

Given,	<u>instructions</u>	<u>clock cycles</u>
	ALU operations	1.0
	Loads	3.5
	Stores	2.8
	Branches	
	Taken.	4.0
	not taken	2.0
	Jumps	2.4

$\frac{9+9}{4+2+2} = \frac{18}{8} = 2.25$

$\frac{9+9}{4+2} = \frac{18}{6} = 3.0$

need to find effective CPI for RISC-V
 → obtain average instruction frequencies by considering. Perlbench & sjeng.

<u>instruction</u>	<u>Perlbench</u>	<u>sjeng</u>	<u>Avg. freq.</u>
loads	25%	19%	22%
stores	14%	7%	10.5%
Branches	15%	15%	15%
Jumps	7%	3%	5%
ALU	39%	56%	47.5%

for eff. CPI.

$$\text{CPI} = (22 * 3.5 + 10.5 * 2.8 + 15 * 1.0 + 5 * 2.4 + 47.5 * 1.0) * \frac{1}{100}$$

$$= (0.77 + 0.294 + 0.45 + 0.12 + 0.475)$$

$$= 2.109_{11}$$

\Rightarrow effective CPT for RISCV mentioned

is 2.109_{11}

2. Given,

- a) So we need to see the execution time in case of without enhancement. \therefore of E.T. also, the computer worked. % of E.T. enhanced mode \rightarrow 50%. non-enhanced mode \rightarrow 50%.

Without enhancement, non-enhanced mode takes 50% but enhanced takes 10 times longer i.e., 500%.

So, this means without enhancement, relative Execution time = 50% + 500% = 550%.

$$\text{Speed UP} = \frac{E.T(\text{non-enh})}{E.T(\text{enh})}$$

$$= 550 / 100 = 5.5$$

- b. % of original E.T has been converted to fast mode.

By Amdahl's law,

(Speed UP \rightarrow S.U)

$$\text{fraction vector} = \frac{S.U_{\text{overall}} * S.U_{\text{fast}} - S.U_{\text{non-fast}}}{S.U_{\text{overall}} * S.U_{\text{fast}} - S.U_{\text{overall}}}.$$

$$\Rightarrow \text{Fraction Vected} = \frac{5.5 * 10 - 10}{5.5 * 10 - 5.5}$$

$$= 55 - 10$$

$$= 45.5 - 5.5$$

$$= \underline{45}$$

$$= 0.9090$$

$$\Rightarrow 90.90\%$$

4. Given,

Consider basic RISC-V, 5-stage Pipeline.
a) Full forwarding.

(i) How pipelining can improve the performance of a given instruction mix.

→ Pipelining is a process used to improve performance of a given instruction mix, by providing parallelism to the set of instructions executed. Although they execute in sequence, certain parts of execution overlaps and improves performance via higher throughput.

v)

	1	2	3	4	5	6	7	8	9	10	11
lw x10, 0(x11)	F	D	E	M	w						
add x9, x11, x11	F	D	E	M	w						
sub x8, x10, x9	F	D	E	M	w						
lw x7, 0(x9)	F	D	E	M	w						
sw x9, 4(x9)	F	D	E	M	w						

(iii) where SW get data from? Here data is present in 'lw' in the MEM1 WB register. Here we can also have a feedback path from this register to one of the IIP's to the data memory.

Control signal could select the appropriate IIP to data memory to support this. If - SW IIP combination

B.

	1	2	3	4	5	6	7	8	9	10	11	12	13
beg X1,X2,X	F	D	E										
lw X10,0(X11)		F	D	E									
sub X14,X10,X10			F	D	E	M	W						
X1: add X4,X1,X2				F	D	E	M	W					
lw X1,0(X4)					F	D	E	M	W				
sub X1,X1,X1						F	D	E	M	W			
add X1,X1,X1							F	D	E	M	W		

→ add dependency Sub instruction before it
∴ Produces the data which was consumed by add instruction.

5. For cache organized in two ways mentioned,

(i) more data elements per block & fewer blocks,

Cons:

→ There is higher potential for conflict misses as there will be fewer blocks.

Pros:

→ Higher performance can be achieved due to spatial locality due to larger block size.

