

# README

## Im2col

### Simple Usage

See `main.cpp`.

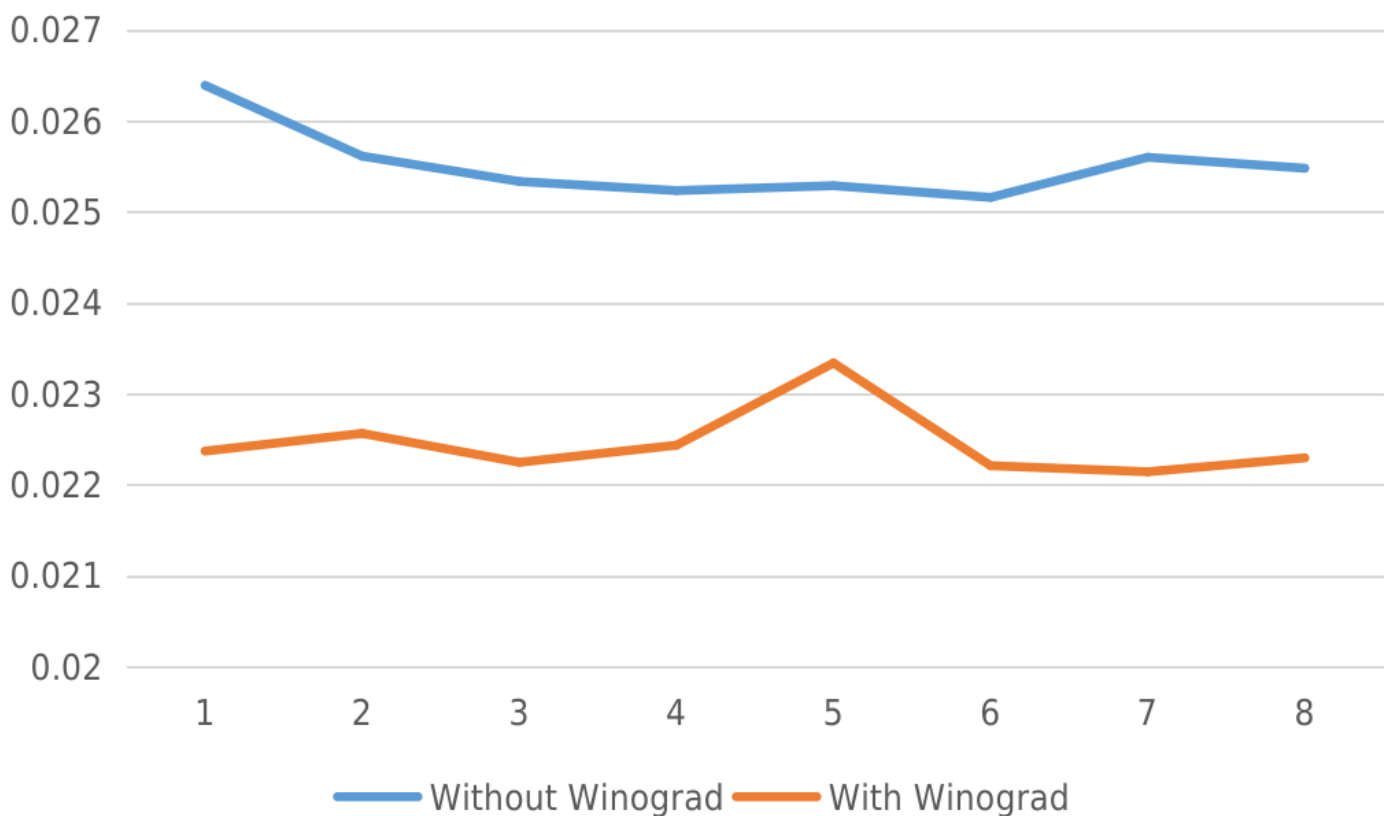
### Files and Directories

```
1 Im2col
2 |— acc_function
3 |   |— AcceleratorFunction.h      // The base class of accelertor functions.
4 |   |— winograd
5 |       |— WinogradFunction_1D.cpp // Winograd 1D. (***** HERE!!!!!! *****)
6 |       |— WinogradFunction_1D.h
7 |       |— WinogradFunction.h     // The base classWinograd accelertor funct
8 |— CMakeLists.txt
9 |— feature_map
10 |   |— DirectFeatureMap.cpp
11 |   |— DirectFeatureMap.h
12 |   |— FeatureMap.h
13 |   |— Im2colFeatureMap.cpp      // Im2col feature map that allows
14 |   |                           convolution with a accelerator
15 |   |                           function. (***** HERE!!!!!! *****)
16 |   |— Im2colFeatureMap.h
17 |   |— Im2colFeatureMap_OMP.cpp
18 |   |— Im2colFeatureMap_OMP.h
19 |— kernel
20 |   |— DirectKernel.cpp
21 |   |— DirectKernel.h
22 |   |— Im2colKernel.cpp
23 |   |— Im2colKernel.h
24 |   |— Kernel.h
25 |— main.cpp
26 |— output_map
27 |   |— OutputMap.cpp
28 |   |— OutputMap.h
29 |— util
30 |   |— GetTime.h
```

## Analysis

```
Direct conv: 0.03022
Im2col conv: 0.0263898
Im2col conv with Winograd: 0.0223682
Im2col conv with Winograd (OpenMP enabled): 0.00318503
```

### Steady Improvement is Witnessed when Using Winograd



In Lab2 settings, my implementation is around 13% faster on average when using Winograd(2, 3), which is not beyond expectation because, in fact, Winograd does not break much spatial locality. As the size of a row is small under our settings, although it loads 2 rows at the same time, we can infer that it still need not reload data from memory to cache.

## Parallelism

Still we can use OpenMP to parallel convolution even if Winograd is used, yet this time we cannot simply add the pragma in the innermost for loop because the innermost for loop will not loop for many times when we use Winograd (e.g., loop for  $K * R * S / 3$  times when using Winograd(2, 3)). Here we'd better let one thread responsible for every two rows, which guarantees little parallel overhead (and also better locality actually).

thread 1  


$d_{00}$	$d_{01}$	$d_{02}$	$d_{10}$	$d_{11}$	$d_{12}$	$d_{20}$	$d_{21}$	$d_{22}$
$d_{01}$	$d_{02}$	$d_{03}$	$d_{11}$	$d_{12}$	$d_{13}$	$d_{21}$	$d_{22}$	$d_{23}$
$d_{10}$	$d_{11}$	$d_{12}$	$d_{20}$	$d_{21}$	$d_{22}$	$d_{30}$	$d_{31}$	$d_{32}$
$d_{11}$	$d_{12}$	$d_{13}$	$d_{21}$	$d_{22}$	$d_{23}$	$d_{31}$	$d_{32}$	$d_{33}$

 thread 2

$$\begin{bmatrix} k_{00} \\ k_{01} \\ k_{02} \\ k_{10} \\ k_{11} \\ k_{12} \\ k_{20} \\ k_{21} \\ k_{22} \end{bmatrix} = \begin{bmatrix} r_{00} \\ r_{01} \\ r_{10} \\ r_{11} \end{bmatrix}$$