Discussed is a capacity planning procedure called USAGE. Various business elements to be individually measured and tracked are presented. Outlined are methods for estimating workload growth. Separate limits of capacity for on-line workload and batch workload demand are discussed. A simple graphic presentation procedure is included to communicate the results of a study to those who need the information for making business decisions.

# A capacity planning methodology

### by J. C. Cooper

The crucial problem for installation managers is that of understanding their application and growth environments sufficiently to gain an understanding of the current workload, to forecast the future workload, and to plan future data processing capacity to meet the needs of the business. Although the CPU is not the only resource that requires planning, it is generally the most critical. The methodology discussed in this paper focuses on capacity planning for the CPU.

To do this, the current workload is broken down into its most important components, which are called business elements. The business plan is also analyzed, and sources of increased workload and new workload are identified. This information is used to quantify future workload requirements for the planning period of interest. Both current and future CPU capacities are quantified and compared with the workload requirements for the planning period. Separate guidelines are used to track the on-line component of the workload and the total workload requirements. This comparison of CPU capacity and workload requirements allows the data processing manager to have an organized view of the growth of the data processing installation, and to plan for the timely acquisition of new hardware to meet the objectives of the business. We call this methodology Understanding Your Application and Growth Environment (USAGE). Like other management methodologies, a USAGE capacity plan can be tracked by measuring the

Copyright 1980 by International Business Machines Corporation. Copying is permitted without payment of royalty provided that (1) each reproduction is done without alteration and (2) the *Journal* reference and IBM copyright notice are included on the first page. The title and abstract may be used without further permission in computer-based and other information-service systems. Permission to *republish* other excerpts should be obtained from the Editor.

future system and comparing those measurements with the forecast. Such a capacity plan can be updated periodically.

The technique was developed in 1974 in Canada, and was based on the experiences of some forty large data processing installations. The technique was brought to the U.S.A. in 1976 for further testing, and the resulting program has been operating since January 1977. It has been widely used to understand current workload and to forecast future workload requirements.

### General concepts

USAGE focuses on the processor resource consumption as recorded in System Measurement Facility (SMF). USAGE is aimed at systems whose operating system is MVT, VS1, SVS, MVS, or MVS/SE for which SMF is available. The input to a study is one month's SMF data. (The USAGE concepts have been applied to other operating systems, such as DOS, although that is not discussed in this paper. The term MVS/SE means MVS/System Extensions.)

A USAGE study begins by analyzing the current system workload from one month's SMF measurements made while the system is operating. Other measures, such as Resource Measurement Facility (RMF) information for long-running jobs, can be used to augment these data.

The second major portion of a study is to forecast the future workload on the system. This requires knowledge of the business plan for growth in the data processing installation. A forecast should show underlying growth of business volumes as reflected in current applications plus planned new applications. There should be an attempt to quantify the latent demand in the system, i.e., the workload demand not presently running due to the current service level, but which will be run when new capacity increases the service level. Thus knowledge of the future business plan is converted into future system workload. This future workload information, when combined with measurements of the current system, allows a forecast of the total future workload to be made. Typically, such a forecast is made every six months.

Each capacity planning study is part of a continuing series of such studies. As with any management tool, feedback from reality is important. In this way, those who measure the existing system and those who forecast future workload become more proficient in these skills.

It is assumed that the system is reasonably well tuned and will continue in that state. It is also assumed that there are no serious bottlenecks other than the CPU, and that DASD devices, control assumptions

units, and channels will be added as required. It is further assumed that SMF is running continuously and recording wait time, elapsed time, and CPU time for each job and TSO user. If MVS is the operating system, Measurement Facility 1 (MFI) or Resource Measurement Facility (RMF) is also installed to record wait times.

auidelines

Guidelines help in all phases of a study, since they provide default values in categories for which real measurements do not yet exist. Guidelines are averages over many installations, although any given installation may deviate from those averages. It is important to understand these deviations from average on the basis of supporting data, so that corrective steps can be taken if the deviations are excessive or unexplained. Guidelines are intended only for a USAGE study, and are not intended to be used out of the USAGE context. There are guidelines for several installation management categories, and part of the USAGE process is to determine whether each category is required or is applicable to the installation being analyzed. Thus the guidelines can be used as reasonableness checks for existing data. If actual data are not available in a category, the guideline value can be used as a default and the USAGE study can be completed. When important guideline parameters are accepted blindly from the default values, the risk is that the values may not represent the data processing installation being studied. A current study even with incomplete information is an impetus to do another study at a later date when more complete information is available.

