

# Stack Metrics: A Taxonomy of Metrics Supporting the Capacity Planning Stack

Richard Gimarc  
CA Technologies  
richard.gimarc@ca.com

Amy Spellmann  
The 451 Group  
amy.spellmann@451research.com

Adrian Johnson  
CA Technologies  
Adrian.Johnson@ca.com

*The Capacity Planning Stack was introduced as a way to organize, discuss and implement a capacity plan that spans the Digital Infrastructure. This paper will examine the fundamental metric categories that support the Capacity Planning Stack. Stack metrics are partitioned into three categories; Business, IT and Facilities. Each category supports different levels within the Stack. Collectively, the metric categories bridge the gap between the high-level business requirements and the lower-level data center that hosts the IT infrastructure.*

*Viewing the set of metrics that support the Stack in terms of the taxonomy enables the practitioner to more easily organize and delineate the data requirements for Digital Infrastructure capacity planning. A constructive side-effect of the taxonomy is that it formally defines the lines of communication required for the successful development of a federated capacity plan that spans the Digital Infrastructure.*

## 1 Introduction

The scope of capacity planning has changed over the past few years. Instead of a traditional focus that ends at the IT infrastructure, today's capacity planners are seeing the breadth of their domain extend into the supporting data center. In other words, today's capacity planner is responsible for planning across the entire Digital Infrastructure which includes the business, applications, IT infrastructure, and the data center. The Capacity Planning Stack (a.k.a. the Stack) was introduced in [SPEL2013] as a way to simplify, structure and focus the practice of capacity planning for today's Digital Infrastructure. The Stack organizes the task of capacity planning for the full breadth and depth of the Digital Infrastructure into a multi-level hierarchy. The hierarchy starts at the Business and progresses through the supporting components. The case was made that the Stack supports a methodology for capacity planning that provides better coverage for today's Digital Infrastructure than the tried-and-true traditional methods that have evolved over the past 30+ years.

In this paper we take a closer look at how to use the Stack. Specifically, we focus our attention on the supporting Digital Infrastructure metrics. This paper introduces a taxonomy that provides structure to the job of identifying and communicating the data needs of the Stack. It will be shown how the taxonomy supports the capacity planner's utilization of the Stack by structuring, organizing and identifying the required metrics.

This paper is organized as follows:

- Section 2 contains a brief review of the Capacity Planning Stack.
- Section 3 introduces the Stack taxonomy.
- Section 4 demonstrates how the taxonomy is used to identify, structure and organize the metrics required for each level in the Capacity Planning Stack.
- Section 5 summarizes the use of the taxonomy to support capacity planning with the Stack.

## 2 The Capacity Planning Stack

The Capacity Planning Stack views today's Digital Infrastructure in a multi-level hierarchy that supports capacity planning workflow from the Business to Facilities (i.e., the data center). The Stack hierarchy describes dependencies and communication between levels (demand and feedback). In addition, each Stack level has a set of efficiency metrics that can be used for long term tracking (measures of success). The entirety of the Stack supports a straightforward and transparent mechanism for the complete assessment of the Digital Infrastructure.

The workflow through the Stack is described in terms of demand, feedback and efficiency metrics:

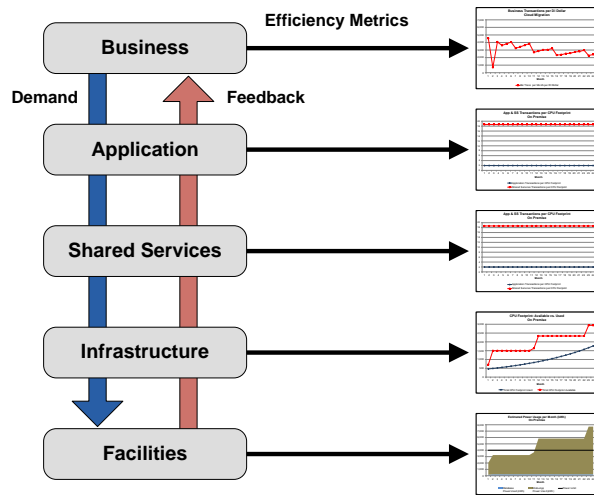


Figure 1. Capacity Planning Stack.

- **Demand flows down the Stack.** Capacity planning starts at the business. Business owners provide application planners with their expected business volumes. Application planners translate the business demand into application-level resource requirements (demand) that are passed downstream to the infrastructure level. It may also be necessary for shared services planners supporting shared database or message queuing tiers to translate application-level requirements into infrastructure-level resource requirements. This level-to-level workflow continues through the Stack. The last demand flow is from the infrastructure level to facilities. This final step is required to ensure that the hosting data center can support the IT equipment required to satisfy the business demand.
- **Feedback flows up the Stack.** A feedback loop communicates results up the Stack to ensure alignment to higher level plans, assist in future planning and potentially refine upstream demands. For example, suppose the infrastructure planners determine that they need 100 additional mid-range servers. What happens if the facilities planners estimate that there is insufficient data center power capacity to support the additional servers? The feedback mechanism provides a means to formally convey this message back up the Stack. In this example, optimization or a change in delivery option at any higher level can potentially alleviate a large facility upgrade expense.
- **Efficiency metrics at each Stack level.** Each level in the Stack has its own set of efficiency metrics that are used to track long term trends. Efficiency metrics serve as a “measure of success” or “report card” for each level in the Stack. For example, the application planners can generate an efficiency metric for their application that describes the number of transactions

processed per unit of resource (similar to miles per gallon for an automobile) and/or report performance of transactions against SLAs. Facilities planners would use PUE as one of their efficiency metrics [TGG2012].

There are a number of metrics used at each level of the Stack to quantify and describe demand, feedback and efficiency. The following table lists some examples [SPEL2013].

