

Research Report

Data Centers in the Wild: A Large Performance Study

R. Birke*, L. Y. Chen*, E. Smirni[#]

*IBM Research – Zurich
8803 Rüschlikon
Switzerland

[#]College of William and Mary
Virginia, USA

LIMITED DISTRIBUTION NOTICE

This report will be distributed outside of IBM up to one year after the IBM publication date.
Some reports are available at <http://domino.watson.ibm.com/library/Cyberdig.nsf/home>.



Research
Almaden • Austin • Brazil • Cambridge • China • Haifa • India • Tokyo • Watson • Zurich

Data Centers in the Wild: A Large Performance Study

Robert Birke
IBM Research Zurich Lab
Rueschlikon, Switzerland
bir@zurich.ibm.com

Lydia Y. Chen
IBM Research Zurich Lab
Rueschlikon, Switzerland
yic@zurich.ibm.com

Evgenia Smirni
College of William and Mary
Virginia, USA
esmirni@cs.wm.edu

ABSTRACT

With the advancement of virtualization technologies and the benefit of economies of scale, industries are seeking scalable IT solutions, such as data centers hosted either in-house or by a third party. In spite of the ubiquity of data centers, little is known about their in-production performance. This study fills this gap by conducting a large scale performance survey on several thousands of data center servers within a time period that spans two years. We provide in-depth analysis on the diversity and time evolution of existing data centers by statistically characterizing typical data center server workloads, highlighting similarities and differences in the usage of basic resource components, including CPU, memory, disk, and file system. In addition, we quantify the time variability and seasonality of resource demands and how they are changing according to different geographical locations as well as virtual and physical operating systems. This survey provides a baseline for workload calibration, which is critical for the development of scalable and efficient resource management and capacity planning of future data centers.

1. INTRODUCTION

Data centers are becoming the standard IT solution to a host of businesses due to their great potential in reducing operation costs and management overheads. Data centers have the clear advantage of economies of scale from the perspective of abundant deployment of resources. Further powered by virtualization technologies [23, 24, 37, 38], they enable multiple resources being multiplexed/shared among a large number of users with diverse time-varying access patterns [4, 11]. Effectively managing a large amount of resources under highly dynamic workloads is no mean feat.

Data center management can be divided into two categories: (i) resource management, that focuses on dynamically controlling workloads given a resource pool, and (ii) capacity planning, that focuses on resource provisioning. In recent years, a large number of studies [7, 14, 24] aim at leveraging virtualization technologies to improve data center efficiency by consolidating workloads. The evaluation of existing consolidation strategies relies heavily on using "representative" workloads, which can be further adopted in simulation, prototypes, and even modeling studies of data centers. Capacity planning on the other hand often relies on time series methodologies [12, 14, 20] for workload forecasting in order to dimension resource capacity. Evolution of workload resource demands is very crucial for short, medium, and long term capacity planning.

There have been large scale studies in the literature that focus on specific workloads or specific resources [4, 5, 28, 30, 35] providing significant insights, e.g., the discovery of self-similarity in web traffic, and disk and DRAM failure patterns. In addition, most detailed workload analysis encompassing major server resources

[22, 25], e.g., CPU, memory, and disk, are based on executing various benchmark suites on different architectures. To the best of our knowledge little is known about "real" workloads placing demands on data centers and how their combined demands on different resources evolve across time.

Motivated by the gap between the needs of data center management and the lack of related studies, we conduct a detailed performance survey across a random selection of *several thousands of servers* hosted in different data centers, located in five continents, during the period of June 2009 to May 2011. We collect the utilization values of CPU, memory, and disk in different time scales, i.e., hourly, daily, weekly, and monthly. We first aim at characterizing the consolidated server workloads in terms of resource demands. In particular, we provide statistics that quantify demands across different resources and derive implications on the potential server consolidation strategies. We also characterize the time evolution of resource demands with respect to geographical locations, time shifts and seasonality, resource interference, and virtual vs. physical machines.

Due to the nature of data that are available to us, our study suffers from some unfortunate limitations. First, because of the large scale and complexity of data, there is no availability of performance metrics in fine time granularities. The smallest resolution we present here is per minute performance but this information is available only for a short period of observation times. For most of analysis presented in this study, we focus on a coarser time granularities, i.e., day, week and month, which are obtained by aggregating performance metrics collected every 15 minutes. A second limitation is that the main metric that is available (either directly or after trace manipulation) is the utilization of different resources, i.e., CPU, memory, disk and file system. CPU utilization describes the amount of time CPU is actively used; whereas we derive the space utilization of memory, disk, and file systems from raw data. Our study sheds no light in the arrival patterns or possible burstiness of requests into the data center(s). In addition, there are no data on the distribution of response times of the various data center users or specific applications.

Despite these limitations, the data presented here can well represent the performance characteristics in today's large scale data centers. Our survey and analysis provide the essential information, i.e., workloads and trends of resource utilizations and usage patterns. Such information can be directly used for resource provisioning and capacity planning or can be used indirectly to guide the parameterizations of simulation studies such that the workloads and scenarios are realistic and relevant. Due to the nature of our data, we are able to present information in different levels of granularity, i.e., across all data centers, across data centers on specific continents, countries, industries, and enterprises.

The summary of our contributions and outline of this work are as follows. In Section 2, we provide an overview of a typical server workload in data centers. In Sections 3 and 4, we analyze the evolution of resource supplies and demands. In particular:

- We present overall characteristics of main resources, i.e., mean, standard deviation, and empirical distribution of CPU, memory, disk, and file system utilization, across all servers (Section 3.1).
- We present the correlations of monthly resource utilizations. We found that CPU, memory, and disk are moderately and positively correlated, and disk and memory are negatively correlated (Section 3.2).
- We present the distribution of CPU utilization across servers. Using elementary fitting, we conclude that surprisingly, a simple truncated exponential distribution can well capture the CPU utilization (Section 3.3).
- We present the maximum and median memory paging rate [KB/s] across all servers. We found that roughly 20 percent of servers experience a noticeable paging rate for a short period of time and that the majority of servers have negligible paging activities for most of the time (Section 3.4).

2. DATA STATISTICS

We collect resource utilization statistics from several thousands of heterogeneous servers at in-production data centers from June 2009 till May 2011. The geographic distribution covers all continents and a wide range of countries, e.g., developed and developing. These systems are used by different industries, including banking, pharmaceutical, IT, consulting, and retail, and are based on various UNIX-like OSs, i.e., AIX, HP-UX, Linux, and Solaris. In summary, our collected samples contain a huge number of representative server statistics reflecting the current practice of resource management in data centers.

In particular, we collect resource utilization values, e.g. the CPU and disk, for their multiple implications of current data center performance. Intuitively, resource utilization values show the *efficient* use of resources in data center servers. Secondly, resource utilization values can be viewed as the *work load* on the different resources. Finally, the trend of utilization values indicate the growth of *demand* and *supply* of resources in data centers. Since utilization values are normalized, it eases comparisons across different servers. Moreover, the utilization values are used as workload inputs to evaluate various power management and capacity planings in data centers [10, 24, 29]

We focus on three main physical resources per server: CPU, memory, and disk, plus the file system. A server can have multiple disks and the disk utilization is the total over all the attached disks. The file system includes both local and remote data storage, which can be on the media of disks and memory. Similar to the disk, when there are multiple file systems, the utilization is taken from the total of all file systems. The CPU utilization is defined by the percentage of time the CPU is active over an observation period (we use 1 and 15 minutes as base periods); whereas utilization values of memory, disk, and file system are defined by the volume usage, i.e., used space divided by the total available space. The average utilization values over base periods are collected via prevailing `vmstat`, `iostat` and `df` utilities (or other equivalent tools) and stored in a round-robin like database so that the space footprint of the database (per server) is constant. Essentially, older data points are aggregated over larger time periods, i.e., days, weeks

and months. Recent data is available at a higher time resolution than older data. In particular we consider monthly data from June 2009 to May 2011, weekly data from June 2010 to May 2011, and daily data for May 2011.

The base utilization values are further divided by different shifts of a day and of a week, i.e., prime working time (8 AM - 5 PM) vs. off-prime working time, and week days vs. weekends. Note that all selected servers have all the data points in the observation period in order to consider a constant server population. As such, we can rule out the explanations of presented data, due to the fluctuation of adding or removing servers. Most of our analysis here is based on the monthly average resource utilizations over a two years span, computed from data points of a 15 minutes long base period. We also present daily and weekly in some of the subsections, based on 15 minutes base values.

