

VELOSTRATA

Decoupling Compute (VMs) from Storage (VMDKs)

Table of Contents

Whitepaper Snapshot	3
The Challenge with Hybrid Cloud	3
Velostrata Overview	4
Velostrata Installation	4
Intelligent Streaming	4
Bi-Directional Optimization and Multi-Tier Caching	5
Resiliency	5
Simple and Transparent Management	5
Performance Results: A Real-World Example	6
Example Use Cases	8
Summary and Conclusion	9
About Velostrata	9

Whitepaper Snapshot

This whitepaper describes Velostrata's unique technology and architecture, as well as several use cases, and a real-world example of the system performance based on a customer deployment. The paper begins by providing an introduction to the challenges typically associated with hybrid cloud deployments when production workloads are involved. This paper is intended for a more technical audience whose responsibilities include data center infrastructure and hybrid cloud initiatives. More information about Velostrata is available on our website at www.velostrata.com.

The Challenge with Hybrid Cloud

The typical data center fluctuates between two states. At times, the data center is over-provisioned and underutilized, and at other times it is under-provisioned and overcommitted. Why? Because enterprise IT organizations are faced with two difficult choices as they try to match data center supply with the growing and often variable demands of the business. IT organizations can over-provision the data center to have enough capacity for peak loads, but that inevitably results in a significantly underutilized infrastructure during non-peak times, and hence leads to wasted resources, time, and money. And the reality is that, even with over-provisioning there are still times when the data center is not be able to meet the needs of the business. Alternatively, IT organizations can provision the data center for average usage, but that impacts the performance and availability of workloads when above-average capacity is needed, leading to business disruption.

The vision for hybrid cloud is to address these data center challenges by providing a third option, namely to leverage the public cloud to augment existing data center capacity with agile and cost-effective "pay as you go" cloud computing resources. However, achieving the hybrid cloud vision, particularly for production workloads, has been elusive for a number of reasons. As is evident from the list to the right, many of the key barriers to hybrid cloud adoption for production workloads center around the myriad of issues involved in migrating storage (data, drivers, boot images, etc.) to and from the public cloud. Velostrata addresses these hybrid cloud barriers in a uniquely innovative way – by decoupling compute from storage.



Risk of storing production data in the cloud

Many enterprise IT organizations are concerned about storing production data in the public cloud due to security, compliance, and in some cases, geographic restrictions.



Cost of block storage in the cloud

Unlike compute, IT organizations must pay for block storage 24x7 in the cloud, regardless of whether that data is accessed or not. And unlike cloud storage used for archival or backup purposes, typically using lower-cost object storage with a small number of transactions, production workloads require high IOPS block storage that provides the necessary performance, and is hence much more expensive.



Time to migrate production workloads to the cloud

Migrating just 10 TBs of data over a 20 Mbps network link takes 50 days with the link 100% utilized. It is just not feasible to migrate production workloads with any significant data to the public cloud.



Migration and management complexity

In order to migrate production workloads to the public cloud, significant changes to the applications, images, storage, and drivers are typically required. In addition, managing storage in the cloud to ensure data resiliency (e.g. backup, replication, HA, DR) and performance, requires adoption of new tools and processes– a major undertaking. Finally, IT organizations need to work with both on-premises and cloud-based management consoles to manage their infrastructure and storage, dramatically increasing complexity.



Vendor Lock-In

Due to all of the afore-mentioned challenges, migrating workloads to the cloud is a one-way trip. There is no going back and there is no safety net should issues arise.

Velostrata Overview

Velostrata is the first hybrid cloud software-based system that streams production workloads to and from the cloud in minutes, keeps the storage and boot images on-premises, and optimizes performance end-to-end. Our patent-pending technology decouples compute (VMs) from storage (VMDKs) and provides intelligent streaming, optimization, multi-tier caching, and data pre-fetching capabilities to ensure optimal performance despite the WAN latency between on-premises storage and compute in the cloud. No manual changes to the applications, images, networks, or storage are required and IT organizations can leverage the same management tools and processes they use today. With Velostrata, after installation and a one-time network set up (described below), streaming production workloads is as simple as a click of a button in our vCenter plug-in. Velostrata also provides extensive monitoring capabilities as well as APIs for simple integration with 3rd-party management solutions.

