

## NBS TECHNICAL NOTE 851

U.S. DEPARTMENT OF COMMERCE / National Bureau of Standards

# Computer System Capacity Fundamentals

QC 100 .45753 no.851 1974 C. 2

#### NATIONAL BUREAU OF STANDARDS

The National Bureau of Standards<sup>1</sup> was established by an act of Congress March 3, 1901. The Bureau's overall goal is to strengthen and advance the Nation's science and technology and facilitate their effective application for public benefit. To this end, the Bureau conducts research and provides: (1) a basis for the Nation's physical measurement system, (2) scientific and technological services for industry and government, (3) a technical basis for equity in trade, and (4) technical services to promote public safety. The Bureau consists of the Institute for Basic Standards, the Institute for Materials Research, the Institute for Applied Technology, the Institute for Computer Sciences and Technology, and the Office for Information Programs.

THE INSTITUTE FOR BASIC STANDARDS provides the central basis within the United States of a complete and consistent system of physical measurement; coordinates that system with measurement systems of other nations; and furnishes essential services leading to accurate and uniform physical measurements throughout the Nation's scientific community, industry, and commerce. The Institute consists of a Center for Radiation Research, an Office of Measurement Services and the following divisions:

Applied Mathematics — Electricity — Mechanics — Heat — Optical Physics — Nuclear Sciences<sup>2</sup> — Applied Radiation<sup>2</sup> — Quantum Electronics<sup>3</sup> — Electromagnetics<sup>3</sup> — Time and Frequency<sup>5</sup> — Laboratory Astrophysics<sup>3</sup> — Cryogenics<sup>5</sup>.

THE INSTITUTE FOR MATERIALS RESEARCH conducts materials research leading to improved methods of measurement, standards, and data on the properties of well-characterized materials needed by industry, commerce, educational institutions, and Government; provides advisory and research services to other Government agencies; and develops, produces, and distributes standard reference materials. The Institute consists of the Office of Standard Reference Materials and the following divisions:

Analytical Chemistry — Polymers — Metallurgy — Inorganic Materials — Reactor Radiation — Physical Chemistry.

THE INSTITUTE FOR APPLIED TECHNOLOGY provides technical services to promote the use of available technology and to facilitate technological innovation in industry and Government; cooperates with public and private organizations leading to the development of technological standards (including mandatory safety standards), codes and methods of test; and provides technical advice and services to Government agencies upon request. The Institute consists of a Center for Building Technology and the following divisions and offices:

Engineering and Product Standards — Weights and Measures — Invention and Innovation — Product Evaluation Technology — Electronic Technology — Technical Analysis — Measurement Engineering — Structures, Materials, and Life Safety <sup>4</sup> — Building Environment <sup>4</sup> — Technical Evaluation and Application <sup>4</sup> — Fire Technology.

THE INSTITUTE FOR COMPUTER SCIENCES AND TECHNOLOGY conducts research and provides technical services designed to aid Government agencies in improving cost effectiveness in the conduct of their programs through the selection, acquisition, and effective utilization of automatic data processing equipment; and serves as the principal focus within the executive branch for the development of Federal standards for automatic data processing equipment, techniques, and computer languages. The Institute consists of the following divisions:

Computer Services — Systems and Software — Computer Systems Engineering — Information Technology.

THE OFFICE FOR INFORMATION PROGRAMS promotes optimum dissemination and accessibility of scientific information generated within NBS and other agencies of the Federal Government; promotes the development of the National Standard Reference Data System and a system of information analysis centers dealing with the broader aspects of the National Measurement System; provides appropriate services to ensure that the NBS staff has optimum accessibility to the scientific information of the world. The Office consists of the following organizational units:

Office of Standard Reference Data — Office of Information Activities — Office of Technical Publications — Library — Office of International Relations.

<sup>&</sup>lt;sup>1</sup>Headquarters and Laboratories at Gaithersburg, Maryland, unless otherwise noted; mailing address Washington, D.C. 20234.

<sup>&</sup>lt;sup>2</sup> Part of the Center for Radiation Research. <sup>3</sup> Localed at Boulder, Colorado 80302.

<sup>&</sup>lt;sup>4</sup> Part of the Center for Building Technology.

mai bureau or Standards

1 1974

VC face. 2100 5753 0.851

974

1.2

### **Computer System Capacity Fundamentals**

D. J. Kuck

**Department of Computer Science University of Illinois** Urbana, Illinois

Prepared for the

Systems and Software Division Institute for Computer Sciences and Technology U.S. National Bureau of Standards Washington, D.C. 20234

t. Technical note no. 851



U.S. DEPARTMENT OF COMMERCE, Frederick B. Dent, Secretary NATIONAL BUREAU OF STANDARDS, Richord W. Roberts, Director

Issued October 1974

#### Library of Congress Catalog Card Number: 74-600148

National Bureau of Standards Technical Note 851 Nat. Bur. Stand. (U.S.), Tech. Note 851, 25 pages (Oct. 1974) CODEN: NBTNAE

> U.S. GOVERNMENT PRINTING OFFICE WASHINGTON: 1974

#### FOREWORD

The Institute for Computer Sciences and Technology at the National Bureau of Standards, U.S. Department of Commerce, and the Association for Computing Machinery, the nation's largest technical society for computing professionals, have been jointly sponsoring a series of workshops and action conferences on national issues. These workshops were designed to bring together the best talents in the country in their respective areas to establish a consensus on (1) current state of the art, (2) additional action required, and (3) where the responsibility for such action lies.

An ACM/NBS Workshop Planning Committee on Computer Capacity meeting was held at NBS in February, 1974. One of the discussions resulted in a new theoretical definition of capacity. Those participating in the discussion were: Dr. George Dodd, General Motors Research Laboratory, Chairman; Professor James C. Browne, University of Texas; Mr. Walter M. Carlson, IBM Corporation; Dr. Sidney Fernbach, Lawrence Livermore Laboratory; Dr. Aaron Finerman, State University of New York at Stoneybrook; Dr. Ronald A. Finkler, Institute for Defense Analyses; Professor Michael J. Flynn, The Johns Hopkins University; Mr. Henry S. Forrest, Control Data Corporation; Mr. S. Jeffery, National Bureau of Standards; Professor David Kuck, University of Illinois; Mr. George Lindamood, National Bureau of Standards; Dr. Michael Muntner, General Services Administration; Mr. James Pomerene, IBM Corporation; Mr. Paul Roth, National Bureau of Standards; Mr. Jules Aronson, National Bureau of Standards.

