

Fast Convolution Operations on Many-Core Architectures

Reporter: Shigang Li
Affiliation: SKL of Computer Architecture Institute of Computing Technology, CAS
Email: shigangli.cs@gmail.com

17th International Conference on High Performance Computing and Communications (HPCC'15)





 In the context of image processing, convolution is the scalar product of the filter weights with the input pixels within a window surrounding each of the output pixels, such as 1D, 2D and multichannel 2D convolutions.



• They are widely used in image processing, such as Gaussian filter; today's Deep Learning, such as Convolution Neural Network.



- Convolution operations have high computational and memory access complexity.
- Many-core architectures evolve quickly, which feature high computational throughput and memory bandwidth, like NVIDIA and AMD GPUs, and Intel Xeon Phi.
- This paper aims at providing some insights for performance optimizations of convolution operations on these emerging many-core architectures.



1D and 2D convolutions





We discuss the optimizations in the context of OpenCL:

- A naïve implementation of 1D and 2D convolutions is each thread is responsible for one module and accesses input pixels from global memory directly.
- 2. A common optimization is utilizing local memory tiling to cache the input pixels, which reduces the total global memory accesses.

Register tiling for 1D and 2D convolutions



In order to exploit the register-level data reuse, register tiling by coarsening the workload of each thread is used.



Register tiling affects memory efficiency









Basically, the serial algorithm contains total 7 nested *for* loops.

The outermost 2-level *for* loops iterate on NumImages (i.e. batchsize) and NumFilters (number of filters), respectively.

The inner 5-level for loops, from outside to inside, iterate on image Height and Width, Channels, FilterSize and FilterSize, respectively

A naive method to implement this algorithm in OpenCL is to let each thread calculate for one pixel of the output images.





1. Each threads group is responsible for calculating NumFltrsPerThrd*

NumImgsPerThrd blocks of the output images, which exploits the data reuse of the outermost 2-level *for* loops, i.e. loop tiling.

2. Each block of input image is loaded into local memory before calculating, i.e. local memory tiling, which exploits inter-module data reuse.

3. Filters are put in constant memory to utilize constant cache.

Optimization details at thread level





For each channel, each thread streams in the pixels of NumImgsPerThrd images and the values of NumFltrsPerThrd filters to registers, and then execute the convolution on the current channel. Register tiling for Multi-channel 2D convolution



Register tiling can also be used to further exploit the inter-module date reuse. As we can see in the figure, the workload of each thread is coarsened.



Performance results of 1D convolution



All of the codes are implemented in OpenCL. We use AMD W8000 (3.23TFlops), GTX TITAN (4.5TFlops) and Intel MIC (2.0496TFlops to do the experimental analysis.



Hoversticeal 11D convolution

- 1. On AMD W8000 and GTX TITAN, the performance of the naive method and local memory tiling is almost equal. This is because that the naïve method may benefit from hardware cache.
- 2. Register tiling greatly improves the performance for both horizontal and vertical 1D convolutions on AMD W8000 and GTX TITAN
- 3. On Intel MIC, local memory tiling performs much worse, as no real hardware support for local memory.

2D convolution





- 1. Register tiling, with tiling size of 4*4, usually performs best for AMD W8000 and GTX TITAN.
- On Intel MIC, register tiling along the vertical dimension with tiling size of 4*1, which guarantees consecutive memory access, always performs best.



	image size	batch size	channels	filters	filter size
Scale ₀	128*128	128	3	96	13*13
Scale ₁	128*128	64	8	96	13*13
Scale ₂	128*128	64	16	96	13*13
Scale ₃	128*128	128	3	96	11*11
Scale ₄	128*128	64	8	96	11*11
Scale ₅	128*128	64	16	96	11*11
Scale ₆	128*128	128	96	128	9*9
Scale ₇	64*64	64	96	128	9*9
Scale ₈	32*32	64	96	128	9*9
Scale ₉	64*64	128	128	128	7*7
Scale ₁₀	32*32	64	128	128	7*7
Scale ₁₁	16*16	64	128	128	7*7
Scale ₁₂	32*32	128	128	384	3*3
Scale ₁₃	16*16	64	128	384	3*3
Scale ₁₄	13*13	64	128	384	3*3

Optimizations comparison of multi-channel 2D convolution





Optimizations comparison of convolution layer on AMD W8000

Using constant memory for filters improves the performance greatly on AMD W8000.

Optimizations comparison of multi-channel 2D convolution





Optimizations comparison of convolution layer on GTX TITAN

We can see that the best performance is obtained after local memory tiling. However, the constant memory defined in OpenCL doesn't work for NVIDIA GPUs.

Optimizations comparison of multi-channel 2D convolution





Optimizations comparison of convolution layer on Intel MIC

The best performance is obtained when the register tiling is conducted only along the vertical dimension, which guarantees consecutive memory access. However, local memory tiling and constant memory for filters don't work on Intel MIC.

Compared with cuDNN-v2





Performance comparison with cuDNN-v2 on GTX TITAN

- For the large filter size, our solution performs apparently better than cuDNNv2, with up to 33% performance improvement. Because our solution fully exploits inter-module date reuse.
- 2. However, our solution performs worse than cuDNN-v2 when the filter size is small, which is caused by the low inter-module data reuse rate of small filter size.

Compared with clBLAS+unfold





Compared with clBLAS+unfold on AMD W8000

clBLAS+unfold follows the same method with the well-known deep learning library—CAFFE.

- 1. Our solution performs apparently better than clBLAS+unfold for large filter size, with up to 28% performance improvement.
- 2. However, our solution performs worse when the filter size is small.





- For 1D and 2D convolutions, register tiling rather than local memory tiling is critical to achieve good performance. And different register tiling modes affect the memory access efficiency.
- Our solution for multi-channel 2D convolution fully exploits inter-module data reuse, and thus gets good performance for large filter sizes, up to 33% over cuDNN-v2 and up to 28% over clBLAS.
- To improve the performance of multi-channel 2D convolution for small filter size is our future work.





Thank You

Reporter: Shigang Li Email: shigangli.cs@gmail.com