Velostrata Installation

Velostrata is a software system that is prepackaged and deployed as virtual appliances and installation requires just a few easy steps. A Velostrata Data Center (DC) Edge virtual appliance is deployed in the data center and a Cloud Edge (CE) virtual appliance is deployed in a dual-node active/active configuration in the cloud for scale and high availability. The Cloud Edge is deployed within a customer-owned VPC, and connected to on-premises via a VPN. In addition, all traffic between the data center and the cloud is encrypted end-to-end, in flight and at rest. Each CE virtual appliance supports 50 concurrent VMs, and CE virtual appliances may be added as needed to scale out linearly. Velostrata includes a vCenter plug-in that, once registered, provides a number of

additional operations within the standard vCenter interface to run VMs in the cloud, run them back on-premises, as well as start and stop VMs when they are in “run-in-cloud” mode. Once the plug-in is installed and network setup (including VPN and VPC) is complete, streaming production workloads to and from the cloud is as simple as right-clicking VMs, selecting “run in cloud” or “run on-premises,” and selecting a few additional deployment options. Velostrata handles everything else, automatically and transparently. Velostrata currently supports vSphere in the data center and AWS as the cloud target, but support for additional clouds and hypervisors will be coming in the future.

Velostrata’s unique technology includes innovations in several key areas, including: Intelligent Streaming, Bi-Directional Optimization and Multi-Tier Caching, Resiliency, and Management. Each one of these areas of innovation will be covered in the section below, followed by real-world performance examples.

Intelligent Streaming

All other solutions that migrate VMs to the cloud replicate the images to the cloud, convert them to meet the format used by the cloud vendor (e.g. AMI for machine templates in AWS), and then instantiate an image and boot it locally. Velostrata differentiates from all other solutions by eliminating the need to replicate the image to the cloud and create a new AMI. Instead, Velostrata utilizes its own generic AMI and performs a native boot of an on-premises operating system over the WAN, in just a few minutes. While the image boots, it is adapted on the fly to meet the requirements of the target hypervisor and cloud environment automatically and transparently, without any user intervention.

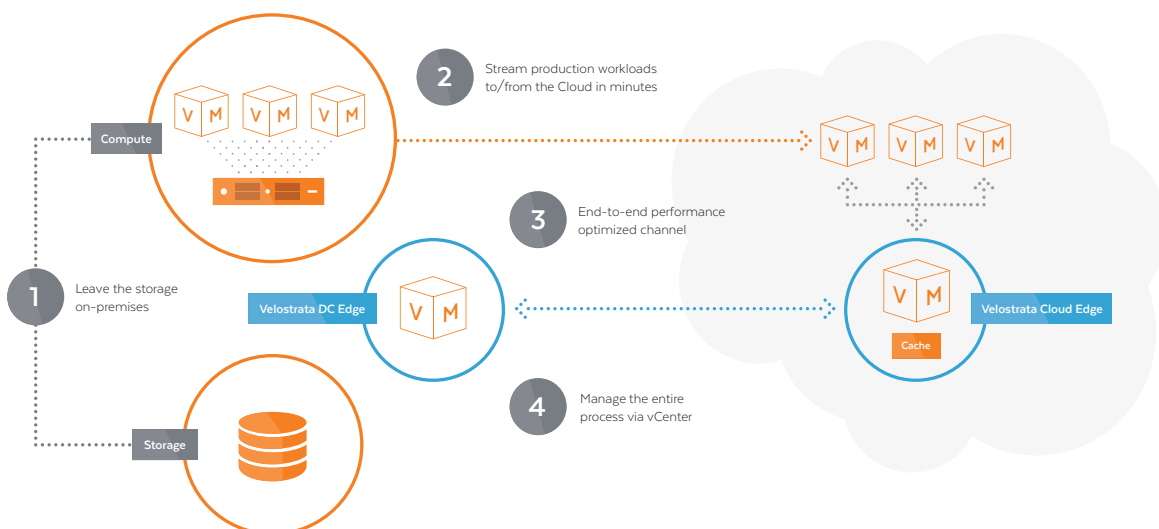


Fig. 1 » Velostrata Architecture Decouples Compute (VMs) from Storage (VMDKs)

Since the image is not replicated, Velostrata streams only that fraction of the image that is required for the workload to run and be available in the cloud. The net result of these capabilities is that production workloads are streamed back and forth between the data center and the cloud in minutes, regardless of image size. This approach is somewhat analogous to the way Netflix or other streaming video services work. Despite the length of the movie, you can start watching in minutes. Note, however, that in the case of a boot image, many applications might not be used by the server, in which case they are never streamed.

