



ScaleOut Software

Implementing Operational Intelligence
Using
In-Memory Computing



William L. Bain (wbain@scaleoutsoftware.com)

June 29, 2015

- What is Operational Intelligence?
- Example: Tracking Set-Top Boxes
- Using an In-Memory Data Grid (IMDG) for Operational Intelligence
 - Tracking and analyzing live data
 - Comparison to Spark
- Implementing OI Using Data-Parallel Computing in an IMDG
- A Detailed OI Example in Financial Services
 - Code Samples in Java
- Implementing MapReduce on an IMDG
- Optimizing MapReduce for OI
- Integrating Operational and Business Intelligence

About ScaleOut Software

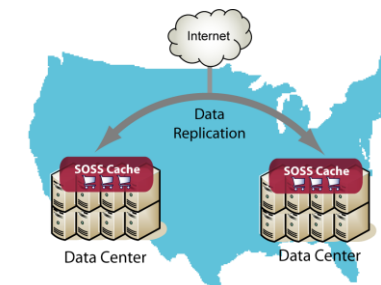
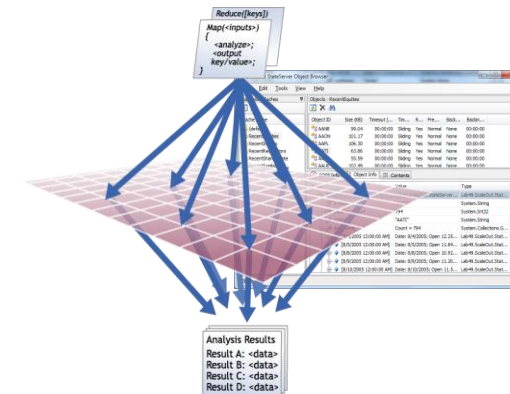
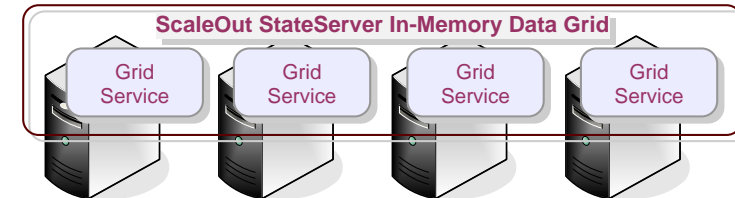


- Develops and markets **In-Memory Data Grids**, software middleware for:
 - **Scaling application performance** and
 - **Providing operational intelligence** using
 - **In-memory data storage and computing**
- Dr. William Bain, Founder & CEO
 - Career focused on parallel computing – Bell Labs, Intel, Microsoft
 - 3 prior start-ups, last acquired by Microsoft and product now ships as Network Load Balancing in Windows Server
- Ten years in the market; 400+ customers, 10,000+ servers
- Sample customers:



- **ScaleOut StateServer® (SOSS)**
 - In-Memory Data Grid for Windows and Linux
 - Scales application performance
 - Industry-leading performance and ease of use
- **ScaleOut ComputeServer™** adds
 - Operational intelligence for “live” data
 - Comprehensive management tools
- **ScaleOut hServer®**
 - Full Hadoop Map/Reduce engine (>40X faster*)
 - Hadoop Map/Reduce on live, in-memory data
- **ScaleOut GeoServer®**
 - WAN based data replication for DR
 - Global data access and synchronization

*in benchmark testing



In-Memory Computing Is Not New

- 1980's: SIMD Systems, Caltech Cosmic Cube



**Thinking Machines
Connection Machine 5**

- 1990's: Commercial Parallel Supercomputers



**Intel
IPSC-2**



**IBM
SP1**

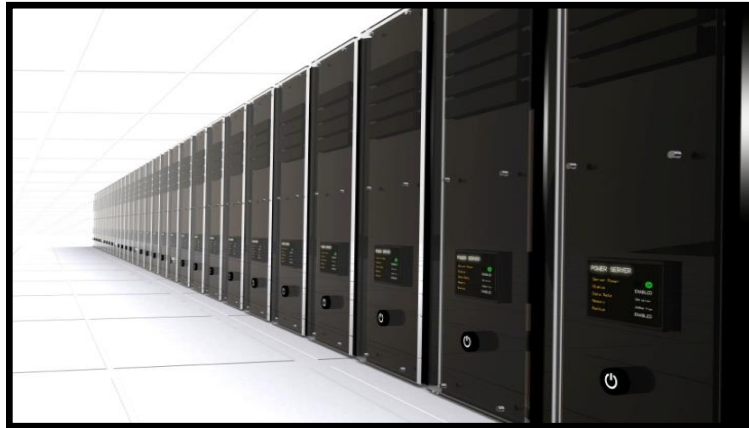
What's New: IMC on Commodity Hardware

- 1990's – early 2000's: HPC on Clusters



**HP
Blade
Servers**

- Since ~2005: Public Clouds



**Amazon EC2,
Windows Azure**

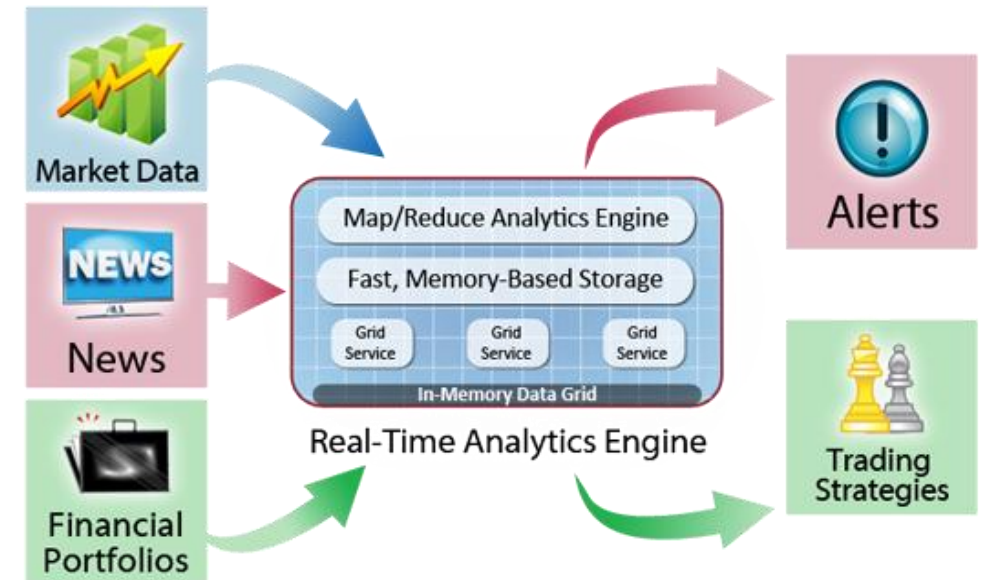
Introductory Video: What is Operational Intelligence

<https://www.youtube.com/watch?v=H6OFzdIEy-g&feature=youtu.be>

Goal: Provide *immediate* (sub-second) feedback to a system handling live data.

A few example use cases requiring immediate feedback within a live system:

- **Ecommerce:** personalized, real-time recommendations
- **Healthcare:** patient monitoring, predictive treatment
- **Equity trading:** minimize risk during a trading day
- **Reservations systems:** identify issues, reroute, etc.
- **Credit cards & wire transfers:** detect fraud in real time
- **IoT, Smart grids:** optimize power distribution & detect issues



Operational vs Business Intelligence

Operational Intelligence



Real-time

Live data sets

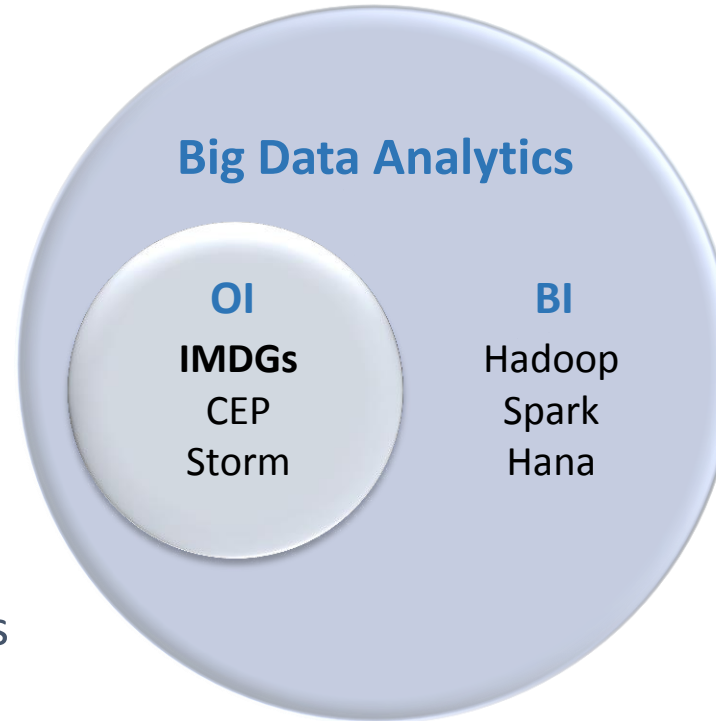
Gigabytes to terabytes

In-memory storage

Sub-second to seconds

Best uses:

- Tracking live data
- Immediately identifying trends and capturing opportunities
- Providing immediate feedback



Business Intelligence

Batch

Static data sets

Petabytes

Disk storage

Minutes to hours

Best uses:

- Analyzing warehoused data
- Mining for long-term trends

Example: Enhancing Cable TV Experience

- **Goals:**

- Make real-time, personalized upsell offers
- Immediately respond to service issues
- Detect and manage network hot spots
- Track aggregate behavior to identify patterns, e.g.:
 - Total instantaneous incoming event rate
 - Most popular programs and # viewers by zip code



©2011 Tammy Bruce presents LiveWire

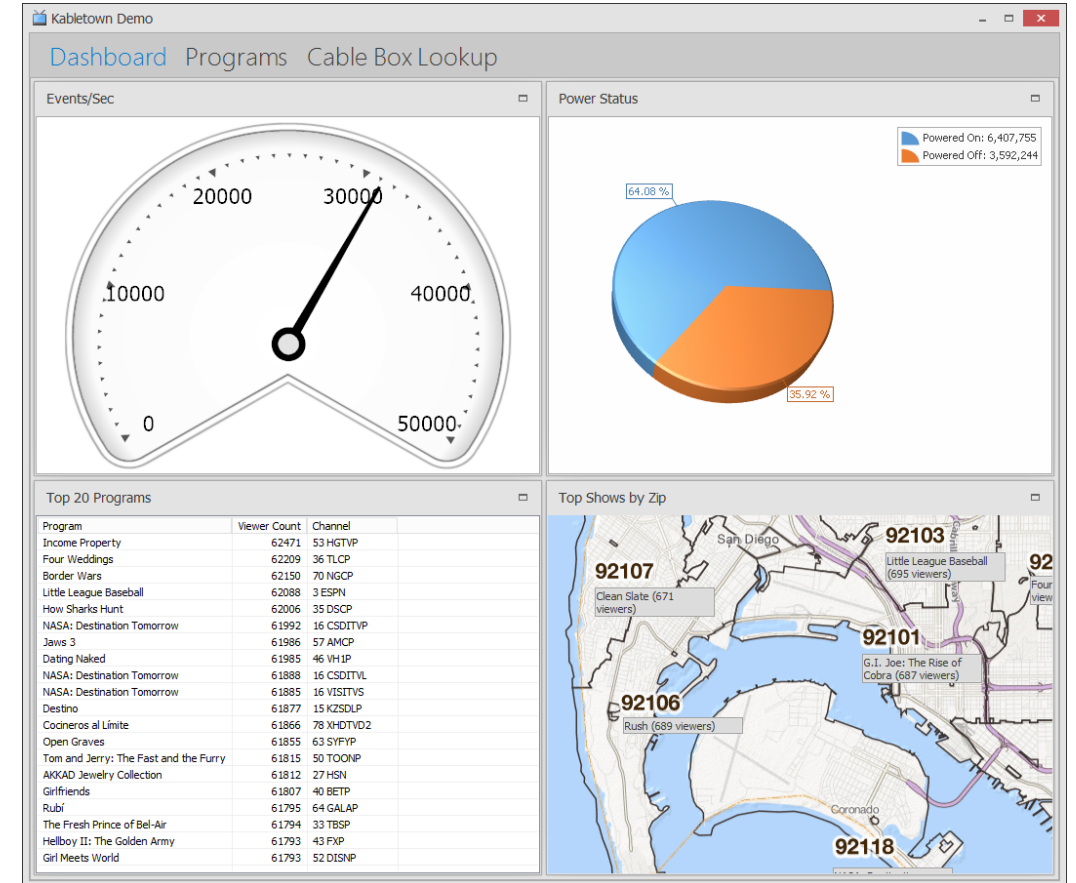
- **Requirements:**

- Track events from 10M set-top boxes with 25K events/sec (2.2B/day)
- Correlate, cleanse, and enrich events per rules (e.g. ignore fast channel switches, match channels to programs)
- Be able to feed enriched events to recommendation engine within 5 seconds
- Immediately examine any set-top box (e.g., box status) & track aggregate statistics

The Result: An OI Platform

Based on a simulated workload for San Diego metropolitan area:

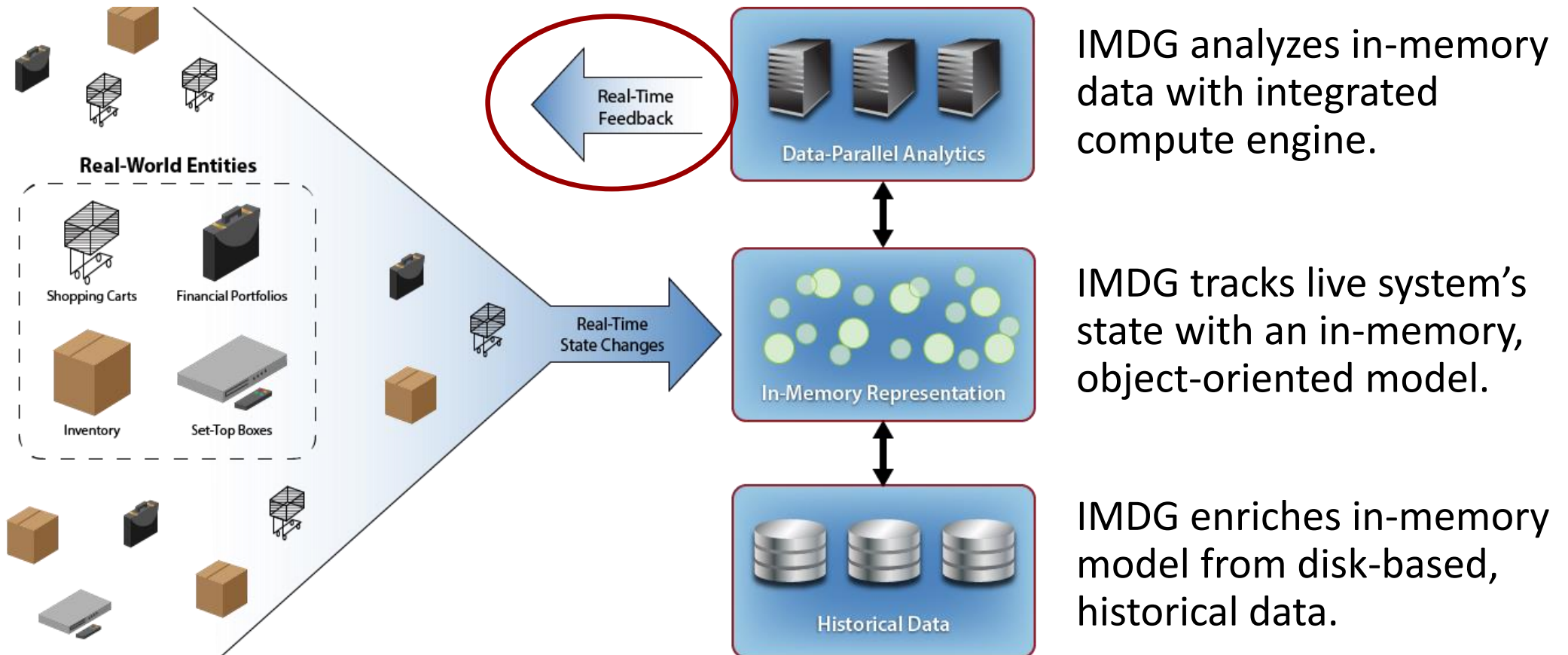
- Continuously correlates and cleanses telemetry from 10M simulated set-top boxes (from synthetic load generator)
- Processes more than 30K events/second
- Enriches events with program information every second
- Tracks aggregate statistics (e.g., top 10 programs by zip code) every 10 seconds



