

Using In-Memory Data Grids for Global Data Integration

by Dr. William Bain, ScaleOut Software, Inc.

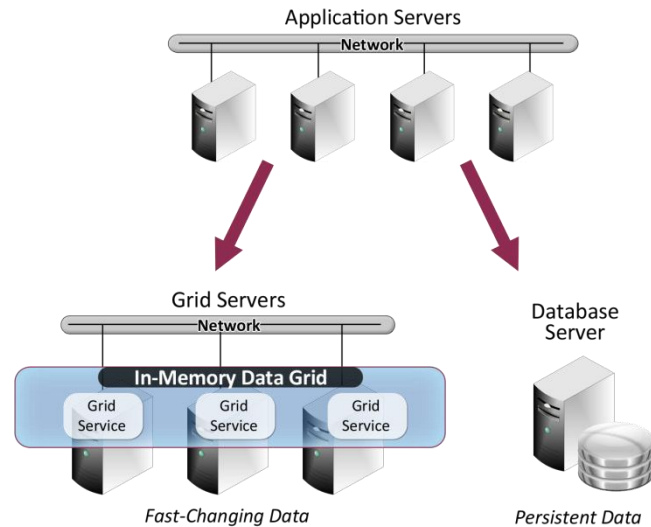


By enabling extremely fast and scalable data access even under large and growing workloads, in-memory data grids (IMDGs) have proven their value in storing fast-changing application data. For example, Web server farms use IMDGs to hold and share large volumes of shopping carts under heavy Web loads. Applications in financial services use IMDGs to hold fast-changing stock trading data for processing orders or for quickly analyzing and responding to emerging market trends.

An increasing number of companies employ multiple data centers to distribute their workloads and mitigate the impact of catastrophic events such as earthquakes and floods. IMDGs can be used to complement disaster recovery strategies by continuously replicating changes to fast-changing grid-based data to remote sites. This enables fast recovery and resumption of processing without data loss after a disaster strikes.

The use of in-memory data grids has also created the opportunity for organizations to employ even more powerful global strategies for data sharing. As organizations work to efficiently access fast-changing data across multiple sites or scale their processing into the cloud, the need to quickly and seamlessly migrate data on demand has grown rapidly. For example, organizations that produce and store fast-changing data in multiple data centers need to be able to access and analyze data without regard to where it originates. Likewise, organizations that access the highly elastic resources of public clouds need an efficient way to restage data in the cloud for processing.

Because IMDGs are specifically designed to store fast-changing data, federating IMDGs across multiple sites and enabling seamless access to data among all federated sites provide an ideal solution to the challenge of global data access. The benefits are twofold. First, applications can efficiently access and update data simply by using the IMDG's data access mechanisms without modification; the federated IMDGs handle all of the details of remote data access and coherent updating. Second, IMDGs provide the scalability and low latency required to enable applications to handle large workloads with fast responsiveness.



Using an IMDG to Store Fast-Changing Data

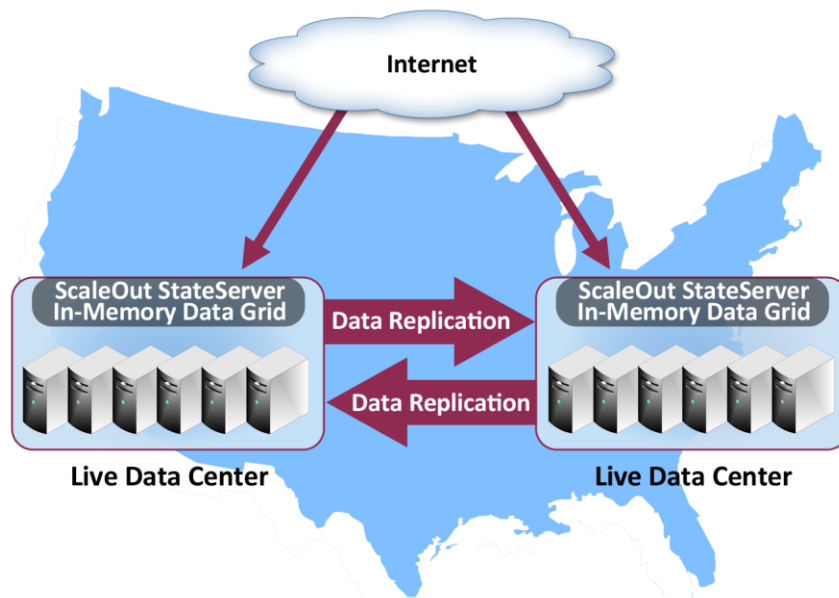
We describe the combined scenarios for data replication and sharing as *global data integration*. This article outlines how in-memory data grids easily can be deployed to implement key strategies for global data integration, and it describes the important benefits this technology brings to organizations with global reach.

