



گزارش پروژه پایانی ساختار زبان و کامپیوتر

طراحی برنامه Convolution-2D

علی فتحی

فهرست مطالب

۱. [مقدمه](#)
۲. [ساختار کلی برنامه](#)
۳. [نحوه اجرای برنامه \(Interactive و CLI\)](#)
۴. [ورودی‌های برنامه](#)
۵. [طراحی Kernel Factory](#)
۶. [خواندن و ذخیره تصویر با OpenCV](#)
۷. [ساختار داده‌ها \(Struct ها\)](#)
۸. [موتورهای Convolution](#)
۹. [توضیح پیاده‌سازی SIMD](#)
۱۰. [تصاویری از اجرای برنامه](#)
۱۱. [معرفی نوت‌بوک‌های مکمل](#)
 - [speed-comparison](#)
 - [object-detection](#)

• مقدمه:

در این گزارش به بررسی برنامه Conv2D می‌پردازیم، به ساختار برنامه، نحوه اجرا، امکانات و محدودیت‌های آن اشاره می‌کنیم و در آخر چند نوت‌بوک تکمیلی جهت بررسی بیشتر نتایج حاصل شده و استفاده از این برنامه در زبان‌های سطح بالا، معرفی می‌کنیم.

• ساختار کلی برنامه:

این برنامه، می‌تواند به صورت تعاملی و یا از طریق آرگومان‌های رابط خط فرمان، از کاربر ورودی‌های مورد نیاز، از جمله آدرس عکس ورودی، فیلتر انتخابی، نوع پردازش رنگ تصویر (grayscale و یا rgb) و ... را گرفته و عمل convolution را به کمک یکی از ۳ موتور (engine) اصلی خود انجام می‌دهد. همچنین در صورت درخواست کاربر، خروجی را در آدرس تعیین شده ذخیره می‌کند.

موتورهای موجود در برنامه عبارتند از:

- Baseline: کد ساده C++ که به صورت scalar عمل می‌کند
 - SSE: استفاده از رجیسترهای SSE (۱۲۸ بیتی) پردازنده x86
 - AVX: استفاده از رجیسترهای AVX (۲۵۶ بیتی) پردازنده x86
- از آنجایی که این برنامه با داده‌هایی از نوع float32 کار می‌کند، رجیسترهای SSE می‌توانند ۴ پردازش، و رجیسترهای AVX، ۸ پردازش را به طور همزمان انجام دهند.

• نحوه اجرای برنامه:

در این بخش توابعی را مشاهده می‌کنید که در ابتدای اجرای برنامه، نوع اجرای برنامه را (Interactive و یا CLI) مشخص کرده و طبق آن توابع مربوطه بعدی را صدا می‌زنند.

سپس اگر برنامه به صورت تعاملی اجرا شده باشد، از کاربر درخواست می‌شود مقادیر ورودی‌های موردنظر را وارد کند، و در غیر این صورت از آرگومان‌های رابط خط فرمان استفاده می‌شود.