Guidelines used in subsequent USAGE studies of an installation should be refined to better describe the actual environment. Thus the installation should track either the parameters that are important or those that are in doubt so as to provide feedback that allows measured data to replace the guideline values.

the USAGE day The USAGE day is a twenty-four-hour period. For the purposes of a study the day is divided into homogeneous units, depending on the work being processed. These units are called *production time periods*. There is usually an on-line (or prime) period during which the system is driven in large part by interactive users. Typically this is from 8 o'clock in the morning to 5 o'clock in the afternoon, for a single-time-zone operation. The rest of the day can usually be viewed as batch operation.

The twenty-four-hour day is further analyzed to determine the times the CPU is available and is not available for normal production. Unavailable time is the sum of such things as scheduled preventive maintenance, idle time, and time lost by hardware or software failures. USAGE assumes that the system will be unavailable an average of two hours per day. Thus twenty-two hours are available for normal production, i.e., the upper bound for capacity planning.

The standard month is 20 working days. If the month being analyzed in a particular data processing installation has a different number of weekdays, appropriate guidelines should be altered.

the **USAGE** month

The average maximum CPU utilization for an MVS system expected in a full month is 90 percent of the elapsed time. Total elapsed time in a 20-day month of 22-hour days is 440 hours. The average maximum CPU utilization is 90 percent of the 440 hours, or 396 hours. As a guideline, we have rounded the 396 hours upward to 400 hours per month. The same procedure is used to calculate the average maximum CPU utilization for installations that count a greater number of work days per month.

For the eight-hour on-line period, we assume that there is no preventive maintenance and no stand-alone time. Thus for the prime period there are 8 hours of elapsed time and, at 90% utilization. 7.2 hours of CPU busy time available.

Elapsed time can be measured by using SMF Type 1 or Type 70 records. MVT, MFT, VS1, and SVS produce Type 1 records, and MVS produces Type 70 records. The elapsed time can be segmented into CPU busy time and wait time as reported in the SMF Type 1 and 70 records. The wait time as accumulated by SMF has been validated using a hardware monitor and therefore does not require adjusting. MVS requires that MFI or RMF be running in order to produce Type 70 records. The CPU busy time is then the difference between the measured elapsed time (wall clock time) and the wait time. CPU busy time divided by elapsed time yields the CPU utilization rate or percent CPU busy.

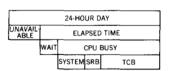
SMF does not account for all the CPU time used in the computer system. Rather, it is intended for job accounting, and as such, it is biased toward repeatability. The price for this repeatability is that all the CPU time spent is not allocated by SMF to each job, which leaves a pool of CPU time called system time. This system time, or uncaptured time, contains CPU time spent for the following types of activities:

- Operator commands
- Job entry and output
- Job scheduler
- TSO control and TCAM
- Supervisor multiprogramming support
- Supervisor virtual storage support
- Other SCP-dependent functions

Task Control Block (TCB) time is allocated to each job by SMF in all systems. In MVS, System Request Block (SRB) time is also allocated to each job and reported by SMF. The balance of the CPU busy time is not directly allocated by SMF to any specific job or task. Figure 1 illustrates the time relationships explained for MVS.

SMF times

Figure 1 USAGE time relationships for MVS



CPU busy time If the total CPU busy time is less than 400 hours the system should be investigated to find the reason. If the system is not used on a round-the-clock or five-days-per-week basis, there is reserve capacity in the off hours. Possibly the system has a very light online load during the day and a heavy batch load in the off-period. Another possible explanation is that the system has availability problems. In this case, a study should be started by the installation to determine the reason for each failure, and to put in place procedures to minimize outage durations and the repetition of known failures.

