# IBM Research Report

# Estimate of Resource Consumption Using Work Index Matrix

**G. Grabarnik, L. Kozakov, Sheng Ma**
IBM Research Division
Thomas J. Watson Research Center
P.O. Box 218
Yorktown Heights, NY 10598

**Research Division**
**Almaden - Austin - Beijing - Haifa - India - T. J. Watson - Tokyo - Zurich**

# Estimate of resource consumption using work index matrix.

G. Grabarnik, L. Kozakov, Sheng Ma
IBM Thomas J. Watson Research Center
Hawthorne, New York
{genady, kozakov, shengma }@us.ibm.com

## Abstract

*Accurate initial estimate of resource consumption is crucial for managing and planning resources and tasks in various computer systems. This short note suggests a new approach that provides initial estimate of the resource consumption for a given task with any background loading of the computer system. We introduce a notion of work index matrix, which is an invariant for specific computer system, and consider any work done by a computer system as a path in the potential field of resources.*

## 1. Motivation

Resource utilization managing and planning is an important issue in various situations, for instance, in autonomic computing systems, in distributed computing environments, in real time systems, etc. During planning phase of any computer job it is essential to estimate what are the necessary resources to do the job. Following are some possibilities for getting that info:

- Using default estimates (guessed ).

- Using historical data stored in some sort of repository.

In both cases, one of the issues is low accuracy of the estimate, which leads to a wrong planning. It is also not clear how to use the previous experience. One of the useful approaches to the planning is to correct a plan by looking how fast appropriate job is executed compared to the other executions or guesses. In this case, the planning gets reduced to poor initial approximation and good reactive correction that requires additional efforts, like dynamic monitoring of resource utilization.

There is a clear need for more accurate initial estimate of the task completion time or resources that would essentially depend on the configuration and other characteristics of a target computer.
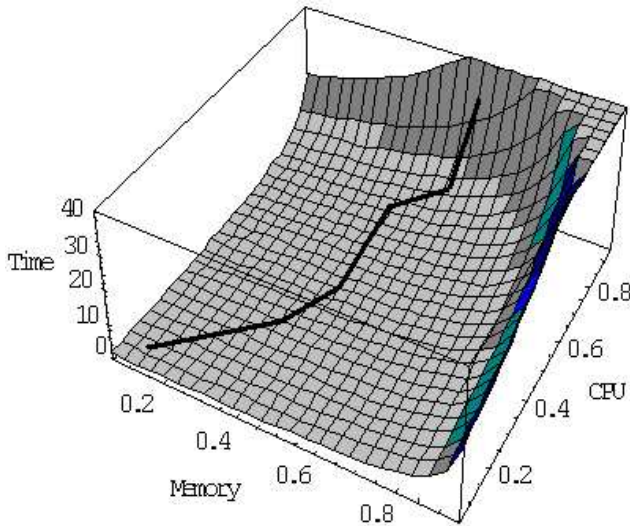
## 2. Main notions

In this section we introduce the main conceptions used in this paper, and present some intuitive analogies. First, we introduce the measure that will be used to characterize the elapse time required to complete a given work unit. We call this measure Rate of Execution (RoE). We can distinguish Computer RoE (CRoE), which characterizes the given individual computer configuration, and System RoE (SRoE), which may be related to a system of computers. RoE for a given work unit depends on two characteristics: the amount of resources that are required to complete the work unit, and the amount of resources that are already consumed by current background loading. We call them Work Unit Resource Consumption (WURC), and System Resource Consumption (SRC) correspondingly.

In the 3-dimensional space (time, memory, CPU) the RoE can be represented as a surface (see Fig. 1), where each point (T, Mem, Cpu) represents time (T) required to complete a given work unit under the background loading that consumed certain portion of the total memory (Mem) and CPU time (Cpu).

Based on the normalized RoE surface (built for a normalized work unit), we can construct a matrix that contains linear distortions to the RoE surface in each point of the resource space grid. We call this construction Work Index Matrix (WIM).

Note that WIM, in fact, is an invariant of the computer system configuration that provides quantitative measure of the system ability to perform any given work unit under any possible background loading. The notion of RoE is analogous to the notion of potential field in physics, and the WIM elements are analogous to the normalized potential values in selected points. Within the physical analogy, the WURC is equivalent to the physical work in the potential field.

**Figure 1. RoE surface**

## 3. Measuring RoE

Now, we explain how the RoE can be measured in a real computer system. RoE is essentially a measure in the multidimensional resource space. Each RoE value is associated with a given work unit performed under a given background loading, in other words, in a given SRC point. To measure RoE value in one SRC point we create appropriate loading levels in each dimension of the resource space, and measure the time required to complete 'normalized' work unit. As shown in Fig. 1, for lower levels of SRC there is almost no dependency between SRC and RoE values. The closer is SRC to 1, the higher is the dependency. In [1] this effect is described as 'saturation'. There may be also a dependency between different resources, and it makes sense to find most influential and less correlated resources (reduce dimensionality of the problem).

## 4. Task resource consumption estimate

This section describes usage of WIM values to estimate task resource consumption (TRC). Following information serves as an input for the method: computer identification, WIM, task as a sequence of work unit elements, measured with respect to some SRC value (say, zero load), target SRC. Method outputs estimated time as well as the sequence of TRC values. We draw the given task as a path on RoE starting at a point corresponding to a given initial load. Then, each element of the task is normalized by the associated WIM value at the element's starting point, and transformed using WIM to its new value under the new load computed consequently. As a result we get a new path starting from the point of target background resource consumption. Further, for low levels of WURC, since dependency on SRC is negligible, linear approximation may be used. For higher levels of WURC one gets linear PDE which should be solved to get estimated task resource consumption. Resulting path may be used for calculation of, say, max required resources, or average required resources etc.

## 5. Summary and Related Work

Initial time and resource consumption estimate is an important problem for provisioning systems. We outline a method that provides much more accurate estimate than existing methods. Similar method may be used in order to produce initial estimate for different target computers.

We refer for all standard definitions to the book [1]. This book contains very broad description of possible measurements related to the computer systems. In order to figure out the resource (time) consumption for specific system one needs both the system characteristics and characteristics of the workload. Characteristics of the workload are called workload indices. For example, in [2] automated estimate of the workload indices is done by using neural networks. In order to do that one needs to run similar workload on the computer system to teach the neural network. In this case, computer related characteristics are relearned by neural network for every new type of job. Number of works on measuring specific application behavior under different system loads were done by the school of prof. M. Seltzer (see, for example, [3][4]). Other area of usage is presented in [5], which deals with CPU usage prediction in heterogenous active networks.

## References

[1] Measurement and Tuning of Computer Systems - Ferrari, Serazzi et al. - 1983

[2] Automated Learning Of Workload Measures For Load Balancing On A Distributed System - Pankaj Mehra, Benjamin W. Wah - ICPP, 1993

[3] Evaluating Windows NT terminal server performance - Wong A. Y. , M. Seltzer - USENIX, 1999

[4] Operating System Benchmarking in the Wake of Lmbench - Brown, A., Seltzer, M. - SIGMETRIX, 1997

[5] Predicting and Controlling Resource Usage in a Heterogeneous Active Network - V. Galtier, et all - 2001