We present an overview of the data set in Table 1. This table summarizes the statistics of servers by different categories, i.e., countries, physical/virtual, and OS. Due to lack of space, we only present detailed data for five countries rather than all. We anonymize and refer to those as country A, B, C, D, and E. Note that this subset however still covers different geographical areas and developed vs. developing areas.

We adopt two perspectives to analyze the data set:

- **Diversity of server workloads:** here we aim to study and quantify the diversity of server workloads, i.e., how server workloads at *different* data centers vary. We capture the statistics of resource demands of servers, i.e., distributions of average resource utilizations and workload variations, across all servers, and over the entire observation period. In addition to in-depth analysis on each resource type, we quantify the dependency among resources. Future data centers can leverage this analysis as a base line of workloads for their system design and resource management, via mathematical models, simulation, emulation, prototyping, and implementation.
- **Time evolution:** here we aim to explore the time series of data center resource demands over the past two years, by adopting a *temporal* perspective. We present time series of resource utilization at different time scales and categories. To identify the representative time series model for predicting resource demands, we analyze the distribution of the auto-correlation function (ACF) in time units of months and days. We capture the time variability of CPU workloads at a minute granularity. In particular, we assess the long term capacity planning for data center resources.

To ease readability, we first introduce some notation used in this study. Let $U_{i,j}(t)$ denote the utilization of resource i at server j during time t . We use the convention that i is subscript for resources, $i = \{c, m, d, f\}$, j is the subscript for server, $j = \{1 \dots J\}$, and t is the index for the time window defined in units of months, weeks, or days, $t = \{1 \dots T\}$. We specifically compute the following statistics:

- $\overline{U_{i,j}}$: the mean utilization of resource i at server j over all time windows t .
- μ_i : the mean utilization of resource i over all servers, i.e., the mean of all $\overline{U_{i,j}}$.
- $\mu_i(t)$: the average utilization of resource i at time window t for all servers.

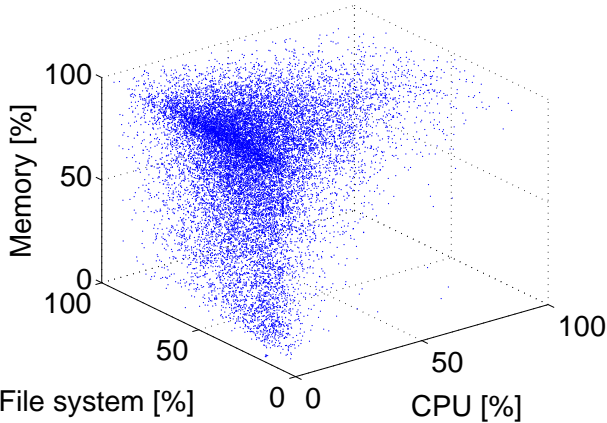
Table 1: Overview of resource utilization by different categories

All	CPU [%]			Memory [%]			Disk [%]			File system [%]		
	mean	std	CoV	mean	std	CoV	mean	std	CoV	mean	std	CoV
All	17.76	18.16	1.02	77.93	23.18	0.30	75.28	24.34	0.32	45.17	19.27	0.43

Country	CPU [%]			Memory [%]			Disk [%]			File system [%]		
	mean	std	CoV	mean	std	CoV	mean	std	CoV	mean	std	CoV
Country A	24.91	18.30	0.73	80.83	16.09	0.20	81.15	18.37	0.23	55.04	18.35	0.33
Country B	12.89	12.00	0.93	81.73	15.57	0.19	64.10	20.08	0.31	46.33	19.50	0.42
Country C	7.25	9.76	1.35	71.87	25.09	0.35	63.59	27.04	0.43	38.89	18.72	0.48
Country D	14.55	12.93	0.89	84.44	19.02	0.23	70.20	23.98	0.34	48.21	17.95	0.37
Country E	19.28	19.45	1.01	78.57	22.22	0.28	72.33	24.96	0.35	44.13	18.33	0.42

Operating System	CPU [%]			Memory [%]			Disk [%]			File system [%]		
	mean	std	CoV	mean	std	CoV	mean	std	CoV	mean	std	CoV
AIX	21.05	19.20	0.91	84.48	16.80	0.20	69.37	23.16	0.33	47.25	18.32	0.39
HP-UX	19.34	16.23	0.84	68.16	20.53	0.30	81.49	17.15	0.21	56.15	16.51	0.29
Linux	6.76	9.73	1.44	80.99	21.29	0.26	85.54	27.17	0.32	36.62	20.24	0.55
Solaris	10.13	11.45	1.13	44.52	23.56	0.53	95.10	12.20	0.13	39.53	20.03	0.51

Type	CPU [%]			Memory [%]			Disk [%]			File system [%]		
	mean	std	CoV	mean	std	CoV	mean	std	CoV	mean	std	CoV
Server	12.15	13.87	1.14	71.93	25.96	0.36	78.66	25.77	0.33	44.91	20.47	0.46
LPAR	49.01	20.10	0.85	85.18	16.23	0.19	71.21	22.17	0.31	45.96	17.76	0.39
Solaris Zones	6.20	8.78	1.42	36.65	20.80	0.57	98.42	7.33	0.07	26.63	16.92	0.64

**Figure 1: Mean CPU, memory, and file system utilization of all servers ($\bar{U}_{i,j}$).**

3. TYPICAL AND REPRESENTATIVE DATA CENTER WORKLOADS

To characterize the data center workloads on all resources, we start from simple visual analysis. We present the mean CPU, memory, and file system utilization values $\bar{U}_{i,j}$ of all servers in Figure 1. Clearly, servers appear very diverse in terms of resource utilizations, although some clustering behavior is observed. As indicated by the large number of points along the north west edge in Fig. 1, there is a significant number of systems with memory that is more than 80% utilized.

We list in Table 1 the mean, standard deviation, and coefficient of variation (CoV) of the mean resource utilizations $\bar{U}_{i,j}$ over all servers by different categories. The average values indicate the

typical resource workloads running on today’s servers: CPU, memory, disk, and file system have mean utilization values at 17.76%, 77.93%, 75.28%, and 45.17% respectively. We observe that all resources have relatively low standard deviation compared to the mean value as reflected also in the CoV. This can be partially explained by the bounded utilization values, i.e., 0–100%. In general the degree of diversity in server resource utilizations is rather low, especially for the categorized servers. We combine the in-depth analysis of the diversity and time evolution of servers for different categories in Section 4.

Clearly, the average CPU utilization is low for most. However one has to keep in mind that this data is averaged over the entire observation period and therefore CPU utilization peaks are hidden. In contrast to the CPU, disk and memory are both highly utilized. The high disk utilization can be easily explained by the definition, i.e., usage of allocated disk spaces, while high memory values by the fact that some OS, e.g., Linux and AIX, use memory as cache to speed up IO operations. We provide detailed statistics in Section 4.5. The file system is moderately utilized leaving space to store new data. Due to the fact that CPU utilization is defined by the percentage of "active" time, one can expect CPU utilization of a server to be the most varying, compared to other resources. Consequently, the CPU has the highest CoV values (~ 1), while the most homogeneous resources in terms of usage are memory and disk (CoV ~ 0.3) with the file system somewhere in between (CoV ~ 0.4). In the following subsections, we further focus on CPU variability, presenting results using finer time granularities and/or on smaller sets of servers.

In summary, statistics listed in Table 1 provide an overview of the server workload at data centers. Essentially, a typical server workload has rather low and varying CPU utilization, medium loaded file systems, and very highly utilized memory and disk. Such information could be used as the basis for designing resource management strategies and workload benchmarks.

3.1 Resource Utilization Distributions

To give more in depth information, we use first and second moment analysis and obtain an empirical distribution. We plot the cumulative density function (CDF) of mean utilizations of each resource, μ^i , over all servers in Fig. 2. Additionally, we fit the empirical distributions with known statistical distributions:

- **CPU** can be approximated with a truncated exponential distribution with $\mu = 17.64$.
- **Memory** can be approximated with a generalized extreme value distribution with $k = -1.14$, $\sigma = 27.85$, and $\mu = 75.66$.
- **Disk** can also be approximated by a generalized extreme value distribution with $k = -1.22$, $\sigma = 25.83$, and $\mu = 78.87$.
- **File system** fits quite well to a Weibull distribution with $a = 50.74$ and $b = 2.51$.

