
REPRODUCTION INSTRUCTIONS FOR THE FASTER R-CNN

Wenjun Ma
Yunnan University
mwj22@stu.ynu.edu.cn

June 11, 2025

ABSTRACT

This article provides an explanation of the second stage of the assessment. Here, I will elaborate on the initial design before the replication process, and describe my file organization structure and its contents. I will also introduce and analyze the problems I encountered during the replication process. Finally, I will summarize this replication and outline my next steps.

1 Introduction

The organizational structure of this explanation report is as follows: Section 2 introduces the preparations before the project implementation as well as the logic of code writing and testing. Section 3 presents the file structure of the project and briefly explains the functions of each file. Section 4 briefly introduces the related packages used in the project as well as the software and hardware environment. Section 5 elaborates and analyzes the problems I encountered. Section 6 provides a summary and next steps for this reproduction.

2 Initial design

Before starting the writing of this project, I decided to proceed in the following sequence: 1) Re-read the paper; 2) Read and study the relevant code (Faster R-CNN and Jittor); 3) Refer to the code and implement it by hand. Therefore, I found the following repositories and documents as the reference materials for my learning:

1. Jittor: a Just-in-time(JIT) deep learning framework
<https://github.com/Jittor/jittor>
2. JRS
<https://github.com/NK-JittorCV/nk-remote>
3. A Simple and Fast Implementation of Faster R-CNN
<https://github.com/chenyuntc/simple-faster-rcnn-pytorch>
4. Documents of Jittor
<https://cg.cs.tsinghua.edu.cn/jittor/assets/docs/index.html>
5. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks
https://github.com/shaoqingren/faster_rcnn
6. Detectron2
<https://github.com/facebookresearch/detectron2>

Among them, 1, 2 and 3 are the contents that I mainly referred to. Due to the absence of certain parts in the 4 (for example, `jt.concat` is not included in the document as shown in the Figure 1), 1 and 2 is basically used as a reference for the document.

The strategy for writing code For most of the code, Jittor is used for implementation (except when reading images, where `np.asarray` is still used in `TwoStageDectection/Utils/VOC_tools/data_utils.py`). In order to save the time for writing code and fully experience the usage of the Jittor framework, my code writing logic is as follows:



Figure 1: The result of the concat operation in the Jittor documentation search.

1. Directly invoke Jittor's API (such as `jt.ones`) or use existing code (such as `TwoStageDetection/Networks/Parts/RoiPool.py`.)
2. Although jittor has been integrated into the system, it is still recommended to use jittor for simple programming. For example, the network model can be implemented using jittor. For details, please refer to `TwoStageDectection/Networks/Parts/Backbone_VGG.py`.
3. Using `jt.code + cuda` for a simple implementation, as Listing 1

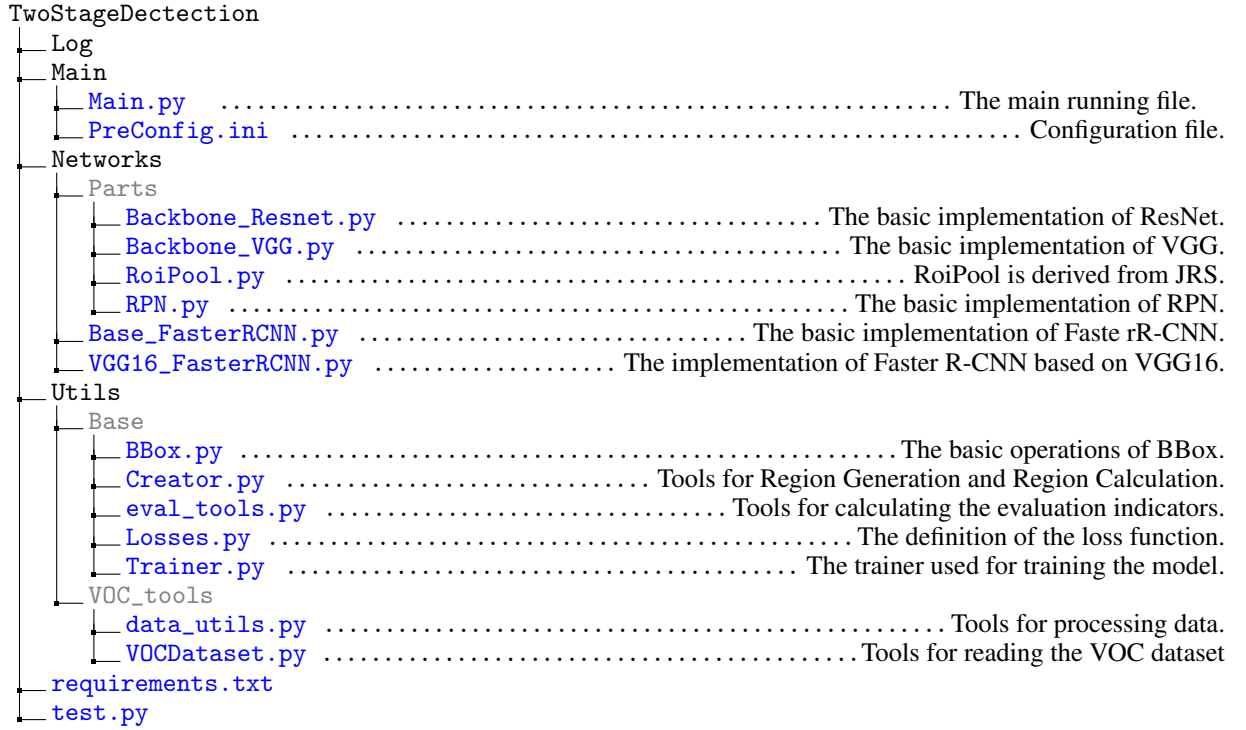
```
def nanmean(inputs):
    CUDA_HEADER = r'''
#include <cmath>
using namespace std;
'''
    CUDA_SRC = r'''
__global__ static void kernel1(@ARGS_DEF) {
    @PRECALC
    float sum = 0;
    int ing_sum = 0;
    for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < in0_shape0;
        i += blockDim.x * gridDim.x){
        __isnanf(@in0(i)) ? ing_sum += 1 : sum += @in0(i);
    }
    int all_len = in0_shape0 - ing_sum;
    all_len == 0 ? @out(0) = NAN : @out(0) = sum / all_len;
}
const int total_count = in0_shape0;
const int thread_per_block = 1024;
const int block_count = (total_count + thread_per_block - 1) / thread_per_block;
kernel1<<<block_count, thread_per_block>>>(@ARGS);
'''
    return jt.code(shape=(1,), dtype='float32', inputs=[inputs],
        cuda_header=CUDA_HEADER, cuda_src=CUDA_SRC)
```