Disaster Recovery

A solid disaster recovery strategy requires that if one data center goes offline, its workload can be handled by another healthy data center to avoid service interruptions. For this recovery strategy to be effective, changes to fast-changing application data must be continuously replicated to a remote site so that the site is immediately ready to handle the workload. An IMDG that includes site-to-site data replication to one or more IMDGs at remote sites can provide this important capability and thereby complement the data center's other replication and recovery strategies. In addition, all data centers can be operated in a "live-live" configuration under normal operating conditions to make full use of all computing resources and avoid the need for an idle "stand-by" data center.

Carefully integrating data replication technology into an IMDG's software architecture enables it to deliver the performance and reliability needed to handle large, fast-changing workloads. It also enables this capability to be easily deployed and managed by IT administrators. ScaleOut GeoServer® DR from ScaleOut Software is an example of a

technology that provides these capabilities. Because it is designed to extend the scalable, highly available architecture of its underlying IMDG, ScaleOut StateServer® (SOSS), it automatically scales replication bandwidth as grid servers are added to handle growing workloads, and it automatically tolerates server failures without interrupting operations. Additionally, it provides management tools that allow IT staff to easily establish and monitor connections to remote sites.



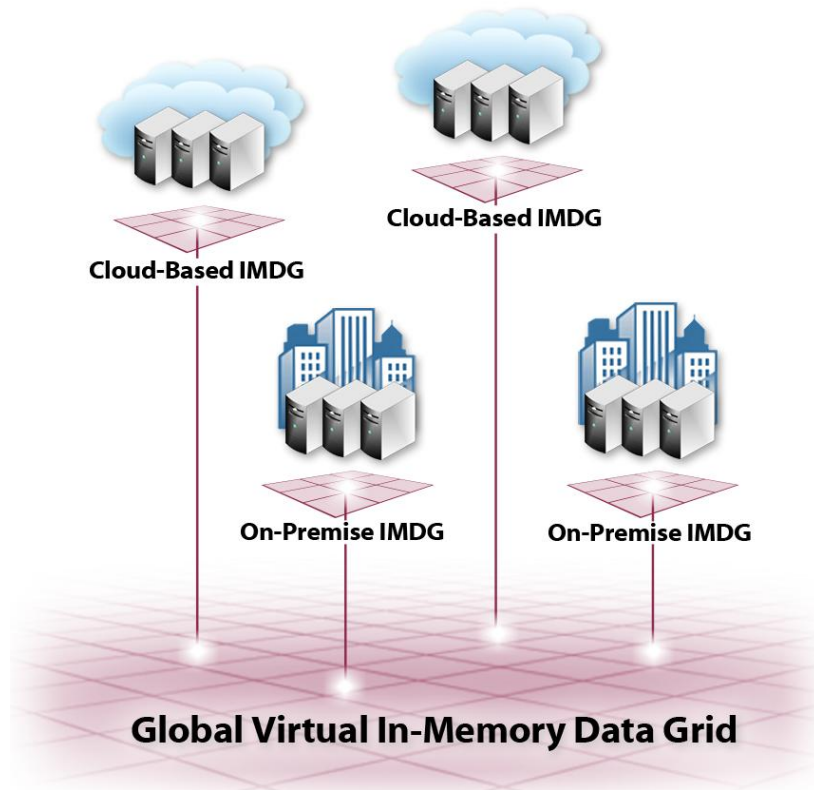
Example of Data Replication Between Two IMDGs for Disaster Recovery

Global Data Access

Beyond data replication for disaster recovery, global data integration provides a range of choices for federating data stored in IMDGs at multiple data centers and cloud sites. For example, multiple data centers can be integrated into a single virtual data grid to provide seamless access to data, regardless of where it is stored and where the access request originates. Also, multiple grids can be interconnected to provide automatic data migration and elastic scaling when needed.

To ensure that global data access can easily be integrated into applications, IMDGs can seamlessly incorporate global access into their data access mechanisms. This simplifies application design by making remote data access transparent and automatic. It also eliminates the need for applications to track where data is located and manually restage it for local access. As an example, ScaleOut GeoServer follows this approach by extending the

APIs provided by ScaleOut StateServer to transparently access data on demand at a configured set of remote sites; all grid accesses proceed as if data were located in the application's local IMDG. ScaleOut GeoServer automatically searches remote IMDGs for missing data and copies it into the local IMDG as needed.