Overall, we observe that it is nearly trivial to find a statistical distribution that fits the empirical distribution of CPU and file system with negligible errors. Fitting disk and file system data into distributions results is more challenging.

Results in Figure 2 can be used for workload verification of large scale studies, such as those that focus on cloud data centers. One can use these results to calibrate simulation, prototyping, and benchmarking. To further clarify how representative the empirical distributions across all servers are, we further analyze resource usage variability across servers that serve five different enterprises.

We select five enterprises, anonymized by Ent 1-5, each of which are in different industries, e.g., banking, communication, e-business, production, and telecommunication, and use more than 100 servers each. Furthermore, we use finer-grained statistics, i.e., $\overline{U}_{i,j}$ is computed from the weekly average over a year span. Our aim here is to validate the fitted resource distributions derived from a much larger server population, when applying to a small number of servers that serve a specific enterprise. In Figure 3 (a)-(d), we summarize the empirical CDFs of those selected enterprise. For all enterprises, one can see that CDFs of memory and file system are very close to CDFs computed from all servers, i.e., the fitted distributions close to the empirical values.¹ The variability of server CPU utilization in enterprises 2-4 can be well captured by exponential distributions; whereas the empirical CDF of enterprise 1 and 5 can be fitted by a Weibull distribution. One can easily see that average CPU loads on enterprise 1 and 5 are much higher.

3.2 Correlation Among Server Resources

We compute the correlation coefficient between different resources, based on monthly utilization, and summarize results in Fig. 4. In general all resources are correlated with absolute coefficient values between 0.26 and 0.1 for the monthly data. As the highest correlations are between CPU/memory and CPU/file system, one expects them to depend on each other, i.e., resource demands have similar peaks and dips. It is surprising to observe that the file system is more correlated with the memory than with the disk, and that the disk is only moderately correlated with CPU and file system. Disk appears to be the least correlated with other resources. Furthermore, all server resources are positively correlated, except memory and disk. The negative correlation between memory and disk suggests that servers with a high memory utilization probably have a low disk space and vice versa. The correlation coefficients computed over the weekly and daily data follow the same trends

¹For the sake of clarity, we do not present fitted data on this graph.

with similar absolute values for the weekly data and smaller absolute values for the daily data. Due to the lack of space, we omit correlation coefficient figures of weekly and daily data.

Typically, a negative correlation indicates a potential benefit of consolidation as the peak-times and off-times of resources of server utilizations are different. Resource management factoring on the negative dependency among resources can reach a more effective and economical operation in data centers, e.g., virtual machine consolidation [24]. On the contrary, the positive correlation implies that the correlated resources can be bundled together for capacity planning. Moreover, neglecting the positive dependency easily leads to under provisioning of resources. From our correlation analysis, we conclude that resource demands of CPU, memory, and file system at typical data centers are dependent on each other with correlation coefficients around 0.25, while disk can be approximately analyzed as an independent resource.

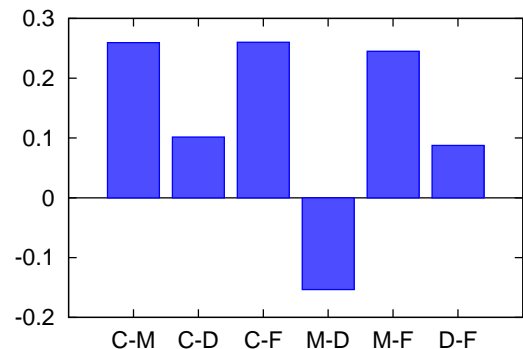


Figure 4: Correlation coefficients between CPU (C), memory (M), disk (D), and file system (F) on all servers.

To better present and validate the correlation coefficients among resources over a closely "related" set of servers, we select from our trace data those that correspond to the five different enterprises presented in Figure 3. We assume that the servers in each enterprise tend to be similar in terms of workload characteristics. In Figure 5, we compute the correlation coefficients among resources for each enterprise. On one hand, one can see that enterprise 1 and 2 have very similar correlation coefficients among resources as the ones computed from the all servers, except of the coefficient between memory and disk utilization. In general, the negative correlation shown in Figure 4 only appears on Enterprise 3. On the other hand, enterprise 4 and 5 have quite different correlation coefficient values from Figure 4. Enterprise 4 even has negative correlation between (1)disk and file system and (2) disk and CPUs. Combining results based on all servers and enterprise servers, we conclude that the correlation coefficient analysis based on all server is valid for some servers in a particular enterprise, but the correlation among memory and disk can be very different or unpredictable.

3.3 Temporal Variability of CPU Utilization

The $\overline{U}_{i,j}$ presented previously are aggregates over all the servers and thus overlooks the temporal variabilities. Such an effect is negligible for the disk, memory, and file system utilizations which represent rather static space demands. In comparison, the CPU utilization is a measure of activity and thus temporal variability is crucial in evaluating workloads [10]. Here, we take a micro approach to characterize the server diversity in terms of temporal variability. Between May 1 and May 7 2011, we present the CPU utilization at time granularities of a minute.

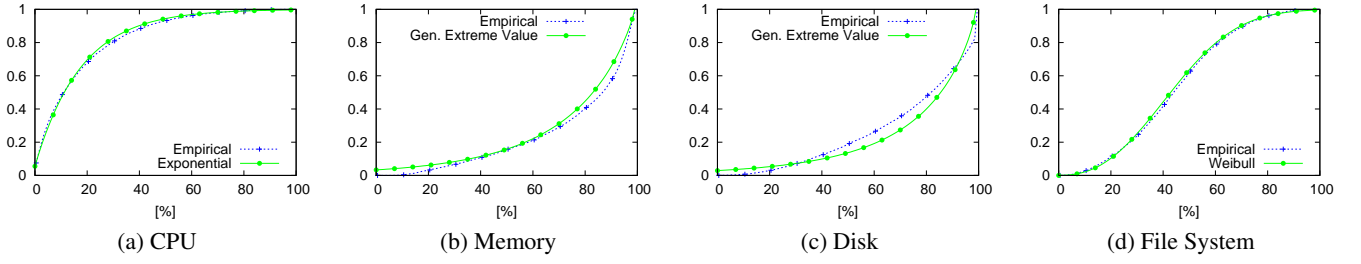


Figure 2: CDF of μ_i over two years monthly averages, across all servers.

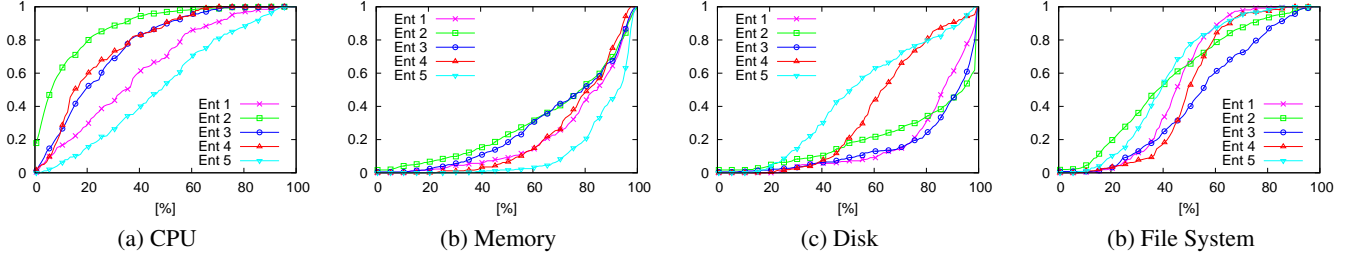


Figure 3: CDF of μ_i over one year weekly averages, across servers dedicated to selected enterprises.

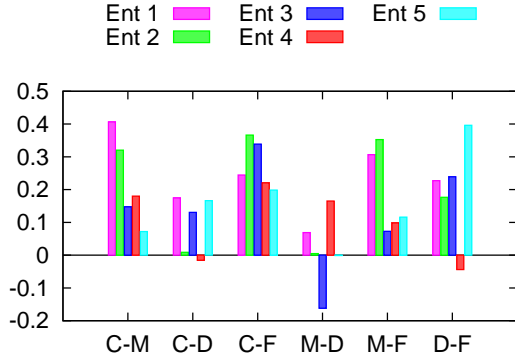


Figure 5: Correlation coefficients between CPU (C), memory (M), disk and file system (F) on servers of selected enterprises.