همانطور که مشاهده می‌کنید، آرگومان‌های CLI باید به صورت زیر به برنامه داده شوند:

```
1 static int read_input_and_run_functional_test() {
2     int res = CODE_SUCCESS;
3
4     FunctionalTestParams params;
5
6     res = read_functional_test_input(params);
7     if (res != CODE_SUCCESS)
8         return res;
9
10    res = run_functional_test(params);
11    return res;
12 }
13
14 static int read_input_and_run_speed_test() {
15     int res = CODE_SUCCESS;
16
17     SpeedTestParams params;
18
19     res = read_speed_test_params(params);
20     if (res != CODE_SUCCESS)
21         return res;
22
23     res = run_speed_test(params);
24     return res;
25 }
26
27 static int run_interactive() {
28     int res = CODE_SUCCESS;
29
30     int option_def = RUN_MODE_SPEED_TEST;
31
32     int option;
33
34     OptionEntry options[2];
35     options[0].option_number = RUN_MODE_FUNCTIONAL_TEST;
36     options[0].option_name = RUN_MODE_FUNCTIONAL_TEST_STR;
37     options[1].option_number = RUN_MODE_SPEED_TEST;
38     options[1].option_name = RUN_MODE_SPEED_TEST_STR;
39
40     res = read_option("Option", options, 2, stdin, &option_def, &option);
41     if (res != CODE_SUCCESS) {
42         print_err("Failed to read option", res);
43         return res;
44     }
45
46     if (option == RUN_MODE_FUNCTIONAL_TEST)
47         res = read_input_and_run_functional_test();
48     else if (option == RUN_MODE_SPEED_TEST)
49         res = read_input_and_run_speed_test();
50
51     return res;
52 }
53
54 static int run_cli(int argc, char **argv) {
55     int res = CODE_SUCCESS;
56
57     CLIArgs args;
58
59     res = parse_cli(argc, argv, args);
60     if (res != CODE_SUCCESS || args.help) {
61         print_help();
62         return res;
63     }
64
65     res = validate_cli(args);
66     if (res != CODE_VALIDATION_OK) {
67         return res;
68     }
69
70     if (args.run_mode == RUN_MODE_FUNCTIONAL_TEST) {
71         FunctionalTestParams params = {
72             args.engine_mode,
73             args.kernel_type,
74             args.kernel_size,
75             args.color_mode,
76             args.input,
77             args.output,
78             args.save_output
79         };
80         res = run_functional_test(params);
81     }
82     else if (args.run_mode == RUN_MODE_SPEED_TEST) {
83         SpeedTestParams params = {
84             args.engine_mode,
85             args.kernel_type,
86             args.kernel_size,
87             args.color_mode,
88             args.input,
89             args.output,
90             args.save_output
91         };
92         res = run_speed_test(params);
93     }
94
95     return res;
96 }
97
98 int main(int argc, char **argv) {
99     int res = CODE_SUCCESS;
100
101     if (argc == 1) res = run_interactive();
102     else
103         res = run_cli(argc, argv);
104
105     return res;
106 }
```

```
1 void print_help() {
2     std::cout <<
3     "Usage: conv2d [OPTIONS]\n\n"
4     "Modes:\n"
5     "  -m --mode functional | speed\n\n"
6     "Options:\n"
7     "  -e, --engine baseline | sse | avx\n"
8     "  -k, --ktype sharpen | box_blur | gaussian_blur | sobel_x | sobel_y\n"
9     "  -s, --ksize kernel size (e.g. 3)\n"
10    "  -i, --input input file or directory\n"
11    "  -o, --output output file or directory (optional)\n"
12    "  -c, --color grayscale | rgb (optional) (default = rgb)\n"
13    "  -h, --help show this help\n";
14 }
```

• ورودی‌های برنامه:

- **Mode:** در حالت functional تنها یک عکس به عنوان ورودی داده می‌شود. هدف این مود تست عملکرد برنامه و مشاهده تاثیر فیلتر انتخاب شده روی عکس ورودی می‌باشد.
مود speed برای تست سرعت برنامه است.
 - **Engine:** موتور اصلی برنامه که نحوه پردازش عمل convolution را مشخص می‌کند. می‌تواند به صورت baseline، و یا با استفاده از رجیسترهای SSE و AVX پردازنده x86 باشد.
 - **Kernel Type:** این آرگومان نوع کرنل (فیلتر) را انتخاب می‌کند. برای مثال می‌توان از فیلترهای blur، sharpen و ... استفاده کرد.
 - **Kernel Size:** این آرگومان سایز کرنل را انتخاب می‌کند. در حال حاضر این مقدار می‌تواند ۳، ۵ و یا ۷ باشد (به نوع کرنل نیز وابسته است)
 - **Input/Output:** آدرس فایل یا پوشه ورودی/خروجی. در مود functional آدرس فایل ورودی و آدرس فایل خروجی برای ذخیره عکس داده می‌شود و در مود speed، یک پوشه حاوی عکس به عنوان ورودی داده می‌شود و یک پوشه برای ذخیره عکس‌های خروجی.
 - **Color:** تعداد کانال‌های پردازش را مشخص می‌کند، در حالت grayscale، عکس‌ها به صورت تک‌کاناله بارگذاری و پردازش می‌شوند و در حالت rgb، به صورت ۳ کاناله.
- در اجرای Interactive نیز همین ورودی‌ها درخواست می‌شوند.

• Kernel Factory:

