

White Paper

# Overcoming Hidden Costs in the Data Center

The undeniable appeal of free cooling  
After the power draw from IT equipment,  
cooling the data center frequently  
consumes the largest amount of energy.  
As such, it represents a significant  
opportunity to improve efficiency.  
Depending on the regional climate, free  
cooling has proven to reduce costs by 10  
to 50 percent or more.

Anxieties about outside humidity and pollutants that once stymied wide adoption have faded. Studies conducted by Lawrence Berkeley National Laboratory in 2007 examined concerns about the impact on IT assets from particulate contamination and humidity controls associated with air-side economizers. Early tests concluded proper design and maintenance results in energy savings without compromising IT assets. Although initial results are favorable, LBNL recommended further studies to examine longer-term effects.

With early data center studies dispelling fears of air-side economization, and with corporate and social sentiment pressing for both efficient and clean energy consumption, today it is almost unthinkable to build a new data center without free cooling. Dr. Kenny Gross, Distinguished Engineer at Oracle and team leader for the System Dynamics Characterization and Control team in Oracle's Physical Sciences Research Center in San Diego cautions, "Data center operators should not naively assume that with every degree they raise the ambient temperature in the data center, they will save money. New challenges with data center operations have emerged with free cooling."

This paper examines the hidden costs linked to free cooling and offers some suggestions about how data center operators can meet these challenges.

## Section one

### **The quest for sustainable growth**

A key tenet to sustainable data center practice is to capitalize on the energy savings from efficient IT equipment designed to reduce energy loads, which in turn reduces the energy required to cool the equipment. For example, for every kilowatt (kW) saved by a server, an additional kW is saved in removing the heat from that server (assuming a PUE of 2). For data center operators, this seemingly straightforward demand for more efficient IT equipment to curb energy costs is offset by the demand for higher density servers. Such servers offer more processing power, thereby increasing productivity (i.e. the number of computations) in the same space. Although this offers a sustainability advantage because the enterprise can increase productivity while deferring the cost of data center expansion, the higher density ultimately requires more energy, even though the kW per computation is lower.

Despite conflicting business pressures to "do more with less," data center professionals are finding solutions by adopting highly efficient practices, which extend to data center design as well as equipment purchase decisions. In response, equipment manufacturers have made significant strides to meet customer demands for sustainable growth. Server manufacturers are now taking advantage of energy-efficient components, such as more efficient processors, fans and power supplies. Variable speed fans offer one example. Slowing fan speeds when CPU and memory workloads are lower can reduce the fan's

power draw significantly. This is due to the fact that power varies with the cube of fan speed. In other words, a fan operating at half speed is only using one eighth the power it would at full speed. This may not seem like much on an individual basis, but there are a number of fans inside of a server and the incremental energy savings multiplied across a large number of servers adds up over time.

In another example, manufacturers have developed blade servers to meet the demand for more processing power in less space. Blade servers are practically stripped down to bare metal to enable more units to fit in a rack, a design change that also included the removal of disposable air filters. However, these air filters were precautionary in server design and were ultimately deemed non-essential. It was a rare customer, indeed, who had changed an air filter in a server.

As server manufacturers advanced technologies to meet new demands, data center operators themselves were being driven by economic pressures to lower their operating costs. Use of efficiency metrics such as PUE and WUE spread throughout the industry as operators began scrutinizing how to reliably extend capacity while reducing energy usage. One of the most popular methods adopted to help improve a data center's efficiency benchmarks is air-side economization, or "free cooling."

Free cooling emerged as an effective passive design<sup>1</sup> method to reduce the energy consumed by the facility infrastructure to replace the warm air exhausted by IT assets. Essentially, data centers located in areas with cool climates could offset the costs of traditional cooling methods (e.g., refrigeration) by using the cooler outside air. Results were undeniably appealing: free cooling has proven to reduce power draw by 10 to 50 percent or more, depending on geographic location.

## Challenges with warmer air and free cooling

### Warmer air gets the cold shoulder

It was understood early on that raising the temperature of the data center would yield diminishing returns in terms of energy savings, but in the absence of real-world experience the actual costs and circumstances were unknown. The System Dynamics Characterization and Control (SDCC) group at Oracle discovered and then quantitatively characterized hidden costs that occur inside the server as a consequence of raising ambient temperatures too high. Following are a few examples:

1. All enterprise servers are now equipped with variable-speed fans inside. The fan motor power has become a significant proportion of the overall server power budget. In many cases the energy consumed by the fan motors exceeds that of the CPU chips. As noted above, the fan motor power draw increases with the cube of the fan RPMs. Consequently, as the ambient air temperature increases, the fan motor consumes significantly more energy to ensure the critical CPUs and memory modules stay within the acceptable operating temperature range.