For the purposes of a USAGE study, the SMF TCB time used is that found in the Type 4 or 5 records for jobs and Type 34 or 35 records for TSO. The total SMF TCB time for the month should sum to about one-half the total CPU busy time. If the total SMF TCB time collected is much less than one-half the total CPU busy time the system should be investigated further. It might be that there are started tasks, such as IMS or CICS, for which SMF is not recording TCB time. (SMF in MVS/Systems Extension 2 measures started tasks.) In this case, special effort must be made to measure or estimate the amount of CPU time consumed by these started tasks. Another possibility is that the system has a large percentage of TSO. The amount of TCB time accumulated by SMF for TSO is considerably less than for the same work accomplished via batch.

It is possible that a job is running with TIME-1440, and the accumulated TCB time is not being recorded for the job by SMF. (SMF in MVS/SE2 measures these jobs.) If this is the case, such jobs should be reviewed to determine whether TIME=1440 is still valid. Yet another reason might be that the CPU is the global processor running Job Entry System 3 (JES3). In this case, the JES3 portion of the work being done by the CPU is not being reported by SMF.

It is possible that the total SMF TCB time is much greater than one-half the total CPU busy time. Such a system is probably a batch environment with a heavy compute-bound workload.

capture ratios

As stated previously in this paper, SMF does not record all the CPU time expended on behalf of a particular job. Its purpose is to be consistent rather than complete, in order to satisfy requirements of charge-back systems. In addition, the portion of time captured by SMF varies with the workload type and the specific SMF implementation for the various SCPs. The *capture ratio* is the proportion of the total CPU time that is captured by the measurement tool. In this case, the capture ratio is the proportion of the total CPU time that is captured as SMF CPU time.

If only one application is running in a system, the capture ratio can be measured. The capture ratio is defined to be the total SMF TCB time divided by the total CPU busy time less the *demand pag-*

Table 1 Capture ratio guidelines

Workload type	Capture ratios				
	MVT/VS1	SVS	MVS TCB	MVS TCB SRB	
	TSO trivial (data entry)	0.30			0.25
TSO program development	0.35	0.30	0.32	0.40	
TSO foreground work/SPF	0.40	0.35	0.37	0.47	
Batch testing	0.50	0.45	0.47	0.60	
Data center operations	0.55	0.50	0.52	0.67	
Commercial production	0.65	0.60	0.62	0.80	
Long scientific or emulation	0.85	0.80	0.82	0.97	

ing time. Demand paging time is the CPU time expended to satisfy page faults. If the SCP's MVS and SRB time is being measured and used, the capture ratio is the TCB time plus SRB time divided by the total CPU busy time less the demand paging time.

The capture ratio varies between approximately 25% and 97% captured. Thus the capture ratio can help allocate uncaptured CPU time to specific applications. Capture ratio guidelines are given in Table 1. These guidelines are averages; actual values for the particular workloads and particular data processing installations should be used where possible.

True CPU hours are the real measure of job CPU utilization. For example, if batch and TSO both have the same amount of SMF hours it is possible for the TSO load to be twice as large in true CPU hours. Assume that commercial production batch has a capture ratio of 0.8 and that TSO program development has a capture ratio of 0.4. Then, for example, if the SMF CPU time for each is 10 hours, the true TSO CPU hours are 25 and the true batch CPU hours are 12.5. Thus TSO has twice as much true CPU time as batch.

The true CPU hours, after adjustment via the capture ratios, more accurately describe the relative CPU consumption of the two types of workload. Consequently, forecasting based on the true CPU times more accurately reflects the CPU resource utilization of a changing workload profile and leads to more realistic conclusions.