Real-Time Dashboard

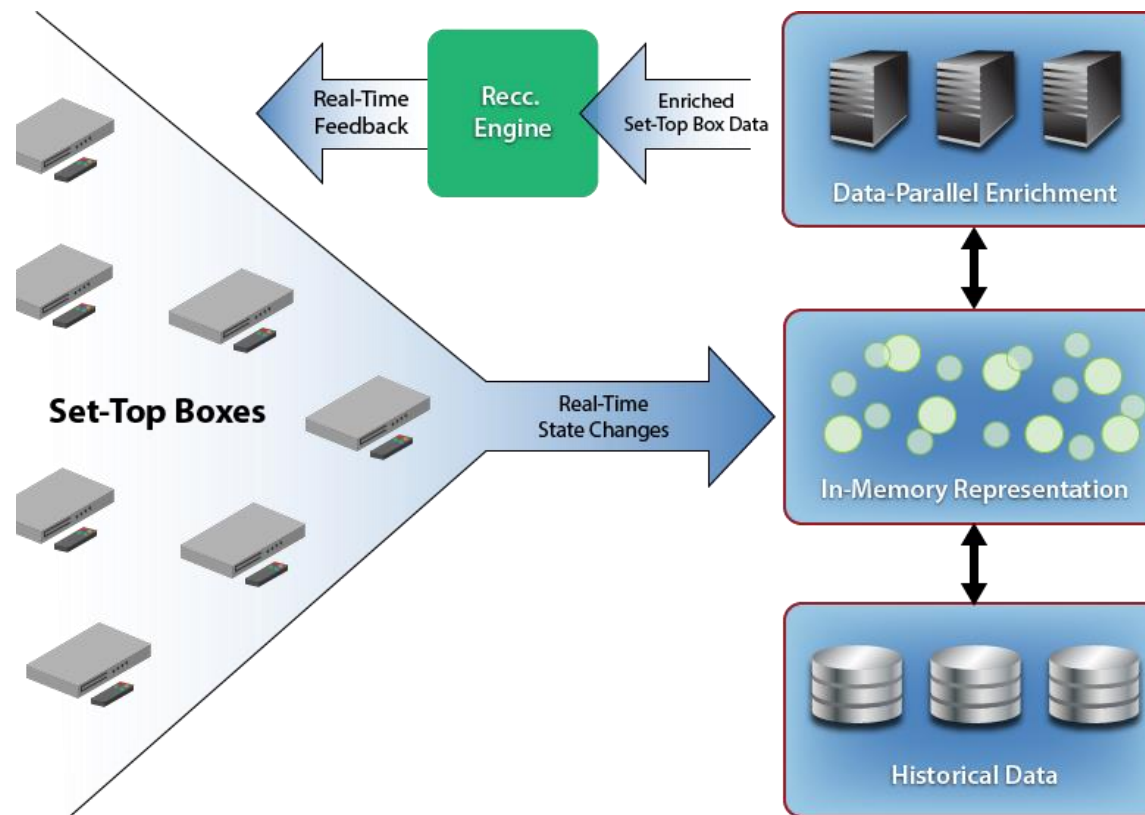
Using an IMDG to Implement OI

- IMDG models and tracks the state of a “live” system.
- IMDG analyzes the system’s state in parallel and provides real-time feedback.

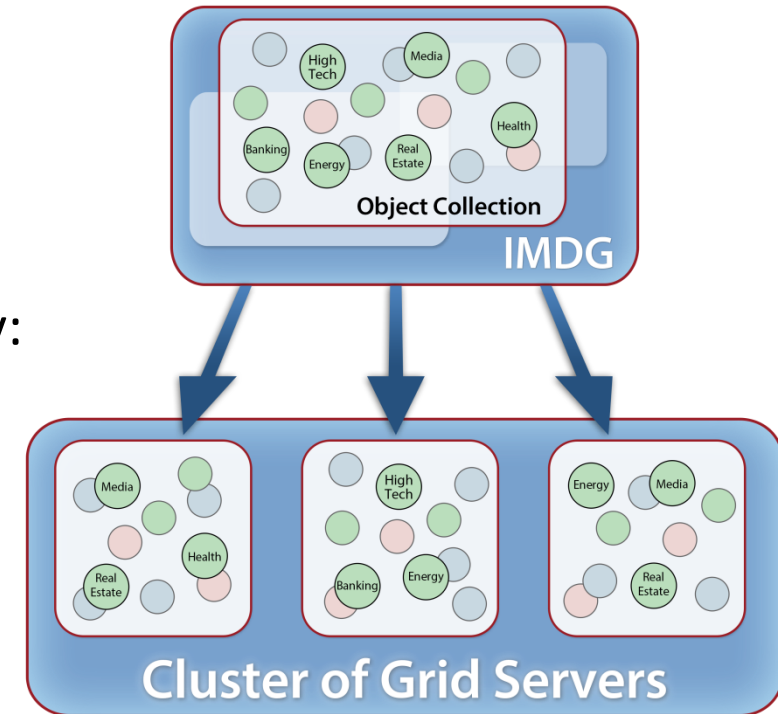


Example: Tracking Set-TopBoxes

- Each set-top box is represented as an object in the IMDG
- Object holds raw & enriched event streams, viewer parameters, and statistics
- IMDG captures incoming events by updating objects
- IMDG uses data-parallel computation to:
 - immediately enrich box objects to generate alerts to recommendation engine, and
 - continuously collect and report global statistics



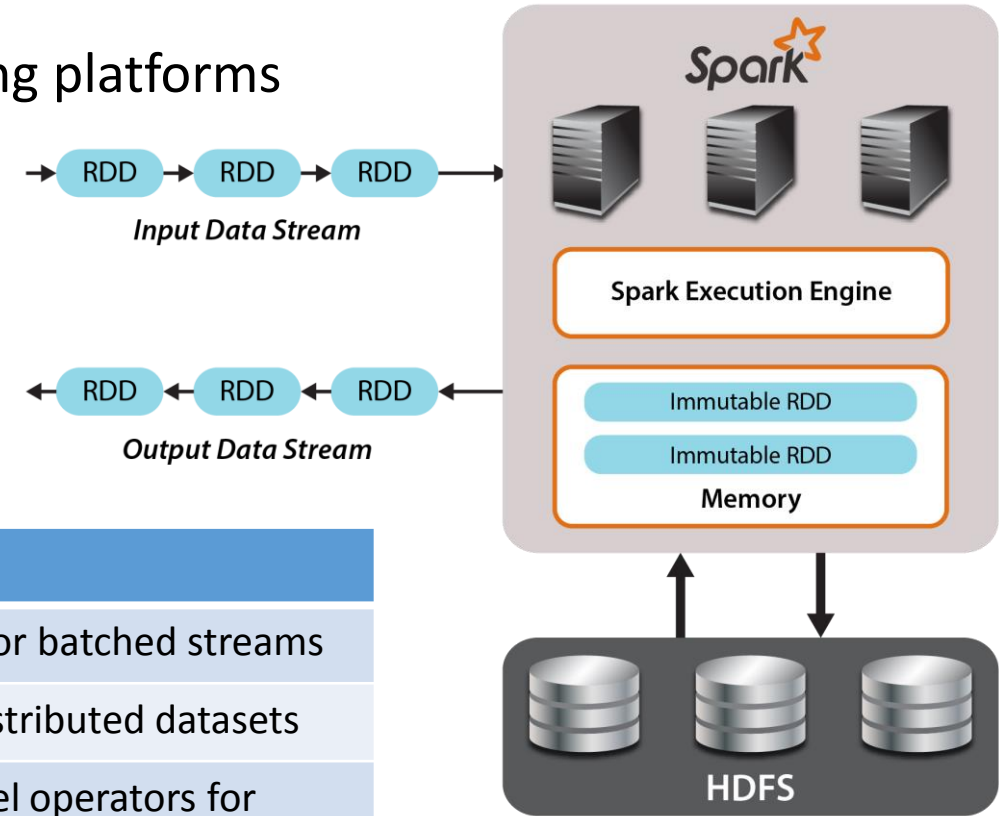
- In-memory data grid (IMDG) provides scalable, hi av storage for live data:
 - Designed to manage business logic state:
 - Object-oriented collections by type
 - Create/read/update/delete APIs for Java/C#/C++
 - Parallel query by object properties
 - Data shared by multiple clients
 - Designed for transparent scalability and high availability:
 - Automatic load-balancing across commodity servers
 - Automatic data replication, failure detection, and recovery
- IMDGs provide ideal platform for operational intelligence:
 - Easy to track live systems with large workloads
 - Appropriate availability model for production deployments



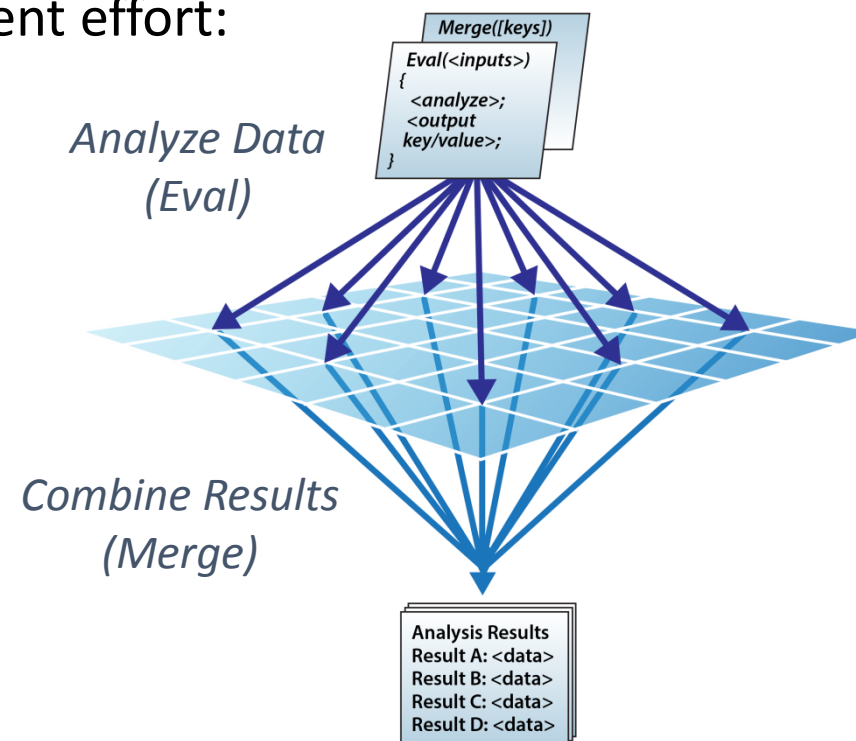
Comparing IMDGs to Spark

- On the surface, both are surprisingly similar:
 - Both designed as scalable, in-memory computing platforms
 - Both implement data-parallel operators
 - Both can handle streaming data
- But there are key differences that impact use for operational intelligence:

| | IMDGs | Spark |
|-----------------------------|---|---------------------------------------|
| Best use | Live, operational data | Static data or batched streams |
| In-memory model | Object-oriented collections | Resilient distributed datasets |
| Focus of APIs | CRUD, eventing, data-parallel computing | Data-parallel operators for analytics |
| High availability tradeoffs | Data replication for fast recovery | Lineage for max performance |



