

# ConvBench

A Comprehensive Convolution Performance Evaluation Benchmark

**Lucas F. A. e Silva** <[lucas.silva@ic.unicamp.br](mailto:lucas.silva@ic.unicamp.br)>

Computer Systems Laboratory (LSC)  
Institute of Computing (IC)  
University of Campinas (UNICAMP)

February, 2024



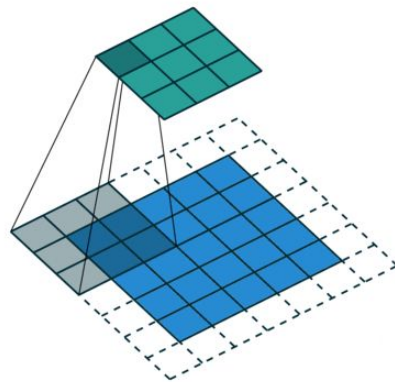
# Table of Contents

- Overview
- 1. Operation Set Construction and Filtering
- 2. Convolution Algorithm Definition and ConvBench Structure
- 3. Results Summary, Dataframe Visualization and Plotting
- Conclusion

## SECTION

# Overview

- Main operation on most DL models
  - Largest share of a CNN execution (~90%).
  - Convolution is a computationally expensive operation.
- Different Convolution Algorithms have been developed:
  - Naive
  - IM2COL transformation + GEMM
  - Winograd
  - Direct Convolution: SConv, YAConv, and so on...



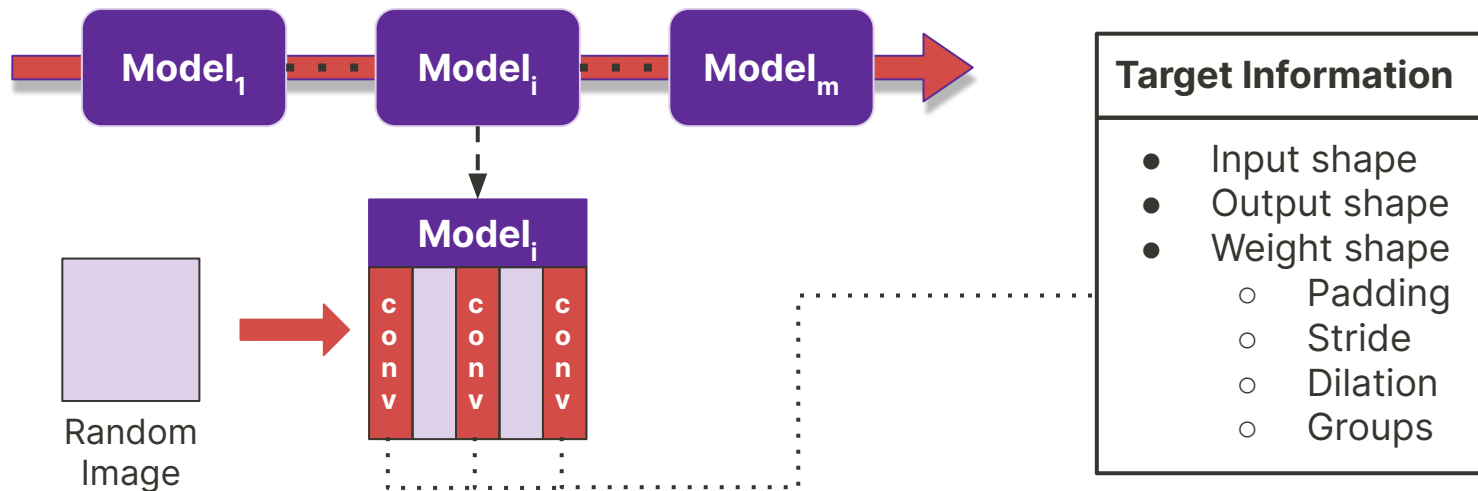
**Figure:** Convolution Animation, obtained from:  
[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

## SECTION

# Operation Set Construction and Filtering

# Operation set construction

- 1<sup>st</sup> step of the Convolution Benchmark: **Operation Set Construction**
- Hugging Face TIMM's DL model collection – 1017 models
  - CNNs, ViTs, and so on.



- Unique keys identifying each convolutions
- **9011 different convolution operations (December/2023).**
  - 5481 elementwise convolutions.
  - 3530 not elementwise.
    - 2269 grouped convolution
    - 17 dilated convolution
    - 93 rectangular convolution (only filters).
    - 2 rectangular convolution (filter and input)
    - 1149 regular convolution

## Operation Set Description:

- Input spatial size: **4×4 up to 1024×1024**
- Output spatial size: **1×1 up to 400×400**
- Channels: **3 up to 2048**
- Kernel number of filters: **8 up to 4096**

## Number of Float Point Operations (#FLOPS):

- ranges from **903 #MFLOPS** up to **43 #GFLOPS**



Filtering convolution operation set.