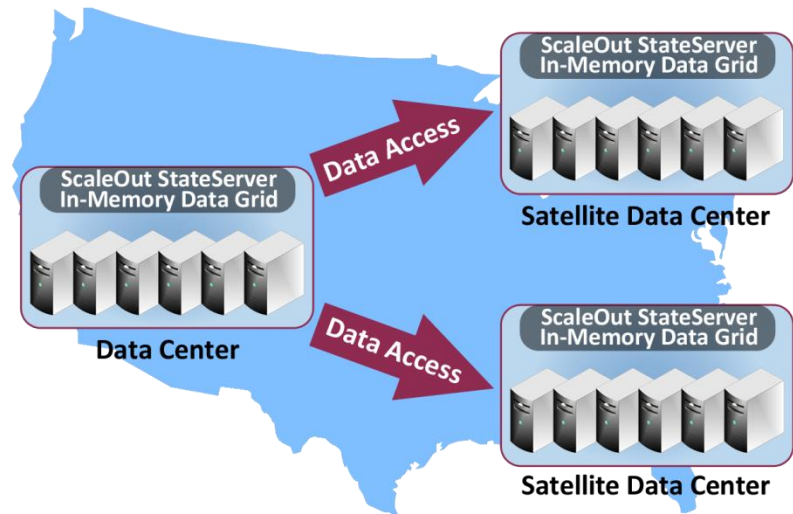
Remote Data Access Among On-Premise and Cloud-Hosted IMDGs

“Mostly Read” Access

ScaleOut GeoServer gives applications fine-grained control over data sharing to ensure efficient use of wide area networks (WANs) and to support various usage models. In one important use case, described as "mostly read" access, applications primarily need to access certain remote data but not perform updates on that data. This type of remotely accessed data is typically static or slowly changing so that local copies only need infrequent refresh over the WAN. Examples could include product pricing information for Web sites or portfolio holdings in financial services.

ScaleOut GeoServer implements mostly-read access by creating a local copy of remotely accessed data and allowing the application to specify a policy for refreshing it. The use of a local copy keeps local reads fast and minimizes WAN usage. Individual data objects can be marked by the application either to be updated periodically or to be updated when a change occurs at the remote site. Called *coherency* policies, these rules allow applications to tailor WAN usage to the characteristics of the data being remotely accessed.

An example of mostly read access, consider a wealth management application that needs to update its portfolios with periodic price changes; prices for different investments are held in multiple data grids around the world. The application can use global data access to obtain and efficiently track prices, with updates flowing into its local IMDG at the frequency required by the application. Also, to minimize WAN usage, only the prices of investments specifically needed by the application are retrieved over the WAN.



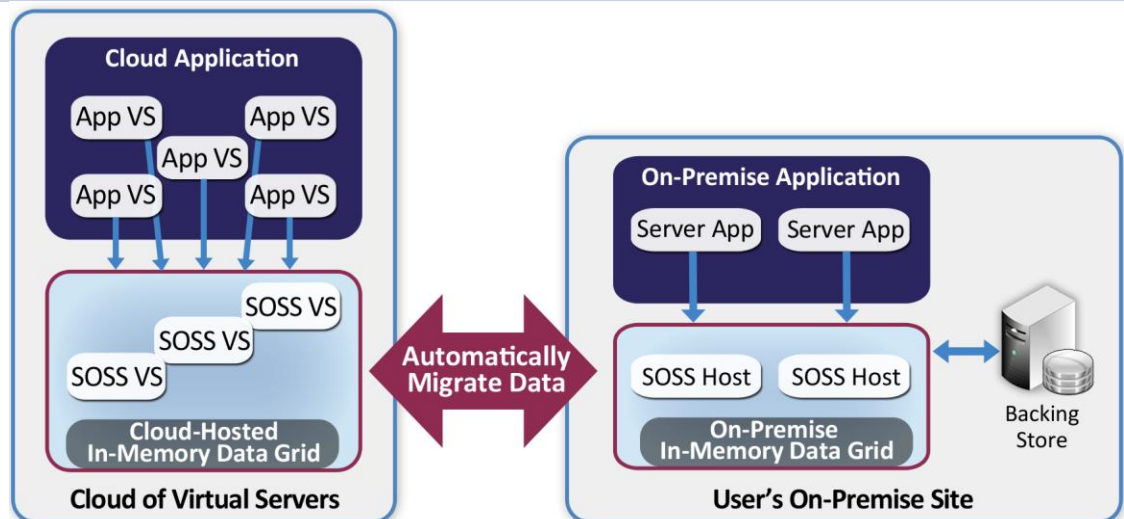
Remote Data Access From Two Remote IMDGs

“Read/Write” Access

In a second important use case called "read/write" access, remotely accessed data needs to be accessed and then updated, and updates by different sites need to be carefully synchronized. Examples include shopping carts in a Web site or financial portfolios being managed (not just examined) at remote sites. These data types can be fast-changing, and it is imperative to synchronize updates to avoid corrupting vital application data.