- IMDGs provide powerful, cost-effective platform for data-parallel computing:
 - Enable integrated computing with data storage:
 - Take advantage of cluster's commodity servers and cores.
 - Avoid delays due to data motion (both to/from disk and across network).
 - Leverage object-oriented model to minimize development effort:
 - Easily define data-parallel tasks as class methods.
 - Easily specify domain as object collection.
- Example: “Parallel Method Invocation” (PMI):
 - Object-oriented version of standard HPC model
 - Runs class methods in parallel across cluster.
 - Selects objects using parallel query of obj. collection.
 - Serves as a platform for implementing MapReduce and other data-parallel operators

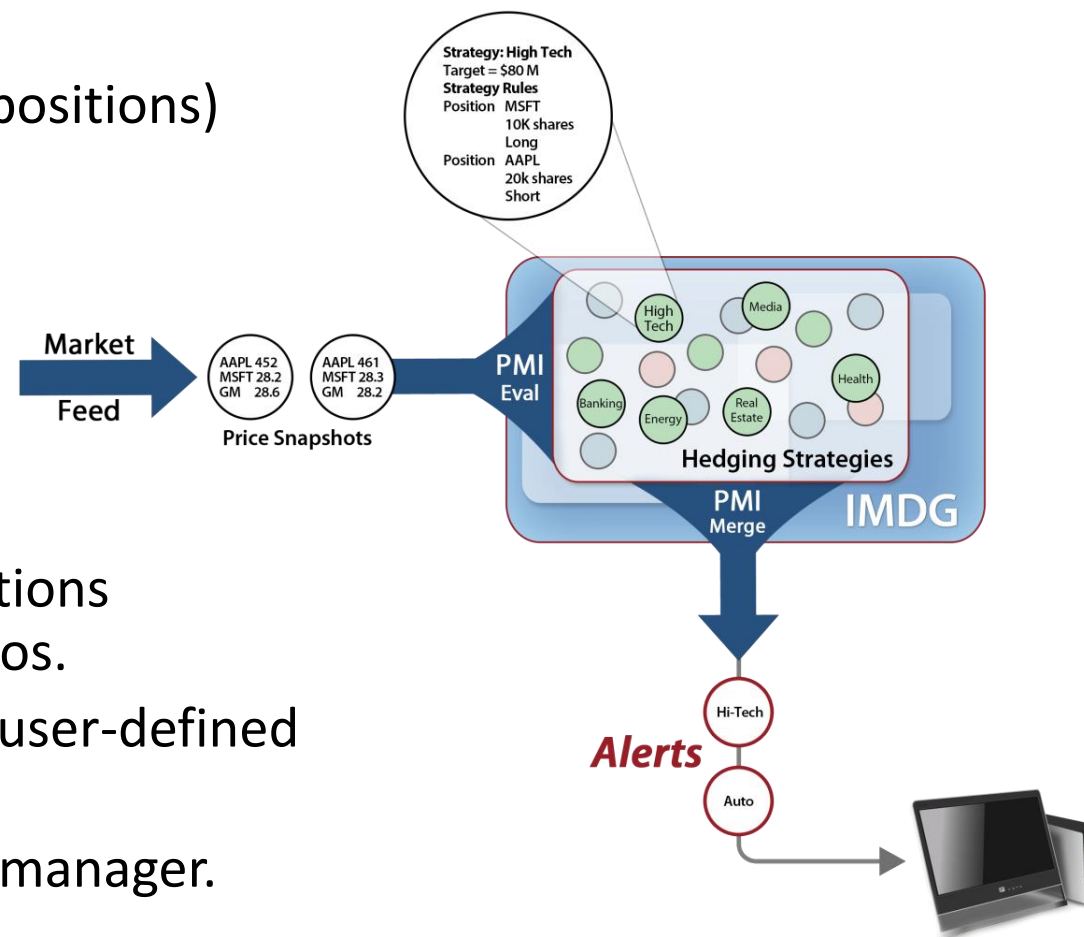


PMI Example: OI in Financial Services

- **Goal:** track market price fluctuations for a hedge fund and keep portfolios in balance.

- **How:**

- Keep portfolios of stocks (long and short positions) in object collection within IMDG.
- Collect market price changes in one-second snapshots.
- Define a method which applies a snapshot to a portfolio and optionally generates an alert to rebalance.
- Perform repeated parallel method invocations on a selected (i.e., queried) set of portfolios.
- Combine alerts in parallel using a second user-defined method.
- Report alerts to UI every second for fund manager.



- Simplified example of a portfolio class (Java):
 - Note: some properties are made query-able.
 - Note: the evalPositions method analyzes the portfolio for a market snapshot.

```
public class Portfolio {  
    private long          id;  
    private Set<Stock>    longPositions;  
    private Set<Stock>    shortPositions;  
    private double        totalValue;  
    private Region        region;  
    private boolean       alerted;    // alert for trading  
  
    @SossIndexAttribute      // query-able property  
    public double getTotalValue() {...}  
    @SossIndexAttribute      // query-able property  
    public Region getRegion() {...}  
  
    public Set<Long> evalPositions(MarketSnapshot ms) {...};  
}
```



Strategy: High Tech
Target = \$80 M
Strategy Rules
Position MSFT
10K shares
Long
Position AAPL
20k shares
Short

- Implement PMI interface to define methods for analyzing each object and for merging the results:

```
public class PortfolioAnalysis implements
    Invokable<Portfolio, MarketSnapshot, Set<Long>>
{
    public Set<Long> eval(Portfolio p, MarketSnapshot ms)
        throws InvokeException {

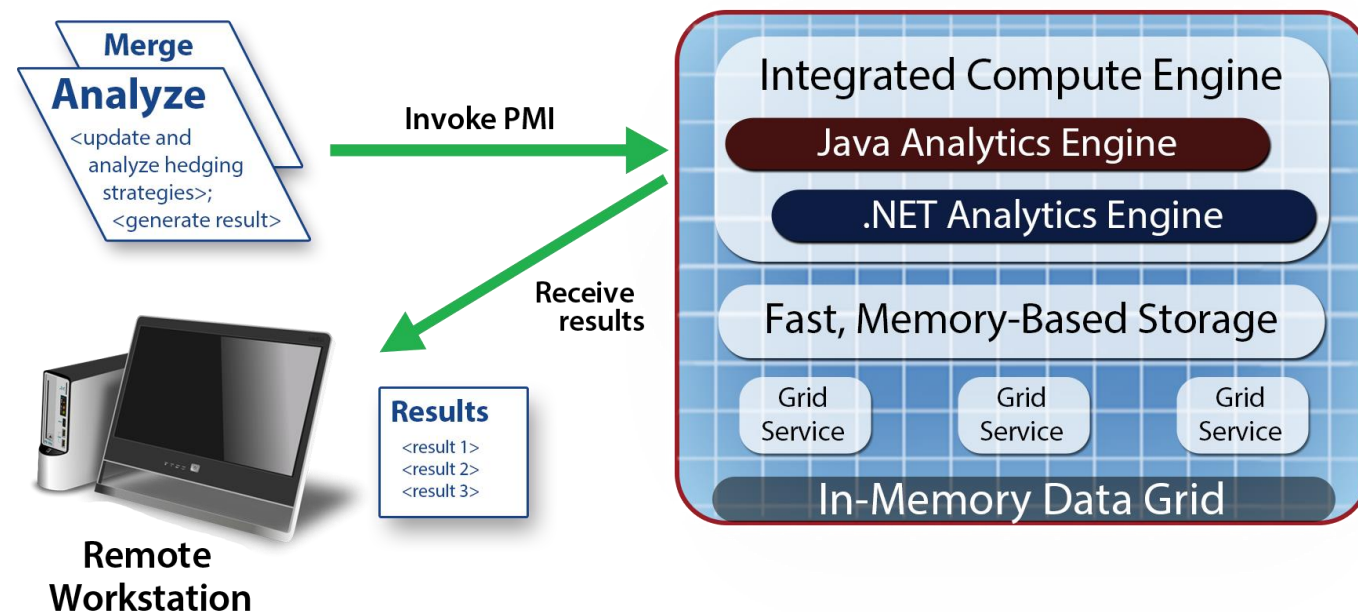
        // update portfolio and return id if alerted:
        return p.evalPositions(ms);
    }

    public Set<Long> merge(Set<Long> set1, Set<Long> set2)
        throws InvokeException {

        set1.addAll(set2);
        return set1; // merged set of alerted portfolio ids
    }
}
```

Running the Analysis

- PMI can be run from a remote workstation.
- IMDG ships code and libraries to cluster of servers:
 - Execution environment can be pre-staged for fast startup.
- In-line execution minimizes scheduling time.
 - Avoids batch scheduling delays.
- PMI automatically runs in parallel across all grid servers:
 - Uses software multicast to accelerate startup.
 - Passes market snapshot parameter to all servers.
 - Uses all servers and cores to maximize throughput.



