# Building the Functional Performance Model of a Processor

Alexey Lastovetsky           Ravi Reddy           Robert Higgins

School of Computer Science and
Informatics
University College Dublin,
Belfield, Dublin 4, Ireland

alexey.lastovetsky@ucd.ie     manumachu.reddy@ucd.ie     robert.higgins@ucd.ie

## ABSTRACT

In this paper, we present an efficient procedure for building a piecewise linear function approximation of the speed function of a processor with hierarchical memory structure. The procedure tries to minimize the experimental time used for building the speed function approximation. We demonstrate the efficiency of our procedure by performing experiments with a matrix multiplication application and a Cholesky Factorization application that use memory hierarchy efficiently and a matrix multiplication application that uses memory hierarchy inefficiently on a local network of heterogeneous computers.

## Categories and Subject Descriptors

B.8.2 [**Hardware**]: Performance and Reliability – *Performance Analysis and Design Aids*; C.4 [**Computer Systems Organization**]: Performance of Systems – *Measurement Techniques, Modeling Techniques, Performance Attributes*; I.6.3 [**Computing Methodologies**]: Simulation and Modeling – *Applications*; I.6.5 [**Computing Methodologies**]: Simulation and Modeling – *Model Development*.

## General Terms

Algorithms, Measurement, Performance, Experimentation.

## Keywords

Performance modeling, processor performance, memory hierarchy, heterogeneous computing, parallel computing, distributed computing, Grid computing.
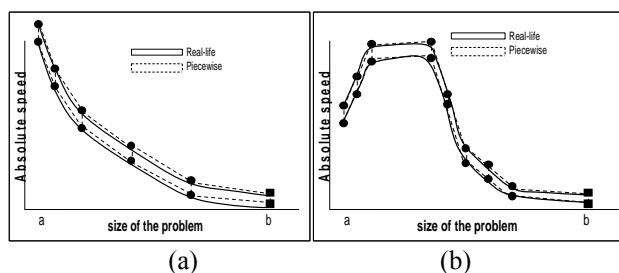
## 1. INTRODUCTION

The quality of problem partitioning and task scheduling for high performance computing on common, heterogeneous networks strongly depends on the accuracy and adequateness of the performance model of such networks. The following two issues have only been recently taken into account and incorporated in such performance models.

The first issue is that the optimal partitioning of the problem or the optimal scheduling of the given tasks can sometimes have

some sub-problems or tasks not fit into the main memory of the processors causing the paging. As a result, relative speeds of processors cannot be accurately approximated by constants independent of the problem size. The more realistic model proposed in [5] represents each processor by a continuous and relatively smooth function of the problem size. This model is application-centric meaning that different applications will characterize the speed of the processor by different functions. Efficient algorithms of problem partitioning and task scheduling with this model have been also designed in [5].

The second issue is that a computer in a common network typically experiences constant and stochastic fluctuations in the workload. This changing transient load will cause a fluctuation in the speed of the computer in the sense that the execution time of the same task of the same size will vary for different runs at different times. The natural way to represent the inherent fluctuations in the speed is to use a speed band rather than a speed function. The width of the band characterizes the level of fluctuation in the performance due to changes in load over time. The band model was proposed in [6], and some algorithms of problem partitioning with this model were designed in [4].

In this paper, we present a procedure for building a piecewise linear function approximation of the speed band of a processor with hierarchical memory structure. Any speed function representing the performance of the processor should fit into the speed band. The procedure tries to minimize the experimental time used to build the approximation of the speed band.



**Figure 1. Using piecewise linear approximation to build speed bands for 2 processors. Circular points are experimentally obtained, square points are calculated using heuristics. (a) The speed band is built from 5 experimental points, application uses the memory hierarchy inefficiently. (b) The speed band is built from 8 experimental points, application uses memory hierarchy efficiently.**

Sample piece-wise linear function approximations of the speed band of a processor are shown in Figure 1. Each approximation is built using a set of few experimentally obtained points. The more points used to build the approximation, the more accurate the

approximation. However it is prohibitively expensive to use large number of points. Hence an optimal set of few points needs to be chosen to build an efficient piecewise linear function approximation of the speed band. Such an approximation gives the speed of the processor for any problem size with certain accuracy within the inherent deviation of the performance of computers typically observed in the network.

The rest of the paper is organized as follows. In section 2, we formulate the problem of building a piecewise linear function approximation of a processor and present an efficient and a practical procedure to solve the problem. To demonstrate the efficiency of our procedure, we perform experiments using a matrix multiplication application and a Cholesky Factorization application that use memory hierarchy efficiently and a matrix multiplication application that uses memory hierarchy inefficiently on a local network of heterogeneous computers.

## 2. PROCEDURE FOR BUILDING SPEED FUNCTION APPROXIMATION

This section is organized as follows. We start with the formulation of the speed band approximation building problem. This is followed by a section on obtaining the load functions characterizing the level of fluctuation in load over time. The third section presents the assumptions adopted by our procedure. We then present some operations and relations related to the piecewise linear function approximation of the speed band. And finally we explain our procedure to build the piecewise linear function approximation.

### 2.1 Problem Formulation

For a given application in a real-life situation, the performance demonstrated by the processor is characterized by a speed band representing the speed function of the processor with the width of the band characterizing the level of fluctuation in the speed due to changes in load over time.

