

Contents

1 Requirement Analysis	2
1.1 Background Knowledge	2
1.2 Design Data Structure	3
1.3 Convolution Function	3
1.4 Testing the correctness	3
1.5 Improve the efficiency	4
2 Code Implementation	4
2.1 Image Operation	4
2.1.1 Input the Image	4
2.1.2 Output the Image	6
2.2 Kernel Initialization	7
2.3 Padding	7
2.4 Convolution Operation	7
2.4.1 Ordinary Function	7
2.4.2 Reduced Loops Function	8
2.4.3 Unloop Function	8
2.4.4 SIMD Function	9
2.4.5 SIMD with OpenMP Function	10
2.5 Convolution Function	11
2.6 Auxiliary Functions	12
2.6.1 Get Random Data	12
2.6.2 Print the Data	13
2.6.3 Calculate the Time Cost	13
2.6.4 Write data into a CSV file	13
2.6.5 Package the Testing	14
2.6.6 Free the memory	14
2.7 Compile the file	15
3 Experiments & Analysis	15
3.1 Efficiency Analysis	15
3.1.1 The Dimension of Loops	15
3.1.2 Unloop	17
3.1.3 SIMD Analysis	17
3.1.4 OpenMP Analysis	20
3.1.5 -O3	21

3.2	Correctness	21
3.2.1	Qualitative Method	21
4	Difficulties & Solutions	22
4.1	Difficulties in input and output the image	22
4.2	Difficulties in Checking the correctness of the convolution operation	22
4.2.1	Attempt Intel Math Kernel Library	23
4.2.2	Attempt to check the correctness quantitatively	23
4.3	Difficulties in compiling the file using CMakeLists	23
5	Brief Summary	24

1 Requirement Analysis

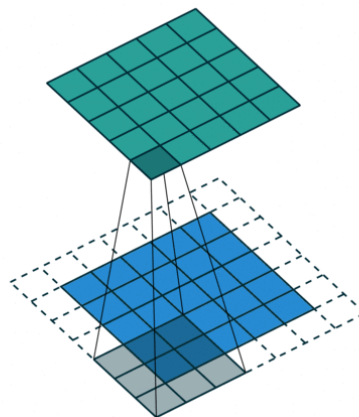
The concise purpose of the project is to design a function that implements the Convolutional Layer in Deep learning using C and to improve the efficiency of the program.

1.1 Background Knowledge

In a convolutional layer, multiple convolution kernels are typically set up. Each kernel is a small matrix (contains data) used for performing a convolution operation on the input data. The size of the kernel is the size of the matrix.

Padding refers to the process of adding extra "padding" pixels (typically added with a value of zero) around the input image before performing convolution. (The dashed squares in the figure) The purpose of padding is to control the size of the output feature maps and preserve spatial information at the borders of the input image.

Stride refers to the number of pixels the kernel is moved in each step while performing convolution operation on the input image. The use of stride helps in reducing the size of the output feature map and also helps in controlling the computational complexity of the model. A stride of 1 means the kernel is shifted by one pixel at a time, and hence the output size will be the same as the input size. A stride of 2 means the kernel is shifted by two pixels at a time, resulting in a smaller output size.



A convolution operation can be seen as a sliding window that moves the kernel over the input data and calculates the dot product between the kernel and the data within the window, resulting in an output value. By moving the kernel over the input data, multiple output values can be obtained, which constitute an output feature map.

1.2 Design Data Structure

The Data structures named Image and Kernel are designed.

Image: The Image data structure contains 4 parts. The height, width, channels, and the data of an image. Considering using SIMD to improve the convolution operation of the images of channel 3, three float pointers storing the image data are set in the data structure. SIMD requires to get continuous data stored in the memory. If we store data of RGB adjacent to each other, it costs a lot to continuously restore the data during the convolution operation. (If the Channel is 1, the data2 and data3 will be set as NULL)

```
1 typedef struct {
2     size_t width;
3     size_t height;
4     size_t channels;
5     float* data1;
6     float* data2;
7     float* data3;
8 } Image;
```

Kernel: The Kernel data structure contains also 4 parts. The height, width, channels, and the data of the kernel. When implementing the convolution operation of the images of channel 3, the kernel data used to operate each channel are different in most applications. Thus, three float pointers storing the kernel data are packaged in a data structure, making it easy to find the specific data during the convolution operation. (If the Channel is 1, the data2 and data3 will be set as NULL)

```
1 typedef struct {
2     size_t width;
3     size_t height;
4     size_t channels;
5     float* data1;
6     float* data2;
7     float* data3;
8 } Kernel;
```

1.3 Convolution Function

First of all, we need to read an image (In this project, we choose to read jpg image). Secondly, we initialize the kernel. Thirdly, we pad the input image according to the size of the kernel. (The purpose is to make the output the same size as the input) Fourthly, we choose the relatively efficient way to implement the Convolution operation. Finally, we output the image in jpg form.

Convolution Function

- Input the Image
 - `input_jpeg_image`
- Kernel Initialization
 - `kernel_init`
- Padding
 - `padding`
- Convolution operation
 - 1x1 kernel
 - `cnn_avx2_kernel1x1`
 - 3x3 kernel
 - `cnn_avx2_omp_kernel3x3` (multiple of 8)
 - `cnn_avx2_kernel3x3`
 - 5x5 kernel
 - `cnn_avx2_omp_kernel5x5` (multiple of 8)
 - `cnn_avx2_kernel5x5`
- Output the Image
 - `output_jpeg_image`

1.4 Testing the correctness

In this project, we use the convolution operation in MATLAB to test the correctness of my own convolution operation qualitatively.

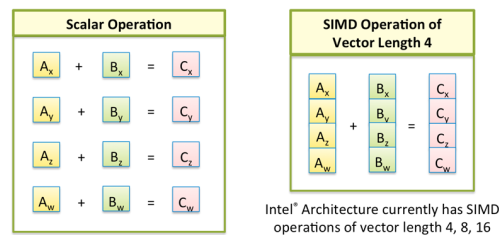
In MATLAB, we can use the `imshow()` function to show the image after the convolution operation. Then we can see if that image is as same as our own output image.

1.5 Improve the efficiency

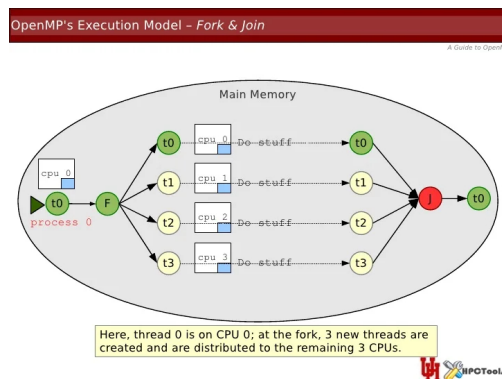
This project mainly improves the efficiency of the convolution operation in three ways: Reducing the number of looping, SIMD (single instruction multiple Data), and OpenMP (Open Multi-processing).

Loop costs a lot. Each loop requires repeated calculations and comparisons, operations that take up processor time and slow down program execution. In addition, loops need to store multiple variables and intermediate results in memory, which can increase memory usage and reduce program efficiency. Thus, lowering the dimension of the looping definitely improves the efficiency of the program.

In SIMD architecture, a single instruction can operate multiple data elements at the same time, making full use of register space. It can greatly reduce the number of instruction executions and improves the running efficiency of the program.



On a multi-core computer, OpenMP can map multiple threads to different CPU cores so that different threads can execute different tasks at the same time. This parallel execution can significantly improve the running efficiency of programs because it can fully utilize the multi-core processing power of the computer.