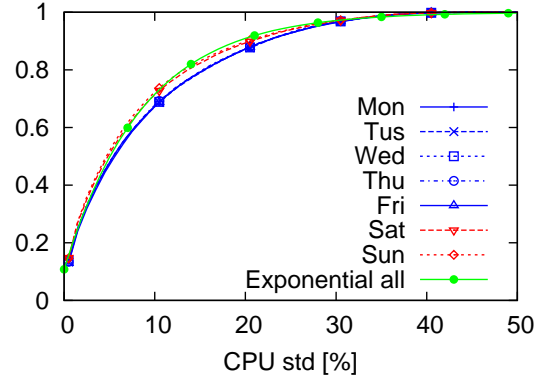


Figure 6: CDF of CPU utilization standard deviation over different days based on minute utilization data).

3.3.1 Standard Deviation of CPU Utilization

From the minute data we compute for each server the standard deviation over the entire day in order to quantify the utilization variability per day. We then plot the empirical CDFs of the standard deviations across servers in Fig. 6. One can see that the CDFs of week days overlap with each other, so do the weekend days and the difference between the two groups is quite small. Moreover we can fit the data reasonably well with an exponential distribution. We use parameters $\mu = 8.75$, $\mu = 8.98$, and $\mu = 8.16$, for all, week, and weekend days, respectively. On the average, the variability of CPU utilization during the weekend is lower than during the week. It is a welcome surprise that an exponential distribution can capture the CPU variability across servers. Combining with the analysis in Section 3.1, we conclude that a typical data center server has both average CPU utilization and temporal standard deviation of CPU utilization following exponential distributions with mean $\mu = 17.64$ and $\mu = 8.75$ respectively. One can refer to this data to validate large scale data center studies.

3.3.2 CPU Temporal Minimum, Median, and Maximum

One common practice of system resource provisioning is based on the the peak usage [35], especially CPU. Therefor, we extend the study of the minute data also to the minimum, median, and maximum CPU utilization per day. More in particular, we compute for each server the 5^{th} , 50^{th} , and 95^{th} percentile of the CPU utilization distribution per day, and use them as a representation of the minimum, median and maximum CPU utilization of a server.

We plot the CDFs of minimum, median, and maximum CPU utilizations in Figure 7. In Section 3.3.1, we already hinted to the fact that the CPU utilization distributions over week days and weekend days are very similar. We observed the same overlapping also for CDFs of the minimum, median and maximum CPU utilizations. Hence for space reasons, we plot only the distributions related to Sunday May 1st, 2011 (green/light lines) and Wednesday May 4th, 2011 (pink/dark lines) representing respectively a weekend day and a week day. The steeper the lines, the more the servers are concen-

trated towards low CPU utilization values, whereas the flatter the lines the more the servers are concentrated towards high CPU utilization values. One can expect that the minimum utilization CDF lines to be higher than the median utilization CDF lines which in turn to be higher than the maximum utilization CDF lines.

In Figure 7(a), we plot the CDFs based on all servers. The differences between a weekday and weekend day grows with the percentile. For the min CDF, the difference is negligible, whereas the difference is most noticeable for the max CDF. Also in agreement with the previous results, week day samples are more spread out than weekend days. The minimum CPU utilization is concentrated towards the 0 value. Specifically, one can see that roughly 37 percent of servers have minimum CPU utilization values less than 1%, and still another roughly 73 percent of servers have minimum CPU utilization values between 1% and 10%. Similarly from the median CPU utilization, one can see that roughly 80 percent of servers have their medium CPU utilization values between 0% and 30%, i.e. half of the time a majority of servers have their CPU utilization between 0% and 30%. Such an observation matches with the previous analysis listed in Table 1 that the majority of servers are under-utilized.

However, from the maximum CPU utilizations, one can observe that servers indeed have some high CPU utilization values. One can see that more than 15 percent of servers have their maximum CPU utilization higher than 80%. In other words, there is a non-negligible share of servers that have a high peak load. Combining with the observation from minimum and median CDF, we further infer that there is a big gap between peak and off-peak CPU usage. As such, when systems are built to satisfy peak loads, the resources tend to be over-provisioned. A qualitative idea is given by the distance between the CDF lines. The closer the lines the smaller the gap, the further the lines the bigger the gap. To give a more precise idea of this gap we plot in Figure 8 the PDF of the difference between the per day minimum and maximum CPU utilization across all servers. The PDF appears similar to the uniform distribution and implies it is very difficult to have good prediction on the max and min CPU loads.

Finally we present the same CDFs for enterprise 1 and 2 in Figure 7(b) and 7(c) to show how well/bad the overall picture can match the global statistics. Client 1 has higher CPU workloads, as indicated by the flatter initial slopes of the CDF lines, whereas Enterprise 2 has a lower CPU workloads, as indicated by the steep initial slopes of the CDF lines. For a given percentile CDF, the difference between week days (supposedly high load) and weekend days (supposedly low load) is bigger for individual Enterprise than for all servers. This can be explained by the smoothing effect due to aggregating of a larger number of servers.

3.4 Peek into Memory Activities

Our aim here is to show how servers use memory within a day and quantify the memory variability across time and across servers. As such, we present fine-grained memory related statistics in terms of space utilization and paging activities. Similarly to Section 3.3 we collect the per minute memory statistics on Sunday May 1, 2011 and Wednesday May 4, 2011 and compute the empirical CDFs of the minimum, median, and maximum.

3.4.1 Memory Space Utilization

In Figure 9(a), we first present the empirical CDFs of the space memory utilization of all servers. All CDFs have similar shapes and the differences among a weekday and weekend is negligible. Moreover as indicated by how close the lines are, the differences among minimum, median and maximum memory utilization values are

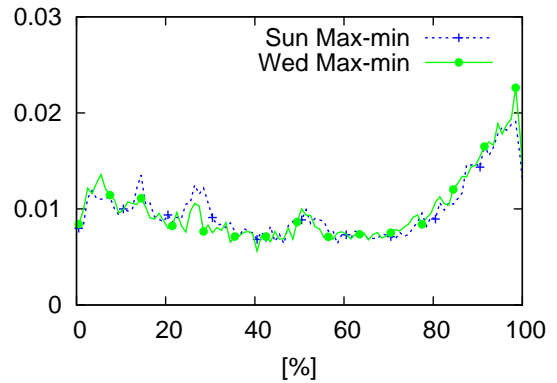


Figure 8: PDF of the difference between max and min CPU utilization within a day, across all servers.

quite small compared to the CPU utilization in Figure 7(a). Both indicate that the memory usage of a server indeed varies somewhat, but overall it is quite static across time.

3.4.2 Memory Paging Rate

Memory space utilization provides only a partial information regarding to the memory performance. In addition to memory usage statistics, we present the memory paging activity in KB/s which can be a good indicator of memory related problems. We express this as the sum of the page-in and page-out rates as reported by `vmstat`. We plot the CDFs of the memory paging rates across all servers in Figure 9(b) (maximum) and 9(c) (median). We skip the minimum since it is always very close to zero. Due to the long tails of the distributions, we add small inset boxes with a bigger scale.

One can see that the paging rate on Wednesday is slightly higher than Sunday, as indicated by the lower CDF lines of max and median. For the max paging rates, roughly 30 percentage of servers have values at zero. This implies that roughly 30 percent of servers have no paging activities. Moreover, roughly 50 percent of servers have their maximum paging rate ranging between 0 to $2000KB/s$. However the distribution clearly shows a long tail and a 20 percent of servers have the maximum paging rates ranging up to $40000KB/s$ or even beyond, as indicated by the small box in Figure 9(b).

Switching to Figure 9(c), one can see that a 70 percent of servers have their median paging rate at zero. It implies that 70 percent of servers have no paging activities for half of the time and 20 percent of servers have their median paging rates up to $10KB/s$, which is a fairly low value. In other words, roughly 90 percent of servers has negligible paging activities for half of their execution. We conclude that most of servers do not suffer from paging for most of the time, and 20 percent of servers experience a non-negligible paging activities for a very small fraction of execution time.