The problem is to find experimentally an approximation of the speed band of the processor that can represent the speed band with sufficient accuracy and at the same time spend minimum experimental time to build the approximation. One such approximation is a piecewise linear function approximation which accurately represents the real-life speed band with a finite number of points. This is shown in Figure 2.
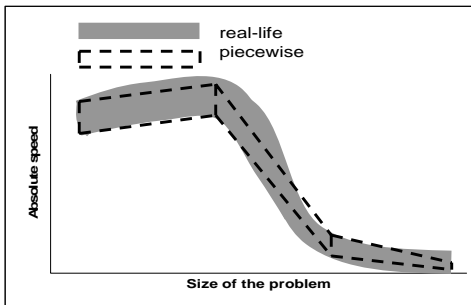


**Figure 2. Real-life speed band of a processor and a piecewise linear function approximation the speed band.**

The piecewise linear function approximation of the speed band of the processor is built using a set of experimentally obtained points for different problem sizes. To obtain an experimental point for a problem size $x$ (we define the size of the problem to be the amount of data stored and processed by the application), we execute the application for this problem size. We measure the ideal execution time $t_{ideal}$. We define $t_{ideal}$ as the time it would require to solve the problem on a completely idle processor. For example on UNIX platforms, this information can be obtained by using the ***time*** utility or the ***getrusage()*** system call, summing the reported user and system cpu seconds a process has consumed. The ideal speed of execution $s_{ideal}$ is then equal to the volume of computations divided by $t_{ideal}$. We assume we have the load functions: $l_{max}(t)$ and $l_{min}(t)$, which are the maximum and minimum load averages observed over increasing time periods. The load average is the number of active processes running on the processor at any time. We make a prediction of the maximum and minimum average load, $l_{max,predicted}(x)$ and $l_{min,predicted}(x)$ respectively, that would occur during the execution of the application for the problem size $x$. The creation of the functions $l_{max}(t)$ and $l_{min}(t)$ and predicting the load averages are explained in detail in the next section. Using $s_{ideal}$ and the load averages predicted, we calculate $s_{max}(x)$ and $s_{min}(x)$ for a problem size $x$:

$$s_{max}(x) = s_{ideal}(x) - l_{min,predicted}(x) \times s_{ideal}(x)$$

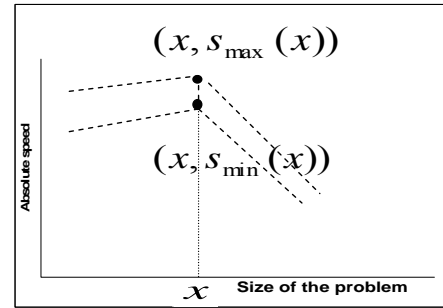$$s_{min}(x) = s_{ideal}(x) - l_{max,predicted}(x) \times s_{ideal}(x)$$



**Figure 3(a). The speeds $s_{max}(x)$ and $s_{min}(x)$ representing a cut of the real band used to build the piecewise linear approximation.**
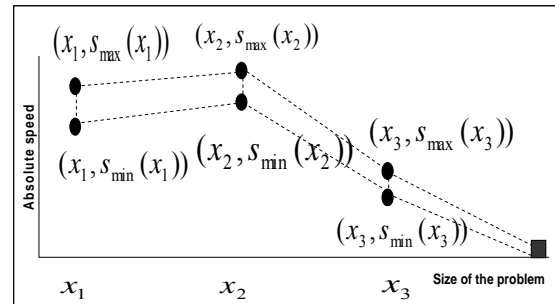


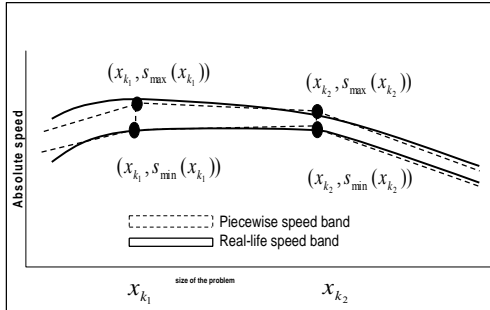**Figure 3(b). Piecewise linear approximation built by connecting the cuts.**

The experimental point is then given by a vertical line connecting the points $(x, s_{max}(x))$ and $(x, s_{min}(x))$. We call this vertical line the "cut" of the real band. This is illustrated in Figure 3(a). The difference between the speeds $s_{max}(x)$ and $s_{min}(x)$ represents the level of fluctuation in the speed due to changes in load during the execution of the problem size $x$. The piecewise linear approximation is obtained by connecting these experimental points as shown in Figure 3(b). So the problem of building the piecewise linear function approximation is to find a

set of such experimental points that can represent the speed band with sufficient accuracy and at the same time spend minimum experimental time to build the piecewise linear function approximation.

Mathematically the problem of building piecewise linear function approximation can be formulated as follows:

**Definition 1**. *Piecewise Linear Function Approximation Building Problem PLFABP($l_{min}(t),l_{max}(t)$):* Given the functions $l_{min}(t)$ and $l_{max}(t)$ ($l_{min}(t)$ and $l_{max}(t)$ are functions of the time characterizing the level of fluctuation in load over time), obtain a set of n experimental points representing the piecewise linear function approximation of the speed band of a processor, each point representing a cut given by $(x_i, s_{max}(x_i))$ and $(x_i, s_{min}(x_i))$ where $x_i$ is the size of the problem and $s_{max}(x_i)$ and $s_{min}(x_i)$ are speeds calculated based on the functions $l_{min}(t)$ and $l_{max}(t)$ and ideal speed $s_{ideal}$ at point $i$, such that:

- The non-empty intersectional area of piecewise linear function approximation with the real-life speed band is a simply connected surface (A surface is said to be connected if a path can be drawn from every point contained within its boundaries to every other point. A topological space is simply connected if it is path connected and it has no holes. This is illustrated in Figure 4), and

- The sum $\displaystyle\sum_{i=1}^{n} t_i$ of the times is minimal where $t_i$ is the experimental time used to obtain point $i$.