The summation of all true CPU time after the capture ratios are applied should approximate the measured total CPU busy time. If the calculated amount is within 5% of the measured amount, the correspondence is good. If it is within 10%, it may be acceptable. If the error is larger than 10%, further analysis is required.

true **CPU** hours

If the calculated total true CPU time is low, SMF may be missing some portion of the workload, or the capture ratios may not adequately reflect the workload characteristics. If the calculated total true CPU time is high, the capture ratios should be investigated. After capture ratios have been chosen, they should not be changed unless there is a very good reason. If the capture ratio is changed it should be carefully noted because changes make it difficult to compare future USAGE studies with past ones.

demand paging CPU Demand paging is the paging the system does to dynamically manage real storage using less real storage than virtual. It does not include functional paging caused by TSO swapping and/or Virtual I/O (VIO). The demand page rate in MVS is the non-swap, non-VIO page rate. RMF measures these categories of paging and can be used to differentiate between them for the MVS user. The demand paging requirement is not only for the path length of the paging mechanism, but also for other overhead items, such as more dispatching and more queue management. For example, the guidelines for a System/370 Model 158 running MVS 3.7 with performance-selectable units is that one page per second requires approximately 0.4% of the CPU.

The CPU cost of demand paging is not included in the capture ratio. Generally, demand paging is a separate line item in the US-AGE study. Demand paging should be measured for each period since it varies by workload mix and normally differs from period to period.

## The current system

production time periods To tailor the USAGE study to the current data processing installation, the first step is to define the production time periods that characterize the installation. Most typically these are an on-line period during the day and a batch period during the evening and night. The weekend is generally a separate period with its own characteristics.

business elements The second step is to define the business elements that characterize the data processing installation. A business element is a grouping of work that is similar from a data processing point of view and can be mapped back into the business in a meaningful manner. Examples of business elements might be major applications such as payroll, the workload submitted by a major department such as engineering, or a major system such as on-line banking. This requires intimate knowledge of the work being done by the system and may require contacting members of user departments. It is important to define enough business elements so that the growth of these categories of work can be traced to some natural portion of the business, yet not so many as to cause

misunderstandings among those who are to be informed about the forecasting process.

The business elements are defined to the program that analyzes the SMF data. This allows all the CPU time consumed by the jobs in that business element to be summed. Thus, another way of looking at business elements is that they consist of work that can be characterized by one capture ratio. In addition, in the forecasting phase, the business element is a portion of the workload that is expected to grow at a common rate.

Business elements are grouped under the following major headings: Production, Testing, and Operations Support. Production includes the batch and on-line jobs that support business operations. Testing represents investment in new applications and is a major source of future growth. Operations Support is the ongoing cost of doing business in a data processing installation.

Certain types of business elements should always be included in the USAGE study. Each major on-line system should be included, even if it is currently small. Typically this portion of the workload is growing and, therefore, should not be overlooked. All testing should be broken out into separate business elements. For example, it is preferable to separate new development from maintenance testing. Operations support should be a separate business element and should include systems programming, backup of data bases, monitoring and measuring, and reruns. The Operations Support normally requires between 8% and 13% of the total CPU time. Values outside this range may indicate that all support functions are not being accumulated properly or that a review is required of all activities in the grouping.

After the data processing installation has been characterized by defining the production time periods and the business elements, the SMF tape can be analyzed. This is usually done by a program that sums the CPU time for each job that is defined to be in a particular business element. This sum is the amount of CPU time captured during the month for that business element. Other sources of system measurement can also be used to characterize the workload. For example, the IMS log tape can be used to characterize the IMS load on the system. The measurements can be summarized on a Measured Data Worksheet such as is shown in Table 2, where P1, P2, and P3 are three production time periods.

Measurement analysis starts by collecting all measured CPU time by business element and by production time period from the SMF tape and other sources of measured CPU time. This is usually done by a computer program. This measured SMF time, however, does not completely characterize the true CPU time. The Current Worksheet, shown in Table 3, can be used to summarize the meameasurements

measurement analysis

CPU/SYSTEM ID		MONTH			
Business elements	Measured CPU time				
	P1	P2	Р3		
Production/Revenue					
On-line					
Time share					
Batch					
Subtotal					
Subtotal			The state of the s		
Testing/Investment					
Development TSO					
Batch		***************************************	/		
Enhancement					
TSO					
Batch		AWA			
Other	7 7778 963 5644				
Subtotal					
Operations Support Development Cost					
Maintenance					
TSO	-				
Batch					
System programming					
Installation support					
Reruns, etc.					
Subtotal			,		
System Summary					
Paging %					
Paging time			The same of the sa		
True CPU time		* No			

surement analysis. It also has space for three production time periods, P1, P2, and P3, and allows the capture ratios (CR) for each business element to change with the production time periods.