	<b>Demand Factors (↓)</b>	<b>Feedback (↑)</b>	<b>Efficiency Metrics (→)</b>
<b>Business</b>	<ul style="list-style-type: none"> <li>- Business volumetrics &amp; priorities</li> <li>- Performance requirements &amp; SLAs</li> </ul>	<ul style="list-style-type: none"> <li>- Total cost</li> <li>- Total time to satisfy requirements</li> <li>- Expected performance</li> </ul>	<ul style="list-style-type: none"> <li>- Business transactions per Digital Infrastructure dollar</li> <li>- Total cost (cumulative from all lower levels)</li> </ul>
<b>Application</b>	<ul style="list-style-type: none"> <li>- Map Business volumetrics to application architecture</li> <li>- Estimated volume of Shared Service and/or Infrastructure requests</li> <li>- Estimate required Application resource footprint and instances</li> <li>- Application-level performance requirements</li> </ul>	<ul style="list-style-type: none"> <li>- Cumulative cost from all lower levels</li> <li>- Application requirements (software licenses and hardware).</li> <li>- Expected performance.</li> <li>- Time to deploy,</li> <li>- Staffing requirements.</li> </ul>	<ul style="list-style-type: none"> <li>- Transactions/minute per resource footprint</li> <li>- Cost per transaction (\$)</li> <li>- Cumulative from lower levels</li> <li>- Performance (e.g., response time) to demonstrate SLA achievement</li> </ul>
<b>Shared Services</b>	<ul style="list-style-type: none"> <li>- Map Shared Service requests to Infrastructure requests</li> <li>- Estimate required Shared Service resource footprint and instances</li> <li>- Shared Services performance requirements</li> </ul>	<ul style="list-style-type: none"> <li>- Cumulative cost from all lower levels</li> <li>- Shared Services requirements (software licenses and hardware).</li> <li>- Expected performance.</li> <li>- Time to deploy</li> <li>- Staffing requirements.</li> </ul>	<ul style="list-style-type: none"> <li>- Transactions/minute per resource footprint</li> <li>- Cost per transaction (\$)</li> <li>- Cumulative from lower levels</li> <li>- Performance (e.g., response time) to demonstrate SLA achievement</li> </ul>
<b>Infrastructure</b>	<ul style="list-style-type: none"> <li>- Translate Application &amp; Shared Service resource footprint and instance requirements to Infrastructure requirements</li> <li>- Determine physical hardware requirements</li> <li>- Initiate procurement process</li> <li>- Evaluates expected performance, headroom and SLAs</li> </ul>	<ul style="list-style-type: none"> <li>- Cumulative cost for infrastructure and facilities</li> <li>- Infrastructure requirements (e.g., servers, storage, network)</li> <li>- Time to procure &amp; deploy</li> </ul>	<ul style="list-style-type: none"> <li>- Count of IT components (servers, storage, network)</li> <li>- Processing capacity per IT component category</li> <li>- Headroom for each IT component category</li> <li>- Cumulative cost of Infrastructure and Facilities</li> </ul>

	Demand Factors (↓)	Feedback (↑)	Efficiency Metrics (→)
Facilities	<ul style="list-style-type: none"> <li>- Estimate required space, power &amp; cooling</li> <li>- Uptime SLA requirements</li> </ul>	<ul style="list-style-type: none"> <li>- Cost for facilities</li> <li>- Data center facilities requirements</li> <li>- Time to satisfy and/or build out</li> </ul>	<ul style="list-style-type: none"> <li>- Power, cooling, space per IT Load</li> <li>- PUE</li> <li>- Facilities headroom</li> <li>- Total Cost (OPEX)</li> </ul>

The above table is just a sampling of the different types of metrics required for a full Digital Infrastructure capacity plan. It is clear that there is a large and diverse set of metrics. Is there a way for the capacity planner to organize the universe of metrics into a more manageable collection?

The intent of the taxonomy introduced in the next section is to provide a structured way to organize, think about and utilize metrics that support the Capacity Planning Stack.

### 3 Stack Taxonomy – Introduction

A taxonomy is defined as a classification of items into a set of ordered categories (see [CAMB2014], [DICT2014], [MERR2014] and [OXFO2014]). Characteristics of a taxonomy classification include the following:

- Natural relationships are preserved
- Classification categories share similar qualities
- Categories have a clear description and means for identification

The following diagram illustrates the Capacity Planning Stack taxonomy.

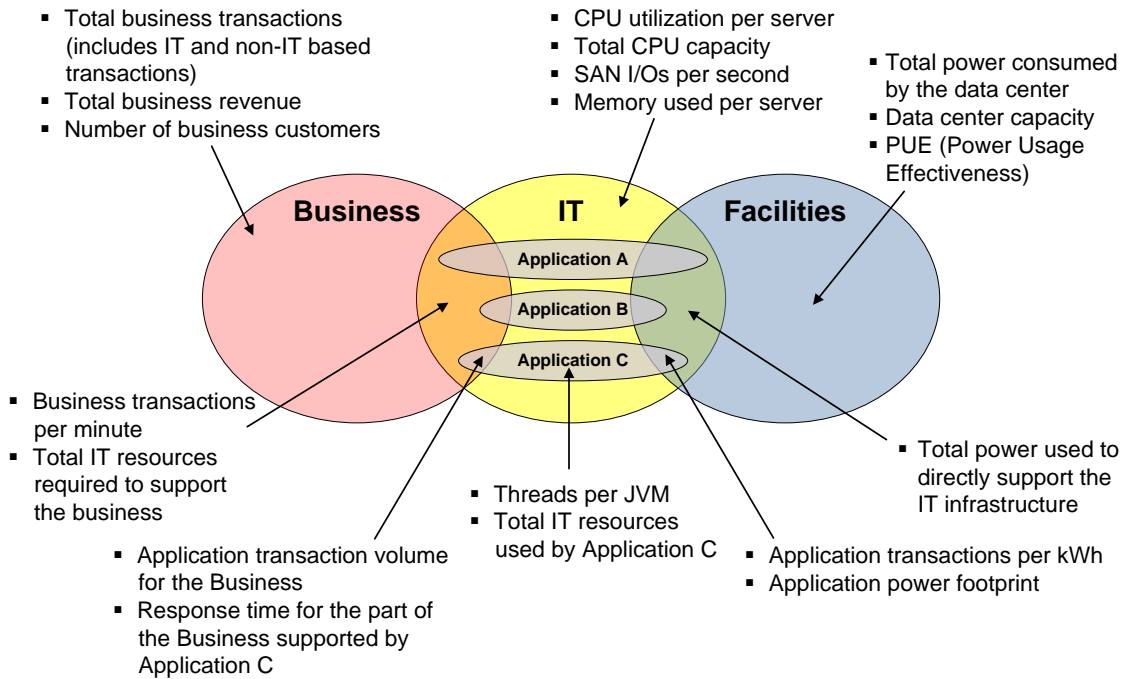


Figure 2. Capacity Planning Stack taxonomy.

The three main categories of metrics in the taxonomy are:

- **Business** - Metrics that describe the business as a whole (e.g., customer count and revenue) and the workload placed on the Digital Infrastructure (e.g., view account balance transactions). There may be one or more applications that directly support each Business.
- **IT** - These are the IT infrastructure metrics that capacity planners are most familiar with; for example, CPU and memory utilization, I/O rate and network bandwidth usage. Note that IT contains subcategories that represent individual Applications.
- **Facilities** - The data center hosting the IT equipment is characterized by a set of Facilities metrics that includes total power usage and PUE.

