

Lecture No 20

Comparison of Memory Models



Comparison of Memory Models

- The δ binomial model is suitable for simple pipelined processors where n requests per T_c are each made with probability δ .

Review and Selection of Queuing Models

- There are basically three dimensions to simple (single) server queuing models.
- These three represent the statistical characterization of arrival Rate, Service rate and amount of buffering present before system saturates.
- For arrival rate, if the source always requests service during a service interval, Use M_B or simple binomial model.

Review and Selection of Queuing Models

- If the particular requestor has diminishingly small probability of making a request during a particular service interval, use poisson arrival.
- For service rate if service time is fixed , use constant (D) service distribution.
- If service time varies but variance is unknown, (choose $c^2=1$ for ease of analysis) use exponential (M) service distribution.

Review and Selection of Queuing Models

- If variance is known and C^2 can be calculated use M/G/1 model.
- The third parameter determining the simple queuing model is amount of buffering available to the requestor to hold pending requests.

Processors with Cache

- The addition of a cache to a memory system complicates the performance evaluation and design.
- For CBWA caches, the requests to memory consists of *line read* and *line write* requests.
- For WTNWA caches, its *line read* requests and *word write* requests.
- In order to develop models of memory systems with caches two basic parameters must be evaluated

Processors with Cache

1. $T_{\text{line access}}$, time it takes to access a line in memory.
2. T_{busy} , potential contention time (when memory is busy and processor/cache is able to make requests to memory)

Accessing a Line $T_{\text{line access}}$

- Consider a pipelined single processor system using interleaving to support fast line access.
- Assume cache has line size of L physical words(bus word size) and memory uses low order interleaving of degree m .
- Now if $m \geq L$, the total time to move a line (for both read and write operations)

$$T_{\text{line access}} = T_a + (L-1) T_{\text{bus.}}$$

Where T_a is word access time & T_{bus} is bus cycle time.

Accessing a Line $T_{\text{line access}}$

- If $L > m$, a module has to be accessed more than once so module cycle time T_c plays a role.
- If $T_c \leq (m \cdot T_{\text{bus}})$, module first used will recover before it is to be used again so even for $L > m$

$$T_{\text{line access}} = T_a + (L-1)T_{\text{bus}}$$

- But for $L > m$ and $T_c \geq (m \cdot T_{\text{bus}})$, memory cycle time dominates the bus transfer

Accessing a Line $T_{\text{line access}}$

- The line access time now depends on relationship between T_a and T_c and we can now use.

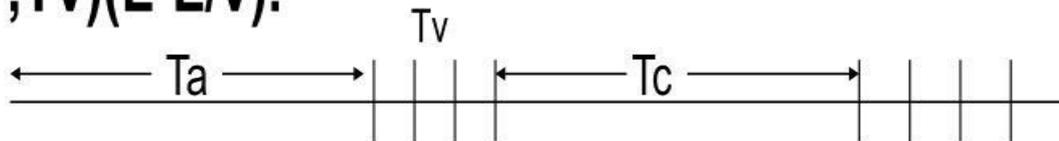
$$T_{\text{line access}} = T_a + T_c \cdot ((L/m) - 1) + T_{\text{bus}} \cdot ((L - 1) \bmod m).$$

- The first word in the line is available in T_a , but module is not available again until T_c . A total of L/m accesses must be made to first module with first access being accounted for in T_a . So additional $(L/m - 1)$ cycles are required.