در این تابع، با توجه به نوع و اندازه کرنل، یک کرنل پیش فرض انتخاب شده و یا یک کرنل مناسب ساخته می شود. اگر در آینده قصد اضافه کردن فیلترهای بیشتری به برنامه داشته باشیم، می توانیم از طریق همین تابع، فیلترهای دلخواه خود را اضافه کنیم.

```
1 int get_kernel(int kernel_type, int kernel_size, Kernel& kernel) {
2
3     if (!validate_odd_kernel(kernel_size)) {
4         print_err("Kernel size must be odd and >= 3", CODE_FAILURE_NOT_SUPPORTED);
5         return CODE_FAILURE_NOT_SUPPORTED;
6     }
7
8     kernel.size = kernel_size;
9     kernel.type = kernel_type;
10    kernel.data = new float[kernel_size * kernel_size];
11
12    int status = CODE_SUCCESS;
13
14    switch (kernel_type) {
15
16        case KERNEL_TYPE_SHARPEN:
17            status = create_sharpen(kernel);
18            break;
19
20        case KERNEL_TYPE_SOBEL_X:
21            status = create_sobel_x(kernel);
22            break;
23
24        case KERNEL_TYPE_SOBEL_Y:
25            status = create_sobel_y(kernel);
26            break;
27
28        case KERNEL_TYPE_BOX_BLUR:
29            status = create_box_blur(kernel);
30            break;
31
32        case KERNEL_TYPE_GAUSSIAN_BLUR:
33            status = create_gaussian_blur(kernel);
34            break;
35
36        default:
37            print_err("Unsupported kernel type", CODE_FAILURE_NOT_SUPPORTED);
38            status = CODE_FAILURE_NOT_SUPPORTED;
39            break;
40    }
41
42    if (status != CODE_SUCCESS) {
43        delete[] kernel.data;
44        kernel.data = nullptr;
45    }
46
47    return status;
48 }
```

- خواندن و ذخیره عکس (OpenCV2):

برای خواندن عکس‌ها و تبدیل آن‌ها به فرمت * float، و همچنین تبدیل آرایه float به عکس و ذخیره آن، از کتابخانه OpenCV استفاده شده است. اگر عکس به صورت grayscale باشد، تنها یک آرایه از float (یک کانال) توسط توابع این بخش ایجاد و یا دریافت می‌شود، و اگر به صورت rgb باشد، ۳ آرایه (البته می‌توانیم مقادیر هر ۳ آرایه را در یک آرایه و به صورت scalar پاس بدهیم).

در اینجا قسمتی از توابع مربوط به خواندن و نوشتن عکس‌ها را مشاهده می‌کنید:

```
1 int load_grayscale_image(  
2     const char *filename,  
3     Image& image  
4 ) {  
5     cv::Mat img = cv::imread(filename, cv::IMREAD_GRAYSCALE);  
6  
7     if (img.empty()) {  
8         print_err("Failed to load image", CODE_FAILURE_READ_INPUT);  
9         return CODE_FAILURE;  
10    }  
11  
12    cv::Mat img_f;  
13    img.convertTo(img_f, CV_32F, 1.0 / 255.0);  
14  
15    image.height = img_f.rows;  
16    image.width = img_f.cols;  
17    image.channels = CHANNELS_GRAYSCALE;  
18  
19    image.data = new float[image.height * image.width];  
20  
21    for (int i = 0; i < image.height; i++) {  
22        const float *row_ptr = img_f.ptr<float>(i);  
23        std::memcpy(  
24            image.data + i * image.width,  
25            row_ptr,  
26            image.width * sizeof(float)  
27        );  
28    }  
29  
30    return CODE_SUCCESS;  
31 }
```

```
1 int save_float_array_as_grayscale_image(  
2     const char *filename,  
3     const Image& output  
4 ) {  
5     cv::Mat img_f(output.height, output.width, CV_32F, const_cast<float*>(output.data));  
6  
7     cv::Mat img_clamped;  
8     cv::min(img_f, 1.0f, img_clamped);  
9     cv::max(img_clamped, 0.0f, img_clamped);  
10  
11    cv::Mat img_u8;  
12    img_clamped.convertTo(img_u8, CV_8U, 255.0);  
13  
14    if (!cv::imwrite(filename, img_u8)) {  
15        print_err("Failed to write output image", CODE_FAILURE_WRITE_OUTPUT);  
16        return CODE_FAILURE;  
17    }  
18  
19    return CODE_SUCCESS;  
20 }
```


- Struct های به کار رفته در برنامه:

قبل از ورود به بخش محاسبات، نگاهی به چند استراکت مهم که برای خوانایی بیشتر کد در برنامه به کار برده شده‌اند می‌اندازیم:

- **Image:**

حاوی اطلاعات عکس، از جمله طول و عرض، تعداد کانال و اعداد نمایانگر پیکسل‌ها. بدیهی است که تعداد float32 هایی که در data قرار می‌گیرد از رابطه $\text{width} * \text{height} * \text{channels}$ حاصل می‌شود.

- **Kernel:**

مانند Image با این تفاوت که در کاربرد این برنامه از کرنل‌های تک کاناله استفاده می‌کنیم.