2 Code Implementation

2.1 Image Operation

In this project, we choose to use a library called libjpeg in C to input and output an jpg image.

```
1 #include <jpeglib.h>
```

2.1.1 Input the Image

Firstly, we open the file, create `jpeg-decompress-struct` (decompress the image), and create `jpeg-error-mgr` (manage error handling during JPEG decompression). Secondly, we set some basic parameters to the `jpeg-decompress-struct` and start the compression. Thirdly, we store the data row by row and set some basic information to the image structure. Finally, we finish the decompression, close the file and free some memories.

```

1  bool input_jpeg_image(const char* filename, Image* image)
2  {
3      if(image==NULL) {
4          printf("Error Image(structure) input!");
5          return false;
6      }
7      //Open JPEG file, create JPEG decoder object and JPEG error object
8      FILE *infile;
9      if ((infile = fopen(filename, "rb")) == NULL) {
10         printf("Error opening JPEG file!");
11         return false;
12     }
13     struct jpeg_decompress_struct cinfo;
14     struct jpeg_error_mgr jerr;
15
16     //Set basic parameters
17     cinfo.err = jpeg_std_error(&jerr);
18     jpeg_create_decompress(&cinfo);
19     jpeg_stdio_src(&cinfo, infile);
20     jpeg_read_header(&cinfo, TRUE);
21     //Start decompression
22     jpeg_start_decompress(&cinfo);
23
24     //Passing data into the Image structure
25     int row_stride = cinfo.output_width * cinfo.output_components;
26     unsigned char *buffer = (unsigned char*)malloc(row_stride*sizeof(unsigned char));
27     if(cinfo.output_components==1)
28     {
29         cinfo.out_color_space = JCS_GRAYSCALE;
30         cinfo.quantize_colors = FALSE;
31         float *Data = (float*)aligned_alloc(256,row_stride*cinfo.output_height*sizeof(float));
32         while (cinfo.output_scanline < cinfo.output_height) {
33             int i = cinfo.output_scanline;
34             jpeg_read_scanlines(&cinfo, &buffer, 1);
35             for(int j=0;j<row_stride;j++)
36                 Data[i*row_stride+j]=(float)buffer[j];
37         }
38         image->data1 = Data;
39         image->data2 = NULL;
40         image->data3 = NULL;
41     }
42     else if(cinfo.output_components==3)
43     {
44         cinfo.out_color_space = JCS_RGB;
45         float* red_channel = (float*)aligned_alloc(256,cinfo.output_width * cinfo.output_height * sizeof(float));
46         float* green_channel = (float*)aligned_alloc(256,cinfo.output_width * cinfo.output_height * sizeof(float));
47         float* blue_channel = (float*)aligned_alloc(256,cinfo.output_width * cinfo.output_height * sizeof(float));
48         while (cinfo.output_scanline < cinfo.output_height)
49         {
50             jpeg_read_scanlines(&cinfo, &buffer, 1);
51             int offset = (cinfo.output_scanline - 1) * cinfo.output_width;
52             for (int i = 0, j = 0; i < row_stride; i += 3, j++) {
53                 red_channel[offset + j] =(float) buffer[i];
54                 green_channel[offset + j] =(float) buffer[i + 1];
55                 blue_channel[offset + j] =(float) buffer[i + 2];
56             }
57         }
58         image->data1 = red_channel;
59         image->data2 = green_channel;
60         image->data3 = blue_channel;
61     }
62     else
63     {
64         printf("The channels of the image must be 1 or 3");
65         return false;
66     }
67     image->width = cinfo.output_width;
68     image->height = cinfo.output_height;
69     image->channels = cinfo.output_components;
70
71     //Complete the decompression process and release all related resources
72     jpeg_finish_decompress(&cinfo);
73     jpeg_destroy_decompress(&cinfo);
74     fclose(infile);
75     free(buffer);
76     return true;
77 }

```

2.1.2 Output the Image

Firstly, we create jpeg-compress-struct (compress the image), create jpeg-error-mgr (manage error handling during JPEG decompression) and open a file. Secondly, we connect the output stream with a compressed object, set some basic parameters to the jpeg-compress-struct, and start the compression. Thirdly, we transfer the data row by row and set some basic information to the compressor. Finally, we finish the compression, output the image, close the file and free some memories.

```
1 bool output_jpeg_image(const char* filename, Image* image)
2 {
3     if(image==NULL) {
4         printf("Error Image(structure) input!");
5         return false;
6     }
7     //JPEG, JPEG
8     struct jpeg_compress_struct cinfo;
9     struct jpeg_error_mgr jerr;
10    cinfo.err = jpeg_std_error(&jerr);
11    jpeg_create_compress(&cinfo);
12    FILE* outfile;
13    if ((outfile = fopen(filename, "wb")) == NULL) {
14        printf("Error opening JPEG file!");
15        return false;
16    }
17    //
18    jpeg_stdio_dest(&cinfo, outfile);
19    //
20    cinfo.image_width = image->width;
21    cinfo.image_height = image->height;
22    cinfo.input_components = image->channels;
23    jpeg_set_defaults(&cinfo);
24    jpeg_set_quality(&cinfo, 90, TRUE);
25    //
26    jpeg_start_compress(&cinfo, TRUE);
27    //image
28    JSAMPROW row_pointer[1];
29    int row_stride = cinfo.image_width * cinfo.input_components;
30    unsigned char *buffer = (char*)malloc(row_stride*sizeof(unsigned char));
31    if(image->channels==1)
32    {
33        cinfo.in_color_space = JCS_GRAYSCALE;
34        while (cinfo.next_scanline < cinfo.image_height) {
35            int i = cinfo.next_scanline;
36            for(int j=0; j<row_stride; j++)
37                buffer[j]=(unsigned char)image->data1[i*row_stride+j];
38            row_pointer[0] = &buffer[0];
39            jpeg_write_scanlines(&cinfo, row_pointer, 1);
40        }
41    }
42    else if(image->channels==3)
43    {
44        cinfo.in_color_space = JCS_RGB;
45        while (cinfo.next_scanline < cinfo.image_height) {
46            int offset = cinfo.next_scanline * cinfo.image_width;
47            for (int i = 0, j = 0; i < row_stride; i += 3, j++) {
48                buffer[i] = (unsigned char)image->data1[offset + j];
49                buffer[i + 1] = (unsigned char)image->data2[offset + j];
50                buffer[i + 2] = (unsigned char)image->data3[offset + j];
51            }
52            row_pointer[0] = &buffer[0];
53            jpeg_write_scanlines(&cinfo, row_pointer, 1);
54        }
55    }
56    else
57    {
58        printf("The channels of the image must be 1 or 3");
59        return false;
60    }
61    //
62    jpeg_finish_compress(&cinfo);
63    fclose(outfile);
64    jpeg_destroy_compress(&cinfo);
65    return true;
66 }
67
68
69
70
```

2.2 Kernel Initialization

We transfer the basic information of a kernel to the function and initialize the kernel through the function.

```
1 bool kernel_init(Kernel* kernel, const size_t height, const size_t width, const size_t ch,  
2 const float* Data1, const float* Data2, const float* Data3)  
3 {  
4     if(ch==1&&(Data1!=NULL))  
5     {  
6         kernel->data1 = Data1;  
7         kernel->data2 = NULL;  
8         kernel->data3 = NULL;  
9     }  
10    else if(ch==3&&(Data1!=NULL)&&(Data2!=NULL)&&(Data3!=NULL))  
11    {  
12        kernel->data1 = Data1;  
13        kernel->data2 = Data2;  
14        kernel->data3 = Data3;  
15    }  
16    else  
17    {  
18        printf("Error Input!");  
19        return false;  
20    }  
21    kernel->width = width;  
22    kernel->height = height;  
23    kernel->channels = ch;  
24    return true;  
25 }
```

2.3 Padding

We transfer the data of the input image, the size of the input image, and the size of the kernel into the function. In order to let the output image maintains the same size as the input image, the function will pad the edges of the data with zeros. If the kernel is nxn we wrap (n-1)/2 circle(s) of zeros around the data.

```
1 float* padding(const float* input_data, const size_t input_height,  
2 const size_t input_width, const size_t kernel_height, const size_t kernel_width)  
3 {  
4     if(input_data==NULL)  
5     {  
6         printf("NULL pointer!");  
7         return FALSE;  
8     }  
9     size_t pad_width = input_width + (kernel_width - 1);  
10    size_t pad_height = input_height + (kernel_height - 1);  
11    float* pad_data=(float*)calloc(pad_width*pad_height,sizeof(float));  
12    if(pad_data==NULL)  
13    {  
14        printf("Error Calloc for the padding data!");  
15        return NULL;  
16    }  
17    for(int i = ((kernel_width - 1) / 2); i < input_height + ((kernel_width - 1) / 2); i++){  
18        memcpy(pad_data + i * pad_width + ((kernel_width - 1) / 2) , input_data + (i - ((kernel_width - 1) / 2))  
19    }  
20    return pad_data;  
21 }
```

2.4 Convolution Operation

Five types of functions are mentioned in this part. They are the Ordinary Function (3 dimensions of loops), the Reduced Loops Function (2 dimensions of loops), the Unloop Function, the SIMD Function, and the SiMD with OpenMP Function.