Bi-Directional Optimization and Multi-Tier Caching

Velostrata software consists of a multi-tier, read-write cache in the cloud. For read purposes, the cache includes the “working set”—the data that is frequently accessed by workloads and hence provides LAN-like performance in access to the cached data. In addition, Velostrata’s cache includes sophisticated pre-fetching algorithms that predict data most likely to be needed next, hence further improving response times. For write purposes, Velostrata acknowledges the write operations locally, hence with minimal latency overhead, and then sends the updates back to on-premises asynchronously. As a result, Velostrata must provide high resiliency (detailed in the section below) and no data-loss while the data is in transit to on-premises. Finally, Velostrata also provides strong bi-directional block level data de-duplication and compression algorithms, to minimize the amount of data that needs to traverse the WAN, and hence accelerate access to on-premises data. The result is that for most enterprise applications, performance is on par with performance in the data center, despite the fact that storage and compute are now decoupled. In fact, in some cases performance actually improves with Velostrata’s approach because larger compute instances may be instantiated in the cloud and IOPs are now a shared resource that can scale out linearly, on-demand.

Resiliency

All deployments of Velostrata include a Cloud Edge (CE) virtual appliance deployed in a dual-node, active/active configuration. Put simply, one CE instance is deployed in one AWS Availability Zone (AZ) and a second CE instance is deployed in a separate AWS AZ for redundancy and high availability. As previously mentioned, Velostrata acknowledges write operations and ensures data resiliency by performing the write operations across two availability zones. In addition, Velostrata stores the journal of write operations in an object store (S3 in AWS) to maintain a transient resilient data store while the data is written back to the VMDKs on-premises. According to Amazon, the annual uptime SLA for EC2 (dual-AZ) is 99.95%, and for S3 it is 99.99% (per year). S3

durability is 99.99999999% (11 9s). Velostrata keeps a maximum 30 seconds of write journal on the dual-AZ nodes before that data is committed to the higher resiliency S3 object store. Velostrata also includes the option to have data writes in the cloud “persist” only in the cloud. Velostrata’s architecture ensures that there is never any data loss related to a single Cloud Edge failure or data consistency issues. Further, Velostrata’s architecture ensures a 30-second RPO for sync in S3 (resilient to dual AZ crash, which is very rare) and a 30-minute RPO for sync on-premises (resilient to whole cloud outage, which is extremely rare).

Simple and Transparent Management

With Velostrata, no changes to the applications, images, networks, storage, or drivers are required and there is no need to learn new tools or processes. Storage can be managed just as it is today, regardless of whether production workloads have been streamed to the cloud or not. Velostrata also handles all image adaption from vSphere to AWS (and back) automatically and transparently. Velostrata extends the actions on an existing VM object without replication or cloning, thus providing administrators with management context, continuity, and coherency. With Velostrata, there is no change to the managed object, no ambiguity, and no loss of administrative history or operational context. Management and reporting is integrated into the vCenter console through Velostrata’s vCenter plug-in. Velostrata is also designed with a REST API for simple integration into 3rd party management tools. Streaming production workloads to and from the cloud with Velostrata involves right-clicking VMs and selecting “run in cloud” or “run on-premises.” Even if the workload is in the cloud, Velostrata’s vCenter plug-in may be used to manage that workload. A few additional deployment options consist of selecting the AWS instance type (larger instances may be selected to further improve performance), storage policy (cloud persist or write-back), security groups, networking, and execution options.