- First obtain a reference to the IMDG's object collection of portfolios:

```
NamedCache pset = CacheFactory.getCache("portfolios");
```

- Create an “invocation grid,” a re-usable compute engine for the application:
 - Spawns a JVM on all grid servers and connects them to the in-memory data grid.
 - Stages the application code on all JVMs.
 - Associates the invocation grid with an object collection.

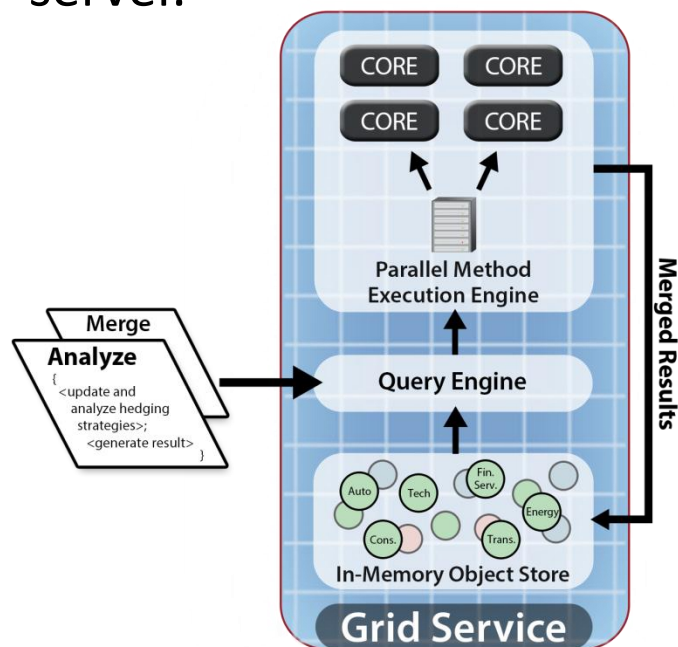
```
InvocationGrid grid = new InvocationGridBuilder("grid")  
    .addClass(DependencyClass.class)  
    .addJar("/path/to/dependency.jar")  
    .setJVMParameters("-Xmx2m")  
    .load();  
  
pset.setInvocationGrid(grid);
```

- Run the PMI on a queried set of objects within the collection:
 - Multicasts the invocation and parameters to all JVMs.
 - Runs the data-parallel computation.
 - Merges the results and returns a final result to the point of call.

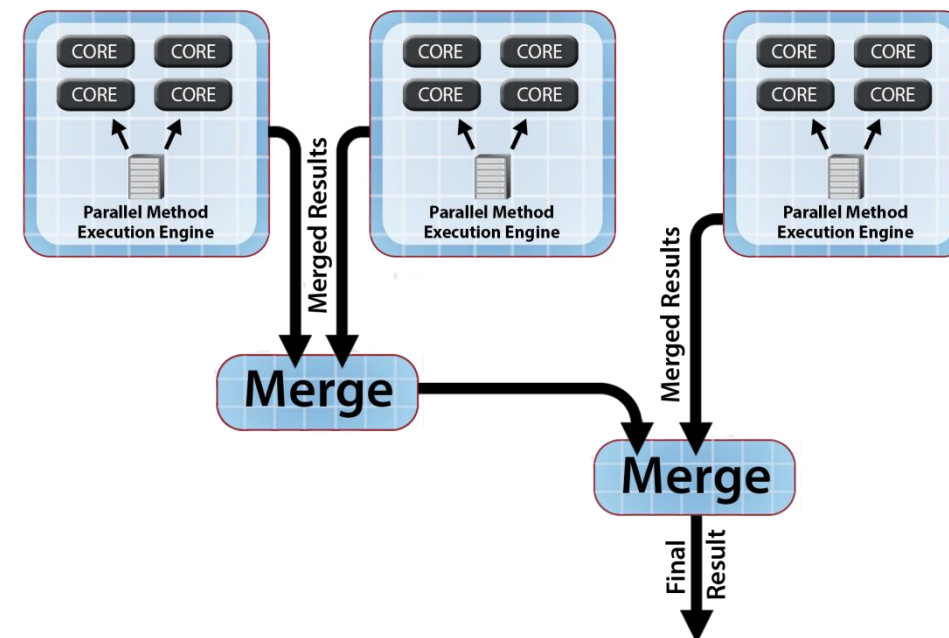
```
InvokeResult alertedPortfolios = pset.invoke(  
    PortfolioAnalysis.class,  
    Portfolio.class,  
    and(greaterThan("totalValue", 1000000), // query spec  
        equals("region", Region.US)),  
    marketSnapshot, // parameters  
    ...  
);  
  
System.out.println("The alerted portfolios are" +  
    alertedPortfolios.getResult());
```

Execution Steps

- **Eval phase:** each server queries local objects and runs eval and merge methods:
 - Note: Accessing local data avoids networking overhead.
 - Completes with one result object per server.

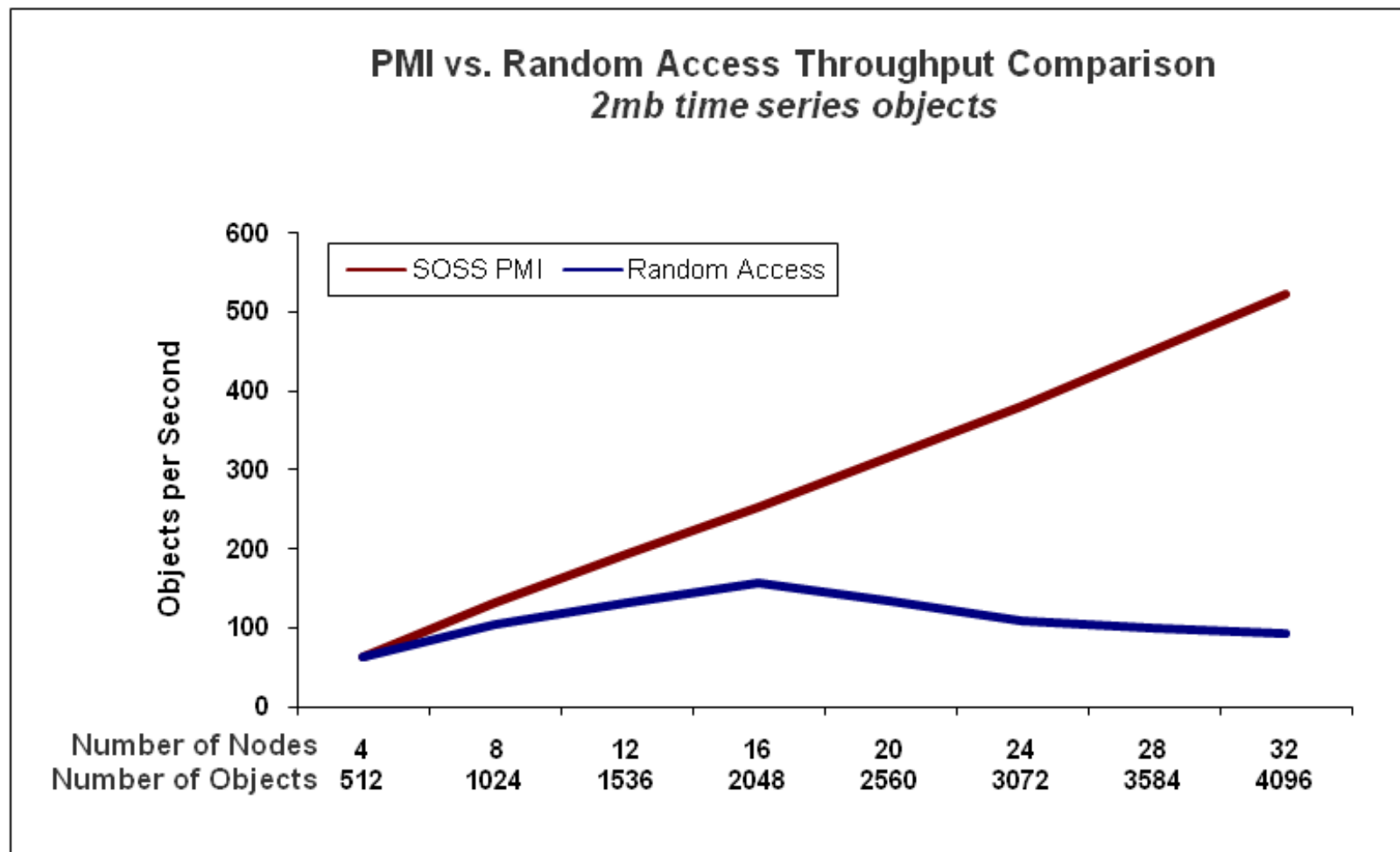


- **Merge phase:** all servers perform distributed merge to create final result:
 - Merge runs in parallel to minimize completion time.
 - Returns final result object to client.