- **Conv2DParams:**

حاوی اطلاعاتی که برای عمل convolution به آن نیاز داریم. شایان ذکر است که در برنامه فعلی ما از stride ثابت 1 استفاده می‌شود.

```
1 struct Image {
2     float *data;
3     int height;
4     int width;
5     int channels;
6 };
7
8 struct Kernel {
9     float *data;
10    int type;
11    int size;
12 };
13
14 struct Conv2DParams {
15     Image image;
16     Kernel kernel;
17     int stride;
18 };
```

- **موتورهای Convolution:**

تابع `conv2d_channels`، همان تابعی است که برای اجرای عملیات کانولوشن آن را صدا می‌زنیم. با توجه به اینکه ممکن است تعداد متغیری کانال داشته باشیم، این تابع، عملیات را به طور جداگانه روی هر کانال اجرا می‌کند:

```
1 int conv2d_channels(  
2     int engine_mode,  
3     const Conv2DParams& params,  
4     Image& output  
5 ) {  
6  
7     int res = CODE_SUCCESS;  
8  
9     if (  
10        params.kernel.size != KERNEL_SIZE_3 &&  
11        engine_mode != ENGINE_MODE_BASELINE  
12    ) {  
13        engine_mode = ENGINE_MODE_BASELINE;  
14        print_warn("Only 3x3 kernels are supported in SSE/AVX engines. Falling back to baseline engine.");  
15    }  
16  
17    Image image = params.image;  
18    Kernel kernel = params.kernel;  
19  
20    int out_height = (image.height - kernel.size) / params.stride + 1;  
21    int out_width = (image.width - kernel.size) / params.stride + 1;  
22  
23    output.height = out_height;  
24    output.width = out_width;  
25    output.channels = image.channels;  
26    output.data = new float[output.height * output.width * output.channels];  
27  
28    for (int c = 0; c < output.channels; c++) {  
29        Conv2DParams ch_params = params;  
30        ch_params.image.data = channel_ptr(image, c);  
31  
32        Image ch_out;  
33        res = conv2d(engine_mode, ch_params, ch_out);  
34  
35        if (res != CODE_SUCCESS) {  
36            std::string err_msg = "Operation failed on channel " + std::to_string(c);  
37            print_err(err_msg.c_str(), res);  
38            return res;  
39        }  
40  
41        std::memcpy(  
42            channel_ptr(output, c),  
43            ch_out.data,  
44            output.height * output.width * sizeof(float)  
45        );  
46        delete[] ch_out.data;  
47    }  
48  
49    return res;  
51 }
```

در شبکه‌های عصبی پیچشی (CNN) معمولاً باید تعداد کانال‌های کرنل با تعداد کانال‌های ورودی آن لایه برابر باشد و خروجی هر عمل convolution تک‌کاناله است، در این کاربرد اگر بخواهیم n کانال

خروجی داشته باشیم (یا به عبارتی n ویژگی از عکس استخراج کنیم) باید n بار عمل convolution را با کرنل‌های متفاوت اجرا کنیم. اما از آنجایی که برنامه ما یک برنامه convolution کلاسیک است که تنها می‌خواهیم اثر فیلترهای گوناگون را روی عکس‌ها مشاهده و سرعت engine های مختلف را مقایسه کنیم، یک کرنل مشخص را روی هر n کانال عکس ورودی اعمال کرده و n کانال خروجی دریافت می‌کنیم. در تابع `conv2d`، بر اساس ورودی `engine_mode` انتخاب می‌کنیم که پردازش با کدام موتور صورت گیرد:

```
1 static int conv2d(  
2     int engine_mode,  
3     const Conv2DParams& params,  
4     Image& output  
5 ) {  
6     int res = CODE_SUCCESS;  
7  
8     res = conv2d_validation(  
9         engine_mode,  
10        params);  
11  
12    if (res != CODE_VALIDATION_OK) {  
13        print_err("Invalid arguments to function", res);  
14        return CODE_FAILURE_INVALID_ARG;  
15    }  
16  
17    switch(engine_mode) {  
18        case ENGINE_MODE_BASELINE:  
19            res = conv2d_baseline(  
20                params,  
21                output);  
22            break;  
23  
24        case ENGINE_MODE_SSE:  
25            res = conv2d_sse(  
26                params,  
27                output);  
28            break;  
29  
30        case ENGINE_MODE_AVX:  
31            res = conv2d_avx(  
32                params,  
33                output);  
34            break;  
35  
36        default: res = CODE_FAILURE;  
37    }  
38  
39    return res;  
40 }
```