Next a first set of capture ratios are assigned to the business elements in order to calculate the true CPU time for each business element by time period. These values are summed for each production time period, and the result is compared with the known total CPU time calculated from the elapsed time and the wait record. If the sum of the calculated true CPU time is within 10% of the total CPU time, the set of capture ratios may be satisfactory. It is the objective of this study to balance the sum of the calculated true CPU times and the total CPU time within 5% or at worst 10%.

The goal of such a study is to limit the risk in decision making. The accuracy of the measurements is valuable only to the extent that it matches the accuracy of the *forecast* information provided

Table 3 Worksheet for the current state of the system

CPU/SYSTEM ID	MONTH						
Business elements	Capture ratio and true CPU time						
	CR	Pl	CR	P2	CR	P3	
Production/Revenue							
On-line							
Time share							
Batch							
Subtotal							
Testing/Investment							
Development							
TSO		477.7					
Batch	-						
Enhancement							
TSO							
Batch							
Other							
Subtotal							
Operations Support/							
Development Cost							
Maintenance							
TSO							
Batch							
System programming							
Installation support							
Reruns, etc.							
Subtotal							
C							
System Summary Paging %							
Paging time		_					
True CPU time							
Utilization %							
Capacity %							

by the user groups. This procedure should be repeated until there is a satisfactory balance between the summation of true CPU time and the total CPU time. This may be done by changing the capture ratios. Any wide deviation from the suggested guideline values should be well understood, since it may be pointing to a loss of data in the measurements. As a result of this analysis it becomes possible to state, in terms of CPU hours per month, how much each business element costs. These data are given in the subtotals in Table 3. When complete, Table 3 represents the present state of the system.

### The forecast

The first step in making a forecast is to understand the business plan of the data processing installation and to translate that into new workload to be done by the system at specific future dates. The forecast workload includes both the growth in the current system workload and any new workload planned or anticipated. A forecast based on a USAGE study is projected two years into the future. A forecast worksheet is a replication of Table 2 for the current state of the system plus each of four half-year periods (at 6, 12, 18, and 24 months) for two production time periods (P1 and P2), since a USAGE study usually generates forecast data at sixmonth intervals. Other intervals can be chosen if necessary. For example, a large IMS application might be going into volume production nine months hence. The first milestone in the study occurs in six months, but the second should be in nine months to reflect major change. Thus, growth is projected in two ways: (1) as a percentage increase in existing CPU hours, and/or (2) as a fixed number of hours added into the business element at some future time.

### the USAGE model

USAGE uses a simple model to forecast future CPU utilization. The growth identified in the forecast portion of the study is applied to the measured CPU hours for each business element, and the result is the amount of CPU time predicted for some future time. Although this is a simple model, it has proved to be a good match for the accuracy of forecast data. When combined with an ongoing tracking procedure to refine the forecast, USAGE has been very useful to many data processing installations in getting them started in capacity planning.

Other models may be appropriate if the installation has exceeded the USAGE guidelines, has known bottlenecks, or has very critical response-time requirements. In such cases, one might try an analytic model<sup>2</sup> or a discrete event simulator, as described in Reference 3. Benchmarking is another approach to modeling for capacity planning purposes.

# current application growth

The first area to investigate for load growth is volume growth (or natural growth) in existing applications due to growth of the enterprise. Volume growth is usually an increase in the number of transactions per unit time. In on-line applications, this growth translates simply into transactions per second. In batch systems, volume growth can be in the number of jobs or in the complexity of each job, depending on how the application is organized. In a USAGE study, this growth in CPU hours per month is usually expressed as a percentage increase for each business element.

# new application growth

It is important to identify each new application that will go into production during the planning period. If the new application is in an early stage of development and little is known about it, a satisfactory estimate of CPU hours can be based on data for a similar application. Other applications that are closer to realization should have better resource consumption information available

as part of the application design, implementation, and test processes. For example, in an on-line application the planning phase should include an estimate of the frequency of occurrence of particular types of transactions and their associated path lengths. This will allow the calculation of the percent of CPU time consumed at some future time.