2. Until just a few years ago, enterprise computing chips consumed energy only for "switching power" in flipping gates as computations are performed. Leakage power, which is considered wasted energy because it does not support the computational workload, had traditionally been negligible. When the ambient temperature is cool enough to keep the CPU chips lower than the threshold for leakage power, all power consumed by the CPUs goes to computation. But as soon as the ambient temperature rises to the point where the CPU chips enter the "leakage zone," leakage power grows exponentially with further increases in ambient temperature. In the most recent enterprise servers, because of relentless miniaturization of microprocessors, leakage power has become significant, using up to 30% of the server's power budget. This is projected to continue growing for future generations of servers.

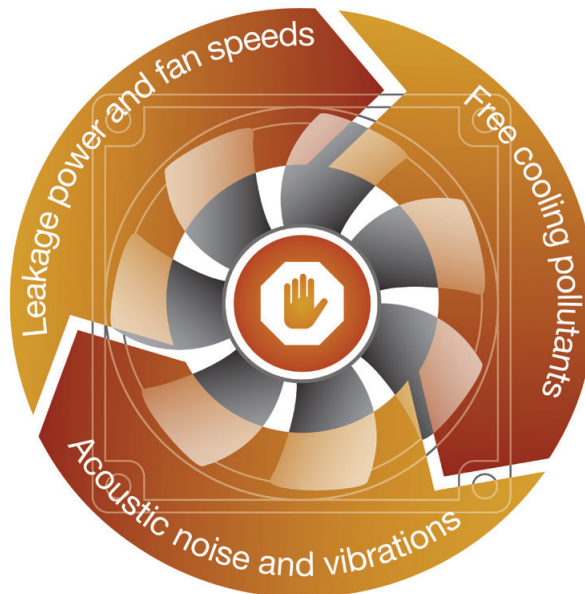
3. An additional challenge with increasing ambient temperatures arises from acoustics. The acoustic noise from servers increases with the 5th power of the fan RPMs.

<sup>1</sup> Passive design implies that energy-consuming mechanical components, like pumps and fans, are not used. Source: [http://en.wikipedia.org/wiki/Passive\\_cooling#Passive\\_cooling](http://en.wikipedia.org/wiki/Passive_cooling#Passive_cooling)



This creates the likelihood of significant customer dissatisfaction for recent generations of servers. Acoustic noise in data centers is ultimately constrained by OSHA standards to assure “a safe and healthful workplace”<sup>2</sup> for the systems engineers in the data centers. Although humans can (and now must) wear protective passive and active noise mitigation, this is not without workplace challenges because it becomes difficult for two or more engineers to communicate with one another and with colleagues outside of the data center. Moreover, the SDCC group at Oracle has discovered and quantified that the acoustic noise introduces an additional energy penalty via its impact on disk performance when the servers inside the data center are running IO-intensive workloads.

Oracle’s Dr. Gross explained there is a three-way relationship between ambient air temperature outside the server, fan motor power and chip leakage power that contributes to diminishing returns on overall data center energy savings as the ambient temperature rises.



#### Cost of failure

The bottom line is not that free cooling is not effective, but simply that employing it requires the most efficient real-time optimization tactics, which in turn rely on fine-grained visibility into device and system performance that accounts for both environmental and server thermal conditions. The above examples illustrate the impact to servers (and operating costs) by increasing the ambient air temperature of the data center. Equally significant is the impact of ambient air quality and humidity.

#### The new threat of silver sulfide

A new risk to “going green” in the data center is the unexpected consequence of “turning silver.” Silver sulfide ( $\text{Ag}_2\text{S}$ ) is a pollutant born of reduced sulfur particles in the air. When it combines with a certain level of moisture in the atmosphere, then the silver contained in the soldering inside the electronics begins to corrode. The corrosion eventually sprouts metal “whiskers,” actually thin oxidized particles that protrude from the soldering surface. These whiskers are electrically conductive and cause shorts and arcing in electrical equipment. Whiskers take time to grow and accumulate, and by themselves are not a hazard. The hazard occurs under certain environmental conditions—when there is enough moisture in the air to cause conduction between the particles. According to numerous reports, including Energy Efficient Thermal Management of Data Centers, “failures can be a rapid event, on the timescale of minutes or seconds.”<sup>3</sup>

The whisker phenomenon is not new. In fact, whiskers in electronics were reported as early as the 1940s when the discovery was coined “tin whiskers” because

<sup>2</sup> OSHA Law & Regulations <http://www.osha.gov/law-regs.html>

<sup>3</sup> Energy Efficient Thermal Management of Data Centers, Yogendra Joshi and Parmod Kumar, 2012, p 113.

of their growth from pure tin coatings in electronic hardware. Tin whiskers were abated by introducing lead to create an alloy that would resist corrosion. In addition, occurrences of “zinc whiskers” in data centers were reported around the turn of the millennium, appearing galvanized surfaces, such as the steel plates on the underside of floor tiles.

