# sVC32 Collaboration and Virtualization using Active Cloud Engine

## Stephen Edel

**IBM SONAS Product Manager**

IBM Storage

# Edge2012

Expect More from Your Storage

**June 4 – 8**  Orlando, Florida

This presentation explores how IBM Active Cloud Engine (ACE) in IBM Scale Out Network Attached Storage (SONAS) enables customers to collaborate and share information both within and across globally distributed locations. Discover how the future virtualization capabilities of ACE will enable migration and sharing of data between systems, to open up new opportunities to help customers manage their data. Real-world customer use cases will be discussed
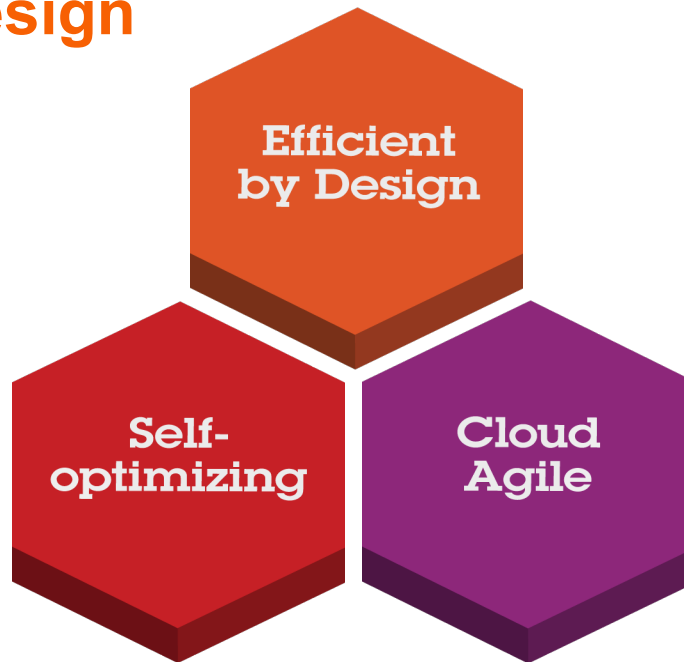
With many thanks to John Sing for Contributions to this presentation