3.5 Discussion and Summary

The following summarizes the gist of our findings:

- 1) The representative average workloads in current data centers are $CPU = 18\%$, $memory = 78\%$, $disk = 75\%$, and $filesystem = 45\%$. The CoV of utilization across servers is more pronounced for CPUs, and less for file systems, memory, and disk.
- 2) Generally speaking, CPU, memory, and file system utilizations in a data center server are moderately correlated, while disk utilization is less tightly correlated with the rest of the resources. Most resources are positively correlated, except memory and disk. Consolidation should factor in the correlations among resources

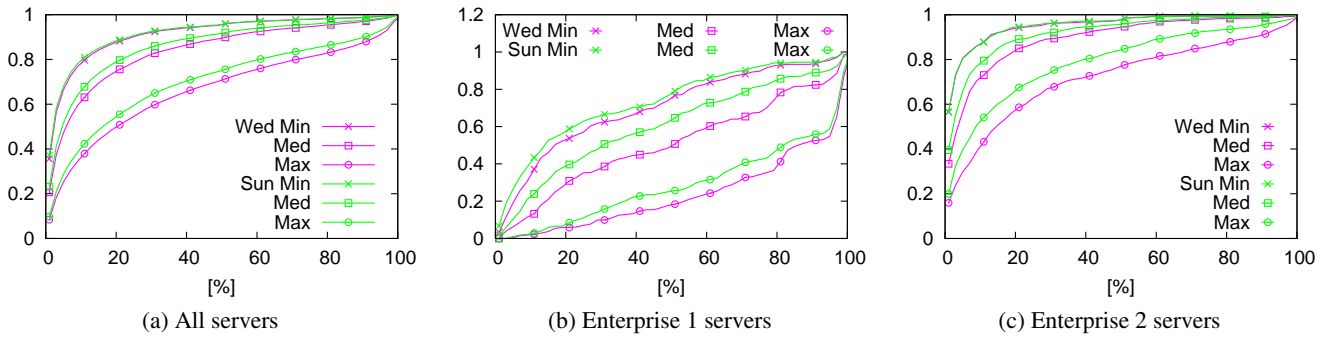


Figure 7: Distribution of CPU minute frequency

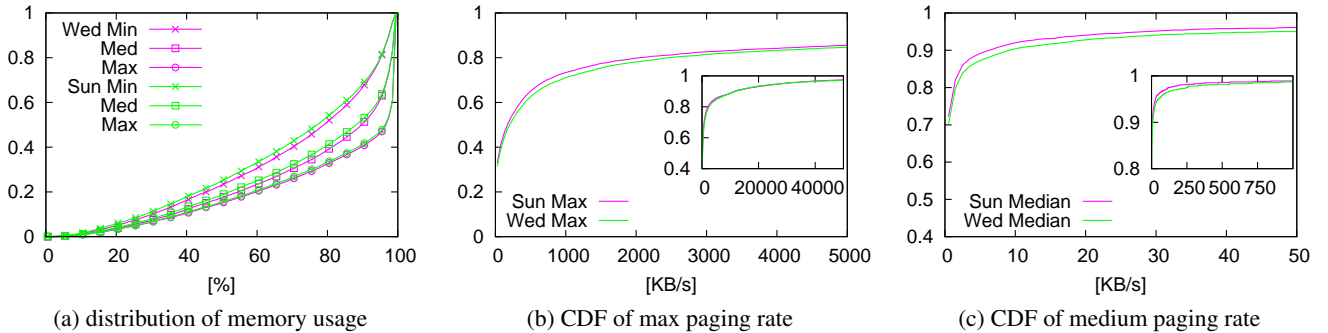


Figure 9: Memory related statistics

and this study provides a qualitative measurement of correlations. We conclude that the provisioning of CPU, memory, and file system should be bundled and that there is a good potential in resource saving in consolidating servers with heavily loaded memory and disk.

3) As CPUs are averagely utilized at 15 – 20% following a truncated exponential distribution, the theoretical upper bound of data center system improvement via workload consolidation is roughly six fold (1/15%), given that the current utilizations are maintained and other resources are not the bottleneck. However, as there are more than forty percent of servers having maximum CPU utilization greater than 50 % and the system needs to fulfill peak loads, efficiency gains thanks to consolidation are immediately limited.

4) Roughly a 20 percent of servers experience a higher paging rate for a short period of time, while most of servers have low memory paging rates.

4. EVOLUTIONARY VIEW

In this section we analyze the time evolution of resource utilizations by category, including geographical location, different time shifts, different seasonality, physical versus virtual machine, and selected enterprises. In addition to the various statistics, we also provide a simple economics analysis that becomes important for capacity planning [13, 31, 36].

4.1 Economics Analysis

The utilization reflects the efficiency as well as the difference between resource demands and supplies at data centers. The average utilization of resource i during month t is defined by the ratio of

the demand $D_i(t)$ over the supply $S_i(t)$:

$$\mu_i(t) = \frac{D_i(t)}{S_i(t)}. \quad (1)$$

Let α_i be the growth rate of resource i utilization. Assuming linear growth of resource utilizations, one can write

$$\mu_i(t+1) = (1 + \alpha_i) \mu_i(t). \quad (2)$$

When resource utilizations show increasing trends $\alpha_i > 0$; otherwise $\alpha_i < 0$. Combining Eq. 1 and Eq. 2, one can obtain:

$$(1 + \alpha_i) = \frac{D_i(t+1)/D_i(t)}{S_i(t+1)/S_i(t)} \quad (3)$$

Essentially, the utilization growth indicates the relative difference between the demand growth ($D_i(t+1)/D_i(t)$) and the supply growth ($S_i(t+1)/S_i(t)$). When $\alpha_i > 0$, it implies that the average demand growth is greater than the supply growth at data centers. On the contrary, the drop of utilization values implies that supply growth is greater than demand growth. In real terms, we can infer that a particular resource has been greatly upgraded. We note that resource demands and supplies at data centers highly depend on the regional economic growth.

We plot the evolution of resource utilizations in Fig. 10. Focusing on the *all* curve (that does not separate the shifts) CPU, disk, and file system show increasing trends ($\alpha_c = 0.22 > \alpha_d = 0.14 \approx \alpha_f = 0.15 > 0$) indicating that their demand growths are stronger than their supply growths. One can observe more frequent dips in CPU and memory utilizations and infer that CPU and memory are upgraded more frequently. Even if in Table 1 CPU utilization is low, the CPU utilization shows a strong increasing trend ($\alpha_c = 0.22$) over time, justifying an initial partly over provisioning of CPU resources. Furthermore this observation matches with

the free market mechanism where the CPU is moving towards the market equilibrium in terms of provisioning.

Overall, with a small $\alpha_d = 0.09$, memory utilization is rather constant at about the common upper limit employed by OSs in managing memory subsystems. As a result, one can infer that the memory demand is roughly the same as the memory supply and at equilibrium. The constant values further indicate that the demand growth is the same as the supply at data centers. Compared to other resources, memory is upgraded with a higher frequency and this can point to issues where memory is the bottleneck resource for many of today's applications.

4.2 Impact of Time Shifts

We separate monthly utilization values by different shifts, i.e., prime, off-prime, weekend, and all in Fig. 10. In general, CPU, file system, and disk utilizations increased by 4%, 3%, and 2% respectively in the past two years, whereas memory utilization appears more constant.

One can observe that the differences among shifts are not much (roughly 1%) especially for memory, disk, and file system utilizations. For CPUs, the difference between prime and weekend shifts is roughly 3%, and the difference between prime and off-prime shifts is roughly 1%. Today's CPUs during the prime shifts are loaded with 1.08 and 1.2 times more workload than on off-prime and weekend respectively. Such an observation leads to the speculation that there might be little resource efficiency to gain by exploring CPU on different shifts, given the criterion of maintaining performance. Similar observations held also for the memory utilization, i.e., prime shift has a higher value than off-prime, whose value is higher than weekend. On the contrary, the shifts have a different effect on disks, i.e., the utilization order for shifts is: weekend, off-prime, all, prime. This observation matches with current common practice of data backup during the off-peak time. Finally, the file system gives mixed signals, i.e., higher utilization at the weekend (matching with the higher disk utilization and therefore most probably related to the backup operations) and lower utilization during the off-prime shift with respect to the prime shift. This last observation could be due to temporary files created by programs during their operation and then deleted.

Considering all shifts, we plot the resource evolutions of the selected enterprises in Figure 11. Compared to the previous subsection, one can see not only the resource evolutions but also their variability among enterprises. Using these specific enterprises, we illustrate how such time series analysis can facilitate the IT capacity planning for enterprises. Enterprise 1 has an obvious CPU workload increase and rather mild workload increases in other resources. One can speculate that for the workloads executed on enterprise 1 an upgrade on CPU might be needed in the near future. Similar observation can be made to the enterprise 5. For enterprises 2 and 3, all resource utilization values are very flat and this can be possibly explained by that its demand growth is very similar to its supply growth. Enterprise 4 has a lot of fluctuations in all resources, especially the memory. This can be indicated by the frequent upgrades. In general, one can use such time series analysis to predict the computational and storage demands for those enterprises.