To synchronize updates, data must migrate from site to site on demand and avoid the use of local copies which could become out of date. ScaleOut GeoServer implements data migration and read/write access by transparently incorporating it into the IMDG's existing distributed locking mechanism, which has been extended to span multiple sites. The IMDG automatically migrates ownership of data from a remote site when it is locked for reading by the application. This ensures that updates are always performed locally and at exactly one site at a time. The application does not have to manually restage data across sites nor provide its own mechanism for global data synchronization.

As an example, consider a premise-hosted ecommerce Web farm that needs to scale into the cloud to handle high seasonal demand. To accomplish this, the Web site's administrator reconfigures the IP load-balancer to distribute Web requests across both on-premise and cloud-based Web servers; this procedure is sometimes called "cloud bursting." By using an IMDG capable of global data integration, all Web servers transparently and coherently retrieve and update shopping carts within a single, virtualized IMDG spanning both sites. The following diagram illustrates this scenario using ScaleOut StateServer ("SOSS") IMDGs at both sites and ScaleOut GeoServer to provide automatic data migration.

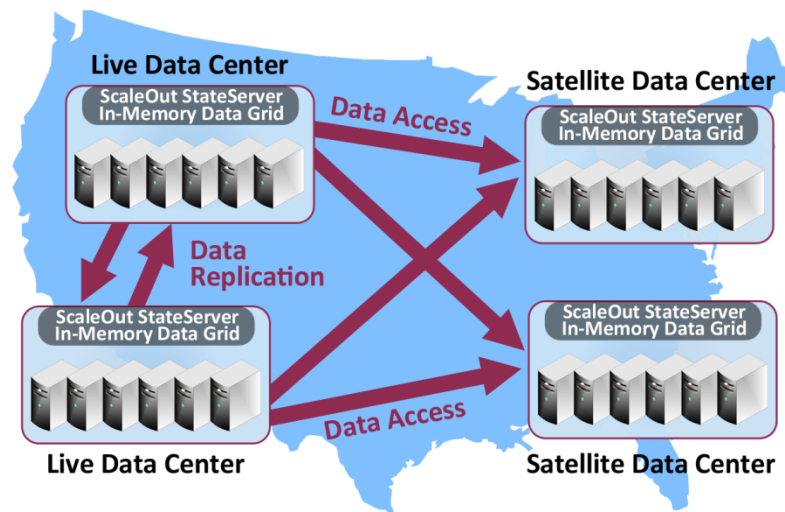


Automatic Data Migration Between On-Premise and Cloud-Hosted IMDGs

Combining Data Replication and Global Data Access

It is often useful to combine the capabilities described above for global data integration to simultaneously address multiple requirements. For example, two central data centers

which hold data accessed by satellite data centers can use data replication for disaster recovery purposes. Both could handle live traffic as described above, but in the case of a data center failure all traffic is routed to the healthy data center. Applications running in satellite data centers can use global data access to retrieve and/or update data held in the two central data centers. These applications can access data from either data center and transparently receive it even if one of the central data centers goes down. As illustrated in the following diagram, this configuration demonstrates the power and flexibility of global data integration.



Combining Data Replication and Remote Data Access

Benefits of a Virtual Data Grid

As we have seen, the goals of global data integration are to replicate data for disaster recovery and to enable applications to transparently access data across multiple sites as needed. ScaleOut GeoServer's implementation of global data integration accomplishes these goals by creating a virtual data grid that seamlessly federates in-memory data grids across multiple sites. This enables application developers to write programs which access all shared data from a single (local) IMDG, leaving the IMDG to implement the details of remote access and synchronization. After a minimal amount of configuration to connect to remote sites, changes to add or remove grid servers in any data center do not affect configuration of the virtual data grid. The virtual data grid is able to withstand and recover from WAN interruptions and other failure conditions without affecting applications.

This article has illustrated the power of global data integration to extend the reach of applications that manage data spanning multiple data centers. As we have seen, in-memory data grids (IMDGs) provide a fast, scalable storage repository for application data. Their mechanisms can be transparently extended to enable data replication for disaster recovery and global access to data held at remote sites. These capabilities open up important new scenarios for globally distributed applications and simplify their implementation. Now applications can seamlessly access data worldwide and extend their processing into the cloud to handle peak workloads. Managing geographically distributed data has never been easier.

Dr. William L. Bain is Founder and CEO of ScaleOut Software, Inc. Bill has a Ph.D. in electrical engineering/parallel computing from Rice University, and he has worked at Bell Labs research, Intel, and Microsoft. Bill founded and ran three start-up companies prior to joining Microsoft. In the most recent company (Valence Research), he developed a distributed Web load-balancing software solution that was acquired by Microsoft and is now called Network Load Balancing within the Windows Server operating system. Dr. Bain holds several patents in computer architecture and distributed computing. As a member of the Seattle-based Alliance of Angels, Dr. Bain is actively involved in entrepreneurship and the angel community.

■ ■ ■

www.scaleoutsoftware.com