# Introducing IBM Smarter Storage for Smarter Computing

**Efficient by Design**

Efficient by Design

Self-optimizing

Cloud Agile

**Self-optimizing**

**Cloud Agile**
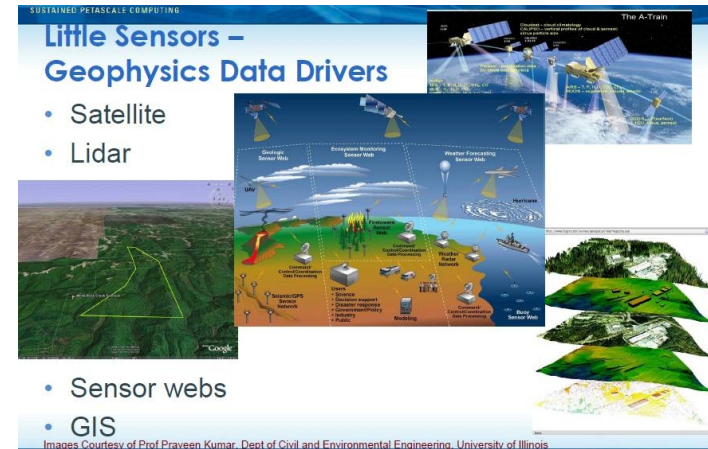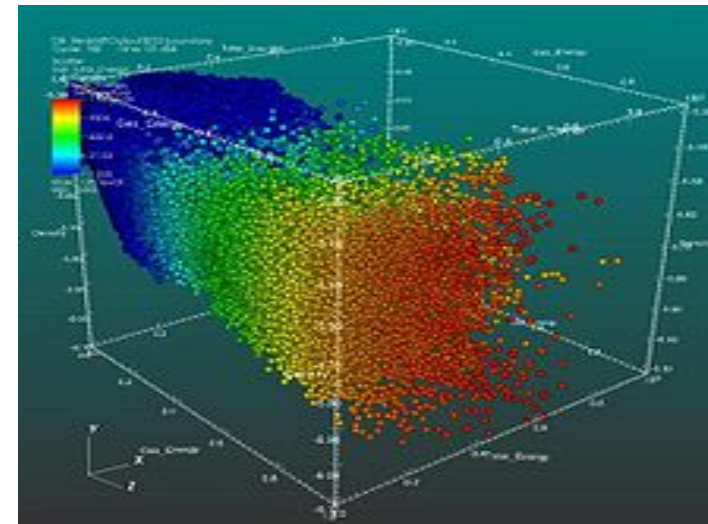
Edge2012

# Business Concerns of Unprecedented File System Data Growth

- **New and larger digital content types continuously evolving**
  - File sizes are growing to TBs
  - Moving from still images to HD video
  - Millions of files to Gazillions of objects
  - More kinds of files and objects
  - Volumes are reaching PBs per day

- **Greater collaboration requirements**
  - Creations points are many and everywhere
  - Collection points used to be few but now many, spread globally

- **New requirements to meet mission deployment / time to delivery**
  - Analysis needs information instantly
  - Access expected anywhere, anytime
  - Analysts no longer only in Washington, are deployed everywhere in the world

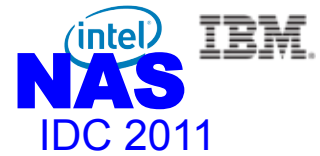- **Mandates for longer retention periods and requirements for Archive & Compliance**



IEEE Massive File Storage presentation, author: Bill Kramer, NCSA: http://storageconference.org/2010/Presentations/MSST/1.Kramer.pdf:



Source: Wikibon March 2011

Edge2012

# Fundamental, common elements to consider when evaluating a scale-out NAS solution:

IDC 2011

- Multidimensional scalability**:** Scalability isn't just about expanding hardware capacity …The best solutions will allow you to **scale each dimension** independently….

- Storage efficiency. Solutions that enable easy adoption of a continually expanding range of storage asset optimization technologies … The best solutions will allow you to introduce or recalibrate these capabilities **without major disruptions or system migrations.**

- Intelligent Information Management. …The best solutions will provide tight integration with **advanced metadata, index, and analytics solutions**… (that enable more intelligent information management as volumes of data sets grow and new use cases proliferate.

Edge2012

# Top requirements for Scale-Out NAS solution

ESG analyst: Terri McClure
December 2011

1. **Global namespace**
2. **Clustered, managed as a single entity**
3. **Ability to scale bandwidth, processors, and storage independently**
4. **Power-efficient**
5. **Self-managing**
6. **Self-healing**
7. **Transparent data mobility**
8. **Tiered storage support**
9. **Feature/functions :**
   **If going after enterprise IT then snapshots/clones/remote mirrors**
   **If going after archive then deduplication/compression**

Edge2012

# Gartner's recommendations to Scale Out NAS vendors

**Companies selling high-end scale-out NAS should:**

- **Continue investing in building the ecosystem and complete solutions for managing large amounts of** unstructured data

- **Invest in R&D to develop practical and cost-effective** backup **and** recovery **solutions**

- **Enhance** global access **capabilities to address the increasing need for** global sharing and collaboration**, as well as** cloud **initiatives**

- **Publish SPECsfs Network File System (NFS) and Common Internet File System (CIFS) benchmark results based on real-world configurations and workload types to make the benchmark more meaningful**

- **Develop** more-affordable solutions **for long-term** retention **of unstructured data**

Gartner analyst: Pushan Rinnen
November 2011

Edge2012

# Why IBM ACE is important to your clients – and You !

*Let's Take a Look at the Recent Evolution of Storage Marketplace*

Mid-1990's **– Data on IP Network =** NAS (Network Attached Storage)

Early 2000's **– Block Level Virtualization =** Storage Virtualization

Mid/Late- 2000's – **Global Namespace = SONAS (Scale Out NAS)**

2011+ – **File Virtualization / Cloud Storage** = **Active Cloud Engine**

Edge2012

# IBM Active Cloud Engine™ (ACE)

- **What is IBM Active Cloud Engine?**
  - **Policy-driven engine** that helps improve storage efficiency by automatically
    - Distributing files, images, and application updates to multiple locations  *
    - Identifying files for backup or replication to a DR location
    - Moving desired files to the right tier of storage including tape in a TSM hierarchy
    - Deleting expired or unwanted files
  - High-performance: can scan billions of files in minutes with scale out nodes

