

# GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal

Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí  
Universitat Politècnica de València, Spain  
Universidad de Córdoba, Spain.

**Presenter: Jie Lei**  
**Universitat Politècnica de València, Spain**



The authors gratefully acknowledge funding from  
- European Union's Horizon2020 Research and Innovation program under the Marie Skłodowska Curie Grant Agreement No. 956090 (APROPOS, <http://www.apropos-itn.eu/>).  
- European High-Performance Computing Joint Undertaking (JU) under grant agreement No 955558.  
- Junta de Andalucía research project PID2020-113656RBC22 of MCIN/AEI/10.13039/501100011033.



This presentation contain content written by Enrique S. Quintana-Ortí

1

Some images are from: Flaticon.com

## Motivations

➤ General Matrix multiplication accelerations

High Performance

Intel Xeon, AMD HPC

High Commodity

NVIDIA GPU+ ARM

Low Power

Embedded Edge AI

Accelerators

SIMD Processor  
Xilinx Versal



2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023



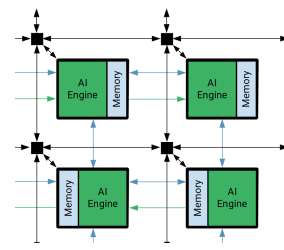
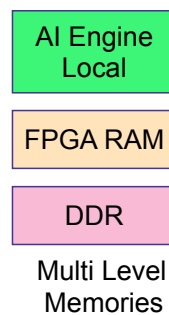
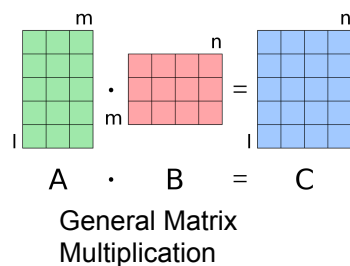
2

Previous work:

**Matrix Multiplication for Deep Learning Inference on the Xilinx Versal,**  
Jie Lei, Jose Flich, Enrique S. Quintana-Ortí,  
31st Euromicro International Conference on Parallel, Distributed, and Network-Based Processing  
03/2023



[arxiv.org/abs/2302.07594](https://arxiv.org/abs/2302.07594)



INT16 Based GEMM,  
28 MAC/Cycle out of 32

Mitigate Communication  
Overhead by Matrix reuse

3

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

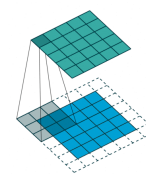


Current work:



GEMM-Like **Convolution** for Deep Learning **Inference** on the **Xilinx Versal**  
Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí

[https://jieggh.github.io/about/H32023/H3\\_PAGE.html](https://jieggh.github.io/about/H32023/H3_PAGE.html)



Extend the use case:  
applicable to direct convolution

16 BIT  
↓  
8 BIT

Compute throughput leap



Evaluating throughput  
with layers of Resnet

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

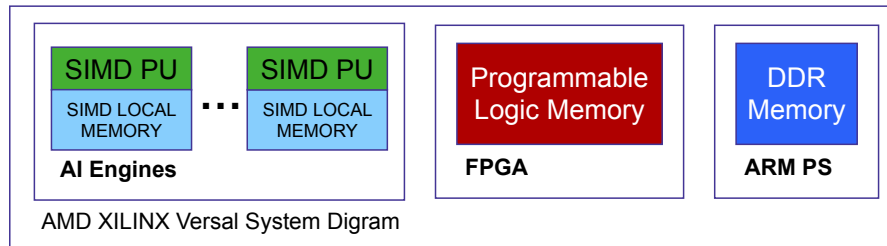
GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023