4.3 Seasonality of CPU

Time-varying workloads are the key for short term as well as long-term capacity planning of computer and communication systems [4, 8, 20]. A large number of studies [7, 10, 11, 33] have identified time series of the family of autoregressive moving average (ARMA) models [12], which can capture well the seasonal effects, for some specific systems and smaller scale datasets. Here, we aim

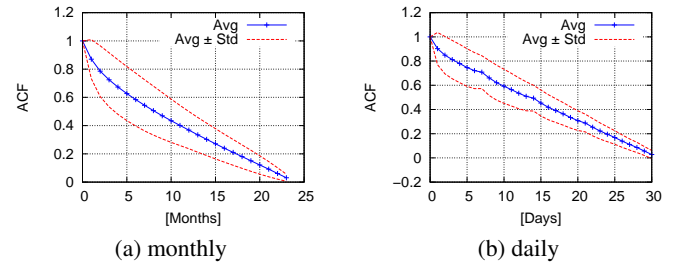


Figure 12: ACF of CPU utilization at different scales: average \pm standard deviation across all servers.

to present the representative ACFs from the time series of CPU utilization at different scales, i.e., monthly, and daily, such that one can obtain the important parameters for ARMA models via our analysis.

For each server, we compute its ACF from the CPU utilization $U_{c,j}(t)$ where t is index for one of 24 months or one of 30 days. For each lag in the ACFs, we further compute the average and standard deviation from all servers. Note that the standard deviation represents the heterogeneity of ACF across different servers. We summarize the curves of the average ACF plus/minus its standard deviation in Fig. 12. The region between the upper and lower curves shows where the majority of data center servers ACF is located. The shapes of ACFs indicate that the time series of CPUs can be captured and predicted well by an autoregressive model. Specifically, the monthly ACF decays faster than than the daily ACF. Furthermore, on the daily ACF it is possible to identify peaks in 7, 14 and 21 days indicating a clear weekly periodicity. One can take these ACF curves to forecast the CPU workloads for their yearly, monthly, weekly, and even daily capacity planning.

4.4 Is CPU waiting for I/O?

We use our data sets to verify such a conjecture by breaking down the CPU utilization into the three modes given by vmstat: system, user, and I/O wait. For each server, we collect the average daily percentage of CPU utilization in these three modes during the whole month of May, 2011. For each day, we compute the mean of all modes across all servers. In Figure 13, we report the CPU utilization in system, user, and I/O wait modes. One can see that most of CPU active time is spent in user and system modes; whereas only a negligible CPU time is spent on I/O wait. Moreover as already seen previously, there is a clear seasonality effect, i.e., there are periodic dips during the weekends. However, the seasonality effect is less visible for the CPU time spent in the system mode than in the other two modes. This indicates that the bookkeeping by the OS is rather constant over time.

Unfortunately the I/O wait time does not include the time spent in network I/O other than file I/O over network, e.g., in the case of network file system. Therefore, we can not directly relate our observations with many DC studies, e.g., [16], pointing out that CPU suffers from relatively low utilization due to not-that-great network performance.

4.5 Memory Breakdown

Similar to CPU breakdown analysis, one can also break down the memory usage into several categories, depending on its use. In general, the vmstat command breaks down the used memory space in buffers, cache, and other, which includes anything which is not buffers or cache. For each server, we collect these three statistics during the entire month of May, 2011. Further, we compute the

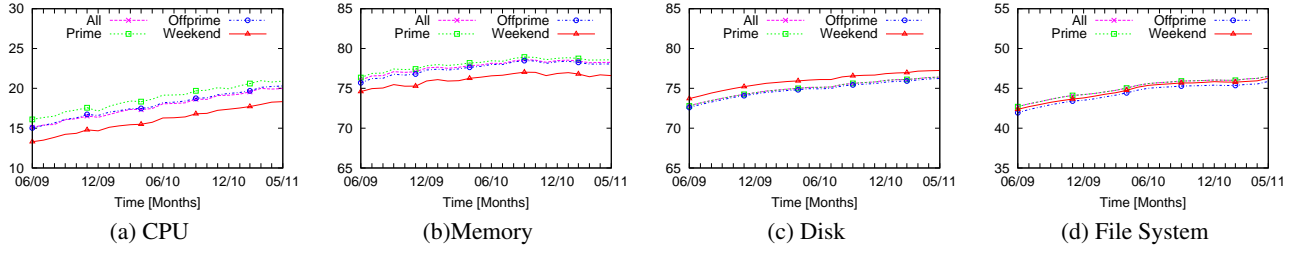


Figure 10: Time series of the average resource utilization, μ_i , on different shifts.

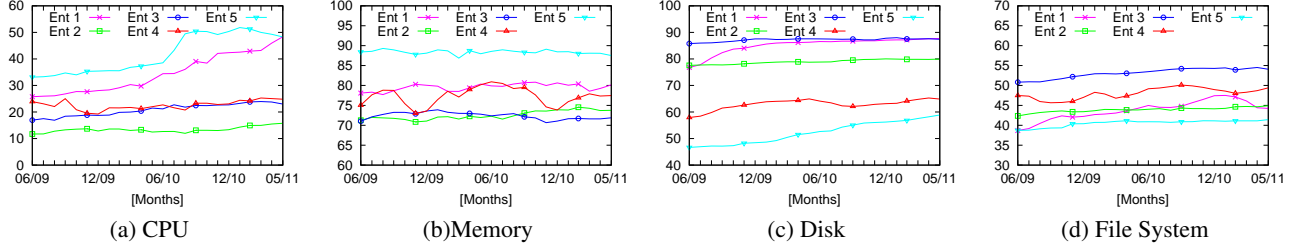


Figure 11: Time series of the average resource utilization, μ_i , on selected enterprises.

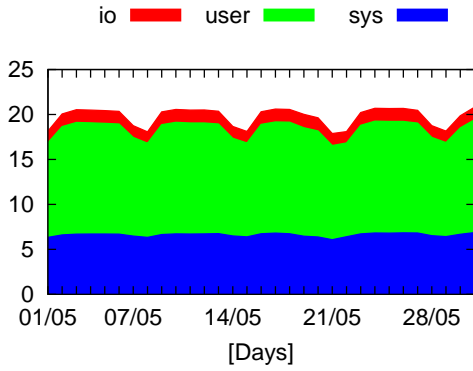


Figure 13: Mean CPU utilization break down across all servers: system, user, and IO wait.

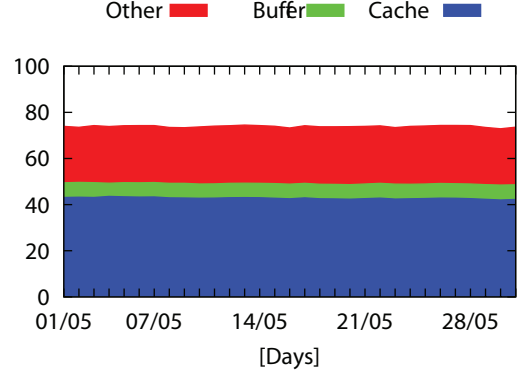


Figure 14: Memory usage breakdown across all servers: buffer, cache, other and free.

average of those values across all servers, and summarize the results in Figure 14. In contrast to the CPU break down, there is no temporal variability in the overall memory usages nor in each category. The total used memory is roughly 77%, which corresponds to the values presented in Table 1. As pointed out earlier, many OSs keep memory roughly utilized around this value to avoid fragmented memory space and heavy memory swapping. A very interesting observation is that on average the applications uses only around of 30% of the total memory, whereas more than half of the memory space is used for system buffer and cache to speed up I/O operations. It shows that operating systems well utilize the available memory space.

4.6 Geographical Locations

Our objective in this subsection is to answer the questions if servers at different geographical locations are the same and how the resource demands deviate from the general observations explained in the previous subsection. We summarize all resource utilizations of some selected countries in Fig. 15. The corresponding averages and standard deviations across all time windows are summarized

in Table 1. One can see that CPU, memory, disk and file system utilizations lie between 7 – 25%, 72 – 84%, 64 – 81%, 39 – 55% respectively for the selected countries. Nevertheless they show various trends in different geographic areas.

We observe a significant growth of supply at data centers, i.e., resource utilization dips, at developing countries and areas. Instead, developed countries have rather steady and static utilization trends. In particular, country C has the lowest CPU, memory, disk and file system utilizations among all the selected countries, whereas country A has the highest resource utilizations at almost all the times. Country E has very steady trends on all resource utilizations, as the number of servers under survey is significantly large, compared to the other countries. In the following, we describe some observations with respect to each resource.