- **What client value does Active Cloud Engine deliver?**
  - Enables ubiquitous **access to files** from **across the globe**  *
  - Reduces networks costs and helps improve application performance by **distributing files closer to users**  *
  - Improves data protection by identifying candidates for backup or DR
  - Lowers storage cost by **moving files transparently** to the most appropriate tier of storage
  - Controls storage growth by transparently **moving older files to tape** and deleting unwanted or expired files
  - Enhances administrator productivity by **automating file management**

Edge2012

* Available only with Active Cloud Engine Global Management

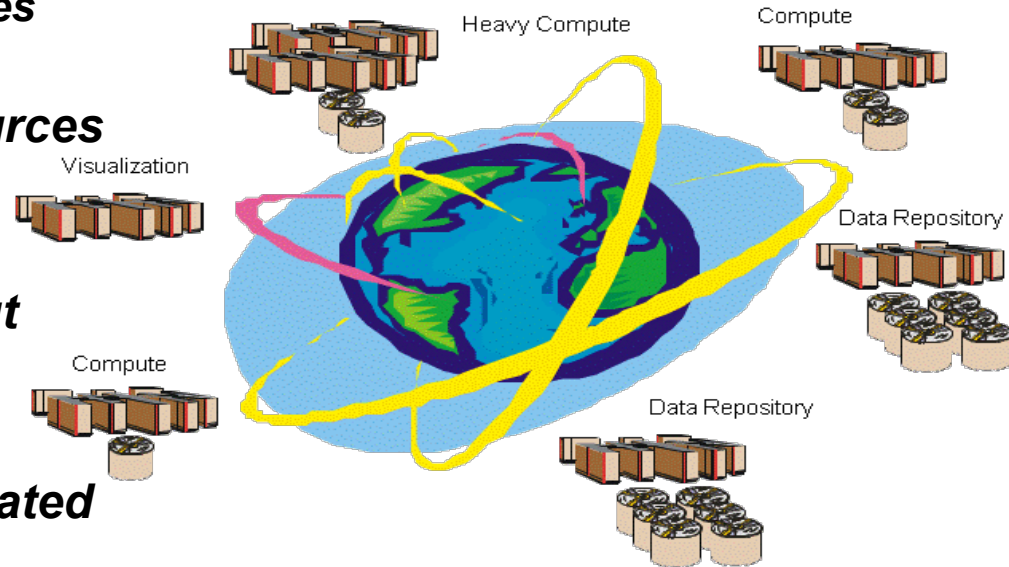# IBM Active Cloud Engine Local Management Capabilities

- **Active Cloud Engine enables two types of policies**

  - *File placement policies*
    - Used to automatically place newly created files in a specific storage pool
    - Files in the same fileset may reside on different storage pools

  - *File management policies*
    - Transparently move files to another internal or external storage pool
    - Delete, backup, incremental backup, restore, replicate
    - Remote caching, intelligent prefetching of subsets of files
    - Virtualize to single global namespace.

Edge2012

# IBM Active Cloud Engine Global Capabilities

Manage file system data with one unified management system for global file access through policies and automation while managing its lifecycle with reduced risk, complexities and cost

- *Break the dependence of an object from its location*
  - *Files (or objects) move around the sites*