The three major categories are easy to recognize; Business, IT and Facilities. A Venn diagram is used to illustrate the taxonomy to highlight the dependencies and logical relationships between the major categories. For example, consider the following:

- $\text{Business} \cap \text{IT}$  - This is the intersection of the Business and IT categories. The intersection contains metrics that describe the IT resources used to support the Business' workload demand.
- $\text{IT} \cap \text{Facilities}$  - This is the intersection of IT and Facilities. Metrics quantifying data center (i.e., Facilities) resources used to directly support IT equipment are contained in this area.

The remainder of this section will develop the taxonomy step by step.

The simplified diagram on the right introduces the three major categories.

The diagram is annotated with sample metrics that describe and characterize each category.

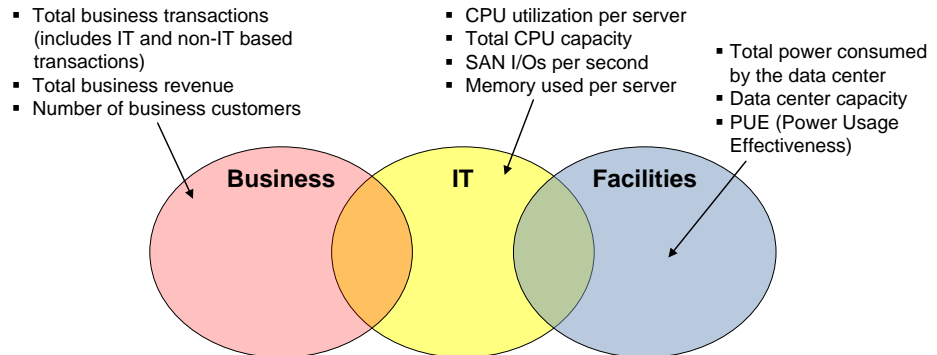
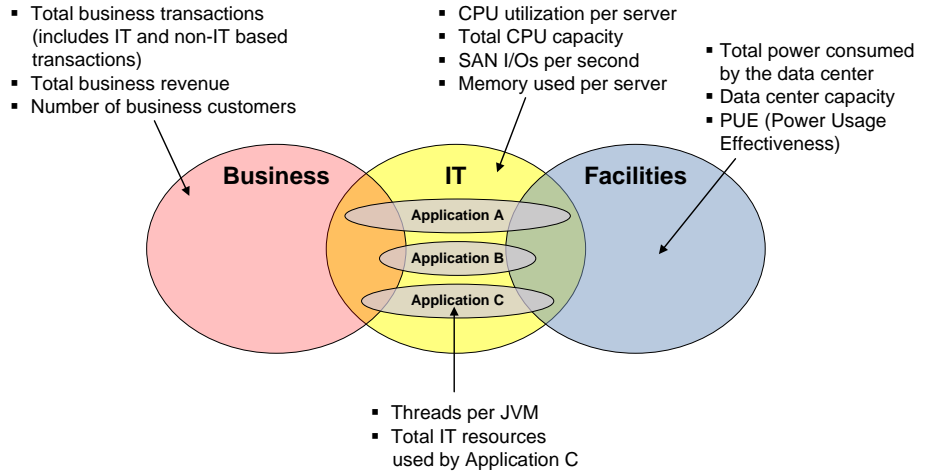


Figure 3. Three major taxonomy categories.

## Stack Metrics: A Taxonomy of Metrics Supporting the Capacity Planning Stack

The next step is to add subcategories for the applications that support the Business. The Application subcategories are fully contained in the IT category and also overlap the Business and Facilities categories.

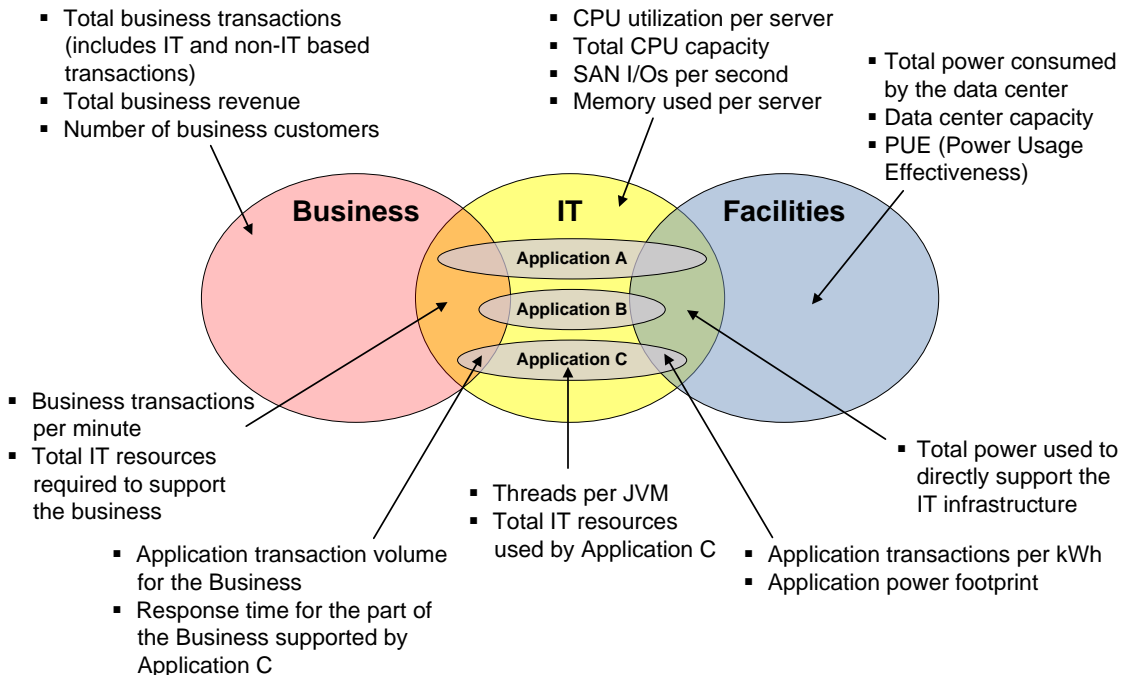


**Figure 4.** Add Application subcategories to IT.

In a fully populated taxonomy, the IT category would contain many Application subcategories; this diagram is simplified for clarity.

The non-Application space in the IT category represents IT overhead (i.e., the “cost of doing business” from the IT perspective).

The final step is to itemize and describe metrics that appear in the intersections of the major and minor categories.



