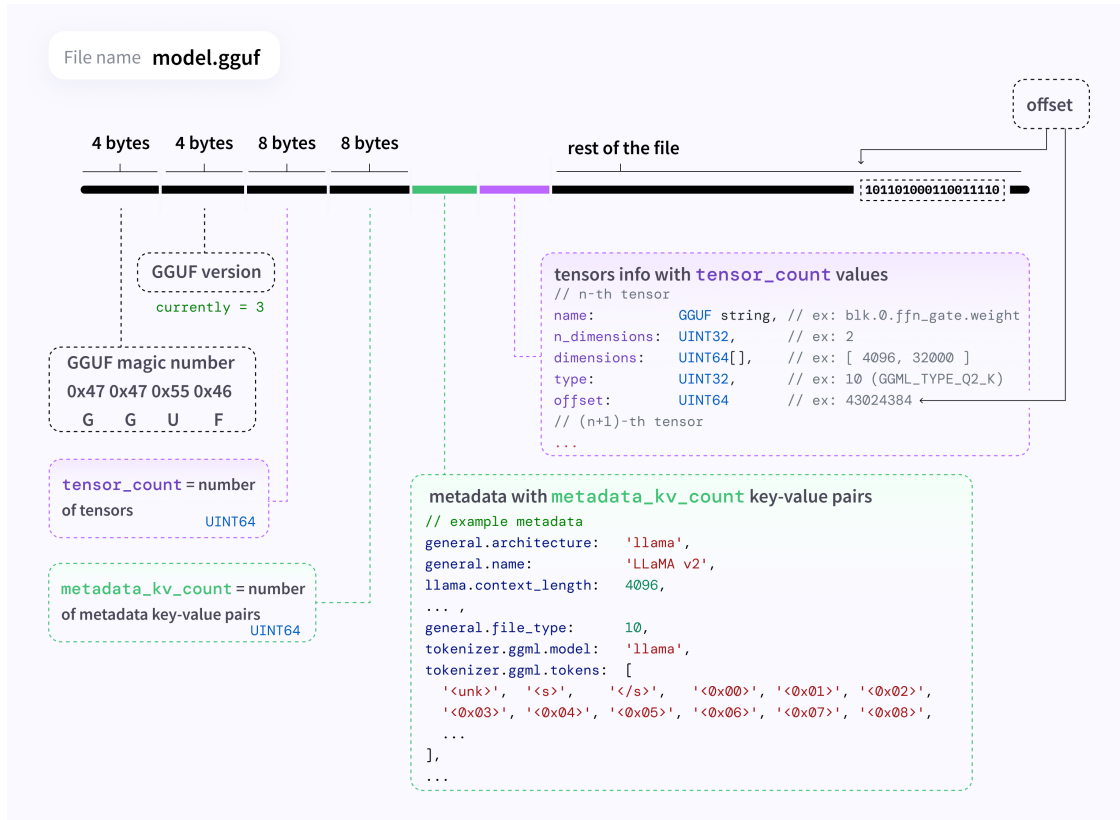


Figure 2: Model visualization: model.gguf

Then, what are `gguf_header_t` and `gguf_tensor_info_t`? Here are the structs that are used internally in GGML, where you can see the real usage of `union` and `enum` that you are familiar with!

```
1 struct gguf_header_t {
2     // Magic number to announce that this is a GGUF file.
3     // Must be `GGUF` at the byte level: `0x47` `0x47` `0x55` `0x46`.
4     // Your executor might do little-endian byte order, so it might be
5     // check for 0x46554747 and letting the endianness cancel out.
6     // Consider being *very* explicit about the byte order here.
7     uint32_t magic;
8     // The version of the format implemented.
9     // Must be `3` for version described in this spec, which introduces big-
10    endian support.
11    //
12    // This version should only be increased for structural changes to the
13    // format.
14    // Changes that do not affect the structure of the file should instead
15    // update the metadata
16    // to signify the change.
17    uint32_t version;
18    // The number of tensors in the file.
19    // This is explicit, instead of being included in the metadata, to ensure
20    // it is always present
```

```

17 // for loading the tensors.
18 uint64_t tensor_count;
19 // The number of metadata key-value pairs.
20 uint64_t metadata_kv_count;
21 // The metadata key-value pairs.
22 gguf_metadata_kv_t metadata_kv[metadata_kv_count];
23 };
24
25 struct gguf_tensor_info_t {
26     // The name of the tensor. It is a standard GGUF string, with the caveat
    that
27     // it must be at most 64 bytes long.
28     gguf_string_t name;
29     // The number of dimensions in the tensor.
30     // Currently at most 4, but this may change in the future.
31     uint32_t n_dimensions;
32     // The dimensions of the tensor.
33     uint64_t dimensions[n_dimensions];
34     // The type of the tensor.
35     ggml_type type;
36     // The offset of the tensor's data in this file in bytes.
37     //
38     // This offset is relative to `tensor_data`, not to the start
39     // of the file, to make it easier for writers to write the file.
40     // Readers should consider exposing this offset relative to the
41     // file to make it easier to read the data.
42     //
43     // Must be a multiple of `ALIGNMENT`. That is, `align_offset(offset) ==
    offset`.
44     uint64_t offset;
45 };
46
47 struct gguf_metadata_kv_t {
48     // The key of the metadata. It is a standard GGUF string, with the
    following caveats:
49     // - It must be a valid ASCII string.
50     // - It must be a hierarchical key, where each segment is `
    lower_snake_case` and separated by a `.`.
51     // - It must be at most 216-1/65535 bytes long.
52     // Any keys that do not follow these rules are invalid.
53     gguf_string_t key;
54
55     // The type of the value.
56     // Must be one of the `gguf_metadata_value_type` values.
57     gguf_metadata_value_type value_type;
58     // The value.
59     gguf_metadata_value_t value;
60 };
61
62 // A string in GGUF.
63 struct gguf_string_t {
64     // The length of the string, in bytes.
65     uint64_t len;
66     // The string as a UTF-8 non-null-terminated string.