- *Allows sites to contribute their resources to the global file system*

- *Allows a single gold copy of data, but copies also at the right place*

- *Allows data to be permanently replicated for fault-tolerance and performance*

- *Supports disconnected operations for network deprived locales*

Heavy Compute

Compute

Visualization

Data Repository

Compute

Data Repository

Edge2012

# IBM Active Cloud Engine - Cache Cluster

## Data Caching

- Data requests are handled locally if the requested data is cached

- Requests for data not in cache are forwarded to the gateway nodes to retrieve the data from the home cluster.
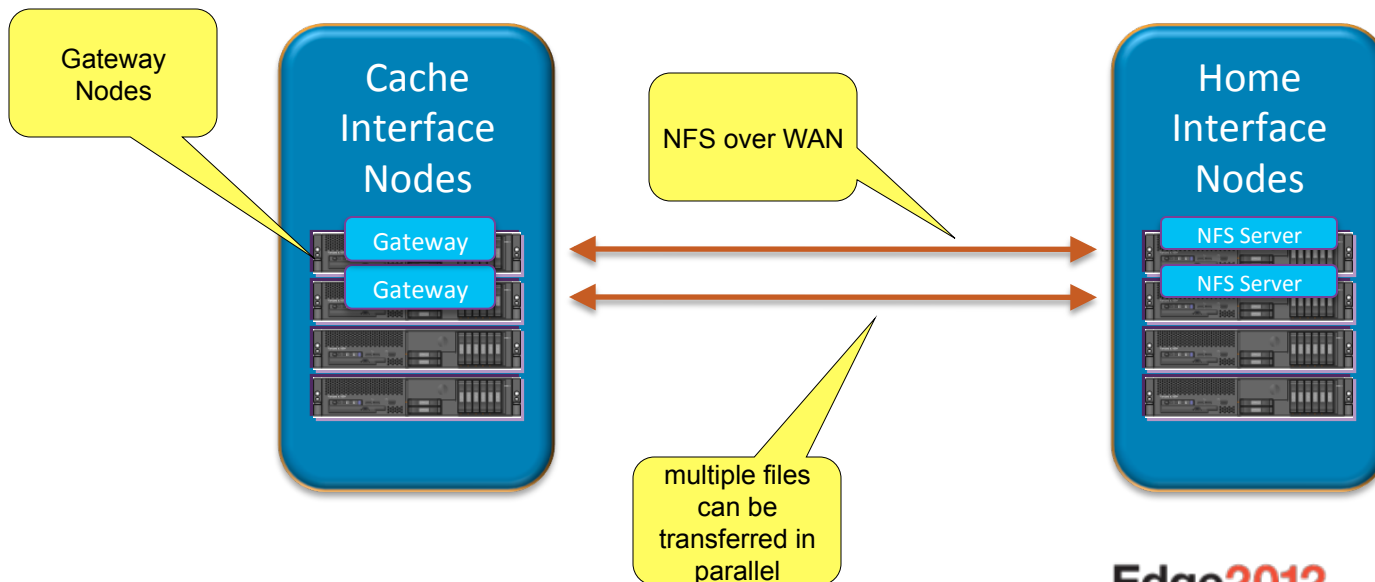
## Gateway Nodes

- IBM SONAS interface nodes take on a personality known as a "**Gateway**" node

- Gateway nodes communicate with the home cluster

- Any interface node can function as a gateway node

  - And still function as an interface node

Additional gateway node personality

Interface Nodes

Gateway          Gateway

InfiniBand Fabric

Storage Pods

Edge2012

# IBM Active Cloud Engine – Home Cluster

- When data is not in the cache, the Gateway nodes request data from the home cluster

- Permissions are checked before the file is copied from the home cluster

- The Home Cluster exports data using the NFS protocol using one or more interface nodes

Gateway Nodes

Cache Interface Nodes

Gateway
Gateway

NFS over WAN

multiple files can be transferred in parallel

Home Interface Nodes

NFS Server
NFS Server

Edge2012

# Active Cloud Engine Features and Advantages

*Extends "common" namespace & ILM/HSM Policy engine to global geo-dispersed file storage*

"Global Namespace, not just "Common" Namespace

Ownership and Relationships all done on a fileset boundary

SONAS/ACE

Protocols
CIFS
NFS
HTTP
FTP
SCP

Management
Central
Administration
Monitoring
File Mgmt

Availability
Data Migration
Replication
Backup

Data Center

Global WAN Caching removes latency effects

**Network**

SONAS/ACE

SONAS/ACE

Geo-dispersed Replicas

SONAS/ACE

SONAS/ACE

Edge2012

# Active Cloud Engine Modes (Fileset level)

- **Read Only Cache**
  - Cache can only read data, no data changes/updates allowed for that fileset at that cache location

- **Exclusive / Single Writer Cache**
  - Cache can write data to each fileset it owns.  Home cluster can't change as it does not own the fileset.  Other peer caches have to be setup as read only caches.