**Figure 5.** Fully populated Stack taxonomy.

The Stack metrics contained in the fully populated sample taxonomy are described below in terms of set theory.

Set Notation	Description	Metrics
BIZ - IT	This is the area of the Business category that does not overlap with IT. Metrics in this area describe the total Business workload, cost and revenue.	<ul style="list-style-type: none"> <li>▪ Total business transaction volume across all supporting applications</li> <li>▪ Total business cost and revenue</li> <li>▪ Number of business customers</li> </ul>
$BIZ \cap IT$	Intersecting Business and IT reveals the metrics that describe the total IT resources used to support the Business.	<ul style="list-style-type: none"> <li>▪ Business transactions per minute</li> <li>▪ Total IT resources required to support the business</li> </ul>
$BIZ \cap APP$	The intersection of Business and Application represents IT metrics that are expressed in terms of the Business and its supporting Applications.	<ul style="list-style-type: none"> <li>▪ Application transaction volume for the Business</li> <li>▪ Response time for the part of the Business supported by the Application</li> </ul>
IT - BIZ - FAC	This is the area of the IT category that does not overlap with Business or IT. These are the metrics that describe the capacity and usage of the IT equipment.	<ul style="list-style-type: none"> <li>▪ CPU utilization per server</li> <li>▪ Total CPU capacity</li> <li>▪ SAN I/Os per second</li> <li>▪ Memory used per server</li> </ul>
$IT \cap FAC$	Intersecting IT and the Facilities shows the metrics that describe the total IT resources that are hosted and supported in the data center.	<ul style="list-style-type: none"> <li>▪ Total IT power consumption</li> <li>▪ Total IT space used (e.g., racks)</li> <li>▪ Total IT cooling requirement</li> </ul>
FAC - IT	Facilities-only metrics are contained in the portion of the Facilities category that does not overlap with IT.	<ul style="list-style-type: none"> <li>▪ Total power consumed by the data center</li> <li>▪ Data center capacity (power, space and cooling)</li> <li>▪ PUE (Power Usage Effectiveness) [TGG2012]</li> <li>▪ CUE (Carbon Usage Effectiveness) [TGG2010]</li> </ul>
$APP \cap FAC$	The intersection of the Application and Facilities contains metrics that describe Facilities resource usage specific to the Application.	<ul style="list-style-type: none"> <li>▪ Application transactions per kWh</li> <li>▪ Application power footprint</li> </ul>
APP - BIZ - FAC	This is the portion of the Application subcategory that does not overlap with either the Business or Facilities categories. This area represents Application metrics are specific to IT.	<ul style="list-style-type: none"> <li>▪ Threads per JVM</li> <li>▪ Total IT resources used by Application C</li> </ul>

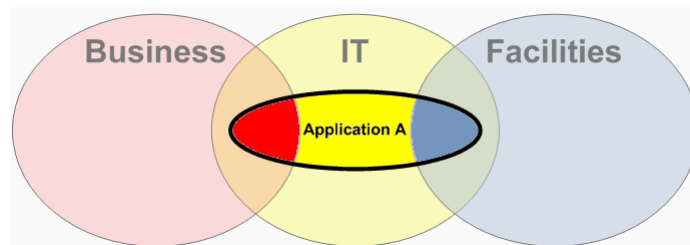
### 3.1 Benefits of the Stack Taxonomy

The four primary benefits derived from the taxonomy are:

1. The taxonomy identifies the three major categories of Digital Infrastructure metrics; Business, IT and Facilities.
2. The dependencies and relationships between the metric categories are clearly illustrated in the taxonomy's Venn diagram.
3. The taxonomy enables the capacity planner to partition metrics by application.
4. The various areas highlighted in the taxonomy's Venn diagram identify the metrics required by the Stack.

Figure 6 illustrates the advantages from Application A's perspective.

- The total IT resources supporting the application are in the Application A subcategory area (APP\_A).
- Application resources stated in terms of the business it supports are contained in the red portion of the Application A subcategory ( $APP\_A \cap BIZ$ ).
- The facilities resources (e.g., power, space and cooling) that directly support Application A are contained in the blue portion of the Application A subcategory ( $APP\_A \cap FAC$ ).



**Figure 6.** Application A metrics partitioned by the Stack taxonomy.

The ability to partition the set of Digital Infrastructure metrics into the categories and subcategories described by the taxonomy gives the capacity planner the complete set of building blocks required to develop a comprehensive plan for the Digital Infrastructure.

The three major categories (Business, IT and Facilities) can be viewed as representing semi-autonomous and de-centrally organized and managed groups within the Digital Infrastructure. The Capacity Planning Stack and metric taxonomy support a federated approach to Digital Infrastructure capacity planning. Instead of relying on a single capacity planning organization that spans Business, IT and Facilities, the Stack and supporting metrics promotes interoperability and sharing of information between the three silos [IRVI2009] [WIKI2014].

### 3.2 Gap Analysis with the Taxonomy

The taxonomy describes the perfect world which normally does not exist for the capacity planner. However, being able to identify the gaps between reality and the taxonomy's perfect world view enables a process for determining what data is missing and identifying the assumptions necessary to account for the discrepancies. It is precisely the gaps identified by the taxonomy that promote the required federated approach to Digital Infrastructure capacity planning.



## 4 Stack Taxonomy – Usage

This section will illustrate how the Stack taxonomy can be used to support a capacity planning exercise using the Capacity Planning Stack.

As an example, suppose a bank is planning to extend their retail banking services to include all states west of the Mississippi over the next two years. This expansion is expected to double their customer base. To make this example tractable, assume that we are only considering their online banking application. The number of new customers is the bank's Natural Forecasting Unit (NFU) (see [LO1986] and [REUL1987]); this is the way they think about and describe their business. The initial step in the capacity planning process is for the business owners to translate the expected increase in their customer base (NFU) into a corresponding increase in the number of business transactions. Based on the staged rollout, the business owners are expecting a 10% per quarter transaction volume increase in their online banking application over the next two years. The capacity planner is now tasked with the following:

- Predict the impact of this new business demand on their Digital Infrastructure
- Determine the most cost effective way to optimize service delivery over the next two years

The following sections will walk through the Stack's demand and feedback workflow to show how the taxonomy organizes and focuses the capacity planner during this planning exercise.