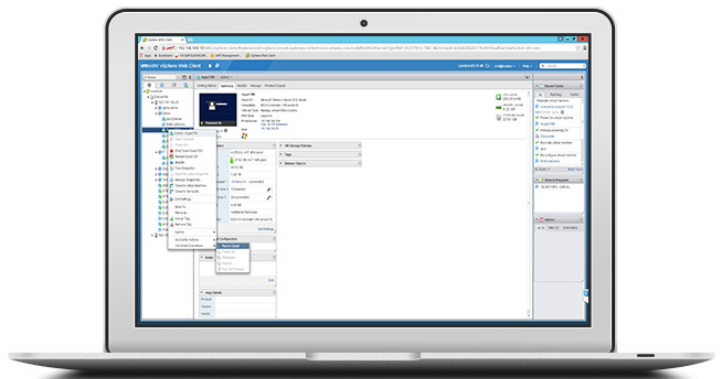


Fig. 2 » Velostrata vCenter Plug-In. Right-click VMs to “run in cloud” or “run on-premises.”

Performance Results: A Real-World Example

The logical question with Velostrata’s approach is—“How is it possible to decouple compute from storage without sacrificing performance?” Previous sections of this whitepaper cover key areas of innovation related to performance (e.g. bi-directional optimization, multi-tier caching, predictive pre-fetching, etc.). However, there is no substitute for a real-world example. To that end, the following section includes an example from customer testing in a real-world production environment. In this example, the customer tested Velostrata with two different use cases. The first use case tested was a “Hybrid Disaster Recovery (DR)” scenario and the customer wanted to ensure that Velostrata would perform as needed should a DR event occur. The second use case involved “Data Center Extension”—augmenting the customer’s currently overcommitted data center by load balancing with cloud computing resources, on-demand. Both use cases are further described later in this whitepaper.

Test Environment

Customer

Global manufacturing company with 30 plants worldwide

Multi-Tier Application

ERP system with 3 TB database server; based on Oracle WebLogic and SQL

Deployment

4 VMs (1 remote session host, 2 WebLogic application servers, 1 MS SQL server)

Network Connectivity

50 Mbps connection to the AWS cloud

Testing Locations

Data Center located in Michigan; AWS Cloud region in Northern Virginia

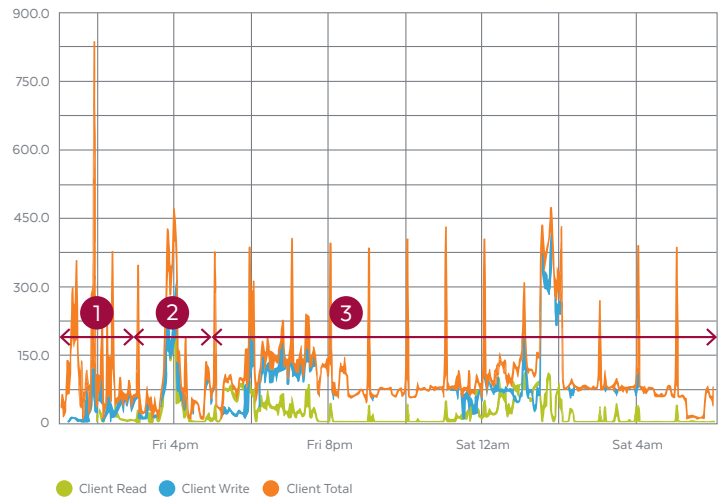
Test Scenario

Tested with real batch processes and an interactive user experience benchmark

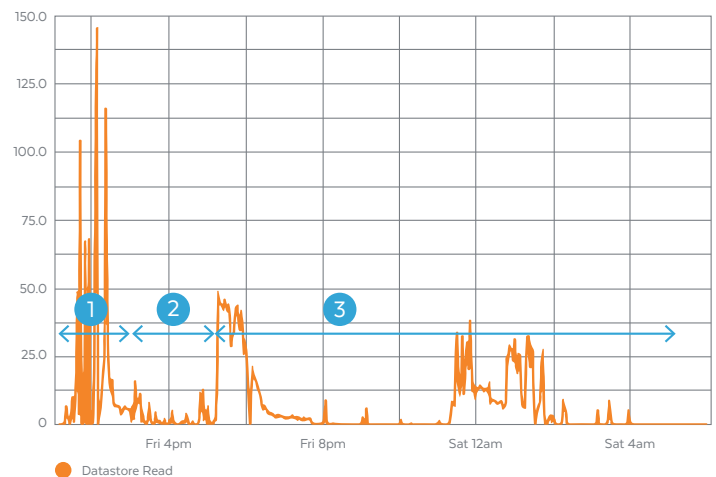
Performance Results

Graphs 1 and 2 below show three phases of the test scenario with Velostrata. The first graph illustrates the IOPs experienced by the workload (VMs) in the cloud. The second graph shows what happens simultaneously to the datastore usage on-premises. The three orange circles in each graphic highlight the following key milestones:

1. Initial boot and start of the application in the cloud
2. First batch process initiated (production line scheduling)
3. Second batch process initiated (cost modeling)



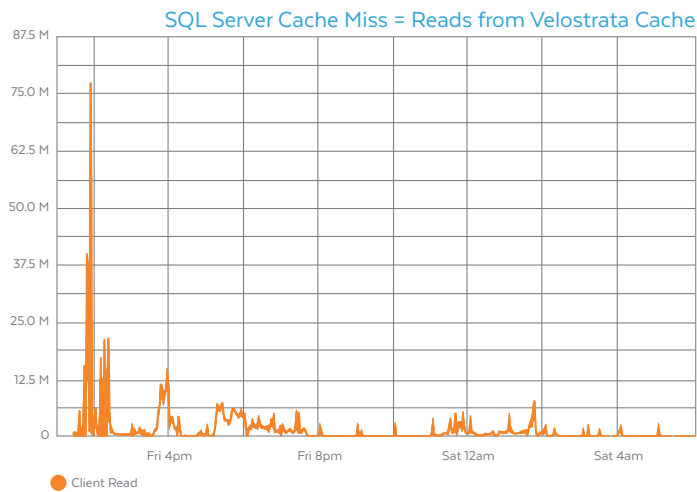
Graph 1 » Workload in Cloud (IOPs)



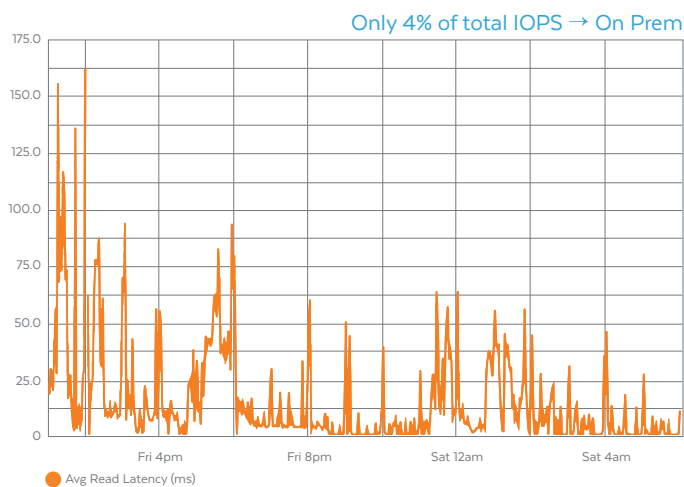
Graph 2 » Datastore usage on-premises

As shown in graphs 1 and 2 above, at each one of the phases, there is a burst of activity but then the application converges on an active working set quickly, the cache becomes “warm” quickly and data is accessed mostly locally, as evidenced by the low activity on the on-premises datastore. What occurs over a short period of time is the application effectively experiences single digit latency due to Velostrata’s caching, pre-fetching, and bi-directional optimization capabilities.

The effects of Velostrata’s read caching is illustrated in graphs 3 and 4 below. Graph 4 illustrates that, after an initial spike, the cache becomes warm and most requests are handled locally, as indicated by the lower response times. It is important to note that, in this example, SQL server had an 80% cache hit rate, and when it accessed the storage for the 20% misses, Velostrata provided an additional 80% cache hit rate. This means that from the application’s perspective, there was effectively a 96% cache hit rate. Said differently, only 4% of the data access requests were fetched from the on-premises datastore. Latency initially is a bit higher than normal, but very quickly converges to 10ms or less. Average read I/O response time for the test period was 6.6ms!