2.4.1 Ordinary Function

```
1 bool cnnFunction_v1(const float* input_data, const size_t input_height, const size_t input_width, const float* kernel_data,  
2 const size_t kernel_width, float* output_data)  
3 {  
4     //Safe Detection  
5     size_t pad_width=input_width+(kernel_width-1);  
6     for(int j = 0; j < input_height; j++)
```

```

7      {
8          for(int i = 0; i < input_width;i++)
9          {
10             for(int m = 0; m<kernel_height;m++)
11             {
12                 for(int n = 0, k = 0; n < kernel_width; n++, k++)
13                 {
14                     *(output_data+(j*input_width)+i) += (*(input_data + (j * pad_width) + i + (m*pad_width) + n))
15                 }
16             }
17         }
18     }
19     return true;
20 }

```

2.4.2 Reduced Loops Function

For the convolution of an element, it is necessary to traverse an entire convolution kernel and the corresponding data of the picture, which we can get the row-index = $(k / \text{kernel-width})$ and the column-index = $(k \% \text{kernel-width})$. For the convolution of the whole image, which we can get the row-index = $(t / \text{input-width})$ and the column-index = $(t \% \text{kernel-width})$. Thus, we can reduce the dimension of the loop from 4 to 2.

```

1  bool cnnFunction_v2(const float* input_data, const size_t input_height, const size_t input_width, const float* kernel_data,
2      const size_t kernel_width, float* output_data)
3  {
4      //Safe Detection
5      size_t pad_width=input_width+(kernel_width-1);
6      for(int t = 0; t < input_height*input_width; t++)
7      {
8          int j = t / input_width;
9          int i = t % input_width;
10         for(int k=0; k<kernel_height * kernel_width;k++)
11         {
12             *(output_data+(j * input_width)+i) += (*(input_data + (j * pad_width) + i + ((k / kernel_width)*pad_width) + (k % kernel_width)))
13         }
14     }
15     return true;
16 }

```

2.4.3 Unloop Function

Based on a three-dimensional loop we calculate 8 elements at one time. In order to reduce the time of calculating the repeated value $(j*\text{input-width})$, we set a value named stride at the top of the first loop. By the way, in order to deal with all the elements in the case that the input width is not a multiple of 8, we can only reduce the loop dimension to 3.

```

1  bool cnn_unloop(const float* input_data, const size_t input_height, const size_t input_width, const float* kernel_data,
2      const size_t kernel_width, float* output_data)
3  {
4      //Safe Detection
5      bool cnn_unloop()
6      {
7          for(int j = 0; j < input_height; j++)
8          {
9              int stride = j*input_width;
10             int i = 0;
11             for(; i < (input_width/8)*8;i+=8)
12             {
13                 for(int m = 0; m<kernel_height;m++)
14                 {
15                     for(int n = 0,k = 0; n < kernel_width; n++,k++)
16                     {
17                         *(output_data+(stride)+i) += (*(input_data + (j * pad_width) + i + (m*pad_width) + n)) * (*(kernel_data+k));
18                         *(output_data+(stride)+(i+1)) += (*(input_data + (j * pad_width) + (i+1) + (m*pad_width) + n)) * (*(kernel_data+k));
19                         *(output_data+(stride)+(i+2)) += (*(input_data + (j * pad_width) + (i+2) + (m*pad_width) + n)) * (*(kernel_data+k));
20                         *(output_data+(stride)+(i+3)) += (*(input_data + (j * pad_width) + (i+3) + (m*pad_width) + n)) * (*(kernel_data+k));
21                         *(output_data+(stride)+(i+4)) += (*(input_data + (j * pad_width) + (i+4) + (m*pad_width) + n)) * (*(kernel_data+k));
22                         *(output_data+(stride)+(i+5)) += (*(input_data + (j * pad_width) + (i+5) + (m*pad_width) + n)) * (*(kernel_data+k));
23                         *(output_data+(stride)+(i+6)) += (*(input_data + (j * pad_width) + (i+6) + (m*pad_width) + n)) * (*(kernel_data+k));
24                     }
25                 }
26             }
27         }
28     }
29 }

```



```

30         (m*pad_width) + n)) * (*(kernel_data+k));
31         *(output_data+(stride)+(i+7)) += (*(input_data + (j * pad_width) + (i+7) +
32         (m*pad_width) + n)) * (*(kernel_data+k));
33     }
34 }
35 }
36 if((input_width % 8) != 0)
37 {
38     for(int w=i;w<(input_width % 8)+i;w++)
39     {
40         for(int m = 0; m<kernel_height;m++)
41         {
42             for(int n = 0,k = 0; n < kernel_width; n++,k++)
43             {
44                 *(output_data+(stride)+w) += (*(input_data + (j * pad_width) + w + m*pad_width + n)) *
45                 (*(kernel_data+k));
46             }
47         }
48     }
49 }
50 }
51 }
52 return true;
53 }

```

2.4.4 SIMD Function

Based on a three-dimensional loop we calculate 8 elements at one time. We use the datatype `_m256` and the function named `_mm256_loadu_ps` to get the continuous 8 floats (8*4 bytes*8 bits) in from the memory into the register. Because the size of the kernel is 3x3, thus we need to get 9 groups of data (each group has 8 floats) at once. As the image data are stored in the datatype `_m256`, so do the kernel data. Moreover, the data of the kernel should also be in the same size group with continuous memory allocated. Having set the data, we can use `_mm256_add_ps` to add and `_mm256_mul_ps` to multiple 2 groups of data (containing 8 floats) at once. After the operation, we use `_mm256_storeu_ps` to store the data in the register back to the memory. (Remember to use `_mm256_setzero_ps` to set the data named `res` to zero at the beginning of every operation).

```

1 bool cnn_avx2_kernel3x3(const float* input_data, const size_t input_height, const size_t input_width,
2 const float* kernel_data, const size_t kernel_height,
3 const size_t kernel_width, float* output_data)
4 {
5     //Safe Detection
6     size_t pad_width=input_width+(kernel_width-1);
7     #ifdef WITH_AVX2
8     _m256 im11,im12,im13,im21,im22,im23,im31,im32,im33;
9     _m256 ker11,ker12,ker13,ker21,ker22,ker23,ker31,ker32,ker33;
10    _m256 res;
11    float*Ker11=(float*)(aligned_alloc(256,8*sizeof(float)));
12    float*Ker12=(float*)(aligned_alloc(256,8*sizeof(float)));
13    float*Ker13=(float*)(aligned_alloc(256,8*sizeof(float)));
14    float*Ker21=(float*)(aligned_alloc(256,8*sizeof(float)));
15    float*Ker22=(float*)(aligned_alloc(256,8*sizeof(float)));
16    float*Ker23=(float*)(aligned_alloc(256,8*sizeof(float)));
17    float*Ker31=(float*)(aligned_alloc(256,8*sizeof(float)));
18    float*Ker32=(float*)(aligned_alloc(256,8*sizeof(float)));
19    float*Ker33=(float*)(aligned_alloc(256,8*sizeof(float)));
20    for(int i=0;i<8;i++)
21    {
22        *(Ker11+i)=*(kernel_data);
23        *(Ker12+i)=*(kernel_data+1);
24        *(Ker13+i)=*(kernel_data+2);
25        *(Ker21+i)=*(kernel_data+3);
26        *(Ker22+i)=*(kernel_data+4);
27        *(Ker23+i)=*(kernel_data+5);
28        *(Ker31+i)=*(kernel_data+6);
29        *(Ker32+i)=*(kernel_data+7);
30        *(Ker33+i)=*(kernel_data+8);
31    }
32    ker11 = _mm256_loadu_ps(Ker11);
33    ker12 = _mm256_loadu_ps(Ker12);
34    ker13 = _mm256_loadu_ps(Ker13);
35    ker21 = _mm256_loadu_ps(Ker21);
36    ker22 = _mm256_loadu_ps(Ker22);
37    ker23 = _mm256_loadu_ps(Ker23);
38    ker31 = _mm256_loadu_ps(Ker31);
39    ker32 = _mm256_loadu_ps(Ker32);
40    ker33 = _mm256_loadu_ps(Ker33);
41    float*temp=(float*)(aligned_alloc(256,8*sizeof(float)));
42    for(size_t j=0; j<input_height;j++)
43    {