### 4.1 Demand: Business & Application

The first step is for the Business to estimate the projected workload volume. This demand is seen as (1) in Figure 7:

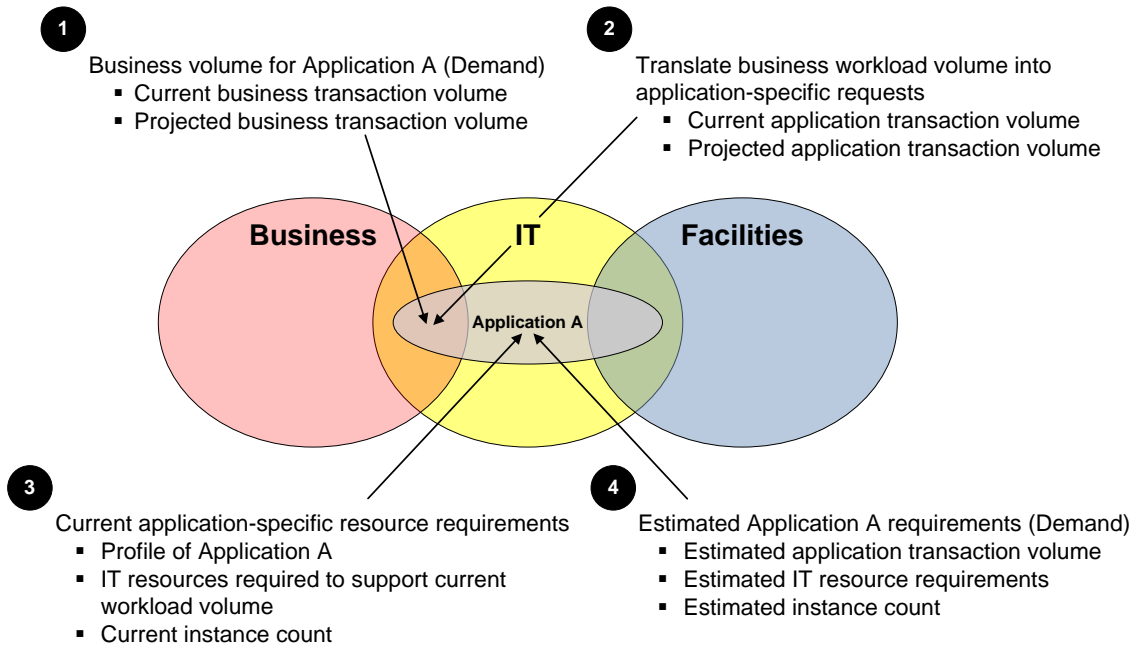
- The application planner will use business-level workload volume measurements to determine the current business volume.
- The new projected business workload volume will be based on a 10% quarterly growth over the next two years.

Given the current and projected business workload volume, the application planner translates business transaction volume into application-level volumes (see (2) in Figure 7). If multiple applications support the business, then application-level volumetrics would be collected from each application.

- Application-level transaction volumetrics are collected and used to determine the ratio of application to business transactions.
- The resulting application-level demand describes the projected application requests that will be passed as demand to the infrastructure level.

The application planner still has work to do: determining the demand that will be passed down to the infrastructure level of the Stack. The application planner must create a profile of the current application(s) supporting the business. The profile will contain the following (see (3) in Figure 7):

- Quantitative description of the physical and logical resources used to support the current application.
- Physical IT resource usage will be expressed in terms of "per application request". Physical resources include those managed by the infrastructure level; for example, CPU, memory, storage and network.
- Logical IT resource usage (e.g., number of application instances, threads or JVMs) is normally expressed in terms of "number of logical resources per active application request).



**Figure 7.** Stack taxonomy supporting Business and Application demand.

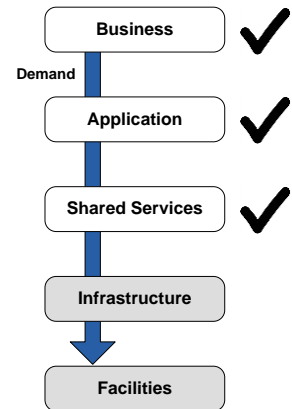
Once the profile is complete, the application planner can use it to create the demand that is passed to the infrastructure level. The application demand includes the following (see (4) in Figure 7):

- Estimated application transaction volume
- Estimated IT physical resource requirements
- Estimated IT logical resources (e.g., instance count)

Application-level modeling is generally used to estimate the required physical and logical resources.

At this point the application planner has transformed the business' projected transaction volume into the resource demand that will be passed to the infrastructure level. For simplicity, we are assuming that there are no Shared Services in this example.

The application demand will be used by the infrastructure level to determine the required IT infrastructure.

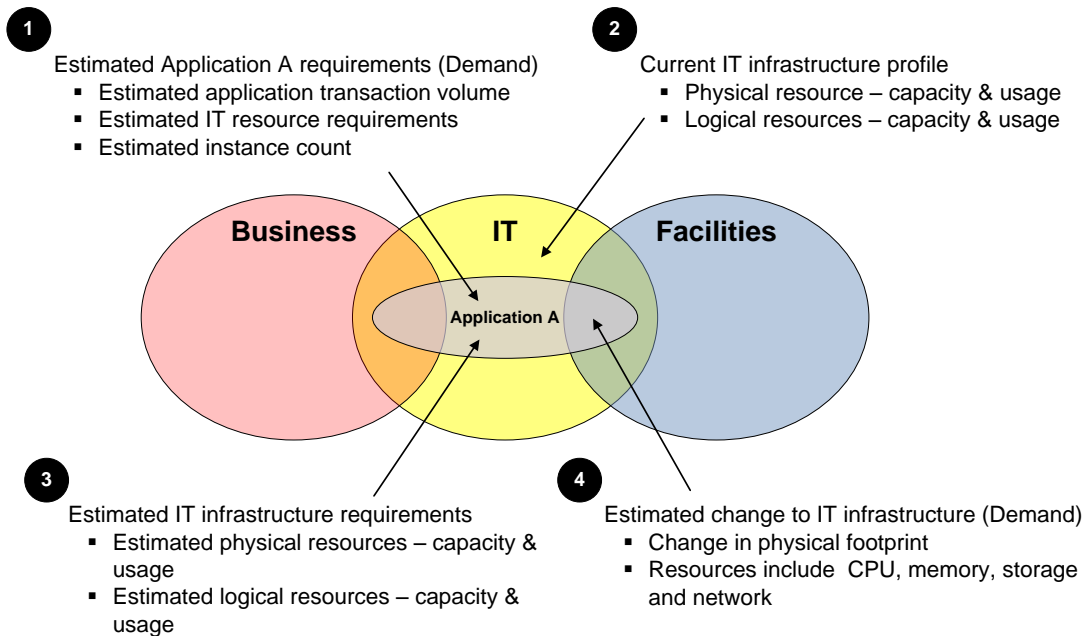


## 4.2 Demand: Infrastructure

The infrastructure planner receives the demand from the application planner (see (1) Figure 8 below). The application demand describes the physical and logical resource requirements for the application. The infrastructure planner is tasked with determining the physical and logical resources that will satisfy the projected application changes.