Example:

In a case where we have high degree of sequential data access.

(ii) Fewer elements per block but more blocks.

Pros:-

→ Here in this case we may see fewer conflicts misses due to more unique mappings unlike other case.

Cons:-

→ we may be more subjected to compulsory misses due to smaller block-size.

Example:

→ In a case where we have more random memory access.

(ii) Given,

Processor has direct mapped cache.

Data words → 8 bits

Physical address → 20 bits.

tag → 11 bits.

As Physical address is 20 bits out of which 11 bits are tag bits,
so 9 bits are left for index & offset.

→ Each block holds 16 bytes of data,
16 words & so, 4 bits are needed
for offset.

\Rightarrow index has 5 bits (9-4).

$$\begin{aligned} \text{cache} &= 2^n n \\ &\rightarrow \text{index} \\ &= 2^5 \\ &= 32 \text{ blocks in the cache.} \end{aligned}$$

(ii) Given,

- 16-way set-associative cache.
- \rightarrow Data words are 64-bit long.
- \rightarrow words are addressed to the half-word.
- \rightarrow Cache holds 2 MB bytes of data
- \rightarrow Each block holds 16 data words
- \rightarrow Physical address are 64-bit long.

word length = 64 bits = 8 bytes.

$$\begin{aligned} \text{cache size in words} - 2 \text{ MB} / 8 &= 256 \text{ kB} \rightarrow 18 \text{ bits} \\ \text{Block size} &= 16 \text{ words} \rightarrow 4 \text{ bits} \\ \text{no of lines} &= CS/BS = 2^{2^n(18-4)} \\ &= 2^{14} \rightarrow 14 \text{ bits.} \end{aligned}$$

but, as we have 16 way set-associative.
set bits = 4 bits.
block bits = 4 bits.

Data word contains half of data word.

$$\Rightarrow 64/2 = 32 = 2^5.$$

\Rightarrow offset \Rightarrow 5 bit.

Total capacity of cache \rightarrow 2 MB \Rightarrow 2²¹ bytes

\Rightarrow total no. of blocks in cache.

$$= 2^{21} \text{ bytes} (1 \text{ block}/16 \text{ words}) * (1 \text{ word}/64 \text{ bits}) \\ * (8 \text{ bits}/1 \text{ byte})$$

$$= 2^{14} \text{ blocks.}$$

Now, we have 2¹⁴ blocks / set.

So there are 2¹⁴ blocks (1 set / 2¹⁴ blocks)

$$= 2^{10} \text{ sets}$$

\Rightarrow index \rightarrow 10 bits

Now, as we know offset & index we can get tag bits.

$$64 - 10 - 5 = 49$$

\Rightarrow tag \rightarrow 49 bits

index \rightarrow 10 bits

offset \rightarrow 5 bits //

3. Given,

we run 2 applications on our dual core.

→ 1st application needs 80% of resources & other needs 20%. So assume when we parallelize a portion of program speed up for that portion is 2.

We know that,

$$\text{Speed Up} = \frac{1}{\text{Time}}$$

$$(1 - \text{Fraction of enhanced program}) + \frac{\text{Fraction of enhanced program}}{\text{Speed Up of enhanced portion}}$$

(i) Given,

40% of 1st application is parallelizable.

$$\text{Speed Up} = \frac{1}{(1 - 0.4) + \frac{0.4}{2}}$$

$$= \frac{1}{0.6 + 0.2} = \frac{1}{0.8} = 1.25 //$$

(ii) Given,

99% of 2nd application is parallelized

$$\text{Speed UP} = \frac{1}{(1-0.99) + \frac{0.99}{2}}$$

$$= \frac{1}{(0.01) + 0.495}$$

$$= \frac{1}{0.505}$$

$$= 1.98 //$$

(iii) Given,

40% of P⁺ application is parallelizable,
overall Speed UP = ?

$$\text{Overall Speed UP} = \frac{1}{1}$$

$$= \frac{1}{(1-0.8) + \frac{0.8}{1.25}}$$

$$= \frac{1}{0.2 + 0.64}$$

$$\text{Overall Speed UP} = 1.19 //$$