```

```

44     size_t i=0;
45     for(; i<(input_width/8)*8; i+=8)
46     {
47         res = _mm256_setzero_ps();
48         im11 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (0*pad_width) + 0);
49         im12 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (0*pad_width) + 1);
50         im13 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (0*pad_width) + 2);
51         im21 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (1*pad_width) + 0);
52         im22 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (1*pad_width) + 1);
53         im23 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (1*pad_width) + 2);
54         im31 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (2*pad_width) + 0);
55         im32 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (2*pad_width) + 1);
56         im33 = _mm256_loadu_ps(input_data + (j * pad_width) + i + (2*pad_width) + 2);
57         res = _mm256_add_ps(res, _mm256_mul_ps(im11, ker11));
58         res = _mm256_add_ps(res, _mm256_mul_ps(im12, ker12));
59         res = _mm256_add_ps(res, _mm256_mul_ps(im13, ker13));
60         res = _mm256_add_ps(res, _mm256_mul_ps(im21, ker21));
61         res = _mm256_add_ps(res, _mm256_mul_ps(im22, ker22));
62         res = _mm256_add_ps(res, _mm256_mul_ps(im23, ker23));
63         res = _mm256_add_ps(res, _mm256_mul_ps(im31, ker31));
64         res = _mm256_add_ps(res, _mm256_mul_ps(im32, ker32));
65         res = _mm256_add_ps(res, _mm256_mul_ps(im33, ker33));
66         _mm256_storeu_ps(temp, res);
67         for(size_t k=0;k<8;k++)
68             *(output_data+(j*input_width)+i+k) = *(temp+k);
69     }
70     if((input_width % 8) != 0)
71     {
72         for(int w=i;w<(input_width % 8)+i;w++)
73         {
74             for(int k =0; k<kernel_height*kernel_width;k++)
75             {
76                 *(output_data+(j*input_width)+w) += (*(input_data +
77                 (j * pad_width) + w + (k / kernel_width)*pad_width +
78                 (k % kernel_width))) * (*(kernel_data+k));
79             }
80         }
81     }
82 }
83 return TRUE;
84 #else
85 printf("AVX2 is not supported");
86 return FALSE;
87 #endif
88 }

```

The example is under the condition of the 3x3 kernel. There are functions written in 1x1 and 5x5 kernels in the source code. (named `cnn_avx2_kernel5x5` and `cnn_avx2_kernel1x1`)

2.4.5 SIMD with OpenMP Function

We use the instruction 'pragma omp parallel for' before the loop to map multiple threads to different CPU cores so that different threads can execute the operation at the same time.

```

1  bool cnn_avx2_omp_kernel3x3(const float* input_data, const size_t input_height, const size_t input_width,
2  const float* kernel_data, const size_t kernel_height,
3  const size_t kernel_width, float* output_data)
4  {
5      //Safe Detection
6      size_t pad_width=input_width+(kernel_width-1);
7  #ifdef WITH_AVX2
8      //Operation in SIMD
9      #pragma omp parallel for
10     for(size_t j=0; j<input_height;j++)
11     {
12         size_t i=0;
13         for(; i<(input_width/8)*8; i+=8)
14         {
15             //operations
16             .....
17         }
18         return TRUE;
19     #else
20     printf("AVX2 is not supported\n");
21     return FALSE;
22     #endif
23 }

```

The example is under the condition of the 3x3 kernel. There are functions written in 1x1 and 5x5 kernels in the source code. (named `cnn_avx2_omp_kernel5x5` and `cnn_avx2_omp_kernel1x1`)

2.5 Convolution Function

This function includes the functions mentioned above. We transfer an image into the input Image structure, the Kernel structure, and the output Image structure to the function. Inside the function, we pad the input data, operate the convolution operation, and package all the information to the output Image structure. Considering the different size of the kernel and the different size of the input data, we choose the relevantly efficient function to operate the convolution operation.

```

1  bool Convolution(const Image* input_image, const Kernel* kernel, Image* output_image)
2  {
3      //Safe Detection
4      if(input_image->channels==1)
5      {
6          //padding
7          float* pad_data = padding(input_image->data1, input_image->height,
8          input_image->width, kernel->height, kernel->width);
9
10         //cnn
11         float* data1 = (float*)malloc(sizeof(float)*input_image
12         ->height*input_image->width);
13         if(kernel->height==1&&kernel->width==1)
14         {
15             cnn_avx2_kernel1x1(pad_data, input_image->height, input_image->width, kernel->data1, kernel->height,
16             kernel->width, data1);
17             output_image->data1 = data1;
18         }
19         else if(kernel->height==3&&kernel->width==3)
20         {
21             if(input_image->width%8==0)
22                 cnn_avx2_omp_kernel3x3(pad_data, input_image->height, input_image->width, kernel->data1,
23                 kernel->height, kernel->width, data1);
24             else
25                 cnn_avx2_kernel3x3(pad_data, input_image->height, input_image->width, kernel->data1,
26                 kernel->height, kernel->width, data1);
27             output_image->data1 = data1;
28         }
29         else if(kernel->height==5&&kernel->width==5)
30         {
31             if(input_image->width%8==0)
32                 cnn_avx2_omp_kernel5x5(pad_data, input_image->height, input_image->width, kernel->data1,
33                 kernel->height, kernel->width, data1);
34             else
35                 cnn_avx2_kernel5x5(pad_data, input_image->height, input_image->width, kernel->data1,
36                 kernel->height, kernel->width, data1);
37             output_image->data1 = data1;
38         }
39         output_image->data2 = NULL;
40         output_image->data3 = NULL;
41         free(pad_data);
42     }
43     else if(input_image->channels==3)
44     {
45         //padding
46         float* pad_data1 = padding(input_image->data1, input_image->height, input_image->width,
47         kernel->height, kernel->width);
48         float* pad_data2 = padding(input_image->data2, input_image->height, input_image->width,
49         kernel->height, kernel->width);
50         float* pad_data3 = padding(input_image->data3, input_image->height, input_image->width,
51         kernel->height, kernel->width);
52
53         //cnn
54         float* data1 = (float*)malloc(sizeof(float)*input_image->height*input_image->width);
55         float* data2 = (float*)malloc(sizeof(float)*input_image->height*input_image->width);
56         float* data3 = (float*)malloc(sizeof(float)*input_image->height*input_image->width);
57         if(kernel->height==1&&kernel->width==1)
58         {
59             cnn_avx2_kernel1x1(pad_data1, input_image->height, input_image->width, kernel->data1,
60             kernel->height, kernel->width, data1);
61             cnn_avx2_kernel1x1(pad_data2, input_image->height, input_image->width, kernel->data1,
62             kernel->height, kernel->width, data2);
63             cnn_avx2_kernel1x1(pad_data3, input_image->height, input_image->width, kernel->data1,
64             kernel->height, kernel->width, data3);
65             output_image->data1 = data1;
66             output_image->data2 = data2;
67             output_image->data3 = data3;
68         }
69         else if(kernel->height==3&&kernel->width==3)

```

```

70 {
71     if(input_image->width%8==0)
72     {
73         cnn_avx2_omp_kernel3x3(pad_data1,input_image->height,input_image->width,kernel->data1,
74         kernel->height,kernel->width,data1);
75         cnn_avx2_omp_kernel3x3(pad_data2,input_image->height,input_image->width,kernel->data1,
76         kernel->height,kernel->width,data2);
77         cnn_avx2_omp_kernel3x3(pad_data3,input_image->height,input_image->width,kernel->data1,
78         kernel->height,kernel->width,data3);
79     }
80     else
81     {
82         cnn_avx2_kernel3x3(pad_data1,input_image->height,input_image->width,kernel->data1,
83         kernel->height,kernel->width,data1);
84         cnn_avx2_kernel3x3(pad_data2,input_image->height,input_image->width,kernel->data1,
85         kernel->height,kernel->width,data2);
86         cnn_avx2_kernel3x3(pad_data3,input_image->height,input_image->width,kernel->data1,
87         kernel->height,kernel->width,data3);
88     }
89     output_image->data1 = data1;
90     output_image->data2 = data2;
91     output_image->data3 = data3;
92 }
93 else if(kernel->height==5&&kernel->width==5)
94 {
95     if(input_image->width%8==0)
96     {
97         cnn_avx2_omp_kernel5x5(pad_data1,input_image->height,input_image->width,kernel->data1,
98         kernel->height,kernel->width,data1);
99         cnn_avx2_omp_kernel5x5(pad_data2,input_image->height,input_image->width,kernel->data1,
100        kernel->height,kernel->width,data2);
101        cnn_avx2_omp_kernel5x5(pad_data3,input_image->height,input_image->width,kernel->data1,
102        kernel->height,kernel->width,data3);
103    }
104    else
105    {
106        cnn_avx2_kernel5x5(pad_data1,input_image->height,input_image->width,kernel->data1,
107        kernel->height,kernel->width,data1);
108        cnn_avx2_kernel5x5(pad_data2,input_image->height,input_image->width,kernel->data1,
109        kernel->height,kernel->width,data2);
110        cnn_avx2_kernel5x5(pad_data3,input_image->height,input_image->width,kernel->data1,
111        kernel->height,kernel->width,data3);
112    }
113    output_image->data1 = data1;
114    output_image->data2 = data2;
115    output_image->data3 = data3;
116 }
117 else
118 {
119     printf("The size of the kernel must be 1x1, 3x3 or 5x5");
120 }
121 free(pad_data1);
122 free(pad_data2);
123 free(pad_data3);
124 output_image->width = input_image->width;
125 output_image->height = input_image->height;
126 output_image->channels = input_image->channels;
127 return true;
128 }
129 }