Importance of Avoiding Data Motion

- Local data access enables linear throughput scaling.
- Network access creates a bottleneck that limits throughput.

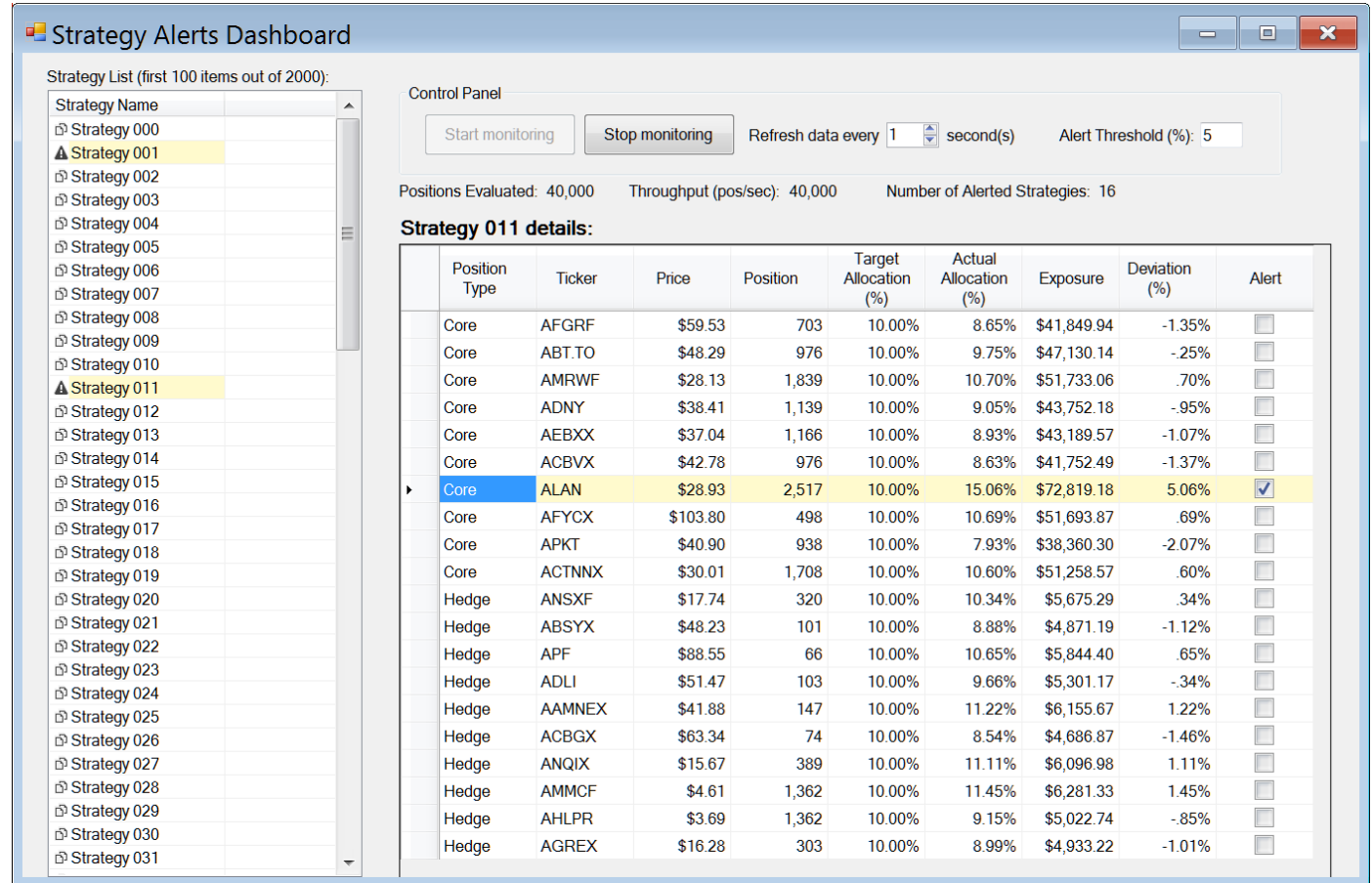


Outputting Continuous Alerts to the UI

- PMI runs every second; it completes in **350 msec.** and immediately refreshes UI.



- UI alerts trader to portfolios that need rebalancing.
- UI allows trader to examine portfolio details and determine specific positions that are out of balance.
- Result: in-memory computing delivers operational intelligence.**



The screenshot shows the "Strategy Alerts Dashboard" window. On the left is a "Strategy List (first 100 items out of 2000):" with a scrollable list of strategies from Strategy 000 to Strategy 031. "Strategy 011" is highlighted with a yellow background. On the right is the "Control Panel" with buttons for "Start monitoring" and "Stop monitoring", a "Refresh data every" dropdown set to "1" second(s), and an "Alert Threshold (%): 5" input field. Below the control panel, summary statistics are shown: "Positions Evaluated: 40,000", "Throughput (pos/sec): 40,000", and "Number of Alerted Strategies: 16". The main section is titled "Strategy 011 details:" and contains a table with 9 columns: Position Type, Ticker, Price, Position, Target Allocation (%), Actual Allocation (%), Exposure, Deviation (%), and Alert. The table lists various positions, with "Core ALAN" highlighted in yellow and its "Alert" checkbox checked.