CPU. The CPU evolution, α_c , of servers differs at different geographic locations. From our data set, we observe that most of the continents have CPU demand growths greater than CPU supply growths. When we look into selected countries, one can also observe a high supply growth in country A. Country B, D, and E all show steady and increasing demand growths and the CPU utiliza-

tions converge to values between 16 – 18%. In general, one can observe periodical dips in CPU utilizations which reflect seasonality and irregular dips reflecting possible upgrades.

Memory. Overall, the evolution of memory utilization is rather flat, but with slightly smaller dips, compared to CPU. One can observe that the evolution of memory at different countries well matches with the CPU evolution, though the relative difference of memory utilization at countries is not the same as for the CPU. The frequent dips in memory utilization for countries B, C, D, and A show that memory is probably often upgraded, whereas country E has a rather steady memory evolution line.

Hard disk. Clearly, one can observe, from the increasing trend of disk utilizations, that disk demand growth is greater than the supply growth in all macro regions.

File system. Most countries show an utilization growth, α_f , around the average value of 0.15 over the 2 years observation period. Country D and B show steady decreasing trends in file system for a short period. We attribute this to innovation in file system technologies.

4.7 Physical vs. Virtual Machines

In this subsection, we focus on the evolution of physical vs. virtual resource utilizations. Virtualization technologies [7, 34] have been greatly employed in the past 10 years to increase resource efficiencies by providing a uniform platform over heterogeneous software and hardware stacks. Moreover, due to the booming of cloud computing [1], a large number of studies have focused on managing virtual resource provisioning [7, 14, 23, 39]. We aim to qualitatively show the efficiency improvement due to the virtualization technologies.

From our data set, we summarize how resource utilizations differ among physical and virtual servers in Fig. 16. For CPU and memory, virtual servers show higher utilizations (with higher increasing trends) than physical ones. The virtual CPU utilization increases from 19.4% to 27.6% in two years, whereas the physical CPU utilization increases only from 10.9% to 12.9%. Clearly, virtualization technologies double the CPU utilization. However, the difference between virtual and physical CPU utilization could be due to the virtualization overhead. Both virtual and physical memories bear steady utilizations around 85% and 70% respectively. One can infer that for both memory demand growth is the same as the memory supply growth. On the contrary, virtual disks have lower utilizations (around 70%) than the physical disks (around 76%). For the file systems, both physical and virtual servers have similar utilizations around 43% with slightly higher utilization growth for virtual ones. In summary one can see that, indeed, the virtual resources show a higher demand growth than the supply growth, especially for CPU.

5. RELATED WORK

In the past few years there have been a host of research efforts focusing on improving data center performance. These studies could be roughly classified at those that focus on power management, capacity planning and resource provisioning, and traffic engineering. A detailed survey of the extensive related work is not possible here, for more details we direct the interested reader in the representative papers that are surveyed here and references therein.

From the perspective of power/cooling delivery and their respective costs, there have been efforts on exploiting the platform heterogeneity that characterizes cloud data centers. They consider differences in power management of the various heterogeneous components and emphasize is on the effectiveness for an allocation mechanism that maps workloads to best fitting platforms [27]. A predic-

tion mechanism that predicts thermal behavior within a data center is proposed in [26] and is used for improving data center power effectiveness. The premise of much of the related work regarding power effectiveness in data centers focuses on either moving workloads from under-utilized servers and shutting them down, e.g., [9] or provide techniques that offer trade-offs on performance with energy reduction, e.g., [11]. Workload-aware statistical multiplexing has been proposed in [15] and evaluated on a small scale prototype of a data center. In general, the evaluation of the proposed data center power saving mechanisms are not easy, as field evaluations are rare and there is no clear workloads and demands placed on data centers. Our evaluation study complements the above works by providing a baseline for real-world data center utilization behavior that can be used to drive the evaluation of different power saving strategies as in [29].

An additional angle of viewing the problem of performance in data centers is that of resource provisioning via multiplexing. Most of these studies have focused on how to multiplex applications and have proposed techniques to best meet service level objectives of different applications, aiming at automated capacity planning and workload management. Virtualization and architectures that effectively support it have been described as key to achieve short- and long-term loads [40]. The non-stationarity of the workload has been viewed as a challenge for allocating server capacity in data centers. Techniques that use clustering algorithms to automatically determine the workload mix have been shown effective in data centers in laboratory settings and have been evaluated experimentally using benchmark applications [32]. The use of online monitoring that is integrated with virtualization technologies within data centers has been proposed in [21]. In general, resource provisioning techniques are evaluated via small scale prototypes at a laboratory setting or simulation. Our large scale characterization study can leverage the above works by providing information on realistic assumptions about the workloads served by data centers.

Last but not least, the need to improve on the data center networking abilities, especially from the geo-diversity and geo-distributing perspective, has been identified in an early position paper as critical [16]. Such work fueled the development of new variants of TCP that are especially tailored for data centers that need to efficiently serve efficiently a diverse mix of short and long flows [3]. In the past years, there is a host of works focusing on the design of networks that are most appropriate for data centers. In turn, such designs focus on certain type of applications, e.g., MapReduce [2, 17, 18] and/or web services [6, 19]. The above works are motivated by measurements of data center traffic. For a comprehensive study on the statistics, topology, and packet-level traffic characteristics of several types of data centers including university, enterprise, and cloud data centers, we direct the interested reader to [6].

To the best of our knowledge, this is the first large scale characterization study that aims to shed some light on the type of real workloads served by today's production data centers. The presented study also offers a glimpse on the performance bottlenecks faced by data centers. Our findings can be used for better parameterization of data center simulation and prototype studies by showing how empirical workloads truly look like.

6. CONCLUSION

In this paper, we surveyed over a period of two years several thousands servers located at different data centers. Using utilization statistics of CPU, memory, disk and file systems, we characterize the typical server workloads and the diversity across servers. In particular, we quantify the average resource utilizations, the tempo-

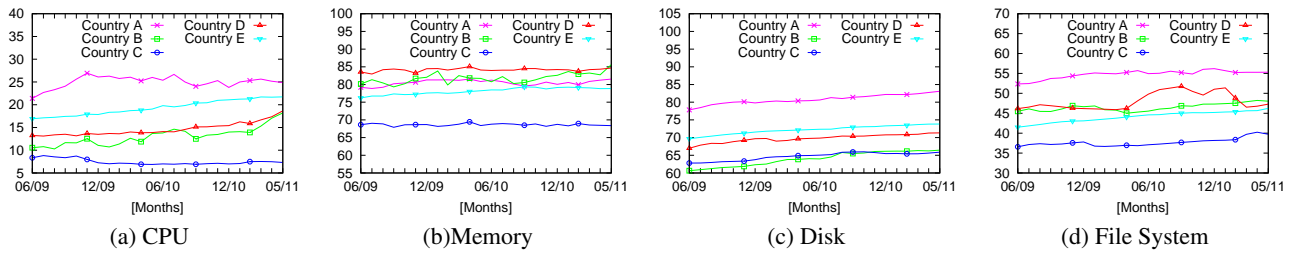


Figure 15: Resource utilizations at selected countries.

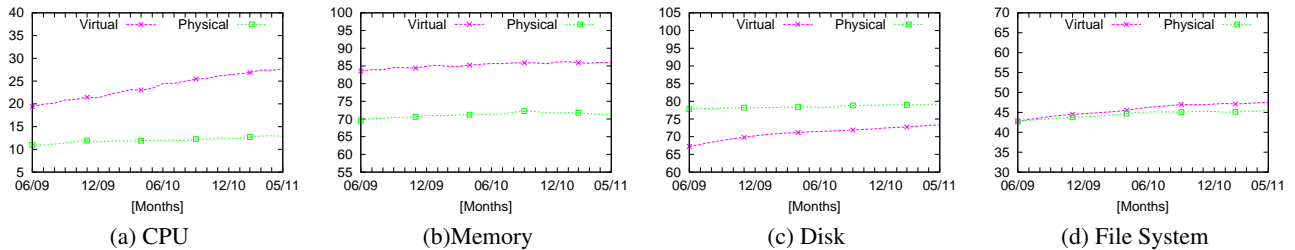


Figure 16: Physical vs. Virtual machines.

ral variability across servers and diversity across resource demands. Moreover, with respect to different categories, e.g., we present the evolution of server resource demands and their economic interpretations. To quantify seasonality, we characterize the autocorrelation across servers at different time scales, i.e., monthly, weekly, and daily. Our analysis on the diversity and evolution of servers provides a basis for resource management and capacity planning at data centers.