```

2.6 Auxiliary Functions

2.6.1 Get Random Data

We obtain random numbers by calculating formulas and obtaining the seeds of random numbers through the system clock. And assign the first place of these numbers to a pointer. (In this project the range of the random float number is from 0 to 255)

```

1 float* Random(const size_t height, const size_t width)
2 {
3     float* data = (float*)malloc(sizeof(float)*height*width);
4     if(data==NULL)
5     {
6         printf("Malloc Error");
7         return NULL;
8     }
9     srand((unsigned)time(NULL));
10    for(int i=0;i<height*width;i++)

```

```

11         *(data + i) = ((rand() / (float)RAND_MAX) + (rand()% 254 ));
12     return data;
13 }

```

2.6.2 Print the Data

When the amount of data is small, we can print out the data in the form of a matrix for us to check if the result of the operation is correct.

```

1 bool printinfo(const float* data, const size_t height, const size_t width, const char* name)
2 {
3     if(data==NULL)
4     {
5         printf("NULL Pointer!");
6         return FALSE;
7     }
8     printf("%s:\n",name);
9     for(int j=0; j<height;j++)
10    {
11        for(int i=0; i<width; i++)
12            printf("%.2f ",*(data + j*width + i));
13        printf("\n");
14    }
15    printf("\n");
16    return true;
17 }

```

2.6.3 Calculate the Time Cost

When dealing with the same data, we use the length of the time cost through this operation to test the efficiency of each function.

Since this program is running under Linux environment, we use the Linux environment timing function gettimeofday() to time the program. The timing function has a high accuracy of microseconds. This function gets the elapsed time and time zone (UTC time) from January 1, 1970 to the present, but according to the official Linux documentation, the time zone is no longer used, so just pass NULL when using it.

The difference between the current time before the convolution operation and the current time after the convolution operation is completed is the difference between them, which is the program's running time. time

```

1 struct timeval start, end;
2 gettimeofday(&start, NULL);
3 //Convolution operation
4 gettimeofday(&end, NULL);
5 printf("Time Cost : %lf s\n", (end.tv_sec - start.tv_sec) + (double)(end.tv_usec - start.tv_usec) / 1000000.0);

```

2.6.4 Write data into a CSV file

When conducting a large number of experiments which collect a large amount of data, we can write our data to a CSV file, which is more convenient for us to process the data.

If we do not have a file, we can create a csv named 'TimeCost.csv' through fopen() using the 'w' mode to write some basic information. (Just an Example)

```

1 FILE *fp;
2 fp = fopen("TimeCost.csv", "w");
3 if(fp == NULL) {
4     printf("Failed to open file\n");
5     return 1;
6 }
7 fprintf(fp, "%s\n", "TimeCost");
8 .....

```

If we already have a file, we can use mode 'a' to append new data after the data in that file.

```

1 FILE *fp;
2 fp = fopen("TimeCost.csv", "a");

```

```

3     if(fp == NULL) {
4         printf("Failed to open file\n");
5         return 1;
6     }
7     fprintf(fp,"%f,%f\n",Time1, Time2);
8     .....

```

CSV file using ',' to separate the data into blocks in a row and using '\n' to change a row.

2.6.5 Package the Testing

Since there are many experiments to be performed, I packaged the experiments into a function test(). The function takes as arguments the function to be tested (using the function pointer), the name of the function to be tested, the data, the height, the width of the image, and the kernel.

This function contains the process of padding, convolution operations, calculating the time costs during the operation, and outputting the result. (printing out the data or writing the data into a CSV file)

```

1 bool test(bool (*cnn)(const float* input_data, const size_t input_height, const size_t input_width,
2 const float* kernel_data, const size_t kernel_height,
3 const size_t kernel_width, float* output_data),const float* input_data, const size_t input_height,
4 const size_t input_width, const float* kernel_data,
5 const size_t kernel_height,const size_t kernel_width,const char* name)
6 {
7     bool flag = false;
8     //Padding
9     float*pad_data = padding(input_data,input_height,input_width,kernel_height,kernel_width);
10    //cnn
11    float* cnn_data = (float*)calloc(input_height*input_width,sizeof(float));
12    if(cnn_data==NULL)
13    {
14        printf("cnn_data calloc error!");
15        return false;
16    }
17    struct timeval start, end;
18    flag = gettimeofday(&start, NULL);
19    if (!flag)
20    {
21        printf("Error: gettimeofday()\n");
22        return false;
23    }
24    flag = cnn(pad_data, input_height, input_width, kernel_data, kernel_height, kernel_width, cnn_data);
25    if(!flag)
26    {
27        printf("cnn Error!");
28        return false;
29    }
30    flag = gettimeofday(&end, NULL);
31    if (!flag)
32    {
33        printf("Error: gettimeofday()\n");
34        return false;
35    }
36    //fprintf(fp,"%lf,", (end.tv_sec - start.tv_sec) +
37    (double)(end.tv_usec - start.tv_usec) / 1000000.0);
38    printf("%s\n",name);
39    printf("Time Cost : %lf s\n", (end.tv_sec - start.tv_sec) + (double)(end.tv_usec - start.tv_usec)/1000000.0);
40    printf("\n");
41    free(cnn_data);
42    free(pad_data);
43 }
44

```

2.6.6 Free the memory

The data of an image is usually very large. So we store the data in dynamic memories (The memory in a heap). Simply using free() to free the structure will not free the memory pointed to by the pointer inside the Image structure. Thus, there is a need for us to design a function to free the memory pointed to by the pointer inside the Image structure. By the way, the data of the kernel is usually not too big. Thus, we store the data of the kernel in static memories (The memory in the stack). There is no need to free the static memory manually, it will be free automatically when the program is finished.

```

1 bool Free_Image(Image* image)
2 {

```

```

3     if(image==NULL)
4     {
5         printf("NULL pointer!");
6         return false;
7     }
8     if(image->data1!=NULL)
9     {
10        free(image->data1);
11        image->data1=NULL;
12    }
13    if(image->data2!=NULL)
14    {
15        free(image->data2);
16        image->data2=NULL;
17    }
18    if(image->data3!=NULL)
19    {
20        free(image->data3);
21        image->data3=NULL;
22    }
23    return true;
24 }

```

2.7 Compile the file

This project includes a library called libjpeg and need to implement SIMD instructions and OpenMP instructions. Thus, there are servals things needed to link when compiling the code.

```

1 gcc project3demo.c Func.c -ljpeg -DWITH_AVX2 -mavx

```

3 Experiments & Analysis

3.1 Efficiency Analysis

3.1.1 The Dimension of Loops

In this experiment, we want to figure out the effect of the number of layers of the loop on the convolution speed.

Four functions are designed. For the convolution of an element, the row-index = (k / kernel-width) and the column-index = (k % kernel-width). For the convolution of the whole image, the row-index = (t / input-width) and the column-index = (t % kernel-width). In this way, we can lower the dimension of the loop.

```

1 bool cnnFunction_v1()
2 {
3     for(int j = 0; j < input_height; j++)
4     {
5         for(int i = 0; i < input_width;i++)
6         {
7             for(int m = 0; m<kernel_height;m++)
8             {
9                 for(int n = 0,k = 0; n < kernel_width; n++,k++)
10                {
11                    *(output_data+(j*input_width)+i) +=
12                    (*(input_data + (j * pad_width) + i + (m*pad_width) + n))
13                    * (*(kernel_data+k));
14                }
15            }
16        }
17    }
18 }
19 bool cnnFunction_v2()
20 {
21     for(int j = 0; j < input_height; j++)
22     {
23         for(int i = 0; i < input_width;i++)
24         {
25             for(int k =0; k<kernel_height * kernel_width;k++)
26             {
27                 *(output_data+(j * input_width)+i) +=
28                 (*(input_data + (j * pad_width) + i + ((k / kernel_width)*pad_width) + (k % kernel_width)))
29                 * (*(kernel_data+k));

```

```

30     }
31 }
32 }
33 }
34 bool cnnFunction_v3()
35 {
36     for(int t = 0; t < input_height*input_width; t++)
37     {
38         int j = t / input_width;
39         int i = t % input_width;
40         for(int m = 0; m<kernel_height;m++)
41         {
42             for(int n = 0,k = 0; n < kernel_width; n++,k++)
43             {
44                 *(output_data+(j*input_width)+i) +=
45                 (*(input_data + (j * pad_width) + i + (m*pad_width) + n))
46                 * (*(kernel_data+k));
47             }
48         }
49     }
50 }
51 bool cnnFunction_v4()
52 {
53     for(int t = 0; t < input_height*input_width; t++)
54     {
55         int j = t / input_width;
56         int i = t % input_width;
57         for(int k =0; k<kernel_height * kernel_width;k++)
58         {
59             *(output_data+(j * input_width)+i) +=
60             (*(input_data + (j * pad_width) + i + ((k / kernel_width)*pad_width) + (k % kernel_width)))
61             * (*(kernel_data+k));
62         }
63     }
64 }

