

# Cache Effect Simulator (CESim) for Virtual Network I/O

Ryota Kawashima  
kawa1983@ieee.org

Here, we explain our mathematical model of packet forwarding within a commodity server. The purpose of the model is to analyze the *pure* effect of CPU caches on the throughput, and to explain performance characteristics of actual packet forwarding from a fundamental viewpoint. Hence, we simplify the model so that minimizing the dependency on actual performance metrics.

TABLE I  
DEFINITION OF PRIMITIVE SYMBOLS

Symbol	Description
<b>Input Parameters</b>	
(from actual performance metrics)	
$T_0$	Maximum throughput [Mpps]
$N_0$	The number of L1 cache accesses per packet (hits + misses)
<b>Variable</b>	
$r_1$	Hit ratio of the L1 cache
<b>Constants</b>	
$f$	Clock frequency of the CPU [GHz]
$c_{L1}$	CPU cycles consumed to access L1 cache
$c_{L2}$	CPU cycles consumed to access L2 cache
$c_{L3}$	CPU cycles consumed to access L3 cache
$l_{mem}$	Access latency to the RAM [ns]
$r_2$	Hit ratio of the L2 cache
$r_3$	Hit ratio of the L3 cache
<b>Others</b>	
$\alpha$	A ratio of pure processing instructions
$\beta$	Acceleration factor

Table I shows the definitions of the primitive symbols used in the model. First, we introduce only two parameters from the actual performance metrics:  $T_0$  and  $N_0$ . These parameters are from a single evaluation program that achieved highest throughput. A hit ratio of the L1 cache ( $r_1$ ) is the only variable whose value is changing during the simulation. In addition, several constant values are introduced in the model, such as access latency to data in each hierarchy level (including the RAM). Note that hit ratios of the L2 and L3 caches are constant in the model because they do not show positive correlations with the actual throughput. Hence, we set 0.70 and 0.79 for the values, respectively. We introduce a parameter  $\alpha$  that expresses the ratio of the execution time for non-data access CPU instructions to the overall ones. Since true value of the parameter cannot be obtained, we examine three values, 0.1, 0.5, and 0.9, in the simulation.

The simulation outputs throughput value in packet-per-second format, corresponding to the hit ratio of the L1 cache. Therefore, the throughput can be derived from the latency of per-packet handling (i.e., *processing* time + *cache/memory access* time) as given in eq. (1).

$$Throughput = \frac{10^3}{Latency} \quad [\text{Mpps}] \quad (1)$$

The value of *Latency* can be obtained by eq. (2-3),

$$Latency = L_{proc} + L_{access} \quad (2)$$

$$= \alpha \frac{1}{T_0} + N L_{L1+} \quad [\text{ns/packet}] \quad (3)$$

where  $L_{proc}$  and  $L_{access}$  denote the per-packet latency of non-data-access (pure processing) and data-access instructions in

total, respectively.  $N$  is the *theoretical* number of accesses to the L1 cache per packet. The value can be calculated from the maximum throughput ( $T_0$ ), a latency of single access to the L1 cache ( $l_{L1}$ ), and the parameter  $\alpha$ , as shown in eq. (4-5).

$$N = \frac{1 - \alpha}{T_0 l_{L1}} \quad [\text{times/packet}] \quad (4)$$

$$l_{L1} = \frac{c_{L1}}{f} \quad [\text{ns}] \quad (5)$$

Equation (6-7) expresses a calculation of the *expected* latency ( $L_{L1+}$ ) to load a single cacheline from the L1 cache, and the value can vary depending on the hit ratio ( $r_1$ ). The remaining  $L_{L2+}$  and  $L_{L3+}$  can also be calculated in a similar way to  $L_{L1+}$ .

$$L_{L1+} = L_{L1_{hit}} + L_{L1_{miss}} \quad (6)$$

$$= r_1 l_{L1} + (1 - r_1) L_{L2+} \quad [\text{ns/packet}] \quad (7)$$

In practice, the number of *actual* accesses to the L1 cache (per packet) is extremely different from the theoretical one ( $N$ ). Then, we introduced an acceleration (parallelization) factor in our model such that the number of apparent access times can vary depending on the hit ratio as given in eq. (8).

$$N^* = r_1^\beta N + (1 - r_1^\beta) N_0 \quad [\text{times/packet}] \quad (8)$$

The degree of parallelization ( $\beta$ ) can be obtained by eq. (9).

$$\beta = \frac{N_0}{N} \quad (9)$$

Finally, we obtain the modified version of the model that supports the parallelization factor as

$$\text{Throughput} = \frac{10^3}{\alpha \frac{1}{T_0} + N^* L_{L1+}} \quad [\text{Mpps}] \quad (10)$$

We demonstrate the model with a dedicated simulator (CESim), and show the results in Fig. 1.

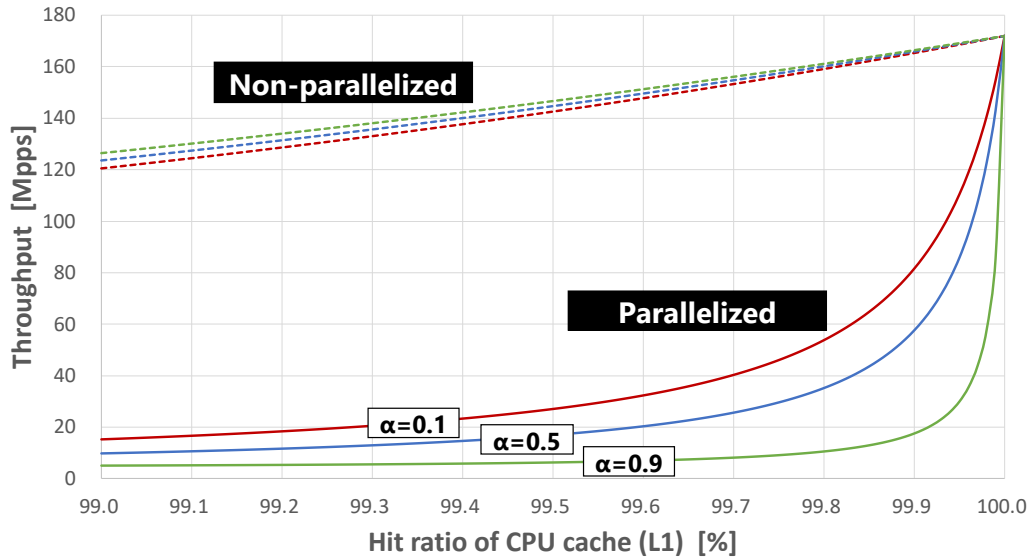


Fig. 1. Throughput vs. L1 cache hit ratio (simulation)