As a starting point, the infrastructure planner leverages a profile of the current IT infrastructure (see (2) in Figure 8). The IT infrastructure profile is a complete inventory and quantification of the physical and logical resources. This current state snapshot includes (at least) the following:

- Complete inventory of physical resources. This includes servers, CPUs, memory, storage and network.
- Quantitative description of physical resources that describes total capacity available and capacity used.
- Logical resource inventory and capacity usage metrics.



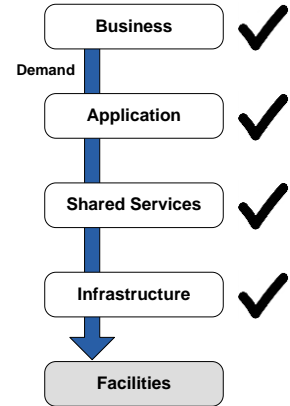
**Figure 8.** Stack taxonomy supporting Infrastructure demand.

Now that the infrastructure planner knows what is currently available (see (2) in Figure 8) and what is required (see (1) in Figure 8), infrastructure planning can begin. Recall the two goals of the capacity planner:

- Predict the impact of this new business demand on their Digital Infrastructure
- Determine the most cost effective way to optimize service delivery over the next two years

The infrastructure planner will use predictive modeling and other analytic tools to determine how to shape the current IT infrastructure into one that will satisfy the application demand (see (3) in Figure 8). This may require the addition of new hardware, repurposing existing hardware or upgrading current hardware.

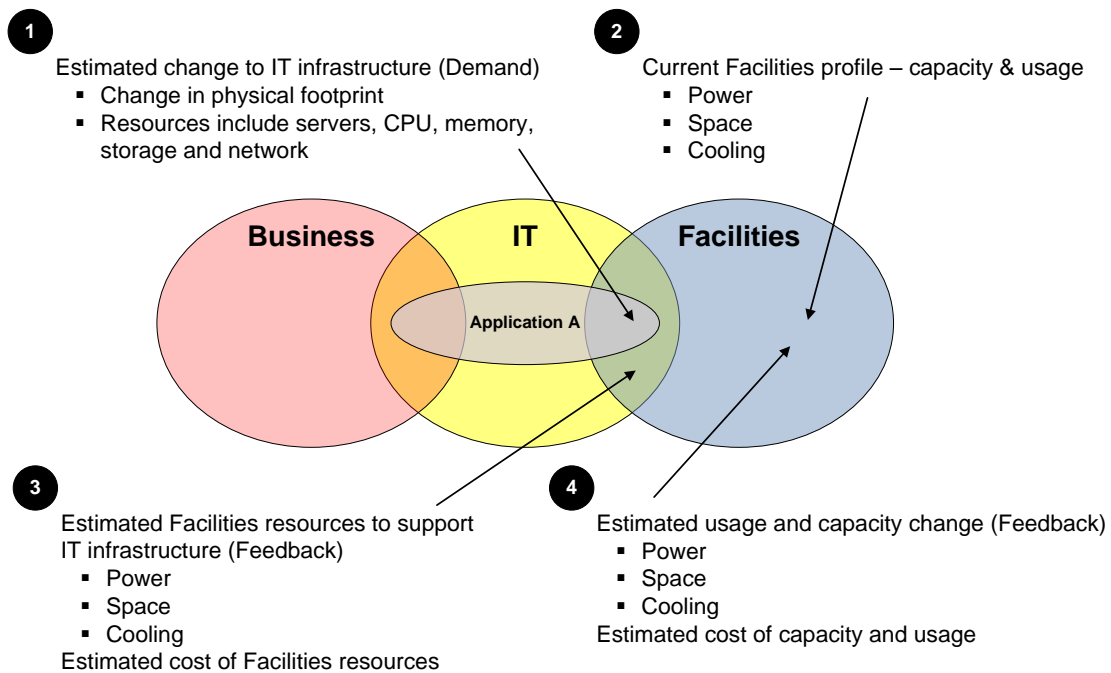
The final step of the infrastructure planner is to pass the infrastructure demand to the Facilities level (see (4) in Figure 8). For expediency, Infrastructure demand is expressed in terms of what is changing.



### 4.3 Demand: Facilities

The last step in the Stack’s demand workflow is for the facilities planner to receive, analyze and process the infrastructure demand (see (1) in Figure 9 below). The facilities demand describes the required change in the physical IT infrastructure. The facilities planner’s job is to determine (a) if the current data center has sufficient capacity to support the change and (b) develop an estimate of the cost to provide that support. If there is insufficient facilities capacity, a feedback loop will be initiated to communicate the discrepancy and estimated cost for remediation.

The facilities planner starts with a profile of the current facilities infrastructure (see (2) in Figure 9). The profile is a complete inventory and quantification of the physical data center resources available to support the IT infrastructure. The facilities planner’s current state profile includes (at least) the total capacity and current usage of the three primary facilities resources: power, space and cooling.



**Figure 9.** Stack taxonomy supporting Facilities demand.

The next step is for the facilities planner to determine if the change in IT infrastructure (see (1) in Figure 9) can be supported by the current facilities infrastructure (see (2) in Figure 9). Once this determination is made (either positive or negative) the facilities planner initiates the feedback workflow that contains an estimate of the facilities resources and cost (see (3) in Figure 9) required to support the “demanded” IT infrastructure.

If the data center does not have sufficient capacity (e.g., cooling cannot support the additional servers), then the facilities feedback will include the cost to add the required capacity (see (4) in Figure 9).

#### 4.4 Demand: Summary

It is interesting to note that each level in the Stack performs the same three basic steps:

1. What is today’s current capacity at level k?
2. How much of today’s level k capacity is used?
3. What do you need to satisfy level k-1’s future demand?

In addition to profiling the total resources available at each level, the Stack taxonomy encourages the partitioning of total resources by application. For example, this can be expressed in set theory notation for the infrastructure level:

Set Notation	Description
IT	This represents the entire set of physical and logical resources in the IT infrastructure.
$IT \cap APP$	The intersection of the IT and Application categories represents the subset of IT infrastructure resources that support a specific application.
$IT \cap FAC$	The intersection of IT and Facilities represents the space, power and cooling required to support the IT equipment.