```

In the experiment, we first use random() to create a certain size of random float image_data and a random 3x3 kernel_data. Then, we use test() to implement the convolution operation, calculate the time duration and store the time interval in the CSV file. The steps mentioned above will be carried out times for a certain size of image data size.

```

1  FILE *fp;
2  //open the file
3  fp = fopen("TimeCost.csv", "a");
4  for(int num = 1; num<=10; num++)
5  {
6      //create random number
7      float* image_data = (float*) malloc(128*128*sizeof(float));
8      image_data = Random(128,128);
9      float* kernel_data = (float*)malloc(3*3*sizeof(float));
10     kernel_data = Random(3,3);
11     fprintf(fp,"%d,",num);
12     test(cnnFunction_v1,image_data,128,128,kernel_data,3,3,fp,"cnn_v1");
13     test(cnnFunction_v2,image_data,128,128,kernel_data,3,3,fp,"cnn_v2");
14     test(cnnFunction_v3,image_data,128,128,kernel_data,3,3,fp,"cnn_v3");
15     test(cnnFunction_v4,image_data,128,128,kernel_data,3,3,fp,"cnn_v4");
16     fprintf(fp,"\n");
17     //free the memory
18 }
19 //close the file

```

The results are as follows.

对于不同的数据量，不同循环层数的卷积所用时间比较						
Function	128x128	256x256	512x512	1024x1024	2048x2048	4096x4096
cnn_v1(s)	0.0005945	0.0024136	0.0095929	0.0380648	0.1511122	0.6040479
cnn_v2(s)	0.0006313	0.0026448	0.0104002	0.0415641	0.1657193	0.6681645
cnn_v3(s)	0.0006612	0.0026523	0.0107056	0.042967	0.1713406	0.6866863
cnn_v4(s)	0.0006731	0.002731	0.011079	0.0443474	0.1773043	0.7116192

The results of the experiment are different from what we just envisioned at the beginning. The reason for the results is that the reduction in consumption through cycle reduction has no way to compensate for the increase in computation time due to the cycle reduction.

3.1.2 Unloop

In this experiment, we want to figure out the effect of the number of loops on the convolution speed.

Two Functions are designed. The first one is the same as the `cnnFunction_v1` mentioned above. The second one calculates 8 elements in one loop. By the way, in order to deal with all the elements in the case that the input width is not a multiple of 8, another block of loops is set to deal with the remaining item(s).

```
1 bool cnnFunction_v1();  
2 bool cnn_unloop();
```

The way of conducting the experiment is the same as the former experiment.

The results are as follows.

对于8的倍数的不同的数据量，解循环与不解循环卷积的时间比较						
Function	128x128	256x256	512x512	1024x1024	2048x2048	4096x4096
cnn_v1(s)	0.0006978	0.0024667	0.0095909	0.0379422	0.152315	0.6049463
Unloop(s)	0.0005326	0.0019161	0.0075812	0.0302822	0.1206239	0.4829237
速度提升率	0.236744053	0.2232132	0.209542379	0.201886027	0.208062896	0.201708152

When the size of the data is a multiple of 8, the speed improvement rate of the unloop function is approximately 21%.

对于8的倍数加7的不同的数据量，解循环与不解循环卷积的时间比较						
Function	15x15	73x73	135x135	263x263	519x519	1031x1031
cnn_v1(s)	0.0000102	0.0002474	0.0007386	0.0025631	0.0098216	0.0385328
Unloop(s)	0.0000083	0.0002042	0.000591	0.0020271	0.007841	0.0307507
速度提升率	0.18627451	0.174616006	0.19983753	0.209121767	0.201657571	0.201960408

When the size of the data is not a multiple of 8, the speed improvement rate of the unloop function is relatively low, especially when the size of the data is small. (approximately 18%) However, as the size of the data increases, the speed improvement gradually approaches 21%. (The average speed improvement rate above)

The reasons for the results are that when we reduce the number of loops, we reduce the cost of calling the loop (The cost of repeated calculations and comparisons). Thus, the speed improves. When the size of the data is not a multiple of 8, it costs time to finish the convolution operation in the normal way and some time to check.

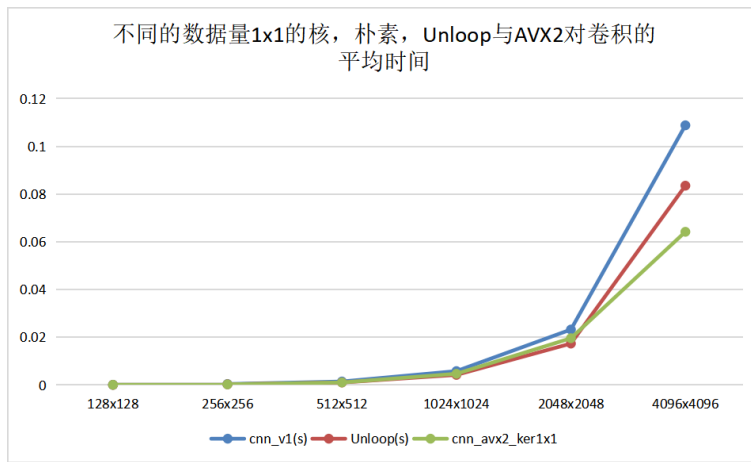
3.1.3 SIMD Analysis

This experiment aims to compare the efficiency between the Unloop function and the SIMD function in the case of 1x1, 3x3, and 5x5 kernels.

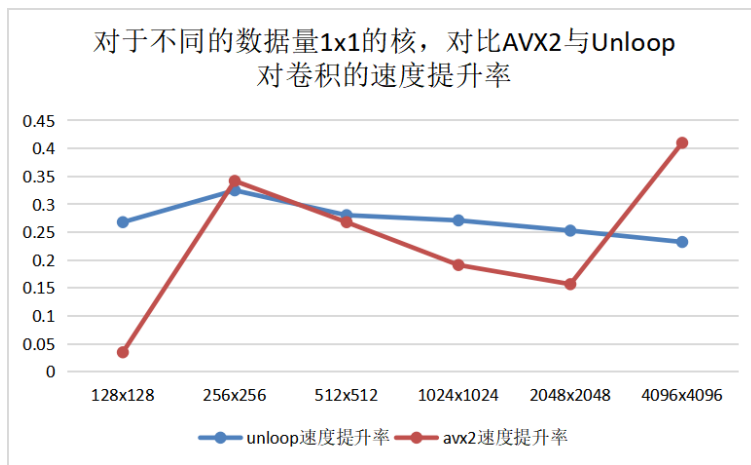
In the case of a 1x1 kernel, we compare the time cost of the following 3 functions.

```
1 bool cnnFunction_v1();  
2 bool cnn_unloop();  
3 bool cnn_avx2_kernel1x1();
```

The results are as follows.



Through the figure above, we can find that, as the size of the data increases, the differences between the time cost of the three functions become apparent.



Through the figure above, we can find that, in the case of a 1x1 kernel, when the size of the data is small, the speed improvement of Unloop function is bigger than that of SIMD. However, when the size of the data is big, the speed improvement of the SIMD function becomes bigger than those of Unloop function.

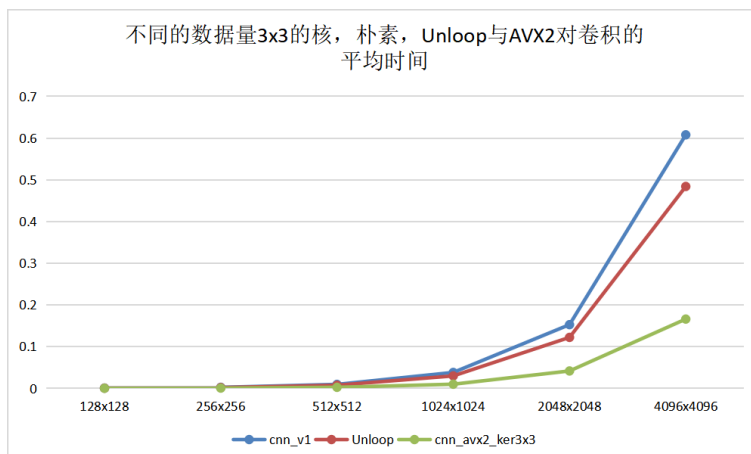
In the case of a 3x3 kernel, we compare the time cost of the following 3 functions.