**Figure 4. The non-empty intersectional area of piecewise linear function approximation with the real-life speed band is a simply connected surface.**

We provide an efficient and a practical procedure to build a piecewise linear function approximation of the speed function of a processor.
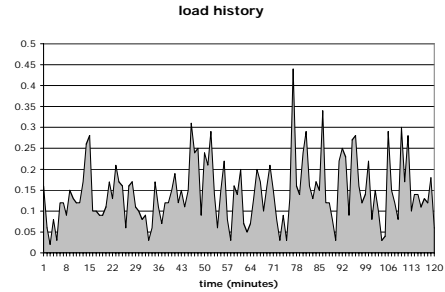
## 2.2  Load Functions

There are a number of experimental methods that can be used to obtain the functions $l_{min}(t)$ and $l_{max}(t)$ (characterizing the level of fluctuation in load over time) input to our procedure for building the piecewise linear function approximation.
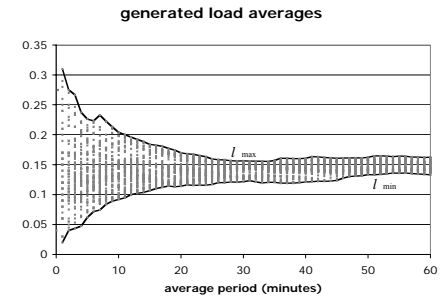
One method is to use the metric of Load Average. Load Average measures the number of active processes at any time. High load averages usually means that the system is being used heavily and the response time is correspondingly slow. A UNIX like operating system maintains three figures for averages over one, five and fifteen minute periods. There are alternative measures available through many utilities on various platforms such as **vmstat** (UNIX), **top** (UNIX), **perfmon** (Windows) or through performance probes and they may be combined to more

accurately represent utilization of a system under a variety of conditions [7]. For this paper we will use the load average metric only.

The load average data is represented by two piecewise linear functions: $l_{max}(t)$ and $l_{min}(t)$. The functions describe load averaged over increasing periods of time up to a limit $w$ as shown in Figure 5(b). This limit should be at most the running time of the largest foreseeable problem, which is the problem size where the speed of the processor can be assumed to be zero (this is given by problem size $b$ discussed in section 2.5). For execution of a problem with a running time greater than this limit, the values of the load functions at $w$ may be extended to infinity. The functions are built from load averages observed every $\Delta$ time units. One, five or fifteen minutes are convenient values for $\Delta$ as statistics for these time periods are provided by the operating system (using a system call **getloadavg()**). Alternate values of $\Delta$ would require additional monitoring of the load average and translation into $\Delta$ time unit load average.



**Figure 5(a). Load history is used to generate a range of load averages.**



**Figure 5(b). $l_{max}(t)$ and $l_{min}(t)$, the maximum and minimum loads calculated from the matrix of load averages A.**

The amount of load observations used in the calculation of $l_{max}(t)$ and $l_{min}(t)$ is given by $h$, the history. A sliding window with a length of $w$ passes over the $h$ most recent observations. At each position of the window a set of load averages is created. The set consists of load averages generated from the observations inside the window. If $\Delta$ were one minute, a one minute average would be given by the first observation in the window, a two minute average would be the average of the first and second observations in the window, and so on. While the window is positioned completely within the history, a total of $w$ load averages would be created in each set, the load averages having periods of $\Delta$, $2\Delta$, … $w\Delta$ time units. The window can move a total of $w$ times, but after the $(h – w)$-th time, its end will slide outside of the history. The sets of averages created at these positions will not range as far as $w\Delta$ but they are still useful. From all of these sets of averages,

maximum and minimum load averages for each time period $\Delta$, $2\Delta$, … $w\Delta$ are extracted and used to create the functions $l_{max}(t)$ and $l_{min}(t)$.

More formally, if we have a sequence of observed loads: $l_1$, $l_2$, …, $l_h$, then the matrix $A$ of load averages created from observations is defined as follows:

$$A = \begin{pmatrix} a_{1,1} & . & . & a_{1,h} \\ . & & & \times \\ . & & \times & \times \\ a_{w,1} & \times & \times & \times \end{pmatrix} \text{ where } a_{ij} = \frac{\sum_{k=j}^{i+j-1} l_k}{i \cdot \Delta},$$

for all $i = 1...h; j = 1...w$ and $i + j < h$

The elements marked as $\times$ in the matrix $A$ are not evaluated as the calculations would operate on observations taken beyond $l_h$. $l_{max}(t)$ and $l_{min}(t)$, are then defined by the maximum and minimum calculated $j$-th load averages respectively, i.e. the maximum or minimum value of a row $j$ in the matrix (see Figure 5(a) and 5(b)). Points are connected in sequence by straight-line segments to give a continuous piecewise function. The points are given by:

$$l_{max}(j) = \max_{i=1}^{h}\left(a_{ij}\right) \quad \text{and} \quad l_{min}(j) = \min_{i=1}^{h}\left(a_{ij}\right)$$

Initial generation of the array has been implemented with a complexity of $h \times w^2$. Maintaining the functions $l_{max}(t)$ and $l_{min}(t)$ after a new observation is made has a complexity of $w^2$. $\Delta$, $h$, and $w$ may be adjusted to ensure the generation and maintenance of the functions is not an intensive task.

When building the speed functions $s_{min}(x)$ and $s_{max}(x)$, we execute the application for a problem size $x$. We then measure the ideal time of execution $t_{ideal}$. We define $t_{ideal}$ as the time it would require to solve the problem on a completely idle processor. On UNIX platforms it is possible to measure the number of CPU seconds a process has used during the total time of its execution. This information is provided by the ***time*** utility or by the ***getrusage()*** system call. We assume that the number of CPU seconds a process has used is equivalent to time it would take to complete execution on a completely idle processor: $t_{ideal}$. We can then estimate the time of execution for the problem running under any load $l$ with the following function:

$$t(l) = \frac{1}{1-l} \times t_{ideal} \qquad (1)$$

This formula assumes that the system is uniprocessor, that no jobs are scheduled if the load is one or greater and that the task we are scheduling is to run as a ***nice***'d process (***nice*** is an operating system call that allows a process to change its priority), only using idle CPU cycles. These limitations fit the target of execution on non-dedicated platforms. If a job is introduced onto a system with a load of, for example, 0.1, the system has a 90% idle CPU, then the formula predicts that the job will take 1/0.9 times longer than the optimal time of execution: $t_{ideal}$.