Perhaps the most important long-term benefit of the taxonomy is to promote interoperability and sharing of information between the major categories within the Digital Infrastructure (Business, IT and Facilities). It is rare to find these three groups working in unison. However, the required sharing of information between Stack levels necessitates that there be some level of cooperation and communication. Furthermore, if information is shared, it must be packaged and structured into a form that can be utilized by at least two adjacent levels of the Stack.

- **Business & Application** These two Stack levels have a track record of working together. Application planners have experience translating business requirements into application-level resource requirements (i.e., demand).
- **Application & Infrastructure** This is another pair of Stack levels that have experience working together. In a number of cases, the source of metrics for these two levels comes from the same set of measurement tools.

- **Infrastructure & Facilities**      Sharing of information between Infrastructure and Facilities is new. Historically, these two organizations have not had much reason to collaborate. The Stack requires that Infrastructure and Facilities create a communication channel and common language that supports Digital Infrastructure capacity planning.

## 4.5 Feedback

The previous section demonstrated the value of the taxonomy for identifying the Digital Infrastructure metrics required to describe the demand that flows down the Stack. In this section we look at how the taxonomy benefits the feedback loop that travels back up the Stack.

In general, the downward flowing demand describes the resources required to support the high-level business demand. For example, the application planners passed the following demand requirements to the infrastructure level (see (4) in Figure 7):

- Estimated application transaction volume
- Estimated IT resource requirements
- Estimated instance count

As a response, the application level expects the following minimum feedback from the infrastructure level:

- Infrastructure resources required to satisfy the application demand
- Cost of infrastructure resources
- Time to deploy infrastructure resources
- Expected application-level transaction response time

In order to get a complete picture that spans the entire Digital Infrastructure, the application planners require the additional feedback that the infrastructure planners received from facilities:

- Description and cost of Facilities resources required to satisfy the application demand
- Time to deploy facilities resources

The feedback requirements for our example are shown in Figure 10.

1. Infrastructure receives feedback from facilities (see (1) in Figure 10) that describes what must be done within the data center to satisfy the infrastructure demand.
2. The application level receives feedback from infrastructure (2) that has a similar description of what was done to satisfy the application demand. In addition, the facilities feedback (1) is forwarded to the application level.
3. Business, sitting at the top of the Stack, will receive feedback from the application level (3) and the cumulative set of feedback from all lower levels ((1) and (2)).

In this example, receiving cumulative feedback from all lower levels in the Stack enables the application planners to understand and quantify the entire set of Digital Infrastructure resources required to satisfy their demand. Note that all feedback is focused on different segments of the application subcategory. As will be shown in the next section, cumulative feedback can be used to formulate efficiency metrics.

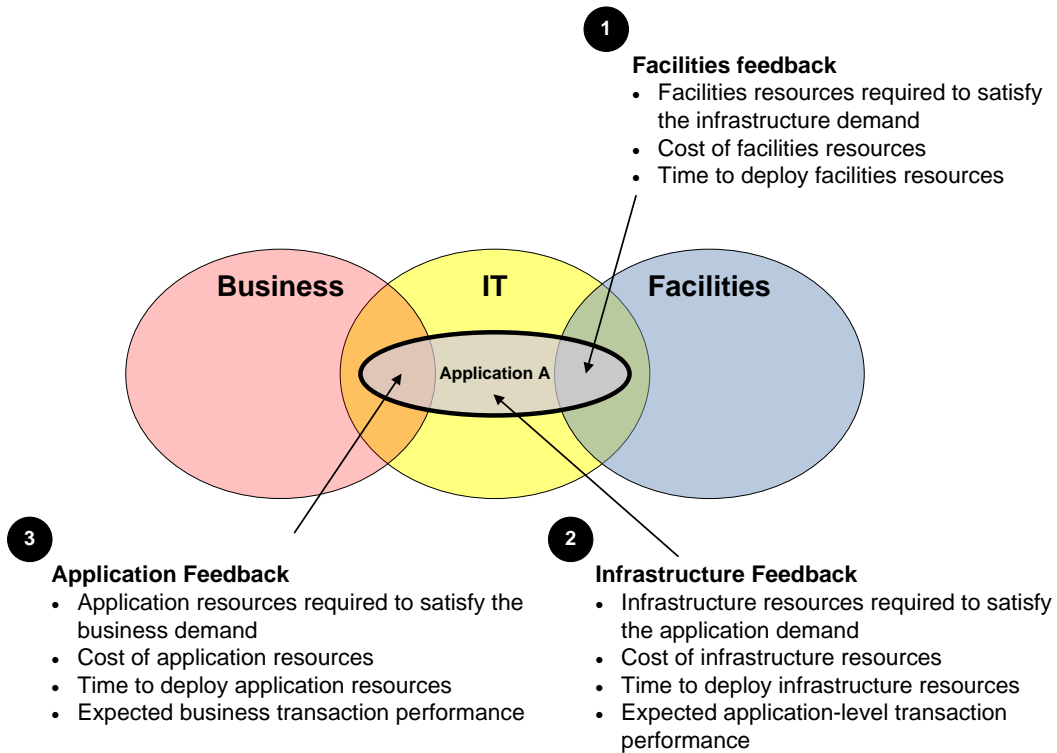


Figure 10. Taxonomy view of Stack feedback.

#### 4.6 Efficiency Metrics

The final set of metrics utilized by the Stack is the efficiency metrics. Efficiency metrics are used to report and track the long term trends and measures of success at each level. There are two types of efficiency metrics:

- **Shared** - These metrics describe the efficiency of the application relative to Business, IT and/or Facilities. Metrics from multiple categories are required to develop these measures of efficiency. Shared efficiency metrics leverage the cumulative feedback passed up through the Stack.
- **Silo** - Each taxonomy category and subcategory has its own set of non-shared (i.e., silo) efficiency metrics that describe how well they are managing their resources.

As an example, consider the computation of the business-level efficiency metric “Business transactions per Digital Infrastructure dollar” for Application A.

This is a shared efficiency metric that leverages the cumulative feedback that is passed up through the Stack. The Business level will have all the information it needs to compute the total cost efficiency of their workload.