- ✓ Operation set organized saved as pickle file.
- ✓ Transform to pandas Dataframe and **Filter** by a just composed filtering string.

Basic Filtering

Advanced Filtering

☐ Elementwise Convolutions

☒ Normal Convolutions (kernel > 1x1)

☐ Grouped Convolutions

☐ Dilated Convolutions

☒ Only Square Filters

☒ Only Square Images

# Operation set construction

10

Filtering

✓ Op

✓ Tra

ng.

Basic Filtering

Advanced Filtering

Enable Advanced Filtering ?

Parameter	Relation	Value
<input type="checkbox"/> Filter by Channels (C)	< ▾	0
<input type="checkbox"/> Filter by Image Height (Hi)	< ▾	0
<input type="checkbox"/> Filter by Image Width (Wi)	< ▾	0
<input type="checkbox"/> Filter by Dimensions (D)	< ▾	0
<input type="checkbox"/> Filter by Output Height (Ho)	< ▾	0
<input type="checkbox"/> Filter by Output Width (Wo)	< ▾	0
<input type="checkbox"/> Filter by Kernel Height (Hk)	< ▾	0
<input type="checkbox"/> Filter by Kernel Width (Wk)	< ▾	0
<input type="checkbox"/> Filter by Horizontal Stride (Hs)	< ▾	0
<input type="checkbox"/> Filter by Vertical Stride (Ws)	< ▾	0
<input type="checkbox"/> Filter by Top Padding (Hpt)	< ▾	0
<input type="checkbox"/> Filter by Bottom Padding (Hpb)	< ▾	0
<input type="checkbox"/> Filter by Lef Padding (Wpl)	< ▾	0
<input type="checkbox"/> Filter by Right Padding (Wpr)	< ▾	0
<input type="checkbox"/> Do it have bias ?		

## SECTION

# Convolution Algorithm Definition and ConvBench Structure

ConvBench was designed as a including **Header file**.

- It Inherits from a *timing.h* class;
- And implements all the benchmark methods in a templated class, **ConvBench**.
  - Instances to be benchmarked must **inherit** from ConvBench class and:
    - Properly **override** the *convolution* function.
    - Properly **override** the *convolution\_baseline* function.
    - And **Implement** a simple *main* function.
  - ConvBench *convset\_exec()* takes care of the convolution algorithm assessment

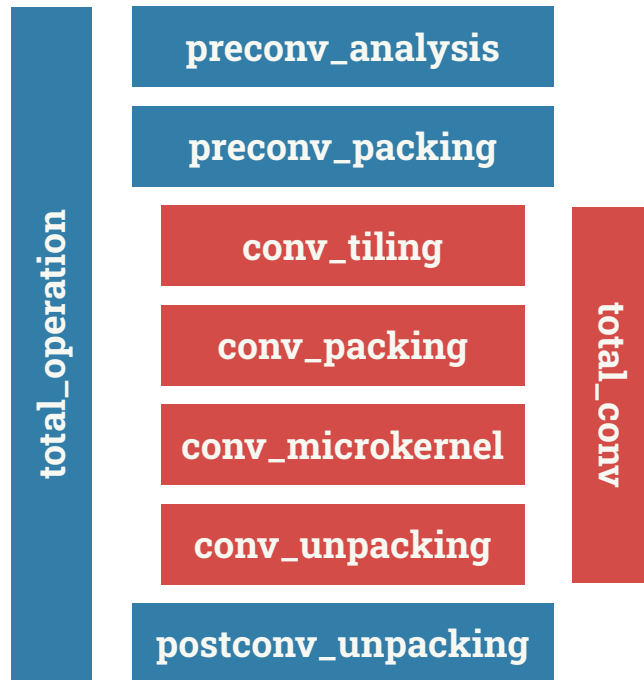
Standardized timing nomenclature.

Just enclose code snippets  
between the:

```
<time_name>_start()
```

```
<time_name>_update()
```

And you are ready to go! At the end  
specialized plots will be available.



ConvBench class is a templated class by  $T$  (Theoretically, accepting any kind of data type).

