

Optimization of the DNN program on the CPU+MIC Platform

University of Electronic Science and Technology of China

1. Preface

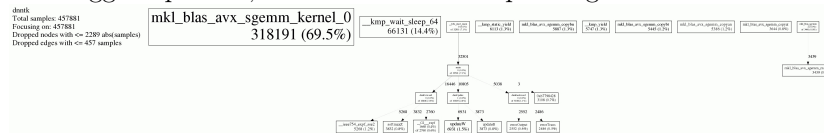
In the section we are required to optimize a DNN(deep neural network) program based on a standalone hybrid CPU+MIC platform. The detailed configuration is as follows:

Item	Name	Configuration	Hosts
Server	Inspur NF5280M4 x 4	CPU : Intel Xeon E5-2680v3 x 2, 2.5Ghz, 12 cores	hostname: mic1, mic2, mic3, mic4
		Memory: 16G x8, DDR4, 2133Mhz	
		Hard disk: 1T SATA x 1	
		Accelerator card: Intel XEON PHI-31S1P (57 cores, 1.1GHz, 1003GFlops, 8GB GDDR5 Memory)	
Network		Infiniband+Ethernet	

Classification	Description	Installation path	Version
OS	GNU/Linux		RHEL 7.1
Compiler	Intel Composer XE Suites	/opt/intel/composer_xe_2015.0.090	2015.0.090
MKL	Intel MKL	/opt/intel/mkl/lib/intel64	
MPI	Intel MPI	/opt/intel/impi/5.0.1.035	5.0.1.035
PBS	Torque	/opt/tsce	3.0.5

2. Analysis of the serial program

First, we generate a call graph by using **Google perftools**, a open source performance profiler, to have a glance though it. Every square represents a function, and the bigger square is, the more time corresponding function cost.



Obviously, the hot spot is something about **MKL**. After googling and searching Intel document we know that MKL provides BLAS routines, which includes

a serial function named “cblas__?gemm” to compute a matrix-matrix product with general matrices.