- **Local Update**
  - Data is cached from home and changes are allowed like Single Writer mode, but changes are not pushed to home.
  - Once data is changed the relationship is broken i.e cache and home are no longer in sync for that file

- **Change of Modes**
  - SW & RO mode caches can be changed to any other mode
    - Requires cli to be run on cache & home
  - LU cache can't be changed

Edge2012

# Use Case: Remote Cache Sites operating in Read Only Mode

Periodic Prefetch

On Demand Pull

Home cluster (Writer)

Remote Cache clusters (Reader)

- Central Data Center is where data is created, maintained, updated / changed.

- Branch locations can periodically prefetch (via policy) or pull on demand files as requested

- Data is revalidated back to Home cluster when accessed by a user at the cache site

- Customers Use Cases:
  - Financial branch offices
  - Dispersed Government locations
  - Media Distribution Centers
  - WW Research Labs
  - Regional Hospitals

Edge2012

# Read Only Scenarios (cache validate/miss)

- **If data is not in cache, pull file attributes and create on demand (lookup, open, …)**

- **Directory Traversal**
  - On first access, entire directory is read remotely (READDIR) and directory entries stored locally

- **On a later data read**
  - Whole file is fetched over NFS and written locally
  - Multiple files can be read in parallel across multiple nodes
  - Application can continue after required data is in cache while the remaining file is being fetched
  - In addition, a cache-site can be pre-populated before any access

- **On a cache hit**
  - Attributes are *revalidated* depend on revalidation delay
    - (done periodically on lookup and open)
  - If data hasn't changed it is read locally

- **On a disconnected mode access**
  - Consider cache valid and allow open without revalidation
  - Files not cached are returned as not existing

Edge2012

# Use Case: Remote Cache Sites having Single Writer Ownership

User A's exclusive fileset (writer)

User B's exclusive fileset (writer)

Home cluster (reader)

Backup site

Remote Cache clusters

- Each site has defined filesets which own the data and serve as their local cluster.

- Remote sites push data back to central Data Center for additional DR processes and for access by others

**Use Cases**

- Federated Institutes of Research with a common repository for access

- Data Repositories or Archives

- Oil/Gas exploration delivering data to a Central Site

- Universities sharing a dedicated home file system for students  but also linking other campuses

- Telco headquarters tracking logs/records

- Genomic sequencing results stored in a Common location

Edge2012

# Single Writer Scenarios Asynchronous Updates (write, create, remove)

- **Updates at cache site are pushed back lazily**
  - Masks the latency of the WAN
  - Data is written to SONAS at cache site synchronously
    - GW node queues the update for later execution
  - Performance identical to a local file system update

- **Write back is asynchronous**
  - Configurable async delay
  - GW nodes queue updates and write back to home as network bandwidth permits
  - Write back can coalesce updates and accommodate out-of-order and parallel writes to files and directories, thus maximizing WAN bandwidth utilization

- **Users can force a sync if needed**

Edge2012

# System states and Expiration of Data

- **Staleness Control**
  - Defined based on time since disconnection occurred
  - Once cache is expired, no access is allowed to cache
  - Manual expire/unexpire option for admin
  - Allowed only for RO mode cache
  - Disabled for SW & LU modes as they are the recognized sources of data

- **Disconnected**
  - No connectivity with home
  - Updates are queued
  - Reads are serviced locally
  - All access stops after expiration

- **Resync Needed**
  - An error or corruption detected at home
  - ACE will try to fix some errors
  - Need to retry if error is temporary
  - Needs admin intervention for data corruption error

Edge2012

# Cache Eviction Scenarios

- When is it needed?
  - When cache cluster is smaller than home cluster
  - If data fills up in cache faster than it can push to home.
  - Need to create space for caching other files or space for incoming writes.
  - Eviction is linked with fileset quotas.