In order to calculate the speed functions $s_{min}(x)$ and $s_{max}(x)$, we need to find the points where the function of performance degradation due to load (Formula 1) intersects with the history of maximal and minimal load $l_{max}(t)$ and $l_{min}(t)$ as shown in Figure 6. For a problem size $x$, the intersection points give the maximum and minimum predicted loads $l_{max,predicted}(x)$ and $l_{min,predicted}(x)$.

Using these loads, the speeds $s_{min}(x)$ and $s_{max}(x)$ for a problem size $x$ are calculated as:

$$s_{max}(x) = s_{ideal}(x) - l_{min,predicted}(x) \times s_{ideal}(x)$$

$$s_{min}(x) = s_{ideal}(x) - l_{max,predicted}(x) \times s_{ideal}(x)$$

where $s_{ideal}(x)$ is equal to the volume of computations involved in solving the problem size $x$ divided by the ideal time of execution $t_{ideal}$.
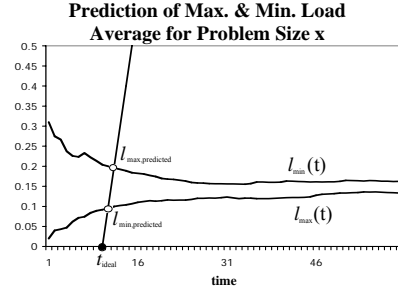


**Figure 6. Intersection of load and running time functions (Formula 1).**

## 2.3 Assumptions

We make some assumptions on the real-life speed band of a processor. Firstly, there are some shape requirements.

(a) We assume that the upper and lower curves of the speed band are continuous functions of the problem size.

(b) The permissible shapes of the real-life speed band are:

- The upper curve and the lower curve are both a non-increasing function of the size of the problem for all problem sizes (Figure 7(a)).

- The upper curve and the lower curve are both a non-decreasing function of the size of the problem followed by a non-increasing function of the size of the problem (Figure 7(b)).

(c) A straight line intersects the upper curve of the real-life speed band in no more than one point between its endpoints and the lower curve of the real-life speed band in no more than one point between its endpoints as shown for applications that use memory hierarchy inefficiently in Figure 7(a) and for applications that use memory hierarchy efficiently in Figure 7(b).

(d) We assume that the width of the real-life speed band, representing the level of fluctuations in speed due to changes in load over time, decreases as the problem size increases.

These assumptions are justified by experiments conducted with a range of applications differently using memory hierarchy presented in [6].

Secondly, we do not take into account the effects on the performance of the processor caused by several users running heavy computational tasks simultaneously. We suppose only one user running heavy computational tasks and multiple users performing routine computations and communications, which are not heavy like email clients, browsers, audio applications, text editors etc.
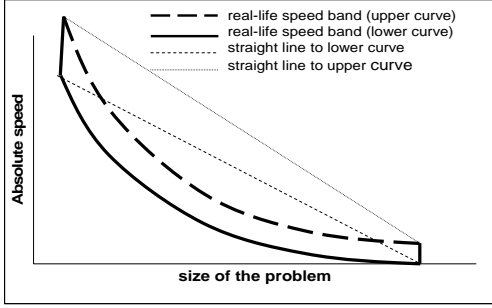
**Figure 7(a). Shape of real-life speed function of processor for applications that use memory hierarchy inefficiently.**
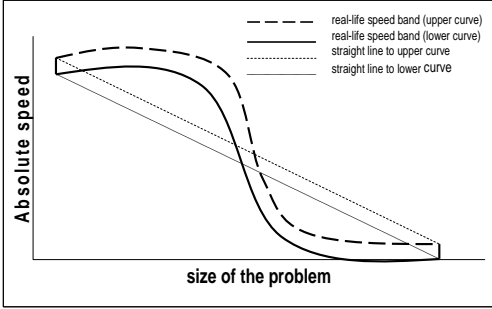


**Figure 7(b). Shape of real-life speed function of processor for applications that use memory hierarchy efficiently.**

## 2.4 Definitions

Before we present our procedure to build a piecewise linear function approximation of the speed band of a processor, we present some operations and relations on cuts that we use to describe the procedure. The piecewise linear function approximation of the speed band of the processor is built by connecting these cuts.

1. We use $I_x$ at problem size $x$ to represent the interval $(s_{min}(x), s_{max}(x))$. $I_x$ is the projection of the cut $C_x$ connecting the points $(x, s_{min}(x))$ and $(x, s_{max}(x))$ on the $y$-axis.

2. $I_x \le I_y$ if and only if $s_{max}(x) \le s_{max}(y)$ and $s_{min}(x) \le s_{min}(y)$.

3. $I_x \cap I_y$ represents intersection between the intervals $(s_{min}(x), s_{max}(x))$ and $(s_{min}(y), s_{max}(y))$. If $I_x \cap I_y = \emptyset$ where $\emptyset$ represents an empty set with no elements, then the intervals are disjoint. If $I_x \cap I_y = I_y$, then the interval $(s_{min}(x), s_{max}(x))$ contains the interval $(s_{min}(y), s_{max}(y))$, that is, $s_{max}(x) \ge s_{max}(y)$ and $s_{min}(x) \le s_{min}(y)$.

4. $I_x = I_y$ if and only if $I_x \le I_y$ and $I_y \le I_x$.