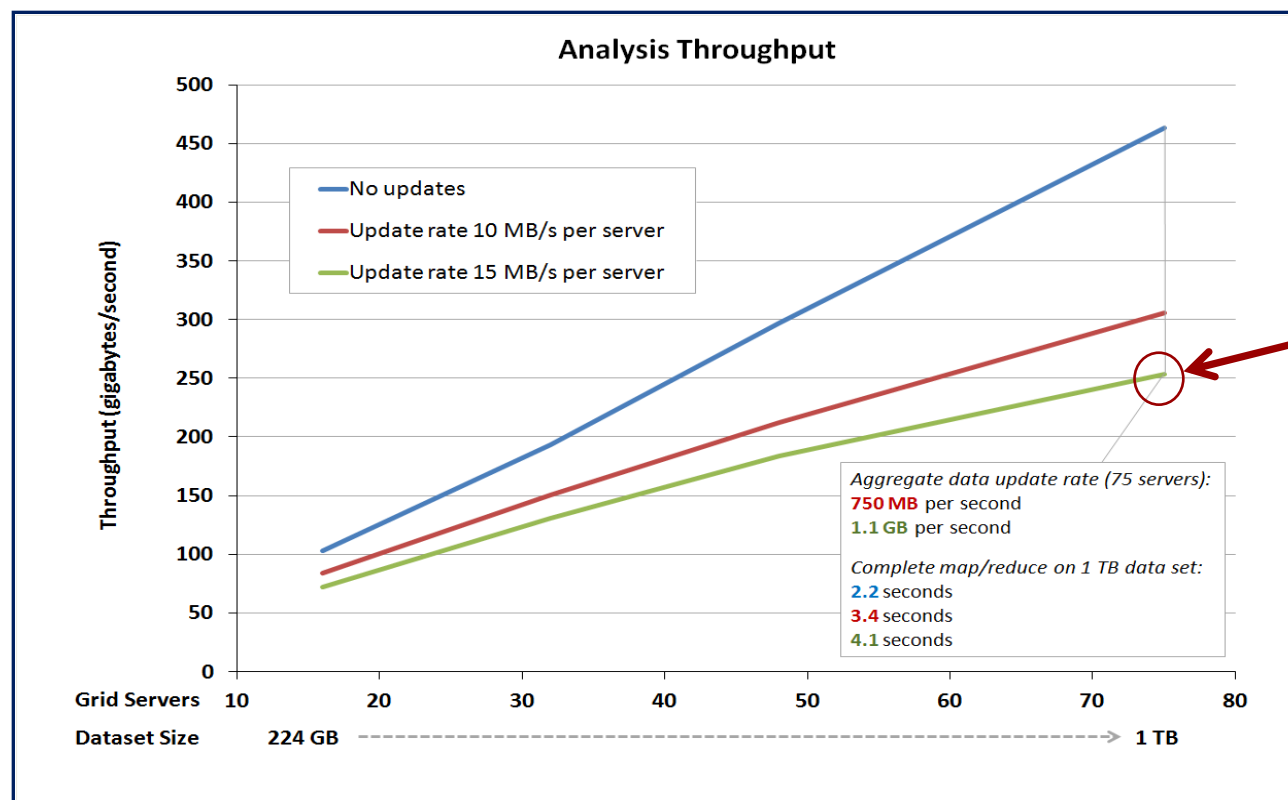
| Position Type | Ticker | Price | Position | Target Allocation (%) | Actual Allocation (%) | Exposure | Deviation (%) | Alert |
|---------------|--------|----------|----------|-----------------------|-----------------------|-------------|---------------|-------------------------------------|
| Core | AFGRF | \$59.53 | 703 | 10.00% | 8.65% | \$41,849.94 | -1.35% | <input type="checkbox"/> |
| Core | ABT.TO | \$48.29 | 976 | 10.00% | 9.75% | \$47,130.14 | -.25% | <input type="checkbox"/> |
| Core | AMRWF | \$28.13 | 1,839 | 10.00% | 10.70% | \$51,733.06 | .70% | <input type="checkbox"/> |
| Core | ADNY | \$38.41 | 1,139 | 10.00% | 9.05% | \$43,752.18 | -.95% | <input type="checkbox"/> |
| Core | AEBXX | \$37.04 | 1,166 | 10.00% | 8.93% | \$43,189.57 | -1.07% | <input type="checkbox"/> |
| Core | ACBVX | \$42.78 | 976 | 10.00% | 8.63% | \$41,752.49 | -1.37% | <input type="checkbox"/> |
| Core | ALAN | \$28.93 | 2,517 | 10.00% | 15.06% | \$72,819.18 | 5.06% | <input checked="" type="checkbox"/> |
| Core | AFYCX | \$103.80 | 498 | 10.00% | 10.69% | \$51,693.87 | .69% | <input type="checkbox"/> |
| Core | APKT | \$40.90 | 938 | 10.00% | 7.93% | \$38,360.30 | -2.07% | <input type="checkbox"/> |
| Core | ACTNNX | \$30.01 | 1,708 | 10.00% | 10.60% | \$51,258.57 | .60% | <input type="checkbox"/> |
| Hedge | ANSXF | \$17.74 | 320 | 10.00% | 10.34% | \$5,675.29 | .34% | <input type="checkbox"/> |
| Hedge | ABSYX | \$48.23 | 101 | 10.00% | 8.88% | \$4,871.19 | -1.12% | <input type="checkbox"/> |
| Hedge | APF | \$88.55 | 66 | 10.00% | 10.65% | \$5,844.40 | .65% | <input type="checkbox"/> |
| Hedge | ADLI | \$51.47 | 103 | 10.00% | 9.66% | \$5,301.17 | -.34% | <input type="checkbox"/> |
| Hedge | AAMNEX | \$41.88 | 147 | 10.00% | 11.22% | \$6,155.67 | 1.22% | <input type="checkbox"/> |
| Hedge | ACBGX | \$63.34 | 74 | 10.00% | 8.54% | \$4,686.87 | -1.46% | <input type="checkbox"/> |
| Hedge | ANQIX | \$15.67 | 389 | 10.00% | 11.11% | \$6,096.98 | 1.11% | <input type="checkbox"/> |
| Hedge | AMMCF | \$4.61 | 1,362 | 10.00% | 11.45% | \$6,281.33 | 1.45% | <input type="checkbox"/> |
| Hedge | AHLPR | \$3.69 | 1,362 | 10.00% | 9.15% | \$5,022.74 | -.85% | <input type="checkbox"/> |
| Hedge | AGREX | \$16.28 | 303 | 10.00% | 8.99% | \$4,933.22 | -1.01% | <input type="checkbox"/> |

Demonstration Video:

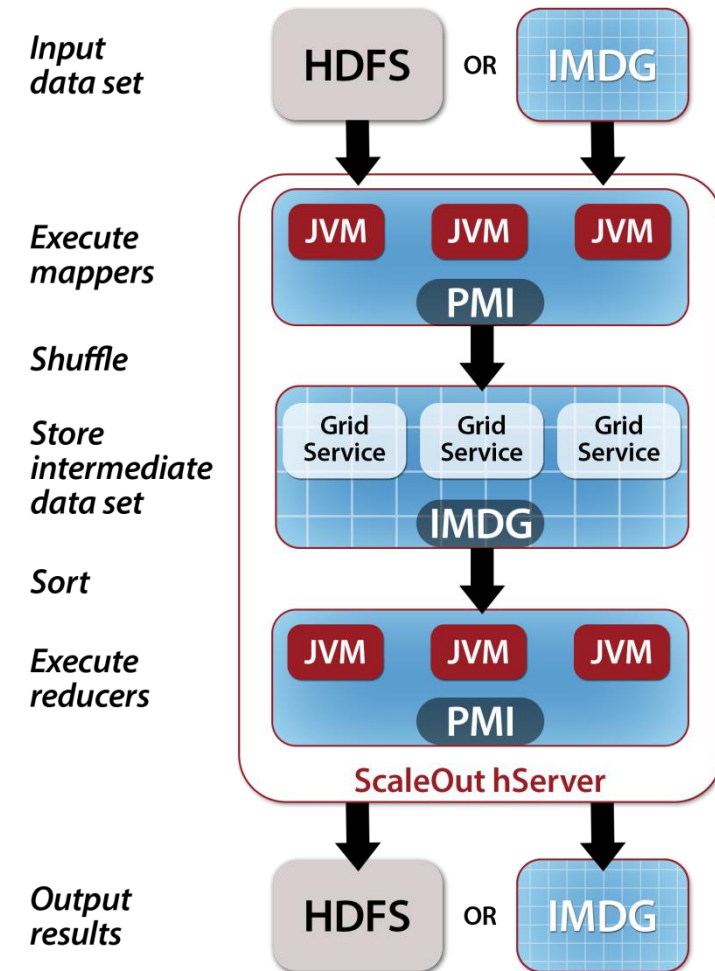
Comparison of PMI to Apache Hadoop

https://www.youtube.com/watch?v=8JTsqp_-Gnw

- Measured a similar financial services application (back testing stock trading strategies on stock histories)
- Hosted IMDG in Amazon EC2 using 75 servers holding **1 TB** of stock history data in memory
- IMDG handled a continuous stream of updates (**1.1 GB/s**)
- Results: analyzed 1 TB in **4.1 seconds** (250 GB/s).
- Observed linear scaling as dataset and update rate grew.

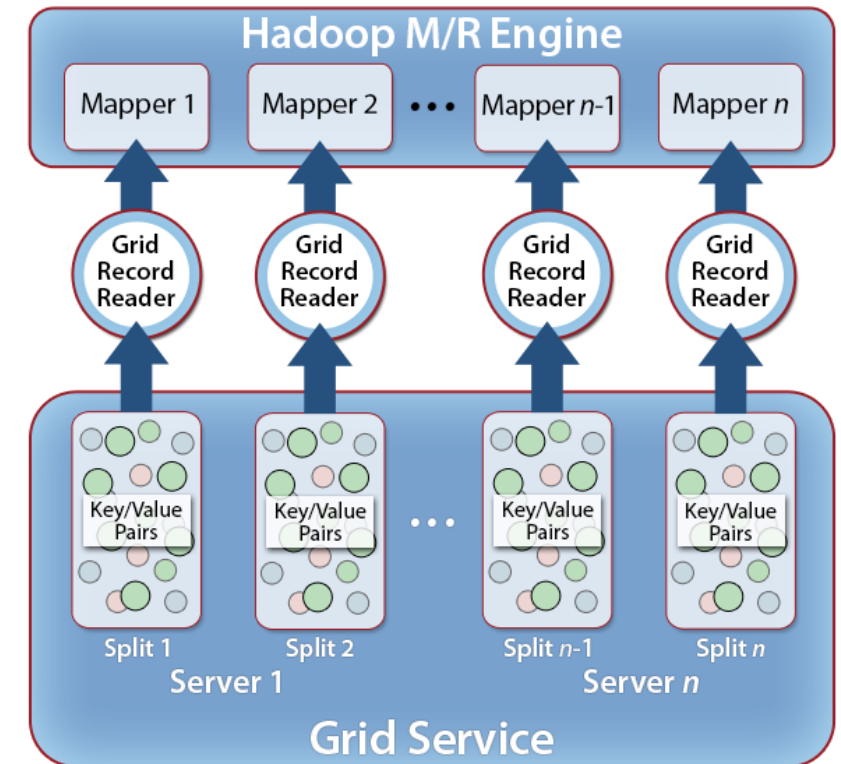


- PMI serves as foundational platform for MapReduce and other parallel operators.
- Implement MapReduce with two PMI phases:
 - Runs standard Hadoop MapReduce applications.
 - Data can be input from either the IMDG or an external data source.
 - Works with any input/output format.
 - IMDG uses PMI phases to invoke the mappers and reducers.
 - Eliminates batch scheduling overhead.
 - Intermediate results are stored within the IMDG.
 - Minimizes data motion in shuffle phase.
 - Allows optional sorting.
 - Note: output of a single reducer/combiner optionally can be globally merged.