- Cache eviction is triggered automatically
  - fileset usage level goes above fileset soft quota limits.
  - Eviction tries to bring down usage below 5-10% of quota limits by de-allocating blocks for files in cache.
  - The dirty files with pending flushes are not looked at by eviction.

- It can be triggered manually also if required

Edge2012

# Use Case: Global Namespace (Mesh)

- ACE can be used as a mesh…that is any fileset in any location can be a cache or a single writer fileset to any other location

- Up to 3000 independent filesets per file system

SONAS2.ibm.com

Clients connect to:
SONAS:/data1
SONAS:/data2
SONAS:/data3
SONAS:/data4
SONAS:/data5
SONAS:/data6

**File System: store2**

Cache Filesets:
/data1
/data2

Local Filesets:
/data3
/data4

Cache Filesets:
/data5
/data6

SONAS1.ibm.com

**File System: store1**

Local Filesets:
/data1
/data2

Cache Filesets:
/data3
/data4

Cache Filesets:
/data5
/data6

Clients connect to:
SONAS:/data1
SONAS:/data2
SONAS:/data3
SONAS:/data4
SONAS:/data5
SONAS:/data6

SONAS3.ibm.com

**File System: store2**

Cache Filesets:
/data1
/data2

Cache Filesets:
/data3
/data4

Local Filesets:
/data5
/data6

Clients connect to:
SONAS:/data1
SONAS:/data2
SONAS:/data3
SONAS:/data4
SONAS:/data5
SONAS:/data6

■Every fileset is accessible from all sites

■Each cache site will export same namespace view

Edge2012

# IBM SONAS with Active Cloud Engine:
# Example use cases

| | |
|---|---|
| **Surveillance and Monitoring** | High resolution image capture, data analysis, and results distribution (e.g. Weather) |
| **Digital Media** | High performance, simplified management alternative for widely varying use cases in digital media environments. |
| **University Collaboration** | Link Institutes and Researchers for Data Sharing and analysis |
| **Energy & Geo-Sciences** | Energy exploration and geo-sciences require huge addressable namespaces for analytics and very high performance. |
| **CAE** | Auto / Aero / Electronics design processes experiencing rapid file-centric storage growth as simulation expands. |

23

Edge2012

# Real Time Image Data Capture/ Surveillance and Monitoring

- High Resolution weather images are recorded from planes or satellites

- Information is directly loaded into SONAS systems at remote locations

- IBM Active Cloud Engine allows transfer of these images to a Central Repository

- Weather experts at other sites can then access the data from the Central Repository, perform analysis and provide guidance and information to be disseminated to others

Edge2012

# Example Applying ACE caching technology to Surveillance / Monitoring environment



Pixel Processing Catalog Server Etc.

**Data is written to SONAS cache**

CIFS/NFS

**User data request at SONAS 1 returns data from cache to user**

**If data is not in local cache, SONAS 1 fetches transparently from SONAS 2**

**File is pushed via ACE (async) to SONAS 2 Home site**

**Async replication copies data to backup system**

Remote
Cache Site

Home site

Centralized (Home)

SONAS 1
SONAS 2
SONAS 3

# Making the Cloud *Active* - Write a File

**Cache cluster user writes a file on February 1st that is copied immediately or via a schedule to Home Site**

Cache

Home

/

/ABC  /DEF  /GHI  /JKL  /Local  /MNO

/Jan  /Feb  /Mar

...

File1

File1
File2
File3
File4
File5
File6

Files in blue represent inodes in the global name space. The files *appears local*, but are **not**

/

/ABC  /DEF  /GHI  /JKL  /Local  /MNO

/Jan  /Feb  /Mar

...