حال، ۳ موتور مختلف برنامه را بررسی می‌کنیم:

```
1 static int conv2d_baseline(
2     const Conv2dParams& params,
3     Image& output
4 ) {
5
6     int res = CODE_SUCCESS;
7
8     Image image = params.image;
9     Kernel kernel = params.kernel;
10
11
12     int out_height = (image.height - kernel.size) / params.stride + 1;
13     int out_width = (image.width - kernel.size) / params.stride + 1;
14     int out_size = out_height * out_width;
15
16     output.height = out_height;
17     output.width = out_width;
18     output.data = new float[out_size];
19
20     for (int i = 0; i < out_height; i++) {
21         for (int j = 0; j < out_width; j++) {
22
23             float sum = 0.0f;
24
25             int base_i = i * params.stride;
26             int base_j = j * params.stride;
27
28             for (int u = 0; u < kernel.size; u++) {
29
30                 int img_row = (base_i + u) * image.width;
31                 int ker_row = u * kernel.size;
32
33                 for (int v = 0; v < kernel.size; v++) {
34                     int img_idx = img_row + (base_j + v);
35                     int ker_idx = ker_row + v;
36
37                     sum += image.data[img_idx] * kernel.data[ker_idx];
38                 }
39
40                 int out_idx = i * out_width + j;
41                 output.data[out_idx] = sum;
42             }
43         }
44     }
45
46     return res;
47 }
```