- ***void convset\_load(csv\_filename):***
  - Load the filtered convolution operation set.
- ***void convdata\_gen(datagen\_strategy):***
  - Populate the data buffers (**Random/Constant**).
- ***void convset\_exec(running\_strategy):***
  - The main execution procedure (**Correctness/Direct/Baseline**).
- ***virtual void convolution(args ...):***
  - The main convolution algorithm to be assessed.
- ***virtual void convolution\_baseline(args ...):***
  - The baseline convolution algorithm to be compared.

# ConvBench Structure • Code example

15

```
#include "convbench.h"

template <T>
class Inherited_Conv : public ConvBench<T>{

    //constructors and destructors
    Inherited_Conv();
    ~Inherited_Conv();

    //Conv algorithms
    void convolution(args ... ) override {
        // Main convolution algorithm implementation
    }

    void convolution_baseline(args ... ) override {
        // baseline convolution algorithm implementation
    }

    //Some other stuff if needed

}

int main() {
    // Object instantiation
    Inherited_Conv<T> bench = Inherited_Conv<T>();
    // Benchmark execution
    bench.convset_exec(args ... );
    return 0;
}
```

# ConvBench Structure • Code example

16

```
#include "convbench.h"

template <T>
class Inherited_Conv : public ConvBench<T>{

    //constructors and destructors
    Inherited_Conv();
    ~Inherited_Conv();

    //Conv algorithms
    void convolution(args ... ) override {
        // Main convolution algorithm implementation
    }

    void convolution_baseline(args ... ) override {
        // baseline convolution algorithm implementation
    }

    //Some other stuff if needed

}

int main() {
    // Object instantiation
    Inherited_Conv<T> bench = Inherited_Conv<T>();
    // Benchmark execution
    bench.convset_exec(args ... );
    return 0;
}
```



# ConvBench Structure • Code example

17

```
#include "convbench.h"

template <T>
class Inherited_Conv : public ConvBench<T>{

    //constructors and destructors
    Inherited_Conv();
    ~Inherited_Conv();

    //Conv algorithms
    void convolution(args ... ) override {
        // Main convolution algorithm implementation
    }

    void convolution_baseline(args ... ) override {
        // baseline convolution algorithm implementation
    }

    //Some other stuff if needed

}

int main() {
    // Object instantiation
    Inherited_Conv<T> bench = Inherited_Conv<T>();
    // Benchmark execution
    bench.convset_exec(args ... );
    return 0;
}
```

# ConvBench Structure • Code example

18

```
#include "convbench.h"

template <T>
class Inherited_Conv : public ConvBench<T>{

    //constructors and destructors
    Inherited_Conv();
    ~Inherited_Conv();

    //Conv algorithms
    void convolution(args ... ) override {
        // Main convolution algorithm implementation
    }

    void convolution_baseline(args ... ) override {
        // baseline convolution algorithm implementation
    }

    //Some other stuff if needed

}

int main() {
    // Object instantiation
    Inherited_Conv<T> bench = Inherited_Conv<T>();
    // Benchmark execution
    bench.convset_exec(args ... );
    return 0;
}
```

## SECTION

# Results Summary, Dataframe Visualization and Plotting

ConvBench returns a CSV file of timing measurements.

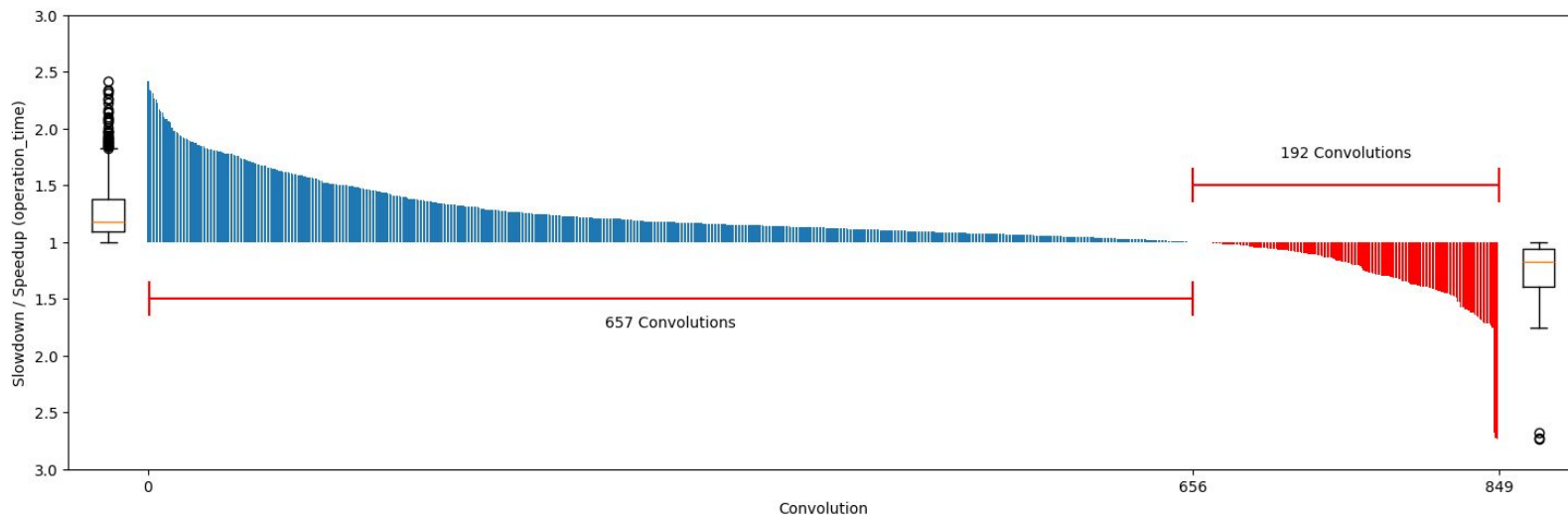
Thus, you can:

- ✓ Simply open it with some spreadsheet software;
- ✓ Open as a **pandas dataframe**;
- ✓ Filter, reshape, inspect.
- ✓ **Plot the results** with matplotlib.
- ✓ And so on ....

For example, we added the following automatically analysis:

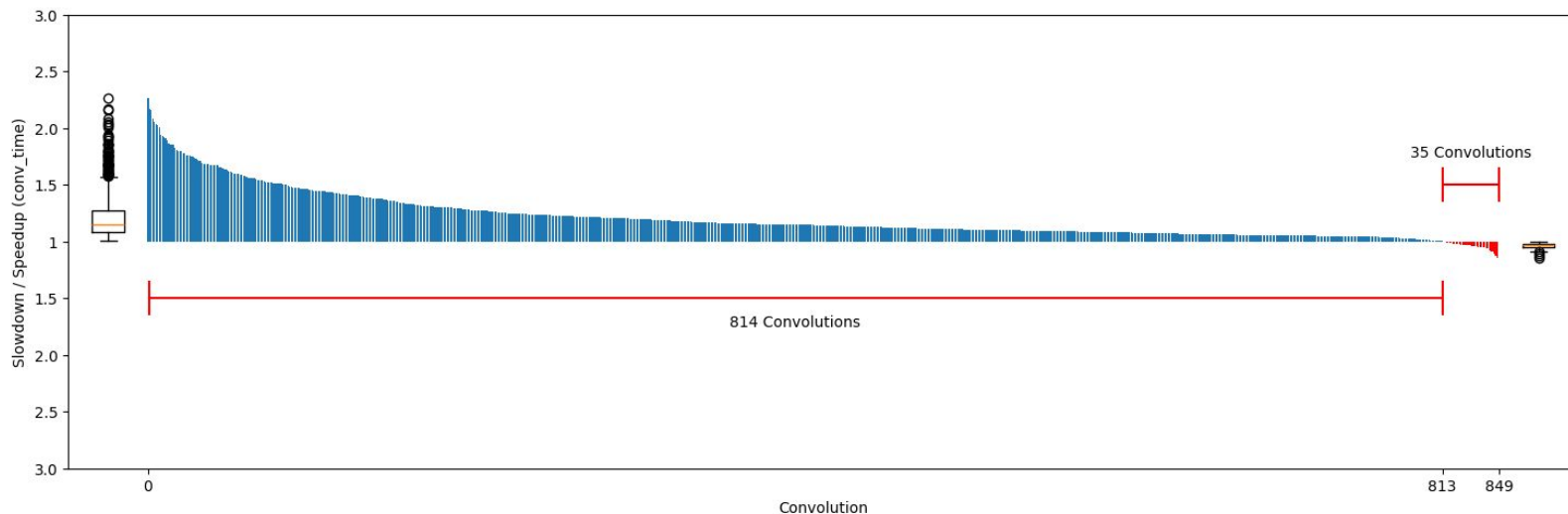
Normal convolution operations (filtering) speedup distribution (analysis)