It was felt that the results of the discussion were significant, and that they should be refined and published. Dr. Kuck was kind enough to undertake this task and produce the following paper.

iii

#### TABLE OF CONTENTS

		Page
1.	Introduction	1
2.	Capacity in Overlapped Machines	2
3.	Capacity in Non-Overlapped Machines	12
4.	Processor-Memory-Disk Systems	15
5.	Primary Memory Size vs. $B_d$	17
6.	Conclusion	19

Computer System Capacity Fundamentals

D. J. Kuck

#### Abstract

A framework for the study of computer capacity is given by means of a definition of capacity in terms of speeds of various parts of a computer as well as memory size. The calculation of theoretical capacity is given for several combinations of processor, memory, and I/O bandwidth for both overlapped and nonoverlapped machines. The tradeoff between primary memory size and I/O bandwidth is discussed in terms of the new definition.

Key words: Capacity; computer; evaluation; measurement; performance.

1. Introduction

On mips and mops and megaflops, and binary capacity.

This report is an attempt to outline a formal structure for the study of computer capacity. Several traditional measures will be discussed and some new measures will be introduced. Our goals for the use of measures of computer capacity include:

- 1. Quantification of upper bounds on a given machine's raw theoretical speed for various kinds of computation.
- 2. Comparisons between diverse computer systems for some set of computations.
- 3. Evaluation of the actual performance of a given machine on some job mix compared with its theoretical capacity.
- 4. Guidelines for improving a given system's cost/performance.

Traditionally, people have often quoted computer speeds in mips (millions of instructions per second). But the execution of an "instruction" yields rather different effects on various machines. The range is from some simple indexing operation on a traditional machine to a vector inner product instruction on a modern pipeline processor. Thus, as computer organizations diverged from one another, mops (millions of operations per second) became a more reasonable measure. But in many numerical calculations, floating-point arithmetic operations are the raison d'etre for the computer and logical operations, shifts, etc., are "overhead." Thus, megaflops (millions of floating-point operations per second) may be the important measure. Quoting megaflops is of course quite irrelevant for most computations performed in the real world every day. In many computations, e.g., data base management, file processing, simulation, etc., almost no floating point arithmetic is performed. The primary memory speed and often input/output speeds are the most important  $t_0$  quote in evaluating or comparing machines. Our formultaion will include consideration of the type of computation being performed in terms of ratios such as primary memory to processor bandwidth used by a computation.

We will attempt to bring together in a uniform way measures of the speeds of various parts of a computer as well as memory size. The two main measures which concern us are speed (of processor, primary and secondary memory) and size of primary memory. By definition, speeds are given in units per second and bits/second is the simplest such measure. It is traditional to call the bit rate of a communication channel its capacity. Similarly, sizes of memories in bits may be thought of as capacities. Since we shall be discussing speeds and sizes together, it seems reasonable to refer to the whole notion as "computer capacity".

In addition to the above machine characteristics, our model will include characteristics of the programs being executed. In particular, we are concerned with the fractions of a computation which use each of the three major parts of a system: processor, primary memory and secondary memory. Thus, our model could be used by independently measuring machine and program characteristics, and relating them through the capacity surfaces we derive.

One difficult question is how to deal with the control unit. It has the potential to allow the several major parts of a computer to operate simultaneously and thereby increase capacity in a major way. We shall briefly discuss "serial" control whereby only one function can be performed at a time. Our major attention will be given to computer systems in which the processor, primary and secondary memory all can operate simultaneously in an overlapped way. The models we discuss can be thought of as assuming a perfect "lookahead" control unit. Alternatively, any idleness due to the control unit may be considered to be lumped together with the processor. Degradations in system capacity due to variously constrained control units could be an interesting area for further study. In fact, the control unit could be treated as a fourth dimension in Figure 6 of Section 4.

#### 2. Capacity in Overlapped Machines

In this section we define processor, memory and system capacity. These definitions are given in terms of machine parameters (our  $\alpha$ 's) and program parameters (our  $\beta$ 's). There is a good deal of symmetry in much of the following, and we illustrate this by displaying a number of equations.

Let us consider a clocked machine with a processor, i.e., an arithmetic and logical unit, operating at maximum bandwidth (i.e,, data rate) B bits/second. Let the primary memory bandwith be B<sub>m</sub> bits/second. We define

$$\alpha_{pm} = \frac{B_m}{B} > 0 \quad and \quad \alpha_{mp} = \frac{B_p}{B_m} > 0.$$

For any given computation, the total available bandwidth of the processor or memory may not be used. Thus, we define  $B_p^u \leq B_p$  to be the bandwidth of the processor which is actually used in some computation. Similarly, we define  $B_m^u \leq B_m$  as the used bandwidth of the memory for a given comuptation. Also, for any given computation we define

$$\beta_{pm} = \frac{B_p^u B_m^u}{B_p^u} \ge 1,$$

and

$$\beta_{mp} = \frac{B_m^u + B_p^u}{B_m^u} \ge 1,$$

so it follows that

$$\beta_{pm} - 1 = \frac{B_m^u}{U} \ge 0,$$

and

$$\beta_{mp} - 1 = \frac{B^{u}}{u} \ge 0.$$

We may interpret  $1/\beta_{pm}$  as the fraction of some computation in which the processor is engaged. Similarly,  $\frac{1}{\beta_{pm}} = 1 - \frac{1}{\beta_{pm}}$  is the fraction of a given computation in which the memory is engaged.