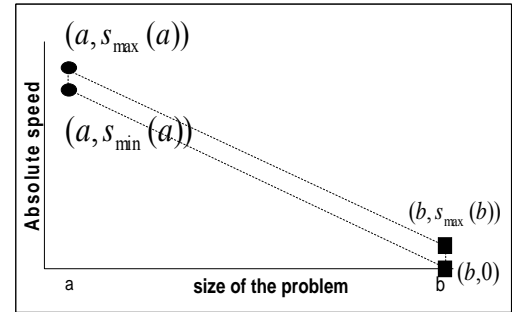## 2.5 Speed Function Approximation Building Procedure

**Procedure** *Geometric Bisection Building Procedure GBBP($l_{max}(t), l_{min}(t)$)*. The procedure to build the piecewise linear function approximation of the speed band of a processor consists of the following steps and is illustrated in Figure 8:
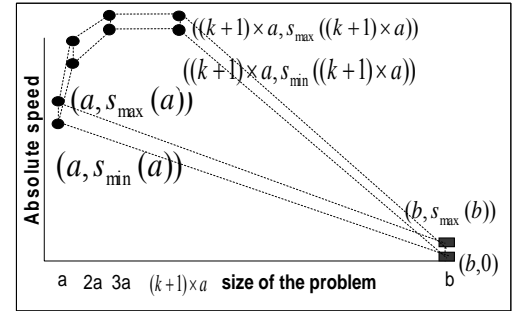
1. We select an interval $[a,b]$ of problem sizes where $a$ is some small size and $b$ is the problem size large enough to make the speed of the processor practically zero. In most cases, $a$ is the problem size that can fit into the top level of memory hierarchy of the computer (L1 cache) and $b$ is the problem size that is obtained based the maximum amount of memory

that can be allocated on heap. To calculate the problem size $b$, run a modified version of the application, which includes only the code that allocates memory on heap. For example consider a matrix-matrix multiplication application of two dense square matrices $A$ and $B$ of size $n \times n$ to calculate resulting matrix $C$ of size $n \times n$, the modified version of the application would just contain the allocation and de-allocation of matrices $A$, $B$, and $C$ on heap. This modified version is then run until the application fails due to exhaustion of heap memory, the problem size at this point gives $b$. It should be noted that finding the problem size $b$ by running the modified version should take just few seconds.
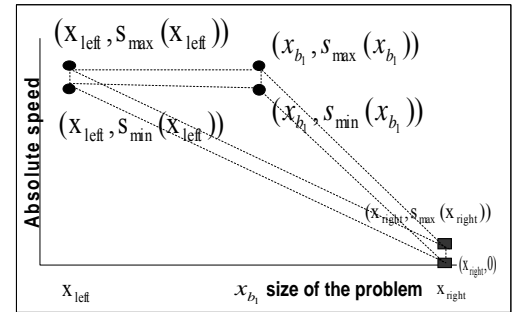
We obtain experimentally the speeds of the processor at point $a$ given by $s_{max}(a)$ and $s_{min}(a)$ and we set the absolute speed of the processor at point $b$ to 0. Our initial approximation of the speed band is a speed band connecting cuts $C_a$ and $C_b$. This is illustrated in Figure 8(a).
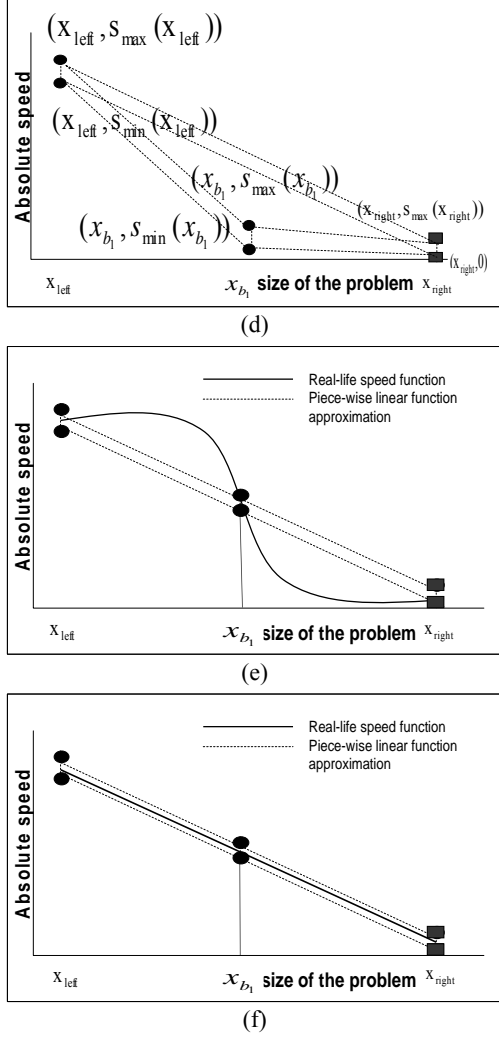


(a)



(b)



(c)

**(d)**

Figure (d) axis label: **Absolute speed** (vertical), horizontal axis: $x_{\text{left}}$, $x_{b_1}$ **size of the problem** $x_{\text{right}}$

Points labeled: $(x_{\text{left}}, S_{\max}(x_{\text{left}}))$, $(x_{\text{left}}, S_{\min}(x_{\text{left}}))$, $(x_{b_1}, S_{\max}(x_{b_1}))$, $(x_{b_1}, S_{\min}(x_{b_1}))$, $(x_{\text{right}}, S_{\max}(x_{\text{right}}))$, $(x_{\text{right}}, 0)$

**(e)**

Legend:
— Real-life speed function
····· Piece-wise linear function approximation

Axis: **Absolute speed** (vertical), $x_{\text{left}}$, $x_{b_1}$ **size of the problem** $x_{\text{right}}$

**(f)**

Legend:
— Real-life speed function
····· Piece-wise linear function approximation

Axis: **Absolute speed** (vertical), $x_{\text{left}}$, $x_{b_1}$ **size of the problem** $x_{\text{right}}$

**Figure 8. (a) to (f) Illustration of the procedure to obtain the piecewise linear function approximation of the speed band for a processor. Circular points are experimentally obtained points. Square points are points of intersection that are calculated but not experimentally obtained. White circular points are experimentally obtained and that fall in the current approximation of the speed band.**

2. We experiment with problem sizes $a$ and $2a$. If $I_{2a} \le I_a$ or $I_{2a} \cap I_a = I_{2a}$, we replace the current approximation of the trapezoidal speed band with two trapezoidal connected bands, the first one connecting the cuts $C_a$ and $C_{2a}$ and the second one connecting the cuts $C_{2a}$ and $C_b$. We then consider the interval $[2a,b]$ and apply step 3 of our procedure to this interval. The speed band in this interval connecting the cuts at problem sizes $2a$ and $b$ is input to step 3 of the procedure. We set $x_{\text{left}}$ to $2a$ and $x_{\text{right}}$ to $b$.