**Total Operation Time** speedup Analysis.

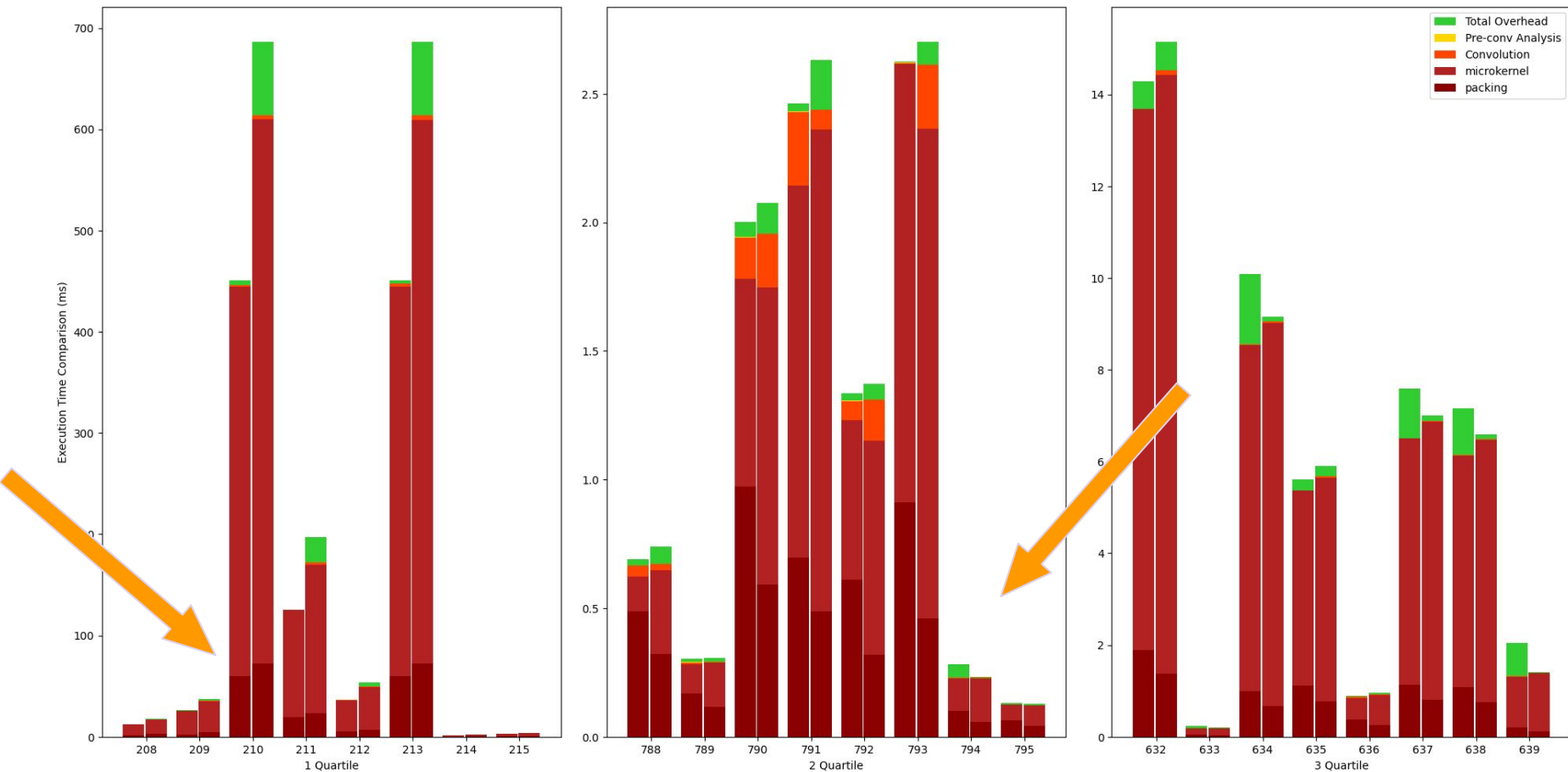


Normal convolution operations (filtering) speedup distribution (analysis)

**Total Convolution Time** speedup analysis.



Execution time breakdown graph.





Execution time breakdown graph.

The conv-packing step is **58,09%** faster than baseline on positive speedup cases.

However, on slowdown cases, the conv-packing step is **68,3%** slower than main algorithm.

## SECTION

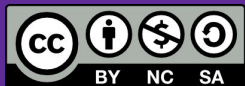
# Conclusion

- ConvBench: an end-to-end platform for evaluating convolution algorithms.
- Fair comparison: timing class.
- A throughout assessment of 2D convolution operations (9K+ operations) + Filtering mechanism
- An automatic experimental procedure + result summarization in tables and plots

# Thank You!

**Email**  
**GitHub**

lucas.silva@ic.unicamp.br  
<https://github.com/LucasFernando-aes>



This work is licensed under CC BY-NC-SA 4.0.