If we assume that each memory cycle and each processor operation require the same amount of time, then the above can be interpreted as follows. For a machine with a control unit which overlaps memory and processor operation,  $1/\beta_{pm}$  is the processor fraction of the total number of instructions executed or the processor fraction of the total bandwidth used for some computation. For a machine with a control unit which allows no overlap of processor and memory operation,  $1/\beta_{pm}$  is the processor fraction of the total number of instructions executed. Obviously, similar statements hold for  $1/\beta_{mp}$ .

Next we consider the notion of the capacity of the processor, the memory and the combination of the two. We shall define capacities in bits/second. Since we are interested in maximum possible data rates, we shall assume that either the memory or the processor bandwidth is saturated in any given computation we discuss. Thus, all our discussion of capacity will assume that for the type of computation under consideration no faster data rate is possible on the machine we are considering.

Let us define

$$\gamma_m = \frac{B_m}{B_m^u} \ge 1,$$

and

$$Y_p = \frac{B_p}{\frac{u}{B_p}} \ge 1,$$

which we call the <u>memory freedom</u> and <u>processor freedom</u>, respectively. When  $\gamma_m = 1$ , the computation is said to be <u>memory bound</u> and when  $\gamma_p = 1$ , the computation is said to be <u>processor bound</u>. As outlined in the preceding paragraph, our subsequent discussions of capacity will assume that either  $\gamma_m = 1$  or  $\gamma_p = 1$ , or both.

We can relate machine parameters ( $\alpha$ 's), program parameters ( $\beta$ 's) and freedoms ( $\gamma$ 's) as follows. Since

$$\frac{\mathbf{v}_m}{\mathbf{v}_p} = \frac{B_m}{B_p} \frac{B^{\mathbf{u}}}{B_p} = \alpha \frac{B^{\mathbf{u}}_p}{B_m}$$

and since

$$\frac{\beta_{mp}}{\beta_{pm}} = \frac{B_p^u}{B_p^u}$$

we have

$$\frac{\gamma_m}{\gamma_p} = \alpha_{pm} \frac{\beta_{mp}}{\beta_{pm}} \,.$$

Since  $\alpha_{pm} = 1/\alpha_{mp}$ , we can derive a similar expression by interchanging m's and p's in this equation.

Now we define, for any given computation on any given machine with overlapped processor and memory, the processor capacity

$$C_{p} = \begin{cases} \frac{\alpha_{pm}}{\beta_{pm}-1}B_{p} & \text{if } \alpha_{pm} < \beta_{pm}-1 \\ B_{p} & \text{otherwise.} \end{cases}$$
(1)

Note that  $\alpha_{pm} \geq \beta_{pm} - 1$  is equivalent to

$$\frac{\frac{B_{m}}{B_{p}}}{\frac{B_{p}}{B_{p}}} \geq \frac{\frac{B_{m}^{u}}{B_{p}}}{\frac{B_{m}^{u}}{B_{p}}}$$

so

$$\frac{\frac{B_m}{u}}{B_m} \geq \frac{\frac{B_p}{u}}{\frac{B_m}{p}}$$

or

$$\gamma_m \ge \gamma_p$$
.

But since we are assuming that either  $\gamma_m = 1$  or  $\gamma_p = 1$ , this implies that  $\gamma_p = 1$ . Thus, in the processor bound situation our definition sets  $C_p = B_p$  which is the maximum processor data rate.

On the other hand, if  $\alpha_{pm} \leq \beta_{pm} - 1$ , it follows that  $\gamma_m \leq \gamma_p$ , but since  $\gamma_m = 1$  or  $\gamma_p = 1$  we conclude that  $\gamma_m = 1$ , and we are memory bound. Thus,  $\mathbf{B}_m = \mathbf{B}_m^u$ . Now the definition of  $C_p$  can be rewritten in this case as

$$C_p = \frac{\alpha_p m^B_p}{\beta_{pm} - 1} = \frac{B_m}{B_m^U / B_p^U} = B_p^U \gamma_m \cdot$$

But since  $\gamma_m = 1$ , we have  $C_p = B_p^u$  in the case of a memory bound computation.

Thus, the processor capacity is defined to be the fraction of the processor bandwidth which can be used for this computation, given the

fact that memory bandwidth is saturated.

If we rewrite processor capacity as

$$C_p = \frac{\gamma_m}{\gamma_p} B_p$$

we can interpret it as B<sub>p</sub> if memory freedom is greater than processor freedom for some computations and as B<sub>p</sub> times the ratio of the freedoms otherwise. We emphasize that the processor only reaches its maximum capacity B<sub>p</sub> when  $\gamma_m \geq \gamma_p$ .

We can derive an expression for memory capacity C with analogous characteristics to processor capacity. Thus, we write

$$C_{m} = \begin{cases} \frac{\alpha_{mp} B_{m}}{\beta_{mp} - 1} & \text{if } \alpha_{mp} \leq \beta_{mp} - 1 \\ B_{m} & \text{otherwise.} \end{cases}$$
(2)

Since  $\alpha_{mp}B_m = B_p$ ,  $\beta_{mp} - 1 = \frac{1}{\beta_{pm}-1}$  and  $B_m = \alpha_{pm}B_p$  we can express  $C_m$  in terms of  $B_p$  as follows

$$C_{m} = \begin{cases} {}^{(\beta_{pm}-1)B_{p}} & \text{if } \beta_{pm}-1 \leq \alpha_{pm} \\ \alpha_{pm}B_{p} & \text{otherwise.} \end{cases}$$
(3)

If we define system capacity  $C_{_S}$  to be the total system bandwidth available for any calculation, by properly adding Equations 1 and 3 we obtain

 $C_{s} = \begin{cases} (1 + \frac{1}{\beta_{pm} - 1})^{\alpha} pm^{B} p & \text{if } \alpha_{pm} \leq \beta_{pm} - 1 \\ (1 + \beta_{pm} - 1)^{B} p & \text{otherwise,} \end{cases}$   $C_{s} = \begin{cases} \frac{\alpha_{pm} \beta_{pm}}{\beta_{pm} - 1}^{B} p & \text{if } \alpha_{pm} \leq \beta_{pm} - 1 \\ \beta_{pm} \beta_{p} & \text{otherwise.} \end{cases}$  (4)

so

This can be expressed in terms of B as

$$C_{s} = \begin{cases} \alpha_{mp} \beta_{mp} \\ \beta_{mp} - 1 \\ \beta_{mp} B_{m} \end{cases} \qquad if \ \alpha_{mp} \leq \beta_{mp} - 1 \\ otherwise. \end{cases}$$
(5)