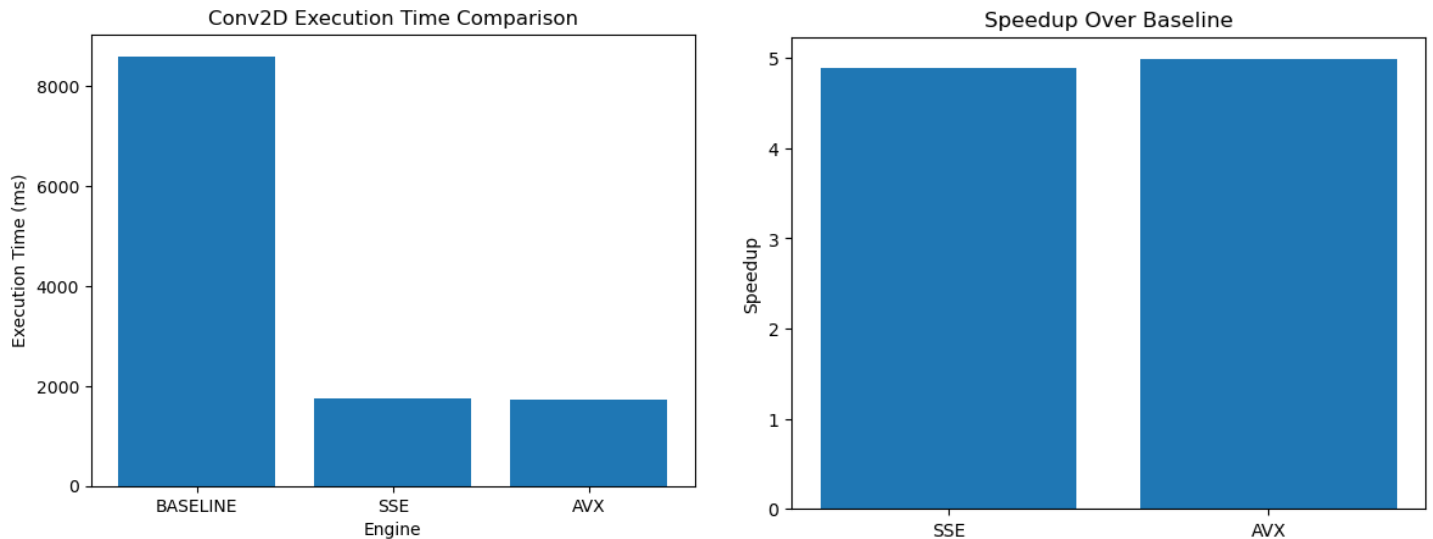
```
1 static int conv2d_sse(
2     const Conv2dParams& params,
3     Image& output
4 ) {
5
6     int res = CODE_SUCCESS;
7
8     Image image = params.image;
9     Kernel kernel = params.kernel;
10
11     int out_height = (image.height - kernel.size) / params.stride + 1;
12     int out_width = (image.width - kernel.size) / params.stride + 1;
13     int out_size = out_height * out_width;
14
15     output.height = out_height;
16     output.width = out_width;
17     output.data = new float[out_size];
18
19     __m128 k00 = _mm_setl_ps(kernel.data[0]);
20     __m128 k01 = _mm_setl_ps(kernel.data[1]);
21     __m128 k02 = _mm_setl_ps(kernel.data[2]);
22
23     __m128 k10 = _mm_setl_ps(kernel.data[3]);
24     __m128 k11 = _mm_setl_ps(kernel.data[4]);
25     __m128 k12 = _mm_setl_ps(kernel.data[5]);
26
27     __m128 k20 = _mm_setl_ps(kernel.data[6]);
28     __m128 k21 = _mm_setl_ps(kernel.data[7]);
29     __m128 k22 = _mm_setl_ps(kernel.data[8]);
30
31     consteval int SSE_FLOATS = 4;
32
33     for (int i = 0; i < out_height; i++) {
34         int base_i = i * params.stride;
35
36         int j = 0;
37         for (; j < out_width; j += SSE_FLOATS) {
38             int base_j = j * params.stride;
39
40             __m128 sum = _mm_setzero_ps();
41
42             // Row 0
43             __m256 r0 = _mm256_loadu_ps(image.data[(base_i + 0) * image.width + base_j + 0]);
44             __m256 r1 = _mm256_loadu_ps(image.data[(base_i + 0) * image.width + base_j + 1]);
45             __m256 r2 = _mm256_loadu_ps(image.data[(base_i + 0) * image.width + base_j + 2]);
46
47             sum = _mm256_add_ps(sum, _mm256_mul_ps(r0, k00));
48             sum = _mm256_add_ps(sum, _mm256_mul_ps(r1, k01));
49             sum = _mm256_add_ps(sum, _mm256_mul_ps(r2, k02));
50
51             // Row 1
52             r0 = _mm256_loadu_ps(image.data[(base_i + 1) * image.width + base_j + 0]);
53             r1 = _mm256_loadu_ps(image.data[(base_i + 1) * image.width + base_j + 1]);
54             r2 = _mm256_loadu_ps(image.data[(base_i + 1) * image.width + base_j + 2]);
55
56             sum = _mm256_add_ps(sum, _mm256_mul_ps(r0, k10));
57             sum = _mm256_add_ps(sum, _mm256_mul_ps(r1, k11));
58             sum = _mm256_add_ps(sum, _mm256_mul_ps(r2, k12));
59
60             // Row 2
61             r0 = _mm256_loadu_ps(image.data[(base_i + 2) * image.width + base_j + 0]);
62             r1 = _mm256_loadu_ps(image.data[(base_i + 2) * image.width + base_j + 1]);
63             r2 = _mm256_loadu_ps(image.data[(base_i + 2) * image.width + base_j + 2]);
64
65             sum = _mm256_add_ps(sum, _mm256_mul_ps(r0, k20));
66             sum = _mm256_add_ps(sum, _mm256_mul_ps(r1, k21));
67             sum = _mm256_add_ps(sum, _mm256_mul_ps(r2, k22));
68
69             __m256_storeu_ps(output.data[(i * out_width + j) * SSE_FLOATS], sum);
70
71             for (; j < out_width; j++) {
72                 int base_j = j * params.stride;
73                 float s = 0.0f;
74                 for (int ky = 0; ky < 3; ky++)
75                     for (int kx = 0; kx < 3; kx++)
76                         s += image.data[(base_i + ky) * image.width + (base_j + kx)] *
77                            kernel.data[ky * 3 + kx];
78                 output.data[(i * out_width + j) * SSE_FLOATS] = s;
79             }
80         }
81     }
82
83     return res;
84 }
```

```
1 static int conv2d_avx(
2     const Conv2dParams& params,
3     Image& output
4 ) {
5
6     int res = CODE_SUCCESS;
7
8     Image image = params.image;
9     Kernel kernel = params.kernel;
10
11     int out_height = (image.height - kernel.size) / params.stride + 1;
12     int out_width = (image.width - kernel.size) / params.stride + 1;
13     int out_size = out_height * out_width;
14
15     output.height = out_height;
16     output.width = out_width;
17     output.data = new float[out_size];
18
19     __m256 k00 = _mm256_setl_ps(kernel.data[0]);
20     __m256 k01 = _mm256_setl_ps(kernel.data[1]);
21     __m256 k02 = _mm256_setl_ps(kernel.data[2]);
22
23     __m256 k10 = _mm256_setl_ps(kernel.data[3]);
24     __m256 k11 = _mm256_setl_ps(kernel.data[4]);
25     __m256 k12 = _mm256_setl_ps(kernel.data[5]);
26
27     __m256 k20 = _mm256_setl_ps(kernel.data[6]);
28     __m256 k21 = _mm256_setl_ps(kernel.data[7]);
29     __m256 k22 = _mm256_setl_ps(kernel.data[8]);
30
31     consteval int AVX_FLOATS = 8;
32
33     for (int i = 0; i < out_height; i++) {
34         int base_i = i * params.stride;
35
36         int j = 0;
37         for (; j < out_width; j += AVX_FLOATS) {
38             int base_j = j * params.stride;
39
40             __m256 sum = _mm256_setzero_ps();
41
42             // Row 0
43             __m256 r0 = _mm256_loadu_ps(image.data[(base_i + 0) * image.width + base_j + 0]);
44             __m256 r1 = _mm256_loadu_ps(image.data[(base_i + 0) * image.width + base_j + 1]);
45             __m256 r2 = _mm256_loadu_ps(image.data[(base_i + 0) * image.width + base_j + 2]);
46
47             sum = _mm256_add_ps(sum, _mm256_mul_ps(r0, k00));
48             sum = _mm256_add_ps(sum, _mm256_mul_ps(r1, k01));
49             sum = _mm256_add_ps(sum, _mm256_mul_ps(r2, k02));
50
51             // Row 1
52             r0 = _mm256_loadu_ps(image.data[(base_i + 1) * image.width + base_j + 0]);
53             r1 = _mm256_loadu_ps(image.data[(base_i + 1) * image.width + base_j + 1]);
54             r2 = _mm256_loadu_ps(image.data[(base_i + 1) * image.width + base_j + 2]);
55
56             sum = _mm256_add_ps(sum, _mm256_mul_ps(r0, k10));
57             sum = _mm256_add_ps(sum, _mm256_mul_ps(r1, k11));
58             sum = _mm256_add_ps(sum, _mm256_mul_ps(r2, k12));
59
60             // Row 2
61             r0 = _mm256_loadu_ps(image.data[(base_i + 2) * image.width + base_j + 0]);
62             r1 = _mm256_loadu_ps(image.data[(base_i + 2) * image.width + base_j + 1]);
63             r2 = _mm256_loadu_ps(image.data[(base_i + 2) * image.width + base_j + 2]);
64
65             sum = _mm256_add_ps(sum, _mm256_mul_ps(r0, k20));
66             sum = _mm256_add_ps(sum, _mm256_mul_ps(r1, k21));
67             sum = _mm256_add_ps(sum, _mm256_mul_ps(r2, k22));
68
69             __m256_storeu_ps(output.data[(i * out_width + j) * AVX_FLOATS], sum);
70
71             for (; j < out_width; j++) {
72                 int base_j = j * params.stride;
73                 float s = 0.0f;
74                 for (int ky = 0; ky < 3; ky++)
75                     for (int kx = 0; kx < 3; kx++)
76                         s += image.data[(base_i + ky) * image.width + (base_j + kx)] *
77                            kernel.data[ky * 3 + kx];
78                 output.data[(i * out_width + j) * AVX_FLOATS] = s;
79             }
80         }
81     }
82
83     return res;
84 }
```