MapReduce for OI Requires New Data Model ScaleOut Software

- IMDGs historically implement a feature-rich data model:
 - Efficiently manages large objects (KBs-MBs).
 - Supports object timeouts, locking, query by properties, dependency relationships, etc.
- MapReduce typically targets very large collections of small key/value pairs:
 - Does *not* require rich object semantics.
 - Does require efficient storage (minimum metadata) and highly pipelined access.
- **Solution:** a new IMDG data model for MapReduce:
 - Uses standard Java named map APIs for access.
 - MapReduce uses standard input/output formats.
 - Stores data in chunks and pipelines to/from engine.
 - Automatically defines splits for mappers and holds shuffled data for reducers.



- Integrate in-memory named map with MapReduce to minimize execution time.
- Use new API (**simpleMR** in Java, C#) to simplify apps and remove Hadoop dependencies.

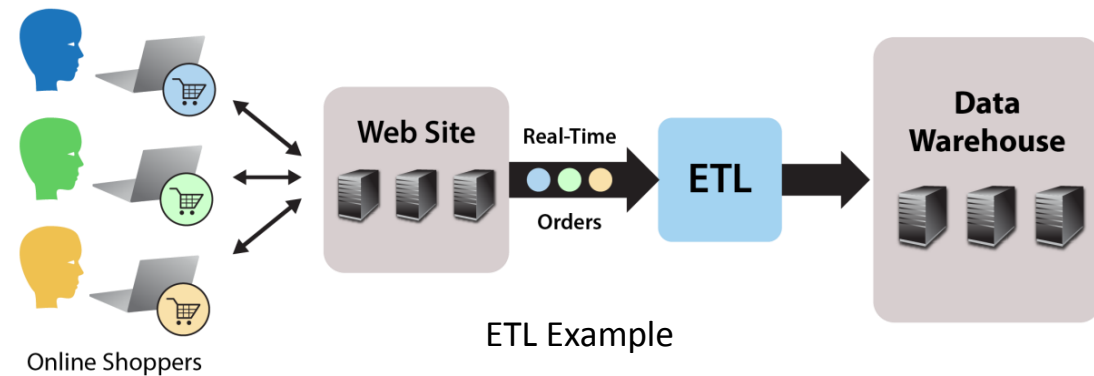
```
public class Mapper : IMapper<int, string, string, int>
{
    void IMapper<int, string, string, int>.Map(int key,
        string value, IContext<string, int> context)
    {
        ...
        context.Emit(Encoding.ASCII.GetString(...), 1);
    }
}

inputMap = new NamedMap<int, string>("Input_Map");
outputMap = new NamedMap<string, int>("Output_Map");

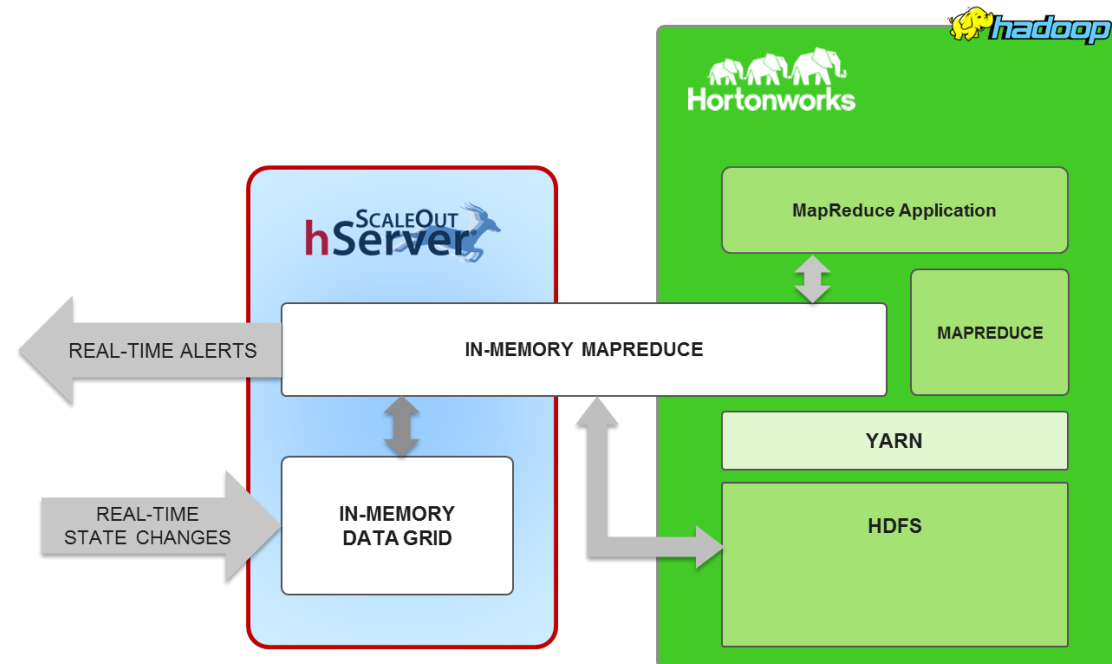
inputMap.RunMapReduce<string, int, string, int>(outputMap,
    new Mapper(), new Combiner(), new Reducer(), ...);
```

Integrating OI and BI in the Data Warehouse ScaleOut Software

- In-memory data grids can add value to a BI platform, e.g.:
 - Transform live data and store in HDFS for analysis.
 - Provide immediate feedback to live system pending deep analysis.

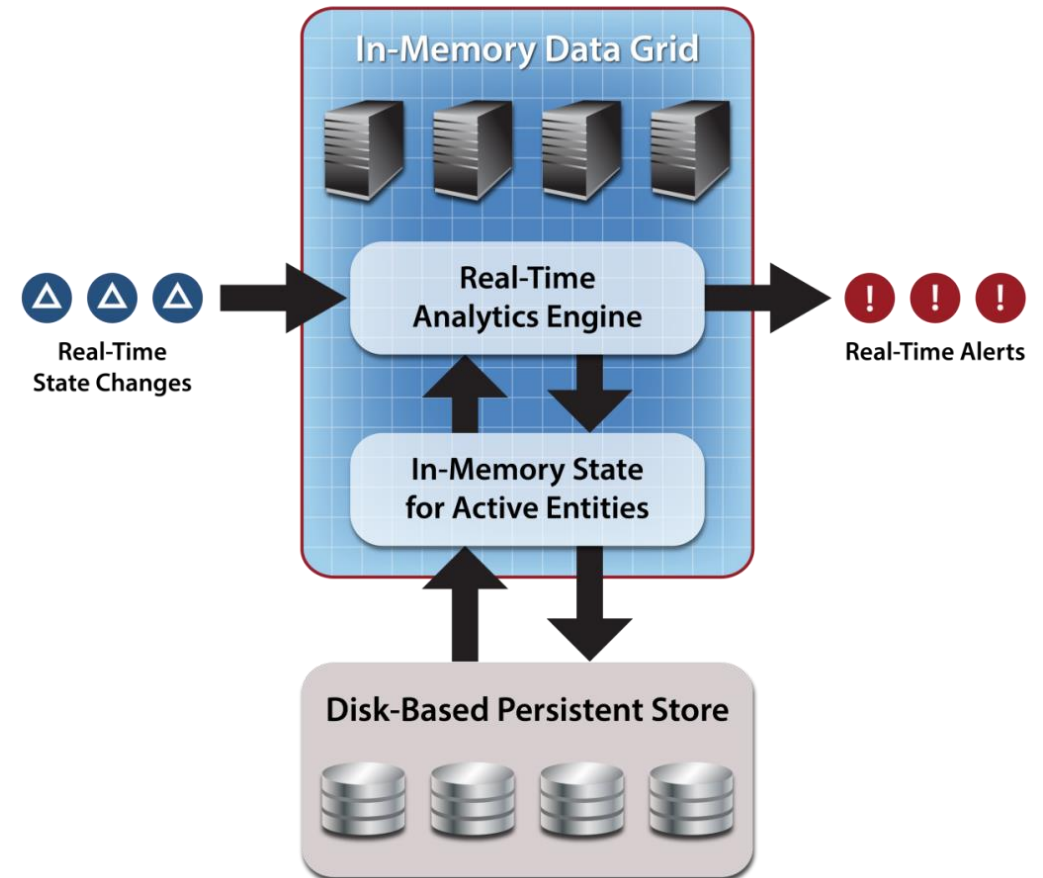


- Using YARN, an IMDG can be directly integrated into a BI cluster:
 - The IMDG holds fast-changing data.
 - YARN directs MapReduce jobs to the IMDG.
 - The IMDG can output results to HDFS.



Recap: In-Memory Computing for OI

- **Online systems need operational intelligence** on “live” data for immediate feedback.
 - Creates important new business opportunities.
- Operational intelligence can be implemented using **standard data-parallel computing techniques**.
- **In-memory data grids provide an excellent platform** for operational intelligence:
 - Model and track the state of a “live” system.
 - Implement high availability.
 - Offer fast, data-parallel computation for immediate feedback.
 - Provide a straightforward, object-oriented development model.





www.scaleoutsoftware.com