Note that maximum system capacity occurs when both the memory and processor are bound, i.e.,  $\gamma_p = \gamma_m = 1$ . Thus, from Equation 4, if  $\alpha_{pm} = \beta_{pm} - 1$  we have

$$C_{s} = \frac{\alpha_{pm}\beta_{pm}}{\beta_{pm}-1} B_{p} = (1 + \frac{1}{\beta_{pm}-1}) \alpha_{pm}B_{p}$$
$$= (1 + \frac{1}{\alpha_{pm}}) B_{m} = (1 + \frac{B_{p}}{B_{m}}) B_{m} = B_{m} + B_{p}.$$

Thus, the maximum system capacity is the sum of the maximum processor and memory bandwidths.

To make matters concrete, we give in Figure 1 examples of capacities for  $\alpha_{pm} = 2$  and various  $\beta_{pm}$  values. In Figure 1 we denote activity by X and inactivity by 0. We show two columns under the label "memory" to denote that the memory bandwidth is twice the processor bandwidth, i.e.,  $\alpha_{pm} = 2$ . The capacities are shown under the columns of activity. Overall results are plotted in Figure 2.

In Figure 3 we plot system and processor capacity for various values of  $\alpha_{pm}$ . Note that the processor can perform at its maximum capacity over a wider range of problems ( $\beta_{pm}$  Values) for larger  $\alpha_{pm}$ . Note also that the memory capacity which is available for memory to memory (or 1/0) operations becomes greater for larger  $\alpha$ . It should be remarked that as  $\beta_{pm}$  approaches 1, reasonable system performance depends on high frequency of register to register operations (or cache to cache operations).

7

$$\beta_{pm} = 4/3$$

processor	memory
x	X 0
Х	00
X	<b>0</b> 0
X	X 0
X	0 0
X	0 0
•	•
•	٠
•	
B + B/3 =	<u>4B</u> 3

$\beta_{Dm} = 3/2$	β <sub>Dm</sub>	=	3/2
--------------------	-----------------	---	-----

processor	memory
X	X 0
X	00
X	X 0
Х	00
•	•
•	•

$$\beta_{pm} = 2$$

processor	memory
X	X 0
X	X 0
X	X 0
•	•
	•
B + B = 2	B

$$B + \frac{B}{2} = \frac{3E}{2}$$

$$\beta_{DM} = 3$$

processor	memory
X	X X
Х	X X
Х	X X
•	•
•	•
•	•
B + 2B =	3B

Figure 1. Overlapped Processor and Memory,  $\alpha_{pm} = 2$ 

$$\beta_{pm} = 4$$
  $\beta_{pm} = 5$ 

processor	memory	processor	memory
0	<i>x x</i>	0	X X
0	<i>x x</i>	<i>X</i>	<u>x x</u>
x	x x	0	x x
X	<i>x x</i>	<i>x</i>	<i>x x</i>
X	x x	0	x x
X	<i>x x</i>	x	x x
0	X X	•	•
0	x x		•
X	x x		•
X	<i>x x</i>		
X	x x		
X	x x	$\frac{B}{2}$ + 2B	$=\frac{5B}{2}$
•	•	2	2
	•		
$\frac{2B}{3}$ + 2B	$=\frac{8B}{3}$		

Figure 1 (continued). Overlapped Processor and Memory,  $\alpha_{pm} = 2$ 



Figure 2. Capacity for  $\alpha_{pm} = 2$ 



Figure 3. Capacities for Various  $\alpha_{pm}$  Values 11

#### 3. Capacity in Non-Overlapped Machines

To contrast the previous section with a simpler machine and demonstrate how capacities vary as a function of machine organization, we now disallow the simultaneous operation of memory and processor. However, we do assume a perfect lookahead control unit. Figure 4 illustrates the situation for  $\alpha = 2$ .

It may be seen that in the case of non-overlapped processor and memory, we have (using the notation of the previous section):

$$C_{p} = \frac{\alpha_{pm}^{B} p}{\alpha_{pm} + \beta_{pm} - 1}, \qquad (6)$$

$$C_m = \frac{\alpha_{pm} (\beta_{pm}^{-1}) B_p}{\alpha_{pm} + \beta_{pm}^{-1}}$$
(7)

and

$$C_{s} = \frac{\alpha_{pm} \beta_{pm} B_{p}}{\alpha_{pm} + \beta_{pm} - 1}$$
(8)

We plot the capacity of a non-overlapped machine for  $\alpha_{pm} = 2$  in Figure 5. Note the contrast with Figure 2, an overlapped machine. Here the processor and memory capacities only reach their maximum bandwidth at the limits of  $1/\beta_{pm}$ . Note also that a good deal less system capacity is left over for 1/0 activities.

We can easily show that an overlapped machine's capacities are all greater than or equal to a non-overlapped machine's. Thus, from Equations 1 and 6 we see that

$$\begin{array}{c} non-\\ overlapped = \frac{\alpha_{pm}^B p}{\alpha_{pm}^A + \beta_{pm}^A - 1} \leq \begin{pmatrix} \alpha_{pm}^B p \\ pm p \\ \beta_{pm}^B - 1 \\ B_p \end{pmatrix} = overlapped C_p$$

since  $\alpha > 0$  in the first case, and  $\beta pm \ge 1$  in the second case. In similar ways we can show that

non-overlapped 
$$C_m \leq overlapped C_m$$

and

non-overlapped  $C_{s} \leq overlapped C_{s}$  .