همانطور که مشاهده می‌کنید، در تابع baseline، پردازش به صورت خطی است و هیچگونه بهینه‌سازی در آن صورت نگرفته. البته compiler بهینه‌سازی و حتی در مواقعی موازی‌سازی را نیز تا حدی انجام می‌دهد. در دو تابع sse و avx، پردازش به صورت برداری و با کمک رجیسترهای برداری ۱۲۸ بیتی و ۲۵۶ بیتی پردازنده x86 انجام شده است.

در نوت‌بوک speed-comparison درباره میزان افزایش سرعت sse و avx نسبت به حالت baseline صحبت می‌کنیم.

تصاویری از نتیجه بررسی سرعت در نوت‌بوک مذکور:



• توضیح پیاده‌سازی SIMD

در این بخش، به جای تمرکز بر نتایج سرعت، نحوه پیاده‌سازی برداری توضیح داده می‌شود.

▪ استفاده از Intrinsics

در موتورهای SSE و AVX از توابع intrinsic مربوط به معماری پردازنده x86 استفاده شده است. این توابع به کامپایلر اجازه می‌دهند مستقیماً دستورهای SIMD متناظر را تولید کند، بدون نیاز به نوشتن کد اسمبلی به صورت جداگانه.

برای مثال:

- بارگذاری داده‌ها به رجیستر برداری
- ضرب برداری مقادیر کرنل و پیکسل‌ها
- جمع برداری نتایج میانی

▪ عرض رجیسترها

- در SSE از رجیسترهای ۱۲۸ بیتی استفاده شده که شامل ۴ مقدار float32 هستند.
 - در AVX از رجیسترهای ۲۵۶ بیتی استفاده شده که شامل ۸ مقدار float32 هستند.
- در نتیجه در هر سیکل پردازشی، به ترتیب ۴ یا ۸ عملیات ضرب-جمع به صورت همزمان انجام می‌شود.

▪ نحوه بردارسازی حلقه‌ها

در پیاده‌سازی SIMD، حلقه داخلی convolution که روی عرض تصویر حرکت می‌کند، به صورت بلوک‌های ۴ یا ۸ تایی پردازش می‌شود. به گونه‌ای که:

- چند پیکسل مجاور به صورت همزمان load می‌شوند.
- با ضرایب کرنل ضرب می‌شوند.
- نتایج در یک رجیستر جمع می‌شوند.
- در نهایت مقدار خروجی ذخیره می‌شود.