Today it is not tin or zinc, but the occurrence of silver whiskers that presents a menacing challenge in the data center. Consider these three factors:

1. Material change – Silver has fairly recently replaced the lead-enhanced alloys in electronics. The Restriction of Hazardous Substances (RoHS) Directive, which took effect in the European Union in 2006 and was followed by similar legislation in China and the United States, prohibits a “covered electronic device from being sold or offered for sale”<sup>4</sup> due to the presence of lead or other certain hazardous substances. Thus, there is more silver in facilities with newer equipment.

2. Silver properties – Silver has the highest electrical and thermal conductivity of any metal.<sup>5</sup> Given its propensity to grow whiskers, one of the main benefits of using silver in electronics (conductivity) is now also a detractor. Furthermore, silver easily reacts with airborne pollutants, and the combination of sulfur gases, higher temperatures and humidity may accelerate its corrosion.

3. New practices – There is a higher probability of silver whiskers occurring in data centers with free air cooling and warmer air, especially if the data center is located in a region prone to increases in hydrogen sulfide (H<sub>2</sub>S), “which comes from organic decay, combustion processes, volcanic activity, and manufacturing sources such as paper mills, sewage plants and high sulfur packaging materials.”<sup>6</sup>

The hidden costs of free cooling pose a formidable threat to system availability and business continuity, but there are ways to mitigate the risks.

David Gallaher, manager of IT Services at National Snow and Ice Data Center (NSIDC), recalls, “Under-floor cooling is the worst for zinc whisker contamination. I’ve seen the underside of floor tiles that are so thick with zinc whiskers that you could run your finger through it and make patterns. Since NSIDC is not using under-floor cooling, we don’t have it. But, I have experienced this before [at another data center]—three times! Each time it had cost as much as US\$200,000 to replace and decontaminate a computer room!”

## Section two

### Beyond the BMS

The state of the art in data center technology in 2007, when the Lawrence Berkley National Laboratory did its studies, is not the state of the art today. Server compute power and rack density continue their unrelenting upward trend. While air pollutants may not have been a significant hazard to data centers in 2007 (as introduced above), advancements to server technology and changes in data center practices may inadvertently heighten sensitivity. In addition, there are changes in environment, such as new or varied concentrations of pollutants and other physical characteristics that pose risks to data centers that rely on significant amounts of outside air:

- **Smoke** from fires in the area can carry high levels of both particulates and corrosive gases.
- **Dust**, while always a concern, can become a larger problem during extremely dry periods, as have been recently experienced in the United States.
- **Volcanic ash** may not seem like a threat to areas where volcanoes are not prominent, but as seen in Northern Europe in 2010 and in New Zealand in 2011, ash can travel great distances and pose a significant hazard to any system that circulates outside air.
- **Gases** in changing concentrations might be indigenous to a region or from man-made sources.

---

<sup>4</sup> Restriction on the use of Certain Hazardous Substances (RoHS) in Electronic Devices  
[http://www.dtsc.ca.gov/Hazardous Waste/RoHS.cfm](http://www.dtsc.ca.gov/Hazardous%20Waste/RoHS.cfm)

<sup>5</sup> Overview of the Use of Silver in Connector Applications, Marjorie Myers, Interconnection & Process Technology, Tyco Electronics, February 5, 2009, page 1.

<sup>6</sup> Overview of the Use of Silver in Connector Applications, Marjorie Myers, Interconnection & Process Technology, Tyco Electronics, February 5, 2009, page 4.

- **Humidity** in and of itself may not be a concern, but in combination with other pollutants, such as dust, is implicated in studies of electronic hardware failures due to copper creep corrosion.
- **Vibration** hinders disk performance and taxes energy consumption.

### New types of monitoring

Awareness of the risks associated with outside air is essential when managing free-cooled data centers, even if the free-cooling occurs only on a part-time basis. While monitoring all possible environmental factors may seem daunting, with the right system in place this proposition is not as extreme as it appears.