 $\beta = 2$ 

processor	memory		processor	memory
0	X X		0	x x
X	0 0		X	0 0
X	0 0		0	X X
0	X X		X	0 0
Х	0 0		0	<i>X X</i>
х	0 0		X	0 0
•	•	•	•	•
•	•		•	
	l .			
2B/3 + 21	$\frac{3}{3} = \frac{4B}{3}$		B/2 + B =	$=\frac{3B}{2}$

 $\beta = 4$ 

 $\beta = 5$ 

processor	memory	processor	memory
0	<i>x x</i>	0	x x
0	<i>x x</i>	0	X X
X	0 0	X	0 0
0	x x	0	X X
x	0 0	0	X X
0	x x	X	0 0
0	X X	•	•
Х	0 0	•	•
0	X X		•
X	0 0		
•	•	$\frac{B}{2} + \frac{4B}{2}$	$=\frac{5B}{2}$
		<u>ر</u> ر	L
•	1.		

 $\frac{2B}{5} + \frac{6B}{5} = \frac{8}{5^B}$ 



Figure 5. Non-overlapped Capacity for  $\alpha = 2$ 

#### 4. Processor-Memory-Disk Systems

Now we turn to a complete system with three components--processor and primary memory as above, together with a secondary memory which we shall refer to as a disk. We shall assume at all times that one of these three components is operating at its highest data rate, i.e., its bandwidth is saturated. We also assume a control unit which overlaps the operation of the processor, the primary memory and the disk. We first give some definitions which are analogous to those of Section 2.

Let B, be the disk or I/O bandwidth.

Then

$$\alpha_{pd} = \frac{B_d}{B_p}$$
,  $\alpha_{md} = \frac{B_d}{B_m}$ 

and

$$\alpha_{dp} = \frac{B_p}{B_d} \, \prime \quad \alpha_{dm} = \frac{B_m}{B_d} \, .$$

We also define

$$\beta_{pd} = \frac{B_p^{u} + B_d^{u}}{B_p^{u}},$$
$$\beta_{md} = \frac{B_m^{u} + B_d^{u}}{B_m^{u}},$$

with  $\beta_{dp}$  and  $\beta_{dm}$  being defined similarly. It follows that processor capacity may be written as:

$$C_{p} = \begin{cases} \alpha_{pm} \stackrel{B}{\xrightarrow{p}} = \frac{B_{m}}{\beta_{pm}^{-1}} & \text{if } \alpha_{pm} \leq \beta_{pm}^{-1} & \text{and } \alpha_{dm} \leq \beta_{dm}^{-1} \\ \\ \alpha_{pd} \stackrel{B}{\xrightarrow{p}} = \frac{B_{d}}{\beta_{pd}^{-1}} & \text{if } \alpha_{pd} \leq \beta_{pd}^{-1} & \text{and } \alpha_{dm} \geq \beta_{dm}^{-1} \\ \\ B_{p} & \text{if } \alpha_{pm} \geq \beta_{pm}^{-1} & \text{and } \alpha_{pd} \geq \beta_{pd}^{-1} & \text{.} \end{cases}$$

Similarly, we have for memory capacity:

$$C_{m} = \begin{cases} \frac{\alpha_{mp} B_{m}}{\beta_{mp} - 1} = \frac{B_{p}}{\beta_{mp} - 1} & \text{if } \alpha_{pm} \ge \beta_{pm} - 1 \text{ and } \alpha_{pd} \ge \beta_{pd} - 1 \\ \frac{\alpha_{md} B_{m}}{\beta_{md} - 1} = \frac{B_{d}}{\beta_{md} - 1} & \text{if } \alpha_{md} \le \beta_{md} - 1 \text{ and } \alpha_{pd} \le \beta_{pd} - 1 \\ B_{m} & \text{if } \alpha_{pm} \le \beta_{pm} - 1 \text{ and } \alpha_{md} \ge \beta_{md} - 1 \end{cases}$$

and for disk capacity:

$$C_{d} = \begin{pmatrix} \alpha_{dm} & B_{d} & = & B_{m} \\ \overline{\beta_{dm}}^{-1} & = & \overline{\beta_{dm}}^{-1} & \text{if } \alpha_{md} \geq \beta_{md}^{-1} \text{ and } \alpha_{pm} \leq \beta_{pm}^{-1} \\ \frac{\alpha_{dp}}{\beta_{dp}}^{-1} & = & \frac{B_{p}}{\beta_{dp}^{-1}} & \text{if } \alpha_{dp} \leq \beta_{dp}^{-1} \text{ and } \alpha_{pn} \geq \beta_{pm}^{-1} \\ B_{d} & \text{if } \alpha_{md} \leq \beta_{md}^{-1} \text{ and } \alpha_{pd}^{-1} \leq \beta_{pd}^{-1} & \text{.} \end{cases}$$

By summing these capacities for consistent conditions, we obtain saturated system capacities as follows:

$$C_{s} = \begin{pmatrix} (1 + \frac{1}{\beta_{mp}-1} + \frac{1}{\beta_{dp}-1}) B_{p} & \text{if } \alpha_{pm} \stackrel{>}{=} \beta_{pm}^{-1} \\ \text{and } \alpha_{dp} \stackrel{<}{=} \beta_{dp}^{-1} \\ (1 + \frac{1}{\beta_{pm}-1} + \frac{1}{\beta_{dm}-1}) B_{m} & \text{if } \alpha_{pm} \stackrel{\leq}{=} \beta_{pm}^{-1} \\ (1 + \frac{1}{\beta_{pd}-1} + \frac{1}{\beta_{md}-1}) B_{d} & \text{if } \alpha_{pd} \stackrel{>}{=} \beta_{pd}^{-1} \\ (1 + \frac{1}{\beta_{pd}-1} + \frac{1}{\beta_{md}-1}) B_{d} & \text{if } \alpha_{pd} \stackrel{\leq}{=} \beta_{pd}^{-1} \\ \text{and } \alpha_{md} \stackrel{\leq}{=} \beta_{md}^{-1} \\ \text{and } \alpha_{md} \stackrel{\leq}{=} \beta_{md}^{-1} \\ \end{pmatrix}$$

It should be noted that in each of these three cases, if the conditions are written as equalities, then the maximum capacity is obtained. In each case this reduces to

$$\max C = B + B + B_{m} + B_{d}.$$

To make matters concrete, in Figure 6 we sketch a surface for  $B_p = B$ ,  $B_m = 2B$ , and  $B_d = \frac{B}{2}$ . The processor capacity is shown as a plateau of height B which runs off to 0 along the memory-disk axis. The top surface is the system capacity. In the region labelled I, the system is processor bound, and in II and III it is memory and disk bound, respectively. Where these three regions meet, the max  $S_c = 3.5B$ point is located.