One technique for estimating future CPU workload requirements is based on CPU time consumed during the development of a new application. USAGE allows the data processing installation to define its categories and to measure the machine time currently used for maintenance, enhancement of existing applications, and development of new applications. Given these measurements, a transfer ratio is needed to allow the projection of future production CPU time from present measured development CPU time. This transfer ratio usually falls in the range of 0.8 to 1.2. A ratio of 1 implies that for every hour of development work measured on the CPU today there will be one hour of additional production work one year hence. Enhancement and new development produce equal amounts of new production CPU time in the future. That is, the same transfer ratio usually holds for both. An exception to this range of transfer ratios is in the case of conversion of an existing program where the transfer ratio is approximately 2. These guidelines for the division of work type in the application development installation and the transfer ratio should be tailored to the individual data processing installation.

Another source of load growth is new services that enhance the way of doing business. A new service is a package that the data processing installation can purchase and, with a minimum effort by the programming staff, deliver a new service to the end users. Examples are personal computing for the professional user, such as A Programming Language (APL) and text processing. These new services are a source of load growth that is not based on the number of programmers in the development installation, but can be based on the number of end users.

If a professional is to use the system, he must use it frequently enough to maintain his skill. Otherwise he spends too much time relearning the procedure necessary to access the system and little time doing professional work. Other criteria for successful use of a system by a professional are high availability and consistently good performance. Availability starts with adequate access to a terminal. Most needs can be met where there is one terminal for each five professionals. As a guideline, each professional consumes 0.2 hours of System/370 Model 168 CPU time or 0.5 hours of Model 158 CPU time per month. If the professionals are using more time than this, they probably need more terminals. If it is projected that the professional uses substantially less CPU time than that indicated here, the actual use may approach zero due to

new services

Table 4 Transaction volumes in an unconstrained environment per eight-hour day

Application	Users/ terminal	Transactions/ day
Data entry	1	1000
Inquiry	1	300
Production	1	500
Programmer	2.5-3	1200
Personal computing	5	500

lack of skill in using the facility. On the other hand, if the use represents a real need, the CPU time required will climb to the indicated figure.

### latent demand

Latent demand is work that is not being submitted to the system due to some constraint. The constraint may be one of policy, e.g., a certain class of work cannot be processed on first shift, or the constraint may be due to poor response time or to lack of an adequate number of terminals in an interactive situation. In any case, if the constraint is removed, additional work which is not currently being processed will appear in the system. Although there is no exact way of calculating latent demand, Table 4 allows an estimate to be made. Measurements are made of the path length of the average transaction for a class of work and the current transaction rate of the application per terminal, and the total CPU consumption for the month is calculated. Note that if this category of work corresponds to a business element, the CPU hours can be gotten directly from the USAGE study. This amount of CPU time can be compared to a theoretical amount of CPU time using the transaction volume per day in Table 4. The difference in CPU hours represents an estimate of the latent demand expected to appear when the constraint is removed.

For example, if a data entry transaction takes 200 ms of CPU time on average, and there are 750 transactions per day, in a 20-day month each terminal consumes 0.83 hours of CPU time. In addition, if it is felt that the response time for data entry is not as good as it could be, and a larger CPU is added to improve the response time, the number of transactions per day per terminal may go to 1000, as indicated in Table 4. If all other factors remain the same, and the data entry work is available, each terminal is estimated to consume 1.1 CPU hours in a 20-day month. Thus the estimated latent demand is the difference, or 0.27 CPU hours per terminal per month. See Reference 4 for a further discussion of latent demand.

### forecast worksheet

Probably the major result of the USAGE study is that the capacity planner has a much better understanding of the data processing installation, both from a current point of view and as it may look in the future. The projection of the future data processing installation is characterized on a Forecast Worksheet, which, as was mentioned earlier in this paper, is a five-fold replication of Table 2. This is a form for quantifying the capacity and the requirements of the data processing installation over the next two years at sixmonth intervals. With the Forecast Worksheet information, there is sufficient information to form a basis for a capacity plan.