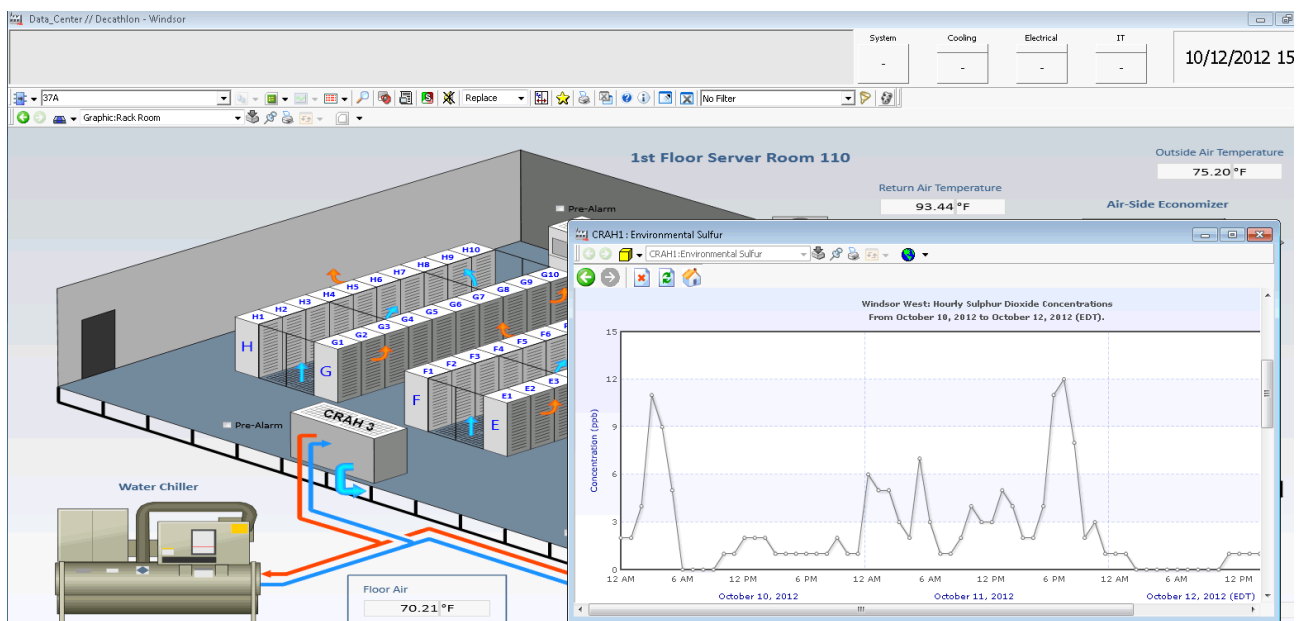
Data center facility monitoring traditionally centers on a Building Management System (BMS). These systems are designed to communicate with and manage a well-defined set of environmental and safety equipment inside a typical building.

A BMS may be limited in incorporating new sources of data, or modeling the entire data center environment.

Increasingly, data center managers are looking to Data Center Infrastructure Management (DCIM) systems to monitor, measure, model and manage their data center operations. DCIM systems offer several advantages over a traditional BMS. Data is collected from several sources including energy resource, power quality and distribution, mechanical equipment, electrical equipment, IT and network systems, and includes all types of physical and environmental factors. In some cases, data collection from these sources is limited by the communication protocols the DCIM system supports. Protocol flexibility is essential to future-proofing any DCIM system, thus those built on open platforms that support protocol plug-ins offer a significant advantage.

A DCIM system with an extensive and flexible protocol interface can collect data from a wide range of non-traditional data sources to perform environmental and physical monitoring beyond a BMS. For example:

- **Air quality analyzers.** Specialized analyzers are available for a number of gaseous and particulate analyses. Also, general purpose spectrographic analyzers can be connected via sophisticated interfaces such as OPC-ADI.
- **Weather stations.** With free cooling, environmental factors outside of the data center are as critical as the environmental factors inside the facility. In addition to outside air temperature and humidity, wind direction and speed should be incorporated into the DCIM inputs.



Decathlon™ IT Analytics

- **Noise and vibration data.** Physical attributes, such as noise and vibration within a server, rack or entire data center, can impact hard drive performance and energy efficiency.

To illustrate, consider the ease of integration with external sources of data for a DCIM system. For example, weather information can be read directly from feeds from the National Weather Service. Even pollution quantities can be read online in many regions.