5. Primary Memory Size vs. B

It is well known that there exists a trade-off between primary memory size and I/O bandwidth. Our purpose here is to sketch an analysis of this trade-off and to relate it to our previous discussion of capacity.

Let the primary memory size be N words of w bits each, for a total of wN bits. The time required to fill this memory from a disk of bandwidth  $B_d$  (assuming  $B_d < B_m$ ) is wN/ $B_d$  sec.

For simplicity, assume a given computation operates on the entire memory. Assume the computation requires  $N^{\alpha}$  time steps. For example, given an nxn matrix, an  $n^3$  step algorithm would give  $\alpha = 3/2$ , since  $n^2 = N$  if the matrix (or a single nxn partition) fills primary memory. Now the time required for the entire computation would be  $wN^{\alpha}/C_{n}$  secs.

17



On the average, the system would be balanced if the processing time were equal to the input time (assuming no output), that is:

or

 $wN^{\alpha}/C_{p} = wNB_{d}$  $N^{\alpha-1} = \frac{C}{B_{d}}$ 

which gives us

$$N = \begin{pmatrix} C \\ D \\ B \\ d \end{pmatrix}^{\frac{1}{\alpha - 1}}$$
(9)

as a relationship between memory size N, I/O bandwidth  $\mathbf{B}_d$  , and processor capacity  $\mathbf{C}_n$  .

The above model can be easily refined in various ways to provide for input and output of data arrays, to provide for multiple buffering, and so on.

#### 6. Conclusion

The point of this report is to provide a framework for the study of computer capacity. We have explored several aspects of the question and Figure 6 shows a system capacity surface as a function of processor, memory and disk bandwidth. For a given class of computations, this surface corresponds to a memory size given by Equation 9 in Section 5.

While we have glossed over many details, the model described here could be useful in the various ways mentioned in the Introduction.

For example, if we were given a set of computations and a machine configuration we could easily determine a Figure 6 type surface from the machine parameters. From the computational algorithms, we could estimate the various  $\beta$  values as discussed in Section 2. This would allow a determination of our operating point in capacity space. While the ideal point is where  $C_s = C_p + C_m + C_d$ , a prudent region is probably somewhere between that point and the processor corner of Figure 6 for "numerical" problems. For "business"-type problems it may be between there and the memory corner of Figure 6. For the class of algorithms under consideration, Equation 9 could be used to make memory size trade-offs.

Given some qualitative idea of the operating rules a user prefers, one could use this model to make quantitative sensitivity studies of capacity as a function of bandwidth and memory size. This could lead to improved system cost/effectiveness. Note that for any given capacity surface, degradation due to operating system overhead, etc., can be quantified by plotting actual performance data in capacity space. In this case, the surfaces shown will serve as theoretical upper bounds on system performance. NBS-114A (REV. 7-73)