If $I_a \le I_{2a}$, we recursively apply this step until $I_{(k+1) \times a} \le I_{ka}$ or $I_{(k+1) \times a} \le I_{ka} = I_{(k+1) \times a}$. We replace the current approximation of the speed band in the interval $[k \times a, b]$ with two connected bands, the first one connecting the cuts $C_{ka}$ and $C_{(k+1) \times a}$ and the second one connecting the cuts $C_{(k+1) \times a}$ and $C_b$. We then consider the interval $[(k+1) \times a, b]$ and apply the step 3 of our

procedure to this interval. The speed band in this interval connecting the cuts $C_{(k+1) \times a}$ and $C_b$ is input to step 3 of the procedure. We set $x_{\text{left}}$ to $(k+1) \times a$ and $x_{\text{right}}$ to $b$. This is illustrated in Figure 8(b).
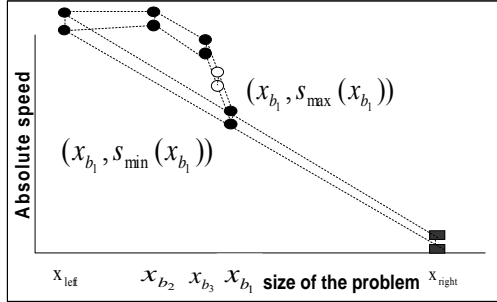
It should be noted that the time taken to obtain the cuts at problem sizes $\{a, 2a, 3a,…,(k+1) \times a\}$ is relatively small (usually milliseconds to seconds) compared to that for larger problem sizes (usually minutes to hours).

3. We bisect this interval $[x_{\text{left}}, x_{\text{right}}]$ into sub-intervals $[x_{\text{left}}, x_{b_1}]$ and $[x_{b_1}, x_{\text{right}}]$ of equal length. We obtain experimentally the cut $Cx_{b_1}$ at problem size $x_{b_1}$. We also calculate the cut of intersection of the line $x = x_{b_1}$ with the current approximation of the speed band connecting the cuts $Cx_{\text{left}}$ and $Cx_{\text{right}}$. The cut of intersection is given by $C'x_{b_1}$.

   a. If $Ix_{\text{left}} \cap Ix_{b_1} \ne \emptyset$, we replace the current approximation of the speed band with two connected bands, the first one connecting the cuts $Cx_{\text{left}}$ and $Cx_{b_1}$ and the second one connecting the cuts $Cx_{b_1}$ and $Cx_{\text{right}}$. This is illustrated in Figure 8(c). We stop building the approximation of the speed band in the interval $[x_{\text{left}}, x_{b_1}]$ and recursively apply step 3 for the interval $[x_{b_1}, x_{\text{right}}]$. We set $x_{\text{left}}$ to $x_{b_1}$.

   b. If $Ix_{\text{left}} \cap Ix_{b_1} = \emptyset$ and $Ix_{\text{right}} \cap Ix_{b_1} \ne \emptyset$, we replace the current approximation of the speed band with two connected bands, the first one connecting the cuts $Cx_{\text{left}}$ and $Cx_{b_1}$ and the second one connecting the cuts $Cx_{b_1}$ and $Cx_{\text{right}}$. This is illustrated in Figure 8(d). We stop building the approximation of the speed band in the interval $[x_{b_1}, x_{\text{right}}]$ and recursively apply step 3 for the interval $[x_{\text{left}}, x_{b_1}]$. We set $x_{\text{right}}$ to $x_{b_1}$.

   c. If $Ix_{\text{left}} \cap Ix_{b_1} = \emptyset$ and $Ix_{\text{right}} \cap Ix_{b_1} = \emptyset$ and $Ix_{b_1} \cap I'x_{b_1} \ne \emptyset$, then we have two scenarios illustrated in Figures 8(e) and 8(f) where experimental point at the first point of bisection falls in the current approximation of the speed band just by accident.
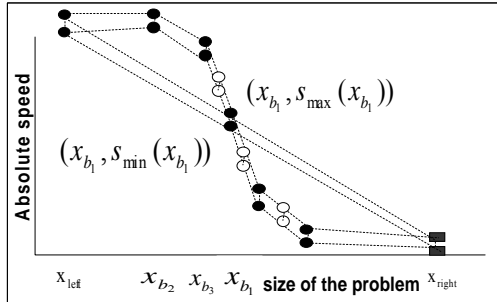
      Consider the interval $[x_{\text{left}}, x_{b_1}]$. It is bisected at the point $x_{b_2}$. We obtain experimentally the cut $Cx_{b_2}$ at problem size $x_{b_2}$. We also calculate the cut of intersection $C'x_{b_2}$ of the line $x = x_{b_2}$ with the current approximation of the speed band. If $Ix_{b_2} \cap I'x_{b_2} \ne \emptyset$, we stop building the approximation of the speed function in the interval $[x_{\text{left}}, x_{b_1}]$ and we replace the current approximation of the trapezoidal speed band in the interval $[x_{\text{left}}, x_{b_1}]$ with two connected bands, the first one connecting the cuts $Cx_{\text{left}}$ and $Cx_{b_2}$ and the second one connecting the points $Cx_{b_1}$ and $Cx_{b_2}$. Since we have obtained the cut at problem size $x_{b_2}$ experimentally, we use it in our approximation. This is chosen as our final piece of our piece-wise linear function approximation in the interval $[x_{\text{left}}, x_{b_1}]$. If $Ix_{b_2} \cap I'x_{b_2} = \emptyset$, the intervals $[x_{\text{left}}, x_{b_2}]$ and $[x_{b_2}, x_{b_1}]$ are recursively bisected using step 3. Figure 9(a) illustrates the procedure.