Listing 1: A simple implementation of nanmean

The strategy for testing code The strategy of the code testing is relatively simpler. Some independent arithmetic operations are tested separately, while those strongly related to model training and testing are checked simultaneously during the overall operation. (When considering this, I thought the overall project was not difficult, so I made this somewhat hasty decision.)

3 File structure

This section presents the file structure of the entire project and briefly describes the contents of each file.



4 Environment configuration

This section presents the environment in which the project operates and the logic for the model's operation.

4.1 Environment configuration

Because the Conda environment is rather messy, we did not directly generate environment.yaml. Instead, we used pipreqs to generate requirements.txt as shown below.

```

jittor==1.3.9.14
numpy==2.3.0
Pillow==11.2.1
tqdm==4.66.2

```

Listing 2: requirements.txt

System environment and hardware:

- CUDA 12.4
- GPU RTX4090D
- ubuntu22.04

5 The problems encountered and the solutions.

5.1 The problems of environment configuration

5.1.1 Absence of files

5.1.2 Other problems

1. Constantly try out different versions of Jittor. (Successfully configured the environment on a certain Windows system using version 1.3.6.6.)
2. Constantly try out different versions of CUDA.
3. Reinstall the operating system.

5.2 The problems of code

This section presents some of the problems I encountered. For each problem, I will provide a wrong example, my provisional solution, the current situation of the problem and the analysis process, and I haven't had the chance to delve into the content of the source code that is built using C yet, so the source code only shows the part implemented in Python.

In order to make the problem occur more consistently, the data shuffling has been disabled here.

5.2.1 Problem0

Example & (provisional) Solution This issue originates from line 120 of TwoStageDetection/Utils/Base/Creator.py. The error code and the corrected code are shown as follows.

```
wrong code:
    rand_index = jt.randperm(len(neg_index) - n_neg)
right code:
    rand_index = jt.randperm(len(neg_index) - n_neg.item())
```

Listing 3: Problem0

The error message is shown in Figure 3.

Reproducibility True

Problem Analysis This issue is named in this way because it was originally functional, but suddenly it stopped working. Now, let's analyze this situation. First, I viewed the source code of the relevant function, shown as Listing 4.

```
def randperm(n, dtype="int32"):
    key = jt.random((n,))
    index, _ = jt.argsort(key)
    return index.cast(dtype)

def random(shape, dtype="float32", type="uniform"):
    for dim in shape:
        if dim < 0:
            raise RuntimeError(f"Trying to create tensor with \
negative dimension {dim}: {shape}")
    ret = ops.random(shape, "float32", type)
    amp_reg = jt.flags.amp_reg
    if amp_reg:
        if amp_reg & 16:
            if amp_reg & 1:
                if ret.dtype != "float32":
                    return ret.float32()
            elif amp_reg & 2:
                if ret.dtype != "float16":
                    return ret.float16()
    return ret

def subtract(self, y: Var)-> Var:
```

Listing 4: Source0

Then we will print out the inputs for both scenarios, shown as Listing 5.

```
wrong code:
    AssertionError: jt.Var([24], dtype=int32)
right code:
    AssertionError: 277
```

Listing 5: Pttout0

We can observe that the output of both methods is of the int type. According to the error message, the type of `jt.Var` is also `int32`. Referring to the error information, there is an `int64` that does not match the basic Python type `int`. This

[illegible]

Figure 3: Error message1_0.

error may occur because the source code failed to correctly read the data of var here, and instead mistakenly converted the entire type into a type similar to `np.int64`. Another supporting basis for this inference is that the dependency package of jittor contains numpy, which may have performed related operations.