Accessing a Line $T_{\text{line access}}$

- Finally $((L-1) \bmod m)$ bus cycles are required for other modules to complete the line transfer.
- If we have single module memory system ($m=1$), with nibble mode or FPM enabled module. Assume v is the no of fast sequential accesses and T_v is the time between each access

$$T_{\text{line access}} = T_a + T_c \left(\left(\frac{L}{v} \right) - 1 \right) + \left(\max(T_{\text{bus}}, T_v) \right) (L - L/v).$$



Accessing a Line $T_{\text{line access}}$

- Now consider a mixed case ie $m > 1$ and nibble mode or FPM mode.

$$T_{\text{line access}} = T_a + T_c((L/m.v) - 1) +$$

$$T_{\text{bus}}(L - (L/m.v))$$

Computing $T_{\text{line access}}$

- Case 1: $T_a = 300\text{ns}$, $T_c = 200\text{ns}$, $m = 2$,
 $T_{\text{bus}} = 50\text{ ns}$ and $L = 8$.

Here we have $L > m$ and $T_c > m \cdot T_{\text{bus}}$

So $T_{\text{line access}} = T_a + T_c \left(\frac{L}{m} - 1 \right) + T_{\text{bus}} \left((L-1) \bmod m \right)$.

$$= 300 + 200(4-1) + 50(1) = 950\text{ns}$$

Computing $T_{\text{line access}}$

- Case 2: $T_a=200\text{ns}$, $T_c=150\text{ns}$, $T_v=40\text{ns}$, $T_{\text{bus}}=50\text{ ns}$, $L=8$, $v=4$, $m=1$.

$$T_{\text{line access}} = T_a + T_c((L/v)-1) + \max(T_{\text{bus}}, T_v)(L-L/v).$$

$$= 200 + 150((8/4) - 1) + 50(8 - (8/4))$$

$$= 200 + 150 + 300$$

$$= 650 \text{ ns}$$

Computing $T_{\text{line access}}$

- Case 3: $T_a=200\text{ns}$, $T_c=150\text{ns}$, $T_v=50\text{ns}$, $T_{\text{bus}}=25\text{ ns}$, $L=16$, $v=4$, $m=2$.

$$T_{\text{line access}} = T_a + T_c((L/m.v)-1) + (T_{\text{bus}})(L-L/m.v).$$

$$=200+ 150((16/2.4)-1)+ 25(16-(16/2.4))$$

$$=200+ 150 +350$$

$$=700\text{ ns}$$

Contention Time & Copy back Caches

- In a simple copy back cache processor ceases on cache miss and does not resume until dirty line (w = probability of dirty line) is written back to main memory and new line read into the cache.

The Miss time penalty thus is

$$T_{\text{miss}} = (1+w) T_{\text{line access}}$$

Contention Time & Copy back Caches

- Miss time may be different for cache and main memory.
 - $T_{c.miss}$ = Time processor is idle due to cache miss.
 - $T_{m.miss}$ = Total time main memory takes to process a miss.
 - $T_{busy} = T_{m.miss} - T_{c.miss}$: Potential Contention time.
 - T_{busy} is =0 for normal CBWA cache

Contention Time & Copy back Caches

- Consider a case when dirty line is written to a write buffer when new line is read into cache. When processor resumes dirty line is written back to memory from buffer.

$T_{m.miss} = (1+w) T_{line\ access}$.

$T_{c.mis} = T_{line\ access}$

So $T_{busy} = w \cdot T_{line\ access}$.

- In case of wrap around load.

$T_{busy} = (1+w) T_{line\ access} - T_a$

Contention Time & Copy back Caches

- If processor creates a miss during T busy we call additional delay as T interference.

T interference = Expected number of misses during T busy.

= No of requests during T busy x prob of miss.

= $\lambda_p \cdot T_{\text{busy}} \cdot F$: where λ_p is processor request rate.

The delay factor given a miss during Tbusy is simply estimated as $T_{\text{busy}} / 2$

So T interference = $\lambda_p \cdot T_{\text{busy}} \cdot F \cdot T_{\text{busy}} / 2$

Contention Time & Copy back Caches

T interference = $\lambda_p \cdot f \cdot (T_{\text{busy}})^2 / 2$ and
total miss time seen from processor

T miss = T c.miss + T interference. And
Relative processor performance

Perf_{rel} = $1 / (1 + f \lambda_p T_{\text{miss}})$