U.S. DEPT. OF COMM.				
BIBLIOGRAPHIC DATA SHEET	1. PUBLICATION OR REPORT NO. NBS TN-851	2. Gov't Accession No.	3. Recipient'	s Accession No.
4. TITLE AND SUBTITLE			S. Publication	n Date
			Octob	r 1974
			6. Performing	Organization Code
Computer Sys	stem Capacity Fundamentals		640.00	, organization code
7. AUTHOR(S) David J.	Kuck		8. Performing	g Organ. Report No.
9. PERFORMING ORGANIZAT	ION NAME AND ADDRESS		10. Project/T	ask/Work Unit No.
NATIONAL E	SUREAU OF STANDARDS			
DEPARTMEN	T OF COMMERCE		11. Contract/	Grant No.
WASHINGTON	N, D.C. 20234			
12. Sponsoring Organization Nat	me and Complete Address (Street, City, S	tate, ZIP)	13. Type of R	leport & Period
NBS			- or or or o	
			Final	
			14. Sponsorin	g Agency Code
15. SUPPLEMENTARY NOTES	· · · · · · · · · · · · · · · · · · ·		L.,	
Library of Cong	ress Catalog Number: 74-6	500148		
16 ABSTRACT (A 200 word or	loss factual summary of most sidnificant	information If document	at includes a si	Idnificant
bibliography or literature su	rvey, mention it here.)	intomation. It document	n menudes a si	igniticant
A framewo	ork for the study of compute	er capacity is g.	iven by mea	ans
of a defi	inition of capacity in terms	of speeds of va	arious part	ts
of a com	puter as well as memory size	. The calculat.	ion of	
theoretic	cal capacity is given for se	veral combination	ons of	
	Jur cupuoreg ro grion lor oc	FOILLE GONDELING FE		
processo	r, memory, and I/O bandwidth	for both overla	apped and	
processon nonoverla	r, memory, and I/O bandwidth apped machines. The tradeof	for both overla f between prima	apped and ry memory	
processor nonoverla size and	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed	for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverla size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	for both overla f between prima. in terms of the	apped and ry memory new	
processo nonoverli size and definitio	r, memory, and I/O bandwidth apped machines. The tradeout I/O bandwidth is discussed Dn.	for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverla size and definitio	r, memory, and I/O bandwidth apped machines. The tradeout I/O bandwidth is discussed on.	n for both overla f between prima. in terms of the	apped and ry memory new	
processo nonoverla size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverli size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima. in terms of the	apped and ry memory new	
processo nonoverli size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima. in terms of the	apped and ry memory new	
processo nonoverli size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverli size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverla size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverla size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverla size and definitio	r, memory, and I/O bandwidth apped machines. The tradeof I/O bandwidth is discussed on.	n for both overla f between prima in terms of the	apped and ry memory new	
processo nonoverla size and definitio	entries; alphabetical order: capitalize onion	n for both overla f between prima. in terms of the	apped and ry memory new	unless a proper
processo nonoverla size and definitio 17. KEY WORDS (six to twelve name; separated by semicolo	entries; alphabetical order; capitalize onions)	n for both overla f between prima. in terms of the	apped and ry memory new	unless a proper
processo nonoverla size and definitio 17. KEY WORDS (six to twelve name; separated by semicolo Capacity	entries; alphabetical order; capitalize onions)	n for both overla f between prima. in terms of the ly the first letter of the	apped and ry memory new	unless a proper
processon nonoverla size and definition 17. KEY WORDS (six to twelve name; separated by semicolo Capacity	entries; alphabetical order; capitalize oni ons) ;; computer; evaluation; mea	n for both overla f between prima. in terms of the y the first letter of the surement; perfor	apped and ry memory new first key word to cmance.	unless à proper
processon nonoverla size and definition 17. KEY WORDS (six to twelve name; separated by semicolo Capacity	entries; alphabetical order; capitalize oni ons) ; computer; evaluation; mea	ty the first letter of the surement; perfor	apped and ry memory new first key word to cmance.	unless à proper
processon nonoverla size and definition 17. KEY WORDS (six to twelve name; separated by semicolo Capacity 18. AVALABILITY	entries; alphabetical order; capitalize oni ons) ; computer; evaluation; mea	n for both overla f between prima. in terms of the y the first letter of the surement; perfor [19. SECURIT (THIS RE	apped and ry memory new first key word to cmance.	unless a proper 21. NO. OF PAGES
<ul> <li>processor nonoverla size and definition</li> <li>17. KEY WORDS (six to twelve name; separated by semicolo Capacity</li> <li>18. AVAILABILITY</li> </ul>	entries; alphabetical order; capitalize oni ons) ; computer; evaluation; mea	n for both overla f between prima. in terms of the y the first letter of the surement; perfor [19. SECURIT (THIS RE	apped and ry memory new first key word of cmance. Y CLASS PORT)	unless a proper 21. NO. OF PAGES
processon nonoverla size and definition 17. KEY WORDS (six to twelve name; separated by semicolo Capacity 18. AVAILABILITY For Official Distribution	entries; alphabetical order; capitalize onions) g; computer; evaluation; mea [X Unlimited n. Do Not Release to NTIS	n for both overla f between prima. in terms of the y the first letter of the surement; perfor (THIS RE UNCLASS	apped and ry memory new first key word to cmance. Y CLASS PORT) SIFIED	unless a proper 21. NO. OF PAGES 20
<ul> <li>processor nonoverla size and definition</li> <li>17. KEY WORDS (six to twelve name; separated by semicolo Capacity</li> <li>18. AVAILABILITY         <ul> <li>For Official Distribution</li> <li>X Order From Sup. of Doc</li> </ul> </li> </ul>	entries; alphabetical order; capitalize onions) g; computer; evaluation; mea [x] Unlimited n. Do Not Release to NTIS US Government Priving Office	n for both overla f between prima. in terms of the surement; perfor [19. SECURIT (THIS RE UNCLASS 20. SECURIT	apped and ry memory new first key word of cmance . Y CLASS PORT) SIFIED	unless a proper 21. NO. OF PAGES 20 22. Price
<ul> <li>processon nonoverla size and definition</li> <li>17. KEY WORDS (six to twelve name; separated by semicolo Capacity</li> <li>18. AVAILABILITY <ul> <li>For Official Distribution</li> <li>Tor Official Distribution</li> <li>Order From Sup. of Doc Washington, D.C. 20402</li> </ul> </li> </ul>	entries; alphabetical order; capitalize onlogo g; computer; evaluation; mea [x] Unlimited n. Do Not Release to NTIS y, SD Cat. No. C13. 40:851	n for both overla f between prima. in terms of the surement; perfor [19. SECURIT (THIS RE UNCLASS 20. SECURIT (THIS PA	apped and ry memory new first key word of cmance . Y CLASS PORT) SIFIED TY CLASS AGE	unless a proper 21. NO. OF PAGES 20 22. Price
<ul> <li>processon nonoverla size and definition</li> <li>17. KEY WORDS (six to twelve name; separated by semicolo Capacity</li> <li>18. AVAILABILITY <ul> <li>For Official Distribution</li> <li>For Official Distribution</li> <li>Order From Sup. of Doc. Washington, D.C. 20402</li> <li>Order From National Te</li> </ul> </li> </ul>	entries; alphabetical order; capitalize onions) g; computer; evaluation; mea <u>x</u> Unlimited n. Do Not Release to NTIS ., U.S. Government Printing Office ., SD Cat. No. C13. 40:851 chnical Information Service (NTIS)	n for both overla f between prima in terms of the y the first letter of the surement; perfor [19. SECURIT (THIS RE UNCL ASS 20. SECURIT (THIS P/	apped and ry memory new first key word to cmance. Y CLASS PORT) SIFIED TY CLASS AGE)	unless a proper 21. NO. OF PAGES 20 22. Price 80 cents
<ul> <li>processon nonoverla size and definition</li> <li>17. KEY WORDS (six to twelve name; separated by semicolo Capacity</li> <li>18. AVAILABILITY <ul> <li>For Official Distribution</li> <li>For Official Distribution</li> <li>Order From Sup. of Doc. Washington, D.C. 20402</li> <li>Order From National Te Springfield, Virginia 22</li> </ul> </li> </ul>	entries; alphabetical order; capitalize onlong) g; computer; evaluation; mea [X] Unlimited n. Do Not Release to NTIS ., U.S. Government Printing Office ., SD Cat. No. C13. 46:851 [S] [S]	n for both overla f between prima. in terms of the y the first letter of the surement; perfor [19. SECURIT (THIS RE UNCL ASS 20. SECURIT (THIS P./ UNCLASS	apped and ry memory new first key word of cmance. Y CLASS PORT) SIFIED TY CLASS AGE)	unless a proper 21. NO. OF PAGES 20 22. Price 80 cents

#### NBS TECHNICAL PUBLICATIONS

#### PERIODICALS

JOURNAL OF RESEARCH reports National Bureau of Standards research and development in physics, mathematics, and chemistry. Comprehensive scientific papers give complete details of the work, including laboratory data, experimental procedures, and theoretical and mathematical analyses. Illustrated with photographs, drawings, and charts, Includes listings of other NBS papers as issued.