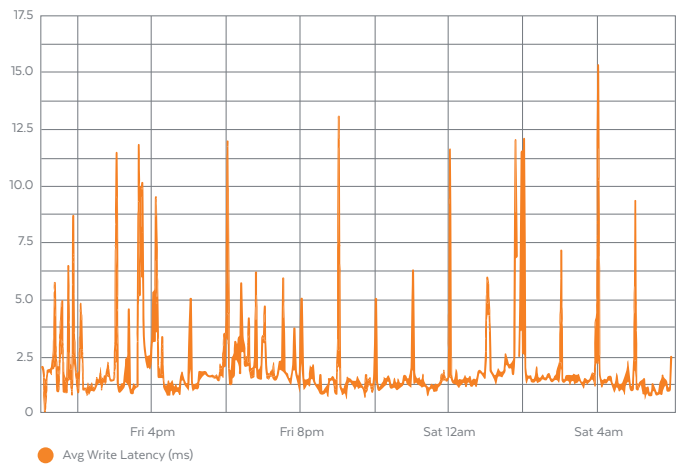


Graph 3 » Observed throughput from cache (bytes/sec).

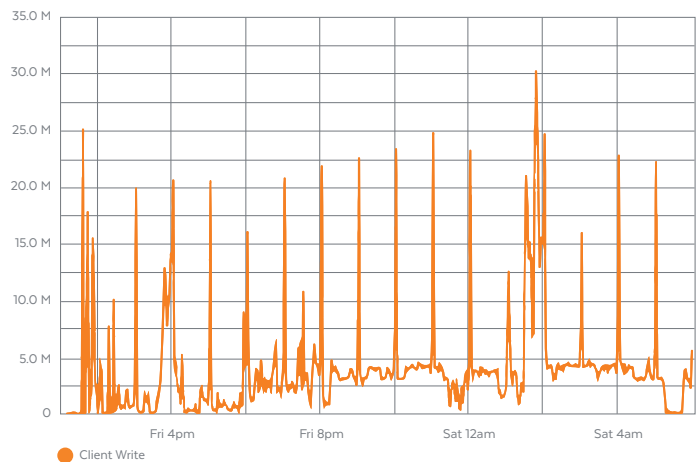


Graph 4 » Latency on SQL Server miss.

With Velostrata, all writes are acknowledged locally and write journaling is effectively mirrored across multiple Availability Zones for high availability and redundancy. As shown in the two graphs below, write latency is consistently in the single-digit (ms) range, and much of the time in the sub 2 ms range (including the time for cross AZ mirroring). Average write latency for the test period is less than 5 ms!



Graph 5 » Data writes average less than 5ms.



Graph 6 » Observed throughput comparable to on-premises (bytes/sec).

In the example above, performance with Velostrata was very quickly comparable to performance prior to deploying the solution, despite the fact that Velostrata decouples compute from storage. In this example, performance with Velostrata would have exceeded prior on-premises performance had larger AWS compute instances been selected. In fact, in several test cases, customers experienced better performance at steady-state with Velostrata than they previously experienced with compute and storage co-located.

The next logical question is—“So, what types of workloads are not a good fit with Velostrata?” It’s worth mentioning that the vast majority of workloads perform extremely well thanks to Velostrata’s extensive IP in this area. Production workloads such as SAP, Oracle Weblogic, ERP, WebSphere, MySQL, SQL server, CRM, and many others have all been tested extensively and performance is on par to performance without Velostrata. As is evident from the list above, even relational databases with a mix of random read/write IO and high performance requirements perform well with Velostrata. Examples of workloads that are not a good fit with Velostrata would be some data warehouse applications that scan the entire record set before reading a record (as an example). In that scenario, cache hits will not be sufficient to provide consistently good performance. Another obvious example includes SaaS applications that are already hosted in the public cloud. Based on real-world testing with many customers in production environments, the vast majority of workloads perform extremely well with Velostrata.

Example Use Cases

In addition to the test example included in the section above, customers deploy Velostrata for a number of different use cases. Example use cases include:

Data Center Extension

Data Center Extension refers to the concept of extending existing data center capacity with cloud computing, typically for planned events. Examples include moving VMs from an overloaded host to the cloud as a load balancing technique, in order to restore performance and stability. The flip side to this example is also appropriate for Data Center Extension, namely moving low-utilization VMs to the cloud in order to improve data center costs and efficiency.

Cloud Bursting

Cloud Bursting refers to the concept of provisioning the data center for average workloads and bursting to the cloud to accommodate peak load times. When customers talk about Cloud Bursting, they typically refer to two types of events—one planned and one unplanned. An example of a planned event might be seasonality that is typical in industries such as retail. In that example, it makes little financial sense to provision the data center for peak load when most of the year experiences much lower volume. An example of an unplanned event might be a boot storm in the middle of the day where additional capacity is needed immediately to accommodate a surprise burst in demand.