File1
File2
File3
File4
File5
File6

Edge2012

# Read a File from "Home"

**Cache user "double clicks" on /JKL/Jan/file3 in his directory that appears local but is actually at the Home site**



Cache

Home

File "is" now local for future access

Edge2012

# Cache Cluster Pre-population

**Policy change to pre-fill the Cache with new version of data for "MNO"**

# Clearing the Cache on policy

**Clear the Cache of all files not modified in the last 30 days**



Cache

Home

ACE scans mtime metadata for files

© 2012 IBM Corporation

Edge2012

# Archive to tape

**Archive data 30 days old to tape by policy**

# Demand File Pull

**Cache User opens up /ABC/File6**



File 6 "appears" local

© 2012 IBM Corporation

Edge2012

# Peer Snapshot Consistent Replication

**(1)** Take a fileset snapshot at cache cluster site.
Mark the point in time in the "write-back queue"

**(2)** Push all updates up to point in time marker

**(3)** Take a snapshot of fileset at the home
Update management tool state of last snapshot time and ids

SONAS Home Cluster

Multi site snapshot scripts

**Push all updates asynchronously Continuous replication with snapshot support**

SONAS Cache Cluster

Edge2012

# IBM SONAS Global Policy Management

| Filename | type | Size | Owner | Creation Timestamp | Last Access Timestamp | Last modified Timestamp | Permissions | World Coordinates | Acquisition time | Country code |
|---|---|---|---|---|---|---|---|---|---|---|
| 783173 | PNG | 10M | User1 | Apr 1, 2009 | Jun 1, 2009 | Apr 1, 2009 | RXX | Latitude/ Longitude | 2009-3-30 T 10:45 | Brazil |
| 273653 | RAW | 100G | User9 | Jan 2, 2011 | Feb 6, 2011 | Jan 15, 2011 | RWX | Latitude/ Longitude | 2010-10-30 T 10:45 | Canada |

Find all the files of "filetype" PNG where "creation timestamp" > than 1 year old and "last access" > 6 months" and move them to tier 3

Specific Metadata



Tier 1    Tier 2    Tier 3

Edge2012

# IBM SONAS Application Specific Storage

| Filename | type | Size | Owner | Creation Timestamp | Last Access Timestamp | Last modified Timestamp | Permissions | World Coordinates | Acquisition time | Country code |
|----------|------|------|-------|--------------------|-----------------------|-------------------------|-------------|-------------------|------------------|--------------|
| 783173 | PNG | 10M | User1 | Apr 1, 2009 | Jun 1, 2009 | Apr 1, 2009 | RXX | Latitude/ Longitude | 2009-3-30 T 10:45 | Brazil |
| 273653 | RAW | 100G | User9 | Jan 2, 2011 | Feb 6, 2011 | Jan 15, 2011 | RWX | Latitude/ Longitude | 2010-10-30 T 10:45 | Canada |

Find all the files where country code="Canada" and move to tier 1

Specific Metadata

With global namespace, **pathnames never change**

Tier 1      Tier 2      Tier 3

273653

783173

Edge2012

# Global Policy Management using IBM Active Cloud Engine

| Filename | type | Size | Owner | Creation Timestamp | Last Access Timestamp | Last modified Timestamp | Permissions | World Coordinates | Acquisition time | Country code |
|----------|------|------|-------|--------------------|-----------------------|-------------------------|-------------|-------------------|------------------|--------------|
| 783173 | PNG | 10M | User1 | Apr 1, 2009 | Jun 1, 2009 | Apr 1, 2009 | RXX | Latitude/ Longitude | 2009-3-30 T 10:45 | Brazil |
| 273653 | RAW | 100G | User9 | Jan 2, 2011 | Feb 6, 2011 | Jan 15, 2011 | RWX | Latitude/ Longitude | 2010-10-30 T 10:45 | Canada |

Replicate all the files where "Country Code"= Canada where "creation timestamp" < than 1 year old and push to "Brazil"

Specific Metadata

With global namespace, **pathnames never change**

Tier 1     Tier 2     Tier 3

783173

273653

Edge2012

# State Government Consolidation

- Customer looking to link all key agencies across the state to capture key data at each city

- Goal is to consolidate local servers in each city.

- Lots of office documents in each agency and city



State Government

**Solution**

- Solution consists of 2 SONAS gateways with V7000 storage located at the State Capitol

- ACE will be used in an exclusive writer mode to move files to the Capital

- Eliminates repetitive "auto-save" operations creating delays since ACE provides local access speeds

- Eventually, V7000 Unified will be the future remote site deployments at key cities for cost optimization

Edge2012

# Financial – Workload Optimization

- Interested in building a file-based Content Distributed Network to enhance workload distribution and placement capabilities

- Ability to distribute workload in a data center today and to various compute farms - move data to applications

Financial Institution

- Meta data (date, access time, permissions, and support of policies for data placement and data migration and replication.