```

```

67     char string[len];
68 }
69
70 union gguf_metadata_value_t {
71     uint8_t uint8;
72     int8_t int8;
73     uint16_t uint16;
74     int16_t int16;
75     uint32_t uint32;
76     int32_t int32;
77     float float32;
78     uint64_t uint64;
79     int64_t int64;
80     double float64;
81     bool bool_;
82     gguf_string_t string;
83     struct {
84         // Any value type is valid, including arrays.
85         gguf_metadata_value_type type;
86         // Number of elements, not bytes
87         uint64_t len;
88         // The array of values.
89         gguf_metadata_value_t array[len];
90     } array;
91 };
92
93 enum ggml_type: uint32_t {
94     GGML_TYPE_F32 = 0,
95     GGML_TYPE_F16 = 1,
96     GGML_TYPE_Q4_0 = 2,
97     GGML_TYPE_Q4_1 = 3,
98     // GGML_TYPE_Q4_2 = 4, support has been removed
99     // GGML_TYPE_Q4_3 (5) support has been removed
100     GGML_TYPE_Q5_0 = 6,
101     GGML_TYPE_Q5_1 = 7,
102     GGML_TYPE_Q8_0 = 8,
103     GGML_TYPE_Q8_1 = 9,
104     // k-quantizations
105     GGML_TYPE_Q2_K = 10,
106     GGML_TYPE_Q3_K = 11,
107     GGML_TYPE_Q4_K = 12,
108     GGML_TYPE_Q5_K = 13,
109     GGML_TYPE_Q6_K = 14,
110     GGML_TYPE_Q8_K = 15,
111     GGML_TYPE_I8,
112     GGML_TYPE_I16,
113     GGML_TYPE_I32,
114     GGML_TYPE_COUNT,
115 };
116
117 enum gguf_metadata_value_type: uint32_t {
118     // The value is a 8-bit unsigned integer.
119     GGUF_METADATA_VALUE_TYPE_UINT8 = 0,
120     // The value is a 8-bit signed integer.

```

```

121 GGUF_METADATA_VALUE_TYPE_INT8 = 1,
122 // The value is a 16-bit unsigned little-endian integer.
123 GGUF_METADATA_VALUE_TYPE_UINT16 = 2,
124 // The value is a 16-bit signed little-endian integer.
125 GGUF_METADATA_VALUE_TYPE_INT16 = 3,
126 // The value is a 32-bit unsigned little-endian integer.
127 GGUF_METADATA_VALUE_TYPE_UINT32 = 4,
128 // The value is a 32-bit signed little-endian integer.
129 GGUF_METADATA_VALUE_TYPE_INT32 = 5,
130 // The value is a 32-bit IEEE754 floating point number.
131 GGUF_METADATA_VALUE_TYPE_FLOAT32 = 6,
132 // The value is a boolean.
133 // 1-byte value where 0 is false and 1 is true.
134 // Anything else is invalid, and should be treated as either the model
    being invalid or the reader being buggy.
135 GGUF_METADATA_VALUE_TYPE_BOOL = 7,
136 // The value is a UTF-8 non-null-terminated string, with length prepended.
137 GGUF_METADATA_VALUE_TYPE_STRING = 8,
138 // The value is an array of other values, with the length and type
    prepended.
139 ///
140 // Arrays can be nested, and the length of the array is the number of
    elements in the array, not the number of bytes.
141 GGUF_METADATA_VALUE_TYPE_ARRAY = 9,
142 // The value is a 64-bit unsigned little-endian integer.
143 GGUF_METADATA_VALUE_TYPE_UINT64 = 10,
144 // The value is a 64-bit signed little-endian integer.
145 GGUF_METADATA_VALUE_TYPE_INT64 = 11,
146 // The value is a 64-bit IEEE754 floating point number.
147 GGUF_METADATA_VALUE_TYPE_FLOAT64 = 12,
148 }