As a result of the forecast analysis, it is possible to estimate the cost in terms of CPU hours per month of the business elements by production time period at the future time period of interest. After the workload has been forecast, a new system to meet this requirement may be characterized. Then the capacity of the future system versus the future workload requirements can be calculated.

In the USAGE study, we are tracking the total CPU hours by business element and production time period for an entire month. Generally, we assume a guideline of a maximum sustained available capacity for an MVS system to be 90%. Thus, for planning purposes, we choose to limit the average CPU utilization during an 8-hour prime production time period to 7.2 hours of CPU time. That is the limit of available sustained capacity for the total workload. In the off-prime production time period there are 16 hours less the 2 hours of unavailable time (or 14 hours) at 90%, which yields 12.6 hours off-prime CPU time per day of available sustained capacity for the total workload.

On-line workload—such as Information Management System (IMS), Customer Information Control System (CICS), or Time Sharing Option (TSO)—should be tracked separately, since it has the quality of instantaneous demand for service, which is measured by response time.

Normally the work of an on-line system has peaks and valleys in the demand for service that follow the habits of the users. Typically, there is a peak around ten o'clock in the morning, a valley at lunch time, and another peak around two o'clock in the afternoon. This demand profile may be modified by the number of geographical time zones covered and the type of work being done by the system. If three time zones are covered, the variations in workload may be less extreme, since lunch time occurs over a three-hour period. If the system is controlling a steel rolling mill, the demand profile will be very flat day and night, twenty-four hours a day, seven days a week. If the workload is a programming installation that has enough terminals and capacity, the peaks and valleys may be pronounced.

Although we are tracking average CPU time per month, the system must provide for peak demand of on-line services. Thus the aver-

total sustained capacity

demand profile

age amount of CPU available to some applications has to be limited for planning purposes in order to meet the need of the peak period.

peak-toaverage ratio We account for peak service requirements with the concept of peak-to-average ratio. This can be measured by finding the online consumption of CPU time for the peak half hour of the peak day of the month. This is generally 1.5 times the average CPU consumption for the on-line applications on the peak day. Then find the ratio of the average CPU consumption on an average day of the month to the average CPU consumption of the peak day. This ratio is also about 1.5. To find the peak-to-average ratio for the entire month, the two measured ratios are multiplied, which yields 2.25. That is, the expected CPU consumption during the peak half hour is 2.25 times the average for the month for the online applications.

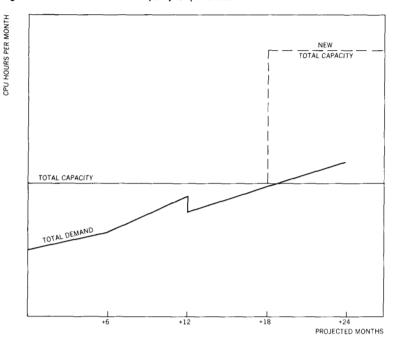
For an overall capacity limit of 90%, and if the on-line applications during the peak half hour of the month use the entire 90%, the average utilization for the month is the peak utilization divided by the peak-to-average ratio. In the example, the peak utilization of 90% is divided by 2.25 to yield 40%. This is the maximum average utilization for all on-line applications. If this maximum is exceeded, one can expect longer response times in on-line applications. It is key to the study that the on-line applications be tracked against this capacity limit as a separate item.

Thus, there are separate limits for the total CPU demand in a production time period and the on-line demand for that period. In a USAGE study, both of these limits are used in the forecasting phase, and the workload demand for each category is tracked against them. This is the basis for a trend toward systems running out of capacity in the on-line category before the total capacity is exhausted in the prime shift.

