

PPW with Dark Silicon and Thermal Model

I. NEW METHOD

In order to estimate the performance-per-watt of a multi-core chip for different number of active cores and various active core distributions, we first propose an iteration based full-chip power estimation method, which can apply the thermal model into the analysis of performance-per-watt of multi-core chips. Furthermore, a non-iteration based method is proposed, which can implement local linearization to avoid time-consuming iterations. Additionally, a greedy based method can be integrated into the non-iteration based power estimation method to achieve further acceleration.

A. Iteration based leakage-aware power estimation

Because of the dependency of static power on temperature, it's not straightforward to compute the static power and temperature of next steady state based on current static power and temperature. Iteration method can be implemented to solve such problem, the computation flow is shown in Fig.x.

The initial value of $P_s^0(T, t + h)$ is a guess we provide based on the process technology. The temperature distribution $T(t + h)^0$ can be calculated with such guess. $P_s^1(T, t + h)$, the static power of next time step is updated with $T(t + h)^0$. Next, the temperature distribution $T(t + h)^1$ can be derived from $P_s^1(T, t + h)$, which concludes one iteration loop. Such iteration goes on until the convergence test is satisfied as $\|P_s^i(T, t + h) - P_s^{i-1}(T, t + h)\| < \epsilon$. Finally, the static power and temperature of steady state is outputted.

The iteration based method can produce an accurate outcome providing the ϵ is chosen to be small enough, yet the computing time is a serious problem, especially when the number of cores is large enough.

B. Local linearized thermal model

C. Non-iteration based power estimation

D. Greedy based acceleration of power estimation

In the coming many-core era, due to the tight power budget, power efficiency is critical for many-core processor design. In the previous work by Dong Hyuk Woo, evaluation of energy efficiency on the basis of performance and power (PPW) models is developed, which shows the tendency of PPW with number of cores. We implement the dark silicon and thermal model into the evaluation of PPW to gain a better understanding of PPW in the dark silicon era.

It is widely acknowledged that the total power of a chip is composed of dynamic and static power. The dynamic power is dependent on the activities of the chip, therefore it's easily estimated by methods such as performance counter. Yet the

static power p_s of the chip is mainly affected by temperature, for it's caused by leakage current I_{leak} as

$$p_s = V_{dd}I_{leak} \quad (1)$$

Due to the non-linear relationship between I_{leak} and temperature, p_s cannot be calculated directly. The iteration method is traditionally implemented to solve such problems.

To reduce the long computing time caused by iteration method, the leakage current I_{leak} can be linearized to eliminate the non-linearity between p_s and temperature, thus accelerating the computation. Taylor expansion is performed on the original I_{leak} model at a expansion point T_0 . Thus the linearized relation of I_{leak} and temperature is obtained. The relation between static power and temperature in linear form can also be achieved as:

$$\begin{aligned} p_s &= V_{dd}I_{leak} \\ &= V_{dd} \times (I_{lin} + I_{gate}) \\ &= V_{dd} \times (I_{lin}(T_p) + I_{const}) \end{aligned} \quad (2)$$

The linearized static power equation in matrix form is

$$P_s = P_0 + A_s T \quad (3)$$

The above mentioned Taylor expansion based method still costs too much computing time, it's especially impractical when the number of cores is too big. To further accelerate the processing, a greedy based method can be applied to handle this problem.

The average power consumption of the many-core processor is as follows:

$$W = \frac{P_1 \times (1 - f) + P_n \times f/n}{(1 - f) + f/n} \quad (4)$$

P_1 is the power consumption during the sequential computation phase, P_n is the power consumption during the parallel computation phase.

$Perf/W$ of a many-core processor is expressed as

$$\begin{aligned} \frac{Perf}{W} &= \frac{1}{(1 - f) + f/n} \times \frac{(1 - f) + f/n}{P_1 \times (1 - f) + P_n \times f/n} \\ &= \frac{1}{P_1 \times (1 - f) + P_n \times f/n} \end{aligned} \quad (5)$$

E. Power consumption with thermal constraints

P_n is made up of the power consumption of n active cores. Please note that in dark silicon era, $P_n = nk$ is not the correct expression. Due to thermal constraint, Dynamic Voltage and Frequency Scaling(DVFS) is necessary for Thermal Design Power(TDP) consideration. As a result, the power consumption of each of the core may decrease with the increase of the number of active cores. From xx, we have

$$(G - B_c A_s)T(t) + C \frac{dT(t)}{dt} = B_c(P_d(t) + P_0) \quad (6)$$

This work is supported in part by National Natural Science Foundation of China under grant No. 61404024, in part by the Scientific Research Foundation for the Returned Overseas Chinese Scholars, State Education Ministry.

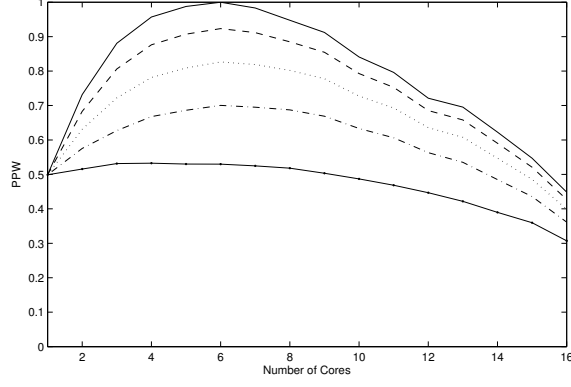


Fig. 1. The relation between PPW and the number of cores.

By applying thermal model, the effect of temperature on cores can be introduced. Through ergodic method, the distribution of light core with the maximum PPW for number of light core from 1 to n can be specified.

In (2) and (3), P_1 is consisted of the power consumption of $n-1$ idle cores and 1 active core. The expression of $P_1 = 1 + (n-1)k$ is not implemented, for the power consumption of idle cores and active core is not constant, due to the influence of the temperature.

By appending idle core model into the above mentioned thermal model, the distribution of idle core and active core with the maximum PPW for number of light core from 1 to n can be specified.

F.

II. CONCLUSION