4

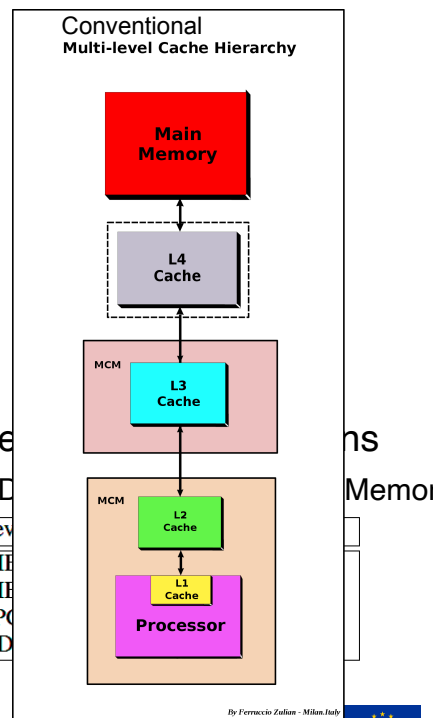


## Differences of our Versal based design compares



- No conventional multi-level caches system
  - Yet, 3 types with different speed and capacity
- Most of the data fetching/storing are made explicitly.
  - Streaming interface
  - Global memory access IO interface
- Matrix F repacked into FPGA memory

➤ Diver



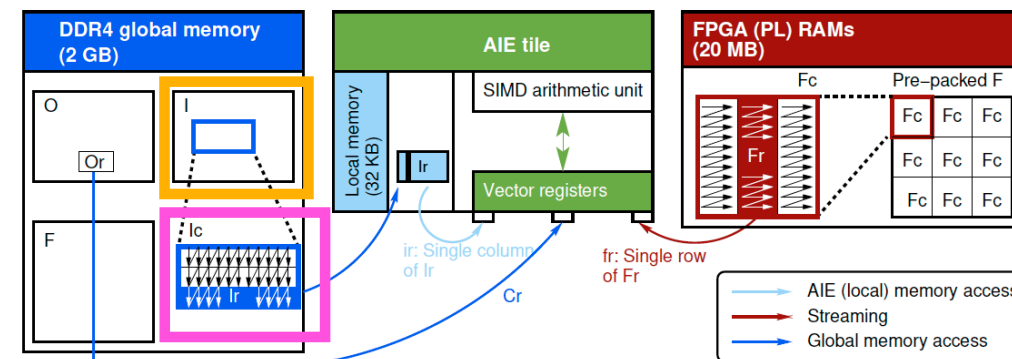
2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

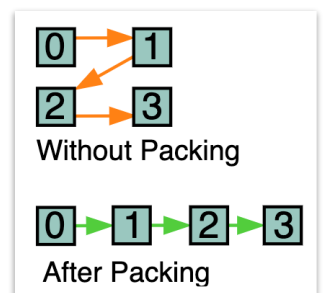
5



## Approach: Matrix tiling and packing



- Partition Matrix into smaller blocks
- Packing into a continuous memory location



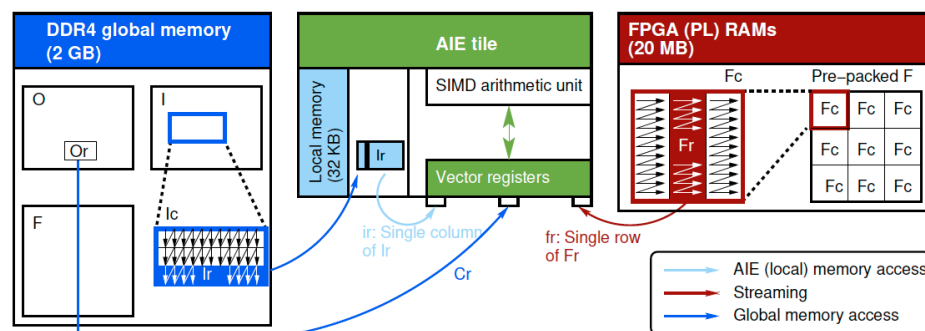
2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

6



## Approach: Tiled matrix transferring



- For an  $O = I * F$ , data movement
  - Output matrix O: From **DDR** -> **Vector register**
  - Input image matrix I: From **DDR** -> **SIMD Local memory** -> **Vector register**
  - Trained filter matrix F/Fc: From **FPGA Memory** -> **Vector register**

Capacities of different memory units in Versal

Level	Capacity
AIE vector registers	2 KB
AIE tile local memory	32 KB
FPGA RAMs	20 MB
DDR4 (global) memory	2 GB