### **Projecting system configurations**

For estimating future systems represented by CPUs not yet installed, we use a table of Average Relative Internal Performance (ARIP) for various CPUs. (Reference 5 has a table of such relative power factors for CICS.) This is not a MIPS (Millions of Instructions Per Second) value, since, with the advent of firmware assists such as MVS/SE, it is possible for the measured instruction rate to decrease while the relative internal performance increases. For projecting future systems on the basis of a USAGE study, a demand versus capacity plot, such as is shown in Figure 2, is used. The figure shows the total demand and the total processor capacity over a projected two years. The jog in the demand curve indicates a change when new software is introduced that

Figure 2 Total demand versus capacity for prime shift



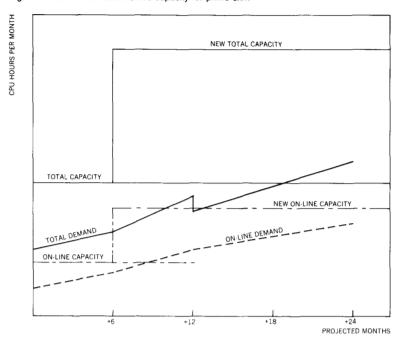
can lead to a lower total demand. Such a drop in demand could also be caused by the migration of some portion of the workload away from this system or out of this production time period.

The dashed line indicates new capacity added to the system in time to provide a capacity upgrade before total demand exceeds total capacity. As previously discussed, total demand is not the only consideration in a capacity planning study. The on-line workload demand should also be tracked independently. In Figure 3, on-line demand and on-line capacity are compared with total demand and total capacity. On this basis, the processor upgrade must be planned much earlier so that on-line demand can be accommodated by new on-line capacity. Clearly, in conjunction with such an analysis, there must be a financial analysis to put a business value on the various possible alternatives of processor upgrades and on the timing of those upgrades.

## Concluding remarks

USAGE is an orderly procedure that starts by measuring an existing system. Like any effective management process, this should be an ongoing process. Projections should be tracked and reworked periodically. The results of a study are useful for planning

Figure 3 On-line demand versus capacity for prime shift



purposes, and as in any good management procedure, projections should be checked against reality. In capacity planning with USAGE, the results should be tracked against measured data. Discrepancies should be analyzed so that later projections can be made more accurately. The information should be fed back to the personnel who forecast growth rates, so that their understanding of the relationship between their portion of the business and the data processing installation can be updated and their future forecasts can be more accurate.

In summary, USAGE is a good place to start in the capacity planning of an installation. It allows the various business elements to be individually measured and tracked. It outlines methods of estimating the workload growth, and sets separate limits of available capacity for on-line workload and batch workload. The simple graphical presentation procedure facilitates communication of results of a study to those who need them for planning purposes.

### ACKNOWLEDGMENTS

USAGE was originally developed in Canada by Prem C. Agrawal and Edward C. Turgeon. The original work in developing the procedures and the guidelines has proved remarkably durable in a changing environment. I wish to acknowledge their assistance in supporting the USAGE program in the United States and their assistance with the original documentation.

### CITED REFERENCES

- R. M. Schardt, "An MVS tuning approach," IBM Systems Journal 19, No. 1, 102-119 (1980, this issue).
- 2. D. C. Schiller, "System capacity and performance evaluation," *IBM Systems Journal* 19, No. 1, 46-67 (1980, this issue).
- 3. H. M. Stewart, "Performance analysis of complex communications systems," *IBM Systems Journal* 18, No. 3, 356-373; and H. C. Nguyen, A. Ockene, R. Revell, and W. J. Skwish, "The role of detailed simulation in capacity planning," *IBM Systems Journal* 19, No. 1, 81-101 (1980, this issue).
- 4. L. Bronner, "Overview of the capacity planning process for production data processing," IBM Systems Journal 19, No. 1, 4-27 (1980, this issue). See also L. Bronner, Capacity Planning Implementation, order number GG22-9015 (January, 1979) available through IBM branch offices, and L. Bronner, Capacity Planning, An Introduction, order number GG22-9001 (January, 1977) available through IBM branch offices.
- Customer Information Control System/Virtual Storage (CICS/VS), System
  Programmer's Reference Manual, order number SC30-0069, available through
  IBM branch offices.

The author is an advisory market support representative at the IBM Washington Systems Center, 18100 Frederick Pike, Gaithersburg, MD 20760.

Reprint Order No. G321-5114.