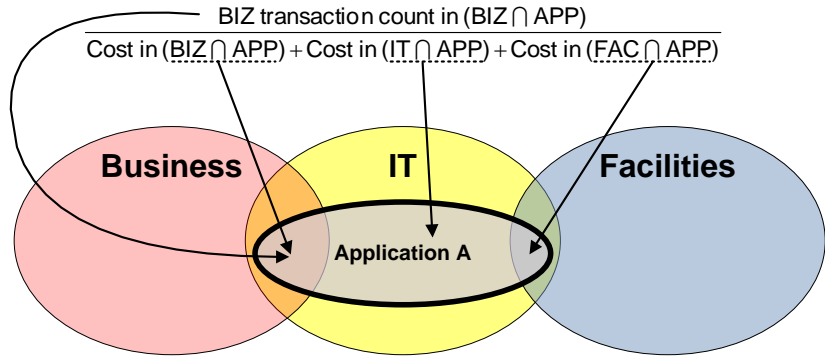


Figure 11. Stack taxonomy supporting efficiency metrics.

To push this example one step further, suppose there are N applications that support the business. The following could be used to compute the efficiency metric “Cost per Business Customer”:

$$\frac{\text{BIZ customer count in (BIZ)}}{\sum_1^N [\text{Cost in } (BIZ \cap APP_i) + \text{Cost in } (IT \cap APP_i) + \text{Cost in } (FAC \cap APP_i)]}$$

The following table lists a number of sample efficiency metrics per Stack level.

Stack Level	Sample Efficiency Metrics	Type
<b>Business</b>	Total business transactions per Digital Infrastructure dollar	Shared
	Business transactions per Digital Infrastructure dollar by application	Shared
	Total Digital Infrastructure cost per business customer	Shared
	Total number of business transactions	Silo
	Average business transaction response time versus SLA	Silo
<b>Application</b>	Application transactions per kWh	Shared
	Application power footprint	Shared
	Average application response time versus SLA	Silo
<b>Shared Services</b>	Average database resources used per application request	Silo
	Average database response time versus SLA	Silo
<b>Infrastructure</b>	Average CPU headroom (i.e., unused capacity)	Silo
	Average power use per server	Shared
<b>Facilities</b>	Power Usage Effectiveness (PUE)	Silo
	Carbon Usage Effectiveness (CUE)	Silo

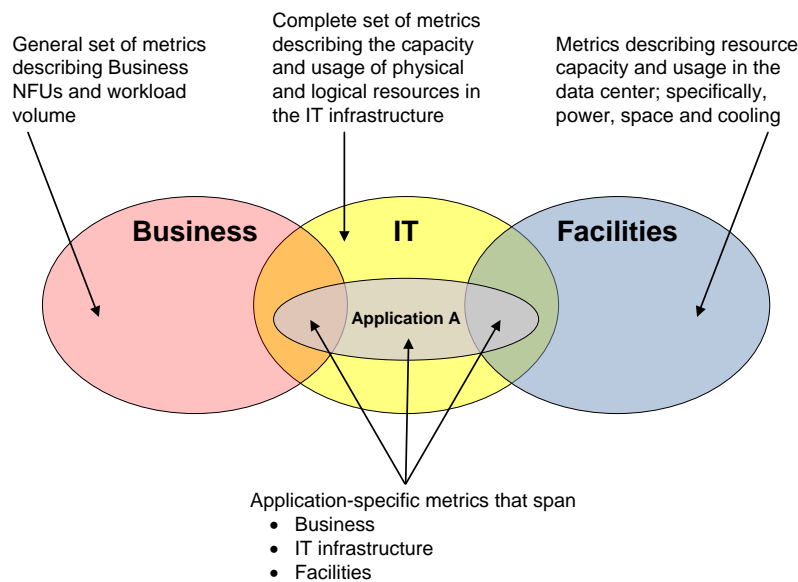


## 5 Summary

This paper introduced a taxonomy for characterizing the metrics required to support the Capacity Planning Stack. As capacity planners, we are faced with a multitude of metrics that come from a wide variety of data sources. The intent of the taxonomy is to organize and focus the capacity planner's selection and use of those metrics during a planning exercise.

As shown in Figure 12, the taxonomy has three major categories; Business, IT and Facilities. In addition, the IT category has subcategories that represent individual applications.

The value of the taxonomy is that it organizes and illustrates the complete set of Stack metrics. The three major categories are straightforward and familiar. Complexity is introduced by the need to describe the logical relationships and dependencies between Business, IT and especially Facilities.



**Figure 12.** Stack taxonomy categories and subcategories

The full set of metrics can be difficult to collect and correlate, especially since the major Stack categories are managed by separate organizations. The taxonomy clearly identifies the touch points between Business, IT and Facilities and shows precisely where sharing and interoperability must occur. The most challenging touch point is often between IT and Facilities. Historically, it has not been necessary for these two organizations to work together. However, with today's pressing need for capacity planning to include the entirety of the Digital Infrastructure, it is imperative for these organizations to not only learn how to work together, but also communicate effectively.

The Stack taxonomy equips the capacity planner with a structured way to think about, organize, communicate and collect the metrics that provide the most value for the Digital Infrastructure. Additionally, the requisite interoperability and sharing of data between Business, IT and Facilities paves the way for federated capacity planning across today's Digital Infrastructure.

## 6 References

- [CAMB2014] Cambridge Dictionaries Online, <http://dictionary.cambridge.org/us/>.
- [DICT2014] Dictionary.com, [dictionary.reference.com/](http://dictionary.reference.com/).
- [IRVI2009] Reed E. Irvin, "*Getting From Point A to Point B: What it Means to Take a Federated Approach*", Information Management Journal, May/June 2009.
- [LO1986] T. L. Lo and J. P. Elias, "*Workload Forecasting Using NFU: A Capacity Planner's Perspective*", CMG 1986 International Conference, December 1986.
- [MERR2014] Merriam-Webster, [www.merriam-webster.com/](http://www.merriam-webster.com/).
- [OXFO2014] Oxford Dictionaries, [www.oxforddictionaries.com/us/](http://www.oxforddictionaries.com/us/).
- [REUL1987] John M. Reyland, "*The Use of Natural Forecasting Units*", CMG 1987 International Conference, December 1987.
- [SPEL2013] Amy Spellmann and Richard Gimarc, "*Capacity Planning: A Revolutionary Approach for Tomorrow's Digital Infrastructure*", CMG 2013 International Conference, November 2013.
- [TGG2010] The Green Grid, "*Carbon Usage Effectiveness (CUE): A Green Grid Data Center Sustainability Metric*", 2010, [www.thegreengrid.org](http://www.thegreengrid.org).
- [TGG2012] The Green Grid, "*PUE: A Comprehensive Examination of the Metric*", 2012, [www.thegreengrid.org](http://www.thegreengrid.org).
- [WIKI2014] Wikipedia, "*Federated Architecture*", 2014, [http://en.wikipedia.org/wiki/Federated\\_Architecture](http://en.wikipedia.org/wiki/Federated_Architecture)