d. Consider the interval $[x_{b_1}, x_{\text{right}}]$. This interval is recursively bisected using step 3. We set $x_{\text{left}}$ to $x_{b_1}$. Figure 9(b) illustrates the procedure. If $Ix_{\text{left}} \cap Ix_{b_1} = \emptyset$ and $Ix_{\text{right}} \cap Ix_{b_1} = \emptyset$ and $Ix_{b_1} \leq I'x_{b_1}$ and $Ix_{b_1} \cap I'x_{b_1} = \emptyset$, we replace the current approximation of the speed band with two connected bands, the first one connecting the cuts $Cx_{\text{left}}$ and $Cx_{b_1}$ and the second one connecting the cuts $Cx_{b_1}$ and $Cx_{\text{right}}$. This is illustrated in Figure 9(c).
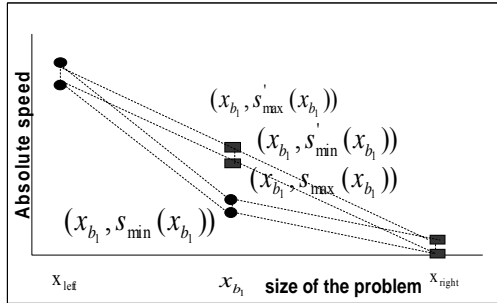
The intervals $[x_{\text{left}}, x_{b_1}]$ and $[x_{b_1}, x_{\text{right}}]$ are recursively bisected using step 3. Figure 9(d) illustrates the procedure.
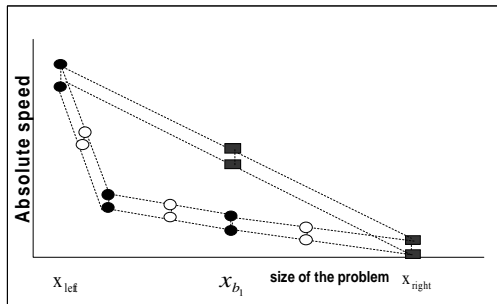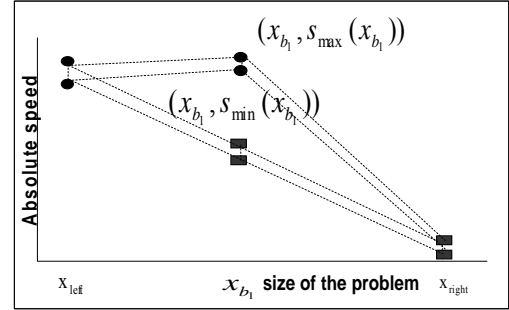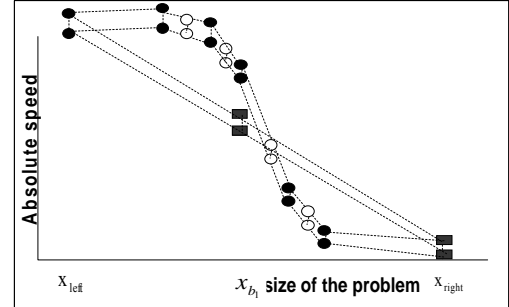


(a)



(b)



(c)



(d)



(e)



(f)

**Figure 9. (a) to (f) Illustration of the procedure to obtain the piecewise linear function approximation of the speed band for a processor. Circular points are experimentally obtained. Square points are calculated. White circular points are experimentally obtained and fall in the current approximation of the speed band.**

e. If $Ix_{\text{left}} \cap Ix_{b_1} = \emptyset$ and $Ix_{\text{right}} \cap Ix_{b_1} = \emptyset$ and $I'x_{b_1} \leq Ix_{b_1}$ and $Ix_{b_1} \cap I'x_{b_1} = \emptyset$, we replace the current approximation of the speed band with two connected bands, the first one connecting the cuts $Cx_{\text{left}}$ and $Cx_{b_1}$ and the second one connecting the cuts $Cx_{b_1}$ and $Cx_{\text{right}}$. This is illustrated in Figure 9(e).

The interval $[x_{\text{left}}, x_{b_1}]$ and $[x_{b_1}, x_{\text{right}}]$ are recursively bisected using step 3. Figure 9(f) illustrates the procedure.

4. The stopping criterion of the procedure is satisfied when we don't have any sub-interval to divide. Figures 1(a) and 1(b) show the final piecewise linear function approximation of the speed band of the processor for an application that uses memory hierarchy efficiently and an application using memory hierarchy inefficiently.
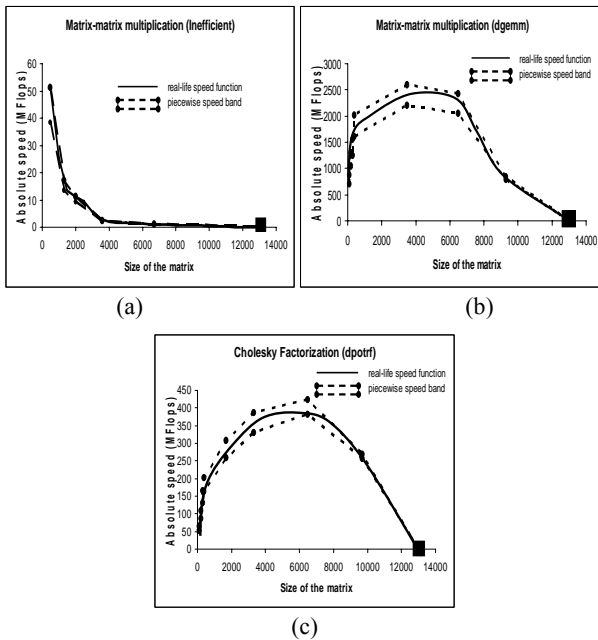
## 3. Experimental Results