```

The output format matters. For some long fields, please ensure they fit into an **80-column width**.

```

1 $ ./hw0205
2 GGUF: true
3 Parameters: 1,100,048,384
4
5 Metadata                                Value
6 version                                2
7 tensor_count                            201
8 kv_count                                21
9 general.architecture                    llama
10 general.name                            models
11 general.file_type                       2
12 general.quantization_version            2
13 llama.context_length                    2048
14 llama.embedding_length                  2048
15 llama.block_count                       22
16 llama.feed_forward_length               5632
17 llama.rope.dimension_count              64
18 llama.rope.freq_base                    10000
19 llama.attention.head_count              32
20 llama.attention.head_count_kv           4

```



```

21 llama.attention.layer_norm_rms_epsilon 0.000009999999747378752
22 tokenizer.ggml.model llama
23 tokenizer.ggml.tokens [<unk>, <s>, </s>, <0x00>, <0x01>,
    ...]
24 tokenizer.ggml.scores [0, 0, 0, 0, 0, ...]
25 tokenizer.ggml.token_type [2, 3, 3, 6, 6, ...]
26 tokenizer.ggml.bos_token_id 1
27 tokenizer.ggml.eos_token_id 2
28 tokenizer.ggml.padding_token_id 2
29
30 Tensors Shape Precision
31 blk.0
32   .attn_k.weight [2048,256] Q4_0
33   .attn_norm.weight [2048] F32
34   .attn_output.weight [2048,2048] Q4_0
35   .attn_q.weight [2048,2048] Q4_0
36   .attn_v.weight [2048,256] Q4_0
37   .ffn_down.weight [5632,2048] Q4_0
38   .ffn_gate.weight [2048,5632] Q4_0
39   .ffn_norm.weight [2048] F32
40   .ffn_up.weight [2048,5632] Q4_0
41 blk.1
42   .attn_k.weight [2048,256] Q4_0
43   .attn_norm.weight [2048] F32
44   .attn_output.weight [2048,2048] Q4_0
45   .attn_q.weight [2048,2048] Q4_0
46   .attn_v.weight [2048,256] Q4_0
47   .ffn_down.weight [5632,2048] Q4_0
48   .ffn_gate.weight [2048,5632] Q4_0
49   .ffn_norm.weight [2048] F32
50   .ffn_up.weight [2048,5632] Q4_0
51 ...
52 blk.21
53   .attn_k.weight [2048,256] Q4_0
54   .attn_norm.weight [2048] F32
55   .attn_output.weight [2048,2048] Q4_0
56   .attn_q.weight [2048,2048] Q4_0
57   .attn_v.weight [2048,256] Q4_0
58   .ffn_down.weight [5632,2048] Q4_0
59   .ffn_gate.weight [2048,5632] Q4_0
60   .ffn_norm.weight [2048] F32
61   .ffn_up.weight [2048,5632] Q4_0
62 output.weight [2048,32000] Q6_K
63 output_norm.weight [2048] F32
64 token_embd.weight [2048,32000] Q4_0

```

If the model.gguf doesn't exist or is not a valid GGUF file (verified by checking the magic number), you only need to print the first line:

```

1 $ ./hw0205
2 GGUF: false

```

You can find more [GGUF example files on HuggingFace](#). Additionally, you can use the tensor viewer to check the correctness of your work. Here is an [example for the above-mentioned TinyLLaMa model](#).

Good Luck!

PS: The source of the structures used in this problem is in the following link.

<https://github.com/ggerganov/ggml/blob/f5c9599cdba3133da0158dce061b33413b49f6fd/src/ggml.c#L20262-L20366>

Follow up: Some models are VERY large (>20GB per file), can your program read them without OOM?

6 Bonus: Where is `errno` (5 pts)

In this class, I have shown you how error number works. As you can see, this is an **extern** variable. Please answer the following questions.

1. Where is **`errno`**?
2. Please find all **known** error numbers. Do not search them from Google but find them from **codes**.
3. Please define a new error number and use **`perror`** to show it.