Hybrid DR

Hybrid DR refers to the concept of providing site-to-site DR without the cost and complexity of duplicating compute resources at the secondary site. Storage must still be replicated to the secondary site as in conventional site-to-site DR. However, in a conventional DR solution the compute infrastructure needs to be

fully provisioned at the secondary site too, and is idle 99% of the time waiting for a DR event to occur. With Velostrata, since workloads are streamed to and from the cloud in minutes, there is no need to keep any compute (other than a very small footprint needed to run Velostrata) at the secondary site. Velostrata is the only solution that provides this unique approach to disaster recovery.

Data Center Consolidation

Data Center Consolidation refers to the concept of consolidating storage to a centralized data center while running VMs in cloud regions geographically close to regional branch offices. Velostrata is the only solution that enables this use case via our unique streaming, caching, and optimization technologies. Despite the WAN latency involved between the now distant compute and storage, performance remains optimal with Velostrata.

Dev/Test

Dev/Test is a well-known use case, but Velostrata’s approach is entirely unique. With Velostrata, there is no need to migrate storage or convert images. Velostrata is the only solution that leaves storage on-premises while streaming compute workloads to and from the cloud in minutes. Data writes in the cloud may “cloud persist” or optionally be configured for write-back to the data center.

In addition to the examples listed above, customers continue to find additional use cases that are valuable to their organizations. More details about each use case may be found on the Solutions page of the Velostrata website, <http://velostrata.com/solutions>.

Summary and Conclusion

Velostrata's unique technology decouples compute from storage and resolves all previous barriers to hybrid cloud adoption. Velostrata mitigates security and compliance risks and reduces costs by 50% or more since authoritative data remains on-premises. There is no longer a need to over-provision the data center or pay for expensive block storage 24x7 in the cloud. Business agility improves and vendor lock-in is prevented since production workloads are streamed to and from the cloud in minutes while performance is optimized end-to-end, despite the WAN latency between compute and storage. Velostrata's integration with vCenter, simple and transparent approach to streaming workloads, and REST API eliminate the complexity typically involved with hybrid cloud deployments. A summary of Velostrata's benefits is included on the right.

About Velostrata

Velostrata's mission is to enable frictionless, no-compromise hybrid clouds. "Frictionless" means enterprises have the speed and flexibility to stream production workloads to and from the cloud in minutes with the click of a button. "No-compromise" means enterprises leave the storage on-premises, and yet performance is optimized end-to-end. With Velostrata, hybrid cloud is finally low risk, cost-effective, fast and simple. Velostrata is backed by Norwest Venture Partners and Greylock Partners IL (83North) and is headquartered in San Jose, California with R&D in Israel. For more information, visit: <http://www.velostrata.com>.



Mitigate security and compliance risks

With Velostrata, authoritative data is left on-premises. The fractional data that is included in the cloud cache is essentially transitory (unless the "cloud persist" option is selected) and all data is encrypted at rest and in motion.



Reduce costs

Reduce costs by 50% or more with Velostrata. Leverage the cloud for compute, on-demand, instead of over-provisioning the data center and paying for server, storage, networking, licensing, and power & cooling costs. With Velostrata there is also no need to pay for expensive cloud storage 24x7 in the cloud since only a fraction of the data ever resides in the cloud.



Improve business agility

With Velostrata, there is no need to migrate storage or boot images to the cloud. Production workloads are streamed to and from the cloud in minutes, and Velostrata handles all image adaptation automatically and transparently. New projects can be brought online in minutes and resources allocated dynamically, on-demand.



Improve operational efficiency

Streaming workloads to and from the cloud requires just a click of a button, and there is no need to manage storage in the cloud. No changes to the applications, images, storage, or drivers are required, and no new tools or processes need to be learned. Velostrata provides extensive monitoring capabilities as well as a REST API for simple 3rd-party integration.



Eliminate vendor lock-in

With Velostrata, it's no longer a one-way, time-consuming migration to the cloud. Workloads may be streamed to and from the cloud on-demand. There is now a safety net if something goes wrong.