7. REFERENCES

- [1] Amazon ec2. <http://aws.amazon.com/ec2/>.
- [2] M. Al-Fares, A. Loukissas, and A. Vahdat. A scalable, commodity data center network architecture. In *SIGCOMM*, pages 63–74, 2008.
- [3] M. Alizadeh, A. G. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan. Data center tcp (dctcp). In *SIGCOMM*, pages 63–74, 2010.
- [4] M. Arlitt and T. Jin. Workload characterization of the 1998 world cup web site. *IEEE Transaction on Network*, 14(3):30–37, 2000.
- [5] L. N. Bairavasundaram, G. R. Goodson, B. Schroeder, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. An analysis of data corruption in the storage stack. In *FAST*, pages 223–238, 2008.
- [6] T. Benson, A. Akella, and D. A. Maltz. Network traffic characteristics of data centers in the wild. In *Internet Measurement Conference*, pages 267–280, 2010.
- [7] N. Bobroff, A. Kochut, and K. Beaty. Dynamic placement of virtual machines for managing SLA violations. In *Integrated Network Management*, pages 119–128, 2007.
- [8] S. C. Borst, A. Mandelbaum, M. I. Reiman, and M. Centrum. Dimensioning large call centers. *Operations Research*, 52:17–34, 2000.
- [9] G. Chen, W. He, J. Liu, S. Nath, L. Rigas, L. Xiao, and F. Zhao. Energy-aware server provisioning and load dispatching for connection-intensive internet services. In *NSDI*, pages 337–350, 2008.
- [10] L. Y. Chen, A. Das, A. Sivasubramaniam, Q. Wang, R. Harper, and M. Bland. Consolidating clients on back-end servers with co-location and frequency control. In *SIGMETRICS/Performance*, pages 383–384, 2006.
- [11] Y. Chen, A. Das, W. Qin, A. Sivasubramaniam, Q. Wang, and N. Gautam. Managing server energy and operational costs in hosting centers. In *SIGMETRICS*, pages 303–314, 2005.
- [12] G. Edward, P. Box, and G. M. Jenkins. *Time Series: Analysis: Forecasting and Control*. Wiley, 2008.
- [13] D. Gmach, J. Rolia, L. Cherkasova, and A. Kemper. Capacity management and demand prediction for next generation data centers. In *ICWS*, pages 43–50, 2007.
- [14] Z. Gong and X. Gu. PAC: Pattern-driven application consolidation for efficient cloud computing. In *MASCOTS*, pages 24–33, 2010.
- [15] S. Govindan, J. Choi, B. Urgaonkar, A. Sivasubramaniam, and A. Baldini. Statistical profiling-based techniques for effective power provisioning in data centers. In *EuroSys*, pages 317–330, 2009.
- [16] A. G. Greenberg, J. R. Hamilton, D. A. Maltz, and P. Patel. The cost of a cloud: research problems in data center networks. *Computer Communication Review*, 39(1):68–73, 2009.
- [17] C. Guo, H. Wu, K. Tan, L. Shi, Y. Zhang, and S. Lu. Dcell: a scalable and fault-tolerant network structure for data centers. In *SIGCOMM*, pages 75–86, 2008.
- [18] D. Halperin, S. Kandula, J. Padhye, P. Bahl, and D. Wetherall. Augmenting data center networks with multi-gigabit wireless links. In *SIGCOMM*, pages 38–49, 2011.
- [19] B. Heller, S. Seetharaman, P. Mahadevan, Y. Yiakoumis, P. Sharma, S. Banerjee, and N. McKeown. Elastictree: Saving energy in data center networks. In *NSDI*, pages 249–264, 2010.
- [20] O. B. Jennings, A. Mandelbaum, W. A. Massey, and W. Whitt. Server staffing to meet time-varying demand. *Management Science*, 42:1383–1394, 1996.

- [21] M. Kutare, G. Eisenhauer, C. Wang, K. Schwan, V. Talwar, and M. W. Matthew. Monalytics: online monitoring and analytics for managing large scale data centers. In *Proceeding of the 7th international conference on Autonomic computing*.
- [22] S. Malkowski, M. Hedwig, and C. Pu. Experimental evaluation of n-tier systems: Observation and analysis of multi-bottlenecks. In *IEEE International Symposium on Workload Characterization (IISWC)*, pages 118–127, 2009.
- [23] S. Mehta and A. Neogi. ReCon: A tool to recommend dynamic server consolidation in multi-cluster data centers. In *NOMS*, pages 363–370, 2008.
- [24] X. Meng, C. Isci, J. Kephart, L. Zhang, E. Bouillet, and D. Pendarakis. Efficient resource provisioning in compute clouds via VM multiplexing. In *ICAC*, pages 11–20, 2010.
- [25] A. Mishra, J. Hellerstein, W. Cirneand, and C. Das. Towards characterizing cloud backend workloads: Insights from Google compute clusters. *SIGMETRICS Perform. Eval. Rev.*, 37:34–41, 2010.
- [26] J. D. Moore, J. S. Chase, and P. Ranganathan. Weatherman: Automated, online and predictive thermal mapping and management for data centers. In *ICAC*, pages 155–164, 2006.
- [27] R. Nathuji, C. Isci, and E. Gorbato. Exploiting platform heterogeneity for power efficient data centers. In *ICAC*, page 5, 2007.
- [28] E. B. Nightingale, J. R. Douceur, and V. Orgovan. Cycles, cells and platters: an empirical analysis of hardware failures on a million consumer pcs. In *EuroSys*, pages 343–356, 2011.
- [29] S. Pelley, D. Meisner, P. Zandevakili, T. F. Wenisch, and J. Underwood. Power routing: dynamic power provisioning in the data center. In *ASPLOS*, pages 231–242, 2010.
- [30] B. Schroeder, E. Pinheiro, and W.-D. Weber. DRAM errors in the wild: A large-scale field study. In *SIGMETRICS/Performance*, pages 193–204, 2009.
- [31] U. Sharma, P. J. Shenoy, S. Sahu, and A. Shaikh. Kingfisher: Cost-aware elasticity in the cloud. In *INFOCOM*, pages 206–210, 2011.
- [32] R. Singh, U. Sharma, E. Cecchet, and P. Shenoy. Autonomic mix-aware provisioning for non-stationary data center workloads. In *International Conference on Autonomic Computing (ICAC)*, pages 21–30, 2010.
- [33] J. Tian, P. Dube, X. Meng, and L. Zhang. Exploiting Resource Usage Patterns for Better Utilization Prediction. In *First International Workshop on Data Center Performance (DCperf)*, 2011.
- [34] B. Urgaonkar, P. Shenoy, and T. Roscoe. Resource overbooking and application profiling in shared hosting platforms. *SIGOPS Oper. Syst. Rev.*, 36:239–254, 2002.
- [35] A. Verma, G. Dasgupta, T. Nayak, P. De, and R. Kothari. Server workload analysis for power minimization using consolidation. In *Proceedings of Usenix Annual Technical Conference*, 2009.
- [36] A. Verma, U. Sharma, R. Jain, and K. Dasgupta. Compass: Cost of migration-aware placement in storage systems. In *Integrated Network Management*, pages 50–59, 2007.
- [37] T. Wood, K. K. Ramakrishnan, P. J. Shenoy, and J. E. van der Merwe. CloudNet: Dynamic pooling of cloud resources by live WAN migration of virtual machines. In *VEE*, pages 121–132, 2011.
- [38] T. Wood, G. Tarasuk-Levin, P. J. Shenoy, P. Desnoyers, E. Cecchet, and M. D. Corner. Memory buddies: Exploiting page sharing for smart colocation in virtualized data centers. In *VEE*, pages 31–40, 2009.
- [39] W. Zhengand, R. Bianchini, J. Janakiraman, J. Santos, and Y. Turner. JustRunIt: experiment-based management of virtualized data centers. In *USENIX Annual technical conference*, pages 18–18, 2009.
- [40] X. Zhu, D. Young, B. J. Watson, Z. Wang, J. Rolia, S. Singhal, B. McKee, C. Hyser, D. Gmach, R. Gardner, T. Christian, and L. Cherkasova. 1000 islands: an integrated approach to resource management for virtualized data centers. *Cluster Computing*, 12(1):45–57, 2009.