Armed with contextually relevant information about current and forecasted weather and pollution, data center operators can make prudent decisions about using free cooling over a given time period. A good DCIM system provides much or all of the analysis to make this decision. The analysis can be presented as decision support, or some operations can be automated. For example, such a system would look for the pre-defined combinations of humidity and air quality that would minimize the risk of equipment damage. Ideally, it would continuously monitor the available data in case conditions change, and alert the operator if a switch back to closed-loop cooling is warranted.

### Optimizing energy usage

Comprehensive, fine-grained data can help address concerns about the environmental impact to data center operations, but what about the effects of running at a higher temperature? As stated earlier, increasing the ambient air temperature in a data center may yield energy savings, but those savings may start to reverse themselves if server response to the temperature changes is not understood and accounted for.

Some DCIM systems can read data directly from the servers themselves. Most modern servers can provide data such as inlet temperature, fan speed, CPU utilization, and other system-related parameters. Collecting and storing server- or device-level data, along with environmental data for root cause analysis is a good start. A good DCIM system will also provide the expert decision support that will help turn all types of data captured and stored in a time series format into contextually relevant information that provides real insight.

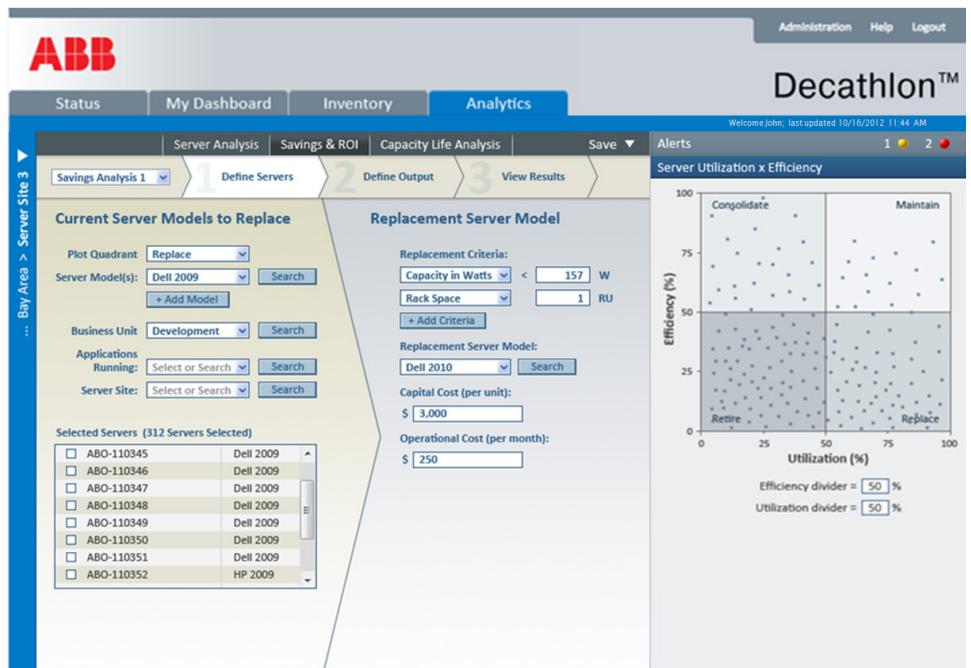
Total cooling power is a combination of facility cooling and internal server fan power. There is a “sweet spot” at which the ambient (inlet) temperature of the data center is high enough to stay within ASHRAE guidelines but low enough not to increase the draw from the servers’ internal fans. This sweet spot will change from data center to data center, due to the different hardware makeup of each facility.

In theory, given enough flexibility, data center operators could find this spot by varying the temperature and monitoring power usage. In practice, however, this is difficult to do systematically in a production facility. The current IT load on the data center will also impact the results, since lightly loaded servers may not have a large delta-T from inlet to outlet, meaning that the inlet temperature can be higher before the fans kick in. This means that the sweet spot may vary from hour to hour.

Given these many variables, it is not enough simply to have real-time data. Rather, access to expert decision support, that is, contextually relevant information based on holistic, high resolution data and analysis, is crucial to determining optimal temperature settings that will yield the highest ROI for each data center environment.

### Finding the outliers

Most data centers contain a variety of types and vintages of assets that give a data center its unique temperature characteristics. It would be very unlikely that all servers in a data center have exactly the same power usage characteristics. At the very least they will vary by model.