We consider a 2Ghz Pentium Xeon Linux workstation with 1GB of RAM (X1) and a 440Mhz UltraSparc IIi Solaris workstation with 512MB of RAM (X2). They are integrated into local departmental network in the experiments. Each machine has roughly 400MB of available RAM, the rest being consumed by the operating system processes and other user applications. These extra processes perform routine computations and communications, applications such as email clients, browsers, text editors, audio applications use a constant percentage of CPU.

There are three applications used to demonstrate the efficiency of our procedure to build the piecewise linear function approximation of the speed band of a processor. The first

application is Cholesky Factorization of a dense square matrix employing the LAPACK [1] routine **dpotrf**. The second application is matrix-matrix multiplication of two dense matrices using memory hierarchy inefficiently. The third application is based on matrix-matrix multiplication of two dense matrices employing the level-3 BLAS routine **dgemm** [2] supplied by Automatically Tuned Linear Algebra Software (ATLAS) [3].

Figures 10(a) to 10(e) show the real-life speed function and the piecewise linear function approximation of the speed band of the processors X1 and X2 for the matrix multiplication and Cholesky Factorization applications. The real-life speed function for a processor is built using a set of experimentally obtained points $(x,s)$ . To obtain an experimental point for a problem size $x$, we execute the application for the problem size at that point. The absolute speed of the processor $s$ for this problem size is obtained by dividing the total volume of computations by the real execution time (and not the ideal execution time).



(a)                                        (b)



(c)

**Figure 10. Piecewise linear approximation of the speed band against the real-life speed function. Circular points are experimentally obtained points. Square points are calculated but not experimentally obtained. (a) Matrix-matrix multiplication using memory hierarchy inefficiently on X1. (b) Matrix-matrix multiplication using ATLAS on X1. (c) Cholesky Factorization using ATLAS on X2.**

Table 1 shows the speedup of Geometric Bisection Building Procedure (GBBP) over a naïve procedure. The naïve procedure divides the interval [a,b] of problem sizes equally into $n$ points. The application is executed for each of the problem sizes $\{a,(a+(b-a)/n),(a+2\times(b-a)/n),\ldots,b\}$ to obtain the experimental points to build the piecewise linear function approximation of the speed band. In our experiments, we have used 20 points. The speedup calculated is equal to the ratio of the experimental time taken to build the piecewise linear function approximation of the speed band using the naïve procedure over the experimental time taken to build the piecewise linear function approximation of the speed band.

We measured the accuracy of the load average functions $l_{max}(t)$ and $l_{min}(t)$ by counting how often a future load was found to be within the bounds of the curves and by measuring the area between the curves. A very wide band will encompass almost all future loads but the prediction of maximum and minimum load will be poor. We fixed $w$, the window size, and varied $h$ to examine how the hit ratio and area of the band changed. Predictions of load X1, a machine operating as a desktop with constant minor fluctuations in load, using a 60 minute window, were most accurate with 4 hours worth of historical load data.

**Table 1. Speedup of GBBP procedure over naïve procedure.**

| CPU | Naïve Matrix Multiplication | dgemm Matrix Multiplication | Cholesky Factorization |
|---|---|---|---|
| | Speedup (Number of points to build using GBBP) | | |
| **X1** | 5.9 (5) | 8.5 (7) | 6.5 (19) |
| **X2** | 5.7 (5) | 5.7 (5) | 15 (8) |

## 4. CONCLUSION

In this paper, we have presented an efficient and practical procedure to build a speed function approximation of a processor. We have demonstrated the efficiency of our procedure by performing experiments with a matrix multiplication application and a Cholesky Factorization application that use memory hierarchy efficiently and a matrix multiplication application that uses memory hierarchy inefficiently on a local network of heterogeneous computers.

## 5. REFERENCES

[1] Anderson, E., Bai, Z., Bischof, C., Demmel, J., Dongarra, J., Croz, J. D., Greenbaum, A., Hammarling, S., McKenney, A., Ostrouchov, S., and Sorensen, D. *LAPACK Users' Guide, Release 3.0*, SIAM, 1999.

[2] Croz, J. D., Dongarra, J., Duff, I. S. and Hammarling, S. A set of level-3 basic linear algebra subprograms. In *ACM Transactions on Mathematical Software*, *16, 1*(1990), 1-17.

[3] Dongarra, J., Whaley, R., C., Petitet, A. *Automated empirical optimizations of software and the atlas project.* Technical Report, Department of Computer Sciences, University of Tennessee, Knoxville, 2000.

[4] Higgins, R. and Lastovetsky, A. Scheduling for Heterogeneous Networks of Computers with Persistent Fluctuation of Load. In Proceedings of ParCo 2005.

[5] Lastovetsky, A and Reddy, R., Data Partitioning with a Realistic Performance Model of Networks of Heterogeneous Computers, In *Proceedings of IPDPS 2004*.

[6] Lastovetsky, A. and Twamley, J. Towards a Realistic Performance Model for Networks of Heterogeneous Computers. In *High Performance Computational Science and Engineering* (IFIP TC5 Workshop), pp.39-58, Springer, 2005.

[7] Spring, J., Spring, N., and Wolski, R. Predicting the CPU Availability of Time-shared Unix Systems on the Computational Grid. In *Cluster Computing, 3, 4*(2000), 293-301.