

① $s = 10$

speedup = 1.3

$$\text{speedup} = \frac{ET_{\text{orig}}}{ET_{\text{new}}} = \frac{1}{(1-f) + f/s}$$

$$1.3 = \frac{1}{(1-f) + f/10}$$

$$1-f + f/10 = 1 - 9f/10$$

$$1.3 = \frac{1}{1 - 9f/10} \Rightarrow 1.3(1 - 9f/10) = 1$$

$$1.3 - 1.17f = 1$$

$$f = 0.26 = \boxed{26\%}$$

$$\text{② speedup} = \frac{1}{(1 - 0.26) + \underbrace{0.26}_{\infty}/\infty} = \frac{1}{(1 - 0.26)} = \boxed{1.35}$$

$$\text{③ average cpi} = 10(0.03) + 5(0.30) + 5(0.20) + 1(0.47) = \boxed{3.27}$$

④ data hazards: arise because of the nature of computation, that is when you need to use a value computed in one instruction in a later instruction (when does that value become available)

structural hazards: when there is contention for a shared hardware resource

control hazards: arise when a jump occurs, either conditional or unconditional, and you have to calculate whether to take the jump and where the target address is

⑤ DCA: 3000

~~00000000~~ ~~00000000~~

011 000 001 000

⑥ Superscalar processors generate conventional sequential machine code in the compiler, then at runtime, hardware figures out dependencies and issues instructions accordingly. The dependencies it works out of the code are generally anti and output dependencies. It does this by static single assignment ~~or register tracking~~ (trying not to reuse register names). What's left are ~~the~~ usually just data dependencies. This makes room for more instructions to run at the same time. Instructions are then fed into reservation stations according to instruction type. These reservation stations maintain information (including dependencies) and let instructions execute when their ~~dependencies~~ operands become available. This ~~keeps~~ ~~respects~~ method respects dependencies and allows instructions to run parallel as long as these dependencies are respected.

⑦ No address or immediate : 1 dx-word

Address : 2 dx words

Immediate : 2 dx words

⑧ instruction 2 \rightarrow instruction 3 (RAW with R4)

instruction 1 \rightarrow instruction 3 (WAR with R2)

instruction 1 \rightarrow instruction 3 (RAW with R1)

instruction 3 \rightarrow instruction 4 (RAW with R2)

instruction 1 \rightarrow instruction 4 (RAW with R1)

instruction 3 \rightarrow instruction 5 (WAR with R4)

instruction 1 \rightarrow instruction 5 (RAW with R1)

instruction 3 \rightarrow instruction 5 (RAW with R2)

instruction 2 \rightarrow instruction 5 (WAW with R4)

⑨ False

⑩ 36 bit phys addr, cache BLOCK 2k bytes, cache is 1MB, dmr
 \downarrow 2^{11} bytes \downarrow 2^{20} bytes

$$\frac{2^{20} \text{ bytes}}{\text{cache}} \div \frac{2^{11} \text{ bytes}}{\text{block}} = 2^9 = 512 \text{ blocks} \quad \downarrow \quad \text{q-bit index}$$

$$\text{⑪ } \frac{2k \text{ bytes}}{\text{block}} - 2^{11} \text{ bytes} = \boxed{11\text{-bit offset}}$$

$$\text{⑫ } 36 - 9 - 11 = \boxed{16\text{-bit tag}}$$

$$\text{⑬ avg mem access time} = \text{hit time} + \text{miss rate} \times \text{miss penalty}$$

" instruction	1	+	0.01	x	10	=	1.1
" data	1	+	0.05	x	10	=	1.5

$$\frac{1}{1.30} \cdot 1.1 + \frac{0.30}{1.30} \cdot 1.5 = \boxed{1.19}$$

⑭ 2-bit branch predictor: with a 1-bit predictor the nested loops will continuously make two wrong predictions: ~~when the loop is meant to stop and keep it~~ when the loop is on its last iteration it will predict to take the loop but will be wrong and when the loop has to run again it will predict not to take it but be wrong. With a 2-bit, the nested loops will only make one wrong prediction: when they are on ~~the~~ their current last iteration. But since it requires ~~to~~ two wrong predictions to change, when the loop has to run again, it will still choose to take the loop and be ~~was~~ right