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

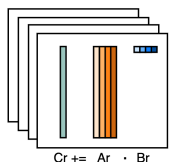
7



## Approach: Micro-kernel design

```
1 for (jc=0; jc<n; jc+=nc) // Loop L1
2 for (pc=0; pc<k; pc+=kc){ // L2
3   // Pack B
4   Bc:=B(pc:pc+kc-1, jc:jc+nc-1);
5   for (ic=0; ic<m; ic+=mc){ // L3
6     // Pack A
7     Ac:=A(ic:ic+mc-1, pc:pc+kc-1);
8     for (jr=0; jr<nc; jr+=nr) // L4
9       for (ir=0; ir<mc; ir+=mr) // L5
10        // Micro-kernel
11        C(ic+ir:ic+ir+mr-1,
12          jc+jr:jc+jr+nr-1)
13        += Ac(ir:ir+mr-1, 0:kc-1)
14        * Bc(0:kc-1, jr:jr+nr-1);
15  }}
```

- In platform code
  - Utilize intrinsic mac16( ... )
    - Performing 128 MAC (multiply and accumulate) operations per clock cycle.
  - Loop L5 unrolls as a factor of 16
    - Split compute into smaller sections, better compute and communication overlaps
    - Utilise more vector and accumulator registers, avoid register spilling



2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

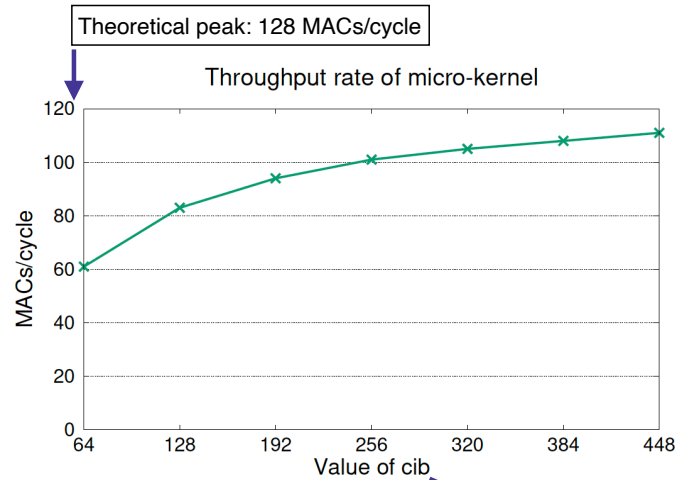
GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

8





## Identify bottleneck, micro-kernel at stationary mode



At stationary mode

Size of matrix problem within the micro-kernel

- In this loop:
  - Ir transfer is not considered
  - Fr transfer happen within the loop
- Throughput profiling at stationary mode
  - Only considering the computations at the micro-kernel level
  - The theoretical peak performance is 128 MACs/cycle
- A smaller workload may result in poor loop pipelining

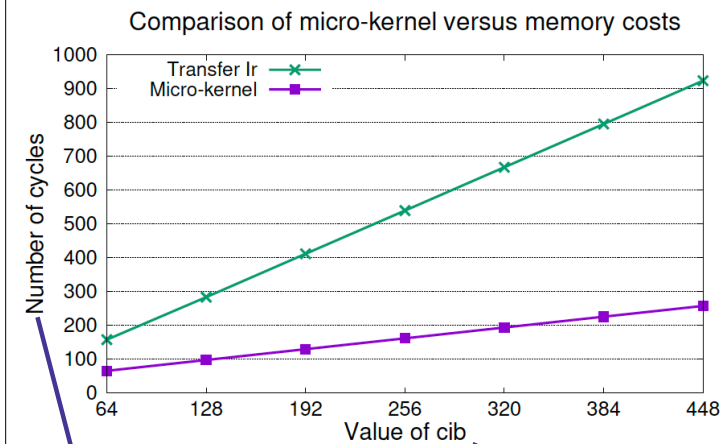
2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

9



## Identify bottleneck, micro-kernel considering some data transfer



Measuring the cycles costs

Size of matrix problem within the micro-kernel

- Transfer of Ir:
  - From DDR. -> SIMD Local memory.
  - > Vector register

Outside MAC Loop:

- Transferring of Ir from DDR into the local memory of the AI Engine

Facts:

- Slow
- Difficult to avoid

To amortize this cost::

- Reuse Ir more, key approach
- > Change the partition parameter of Ir

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

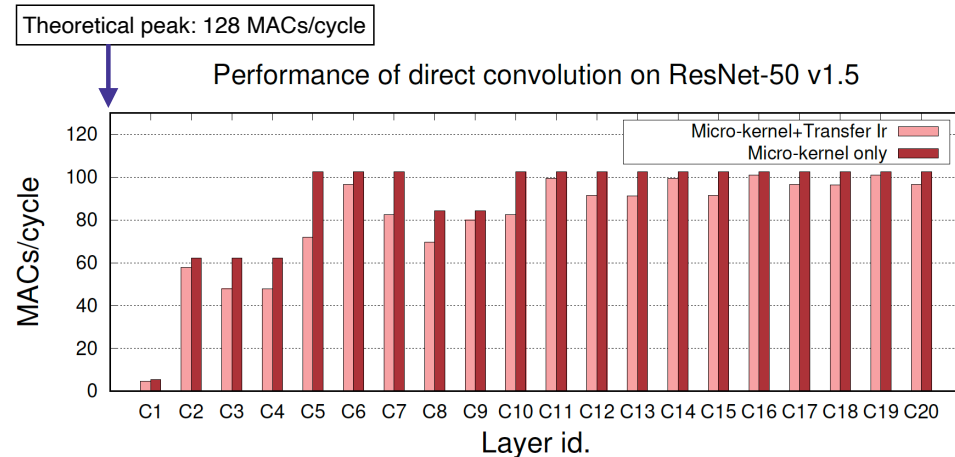
GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

10



- The theoretical peak performance is 128 MACs/cycle

## Evaluating computation throughput with layers in ResNet-50



> With larger layers: deliver optimal performance

> Low number of input channel (Ci) cause poor overall throughput

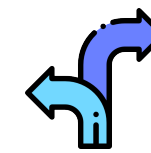
2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

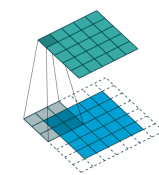
11



## Recap



Alternative GEMM approach for AMD Versal



Extend the use case: applicable to direct convolution

16 BIT  
↓  
8 BIT

Compute throughput leap



Tackle bottleneck: Throughput analyse



70% of peak performance for some layers of the ResNet

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

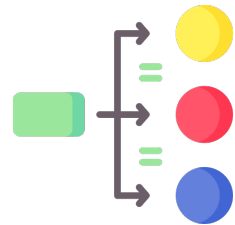
GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

12

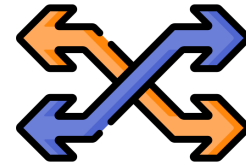


# APROPOS

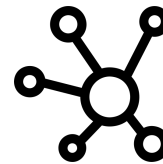
Future work



Utilise Multiple AI Engines



Mixed precisions support



Exploring heterogeneous architecture

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

13



# APROPOS

<http://apropos-itn.eu>

This project has received funding from the European Union's Horizon 2020 (H2020) Marie Skłodowska-Curie Innovative Training Networks H2020-MSCA-ITN-2020 call, under the Grant Agreement no 956090.



14

2023 WORKSHOP: HPC ON HETEROGENEOUS HARDWARE (H3)

GEMM-Like Convolution for Deep Learning Inference on the Xilinx Versal, Jie Lei, Héctor Martínez, José Flich, Enrique S. Quintana-Ortí 05/2023

# APROPOS

The authors gratefully acknowledge funding from

- European Union's Horizon2020 Research and Innovation program under the Marie Skłodowska Curie Grant Agreement No. 956090 (APROPOS, <http://www.apropos-itn.eu/>).
- European High-Performance Computing Joint Undertaking (JU) under grant agreement No 955558.
- Junta de Andalucía research project PID2020-113656RBC22 of MCIN/AEI/10.13039/501100011033.



15