Decathlon™ IT Analytics

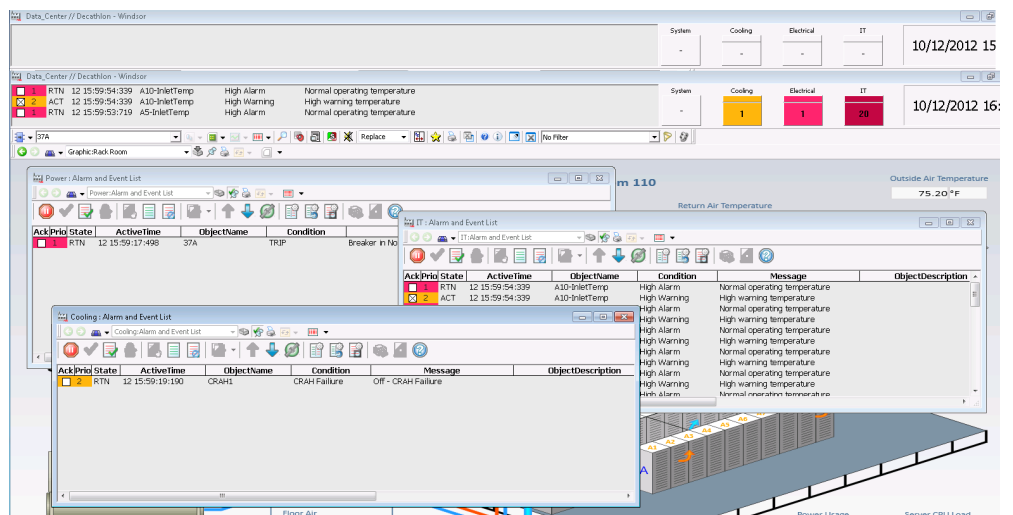
In such a scenario, some servers will be the “worst offenders” in terms of energy use. Finding and eliminating these servers could pay significant dividends over time. One way to do this is to identify potential savings based on static measurements of individual data center models using analytic tools and metrics such as PAR4.

While this technique is extremely powerful in identifying the best optimization targets, capturing real temperature, power, CPU usage and fan speed data will allow further analysis of the performance of individual servers, and help identify further potential savings.

## Section three

### Dynamic data center

Today's data center is a complex production facility, more similar to a chemical plant, power plant or other mission critical facility, than a simple building with a room full of servers. Even some of the current “smart building” technologies do not fully address the high reliability and performance needs that data centers increasingly demand.



Decathlon™ Alarm Management



The drive to reduce energy costs and minimize environmental impact leads to new designs that require data centers to straddle a fine line between energy efficiency and fault tolerance. Systems used to manage these facilities must be highly responsive, highly intelligent, and highly extensible. The support given to the data center operator must be quick and precise. The controls they use must also be secure and structured in a way to minimize response time and maximize effectiveness.

### Riding through a failure

Recent data center surveys show that approximately one-third of data centers asked had experienced a power outage within the year prior to the survey. The time between a failure of the data center cooling system and the IT equipment reaching critical temperatures is called the ride-through time and is typically in the range of about 2 minutes or less. Ride-through time decreases as the starting ambient temperature of the data center is increased. This means that the time to respond to the loss of the cooling system becomes increasingly important as data centers look for ways to save energy.

When a power failure causes the loss of cooling, systems are usually already in place to activate backup generators within the required time. But what if the failure is due to a device failure? What if free air cooling is suddenly not available because the outside air vents have unexpectedly closed due to a fire alarm or some unexpected environmental circumstances?

Typically, the options to handle these types of emergency situations have been few; either wait for the equipment to overheat and shut down (and hope that no warranty problems will be incurred), or perhaps start powering equipment down manually. Neither of these approaches allows for a systematic reduction in load that could potentially increase ride-through time, or for the transfer of IT load to an alternate location prior to a complete data center failure.

Data center automation enables these capabilities by predefining response scenarios that take the guesswork out of managing emergency situations. For example, if a critical cooling component fails (e.g., an evaporative cooling unit) due to a breaker trip, it may not be clear right away if the problem can be easily rectified or whether the backup cooling can be brought online quickly.

