

Part A

Question 1

Posix_memalign is a function in the POSIX standard library that is used to align the memory to ensure that data is stored in a way that corresponds to the cache's line size. If the data is not aligned, it may span multiple cache lines, resulting in more cache misses and slower performance.

Question 2

In naive matrix multiplication, its for loop for row cannot well utilize cache. When the dimension of the matrix is large, cache cannot store all elements in one row. Then, it keeps generating cache misses when computing the inner product between row and column. This is because after it iterates to the end of a row, the values previously fetched has been removed.

In addition, for each row, it computes the inner product with each columns. However, after computing inner product of row i with all columns, it needs to fetch all column values from the main memory again when computing the inner product with row $i+1$.

Starting from matrix with size 256, the results show many noticeable differences. And this difference increases when size of matrix gets larger.

Question 3

When tile size is small, the megaflops are small and increases as tile size increases. When tile size is large, the megaflops are also small and decreases as tile size increases. The optimal tile size is 16. Details could see the following results:

Tile	Matrix	MFLOPS
=====		
Naive	128	497.6988826000000000000000
Naive	256	387.7734726000000000000000
Naive	512	312.2995818000000000000000
Naive	1024	331.3849300000000000000000
Naive	2048	99.9356756000000000000000
=====		
4	128	336.5923542000000000000000
4	256	323.6110706000000000000000
4	512	321.1134044000000000000000
4	1024	322.9726508000000000000000
4	2048	235.8644074000000000000000
=====		
8	128	362.8427200000000000000000
8	256	373.3907060000000000000000
8	512	373.4118380000000000000000
8	1024	372.0156818000000000000000
8	2048	350.1731590000000000000000
=====		
16	128	756.8243124000000000000000
16	256	764.8316772000000000000000
16	512	656.7675540000000000000000
16	1024	658.9617510000000000000000
16	2048	627.1330214000000000000000
=====		
32	128	868.6356268000000000000000
32	256	638.3576886000000000000000
32	512	540.7361550000000000000000
32	1024	541.2610866000000000000000
32	2048	434.8541128000000000000000
=====		
64	128	605.7598694000000000000000
64	256	576.3511986000000000000000
64	512	483.5854078000000000000000
64	1024	415.2916810000000000000000
64	2048	398.2978626000000000000000

Question 4

Yes it somehow changes the performance. But the differences are not significantly large.

Metric: MegaFlops/s.

Row-Column Order

16*16: 592.40

32*32: 268.43

512*512: 264.87

1024*1024: 316.66

2048*2048: 93.95

Column-Row Order

16*16: 241.97

32*32: 289.95

512*512: 284.29

1024*1024: 309.01

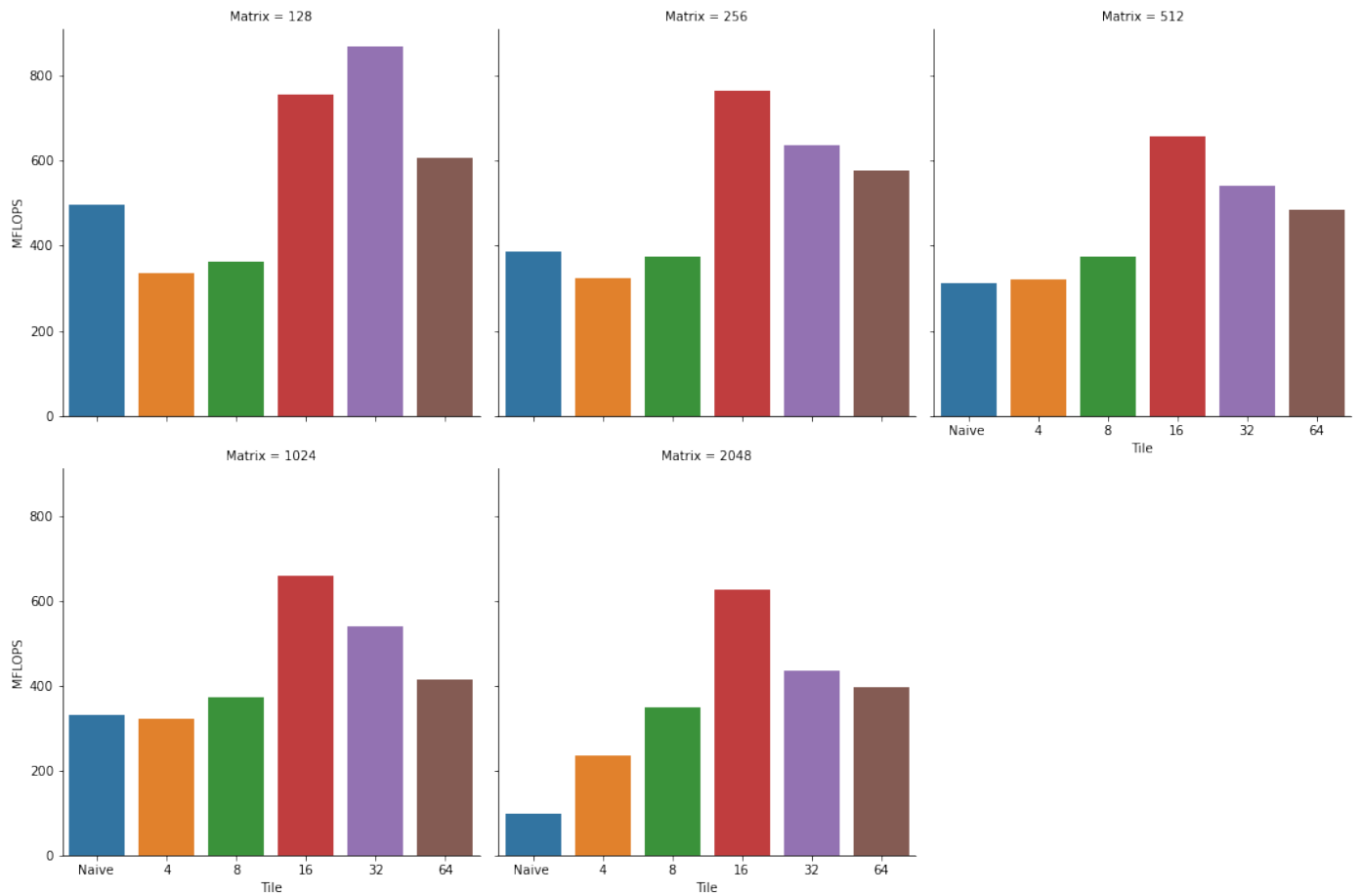
2048*2048: 86.34

Except when matrix is 16*16, row column-order is significantly faster than column-row order, in other small matrix multiplication, column-row order is slightly faster than row-column order.

However, as dimension of matrix is large, row column-order is consistently faster than column-row order.

I think this might due to spatial locality. When matrix dimension is small a cache line could store entire row and thus, fix column at outer loop is more efficient than fix row at outer loop. However, when dimension of matrix is large this advantage diminishes, because it cannot store the entire row in a cache block.

Part B



The above subplots are generated by `seaborn.catplot` function, which compare the MFLOP for different variants and tile number with the matrix size.