- Central point of management, but with multiple copies in different locations

- Data resides at a central Data Center but seamlessly distributes data to edge cache sites

- Supports copying of data at remote sites back to home cluster

## Solution

- Solution consists of 2 SONAS gateways currently in an ACE relationship

- Testing underway demonstrating how ACE provides the key capabilities listed above

Edge2012

# Statements of Direction* for IBM Smarter Storage - An integral part of IBM Smarter Computing

- IBM Active Cloud Engine™ will provide NAS Virtualization capabilities to enable transparent file migrations from selected NAS systems into SONAS or Storwize V7000 Unified.

- IBM Active Cloud Engine will allow users at geographically dispersed locations to both share and modify the same file (SONAS/ Storwize V7000 Unified)

- Using cloud enhancements and self-service cloud portal for SONAS and Storwize V7000 Unified, enterprises can implement a private cloud storage service where users, with a few clicks, can request & receive storage capacity, share files with other users, and administrators can easily monitor & report usage.

*Statements of IBM's future plans and direction are provided for informational purposes only. Plans and directions are subject to change without notice

Edge2012

# IBM Active Cloud Engine
## Statements of Direction*

## IBM Active Cloud Engine<sup>TM</sup> will help provide:

- **Virtualization capabilities while facilitating transparent file migrations from selected NAS systems to IBM SONAS and Storwize V7000 Unified**

- **Investment protection as clients can continue to leverage their existing NAS assets**

- **Improved productivity and true collaboration where users at geographically dispersed locations can both share and modify the same file.**



**Application Servers**

SONAS
or
Storwize
V7000 Unified

**Global namespace**

| Protocols | Active Cloud Engine | Availability |
|---|---|---|
| CIFS NFS HTTP FTP SCP | | Data Migration Replication Backup |

Automated data movement
Central policy engine

• **Virtualize selected NAS devices**

**Existing NAS filers**

Edge2012

# Statement of Direction* - Migration of Legacy NAS data

**Local Use Case**

Legacy storage serves as repository – mounted to SONAS

1) Customer requests a file from SONAS storage

2) If Data not already in SONAS, retrieve it from Server/system

3) File is now stored (migrated to) on SONAS for future access

**Legacy NAS**

Celerra NS-960
Celerra NS-480
Celerra NS-120

**SONAS**
Cache Fileset /data1

Local Fileset/data1

**Data Center**

## Definition

- Legacy NAS system is a "home site" to a SONAS "cache cluster"
- Migrate or access any data that Legacy NAS exports using NFSv3 protocol
- Will initially support read operations only

## Benefits

- Can provide transparent migration of data to SONAS or Storwize V7000 Unified while still allowing normal operations
- Reduces "down time" for users during migration period to just migrating ACLs, permissions, etc.

Edge2012

*Statements of IBM's future plans and direction are provided for informational purposes only. Plans and directions are subject to change without notice

# Customer Example*: Migration of Legacy NAS data using Active Cloud Engine

## Research Institute

- Existing NFS & CIFS environment, approximately 1PB and growing rapidly

- 10,000 students, employees/professors and associated research projects

- Mixed environment of NAS devices as well as approximately mixed protocols access via NFS only, CIFS only, and Mixed

- Original estimates: Over 6 months to migrate using current standard industry migration tools

- ACE Based Proposal for data migration

  – ACE can migrate data faster than industry standard tools

  – ACE can help migrate data from legacy NAS and servers via NFSv3 protocol

  – Goal is reduced overall time to migrate data and reduction of system down time

- **Future Lab Based Services Engagements required**

**Local Use Case**



Legacy NAS Server

SONAS Cache Fileset /data1

Local Fileset          Data Center

*Statements of IBM's future plans and direction are provided for informational purposes only. Plans and directions are subject to change without notice

# Thank You!

Session:     sVC32
Presenter:   Stephen Edel

# Trademarks and disclaimers

ZSP03490-USEN-00