5.2.2 Problem1

Example & (provisional) Solution This issue has not been resolved yet. Here, only the problem description and error information are provided. This problem is also the biggest obstacle I encountered during the implementation process.

Reproducibility True

Problem Analysis The problem directly related to the error message is caused by `TwoStageDetection/Networks/VGG16_FasterRCNN.py`. There is no data in the input pool. After conducting a thorough investigation, I found that this step shown as Listing 6 in line 202 of `TwoStageDetection/Utils/Base/Creator.py` resulted in the data being empty. There were no data that met the filtering requirements.

```
keep = jt.where((hs >= min_size) & (ws >= min_size))[0]
```

Listing 6: Code1

When data shuffling was enabled, I printed the contents of the data multiple times. On one occasion, I discovered that there were `nan` values in the data.

Regarding the methods for handling nan values in the data, I made the following attempts for processing, but all failed.

```

[1 0008 22:37:50.819071 84 cuda_flags.cc:49] CUDA enabled.
load data
model construct completed
log
Epoch: 1/20 : 0%|
Caught segfault at address 0x10, thread_name: 'C5', flush log...
Segfault, exit
[e 0608 22:37:51.791066 84 parallel_compiler.cc:324] Segfault happen, main thread exit

```

Figure 4: Error message1_1.

1. Modified the initialization method of the convolution layer.
2. Adjusted the learning rate.
3. Applied gradient clipping.
4. Disabled data shuffling and swapped or deleted the data.
5. Limited and adjusted the numerical range of all bounding box calculation values.
6. Used the official export `JT_CHECK_NAN=1`, `export trace_py_var=3` provided by the Jittor to locate the error position.

Based on this, I plan to use `assert not jt.isnan(features).any()`, `"!!!"` to locate and find out exactly where the problem lies. However, a segmentation fault occurred, as shown in the Figure 4.

The official debugging method(`export debug=1`, `export gdb_attach=1`) provided by the Jittor was used to detect the segmentation fault, then the investigation was carried out. Once before, there was a prompt indicating that the accessed memory was already released, but no screenshot was taken at that time. From now on, if we test again, the only possible outcomes would be either a system crash or no additional error messages.

Based on this, I can infer that there is a problem with Jittor's memory management. It is very likely that there is an issue with the counting or marking during garbage collection, which leads to the incorrect release of memory.

```

wrong code:
height = src[:, 2] - src[:, 0]
width = src[:, 3] - src[:, 1]
ctr_y = src[:, 0] + 0.5 * height
ctr_x = src[:, 1] + 0.5 * width

basey0 = dst[:, 0]
basey1 = dst[:, 2]
basex0 = dst[:, 1]
basex1 = dst[:, 3]
right code:
height = src[:, 2] - src[:, 0]
width = src[:, 3] - src[:, 1]
ctr_y = src[:, 0] + 0.5 * height
ctr_x = src[:, 1] + 0.5 * width

basey0 = dst[... , 0]
basey1 = dst[... , 2]
basex0 = dst[... , 1]
basex1 = dst[... , 3]

```

Listing 7: Problem2.

5.2.3 Problem2

Example & (provisional) Solution This issue is quite strange and cannot be reproduced now. This problem is presented in `TwoStageDectection/Utils/Base/BBox.py`, the detailed information shown in listing 7. The above

```

>>>
>>> a = jt.array([5, 5])
>>> print(a)
jt.Var([0 0], dtype=int32)
>>>

```

Figure 5: Problem3.

code snippet would trigger a segmentation fault, while the below one runs correctly. Just now, when I attempted to reproduce it, I couldn't get the error message again. I'm just making a note here.

Reproducibility False

5.2.4 Problem3

Example & (provisional) Solution This was a problem that occurred during a simple test after installing Jittor on Windows. This issue was somewhat beyond my comprehension. The detailed situation of the problem is shown in Figure 5. After seeing this problem, I decided to stop configuring the environment on the Windows platform and reinstalled the system as Ubuntu22 to continue the environment configuration.

Reproducibility False

6 Conclusion

Let me briefly talk about my gains and my next goals. The most direct gain is a deeper and broader understanding of Python and Jittor. Indirectly, it has enriched my experience in reading source code and improved my speed and level of reading source code. The next plan is to learn more about Jittor and CUDA programming, try to refactor my previous code using Jittor, and do more learning.