

1.
(a)

we could draw the graph:

Inst \ cycle	1	2	3	4	5	6	...	500	501	502	503	504	505
1	/	2	3	4	5	6	...						
2		/	2	3	4	5	...						
3			/	2	3	4	...						
4				/	2	3	...						
5					/	2	...						
6						/	...						
:							...						
498							...	3	4	5	6		
499							...	2	3	4	5	6	
500							...	1	2	3	4	5	6

From the graph, we see we need 505 cycles to finish (500 cycles for the instruction #500 to start, and additional 5 cycles for it to finish)

(b) Total time = Cycles / Frequency
= $505 / 2 \text{ GHz}$
 $= 2.525 \cdot 10^{-7} \text{ s} = 252.5 \text{ ns}$

2.

(a) the line $A[i] += B[i] + B[i+1]$

is actually $A[i] = A[i] + B[i] + B[i+1]$

So, for this single line, we need to have

$A[i]$, $B[i]$, $B[i+1]$ in our instruction.

We could draw a table for the memory we need for each instruction: (Read directly from memory is marked purple, read from register is marked green)

inst #

0	$A[0]$	$B[0]$	$B[1]$
1	$A[1]$	$B[1]$	$B[2]$
2	$A[2]$	$B[2]$	$B[3]$
:	:	:	:
255	$A[255]$	$B[255]$	$B[256]$

From the table, we could see that :

from memory : $256 \times 2 + 1 = 513$

from register: 255

Time it takes to read from memory :

$$513 \times 100 \text{ ns} = 51300 \text{ ns} = 51.3 \text{ ms}$$

(b) The hit ratio is:

$$\frac{255}{255+513} = \frac{255}{768} \approx 0.332 = 33.2\%$$

3. Again, we could draw a table

purple = read from memory

blue = read from cache

green = read from register

inst #	register			cache
0	A[0]	B[0]	B[1]	Cache A[0:3] Cache B[0:3]
1	A[1]	B[1]	B[2]	
2	A[2]	B[2]	B[3]	
3	A[3]	B[3]	B[4]	Cache B[4:7]
4	A[4]	B[4]	B[5]	
:	:	:	:	(this means the memory from base_ptr + 256 to base_ptr + 259, inclusive. B itself might not be that long)
255	A[255]	B[255]	B[256]	Cache B[256:259]

From the table, we could see that;

$$\# \text{ from memory} : 64 \text{ (from A)} + 65 \text{ (from B)} = 129$$

$$\# \text{ from cache} : 192 \text{ (from A)} + 192 \text{ (from B)} = 384$$

$$\# \text{ from register} : 255$$

Thus, time takes:

$$t = 129 \times 100 \text{ ns} + 384 \text{ ns} \times 2 \text{ ns} = 13668 \text{ ns}$$

$$= 13.668 \text{ ms}$$

4. We could draw a table

block \ index	0	1	2	3	4	5	6	7
0	64 (f)	65	66	67	68	69	70	71 (loaded (f))
1								
2	16 (b)	17 (b)	18	19 (c)	20 (a)	21	22	23 (loaded (a))
3								
4								
5	104	105	106	107	108 (e)	109	110	111 (loaded (e))
6								
7	248	249	250 (d)	251	252	253	254	255 (loaded (d))

(a) line : $\lfloor 20/8 \rfloor = 2$

block : $2 \% 8 = 2$

since there is no line currently loaded. This is a cold miss. Load cache line #2 to block 2.

(b) line : $\lfloor 17/8 \rfloor = 2$

block : $2 \% 8 = 2$

The line #2 is already in cache (block 2).

this is a hit, read from the block 2

(c) line: $\lfloor 19/8 \rfloor = 2$

block: $2 \% 8 = 2$

The line #2 is already in cache (block 2).

this is a hit, read from the block 2

(d) line: $\lfloor 250/8 \rfloor = 31$

block: $31 \% 8 = 7$

Line #31 currently is not loaded in cache. This is a cold miss. Load cache line #31 to block 7.

(e) line: $\lfloor 108/8 \rfloor = 13$

block: $13 \% 8 = 5$

Line #13 currently is not loaded in cache. This is a cold miss. Load cache line #13 to block 5.

(f) line: $\lfloor 64/8 \rfloor = 8$

block: $8 \% 8 = 0$

Line #8 currently is not loaded in cache. This is a cold miss. Load cache line #8 to block 0.

5.

(1) **False.** It's not always useful. Consider a cache line of size 4. If CPU asks data on address 0, 4, 8, 12... It still misses everytime, and need to go to the memory. In this case, the caching hardware need to do extra works (load cache) that has no benefits.

(2) **True.** If the program has spatial/temporal locality. Prefetching could reduce compulsory misses. Yet, it's not guaranteed to reduce such misses if the program's memory access doesn't have certain patterns (say the program randomly accesses the memory, or program traverses the matrix in column major).

(3) **True**

(4) **False.** In the Q1, we see it's not $N \cdot M$ steps for N independent operation. The steps (or cycles) is actually $N + (M - 1)$.

(5) **False.** The pipeline is designed to speed up multiple instructions. For single instruction,

It still requires same cycles. Say for an instruction that takes m cycles to execute, if you split the instruction into M substeps each of which takes 1 cycle to finish, you still need $1 + (M-1) = m$ cycles to execute, and in real situations, implementing the pipeline might make each substep slower. So, pipeline could not make single instruction faster. It speeds up multiple instruction, though.