With automation, data center operators can choose from a selection of options that will bring the data center to a lower power consumption state, thus increasing the ride-through time and potentially avoiding a full shutdown. These might include any of the following:

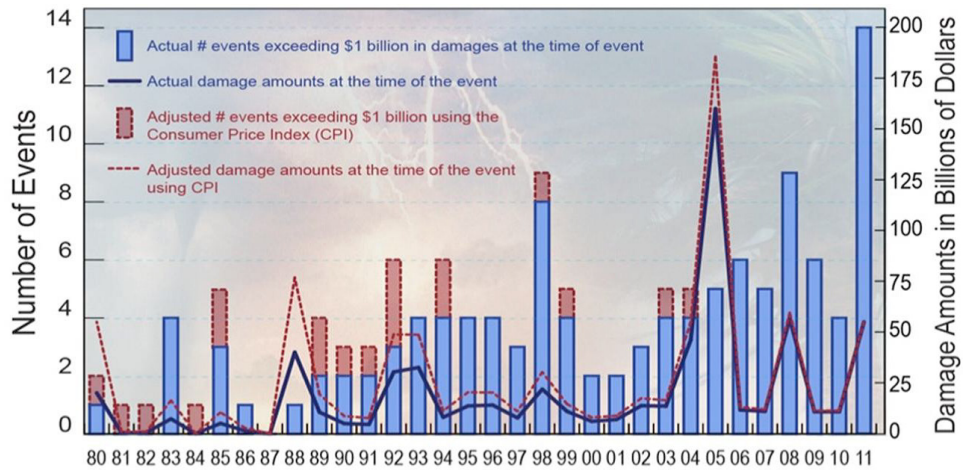
- Selectively shut down non-essential equipment
- Immediately initiate power-capping in all essential servers
- Initiate transfer of some IT load to a D/R site (if the system is highly virtualized)

Once the system has begun to recover, the DCIM can reverse the changes automatically.

In the case of a lights-out facility, the DCIM system can automatically initiate these steps after a cooling failure. Sequences may be triggered based on temperature and/or time thresholds.

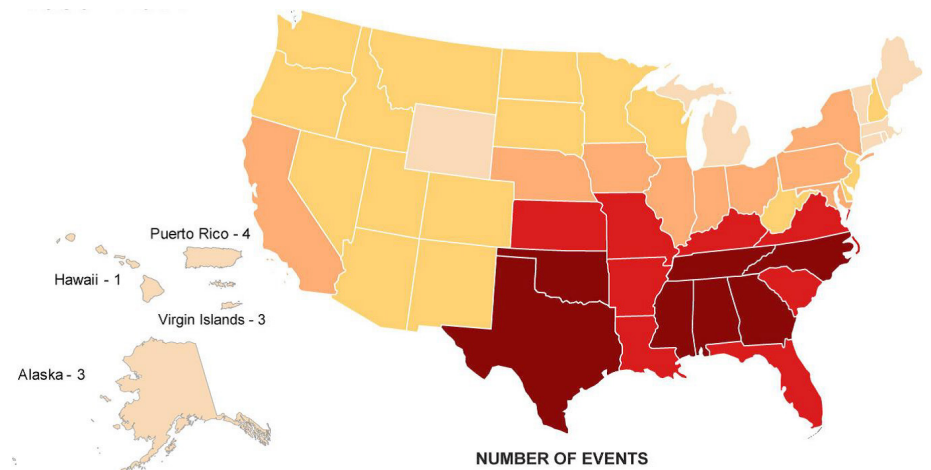
## Preparing for disaster

According to the National Climatic Data Center, the number of “billion dollar” natural disasters in the US averages more than four per year.



Source: National Climatic Data Center, NOAA, USA (2012)

The causes of these disasters are varied – hurricanes, tornados, floods, etc. – but the impact is always the same: major disruption of infrastructure and people’s lives. Certain regions have a higher likelihood of experiencing these events than others— as much as 30 times higher than in less-affected regions.



Billion Dollar Weather / Climate Disasters 1980 – 2011  
Source: National Climatic Data Center, NOAA, USA (2012)

Data centers can prepare for some climate-related events, but the severity of any given event is rarely known in advance. Because of the time it takes to initiate emergency action, it may be that in many cases, the only option available is a full data center shutdown. The point at which this decision needs to be made depends entirely on how quickly such action can be executed.

Even if the data center in question is outside the impacted area of a climate-related disaster, it can still be affected. Flooding can damage electrical and communication infrastructure used by a data center located miles from the flood itself. Roads may be blocked, making it difficult to keep diesel generators fueled or for staff to get to their jobs. Worse, staff may be dealing with the impact of the disaster on their own lives.

All of these scenarios point to the need for automated data center disaster response. A truly dynamic data center, using DCIM automation, can shorten the amount of time it takes to move IT load to a disaster recovery data center and thus allow the local site to delay critical decisions and possibly avoid taking some actions altogether.