▪ مدیریت بخش باقیمانده (Tail Processing)

اگر عرض تصویر مضرب ۴ یا ۸ نباشد، بخش انتهایی تصویر که در قالب بردار کامل قرار نمی‌گیرد، به صورت scalar یا خطی پردازش می‌شود.

داده‌های تصویر به صورت آرایه خطی ذخیره شده‌اند، که باعث می‌شود دسترسی حافظه در حلقه داخلی به صورت متوالی باشد. این موضوع نقش مهمی در عملکرد SIMD ایفا می‌کند.

نکته قابل توجه این است که در حالی که عرض رجیسترهای AVX دو برابر SSE است، افزایش سرعت دقیقاً دو برابر مشاهده نمی‌شود. تحلیل کمی و تجربی این موضوع در نوت‌بوک speed-comparison صورت گرفته است.

- تصاویری از اجرای برنامه:

در ادامه تصاویری از اجرای برنامه، عکس‌های ورودی و خروجی متناظر با آن‌ها قرار داده شده است.

تصویری که فیلترهای مختلف را روی آن تست خواهیم کرد:




```
Options:
  (1) Functional Test
  (2) Speed Test
Enter Option [2] > 1
Options:
  (1) Baseline
  (2) SSE
  (3) AVX
Enter Engine Mode [1] >
Options:
  (1) Sharpen
  (2) Box Blur
  (3) Gaussian Blur
  (4) Sobel X
  (5) Sobel Y
Enter Kernel Type [1] >
Options:
  (3) 3 x 3
  (5) 5 x 5
  (7) 7 x 7
Enter Kernel Size [3] >
Enter Image Path [./images/normal-small/01.jpeg] >
Enter Output Filename [./images/output/01.jpeg] >
Options:
  (1) RGB
  (2) Grayscale
Enter Color Mode [1] >
[TIMING] Engine: Baseline took 86.699823 ms
```



```
Options:
  (1) Functional Test
  (2) Speed Test
Enter Option [2] > 1
Options:
  (1) Baseline
  (2) SSE
  (3) AVX
Enter Engine Mode [1] >
Options:
  (1) Sharpen
  (2) Box Blur
  (3) Gaussian Blur
  (4) Sobel X
  (5) Sobel Y
Enter Kernel Type [1] > 2
Options:
  (3) 3 x 3
  (5) 5 x 5
  (7) 7 x 7
Enter Kernel Size [3] > 7
Enter Image Path [./images/normal-small/01.jpeg] >
Enter Output Filename [./images/output/01.jpeg] >
Options:
  (1) RGB
  (2) Grayscale
Enter Color Mode [1] >
[TIMING] Engine: Baseline took 212.983973 ms
```




```
Options:
  (1) Functional Test
  (2) Speed Test
Enter Option [2] > 1
Options:
  (1) Baseline
  (2) SSE
  (3) AVX
Enter Engine Mode [1] > 3
Options:
  (1) Sharpen
  (2) Box Blur
  (3) Gaussian Blur
  (4) Sobel X
  (5) Sobel Y
Enter Kernel Type [1] > 4
Options:
  (3) 3 x 3
  (5) 5 x 5
  (7) 7 x 7
Enter Kernel Size [3] >
Enter Image Path [./images/normal-small/01.jpeg] >
Enter Output Filename [./images/output/01.jpeg] >
Options:
  (1) RGB
  (2) Grayscale
Enter Color Mode [1] >
[TIMING] Engine: AVX took 22.451762 ms
```



لازم به ذکر است که در این تصاویر از مود functional برنامه و در نوتبوک speed-comparison برای مقایسه سرعت، از مود speed استفاده شده است.

- معرفی نوتبوک‌های تکمیلی:

- نوتبوک speed-comparison:

در این نوتبوک سرعت ۳ موتور مختلف با هم مقایسه شده.

- نوتبوک object-detection:

در این نوتبوک تلاش شده است با استفاده از فیلترهای پیاده‌سازی شده، اشکال ساده دیتاست‌های ساختگی را شناسایی و دقت تشخیص را بیازماییم.

- نوتبوک cnn-pneumonia:

در این نوتبوک با کمک کتابخانه pytorch، یک مدل بسیار ساده روی داده‌های پزشکی آموزش داده شده و سپس وزن‌های آن را به برنامه C++ خود انتقال می‌دهیم و با ۳ موتور اصلی برنامه آن دقت و سرعت پردازش روی داده‌های test را بررسی می‌کنیم.