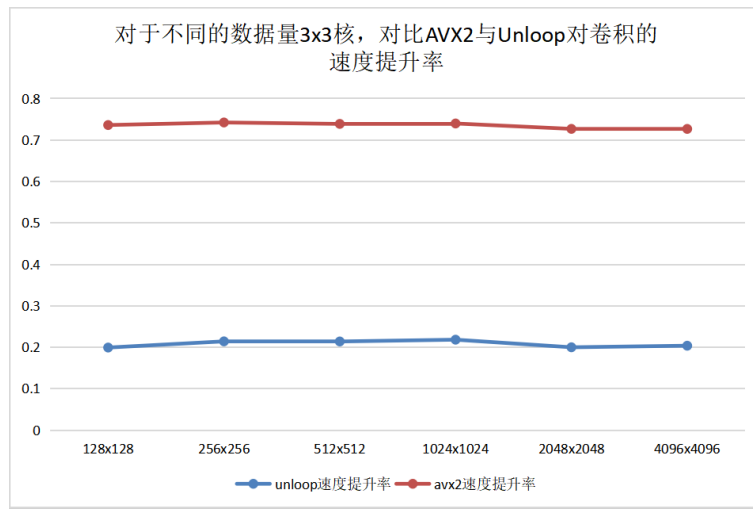
```

1 bool cnnFunction_v1();
2 bool cnn_unloop();
3 bool cnn_avx2_kernel3x3();

```

The results are as follows.



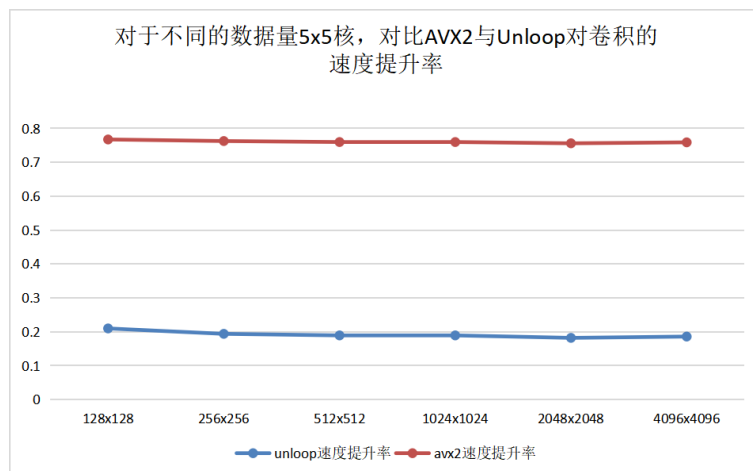
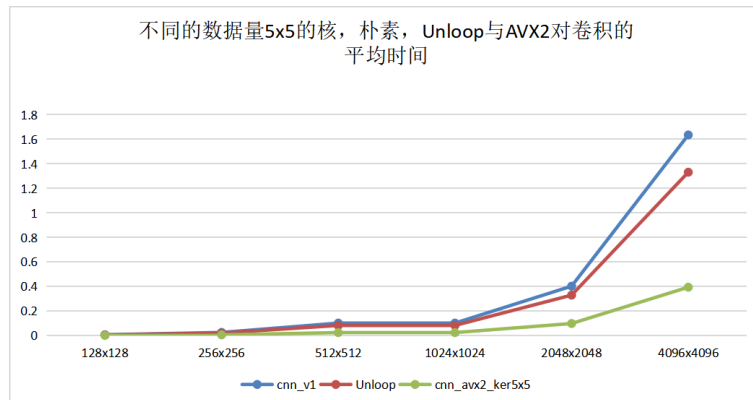


Through the figures above, we can find that, in the case of a 3x3 kernel, the speed improvement rate of the SIMD function is constantly bigger than that of the Unloop function. And the speed improvement rate of both SIMD and Unloop remain relatively constant among the different size of the data. (about 73% and 20% separately)

In the case of a 5x5 kernel, we compare the time cost of the following 3 functions.

```
1 bool cnnFunction_v1();
2 bool cnn_unloop();
3 bool cnn_avx2_kernel5x5();
```

The results are as follows.



Through the figures above, we can find that the result in the case of a 5x5 kernel is similar to the result in the case of a 3x3 kernel. (about 76% and 20% separately)

SIMD is a single instruction that can operate multiple data elements at the same time, making full use of register space. It saves processor time and resources by reducing the number of data reads and thus the number of memory accesses. However, implementing a SIMD has some costs in the preparation.

When the size of the kernel is small, the speed improvement due to the SIMD will not compensate for the costs of the preparation, resulting in the speed improvement rate of SIMD being smaller than Unloop at a small data size. When the size of the data gets bigger, the advantage of SIMD will be apparent.

When the size of the kernel is big, it is apparent that SIMD has a great advantage in efficiency compared to Unloop.

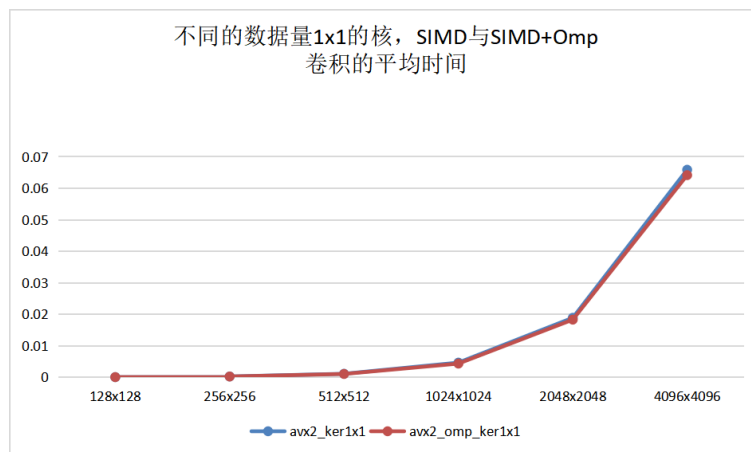
3.1.4 OpenMP Analysis

This experiment aims to compare the efficiency between SIMD and SIMD with OpenMP in the case of 1x1, 3x3, and 5x5 kernels.

In the case of a 1x1 kernel, we compare the time cost of the following 2 functions.

```
1 bool cnn_avx2_kernel1x1();
2 bool cnn_avx2_omp_kernel1x1();
```

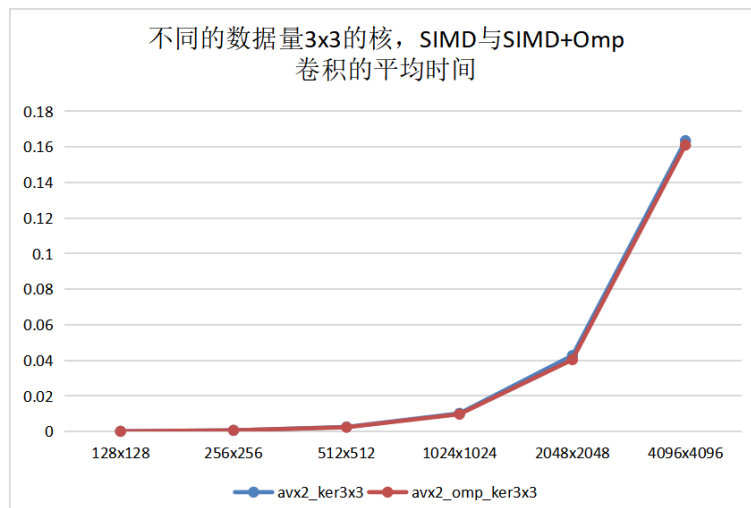
The results are as follows.



In the case of a 3x3 kernel, we compare the time cost of the following 2 functions.

```
1 bool cnn_avx2_kernel3x3();
2 bool cnn_avx2_omp_kernel3x3();
```

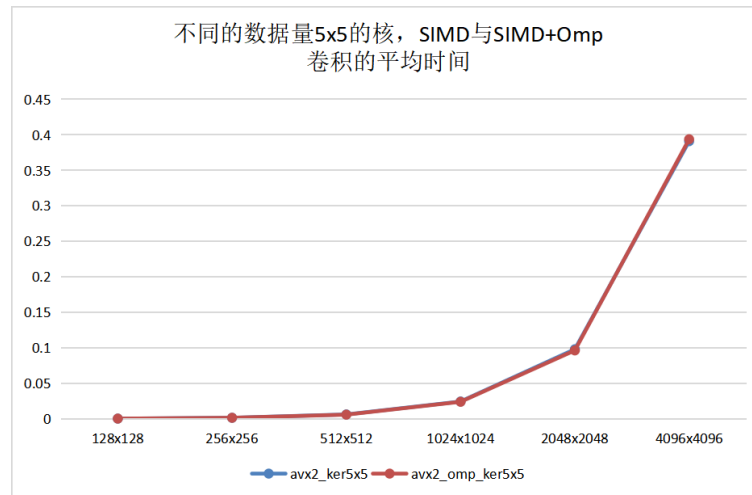
The results are as follows.