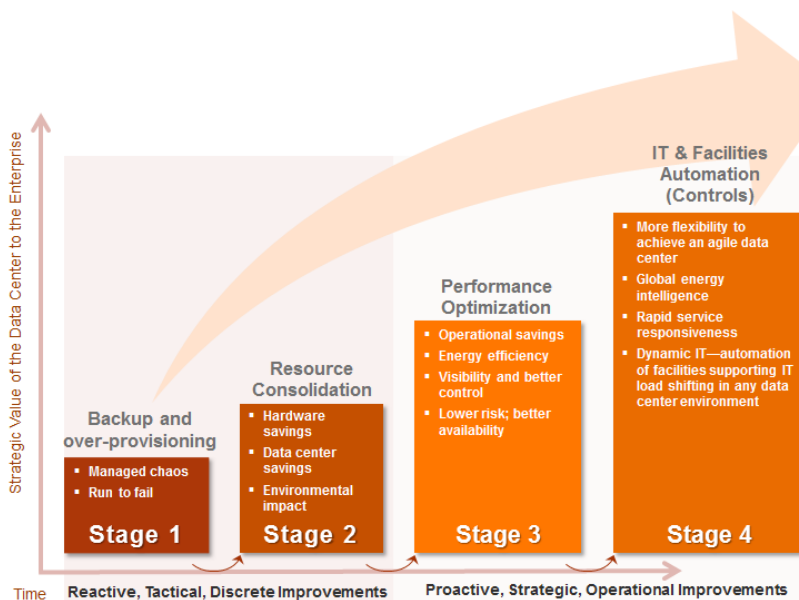
As an example, suppose you have built a 100% free-air cooled data center, and several miles away there is a forest fire. Using the flexible protocol interface defined in the previous section, you would have access to weather and wind-related data indicating that the wind, although blowing away from the data center at the moment, is forecast to shift. At this point, it is prudent to begin moving IT load to an alternate site, which you do using your DCIM automation.

A short time later, the weather data says that the wind has indeed shifted. Using DCIM automation, you take steps to shut down all non-critical load and shift the remaining load to an alternate site. Once all load has been shifted and all servers switched off, the data center vents can be closed.

Vents can be opened again once the connected air-quality analyzer indicates that the smoke and gases associated with the fire have dissipated. At that point, the DCIM system will automatically reverse the process.

### Achieving sustainable growth—reliably with DCIM

The advent of free cooling is just one example of innovative efficiency solutions that enable data center operators to meet business demands for growth while at the same time extending the productivity of existing resources and curbing costs. Sometimes with new solutions there are unintended consequences that only surface over time. A DCIM system provides the tools that enable data center professionals to mitigate risks and capitalize on new technologies and best practices.



Data Center Operational Maturity Model

A good DCIM system improves the strategic value of the data center by enabling operations to mature in stages. When employed correctly, the data center evolves from reactive, discrete improvements such as gaining visibility to data center assets and systems for managing capacity and growth, to proactive, strategic improvement such as maximum performance optimization and data center automation. Ultimately, a full-suite, open platform DCIM system enables an agile data center that readily supports business continuity.

### For more information

ABB is a global leader in power and automation technologies that enable utility, industrial and enterprise customers to improve performance while lowering environmental impact. To learn more about how ABB can help you take a strategic leap forward to the “future state” of your data center operations, contact your ABB representative, or visit <http://www.abb.com/decaathlon>.

# Contact us

## AMER

Richard Ungar

Tel: +1 905 333 7532

Email: richard.t.ungar@ca.abb.com

## EMEA

Jim Shanahan

Tel: +353 1 405 7300

Email: Jim.Shanahan@ie.abb.com

## APAC

Madhav Kalia

Tel: +65 6773 8746

Email: Madhav.Kalia@sg.abb.com

[www.abb.com/decathlon](http://www.abb.com/decathlon)

## Acknowledgements

The authors would like to thank the following for their invaluable expertise and support: Dr. Kenny Gross of Oracle, Steve Uhlir of Engineering Serendipity, David Gallaher of the National Snow and Ice Data Center, Clemens Pfeiffer of Power Assure, and from ABB: Dennis McKinley, Adrian Timbus, Chris Anthony, Henry Buijs and Juergen Kappler.

## Authors

### ABB Inc.

#### Richard Ungar

3450 Harvester Road

Burlington, Ontario L7N 3W5

Tel: +1 905 333 7532

Email: richard.t.ungar@ca.abb.com

#### Marina Thiry

730 University Drive

Menlo Park, CA 94025

Tel: +1 650 799 8299

Email: marina.thiry@us.abb.com