#### Published in two sections, available separately:

#### • Physics and Chemistry (Section A)

Papers of interest primarily to scientists working in these fields. This section covers a broad range of physical and chemical research, with major emphasis on standards of physical measurement, fundamental constants, and properties of matter. Issued six times a year. Annual subscription: Domestic, \$17.00; Foreign, \$21.25.

#### • Mathematical Sciences (Section B)

Studies and compilations designed mainly for the mathematician and theoretical physicist. Topics in mathematical statistics, theory of experiment design, numerical analysis, theoretical physics and chemistry, logical design and programming of computers and computer systems. Short numerical tables. Issued quarterly, Annual subscription: Domestic, \$9.00; Foreign, \$11.25.

DIMENSIONS/NBS (formerly Technical News Bulletin)—This monthly magazine is published to inform scientists, engineers, businessmen, industry, teachers, students, and consumers of the latest advances in science and technology, with primary emphasis on the work at NBS.

DIMENSIONS/NBS highlights and reviews such issues as energy research, fire protection, building technology, metric conversion, pollution abatement, health and safety, and consumer product performance. In addition, DIMENSIONS/NBS reports the results of Bureau programs in measurement standards and techniques, properties of matter and materials, engineering standards and services, instrumentation, and automatic data processing.

Annual subscription: Domestic, \$6.50; Foreign, \$8.25.

#### NONPERIODICALS

Monographs—Major contributions to the technical literature on various subjects related to the Bureau's scientific and technical activities.

Handbooks—Recommended codes of engineering and industrial practice (including safety codes) developed in cooperation with interested industries, professional organizations, and regulatory bodies.

Special Publications—Include proceedings of high-level national and international conferences sponsored by NBS, precision measurement and calibration volumes, NBS annual reports, and other special publications appropriate to this grouping such as wall charts and bibliographies.

Applied Mathematics Series—Mathematical tables, manuals, and studies of special interest to physicists, engineers, chemists, biologists, mathematicians, computer programmers, and others engaged in scientific and technical work. National Standard Reference Data Series—Provides quantitative data on the physical and chemical properties of materials, compiled from the world's literature and critically evaluated. Developed under a world-wide program coordinated by NBS. Program under authority of National Standard Data Act (Public Law 90-396). See also Section 1.2.3.

Building Science Series—Disseminates technical information developed at the Bureau on building materials, components, systems, and whole structures. The series presents research results, test methods, and performance criteria related to the structural and environmental functions and the durability and safety characteristics of building elements and systems.

Technical Notes—Studies or reports which are complete in themselves but restrictive in their treatment of a subject. Analogous to monographs but not so comprehensive in scope or definitive in treatment of the subject area. Often serve as a vehicle for final reports of work performed at NBS under the sponsorship of other government agencies.

Voluntary Product Standards—Developed under procedures published by the Department of Commerce in Part 10, Title 15, of the Code of Federal Regulations. The purpose of the standards is to establish nationally recognized requirements for products, and to provide all concerned interests with a basis for common understanding of the characteristics of the products. The National Bureau of Standards administers the Voluntary Product Standards program as a supplement to the activities of the private sector standardizing organizations.

Federal Information Processing Standards Publications (FIPS PUBS)—Publications in this series collectively constitute the Federal Information Processing Standards Register. The purpose of the Register is to serve as the official source of information in the Federal Government regarding standards issued by NBS pursuant to the Federal Property and Administrative Services Act of 1949 as amended, Public Law 89-306 (79 Stat. 1127), and as implemented by Executive Order 11717 (38 FR 12315, dated May 11, 1973) and Part 6 of Title 15 CFR (Code of Federal Regulations). FIPS PUBS will include approved Federal information processing standards information of general interest, and a complete index of relevant standards publications.

Consumer Information Series—Practical information, based on NBS research and experience, covering areas of interest to the consumer. Easily understandable language and illustrations provide useful background knowledge for shopping in today's technological marketplace.

NBS Interagency Reports—A special series of interim or final reports on work performed by NBS for outside sponsors (both government and non-government). In general, initial distribution is handled by the sponsor; public distribution is by the National Technical Information Service (Springfield, Va. 22151) in paper copy or microfiche form.

Order NBS publications (except Bibliographic Subscription Services) from: Superintendent of Documents, Government Printing Office, Washington, D.C. 20402.

#### BIBLIOGRAPHIC SUBSCRIPTION SERVICES

The following current-awareness and literature-survey bibliographies are issued periodically by the Bureau:

- Cryogenic Data Center Current Awareness Service (Publications and Reports of Interest in Cryogenics). A literature survey issued weekly. Annual subscription: Domestic, \$20.00; foreign, \$25.00.
- Liquefied Natural Gas. A literature survey issued quarterly. Annual subscription: \$20.00.
- Superconducting Devices and Materials. A literature survey issued quarterly. Annual subscription: \$20.00. Send subscription orders and remittances for the pre-

ceding bibliographic services to the U.S. Department of Commerce, National Technical Information Service, Springfield, Va. 22151.

Electromagnetic Metrology Current Awareness Service (Abstracts of Selected Articles on Measurement Techniques and Standards of Electromagnetic Quantities from D-C to Millimeter-Wave Frequencies). Issued monthly. Annual subscription: \$100.00 (Special rates for multi-subscriptions). Send subscription order and remittance to the Electromagnetic Metrology Information Center, Electromagnetics Division, National Bureau of Standards, Boulder, Colo. 80302. DFFICIAL BUSINESS

Penalty for Private Use, \$300

POSTAGE AND FEES PAID U.S. DEPARTMENT DF COMMERCE COM-215