In the case of a 5x5 kernel, we compare the time cost of the following 2 functions.

```
1 bool cnn_avx2_kernel5x5();
2 bool cnn_avx2_omp_kernel5x5();
```

The results are as follows.



Through the figures listed above, we find that no matter the size of the data or the size of the kernel, combining OpenMP with SIMD will not bring extra efficiency compared to the single SIMD operation.

The reason for the result is that when performing convolution operations via SIMD, it is necessary to chunk the input data, which already adds additional copy and reorder overhead. In addition, in OpenMP, the input data needs to be copied again to avoid competition between threads, which may take up a lot of memory space and thus affect the efficiency of the program.

3.1.5 -O3

We compare the speed of the convolution operation with and without the -O3 instruction.

The results are as follows

	O3的优化			
	unloop	omp_unloop	SIMD	SIMD_omp
加-O3	0.0596392	0.0600104	0.0191133	0.0166557
不加-O3	0.4940981	0.4931542	0.1663743	0.1611616
速度提升率	0.879296844	0.878313112	0.885118675	0.89665218

The speed improvement rate of adding the -O3 command is very high, larger than that of the SIMD. (around 70%)

So if you want to improve the speed of your program, maybe you can try using this command.

3.2 Correctness

3.2.1 Qualitative Method

We output the images after the convolution operation in C and Matlab to see if they are the same.

Convolution layer in C. (The example is an image of one channel)

```
1 //Input the Image
2 Image input_image;
3 input_jpeg_image("gray.jpg",&input_image);
```

```

4 //Initialize the kernel
5 Kernel kernel;
6 float kernel_data[9] = {-1, 0, 1, -2, 0, 2, -1, 0, 1};
7 flag=kernel_init(&kernel,3,3,input_image.channels,kernel_data,NULL,NULL);
8 //Create the output image
9 Image output_image;
10 //Convolution operation
11 Convolution(&input_image,&kernel,&output_image);
12 //output the image
13 flag = output_jpeg_image("grayoutput.jpg",&output_image);
14 //free the memory
15 Free_Image(&input_image);

```

Convolution layer in MATLAB. (The example is an image of one channel)

```

1 img = imread('gray.jpg');
2 kernel = [-1 0 1; -2 0 2; -1 0 1];
3 result = conv2(img, kernel, 'same');
4 imshow(result);

```

The result in C and the result in MATLAB are as follows.



Figure 1: Output from C

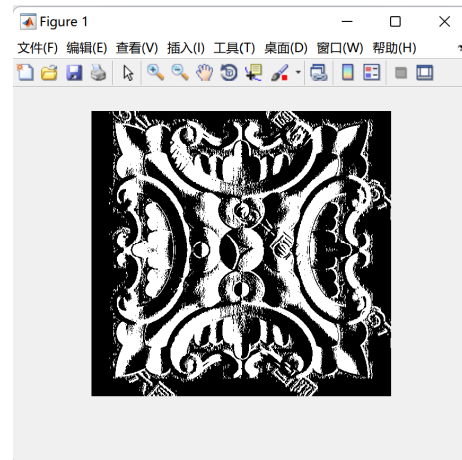


Figure 2: Output from MATLAB

From the result, it is likely that the Convolution layer in C is correct.

4 Difficulties & Solutions

4.1 Difficulties in input and output the image

Reading images in C is difficult as it does not contain any related function. This project uses a library called 'libjpeg' to read the data of a jpg image.

First, we need to download the library in Ubuntu and include the head file when running the program. (The instruction and the head file are as follows.)

```

1 sudo apt-get install libjpeg-dev
2 #include <jpeglib.h>

```

Then, we need to learn the grammar of decompressing and compressing the jpg. (What I have learned is mentioned in the 'input and output image' of the Code Implementation part)

4.2 Difficulties in Checking the correctness of the convolution operation

Checking if our convolution operation in C is difficult. Not only does it have no function of convolution operation, but it is also very difficult to find a library linked by C that can operate the convolution operation.

4.2.1 Attempt Intel Math Kernel Library

By searching on the internet, I found that the Math Kernel Library can conduct the convolution operation in C. However, having spent two to three hours, I still failed to install the library in my Ubuntu. Considering the time limit, I quit this method.

4.2.2 Attempt to check the correctness quantitatively

By searching on the internet, I found that we can operate convolution operations in MATLAB and can write the data of the images after the convolution operation in C and Matlab into Txt files, and use VSCode's file-matching feature to check if the results are correct.

Output data in C.

```
1 FILE *fp;
2 fp = fopen("C_output.txt", "w");
3 for(int i=0;i<input_image.height*input_image.width*input_image.channels;i++)
4     fprintf(fp,"%f\n",*(input_image.data1+i));
```

Output data in MATLAB.

```
1 fileID = fopen('M_output.txt', 'w');
2 fprintf(fileID, '%d\n', result);
3 fclose(fileID);
```

However, the datatype of storing the image data in C through libjpeg differs from MATLAB. The original datatype of the jpg image reading by libjpeg is unsigned char (I make a type conversion from unsigned char to float in order to satisfy the requirement) and the datatype of the jpg image reading by MATLAB is an integer (which contains the negative numbers). Without the relevant knowledge of the Algorithm to read the image in both libjpeg and MATLAB and without plenty of time to learn, I quit this method.

4.3 Difficulties in compiling the file using CMakeLists

In this project, we invoke a library called 'libjpeg' and use several extra commands to implement the SIMD operation. Thus, it is a little bit of trouble to write all the commands in the command window by ourselves. Writing a CMakeLists.txt is a good way to solve this problem.

However, although I successfully compile my CMakeLists.txt, it still can not work correctly, I failed to fix the problem before the deadline.

Here is my CMakeLists.txt.

```
1 cmake_minimum_required(VERSION 3.12)
2
3 add_definitions(-DWITH_AVX2)
4 add_definitions(-mavx)
5 add_definitions(-ljpeg)
6 add_definitions(-O3)
7
8 set(CMAKE_CXX_STANDARD 11)
9
10 project(pro)
11
12 ADD_EXECUTABLE(pro project3demo.c Func.c)
13
14 include_directories(${JPEG_INCLUDE_DIR})
15 find_package(JPEG REQUIRED)
16 target_link_libraries(pro PUBLIC ${JPEG_LIBRARIES})
17
18 find_package(OpenMP)
19 if(OpenMP_CXX_FOUND)
20     message("OpenMP found.")
21     target_link_libraries(pro PUBLIC OpenMP::OpenMP_CXX)
22 endif()
```

5 Brief Summary

This project focused on the implementation of a convolutional layer through C language and improving the efficiency of the program. Through the process of making the project, I have some gains and many regrets as well.

The gains from this project.

First of all, this project gave me a preliminary understanding of neural convolutional networks (CNN), and a preliminary exposure to the field of image processing, which broadened my horizon.

Secondly, I gained a deeper understanding of the underlying nature of the language C and the underlying nature of computer programming. I learned that the process of computing is done by moving the data to be computed out of memory into registers and then computing in the registers. The operations in the registers are implemented using different instructions depending on the computer chip (SIMD), and multiple data can be manipulated at once to achieve acceleration of operations. For our modern computers, which are multi-core systems, we can use OpenMP to call multiple cores to perform operations at the same time. (be careful about conflicts when storing data)

The project leaves a lot to be desired.

Firstly, no explore additional ways to achieve speedup. By talking with others, I found that it is actually possible to do the convolution of a whole image as if it were the multiplication of a matrix. But this involves rearranging the data. It's actually a bit like the method used in my SIMD acceleration. If it could be converted into a multiplication, it could be verified correctly by matrix multiplication in the openblas library.

Secondly, no good way to verify correctness has been explored. I wanted to download an MKL library, but there was a problem and I couldn't download it. (But the problem is a package, maybe I can try to find the individual library and download it, and maybe it will work)

Thirdly, failed to write CMakeLists.txt correctly. There are some unknown bugs in it.

Fourthly, did not try when the size of 512 and 128 register variables to compare the speed change with the size of 256 register variables, is a pity

Fifthly, did not bother with the development board. Because no time is left at all.

Some bugs were recorded.

The first one is that when dealing with remainders for arrays that are not multiples of 8, I didn't see where the brackets were when I was accessing the data through the pointer plus the offset, which had a lot of brackets, resulting in an error in the data accessed and never found.

The second is when writing the return type function for bool, there is a function of return forgot to write. Then, when I called the function (return type is also bool) written just behind it, no matter what the result is, the latter function always returns a FALSE.

The third is that when writing formulas to calculate data, due to the failure to consider the priority correctly, which results in the wrong results, adding brackets can solve this problem.